

1 **Validation of an unbiased metagenomic detection assay for RNA viruses in**
2 **viral transport media and plasma**

3 Anthony D. Kappell (0000-0003-3511-9207)*, Kathleen Q. Schulte (0000-0002-0558-0660), Elizabeth
4 A. Scheuermann (0000-0002-7498-7130), Matthew B. Scholz (0000-0003-3686-1227), Nicolette C.
5 Keplinger, Amanda N. Scholes (0000-0003-4964-6139), Taylor A. Wolt, Viviana M. June (0000-0001-
6 5519-4971), Cole J. Schulte, Leah W. Allen, Krista L. Ternus (0000-0003-1138-5308) and F. Curtis
7 Hewitt (0000-0001-9546-5625)

8 Signature Science, LLC, Austin, Texas, United States of America

9 *Corresponding author

10

11 Corresponding author contact information:

12 Anthony D. Kappell

13 akappell@signaturescience.com

14 **Abstract (250 words)**

15 Unbiased long read sequencing holds enormous potential for the detection of pathogen sequences in
16 clinical samples. However, the untargeted nature of these methods precludes conventional PCR
17 approaches, and the metagenomic content of each sample increases the challenge of bioinformatic
18 analysis. Here, we evaluate a previously described novel workflow for unbiased RNA virus sequence
19 identification in a series of contrived and real-world samples. The novel multiplex library preparation
20 workflow was developed for the Oxford Nanopore Technologies (ONT) MinION™ sequencer using
21 reverse transcription, whole genome amplification, and ONT's Ligation Sequencing Kit with Native
22 Barcode Expansion. The workflow includes spiked MS2 Phage as an internal positive control and
23 generates an 8-plex library with 6 samples, a negative control and a *gfp* transcript positive control.
24 Targeted and untargeted data analysis was performed using the EPI2ME Labs framework and open access
25 tools that are readily accessible to most clinical laboratories. Contrived samples composed of common
26 respiratory pathogens (Influenza A, Respiratory Syncytial Virus and Human Coronavirus 229E) in viral
27 transport media (VTM) and bloodborne pathogens (Zika Virus, Hepatitis A Virus, Yellow Fever Virus
28 and Chikungunya Virus) in human plasma were used to establish the limits of detection for this assay. We
29 also evaluated the diagnostic accuracy of the assay using remnant clinical samples and found that it
30 showed 100% specificity and 62.9% clinical sensitivity. More studies are needed to further evaluate
31 pathogen detection and better position thresholds for detection and non-detection in various clinical
32 sample metagenomic mixtures.

33 **Keywords (3-10 words)**

34 Nanopore Sequencing; MinION; RNA Virus; Agnostic Sequencing; Metagenomics; Respiratory
35 Pathogens; Blood Pathogens; coronavirus; SARS-CoV-2; COVID-19

36

37 **Introduction**

38 RNA viruses pose a significant threat to global public health. The prime example of this is the
39 COVID-19 pandemic, but it is far from the only virus that has caused outbreaks in recent years. There are
40 seasonal outbreaks every year of respiratory RNA viruses, such as paramyxoviruses (e.g., respiratory
41 syncytial virus and human metapneumovirus) and alphainfluenzaviruses (e.g., Influenza A). There have
42 also been recurring regional outbreaks of mosquito borne diseases caused by RNA viruses such as
43 alphaviruses (e.g., Chikungunya virus) and flaviviruses (e.g., Zika virus, Dengue virus and Yellow Fever
44 virus). These viruses pose a serious threat to human health – in the United States alone, hundreds of
45 thousands of people are hospitalized each year from complications of respiratory viral infections
46 (Thompson et al. 2004). Respiratory viral illnesses are also associated with a steep economic burden; they
47 are collectively estimated to cost over \$127 billion each year (Fendrick et al. 2003; Young-Xu et al.
48 2017).

49 Unbiased, metagenomic sequence-based approaches could lead to early detection of both
50 previously characterized and novel RNA viruses and serve as a universal assay to detect infectious
51 disease agents (Bibby 2013; Miller et al. 2013; Schlaberg et al. 2017). Established clinical tests largely
52 rely on PCR, culturing or antibody specificity to detect pathogens. These approaches are limited in that
53 they require prior assumptions about the pathogens that might be present. Furthermore, many diagnostic
54 tests rely on culturing the pathogen, which can create delays of several days to diagnosis. Hypothesis-free
55 approaches have the distinct advantage of being able to survey the presence of multiple infectious agents
56 at once in a single assay, which allows pathogens to be identified more quickly and without collecting
57 further samples for a series of tests.

58 Nanopore sequencing is a cost-effective third-generation sequencing method that allows long
59 reads to be generated on a hand-held device. It has a shorter turn-around time than other next-generation
60 sequencing technologies (Petersen et al. 2019; Miller and Chiu 2022) and has been successfully employed
61 to identify viral (Russell et al. 2018; Arévalo et al. 2022), fungal (Ohta et al. 2023) and bacterial agents

62 (Hewitt et al. 2017; Charalampous et al. 2019; Bouchiat et al. 2022). Nanopore sequencing, unlike
63 traditional RT-PCR methods to diagnose viral illness, allows for whole genome sequencing of viruses,
64 allowing the identification of viral variants to improve disease surveillance efforts. This was used in the
65 recent COVID-19 pandemic to monitor emerging SARS-CoV-2 variants (Yakovleva et al. 2022; Centers
66 for Disease Control and Prevention 2022). To effectively bring this sequencing to a clinical diagnostic
67 setting, there is a need to effectively validate the use of sequencing assays on clinical samples and
68 develop bioinformatics pipelines that can generate reports that provide clinically actionable insights
69 (Miller et al. 2019).

70 Here, we validate the use of a metagenomic sequencing workflow that rapidly and accurately
71 detects RNA viruses in multiple clinical sample types. This workflow uses unbiased amplification to
72 increase assay sensitivity and multiplexed sequencing to increase throughput and decrease sample
73 analysis cost. Several performance metrics were collected to validate the use of this assay for clinical
74 metagenomics. The limit of detection (LoD) and precision of this assay were determined for seven RNA
75 viruses using contrived samples. We also evaluated the degree to which common contaminants in clinical
76 samples, such as EDTA, interfered with pathogen detection in this workflow. Furthermore, this workflow
77 was tested on commercially purchased remnant clinical samples in addition to contrived samples created
78 in our laboratory, allowing evaluation of assay performance directly against existing clinical assays to
79 determine accuracy.

80

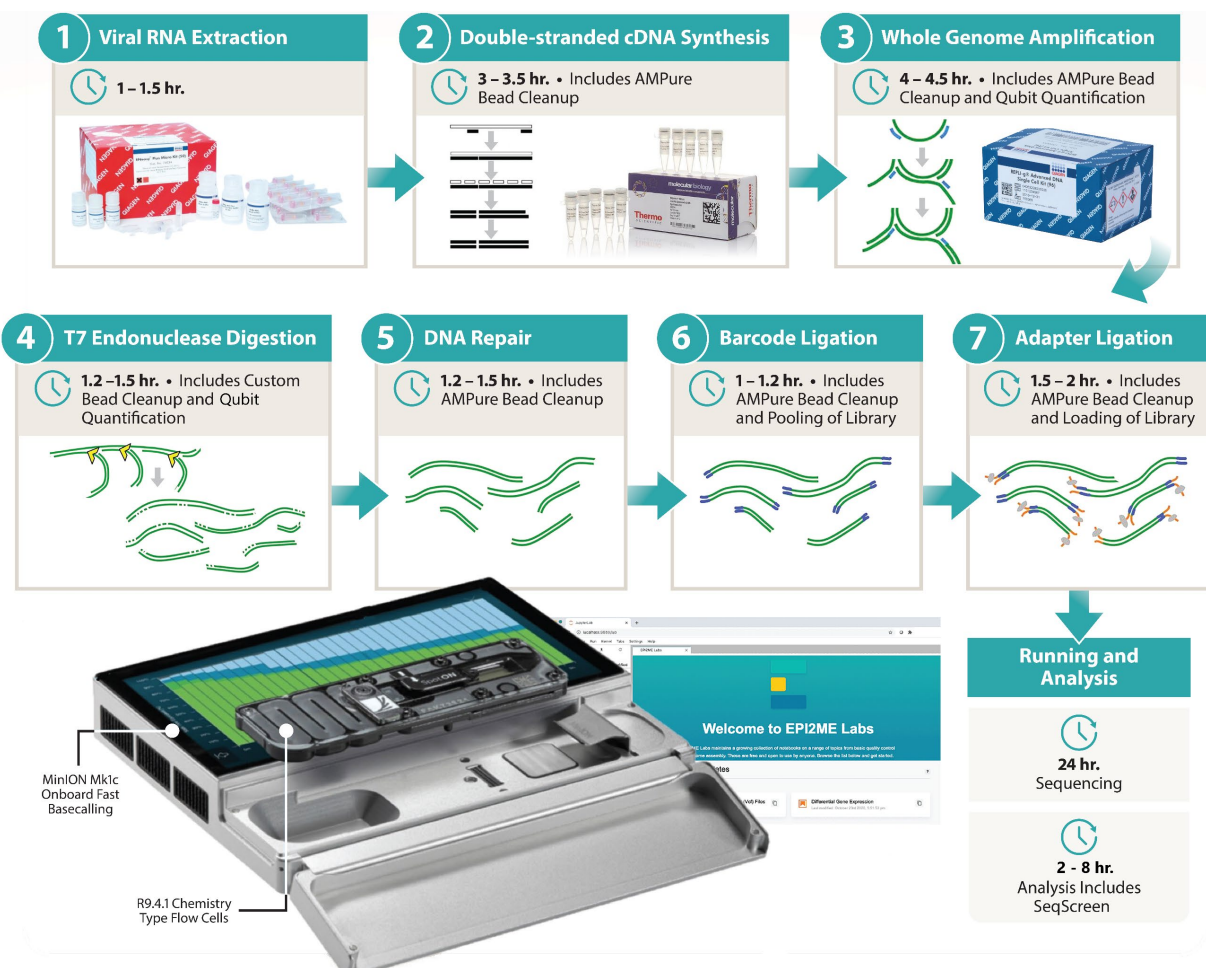
81 **Results**

82 **Sample Processing and Bioinformatics Analysis**

83 We developed an untargeted third-generation sequencing assay for RNA virus pathogen identification
84 from nasopharyngeal swabs in viral transport media (VTM) and plasma (Figure 1). This assay includes
85 library preparation, sequencing, and bioinformatics analysis. We then evaluated the performance of the

86 assay using contrived and clinically relevant remnant samples. For each sequencing run, internal
87 sequencing controls consisting of MS2 Phage were spiked into each of the six clinically relevant samples
88 and a NTC (“no template control”) of sterile phosphate buffered saline processed in parallel through RNA
89 extraction. Each sequencing run also had a PC (positive control) consisting of an RNA transcript of the
90 *gfp* gene which was processed through the RNA sequencing library preparation in parallel with the RNA
91 extracted clinically relevant samples and NTC. The sequencing library preparation (Figure 1) for RNA for
92 the untargeted sequencing assay included (1) double-stranded complimentary DNA synthesis (dscDNA
93 synthesis), (2) whole-genome amplification (WGA), (3) ONT native barcoding kit, (4) library pooling in
94 equal molar concentration, (5) ONT ligation sequencing kit, and (6) sequencing on the MinION™ (Mk1B
95 or Mk1C). Raw reads, ‘squiggles’ in the fast5 format, were basecalled and demultiplexed either ‘live’
96 while sequencing or post-run basecalled using ONT’s official ‘guppy’ basecaller. The resulting
97 demultiplexed and quality filtered reads were analyzed through a bioinformatic pipeline using Epi2Me-
98 Labs as a framework. The SigSciDx bioinformatic pipeline consisted of mapping each of the barcoded
99 read bins to the human genomic sequence using ‘minimap2’ (Li 2021, 2018) to remove human aligned
100 reads. The remaining reads were then mapped to the positive control gene sequence (*gfp*) and internal
101 sequencing control (MS2 Phage) using ‘minimap2’ and alignment statistics were generated with
102 ‘samtools’ (Danecek et al. 2021; Li et al. 2009) for assessment of quality of the sequencing run and for
103 the individual samples. The barcoded read bins were also analyzed for taxonomical determination through
104 the ‘reference inference’ module from SeqScreen. Results from the SigSciDx bioinformatic pipeline were
105 published in a html formatted report for clinician review and interpretation. All remnant samples, positive
106 or negative, were blinded to the analysts performing the extraction and library preparation, and to the
107 analysts performing initial bioinformatic evaluation and interpretation.

108



109

110

111

Figure 1. Untargeted sequencing workflow

112 Establishing thresholds for reporting detected pathogens

113 To minimize false-positive results from cross-contamination caused by low-level barcode crosstalk (Xu et
114 al. 2018), we examined the use of ONT’s current official basecaller ‘guppy’ on 1) a setting for single-end
115 barcode binning of reads, requiring detection of only one barcode per sequence read, or 2) both-end
116 barcode binning, requiring matching barcodes are present at both ends of each read. We examined
117 crosstalk from samples in the NTC (Figure 2, Table S1), indicating potential false positives and the loss
118 of sensitivity (Table 1).

119

120

Table S1. Indication of Barcode Crosstalk by Detection of RNA Virus in TNC

Media		Viral Transport Media						Human Plasma							
Organism	Influenza A	Human Respiratory Syncytial Virus		Human Coronavirus 229E		Zika Virus		Hepatitis A Virus		Yellow Fever Virus		Chikungunya Virus			
		S	B	S	B	S	B	S	B	S	B	S	B		
Highest Load Levels on Run	A	2/2	0/2	2/2	0/2	2/2	2/2	1/1	1/1	1/1	0/1	1/1	1/1	1/1	1/1
	B	2/2	0/2	2/2	0/2	2/2	1/2	NA	NA	NA	NA	NA	NA	NA	NA
	C	0/1	0/1	1/1	0/1	1/1	0/1	2/2	2/2	2/2	0/2	2/2	0/2	2/2	0/2
	D	0/1	0/1	0/1	0/1	0/1	0/1	2/2	0/2	0/2	0/2	2/2	0/2	2/2	0/2
	E	NA	NA	NA	NA	NA	NA	1/3	0/3	0/3	0/3	0/3	0/3	2/3	0/3

* S: Single-barcode binning, B: Both-barcode binning

121

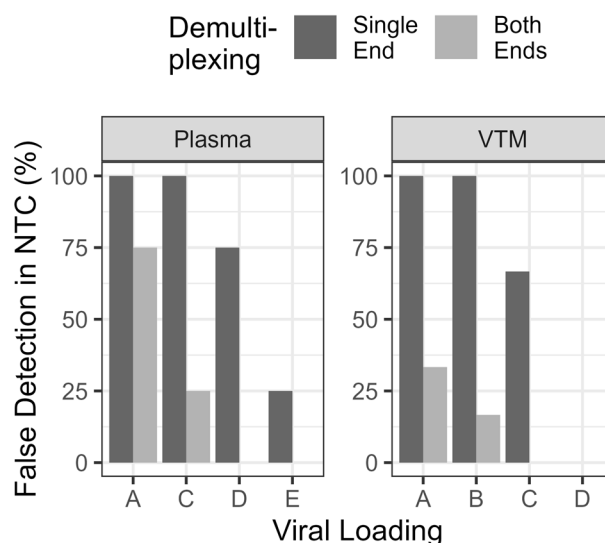


Figure 2. Indication of Barcode Crosstalk by Detection of RNA Viruses in NTC of a Run. Percent of RNA viruses detected in contrived samples at different loadings (Table 4) across run using single-end or both-ends demultiplexing methods.

122 Utilizing the both-end barcode binning setting on the basecaller significantly reduced inherent
 123 potential false positive calls within the NTC with only a minor reduction in sensitivity of between 0.3 and
 124 1-Log. Total read counts of binned barcodes reduced crosstalk by 38.9 ± 5.5 % between single-end and
 125 both-end binning settings.

126 **Table 1. Comparison of Limit of Detection (LoD) Differences**
 127 **between Single and Both-end Barcoding Binning**

Organism	Both-end Barcode Binning	Single-end Barcode Binning	LoD Log Differences
Human Respiratory Virus	2.33	1.95	0.38
Influenza A	4.43	4.05	0.38
Human coronavirus 229E	-0.01	NA ^	NA
Zika Virus	2.74	2.45	0.29
Chikungunya Virus	2.06	1.55	0.51

Yellow Fever Virus	1.88	0.88	1.01
Hepatitis A Virus	3.12	NA ^A	NA

^A Overestimated due to fewer samples at lower load levels

128

129 The sequence analysis results from contrived and remnant clinical samples informed the
130 SigSciDx bioinformatic workflow thresholds as described in the Materials and Methods. This included
131 limiting the number of taxonomical ids initially entering the ‘reference inference module’ and decreasing
132 the stringency of which reference genomic sequences proceed from the first round of mapping to the
133 second round of mapping within the ‘reference inference module’. The criteria for analyst interpretation
134 of detection results were based on replicate sequencing results and inter-run information from contrived
135 and remnant clinical samples, including those from PC, NTC, and internal control. Total number of reads
136 per a barcode with a quality score greater than Q8 was expected to be greater than 40,000. The MinION™
137 run was expected to achieve greater than 1.5 million long-reads representing greater than 4 Gb estimated
138 bases with an estimated N50 greater than 2.5 kb, usually 3.5 kb. The minimum number of *gfp* gene reads
139 in the PC representative barcode must have been greater than 5,000 for a successful run to be considered
140 for further interpretation. MS2 Phage was accessed on two separate criteria based on mapped read count
141 or percentage, depending on which was greater. MS2 Phage in each sample was expected to have mapped
142 greater than 1000 reads, especially in samples with many reads remaining after host removal, or at 50% of
143 the reads remaining after host removal, especially if fewer reads remain after host removal. The NTC and
144 *gfp* gene barcode samples were used to interpret if high barcode crosstalk or cross-contamination exist. It
145 was expected that only MS2 Phage, or other phages, may be detected within the NTC. If an NTC detected
146 a virus that was also detected in at least one other sample of the run, this dictated a deeper review.
147 Generally, the preferred remedy was determining the sample with a high read count which was positive
148 for the detected virus and re-sequence using the cleaned WGA DNA material of the other samples to
149 verify their positive calls. Otherwise, resequencing all the samples again was also viable. There was little
150 ability to access if a statistical cut-off could be employed such as requiring a 2 or 3 -fold higher number of
151 reads in a sample compared to the NTC for that virus to indicate a detection, as this appeared in only 1 of

152 the runs performed on the remnant samples. Further study of remnant or clinical samples could inform a
153 statistical required read count in the presence of barcode crosstalk or cross-contamination.

154

155 **Limit of Detection**

156 A 95% limit of detection (LoD) was determined for each of the three and four representative RNA viruses
157 in VTM and plasma contrived samples, respectively. We evaluated each RNA virus over a minimum 6-
158 Log dilution range, testing 2 to 14 replicates at each concentration. The 95% LoD, defined as the lowest
159 concentration at which 95% of positive samples are detected, was determined for each of the seven RNA
160 viruses using probit analysis (Table 2).

161

Table 2. Performance characteristics for the untargeted sequencing assay

Performance metric	Method	Result
Limit of detection (LoD)	Qualitative detection of RNA virus dilution replicates by probit analysis	Pathogen LoD (Estimated PFU ¹)
		Influenza A Virus 26,915 TCID ₅₀ /mL (18,840 PFU/mL)
		Human Respiratory Syncytial Virus 214 TCID ₅₀ /mL (150 PFU/mL)
		Human Coronavirus 229E 1 TCID ₅₀ /mL (1 PFU/ mL)
		Zika Virus 550 TCID ₅₀ /mL (385 PFU/ mL)
		Hepatitis A Virus 1318 TCID ₅₀ /mL (923 PFU/ mL)
		Yellow Fever Virus 68 TCID ₅₀ /mL (47 PFU/ mL)
		Chikungunya Virus 115 TCID ₅₀ /mL (80 PFU/ mL)
Precision	Qualitative detection over 2 to 7 contrived sample runs of each organism (inter-assay) ²	100 % concordance
	Qualitative detection of duplicate contrived samples on the same run (intra-assay) ²	100 % concordance
	Qualitative detection over 21 remnant sample (duplicate or triplicate) runs (inter-assay)	100% concordance
Interference	Quantitative read count of viruses with spiked blood (5%, 2%)	Decrease of number of reads by 1.3-Logs to 1.5-Logs mapped to viruses
	Quantitative read count of viruses with spiked EDTA (100 mM, 10 mM, 5 mM)	Decrease of number of reads by 0.6-Logs and 0.1-Logs with addition of 100 mM and 10 mM EDTA, respectively
	Quantitative read count of viruses with spiked bacteria (<i>Micrococcus luteus</i> or <i>Staphylococcus epidermidis</i>)	No significant changes in number of reads mapped to RNA viruses
Accuracy	18 of single-detection remnant sample runs, results comparisons (3 for each organism).	
	<u>Sensitivity</u>	<u>Specificity</u>
	Human metapneumovirus 33.3% ³	100%
	Parainfluenza IV 33.3%	100%
	SARS-CoV-2 100%	100%
	RSV 100%	100%
	Influenza A 100%	100%
	Enterovirus 100%	100%

¹ Estimated PFU calculated by multiplying the TCID₅₀ /mL by 0.7² Nearest higher loading (Table 4) to that of the LOD (within 1-Log)³ Includes two of the three samples that did not pass internal control QC

164 **Precision**

165 We demonstrated inter-assay reproducibility of the untargeted sequencing assay by testing of the NTC
166 and contrived samples at the nearest loading above the calculated LoD for each organism across at least 4
167 sequencing runs and intra-assay reproducibility by testing of at least 4 independently generated sets of
168 duplicate contrived samples on the same run (and negative remnant samples in two runs). The *gfp* gene
169 mRNA transcript positive control passed QC for every completed sequencing run (36 total runs: 20
170 contrived sample runs and 16 remnant sample runs), indicating successful sequencing library generation
171 from cDNA synthesis through sequencing. Internal spiked MS2 phage controls passed QC for all
172 sequenced replicate contrived samples (216 total), indicating conditions across all processes of the
173 workflow were successful, including the RNA extraction process. For remnant samples, 2 samples
174 consistently failed internal spiked phage control QC after 3 independent runs were performed (6 total
175 remnant replicate samples). The remaining 90 of the 96 replicate remnant samples across the 16 runs
176 passed QC for the internal spiked MS2 phage control. The 6 failed internal controls in the replicates of the
177 2 remnant samples occurred in samples that were Human metapneumovirus-positive by original clinical
178 assay and had relative low human sequence contributions compared to the other VTM samples. One of
179 the samples, despite the failed QC, detected the presence of Human metapneumovirus while the other
180 consistently failed to detect. All seven organisms in the contrived samples were detected using pre-
181 established threshold criteria for the duplicate intra-assay samples on the same run and each replicate
182 inter-assay run (Table 2). All 9 negative remnant samples, determined to be absent of viral pathogens by
183 previous clinical testing, had no detection by untargeted sequencing of pathogenic RNA viruses within
184 the same run which were divided between 2 individual runs (intra-assay testing). From across all 36
185 sequencing runs (inter-assay testing), 35 NTCs showed no pathogenic RNA virus detections. The single
186 run that showed a positive detection within the NTC, presence of Influenza A, was a remnant run
187 containing a patient sample with a high-load of Influenza A and subsequent two repeats of the run using
188 that sample did not show detection in the NTCs. The read count was at 14 for Influenza A within the NTC

189 and had enough coverage and depth to be determined as present despite a lower calculated statistics for
190 the mapping. Interpretation was not impacted with the other remnant samples in the run based on pre-
191 established interpretation guidelines and criteria.

192 Accuracy

193 A total of 31 remnant samples consisting of 18 samples with original clinical positive single detections
194 and 13 samples with negative clinical detections were tested using the untargeted sequencing assay.
195 Concordance was determined by comparing the assay results to original clinical test results (Table 2 and
196 Table S2). Two clinical positive samples failed internal control QC in 3 independent runs. Of the other 16
197 clinical positive samples, 13 showed positive concordant detections or true-positives and 3 samples
198 showed unexpected negative detections and were classified as false-negatives for detection by the
199 untargeted assay. All 13 negative clinical samples were also negative for the untargeted sequencing assay.

200

Table S2. Single-detection remnant sample calls

Sample	Original Detection Method	Untargeted Sequencing Call
Human Metapneumovirus	Genmark Eplex PCR	Not Detected, Internal Control (IC) Quality Control (QC) failed
Human Metapneumovirus	Genmark Eplex PCR	Not Detected
Human Metapneumovirus	Genmark Eplex PCR	Human metapneumovirus, IC QC failed
Parainfluenza IV	Genmark Eplex PCR	Not Detected
Parainfluenza IV	Genmark Eplex PCR	Human parainfluenza virus 4a
Parainfluenza IV	Genmark Eplex PCR	Not Detected
SARS-CoV-2	Gene Xpert Infinity	Severe acute respiratory syndrome coronavirus 2
SARS-CoV-2	Gene Xpert Infinity	Severe acute respiratory syndrome coronavirus 2
SARS-CoV-2	Gene Xpert Infinity	Severe acute respiratory syndrome coronavirus 2
RSV	Biofire	Respiratory syncytial virus
RSV	Biofire	Respiratory syncytial virus
RSV	Biofire	Respiratory syncytial virus
Influenza A	Biofire	Influenza A
Influenza A	Diasorin Integrated Cyclor	Influenza A
Influenza A	Biofire	Influenza A
Enterovirus	Diasorin Integrated Cyclor	Rhinovirus A1
Enterovirus	Diasorin Integrated Cyclor	Rhinovirus A1
Enterovirus	Diasorin Integrated Cyclor	Rhinovirus A1

201

202 Overall, the untargeted sequencing assay showed 81.25 % sensitivity (13/16) and 100%
203 specificity (29/29) compared to original clinical positive sample results excluding the internal control QC
204 failures (Table 2). No cases, excluding the internal control QC failed samples, were classified as

205 untargeted sequencing assay false positives and two samples originally identified as Parainfluenza IV
206 were classified as false negatives; discrepancy testing was not performed.

207 Interference

208 We evaluated the effects of interference from blood, EDTA, and increasing additions of *Micrococcus*
209 *luteus* or *Staphylococcus epidermidis* to contrived plasma samples prior to extraction on the untargeted
210 sequencing assay. Exogenous addition of blood to plasma samples at 2% and 5% of the sample prior to
211 RNA extraction resulted in a significant 1.34 ± 0.04 -Log and 1.54 ± 0.11 -Log reduction in the number
212 of RNA virus reads, respectively (Figure 3) ANCOVA and TukeyHSD, $p < 0.001$). The MS2 Phage
213 internal control for the samples indicated only a 0.83-Log reduction in reads (Figure 3).

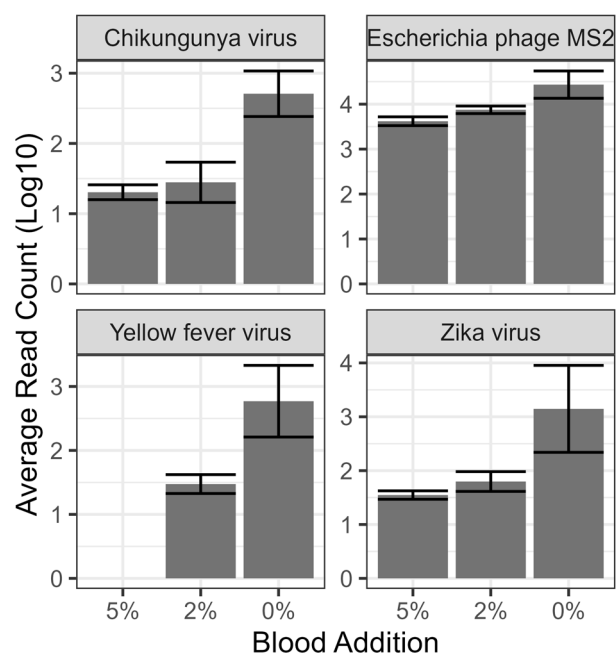


Figure 3. Interference of RNA virus reads with increasing blood concentration. Number of reads mapped to the individual viruses within contrived samples containing increasing additions of human whole blood.

214

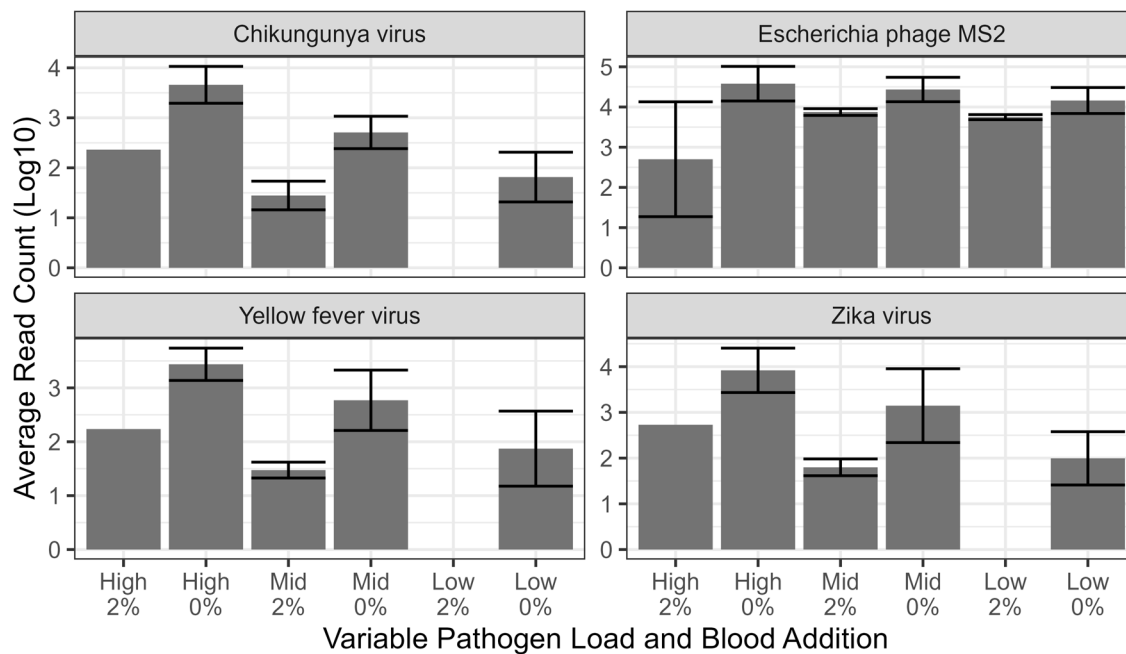


Figure 4. Interference from blood with increasing RNA virus concentration. Number of reads associated with individual viruses at different loading in contrived samples in the presence of 2% blood. Threshold for detection was a minimum of 10 reads leading to absence in graph and subsequent statistical analysis.

215

216 The Log reduction in the number of RNA virus reads was also consistent with increasing
217 concentrations of virus within the sample containing 2% blood (Figure 4). The addition of 2% blood to
218 samples with increasing concentrations of viral pathogens had a significant reduction of 1.31 ± 0.05 -Log
219 reads at all pathogen concentrations, with the exception of 'low' pathogen concentrations where the read
220 count dropped below the 10 read threshold and was removed from analysis (ANCOVA and TukeyHSD,
221 $p < 0.001$). The MS2 Phage internal control had only a 0.5-Log reduction in reads with greater variability
222 in the 'high' load viral pathogens in 2% blood.

223 Interference caused by the addition of increasing concentration of EDTA was also established
224 (Figure 5). The addition of low concentrations of EDTA as sodium salt at 5 mM to the existing
225 approximate 4.5 mM potassium EDTA for plasma collection may have increased the extraction efficiency
226 of some RNA viruses with an approximate increase in viral reads by 0.3-Log, yet was not statistically

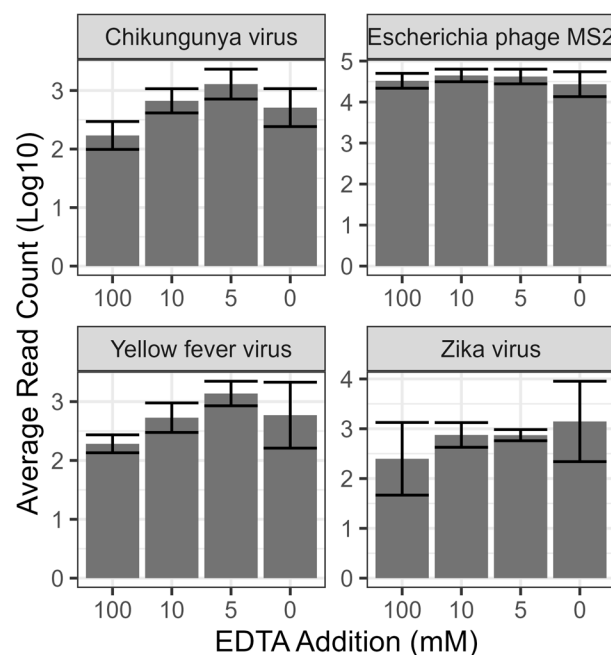


Figure 5. Interference in the presence of EDTA. Number of reads associated with different viruses within contrived samples with increasing concentrations of EDTA.

227 significant compared to no EDTA addition (ANCOVA and TukeyHSD, $p=0.69$). This potential increase
228 was eliminated at addition of 10 mM EDTA levels ($p=0.99$ compared to no addition) and negatively
229 impacted RNA virus reads by approximately 0.60 ± 0.13 -Log reduction at additions of 100 mM EDTA
230 compared to no addition ($p=0.005$).

231 Finally, we evaluated the addition of *Staphylococcus epidermidis* or *Micrococcus luteus* to known
232 positive plasma samples, neither of which negatively impacted viral read counts (Figure 6, ANCOVA:
233 $p=0.239$).

234

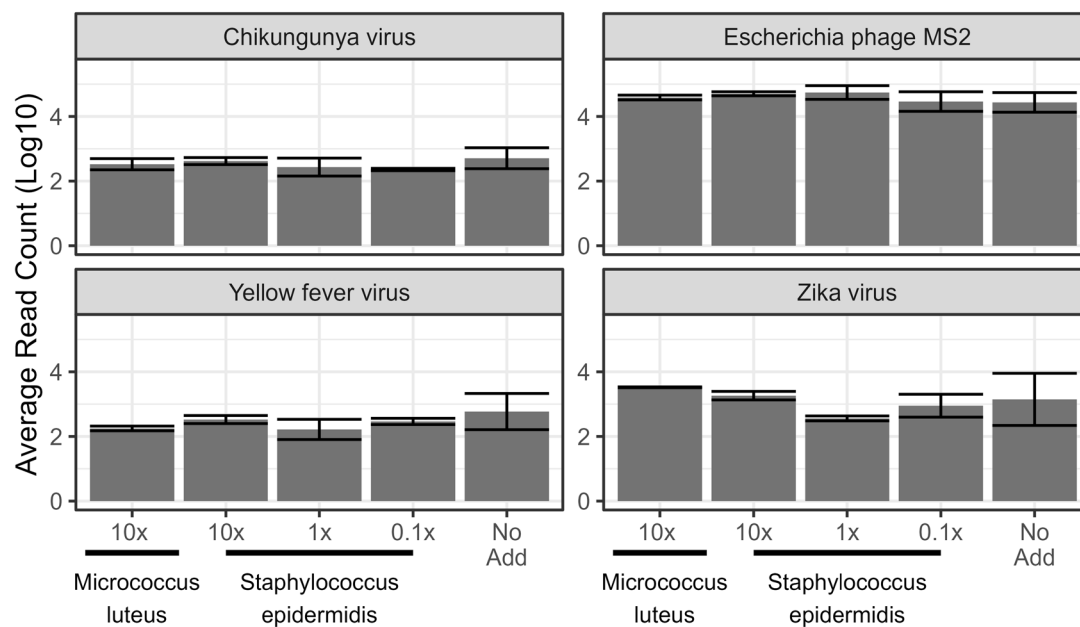


Figure 6. Interference by bacteria addition. Number of reads associated with different viruses in contrived samples with increasing concentrations of *Staphylococcus epidermidis* or *Micrococcus luteus*

235

236

237

238 The relative selectivity of sequences within the different clinical sample matrices, VTM and
239 Plasma, and the absence of a complex matrix, NTC, was examined using the MS2 Phage IC used between
240 the samples (Figure 7). There was no significant difference in the number of MS2 Phage reads between

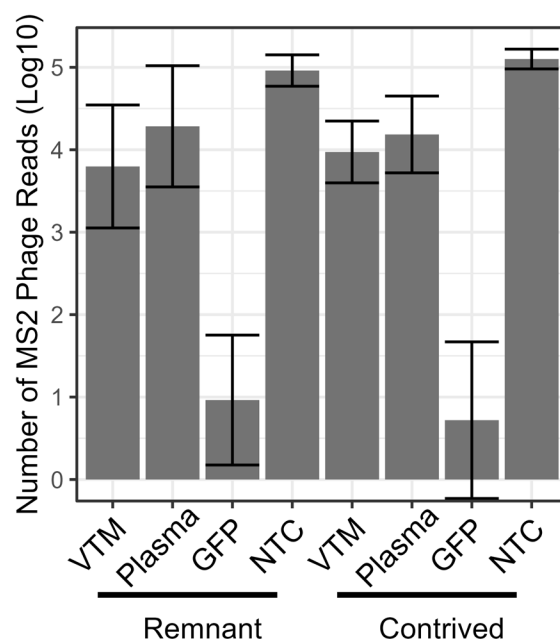


Figure 7. Number of MS2 Phage Reads in Different Sample Matrices.

241 the remnant and contrived samples runs (ANCOVA, $p=0.537$). There was also no significant difference in
242 the number of MS2 Phage reads between the plasma and VTM in either the remnant (TukeyHSD, $p=0.18$)
243 or contrived (TukeyHSD, $p=0.78$) runs. The number of MS2 Phage reads in VTM samples had a
244 significant decrease of 1.13 and 1.16-Log compared to NTC (TukeyHSD, $p<0.001$) in contrived and
245 remnant samples respectively. Plasma had a significant decrease in MS2 Phage reads of 0.92-Log in
246 contrived sample runs (TukeyHSD, $p<0.001$) but not in remnant sample runs (0.68-Log; TukeyHSD
247 $p=0.086$) compared to NTC. There was no significant difference between the number of human reads
248 between the contrived and remnant samples (ANCOVA, $p=0.66$) nor between VTM and plasma
249 (ANCOVA, $p=0.99$). The presence of human reads in the NTC and PC of approximately 300 reads
250 indicates how the majority reads from 6 samples can impact the contamination of reads in samples
251 expected not to contain human reads.

252

253 Challenge remnant samples

254 We evaluated remnant samples containing multiple virus detections by current molecular methods or non-
255 overlapping single calls (Table 3). Of the 17 diagnosed viruses present by current molecular diagnostics
256 in the 9 samples, 8 of those viruses were confirmed through untargeted sequencing (47%) without
257 incorrect identification of other viruses. Discrepancy testing was not performed.

258

Table 3. Multiple-detection remnant samples

Sample	Original Detection Method	Untargeted Sequencing Call
Human Metapneumovirus Rhinovirus	Genmark Eplex PCR	Not Detected Not Detected
Parainfluenza IV Influenza A	Genmark Eplex PCR	Not Detected Influenza A
Human Metapneumovirus SARS-CoV-2	Genmark Eplex PCR	Not Detected SARS-CoV-2
Human Metapneumovirus Parainfluenza IV	Genmark Eplex PCR	Not Detected Not Detected
Coronavirus OC43 Human Metapneumovirus	Biofire	Human coronavirus OC43 Human metapneumovirus
Parainfluenza II	Biofire	Human orthorubulavirus 2
Parainfluenza I	Biofire	Not Detected
Adenovirus Influenza A Parainfluenza III	Biofire	Not Detected Influenza A Simian Agent 10
Rhinovirus/Enterovirus Influenza A	Biofire	Not Detected Influenza A

259

260

261 Discussion

262 We have previously developed and evaluated a clinical untargeted sequencing assay intended to aid in the
263 diagnosis of respiratory and blood-borne infections by RNA viruses (Kappell et al. 2023). Here, we
264 rigorously evaluate the best performing of the sequencing methods for different RNA viruses within viral
265 transport media and human plasma through clinically relevant contrived samples and remnant patient
266 samples. We incorporated multiple QC materials, including an external positive control (*gfp* gene
267 transcript RNA), NTC samples that are run in parallel with clinical samples, and spiked MS2 Phage as an
268 internal process control. Reproducible threshold metrics were established and evaluated to enable
269 identification of pathogens from long-read nanopore sequencing data above background noise and to
270 minimize potential barcode crosstalk within multiplexed analyses. The final sequencing protocol

271 incorporated: (1) RNA extraction of MS2 Phage spiked samples and NTC, (2) reverse transcription of
272 RNA into double-stranded cDNA, (3) unbiased amplification through whole genome amplification, (4)
273 ligation of barcodes and adapters using ONT kits, (5) sequencing, basecalling, and demultiplexing using
274 ONT MinION™ (Mk1B or Mk1C), and (6) SigSciDx open-source bioinformatics pipeline for QC and
275 pathogen detection.

276 Our untargeted sequencing library preparation protocol uses a whole genome amplification
277 (WGA) approach after cDNA synthesis. This step increases assay specificity without requiring targeted
278 amplification. The WGA step has the potential to be further optimized to shorten the time of reaction, but
279 the material generated was sufficient for approximately 3 reactions after a 2-hour incubation, starting with
280 nearly undetectable quantities of cDNA. The amplified material serves as an ideal reference point for
281 additional sequencing or molecular analysis, if needed. WGA was performed using Qiagen REPLI-g kit
282 which has been shown to introduce the lowest amplification bias compared to other multiple displacement
283 amplification (MDA) methods (Pinard et al. 2006) inducing the least distortion in read counts per bin
284 across the length of a genome compared to unamplified controls. While sequencing independent single
285 primer amplification (SISPA) methods are generally quicker than MDA methods, MDA methods are
286 known to perform better than SISPA methods in sequencing output, taxonomic assignments, diversity,
287 and assembly statistics (Kallies et al. 2019; Parras-Moltó et al. 2018). The simplified execution of a MDA
288 based method in the Qiagen kit format has led to our adoption of the method for this untargeted
289 sequencing assay.

290 Contrived VTM and plasma samples in this study that contained 3 and 4 RNA viruses
291 respectively, indicated a range of limits of detection for the untargeted sequencing assay ranging from 1
292 TCID₅₀ /mL (est. 1 PFU/ mL) for Human coronavirus 229E to 26,915 TCID₅₀ /mL (est. 18,840 PFU /
293 mL) for Influenza A virus. Untargeted sequencing assay sensitivity for detection of a given organism is
294 dependent on multiple factors, including extraction efficiency, size of the genome, complexity of the
295 genome, the variability of genomic equivalences (sequenceable material) vs infectious particles, library

296 preparation bias, and availability of matching reference genomes in the database. We think the higher
297 levels of detection for Influenza A are likely due to genome complexity and size. While Influenza A has
298 an approximately 13.5 kb genome, it is fragmented into 8 segments ranging in size from 890 to 2,341 bp.
299 The subset of Influenza A genomes available in the NCBI Genbank of complete genomes is also limited,
300 preventing accurate strain level detection. We speculate that Human coronavirus 229E was detected at
301 lower LOD levels due to its large unsegmented 27.5 kb genome allowing multiple targets and ease of
302 progressivity for reverse-transcription and WGA. Another important context of the LOD for the
303 individual viruses used in the study is the variability of genome equivalents infectious particles in a viral
304 stock, similar to particle-to-PFU ratios (Bhat et al. 2022). For example, influenza has a particle-to-PFU
305 ratio of 20-to-50 and SARS-CoV-2 has a range of 10^4 -to- 10^6 for genomic RNA to PFU. As viral stocks
306 used for this study were reported in TCID₅₀ related to PFU, correcting for the particle-to-PFU ratio
307 consistent with these values would suggest a more consistent LOD between these viruses.

308 A key limitation for infectious disease diagnostics using untargeted sequencing assays is
309 background interference. The presence of blood as little as 2% volume of the sample led to a greater than
310 10-fold decrease in the number of reads, indicating a potential 1-Log decrease in analytical sensitivity of
311 the assay. This highlights the importance of the plasma capture step from a blood sample during initial
312 sample preparation to maximize assay sensitivity. Additionally, the presence of 100 mM EDTA caused a
313 0.6-Log reduction in reads and potential loss in analytical sensitivity. The presence of such an excess of
314 EDTA is unlikely, but a low volume of blood captured in a K₂-EDTA vacutainer tube would increase the
315 concentration of EDTA present in the extracted sample could be indicated as a potential cause for
316 sequencing failure. The spiked MS2 phage IC was useful for assessing the decreased assay sensitivity
317 caused by interferents or due to matrix effects. Notably, matrix effects on read counts from the MS2
318 Phage IC were more variable in the remnant samples compared to the NTC. While these matrix effects
319 had an impact on total read output, the read count following host removal was most noteworthy. Remnant
320 samples that had fewer mapped MS2 reads following host removal typically showed fewer or no

321 pathogen reads. This indicated potential interference from high human background or potentially other
322 confounding inhibitors. Accordingly, repeated failures or low read abundance of the MS2 Phage IC may
323 indicate poor assay performance. In this case, other diagnostic tests that are less sensitive to interference
324 should be considered (e.g., RT-PCR).

325 Of the 35 RNA viruses detected by conventional clinical molecular assays across 27 total remnant
326 samples, untargeted sequencing detected 22 concordant positives. Concordant negative detections
327 occurred in 13 out of 13 negative remnant samples. These findings show the unbiased sequencing assay
328 has 62.9% sensitivity and 100% specificity. The calculated 62.9% sensitivity refers to clinical sensitivity
329 in diagnosis of infection and not analytical sensitivity or detection of pathogen nucleic acid. Multiple
330 factors likely play a role in limiting the clinical sensitivity of untargeted sequencing, including: (1) the use
331 of remnant samples as a “gold standard”, with the potential that some samples may represent false
332 positives due to PCR artifacts or incorrect handling (e.g., contamination), (2) the potential that remnant
333 samples have undergone nucleic acid degradation from prior freeze-thaw steps, (3) the necessary use of
334 both-end barcode demultiplexing and robust pre-established thresholds to minimize false-positive
335 detections that could filter appropriate pathogen sequencing data, and (5) potential patient medications,
336 treatments, or other interferents that cause assay inhibition. Discrepancy testing was not performed to
337 determine the potential causes for the false-negative cases.

338 Approximately 14 of the 40 total remnant samples used in this study, including negative remnant
339 samples, had additional virus species or genera detected above pre-established thresholds consisting of
340 environmental or normal flora. Most of the additional detections were of phage, most likely due to the
341 mappings of the MS2 Phage to additional phage sequences and small amounts of contamination in
342 workflow reagents. Torque teno viruses, Torque teno midi virus, and TTV-like mini viruses were also
343 identified in several of the remnant VTM samples. These anelloviruses have small circular, negative-
344 sense, single-stranded DNA genomes which may have escaped the DNA removal column-based step in
345 the RNA extraction process. These anelloviruses are not directly concerning as they are diverse and

346 commensal members of the human virome (Bendinelli and Maggi 2010). However, their increase leading
347 to detection may indicate host immunosuppression (del Rosal et al. 2023; Prasetyo et al. 2015).
348 Determining the clinical significance of detecting viruses that may be contaminants or normal flora
349 remains a classical problem in clinical microbiology and often requires clinical context for interpretation.
350 While plasma is normally a sterile sample, manipulation and handling of the sample may introduce low
351 levels of contamination which ultimately gets amplified by our use of WGA within the untargeted
352 sequencing protocol. Additionally, the collection using nasopharyngeal swabs into VTM in a non-sterile
353 site collection environment presents a reasonable risk of several potential environmental and normal flora
354 being captured along with viral pathogens. Thus, detection of multiple viruses as potential contamination
355 or normal flora were noted and were considered negative in our evaluation of assay performance.

356 Untargeted sequencing testing can provide broad-spectrum RNA virus detection and
357 identification, however assessment of the clinical significance of the reported findings may require
358 interpretation. Clinical context of the results from untargeted sequencing is important in discussing and
359 reviewing patient cases with the treating clinicians. The potential for untargeted sequencing for RNA
360 virus identification includes characterization of antiviral or vaccination escape mutations, genotyping or
361 strain-level identification, and presence of reads from potential pathogens below formal threshold that can
362 inform follow-up targeted testing.

363

364 **Methods**

365 **Remnant Clinical Samples**

366 Clinically diagnosed, de-identified remnant samples were obtained from Discovery Life Sciences and
367 Precision for Medicine. Of the 40 remnant samples evaluated, 9 VTM samples were negative based on
368 molecular respiratory panels, 4 were negative serum or plasma samples, and 27 VTM samples were
369 positive for single or multiple respiratory viruses. Of the remnant samples, 6 different RNA virus

370 pathogens as determined by standard clinical testing were represented by 3 samples each. These 18
 371 remnant VTM samples were used for the initial accuracy study. The other 9 VTM samples containing
 372 mixed or singlet samples with one diagnosed pathogen were kept for additional challenge samples.
 373 Remnant samples were pipetted into 140 μ L single use aliquots for extraction. The number of prior freeze
 374 thaws of remnant samples was unknown before receipt. Samples were stored at -80 $^{\circ}$ C upon receipt and
 375 again after aliquots were prepared.

376 **Contrived Sample Preparation**

377 Aliquots of negative diagnosed remnant VTM or plasma were pooled to generate a matrix for contrived
 378 RNA virus samples. Contrived samples in VTM were prepared using Influenza A virus (BEI Resources
 379 NR-42007), Respiratory Syncytial virus (BEI Resources NR-28529), and Human coronavirus (HCoV)
 380 229E (NR-52726). Contrived samples in plasma were prepared using Chikungunya virus (NR-56523),
 381 Hepatitis A virus (BEI Resources NR-137), Yellow Fever virus (BEI Resources NR-116), and Zika virus
 382 (BEI Resources NR-50065). To determine LoD, contrived samples were prepared at a starting dilution of
 383 25 μ L/mL stock material in pooled VTM or plasma. Tenfold serial dilutions were prepared in respective
 384 pre-screened negative matrices to a final concentration of 0.00025 μ L/mL. The viral load levels in
 385 TCID₅₀/ mL at the different serial dilutions are indicated in Table 4. All samples were prepared into single
 386 use 140 μ L sample aliquots and frozen at -80 $^{\circ}$ C until extraction.

387

Table 4 Contrived Sample Loading

Media:		Viral Transport Media			Human Plasma			
Virus:		Influenza A Virus (A/Wisconsin/15/ 2009 (H3N2))	Human Respiratory Syncytial Virus (A 1998/ 3-2)	Human Coronavirus (229E)	Zika Virus (MR 766)	Hepatitis A Virus (HM175/18f)	Yellow Fever Virus (17D)	Chikungunya Virus (181/25)
BEI Cat#		NR-42007	NR-28529	NR-52726	NR-50065	NR-137	NR-116	NR-56523
Viral Load Levels Log ₁₀ (TCID ₅₀ / mL)	A*	7.10	5.60	3.85	7.35	6.85	6.60	6.35
	B	6.10	4.60	2.85	6.35	5.85	5.60	5.35
	C	5.10	3.60	1.85	5.35	4.85	4.60	4.35
	D	4.10	2.60	0.85	4.35	3.85	3.60	3.35
	E	3.10	1.60	-0.15	3.35	2.85	2.60	2.35
	F	2.10	0.60	-1.15	2.35	1.85	1.60	1.35
	G	NA	NA	NA	1.35	0.85	0.60	0.35

*The letters in column 1 provide a designation for the specific viral load of each virus in the corresponding row. Viral loads decrease as the letters continue, with "A" representing the highest load level.

388

389 **Untargeted Third-Generation Sequencing**

390 RNA from contrived and remnant clinical samples was extracted using the Qiagen RNeasy Plus Micro
391 Kit (Qiagen #74034) with minor changes to the manufacturer's standard protocol adapting to the specific
392 sample matrix and input volume of 140 μ L sample material augmented with 5 μ L of MS2 Phage internal
393 control (ZeptoMetrix, 0810274; 145 μ L total input). Extracted RNA was immediately reverse-transcribed
394 using random hexamers and Maxima H Minus Double-Stranded cDNA Synthesis Kit (ThermoFisher
395 K2561) following manufacturer's instructions including the recommended changes to increase the
396 addition of 1st strand enzyme mix to 2 μ L and the incubation temperature to 55 $^{\circ}$ C. In addition to the
397 extracted RNA, an RNA transcript of the *gfp* gene (System Biosciences, LLC MR700A-1) was included
398 as an additional positive control sample at 6 pg. If necessary, DNA forms of samples were stored at -20
399 $^{\circ}$ C only after bead clean-up of the previous reaction steps. Double-stranded cDNA was cleaned using
400 AMPure XP beads before whole genome replication using the Qiagen REPLI-g Advanced DNA Single
401 Cell Kit (Qiagen, 150363) following the Qiagen supplementary protocol "Whole genome amplification
402 from genomic DNA using the REPLI-g Advanced DNA Single Cell Kit with increased sample volumes."
403 Briefly, 15 μ L of double-stranded cDNA was denatured with 2 μ L of Advanced Buffer DLB for 3
404 minutes at 25 $^{\circ}$ C. Stop Solution was added at 3 μ L and mixed. The 29 μ L REPLI-g sc Advanced Reaction
405 Buffer with 2 μ L REPLI-g sc DNA Polymerase was added directly to the denatured cDNA and incubated
406 for 2 hours at 30 $^{\circ}$ C. The REPLI-g sc DNA polymerase was inactivated at 65 $^{\circ}$ C for 3 minutes. The
407 resulting amplified DNA was purified from the reaction by AMPure XP bead clean up and digested using
408 T7 Endonuclease I (New England Biolabs (NEB), M0302L). For effective removal of T7 digested
409 fragments, a custom AMPure XP bead solution was made using PEG 8000 50% (w/v) (Fisher Scientific,
410 NC1017553) as described in Oxford Nanopore Technology's (ONT) whole genome amplification
411 protocol for LSK110 (v110, rev. 10 Nov 2020, Oxford Nanopore Technologies 2023). Once purified, the
412 sequencing library was prepared by using end repair and ligation following the ONT's protocol for the

413 Ligation sequencing kit (SQK-LSK110). Native barcode expansion kit 1-12 (EXP-NBD104) was used for
414 multiplexing samples. All bead cleanups were done on a microfuge tube magnetic separation stand
415 (Permagen). Sequencing was performed on Oxford Nanopore Technologies MinION™ Mk1B or Mk1C
416 using R9.4.1 flow cells (FLO-MIN106D). Each flow cell was primed and loaded using manufacturer's
417 instructions. Each run used 24-hour run default settings except no reserve pores. 'Live' basecalling was
418 performed in 'Fast basecalling mode' and single-end barcode demultiplexing. The 'live' basecalling was
419 performed with the latest MinKNOW (22.12.5 to 23.04.03) running the Mk1C or Mk1B which did not
420 impact 'Fast basecalling' (guppy 6.4.6 - 6.5.7) based on change logs and no changes on re-basecalling of
421 the earliest run. Alternative basecalling and demultiplexing including both-end barcode demultiplexing of
422 the sequencing runs were performed using standalone 'guppy' basecaller version 6.2.1. All data analyzed
423 was basecalled with version 6.2.1 in both-end barcode demultiplexing as a final workflow, with the
424 exception for initial comparison of single-end and double-end barcode demultiplexing (Figure 2, Table 1).

425 **Bioinformatics Analysis**

426 Analyses were developed using EPI2ME Labs (<https://github.com/epi2me-labs>), a Jupyter-notebook
427 platform allowing users to run pre-written python, R, or BASH based analyses in a web-browser based
428 environment. For each sequencing run, passing reads (default Q8 threshold) were concatenated for each
429 barcode, and quality of run statistics (e.g., read length and quality distribution) were reported and
430 visualized in the web-based environment. Human reads were removed and reported after mapping to the
431 human genome using 'minimap2' (Li 2021, 2018) with default 'ont' parameters. The remaining human
432 depleted reads were then mapped to the control sequences of the *gfp* gene and MS2 Phage internal
433 controls using minimap2, while removing the *gfp* mapped reads. The number and percentage of reads
434 aligned to the human genome, *gfp* gene, and MS2 Phage were analyzed using SAMtools (Danecek et al.
435 2021; Li et al. 2009) and reported in the Quality Control (QC) portion of the output for review.
436 Untargeted sequence analysis was performed using Centrifuge (Kim et al. 2016) with a customized virus

437 database to identify putative NCBI taxonomy identifiers (taxids) and input those taxids into the ‘reference
438 inference module’ of SeqScreen-Nano (Balaji et al. 2023b, 2023a).

439 The customized Centrifuge database is formed of the ‘Complete Genome’ or ‘Chromosome’ from the
440 viral genomes from NCBI’s GenBank database, excluding all partial or incomplete genome sequences,
441 using the centrifuge-download command. The centrifuge-build command was used as described by
442 Centrifuge developers (Centrifuge Developers 2023; Kim et al. 2016) to generate the custom virus
443 sequence database to use with the Centrifuge taxonomical classifier. The taxids reported by Centrifuge
444 were filtered to include only those taxids with greater than 10 reads assigned to a species or leaf
445 taxonomical level. The ‘reference inference module’ from SeqScreen-Nano (SeqScreen version 4.0) then
446 used the filtered taxids from the Centrifuge to download the corresponding reference genomes from
447 NCBI. All reads were then independently mapped using 'minimap2' to each of those individual reference
448 genomes, and mapping metrics (e.g., depth, breadth of coverage) were calculated by SeqScreen-Nano.
449 High-coverage reference genomes were concatenated by SeqScreen-Nano's ‘reference inference module’
450 before the second alignment step, where the reads were remapped to the concatenated reference to allow
451 each read to competitively map to the single best reference. Mapping metrics were re-calculated and
452 further statistical analysis was performed to determine the likelihood of presence or absence of individual
453 genomes based on a rubric. The SeqScreen-Nano ‘reference inference module’ analysis, metrics, and
454 rubrics are originally described in Balaji et al. (2023a). Modification to the threshold criteria was made to
455 be more inclusive of the reference genomes taken from the first mapping of individual reference genomes
456 to the second mapping of a concatenated reference. This included reducing the original calculated in
457 SeqScreen-Nano ‘coverage score’ threshold, which is the ratio of the observed ‘breadth of coverage’ and
458 ‘expected breadth of coverage’ for each taxon in a sample (Balaji et al. 2023a), from a minimum of 0.70
459 to 0.10. The large ratio of 0.7 for the ‘coverage score’ was not ideal for smaller genomes for viruses,
460 which had inflated expected coverages compared to genomes sizes of bacteria. This allowed for more
461 reference viral genomes within the sample to be used in second stage mapping, while still removing

462 potential contaminants. Additionally, signals from low abundance or partially covered genomes of close
463 neighboring viral taxons were successfully reduced by requiring a greater than 0.15 ‘breadth of coverage’
464 for a reference genome to move forward in the second stage of mapping. These changes increased the
465 inclusion of viral genomes known to be present within a sample while also reducing background and
466 known contamination. The rubrics within SeqScreen v4.0 for “Present” (Species) and “Genus Present”
467 calls were not modified.

468

469 **Evaluation of Untargeted Sequencing and Analysis**

470 Statistical analysis and visualizations were performed using R version 4.3.1 (R Core Team 2023)
471 and RStudio version 2023.06.01.524 (Posit team 2023). Visualization of the data was performed using the
472 ‘ggplot2’ R package (Wickham 2016).

473 Limits of detection were determined for each of the seven representative RNA viruses in their
474 respective clinical matrices by probit analysis using a series of dilutions across a minimum of 6-Log range
475 (Table 4). The contrived samples in the plasma matrix were sequenced across a total of 8 multiplex runs
476 and samples in VTM were sequenced across a total of 6 multiplex runs. Assay LoD were calculated in R
477 and RStudio using probit regression model using the ‘glm’, generalized linear model, function as the
478 concentration at which RNA virus was detected in 95% of replicates with at least 2 samples performed at
479 each tested concentration (Table S3).

480

Table S3. Number of contrived samples at virus load sequenced

Load Level (See Table 5)	Plasma	VTM
A	2	4
B	2	8
C	6	10
D	8	8
E	14	4
F	10	2
G	6	NA

481

482 Precision was determined using repeat analysis of the contrived samples for RNA viruses at the
483 load above and nearest their individual LoD as determined by probit analysis (LoD load) across 4 or 5
484 runs with VTM contrived samples and 4 or 7 runs with plasma contrived samples for inter-assay
485 reproducibility. Precision was also determined using repeat analysis of NTC samples across all 14 runs
486 for inter-assay reproducibility. Replicate remnant samples were also used for inter-assay reproducibility
487 ranging from 2 to 3 independent runs performed for each sample performed across a total of 16 runs.
488 Precision for intra-assay reproducibility was determined for each individual RNA virus for duplicate
489 contrived samples at LoD load performed on 4 to 7 runs. Intra-assay reproducibility was also determined
490 for negative detection remnant samples representing 3 (plasma/serum) or all 6 (VTM) samples across a
491 total of 2 runs.

492 Accuracy was determined using 24 remnant VTM samples comprised of 18 positive samples
493 containing 6 detected organisms (three samples for each organism) and 6 negative samples. Sensitivity
494 and specificity of the untargeted sequencing assay were calculated relative to prior ‘gold standard’ clinical
495 molecular method, which was information provided with the obtained remnant samples.

496 To evaluate the effects of potential interfering substances, EDTA, human blood, and microbial
497 contamination simulated by *Micrococcus luteus* or *Staphylococcus epidermidis* were added to contrived
498 plasma samples. The addition of interfering substances was performed on plasma samples with RNA viral
499 load level E (Table 4). EDTA was added to an additional 5 mM, 10 mM, and 100 mM EDTA disodium
500 salt to the plasma samples prior to RNA extraction and downstream workflow. Blood was added to
501 represent 2% or 5% of the total volume of the plasma sample. Additional analysis of the effects of human
502 blood on the untargeted sequencing assay was performed on plasma samples with RNA virus loads of D,
503 E, and F with the addition of 2% blood. For addition of *Micrococcus luteus* or *Staphylococcus*
504 *epidermidis*, an overnight culture was washed twice with sterile molecular grade water. The cells were
505 resuspended at high concentration and quantified by OD₆₀₀. The bacteria were added to contrived plasma
506 samples to generate bacterial concentrations at 10⁷ (high), 10⁶ (medium), and 10⁵ (low) cells per mL. The

507 number of reads were determined by the Centrifuge results mapping to the individual RNA viruses. The
508 read counts were normalized by each organism prior to performing ANCOVA in R to determine if
509 significant differences existed between the concentration of interfering substances and no addition. Post-
510 hoc comparisons between the concentration and no addition were made using TukeyHSD (honestly
511 significant difference) testing.

512 An additional examination of remnant samples containing multiple viral detections or single
513 sample representatives of a virus was performed. These 9 samples represented 10 different viruses and
514 were sequenced using the untargeted sequencing assay and concordance with ‘gold standard’ clinical
515 molecular testing was examined.

516 **Software availability**

517 The “SigSciDx” computational software used in this study, along with the custom Centrifuge database, is
518 publicly available as the following docker image on Docker Hub: sigsci/sigscidx (Signature Science
519 Team 2023b). The version used for analysis in this publication is 1.0.1. A bash command to initiate the
520 EPI2ME Labs framework is available on GitHub (Signature Science Team 2023c). Note that the database
521 used for the custom viral Centrifuge database was from the June 2023 distribution of GenBank genome
522 assemblies.

523 **Data Access**

524 Metagenomic reads from contrived and remnant samples from this study were depleted of human host
525 sequences and will be submitted to the NCBI BioProject database (Signature Science Team 2023a). The
526 MS2 Phage genome sequence used was NC001417.2 for mapping. The *gfp* gene sequence used for
527 mapping was obtained through direct communication with the vendor, System Biosciences, LLC
528 (MR700A-1). Both MS2 Phage and *gfp* gene sequences from the specific vendors are included in the
529 SigSciDx docker image.

530 **Competing interests**

531 The authors declare that they have no competing interests.

532 **Acknowledgements**

533 The authors would like to thank Leslie Parke for oversight and project management for this effort. The
534 authors would also like to thank Jim Gibson for his assistance creating the workflow figure. The authors
535 would also like to thank Leslie Parke for her review. This work was supported by the Centers for Disease
536 Control and Prevention under contract number 75D30122C15359.

537 **Authors' contributions**

538 All authors have read and approved the manuscript. Conceptualization: ADK, KQS, MBS, FCH. Data
539 curation: ADK, MBS, NCK. Laboratory analysis: ANS, NCK, KQS, EAS, CJS, ADK, LWA, TAW.
540 Funding acquisition: ADK, FCH, KLT, KQS. Data analysis: ADK, MBS, NCK, KLT, FCH. Project
541 administration: ADK, FCH. Writing – original draft: ADK, NCK, MBS, KQS, VMJ, FCH, KLT. Writing
542 – review & editing: ADK, KQS, MBS, LWA, KLT, FCH.

543 **References**

- 544 Arévalo MT, Karavis MA, Katoski SE, Harris JV, Hill JM, Deshpande SV, Roth PA, Liem AT,
545 Bernhards RC. 2022. A Rapid, Whole Genome Sequencing Assay for Detection and
546 Characterization of Novel Coronavirus (SARS-CoV-2) Clinical Specimens Using Nanopore
547 Sequencing. *Front Microbiol* **13**. <https://www.frontiersin.org/articles/10.3389/fmicb.2022.910955>
548 (Accessed September 8, 2023).
- 549 Balaji A, Liu Y, Nute MG, Hu B, Kappell A, LeSassier DS, Godbold GD, Ternus KL, Treangen TJ.
550 2023a. SeqScreen-Nano: a computational platform for rapid, in-field characterization of
551 previously unseen pathogens. 2023.02.10.528096.
552 <https://www.biorxiv.org/content/10.1101/2023.02.10.528096v2> (Accessed September 7, 2023).
- 553 Balaji A, Liu Y, Nute MG, Treangen TJ. 2023b. SeqScreen's reference inference module.
554 https://gitlab.com/treangenlab/seqscreen/-/blob/master/scripts/reference_inference.py.
- 555 Bendinelli M, Maggi F. 2010. TT Virus and Other Anelloviruses. In *Topley & Wilson's Microbiology and*
556 *Microbial Infections*, John Wiley & Sons, Ltd
557 <https://onlinelibrary.wiley.com/doi/abs/10.1002/9780470688618.taw0263> (Accessed September
558 7, 2023).
- 559 Bhat T, Cao A, Yin J. 2022. Virus-like Particles: Measures and Biological Functions. *Viruses* **14**: 383.

- 560 Bibby K. 2013. Metagenomic identification of viral pathogens. *Trends Biotechnol* **31**: 275–279.
- 561 Bouchiat C, Ginevra C, Benito Y, Gaillard T, Salord H, Dauwalder O, Laurent F, Vandenesch F. 2022.
562 Improving the Diagnosis of Bacterial Infections: Evaluation of 16S rRNA Nanopore
563 Metagenomics in Culture-Negative Samples. *Front Microbiol* **13**.
564 <https://www.frontiersin.org/articles/10.3389/fmicb.2022.943441> (Accessed September 8, 2023).
- 565 Centers for Disease Control and Prevention. 2022. SARS-CoV-2 Sequencing Resources.
566 https://github.com/CDCgov/SARS-CoV-2_Sequencing (Accessed August 22, 2022).
- 567 Centrifuge Developers. 2023. *Manual for Centrifuge, Database download and index building*.
568 <https://ccb.jhu.edu/software/centrifuge/manual.shtml#database-download-and-index-building>.
- 569 Charalampous T, Kay GL, Richardson H, Aydin A, Baldan R, Jeanes C, Rae D, Grundy S, Turner DJ,
570 Wain J, et al. 2019. Nanopore metagenomics enables rapid clinical diagnosis of bacterial lower
571 respiratory infection. *Nat Biotechnol* **37**: 783–792.
- 572 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T, McCarthy
573 SA, Davies RM, et al. 2021. Twelve years of SAMtools and BCFtools. *GigaScience* **10**: giab008.
- 574 del Rosal T, García-García ML, Casas I, Iglesias-Caballero M, Pozo F, Alcolea S, Bravo B, Rodrigo-
575 Muñoz JM, del Pozo V, Calvo C. 2023. Torque Teno Virus in Nasopharyngeal Aspirate of
576 Children With Viral Respiratory Infections. *Pediatr Infect Dis J* **42**: 184.
- 577 Fendrick AM, Monto AS, Nightengale B, Sarnes M. 2003. The Economic Burden of Non-Influenza-
578 Related Viral Respiratory Tract Infection in the United States. *Arch Intern Med* **163**: 487–494.
- 579 Hewitt FC, Guertin SL, Ternus KL, Schulte K, Kadavy DR. 2017. Toward Rapid Sequenced-Based
580 Detection and Characterization of Causative Agents of Bacteremia. 162735.
581 <https://www.biorxiv.org/content/10.1101/162735v1> (Accessed August 22, 2022).
- 582 Kallies R, Hölzer M, Brizola Toscan R, Nunes da Rocha U, Anders J, Marz M, Chatzinotas A. 2019.
583 Evaluation of Sequencing Library Preparation Protocols for Viral Metagenomic Analysis from
584 Pristine Aquifer Groundwaters. *Viruses* **11**: 484.
- 585 Kappell AD, Scholes AN, Scholz MB, Keplinger NC, Allen LW, Murray MC, Ternus KL, Hewitt FC.
586 2023. Unbiased metagenomic detection of RNA viruses for rapid identification of viral pathogens
587 in clinical samples. *Prepration*.
- 588 Kim D, Song L, Breitwieser FP, Salzberg SL. 2016. Centrifuge: rapid and sensitive classification of
589 metagenomic sequences. *Genome Res* **26**: 1721–1729.
- 590 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–3100.
- 591 Li H. 2021. New strategies to improve minimap2 alignment accuracy. *Bioinformatics* **37**: 4572–4574.
- 592 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The
593 Sequence Alignment/Map format and SAMtools. *Bioinforma Oxf Engl* **25**: 2078–9.
- 594 Miller RR, Montoya V, Gardy JL, Patrick DM, Tang P. 2013. Metagenomics for pathogen detection in
595 public health. *Genome Med* **5**: 81.

- 596 Miller S, Chiu C. 2022. The Role of Metagenomics and Next-Generation Sequencing in Infectious
597 Disease Diagnosis. *Clin Chem* **68**: 115–124.
- 598 Miller S, Naccache SN, Samayoa E, Messacar K, Arevalo S, Federman S, Stryke D, Pham E, Fung B,
599 Bolosky WJ, et al. 2019. Laboratory validation of a clinical metagenomic sequencing assay for
600 pathogen detection in cerebrospinal fluid. *Genome Res* **29**: 831–842.
- 601 Ohta A, Nishi K, Hirota K, Matsuo Y. 2023. Using nanopore sequencing to identify fungi from clinical
602 samples with high phylogenetic resolution. *Sci Rep* **13**: 9785.
- 603 Oxford Nanopore Technologies. 2023. *LIGATION SEQUENCING GDNA - WHOLE GENOME*
604 *AMPLIFICATION (SQK-LSK110)*.
605 [https://community.nanoporetech.com/docs/prepare/library_prep_protocols/premium-whole-](https://community.nanoporetech.com/docs/prepare/library_prep_protocols/premium-whole-genome-amplification-sqk-lsk110/v/wal_9115_v110_revh_10nov2020/whole-genome-amplification)
606 [genome-amplification-sqk-lsk110/v/wal_9115_v110_revh_10nov2020/whole-genome-](https://community.nanoporetech.com/docs/prepare/library_prep_protocols/premium-whole-genome-amplification-sqk-lsk110/v/wal_9115_v110_revh_10nov2020/whole-genome-amplification)
607 [amplification](https://community.nanoporetech.com/docs/prepare/library_prep_protocols/premium-whole-genome-amplification-sqk-lsk110/v/wal_9115_v110_revh_10nov2020/whole-genome-amplification).
- 608 Parras-Moltó M, Rodríguez-Galet A, Suárez-Rodríguez P, López-Bueno A. 2018. Evaluation of bias
609 induced by viral enrichment and random amplification protocols in metagenomic surveys of
610 saliva DNA viruses. *Microbiome* **6**: 119.
- 611 Petersen LM, Martin IW, Moschetti WE, Kershaw CM, Tsongalis GJ. 2019. Third-Generation
612 Sequencing in the Clinical Laboratory: Exploring the Advantages and Challenges of Nanopore
613 Sequencing. *J Clin Microbiol* **58**: 10.1128/jcm.01315-19.
- 614 Pinard R, de Winter A, Sarkis GJ, Gerstein MB, Tartaro KR, Plant RN, Egholm M, Rothberg JM,
615 Leamon JH. 2006. Assessment of whole genome amplification-induced bias through high-
616 throughput, massively parallel whole genome sequencing. *BMC Genomics* **7**: 216.
- 617 Posit team. 2023. RStudio: Integrated Development Environment for R. <http://www.posit.co/>.
- 618 Prasetyo AA, Desyardi MN, Tanamas J, Suradi, Reviono, Harsini, Kageyama S, Chikumi H, Shimizu E.
619 2015. Respiratory Viruses and Torque Teno Virus in Adults with Acute Respiratory Infections.
620 *Intervirology* **58**: 57–68.
- 621 R Core Team. 2023. R: A Language and Environment for Statistical Computing. [https://www.R-](https://www.R-project.org/)
622 [project.org/](https://www.R-project.org/).
- 623 Russell JA, Campos B, Stone J, Blosser EM, Burkett-Cadena N, Jacobs JL. 2018. Unbiased Strain-Typing
624 of Arbovirus Directly from Mosquitoes Using Nanopore Sequencing: A Field-forward
625 Biosurveillance Protocol. *Sci Rep* **8**: 5417.
- 626 Schlaberg R, Chiu CY, Miller S, Procop GW, Weinstock G, the Professional Practice Committee and
627 Committee on Laboratory Practices of the American Society for Microbiology, the Microbiology
628 Resource Committee of the College of American Pathologists. 2017. Validation of Metagenomic
629 Next-Generation Sequencing Tests for Universal Pathogen Detection. *Arch Pathol Lab Med* **141**:
630 776–786.
- 631 Signature Science Team. 2023a. NCBI BioProject Database Respository.
632 <https://www.ncbi.nlm.nih.gov/bioproject>.

- 633 Signature Science Team. 2023b. SigSciDx.
634 <https://hub.docker.com/repository/docker/sigsci/sigscidx/general>.
- 635 Signature Science Team. 2023c. SigSciDx Docker Start Command.
636 https://github.com/signaturescience/lrn_linux_epi2me_start.
- 637 Thompson WW, Shay DK, Weintraub E, Brammer L, Bridges CB, Cox NJ, Fukuda K. 2004. Influenza-
638 Associated Hospitalizations in the United States. *JAMA* **292**: 1333–1340.
- 639 Wickham H. 2016. *ggplot2*. Springer International Publishing, Cham
640 <http://link.springer.com/10.1007/978-3-319-24277-4> (Accessed September 7, 2023).
- 641 Xu Y, Lewandowski K, Lumley S, Pullan S, Vipond R, Carroll M, Foster D, Matthews PC, Peto T, Crook
642 D. 2018. Detection of Viral Pathogens With Multiplex Nanopore MinION Sequencing: Be
643 Careful With Cross-Talk. *Front Microbiol* **9**: 2225.
- 644 Yakovleva A, Kovalenko G, Redlinger M, Liulchuk MG, Bortz E, Zadorozhna VI, Scherbinska AM,
645 Wertheim JO, Goodfellow I, Meredith L, et al. 2022. Tracking SARS-COV-2 variants using
646 Nanopore sequencing in Ukraine in 2021. *Sci Rep* **12**: 15749.
- 647 Young-Xu Y, Aalst R van, Russo E, Lee JKH, Chit A. 2017. The Annual Burden of Seasonal Influenza in
648 the US Veterans Affairs Population. *PLOS ONE* **12**: e0169344.
- 649