

1 **Biology and taxonomy of crAss-like bacteriophages, the most abundant virus in the**  
2 **human gut**

3

4 Emma Guerin<sup>a,b,1</sup>, Andrey Shkoporov<sup>a,1</sup>, Stephen R. Stockdale<sup>a,c,1</sup>, Adam G. Clooney<sup>a</sup>, Feargal  
5 J. Ryan<sup>a</sup>, Lorraine A. Draper<sup>a</sup>, Enrique Gonzalez-Tortuero<sup>a</sup>, R. Paul Ross<sup>a,b,c</sup>, Colin Hill<sup>a,b,2</sup>

6

7 <sup>a</sup>APC Microbiome Ireland, University College Cork, Co. Cork, Ireland

8 <sup>b</sup>School of Microbiology, University College Cork, Co. Cork, Ireland

9 <sup>c</sup>Teagasc Food Research Centre, Moorepark, Fermoy, Co. Cork, Ireland

10

11 <sup>1</sup>These authors contributed equally to this work.

12 <sup>2</sup>Corresponding author.

13

14 Keywords: Bacteriophages, crAssphage, gut microbiota, human microbiome, phage  
15 taxonomy, phage characterisation

16

17     **Abstract**

18             CrAssphage is yet to be cultured even though it represents the most abundant virus in  
19     the gut microbiota of humans. Recently, sequence based classification was performed on  
20     distantly related crAss-like phages from multiple environments, leading to the proposal of a  
21     familial level taxonomic group [Yutin N, et al. (2018) Discovery of an expansive  
22     bacteriophage family that includes the most abundant viruses from the human gut. Nat  
23     Microbiol 3(1):38–46]. Here, we assembled the metagenomic sequencing reads from 702  
24     human faecal virome/phageome samples and obtained 98 complete circular crAss-like phage  
25     genomes and 145 contigs  $\geq 70$ kb. *In silico* comparative genomics and taxonomic analysis was  
26     performed, resulting in a classification scheme of crAss-like phages from human faecal  
27     microbiomes into 4 candidate subfamilies composed of 10 candidate genera. Moreover,  
28     laboratory analysis was performed on faecal samples from an individual harbouring 7 distinct  
29     crAss-like phages. We achieved propagation of crAss-like phages in *ex vivo* human faecal  
30     fermentations and visualised *Podoviridae* virions by electron microscopy. Furthermore,  
31     detection of a crAss-like phage capsid protein could be linked to metagenomic sequencing  
32     data confirming crAss-like phage structural annotations.

33

34

35

36

37

38

39

40

41

42 **Significance**

43 CrAssphage is the most abundant biological entity in the human gut, but it remains  
44 uncultured in the laboratory and its host(s) is unknown. CrAssphage was not identified in  
45 metagenomic studies for many years as its sequence is so different from anything present in  
46 databases. To this day, it can only be detected from sequences assembled from metagenomics  
47 or viromic datasets (**crAss** – **cross Assembly**). In this study, we identified 243 new crAss-like  
48 phages from human faecal metagenomic studies. Taxonomic analysis of these crAss-like  
49 phages highlighted their extensive diversity within the human microbiome. We also present  
50 the first propagation of crAssphage in faecal fermentations and provide the first electron  
51 micrographs of this extraordinary bacteriophage.

52

## 53 **Introduction**

54 In recent years, increasing numbers of bacteria, archaea, fungi, protists and viruses  
55 residing on and within the human body have been associated with various states of human  
56 health and disease, including diet, age, weight, inflammatory bowel disease (IBD), diabetes,  
57 and cognition (1–7). A relatively small number of eukaryote viruses present in the  
58 gastrointestinal tract can target the human host, however, much larger and much more  
59 complex populations of viruses that target bacteria (bacteriophages) also reside there. The  
60 role of phages in the gut has been a subject of increased interest as initial investigations have  
61 revealed substantial differences in bacteriophage populations between healthy and diseased  
62 cohorts (7–11). It is likely that phages have an important role in shaping our gut microbiome,  
63 but their precise role remains poorly understood.

64 In 2014, metagenomic studies of the viral fraction of the human gut microbiota  
65 identified a DNA phage, crAssphage, detectable in approximately 50% of individuals from  
66 specific human populations and reaching up to 90% of the total viral DNA load in faeces of  
67 certain individuals (12). Dutilh and colleagues noted that crAssphage had been overlooked in  
68 previous metagenomic studies as the vast majority of its genes do not match known  
69 sequences present in databases. It has been predicted, based on indirect evidence using host  
70 co-occurrence profiling, that prototypical crAssphage infects *Bacteroides*, an abundant genus  
71 of bacteria important for the normal gut function of humans. However, since crAssphage has  
72 never been isolated in culture, its host range, replication strategy, virion morphology and  
73 impact on the human gastrointestinal microbiota remains unknown. Thus, a better  
74 understanding of crAssphage is crucial to understanding phage host dynamics in the human  
75 gut microbiota.

76 Originally crAssphage was published as an individual phage following cross-  
77 assembly of several metagenomic samples (12). Analysis by Manrique *et al.*, of the healthy

78 human gut phageome identified 4 circular crAssphage genomes and several related  
79 incomplete contigs (10). PCR amplification and sequencing of the crAssphage polymerase  
80 gene by Liang and colleagues similarly demonstrated diversity amongst crAssphage-  
81 positive faecal samples (13). Recently, Cinek *et al.* described updated PCR primer  
82 sequences for the detection and evaluation of crAssphage diversity, while Stachler *et al.*  
83 developed their own primers targeting conserved genomic regions to evaluate the  
84 abundance of crAssphage as an indicator of human faecal pollution (14, 15). Finally, an  
85 epidemiological survey of crAssphages conducted by Dutilh, Edwards and colleagues has  
86 suggested crAssphage is associated with humans and primates globally with significant  
87 diversity (manuscript currently in preparation).

88 A recent study provided the first detailed sequence-based taxonomic categorisation  
89 of crAss-like phages, proposing a novel familial level taxonomic group that would include  
90 crAssphage itself, as well as various related bacteriophages, from multiple environments  
91 (16). However, the authors noted that this classification is in contrast with the classical  
92 viral taxonomy scheme currently in use. Such taxonomy strictly categorises crAssphage as  
93 a member of the *Podoviridae* family. Previous attempts to reconcile sequence-based and  
94 classical viral taxonomy have proposed *Podoviridae* sharing >40% orthologous protein-  
95 coding genes be grouped at the taxonomic rank of genus, while phages sharing only 20-  
96 40% orthologous protein-coding genes should be grouped at the higher taxonomic rank of  
97 subfamily (17). Other reports describe a phage genus as a cohesive group of viruses  
98 sharing >50% nucleotide sequence similarity (18). As crAssphage is not a single entity, but  
99 rather a group of crAss-like phages that share similarity with the prototypical crAssphage  
100 at various levels, a comparative analysis of crAss-like phage sequences is required to  
101 enable detailed taxonomic characterisation.

102           In this study, we combine several *in silico* and *in vitro* approaches to further explore  
103 the diversity of crAss-like phages in the human gut, and better understand their biological  
104 properties. We performed an in-depth analysis of crAss-like sequences from a number of  
105 previously published and unpublished human faecal virome datasets (1, 7, 9, 10, 19).  
106 Subsequent to the assembly of metagenomic sequencing reads, crAss-like phage contigs were  
107 identified using conserved genetic signatures. In total, 98 complete circular and 145 near-  
108 complete ( $\geq 70$ kb) linear contigs of crAss-like phages were identified for genomic and  
109 taxonomic analyses. Laboratory analysis of crAss-like phages was focused on a human donor  
110 identified as a stable carrier of several highly predominant crAssphage-like DNA sequences,  
111 including one closely related to the prototypical crAssphage. *Ex vivo* faecal fermentations  
112 enabled the amplification of a virus highly related to the prototypical crAssphage, with  
113 electron micrographs supporting the proposal that crAss-like phages are members of the  
114 *Podoviridae* family. These results represent the first example of biological characterisation of  
115 this highly prevalent and, potentially, very important human microbiome virus.

## 116 **Results**

117 **Detection of crAss-like phage contigs.** Following the assembly of 702 human faecal  
118 virome/phageome metagenomic samples listed in Supplementary Table 1, contigs were  
119 screened for relatedness to the prototypical crAssphage virus, henceforth referred to as  
120 crAssphage *sensu stricto*. Initially, the polymerase of crAssphage *sensu stricto* (UGP\_018,  
121 NC\_024711.1) was used for crAss-like phage detection due to its use in several studies as a  
122 genetic signature to determine diversity of crAss-like phages (13, 20, 21). However, we  
123 extended our criteria in order to include partial genomes ( $\geq 70\text{kb}$ ) that may not have included  
124 the polymerase gene in the assembly. Therefore, after an initial detection of crAss-like phages  
125 using the polymerase sequence, we identified the most conserved crAss-like phage protein in  
126 our dataset as the terminase protein, encoded by crAssphage *sensu stricto* UGP\_092. The  
127 terminase was subsequently used as a second genetic signature for identifying crAss-like  
128 phage contigs.

129 Initially, 239 contigs  $\geq 70\text{kb}$  were detected with similarity to crAssphage *sensu stricto*  
130 polymerase sequence. An additional 59 contigs  $\geq 70\text{kb}$  were subsequently detected with  
131 relatedness to crAssphage *sensu stricto* terminase sequence. Following an initial examination  
132 of the contig sequences retrieved, more stringent parameters were implemented. Only contigs  
133 whose polymerase and/or terminase sequence(s) aligned with greater than 350bp were  
134 considered for further analysis as crAss-like phages. This reduced the total number of crAss-  
135 like phages to 256. In addition, as several assembled metagenomic samples were from the  
136 same person sequenced at multiple time points, redundant contigs were removed from further  
137 analysis. When two or more contigs aligned with 100 percent identity, the longer contig or the  
138 contig with the highest coverage was retained. This resulted in a total of 244 crAss-like  
139 contigs (including crAssphage *sensu stricto*), with 143 contigs containing both a polymerase  
140 and terminase, 60 a polymerase only and 40 a terminase only. Of the 244 crAss-like phage

141 contigs, metadata was available for the majority of their originating faecal samples. CrAss-  
142 like phages were detected in healthy individuals across a wide age range (including infants 1  
143 year of age and individuals  $\geq 65$  years of age) and individuals suffering from Crohn's disease,  
144 ulcerative colitis, cystic fibrosis, kwashiorkor and marasmus.

145 **Taxonomy of crAss-like phages.** In order to compare the phylogeny of the more  
146 distantly related phages proposed to be included into a crAss-like familial level taxon by  
147 Yutin *et al.* (16) with those identified in this study, a phylogenetic tree of conserved crAss-  
148 like phage terminase sequences was constructed (Supplementary Figure 1). Amino acid  
149 terminase sequences were used to generate mid-point rooted phylogenetic trees.  
150 Predominantly, the terminase sequences of very distant crAss-like phage relatives identified  
151 by Yutin *et al.* from various environmental sources were distinct from the various candidate  
152 genera of crAss-like phages observed in the phylogram. However, the human gut microbiome  
153 phage, IAS virus (16), characterised by Yutin *et al.* as crAss-like, clustered closely with  
154 candidate genus VI crAss-like phages identified in this study.

155 Previously, studies have used the percentage of shared homologous proteins as a means  
156 of defining phage taxonomic ranks (17). Therefore, clusters of phages sharing between 20-  
157 40% of their protein-coding genes were categorised as related at the subfamily level, while  
158 phages sharing >40% protein-coding genes were grouped at the genus level. A heatmap based  
159 on the percentages of shared orthologous proteins suggests that crAss-like phages form 4  
160 candidate subfamilies. The four subfamilies were assigned the nomenclature  
161 *alphacrAssvirinae* (which contains crAssphage *sensu stricto*), *betacrAssvirinae* (which  
162 contains IAS virus), *gammacrAssvirinae* and *deltacrAssvirinae* (Figure 1). These subfamilies  
163 can be further subdivided into 10 candidate genera, with Candidate Genus I containing  
164 crAssphage *sensu stricto* and Candidate Genus VI containing the IAS virus. Metadata of all



165 crAss-like phages analysed in this study, including their categorisation into the various  
166 taxonomic divisions, is available in Supplementary Table 2.

167 An alternative approach for characterising the encoded proteome of crAss-like phages  
168 was performed by visualisation of genome clusters using the t-SNE machine learning  
169 algorithm with Euclidean distances of orthologous genes distribution between genomes as an  
170 input. Applying the previously determined 10 crAss-like phage candidate genera  
171 classifications to the t-SNE two-dimensional ordination demonstrated that some clusters  
172 showed uniformity while others groups were quite dispersed, such as Candidate Genus II and  
173 VII, respectively (Figure 2A). In addition, no single cluster of crAss-like phages is  
174 exclusively associated with healthy or diseased individuals.

175 Groups of crAss-like phages with a similar G+C nucleotide content would be expected  
176 to infect related bacteria, since phage G+C content often aligns to that of its host (22, 23).  
177 Therefore, several groups of crAss-like phages, such as candidate genera II, IV, V, VII and X,  
178 are likely infect closely related bacterial taxa within the human microbiome (Figure 2B).  
179 Candidate genus I is the most homogenous group of crAss-like phages containing crAssphage  
180 *sensu stricto* and 30 additional complete circular genomes and 29 linear contigs  $\geq 70$ kb with a  
181 distinct G+C nucleotide content ( $29.11 \pm 0.14\%$ ). Candidate genera III and VI display the  
182 greatest heterogeneity, with G+C contents of  $28.94 \pm 3.03\%$  and  $35.81 \pm 2.56\%$ , respectively.

183 **Nucleotide comparison of crAss-like phages.** To further investigate the relatedness of  
184 crAss-like phages, a more detailed comparison at the nucleotide level was performed by  
185 calculating their average nucleotide identity (Figure 3). Candidate genera III and VI of crAss-  
186 like phages, as defined by the percentage of their shared encoded proteins, also do not cluster  
187 into clearly definable groups based on nucleotide composition. Candidate Genus I, containing  
188 crAssphage *sensu stricto*, forms a well-defined homogenous taxonomic group even when  
189 analysed at the higher resolution of nucleotide composition. This is to be expected as

190 crAssphage *sensu stricto* was the starting point for finding all crAss-like phages examined in  
191 this study and thus has the most sequences available for analysis.

192 Interestingly, the majority of crAss-like candidate genera demonstrate the same type of  
193 genomic organization (Supplementary Figure 2). Prominent features were shared between  
194 candidate genera I – V, IX, and X. These include; circular genomes with size ranging from 92  
195 to 104kb, two clearly separated genome regions with opposite gene orientation and inverted  
196 G+C skew (the smaller region encodes proteins involved in replication, the bigger region  
197 coding for proteins involved in transcription and virion assembly, as suggested by Yutin *et*  
198 *al.*), the presence of giant open reading frames with sizes up to 15kb (UGP\_052, UGP\_053,  
199 UGP\_052 in the genome of crAssphage *sensu stricto*), possibly coding for fused subunits of  
200 RNA polymerase (16), as well as an absence or scarcity of tRNA genes. By contrast,  
201 members of candidate genus VI had two genome regions of approximately equal size with  
202 opposite gene orientation and G+C skew and large sets of tRNA genes (up to 27;  
203 Supplementary Table 2). A prominent common feature of the members of candidate genera  
204 VII and VIII was absence of the giant open reading frames.

205 In order to further demonstrate the homogeneity of the candidate genus I of crAss-like  
206 phages, comparative genomic analysis was performed on complete genomes. We  
207 characterised crAss-like phages as having pac-type circularly permuted genomes (24, 25);  
208 therefore, only genomes determined as circular were considered for this analysis. The  
209 genomic start coordinates of circular Candidate Genus I crAss-like phages were altered to  
210 match that of the published prototypical crAssphage *sensu stricto*. Candidate Genus I crAss-  
211 like phages showed high levels of synteny and strong homology across their entire genomes.  
212 However, the most notable area of diversity is observed in the crAss-like phage putative  
213 receptor binding protein (UGP\_074), which likely targets the different crAssphage strains  
214 towards their specific bacterial hosts (Supplementary Figure 3).

215 **Prevalence of crAss-like phages in human faecal virome samples.** To get insights  
216 into relative abundance of different crAss-like phages in various human populations we  
217 aligned quality filtered reads, representing 532 human faecal samples from the same datasets  
218 as used for assembly of crAss-like genomes, to a database of 93 nonredundant crAss-like  
219 phage genomic sequences (with <90% of homology and/or <90% overlap between them)  
220 representing all 10 candidate genera.

221 CrAss-like phage colonization rates varied from 51-58% in Malawian infants to 98-  
222 100% of healthy individuals of various ages in the Western cohorts. While relative crAss-like  
223 phage content ranged from 0 to 87% of the reads per sample, and depended significantly on  
224 the country of residence ( $p = 6.5E-09$  in Kruskal-Wallis test) and age group of the donor ( $p =$   
225  $1.6E-10$ ). In ~8% of all virome samples, >50% of reads aligned to crAss-like phage genomes.  
226 Lowest overall crAss-like phage counts were seen in healthy Irish and Malawian infants and  
227 in USA adults with IBD (Figure 4A). On a global scale, crAss-like candidate genera I, III,  
228 and VIII seem to be the most prevalent ones (Figure 4B).

229 The specific composition of crAss-like phages in faeces partly separated a cohort of  
230 healthy and malnourished infants living in rural areas of Malawi from the healthy and  
231 diseased urban Western cohorts (Figure 4C). PERMANOVA analysis suggested that crAss-  
232 like phage composition was mostly driven by place of residence ( $R^2 = 0.24$ ,  $p = 0.001$ ) with  
233 condition and age group also having significant impact ( $R^2 = 0.05$  and  $0.01$  respectively,  $p =$   
234  $0.001$ ). This observation is further supported by a clear difference in the distribution of  
235 specific crAss-like candidate genera across different populations (Figure 5). Specifically,  
236 Candidate Genus I, which includes crAssphage *sensu stricto* is by far the most prevalent type  
237 of crAss-like phages in Western population regardless of age. At the same time, same genus  
238 was extremely scarce in Malawian cohort where Candidate Genus III and VIII were the most  
239 common ( $p = 6.7E-03$  and  $1.4E-06$ , respectively).

240 **Faecal fermentations of a crAssphage rich sample.** During an ongoing longitudinal  
241 study of faecal viromes in healthy adults we identified one individual (subject ID 924), in  
242 which crAssphage *sensu stricto* was consistently contributing >30% of virome metagenomic  
243 reads over a 12 month period. Thus, this donor was selected in order to investigate if  
244 crAssphage *sensu stricto* could be propagated in a batch faecal fermentation system.  
245 Quantitative PCR (qPCR) detection of a conserved fragment of the crAssphage *sensu stricto*  
246 DNA polymerase gene in the viral nucleic acid fractions throughout the fermentation revealed  
247 that crAssphage *sensu stricto* was effectively propagated. CrAssphage *sensu stricto* was  
248 found to increase in titre by 89 fold for up to 21 hours into the fermentation (Figure 6A).

249 Interestingly, shotgun metagenomic sequencing of the viral enriched DNA from the  
250 fermentation supernatants showed the presence of six other crAss-like phages in the study  
251 subject, in addition to crAssphage *sensu stricto* (Supplementary Table 2). These crAss-like  
252 phage contigs were all  $\geq 70$ kb and grouped into five of the candidate genera (Figure 6B), four  
253 of which contributed to  $\geq 1\%$  of the reads per sample. The most abundant crAss-like contig of  
254 subject ID 924, designated as Fferm\_ms\_6 (linear, 90.4kb), is a member of proposed  
255 Candidate Genus I and closely related to crAssphage *sensu stricto*. Contig Fferm\_ms\_2  
256 (linear, 88.8 kb) is the second most abundant in the sample and belongs to Candidate Genus  
257 V. Other crAss-like phages showed varying degrees of similarity at the amino acid level to  
258 different crAss-like phage at the genus-level taxonomic groups. Analysis of bacterial  
259 microbiota in the fermentation vessel using compositional 16S rRNA gene amplicon  
260 sequencing revealed a concomitant increase in the course of fermentation of a number of  
261 *Bacteroides* species, including; *B. dorei*, *B. uniformis*, *B. fragilis*, *B. xylanisolvens*, *B. nordii*,  
262 *Parabacteroides distasonis* and *Parabacteroides chinchillae* (Supplementary Figure 4).

263 **Biological characterisation of crAss-like phages.** Transmission electron microscopy  
264 (TEM) of a crAssphage *sensu stricto* rich faecal filtrate showed a significant presence of

265 short-tailed or non-tailed viral particles with icosahedral or isometric heads (53% of  
266 *Podoviridae* type and 29% of *Microviridae* or a smaller type of *Podoviridae*), with lower  
267 levels of tailed bacteriophages of the family *Siphoviridae* (15%; Figure 7A). *Podoviridae*-  
268 type virions could be further classified into two types: type I, with head diameters of ~76.5  
269 nm and short tails; and type II, with a similar head size but head-tail collar structures and  
270 slightly longer tails (Figure 7B). Sequencing of the same fraction as used for the TEM  
271 showed that approximately 40% of reads aligned to crAss-like genomic contigs (Figure 7D).  
272 Based on the size of crAss-like genomic contigs assembled from subject ID 924 samples  
273 (88.8-97.3 kb), it seems likely that the predominant *Podoviridae* morphology observed  
274 corresponds to the crAss-like group of bacteriophages. For comparison, *Microviridae* phages  
275 have genomes 4.4-6.1 kb and icosahedral capsids of approx. 15-30 nm in diameter (26, 27).

276 The same CsCl fraction that was subjected to metagenomic sequencing and TEM  
277 visualisation was also analysed by SDS-PAGE followed by identification of major bands  
278 using MALDI-TOF mass spectrometry. A major structural protein of a crAss-like phage,  
279 denoted as Fferm\_ms\_2\_MCP, was detected following MALDI-TOF analysis of a band  
280 excised from the ~55kDa area on a SDS-PAGE gel (Figure 7C). The obtained peptide profile  
281 corresponded to a protein of 490 amino acids and 55.4 kDa, encoded by Fferm\_ms\_2. Further  
282 analyses using BLASTp showed the protein to have 37% identity with UGP\_086, predicted  
283 as the major capsid protein of the prototypical crAssphage (16).

284 In addition, we attempted to independently establish the size of crAss-like phage  
285 virions by passing faecal filtrates through a series of filters with gradually decreasing pore  
286 sizes (Supplementary Figure 5). Filtration through 0.1  $\mu\text{m}$  pores (equivalent to 100 nm)  
287 resulted in partial retention of crAss-like phages while pores of 0.02  $\mu\text{m}$  size completely  
288 removed crAssphage from the filtrate, as judged by the qPCR assay.

289

## 290 Discussion

291 The overall objective of this study was to gain a more in depth insight into one of the  
292 most enigmatic phages discovered to date, crAssphage. This phage is highly abundant in the  
293 human microbiome on a global scale; however, it remains poorly understood. One reason  
294 why crAssphage has remained such a mystery is due to the lack of available genome  
295 sequences for comparison. When crAssphage was assigned a specific nomenclature and  
296 uploaded to a public repository by Dutilh and colleagues (12), it became a template for other  
297 studies to compare against. This highlights the need for researchers to upload both the  
298 sequencing reads and assembled contigs following metagenomic studies.

299 CrAssphage is a representative of an expanding group of human gut-associated  
300 bacteriophages. While previous studies have proposed a sequence-based classification of  
301 crAss-like viruses at the familial level (16), our *in silico* analysis fits within classical familial  
302 taxonomic assignments whereby crAss-like phages are categorised as *Podoviridae*. In this  
303 study, we present 243 new crAss-like phage genomes from various metagenomic studies.  
304 Comparative genomics of the 244 available crAss-like phages demonstrates an extensive  
305 degree of diversity among these phages, including the potential identification of four crAss-  
306 like phage subfamilies. While the *alphacrAssvirinae* subfamily is currently the largest of the  
307 4 subfamilies, future studies looking for additional homologues of *betacrAssvirinae*,  
308 *gammacrAssvirinae* and *deltacrAssvirinae* members will refine these taxonomic categories.

309 Assigning phage taxonomy, in the absence of a universal genetic marker such as 16S  
310 rRNA, is a difficult and potentially erroneous process. In our study, we adopted a method  
311 previously employed to assign taxonomic ranks to *Podoviridae* based on the percentage of  
312 shared homologous proteins (17). This categorisation strategy identified 10 candidate genera,  
313 with crAss-like phages in each genera originating from the faeces of putatively healthy  
314 individuals and people suffering from various diet and bowel-related disorders. Alternative

315 proposed methods for defining phage genera include grouping phages with >50% nucleotide  
316 similarity identity together (18). Noteworthy, the 10 proposed crAss-like phage genera as  
317 determined by percentage of shared homologous proteins closely resembles that observed for  
318 crAss-like phage groups when characterised by >50% shared average nucleotide identity .

319 Several crAss-like phage genera proposed in this study have distinct nucleotide G+C  
320 compositions. The nucleotide composition of obligate parasites, such as phages, likely  
321 evolves in close association with the host bacterium (23, 28–30). Thus, Candidate Genera III  
322 and VI with diverse G+C compositions are either heterogeneous groups of crAss-like phages  
323 that require further sequences to refine their taxonomic structure, or they are potentially  
324 capable of infecting across a broad host range.

325 Quantitative analysis of crAss-like phage content in several cohorts revealed that in  
326 agreement with the previous studies the vast majority of faecal viral metagenomic samples  
327 contained varied amounts of crAssphage DNA. CrAssphage *sensu stricto* (Candidate Genus  
328 I) is by far most predominant type in Western populations, co-existing with other crAss-like  
329 phages in the majority of samples. By contrast, in the cohort of malnourished and healthy  
330 Malawian infants (9, 31), other candidate genera such as III, VIII and IX seem to play the  
331 leading role. It is well known that non-Western rural populations, which mostly consume high  
332 fibre, low fat and low animal protein diet are predominantly associated with high  
333 *Prevotella*/low *Bacteroides* type of gut microbiota (known as enterotype II (32)), as opposed  
334 to *Bacteroides*/*Clostridia*-dominated microbiota (enterotype I) in urban populations  
335 consuming western diet (33, 34). Indeed, our analysis of the Reyes et al. (2015) 16S rRNA  
336 gene sequencing data confirmed high prevalence of *Prevotella* in Malawian samples  
337 (Supplementary Figure 6). One can hypothesize that members of candidate genera III, VIII  
338 and IX might be associated with *Prevotella* or other members of the order *Bacteroidales* apart  
339 from *Bacteroides sensu stricto*.



340 The *in vitro* analysis of samples obtained from subject ID 924 was particularly  
341 intriguing. By mapping metagenomic sequencing reads against crAssphage *sensu stricto*, it  
342 was initially thought that this donor only carried the prototypical crAssphage at levels  
343 exceeding 30% of total viral reads for a 1 year period. A subsequent mining for phages  
344 related to crAssphage *sensu stricto* using metagenomic sequencing at later time points, with  
345 and without multiple displacement amplification resulted in 5 additional crAss-like phages  
346 being simultaneously detected from a single donor. However, the initial screening and  
347 inclusion criteria for bioinformatic detection of crAss-like phages resulted in a fragmented  
348 crAss-like phage contig being missed. The overlooked crAss-like phage, Fferm\_ms\_2  
349 (Candidate Genus V), turned out to be extremely important during the *in vitro* biological  
350 characterisation experiment. Therefore, it is possible many additional crAss-like phage  
351 genomes could be present within the metagenomic datasets that were examined in this study,  
352 but they were not included in our analysis because of the inclusion criteria chosen or even the  
353 choice of assembly program.

354 In total, subject ID 924 consistently carried 7 crAss-like phages, which resolved in our  
355 taxonomic analysis into 5 candidate genera. Three of the crAss-like phages were identified in  
356 Candidate Genus VI, supporting the notion this is a heterogeneous group and not simply  
357 composed of broad host range infecting phages. It is possible that there are potentially more  
358 than 7 crAss-like phages within subject ID 924. However, we believe that only a single  
359 representative of each candidate crAss-like phage genus (with the exception of the  
360 heterogeneous candidate genus VI) could assemble correctly, with two or more highly  
361 identical phages amalgamating their single nucleotide polymorphisms into a single consensus  
362 representative sequence (Supplementary Figure 7).

363 This study demonstrates the proliferation of crAss-like phages in a faecal fermenter, the  
364 first evidence of crAss-like phage propagation in the laboratory. Furthermore, following our



365 ability to propagate faecal crAss-like phages, we conducted the first transmission electron  
366 micrographs (TEMs) of these phages. Indeed, the most abundant faecal viruses present in  
367 samples used to inoculate faecal fermentation were *Podoviridae*. This is in agreement with  
368 the predictions made by Yutin *et al.*, following their detailed genome annotation of two  
369 crAss-like phages (16). Interestingly, however, our TEMs suggest presence of two types of  
370 virions with short non-contractile tails (Figure 7C). Presumably, the more abundant type I  
371 virions with shorter tail can belong to members of Candidate Genus I, also found as the most  
372 abundant crAss-like phage group in subject ID 924 by means of metagenomic sequencing  
373 (Figure 6B). Whereas type II virions with slightly longer tails and visible head-tail collar  
374 structures may correspond to Candidate Genus VI, found as the second most abundant crAss-  
375 like phage subfamily in shotgun metagenomics. But without isolating these phages in pure  
376 culture, it is not possible to accurately assign which *Podoviridae* tail corresponds to which  
377 specific crAss-like phage subfamily or genera.

378 This work provides the first *in vitro* evidence confirming that crAss-like phages are  
379 members of the *Podoviridae* family. This is shown from three levels of experimentation using  
380 the same CsCl fraction purified from crAssphage rich faeces of a healthy human donor. The  
381 TEM images produced from the CsCl fraction showed an abundance of the signature  
382 *Podoviridae* morphology. Other phage capsids present, predominantly *Microviridae*, would  
383 typically be associated with smaller genome sizes than that of crAss-like phages (26).  
384 Sequencing of the same fraction identified that almost 40% of the reads aligned to crAss-like  
385 phages. This is consistent with the percentage of *Podoviridae* identified in the TEM images.  
386 Furthermore, a highly predominant protein denoted as Fferm\_ms\_2\_MCP, was isolated from  
387 the fraction and was found to have significant similarity to crAss-like phages of (Candidate  
388 Genera V) as well as a moderate degree of similarity to crAssphage *sensu stricto* (Candidate

389 Genera I). This *in vitro* evidence, in line with the taxonomic analysis performed by Yutin *et*  
390 *al.*, proves that crAss-like phages do indeed belong to the *Podoviridae* family.

391 Identifying a means of propagating crAss-like phages is of particular importance.  
392 However, it was also observed that the primers applied in the qPCR analyses of viral nucleic  
393 acids were not suitable for targeting crAss-like phages associated with the various  
394 subfamilies and candidate genera that differed significantly from crAssphage *sensu stricto*.  
395 With the availability of more crAss-like phage sequences, broad and narrow spectrum  
396 primers can now be designed and applied in the analysis of these phages. The choice of  
397 primers for detecting crAss-like phages was also discussed in the recent work of Cinek *et al.*  
398 (14). This will be an important part of further work.

399 It also has to be considered that human gut crAssphage is not one single entity, but  
400 rather a group of diverse viruses, sharing certain signature genomic traits. It is most likely  
401 that these diverse phages target multiple bacterial taxa. Previously, a member of the  
402 *Bacteroides* genus was hypothesised as being the host for crAssphage (12). In a study prior to  
403 the discovery of crAssphage (35), a 95.9kb contig corresponding to a putative virus  $\phi$ HSC05  
404 was shown to be stably engrafted after transplantation of human faecal virus fraction into  
405 germ-free mice colonized with an artificial defined community of 15 bacterial species. The  
406 artificial bacterial community, among others, included: *Bacteroides thetaiotaomicron* (2  
407 strains), *B. caccae*, *B. ovatus*, *B. vulgatus*, *B. cellulosilyticus* and *B. uniformis*. One might  
408 conclude that one of the above mentioned 7 strains of the genus *Bacteroides*, more likely than  
409 the remaining 8 strains of Gram-positive anaerobic bacteria used in that study, must have  
410 served as a host for crAssphage propagation. The retrospective analysis of contigs from that  
411 study conducted by ourselves showed that the  $\phi$ HSC05 contig was 91.73% identical by its  
412 nucleotide sequence to crAssphage *sensu stricto*. Since crAssphage had not been described at

413 the time the article was published, this very interesting observation was never made by the  
414 authors of the original work.

415 With more divergent sequences, we could assume that different members of the  
416 *Bacteroides* genus, or even *Bacteroidetes* phylum for example, may serve as hosts for  
417 different crAss-like phages. One host that has been hypothesised for prototypical crAss-like  
418 phages is *B. dorei*. This was inferred following the analysis of a dataset generated from  
419 infants and toddlers with islet autoimmunity. It was correlated that crAssphage was only  
420 present when *B. dorei* also was detected within the samples. This was not true for other  
421 *Bacteroides* members tested, including *B. vulgatus* which is highly related to *B. dorei*. This  
422 correlation is compelling; however, it should be noted that there was no confirmation that  
423 crAssphage has any role in causing bacteriome alterations that lead to islet autoimmunity  
424 (36). Interestingly, one of the key *Bacteroides* species detected from our faecal fermentation  
425 16S rRNA analysis was *B. dorei*. Its levels were inversely proportional to that of crAssphage.  
426 Therefore, this possible phage-host pair should be investigated further.

427 CrAss-like phages have also been defined as a part of the core human gut phageome  
428 (10). This emphasises the importance of identifying hosts for diverse crAss-like phages  
429 belonging to different candidate genera proposed in this study. Such knowledge along with  
430 the ability to propagate crAss-like phages *in vitro* will provide an insight into its biological  
431 significance including their possible role in shaping the bacterial composition of the human  
432 gut microbiome in a positive or negative manner, in context of various disease states, such as  
433 inflammatory bowel disease, cancer, and obesity among others. Thus far, only a few studies  
434 has attempted to correlate crAss-like phages with a gastrointestinal disorder (7, 13, 36).  
435 Exploring this aspect of crAss-like phages further will be a key part of future work.

436 In conclusion, our results expand the repertoire of known crAss-like phages  
437 significantly, providing a path towards the identification of further crAss-like phages and their

438 hosts. This will lead to a better understanding of their role, if any, in human health and  
439 disease. Our work also provides an interesting insight into the diversity of these human gut-  
440 associated phages in various populations through *in silico* and *in vitro* methods. In addition,  
441 we also demonstrate that these enigmatic phages can be efficiently propagated *in vitro* in a  
442 mixed culture as well as present the first TEMs of crAss-like phages, giving an insight into  
443 their morphology. CrAss-like phages appear to be universally present in human populations,  
444 including various disease states. Due to the specificity of phage-host interactions, the  
445 diversity of crAss-like phages suggests they infect multiple diverse bacteria of the human  
446 gastrointestinal microbiota. However, more studies will be required to determine the  
447 biological significance and role of crAss-like phages in the human gut and determine if its  
448 presence positively or negatively impacts human gastrointestinal health.

## 449 **Methods**

450 **Metagenomic datasets and contig assemblies.** Sequencing reads from publicly  
451 available metagenomic datasets were downloaded from NCBI Sequence Read Archive (SRA)  
452 database. All published and unpublished metagenomic datasets that yielded crAss-like phage  
453 contigs, the DNA preparation protocol, the sequencing technology, the assembly program,  
454 and information related to contig nomenclature, are briefly described in Supplementary Table  
455 1. All reads were processed using Trimmomatic v0.32 to remove adaptor sequences and to  
456 trim reads when the Phred quality score dropped below 30 for a 4bp sliding window.  
457 Trimmed reads were assembled using either SPAdes v3.6.2 (37) or metaSPAdes v3.10.0 (38).  
458 Contigs from the assembly of 702 metagenomic samples were assigned a specific  
459 nomenclature, representing: [1] study/sample description, [2] SPAdes or metaSPAdes  
460 assembly, and [3] numerical rank of largest-to-smallest assembled contigs. The full list of  
461 contigs assembled in this study, the available associated metadata, and contig accession  
462 numbers, are detailed in Supplementary Table 2.

463 **Detection and curation of crAss-like phages.** The detection of crAss-like phage  
464 contigs was performed as follows. The amino acid polymerase sequence of prototypical  
465 crAssphage (UGP\_018, NC\_024711.1) was queried using BLAST v2.2.28+ (39) against a  
466 translated nucleotide database consisting of assembled metagenome contig sequences. The  
467 most conserved orthologous protein group detected in our initial putative crAss-like phage  
468 screening included prototypical crAssphage protein UGP\_092, which was annotated through  
469 the HHPred homology and structural prediction web server (40) as a phage terminase. This  
470 was then used as a second genetic signature of crAss-like phages and used in an additional  
471 BLAST search. All putative crAss-like phages selected for analysis met the following  
472 criteria: [1] a BLAST hit against either prototypical crAssphage polymerase or terminase

473 with an e-value less than 1e-05, [2] a BLAST query alignment length  $\geq 350$ bp, and [3] a  
474 minimum contig length of 70kb (representing near-complete crAss-like phage contigs).

475 **Identification of crAss-like phage orthologous proteins and clusters.** The encoded  
476 proteins of crAss-like phages were predicted using Prodigal v2.6.3 (41). Orthologous proteins  
477 shared between crAss-like phages were detected using OrthoMCL v2.0 using default  
478 parameters (42). The presence/absence of orthologous proteins between crAss-like phages was  
479 initially converted into a binary count matrix where the percentage of shared orthologous  
480 proteins was calculated (Figure 1B). The optimum number of phage clusters was calculated  
481 using the percentage of shared homologous proteins using the NbClust v3.0 package for R  
482 (43). Hierarchical clustering was performed on the count matrix of percentage shared crAss-  
483 like phage orthologous proteins using Ward's minimum variance method ['Ward.D2'  
484 algorithm in R (44)]. The resulting dendrogram was cut at  $k = 10$  based on the estimation of  
485 the number of crAss-like phage clusters (Figure 1A).

486 As a verification of the 10 predicted crAss-like phage clusters, the original abundance  
487 matrix of crAss-like phage orthologous proteins was used to calculate Euclidean distances  
488 between samples. These distance variations were calculated using the t-SNE machine  
489 learning algorithm ['tsne' v0.1-3 for R; (45)] and plotted using ggplot v2.2.1 (Figure 2). The  
490 presence or absence of orthologous protein groups was used to determine the core proteome  
491 of crAss-like phage clusters (Supplementary Figure 8).

492 **Phylogeny of crAss-like phage terminase sequences.** Following the work of Yutin *et*  
493 *al.*, (16) all publically available crAss-like phage terminase sequences were included in an  
494 additional phylogenetic analysis (Supplementary Figure 2). The terminase amino acid  
495 sequences of crAss-like phages were aligned using Muscle v3.8.31 (46). The resultant  
496 alignment was converted to Phylip format and phylogeny was determined by PhyML using a  
497 JTT amino acid substitution model (47). The phylogenetic tree was visualised using FigTree

498 v1.4.3. The phylogenetic tree is coloured based on the crAss-like phage clustering analysis  
499 with node support values displayed.

500 **Genomic comparisons of crAss-like phages.** The average nucleotide identity between  
501 crAss-like phage contigs was calculated using Pyani v0.2.3 by the ANIm method with a  
502 500bp fragment size. Pairwise comparisons of complete crAss-like phage genomes belonging  
503 to Candidate Genera I was performed using Easyfig v2.2.2. Genomic start coordinates and  
504 contig orientations were altered to match the published GenBank sequence of prototypical  
505 crAssphage NC\_024711.1. The order of crAss-like phages in the Easyfig image was adjusted  
506 to match to the order they appear in the average nucleotide identity analysis (Figure 3). The  
507 Easyfig image was generated using tBLASTx comparisons, with a minimum BLAST length  
508 of 50bp and identity of 30bp (Supplementary Figure 3). The presence of crAss-like phage  
509 tRNA-encoding sequences were detected using ARAGORN v1.2.36 (48). To determine the  
510 genomic packaging mechanism of crAss-like phages, metagenomic sequencing reads from a  
511 TruSeq (Illumina) manually fragmented DNA library were analysed using PhageTerm (25).  
512 Single nucleotide polymorphisms (SNPs) of crAss-like phages were observed by aligning  
513 metagenomic sequencing reads to the consensus assembled contig sequence using Bowtie2  
514 and Samtools, and visualising SNPs using Tablet v1.17.08.17 (49).

515 **Alignment of virome metagenomic reads to crAss-like contigs.** The quality filtered  
516 reads from 532 human faecal viromes (as subset of 701 viromes selected based on availability  
517 of sufficient metadata) were then aligned to the set of 93 nonredundant crAss-like phage  
518 genomic (with <90% of homology and/or <90% overlap between them) using Bowtie2 v2.3.0  
519 (50) using the end-to-end alignment mode. A count table was generated with Samtools  
520 v0.1.19 which was then imported into R v3.3.1 for statistical analysis.  $\beta$ -diversity of crAss-  
521 like viral populations in human cohorts was visualized using PCoA plot based on Spearman  
522 rank distances ( $D = 1 - \rho$ , where  $\rho$  is Spearman rank correlation coefficient of relative



523 abundance of different crAss-like contigs between samples). Statistical analysis was  
524 performed using permutational multivariate analysis of variance (PERMANOVA)  
525 implemented in Vegan v2.4.3 package for R (51) and non-parametric Kruskal-Wallis test.

526 **Recruitment of a crAssphage faecal donor and faecal fermentations.** Human faecal  
527 viromes from a number of ongoing studies sequenced using Illumina HiSeq and MiSeq  
528 platforms were screened for crAss-like phages by aligning the obtained sequencing reads  
529 against prototypical crAssphage NC\_024711.1 using Bowtie2 v2.3.0. One individual (subject  
530 ID 924) was found to carry crAssphage consistently at levels exceeding 30% of the total  
531 number of reads over a one year period. The recruited individual is an adult female that  
532 suffers from gastritis and is vitamin B12 deficient. A frozen standard inoculum (FSI) sample  
533 was processed as described by (52) with the following modification: the sample was  
534 resuspended in 1X phosphate buffered saline (37 mM NaCl, 2.7 mM KCl, 8 mM Na<sub>2</sub>HPO<sub>4</sub>,  
535 and 2 mM KH<sub>2</sub>PO<sub>4</sub>), 0.05% (w/v) L-cysteine (Sigma Aldrich, Ireland) and (1 mg/L)  
536 resazurin (Sigma Aldrich, Ireland). The crAssphage-rich FSI was inoculated into 400 ml  
537 YCFA-GSCM broth in a 500 ml fermenter vessel at 5% (v/v). Fermentation media was  
538 prepared exactly as described by (53) with the addition of glucose (2 g/L), soluble starch (2  
539 g/L), cellobiose (2 g/L) and maltose (2 g/L). Fermentation was performed in batch format at  
540 approximately 37°C for 51 hours. Dissolved oxygen was sustained at <0.1% by constantly  
541 sparging the vessel with anaerobic gas mix (80% (v/v) N<sub>2</sub>, 10% (v/v) CO<sub>2</sub>, 10% (v/v) H<sub>2</sub>) and  
542 stirring at 200 rpm. Both 2M NaOH and HCl solutions were used to maintain pH at ~7.  
543 Samples were collected at the following time points; 0, 4, 21, 28, 45 and 51 hours. Collected  
544 samples were centrifuged at 4,700 rpm at +4°C for 10 minutes. The resulting supernatants  
545 were filtered once through a 0.45 µm pore syringe filter and stored at +4°C. Resultant pellets  
546 were stored at -80°C.



547           **Extraction of viral nucleic acids and sequencing library preparation.** Total virome  
548           extractions were performed on 0.45  $\mu$ M pore filtered fermentation supernatants. Solid NaCl  
549           and polyethylene glycol 8000 were added to the filtrates to give a final concentration of 0.5M  
550           and 10% (w/v), respectively. After overnight incubation at +4°C samples were centrifuged at  
551           4,700 rpm and +4°C for 20 minutes. The pellets were then resuspended in 400 $\mu$ l of SM buffer  
552           (1M Tris-HCl pH 7.5, 5M NaCl, 1M MgSO<sub>4</sub>) and briefly vortexed with an equal volume of  
553           chloroform. This mixture was then centrifuged at 2,500g for 5 minutes using a standard  
554           desktop centrifuge. The resultant aqueous phase was then transferred into an Eppendorf to  
555           which 40 $\mu$ l DNase buffer (10mM CaCl<sub>2</sub> and 50mM MgCl<sub>2</sub>) and 8U and 4U TURBO DNase  
556           (Ambion/ThermoFisher Scientific) and RNase I (ThermoFisher Scientific) were added,  
557           respectively. This was incubated at 37°C for 1 hour followed by an enzyme inactivation step  
558           at 70°C for 10 minutes. This was followed by the addition of 2 $\mu$ l proteinase K and 10% SDS  
559           and further incubation at 56°C for 20 minutes. Lastly, 100 $\mu$ l phage lysis buffer (4.5 M  
560           guanidinium isothiocyanate, 44 mM sodium citrate pH 7.0, 0.88% sarkosyl, 0.72% 2-  
561           mercaptoethanol) was added to lyse the viral particles. The final incubation was carried out at  
562           65°C for 10 minutes. The resulting lysates were lightly vortexed with an equal volume of  
563           phenol/chloroform/isoamyl alcohol 25:24:1 (Fisher Scientific) and were centrifuged at room  
564           temperature for 5 minutes at 8,000g. This was again repeated with the resulting aqueous  
565           phase. Following the second extraction, the aqueous phase was passed through a DNeasy  
566           Blood and Tissue Kit (Qiagen) for final lysate purification. The wash steps were each  
567           repeated twice and the final elution was carried out in 50 $\mu$ l elution buffer. Viral DNA  
568           quantification was carried out with the Qubit HS DNA Assay Kit (Invitrogen/ThermoFisher  
569           Scientific) in a Qubit 3.0 Fluorometer (Life Technologies). The viral nucleic acids were then  
570           subjected to reverse transcription using SuperScript IV Reverse Transcriptase (RT) kit  
571           (Invitrogen/ThermoFisher Scientific). The protocol was carried out exactly as described in

572 the manufacturer's protocol for random hexamer primers. Following this, 1µl of the reversed  
573 transcribed viral DNA was subjected to GenomiPhi V2 (GE Healthcare) Multiple  
574 Displacement Amplification (MDA). Finally, MDA and non-MDA viral DNA was prepared  
575 for sequencing using TruSeq DNA Library Preparation Kit (Illumina, Ireland). All steps were  
576 performed as per the manufacturer's instructions. Prepared libraries were sequenced on an  
577 Illumina HiSeq platform (Illumina, San Diego, California) with 2x300bp paired-end  
578 chemistry at GATC Biotech AG, Germany. Reads were filtered, trimmed and assembled into  
579 contigs as described above. A count matrix was created by aligning quality-filtered reads back  
580 to contigs using Bowtie2 and Samtools.

581 **CrAssphage PCR detection.** Two oligonucleotide primer pairs were designed based  
582 on the prototypical crAssphage DNA polymerase sequence UGP\_018 (1) using PerlPrimer  
583 software (54). Primer sequences are as follows: CrAss-Pol-F5 5'-  
584 GCCTATTGTTGCTCAAGCTATTGAA-3', CrAss-Pol-R5 5'-  
585 ACAACAGAACCAGCTGCCAT-3', CrAss-Pol-F6 5'-  
586 AGTGGTCTTGCTCCNGAACAATGG-3' and CrAss-Pol-R6 5'-  
587 AACCTCCAGTTGCAACAGTATAAGT-3'. PCR products were cloned into pCR2.1-TOPO  
588 TA vector (ThermoFisher Scientific) and obtained plasmids at known concentrations were  
589 used to establish calibration curves through serial two-fold dilutions. Subsequently, qPCR  
590 were run in 15µl reaction volumes using SensiFAST SYBR No-ROX mastermix and  
591 LightCycler 480 thermocycler with the following conditions: initial denaturation at 95°C for  
592 5 minutes, then 35 cycles of 94°C for 20 seconds, 55°C for 20 seconds and 72°C for 20  
593 seconds, with a final extension at 72°C for 7 minutes. All samples were run in triplicate and  
594 the standard error was determined following calculation of DNA concentration based on the  
595 above standard curve.

596           **Electron microscopy and detection of crAssphage proteins.** A virus-enriched  
597 fraction of the crAssphage positive faecal sample, collected from subject ID 924, was  
598 prepared for electron microscopy imaging as follows. A 1:20 suspension (w/v) of faeces was  
599 prepared in SM buffer followed by vigorous vortexing until homogenised. The homogenised  
600 sample was chilled on ice for 5 minutes prior to centrifugation twice at 4,700 rpm for 10  
601 minutes at +4°C. The resulting supernatant was then filtered twice through a 0.45 µM pore  
602 syringe filters. The filtrate was ultra-centrifuged at 120,000g for 3 hours using a F65L-6x13.5  
603 rotor (ThermoScientific). The resulting pellets were resuspended in 5 ml SM buffer. The viral  
604 suspensions were ultracentrifuged again by overlaying them onto a caesium chloride (CsCl)  
605 step gradient of 5M and 3M, followed by centrifugation at 105,000g for 2.5 hours. A band of  
606 viral particles visible under side illumination was collected and buffer-exchanged using 3  
607 sequential rounds of 10-fold diluting and concentrating to the original volume by ultra-  
608 filtration using Amicon Centifugal Filter Units 10,000 MWCO (Merck). The purified fraction  
609 was then analysed by qPCR for the presence of crAssphage as described above. Following  
610 this, 5µl aliquots of the viral fraction were applied to Formvar/Carbon 200 Mesh, Cu grids  
611 (Electron Microscopy Sciences) with subsequent removal of excess sample by blotting. Grids  
612 were then negatively contrasted with 0.5% (w/v) uranyl acetate and examined at UCD  
613 Conway Imaging Core Facility (University College Dublin, Dublin, Ireland) by transmission  
614 electron microscope. The faecal viral fraction from subject ID 924 was further concentrated  
615 using Amicon Ultra-0.5 Centrifugal Filter Unit with 3 kDa MWCO membrane (Merck,  
616 Ireland). This concentrated fraction was loaded onto a premade Bolt 4-12% Bis-Tris Plus  
617 reducing SDS-PAGE gel (Invitrogen) and separated at 200 V for 30 minutes using 1X  
618 NuPAGE MOPS SDS Running Buffer. Six brightest bands with approximate molecular  
619 weights of 28, 35, 45, 55, 120 and 200 kDa were excised and subjected to MALDI-TOF/TOF

620 (Bruker ultraflex III) protein identification following in-gel trypsinization, at Metabolomics  
621 & Proteomics Technology Facility (University of York, York, UK).

622 **16S rRNA gene library preparations.** Total DNA was extracted from the pellets  
623 formed following centrifugation of fermentation samples. This was carried out using the  
624 QIAamp Fast DNA Stool Mini Kit (Qiagen, Hilden, Germany). All steps were carried out as  
625 per the manufacturer's protocol with the addition of a bead-beating step to aid total DNA  
626 extraction from the bacterial cells. Approximately 200mg of each pellet was placed in a 2ml  
627 screw-cap tube containing a mixture of one 3.5 mm glass bead, a 200µl scoop of 1mm  
628 zirconium beads and a 200µl scoop of 0.1mm zirconium beads (ThistleScientific) with 1ml of  
629 InhibitEX Buffer. Bead-beating was carried out three times for 30 seconds using the  
630 FastPrep-24 benchtop homogeniser (MP Biomedicals). Between each bead-beating the  
631 samples were cooled on ice for 30 seconds. The samples were then lysed at 95°C for 5  
632 minutes. All other steps were carried out as per the manufacturer's protocol. Following  
633 extraction of total bacterial DNA, the hypervariable regions of V3 and V4 16S ribosomal  
634 RNA genes were amplified from 15ng of the DNA using Phusion High-Fidelity PCR Master  
635 Mix (ThermoFisher Scientific) and 0.2µM of each of the following primers, containing  
636 Illumina-compatible overhang adapter sequences: 16S-FP: 5'-  
637 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCCTACGGGNGGCWGCAG-3' and  
638 16S-RP: 5'-  
639 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGACTACHVGGGTATCTAATCC-  
640 3'. The PCR program was run as follows: 98°C for 30 seconds, 25 cycles of 98°C for 10  
641 seconds, 55°C for 15 seconds and 72°C for 20 seconds, with a final extension of 72°C for 5  
642 minutes. The amplicons were then purified using Agencourt AMPure XP magnetic beads  
643 (Beckman-Coulter) followed by a second PCR to attach dual Illumina Nextera indices using  
644 the Nextera XT index kit v2 (Illumina). Purification was performed once again and the

645 libraries were quantified using a Qubit dsDNA HS Assay Kit. The libraries were then pooled  
646 in equimolar concentration and sent for sequencing on an Illumina MiSeq platform (Illumina,  
647 San Diego, California) at GATC Biotech AG, Germany. The quality of the raw reads were  
648 assessed with FastQC (v11.5) and initial quality filtering was performed using Trimmomatic  
649 v0.36. Filtered reads were imported into R (v3.4.3) for analysis with DADA2 v1.6.0. (55)  
650 Further quality filtering and trimming (maxN of 0 and a maxEE of 2) was carried out on both  
651 the forward and reverse reads with only retention in cases of pairs being of sufficient high  
652 quality. Error correction was performed on forward and reverse reads separately and  
653 following this, reads were merged. The resulting unique Ribosomal Variant Sequences  
654 (RSVs) were subjected to further chimera filtering using USEARCH v8.1 (56) with the  
655 Chimera-Slayer gold database v20110519. The retained, high quality, chimera-free, RSVs  
656 were classified with the RDP-classifier in mothur v1.34.4 (57) against the RDP database  
657 v11.4 (phylum to genus) and SPINGO (58) for species assignment. Plots were generated  
658 using the R package ggplot2 v2.2.1.

659

660 **Acknowledgements**

661 We would like to additionally thank Tom Sutton for his bioinformatic expertise and insightful  
662 discussions. This publication has emanated from research conducted with the financial  
663 support of Science Foundation Ireland (SFI) under Grant Numbers SFI/12/RC/2273,  
664 SFI/15/ERCD/3189 and SFI/14/SP APC/B3032, and a research grant from Janssen Biotech,  
665 Inc.

666

667 **Author contributions:**

668 EG and SRS performed the laboratory and bioinformatic work, respectively. AS assisted in  
669 both the laboratory and bioinformatic analyses. AGC performed the 16S analysis. FJR, LAD  
670 and EGT assisted in the design, implementation and interpretation of experiments. EG, AS  
671 and SRS wrote the paper and generated the figures. AGC, FJR, LAD and EGT reviewed  
672 drafts of the manuscript and provided constructive criticism for its improvement. PR and CH  
673 secured the funding and wrote the paper. All authors contributed to the analysis of the data.

674

675 **Conflict of interest:**

676 The authors declare no conflict of interest.

677

678 **Data deposition:**

679 The 244 crAss-like phage contigs analysed in this study have been submitted to GenBank and  
680 are currently under revision. Contigs are currently accessible at:

681 [https://figshare.com/articles/crAss-like\\_contigs\\_fasta\\_tar\\_gz/6098321](https://figshare.com/articles/crAss-like_contigs_fasta_tar_gz/6098321)

682 **References**

- 683 1. Reyes A, et al. (2010) Viruses in the fecal microbiota of monozygotic twins and their  
684 mothers. *Nature* 466(7304):334–338.
- 685 2. Frank DN, et al. (2011) Disease phenotype and genotype are associated with shifts in  
686 intestinal-associated microbiota in inflammatory bowel diseases. *Inflamm Bowel Dis*  
687 17(1):179–184.
- 688 3. Tremaroli V, Bäckhed F (2012) Functional interactions between the gut microbiota and  
689 host metabolism. *Nature* 489(7415):242.
- 690 4. Claesson MJ, et al. (2012) Gut microbiota composition correlates with diet and health in  
691 the elderly. *Nature* 488(7410):178–184.
- 692 5. Cryan JF, Dinan TG (2014) Mind-altering microorganisms: the impact of the gut  
693 microbiota on brain and behaviour. *Nat Rev Neurosci* 13(10):701–712.
- 694 6. Everard A, Cani PD (2013) Diabetes, obesity and gut microbiota. *Best Pract Res Clin*  
695 *Gastroenterol* 27(1):73–83.
- 696 7. Norman JM, et al. (2015) Disease-specific Alterations in the Enteric Virome in  
697 Inflammatory Bowel Disease. *Cell* 160(3):447–460.
- 698 8. Mills S, et al. (2013) Movers and shakers: influence of bacteriophages in shaping the  
699 mammalian gut microbiota. *Gut Microbes* 4(1):4–16.
- 700 9. Reyes A, et al. (2015) Gut DNA viromes of Malawian twins discordant for severe acute  
701 malnutrition. *Proc Natl Acad Sci U S A* 112(38):11941–11946.

- 702 10. Manrique P, et al. (2016) Healthy human gut phageome. *Proc Natl Acad Sci U S A*  
703 113(37):10400–10405.
- 704 11. Manrique P, Dills M, Young MJ (2017) The Human Gut Phage Community and Its  
705 Implications for Health and Disease. *Viruses* 9(6):141.
- 706 12. Dutilh BE, et al. (2014) A highly abundant bacteriophage discovered in the unknown  
707 sequences of human faecal metagenomes. *Nat Commun* 5:ncomms5498.
- 708 13. Liang YY, Zhang W, Tong YG, Chen SP (2016) crAssphage is not associated with  
709 diarrhoea and has high genetic diversity. *Epidemiol Amp Infect* 144(16):3549–3553.
- 710 14. Cinek O, et al. (2018) Quantitative CrAssphage real-time PCR assay derived from data of  
711 multiple geographically distant populations. *J Med Virol* 90(4):767–771.
- 712 15. Stachler E, et al. (2017) Quantitative CrAssphage PCR Assays for Human Fecal Pollution  
713 Measurement. *Environ Sci Technol* 51(16):9146–9154.
- 714 16. Yutin N, et al. (2018) Discovery of an expansive bacteriophage family that includes the  
715 most abundant viruses from the human gut. *Nat Microbiol* 3(1):38–46.
- 716 17. Lavigne R, Seto D, Mahadevan P, Ackermann H-W, Kropinski AM (2008) Unifying  
717 classical and molecular taxonomic classification: analysis of the Podoviridae using  
718 BLASTP-based tools. *Res Microbiol* 159(5):406–414.
- 719 18. Adriaenssens E, Brister JR (2017) How to Name and Classify Your Phage: An Informal  
720 Guide. *Viruses* 9(4):70.
- 721 19. Minot S, et al. (2011) The human gut virome: Inter-individual variation and dynamic  
722 response to diet. *Genome Res* 21(10):1616–1625.



- 723 20. García-Aljaro C, Ballesté E, Muniesa M, Jofre J (2017) Determination of crAssphage in  
724 water samples and applicability for tracking human faecal pollution. *Microb Biotechnol*  
725 10(6):1775–1780.
- 726 21. Liang Y, Jin X, Huang Y, Chen S (2018) Development and application of a real-time  
727 polymerase chain reaction assay for detection of a novel gut bacteriophage (crAssphage).  
728 *J Med Virol* 90(3):464–468.
- 729 22. Edwards RA, McNair K, Faust K, Raes J, Dutilh BE (2016) Computational approaches to  
730 predict bacteriophage–host relationships. *FEMS Microbiol Rev* 40(2):258–272.
- 731 23. Lucks JB, Nelson DR, Kudla GR, Plotkin JB (2008) Genome Landscapes and  
732 Bacteriophage Codon Usage. *PLOS Comput Biol* 4(2):e1000001.
- 733 24. Casjens SR, Gilcrease EB (2009) Determining DNA Packaging Strategy by Analysis of  
734 the Termini of the Chromosomes in Tailed-Bacteriophage Virions. *Bacteriophages,*  
735 *Methods in Molecular Biology*<sup>TM</sup>. (Humana Press), pp 91–111.
- 736 25. Garneau JR, Depardieu F, Fortier L-C, Bikard D, Monot M (2017) PhageTerm: a tool for  
737 fast and accurate determination of phage termini and packaging mechanism using next-  
738 generation sequencing data. *Sci Rep* 7(1):8292.
- 739 26. Zhong X, Guidoni B, Jacas L, Jacquet S (2015) Structure and diversity of ssDNA  
740 Microviridae viruses in two peri-alpine lakes (Annecy and Bourget, France). *Res*  
741 *Microbiol* 166(8):644–654.
- 742 27. Roux S, Krupovic M, Poulet A, Debroas D, Enault F (2012) Evolution and Diversity of  
743 the Microviridae Viral Family through a Collection of 81 New Complete Genomes  
744 Assembled from Virome Reads. *PLOS ONE* 7(7):e40418.

- 745 28. Pride DT, Wassenaar TM, Ghose C, Blaser MJ (2006) Evidence of host-virus co-  
746 evolution in tetranucleotide usage patterns of bacteriophages and eukaryotic viruses.  
747 *BMC Genomics* 7:8.
- 748 29. Roux S, Hallam SJ, Woyke T, Sullivan MB (2015) Viral dark matter and virus–host  
749 interactions resolved from publicly available microbial genomes. *eLife* 4:e08490.
- 750 30. Mavrich TN, Hatfull GF (2017) Bacteriophage evolution differs by host, lifestyle and  
751 genome. *Nat Microbiol* 2(9):17112.
- 752 31. Smith MI, et al. (2013) Gut Microbiomes of Malawian Twin Pairs Discordant for  
753 Kwashiorkor. *Science* 339(6119):548–554.
- 754 32. Arumugam M, et al. (2011) Enterotypes of the human gut microbiome. *Nature*  
755 473(7346):174–180.
- 756 33. Filippo CD, et al. (2010) Impact of diet in shaping gut microbiota revealed by a  
757 comparative study in children from Europe and rural Africa. *Proc Natl Acad Sci*  
758 107(33):14691–14696.
- 759 34. Gorvitovskaia A, Holmes SP, Huse SM (2016) Interpreting Prevotella and Bacteroides as  
760 biomarkers of diet and lifestyle. *Microbiome* 4:15.
- 761 35. Reyes A, Wu M, McNulty NP, Rohwer FL, Gordon JI (2013) Gnotobiotic mouse model  
762 of phage–bacterial host dynamics in the human gut. *Proc Natl Acad Sci U S A*  
763 110(50):20236–20241.
- 764 36. Cinek O, et al. (2017) Imbalance of bacteriome profiles within the Finnish Diabetes  
765 Prediction and Prevention study: Parallel use of 16S profiling and virome sequencing in

- 766 stool samples from children with islet autoimmunity and matched controls. *Pediatr*  
767 *Diabetes* 18(7):588–598.
- 768 37. Bankevich A, et al. (2012) SPAdes: A New Genome Assembly Algorithm and Its  
769 Applications to Single-Cell Sequencing. *J Comput Biol* 19(5):455–477.
- 770 38. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA (2017) metaSPAdes: a new versatile  
771 metagenomic assembler. *Genome Res* 27(5):824–834.
- 772 39. Altschul SF, et al. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein  
773 database search programs. *Nucleic Acids Res* 25(17):3389–3402.
- 774 40. Söding J, Biegert A, Lupas AN (2005) The HHpred interactive server for protein  
775 homology detection and structure prediction. *Nucleic Acids Res* 33(suppl\_2):W244–  
776 W248.
- 777 41. Hyatt D, et al. (2010) Prodigal: prokaryotic gene recognition and translation initiation site  
778 identification. *BMC Bioinformatics* 11:119.
- 779 42. Li L, Stoeckert CJ, Roos DS (2003) OrthoMCL: Identification of Ortholog Groups for  
780 Eukaryotic Genomes. *Genome Res* 13(9):2178–2189.
- 781 43. NbClust: An R Package for Determining the Relevant Number of Clusters in a Data Set |  
782 Charrad | Journal of Statistical Software doi:10.18637/jss.v061.i06.
- 783 44. Ward JHJ (1963) Hierarchical Grouping to Optimize an Objective Function. *J Am Stat*  
784 *Assoc* 58(301):236–244.
- 785 45. Maaten L van der, Hinton G (2008) Visualizing Data using t-SNE. *J Mach Learn Res*  
786 9(Nov):2579–2605.

- 787 46. Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high  
788 throughput. *Nucleic Acids Res* 32(5):1792–1797.
- 789 47. Guindon S, et al. (2010) New Algorithms and Methods to Estimate Maximum-Likelihood  
790 Phylogenies: Assessing the Performance of PhyML 3.0. *Syst Biol* 59(3):307–321.
- 791 48. Laslett D, Canback B (2004) ARAGORN, a program to detect tRNA genes and tmRNA  
792 genes in nucleotide sequences. *Nucleic Acids Res* 32(1):11–16.
- 793 49. Milne I, et al. (2010) Tablet—next generation sequence assembly visualization.  
794 *Bioinformatics* 26(3):401–402.
- 795 50. Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nat*  
796 *Methods* 9(4):357–359.
- 797 51. Anderson MJ (2001) A new method for non-parametric multivariate analysis of variance.  
798 *Austral Ecol* 26(1):32–46.
- 799 52. O’Donnell MM, et al. (2016) Preparation of a standardised faecal slurry for ex-vivo  
800 microbiota studies which reduces inter-individual donor bias. *J Microbiol Methods*  
801 129:109–116.
- 802 53. Duncan SH, Hold GL, Harmsen HJM, Stewart CS, Flint HJ (2002) Growth requirements  
803 and fermentation products of *Fusobacterium prausnitzii*, and a proposal to reclassify it as  
804 *Faecalibacterium prausnitzii* gen. nov., comb. nov. *Int J Syst Evol Microbiol* 52(6):2141–  
805 2146.
- 806 54. Marshall OJ (2004) PerlPrimer: cross-platform, graphical primer design for standard,  
807 bisulphite and real-time PCR. *Bioinformatics* 20(15):2471–2472.

- 808 55. Callahan BJ, et al. (2016) DADA2: High-resolution sample inference from Illumina  
809 amplicon data. *Nat Methods* 13(7):581–583.
- 810 56. Edgar RC (2010) Search and clustering orders of magnitude faster than BLAST.  
811 *Bioinforma Oxf Engl* 26(19):2460–2461.
- 812 57. Schloss PD, et al. (2009) Introducing mothur: open-source, platform-independent,  
813 community-supported software for describing and comparing microbial communities.  
814 *Appl Environ Microbiol* 75(23):7537–7541.
- 815 58. Allard G, Ryan FJ, Jeffery IB, Claesson MJ (2015) SPINGO: a rapid species-classifier for  
816 microbial amplicon sequences. *BMC Bioinformatics* 16:324.

817 **Figure Legends**

818 **Figure 1.** Determination of crAssphage candidate subfamilies and genera based on the  
819 percentage of shared protein-encoding genes. **(A)** The 4 red lines cut the hierarchical  
820 clustering dendrogram of crAss-like phage contigs, with Euclidean distances calculated  
821 between the percentages of shared protein-encoding genes, into the 4 proposed candidate  
822 subfamilies of crAss-like phages. The histogram insert (top-right) represents the calculated  
823 optimal number of crAss-like phage clusters. The 10 optimal crAss-like phage clusters  
824 represent the putative candidate genera, and are assigned specific colours. **(B)** Heatmap  
825 showing the percentage of shared protein-coding genes between crAss-like phage genomes.  
826 CrAss-like phages with 20-40% shared protein encoding genes are considered related at the  
827 subfamily level while phages with >40% similarity are believed to be related at the genus  
828 level, consistent with the calculated number of crAss-like phage clusters.

829 **Figure 2.** Two-dimensional ordination of crAss-like phages based on the abundance of their  
830 protein-encoded orthologous sequences was performed using t-SNE machine learning  
831 algorithm. **(A)** CrAss-like phages are coloured by candidate genus annotations and shape is  
832 determined by their origin. CrAss-like phages originating from individuals with kwashiorkor  
833 and marasmus, or lacking metadata, are grouped together as 'Other/Unknown'. **(B)** CrAss-  
834 like phages are coloured by the percentage G+C nucleotide composition of their contig, while  
835 shape represents complete (circular) or partial (linear) genomes.

836 **Figure 3.** Average nucleotide identity of crAss-like phage contigs. The column annotation  
837 colour scheme highlights the predicted crAss-like phage candidate genus annotations, while  
838 the coloured row annotation represents the origin of the respective crAss-like phage contig.

839 **Figure 4.** Prevalence of crAss-like phage in human faecal viromes. **(A)** Relative abundance  
840 of total crAss-like phage in several cohorts differing in age, health status and country of  
841 origin, based on the fraction of metagenomic reads aligned. Bars represent median relative

842 abundances, the values within boxes represent percentage of positive samples. **(B)** Relative  
843 abundance of specific crAss-like candidate genera in total human populations analysed. **(C)**  
844 PCoA plot of crAss-like phages based on Spearman rank distances.

845 **Figure 5.** Relative abundance of the ten candidate genera of crAss-like phages in six different  
846 human cohorts based on the fraction of metagenomic reads aligned. Bars represent median  
847 relative abundances, while values within boxes represent percentage of positive samples.

848 **Figure 6.** Analysis of crAss-like phage dynamics in a faecal fermenter. **(A)** Evidence of  
849 crAssphage *sensu stricto* propagation following *in vitro* fermentations (standard error, n=3).  
850 The level of crAssphage *sensu stricto* propagation was determined by qPCR analysis of viral-  
851 enriched DNA, respectively, using primers specific to a segment of the crAssphage *sensu*  
852 *stricto* DNA polymerase gene. **(B)** Six additional crAss-like phages, that group into five of  
853 the candidate genera, were identified following sequencing of the same viral-enriched DNA  
854 from the fermenter. The relative abundance of each of these crAss-like phages is skewed due  
855 to the biased amplification of other components of the viral-enriched DNA fraction that is  
856 associated with multiple displacement amplification.

857 **Figure 7.** CrAss-like phage morphology was examined using a CsCl fraction purified from a  
858 crAssphage rich faecal filtrate of donor subject ID 924. **(A)** Analysis of the fraction through  
859 transmission electron microscopy (TEM) was performed. The TEM images are largely  
860 dominated by *Podoviridae* (53%), *Microviridae* (29%), *Siphoviridae* (15%) and other phage  
861 morphologies (3%). **(B)** Further examination of the observed *Podoviridae* identifies two  
862 variants with differing tail morphologies. Both variants have head diameters of ~76.5 nm. **(C)**  
863 SDS-PAGE gel of the CsCl fraction. Six bands containing possible crAssphage proteins were  
864 excised and analysed by mass spectrometry. A protein, denoted as Fferm\_ms\_2\_MCP,  
865 isolated from the ~55 kDa (\*) band was found to have high sequence similarity with  
866 Candidate Genus V crAss-like phages. **(D)** Sequencing of the CsCl purified viral fraction,

867 without multiple displacement amplification, showed that approximately 40% the reads  
868 aligned to crAss-like phages.

869

### 870 **Supplementary Figure Legends**

871 **Supplementary Figure 1.** Phylogeny of crAss-like phage terminase protein sequences,  
872 including publically available terminase sequences from the Yutin *et al.* (2017)  
873 characterisation of familial-related crAss-like phages. The figure legend insert corresponds to  
874 the colour scheme of the 10 proposed candidate genera groupings. NC\_024711 crAssphage  
875 and IAS virus, discussed in the main text, are highlighted in red. Bootstrapping node support  
876 values are shown.

877 **Supplementary Figure 2.** Comparison of general structural feature of representative  
878 complete circular genomes of the 10 proposed genera of crAss-like bacteriophages.  
879 Innermost circle (green/blue), G+C skew; middle circle, G+C content deviation from mean  
880 value; outermost circle, protein-coding genes (CDS) located on positive (red) and negative  
881 (blue) DNA strands, respectively; and tRNA genes (orange).

882 **Supplementary Figure 3.** Comparison of circular Candidate Genus I crAss-like phage  
883 genomes. Start co-ordinates of crAss-like phage genomes were adjusted to match crAssphage  
884 *sensu stricto*. The order of crAss-like phage genomes was determined by the average  
885 nucleotide identity comparisons. Open reading frames corresponding to specific predicted  
886 phage structural proteins are highlighted.

887 **Supplementary Figure 4.** The relative abundance of 16S rRNA throughout the crAssphage-  
888 rich frozen standard inoculum initiated faecal fermentation. **(A)** The relative abundance of the  
889 major genera detected throughout the fermentation. *Bacteroides* (\*), the genus hypothesised  
890 to be associated with crAssphage, can be seen to decrease between time points 0 and 4 of the  
891 fermentation after which levels gradually begin to increase again. **(B)** The relative abundance



892 of total *Bacteroides* at each time point. (C) Abundances of individual *Bacteroides* species  
893 detected. *B. dorei* is found to be particularly abundant and seemingly inversely proportional  
894 to the detected crAssphage levels.

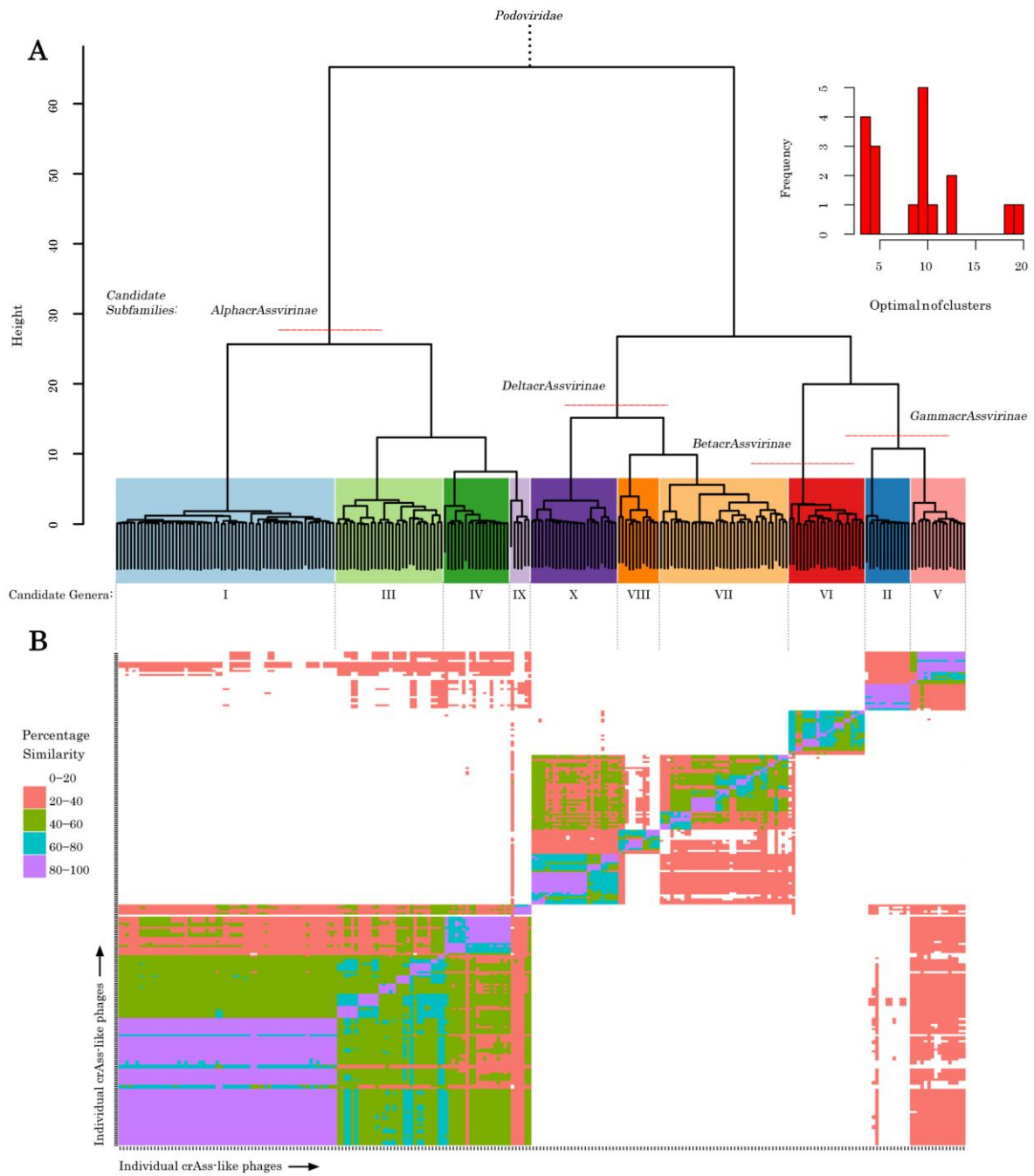
895 **Supplementary Figure 5.** Quantitative PCR analysis of filtrates obtained with different pore  
896 sizes from a crAssphage-rich faecal sample collected from subject ID 924.

897 **Supplementary Figure 6.** Comparison of 16S rRNA *Prevotella* abundances in healthy Irish  
898 adults and infants with Malawian infants.

899 **Supplementary Figure 7.** Visualisation of an example of metagenomic read-specific single  
900 nucleotide polymorphisms within the assembled of crAss-like phage contig, Fferm\_ms\_2,  
901 highlighting within sample species and/or strain level diversity of crAss-like phages are not  
902 resolved.

903 **Supplementary Figure 8.** Visualisation of the core proteome of the 10 crAss-like phage  
904 candidate genera.

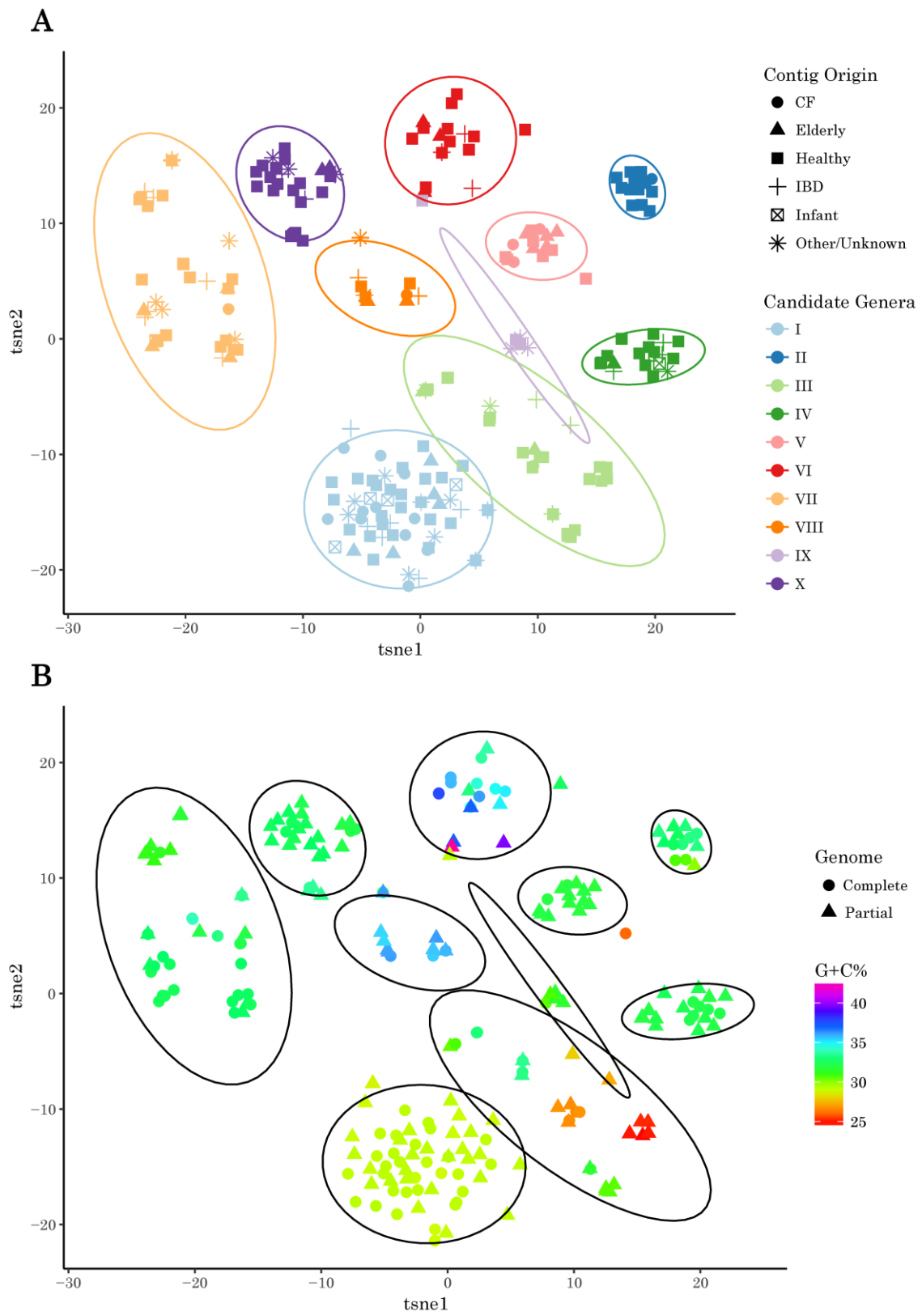
905 **Figure 1.**



906

907

908 **Figure 2.**

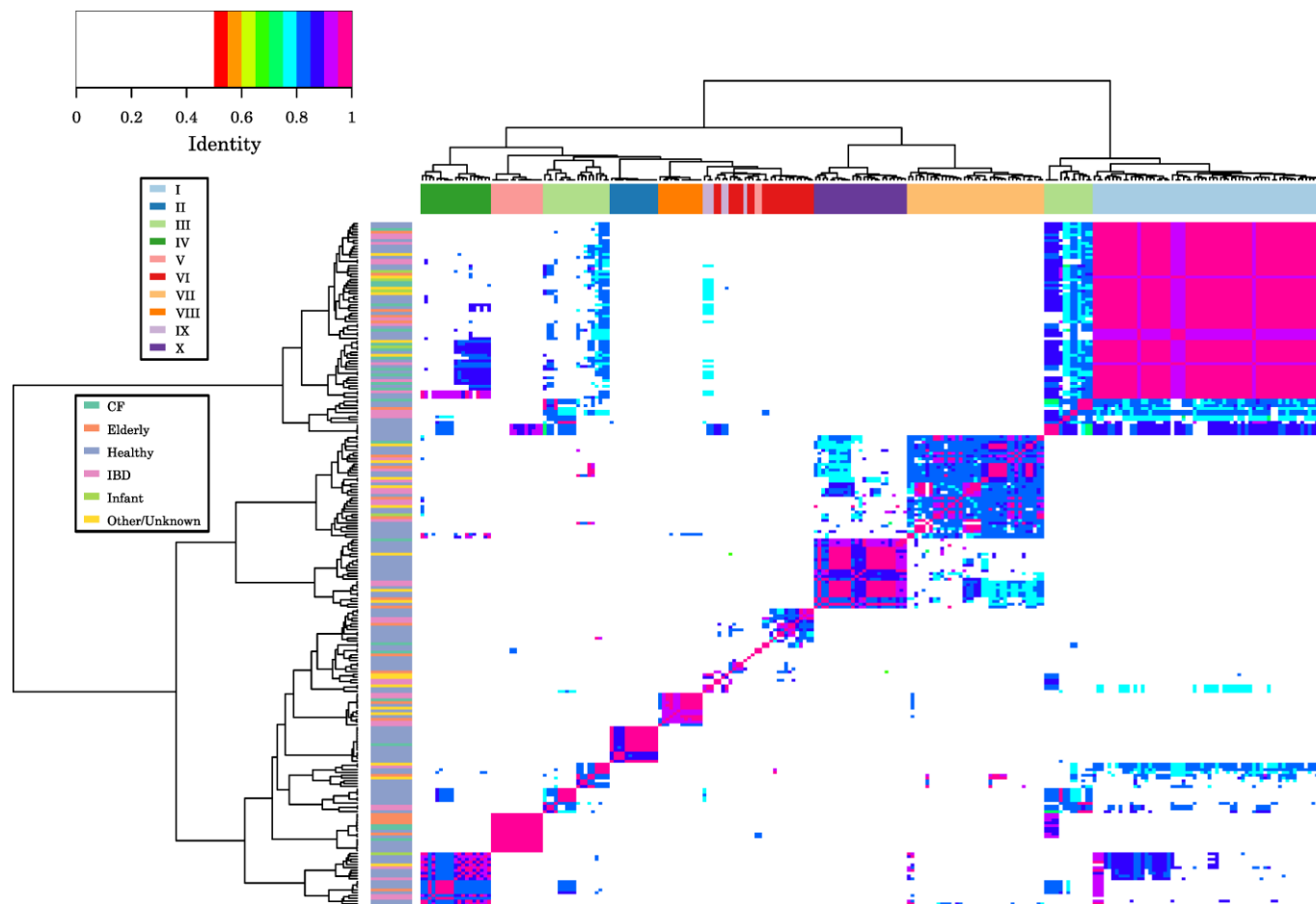


909

910

911

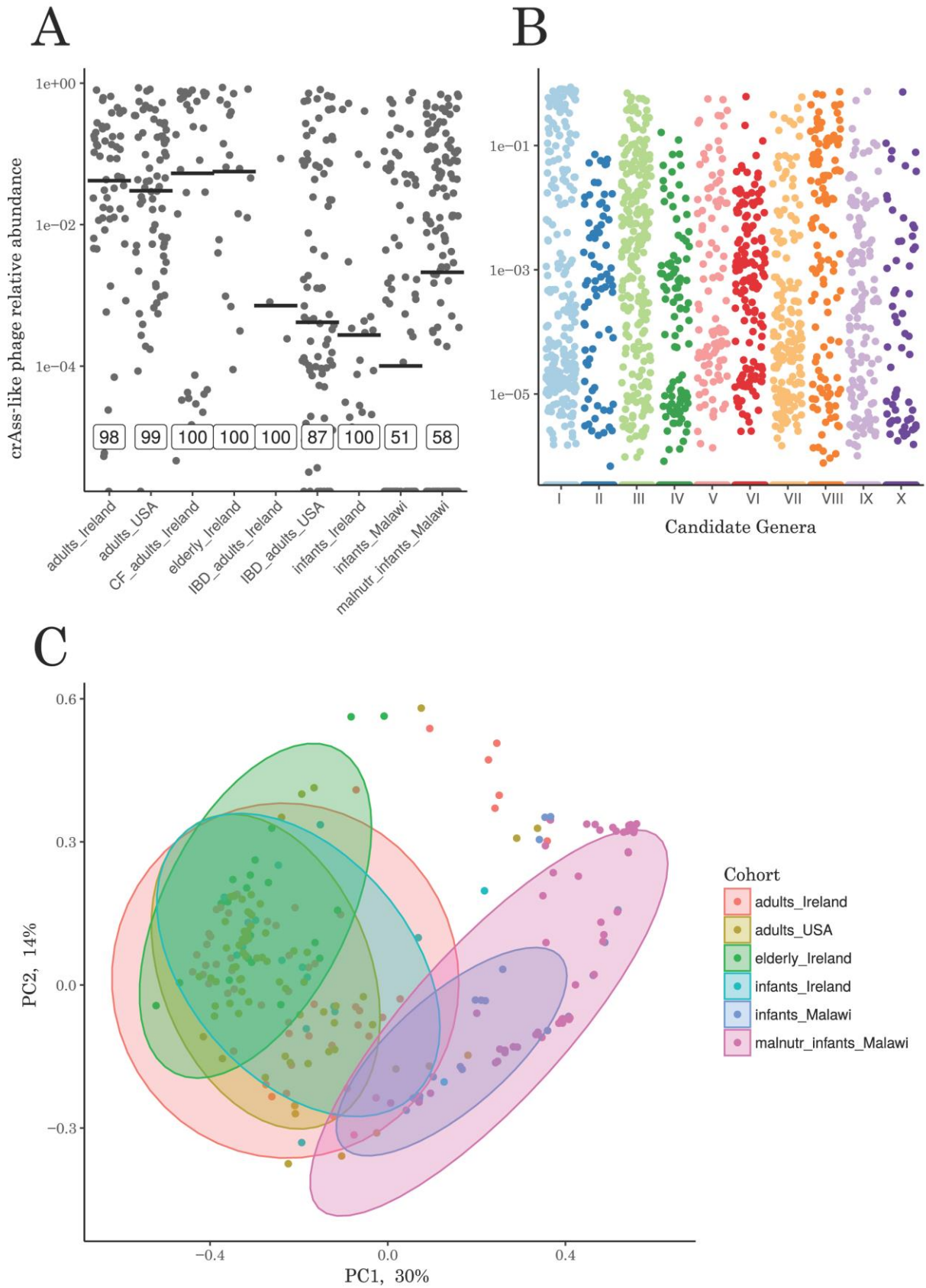
Figure 3.



912

913

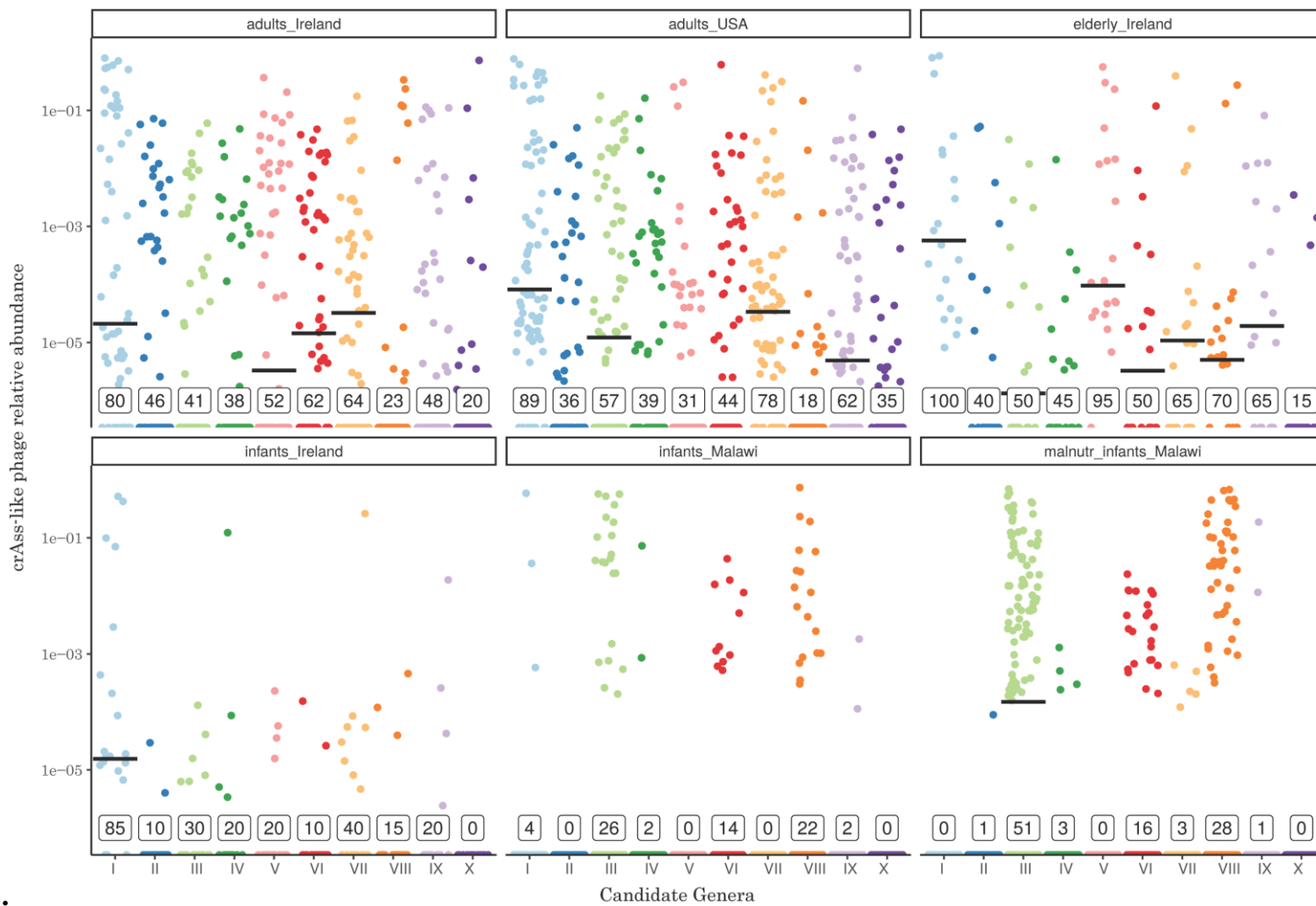
914 **Figure 4.**



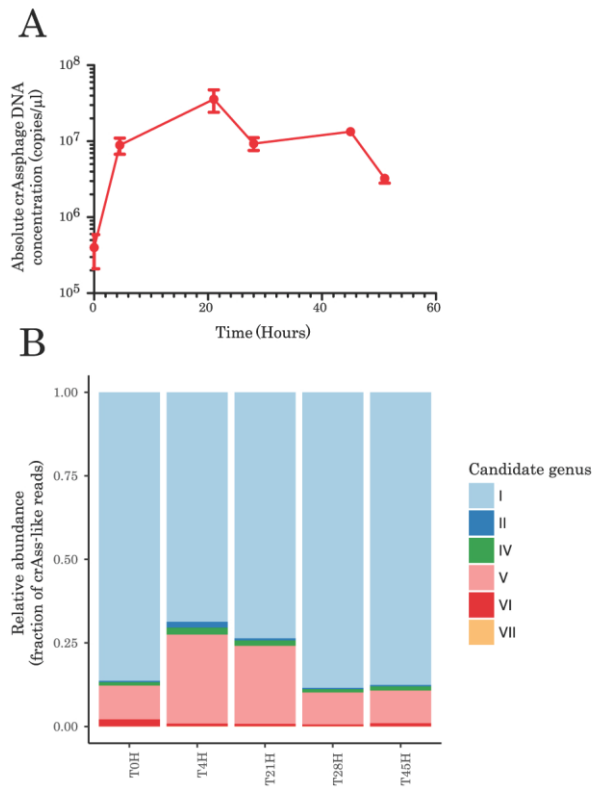
915

916

Figure 5



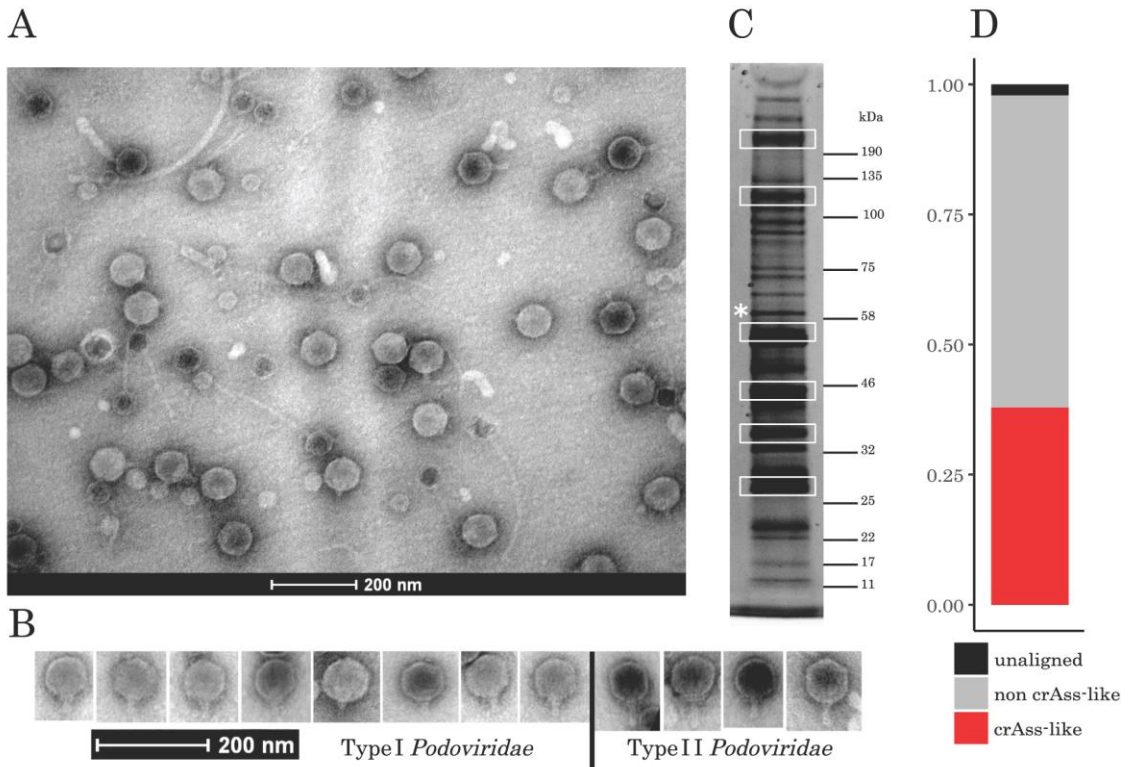
919 **Figure 6.**



920

921

922 **Figure 7.**

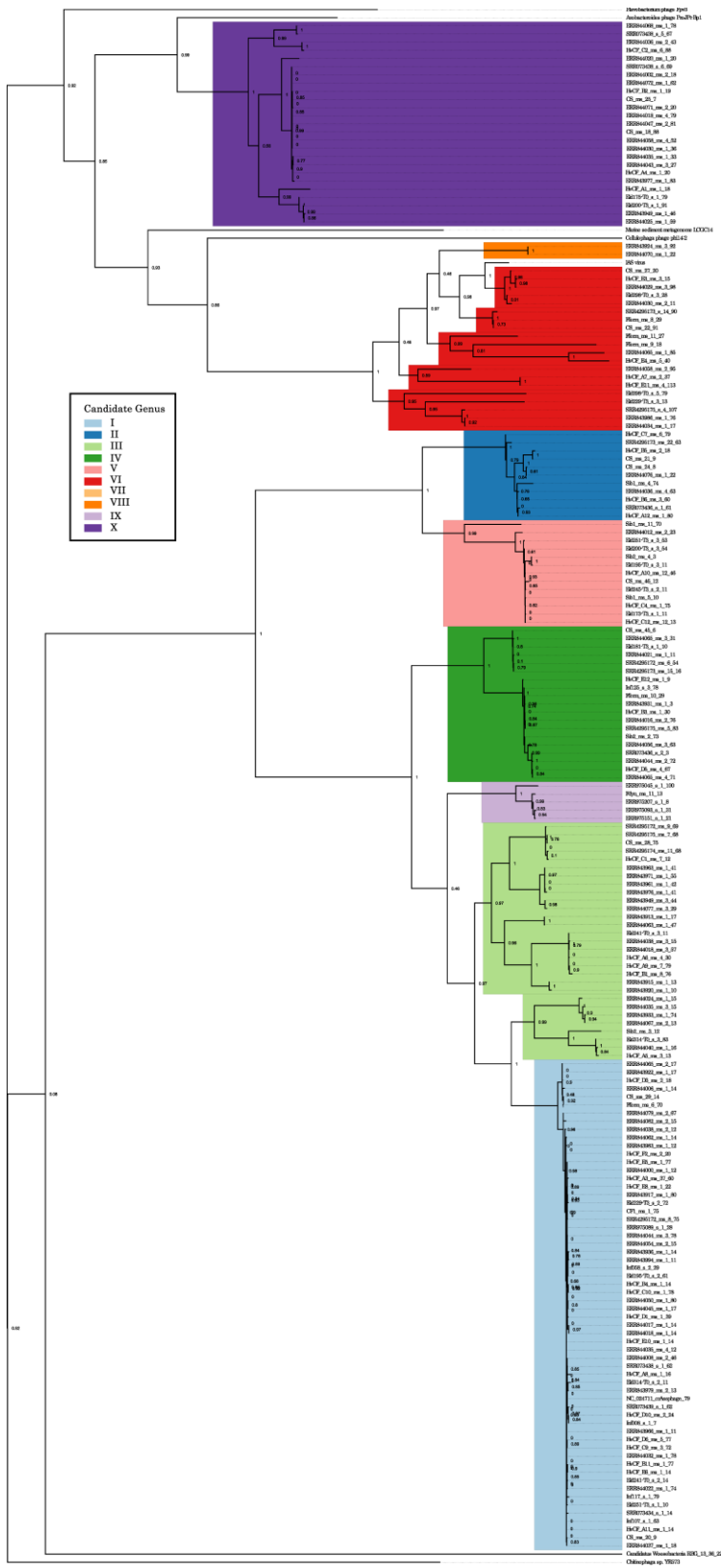


923

924



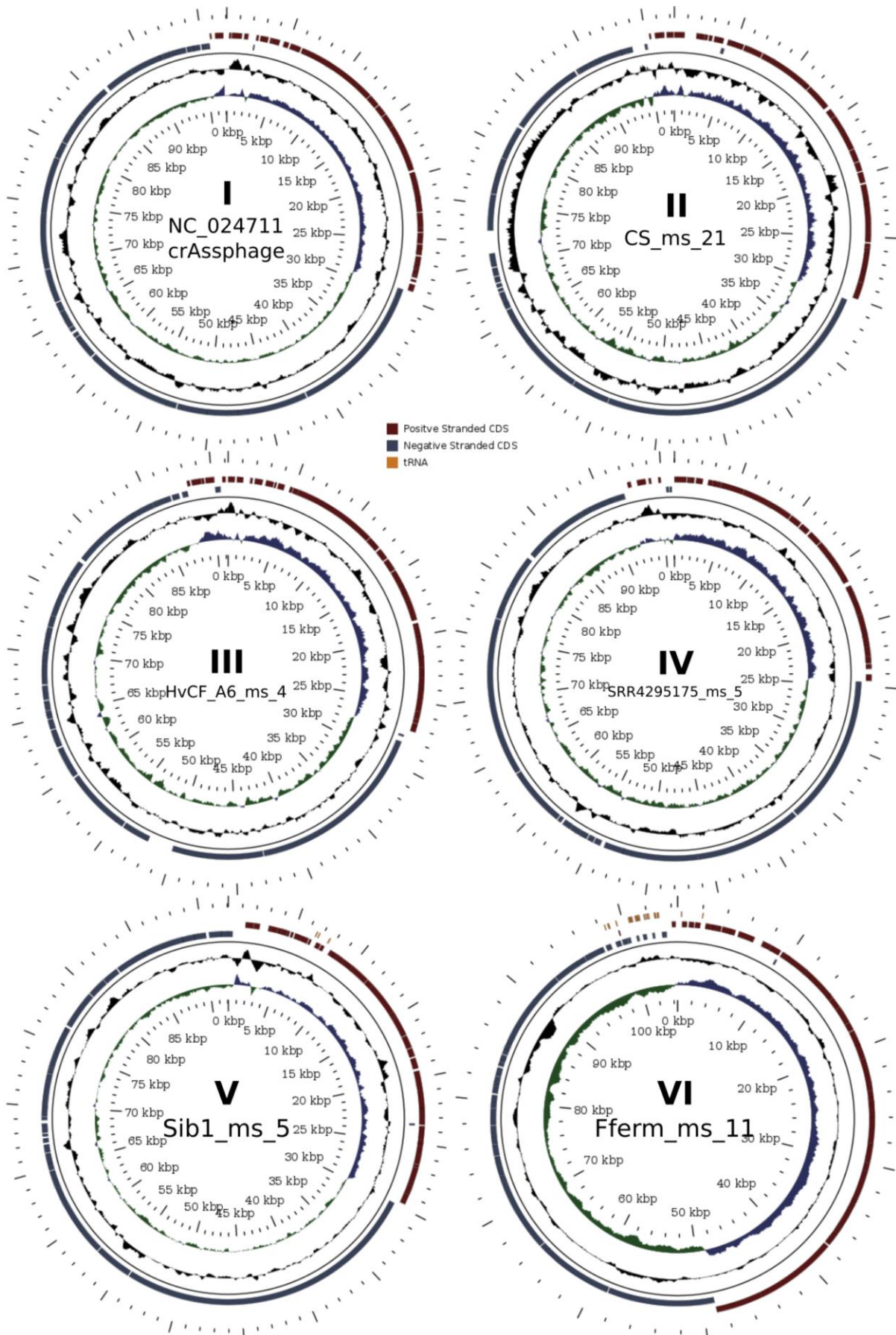
925 **Supplementary Figure 1.**



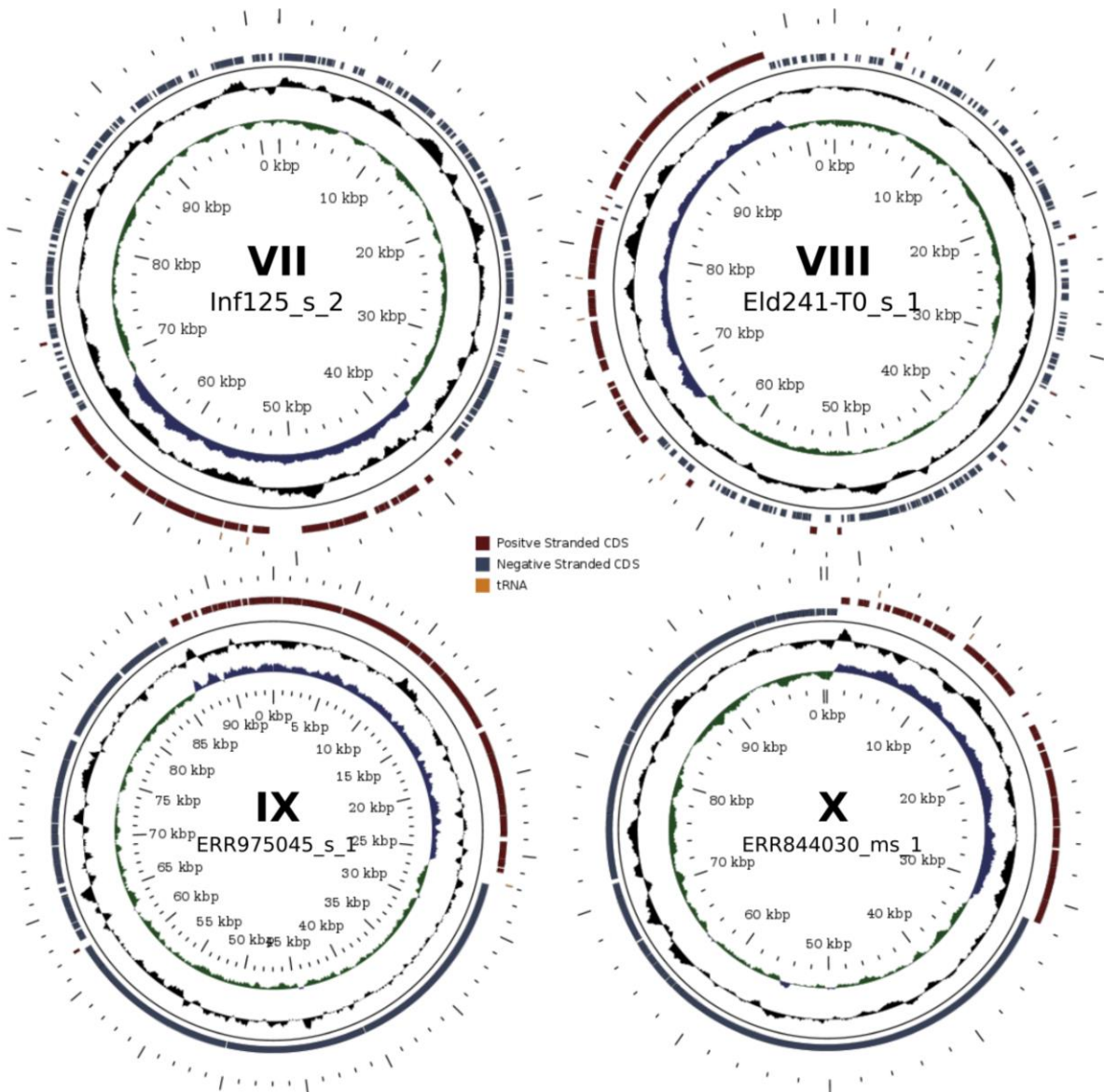
926

927

928 **Supplementary Figure 2.**



929

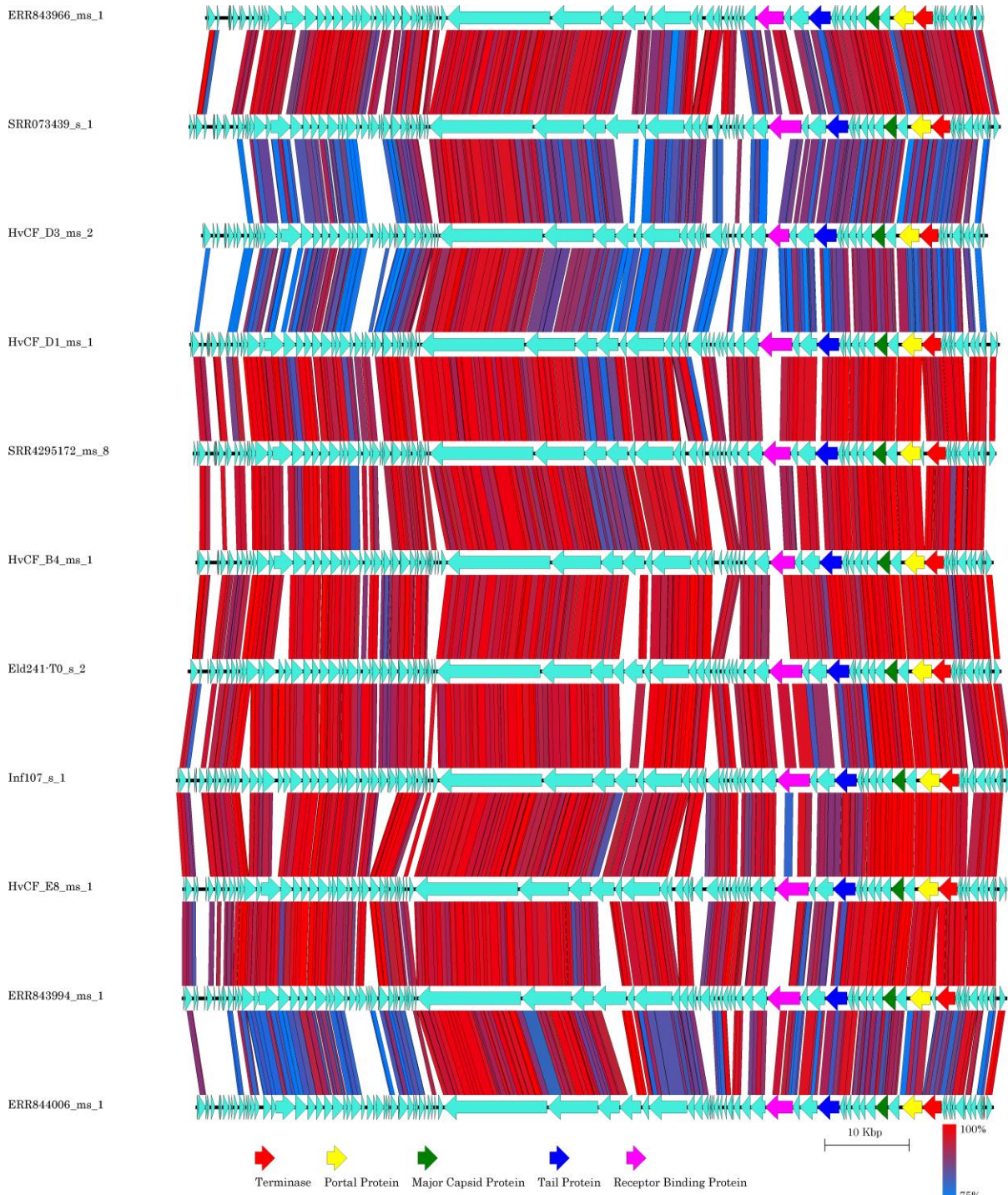


930

931

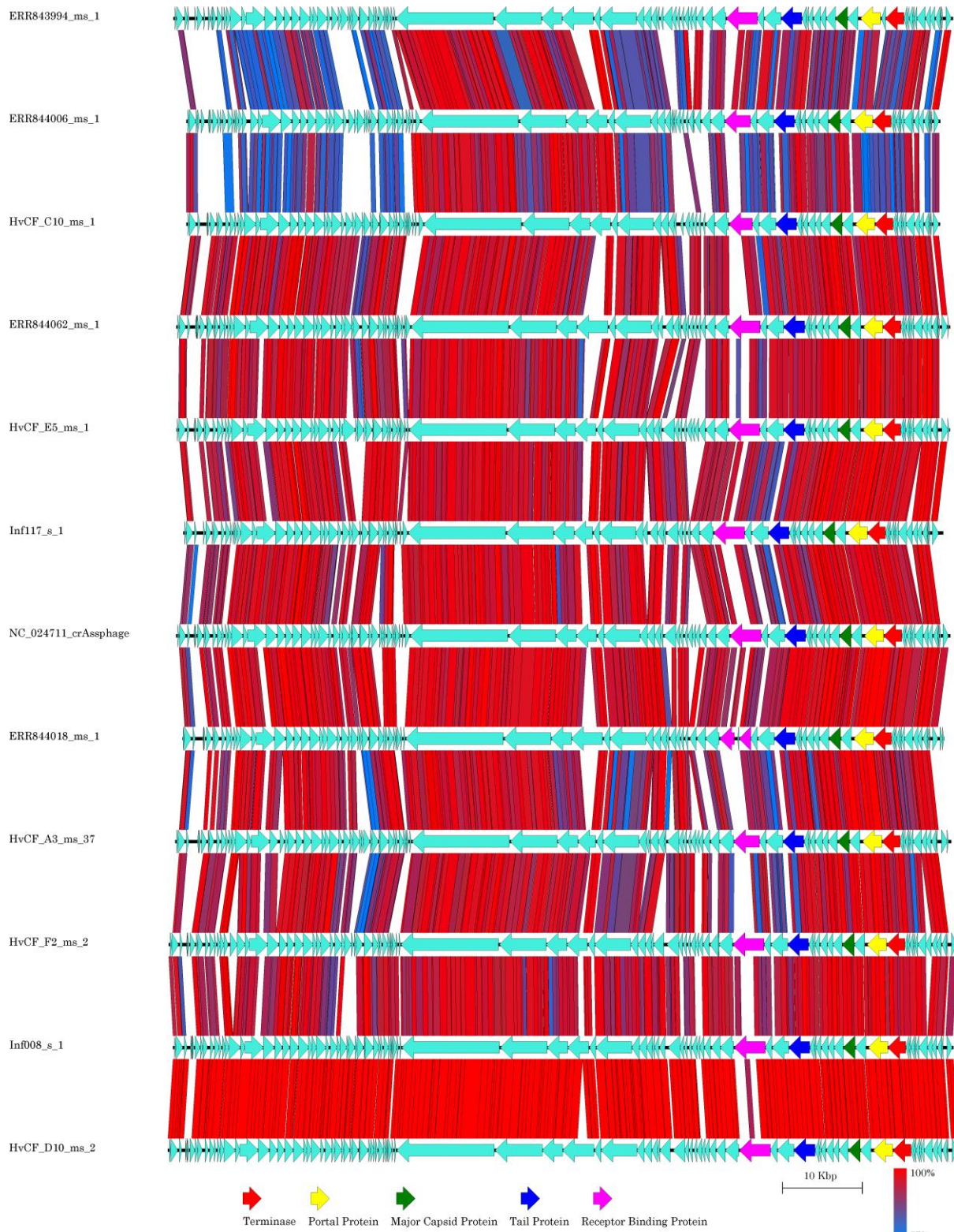


932 **Supplementary Figure 3.**



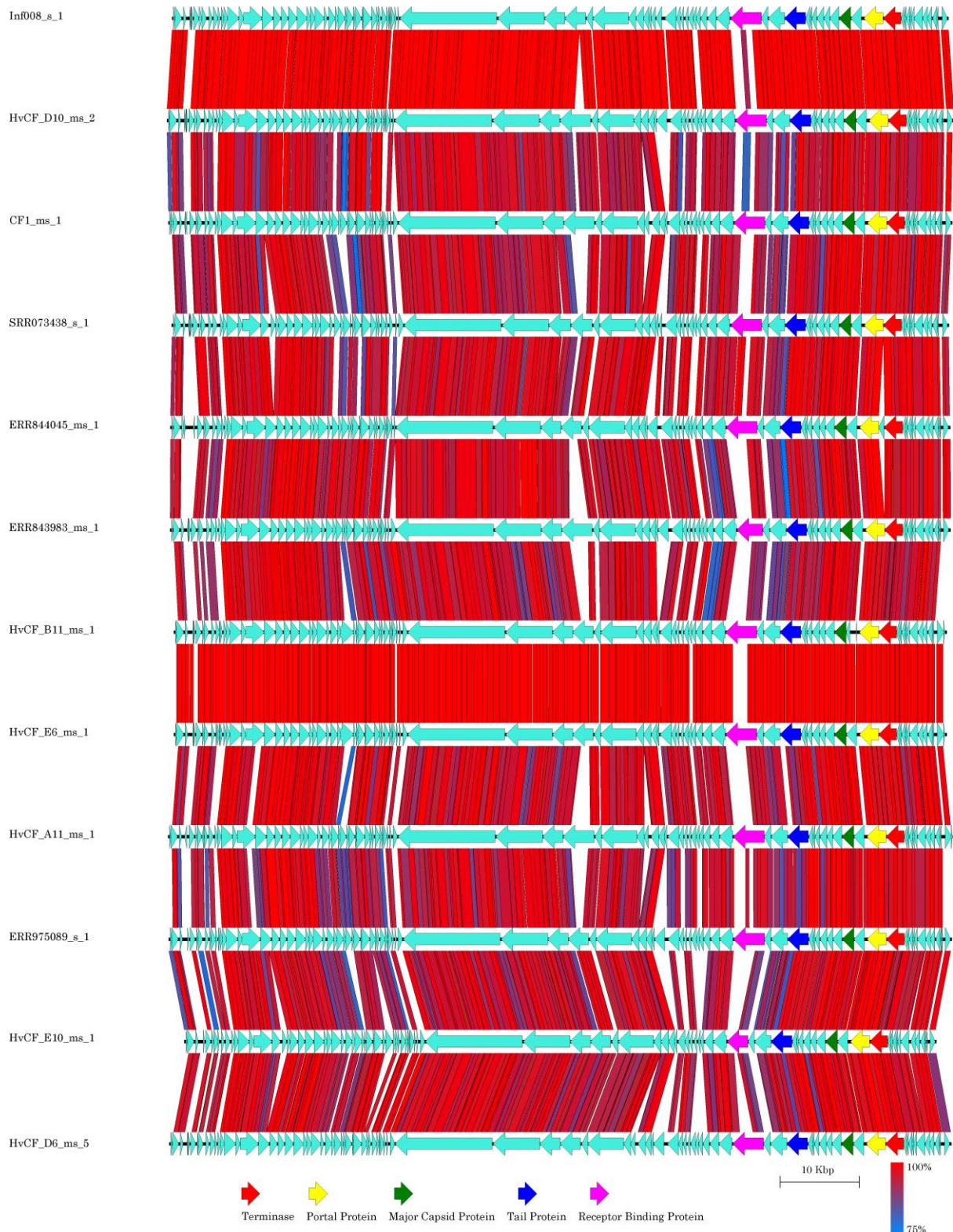
933



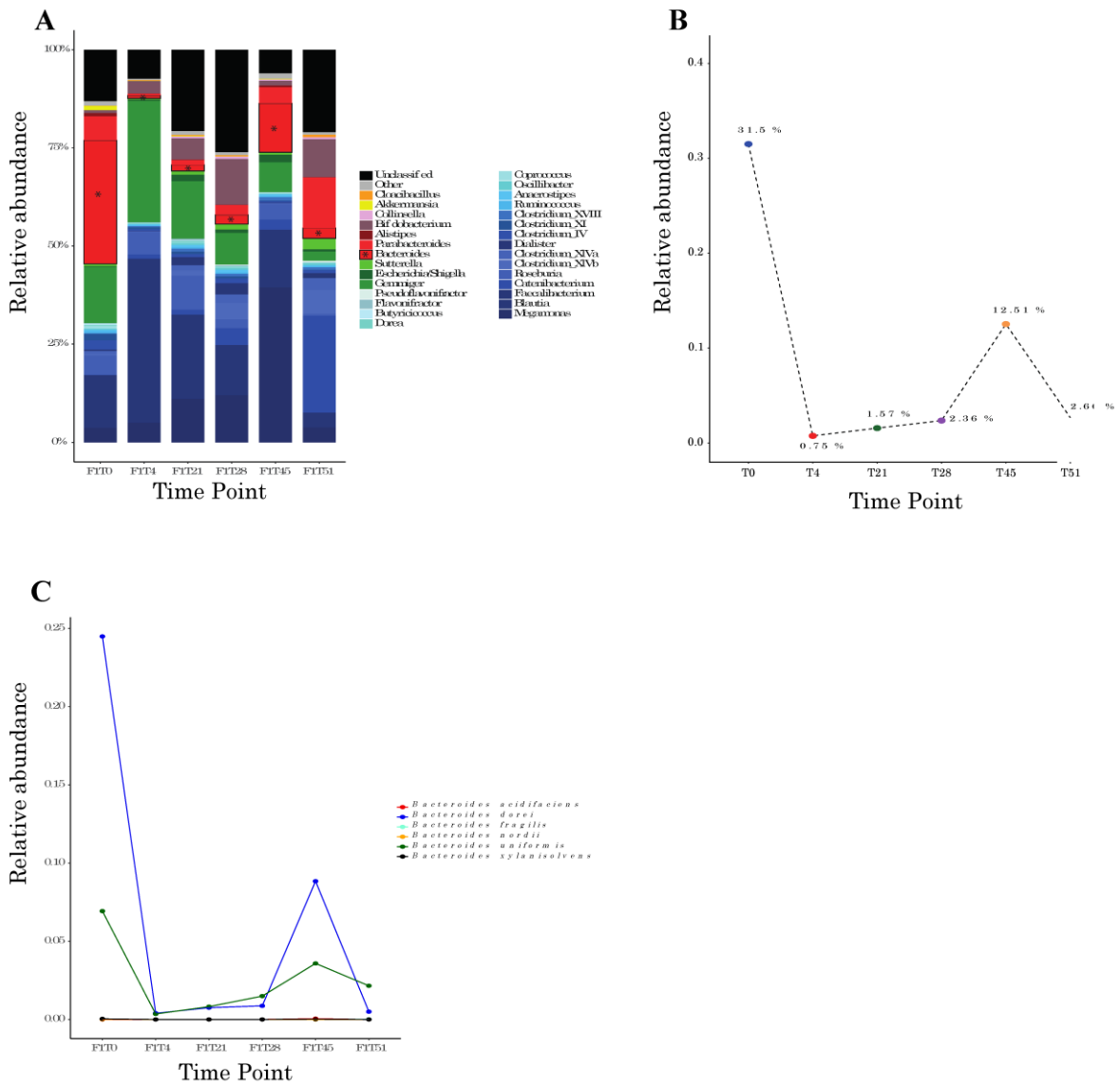


934





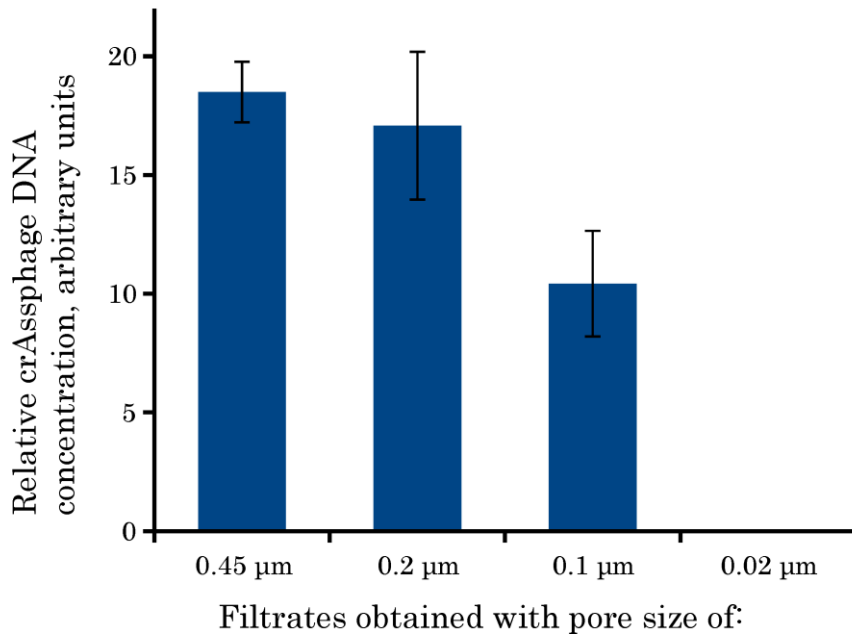
937 **Supplementary Figure 4.**



938

939

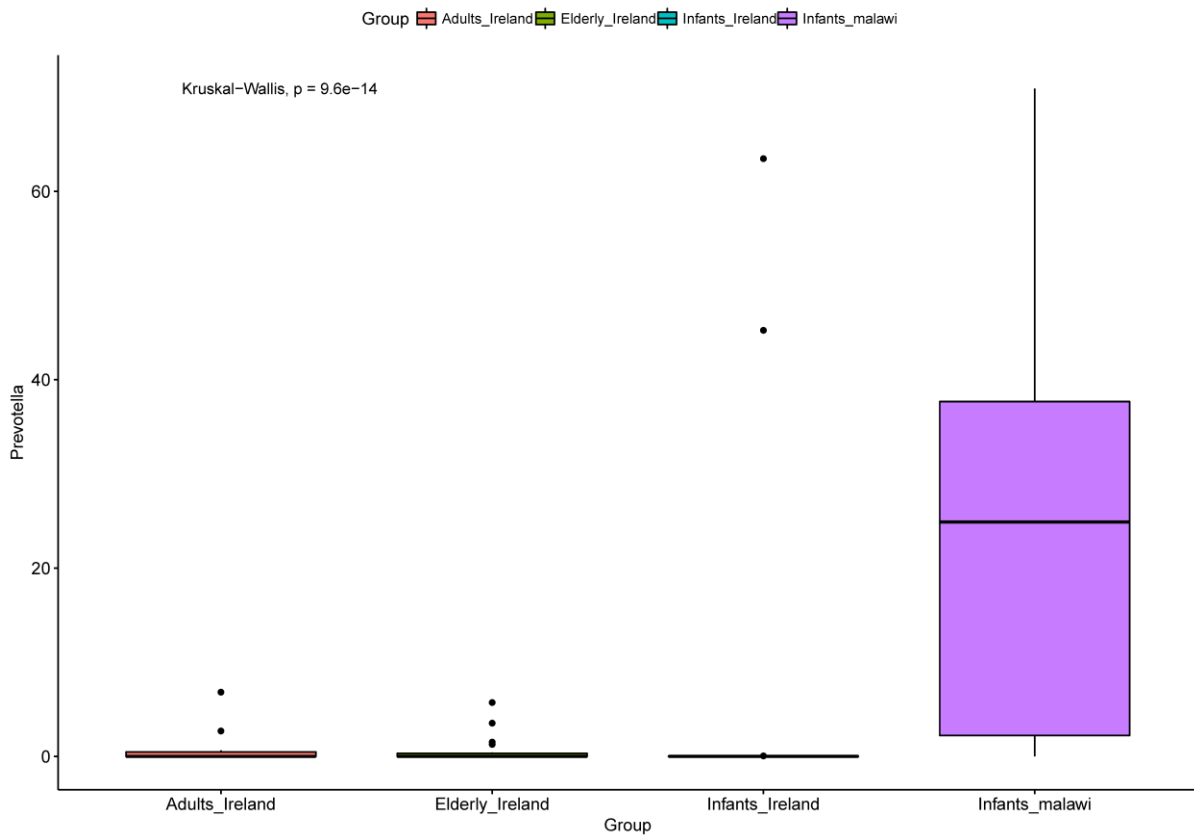
940 **Supplementary Figure 5.**



941

942

943 **Supplementary Figure 6.**



944

945



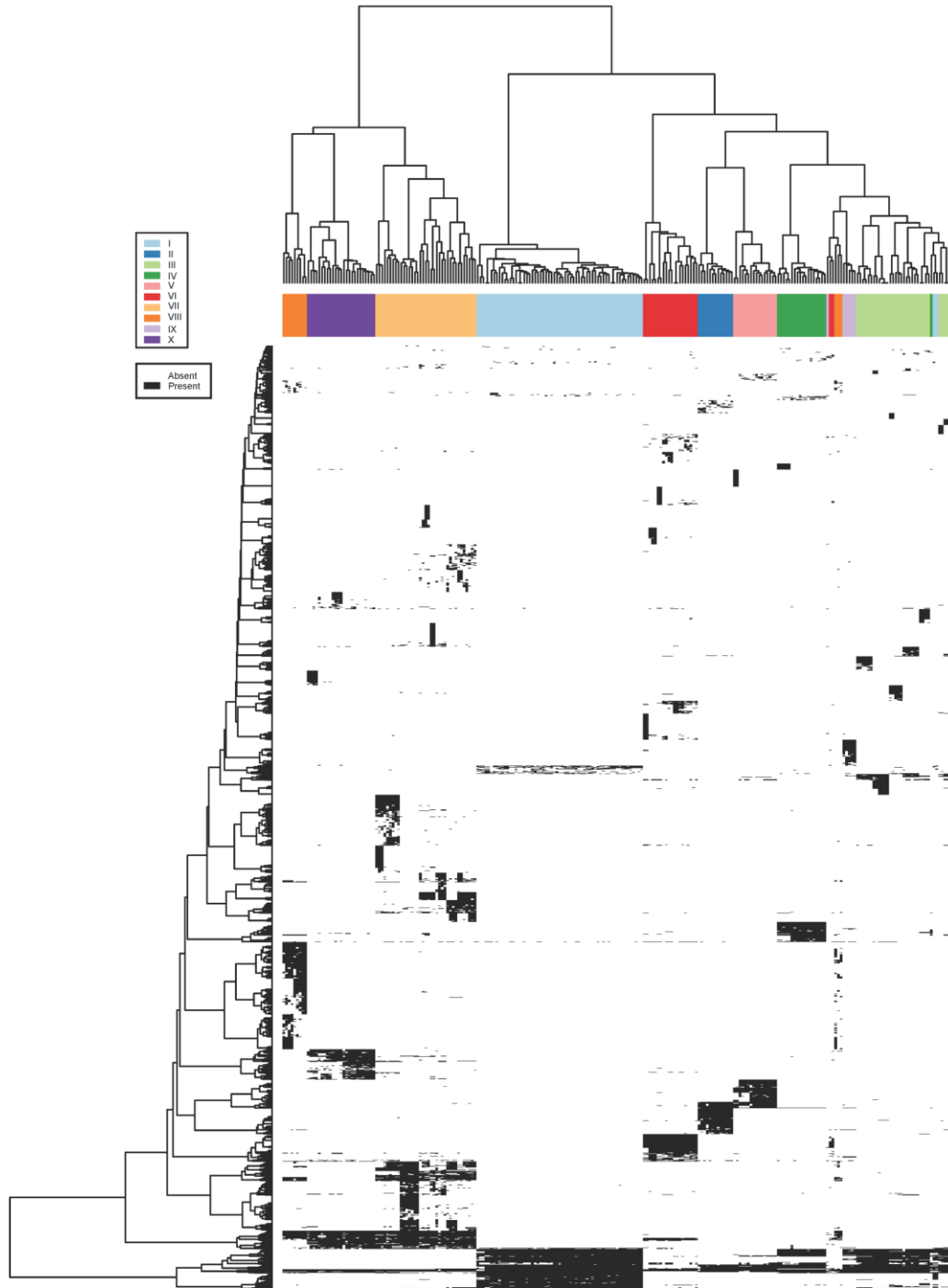
946 **Supplementary Figure 7.**



947

948

949 **Supplementary Figure 8.**



950