# Viromic Analysis of Wastewater Input to a River Catchment Reveals a Diverse Assemblage of RNA Viruses

Evelien M. Adriaenssens,[a] Kata Farkas,[b] Christian Harrison,[a] David L. Jones,[b] Heather E. Allison,[a] Alan J. McCarthy[a]

[a]Microbiology Research Group, Institute of Integrative Biology, University of Liverpool, Liverpool, United Kingdom
[b]School of Environment, Natural Resources and Geography, Bangor University, Bangor, United Kingdom

**ABSTRACT** Detection of viruses in the environment is heavily dependent on PCR-based approaches that require reference sequences for primer design. While this strategy can accurately detect known viruses, it will not find novel genotypes or emerging and invasive viral species. In this study, we investigated the use of viromics, i.e., high-throughput sequencing of the biosphere's viral fraction, to detect human-/animal-pathogenic RNA viruses in the Conwy river catchment area in Wales, United Kingdom. Using a combination of filtering and nuclease treatment, we extracted the viral fraction from wastewater and estuarine river water and sediment, followed by high-throughput RNA sequencing (RNA-Seq) analysis on the Illumina HiSeq platform, for the discovery of RNA virus genomes. We found a higher richness of RNA viruses in wastewater samples than in river water and sediment, and we assembled a complete norovirus genotype GI.2 genome from wastewater effluent, which was not contemporaneously detected by conventional reverse transcription-quantitative PCR (qRT-PCR). The simultaneous presence of diverse rotavirus signatures in wastewater indicated the potential for zoonotic infections in the area and suggested runoff from pig farms as a possible origin of these viruses. Our results show that viromics can be an important tool in the discovery of pathogenic viruses in the environment and can be used to inform and optimize reference-based detection methods provided appropriate and rigorous controls are included.

**IMPORTANCE** Enteric viruses cause gastrointestinal illness and are commonly transmitted through the fecal-oral route. When wastewater is released into river systems, these viruses can contaminate the environment. Our results show that we can use viromics to find the range of potentially pathogenic viruses that are present in the environment and identify prevalent genotypes. The ultimate goal is to trace the fate of these pathogenic viruses from origin to the point where they are a threat to human health, informing reference-based detection methods and water quality management.

**KEYWORDS** RNA viruses, norovirus, pathogen detection, rotavirus, viromics, wastewater

Pathogenic viruses in water sources are likely to originate primarily from contamination with sewage. Classic marker bacteria used for fecal contamination monitoring, such as *Escherichia coli* and *Enterococcus* spp., are not, however, good indicators for the presence of human enteric viruses (1). The virus component is often monitored using quantitative PCR (qPCR) approaches, which can give information on the abundance of specific viruses and their genotypes, but only those that are both known and characterized (2). Viruses commonly targeted in sewage contamination assays include noroviruses (NoV) (3), hepatitis viruses (4), enteroviruses (5), and various adenoviruses (6, 7). Viral monitoring in sewage has previously yielded positive results for norovirus,

**FIG 1** Map of the sampling locations, indicated with blue arrows. WWTP, wastewater treatment plant. Data in the left panel were taken from Google Maps (Map data ©Google 2017).

sapovirus (SaV), astrovirus, and adenovirus, indicating that people are shedding viruses that are not necessarily detected in a clinical setting (8). This same study found a spike in norovirus genogroup GII sequence signatures in sewage 2 to 3 weeks before the outbreak of associated disease was reported in hospitals and nursing homes. The suggestion, therefore, is that environmental viromics can provide an early warning of disease outbreaks, in addition to the monitoring of virus dissemination in watercourses.

Recent reviews have proposed the use of viral metagenomics or viromic approaches as an alternative method to test for the presence of pathogenic viruses in the environment, offering the potential to detect novel genotypes or even entirely novel viruses (2, 9, 10). Potential new viral markers for fecal contamination have already been revealed, such as pepper mild mottle virus and crAssphage (11, 12), among the huge diversity of human viruses found in sludge samples (13–16).

In this pilot study, we have used viromics to investigate the presence of human-pathogenic RNA viruses in wastewater and estuarine surface water and sediment in a single catchment. The water and sediment samples were collected at the wastewater treatment plant (Llanrwst, Wales, United Kingdom) and downstream from it at the estuary of the river Conwy near a bathing water beach (Morfa, Wales, United Kingdom) (Fig. 1). To our knowledge, this is the first study to use unamplified environmental viral RNA for sequencing library construction, sequence data set production, and subsequent analysis. Because we used a directional library sequencing protocol on RNA, rather than amplifying to cDNA, we were able to distinguish single-stranded from double-stranded RNA genome fragments.

**TABLE 1** Summary of viromic and qRT-PCR detection of specific RNA viruses across sewage, estuarine water, and sediment samples

| Sample[a] | Sample vol or mass | Location | No. of contigs (curated) | Target RNA viruses detected in contigs[b] | qRT-PCR results (gc/liter)[c] |
|---|---|---|---|---|---|
| LI_13-9 | 1 liter | Llanrwst WWTP[d] | 5,721 | RVA, RVC, PBV, SaV | NoVGII (1,457) |
| LE_13-9 | 1 liter | Llanrwst WWTP | 2,201 | RVA, RVC, PBV | NoVGII (1,251) |
| LI_11-10 | 1 liter | Llanrwst WWTP | 859 | PBV | NoVGII (detected) |
| LE_11-10 | 1 liter | Llanrwst WWTP | 5,433 | NoVGI, RVA, RVC, PBV, AsV | NoVGII (50,180) |
| SW | 50 liters | Morfa beach | 243 | | |
| Sed1 | 60 g | Morfa beach | 550[e] | | |
| Sed2 | 60 g | Morfa beach | 550[e] | | |

[a]LI, sewage influent; LE, sewage effluent; SW, estuarine surface water; Sed, estuarine sediment.
[b]RVA, rotavirus A; RVB, rotavirus B; PBV, picobirnavirus; SaV, sapovirus; NoVGI, norovirus genogroup I; AsV, astrovirus.
[c]Samples were tested with qRT-PCR for the following targets: NoVGI, NoVGII, SaV, HAV, and HEV. Results are reported in genome copies per liter (gc/liter). NoVGII below the limit of quantification (approximately 200 gc/liter) was detected in sample LI_11-10. NoVGII was the only target virus detected by qRT-PCR.
[d]WWTP, wastewater treatment plant.
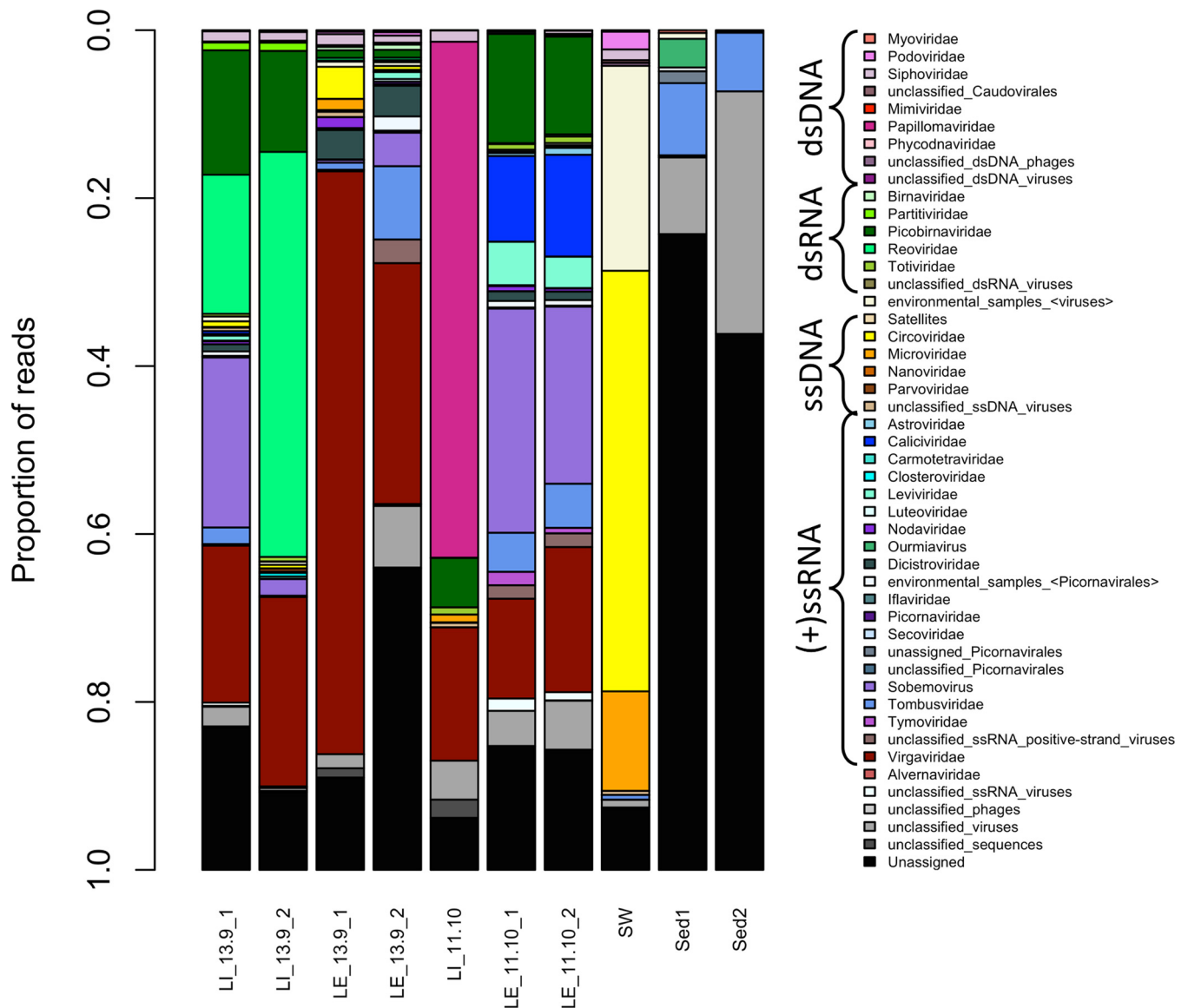[e]Samples Sed1 and Sed2 were assembled together into the contig data set Sed.

## RESULTS

**Sample overview.** Wastewater influent and effluent samples were collected from the Llanrwst wastewater treatment plant (53°08′24.4″N, 3°48′12.8″W) (Fig. 1) in September and October 2016, resulting in four different samples, LI_13-9 (Llanrwst influent September 2016), LE_13-9 (Llanrwst effluent September 2016), LI_11-10 (Llanrwst influent October 2016), and LE_11-10 (Llanrwst effluent October 2016). Estuarine surface water (SW) was collected from Morfa beach (53°17′37.7″N, 3°50′22.2″W; Conwy, Wales) (Fig. 1) in November 2016 and sediment from the same site in October and November 2016 (Sed1 and Sed2, respectively).

As an initial assessment, samples were tested for the presence of a subset of locally occurring enteric RNA viruses using reverse transcription-quantitative PCR (qRT-PCR) (Table 1). Only norovirus (NoV) genogroup GII signatures were detected in the wastewater samples. In the samples collected in September 2016, $10^3$ genome copies (gc)/liter of norovirus GII was observed in both the influent (LI_13-9) and the effluent (LE_13-9). In the samples collected in October 2016, approximately $10^2$ gc/liter (below the limit of quantification, which was approximately 200 gc/liter) was observed in the influent (LI_11-10) and a considerably higher concentration of $5 \times 10^4$ gc/liter was noted in the effluent (LE_11-10). All qRT-PCRs were negative for the presence of sapoviruses (SaV) and hepatitis A and E viruses (HAV and HEV, respectively). None of the target enteric viruses were found in the surface water and sediment samples.

**Summary of viral diversity.** The virus taxonomic diversity present in each sample was assessed by comparison of curated read and contig data sets with both the RefSeq viral protein database and the nonredundant protein database of NCBI, using Diamond blastx (17) and lowest-common-ancestor taxon assignment with MEGAN6 (18). For wastewater samples LI_13-9, LE_13-9, and LE_11-10, two libraries were processed (indicated with _1 and _2 in the data set names), and for the wastewater influent sample LI_11-10, the surface water sample (SW), and the two sediment samples (Sed1 and Sed2), one library was processed for each. This section focuses on those reads and contigs that have been assigned to the viral fraction exclusively, disregarding sequences of cellular or unknown origin.

The wastewater samples showed a greater richness of known viruses and had a larger number of curated contigs than the surface water and sediment samples (Fig. 2 and 3). At the virus family level, between 14 and 34 groups, including the unclassified levels, were observed for wastewater influent and effluent samples, 12 for the surface estuarine water sample, and 11 and 5 for the sediment samples Sed1 and Sed2, respectively. The unclassified viruses and unassigned bins are indicated in gray and black in Fig. 2 and made up the majority of known reads in the estuarine sediment samples. In most of the viromes, double-stranded DNA (dsDNA) and single-stranded DNA (ssDNA) virus families were present, despite a DNase treatment having been performed after viral nucleic acid extraction (Fig. 2 and 3). These families represented

**FIG 2** Taxonomic distribution of curated read data (relative abundance) at the virus family level. Reads were assigned to a family or equivalent group by MEGAN6 using a lowest-common-ancestor algorithm, based on blastx-based homology using the program Diamond with the RefSeq Viral protein database (January 2017 version) and the nonredundant protein database (May 2017 version). Only viral groupings are shown. LI, sewage influent; LE, sewage effluent; SW, estuarine surface water; Sed, estuarine sediment.

only a minor (<5%) proportion of the total assigned reads, with a few exceptions. In wastewater influent sample LI_11-10, reads assigned to the dsDNA family *Papillomaviridae* accounted for 61% of the total (Fig. 2, dark pink), and these reads were assembled into a single contig representing a nearly complete *Betapapillomavirus* genome. In the surface water sample, reads assigned to the ssDNA families *Circoviridae* and *Microviridae* represented 50% and 12% of the total, respectively (Fig. 2, yellow and orange), assembling into contigs representing a significant proportion of the genome. The presence of both ssDNA and dsDNA virus signatures in all data sets is most likely due to incomplete digestion of the viral DNA with the DNase Max kit.

The families of dsRNA viruses present in these data sets were *Totiviridae* (fungal and protist hosts), *Reoviridae* (invertebrate, vertebrate, and plant hosts), *Picobirnaviridae* (mammals), *Partitiviridae* (fungi and protists), and *Birnaviridae* (vertebrates and invertebrates), with a small number of reads and contigs recognized as unclassified dsRNA viruses (Fig. 2 and 3). None of these groups were present in all libraries, but totivirus

**FIG 3** Heatmap of viral richness at the family level per sample. Heatmap colors denote relative abundances per sample. Contigs larger than 300 nucleotides (nt) were assigned to a family or grouping by MEGAN6 using a lowest-common-ancestor algorithm, based on blastx-based homology using the program Diamond with the RefSeq viral protein database (version January 2017) and the nonredundant protein database (May 2017 version). Only those families/groups comprising large contigs (>1,000 nt) or with contigs mapping to viral signature genes (e.g., capsid and RNA-dependent RNA polymerase genes) were retained. LI, sewage influent; LE, sewage effluent; SW, estuarine surface water; Sed, estuarine sediment.

and picobirnavirus (Fig. 2, dark green) signatures were present in all wastewater samples and reoviruses (Fig. 2, bright green) were found in three of the four wastewater samples. *Partitiviridae* signatures were only found in wastewater samples LE_11-10 and LI_13-9, while *Birnaviridae* reads were only present in the wastewater LE_13-9 libraries. The sediment and surface water samples did not have detectable levels of dsRNA virus sequences.

Positive-sense ssRNA viruses were the most diverse class of viruses present in these data sets. The family *Tombusviridae*, which groups plant viruses with monopartite or bipartite linear genomes (19), was present in all samples with the sole exception of the wastewater influent sample LI_11-10 (Fig. 2, cornflower blue, and Fig. 3). Virus signatures belonging to the family *Virgaviridae*, representing plant viruses, were present in all wastewater samples at high relative abundances (Fig. 2, dark red). Other highly represented families or groupings were the families *Dicistroviridae* (invertebrate hosts) and *Nodaviridae* (invertebrate and vertebrate hosts) and the bacteriophage family *Leviviridae*, the plant virus genus *Sobemovirus* (Fig. 2, medium purple), and the groupings of "unclassified ssRNA positive-strand viruses" and several unclassified/unassigned/ environmental members of the order *Picornavirales*. Sediment sample Sed1 was the only sample with signatures of the family *Alvernaviridae*, which has as its sole member the dinoflagellate virus *Heterocapsa circularisquama RNA virus 01*. The wastewater effluent sample LE_11-10 and influent sample LI_13-9 were the only samples with

*Calicivirus* signatures (Fig. 2, medium blue), and sample LE_11-10_1 and LE_1-10_2 were the only samples with *Astroviridae* reads (vertebrate host). Several families of the order *Picornavirales* were detected in the wastewater samples at different levels in different samples, and a small number of unassigned picornaviruses were detected in the surface water sample (SW).

We did not observe any known negative-sense ssRNA [(−)ssRNA] viruses in any of the sequencing libraries, but it is possible that some of the unaffiliated viral contigs belong to this class. The known human-pathogenic (−)ssRNA viruses are enveloped (19) and predicted to degrade more rapidly than the nonenveloped enteric viruses, especially in wastewater (20, 21). We cannot rule out the possibility that (−)ssRNA viruses were present but were removed by our sampling protocol.
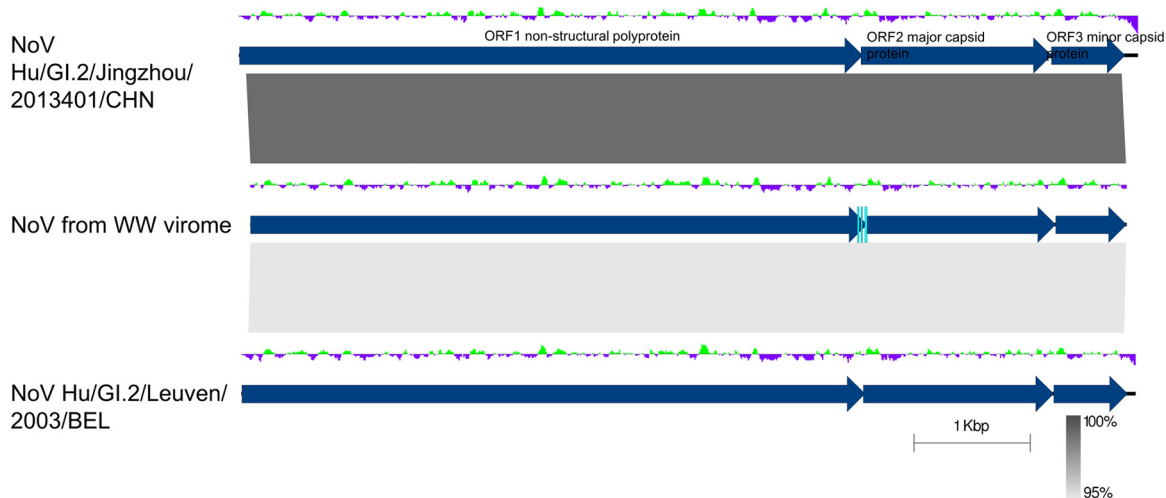
The general wastewater viral diversity found here is similar to that reported previously. Those studies that investigated RNA viruses found both bacterial and eukaryotic viruses, with a high abundance of plant viruses of the family *Virgaviridae*, which includes the *Tobamovirus Pepper mild mottle virus* (11, 14). The families of viruses with potential human hosts found in previous metagenomics studies of sewage include *Astroviridae*, *Caliciviridae*, *Picobirnaviridae*, and *Picornaviridae* (13–16), of which only picobirnaviruses were recovered in all wastewater viromes in this study. In contrast, members of the family *Reoviridae*, represented by the genus *Rotavirus*, were found in three of our four wastewater samples but were not detected in many of the previous studies (14–16).

**Potential human-pathogenic viruses.** An important aim of this study was to investigate the presence and genomic diversity of potential human-pathogenic RNA viruses in different sample types within the river catchment area. To minimize misassignments of short sequences to taxa, we used the assembled, curated contig data set and looked for contigs representing nearly complete viral genomes.

**Presence of a norovirus GI.2 genome.** We were particularly interested in finding norovirus genomes in order to explore the genomic diversity of these important and potentially abundant pathogens originating from sewage and disseminated in watercourses, with implications for shellfisheries and recreational waters. This is of relevance due to known issues of sewage contamination in the region (22). Members of the genus *Norovirus* (family *Caliciviridae*) are nonenveloped, icosahedral (+)ssRNA viruses with a linear, unsegmented, ~7.6-kb genome encoding three open reading frames (ORFs) (19). These viruses are divided into different genogroups, of which GI and GII are associated with human gastroenteritis (23, 24). Noroviruses are identified routinely by qRT-PCR, providing an opportunity here to examine correlations between qRT-PCR and metaviromic data.

We only found norovirus signatures in the libraries of wastewater effluent sample LE_11-10. These reads assembled into a single contig of 7,542 bases, representing a nearly complete norovirus genome (GenBank accession number MG599789). Read mapping showed uneven coverage over the genome length between 18× and 745× (13,165 reads of library 1 and 8,986 reads of library 2). Based on this mapping, we performed variant calling and corrected the consensus sequence in cases where the variant was present in more than 85% of the reads.

A BLASTN search revealed two close relatives to our wastewater-associated norovirus genome, norovirus Hu/GI.2/Jingzhou/2013401/CHN (KF306212), which is 7,740 bases in length (25), displaying a nucleotide sequence identity of 99% over 99% of the genome length, and norovirus Hu/GI.2/Leuven/2003/BEL (FJ515294), at 95% sequence identity over 99% of the alignment length (Fig. 4). From the 5′ end of our norovirus contig, 62 bases were missing compared with the sequence of Hu/GI.2/Jingzhou/2013401/CHN, and from the 3′ end, 165 bases and the poly(A) tail were not present. We compared the sequence of our norovirus with that of Hu/GI.2/Jingzhou/2013401/CHN base by base and observed 81 single-nucleotide polymorphisms (SNPs) and no other forms of variation. Of the SNPs, only eight were nonsynonymous, resulting in five different amino acids incorporated in the nonstructural polyprotein (ORF1), one in the

**FIG 4** Pairwise genome comparison between the virome's norovirus genome (middle) and its closest relatives, norovirus Hu/GI.2/Jingzhou/2013401/CHN and norovirus Hu/GI.2/Leuven/2003/BEL. BLASTN similarity is indicated in shades of gray. ORFs are delineated by dark-blue arrows. Deviations from the average GC content are indicated above the genomes in a green and purple graph. The qRT-PCR primer binding sites for the wastewater (WW)-associated genome are indicated by light-blue rectangles. The figure was created with Easyfig (92).
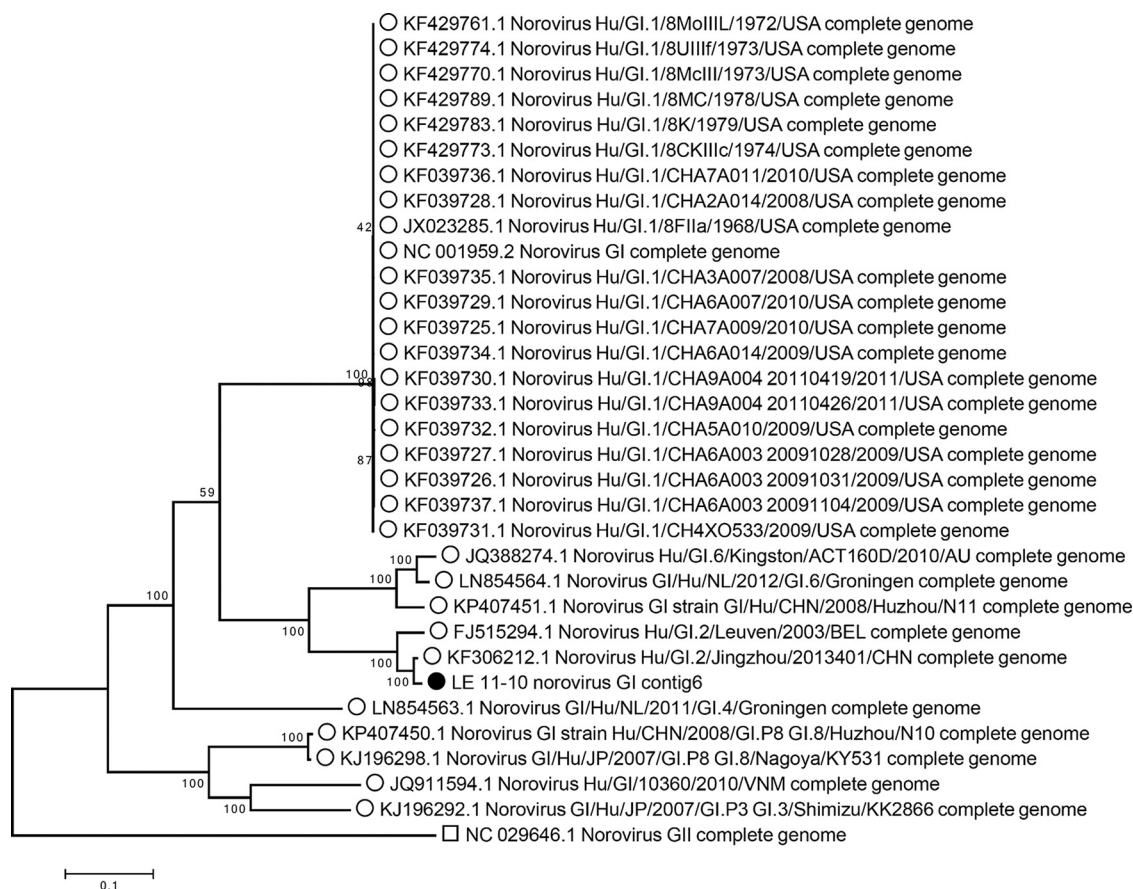
major capsid protein (MCP) (ORF2), and two in the minor structural protein (ORF3). According to the current classification criteria, this level of similarity places our assembled genome in genogroup GI, genotype GI.2, with only a single amino acid differing between the major capsid protein of Hu/GI.2/Jingzhou/2013401/CHN and the genome assembled here.

We tested the genotype grouping of our genome in a whole-genome phylogeny with all complete genome sequences of genogroup I available in GenBank. The phylogenomic tree clearly delineated the different genotypes within genogroup GI, placing the newly assembled genome within genotype GI.2, with the reference isolate for GII used as an outgroup (Fig. 5).

For further validation, the full genome of the novel norovirus GI was recovered using RT-PCR. However, the amplicon could not be ligated into a plasmid and, hence, was not fully sequenced.

**Presence of diverse rotavirus segments in wastewater samples.** Rotaviruses are segmented dsRNA viruses belonging to the family *Reoviridae* that cause gastroenteric illness in vertebrates and are transmitted through the fecal-oral route (19). Read signatures assigned to the genus *Rotavirus* were found in three of the four wastewater samples (all but LI_11-10). Wastewater influent sample LI_13-9 contained the most signatures, with approximately 75,000 reads, assembled into 120 contigs, representing genome fragments of 10 of the 11 rotavirus segments. At the species level, these genome fragments were assigned to either species *Rotavirus A* or *Rotavirus C*. Comparing the amino acid sequences of the predicted proteins, some contigs showed high levels of identity (>88%) with the segments of either rotavirus A (RVA) or rotavirus C (RVC) reference genomes as available in the RefSeq database (26, 27), while others showed lower identities with a variety of RVC isolates only. The segmented-genome nature and the possibility of segment exchange make it difficult to confidently identify the number of rotavirus types present in this sample. Given the amino acid similarities with both RVA and RVC types, we suggest there are at least two and possibly three types present here.

Using the RotaC 2.0 typing tool for RVA and Blast-based similarity to known genotypes, we have typed the rotavirus genome segments found here (Table 2). The combined genomic makeup of the RV community in sample LI_13-9 was G8/G10/Gx-P[1]/P[14]/P[41]/P[x]-I2/Ix-R2/Rx-C2/Cx-M2/Mx-A3/A11/Ax-Nx-T6/Tx-E2/Ex (28, 29).

**FIG 5** Maximum-likelihood phylogenetic tree of norovirus genomes belonging to genogroup GI, with the norovirus GII reference genome as an outlier. The nucleotide sequences were aligned with MUSCLE, and the alignment was trimmed to the length of contig 6 of the LE_11-10 virome sequence, resulting in 7,758 positions analyzed for tree building. The maximum-likelihood method was used, with the Tamura-Nei model for nucleic acid substitution. The percentages of trees in which the associated taxa clustered together are shown next to the branches. The scale bar represents the number of substitutions per site.

**Partial genomes of other potentially pathogenic RNA viruses.** In sample LI_13-9, a small contig of 347 bases was found that was 94% identical at the nucleotide level to ORF1 of sapovirus Mc2 (AY237419) in the family *Caliciviridae*. We also identified four contigs of approximately 500 bases in sample LE_11-10 that resembled most closely the astrovirus MLB2 isolates MLB2/human/Geneva/2014 (KT224358) and MLB2-LIHT (KX022687), at 99% nucleotide identity. In addition, we identified several reads and contigs assigned to the family *Picornaviridae*, which comprises a diverse set of enteric viruses, but the closest relatives in the databases were metagenomically assembled or unidentified picornaviruses.

**Picobirnaviruses showed a high prevalence in wastewater.** All the wastewater virome libraries contained signatures assigned to the dsRNA family *Picobirnaviridae*, genus *Picobirnavirus* (Fig. 2), and these reads assembled into between 42 (LE_13-9) and 510 (LI_13-9) contigs. Both picobirnavirus genome segments, segment 1 encoding two hypothetical proteins and segment 2 on which the RNA-dependent RNA polymerase (RdRP) is encoded, were observed in the samples. The contigs showed little sequence similarity with the *Human picobirnavirus* reference genome (RefSeq GenBank accession numbers NC_007026.1 and NC_007027.1). Phylogenetic analysis of a partial region of the predicted RdRPs in the virome contigs was not able to resolve any cluster or evolutionary origin (Fig. 6A). Picobirnavirus RdRPs from human, animal, and environmental isolates, as well as the majority of the virome sequences, were grouped in one large, unsupported cluster that showed relatively little genomic diversity. While many picobirnaviruses have been isolated from humans with gastroenteritis, a review of the

**TABLE 2** Rotavirus A and C genome information and detection in the LI_13-9 sample data set

| Virus, genome segment | Length (nt) | Protein(s) encoded | Predicted function | No. of contigs (no. of RVCX contigs[a]) | Putative genotype(s) | Potential host(s)[b] |
|---|---|---|---|---|---|---|
| RVA |
| 1 | 3,302 | VP1 | RNA-dependent RNA polymerase | 7 | R2 | Human, cow |
| 2 | 2,693 | VP2 | Core capsid protein | 1 | C2 | Human |
| 3 | 2,591 | VP3 | RNA capping protein | 1 | M2 | Human, sheep |
| 4 | 2,363 | VP4 | Outer capsid spike protein | 3 | P[1], P(41), P[14] | Human, pig, alpaca, monkey |
| 5 | 1,614 | NSP1 | Interferon antagonist protein | 6 | A3, A11 | Human, cow, pig, deer |
| 6 | 1,356 | VP6 | Inner capsid protein | 1 | I2 | Human |
| 7 | 1,105 | NSP3 | Translation effector protein | 4 | T6 | Human, dog, cow |
| 8 | 1,059 | NSP2 | Viroplasm RNA binding protein | 0 | | |
| 9 | 1,062 | VP7 | Outer capsid glycoprotein | 2 | G10, G8 | Cow, human |
| 10 | 751 | NSP4 | Enterotoxin | 1 | E2 | Human, cow |
| 11 | 667 | NSP5 and -6 | Phosphoprotein, nonstructural protein | 0 | | |
| RVC |
| 1 | 3,309 | VP1 | RNA-dependent RNA polymerase | 7 (0) | Rx | Pig, cow |
| 2 | 2,736 | VP2 | Core capsid protein | 4 (2) | Cx | Pig, dog |
| 3 | 2,283 | VP4 | Outer capsid protein | 2 (4) | P[x] | Pig |
| 4 | 2,166 | VP3 | Guanylyl transferase | 6 (0) | Mx | Pig |
| 5 | 1,353 | VP6 | Inner capsid protein | 1 (0) | Ix | Pig |
| 6 | 1,350 | NSP3 | | 0 (1) | Tx | Human |
| 7 | 1,270 | NSP1 | | 0 (2) | Ax | Pig, dog |
| 8 | 1,063 | VP7 | Outer capsid glycoprotein | 0 (2) | Gx | Pig |
| 9 | 1,037 | NSP2 | | 2 (0) | Nx | Pig |
| 10 | 730 | NSP5 | | 0 (0) | | |
| 11 | 613 | NSP4 | Enterotoxin | 0 (4) | Ex | Pig |

[a]Number of predicted RVCX contigs are in parentheses, i.e., contigs with only limited amino acid similarity to RVC.
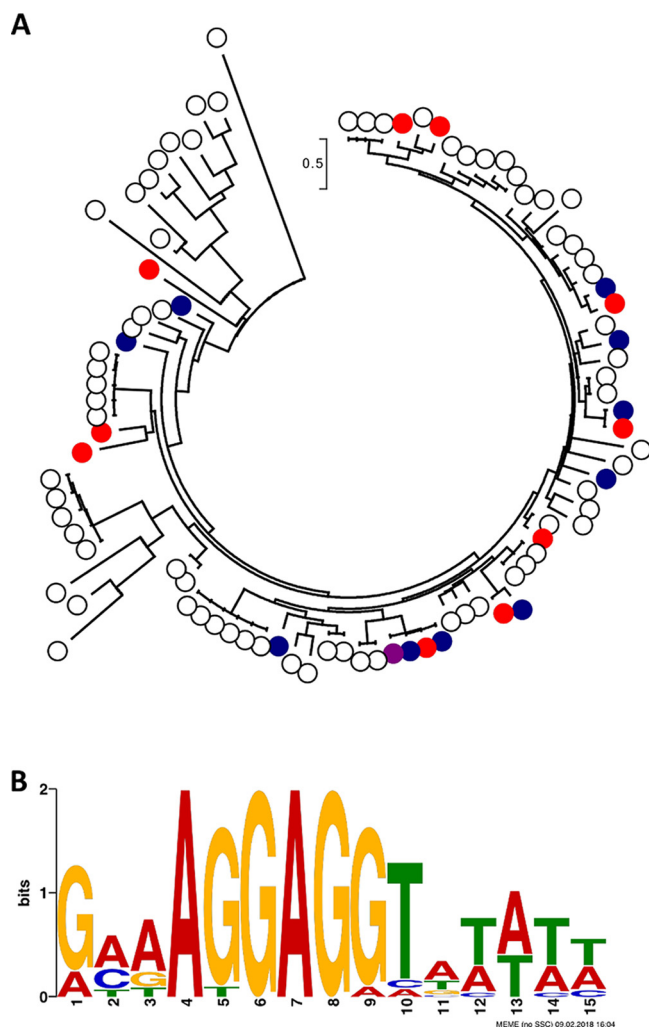[b]Potential hosts are defined as the hosts of the reference rotavirus sequence with the highest similarity to the contigs found in the virome sample LI_13-9.

known cases suggested that picobirnaviruses are probably not the main cause of acute diarrhea and are secondary pathogens with potentially synergistic effects (30). A qRT-PCR-based investigation into the suitability of human picobirnaviruses as indicators of human fecal contamination showed that they were not present in a sufficient proportion of tested samples to be good water quality indicators (31), but their high relative abundance in our sample set warrants further investigation for their use as water quality markers using metaviromic methods.

A recent study of picobirnaviruses gave rise to the hypothesis that these viruses do not infect mammals but are a new family of RNA bacteriophages, based on the presence of bacterial ribosome binding sites (RBS) upstream from the coding sequences (CDSs) (32). To test this hypothesis, we extracted all contigs with amino acid similarity to RdRPs or capsid proteins of known picobirnaviruses, annotated the CDSs, and extracted the 21 nucleotides upstream from the transcription start site. In the 233 contigs found, 71 partial CDSs were predicted, from which we extracted 17 5′ UTRs (untranslated regions), discarding those partially annotated CDSs missing the transcription start site. We discovered the 6-mer motif AGGAGG (Fig. 6B) in 100% of the upstream sequences, similar to the frequency reported by Krishnamurthy and Wang (32), who found at least a 4-mer RBS in 100% of the 98 picobirnavirus 5′ UTRs investigated. In contrast, the different families of eukaryotic viruses analyzed in that study only showed a low incidence of RBSs, which were mostly 4-mers. Our findings, therefore, support the hypothesis that picobirnaviruses are bacteriophages, and we suggest that they belong to a novel RNA bacteriophage family with a high level of genomic diversity.

## DISCUSSION

We set out to explore the possibility of using viromics to find human-pathogenic RNA viruses in the environment. We have been successful in identifying several potentially human-pathogenic, including potentially zoonotic, viral genomes from the wastewater samples but did not find any in the surface estuarine water and sediment

**FIG 6** Picobirnavirus diversity. (A) Maximum-likelihood phylogenetic tree of RdRP amino acid sequences of isolated and virome picobirnaviruses. Sequences from isolates are indicated with white dots and virome-derived sequences with filled colored dots, as follows: sample LI_11-10 in purple, sample LE_11-10 in blue, and sample LI_13-9 in red. Sequences were aligned using MUSCLE, providing 114 amino acid positions for tree generation. The maximum-likelihood method was used, with the JTT matrix-based model. The scale bar represents the number of substitutions per site. The bootstrap values of all branches were low. (B) Predicted ribosome binding site consensus sequence from extracted 5' UTRs. The logo was produced using the MEME Suite.

samples. The absence of signatures does not necessarily mean that there are no pathogenic viruses present in water or sediment but only that their levels could be below our limit of quantification for qPCR (approximately 200 gc/liter).

It is important to note here that during the RNA extraction process, many biases could have been introduced, leading to a lower recovery of input viruses. Samples were first concentrated from volumes of 1 liter (wastewater) or 50 liters (surface water) down to 50 ml, using tangential flow filtration (TFF) at a molecular-mass cutoff of 100 kDa, followed by polyethylene glycol (PEG) 6000 precipitation. These samples were diluted in fresh buffer, filtered through syringe filters of 0.22-$\mu$m pore size, and then treated with nuclease to remove free DNA and RNA. Previous research has shown that, while any enrichment method aimed at fractionating the viral and cellular components will decrease the total quantity of viruses, a combination of centrifugation, filtration, and nuclease treatment increases the proportion of viral reads in sequencing data sets (33). After implementing these steps, we performed viral RNA extraction using the Mo Bio PowerViral environmental DNA/RNA extraction kit, which has previously been shown to

perform best overall in spiking experiments with murine norovirus, in terms of extraction efficiency and removal of inhibitors (34). The kit has, however, given low recoveries of viruses from sediment (35).

We did not perform an amplification step before library construction with the NEBNext Ultra directional RNA library preparation kit for Illumina, to retain the genome sense and strand information. Instead, we increased the number of cycles of random PCR during library preparation from 12 to 15 to counteract the low input quantity of RNA (<1 ng). The random amplification during library construction led to a trade-off in which genome strand information was gained for a loss of quantitative power, making it difficult to compare abundances of viral types within and across libraries. This random-PCR-based bias has been highlighted before, but the proposed solution of using library preparation protocols which limit the use of PCR is only feasible with large amounts of input nucleic acid (36), which we have not found to be possible when processing environmental/wastewater samples to generate RNA metaviromes.

A critical issue to highlight here is the inclusion of controls in our sequencing libraries in order to identify potential contaminants and their origins, as has been suggested previously (37, 38). There have been multiple reports of false-positive genome discoveries, in particular the novel parvovirus-like hybrid in hepatitis patients that was later revealed to originate from the silica-based nucleic acid extraction columns (39–41). In this study, we included a positive control that comprised bacterial cells (*Salmonella enterica* serovar Typhimurium isolate D23580; GenBank RefSeq accession number NC_016854) and mengovirus (36), an RNA virus that serves as a process control, and two negative controls, an extraction control and a library preparation control. Analysis of the control libraries showed that while the *Salmonella* cells and DNA were successfully removed from the positive-control sample by the enrichment protocol, the mengovirus was not recovered. Subsequent qRT-PCR analysis revealed that the mengovirus remained detectable in the preprocessing stages of the extraction but was lost after RNase treatment (data not shown). The inclusion of an inactivation step of the DNase at 75°C potentially exacerbated the effect of the RNase step. Consequently, it is likely that we have missed viral types during the extraction process despite having still managed to recover an RNA metavirome harboring substantial diversity.

Further examination of the HiSeq and MiSeq control data sets revealed a wide range of contaminant signatures of prokaryotic, eukaryotic, and viral origin, making up 45 million read pairs per control on the HiSeq platform and 1 million read pairs for the MiSeq, even though the 16S and 18S rRNA PCR and RT-PCR reactions produced no visible bands on an agarose gel. Most bacterial contaminant reads belonged to the phyla *Proteobacteria*, *Actinobacteria*, and *Firmicutes*. The most abundant genera included *Corynebacterium*, *Propionibacterium*, *Sphingomonas*, *Ralstonia*, *Pseudomonas*, *Streptomyces*, *Staphylococcus*, and *Streptococcus*, members of which have been identified as common laboratory contaminants in the past (42). Within the eukaryotic signatures, human-derived, *Beta vulgaris*, and *Anopheles* reads were the most prevalent, pointing toward potential cross-contamination of the sequencing libraries. A small number of virus signatures were also identified, with the most prominent being feline calicivirus and dengue virus. The presence of the calicivirus was traced back to the library preparation kits after the libraries were reconstructed and resequenced. The dengue virus signature was a <100-nt sequence which was coextracted in all the samples and potentially originated in one of the reagents or the spin extraction column. All sequences present in the controls were carefully removed from the sample data sets during the quality control stage of the bioinformatics processing before further analysis. For future experiments, we will omit the RNase treatment step during extraction and filter out any contaminating rRNA or cellular-derived mRNA sequences as part of the bioinformatic quality control workflow.

Our results show that while contamination is an issue when dealing with low-biomass samples, the combination of increased random PCR cycles during library preparation, deep sequencing (i.e., HiSeq rather than MiSeq), and computational

subtraction of control sequences provides data of sufficient quantity and quality to assemble nearly complete RNA virus genomes *de novo*.

**Norovirus.** Noroviruses are one of the most common causes of gastrointestinal disease in the developed world, with an incidence in the United Kingdom estimated as approaching 4 million cases per annum (43). The genotype most commonly associated with disease is GII.4 (44–46), which was not detected in the metaviromes generated here.

We retrieved one norovirus GI genome, assembled from 22,151 reads, in wastewater effluent sample LE_11-10. This finding was in direct conflict with the qRT-PCR analysis of this sample, which did not detect any NoV GI signatures (Table 1). In contrast, NoV GII signatures were detected by qRT-PCR, but no NoV GII genomes or genome fragments were observed in the virome libraries. One hypothesis to explain the discrepancy between PCR and viromics approaches lies in the differences in extraction protocol. For qRT-PCR, no viral enrichment step was performed and RNA was not extracted with the PowerViral kit. Therefore, NoV GII could have been lost before virome sequencing, as was the process control mengovirus. An alternative hypothesis is that the NoV GII signatures detected during qRT-PCR were derived from fragmented RNA or from particles with a compromised capsid. In both these cases, the RNA would not be detected in the virome data because of the RNase preprocessing steps implemented in the enrichment/extraction protocol. This calls into question the reliance of qRT-PCR for NoV detection and whether the detected viruses are infectious or merely remnants of previous infections. Further research using, for example, capsid integrity assays combined with infectious particle counts will need to be conducted to assess the validity of qRT-PCR protocols for norovirus detection.

The inability to identify NoV GI with qRT-PCR might be related to the mismatched base present in the forward primer sequence used for detection, but even without mismatches, primer-probe pairs can be improved to provide better detection. In a recent study, researchers designed an improved probe for NoV GI.2 strains, lowering the limit of detection for these strains from waterborne samples (47). It is, therefore, possible that the NoV GI.2 detected here with viromics methods was present below the limit of detection of the ISO standard primer/probe combination (48) used in our study. Viromics as a means of investigating water samples for the presence of norovirus does have the advantage of demonstrating the presence of an undegraded genome, provided the sample processing requirements do not lead to excessive loss of virus particles, resulting in false negatives. Certainly, time and cost permitting, viromics is a useful adjunct to qPCR for samples for which it is deemed particularly important or critical that the presence of intact viral genome be determined.

While there have been recent breakthroughs in growing human NoV (49, 50), its culture remains very difficult and not yet suitable for routine testing for NoV in the environment. Hence, studies using male-specific coliphages, such as MS2 and GA, which are ssRNA phages belonging to the family *Leviviridae*, are still worthwhile as alternative model systems (51, 52). Interestingly, while some levivirus signatures were present in all wastewater samples (<500 reads), we observed a striking cooccurrence of these viruses with norovirus signatures in both libraries of sample LE_11-10 (>2,500 reads). The most commonly observed viruses in this sample were *Pseudomonas* phage PRR1, an unclassified levivirus, and *Escherichia* phages FI and M11 in the genus *Allolevivirus*. Further studies with more samples and replicates will indicate whether there is a significant correlation between the presence of leviviruses and noroviruses in water samples. Furthermore, the higher abundance of alloleviviruses than of MS2-like viruses could indicate that the former might be more relevant as model systems for noroviruses.

**Rotavirus.** Rotaviruses are, like noroviruses, agents of gastroenteritis, but the disease is commonly associated with children under the age of 5, where severe diarrhea and vomiting can lead to over 10,000 hospitalizations per year in England and Wales (53). Since the introduction of the live attenuated vaccine Rotarix, the incidence of

gastroenteritis in England has declined, specifically for children aged <2 and during peak rotavirus seasons (54–56). Therefore, the discovery of a diverse assemblage of rotavirus genome segments in the wastewater samples here was less expected than the norovirus discovery. While we were unable to recover the genome of the vaccine strain, our genomic evidence suggests that at least one RVA and one RVC population were circulating in the Llanrwst region in September 2016.

The genome constellation for the RVA segments in sample LI_13-9, G8/G10-P[1]/P[14]/P[41]-I2-R2-C2-M2-A3/A11-(N)-T6-E2-(H) (N and H segments were not recovered in this study), is distinctly bovine in origin (28). The closest genome segment relatives based on nucleic acid similarity, however, have been isolated from humans (Table 2), possibly pointing toward a bovine-human zoonotic transmission of this virus (57). The same genomic constellation has been found recently when unusual G8P[14] RVA isolates were recovered from human strain collections in Hungary (58) and Guatemala (59) and isolated from children in Slovenia (60) and Italy (61). Cook and colleagues calculated that there would be approximately 5,000 zoonotic human infections per year in the United Kingdom from livestock transmission, but many would be asymptomatic (62). The presence of RVA in the wastewater of Llanrwst could be from zoonotically infected individuals shedding into the wastewater, but it is equally likely that RVA from cattle farms in the area spilled over into the sewage system.

The origins of the RVC genome segments are more difficult to trace, because of lower similarity scores with known RVC isolates. The majority of the segments were similar to porcine RVC genomes, while others showed no nucleotide similarity at all, only a low degree of amino acid similarity. An explanation for the presence of pig-derived rotavirus signatures could be farm runoff. While farm waste is not supposed to end up in the sewage treatment plant, it is likely that the RVC segments originate directly from pigs, not through zoonotic transfer. Runoff from fields onto public roads, broken farm sewer pipes, or polluted small streams might lead to porcine viruses entering the human sewerage network, but we cannot provide formal proof from the data available. Based on the evidence, we hypothesize that there are one or possibly two divergent strains of RVC circulating in the pig farms in the Llanrwst area.

**Conclusion.** In this study, we investigated the use of metagenomics for the discovery of RNA viruses circulating in watercourses. We have found RNA viruses in all samples tested, but potential human-pathogenic viruses were only identified in wastewater. The recovery of plant viruses in most samples points toward potential applications in crop protection, for example, the use of metaviromics in phytopathogen diagnostics. However, technical limitations, including the amount of input material necessary and contamination of essential laboratory consumables and reagents, are currently the main bottleneck for the adoption of fine-scale metagenomics in routine monitoring and diagnostics. The discovery of a norovirus GI and a diverse set of rotavirus segments in the corresponding metaviromes indicates that qPCR-based approaches can miss a significant portion of relevant pathogenic RNA viruses present in water samples. Therefore, metagenomics can, at this time, best be used for exploration, to design new diagnostic markers/primers targeting novel genotypes, and to inform diagnostic surveys on the inclusion of specific additional target viruses.

## MATERIALS AND METHODS

**Sample collection and processing.** Wastewater samples were collected as part of a viral surveillance study described elsewhere (63). Wastewater influent and effluent samples, 1 liter each, were collected at the Llanrwst wastewater treatment plant by Welsh Water (Wales, United Kingdom) (Fig. 1) on 12 September 2016 (processed on 13 September; sample designations LI_13-9 and LE_13-9) and 10 October 2016 (processed on 11 October; sample designations LI_11-10 and LE_11-10). The wastewater treatment plant uses filter beds for secondary treatment and serves approximately 4,000 inhabitants. The estuarine surface water sample (SW; 50 liters) was collected at Morfa Beach (Conwy, Wales) (Fig. 1), approximately 22 km downstream from the Llanrwst wastewater treatment plant, on 19 October and 2 November 2016 at low tide (only the sample from November was used for sequencing, as the October sample extract failed quality control). Together with the surface water sample, 90 g of the top 1- to 2-cm layer of the sediment was also collected (sample designations Sed1 for the October sample and Sed2 for the November sample).

The wastewater and surface water samples were processed using a two-step concentration method as described elsewhere (63). In brief, the 1-liter (wastewater) and 50-liter (surface water) samples were first concentrated down to 50 ml using a KrosFlo research IIi tangential flow filtration (TFF) system (Spectrum Labs, United States) with a 100-PES (polyethersulfone) membrane. Particulate matter was then eluted from solid matter in the concentrates, using beef extract buffer, and then viruses were precipitated using polyethylene glycol (PEG) 6000. The viruses from the sediment samples were eluted and concentrated using beef extract elution and PEG precipitation as described elsewhere (35). The precipitates were eluted in 2- to 10-ml phosphate saline buffer (PBS, pH 7.4) and stored at −80°C.

**Detection and quantification of enteric viruses with qRT-PCR.** Total nucleic acids were extracted from a 0.5-ml aliquot of the concentrates using the NucliSENS miniMag nucleic acid purification system (BioMérieux, France). The final volumes of the nucleic acid solution were 0.05 ml (surface water and sediment) and 0.1 ml (wastewater samples). Norovirus GI (64, 65) and GII (66, 67), sapovirus GI (68), and hepatitis A and E viruses (69, 70) were targeted in qRT-PCR assays as described elsewhere (71).

**Viral RNA extraction for metaviromic sequencing.** Viral particles were extracted from the concentrated samples by filtration. In a first step, the samples were diluted in 10 ml of sterile 0.5 M NaCl buffer and incubated at room temperature (20°C) with gentle shaking for 30 min to disaggregate particles. The suspension was then filtered through a sterile, 0.22-$\mu$m pore-size syringe filter (PES membrane; Millex). The sample was desalted by centrifugation (3,200 $\times$ $g$, between 1 and 6 h for different samples) in a sterilized spin filter (Vivaspin 20, 100-kDa molecular-mass cutoff) and replacement of the buffer solution with 5 ml of a Tris-based buffer (10 mM Tris-HCl, 10 mM MgSO$_4$, 150 mM NaCl, pH 7.5). The buffer exchange was performed twice, and the volume retained after the final spin was <500 $\mu$l. The samples were then treated with Turbo DNase (20 units; Ambion) and incubated for 30 min at 37°C, followed by inactivation at 75°C for 10 min. In a next step, all samples were treated with 80 $\mu$g RNase A (Thermo Fisher Scientific) and incubated at 37°C for 30 min. The RNase was inactivated with RiboLock RNase inhibitor (Thermo Fisher Scientific), the inactivated complex was removed by spin filtration (Vivaspin 500, 100-kDa molecular-mass cutoff), and the samples centrifuged until the volume was approximately 200 $\mu$l. Viral DNA and RNA were coextracted using the PowerViral environmental DNA/RNA kit (Mo Bio Laboratories) according to the manufacturer's instructions. In this protocol, buffer PV1 was supplemented with 20 $\mu$l/ml betamercaptoethanol to further reduce RNase activity. The nucleic acid was eluted in 100 $\mu$l RNase-free water. The extracted viral DNA was degraded using the DNase Max kit (Mo Bio Laboratories) according to the manufacturer's instructions. The remaining viral RNA was further purified and concentrated by ethanol precipitation using 2.5$\times$ the sample volume of 100% ethanol and 1/10 volume of diethyl pyrocarbonate (DEPC)-treated Na-acetate (3 M). The quantity and quality of RNA were determined with Bioanalyzer Pico RNA 6000 capillary electrophoresis (Agilent Technologies). Positive and negative extraction control samples were processed alongside the main samples. The positive-control samples contained *Salmonella enterica* serovar Typhimurium strain D23580, which is not found in the United Kingdom (72), and mengovirus as a process control virus (71, 73).

The viral RNA extracts were tested for bacterial and eukaryotic cellular contamination using 16S and 18S rRNA gene PCR and RT-PCR, with primers e9F (74) and 519R (75) for the 16S rRNA gene and primers 1389F and 1510R (76) for the 18S rRNA gene. Complementary DNA was created using SuperScript III reverse transcriptase (Invitrogen) with random hexamer primers according to the manufacturer's instructions. RT-PCR was performed with MyTaq red mix (Bioline) for 35 cycles (95°C for 45 s, 50°C for 30 s, and 72°C 1 min 40 s) and visualized on a 1% agarose gel. Samples were considered suitable for sequencing if no DNA bands were visible on the gel.

**Library preparation and sequencing.** The library preparation and sequencing were performed at the University of Liverpool Centre for Genomics Research (CGR). Twelve dual-indexed, strand-specific libraries were created using the NEBNext Ultra directional RNA library preparation kit for Illumina, according to the manufacturer's instructions. These libraries were pooled and sequenced at 2 $\times$ 150-bp read lengths on the Illumina HiSeq 4000 platform. This generated between 10 and 110 million paired reads per sample.

To confirm our results, a second set of libraries was constructed from new kits and a Milli-Q water sample was included as a library preparation control. The 13 resulting libraries were sequenced at 2 $\times$ 150-bp read lengths on the Illumina MiSeq platform at the CGR, University of Liverpool. These data were used for verification and control purposes only, as the sequencing depth was insufficient for the bioinformatics analyses described in the rest of the study.

**Bioinformatics.** All command line programs for data analysis were run on the bioinformatics cluster of CGR (University of Liverpool) in a Debian 5 or 7 environment.

Raw fastq files were trimmed to remove Illumina adapters using Cutadapt version 1.2.1 with option -O 3 (77) and Sickle version 1.200 with a minimum quality score of 20 (78). Further quality control was performed with Prinseq-lite (79) with the following parameters: minimum read length of 35, GC percentage between 5 and 95%, minimum mean quality of 25, dereplication (removal of identical reads, leaving 1 copy), and removal of tails of a minimum of 5 poly(N) sequences from 3′ and 5′ ends of reads.

The positive- and negative-control libraries described earlier were used for contaminant removal. The reads of the control samples were analyzed using Diamond blastx (17) against the nonredundant protein database of NCBI (nr, November 2015 version). The blast results were visualized using MEGAN6 Community Edition (18). An extra contaminant file was created with the complete genomes of species present at over 1,000 reads in the positive- and negative-control samples. Then, bowtie2 (80) was used for each sample to subtract the reads that mapped to the positive-control, negative-control, or contaminant file. The unmapped reads were used for assembly with SPAdes version 3.9.0, with $k$-mer values of 21, 31, 41, 51, 61, and 71 and the options --careful and a minimum coverage of 5 reads per contig (81).

The contig files of each sample were compared with the contigs of the controls (assembled using the same parameters) using blastn of the BLAST+ suite (82). Contigs that showed significant similarity with control contigs were manually removed, creating a curated contig data set. The unmapped read data sets were then mapped against this curated contig data set with bowtie2, and only the reads that mapped were retained, resulting in a curated read data set.

The curated contig and read data sets were compared to the RefSeq viral (January 2017 release) and nonredundant protein (nr, May 2017 release) reference databases using Diamond blastx at an e value of 1e−5 for significant hits (17, 83, 84). Taxon assignments were made with MEGAN6 Community Edition according to the lowest-common-ancestor algorithm with default settings (18). We chose the family level taxon assignments to represent the overall viral diversity because there is generally little amino acid identity between viral families. The taxon abundance data were extracted from MEGAN6 and imported into RStudio for visualization (85). Genes on the assembled contigs were predicted with Prokka (86) using the settings --kingdom Viruses and an e value of 1e−5. Multiple alignments of genes and genomes were made in MEGA7 using the MUSCLE algorithm with default settings (87, 88). The alignments were manually trimmed, and phylogenetic trees were built using the maximum-likelihood method in MEGA7 with the default settings. Sequences upstream from potential CDSs of Prokka-annotated picobirnaviruses were extracted using extractUpStreamDNA (https://github.com/ajvilleg/extractUpStreamDNA), and all 5′ UTRs and transcription start sites were manually verified in UGene (89). These extracted sequences were then subjected to a motif search using the MEME Suite (90, 91).

**Accession numbers.** Read and contig data sets are available from NCBI under the following BioProject accession numbers: PRJNA421889 (wastewater data), PRJNA421892 (sediment data), and PRJNA421894 (estuarine water data). The NoV GI genome isolate was deposited in GenBank under accession number MG599789.

## REFERENCES

1. Lin J, Ganesh A. 2013. Water quality indicators: bacteria, coliphages, enteric viruses. Int J Environ Health Res 23:484–506. https://doi.org/10.1080/09603123.2013.769201.
2. Girones R, Ferrús MA, Alonso JL, Rodriguez-Manzano J, Calgua B, de Abreu Corrêa A, Hundesa A, Carratala A, Bofill-Mas S. 2010. Molecular detection of pathogens in water—the pros and cons of molecular techniques. Water Res 44:4325–4339. https://doi.org/10.1016/j.watres.2010.06.030.
3. Laverick MA, Wyn-Jones AP, Carter MJ. 2004. Quantitative RT-PCR for the enumeration of noroviruses (Norwalk-like viruses) in water and sewage. Lett Appl Microbiol 39:127–136. https://doi.org/10.1111/j.1472-765X.2004.01534.x.
4. Rodriguez-Manzano J, Miagostovich M, Hundesa A, Clemente-Casares P, Carratala A, Buti M, Jardi R, Girones R. 2010. Analysis of the evolution in the circulation of HAV and HEV in Eastern Spain by testing urban sewage samples. J Water Health 8:346–354. https://doi.org/10.2166/wh.2009.042.
5. Schvoerer E, Ventura M, Dubos O, Cazaux G, Serceau R, Gournier N, Dubois V, Caminade P, Fleury HJA, Lafon ME. 2001. Qualitative and quantitative molecular detection of enteroviruses in water from bathing areas and from a sewage treatment plant. Res Microbiol 152:179–186. https://doi.org/10.1016/S0923-2508(01)01190-1.
6. Fong TT, Phanikumar MS, Xagoraraki I, Rose JB. 2010. Quantitative detection of human adenoviruses in wastewater and combined sewer overflows influencing a Michigan river. Appl Environ Microbiol 76:715–723. https://doi.org/10.1128/AEM.01316-09.
7. Bofill-Mas S, Albinana-Gimenez N, Clemente-Casares P, Hundesa A, Rodriguez-Manzano J, Allard A, Calvo M, Girones R. 2006. Quantification and stability of human adenoviruses and polyomavirus JCPyV in wastewater matrices. Appl Environ Microbiol 72:7894–7896. https://doi.org/10.1128/AEM.00965-06.
8. Hellmér M, Paxéus N, Magnius L, Enache L, Arnholm B, Johansson A, Bergström T, Norder H. 2014. Detection of pathogenic viruses in sewage provided early warnings of hepatitis A virus and norovirus outbreaks. Appl Environ Microbiol 80:6771–6781. https://doi.org/10.1128/AEM.01981-14.
9. Nieuwenhuijse DF, Koopmans MPG. 2017. Metagenomic sequencing for surveillance of food- and waterborne viral diseases. Front Microbiol 8:230. https://doi.org/10.3389/fmicb.2017.00230.
10. Symonds EM, Breitbart M. 2015. Affordable enteric virus detection techniques are needed to support changing paradigms in water quality management. Clean Soil Air Water 43:8–12. https://doi.org/10.1002/clen.201400235.
11. Rosario K, Symonds EM, Sinigalliano C, Stewart J, Breitbart M. 2009. Pepper mild mottle virus as an indicator of fecal pollution. Appl Environ Microbiol 75:7261–7267. https://doi.org/10.1128/AEM.00410-09.
12. Stachler E, Bibby K. 2014. Metagenomic evaluation of the highly abundant human gut bacteriophage CrAssphage for source tracking of human fecal pollution. Environ Sci Technol Lett 1:405–409. https://doi.org/10.1021/ez500266s.
13. Bibby K, Peccia J. 2013. Identification of viral pathogen diversity in sewage sludge by metagenome analysis. Environ Sci Technol 47:1945–1951. https://doi.org/10.1021/es305181x.
14. Cantalupo PG, Calgua B, Zhao G, Hundesa A, Wier AD, Katz JP, Grabe M,

Hendrix RW, Girones R, Wang D, Pipas JM. 2011. Raw sewage harbors diverse viral populations. mBio 2:e00180-11. https://doi.org/10.1128/mBio.00180-11.

15. Ng TFF, Marine R, Wang C, Simmonds P, Kapusinszky B, Bodhidatta L, Oderinde BS, Wommack KE, Delwart E. 2012. High variety of known and new RNA and DNA viruses of diverse origins in untreated sewage. J Virol 86:12161–12175. https://doi.org/10.1128/JVI.00869-12.

16. Fernandez-Cassi X, Timoneda N, Martínez-Puchol S, Rusiñol M, Rodriguez-Manzano J, Figuerola N, Bofill-Mas S, Abril JF, Girones R. 2018. Metagenomics for the study of viruses in urban sewage as a tool for public health surveillance. Sci Total Environ 618:870–880. https://doi.org/10.1016/j.scitotenv.2017.08.249.

17. Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using DIAMOND. Nat Methods 12:59–60. https://doi.org/10.1038/nmeth.3176.

18. Huson DH, Weber N. 2013. Microbial community analysis using MEGAN. Methods Enzymol 531:465–485. https://doi.org/10.1016/B978-0-12-407863-5.00021-6.

19. King AMQ, Adams MJ, Carstens EB, Lefkowitz EJ (ed). 2012. Virus taxonomy. Classification and nomenclature of viruses. Ninth report of the International Committee on Taxonomy of Viruses. Elsevier Academic Press, San Diego, CA.

20. Ye Y, Ellenberg RM, Graham KE, Wigginton KR. 2016. Survivability, partitioning, and recovery of enveloped viruses in untreated municipal wastewater. Environ Sci Technol 50:5077–5085. https://doi.org/10.1021/acs.est.6b00876.

21. Aquino De Carvalho N, Stachler EN, Cimabue N, Bibby K. 2017. Evaluation of Phi6 persistence and suitability as an enveloped virus surrogate. Environ Sci Technol 51:8692–8700. https://doi.org/10.1021/acs.est.7b01296.

22. Winterbourn JB, Clements K, Lowther JA, Malham SK, McDonald JE, Jones DL. 2016. Use of Mytilus edulis biosentinels to investigate spatial patterns of norovirus and faecal indicator organism contamination around coastal sewage discharges. Water Res 105:241–250. https://doi.org/10.1016/j.watres.2016.09.002.

23. Patel MM, Hall AJ, Vinjé J, Parashar UD. 2009. Noroviruses: a comprehensive review. J Clin Virol 44:1–8. https://doi.org/10.1016/j.jcv.2008.10.009.

24. Zheng DP, Ando T, Fankhauser RL, Beard RS, Glass RI, Monroe SS. 2006. Norovirus classification and proposed strain nomenclature. Virology 346:312–323. https://doi.org/10.1016/j.virol.2005.11.015.

25. Huo Y, Cai A, Yang H, Zhou M, Yan J, Liu D, Shen S. 2014. Complete nucleotide sequence of a norovirus GII.4 genotype: evidence for the spread of the newly emerged pandemic Sydney 2012 strain to China. Virus Genes 48:356–360. https://doi.org/10.1007/s11262-013-1018-8.

26. Small C, Barro M, Brown TL, Patton JT. 2007. Genome heterogeneity of SA11 rotavirus due to reassortment with "O" agent. Virology 359:415–424. https://doi.org/10.1016/j.virol.2006.09.024.

27. Chen Z, Lambden PR, Lau J, Caul EO, Clarke IN. 2002. Human group C rotavirus: completion of the genome sequence and gene coding assignments of a non-cultivatable rotavirus. Virus Res 83:179–187. https://doi.org/10.1016/S0168-1702(01)00442-7.

28. Matthijnssens J, Ciarlet M, Heiman E, Arijs I, Delbeke T, McDonald SM, Palombo EA, Iturriza-Gómara M, Maes P, Patton JT, Rahman M, Van Ranst M. 2008. Full genome-based classification of rotaviruses reveals a common origin between human Wa-like and porcine rotavirus strains and human DS-1-like and bovine rotavirus strains. J Virol 82:3204–3219. https://doi.org/10.1128/JVI.02257-07.

29. Matthijnssens J, Ciarlet M, Rahman M, Attoui H, Bányai K, Estes MK, Gentsch JR, Iturriza-Gómara M, Kirkwood CD, Martella V, Mertens PPC, Nakagomi O, Patton JT, Ruggeri FM, Saif LJ, Santos N, Steyer A, Taniguchi K, Desselberger U, Van Ranst M. 2008. Recommendations for the classification of group a rotaviruses using all 11 genomic RNA segments. Arch Virol 153:1621–1629. https://doi.org/10.1007/s00705-008-0155-1.

30. Ganesh B, Bányai K, Martella V, Jakab F, Masachessi G, Kobayashi N. 2012. Picobirnavirus infections: viral persistence and zoonotic potential. Rev Med Virol 22:245–256. https://doi.org/10.1002/rmv.1707.

31. Hamza IA, Jurzik L, Überla K, Wilhelm M. 2011. Evaluation of pepper mild mottle virus, human picobirnavirus and torque teno virus as indicators of fecal contamination in river water. Water Res 45:1358–1368. https://doi.org/10.1016/j.watres.2010.10.021.

32. Krishnamurthy SR, Wang D. 2018. Extensive conservation of prokaryotic ribosomal binding sites in known and novel picobirnaviruses. Virology 516:108–114. https://doi.org/10.1016/j.virol.2018.01.006.

33. Hall RJ, Wang J, Todd AK, Bissielo AB, Yen S, Strydom H, Moore NE, Ren X, Huang QS, Carter PE, Peacey M. 2014. Evaluation of rapid and simple techniques for the enrichment of viruses prior to metagenomic virus discovery. J Virol Methods 195:194–204. https://doi.org/10.1016/j.jviromet.2013.08.035.

34. Iker BC, Bright KR, Pepper IL, Gerba CP, Kitajima M. 2013. Evaluation of commercial kits for the extraction and purification of viral nucleic acids from environmental and fecal samples. J Virol Methods 191:24–30. https://doi.org/10.1016/j.jviromet.2013.03.011.

35. Farkas K, Hassard F, McDonald JE, Malham SK, Jones DL. 2017. Evaluation of molecular methods for the detection and quantification of pathogen-derived nucleic acids in sediment. Front Microbiol 8:53. https://doi.org/10.3389/fmicb.2017.00053.

36. Van Dijk EL, Jaszczyszyn Y, Thermes C. 2014. Library preparation methods for next-generation sequencing: tone down the bias. Exp Cell Res 322:12–20. https://doi.org/10.1016/j.yexcr.2014.01.008.

37. Weiss S, Amir A, Hyde ER, Metcalf JL, Song SJ, Knight R. 2014. Tracking down the sources of experimental contamination in microbiome studies. Genome Biol 15:564. https://doi.org/10.1186/s13059-014-0564-2.

38. Strong MJ, Xu G, Morici L, Splinter Bon-Durant S, Baddoo M, Lin Z, Fewell C, Taylor CM, Flemington EK. 2014. Microbial contamination in next generation sequencing: implications for sequence-based analysis of clinical samples. PLoS Pathog 10:e1004437. https://doi.org/10.1371/journal.ppat.1004437.

39. Zhi N, Hu G, Wan Z, Zheng X, Liu X, Wong S, Kajigaya S, Zhao K, Young NS, Africa S. 2014. Correction for Xu et al., Hybrid DNA virus in Chinese patients with seronegative hepatitis discovered by deep sequencing. Proc Natl Acad Sci U S A 111:4344–4345. https://doi.org/10.1073/pnas.1402288111.

40. Zhi N, Hu G, Wong S, Zhao K, Mao Q, Young NS. 2014. Reply to Naccache et al: Viral sequences of NIH-CQV virus, a contamination of DNA extraction method. Proc Natl Acad Sci U S A 111:E977. https://doi.org/10.1073/pnas.1318965111.

41. Naccache SN, Greninger AL, Lee D, Coffey LL, Phan T, Rein-Weston A, Aronsohn A, Hackett J, Delwart EL, Chiu CY. 2013. The perils of pathogen discovery: origin of a novel parvovirus-like hybrid genome traced to nucleic acid extraction spin columns. J Virol 87:11966–11977. https://doi.org/10.1128/JVI.02323-13.

42. Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, Turner P, Parkhill J, Loman NJ, Walker AW. 2014. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. BMC Biol 12:87. https://doi.org/10.1186/s12915-014-0087-z.

43. Harris JP, Iturriza-Gomara M, O'Brien SJ. 2017. Re-assessing the total burden of norovirus circulating in the United Kingdom population. Vaccine 35:853–855. https://doi.org/10.1016/j.vaccine.2017.01.009.

44. Siebenga JJ, Vennema H, Zheng D, Vinjé J, Lee BE, Pang X, Ho ECM, Lim W, Choudekar A, Broor S, Halperin T, Rasool NBG, Hewitt J, Greening GE, Jin M, Duan Z, Lucero Y, O'Ryan M, Hoehne M, Schreier E, Ratcliff RM, White PA, Iritani N, Reuter G, Koopmans M. 2009. Norovirus illness is a global problem: emergence and spread of norovirus GII.4 variants, 2001–2007. J Infect Dis 200:802–812. https://doi.org/10.1086/605127.

45. Eden JS, Tanaka MM, Boni MF, Rawlinson WD, White PA. 2013. Recombination within the pandemic norovirus GII.4 lineage. J Virol 87:6270–6282. https://doi.org/10.1128/JVI.03464-12.

46. Cannon JL, Barclay L, Collins NR, Wikswo ME, Castro CJ, Magaña LC, Gregoricus N, Marine RL, Chhabra P, Vinjé J. 2017. Genetic and epidemiologic trends of norovirus outbreaks in the United States from 2013 to 2016 demonstrated emergence of novel GII.4 recombinant viruses. J Clin Microbiol 55:2208–2221. https://doi.org/10.1128/JCM.00455-17.

47. Cho HG, Lee SG, Mun SK, Lee MJ, Park PH, Jheong WH, Yoon MH, Paik SY. 2017. Detection of waterborne norovirus genogroup I strains using an improved real time RT-PCR assay. Arch Virol 162:3389–3396. https://doi.org/10.1007/s00705-017-3505-z.

48. International Organization for Standardization. 2013. Microbiology of food and animal feed—horizontal method for determination of hepatitis A virus and norovirus in food using real-time RT-PCR. Part 2. Method for qualitative detection. ISO/TS 15216 2:2013. ISO, Geneva, Switzerland. https://www.iso.org/standard/60297.html.

49. Jones MK, Grau KR, Costantini V, Kolawole AO, De Graaf M, Freiden P, Graves CL, Koopmans M, Wallet SM, Tibbetts SA, Schultz-Cherry S, Wobus CE, Vinjé J, Karst SM. 2015. Human norovirus culture in B cells. Nat Protoc 10:1939–1947. https://doi.org/10.1038/nprot.2015.121.

50. Ettayebi K, Crawford SE, Murakami K, Broughman JR, Karandikar U, Tenge VR, Neill FH, Blutt SE, Zeng X, Qu L, Kou B, Opekun AR, Burrin D,

Graham DY, Ramani S, Atmar RL, Estes MK. 2016. Replication of human noroviruses in stem cell-derived human enteroids. Science 353: 1387–1393. https://doi.org/10.1126/science.aaf5211.

51. Dunkin N, Weng S, Coulter CG, Jacangelo JG, Schwab KJ. 2017. Reduction of human norovirus GI, GII, and surrogates by peracetic acid and monochloramine in municipal secondary wastewater effluent. Environ Sci Technol 51:11918–11927. https://doi.org/10.1021/acs.est.7b02954.

52. Arredondo-Hernandez LJR, Diaz-Avalos C, Lopez-Vidal Y, Castillo-Rojas G, Mazari-Hiriart M. 2017. FRNA bacteriophages as viral indicators of faecal contamination in Mexican tropical aquatic systems. PLoS One 12: e0170399. https://doi.org/10.1371/journal.pone.0170399.

53. Harris JP, Jit M, Cooper D, Edmunds WJ. 2007. Evaluating rotavirus vaccination in England and Wales. Part I. Estimating the burden of disease. Vaccine 25:3962–3970. https://doi.org/10.1016/j.vaccine.2007.02.072.

54. Bawa Z, Elliot AJ, Morbey RA, Ladhani S, Cunliffe NA, O'Brien SJ, Regan M, Smith GE. 2015. Assessing the likely impact of a rotavirus vaccination program in England: the contribution of syndromic surveillance. Clin Infect Dis 61:77–85. https://doi.org/10.1093/cid/civ264.

55. Thomas SL, Walker JL, Fenty J, Atkins KE, Elliot AJ, Hughes HE, Stowe J, Ladhani S, Andrews NJ. 2017. Impact of the national rotavirus vaccination programme on acute gastroenteritis in England and associated costs averted. Vaccine 35:680–686. https://doi.org/10.1016/j.vaccine.2016.11.057.

56. Hungerford D, Read JM, Cooke RPD, Vivancos R, Iturriza-Gómara M, Allen DJ, French N, Cunliffe N. 2016. Early impact of rotavirus vaccination in a large paediatric hospital in the UK. J Hosp Infect 93:117–120. https://doi.org/10.1016/j.jhin.2015.12.010.

57. Wilhelm B, Waddell L, Greig J, Rajić A, Houde A, McEwen SA. 2015. A scoping review of the evidence for public health risks of three emerging potentially zoonotic viruses: hepatitis E virus, norovirus, and rotavirus. Prev Vet Med 119:61–79. https://doi.org/10.1016/j.prevetmed.2015.01.015.

58. Marton S, Dóró R, Fehér E, Forró B, Ihász K, Varga-Kugler R, Farkas SL, Bányai K. 2017. Whole genome sequencing of a rare rotavirus from archived stool sample demonstrates independent zoonotic origin of human G8P[14] strains in Hungary. Virus Res 227:96–103. https://doi.org/10.1016/j.virusres.2016.09.012.

59. Gautam R, Mijatovic-Rustempasic S, Roy S, Esona MD, Lopez B, Mencos Y, Rey-Benito G, Bowen MD. 2015. Full genomic characterization and phylogenetic analysis of a zoonotic human G8P[14] rotavirus strain detected in a sample from Guatemala. Infect Genet Evol 33:206–211. https://doi.org/10.1016/j.meegid.2015.05.004.

60. Steyer A, Naglič T, Jamnikar-Ciglenečki U, Kuhar U. 2017. Detection and whole-genome analysis of a zoonotic G8P[14] rotavirus strain isolated from a child with diarrhea. Genome Announc 5:e01053-17. https://doi.org/10.1128/genomeA.01053-17.

61. Medici MC, Tummolo F, Bonica MB, Heylen E, Zeller M, Calderaro A, Matthijnssens J. 2015. Genetic diversity in three bovine-like human G8P[14] and G10P[14] rotaviruses suggests independent interspecies transmission events. J Gen Virol 96:1161–1168. https://doi.org/10.1099/vir.0.000055.

62. Cook N, Bridger J, Kendall K, Gomara MI, El-Attar L, Gray J. 2004. The zoonotic potential of rotavirus. J Infect 48:289–302. https://doi.org/10.1016/j.jinf.2004.01.018.

63. Farkas K, Cooper DM, McDonald JE, Malham SK, de Rougemont A, Jones DL. 2018. Seasonal and spatial dynamics of enteric viruses in wastewater and in riverine and estuarine receiving waters. Sci Total Environ 634: 1174–1183. https://doi.org/10.1016/j.scitotenv.2018.04.038.

64. Da Silva AK, Le Saux JC, Parnaudeau S, Pommepuy M, Elimelech M, Le Guyader FS. 2007. Evaluation of removal of noroviruses during wastewater treatment, using real-time reverse transcription-PCR: different behaviors of genogroups I and II. Appl Environ Microbiol 73:7891–7897. https://doi.org/10.1128/AEM.01428-07.

65. Svraka S, Duizer E, Vennema H, De Bruin E, Van Der Veer B, Dorresteijn B, Koopmans M. 2007. Etiological role of viruses in outbreaks of acute gastroenteritis in the Netherlands from 1994 through 2005. J Clin Microbiol 45:1389–1394. https://doi.org/10.1128/JCM.02305-06.

66. Loisy F, Atmar RL, Guillon P, Le Cann P, Pommepuy M, Le Guyader FS. 2005. Real-time RT-PCR for norovirus screening in shellfish. J Virol Methods 123:1–7. https://doi.org/10.1016/j.jviromet.2004.08.023.

67. Kageyama T, Kojima S, Shinohara M, Uchida K, Fukushi S, Hoshino FB, Takeda N, Katayama K. 2003. Broadly reactive and highly sensitive assay for Norwalk-like viruses based on real-time quantitative reverse transcription-PCR. J Clin Microbiol 41:1548–1557. https://doi.org/10.1128/JCM.41.4.1548-1557.2003.

68. Chan MCW, Sung JJY, Lam RKY, Chan PKS, Lai RWM, Leung WK. 2006. Sapovirus detection by quantitative real-time RT-PCR in clinical stool specimens. J Virol Methods 134:146–153. https://doi.org/10.1016/j.jviromet.2005.12.013.

69. Costafreda MI, Bosch A, Pintó RM. 2006. Development, evaluation, and standardization of a real-time TaqMan reverse transcription-PCR assay for quantification of hepatitis A virus in clinical and shellfish samples. Appl Environ Microbiol 72:3846–3855. https://doi.org/10.1128/AEM.02660-05.

70. Jothikumar N, Cromeans TL, Robertson BH, Meng XJ, Hill VR. 2006. A broadly reactive one-step real-time RT-PCR assay for rapid and sensitive detection of hepatitis E virus. J Virol Methods 131:65–71. https://doi.org/10.1016/j.jviromet.2005.07.004.

71. Farkas K, Peters DE, McDonald JE, de Rougemont A, Malham SK, Jones DL. 2017. Evaluation of two triplex one-step qRT-PCR assays for the quantification of human enteric viruses in environmental samples. Food Environ Virol 9:342–349. https://doi.org/10.1007/s12560-017-9293-5.

72. Kingsley RA, Msefula CL, Thomson NR, Kariuki S, Holt KE, Gordon MA, Harris D, Clarke L, Whitehead S, Sangal V, Marsh K, Achtman M, Molyneux ME, Cormican M, Parkhill J, MacLennan CA, Heyderman RS, Dougan G. 2009. Epidemic multiple drug resistant Salmonella Typhimurium causing invasive disease in sub-Saharan Africa have a distinct genotype. Genome Res 19:2279–2287. https://doi.org/10.1101/gr.091017.109.

73. Hennechart-Collette C, Martin-Latil S, Guillier L, Perelle S. 2015. Determination of which virus to use as a process control when testing for the presence of hepatitis A virus and norovirus in food and water. Int J Food Microbiol 202:57–65. https://doi.org/10.1016/j.ijfoodmicro.2015.02.029.

74. Reysenbach A, Pace N. 1995. Reliable amplification of hyperthermophilic archaeal 16S rRNA genes by the polymerase chain reaction, p 101–107. In Robb F, Place A (ed), Archaea: a laboratory manual. Cold Spring Harbor Laboratory Press, New York, NY.

75. Turner S, Pryer KM, Miao VP, Palmer JD. 1999. Investigating deep phylogenetic relationships among cyanobacteria and plastids by small subunit rRNA sequence analysis. J Eukaryot Microbiol 46:327–338. https://doi.org/10.1111/j.1550-7408.1999.tb04612.x.

76. Amaral-Zettler LA, McCliment EA, Ducklow HW, Huse SM. 2009. A method for studying protistan diversity using massively parallel sequencing of V9 hypervariable regions of small-subunit ribosomal RNA Genes. PLoS One 4:e6372. https://doi.org/10.1371/journal.pone.0006372.

77. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. EMBnet.journal 17:10–12. https://doi.org/10.14806/ej.17.1.200.

78. Joshi N, Fass J. 2011. Sickle: a sliding-window, adaptive, quality-based trimming tool for FastQ files (version 1.33). https://github.com/najoshi/sickle.

79. Schmieder R, Edwards R. 2011. Quality control and preprocessing of metagenomic datasets. Bioinformatics 27:863–864. https://doi.org/10.1093/bioinformatics/btr026.

80. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nat Methods 9:357–359. https://doi.org/10.1038/nmeth.1923.

81. Nurk S, Bankevich A, Antipov D, Gurevich A, Korobeynikov A, Lapidus A, Prjibelsky A, Pyshkin A, Sirotkin A, Sirotkin Y, Stepanauskas R, McLean J, Lasken R, Clingenpeel SR, Woyke T, Tesler G, Alekseyev MA, Pevzner PA. 2013. Assembling genomes and mini-metagenomes from highly chimeric reads, p 158–170. In Deng M, Jiang R, Sun F, Zhang X (ed), Research in computational molecular biology. RECOMB 2013. Lecture notes in computer science. Springer, Berlin, Heidelberg.

82. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: architecture and applications. BMC Bioinformatics 10:421. https://doi.org/10.1186/1471-2105-10-421.

83. Brister JR, Ako-adjei D, Bao Y, Blinkova O. 2015. NCBI viral genomes resource. Nucleic Acids Res 43:D571–D577. https://doi.org/10.1093/nar/gku1207.

84. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O, Brover V, Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D, Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, McGarvey KM, Murphy MR, O'Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD, Schoch C, Shkeda A, Storz SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C, Webb D, Wu W, Landrum MJ, Kimchi A, Tatusova T, DiCuccio M, Kitts P, Murphy TD, Pruitt KD. 2016. Reference sequence (RefSeq) database at

NCBI: current status, taxonomic expansion, and functional annotation. Nucleic Acids Res 44:D733–D745. https://doi.org/10.1093/nar/gkv1189.

85. Racine JS. 2012. RStudio: a platform-independent IDE for R and Sweave. J Appl Econ 27:167–172. https://doi.org/10.1002/jae.1278.

86. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics 30:2068–2069. https://doi.org/10.1093/bioinformatics/btu153.

87. Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Res 32:1792–1797. https://doi.org/10.1093/nar/gkh340.

88. Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets. Mol Biol Evol 33:msw054. https://doi.org/10.1093/molbev/msw054.

89. Okonechnikov K, Golosova O, Fursov M, UGENE Team. 2012. Unipro UGENE: a unified bioinformatics toolkit. Bioinformatics 28:1166–1167. https://doi.org/10.1093/bioinformatics/bts091.

90. Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS. 2009. MEME Suite: tools for motif discovery and searching. Nucleic Acids Res 37:W202–W208. https://doi.org/10.1093/nar/gkp335.

91. Bailey TL, Johnson J, Grant CE, Noble WS. 2015. The MEME Suite. Nucleic Acids Res 43:W39–W49. https://doi.org/10.1093/nar/gkv416.

92. Sullivan MJ, Petty NK, Beatson SA. 2011. Easyfig: a genome comparison visualizer. Bioinformatics 27:1009–1010. https://doi.org/10.1093/bioinformatics/btr039.