

1 **Title**

2 Tomato brown rugose fruit virus Mo gene is a novel microbial source tracking marker

3

4 **Authors**

5 Aravind Natarajan\*<sup>1,2</sup>, Brayon J. Fremin\*<sup>1</sup>, Danica T. Schmidtke<sup>3</sup>, Marlene K. Wolfe<sup>4</sup>, Soumaya  
6 Zlitni<sup>1,2</sup>, Katherine E. Graham<sup>4#</sup>, Erin F. Brooks<sup>2</sup>, Christopher J. Severyn<sup>5</sup>, Kathleen M.  
7 Sakamoto<sup>5</sup>, Norman J. Lacayo<sup>5</sup>, Scott Kuersten<sup>6</sup>, Jeff Koble<sup>6</sup>, Glorianna Caves<sup>6</sup>, Inna Kaplan<sup>7</sup>,  
8 Upinder Singh<sup>8</sup>, Prasanna Jagannathan<sup>8,9</sup>, Andrew R. Rezvani<sup>7</sup>, Ami S. Bhatt<sup>1,2,\*\*</sup>, Alexandria B.  
9 Boehm<sup>4,\*\*</sup>

10 *\*Aravind Natarajan and Brayon J Fremin contributed equally to this work. Order of co-first*  
11 *authors was determined by contribution to writing the manuscript.*

12

13 **Affiliations**

14 <sup>1</sup> Department of Genetics, Stanford University, Stanford, CA, USA.

15 <sup>2</sup> Department of Medicine (Hematology, Blood and Marrow Transplantation), Stanford  
16 University, Stanford, CA, USA.

17 <sup>3</sup> Department of Microbiology and Immunology, Stanford University, Stanford, CA, USA.

18 <sup>4</sup> Department of Civil and Environmental Engineering, Stanford University, Stanford, CA, USA.

19 <sup>5</sup> Department of Pediatrics, (Hematology/Oncology/Stem Cell Transplant & Regenerative  
20 Medicine), Stanford University, Stanford, CA, USA.

21 <sup>6</sup> Illumina, Inc., San Diego, CA

22 <sup>7</sup> Department of Medicine (Blood and Marrow Transplantation and Cellular Therapy), Stanford  
23 University, Stanford, CA, USA.

24 <sup>8</sup> Department of Medicine (Infectious Diseases and Geographic Medicine), Stanford University,  
25 Stanford, CA, USA

26 <sup>9</sup> Department of Microbiology and Immunology, Stanford University, Stanford, CA, USA

27

28 **Present address**

29 # Katherine E. Graham: School of Civil and Environmental Engineering, Georgia Institute of  
30 Technology, Atlanta , GA 30332, USA.

31

32 **Corresponding authors**

33 \*\*

34 Ami S. Bhatt, Center for Clinical Sciences Research RM. 1155b, Stanford University, Stanford,  
35 CA, 94305. Tel: (650) 498-4438; Email: [asbhatt@stanford.edu](mailto:asbhatt@stanford.edu).

36 Alexandria Boehm, Jerry Yang & Akiko Yamazaki Environment & Energy Building RM. 189,  
37 Stanford University, Stanford, CA 94305. Tel: (650) 724-9128; Email: [aboehm@stanford.edu](mailto:aboehm@stanford.edu).

38

39

40

41 **Abstract**

42           Microbial source tracking (MST) identifies sources of fecal contamination in the  
43 environment using fecal host-associated markers. While there are numerous bacterial MST  
44 markers, there are few viral markers. Here we design and test novel viral MST markers based  
45 on tomato brown rugose fruit virus (ToBRFV) genomes. We assembled eight nearly complete  
46 genomes of ToBRFV from wastewater and stool samples from the San Francisco Bay Area in  
47 the United States of America. Next, we developed two novel probe-based RT-PCR assays  
48 based on conserved regions of the ToBRFV genome, and tested the markers' sensitivities and  
49 specificities using human and non-human animal stool as well as wastewater. The ToBRFV  
50 markers are sensitive and specific; in human stool and wastewater, they are more prevalent and  
51 abundant than a currently used marker, the pepper mild mottle virus (PMMoV) coat protein (CP)  
52 gene. We applied the assays to detect fecal contamination in urban stormwater samples and  
53 found that the ToBRFV markers matched cross-assembly phage (crAssphage), an established  
54 viral MST marker, in prevalence across samples. Taken together, ToBRFV is a promising viral  
55 human-associated MST marker.

56

57

58 **Importance.** Human exposure to fecal contamination in the environment can cause  
59 transmission of infectious diseases. Microbial source tracking (MST) can identify sources of  
60 fecal contamination so that contamination can be remediated and human exposures can be  
61 reduced. MST requires the use of fecal host-associated MST markers. Here we design and test  
62 novel MST markers from genomes of tomato brown rugose fruit virus (ToBRFV). The markers

63 are sensitive and specific to human stool, and highly abundant in human stool and wastewater  
64 samples.

65

66 **Keywords:** Microbial source tracking, tomato brown rugose fruit virus, human-associated  
67 marker, MST, virus, ToBRFV, wastewater

## 68 **Introduction**

69           Across the world, water quality is assessed for human fecal contamination using  
70 microbial indicators, including total coliforms like *Escherichia coli* and enterococci (1–3). Using  
71 these organisms to assess water quality is advantageous because they are abundant in human  
72 stool, which enables detection of even trace contamination of waterbodies. Additionally, their  
73 presence may indicate the potential contamination of waterbodies by other sparser human  
74 pathogens that may be harder to detect. However, there are limitations to their utility. These  
75 microbial indicators of human fecal contamination are also found in non-human stool (4).  
76 Additionally, they can be present and even grow in the environment, including in decaying plant  
77 material (1, 5), and in soils and sands (6, 7). Therefore, there is a need to identify new microbial  
78 indicator targets that can be used to specifically assess the presence of human fecal  
79 contamination.

80           The process of detecting microbes and identifying sources of microbial contamination in  
81 the environment is known as Microbial source tracking (MST). MST targets have also been used  
82 in COVID-19 wastewater-based epidemiology applications as “fecal strength” and endogenous  
83 extraction controls (8). Over the last decade, sensitive and specific molecular MST markers  
84 have been developed for various animal stools including human (9), cow (10), and birds (11).  
85 Most of these MST markers target conserved regions of bacterial genomes (9), with the  
86 exception of two that target the cross-assembly phage (crAssphage) (12) and pepper mild  
87 mottle virus (PMMoV) (13). crAssphage, a phage of *Bacteroidetes*, is a DNA virus that is highly  
88 abundant in the human gut (14). PMMoV is a plant RNA virus found at high concentrations in  
89 the human gut given its presence in popular spices, hot sauces, and other food products (15).  
90 The performance of MST targets is evaluated in terms of sensitivity and specificity for the host  
91 stool. For instance, a sensitive target for human stool is present at high concentrations in nearly

92 all human fecal samples so that dilute human stool can be detected in the environment.  
93 Meanwhile, a specific target is absent in nearly all non-human fecal samples. A previous study  
94 defined an MST assay as being sensitive and specific if the true positive and true negative rates  
95 were greater than 80% (9).

96 In this study, we present a new human-associated, RNA-based, viral MST target that is  
97 highly abundant in human stool and wastewater, tomato brown rugose fruit virus (ToBRFV).  
98 ToBRFV was first identified in Israel in 2014 and has since been detected across the world. To  
99 date, ToBRFV has been found across four continents, in at least 35 countries; this is likely an  
100 underestimate (16). We assembled eight nearly complete genomes of ToBRFV from wastewater  
101 and stool samples from the San Francisco Bay Area (Bay Area) in California in the United  
102 States of America (U.S.A), representing some of the first complete genomes from stool and  
103 wastewater in the area. Using these complete genomes, and other publicly available genomes,  
104 we developed two novel hydrolysis probe-based RT-PCR assays based on conserved regions  
105 of its RNA genome, and tested their sensitivity and specificity using stool and wastewater  
106 samples. Finally, we apply this assay for MST in stormwater samples collected from an urban  
107 environment. With the finding that ToBRFV is a reliable RNA-virus based MST marker, this  
108 study makes a valuable contribution to detecting human fecal contamination of the environment,  
109 and to waste-water based epidemiology.

110

111 **Materials and methods**

112 **Assembly and analysis of ToBRFV genomes, and design of hydrolysis-probe RT-PCR**  
113 **assays.**

114 In order to design ToBRFV-specific primers and probes for hydrolysis-probe RT-PCR  
115 assays, all ToBRFV genomes available in February 2021 were obtained. These were  
116 supplemented with new genomes assembled from stool samples processed and sequenced in  
117 this study (Table 1).

118 In February 2021, all near complete genomes (n = 70) of ToBRFV were downloaded  
119 from NCBI Genbank. In the same month, raw reads from the only publicly available wastewater  
120 metatranscriptomics dataset (obtained from wastewater in the Bay Area, collected between May  
121 and July 2020; Bioproject accession PRJNA661613) were also downloaded. Using these reads,  
122 five ToBRFV genomes were assembled as outlined below.

123 In addition to using existing sequencing data and genomes, RNA from three human stool  
124 samples obtained longitudinally from one individual were also sequenced; the first two samples  
125 were collected 10 days apart, and the third was collected 93 days after the second sample. The  
126 samples were obtained from an individual with laboratory-confirmed COVID-19 and were  
127 collected under an Institutional Review Board-approved protocol (Stanford IRB protocol  
128 #55619). Total RNA was extracted from these samples, rRNA was depleted, libraries were  
129 prepared and sequenced using NextSeq 550 as outlined in Note S1.

130 The following bioinformatic methods were used to assemble genomes from both the  
131 existing (from wastewater) and newly obtained (from stool) metatranscriptomics reads. Reads  
132 were trimmed with Trim Galore (version 0.4.0) using Cutadapt (version 1.8.1) (17) set to flags -q  
133 30 and -illumina. SPAdes (version 3.14.1) set to -meta was used to assemble genomes *de*  
134 *novo* (17, 18). Contigs belonging to ToBRFV were classified using One Codex (19). Genes were

135 annotated using Prodigal (version 2.6.3) set to -meta (20). If all genes were predicted on the  
136 negative strand of the contig, the entire contig was reverse complemented. The completeness of  
137 potential ToBRFV genomes was assessed using CheckV (version 1.0.1) (21) and genomes that  
138 were >90.0% complete were selected for subsequent analyses.

139 To assess strain diversity of ToBRFV in the longitudinal stool samples, RNA sequencing  
140 reads from stool samples were aligned to the ToBRFV reference genome (NCBI accession  
141 number NC\_028478) using Bowtie (version 2.4.2) (22). The resulting bam files were used as  
142 input to inStrain (version 1.0.0) (23) to calculate popANI between genomes.

143 To assess abundance of ToBRFV relative to other viruses in the RNA-Seq data, reads  
144 were classified against the Viral Kraken2 database ([https://benlangmead.github.io/aws-](https://benlangmead.github.io/aws-indexes/k2)  
145 [indexes/k2](https://benlangmead.github.io/aws-indexes/k2)) (24) using default parameters. Counts from the classification were used to calculate  
146 relative abundance of viral reads.

147 A multiple sequence alignment of all near complete genomes of ToBRFV including  
148 genomes downloaded from NCBI Genbank in February 2021 (70 genomes) and those we  
149 assembled from wastewater and stool (8 genomes) was performed using Geneious Alignment  
150 (Geneious Prime version 2021.0.3) (25) with default settings, global alignment with free end  
151 gaps, cost similarity matrix set to 65.0%. SNPs were called from the multiple sequence  
152 alignment using SNP-Sites (version 2.5.1) (26). A phylogenetic tree was built using Geneious  
153 Tree Builder (version 2021.0.3) with default settings, Tamura-Nei genetic distance model with  
154 the neighbor-joining method. Primers and probes were designed to be specific for ToBRFV  
155 using Geneious Primer (version 3 2.3.7) (27) based on the 78 genomes we had access to in  
156 February 2021 with near default settings, requiring product size to be between 95-125 base  
157 pairs in length and primers to be based on consensus with 100.0% identity across all ToBRFV



158 genomes. Primers and probe sequences were screened for specificity, *in silico*, using NCBI  
159 Blast.

160

161 **New genomes available in November 2022.**

162 New ToBRFV genomes became available on public databases within the last six  
163 months. An additional 113 (total 183) from NCBI (28) and 250 assembled ToBRFV genomes  
164 from a study of wastewater from Southern California (29) were downloaded (Table 1).

165 As Geneious alignment and tree building is computationally intensive, a phylogenetic  
166 tree of all 441 near complete genomes of ToBRFV was built using ViPTree (30), visualized and  
167 color coded by region using Iroki (31). In addition, the applicability of the primers and probes  
168 designed in this study was tested *in silico* using NCBI Blast.

Sample type	Sample source	Number of ToBRFV genomes available in		Reference	
		Feb 2021	Nov 2022	Sequence data	Assembled genomes
Stool	Bay Area, CA, U.S.A.	3	0	(current study)	(current study)
Tomatoes	Global	70	183	N/A	<a href="https://www.ncbi.nlm.nih.gov/nucleotide">https://www.ncbi.nlm.nih.gov/nucleotide</a>
Wastewater	Southern CA, U.S.A.	0	250	N/A	(29)
Wastewater	Bay Area, CA, U.S.A.	5	0	NCBI Bioproject (PRJNA661613)	(current study)

169 **Table 1: Source of genomes analyzed**

170

171 **Processing of animal stool samples for RNA quantification.**

172 One stool sample was either collected from a) a single animal (cat, dog, horse, pig,  
173 rabbit) raised as a pet, b) a group of cohabiting animals of a single kind (chicken, cow, goat,  
174 mouse, sheep) from the Deer Hollow Farms (California, U.S.A.), c) a group of cohoused animals  
175 at the Deer Hollow Farms in the case of the ducks and geese, or d) from the wild (bear, deer).  
176 Samples were collected wearing gloves, using a spatula and in a sterile clinical stool collection  
177 container. Samples were transported at room temperature, aliquoted into cryovials and stored at  
178  $-80^{\circ}\text{C}$  within 12 hours from collection. Samples were further processed within a month of  
179 storage.

180 A single, defined solid volume of sample of each animal stool was acquired using the  
181 Integra Miltex Biopsy Punches with Plunger System (Thermo Fisher Scientific; Catalog # 12-  
182 460-410) to independent microcentrifuge tubes. 500  $\mu\text{l}$  of RNALater (Ambion; Catalog #  
183 AM7023M) was added and samples were processed using a previously validated methodology  
184 (32) as follows. A stock Bovine Coronavirus (BCoV) vaccine was prepared by adding 3 mL of 1X  
185 Phosphate Buffered Saline (PBS; Fisher Scientific; Catalog # BP399-500) to one vial of  
186 lyophilized Zoetis Calf-Guard Bovine Rotavirus-Coronavirus Vaccine (Catalog # VLN 190/PCN  
187 1931.20) to create an undiluted reagent as per manufacturer's instructions. 10  $\mu\text{l}$  of this  
188 attenuated BCoV vaccine was added to every sample as an external control and vortexed for 15  
189 minutes. BCoV is an RNA virus that was previously found to be a reliable positive control for  
190 RNA extraction from stool (32) and helps identify instances of PCR inhibition. Samples were  
191 processed immediately after addition of the BCoV control.

192 **Collection and processing of human stool samples used for RNA quantification.**

193 Human stool samples were previously collected and biobanked in RNALater solution as  
194 part of Stanford Institutional Review Board-approved protocols #8903 (“Blood and Bone Marrow  
195 Grafting for Leukemia and Lymphoma”), #11062 (“Genome, Proteome and Tissue Microarray  
196 Studies in Childhood malignant and Non-Malignant Hematologic Disorders”), and #48548  
197 (“Hematopoietic Recovery During Induction Chemotherapy in Pediatric Leukemia”). From these  
198 biobanks, 194 and 28 samples collected over the span of a year from November 2019 to  
199 October 2020 from 125 adult and 4 pediatric participants respectively were used in this study.  
200 These samples had been stored for between 1 - 12 months depending on date of collection. All  
201 samples were spiked with 10 µl of attenuated Bovine Coronavirus (BCoV) vaccine as control  
202 and processed similar to the animal stool samples.

203 **RNA extraction from all stool samples used for RNA quantification.**

204 RNA was extracted from these stool samples using the QIAamp Viral RNA Mini Kit  
205 (Qiagen; Catalog # 52906) as previously optimized (32). Briefly, the prepared stool samples  
206 were spun down at 10,000x g for 2 minutes to acquire 140 µL of clarified supernatant. RNA was  
207 extracted from this supernatant using the QIAamp Viral RNA Mini Kit (Qiagen; Catalog # 52906)  
208 as per the manufacturer’s instructions. Finally RNA was eluted in 100 µL of the elution buffer  
209 and stored in a 96-well plate at -80°C for up to 12 months. In previous work on BCoV and  
210 SARS-CoV-2 RNA (32) we found that RNA extracted using this method did not result in RT-  
211 PCR inhibitors. Therefore, we assume that samples extracted here also do not have any RT-  
212 PCR inhibitors.

213 **Augmenting analysis of stool with metatranscriptomic data from healthy individuals.**

214 As described below, we assessed the prevalence and abundance of MST markers in stool  
215 acquired from participants with hematologic disorders. This presented a caveat to the

216 generalizability of our work. Therefore, we acquired metatranscriptomics data from stool  
217 samples from 10 healthy participants presented in a previous study (33). Though many human  
218 stool metatranscriptomic datasets exist, this was the most recent dataset we had access to.

#### 219 **Collection and processing of wastewater samples used for RNA quantification.**

220 Settled solids were obtained from 15 wastewater treatment plants across the US (Table  
221 S2). Solids were collected from the primary clarifier, or settled from a 24 hour composited  
222 influent sample using Imhof cones. Samples were collected in sterile containers and transported  
223 to the lab. Samples from the Bay Area were processed immediately, while other samples were  
224 stored at -80°C until analysis (between 5 and 20 months).

225 Solids were dewatered using centrifugation and then an aliquot of the dewatered solids  
226 was set aside for dry weight analysis. Solids were then suspended in a buffer (approximately 75  
227 mg/ml), homogenized and centrifuged. This suspension of solids in buffer was found to alleviate  
228 inhibition of RT-PCR (35). An aliquot of the supernatant was processed for total nucleic-acid  
229 extraction using Chemagic 360 (Perkin Elmer). Nucleic acid preparations from wastewater  
230 samples are known to contain PCR inhibitors that interfere with their accurate quantification  
231 using PCR-based methods. Therefore, inhibitors were removed using the OneStep PCR  
232 Inhibitor Removal Kit (Zymo research; Catalog # D6035). These methods have been published  
233 in detail (36) and step-by-step protocols are available on protocols.io (37, 38).

#### 234 **Source of RNA extracted from stormwater samples used here for RNA quantification.**

235 RNA extracted from stormwater samples was derived from a previous study from our  
236 group (39). Briefly, nine stormwater samples from the Bay Area - one from Guadalupe River,  
237 Pilarcitos Creek, San Francisquito creek and San Pedro Creek, two from Stevens Creek, and  
238 three from Lobos Creek - collected between October 2018 and March 2019 were used to extract

239 RNA (Table S3). Specifically, viruses were concentrated from 1 - 5.5 L stormwater samples  
240 using electronegative filtration using MgCl<sub>2</sub>. The filtration membranes were preserved in 250 µl  
241 of RNALater (Qiagen; Catalog #76104) for 5 minutes prior to storage at -80°C. Nucleic acids  
242 were extracted from the stored filtration membrane using the Qiagen DNA/RNA AllPrep  
243 PowerViral Kit using the protocol including β-mercaptoethanol and bead-beating, and stored in  
244 microcentrifuge tubes at -80°C. Previous work suggested that RT-PCR inhibitors from the  
245 samples were not co-extracted in this RNA extraction process (39). These extracts were used in  
246 the current study after 30 months of storage.

#### 247 **Quantification of viral RNA sequences by ddRT-PCR.**

248 The CP gene encoding the coat protein from PMMoV, Mo gene encoding the movement  
249 protein and RdRP gene encoding the RNA-dependent RNA polymerase from ToBRFV, and M  
250 gene encoding the membrane protein from BCoV were quantified using droplet digital reverse  
251 transcription-polymerase chain reaction (ddRT-PCR). Human participants in this study were  
252 enrolled and hospitalized during the first year of the COVID-19 pandemic. We tested their stool  
253 for genes encoding the envelope (E) and a nucleocapsid (N2) protein from the SARS-CoV-2  
254 genome as previously described (32), in order to assess occurrence of COVID-19 during  
255 hospitalization at Stanford Hospital. However, we did not find any presence of COVID-19 RNA  
256 in these samples. Sequences of the newly designed primers and probes targeting ToBRFV Mo  
257 and RdRp genes are listed in Table 2. Previously published primers and probes targeting  
258 BCoV, PMMoV and SARS-CoV-2 RNAs are listed in Table S4. We chose ddRT-PCR instead of  
259 RT-qPCR for nucleic acid detection and quantification because of its superior sensitivity and  
260 resistance to PCR inhibitors (32, 40).

Primer	Description	Sequence (5' to 3')	Amplicon length (bps)
--------	-------------	---------------------	-----------------------

ToBRFV_Mo_F	ToBRFV Mo gene; forward primer	TCA GTG TCT GTT TGG TCG ATA A	105
ToBRFV_Mo_R	ToBRFV Mo gene; reverse primer	GGA ACG ACT TTG AAC TGA AAC C	
ToBRFV_RdRP_F	ToBRFV RdRP gene; forward primer	AGC CAC AAG AGA TAA TGT TCG TA	103
ToBRFV_RdRP_R	ToBRFV RdRP gene; reverse primer	ACA TCA GAC CTT CGT CGA TAA AT	
<b>Probe</b>	<b>Description</b>	<b>Sequence (5' to 3')</b>	<b>Modifications</b>
ToBRFV_Mo_P	ToBRFV Mo gene; probe	AGA GCG GAC GAG GCA ACT CTT G	FAM/ZEN/IBH Q
ToBRFV_RdRP_P	ToBRFV RdRP gene; probe	ACG GTA AAG GAA CAC GCT GTC AGT	FAM/ZEN/IBH Q

261 **Table 2: Sequences of primers and probes designed in this study to quantify ToBRFV M**  
 262 **ands RdRP genes.**

263 The Droplet Digital PCR Applications Guide on QX200 machines (BioRad) (41) and  
 264 digital Minimum Information for Publication of Quantitative Real-Time PCR Experiments  
 265 (dMIQE) guidelines (42) inform this methodology. The experimental checklist recommended by  
 266 dMIQE is available at the Stanford Digital Repository (<https://purl.stanford.edu/nf771cs9443>). A  
 267 Biomek FX liquid handler (Beckman Coulter) was used to prepare the ddRT-PCR reaction by  
 268 adding 5.5 µL of eluted RNA to 5.5 µL Supermix, 2.2 µL reverse transcriptase, 1.1 µL of 300 nM  
 269 Dithiothreitol (DTT), 1.1 µL of each of the 20x Custom ddPCR Assay Primer/Probe Mix (BioRad,  
 270 Catalog # 10031277) and 5.5 µL of nuclease-free water (Ambion, Catalog # AM9937, Lot  
 271 2009117). The Supermix, reverse transcriptase and DTT were from the One-Step ddRT-PCR  
 272 Advanced Kit for Probes (BioRad, Catalog # 1864021). A QX200 AutoDG Droplet Digital PCR

273 System (BioRad) was used to partition the samples into droplets of roughly 1 nl using the default  
274 settings and the template was amplified using a BioRad T100 thermocycler with the following  
275 thermocycling program: 50°C for 60 min, 95°C for 10 min, 40 cycles of 94°C for 30 sec and  
276 55°C for 1 min, followed by 1 cycle of 98°C for 10 min and 4°C for 30 min with ramp speed of  
277 1.6°C per second at each step (43).

278 A multistep approach was adopted to calculate the raw RNA concentrations, as  
279 previously described (32). Every plate of ddRT-PCR assays included appropriate positive and  
280 negative controls including synthetic target genes (PMMoV CP gene PMMoV, ToBRFV Mo and  
281 RdRP genes) cloned in the pIDT vector, RNA extracted from reconstituted attenuated BCoV  
282 vaccine, water and RNA later. The signal threshold corresponding every plate was manually set  
283 between the mean positive and negative amplitudes of these controls such that the number of  
284 detected copies in the negative controls was minimal and those from the relevant positive  
285 controls most closely matched the expected RNA concentration. Next, the difference between  
286 the mean negative amplitude and the threshold amplitude in the negative control reactions was  
287 calculated and added to the mean negative amplitude for every sample on that plate. Applying  
288 this threshold yielded the raw RNA concentrations.

289 In order to derive the limit of blank (LoB) and limit of detection (LoD) of our assays to  
290 further process the raw RNA concentrations we adopted the following steps: a) LoB indicates  
291 the highest background RNA concentration registered from control samples that are confidently  
292 negative for the relevant gene targets. In order to determine the LoB, water, RNA later and  
293 synthetic genes discordant with the target gene (e.g. ToBRFV Mo gene is a negative control on  
294 an assay of ToBRFV RdRP gene) were assayed in duplicate. The highest RNA concentration  
295 measured in these LoB samples for each of the primer/probe sets was set as the relevant LoB.  
296 All samples that we detected an RNA concentration equal to or less than the LoB were zeroed.

297 b) LoD is defined as the lowest concentration of RNA that can be reliably detected. To  
298 determine the LoD, duplicate serial dilution series of the synthetic target genes at the following  
299 concentrations - 1, 2, 5, 10, 100, 1000 copies/  $\mu$ L of template - were assayed for the  
300 corresponding target gene (Fig. S1). The synthetic target genes were acquired from Integrated  
301 DNA Technologies and cloned in their standard backbone, pIDTSmart. These plasmids were  
302 transformed into *E. coli*, isolated using the QIAprep Spin miniprep kit (Qiagen; Catalog # 27104)  
303 and quantified using Qubit. LoD for a primer/probe set was defined as the least concentration of  
304 the standard at which both replicates had a detectable RNA concentration. All viral RNA  
305 concentrations below the LoD were zeroed.

306 Finally, after applying these data processing and analysis steps, the samples were  
307 assigned a final viral RNA concentration in copies/  $\mu$ L of template. Eluate refers to the 100  $\mu$ L of  
308 sample acquired from the RNA extraction. Viral RNA concentrations from animal and human  
309 stool samples are expressed as copies/  $\mu$ L of template, from wastewater samples as copies/ g  
310 of wastewater, and from stormwater samples as copies/ liter of stormwater.

311 In the case of all non-human animal stool, wastewater and stormwater samples, RNA  
312 was quantified using singleplex reactions. For the human stool samples, which were limited in  
313 quantity, the detection of the BCoV M gene and PMMoV CP gene were multiplexed with the  
314 detection of the SARS-CoV-2 E and N2 genes using orthogonal fluorescent probes. After  
315 extensive optimization (outlined in the Note S1), we paired the detection of the E gene (SARS-  
316 CoV-2) with the CP gene (PMMoV), and N2 gene (SARS-CoV-2) with the M gene (BCoV) in two  
317 independent reactions using the carboxyfluorescein (FAM) and hexachlorofluorescein (HEX)  
318 fluors, respectively.



319 **Data analysis and generation of plots.**

320 Data was analyzed using RStudio (ver 1.2.5042), using packages cowplot (ver 1.1.1),  
321 dplyr (ver 1.0.8), eulerr (ver 6.1.1), ggplot2 (ver 3.3.6), and UpSetR (ver 1.4.0).

322 **Data availability.**

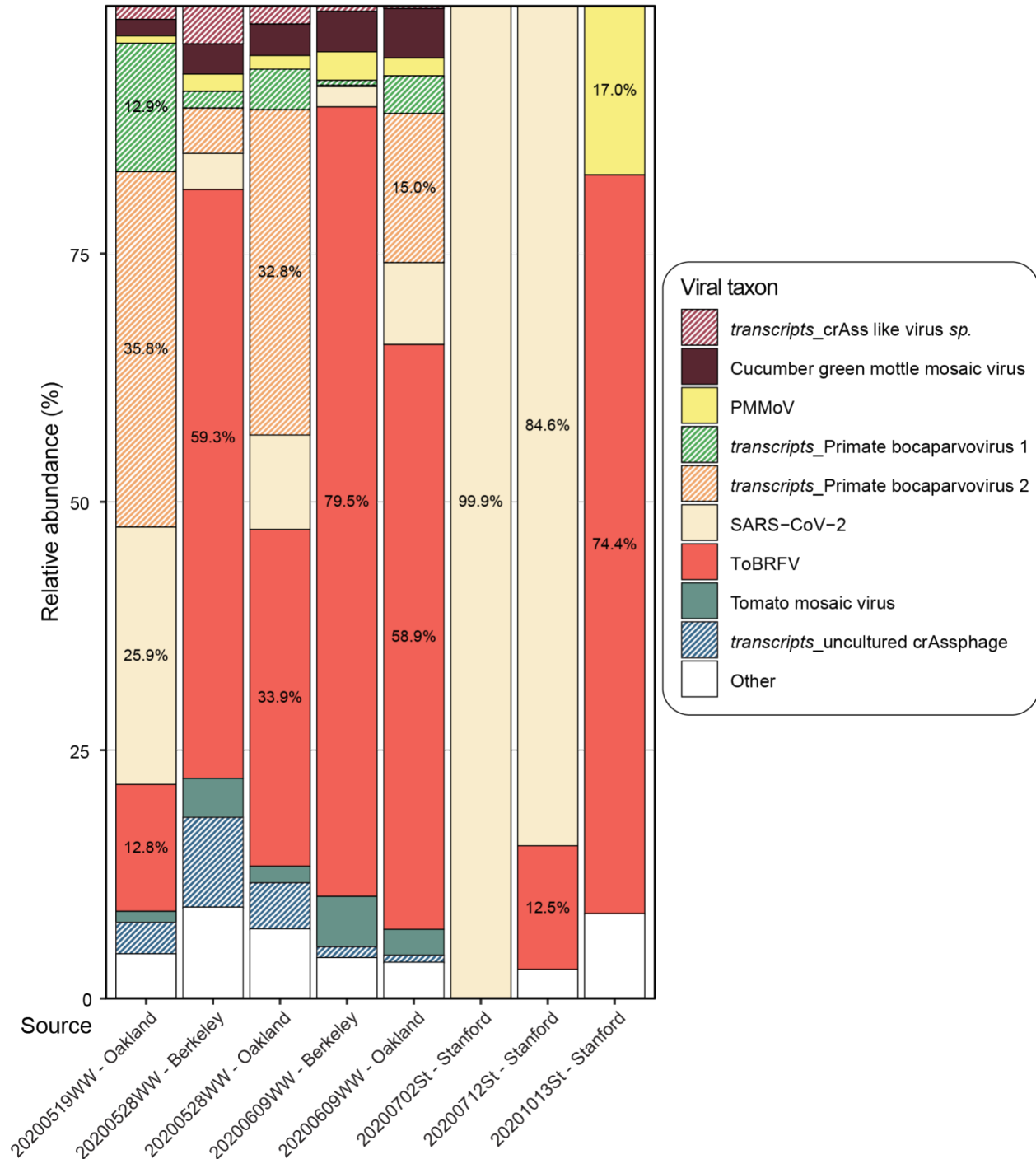
323 Newly generated genomes and raw sequencing reads from stool samples are available  
324 on NCBI's Sequence Read Archive (SRA) database at accession #PRJNA917455. All other  
325 relevant data are included in this manuscript and available through the Stanford Digital  
326 Repository (<https://purl.stanford.edu/nf771cs9443>).

327 **Results and discussion**

328 **ToBRFV is widely prevalent and abundant in sequence data from stool and wastewater**  
329 **samples.**

330 Tracking the presence of human feces in the environment, and the identification of  
331 internal controls for the processing of stool and wastewater samples requires marker genes that  
332 are: a) prevalent - consistently present across samples, and b) abundant - available in high  
333 enough concentration for reliable detection. crAssphage (12) presents one such DNA-based  
334 marker and PMMoV is another RNA-based marker (13). We sought to test whether RNA-based  
335 markers from ToBRFV also meet these criteria.

336 We isolated and sequenced RNA from three longitudinal stool samples from one human  
337 participant, who had tested positive for SARS-CoV-2. In parallel, we acquired publicly available  
338 transcriptomics data from five wastewater samples that had been collected and sequenced from  
339 the Bay Area (44). Using these sequence data from eight samples, we identified all represented  
340 RNA viruses and their relative abundances (Fig. 1). The tomato brown rugose fruit virus  
341 (ToBRFV) was the most widely prevalent RNA virus, present in all five wastewater samples and  
342 three stool samples. It was detected at very low relative abundance (0.077% of viral reads) in  
343 one of the stool samples during the time of active SARS-CoV-2 infection, in which 99.9% of  
344 reads belonged to SARS-CoV-2. In the other seven samples where it is detected, it is the only  
345 viral RNA that is consistently over 10.0% in relative abundance of viral reads, often making up  
346 over 50.0% of the reads. Notably, the relative abundance of ToBRFV was consistently greater  
347 than PMMoV, which is a well established microbial source tracking marker and known to be  
348 highly abundant in wastewater (8). This is consistent with reports from studies carried out prior  
349 to (44), and in parallel with (29, 45) ours, that also show that ToBRFV is a highly prevalent virus  
350 in wastewater.



351 **Fig. 1: Relative abundance of viral RNA from sequencing wastewater and stool samples.**

352 The x-axis represents the source of the eight sequencing datasets analyzed here. Five  
 353 wastewater samples are marked by the date of collection in 'YYYYMMDD' format followed by  
 354 "WW" and the location of collection. The three stool samples are marked by the date of  
 355 collection in 'YYYYMMDD' format followed by "St" and the location of collection. The y-axis  
 356 indicates the relative abundance of each taxon. The color scheme represents specific taxon as  
 357 light brown for crAss like virus *sp.*, dark brown for cucumber green mottle mosaic virus, yellow

358 for PMMoV, green for primate bocaparvovirus 1, orange for primate bocaparvovirus 2, cream for  
359 SARS-CoV-2, salmon for ToBRFV, light blue for tomato mosaic virus, dark blue for uncultured  
360 crAssphage, and white for other minor taxa. Patterned bars highlight sequence reads from  
361 transcripts from taxa that have DNA genomes. Taxa with >10.0% relative abundance also list  
362 the percentage of abundance in the histogram.

363

### 364 **Novel ToBRFV genomes and sequence analysis reveal suitable RNA-borne marker genes.**

365 Having identified that ToBRFV is a prevalent and abundant RNA virus in sequence data,  
366 we next set out to identify genomic regions suitable as a target for primer/probe for its reliable  
367 molecular detection.

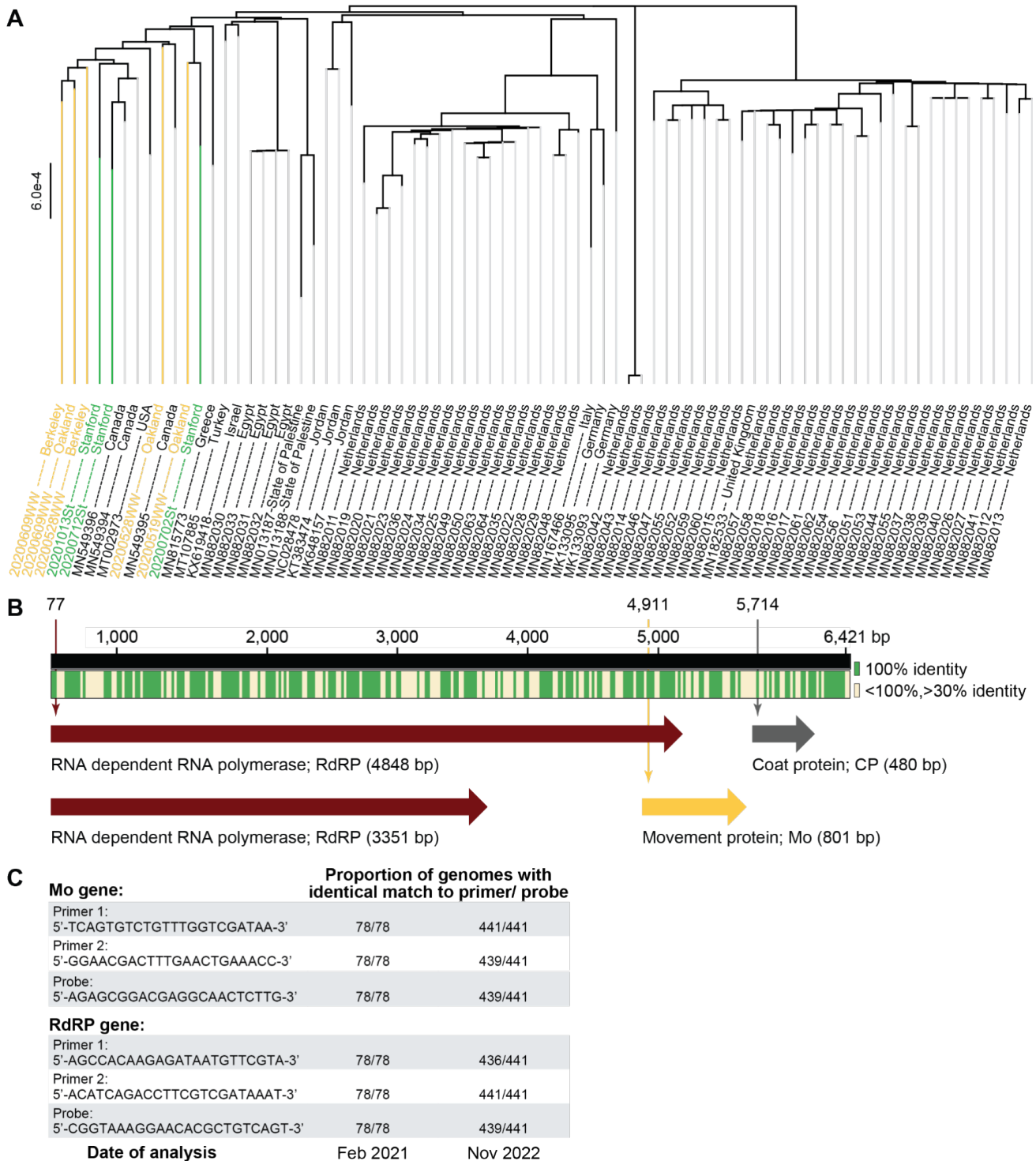
368 In February 2021, at the start of this study, only 70 near complete ToBRFV genomes  
369 were known. 50 of these were from the Netherlands. None had been sequenced from human  
370 stool or wastewater samples, and only one sequence was derived from the U.S.A. In order to  
371 ensure that the assay we developed was universal, we first decided to augment the number of  
372 ToBRFV genomes and the diversity of its sources. Therefore, we assembled near complete  
373 genomes of ToBRFV using sequence data generated in this study from stool samples and using  
374 existing data from wastewater samples (44), both collected in the Bay Area. The eight newly  
375 assembled genomes had a mean completeness of 98.8% (range 93.6% - 100.0%; median  
376 99.4%) (Table S5). The longitudinally acquired stool samples yielded ToBRFV genomes with  
377 SNPs in 27 positions, suggesting possible strain variation over time. Looking more broadly,  
378 across all 78 near complete ToBRFV genomes, we identified 2,808 positions containing SNPs  
379 (across an average contig length of 6,366 bp), and the 12 North American strains form their  
380 own, distinct cluster (Fig. 2B).

381 Multiple sequence analysis across all 78 ToBRFV genomes highlights regions that are  
382 100.0% conserved (Fig. 2C). Among these, gene annotation reveals a) two variants of the RNA  
383 dependent RNA polymerase (RdRP) encoding gene at 2,700 bp on the chromosome, that differ

384 by whether an internal stop codon is read through (size 3,351 bp or 4,848 bp), b) the movement  
385 protein (Mo) encoding gene (size 480 bp) at 5,166 bp, and c) the coat protein (CP) encoding  
386 gene (size 801 bp) at 5,166 bp (Fig. 2C). Among these we designed primer/probe sets targeting  
387 the 5' end of the RdRP gene, and the Mo gene. We were unable to identify a suitable primer set  
388 for the CP gene for ddRT-PCR. Notably, the primer/probe sets designed here (Table 2) were  
389 conserved across all 78 genomes (Fig. 2D).

390           Within the last six months, the number of near complete ToBRFV genomes have  
391 increased to 441 (Table 1), with additional genomes from Belgium, France, Mexico, Switzerland,  
392 and the U.S.A. (summarized in Fig. 2b). Therefore, we repeated the phylogenetic analysis of the  
393 novel genomes generated in the current study in the context of all 441 currently known genomes  
394 (Fig. S3). Again, we find that the genomes derived from North American cluster distinctly.  
395 Finally, we analyzed whether the primer/probe sets proposed here continue to be universal and  
396 found that the oligonucleotides targeting Mo are a perfect sequence match in 439/441 genomes,  
397 while those targeting RdRP are a perfect match in 436/441 genomes (Fig. 2D).

398  
399  
400  
401  
402  
403  
404  
405  
406



407 **Fig. 2: Analysis of newly assembled ToBRFV genomes and generation of primer/probe**  
 408 **sets for ddRT-PCR.**  
 409 (A) Phylogenetic tree of 78 near complete genomes of ToBRFV, including eight genomes  
 410 generated in the current study from wastewater and stool. All genomes are listed by their NCBI  
 411 accession number and source location. 70 pre-existing genomes are listed in black font, five  
 412 genomes derived from wastewater samples are listed in yellow, and three from stool samples in  
 413 green. (B) Summary of multiple sequence alignment and gene annotation across the 78

414 ToBRFV genomes. Green indicates regions that are 100.0% conserved across all genomes,  
415 while cream marks those that are greater than 30.0% but less than 100.0% conserved. Two  
416 variants of the RNA dependent RNA polymerase (RdRP) encoding gene are found at 77 bp and  
417 are of either 4,848 bp or 3,351 bp in size. The movement protein (Mo) encoding gene is found at  
418 4,911 bp and is 801 bp in size. The coat protein (CP) encoding gene is found at 5,714 bp and is  
419 480 bp in size. Genomic locations are based on genome ID NC\_028478. Sequences of  
420 primer/probe sets generated in the current study, in Feb 2021, aimed at targeting the Mo and  
421 RdRP genes across all known genomes. Since the number of known genomes grew from Feb  
422 2021 to Nov 2022, the final column indicates the proportion of the 441 current genomes bearing  
423 sequences identical to the designed primer/probe sets.

424

#### 425 **ToBRFV targeting primer/probe sets have low LoB and LoD.**

426 Having newly designed primer/probe sets targeting the Mo and RdRP genes in ToBRFV,  
427 we aimed to validate these oligonucleotides and establish the limits of their reliable utility.

428 To this end, we acquired synthetic DNA constructs featuring regions of the ToBRFV Mo  
429 and RdRP genes targeted by hydrolysis-probe RT-PCR assays from Integrated RNA  
430 Technologies (IDT) cloned into the pIDT plasmid. We also acquired a similar plasmid containing  
431 the PMMoV CP gene. Using ddRT-PCR, we assayed a dilution series of these synthetic plasmid  
432 constructs at 1, 2, 5, 10, 100, 1000 copies /  $\mu\text{L}$  of template in triplicate and found that all the  
433 primer/probe sets detected the target gene at all concentrations (Fig. S4). Next, we focused our  
434 attention on the negative controls included in the assays to identify the limit of detection for each  
435 primer/probe set. The negative controls included two no template controls, water and RNALater,  
436 and two mismatched controls that were the synthetic pIDT plasmids bearing targets orthogonal  
437 to the primer/probe sets. Therefore, theoretically, all the negative controls would have no  
438 detectable gene target. For each primer/probe set, among the negative controls, we identified  
439 the highest concentration of target detected and set this value as the limit of blank (LoB). This  
440 means any concentration below  $-0.552 \log_{10}$  copies/  $\mu\text{L}$  of template for primer/probe set  
441 targeting PMMoV CP gene,  $-0.590 \log_{10}$  copies/  $\mu\text{L}$  of template for ToBRFV Mo gene, and  $0.407$   
442  $\log_{10}$  copies/  $\mu\text{L}$  of template for ToBRFV RdRP gene is not reliable (Fig. S4). After converting all

443 concentrations of gene targets below the LoB to zero, we focused our attention on the triplicate  
444 dilution series to identify the lowest concentration of template at which all three reactions had a  
445 detectable target concentration (Fig. S4). We set this concentration as the limit of detection  
446 (LoD), i.e. the least concentration at which a gene target can be reliably detected. The LoD for  
447 the primer/probe set targeting PMMoV CP gene was 1 copies/  $\mu$ L of template, for ToBRFV Mo  
448 gene was 5 copies/  $\mu$ L of template, and for ToBRFV RdRP gene was 5 copies/  $\mu$ L of template.  
449 All gene target concentrations below the LoD were set to zero.

450

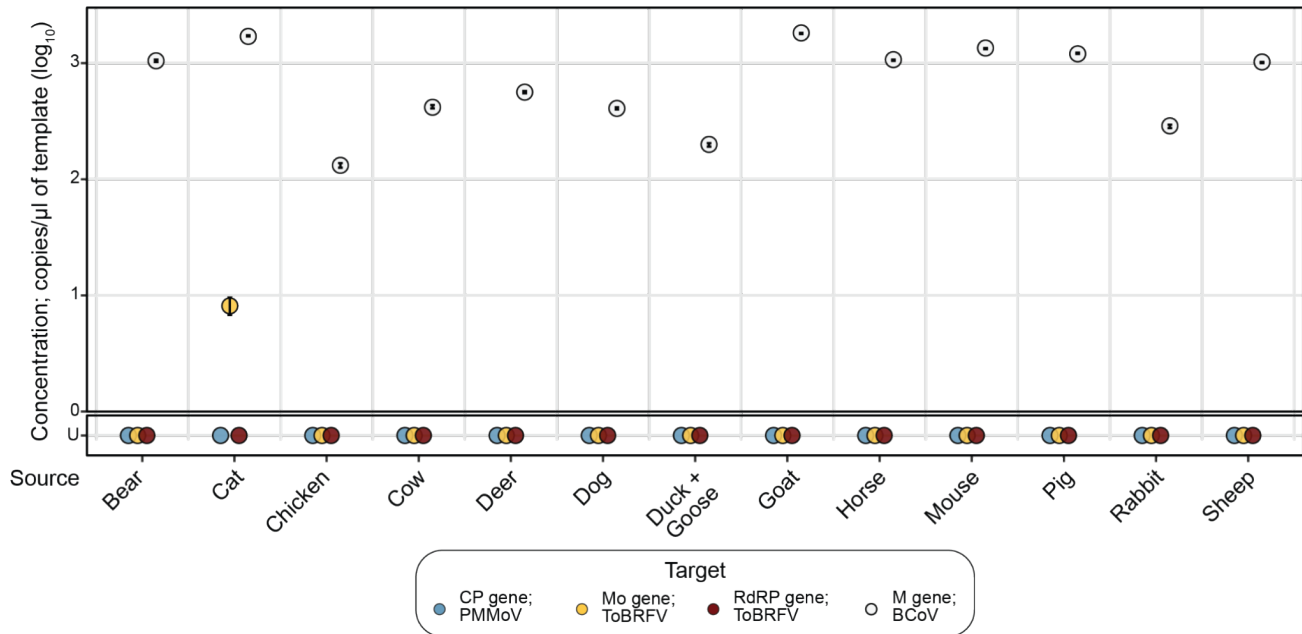
#### 451 **ToBRFV is not detected in stool from non-human animals.**

452 MST targets should be specific, meaning they are mostly absent in stool from other  
453 common animals. Therefore, having established that our primer/probe sets are functional, we  
454 tested them against stool collected from 14 different animals including wild bear and deer,  
455 chicken, cow, duck, goose, goat and sheep from a farm, horse and pig from a barn, a household  
456 cat, dog and rabbit, and laboratory mice. Notably, these animals are rather diverse and are fed a  
457 wide variety of foods. While RNA extracted from all of these animal samples had a detectable  
458 concentration of the M gene target from the spiked-in BCoV used as a control, none of them had  
459 RNA containing either the PMMoV CP gene or the ToBRFV RdRP gene. The ToBRFV Mo gene  
460 was detectable only in the sample derived from the domesticated cat, perhaps due to inclusion  
461 of tomatoes in its processed kibble or cross contamination of diet with its human cohabitant.  
462 Therefore, all three sets of primer/probe to detect RNA from PMMoV and ToBRFV are not  
463 detected in most animal feces, except for the ToBRFV Mo gene in a cat, indicating they are  
464 specific for human stool.

465

466





467 **Fig. 3: Concentrations of PMMoV and ToBRFV target genes in animal stool samples.**  
468 Dot plot marking the concentrations of the PMMoV CP (blue), ToBRFV Mo (yellow) and RdRP  
469 (red) genes. Concentrations of BCoV M gene, used as a control, are marked by white dots.  
470 Error bars marking the standard deviation are plotted along with the dots, and are mostly  
471 subsumed within the dot. The x-axis lists 13 samples from 14 different animals, where a single  
472 sample is derived from cohoused ducks and geese. The y-axis lists concentrations of the genes  
473 in log<sub>10</sub> copies/ μl of template; U stands for “Undetermined” and marks samples with no  
474 detectable gene

475

### 476 **Description of participants who provided human stool samples used for RNA** 477 **quantification.**

478 Analyzing sequence information from three stool samples collected from one human  
479 participant revealed ToBRFV to be abundantly present. To further test the sensitivity of the  
480 assays to human stool, we relied on a stool biobank including 194 stool samples from 125  
481 adults and 28 samples from four children, all of whom were undergoing hematopoietic cell  
482 transplantation (HCT), cell therapy (CAR-T) or induction chemotherapy for the treatment of  
483 underlying hematologic disorders.

484 Of the adult participants, 79 are male, 45 are female and 1 participant did not provide  
485 information on their sex. The median age of the adult participants is 60 years (range 19 - 82  
486 years) and pediatric participants is 6 years (range 3 - 16 years). Among the adult participants,  
487 61.6% self-identify as White. Age, race and ethnicity information on pediatric participants are  
488 withheld since these can be used to identify the participants. Timeline of stool collection is  
489 summarized in Fig. S1. Demographic information is summarized in Fig. S2 and Table S1.

490

#### 491 **ToBRFV is more prevalent in human stool samples compared to PMMoV.**

492 220 out of 222 RNA extracts derived from 129 participants had detectable BCoV RNA.  
493 This suggests that two of the RNA extractions failed; those samples are therefore excluded from  
494 further analysis, altering our study cohort size to 127 (123 adult; 4 pediatric). Among the  
495 remaining samples, 126/220 (57.3%) of stool samples had detectable levels of the PMMoV CP  
496 gene, while 143/220 (65.0%) had the ToBRFV Mo gene and 108/220 (49.1%) had the ToBRFV  
497 RdRP gene; ToBRFV Mo gene was the most prevalent target gene. This prevalence varied in  
498 the two patient cohorts (Fig. 4A); 127/192 (66.2%) stool samples from adult participants had  
499 detectable amounts of the ToBRFV Mo gene, more than in the case of PMMoV CP gene  
500 (103/192; 54.7%), but only 16/28 (57.1%) stool samples from pediatric patients had detectable  
501 amounts of the ToBRFV Mo gene, fewer than in the case of PMMoV CP gene (23/28; 82.1%).

502 In analyzing the prevalence of the three gene targets of interest in the stool samples, we  
503 detected all three gene targets in 70 (31.8%) of the samples, while we detected none of the  
504 three gene targets in 43 (19.6%) (Fig. 4B, Fig. S5). Notably, in 34 (15.5%) of the samples, we  
505 only detected the PMMoV CP gene, and in 13 (5.9%), we only detected the ToBRFV Mo gene.  
506 In all samples that we detected the ToBRFV RdRP gene we also detected the ToBRFV Mo  
507 gene. This analysis suggests that while the ToBRFV Mo gene is the most prevalent RNA-based

508 marker of human stool, combining this with the detection of the PMMoV CP gene will provide  
509 the most coverage of more than 80.0% of stool samples.

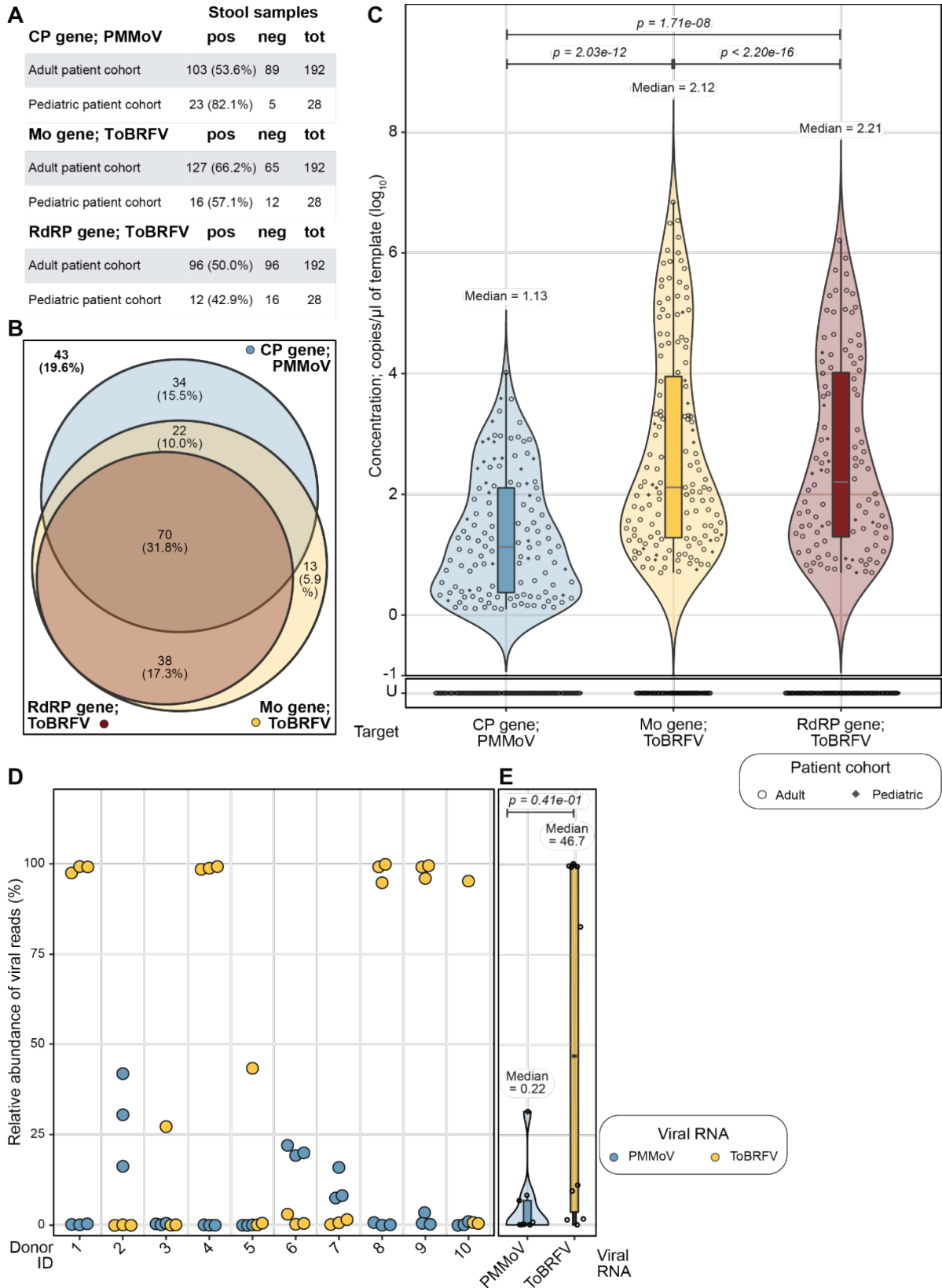
510 Next we analyzed the abundance of each of these gene targets in stool samples. The  
511 median detected concentration of the PMMoV CP gene is lower than the ToBRFV Mo gene  
512 ( $1.13 \log_{10}$  copies/  $\mu\text{l}$  vs.  $2.12 \log_{10}$  copies/  $\mu\text{l}$ , Wilcoxon signed rank test  $p = 2.03\text{e-}12$ ) and the  
513 ToBRFV RdRP gene ( $1.13 \log_{10}$  copies/  $\mu\text{l}$  vs.  $2.21 \log_{10}$  copies/  $\mu\text{l}$ ,  $p = 1.17\text{e-}8$ ; Fig. 4C). These  
514 stool samples were derived from participants undergoing different treatments for underlying  
515 hematologic disorders. Therefore, we investigated whether the nature of treatment was a  
516 confounding factor. Here, again, we find that the median abundances of both target genes from  
517 ToBRFV are higher than the PMMoV CP gene, even when the samples are separated by  
518 treatment cohort (Fig. S6A). Further, a paired comparison of target gene abundances validates  
519 the previous observation that all samples that tested positive for the ToBRFV RdRP gene also  
520 tested positive for the ToBRFV Mo gene (Fig. S6B).

521 While the concentration of the various gene targets has so far been reported in copies  
522 per  $\mu\text{l}$  of template, we recognize that studies also measure molecular targets in units per g dry  
523 weight of stool sample. Therefore, we chose five samples per cohort at random, dried two  
524 biopsy punches from each sample and found the mean percent dry weight of the samples from  
525 adults undergoing HCT treatment as 23.6% (range 18.2 - 33.9%), CAR-T treatment as 27.5%  
526 (range 22.2 - 31.6%), and those of pediatric patients undergoing induction chemotherapy as  
527 32.4% (range 23.8 - 40.3%). We used the average percent dry weight to convert gene target  
528 concentrations to copies per gram (g) dry weight of stool samples in Fig. S6C. In brief, the  
529 median concentrations of ToBRFV RdRP and Mo genes were  $6.45$  and  $6.32 \log_{10}$  copies/ g dry  
530 weight of stool, the PMMoV CP gene was  $5.36 \log_{10}$  copies/ g (Fig. S6C).

531           To determine if our findings are generalizable to applications beyond a cohort of patients,  
532 we looked at an alternate dataset recently generated in our group that sequenced RNA from  
533 stool collected and frozen from 10 healthy individuals in triplicate (33). In this dataset also, the  
534 relative abundance of ToBRFV was consistently greater than PMMoV, as reflected by their  
535 median relative abundance of 46.7 vs. 0.22% viral RNA reads (Fig. 4D). Taken together, the  
536 abundance of ToBRFV is greater than PMMoV in human stool samples, and ToBRFV Mo gene  
537 may thus be a preferable MST marker compared to the the PMMoV CP gene.

538

539



540 **Fig.4: Prevalence of PMMoV and ToBRFV target genes in human stool samples.**

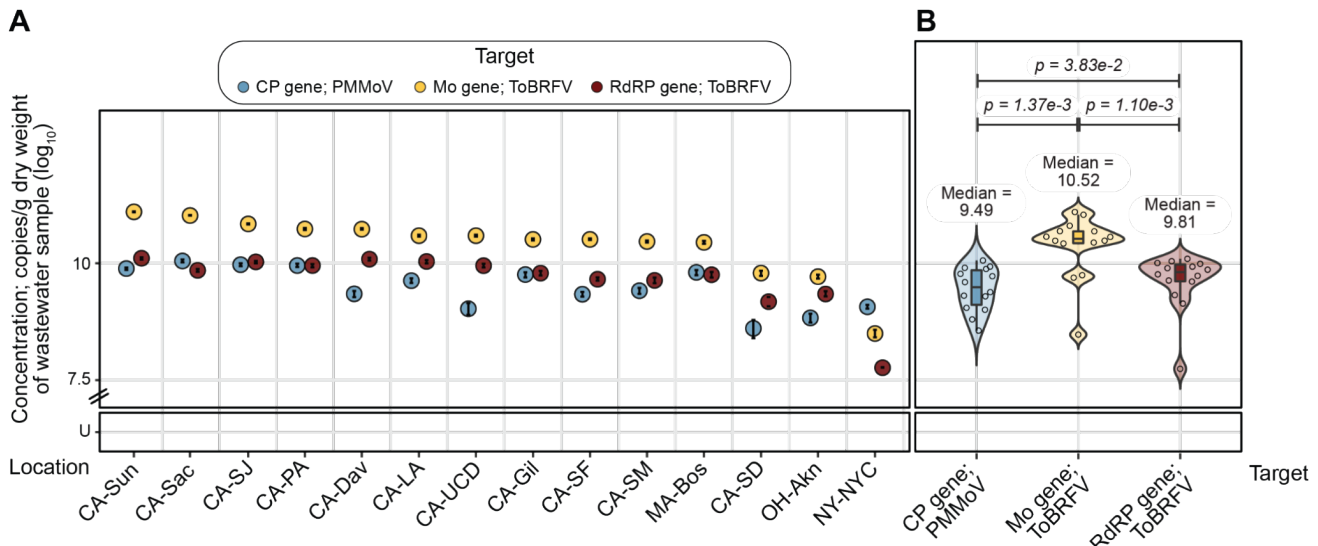
541 (A) Tabular summary of detection of the three gene targets in samples from adult and pediatric  
542 cohorts. The first column lists the name of the target gene and cohort, followed by the number  
543 and percent of samples that were positive (pos) or number of samples that were negative (neg)  
544 for that target gene, and the total number of samples tested (tot). (B) Venn diagram  
545 summarizing the detection of the PMMoV CP gene (blue), and ToBRFV Mo (yellow) and RdRP  
546 (red) genes across 220 human stool samples. In 34 (15.5%) we only detected the PMMoV CP  
547 gene, while 13 (5.9%) only the ToBRFV Mo gene. In 38 (17.3%) samples we detected both  
548 ToBRFV target genes, while 22 (10.0%) had both the PMMoV CP gene and ToBRFV Mo gene.  
549 In 70 (31.8%) of the samples we detected all three gene targets, while 43 (19.6%) had none of  
550 them. (C) Dot plot marking the concentrations of PMMoV CP (blue), ToBRFV Mo (red) and  
551 RdRP (yellow) genes, with violin and box plots summarizing their distributions, in RNA extracted  
552 from stool samples collected from humans. The x-axis marks the target genes, and the y-axis  
553 lists their concentrations in  $\log_{10}$  copies/  $\mu$ l of template; U stands for “Undetermined” and marks  
554 samples with no detectable gene target above LoB. Adult samples are marked by unfilled circles  
555 and pediatric samples are marked by a filled diamond. The concentration of PMMoV CP gene is  
556 a median of 1.13 with a standard deviation of 1.00 and IQR of 1.74  $\log_{10}$  copies/  $\mu$ l, ToBRFV Mo  
557 gene is a median of 2.12 with a standard deviation of 1.69 and IQR of 2.67  $\log_{10}$  copies/  $\mu$ l, and  
558 ToBRFV RdRP gene is a median of 2.20 with a standard deviation of 1.56 and IQR of 2.72  $\log_{10}$   
559 copies/  $\mu$ l. *p* values derived from paired Wilcoxon signed-rank tests with continuity correction  
560 and excluding samples with undetermined concentration across all combinations of the three  
561 gene targets are listed at the top of the plot. (D) Dot plot marking the relative abundance of viral  
562 reads of PMMoV (blue) and ToBRFV (yellow) from previously published metatranscriptomics  
563 data derived from healthy stool samples. The x-axis lists the 10 donors who provided samples,  
564 and each sample provided RNA sequences in biological triplets; each dot denotes a single  
565 replicate. The y-axis lists relative abundance in percent. (E) Dot plot summarizing data from  
566 panel D, now including violin and box plots to highlight distribution of viral RNA concentrations  
567 and associated statistics. The x-axis marks the target viral RNA, and the y-axis lists their relative  
568 abundance in percent. Dots represent the average of data from three biological replicates.  
569 PMMoV (blue) is present at a median relative abundance of 0.217% with a standard deviation of  
570 9.83% and IQR of 5.19%, ToBRFV (yellow) is present at a median relative abundance of 46.7  
571 with a standard deviation of 48.5% and IQR of 95.4%. *p* value derived from a Wilcoxon signed-  
572 rank test of pairwise differences in relative abundance with continuity correction and excluding  
573 samples with undetermined concentration is listed at the top of the plot.

#### 574 **ToBRFV Mo gene is prevalent and abundant in wastewater samples.**

575 Wastewater is a complex matrix containing human stool and other biological excretions,  
576 in addition to food waste, industrial wastes, and infiltrated stormwater in some cases. We next  
577 validate the molecular detection test developed here for testing this sample type. We acquired  
578 wastewater solids samples from 15 cities in the U.S.A, extracted RNA and assayed these for the  
579 presence and abundance of the gene targets of interest.

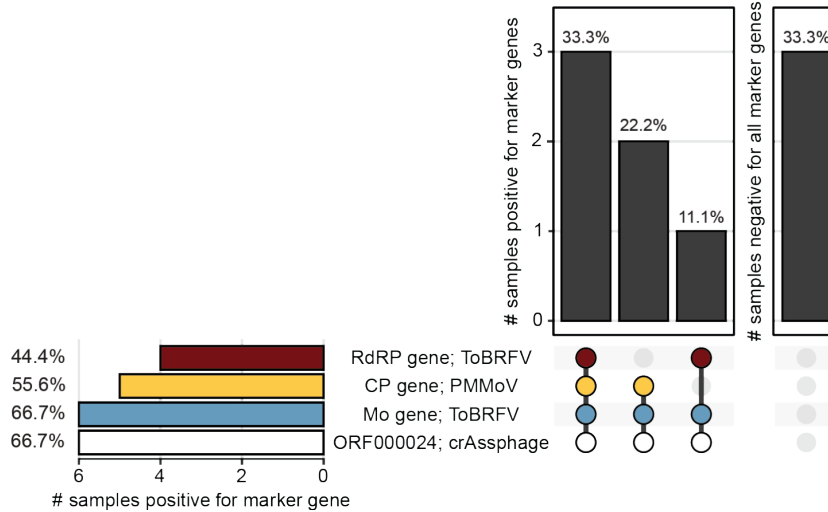
580 The extracted RNA from Wisconsin failed to have detectable amounts of any of the gene  
581 targets; this matches unpublished data generated using this sample by a different group and is

582 excluded from further analysis reducing our sample size to 14. 13 of these samples have more  
 583 ToBRFV Mo gene than the other two molecular markers, with the sample from New York being  
 584 the exception, having PMMoV CP gene in the highest concentration (Fig. 5A). Looking at the  
 585 data in aggregate, the samples have a median concentration of 10.5 log<sub>10</sub> copies/ g dry weight  
 586 of wastewater solids with a standard deviation of 0.67 and IQR of 0.26 log<sub>10</sub> copies/ g of the  
 587 ToBRFV Mo gene, followed by 9.81 log<sub>10</sub> copies/ g, standard deviation of 0.60 and IQR of 0.36  
 588 log<sub>10</sub> copies/ g of the ToBRFV RdRP gene, and 9.49 log<sub>10</sub> copies/ g, standard deviation of 0.46  
 589 and IQR of 0.74 log<sub>10</sub> copies/ g of the PMMoV CP gene. Pair-wise comparison of gene target  
 590 concentrations across samples using the Wilcoxon signed-rank test reveals that the increased  
 591 detection of the ToBRFV Mo gene is statistically significant in comparison to the PMMoV CP  
 592 gene ( $p = 1.37e-3$ ) and the ToBRFV RdRP gene ( $p = 1.10e-3$ ).



593 **Fig. 5: Concentrations of PMMoV and ToBRFV target genes in wastewater samples from**  
 594 **across the U.S.A.**  
 595 (A) Dot plot marking the concentrations of PMMoV CP (blue), ToBRFV Mo (yellow) and  
 596 ToBRFV RDRP (red) genes across samples. Error bars marking the standard deviation are  
 597 plotted along with the dots, and are mostly subsumed within the dot. The x-axis lists the 15 cities  
 598 from where samples were sourced in decreasing concentration of the Mo gene; abbreviations  
 599 for state and cities are expanded in Table S2. The y-axis lists concentrations of the genes in  
 600 log<sub>10</sub> copies/ g dry weight. (B) Dot plot marking the concentrations of PMMoV CP (blue),  
 601 ToBRFV Mo (red) and RdRP (yellow) genes, with violin and box plots summarizing their  
 602 distributions, in RNA extracted from wastewater samples collected from across the U.S.A. The

603 x-axis marks the target genes, and the y-axis lists their concentrations in log<sub>10</sub> copies/ g of  
 604 wastewater; U stands for “Undetermined” and marks samples with no detectable gene target  
 605 above LoB. The PMMoV CP gene has a median of 9.49 with a standard deviation of 0.46 and  
 606 IQR of 0.74 log<sub>10</sub> copies/ g dry weight of wastewater sample, ToBRFV Mo gene has a median of  
 607 10.5 with a standard deviation of 0.67 and IQR of 0.26 log<sub>10</sub> copies/ g, and ToBRFV RdRP gene



608 has a median of 9.81 with a standard deviation of 0.60 and IQR of 0.36 log<sub>10</sub> copies/ g. *p* values  
 609 derived from paired Wilcoxon signed-rank tests with continuity correction across all  
 610 combinations of the three gene targets are listed at the top of the plot. U stands for  
 611 “Undetermined” and marks samples with no detectable gene target above LoB.

612 **ToBRFV Mo gene matches crAssphage ORF000024 as an indicator of fecal contamination**  
 613 **of stormwater.**

614 crAssphage ORF000024 is a well established human-associated microbial source  
 615 tracking marker (12). We compared concentrations of PMMoV and ToBRFV RNA targets to  
 616 those of this crAssphage DNA target in stormwater draining from urbanized watersheds in the  
 617 Bay Area. crAssphage ORF000024 was previously quantified in these samples and is reported  
 618 in Graham *et. al.* (39).

619 We found that in the nine stormwater samples, crAssphage ORF000024 had the highest  
 620 median concentration of 4.65 with a standard deviation of 0.56 and IQR of 0.66 log<sub>10</sub> copies/ liter  
 621 of stormwater, followed by the ToBRFV RdRP gene with a median of 3.48, standard deviation of  
 622 0.97 and IQR of 1.24 log<sub>10</sub> copies/ liter of stormwater, ToBRFV Mo gene with a median of 3.34,  
 623 standard deviation of 0.99 and IQR of 1.36 log<sub>10</sub> copies/ liter of stormwater, and finally the



624 PMMoV CP gene with a median of 3.02, standard deviation of 0.54 and IQR of 0.44 log<sub>10</sub>  
625 copies/ liter of stormwater (Fig. S7). Pair-wise comparison of gene target concentrations across  
626 samples using the Wilcoxon signed-rank test reveals that differences in concentrations are not  
627 statistically significant and gene targets are similarly abundant. The concentration of gene  
628 targets in each of the samples is presented in Fig. S7. Notably, the ToBRFV Mo gene is  
629 detected in as many samples (6/9) as crAssphage ORF000024 (Fig. 6). This result suggests  
630 using an RNA based marker from ToBRFV to detect human stool contamination of storm water  
631 may be as useful as using the DNA marker from crAssphage ORF000024.

632 **Fig. 6: Prevalence of PMMoV, ToBRFV and crAssphage target genes in stormwater**  
633 **samples from across California.**  
634 UpSet plot summarizing the number of stormwater samples (total n = 9) that are either positive  
635 for multiple marker genes (left), or negative for all marker genes (right) in the vertical bar plots.  
636 Marker genes are listed under the plots, with colored dots representing positive presence and  
637 grey dots representing absence. Marker genes present in samples represented in the vertical  
638 bar are also connected by a thick line. Prevalence of independent marker genes are also  
639 summarized in the horizontal bar plot. All bars list data in percentage units. PMMoV CP gene is  
640 marked by blue, ToBRFV Mo and RdRP genes by yellow and red respectively, and crAssphage  
641 ORF000024 by white. Data on crAssphage are derived from a previous study (39).

642

### 643 **Conclusions and limitations**

644 In this study, we generate eight nearly complete genomes of ToBRFV from wastewater  
645 and stool from the Bay Area. We catalog SNPs in all existing genomes, including in those that  
646 we assembled here, and note variations in viral genomes isolated from the same individual over  
647 ~ 100 days. We then went on to identify two sets of primers and probes that can universally  
648 detect ToBRFV across the world.

649 Assays developed using these primer, probe sequences are sensitive and specific for  
650 human stool and wastewater as it was present in a wide range of wastewaters and stool  
651 samples, and not present in any tested animal stool aside from one cat. Like the established  
652 viral MST target PMMoV (8), the ToBRFV target is derived from the genome of a plant virus

653 likely present in the human gut owing to dietary intake of diseased plants. Concentrations of  
654 ToBRFV Mo and RdRP gene targets were as high or higher than PMMoV CP gene in  
655 wastewaters and stormwater known to contain sewage. The ToBRFV targets' high  
656 concentrations in wastewater, as well as in human stool samples, suggests that they may be  
657 useful as endogenous fecal strength controls for wastewater-based epidemiology applications  
658 (46), as well as an endogenous positive extraction control during nucleic acid extractions in  
659 studies seeking to quantify rare infectious disease targets (8).

660         There are several limitations to this work. First, the specificity of the ToBRFV Mo and  
661 RdRP gene targets was tested using just one representative sample of various non-human,  
662 animal stools. Additional work to test more animal stool samples would be helpful to further  
663 characterize the assays' specificity for human stool. Second, the sensitivities of the various  
664 assays were tested using only human stool samples from individuals residing in the Bay Area. It  
665 is possible that the distribution of the targets in individuals from other locations may differ from  
666 those studied here and more work is encouraged to document the ToBRFV prevalence and  
667 abundance in samples globally. Third, we assayed wastewater solids sampled from around the  
668 U.S.A, from New York to California, and they contained high concentrations of the ToBRFV  
669 targets. Further work from samples around the world will be valuable to testing the  
670 generalizability of the assays. Notably, the presence of ToBRFV genomes from this study and  
671 others collected from many countries reassures us that ToBRFV is likely to be a universal global  
672 MST marker. Finally, As more ToBRFV genomes become available, it will be important to test  
673 whether the primers and probes developed herein continue to overlap with conserved regions of  
674 the genomes.

675

676 **Acknowledgements**

677           This work was supported by a ChemH-IMA grant (to A.S.B.), a gift from the CDC-  
678 Foundation (to A.B.B.), Stanford MCHRI and NIH T32 DK098132 (C.J.S.), DGE-1656518 and  
679 T32GM007276 (to D.T.S.), and the Stanford Dean's Postdoctoral Fellowship (to A.N.). A.S.B.  
680 laboratory is supported by NIH R01 AI148623 and R01 AI143757. We acknowledge Aaron Behr,  
681 Alvin Han, David Miklos, David Solow-Cordero, Dhananjay Wagh, Jennifer Estes, Isabel Delwel,  
682 Luisa Jiminez, Said Attiya, Sopheak Sim, and Summer Vance for technical assistance, and  
683 sharing of samples and data for use in this study. A.N., B.J.F, M.W. and A.B.B., are co-inventors  
684 on a U.S. provisional patent application #63/387,657 that has been filed and relates to the  
685 methods presented in this manuscript. The other authors declare no competing interests. This  
686 study was performed on the ancestral and unceded lands of the Muwekma Ohlone people. We  
687 pay our respects to them and their Elders, past and present, and are grateful for the opportunity  
688 to live and work here.

689

690 **References**

- 691 1. Whitman RL, Shively DA, Pawlik H, Nevers MB, Byappanahalli MN. 2003. Occurrence of  
692 *Escherichia coli* and enterococci in *Cladophora* (Chlorophyta) in nearshore water and  
693 beach sand of Lake Michigan. *Appl Environ Microbiol* 69:4714–4719.
- 694 2. Us Epa OW. 2013. Recreational Water Quality Criteria and Methods.
- 695 3. Us Epa OW. 2015. Drinking Water Regulations.
- 696 4. Layton BA, Walters SP, Lam LH, Boehm AB. 2010. Enterococcus species distribution  
697 among human and animal hosts using multiplex PCR. *J Appl Microbiol* 109:539–547.
- 698 5. Imamura GJ, Thompson RS, Boehm AB, Jay JA. 2011. Wrack promotes the persistence of  
699 fecal indicator bacteria in marine sands and seawater. *FEMS Microbiol Ecol* 77:40–49.
- 700 6. Yamahara KM, Walters SP, Boehm AB. 2009. Growth of enterococci in unaltered,  
701 unseeded beach sands subjected to tidal wetting. *Appl Environ Microbiol* 75:1517–1524.
- 702 7. Byappanahalli MN, Whitman RL, Shively DA, Sadowsky MJ, Ishii S. 2006. Population  
703 structure, persistence, and seasonality of autochthonous *Escherichia coli* in temperate,  
704 coastal forest soil from a Great Lakes watershed. *Environ Microbiol* 8:504–513.
- 705 8. McClary-Gutierrez JS, Aanderud ZT, Al-Faliti M, Duvallet C, Gonzalez R, Guzman J, Holm  
706 RH, Jahne MA, Kantor RS, Katsivelis P, Kuhn KG, Langan LM, Mansfeldt C, McLellan SL,  
707 Grijalva LMM, Murnane KS, Naughton CC, Packman AI, Paraskevopoulos S, Radniecki TS,  
708 Roman FA Jr, Shrestha A, Stadler LB, Steele JA, Swalla BM, Vikesland P, Wartell B,  
709 Wilusz CJ, Wong JCC, Boehm AB, Halden RU, Bibby K, Vela JD. 2021. Standardizing data  
710 reporting in the research community to enhance the utility of open data for SARS-CoV-2  
711 wastewater surveillance. *Environ Sci* 9.

- 712 9. Boehm AB, Van De Werfhorst LC, Griffith JF, Holden PA, Jay JA, Shanks OC, Wang D,  
713 Weisberg SB. 2013. Performance of forty-one microbial source tracking methods: a twenty-  
714 seven lab evaluation study. *Water Res* 47:6812–6828.
- 715 10. Shanks OC, White K, Kelty CA, Hayes S, Sivaganesan M, Jenkins M, Varma M, Haugland  
716 RA. 2010. Performance assessment PCR-based assays targeting bacteroidales genetic  
717 markers of bovine fecal pollution. *Appl Environ Microbiol* 76:1359–1366.
- 718 11. Green HC, Dick LK, Gilpin B, Samadpour M, Field KG. 2012. Genetic markers for rapid  
719 PCR-based identification of gull, Canada goose, duck, and chicken fecal contamination in  
720 water. *Appl Environ Microbiol* 78:503–510.
- 721 12. García-Aljaro C, Ballesté E, Muniesa M, Jofre J. 2017. Determination of crAssphage in  
722 water samples and applicability for tracking human faecal pollution. *Microb Biotechnol*  
723 10:1775–1780.
- 724 13. Rosario K, Symonds EM, Sinigalliano C, Stewart J, Breitbart M. 2009. Pepper mild mottle  
725 virus as an indicator of fecal pollution. *Appl Environ Microbiol* 75:7261–7267.
- 726 14. Edwards RA, Vega AA, Norman HM, Ohaeri M, Levi K, Dinsdale EA, Cinek O, Aziz RK,  
727 McNair K, Barr JJ, Bibby K, Brouns SJJ, Cazares A, de Jonge PA, Desnues C, Díaz Muñoz  
728 SL, Fineran PC, Kuriishikov A, Lavigne R, Mazankova K, McCarthy DT, Nobrega FL, Reyes  
729 Muñoz A, Tapia G, Trefault N, Tyakht AV, Vinuesa P, Wagemans J, Zhernakova A,  
730 Aarestrup FM, Ahmadov G, Alassaf A, Anton J, Asangba A, Billings EK, Cantu VA, Carlton  
731 JM, Cazares D, Cho G-S, Condeff T, Cortés P, Cranfield M, Cuevas DA, De la Iglesia R,  
732 Decewicz P, Doane MP, Dominy NJ, Dziewit L, Elwasila BM, Eren AM, Franz C, Fu J,  
733 Garcia-Aljaro C, Ghedin E, Gulino KM, Haggerty JM, Head SR, Hendriksen RS, Hill C,  
734 Hyöty H, Iliina EN, Irwin MT, Jeffries TC, Jofre J, Junge RE, Kelley ST, Khan Mirzaei M,

- 735 Kowalewski M, Kumaresan D, Leigh SR, Lipson D, Lisitsyna ES, Llagostera M, Maritz JM,  
736 Marr LC, McCann A, Molshanski-Mor S, Monteiro S, Moreira-Grez B, Morris M, Mugisha L,  
737 Muniesa M, Neve H, Nguyen N-P, Nigro OD, Nilsson AS, O'Connell T, Odeh R, Oliver A,  
738 Piuri M, Prussin AJ II, Qimron U, Quan Z-X, Rainetova P, Ramírez-Rojas A, Raya R,  
739 Reasor K, Rice GAO, Rossi A, Santos R, Shimashita J, Stachler EN, Stene LC, Strain R,  
740 Stumpf R, Torres PJ, Twaddle A, Ugochi Ibekwe M, Villagra N, Wandro S, White B,  
741 Whiteley A, Whiteson KL, Wijmenga C, Zambrano MM, Zschach H, Dutilh BE. 2019. Global  
742 phylogeography and ancient evolution of the widespread human gut virus crAssphage. *Nat*  
743 *Microbiol* 4:1727–1736.
- 744 15. Colson P, Richet H, Desnues C, Balique F, Moal V, Grob J-J, Berbis P, Lecoq H, Harlé J-R,  
745 Berland Y, Raoult D. 2010. Pepper mild mottle virus, a plant virus associated with specific  
746 immune responses, Fever, abdominal pains, and pruritus in humans. *PLoS One* 5:e10041.
- 747 16. Zhang S, Griffiths JS, Marchand G, Bernards MA, Wang A. 2022. Tomato brown rugose  
748 fruit virus: An emerging and rapidly spreading plant RNA virus that threatens tomato  
749 production worldwide. *Mol Plant Pathol* 23:1262–1277.
- 750 17. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing  
751 reads. *EMBnet.journal* <https://doi.org/10.14806/ej.17.1.200>.
- 752 18. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. 2017. metaSPAdes: a new versatile  
753 metagenomic assembler. *Genome Res* 27:824–834.
- 754 19. Minot SS, Krumm N, Greenfield NB. One Codex: A Sensitive and Accurate Data Platform  
755 for Genomic Microbial Identification <https://doi.org/10.1101/027607>.
- 756 20. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal:

- 757 prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*  
758 11:119.
- 759 21. Nayfach S, Camargo AP, Schulz F, Eloë-Fadrosch E, Roux S, Kyrpides NC. 2020. CheckV  
760 assesses the quality and completeness of metagenome-assembled viral genomes. *Nat*  
761 *Biotechnol* <https://doi.org/10.1038/s41587-020-00774-7>.
- 762 22. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
763 9:357–359.
- 764 23. Olm MR, Crits-Christoph A, Bouma-Gregson K, Firek BA, Morowitz MJ, Banfield JF. 2021.  
765 inStrain profiles population microdiversity from metagenomic data and sensitively detects  
766 shared microbial strains. *Nat Biotechnol* <https://doi.org/10.1038/s41587-020-00797-0>.
- 767 24. Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2.  
768 *Genome Biol* 20:257.
- 769 25. 2019. Geneious. Geneious. <https://www.geneious.com>. Retrieved 21 July 2022.
- 770 26. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. : rapid  
771 efficient extraction of SNPs from multi-FASTA alignments. *Microb Genom* 2:e000056.
- 772 27. Rozen S, Skaletsky H. 2000. Primer3 on the WWW for general users and for biologist  
773 programmers. *Methods Mol Biol* 132:365–386.
- 774 28. Sayers EW, Bolton EE, Brister JR, Canese K, Chan J, Comeau DC, Farrell CM, Feldgarden  
775 M, Fine AM, Funk K, Hatcher E, Kannan S, Kelly C, Kim S, Klimke W, Landrum MJ, Lathrop  
776 S, Lu Z, Madden TL, Malheiro A, Marchler-Bauer A, Murphy TD, Phan L, Pujar S,  
777 Rangwala SH, Schneider VA, Tse T, Wang J, Ye J, Trawick BW, Pruitt KD, Sherry ST.

- 778 2022. Database resources of the National Center for Biotechnology Information in 2023.  
779 Nucleic Acids Res <https://doi.org/10.1093/nar/gkac1032>.
- 780 29. Rothman JA, Whiteson KL. 2022. Sequencing and Variant Detection of Eight Abundant  
781 Plant-Infecting Tobamoviruses across Southern California Wastewater. *Microbiol Spectr*  
782 e0305022.
- 783 30. Nishimura Y, Yoshida T, Kuronishi M, Uehara H, Ogata H, Goto S. 2017. ViPTree: the viral  
784 proteomic tree server. *Bioinformatics* 33:2379–2380.
- 785 31. Moore RM, Harrison AO, McAllister SM, Polson SW, Wommack KE. 2020. Iroki: automatic  
786 customization and visualization of phylogenetic trees. *PeerJ* 8:e8584.
- 787 32. Natarajan A, Han A, Zlitni S, Brooks EF, Vance SE, Wolfe M, Singh U, Jagannathan P,  
788 Pinsky BA, Boehm A, Bhatt AS. 2021. Standardized preservation, extraction and  
789 quantification techniques for detection of fecal SARS-CoV-2 RNA. *Nat Commun* 12:5753.
- 790 33. Maghini D, Dvorak M, Dahlen A, Roos M, Kuersten S, Bhatt AS. 2022. Achieving  
791 quantitative and accurate measurement of the human gut microbiome. *bioRxiv*.
- 792 34. Franzosa EA, Morgan XC, Segata N, Waldron L, Reyes J, Earl AM, Giannoukos G, Boylan  
793 MR, Ciulla D, Gevers D, Izard J, Garrett WS, Chan AT, Huttenhower C. 2014. Relating the  
794 metatranscriptome and metagenome of the human gut. *Proceedings of the National*  
795 *Academy of Sciences* <https://doi.org/10.1073/pnas.1319284111>.
- 796 35. Huisman JS, Scire J, Caduff L, Fernandez-Cassi X, Ganesanandamoorthy P, Kull A,  
797 Scheidegger A, Stachler E, Boehm AB, Hughes B, Knudson A, Topol A, Wigginton KR,  
798 Wolfe MK, Kohn T, Ort C, Stadler T, Julian TR. 2022. Wastewater-Based Estimation of the  
799 Effective Reproductive Number of SARS-CoV-2. *Environ Health Perspect* 130:57011.



- 800 36. Wolfe MK, Topol A, Knudson A, Simpson A, White B, Vugia DJ, Yu AT, Li L, Balliet M,  
801 Stoddard P, Han GS, Wigginton KR, Boehm AB. 2021. High-Frequency, High-Throughput  
802 Quantification of SARS-CoV-2 RNA in Wastewater Settled Solids at Eight Publicly Owned  
803 Treatment Works in Northern California Shows Strong Association with COVID-19  
804 Incidence. *mSystems* 6:e0082921.
- 805 37. Topol A, Wolfe M, White B, Wigginton K, Boehm AB. 2021. High Throughput pre-analytical  
806 processing of wastewater settled solids for SARS-CoV-2 RNA analyses. *protocols.io*.  
807 [https://www.protocols.io/view/high-throughput-pre-analytical-processing-of-waste-](https://www.protocols.io/view/high-throughput-pre-analytical-processing-of-waste-kxygxpod4l8j/v2)  
808 [kxygxpod4l8j/v2](https://www.protocols.io/view/high-throughput-pre-analytical-processing-of-waste-kxygxpod4l8j/v2). Retrieved 6 September 2022.
- 809 38. Topol A, Wolfe M, Wigginton K, White B, Boehm A. 2021. High Throughput RNA Extraction  
810 and PCR Inhibitor Removal of Settled Solids for Wastewater Surveillance of S. *protocols.io*.  
811 [https://www.protocols.io/view/high-throughput-rna-extraction-and-pcr-inhibitor-r-](https://www.protocols.io/view/high-throughput-rna-extraction-and-pcr-inhibitor-r-81wgb72bovpk/v2)  
812 [81wgb72bovpk/v2](https://www.protocols.io/view/high-throughput-rna-extraction-and-pcr-inhibitor-r-81wgb72bovpk/v2). Retrieved 6 September 2022.
- 813 39. Graham KE, Anderson CE, Boehm AB. 2021. Viral pathogens in urban stormwater runoff:  
814 Occurrence and removal via vegetated biochar-amended biofilters. *Water Res* 207:117829.
- 815 40. Kuypers J, Jerome KR. 2017. Applications of Digital PCR for Clinical Microbiology. *J Clin*  
816 *Microbiol* 55:1621–1628.
- 817 41. Droplet Digital PCR Applications Guide (6407 Ver B). Bio-Rad.
- 818 42. dMIQE Group, Huggett JF. 2020. The Digital MIQE Guidelines Update: Minimum  
819 Information for Publication of Quantitative Digital PCR Experiments for 2020. *Clin Chem*  
820 66:1012–1029.
- 821 43. Loeb S. 2020. One-Step RT-ddPCR for Detection of SARS-CoV-2, Bovine Coronavirus,

- 822 and PMMoV RNA in RNA Derived from Wastewater or Primary Settled Solids.
- 823 44. Crits-Christoph A, Kantor RS, Olm MR, Whitney ON, Al-Shayeb B, Lou YC, Flamholz A,  
824 Kennedy LC, Greenwald H, Hinkle A, Hetzel J, Spitzer S, Koble J, Tan A, Hyde F, Schroth  
825 G, Kuersten S, Banfield JF, Nelson KL. 2021. Genome Sequencing of Sewage Detects  
826 Regionally Prevalent SARS-CoV-2 Variants. *MBio* 12.
- 827 45. Rothman JA, Loveless TB, Kapcia J 3rd, Adams ED, Steele JA, Zimmer-Faust AG,  
828 Langlois K, Wanless D, Griffith M, Mao L, Chokry J, Griffith JF, Whiteson KL. 2021. RNA  
829 Viromics of Southern California Wastewater and Detection of SARS-CoV-2 Single-  
830 Nucleotide Variants. *Appl Environ Microbiol* 87:e0144821.
- 831 46. McClary-Gutierrez JS, Mattioli MC, Marcenac P, Silverman AI, Boehm AB, Bibby K, Balliet  
832 M, de Los Reyes FL 3rd, Gerrity D, Griffith JF, Holden PA, Katehis D, Kester G, LaCross N,  
833 Lipp EK, Meiman J, Noble RT, Brossard D, McLellan SL. 2021. SARS-CoV-2 Wastewater  
834 Surveillance for Public Health Action. *Emerg Infect Dis* 27:1–8.
- 835