
1 **Meta-transcriptomic analysis reveals the gene expression and novel conserved**
2 **sub-genomic RNAs in SARS-CoV-2 and MERS-CoV**

3 Lin Lyu^{1†}, Ru Feng^{1†}, Mingnan Zhang^{1†}, Qiyu Gong¹, Yinjing Liao², Yanjiao Zhou³, Xiaokui
4 Guo¹, Bing Su¹, Yair Dorsett^{3*}, Lei Chen^{1*}

5

6 1 Shanghai Institute of Immunology, Shanghai Jiao Tong University School of Medicine,
7 Shanghai 200025, China

8 2 College of Chemistry, Sichuan University, Chengdu 610064, China

9 3 Department of Medicine, University of Connecticut Health, Farmington, CT, USA

10

11 * **Correspondence:**

12 Lei Chen: lei.chen@sjtu.edu.cn;

13 Yair Dorsett: dorsett@uchc.edu

14

15 [†] These authors contributed equally.

16

17

18

19

20

21

22

23 **ABSTRACT**

24 **Background:** Fundamental to viral biology is identification and annotation of viral genes and
25 their function. Determining the level of coronavirus gene expression is inherently difficult
26 due to the positive stranded RNA genome and the identification of sub-genomic RNAs
27 (sgRNAs) that are required for expression of most viral genes. In the COVID-19 epidemic so
28 far, few genomic studies have looked at viral sgRNAs and none have systematically
29 examined the sgRNA profiles of large numbers of SARS-CoV2 datasets in conjunction with
30 data for other coronaviruses.

31 **Results:** We developed a bioinformatic pipeline to analyze the sgRNA profiles of
32 coronaviruses and applied it to 588 individual samples from 20 independent studies, covering
33 more than 10 coronavirus species. Our result showed that SARS-CoV, SARS-CoV-2 and
34 MERS-CoV each had a core sgRNA repertoire generated via a canonical mechanism. Novel
35 sgRNAs that encode peptides with evolutionarily conserved structures were identified in
36 several coronaviruses and were expressed *in vitro* and *in vivo*. Two novel peptides may have
37 direct functional relevance to disease, by alluding interferon responses and disrupting IL17E
38 (IL25) signaling. Relevant to coronavirus infectivity and transmission, we also observed that
39 the level of Spike sgRNAs were significantly higher *in-vivo* than *in-vitro*, while the opposite
40 held true for the Nucleocapsid protein.

41 **Conclusions:** Our results greatly expanded the predicted number of coronavirus proteins
42 and identified potential viral peptide suggested to be involved in viral virulence. These
43 methods and findings shed new light on coronavirus biology and provides a valuable resource
44 for future genomic studies of coronaviruses.

45

46 **Key words: coronavirus, SARS-CoV-2, MERS-CoV, meta-transcriptome, sub-genomic**

47 **RNA**

48

49 **BACKGROUND**

50 Corona virus disease 2019 (COVID-19) reached pandemic levels beginning March 2020
51 and brought unprecedented devastation to human lives and the global economy [1]. The
52 causative agent is Severe Acute Respiratory Syndrome – Corona Virus - 2 (SARS-COV-2), a
53 beta coronavirus similar to MERS-CoV, the only other active virulent beta-coronavirus.
54 MERS-CoV is the causative agent of Middle Eastern Respiratory Syndrome (MERS) and is
55 more virulent but less infectious than SARS-CoV-2 and is phylogenetically different from
56 SARS-CoV-2 (less than 90% amino acid sequence homology). Both viruses have a positive
57 single-stranded RNA genome of approximately 30 kilobases that is polyadenylated that
58 encodes 4 structural proteins (spike (S), membrane (M), envelope (E) and nucleocapsid (N))
59 that play similar roles within each virus. The two viruses diverge with respect to the receptor
60 used for cell entry, their virulent accessory proteins and the specific function(s) of the 16 non
61 structural proteins (nsp1 to nsp16). Nsp's are produced by viral proteinase cleavage of two
62 large polyproteins encoded by ORF1a and ORF1b. ORF1 is closest to the 5' end and is
63 directly translated from genomic RNA upon entrance into host cells and a ribosome skipping
64 mechanism divides it into ORF1a and ORF1b [2]. While MERS-CoV encodes at least 5
65 accessory proteins (ORF3, ORF4a, ORF4b, ORF5 and ORF8b), SARS-CoV-2 encodes at
66 least 6 (ORF3a, ORF6, ORF7a, ORF7b, ORF8, and ORF10 [3]. All proteins not encoded by
67 ORF1a or ORF1b, must be translated from sub-genomic RNAs (sgRNAs) [4, 5]. SgRNAs are
68 generated via a mechanism termed discontinuous extension that uses short sequences of
69 varying length (usually 6 to 12 nucleotides (nts)) termed Transcription Regulatory Sequences
70 (TRS's) spaced between genes to pair a 3' portion of the negative viral strand to a

71 complementary 5' leader sequence of around 70 nts. This is followed by extension of the
72 negative strand to the 5' end of the positive strand, generating a short negative strand sgRNA
73 intermediat. The RNA intermediary is then replicated to generate a positive strand sgRNA
74 that encodes viral protein(s) [6] (**Fig. 1A**).

75 Annotating viral transcriptomes is fundamental to understanding virus biology, which is
76 a key aspect in combating viral transmission, replication and pathogenesis. Prior coronavirus
77 outbreaks, such as the Severe Acute Respiratory Syndrome (SARS) outbreak in 2003 and the
78 MERS outbreak that began in 2012 and is still ongoing [7, 8], has increased research on these
79 viruses as well as coronaviruses of zoonotic origin from which human coronaviruses are
80 thought to originate. Comparing transcriptional variation of different coronaviruses may
81 reveal mechanisms behind their distinct pathogenicity and infectivity, and potentially explain
82 the molecular etiology behind how species barriers are crossed. Systematically annotating
83 differences in the transcriptional profiles of virulent coronaviruses that is buried within
84 numerous metatranscriptomic data sets may shed new light on viral transmissibility and
85 virulence. However, even a simple systematic comparison of their *in-vitro* transcriptional
86 profiles is lacking.

87 For newly emerged SARS-CoV-2 virus, sequencing plays an essential role in diagnosis
88 and monitoring of strain evolution [3, 9]. However, in general, sequencing data sets for
89 SARS-CoV-2 and MERS-CoV are limited to the description of both viral and host transcripts
90 generated during infection of *in-vitro* cell lines as well as model organisms. Analysis of viral
91 transcriptomes originating from different viral strains in humans is overlooked as suitable
92 analysis tools are lacking.

93 Sequence homology plays an essential role in the functional annotation of viral genes.
94 However, sequence homology alone does not guarantee protein expression as rapidly
95 mutating RNA viruses can harbor sequence alterations that result in novel or mutated ORFs
96 that are not transcribed nor expressed. Therefore, direct profiling of viral RNAs is the step
97 toward understanding which viral products can actually be generated. For SARS-CoV-2,
98 direct profiling of viral RNAs produced in a cultured cell line was recently conducted using
99 Oxford nanopore technology and identified the existence of a canonical and non-canonical
100 viral transcriptome. All three of these studies used isolated virus strains to infect the Vero cell
101 line isolated from kidney epithelial cells of the African green monkey that does not initiate an
102 interferon (IFN) response upon infection. Although these studies establish a basic
103 characterization of virus transcription, individual studies only characterize viral gene
104 expression of a single viral strain and are unable to determine if viral transcriptional
105 responses are altered in response to even the most basic of immune responses (e.g. IFN γ)
106 [10-12].

107 We developed a bioinformatics pipeline CORONATATOR (CORONAvirus annoTATOR)
108 to quantify viral gene expression and identify bona-fide sgRNAs in numerous publicly
109 available meta-transcriptomic data sets. Beyond outlining the variation in sgRNA profiles and
110 their relative expression, our analysis identified novel sgRNAs for several different
111 coronaviruses. It also revealed the presence of a core sgRNA repertoire that is shared between
112 SARS and SARS-CoV-2 and one that is unique to MERS-CoV. A subset of novel sgRNAs for
113 SARS-CoV-2 and MERS-CoV appear to be evolutionarily conserved in related coronaviruses
114 found in bat and pangolin. Finally, we show that the transcription of specific sgRNAs differs

115 significantly in-vitro and in-vivo as well as between different coronaviruses.

116

117 **RESULTS**

118 **CORONATATOR profiles viral sgRNAs via alignment breakpoint analysis**

119 To systematically identify and compare coronavirus sgRNAs, we sought to identify
120 publicly available coronavirus transcriptomic data sets. As of 2021/09/10, more than 3410427
121 viral genome sequences were submitted to the Global Initiative on Sharing All Influenza Data
122 (GISAID) [13]. However, few data sets contain the raw sequencing reads. Using the search
123 term “coronavirus” along with manual curation, we located raw reads in a total of 19
124 Bioprojects within the NCBI Short Read Archive that contain 588 samples for SARS-CoV-2
125 as well as related coronaviruses, such as SARS and MERS (**Table S1**). We also used an
126 additional dataset with a single sample that was recently published [10].

127 To profile the sgRNAs present within these data sets, we developed an informatics
128 pipeline (CORONATATOR). CORONATATOR is designed for the utilization of sequences
129 produced by highly accurate second generation sequencing technology that permits
130 identification of TRS sequences from individual reads. Direct RNA Sequencing on the
131 Oxford Nanopore platform can also be used to profile viral sgRNAs but is currently not
132 supported by CORONATATOR due to the limited data set availability as well as its
133 restrictions in terms of sequencing accuracy and read length bias (see **Methods**).

134 Briefly, raw reads were first aligned to their respective viral references, i.e.
135 SARS-CoV-2 (GeneBank ID NC_045512.2), SARS-CoV (GeneBank ID NC_004718.3),

136 MERS-CoV (GeneBank ID NC_019843.3) or reference for other species of
137 coronaviruses(**Table S1**). Specific sgRNAs were inferred from alignment breakpoint analysis
138 that identified reads that spanned the junctions between the 5' leader sequence and more
139 distal genomic sequence (**Fig. 1A and Supplementary Methods**). The relative abundance of
140 a specific sgRNA to all sgRNA's in a particular sample is analogous to relative gene
141 expression. We constructed a heatmap to determine how viral genotype and viral origin
142 (e.g.-in-vivo vs in-vitro) influences viral gene expression (**Fig. 3** and discussion below).

143 CORONATATOR was designed to profile all possible breakpoints. However, to obtain
144 bona-fide sgRNAs, we removed both rare breakpoints and breakpoints that were inconsistent
145 across samples. A complete breakpoint consists of two separate genomic positions (**Fig. 1A**).
146 We also analyzed non-sgRNA breakpoints, for which the 5' position does not encompass the
147 leader TRS. Our data suggested that non-sgRNA breakpoints are very rare (usually below
148 0.05% of total sgRNA breakpoints) and inconsistent, as these breakpoints were never
149 identified in more than a single study. We therefore focused on sgRNAs formed with the
150 canonical 5' leader and a 3' body part.

151

152 **Most predicted coronavirus ORFs can be validated by sgRNA analysis**

153 Many ORFs are annotated for SARS-CoV-2 based on consensus sequence annotation
154 and the existence of some are disputed by proteomics as well as sequencing studies [3, 11].
155 Only after examination of a large number of data sets from multiple studies were we able to
156 confidently assign commonly annotated ORFs into one of three categories (core, low support
157 and no support) (**Fig. 1B**). Identifying bonafide sgRNAs requires multistudy and multisample

158 analysis as unique artifactual sequences are often generated during sequence library
159 preparation or sequencing([14, 15]). Additionally, many non-canonical sgRNAs found in low
160 abundance may be random aberrant transcripts without dedicated function (Kim et al, 2020).
161 Therefore, only sgRNAs that are present in multiple studies and data sets are true sgRNA
162 candidates. To classify each viral gene we considered factors such as sgRNA relative
163 abundance, TRS conservation and the potential for leaky ribosome scanning that can be
164 affected by start-codon hijacking[16].

165 We first validated the commonly annotated ORFs for SARS-CoV-2, SARS-CoV and
166 MERS-CoV by looking at the extent of sequencing evidence that supports the existence of
167 specific sgRNAs. For SARS-CoV-2, thirty four samples were kept after removing those with
168 less than 20 sgRNA reads. To identify robust and consistent sgRNAs that represent the “core”
169 repertoire, which we assign to our first sgRNA category, we pooled all sgRNAs identified for
170 a specific virus using a weighted average approach (see Methods) and noted their relative
171 abundance. At a relative abundance of 0.5%, 8 canonical breakpoints emerged corresponding
172 to 8 sgRNA species that harbor 8 well-described ORFs for SARS-CoV-2: S, E, M, N, ORF3a,
173 ORF6, ORF7a and ORF8 (**Fig.1B-C, Fig. S2A, Table S2, Table S3**). The sgRNA breakpoints
174 for these ORFs are situated between 9 and 162 nt upstream of the start codon. N is the most
175 abundant core sgRNA, representing 54% of the core sgRNAs identified in all samples. The E
176 sgRNA is the least abundant at 1.5%, and the only core protein not identified in recent
177 proteomics studies [11, 17]. ORF7a, M, ORF3a, S, ORF8 and ORF6 are present at 10.6%,
178 8.4%, 6.9%, 6.1%, 5.9 and 2.7% respectively. Together, these 8 core sgRNAs account for 70%
179 to 100% of the total sgRNAs depending on sample type (e.g. *in-vivo* vs *in-vitro*), viral strain

180 and read coverage (**Fig. 2, Table. S3**).

181 Beside their high relative abundance, these 8 core sgRNA are also defined by a shared
182 canonical body TRS with a converseved core sequence of “ACGAAC”, which is unique to
183 this group of sgRNAs. This core sequence could be necessary and sufficient for sgRNA
184 formation. Futhermore, the same 8 core sgRNAs, as well the core TRS sequence, were shared
185 by SARS (Fig. S2). The 7 core sgRNAs for MERS following (S, E, M, N, ORF3, ORF4a and
186 ORF5) (Fig. 1C) also utilize this core sequence, with the exception of N that has a TRS
187 which contains “ACGAA”.

188 A second category of sgRNA was generally present at low relative abundance and does
189 not use this core sequence a conserved core TRS sequence. This category include ORF7b in
190 SARS-CoV-2 and SARS-CoV, ORF3b in SARS and ORF4b and ORF8b in MERS-CoV. For
191 SARS-CoV-2, E has an average relative abundance of 1.5% which is the lowest amongst the
192 core ones, while ORF7b’s is only 0.02% . This low abundance or low efficiency in sgRNAs
193 formation may result from the use of noncanonical TRSs. This group of sgRNAs do not use
194 the conserved core TRS sequence as core sgRNAs do, meaning the sequence homology they
195 rely on for recombination is always shifted a few bases from the core and quite often they
196 contain mismatches between leader and body TRS.

197 Other predicted ORFs fell into the third category with no sgRNA support, at least in the
198 data set we examined. When factor in evidence beyond sgRNA support. This category can be
199 futher divided into two sub-categories. The first would be no sgRNA support but can
200 potentially be translated. It has been observed that some coronavirus ORFs can be expressed
201 via a leaky ribosome scanning mechanism [16]. ORF9b of SARS-CoV-2 falls into this

202 sub-category. Indeed, multiple recent proteomics studies showed support for the ORF9b
203 protein product in SARS-CoV-2[17, 18]. Its homolog in SARS-CoV, also named ORF9b,
204 falls in the same category. Interestingly, the ORF7b of SARS-CoV-2 and SARS-CoV were
205 mentioned in previous studies to be in this category (Schaecher et al., 2007), and indeed the
206 long stretch (362 nt in SARS-CoV-2 and 365 nt in SARS) between start codons of ORF7b
207 and preceding ORF7a are void of additional start codons. Yet, these gene products still form
208 their own sgRNAs at low abundance.

209 The second sub-category would contain the most suspicious ORFs, where sgRNA
210 support cannot be found and intervening start codon between them and the closest sgRNA
211 breakpoint would make their expression very unlikely. This category includes a few
212 commonly annotated ORFs: ORF3b, ORF9c and ORF10 of SARS-CoV-2, ORF8b in
213 SARS-CoV. The several out of frame start codons between these ORFs and preceding ones,
214 along with the absence of corresponding sgRNAs and its absence from proteomic studies [10,
215 17, 18], strongly argues that these proteins are not generated. Indeed, the existence ORF10
216 was recently debated in recent manuscripts [11, 12]. The evidence described above indicates
217 the potential pitfalls of conducting experiments on viral products from putative ORFs with no
218 sgRNA or proteomic support. For example, a recent study that generated a synthetic version
219 of the predicted truncated version of ORF3b in SARS-CoV-2 speculated that the putative
220 truncated version in SARS-CoV-2 had a stronger anti-IFN activity than the SARS version
221 [19].

222

223 **Identification of novel sgRNAs with non-canonical TRSs in SARS-CoV-2, MERS-CoV**

224 **and SARS-CoV.**

225 As mentioned before, during formation of the core sgRNA repertoire, a body TRS that
226 contains a minimal core sequence will pair with the leader TRS. For each particular core
227 sgRNA, the two TRS's used must be of the same length and sequence, although the length
228 can vary between sgRNAs (**Fig. 1C**). We found the average length of these canonical TRS's
229 for SARS-CoV-2 was ~9.6 nts. Interestingly, the same core sequence is used in SARS, while
230 MERS also uses a six nucleotide TRS with a different core sequence (**Fig. S2A,B**).

231 When we looked for sgRNAs that composed more than 0.2% of sgRNA transcripts, we
232 identified three additional sgRNAs that were present in at least two separate samples and
233 studies (**Fig. 2A**). All three novel sgRNAs contained breakpoints that did not utilize canonical
234 TRS sequences that are present in core sgRNAs. The three breakpoints support the
235 discontinuous extension model of sgRNA formation, as the sequence from the body strand
236 was found in the TRS sequences of the final transcript (**Fig. 2B, Fig. S3A-C**). On a separate
237 note, sequence analysis of stranded RNA library preps identified the presence of negative
238 strand sgRNAs, which were not described in the previous Nanopore sequencing manuscripts
239 [10-12]. As previously noted for artificial TRS's, analysis of these non-canonical breakpoint
240 sequences revealed that TRS's without perfect complementarity may pair, and/or that large
241 regions of complementarity around a core TRS between the body to itself, maybe used for the
242 formation of sgRNAs (**Fig. 2B**). Our analysis confirmed that TRS sequences can vary
243 significantly between distantly related viruses and find that canonical TRS sequences can be
244 more than 30 nt in length in some coronaviruses (**Fig. S2D**).

245 The three novel TRS's generated three novel sgRNAs that we have termed putative

246 ORF2b (pORF2b), alternative M (aM) and truncated ORF7b (tORF7b). The longest novel
247 sgRNA, pORF2b, is within the S gene and has two alternative TRS's positioned around
248 22501. Interestingly, it encodes a novel peptide that has a domain structure that is conserved
249 in closely related coronaviruses, with at least one virus harboring and extended ORF (**Fig.**
250 **2B,C**). The second novel breakpoint is located at 26494, 31 nt downstream of the canonical
251 breakpoint for M. The sgRNA would support M expression, but with an alternative 5' UTR
252 (**Fig. S3A**). The shortest of the three novel sgRNA's has its breakpoint positioned at 27761
253 and codes for a truncated version of ORF7b (tORF7b). The truncation removes the
254 extracellular domain and 14 of the 24 amino acids that comprise the transmembrane domain
255 (**Fig. 2A, Fig. S3B**). This sgRNA is expressed at relatively high levels both *in-vivo* and
256 *in-vitro* and likely harbors novel functions (see discussion below).

257 Translation of pORF2b results in a 36 amino acid peptide. It was predicated by PSIPRED
258 [20] to have a intracellular protein binding coil and two short alpha-helices that overlap a
259 transmembrane domain, with the second alpha helix partially extracellular (**Fig. 2A, C**).
260 pORF2b was present in 4 samples in two separate studies. The highest expression of pORF2b
261 was observed in a patient derived sample from Washington State in the US (SRX7884411),
262 where it accounted for a substantial 11.1% of the total sgRNAs. In a separate patient sample
263 (SRX7884409) from the same bioproject, the novel ORF represented 1.1% of the sgRNAs
264 identified (**Table. S3**). The virus strains infecting these two patients differed by one
265 nucleotide. Five other patient samples from the same study with different viral strains (**Table.**
266 **S3**) did not yield sgRNAs for pORF2b. The low breakpoint read numbers for these samples
267 as well as viral strain may contribute to the variable detection of pORF2b *in-vivo*. This

268 indicates that the level of pORF2b transcripts maybe loosely correlated with viral strain and
269 further demonstrates that samples within this bioproject are not cross contaminated with an
270 artifactual pORF2b sgRNA. SgRNA pORF2b was also identified in a separate study
271 (PRJNA615032), in two in-vitro samples that used a different viral strain than any of those
272 identified in the *in-vivo* study (**Table. S3**).

273 We searched for sequence conservation of pORF2b in other related Sarbecoviruses,
274 including SARS-CoV, HKU3 (bat coronavirus), RaTG13 (a bat coronavirus proposed to be
275 directly related to SARS-CoV-2) and a coronavirus infecting pangolin (SRX7732088)[21]. A
276 corresponding ORF was identified in all four viruses, with the highest level of homology
277 found in RaTG13, with 91.89% nucleotide identity (**Table S4 & Fig. 2C**). Interestingly,
278 pORF2b and more so the pangolin version which has a C terminal extension, share high
279 similarity with the ligand binding domain of human IL17RB (**Fig. 2D and see section**
280 **“Discussion”**).

281 The third novel breakpoint was located at position 27761, within ORF7b, and encodes a
282 truncated version of ORF7b (tORF7b). We identified this transcript and its relative abundance
283 *in-vivo* and *in-vitro* in two separate bioprojects that included more than one viral strain. This
284 transcript was also recently identified in a VERO cell line infected by a single viral strain
285 [10]. Interestingly, this novel sgRNA was expressed at relatively high levels both *in-vivo* and
286 *in-vitro* (**Fig. 3**), and a SARS-CoV homolog of this sgRNA was also present in several
287 samples across two studies. This truncated version of ORF7b is missing the intracellular
288 domain and more than half of its transmembrane domain, while retaining its hydrophilic
289 extracellular domain (**Fig. S3D**). ORF7b is present in the SARS-CoV virion particle and is

290 homologous ORF7b encoded by SARS-CoV-2 [16]. The portion of ORF7b encoded by
291 tORF7b is highly conserved in SARS (**Table S4 & Fig. S3D**). Intriguingly, a previous study
292 observed that a 45 nt deletion in SARS ORF7b that removes much of the transmembrane
293 domain lost in tORF7b, attenuated the induction of interferon-beta, provides a replicative
294 advantage *in-vitro* and *in-vivo* as well as to cells pretreated with interferon-beta [22]. Future
295 research will reveal if this novel sgRNA encodes a novel virulent peptide that has function(s)
296 antagonistic to IFN while subverting the initiation of an interferon response.

297 We also obtained a significant amount of *in-vivo* and *in-vitro* sequence data sets for
298 MERS-CoV, allowing us to identify abundant non-canonical sgRNAs (**Fig. S2B**). This novel
299 sgRNA (putative ORF8c or pORF8c), is predicted to encode a ORF that translate into a novel
300 51 amino acid peptide. This novel sgRNA was identified in 5 separate studies, both *in-vivo*
301 and *in-vitro*, ranging in abundance from 0.03% to 1.0% of total sgRNAs. PSIPRED suggest
302 this novel peptide has a transmembrane domain connected to a cytoplasmic helix domain. We
303 also looked for its conservation in other Merbecoviruses, including HKU4, HKU5 and an
304 Erinaceus coronavirus. pORF8c could be found in all 3 with varying conservation (**Fig. S3E,**
305 **Table. S4**). The cytoplasmic N terminal was the most conserved across Merbecovriuses and
306 C terminal elongated versions were observed in HKU5 and Erinaceus (**Fig. S3E**).

307 To exhaust our search for novel sgRNAs, we lowered our threshold value to a relative
308 abundance of 0.01%, while maintaining our other criteria. This analysis identified additional
309 novel sgRNAs that appeared in more than one study for SARS-CoV-2, SARS-CoV and
310 MERS-CoV (**Table S5**). Additional sequencing and future experiments will determine the
311 significance of pORF2b, tORF7B and aM as well as the numerous other novel sgRNAs

312 present at extremely low abundance.

313

314 **CORONATOR detects experimentally induced alteration of novel pORF8c relative**
315 **abundance.**

316 We next wished to validate the experimental utility of our pipeline and validate that a
317 novel sgRNA responds to experimental stimuli in a manner similar to other established viral
318 genes. To accomplish this, we utilized an experimental data set that tested the effects of
319 Gleevec and IFN- β on host gene expression during treatment of MERS-CoV infection
320 in-vitro (PRJNA233943 & PRJNA233944) (**Fig. S5**). Specifically, we analyzed the effects on
321 viral load, viral gene expression and the expression of novel pORF8c. Initial analysis
322 demonstrated that decreased viral load broadened the expression of individual viral genes
323 (**Fig. S5**). Importantly, even at low viral loads, the ratio of N to S remained high,
324 demonstrating that this ratio is not influenced by viral abundance, but by *in-vitro* and *in-vivo*
325 context. The effect of both Gleevec and IFN- β on viral gene expression was not uniform,
326 having different effects on different viral genes. Interestingly, the expression profile of novel
327 pORF8c followed the same trend of N and E with respect to viral load in response to IFN- β
328 and Gleevec. This demonstrates that pORF8c has the same biological response in terms of
329 gene expression as some “core” sgRNAs in this context.

330

331 **The relative abundance of Spike sgRNAs is elevated for SARS-CoV-2 *in-vivo*.**

332 When processing the data sets, we noticed two distinct patterns of read coverage along
333 the SARS-CoV-2 reference genome that suggested that viral reads originate from two sources.

334 Upon further examination, it was revealed the two sources were *in-vivo* and *in-vitro* derived
335 samples (**Fig. S1**). The former is composed of extracellular virion particles and infected host
336 cells present in BALF (human) and nasal washes (Ferret) or lung homogenate (MERS), while
337 the latter is composed of infected cells that are not subject to systemic or sometimes innate
338 (e.g. VERO cells do not produce IFN) anti-viral responses. *In-vivo* derived viral sequences
339 obtained primarily from BALF for SARS-CoV-2 (primarily BALF) generally covered the
340 entire viral reference length, with little bias towards the sgRNA containing 3' end. In contrast,
341 highly elevated coverage at the 3' end of the viral genome was observed in the *in-vitro*
342 samples due to the formation of nested sgRNAs during viral transcription.

343 SARS-CoV-2 and MERS-CoV are the only active virulent coronaviruses and are present
344 in both *in-vivo* and *in-vitro* derived metatranscriptomic data sets. We analyzed the relative
345 abundance of sgRNAs generated *in-vivo* and *in-vitro* for both SARS-CoV-2 and MERS-CoV.
346 When comparing the relative abundance of viral sgRNAs generated *in-vivo* to those
347 generated *in-vitro*, it was evident that the ratio of S sgRNAs to N sgRNAs was significantly
348 higher *in-vivo*, especially for SARS-CoV-2 (0.04 *in-vitro* vs 0.69 *in-vivo* for SARS-CoV-2, p
349 value 0.0012 with Wilcoxon ranksum test) (**Fig. 4A-B, Fig. S4**). The differences in
350 environmental pressures that influence the requirement for these sgRNAs for viral
351 replication, provides a general explanation for this striking variation in sgRNA levels. The
352 selective pressures may alter viral transcriptional responses that promote viral propagation.
353 For example, the primary function of the S protein centers around host cell recognition and
354 invasion while the primary function of the N protein centers around the regulation of viral
355 RNAs to promote viral replication. This is mediated by direct binding of the 3' end of the

356 viral genome, the viral packaging signal as well as TRS's [23-25]. Another explanation,
357 which is not mutually exclusive, is that the increased S/N ratio *in-vivo* is due to an altered
358 viral Replication Transcription Complex (RTC) that favors TRS read-through, preferentially
359 generating longer sgRNAs. Such a “global” alteration of viral transcription likely involves
360 host factors, as observed for Infectious bronchitis virus (IBV), a gamma coronavirus in which
361 the N protein is phosphorylated by cellular GSK3 to recruit the helicase DDX1 to promote
362 TRS read through during the formation of long sgRNAs [26]. In this regard, the N protein is
363 generated from the shortest sgRNA while the S protein is generated from the longest. Future
364 electron microscopy studies on *in vivo* and *in-vitro* viron particles will determine if Spike
365 sgRNA abundance in SARS-CoV-2 correlates with spike protein levels on viron surfaces.
366 Other examples of sgRNAs that are significantly differentially expressed *in-vitro* and *in-vivo*
367 include the overall increase in the levels of accessory sgRNAs that act via multiple pathways
368 to quell the immune response in both SARS-CoV-2 and MERS-CoV (**Fig. S4**, [27]).

369 To obtain a clearer perspective on how the relative abundance of SARS-CoV-2 sgRNAs
370 compares to other coronaviruses *in-vivo* and *in-vitro* as well as determine if additional novel
371 sgRNAs have been overlooked, CARONTATOR was utilized to analyze additional
372 coronaviruses. This analysis included OC43, NL63, HKU1 as well as bat and pangolin
373 viruses with high sequence homology to SARS-COV-2 [9, 21] (**Fig 4C**, **Table S1** and **Fig.**
374 **S2**). Some datasets did not yield enough breakpoint reads to be informative. For example,
375 analysis of the the bat virus RaTG13, with the highest homology to SARS-COV-2, yielded
376 only 1 break point read and was therefore omitted from **Fig. 4C**.

377 Of the different coronaviruses profiled, SARS-COV-2 stands out as having the highest

378 levels of S sgRNAs, especially *in-vivo* (see discussion below and **Fig. 4C**). Our analysis
379 indicates that this is independent of viral strain as it is present at high levels in different
380 strains identified *in-vivo* (**Fig. 2**). The high levels of Spike protein may play a role in the
381 viruses ability to cross the species barrier (see discussion below) and it's high rate of
382 infectivity. In agreement, we noted that the relative levels of the Spike sgRNA is positively
383 correlated with coronavirus infectivity. Viral infectivity and levels of S sgRNAs *in-vivo* are as
384 follows: SARS-COV-2> HKU1> MERS [28]. However, S protein levels alone are not
385 sufficient to cause high levels of SARS-CoV-2 transmissibility, as factors such as Spike
386 protein stability, receptor avidity [29] and viron stability[30], also contribute to viral
387 transmissibility.

388

389 **Mutations in the RTC reverse the expression of N and S sgRNAs in vitro and in-vivo.**

390 We also observed mutations in viral RTC components that altered the expression profile
391 of S to N. Specifically, the viral strain Kim 2020 had one unique non-synonymous mutation
392 in the RTC component nsp3, a papain protease that binds the N and M protein (**Fig. 3**). The
393 transcriptome generated *in-vitro* for this viral strain showed a dramatic increase in the S to N
394 ratio, mimicking the expression profile of viruses found *in-vivo* (**Fig. 3, Fig. S6**). Interestingly,
395 a viral strain identified *in-vivo* (SRX7852918), had two non-synonymous mutations in nsp3,
396 as well as nsp6 and nsp12 and had an *in-vitro* like transcription profile, with a decreased S to
397 N ratio (**Fig. S6**).

398 The observation that mutations in nsp3 occur in the two viruses with altered gene
399 expression is thought provoking. Nsp3 is reported to to bind TRS's, the 3' end of the viral

400 genome, the global viral RNA packaging signal as well as the N and M proteins [23, 31, 32].
401 Additionally, phosphorylation of the N protein has been reported to alter it's conformation to
402 preferentially bind viral RNA and as mentioned above for IBV, promote TRS readthrough
403 during the generation of long sgRNAs [26, 33]. This observation tentatively implies that
404 mutations within nsp3 affect the relative abundance of sgRNAs by acting in a global
405 mechanism that influences overall viral structure and may act in concert with the mechanism
406 described above for IBV. Additionally, the altered relative abundance of N sgRNAs *in-vivo*
407 and *in-vitro* due to mechanisms discussed above, may feedback on it's interaction with Nsp3
408 and influence the function of mutations in Nsp3 *in-vivo* and *in-vitro* (**Fig. 4A-B**).

409

410 **DISCUSSION**

411 The vast amount of sequence data generated for SARS-CoV-2 thus far has primarily been
412 used for the typing and following of emerging viral strains. Although this is important, we
413 felt such a focus could be an under-utilization of a valuable information. By developing the
414 Coronator informatics pipeline, we took a step beyond the characterization of viral strains
415 and described coronavirus viral sgRNA expression and uncovered novel and conserved
416 sgRNAs with unknown function that are generated via a non-canonical TRS pairing
417 mechanism (**Fig. 2**). Functional prediction for some of these novel putative proteins is still
418 ongoing. We tentatively show that a homolog of SARS-CoV-2 pORF2b in pangolin virus
419 shares extensive similarity with human IL17RB's ligand binding domain (**Fig. 2D**). It is
420 curious that a coronavirus may generate a peptide that could theoretically disrupt IL17B and

421 IL17E (IL25) signaling as they are generally associated with promoting or inhibiting
422 inflammatory responses in specific contexts. Future proteomic studies and/or ribosome
423 sequencing studies will be required to verify the production of the protein products encoded
424 by the novel sgRNAs identified here.

425 The analysis presented here also implicates that different strains of SARS-CoV-2
426 express sgRNAs at different levels (**Table S3, Fig. 3**), especially for the newly discovered
427 sgRNAs. Our findings underscore that a true understanding of viral pathogenesis in terms of
428 sgRNA expression can only come from thorough sequencing of patient samples in which the
429 virus is under selective pressure. This begs for in-depth case examination, in which thorough
430 sequencing and analysis is conducted for different stages of COVID-19 on a strain by strain
431 basis. This would result in truly individualized patient care.

432 Although other zoonotic viruses may share extensive sequence similarity to
433 SARS-CoV-2 at the gene or genomic level, similarity alone is not sufficient for the generation
434 of pathogenic human viruses. Generally not considered during discussion of zoonotic viral
435 origins, the specific expression level of viral genes, such as the Spike protein, are likely
436 important for crossing the species barrier. For example, considering the vast number of
437 un-sampled zoonotic viruses, it is likely Spike proteins capable of crossing the species barrier
438 already exist, yet are not expressed at sufficiently high levels to enable sustainable
439 inter-human transmission. However, low level Spike protein expression would allow sporadic
440 transmission from bat to human, yet would not be sustainable as human to human transmission
441 would be low due to low S protein expression as well sanitary environments that do not exist
442 for bats. In agreement, it has been observed that people living in proximity to bat caves

443 harbor virus specific antibody without ever experiencing severe disease [34].

444 Our analysis of the meta-transcriptomic data sets identified numerous sources of RNA,
445 such as host RNA as well as microbial RNA (although not optimally captured). In a time
446 when it is unclear why some people succumb to SARS-CoV-2 infection while others do not,
447 these valuable sequences should not be wasted and could be made more useful if more
448 clinical information is shared for these data sets. Most GISAID entries for SARS-CoV-2 have
449 a meta-transcriptomic dataset that supports it. However, current GISAID entries that simply
450 outline the viral genome sequence and strain far outnumber the raw read entries we
451 identified in SRA. Sharing the raw read information will greatly help researchers study this
452 virus and ultimately curb it.

453 **CONCLUSIONS**

454 We developed a bioinformatics pipeline CORONATATOR that can take meta-transcriptomic
455 sequencing reads generated from coronavirus samples and analyze the sub-genomic RNA
456 profiles of the underlying virus, akin to a transcriptome for the virus. For emergent viruses,
457 as in the case of SARS-CoV2, homology search was usually the first and only choice of
458 predicting viral ORFs after sequencing was done. Now our tool can provide additional
459 evidence. By applying it to large number of SARS-CoV2 and related viral datasets,
460 interesting biology about these coronaviruses were revealed. In addition to define core and
461 predict novel ORFs, our results suggested, for beta-coronaviruses, the spike to nucleocapsid
462 ratio to be a potential tunable in adjusting viral life style and the elevation of this ratio in
463 SARS-CoV2 may contribute to its strong transmissibility. The methods and findings

464 presented here provides a valuable resource for future genomic studies of coronaviruses.

465 **METHODS**

466 **Data collection**

467 All sequencing data used were collected from NCBI Short Reads Archive (SRA). Some
468 nanopore datasets were downloaded from online repository described in their respective
469 manuscripts [10]. The bioprojects were located by searching with key words “coronavirus”
470 and with manual curation, only meta-transcriptomic data were kept Raw reads files were
471 downloaded from SRA using wget with a customized script, SRAtoolkit were used to
472 generate compressed fastq files from downloaded sra files. After initial sequence alignment
473 using bwa with reference genome sequences of SARS-CoV , SARS-CoV-2 or MERS-CoV,
474 samples with too few viral reads were filtered out. CORONATATOR only uses reads
475 generated from second generation technologies (Illumina), nanopore data were used for
476 comparison.

477 **Coronator**

478 CORONATATOR were a series of perl and bash scripts developed for profiling and analysis
479 of RNA-Seq data from coronavirus. It consists of 3 major steps, including preprocessing,
480 breakpoint identification, sgRNA calling and profiling, details below.

481 **Preprocessing**

482 BAM files were generated from sequence alignment with reference genomes of SARS-CoV ,
483 SARS-CoV-2 or MERS-CoV, for viruses from bat and pangolin, responsive genome
484 assemblies were obtained from NCBI as references. SNPs were called and filtered with

485 bcftools [35] and annotated with vcf-annotator [36]. In addition, consensus genome
486 sequences were also generated with filtered SNPs for further analysis.

487 **Breakpoint identification**

488 Breakpoints were identified from alignments with soft or hard clips, these alignments were all
489 partial alignments largely caused by reads with recombination joints, which was generated by
490 the mechanism through which coronavirus produce their sgRNA. In this step, a matrix of
491 reads' information, breakpoint sites, CIGAR strings together with possible TRS sequences
492 was generated.

493 **sgRNA calling and profiling**

494 Typical sgRNAs were identified and defined by two breakpoint coordinates on a reference
495 genome sequence, these sites were obtained by extracting breakpoints from partial alignments,
496 i.e. one from primary alignment and the other from supplementary alignment. To recognize
497 possible TRS pattern, sequences between breakpoint pairs were extracted from previous
498 generated consensus genome sequences. After that, corresponding genes of called sgRNAs
499 were identified by manually comparing the distances between start codons of known viral
500 genes and their breakpoints. Biosamples with more than 20 sgRNAs were used for further
501 analysis, in these samples, sgRNAs were counted by genes and normalized by total sgRNA
502 count to obtain a transcription profile matrix.

503 **Novel ORF identification**

504 Potential ORFs were predicted using Prodigal [37] with -s arguments to write all potential
505 genes. An in-house python script was also used to identify very short ORFs. Then for
506 sgRNAs with multiple bioproject support, we calculated and sorted the distances between

507 their breakpoints and all identified start codon sites. ORFs that start closest to upstream
508 breakpoints were bookmarked and manually checked for verification.

509 **Sequence alignment and phylogenetic analysis**

510 Consensus genome sequences of SARS-CoV-2, SARS-CoV, MERS-CoV and biosamples
511 from bat or pangolin or other human coronavirus with more than 20 sgRNAs were used for
512 phylogenetic analysis. Multi-sequence alignment were performed with MAFFT [38],
513 Maximum likelihood consensus trees were constructed using IQ-TREE [39] with 1000
514 bootstrap times.

515 **Converting Nanopore sgRNA proportion to short reads'**

516 Kim et al included both nanopore data and short read data. The ratios between the two were
517 used to convert the other nanopore data sets to proportions comparable with others in this
518 study.

519 **Plots and statistical analysis**

520 Heatmaps showing gene expression profile were produced using 'heatmap.plus' package.
521 SgRNA expression dot plots and boxplots were made with 'ggplot2' package to compare
522 difference between gene expression among different sample origin, T-test and wilcoxon test
523 were used for statistical analysis.

524 **Function annotation**

525 Novel peptide sequences were aligned with EMBL online tool FASTA
526 (<https://www.ebi.ac.uk/Tools/sss/fasta/>) against UniProtKB/Swiss-Prot database with default
527 arguments. NCBI CD Blast online service was used to identify protein domains.

528 **Sequence conservation**

529 To check for sequence conservation of putative peptides in related viral species, we generated
530 a reference database containing all predicted ORFs from related viral genomes. DC
531 MegaBlast (DisContinuous MegaBlast) was used to search for inter-species homologs.
532 Arguments were set as follows: window_size 0, gapopen 0, gapextend 2, penalty -1, reward 1,
533 num_alignments 1. A group of homologous ORFs were then subjected to multiple sequence
534 alignment (MSA) using MAFFT. After that CLUSTAO (Clustal Omega) was used to
535 calculate an identity matrix for the MSA result. The same procedure was performed for both
536 nucleotide and amino acid sequences.

537 **ABBREVIATIONS**

538 sgRNA: sub-genomic RNA

539 ORF: open reading frame

540 TRS: Transcription Regulatory Sequences

541 COVID: Corona virus disease

542 MERS: Middle Eastern Respiratory Syndrome

543 SARS: Severe Acute Respiratory Syndrome

544 CORONATATOR: CORONAvirus annoTATOR

545 **DECLARATIONS**

546 **Ethics approval and consent to participate:** Not applicable.

547 **Consent for publication:** Not applicable.

548 **Availability of Data and Materials:** All used sequencing data are accessible with accession

549 number provided in supplementary table 1, code of CORONATATOR is accessible at:
550 <https://github.com/15274972986/CORONATATOR>.

551 **Competing interests:** The authors declare that they have no competing interests.

552 **Funding:** This study was supported by Shanghai Institute of Immunology COVID-19 Special
553 Fund.

554 **Authors' contributions:** Conceptualization: Lei Chen; Methodology: Lei Chen, Lyu Lin;
555 Investigation: Lyu Lin, Ru Feng, Mingnan Zhang, Yinjing Liao; Visualization: Lyu Lin, Qiyu
556 Gong; Supervision: Lei Chen, Xiaokui Guo, Bing Su, Yanjiao Zhou; Writing—original draft:
557 Lei Chen, Yair Dorsett, Lyu Lin; Writing—review & editing: Lei Chen, Yair Dorsett, Lyu Lin,
558 Ru Feng.

559 **Acknowledgements:** We thank Dr. Qiming Liang of Shanghai Institute of Immunology for
560 his insightful suggestion.

561

562

563 REFERENCES

564 1. WHO: **World Health Organization: Coronavirus disease 2019 (COVID-19) Situation**
565 **Report.** In., vol. 2020: World Health Organization; 2020.

566 2. Knipe DM, Howley PM: **Coronaviridae.** In: *Fields virology.* vol. 28, sixth edn.
567 Philadelphia: Wolters Kluwer Health/Lippincott Williams & Wilkins; 2013: 825-859.

568 3. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH, Pei YY *et al*:

-
- 569 **A new coronavirus associated with human respiratory disease in China.** *Nature* 2020,
570 **579**(7798):265-269.
- 571 4. Brian DA, Baric RS: **Coronavirus genome structure and replication.** *Coronavirus*
572 *Replication and Reverse Genetics* 2005, **287**:1-30.
- 573 5. Yount B, Curtis KM, Fritz EA, Hensley LE, Jahrling PB, Prentice E, Denison MR,
574 Geisbert TW, Baric RS: **Reverse genetics with a full-length infectious cDNA of severe**
575 **acute respiratory syndrome coronavirus.** *Proceedings of the National Academy of*
576 *Sciences of the United States of America* 2003, **100**(22):12995-13000.
- 577 6. Sola I, Almazan F, Zuniga S, Enjuanes L: **Continuous and Discontinuous RNA**
578 **Synthesis in Coronaviruses.** *Annual Review of Virology, Vol 2* 2015, **2**:265-288.
- 579 7. Peiris JSM, Yuen KY, Osterhaus ADME, Stohr K: **Current concepts: The severe acute**
580 **respiratory syndrome.** *New England Journal of Medicine* 2003, **349**(25):2431-2441.
- 581 8. Assiri A, McGeer A, Perl TM, Price CS, Al Rabeah AA, Cummings DAT, Alabdullatif
582 ZN, Assad M, Almulhim A, Makhdoom H *et al*: **Hospital Outbreak of Middle East**
583 **Respiratory Syndrome Coronavirus.** *New England Journal of Medicine* 2013,
584 **369**(5):407-416.
- 585 9. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL *et*
586 *al*: **A pneumonia outbreak associated with a new coronavirus of probable bat origin.**
587 *Nature* 2020, **579**(7798):270-273.
- 588 10. Kim D, Lee J-Y, Yang J-S, Kim JW, Kim VN, Chang H: **The Architecture of**
589 **SARS-CoV-2 Transcriptome.** *Cell* 2020, **181**(4):914-921.e910.
- 590 11. Davidson AD, Williamson MK, Lewis S, Shoemark D, Carroll MW, Heesom K, Zambon

-
- 591 M, Ellis J, Lewis PA, Hiscox JA *et al*:
592 2020:<https://www.biorxiv.org/content/10.1101/2020.1103.1122.002204v002201>.
- 593 12. Taiaroa G, Rawlinson D, Featherstone L, Pitt M, Caly L, Druce J, Purcell D, Harty L,
594 Tran T, Roberts J *et al*:
595 2020:<https://www.biorxiv.org/content/10.1101/2020.1103.1105.976167v976162>.
- 596 13. Shu YL, McCauley J: **GISAID: Global initiative on sharing all influenza data - from**
597 **vision to reality**. *Eurosurveillance* 2017, **22**(13):2-4.
- 598 14. Lebrigand K, Magnone V, Barbry P, Waldmann R: **High throughput error corrected**
599 **Nanopore single cell transcriptome sequencing**. *Nat Commun* 2020, **11**(1):4025.
- 600 15. Peng Q, Vijaya Satya R, Lewis M, Randad P, Wang Y: **Reducing amplification artifacts**
601 **in high multiplex amplicon sequencing by using molecular barcodes**. *BMC Genomics*
602 2015, **16**:589.
- 603 16. Schaecher SR, Mackenzie JM, Pekosz A: **The ORF7b protein of severe acute**
604 **respiratory syndrome coronavirus (SARS-CoV) is expressed in virus-infected cells**
605 **and incorporated into SARS-CoV particles**. *Journal of Virology* 2007, **81**(2):718-731.
- 606 17. Bojkova D, Klann K, Koch B, Widera M, Krause D, Ciesek S, Cinatl J, Münch C:
607 **Proteomics of SARS-CoV-2-infected host cells reveals therapy targets**. *Nature* 2020,
608 **583**:469-472.
- 609 18. Gordon DE, Jang GM, Bouhaddou M, Xu J, Obernier K, White KM, O’Meara MJ,
610 Rezelj VV, Guo JZ, Swaney DL *et al*: **A SARS-CoV-2 protein interaction map reveals**
611 **targets for drug repurposing**. *Nature* 2020, **583**:459-468.
- 612 19. Konno Y, Kimura I, Uriu K, Fukushi M, Irie T, Koyanagi Y, Nakagawa S, Sato K:

- 613 2020:<https://www.biorxiv.org/content/10.1101/2020.1105.1111.088179v088171>.
- 614 20. Buchan DWA, Jones DT: **The PSIPRED Protein Analysis Workbench: 20 years on.**
- 615 *Nucleic Acids Research* 2019, **47**(W1):W402-W407.
- 616 21. Lam TT-Y, Shum MH-H, Zhu H-C, Tong Y-G, Ni X-B, Liao Y-S, Wei W, Cheung WY-M,
- 617 Li W-J, Li L-F *et al*: **Identifying SARS-CoV-2 related coronaviruses in Malayan**
- 618 **pangolins.** *Nature* 2020, **583**:282-285.
- 619 22. Pfefferle S, Krähling V, Ditt V, Grywna K, Mühlberger E, Drosten C: **Reverse genetic**
- 620 **characterization of the natural genomic deletion in SARS-Coronavirus strain**
- 621 **Frankfurt-1 open reading frame 7b reveals an attenuating function of the 7b protein**
- 622 **in-vitro and in-vivo.** *Virology journal* 2009, **6**:131-131.
- 623 23. Liang Y, Wang M-L, Chien C-S, Yarmishyn AA, Yang Y-P, Lai W-Y, Luo Y-H, Lin Y-T,
- 624 Chen Y-J, Chang P-C *et al*: **Highlight of Immune Pathogenic Response and**
- 625 **Hematopathologic Effect in SARS-CoV, MERS-CoV, and SARS-Cov-2 Infection.**
- 626 *Frontiers in Immunology* 2020, **11**:1022.
- 627 24. Molenkamp R, Spaan WJ: **Identification of a specific interaction between the**
- 628 **coronavirus mouse hepatitis virus A59 nucleocapsid protein and packaging signal.**
- 629 *Virology* 1997, **239**(1):78-86.
- 630 25. Fan H, Ooi A, Tan YW, Wang S, Fang S, Liu DX, Lescar J: **The nucleocapsid protein of**
- 631 **coronavirus infectious bronchitis virus: crystal structure of its N-terminal domain**
- 632 **and multimerization properties.** *Structure (London, England : 1993)* 2005,
- 633 **13**(12):1859-1868.
- 634 26. Wu C-H, Chen P-J, Yeh S-H: **Nucleocapsid Phosphorylation and RNA Helicase DDX1**

-
- 635 **Recruitment Enables Coronavirus Transition from Discontinuous to Continuous**
636 **Transcription.** *Cell Host & Microbe* 2014, **16**(4):462-472.
- 637 27. Canton J, Fehr AR, Fernandez-Delgado R, Gutierrez-Alvarez FJ, Sanchez-Aparicio MT,
638 Garcia-Sastre A, Perlman S, Enjuanes L, Sola I: **MERS-CoV 4b protein interferes with**
639 **the NF-kappaB-dependent innate immune response during infection.** *PLoS Pathog*
640 2018, **14**(1):e1006838.
- 641 28. Kissler SM, Tedijanto C, Goldstein E, Grad YH, Lipsitch M: **Projecting the**
642 **transmission dynamics of SARS-CoV-2 through the postpandemic period.** *Science*
643 2020, **368**:860-868.
- 644 29. Wrobel AG, Benton DJ, Xu P, Roustan C, Martin SR, Rosenthal PB, Skehel JJ, Gamblin
645 **SJ: SARS-CoV-2 and bat RaTG13 spike glycoprotein structures inform on virus**
646 **evolution and furin-cleavage effects.** *Nat Struct Mol Biol* 2020, **27**(8):763-767.
- 647 30. Aboubakr HA, Sharafeldin TA, Goyal SM: **Stability of SARS-CoV-2 and other**
648 **coronaviruses in the environment and on common touch surfaces and the influence of**
649 **climatic conditions: A review.** *Transbound Emerg Dis* 2020, **68**:296-312.
- 650 31. Hurst KR, Koetzner CA, Masters PS: **Characterization of a critical interaction**
651 **between the coronavirus nucleocapsid protein and nonstructural protein 3 of the**
652 **viral replicase-transcriptase complex.** *Journal of virology* 2013, **87**(16):9159-9172.
- 653 32. Lei J, Kusov Y, Hilgenfeld R: **Nsp3 of coronaviruses: Structures and functions of a**
654 **large multi-domain protein.** *Antiviral Research* 2018, **149**:58-74.
- 655 33. Chang C-k, Hou M-H, Chang C-F, Hsiao C-D, Huang T-h: **The SARS coronavirus**
656 **nucleocapsid protein – Forms and functions.** *Antiviral Research* 2014, **103**:39-50.

-
- 657 34. Wang N, Li S-Y, Yang X-L, Huang H-M, Zhang Y-J, Guo H, Luo C-M, Miller M, Zhu G,
658 Chmura AA *et al*: **Serological Evidence of Bat SARS-Related Coronavirus Infection in**
659 **Humans, China.** *Virologica Sinica* 2018, **33**(1):104-107.
- 660 35. Li H: **A statistical framework for SNP calling, mutation discovery, association**
661 **mapping and population genetical parameter estimation from sequencing data.**
662 *Bioinformatics* 2011, **27**(21):2987-2993.
- 663 36. **vcf-annotator** [<https://github.com/rpetit3/vcf-annotator>]
- 664 37. Hyatt D, Chen GL, Locascio PF, Land ML, Larimer FW, Hauser LJ: **Prodigal:**
665 **prokaryotic gene recognition and translation initiation site identification.** *BMC*
666 *Bioinformatics* 2010, **11**:119.
- 667 38. Katoh K, Misawa K, Kuma K, Miyata T: **MAFFT: a novel method for rapid multiple**
668 **sequence alignment based on fast Fourier transform.** *Nucleic Acids Res* 2002,
669 **30**(14):3059-3066.
- 670 39. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A,
671 Lanfear R: **IQ-TREE 2: New Models and Efficient Methods for Phylogenetic**
672 **Inference in the Genomic Era.** *Mol Biol Evol* 2020, **37**(5):1530-1534.
- 673

674 **FIGURE LEGEND**

675 **Fig. 1: Study overview and the sgRNA profile of SARS-CoV-2.** (A) Study overview; Top
676 panel, datasets used in this study came from different hosts infected by different
677 coronaviruses , for in-vitro studies, a cell line culturing step was added. These samples were
678 subjected to meta-transcriptomic sequencing and reads were collected from Short Read
679 Archive. Bottom panel, reads mapped to viruses' genome. Alignment of reads spanning the
680 genome shows breakpoints sites. As viral sgRNAs were formed via the recombination of
681 transcript body and a fixed 5' leader via TRS homology, breakpoints with 5' position close to
682 leader TRS were sent to sub-genomic RNA profiling. (B) The canonical breakpoints plot of
683 SARS-CoV-2. The ratio of putative sub-genomic RNA, black bar indicates relative
684 abundance of profiled sgRNA. (C) ORF annotation and comparison among 3 coronaviruses.
685 Different color of blocks represent sgRNA supportive stat of that gene. Specially, red blocks
686 demonstrate novel sgRNA the algorithm identified.

687 **Fig. 2: Novel sgRNAs and responsive translational product for SARS-CoV-2.** (A)
688 Breakpoints plot for SARS-CoV-2 showing the three novel breakpoints at relative abundance
689 cut-off of 0.1%, putative TRS sequences were shown below. Count of classical breakpoints
690 were shown in grey as background. Peptides of novel ORFs, i.e. putative ORF2b (pORF2b)
691 and truncated ORF7b (tORF7b), were shown in inlets, secondary structures of these peptides
692 were predicted and shown in different color. Specially, a complete ORF7b peptide was shown
693 in grey as a reference for the truncated one; (B) Sequence homology between leader TRS
694 (top), sgRNA (middle) and body TRS (bottom) for novel sgRNAs, TRS core were shown in
695 blue. (C) Structural conservation of novel peptide translated from newly discovered sgRNA

696 putative ORF2b, left panel demonstrates consensus phylogenetic tree of responsive
697 coronavirus determined by genomic sequences, right panel compared second structure of
698 novel peptide predicted by PSIPRED Workbench. (D) Putative ORF2b in pangolin CoV
699 shows homology with IL17RB's fibronectin III like domain, which is also a ligand binding
700 domain.

701 **Fig. 3: Heatmap of sgRNA expression profile of SARS-CoV-2 with SNP annotation.** Left
702 panel shows sgRNA expression profile of SARS-CoV-2 in the transcriptomic or
703 meta-transcriptomic dataset profiled. Right panel shows all the SNP sites with annotation of
704 the responsive biosamples. Interestingly, virus strains from SRX7852918 and Kim et al had
705 distinctive SNP pattern as well as characteristic expression profiles.

706 **Fig.4: Comparison of *in-vivo* and *in-vitro* sgRNA expression.** (A) and (B) Expression
707 profile of SARS-CoV-2 and MERS-CoV, both *in-vivo* and *in-vitro* datasets were included. It
708 should be noticed that two third generation sequencing technology data were added as
709 complementary datasets to SARS-CoV-2 *in-vitro* plot, a math model was applied to adjust
710 long read expression ratio into an adapted version which was comparable with short read
711 archive datasets. Interestingly, higher levels of S and M expression ratio and lower level N
712 expression ratio were observed in *in-vivo* sample versus *in-vitro* sample in these two
713 coronaviruses. (C) Phylogenetic tree of involved coronaviruses (left), scale bar indicates
714 phylogenetic distance which were calculated as the ratio of nonidentical base positions to all
715 base positions, taxonomic classification at genus level were indicated at left part. Expression
716 ratio of Spike (S) genes in vivo and in vitro in different coronaviruses (right), each dot
717 represent a biosample, black bars indicate average expression level of responsive virus.

718

719 **SUPPLEMENTARY INFORMATION**

720 **Additional file 1: Fig. S1:** Coverage plot for in-vivo and in-vitro datasets. **(A)** Coverage plot
721 for SRX8089279, which is a representative of in-vitro sample. In-vivo RNA-Seq reads
722 relatively evenly mapped to viral genome, indicating a genomic RNA dominated sample. **(B)**
723 Coverage plot for SRX7736886, which is a representative of in-vivo sample. In-vitro reads
724 resulted a dense mapping at 3' and 5' end of the genome, which revealed active viral
725 transcription and replication in cultured cells.

726 **Additional file 2: Fig. S2 :** Breakpoint profile of 4 coronaviruses. **(A)** Breakpoints profile of
727 SARS-CoV; **(B)** Breakpoints profile of MERS-CoV; **(C)** Breakpoints profile of
728 Pangolin-CoV; **(D)** Breakpoints profile of HKU1, interestingly, this virus uses long TRS for
729 discontinuous sgRNA production.

730 **Additional file 3: Fig. S3:** Novel sgRNA breakpoint and TRS sequence. **(A)** Alternative TRS
731 of M in SARS-CoV-2. The canonical TRS region (upper panel) has 12 bases while the novel
732 one has only 6 (lower panel), start codon of M was shown in red. **(B)** TRS for tORF7b in
733 SARS-CoV-2, which has 7 bases overlapped with leader TRS. **(C)** TRS for pORF8c in
734 MERS-CoV. **(D)** and **(E)** Structural conservation of novel peptide translated from newly
735 discovered sgRNA truncated ORF7b and putative ORF8c, left panel demonstrates consensus
736 phylogenetic tree of responsive coronavirus determined by genomic sequences, right panel
737 compared second structure of novel peptide predicted by PSIPRED Workbench.

738 **Additional file 4: Fig. S4:** Detailed gene expression profile of SARS-CoV-2 and MERS-CoV.

739 **(A)** and **(B)** Detailed accessory gene expression profile of SARS-CoV-2 and MERS-CoV,
740 between in-vivo and iv-vitro datasets. Remarkably, MERS-CoV had higher ORF4a in-vivo
741 expression level while lower in-vivo ORF5 expression level. **(C)** and **(D)** Gene expression
742 level of SARS-CoV-2 and MERS-CoV among different hosts.

743 **Additional file 5: Fig. S5:** Detailed gene expression profile of MERS-CoV in PRJNA233943
744 & PRJNA233944. In this study, cells infected with MERS-CoV were treated with different
745 drugs, i.e. Gleevec and IFN- β , after 24 or 48h post-infection, distinct pattern can be observed
746 from viral gene expression profile as condition alters, especially for IFN- β , treated 24 hpi and
747 48 hpi resulted in distinct expression levels among several structural and accessory genes. It
748 also indicates that our analytical pipeline CORONATATOR is a powerful and sensitive tool
749 for analyzing how experimental manipulation effects the relative expression of specific
750 sgRNAs.

751 **Additional file 6: Fig. S6 :** Outliers in expression profiles harbour interesting SNPs. **(A)** and
752 **(B)** The in-vivo sample from Kim et al 2020 had S expression level similar to that of in-vitro
753 samples. While also have a few mutations in ORF1a that's not found in other viral strains.
754 Mirroring this, the in-vitro sample SRX7852918 have S expression level similar to that of
755 in-vivo samples, and hold several private mutations in ORF1a as well.

756 **Additional file 7: Table S1 :** Meta information of samples collected.

757 **Additional file 8: Table S2 :** Annotation of SARS-CoV-2, SARS-CoV and MERS-CoV.

758 **Additional file 9: Table S3:** SARS-CoV-2 sgRNA abundance across samples.

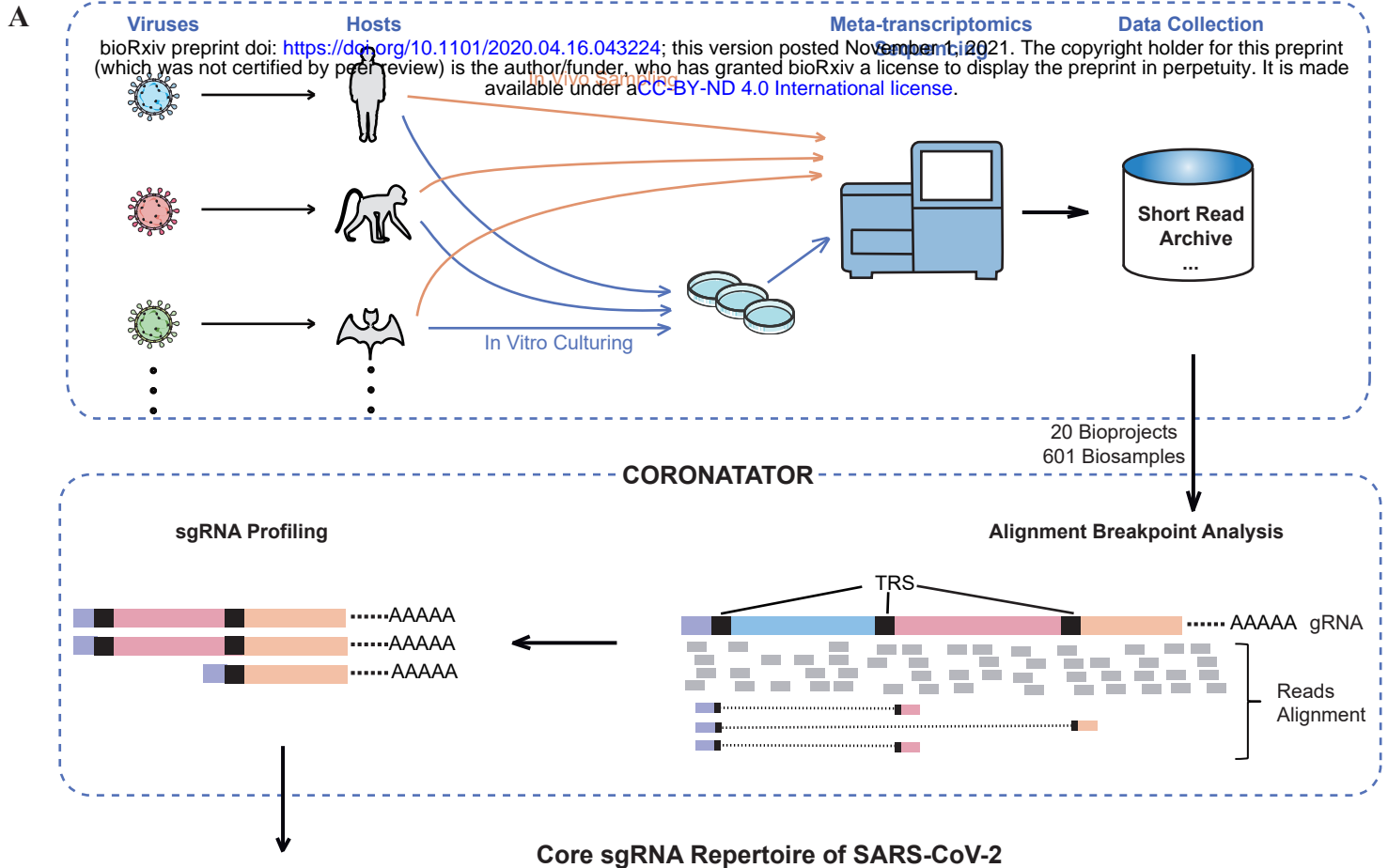
759 **Additional file 10: Table S4:** Conservation of selected novel proteins.

760 **Additional file 11: Table S5:** List of novel sgRNAs.

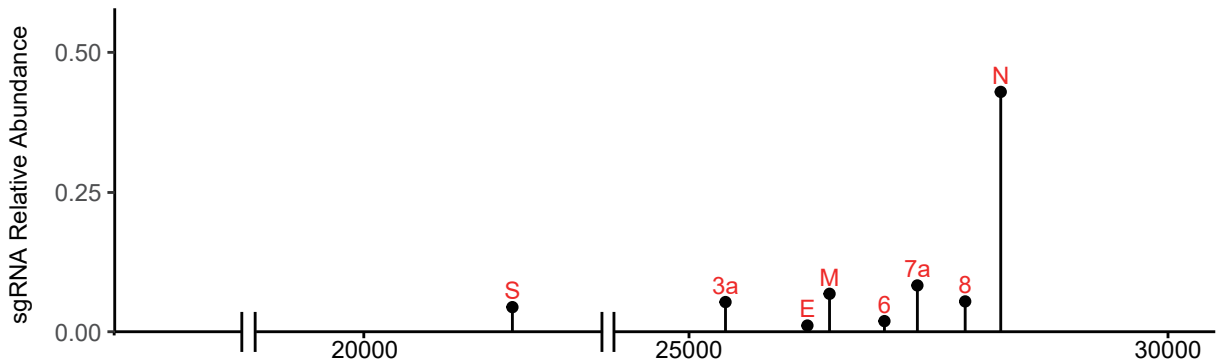
761

762

Fig. 1



B



C

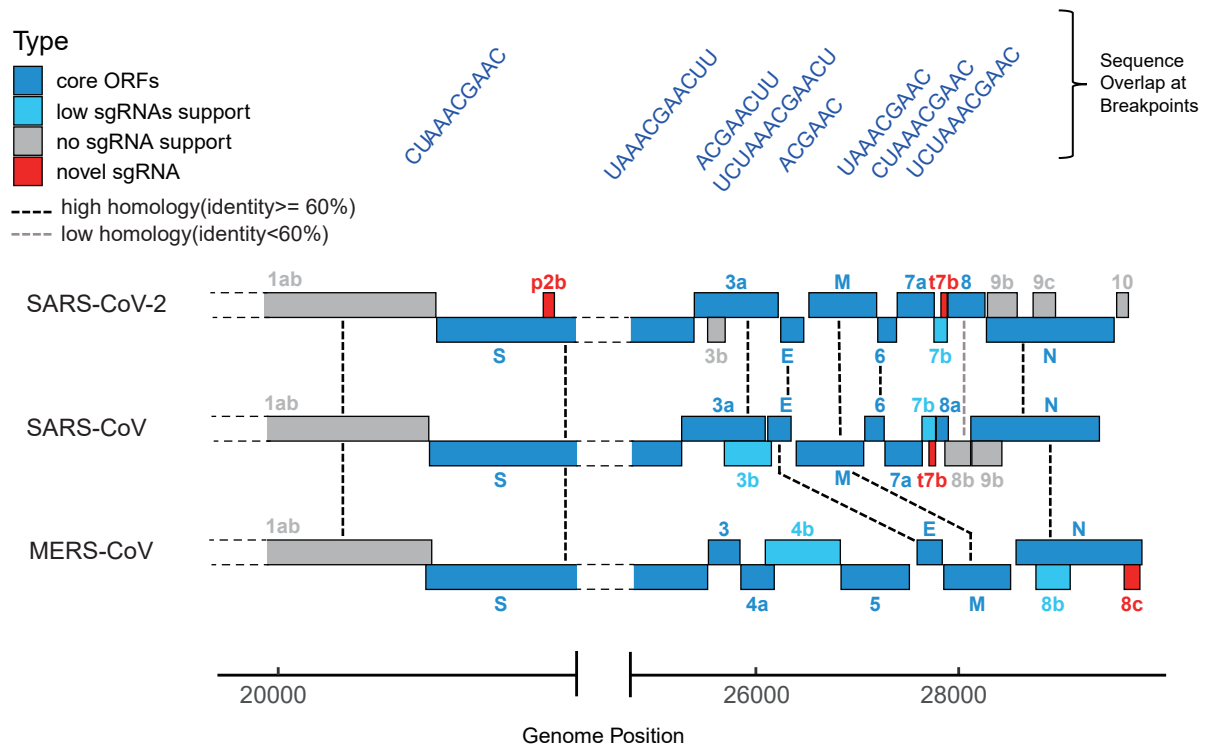
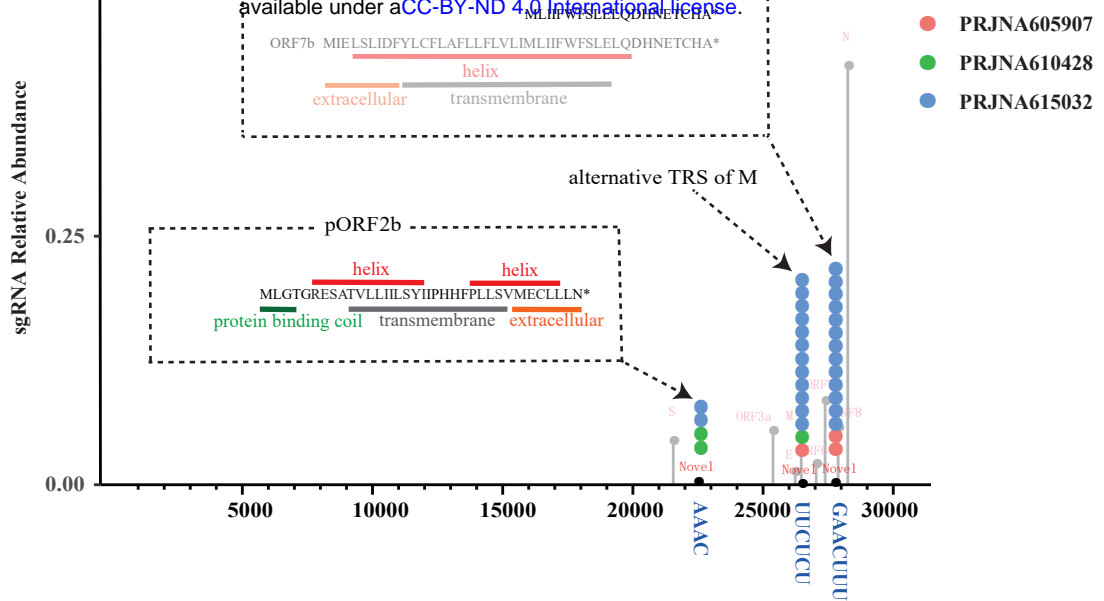


Fig. 2

Novel sgRNA in SARS-CoV-2

A

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.16.043224>; this version posted November 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.



B

Breakpoint and TRS of SARS-CoV-2 pORF2b(two alternatives)

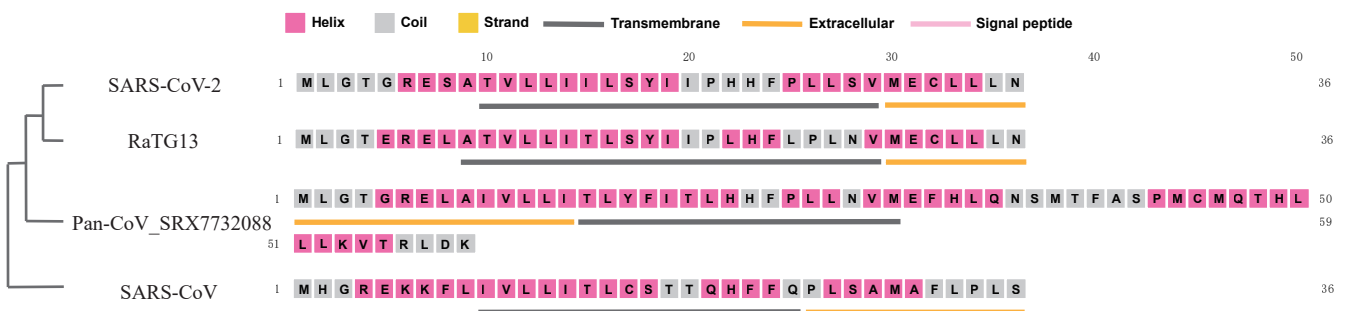
TRS core

Leader 51 UUGUAGAUCUGUUCUCUAA**ACGAAC**UUUAAAUCUGUGUGG 91
 |||
 sgRNA1 UUGUAGAUCUGUUCUCUAAACUUCUAAACUUUAGAGTCCAAC
 |||
 Body 22486 AGAAAAGGAAUCUAUCAAAACUUCUAAACUUUAGAGUCCAAC 22526

Leader 51 UUGUAGAUCUGUUCUCUAA**ACGAAC**UUUAAAUCUGUGUGG 91
 |||
 sgRNA2 UUGUAGAUCUGUUCUCUAAACUAAACUUUAGAGUCCAACCAA
 |||
 Body 22489 AAAAGGAAUCUAUCAAAACUUCUAAACUUUAGAGUCCAACCAA 22529

C

Structural conservation of the SARS-CoV-2 pORF2b amongst selected Sarbecoviruses



D

Putative ORF2b in pangolin showing homology with IL-17 receptor B

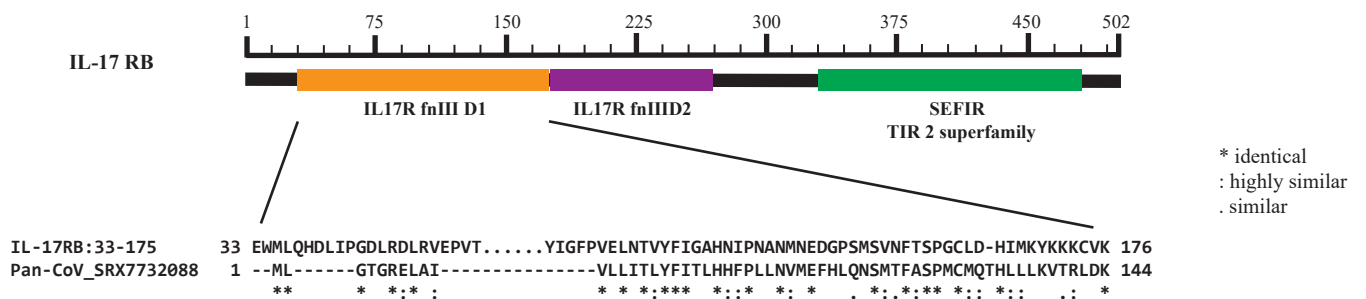
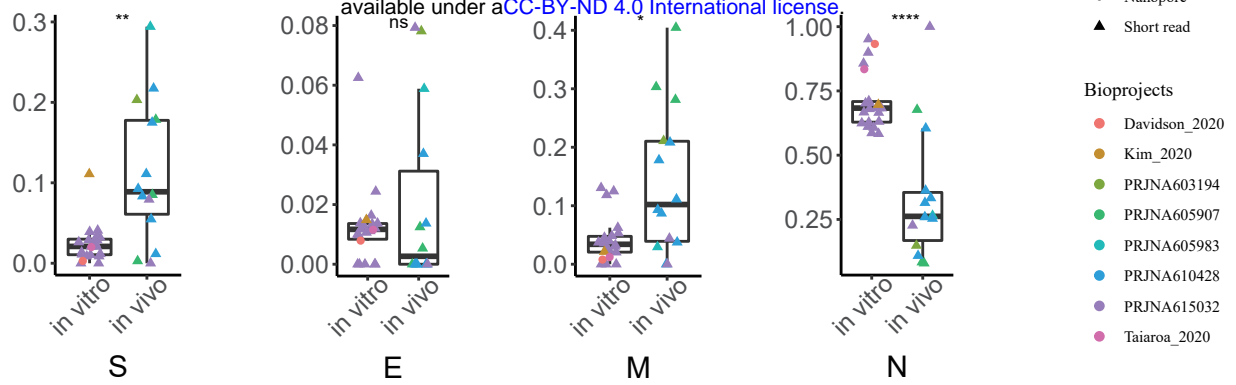


Fig 4.

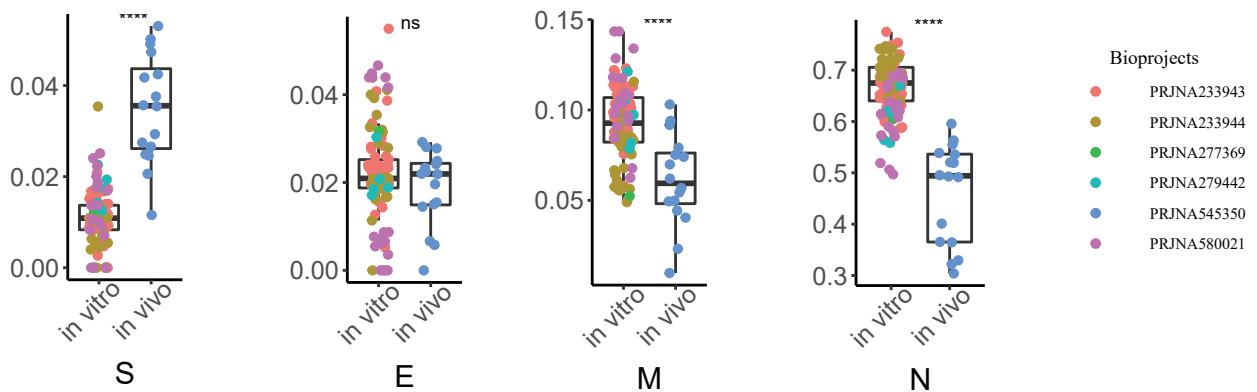
In-vivo versus in-vitro, SARS-CoV-2

A bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.16.043224>; this version posted November 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.



B

In-vivo versus in-vitro, MERS-CoV



C

Phylogenetic Tree of Involved Coronaviruses

Expression of S in Vivo Across Coronaviruses

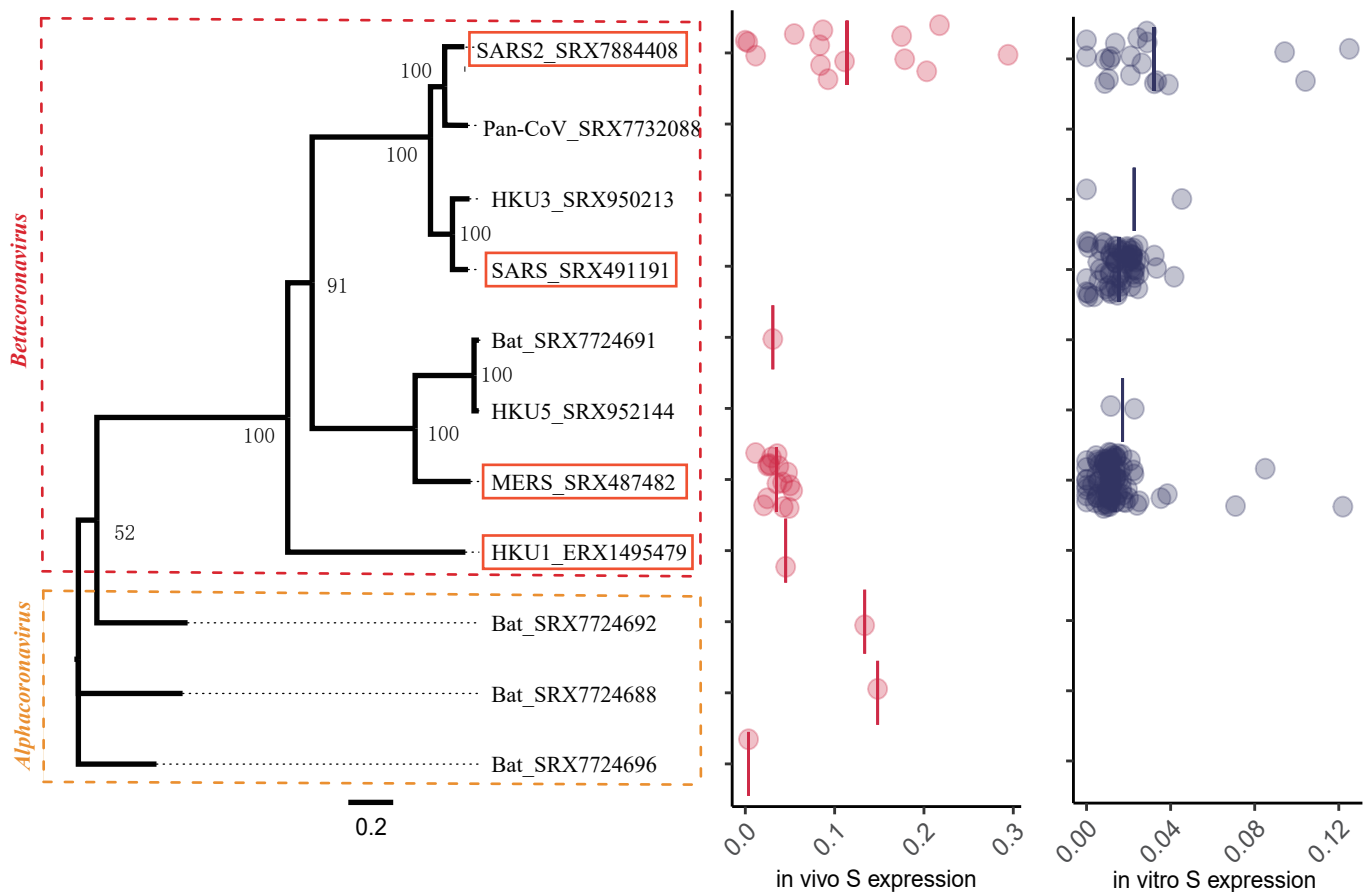
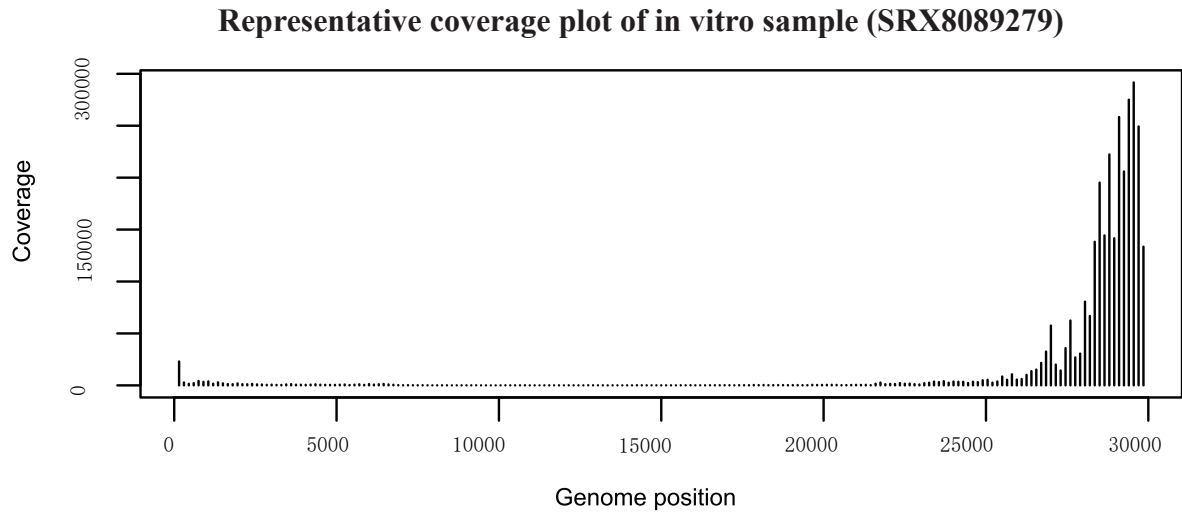


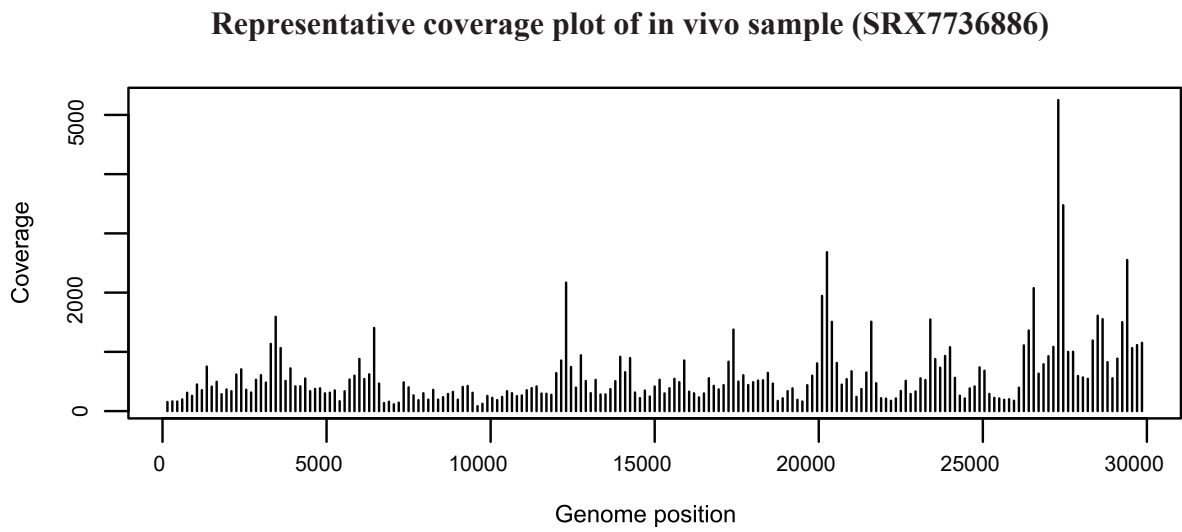
Fig. S1

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.16.043224>; this version posted November 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

A

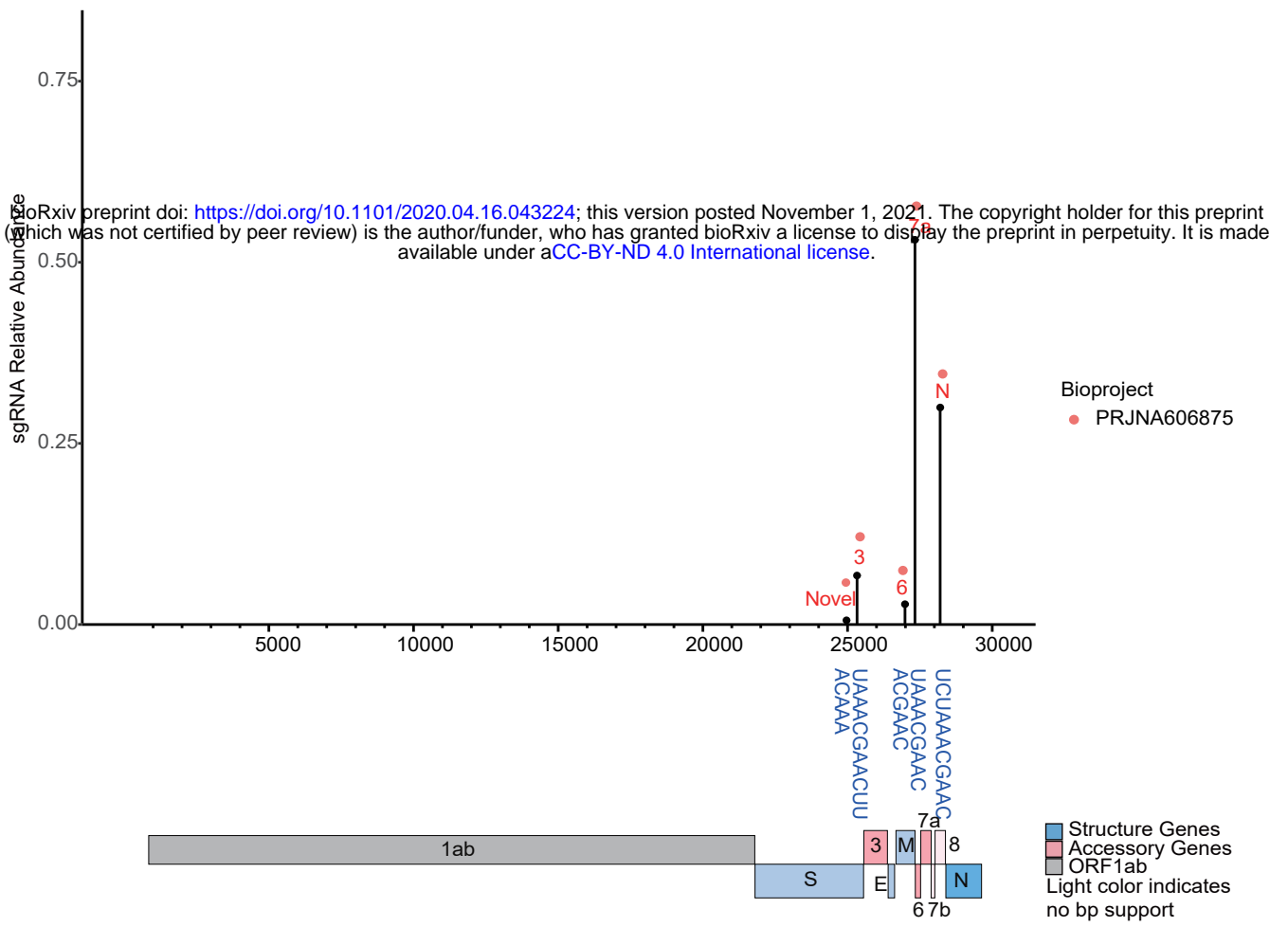


B



C

The sgRNA Profile of Pan-CoV



D

The sgRNA Profile of HKU1

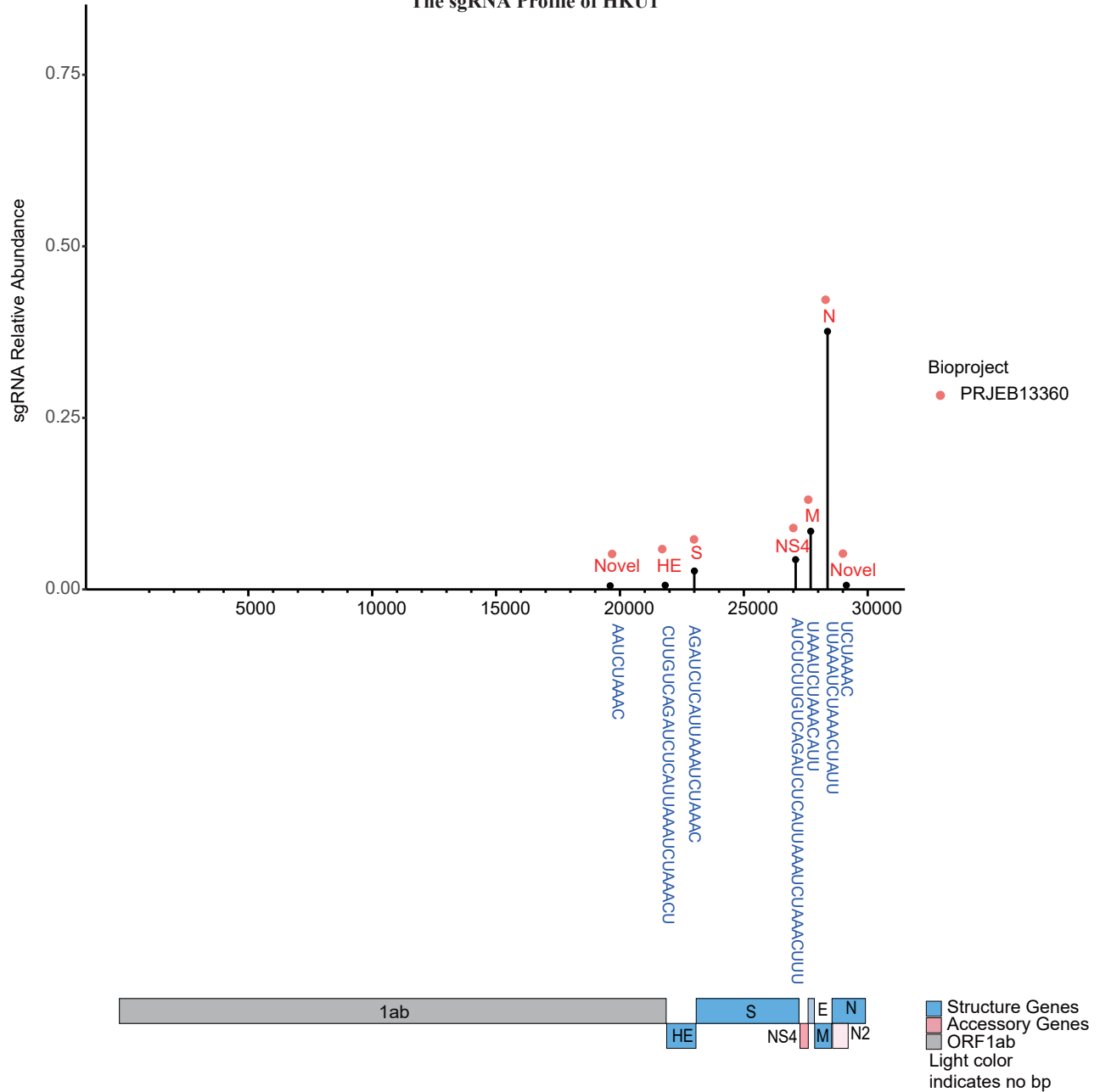


Fig. S3

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.16.043224>; this version posted November 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

A

```

Leader 51 UUGUAGAUCUGUUCUCUAAACGAACUUUAAAAUCUGUGUGG 91
          |||
sgRNA.cano UUGUAGAUCUGUUCUCUAAACGAACUAAUUAUUUAUUAGU
          |||
Body 26454 UCCUGAUCUUCUGGUCUAAACGAACUAAUUAUUUAUUAGU 26494

Leader 51 UUGUAGAUCUGUUCUCUAAACGAACUUUAAAAUCUGUGUGG 91
          |||
sgRNA.alt UUGUAGAUCUGUUCUGU-UUGGAACUUUAAUUUUAGCCAUG
          |||
Body 26486 UAUUAUAGUUUUUCUGU-UUGGAACUUUAAUUUUAGCCAUG 26525
                                     M start
    
```

B

SARS-CoV-2 tORF7b

```

Leader 51 UUGUAGAUCUGUUCUCUAAACGAACUUUAAAAUCUGUGUGG 91
          |||
sgRNA    UUGUAGAUCUGUUCUCUAAACGAACUUUCAUUAUUUGACUU
          |||
Body 27741 CAAAAGAAAGACAGAAUGAUUGAACUUUCAUUAUUUGACUU 27781
    
```

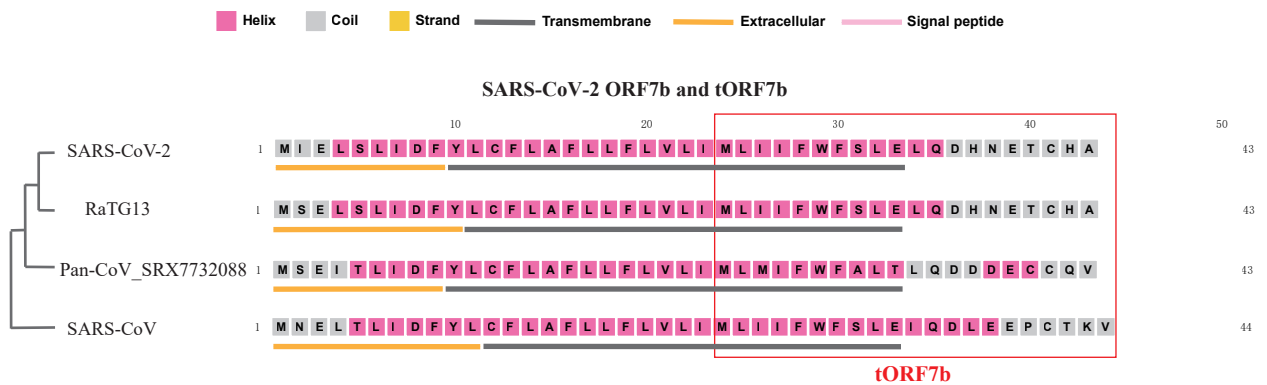
C

MERS-CoV pORF8c

```

Leader 51 ACUUUGAUUUUUAACGAACUUAAAUAAGCCUGUUGUUUA 91
          |||
sgRNA    ACUUUGAUUUUUAACGAACUUAAAAGAUCCAACUACAAUAA
          |||
Body 29562 GGAGCCAUUAACUUGACCCAAAGAAUCCAACUACAAUAA 29602
    
```

D



E

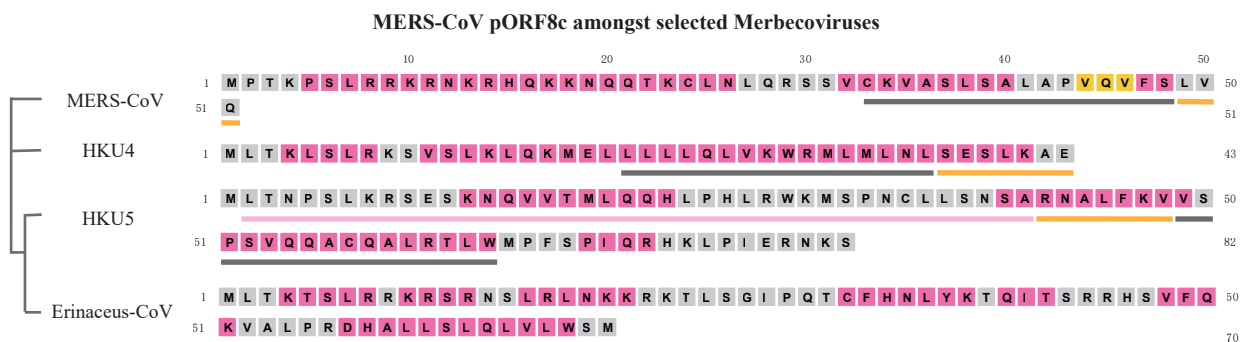
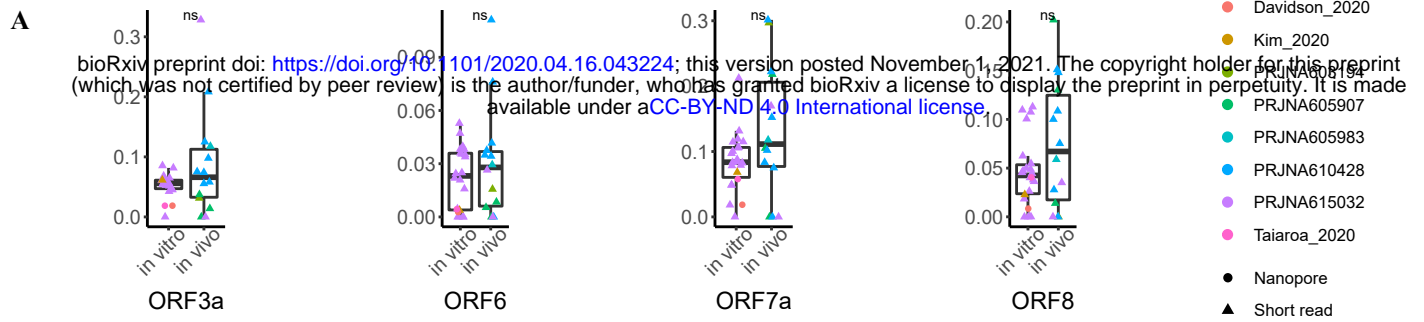
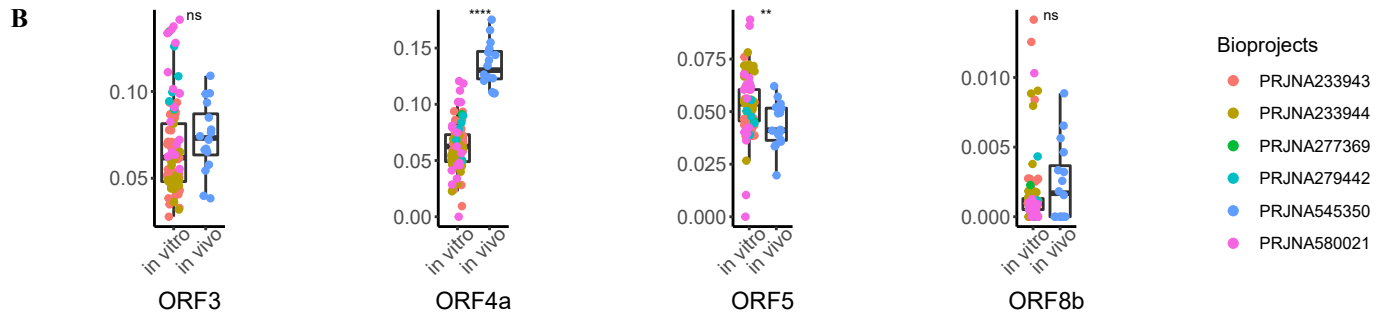
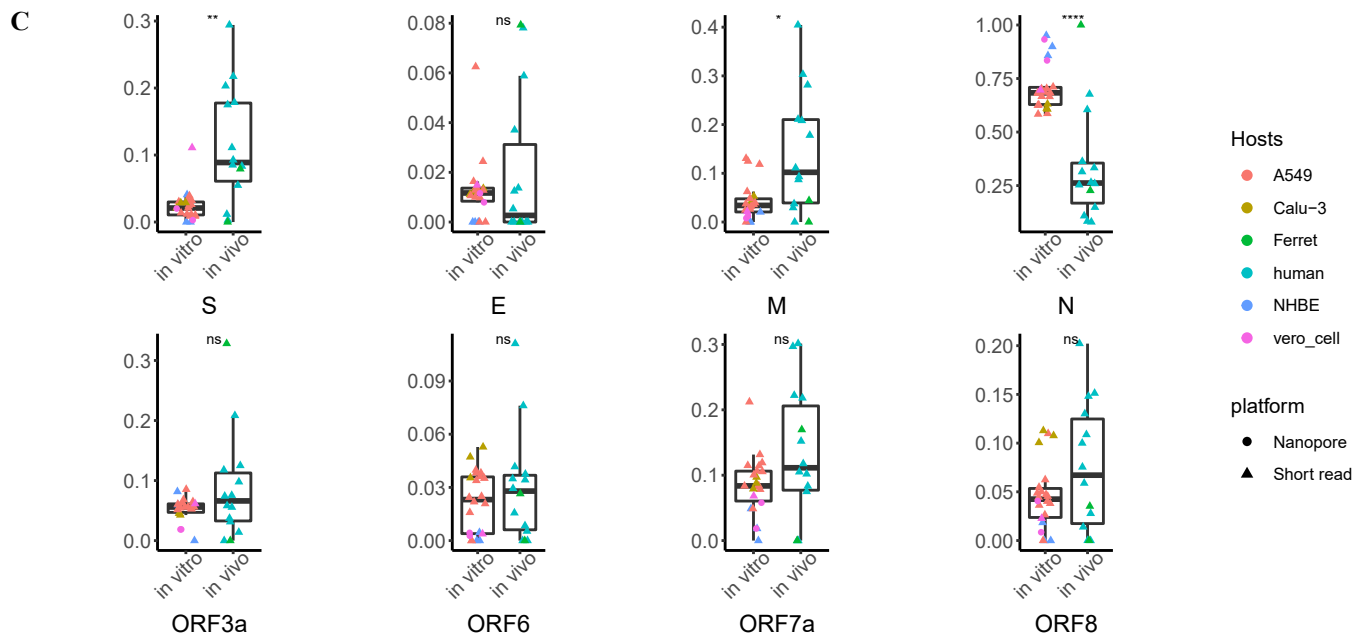
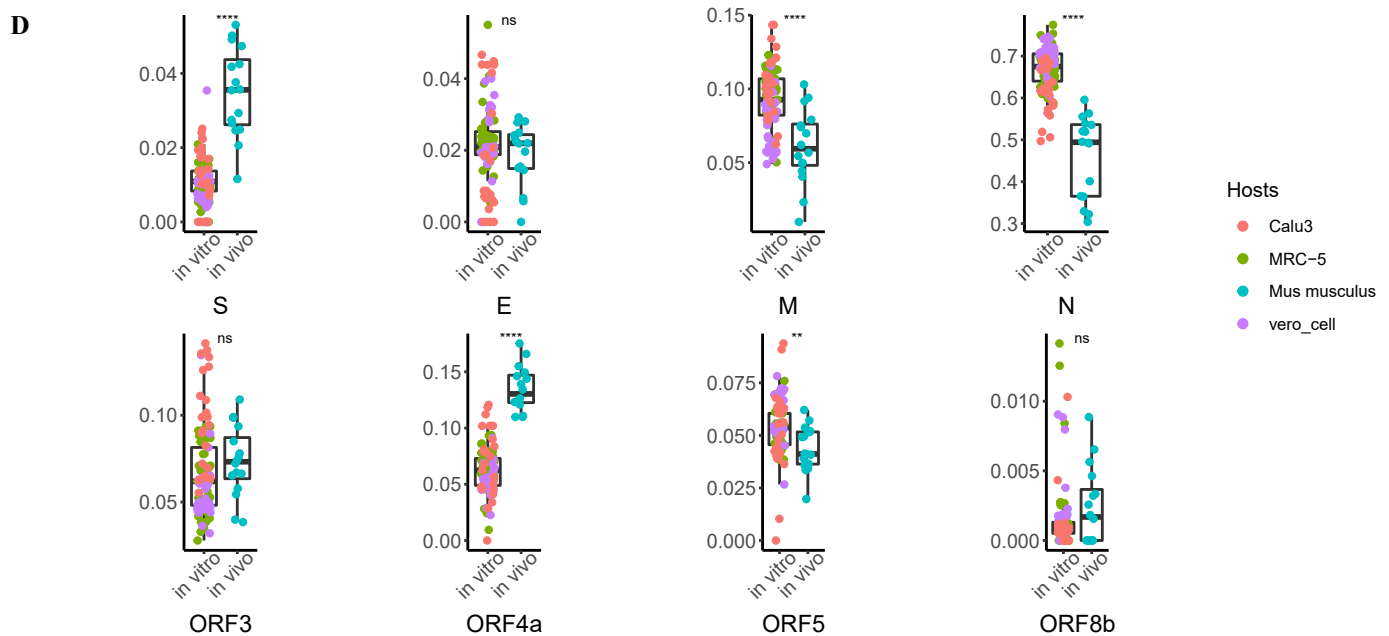


Fig. S4**Expression of Accessory proteins, In-vivo versus in-vitro, SARS-CoV-2****Expression of Accessory proteins, In-vivo versus in-vitro, MERS-CoV****In-vivo versus in-vitro, SARS-CoV-2, Colored by hosts****In-vivo versus in-vitro, MERS-CoV, Colored by hosts**

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.16.043224>; this version posted November 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

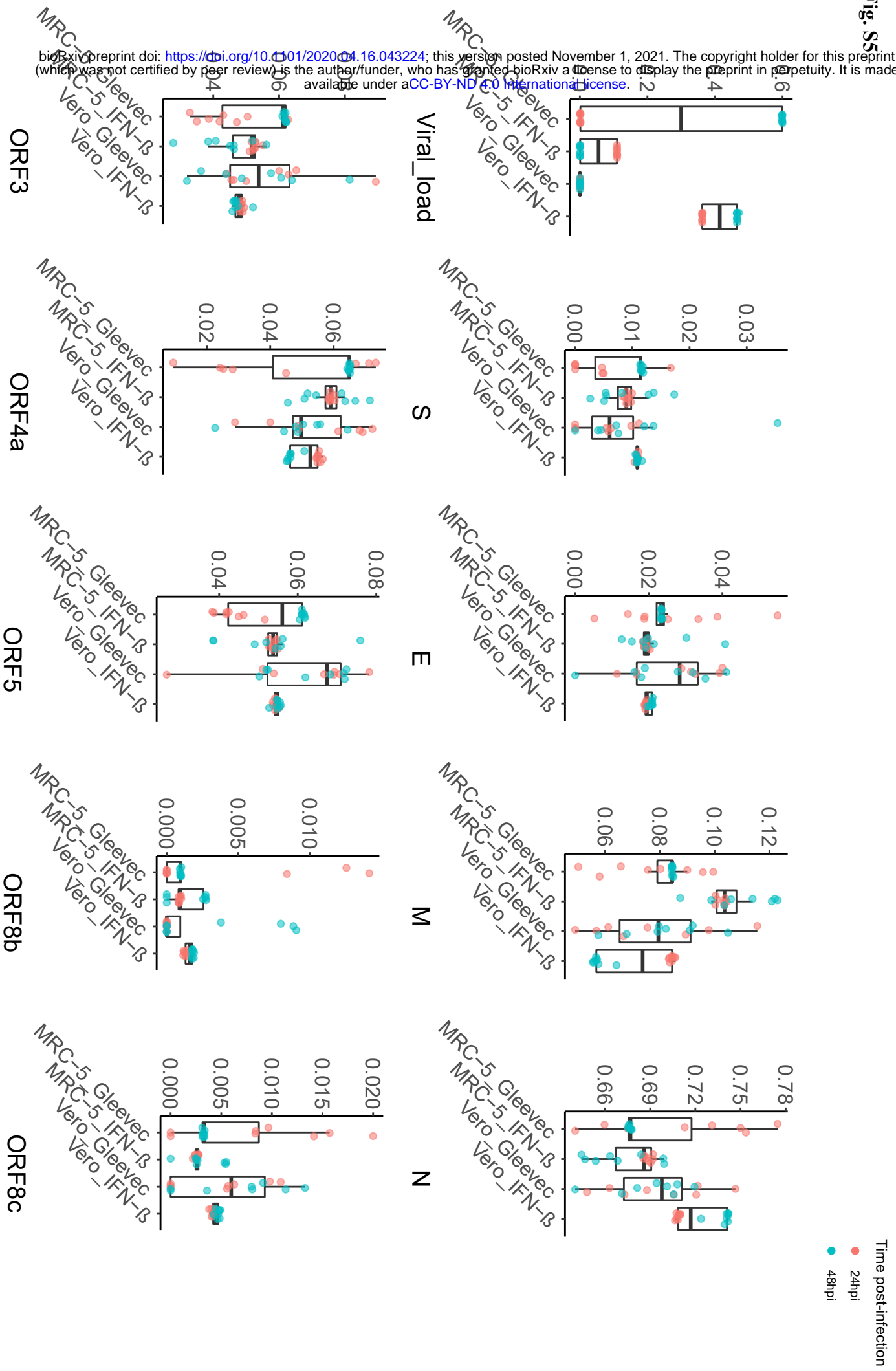


Fig. S6

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.16.043224>; this version posted November 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

