

## 1 **The human gut virome database**

2 Ann C. Gregory<sup>1,\*</sup>, Olivier Zablocki<sup>1,\*</sup>, Allison Howell<sup>1</sup>, Benjamin Bolduc<sup>1</sup> & Matthew B. Sullivan<sup>1,2,#</sup>

3

4 <sup>1</sup>Department of Microbiology, Ohio State University, Columbus, OH, United States

5 <sup>2</sup>Department of Civil, Environmental and Geodetic Engineering, Ohio State University, Columbus, OH,  
6 United States

7 # Correspondence to: Matthew Sullivan, [sullivan.948@osu.edu](mailto:sullivan.948@osu.edu)

8 \* These authors contributed equally

9

## 10 **ABSTRACT**

11 The gut microbiome profoundly impacts human health and disease, but viruses that infect these microbes  
12 are likely also important. Problematically, viral sequences are often missed due to insufficient reference  
13 viral genomes. Here we (i) built a human gut virome database, GVD, from 648 viral particle  
14 metagenomes or microbial metagenomes from 572 individuals previously searched for viruses, (ii)  
15 assessed its effectiveness, and (iii) conducted meta-analyses. GVD contains 13,203 unique viral  
16 populations (approximately species-level taxa) organized into 702 novel genera, which roughly doubles  
17 known phage genera and improves viral detection rates over NCBI viral RefSeq nearly 60-fold. Applying  
18 GVD, we assessed and rejected the idea of a ‘core’ gut virome in healthy individuals, and found through  
19 meta-analyses that technical artifacts are more impactful than any ‘treatment’ effect across the entire  
20 meta-study dataset. Together, this foundational resource and these findings will help human microbiome  
21 researchers better identify viral roles in health and disease.

## 22 **Main text**

23 The human gut microbiome is now thought to play an integral role in health and disease<sup>1-4</sup>.  
24 Persistent alterations in the structure, diversity and function of gut microbial communities—dysbiosis—  
25 are increasingly recognized as key contributors in the establishment and maintenance of a growing  
26 number of disease states<sup>5-7</sup>, including obesity<sup>8</sup> and cancer<sup>9</sup>. Gut dysbiosis can develop from complex

27 interplays between host, cognate microbiota and external environmental factors<sup>10,11</sup>. Within the gut  
28 microbial consortium, the bacteriome has been the most extensively studied, where significant shifts in  
29 population dynamics have been observed between healthy and diseased individuals<sup>12</sup>. However,  
30 emerging views<sup>10,13,14</sup> suggest that the gut virome plays an important role in homeostatic regulation and  
31 disease progression through multiple interaction paths with the co-occurring bacteriome, and even  
32 directly with human immune system components<sup>15</sup>.

33         The first step in studying viruses in complex communities is to “see” them. Problematically,  
34 identifying viral sequences in large datasets is notoriously challenging. Because viruses lack a universal  
35 viral marker<sup>16</sup>, as opposed to bacterial 16S rRNA for example, researchers often resort to sequence  
36 homology searches against reference databases (e.g. NCBI viral RefSeq). Such searches are variably  
37 successful with anywhere from 14% to 87% of the observed gut viral genomes having detectable  
38 similarity to viruses in such databases<sup>10</sup>. This large range stems from several factors that are not  
39 mutually-exclusive including the following: (i) broad under-representation of viral genome space in  
40 databases, (ii) non-standardized database usage per study, (iii) overrepresentation of certain virus groups  
41 due to sample preparation and cultured host availability, and (iv) natural sample variation. In addition,  
42 although viral reference datasets are being generated at unprecedented rates<sup>17</sup>, these new data are rarely  
43 incorporated for cross-comparisons, which inflates virus novelty in new datasets and/or leaves many virus  
44 sequences undetected. Therefore, given the rapid accrual of so many studies, there is a need to aggregate  
45 their findings into a central gut-specific database to improve gut virome inference capabilities.

46         Here we collected and curated 648 gut metagenomes from 21 datasets (i.e., any metagenomic  
47 dataset that looked at gut viruses published before 2018), consistently processed them to map known and  
48 unknown viral populations, and used this in multiple meta-analyses to assess improvement and reveal  
49 new biology. The resulting Gut Virome Database (GVD) was born by (i) collecting 648 gut metagenomes  
50 from 572 individuals, (ii) extensive metadata curation through literature mining and, as needed, direct  
51 communication with the original researchers, and (iii) re-analysis of the virome data to establish  
52 consistent processing and extensive virus identification. The value of GVD was assessed for performance

53 against the best currently available databases (NCBI viral RefSeq and IMG/VR<sup>18</sup>), and then used to re-  
54 evaluate global diversity patterns and the relationship between gut virome diversity and diet.

55

## 56 **RESULTS AND DISCUSSION**

57

58 *GVD contains 13,204 viral populations, dominated by phages*

59 To build a collection of the commensal human gut virome, 648 metagenomic samples from 572  
60 individuals were processed from all datasets publicly available as of December 2017 (n=19), along with 2  
61 unpublished datasets where access was granted prior to publication. These studies represented a total of  
62 1.28 Tbp of sequence data derived from a spectrum of gut virome study areas including: (i) healthy gut  
63 viromes of infants<sup>19,20</sup> and adults<sup>21-26</sup>, as well as individuals experiencing (ii) fecal matter transplant, or  
64 FMT<sup>27-31</sup>, (iii) inflammatory bowel disease, or IBD<sup>32,33</sup>, (iv) HIV infection<sup>34</sup>, (v) Type I diabetes<sup>35,36</sup>,  
65 (vi) malnutrition<sup>37</sup>, or (vii) chronic fatigue syndrome<sup>38</sup> (see **Supplementary Table 1**). Datasets had a  
66 worldwide distribution, though most originated from the United States (48.4%; **Fig. 1a**). All reads were  
67 processed consistently, assembled into contigs and viral-like sequence were identified using three  
68 independent methods and validated by cross-comparisons between methods (**Fig. 1b**, see Methods). To  
69 avoid duplicate viral fragments/partial virus genomes across the datasets, contigs were de-replicated by  
70 clustering sequences according to percentage of average nucleotide identity (ANI) and sequence length.  
71 Multiple reports<sup>17,39-43</sup> have revealed that > 95% ANI was a suitable threshold for defining a set of  
72 closely-related discrete ‘viral populations’, with follow-on studies suggesting that this cut-off establishes  
73 populations that are largely concordant with a biologically relevant viral species definition<sup>39,41,44</sup>. Using  
74 this clustering strategy, we identified highly variable numbers of unique viral populations per study  
75 (range: 0 - 3596; mean = 670) (**Supplementary Fig. 1a**). GVD comprises 13,203 viral populations (N50  
76 = 34,220 bp ; L50 = 2,066 bp). For context, NCBI’s viral RefSeq v88 (released May 2018) database holds  
77 8,013 viruses of eukaryotes, bacteria and archaea from all environments, combined. Moreover, if only  
78 comparing phage genomes to the same database, GVD contains 7 times more phages compared to the

79 entire set of cultured phage isolates in viral RefSeq to date. Thus, GVD greatly augments the repertoire of  
80 known viruses in the human gut.

81 Taxonomically, 96.1% of GVD viral populations are bacterial viruses (i.e., phages), with a  
82 minority of GVD viral populations more likely to represent eukaryotic viruses (3.8%) and archaeal  
83 viruses (0.1%) (**Fig. 2a**). Though in the minority, the 505 eukaryotic viruses were taxonomically diverse  
84 (14 families), dominated by ssDNA families *Anelloviridae* (72%), *Genomoviridae* (10%) and  
85 *Circoviridae* (8%). All, with the exception of Genomoviruses, have been reported previously in the  
86 datasets underlying GVD<sup>34</sup>. Among the phages, 82% did not have ICTV classification, with the  
87 remaining fraction comprised of dsDNA tailed phage families (*Siphoviridae*, *Myoviridae* and  
88 *Podoviridae*), *Microviridae* and *Inoviridae* (**see Supplementary Table 2**). Twelve unknown archaeal  
89 viral populations were detected, with no close genome/gene homology to any of the classified archaeal  
90 viruses. The high number of unclassified phages likely results from underrepresentation of gut phages in  
91 the database, coupled to unresolved and/or missing taxonomic assignments for ~ 60% of reference phage  
92 genomes in RefSeq, with the currently classified fraction organized into ~250 genera<sup>45</sup>. To fill this phage  
93 and archaeal virus taxonomic classification gap, we used a genome-based, gene-sharing network strategy  
94<sup>46,47</sup> that *de novo* predicts genus-level groupings ('viral clusters' or 'VCs') from viral population data. A  
95 network was computed from 6,373 GVD phage genomes (only those  $\geq 10$  kb in length; 48% of GVD),  
96 combined with 2,304 curated reference phage genomes from NCBI Viral RefSeq (version 88). The  
97 resulting gene-sharing network (**Fig. 2b**) revealed 957 VCs, 702 of which were novel and exclusively  
98 composed of GVD genomes (3,220 viral genomes or ~51% of GVD genomes). This would roughly  
99 double the current number of ICTV-recognized phage genera. Though not explored here, as our goals  
100 focused on taxonomic classification, the shared protein content within and between VCs calculated in our  
101 network analyses could be used to guide qPCR assays for NGS validation<sup>48</sup> and/or tracking of viruses at  
102 either the viral population- or genera- level under changing conditions<sup>35</sup>.

103 Next, we sought to link phage populations to their hosts using *in silico* strategies (see **Methods**).  
104 The most common identifiable phage hosts (**Fig. 2c**) in GVD belonged the bacterial phylum Firmicutes

105 (38%), about 2-fold more than the next most abundantly identified host phyla (Bacteroidetes and  
106 Proteobacteria; **see Supplementary Table 2**). Though Firmicutes and Bacteroides are the most prominent  
107 bacterial phyla in the human gastrointestinal tract<sup>49</sup>, Firmicutes typically outnumber Bacteroidetes in  
108 unhealthy individuals with metabolic and digestive disorders<sup>50-52</sup>. GVD metagenomes originated from  
109 ~16% healthy individuals and ~84% unhealthy individuals, many of which have metabolic and digestive  
110 disorders. Thus, it is perhaps not surprising that most of the annotated viral populations were linked to the  
111 phylum Firmicutes.

112

### 113 *GVD significantly improves virus detection in all gut datasets*

114 We then quantitatively evaluated virus identification sensitivity (through read mapping) between  
115 multiple databases by comparing the number of identified viral populations in each study detected by  
116 GVD, viral RefSeq v88, IMG/VR 1.1 (2018 release) and the individual virome datasets ('IV') from each  
117 study (**Fig. 3**). For the latter, IV reads were mapped against viral populations (predicted in this study)  
118 derived exclusively from its matching IV. In all datasets, GVD surpassed viral RefSeq (mean increase:  
119 59-fold  $\pm$  95-fold) and IVs (mean increase: 3.2-fold  $\pm$  6.6-fold). In 5 of 18 studies (28%), GVD  
120 outperformed IMG/VR (mean increase: 1.1-fold  $\pm$  2-fold), with the remaining studies finding no  
121 significant difference between or too low of a sample size to compare GVD and IMG/VR. After GVD,  
122 IMG/VR was the next best performing database for viral detection in the gut, as our tests showed an  
123 average of 49-fold ( $\pm$  87-fold) increase over viral RefSeq. IMG/VR was expected to surpass viral RefSeq,  
124 as it aggregates both cultivated reference virus genomes, >12,000 prophages and >700,000 uncultivated  
125 virus genomes/fragments from many environments, including multiple human body sites<sup>53</sup>. Moreover,  
126 given the high performance of IMG/VR in our tests, we wondered about the extent of viral population  
127 overlap with GVD (**Fig. 3b**). There were 1,730 viral populations shared between the two databases, but  
128 still each database is overwhelmingly unique (82% and 69% unique to GVD and IMG/VR, respectively).  
129 This is because IMG/VR includes human gut studies that did not explore the viral fraction as well.  
130 Overall, the significant increase in virus detection by GVD over other databases (two-tailed Mann-

131 Whitney U-tests;  $p$ -value < 0.05) highlights the low representation of gut viruses recorded in RefSeq and  
132 thus demonstrates the value of GVD for sequence-based virus identification in human gut microbiome  
133 datasets. Because the datasets used to compile GVD were originally analyzed most often (55% of the  
134 studies) using viral RefSeq as the primary source to identify viruses (**Supplementary Table 1**), we  
135 wondered whether significant fractions of viruses could have been missed, and whether a possibly  
136 reduced viral “signal” would influence previous conclusions.

137

### 138 *MDA amplification skews diversity and prohibits quantitative analysis of gut viromes*

139 To evaluate this possible reduced viral “signal”, we first examined the role of methodological  
140 approaches in influencing inferences about ssDNA viruses. This is because we noticed that the bulk of  
141 ssDNA eukaryotic viruses (Anelloviruses, Circoviruses, Genomoviruses, Geminiviruses) and phages  
142 (Microviruses) originated from only 4 of the 21 studies gathered in this work (**Fig. 4 a,b**). These studies  
143 evaluated 2 infant gut viromes<sup>19,37</sup> and 2 adult inflammatory bowel disease viromes<sup>31,32</sup>, and they  
144 reported relative abundance shifts of ssDNA and dsDNA phages within these viromes. From this  
145 observation, these studies concluded that such shifts could discriminate between healthy and disease  
146 states associated with virome development in early life.

147 However, the abundance of ssDNA viruses can also be enriched from methodologies used in  
148 making the viromes, even if all samples are processed consistently. Specifically, early virome studies  
149 where limiting viral nucleic acids were obtained, often used whole genome amplification kits that  
150 leverage a DNA polymerase from the phi29 ssDNA virus to obtain many-fold increases in DNA via  
151 multiple displacement amplification or MDA<sup>54</sup>. Though attractive at first, MDA is now known to have  
152 stochastic biases (e.g., 100s –10,000s-fold biases in coverage,<sup>55,56</sup>), which result from randomized initial  
153 template interactions and can induce chimera formation and uneven amplification of linear genomic  
154 sections (whether ssDNA or dsDNA templates), as well as systematic biases resulting from preferential  
155 amplification of small, circular and ssDNA genomes<sup>57–61</sup>. Taken together, MDA-associated artifacts skew  
156 the taxonomic representation of a community in non-repeatable ways and preclude quantitative analysis

157 of viromes<sup>57</sup>. Although non-quantitative, MDA-amplified viromes do still have value enriching for  
158 ssDNA viruses, as well as estimating presence of viruses.

159 Consistent with the idea that these ssDNA viruses are methodologically enriched in the MDA  
160 libraries, we found that non-MDA amplified gut viromes contained significantly less ssDNA viruses than  
161 MDA amplified gut viromes (range: 0% - 4% versus 0-42%; Mann-Whitney U-test;  $p$ -value = 0.0083),  
162 though sample size was quite low. Further, while we see a strong linear relationship ( $R^2 = 0.86$ ) between  
163 sequencing depth and the number of viral populations sequenced in non-MDA viromes, this relationship  
164 is weak in MDA viromes ( $R^2 = 0.39$ ), suggesting that MDA can skew the number of assembled viral  
165 contigs in datasets (**Supplementary Fig. 1b**). Critically, 14 of the 21 studies gathered in this work  
166 employed MDA, which calls into question the quantitative nature of these datasets. Fortunately, viral  
167 nucleic acid extraction from feces often yield sufficient quantities for high throughput sequencing<sup>26</sup>, and  
168 in cases where they do not there are now several viable alternative methods to more quantitatively  
169 establish viromes with as little as 1pg of DNA<sup>61,62</sup>. Problematically, current established gut virome  
170 protocols recommend an MDA step<sup>48,63</sup>. If a researcher's goal is to provide quantitative datasets, then we  
171 strongly advocate against this recommendation and instead suggest that alternative methods<sup>61,62</sup> be used to  
172 generate gut viromes.

173

174 *Human gut virome study conclusions are more impacted by methodology than disease state*

175 Given a systematically processed GVD, we next sought to determine whether global clustering patterns  
176 would emerge between study themes between all dataset used to build GVD. To this end, viral  
177 populations identified in this study were matched back to their respective datasets, and used in a co-  
178 occurrence network analysis (see **Methods**) to assess co-variation at two levels: between study datasets  
179 (**Fig. 4c**), and between viromes across all datasets (**Fig. 4d**). Between datasets, the fraction of shared viral  
180 populations was low (mean: 3%  $\pm$ 3%; **Fig. 4c**), except for 6 datasets that clustered together (hierarchical  
181 clustering bootstrap = 100%; **Fig. 4c**) and had a higher level of shared viral populations (>4-fold  
182 increase). Presumably, these elevated similarities across the 6 datasets may be due to deeper sequencing

183 (Fig. 4d, top panel) that allowed deeper sequencing into the rare tail of viral populations among samples.  
184 A similar trend was observed when looking at the level of individuals within each study (**Fig. 4d**), where  
185 the co-occurrence network revealed close clustering between individuals derived from the same study,  
186 irrespective of geographical origin, health status and/or diet. This per study clustering implies that, taken  
187 together, these studies are not comparable likely due inconsistent sampling and extraction methodologies.  
188 We then investigated the prevalence of gut viral populations amongst all samples, so as to establish  
189 whether any viral populations were detected in all samples (i.e., a ‘core’ gut virome<sup>22</sup>). On average,  $138 \pm$   
190  $170$  (average  $\pm$  SD; range: 0 to 849) viral populations were detected per sample, but not one viral  
191 population was found across all samples. We then explored deeper to detect whether subsets of the  
192 samples would reveal shared viral populations. We found that only 28 viral populations occurred in over  
193 20% of the GVD samples. Most viral populations were detected in very few samples. In fact, >40% of  
194 the viral populations occurred in <0.5% of the samples and 98% of the viral populations occurred in  
195 <0.1% of the samples in GVD (**Fig. 5 a, b** and **Supplementary Table 3**). Further, we specifically looked  
196 at the prevalence of crAssphages, a well-recognized, multi-genera group of phages known to be  
197 widespread in gut viromes<sup>64</sup> (**Fig. 5 b, c**). While crAssphages are ubiquitous across the GVD samples,  
198 there was not one crAssphage viral population found universally, with the most widespread crAssphage  
199 population occurring in only 38% of samples. Importantly, when we looked at all healthy samples and  
200 healthy western samples specifically, still no shared viral populations were identified in all samples.  
201 (**Supplementary Fig. 2a, b**). Assuming samples were sufficiently sequenced, this may be indicative that  
202 individuals carry a unique ‘gut virome fingerprint’, even between twins, which is perhaps not surprising  
203 given recent suggestions of a similar ‘fingerprint’ for gut microbes (the ‘personal’ microbial microbiomes  
204 <sup>65</sup>). This apparent lack of core gut virome among individuals contrasts with a recent report <sup>22</sup>, in which  
205 overlapping patterns of phage genomes between 2 unrelated healthy individuals, as well as within a re-  
206 analyzed larger cohort <sup>66</sup> revealed three levels of sharing patterns: (i) core (phage found in >50% of  
207 samples, (ii) common (phage found in >20-50% of samples), and (iii) unique (phage found in <20% of  
208 samples). Our analyses showed no viral populations shared above >50% of samples, thus bringing into



209 question the presence of a ‘core’ virome as previously defined<sup>22</sup>, as well as a very limited ‘common’  
210 virome (20-50% sharing across samples), in which we observed either 1% (all healthy; n=132) or 0.1%  
211 (all healthy Westerners; n=18) of GVD viral populations, similar to the 3% previously reported<sup>22</sup> (see  
212 **Supplementary Table 4**). Likely, this discrepancy with our results could be attributed to how viruses  
213 were identified through read mapping. In the initial study reporting a core virome<sup>22</sup>, a virus was  
214 considered present if a single read mapped to a genome, a very permissive cut-off which does not take  
215 into account shared homologous regions between distinct viral populations. In this study, we considered  
216 a virus present if reads mapped 70% of the genome length (if genome is <5kb) or reads mapped at least  
217 5kb of the genome (for genome >5kb in length) (see Methods). While our cut-off is more conservative, it  
218 better ensures that we are detecting the same viral population. Nonetheless, the idea of a core virome  
219 might still be an open question.

220

221 *Re-evaluation of a previous study: the virome across different geographic regions and lifestyles*

222 Due to the high level of sample clustering per study (**Fig. 4c**), we were unable to conduct cross-  
223 study analyses. Instead, we sought to assess if the virome community patterns between populations of  
224 varying lifestyles (industrialized versus semi-industrialized versus hunter-gatherer) would vary between  
225 the initial study<sup>26</sup> or GVD-based, to test whether there were geographic biases around GVD viral  
226 populations, and how well sampled are the different geographic regions. This initial study encompassed a  
227 globally-distributed dataset (USA, Italy, Tanzania and two Peruvian populations: Tunapuco and Matses;  
228 **Fig. 6a**), and explored the impact of geography and diet on eukaryotic gut viruses (but did not include  
229 phages) and found that the hunter-gatherers (Hadza in Tanzania and the Matses in Peru) had the highest  
230 eukaryotic viral richness<sup>26</sup>.

231 In this re-analysis, however, we included phages in addition to eukaryotic viruses, and focused on  
232 how the virome diversity varied along the dataset. We first evaluated whether per-region GVD-mediated  
233 detection of viruses would incur biases, potentially stemming from underrepresented viral populations  
234 from less-sampled geographical regions. This did not appear to be the case, as significant increases in

235 virus detection were observed across 4 out of the 5 regions sampled (**Fig. 6b**). We next calculated  
236 diversity indices (**Fig. 6c** and **Supplementary Fig. 3**) for each regional dataset, and looked at the number  
237 of viral populations mapped with GVD. Overall, we reached a similar conclusion to the initial study (even  
238 when considering phages), in which the hunter-gatherers (Peru Matses) generally contained higher viral  
239 richness (**Fig. 6c - left**) and biodiversity (Shannon's H, **Fig. 6c - middle**), but not higher evenness  
240 (Peilou's J, **Fig. 6c -right**). Collector's curves revealed that we have not saturated the human gut viral  
241 diversity among individuals globally (**Supplementary Fig. 4**) or even among just among American  
242 samples (**Supplementary Fig. 2, inset**). Thus, it appears much more viral diversity remains to be  
243 discovered across all geographic regions.

244 We next wondered whether the addition of phage in our analysis would reflect on overall viral  
245 community similarities by using Bray-Curtis distances between individuals across these geographic and  
246 lifestyle gradients (**Fig. 6d**). While unequal database representation can have an impact on alpha-  
247 diversity, beta-diversity is often less impacted<sup>67</sup>. Principal coordinate analyses (PCoA) of Bray-Curtis  
248 distances derived from using the individual Rampelli et al., 2017 virome database (**Fig. 6d**, left panel) and  
249 GVD (**Fig. 6d**, right panel) revealed no significant differences (Mantel's test;  $R = 0.95$ ,  $p = 0.001$ ).  
250 However, analysis of the GVD-referenced PCoA revealed individuals with the same lifestyle and from  
251 the same region clustered together (PERMANOVA;  $p \leq 0.001$ ) and provided better resolution of the  
252 clustering in comparison to the IV-referenced PCoA. However, lifestyle alone may not account for the  
253 observed clustering patterns. The viromes of the Hadza in Tanzania and semi-industrialized, agrarian  
254 Tunapuco population in Peru strongly overlapped (hierarchical clustering bootstrap = 100%; **Fig. 6d**),  
255 most likely driven by their diets rich in root vegetables<sup>68-70</sup>. Nonetheless, when we look at differences  
256 between dominant viral populations (found in >50% individuals) across these geographic and lifestyle  
257 gradients, we see that there are key viruses missing from Western, industrialized gut viromes (**Fig. 6e**),  
258 specifically viruses that infect the genus *Prevotella* spp. This parallels the bacterial analyses that show  
259 that *Prevotella* spp. are enriched in non-Western gut microbiomes and many species are missing from

260 Western, industrial gut microbiomes<sup>69–71</sup>. Overall, this suggests that lifestyle and diet has an impact not  
261 only on the bacterial community, but also on the viral community in the gut.

## 262 CONCLUSIONS

263 The lack of a curated database for the detection of viral sequences in the human gut has been  
264 identified as the most critical shortcoming of applying metagenomic approaches to studying the human  
265 gut virome<sup>72</sup>. Although GVD is geared towards filling this gap and performs well (increasing viral  
266 detection 59-fold over the most commonly used database, NCBI viral RefSeq), there are limitations. *First*,  
267 the geographic and ethnic representation across the dataset is not very broad. Meta-analyses will benefit  
268 from more broadly representative datasets. *Second*, GVD was built using all datasets available by the end  
269 of 2017. Since then, as of May 2019, there are 11 additional datasets that study the gut virome, 8 of which  
270 use viral particle-enriched metagenomes (**Supplementary Table 5**). Further, there are many more human  
271 gut microbial metagenomic datasets and these could be a rich source for virus reference genomes as found  
272 for soils<sup>73</sup> and the large-scale Earth Virome study<sup>74</sup>. To maintain significance as a resource, we will  
273 update GVD annually by extracting the viral signal from such gut-related datasets, as well as monitoring  
274 IMG/VR for gut-related viruses that should be integrated. *Third*, GVD is accessible through direct  
275 download as a single fasta file containing all GVD viral populations (see link in the ‘Data availability’  
276 statement below), and is likely best paired with IMG/VR to maximize viral signal recovery. Future GVD  
277 updates and development will be required to improve the user experience for those not comfortable at  
278 command-line interfaces, but these are likely best integrated with large-scale standardizing efforts like the  
279 National Microbiome Initiative.

280 Given the relatively minimal value added via non-quantitative MDA-based approaches and the  
281 availability now of low-input quantitative approaches pioneered studying ocean viruses<sup>61,75</sup> suggest that  
282 gut virome studies should move away from the former towards the latter. GVD, combined with the means  
283 to classify uncultivated virus genomes<sup>47</sup>, are prime starting requirements for enabling ecosystem-wide  
284 examinations<sup>76</sup> of the dynamics and impacts of the virome within the human gut. Other environmental  
285 advances also invite such studies to include assessing the role of micro- and macro-diversity on virus

286 persistence<sup>41</sup>, and metabolic reprogramming via virus-encoded auxiliary metabolic genes<sup>73,76</sup>. These  
287 combined efforts are critical to enable studies of the human gut virome to advance from ‘stamp  
288 collecting’ diversity studies towards the kinds of comprehensive efforts needed to incorporate viruses into  
289 mechanistic, predictive models. Such efforts, with future viral mapping outside the gut to parallel efforts  
290 for the ‘non-gut’ human microbiome<sup>77</sup>, should help transform personalized medicine and lead to a better  
291 understanding of human ecosystems.

292

293 **FIGURE LEGENDS:**

294

295 **Figure 1. Overview of studies and meta-analyses comprising the Gut Viral Database (GVD). (a)**

296 Global heatmap of the world showing the number of individual's gut viromes coming from different  
297 countries within the GVD. Importantly, individual's viromes coming from the Cameroon were pooled  
298 based on their location, age, and contact with bats. The pools were counted as a single individual's virome  
299 for our analyses. **(b)** Pipeline for the selection and processing of human gut virome datasets (see  
300 **Methods**). Datasets were processed individually and, within each dataset, viromes were pooled by  
301 individual, except for fecal microbiota transfer (FMT) studies and data that was given to us prior to  
302 publication (Yinda et al., 2019; Neto et al. (unpublished)). Reads were filtered for quality and trimmed  
303 and reads that mapped to  $\Phi$ x174 and the human genome were removed. The remaining reads were  
304 assembled into scaffolds, filtered for lengths  $\geq 1.5$ kb, and run through tools that collectively utilize  
305 homology to viral reference databases, probabilistic models on viral genomic features, and viral  $k$ -mer  
306 signatures to identify viral contigs. Viral contigs were then deduplicated to get a total of 13,203 viral  
307 populations.  
308

309 **Figure 2. The Gut Viral Database (GVD). (a)** Pie charts showing the number of bacteriophages,  
310 eukaryotic viruses, and archaeal viruses in the GVD (center) and their familial taxonomic composition by  
311 the bacteriophages (left) and the eukaryotic viruses (right). **(b)** Gene-sharing taxonomic network of the  
312 GVD, including viral RefSeq viruses v88. RefSeq viruses are highlighted in red. Every node represent a  
313 virus genome, while connecting edges identify significant gene-sharing between genomes, which form the  
314 basis for their clustering in genus-level taxonomy. **(c)** Bar chart showing the number of bacterial host  
315 phyla of the GVD bacteriophages, with an inset providing resolution for the low frequency bacteria host  
316 phyla. Putative host phyla per each bacteriophage population are in **Supplementary Table 3**  
317

318 **Figure 3. GVD as a reference database increases viral population detection. (a)** Boxplots showing  
319 median and quartiles of the number of viral populations detected per study using the IV, Viral Refseq  
320 v88, JGI IMG/VR, or GVD databases. Studies where the reads were given to us prior to publication are  
321 excluded from this analysis (Yinda et al., 2019; Neto et al. (unpublished)). **(b)** Venn diagram showing the  
322 number of viral populations unique and shared between the different databases. Importantly, we only  
323 compared dereplicated viral populations from IMG/VR that came directly from human gut samples or had  
324 reads mapping to them from GVD gut samples.  
325

326 **Figure 4. Individual Viromes (IV) Study Databases and Cross-Study Comparisons. (a)** Barplot  
327 showing the proportion of those viruses that are bacteriophages, archaeal viruses, or eukaryotic viruses.  
328 The total number of assembled viral contigs and viral populations per study are available in  
329 **Supplementary Fig. 1a. (b)** Barplot showing the proportion of those viruses that are dsDNA, ssDNA, or  
330 RNA viruses. Studies where multiple displacement amplification (MDA) was used show a higher  
331 prevalence of ssDNA viruses. No viral contigs  $\geq 1.5$ kb were assembled from the Reyes *et al.* 2010 study.  
332 **(c)** Hierarchically clustered heatmap showing the number of viral populations shared within and between  
333 studies. The barplot on top of the heatmap shows the total number of sequenced base pairs following  
334 quality control within each study. **(d)** Viral population co-occurrence network per individual within each  
335 study shows that individuals within a study cluster together regardless of health status. The squares  
336 represent the healthy individuals within each study.  
337

338 **Figure 5. There are no core viral populations across GVD samples. (a)** Histogram showing the  
339 number of viral populations present in different percentages of GVD samples. The vast majority of viral  
340 populations are found in  $<10\%$  of the individuals. **(b)** Hive plot showing the percentage of GVD sample  
341 each viral population is detected within. The dots on the x-axis represent each GVD viral population in  
342 ascending order of the percentage of GVD samples that they are found within. The y-axis is the  
343 percentage of GVD samples that each viral population is detected within. CrAssphage viral populations

344 are highlighted in red. (c) Heatmap showing the presence or absence of each crAssphage viral population  
345 across the different GVD samples.

346  
347 **Figure 6. Diet and geography widely influence gut virome.** (a) World map showing the geographical  
348 distribution of the Rampelli *et al.*, 2017 dataset. (b) Boxplots showing median and quartiles of the number  
349 of viral populations detected using the GVD database and the Rampelli *et al.*, 2017 viral database alone  
350 (IV) within each geographic group. (c) Boxplots showing median and quartiles of the  $\alpha$ -diversity metrics  
351 – richness, Shannon’s H and Peilou’s J – across the different geographic groups using the GVD database  
352 (see **Fig. S3** for  $\alpha$ -diversity metrics using the IV database). (d) Principal coordinate analysis (PCoA) of a  
353 Bray-Curtis dissimilarity matrix calculated from mapping the Rampelli *et al.*, 2017 dataset against the IV  
354 (left) and GVD (right) databases. Analyses show that the viromes significantly (Permanova  $p < 0.05$ )  
355 structure into based on the geographic groups, with mapping to the GVD showing revealing much  
356 stronger clustering based on geography. Ellipses in the PCoA plot are drawn around the centroids of each  
357 group at a 95% confidence interval. The dashed lines connecting the different points reveal the  
358 connections determined by hierarchically clustering between the different samples. (e) Heatmap of the  
359 abundances of the viral populations found across >50% of individuals within the study. Individuals on a  
360 Western diet (from the USA and Italy) lack phages that infect Bacteroidetes, specifically those that infect  
361 *Prevotella* sp. All pairwise comparisons were performed using a two-tailed Mann-Whitney U-tests.

362  
363 **Supplementary Figure 1. Number of assembled viral contigs and populations.** (a) Barplot showing  
364 the number of assembled viral contigs versus the number of deduplicated viral populations per study. (b)  
365 Scatterplots with linear regressions showing the impact of increased sequencing on the number of  
366 assembled contigs per study divided by studies that did not have multiple displacement amplification  
367 (MDA; **top**) and those that did have MDA (**bottom**).

368 **Supplementary Figure 2. There are no core viral populations across healthy samples and across**  
369 **healthy western samples.** Hive plots showing the percentage of GVD samples each viral population is  
370 detected within across (A) all healthy individuals and (B) across only healthy western adults. The dots on  
371 the x-axis represent each GVD viral population in ascending order of the percentage of GVD samples that  
372 they are found within. The y-axis is the percentage of GVD samples that each viral population is detected  
373 within.

374  
375 **Supplementary Figure 3.** Boxplots showing median and quartiles of the  $\alpha$ -diversity metrics – richness,  
376 Shannon’s H and Peilou’s J – across the different geographic groups in the Rampelli *et al.* 2017 study  
377 using the IV and GVD databases.

378  
379 **Supplementary Figure 4. The number of gut viral populations will still increase with more samples**  
380 **added to GVD.** Collector’s curve for gut viral populations in the GVD. (**inset**) Collector’s curve for just  
381 the viromes from samples from the USA.

382  
383  
384 **Supplementary Table 1.** Origin of datasets and associated metadata used to create the gut virome  
385 database.

386  
387 **Supplementary Table 2.** Gut Viral Database contigs family-level taxonomy and putative hosts.

388  
389 **Supplementary Table 3.** Distributions of viral populations across GVD samples.

390  
391 **Supplementary Table 4.** Core, common, low-overlap, and unique GVD viral populations

392  
393 **Supplementary Table 5.** Human gut virome studies since the end of 2017

## 394 METHODS

395 **Experimental Model and Subject Details.** Gut virome database (GVD) studies were selected by doing a  
396 thorough and manually curated search of the Web of Science Core Collection of Thomson Reuters for  
397 studies looking at viruses in the gut published prior to 2018. All studies that used next-generation  
398 sequencing and looked for viruses within the gut microbiome were selected to be part of GVD (see full  
399 list of studies in **Supplementary Table 1**). Additionally, we were given access to the reads of two studies  
400 that were unpublished at the time. One of the studies, however, is now published (Yinda et al., 2019).

401  
402 **Viral contig assembly, identification, and dereplication.** Previously published GVD reads were  
403 downloaded from their respective hosting databases (e.g. SRA, iVirus, or MG-RAST). Prior work  
404 revealed that an individual's gut virome is stable across time (Minot et al., 2013), so reads were pooled  
405 per individual regardless of the number of time points, with a few exceptions (**Fig. 1**). These exceptions  
406 included studies with fecal microbiota transfers and studies whose reads were given to us prior to  
407 publication. For fecal microbiota transfers, all time points per individual were kept separate and processed  
408 independently. Read sets from two studies were given to us prior to publication (Yinda et al., 2019; Neto  
409 et al., unpublished). For the Yinda *et al.*, 2019 study, individual's reads were pooled based on their  
410 location, age, and contact with bats. The pools were counted as a single individual's virome for our  
411 analyses. For the Nadia et al., (unpublished), all reads from all individuals were pooled together. A global  
412 map showing the number of individuals (or pooled read sets) originating from each country was created  
413 using the R packages 'rworldmap.'

414 Pooled reads were then assembled using metaSPAdes 3.11.1<sup>78</sup>. Following assembly, contigs  
415  $\geq 1.5$ kb were piped through VirSorter<sup>79</sup> and VirFinder<sup>80</sup> and those that mapped to the human, cat or dog  
416 genomes were removed. For viral-enriched metagenomes (i.e. viromes), contigs  $\geq 5$ kb or  $\geq 1.5$ kb and  
417 circular that were sorted as VirSorter categories 1-6 and/or VirFinder score  $\geq 0.7$  and  $p < 0.05$  were pulled  
418 for further investigation. Of these contigs, those sorted as VirSorter categories 1 and 2, VirFinder score  
419  $\geq 0.9$  and  $p < 0.05$  or were identified as viral by both VirSorter (categories 1-6) and VirFinder (score  $\geq 0.7$   
420 and  $p < 0.05$ ) were classified as viral. The remaining contigs were run through CAT<sup>81</sup> and those with  
421  $< 40\%$  (based on an average gene size of 1000) of the genome classified as bacterial, archaeal, or  
422 eukaryotic were considered viral. For the microbial metagenomes, we took a more conservative approach  
423 with only contigs  $\geq 5$ kb or  $\geq 1.5$ kb and circular that were sorted as VirSorter categories 1-2 and VirFinder  
424 score  $\geq 0.6$  and  $p < 0.05$  were considered viral. Across the both the viral-enriched and microbial  
425 metagenomes, contigs  $\geq 5$ kb or  $\geq 1.5$ kb and circular that were classified as eukaryotic viral contigs by  
426 CAT were also considered viral. In total, 29,345 viral contigs were identified.

427 Viral contigs that were from known ssDNA or RNA viral families using CAT were grouped into  
428 populations if they shared  $\geq 95\%$  nucleotide identity across  $\geq 100\%$  of the genome. Because there are no  
429 benchmarked metagenomic population boundaries for ssDNA and RNA viral families, we chose to not  
430 use stringent dereplication. All other contigs were considered double-stranded DNA and were grouped  
431 into populations if they shared  $\geq 95\%$  nucleotide identity across  $\geq 70\%$  of the genome (*sensu*<sup>82</sup>) using  
432 nucmer<sup>83</sup>. All the viral contigs that were assembled were dereplicated per study to create the individual  
433 virome (IV) databases and across all of GVD (see **Supplementary Fig. 1**). For GVD, this resulted in  
434 13,203 total viral populations found in GVD (see **Supplementary Table 3** for VirSorter, VirFinder, and  
435 CAT results), of which 6,373 were  $\geq 10$ kb in length.

436

437 **Core Viral Population Analyses.** To explore if there were any core viral populations, the abundance  
438 table was turned into a binary presence-absence matrix. The number of GVD samples that each viral  
439 population was detected within was then calculated using R and divided by the total number (648) to get  
440 the percentage of samples. Each viral population's percentage was plotted in hive plot using  
441 'geom\_curve' in ggplot2<sup>84</sup>. This process was repeated on subsets of the matrix including all healthy  
442 individuals and only the healthy western adults. The number of viral populations that were present across  
443 different percentages were calculated using R and their distributions plotted using 'geom\_histogram' in  
444 ggplot2<sup>84</sup>. CrAssphage viral populations in GVD were identified using CAT results and by dereplicating  
445 GVD viral populations with the crAssphage genomes identified in Guerin *et al.*<sup>64</sup> and seeing which GVD  
446 genomes cluster. In total, there were 95 unique crAssphage populations. The binary presence-absence  
447 data for the crAssphage populations were plotted using pheatmap in R.

448  
449 **Viral taxonomy.** For each viral population, ORFs were called using Prodigal<sup>85</sup> and the resulting protein  
450 sequences were used as input for vConTACT2<sup>47</sup> and for BLASTp. Double-stranded DNA viral  
451 populations represented by contigs >10kb were clustered with Viral RefSeq release 88 viral genomes  
452 using vConTACT2. Those that clustered with a virus from RefSeq based on amino acid homology based  
453 on DIAMOND<sup>86</sup> alignments were able to be assigned to a known viral taxonomic genera. For viral  
454 dsDNA populations that could not be assigned taxonomy or were <10kb, family level taxonomy was  
455 assigned using a majority-rules approach, where if >50% of a genome's proteins were assigned to the  
456 same viral family using a blastp bitscore  $\geq 50$  with a Viral RefSeq virus, it was considered part of that  
457 viral family (see **Supplementary Table 3** for family-level taxonomy). For ssDNA and RNA viruses,  
458 CAT was used to assign the viral family (see **Supplementary Table 3** for family-level taxonomy).

459  
460 **Viral Host Prediction.** Bacteriophage hosts were predicted using a variety of bioinformatic methods  
461 including: (i) CRISPR-spacer matches, (ii) prophage blasts, (iii) tRNA genes matches, and (iv) WiSH  
462 matches<sup>87</sup> against Bacterial Refseq v88. CRISPR spacers were predicted using MinCED  
463 (<https://github.com/ctSkennerton/minced>) and the CRISPR Recognition Tool (CRT<sup>88</sup>) and a BLASTn (-  
464 task blastn-short -word\_size 5) was used to assess matches between the CRISPR spacers and viral  
465 populations in GVD. Those with 1 mismatch were considered a match. For prophage blasts, a blastn of  
466 the viral population against Bacterial RefSeq was performed. A bacterial genome with  $\geq 2500$ bp regions  
467 of their genome matching at 95%ID with a viral population genome were considered putative hosts of that  
468 viral population (see<sup>76</sup>). Viral tRNA genes and Bacterial RefSeq tRNA genes were predicted using  
469 tRNA-scan<sup>89</sup> and then a blastn was performed between the viral and bacterial tRNA genes. Bacterial  
470 tRNA genes that matched viral tRNA genes at 95% ID across 100% of the length were considered  
471 putative bacterial hosts. Lastly, WiSH was used to predict hosts according to default settings<sup>87</sup>. Priority  
472 host assignment was given to CRISPR, then prophage, WiSH and tRNA results. Viruses with putative  
473 archaeal hosts were predicted using MarVD<sup>90</sup>. Viruses with predicted eukaryotic hosts were assigned  
474 based on their assigned taxonomic viral family.

475  
476 **Detecting viral populations and calculating their raw abundances.** To calculate the raw abundances of  
477 the different viral populations in each sample, reads from each GVD pooled read set were first non-  
478 deterministically mapped to all GVD viral population genomes using bowtie2. Further, reads from each  
479 GVD pooled read set per study were mapped to their respective IV databases. BamM  
480 (<https://github.com/ecogenomics/BamM>) was used to remove reads that mapped at <95% nucleotide



481 identity to the contigs, bedtools genomecov<sup>91</sup> was used to determine how many positions across each  
482 genome were covered by reads, and custom Perl scripts were used to further filter out contigs without  
483 enough coverage across the length of the contig. All contigs  $\leq 5\text{kb}$  in length with  $>70\%$  of the contig  
484 covered were considered detected in the sample. Contigs  $>5\text{kb}$  in length with  $\geq 5\text{kb}$  in length covered were  
485 also considered detected in the sample<sup>92</sup>. BamM was used to calculate the average read depth ('tpmean' -  
486 minus the top and bottom 10% depths) across each detected contig. For the alpha-diversity calculations,  
487 the average read depth was used as a proxy for abundance and normalized by total read number per  
488 metagenome to allow for sample-to-sample comparison. However, because most of the studies in GVD  
489 involved MDA, which can skew abundances, we chose to use only a presence-absence statistic (richness)  
490 for most of our  $\alpha$ -diversity calculations. Collector's curves and the whole GVD and across only American  
491 samples were calculated using the function 'specaccum' in the R 'vegan' package<sup>93</sup>.

492  
493 **Comparisons to IMG/VR, Viral RefSeq v88, and IV databases.** The IMG/VR (1.1.2018 release)  
494 included all viral contigs assembled from different datasets. All of the viral contigs in GVD, Viral Refseq  
495 v88, and IV databases are dereplicated at the population level. In order to make IMG/VR comparable to  
496 GVD, Viral Refseq and IV databases, we needed to dereplicate the IMG/VR database. IMG/VR (1.1.2018  
497 release) is composed of 715,672 contigs. Because dereplication is extremely computationally intensive,  
498 we decided to only focus on dereplicating viral contigs that originated from the human gut and had at  
499 least 1 read from a GVD metagenome map. These IMG/VR viral contigs were then dereplicated using the  
500 same methodology as previously described in the methods section. In total, 29,378 IMG/VR viral contigs  
501 were dereplicated into 6,652 viral populations. GVD pooled read sets were mapped to this IMG/VR  
502 human gut viral population database, Viral RefSeq v88, and the IV databases for each individual study in  
503 GVD. The raw abundances of the different IMG/VR and Viral RefSeq viral populations in each sample  
504 were calculated the same way as described in the previous section. The total number of viral populations  
505 detected per sample per study using the different databases were then plotted and comparative statistics  
506 using the 'ggboxplot' function from the 'ggpubr' package in R.

507 All of the viral populations from GVD, the dereplicated IMG/VR gut-specific dataset, and Viral  
508 Refseq were then dereplicated to see how many viral populations overlapped between databases. The  
509 results were then plotted using the 'VennDiagram' package in R. Importantly, in the dereplication  
510 process, some of the original viral populations in each database may be dereplicated down due to the  
511 presence of a longer viral contig from the same population that links the two together into the same  
512 population. Across the databases, 329, 177, and 459 viral populations were dereplicated in GVD,  
513 IMG/VR, and Viral Refseq, respectively. This is why the total number of populations displayed in the  
514 Venn diagram does not add up to the total number of viral populations in each database.

515  
516 **Clustering studies based on shared viral populations.** To test how studies clustered together, the viral  
517 population presence-absence data from individuals (or pooled read sets) within a study were merged. In  
518 Study 1, individual A had viral population 1, 2, 4, 5 and individual B had viral population 3, then Study 1  
519 had viral populations 1, 2, 3, 4, and 5. The different studies were then assessed for the number of shared  
520 viral populations that were present in both studies. These values were then displayed and hierarchically  
521 clustered using the R 'pheatmap' package and the stability of the hierarchical clusters were assess using  
522 the R 'pvclust' package. The number of shared viral populations between individuals (or pooled read sets  
523 within a sample) were clustered using the R 'SPIEC-EASI' package<sup>94</sup> using the Meinshausen and

524 Bühlmann (MB) method to infer associations between samples based on the shared number of viral  
525 populations. The network was plotted using the R ‘igraph’ package.

526  
527 **Alpha- and Beta-Diversity calculations.** The  $\alpha$ - (Richness, Shannon’s  $H$ , and Peilous’  $J$ ) and  $\beta$ - (Bray-  
528 Curtis dissimilarity) diversity statistics were performed using VEGAN<sup>93</sup> in R. For all studies, except for  
529 Rampelli et al.<sup>26</sup>, only richness was calculated for both abundances based on read mapping to IMG/VR,  
530 Viral Refseq, the IV databases and GVD. Comparisons were plotted using ‘ggboxplot’ function in the R  
531 ‘ggpubr’ package. The Rampelli et al.<sup>26</sup> did not use MDA, so we went ahead with scaling the raw  
532 abundances based on the number of quality controlled base pairs sequenced to normalize the data. All  $\alpha$ -  
533 diversity statistics were calculated and  $\beta$ -diversity was used to look at community structure using both the  
534 IV and GVD databases. Principal Coordinate analysis (function capscale of VEGAN package with no  
535 constraints applied) was used as the ordination method to plot the Bray-Curtis dissimilarity matrices  
536 (function vegdist; method “bray”) after a cube root transformation (function nthroot; n = 3). To  
537 determine if the Rampelli et al. samples clustered by geographic region, a permanova test (function  
538 “adonis”) and the 95% confidence interval were plotted using function “ordiellipse.” Further, the samples  
539 were hierarchically clustered and plotted within the PCoA. To specifically look at abundance differences  
540 in the most abundant viral populations in the Rampelli et al.<sup>26</sup> study, viral populations that were present  
541 in 50% study individuals and their hosts information were plotted using the R ‘pheatmap’ package.

542  
543 **Code availability.** Scripts used in this manuscript are available on the Sullivan laboratory bitbucket under  
544 Gut\_Virome\_Database.

545  
546 **Data availability.** All raw reads are available through SRA, iVirus, or MG-RAST using the identifiers

547 listed in **Supplementary Table 1**. GVD viral populations can be downloaded directly from iVirus

548 through the following link: <https://de.cyverse.org/dl/d/E83EFBFF-2A23-4794-8819->

549 [ADD34160D018/FINAL\\_Gut\\_Viral\\_Database\\_GVD\\_1.7.2018.fna](https://de.cyverse.org/dl/d/E83EFBFF-2A23-4794-8819-ADD34160D018/FINAL_Gut_Viral_Database_GVD_1.7.2018.fna)

550

## 551 REFERENCES

- 552 1. Clemente, J. C., Ursell, L. K., Parfrey, L. W. & Knight, R. The impact of the gut microbiota on  
553 human health: An integrative view. *Cell* (2012). doi:10.1016/j.cell.2012.01.035
- 554 2. Lynch, S. V. & Pedersen, O. The Human Intestinal Microbiome in Health and Disease. *N. Engl. J.*  
555 *Med.* **375**, 2369–2379 (2016).
- 556 3. Gilbert, J. A. *et al.* Current understanding of the human microbiome. *Nat. Med.* (2018).  
557 doi:10.1038/nm.4517
- 558 4. Schmidt, T. S. B., Raes, J. & Bork, P. The Human Gut Microbiome: From Association to  
559 Modulation. *Cell* (2018). doi:10.1016/j.cell.2018.02.044
- 560 5. Frank, D. N. *et al.* Molecular-phylogenetic characterization of microbial community imbalances in

- 561 human inflammatory bowel diseases. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 13780–5 (2007).
- 562 6. Huttenhower, C. & Human Microbiome Project Consortium. Structure, function and diversity of  
563 the healthy human microbiome. *Nature* **486**, 207–14 (2012).
- 564 7. Qin, J. *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*  
565 **490**, 55–60 (2012).
- 566 8. Turnbaugh, P. J. *et al.* An obesity-associated gut microbiome with increased capacity for energy  
567 harvest. *Nature* **444**, 1027–1031 (2006).
- 568 9. Yoshimoto, S. *et al.* Obesity-induced gut microbial metabolite promotes liver cancer through  
569 senescence secretome. *Nature* **499**, 97–101 (2013).
- 570 10. Mirzaei, M. K. & Maurice, C. F. M $\acute{e}$ nage  $\grave{a}$  trois in the human gut: Interactions between host,  
571 bacteria and phages. *Nature Reviews Microbiology* **15**, 397–408 (2017).
- 572 11. Schreiner, Andrew B., John Y. Kao, and V. B. Y. The gut microbiome in health and in disease.  
573 *Curr. Opin. Gastroenterol.* **31**, 69–75 (2015).
- 574 12. Zhang, Y.-J. *et al.* Impacts of Gut Bacteria on Human Health and Diseases. *Int. J. Mol. Sci.* **16**,  
575 7493–7519 (2015).
- 576 13. Ogilvie, L. A. & Jones, B. V. The human gut virome: A multifaceted majority. *Frontiers in*  
577 *Microbiology* **6**, (2015).
- 578 14. Tetz, G. V. *et al.* Bacteriophages as potential new mammalian pathogens. *Sci. Rep.* **7**, (2017).
- 579 15. Keen, E. C. & Dantas, G. Close Encounters of Three Kinds: Bacteriophages, Commensal Bacteria,  
580 and Host Immunity. *Trends Microbiol.* **26**, 943–954 (2018).
- 581 16. Rohwer, F. & Edwards, R. The phage proteomic tree: A genome-based taxonomy for phage. *J.*  
582 *Bacteriol.* **184**, 4529–4535 (2002).
- 583 17. Roux, S. *et al.* Minimum Information about an Uncultivated Virus Genome (MIUViG): a  
584 community consensus on standards and best practices for describing genome sequences from  
585 uncultivated viruses. *Nat. Biotechnol.* **37**, 29–37 (2018).
- 586 18. Paez-Espino, D. *et al.* IMG/VR: a database of cultured and uncultured DNA Viruses and  
587 retroviruses. *Nucleic Acids Res.* **45**, D457–D465 (2017).
- 588 19. Lim, E. S. *et al.* Early life dynamics of the human gut virome and bacterial microbiome in infants.  
589 *Nat. Med.* **21**, 1228–1234 (2015).
- 590 20. Reyes, A. *et al.* Viruses in the faecal microbiota of monozygotic twins and their mothers. *Nature*  
591 **466**, 334–338 (2010).
- 592 21. Ly, M. *et al.* Transmission of viruses via our microbiomes. *Microbiome* **4**, 64 (2016).
- 593 22. Manrique, P. *et al.* Healthy human gut phageome. *Proc. Natl. Acad. Sci.* **113**, 10400–10405  
594 (2016).

- 595 23. Minot, S. *et al.* The human gut virome: Inter-individual variation and dynamic response to diet.  
596 *Genome Res.* **21**, 1616–1625 (2011).
- 597 24. Minot, S., Grunberg, S., Wu, G. D., Lewis, J. D. & Bushman, F. D. Hypervariable loci in the  
598 human gut virome. *Proc. Natl. Acad. Sci.* **109**, 3962–3966 (2012).
- 599 25. Minot, S. *et al.* Rapid evolution of the human gut virome. *Proc. Natl. Acad. Sci.* **110**, 12450–  
600 12455 (2013).
- 601 26. Rampelli, S. *et al.* Characterization of the human DNA gut virome across populations with  
602 different subsistence strategies and geographical origin. *Environ. Microbiol.* **19**, 4728–4735  
603 (2017).
- 604 27. Broecker, F., Klumpp, J. & Moelling, K. Long-term microbiota and virome in a Zürich patient  
605 after fecal transplantation against *Clostridium difficile* infection. *Ann. N. Y. Acad. Sci.* **1372**, 29–  
606 41 (2016).
- 607 28. Broecker, F., Russo, G., Klumpp, J. & Moelling, K. Stable core virome despite variable  
608 microbiome after fecal transfer. *Gut Microbes* **8**, 214–220 (2017).
- 609 29. Chehoud, C. *et al.* Transfer of viral communities between human individuals during fecal  
610 microbiota transplantation. *MBio* **7**, 1–8 (2016).
- 611 30. Kang, D. W. *et al.* Microbiota Transfer Therapy alters gut ecosystem and improves gastrointestinal  
612 and autism symptoms: An open-label study. *Microbiome* **5**, (2017).
- 613 31. Zuo, T. *et al.* Bacteriophage transfer during faecal microbiota transplantation in *Clostridium*  
614 *difficile* infection is associated with treatment outcome. *Gut* **67**, 634–643 (2017).
- 615 32. Norman, J. M. *et al.* Disease-specific alterations in the enteric virome in inflammatory bowel  
616 disease. *Cell* **160**, 447–460 (2015).
- 617 33. Pérez-Brocá, V. *et al.* Study of the viral and microbial communities associated with Crohn’s  
618 disease: A metagenomic approach. *Clin. Transl. Gastroenterol.* **4**, (2013).
- 619 34. Monaco, C. L. *et al.* Altered Virome and Bacterial Microbiome in Human Immunodeficiency  
620 Virus-Associated Acquired Immunodeficiency Syndrome. *Cell Host Microbe* **19**, 311–322 (2016).
- 621 35. Kramná, L. *et al.* Gut virome sequencing in children with early islet autoimmunity. *Diabetes Care*  
622 **38**, 930–933 (2015).
- 623 36. Zhao, G. *et al.* Intestinal virome changes precede autoimmunity in type I diabetes-susceptible  
624 children. *Proc. Natl. Acad. Sci.* 201706359 (2017). doi:10.1073/pnas.1706359114
- 625 37. Reyes, A. *et al.* Gut DNA viromes of Malawian twins discordant for severe acute malnutrition.  
626 *Proc. Natl. Acad. Sci.* **112**, 11941–11946 (2015).
- 627 38. Giloteaux, L., Hanson, M. R. & Keller, B. A. A pair of identical twins discordant for myalgic  
628 encephalomyelitis/chronic fatigue syndrome differ in physiological parameters and gut

- 629 microbiome composition. *Am. J. Case Rep.* **17**, 720–729 (2016).
- 630 39. Gregory, A. C. *et al.* Genomic differentiation among wild cyanophages despite widespread  
631 horizontal gene transfer. *BMC Genomics* **17**, (2016).
- 632 40. Brum, J. R. *et al.* Ocean Viral Communities. *Science (80-. )*. **348**, 1261498-1–11 (2015).
- 633 41. Gregory, A. *et al.* Marine viral macro- and micro-diversity from pole to pole. *Cell* (2019).
- 634 42. Duhaime, M. B. & Sullivan, M. B. Ocean viruses: Rigorously evaluating the metagenomic  
635 sample-to-sequence pipeline. *Virology* **434**, 181–186 (2012).
- 636 43. Duhaime, M. B. *et al.* Comparative omics and trait analyses of marine Pseudoalteromonas phages  
637 advance the phage OTU concept. *Front. Microbiol.* (2017). doi:10.3389/fmicb.2017.01241
- 638 44. Bobay, L. & Ochman, H. Biological species in the viral world. **115**, (2018).
- 639 45. Jang, H. Bin *et al.* Gene sharing networks to automate genome-based prokaryotic viral taxonomy.  
640 *bioRxiv* (2019). doi:<https://doi.org/10.1101/533240>
- 641 46. Bolduc, B. *et al.* vConTACT: an iVirus tool to classify double-stranded DNA viruses that infect  
642 *Archaea* and *Bacteria*. *PeerJ* **5**, e3243 (2017).
- 643 47. Bin Jang, H. *et al.* Taxonomic assignment of uncultivated prokaryotic virus genomes is enabled by  
644 gene-sharing networks. *Nat. Biotechnol.* (2019). doi:10.1038/s41587-019-0100-8
- 645 48. Monaco, C. & Kwon, D. Next-generation Sequencing of the DNA Virome from Fecal Samples.  
646 *BIO-PROTOCOL* (2017). doi:10.21769/bioprotoc.2159
- 647 49. Eckburg, P. B. *et al.* Microbiology: Diversity of the human intestinal microbial flora. *Science (80-.*  
648 *)*. (2005). doi:10.1126/science.1110591
- 649 50. Ott, S. J. *et al.* Reduction in diversity of the colonic mucosa associated bacterial microflora in  
650 patients with active inflammatory bowel disease. *Gut* (2004). doi:10.1136/gut.2003.025403
- 651 51. Ley, R. E. *et al.* Obesity alters gut microbial ecology. *Proc. Natl. Acad. Sci. U. S. A.* (2005).  
652 doi:10.1073/pnas.0504978102
- 653 52. Nicholson, J. K. *et al.* Host-gut microbiota metabolic interactions. *Science* (2012).  
654 doi:10.1126/science.1223813
- 655 53. Paez-Espino, D. *et al.* IMG/VR v.2.0: an integrated data management and analysis system for  
656 cultivated and environmental viral genomes. *Nucleic Acids Res.* gky1127–gky1127 (2018).
- 657 54. Angly, F. E. *et al.* The marine viromes of four oceanic regions. *PLoS Biol.* **4**, 2121–2131 (2006).
- 658 55. Zhang, K. *et al.* Sequencing genomes from single cells by polymerase cloning. *Nat. Biotechnol.*  
659 (2006). doi:10.1038/nbt1214
- 660 56. Woyke, T. *et al.* Assembling the marine metagenome, one cell at a time. *PLoS One* (2009).  
661 doi:10.1371/journal.pone.0005299
- 662 57. Yilmaz, S., Allgaier, M. & Hugenholtz, P. Multiple displacement amplification compromises

- 663 quantitative analysis of metagenomes. *Nat. Methods* **7**, 943–944 (2010).
- 664 58. Kim, K.-H. & Bae, J.-W. Amplification methods bias metagenomic libraries of uncultured single-  
665 stranded and double-stranded DNA viruses. *Appl. Environ. Microbiol.* **77**, 7663–8 (2011).
- 666 59. Marine, R. *et al.* Caught in the middle with multiple displacement amplification: The myth of  
667 pooling for avoiding multiple displacement amplification bias in a metagenome. *Microbiome* **2**, 1–  
668 8 (2014).
- 669 60. Kim, K. H. *et al.* Amplification of uncultured single-stranded DNA viruses from rice paddy soil.  
670 *Appl. Environ. Microbiol.* **74**, 5975–5985 (2008).
- 671 61. Roux, S. *et al.* Towards quantitative viromics for both double-stranded and single-stranded DNA  
672 viruses. *PeerJ* **4**, e2777 (2016).
- 673 62. Duhaime, M. B., Deng, L., Poulos, B. T. & Sullivan, M. B. Towards quantitative metagenomics of  
674 wild viruses and other ultra-low concentration DNA samples: A rigorous assessment and  
675 optimization of the linker amplification method. *Environ. Microbiol.* **14**, 2526–2537 (2012).
- 676 63. Mayo-Muñoz, D. Viral genome isolation from human faeces for succession assessment of  
677 the human gut virome. in *Methods in Molecular Biology* (2018). doi:10.1007/978-1-4939-8682-  
678 8\_8
- 679 64. Guerin, E. *et al.* Biology and Taxonomy of crAss-like Bacteriophages, the Most Abundant Virus  
680 in the Human Gut. *Cell Host Microbe* (2018). doi:10.1016/j.chom.2018.10.002
- 681 65. Franzosa, E. A. *et al.* Identifying personal microbiomes using metagenomic codes. *Proc. Natl.*  
682 *Acad. Sci.* (2015). doi:10.1073/pnas.1423854112
- 683 66. Norman, J. M. *et al.* Disease-specific alterations in the enteric virome in inflammatory bowel  
684 disease. *Cell* **160**, 447–460 (2015).
- 685 67. Lemos, L. N., Fulthorpe, R. R., Triplett, E. W. & Roesch, L. F. W. Rethinking microbial diversity  
686 analysis in the high throughput sequencing era. *J. Microbiol. Methods* (2011).  
687 doi:10.1016/j.mimet.2011.03.014
- 688 68. Yatsunenکو, T. *et al.* Human gut microbiome viewed across age and geography. *Nature* (2012).  
689 doi:10.1038/nature11053
- 690 69. De Filippo, C. *et al.* Impact of diet in shaping gut microbiota revealed by a comparative study in  
691 children from Europe and rural Africa. *Proc. Natl. Acad. Sci. U. S. A.* (2010).  
692 doi:10.1073/pnas.1005963107
- 693 70. Schnorr, S. L. *et al.* Gut microbiome of the Hadza hunter-gatherers. *Nat. Commun.* (2014).  
694 doi:10.1038/ncomms4654
- 695 71. Yatsunenکو, T. *et al.* Human gut microbiome viewed across age and geography. *Nature* **486**, 222–  
696 227 (2012).

- 697 72. Shkoporov, A. N. & Hill, C. Bacteriophages of the Human Gut: The “Known Unknown” of the  
698 Microbiome. *Cell Host Microbe* **25**, 195–209 (2019).
- 699 73. Emerson, J. B. *et al.* Host-linked soil viral ecology along a permafrost thaw gradient. *Nat.*  
700 *Microbiol.* (2018). doi:10.1038/s41564-018-0190-y
- 701 74. Paez-Espino, D. *et al.* Uncovering Earth’s virome. *Nature* **536**, 425–430 (2016).
- 702 75. Solonenko, S. a *et al.* Sequencing platform and library preparation choices impact viral  
703 metagenomes. *BMC Genomics* **14**, 320 (2013).
- 704 76. Roux, S. *et al.* Ecogenomics and potential biogeochemical impacts of globally abundant ocean  
705 viruses. *Nature* **537**, 689–693 (2016).
- 706 77. Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000  
707 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* (2019).  
708 doi:10.1016/j.cell.2019.01.001
- 709 78. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. MetaSPAdes: A new versatile  
710 metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
- 711 79. Roux, S., Enault, F., Hurwitz, B. L. & Sullivan, M. B. VirSorter: mining viral signal from  
712 microbial genomic data. *PeerJ* **3**, e985 (2015).
- 713 80. Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A. & Sun, F. VirFinder: a novel k-mer based tool  
714 for identifying viral sequences from assembled metagenomic data. *Microbiome* **5**, 69 (2017).
- 715 81. Cambuy, D. D., Coutinho, F. H. & Dutilh, B. E. Contig annotation tool CAT robustly classifies  
716 assembled metagenomic contigs and long sequences. *bioRxiv* (2016). doi:10.1101/072868
- 717 82. Brum, J. R. *et al.* Patterns and ecological drivers of ocean viral communities. *Science* (80-. ). **348**,  
718 (2015).
- 719 83. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.* (2004).  
720 doi:10.1186/gb-2004-5-2-r12
- 721 84. Ginevet, C. ggplot2: Elegant Graphics for Data Analysis. *J. R. Stat. Soc. Ser. A (Statistics Soc.*  
722 (2011). doi:10.1111/j.1467-985x.2010.00676\_9.x
- 723 85. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification.  
724 *BMC Bioinformatics* **11**, 119 (2010).
- 725 86. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat.*  
726 *Methods* **12**, 59–60 (2015).
- 727 87. Galiez, C., Siebert, M., Enault, F., Vincent, J. & Söding, J. WISH: who is the host? Predicting  
728 prokaryotic hosts from metagenomic phage contigs. *Bioinformatics* **33**, 3113–3114 (2017).
- 729 88. Bland, C. *et al.* CRISPR recognition tool (CRT): a tool for automatic detection of clustered  
730 regularly interspaced palindromic repeats. *BMC Bioinformatics* **8**, 209 (2007).

- 731 89. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: A program for improved detection of transfer RNA  
732 genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).
- 733 90. Vik, D. R. *et al.* Putative archaeal viruses from the mesopelagic ocean. *PeerJ* **5**, e3428 (2017).
- 734 91. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing genomic  
735 features. *Bioinformatics* (2010). doi:10.1093/bioinformatics/btq033
- 736 92. Ann Gregory, A. C. *et al.* Marine DNA Viral Macro-and Microdiversity from Pole to Pole. *Cell*  
737 (2019). doi:10.1016/j.cell.2019.03.040
- 738 93. Dixon, P. VEGAN, a package of R functions for community ecology. *Journal of Vegetation*  
739 *Science* (2003). doi:10.1111/j.1654-1103.2003.tb02228.x
- 740 94. Kurtz, Z. D. *et al.* Sparse and Compositionally Robust Inference of Microbial Ecological  
741 Networks. *PLoS Comput. Biol.* (2015). doi:10.1371/journal.pcbi.1004226

742

743 **ACKNOWLEDGEMENTS.** Computational support was provided by an award from the Ohio  
744 Supercomputer Center (OSC) to MBS. Study design and manuscript comments from Shini Sunagawa,  
745 Miguelangel Cuenca Vera, Bas E. Dutilh, Ksenia Arkhipova, Pedro Meirelles and Simon Roux are  
746 gratefully acknowledged. Funding was provided by the Gordon and Betty Moore Foundation (#3790) to  
747 MBS and an NIH T32 training grant fellowship (AI112542) to ACG.

748

749 **AUTHOR CONTRIBUTIONS.** A.C.G. collected all datasets and metadata for the study. A.C.G. and  
750 A.H. curated metadata for the study. A.C.G., O.Z., A.H., B.B., M.B.S created the study design, analyzed  
751 the data, and wrote the manuscript. All authors approved the final manuscript.

752

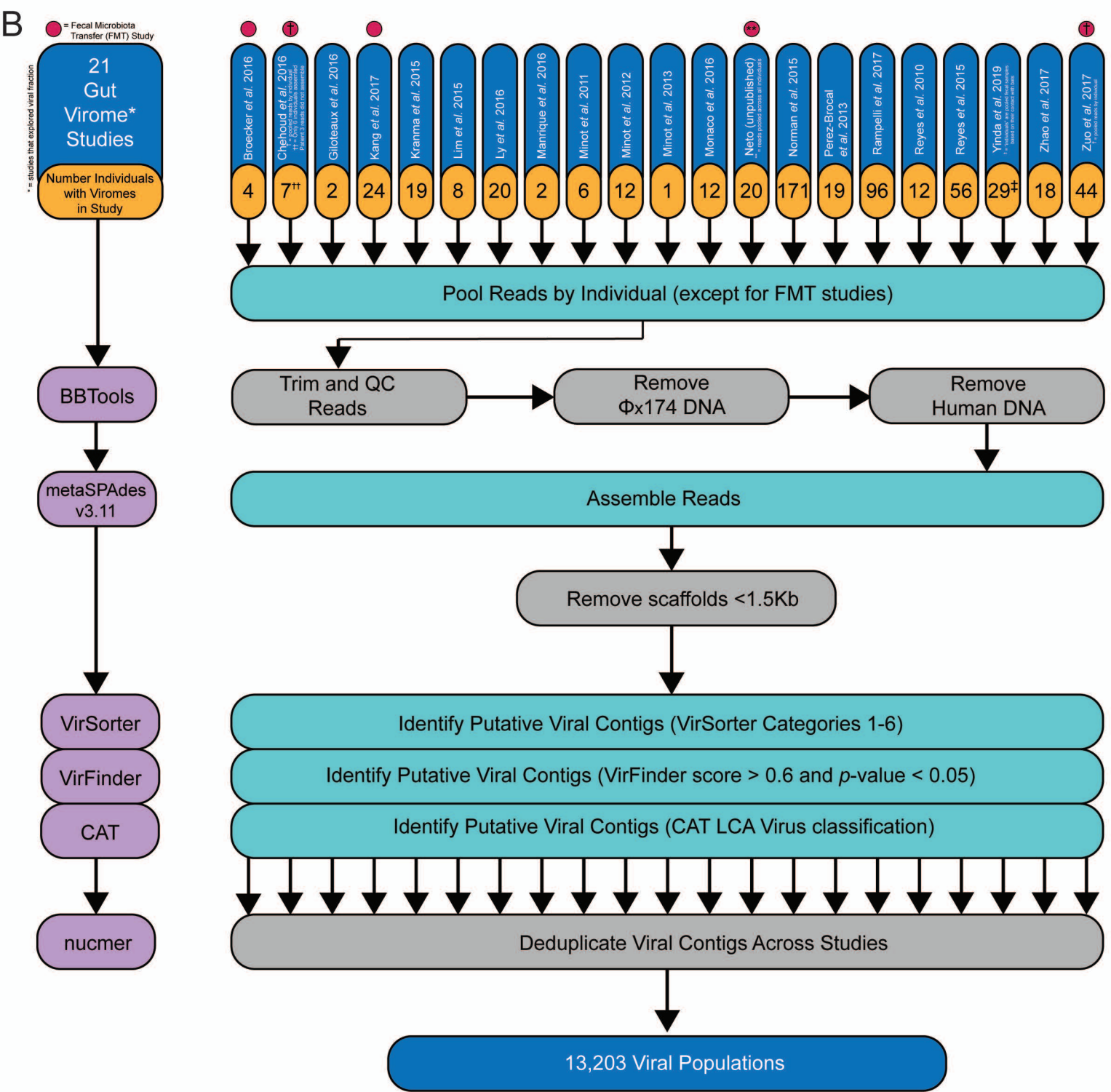
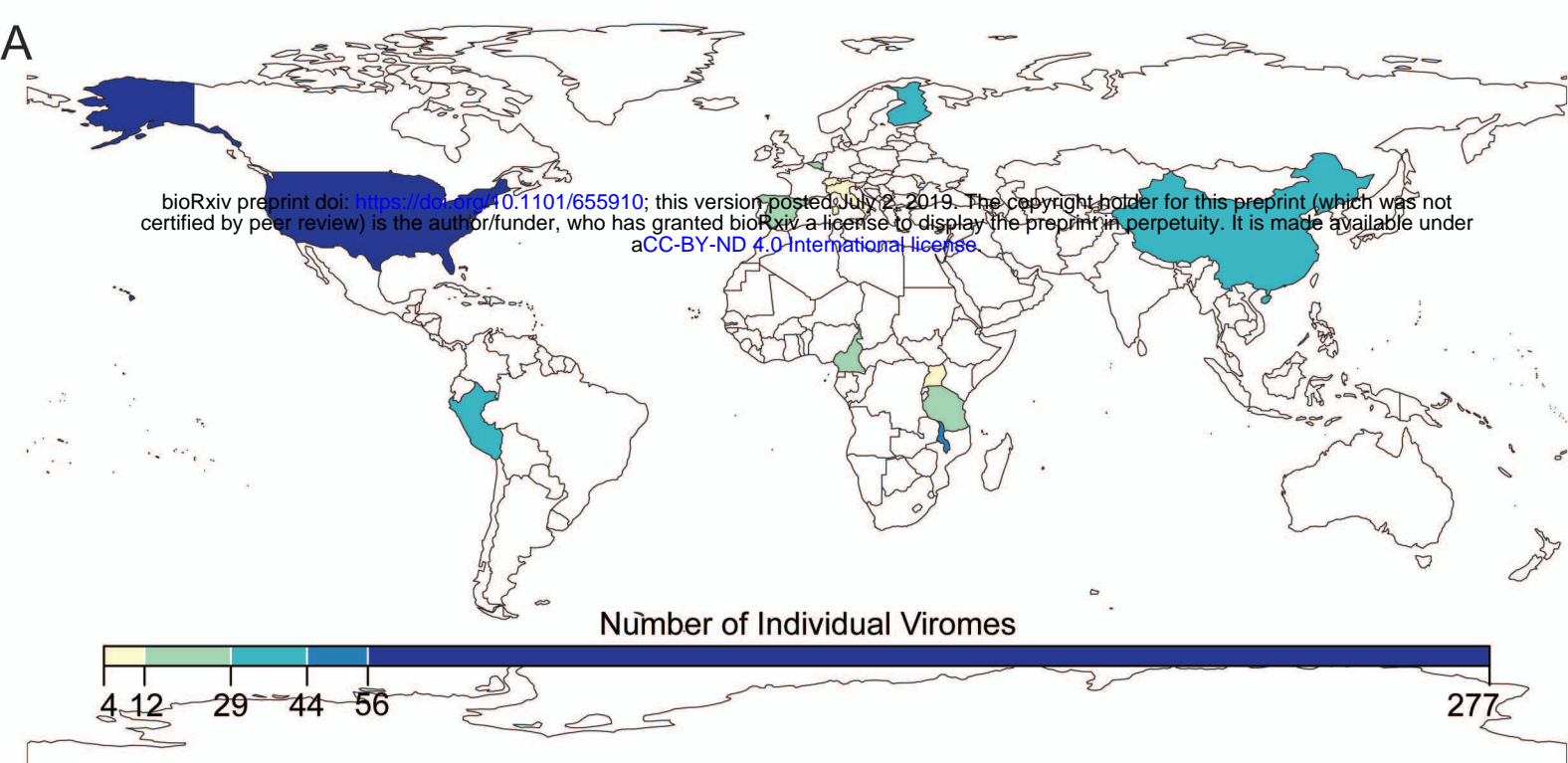
753 **COMPETING INTERESTS.** The authors declare no competing interests.

754

755 **MATERIALS & CORRESPONDENCE.** Correspondence and material requests should be addressed to  
756 Matthew B. Sullivan at [sullivan.948@osu.edu](mailto:sullivan.948@osu.edu).

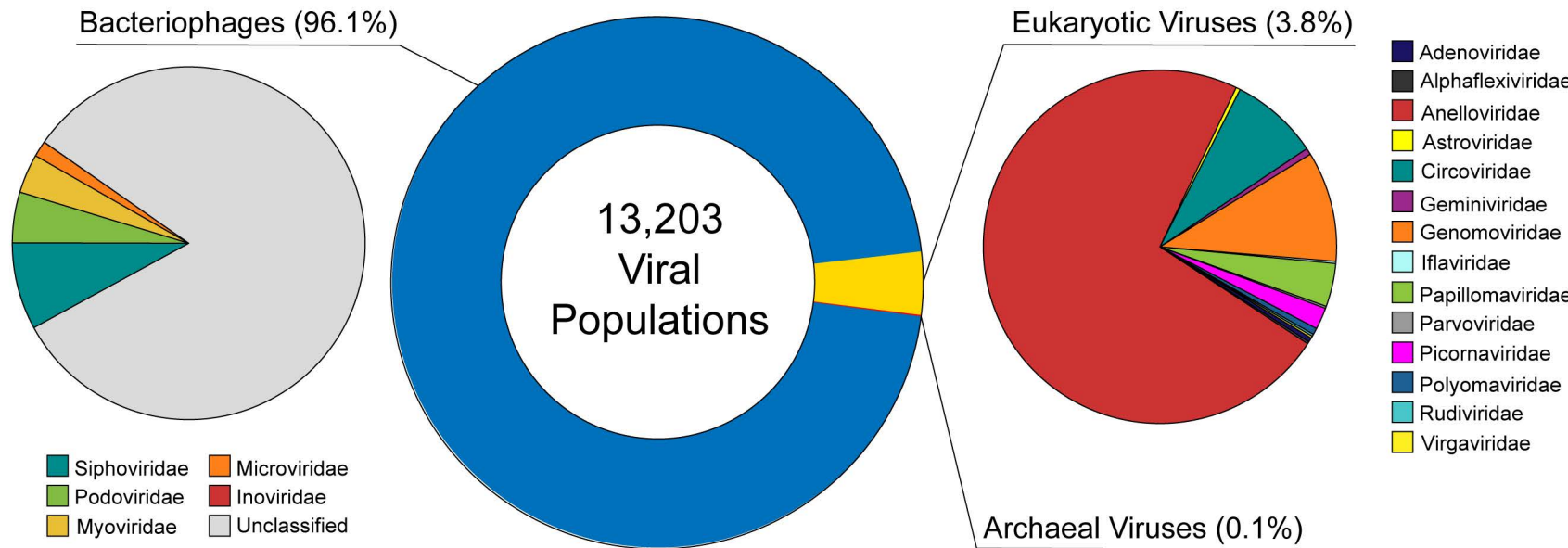
757



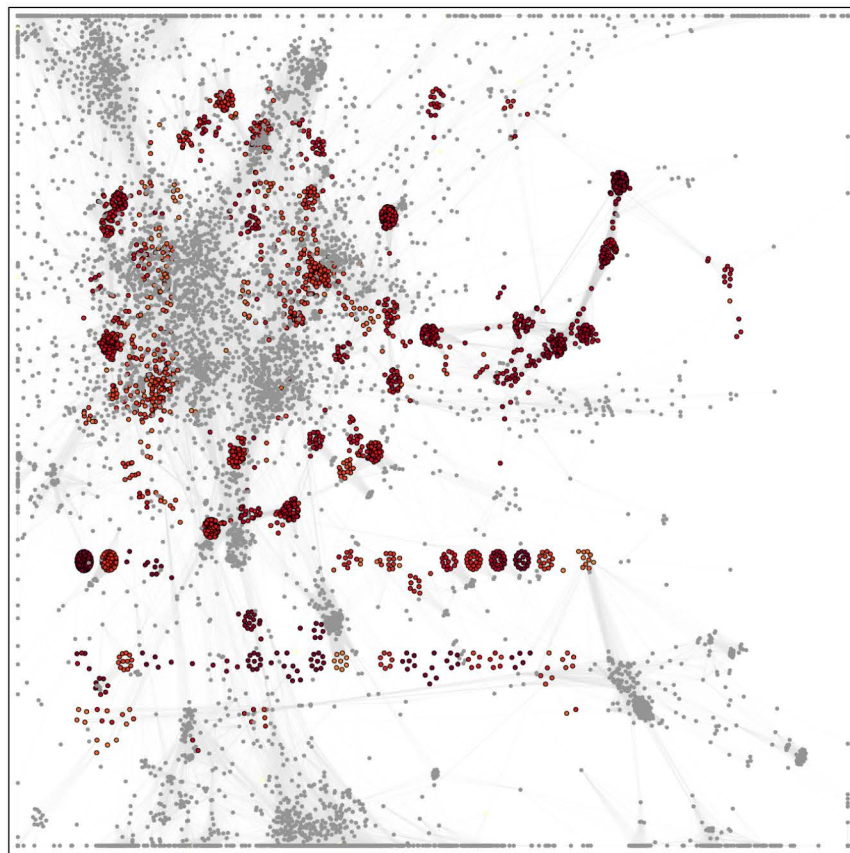


A

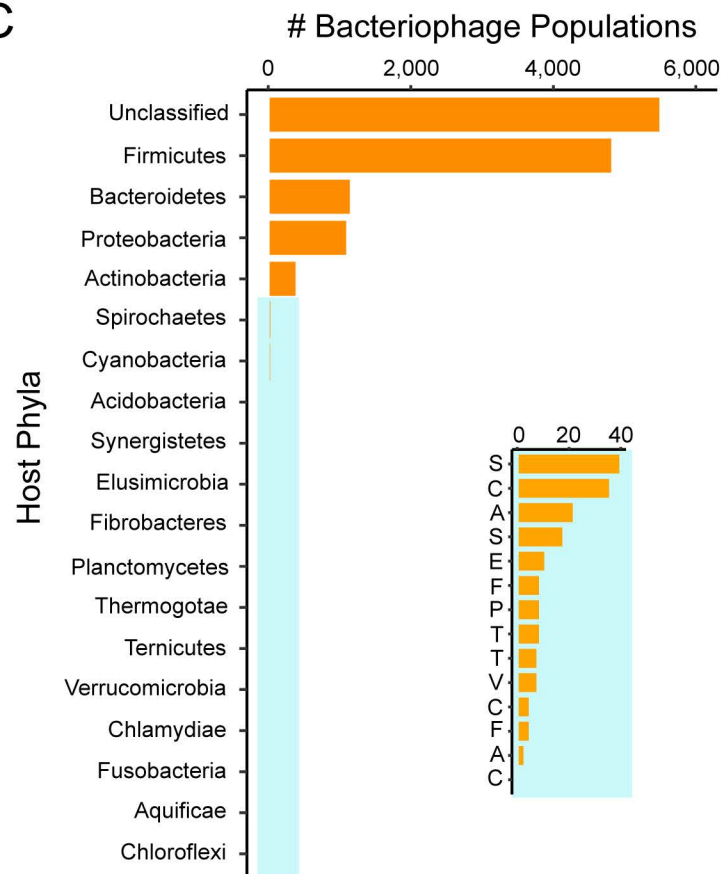
## Gut Viral Database

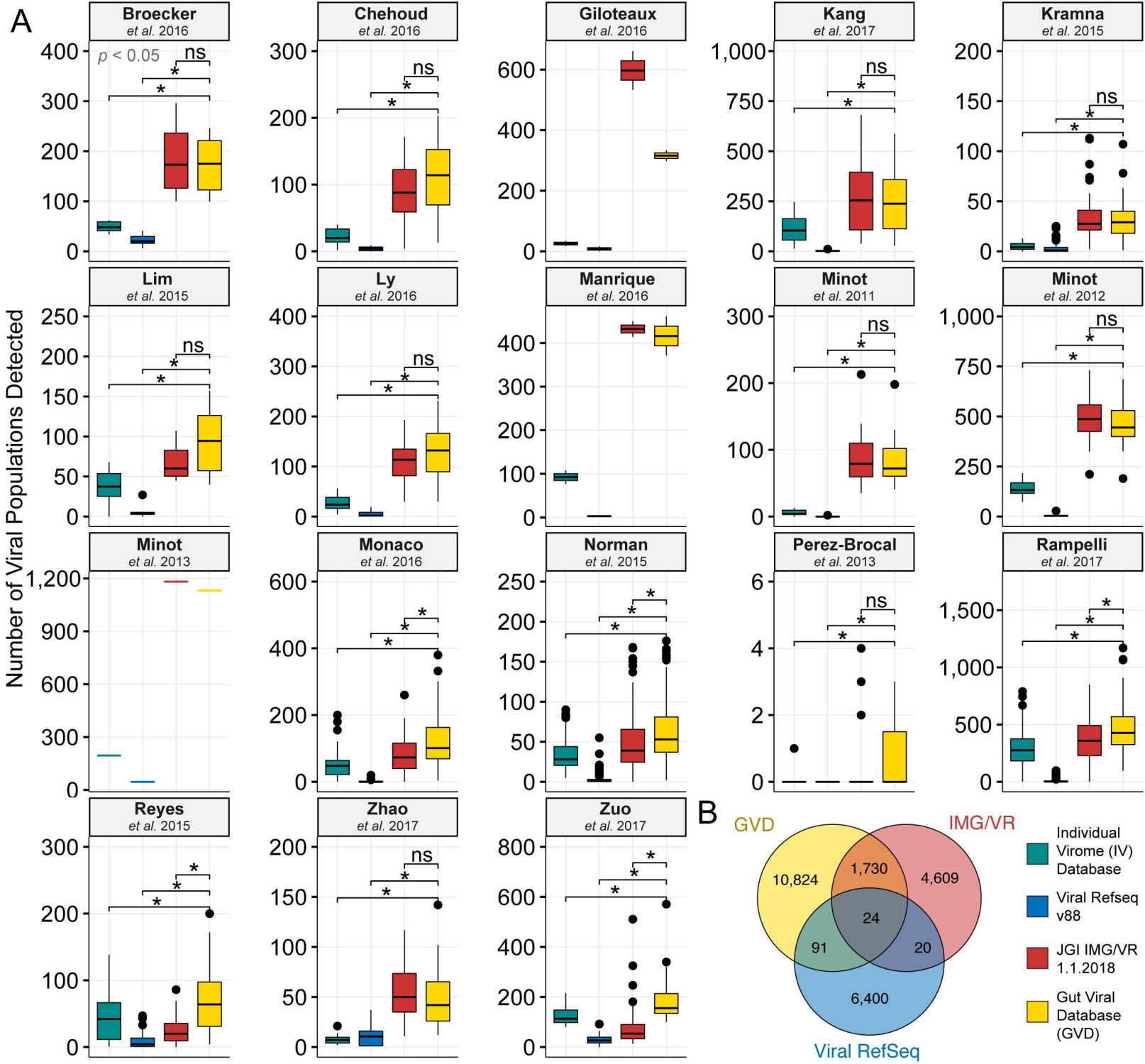


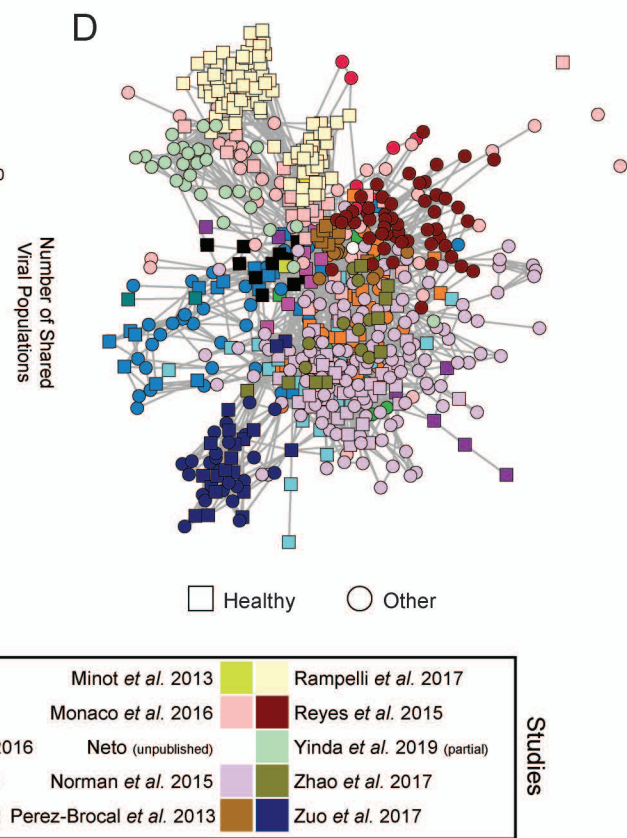
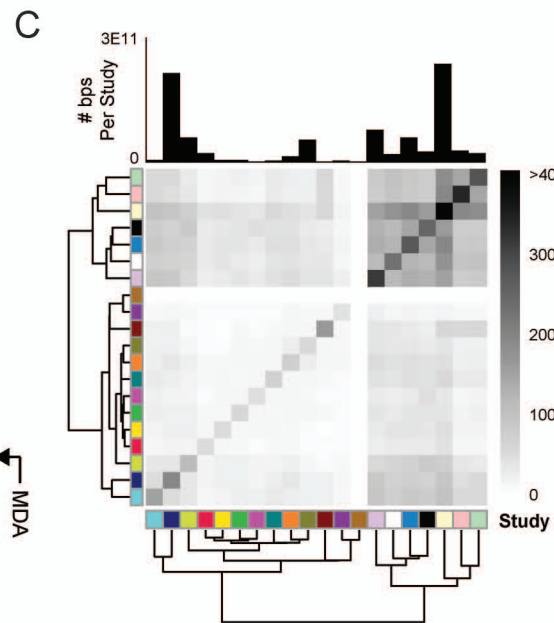
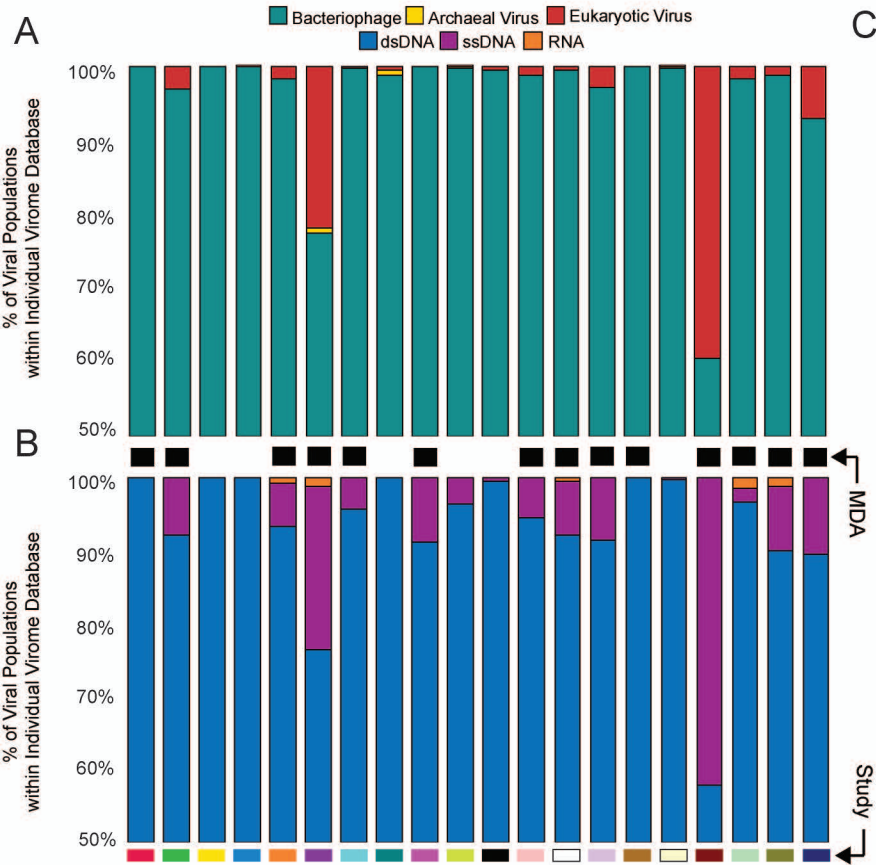
B

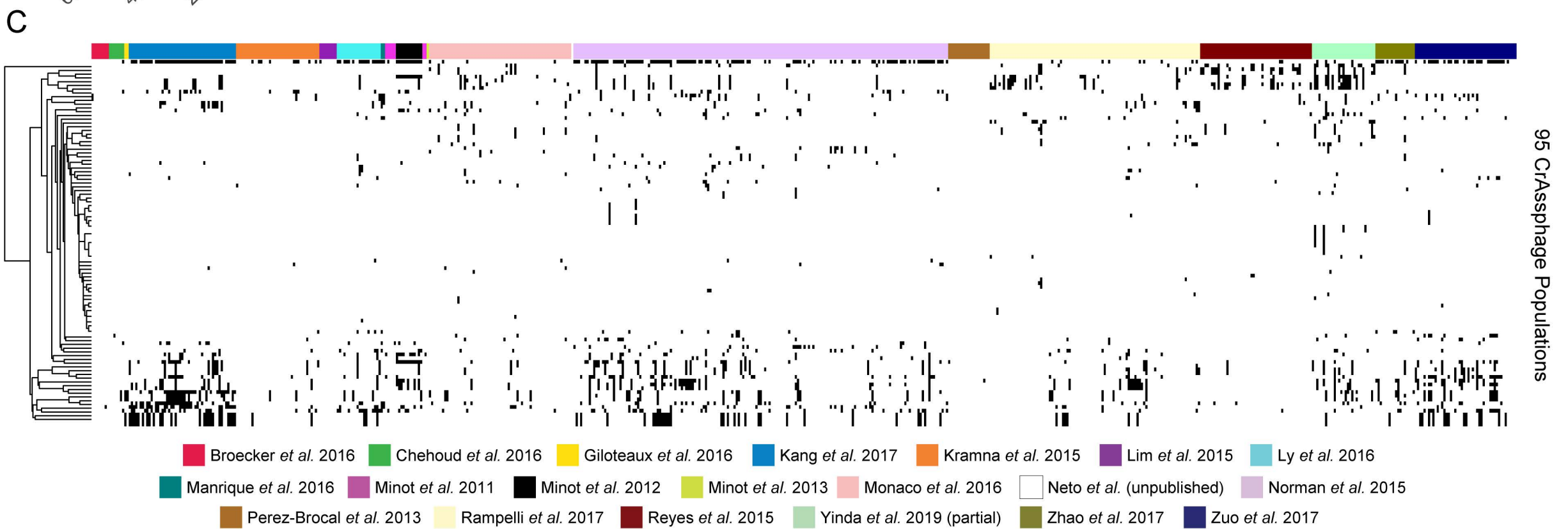
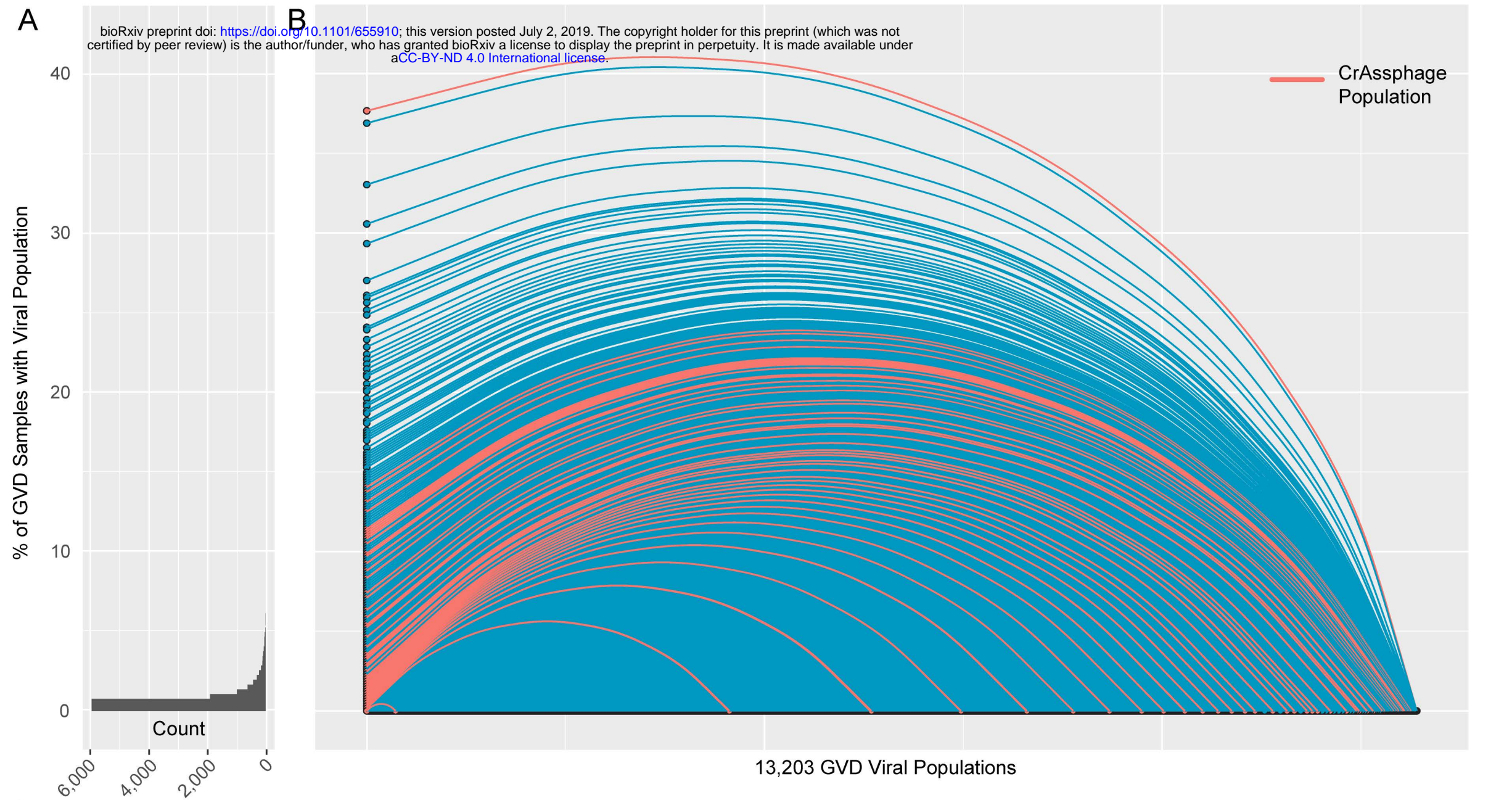


C

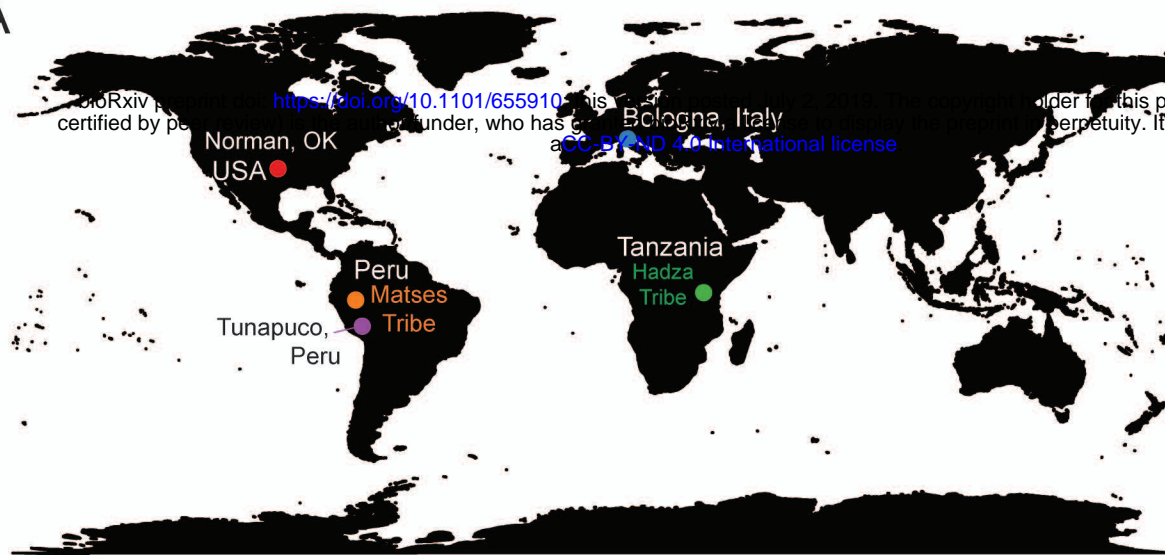




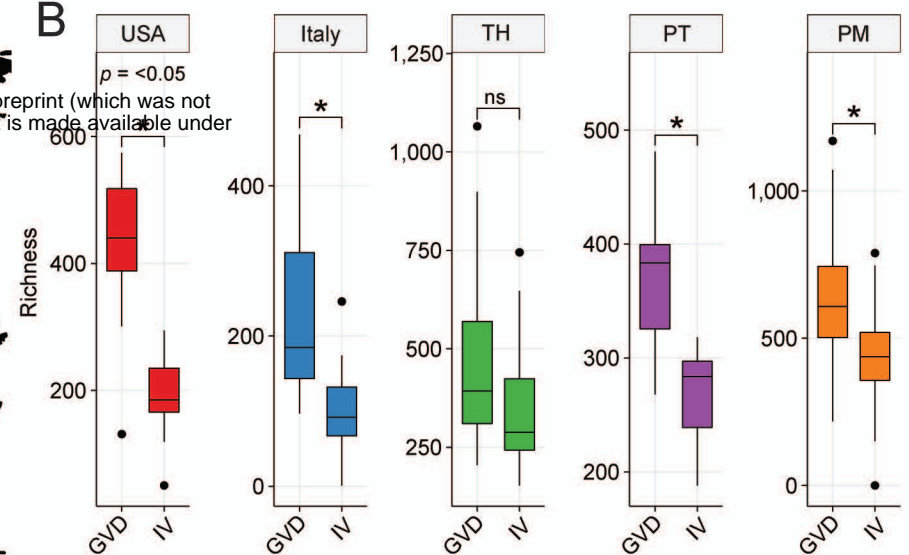




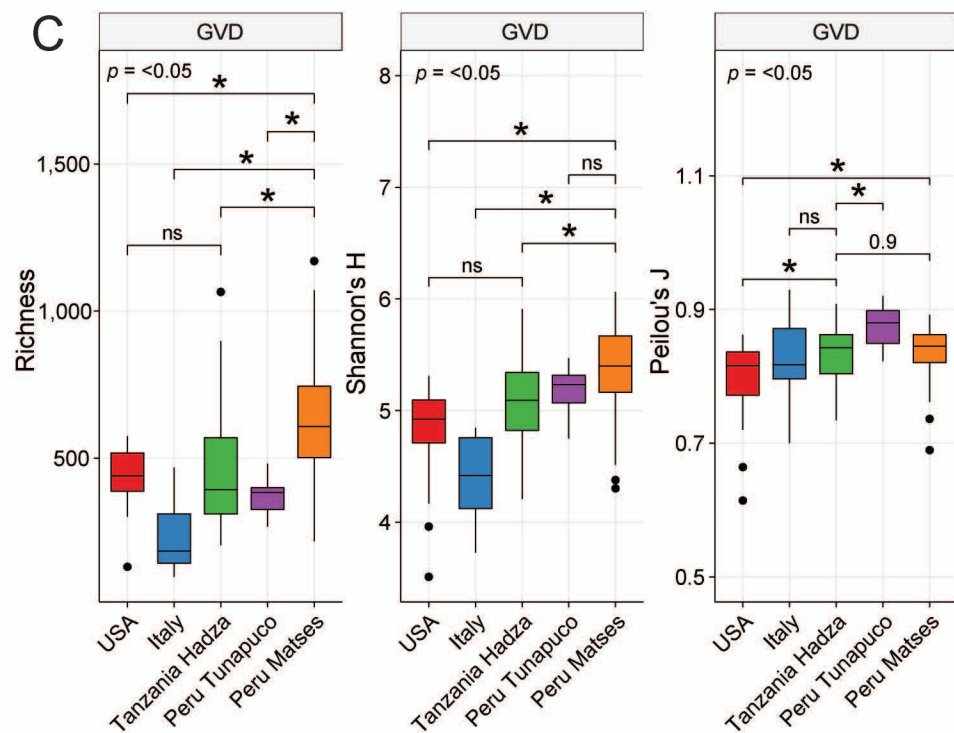
A



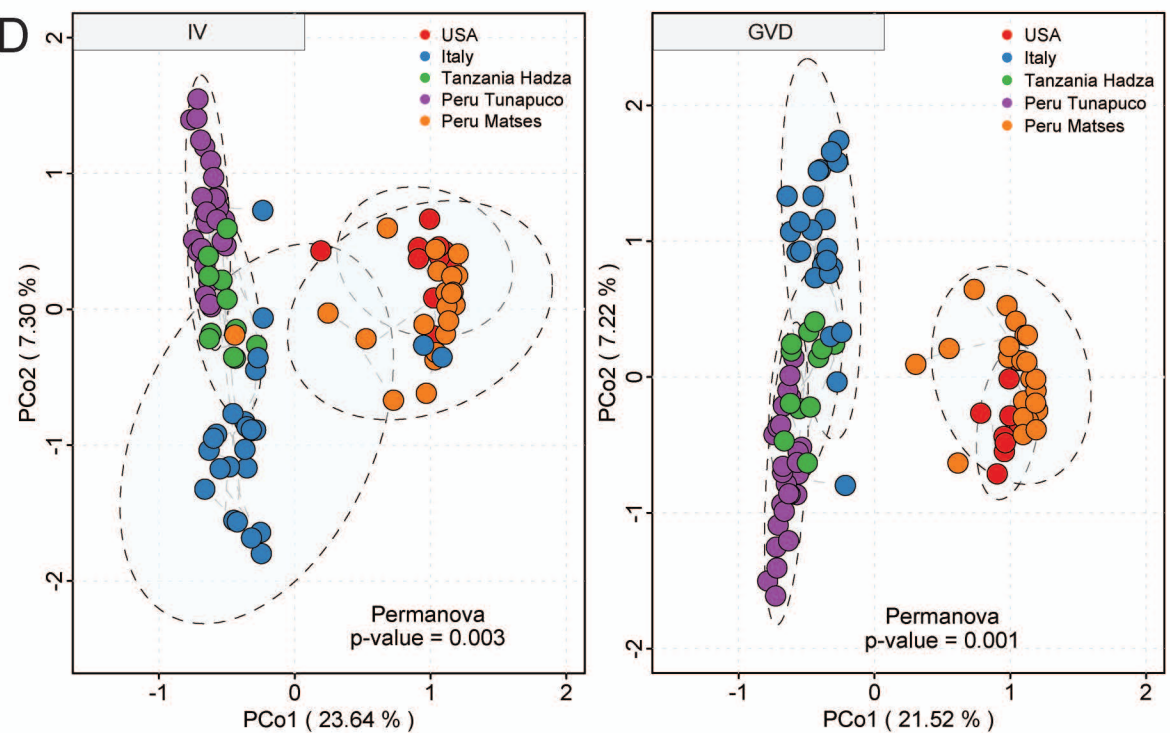
B



C



D



E

