1  # Members of a highly widespread bacteriophage family are hallmarks
2  # of metabolic syndrome gut microbiomes

3

4  Patrick A. de Jonge[1†], Koen Wortelboer[1†], Torsten P.M. Scheithauer[1], Bert-Jan H. van
5  den Born[1], Aeilko H. Zwinderman[2], Franklin L. Nobrega[3], Bas E. Dutilh[4], Max
6  Nieuwdorp[1], Hilde Herrema[1*]

7

8  [1]: Departments of Internal and Experimental Vascular Medicine, Amsterdam University
9  Medical Centers, Location AMC, Amsterdam, The Netherlands
10 [2]: Department of Clinical Epidemiology, Biostatistics and Bioinformatics, Amsterdam
11 University Medical Centers, Location AMC, University of Amsterdam, Amsterdam, The
12 Netherlands
13 [3]: School of Biological Sciences, Faculty of Environmental and Life Sciences,
14 University of Southampton, Southampton, United Kingdom
15 [4]: Theoretical Biology and Bioinformatics, Science for Life, Utrecht University, Utrecht,
16 the Netherlands
17 [†]: These authors contributed equally to this work.

18

19 [*]Correspondence: h.j.herrema@amsterdamumc.nl

20

23

24 ## Summary

25

26 There is significant interest in altering the course of cardiometabolic disease
27 development via the gut microbiome. Nevertheless, the highly abundant phage
28 members -which impact gut bacteria- of the complex gut ecosystem remain
29 understudied. Here, we characterized gut phageome changes associated with
30 metabolic syndrome (MetS), a highly prevalent clinical condition preceding
31 cardiometabolic disease. MetS gut phageome populations exhibited decreased
32 richness and diversity, but larger inter-individual variation. These populations were
33 enriched in phages infecting *Bacteroidaceae* and depleted in those infecting

34   *Ruminococcaeae*. Differential abundance analysis identified eighteen viral clusters

35   (VCs) as significantly associated with either MetS or healthy phageomes. Among these

36   are a MetS-associated *Roseburia* VC that is related to healthy control-associated

37   *Faecalibacterium* and *Oscillibacter* VCs. Further analysis of these VCs revealed the

38   *Candidatus Heliusviridae*, a highly widespread gut phage lineage found in 90+% of the

39   participants. The identification of the temperate *Ca. Heliusviridae* provides a novel

40   starting point to a better understanding of the effect that phages have on their bacterial

41   hosts and the role that this plays in MetS.

42

43   **Introduction**

44

45   The human gut microbiome influences many (metabolic) processes, including

46   digestion, the immune system[1], and endocrine functions[2]. It is also involved in

47   diseases such as type 2 diabetes[3], fatty liver disease[4] and inflammatory bowel

48   disease[5]. Though studies of these gut microbiome effects on health and disease mostly

49   focus on bacteria, increasing attention is devoted to bacteriophages (or phages).

50          Phages are viruses that infect bacteria. By infecting bacteria, they can

51   significantly alter gut bacterial communities, mainly by integrating into bacterial

52   genomes as prophages (lysogeny) or killing bacteria (lysis). Such alterations to

53   bacterial communities in turn affect the interactions between bacteria and host, making

54   phages part of an interactive network with bacteria and hosts. For example, an

55   increase in phage lytic action is linked to decreased bacterial diversity in inflammatory

56   bowel disease[6,7], prophage integration into *Bacteroides vulgatus* modifies bacterial bile

57   acid metabolism[8], and dietary fructose intake prompts prophages to lyse their bacterial

58   hosts[9].

59          Gut virome alterations have been linked to several disease states like

60   inflammatory bowel diseases[6,7], malnutrition[10], and type 2 diabetes[11]. But many such

61   studies have not been able to identify specific viral lineages that are involved in such

62   diseases, mainly due to the lack of viral marker genes[12,13] and high phage diversity

63   due to the rapid evolution[14]. Consequently, human gut phage studies are limited to

64   relatively low taxonomic levels. While recent efforts uncovered viral families that are

65   widespread in human populations, such as the crAss-like[15,16] phages, these have not

66   been successfully linked to disease states. In order to develop microbiome-targeted

67 interventions to benefit human health, it is pivotal to study such higher-level phage
68 taxonomies in the gut among relevant cohorts.

69      Here, we report on gut phageome alterations in metabolic syndrome (MetS)[17]
70 among 196 people. MetS is a collection of clinical manifestations that affects about a
71 quarter of the world population, and is a major global health concern because it can
72 progress into cardiometabolic diseases like type 2 diabetes, cardiovascular disease,
73 and non-alcoholic fatty liver disease[18,19]. As gut bacteria are increasingly seen as
74 contributing agents of MetS[20–22], it stands to reason that the phages which infect these
75 bacteria exhibit altered population compositions in MetS. For our analysis, we focused
76 on dsDNA phages, which form a large majority of gut phages in particular and gut
77 viruses in general[14,23]. We found MetS-connected decreases in phageome richness
78 and diversity, which are correlated to bacterial population patterns. Such correlations
79 extended to the relative abundance of phages and their particular bacterial hosts. We
80 further identified eighteen viral clusters (VCs) that are significantly correlated with
81 either MetS or healthy controls. We found that sequences contained in three of these
82 VCs, one VC correlated with MetS and two VCs with controls, are related. These
83 contain members of the *Candidatus Heliusviridae*, a previously unstudied lineage of
84 temperate *Clostridiales* phages that is present in over 90% of the participants.
85 Phylogenetic and taxonomic classification revealed at least six distinct *Ca.*
86 *Heliusviridae* sub-groups, two of which are significantly more abundant in MetS. As the
87 *Ca. Heliusviridae* include both phages which are associated with MetS and with healthy
88 controls, this extremely widespread lineage is an interesting target for research into
89 the human gut phageome.

90

91 **<u>Results</u>**

92

93 <u>Metagenomic sequencing identifies high divergence in MetS phageomes</u>

94

95 We performed bulk metagenome sequencing on fecal samples of 97 MetS and 99
96 healthy participants from the Healthy Life in an Urban Setting (HELIUS) cohort[24], a
97 large population study in Amsterdam, the Netherlands (**Supplementary Table 1**). This
98 yielded an average of 42.1 ± 6.7 million reads per sample (median: 40.6 million reads).
99 We assembled reads and selected 6,780,412 contigs longer than 1,500 bp or shorter
100 but circular, among which we predicted phage sequences that we clustered at 95%

101  nucleotide identity. This produced a database of 25,893 non-redundant phage contigs,

102  which we grouped by shared protein content[25] into 2,866 viral clusters (VCs). These

103  comprised 14,433 contigs, while the remainder were singletons too distinct to

104  confidently cluster with other phage contigs. Treating such singletons as VCs with one

105  member gave a final dataset of 14,325 VCs.

106  Analysis of the read depth per VC across participants (**Supplementary Table**

107  **2**) underscored the extremely high inter-individual diversity in gut phageomes, as 8,799

108  VCs (61.4% of total VCs) were specific to a single individual and 5,122 VCs (35.8% of

109  total VCs) were found in two to twenty participants (*i.e.*, fewer than 10% of participants,

110  **Supplementary Figure 1a**). Due to being so common, these two sets of VCs also

111  comprised a mean of 92.6 ± 4.4% (median: 93.5%) of phage relative abundance

112  (**Supplementary Figure 1b**). The remaining 241 VCs (1.7%) were present in over 10%

113  of participants and represented 7.4 ± 4.4% (median: 6.6%) of phage relative

114  abundance. Of these, 27 (0.2%) were found in over 30% of participants and may be

115  part of the core human gut phageome[26].

116  Next, we examined the relative abundance the of four VC sets (*i.e.*, individual-

117  specific, present in ≤10, 10-30% and ≥30% of participants) in individual participants.

118  For all four sets both the participant with the highest and lowest relative abundance of

119  that VC set had MetS (**Supplementary Figure 1c**). This suggested greater β-diversity

120  variation among MetS viromes than those in healthy controls, which we confirmed with

121  a permutational analysis of multivariate dispersions (permdisp) on Bray-Curtis

122  dissimilarities ($p = 0.003$, **Supplementary Figure 1d**). In conclusion, while we found

123  high inter-individual diversity among phageomes in all participants, MetS phageomes

124  exhibited greater β-diversity variation than controls.

125

126  <u>Gut phage and bacterial populations exhibit altered richness and diversity measures</u>

127  <u>in MetS</u>

128

129  To gain a deeper understanding of MetS phageome community dynamics, we first

130  examined total read fractions that mapped to VCs. This was significantly lower in MetS

131  compared to controls, implying that MetS participants either had lower phage loads or

132  had higher absolute bacterial abundance than controls, or both (Wilcoxon signed rank

133  test, $p = 0.013$, **Supplementary Figure 2a**). This pattern extended to prophage-

134  carrying bacterial contigs, which likewise had lowered relative abundance among MetS

135   participant than controls (Wilcoxon signed rank test, $p$ = 1.4 x 10$^{-4}$, **Supplementary**

136   **Figure 2b**). This notably reflects a decrease in relative abundance of prophage-

137   containing bacteria, not a decrease in that of temperate phages, as the relative

138   abundances of VCs that encode the integrases used in prophage formation were

139   unaltered (Wilcoxon signed rank test, $p$ = 0.47, **Supplementary Figure 2c**). We

140   hypothesize that prophage formation rates among MetS phageomes decrease and that

141   phages possibly more commonly utilize the lytic phage life cycle. Furthermore,

142   combined relative abundance of temperate VCs was a mean 34.8 ± 11.3% total phage

143   relative abundance (median: 32.7%). Thus, while our bulk sequencing approach might

144   be expected to bias in favor of prophages, the majority of phage relative abundance

145   was composed of non-temperate phages.

146        For further analysis of phage communities, we examined phageome richness

147   and diversity. We determined phage richness by measuring the number of VCs that

148   were present (*i.e.*, had a relative abundance above 0) in each participant, using a

149   horizontal coverage cutoff of 75%[27]. This showed that besides lower total phage

150   relative abundance, MetS phageomes also had lower phage richness than controls,

151   but equal evenness (Wilcoxon signed rank test, richness $p$ = 8.7 x 10$^{-8}$, Pielou

152   evenness $p$ = 0.79, **Figure 1a** and **b**). Nevertheless, due to the strong differences in

153   species richness, phage α-diversity was significantly decreased among MetS

154   participants (Shannon H' $p$ = 1.3 x 10$^{-3}$, **Figure 1c**). This suggested that MetS gut

155   phageomes are distinct from healthy communities. Indeed, MetS and control

156   phageomes displayed significant separation when assessed by principal covariate

157   analyses (PCoA) of β-diversity based on Bray-Curtis dissimilarities (permanova $p$ =

158   0.001, **Figure 2d**).

159        Because phages are obligate parasites of bacteria, we also studied bacterial

160   16s rRNA amplicon sequencing data. This showed that MetS phageomes mirror

161   bacterial communities in species richness and α-diversity, but not evenness, which

162   was significantly lowered in MetS bacterial populations (Wilcoxon signed rank test,

163   Chao1 richness $p$ = 9.1 x 10$^{-4}$, Shannon $H'$ $p$ = 1.5 x 10$^{-15}$, Pielou evenness $p$ =1.8 x

164   10$^{-14}$, **Supplementary Figure 3a-c**). Additionally, bacterial communities separated in

165   PCoA analysis in similar fashion to phageomes (permanova $p$ = 0.001 for both bacteria

166   and phages, **Supplementary Figure 3d**). Population-level phageome changes in

167   MetS are thus directly related to a depletion of host bacteria populations, an assertion

168   strengthened by significant direct correlations between phage and bacterial

169  communities in richness (Spearman ρ = 0.42, $p$ = 1.3 x 10$^{-9}$, **Figure 2a**), evenness

170  (Spearman ρ = 0.24, $p$ = 5.7 x 10$^{-4}$, **Figure 2b**).

171       As the bacterial and phage populations did not equally decrease in richness and

172  evenness, they also did not equally correlate with MetS clinical parameters. Rather,

173  phage richness was significantly negatively correlated with obesity, blood glucose

174  levels, blood pressure, and triglyceride concentrations but bacterial richness was not

175  ($p$ < 0.05, **Figure 2c**). Bacterial evenness, meanwhile, did significantly negatively

176  correlate with these clinical parameters while phage evenness did not ($p$ < 0.05, **Figure**

177  **2d**). Increasingly severe MetS phenotypes thus result in stronger decreases in

178  bacterial evenness than richness, while phage populations exhibit stronger decreases

179  in richness than evenness. The decreasing bacterial evenness could be caused by

180  depletion of certain bacterial species in MetS, which results in the viruses infecting

181  these depleted bacteria to become undetectable, thereby decreasing richness more

182  than evenness. Otherwise, the success of certain bacterial species could also

183  decrease evenness. In the process this could obfuscate rare phage species, which

184  could cause the decreased phage richness. Combined with the results showing MetS-

185  associated reduction in total phage abundance but no increase in lysogeny

186  (**Supplementary Figure 2**), our findings indicate that certain phages are either

187  completely removed from the gut or become too rare to detect in MetS.

188

189  Phages infecting select bacterial families are more abundant in MetS phageomes

190

191  We next studied individual bacterial lineages and the phages that infect them. To do

192  this, we linked viral contigs to bacterial hosts by determining CRISPR protospacer

193  alignments, taxonomies of prophage-containing bacterial sequences, and hosts of

194  previously isolated phages co-clustered in VCs  (see methods for details). We found

195  2,621 host predictions between 2,575 VCs (18% of all VCs) and eleven bacterial phyla,

196  most commonly *Firmicutes* (2,067 VCs) and *Bacteroidetes* (234 VCs, **Supplementary**

197  **Table 3**). We also identified 43 VCs with multi-phyla host range predictions, similar to

198  previous works[28].

199       To increase statistical accuracy, we selected 1,744 host predictions between

200  1,514 VCs (10.6% of all VCs) and the twelve most commonly occurring host families.

201  We then performed an analysis of compositions of microbiomes with bias correction

202  (ANCOM-BC)[29], which showed higher relative abundances in controls for

203  *Ruminococcaceae* VCs ($q = 9.1 \times 10^{-3}$), and in MetS for *Bacteroidaceae* VCs ($q = 2.2$

204  $\times 10^{-4}$), plus marginally for *Acidaminococcaceae* and *Tannerellaceae* VCs ($q = 0.04$,

205  **Figure 3a**). ANCOM-BC on 16s rRNA gene sequencing data found that multiple

206  *Ruminococcaeae* ASVs were significantly differentially abundant in controls

207  (**Supplementary Figure 5**). Of *Ruminococcaceae* VCs with host predictions at the

208  species level, those linked to *Faecalibacterium sp. CAG:74*, *Ruthenibacterium*

209  *lactatiformans,* and *Subdoligranulum sp. APC924/74* had higher relative abundance in

210  controls (Wilcoxon signed rank test with Benjamini and Hochberg adjustment, $q \le 0.05$,

211  **Figure 3b**). These results are congruent with *Ruminococcaeae* being commonly linked

212  to healthy gut microbiomes[3,30,31].

213      ANCOM-BC on 16s rRNA sequencing data identified *Bacteroidaceae* bacteria

214  as significantly differentially abundant among MetS participants ($q = 1.03 \times 10^{-13}$).

215  Since some widespread crAss-like gut phages infect *Bacteroidaceae* hosts[32,33], we

216  ascertained whether this phage family was more abundant among MetS participants.

217  We did not find significant relations between crAss-like phage VC relative abundance

218  and MetS (**Supplementary Figure 4a**), although VCs containing such phages were

219  present more often among control (70/99) than MetS (57/97) participants (Fisher's

220  exact test, $p = 0.1$, **Supplementary Figure 4b**). Next to being absent more often, the

221  participant with the highest relative abundance of crAss-like phage VCs belonged to

222  the MetS group (17.2% of total phage relative abundance, **Supplementary Figure 4a**),

223  which was indicative of greater variation in crAss-like phage relative abundance among

224  MetS (mean 1.29±2.62%) than controls (mean 0.830±1.44%). MetS-associated

225  alterations to crAss-like phage composition may thus occur at an individual level.

226

227  <u>*Bacteroidaceae* VCs are markers of the MetS phageome</u>

228

229  The above results all indicate that MetS gut phageomes are distinct from those in

230  healthy individuals. In light of this, we surveyed our cohort for individual VCs that were

231  correlated with either MetS or healthy control phageomes. ANCOM-BC uncovered

232  thirty-nine VCs that were more abundant in MetS participants, and eight more in

233  controls ($q \le 0.05$, **Figure 4a**).

234      In line with the above findings that *Bacteroidaceae* VCs are hallmarks of the

235  MetS phageome, three MetS-associated VCs infected *Bacteroides* bacteria. The first

236  (VC_1180_0) contained a non-prophage contig (*i.e.*, no detected bacterial

237 contamination) of 34,170 bp with a checkV[34] completion score of 100%. It further co-
238 clustered with a contig that checkV identified as a complete prophage flanked by
239 bacterial genes. Analysis with the contig annotation tool (CAT[35]) identified this contig
240 as *Bacteroides fragilis*. Additionally, the most complete VC_1180_0 contig shared 6/69
241 (8.7%) ORFs with *Bacteroides uniformis Siphoviridae* phage Bacuni_F1[36] (BLASTp
242 bit-score ≥ 50). The second MetS-associated *Bacteroides* VC (VC_786_0) contained
243 one contig with CRISPR spacer hits to *Bacteroides*. Its most complete contig had a
244 CheckV completeness score of 98.94% and was classified by the contig annotation
245 tool (CAT[35]) as *Phocaeicola vulgatus* (formerly *Bacteroides vulgatus*[37]). This near-
246 complete contig furthermore shares 11/77 ORFs (14.3%) with *Riemerella Siphoviridae*
247 phage RAP44 (BLASTp bit scores >50). This last finding was notable because the third
248 and final MetS-associated *Bacteroides* VC (VC_775_1) also contained a near-
249 complete genome (CheckV: 90.32% complete) that shared 16/81 ORFs (19.8%) with
250 RAP44. Comparison of the most complete VC_786_0 and VC_775_1 contigs indicated
251 that they share nine ORFs, revealing that they are part of an extended family of
252 *Bacteroidetes Siphoviridae* phages of which members are hallmarks of MetS.

253

254 <u>A widespread phage family contains markers for healthy and MetS phageomes</u>

255

256 Besides the above-mentioned *Bacteroidaceae* VCs, all other differentially abundant
257 VCs with host links, two MetS- and four control-associated, infected *Firmicutes*,
258 particularly in the *Clostridiales* order. The sole exception to this (VC_745_0)
259 remarkably had CRISPR protospacer matches to *Firmicutes* bacteria *Faecalibacterium*
260 *sp.* CAG: 74_58_120 and *Ruthenibacterium lactatiformans,* as well as to
261 *Actinobacteria* bacterium *Parascardovia denticolens*. As this VC included genome
262 fragments with simultaneous CRISPR protospacer hits to both phyla, VC_745_0
263 seemingly contains phages with an extraordinarily broad host range.

264 Besides this broad host range VC, our attention was drawn to the two MetS-
265 associated *Clostridiales* VCs. Both were predicted to infect hosts that are usually
266 associated with healthy gut microbiomes: *Roseburia*[3] for VC_659_0, and
267 *Oscillospiraceae*[38] and *Faecalibacterium prausnitzii*[39] for VC_1040_0. Further
268 examination of their largest genomes revealed that MetS-associated VC_659_0 was
269 remarkably similar to two VCs that were significantly associated with healthy controls:

270 the above-mentioned broad host-range VC_745_0 and the less broad host-range

271 *Oscillibacter/Ruminococcaceae* VC_643_0 (**Figure 4b**).

272 Intrigued by this apparent relatedness of VCs that included markers of MetS and

273 healthy controls among our cohort, we sought to identify additional related sequences

274 among our cohort. For this we first determined the exact length of VC_659_0 genomes

275 by analyzing read coverage plots of a prophage flanked by bacterial genes (**Figure**

276 **4b**). By analyzing coverage of the contig in subjects where bacterial genes were highly

277 abundant but viral genes were absent, we extracted a genome of 68,665 bp long.

278 Homology searches of all 74 ORFs encoded by this prophage against all ORFs from

279 all phage contigs in the cohort identified 249 contigs of over 30,000 bp that all shared

280 nine genes (BLASTp bit score ≥ 50, **Figure 4b**). Additionally, we identified 51

281 *Siphoviridae* phage genomes in the National Center for Biotechnology Information

282 (NCBI) nucleotide database that also shared these nine genes. With one exception,

283 these were *Streptococcus* phages, the exception being *Erysipelothrix* phage phi1605.

284 The genes shared by all these phage genomes formed three categories. First

285 are genes encoding structural functions: a major capsid protein, portal protein, CLP-

286 like prohead maturation protease, and terminase. The second group are transcription-

287 related genes encoding a DNA polymerase I, probable helicase, and nuclease. Finally,

288 there are two genes that encode domains of unknown function, but which given their

289 adjacency to the second group are likely transcription-related.

290 Earlier studies have used a cutoff of 10% gene similarity for phages that are in

291 the same families, 20% for sub-families, and 40% for genera[40,41]. The nine shared

292 genes form 10-25% of ORFs found on both the characterized phages and non-provirus

293 contigs with checkV 'high-quality' designations. Thus, these phages form a family,

294 which we dubbed the *Candidatus Heliusviridae*. Next, we further studied the *Ca.*

295 *Heliusviridae* interrelatedness by calculating the pairwise percentages of shared

296 protein clusters and hierarchically clustering the results (**Figure 5a**. This showed that

297 the *Ca. Heliusviridae* form six groups, which we designated as *Ca. Heliusviridae* group

298 alpha, beta, gamma, delta, epsilon, and zeta *Heliusvirinae*. A concatenated

299 phylogenetic tree made from alignments of nine conserved *Ca. Heliusviridae* genes

300 largely confirmed the hierarchical clustering (**Supplementary Figure 6**).

301 The *Ca. Heliusviridae* group alpha solely contained previously isolated

302 *Streptococcus* phages, which both in the hierarchical clustering and the phylogenetic

303 tree were distinct from the other genomes. Meanwhile, all three VCs that were

304    significantly associated with either MetS or healthy controls where part of the *Ca.*
305    *Heliusviridae* group zeta, by far the largest and most diverse group. Two of these,
306    VC_659_0 and VC_745_0, formed distinct sub-clades in both hierarchical clustering
307    and phylogenetic tree, while VC_643_0 conversely was spread out over multiple
308    clades.

309         The *Ca. Heliusviridae* were present in 181 participants (92.3%), 94 controls and
310    87 MetS participants (**Figure 5b**). We also tested this finding in two cohorts in which
311    the gut phageome was studied earlier, in the context of hypertension[42] and type 2
312    diabetes[11]. To allow for incomplete assemblies, we searched for contigs in these two
313    cohorts that contain the four conserved *Ca. Heliusviridae* structural genes. A
314    phylogenetic tree containing concatenated alignments of the structural genes clearly
315    showed that both validation cohorts contained sequences from all *Ca. Heliusviridae*
316    groups (**Supplementary Figure 7**). Only a small minority of 47 sequences, largely
317    from the hypertension cohort, formed a separate and distant clade of which the relation
318    to the remainder of *Ca. Heliusviridae* is unclear. Among the two cohorts, *Ca.*
319    *Heliusviridae* were present in 140/196 (71.4%, hypertension) and 112/145 (77.2%,
320    T2D) participants. Finally, as this study and the two validation cohorts all utilized whole
321    genome shotgun sequencing, the phages identified here might be inactive prophages.
322    Thus, we used datasets of fecal virus-like particle (VLP) sequencing from ten people
323    that were published earlier[43]. Cross-assembly of the ten VLP sequence datasets
324    identified one contig of 43,244 bp (70.68% checkV completeness) and eight contig
325    fragments that contained four or more conserved *Ca. Heliusviridae* genes. Thus,
326    phages in this family are also found in VLP fractions, implying that they are inducible.

327         Of the *Heliusviridae* groups, the zeta was by far the most abundant, being
328    present in 88 controls and 72 MetS participants. This meant healthy control
329    phageomes were significantly more likely to contain *Heliusviridae* group zeta (Fisher
330    exact test, $p = 0.0096$), though they were not significantly more abundant. Of the other
331    candidate sub-families, the groups delta and epsilon were in significantly higher
332    relative abundance (Wilcoxon signed rank test, $p = 0.0043$ and $0.0063$, respectively)
333    among MetS participants. The *Ca. Heliusviridae* group delta infects *Lachnospiraceae*,
334    in particular *Butyrivibrio* sp. CAG:318 and *Lachnoclostridium* sp. An181. Meanwhile,
335    the *Heliusviridae* group epsilon were distinct among the *Heliusviridae* in that they infect
336    *Negativicutes* rather than *Clostridia*, specifically *Acidominacoccus* and several other
337    genera in the *Veillonellaceae* (**Supplementary Table 3**). These results, combined with

338    the fact that group zeta VC_659_0 is strongly correlated with MetS (**Figure 4a**), show

339    that Ca. *Heliusviridae* are part of the core human gut phageome, where they may affect

340    intricate relations with human health.

341

342    MetS-associated group zeta *Ca. Heliusviridae* prophages encode possible metabolic

343    genes

344

345    The *Ca. Heliusviridae* are generally linked to bacteria that are associated with healthy

346    human gut microbiomes. Therefore, it is an apparent contradiction that MetS-

347    associated *Ca. Heliusviridae* group zeta VC_659_0 infected *Roseburia*, a short chain

348    fatty acids producer that is often abundant in healthy microbiomes[44]. Due to this

349    contradiction, we explored this VC further. It contained two additional prophages, which

350    where both incomplete (**Figure 6a**). Whole genome alignment showed that all three

351    prophages shared their phage genes, and that the two incomplete ones also shared

352    host-derived genes. This indicated that the incomplete prophages integrated into highly

353    similar host bacteria which were distinct from *Roseburia*. To confirm this, we performed

354    homology searches of the bacterial host ORFs found on these contigs against the

355    NCBI nr database (BLASTp, bit-score ≥50). In both incomplete prophages, the majority

356    of ORFs had *Blautia* as their top hit, which for a plurality of ORFs involved *Blautia*

357    *wexlerae* (**Figure 6a**)*.* VC_659_0 thus contains MetS-associated phages that integrate

358    into at least two genera (*Roseburia* and *Blautia*) within the *Lachnospiraceae*.

359        To examine if the hosts infected by VC_659_0 were more abundant in MetS

360    participants, we determined mean coverage of bacterial genes found adjacent to the

361    prophages. We thus assured that we analyzed the particular host strains infected by

362    these phages, rather than unrelated strains in the same genera. This showed that both

363    the *Blautia* and the *Roseburia* host genes were more abundant among MetS

364    participants (Wilcoxon signed rank test, *Blautia* $p = 5.1 \times 10^{-4}$, *Roseburia* $p = 0.042$,

365    **Figure 6b** and **c**). The specific *Lachnospiraceae* strains infected by VC_659_0 phages

366    thus seem to thrive in MetS microbiomes. This could in part be due to functions

367    conferred upon these bacteria by these prophages, as particularly the *Roseburia*

368    prophage which carried several virulence and metabolism-related genes, including

369    ones    encoding    a    chloramphenicol    acetyltransferase    3    (2.3.1.28),

370    Glyoxalase/Bleomycin resistance protein (IPR004360), multi antimicrobial extrusion

371    protein    (IPR002528),    2-succinyl-6-hydroxy-2,4-cyclohexadiene-1-carboxylate

372  synthase (4.2.99.20), and NADPH-dependent FMN reductase (PF03358). The latter

373  two in particular are both associated with vitamin K (menaquinone) metabolism, which

374  is part of (an)aerobic respiration in bacteria[45]. We speculate that this opens up the

375  possibility that this *Roseburia* prophage aids its host bacterium, which in turn may

376  contribute to MetS phenotypes.

377

378  **Discussion**

379

380  This is the first study of the gut phages in the context of MetS, a widespread global

381  health concern to which the gut bacteria targeted by phages are believed to be a main

382  contributor[18]. We have shown that MetS is associated with decreases in gut phageome

383  total relative abundance and richness, but no change in evenness. Due to their

384  compositional nature, these phageome alterations could be bacterially driven, as

385  phage total relative abundance decreases could be caused by bacterial counts

386  increasing rather than phage counts decreasing. But since we measured decreased

387  bacterial richness and evenness, MetS gut metagenomes would need to have larger

388  numbers of bacterial cells that are distributed among fewer strains that are more

389  unevenly divided than in healthy individuals. Conversely, total phage relative

390  abundances could be lower in MetS due to lower viral loads, which would be in line

391  with decreased phage richness and is in agreement with recently reported direct

392  correlations between gut viral and bacterial populations in healthy individuals[46]. Future

393  confirmation of this would necessitate counts of viable bacterial cells and VLP. In either

394  case, we surmise that the main driver of these effects is diet, which affects bacterial[47−]

395  [49] as well as viral[50] populations. It is also possible that phage populations as described

396  here may further exacerbate bacterial diversity losses, as low phage abundance may

397  decrease their positive effects on bacterial diversity[51,52].

398        One aspect in which our study separated itself from some other gut virome

399  studies was its usage of whole genome rather than VLP sequencing. We believe our

400  approach has its advantages, such as a lack of pre-sequencing amplification and the

401  biases this introduces, and a greater emphasis on the prophage community significant

402  to functioning of both bacterial and phage communities in the gut[8,9,12,53]. It must

403  nevertheless be noted that our sequencing approach may underestimate the virulent

404  proportion of the human gut phageome, likewise to VLP sequencing approaches

405  probably underestimating gut prophage communities. Indeed, analysis of both VLP

406 and bulk communities is likely needed to gain a full reckoning of human gut
407 phageomes[14]. This approach could furthermore distinguish between inducible and
408 defective prophages[54], which we were unable to do.

409 Like other studies, we found highly individual-specific[26,43] human gut
410 phageomes in which singular widespread phage VCs are rare[14,15]. Despite the
411 universally high inter-individual phageome diversity, we found larger within-group
412 variation among MetS phageomes than healthy controls. This is consistent with the
413 Anna Karenina principle (AKP), which holds that "all healthy microbiomes are similar;
414 each dysbiotic microbiome is dysbiotic in its own way"[55]. Such AKP-dynamics mirror
415 previous findings of obesity-related alterations in the gut bacterial populations[56].
416 Particularly, we hypothesize that stressors inherent to MetS perturb gut bacterial
417 populations in a stochastic fashion, the effects of which reverberate to the phage
418 populations and result in their increased within-group variation.

419 We further found strong negative correlations between the risk factors that
420 constitute MetS and phage richness but not evenness. This likely results from the
421 whole-genome sequencing approach that we took, which better captures intracellular
422 phages (*e.g.*, actively replicating or integrated prophages) than extracellular phages[14].
423 The phage VC richness that we report here thus represents phages that are actively
424 engaging with their hosts or are highly abundant extracellularly. As phages that target
425 depleted bacteria are more likely to be low in abundance and extracellular, our
426 approach does not capture them. Thus, the apparent species richness drops because
427 low abundant extracellular phages are below the detection limit of our sequencing
428 approach. This removal of rare phages in turn prohibits significant drops in species
429 evenness in MetS. It could also be that bacteria depleted in MetS reside in phage-
430 inaccessible locales within the gut[57], which perhaps results in removal of the
431 corresponding phages from the gut to below detectable levels. This would explain the
432 stronger correlation between bacterial evenness than richness to MetS risk factors.

433 As most (gut) phages remain unstudied[14,58], it is often difficult to link phages to
434 host bacteria[59]. Here, we linked roughly one tenth of all VCs to a bacterial host. The
435 remaining majority of VCs likely represent phages that infect bacterial lineages lacking
436 CRISPR systems[60], are exclusively lytic, or that integrate into hosts which we could
437 not taxonomically classify. Whichever is the case, our study underscores the great
438 need for methods that link phages to hosts with high accuracy[61,62]. From the phage-
439 host linkages that we obtained, we found that VCs containing phages infecting specific

440    bacterial families tend to be either depleted (*Bifidobacteriaceae, Ruminococcaeae*, and

441    *Oscillospiraceae)* or enriched (*Bacteroidaceae*) in tandem to their hosts. For the

442    *Bifidobacteriaceae* and *Oscillospiracaeae*, this is in line with established studies that

443    show depletion of these families in MetS[3] and MetS-associated disease states[30,38].

444          For *Ruminococcaceae* bacteria, associations with MetS and MetS-related

445    diseases are less clear, with reports of both positive[3,31] or negative[39,63] associations.

446    The specific *Ruminococcaceae* of which we link the phage-containing VCs to healthy

447    controls,    most    notably    *Faecalibacterium    prausnitzii*    and    *Ruthenibacterium*

448    *lactatiformans*, are both considered hallmarks of healthy human gut microbiomes[3,30,39].

449    Interestingly,    we    also    succeeded    in    linking    specific    viral    clusters    that    infect

450    *Faecalibacterium* to both healthy controls (VC_745_0) and MetS participant groups

451    (VC_1040_0). This contradiction may indicate that infections of *Faecalibacteria* could

452    result in differing outcomes for the bacterium depending on the phage species. As both

453    VCs contain sequences with integrase-like genes, they contain temperate phages. It

454    could be that VC_745_0 prophages augment *Faecalibacterium* growth, as prophages

455    are known to do[64]. Meanwhile VC_1040_0 prophages could be detrimental to the same

456    bacteria, for example by becoming lytic in the presence of MetS-associated dietary

457    components, likewise to *Lactobacillus reuteri* prophages lysing their hosts in the

458    presence of dietary fructose[9]. Such behavior can lead to the collapse of the bacterial

459    population[44], and may thereby be a contributing factor to depletion of *Faecalibacterium*

460    in MetS.

461          As mentioned, the *Bacteroidaceae* were the only bacterial family that are

462    infected by phages of which the VCs were significantly more abundant among MetS

463    participants. Concordantly, we found several individual *Bacteroides* VCs that were

464    MetS-associated. The *Bacteroides* are often positively associated with high-fat and

465    high-protein diets[65,66]. Simultaneously, however, reports disagree on individual

466    *Bacteroides* species and their associations with MetS-related diseases like obesity,

467    type 2 diabetes, and non-alcoholic fatty liver disease[30]. Such conflicting reports likely

468    reflect the large diversity in metabolic effects at strain level among these bacteria[67].

469    Based on our results, we drew two conclusions. First, that *Bacteroidaceae*-linked VCs

470    mirror their hosts in MetS-associated relative abundance increase, and second that

471    *Bacteroidaceae*-linked VCs are of significant interest to studies of the MetS

472    microbiome. The latter conclusion is strengthened by findings that *Bacteroides*

473    prophages can alter bacterial metabolism in the gut[8].

474    While *Bacteroidaceae* VCs at large were thus seemingly associated with MetS

475    phenotypes, we uncovered larger variation of crAss-like phage-containing VC

476    abundance, which suggest at individual-specific alterations to this gut phage family

477    among MetS phageomes. This widespread and often abundant human gut phage

478    family infects *Bacteroidetes*, including members of the *Bacteroidaceae*[68,69]. As these

479    phages are commonly linked to healthy gut microbiomes[41,69,70], it is conceivable that

480    they would be negatively correlated with MetS phageomes. That this is not the case

481    among the entire cohort is likely due to the great variety within this family[69], and

482    perhaps additionally to the hypothesized aptitude of crAss-like phages for host

483    switching through genomic recombination[69].

484

485    Finally, our study revealed the *Candidatus Heliusviridae*, a highly widespread

486    family of gut phages that largely infect *Clostridiales* hosts. This prospective family is

487    also expansive, and includes at least six distinct candidate subfamilies. Our uncovering

488    of this novel human gut phage family underscores the usefulness of database-

489    independent *de novo* sequence analyses[25,27,71], as well as the need for a wider view

490    on viral taxonomy than has presently been exhibited in the field of gut viromics.

491    The *Ca. Heliusviridae* are of particular interest to studies of MetS and related

492    illnesses because its member phages include some associated with MetS and others

493    with healthy controls. Most striking is the fact that most of the bacteria infected by

494    MetS-associated *Ca. Heliusviridae,* are generally producers of short chain fatty acids

495    (SCFA) such as butyrate and commonly depleted in MetS[30]. Such SCFA-producing

496    bacteria are commonly positively associated with healthy microbiomes, as SCFAs that

497    result from microbial digestion of dietary fibers have a role in the regulation of

498    satiation[72,73]. The exception to this are the *Veillonellaceae* that are infected by the

499    *Heliusviridae* group epsilon, which are found at elevated abundance in non-alcoholic

500    fatty liver disease[30]. While higher abundance of some of the other butyrate-producers

501    infected by *Ca. Heliusviridae* is associated with metformin use[74], this is used to treat

502    type 2 diabetes rather than MetS.

503    Particularly interesting are the *Roseburia/Blautia* phages in VC_659_0, which

504    was the most strongly correlated with MetS out of all VCs. The positive correlation

505    between the relative abundance of these phages and that of their hosts indicates that

506    they have a stable relation with their hosts in the MetS microbiome. This is to be

507    expected, as large-scale prophage induction is generally associated with sudden

508 alterations to the microbiome, such as the addition of a specific food supplement that
509 acts as an inducer of prophages[9]. Such sudden alterations in phage behavior are
510 unlikely to be captured in large cohorts with single measurements. In fact, as phages
511 are strongly dependent on their host, one might expect the abundance of many gut
512 phages to be positively correlated to that of their particular hosts under the relatively
513 temporally stable conditions of MetS. The strong correlation of VC_659_0 to MetS
514 phenotypes, coupled to the commonly found correlation to healthy microbiomes of
515 VC_659_0 host bacteria, and the presence of potential auxiliary metabolic genes in
516 VC_659_0 phage sequences combined introduce the possibility that prophage
517 formation of these *Ca. Heliusviridae* phages alters the metabolic behavior of their host
518 bacteria, as is known to happen in marine environments[75,76]. This could make these
519 bacteria detrimental to health. Proving this hypothesis necessitates future isolation of
520 VC_659_0 phages.

521 Despite efforts to catalog the human gut phageome[14,28], taxonomically higher
522 structures are still largely absent. This study shows the worth of analyzing phages at
523 higher taxonomic levels than genomes or VCs, similarly to what has been shown in
524 recent years regarding the crAss-like phage family[15,16]. Unlike the crAss-like phage
525 family, however, the *Ca. Heliusviridae* seem to be strongly correlated with human
526 health. We hope that further research will provide a deeper understanding of the effect
527 that these phages have on their bacterial hosts and the role that this plays in MetS, as
528 well as a refinement of their taxonomy.

529

542 (Hartstichting; 2010T084), the Netherlands Organization for Health Research and
543 Development (ZonMw; 200500003), the European Integration Fund (EIF;
544 2013EIF013), and the European Union (Seventh Framework Programme, FP-7;
545 278901). We gratefully acknowledge the AMC Biobank for their support in biobank
546 management and high-quality storage of collected samples. We are most grateful to
547 the participants of the HELIUS study and the management team, research nurses,
548 interviewers, research assistants and other staff who have taken part in gathering the
549 data of this study.

550

551 **Author contributions**
552 PAdJ and KW conducted data analysis; TPMS, BJvdB, AHZ, FLN, BED, and MN
553 assisted with experimental design and data interpretation; PAdJ and HH designed the
554 study and wrote the manuscript. All authors read and provided input on the manuscript.

555

556 **Declaration of Interests**
557 MN owns stock in, consults for, and has intellectual property rights in Caelus Health.
558 He consults for Kaleido. None of these are directly relevant to the current paper.

559

560 **Methods**

561

562 Sequencing and contig assembly

563

564 The Healthy Life in an Urban Setting (HELIUS) cohort includes some 25,000 ethnically
565 diverse participants from Amsterdam, the Netherlands. The cohort details were
566 published previously[77]. The HELIUS cohort conformed to all relevant ethical
567 considerations. It complied with the Declaration of Helsinki (6th, 7th revisions), and was
568 approved by the Amsterdam University Medical Centers Medical Ethics Committee.
569 For details on stool sample collection from among the participants, their storage, and
570 DNA extraction, see Deschasaux, *et al*[24]. In summary, participants were asked to
571 deliver stool samples to the research location within 6 hours after collection with pre-
572 provided kit consisting of a stool collection tube and safety bag. If not possible, they
573 were instructed to store their sample in a freezer overnight. Samples were stored at
574 the study visit location at -20°C until daily transportation to a central -80°C freezer.

575 Total genomic DNA was extracted using a repeated bead beating method described
576 previously[24,78]. Libraries for shotgun metagenomic sequencing were prepared using a
577 PCR-free method at Novogene (Nanjing, China) on a HiSeq instrument (Illumina Inc.
578 San Diego, CA, USA) with 150 bp paired-end reads and 6 Gb data/sample. All
579 bioinformatics software was run using standard settings, unless otherwise stated. All
580 sequencing reads are available at the European Nucleotide Archive under project
581 PRJEB42542. Samples have accession numbers ERS5585222-ERS5585321, all
582 phage contigs are under accession number ERZ1762427.

583 Following previously set definitions[79], participants were classified in the
584 MetS group if three of the following five health issues occurred: abdominal obesity
585 measured by waist circumference, insulin resistance measured by elevated fasting
586 blood glucose, hypertriglyceridemia, low serum high-density lipoprotein (HDL), and
587 high blood pressure[79]. All participants of the HELIUS cohort reside in Amsterdam, the
588 Netherlands. Participants were roughly evenly divided by ethnicity, with European
589 Dutch comprising 49 controls and 49 MetS participants, and African Surinamese 50
590 controls and 49 MetS participants. The MetS group contained 55 women and had a
591 median age of 58 (mean 56.8±8.09), and the controls 71 and had a median age of 50
592 (mean 49.1±12). Of the 196 participants, 26 used metformin, of whom 2 were controls
593 who did not concur to the MetS criteria. Analysis of sequencing output started with
594 assembly of the sequencing reads per sample (*i.e.,* 196 individual assemblies) into
595 contigs using the metaSPAdes v3.14.1 software[80]. For each sample, we selected
596 contigs of more than 5,000 bp for further analysis. In addition, among contigs between
597 1,500 and 5,000 bp we identified circular contigs by checking for identical terminal ends
598 using a custom R script that employed the Biostrings R package v3.12[81]. All 6,780,412
599 circular contigs and contigs over 5,000 bp were then pooled before phage sequence
600 prediction.

601

602 <u>Phage and bacterial sequence selection</u>

603

604 We predicted phage sequences as described previously[82]. In short, we first analyzed
605 contigs using VirSorter v1.0.6[83] and selected those in category 1, 2, 4, and 5. In
606 parallel, contigs were analyzed using VirFinder v1.1, after which we selected those
607 with a score above 0.9 and a p-value below 0.05. We additionally classified contigs as
608 phage if (I) they were both in VirSorter categories 3 or 6 and had VirFinder scores

609    above 0.7 with p-values below 0.05, and (II) annotation with the contig annotation tool

610    (CAT) v5.1.2[35] was as "Viruses" or "unclassified" at the superkingdom level. After

611    removing those with CAT classifications as Eukaryotic viruses, this resulted in a

612    database of 45,568 phage contigs. Bacterial sequences were predicted by selecting

613    all contigs that CAT annotated in the "Bacteria" at the superkingdom level, and

614    removing contigs that were also found in the phage dataset. An exception was made

615    for prophage contigs in VirSorter category 4, 5, and 6, which were left among the

616    bacterial dataset (see "Phage-host linkage prediction"). This resulted in a total of

617    1,579,361 bacterial contigs. The 1,624,929 bacterial and phage datasets were then

618    concatenated and deduplicated using dedupe from BBTools v38.84 with a minimal

619    identity cutoff of 90% (option minidentity=90). This identified 759,403 duplicates and

620    resulted in 829,633 non-redundant bacterial sequences and 25,893 non-redundant

621    phage sequences. While the bacterial sequences were used for host prediction (see

622    "Phage-host linkage prediction"), we subsequently predicted open reading frames

623    (ORFs) in phage contigs using Prodigal v2.6.2[84] (option -p meta). These ORFs were

624    then used to group phage sequences in viral clusters (VCs) using vContact2 v0.9.18[25].

625    This resulted in 2,866 VCs comprising 14,433 phage contigs and 11,460 singletons

626    and outliers, which we treated as VCs with one member. This resulted in 14,325 VCs.

627    For a full accounting of phage contigs, see **Supplementary Table 2 and 4.**

628

629    <u>Read mapping and community composition</u>

630

631    For bacterial community composition, we used sequencing data targeting the V4

632    region of the 16s rRNA gene that had been performed previously[24,85]. Details on ASV

633    construction from these samples was described previously in Verhaar, et al[85]. As part

634    of this previous analysis, samples with fewer than 5000 read counts had been

635    removed, and samples had been rarified to 14932 counts per sample.

636         To determine phage community composition, we mapped reads from each

637    sample to the non-redundant contig dataset using bowtie2 v2.4.0[86]. As previously

638    recommended[27], we removed spurious read mappings at less than 90% identity using

639    coverM filter v0.5.0 (unpublished; https://github.com/wwood/CoverM, option -min-read-

640    percent-identity 90). The number of reads per contig was calculated using samtools

641    idxstats v1.10[87]. As was also recommended[27], contig coverage was calculated with

642    bedtools genomecov v2.29.2[88], and read counts to contigs with a coverage of less than

643   75% were set to zero. Read counts for each sample were finally summed per VC. All

644   contigs were analyzed for completion with CheckV v 0.7.0-1[34].

645

646   Ecological measures

647

648   In all boxplots, we tested statistical significance using the Wilcoxon rank sum test as it

649   is implemented in the ggpubr v0.4.0 R package (available from: https://cran.r-

650   project.org/web/packages/ggpubr/index.html). Unless stated otherwise, all plots were

651   made using either ggpubr or the ggplot2 v3.3.2 R package (available from:

652   https://cran.r-project.org/web/packages/ggplot2/index.html). Alpha diversity measures

653   (observed VCs and Shannon H' for phages and Chao1 and Shannon H' for bacteria)

654   were calculated using read count tables with the plot_richness function in the phyloseq

655   R package v1.33.0[89]. For β-diversity, we converted read counts to relative abundances

656   using the transform function from the microbiome v1.11.2 R package. We then used

657   the phyloseq package to calculate pairwise Bray-Curtis dissimilarities and construct a

658   principal coordinates analysis (PCoA). Statistical significance of separation in the

659   PCoA analysis was determined with a permutational multivariate analysis of variance

660   (permanova) using the adonis function from the vegan R package[90]. For this analysis,

661   we adjusted for smoking, sex, age, alcohol use, and metformin use. Direct correlation

662   coefficients between richness and diversity were calculated using the stat_cor function

663   in the ggpubr R package.

664

665   Phage-host linkage prediction

666

667   We predicted VC-bacterium links in three ways: (i) CRISPR protospacers, (ii) prophage

668   similarity, and (iii) characterized phage similarity.

669          We predicted CRISPR arrays among the bacterial contigs using CRISPRdetect

670   v2.4[91] (option array_quality_score_cutoff 3) and used these to match bacterial contigs

671   and phage contigs. In addition, we used a dataset of 1,473,418 CRISPR spacers that

672   had previously been predicted[62,92] in genomes contained in the Pathosystems

673   Resource Integration Center (PATRIC)[93] database to match to phage contigs with

674   spacePharer v2-fc5e668[94] using standard settings and cutoffs. This process resulted

675   in 3,727 spacer hits, of which 2,244 hits were either to PATRIC genomes or to bacterial

676 contigs from this study with definite CAT classifications at the phylum level
677 (**Supplementary Table 3**).

678      To identify predicted phage contigs with high sequence similarity to prophages,
679 we analyzed which viral clusters contained on of the 7,691 bacterial contigs with
680 VirSorter prophage predictions in category 4 or 5. CAT was subsequently used to
681 determine the taxonomy of bacterial contigs with prophage regions. In total, we linked
682 1,102 VCs to prophages with this approach.

683      Finally, VCs were linked to bacterial hosts by vContact2 clustering with
684 characterized phages from the viral RefSeq V85 database[95] with a known host. To
685 achieve this, we selected all VCs from the vContact2 output that contained both
686 characterized genomes and phage contigs. If all characterized phages infected hosts
687 within the same bacterial family, we took that to mean that the whole VC infects hosts
688 from that family. This approach linked 44 VCs to hosts.

689

690 Differential abundance analysis

691

692 To determine which bacteria and VCs were differentially abundant between MetS and
693 control subjects, we employed the analysis of composition of microbiomes with bias
694 correction (ANCOM-BC)[29]. This novel method, unlike other similar methods like
695 DeSeq2, takes into account the compositional nature of metagenomics sequencing
696 data[96]. To implement this method, we applied the ANCOM-BC v1.0.2 R package to
697 raw read count tables, as ANCOM-BC employs internal corrections for library size and
698 sampling biases[29]. Significance cutoff was set at an adjusted p-value of 0.05, p-values
699 were adjusted using the Benjamini-Hochberg method, and all entities (bacteria
700 taxa/VCs) that were present in more than 10% of the samples were included (options
701 p_adj_method = "BH", zero_cut = 0.9, lib_cut = 0, struc_zero = T, neg_lb = F, tol = 1e-
702 5, max_iter = 100, alpha = 0.05). For this analysis, we adjusted for smoking, sex, age,
703 alcohol use, and metformin use.

704

705 crAss-like phages

706

707 To identify crAss-like phages, we employed a methodology described earlier[41]. Shortly,
708 a BLAST database was made containing all ORFs from all phage contigs (predicted
709 before viral clustering, see "Viral and bacterial sequence selection") using BLAST

710 v2.9.0+[97]. We then performed two BLASTp searches in this database, one with the
711 terminase (YP_009052554.1) and one with the polymerase (YP_009052497.1) of
712 crAssphage (NC_024711.1), with a bitscore cutoff of 50. All phage contigs that had (i)
713 a hit against both crAssphage terminase and polymerase and a query alignment of
714 ≥350 bp, and (ii) a contig length of ≥70 kbp were considered crAss-like phages. This
715 resulted in 146 crAss-like phage contigs, which were contained in 29 VCs.

716

717 <u>Candidatus *Heliusviridae* analysis</u>
718 To detect pairwise similarity, whole genome analyses were constructed with Easyfig
719 v2.2.5[98]. The prophage borders in NODE_38_length_205884_cov_102.806990 were
720 determined by determining the read depth along the entire contig from the bam files
721 with read mapping data ("Read mapping and community composition") using bedtools
722 genomecov v2.29.2[88] with option -bg. Resultant output was parsed and plotted in R.
723 Other related phages among the cohort were detected by performing a BLASTp search
724 with all phage ORFs of NODE_38_length_205884_cov_102.806990 against all phage
725 ORFs of the cohort with Diamond v2.0.4. This identified nine genes that were present
726 in 249 contigs. The ORFs on these contigs were annotated using PROKKA v1.14.6[99]
727 and InterProScan v5.48-83.0[100]. To identify isolated phages that share these nine
728 contigs, we performed a BLASTp against the NCBI nr-database using the NCBI
729 webserver[101] on February 26 2021 and collected all genomes with hits against all nine
730 genes (bit score ≥50).
731      The phages sharing all nine genes were clustered by analyzing them with
732 vContact2 v0.9.18[25], extracting the protein clustering data and calculating the number
733 of shared clusters between each pair of contigs. Contigs were clustered in R based on
734 Euclidean distances with the average agglomeration method.
735      To build a taxonomic tree, the nine genes were separately aligned using Clustal
736 Omega v1.2.4[102], positions with more than 90% gaps were removed with trimAl
737 v1.4.rev15[103] and alignments were concatenated. From the concatenated alignment,
738 an unrooted phylogenetic tree was built using IQ- Tree v2.0.3[104] using model finder[105]
739 and performing 1000 iterations of both SH-like approximate likelihood ratio test and the
740 ultrafast bootstrap approximation (UFBoot)[106]. In addition, ten iterations of the tree
741 were separately constructed, as has been recommended[107] (Zhou et al., 2018) (IQ-
742 Tree options -bb 1000, -alrt 1000, and --runs 10).

743

744 <u>Validation of *Ca. Heliusviridae* in other cohorts</u>

745 We used three additional studies to analyze prevalence of the *Ca. Heliusviridae*; one

746 composing of 145 participants used to study the gut virome in type 2 diabetes[11], a

747 second containing 196 participants and used to study the gut virome in hypertension[42],

748 and a final one containing ten healthy participants studied by VLP sequencing[43]. Reads

749 belonging to all studies were downloaded from the NCBI sequencing read archive

750 (SRA) and assembled as described above. The ten-patient VLP cohort was cross-

751 assembled, while the other two cohorts were assembled separately. After assembly,

752 ORFs were predicted using Prodigal v2.6.2[84]. *Ca. Heliusviridae* members were

753 identified by blastp using Diamond v2.0.4[108] against ORFs from each study, in which

754 the terminase, portal protein, Clp-protease, and major capsid protein of

755 NODE_38_length_205884_cov_102.806990 were used as queries. This was done

756 instead of all nine signature *Ca. Heliusviridae* genes to better allow for incomplete

757 assemblies. Contigs containing all four genes were selected, and a concatenated

758 alignment was made of the four head genes found in the T2D and hypertension

759 cohorts, plus all *Ca. Heliusviridae* in the tree depicted in Supplementary Figure 5.

760 These were then used to build a phylogenetic tree. The concatenated alignment and

761 phylogenetic tree were constructed as described above under "Candidatus

762 Heliusviridae analysis".

763

764 **<u>References</u>**

765

766 1.   Belkaid, Y. & Hand, T. W. Role of the microbiota in immunity and inflammation.
767      *Cell* **157**, 121–141 (2014).

768 2.   Rastelli, M., Cani, P. D. & Knauf, C. The Gut Microbiome Influences Host
769      Endocrine Functions. *Endocr. Rev.* **40**, 1271–1284 (2019).

770 3.   Gurung, M. *et al.* Role of gut microbiota in type 2 diabetes pathophysiology.
771      *EBioMedicine* **51**, 102590 (2020).

772 4.   Lang, S. & Schnabl, B. Microbiota and Fatty Liver Disease—the Known, the
773      Unknown, and the Future. *Cell Host Microbe* **28**, 233–244 (2020).

774 5.   Frank, D. N. *et al.* Molecular-phylogenetic characterization of microbial
775      community imbalances in human inflammatory bowel diseases. *Proc. Natl. Acad.*
776      *Sci.* **104**, 13780–13785 (2007).

777 6.   Clooney, A. G. *et al.* Whole-Virome Analysis Sheds Light on Viral Dark Matter in

778    Inflammatory Bowel Disease. *Cell Host Microbe* **26**, 764-778.e5 (2019).

779  7.  Norman, J. M. *et al.* Disease-Specific Alterations in the Enteric Virome in
780      Inflammatory Bowel Disease. *Cell* **160**, 447–460 (2015).

781  8.  Campbell, D. E. *et al.* Infection with Bacteroides Phage BV01 Alters the Host
782      Transcriptome and Bile Acid Metabolism in a Common Human Gut Microbe. *Cell*
783      *Rep.* **32**, 108142 (2020).

784  9.  Oh, J.-H. *et al.* Dietary Fructose and Microbiota-Derived Short-Chain Fatty Acids
785      Promote Bacteriophage Production in the Gut Symbiont Lactobacillus reuteri.
786      *Cell Host Microbe* **25**, 273-284.e6 (2019).

787  10. Reyes, A. *et al.* Gut DNA viromes of Malawian twins discordant for severe acute
788      malnutrition. *Proc. Natl. Acad. Sci.* **112**, 11941–11946 (2015).

789  11. Ma, Y., You, X., Mai, G., Tokuyasu, T. & Liu, C. A human gut phage catalog
790      correlates the gut phageome with type 2 diabetes. *Microbiome* **6**, 1–12 (2018).

791  12. De Sordi, L., Lourenço, M. & Debarbieux, L. The Battle Within: Interactions of
792      Bacteriophages and Bacteria in the Gastrointestinal Tract. *Cell Host Microbe* **25**,
793      210–218 (2019).

794  13. Paez-Espino, D. *et al.* Uncovering Earth's virome. *Nature* **536**, 425–430 (2016).

795  14. Gregory, A. C. *et al.* The Gut Virome Database Reveals Age-Dependent
796      Patterns of Virome Diversity in the Human Gut. *Cell Host Microbe* **28**, 724-740.e8
797      (2020).

798  15. Dutilh, B. E. *et al.* A highly abundant bacteriophage discovered in the unknown
799      sequences of human faecal metagenomes. *Nat. Commun.* **5**, 4498 (2014).

800  16. Yutin, N. *et al.* Discovery of an expansive bacteriophage family that includes the
801      most abundant viruses from the human gut. *Nat. Microbiol.* **3**, 38–46 (2018).

802  17. O'Neill, S. & O'Driscoll, L. Metabolic syndrome: A closer look at the growing
803      epidemic and its associated pathologies. *Obes. Rev.* **16**, 1–12 (2015).

804  18. Dabke, K., Hendrick, G. & Devkota, S. The gut microbiome and metabolic
805      syndrome. *J. Clin. Invest.* **129**, 4050–4057 (2019).

806  19. Mazidi, M., Rezaie, P., Kengne, A. P., Mobarhan, M. G. & Ferns, G. A. Gut
807      microbiome and metabolic syndrome. *Diabetes Metab. Syndr. Clin. Res. Rev.*
808      **10**, S150–S157 (2016).

809  20. Fujisaka, S. *et al.* Diet, Genetics, and the Gut Microbiome Drive Dynamic
810      Changes in Plasma Metabolites. *Cell Rep.* **22**, 3072–3086 (2018).

811  21. Ussar, S. *et al.* Interactions between gut microbiota, host genetics and diet

812      modulate the predisposition to obesity and metabolic syndrome. *Cell Metab.* **22**,
813      516–530 (2015).

22. Haro, C. *et al.* The gut microbial community in metabolic syndrome patients is modified by diet. *J. Nutr. Biochem.* **27**, 27–31 (2016).

23. Shkoporov, A. N. & Hill, C. Bacteriophages of the Human Gut: The "Known Unknown" of the Microbiome. *Cell Host Microbe* **25**, 195–209 (2019).

24. Deschasaux, M. *et al.* Depicting the composition of gut microbiota in a population with varied ethnic origins but shared geography. *Nat. Med.* **24**, 1526–1531 (2018).

25. Bin Jang, H. *et al.* Taxonomic assignment of uncultivated prokaryotic virus genomes is enabled by gene-sharing networks. *Nat. Biotechnol.* **37**, 632–639 (2019).

26. Manrique, P. *et al.* Healthy human gut phageome. *Proc. Natl. Acad. Sci.* **113**, 10400–10405 (2016).

27. Roux, S., Emerson, J. B., Eloe-Fadrosh, E. A. & Sullivan, M. B. Benchmarking viromics: An in silico evaluation of metagenome-enabled estimates of viral community composition and diversity. *PeerJ* **2017**, 1–26 (2017).

28. Camarillo-Guerrero, L. F., Almeida, A., Rangel-Pineros, G., Finn, R. D. & Lawley, T. D. Massive expansion of human gut bacteriophage diversity. *Cell* **184**, 1098-1109.e9 (2021).

29. Lin, H. & Peddada, S. Das. Analysis of compositions of microbiomes with bias correction. *Nat. Commun.* **11**, 1–11 (2020).

30. Aron-Wisnewsky, J. *et al.* Gut microbiota and human NAFLD: disentangling microbial signatures from metabolic disorders. *Nat. Rev. Gastroenterol. Hepatol.* **17**, 279–297 (2020).

31. Liu, R. *et al.* Gut microbiome and serum metabolome alterations in obesity and after weight-loss intervention. *Nat. Med.* **23**, 859–868 (2017).

32. Hryckowian, A. J. *et al.* Bacteroides thetaiotaomicron-Infecting Bacteriophage Isolates Inform Sequence-Based Host Range Predictions. *Cell Host Microbe* **28**, 371-379.e5 (2020).

33. Yutin, N. *et al.* Analysis of metagenome-assembled viral genomes from the human gut reveals diverse putative CrAss-like phages with unique genomic features. *Nat. Commun.* **12**, 1044 (2021).

34. Nayfach, S. *et al.* CheckV assesses the quality and completeness of

846    metagenome-assembled viral genomes. *Nat. Biotechnol.* (2020).
847    doi:10.1038/s41587-020-00774-7

848    35.  Von Meijenfeldt, F. A. B., Arkhipova, K., Cambuy, D. D., Coutinho, F. H. & Dutilh,
849         B. E. Robust taxonomic classification of uncharted microbial sequences and bins
850         with CAT and BAT. *Genome Biol.* **20**, 1–14 (2019).

851    36.  Hedzet, S., Accetto, T. & Rupnik, M. Lytic Bacteroides uniformis bacteriophages
852         exhibiting host tropism congruent with diversity generating retroelement. *bioRxiv*
853         2020.10.09.334284 (2020). doi:10.1101/2020.10.09.334284

854    37.  García-López, M. *et al.* Analysis of 1,000 Type-Strain Genomes Improves
855         Taxonomic Classification of Bacteroidetes. *Front. Microbiol.* **10**, (2019).

856    38.  Maya-Lucas, O. *et al.* The gut microbiome of Mexican children affected by
857         obesity. *Anaerobe* **55**, 11–23 (2019).

858    39.  Miquel, S. *et al.* Faecalibacterium prausnitzii and human intestinal health. *Curr.
859         Opin. Microbiol.* **16**, 255–261 (2013).

860    40.  Lavigne, R., Seto, D., Mahadevan, P., Ackermann, H. W. & Kropinski, A. M.
861         Unifying classical and molecular taxonomic classification: analysis of the
862         Podoviridae using BLASTP-based tools. *Res. Microbiol.* **159**, 406–414 (2008).

863    41.  Guerin, E. *et al.* Biology and Taxonomy of crAss-like Bacteriophages, the Most
864         Abundant Virus in the Human Gut. *Cell Host Microbe* **24**, 653-664.e6 (2018).

865    42.  Han, M., Yang, P., Zhong, C. & Ning, K. The Human Gut Virome in Hypertension.
866         *Front. Microbiol.* **9**, 1–10 (2018).

867    43.  Shkoporov, A. N. *et al.* The Human Gut Virome Is Highly Diverse, Stable, and
868         Individual Specific. *Cell Host Microbe* **26**, 527-541.e5 (2019).

869    44.  Cornuault, J. K. *et al.* The enemy from within: a prophage of Roseburia
870         intestinalis systematically turns lytic in the mouse gut, driving bacterial
871         adaptation by CRISPR spacer acquisition. *ISME J.* **14**, 771–787 (2020).

872    45.  Walther, B., Karl, J. P., Booth, S. L. & Boyaval, P. Menaquinones, Bacteria, and
873         the Food Supply: The Relevance of Dairy and Fermented Food Products to
874         Vitamin K Requirements. *Adv. Nutr.* **4**, 463–473 (2013).

875    46.  Moreno-Gallego, J. L. *et al.* Virome Diversity Correlates with Intestinal
876         Microbiome Diversity in Adult Monozygotic Twins. *Cell Host Microbe* **25**, 261-
877         272.e5 (2019).

878    47.  Zmora, N., Suez, J. & Elinav, E. You are what you eat: diet, health and the gut
879         microbiota. *Nat. Rev. Gastroenterol. Hepatol.* **16**, 35–56 (2019).

48. Falony, G. *et al.* Population-level analysis of gut microbiome variation. *Science (80-. ).* **352**, 560–564 (2016).

49. Zhernakova, A. *et al.* Population-based metagenomics analysis reveals markers for gut microbiome composition and diversity. *Science (80-. ).* **352**, 565–569 (2016).

50. Minot, S. *et al.* The human gut virome: Inter-individual variation and dynamic response to diet. *Genome Res.* **21**, 1616–1625 (2011).

51. Rodriguez-Valera, F. *et al.* Explaining microbial population genomics through phage predation. *Nat. Rev. Microbiol.* **7**, 828–836 (2009).

52. Koskella, B. & Brockhurst, M. A. Bacteria–phage coevolution as a driver of ecological and evolutionary processes in microbial communities. *FEMS Microbiol. Rev.* **38**, 916–931 (2014).

53. Silveira, C. B. & Rohwer, F. L. Piggyback-the-Winner in host-associated microbial communities. *npj Biofilms Microbiomes* **2**, 16010 (2016).

54. Fujimoto, K. *et al.* Metagenome Data on Intestinal Phage-Bacteria Associations Aids the Development of Phage Therapy against Pathobionts. *Cell Host Microbe* **28**, 380-389.e9 (2020).

55. Zaneveld, J. R., McMinds, R. & Vega Thurber, R. Stress and stability: applying the Anna Karenina principle to animal microbiomes. *Nat. Microbiol.* **2**, 17121 (2017).

56. Holmes, I., Harris, K. & Quince, C. Dirichlet Multinomial Mixtures: Generative Models for Microbial Metagenomics. *PLoS One* **7**, e30126 (2012).

57. Lourenço, M. *et al.* The Spatial Heterogeneity of the Gut Limits Predation and Fosters Coexistence of Bacteria and Bacteriophages. *Cell Host Microbe* **28**, 390-401.e5 (2020).

58. Hatfull, G. F. Dark Matter of the Biosphere: the Amazing World of Bacteriophage Diversity. *J. Virol.* **89**, 8107–8110 (2015).

59. Edwards, R. A., McNair, K., Faust, K., Raes, J. & Dutilh, B. E. Computational approaches to predict bacteriophage-host relationships. *FEMS Microbiol. Rev.* **40**, 258–272 (2016).

60. Burstein, D. *et al.* Major bacterial lineages are essentially devoid of CRISPR-Cas viral defence systems. *Nat. Commun.* **7**, 10613 (2016).

61. Džunková, M. *et al.* Defining the human gut host–phage network through single-cell viral tagging. *Nat. Microbiol.* **4**, 2192–2203 (2019).

62. de Jonge, P. A. *et al.* Adsorption Sequencing as a Rapid Method to Link Environmental Bacteriophages to Hosts. *iScience* **23**, 101439 (2020).

63. Walters, W. A., Xu, Z. & Knight, R. Meta-analyses of human gut microbes associated with obesity and IBD. *FEBS Lett.* **588**, 4223–4233 (2014).

64. Drulis-Kawa, Z., Majkowska-Skrobek, G., Maciejewska, B., Delattre, A.-S. & Lavigne, R. Learning from Bacteriophages - Advantages and Limitations of Phage and Phage-Encoded Protein Applications. *Curr. Protein Pept. Sci.* **13**, 699–722 (2012).

65. Ridaura, V. K. *et al.* Gut Microbiota from Twins Discordant for Obesity Modulate Metabolism in Mice. *Science (80-. ).* **341**, 1241214 (2013).

66. David, L. A. *et al.* Diet rapidly and reproducibly alters the human gut microbiome. *Nature* **505**, 559–563 (2014).

67. De Filippis, F. *et al.* Distinct Genetic and Functional Traits of Human Intestinal Prevotella copri Strains Are Associated with Different Habitual Diets. *Cell Host Microbe* **25**, 444-453.e3 (2019).

68. Shkoporov, A. N. *et al.* ΦCrAss001 represents the most abundant bacteriophage family in the human gut and infects Bacteroides intestinalis. *Nat. Commun.* **9**, 4781 (2018).

69. Koonin, E. V. & Yutin, N. The crAss-like Phage Group: How Metagenomics Reshaped the Human Virome. *Trends Microbiol.* **28**, 349–359 (2020).

70. Edwards, R. A. *et al.* Global phylogeography and ancient evolution of the widespread human gut virus crAssphage. *Nat. Microbiol.* **4**, 1727–1736 (2019).

71. Garmaeva, S. *et al.* Studying the gut virome in the metagenomic era: Challenges and perspectives. *BMC Biol.* **17**, 1–14 (2019).

72. Zhao, L. *et al.* Gut bacteria selectively promoted by dietary fibers alleviate type 2 diabetes. *Science (80-. ).* **359**, 1151–1156 (2018).

73. Narita, M. The gut microbiome as a target for prevention of allergic diseases. *Japanese J. Allergol.* **69**, 19–22 (2020).

74. De La Cuesta-Zuluaga, J. *et al.* Metformin is associated with higher relative abundance of mucin-degrading akkermansia muciniphila and several short-chain fatty acid-producing microbiota in the gut. *Diabetes Care* **40**, 54–62 (2017).

75. Gazitúa, M. C. *et al.* Potential virus-mediated nitrogen cycling in oxygen-depleted oceanic waters. *ISME J.* **15**, 981–998 (2021).

76. Sharon, I. *et al.* Photosystem I gene cassettes are present in marine virus

948      genomes. *Nature* **461**, 258–262 (2009).

949    77.   Snijder, M. B. *et al.* Cohort profile: The Healthy Life in an Urban Setting (HELIUS)
950      study in Amsterdam, the Netherlands. *BMJ Open* **7**, 1–11 (2017).

951    78.   Mobini, R. *et al.* Metabolic effects of Lactobacillus reuteri DSM 17938 in people
952      with type 2 diabetes: A randomized controlled trial. *Diabetes, Obes. Metab.* **19**,
953      579–589 (2017).

954    79.   Alberti, K. G. M. M. *et al.* Harmonizing the metabolic syndrome: A joint interim
955      statement of the international diabetes federation task force on epidemiology and
956      prevention; National heart, lung, and blood institute; American heart association;
957      World heart federation; International . *Circulation* **120**, 1640–1645 (2009).

958    80.   Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. MetaSPAdes: A new
959      versatile metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).

960    81.   Pagès H, Aboyoun P, Gentleman R, D. S. Biostrings: Efficient manipulation of
961      biological strings. (2020).

962    82.   Gregory, A. C. *et al.* Marine DNA Viral Macro- and Microdiversity from Pole to
963      Pole. *Cell* **177**, 1109-1123.e14 (2019).

964    83.   Roux, S., Enault, F., Hurwitz, B. L. & Sullivan, M. B. VirSorter: mining viral signal
965      from microbial genomic data. *PeerJ* **3**, e985 (2015).

966    84.   Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation
967      site identification. *BMC Bioinformatics* **11**, 119 (2010).

968    85.   Verhaar, B. J. H. *et al.* Associations between gut microbiota, faecal short-chain
969      fatty acids, and blood pressure across ethnic groups: the HELIUS study. *Eur.*
970      *Heart J.* **41**, 4259–4267 (2020).

971    86.   Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat.*
972      *Methods* **9**, 357–359 (2012).

973    87.   Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
974      **25**, 2078–2079 (2009).

975    88.   Quinlan, A. R. *BEDTools: The Swiss-Army tool for genome feature analysis.*
976      *Current Protocols in Bioinformatics* **2014**, (2014).

977    89.   McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible
978      Interactive Analysis and Graphics of Microbiome Census Data. *PLoS One* **8**,
979      e61217 (2013).

980    90.   Dixon, P. VEGAN, a package of R functions for community ecology. *J. Veg. Sci.*
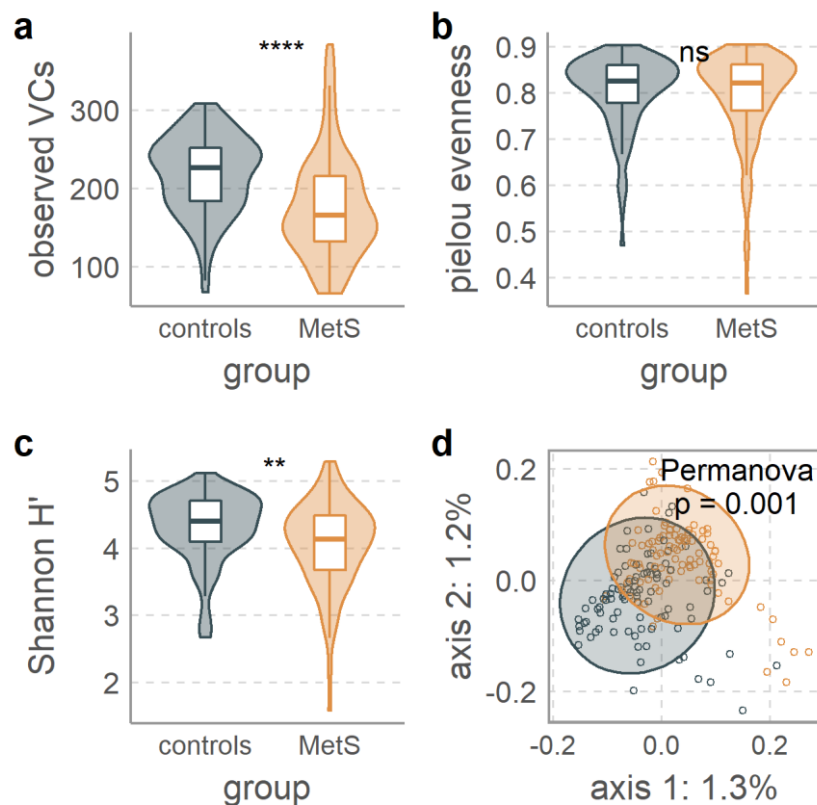981      **14**, 927–930 (2003).

91. Biswas, A., Staals, R. H. J., Morales, S. E., Fineran, P. C. & Brown, C. M. CRISPRDetect: A flexible algorithm to define CRISPR arrays. *BMC Genomics* **17**, 1–14 (2016).

92. Nobrega, F. L., Walinga, H., Dutilh, B. E. & Brouns, S. J. J. J. Prophages are associated with extensive CRISPR–Cas auto-immunity. *Nucleic Acids Res.* **48**, 12074–12084 (2020).

93. Wattam, A. R. *et al.* PATRIC, the bacterial bioinformatics database and analysis resource. *Nucleic Acids Res.* **42**, 581–591 (2014).

94. Zhang, R. *et al.* SpacePHARER: Sensitive identification of phages from CRISPR spacers in prokaryotic hosts. *bioRxiv* 2020.05.15.090266 (2020). doi:10.1101/2020.05.15.090266

95. Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): A curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* **35**, 61–65 (2007).

96. Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome datasets are compositional: And this is not optional. *Front. Microbiol.* **8**, 1–6 (2017).

97. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).

98. Sullivan, M. J., Petty, N. K. & Beatson, S. A. Easyfig: A genome comparison visualizer. *Bioinformatics* **27**, 1009–1010 (2011).

99. Seemann, T. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).

100. Jones, P. *et al.* InterProScan 5: Genome-scale protein function classification. *Bioinformatics* **30**, 1236–1240 (2014).

101. Johnson, M. *et al.* NCBI BLAST: a better web interface. *Nucleic Acids Res.* **36**, 5–9 (2008).

102. Sievers, F. & Higgins, D. G. Clustal Omega for making accurate alignments of many protein sequences. *Protein Sci.* **27**, 135–145 (2018).

103. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: A tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973 (2009).

104. Nguyen, L. T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: A fast and effective stochastic algorithm for estimating maximum-likelihood

1016      phylogenies. *Mol. Biol. Evol.* **32**, 268–274 (2015).

1017  105. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermiin,

1018      L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat.*

1019      *Methods* **14**, 587–589 (2017).

1020  106. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S.

1021      UFBoot2: Improving the Ultrafast Bootstrap Approximation. Molecular biology

1022      and evolution. *Mol. Biol. Evol.* **35**, 518–522 (2018).

1023  107. Zhou, X., Shen, X. X., Hittinger, C. T. & Rokas, A. Evaluating fast maximum

1024      likelihood-based phylogenetic programs using empirical phylogenomic data

1025      sets. *Mol. Biol. Evol.* **35**, 486–503 (2018).

1026  108. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using

1027      DIAMOND. *Nat. Methods* **12**, 59–60 (2014).

1028  109. Charrad, M., Ghazzali, N., Boiteau, V. & Niknafs, A. NbClust : An R Package for

1029      Determining the Relevant Number of Clusters in a Data Set. *J. Stat. Softw.* **61**,

1030      11744–11750 (2014).

1031

1032  **Figures**

1033



1034

1035 **Figure 1: Gut phage populations are altered in MetS. a** MetS-associated decreased phage

1036 species richness is evidenced by the number of unique VCs observed per sample. **b** No

1037 change in phage pielou evenness measurements. **c** significantly decreased phage α-diversity

1038 measured by Shannon diversity. **d** clear separation between phageomes of MetS and control

1039 participant as shown by β-diversity depicted in a principal coordinates analysis (PCoA) of Bray-

1040 Curtis dissimilarities. Permanova test was adjusted for smoking, age, sex, alcohol use, and

1041 metformin use. Statistical significance in A-C is according to the Wilcoxon signed rank test,

1042 where p-values are denoted as follows: ns not significant, * ≤ 0.05, ** ≤ 0.01, *** ≤ 0.001, ****

1043 ≤ 0.0001. Box plots show the median, 25$^{th}$, and 75$^{th}$ percentile, with upper and lower whiskers

1044 to the 25$^{th}$ percentile minus and the 75$^{th}$ percentile plus 1.5 times the interquartile range.

1045



1046

1047 **Figure 2: Correlations between phage and bacterial populations as well as between**

1048 **population measures and MetS clinical parameters.** Strong correlations between **a** phage

1049 richness (observed VCs) and bacterial richness (Chao1 index), as well as between **b** phage

1050 and bacterial evenness (Pielou's index), both with significant positive Spearman's rank

1051 correlation coefficient. Both of these measures were correlated to MetS clinical parameters.

1052   Plotted are the Spearman's rank correlation coefficients between the five MetS risk factors and

1053   **c** richness and **d** evenness. Points with p-values below 0.05 are colored in and labeled.

1054



1055

**Figure 3: Phages infecting selected bacterial families are differentially abundant in MetS or healthy controls. a** ANCOM-BC[29] analysis of VCs that infect the twelve bacterial families to which the most VCs were linked shows significant association between *Bacteroidaceae* VCs and MetS, as well as between *Ruminococcaceae, Acidominacoccaceae, and Tannerellaceae* VCs and healthy controls. ANCOM-BC was adjusted for smoking, age, sex, alcohol use, and metformin use. **b** relative abundance comparisons between MetS and control participants of VCs infecting *Faecalibacterium sp.* CAG:74, *Ruthenibacterium lactatiformans,*

1063 *Subdoligranulum sp.* APC924/74. Stars denote significance according to the Wilcoxon signed

1064 rank test, with p-values adjusted with the Benjamini and Hochberg procedure (*q*). * ≤ 0.05, **

1065 ≤ 0.01, *** ≤ 0.001, **** ≤ 0.0001. Box plots show the median, 25th, and 75th percentile, with

1066 upper and lower whiskers to the 25th percentile minus and the 75th percentile plus 1.5 times

1067 the interquartile range. Error bars in **a** denote the standard error adjusted by the Benjamini-

1068 Hochberg procedure for multiple testing.



1069

1070 **Figure 4: Among significantly differentially abundant VCs some are related. a** VCs

1071 identified by ANCOM-BC analysis as significantly abundant (q ≤ 0.05 after implementing the

1072 Benjamini-Hochberg procedure for multiple testing). Error bars denote ANCOM-BC-supplied

1073 standard error. The analysis was adjusted for smoking, age, sex, alcohol use, and metformin

1074 use. Taxonomic names to the right of the plot denote host predictions, which are colored as

1075 follows: *Firmicutes;* grey, *Bacteroidetes;* red, *Actinobacteria*; green. The full taxonomies are

1076 listed in Supplementary Tables 2 and 4. For brevity, only the ten VCs most significantly

1077 associated with MetS (out of 38) are shown. See Supplementary table 7 for a full reckoning of

1078 significant VCs and the full names of the two singletons. **b** Whole genome analysis of three

1079 contigs that belong to VC_659_0, VC_745_0 and VC_643_0. The VC_659_0 contig is zoomed

1080 in on the prophage region, for the entire contig, see Figure 6. The read coverage depth of this

1081    contig in two samples is displayed at the top, on in which the prophage is present (S194) and

1082    one in which it is absent (S095). The nine genes shared by all *Candidatus Heliusviridae* are

1083    colored red, and annotated at the bottom.

1084



1085

**Figure 5: Three VCs that are hallmarks for either MetS or healthy control phageomes are part of the widespread *Candidatus Heliusviridae* family of gut phages. a** heatmap and hierarchical clustering of pairwise shared protein cluster values for 249 contigs from the current study and 51 previously isolated phages. The line graph shows the optimal number of clusters as determined using the NbClust R package[109], and the dendrogram is cut to form six clusters. These six clusters are labeled as alpha, beta, gamma, delta, epsilon, and zeta subfamilies. The top row of colors beneath the dendrogram denote the differentially abundant VCs, from left to right: VC_745_0 (red), VC_659_0 (green), and VC_643_0 (purple). The bottom colors are according to the candidate subfamilies. **b** the prevalence of the *Candidatus Heliusviridae* (left) and the separate candidate subfamilies (right). **c** the relative abundances of the candidate subfamilies (the whole family was not significantly more abundant in either group and is thus not depicted). q-values are denoted as follows * ≤ 0.05, ** ≤ 0.01, *** ≤ 0.001, **** ≤ 0.0001. Box plots show the median, 25th, and 75th percentile, with upper and lower whiskers to the 25th percentile minus and the 75th percentile plus 1.5 times the interquartile range.



**Figure 6: VC_659_0 infects *Roseburia* and *Blautia*, and carries possible auxiliary metabolic genes. a** Whole genome alignment of three prophages contained within VC_659_0, with pie charts denoting the top BLASTp hit of all host genes on the contigs. **b** and **c** the mean coverage of host-derived regions in NODE_38 (**b**) and NODE_192 (**c**). Significance according to Wilcoxon signed rank test, p-values are denoted as follows * ≤ 0.05, ** ≤ 0.01, *** ≤ 0.001,

1107 **** ≤ 0.0001. Box plots show the median, 25th, and 75th percentile, with upper and lower

1108 whiskers to the 25th percentile minus and the 75th percentile plus 1.5 times the interquartile
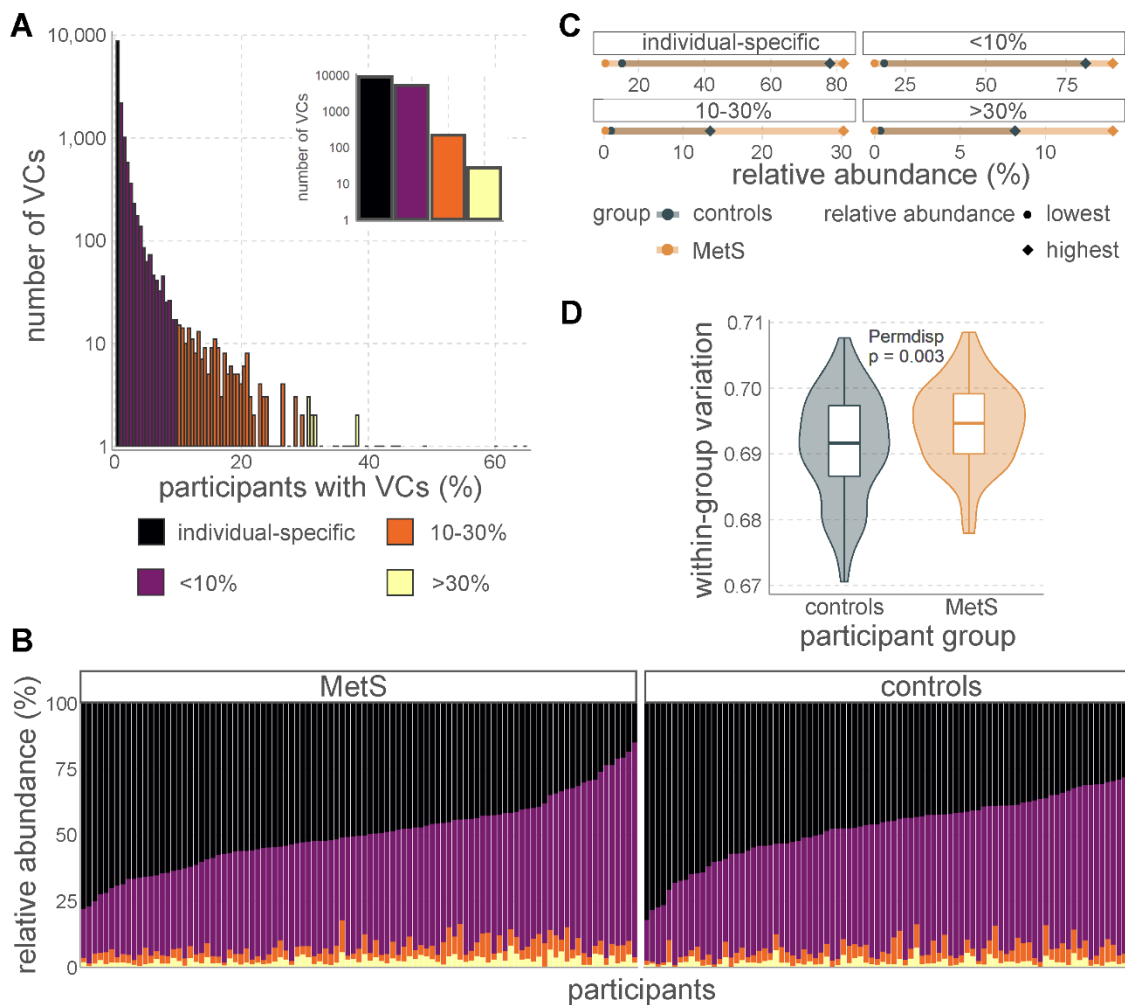
1109 range.

1110

1111

1112 **Supplementary Figures**

1113

1114



1115

1116 **Supplementary Figure 1: Overview of the phageomes show more variability among**

1117 **MetS participants. a** Histogram of VCs by number of participants that they are found in shows

1118 most VCs are individual-specific. The inset is the same dataset with one bar for each category

1119 shown in the legend. **b** Stacked bar charts of community composition show high inter-individual

1120 phageome diversity. Color legend is identical to (A). **c** Comparisons show that participants with

1121 the highest and lowest relative abundance in each VC category all belong to the MetS group.

1122 **d** MetS participants have significantly higher within-group variation, as measured by permdisp

1123 on Bray-Curtis dissimilarities. Box plot shows the median, 25th, and 75th percentile, with upper

1124 and lower whiskers to the 25th percentile minus and the 75th percentile plus 1.5 times the
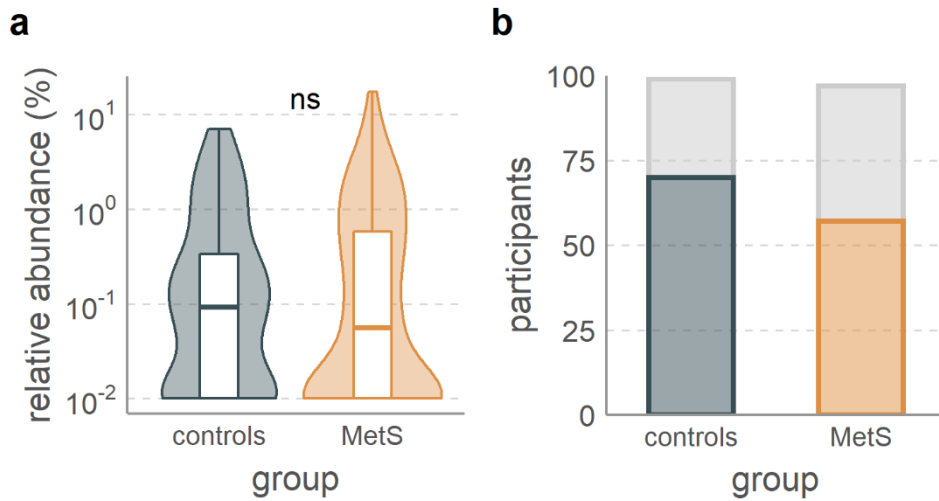
1125 interquartile range.

1126



1127

**Supplementary Figure 2: Differences in total phage abundance and potential temperate phage abundance. a** total phage abundance, as shown by the percentage of reads that map to phage sequences. **b** significantly more reads map to bacterial contigs that contain prophage-like sequences. **c** no significant difference in relative abundance of VCs that carry integrase genes. Stars denote significance according to the Wilcoxon signed rank test. * ≤ 0.05, ** ≤ 0.01, *** ≤ 0.001, **** ≤ 0.0001. Box plots show the median, 25th, and 75th percentile, with upper and lower whiskers to the 25th percentile minus and the 75th percentile plus 1.5 times the interquartile range.

1136

**Supplementary Figure 3: Gut bacterium populations are altered in MetS. a** MetS-associated decreased bacterial species richness is evidenced by the Chao1 index. **b** decreased bacterial pielou evenness measurements. **c** significantly decreased bacterial α-diversity measured by Shannon diversity. **d** clear separation between bacterial populations of MetS and control participant as shown by β-diversity depicted in a principal coordinates analysis (PCoA) of Bray-Curtis dissimilarities. Permanova test was adjusted for smoking, age, sex, alcohol use, and metformin use. Statistical significance in A-C is according to the Wilcoxon signed rank test, where p-values are denoted as follows: ns not significant, * ≤ 0.05, ** ≤ 0.01, *** ≤ 0.001, **** ≤ 0.0001. Box plots show the median, 25th, and 75th percentile, with upper and lower whiskers to the 25th percentile minus and the 75th percentile plus 1.5 times the interquartile range.

1148

1149

1150



1151

1152 **Supplementary Figure 4: Non-significant differences in crAss-like phage populations. a**
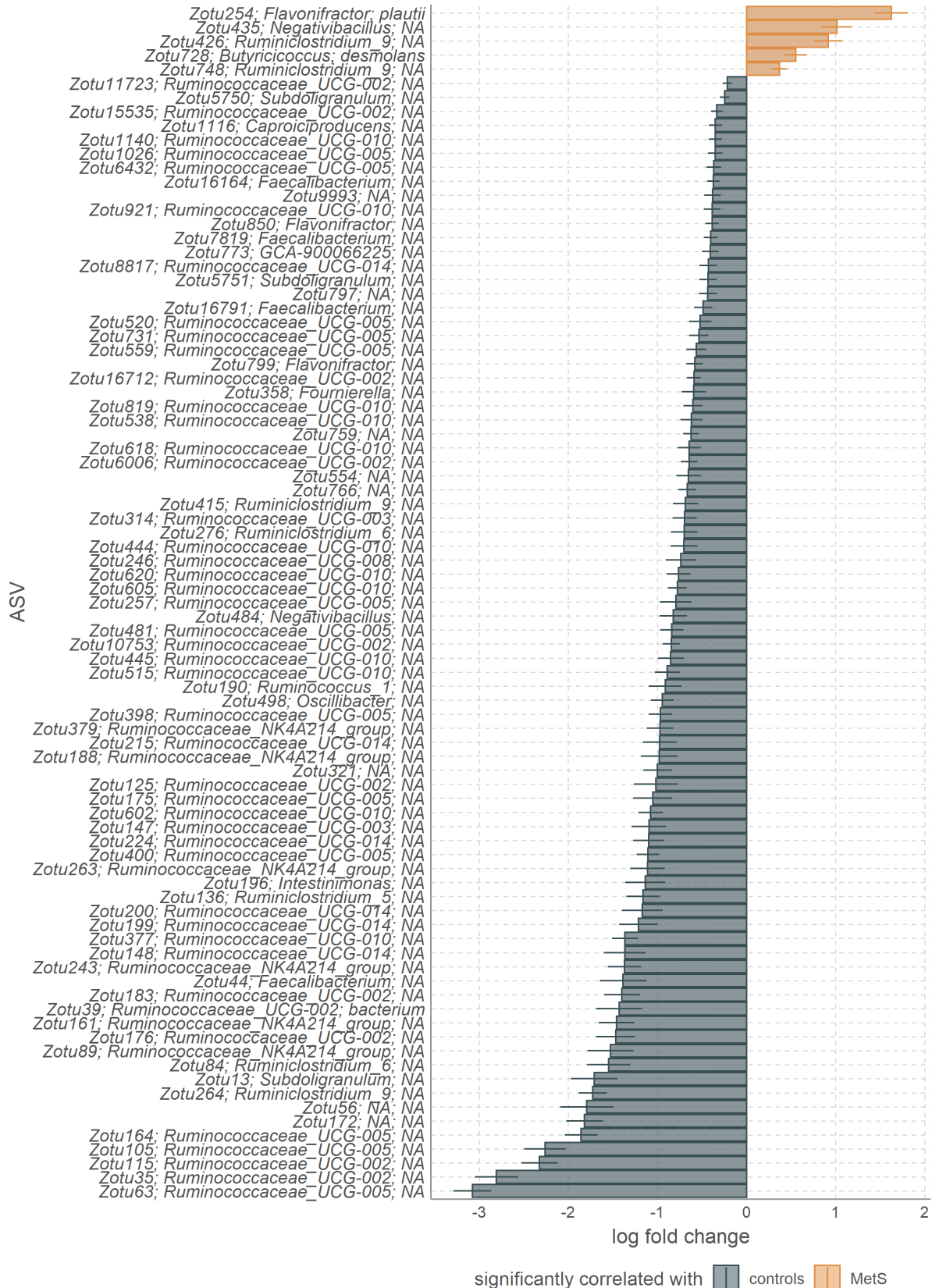
1153 relative abundance of crAss-like phages in controls and MetS. **b** the number of participants in

1154 which crAss-like phages were present. Box plots show the median, 25th, and 75th percentile,

1155 with upper and lower whiskers to the 25th percentile minus and the 75th percentile plus 1.5
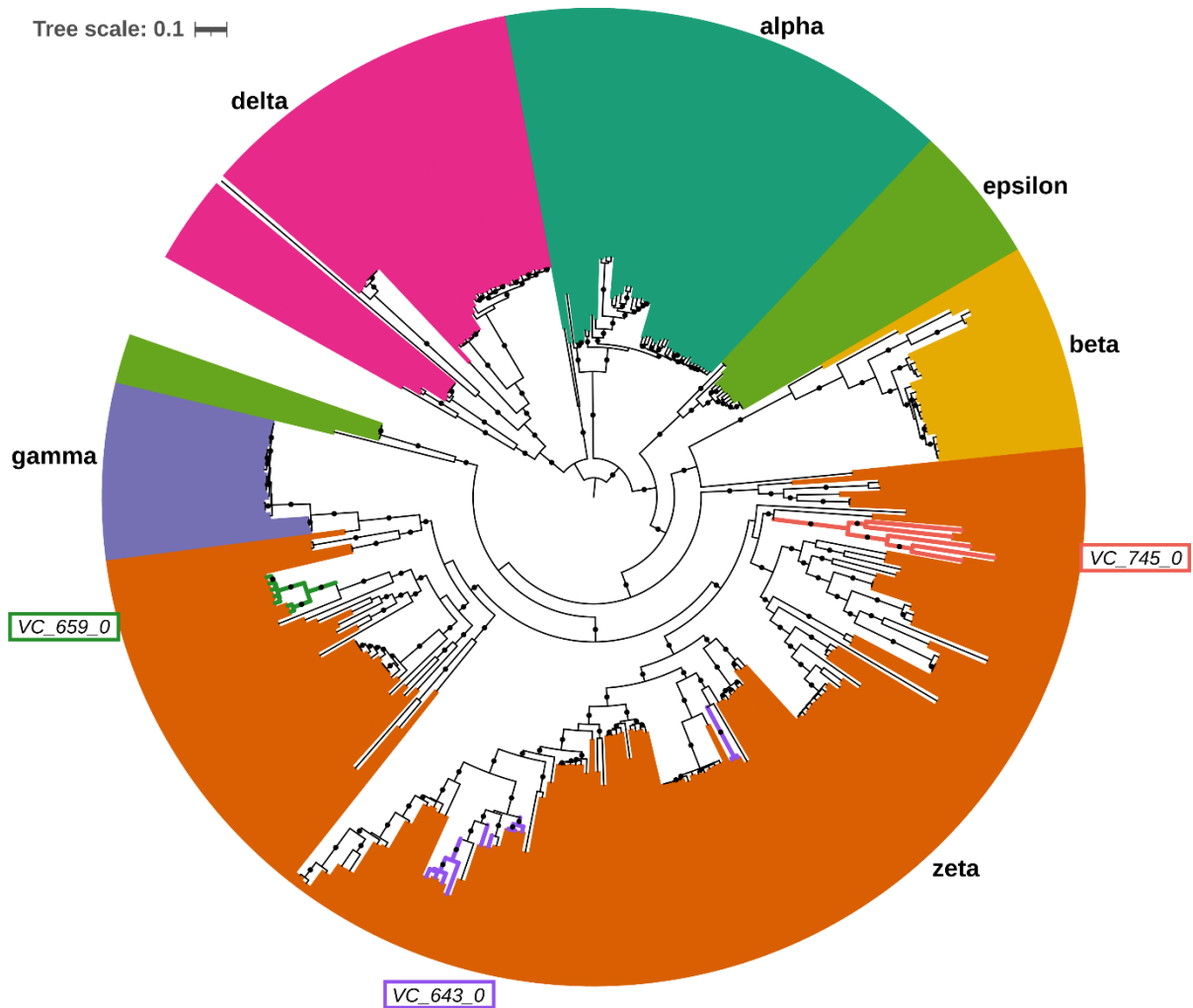
1156 times the interquartile range.

**Supplementary Figure 5:** ANCOM-BC analysis results of significantly differentially abundant *Ruminococcaceae* ASVs. Error bars denote the standard error with Holm adjustment for multiple testing.
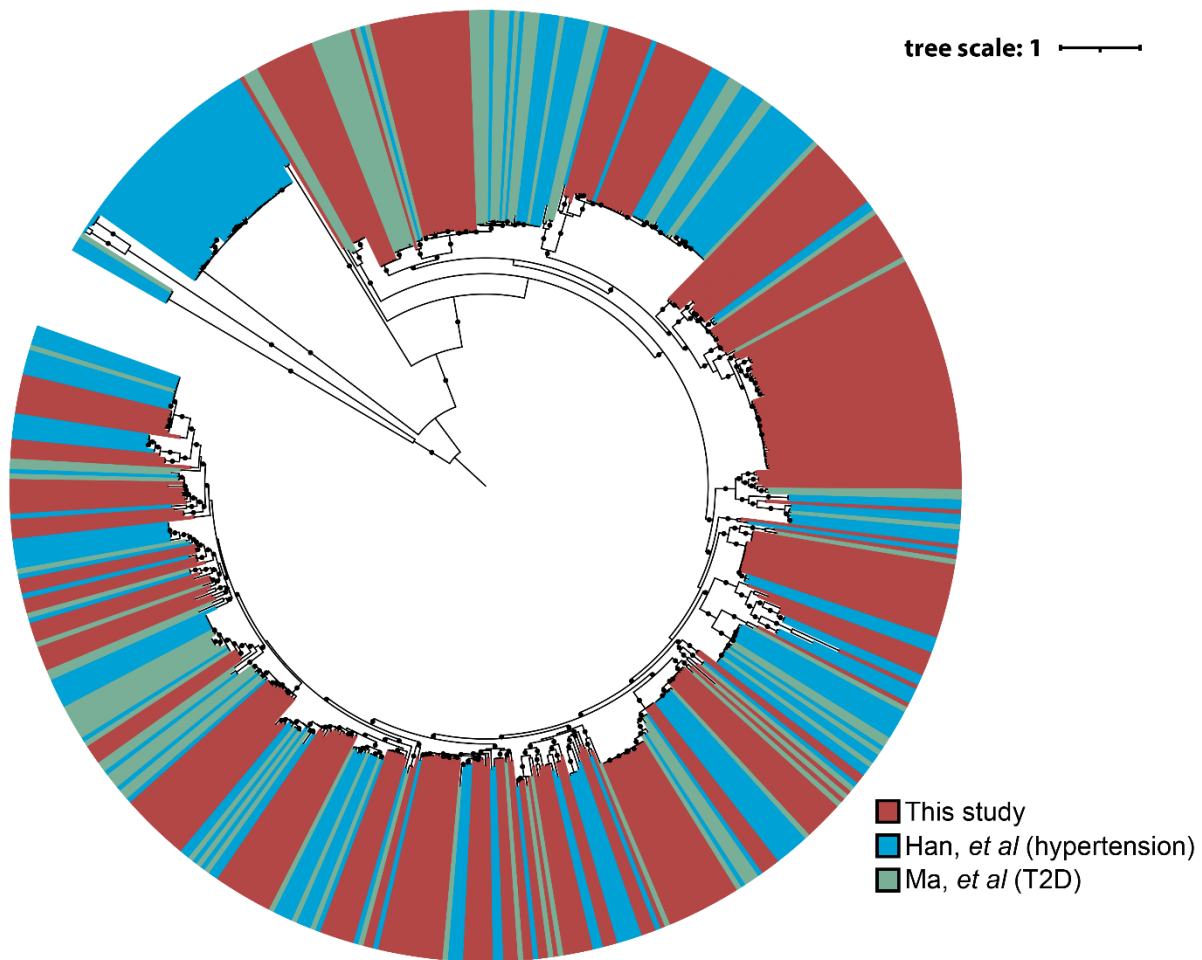
1161



1162

**Supplementary Figure 6: A midpoint-rooted approximate maximum likelihood tree made from the concatenated alignments of the nine universally shared *Candidatus Heliusviridae* genes, with colors denoting the groups. Dots represent bootstrap values of ≥95. Branch colors show contigs that belong to the three *Ca. Heliusviridae* VCs that are significantly differentially abundant in either controls or MetS participants.**

1168

**tree scale: 1**
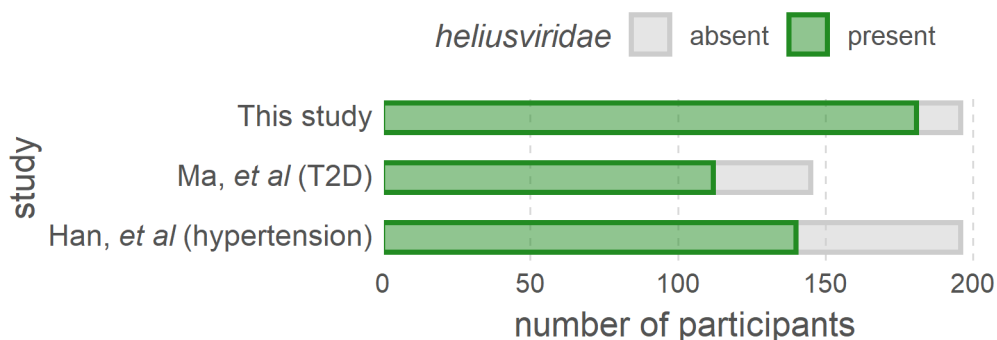
This study
Han, *et al* (hypertension)
Ma, *et al* (T2D)

1169

1170 **Supplementary Figure 7: A midpoint-rooted approximate maximum likelihood tree made**

1171 **from the concatenated alignments of the four structural *Candidatus Heliusviridae* genes**

1172 **in contigs from this study and two cohorts in which the phageome was analyzed before,**

1173 **with colors denoting the study. Dots represent bootstrap values of ≥95.**

1174



1175

1176 **Supplementary Figure 8: Occurrence of *Candidatus Heliusviridae* in this study and two**

1177 **validation cohorts. To circumvent incomplete assemblies, contigs were identified as**

1178 ***Candidatus Heliusviridae* if they 1) contained the terminase, portal protein, major capsid**

1179 **protein, and clp-proteas, and 2) were located in the same clade as *Candidatus***

1180     *Heliusviridae* **from this study in the phylogenetic tree depicted in Supplementary Figure**

1181     **7.**