

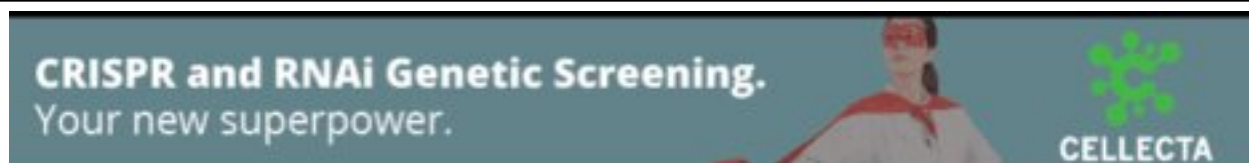


Enhanced virome sequencing through solution-based capture enrichment

Todd N Wylie, Kristine M Wylie, Brandi N Herter, et al.

Genome Res. published online September 22, 2015
Access the most recent version at doi:[10.1101/gr.191049.115](https://doi.org/10.1101/gr.191049.115)

P<P	Published online September 22, 2015 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Enhanced Virome Sequencing Using Targeted Sequence Capture**

2

3 Todd N. Wylie^{1,2,*}, Kristine M. Wylie^{1,2,*}, Brandi N. Herter¹, and Gregory A. Storch¹

4

5 The Department of Pediatrics¹, Washington University School of Medicine, Campus Box 8208,
6 660 S. Euclid Avenue, St Louis, MO 63110, USA

7 McDonnell Genome Institute², Washington University School of Medicine, 4444 Forest Park
8 Avenue, St Louis 63108, MO, USA

9 *These authors contributed equally to this work.

10

11 Corresponding Author

12 Todd N. Wylie

13 wylie_t@kids.wustl.edu

14 Campus Box 8208

15 Washington University School of Medicine

16 660 S. Euclid Avenue

17 St. Louis, MO 63110

18 314-747-4069

19

20 Running Title:

21 Enhanced virome sequencing using sequence capture

22

23 **Keywords:** targeted sequence capture; infectious disease; virus; virome; microbiome;
24 metagenomics; genomics; next-generation sequencing; deep sequencing; ViroCap

25

26

27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52

ABSTRACT

Metagenomic shotgun sequencing (MSS) is an important tool for characterizing viral populations. It is culture-independent, requires no *a priori* knowledge of the viruses in the sample, and may provide useful genomic information. However, MSS can lack sensitivity and may yield insufficient data for detailed analysis. We have created a targeted sequence capture panel, ViroCap, designed to enrich nucleic acid from DNA and RNA viruses from 34 families that infect vertebrate hosts. A computational approach condensed nearly 1 billion base pair (bp) of viral reference sequence into less than 200 million bp of unique, representative sequence suitable for targeted sequence capture. We compared the effectiveness of detecting viruses in standard MSS versus MSS following targeted sequence capture. First, we analyzed two sets of samples, one derived from samples submitted to a diagnostic virology laboratory and one derived from samples collected in a study of fever in children. We detected 14 and 18 viruses in the two sets, comprising 19 genera from 10 families, with dramatic enhancement of genome representation following capture enrichment. The median fold-increases in percent viral reads post-capture were 674 and 296. Median breadth of coverage increased from 2.1% to 83.2% post-capture in the first set and from 2.0% to 75.6% in the second set. Next, we analyzed samples containing a set of diverse anellovirus sequences and demonstrated that ViroCap could be used to detect viral sequences with up to 58% variation from the references used to select capture probes. ViroCap substantially enhances MSS for a comprehensive set of viruses and has utility for research and clinical applications.

53 INTRODUCTION

54 High-throughput, massively parallel nucleotide sequence analysis has made in-depth studies of
55 the human microbiome feasible. Thus far, most microbiome studies have focused on bacteria
56 (Arumugam et al. 2011; Human Microbiome Project Consortium 2012; Turnbaugh et al. 2009;
57 Gajer et al. 2012), although some include fungi (Cleland et al. 2014; Paulino et al. 2006; Findley
58 et al. 2013; Willger et al. 2014) and viruses (Wylie et al. 2014; 2012; De Vlaminck et al. 2013;
59 Minot et al. 2011; Reyes et al. 2010). Viruses are particularly understudied, in part due to the
60 challenges of assessing their presence in clinical samples. Viruses as a group have highly
61 variable genomes, with no gene shared among all viruses that can be surveyed by an amplicon-
62 based sequencing strategy. Therefore, studies of viruses based on nucleotide sequencing
63 require a metagenomic approach. Metagenomic shotgun sequencing (MSS) is a relatively
64 unbiased, culture-independent method in which nucleic acid extracted from a sample is
65 sequenced. Sequence reads are classified based on similarity to reference genomes. This
66 approach allows comprehensive study of the viral component of the microbiome (the virome)
67 and has led to the discovery of novel viruses (reviewed in (Chiu 2013)) and the characterization
68 of viruses present in healthy and sick people (Wylie et al. 2012; 2014; Reyes et al. 2010; Holtz
69 et al. 2014; Oh et al. 2014; Lysholm et al. 2012; Minot et al. 2011; Young et al. 2014). When
70 adequate numbers of sequence reads are generated, viruses can be characterized with regard
71 to taxonomy and the presence of genes associated with virulence and resistance to antiviral
72 drugs.

73
74 A limitation of MSS as employed to date for virus detection is that the amount and proportion of
75 viral nucleic acid in samples from humans may be very low, and in these cases few viral
76 sequences are generated. In our experience using MSS, we have detected fewer than 10 viral
77 sequences per 25 million sequence reads generated for a virus that was detected in a sample
78 by a molecular assay (Wylie et al. 2012). In other instances we have failed to detect viruses

79 known to be present based on molecular assays (Wylie et al. 2012). These difficulties may
80 reflect the small genome size of some viruses and/or low levels of virus in the sample. This can
81 be a particular problem for studies of the virome of healthy, asymptomatic individuals (Wylie et
82 al. 2014; 2012), in whom virus levels may be low. In efforts to increase the sequence yield,
83 purification or enrichment procedures have been employed, including low-speed centrifugation
84 and/or filtration to remove bacterial and host cells, sample treatment with nucleases to digest
85 nucleic acid not protected within virions (Allander et al. 2001), or concentration of viral particles
86 by high-speed gradient centrifugation (reviewed in (Duhaime and Sullivan 2012)). Each of these
87 procedures may bias against detection of some viruses (Young et al. 2014; Breitbart and
88 Rohwer 2005).

89

90 An alternative method for enrichment of viral sequences in a metagenomic sample prior to
91 sequencing is targeted sequence capture, a well-established approach for targeted enrichment
92 of specific nucleic acids. Targeted sequence capture has been used extensively to assess the
93 human exome as well as specific gene targets (Lovett et al. 1991; Okou et al. 2007; Hodges et
94 al. 2007; Albert et al. 2007). Sequence capture has also been applied to the study of specific
95 viruses (Duncavage et al. 2011; Depledge et al. 2011; Koehler JW et al. 2014). Our aim was to
96 develop a comprehensive viral targeted sequence capture panel that could be used to (1)
97 assess all viruses known to infect vertebrate cells and (2) detect divergent viruses. To this end,
98 we created ViroCap, a targeted sequence capture panel that enhances the detection of a
99 comprehensive set of viruses with vertebrate hosts. Here we describe the first application of
100 ViroCap to enrich a broad range of viruses from human clinical samples.

101

102 **RESULTS**

103 ViroCap includes targets from 34 viral families, comprising 190 annotated viral genera and 337
104 species (**Fig. 1**). Included viruses represent all DNA and RNA viruses with sequenced genomes

105 from vertebrate hosts, except human endogenous retroviruses, which were excluded due to
106 their prevalence within the human genome. Nearly 1 billion base pairs (bp) of viral genome
107 sequences were condensed into less than 200 million bp of targets (**Supplemental Table S1**)
108 using *k*-mer and clustering analyses to define a unique set of reference sequences, as
109 described in Methods.

110

111 **Analysis of Clinical and Research Samples with ViroCap**

112 We evaluated the effectiveness of detecting DNA and RNA viruses in MSS data compared with
113 ViroCap targeted sequence capture data in two sets of human samples. In experiment 1, the
114 sample set consisted of clinical samples that had been found to be positive by molecular tests in
115 the Diagnostic Virology Laboratory at St. Louis Children's Hospital. Nucleic acid extracts
116 available in the Virology Laboratory were pooled, and a sequencing library was prepared from
117 this pooled nucleic acid (see Methods). In experiment 2, 8 patient samples from a research
118 study of young children with fever (Colvin et al. 2012; Wylie et al. 2012) were selected for use in
119 the present study because each had been found to be positive for one or more viruses when
120 tested by batteries of PCR assays used in that study. Individual sequencing libraries were
121 prepared from each of the 8 samples as described in Methods and pooled for sequencing.
122 Experiments 1 and 2 were analyzed in separate sequencing runs. In each experiment,
123 sequencing libraries were divided and the same library was sequenced without targeted
124 sequence capture (pre-capture) and following targeted sequence capture using ViroCap (post-
125 capture).

126

127 In experiment 1, we detected 10 viruses in the pre-capture MSS data (**Table 1**). After targeted
128 sequence capture using the same sequencing library, we detected the same 10 viruses plus 4
129 additional viruses. Targeted sequence capture resulted in dramatic improvements in all
130 sequence coverage metrics (**Table 1 and Supplemental Table S2**), including number and

131 percent viral reads, breadth and depth of coverage, and coverage gaps. In experiment 1, the
132 median increase in percent viral reads was 674 (range >13 to 9335), and the median breadth of
133 coverage increased from 2.1% (range 0 to 89.8%) to 83.2% (range 0.8 to 100%). Illustrative
134 examples are shown in **Figure 2 A-D**.

135
136 In experiment 2, 11 viruses were detected in the pre-capture MSS data (**Table 2**). After targeted
137 sequence capture with ViroCap using the same sequencing libraries, we detected those 11
138 viruses plus 7 additional viruses. Thus, in the two experiments together, the number of viruses
139 detected went from 21 to 32, a 52% increase. All of the viruses detected in both experiments
140 were confirmed by PCR assays except for a torque teno virus in the clinical pool, which was not
141 evaluated by PCR (**Table 1, Supplemental Tables S8 and 9**). Viruses detected encompassed
142 19 genera from 10 families (**Supplemental Figure S1**). In experiment 2, we again found that
143 targeted sequence capture resulted in dramatic improvements in sequencing parameters. In
144 experiment 2, the median fold increase in percent viral reads was 296 (range >56 to 2,722), and
145 the median breadth of coverage increased from 2.0% (range 0 to 99.9%) to 75.6% (range 13.5
146 to 100%). Illustrative examples are shown in **Figure 2 E-H**.

147
148 Using targeted sequence capture, >80% breadth of coverage of the viral genomes was obtained
149 for 16 of 32 viruses, including diverse DNA and RNA genomes of sizes ranging from 5 to 161
150 kilobases (**Tables 1 and 2 and Supplemental Tables S2 and S3**). Greater than 90% breadth
151 of coverage was obtained for 12 of 32 viruses, and 8 viruses had 100% coverage. Pre-capture,
152 the median gap size in genome coverage was 1704 bp (range 4 to 152,261 bp) and post-
153 capture the median gap size was 82 bp (range 0 to 13,734) (**Supplemental Tables S2 and S3**).
154 High genome representation was obtained for multiple viruses in the same capture reaction, as
155 experiments 1 and 2 were each single, independent capture reactions encompassing multiple
156 samples (see Methods).

157

158 **Targeted sequence capture identifies divergent viral sequences**

159 To determine whether or not divergent sequences could be identified using targeted sequence
160 capture, we tested ViroCap on samples containing anelloviruses, a highly divergent group of
161 ssDNA viruses which have a common genome structure, but may have up to 30 to 50%
162 nucleotide sequence diversity among separate species (Ninomiya et al. 2007; de Villiers et al.
163 2011). We selected anellovirus-positive samples that we had previously characterized using
164 multi-strand displacement amplification followed by high-throughput sequencing. After
165 assembling the pre-capture sequences to generate contiguous sequences (contigs), we
166 identified anellovirus contigs greater than 1 kilobase in length. The contigs had varying degrees
167 of similarity to the reference genomes used in the ViroCap panel based on BLAST alignments,
168 ranging from 58% to 98% nucleotide sequence identity for the top high-scoring segment pair
169 (HSP) alignment (**Fig. 3A, Supplemental Table S4**). All of the contigs assembled using the pre-
170 capture sequence data were also detected post-capture. The contig with 58% identity to the
171 reference database was missing 13% of its length post-capture (**Fig. 3A**). The contig with the
172 next lowest percent identity to the reference database (62%) was fully sequenced (i.e. 100%
173 breadth-of-coverage) (**Figs. 3A and B**). Figure 3B illustrates the nucleotide sequence
174 matches/mismatches between the contig and the most similar reference genome in the
175 sequences used for the ViroCap design. These results demonstrate that targeted sequence
176 capture using the ViroCap panel allows us to identify variant virus sequences having as low as
177 58% nucleotide sequence identity.

178

179 **Specificity of targeted sequence capture**

180 In order to determine whether ViroCap systematically enriched off-target sequences, we
181 compared the filtering and classification statistics of the non-viral sequences in the pre-capture
182 MSS and targeted sequence capture data (**Supplemental Table S5, Supplemental Figures**

183 **S2 and S3**). If our probes were specific, we would not observe any systematic enrichment of
184 specific human chromosomes or bacterial genomes post-capture. However, we anticipated a
185 small amount of variation because the targeted sequence capture library had been through
186 more sample handling in the form of incubations, dilutions, and amplifications. We found that the
187 proportions of the non-viral sequences were strongly correlated (Pearson's correlation value:
188 $r=0.9881$ to 0.9996) (**Supplemental Table S5 and Supplemental Figure S2**). A slightly higher
189 percentage (mean 5.8%, median 5%, range 0 to 10.7%) of reads aligned to non-viral reference
190 genomes in the post-capture data compared with pre-capture in all but one of the samples.
191 However, the distribution of sequences among reference genomes did not show a systematic
192 bias. This can be seen in the conserved distribution of sequences among human chromosomes
193 (**Supplemental Figure S3**).

194

195 **DISCUSSION**

196 We designed the ViroCap panel to enhance the sensitivity of metagenomic shotgun sequencing
197 for comprehensive detection of known vertebrate viruses as well as detection of divergent
198 viruses that have nucleotide sequence similarities to known viruses. Here we have
199 demonstrated that targeted sequence capture using ViroCap dramatically increases the amount
200 of viral sequence obtained from human samples compared to conventional MSS, greatly
201 enhancing the resolution of genomic characterization and increasing the number of viruses
202 detected by more than 50%. Enhancement was demonstrated for DNA and RNA viruses from
203 multiple diverse families. The increased sensitivity will be valuable in multiple research
204 applications including descriptions of the human virome and will also improve the potential for
205 MSS as a diagnostic tool in human and animal health.

206

207 The dramatic enrichment of viral nucleic acids present within the targeted sequence capture
208 libraries offers important advantages. First, as we demonstrate, MSS with ViroCap can be used

209 to generate complete or nearly complete genome sequences directly from clinical samples,
210 including those with very low proportions of viral nucleic acid, without culturing the viruses.
211 Availability of extensive sequence data provides the opportunity to distinguish among closely
212 related virus subtypes or even among viral strains, which might not be distinguished by other
213 types of assays. In the data set presented here, we demonstrated the ability to type
214 rhinoviruses, and distinguish between human herpesvirus 6B and 6A, adenovirus types A and
215 C, and polyomaviruses JC and BK. Notably, influenza A virus was identified pre-capture but
216 could only be typed as an H3N2 virus post-capture. Elsewhere, we used ViroCap to sequence
217 the enterovirus D68 genome directly from clinical samples (Wylie et al. 2015), and in that work,
218 the extensive sequence data that we obtained allowed us carry out detailed comparative
219 analysis of closely related strains that differed at a limited number of nucleotide positions.
220 Second, the use of ViroCap can reduce the depth of sequencing needed to detect viruses in
221 samples. Because targeted sequence capture results in a large increase in the percentage of
222 sequencing reads that are viral (**Fig. 2, Tables 2 and 3, Supplemental Tables 2-4**), ViroCap
223 achieves better viral coverage while requiring the generation of fewer total sequence reads. This
224 increased efficiency has the potential to lower sequencing costs.

225
226 An important feature of ViroCap is the tiling of capture probes across genomes, including highly
227 conserved regions which may allow detection of genomic fragments of divergent viruses that
228 share little overall sequence homology with known viruses. We illustrated such capability using
229 anelloviruses containing divergent nucleotide sequence (**Fig. 3**). In addition, the inclusion of
230 Genome Neighbor targets enhanced our design not only by expanding beyond the tiled RefSeq
231 viruses, but also by adding sensitivity for genomic regions where RefSeq capture probes alone
232 might not have captured divergent strains (see Methods). ViroCap cannot detect viruses that do
233 not share any nucleotide sequence similarity to known viruses; however, we note that because
234 the enrichment of viral nucleic acids occurs after sequence library construction, the uncaptured

235 portion of the library could subsequently be sequenced for additional attempts at pathogen
236 discovery. Furthermore, the ViroCap panel is extensible and will be updated periodically with
237 new viral sequences as they are added to RefSeq and the Genome Neighbors databases.
238 Updates will be publicly available through our GitHub repository (see Methods: Data access).

239
240 There were a few genomes (<10) in the NCBI reference databases that had been cloned into
241 bacterial vectors prior to sequencing, and the deposited viral genome sequences contained
242 bacterial vector sequence. We were not aware of this prior to probe design, so ViroCap includes
243 capture probes that target these sequences. This resulted in enrichment of some sequences (on
244 average 1.1% of total non-viral reads) that were subsequently recognized by our analysis
245 pipeline as bacterial based on nucleotide sequence alignment. In subsequent versions of
246 ViroCap, we will filter out these bacterial vector sequences.

247
248 In the experiments reported here, we pooled sequencing libraries prior to targeted sequence
249 capture in order to reduce cost, but still achieved enhanced detection of multiple viruses of
250 varying abundance. As has been reported for strategies that involve sequencing indexed,
251 pooled libraries (Kircher et al. 2012), we observed some sample cross-contamination. This
252 cross-contamination is recognizable when a high number of viral sequences are detected in the
253 truly positive sample, while few sequences ($\leq 0.05\%$ of the viral sequences in the truly positive
254 sample) of the same virus are detected in other samples in the pool. In a clinical setting, each
255 sample would optimally be captured and sequenced independently to reduce the possibility of
256 sample cross-contamination. However, future methodological improvements could allow pooling
257 of clinical specimens.

258
259 The success of viral targeted sequence capture is affected by the representation of the virus in
260 the sequencing library. In our sample preparation, total nucleic acid extracted from the sample

261 was reverse transcribed and randomly amplified prior to library construction (Wang et al. 2003),
262 allowing detection of DNA and RNA viral genomes within the same sequencing experiment. The
263 uneven sequence representation observed for some genomes (**Fig. 2**) is likely due in part to
264 detection of messenger RNA, whose abundance reflects patterns of gene expression, as well as
265 to primer biases during the reverse transcription and amplification steps. Capture hybridization
266 may also induce bias, in that sequences that diverge from target probe sequences may be
267 captured less efficiently than those with exact or nearly exact matches to the probe. Taken
268 together, these data suggest that further improvement in performance of viral targeted
269 sequence capture may be achievable by improving efficiency of reverse transcription,
270 amplification, and library construction, while continuing to update the ViroCap panel as new,
271 divergent genome sequences become available.

272
273 Methods other than genome sequencing have been used for virus characterization and
274 discovery, including Virochip, a microarray-based method for detection/genotyping of viral
275 pathogens (Chen et al. 2011; Wang et al. 2002), and PathoChip, a microarray designed to
276 detect viruses and other microbial pathogens (Baldwin DA et al. 2014). While designed to detect
277 known viruses by means of microarray probe spotting, this technology has also shown the
278 ability to detect emerging viruses (Yu et al. 2012; Wang et al. 2003). The primary difference
279 between the designs of these microarrays and ViroCap targeted sequence capture is that the
280 latter approach targets complete viral genomes while the microarrays target smaller, discrete
281 genomic regions. The results obtained from each approach also differ significantly. The
282 microarray approach detects the presence of a virus but does not directly provide sequence
283 information. In contrast MSS enhanced by ViroCap targeted sequence capture provides
284 sequence data, sometimes covering the entire genome.

285

286 In conclusion, ViroCap greatly enhances the sensitivity of MSS for nucleotide sequence-based
287 virus detection. To our knowledge, ViroCap represents the first effort to apply a targeted
288 sequence capture approach to the detection of a comprehensive set of viruses. Its research
289 applications are far reaching, allowing a new, higher resolution view of eukaryotic DNA and
290 RNA viruses in the microbiome. ViroCap should also help realize the potential of MSS as a
291 clinical diagnostic tool that can simultaneously detect viruses as well as provide immediate
292 characterization including taxonomic assignment, strain typing, virulence characteristics, and
293 anti-viral drug resistance genotyping. ViroCap could also be modified into a tool for broader
294 pathogen identification, which might include a comprehensive set of human pathogens: genes
295 from viruses, bacteria (e.g. toxin genes, antibiotic resistance genes), fungi, protists, and other
296 microbes.

297

298 **METHODS**

299 **Taxonomy selection.** At the time we designed the ViroCap panel, NCBI GenBank had
300 available for download a total of ~1 Gb of sequence representing 440 viral families, well beyond
301 the 200 megabase (Mb) of target space supported by the custom SeqCap EZ library format
302 (NimbleGen, Madison, WI, USA). Therefore, we developed the following approach for selecting
303 representative targeted sequence capture probes. Because we were interested in studying viral
304 diseases of humans, we excluded bacteriophages and endogenous human retroviruses. We
305 also specifically did not include references from the following NCBI viral reference genome
306 database host categories: algae, archaea, bacteria, diatom, environment, fungi, invertebrates,
307 plants, and protozoa. After filtering, our target list contained reference sequences from the
308 following host categories: human, vertebrates, and “unknown”. This list included viruses that
309 could have both vertebrate and invertebrate hosts, such as vertebrate viruses with insect
310 vectors. Based on these broad viral-host categories, we downloaded all of the associated viral
311 reference sequences in each chosen category from NCBI (accessed February 3, 2014). These

312 sequences comprise the core reference database from which our capture library is designed.
313 Our capture library includes targets from 34 viral families composed of 190 annotated viral
314 genera and 337 species (**Fig. 1, Supplemental Tables S6 and 7**). Sources of viral sequences
315 include complete representation of the viral genomes from NCBI's Reference Sequence
316 (RefSeq) collection, complementary representation of unique regions from Genome Neighbor
317 targets, selected representation of NCBI Influenza Virus Resource sequences, and the entirety
318 of the probe space represented on the Virochip microarray (Yu G et al. 2012), GEO accession
319 number GPL15905. The methods used to consolidate these database sequences follow.

320
321 **RefSeq.** NCBI's RefSeq (<http://www.ncbi.nlm.nih.gov/refseq/>) genome collection is a database
322 of taxonomically diverse entries representing comprehensive, well-annotated genome
323 sequences (Pruitt et al. 2014; Tatusova et al. 2014). As RefSeq entries are the most complete
324 sequence representatives in terms of annotation and metadata consistency, we targeted
325 selected viral RefSeqs by tiling of targeted sequence capture probes across the entire length of
326 each RefSeq's genome, with the intention of capturing the entire viral genome. For our capture
327 library, RefSeq nucleotide FASTA sequences were downloaded for desired viral-host categories
328 (human; vertebrates; vertebrates, human; vertebrates, invertebrates; vertebrates, invertebrates,
329 human; invertebrates, vertebrates; unknown) using both the online NCBI taxonomy viewer
330 (<http://www.ncbi.nlm.nih.gov/genomes/GenomesGroup.cgi?opt=virus&taxid=10239>) as well as
331 the RefSeq specific FTP site (<ftp://ftp.ncbi.nlm.nih.gov/refseq/release/viral>). Entries were
332 merged to avoid redundancy. RefSeq targets were pooled with the other sequence candidates
333 (see: **Design consolidation**). A total of 1,456 RefSeq FASTA entries (26.9 Mb) representing
334 190 viral genera were completely tiled for inclusion in the ViroCap library, accounting for 13.5%
335 of the total capture library's target-space.

336

337 **Genome Neighbors.** While RefSeq entries are single, canonical species representations, other
338 complete or partial viral sequences also exist in DDBJ/EMBL/GenBank. In the case of viral
339 sequences, there is extensive redundancy in these databases due to the large number of similar
340 viral strains, isolates, and mutants. Therefore, non-RefSeq (e.g. DDBJ, EMBL, GenBank)
341 nucleotide sequences of complete viral genomes that belong to the same species as a RefSeq
342 are classified as Genome Neighbors for that reference sequence, provided that they match all of
343 the criteria that were used to select complete genomic sequences (Bao et al. 2004). At the time
344 of our ViroCap panel design, Genome Neighbors (sequences downloaded from Entrez Genome
345 link "Other genomes for species"; accessed February 3, 2014) in total represented an additional
346 56,314 entries and 507.1 Mb of sequence, more than 2.5 times our SeqCap EZ capture target
347 sequence space limit. Therefore, an alternative target selection approach was employed to add
348 diversity to our RefSeq selections by selecting unique, complementary Genome Neighbor
349 sequences.

350

351 **RefSeq and Genome Neighbor sequence association.** We began the process of variant
352 sequence selection by identifying conserved regions in Genome Neighbors already represented
353 by completely tiled RefSeq capture probes. First, we associated our viral RefSeq selections with
354 corresponding Genome Neighbors. This was performed by downloading Genome Neighbor
355 annotation files from NCBI
356 (<http://www.ncbi.nlm.nih.gov/genomes/GenomesGroup.cgi?opt=virus&taxid=10239>) and
357 associating the information with our RefSeq annotation files, by means of *ad hoc* Perl
358 parsing/coupling scripts (results in **Supplemental Tables S6 and S7**). Once associated, the
359 parent RefSeq sequences could be compared to related Genome Neighbor sequences to
360 determine conserved and divergent nucleotide regions. Each viral RefSeq entry was individually
361 reviewed, along with associated Genome Neighbor entries. FASTA sequences were collected
362 for each RefSeq entry and its related Genome Neighbors for subsequent *k*-mer analysis.

363

364 **K-mer analysis.** Each of the Genome Neighbor sequences was split into 100 bp *k*-mers by
365 means of an exhaustive 1 bp sliding window algorithm, as depicted in **Supplemental Figure**
366 **S4**. The resultant output thus included all possible 100 bp sequences based on the combined
367 Genome Neighbor sequence space. As our SeqCap EZ targeted sequence capture probe
368 lengths are 100 bp, the sequences generated by the sliding window algorithm represent the
369 total number of possible probe combinations based on the aggregate of Genome Neighbor
370 sequences. Based on our conservative expectation of hybridization/homology at the capture
371 probe level, we then clustered all of the Genome Neighbor 100-mers back to the parent RefSeq
372 sequence at $\geq 90\%$ sequence identity using length-sorted FASTA entries and the UCLUST
373 (Edgar 2010) package (version 1.1.579; parameters: --rev --id .90). Given that all of our
374 candidate sequences were 100 bp in length, and all RefSeq entries are >100 bp, UCLUST
375 always used the longer RefSeq as the first seed (*centroid*) in which to attempt folding of other
376 sequences. As the parent RefSeq had complete probe tiling in our design, any Genome
377 Neighbor 100-mer with $\geq 90\%$ identity was considered already represented in our capture
378 library, and therefore discarded. Genome Neighbor 100-mers with less than 90% identity were
379 chosen for inclusion in the capture library. As the sliding window approach produces 100-mers
380 that overlapped one another, we merged overlapping 100-mers based on their Genome
381 Neighbor genomic coordinates into single contiguous spans using BEDTools (Quinlan and Hall
382 2010) functions.

383

384 **Genome Neighbor sub-sequences.** Resultant sub-sequences were excised as FASTA entries
385 from corresponding Genome Neighbor references using WU-BLAST's (<http://blast.wustl.edu>)
386 xdget application and added to the ViroCap panel. These supplementary entries are easily
387 identifiable in our final target design, as the FASTA headers for the entries list the original
388 parent sequence ID with the excised span indicated in curly braces (e.g.

389 gi|1249624|emb|A28090.1| HPV42 (partial) genomic sequence {SQ 2444-2644}). In this
390 manner, for each RefSeq species, we generated Genome Neighbor sub-sequences from 100
391 bp to 21 kb in length to add to our capture panel.

392

393 These processing steps reduced the aggregate input Genome Neighbors targeted sequence
394 space from 507.1 Mb to 153.2 Mb (**Supplemental Table S1**), and these sequences were
395 pooled with our other targeted capture sequence targets (see: **Design consolidation**). A total
396 of 130,808 partial Genome Neighbor FASTA entries (153.2 Mb) were added for capture in our
397 ViroCap library, accounting for 77.1% of the total capture library's target-space.

398

399 **Influenza Virus Resource.** We obtained reference sequences from NCBI's Influenza Virus
400 Resource database (<http://www.ncbi.nlm.nih.gov/genomes/FLU/FLU.html>), which contains
401 sequence data from the NIAID Influenza Genome Sequencing Project as well as from GenBank.
402 At the time of our capture panel design, the NCBI Influenza Virus Resource contained 305,524
403 influenza entries, representing 458.1 Mb of sequence. This is 17 times the size of our viral
404 RefSeq selections and 3 times the size of our collapsed Genome Neighbor targets. Our
405 selected RefSeqs included 29 influenza RefSeq entries (each influenza virus segment is
406 represented as a separate entry), targeted in its entirety. These sequences served as the core
407 of influenza reference genomes against which all other influenza sequences were compared.
408 We directly clustered the long influenza sequences using length-sorted FASTA and the
409 UCLUST package (version 1.1.579; parameters: --rev --id .90). In UCLUST, a cluster is defined
410 by one sequence, known as the *centroid* or representative sequence. To lessen the
411 computational burden and ensure that our core influenza RefSeq genomes were always the
412 longest first seeds (*centroids*) in UCLUST's clustering process, we artificially concatenated the
413 29 parent RefSeq sequences into one linear sequence representation, and then split this
414 representation into 6 segments ranging in size from 18-26 kb. UCLUST preferentially seeded

415 with the long RefSeq construct segments when clustering, ensuring that clustering was first
416 attempted within the longer, canonical references. ULCUST was run with a requirement of
417 $\geq 90\%$ sequence identity to fold into a parent influenza RefSeq segment. Therefore, only
418 sequences that (1) had $< 90\%$ identity to influenza RefSeqs and (2) were subsequent centroids
419 in non-RefSeq clusters were chosen for inclusion in our capture panel. This process reduced
420 the aggregate input Influenza Resource Database reference sequence from 458.1 Mb to 15.7
421 Mb (**Supplemental Table S1**). Finally, supplementary influenza targets were pooled with the
422 other sequence candidates (see: **Design consolidation**). A total of 9,759 influenza FASTA
423 entries (15.7 Mb) were added for targeted sequence capture in our ViroCap library, accounting
424 for 7.9% of the total capture library's target-space.

425

426 **Virochip microarray.** Considering the biologically important short sequence signatures
427 represented on the Virochip panel (Yu G et al. 2012), as well as the comparatively small
428 footprint, we subsumed these sequences within our targeted sequence capture panel design.
429 The probe sequences for the microarray are publicly available at NCBI's Gene Expression
430 Omnibus (GEO) repository (Edgar et al. 2002). We downloaded this information for Platform
431 GPL15905 (Viro5AG-60k) as a text file
432 (<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GPL15905>). This platform included more
433 than 60,000 oligonucleotides of length 60-70 bp, corresponding to 3.1 Mb of probes
434 (**Supplemental Table S1**). Virochip targets were pooled with the other sequence candidates
435 (see: **Design consolidation**). Upon review, 1.3 Mb of the probes were already directly
436 represented by RefSeq, Genome Neighbor, and Influenza Viral Resource sequences during
437 capture library design and synthesis. Therefore, the remaining 25,749 (60-70 bp) Virochip
438 FASTA entries of 1.8 Mb total size were added to ViroCap, accounting for $< 1\%$ of the total
439 targeted sequence capture panel.

440

441 **Design consolidation.** All of our selected candidate target sequences from RefSeq, Genome
442 Neighbors, Influenza Virus Resource, and the Virochip microarray were combined into a single
443 FASTA sequence file. Human endogenous retroviruses were removed from inclusion by means
444 of a two-part process: 1) entries were filtered by taxonomic annotation indicating human
445 endogenous retrovirus identity; 2) the remaining entries were BLAST-aligned to the GRCh37-lite
446 build of the human reference genome to remove sequences with high percent identity ($\geq 75\%$)
447 at the 100 bp probe level. Finally, sequences were hard-masked (i.e. bases converted to N's) in
448 low complexity regions using the DUST (R. Tatusov and D.J. Lipman, unpublished) software
449 module. The final ViroCap targeted sequence capture panel consists of 185,835 FASTA
450 sequences totaling 198.9 Mb (see: **Taxonomy selection**).

451
452 **NimbleGen Sequence Capture design.** Our consolidated target sequences were submitted to
453 Roche NimbleGen (Madison, WI, USA) on March 31, 2014 for capture library design and
454 synthesis. As our final ViroCap design required 198.9 Mb, manufacturing was implemented
455 under the custom NimbleGen SeqCap EZ Developer Library format, which has a maximum
456 capture space of 200 Mb of non-human sequence. NimbleGen's Sequence Capture design
457 offered up to 2.1 million of 50-105mer sequence probes. It was at the discretion of NimbleGen,
458 based on proprietary algorithms, to redistribute probes for better capture uniformity,
459 redundancy, and comprehensive target base coverage. NimbleGen provided us with a proposed
460 capture design accompanied by coordinate (GFF, BED) files and associated sequence
461 coverage metrics on April 14, 2014. The design set contained probe representation generated
462 by first masking all but one exact copy of each 100-mer in our original FASTA file, tiling the
463 unmasked regions, screening the resulting probes against the (hg19) human genome, and
464 finally selecting only those probes that had no matches in the human genome as determined by
465 the SSAHA (Ning et al. 2001) algorithm. NimbleGen provides two metrics for assessing in silico
466 targeted sequence capture design coverage: 1) target bases covered with 0-bp-offset are

467 determined by counting target bases directly represented in probe sequences; 2) target bases
468 covered with 100-bp-offset are determined by counting all target bases within 100 bp of a probe.
469 The capture design provided 95.9% 0-bp-offset coverage and 99.6% 100-bp-offset coverage of
470 our initial 198.9 Mb target request. We approved the design on April 17, 2014 for capture library
471 synthesis and received our first 12 SeqCap EZ Library reactions for in-house Illumina
472 sequencing and analysis on April 28, 2014.

473

474 **Human subjects approval and sample selection.** Samples were collected under protocols
475 approved by the Human Research Protection Office at Washington University School of
476 Medicine (IRB protocol numbers 201106177, 201102561, and 201102045). Samples were
477 selected to represent a broad range of viruses that are commonly encountered in the clinical
478 laboratory and in our research studies. Viruses were identified in samples based on clinical
479 laboratory test results in the Diagnostic Virology Laboratory at St. Louis Children's Hospital or
480 by PCR assays and sequencing results carried out in previous studies (Colvin et al. 2012;
481 McElvania TeKippe et al. 2012; Wylie et al. 2012). Specimens of nasopharyngeal secretions,
482 plasma, and stool were included.

483

484 **Sequencing.** Total nucleic acid was extracted from clinical samples using the EasyMag
485 NucliSENS instrument (bioMerieux, Marcy l'Etoile, France). Samples were processed in one of
486 two ways. In experiment 1, nucleic acids from clinical specimens from the Diagnostic Virology
487 Laboratory were combined, and the resulting pooled nucleic acid was used as input for a single
488 sequencing library (constructed as described below). These samples are designated with a
489 sample identification prefix of "P" in the various tables and figures. Alternatively, in experiment
490 2, individual sequencing libraries were made from each set of 8 different specimens prior to
491 combining the libraries for sequencing. These samples are designated with a sample
492 identification prefix of "S" in the various tables.

493

494 For sequencing libraries, DNA and RNA viruses were assessed in the same assay as described
495 (Wang et al. 2003). Specifically, the RNA in the total nucleic acid was reverse transcribed with
496 reverse transcriptase (Promega, Fitchburg, Wisconsin) and random nonomers tagged with a
497 conserved sequence (5' **GTTTCCCAGTCACGATA** 3') to be used for subsequent amplification
498 (Integrated DNA Technologies, Coralville, Iowa), and second strand synthesis was carried out
499 with Sequenase V2.0 DNA polymerase (Affymetrix, Santa Clara, California). DNA and RNA
500 were subsequently amplified with Accuprime Taq (Life Technologies, Grand Island, New York)
501 using the conserved sequence on the ends of the random primers, and the DNA/cDNA mixture
502 was sheared using the Qsonica Q800R instrument (Qsonica, Newtown, Connecticut) to
503 generate fragments with an average length of 500 basepairs. Dual-indexed sequencing libraries
504 were constructed using the KAPA Low Throughput Library Construction Kit (KAPA Biosystems,
505 Wilmington, Massachusetts).

506

507 For the anellovirus samples, DNA was amplified with the Illustra GenomiPhi V2 DNA
508 Amplification Kit (GE Healthcare Life Sciences, Pittsburgh, Pennsylvania); RNA was not
509 assessed. DNA was sheared and libraries were constructed from each sample as described
510 above. Sequencing libraries were pooled (**Tables 2 and 3, column "Sample ID"**).

511

512 In each case, the libraries were divided, and part was directly sequenced (pre-capture) and part
513 was subjected to targeted sequence capture with the custom ViroCap probes prior to
514 sequencing (post-capture). Targeted sequence capture was carried out according to the
515 manufacturer's specifications. We carried out 10, 10, and 16 cycles of post-capture linker-
516 mediated PCR for experiments 1 (pooled clinical samples), 2 (individual samples from the
517 research study) and 3 (anellovirus samples), respectively, prior to sequencing. The number of
518 cycles was empirically determined to be the minimum number needed to obtain a 5 nanomolar

519 dilution of library material for qPCR and loading. Libraries were sequenced on the Illumina
520 HiSeq 2000 or HiSeq 2500 instrument, generating 100 bp, paired-end reads.

521
522 **Sequence analysis.** Viral sequences were identified based on nucleotide and translated protein
523 sequence alignment against reference genomes. The pipeline is adapted from previously
524 described methods (Wylie et al. 2014), except that nucleotide alignments were carried out using
525 BWA-MEM with default settings (Li and Durbin 2009). Because many similar genomes are
526 included in the reference database, we used the initial alignment statistics for each sample to
527 choose a single reference from each species to calculate and report coverage statistics.
528 References were chosen based on having the highest number of reference bases covered.
529 Sequences were re-aligned to the selected references with BWA-MEM for calculation of
530 coverage statistics and comparison of samples pre- and post-capture. Sequence alignments
531 were evaluated with samtools (Li 2011), and sequence coverage was determined with RefCov
532 (<http://gmt.genome.wustl.edu/packages/refcov/>) and visualized with Plot2
533 (<http://plot2doc.micw.eu>). For illustrative purposes, the genome coverage panels in **Figures 2**
534 **and 3** were normalized by removing (*deduplicating*) reads based on identical alignment start-
535 sites using the samtools *rmdup* command. For each alignment start-site, only the highest quality
536 read was retained for forward and reverse alignment orientations. Therefore, for the 100 bp read
537 data shown in each coverage panel the theoretical maximum depth is 200X.

538
539 Anellovirus contigs were assembled from the pre-capture sequence data using IDBA-UD (Peng
540 et al. 2012). Contigs were aligned against the sequence database used to design the ViroCap
541 panel using BLAST (Altschul et al. 1997) with the following parameters to detect low-similarity
542 sequences: -G 5 -E 2 -r 1 -q -1. The percent identity of the highest scoring HSP is reported in
543 **Supplemental Table S3.**

544

545 **DATA ACCESS**

546 Data files associated with our ViroCap panel are publicly available through the
547 <https://github.com/WashU-PMG/ViroCap> GitHub repository. Hosted files include 198.9 Mb of
548 ViroCap target sequence in FASTA format, taxonomy information, corresponding RefSeq and
549 Genome Neighbor associations, and NimbleGen's target design coverage metrics for 0- and
550 100-bp-offset intervals. MSS data used for ViroCap evaluation have been deposited (with
551 potentially identifiable human sequences removed) in the NCBI Sequence Read Archive
552 (PRJNA273884).

553

554 **ACKNOWLEDGMENTS**

555 We thank Richard Hotchkiss and Andrew Walton for providing samples for the anellovirus work;
556 Maria Cannella and Richard Buller for carrying out the virus-specific PCR assays on the
557 samples from the research study; Stephanie Bledsoe and the staff of the Virology Laboratory at
558 St. Louis Children's Hospitals for supplying remnant specimen material; Elena Deych and
559 William Shannon for assisting with statistical analysis; and The McDonnell Genome Institute for
560 generating sequence data. We thank Efreem Lim and David Wang for performing the PCR
561 validation test on the HPyV10 virus.

562

563 **AUTHOR CONTRIBUTIONS**

564 TNW and KMW developed and benchmarked the methodology. TNW developed computational
565 tools for targeted sequence capture selection. BNH and KMW generated MSS libraries,
566 optimized molecular methods, and validated results. TNW, KMW, and GAS performed the MSS
567 analysis and virus review. KMW and GAS were involved with study design. GAS provided
568 patient samples representing a broad range of viruses that are commonly encountered in the
569 clinical laboratory and provided the samples from the research study. TNW, KMW, and GAS
570 wrote the manuscript.

571

572 **FUNDING**

573 This study was supported by grant number R01AI097213 from the National Institute of Allergy
574 and Infectious Diseases, awarded to GAS. The funders had no role in study design, data
575 collection and analysis, decision to publish, or preparation of the manuscript. The authors have
576 no conflicts of interest to report.

577

578 **FIGURE LEGENDS**

579 **Figure 1. Taxonomic distribution of target genomes included in ViroCap.** Shown are the
580 viral classes, families, and genera included in the ViroCap targeted sequence capture panel.
581 Taxonomic assignments were obtained from the NCBI Taxonomy Viewer
582 (<http://www.ncbi.nlm.nih.gov/genomes/GenomesGroup.cgi?opt=virus&taxid=10239>).

583

584 **Figure 2. Targeted sequence capture enrichment.** Examples are given showing the impact of
585 targeted sequence capture on breadth and depth of genome coverage for 8 representative viral
586 genomes. For illustrative purposes, all of the coverage panels in this figure have been
587 normalized by removing (*deduplicating*) reads based on identical alignment start-sites.
588 Nucleotide positions along the reference genome are shown on the X-axis. The depth of
589 deduplicated reads is shown on the Y-axis. The shaded portion indicates the sequence
590 coverage (breadth and depth) for each virus. Post-capture sequence coverage is represented in
591 the larger panels in blue; pre-capture sequence coverage is shown in the insets in red. Note that
592 Y-axis ranges are different for each panel. At the top of each panel is shown the breadth of
593 coverage (BoC) for the sample. The header of each panel includes breadth of coverage gain
594 (BoC gain), sample id, and reference genome name and NCBI version number. BoC gain is
595 calculated by subtracting the percent of the length of the reference genome that was covered by

596 sequence reads in pre-capture MSS from the percent of the length of the reference genome
597 covered by post-capture sequence reads.

598
599 **Figure 3. Targeted sequence capture identifies divergent sequences.** (A) The percent
600 identity of the top HSP identified from the BLAST alignment of anellovirus contig sequences to
601 the references used to design ViroCap is plotted on the Y-axis. The X-axis represents the
602 percent of the length of the anellovirus contig covered after targeted sequence capture. (B) This
603 coverage plot represents the sequence coverage of a divergent anellovirus contig sequence.
604 The figure is designed as described in the figure legend for **Figure 2**, with the following addition:
605 the post-capture coverage plot is shaded to show regions of nucleotide sequence variation
606 between the anellovirus contig and the most similar reference genome in the ViroCap panel.
607 Dark shading represents areas of identical sequence, and each position with nucleotide
608 mismatch between aligned sequences is shown in the lighter color. All of the HSPs are shown,
609 rather than just the top HSP.

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624 **REFERENCES**

625 Albert TJ, Molla MN, Muzny DM, Nazareth L, Wheeler D, Song X, Richmond TA, Middle CM,
626 Rodesch MJ, Packard CJ, et al. 2007. Direct selection of human genomic loci by microarray
627 hybridization. *Nat Methods* **4**: 903–905.

628 Allander T, Emerson SU, Engle RE, Purcell RH, Bukh J. 2001. A virus discovery method
629 incorporating DNase treatment and its application to the identification of two bovine
630 parvovirus species. *Proc Natl Acad Sci USA* **98**: 11609–11614.

631 Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped
632 BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic
633 Acids Res* **25**: 3389–3402.

634 Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, Fernandes GR, Tap J,
635 Bruls T, Batto J-M, et al. 2011. Enterotypes of the human gut microbiome. *Nature* **473**: 174–
636 180.

637 Baldwin DA, Feldman M, Alwine JC, Robertson ES. 2014. Metagenomic assay for
638 identification of microbial pathogens in tumor tissues. *MBio* **5**: e01714–14.

639 Bao Y, Federhen S, Leipe D, Pham V, Resenchuk S, Rozanov M, Tatusov R, Tatusova T. 2004.
640 National center for biotechnology information viral genomes project. *J Virol* **78**: 7291–7298.

641 Breitbart M, Rohwer F. 2005. Method for discovering novel DNA viruses in blood using viral
642 particle selection and shotgun sequencing. *BioTechniques* **39**: 729–736.

643 Chen EC, Miller SA, DeRisi JL, Chiu CY. 2011. Using a pan-viral microarray assay
644 (Virochip) to screen clinical samples for viral pathogens. *J Vis Exp* e2536–e2536.

645 Chiu CY. 2013. Viral pathogen discovery. *Curr Opin Microbiol* **16**: 468–478.

646 Cleland EJ, Bassioni A, Boase S, Dowd S, Vreugde S, Wormald P-J. 2014. The fungal
647 microbiome in chronic rhinosinusitis: richness, diversity, postoperative changes and patient
648 outcomes. *Int Forum Allergy Rhinol* **4**: 259–265.

649 Colvin JM, Muenzer JT, Jaffe DM, Smason A, Deych E, Shannon WD, Arens MQ, Buller RS,
650 Lee WM, Weinstock EJS, et al. 2012. Detection of viruses in young children with fever
651 without an apparent source. *Pediatrics* **130**: e1455–62.

652 de Villiers E-M, Borkosky SS, Kimmel R, Gunst K, Fei J-W. 2011. The diversity of torque teno
653 viruses: in vitro replication leads to the formation of additional replication-competent subviral
654 molecules. *J Virol* **85**: 7284–7295.

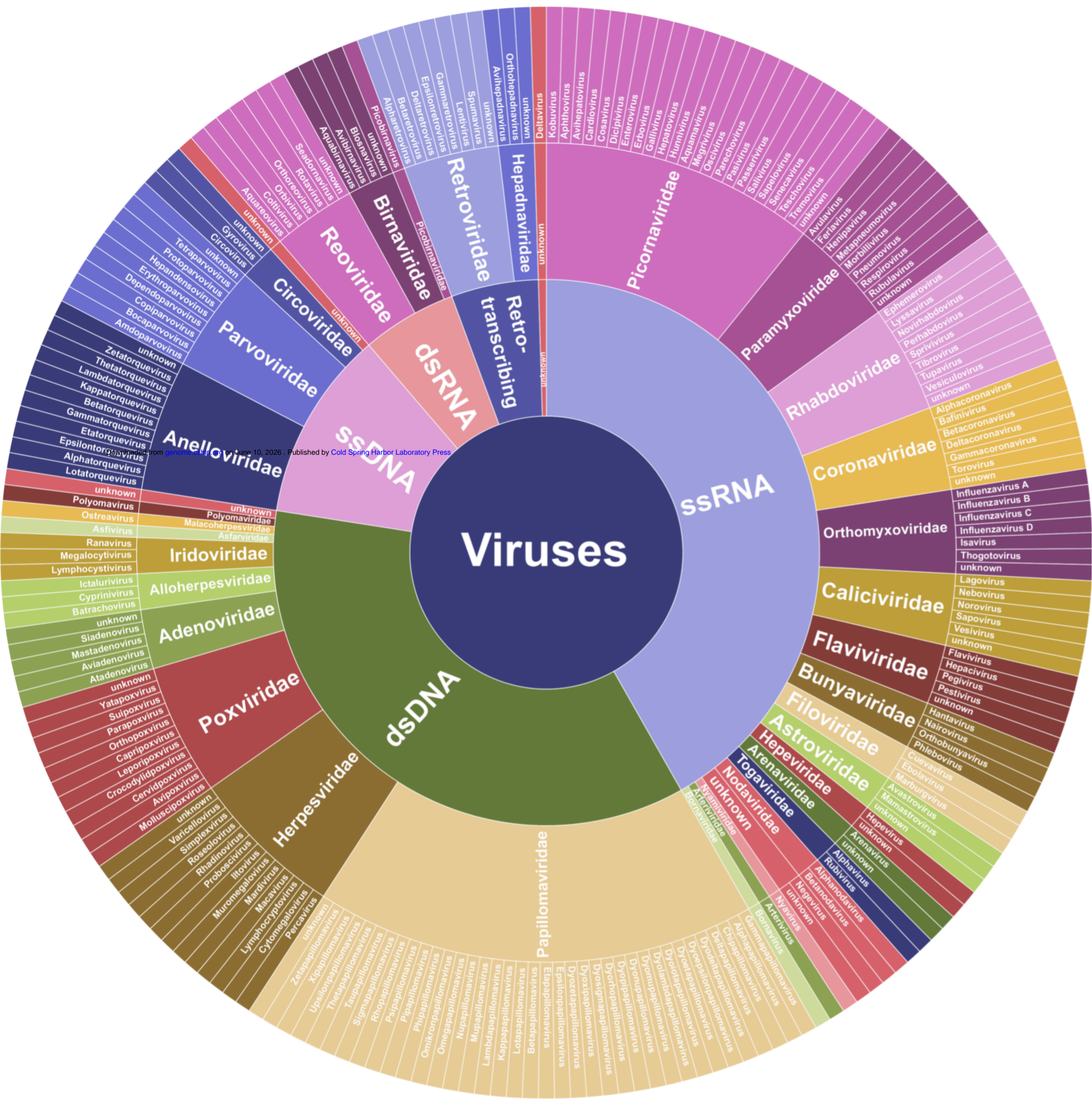
655 De Vlaminc I, Khush KK, Strehl C, Kohli B, Luikart H, Neff NF, Okamoto J, Snyder TM,
656 Cornfield DN, Nicolls MR, et al. 2013. Temporal response of the human virome to
657 immunosuppression and antiviral therapy. *Cell* **155**: 1178–1187.

- 658 Depledge DP, Palser AL, Watson SJ, Lai IY-C, Gray ER, Grant P, Kanda RK, Leproust E,
659 Kellam P, Breuer J. 2011. Specific capture and whole-genome sequencing of viruses from
660 clinical samples. ed. R. Jhaveri. *PLoS ONE* **6**: e27805.
- 661 Duhaime MB, Sullivan MB. 2012. Ocean viruses: rigorously evaluating the metagenomic
662 sample-to-sequence pipeline. *Virology* **434**: 181–186.
- 663 Duncavage EJ, Magrini V, Becker N, Armstrong JR, Demeter RT, Wylie T, Abel HJ, Pfeifer JD.
664 2011. Hybrid capture and next-generation sequencing identify viral integration sites from
665 formalin-fixed, paraffin-embedded tissue. *J Mol Diagn* **13**: 325–333.
- 666 Edgar R, Domrachev M, Lash AE. 2002. Gene Expression Omnibus: NCBI gene expression
667 and hybridization array data repository. *Nucleic Acids Res* **30**: 207–210.
- 668 Edgar RC. 2010. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*
669 **26**: 2460–2461.
- 670 Findley K, Oh J, Yang J, Conlan S, Deming C, Meyer JA, Schoenfeld D, Nomicos E,
671 Park M, NIH Intramural Sequencing Center Comparative Sequencing Program, et
672 al. 2013. Topographic diversity of fungal and bacterial communities in human skin.
673 *Nature* **498**: 367–370.
- 674 Gajer P, Brotman RM, Bai G, Sakamoto J, Schütte UME, Zhong X, Koenig SSK, Fu L, Ma ZS,
675 Zhou X, et al. 2012. Temporal dynamics of the human vaginal microbiota. *Sci Transl Med* **4**:
676 132ra52–132ra52.
- 677 Hodges E, Xuan Z, Balija V, Kramer M, Molla MN, Smith SW, Middle CM, Rodesch MJ, Albert
678 TJ, Hannon GJ, et al. 2007. Genome-wide in situ exon capture for selective resequencing.
679 *Nat Genet* **39**: 1522–1527.
- 680 Holtz LR, Cao S, Zhao G, Bauer IK, Denno DM, Klein EJ, Antonio M, Stine OC, Snelling TL,
681 Kirkwood CD, et al. 2014. Geographic variation in the eukaryotic virome of human diarrhea.
682 *Virology* **468-470**: 556–564.
- 683 Human Microbiome Project Consortium. 2012. Structure, function and diversity of the healthy
684 human microbiome. *Nature* **486**: 207–214.
- 685 Kircher M, Sawyer S, Meyer M. 2012. Double indexing overcomes inaccuracies in multiplex
686 sequencing on the Illumina platform. *Nucleic Acids Res* **40**: e3–e3.
- 687 Koehler JW, Hall AT, Rolfe PA, Honko AN, Palacios GF, Fair JN, Muyembe J-J,
688 Mulembekani P, Schoepp RJ, Adesokan A, et al. 2014. Development and evaluation
689 of a panel of filovirus sequence capture probes for pathogen detection by next-
690 generation sequencing. ed. J.H. Kuhn. *PLoS ONE* **9**: e107007.
- 691
- 692 Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and
693 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**: 2987–
694 2993.

- 695 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
696 *Bioinformatics* **25**: 1754–1760.
- 697 Lovett M, Kere J, Hinton LM. 1991. Direct selection: a method for the isolation of cDNAs
698 encoded by large genomic regions. *Proc Natl Acad Sci USA* **88**: 9628–9632.
- 699 Lysholm F, Wetterbom A, Lindau C, Darban H, Bjerkner A, Fahlander K, Lindberg AM, Persson
700 B, Allander T, Andersson B. 2012. Characterization of the viral microbiome in patients with
701 severe lower respiratory tract infections, using metagenomic sequencing. ed. S.K.
702 Highlander. *PLoS ONE* **7**: e30875.
- 703 McElvania TeKippe E, Wylie KM, Deych E, Sodergren E, Weinstock G, Storch GA. 2012.
704 Increased prevalence of anellovirus in pediatric patients with fever. ed. R. Jhaveri. *PLoS*
705 *ONE* **7**: e50937.
- 706 Minot S, Sinha R, Chen J, Li H, Keilbaugh SA, Wu GD, Lewis JD, Bushman FD. 2011. The
707 human gut virome: inter-individual variation and dynamic response to diet. *Genome Res* **21**:
708 1616–1625.
- 709 Ning Z, Cox AJ, Mullikin JC. 2001. SSAHA: a fast search method for large DNA databases.
710 *Genome Res* **11**: 1725–1729.
- 711 Ninomiya M, Takahashi M, Shimosegawa T, Okamoto H. 2007. Analysis of the entire genomes
712 of fifteen torque teno midi virus variants classifiable into a third group of genus Anellovirus.
713 *Arch Virol* **152**: 1961–1975.
- 714 Oh J, Byrd AL, Deming C, Conlan S, NISC Comparative Sequencing Program, Kong HH, Segre
715 JA. 2014. Biogeography and individuality shape function in the human skin metagenome.
716 *Nature* **514**: 59–64.
- 717 Okou DT, Steinberg KM, Middle C, Cutler DJ, Albert TJ, Zwick ME. 2007. Microarray-based
718 genomic selection for high-throughput resequencing. *Nat Methods* **4**: 907–909.
- 719 Paulino LC, Tseng C-H, Strober BE, Blaser MJ. 2006. Molecular analysis of fungal microbiota in
720 samples from healthy human skin and psoriatic lesions. *J Clin Microbiol* **44**: 2933–2941.
- 721 Peng Y, Leung HCM, Yiu SM, Chin FYL. 2012. IDBA-UD: a de novo assembler for single-cell
722 and metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**: 1420–
723 1428.
- 724 Pruitt KD, Brown GR, Hiatt SM, Thibaud-Nissen F, Astashyn A, Ermolaeva O, Farrell CM, Hart
725 J, Landrum MJ, McGarvey KM, et al. 2014. RefSeq: an update on mammalian reference
726 sequences. *Nucleic Acids Res* **42**: D756–63.
- 727 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
728 features. *Bioinformatics* **26**: 841–842.
- 729 Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, Gordon JI. 2010. Viruses in
730 the faecal microbiota of monozygotic twins and their mothers. *Nature* **466**: 334–338.
- 731 Tatusova T, Ciufu S, Fedorov B, O'Neill K, Tolstoy I. 2014. RefSeq microbial genomes

- 732 database: new representation and annotation strategy. *Nucleic Acids Res* **42**: D553–9.
- 733 Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, Sogin ML, Jones
734 WJ, Roe BA, Affourtit JP, et al. 2009. A core gut microbiome in obese and lean twins.
735 *Nature* **457**: 480–484.
- 736 Wang D, Coscoy L, Zylberberg M, Avila PC, Boushey HA, Ganem D, DeRisi JL. 2002.
737 Microarray-based detection and genotyping of viral pathogens. *Proc Natl Acad Sci USA* **99**:
738 15687–15692.
- 739 Wang D, Urisman A, Liu Y-T, Springer M, Ksiazek TG, Erdman DD, Mardis ER, Hickenbotham
740 M, Magrini V, Eldred J, et al. 2003. Viral discovery and sequence recovery using DNA
741 microarrays. *PLoS Biol* **1**: E2.
- 742 Willger SD, Grim SL, Dolben EL, Shipunova A, Hampton TH, Morrison HG, Filkins LM, O'Toole
743 GA, Moulton LA, Ashare A, et al. 2014. Characterization and quantification of the fungal
744 microbiome in serial samples from individuals with cystic fibrosis. *Microbiome* **2**: 40.
- 745 Wylie KM, Mihindukulasuriya KA, Sodergren E, Weinstock GM, Storch GA. 2012. Sequence
746 analysis of the human virome in febrile and afebrile children. ed. C. Zhang. *PLoS ONE* **7**:
747 e27735.
- 748 Wylie KM, Mihindukulasuriya KA, Zhou Y, Sodergren E, Storch GA, Weinstock GM. 2014.
749 Metagenomic analysis of double-stranded DNA viruses in healthy adults. *BMC Biol* **12**: 71.
- 750 Wylie KM, Wylie TN, Orvedahl A, Buller RS, Herter BN, Magrini V, Wilson RK, Storch GA. 2015.
751 Genome Sequence of Enterovirus D68 from St. Louis, Missouri, USA. *Emerging Infect Dis*
752 **21**(1), 184-186.
- 753 Young JC, Chehoud C, Bittinger K, Bailey A, Diamond JM, Cantu E, Haas AR, Abbas A, Frye L,
754 Christie JD, et al. 2014. Viral Metagenomics Reveal Blooms of Anelloviruses in the
755 Respiratory Tract of Lung Transplant Recipients. *Am J Transplant* **15**: 200-209.
- 756 Yu G, Greninger AL, Isa P, Phan TG, Martínez MA, la Luz Sanchez de M, Contreras JF,
757 Santos-Preciado JI, Parsonnet J, Miller S, et al. 2012. Discovery of a novel polyomavirus in
758 acute diarrheal samples from children. ed. T. Ramqvist. *PLoS ONE* **7**: e49449.
- 759

Figure 1



Downloaded from academic.oup.com/guidelines/advance-article-abstract/doi/10.1093/guides/gaa011/5788888 on June 10, 2026 . Published by Cold Spring Harbor Laboratory Press

Figure 2

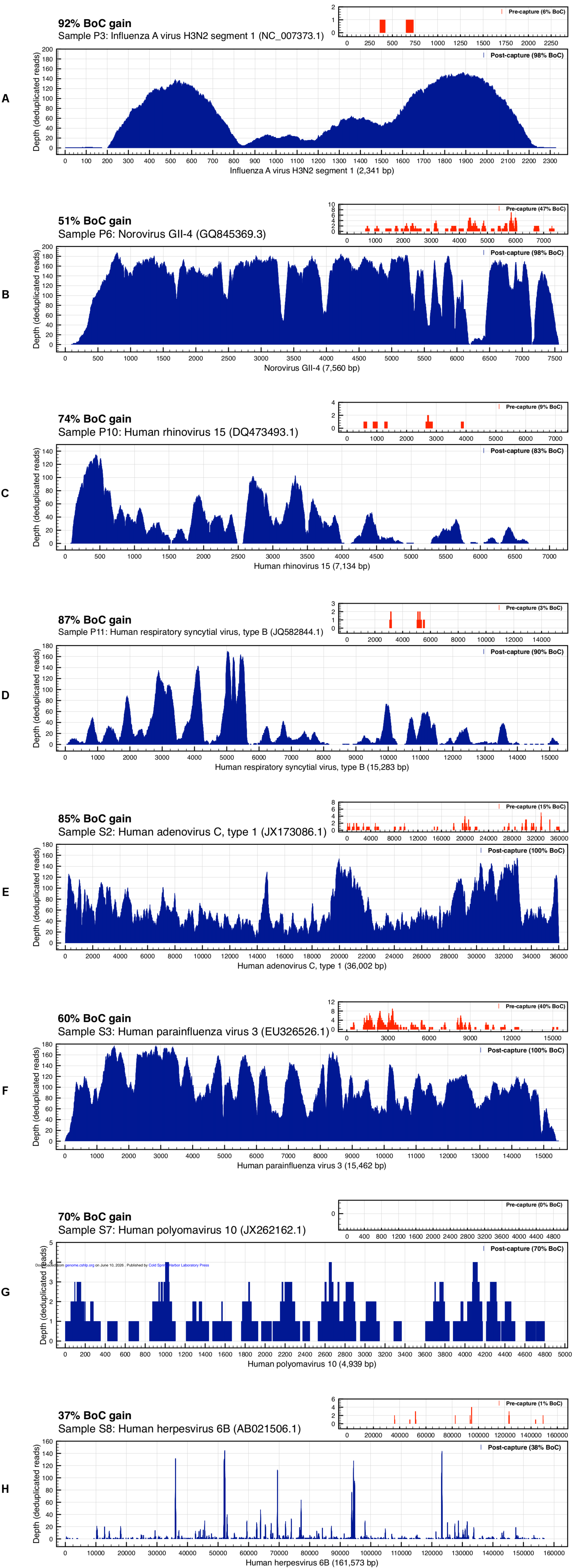
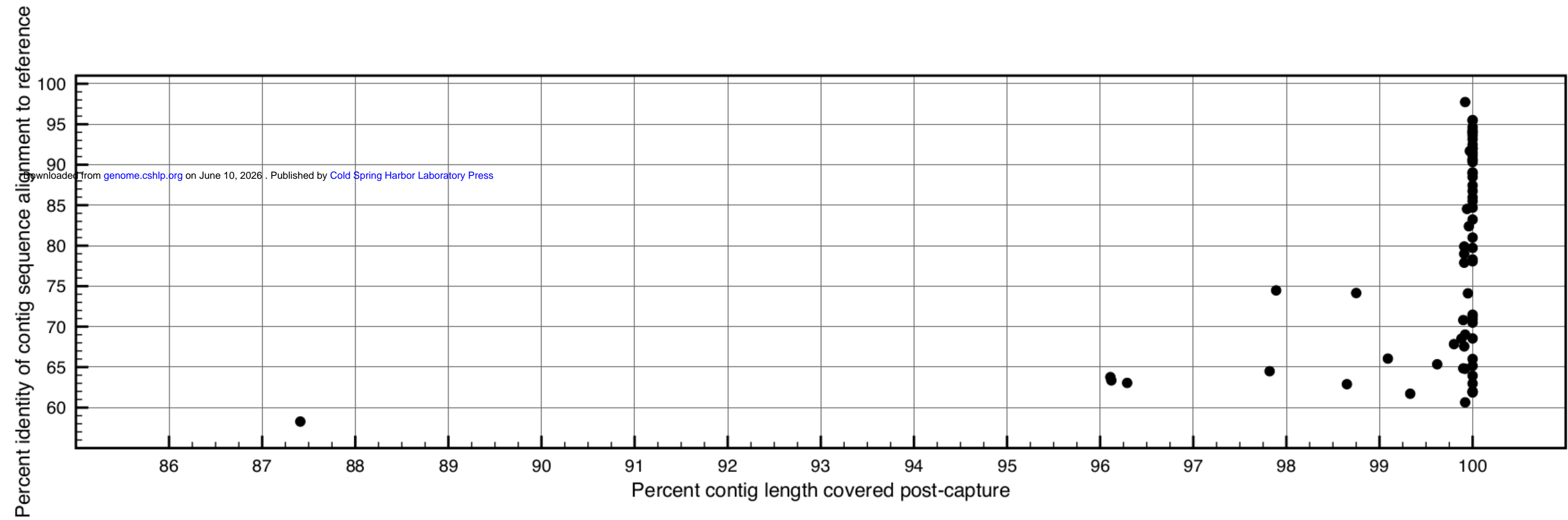


Figure 3



Anellovirus: Torque teno virus (FR751499.1)
Contig A170

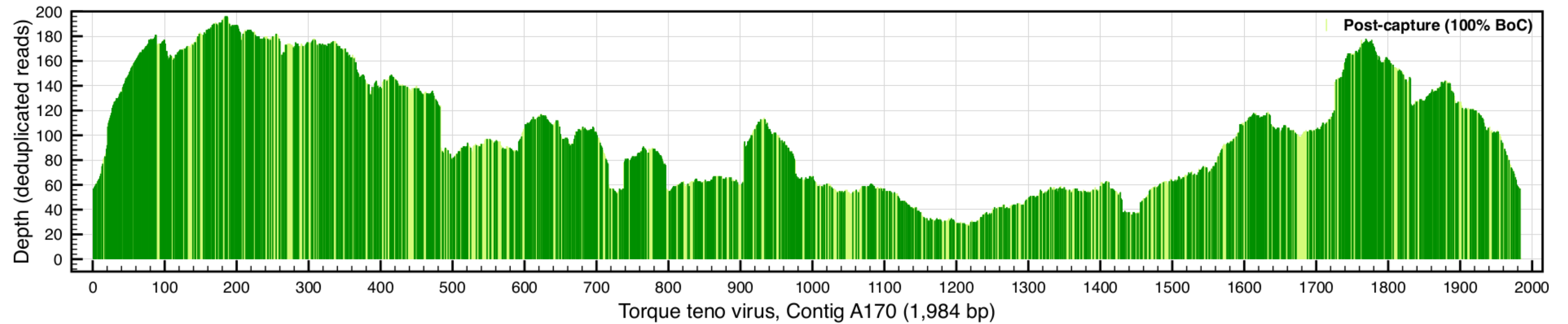
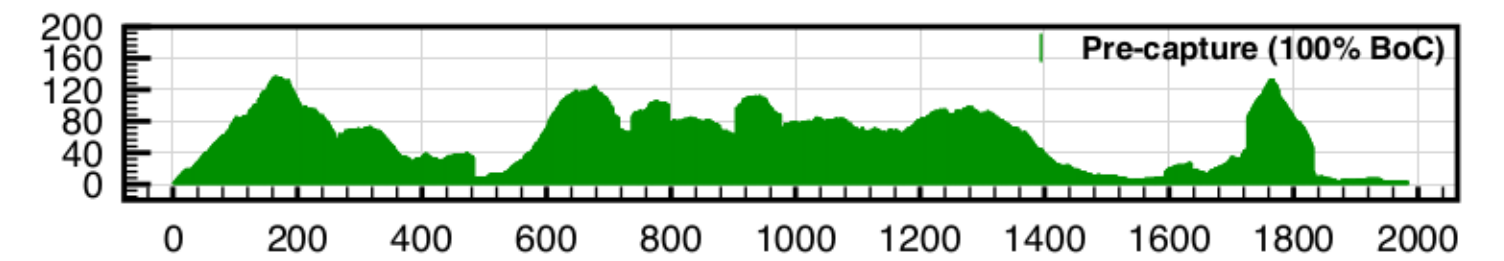


Table 1: Results of metagenomic shotgun sequencing for pooled specimens before and after viral targeted sequence capture

Virus ^a	Sample ID	Virus length (bp)	Viral sequence reads			Genome coverage			
			Viral read count [*]		PVR fold increase ^d	Breadth of coverage: percent		Depth of coverage: mean (SD)	
			Pre-capture ^b	Post-capture ^c		Pre-capture	Post-capture	Pre-capture	Post-capture
Human adenovirus B, type 35	P1	34,794	5	8,103	1,300	0.9	83.6	<0.1 (0.1)	19.3 (52.7)
Human bocavirus 1	P2	5,299	110	15,277	111	68.6	100.0	1.8 (2.0)	251.1 (122.1)
Influenza A virus (H3N2)	P3	13,267	4	46,540	9,335	2.4	74.0	<0.1 (0.3)	263.0 (394.4)
Influenza B virus	P4	14,452	0	513	>513 [‡]	0	9.8	0	2.6 (0.7)
Human parvovirus B19	P5	5,596	1,867	474,849	204	89.8	100.0	28.8 (17.9)	7,367.9 (4,193.0)
Norovirus GII-4	P6	7,560	72	527,656	5,880	46.5	98.4	0.8 (1.1)	5,870.4 (8,194.3)
Parechovirus 1	P7	7,339	0	13	>13 [‡]	0	8.0	0	0.2 (0.6)
BK polyomavirus	P8	5,142	1	1,520	1,220	1.6	88.6	<0.1 (0.1)	24.9 (32.2)
JC polyomavirus	P9	5,121	5	2,760	443	8.6	98.5	0.1 (0.3)	46.3 (54.0)
Human rhinovirus 15	P10	7,134	8	9,624	965	8.8	82.9	0.1 (0.3)	115.9 (170.3)
Human respiratory syncytial virus, type B	P11	15,283	9	67,778	6,042	3.0	89.5	<0.1 (0.3)	350.1 (1,667.2)
Human herpesvirus 1	P12	152,261	0	14	>14 [‡]	0	0.8	0	<0.1 (0.1)
Torque teno virus	P13	3,260	1	447	447	1.78	60.74	0.02 (0.13)	10.4 (23.4)
Human herpesvirus 3	P14	125,030	0	834	>834 [‡]	0	8.33	0	0.6 (3.5)

Pre-capture: Metagenomic shotgun sequencing without targeted sequence capture.

Post-capture: Metagenomic shotgun sequencing using (ViroCap) targeted sequence capture.

^{*} Viral reads per million sequences generated pre-capture and post-capture are statistically different ($P < 0.0001$, Wilcoxon test).

^a Viruses listed were incorporated into a viral pool that was subjected to MSS without and with targeted sequence capture, as described in the text.

^b MSS of the virus pool without targeted sequence capture yielded 7,458,192 total reads.

^c MSS of the virus pool after targeted sequence capture yielded 9,295,438 total reads.

^d PVR (percent viral reads) fold increase: percent of post-capture viral reads divided by percent of pre-capture viral reads.

[‡] PVR fold increase could not be calculated because the number of pre-capture reads was 0.

Table 2: Results of metagenomic shotgun sequencing for individual specimens before and after viral targeted sequence capture

Sample type	Sample ID	Total reads		Virus(es) detected	Virus length (bp)	Viral sequence reads			Genome coverage			
		Pre-capture	Post-capture			Viral read count [*]		PVR fold increase ^a	Breadth of coverage: percent		Depth of coverage: mean (SD)	
						Pre-capture	Post-capture		Pre-capture	Post-capture	Pre-capture	Post-capture
Nasopharyngeal swab	S1	4,202,474	16,080,640	Torque teno virus	3,736	0	231	>231 [‡]	0	74.4	0	5.2 (7.6)
				Human adenovirus B, type 3	35,269	15,116	7,703,787	133	89.6	100.0	36.1 (213.8)	19,097.4 (17,285.5)
Nasopharyngeal swab	S2	5,087,234	3,713,962	TTV-like mini virus isolate	2,912	0	249	>249 [‡]	0	37.5	0	6.4 (18.1)
				Human adenovirus C, type 1	36,002	87	116,502	1,834	15.3	100.0	0.2 (0.6)	283.1 (422.3)
Nasopharyngeal swab	S3	5,139,462	4,094,568	Torque teno mini virus	2,908	0	256	>256 [‡]	0	22.7	0	6.4 (21.8)
				Human parainfluenza virus 3	15,462	172	262,806	1,918	39.5	99.9	1.0 (1.6)	1,462.1 (1,767.0)
Nasopharyngeal swab	S4	6,124,424	5,659,554	TTV-like mini virus isolate	2,915	1	887	960	2.8	13.5	<0.1 (0.2)	23.9 (105.7)
				Human bocavirus 1	5,543	4,419	460,617	113	95.4	100.0	68.3 (94.1)	7,294.7 (3,443.1)
Nasopharyngeal swab	S5	9,152,970	5,834,446	Torque teno virus	3,741	0	56	>56 [‡]	0	14.9	0	1.0 (3.2)
Nasopharyngeal swab	S6	10,179,884	12,064,068	Human adenovirus B, type 3A	35,264	0	1,071	>1,071 [‡]	0	75.3	0	2.6 (2.8)
				KI polyomavirus	5,040	1,034	411,173	336	83.4	100.0	18.4 (20.5)	7,373.2 (7,123.3)
				Human rhinovirus 80	7,138	19,232	3,120,184	137	65.0	75.9	238.2 (313.2)	39,218.6 (102,545.7)
Stool	S7	3,691,496	4,104,534	Human adenovirus C, type 1	36,006	3	9,081	2,722	0.8	99.9	<0.1 (0.1)	22.9 (25.7)
				Sapovirus	7,429	33	61,982	1,689	23.3	100.0	0.4 (0.8)	740.1 (437.6)
				Human astrovirus 1	6,816	6,641	1,140,900	155	99.9	100.0	85.1 (135.6)	14,725.6 (13,370.2)
Plasma	S8	10,875,448	7,088,360	Human polyomavirus 10	4,939	0	81	>81 [‡]	0	70.4	0	1.6 (1.6)
				Torque teno virus	3,880	0	1,817	>1,817 [‡]	0	38.2	0	37.0 (105.0)
				Human herpesvirus 6B	161,573	38	27,523	1,111	1.1	38.4	<0.1 (0.2)	14.2 (112.9)

Pre-capture: Metagenomic shotgun sequencing without targeted sequence capture; Post-capture: Metagenomic shotgun sequencing using (ViroCap) targeted sequence capture.

* Viral reads per million sequences generated pre-capture and post-capture are statistically different (P=0.0002, Wilcoxon test).

^a PVR (percent viral reads) fold increase: percent of post-capture viral reads divided by percent of pre-capture viral reads.

[‡] PVR fold increase could not be calculated because the number of pre-capture reads was 0.