

1 To be submitted to [Genome Biology](#)

2 **Enhancing Infectious Intestinal Disease diagnosis through metagenomic and**
3 **metatranscriptomic sequencing of over 1000 human diarrhoeal samples**

4 Edward Cunningham-Oakes^{1,2}, Blanca M. Perez-Sepulveda³, Yan Li³, Jay C. D. Hinton³,
5 Charlotte A. Nelson⁴, K. Marie McIntyre⁵, Maya Wardeh^{6,7,8}, Sam Haldenby⁴, Richard
6 Gregory⁴, Miren Iturriza-Gómara^{3,9}, Christiane Hertz-Fowler⁴, Sarah J. O'Brien⁵, Nigel A.
7 Cunliffe^{2,3}, Alistair C. Darby^{1,2,4*}; on behalf of the INTEGRATE consortium

8 ¹Institute of Infection, Veterinary & Ecological Sciences, University of Liverpool, Liverpool, UK

9 ²NIHR Health Protection Research Unit in Gastrointestinal Infections, Liverpool, UK

10 ³Department of Clinical Infection, Microbiology and Immunology, Institute of Infection,
11 Veterinary and Ecological Sciences, University of Liverpool, Liverpool, UK

12 ⁴Centre for Genomic Research, Institute of Systems, Molecular and Integrative Biology,
13 University of Liverpool, Liverpool, UK

14 ⁵School of Natural and Environmental Sciences, Newcastle University, Newcastle, UK

15 ⁶Department of Computer Science, University of Liverpool, Liverpool, UK

16 ⁷Department of Livestock and One Health, Institute of Infection, Veterinary and Ecological
17 Sciences, University of Liverpool, Liverpool, UK

18 ⁸Department of Mathematics, University of Liverpool, Liverpool, UK

19 ⁹Centre for Vaccine Innovation and Access, Program for Appropriate Technology in Health
20 (PATH), Geneva, 1218, Switzerland

21 *For correspondence:

22 E-mail: Alistair.Darby@liverpool.ac.uk

23 **Abstract**

24 Fundamental issues in the traditional surveillance of diarrhoeal disease need to be addressed.
25 The limitations of traditional microbiological diagnostic methods often mean that the cause of
26 diarrhoea remains unknown, especially for novel or difficult-to-isolate pathogens. Sequencing

27 samples directly, without isolating pathogens, would address this issue. However, we must
28 ensure that sequencing also captures pathogens that are detectable using current diagnostic
29 methods.

30 We show that metagenomic and metatranscriptomic approaches can effectively detect nine
31 gastrointestinal pathogens in the UK. Metatranscriptomics shows increased sensitivity of
32 detection for pathogens like *Campylobacter*, *Clostridioides difficile*, *Cryptosporidium* and
33 *Giardia*, while metagenomics is more effective for detecting pathogens such as *Adenovirus*,
34 pathogenic *Escherichia coli*, *Salmonella*, *Shigella*, and *Yersinia enterocolitica*. Certain
35 pathogens were detected by both metagenomic and metatranscriptomic sequencing.
36 Metatranscriptomics gave near-complete genome coverage for Human mastadenovirus F and
37 detected *Cryptosporidium* via capture of *Cryptosporidium parvum* virus (CSpV1). A
38 comprehensive transcriptomic profile of *Salmonella* Enteritidis was recovered from the stool
39 of a patient with a laboratory-confirmed *Salmonella* infection.

40 This study highlights the power of direct sequencing of human samples to augment GI
41 pathogen surveillance and clinical diagnostics. Metatranscriptomics was best for capturing a
42 wide breadth of pathogens and was more sensitive for this purpose. We propose that
43 metatranscriptomics should be considered for future surveillance of gastrointestinal
44 pathogens. This study has generated a rich data resource of paired metagenomic and
45 metatranscriptomic datasets, direct from over 1000 patient stool samples. We have made
46 these data publicly available to promote the improved understanding of pathogens associated
47 with infectious intestinal diseases.

48 **Keywords**

49 microbiome, culture-independent, metagenome, metatranscriptome, diagnostics, genomics,
50 pathogens.

51 **Background**

52 The incidence of infectious intestinal disease (or acute gastroenteritis) is estimated to be 18
53 million cases each year in the United Kingdom (UK)[1]. About 25% of infected people
54 experience diarrhoeal and related gastrointestinal symptoms. The current mainstay for
55 identifying gastrointestinal pathogens in faecal specimens in the UK are conventional
56 laboratory techniques, including microscopy and antigen detection, and increasingly,
57 molecular assays such as nucleic acid amplification[2].

58 Although conventional and polymerase chain-reaction (PCR)-based approaches (such as
59 BioFire Panels) are validated for clinical laboratory use[2], both focus on a single gene or set
60 of characteristics providing limited information about pathogens[3]. In the case of bacterial
61 culture, the time required for growth, lack of sensitivity, and the challenge of culturing fastidious
62 organisms cause diagnostic delays[3]. Current methods lack the sensitivity required to detect
63 pathogens that are present intermittently or in low numbers[4]. In contrast, PCR-based
64 methods use target sequences for organism detection, resulting in increased sensitivity and
65 no strict requirement for the prior growth of organisms[3].

66 Whilst PCR-based methods are more sensitive than conventional (traditional) methods, they
67 lack resolution[5] and are unable to achieve the strain-level discrimination required for
68 outbreak monitoring[6]. Inevitably, molecular assays target known genes from well-
69 characterised organisms [5], meaning that unexpected pathogens and unique genes will be
70 missed. Whole-genome sequencing partly overcomes this, but still requires the isolation of a
71 pure culture of pathogenic organisms.

72 The speed and sensitivity of metagenomic and metatranscriptomic data analysis[7] has been
73 significantly enhanced by k-mer-based methods, an approach that has been widely adopted
74 in many popular workflows[8,9] to identify pathogens in metagenomic samples through
75 database matching. The computational efficiency of k-mers is ideal for high-throughput
76 sequencing applications[10]. However, it is important to note that sequencing errors and the
77 comprehensiveness of the databases used[11] can influence the effectiveness of k-mer-based
78 approaches.

79 It has been proposed that DNA and RNA sequencing of clinical samples could be a valuable
80 future approach[5]. To ensure that the presence of pathogens is recorded accurately, it is
81 essential to understand the strengths and limitations of metagenomic and metatranscriptomic
82 approaches.

83 The INTEGRATE study[12] compared traditional (culture, ELISA, microscopy and PCR)
84 diagnostic methods with state-of-the-art, sensitive molecular and genome-based
85 microbiological methods for identifying and characterising causative pathogens[12]. Here, we
86 present data generated by next-generation sequencing of the stool microbiomes of 1,067
87 patients with symptoms of gastroenteritis. This dataset represented a unique opportunity to
88 explore the effectiveness of k-mer-based analyses using a large number of samples that were
89 also characterised with gold-standard clinical laboratory diagnostics. We considered the
90 comparative benefits of different sequencing types in various scenarios (right test, right time,
91 right patient).

92 We use these data to show that both metagenomic (DNA) and metatranscriptomic (RNA)
93 sequencing directly from stool can detect the major community-associated gastrointestinal
94 (GI) pathogens in the United Kingdom. We found that metagenomic and metatranscriptomic

95 sequencing have distinctive features for pathogen detection, and discovered that
96 metatranscriptomics offers unexpected benefits for pathogen surveillance.

97 All of these data have been made publicly available (PRJEB62473) to provide a rich data
98 source for researchers to foster a deeper understanding of the pathogens associated with
99 infectious intestinal diseases.

100 **Results**

101 ***Metagenomics and metatranscriptomics show different levels of sensitivity for GI*** 102 ***pathogens***

103 The DNA and RNA from a total of 1,067 samples were sequenced, with 985 providing both
104 metagenomic and metatranscriptomic data (see Supplementary File 1 for all k-mer counts and
105 associated taxonomy from these samples). For *Campylobacter*, *Cryptosporidium*, and *Giardia*
106 (Figure 1), metatranscriptomics showed greater sensitivity than metagenomics (see Methods
107 for definition of sensitivity). In contrast, metagenomics displayed greater sensitivity than
108 metatranscriptomics for the *Adenoviridae*, *Clostridium difficile*, pathogenic *Escherichia coli*,
109 *Salmonella*, *Shigella* and *Yersinia enterocolitica* (Figure 1). *Entamoeba histolytica* were not
110 detected in either the metagenomic or metatranscriptomic datasets.

111 ***The detection of GI pathogens in metagenomic and metatranscriptomic data mirrors*** 112 ***clinical laboratory results***

113 Our analysis showed that the pathogens detected in sequencing reads closely match results
114 generated by laboratory diagnostics for Adenovirus, *C. difficile*, *Campylobacter*,
115 *Cryptosporidium*, Norovirus, Rotavirus, *Salmonella*, Sapovirus, *Shigella*, and *Y. enterocolitica*
116 (Figure 2). Most major GI community pathogens in the United Kingdom were detected in both
117 metagenomic and metatranscriptomic data, but RNA viruses could only be detected by

118 metatranscriptomics. A summary of the “traditional” methods used for pathogen diagnosis in
119 the INTEGRATE study is presented in Table 1.

120 ***Viral pathogens***

121 DNA viruses such as Adenovirus were detected in both the metagenomic and
122 metatranscriptomic datasets. For Adenovirus, positive correlations were observed between
123 detection in metagenomic reads, metatranscriptomic reads, and Luminex xTAG
124 Gastrointestinal Pathogen Panel (Luminex) results ($p < 0.001$). The metatranscriptomic results
125 correlated positively with Rotavirus ($p < 0.001$). The detection of Norovirus and Sapovirus by
126 metatranscriptomics was significantly correlated ($p < 0.001$) with the Luminex results.
127 Metagenomic and metatranscriptomic results did not correlate with the detection of Astrovirus
128 using Traditional and Luminex methods.

129 ***Protists***

130 Protists were detected by both metagenomics and metatranscriptomics. However, the
131 metatranscriptomic results had a much higher sensitivity for the detection of parasites than
132 metagenomics. Positive correlations between the detection of *Cryptosporidium* in
133 metatranscriptomic data and laboratory data were highly significant ($p < 0.001$). No
134 associations were observed between metagenomic data and laboratory results for
135 *Cryptosporidium*. There was no correlation between the detection of *Giardia* using Traditional
136 or Luminex methods, and detecting *Giardia* using metagenomics or metatranscriptomics.

137 ***Bacterial pathogens***

138 The identification of bacterial pathogens from sequencing data is challenging, as commensal
139 organisms and pathogens can have extremely high levels of genomic similarity. Laboratory
140 diagnostics tend to differentiate commensal and pathogenic organisms using genes or

141 phenotypes associated with pathogenicity. Our results show that metagenomics and
142 metatranscriptomics can both identify bacterial pathogens with differing sensitivities. For
143 *Campylobacter*, positive correlations were observed between direct sequencing and all
144 laboratory results ($p < 0.001$). *Salmonella* displayed positive correlations between sequencing
145 data and both Traditional ($p < 0.001$) and Luminex ($p < 0.25$) diagnostics. *C. difficile*
146 metatranscriptomic sequencing data positively correlated with both Traditional ($p < 0.25$) and
147 Luminex ($p < 0.001$) diagnostics. *Y. enterocolitica* sequencing data positively correlated with
148 Luminex results as follows: *C. difficile* metatranscriptomic reads ($p < 0.001$), *Y. enterocolitica*
149 metatranscriptomic reads ($p < 0.01$), and *Y. enterocolitica* metagenomic reads ($p < 0.001$).

150 *E. coli* and *Shigella* are closely-related species; detection of *Shigella* in metagenomic data
151 correlated positively with traditional and Luminex diagnostics ($p < 0.25$), while *E. coli* showed a
152 non-significant correlation ($p > 0.25$). *Vibrio cholerae* were not be detected in either
153 metagenomic or metatranscriptomic data, consistent with laboratory diagnostics, which
154 identified no *V. cholerae* infections.

155 A summary of all correlations between the detection of GI pathogens in sequencing reads and
156 laboratory data and their significance is provided in Supplementary Files 2 and 3.

Test Parameter	Liverpool	Manchester	Preston
Adenovirus 41/42	PCR	PCR	Immunoassay
Rotavirus A	PCR	PCR	Immunoassay
Norovirus GI/GII	PCR	PCR	Immunoassay & PCR
Sapovirus	PCR	PCR	Not Available
<i>Clostridioides difficile</i> toxin A/B & GDH	Immunoassay	Immunoassay	Immunoassay
<i>Salmonella</i>	Culture	Culture	Culture
<i>Shigella</i>	Culture	Culture	Culture
<i>Campylobacter</i> (<i>C. jejuni</i> , <i>C. coli</i> , <i>C. lari</i>)	Culture	Culture	Culture
<i>E. coli</i> O157	Culture	Culture	Culture
Enterotoxigenic <i>E. coli</i> (EPEC) LT/ST	Not Available	Not Available	Not Available
Enterotoxigenic <i>E. coli</i> (EPEC) EA/EC	Not Available	Not Available	Not Available
Enterotoxigenic <i>E. coli</i> (EPEC) EA/EC	Not Available	Not Available	Not Available
<i>Yersinia enterocolitica</i>	Culture	Culture	Not Available
<i>Vibrio cholerae</i>	Culture	Culture	Culture
Shigella like toxin producing <i>E. coli</i> (STEC)	Not Available	Not Available	Not Available
<i>Giardia lamblia</i>	Microscopy	Immunoassay	Immunoassay
<i>Cryptosporidium</i>	Microscopy	Immunoassay	Immunoassay
<i>Entamoeba histolytica</i>	Microscopy	Microscopy	Microscopy

Table 1: Summary of “traditional” methods used at clinical laboratories during the INTEGRATE study. Samples were processed via routine diagnostic pathways at each laboratory involved in the study (see Supplementary File 5). Traditional assays for Enterotoxigenic and Enterotoxigenic *E. coli*, as well as *E. coli* O157, were not available (only available in the Luminex xTAG GPP panel).

157 **Case-studies for the use of metatranscriptomics in pathogen surveillance**

158 ***Complete genomes from diarrhoeal-associated Adenovirus can be detected in both***
159 ***metagenomic and metatranscriptomic data***

160 Whilst Adenovirus is a DNA virus, it was surprising to see that Adenovirus could also be
161 detected in RNA. There was a strong correlation between the detection of Adenovirus in
162 metagenomic and metatranscriptomic data, and detection using Luminex methods (see
163 Supplementary File 2 and Figure 2). Mapping of Adenovirus-associated reads to the Human
164 mastadenovirus F genome showed that, in 7 out of 9 samples, metagenomics generated more
165 complete genes, at a higher depth (Supplementary File 4). However, in 2 out of 9 samples
166 (Samples 5638 and 6985) near-complete genome coverage was achieved using both
167 metagenomics and metatranscriptomics (83.9-100% - Supplementary File 4 and Figure 3).
168 These results demonstrate the potential of metatranscriptomics to directly capture the virome
169 from clinical samples, including DNA viruses relevant to the condition of interest.

170 ***Cryptosporidium-associated RNA viruses facilitate detection directly from stool***

171 Another interesting observation was the correlation ($p < 0.001$, see Figure 2) between the
172 detection of *Cryptosporidium* using metatranscriptomics and the detection of *Cryptosporidium*
173 in the laboratory. In contrast, detecting *Cryptosporidium* using metagenomics did not correlate
174 with laboratory results. Mapping revealed that *Cryptosporidium* was accurately identified in
175 metatranscriptomic data due to the presence of *Cryptosporidium parvum* virus (CSpV1), which
176 is an RNA virus. CSpV1 was identified in 33 metatranscriptomic samples (Table 2). Of these
177 33 samples, 9 received a positive result using Traditional methods, whilst 16 were positive by
178 Luminex. CSpV1 received a high-confidence score (0.995) in 21 out of the 33 samples, with
179 percentage breadth of genome coverage ranging from 57.1-100%. This illustrates the potential
180 of CSpV1 to be used as a reliable biomarker for human *Cryptosporidium* infection.

Study ID	Traditional result	Luminex result	TAXID	Read count	Coverage breadth (%)	Coverage depth (fold)	PanGIA score
238	0	1	675060.1	24334	0.9937	1057.6231	0.995
299	0	1	675060.1	17324	0.9907	742.9534	0.995
347	0	1	675060.1	18914	0.9513	820.9334	0.995
1530	NA	0	675060.1	132778	0.9997	5789.0705	0.995
1730	0	1	675060.1	277178	1	12099.0436	0.995
1868	1	1	675060.1	25942	0.9836	1121.2385	0.995
1996	NA	0	675060.1	490	0.9659	21.2555	0.995
2237	0	0	675060.1	54	0.5705	2.3751	0.995
2270	1	1	675060.1	2481486	1	108131.0308	0.995
4580	0	0	675060.1	136	0.8757	5.711	0.995
4667	1	1	675060.1	45590	1	1975.986	0.995
4922	1	1	675060.1	11154	0.9438	474.8819	0.995
5019	1	1	675060.1	27990	0.997	1230.6458	0.995
5195	0	0	675060.1	526	0.8951	22.775	0.995
5215	1	1	675060.1	43408	0.9997	1895.0209	0.995
5563	0	1	675060.1	85048	0.9988	3700.396	0.995
5675	0	1	675060.1	8338	0.9913	360.4002	0.995
6446	1	1	675060.1	2200	0.9668	94.7941	0.995
6602	0	1	675060.1	5770	0.9949	250.243	0.995
6912	1	1	675060.1	21542	0.8655	932.6551	0.995
7233	1	1	675060.1	22898	1	1005.3114	0.995
1817	NA	0	675060.1	14	0.2624	0.601	0.601
1548	0	0	675060.1	12	0.2585	0.535	0.535
1111	NA	0	675060.1	10	0.2409	0.4471	0.4471
769	NA	0	675060.1	10	0.2445	0.439	0.439
1436	NA	0	675060.1	6	0.1491	0.2687	0.2687
54	0	0	675060.1	6	0.1626	0.2469	0.2469
127	0	NA	675060.1	4	0.0819	0.1793	0.1793
6890	NA	0	675060.1	4	0.1136	0.179	0.179
5734	0	0	675060.1	4	0.0562	0.1787	0.1787
360	0	0	675060.1	4	0.1482	0.1781	0.1781
2279	0	0	675060.1	4	0.0855	0.1614	0.1614
369	0	0	675060.1	4	0.1384	0.1599	0.1599

Table 2: Identification of CSpV1 in metatranscriptomic data in comparison to results from *Cryptosporidium* laboratory diagnostics. For both Traditional and Luminex results, NA represents instances where a diagnostic test could not be performed.

181 **Generation of a complete transcriptomic profile for *Salmonella***

182 Metatranscriptomic analysis of stool from a patient with a laboratory-confirmed *Salmonella*
183 infection yielded functional insights that cannot be achieved with Traditional and Luminex
184 diagnostics. Transcriptomic analysis reveals gene expression patterns that key biological
185 processes in bacteria. The high-quality transcriptomic profile was generated from 12.7 million
186 sequence reads that mapped to the genome of *S. enterica* serovar Enteritidis PT4 strain
187 P125109. The *S. Enteritidis* transcripts from this novel gene expression data can be visualised
188 and interrogated in a bespoke genome browser (https://s.hintonlab.com/study_74).

189 A variety of environmentally responsive *Salmonella* genes were highly expressed (as defined
190 by Kröger *et al.* 2013; *Cell Host Microbe* [13]), likely reflecting physicochemical stresses the
191 bacteria had been exposed to in the stool sample. Examples include *ahpC* (oxidative stress),
192 *hmpA* (nitrosative stress), *phoH* (phosphate starvation), *pspA* (extracytoplasmic stress), and
193 the *rpoE* and *rpoS* transcription factor genes, as can be seen with the [SalComMac data](#)
194 [visualisation tool](#). The unexpected discovery that the metatranscriptomic analysis of a human
195 stool sample can generate a comprehensive gene expression profile of a *Salmonella* pathogen
196 is worthy of future exploitation.

197 **Discussion**

198 We have demonstrated that metagenomic and metatranscriptomic approaches provide
199 agnostic detection of important UK GI pathogens from human stool. The primary impact of this
200 work lies within GI pathogen diagnostics. Our findings demonstrate the potential for improving
201 current GI pathogen diagnostics and the bridging of gaps not addressed by standard
202 approaches.

203 ***Improvements within the scope of current diagnostics***

204 Sequencing directly from stool could minimise the time required for pathogen detection,
205 allowing more laborious detection methods such as cultivation to be appropriately tailored to
206 confirm the presence of the suspected pathogens.

207 The metatranscriptomic strategy displays increased sensitivity for *Campylobacter*, *C. difficile*,
208 *Cryptosporidium* and *Giardia*, whilst metagenomics displayed increased sensitivity for other
209 GI pathogens including Adenovirus, pathogenic *E. coli*, *Salmonella*, *Shigella*, and *Y.*
210 *enterocolitica*. Direct extraction of RNA from stool represents a single sample format and
211 cultivation-independent process for detecting a broad range of GI pathogens, including
212 unexpected aetiological agents and those that cannot be detected by metagenomic
213 sequencing, such as RNA viruses. The observation of near-complete genome coverage for
214 Human mastadenovirus F in both the metagenome and metatranscriptome highlights the
215 potential to optimise metatranscriptomic sequencing from stool to capture the virome,
216 including DNA virus transcriptomes relevant to clinical conditions. This finding is supported by
217 previous clinical studies, which used metatranscriptomics to simultaneously measure the
218 virome, microbiome, and host response[14]. Our data and previous studies[15] confirm the
219 ability to characterise disease-related microbiomes with increased sensitivity via
220 metatranscriptomics.

221 Increased sensitivity for the detection of protists of concern in GI infections was also
222 demonstrated. Our visualisations of metagenomic and metatranscriptomic reads (Figure 1)
223 showed that metatranscriptomic data provide greater sensitivity for detecting *Cryptosporidium*
224 and *Giardia* (protists). Finally, our multivariable model demonstrated the strong correlation and
225 high significance between the detection of *Cryptosporidium* in the laboratory, and in
226 metatranscriptomic data, a finding that was supported by a previous study that detected 23%
227 more blood infections than traditional methods[16]. In these data, the presence of
228 *Cryptosporidium*-associated viruses increases the sensitivity of detection for this particular

229 protist. In contrast, the detection of *Cryptosporidium* using metagenomics appears to be
230 spurious. This virus has recently been reported in various subtypes of *Cryptosporidium*
231 *parvum* from diarrhoeic farm animals[17,18], but it is not currently used as a diagnostic marker
232 in humans. These results highlight the advantages of metatranscriptomics for *Cryptosporidium*
233 surveillance, where the use of metagenomics alone could result in missed identification. This
234 suggests that RNA viruses could be considered sensitive biomarkers for *Cryptosporidium* and
235 other protists.

236 Overall, our findings reveal that RNA is a valuable diagnostic target for the detection of
237 pathogens of low abundance and reduces false-positive signals from commensals. Our
238 approach could influence the future allocation of resources for reference laboratory
239 diagnostics.

240 ***Bridging gaps not addressed by current diagnostics***

241 Metatranscriptomic data could fill gaps in areas of clinical relevance that are not fulfilled by
242 routine clinical diagnostics. Firstly, metagenomic and metatranscriptomic data permits the
243 identification of multiple species and strains within a sample (Supplementary Figure 1,
244 Supplementary File 4), including novel pathogens. Such analysis is beyond the scope of our
245 study, but has been used to successfully identify novel pathogens from the stools of various
246 mammalian species[19,20]. Additionally, we have demonstrated the ability to rapidly generate
247 gene expression profiles for pathogens of concern, without prior enrichment. Finally, we have
248 generated illuminating metatranscriptomic data from a human diarrhoeal sample. Future
249 studies could generate true disease-state expression profile by using appropriate
250 methodology. From a clinical perspective, the use of metagenomic and metatranscriptomic
251 sequencing has the potential to reveal the effects of interventions[21] and to accurately
252 investigate host-pathogen dynamics during genuine human infections[14].

253 **Limitations**

254 In certain scenarios, metagenomic sequencing captures more information than
255 metatranscriptomic sequencing. For DNA viruses, while it is possible to capture expression
256 profiles, optimisation is needed to improve this process. Our data demonstrate that the
257 underlying biology can be captured, but further refinement is necessary. Additionally, *E.*
258 *histolytica* was not captured by metagenomic or metatranscriptomic approaches, a finding that
259 requires further investigation.

260 Future adaptation of our workflow is needed for the accurate identification of *E. coli*
261 pathovariants from sequencing data. *Shigella* and *E. coli* pathovariants are extremely similar
262 on a genome-wide (and taxonomic) level[22], and are currently distinguished using specific
263 gene-based assays[23]. In contrast, our study drew correlations between pathogens in reads
264 and laboratory tests based on taxonomy. Due to this approach, and the ubiquitous presence
265 of *E. coli* in all stool samples, it was not possible to associate the presence of *E. coli*, and
266 gene-based assays used for *E. coli* pathovariant identification (Supplementary Table S1,
267 Supplementary File 4). This may explain the limited overlap between *Shigella* sequencing
268 reads and laboratory tests (Figure 2), due to the conflation of *Shigella* with *E. coli* (and *vice*
269 *versa*) in our current analysis. Future work should also validate this approach on a range of
270 sample types (beyond stool) to ensure robustness and reliability across different clinical
271 scenarios.

272 **Perspective**

273 With sufficient benchmarking, the diagnosis of various GI pathogens can be achieved from
274 clinical samples without culturing. Metatranscriptomics can detect active DNA viruses and
275 enhance sensitivity for protists by using RNA viruses as biomarkers. Perhaps the value of
276 clinical metagenomics has been overstated, and metatranscriptomics could offer a
277 comprehensive approach to both detect disease-relevant pathogens and understand their
278 biology.

279 To our knowledge, this study is the first to demonstrate and quantify the benefits that
280 metatranscriptomics could bring to gastrointestinal surveillance in the United Kingdom by
281 direct comparison of all major community pathogens to validated diagnostics. Our study lays
282 the groundwork for the implementation of sequencing-based diagnostics in clinical settings,
283 with the potential to detect a broader range of organisms than current approaches and to
284 identify novel pathogens.

285 **Materials and methods**

286 ***Patient recruitment and sample collection***

287 Recruitment and sample collection was described previously[12]. Briefly, stool was collected
288 from 1,067 members of the public with symptoms of acute gastroenteritis via practices in the
289 Royal College of General Practitioners Research and Surveillance Centre National Monitoring
290 Network (RCGP RSC NMN). Patients meeting inclusion criteria were invited to submit a stool
291 sample for microbiological analysis. Consent was obtained for this procedure, as stool
292 sampling is usually only performed if a case is severe, or persistent. Patients who provided a
293 stool sample were then recruited into the study.

294 ***Sample processing***

295 Faecal samples were received by one of three clinical laboratories (Royal Liverpool and
296 Broadgreen University Hospitals NHS Trust, Central Manchester University Hospitals NHS
297 Foundation Trust, or Lancashire Teaching Hospitals NHS Foundation Trust), and divided into
298 two aliquots. One part of the sample was processed using traditional methods (culture, ELISA,
299 microscopy or PCR with no additional hybridisation probe – see Supplementary File 5) at each
300 laboratory; the other was processed using a combined molecular multiplex real-time
301 polymerase chain reaction (PCR) and target-specific hybridisation probe [Luminex xTAG
302 Gastrointestinal Pathogen Panel, Luminex, I032C0324], supplemented with targets for
303 Enterococcal *Escherichia coli* and Sapovirus. Nucleic acid extraction from faeces was

304 performed using QIASymphany and EasyMag automated nucleic acid extraction platforms.
305 Further details can be found in the primary study protocol[12]. Samples that returned a positive
306 result according to routine clinical practice were designated as “clinical positive”. Those that
307 returned a positive result by Luminex were designated as “molecular positive”.

308 ***Metagenomic and metatranscriptomic sequencing***

309 Illumina fragment libraries from DNA were prepared using NEBNext DNA Ultra kits. For the
310 generation of dual-indexed, strand-specific RNASeq libraries, RiboZero rRNA depleted RNA
311 samples were prepared using NEBNext Ultra Directional RNA kits. For all libraries, paired-
312 end, 150-bp sequencing was then performed on an Illumina HiSeq 4000, generating data from
313 >280 million clusters per lane.

314 ***Quality control for second-generation sequencing reads***

315 Modules from the MetaWRAP[24] (v1.3.2) pipeline were used to standardise metagenome
316 analysis. The pipeline was deployed in a dedicated Conda environment, using the “manual
317 installation” guide (see [GitHub](#)). All paired-end reads underwent quality-control using the
318 MetaWRAP “read_qc” module to remove low-quality, adapter, and human sequence reads.
319 The T2T consortium complete human genome, (GCF_009914755.1) and human
320 mitochondrial genome (NC_012920.1) were used as references for the removal of human
321 reads.

322 ***Assigning taxonomy to genomic DNA and RNA reads and assessing microbiome*** 323 ***diversity***

324 DNA and RNA reads were used for taxonomic assignments with Kraken2[8] (v2.1.2), using a
325 custom database, which included all RefSeq complete genomes and proteins for archaea,
326 bacteria, fungi, viruses, plants, protozoa, as well as all complete RefSeq plasmid nucleotide
327 and protein sequences, and a false-positive minimised version of the NCBI UniVec database.

328 A confidence threshold of 0.1 was set for read assignments, and reports were generated for
329 downstream biom file generation. For DNA sequencing data, read counts assigned to
330 taxonomies in each sample were then re-estimated using the average read length of that
331 sample, using Bracken[25] (v2.0). Kraken-biom (v1.0.1) was then used to generate biom file
332 in json format, using initial Kraken reports for RNA samples, and Bracken reports for DNA
333 samples. Biom (v2.1.6) was then used to assign tabulated metadata to this biom file.

334 ***Visualisation and comparison of taxa of interest in RNA and DNA***

335 A taxonomy table was generated from the biom file in R (v4.2.2) using Phyloseq[26] (v1.42.0)
336 and MicrobiotaProcess[27] (v1.10.3). Read-assigned taxonomy counts were parsed from this
337 table for any samples with both metagenomic (DNA) and metatranscriptomic data (n=985).
338 Counts were extracted for the following taxa: *Adenoviridae*, *Campylobacter*, *Clostridioides*
339 *difficile* (*C. difficile*), *Cryptosporidium*, *Escherichia coli* (*E. coli*), Norovirus, Rotavirus,
340 *Salmonella*, *Shigella*, Sapovirus, *Vibrio cholerae* (*V. cholerae*) and *Yersinia enterocolitica* (*Y.*
341 *enterocolitica*). These taxa were chosen to reflect the pathogen panels used during this study.
342 RNA virus (Astrovirus, Norovirus, Rotavirus, and Sapovirus) read counts could not be
343 extracted for this part of the analysis, as visualisations relied on the presence of DNA reads.
344 DNA and RNA counts were log-transformed and plotted against one another as a line graph
345 using standard functions in ggplot2[28] (v3.4.0). Visualisations were then used to assess the
346 sensitivity of metagenomics and metatranscriptomics for the selected taxa, where we define
347 sensitivity as the skew of data points towards either metagenomics (x-axis) or
348 metatranscriptomics (y-axis). A 0,0 intercept line was included in each line graph to assist in
349 illustrating sensitivity differences.

350 ***Correlation of genomic reads assigned to taxa of interest with number of observed taxa*** 351 ***and results from laboratory diagnostics***

352 Associations between read counts and laboratory results for organisms of interest were
353 assessed using a multivariable linear regression model in MaAsLin 2[29] (v1.6.0) under default

354 settings. The introduction of another variable into the model (laboratory results) provided a
355 point of reference. This allowed us to determine the relationship between any sample with
356 sequencing data and laboratory results. As such, for this analysis, all sequenced patient
357 samples (n = 1,067) were used, even if they did not contain both metagenomic and
358 metatranscriptomic data. Our approach allowed RNA virus read counts from
359 metatranscriptomic data to be included in this analysis. To visualise the strength of correlations
360 between laboratory results and pathogen-assigned sequencing reads, correlation coefficients
361 and adjusted p-values from the model were tabulated and used to generate a heatmap with
362 corrplot (v0.9.2). Adjusted p-values were generated using the Benjamini-Hochberg Procedure.

363 ***Comparative analysis of Adenovirus-associated k-mers in DNA and RNA***

364 The extract_kraken_reads.py utility from Kraken-tools (v1.2) was used alongside Kraken2
365 reports to extract reads with k-mer profiles associated with the family *Adenoviridae* for samples
366 that tested positive using either traditional or Luminex methods. Samples where sequencing
367 was not successful for both DNA and RNA were excluded from this analysis.

368 These reads were then mapped to the Human Mastadenovirus F genome (Accession:
369 GCF_000846685.1) using HISAT2[30] (v2.2.1) for splice-aware mapping. Coverage statistics
370 were then generated using samtools coverage. Coverage statistics for each sample were
371 compiled into a single table and visualised as a bar chart using ggplot2 (v3.4.0).

372

373 **Identification of CSpV1 as a biomarker of *Cryptosporidium* infection**

374 To understand why *Cryptosporidium*-associated k-mers showed a positive correlation with
375 using gold-standard diagnostics in metatranscriptomic but not metagenomic sequencing data,
376 we employed competitive mapping using PanGIA[31] (v1.0.0-RC6.2). We mapped quality-
377 controlled reads from all INTEGRATE samples against a database containing representative
378 and reference genomes of bacteria, archaea, and viruses in NCBI RefSeq (release 89). This
379 helped to validate our k-mer-based results and offers a less computationally intensive
380 alternative to mapping-based approaches for future users of k-mer-based databases. By
381 aligning the reads to these genome sequences, we obtained a read count and depth of
382 coverage for each organism. We then extracted entries associated with the term
383 '*Cryptosporidium*' along with their corresponding scores and mapping information. PanGIA
384 also accounts for many reads mapped equally well to other organisms and the percentage of
385 identity of these hits and derived a confidence score from this, ranging from 0 to 1 for each
386 query sequence at each taxonomy level. This allowed us to determine the certainty that the
387 organism is truly present in the sequencing data.

388 **Visualisation of a *Salmonella* transcriptome directly from stool**

389 Metatranscriptomic reads from a sample of a patient with a later-confirmed (culture positive)
390 *Salmonella* spp. infection underwent quality control, alignment, and quantification using the
391 Bacpipe RNA-seq processing pipeline (v0.6.0). The GFF annotation[32] for the *Salmonella*
392 *enterica* subsp. *enterica* serovar Enteritidis PT4 strain P125109 (Accession:
393 GCA_015240635.1) was used in this analysis. Coverage tracks and annotation were
394 visualised using JBrowse (v1.16.8). This visualisation can be found here:
395 https://s.hintonlab.com/study_74.

396 **Declarations**

397 ***Ethics approval and consent to participate***

398 Members of the public with symptoms of acute gastroenteritis, including a case definition of
399 vomiting and diarrhoea, who sought health advice from general practices in the RCGP RSC
400 NMN were invited to submit a stool sample for microbiological examination. Their consent for
401 this procedure was sought because normal care would not necessarily entail stool sampling
402 for most patients unless their symptoms were severe or had persisted for a long time. The
403 North West - Greater Manchester East Research Ethics Committee (REC reference:
404 15/NW/0233) and NHS Health Research Authority (HRA) Confidential Advisory Group (CAG)
405 (CAG reference: 15/CAG/0131) granted a favourable ethics opinion for the INTEGRATE
406 project. Approval was also granted by NHS Research Management and Governance
407 Committees (including Royal Liverpool and Broadgreen University Hospital Trust, Lancashire
408 Teaching Hospitals NHS Foundation Trust, Central Manchester University Hospitals NHS
409 Foundation Trust, and the University of Liverpool Sponsor), the Lancaster University Faculty
410 of Health and Medicine Ethics Committee, and the University of Liverpool Ethics Sub-
411 Committees. An Information Governance Toolkit (IGT) from the Department of Health hosted
412 by the Health and Social Care Information Centre (HSCIC) was also completed for the project,
413 and all project research staff obtained Honorary NHS contracts, research passports, and
414 letters of access, as necessary.

415 ***Consent for publication***

416 The publication was approved by the National Institute for Health and Care Research on 29th
417 March 2023.

418 ***Availability of data and materials***

419 Illumina sequence reads with human data removed have been deposited in the European
420 Nucleotide Archive (ENA) under ENA project accession number PRJEB62473.

421 ***Competing interests***

422 M.I.G. has received research grants from GSK and Merck, and has provided expert advice to
423 GSK. M.I.G. has been an employee of GSK since January 2023, although the work presented
424 here was completed prior to this date.

425 ***Funding***

426 This publication presents independent research supported by the Health Innovation Challenge
427 Fund (WT096200, HICF-T5-354), a parallel funding partnership between the Department of
428 Health and Wellcome Trust. The views expressed in this publication are those of the author(s)
429 and not necessarily those of the Department of Health or Wellcome Trust. This study is also
430 funded by the National Institute for Health Research (NIHR) Health Protection Research Unit
431 in Gastrointestinal Infections at University of Liverpool, in partnership with the UK Health
432 Security Agency (UKHSA), in collaboration with University of Warwick. E.C.-O., N.A.C. and
433 A.C.D. are based at The University of Liverpool. The views expressed are those of the
434 author(s) and not necessarily those of the NIHR, the Department of Health and Social Care or
435 the UK Health Security Agency. N.A.C. is a NIHR Senior Investigator (NIHR203756).

436 This work was supported by a Wellcome Trust Investigator award (grant number
437 222528/Z/21/Z) to J.C.D.H.

438 ***Authors' contributions***

439 Conceptualisation: E.C.-O. and A.C.D. Data curation: E.C.-O., B.P-S., M.W., C.A.N., K.M.M.,
440 S.H. and R.G. Formal analysis: E.C.-O. and Y.L. Funding acquisition: S.J.O'B. and N.A.C.
441 Investigation: E.C.-O. and A.C.D. Project administration: E.C.-O., M.I.G., C.H.-F., S.J.O'B.,
442 N.A.C. and A.C.D. Resources: Y.L., B.P-S., J.C.D.H., C.A.N., C.H.-F. and A.C.D. Supervision:
443 N.A.C. and A.C.D. Validation: E.C.-O. Visualisation: E.C.-O. and Y.L. Writing – original draft:
444 E.C.-O., B.P-S., J.C.D.H., N.A.C. and A.C.D. Writing – review and editing: all authors.

445 ***Acknowledgements***

446 The authors thank the members of the INTEGRATE Consortium. The INTEGRATE
447 Consortium investigators in the United Kingdom are Sarah J O'Brien (principal investigator),
448 Frederick J Bolton, Rob M Christley, Helen E Clough, Nigel A Cunliffe, Susan Dawson,
449 Elizabeth Deja, Ann E Durie, Sam Haldenby, Neil Hall, Christiane Hertz-Fowler, Debbie
450 Howarth, Lirije Hyseni, Miren Iturriza-Gómara, Kathryn Jackson, Lucy Jones, Trevor Jones, K
451 Marie McIntyre, Charlotte A Nelson, Lois Orton, Jane A Pulman, Alan D Radford, Danielle
452 Reaves, Helen K Ruddock, Darlene A Snape, Debbi Stanistreet, Tamara Thiele, Maya
453 Wardeh, David Williams, and Craig Winstanley (University of Liverpool), Kate Dodd (NIHR
454 Clinical Research Network: North West Coast), Peter J Diggle, Alison C Hale, Barry S
455 Rowlingson (Lancaster University), Jim Anson, Caroline E Corless, Viki Owen (Royal
456 Liverpool and Broadgreen University Hospitals NHS Trust), Malcolm Bennett (University of
457 Nottingham), Lorraine Bolton, John Cheesbrough, Katherine Gray, David Orr, Lorna Wilson
458 (Lancashire Teaching Hospitals NHS Foundation Trust), Andrew R Dodgson, Ashley McEwan
459 (Manchester University NHS Foundation Trust), Paul Cleary, Alex J Elliot, Ken H Lamden,
460 Lorraine Lighton, Catherine M McCann, Matthieu Pegorie, Nicola Schinaia, Anjila Shah, Gillian
461 E Smith, Roberto Vivancos, Bernard Wood (PHE), Rikesh Bhatt, Dyfrig A Hughes (Bangor
462 University), Rob Davies (APHA); Simon de Lusignan, Filipa Ferreira, Mariya Hriskova, Sam
463 O'Sullivan, Stacy Shinneman and Ivelina Yonova (University of Surrey/Royal College of
464 General Practitioners).

465 References

- 466 1. Foodborne Disease Estimates for the United Kingdom in 2018. Food Standards Agency;
467 2020. Available from:
468 [https://webarchive.nationalarchives.gov.uk/ukgwa/20200803160512/https://www.food.gov.uk](https://webarchive.nationalarchives.gov.uk/ukgwa/20200803160512/https://www.food.gov.uk/sites/default/files/media/document/foodborne-disease-estimates-for-the-united-kingdom-in-2018.pdf)
469 [/sites/default/files/media/document/foodborne-disease-estimates-for-the-united-kingdom-in-](https://webarchive.nationalarchives.gov.uk/ukgwa/20200803160512/https://www.food.gov.uk/sites/default/files/media/document/foodborne-disease-estimates-for-the-united-kingdom-in-2018.pdf)
470 [2018.pdf](https://webarchive.nationalarchives.gov.uk/ukgwa/20200803160512/https://www.food.gov.uk/sites/default/files/media/document/foodborne-disease-estimates-for-the-united-kingdom-in-2018.pdf). Accessed 09 Sep 2024.
471
- 472 2. Public Health England, NHS. UK Standards for Microbiology Investigations -
473 Gastroenteritis. UK Government; 2020. Available from:
474 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_da](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/930517/S_7i2_FINAL-UKSMI.pdf)
475 [ta/file/930517/S_7i2_FINAL-UKSMI.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/930517/S_7i2_FINAL-UKSMI.pdf). Accessed 09 Sep 2024.
- 476 3. Foddai ACG, Grant IR. Methods for detection of viable foodborne pathogens: current
477 state-of-art and future prospects. *Appl Microbiol Biotechnol*. 2020;104:4281–8.
- 478 4. Kanagarajah S, Waldram A, Dolan G, Jenkins C, Ashton PM, Carrion Martin AI, *et al*.
479 Whole genome sequencing reveals an outbreak of *Salmonella* Enteritidis associated with
480 reptile feeder mice in the United Kingdom, 2012-2015. *Food Microbiology*. 2018;71:32–8.
- 481 5. Chiu CY, Miller SA. Clinical metagenomics. *Nat Rev Genet*. 2019;20:341–55.
- 482 6. Buytaers FE, Saltykova A, Mattheus W, Verhaegen B, Roosens NHC, Vanneste K, *et al*.
483 Application of a strain-level shotgun metagenomics approach on food samples: resolution of
484 the source of a *Salmonella* food-borne outbreak. *Microb Genom*. 2021;7:000547.
- 485 7. Moeckel C, Mareboina M, Konnaris MA, Chan CSY, Mouratidis I, Montgomery A, *et al*. A
486 survey of k-mer methods and applications in bioinformatics. *Comput Struct Biotechnol J*.
487 2024;23:2289–303.
- 488 8. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome*
489 *Biol*. 2019;20:257.
- 490 9. Juul S, Izquierdo F, Hurst A, Dai X, Wright A, Kulesha E, *et al*. What's in my pot? Real-
491 time species identification on the MinION™. *bioRxiv*. 2015;030742.
- 492 10. Breitwieser FP, Lu J, Salzberg SL. A review of methods and databases for metagenomic
493 classification and assembly. *Brief Bioinform*. 2019;20:1125–36.
- 494 11. Van Etten J, Stephens TG, Bhattacharya D. A k-mer-Based Approach for Phylogenetic
495 Classification of Taxa in Environmental Genomic Data. *Systematic Biology*. 2023;72:1101–
496 18.
- 497 12. McIntyre KM, Bolton FJ, Christley RM, Cleary P, Deja E, Durie AE, *et al*. A Fully
498 Integrated Real-Time Detection, Diagnosis, and Control of Community Diarrheal Disease
499 Clusters and Outbreaks (the INTEGRATE Project): Protocol for an Enhanced Surveillance
500 System. *JMIR Res Protoc*. 2019;8:e13941.
- 501 13. Kröger C, Colgan A, Srikumar S, Händler K, Sivasankaran SK, Hammarlöf DL, *et al*. An
502 Infection-Relevant Transcriptomic Compendium for *Salmonella enterica* Serovar
503 Typhimurium. *Cell Host & Microbe*. 2013;14:683–95.
- 504 14. Rajagopala SV, Bakhoun NG, Pakala SB, Shilts MH, Rosas-Salazar C, Mai A, *et al*.
505 Metatranscriptomics to characterize respiratory virome, microbiome, and host response
506 directly from clinical samples. *Cell Rep Methods*. 2021;1:100091.

- 507 15. Feng Y, Ramnarine VR, Bell R, Volik S, Davicioni E, Hayes VM, *et al.* Metagenomic and
508 metatranscriptomic analysis of human prostate microbiota from patients with prostate
509 cancer. *BMC Genomics*. 2019;20:146.
- 510 16. Galen SC, Borner J, Williamson JL, Witt CC, Perkins SL. Metatranscriptomics yields new
511 genomic resources and sensitive detection of infections for diverse blood parasites. *Mol Ecol*
512 *Resour*. 2020;20:14–28.
- 513 17. Adjou KT, Chevillot A, Lucas P, Blanchard Y, Louifi H, Arab R, *et al.* First identification of
514 *Cryptosporidium parvum* virus 1 (CSpV1) in various subtypes of *Cryptosporidium parvum*
515 from diarrheic calves, lambs and goat kids from France. *Vet Res*. 2023;54:66.
- 516 18. Chae J-B, Shin S-U, Kim S, Jo Y-M, Roh H, Chae H, *et al.* The First Identification of
517 *Cryptosporidium parvum* Virus-1 (CSpV1) in Hanwoo (*Bos taurus coreanae*) Calves in
518 Korea. *Vet Sci*. 2023;10.
- 519 19. Geldenhuys M, Mortlock M, Weyer J, Bezuidt O, Seamark ECJ, Kearney T, *et al.* A
520 metagenomic viral discovery approach identifies potential zoonotic and novel mammalian
521 viruses in *Neoromicia* bats within South Africa. *PLoS One*. 2018;13:e0194527.
- 522 20. Vibin J, Chamings A, Klaassen M, Bhatta TR, Alexandersen S. Metagenomic
523 characterisation of avian parvoviruses and picornaviruses from Australian wild ducks. *Sci*
524 *Rep*. 2020;10:12800.
- 525 21. Qin H, Lo NW-S, Loo JF-C, Lin X, Yim AK-Y, Tsui SK-W, *et al.* Comparative
526 transcriptomics of multidrug-resistant *Acinetobacter baumannii* in response to antibiotic
527 treatments. *Sci Rep*. 2018;8:3515.
- 528 22. Parks DH, Chuvochina M, Reeves PR, Beatson SA, Hugenholtz P. Reclassification of
529 *Shigella* species as later heterotypic synonyms of *Escherichia coli* in the Genome Taxonomy
530 Database. *bioRxiv*. 2021;2021.09.22.461432.
- 531 23. Devanga Ragupathi NK, Muthuirulandi Sethuvel DP, Inbanathan FY, Veeraraghavan B.
532 Accurate differentiation of *Escherichia coli* and *Shigella* serogroups: challenges and
533 strategies. *New Microbes New Infect*. 2018;21:58–62.
- 534 24. Uritskiy GV, DiRuggiero J, Taylor J. MetaWRAP—a flexible pipeline for genome-resolved
535 metagenomic data analysis. *Microbiome*. 2018;6:158.
- 536 25. Lu J, Breitwieser F, Thielen P, Salzberg S. Bracken: Estimating species abundance in
537 metagenomics data. *PeerJ Comput Sci*. 2017;3:e104.
- 538 26. McMurdie PJ, Holmes S. phyloseq: An R Package for Reproducible Interactive Analysis
539 and Graphics of Microbiome Census Data. *PLoS One*. 2013;8:e61217.
- 540 27. Xu S, Zhan L, Tang W, Wang Q, Dai Z, Zhou L, *et al.* MicrobiotaProcess: A
541 comprehensive R package for deep mining microbiome. *The Innovation*. 2023;4:100388.
- 542 28. Valero-Mora PM. ggplot2: Elegant Graphics for Data Analysis. *J Stat Soft*. 2010;35:1–3.
- 543 29. Mallick H, Rahnavard A, Mclver LJ, Ma S, Zhang Y, Nguyen LH, *et al.* Multivariable
544 association discovery in population-scale meta-omics studies. *PLoS Comput Biol*.
545 2021;17:e1009442.

546 30. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and
547 genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol.* 2019;37:907–15.

548 31. Li P-E, Russell JA, Yarmosh D, Shteyman AG, Parker K, Wood H, *et al.* PanGIA: A
549 Metagenomics Analytical Framework for Routine Biosurveillance and Clinical Pathogen
550 Detection. *bioRxiv.* 2020;2020.04.20.051813.

551 32. Perez-Sepulveda BM, Predeus AV, Fong WY, Parry CM, Cheesbrough J, Wigley P, *et*
552 *al.* Complete Genome Sequences of African *Salmonella enterica* Serovar Enteritidis Clinical
553 Isolates Associated with Bloodstream Infection. *Microbiol Resour Announc.* 2021;10.

554 **Figure legends**

555 **Figure 1: Visual overview and comparison of DNA (metagenomic) and RNA**
556 **(metatranscriptomic) sequencing reads** assigned to GI pathogens of relevance to the UK
557 setting. For all graphs, the dashed (black) intercept line is provided to highlight the skew of
558 sensitivity towards either DNA or RNA.

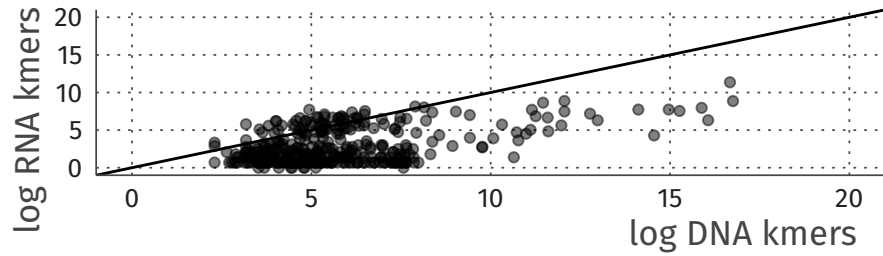
559 **Figure 2: Statistically significant correlations were observed between sequencing data**
560 **and laboratory tests for 10 out of 14 major GI community pathogens in the United**
561 **Kingdom.** Results where at least one statistically significant correlation was observed are
562 shown. All correlations, whether significant or not, are displayed in Supplementary Figure 1.
563 No statistically significant correlation was found between the sequencing and diagnostic test
564 for Astrovirus, *E. histolytica*, *Giardia* or *V. cholerae*.

565 The darker the colour of a quadrant in a heatmap, the stronger the correlation (coefficient)
566 between the detection of a pathogen in sequencing data (metagenomic or metatranscriptomic)
567 and a laboratory result (Luminex or Traditional). Asterisks in quadrants indicate the statistical
568 significance of correlations as follows: *: $p < 0.25$; **: $p < 0.05$; ***: $p < 0.01$; ****: $p < 0.001$.

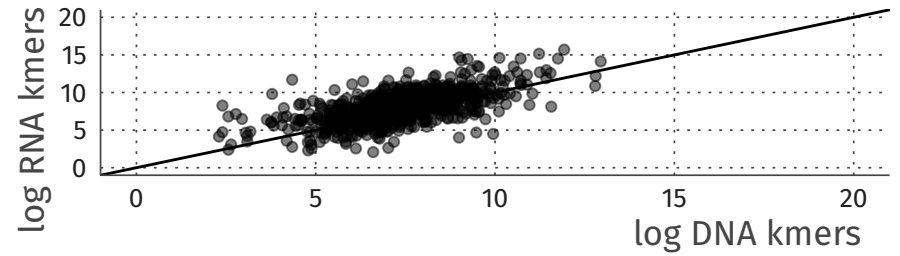
569
570 **Figure 3: Adenovirus can be detected through its genomic material and the expression**
571 **of transcript, directly from stool.** Graphs display the breadth of coverage (%) for both DNA
572 (purple) and RNA (orange) across nine samples, chosen on the basis of positive results
573 through gold-standard laboratory methods. Coverage values were generated via mapping to
574 a representative Human mastadenovirus F genome (GCF_000846685.1, International
575 Committee on Taxonomy of Viruses species exemplar).

576 **Supplementary Figure 1:** Complete overview of correlations observed between sequencing
577 data and laboratory tests for major GI community pathogens in the United Kingdom
578 (alternative to Figure 2).

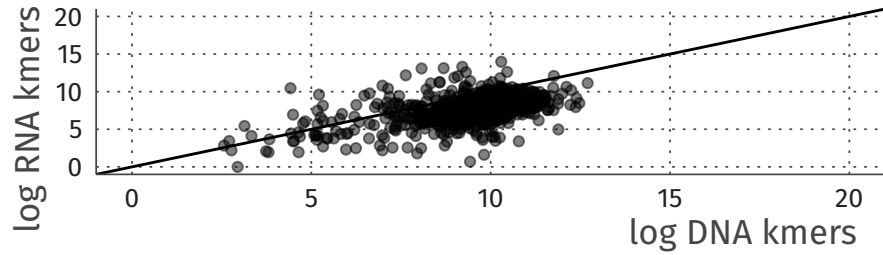
Adenovirus



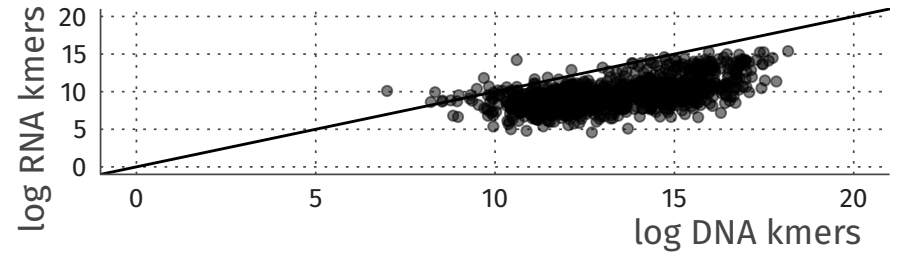
Campylobacter



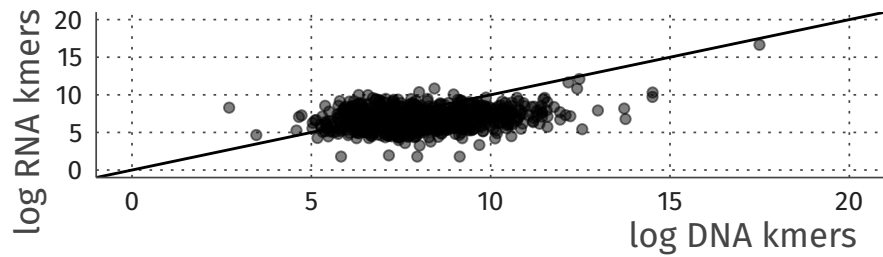
Clostridioides difficile



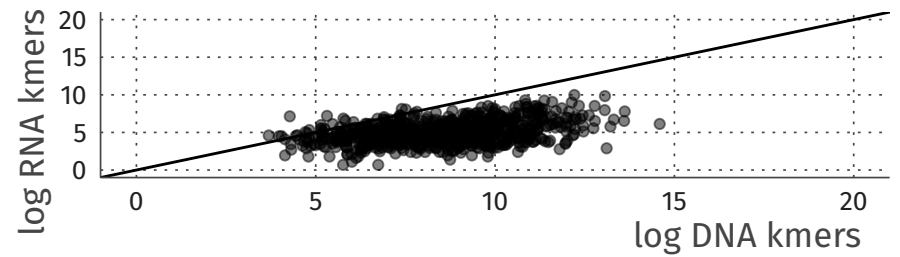
Escherichia coli



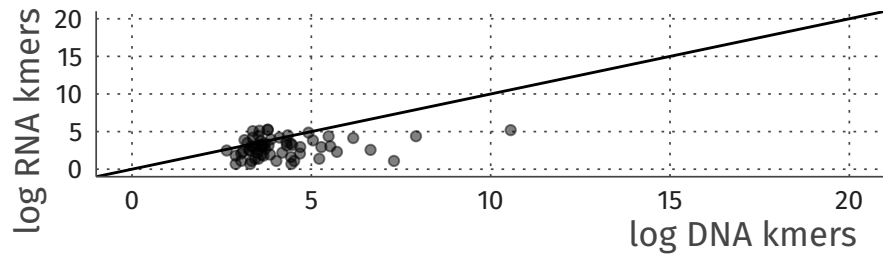
Salmonella



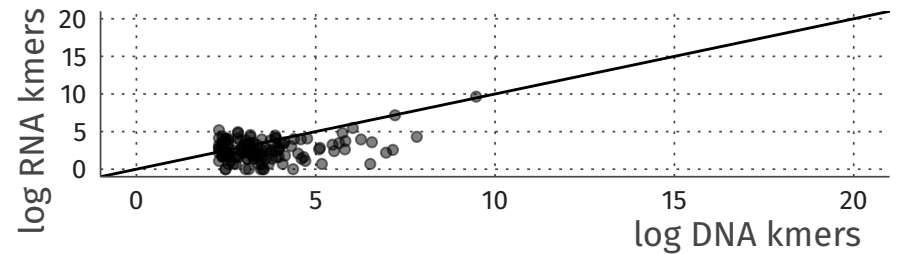
Shigella



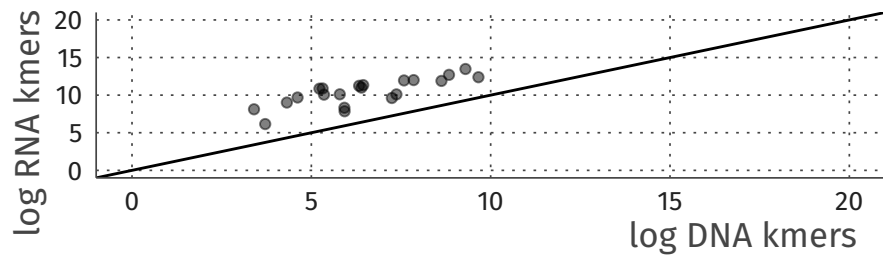
Vibrio cholerae



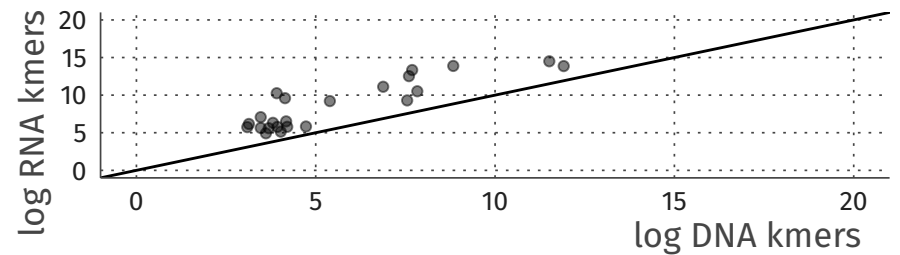
Yersinia enterocolitica

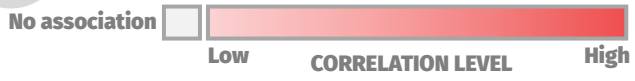
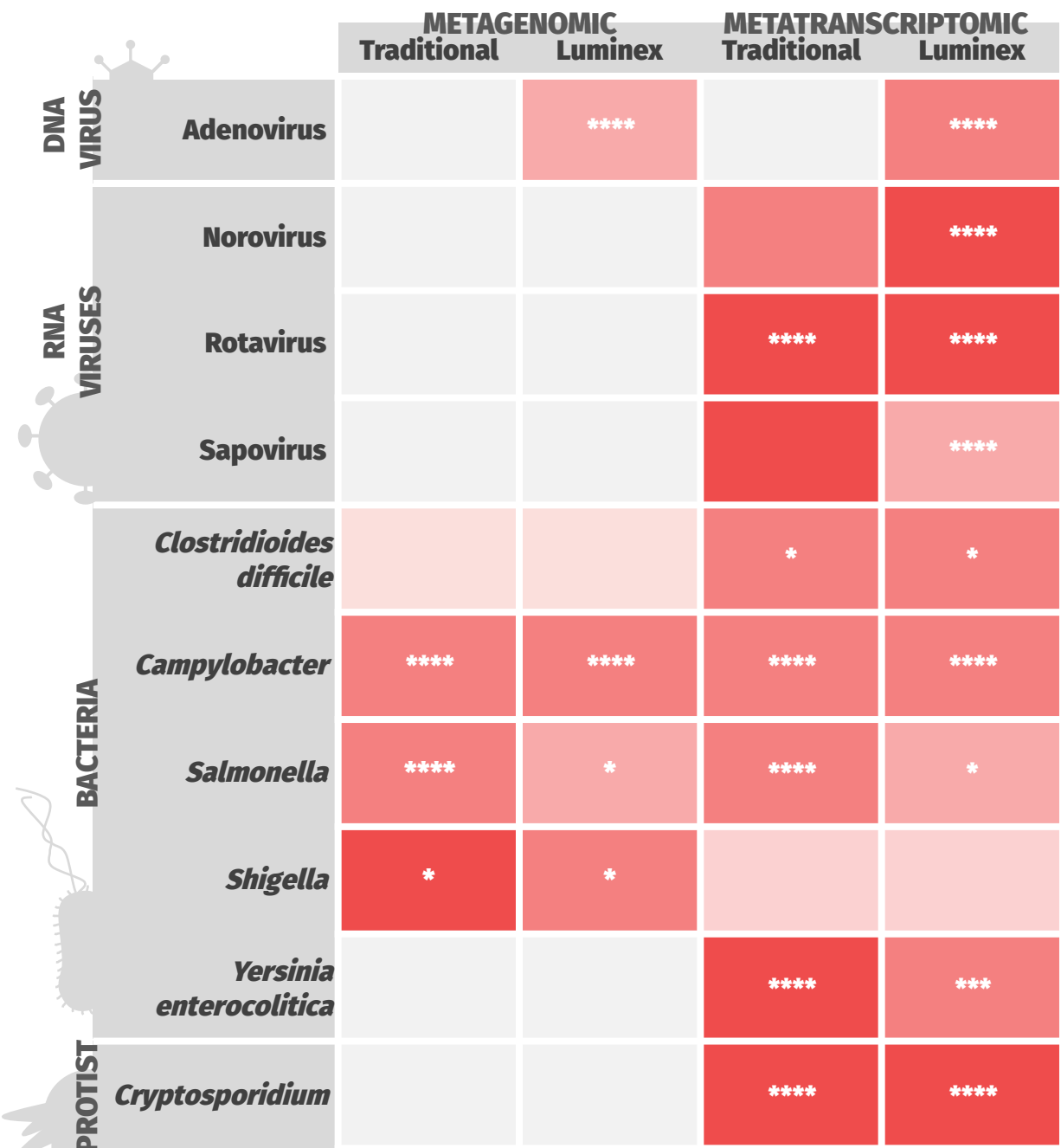


Cryptosporidium



Giardia





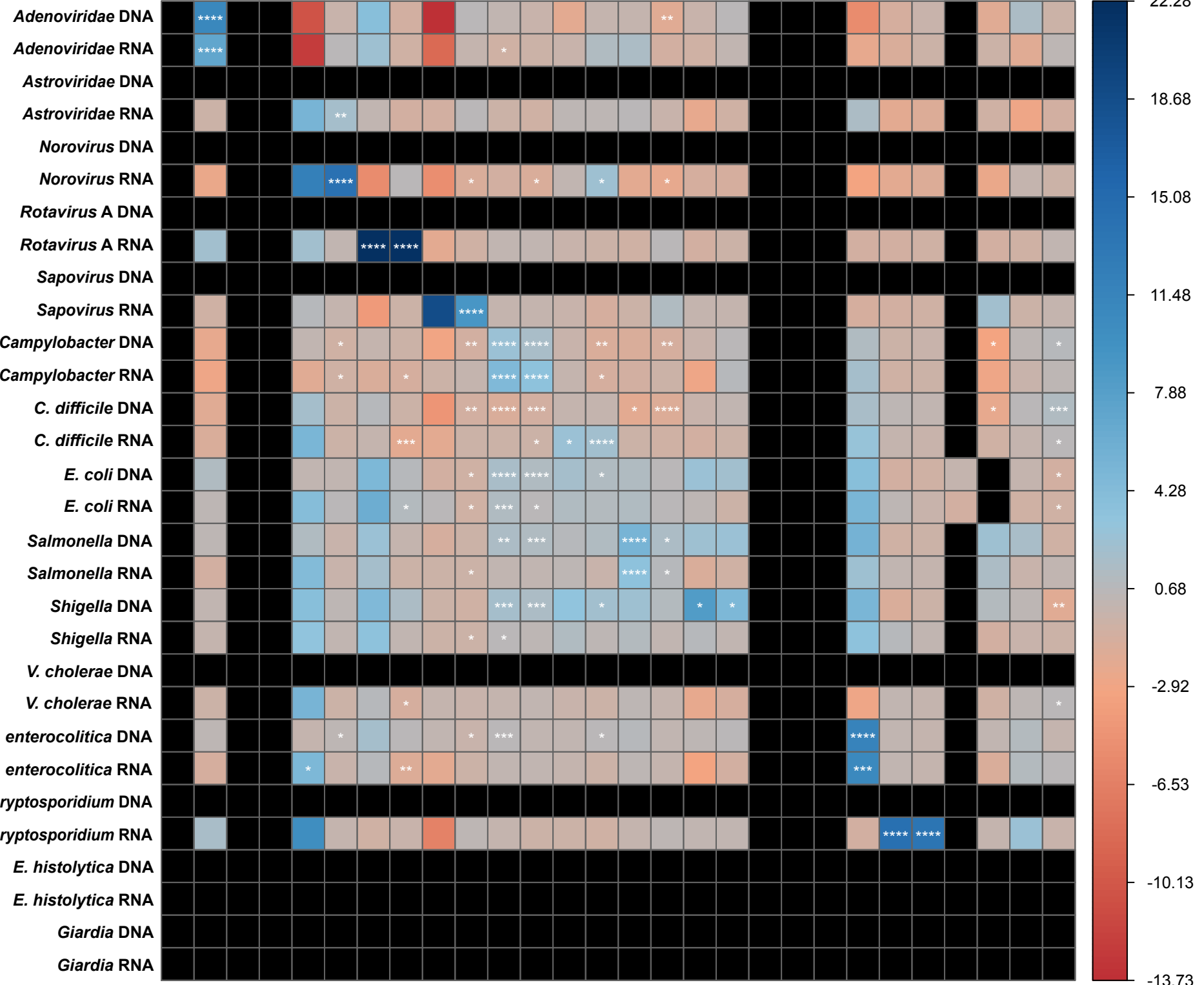
* p < 0.25 **p < 0.05 ***p < 0.01 ****p < 0.001

Human mastadenovirus F

GCF_000846685.1



Adenovirus culture-based
 Adenovirus PCR-based
 Astrovirus culture-based
 Astrovirus PCR-based
 Norovirus GI culture-based
 Norovirus GI PCR-based
 Rotavirus A culture-based
 Rotavirus A PCR-based
 Sapovirus culture-based
 Sapovirus PCR-based
 Campylobacter culture-based
 Campylobacter PCR-based
 C. difficile culture-based
 C. difficile PCR-based
 Salmonella culture-based
 Salmonella PCR-based
 Shigella culture-based
 Shigella PCR-based
 V. cholerae culture-based
 V. cholerae PCR-based
 Y. enterocolitica culture-based
 Y. enterocolitica PCR-based
 Cryptosporidium culture-based
 Cryptosporidium PCR-based
 E. histolytica culture-based
 E. histolytica PCR-based
 Giardia culture-based
 Giardia PCR-based



Key for significance values:
 * = <0.25, ** = <0.05, *** = <0.01, **** = <0.001