1    **Bias of library preparation for virome characterization in untreated and treated**

2    **wastewaters**

3    Alba    Pérez-Cataluña[1,2*],    Enric    Cuevas-Ferrando[1],    Walter    Randazzo[1,2,3],    Gloria

4    Sánchez[1,2*]

5    [1]Department of Preservation and Food Safety Technologies, IATA-CSIC, Avda.

6    Agustin Escardino 7, 46980 Paterna, Valencia, Spain

7    [2]European Virus Bioinformatics Center, Leutragraben 1, 07743 Jena, Germany

8    [3]Department of Microbiology and Ecology, University of Valencia, Av. Dr. Moliner,

9    50, 46100 Burjassot, Valencia, Spain

10

11    *Corresponding authors:

12    Gloria Sánchez. Department of Preservation and Food Safety Technologies (IATA-

13    CSIC). Av. Agustín Escardino 7. 46980 Paterna. Valencia. Spain.

14    Tel.: + 34 96 3900022; Fax: + 34 96 3939301; E-mail: gloriasanchez@iata.csic.es

15    Twitter: @FoodViruses

16    Alba Pérez-Cataluña. Department of Preservation and Food Safety Technologies

17    (IATA-CSIC). Av. Agustín Escardino 7. 46980 Paterna. Valencia. Spain.

18    Tel.: + 34 96 3900022; Fax: + 34 96 3939301; E-mail: alba.perez@iata.csic.es Twitter:

19    Twitter: @AlbaPerezCat

20

21

22

23

24

1 **Abstract**

2 The use of metagenomics for virome characterization and its implementation for

3 wastewater analyses, including wastewater-based epidemiology, has increased in the

4 last years. However, the lack of standardized methods can led to highly different results.

5 The aim of this work was to analyze virome profiles in upstream and downstream

6 wastewater samples collected from four wastewater treatment plants (WWTPs) using

7 two different library preparation kits. Viral particles were enriched from wastewater

8 concentrates using a filtration and nuclease digestion procedure prior to total nucleic

9 acid (NA) extraction. Sequencing was performed using the ScriptSeq v2 RNA-Seq (LS)

10 and the NEBNext® Ultra™ II RNA (NB) library preparation kits. Cleaned reads and

11 contigs were annotated using a curated *in-house* database composed by reads assigned

12 to viruses at NCBI. Significant differences in viral families and in the ratio of detection

13 were shown between the two library kits used. The use of LS library showed

14 *Virgaviridae*, *Microviridae* and *Siphoviridae* as the most abundant families; while

15 *Ackermannviridae* and *Helleviridae* were highly represented within the NB library.

16 Additionally, the two sequencing libraries produced outcomes that differed in the

17 detection of viral indicators. These results highlighted the importance of library

18 selection for studying viruses in untreated and treated wastewater. Our results underline

19 the need for further studies to elucidate the influence of sequencing procedures in

20 virome profiles in wastewater matrices in order to improve the knowledge of the virome

21 in the water environment.

22

23 **Keywords:** Wastewater, Metagenomics, Enteric viruses, viability RT-qPCR.

24

## 1. Introduction

The reuse of water, including for irrigation, cooling, and other non-potable applications is an emerging topic due to climate change and water scarcity. Treatment and regeneration of household sewage water in urban regions are usually performed by wastewater treatment plants (WWTPs); however, they are not always able to completely eliminate the microbiological risks present in treated wastewaters (Chalmers et al., 2010; Randazzo et al., 2019; Sano et al., 2016). Fecal bacteria have traditionally been used as indicators for the presence of pathogenic microorganisms even though they fail to detect the presence of human pathogenic enteric viruses (Eslamian, 2016; Gerba et al., 2013; Kitajima et al., 2014). Thus, several viruses (i.e. crAssphage, Pepper mild mottle virus, adenovirus, polyomavirus, …) have been proposed as indicators because of their similarity to pathogenic viruses in terms of environmental stability and resistance to wastewater sanitation treatments (Farkas et al., 2020). The presence of human enteric viruses in treated wastewaters has been well documented (Gerba et al., 2018; Randazzo et al., 2019; Sano et al., 2016), posing public health risk-related concerns also because of their stability into the environment. Thus far, nearly one hundred different types of human enteric viruses are known, which cause a variety of illnesses and diseases in humans (Fong and Lipp, 2005), primarily gastroenteritis and hepatitis, and new pathogenic strains and species continue to be discovered. Among others, the viruses most commonly detected in untreated and treated wastewaters include human norovirus, adenovirus (AdV), enterovirus (EV), sapovirus (SaV), astrovirus (HAstV), rotavirus A (RV), and hepatitis A and E viruses (HAV and HEV) (Haramoto et al., 2018). Surveillance of human enteric viruses in untreated and treated wastewaters is performed by molecular procedures (e.g., real time PCR (qPCR) or digital PCR (dPCR)) (Haramoto et al., 2018). Currently, a wastewater-based

26 epidemiology surveillance has been globally implemented to monitor COVID-19

27 disease, caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)

28 with notable implications for public health response in local settings (Bivins et al.,

29 2020; Polo et al., 2020). These approaches require reference sequences for primer and

30 probe design which limit the number and variety of viruses to be analyzed.

31 Alternatively, recent shotgun or untargeted metagenomic approaches enable the

32 simultaneous identification of viral sequences from a sample, referred to as 'virome',

33 which is a diverse community of mainly eukaryotic RNA and DNA eukaryotic viruses

34 and bacteriophages. Virome characterization in wastewater provides a potential solution

35 to the challenges associated with the traditional surveillance of viruses in sewage

36 (Nieuwenhuijse and Koopmans, 2017).

37 In this pilot study, we have used metagenomics analyses using two different library

38 preparation kits for metagenomic sequencing to characterize the virome composition in

39 influent and effluent samples from four different WWTPs. Thus, the objectives of this

40 study were to: 1) evaluate different sequencing libraries for virome characterization; and

41 2) investigate virome distribution and diversity in influent and effluent samples.

42 **2. Materials and Methods**

43 **2.1. Sample processing**

44 Five-hundred mL of influent (IW) and effluent (EW) grab samples from four different

45 WWTPs were collected in November 2018 in Valencia (Spain). Treatment plants

46 differed in the number of equivalent population, the volume of treated wastewater and

47 the disinfection treatments (Table S1). *Escherichia coli* counts, expressed as Most

48 Probably Number (MPN), were performed using the Colilert® kit (IDDEX Laboratories,

49 Spain) following the ISO 9308-2:2012 standard on the same sampling day. Samples

50 were kept for further analyses at -20°C and thawed for 12 h at approximately 20 °C

before processing. After thawing, 200 mL of each sample were inoculated with 7 log PCRU/L of mengovirus (MgV) vMC$_0$ (CECT 100000), used as a process control. Samples were processed using the aluminum-based precipitation protocol described elsewhere (AAVV, 2018; Randazzo et al., 2019). Briefly, 200 mL of sample was adjusted to pH 6.0. The Al(OH)$_3$ precipitate was performed mixing 1 part of AlCl$_3$ 0.9N per 100 parts of sample and the solution was mixed at 150 rpm for 15 min. Then, samples were centrifuged at 1,700 x $g$ for 20 min and the pellet was resuspended in 10 mL of 3% beef extract (pH 7.4) and shaken at room temperature (RT) for 10 min at 150 rpm. Finally, samples were centrifuged for 30 min at 1,900 x $g$ and the resulting pellet was resuspended in 1 mL phosphate saline buffer (PBS, pH 7.4) and stored at -80°C.

### 2.2. Sample processing for metagenomics

Viral particles were enriched from sample concentrates (n=8) following the NetoVIR protocol, which includes both filtration and nuclease digestion steps (Conceição-Neto et al., 2015). In brief, 500 µL of concentrates were homogenized using the MP FastPrep24 5G (MP Biomedicals, Spain) for 40 seconds at a speed of 6.0. The homogenate was centrifuged at 16,000 × $g$ for 3 min and 200 µL of the supernatant was filtered through a 0.8 µm PES filter (Sartorius, UK) to remove large particles. The filtrate was incubated with benzonase (Millipore, Spain) and microccocal nuclease (New England Biolabs, USA) enzymes at 37°C for 2 h to degrade free nucleic acids. Capsid protected viral nucleic acids were extracted with the NucleoSpin®RNA virus kit (Macherey-Nagel GmbH & Co., Germany), according to the manufacturer's instructions, without adding carrier RNA. Thus, both DNA and RNA viral nucleic acids were concomitantly extracted. Nucleic acids were eluted in 50 µL RNase-free water. Libraries were generated from 1 to 50 ng of a DNA-RNA sample using two different library preparation kits. The first library preparation kit was the ScriptSeq v2 RNA-Seq Library

76  Preparation Kit (Illumina, USA), referenced as LS, with slight modifications. An initial

77  denaturation step (95 °C for 5 min) was added to the protocol, and PCR cycles were

78  increased to 20 to obtain enough library concentration to sequence. Additionally, the RT

79  enzyme from the original library preparation kit was substituted by the AMV Reverse

80  Transcriptase (Promega, Spain). The second library preparation kit was the NEBNext®

81  Ultra™ II RNA Library Prep Kit (New England BioLabs Inc., Ipswich, UK) (referenced

82  as NB) following manufacturer's instructions. The two libraries compared in this study

83  differ in terms of fragmentation times, enzymes, cDNA synthesis conditions, primers

84  used in the PCR, as well as the conditions for aforesaid amplification (Table S2).

85  Libraries were normalized, pooled, and sequenced using the NextSeq™ 500 platform

86  (Illumina), following the manufacturer's protocol, with a configuration of 150 cycles

87  paired-end reads. Sequencing was performed by Lifesequencing S.L. (Valencia, Spain).

88    **2.3. Data analyses**

89  Obtained reads were cleaned for adaptor removal using cutadapt software (Martin,

90  2011) with a minimum overlap of 5 nucleotides between read and adaptor and a

91  maximum error rate of 0.1. Reads were cleaned with the *reformat.sh* script from

92  BBMap software (sourceforge.net/projects/bbmap/) in order to remove nucleotides from

93  both ends with Phred scores lower than 20 and reads shorter than 50 bp. Cleaned reads

94  were merged in to single reads with FLASH v1.2.11 (Magoč and Salzberg, 2011)

95  allowing outies. Additionally, cleaned reads were assembled with Ray 2.3.1. (Boisvert

96  et al., 2012) using 31-mers.

97  Merged reads and contigs were taxonomically annotated using BLASTn algorithm

98  (Boratyn et al., 2013) with a manually curated *in-house* database constructed with all

99  the viral sequences (NCBI:txid10239; release May, 5 2020) available at GenBank

100  (https://www.ncbi.nlm.nih.gov/nuccore/?term=viruses%5Borganism%5D).     For     the

BLASTn analysis of viral reads against this curated *in-house* database, a cut off of 70% of query sequence coverage and 80% of identity was used, respectively. Rarefaction curves and diversity indexes Shannon and Simpson were calculated with R package *vegan* v2.5-6.

## 2.4. Virus quantification

For virus quantification an optimized viability RT-qPCR was applied as previously described (Randazzo et al., 2019). In brief, 150 μL sample concentrates were added to 50 μM PMAxx (Biotium, USA) and 0.5% Triton 100-X (Thermo Fisher Scientific, Spain) and incubated in the dark at RT for 10 min at 150 rpm. Then, samples were exposed to photo-activation using a photo-activation system (Led-Active Blue, GenIUL, Spain) for 15 min. RNA was extracted using the NucleoSpin® RNA virus kit (Macherey-Nagel GmbH & Co.) according to the manufacturer's instructions including the Plant RNA Isolation Aid (Ambion, Spain) pretreatment. Primers, probes and RT-qPCR conditions for norovirus GI, norovirus GII, RV, HAV, HEV, mengovirus and HAstV quantification have been previously reported (Randazzo et al., 2019, Cuevas-Ferrando et al., 2020).

For crAssphage quantification by qPCR, the primer set CPQ_064 described by Stachler et al. (2017) was used. PCR conditions were an initial denaturation step of 30 seconds at 95˚C followed by 45 cycles of 5s at 95˚C and 30s at 60˚C. The Premix Ex Taq master mix for probe-based real-time PCR kit (Takara, France) was used for the reaction. For the crAssphage quantification, the standard curve was performed with a customized gBlock® fragment (Integrated DNA Technologies, Spain) of 228 bp that contained the crAssphage sequences used for amplification.

5

125     Limit of quantification, qPCR efficiency and standard curve $R^2$ values for all the tested

126     genes are displayed in Table S3. For all RT-qPCR assays, undiluted and ten-fold diluted

127     RNA was tested to check for RT-qPCR inhibitors.

128     **2.5. Correlation and similarity analyses**

129     Correlation analyses were carried out between data sets obtained by both libraries at

130     family level, and between metagenomics and RT-qPCR results using the R package

131     *Hmisc* v4.2-0 (https://CRAN.R-project.org/package=Hmisc) and applying the Spearman

132     method (ρ). Significance was set at 0.05. Representation of correlation matrix values

133     was performed with the R library *corrplot* v0.84 (https://CRAN.R-

134     project.org/package=corrplot).

135     For each individual sample, the Jaccard index was used to analyze the similarity among

136     results obtained with both libraries. Calculations were performed using R package

137     *betapart* v1.5.2 (Baselga, 2010) taking into account the beta.JAC values representing

138     the overall beta diversity for each sample pair.

139     **3.   Results**

140     **3.1. Overview of bias due to library preparation**

141     Each concentrated sample was sequenced using two sequencing libraries: the ScriptSeq

142     v2 RNA-Seq Library Preparation Kit (LS) and the NEBNext® Ultra™ II RNA Library

143     Prep Kit (NB). The average number of reads was 3.2 and 11.5 million for LS and NB

144     libraries, respectively. Rarefaction analyses showed that 5 out of 8 samples sequenced

145     by the LS library reached the plateau, while 2 out of 8 samples sequenced by NB library

146     reached it. Despite that, remaining samples were close to stabilization with both

147     libraries (Fig. S1). Merged viral reads were annotated through a BLASTn comparison

148     with the curated *in-house* database that comprised all the viral sequences (CDS and

149    complete genomes) available at GenBank. For the LS library, the percentage of viral

150    reads ranged from 0.6% to 2.4% in influent and from 0.4% to 4.4% in effluent samples.

151    For NB library, the BLASTn analysis showed a high number of sequences ascribed to

152    the same taxon, suggesting an overrepresentation due to sequencing bias, representing

153    between 33 and 60% of the total viral reads. For that reason, the relative calculations of

154    subsequent analyses were made also taking into account this overrepresentation. These

155    corrected calculations will be called NB-corrected. For the NB library, viral reads

156    ranged from 38% to 58% in influent and from 14% to 24% in effluent samples. Taking

157    into account the results of NB-corrected, these percentages ranged from 9% to 12% and

158    from 7% to 12% in influent and effluent samples, respectively. Shannon and Simpson

159    diversity indexes were calculated for each type of sample (influent or effluent) and for

160    each library (LS or NB) (Fig. 1). Shannon indexes were higher in influent samples

161    sequenced with LS library (mean values of 3.85±0.33 for LS and 1.30±0.12 for NB);

162    however, for effluent samples both indexes showed similar means (1.81±1.21 for LS

163    and 1.90±0.2 for NB), being effluent samples sequenced with LS library more variable

164    (0.36-3.07). Similar results were obtained for Simpson index, even though the mean

165    values for influent samples were highly different (0.93±0.02 for LS and 0.62±0.05 for

166    NB).

167    Raw data was deposited at SRA under the Bioproject PRJNA67378 with the following

168    accession numbers: SAMN16633937-SAMN16633944 for ScriptSeq v2 RNA-Seq

169    Library Preparation Kit samples and SAMN16634071-SAMN16634078 for NEBNext®

170    Ultra™ II RNA Library Prep Kit samples.

171    **3.2. Mengovirus recovery**

172    Mengovirus (MgV) was used as a process control to analyze the performance of each

173    library to recover reads and the entire genome of MgV. Its recovery, represented as the

174  percentage of viral reads and the percentage of MgV isolate M genome (L22089.1)

175  obtained for each sample with each library, was different depending on the library used.

176  For LS library, the percentage of viral reads of MgV ranged from 0.05% to 0.79% in

177  influent and from 0.35% to 3.68% in effluent samples. For NB libraries, these values

178  ranged from 0.01% to 0.16% in influent samples, and from 0.63% to 5.77% in effluent

179  samples. However, the percentage of MgV reads with the NB-corrected values were

180  higher in effluent samples, ranging from 1.38% to 11.38%. For the analysis of the

181  recovery of MgV genome, assembled contigs belonging to this species were compared

182  with the genome of Mengovirus isolate M (L22089.1). LS library genome recovery

183  ranges from 6.0% to 95.1%. The highest recovery was obtained in the sample IW3. On

184  the other hand, the coverage of this genome by the NB library ranged from 98.4% to

185  100% (EW3).

186  **3.3. Virome comparison**

187  Regarding the virome composition for each library at family levels, results showed high

188  differences between the two approaches (Fig. 2). While the most represented families

189  with the LS library were *Virgaviridae*, *Microviridae* and *Siphoviridae*; the most

190  abundant families with the NB library were *Ackermannviridae* and *Helleviridae*. These

191  differences allowed the detection of some viral families depending on the library used.

192  For example, families as *Rhabdoviridae*, *Pospiviroidae* or *Mitoviridae* were only

193  detected when the NB library was used for sequencing. Also the taxa uncultured human

194  fecal virus (NCBI:txid239364) and uncultured marine virus (NCBI:txid186617) were

195  only detected with the NB approach. Regarding the families detected with both

196  libraries, only *Genomoviridae* showed total correlation ($\rho=1$) between values obtained

197  with both libraries in influent and effluent samples. For influent wastewater samples,

198  families *Nairoviridae* and *Virgaviridae* showed total correlations; however, this

correlation was only observed for *Parvoviridae* in effluent samples. Other families

showed high correlations (ρ=0.8) in influent wastewaters, as *Peribunyaviridae* and

*Picornaviridae*; while *Podoviridae*, *Poxviridae*, *Reoviridae* and *Virgaviridae* families

showed high correlations (ρ=0.8) in effluent samples. Jaccard indexes showed

similarities between the same sample sequenced with each library that ranged from 0.76

(IW1) to 0.91 (IW2), with mean values of 0.83±0.09 for IWs and 0.81±0.04 for EWs.

## 3.4. Analyses of viral fecal indicators by NGS and correlation with enteric viruses detected by RT-qPCR

Each of the libraries used in this study showed different power for detecting fecal

indicators. Similarly, influent and effluent samples showed different detections rates

(Fig. 3A). For example, LS library detected CrAssphage with read percentages higher

than 1% but these percentages decreased to less than 0.01% with NB library. Most

importantly, LS library was unable to detect the fecal indicator adenovirus in effluent

wastewaters, while the NB library detected adenoviruses in percentages between 0.1-

1%. The same scenario was observed for the Picobirnavirus indicator. However,

indicators families as *Inoviridae*, *Microviridae*, *Myoviridae* and *Podoviridae* showed a

better detection with the LS library. *Siphoviridae* family detection did not show

differences in its detection capacity between the two different tested libraries, with the

exception of sample EW3 (Fig. 3A).

Correlation between the number of reads of proposed viral indicators obtained with both

libraries and the quantifications obtained by RT-qPCR for enteric viruses along with the

*E. coli* counts were calculated. Figure 3B shows the Spearman values of correlation (ρ)

calculated with 95% confidence. Norovirus GI and GII showed high correlation values

(ρ>0.8) with the indicators crAssphage, Picobirnavirus and *Inoviridae*. The highest

correlation values (ρ=0.7) between RV and indicators reads were crAssphage,

9

*Inoviridae* and *Microviridae*. For HAstrV, high correlations ($\rho$>0.8) only occurred with

crAssphage and *Myoviridae*. HAV and HEV did not correlate with any of the indicators.

Interestingly, the proposed indicator AdV showed negative correlations with all the

enteric viruses analyzed ($\rho$>-0.7), with the exception of HAV and HEV. Results

obtained by NGS for the pepper mild mottle virus (PMMoV), also proposed as viral

indicator, showed no correlation with any of the enteric viruses.

### 3.5. **Virome comparison between influent and effluent wastewaters**

Clean reads obtained from each library and each sample were assembled and contigs

longer than 200 bp were used for taxonomical classification by using BLASTn

algorithm and the *in-house* database. Due to the different results observed at reads level

for each library, contigs classification obtained with both libraries for each sample were

merged for results representation and virome analysis. The relative abundances of

different taxa are shown in Fig. 4. As observed in the heatmap graphic, the most

abundant viruses were bacteriophages, as Dickeya phage or Listeria phage WIL-3, even

with higher percentages in effluent samples. The higher detection of these phages in

treated samples, as occurred with other species (i.e. cucumber green mottle mosaic

virus, EBPR podovirus 2, PMMoV, Stealth virus 1, or Tobacco and tomato mosaic

viruses), can be due to the decrease of other viruses after wastewater treatment that

allows the its detection. Similarly, this effect could be the responsible of the detection of

some viruses in effluent samples that were not detected in influent samples, as the case

of human adenovirus and human gammaherpesvirus. Some viruses were only found in

high percentages in influent samples, as the indicator crAssphage, some Aeromonas

phages, Escherichia phages or viruses belonging to the *Microviridae* family.

Wastewater treatments could be the factor that produce this decrease; however, more

248  studies along time from the same WWTPs must be performed in order to ensure the

249  effect of performed treatments.

## 4. Discussion

251  The virome of wastewaters have been previously characterized from samples collected

252  around the world (Adriaenssens et al., 2018; Aw et al., 2014; Cantalupo et al., 2011;

253  Fernandez-Cassi et al., 2018; Furtak et al., 2016; Nieuwenhuijse et al., 2020; Rusiñol et

254  al., 2020; Strubbia et al., 2019a, 2019b; Wang et al., 2018); however, much less is

255  known about the virome of effluent samples as only one study has analyzed two effluent

256  samples collected in the UK (Adriaenssens et al., 2018). As far as we know, this is the

257  first study that concomitantly analyzes the RNA and DNA viruses present in influent

258  and effluent samples besides providing a comparison of viruses profiles detected using

259  different sequencing library kits.

260  Results obtained in our study showed high differences regarding not only viruses, but

261  also the power of detection of viral fecal indicators. Both aspects are important for the

262  use of random metagenomics as tool for specific detection. Our results evidenced the

263  influence of the library used for virome studies together with their variability.

264  Additionally, by using MgV as process control for both metagenomic and RT-qPCR

265  analyses, we further assessed the sensitivity of each library, being higher when using

266  NB library. Recoveries of MgV complete genome were between 6.0 and 95.1% for LS

267  library and between 98.4 and 100% for NB library. In contrast, in a recent study, MgV

268  reads were not recovered from spiked water and sediment samples (Adriaenssens et al.,

269  2020). According to the authors, this was likely due to an inclusion of an inactivation

270  step of the DNase at 75°C, which potentially exacerbated the effect of the RNase step

271  (Adriaenssens et al., 2018). The use of models of a virus of interest when comparing

272  sequencing libraries can be an excellent tool for the library selection.

273  For the analysis of the virome of influent and effluent wastewaters, results obtained by

274  both libraries were merged. Phages as crAssphage, *Aeromonas* phages, *Escherichia*

275  phages or viruses belonging to the *Microviridae* family were found in high percentages

276  in influent wastewaters. The absence of these viruses in effluent samples can be due to

277  the sanitation treatments applied in WWTPs, even though further analysis that includes

278  a wider sampling design needs to be performed. These results are in line with previous

279  studies showing a high abundance of bacteriophages families (Aw et al., 2014;

280  Cantalupo et al., 2011; Fernandez-Cassi et al., 2018; Rusiñol et al., 2020; Wang et al.,

281  2018) in influent sewage samples. Nevertheless, other studies showed *Virgaviridae* as

282  the most represented viruses (Furtak et al., 2016). Differences in virome profiling with

283  other studies might be due to the influence of library sequencing and the intrinsic

284  characteristics of the virome related to the sample itself and the area of study. On the

285  other hand, the higher presence of some viruses or even its detection only in effluent

286  samples could be produced by the decrease of other viruses that allowed its detection.

287  The presence of pathogenic viruses is an important aspect for defining the final use of

288  treated waters as it may be the case of irrigation. Due to their high environmental

289  resistance, the presence of human enteric viruses has been reported in treated

290  wastewaters (Adriaenssens et al., 2018). However, some of these pathogenic viruses are

291  not always detected by metagenomics analyses. For instance, in the study by Fernández-

292  Cassi et al. (2018), human adenoviruses (HAdV) reads were not detected in samples

293  concentrated from 10 liters of wastewater. *Adenoviridae* was also not detected in the

294  study of Adrianssens et al. (2018), in which the sample was concentrated from 1 liter of

295  wastewater. In our study, concentrating 200 mL of effluent samples, we were able to

296  detect HAdV in percentages between 0.16% and 0.35%. In contrast, percentages of

297  HAdV in influent wastewaters were lower than 0.01%. Overall, the majority of the

298    annotated virome belonged to bacteriophages. This indicates that metagenomics is poor

299    in sensitivity when used to detect a low abundance of viral pathogens against a large

300    background of bacteriophages, as occurred for the enteric viruses detected by viability

301    RT-qPCR. For example, in the present study, norovirus genomes could not be retrieved

302    from the reads as reported elsewhere (Adriaenssens et al., 2018; Fernández-Cassi et al.,

303    2018; Strubbia et al., 2019b). In the current study, the number of generated paired reads

304    per sample was 3.2 and 11.5 million for LS and NB, respectively; while Adriaenssens et

305    al., (2018) reported between 10 and 110 million, increasing significantly the probability

306    to retrieve full or partial viral genomes. Alternatively, methods to detect and

307    characterize specific viruses have been described and rely on the selection of target

308    RNA prior to library preparation through a capture using VirCapSeq-VERT target

309    enrichment, as reported for norovirus (Strubbia et al., 2019b).

310    **5. Conclusion**

311    The use of metagenomics for virome characterization and its implementation for

312    wastewater analyses has increased in the last years ( Nieuwenhuijse and Koopmans,

313    2017). However, the major problem of this approach is the lack of standardized

314    procedures and the substantial differences among studies; thus, available data must be

315    interpreted with caution. The present study showed a procedure that allows the detection

316    and the characterization of viral populations in untreated and treated wastewater

317    samples. Overall, this study sheds light on the diversity of the viral communities in

318    untreated and treated wastewaters providing valuable information also in terms of viral

319    fecal indicators. The study also evidences the bias on virome profiles obtained by tested

320    sequencing libraries. Our results underline the need for further studies to elucidate the

321    influence of sequencing procedures in virome profiles in wastewater matrices in order to

322    improve the knowledge of the virome in the environment.

**Declaration of competing interest**

**Acknowledgements**

**Contributions**

GS designed the work. AP-C and EC-F processed the samples. AP-C performed the bioinformatic work and data analysis. AP-C, EC-F, WR and GS wrote the paper. All authors have read and agreed to the published version of the manuscript.

**References**

AAVV (2018). "Section 9510 D. Virus concentration by aluminum hydroxide adsorption-precipitation, chapter detection of enteric viruses," in Standard Methods for the Examination of Water and Wastewater, 23rd Edn, eds

346    E. W. Rice, R. B. Baird, and A. D. Eaton (Denver, CO: American Water Works

347    Association)

348    Adriaenssens, E.M., Farkas, K., Harrison, C., Jones, D.L., Allison, H.E., McCarthy,

349        A.J., 2018. Viromic Analysis of Wastewater Input to a River Catchment Reveals a

350        Diverse Assemblage of RNA Viruses. mSystems 3, 1–18.

351        https://doi.org/10.1128/mSystems.00025-18

352    Aw, T.G., Howe, A., Rose, J.B., 2014. Metagenomic approaches for direct and cell

353        culture evaluation of the virological quality of wastewater. J. Virol. Methods 210,

354        15–21. https://doi.org/10.1016/j.jviromet.2014.09.017

355    Baselga, A. 2010. Partitioning the turnover and nestedness components of beta

356        diversity. Global Ecology and Biogeography 19:134-143.

357        https://doi.org/10.1111/j.1466-8238.2009.00490.x

358    Boisvert, S., Raymond, F., Godzaridis, É., Laviolette, F., Corbeil, J., 2012. Ray Meta:

359        scalable de novo metagenome assembly and profiling. Genome Biol. 13, R122.

360        https://doi.org/10.1186/gb-2012-13-12-r122

361    Boratyn, G.M., Camacho, C., Cooper, P.S., Coulouris, G., Fong, A., Ma, N., Madden,

362        T.L., Matten, W.T., McGinnis, S.D., Merezhuk, Y., Raytselis, Y., Sayers, E.W.,

363        Tao, T., Ye, J., Zaretskaya, I., 2013. BLAST: a more efficient report with usability

364        improvements. Nucleic Acids Res. 41, W29-33. https://doi.org/10.1093/nar/gkt282

365    Bivins, A., North, D., Ahmad, A., Ahmed, W., Alm, E., Been, F., Bhattacharya, P.,

366        Bijlsma, L., Boehm, A.B., Brown, J., Buttiglieri, G., Calabro, V., Carducci, A.,

367        Castiglioni, S., Cetecioglu Gurol, Z., Chakraborty, S., Costa, F., Curcio, S., De Los

368        Reyes, F.L., Delgado Vela, J., Farkas, K., Fernandez-Casi, X., Gerba, C., Gerrity,

369    D., Girones, R., Gonzalez, R., Haramoto, E., Harris, A., Holden, P.A., Islam, M.T.,

370    Jones, D.L., Kasprzyk-Hordern, B., Kitajima, M., Kotlarz, N., Kumar, M., Kuroda,

371    K., La Rosa, G., Malpei, F., Mautus, M., McLellan, S.L., Medema, G., Meschke,

372    J.S., Mueller, J., Newton, R.J., Nilsson, D., Noble, R.T., Van Nuijs, A., Peccia, J.,

373    Perkins, T.A., Pickering, A.J., Rose, J., Sanchez, G., Smith, A., Stadler, L.,

374    Stauber, C., Thomas, K., Van Der Voorn, T., Wigginton, K., Zhu, K., Bibby, K.,

375    2020. Wastewater-Based Epidemiology: Global Collaborative to Maximize

376    Contributions in the Fight against COVID-19. Environ. Sci. Technol.

377    https://doi.org/10.1021/acs.est.0c02388

378  Cantalupo, P.G., Calgua, B., Zhao, G., Hundesa, A., Wier, A.D., Katz, J.P., Grabe, M.,

379    2011. Raw Sewage Harbors Diverse Viral Populations. MBio 2, e00180-11.

380    https://doi.org/10.1128/mBio.00180-11.Editor

381  Chalmers, R.M., Robinson, G., Elwin, K., Hadfield, S.J., Thomas, E., Watkins, J.,

382    Casemore, D., Kay, D., 2010. Detection of Cryptosporidium species and sources of

383    contamination with Cryptosporidium hominis during a waterborne outbreak in

384    north west Wales. J. Water Health 8, 311–25. https://doi.org/10.2166/wh.2009.185

385  Conceição-Neto, N., Zeller, M., Lefrère, H., De Bruyn, P., Beller, L., Deboutte, W.,

386    Yinda, C.K., Lavigne, R., Maes, P., Van Ranst, M., Heylen, E., Matthijnssens, J.,

387    2015. Modular approach to customise sample preparation procedures for viral

388    metagenomics: a reproducible protocol for virome analysis. Sci. Rep. 5, 16532.

389    https://doi.org/10.1038/srep16532

390  Cuevas-Ferrando, E., Randazzo, W., Pérez-Cataluña, A., Sánchez, G., 2020. HEV

391    Occurrence in Waste and Drinking Water Treatment Plants. Front. Microbiol. 10.

392    https://doi.org/10.3389/fmicb.2019.02937

393    Eslamian, S. (Ed.), 2016. Urban Water Reuse Handbook. CRC Press.

394        https://doi.org/10.1201/b19646

395    Farkas, K., Walker, D.I., Adriaenssens, E.M., McDonald, J.E., Hillary, L.S., Malham,

396        S.K., Jones, D.L., 2020. Viral indicators for tracking domestic wastewater

397        contamination in the aquatic environment. Water Res.

398        https://doi.org/10.1016/j.watres.2020.115926

399    Fernandez-Cassi, X., Timoneda, N., Martínez-Puchol, S., Rusiñol, M., Rodriguez-

400        Manzano, J., Figuerola, N., Bofill-Mas, S., Abril, J.F., Girones, R., 2018.

401        Metagenomics for the study of viruses in urban sewage as a tool for public health

402        surveillance. Sci. Total Environ. 618, 870–880.

403        https://doi.org/10.1016/j.scitotenv.2017.08.249

404    Fong, T.-T.,  Lipp, E. K. 2005. Enteric viruses of humans and animals in aquatic

405        environments: health risks, detection, and potential water quality assessment tools.

406        MMBR, 69(2), 357–371. https://doi.org/10.1128/MMBR.69.2.357-371.2005

407    Furtak, V., Roivainen, M., Mirochnichenko, O., Zagorodnyaya, T., Laassri, M., Zaidi,

408        S.Z., Rehman, L., Alam, M.M., Chizhikov, V., Chumakov, K., 2016.

409        Environmental surveillance of viruses by tangential flow filtration and

410        metagenomic reconstruction. Eurosurveillance 21, 30193.

411        https://doi.org/10.2807/1560-7917.ES.2016.21.15.30193

412    Gerba, C.P., Betancourt, W.Q., Kitajima, M., Rock, C.M., 2018. Reducing uncertainty

413        in estimating virus reduction by advanced water treatment processes. Water Res.

414        133, 282–288. https://doi.org/10.1016/j.watres.2018.01.044

415    Gerba, C.P., Kitajima, M., Iker, B., 2013. Viral presence in waste water and sewage and

416    control methods. Viruses Food Water 293–315.

417    https://doi.org/10.1533/9780857098870.3.293

418    Haramoto, E., Kitajima, M., Hata, A., Torrey, J.R., Masago, Y., Sano, D., Katayama,

419    H., 2018. A review on recent progress in the detection methods and prevalence of

420    human enteric viruses in water. Water Res. 135, 168–186.

421    https://doi.org/10.1016/j.watres.2018.02.004

422    ISO 9308-2:2012. Water quality -- Enumeration of Escherichia coli and coliform

423    bacteria -- Part 2: Most probable number method.

424    Kitajima, M., Iker, B.C., Pepper, I.L., Gerba, C.P., 2014. Relative abundance and

425    treatment reduction of viruses during wastewater treatment processes —

426    Identification of potential viral indicators. Sci. Total Environ. 488–489, 290–296.

427    https://doi.org/10.1016/J.SCITOTENV.2014.04.087

428    Lamori, J.G., Xue, J., Rachmadi, A.T., Lopez, G.U., Kitajima, M., Gerba, C.P., Pepper,

429    I.L., Brooks, J.P., Sherchan, S., 2019. Removal of fecal indicator bacteria and

430    antibiotic resistant genes in constructed wetlands. Environ. Sci. Pollut. Res. 26,

431    10188–10197. https://doi.org/10.1007/s11356-019-04468-9

432    Magoč, T., Salzberg, S.L., 2011. FLASH: fast length adjustment of short reads to

433    improve genome assemblies. Bioinformatics 27, 2957–63.

434    https://doi.org/10.1093/bioinformatics/btr507

435    Martin, M., 2011. Cutadapt removes adapter sequences from high-throughput

436    sequencing reads. EMBnet.journal 17, 10. https://doi.org/10.14806/ej.17.1.200

437    Nieuwenhuijse, D.F., Koopmans, M.P.G., 2017. Metagenomic sequencing for

438    surveillance of food- and waterborne viral diseases. Front. Microbiol. 8, 1–11.

439 https://doi.org/10.3389/fmicb.2017.00230

440 Nieuwenhuijse, D.F., Oude Munnink, B.B., Phan, M.V.T., Hendriksen, R.S., Bego, A.,

441  Rees, C., Neilson, E.H., Coventry, K., Collignon, P., Allerberger, F., Rahube, T.O.,

442  Oliveira, G., Ivanov, I., Sopheak, T., Vuthy, Y., Yost, C.K., Tabo, D. adjim,

443  Cuadros-Orellana, S., Ke, C., Zheng, H., Baisheng, L., Jiao, X., Donado-Godoy,

444  P., Coulibaly, K.J., Hrenovic, J., Jergović, M., Karpíšková, R., Elsborg, B.,

445  Legesse, M., Eguale, T., Heikinheimo, A., Villacis, J.E., Sanneh, B., Malania, L.,

446  Nitsche, A., Brinkmann, A., Saba, C.K.S., Kocsis, B., Solymosi, N.,

447  Thorsteinsdottir, T.R., Hatha, A.M., Alebouyeh, M., Morris, D., O'Connor, L.,

448  Cormican, M., Moran-Gilad, J., Battisti, A., Alba, P., Shakenova, Z., Kiiyukia, C.,

449  Ng'eno, E., Raka, L., Bērziņš, A., Avsejenko, J., Bartkevics, V., Penny, C.,

450  Rajandas, H., Parimannan, S., Haber, M.V., Pal, P., Schmitt, H., van Passel, M.,

451  van de Schans, M.G.M., Zuidema, T., Jeunen, G.J., Gemmell, N., Fashae, K.,

452  Wester, A.L., Holmstad, R., Hasan, R., Shakoor, S., Rojas, M.L.Z., Wasyl, D.,

453  Bosevska, G., Kochubovski, M., Radu, C., Gassama†, A., Radosavljevic, V., Tay,

454  M.Y.F., Zuniga-Montanez, R., Wuertz, S., Gavačová, D., Trkov, M., Keddy, K.,

455  Esterhuyse, K., Cerdà-Cuéllar, M., Pathirage, S., Larsson, D.G.J., Norrgren, L.,

456  Örn, S., Van der Heijden, T., Kumburu, H.H., de RodaHusman, A.M., Njanpop-

457  Lafourcade, B.M., Bidjada, P., Nikiema-Pessinaba, S.C., Levent, B., Meschke,

458  J.S., Beck, N.K., Van Dang, C., Tran, D.M.N., Do Phuc, N., Kwenda, G., Munk,

459  P., Venkatakrishnan, S., Aarestrup, F.M., Cotten, M., Koopmans, M.P.G., 2020.

460  Setting a baseline for global urban virome surveillance in sewage. Sci. Rep. 10, 1–

461  13. https://doi.org/10.1038/s41598-020-69869-0

462 Polo, D.D., Quintela-Baluja, M., Corbishley, A., Jones, D.L., Singer, A.C., Graham,

463  D.W., Romalde, P.J.L., 2020. Making waves: Wastewater-based epidemiology for

464      SARS-CoV-2 – Developing robust approaches for surveillance and prediction is

465      harder than it looks. Water Res. 186, 116404.

466      https://doi.org/10.1016/j.watres.2020.116404

467  Randazzo, W., Piqueras, J., Evtoski, Z., Sastre, G., Sancho, R., Gonzalez, C., Sánchez,

468      G., 2019. Interlaboratory Comparative Study to Detect Potentially Infectious

469      Human Enteric Viruses in Influent and Effluent Waters. Food Environ. Virol.

470      https://doi.org/10.1007/s12560-019-09392-22

471   Rusiñol, M., Martínez-Puchol, S., Timoneda, N., Fernández-Cassi, X., Pérez-Cataluña,

472      A., Fernández-Bravo, A., Moreno-Mesonero, L., Moreno, Y., Alonso, J.L.,

473      Figueras, M.J., Abril, J.F., Bofill-Mas, S., Girones, R., 2020. Metagenomic

474      analysis of viruses, bacteria and protozoa in irrigation water. Int. J. Hyg. Environ.

475      Health 224, 113440. https://doi.org/10.1016/j.ijheh.2019.113440

476  Sano, D., Amarasiri, M., Hata, A., Watanabe, T., Katayama, H., 2016. Risk

477      management of viral infectious diseases in wastewater reclamation and reuse:

478      Review. Environ. Int. 91, 220–229.

479      https://doi.org/10.1016/J.ENVINT.2016.03.001

480  Stachler, E., Kelty, C., Sivaganesan, M., Li, X., Bibby, K., Shanks, O.C., 2017.

481      Quantitative CrAssphage PCR Assays for Human Fecal Pollution Measurement.

482      Environ. Sci. Technol. 51, 9146–9154. https://doi.org/10.1021/acs.est.7b02703

483  Strubbia, S., Phan, M.V.T., Schaeffer, J., Koopmans, M., Cotten, M., Le Guyader, F.S.,

484      2019a. Characterization of Norovirus and Other Human Enteric Viruses in Sewage

485      and Stool Samples Through Next-Generation Sequencing. Food Environ. Virol.

486      https://doi.org/10.1007/s12560-019-09402-3

487  Strubbia, S., Schaeffer, J., Oude Munnink, B.B., Besnard, A., Phan, M.V.T.,

488      Nieuwenhuijse, D.F., de Graaf, M., Schapendonk, C.M.E., Wacrenier, C., Cotten,

489      M., Koopmans, M.P.G., Le Guyader, F.S., 2019b. Metavirome Sequencing to

490      Evaluate Norovirus Diversity in Sewage and Related Bioaccumulated Oysters.

491      Front. Microbiol. 10, 2394. https://doi.org/10.3389/fmicb.2019.02394

492  Wang, Y., Jiang, X., Liu, L., Li, B., Zhang, T., 2018. High-Resolution Temporal and

493      Spatial Patterns of Virome in Wastewater Treatment Systems. Environ. Sci.

494      Technol. 52, 10337–10346. https://doi.org/10.1021/acs.est.8b03446

495

496

497

498

499

500

501

502

503

504

505

506

507

**Figure legends**

509    Figure 1: Shannon and Simpson diversity indexes for viral species in influent (IW) and

510    effluent (EW) samples processed by using ScriptSeq v2 RNA-Seq Library Preparation

511    kit (LS) and NEBNext® Ultra™ II RNA Library Prep Kit libraries (NB) for

512    metagenomics characterization.

513    Figure 2. Relative abundance at family level of the viral population detected in influent

514    and effluent samples from four different WWTPs by metagenomics with ScriptSeq v2

515    RNA-Seq Library Preparation kit (LS, green) and NEBNext® Ultra™ II RNA Library

516    Prep Kit (NB, orange).

517    Figure 3. Viral indicators analysis in influent (IW) and effluent (EW) samples. Panel A,

518    Detection of viral indicators with ScriptSeq v2 RNA-Seq Library Preparation kit (LS)

519    and NEBNext® Ultra™ II RNA Library Prep Kit (NB and NB-corrected). Panel B,

520    Correlation matrix between the reads of viral indicators obtained by NGS and the load

521    of enteric viruses (RT-qPCR) and E. coli counts.

522    Figure 4. Heatmap showing the virome composition at species level obtained by

523    merging the results of ScriptSeq v2 RNA-Seq Library Preparation kit (LS) and

524    NEBNext® Ultra™ II RNA Library Prep Kit (NB). Only species with percentages

525    higher than 1% are shown.

526