

1 **Short Note**

2 **The anomalous nature of the fecal swab sample used for RaTG13**
3 **genome assembly as revealed by NGS data analysis**

4

5

6 **Monali C. Rahalkar^{1*} and Rahul A. Bahulikar²**

7 ¹C2, Bioenergy group, MACS Agharkar Research Institute, G.G. Agarkar Road,

8 Pune 411004, Maharashtra, India

9 ²BAIF Development Research Foundation, Central Research Station,

10 Urulikanchan, Pune 412202

11 *Corresponding author: monalirahalkar@aripune.org

12

13

14

15

16 Abstract:

17 RaTG13, a SARS-like beta coronavirus, which exists in the form of a genome sequence, is
18 the closest relative of SARS-CoV-2 reported till date. The sample from which RaTG13 virus
19 was sequenced was a bat fecal swab collected in 2013 from Tongguan, Mojiang, Yunnan
20 province, China. The genome data for RaTG13, MN996532.1, was deposited on 27th Jan
21 2020 and the raw data (Illumina reads) was deposited a fortnight later on 13th Feb 2020
22 [https://www.ncbi.nlm.nih.gov/sra/SRX7724752\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX7724752[accn]). Comparison of the RNA Seq data of
23 RaTG13 fecal swab sample to the corresponding data from the bat fecal swabs deposited by
24 the same working group indicated that the raw data seemed to be anomalous in several
25 aspects. Thirty percent of the reads did not match with anything. From the rest of the 70%, an
26 abnormal high proportion was contributed by reads derived from eukaryotes (~68%). These
27 matched with the sequences of not one but various bat species (round leaf bats, fruit bats and
28 other bats) and animal species (squirrels, foxes, etc.) as per Krona analysis included with the
29 SRA data. The proportion of the bacterial reads in the swab was exceptionally low, i.e. 0.7%,
30 which is abnormal, compared to the 70-90% bacterial abundance in other bat fecal swabs.
31 Furthermore, we also found another set of raw data associated with RaTG13, amplicon
32 sequencing of the genome (SRX8357956), which was submitted in May 2020. Analysis of
33 the amplicons by BLAST showed that these collectively do not cover the whole genome
34 (MN996532.1). On closer inspection, the dates mentioned in the files of the sequenced
35 amplicons were also found to be older (2017, 2018). Collectively, the anomalies in the raw
36 data of RaTG13 pose an important question about the overall authenticity of the RaTG13
37 genome sequence.

38 **Key words:** RaTG13; SARS-COV-2; Illumina sequencing, amplicon sequencing, NGS; fecal
39 swab

40

41 Covid-19 has been a devastating pandemic affecting more than nineteen million people and
42 killing about three quarter million people till date (8th August 2020). SARS-CoV2, the virus
43 responsible for the pandemic is most similar to RaTG13 (a bat derived coronavirus) on the
44 genomic level. RaTG13 has been known as the sister virus of SARS-CoV-2 as its shows the
45 closest overall genomic similarity (96.2%) to SARS-CoV-2 genome (Zhou et al., 2020).
46 RaTG13 has been used for various comparative experiments with SARS-CoV-2. These
47 include: the capacity of its spike to bind to human ACE-2 and the capacity to cause human
48 infections (Wrobel et al., 2020), the evolutionary analysis and prediction of a common
49 ancestor of SARS-CoV-2 and RaTg13 (Boni et al., 2020), and many more upcoming papers.

50 RaTG13 is a beta SARS-like corona virus and the name was introduced to us in 2020 (Zhou
51 et al., 2020). The sequence of RaTG13 was retrieved by RNA sequencing using a next
52 generation sequencing approach, though the sample was collected in 2013 (Zhou et al.,
53 2020). The RNA sample was that of a bat fecal swab collected in July 2013, from Tongguan
54 mineshaft in Yunnan. The details of the location were predicted earlier (Arbuthnott et al.,
55 2020, Rahalkar and Bahulikar, 2020). However, in a recent reply to the Science magazine it
56 has been clarified by Dr. Zheng-Li Shi that the TG in RaTG13 is for Tongguan, Mojiang,
57 Yunna, China (Cohen, 2020). She has also confirmed that the old name of RaTG13 virus was
58 BtCoV/CoV4991, described earlier (Ge et al., 2016). However, the sample appears to be
59 finished or disintegrated or not available to the scientific community as per a recent media
60 investigation (Arbuthnott et al., 2020). Also, Zheng-Li Shi has confirmed that her lab has
61 never cultured this virus and it is not in live condition in her lab (Cohen, 2020). Therefore,
62 the entire scientific community has to rely on the RaTG13 genome available. And if the
63 RaTG13 genome sequence is to be used in all the bioinformatics and model experiments, the
64 pre-supposition is that the sequence of this virus should be accurate and based on a reliable
65 raw data. Therefore, we have looked at the RaTG13 raw data in a closer manner in this paper.

66 The name RaTG13 was first mentioned in 2020 (Zhou et al., 2020) and the full genome
67 sequence was not available before January 2020 on any of the databases, to the best of our
68 knowledge. The Illumina based NGS sequence of RaTG13 **MN996532.1** was deposited on
69 27th Jan 2020 and the raw data was available a little later on 13th Feb 2020
70 [https://www.ncbi.nlm.nih.gov/sra/SRX7724752\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX7724752[accn]).

71 The older name of RaTG13 was BtCoV/4991, where the names were based on RdRP
72 sequences and sample numbers (Cohen, 2020). A 370 base RdRp fragment of BtCoV/4991
73 (KP378696.1) or RaTG13 shows the highest similarity to SARS-CoV-2 with only 3-5 bases
74 difference (blast comparison with the SARS-CoV-2 sequences deposited till date). Also,
75 BtCoV/4991 or RaTG13 has a great significance as it had been recovered from the same
76 location, a mineshaft in Tongguan, where six miners were afflicted with a suspiciously
77 COVID-19 like pneumonia in 2012, and three succumbed to the infection and died
78 (Arbuthnott et al., 2020, Rahalkar and Bahulikar, 2020). Thus, BtCoV/4991 or RaTG13 is
79 also the first and the only beta SARS-like CoV known so far associated with Tongguan
80 mineshaft where lethal human pneumonia cases were reported in 2012 (Arbuthnott et al.,
81 2020, Rahalkar and Bahulikar, 2020).

82 **Problems seen in the RAW DATA of RaTG13: Illumina sequence SRX7724752**

83 **Here are the basic discrepancies encountered after the analysis of the RaTG13 fecal**
84 **swab Illumina data** [https://www.ncbi.nlm.nih.gov/sra/SRX7724752\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX7724752[accn]):

85 1. The genome of RaTG13 (MN996532.1) is derived from a fecal swab sample collected in
86 2013 as per the description. However, the Illumina sequencing entry of SRX7724752, the
87 sample is recorded as being extracted from a BAL fluid (broncho alveolar lavage) (Fig. 1a).

88 2. The total raw data is 3.3 gigabases (Fig. 1b). After the Krona analysis it is seen that ~30%
89 reads are unidentified (no matches) and only ~ 70% reads are identified. Out of the 70%, a
90 vast majority i.e. 68% is seen to be contributed by eukaryotic sequences (Fig. 1b). This is
91 highly unusual as it is a fecal swab and the analysis of other bat fecal or anal swabs usually
92 do not show such high proportion of eukaryotic RNA.

93 3. Within the 68% of the eukaryotic sequences, the bat sequences are about 36-40% (Fig 1a.).
94 The rest of the 30% sequences are contributed by squirrels (*Marmota*), flying foxes, foxes,
95 and other types of animals (Fig.1 b). In a bat fecal swab, the majority of the eukaryotic RNA
96 should be arising from the same species, e.g. in this case it should be of *Rhinoplohus affinis*.
97 Given the fact that *Rhinolophus affinis* bat has not been sequenced completely, it could show
98 similarities with other bat species, such as *Hipposideros* (round leaf bats) or *Rousettus*
99 (egyptian fruit bats). However, the rest of the hits are from totally unrelated taxa, such as
100 *Marmota* (squirrels), *Vulpus* (foxes), *Pteropus* (mega-bats), etc. Similar discrepancies in the
101 raw data have been pointed out recently (Zhang, 2020).

102 4. Another major discrepancy is that the RNA sequencing data shows extremely less
103 abundance of bacteria, only 0.65%. This is far too less in comparison with other fecal or anal
104 swab of bats, which show a very high proportion of bacterial sequences ~76-90% (Fig.2 and
105 Fig.3). SRA data of six other bat fecal swabs submitted by the same group also showed a
106 high abundance of bacterial reads (details not shown). Bacteria are usually the highest
107 constituents of gut flora and hence contribute to a high extent to a fecal sample.

108 5. The coronavirus sequence (RaTG13) contributes to ~0.003% of the total sequence reads.
109 These raw reads were used to build an almost complete assembly, though the overall
110 coverage is less ~7-8X. Though there were less overlaps in a few regions, there are only 3
111 gaps as per our analysis using RaTG13 as the reference genome. The Wuhan Institute of

112 Virology has recently described methods like probe-capture for getting the whole genome of
113 viruses from samples like bat feces (Li et al 2019). In this case, without the use of any other
114 methods, and after using a seven year old fecal swab or fecal swab RNA it is surprising that
115 how the viral reads were of a better quality.

116 6. The assembly method and the actual assembly accession for RaTG13 is not described or
117 linked to the whole genome of RaTG13, i.e. MN669532.1 and also no assembly method is
118 specified in the raw data SRX7724752. Also, no assembly data accession number is available
119 for RaTG13 genome as per our information and searches.

120 7. After blasting the RaTG13 genome against the SRA, ~1700 reads can be retrieved which
121 covers only a small portion, i.e. 252 kb of the total 3.3 Gbases. The genome size of RaTG13
122 is known to be ~30 kb. According to our knowledge, this is ~8x coverage is low and may be
123 insufficient to arrive at a definitive genome assembly.

124 8. We also compared the fecal/anal swab RNA Seq data deposited for the same bat species,
125 i.e. *Rhinolophus affinis* (Fig.2) and fecal swab from another bat (Fig. 3). It is clearly seen that
126 the other two swabs showed normal findings, with 70-90% bacterial reads and very few reads
127 associated with the host. Also these swabs do not show sequences affiliated with other
128 animals.

129 9. Similar findings have been documented in a latest preprint by Zhang, D. (Zhang, 2020)
130 <https://zenodo.org/record/3969272#.Xypwfn5S-Un>.

131 **Problems in the Amplicon sequencing data:**

132 We found that some amplicon sequencing data for RaTG13 (SRX8357956) was submitted in
133 May 2020.

134 1. No indications of amplicon sequencing has been given by Zhou et al 2020, or in any of the
135 recent publications by the WIV workgroup.

136 2. A total of 33 spots are present in the raw data (Fig.4). The sequencing file names indicate
137 that the dates are from 2017 and 2018. However, the submission has been done in May 2020.

138 3. There are two contrasting sequences for a single patch (spot 23 and spot 24), e.g. shows
139 94-96% similarity to that of MN669532.1. However, two spots (22 and 25) covering the same
140 area showed 99% similarity to the described RaTG13 consensus MN669532.1.

141 4. In general, most of the amplicons showed 97-99% similarity with that of MN669532.1.
142 However, collectively, the spots do not cover the entire genome and major gaps are seen in
143 various regions.

144 5. RdRp derived from the amplicon sequencing is incomplete (spots 31 and 32) and does not
145 match with RdRp of BtCoV/4991 KP876546.1. Around 170 bases from 370 base sequences
146 are missing and it shows 2 base mismatches.

147 **Conclusions:**

148 a. Our main grievance is that the fecal swab from which RaTG13 sequence has been derived
149 appears as an anomalous fecal swab as pointed above with respect to its community
150 composition. The swab shows 70% of eukaryotic sequences from sources which should not
151 have been detected in *Rhinolophus* bat feces such as mexican bats, squirrels, flying foxes, red
152 foxes, etc. And most importantly, there is extreme low representation of bacteria. Bacteria
153 constitute a major part of feces from any eukaryotic organism.

154

155 b. RaTG13 genome sequence has been used extensively for in various evolutionary
156 calculations, simulation experiments and would be used in future for bioinformatics and
157 experimental comparisons. And therefore, all the data associated with RaTG13 should be
158 inspected properly.

159 c. The reads from which the viral sequence of RaTG13 was assembled appears not to be
160 affected by the anomalous nature of the RNA Seq data of the fecal swab. An almost complete
161 assembly is assumed to have been built from this raw data (Illumina reads). An important
162 question is how did this considerably good data related to the virus come from this swab?
163 And if it was not degraded, what are the reasons of its anomalous composition?

164 d. The amplicon data is incomplete and does not help us in further confirming the RaTG13
165 whole genome sequence.

166 e. Considering the anomalous nature of the raw data presented for RaTG13 (both Illumina
167 and amplicon sequence) it would be a real question can the scientific community rely on the
168 integrity of the RaTG13 genome sequence MN996532.1? Moreover, with the discrepancies
169 pointed out in the raw data of RaTG13, we suggest that RaTG13 genome sequence should be
170 interpreted with caution.

171

172

173 **Figures:**174 **Fig.1** RNA-Seq of *Rhinolophus affinis*:Fecal swabTaxonomy Analysis (RaTG13)

Full ▾ Send to: ▾

SRX7724752: RNA-Seq of *Rhinolophus affinis*:Fecal swab
1 ILLUMINA (Illumina HiSeq 3000) run: 11.6M spots, 3.3G bases, 1.7Gb downloads

Design: Total RNA was extracted from bronchoalveolar lavage fluid using the QIAamp Viral RNA Mini Kit following the manufacturers instructions. An RNA library was then constructed using the TruSeq Stranded mRNA Library Preparation Kit (Illumina, USA). Paired-end (150 bp) sequencing of the RNA library was performed on the HiSeq 3000 platform (Illumina).

Submitted by: Wuhan Institute of Virology, Chinese Academy of Sciences

Study: Bat coronavirus RaTG13 Genome sequencing
[PRJNA606165](#) • [SRP249482](#) • [All experiments](#) • [All runs](#)
[show Abstract](#)

Sample:
[SAMN14082201](#) • [SRS6146537](#) • [All experiments](#) • [All runs](#)
Organism: [unidentified coronavirus](#)

Library:
Name: RaTG13
Instrument: Illumina HiSeq 3000
Strategy: RNA-Seq
Source: METAGENOMIC
Selection: RANDOM
Layout: PAIRED

Runs: 1 run, 11.6M spots, 3.3G bases, [1.7Gb](#)

Run	# of Spots	# of Bases	Size	Published
SRR11085797	11,604,666	3.3G	1.7Gb	2020-02-13

ID: 10102765

175

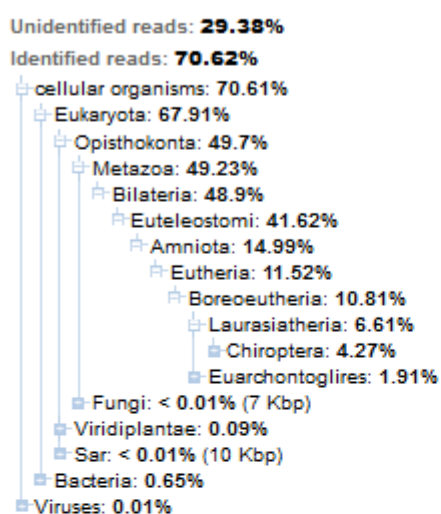
176 **Fig1a.** RNA-Seq of *Rhinolophus affinis*:Fecal swab (**RaTG13**)

177

RNA-Seq of *Rhinolophus affinis*:Fecal swab (SRR11085797)

Metadata Analysis Reads Data access

Taxonomy Analysis



[View in Krona](#)

Strong signals

SuperKingdom	Organism	Rank	%%	Kbp	Coverage
Eukaryota	<i>Hipposideros armiger</i>	species	31.8	1,048,945	
Eukaryota	<i>Rousettus aegyptiacus</i>	species	4.6	151,010	
Eukaryota	<i>Marmota marmota marmota</i>	subspecies	4.6	150,069	
Eukaryota	<i>Vulpes vulpes</i>	species	4.0	131,805	
Eukaryota	<i>Marmota flaviventris</i>	species	3.6	118,361	
Eukaryota	<i>Pteropus</i>	genus	3.0	100,495	
Eukaryota	Odontoceti	parvorder	3.0	98,516	
Eukaryota	<i>Myotis</i>	genus	3.0	97,335	
Eukaryota	Miniopterinae	subfamily	2.6	86,145	

178

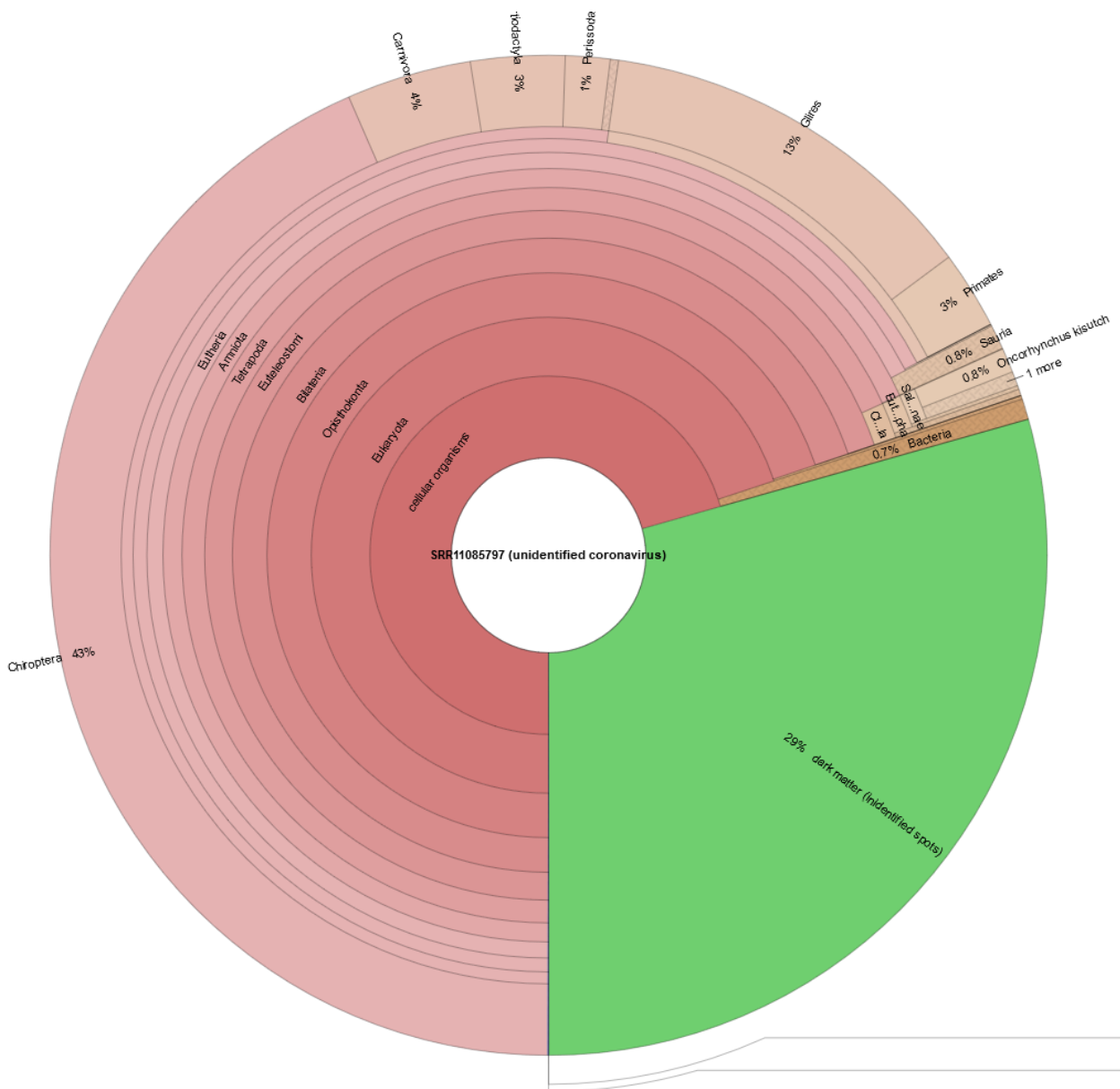
179

180

181 Fig. 1b. Distribution of the reads in the raw data. The individual distribution is given and in the
 182 second part, the reads with strong signals, i.e. which contribute to a higher extent are given in
 183 decreasing order.

184

185



186

187

188 Fig.1 c. Krona chart of RaTG13 raw data, 29% unidentified reads, 43% Chiroptera, 13% Gileres, 3%

189 Primates, 0.7% bacteria and 0.024% RaTG13 reads

190

191 Fig 2. RNA-Seq of Rhinolophus affinis: Fecal swab Taxonomy Analysis

192 [https://www.ncbi.nlm.nih.gov/sra/SRX7724693\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX7724693[accn])

Full ▾

Send to: ▾

SRX7724693: RNA-Seq of Rhinolophus affinis: Anal swab

1 ILLUMINA (Illumina HiSeq 3000) run: 11.9M spots, 3.5G bases, 1.6Gb downloads

Design: Total RNA was extracted from bronchoalveolar lavage fluid using the QIAamp Viral RNA Mini Kit following the manufacturers instructions. An RNA library was then constructed using the TruSeq Stranded mRNA Library Preparation Kit (Illumina, USA). Paired-end (150 bp) sequencing of the RNA library was performed on the HiSeq 3000 platform (Illumina).

Submitted by: Wuhan Institute of Virology, Chinese Academy of Sciences

Study: Discovery of Bat Coronaviruses through Surveillance and Probe Capture-Based Next-Generation Sequencing.

[PRJNA606159](#) • [SRP249478](#) • [All experiments](#) • [All runs](#)

[show Abstract](#)

Sample:

[SAMN14086235](#) • [SRS6146479](#) • [All experiments](#) • [All runs](#)

Organism: [unclassified Rhinacovirus](#)

Library:

Name: 160660

Instrument: Illumina HiSeq 3000

Strategy: RNA-Seq

Source: METAGENOMIC

Selection: RANDOM

Layout: PAIRED

Runs: 1 run, 11.9M spots, 3.5G bases, [1.6Gb](#)

Run	# of Spots	# of Bases	Size	Published
SRR11085736	11,924,182	3.5G	1.6Gb	2020-02-13

ID: 10102706

193

194 **Fig. 2a. RNA-Seq of Rhinolophus affinis: Anal swab (SRR11085736)**

195

Taxonomy Analysis

Unidentified reads: **0.86%**

Identified reads: **99.14%**

cellular organisms: 99.11%

Bacteria: 91.07%

Eukaryota: 4.36%

Viruses: 0.03%

[View in Krona](#)

Strong signals

SuperKingdom	Organism	Rank	%%	Kbp	Coverage
Bacteria	Clostridium	genus	37.3	1,288,845	
Bacteria	Niameybacter massiliensis	species	24.6	849,347	
Bacteria	Pasteurellaceae	family	11.7	404,812	
Bacteria	Clostridioides difficile	species	5.8	199,353	47.6
Eukaryota	Boreoeutheria		4.2	145,969	
Bacteria	Romboutsia lituseburensis	species	3.7	126,405	
Bacteria	Escherichia coli	species	3.2	110,843	21.5
Bacteria	Paenibacillus	genus	1.4	47,848	
Bacteria	Helicobacter	genus	1.1	38,581	
Bacteria	Paeniclostridium sordellii	species	0.8	28,640	8.2
Bacteria	Enterococcus faecalis	species	0.4	14,079	4.7
Bacteria	Staphylococcus aureus	species	0.3	11,072	3.9
Bacteria	Enterococcus faecium	species	0.3	10,030	3.4

196

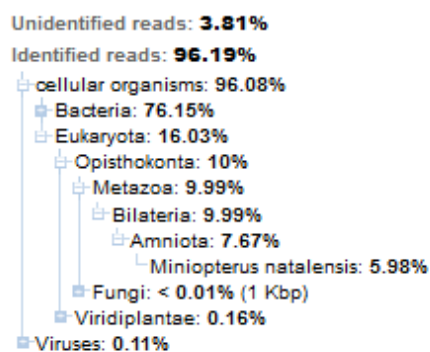
197 Fig. 2b. Distribution of the reads in the raw data. The individual distribution is given and in the
 198 second part, the reads which contribute to a higher extent are given.

199

204 **Fig 3** RNA-Seq of *Miniopterus schreibersii*: Fecal swab Taxonomy AnalysisRNA-Seq of *Miniopterus schreibersii*: Anal swab (SRR11085734)

Metadata Analysis Reads Data access

Taxonomy Analysis

[View in Krona](#)

Strong signals

SuperKingdom	Organism	Rank	%%	Kbp	Coverage
Bacteria	Paenibacillus	genus	77.2	2,124,474	
Eukaryota	<i>Miniopterus natalensis</i>	species	16.3	449,835	
Bacteria	unclassified Paenibacillus		5.9	162,892	
Bacteria	Paenibacillus sp. FSL R5-0345	species	1.6	44,539	
Bacteria	Paenibacillus odorifer	species	1.2	33,308	4.8
Bacteria	Clostridium perfringens	species	0.8	22,508	6.3
Bacteria	unclassified Massilia		0.5	14,457	
Bacteria	Mycoplasma	genus	0.2	5,115	
Bacteria	<i>Kluyvera ascorbata</i>	species	0.2	4,588	

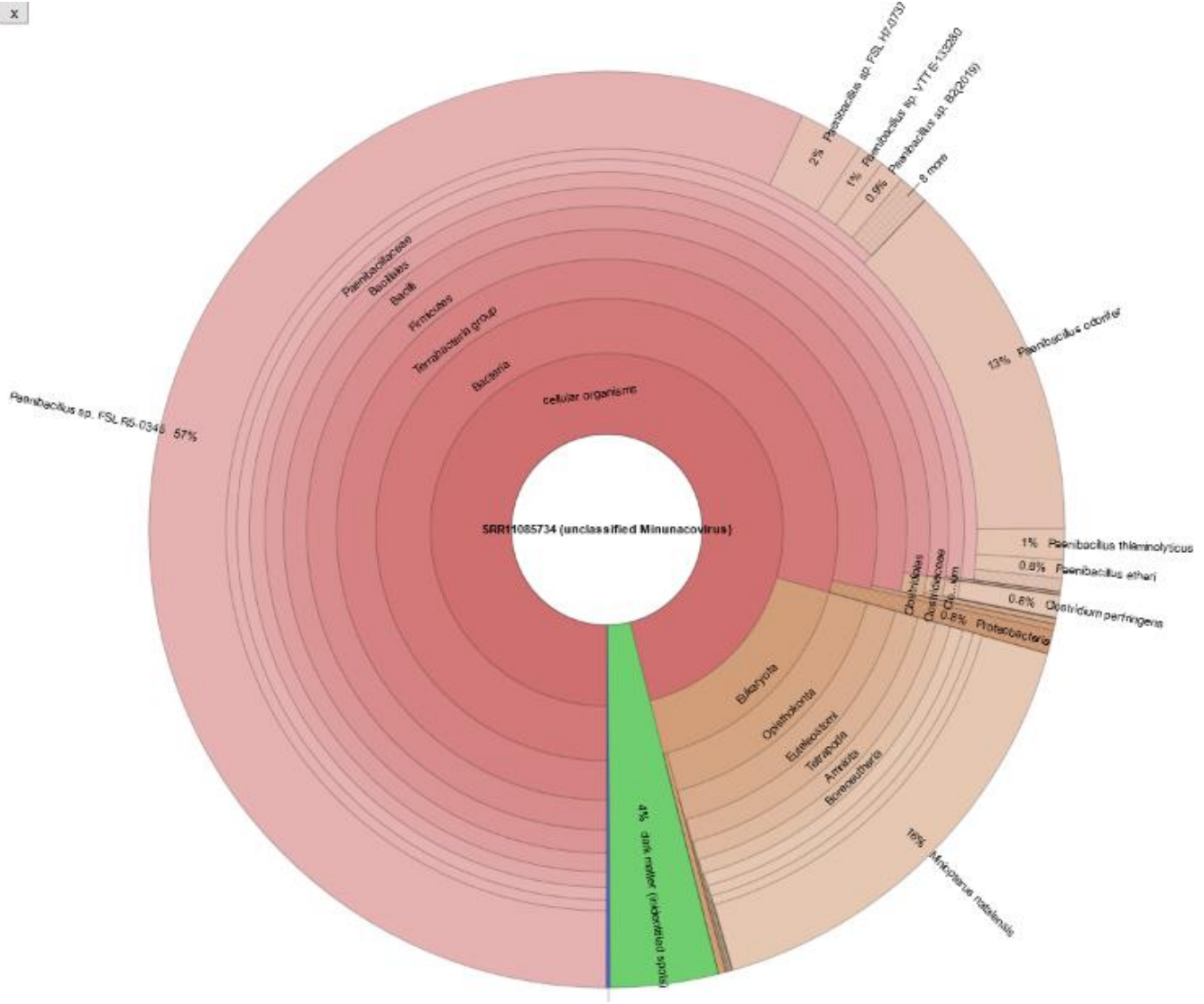
205

206 Fig. 3a. RNA-Seq of fecal swab *Miniopterus schreibersii*

207

208 Fig. 3c. Krona chart of *Miniopterus schreibersii*: Fecal swab Taxonomy

x



209
210

211 Fig. 4

Full ▾

Send to: ▾

SRX8357956: amplicon_sequences of RaTG13

1 CAPILLARY (AB 310 Genetic Analyzer) run: 33 spots, 30,576 bases, 1.1Mb downloads

Design: Primer-based amplicon sequences**Submitted by:** Wuhan Institute of Virology, Chinese Academy of Sciences**Study:** Bat coronavirus RaTG13 Genome sequencing[PRJNA606165](#) • [SRP249482](#) • [All experiments](#) • [All runs](#)[show Abstract](#)**Sample:**[SAMN14082201](#) • [SRS6146537](#) • [All experiments](#) • [All runs](#)**Organism:** [unidentified coronavirus](#)**Library:****Name:** RaTG13_amplicon_sequences**Instrument:** AB 310 Genetic Analyzer**Strategy:** AMPLICON**Source:** METAGENOMIC**Selection:** PCR**Layout:** SINGLE**Runs:** 1 run, 33 spots, 30,576 bases, [1.1Mb](#)

Run	# of Spots	# of Bases	Size	Published
SRR11806578	33	30,576	1.1Mb	2020-05-19

212 ID: 10870921

213

214 **References:**

215 Arbuthnott, George, Calvert, Jonathan & Sherwell, Philip. 2020.
 216 <https://www.thetimes.co.uk/article/seven-year-covid-trail-revealed-l5vxt7jqp>. *The Sunday*
 217 *Times*.

218 Boni, M.F., Lemey, P & Jiang, X. et al. 2020. Evolutionary origins of the SARS-CoV-2 sarbecovirus
 219 lineage responsible for the COVID-19 pandemic. *Nature Microbiology*,
 220 <https://doi.org/10.1038/s41564-020-0771-4>.

221 Cohen, Jon 2020. Wuhan coronavirus hunter Shi Zhengli speaks out. *Science*, 369, 487-488.

222 Ge, X. Y., Wang, N., Zhang, W., Hu, B., Li, B., Zhang, Y. Z., Zhou, J. H., Luo, C. M., Yang, X. L., Wu, L. J.,
 223 Wang, B., Zhang, Y., Li, Z. X. & Shi, Z. L. 2016. Coexistence of multiple coronaviruses in
 224 several bat colonies in an abandoned mineshaft. *Virol. Sin.*, 31, 31-40.

225 Rahalkar, Monali C. & Bahulikar, Rahul A. 2020. Understanding the origin of 'BatCoV RaTG13', a virus
 226 closest to SARS-CoV-2.

227 Wrobel, A.G., Benton, D.J. & Xu, P. et al. 2020. SARS-CoV-2 and bat RaTG13 spike glycoprotein
 228 structures inform on virus evolution and furin-cleavage effects. *Nat. Struct. Mol. Biol.*, 27,
 229 763-767.

230 Zhang, Daoyu 2020. Anomalies in BatCoV/RaTG13 sequencing and provenance.
 231 https://zenodo.org/record/3969272#.Xy0m5jVS_Ix.

232 Zhou, P., Yang, X. L., Wang, X. G., Hu, B., Zhang, L., Zhang, W., Si, H. R., Zhu, Y., Li, B., Huang, C. L.,
 233 Chen, H. D., Chen, J., Luo, Y., Guo, H., Jiang, R. D., Liu, M. Q., Chen, Y., Shen, X. R., Wang, X.,
 234 Zheng, X. S., Zhao, K., Chen, Q. J., Deng, F., Liu, L. L., Yan, B., Zhan, F. X., Wang, Y. Y., Xiao, G.
 235 F. & Shi, Z. L. 2020. A pneumonia outbreak associated with a new coronavirus of probable
 236 bat origin. *Nature*, 579, 270-273.

237