1  **Individuals at risk for developing rheumatoid arthritis harbor differential intestinal**

2  **bacteriophage communities with distinct metabolic potential**

3

4

5

6  Mihnea R. Mangalea[1], David Paez-Espino[2,3], Kristopher Kieft[4], Anushila Chatterjee[1], Jennifer A.

7  Seifert[5], Marie L. Feser[5], M. Kristen Demoruelle[5], Meagan E. Chriswell[5], Alexandra Sakatos[3],

8  Karthik Anantharaman[4], Kevin D. Deane[5], Kristine A. Kuhn[5], V. Michael Holers[5], and Breck A.

9  Duerkop[1,6,*]

10

11  [1]Department of Immunology and Microbiology, University of Colorado School of Medicine,

12  Aurora, CO, USA

13  [2]Mammoth Biosciences, San Francisco, CA, USA

14  [3]Anicilia Therapeutics, New York, NY, USA

15  [4]Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA

16  [5]Division of Rheumatology, University of Colorado School of Medicine, Aurora, CO, USA

17  [6]Lead Contact

18  [*]Correspondence: breck.duerkop@cuanschutz.edu

19

20

21

22

23

24

25

26

**SUMMARY**

Rheumatoid arthritis (RA) is an autoimmune disease characterized in seropositive individuals by the presence of anti-cyclic citrullinated protein (CCP) antibodies. RA is linked to the intestinal microbiota, yet the association of microbes with CCP serology and their contribution to RA is unclear. We describe intestinal phage communities of individuals at risk for developing RA, with or without anti-CCP antibodies, whose first degree relatives have been diagnosed with RA. We show that at-risk individuals harbor intestinal phage compositions that diverge based on CCP serology, are dominated by Lachnospiraceae phages, and originate from disparate ecosystems. These phages encode unique repertoires of auxiliary metabolic genes (AMGs) which associate with anti-CCP status, suggesting that these phages directly influence the metabolic and immunomodulatory capability of the microbiota. This work sets the stage for the use of phages as preclinical biomarkers and provides insight into a possible microbial-based causation of RA disease development.

**KEYWORDS**

Bacteriophages, rheumatoid arthritis, autoimmune disease, microbiome, phage-host interaction, phage-host metabolism

**INTRODUCTION**

Rheumatoid arthritis (RA) is a systemic autoimmune disease with a global prevalence of approximately 1%. The development of RA in at-risk individuals is dependent on a combination of genetics, epidemiological factors, and systemic immune dysregulation [1]. The heritability of RA is estimated to be 40–60%, with increased familial risk evident among first-degree relatives (FDRs) of individuals with diagnosed RA [2, 3]. Analyses of at-risk FDRs, even those without serum RA-related autoantibodies, have identified patterns of mucosal inflammation whereby anti-cyclic citrullinated peptide (anti-CCP) antibodies and rheumatoid factors (RF), as well as cytokines and chemokines, are expressed locally in a subset of individuals [4-6]. In addition, anti-CCP and RF are present in the blood for years prior to the onset of RA, and their presence as well as circulating cytokine and chemokine biomarkers, are predictive of future RA development [7-9]. To probe the mucosal origins hypothesis [1] and the mounting evidence implicating intestinal microbiota perturbations in RA etiopathogenesis [10], it is necessary to characterize the ecological associations of the microbiota in at-risk individuals susceptible to RA.

Studies linking the role of the intestinal microbiota to systemic autoimmune diseases predominantly rely on 16S ribosomal gene analyses of bacteria within the microbiome, and have expanded our understanding of dysbiosis in the RA intestine. Individuals with established RA harbor a microbiota dominated by *Prevotella copri* [11, 12], enriched with Gram-positive bacteria [13], and decreased carriage of bifidobacteria [14], Gram-negative *Bacteroides,* and Firmicutes [13, 15]. The association of enriched Prevotellaceae, including *P. copri*, has also been described in individuals with preclinical RA [16], indicating that intestinal *P. copri* is immune-relevant to the pathogenesis of RA [17]. The presence of *P. copri* may therefore represent a biological indicator and additional risk factor for RA development and progression [18]. However, associating a single organism to RA etiology neglects the interactions of bacteria with

3

78    their surrounding environment and other bacterial community members whose populations can

79    be influenced by predatory bacteriophages (phages).

80         In contrast to the recent enthusiasm for characterizing microbial links to the etiology of

81    RA, relatively little is known concerning the composition of phage communities in the intestine

82    as it relates to RA disease risk. Phages of the intestinal microbiota can fluctuate in community

83    composition in response to immune system function and disease, which suggests that they

84    could be exploited as biomarkers for early disease detection [19]. Metagenomic sequencing

85    strategies have revealed extensive and diverse populations of phages in the human intestine

86    [20-22], in which phage community dynamics correlate with distinct disease states [23-25].

87    Specific intestinal phage genomic signatures precede autoimmunity development of type 1

88    diabetes in a cohort of diabetes-susceptible children, with disease-associated phages

89    correlating to the bacterial component of the microbiota [26]. In addition to the direct impact of

90    intestinal phages on bacterial community composition via classical predation and prophage

91    mediated bacterial competition and metabolism, phages also adhere to mucosal surfaces,

92    significantly impacting microbial colonization [27] and host mucosal immunity development [28].

93    Evidence is emerging that phages are also immunomodulatory through intrinsic anti-

94    inflammatory properties, and are capable of direct lymphocyte regulation through the ability to

95    translocate to multiple tissues and organs [29]. Despite these observations and potential

96    implications for systemic autoimmune diseases like RA, evaluation of intestinal phages in the

97    context of RA disease risk has yet to be described.

98         The interplay between intestinal bacteria, their phages, and the host immune system,

99    whose interactions have consequences not only for compositional dysbiosis but

100    immunomodulation, must be considered in the etiopathogenesis of RA. The microbiome, and

101    more recently the virome, have been implicated in a range of human diseases including cancers

102    [30, 31], inflammatory bowel diseases [32, 33], and arthritis [11, 34]. By characterizing the

103    phage populations in an at-risk RA FDR cohort; further sub-grouped with regard to autoantibody

4

104    status as defined by the presence of anti-CCP antibodies and compared to healthy controls, we

105    have begun to address this question. The cohort contains individuals that do not have

106    inflammatory arthritis or established RA disease but are FDRs to an individual with diagnosed

107    RA, which alone increases RA risk. Studying the microbiomes of at-risk individuals in the

108    preclinical RA state could lead to the identification of biomarkers and therapeutic targets

109    independent of confounding by the use of drugs in subjects with active arthritis.

110        We used metagenomics to define intestinal phage populations of anti-CCP positive

111    (CCP+) and negative (CCP-) individuals in an at-risk FDR cohort. Phage matching to bacterial

112    hosts showed divergent intestinal phage communities dependent on anti-CCP serology status.

113    We observed an overabundance of phages targeting Bacteroidaceae and Sreptococcaceae

114    bacteria in CCP+ at-risk FDRs as well as phages targeting Bacteroidaceae bacteria in CCP- at-

115    risk FDRs. Importantly, analysis of the metabolic traits encoded in phage metagenomes

116    revealed intra-cohort profiles reflecting distinct immunomodulatory potential. Phages with

117    auxiliary metabolic genes (AMGs) that modify lipopolysaccharide and other outer membrane

118    glycans of host bacteria were differentially abundant, implicating modifications to bacterial

119    antigenicity [35] and bacterial fitness [36] in RA-associated communities. Core phage metabolic

120    genes, including 14 genes which are globally conserved among phages from multiple diverse

121    environments [37], as well as bacterial surface modifying enzymes, were associated with

122    phages targeting *Flavonifractor* sp. in the CCP+ cohort and *Bacteroides* sp. in the CCP- cohort.

123    Phages targeting Lachnospiraceae (*Clostridium scindens*) and Actinomyces (*A. oris*), including

124    several AMGs, were over-abundant among CCP+ and CCP- individuals, respectively, compared

125    to healthy controls. Our data show that there are unique and abundant intestinal phages specific

126    to RA-susceptibility status, and this highlights their potential as biomarkers for preclinical RA

127    and the need for further pursuit of community-level bacteria-phage interactions during the

128    development and progression of RA.

129

## RESULTS

**First-degree relatives to individuals with rheumatoid arthritis.**

A total of 25 human subjects were identified from the Studies of the Etiology of Rheumatoid Arthritis (SERA) [38], including 16 FDRs of individuals with RA and 9 age and sex matched healthy controls (HC). FDR subjects for which a detectable level of anti-CCP autoantibody was present (defined by a value of ≥ 20 units/mL in either ELISA assay for anti-CCP3.1 IgA/IgG or anti-CCP3 IgG (Inova Diagnostics) [39]) were designated the CCP+ group (n = 8). FDRs with no anti-CCP detected were designated the CCP- group (n = 8) (Table 1). Mean ages for the three groups in this study were 61.3 ± 11.0 for CCP+, 49.0 ± 15.7 for CCP-, and 44.4 ± 13.6 for HC. The distribution of sexes for each group is reported as percent female, with 88.9% for CCP+, 62.5% for CCP-, and 66.7% for HC. Among the CCP+ and HC groups, 3/9 and 2/9 of individuals have reported ever smoking (a risk factor associated with RA), respectively (Table 1).

**Generation and curation of *de novo* assembled VLP contigs.**

We used individual fecal samples from the subjects obtained at the time of autoantibody and clinical evaluations, and isolated total genomic DNA for shotgun metagenomic sequencing using an untargeted amplification-independent approach [23, 40]. Samples were physically separated into whole metagenome (M), including all genomic DNA present in the sample, and virus-like particle (VLP) fractions, which were subjected to phage-specific precipitation (Figure S1A). Illumina sequencing resulted in an average of 123.8 ± 32.2, 135.2 ± 40.4, and 104.7 ± 45.9 million (M) paired end reads per sample for CCP+, CCP- and HC whole metagenomes, respectively, and an average of 67.3 ± 29.5, 73.2 ± 33.7, and 89.6 ± 47.8 M paired reads per sample for CCP+, CCP- and HC VLP fractions, respectively (Figure S1B). VLP sequencing reads were used for *de novo* contig assembly of VLP metagenomes. In total, 3.56 M contigs were assembled and pooled from the 25 individual metagenomes, with 80,762 contigs longer

6

155 than 5 kb (Figure 1A). VLP contigs longer than 5 kb were distributed evenly across the three

156 sample groups, totaling 2908.6 ± 1461.3, 3209.0 ± 2573.8, and 3535.7 ± 2826.4 contigs per

157 sample for CCP+, CCP- and HC respectively (Figure S1C).

158  These 80,762 contigs served as a starting point for identifying putative phages using a

159 three-pronged approach of independent phage discovery methods (Figures 1A and S1D). The

160 first method (P/M ratio) employed a previously validated read mapping strategy whereby VLP

161 read sets from all 25 samples were mapped to both whole metagenome (M) and VLP (P)

162 contigs [23]. Using the read-mapping P/M ratio (see Methods), we identified 2,117 unique

163 putative phage contigs after dereplication at 95% sequence identity. Next, we identified an

164 independent set of phage contigs by aligning all open reading frames (ORFs) of the 80,762 VLP

165 contigs against a set of 25,281 curated viral protein families (VPFs) [41]. Using this VPF

166 method, several filters were applied to identify viral contigs; (i) 2,902 contigs were identified as

167 having 5 or more VPF hits and non-viral Pfam hits below 20% of total ORFs on a contig, (ii) 263

168 contigs were identified with 5 or more VPF hits and less than 50% non-viral Pfam hits on a

169 contig, (iii) 644 contigs with 2-4 VPF hits and 0 non-viral Pfams, (iv) 976 contigs with at least 1

170 VPF hit, without considering any non-viral Pfams. In total, after dereplication, the viral contigs

171 arising from all above filters resulted in 4,785 unique viral contigs. For the third and final

172 approach we employed VIBRANT (Virus Identification By iteRative ANnoTation), a sequence-

173 independent algorithm that uses neural networks of viral protein signatures to identify lytic and

174 lysogenic phages [37]. Using VIBRANT, we identified 4,758 unique viral contigs.

175  To consolidate this list, we identified contigs that were shared between all three phage

176 discovery methods, resulting in a curated list of 660 contigs (Figures 1A and 1B). This curated

177 list of putative phage contigs range in size from 5,007 bp to 557,525 bp. To assess host

178 bacterial contamination among these contigs, we employed CheckV, a pipeline for assessing

179 the quality of viral genomes [42]. CheckV analysis revealed a reduced level of host bacterial

180 contamination and an increase of pure viral genomes in the final list of 660 curated contigs as

181 compared to varying levels of contamination among the three separate methods prior to contig

182 overlap identification (Figure 1C). We estimated completeness of our curated contigs using

183 CheckV and determined a greater distribution of "high quality" contigs relative to contig length,

184 in comparison to the three independent methods (Figure S2) [43]. Further, using the VIBRANT

185 platform for integrated provirus prediction, we describe communities of predominantly lytic viral

186 genomes belonging to Siphoviridae morphology (Figure S3). By using a combination of

187 approaches for viral contig discovery and assessing the overlap among these methods, we

188 have extracted a set of 660 predicted phages which are of overall high quality, both in terms of

189 viral contig completeness and lack of bacterial contamination than those from each of the

190 individual methods (Figures 1C and S2), which to date have been used primarily in isolation to

191 identify and characterize viral metagenomes.

192

193 **Clustering of metagenomic viral contigs reveals distinct viral ecological composition.**

194 Next we compared our set of curated contigs to over 2.3 million viral whole genome and

195 metagenome sequences from the IMG/VR database [44]. We used blastn at a threshold of 95%

196 sequence identity over 85% of 1 kb sequence length and Markov clustering to group our contigs

197 with related sequences from IMG/VR. Of the 660 contigs, 346 (52.4%) clustered into 255

198 clusters that contained 7,736 additional metagenomic viral contigs (mVCs) from IMG/VR. The

199 remaining 314 contigs (47.6%) were classified as singletons, with an even distribution among

200 CCP cohorts compared to healthy controls (Figure S4A). Of the curated contigs that were

201 clustered, cluster sizes ranged from 2 to 646 members with 78.4% of the groups containing

202 more than 2 partners and 36.5% containing more than 10 members, and 65.9% between 2 – 10

203 members (Figure S4B). Among these 255 clusters, 14 included reference prophages and lytic

204 phages, and 318 (48.2%) clustered with classified mVCs, thus assigning multiple levels of

205 taxonomy to our contigs (Figures 2A, 2B, and Supplementary Table 1).

206       Although host assignments were made using sequence-based clustering, host

207    specificity was further determined by aligning Clustered Regularly Interspaced Short

208    Palindromic Repeat (CRISPR) spacer sequences to our 660 curated contigs. CRISPR-Cas

209    serves as a snapshot of previous phage infections in the form of acquired spacer sequences

210    that represent invading viral genomes [45], and these sequences can be used for accurate

211    identification of phage-host interactions in intestinal microbiomes [23, 46]. CRISPR spacer host

212    assignments at the family level were present in 207 of 660 contigs (31.4%). All CRISPR spacer

213    queries considered for these analyses, ranging in length from 18 to 70 bp, were matches of

214    93.1–100% identity across the full length of the query and allowing for 0–2 mismatches and up

215    to 1 gap throughout [47] (Supplementary Table 2). Among predicted phages, total assigned

216    CRISPR spacers were evenly distributed, yet CCP+ sample containing phages predicted to

217    target Lachnospiraceae, Ruminococcaceae, Streptococcaceae, and Veillonellaceae bacterial

218    families were disproportionately abundant (Figures 2A and 2B). In total, 21 bacterial families

219    were identified as hosts via CRISPR spacer matching, supplementing the phage-host

220    interactions discerned from sequence-based clustering (Figure 2A). Among all samples in this

221    study, phages were predicted to target Lachnospiraceae, Ruminococcaceae, Clostridiaceae and

222    Bacteroidaceae bacteria with highest frequency of total CRISPR spacers (Figure 2A). Phage-

223    host interactions were also measured in terms of host range specificity, showing that while the

224    majority of the phages were predicted to have narrow host ranges, several spacers were linked

225    to multiple hosts across family level and higher taxa (Figure 2C), consistent with prior

226    observations of diverse viromes [47]. Broad host range phages were found across all cohorts,

227    but particularly among CCP+ sample contigs (Figure 2D) suggesting a more dysbiotic

228    community of host bacteria among these individuals' metagenomes.

229       We further explored the association of sample cohorts to phage hosts using read

230    mapping to determine differential host abundance profiles (Figure 3). Reads from all samples

231    were mapped to assembled phage contigs whose host assignments were deduced using

232    CRISPR-spacer matching and Markov clustering to quantify sequence abundances by

233    measuring cohort-based read recruitment [23, 48-50]. In comparing the differential read

234    recruitment to phages predicted to infect separate bacterial families, we observed differences

235    based on reads originating from either the CCP+ or CCP- groups in relation to the HC cohort

236    (Figure 3). Among the most striking, phage contigs targeting Bacteroidaceae recruited

237    significantly more reads from CCP+ viromes than either HC or CCP- individuals (Figure 3A). In

238    contrast, phages predicted to target Clostridiaceae bacteria were evenly abundant across all

239    three groups (Figure 3B). For Lachnospiraceae bacteria, CCP+ phages recruited were evenly

240    distributed among the groups with a slight elevation in CCP+ individuals that was not statistically

241    significant (Figure 3C). Ruminococcaceae phages were significantly skewed when comparing

242    HC to CCP- individuals (Figure 3D) and a major shift in phage read recruitment abundance was

243    evident for Streptococcaceae phages, as a greater percentage of total CCP+ reads were

244    mapped to these phages in relation to either HC or CCP- virome reads (Figure 3E). This skew

245    among CCP+ individuals is supported by prior works showing elevated Streptococcal phage

246    abundances in intestinal viromes of humans with inflammatory bowel disease [32] and a murine

247    model of colitis [23]. Lastly, no significant differences were observed for read recruitment to

248    Veillonellaceae-targeting phages (Figure 3F). Thus, differences in the host specificities were

249    evident between CCP+, CCP-, and HC groups with respect to read mapping abundance profiles

250    for Bacteroidaceae, Ruminococcaceae, and Streptococcaceae phages.

251

252    **CRISPR spacer host metadata reveal CCP+ phages represent greater variability in**

253    **microbial host ecology.**

254    To further explore the phage ecology from our subject cohort, we analyzed the host and mVC

255    metadata from the Joint Genome Institute's (JGI) Genomes OnLine Database (GOLD) [51]. The

256    JGI GOLD database contains metadata from over 100,000 biosamples and over 350,000

257    sequencing projects involving genomic and metagenomic sequencing data from biological

10

258     isolates worldwide. Moreover, recent work has contributed an additional 52,515 metagenome-

259     assembled genomes from diverse ecologies and geographic distributions [52], further

260     enhancing microbial host ecosystem analysis. Using the GOLD Biosample Ecosystem

261     Classification system, we analyzed the ecosystem distributions for all CRISPR spacers

262     identified in our curated contig list and discovered that the majority of host assigned contigs fell

263     within four distinct ecosystem classification levels; from broad to specific environments: host-

264     associated, human-associated, digestive system, and large intestine (Figure 4). For phages that

265     were previously identified as having CRISPR spacer host assignments, total spacer alignments

266     as identified by blastn ranging from 1 to 825 per contig, were tallied and used to calculate the

267     uniformity of spacer origins per contig (Supplemental Table 3). For each of the four ecosystem

268     categories, the most abundant classifications were used to compare across the study cohorts.

269     At the highest order GOLD Ecosystem distribution, the host-associated (i.e., human, mammal,

270     plant, arthropod, fungi) origin classification per contig was comparable for the HC and CCP-

271     groups but not for the CCP+ group (Figure 4A). A similar pattern was evident at the lower order

272     metadata distributions, with phage contigs derived from CCP+ individuals being more divergent

273     from the other cohorts for contigs of human-associated origin (Figure 4B), digestive system

274     origin (Figure 4C), and large intestine origin for the Ecosystem Subtype (Figure 4D).

275         These compositions of multiple CRISPR spacer ecosystem distributions reveal

276     homogeneity among phages derived from HC and CCP- samples, and indicates more dysbiotic

277     communities across CCP+ samples, suggesting that CCP+ individuals harbor disparate phage

278     communities that are more likely to originate from non-host associated sources. The putative

279     origins of these phages are related to environmental metadata of CRISPR spacers in the JGI

280     GOLD database describing the origin of bacterial DNA samples across ecologically diverse

281     biomes worldwide [52]; and increased heterogeneity in the CCP+ phages suggests a condition-

282     dependent host intestinal environment that maintains diversity. At the highest Ecosystem

283     classification level, with only three unique classification terms, these non-host associated

284  sources that are more prevalent in the CCP+ group, correspond to a higher degree of spacers

285  matching organisms originating from environmental and/or engineered habitats as archived in

286  GOLD (Figure S5). The ecosystem distributions of Category, Type, and Subtype have 43, 126,

287  and 146 unique terms for each classification level respectively, indicating multiple possible

288  combinations for organism habitats. Thus, our analysis of GOLD metadata for all phages with

289  predicted host isolates within our study reveals divergent habitat origins for CCP+ derived

290  contigs.

291

292  **Quantitative read mapping reveals differentially abundant contigs despite sample**

293  **cohesiveness.**

294  We next asked whether certain phage community members are present in different abundances

295  among the members of the cohort at-risk for rheumatoid arthritis compared to healthy controls.

296  To assess differences between phages among the sample groups, we used a viral read

297  recruitment strategy whereby VLP reads from all samples were mapped to the 660 curated

298  contigs [23, 48]. Using read count matrices for all contigs as input in the DEseq2 statistical

299  package for differential analysis of comparative count data [53], we analyzed three pairwise

300  comparisons for over- or under-abundant viral contigs (Figure 5). Initial comparisons of the

301  normalized and log-transformed count matrices were performed to evaluate the experiment-

302  wide trends across all samples. Principal component analyses reveal minimal variance

303  explained by the first two principal components for CCP+ vs HC samples (Figure 5A), CCP- vs

304  HC samples (Figure 5B), and CCP+ vs CCP- samples (Figure 5C), indicating that total sample

305  community signatures cannot be readily differentiated based on at-risk or healthy control

306  cohorts. We further explored the sample similarities by comparing Euclidian sample-to-sample

307  distances of the regularized log-transformed count matrices. Hierarchical clustering of sample-

308  to-sample distances did not reveal any discernable clustering for CCP+ vs HC samples (Figure

309  5D), and only minimal similarities between two CCP- samples when compared to the HC

310    (Figure 5E) and CCP+ (Figure 5F) groups, suggesting general sample cohesiveness between

311    cohorts.

312        We next analyzed specific members of the intestinal phage community, considering the

313    rationale that samples with complex communities are better explored at the level of each unique

314    member [33]. Visualization of the principal components incorporating the viral identification

315    metrics used in the VIBRANT neural network for our 660 curated contigs shows minimal

316    differentiation among phage scaffolds based on scaffold quality (Figure S6A) or predicted phage

317    state (i.e., lytic or lysogenic) (Figure S6B), although fragmentation of smaller sized contigs is

318    evident for both analyses. Further, grouping of contigs at the sample type level does not

319    differentiate any specific cluster (Figure S6C), which is consistent with the minimal variance

320    observed at the sample level (Figures 5A, 5B, and 5C). Finally, we assessed the differential

321    abundance of read recruitment counts for the set of 660 contigs and estimated fold changes

322    based on the negative binomial generalized linear model provided by DESeq2 [53]. Using

323    thresholds of log2-fold change greater than 1 or less than -1 (equivalent to fold change of $\pm$ 2)

324    and Benjamini-Hochberg adjusted $p$-values < 0.001, we identified a total of 178 differentially

325    abundant contigs (27% of the 660 phages) across three pair-wise abundance comparisons. For

326    CCP+ vs HC samples a total of 59 contigs (30 over- and 29 under-abundant) (Figure 5G), for

327    CCP- vs HC a total of 66 contigs (27 over- and 39 under-abundant) (Figure 5H), and for CCP+

328    vs CCP- a total of 53 contigs (27 over- and 21 under-abundant) (Figure 5I) passed our

329    thresholds for significance. This suggests that there are unique changes in select phage

330    abundances from the intestinal viromes of individuals at risk for RA, and that these changes are

331    more nuanced than sample-based community associations can reveal. These data indicate that

332    these cohort groups represent minimal sample-sample variation, but may provide clues related

333    to detection of biomarkers via specific community members. The top phage contigs associated

334    with either CCP+ or CCP- individuals were *Clostridium scindens* (Lachnospiraceae) and

335    *Actinomyces oris* (Actinomycetaceae), respectively, over-abundant at $\log_2$ fold changes of 25.9

336    and 23.5 compared to the healthy control samples.

337           A comparison of the bacterial relative abundances via 16S amplicon sequencing

338    confirmed an expansion of Lachnospiraceae bacteria among samples originating from CCP+

339    individuals (Figure S7A). This confirms, in part, observations of over-abundant

340    Lachnospiraceae-targeting phage contigs for the CCP+ but not CCP- cohorts (Figures 6B and

341    6C). The bacterial composition across all cohorts was relatively even in terms of richness

342    (Figures S7B and S7C), evenness (Figure S7D), and species diversity (Figure S7E).

343    Conversely, phage host abundances in the CCP- cohort relative to healthy controls were not

344    correlated to a family-level differentiation in bacterial taxa relative abundance.

345

346    **Phage auxiliary metabolic gene abundances highlight cohort-associated disparities in**

347    **metabolic potential.**

348    To determine the functional potential and metabolic capabilities within intestinal phages, we

349    quantified AMGs assigned to specific metabolic pathways in the Kyoto Encyclopedia of Genes

350    and Genomes (KEGG) database across at-risk and healthy cohorts. Since their identification as

351    viral drivers of host metabolism [54], phage-encoded AMGs have been recognized as

352    consequential actors that redirect host functional capacities thereby directly influencing local

353    ecology [55, 56]. Analyses of AMGs using VIBRANT and KEGG pathway annotations can

354    provide valuable insights into potentially altered metabolic functions or informative biosignatures

355    for cohort-associated microbial communities [37, 57]. To this end, we assessed our set of

356    curated phage contigs against 2,835 AMGs with KEGG annotations identified as "metabolic

357    pathways" or "sulfur relay system" [37]. Among our 660 phage contigs, 161 (24%) were found to

358    encode at least 1 AMG, with 252 AMGs in total across all samples (Supplemental Table 4).

359    Phages originating from the HC cohort accounted for 131 metabolic signatures, while CCP+ and

360    CCP- had less total AMGs with 77 and 44, respectively (Figure S8A). Among the most

14

361  represented metabolic categories across all phages, amino acid metabolism and the

362  metabolism of cofactors and vitamins contained 121 and 88 AMGs, respectively, with energy

363  metabolism being the next largest category with 22 AMG hits (Figure S8B). These general

364  pathway results indicate that phages in the intestine presumably affect host metabolism through

365  the consumption of metabolic resources needed for their own biogenesis, as described in

366  phage-host infection studies of model pathogens [58-60] and marine virocells [61].

367      To further probe all metabolic phage-encoded functions corresponding to our sample

368  cohorts, we assessed all AMG hits for total KEGG pathway abundances. Hierarchical clustering

369  grouped AMGs into 5 distinct metabolic clusters relative to HC and at-risk CCP cohorts (Figure

370  6A). Among these groups, the gene coding for *phnP* (K06167) stands apart from the others,

371  both in terms of clustering and also for relative pathway abundance (Figure 6A). Among group-

372  associated differences in AMG pathway abundances, there are notable absences among both

373  CCP+ and CCP- individuals. Namely, several clustered transferases such as the mannose-

374  phosphate transferases (*algA, xanB, rfbA, wbpW, pslB*), manno-heptose transferases (*gmhC,*

375  *hldE, waaE, rfaE*), and the *galE* epimerase and *glmS* transaminase (Figure 6A). Considering the

376  impact of such transferases on bacterial cell wall polysaccharides and biofilm formation [62, 63],

377  these results point to a baseline of phage-driven bacterial surface modifications from HC-

378  derived phages. Conversely, AMGs involved in lipopolysaccharide (LPS) biosynthesis such as

379  the *waaL* O-antigen ligase and the *gmhB* phosphatase are only present in CCP+ phages or at

380  greater abundance in CCP+ phages, respectively, indicating a possible role in immune evasion.

381  Within the CCP- cohort, one of the most abundant AMGs, KEGG orthology entry K23144

382  encoding for a polyketide sugar transferase important in peptidoglycan biosynthesis is

383  completely absent from the HC cohort and present at lower levels for CCP+ samples. Thus,

384  phage-encoded bacterial surface modifying enzymes such as the sugar transferases and

385  LPS/peptidoglycan biosynthetic genes, are differentially represented across the cohorts in this

15

386    study, which has implications for bacterial fitness in the intestinal ecosystems and their

387    interactions with the immune system.

388        We next incorporated the AMG characterization of genomes within our curated set of

389    phages to those that were significantly over- or under-abundant in previous differential

390    abundance analyses (Figures 5G, 5H, and 5I). Among the 20 differentially abundant contigs

391    from the CCP+ vs HC pairwise comparison that contained CRISPR spacer-predicted hosts, 8 of

392    these encoded at least one AMG (Figure 6B). The 9 under-abundant phages in this comparison

393    encode 5 AMGs, including manno-heptose transferases (*gmhC, hldE, waaE, rfaE*), mannose-1-

394    phostphate transferases (*algA, xanB, rfbA, wbpW, pslB*) and *ahbD* AdoMet-dependent heme

395    synthase all together on 1 contig, and *cysH* and *iscS* genes on 2 other contigs (Figure 6B).

396    Among the 11 significantly over-abundant contigs, 3 of these encode the *phnP*

397    phosphodiesterase; 3 phages predicted to infect *Flavonifractor* sp. (Ruminococcaceae) and one

398    predicted to infect Clostridiales bacteria. The remaining AMG found in CCP+-associated over-

399    abundant phages encodes for a cobalamin biosynthesis protein *cobS*, found in marine

400    cyanophages [64], viruses of marine archaea [65], and is considered a core component of

401    marine phage genomes [66], but also ubiquitous in phage genomes that infect *E. coli* [67]. Our

402    identification of a CCP+ over-abundant phage contig targeting *Bacteroides fragilis* and carrying

403    the *cobS* AMG (Figure 6B) reinforces the universal nature of this central AMG that is conserved

404    across hosts and environments [37].

405        We also identified 16 unique phage contigs with definitive CRISPR spacer-predicted

406    hosts that were differentially abundant and associated with the CCP- cohort (Figure 6C). Within

407    these contigs, 9 are significantly under-abundant compared to healthy controls, with 3 of these

408    encoding AMGs. CCP- associated phages were identified as carrying *cobS*, *DNMT3A*, *thiF*, and

409    *iscS* metabolic genes (Figure 6C). Thus, in contrast to CCP+ associated contigs which harbored

410    *phnP* and *cobS* on a combination of Lachnospiraceae, Rumminococcaceae, and

411    Bacteroidaceae targeting phages, CCP- associated phages were identified to target primarily

412    Bacteroidaceae and *Actinomyces oris* and harbor a combination of AMGs.

413

414    **DISCUSSION**

415    RA is a complex disease with an unknown etiology that puts a burden on quality of life resulting

416    in a strong societal impact [68, 69]. In addition to multiple epidemiological factors being

417    associated with RA development, including genetic and familial risk, environmental risk factors

418    and biological sex [3], the microbiota remains an important and understudied factor that likely

419    influences RA autoimmunity [70]. Given the widespread occurrence and diversity of phages in

420    the human intestinal microbiota and their impact on intestinal microbial ecology during health

421    and disease [19, 20, 71], we analyzed this previously neglected component of the microbiota as

422    it relates to RA etiopathogenesis. We used shotgun metagenomics to identify intestinal phages

423    of individuals at risk for developing RA and discovered an association of distinct phage

424    communities with RA-specific serology in the at-risk population.

425    Using three separate database-independent approaches, we describe a collection of 660

426    phage genomic sequences, their potential metabolic capability, and their differential abundance.

427    Through a combination of CRISPR spacer matching and Markov clustering with other viral

428    metagenomic sequences from diverse environments, we predicted host assignments for 285 or

429    43.2% of these phages, which is a high level of taxonomic assignments relative to recent

430    reports of approximately 10 – 30% host assignment identification [23, 48, 72]. By analyzing a

431    core set of *de novo* assembled phage contigs paired with taxonomy, we identified differential

432    phage communities associated with the at-risk RA individuals compared to healthy controls, all

433    while adding novel phage-host assignments to previously unidentified intestinal phages [73, 74].

434    Phage-host assignments were dominated by Lachnospiraceae-targeting phages, some

435    of which were over-abundant in CCP+ individuals. This expansion of phages also correlated

436    with increased abundances of Lachnospiraceae bacteria in the CCP+ cohort compared to either

17

437    CCP- or the healthy cohort, suggesting a link to this family of Firmicutes and CCP autoantibody

438    production in the human intestine. Interestingly, increased abundance of Lachnospiraceae has

439    been observed in at least two previous studies of intestinal microbiotas in mice during the

440    course of collagen-induced arthritis (CIA) [75, 76]. Considering the precedence for overlap of

441    identified phage contigs from mouse intestines to human-associated intestinal phages [23], the

442    previously-reported increase in abundance of Lachnospiraceae bacteria during experimental

443    arthritis in mice is supported by our findings of increased Lachnospiraceae phage-host

444    interactions in CCP seropositive individuals at-risk for developing RA. To this end, given that the

445    FDR individuals included in this study do not show clinical signs of established RA, our

446    identification of a preclinical cohort with increased Lachnospiraceae phage-host interactions

447    could serve as a biological indicator of disease. Similarly, an expansion of Bacteroidaceae-

448    targeting phages associated with the CCP- cohort was described, which corresponds to a

449    previously observed expansion of Bacteroidaceae bacteria following CIA in mice [75]. In

450    addition to these phages serving as potential biomarkers of disease in humans at risk for RA,

451    our data indicate that Bacteroidaceae and Lachnospiraceae-targeting phages designate a

452    distinction between CCP serology status that may serve as an additional indicator of disease

453    progression and/or future disease severity [77]. Notably, bacteria in the Lachnospiraceae and

454    Ruminococcaceae families have been linked to the pre-diabetic intestinal microbiota and

455    diabetic pathogenesis, while Bacteroidaceae are associated with disease protection in a murine

456    model of diabetes [78]. The identification of cohort-specific phage-host interactions sheds light

457    on potential preclinical biomarkers connecting specific dysbiotic intestinal microbial communities

458    to possible regulation of microbiota-mediated mucosal inflammation [1, 79].

459    We calculated the differential abundance of curated phages on a contig-to-contig basis

460    to estimate dispersion and fold changes of quantitative read mapping matrices. In doing so, we

461    identified 178 differentially abundant contigs (27% of the total curated list) across three pair-wise

462    cohort comparisons. Among the CCP+ vs HC comparison, we observed over-abundant phages

463    targeting *Clostridium scindens, Flavonifractor* sp*., Actinomyces oris*, as well as other family-level

464    taxonomic assignments. A member of the Lachnospiraceae, *C. scindens* is an intestinal

465    commensal bacterium involved in maintaining homeostatic large intestinal bile acid composition

466    and providing host protection from opportunistic *Clostridioides difficile* blooms [80, 81]. A

467    differential abundance of phage targeting *C. scindens* in the CCP+ at-risk cohort, may have

468    implications for bile acid dysmetabolism in these individuals, which has consequences for

469    inflammatory bowel diseases [82, 83]. Differential abundance of phages in the CCP- cohort

470    revealed several phages targeting Bacteroidaceae and *Bacteroides* species, bacteria involved

471    in multiple reactions of bile acid metabolism promoting host metabolic health [84, 85]. Recent

472    phage-*Bacteroides* interactions have described the influence of phage BV01 in reducing

473    *Bacteroides* bile acid metabolism [86], which has implications for the impact of phages on

474    mammalian gut metabolic function. Our findings suggest individuals at risk for RA harbor

475    divergent communities of phages with potential to alter intestinal metabolic potential through

476    either reduction of key bacterial species and thus reducing endogenous metabolic function, or

477    through the phage-derived introduction of specific AMGs.

478         Changes to the intestinal metabolome can lead to compositional microbiota transitions

479    that in turn impact host nutrient uptake and immune homeostasis [87]. Considering that

480    manipulations of microbial metabolic pathways in the intestine can influence inflammation and

481    dysbiosis [88], our identification of phage communities with differential abundances of encoded

482    AMGs points to divergent metabolic landscapes associated with at-risk RA cohorts. A majority

483    of the AMGs identified in our analysis make up a group of 14 genes conserved across many

484    environments [37], indicating their functional importance in core metabolism. We were surprised

485    to identify three phages that were over-abundant in the CCP+ cohort (3 of 11 in total), three with

486    *Flavonifractor sp.* predicted hosts and one Clostridiales-targeting phage, encoding the *phnP*

487    phosphodiesterase. Encoding a phosphoribosyl 1,2-cyclic phosphate phosphodiesterase, *phnP*

488    accounts for 10% of the total AMGs represented in our phage genomes, and is differentially

19

489    abundant among the CCP+ cohort samples. While *phnP* is one of 14 genes considered to be

490    globally conserved across multiple environments [37], it is the only gene among AMGs in our

491    analysis that is the lone representative of its pathway. The PhnP phosphodiesterase, part of a

492    14-gene operon originally described in *Escherichia coli,* plays a crucial intermediary role in the

493    carbon-phosphorous lyase pathway by degrading a dead-end cyclic phosphonate byproduct

494    [89]. The uniform presence of *phnP* across phages derived from at-risk and healthy cohorts

495    (Figure 6A), suggests phage-driven organophosphonate degradation, which is fundamental for

496    bacteria in diverse environments [90].

497        Phosphonate degradation is important for phosphorus assimilation in enteric bacteria

498    [91], although phosphonate metabolism has not been described for *Flavonifractor* species and a

499    *phnP* homolog is not available for this genus in the KEGG database (K06167). In a recent study

500    characterizing microbiota KEGG orthologs as predictors of methotrexate responsiveness for RA

501    treatment, a gene in the phosphonate transport system, *phnC* (K02041), exhibited high median

502    random forest importance as a predictor of drug response in new-onset RA subjects [92]. The

503    contribution of the phosponate metabolic pathway in bacteria and phages, will require further

504    exploration in the context of RA pathogenesis and treatment. However, it is possible that these

505    phage-encoded metabolic products are supplementing phosphorous uptake among

506    Ruminococcaceae and Lachnospiraceae bacteria that predominate in CCP+ individuals prior to

507    RA clinical symptoms. Our analysis is limited in that we did not measure a longitudinal

508    progression of microbial metabolic pathways in these human samples, yet these metabolic

509    associations warrant further investigations into causality and the potentially cascading effects on

510    interbacterial interactions [93].

511        Our results point to divergent communities of phages with multiple bacterial host targets

512    that group according to anti-CCP serology in individuals predisposed to developing RA. These

513    at-risk individuals who develop seropositive RA, a disease manifestation that is more severe

514    [94] and less responsive to treatment [95], endure a prolonged asymptomatic period before

515   pathological early RA develops in those who are at a higher disease susceptibility in the

516   preclinical RA state [1]. Current approaches for RA diagnosis rely in large part on anti-CCP

517   serology which has up to 93% specificity but as low as 67% sensitivity (for the CCP3.1 assay

518   used here) [39], indicating that a negative result does not preclude current or development of

519   clinically apparent RA. Phage community composition analyses may complement existing

520   diagnoses for RA, considering that intestinal phages can play important roles in immune

521   tolerance, mucosal immunity, and microbial homeostasis [96]. Given that phage community

522   alterations have been shown to precede autoimmunity development in children at risk for

523   developing type 1 diabetes [26], phage community structure should be considered as a

524   biomarker for diseases such as RA that are influenced by non-genetic microbial factors [19]. To

525   that end, we have characterized the intestinal viromes of RA at-risk individuals corresponding to

526   anti-CCP serology status. Furthermore, we calculated species-specific phage-host interactions

527   and identified over-abundances of *C. scindens* and *A. oris* targeting phages in CCP+ and CCP-

528   individuals, respectively. Divergent metabolic profiles evident by differential abundance of AMG-

529   encoding phages in both conditions warrant further interrogation during models of RA-like

530   disease. Future work should investigate the potential of phages in a murine CIA model to

531   determine the influence of RA-associated phages with immunomodulation and inflammatory

532   disease progression. Our multifaced approaches for phage prediction and phage host

533   assignments hold promise to better ascertain the occurrence and diversity of the virome and the

534   identification of key phages influencing the microbiota and individuals at risk for developing RA

535   autoimmune disease. This RA-focused study implicating specific phage populations could open

536   new avenues to assess the basis for phage implication in other microbiota dysbiosis-associated

537   diseases.

538

539   **RESOURCE AVAILABILITY**

540   **Lead Contact**

541     Further information and requests for resources and reagents should be directed to and will be

542     fulfilled by the Lead Contact, Breck A. Duerkop (breck.duerkop@cuanschutz.edu).

543

544     **Materials Availability**

545     This study did not generate new unique reagents.

546

547     **Data and Code Availability**

548     The VLP and whole metagenome DNA sequencing reads as well as the final curated phage

549     contigs generated in this study are available at the European Nucleotide Archive under the

550     Study titled "Intestinal VLP reads and predicted phage contigs for at-risk RA individuals"

551     (accession numbers PRJEB42612 and ERP126498). The VLP and whole metagenome raw

552     unmapped read sets are available for each of the 25 individual samples included in this study

553     and are available under the Study Primary Accession PRJEB42612. The 660 curated contigs

554     are compiled in a multifasta file deposited as Sample SAMEA7856466 under the same Study

555     PRJEB42612.

556

557     **EXPERIMENTAL MODEL AND SUBJECT DETAILS**

558     **Study Subjects and Fecal Samples**

559     Fecal samples were obtained from individuals recruited for the SERA (Studies of the Etiology of

560     Rheumatoid Arthritis) initiative, aimed at understanding the mechanisms that prelude the

561     preclinical development of RA. SERA is a multicenter prospective cohort study that has

562     identified first-degree relative (FDR) probands defined as a parent, full sibling, or offspring of

563     individuals with diagnosed clinical RA [38]. FDR probands were evaluated in extensive clinical

564     research visits, longitudinal follow-ups, and autoantibody testing to determine CCP status [38].

565     FDR probands were split into cohorts dependent on serum CCP levels, with 100% of subjects in

566     the CCP+ cohort positive and 0% of subjects in either CCP- or HC (Healthy Control) cohorts

567    testing positive. Healthy control subjects were recruited and included in this study as described

568    previously [97]. The present study consisted of 25 subjects split into 3 cohorts, of which 8 were

569    CCP+, 8 were CCP-, and 9 were HC. Ethical approval for this study was obtained from the

570    University of Colorado Multiple Institutional Review Board (COMIRB) study numbers 01-675

571    (primary) and also 13-2606 and 14-1751. COMIRB Protocol 01-675 included informed consent

572    with HIPAA authorization for stool sample collections. Stool samples were obtained

573    independently by SERA study participants and returned within 1 week of their original visit.

574    Samples were stored in aliquots at -20°C until processing.

575

576    **METHODS**

577    **Extraction of Fecal Whole Metagenome and VLP DNA, Library Preparation and**

578    **Sequencing**

579    Whole metagenome and VLP fraction DNA were isolated as described previously [98], with

580    some modifications as follows. For all samples, 0.1 g of human stool was homogenized in 8 mL

581    salt magnesium plus (SM+) buffer [99] and 0.5 ml of homogenate was transferred to a

582    BashingBead Lysis tube (Zymo) and designated as the whole metagenome sample. Whole

583    metagenome DNA was extracted using a ZymoBIOMICS DNA kit (Zymo) following the

584    manufacturer recommended protocol. VLPs were clarified from the remaining 7.5 mL of sample

585    by three successive centrifugation steps (3200g for 10 min, 3200g for 10 min, 7800g for 10

586    min), and the supernatant was filtered through a 0.45-µm PVDF filter membrane. VLPs were

587    precipitated by adding 0.5M NaCl and 10% wt/vol PEG8000 and incubating on ice at 4°C

588    overnight, followed by centrifugation (7800g for 20 min). VLP pellets were resuspended in 400

589    µL SM+ buffer and treated with 40 µL DNase buffer (10 mM $CaCl_2$, 50 mM $MgCl_2$), 25 units

590    DNase, and 15 units RNase for 1 hr at 37°C. VLPs were further treated with 50 mg/mL

591    proteinase K and 0.5% SDS for 30 min at 56°C before addition of 100 µL phage lysis buffer (4.5

592    M guanidiniumisothiocyanate, 44 mM sodium citrate pH 7.0, 0.88% sarkosyl, 0.72% 2-

23

593    mercaptoethanol) and incubated for 10 min at 65ºC. VLP DNA was precipitated and extracted

594    with an equal volume of phenol/chloroform/isoamyl alcohol 25:24:1, spun at 7800g for 5 min,

595    and further extracted with an equal volume of chloroform. VLP DNA was precipitated with 0.3M

596    NaOAc (pH 5.2) and an equal volume of isopropanol, washed with ice-cold 70% ethanol, and

597    resuspended in sterile water.

598

599    **Metagenomic DNA Sequencing**

600    VLP and whole metagenomic DNA was sequenced using the Illumina NovaSEQ 6000 platform

601    with paired-end 150-cycle sequencing chemistry. DNA libraries were amplified using the

602    Ovation Ultralow System v2 (Nugen, part no. 0334) library preparation kit including 12 cycles of

603    amplification. TruSeq adapters (Illumina) were used for multiplexing. Libraries were quantified

604    using a Qubit and quality was measured using a Tapestation. All library preparation,

605    quantification, quality assessment and control, were performed by the University of Colorado

606    Anschutz Medical Campus Genomics and Microarray Core.

607

608    **16S rRNA Amplicon Sequencing and Analysis**

609    16S rRNA gene analysis was performed using fecal samples that were processed for isolation

610    of whole metagenomic DNA using a ZymoBIOMICS DNA kit (Zymo) and stored at -80°C.

611    Amplicons of the 16S rRNA gene V4 region were amplified using Earth Microbiome Project

612    primers 515F and 806R [100] with custom barcodes. Samples were sequenced on the Illumina

613    MiSeq platform with paired end 250 bp reads using bTEFAP technology [101] by MR DNA

614    (Molecular Research LP, Shallowater, TX), and processed using mothur v.1.44.0 [102].

615    Sequenced reads, which averaged 607,915 $\pm$ 112,641.7 per sample, were demultiplexed,

616    assembled as contigs, and processed to remove chimeras and erroneous sequences per the

617    Kozich protocol [103] and mothur MiSeq SOP (https://mothur.org/wiki/miseq_sop/ accessed

618    07/16/2020). Sequences were aligned to the Greengenes core reference alignment for

619    taxonomy using the 2013 release (gg_13_8_99) [104]. Sequences were differentiated into

620    amplicon sequence variants (ASVs) using the make.shared command, resulting in a total of

621    8,108,071 sequences. Subsampling was performed using 186,745 sequences, which

622    corresponded to the smallest sample in our dataset. Diversity measurements were performed

623    using mothur calculators to estimate community richness (Chao1 estimator), community

624    evenness (Shannon evenness), and community diversity (inverse Simpson index).

625

626    **Decontamination and Read Processing**

627    Metagenomic reads were decontaminated and trimmed as previously described [23] using

628    BBMap short read aligner v38.56 [105]. Briefly, raw reads were mapped to the internal Illumina

629    phage genome control phiX174 (J02482.1), human reference genome (hg38), and potential

630    laboratory contaminants including mouse genome (mm10), *Enterococcus faecalis* V583

631    genome (NC_004668.1), *E. faecalis* OG1RF genome (NC_017316.1), and *E. faecalis* phage

632    VPE25 (LT546030.1) using the bbsplit algorithm with default settings. Unmapped reads were

633    trimmed of adapter sequences, with low quality reads and reads of insufficient length removed

634    using the bbduk algorithm with the following parameters: ktrim = lr, k = 20, mink = 4, minlength =

635    20, qtrim = f, as previously described [23].

636

637    **Metagenomic Assembly**

638    Decontaminated and trimmed R1 and R2 reads were interleaved using the fq2fa --merge

639    command from the IDBA-UD package [106]. Whole metagenome and VLP assemblies were

640    performed using the MEGAHIT assembler v1.2.7 [107] using the default setting plus the

641    following flags: --presets meta-large (--k-min 27 --k-max 127 --k-step 10) for large and complex

642    metagenomes.

643

644 **Quantitative Read Mapping and Construction of the Curated VLP Contig Database**

645 VLP reads were assembled into 25 individual sample sets, corresponding to the 25 individual

646 fecal samples included in our study. All contigs resulting from MEGAHIT assembly were filtered

647 to a minimum length of 5kb, resulting in a pool of 80,762 total contigs from all samples. Three

648 separate independently published methods were employed to identify putative phages from the

649 pooled set of contigs over 5kb in length. First, the P/M read mapping approach was used

650 whereby each sample's VLP and whole metagenome reads were mapped to their

651 corresponding assembled contigs, using BBMap as previously described [23]. After pooling, the

652 top 100 largest ratios of VLP reads to whole metagenome reads for all 25 read-mapping sets for

653 each sample were identified and pooled. Redundancy was removed using cd-hit-est at an

654 identity threshold of 95% resulting in 2117 unique contig sequences. Next, as a separate

655 method, putative phages from the 80,762 contigs were identified by searching for viral protein

656 family (VPF) hits, as previously described [41]. Separate filters were applied for VPF hits

657 calculated in relation to total genes, microbial genes, and percent non-viral Pfams. 2,902 contigs

658 were identified that contained 5 or more VPF hits and with non-viral Pfam hits below 20%. 263

659 contigs were identified with 5 or more VPF hits, with more viral gene content than microbial

660 genes per contig, and 644 contigs were identified with 2 – 4 VPF hits and 0 microbial gene hits.

661 Finally, 976 contigs were identified with only 1 VPF hit per contig, and were included regardless

662 of microbial gene content. The third and final independent phage contig identification method

663 used was VIBRANT v1.2.1 [37], a neural network machine learning algorithm that identifies viral

664 protein signatures. VIBRANT identified 4,758 unique phage contigs. After combining these three

665 independent approaches used to identify unique sets of phages, all sets were combined and the

666 overlapping 660 contigs were used for analysis as the curated contig set. To assess contig

667 completion and contamination levels, CheckV v0.6.0 was used with standard operating

668 parameters.

669

**Differential Abundance Analyses**

To calculate differential abundance in pairwise analyses, we first generated read mapping count matrices by mapping all VLP reads to the curated contig set of 660 contigs. The bbmap algorithm from the BBMap suite of tools was used with the following parameters: ambiguous = random, qtrim = lr, minid = 0.97. Total raw read counts were aggregated per contig and assembled into 25 count matrices for all samples, which were then used as input for DESeq2 v1.20.0 [53] running in R version 3.6.3 for comprehensive differential abundance analysis. Raw un-normalized read count coverage values were used to compare fold changes across three pairwise comparisons: CCP+ vs. HC, CCP- vs. HC, and CCP+ vs. CCP- groups. The standard workflow for differential analysis within DESeq2 was used, producing logarithmic fold-change values incorporating Wald tests for $p$-value calculations and the Benjamini-Hochberg multiple testing correction for the adjusted $p$-value. In total, 178 phage contigs from our set of 660 were found to be differentially abundant using thresholds of $\log_2$ Fold Change < -1 or > 1 and adjusted p-value < 0.001.

**VLP Clustering, Phage Host Matching, and AMG Identification**

Clustering of all viral contigs within the RA virome described in this study was performed using two lists of contigs, the total 4,785 viral sequences identified by all filters of the VPF method, as well as the final curated set of 660 contigs. First, all 4,785 sequences were screened against the most recent iteration of the public viral database IMG/VR v3.0 [44] using blastn with 95% sequence similarity over 85% of each 1kb region of the contig, which resulted in 19,892 viral sequences. Then, a total of 24,926 sequences were screened against each other using blastn with the same parameters and omitting duplicate hits. Markov clustering of these 9.4 million connections resulted in a total of 1,193 clusters encompassing 22,306 total sequences. Overall, 2,420 of the 4,785 total RA virome sequences were clustered into 1,184 clusters. Of these clusters, 41 contained a reference viral isolate, 1,037 contained another metagenomic viral

27

696    contig from IMG/VR, and 106 were identified as originating from RA metagenomic sequencing

697    projects. Lastly, clustering was also calculated for the 660 curated viral sequences, which

698    resulted in 266 individual clusters containing 336 (roughly 48% of curated set) unique

699    sequences. Phage host assignments were determined via bacterial CRISPR spacer matching

700    as previously described [23], requiring at least 93% sequence identity match over the entire

701    spacer length and allowing for up to 2 mismatches. Of our 660 curated contig list, 207 (31.4%)

702    had CRISPR spacers matching reference isolates therefore leading to host predictions for a

703    third of our final contigs. VIBRANT v1.2.1 was used to identify auxiliary metabolic genes (AMGs)

704    according to KEGG metabolic pathway annotations. VIBRANT annotates using VOG, Pfam, and

705    KEGG databases; therefore, if the best HMM hit is to the KEGG database and also if the

706    annotation is in a metabolic pathway, the hit gets called as an AMG.

707

708    **Data Visualizations**

709    Various R packages were used, including DESEq2, ggplot2, ComplexHeatmap, pheatmap,

710    corrplot, RColorBrewer, and EnhancedVolcano. Graphpad Prism v8.4.3 was used for all

711    supplemental calculations. Lastly, SankeyMATIC (https://github.com/nowthis/sankeymatic) and

712    meta-chart (https://www.meta-chart.com/venn) were used to create the Sankey and Venn

713    diagrams depicted in Figure 1, respectively.

714

715    **FIGURE LEGENDS**

716    **Figure 1. Generation and curation of *de novo* assembled VLP contigs.** Metagenomic

717    sequencing was carried out for 25 samples belonging to 3 cohorts of individuals, FDRs at risk

718    for developing RA later in life with either CCP+ or CCP- serology status, and a Healthy Control

719    (HC) group. (A) Contigs were assembled *de novo* for all samples, ranging from 30,011 to

720    284,689 contigs per sample, and a total of 3,557,500 contigs for the entire sample set. Each

721    node on the Sankey diagram is scaled to the number of contigs it contains. Thresholds of

28

722 minimum contig sizes being greater than 1 and 5 kilobases reduced the total numbers to

723 564,228 and 80,762 contigs respective to the size cut-off. Three independent methods were

724 used to identify putative phages from the list of 80,762 contigs (boxed portion of panel A),

725 resulting in 2,117 contigs from the P/M ratio method, 4,785 contigs from the Viral Protein

726 Families method, and 4,758 contigs using the VIBRANT algorithm. (B) A Venn diagram was

727 created to show the overlap of redundant contigs identified among the three methods. 660

728 unique contigs were identified independently by all phage identification methods. (C) CheckV

729 analysis of the three separate methods as well as the final set of curated contigs revealed a

730 disparity in host contamination, with the set of 660 contigs being relatively free of host bacterial

731 contamination. Colors were assigned to the CheckV categories that account for prophage

732 elements based on their position on the contig sequence, as well as pure viral (green) and pure

733 bacterial (grey) classifications.

734

735 **Figure 2. Clustering with metagenomic viral contigs reveals viral ecological composition.**

736 (A) Host assignments for the set of curated phages based on Markov clustering to the IMG/VR

737 database metagenomic viral clusters or direct match to bacterial CRISPR spacers, based on

738 cohort abundance. Bacteroidaceae, Lachnospiraceae, Ruminococcaceae, and

739 Streptococcaceae hosts are evident to be cumulatively more abundant than other bacterial

740 families, especially for the CCP+ cohort. (B) Cladogram of the complete host phylogeny at the

741 genus level for all spacers identified from total RA virome via the VPF method. The pie chart at

742 the center represents all 958 CRISPR spacers from the family level quantified in panel A that

743 have been color coordinated on this cladogram as well. Total host hits were quantified at the

744 genus level and are represented in relative size by colored circles, indicating host assignments

745 that were discerned via clustering (dark grey) and those that were identified via direct CRISPR

746 spacer matching (light grey). Total CRISPR spacers per contig with family level host taxonomy

747    assignments were tabulated per cohort group (C) and differentiated as narrow or broad phage

748    host ranges (D) based on target uniformity to bacterial family.

749

750    **Figure 3. Phage-host assignments for curated VLP contigs reveal cohort-based**

751    **differential read recruitment among several bacterial families.** Relative abundances were

752    calculated for all VLP reads mapped to phages predicted to target Bacteroidaceae (A),

753    Clostridiaceae (B), Lachnospiraceae (C), Ruminococcaceae (D), Streptococcaceae (E), and

754    Veillonellaceae (F) bacterial families. For these analyses, VLP reads were mapped to predicted

755    phage contigs to which CRISPR spacers were assigned using bbmap at a 97% minimum read-

756    mapping identity level. Scaffold abundances were averaged across all samples and statistics

757    were determined by nonparametric Wilcoxon tests (* $p < 0.05$, ** $p < 0.01$, **** $p < 0.0001$).

758

759    **Figure 4. CRISPR spacer host metadata reveal CCP+ phages represent greater variability**

760    **in microbial host ecology.** Phage host isolate ecology metadata was compiled from

761    JGI/GOLD v7.0 and broken down by Ecosystem, Ecosystem Category, Ecosystem Type, and

762    Ecosystem Subtype distributions accordingly for all CRISPR spacers identified within our list of

763    660 phages. (A) Ecosystem Distribution showing the percent host-associated spacers

764    calculated for each contig based on cohort distribution. (B) Ecosystem Category distribution

765    showing the percent human-associated spacers. (C) Ecosystem Type distribution showing the

766    percent of contigs that contain spacers originating from the digestive system. (D) Ecosystem

767    Subtype showing the percent of contigs that contain spacers originating from the large intestine

768    microenvironment. Cohort distributions based on these metadata revealed a disproportionate

769    distribution of CRISPR spacers among samples originating from CCP+ individuals when

770    compared to CCP- or HC groups. Statistical significance was determined using pairwise

771    Wilcoxon rank sum tests for comparisons between the three groups, using the Benjamini-

772    Hochberg correction for multiple testing comparisons (* $p = 0.023$, **** $p < 2 \times 10^{-16}$).

773

**Figure 5. Quantitative read mapping exposes differentially abundant contigs despite sample cohesiveness.** Quantitative read mapping of all VLP read sets to the final curated 660 phages reveals contig-contig dissimilarities despite minimal sample-sample variance or intra-sample hierarchical clustering. Differential abundance calculations were carried out within the DESeq2 package by way of 3 pairwise comparisons: CCP+ vs. HC, CCP- vs. HC, and CCP+ vs. CCP-. (A, B, C) Analyses of the first and second principal components for sample-to-sample exploratory analyses revealed minimal variance explained across all comparisons. (D, E, F) Euclidian distances for sample-sample read-based coverages were used for hierarchical clustering across all pairwise comparisons reveal minimal clustering based on sample type. (G, H, I) Volcano plots reveal 9%, 10%, and 8% of contigs included in our analysis are differentially abundant respective to CCP+ vs. HC, CCP- vs. HC, and CCP+ vs. CCP- group-based comparisons of specific contig community members.

786

**Figure 6. Phage auxiliary metabolic gene abundances highlight cohort-associated disparities in potential metabolic function.** AMGs were identified within the VIBRANT algorithm, based on screening 2,835 auxiliary metabolic genes with KEGG Orthology annotations (March 2019 KEGG release, ftp://ftp.genome.jp/pub/db/kofam/archives/2019-03-20/). (A) Total counts per KEGG Pathway were used normalize relative abundance of AMGs per sample, which were clustered using the ComplexHeatmap package in R. Areas in black indicate no AMG hits were present for the entire cohort for the 660 contig samples. (B) Differentially abundant contig for the CCP+ to HC pairwise comparison, visualizing only the contigs which had CRISPR spacer-predicted hosts. Color-coded stars belong to a list of AMGs and indicate association with the contig they are adjacent to. (C) Differentially abundant contigs for the CCP- vs HC comparison.

798

31

799 **TABLES**

800 **Table 1. Summary of the Subjects' Characteristics for the Samples Included in the Study**

| VARIABLE | HC | CCP+ | CCP- |
|---|---|---|---|
| Count | 9 | 8 | 8 |
| Age (mean) | 44.4 | 61.3 | 49 |
| Age (SD) | 13.6 | 11 | 15.7 |
| Sex (% female) | 66.7 | 88.9 | 62.5 |
| Serum CCP+ (%) | 0 | 100 | 0 |
| Ever smokers (%) | 22.2 | 33.3 | 0 |

801

802 **SUPPLEMENTAL INFORMATION**

803 **Figure S1. Overview of methods for VLP isolation and phage identification from**

804 **sequencing reads.** (A) Individual stool samples were homogenized and split into P and M

805 subsamples for generating VLP and whole metagenome DNA, respectively. (B) Total

806 sequencing reads generated per sample for each P and M read sets, after quality control and

807 decontamination. (C) Total assembled contigs with length greater than 5kb generated per

808 sample for each P and M read sets. (D) Overview of the computational pipeline used to identify

809 phages; from short-insert pair end read sets averaging approximately 75M read pairs per

810 sample, to the 80,762 *de novo* assembled contigs greater than 5kb in length, and the three

811 independent methods for phage identification (P/M, VPF, VIBRANT).

812

813 **Figure S2. Estimation of contig completeness by CheckV.** Distribution of contig lengths

814 across contig quality categories according to the MIUViG standards. Contigs derived from the

815 (A) P/M ratio method of phage identification, (B) the VPF method, (C) VIBRANT algorithm, and

816 finally (D) the curated contig list. Boxplots depict the following five summary statistics: median,

817 lower and upper hinges corresponding to the first and third quartiles, and two whiskers

818 corresponding to 1.5 times the interquartile range between the first and third quartiles. Points

819 beyond the whiskers correspond to outlier points.

820

32

821 **Figure S3. Lifestyle and morphology distributions of curated phage contigs.** (A) Total

822 contigs per sample among the three groups, divided according to infection mechanism (lytic vs.

823 lysogenic) as determined by the VIBRANT algorithm. (B) Relative abundances of phage

824 lifestyles as determined by the VIBRANT algorithm. In total, for our 660 predicted phages, 467

825 (70.8%) are classified as lytic and 193 (29.2%) are classified as lysogenic by VIBRANT. (C)

826 Viral taxonomy of all contigs per sample including the top four morphotypes: Siphoviridae,

827 Myoviridae, Podoviridae, and Microviridae. (D) Relative abundance of all viral morphotypes

828 identified for all 660 phages. Viral taxonomy was determined using a custom database

829 described in this preprint by Kieft et al. (2020; bioRxiv preprint doi:

830 https://doi.org/10.1101/2020.08.24.253096).

831

832 **Figure S4. Clustering distribution of singletons and viral groups.** (A) Total singletons, viral

833 contigs that did not cluster with any other genome, per sample and RA cohort group. (B)

834 Distribution of total viral clusters in relation to the number of viral genomes clustered within each

835 group.

836

837 **Figure S5. CRISPR spacer host metadata distribution of environmental and engineered**

838 **derived phages per cohort.** Phage host isolate ecology metadata was compiled from

839 JGI/GOLD v7.0 at the highest Ecosystem classification level for all CRISPR spacers identified

840 within our list of 660 phages. Data is presented as percent of spacers per contig whose

841 metadata is designated as originating from (A) environmental or (B) engineered environments,

842 distributed across the three RA cohort groups. Statistical significance was determined using

843 pairwise Wilcoxon rank sum tests for comparisons between the three groups, using the

844 Benjamini-Hochberg correction for multiple testing comparisons (* $p = 0.011$, **** $p < 2 \times 10^{-16}$).

845

**Figure S6. Principal component analyses based on quality, predicted phage lifestyle, and sample cohort.** Principal components for the final curated set of 660 contigs derived from the VIBRANT phage identification program categorized by (A) contig quality, (B) phage lifestyle, and (C) cohort group. Total identified open reading frames were incorporated in analyses in (A) and (B), showing a greater dispersion of smaller sized contigs and a consensus grouping of bigger contigs.

**Figure S7. Analysis of bacterial family diversity from fecal samples based on 16S sequencing and analyzed using mothur.** (A) Relative abundances of bacterial families based on ASV binning reveals a significant difference in Lachnospiraceae bacteria originating from CCP+ fecal DNA samples. Unpaired nonparametric Mann-Whitney tests were used to compare ranks, revealing $p$ values of 0.0464 comparing CCP+ to HC individuals. Community richness was measured by the standard observed richness calculator in mothur (B) as well as the Chao1 richness estimate (C). Community evenness was measured using the Shannon index (D), and community diversity was measured using the inverse Simpson index (E). No statistically significant differences were observed among any of the above calculators using nonparametric tests of significance among the three groups.

**Figure S8. Distribution of auxiliary metabolic genes found on curated contigs.** (A) A total of 252 AMGs were discovered among our 660 phages, distributed across the three cohorts. (B) AMGs were categorized predominantly as belonging to amino acid and cofactor/vitamin metabolism categories.

878

**AUTHOR CONTRIBUTIONS**

880    Conceptualization, M.R.M, D.P., A.C., K.A.K., and B.A.D.; Methodology, M.R.M., D.P., K.K.,

881    A.C., J.A.S., M.L.F., M.K.D., and B.A.D.; Investigation, M.R.M, D.P., K.K., A.C., M.E.C.; Sample

882    Procurement, M.E.C., J.A.S., M.L.F., and M.K.D.; Visualization, M.R.M, D.P. and K.K.; Writing –

883    Original Draft, M.R.M and B.A.D.; Writing – Review & Editing, M.R.M, D.P., K.K., A.C., M.E.C,

884    A.S., K.D.D., V.M.H., K.A.K. and B.A.D.; Funding Acquisition, V.M.H. and B.A.D.; Resources,

885    B.A.D.; Supervision, B.A.D., V.M.H., K.D.D., K.A., A.S. and K.A.K.

886

**DECLARATION OF INTERESTS**

888    D.P.E is an employee of Mammoth Biosciences and co-founder/employee of Ancilia

889    Therapeutics. A.S. is the founder/employee of Ancilia Therapeutics. B.A.D. is a co-founder and

890    shareholder of Ancilia Therapeutics.

891

**REFERENCES**

893    1.    Holers VM, Demoruelle MK, Kuhn KA, Buckner JH, Robinson WH, Okamoto Y, Norris

894    JM, Deane KD. Rheumatoid arthritis and the mucosal origins hypothesis: protection turns to

895    destruction. Nat Rev Rheumatol. 2018;14(9):542-57.

896    2.      MacGregor AJ, Snieder H, Rigby AS, Koskenvuo M, Kaprio J, Aho K, Silman AJ.

897    Characterizing the quantitative genetic contribution to rheumatoid arthritis using data from twins.

898    Arthritis Rheum. 2000;43(1):30-7.

899    3.      Deane KD, Demoruelle MK, Kelmenson LB, Kuhn KA, Norris JM, Holers VM. Genetic

900    and environmental risk factors for rheumatoid arthritis. Best Pract Res Clin Rheumatol.

901    2017;31(1):3-18.

902    4.      Chang HH, Liu GY, Dwivedi N, Sun B, Okamoto Y, Kinslow JD, Deane KD, Demoruelle

903    MK, Norris JM, Thompson PR, Sparks JA, Rao DA, Karlson EW, Hung HC, Holers VM, Ho IC. A

904    molecular signature of preclinical rheumatoid arthritis triggered by dysregulated PTPN22. JCI

905    Insight. 2016;1(17):e90045.

906    5.      Demoruelle MK. Mucosa biology and the development of rheumatoid arthritis: potential

907    for prevention by targeting mucosal processes. Clin Ther. 2019;41(7):1270-8.

908    6.      Demoruelle MK, Bowers E, Lahey LJ, Sokolove J, Purmalek M, Seto NL, Weisman MH,

909    Norris JM, Kaplan MJ, Holers VM, Robinson WH, Deane KD. Antibody responses to citrullinated

910    and noncitrullinated antigens in the sputum of subjects with rheumatoid arthritis and subjects at

911    risk for development of rheumatoid arthritis. Arthritis Rheumatol. 2018;70(4):516-27.

912    7.      Demoruelle MK, Harrall KK, Ho L, Purmalek MM, Seto NL, Rothfuss HM, Weisman MH,

913    Solomon JJ, Fischer A, Okamoto Y, Kelmenson LB, Parish MC, Feser M, Fleischer C, Anderson

914    C, Mahler M, Norris JM, Kaplan MJ, Cherrington BD, Holers VM, Deane KD. Anti-citrullinated

915    protein antibodies are associated with neutrophil extracellular traps in the sputum in relatives of

916    rheumatoid arthritis patients. Arthritis Rheumatol. 2017;69(6):1165-75.

917    8.      Hughes-Austin JM, Deane KD, Derber LA, Kolfenbach JR, Zerbe GO, Sokolove J, Lahey

918    LJ, Weisman MH, Buckner JH, Mikuls TR, O'Dell JR, Keating RM, Gregersen PK, Robinson

919    WH, Holers VM, Norris JM. Multiple cytokines and chemokines are associated with rheumatoid

920    arthritis-related autoimmunity in first-degree relatives without rheumatoid arthritis: Studies of the

921    Aetiology of Rheumatoid Arthritis (SERA). Ann Rheum Dis. 2013;72(6):901-7.

922    9.      Willis VC, Demoruelle MK, Derber LA, Chartier-Logan CJ, Parish MC, Pedraza IF,

923    Weisman MH, Norris JM, Holers VM, Deane KD. Sputum autoantibodies in patients with

924    established rheumatoid arthritis and subjects at risk of future clinically apparent disease.

925    Arthritis Rheum. 2013;65(10):2545-54.

926    10.     Brusca SB, Abramson SB, Scher JU. Microbiome and mucosal inflammation as extra-

927    articular triggers for rheumatoid arthritis and autoimmunity. Curr Opin Rheumatol.

928    2014;26(1):101-7.

929    11.     Scher JU, Sczesnak A, Longman RS, Segata N, Ubeda C, Bielski C, Rostron T,

930    Cerundolo V, Pamer EG, Abramson SB, Huttenhower C, Littman DR. Expansion of intestinal

931    *Prevotella copri* correlates with enhanced susceptibility to arthritis. Elife. 2013;2:e01202.

932    12.     Maeda Y, Kurakawa T, Umemoto E, Motooka D, Ito Y, Gotoh K, Hirota K, Matsushita M,

933    Furuta Y, Narazaki M, Sakaguchi N, Kayama H, Nakamura S, Iida T, Saeki Y, Kumanogoh A,

934    Sakaguchi S, Takeda K. Dysbiosis contributes to arthritis development via activation of

935    autoreactive T cells in the intestine. Arthritis Rheumatol. 2016;68(11):2646-61.

936    13.     Zhang X, Zhang D, Jia H, Feng Q, Wang D, Liang D, Wu X, Li J, Tang L, Li Y, Lan Z,

937    Chen B, Li Y, Zhong H, Xie H, Jie Z, Chen W, Tang S, Xu X, Wang X, Cai X, Liu S, Xia Y, Li J,

938    Qiao X, Al-Aama JY, Chen H, Wang L, Wu QJ, Zhang F, Zheng W, Li Y, Zhang M, Luo G, Xue

939    W, Xiao L, Li J, Chen W, Xu X, Yin Y, Yang H, Wang J, Kristiansen K, Liu L, Li T, Huang Q, Li

940    Y, Wang J. The oral and gut microbiomes are perturbed in rheumatoid arthritis and partly

941    normalized after treatment. Nat Med. 2015;21(8):895-905.

942    14.    Vaahtovuo J, Munukka E, Korkeamaki M, Luukkainen R, Toivanen P. Fecal microbiota

943    in early rheumatoid arthritis. J Rheumatol. 2008;35(8):1500-5.

944    15.    Toivanen P, Vartiainen S, Jalava J, Luukkainen R, Mottonen T, Eerola E, Manninen R.

945    Intestinal anaerobic bacteria in early rheumatoid arthritis (RA). Arthritis Research & Therapy.

946    2002;4.

947    16.    Alpizar-Rodriguez D, Lesker TR, Gronow A, Gilbert B, Raemy E, Lamacchia C, Gabay

948    C, Finckh A, Strowig T. *Prevotella copri* in individuals at risk for rheumatoid arthritis. Ann Rheum

949    Dis. 2019;78(5):590-3.

950    17.    Pianta A, Arvikar S, Strle K, Drouin EE, Wang Q, Costello CE, Steere AC. Evidence of

951    the immune relevance of *Prevotella copri*, a gut microbe, in patients with rheumatoid arthritis.

952    Arthritis Rheumatol. 2017;69(5):964-75.

953    18.    Drago L. *Prevotella copri* and microbiota in rheumatoid arthritis: fully convincing

954    evidence? J Clin Med. 2019;8(11).

955    19.    Duerkop BA. Bacteriophages shift the focus of the mammalian microbiota. PLoS Pathog.

956    2018;14(10):e1007310.

957    20.    Minot S, Sinha R, Chen J, Li H, Keilbaugh SA, Wu GD, Lewis JD, Bushman FD. The

958    human gut virome: inter-individual variation and dynamic response to diet. Genome Res.

959    2011;21(10):1616-25.

960    21.    Manrique P, Bolduc B, Walk ST, van der Oost J, de Vos WM, Young MJ. Healthy human

961    gut phageome. Proc Natl Acad Sci U S A. 2016;113(37):10400-5.

962    22.    Shkoporov AN, Hill C. Bacteriophages of the human gut: the "known unknown" of the

963    microbiome. Cell Host Microbe. 2019;25(2):195-209.

964    23.    Duerkop BA, Kleiner M, Paez-Espino D, Zhu W, Bushnell B, Hassell B, Winter SE,

965    Kyrpides NC, Hooper LV. Murine colitis reveals a disease-associated bacteriophage community.

966    Nat Microbiol. 2018;3(9):1023-31.

967    24.    Clooney AG, Sutton TDS, Shkoporov AN, Holohan RK, Daly KM, O'Regan O, Ryan FJ,

968    Draper LA, Plevy SE, Ross RP, Hill C. Whole-virome analysis sheds light on viral dark matter in

969    inflammatory bowel disease. Cell Host Microbe. 2019;26(6):764-78 e5.

970    25.    Khan Mirzaei M, Khan MAA, Ghosh P, Taranu ZE, Taguer M, Ru J, Chowdhury R, Kabir

971    MM, Deng L, Mondal D, Maurice CF. Bacteriophages isolated from stunted children can

972    regulate gut bacterial communities in an age-specific manner. Cell Host Microbe.

973    2020;27(2):199-212 e5.

974    26.    Zhao G, Vatanen T, Droit L, Park A, Kostic AD, Poon TW, Vlamakis H, Siljander H,

975    Harkonen T, Hamalainen AM, Peet A, Tillmann V, Ilonen J, Wang D, Knip M, Xavier RJ, Virgin

976    HW. Intestinal virome changes precede autoimmunity in type I diabetes-susceptible children.

977    Proc Natl Acad Sci U S A. 2017;114(30):E6166-E75.

978    27.    Barr JJ, Auro R, Furlan M, Whiteson KL, Erb ML, Pogliano J, Stotland A, Wolkowicz R,

979    Cutting AS, Doran KS, Salamon P, Youle M, Rohwer F. Bacteriophage adhering to mucus

980    provide a non-host-derived immunity. Proc Natl Acad Sci U S A. 2013;110(26):10771-6.

981    28.    Quistad SD, Grasis JA, Barr JJ, Rohwer FL. Viruses and the origin of microbiome

982    selection and immunity. ISME J. 2017;11(4):835-40.

983    29.    Gorski A, Dabrowska K, Miedzybrodzki R, Weber-Dabrowska B, Lusiak-Szelachowska

984    M, Jonczyk-Matysiak E, Borysowski J. Phages and immunomodulation. Future Microbiol.

985    2017;12:905-14.

986   30.    Minot SS, Willis AD. Clustering co-abundant genes identifies components of the gut

987        microbiome that are reproducibly associated with colorectal cancer and inflammatory bowel

988        disease. Microbiome. 2019;7(1):110.

989   31.    Yu AI, Zhao L, Eaton KA, Ho S, Chen J, Poe S, Becker J, Gonzalez A, McKinstry D,

990        Hasso M, Mendoza-Castrejon J, Whitfield J, Koumpouras C, Schloss PD, Martens EC, Chen

991        GY. Gut microbiota modulate CD8 T cell responses to influence colitis-associated

992        tumorigenesis. Cell Rep. 2020;31(1):107471.

993   32.    Norman JM, Handley SA, Baldridge MT, Droit L, Liu CY, Keller BC, Kambal A, Monaco

994        CL, Zhao G, Fleshner P, Stappenbeck TS, McGovern DP, Keshavarzian A, Mutlu EA, Sauk J,

995        Gevers D, Xavier RJ, Wang D, Parkes M, Virgin HW. Disease-specific alterations in the enteric

996        virome in inflammatory bowel disease. Cell. 2015;160(3):447-60.

997   33.    Gevers D, Kugathasan S, Denson LA, Vazquez-Baeza Y, Van Treuren W, Ren B,

998        Schwager E, Knights D, Song SJ, Yassour M, Morgan XC, Kostic AD, Luo C, Gonzalez A,

999        McDonald D, Haberman Y, Walters T, Baker S, Rosh J, Stephens M, Heyman M, Markowitz J,

1000       Baldassano R, Griffiths A, Sylvester F, Mack D, Kim S, Crandall W, Hyams J, Huttenhower C,

1001       Knight R, Xavier RJ. The treatment-naive microbiome in new-onset Crohn's disease. Cell Host

1002       Microbe. 2014;15(3):382-92.

1003  34.    Lee JY, Mannaa M, Kim Y, Kim J, Kim GT, Seo YS. Comparative analysis of fecal

1004       microbiota composition between rheumatoid arthritis and osteoarthritis patients. Genes (Basel).

1005       2019;10(10).

1006  35.    Van Belleghem JD, Dabrowska K, Vaneechoutte M, Barr JJ, Bollyky PL. Interactions

1007       between bacteriophage, bacteria, and the mammalian immune system. Viruses. 2018;11(1).

1008    36.    Rodriguez-Valera F, Martin-Cuadrado AB, Rodriguez-Brito B, Pasic L, Thingstad TF,

1009    Rohwer F, Mira A. Explaining microbial population genomics through phage predation. Nat Rev

1010    Microbiol. 2009;7(11):828-36.

1011    37.    Kieft K, Zhou Z, Anantharaman K. VIBRANT: automated recovery, annotation and

1012    curation of microbial viruses, and evaluation of viral community function from genomic

1013    sequences. Microbiome. 2020;8(1):90.

1014    38.    Kolfenbach JR, Deane KD, Derber LA, O'Donnell C, Weisman MH, Buckner JH, Gersuk

1015    VH, Wei S, Mikuls TR, O'Dell J, Gregersen PK, Keating RM, Norris JM, Holers VM. A

1016    prospective approach to investigating the natural history of preclinical rheumatoid arthritis (RA)

1017    using first-degree relatives of probands with RA. Arthritis Rheum. 2009;61(12):1735-42.

1018    39.    Demoruelle MK, Parish MC, Derber LA, Kolfenbach JR, Hughes-Austin JM, Weisman

1019    MH, Gilliland W, Edison JD, Buckner JH, Mikuls TR, O'Dell JR, Keating RM, Gregersen PK,

1020    Norris JM, Holers VM, Deane KD. Performance of anti-cyclic citrullinated peptide assays differs

1021    in subjects at increased risk of rheumatoid arthritis and subjects with established disease.

1022    Arthritis Rheum. 2013;65(9):2243-52.

1023    40.    Kang DW, Adams JB, Gregory AC, Borody T, Chittick L, Fasano A, Khoruts A, Geis E,

1024    Maldonado J, McDonough-Means S, Pollard EL, Roux S, Sadowsky MJ, Lipson KS, Sullivan

1025    MB, Caporaso JG, Krajmalnik-Brown R. Microbiota Transfer Therapy alters gut ecosystem and

1026    improves gastrointestinal and autism symptoms: an open-label study. Microbiome.

1027    2017;5(1):10.

1028    41.    Paez-Espino D, Pavlopoulos GA, Ivanova NN, Kyrpides NC. Nontargeted virus

1029    sequence discovery pipeline and virus clustering for metagenomic data. Nat Protoc.

1030    2017;12(8):1673-82.

1031  42.     Nayfach S, Camargo AP, Eloe-Fadrosh E, Roux S, Kyrpides N. CheckV: assessing the

1032  quality of metagenome-assmebled viral genomes. bioRxiv preprint. 2020.

1033  43.     Roux S, Adriaenssens EM, Dutilh BE, Koonin EV, Kropinski AM, Krupovic M, Kuhn JH,

1034  Lavigne R, Brister JR, Varsani A, Amid C, Aziz RK, Bordenstein SR, Bork P, Breitbart M,

1035  Cochrane GR, Daly RA, Desnues C, Duhaime MB, Emerson JB, Enault F, Fuhrman JA,

1036  Hingamp P, Hugenholtz P, Hurwitz BL, Ivanova NN, Labonte JM, Lee KB, Malmstrom RR,

1037  Martinez-Garcia M, Mizrachi IK, Ogata H, Paez-Espino D, Petit MA, Putonti C, Rattei T, Reyes

1038  A, Rodriguez-Valera F, Rosario K, Schriml L, Schulz F, Steward GF, Sullivan MB, Sunagawa S,

1039  Suttle CA, Temperton B, Tringe SG, Thurber RV, Webster NS, Whiteson KL, Wilhelm SW,

1040  Wommack KE, Woyke T, Wrighton KC, Yilmaz P, Yoshida T, Young MJ, Yutin N, Allen LZ,

1041  Kyrpides NC, Eloe-Fadrosh EA. Minimum Information about an Uncultivated Virus Genome

1042  (MIUViG). Nat Biotechnol. 2019;37(1):29-37.

1043  44.     Roux S, Paez-Espino D, Chen IA, Palaniappan K, Ratner A, Chu K, Reddy TBK,

1044  Nayfach S, Schulz F, Call L, Neches RY, Woyke T, Ivanova NN, Eloe-Fadrosh EA, Kyrpides

1045  NC. IMG/VR v3: an integrated ecological and evolutionary framework for interrogating genomes

1046  of uncultivated viruses. Nucleic Acids Res. 2020.

1047  45.     Barrangou R, Fremaux C, Deveau H, Richards M, Boyaval P, Moineau S, Romero DA,

1048  Horvath P. CRISPR provides acquired resistance against viruses in prokaryotes. Science.

1049  2007;315(5819):1709-12.

1050  46.     Stern A, Mick E, Tirosh I, Sagy O, Sorek R. CRISPR targeting reveals a reservoir of

1051  common phages associated with the human gut microbiome. Genome Res. 2012;22(10):1985-

1052  94.

1053    47.    Paez-Espino D, Eloe-Fadrosh EA, Pavlopoulos GA, Thomas AD, Huntemann M,

1054    Mikhailova N, Rubin E, Ivanova NN, Kyrpides NC. Uncovering Earth's virome. Nature.

1055    2016;536(7617):425-30.

1056    48.    Moreno-Gallego JL, Chou SP, Di Rienzi SC, Goodrich JK, Spector TD, Bell JT,

1057    Youngblut ND, Hewson I, Reyes A, Ley RE. Virome diversity correlates with intestinal

1058    microbiome diversity in adult monozygotic twins. Cell Host Microbe. 2019;25(2):261-72 e5.

1059    49.    Roux S, Emerson JB, Eloe-Fadrosh EA, Sullivan MB. Benchmarking viromics: an in

1060    silico evaluation of metagenome-enabled estimates of viral community composition and

1061    diversity. PeerJ. 2017;5:e3817.

1062    50.    Liang G, Zhao C, Zhang H, Mattei L, Sherrill-Mix S, Bittinger K, Kessler LR, Wu GD,

1063    Baldassano RN, DeRusso P, Ford E, Elovitz MA, Kelly MS, Patel MZ, Mazhani T, Gerber JS,

1064    Kelly A, Zemel BS, Bushman FD. The stepwise assembly of the neonatal virome is modulated

1065    by breastfeeding. Nature. 2020;581(7809):470-4.

1066    51.    Mukherjee S, Stamatis D, Bertsch J, Ovchinnikova G, Katta HY, Mojica A, Chen IA,

1067    Kyrpides NC, Reddy T. Genomes OnLine database (GOLD) v.7: updates and new features.

1068    Nucleic Acids Res. 2019;47(D1):D649-D59.

1069    52.    Nayfach S, Roux S, Seshadri R, Udwary D, Varghese N, Schulz F, Wu D, Paez-Espino

1070    D, Chen IM, Huntemann M, Palaniappan K, Ladau J, Mukherjee S, Reddy TBK, Nielsen T,

1071    Kirton E, Faria JP, Edirisinghe JN, Henry CS, Jungbluth SP, Chivian D, Dehal P, Wood-

1072    Charlson EM, Arkin AP, Tringe SG, Visel A, Consortium IMD, Woyke T, Mouncey NJ, Ivanova

1073    NN, Kyrpides NC, Eloe-Fadrosh EA. A genomic catalog of Earth's microbiomes. Nat Biotechnol.

1074    2020.

1075    53.    Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for

1076    RNA-seq data with DESeq2. Genome Biol. 2014;15(12):550.

1077    54.    Breitbart M, Thompson LR, Suttle CA, Sullivan MB. Exploring the vast diversity of marine

1078    viruses. Oceanography. 2007;20(2):135-9.

1079    55.    Thompson LR, Zeng Q, Kelly L, Huang KH, Singer AU, Stubbe J, Chisholm SW. Phage

1080    auxiliary metabolic genes and the redirection of cyanobacterial host carbon metabolism. Proc

1081    Natl Acad Sci U S A. 2011;108(39):E757-64.

1082    56.    Breitbart M, Bonnain C, Malki K, Sawaya NA. Phage puppet masters of the marine

1083    microbial realm. Nat Microbiol. 2018;3(7):754-66.

1084    57.    Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, Gordon JI. Viruses in

1085    the faecal microbiota of monozygotic twins and their mothers. Nature. 2010;466(7304):334-8.

1086    58.    Chevallereau A, Blasdel BG, De Smet J, Monot M, Zimmermann M, Kogadeeva M,

1087    Sauer U, Jorth P, Whiteley M, Debarbieux L, Lavigne R. Next-generation "-omics" approaches

1088    reveal a massive alteration of host RNA metabolism during bacteriophage infection of

1089    *Pseudomonas aeruginosa*. PLoS Genet. 2016;12(7):e1006134.

1090    59.    De Smet J, Zimmermann M, Kogadeeva M, Ceyssens PJ, Vermaelen W, Blasdel B, Bin

1091    Jang H, Sauer U, Lavigne R. High coverage metabolomics analysis reveals phage-specific

1092    alterations to *Pseudomonas aeruginosa* physiology during infection. ISME J. 2016;10(8):1823-

1093    35.

1094    60.    Chatterjee A, Willett JLE, Nguyen UT, Monogue B, Palmer KL, Dunny GM, Duerkop BA.

1095    Parallel Genomics Uncover Novel Enterococcal-Bacteriophage Interactions. mBio. 2020;11(2).

1096  61.  Howard-Varona C, Lindback MM, Bastien GE, Solonenko N, Zayed AA, Jang H,

1097  Andreopoulos B, Brewer HM, Glavina Del Rio T, Adkins JN, Paul S, Sullivan MB, Duhaime MB.

1098  Phage-specific metabolic reprogramming of virocells. ISME J. 2020;14(4):881-95.

1099  62.  Valvano MA, Messner P, Kosma P. Novel pathways for biosynthesis of nucleotide-

1100  activated glycero-manno-heptose precursors of bacterial glycoproteins and cell surface

1101  polysaccharides. Microbiology (Reading). 2002;148(Pt 7):1979-89.

1102  63.  Nakao R, Senpuku H, Watanabe H. *Porphyromonas gingivalis galE* is involved in

1103  lipopolysaccharide O-antigen synthesis and biofilm formation. Infect Immun. 2006;74(11):6145-

1104  53.

1105  64.  Crummett LT, Puxty RJ, Weihe C, Marston MF, Martiny JBH. The genomic content and

1106  context of auxiliary metabolic genes in marine cyanomyoviruses. Virology. 2016;499:219-29.

1107  65.  Lopez-Perez M, Haro-Moreno JM, de la Torre JR, Rodriguez-Valera F. Novel

1108  caudovirales associated with marine group I Thaumarchaeota assembled from metagenomes.

1109  Environ Microbiol. 2019;21(6):1980-8.

1110  66.  Ignacio-Espinoza JC, Sullivan MB. Phylogenomics of T4 cyanophages: lateral gene

1111  transfer in the 'core' and origins of host genes. Environ Microbiol. 2012;14(8):2113-26.

1112  67.  Breitbart M. Marine viruses: truth or dare. Ann Rev Mar Sci. 2012;4:425-48.

1113  68.  Markenson JA. Worldwide trends in the socioeconomic impact and long-term prognosis

1114  of rheumatoid arthritis. Semin Arthritis Rheum. 1991;21(2 Suppl 1):4-12.

1115  69.  Hunter TM, Boytsov NN, Zhang X, Schroeder K, Michaud K, Araujo AB. Prevalence of

1116  rheumatoid arthritis in the United States adult population in healthcare claims databases, 2004-

1117  2014. Rheumatol Int. 2017;37(9):1551-7.

1118    70.    Scher JU, Abramson SB. The microbiome and rheumatoid arthritis. Nat Rev Rheumatol.

1119    2011;7(10):569-78.

1120    71.    Mirzaei MK, Maurice CF. Menage a trois in the human gut: interactions between host,

1121    bacteria and phages. Nat Rev Microbiol. 2017;15(7):397-408.

1122    72.    Bin Jang H, Bolduc B, Zablocki O, Kuhn JH, Roux S, Adriaenssens EM, Brister JR,

1123    Kropinski AM, Krupovic M, Lavigne R, Turner D, Sullivan MB. Taxonomic assignment of

1124    uncultivated prokaryotic virus genomes is enabled by gene-sharing networks. Nat Biotechnol.

1125    2019;37(6):632-9.

1126    73.    Sutton TDS, Hill C. Gut bacteriophage: current understanding and challenges. Front

1127    Endocrinol (Lausanne). 2019;10:784.

1128    74.    Roux S, Hallam SJ, Woyke T, Sullivan MB. Viral dark matter and virus-host interactions

1129    resolved from publicly available microbial genomes. Elife. 2015;4.

1130    75.    Liu X, Zeng B, Zhang J, Li W, Mou F, Wang H, Zou Q, Zhong B, Wu L, Wei H, Fang Y.

1131    Role of the gut microbiome in modulating arthritis progression in mice. Sci Rep. 2016;6:30594.

1132    76.    Jubair WK, Hendrickson JD, Severs EL, Schulz HM, Adhikari S, Ir D, Pagan JD, Anthony

1133    RM, Robertson CE, Frank DN, Banda NK, Kuhn KA. Modulation of inflammatory arthritis in mice

1134    by gut microbiota through mucosal inflammation andautoantibody generation. Arthritis

1135    Rheumatol. 2018;70(8):1220-33.

1136    77.    Braschi E, Shojania K, Allan GM. Anti-CCP: a truly helpful rheumatoid arthritis test? Can

1137    Fam Physician. 2016;62(3):234.

1138    78.    Krych L, Nielsen DS, Hansen AK, Hansen CH. Gut microbial markers are associated

1139    with diabetes onset, regulatory imbalance, and IFN-gamma level in NOD mice. Gut Microbes.

1140    2015;6(2):101-9.

1141    79.    Chriswell ME, Kuhn KA. Microbiota-mediated mucosal inflammation in arthritis. Best

1142    Pract Res Clin Rheumatol. 2019;33(6):101492.

1143    80.    Studer N, Desharnais L, Beutler M, Brugiroux S, Terrazos MA, Menin L, Schurch CM,

1144    McCoy KD, Kuehne SA, Minton NP, Stecher B, Bernier-Latmani R, Hapfelmeier S. Functional

1145    intestinal bile acid 7alpha-dehydroxylation by *Clostridium scindens* associated with protection

1146    from *Clostridium difficile* infection in a gnotobiotic mouse model. Front Cell Infect Microbiol.

1147    2016;6:191.

1148    81.    Buffie CG, Bucci V, Stein RR, McKenney PT, Ling L, Gobourne A, No D, Liu H,

1149    Kinnebrew M, Viale A, Littmann E, van den Brink MR, Jenq RR, Taur Y, Sander C, Cross JR,

1150    Toussaint NC, Xavier JB, Pamer EG. Precision microbiome reconstitution restores bile acid

1151    mediated resistance to *Clostridium difficile*. Nature. 2015;517(7533):205-8.

1152    82.    Duboc H, Rajca S, Rainteau D, Benarous D, Maubert MA, Quervain E, Thomas G,

1153    Barbu V, Humbert L, Despras G, Bridonneau C, Dumetz F, Grill JP, Masliah J, Beaugerie L,

1154    Cosnes J, Chazouilleres O, Poupon R, Wolf C, Mallet JM, Langella P, Trugnan G, Sokol H,

1155    Seksik P. Connecting dysbiosis, bile-acid dysmetabolism and gut inflammation in inflammatory

1156    bowel diseases. Gut. 2013;62(4):531-9.

1157    83.    Atarashi K, Tanoue T, Oshima K, Suda W, Nagano Y, Nishikawa H, Fukuda S, Saito T,

1158    Narushima S, Hase K, Kim S, Fritz JV, Wilmes P, Ueha S, Matsushima K, Ohno H, Olle B,

1159    Sakaguchi S, Taniguchi T, Morita H, Hattori M, Honda K. Treg induction by a rationally selected

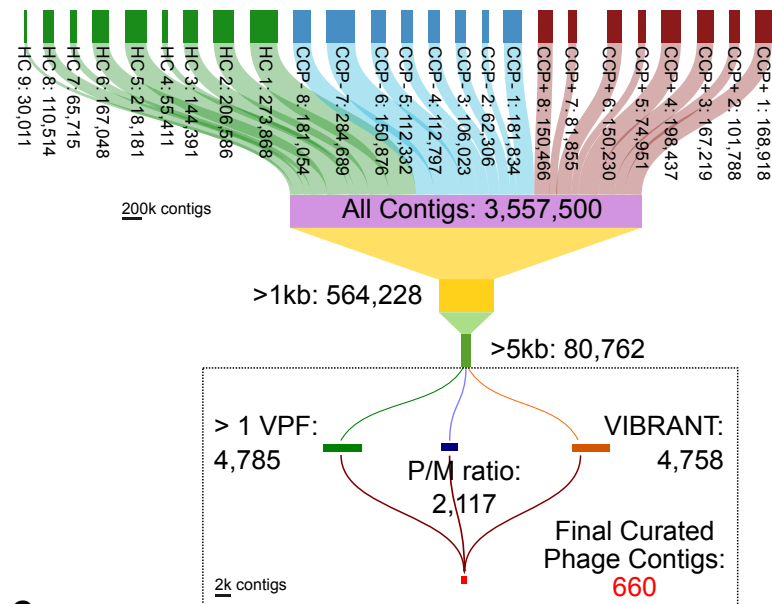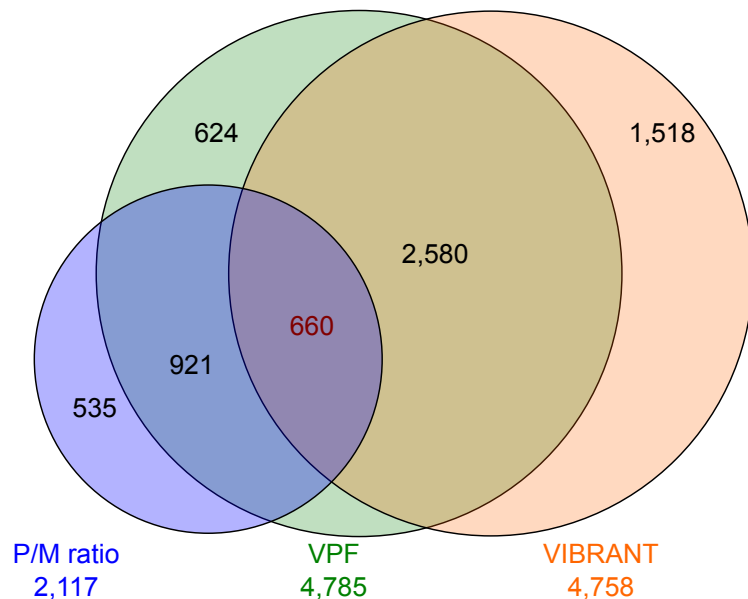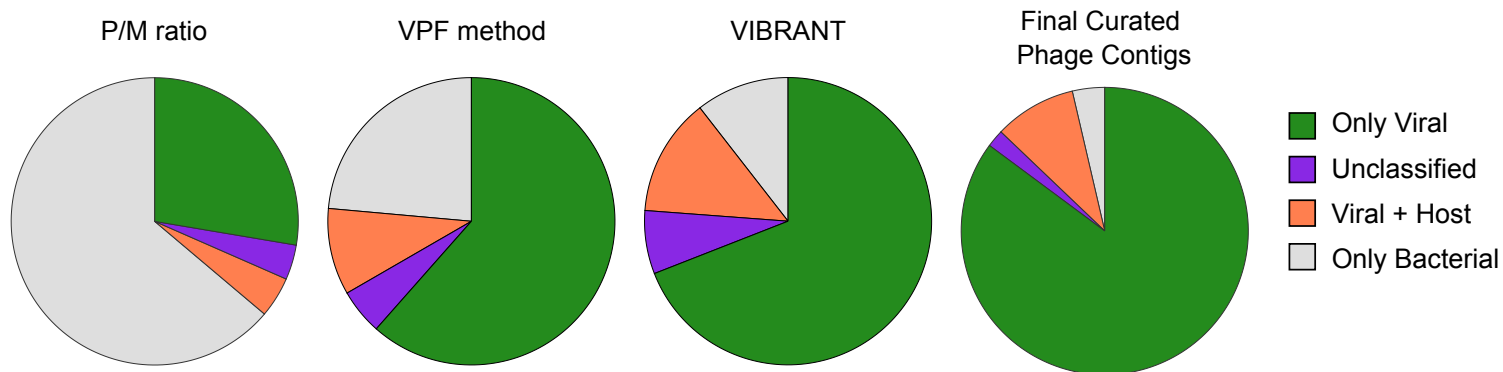1160    mixture of Clostridia strains from the human microbiota. Nature. 2013;500(7461):232-6.

1161    84.    Gerard P. Metabolism of cholesterol and bile acids by the gut microbiota. Pathogens.

1162    2013;3(1):14-24.

1163    85.    Yao L, Seaton SC, Ndousse-Fetter S, Adhikari AA, DiBenedetto N, Mina AI, Banks AS,

1164    Bry L, Devlin AS. A selective gut bacterial bile salt hydrolase alters host metabolism. Elife.

1165    2018;7.

1166    86.    Campbell DE, Ly LK, Ridlon JM, Hsiao A, Whitaker RJ, Degnan PH. Infection with

1167    *Bacteroides* phage BV01 alters the host transcriptome and bile acid metabolism in a common

1168    human gut microbe. Cell Rep. 2020;32(11):108142.

1169    87.    Lozupone CA, Stombaugh JI, Gordon JI, Jansson JK, Knight R. Diversity, stability and

1170    resilience of the human gut microbiota. Nature. 2012;489(7415):220-30.

1171    88.    Zhu W, Winter MG, Byndloss MX, Spiga L, Duerkop BA, Hughes ER, Buttner L, de Lima

1172    Romao E, Behrendt CL, Lopez CA, Sifuentes-Dominguez L, Huff-Hardy K, Wilson RP, Gillis CC,

1173    Tukel C, Koh AY, Burstein E, Hooper LV, Baumler AJ, Winter SE. Precision editing of the gut

1174    microbiota ameliorates colitis. Nature. 2018;553(7687):208-11.

1175    89.    He SM, Wathier M, Podzelinska K, Wong M, McSorley FR, Asfaw A, Hove-Jensen B, Jia

1176    Z, Zechel DL. Structure and mechanism of PhnP, a phosphodiesterase of the carbon-

1177    phosphorus lyase pathway. Biochemistry. 2011;50(40):8603-15.

1178    90.    Martinez A, Ventouras LA, Wilson ST, Karl DM, Delong EF. Metatranscriptomic and

1179    functional metagenomic analysis of methylphosphonate utilization by marine bacteria. Front

1180    Microbiol. 2013;4:340.

1181    91.    Lee KS, Metcalf WW, Wanner BL. Evidence for two phosphonate degradative pathways

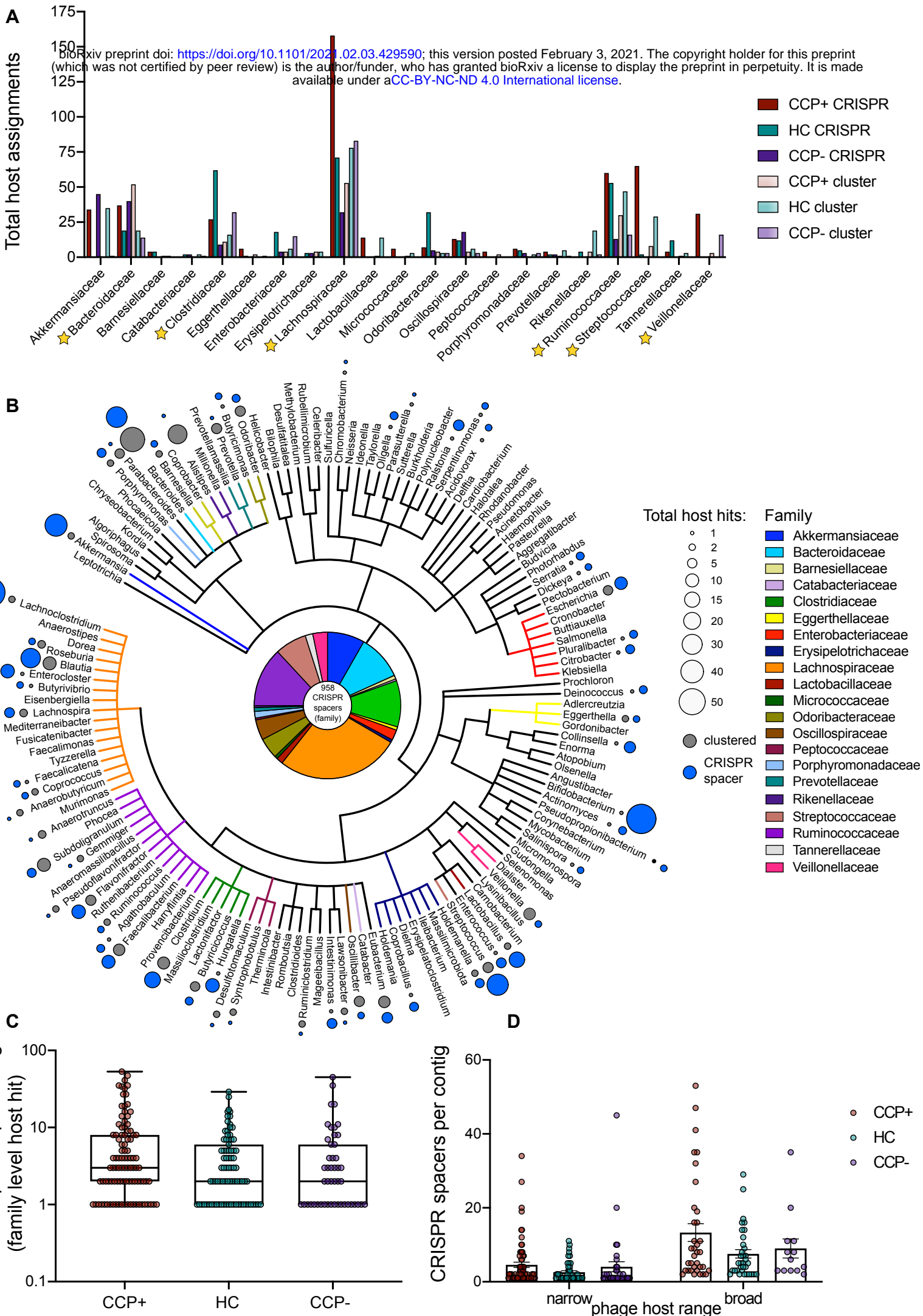1182    in *Enterobacter aerogenes*. J Bacteriol. 1992;174(8):2501-10.

1183    92.    Artacho A, Isaac S, Nayak R, Flor-Duro A, Alexander M, Koo I, Manasson J, Smith PB,

1184    Rosenthal P, Homsi Y, Gulko P, Pons J, Puchades-Carrasco L, Izmirly P, Patterson A,

1185    Abramson SB, Pineda-Lucena A, Turnbaugh PJ, Ubeda C, Scher JU. The pre-treatment gut

1186    microbiome is associated with lack of response to methotrexate in new onset rheumatoid

1187    arthritis. Arthritis Rheumatol. 2020.

1188    93.    Hsu BB, Gibson TE, Yeliseyev V, Liu Q, Lyon L, Bry L, Silver PA, Gerber GK. Dynamic

1189    modulation of the gut microbiota and metabolome by bacteriophages in a mouse model. Cell

1190    Host Microbe. 2019;25(6):803-14 e5.

1191    94.    Kuller LH, Mackey RH, Walitt BT, Deane KD, Holers VM, Robinson WH, Sokolove J,

1192    Chang Y, Liu S, Parks CG, Wright NC, Moreland LW. Determinants of mortality among

1193    postmenopausal women in the women's health initiative who report rheumatoid arthritis. Arthritis

1194    Rheumatol. 2014;66(3):497-507.

1195    95.    Choi S, Lee KH. Clinical management of seronegative and seropositive rheumatoid

1196    arthritis: A comparative study. PLoS One. 2018;13(4):e0195550.

1197    96.    Chatterjee A, Duerkop BA. Beyond bacteria: bacteriophage-eukaryotic host interactions

1198    reveal emerging paradigms of health and disease. Front Microbiol. 2018;9:1394.

1199    97.    Ayyappan P, Harms RZ, Seifert JA, Bemis EA, Feser ML, Deane KD, Demoruelle MK,

1200    Mikuls TR, Holers VM, Sarvetnick NE. Heightened levels of antimicrobial response factors in

1201    patients with rheumatoid arthritis. Front Immunol. 2020;11:427.

1202    98.    Kleiner M, Hooper LV, Duerkop BA. Evaluation of methods to purify virus-like particles

1203    for metagenomic sequencing of intestinal viromes. BMC Genomics. 2015;16:7.

1204    99.    Duerkop BA, Huo W, Bhardwaj P, Palmer KL, Hooper LV. Molecular basis for lytic

1205    bacteriophage resistance in Enterococci. mBio. 2016;7(4).

1206    100.    Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ,

1207    Fierer N, Knight R. Global patterns of 16S rRNA diversity at a depth of millions of sequences

1208    per sample. Proc Natl Acad Sci U S A. 2011;108 Suppl 1:4516-22.

1209    101.    Dowd SE, Sun Y, Wolcott RD, Domingo A, Carroll JA. Bacterial tag-encoded FLX

1210    amplicon pyrosequencing (bTEFAP) for microbiome studies: bacterial diversity in the ileum of

1211    newly weaned *Salmonella*-infected pigs. Foodborne Pathog Dis. 2008;5(4):459-72.

1212    102.    Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,

1213    Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, Weber CF.

1214    Introducing mothur: open-source, platform-independent, community-supported software for

1215    describing and comparing microbial communities. Appl Environ Microbiol. 2009;75(23):7537-41.

1216    103.    Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD. Development of a dual-

1217    index sequencing strategy and curation pipeline for analyzing amplicon sequence data on the

1218    MiSeq Illumina sequencing platform. Appl Environ Microbiol. 2013;79(17):5112-20.

1219    104.    DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber T, Dalevi D,

1220    Hu P, Andersen GL. Greengenes, a chimera-checked 16S rRNA gene database and workbench

1221    compatible with ARB. Appl Environ Microbiol. 2006;72(7):5069-72.

1222    105.    Bushnell B. BBMap short read aligner, and other bioinformatic tools. 38.56 ed.

1223    https://sourceforge.net/projects/bbmap/2019.

1224    106.    Peng Y, Leung HC, Yiu SM, Chin FY. IDBA-UD: a de novo assembler for single-cell and

1225    metagenomic sequencing data with highly uneven depth. Bioinformatics. 2012;28(11):1420-8.

1226    107.    Li D, Luo R, Liu CM, Leung CM, Ting HF, Sadakane K, Yamashita H, Lam TW.

1227    MEGAHIT v1.0: A fast and scalable metagenome assembler driven by advanced methodologies

1228    and community practices. Methods. 2016;102:3-11.

1229

# Figure 1



**A**

HC 9: 30,011
HC 8: 110,514
HC 7: 65,715
HC 6: 167,048
HC 5: 218,181
HC 4: 55,411
HC 3: 144,391
HC 2: 206,586
HC 1: 273,868
CCP-8: 181,054
CCP-7: 284,689
CCP-6: 150,876
CCP-5: 112,332
CCP-4: 112,797
CCP-3: 106,023
CCP-2: 62,306
CCP-1: 181,834
CCP+8: 150,466
CCP+7: 81,855
CCP+6: 150,230
CCP+5: 74,951
CCP+4: 198,437
CCP+3: 167,219
CCP+2: 101,788
CCP+1: 168,918

200k contigs

All Contigs: 3,557,500

>1kb: 564,228

>5kb: 80,762

> 1 VPF: 4,785

P/M ratio: 2,117

VIBRANT: 4,758

Final Curated Phage Contigs: 660

2k contigs

**B**

624
1,518
2,580
660
921
535

P/M ratio 2,117
VPF 4,785
VIBRANT 4,758

**C**

P/M ratio     VPF method     VIBRANT     Final Curated Phage Contigs

- Only Viral
- Unclassified
- Viral + Host
- Only Bacterial

# Figure 2

# Figure 3

**Figure 4**

# Figure 5

# Figure 6