

1 **The ecogenomics of dsDNA bacteriophages in feces of stabled and feral horses.**

2

3 **V. V. Babenko<sup>1,\*</sup>, A. Millard<sup>2,\*</sup>, E. E. Kulikov<sup>3</sup>, N.N. Spasskaya<sup>4</sup>, M. A.**  
4 **Letarova<sup>3</sup>, D. N. Konanov<sup>1</sup>, I. Sh. Belalov<sup>3</sup> A.V. Letarov<sup>3,5,\*</sup>**

5

6 <sup>1</sup> FSC Physico-chemical medicine FMBA, Russia

7 <sup>2</sup> Dept Genetics and Genome Biology, University of Leicester, UK

8 <sup>3</sup> Winogradsky institute of microbiology RC Biotechnology RAS, Moscow, Russia

9 <sup>4</sup>Zoology museum, Faculty of biology, Lomonosov Moscow state university, Russia

10 <sup>5</sup>Faculty of biology Lomonosov Moscow state university, Russia

11

12 \* These authors contributed equally to this work

13 \* corresponding author, letarov@gmail. com

14

15

16 **Abstract**

17 The viromes of the mammalian lower gut were shown to be heavily dominated by  
18 bacteriophages; however, only for humans were the composition and intervariability  
19 of the bacteriophage communities studied in depth. Here we present an ecogenomics  
20 survey of dsDNA bacteriophage diversity in the feces of horses (*Equus caballus*),  
21 comparing two groups of stabled horses, and a further group of feral horses that were  
22 isolated on an island. Our results indicate that the dsDNA viromes of the horse feces  
23 feature higher richness than in human viromes, with more even distribution of  
24 genotypes. No over-represented phage genotypes, such as CrAssphage-related viruses  
25 found in humans, were identified. Additionally, many bacteriophage genus-level  
26 clusters were found to be present in all three geographically isolated populations. The  
27 diversity of the horse intestinal bacteriophages is severely undersampled, and so  
28 consequently only a minor fraction of the phage contigs could be linked with the  
29 bacteriophage genomes. Our study indicates that bacteriophage ecological parameters  
30 in the intestinal ecosystems in horses and humans differ significantly, leading them to  
31 shape their corresponding viromes in different ways. Therefore, the diversity and  
32 structure of the intestinal virome in different animal species needs to be  
33 experimentally studied.

34

35 **Short abstract (150 words) (needed in some journals as eLife)**

36 The viromes of the mammalian gut were shown to be heavily dominated by  
37 bacteriophages; however, only for humans were the composition and intervariability  
38 of the bacteriophage communities studied in depth. Here we present an ecogenomics  
39 survey of dsDNA bacteriophage diversity in the feces of horses (*Equus caballus*),  
40 comparing stabled horses, and feral horses that were isolated on an island. The  
41 viromes equine fecal viromes feature higher richness than in human viromes, with  
42 more even distribution of genotypes. No over-represented phage genotypes were  
43 identified. Additionally, many bacteriophage genus-level clusters were found to be  
44 present in geographically isolated populations. Only a minor fraction of the phage  
45 contigs could be linked with the bacteriophage genomes. Our study indicates that  
46 bacteriophage ecological parameters in the intestinal ecosystems in horses and  
47 humans differ significantly, leading them to shape their corresponding viromes in  
48 different ways.

49

50

51 **Importance. (needed for mBio)** The study presents the first in depth analysis of the  
52 composition and variability of the gut dsDNA bacteriophage community in the  
53 mammalian species, other than humans. The study demonstrates that the  
54 bacteriophage ecology in the gut is substantially different in different animal species.  
55 The results also indicate that the genetic diversity of the equine intestinal  
56 bacteriophages is immense and almost totally unexplored by the moment.

57

58 **Introduction**

59

60 The existence of microbial populations inhabiting different niches of the human and  
61 other animal bodies was first observed by Antony van Leeuwenhoek in 17<sup>th</sup> century  
62 (1) and has since become a commonly accepted paradigm that is mentioned in almost  
63 any microbiology textbook. Remarkable progress has been made in this field over the  
64 last 15 years due to the introduction of the culture-independent tools for the analysis  
65 of the composition and function of the microbial component of the human (2) or  
66 animal holobiont (3). Much emphasis has been placed on the gut microbiome,  
67 representing the largest microbial community associated with humans or other  
68 mammalian bodies. The intestinal microbiome is now considered as a “new organ”,

69 exerting strong and multifaceted influence over the physiology of the macro-host (4).  
70 The gut microbiome is involved in the pathology of numerous conditions, including  
71 Crohn disease (5, 6), obesity (7, 8), cancer (9) and even behavioral alterations (10-  
72 12).

73

74 It is well established that in all vertebrate animals the intestinal microbiome is  
75 associated with the corresponding virome - the community of the viruses infecting or  
76 produced by the microorganisms comprising the bacterial microbiome (13-16).  
77 Although the bulk of the intestinal viromes are comprised of bacteriophages (13, 14,  
78 17), these viral communities are also believed to be involved in multiple  
79 physiological effects and pathological processes *via* alteration of the composition and  
80 activity of the microbial community, and through direct interaction with the macro-  
81 host tissues and immune system (14, 18-20).

82

83 Bacteriophage diversity, biogeography and dynamics in the human gut has been  
84 investigated in depth in numerous studies using metagenomic approaches (16, 17, 21);  
85 see also reviews (17, 22, 23). It has to be mentioned, however, that in almost all of the  
86 studies the viral community of the feces was used as a proxy of the intestinal viromes.  
87 The “normal” composition of the human fecal bacteriophage community has been  
88 established and the “core phageome” composition defined as bacteriophage genotypes  
89 present in more than 50% individuals worldwide was evaluated (24). The first  
90 identified core phage lineage, named CrAssphage, that is highly prevalent in some of  
91 the samples (up to 90% of the viral reads) was initially identified using bioinformatic  
92 approaches (25) and was later cultured and shown to be a large podovirus infecting  
93 *Bacteroides* (26).

94 Therefore, the main characteristics of healthy human viromes have been established  
95 as follows: a diverse community, that is highly stable in time (17, 21), highly  
96 individual with larger inter-individual distances compared to different time points (16,  
97 17, 27) even if the dietary interventions were applied (21). Human viromes are  
98 suggested to be dominated by temperate bacteriophages (27) although the prevalence  
99 of the contigs containing integrases or site-specific recombinases genes is found to  
100 vary greatly (0-68%) between individual viral metagenomes (17).

101 Despite significant progress in the understanding of bacteriophage ecology in the  
102 human gut, the data on other animal species are scarce. Although a significant number

103 of metagenomic datasets from various species have been published (28-31), the vast  
104 majority of these studies focus on detection and interpretation of the animal viruses  
105 sequences, and bacteriophages have not been given significant attention. Only a few  
106 studies give emphasis to bacteriophage diversity in these samples, although this is  
107 limited to identifying differences between health and disease states in rhesus monkeys  
108 (32, 33) or specifically focusing on the diversity of ssDNA viruses (34).

109

110 In the present work we focus on the ecogenomics of dsDNA viruses present in horse  
111 feces. The equine intestinal microbiome plays an essential role in animal nutrition,  
112 allowing the horse to digest cellulose which is the major component of the grass  
113 consumed (35). In contrast to ruminants where microbial cellulose digestion takes  
114 place in the forestomach (rumen), in horses the cellulolytic microbial community  
115 develops in the cecum and large intestine that have a cumulative volume of about 100  
116 L with food retention time about 48-72 hours (36). The large intestine content is not  
117 subjected to any subsequent digestion (such as in ruminants) and is pushed by  
118 peristalsis into the rectum where it is subjected to partial dehydration to form the  
119 feces (37). The average time intervals between food intake or between the defecation  
120 acts in horses are much shorter than the indicated retention time (38). Therefore, the  
121 horse large intestine functions as a natural chemostat with highly stable physical and  
122 chemical conditions and fairly constant flow through.  
123 Adult horses do not show any coprophagy, but at the same time they do not avoid  
124 contact with feces of other individuals or fecally contaminated objects (38),  
125 potentially enhancing the exchange of bacteria and viruses between the individual  
126 viromes.

127 Only a few studies have been dedicated to horse intestinal bacteriophages. A limited  
128 Sanger-based metagenomic analysis of a single sample allowed the estimation of  
129 richness of the viral community, finding 1200 bacteriophage genotypes (39). A more  
130 recent metagenomic study compared fecal microbiomes and viromes of cattle and  
131 horses held on the same farm (40). However, the amount of data for each virome in  
132 this study was limited and no information concerning the specific characteristics of  
133 the bacteriophage communities of the samples was reported.  
134 There have also been several studies of the horse gut bacteriophage community based  
135 on other approaches. In a limited electron-microscopy study of horse feces almost all  
136 VLPs identified were classified as tailed phages, with 69 morphologically distinct



137 types reported out of <200 particles measured, indicating a high level of diversity  
138 (see (13) for review of earlier work). A comprehensive study of coliphage diversity  
139 and dynamics (41) in the feces of four horses held in the same location suggested high  
140 prevalence of virulent coliphages.. The *E. coli* host population was found to be highly  
141 divergent, and represented by hundreds of strains simultaneously present in the same  
142 sample. The overlap of the sensitivity of these strains to co-occurring bacteriophages  
143 was limited (41, 42) with ~1-5 % of the total *E. coli* counts being suitable hosts for  
144 any particular phage isolate. The data of Golomidova et al. (41, 43) and the results of  
145 the longitudinal study of G7C-related bacteriophages persistence and evolution within  
146 the ecosystem of a horse stable (44) indicated the flow of the coliphage genotypes  
147 between the animals. However, *E. coli* and its phages make for only a minor fraction  
148 of the total quine microbiome and viromes respectively. Currently, to the best of our  
149 knowledge, no methods exist that allow the translation of the findings made using  
150 this model system to the total community.

151 Here we present the ecogenomics of horse fecal dsDNA viromes of three separate  
152 horse populations including two groups of stabled horses and one herd of feral  
153 horses isolated on an island. Our data indicate that equine intestinal viromes are  
154 highly diverse communities dominated by the tailed bacteriophages. Although the site  
155 of sampling or/and the life conditions of distinct populations have marked influence  
156 over the composition of the individual viromes, it was possible to identify the equine  
157 intestinal core virome

158

## 159 **Results**

160 *The sampling strategy, workflow and sequencing results.*

161 In order to characterize the virome composition and diversity in horse feces we  
162 collected samples from three populations of horses in Russia. These included two  
163 groups of stabled animals and one feral population. The stable 1(S1) population was  
164 kept at the equestrian center in the city of Moscow. The stable 2 (S2) population was  
165 located in the country side ~90 km from Moscow. The lifestyle and diet of these two  
166 populations differed significantly (see material and methods). In addition, we sampled  
167 from feral horses at Rostovsky national reserve, that have been isolated on the island  
168 for several decades (see material and methods for detail) – population F. In population  
169 S1 we sampled four animals, two of which were sampled twice. In the S2 population  
170 five animals were sampled a single time (October 2018). In the F population, six

171 animals (two harem stallions and two mares belonging to each of the stallions) were  
172 sampled both in May and October 2018 (Table S1)

173 The viromes were extracted from all samples: viral DNA was extracted and  
174 sequenced using IonTorrent technology. It is important to note that the procedure  
175 applied for virome isolation (Fig.1) did not include any gradient centrifugation or  
176 ultrafiltration steps that may selectively remove some types of viral particles. We  
177 also did not use DNA amplification to avoid the biased representation of sequences  
178 that can occur. Based on previously published research (39, 45), see also (40) the bulk  
179 of the horse intestinal virome is composed of tailed bacteriophages, so we decided to  
180 focus on dsDNA viruses.

181 To check for bacterial contamination both virome QC and sortmeRNA were used, and  
182 both methods suggest the samples were highly enriched for viral DNA with minimal  
183 bacterial contamination. A total of 8097 nonredundant viral contigs >5 kb were  
184 identified, and were used for all further analysis. Among these contigs we identified  
185 46 contigs longer than 30 kb. that may represent complete or almost complete phage  
186 genomes.

187

#### 188 *Complexity of the individual viromes*

189 To estimate the alpha-diversity Shannon and Simpson indexes (Fig. 2) were calculated  
190 from relative abundance (how) and revealed that individual virome diversity tends to  
191 be higher in feral horses than in stabled (population S1). The population S2 falls in  
192 between, being closer to the feral populations. The samples ranking by Shannon and by  
193 Simpson indexes (TableS2) were almost identical (the maximum difference in a  
194 sample rank was 1). Shannon index is known to give more weight to species richness  
195 while Simpson index gives more emphasis on evenness (46). High correlation  
196 between these index values in the samples indicate that changes in richness are not  
197 associated with significant alterations of the evenness of the horse viromes, so all the  
198 samples contain high numbers of viral genotypes, none of which are significantly  
199 overrepresented.

200 To estimate the richness we used the approach mimicking that of Torsvik et al. (47)  
201 to estimate bacterial population richness in a soil sample. These authors estimated the  
202 complexity of the bacterial DNA extracted from soil using DNA re-association  
203 kinetics measurements. Knowing the average bacterial genome size, these authors  
204 calculated the approximate number of unique bacterial genotypes present. Instead of

205 experimental determination of the viral DNA re-association we computed the plots of  
206 the cumulative read recruitment against the cumulative length of the contigs ranked  
207 by the abundance (TPM) in the given sample (Fig. 2).

208 The samples appear to differ by abundance of the most prevalent viral genotypes.  
209 This is reflected by different slopes of the initial rise of the curves, though after this  
210 initial rise the curves are almost parallel, indicating the similar law of the genotypes  
211 abundance distribution in the viromes of different animals belonging to different  
212 population. Noteworthy, in none of the samples could we detect the presence of a  
213 over-represented genotypes. If we estimate an average phage genome as 50-100 kbp,  
214 the top represented 20-40 genotypes would account for 2 Mbp of the cumulative non-  
215 redundant DNA sequence. This value corresponds to 1.3 – 8.6% of the total amount  
216 of the viral DNA (Fig. 2). After this initial rise the curves are almost parallel that  
217 indicates the similar law of the genotypes abundance distribution in the viromes of  
218 different animals belonging to different population. We were not able to reveal the  
219 law of the distribution of the genotypes fractions within the communities study and,  
220 therefore, we did not find any reliable function to extrapolate the curves outside of the  
221 available data interval. However, to estimate the lower limit of the richness we used  
222 the function  $f(x) = ax^b \log(cx+1)$ , where  $x$  is the cumulative length of the contigs,  $f(x)$   
223 is the fraction of reads recruited by the most covered contigs with the cumulative  
224 length  $x$ , and  $a$ ,  $b$  and  $c$  are the parameters fitted to minimize the square deviation  
225 from the experimental curves. As shown in the Fig 2X, within the range the modeled  
226 curves run always higher than the experimental curves. Therefore, if the distribution  
227 law remains the same, the real  $x$  values corresponding to any  $(x)$  threshold chosen will  
228 be higher than the values predicted by the function. The calculated lowest estimates  
229 of the non-redundant length of the genomes sequences of the phage particles  
230 comprising 50 % of total community for most (20 out of 24) of the curves were in the  
231 range  $10^9 - 10^{11}$  b.p. This translates into  $10^4 - 10^6$  distinct bacteriophage genotypes  
232 without taking into consideration of possible overlap of the sequences in many  
233 different but still related viral genomes.

234

235

### 236 *Composition of horse intestinal virome*

237 Having established the high level of diversity of the equine gut dsDNA bacteriophage  
238 community, we asked how related are these numerous viral genotypes to known

239 viruses. First, we attempted direct classification of the filtered reads using  
240 centrifuge. However, only ~ 0.1% of reads can be classified this way. Out of them,  
241 97% matched viruses and 94% could be classified as dsDNA containing viruses, 81%  
242 of which were assigned to the order *Caudovirales* (tailed phages). Of these 43% were  
243 *Siphoviridae*, 41% *Myoviridae* and 13% *Podoviridae* (see Supplementary file Fig S1  
244 for the interactive Krona plot this analysis). However, given only a minor fraction of  
245 the reads could be classified, all the subsequent analysis was performed on the  
246 assembled contigs.

247

248 Initially we used pVOGs (48) to annotate all predicted proteins on viral contigs, with  
249 a simple scoring matrix. Out of 8097 contigs, 7483 (92%) had at least one pVOG  
250 detected. Only seven contigs where pVOGs were detected, were found to have a  
251 pVOG not found in the order *Caudovirales*. Further suggesting that the vast majority  
252 of contigs originated from tailed phages. We then analyzed the relatedness at the  
253 protein level using vCONTACT2, including RefSeq genomes plus other available  
254 phage genomes at the time (May 2019). The horse virome contigs were spread across  
255 1156 viral clusters (VCs), but only 31 were found in VCs that contain a known  
256 bacteriophage reference sequence, allowing classification at the genus level (Fig. 3,  
257 Tables S3). A further 2873 virome contigs remained singletons, once more  
258 highlighting the diversity of phages present.

259

260 Due to the inability to link the majority of contigs to any known phage at the  
261 subfamily or genus level, we manually inspected the 10 largest contigs that belonged  
262 to 10 different VC clusters. Gene products were analysed with both BLASTp and  
263 HHpred (49, 50) along with gene order and orientation in the genomes. We confirmed  
264 that even for the large (35-65 kbp) contigs the links to the known viral genomes were  
265 barely detectable and lie beyond the genus or subfamily level (Table S4, Fig.S3) at  
266 which vCONTACT2 is able to operate. Only in one case (the contig 070k255\_67966)  
267 were distant, but reliable relationships to the known *Gordonia* phage Gravy  
268 discovered, which was not in the vCONTACT database at the time of analysis. The  
269 results of the manual analysis further confirmed the viral diversity of the horse gut is  
270 to date very poorly sampled.

271 We also detected pVOGs that may be considered markers of the temperate life style  
272 (transposase, integrase, recombinase, resolvase and excisionases; see Material and

273 methods section for detail). At least one of such pVOGs was detected in 462 (5.7%)  
274 of the 5K contigs (Table S5). Among the contigs detected in the individual samples,  
275 the highest prevalence of the temperate lifestyle markers was observed in the  
276 population S1 (7%), in the populations F and S2 the prevalence of the “temperate”  
277 contigs was about 4%. Taking into consideration that the mean length of viral contigs  
278 was 8.3 kbp, the average number of the temperate lifestyle markers per genome is  
279 three, and estimating the average length of a temperate phage genome as 40-100 kbp,  
280 we can estimate the prevalence of the phage genotypes carrying these markers as 10-  
281 25%.

282

283 *Variability of the fecal viromes between the individuals and between the populations.*

284

285 To compare samples at the read level, we computed the Jacard’s distances between  
286 the datasets using mash (Fig.4). The samples clustered according to the populations  
287 they were collected from. Interestingly, the feral horses cluster tighter than the stable  
288 animals within the populations S1 and S2. At the same time, we did not observe any  
289 cluster formation according to the social (harem) groups of population F.  
290 Interestingly, in population F the samples collected at the different time points from  
291 the same animal were closer to each other than to any sample from the other animals.  
292 This remarkable stability of the individual viromes was observed despite the fact that  
293 between two sampling points the animals lived out a very hot summer that was  
294 associated with severe water deprivation because the debit of the water hole was  
295 decreased about two times for several months because the hole was blocked by the  
296 sand (it was cleaned by the rangers in late September). The clustering of the samples  
297 collected from the same animal at different time points was not observed in the  
298 population S1. However, in this population the period between sequential sampling  
299 was longer (Table S1). The population S2 appears to be much closer to the  
300 population F. This may reflect the fact that the diet of these two populations is much  
301 closer to each other than to the population S1.

302

303 To compare the individual samples and populations at the contig level we mapped  
304 reads from each sample against viral contigs. Contigs were considered present in a  
305 sample if the contig had  $\geq 1x$  coverage  $\geq 75\%$ , when mapping at 95% identity. Contig  
306 abundance

307 abundance was normalized, for both contig length and depth of sequencing. Thus we  
308 used “counts per thousand per million” (CPM) value as a proxy of contig abundance  
309 (52). An average of 913 contigs (range 655 – 1105, Table S3) were detected per  
310 sample.

311

312 The heatmap of the contig abundance in the samples is shown on the Fig 4. One can  
313 see that many contigs are shared by the animals belonging to the same group but  
314 much fewer are shared between the animals. The existence of the core-virome,  
315 defined as the assortment of the viral lineages present in the majority of the sequenced  
316 samples, was recently demonstrated for human feces (24). The crAssphage -like  
317 viruses that were shown to be highly abundant in some of the samples (24, 25, 51)  
318 also belong to the human feces core-virome. In order to reveal a possible horse core-  
319 virome we identified the contigs detected in all the samples or in more than half of the  
320 animals (the samples collected at different time points from the same animal were  
321 thus joined together). These criteria were applied for each population to reveal the  
322 local core-viromes, and for all the samples to retrieve the universal core-virome.

323 Venns diagrams of the local core viromes relatedness are shown on Fig. 5. Only 1  
324 contig was omnipresent in all the samples, however 192 contigs were shared by all  
325 three populations, among them 14 contigs were simultaneously present in 50% of the  
326 samples in each population (Fig. 5). Given the F population was completely isolated  
327 from S1 and S2 due to ca. 1500 km distance and protection by national reserve  
328 regimen (indirect exchange of viruses between the animal of the populations S1 and  
329 S2 is also highly unlikely though could not be completely excluded), these 192  
330 contigs can be considered as potential candidates for a “equine core viromes”. At the  
331 same time no contig exhibited abnormal coverage comparable to the values reported  
332 for human CrAssphage (25). The largest fraction of a single in the sum of all CPMs of  
333 all contigs of a sample was 0.006 (range 0.001 – 0.006). So, no equine analog of the  
334 over-represented and wide spread crAssphage group was detected in our dataset.

335 The distribution of VCs between the populations (Fig. 5B) revealed more of  
336 commonality between locations studied. Out of 1156 VCs, 262 (23%) VCs were  
337 detected in all three populations. This equates to 34-40% of VCs present in any of  
338 these populations (the VC was considered as present in a population if at least 1  
339 contig belonging to this VC passed the detection criteria for at least 1 of the samples  
340 from this population). Interestingly, in each of the populations many VCs were

341 present in 50% or more of the samples (Fig 5.C). The fractions of such prolific VC  
342 are larger in the populations F and S2 (89% and 90%) compared to S1 (23%). The  
343 fraction of contigs present in 50% or more of the samples in each of the populations  
344 are smaller (13%, 11% and 2% for the populations F, S2 and S1 respectively) with  
345 only 30 contigs present in 50% in all three populations simultaneously (Fig. 5D.).  
346 Thus, viromes appear to be highly individual at the level of the viral genotypes, but  
347 they appear to consist similar sets of bacteriophage genera.

348

349

350 *Host-phage relationships.*

351 High prevalence of the common VCs in three geographically isolated populations  
352 may reflect the presence of similar bacterial groups in the gut of horses belonging to  
353 different populations. To test this hypothesis we performed sequencing of bacterial  
354 16S rRNA genes libraries for all samples. We also predicted putative hosts for viral  
355 contigs using WiSH (53). The prevalence of the bacterial genus level OTUs and the  
356 prevalence of the contigs predicted to belong to bacteriophages infecting these host  
357 groups is shown in Fig 6.

358

359 The samples clustered according to 16S pattern still reflects the location (Fig 6.). At  
360 the same time the distribution of the prevalence of the phage genotypes predicted to  
361 infect different hosts was close to uniform. It has to be mentioned that the list of the  
362 genus-level bacterial OTUs and the list of bacterial genera predicted to be the hosts of  
363 the bacteriophage contigs only partially overlap. However, the non-overlapping OTUs  
364 have low prevalence as inferred from 16S sequencing data or from the statistics of the  
365 prevalence of the contigs allocated to the particular host.

366

## 367 **Discussion.**

368 Although the viral component of the intestinal microbiome is now widely believed to  
369 be an important factor in both shaping the microbial community of the gut and  
370 mediating its interactions with the macro-host (14, 18-20). The main component of  
371 the gut viromes – the community of tailed bacteriophages has only been  
372 comprehensively studied in humans. The results of our study provide a basic  
373 understanding of the composition of the dsDNA viromes of one more mammal  
374 species – *Equus caballus*. Given the very long history of domestication (54), until



375 very recently the integral involvement of domestic horses in almost all spheres of  
376 business and military activity, and the multifaceted influence of the relationship with  
377 this species on human culture (54). It would not be an exaggeration to say that this  
378 animal species is the second most important in the development of our civilization  
379 after *Homo sapiens*.

380

381 The equine intestinal bacterial community has been extensively studied, and it is now  
382 considered to influence horse organism homeostasis and health, to the same extent as  
383 bacteria in the human gut (reviewed in (35, 55), see also (56, 57). The equine gut  
384 bacterial community is involved in pathology of specific diseases such as equine  
385 metabolic syndrome and laminitis (58). The horse behavior was also suggested to be  
386 influenced by the gut bacterial community (59).

387 In contrast to the bacterial community, intestinal viromes, in particular, its main  
388 component – the dsDNA containing bacteriophages – were not investigated in any  
389 detail. Our data confirms that tailed bacteriophages (order *Caudovirales*) comprise the  
390 majority of total dsDNA viromes, so for brevity we use here below the term “virome”  
391 to describe the community of dsDNA containing bacteriophages.

392 An individual horse virome appears to be highly diverse including more than 2000  
393 viral genotypes (the extrapolation of the curves Fig. 2 gives estimates of more than  
394  $10^4 - 10^6$  viral genotypes per sample). The viromes richness did not differ much  
395 between the samples analyzed. In contrast to human viromes where the contig number  
396 per sample has been shown to vary more than three orders of magnitude (17), in our  
397 samples the variation was limited to a factor of less than 2.

398 At the same time the evenness of the viral genotype abundance was much higher in  
399 horses as could be seen by comparison of Shannon and Simpson diversity indexes and  
400 also inferred from the reads recruitment analysis (Fig. 2). No analog of human  
401 crAssphage that is hyper-dominant in some human gut samples, was observed. The  
402 most prolific 20-40 viral genotypes, between them only accounted for 1.3-8.6% of the  
403 total number of reads in all the samples. This makes a striking contrast to the situation  
404 described for human viromes where ~2% of the contigs that are so called persistent  
405 personal viromes recruited 92.3 % of VLP reads per sample (17).

406 Most of the bacteriophage genotypes comprising the bulk of equine intestinal dsDNA  
407 viromes are unrelated or very distantly related to known phage genotypes. Only 31  
408 out of 1152 identified VCs contained simultaneously horse virome contigs and known

409 bacteriophage genomes. Moreover, the manual analysis of the 10 largest contigs, did  
410 not allow (with a single exception) to assign these sequences to any known tailed  
411 phage group. Interestingly, we estimated that only 10-25% of this immense phage  
412 diversity are represented by temperate bacteriophages. These values are in marked  
413 contrast to the human fecal viromes that are dominated by the temperate  
414 bacteriophages (27, 60). Noteworthy, the high prevalence of virulent bacteriophages  
415 in horses is in agreement with previous data on the diversity of the coliphages isolated  
416 from horse feces (41) (13). The coliphages isolated from the human feces were  
417 reported to be mainly temperate (13, 60). At the same time, high richness and high  
418 evenness of the viral community as well as the lack of the correlation between the  
419 abundance of the host 16S-based OTUs and abundance of the predicted phages to  
420 these hosts (Fig. XX) indicate that the community most probably contains multiple  
421 viruses for many (if not for all) the bacterial species present in the samples. This  
422 pattern may support the elevated diversity at the strain level, might be maintained by a  
423 kill-the-winner type mechanism (61). Our metagenomic data does not provide any  
424 direct estimates of the bacterial diversity and/or phage-host relationships at the strain  
425 level. However high intraspecies diversity of *E. coli* within horse feces associated  
426 with high diversity of co-occurring coliphages, having relatively narrow host ranges,  
427 was previous demonstrated using the culture-based approaches (41, 42). Additional  
428 culture-based evaluation of the strain-level diversity of a more prevalent species than  
429 *E. coli* combined with characterization of its co-occurring phages, may help to shed  
430 more light over the pattern of the phage-host relationships in the horse gut ecosystem.  
431 Despite the observed high virome diversity, our data suggest that a healthy horse (in  
432 the feral population the animals without visible abnormalities, wounds and marked  
433 anomalies were considered as healthy) intestinal virome includes a certain number of  
434 conserved components. The human core-virome was defined by (24) as a set of viral  
435 genotypes that are present in more than 50% of the human fecal viromes. However,  
436 the limited amount of data (22 samples from 14 animals) makes this criterion less  
437 useful for evaluation of our data. At the same time, we may benefit from the known  
438 history of strict isolation of the population of the feral horses preserved on an island in  
439 Rostovski national reserve (population F) from any contacts with other horses. The  
440 factor of geographical isolation makes direct transfer of the viral genotypes even  
441 between the ancestors of these animals over last 80-100 years unlikely. Nevertheless,  
442 we were able to find significant number of the bacteriophage genotypes present in all

443 three populations. Approximately 3% (192 out of 6438) of contigs detected in the  
444 samples were universally present in all three locations. At the higher taxonomy level  
445 262 out of 1130 VCs (approximately corresponding to the genus or subfamily level of  
446 relatedness) detected in the individual samples were present in all the populations (Fig  
447 5). Moreover, despite the long history of isolation, populations S2 and F shared 552  
448 out of 875 VCs. Thus, a significant fraction of bacteriophage OTUs of species or  
449 genus level are widely present in the horse intestinal viromes, but the fractions of  
450 these common taxa may vary significantly. Higher similarity of the fecal viromes  
451 compositions of the populations F and S2 compared to their distance to the population  
452 S1 (Figs 4 and 5) may be explained by similar diets (grass only or grass and forages  
453 compared to grass and grain diets). Given the fact that the distribution of abundances  
454 of the viral genotypes in the individual viromes is very even (Fig 2 and Table S2 ),  
455 such variations may obscure the commonality of the viromes composition because  
456 many shared components are present below current limits of detection. At the same  
457 time, increasing the detection sensitivity using less stringent criteria may lead to a  
458 high frequency of false detection of the viral genotypes. The deep sequencing of  
459 several viromes from different location using, for example, high output Illumina  
460 sequencing and combined with long-reads single-molecule based sequencing (e.g.  
461 Oxford nanopore) may allow characterization of the repertory of the core components  
462 of equine virome. It is logical to expect that some endemic bacteriophage genotypes  
463 may also exist in certain populations, especially in the isolated animal groups, such as  
464 the population F. However, high diversity of the viromes does not allow the  
465 identification of them at the metagenome sequence coverage levels achieved in our  
466 work.

467 Another remarkable feature of the horse viromes is revealed by clustering of the  
468 individual viromes compositions according to the sampling site (Fig 4). The tightest  
469 cluster was formed by the samples from the population F. The clustering of these  
470 samples did not reflect the social structure of the herd. Only the samples taken from  
471 the same animals always clustered together.

472 Significant fractions of contigs and VCs were found in at least 50% of the samples in  
473 each of the populations, however the percentage of the shared contigs and VCs was  
474 lower in the group of the horses stabled in the city equestrian centre (S1). These  
475 findings may be explained by significant exchange by the viral genotypes between the  
476 animals. Horse in stable 1 (S1) are held in the individual boxes, which is much more

477 restrictive of behavior facilitating virus exchange through a fecal-oral route (see (38)).  
478 In the population from stable 2 (S2) during the summer season, horses spend most of  
479 their time at the pasture, and in the feral population (F) they have no human-imposed  
480 restrictions at all. During our field work we regularly observed behavior that may  
481 allow viral exchange (for example, during the spring time large groups of horses take  
482 the mud baths in the freshwater pools, where some animals may defecate and from  
483 which they also may drink), though detailed recording of the behavior falls out the  
484 scope of this work. In such conditions, the level of all-to-all exposure may erase the  
485 signal from tighter contacts within a harem group.

486

487 The phage genotypes transfer between the horses was earlier observed by the  
488 detection of the highly related coliphage isolates that could be obtained from multiple  
489 animals held in the same stable but could not be discovered in other locations (44, 62,  
490 63). These observations are in good agreement with the metagenomic data indicating  
491 that the transfer of the phages between the individual viromes is not limited to the  
492 minor viromes fraction(s) such as coliphages. So, the individuality and stability of the  
493 intestinal viromes are less pronounced in horses compared to humans (16, 17, 27).

494 Summarizing all the data, we conclude that horse intestinal viromes appear to be a  
495 more open ecological system than has been inferred from the human viromes. The  
496 bacteriophages of equine intestinal viromes represent a large pool of novel viral  
497 groups including the high level taxa such as families or subfamilies (we mean here  
498 new contemporary understanding of phage families, not the old *Siphoviridae* –  
499 *Myoviridae* – *Podoviridae* grouping within *Caudovirales*). This work provides an  
500 essential starting point from which the full genetic diversity for phages can be  
501 explored using long-read sequencing and culture based methods. Additional work is  
502 also required to analyze temporal stability of the horse viromes.

503

504

505

## 506 **Material and methods**

507

508 The horse populations and sampling

509 The sampling was performed in three geographically separated horse populations. The  
510 stable 1 population (S1) was located in a children's equestrian club in Detski park

511 Fili, Moscow, Russia, and represents typical stabled horses. These animals are kept in  
512 the boxes and taken outdoors for a limited time to be exercised (1-4 hours per day)  
513 and to have a rest (ca. 2 hours). These animals diet is typical for sportive horses and is  
514 comprised of foraging, supplemented with, oats, offal and carrots. The animals have  
515 an *ad libitum* access to water. The population from stable 2 (S2) was a group of  
516 horses living in a stable located in the village of Tretyakovo, Klinski district of  
517 Moskovskaya oblast, Russia. Horses are stabled in boxes and fed by forages (carrots  
518 or apples are occasionally given to them), but they spend 8-16 hours (depending on  
519 the season) per day in the field where they are able to graze. Access to water is not  
520 limited for this population, with water provided twice a day during dry weather,  
521 where the horses can drink *ad libitum*. In the population III four animals were  
522 sampled only once (in October 2018). Thus, the living conditions and diet of these  
523 three populations represent almost the whole spectrum of the conditions th  
524 The herd of feral horses (F population) inhabiting Vodny Island in the salty lake  
525 Manych-Gudilo in Rostovskaya oblast, Russia. This island belongs to the core part  
526 national Reserve “Rostovski” and therefore no business activity is allowed there, and  
527 the visits to the island are restricted. The horses on the Vodny island do not receive  
528 any feeding from humans and their diet includes only grass they forage. The access  
529 to drinking water is variable. From late autumn to spring the animals can drink from  
530 the pools or consume snow, the grass is also frequently wet due to rain and/or to dew-  
531 fall. In summer the watering is limited to water piped from the terrestrial beach once  
532 per day, and limited (ca. 3 L per min) output to an old water hole that exists in the  
533 middle of the island (64). The water from the holes is slightly saline. The access to  
534 these water sources differs significantly for different animals depending on their ranks  
535 in the social groups (harems) and on the rank of the harem stallion of their group  
536 among the harem stallions of the herd (at the moment of sampling in May and  
537 October 2017 there were 17-18 harem groups and a group of bachelor stallions).  
538 Samples of freshly voided feces were collected immediately after the natural  
539 defecation and placed in sterile plastic containers. The containers were placed on ice  
540 and transported to the laboratory. The samples from the populations S1 and S2 were  
541 processed within 24 h, the samples from the population F – within 72 h.

542

543 **Extraction of the viromes, DNA isolation and sequencing**

544 The viromes were extracted as it described in (65) with minor modifications. Briefly,  
545 the 10 g of the fecal sample was suspended in 100 ml of the extraction buffer (0.2M  
546 NaCl, 0.1 mg/ml NaN<sub>3</sub>, 1 mg/ml Tween 20 (Sigma-Aldrich, USA)) and extracted on  
547 the planetary shaker at 200 rpm, room temperature for 4 h. The coarse material was  
548 then filtered out using meltblown tissue (Miracle cloth meltblown fabric), then  
549 pelleted by centrifugation at 10000 g for 15 min. The supernatant was carefully  
550 separated. The samples were filtered through the combined filter composed of  
551 Whatman GF-F glass fiber paper and a layer of diatomite (Hyflo Super-Cel). The  
552 DNase was added to the filtrate up to 0.01 mg ml<sup>-1</sup> and the filtrate was incubated for  
553 1 h at room temp. The virus-like particles were then PEG-precipitated by adding dry  
554 NaCl to 0.6M and dry PEG 6000 (Panreac), dissolving both on orbital shaker (45  
555 min) and allowing precipitate to form in refrigerator (+4°C, 5-6 days). Brown-  
556 greenish precipitates containing VLPs were directly extracted with CTAB using a  
557 protocol described in (66).

558 Ion proton shotgun sequencing metavirome DNA (approx. 1000 ng) was fragmented  
559 to a mean size of 200-300 bp using the Covaris S220 System (Covaris, Woburn,  
560 Massachusetts, USA). Then, an Ion XpressPlus Fragment Library Kit (Life Tech-  
561 nologies) was employed to prepare a barcoded shotgunlibrary. Emulsion PCR was  
562 performed using the OneTouch system (Life Technologies). Beads were prepared  
563 using the One Touch 2 and Template Kit v2, and se-quencing was performed using  
564 Ion Proton 200 Sequen-cing Kit v2 and the P1 Ion chip. The reads were deposited to  
565 Sequence read archive (SRA) database, the accession numbers are given

566

### 567 **Bioinformatic analysis**

568 Prior to assembly reads were quality controlled by trimming with Sickle (67) with  
569 default settings. Contaminating of horse DNA was removed by mapping all reads to  
570 EquCab3.0 (GCA\_002863925.1) as reference genome, using bbmap with the  
571 following settings `minid=0.95` with any reads that mapped removed prior to  
572 assembly(68). Mapping suggested minimal contamination of horse DNA with the  
573 highest percentage of reads that mapped from any library at 0.07%. Metagenomes  
574 were assembled with MEGAHITv1.1.2, using the following parameters `k-min 21 --  
575 k-max 255 --k-step 10 -t 30`. Reads were mapped back against resultant contigs  
576 using BBmap `minid=0.95 covstats rpkm` (68). Resultant bam and sam files were



577 processed using Samtools v1.6 (69). To assess the level of bacterial DNA  
578 contamination all samples were processed with SortMeRNA v2.1 to check for  
579 contaminating rRNA reads `sortmerna --ref /usr/local/bioinf/sortmerna-`  
580 `2.1/rRNA_databases/silva-bac-16s-id90.fasta,/usr/local/bioinf/sortmerna-`  
581 `2.1/index/silva_b90:/usr/local/bioinf/s$` (70) and also using viromeQC. For read  
582 based assessment of viral diversity centrifuge was used with default settings and  
583 database of known phage genomes .`

584 For contig based assessment of viral diversity, contigs were first filtered with  
585 DeepVirFinder to remove any contigs that are of likely bacterial origin, only contigs  
586 with a p value <0.05 and were greater than 5 kb in length were considered for further  
587 analysis (71). Contigs were considered to be present within a sample if the average  
588 coverage of mapped reads was  $\geq 1X$  over  $\geq 70\%$  of the sample as recommended in  
589 other studies (4). The relative abundance of contigs within each sample was  
590 determined by counting the number of reads mapped to each contig, divided by the  
591 length of the contig (Kbp) to give RPK. The sum of all RPK values per sample was  
592 divided 1 000 000, with each RPK divided by a 1 000 000. Processing of data was  
593 carried out in R, using the PhyloSeq (72) library to calculate diversity statistics. To  
594 identify circular contigs `lastal -s 1 -x 300 -f 0 -T`` was used to identify the ends of  
595 contigs that overlapped (5).

596  
597 Contigs were annotated automatically using Prokka v1.12 using the following settings  
598 `--meta ``, using a custom database phage proteins (73). This database was  
599 constructed by extracting all the proteins from publically available phage genomes  
600 within the European Nucleotide Archive (5) and then further annotated using the  
601 scripts associated with prokka to do so []. Further annotation was provided by the use  
602 of hmmpfiles using hmmscan with the prokaryotic Viral Orthologous Groups  
603 (pVOG) collection of hmm profiles using a cutoff value of  $1E^{-5}$  (7, 8). To identify  
604 putative temperate phages a method akin to Sh & Hill was used. We utilised the a set  
605 of PFAM hmms (PF07508, PF00589, PF01609, PF03184, PF02914, PF01797,  
606 PF04986, PF00665, PF07825, PF00239, PF13009, PF16795, PF01526, PF03400,  
607 PF01610, PF03050, PF04693, PF07592, PF12762, PF13359, PF13586, PF13610,  
608 PF13612, PF13701, PF13737, PF13751, PF13808, PF13843, and PF13358) that are  
609 specific to bacteriophage transposase, integrase, recombinase, resolvase and  
610 excisionases.



611

612 Putative hosts were predicted using WIsH (53). A database of 9075 complete bacterial  
613 genomes was downloaded from Genbank (Jan 2018) and models were constructed for  
614 each bacterial genome within WIsH. Null parameters were calculated for each  
615 bacterial model using 7000 bacteriophage genomes. Hosts were predicted for each  
616 phage contig in the virome, with only predictions that had a pvalue of  $< 0.05$   
617 considered for further analysis.

618

619 Closest relatives

620 A custom database of all known phages genomes was produced by extraction of ~  
621 10,000 complete phage genomes from genbank as previously described . A MASH  
622 database was produced for using sketch  $-s$  10000. Each contig was queried against  
623 this database using mash dist function, with the top hit that had a distance of  $< 0.05$   
624 assigned as it closest known relative. To cluster contigs at the genus level, vContact2  
625 was used with the following settings “--rel-mode ‘Diamond’ --db  
626 'ProkaryoticViralRefSeq85-Merged' --pcs-mode MCL --vcs-mode ClusterONE”. The  
627 network graphs was visualized in Cytoscape and using Python package  
628 graphviz\_layout.

629

630 Diversity indices were produced by use the R Phyloseq package (72)

631

### 632 **Acknowledgements.**

633 We acknowledge T.Redgwell for help with proofreading and V.N. Filatova for help with  
634 the samples collection. The work was partially supported by RFBR grant #18-29-  
635 13029 (the field work and the pre-sequencing samples processing were performed  
636 before the grant acquisition). AM was funded by MRC grants MR/L015080/1 &  
637 MR/T030062/1. Bioinformatic analysis was in part carried out on infrastructure provided  
638 by MRC-CLIMB.

639

640

641

642

643

644

645 **References**

646

- 647 1. Fred EB. 1933. Antony van Leeuwenhoek: on the three-hundredth anniversary  
648 of his birth. *J Bacteriol* 25:iv 2-18.
- 649 2. Lloyd-Price J, Abu-Ali G, Huttenhower C. 2016. The healthy human  
650 microbiome. *Genome Med* 8:51.
- 651 3. Barko PC, McMichael MA, Swanson KS, Williams DA. 2018. The  
652 gastrointestinal microbiome: A review. *J Vet Intern Med* 32:9-25.
- 653 4. Heintz-Buschart A, Wilmes P. 2018. Human gut microbiome: function  
654 matters. *Trends Microbiol* 26:563-574.
- 655 5. Torres J, Mehandru S, Colombel J-F, Peyrin-Biroulet LJTL. 2017. Crohn's  
656 disease. 389:1741-1755.
- 657 6. Khanna S, Raffals LE. 2017. The microbiome in Crohn's disease: role in  
658 pathogenesis and role of microbiome replacement therapies. *Gastroenterol*  
659 *Clin North Am* 46:481-492.
- 660 7. Chen X, Devaraj S. 2018. Gut microbiome in obesity, metabolic syndrome,  
661 and diabetes. *Curr Diab Rep* 18:129.
- 662 8. Maruvada P, Leone V, Kaplan LM, Chang EB. 2017. The human microbiome  
663 and obesity: moving beyond associations. *Cell Host Microbe* 22:589-599.
- 664 9. Picardo SL, Coburn B, Hansen AR. 2019. The microbiome and cancer for  
665 clinicians. *Crit Rev Oncol Hematol* 141:1-12.
- 666 10. Nitschke A, Deonandan R, Konkle AT. 2020. The link between autism  
667 spectrum disorder and gut microbiota: A scoping review. *Autism*  
668 doi:10.1177/1362361320913364:1362361320913364.
- 669 11. Vuong HE, Yano JM, Fung TC, Hsiao EY. 2017. The microbiome and host  
670 behavior. *Annu Rev Neurosci* 40:21-49.
- 671 12. Warner BB. 2019. The contribution of the gut microbiome to  
672 neurodevelopment and neuropsychiatric disorders. *Pediatr Res* 85:216-224.
- 673 13. Letarov A, Kulikov E. 2009. The bacteriophages in human- and animal body-  
674 associated microbial communities. *J Appl Microbiol* 107:1-13.
- 675 14. Seo SU, Kweon MN. 2019. Virome-host interactions in intestinal health and  
676 disease. *Curr Opin Virol* 37:63-71.
- 677 15. Wang H, Ling Y, Shan T, Yang S, Xu H, Deng X, Delwart E, Zhang WJVe.  
678 2019. Gut virome of mammals and birds reveals high genetic diversity of the  
679 family Microviridae. 5:vez013.
- 680 16. Moreno-Gallego JL, Chou SP, Di Rienzi SC, Goodrich JK, Spector TD, Bell  
681 JT, Youngblut ND, Hewson I, Reyes A, Ley RE. 2019. Virome diversity  
682 correlates with intestinal microbiome diversity in adult monozygotic twins.  
683 *Cell Host Microbe* 25:261-272 e5.
- 684 17. Shkoporov AN, Hill C. 2019. Bacteriophages of the Human Gut: The "Known  
685 Unknown" of the Microbiome. *Cell Host Microbe* 25:195-209.
- 686 18. Neil JA, Cadwell KJTJoI. 2018. The intestinal virome and immunity.  
687 201:1615-1624.
- 688 19. Clooney AG, Sutton TDS, Shkoporov AN, Holohan RK, Daly KM, O'Regan  
689 O, Ryan FJ, Draper LA, Plevy SE, Ross RP, Hill C. 2019. Whole-virome  
690 analysis sheds light on viral dark matter in inflammatory bowel disease. *Cell*  
691 *Host Microbe* 26:764-778 e5.
- 692 20. Emler C, Ruffin M, Lamendella R. 2020. Enteric virome and carcinogenesis in  
693 the gut. *Dig Dis Sci* 65:852-864.

- 694 21. Minot S, Sinha R, Chen J, Li H, Keilbaugh SA, Wu GD, Lewis JD, Bushman  
695 FD. 2011. The human gut virome: inter-individual variation and dynamic  
696 response to diet. *Genome Res* 21:1616-25.
- 697 22. Koonin EV, Yutin N. 2020. The crass-like phage group: how metagenomics  
698 reshaped the human virome. *Trends Microbiol* 28:349-359.
- 699 23. Garmaeva S, Sinha T, Kurilshikov A, Fu J, Wijmenga C, Zhernakova A. 2019.  
700 Studying the gut virome in the metagenomic era: challenges and perspectives.  
701 *BMC Biol* 17:84.
- 702 24. Manrique P, Bolduc B, Walk ST, van der Oost J, de Vos WM, Young MJ.  
703 2016. Healthy human gut phageome. *Proc Natl Acad Sci U S A* 113:10400-5.
- 704 25. Dutilh BE, Cassman N, McNair K, Sanchez SE, Silva GG, Boling L, Barr JJ,  
705 Speth DR, Seguritan V, Aziz RK, Felts B, Dinsdale EA, Mokili JL, Edwards  
706 RA. 2014. A highly abundant bacteriophage discovered in the unknown  
707 sequences of human faecal metagenomes. *Nat Commun* 5:4498.
- 708 26. Shkoporov AN, Khokhlova EV, Fitzgerald CB, Stockdale SR, Draper LA,  
709 Ross RP, Hill C. 2018. PhiCrAss001 represents the most abundant  
710 bacteriophage family in the human gut and infects *Bacteroides intestinalis*. *Nat*  
711 *Commun* 9:4781.
- 712 27. Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, Gordon JL.  
713 2010. Viruses in the faecal microbiota of monozygotic twins and their  
714 mothers. *Nature* 466:334-8.
- 715 28. Duarte AM, Silva JM, Brito C, Teixeira D, Melo F, Ribeiro B, Nagata T, S.  
716 Campos F. 2019. Faecal Virome Analysis of Wild Animals from Brazil.  
717 *Viruses* 11.
- 718 29. Bergner LM, Orton RJ, Benavides JA, Becker DJ, Tello C, Biek R, Streicker  
719 DG. 2019. Demographic and environmental drivers of metagenomic viral  
720 diversity in vampire bats. *Mol Ecol* doi:10.1111/mec.15250.
- 721 30. Chen L, Gu W, Liu C, Wang W, Li N, Chen Y, Lu C, Sun X, Han Y, Kuang  
722 D, Tong P, Dai J. 2019. Characteristics of the tree shrew gut virome. *PLoS*  
723 *One* 14:e0212774.
- 724 31. Lima DA, Cibulski SP, Tochetto C, Varela APM, Finkler F, Teixeira TF,  
725 Loiko MR, Cerva C, Junqueira DM, Mayer FQ, Roehe PM. 2019. The  
726 intestinal virome of malabsorption syndrome-affected and unaffected broilers  
727 through shotgun metagenomics. *Virus Res* 261:9-20.
- 728 32. Li H, Li H, Wang J, Guo L, Fan H, Zheng H, Yang Z, Huang X, Chu M, Yang  
729 F, He Z, Li N, Yang J, Wu Q, Shi H, Liu L. 2019. The altered gut virome  
730 community in rhesus monkeys is correlated with the gut bacterial microbiome  
731 and associated metabolites. *Virology* 16:105.
- 732 33. Zhao G, Droit L, Gilbert MH, Schiro FR, Didier PJ, Si X, Paredes A, Handley  
733 SA, Virgin HW, Bohm RP, Wang D. 2019. Virome biogeography in the lower  
734 gastrointestinal tract of rhesus macaques with chronic diarrhea. *Virology*  
735 527:77-88.
- 736 34. Wang H, Ling Y, Shan T, Yang S, Xu H, Deng X, Delwart E, Zhang W. 2019.  
737 Gut virome of mammals and birds reveals high genetic diversity of the family  
738 Microviridae. *Virus Evol* 5:vez013.
- 739 35. Garber A, Hastie P, Murray JA. 2020. Factors influencing equine gut  
740 microbiota: Current knowledge. *J Equine Vet Sci* 88:102943.
- 741 36. Hintz H, Cymbaluk N. 1994. Nutrition of the horse. *Annual review of*  
742 *nutrition* 14:243-267.
- 743 37. Зеленевский НВ. 2007. Анатомия лошади. ИКЦ.

- 744 38. Waring GH. 2003. Horse behavior, 2nd ed. Noyes Publishing, Norwich, N.Y.
- 745 39. Cann AJ, Fandrich SE, Heaphy S. 2005. Analysis of the virus population  
746 present in equine faeces indicates the presence of hundreds of uncharacterized  
747 virus genomes. *Virus Genes* 30:151-6.
- 748 40. Park J, Kim EB. 2020. Differences in microbiome and virome between cattle  
749 and horses in the same farm. *Asian-Australas J Anim Sci* 33:1042-1055.
- 750 41. Golomidova A, Kulikov E, Isaeva A, Manykin A, Letarov A. 2007. The  
751 diversity of coliphages and coliforms in horse feces reveals a complex pattern  
752 of ecological interactions. *Appl Environ Microbiol* 73:5975-81.
- 753 42. Isaeva AS, Kulikov EE, Tarasyan KK, Letarov AV. 2010. A novel high-  
754 resolving method for genomic PCR-fingerprinting of Enterobacteria. *Acta*  
755 *Naturae* 2:82-8.
- 756 43. Golomidova AK, Kulikov EE, Prokhorov NS, Guerrero-Ferreira RC,  
757 Ksenzenko VN, Tarasyan KK, Letarov AV. 2015. Complete genome  
758 sequences of T5-related Escherichia coli bacteriophages DT57C and DT571/2  
759 isolated from horse feces. *Arch Virol* 160:3133-7.
- 760 44. Babenko VV, Golomidova AK, Ivanov PA, Letarova MA, Kulikov EE,  
761 Manolov AI, Prokhorov NS, Kostrukova ES, Matyushkina DM, Prilipov AG,  
762 Maslov S, Belalov IS, Clokie MRJC, Letarov AV. 2019. Phages associated  
763 with horses provide new insights into the dominance of lateral gene transfer in  
764 virulent bacteriophages evolution in natural systems. doi:10.1101/542787  
765 bioRxiv:542787.
- 766 45. Kulikov EE, Isaeva AS, Rotkina AS, Manykin AA, Letarov AV. 2007.  
767 Diversity and dynamics of bacteriophages in horse feces. *Mikrobiologiya*  
768 76:271-8.
- 769 46. Kim BR, Shin J, Guevarra R, Lee JH, Kim DW, Seol KH, Lee JH, Kim HB,  
770 Isaacson R. 2017. Deciphering Diversity Indices for a Better Understanding of  
771 Microbial Communities. *J Microbiol Biotechnol* 27:2089-2093.
- 772 47. Torsvik V, Goksoyr J, Daae FL. 1990. High diversity in DNA of soil bacteria.  
773 *Appl Environ Microbiol* 56:782-7.
- 774 48. Grazziotin AL, Koonin EV, Kristensen DM. 2017. Prokaryotic Virus  
775 Orthologous Groups (pVOGs): a resource for comparative genomics and  
776 protein family annotation. *Nucleic Acids Res* 45:D491-D498.
- 777 49. Johnson M, Zaretskaya I, Raytselis Y, Merezuk Y, McGinnis S, Madden TL.  
778 2008. NCBI BLAST: a better web interface. *Nucleic Acids Res* 36:W5-9.
- 779 50. Zimmermann L, Stephens A, Nam SZ, Rau D, Kubler J, Lozajic M, Gabler F,  
780 Soding J, Lupas AN, Alva V. 2018. A Completely Reimplemented MPI  
781 Bioinformatics Toolkit with a New HHpred Server at its Core. *J Mol Biol*  
782 430:2237-2243.
- 783 51. Edwards RA, Vega AA, Norman HM, Ohaeri M, Levi K, Dinsdale EA, Cinek  
784 O, Aziz RK, McNair K, Barr JJ, Bibby K, Brouns SJJ, Cazares A, de Jonge  
785 PA, Desnues C, Diaz Munoz SL, Fineran PC, Kurilshikov A, Lavigne R,  
786 Mazankova K, McCarthy DT, Nobrega FL, Reyes Munoz A, Tapia G,  
787 Trefault N, Tyakht AV, Vinuesa P, Wagemans J, Zhernakova A, Aarestrup  
788 FM, Ahmadov G, Alassaf A, Anton J, Asangba A, Billings EK, Cantu VA,  
789 Carlton JM, Cazares D, Cho GS, Condeff T, Cortes P, Cranfield M, Cuevas  
790 DA, De la Iglesia R, Decewicz P, Doane MP, Dominy NJ, Dziewit L, Elwasila  
791 BM, Eren AM, et al. 2019. Global phylogeography and ancient evolution of  
792 the widespread human gut virus crAssphage. *Nat Microbiol* 4:1727-1736.

- 793 52. StatQuest. July 22, 2015 2015. RPKM, FPKM and TPM, clearly explained.  
794 <https://www.rna-seqblog.com/rpkm-fpkm-and-tpm-clearly-explained/>.  
795 Accessed
- 796 53. Galiez C, Siebert M, Enault F, Vincent J, Soding J. 2017. WISH: who is the  
797 host? Predicting prokaryotic hosts from metagenomic phage contigs.  
798 *Bioinformatics* 33:3113-3114.
- 799 54. Kelekna P. 2009. The horse in human history. Cambridge University Press,  
800 Cambridge ; New York.
- 801 55. Murcia PR. 2019. Clinical insights: The equine microbiome. *Equine*  
802 *veterinary journal* 51:714.
- 803 56. Langner K, Blaue D, Schedlbauer C, Starzonek J, Julliand V, Vervuert I.  
804 2020. Changes in the faecal microbiota of horses and ponies during a two-year  
805 body weight gain programme. *PLoS One* 15:e0230015.
- 806 57. Lindenberg F, Krych L, Kot W, Fielden J, Frokiaer H, van Galen G, Nielsen  
807 DS, Hansen AK. 2019. Development of the equine gut microbiota. *Sci Rep*  
808 9:14427.
- 809 58. Patterson-Kane JC, Karikoski NP, McGowan CM. 2018. Paradigm shifts in  
810 understanding equine laminitis. *Vet J* 231:33-40.
- 811 59. Bulmer LS, Murray JA, Burns NM, Garber A, Wemelsfelder F, McEwan NR,  
812 Hastie PM. 2019. High-starch diets alter equine faecal microbiota and increase  
813 behavioural reactivity. *Sci Rep* 9:18621.
- 814 60. Mathieu A, Dion M, Deng L, Tremblay D, Moncaut E, Shah SA, Stokholm J,  
815 Krogfelt KA, Schjørring S, Bisgaard HJNC. 2020. Virulent coliphages in 1-  
816 year-old children fecal samples are fewer, but more infectious than temperate  
817 coliphages. 11:1-12.
- 818 61. Weinbauer MG. 2004. Ecology of prokaryotic viruses. *FEMS Microbiol Rev*  
819 28:127-81.
- 820 62. Golomidova AK, Kulikov EE, Babenko VV, Ivanov PA, Prokhorov NS,  
821 Letarov AV. 2019. Escherichia coli bacteriophage Gostya9, representing a  
822 new species within the genus T5virus. *Arch Virol* 164:879-884.
- 823 63. Golomidova AK, Kulikov EE, Prokhorov NS, Guerrero-Ferreira Rcapital Es  
824 C, Knirel YA, Kostryukova ES, Tarasyan KK, Letarov AV. 2016. Branched  
825 lateral tail fiber organization in T5-like bacteriophages DT57C and DT571/2  
826 is revealed by genetic and functional analysis. *Viruses* 8.
- 827 64. Minoransky VA, Uzdénov VM. 2011. Feral horses of Vodny island (Manych-  
828 Gudilo lake, Rostov region, Russia). *Vesti biosfernogo zapovednika "Askania-  
829 nova"* 13:135-145.
- 830 65. Kulikov EE, Golomidova AK, Prokhorov NS, Ivanov PA, Letarov AV. 2019.  
831 High-throughput LPS profiling as a tool for revealing of bacteriophage  
832 infection strategies. *Sci Rep* 9:2958.
- 833 66. Thomas JC, Houry R, Neeley CK, Akroush AM, Davies ECJBE. 1997. A  
834 fast CTAB method of human DNA isolation for polymerase chain reaction  
835 applications. 25:233-235.
- 836 67. Joshi N, Fass J. 2011. Sickle: A Sliding-Window, Adaptive, Quality-Based  
837 Trimming Tool for FastQ Files (Version 1.21)[Software],
- 838 68. Bushnell B. 2014. BBMap: a fast, accurate, splice-aware aligner. Lawrence  
839 Berkeley National Lab.(LBNL), Berkeley, CA (United States),
- 840 69. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G,  
841 Abecasis G, Durbin R. 2009. The sequence alignment/map format and  
842 SAMtools. *Bioinformatics* 25:2078-2079.



- 843 70. Kopylova E, Noé L, Pericard P, Salson M, Touzet H. Sortmerna 2: ribosomal  
844 rna classification for taxonomic assignation, p. *In* (ed),  
845 71. Ren J, Song K, Deng C, Ahlgren NA, Fuhrman JA, Li Y, Xie X, Sun F. 2018.  
846 Identifying viruses from metagenomic data by deep learning. arXiv preprint  
847 arXiv:180607810.  
848 72. McMurdie PJ, Holmes S. 2013. phyloseq: an R package for reproducible  
849 interactive analysis and graphics of microbiome census data. *PLoS one* 8.  
850 73. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation.  
851 *Bioinformatics* 30:2068-9.  
852

## 853 **Figure legends**

854

855 **Figure 1.** Sample processing workflow

856

857 **Figure S1** Krona plots of reads classification using Centrifuge algorithm.

858

859 **Figure 2.** Individual viromes richness. The curves represent the cumulative fraction  
860 of the reads recruited by the contigs plotted against the cumulative length of the  
861 contigs. Contigs were sorted by CPM, taken as a proxy for their relative abundance.  
862 The dotted lines on the panel with the all the samples combined indicate the  
863 modelling function fitted to each of the curves. The distribution of the Shannon and  
864 Simpson diversity indexes in the samples from each population are plotted as box and  
865 whisker plots as inserts of respective panels.

866

867 **Figure 3.** Network graph produced using vCONTACT 2. Red dots – bacteriophage  
868 reference sequences, retrieved from GenBank, green dots – the horse viromes contigs.

869

870 **Figure S2** Graphs of contigs clustered using vCONTACT 2 algorithm by populations.

871

872 **Figure S3** Genome maps of 11 largest contigs subjected to manual analysis. Genes  
873 coloured had a putative function assigned, those in purple have no known function.

874

875 **Figure 4.** Individuality of horse fecal viromes. A – Jaccard's distances between the  
876 samples. B. Abundance of contigs in different samples.

877

878 **Figure 5.** Composition of viromes in the different horse populations. A – distribution  
879 of the contigs detected in three populations. B. – distribution of the VC clusters

880 between the populations. C. – distribution of the VCs, detected in 50% or more of the  
881 samples at least in one of three populations. D. – distribution of the contigs, detected  
882 in 50% or more of the samples at least in one of three populations

883

884 **Figure 6.** Comparison of the samples by abundance of genus-level 16S bacterial  
885 OTUs (top) and by predicted host of phage contigs. Only the OTUs or genera  
886 overlapping between the 16S sequencing results and viral host prediction are shown  
887 in color. Other OTUs are shown in gray scale.

888

889



Feces in extraction buffer containing sodium azide and Tween 20

Clarification of rough extract by filtration, removal of coarse particles

DNase treatment, centrifugation at 10 000 g

Supernatant vacuum filtering on diatomite bed

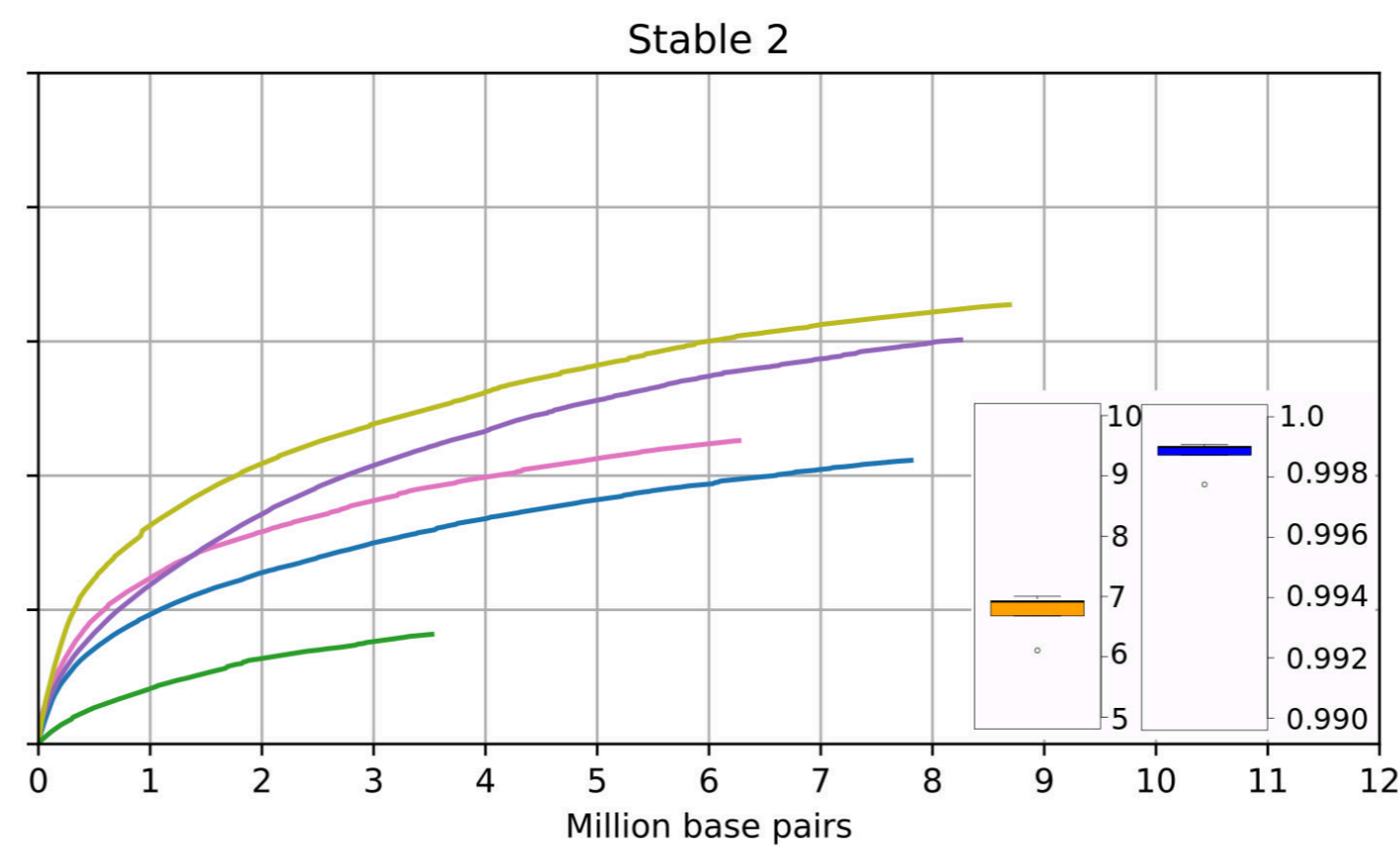
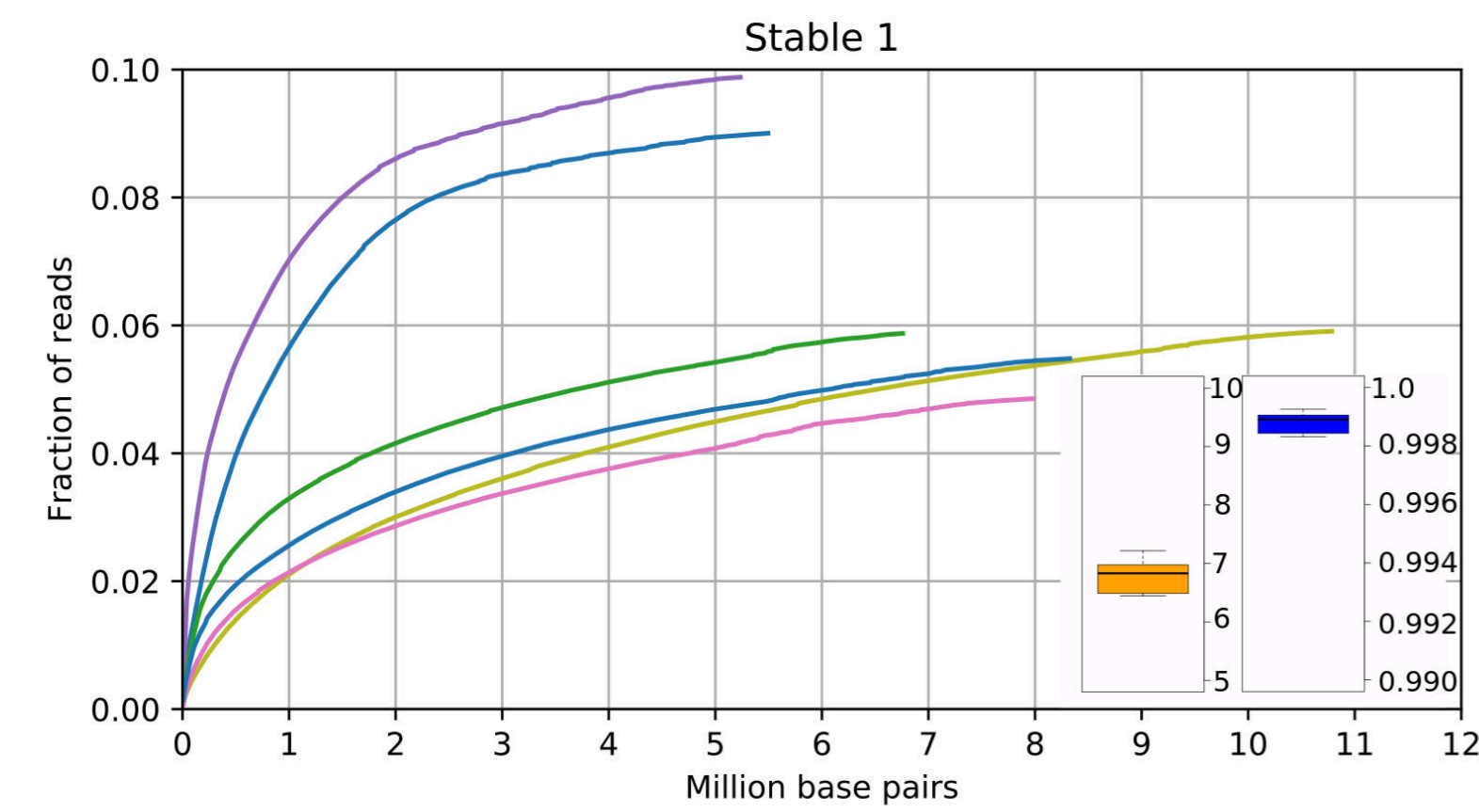
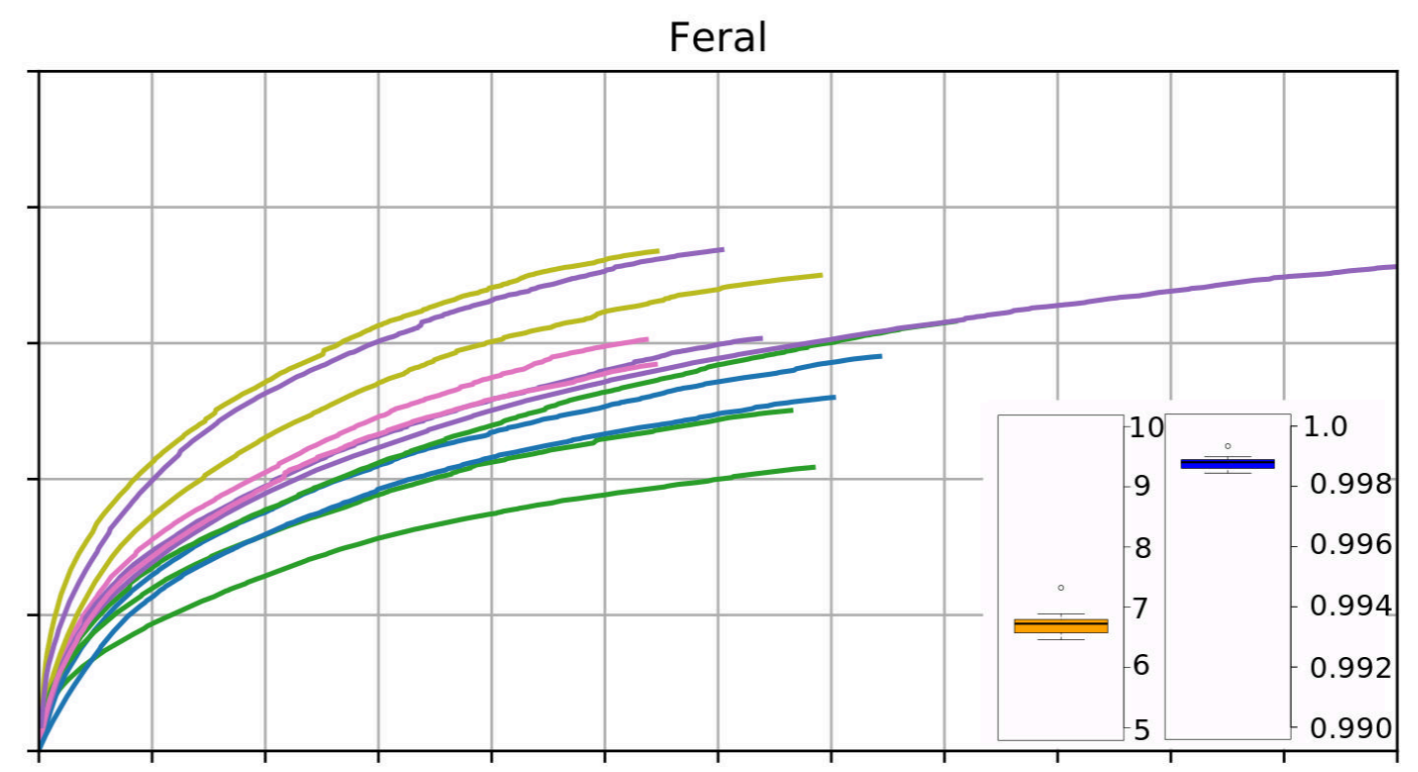
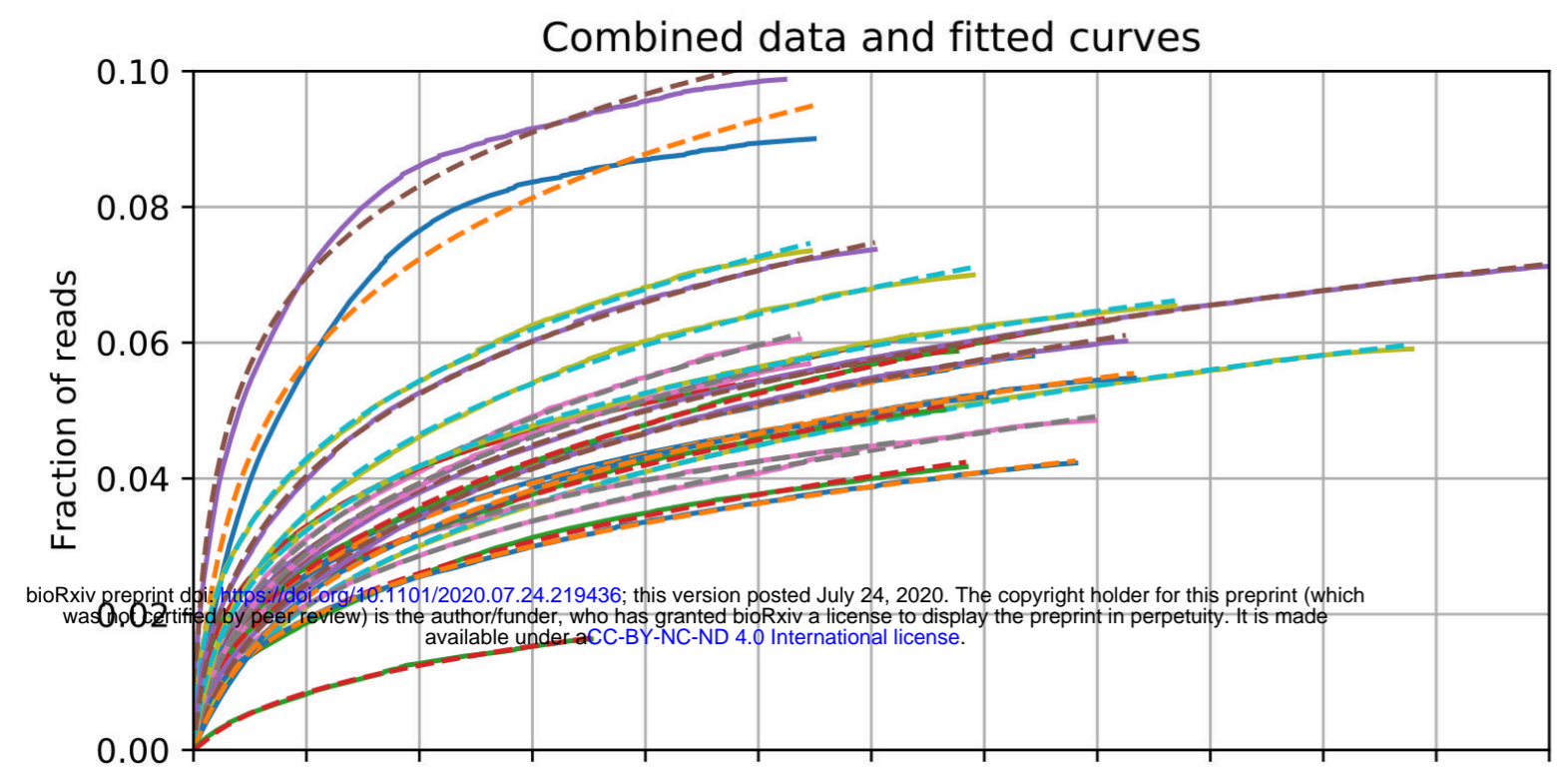
PEG precipitation of virome

Virome DNA extraction using CTAB

DNA precipitation and quantitation

Nextgen sequencing

bioRxiv preprint doi: <https://doi.org/10.1101/2020.07.24.219436>; this version posted July 24, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



- Fe9M11-1
- Fe1M6-1
- Fe1M11-1
- Fe7M6-1
- Fe39M6-1
- Fe2M11-1
- Fe2M11-2
- Fe1M6-2
- Fe9M11-2
- Fe7M6-2
- Fe1M11-2
- Fe39M6-2
- S1Ca-1
- S1Cr-1
- S1Ca-2
- S1Pr
- S1Cr-2
- S1Er
- S2Dt
- S2DI
- S2Mi
- S2Tb

Fig. 2

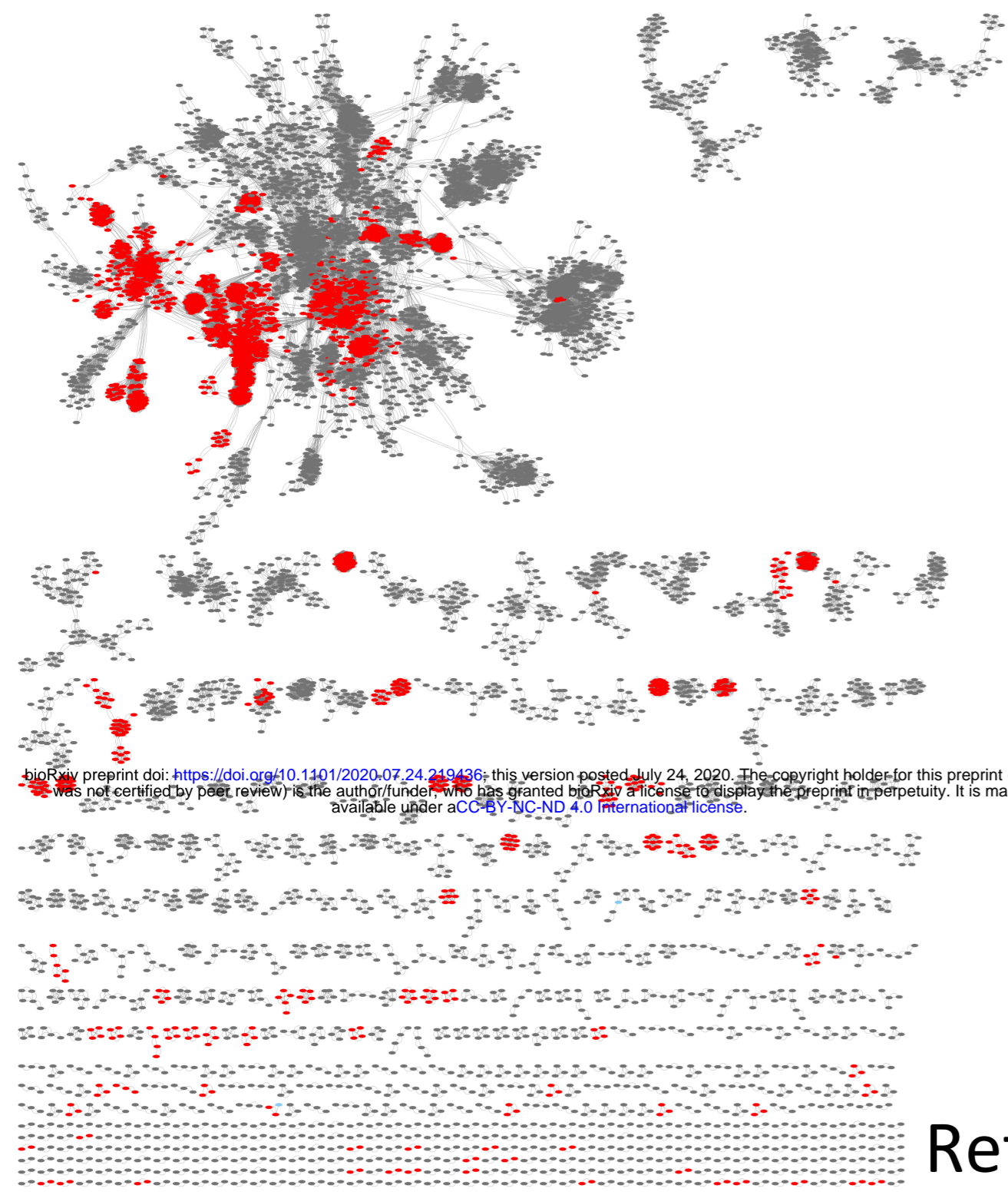


bioRxiv preprint doi: <https://doi.org/10.1101/2020.07.24.219436>; this version posted July 24, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

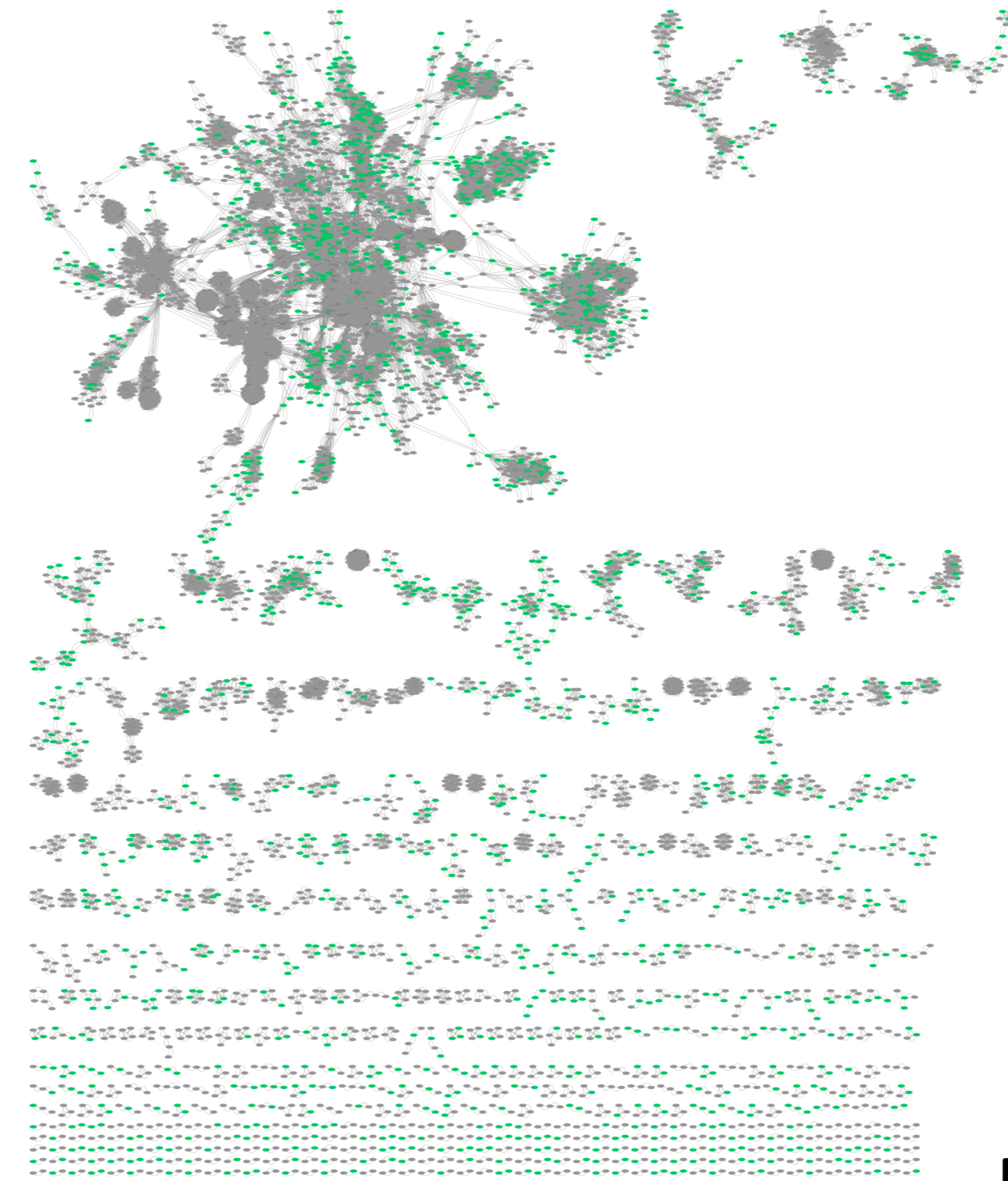


Fig 3

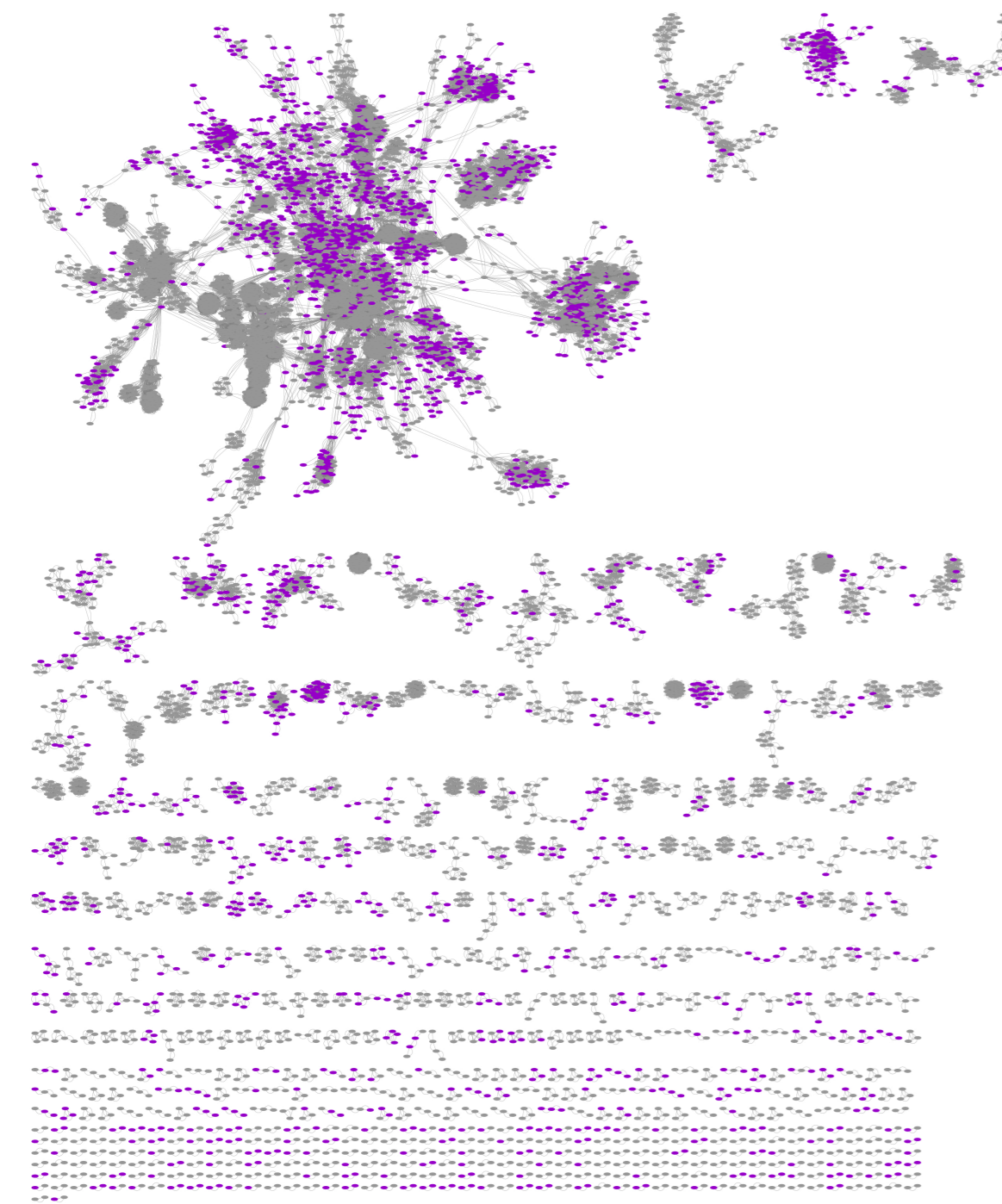




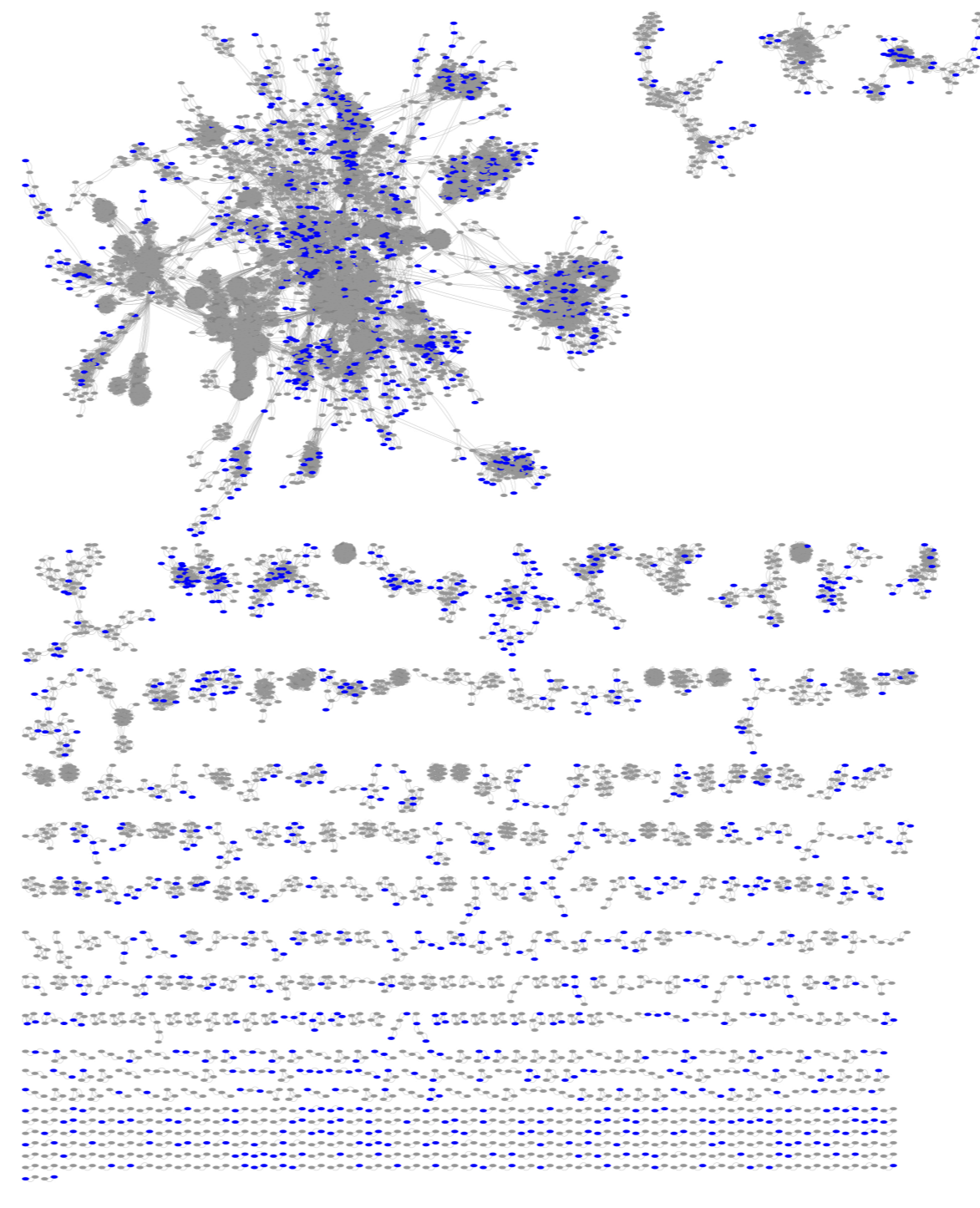
RefSeqs



Feral



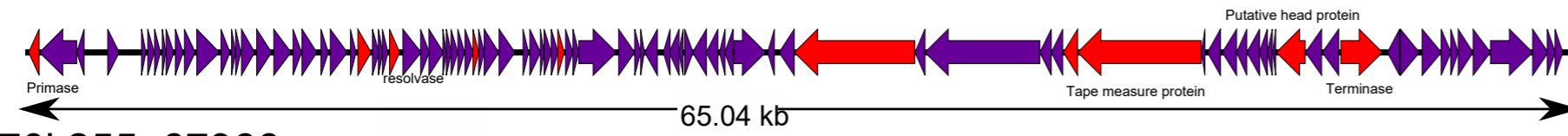
Stable 1



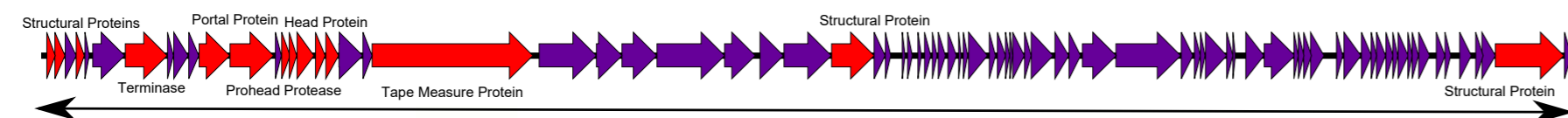
Stable 2



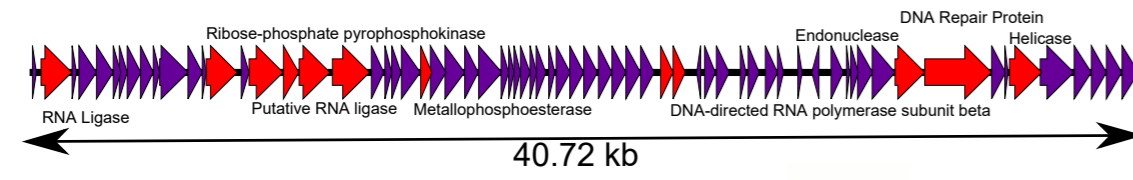
008k255\_111094



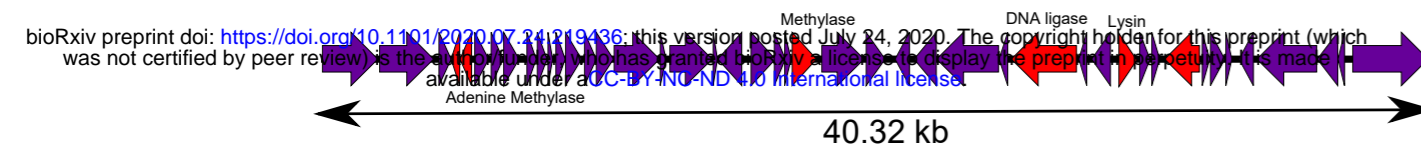
070k255\_67966



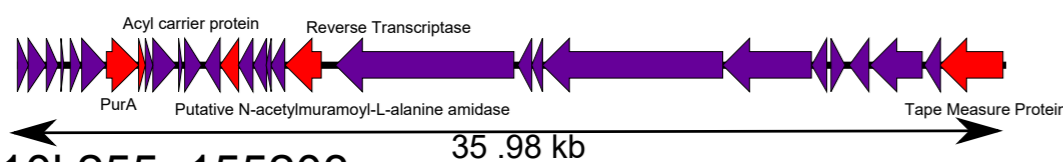
010k255\_55070



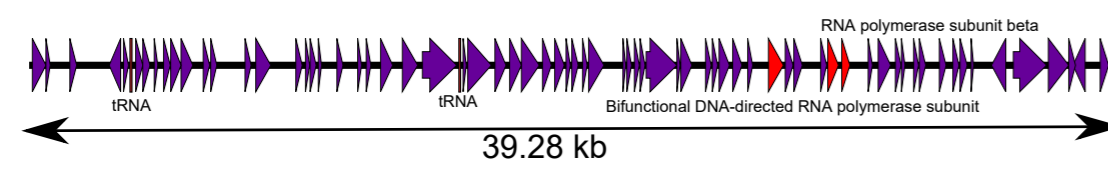
016k255\_18131



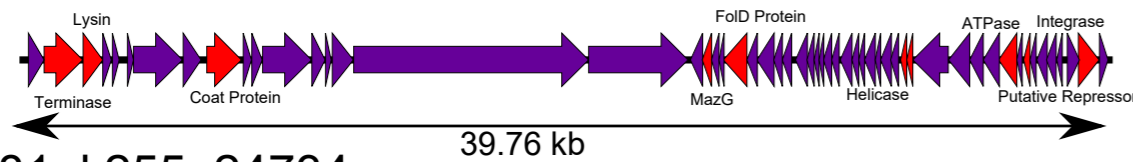
015k255\_41289



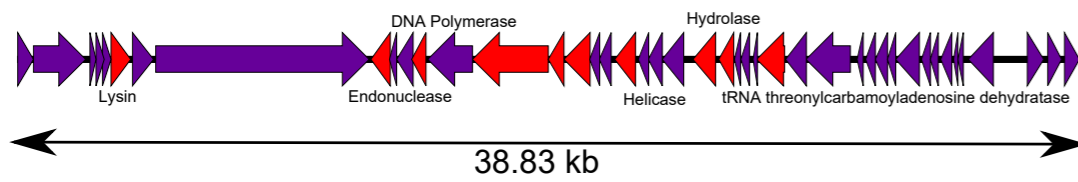
010k255\_155203



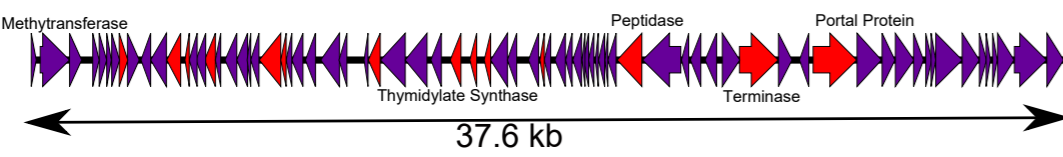
013k255\_49020



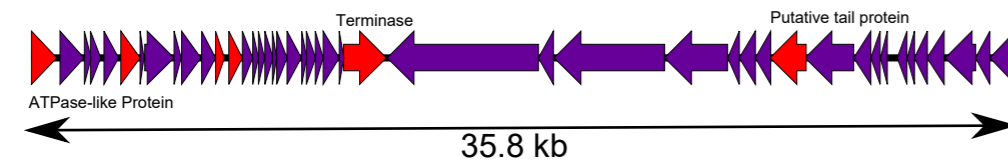
081\_k255\_24734



003k255\_21071



006k255\_28256



011k255\_69359

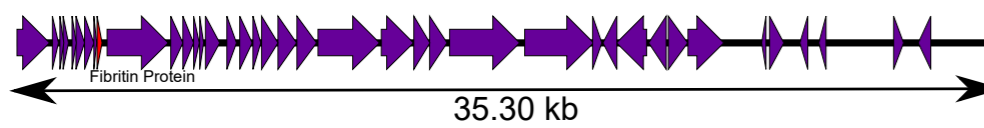


Fig. S3

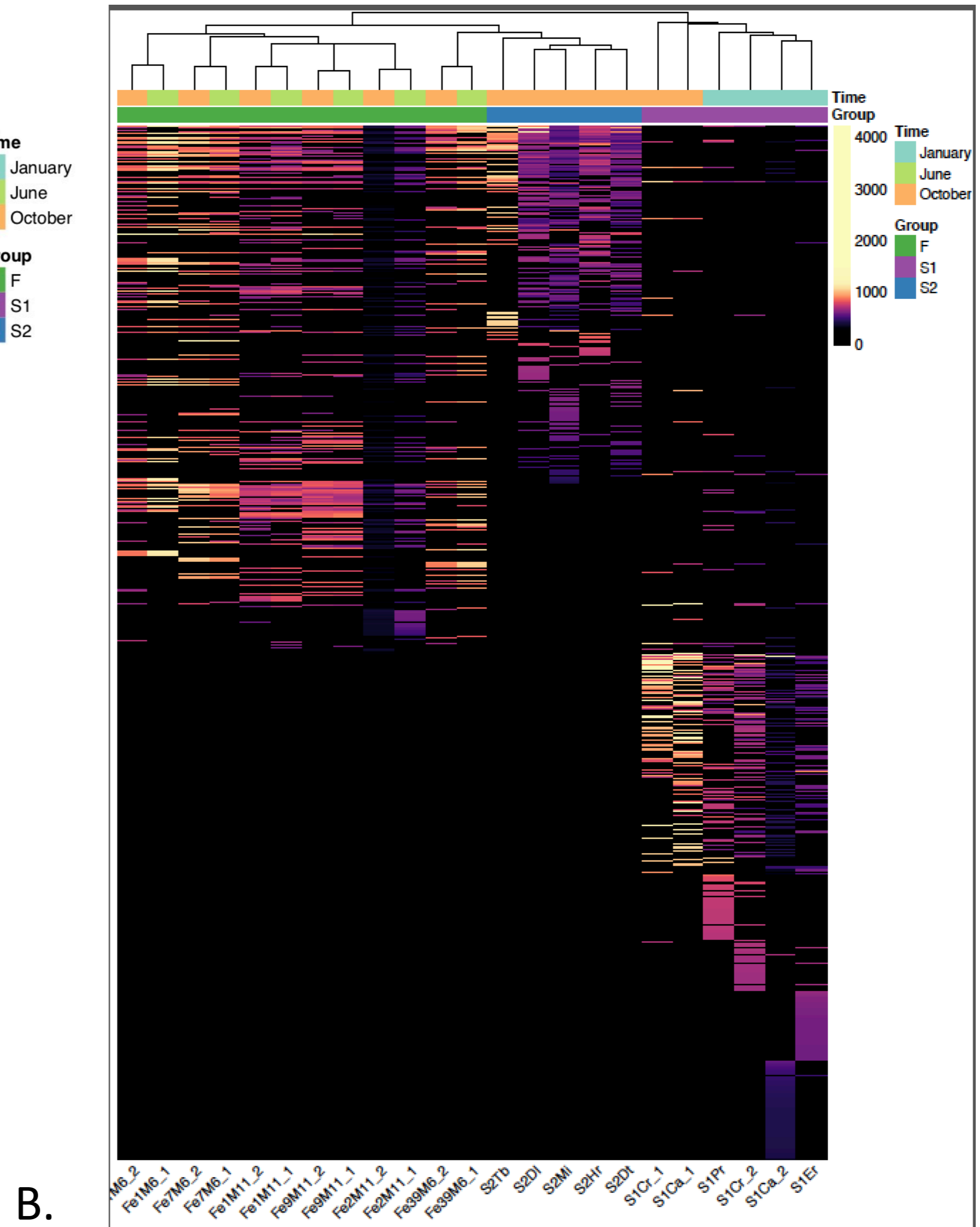
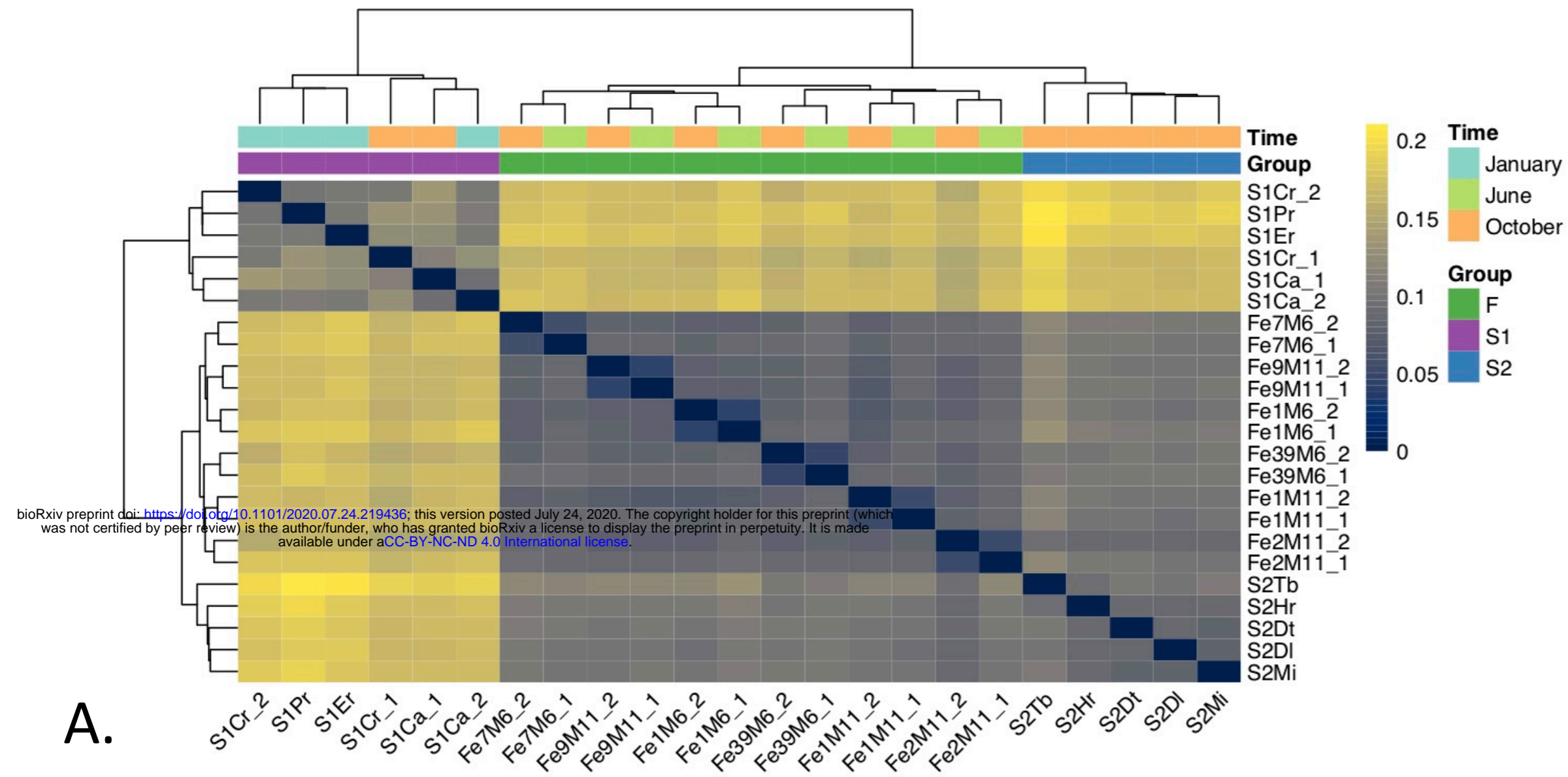


Fig. 4

bioRxiv preprint doi: <https://doi.org/10.1101/2020.07.24.219436>; this version posted July 24, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

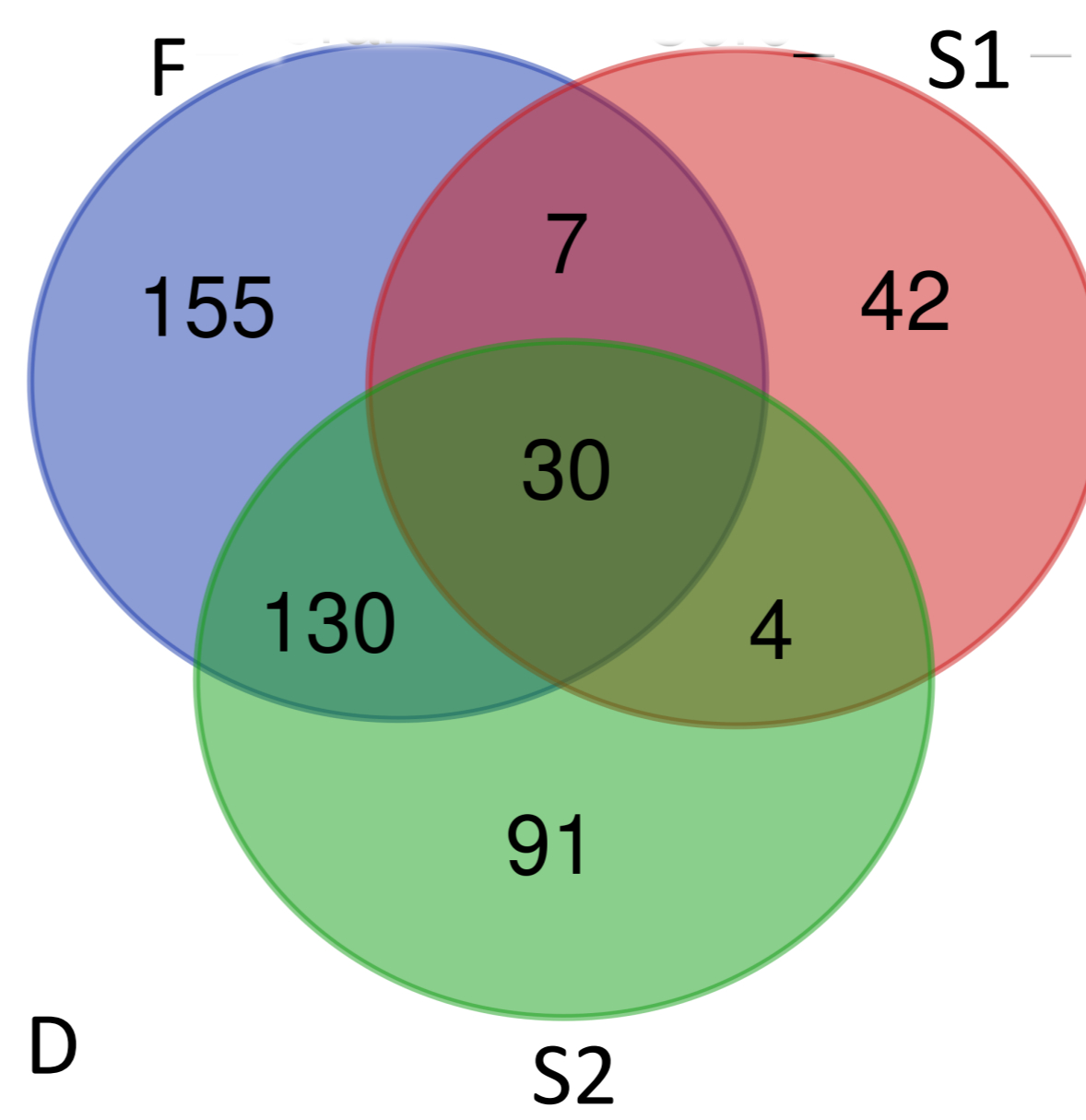
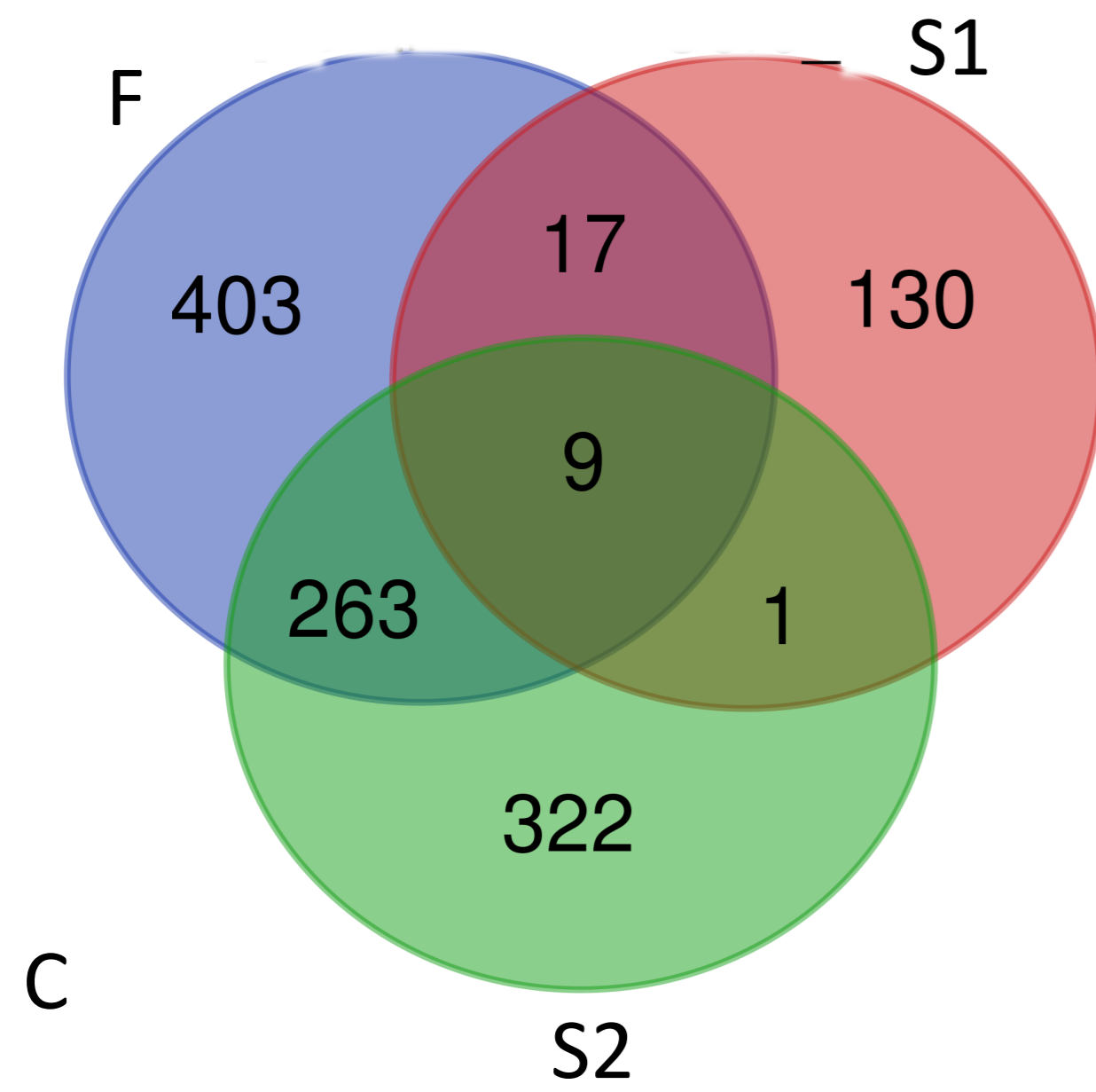
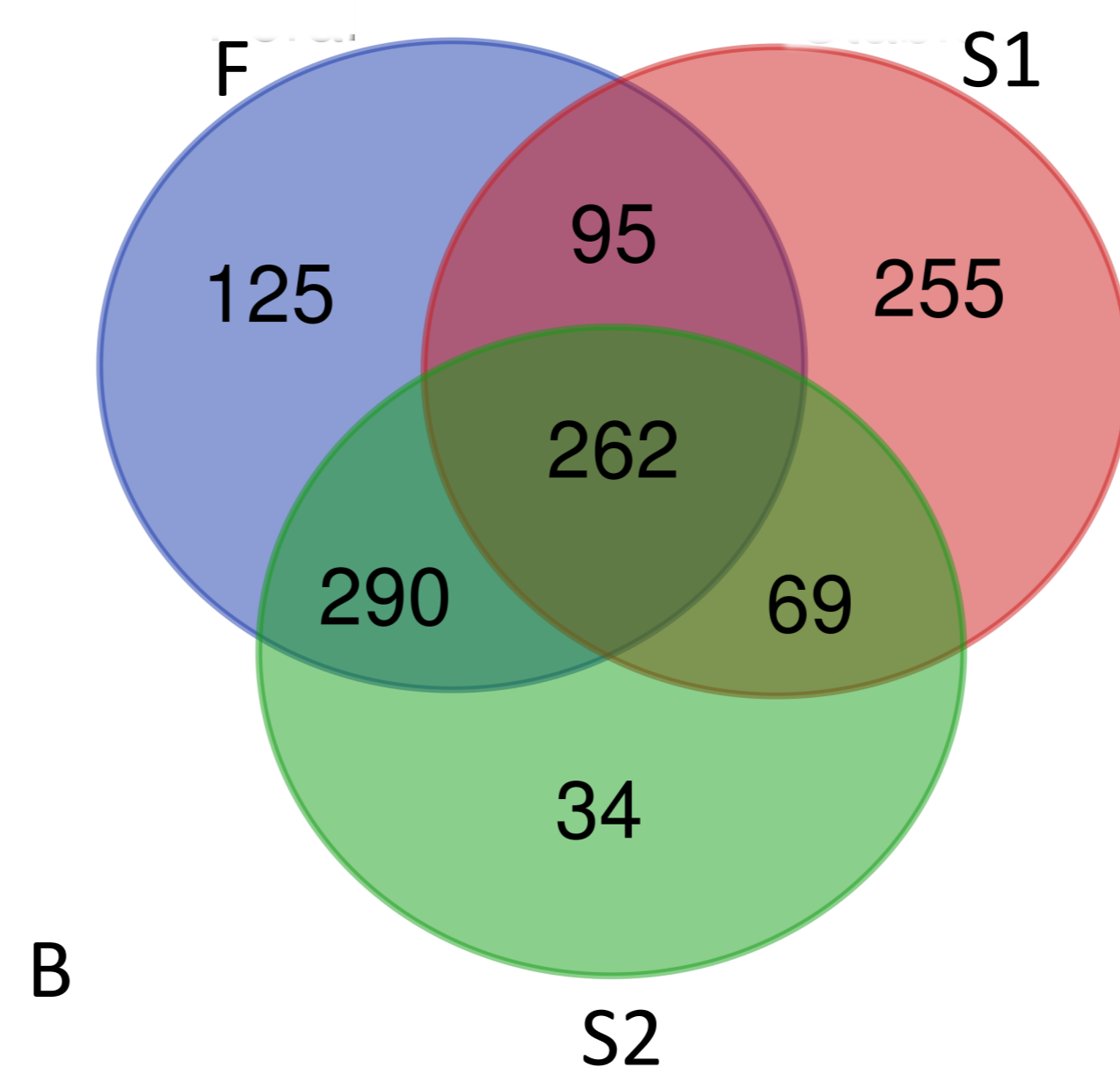
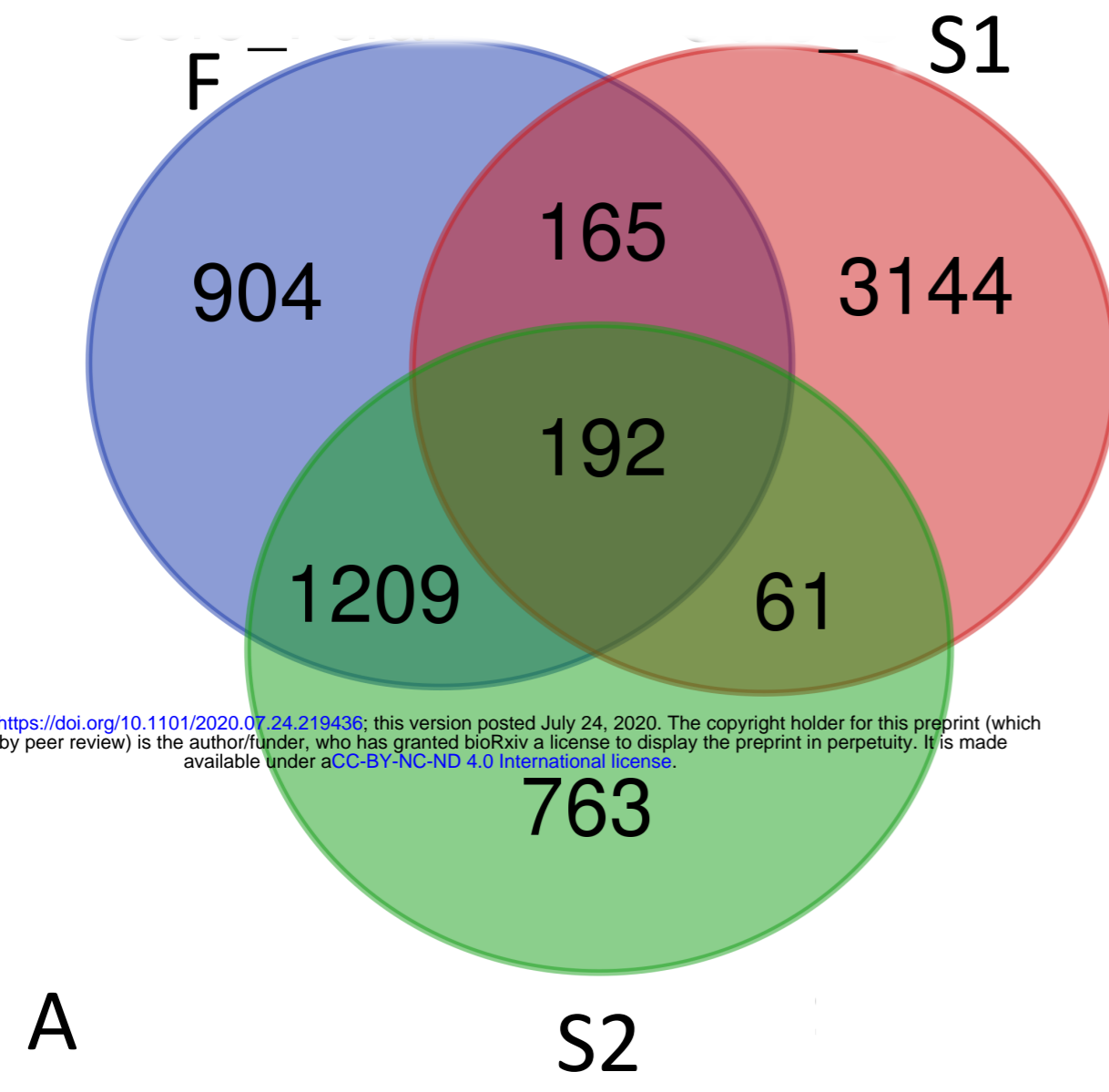


Fig. 4



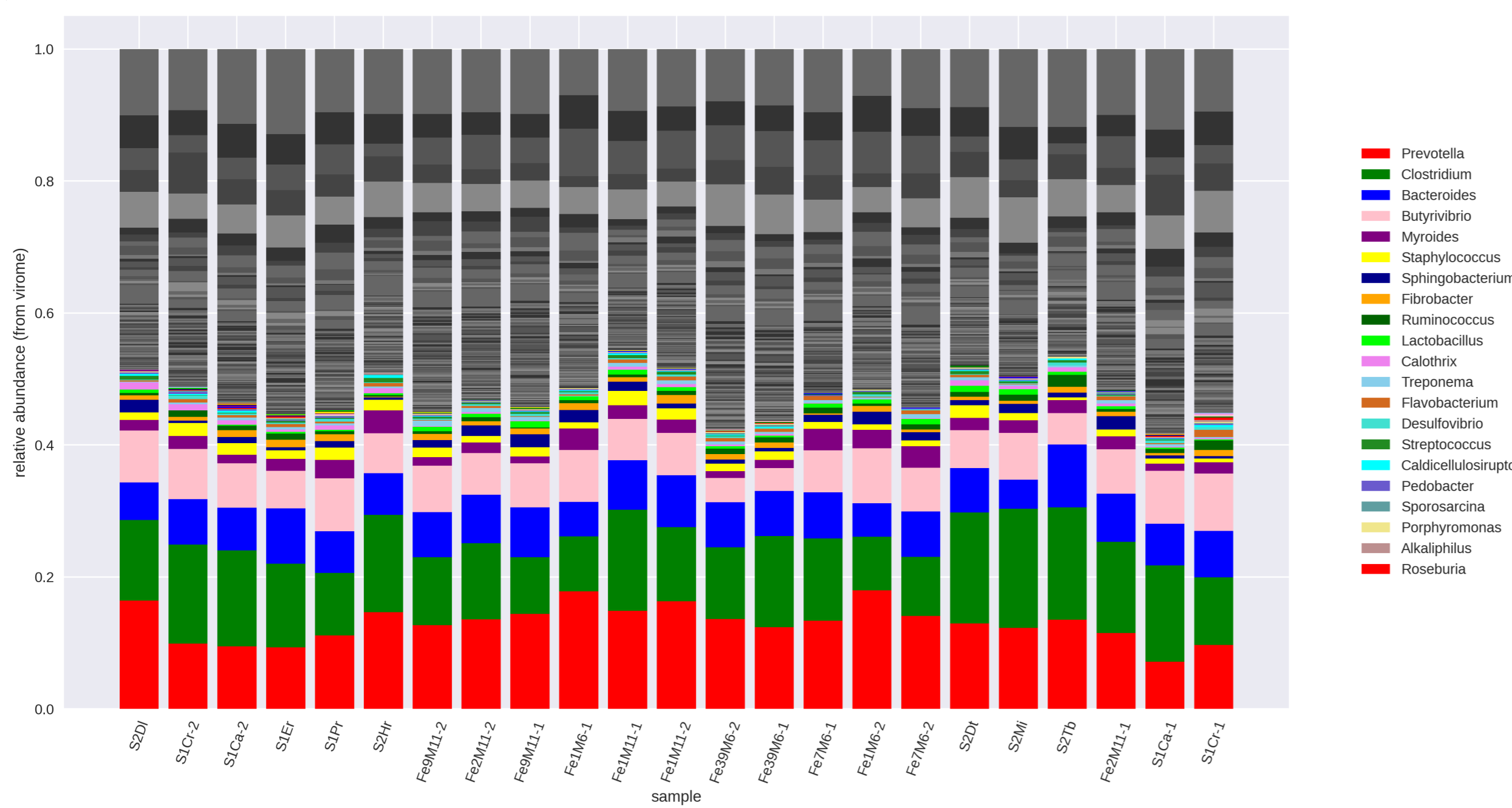
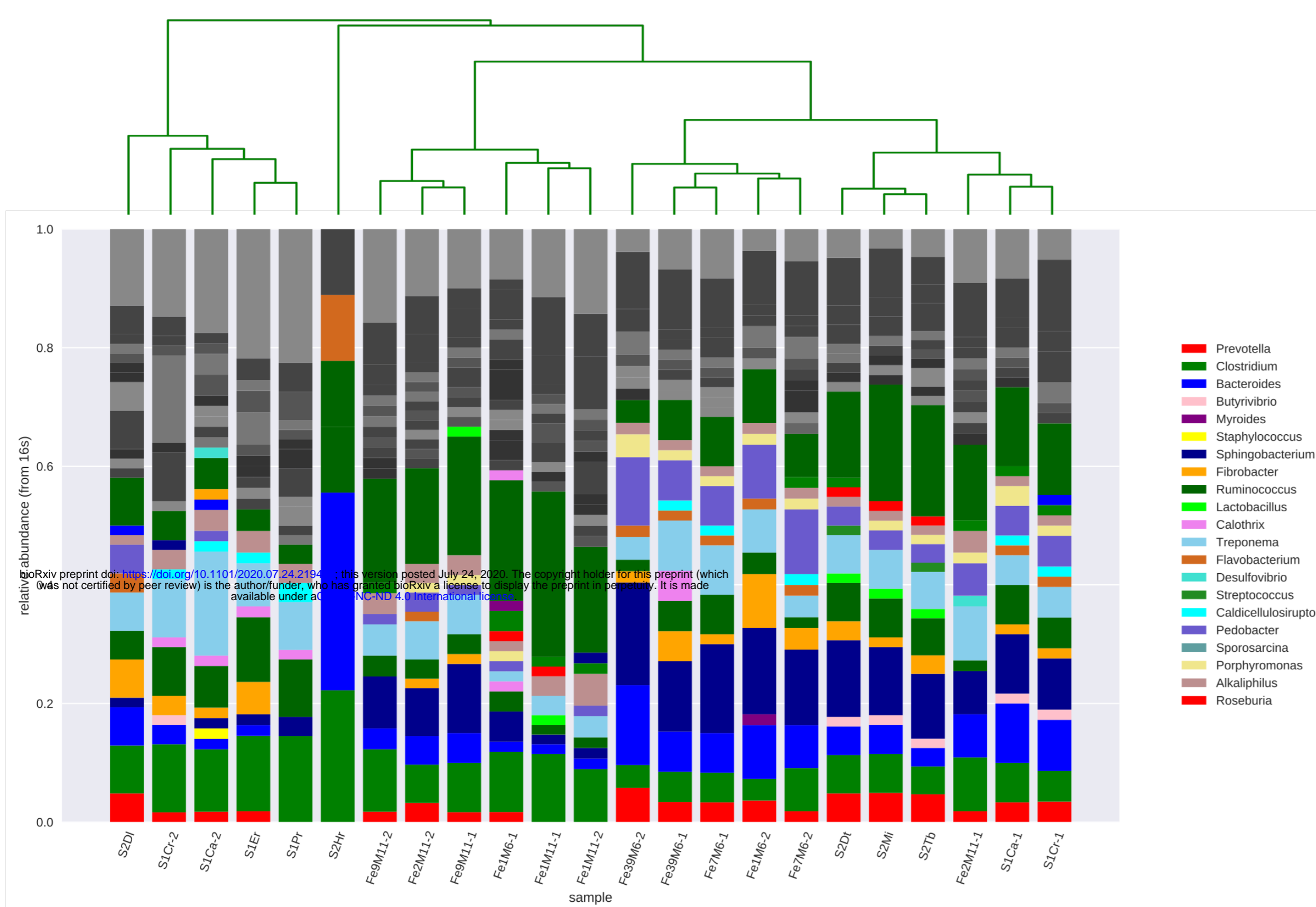


Fig. 6