1    **The virome in adult monozygotic twins with concordant or discordant gut**

2    **microbiomes**

3

4    J. Leonardo Moreno-Gallego[1]*, Shao-Pei Chou[2]*, Sara C. Di Rienzi[2], Julia K. Goodrich[1],

5    Timothy Spector[3], Jordana T. Bell[3], Youngblut[1], Ian Hewson[6], Alejandro Reyes[4,5], and Ruth

6    E. Ley[1]¥

7    Addresses:

8    1. Department of Microbiome Science, Max Planck Institute for Developmental Biology,

9    Tübingen 72076, Germany

10   2. Department of Molecular Biology and Genetics, Cornell University, Ithaca NY 14853, USA

11   3. Department of Twin Research and Genetic Epidemiology, King's College London, London

12   SE1 7EH, UK

13   4. Max Planck Tandem Group in Computational Biology, Department of Biological Sciences,

14   Universidad de los Andes, Bogotá 111711, Colombia

15   5. Center for Genome Sciences and Systems Biology, Washington University School of

16   Medicine, Saint Louis, MO 63108, USA

17   6. Department of Microbiology, Cornell University, Ithaca NY 14853 USA

18   * Co-first

19   ¥ Correspondence: rley@tuebingen.mpg.de

20

21   Keywords: Human gut virome; Human gut microbiome; TwinsUK
22

23

24

25

26

27 **SUMMARY**

28 The virome is one of the most variable components of the human gut microbiome.

29 Within twin-pairs, viromes have been shown to be similar for infants but not for

30 adults, indicating that as twins age and their environments and microbiomes diverge,

31 so do their viromes. The degree to which the microbiome drives the virome's vast

32 diversity is unclear. Here, we examined the relationship between microbiome

33 diversity and virome diversity in 21 adult monozygotic twin pairs selected for high or

34 low microbiome concordance. Viromes derived from virus-like particles were unique

35 to each subject, dominated by Caudovirales and Microviridae, and exhibited a small

36 core that included crAssphage. Microbiome-discordant twins had more dissimilar

37 viromes compared to microbiome-concordant twins, and the richer the microbiomes,

38 the richer the viromes. These patterns were driven by the bacteriophages, not

39 eukaryotic viruses. These observations support a strong role of the microbiome in

40 patterning the virome.

41 **INTRODUCTION**

42      The bulk of the human gut microbiome is composed of a vast diversity of

43 bacterial cells, along with a minority of archaeal and eukaryotic cells. The cellular

44 fraction of the microbiome forms a high density microbial ecosystem ($10^{11}$-$10^{12}$ per

45 gram of feces (Sender et al., 2016). All of these cells are accompanied by a virome

46 estimated to be in about equal proportion (ranging between $10^9$ to $10^{12}$ per gram of

47 feces (Castro-Mejía et al., 2015; Hoyles et al., 2014; Ogilvie and Jones, 2017; Reyes

48 et al., 2010). The viral fraction of the human gut microbiome is primarily composed of

49 bacteriophages and prophages, and it also includes rarer eukaryotic viruses and

50 endogenous retroviruses (Breitbart et al., 2003; Minot et al., 2011; Reyes et al.,

51 2010). Currently, the majority of phages have no matches in databases and their

52 hosts remain to be elucidated. Matching phages to their hosts is challenging: for

53 instance, the host of the most common human gut phage, crAssphage, has only

54 recently been identified as *Bacteroides spp.* (Shkoporov et al., 2018; Yutin et al.,

55 2018). In addition to the identification of hosts, other questions remain as to the

56 factors most important in shaping the virome, and how predictive the cellular fraction

57 of the microbiome can be of the virome.

58      The temporal population dynamics of phages and their hosts might be

59 expected to be linked. Indeed, population oscillations of viruses and their bacterial

60 hosts are described for aquatic systems, where they indicate that viruses play a key

61 role in regulating bacterial populations (Suttle, 2007; Thingstad, 2000; Thingstad et

62 al., 2014; Weitz and Dushoff, 2008). But such patterns of predator/prey dynamics are

63 not typical for the human gut virome and microbiome (for clarity, from here on we

64 use 'microbiome' to refer to cellular fraction of the microbiome, e.g., mostly bacterial

65    cells) (Minot et al., 2011; Reyes et al., 2013; Rodriguez-Brito et al., 2010; Rodriguez-

66    Valera et al., 2009). Nonetheless, the virome and microbiome do display some

67    common patterns of diversity across hosts, such as high levels of interpersonal

68    differences and relative stability over time (Reyes et al., 2010). The microbiome

69    tends to be more similar for related individuals compared to unrelated individuals,

70    possibly due to shared dietary habits, which drive similarity between microbiomes

71    (Cotillard et al., 2013; David et al., 2014). In accord, diet has been associated with

72    virome diversity, quite possibly through diet effects on the microbiome (Minot et al.,

73    2011). In infants, twin comparisons have revealed viromes to be more similar

74    between co-twins than between unrelated individuals (Lim et al., 2015; Reyes et al.,

75    2015). This pattern was not observed in adult twins (Reyes et al., 2010) possibly due

76    to divergence of their microbiomes (Reyes et al., 2010). The degree to which the

77    microbiome itself drives patterns of virome diversity across hosts has been difficult to

78    assess due to confounding factors such as host relatedness.

79         Here, we focus on adult monozygotic (MZ) twin microbiomes to explore

80    further the relationship between microbiome and virome diversity. By studying the

81    viromes of MZ twin pairs, we control for host genetic relatedness. Although MZ twin

82    pairs generally have more similar microbiomes compared to dizygotic (DZ) twin pairs

83    or unrelated individuals, MZ twins nevertheless can display a large range of within-

84    twin-pair microbiome diversity (Goodrich et al., 2014). We previously generated fecal

85    microbiome data for twin pairs from the TwinsUK cohort (Goodrich et al., 2014), and

86    based on this information we selected twin pairs either highly concordant or highly

87    discordant for their microbiomes. We generated viromes from virus-like particles

88    (VLPs) obtained from the same samples from which the microbiomes were derived.

89  Results indicate that microbiome diversity and virome diversity measures are

90  positively associated.

91

92  **RESULTS**

93  **Selection of microbiome-concordant and discordant monozygotic twin**

94  **pairs -** We selected twin pairs with a similar body mass index (BMI), whose

95  microbiomes were either concordant or discordant for microbiome between-sample

96  diversity (β-diversity) based on previously obtained 16S rRNA gene data. The adult

97  co-twins in this study did not share a household and we assume that other

98  environmental variability was similar across twin pairs. We determined the degree of

99  concordance or discordance between co-twins' microbiomes based on three β-

100  diversity distance metrics: Bray-Curtis, weighted UniFrac and unweighted UniFrac

101  **(See Methods)**. As expected, the β-diversity measures were correlated (Pearson

102  pairwise correlation coefficient > 0.4). Based on the distribution of pairwise distance

103  measures, we selected 21 MZ twin pairs from the boundaries of all three distributions

104  **(Figure 1A),** while maintaining a balanced distribution of age and BMI across the set

105  **(Table S1).** Within the 21 selected twin pairs, the microbiomes of microbiome-

106  concordant co-twins were, as expected, more similar to each other than microbiomes

107  of microbiome-discordant co-twins (p = $6.31 \times 10^{-12}$). The microbiomes of the

108  discordant co-twins differed compositionally at all taxonomic levels, particularly at the

109  phylum level, with Firmicutes and Bacteroidetes, the two dominant phyla,

110  contributing the most to the variation between co-twins **(Figure 1B and 1C)**.

111  **Shotgun metagenomes of VLPs -** We isolated virus-like particles (VLPs)

112  from the same fecal samples that had been used for 16S rRNA gene diversity

5

113    profiling **(See Methods)**. DNA extracted from VLPs was used in whole genome

114    amplification followed by shotgun metagenome sequencing **(See Methods)**. A first

115    library ("large-insert-size library") was selected with an average insert size of 500 bp

116    (34,325,116 paired reads in total; 817,265 ± 249,550 paired reads per sample after

117    quality control) and used for *de novo* assembly of viral contigs. Smaller fragments

118    with an average insert size of 300bp were purified in a second library ("small-insert-

119    size library") and sequenced. The resulting pair-end reads were merged into

120    25,324,163 quality filtered longer reads to increase mapping accuracy (602,956 ±

121    595,444 merged reads per sample) **(See Methods) (Table S2)**.

122        **Identification of putative bacterial contaminants -** Viromes prepared and

123    sequenced from VLPs may be contaminated with bacterial DNA (Roux et al., 2013).

124    However, given that phages are major agents of horizontal gene transfer and that

125    temperate viruses often comprise up to 10% of bacterial genomes in a prophage

126    state, removal of potential bacterial contamination risks also removing viral reads. To

127    assess bacterial DNA contamination, we mapped virome reads against a set of

128    8,163 fully assembled bacterial genomes. Our strategy consisted of evaluating the

129    coverage along the length of each genome (in bins of 100Kb), and those genomes

130    with a median coverage greater than 100 were considered contaminants. Reads

131    mapping to short regions were considered to be prophages or horizontally

132    transferred genes and retained **(See Methods) (Figure 2A)**. Reads mapping to

133    genomes determined to be potential contaminants were removed from further

134    analyses.

135        We identified 65 bacterial genomes as contributing to potential contaminant,

136    with 1.006 ± 1.125% (average ± std) reads per sample mapping to those bacterial

137    genomes **(Table S2)**. The majority (37/68) belonged to the Firmicutes phylum; at the

138    species level, *Bacteroides dorei*, *B. vulgatus*, *Ruminococcus bromii*,

139    *Faecalibacterium prausnitzii*, *B. xylanisolvens, Odoribacter splanchnicus and B.*

140    *caecimuris* (in that order) were detectable in at least 50% of the samples **(Table S2)**.

141    If the most abundant bacterial species in the microbiome are the most likely sources

142    of contamination, then the taxonomic composition of the bacterial contaminants

143    should correlate with their corresponding bacterial abundances in the microbiome.

144    However, we observed no significant correlation between the relative abundances of

145    taxa represented in the contaminant DNA and in the microbiomes **(Figure 2B)**.


146       **Functional profiles support viral enrichment in VLP purifications -** To

147    assess the functional content of the viromes, we annotated the "short-insert-size

148    library" raw reads using the KEGG annotation of the Integrated Gene Catalog (IGC)

149    (Li et al., 2014) **(See Methods)**. In line with previous reports (Breitbart et al., 2008;

150    Minot et al., 2011; Reyes et al., 2010), the majority of reads (85.43 ± 5.74%) from

151    our VLP metagenomes mapped to genes with unknown function **(Figure 3A)**.


152       To further verify that sequences were derived from VLPs and not microbiomes

153    generally, we conducted an internal check in which we generated and compared

154    additional metagenomes from VLPs and bulk fecal DNA for an additional 4

155    individuals (2 twin pairs**; Figure 1A)**. As expected, the functional profiles of viromes

156    and microbiome-metagenomes derived from the same samples were dissimilar.

157    Virome reads that mapped to annotated genes were enriched in two categories:

158    Genetic Information Process (48.87 ± 12.12%) and Nucleotide Metabolism (17.59 ±

159    8.81%), compared to 24.31 ± 1.28% and 5.47 ± 0.4% for the microbiome-

160    metagenome, respectively **(Figure 3B)**. Most of the other functional categories

161     present in the bacterial metagenomes were essentially absent from the viromes.

162     Furthermore, the functional annotations of the viromes show greater between-

163     sample variability than the microbiomes and a lower intraclass correlation coefficient

164     (**Figure 3B**).

165          **Viromes are unique to individuals** - We assembled reads from the "large-

166     insert-size library" resulting in a total of 107,307 contigs ≥ 500 nt (max: 79,863 nt;

167     mean 1,186nt ± 1,741; **Figure S1**). To assess the structure and composition of the

168     viromes, a matrix of the recruitment of reads against dereplicated contigs were built

169     **(See Methods)**. The recruitment matrix included 14,584 contigs that were both long

170     (> 1,300 nt) and well covered (> 5X); these are referred to as 'virotypes' **(Figure S1)**.

171     Analysis of the recruitment matrix showed that each individual harbored a unique set

172     of virotypes: 3,415 virotypes (23.41% of total) were present in only one individual;

173     413 virotypes (2.83%) were present in at least 50% of the individuals; only 18

174     virotypes (0.1%) were present in all individuals.

175          **Twins with concordant microbiomes share virotypes -** We checked for

176     virotypes shared between twins and observed that co-twins did not share more

177     virotypes than unrelated individuals (p = 0.074). We then assessed microbiome-

178     concordant and discordant twin pairs separately: twins with a discordant microbiome

179     did not share more virotypes that unrelated individuals (p = 0.254), and twins with a

180     concordant microbiome did share more virotypes than unrelated individuals (p =

181     0.048). Furthermore, we also found that twins with a concordant microbiome shared

182     more virotypes than twins with a discordant microbiome (p = 0.015; **Figure S2)**.

8

183    **Bacteriophage dominance of the gut virome -** In order to characterize the

184    taxonomic composition of the virome, we attempt to annotated all 66,446

185    dereplicated and well covered contigs **(Figure S1)** using a voting system approach

186    that exploited the information in both the assembled contigs and their encoding

187    proteins **(See Methods)**. In addition, we performed a custom annotation on two

188    highly abundant gut-associated bacteriophage families: (i) the crAssphage (Dutilh et

189    al. 2014; Yuting et al. 2018) and (ii) the *Microviridae* families (Székely and Breitbart

190    2016). For this, we used profile Hidden Markov Models (HMMs) to search for

191    crAssphage (dsDNA viruses) and *Microviridae* (ssDNA viruses) contigs **(See**

192    **Methods)**.

193        Using HMMs allowed us to identify distant homologs, which we then

194    incorporated into a phylogenetic tree with known reference sequences to confirm the

195    annotation and better resolve the taxonomy. We annotated 108 contigs (19

196    crAssphage, 90 *Microviridae*), validated the family assignment of 68 contigs, and

197    assigned a subfamily to 97 contigs without previous subfamily assignment. For the

198    *Microviridae*, only 11 contigs had a previous taxonomic assignment, all belonging to

199    the *Gokushovirinae*: we confirmed these and 23 more as *Gokushovirinae*, 54 as

200    Alpavirinae and 1 contig as *Pichovirinae* **(Figure S3A)**. For the crAssphage, 11

201    contigs were clustered with the original crAssphage, 3 contigs grouped with the

202    reference Chlamydia phage, and 5 contig grouped with the reference IAS virus

203    **(Figure S3B)**.

204        After collating the voting system annotation and the HMM annotation, a total

205    of 12,751 contigs (29,62%) were taxonomically assigned **(Figure S1)**. Viromes were

206    dominated by bacteriophages with only 6.42% of contigs annotated as Eukaryotic

9

207    viruses. As expected, most of the contigs (96.98%) were dsDNA viruses, while only

208    2.43% of contigs were annotated as ssDNA viruses. Caudovirales was the most

209    abundant Order, with its three main families represented: Myoviridae (20.22 ±

210    4.83%), *Podoviridae* (10.54 ± 3.27%), and *Siphoviridae* (35.25 ± 7.19%). The

211    crAssphage family constituted on average 13.26% (± 12.24%) of the contigs,

212    reaching a maximum contribution of 55.80% in one virome, and *Microviridae*

213    represented 3.87 ± 2.57% of the viromes. Interestingly, we observed that

214    *Phycodnaviridae* exceeded 1% of average abundance (1.77 ± 1.12%; **Figure 4A)**

215    and that contigs related to any nucleocytoplasmic large DNA viruses (NCLDV) had a

216    mean relative contribution of 3.99 ± 2.22%. The 18 contigs present in all samples

217    included 10 annotated as crAssphage, 2 annotated as "unclassified Myoviridae", 2

218    "unclassified Caudovirales", 1 classified as *Microviridae,* and 3 unclassified. Within a

219    defined taxonomic profile for each sample, we looked for differences in composition

220    between viromes at all taxonomic levels for concordant and discordant twin-pairs.

221    There were no significant differences between groups for any taxa at the Order and

222    Family levels, including crAssphage and *Microviridae* families **(Figure 4B)**.

223        We used CRISPR spacer mapping and the microbe-versus-phage (MVP)

224    database (Gao et al., 2018) to predict hosts for virotypes and taxonomically

225    characterized contigs **(See Methods)**. As host annotation was directed to

226    bacteriophages, we did not gain any information for contigs annotated as Eukaryotic

227    viruses. These approaches allowed us to identify putative hosts for 910 contigs.

228    Within these 910 contigs, only one was previously annotated as crAssphage, and ss

229    expected, its host was inferred to be a member of *Bacteroidetes.* In total we

230    identified 1,280 bacterial putative host strains, including 187 species from 87 genera

231    over several phyla; most of them from Firmicutes (92), followed by Bacteroidetes

232     (41) and Proteobacteria (38). The median number of host for each contig was 1

233     (IQR=1-2) while the median number of phages per host, at the strain level, was 2

234     (IQR=1-3) (Figure S4).

235     **Virome diversity correlates with microbiome diversity -** To assess the

236     relationship between virome and microbiome diversity, we examined the within-

237     samples diversity (α-diversity) and β-diversity of the viromes using three different

238     layers of information that we recovered from the sequence data: i) virotypes, iii)

239     taxonomically annotated contigs, and iii) annotated genes from short reads **(Figure**

240     **S1)**.

241     *Alpha-diversity* - α-diversities of the microbiome and the virome were

242     positively correlated in two of the three layers of information used to test the

243     correlation (virotypes and taxonomy annotated contigs but not genes; **Figure 5A)**.

244     We used annotated contigs to ask about the α-diversity within subgroups of viruses:

245     (ssDNA eukaryotic, dsDNA eukaryotic, ssDNA bacteria and dsDNA bacteria). Our

246     results show that the diversity of eukaryotic viruses does not correlate with the

247     microbiome α-diversity. In contrast, bacteriophages and microbiome α-diversity were

248     positively correlated, for both ssDNA or dsDNA bacterial viruses **(Figure 5B)**.

249     *Beta-diversity -* We observed that concordant twins had lower virome β-

250     diversity compared to discordant twins using Hellinger distances **(Figure 6)**; the

251     mean binary Jaccard distance and Bray-Curtis dissimilarity of viromes also showed

252     the same trend **(Figure S5A and S5B)**. Similar to what we observed with α-diversity,

253     regardless of the layer of information used, the mean Hellinger distance of viromes

254     within MZ twin pairs with concordant microbiomes was significantly lower than that of

255     MZ twin pairs with discordant microbiomes ($p < 0.04$, Mann-Whitney's U test)

256 **(Figure 6)**. Furthermore, a similar significant positive correlation was observed

257 between microbiome and virome β-diversity when using the annotated contigs. This

258 relationship was driven by the bacteriophages ($p = 0.009$, Mann-Whitney's U test),

259 but not the eukaryotic viruses ($p = 0.243$, Mann-Whitney's U test).

260 Finally, we compared the virome and microbiome pairwise distances among

261 related (co-twins) and unrelated individuals. The pairwise distance matrices showed

262 a positive correlation between virome and microbiome β-diversity measures not only

263 within twin pairs (Pearson correlation coefficient > 0.50) but also generally across all

264 individuals (Pearson correlation coefficient > 0.25; $p < 0.003$, Mantel test; **Figure**

265 **S5C)**. These results show that regardless of genetic relatedness between hosts,

266 individuals with more similar microbiomes harbour more similar viromes.

267

268 **DISCUSSION**

269 Co-twins, like other siblings, generally have more similar gut microbiomes

270 within their twinships compared to unrelated individuals (Lee et al., 2011; Palmer et

271 al., 2007; Tims et al., 2013; Turnbaugh et al., 2009; Yatsunenko et al., 2012).

272 Moreover, MZ twins have overall more similar microbiomes than DZ twins, although

273 at a whole-microbiome level this effect is small and primarily driven by a small set of

274 heritable microbiota (Goodrich et al., 2014, 2016). Within a population of MZ twin

275 pairs, however, the range of within-twin pair differences in the microbiomes can be

276 as great as for DZ twins (Goodrich et al., 2014). We took advantage of the large

277 spread in β-diversity for MZ co-twins to select co-twins that were either highly

278 concordant or discordant for their gut microbiomes. Our analysis of their viromes

279 showed that despite the high variation in the gut viromes between individuals, and

12

280    regardless of host relatedness, the more dissimilar their microbiomes, the more

281    dissimilar their viromes. This pattern was driven by the bacteriophage component of

282    the virome.

283         Here, by choosing MZ twins from a distribution of divergence in the

284    microbiome, we removed host genetic relatedness as a variable. Previous studies of

285    the viromes and microbiomes of infant twin pairs showed that the microbiomes and

286    viromes of co-twins were more similar than those of unrelated individuals, suggested

287    shared host genotype and/or environment were key (Lim et al., 2015; Reyes et al.,

288    2015). In contrast, an early study of the virome of adult twins showed that adult co-

289    twins did not have more similar viromes than unrelated individuals (Reyes et al.,

290    2010); however, in light of the current study's results, this was likely a power issue.

291    Indeed, in our dataset we observed that regardless of whether twins were

292    concordant or discordant for their microbiomes, co-twins had more similar viromes

293    (virotypes and taxonomy) than unrelated individuals.

294         The previously reported greater virome similarity in young compared to adult

295    twins has been related to the fact that infants have a greater shared environment

296    compared to adult twins (Lim et al., 2015), particularly in terms of their diet. Minot et

297    al., have also shown that individuals on the same diet have more similar gut viromes

298    than individuals on dissimilar diets (Minot et al., 2011). It is well established that diet

299    is a strong driver of daily microbiome fluctuation (Claesson et al., 2012; David et al.,

300    2014; De Filippo et al., 2010; Wu et al., 2011), so the effect of diet on the virome is

301    likely mediated by the microbiome. However, we did not control for diet, so it is

302    possible that the microbiome discordance that we observe was caused by co-twins

303    eating differently around the time of sampling. Regardless of what underlies the

13

304    variance in microbiome concordance, it is strongly associated with virome

305    concordance.

306         The relationship between virome richness and microbiome richness had not

307    previously been directly addressed in adults. We observed that the α-diversity of the

308    microbiome and the virome were positively correlated using two of the three layers of

309    information describing virome diversity. Specifically, this pattern was observed for

310    virotypes and taxonomy but not for genes. However, since virome genes were

311    observed to be enriched in only two categories, Genetic Information Processing and

312    Nucleotide Metabolism, we would not expect differences in diversity of virome genes

313    between subjects. The taxonomic annotation layer showed that the bacteriophage

314    component of the virome, not the eukaryotic viruses, was driving this α-diversity

315    correlation pattern.

316         The positive relationship between virome and microbiome α-diversity

317    suggests that a greater availability of hosts drives a greater availability of viruses.

318    These observations are in accordance with "(Minot et al., 2013; Reyes et al., 2010),

319    which posits that in a (Minot et al., 2013; Reyes et al., 2010) (Knowles et al., 2016).

320    Indeed, longitudinal studies of the human gut virome have reported genes

321    associated with lysogeny, low mutation rate over time in temperate-like contigs, and

322    long-term stability of the virome, suggesting preference for a lysogenic cycle (Minot

323    et al., 2013; Reyes et al., 2010). Nevertheless, phage predation has been

324    acknowledged as an important factor for the maintenance of highly diverse and

325    efficient ecosystems (Rodriguez-Valera et al., 2009) and may play a role in the

326    maintenance of diversity in a rapidly changing ecosystem as the human gut (David et

327    al., 2014). Short scale time-series analyses of virome-microbiome interactions, along

328  with a better understanding of the lysogenic-lytic switch in viral reproduction, would

329  help to interpret the observed patterns in the human gut virome.

330      The composition of the viromes described here was similar to what has been

331  previously reported for adult fecal viromes (Minot et al., 2011, 2013; Reyes et al.,

332  2010) but stands in contrast to what has been observed in babies (Lim et al., 2015).

333  From the annotated fraction of the virome, the order *Caudovirales* and its families

334  *Siphoviridae*, *Myoviridae*, and *Podoviridae,* along with crAssphage, were the

335  dominant phages in all samples. Manrique *et al.* have summarized the phage

336  colonization of the infant gut as follows: the eukaryotic viruses first dominate the

337  newborn gut, followed by the *Caudovirales*, and by 2.5 years of age the *Microviridae*

338  start to dominate (Manrique et al., 2017). We did observe abundant *Microviridae* in

339  our sample set, but the Caudovirales were the dominant group. Age was not related

340  to patterns of diversity in the set of adult subjects studied here.

341      Despite the high diversity and uniqueness of each virome described here, we

342  nonetheless recovered a core virome among the subjects: 18 contigs were present

343  in all samples. More than half of these contigs were annotated as crAssphage,

344  consistent with recent reports that this phage is widespread (Dutilh et al., 2014;

345  Manrique et al., 2016; Yarygin et al., 2017). Other shared virotypes in our dataset

346  were classified as *Myoviridae* and *Microviridae.* We also recovered contigs mapping

347  to representative families of the nucleocytoplasmic large DNA viruses (NCLDV),

348  *Phycodnaviridae* and *Mimiviridae*. These types of viruses are increasingly reported

349  as members of the human gut virome (Colson et al., 2013; Halary et al., 2016). A

350  core set of bacteriophages consisting of nine representatives, including crAssphage,

351  has previously been reported for the human gut (Manrique et al., 2016). Widely

15

352     shared virotypes may indicate the wide sharing of specific hosts between individuals,

353     or that these viruses have a broad host range within the human microbiome.

354        Our use of the HMMs to annotate viral contigs allowed a deep exploration into

355     the taxonomic content of the virome. We annotated a diversity of contigs beyond

356     what was revealed from comparisons to public databases, and also confirmed those

357     annotations. Because each type of virus (*e.g.*, family) requires its own HMM, we

358     applied this method to a few key groups. When applied to the crAssphage, the HMM

359     retrieved contigs that grouped only with sequences derived from fecal viromes and

360     not with sequences from other environments (e.g., terrestrial or marine). This

361     suggests that although crAssphage is a diverse group of bacteriophages, its diversity

362     in the human gut is restricted to sequences related to the reference crAssphage

363     genome (Dutilh et al., 2014), the IAS virus reference (Shkoporov et al., 2018), or

364     *Chlamydia* bacteriophage (Yutin et al., 2018). We also applied HHM to the family

365     *Microviridae*, which are single strand DNA bacteriophages. We were able to confirm

366     the presence of diverse members of *Gokushovirinae* and Alpavirinae subfamilies.

367     Although there is evidence that described Alpavirinae genomes constitute a third

368     group of the Microviridae family (Krupovic and Forterre 2011; Roux et al. 2012), they

369     correspond to prophages, which makes it difficult to integrate them into the taxonomy

370     of the International Committee on Taxonomy of Viruses (ICTV), thus, no contigs

371     were annotated as Alpavirinae prior to application of the HMM profiles.

372        For each taxonomic group of viruses, there is a corresponding set of bacterial

373     hosts. From the 16S rRNA gene diversity data we used to select the twin pairs, it is

374     clear which bacteria phyla contribute the most to the differences in the microbiomes

375     of concordant and discordant twins. But unlike for bacteria, we were not able to

376    discern such clear patterns by order or family in the virome. Indeed, most of the

377    bacteriophage diversity is grouped in just one order, *Caudovirales,* and its three

378    families *Myoviridae*, *Podoviridae* and *Siphoviridae*. Representatives of these families

379    can infect unrelated hosts (Barylski et al., 2017). As such, we wouldn't necessarily

380    expect specific orders or families of viruses to show the patterns observed in the

381    bacterial phyla.

382        Finally, we noted an interesting pattern of complete bacterial genome

383    coverage for select bacteria in the genomes. As these putative contaminants were

384    not the most abundant members of the microbiome, they are unlikely to represent

385    random contamination of bulk DNA. Why certain bacterial genomes showed such

386    high coverage is unclear. One possibility is that we are observing the host species

387    range of transposable phages. Phages such as the Mu phage randomly integrate

388    into the host genome (Taylor, 1963), amplify by successive rounds of replicative

389    transposition, and then can package any section of their host's genome (Hulo et al.,

390    2011; Toussaint and Rice, 2017). Intriguingly, several of the contaminants detected

391    here (*e.g., B. vulgatus*, *B. dorei*, *F. prausnitzii* and *B. thetaiotaomicron*) have also

392    been reported as contaminants in other human gut virome studies (Minot et al.,

393    2011; Roux et al., 2013), which could indicate host-specificity of Mu phages.

394    Alternative explanations include vesicle production, gene transfer agents and/or

395    generalized transduction processes (Biller et al., 2014; McDaniel et al., 2010; Minot

396    et al., 2011). Further comparisons of whole bacterial genomes recovered in diverse

397    virome datasets may help shed light on their source, particularly if the same bacterial

398    species are recovered across multiple studies.

399      ***Prospectus*** – Our results show that gut microbiome richness and diversity

400      correlate to  virome richness and diversity, and vice-versa. The mechanics

401      underlying this association remain to be resolved for the human gut. That the two are

402      coupled may be useful to take into consideration when designing future studies of

403      the virome and factors affecting. Baseline microbiome diversity may be important to

404      balance between groups, for instance, prior to assessing the diversity of the virome.

405

406      **METHODS**

407      **Selection of concordant and discordant monozygotic twin pairs -** From

408      16S rRNA gene diversity previously measured for 354 monozygotic twin pairs whose

409      fecal samples were received between January 28th 2013 and July 14th 2014

410      (Goodrich et al., 2014), we selected 11 concordant and 13 discordant MZ co-twins

411      based on three microbiota β-diversity distances within twin pairs: unweighted

412      UniFrac, weighted UniFrac (Lozupone et al., 2007) and Bray-Curtis (Bray and Curtis,

413      1957). The twins pairs in the the concordant and the discordant groups were

414      selected to be balanced between those two groups for age, BMI, and BMI difference

415      within a twin pair **(TableS1)**. Twins within the concordant group ranged in age from

416      23 to 77 years old and included 5 men and 4 women, while those in the discordant

417      group ranged in age from 29 to 81 years old with 5 men and 7 women.

418      **Isolation of virus-like particles (VLPs) from human fecal samples -** VLP

419      isolation procedures were based on the protocol described by (Gudenkauf et al.,

420      2014) and Minot *et al. (Minot et al., 2013).* For VLP isolation, ~0.5 g of fecal sample

421      was resuspended by vortexing for 5-10 minutes in 15 ml PBS, previously filtered

422      through 0.02 μm filter (Whatman). The homogenates were centrifuged for 30 min at

423    4,500 *x*g, and the supernatant was filtered through 0.22 μm polyethersulfone (PES)

424    Express Plus Millipore Stericup (150 ml) to remove cell debris and bacterial-sized

425    particles. The filtrate was then concentrated on a Millipore Amicon Ultra-15

426    Centrifugal Filter Unit 100K to ~1 ml. The concentrate was transferred to 5 Prime

427    Phase Lock Gel and incubated with 200 μl chloroform for 10 min at room

428    temperature. After being centrifuged for 1 min at 15,000 *x*g, the aqueous layer was

429    transferred to a new microcentrifuge tube, and was treated with Invitrogen TURBO

430    DNase (14 U), Promega RNase One (20 U) and 1 μl Benzonase Nuclease (E1014

431    Sigma Benzonase® Nuclease) at 37 ℃ for 3 hr (Gudenkauf and Hewson, 2016;

432    Reyes et al., 2012). After incubation, 0.04 volumes 0.5 M EDTA was added to each

433    sample. The sample was then stored at -80 ℃ before further processing.

434        **Viral DNA shotgun sequencing -** The viral DNA was extracted with

435    PureLink® Viral RNA/DNA Mini Kit from Invitrogen™. Each viral DNA sample was

436    then amplified using GenomePlex® Complete Whole Genome Amplification (WGA2)

437    Kit from Sigma-Aldrich (Gudenkauf and Hewson, 2016). Two blank controls were

438    included in this step, but very low yield precluded library construction. The amplified

439    product was then fragmented with Covaris S2 Adaptive Focused Acoustic Disruptor

440    with the parameters set as follows: the duty cycle set at 10%, cycle per burst 200,

441    intensity 4 and duration 60 seconds. Each viral sequencing library was prepared

442    following Illumina TruSeq DNA Preparation Protocol with one unique barcode per

443    sample. All barcoded libraries were pooled together. Half of the pool was size

444    selected by BluePippin (Sage Science, Beverly, MA, USA) to enrich fragments with

445    longer inserts (425 bp to 875 bp including the adapters). Both pools, the "large-

446    insert-size library" and the "short-insert-size library", were sequenced in independent

447     lanes on an Illumina HiSeq 2500 instrument, operating in Rapid Run Mode with 250

448     bp paired-end chemistry at the Cornell Genomics facility.

449     **Whole fecal metagenome shotgun sequencing** - The genomic DNA was

450     isolated from an aliquot of ~100 mg from each sample using the PowerSoil® - htp

451     DNA isolation kit (MoBio Laboratories Ltd, Carlsbad, CA). Each sequencing library

452     was then prepared following Illumina TruSeq DNA Preparation Protocol with 500 ng

453     DNA using the gel-free method, 14 cycles of PCR, and with one unique barcode per

454     sample. Sequencing was performed on an Illumina HiSeq 2500 instrument in Rapid

455     Run mode with 2x150 bp paired-end chemistry at the Cornell Biotechnology

456     Resource Center Genomics Facility.

457     **Assessment of Bacterial Contamination -** A set of 8,163 finished bacterial

458     genomes was retrieved from the NCBI FTP on 21 February 2017. Reads per sample

459     were mapped against this bacterial genomes dataset using Bowtie2 v.2.2.8

460     (Langmead and Salzberg, 2012) with the following parameters: --local --maxins 800 -

461     k=3. Genome coverage per base was calculated considering only reads with a

462     mapping quality above 20 using *view* and *depth* Samtools commands v.1.5 (Li et al.,

463     2009). Next, genome coverage was averaged for 100Kbp bins. We observed that

464     evenly covered genomes had a median bin coverage of at least 100; those genomes

465     with a median bin coverage greater than 100 were considered as contaminants. The

466     reads mapping to those genomes were removed. Bacterial genomes can have one

467     or more prophage(s) in their genomes (Munson-McGee et al., 2018) bursting events

468     of those prophages can occur, generating several VLPs. As a conservative measure

469     to avoid the loss of reads originating from prophages and not the bacterial genome

470     *per se*, bins with a coverage over three standard deviations of the bacterial mean

471     coverage were also identified and catalogued as prophages-like regions. Reads

472     mapping to potential contaminant genomes were tagged as "contaminants" and

473     removed from further analysis while reads mapping to high coverage bins were

474     tagged as "possible prophages".

475          A matrix of the abundance of each potential contaminant per sample was built

476     using an in-house Python script and normalized by RPKM. In parallel, from Goodrich

477     *et al.* data (Goodrich et al., 2014), the relative abundance of each OTU was

478     recovered and summarized at the species level using summarize_taxa.py qiime

479     script. The Spearman rank order correlation between relative abundances of

480     contaminants and their corresponding 16S rRNAs data was calculated for species in

481     both sets.

482          **Functional profiles -** The joined and trimmed reads from the "short-insert-

483     size library" were mapped onto Integrated Gene Catalogs (IGC), an integrated

484     catalog of reference genes in the human gut microbiome (Li et al., 2014) by BLASTX

485     using DIAMOND v.0.7.5 (Buchfink et al., 2015) with maximum e-value cutoff 0.001,

486     and maximum number of target sequences to report set to 25.

487          After the mapping onto IGC, an abundance matrix was generated using an in-

488     house Python script. The matrix was then annotated according to the KEGG

489     annotation of each gene provided by IGC. The annotated abundance matrix was

490     rarefied (subsampling without replacement) to 2,000,000 read hits per sample. The

491     KEGG functional profile was then generated using QIIME 1.9 (Quantitative Insights

492     Into Microbial Ecology) (Caporaso et al., 2010) using the command

493     summarize_taxa_through_plots.py. The Intraclass Correlation Coefficient of the

494     functional profiles for each group (additional microbiomes, additional viromes,

495    viromes of concordant-microbiome samples and viromes of discordant-microbiome

496    samples) was calculated using the Psych R package.

497         ***De-novo* assembly -** Reads from the "large-insert-size library" that remain

498    paired (forward and reverse) after the trimming step were assembled using

499    Integrated metagenomic assembly pipeline for short reads (InteMAP) (Lai et al.,

500    2015) with insert size 325 bp ± 100 bp. Each sample was assembled separately.

501    After the first run of assembly, all clean reads were mapped to the assembled

502    contigs using Bowtie2 v.2.2.8 (Langmead and Salzberg, 2012) with the following

503    parameter: --local --maxins 800. The pairs of reads that aligned concordantly at least

504    once were then submitted for the second run of assemble by InteMAP. Contigs

505    larger than 500 bp from all samples were pooled together and compared all vs all,

506    using an in-house Perl script, on the comparison file it was possible to identify

507    potential circular genomes, and dereplicate contigs that were contained in over 90%

508    of their length within another contig.

509         In order to build an abundance matrix, the recruitment of reads to the

510    dereplicated metagenomic assemblies was used implementing a filter of coverage

511    and length as recommended in Roux *et al.* (Roux et al., 2017). With this in mind,

512    reads (not tagged as contaminants in the previous step) were mapped to

513    dereplicated contigs using Rsubread v.1.28.0 (Liao et al., 2013). Mapping outputs

514    were parsed using an in-house Python script into an abundance matrix that was

515    normalized by reads per kilobase of contig length per million sequenced reads per

516    sample (RPKM) and transformed to $Log_{10}(x+1)$, being *x* the normalized abundance.

517    Contigs with a normalized coverage bellow 5x were excluded. Finally, to virotypes, a

518    filter on contig length was applied. A length threshold was chosen as the elbow of

519  the decay curve generated when plotting the number of contigs as a function of

520  length, which occurred at a length of 1,300 bp.

521  **HMM annotation -** Independent HMM profiles were built to identify crAss-like

522  contigs and Microviridae contigs. To build the HMM-crAsslike profile, sequences for

523  the Major Capsid Protein (MCP) of the proposed crAss-like family (Yutin et al., 2018)

524  were retrieved from ftp.ncbi.nih.gov/pub/yutinn/crassphage_2017/. Multiple

525  sequence alignments (MSA) were done using MUSCLE v.3.8.31(Edgar, 2004) and

526  inspected using UGENE v.1.31.0 (Okonechnikov et al., 2012); positions with more

527  than 30% of gaps were removed. Finally, the HMM-crAsslike profile was built using

528  *hmmbuild* from the HMMER package v.3.1b2 (http://hmmer.org/) (Eddy, 1998). For

529  the Microviridae case, all HMM-profiles for the viral protein 1 (VP1) developed by

530  Alves *et al.* (Alves et al., 2016) were adopted.

531  Predicted proteins of the assembled contigs were queried for matching the

532  HMM-profiles using *hmmsearch* (Eddy, 1998). Matching proteins with an e-value

533  below $1 \times 10^{-5}$ were considered as true homologs but only proteins between the size

534  rank of the reference proteins (crAsslike MCP: 450-510 residues; Microviridae: 450-

535  800 residues), a coverage of at least 50% and a percentage of identity of at least

536  40% to at least one reference sequence were used for further analysis. Coverage

537  and identity percentage were determined making a BLASTp of the true homologues

538  against the reference sequences.

539  True homologues passing the filters mentioned above were used in

540  phylogenetic analysis. Reference and homologous sequences were aligned using

541  MUSCLE v.3.8.31 and sites with at least 30% of gaps were removed using UGENE

542  v.1.31.0. A maximum-likelihood (ML) phylogenetic analysis was done using RAxML

23

543     v.8.2.4 (Stamatakis, 2014), the best evolutive model was obtained with prottest

544     v.3.4.2 (Darriba et al., 2011) and support for nodes in the ML trees were obtained by

545     bootstrap with 100 pseudoreplicates.

546     **Taxonomic profiles -** To infer the taxonomic affiliation of the assembled

547     VLPs, genes were predicted from all assembled contigs larger than 500 bp using

548     GeneMarkS v.4.32 (Besemer et al., 2001). The amino acid sequence of the

549     predicted genes was then used in a BLASTp search against the NR NCBI viral

550     database using DIAMOND v.0.7.5 (Buchfink et al., 2015) with maximum e-value

551     cutoff 0.001 and maximum number of target sequences to report set to 25. Using the

552     BLASTp results, the taxonomy of each gene was assigned by the lowest-common-

553     ancestor algorithm in MEtaGenome ANalyzer (MEGAN5) v.5.11.3 (Huson et al.,

554     2011) with the following parameters: Min Support: 1, Min Score: 40.0, Max Expected:

555     0.01, Top Percent: 10.0, Min-Complexity filter: 0.44. Independently, the taxonomy

556     annotation of each contig was obtained using CENTRIFUGE v.1.0.4 (Kim et al.,

557     2016) against the NT NCBI viral genomes database. The final taxonomic annotation

558     of each contig was then assigned using a voting system where the taxonomic

559     annotation of each protein and the CENTRIFUGE annotation of the contig were

560     considered as votes. With all the possible votes for a contig, an N-ary tree was build

561     and the weight of each node was the number of votes including that node. The

562     taxonomic annotation of a contig will be the result of traverse the tree passing

563     through the heaviest nodes with one consideration: if all children nodes of a node

564     have the same weight the traversing must be stopped. The taxonomic profile was

565     considered as a subset of the recruitment matrix containing all contigs annotated

566     either by the voting system or annotated through the HMM profiles (see above).

24

567    **Prediction of phage-host interaction -** Clustered Regularly Interspaced

568    Short Palindromic Repeats (CRISPRs) were identified using the PilerCR program

569    v.1.06 (Edgar, 2007) from the same set of 8,163 bacterial used to asses the bacterial

570    contamination. Spacers within the expected size of 20 bp and 72 bp (Horvath and

571    Barrangou, 2010) were used as queries against virotypes and taxonomically

572    annotated contigs using BLASTn (v.2.6.0+) with short query parameters (Camacho

573    et al., 2009). Matches covering at least 90% of the spacer and with an e-value <

574    0.001 were considered to be CRISPR spacer-virus associations. Additionally,

575    virotypes and taxonomically annotated contigs were mapped against the

576    representatives genomes of the viral clusters in the MVP database (Gao et al., 2018)

577    using LAST-959 (Kiełbasa et al., 2011). As viral clusters in MVP comprise

578    sequences that have at least 95% identity along at least 80% of their lengths, only

579    matches that fulfill those constraints were kept. The host(s) of a contig was

580    determined from its matching viral cluster.

581    **Diversity indexes -** The Shannon diversity index within-samples (α-diversity)

582    and the Hellinger distance within co-twins (β-diversity) were calculated using

583    *diversity* and *vegdist* functions of Vegan R package for all three abundance matrices

584    generated (function, taxonomy and read recruitment matrices). Correlations between

585    virome α-diversity and microbiome α-diversity were measured using the Pearson

586    correlation coefficient. Correlations between viromes β-diversity and the

587    microbiomes β-diversity was computed with a the Mantel test using the Pearson

588    correlation coefficient. Additionally, the β-diversity between concordant MZ co-twins

589    was compared to the β-diversity between discordant MZ co-twins; p values were

590    calculated using Mann-Whitney U test.

**DATA AND SOFTWARE AVAILABILITY**

591

592     Jupyter notebooks and scripts describing the data analysis process are

593     available on GitHub at https://github.com/leylabmpi/TwinsUK_virome

594     The sequence data have been deposited in the European Nucleotide Archive under

595     the study accession number PRJEB29491.

**ACKNOWLEDGEMENTS**

596

**AUTHOR CONTRIBUTIONS**

606

607     RL and SPC designed the study. TS and JB were involved in sample

608     collection. SPC and IH generated the data. JLM-G, SPC, JKG, NY, AR and RL

609     analyzed the data. JLM-G, SPC, SCD, AR and RL wrote the manuscript. All authors

610     read and approved the final manuscript.

**DECLARATION OF INTERESTS**

611

612     The authors declare no competing interests.

**REFERENCES**

Alves, J.M.P., de Oliveira, A.L., Sandberg, T.O.M., Moreno-Gallego, J.L., de Toledo, M.A.F., de Moura, E.M.M., Oliveira, L.S., Durham, A.M., Mehnert, D.U., Zanotto, P.M. de A., et al. (2016). GenSeed-HMM: A tool for progressive assembly using profile HMMs as seeds and its application in Alpavirinae viral discovery from metagenomic data. Front. Microbiol. *7*, 269.

Barylski, J., Enault, F., Dutilh, B.E., Schuller, M.B.P., Edwards, R.A., Gillis, A., Klumpp, J., Knezevic, P., Krupovic, M., Kuhn, J.H., et al. (2017). Genomic, proteomic, and phylogenetic analysis of spounaviruses indicates paraphyly of the order Caudovirales.

Besemer, J., Lomsadze, A., and Borodovsky, M. (2001). GeneMarkS: a self-training method for prediction of gene starts in microbial genomes. Implications for finding sequence motifs in regulatory regions. Nucleic Acids Res. *29*, 2607–2618.

Biller, S.J., Schubotz, F., Roggensack, S.E., Thompson, A.W., Summons, R.E., and Chisholm, S.W. (2014). Bacterial vesicles in marine ecosystems. Science *343*, 183–186.

Bray, J.R., and Curtis, J.T. (1957). An Ordination of the Upland Forest Communities of Southern Wisconsin. Ecol. Monogr. *27*, 326–349.

Breitbart, M., Hewson, I., Felts, B., Mahaffy, J.M., Nulton, J., Salamon, P., and Rohwer, F. (2003). Metagenomic analyses of an uncultured viral community from human feces. J. Bacteriol. *185*, 6220–6223.

Breitbart, M., Haynes, M., Kelley, S., Angly, F., Edwards, R.A., Felts, B., Mahaffy, J.M., Mueller, J., Nulton, J., Rayhawk, S., et al. (2008). Viral diversity and dynamics in an infant gut. Res. Microbiol. *159*, 367–373.

Buchfink, B., Xie, C., and Huson, D.H. (2015). Fast and sensitive protein alignment using DIAMOND. Nat. Methods *12*, 59–60.

Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L. (2009). BLAST+: architecture and applications. BMC Bioinformatics *10*, 421.

Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K., Fierer, N., Peña, A.G., Goodrich, J.K., Gordon, J.I., et al. (2010). QIIME allows analysis of high-throughput community sequencing data. Nat. Methods *7*, 335–336.

Castro-Mejía, J.L., Muhammed, M.K., Kot, W., Neve, H., Franz, C.M.A.P., Hansen, L.H., Vogensen, F.K., and Nielsen, D.S. (2015). Optimizing protocols for extraction of bacteriophages prior to metagenomic analyses of phage communities in the human gut. Microbiome *3*, 64.

Claesson, M.J., Jeffery, I.B., Conde, S., Power, S.E., O'Connor, E.M., Cusack, S., Harris, H.M.B., Coakley, M., Lakshminarayanan, B., O'Sullivan, O., et al. (2012). Gut microbiota composition correlates with diet and health in the elderly. Nature *488*,

652    178–184.

653    Colson, P., Fancello, L., Gimenez, G., Armougom, F., Desnues, C., Fournous, G.,
654    Yoosuf, N., Million, M., La Scola, B., and Raoult, D. (2013). Evidence of the
655    megavirome in humans. J. Clin. Virol. *57*, 191–200.

656    Cotillard, A., Kennedy, S.P., Kong, L.C., Prifti, E., Pons, N., Le Chatelier, E.,
657    Almeida, M., Quinquis, B., Levenez, F., Galleron, N., et al. (2013). Dietary
658    intervention impact on gut microbial gene richness. Nature *500*, 585–588.

659    Darriba, D., Taboada, G.L., Doallo, R., and Posada, D. (2011). ProtTest 3: fast
660    selection of best-fit models of protein evolution. Bioinformatics *27*, 1164–1165.

661    David, L.A., Materna, A.C., Friedman, J., Campos-Baptista, M.I., Blackburn, M.C.,
662    Perrotta, A., Erdman, S.E., and Alm, E.J. (2014). Host lifestyle affects human
663    microbiota on daily timescales. Genome Biol. *15*, R89.

664    De Filippo, C., Cavalieri, D., Di Paola, M., Ramazzotti, M., Poullet, J.B., Massart, S.,
665    Collini, S., Pieraccini, G., and Lionetti, P. (2010). Impact of diet in shaping gut
666    microbiota revealed by a comparative study in children from Europe and rural Africa.
667    Proc. Natl. Acad. Sci. U. S. A. *107*, 14691–14696.

668    Dutilh, B.E., Cassman, N., McNair, K., Sanchez, S.E., Silva, G.G.Z., Boling, L., Barr,
669    J.J., Speth, D.R., Seguritan, V., Aziz, R.K., et al. (2014). A highly abundant
670    bacteriophage discovered in the unknown sequences of human faecal
671    metagenomes. Nat. Commun. *5*, 4498.

672    Eddy, S.R. (1998). Profile hidden Markov models. Bioinformatics *14*, 755–763.

673    Edgar, R.C. (2004). MUSCLE: multiple sequence alignment with high accuracy and
674    high throughput. Nucleic Acids Res. *32*, 1792–1797.

675    Edgar, R.C. (2007). PILER-CR: fast and accurate identification of CRISPR repeats.
676    BMC Bioinformatics *8*, 18.

677    Gao, N.L., Zhang, C., Zhang, Z., Hu, S., Lercher, M.J., Zhao, X.-M., Bork, P., Liu, Z.,
678    and Chen, W.-H. (2018). MVP: a microbe–phage interaction database. Nucleic Acids
679    Res. *46*, D700–D707.

680    Goodrich, J.K., Waters, J.L., Poole, A.C., Sutter, J.L., Koren, O., Blekhman, R.,
681    Beaumont, M., Van Treuren, W., Knight, R., Bell, J.T., et al. (2014). Human genetics
682    shape the gut microbiome. Cell *159*, 789–799.

683    Goodrich, J.K., Davenport, E.R., Beaumont, M., Jackson, M.A., Knight, R., Ober, C.,
684    Spector, T.D., Bell, J.T., Clark, A.G., and Ley, R.E. (2016). Genetic determinants of
685    the gut microbiome in UK Twins. Cell Host Microbe *19*, 731–743.

686    Gudenkauf, B.M., and Hewson, I. (2016). Comparative metagenomics of viral
687    assemblages inhabiting four phyla of marine invertebrates. Frontiers in Marine
688    Science *3*, 23.

689    Gudenkauf, B.M., Eaglesham, J.B., Aragundi, W.M., and Hewson, I. (2014).

690  Discovery of urchin-associated densoviruses (family Parvoviridae) in coastal waters
691  of the Big Island, Hawaii. J. Gen. Virol. *95*, 652–658.

692  Halary, S., Temmam, S., Raoult, D., and Desnues, C. (2016). Viral metagenomics:
693  are we missing the giants? Curr. Opin. Microbiol. *31*, 34–43.

694  Horvath, P., and Barrangou, R. (2010). CRISPR/Cas, the immune system of bacteria
695  and archaea. Science *327*, 167–170.

696  Hoyles, L., McCartney, A.L., Neve, H., Gibson, G.R., Sanderson, J.D., Heller, K.J.,
697  and van Sinderen, D. (2014). Characterization of virus-like particles associated with
698  the human faecal and caecal microbiota. Res. Microbiol. *165*, 803–812.

699  Hulo, C., de Castro, E., Masson, P., Bougueleret, L., Bairoch, A., Xenarios, I., and Le
700  Mercier, P. (2011). ViralZone: a knowledge resource to understand virus diversity.
701  Nucleic Acids Res. *39*, D576–D582.

702  Huson, D.H., Mitra, S., Ruscheweyh, H.-J., Weber, N., and Schuster, S.C. (2011).
703  Integrative analysis of environmental sequences using MEGAN4. Genome Res. *21*,
704  1552–1560.

705  Kiełbasa, S.M., Wan, R., Sato, K., Horton, P., and Frith, M.C. (2011). Adaptive seeds
706  tame genomic sequence comparison. Genome Res. *21*, 487–493.

707  Kim, D., Song, L., Breitwieser, F.P., and Salzberg, S.L. (2016). Centrifuge: rapid and
708  sensitive classification of metagenomic sequences. Genome Res. *26*, 1721–1729.

709  Knowles, B., Silveira, C.B., Bailey, B.A., Barott, K., Cantu, V.A., Cobián-Güemes,
710  A.G., Coutinho, F.H., Dinsdale, E.A., Felts, B., Furby, K.A., et al. (2016). Lytic to
711  temperate switching of viral communities. Nature *531*, 466–470.

712  Lai, B., Wang, F., Wang, X., Duan, L., and Zhu, H. (2015). InteMAP: Integrated
713  metagenomic assembly pipeline for NGS short reads. BMC Bioinformatics *16*, 1–14.

714  Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie
715  2. Nat. Methods *9*, 357–359.

716  Lee, S., Sung, J., Lee, J., and Ko, G. (2011). Comparison of the gut microbiotas of
717  healthy adult twins living in South Korea and the United States. Appl. Environ.
718  Microbiol. *77*, 7433–7437.

719  Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,
720  Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing Subgroup
721  (2009). The sequence alignment/map format and SAMtools. Bioinformatics *25*,
722  2078–2079.

723  Li, J., Jia, H., Cai, X., Zhong, H., Feng, Q., Sunagawa, S., Arumugam, M., Kultima,
724  J.R., Prifti, E., Nielsen, T., et al. (2014). An integrated catalog of reference genes in
725  the human gut microbiome. Nat. Biotechnol. *32*, 834–841.

726  Liao, Y., Smyth, G.K., and Shi, W. (2013). The Subread aligner: fast, accurate and
727  scalable read mapping by seed-and-vote. Nucleic Acids Res. *41*, e108.

728 Lim, E.S., Zhou, Y., Zhao, G., Bauer, I.K., Droit, L., Ndao, I.M., Warner, B.B., Tarr,
729 P.I., Wang, D., and Holtz, L.R. (2015). Early life dynamics of the human gut virome
730 and bacterial microbiome in infants. Nat. Med. *21*, 1228–1234.

731 Lozupone, C.A., Hamady, M., Kelley, S.T., and Knight, R. (2007). Quantitative and
732 qualitative beta diversity measures lead to different insights into factors that structure
733 microbial communities. Appl. Environ. Microbiol. *73*, 1576–1585.

734 Manrique, P., Bolduc, B., Walk, S.T., van der Oost, J., de Vos, W.M., and Young,
735 M.J. (2016). Healthy human gut phageome. Proc. Natl. Acad. Sci. U. S. A. *113*,
736 10400–10405.

737 Manrique, P., Dills, M., and Young, M.J. (2017). The human gut phage community
738 and its implications for health and disease. Viruses *9*, 10.3390/v9060141.

739 McDaniel, L.D., Young, E., Delaney, J., Ruhnau, F., Ritchie, K.B., and Paul, J.H.
740 (2010). High frequency of horizontal gene transfer in the oceans. Science *330*, 50.

741 Minot, S., Sinha, R., Chen, J., Li, H., Keilbaugh, S.A., Wu, G.D., Lewis, J.D., and
742 Bushman, F.D. (2011). The human gut virome: inter-individual variation and dynamic
743 response to diet. Genome Res. *21*, 1616–1625.

744 Minot, S., Bryson, A., Chehoud, C., Wu, G.D., Lewis, J.D., and Bushman, F.D.
745 (2013). Rapid evolution of the human gut virome. Proc. Natl. Acad. Sci. U. S. A. *110*,
746 12450–12455.

747 Munson-McGee, J.H., Peng, S., Dewerff, S., Stepanauskas, R., Whitaker, R.J.,
748 Weitz, J.S., and Young, M.J. (2018). A virus or more in (nearly) every cell: ubiquitous
749 networks of virus-host interactions in extreme environments. ISME J.

750 Ogilvie, L.A., and Jones, B.V. (2017). The human gut virome: form and function.
751 Emerging Topics in Life Sciences *1*, 351–362.

752 Okonechnikov, K., Golosova, O., Fursov, M., and UGENE team (2012). Unipro
753 UGENE: a unified bioinformatics toolkit. Bioinformatics *28*, 1166–1167.

754 Palmer, C., Bik, E.M., DiGiulio, D.B., Relman, D.A., and Brown, P.O. (2007).
755 Development of the human infant intestinal microbiota. PLoS Biol. *5*, e177.

756 Reyes, A., Haynes, M., Hanson, N., Angly, F.E., Heath, A.C., Rohwer, F., and
757 Gordon, J.I. (2010). Viruses in the faecal microbiota of monozygotic twins and their
758 mothers. Nature *466*, 334–338.

759 Reyes, A., Semenkovich, N.P., Whiteson, K., Rohwer, F., and Gordon, J.I. (2012).
760 Going viral: next-generation sequencing applied to phage populations in the human
761 gut. Nat. Rev. Microbiol. *10*, 607–617.

762 Reyes, A., Wu, M., McNulty, N.P., Rohwer, F.L., and Gordon, J.I. (2013). Gnotobiotic
763 mouse model of phage-bacterial host dynamics in the human gut. Proc. Natl. Acad.
764 Sci. U. S. A. *110*, 20236–20241.

765 Reyes, A., Blanton, L.V., Cao, S., Zhao, G., Manary, M., Trehan, I., Smith, M.I.,

766  Wang, D., Virgin, H.W., Rohwer, F., et al. (2015). Gut DNA viromes of Malawian
767  twins discordant for severe acute malnutrition. Proc. Natl. Acad. Sci. U. S. A. *112*,
768  11941–11946.

769  Rodriguez-Brito, B., Li, L., Wegley, L., Furlan, M., Angly, F., Breitbart, M., Buchanan,
770  J., Desnues, C., Dinsdale, E., Edwards, R., et al. (2010). Viral and microbial
771  community dynamics in four aquatic environments. ISME J. *4*, 739–751.

772  Rodriguez-Valera, F., Martin-Cuadrado, A.-B., Rodriguez-Brito, B., Pasić, L.,
773  Thingstad, T.F., Rohwer, F., and Mira, A. (2009). Explaining microbial population
774  genomics through phage predation. Nat. Rev. Microbiol. *7*, 828–836.

775  Roux, S., Krupovic, M., Debroas, D., Forterre, P., and Enault, F. (2013). Assessment
776  of viral community functional potential from viral metagenomes may be hampered by
777  contamination with cellular sequences. Open Biol. *3*, 130160.

778  Roux, S., Emerson, J.B., Eloe-Fadrosh, E.A., and Sullivan, M.B. (2017).
779  Benchmarking viromics: an in silico evaluation of metagenome-enabled estimates of
780  viral community composition and diversity. PeerJ *5*, e3817.

781  Sender, R., Fuchs, S., and Milo, R. (2016). Are we really vastly outnumbered?
782  revisiting the ratio of bacterial to host cells in humans. Cell *164*, 337–340.

783  Shkoporov, A., Khokhlova, E.V., Brian Fitzgerald, C., Stockdale, S.R., Draper, L.A.,
784  Paul Ross, R., and Hill, C. (2018). ΦCrAss001, a member of the most abundant
785  bacteriophage family in the human gut, infects Bacteroides.

786  Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-
787  analysis of large phylogenies. Bioinformatics *30*, 1312–1313.

788  Suttle, C.A. (2007). Marine viruses--major players in the global ecosystem. Nat. Rev.
789  Microbiol. *5*, 801–812.

790  Taylor, A.L. (1963). Bacteriophage-induced mutation in Escherichia coli. Proc. Natl.
791  Acad. Sci. U. S. A. *50*, 1043–1051.

792  Thingstad, T.F. (2000). Elements of a theory for the mechanisms controlling
793  abundance, diversity, and biogeochemical role of lytic bacterial viruses in aquatic
794  systems. Limnol. Oceanogr. *45*, 1320–1328.

795  Thingstad, T.F., Våge, S., Storesund, J.E., Sandaa, R.-A., and Giske, J. (2014). A
796  theoretical analysis of how strain-specific viruses can control microbial species
797  diversity. Proc. Natl. Acad. Sci. U. S. A. *111*, 7813–7818.

798  Tims, S., Derom, C., Jonkers, D.M., Vlietinck, R., Saris, W.H., Kleerebezem, M., de
799  Vos, W.M., and Zoetendal, E.G. (2013). Microbiota conservation and BMI signatures
800  in adult monozygotic twins. ISME J. *7*, 707–717.

801  Toussaint, A., and Rice, P.A. (2017). Transposable phages, DNA reorganization and
802  transfer. Curr. Opin. Microbiol. *38*, 88–94.

803  Turnbaugh, P.J., Hamady, M., Yatsunenko, T., Cantarel, B.L., Duncan, A., Ley, R.E.,

804    Sogin, M.L., Jones, W.J., Roe, B.A., Affourtit, J.P., et al. (2009). A core gut
805    microbiome in obese and lean twins. Nature *457*, 480–484.

806    Weitz, J.S., and Dushoff, J. (2008). Alternative stable states in host–phage
807    dynamics. Theor. Ecol. *1*, 13–19.

808    Wu, G.D., Chen, J., Hoffmann, C., Bittinger, K., Chen, Y.-Y., Keilbaugh, S.A.,
809    Bewtra, M., Knights, D., Walters, W.A., Knight, R., et al. (2011). Linking long-term
810    dietary patterns with gut microbial enterotypes. Science *334*, 105–108.

811    Yarygin, K., Tyakht, A., Larin, A., Kostryukova, E., Kolchenko, S., Bitner, V., and
812    Alexeev, D. (2017). Abundance profiling of specific gene groups using precomputed
813    gut metagenomes yields novel biological hypotheses. PLoS One *12*, e0176154.

814    Yatsunenko, T., Rey, F.E., Manary, M.J., Trehan, I., Dominguez-Bello, M.G.,
815    Contreras, M., Magris, M., Hidalgo, G., Baldassano, R.N., Anokhin, A.P., et al.
816    (2012). Human gut microbiome viewed across age and geography. Nature *486*, 222–
817    227.

818    Yutin, N., Makarova, K.S., Gussow, A.B., Krupovic, M., Segall, A., Edwards, R.A.,
819    and Koonin, E.V. (2018). Discovery of an expansive bacteriophage family that
820    includes the most abundant viruses from the human gut. Nat Microbiol *3*, 38–46.

821

822

823    **FIGURE TITLES AND LEGENDS**

824    **Figure 1. Microbiome discordance in twin pairs. (A)** The β-diversity

825    measures of the microbiotas of 354 monozygotic twin pairs from a previous study

826    (Goodrich et al., 2014) are shown. Each dot represents the β-diversity of a pair of

827    twins, measured by the weighted UniFrac (x-axis), unweighted UniFrac (z-axis), and

828    Bray-Curtis (y-axis) β-diversity metrics. The three β-diversity metrics are in general

829    correlated (Pearson pairwise correlation coefficient > 0.4). The plane is the least

830    squared fitted plane Bray-Curtis ~ Weighted UniFrac + Unweighted UniFrac. A

831    subset of twin pairs with concordant microbiotas (blue) and discordant microbiotas

832    (orange) were chosen from the two edges. Black dots indicate the samples used for

833    virome and whole fecal metagenome comparison. **(B)** Comparison of the taxonomic

834    profiles (relative abundance) at the Phylum level for the 21 MZ twin pairs concordant

835    (1-9) or discordant (10-21) for their microbiotas. **(C)** Differences in the relative

836    abundances for the major phyla for concordant (blue points, n=9) and discordant

837    (orange points, n=12) twin pairs. Mann-Whitney's U test. *** $p < 0.0005$, * $p = 0.055$

838

839    **Figure 2. Bacterial contamination in VLP preparations. (A)** Heatmap of

840    VLP reads from sample 4A mapping to bacterial genomes before and after the

841    removal of reads determined as contaminants. Genomes are sorted by length and

842    split in bins of 100,000 bp. Bacterial genomes with a median coverage greater than

843    100 were considered as contaminants. **(B)** Cladogram based on the NCBI taxonomy

844    of the 65 genomes identified as contaminants across all VLP extractions. **(Right)**

845    Spearman rank correlation coefficient (rho) between the abundance of the bacterial

846    genomes in the VLP extractions and 16S rRNA gene profile from the microbiome.

847    **(Left)** Total abundance of each bacterial genome added across all individuals.

848

849 **Figure 3. Comparison of the gene content of whole fecal metagenomes**

850 **and viromes.** Relative abundance of KEGG categories in whole fecal metagenomes

851 and viromes. **(A)** The relative abundance of KEGG categories in whole fecal

852 metagenomes and viromes, including all hits to IGC genes, regardless of the

853 annotation. **(B)** Heatmap of the relative abundance of the second level of KEGG

854 categories in whole fecal metagenomes and viromes, excluding the IGC genes with

855 unknown annotation. A.V.: Additional viromes; A.M.: Additional microbiomes (whole

856 genome extractions). Intra-class coefficient (ICC) for A.M. = 0.99; ICC for A.V. =

857 0.85; ICC concordant-microbiome co-twins = 0.69; ICC discordant-microbiome co-

858 twins = 0.68.

859

860 **Figure 4. Virome composition.** Comparison of the taxonomic profiles at the

861 Family level for the 21 MZ twin pairs concordant (1-9) or discordant (10-21) for their

862 viromes. **(A)** The viral family composition of the MZ twins. **(B)** Differences of the

863 relative abundances of each family for concordant (blue points, n=9) and discordant

864 (orange points, n=12) twin pairs.

865

866 **Figure 5. Bacteriophages diversity correlates with microbiome diversity**

867 **but eukaryotic viruses diversity do not. (A)** Correlation of Shannon α-diversity of

868 viromes to Shannon α-diversity of microbiomes (n=42). **i) Virotypes:** Pearson

869 correlation coefficient = 0.406, m = 0.3, p = 0.007, $R^2$ = 0.165; **ii) Taxonomy:**

870 Pearson correlation coefficient = 0.389, m = 0.25, p = 0.010, $R^2$ = 0.151; iii) **Genes:**

871 Pearson correlation coefficient = 0.105, m = 0.11, p = 0.506, $R^2$ = 0.011 **(B)**

872 Correlation of the Shannon α-diversity of the virome, calculated from contigs

873      annotated as ssDNA eukaryotic viruses, ssDNA phages, dsDNA eukaryotic viruses,

874      and dsDNA phages, to Shannon α-diversity of the microbiome (n=42). **ssDNA**

875      **eukaryotic viruses:** Pearson correlation coefficient = 0.027, m = 0.034, p = 0.863,

876      $R^2$ = 0.000751; **ssDNA bacteriophages:** Pearson correlation coefficient = 0.394, m

877      = 0.35, p = 0.009, $R^2$ = 0.155; **dsDNA eukaryotic viruses:** Pearson correlation

878      coefficient = 0.143, m = 0.15, p = 0.368, $R^2$ = 0.020; **dsDNA bacteriophages:**

879      Pearson correlation coefficient = 0.400, m = 0.25, p = 0.008, $R^2$ = 0.16.

880

881      **Figure 6. Virome Beta-diversity patterns mirror microbiome Beta-**

882      **diversity.** Box plots show the distribution of Hellinger distances for microbiomes and

883      viromes, according to the three different layers of information recovered (virotypes,

884      function, and taxonomy), for concordant co-twins (blue, n=9), discordant co-twins

885      (orange, n=12), unrelated samples within the concordant co-twins (blue edges,

886      n=144), and unrelated samples within the discordant co-twins (orange edges,

887      n=264). Significant differences between means (Mann-Whitney's U test, p < 0.020)

888      are denoted with different letters.

889

890      **SUPPLEMENTAL INFORMATION LEGENDS**

891      **Figure S1.** Schematic representation summarizing the procedures applied to

892      **(left)** the "large-insert-size library" and **(right)** the "short-insert-size library" to obtain

893      three different layers of information used to analyze the virome diversity of the

894      microbiome-concordant and microbiome-discordant co-twins.

895

896      **Figure S2.** Box plots showing the distribution of the number of shared

897      virotypes between different groups made from the 21 MZ co-twins. (Up left) All co-

898    twins vs unrelated individuals. (Up right) Microbiome-discordant co-twins vs

899    unrelated individuals in the same group. (Down left) Microbiome-concordant co-twins

900    vs unrelated individuals in the same group. (Down right) Microbiome-concordant co-

901    twins vs microbiome-discordant co-twins. Mann-Whitney's U test. * $p < 0.05$; n.s: not

902    significant difference.

903

904        **Figure S3.** Maximum likelihood phylogenetic analysis of **(A)** the VP1 protein

905    of *Microviridae* phages and **(B)** the MCP protein of crAssphage found in the 42 MZ

906    viromes. Reference sequences are in purple, outgroup sequences are in red while

907    the different MCP or VP1 proteins found in this work are labeled in black. Circles in

908    the nodes indicates bootstrap values above 70%.

909

910        **Figure S4.** Cladogram based on the NCBI taxonomy showing the bacteria

911    identified as hosts. The cladogram is summarized by genus, and clades are colored

912    by Phylum. Blue: Firmicutes; Red: Actinobacteria; Yellow: Tenericutes; Green:

913    Proteobacteria; Purple: Bacteroidetes; Light green: Fusobacteria; Magenta:

914    Verrucomicrobia; Light blue: Euryarchaeota. Red bars indicate the number of

915    species in each genus, and green bars show the dereplicated number of contigs

916    associated to each genus (i.e. if a contig was found associated to two species in that

917    genus, it is only shown one time).

918

919        **Figure S5.** Box plots showing the distribution of **(A)** the Jaccard distances

920    and **(B)** Bray-Curtis distances for microbiomes and viromes, according to the three

921    different layers of information recovered (virotypes, function and taxonomy).

922    Significant differences between means (Mann-Whitney's U test) are denoted with

36

923      different letters. Groups and n values as in Figure 6. **(C)** Correlation between virome

924      β-diversity and microbiome β-diversity (n=840). **i) Virotypes:** Pearson correlation

925      coefficient among all individuals = 0.382 (p = 0.0005, Mantel test), m = 0.167, p = 0,

926      $R^2$ = 0.157; Pearson correlation coefficient among co-twins = 0.522, m = 0.188, p =

927      0.015, $R^2$ = 0.1508 ; **ii) Taxonomy annotated contigs:** Pearson correlation

928      coefficient among all individuals = 0.266 (p = 0.003, Mantel test), m = 0.140, p = 0,

929      $R^2$ = 0.0796; Pearson correlation coefficient among co-twins = 0.512, m = 0.186, p =

930      0.017, $R^2$ = 0.224; **iii) Genes:** Pearson correlation coefficient among all individuals =

931      0.344 (p = 0.0009, Mantel test), m = 0.162, p = 0, $R^2$ = 0.123; Pearson correlation

932      coefficient among co-twins = 0.53, m = 0.182, p = 0.012, $R^2$ = 0.248. Lines describe

933      linear regressions of pairwise distances among all individuals. Triangles indicate

934      concordant-microbiome co-twins and squares indicate discordant-microbiome co-

935      twins.

936

937      **Table S1.** Additional information pertaining to the 21 selected MZ twin pairs

938      (metadata), and counts of viromes reads and contigs per sample.

939

940      **Table S2.** Median bin coverage of bacterial genomes by VLP reads per

941      sample.

# Figure 1.

Figure 2.

**Figure 3.**

**Figure 4.**

# Figure 5.

**Figure 6.**

**Figure S1.**

# Figure S2.

# Figure S3.

**A**

**Microviridae subfamilies**
- Pichovirinae
- Alpavirinae
- Gokushovirinae
- Outgroup
- Known representants
- Viromes VP1

Tree scale: 1

**B**



**crAssphage**
- Chlamydia group
- IAS virus group
- crAssphage group
- outgroup
- Known representants
- Viromes MCP

Tree scale: 1

## Figure S4.

# Figure S5.

**A**



**B**



**C**