

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

Effect of geographic isolation on the nasal virome of indigenous children

Eda Altan^{1,2}, Juan Carlos Dib³, Andres Rojas Gullos³, Duamaco Escribano Juandigua³, Xutao, Deng^{1,2}, Roberta Bruhn^{1,2}, Kristen Hildebrand⁴, Pamela Freiden⁴, Janie Yamamoto^{1,2}, Stacey Schultz-Cherry⁴, Eric Delwart^{1,2}.

¹ Vitalant Research Institute, 270 Masonic Ave, San Francisco, California, United States of America

² University of California San Francisco, Department of Laboratory Medicine, San Francisco, California, United States of America

³ Fundación Salud Para el Trópico-Tropical health Foundation, Santa Marta, Magdalena, Colombia

⁴ Department of Infectious Diseases, St Jude Children's Research Hospital, Memphis, Tennessee, United States of America

*Corresponding author

E-mail: eric.delwart@ucsf.edu

Running title: Nasal virome in isolated villages

22 **Abstract**

23 The influence of living in small remote villages on the diversity of viruses in the nasal
24 mucosa was investigated in three Colombian villages with increasing levels of geographic
25 isolation. Viral metagenomics was used to characterize viral nucleic acids on nasal swabs of
26 63 apparently healthy young children. Sequences from human virus members of the families
27 *Anelloviridae*, *Papillomaviridae*, *Picornaviridae*, *Herpesviridae*, *Polyomaviridae*, *Adenoviridae*,
28 *and Paramyxoviridae* were detected in a decreasing fraction of children. The diversity of human
29 viruses was not reduced in the most isolated indigenous Kogi villages. Multiple viral
30 transmission clusters were also identified as closely related variants of rhinoviruses A or B in 2
31 to 4 children from each of villages. The number of papillomavirus detected was greater in the
32 village most exposed to outside contacts while conversely more anellovirus infections were
33 detected in the more isolated indigenous villages. Genomes of viruses not known to infect
34 humans, including in the family *Parvoviridae* (genus densovirus), *Partitiviridae*, *Dicistroviridae*,
35 *and Iflaviridae* and circular Rep expressing ssDNA genomes (CRESS-DNA) were also detected
36 in nasal swabs likely reflecting environmental contamination from insect, fungal, and unknown
37 sources. Despite the high level of geographic and cultural isolation, the diversity of human
38 viruses in the nasal passages of children was not reduced in indigenous villages indicating
39 ongoing exposure to globally circulating viruses.

40 **Importance**

41 Extreme geographic and cultural isolation can still be found in some indigenous South
42 American villages. Such isolation may be expected to limit the introduction of globally circulating
43 viruses. Very small population size may also result in rapid local viral extinction due to lack of
44 sufficient sero-negative subjects to maintain transmission chains of rapidly cleared viruses. We
45 compared the viruses in the nasal passage of young children in three villages with increasing
46 level of isolation. We find that isolation did not reduce the diversity of viral infections in the most
47 isolated villages. Ongoing viral transmission of rhinoviruses could also be detected within all
48 villages. We conclude that despite their geographic isolation remote villages are continuously
49 exposed to globally circulating respiratory viruses.

50

51

52

53 **Introduction**

54 The impact of geographic isolation in shaping the respiratory virome remains largely
55 unknown. In the pre-agricultural era, people typically lived widely dispersed in small nomadic
56 groups, a lifestyle which may have minimized the spread and maintenance of infectious
57 diseases that did not establish long lasting or chronic infections. Small populations now settled
58 in hard to explore regions may still be relatively isolated from repeated exposures to highly
59 prevalent viruses circulating in larger, more connected, communities. Inhabitants of such highly
60 isolated villages may have therefore lost viruses dependent on large population size of young,
61 seronegative, susceptible hosts found in larger populations [1].

62 Coincident with the arrival of Europeans, native Amerindian populations underwent
63 strong population bottlenecks possibly due to imported airborne epidemics such as small pox,
64 measles, and more recently influenza viruses to which they had no prior exposure [2,3]. To
65 determine whether reduced rate of outside contact coincides with a reduction or even an
66 absence of detectable human viruses, we analyzed and compared the nasal virome of children
67 in two highly isolated Amerindian Kogi villages in a tropical forest of Northern Colombia and of
68 one largely Hispanic village alongside a coastal highway. In order to detect all human viruses,
69 viral metagenomics was applied to nasal swabs collected from children two to nine years old.

70 **Results**

71 **Sample collection and village location**

72 Nasal swabs were collected from 63 children (53.9% female) with a mean age of 5 years
73 (Table 1). The children lived in three Northern Columbia villages that differed in degree of
74 outside contacts. Samples used for comparison were from age, sex and race matched children
75 (Table 1). The first village Calabazo (GPS 11.28448, -74.00195) is located along a major road
76 (highway 90) running alongside the National Natural Park of Tayrona and is frequently visited by
77 tourists. Calabazo has a 2005 census population size of 499 and the main language is Spanish.
78 Seywiaka (GPS 11.2174, -73.5794) is an isolated village with a population size of 250-300
79 accessible only by foot (1.5 hours walk from nearest road) inhabited by Kogi people speaking
80 their indigenous language (Fig 1A). An even more isolated Kogi village Umandita (GPS
81 11.09698,-73.64781) with an estimated population of 350-400 inhabitants is accessible after a
82 9-10 hours walk from Seywiaka (Fig 1B).



83
84

Figure 1A. View of Seywiaka village.



85

86 Figure 1B. View of Umandita village.

87 **Nasal mucosa virome**

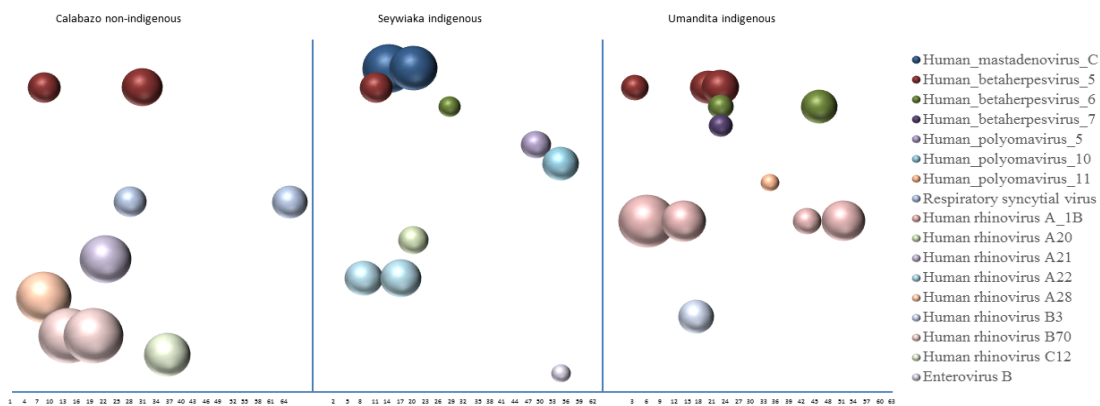
88 Following viral metagenomics enrichment of viral particles-associated nucleic acids in
89 nasal swabs, random RNA and DNA amplification, and deep sequencing of 63 individual nasal
90 swab supernatants a total number of 63 million reads were generated for an average number of
91 reads of approximately one million per sample. The raw sequence data for each pool is
92 available at NCBI's Short Reads Archive under GenBank accession number PRJNA530270. We
93 found 92% of samples (58/63) to be positive for at least one human virus. Human virus
94 belonging to 7 viral families were detected and are listed in decreasing prevalence of detection
95 (*Anelloviridae*, *Papillomaviridae*, *Picornaviridae*, *Herpesviridae*, *Polyomaviridae*, *Pneumoviridae*,
96 *and Adenoviridae*).

97 *Anelloviridae* family members reads were the most commonly detected viral sequences
98 and were found in 49/63 children or 77.7%. 0.16% (n=100,957 sequence reads) of 63 million

99 total reads could be mapped to the *Anelloviridae* family with BLASTx E scores $<10^{-10}$. The
100 second most commonly detected human virus reads belonged to the *Papillomaviridae* family,
101 which were detected in 44.4% (28/63 children) with 0.087% of total reads (n= 55,248). Next, with
102 a prevalence of 23.8% (15/63 children) were reads from the *Picornaviradae* family. Of these,
103 0.094% reads (n=59,819) encoded picornavirus reads from the species Rhinovirus A (10/63
104 children, 15.8%), Rhinovirus B (3/63 children, 4.7%), Rhinovirus C (1/63 children, 1.58%) and
105 Enterovirus B (1/63 children, 1.58%). *Herpesviridae* family members were next in prevalence
106 being detected in 7/63 children (11.1%) including human betaherpesvirus 5 (CMV or HHV5)(
107 6/63 children, 9.52%), human herpesvirus 6 (Roseolovirus or HHV6) (3/63 children, 4.76%), and
108 human betaherpesvirus 7 (Kaposi Sarcoma virus or HHV7) (1/63 children, 1.58%). In the
109 *Polyomaviridae* family, human polyomavirus 5 (Merkel Cell carcinoma virus or HPyV5 [4]) (1/63
110 children, 1.58%), human polyomavirus 10 [5,6], (1/63 children, 1.58%), human polyomavirus
111 11[7] (1/63 children, 1.58%) were detected. Adenovirus C reads were detected in 2/63 children
112 (3.17%). Respiratory syncytial virus (RSV), belonging to the *Paramyxoviridae* family, was found
113 in 2/63 children (3.17%), This was the only viral family detected exclusively in the most exposed
114 Calabazo village. The fraction of total reads from each sample encoding proteins with high-level
115 similarity (E scores $<10^{-10}$) to human viruses are shown in (Fig. 2) with the exception of the
116 papillomaviruses and anelloviruses that are analyzed below.

117

118



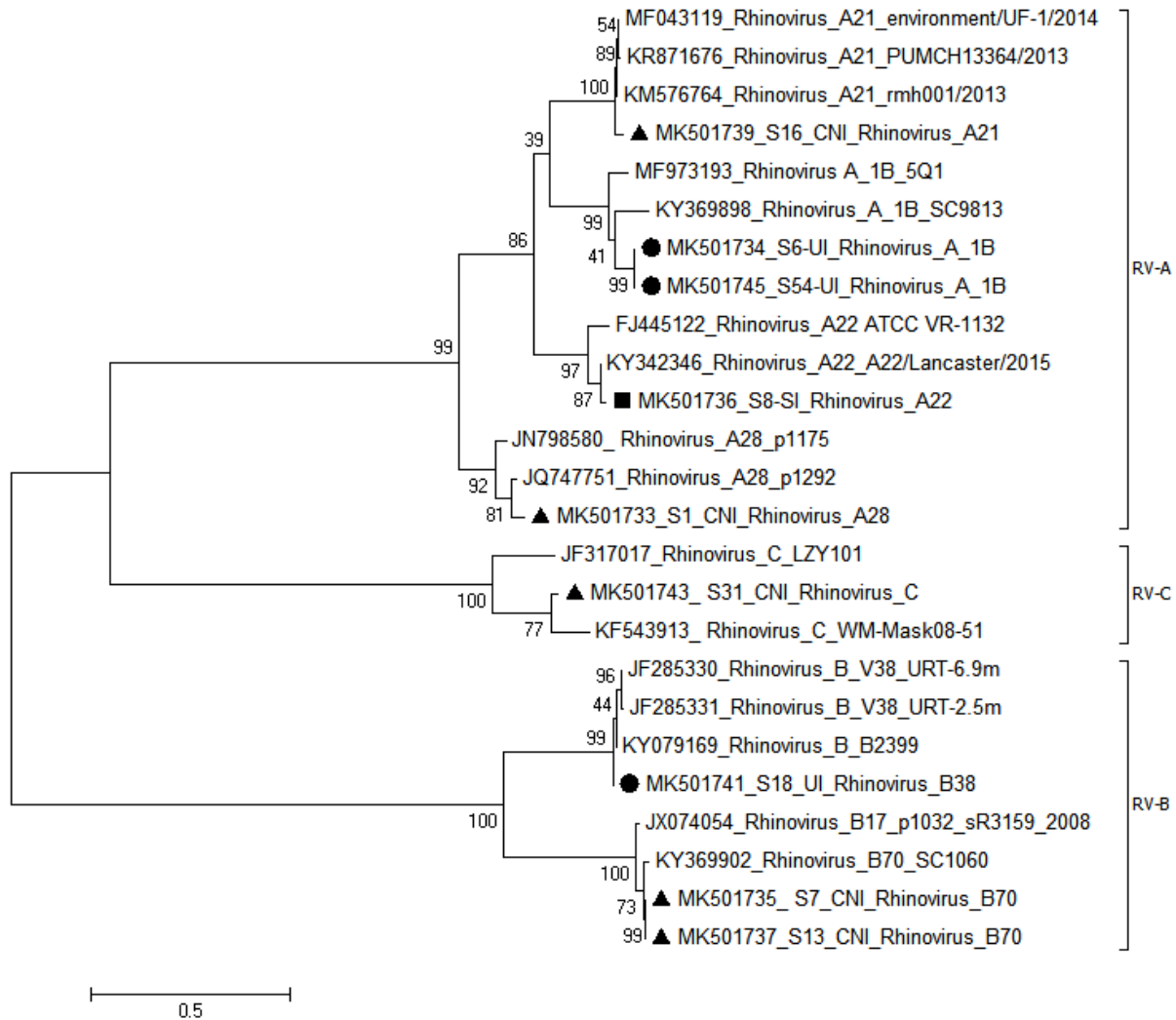
120 Figure 2. Distribution and level of viral reads to human viruses in the three villages.

121

122

123 **Family *Picornaviridae***

124 Fourteen children showed the presence of picornavirus sequences, 13 from rhinoviruses
125 A, B, or C species and 1 from enterovirus B species. Rhinovirus (RV) reads generated contigs
126 ranging in size from 481 to 7,089 nt (GenBank accession no MK501733 - MK501745). In total,
127 8 contigs included complete 5UTR-VP4-VP2, 2 contigs complete 5UTR-VP4, and 1 contigs a
128 complete VP4-VP2 sequences that were used for phylogenetic analysis (Fig 3). Four rhinovirus
129 contigs from Umandita region showed closest nucleotide identity (90 to 92%) to genotype 1B of
130 rhinovirus A (RV-A-1B). Contigs from 3 of these children overlapped over almost the entire
131 genome (6.6Kb without gaps) and showed a nucleotide identity of 99.8-99.9%. These three
132 rhinovirus A-1B contigs clustered tightly together reflecting a recent common origin and an
133 ongoing transmission cluster in the most isolated village, Umandita. Two children from Calabazo
134 were infected with rhinovirus B70 that had 99.9% nucleotide identity indicating another
135 transmission cluster occurring at the time of sampling. Two children from Seywiaka were
136 shedding rhinovirus A22 but did not generate enough sequence reads to be included in
137 phylogenetic analysis. Reads from these two children did overlap by 154 bases showing a
138 single mismatch indicating another possible transmission cluster. Rhinovirus transmission
139 clusters (genotypes A1B in Umandita, B70 in Calabazo, and A22 in Seywiaka) were therefore
140 detected in each village. The enterovirus B reads from a Seywiaka child showed closest amino
141 acid identity (93%) to Echovirus E15 (GenBank AY302541).



143 Figure 3. Phylogenetic analysis of VP4-VP2 region of rhinoviruses.

144 ● Calabazo non-indigenous (CNI), ▲ Seywiaka indigenous (SNI), ■ Umandita indigenous (UI)

145 Polyomaviruses

146 Polyomavirus sequences were also found in the isolated villages of Seywiaka (n=2) and
 147 Umandita (n=1). Two Seywiaka villagers were shedding human polyomavirus 5 (Merkel cell
 148 polyomavirus), or human polyomavirus 10 (MW polyomavirus), and one child from Umandita
 149 was shedding human polyomavirus 11 (STL polyomavirus)(Fig 2).

150 Herpesviruses

151 Sequences of human CMV, roseolovirus, and Kaposi sarcoma virus were identified.
 152 CMV sequences were found in six children (2, 1 and 3 children from Calabazo, Seywiaka and

153 Umandita respectively). Three children shed Roseolovirus (2 and 1 from Seywiaka and
154 Umandita respectively). One Kaposi sarcoma virus infection was detected in a child from
155 Umandita (Fig 2). All contigs showed 98-100% nucleotide identities to genomes in GenBank.

156 **Adenovirus, pneumovirus, and parvovirus**

157 Sequences from human_mastadenoviruses C species (HAdV-C) in the *Adenoviridae*
158 family ranging in size from 250 to 1831 nt, were identified in two children from Seywiaka village
159 (Fig 2). Six different regions (E3, E4, E1a, and L3) showed overlap between children with
160 nucleotide identity of 83 to 97% likely reflecting two independent infections

161 Two respiratory syncytial virus generating contigs of size 363 nt and 888 nt were
162 generated from two children in Calabazo both showing 99% nucleotide identity with respiratory
163 syncytial virus strain A (GB accession number MG793382) (Fig 2). Contigs from these two
164 children overlapped in the G gene (350 bp) showing a nucleotide identity of 99.1%. The close
165 sequence identity of these two RSV strain may also reflect an ongoing transmission cluster
166 within that village.

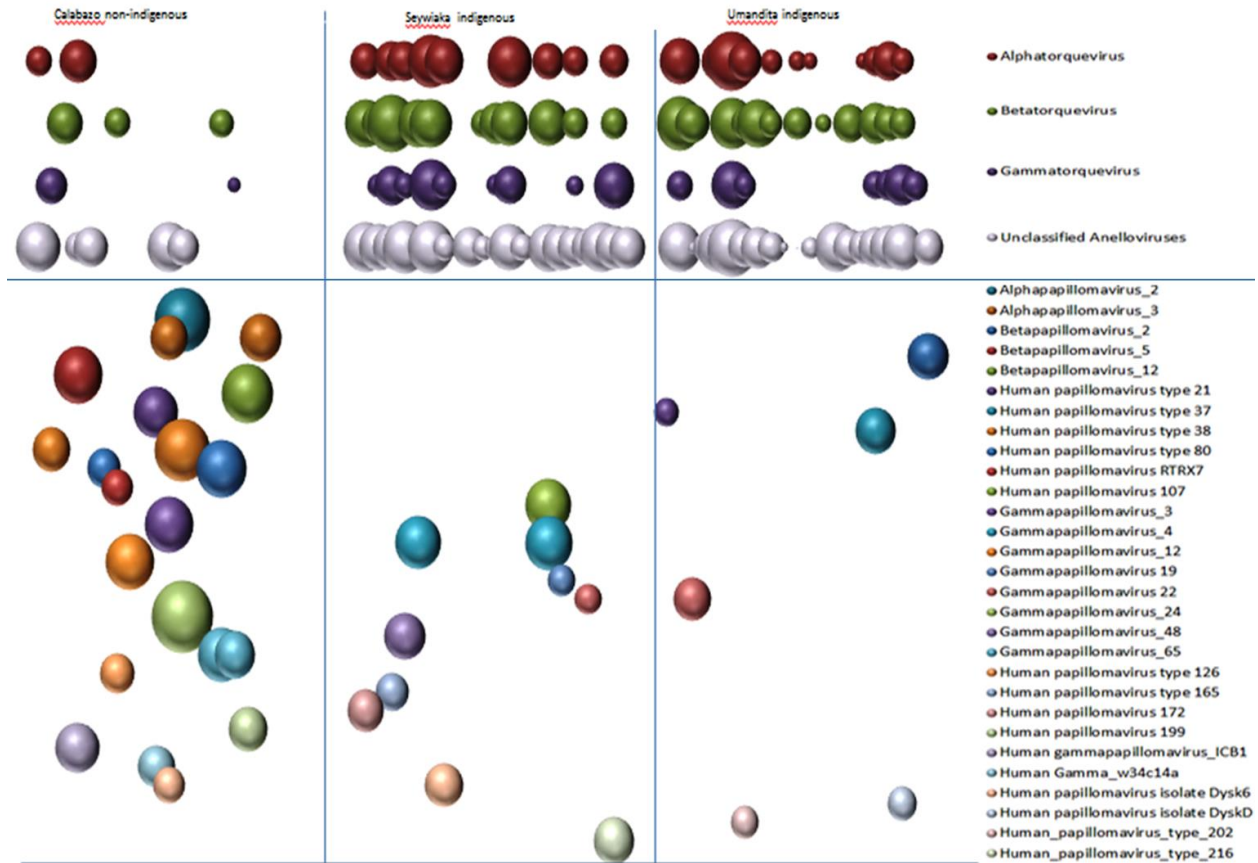
167 Unexpectedly, two reads of canine bocavirus were also detected in one swab sample
168 from Calabazo village showing 92 and 97% aa identity to canine bocavirus NS1 gene region
169 (GB accession number MG025952).

170 **Anelloviruses**

171 Reads matching *Anelloviridae* family viruses were found in 77.7% (49/63) of children.
172 Prevalence of anellovirus detection was 42% (10/21), 90% (19/21) and 95.2% (20/21) in
173 children from Calabazo, Seywiaka and Umandita respectively. The overall fraction of children
174 infected with different anellovirus genera were 34% with alphatorquevirus, 44.4% with
175 betatorquevirus, 28.5% with gammatorquevirus and 65% with unclassified anelloviruses.

176 **Papillomaviruses**

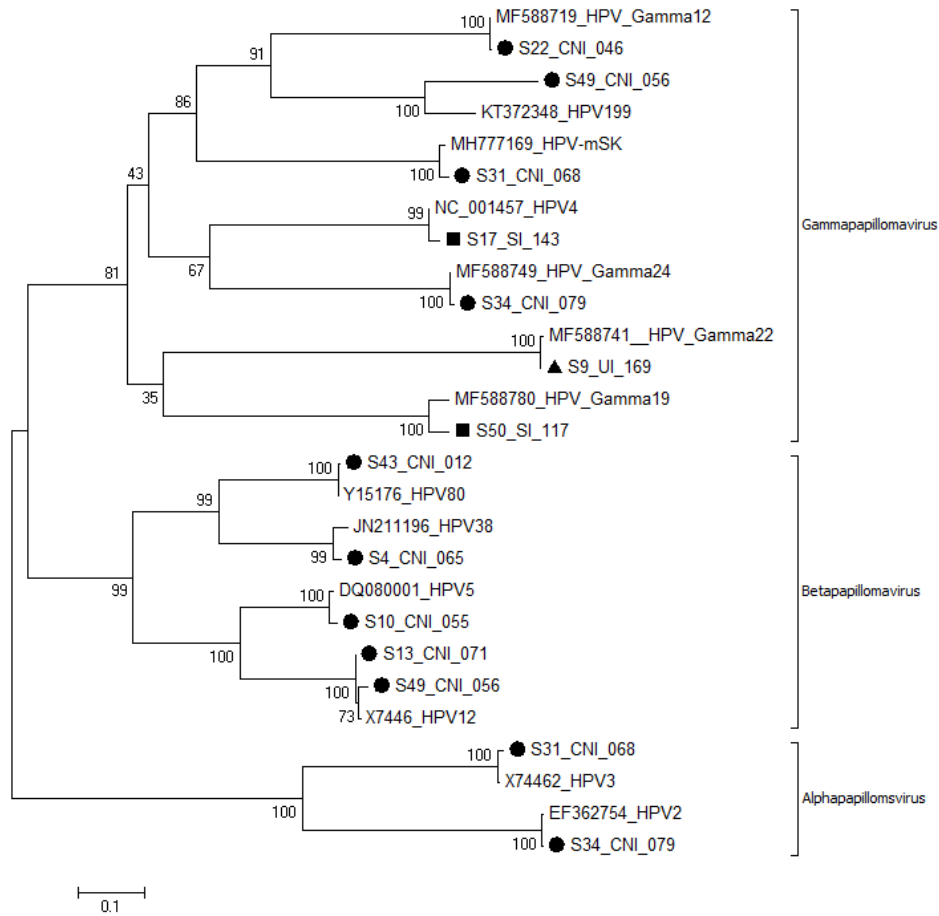
177 Altogether we detected 29 papillomaviruses consisting of 17 genotypes in 13 Calabazo
178 children; 10 papillomaviruses consisting of 9 genotypes in 9 Seywiaka children; and 6
179 papillomaviruses consisting of 6 genotypes in 6 Umandita children (Figure 4B).



180

181 Figure 4A and B. Distribution and level of anellovirus and papillomavirus reads in all three
 182 villages.

183 37 partial papillomavirus genome contigs ranging in size from 261 nucleotides (nt) to
 184 7,392 nt were generated, 14 of which included a partial L1 gene region. Phylogenetic analysis of
 185 these L1 sequences was generated (Fig 4). All papillomavirus contigs showed 97-100% aa
 186 identities to papillomavirus proteins in GenBank. Papillomaviruses (HPV12) in two children from
 187 Calabazo village (S13-CNI, S49-CNI), were closely linked (Fig 5) showing 99% nucleotide
 188 identity.



189

0.1

190 Figure 5. Phylogenetic analysis of major capsid (L1) protein of papillomaviruses.

191 ● Calabazo non-indigenous (CNI), ▲ Seywiaka indigenous (SNI), ■ Umandita indigenous(UI)

192 Virome comparison between villages

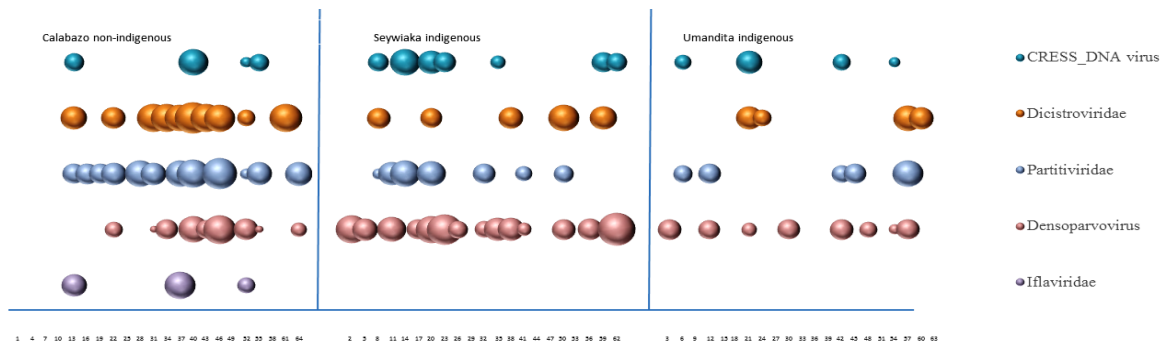
193 We next compared the distribution of the two viral families that yielded the most reads,
194 anelloviruses and papillomaviruses, among the 3 villages (Fig 4). The numbers of anellovirus
195 infections were significantly different among the villages ($p=0.0001$). Inspection of the
196 anellovirus distribution indicated that fewer infections were detected in the most exposed
197 Hispanic village of Calabazo.

198 An analysis of papillomavirus reads distribution also showed that the number of
199 papillomavirus infections were significantly different among the villages ($p=0.043$). As opposed
200 to anelloviruses a greater number of papillomavirus infections were detected in the two most
201 isolated villages relative to Calabazo (Fig 4).

202 **Viral families of non-vertebrate or unknown host tropism**

203 Sequences from viral families not known to infect human (or vertebrates), likely
204 representing air-borne mucosal surfaces contamination, were detected in 44/63 (69.8%) of
205 children. Members of the following viral clades, ranked from highest to lowest prevalence, were
206 detected (*Parvoviridae*-densoviruses, *Partitiviridae*, *Dicistroviridae*, *circular Rep encoding single*
207 *stranded DNA viruses-CRESSS-DNA*, and *Iflaviridae*) were found in 49.2%, 38.09%, 30.1%,
208 23.8%, 4.7% of the swab samples respectively (Fig 6).

209 Densoparvoviruses, dicistroviruses, and iflaviruses are known to infect invertebrates and
210 partitiviruses to infect fungi and plants. Some *CRESS-DNA* viruses can infect fungi, plants, or
211 mammals, but the tropism of most *CRESS-DNA* genomes largely identified through
212 metagenomics of environmental samples (including those detected here) remain unknown.



213

214 Figure 6. Distribution of viral sequences from viral groups not known to infect vertebrates
215 or of unknown tropism.

216

217

218

219 Discussion

220 Viral metagenomics of human respiratory secretions have analyzed different sample
221 types mainly from clinical cases while only a few studies have studied samples from healthy
222 subjects. Early studies of nasopharyngeal aspirates from lower respiratory tract infections in
223 China [8] and Sweden [9] revealed numerous viruses with members of the
224 families *Paramyxoviridae* (respiratory syncytia virus, metapneumovirus, parainfluenza virus),
225 *Picornaviridae* (rhinovirus), *Orthomyxoviridae* (influenza) pre-dominating. A viral metagenomics
226 analysis of nasopharyngeal swabs, aspirates, and sputums from patients with
227 acute respiratory infections could confirm the presence of diverse viruses expected from PCR
228 results [10]. Pediatric nasopharyngeal swabs viral metagenomics results also showed high
229 concordance with the results of a commercial respiratory virus panel as well as reveal the
230 presence of other, non-tested for, viruses [11]. Nasopharyngeal/oropharyngeal swabs from
231 children with pneumonia of unknown etiology and asymptomatic controls, when analyzed by
232 viral metagenomics could identify numerous viruses some of which could be associated with
233 respiratory symptoms [12]. Nasopharyngeal swabs from children with community-acquired
234 pneumonia but negative for common respiratory viruses revealed non-tested for viruses and
235 human parainfluenza 3 virus with a large deletion that may have precluded specific PCR
236 amplification [13] highlighting an advantage of non-specific amplification and deep sequencing.
237 Lung transplant recipients with respiratory symptoms showed human rhinovirus infections
238 recalcitrant to PCR detection as well as frequent HHV7 and anellovirus infections [14]. Diverse
239 viruses could be detected in about a third of immune-compromised children with pulmonary
240 disease including co-infections missed by prior clinical tests together with untested for viruses
241 [15]. Metagenomics analyses therefore shows great promises as a supplement or even
242 replacement for more specific viral genome detection assays although sensitivity issues remain
243 [16-19].

244 A more limited number of metagenomics studies have analyzed the respiratory track
245 virome of healthy children. Double stranded DNA from the anterior nares of healthy individuals
246 in the human microbiome project showed that beta and gamma papillomaviruses were the most
247 common viruses detected, followed by roseolovirus (HHV6) [20]. A PCR study of sinonasal
248 mucosa from sinus surgery patients for 16 common respiratory viruses indicated HHV6 was the
249 most the frequently detected virus [21]. A study of nasopharyngeal swabs from healthy 18
250 months old children showed the presence of human rhinoviruses, adenoviruses, bocaviruses,

251 and parainfluenza virus [22]. Nasopharyngeal swabs from healthy children showed
252 anelloviruses, HHV6, and HHV7 to be the most common infections [12].

253 Here we characterize the nasal mucosal viromes of 63 age and gender matched
254 children from three villages and show that the geographical and cultural isolation of the two
255 indigenous villages did not eliminate or even reduce the diversity of their human viruses. Three
256 different herpesviruses (HHV5-7) and three different polyomaviruses (human polyomavirus 5,
257 10, 11) were detected in the two isolated villages while only one herpesvirus (HHV5) and no
258 polyomaviruses were found in the village with frequent outside contact (Fig 2). Respiratory
259 syncytia virus was the only virus found exclusively in the most exposed village. Four different
260 rhinovirus and one enterovirus B genotypes were found in the two isolated villages while four
261 rhinovirus genotypes were detected in the more exposed village. There was no overlap in the
262 picornavirus genotypes in the different villages. In 0-5 years old children, the rates of
263 asymptomatic rhinovirus detection have been reported from 12.5 to 33% [23-26]. In under 3
264 years old asymptomatic children, a rhinovirus detection rate of 33% did not significantly differ
265 from that found in matched hospitalized children [25]. Here, we found an average 20.6% rate of
266 rhinovirus detection in healthy 2-9 years old children, ranging from 23% in Calabazo and
267 Umandita villages to 14% in Seywiaka.

268 Outbreaks, as reflected by the detection of closely related variants of the same
269 rhinovirus genotypes, were detected in both the most isolated (four cases of rhinovirus A1B in
270 Umandita) and the least isolated (two cases of rhinovirus B70 in Calabazo) villages. Two
271 rhinovirus A22 infections in isolated Seywiaka were also very closely related and likely also
272 epidemiologically linked. Because rhinovirus infections are of short duration it seems likely that
273 each of these clusters resulted from recent introductions within these communities.

274 The origin of sequence reads from viral clades not known to infect humans, namely from
275 the *Parvoviridae* genus densoparvovirus, *Partitiviridae*, *Dicistroviridae*, *Iflaviridae* and
276 CRESS_DNA viruses remains unknown but deposition onto nasal mucosa surfaces from
277 environmental sources such as the ambient air remains a likely possibility. Possible source for
278 such viruses include plants and fungi for the partitiviruses, and insects for the dicistroviruses,
279 iflaviruses, and densoparvoviruses. The origins of CRESS-DNA viral genomes are unknown.
280 The detection of a very few reads of canine bocavirus (n=3), a virus reported in dogs as well as
281 cats [27-31], might similarly reflect environmental contamination from local pets.

282 More frequent infections with anelloviruses were detected in the more isolated villages of
283 Seywiaka and Umandita. Anellovirus concentration in blood are highly dependent on the host's
284 immunocompetence and viral titers have been shown to increase in febrile patients [32],
285 immune-suppressed transplant patients [33-35] and AIDS patients [36-41]. The higher rate of
286 detectable anellovirus infections in the most isolated villages may therefore reflect generally
287 weaker immune systems leading to more readily detectable anelloviruses possibly a result of
288 poorer diet and medical care in these remote locations. The converse relationship was found
289 for papillomaviruses which were more commonly detected in the most exposed village of
290 Calabazo (21 distinct infections) versus the more isolated villages (10 and 6 infections in
291 Seywiaka and Umandita respectively). Carcinogenic papillomaviruses were not detected. The
292 number of papillomavirus infections therefore appears to correlate with the amount of exposure
293 to people from outside their villages and was the only virus family where geographical isolation
294 was associated with reduced viral diversity. Whether papillomavirus infections are consistently
295 lower in prevalence in other small, isolated, villages relative to more connected or larger
296 populations, remains to be further tested and confirmed.

297 We therefore show here that children from both connected and highly isolated villages in
298 Northern Colombia carry diverse human viruses in their nasal mucosa, most frequently
299 anelloviruses and papillomaviruses, that rhinovirus transmission clusters can be readily
300 detected within these small communities, and that extreme geographical and cultural isolation
301 did not result in a general reduction in viral diversity.

302 Closely related picornaviruses (and caliciviruses) have also been described in fecal
303 samples from children of highly isolated Amazonian villages in Venezuela. This observation
304 likely also reflect ongoing transmission chains among epidemiologically linked children within
305 very small communities as described here [42]. This recent study of fecal viromes also showed
306 that extreme isolation did not reduce the diversity of circulating enteric viruses [42] as we show
307 here for the nasal mucosa. These results support our conclusion that the current reach of
308 common human viruses extends to some of the most geographically remote populations.

309 **Material Method**

310 **Study population and Study design**

311 Nasal swab samples were collected from September 2016 to February 2017 from
312 children with no apparent clinical signs enrolled in an influenza surveillance study located in
313 three different villages in the Magdalena Department of Colombia by the Caribbean sea (Figure

314 1). Nasal swabs from 21 children from each village were collected totaling 63 samples from 34
315 girls and 29 boys (Table 1). Dry sterile swabs (Nylon flocked, Fisher) were used in both nostril
316 and stored in 1 ml universal transport medium (Quidel). Samples were kept on ice for 4-7 days
317 and then stored at -80C.

	2 to 5 years old	6 to 9 years old	Girl	Boy
Calabazo	11	10	8	13
Seywiaka	12	9	9	12
Umandita	13	8	17	4

318 Table 1. Age and gender of children analyzed.

319 **Viral metagenomics**

320 To reduce possible batch effects, samples from the 3 locations were processed in an
321 interdigitated manner (First sample from village 1,2,3, then second sample from village 1,2,3,
322 then repeat) using two Illumina MiSeq runs. Individual swab supernatants (150ul) were filtered
323 using a 0.45- μ m filter (Millipore). The filtrates were treated with a mixture of DNases (Turbo
324 DNase [Ambion], Baseline-ZERO [Epicentre], benzonase [Novagen]) and RNase (Fermentas)
325 at 37°C for 90 minutes to enrich for viral capsid-protected nucleic acids were then extracted
326 using a Maxwell 16 automated extractor (Promega)[43]. Random RT-PCR followed by
327 Nextera™ XT Sample Preparation Kit (Illumina) were used to generate a library for Illumina
328 MiSeq (2 × 250 bases) with dual barcoding as previously described[44].

329

330 **Bioinformatic analysis**

331 An in-house analysis pipeline was used to analyze sequence data. Before analyzing raw
332 data was pre-processed by subtracting human and bacterial sequences, duplicate sequences,
333 and low quality reads. Following de novo assembly using the Ensemble program [45], both
334 contigs and singlets viral sequences were then analyzed using translated protein sequence
335 similarity search (BLASTx v.2.2.7) to all annotated viral proteins available in GenBank.
336 Candidate viral hits were then compared to a non-virus non-redundant (nr) protein database to
337 remove false positive viral hits. To align reads and contigs to reference viral genomes from
338 GenBank and to generate complete or partial genome sequences the Geneious R10 program
339 was used. For plotting read numbers to different viruses the number of reads with BLASTx E
340 score $<10^{-10}$ to named viruses was divided by the total number of reads multiplied by 10^4 then
341 log 10 transformed to determine the size of the colored circles using Excel.

342 **Phylogenetic analyses**

343 Phylogenetic trees were constructed from VP4-VP2 nucleotide sequence for
344 rhinoviruses and amino acid sequence for papillomaviruses. Evolutionary analyses were
345 conducted in MEGA6 using the using the Maximum Likelihood method based on the General
346 Time Reversible model [46,47].

347
348 **Statistical methods**

349 To evaluate the proportional distribution of viral types among villages, a nonparametric,
350 oneway, ANOVA was performed using the Kruskal Wallis test with ties and an a priori statistical
351 significance level set at $p < 0.05$. Stata/MP 15.1 (StataCorp, College Station, TX) was used for
352 the statistical analysis. The Kruskal-Wallis equality of population rank test was done using two
353 degree of freedom.

354

355 **Ethics statement**

356 Studies were approved by the Indigenous Health Council, Tropical Health Foundation
357 Ethics Committee, and St. Jude Children's Research Hospital Institutional Review Board. The
358 investigators ensure that this study is conducted in full conformity with the principles set forth in
359 The Belmont Report: Ethical Principles and Guidelines for the Protection of Human Subjects of
360 Research of the US National Commission for the Protection of Human Subjects of Biomedical
361 and Behavioral Research (April 18, 1979) and codified in 45 CFR Part 46, 21 CFR 50, 21 CFR
362 56 and/or the ICH E6; 62 Federal Regulations 25691 (1997), if applicable. The investigator's
363 Institution's hold current Federal Wide Assurances (FWA) issued by the Office of Human
364 Research Protection (OHRP) for federally funded research.

365

366

367

368

369 **Acknowledgments**

370 Funding sources consisted of support from Vitalant Inc. to ED and ALSAC and NIAID contract
371 HHSN272201400006C to SSC.

372

373

374

375

References

- 376 1. Black FL (1975) Infectious diseases in primitive societies. *Science* 187 (4176):515-518
- 377 2. O'Fallon BD, Fehren-Schmitz L (2011) Native Americans experienced a strong population bottleneck
- 378 coincident with European contact. *Proc Natl Acad Sci U S A* 108 (51):20444-20448.
- 379 doi:10.1073/pnas.1112563108
- 380 3. Walker RS, Sattenspiel L, Hill KR (2015) Mortality from contact-related epidemics among indigenous
- 381 populations in Greater Amazonia. *Sci Rep* 5:14032. doi:10.1038/srep14032
- 382 4. Feng H, Shuda M, Chang Y, Moore PS (2008) Clonal integration of a polyomavirus in human Merkel
- 383 cell carcinoma. *Science* 319 (5866):1096-1100. doi:10.1126/science.1152586
- 384 5. Siebrasse EA, Reyes A, Lim ES, Zhao G, Mkakosya RS, Manary MJ, Gordon JI, Wang D (2012)
- 385 Identification of MW polyomavirus, a novel polyomavirus in human stool. *J Virol* 86 (19):10321-10326.
- 386 doi:10.1128/JVI.01210-12
- 387 6. Yu G, Greninger AL, Isa P, Phan TG, Martinez MA, de la Luz Sanchez M, Contreras JF, Santos-Preciado
- 388 JI, Parsonnet J, Miller S, DeRisi JL, Delwart E, Arias CF, Chiu CY (2012) Discovery of a novel polyomavirus
- 389 in acute diarrheal samples from children. *PLoS ONE* 7 (11):e49449. doi:10.1371/journal.pone.0049449
- 390 7. Lim ES, Reyes A, Antonio M, Saha D, Ikumapayi UN, Adeyemi M, Stine OC, Skelton R, Brennan DC,
- 391 Mkakosya RS, Manary MJ, Gordon JI, Wang D (2013) Discovery of STL polyomavirus, a polyomavirus of
- 392 ancestral recombinant origin that encodes a unique T antigen by alternative splicing. *Virology* 436
- 393 (2):295-303. doi:10.1016/j.virol.2012.12.005
- 394 8. Yang J, Yang F, Ren L, Xiong Z, Wu Z, Dong J, Sun L, Zhang T, Hu Y, Du J, Wang J, Jin Q (2011) Unbiased
- 395 parallel detection of viral pathogens in clinical samples by use of a metagenomic approach. *J Clin*
- 396 *Microbiol* 49 (10):3463-3469. doi:10.1128/JCM.00273-11
- 397 9. Lysholm F, Wetterbom A, Lindau C, Darban H, Bjerkner A, Fahlander K, Lindberg AM, Persson B,
- 398 Allander T, Andersson B (2012) Characterization of the viral microbiome in patients with severe lower
- 399 respiratory tract infections, using metagenomic sequencing. *PLoS ONE* 7 (2):e30875.
- 400 doi:10.1371/journal.pone.0030875
- 401 10. Bal A, Pichon M, Picard C, Casalegno JS, Valette M, Schuffenecker I, Billard L, Vallet S, Vilchez G,
- 402 Cheynet V, Oriol G, Trouillet-Assant S, Gillet Y, Lina B, Brengel-Pesce K, Morfin F, Josset L (2018) Quality
- 403 control implementation for universal characterization of DNA and RNA viruses in clinical respiratory
- 404 samples using single metagenomic next-generation sequencing workflow. *BMC Infect Dis* 18 (1):537.
- 405 doi:10.1186/s12879-018-3446-5
- 406 11. Graf EH, Simmon KE, Tardif KD, Hymas W, Flygare S, Eilbeck K, Yandell M, Schlaberg R (2016)
- 407 Unbiased Detection of Respiratory Viruses by Use of RNA Sequencing-Based Metagenomics: a
- 408 Systematic Comparison to a Commercial PCR Panel. *J Clin Microbiol* 54 (4):1000-1007.
- 409 doi:10.1128/JCM.03060-15
- 410 12. Schlaberg R, Queen K, Simmon K, Tardif K, Stockmann C, Flygare S, Kennedy B, Voelkerding K,
- 411 Bramley A, Zhang J, Eilbeck K, Yandell M, Jain S, Pavia AT, Tong S, Ampofo K (2017) Viral Pathogen
- 412 Detection by Metagenomics and Pan-Viral Group Polymerase Chain Reaction in Children With
- 413 Pneumonia Lacking Identifiable Etiology. *J Infect Dis* 215 (9):1407-1415. doi:10.1093/infdis/jix148
- 414 13. Xu L, Zhu Y, Ren L, Xu B, Liu C, Xie Z, Shen K (2017) Characterization of the Nasopharyngeal Viral
- 415 Microbiome from Children with Community-Acquired Pneumonia but Negative for Luminex xTAG
- 416 Respiratory Viral Panel Assay Detection. *J Med Virol*. doi:10.1002/jmv.24895
- 417 14. Lewandowska DW, Schreiber PW, Schuurmans MM, Ruehe B, Zagordi O, Bayard C, Greiner M,
- 418 Geissberger FD, Capaul R, Zbinden A, Boni J, Benden C, Mueller NJ, Trkola A, Huber M (2017)

- 419 Metagenomic sequencing complements routine diagnostics in identifying viral pathogens in lung
420 transplant recipients with unknown etiology of respiratory infection. *PLoS ONE* 12 (5):e0177340.
421 doi:10.1371/journal.pone.0177340
- 422 15. Zinter MS, Dvorak CC, Mayday MY, Iwanaga K, Ly NP, McGarry ME, Church GD, Faricy LE, Rowan CM,
423 Hume JR, Steiner ME, Crawford ED, Langelier C, Kalantar K, Chow ED, Miller S, Shimano K, Melton A,
424 Yanik GA, Sapru A, DeRisi JL (2018) Pulmonary Metagenomic Sequencing Suggests Missed Infections in
425 Immunocompromised Children. *Clin Infect Dis*. doi:10.1093/cid/ciy802
- 426 16. Blauwkamp TA, Thair S, Rosen MJ, Blair L, Lindner MS, Vilfan ID, Kawli T, Christians FC,
427 Venkatasubrahmanyam S, Wall GD, Cheung A, Rogers ZN, Meshulam-Simon G, Huijse L, Balakrishnan S,
428 Quinn JV, Hollemon D, Hong DK, Vaughn ML, Kertesz M, Bercovici S, Wilber JC, Yang S (2019) Analytical
429 and clinical validation of a microbial cell-free DNA sequencing test for infectious disease. *Nature*
430 *Microbiology*. doi:10.1038/s41564-018-0349-6
- 431 17. Gu W, Miller S, Chiu CY (2019) Clinical Metagenomic Next-Generation Sequencing for Pathogen
432 Detection. *Annu Rev Pathol* 14:319-338. doi:10.1146/annurev-pathmechdis-012418-012751
- 433 18. Parize P, Muth E, Richaud C, Gratigny M, Pilmis B, Lamamy A, Mainardi JL, Cheval J, de Visser L,
434 Jagorel F, Ben Yahia L, Bamba G, Dubois M, Join-Lambert O, Leruez-Ville M, Nassif X, Lefort A, Lanternier
435 F, Suarez F, Lortholary O, Lecuit M, Eloit M (2017) Untargeted next-generation sequencing-based first-
436 line diagnosis of infection in immunocompromised adults: a multicentre, blinded, prospective study. *Clin*
437 *Microbiol Infect* 23 (8):574 e571-574 e576. doi:10.1016/j.cmi.2017.02.006
- 438 19. Schlaberg R, Chiu CY, Miller S, Procop GW, Weinstock G, Professional Practice C, Committee on
439 Laboratory Practices of the American Society for M, Microbiology Resource Committee of the College of
440 American P (2017) Validation of Metagenomic Next-Generation Sequencing Tests for Universal
441 Pathogen Detection. *Arch Pathol Lab Med* 141 (6):776-786. doi:10.5858/arpa.2016-0539-RA
- 442 20. Wylie KM, Mihindukulasuriya KA, Zhou Y, Sodergren E, Storch GA, Weinstock GM (2014)
443 Metagenomic analysis of double-stranded DNA viruses in healthy adults. *BMC Biol* 12:71.
444 doi:10.1186/s12915-014-0071-7
- 445 21. Goggin RK, Bennett CA, Bassiouni A, Bialasiewicz S, Vreugde S, Wormald PJ, Psaltis AJ (2018)
446 Comparative Viral Sampling in the Sinonasal Passages; Different Viruses at Different Sites. *Frontiers in*
447 *cellular and infection microbiology* 8:334. doi:10.3389/fcimb.2018.00334
- 448 22. Bogaert D, Keijser B, Huse S, Rossen J, Veenhoven R, van Gils E, Bruin J, Montijn R, Bonten M,
449 Sanders E (2011) Variability and diversity of nasopharyngeal microbiota in children: a metagenomic
450 analysis. *PLoS ONE* 6 (2):e17035. doi:10.1371/journal.pone.0017035
- 451 23. Fry AM, Lu X, Olsen SJ, Chittaganpitch M, Sawatwong P, Chantra S, Baggett HC, Erdman D (2011)
452 Human rhinovirus infections in rural Thailand: epidemiological evidence for rhinovirus as both pathogen
453 and bystander. *PLoS ONE* 6 (3):e17780. doi:10.1371/journal.pone.0017780
- 454 24. Iwane MK, Prill MM, Lu X, Miller EK, Edwards KM, Hall CB, Griffin MR, Staat MA, Anderson LJ,
455 Williams JV, Weinberg GA, Ali A, Szilagyi PG, Zhu Y, Erdman DD (2011) Human rhinovirus species
456 associated with hospitalizations for acute respiratory illness in young US children. *J Infect Dis* 204
457 (11):1702-1710. doi:10.1093/infdis/jir634
- 458 25. Singleton RJ, Bulkow LR, Miernyk K, DeByle C, Pruitt L, Hummel KB, Bruden D, Englund JA, Anderson
459 LJ, Lucher L, Holman RC, Hennessy TW (2010) Viral respiratory infections in hospitalized and community
460 control children in Alaska. *J Med Virol* 82 (7):1282-1290. doi:10.1002/jmv.21790
- 461 26. van Bentem I, Koopman L, Niesters B, Hop W, van Middelkoop B, de Waal L, van Drunen K, Osterhaus
462 A, Neijens H, Fokkens W (2003) Predominance of rhinovirus in the nose of symptomatic and
463 asymptomatic infants. *Pediatric allergy and immunology : official publication of the European Society of*
464 *Pediatric Allergy and Immunology* 14 (5):363-370

- 465 27. Kapoor A, Mehta N, Dubovi EJ, Simmonds P, Govindasamy L, Medina JL, Street C, Shields S, Lipkin WI
466 (2012) Characterization of novel canine bocaviruses and their association with respiratory disease. *J Gen*
467 *Virol* 93 (Pt 2):341-346. doi:10.1099/vir.0.036624-0
- 468 28. Piewbang C, Jo WK, Puff C, Ludlow M, van der Vries E, Banlunara W, Rungsipipat A, Kruppa J, Jung K,
469 Techangamsuwan S, Baumgartner W, Osterhaus A (2018) Canine Bocavirus Type 2 Infection Associated
470 With Intestinal Lesions. *Vet Pathol* 55 (3):434-441. doi:10.1177/0300985818755253
- 471 29. Bodewes R, Lapp S, Hahn K, Habierski A, Forster C, Konig M, Wohlsein P, Osterhaus AD, Baumgartner
472 W (2014) Novel canine bocavirus strain associated with severe enteritis in a dog litter. *Vet Microbiol* 174
473 (1-2):1-8. doi:10.1016/j.vetmic.2014.08.025
- 474 30. Niu J, Yi S, Wang H, Dong G, Zhao Y, Guo Y, Dong H, Wang K, Hu G (2019) Complete genome
475 sequence analysis of canine bocavirus 1 identified for the first time in domestic cats. *Arch Virol* 164
476 (2):601-605. doi:10.1007/s00705-018-4096-z
- 477 31. Lau SK, Woo PC, Yeung HC, Teng JL, Wu Y, Bai R, Fan RY, Chan KH, Yuen KY (2012) Identification and
478 characterization of bocaviruses in cats and dogs reveals a novel feline bocavirus and a novel genetic
479 group of canine bocavirus. *J Gen Virol* 93 (Pt 7):1573-1582. doi:10.1099/vir.0.042531-0
- 480 32. McElvania TeKippe E, Wylie KM, Deych E, Sodergren E, Weinstock G, Storch GA (2012) Increased
481 prevalence of anellovirus in pediatric patients with fever. *PLoS ONE* 7 (11):e50937.
482 doi:10.1371/journal.pone.0050937
- 483 33. Young JC, Chehoud C, Bittinger K, Bailey A, Diamond JM, Cantu E, Haas AR, Abbas A, Frye L, Christie
484 JD, Bushman FD, Collman RG (2015) Viral metagenomics reveal blooms of anelloviruses in the
485 respiratory tract of lung transplant recipients. *American journal of transplantation : official journal of the*
486 *American Society of Transplantation and the American Society of Transplant Surgeons* 15 (1):200-209.
487 doi:10.1111/ajt.13031
- 488 34. De Vlaminck I, Khush KK, Strehl C, Kohli B, Luikart H, Neff NF, Okamoto J, Snyder TM, Cornfield DN,
489 Nicolls MR, Weill D, Bernstein D, Valentine HA, Quake SR (2013) Temporal response of the human
490 virome to immunosuppression and antiviral therapy. *Cell* 155 (5):1178-1187.
491 doi:10.1016/j.cell.2013.10.034
- 492 35. Blatter JA, Sweet SC, Conrad C, Danziger-Isakov LA, Faro A, Goldfarb SB, Hayes D, Jr., Melicoff E,
493 Schechter M, Storch G, Visner GA, Williams NM, Wang D (2018) Anellovirus loads are associated with
494 outcomes in pediatric lung transplantation. *Pediatr Transplant* 22 (1). doi:10.1111/petr.13069
- 495 36. Li L, Deng X, Linsuwanon P, Bangsberg D, Bwana MB, Hunt P, Martin JN, Deeks SG, Delwart E (2013)
496 AIDS alters the commensal plasma virome. *J Virol* 87 (19):10912-10915. doi:10.1128/JVI.01839-13
- 497 37. Sherman KE, Rouster SD, Feinberg J (2001) Prevalence and genotypic variability of TTV in HIV-
498 infected patients. *Digestive diseases and sciences* 46 (11):2401-2407
- 499 38. Touinssi M, Gallian P, Biagini P, Attoui H, Vialettes B, Berland Y, Tamalet C, Dhiver C, Ravaux I, De
500 Micco P, De Lamballerie X (2001) TT virus infection: prevalence of elevated viraemia and arguments for
501 the immune control of viral load. *Journal of clinical virology : the official publication of the Pan American*
502 *Society for Clinical Virology* 21 (2):135-141
- 503 39. Shibayama T, Masuda G, Ajisawa A, Takahashi M, Nishizawa T, Tsuda F, Okamoto H (2001) Inverse
504 relationship between the titre of TT virus DNA and the CD4 cell count in patients infected with HIV. *Aids*
505 15 (5):563-570
- 506 40. Thom K, Petrik J (2007) Progression towards AIDS leads to increased Torque teno virus and Torque
507 teno minivirus titers in tissues of HIV infected individuals. *J Med Virol* 79 (1):1-7. doi:10.1002/jmv.20756
- 508 41. Christensen JK, Eugen-Olsen J, M SL, Ullum H, Gjedde SB, Pedersen BK, Nielsen JO, Krosgaard K
509 (2000) Prevalence and prognostic significance of infection with TT virus in patients infected with human
510 immunodeficiency virus. *J Infect Dis* 181 (5):1796-1799. doi:10.1086/315440

- 511 42. Siqueira JD, Dominguez-Bello MG, Contreras M, Lander O, Caballero-Arias H, Xutao D, Noya-Alarcon
512 O, Delwart E (2018) Complex virome in feces from Amerindian children in isolated Amazonian villages.
513 *Nat Commun* 9 (1):4270. doi:10.1038/s41467-018-06502-9
- 514 43. Victoria JG, Kapoor A, Li L, Blinkova O, Slikas B, Wang C, Naeem A, Zaidi S, Delwart E (2009)
515 Metagenomic analyses of viruses in stool samples from children with acute flaccid paralysis. *J Virol* 83
516 (9):4642-4651. doi:10.1128/JVI.02301-08
- 517 44. Li L, Deng X, Mee ET, Collot-Teixeira S, Anderson R, Schepelmann S, Minor PD, Delwart E (2015)
518 Comparing viral metagenomics methods using a highly multiplexed human viral pathogens reagent. *J*
519 *Virol Methods* 213:139-146. doi:10.1016/j.jviromet.2014.12.002
- 520 45. Deng X, Naccache SN, Ng T, Federman S, Li L, Chiu CY, Delwart EL (2015) An ensemble strategy that
521 significantly improves de novo assembly of microbial genomes from metagenomic next-generation
522 sequencing data. *Nucleic Acids Res.* doi:10.1093/nar/gkv002
- 523 46. S. NMaK (2000) *Molecular Evolution and Phylogenetics*. Oxford University Press, New York.
- 524 47. Tamura K, Stecher G, Peterson D, Filipski A, Kumar S (2013) MEGA6: Molecular Evolutionary Genetics
525 Analysis version 6.0. *Mol Biol Evol* 30 (12):2725-2729. doi:10.1093/molbev/mst197

526

527

528