

1 **Positive selection within the genomes of SARS-CoV-2 and other**
2 **Coronaviruses independent of impact on protein function**

3
4

5 Alejandro Berrio¹, Valerie Gartner¹, Gregory A Wray^{1,2}

6
7
8

9 ¹ Biology, Duke University, Durham, North Carolina, United States

10 ² Center for Genomic and Computational Biology, Duke University, Durham, North Carolina,
11 United States

12

13 Corresponding Author:

14 Alejandro Berrio¹

15

16 130 Science Drive, Durham, North Carolina, 27708, United States

17 Email address: alejo.berrio@duke.edu

18

19 **Abstract**

20 **Background.** The emergence of a novel coronavirus (SARS-CoV-2) associated with severe
21 acute respiratory disease (COVID-19) has prompted efforts to understand the genetic basis for its
22 unique characteristics and its jump from non-primate hosts to humans. Tests for positive
23 selection can identify apparently nonrandom patterns of mutation accumulation within genomes,
24 highlighting regions where molecular function may have changed during the origin of a species.
25 Several recent studies of the SARS-CoV-2 genome have identified signals of conservation and
26 positive selection within the gene encoding Spike protein based on the ratio of synonymous to
27 nonsynonymous substitution. Such tests cannot, however, detect changes in the function of RNA
28 molecules.

29 **Methods.** Here we apply a test for branch-specific oversubstitution of mutations within narrow
30 windows of the genome without reference to the genetic code.

31 **Results.** We recapitulate the finding that the gene encoding Spike protein has been a target of
32 both purifying and positive selection. In addition, we find other likely targets of positive
33 selection within the genome of SARS-CoV-2, specifically within the genes encoding Nsp4 and
34 Nsp16. Homology-directed modeling indicates no change in either Nsp4 or Nsp16 protein
35 structure relative to the most recent common ancestor. Thermodynamic modeling of RNA
36 stability and structure, however, indicates that RNA secondary structure within both genes in the
37 SARS-CoV-2 genome differs from those of RaTG13, the reconstructed common ancestor, and
38 Pan-CoV-GD (Guangdong). These SARS-CoV-2-specific mutations may affect molecular
39 processes mediated by the positive or negative RNA molecules, including transcription,
40 translation, RNA stability, and evasion of the host innate immune system. Our results highlight
41 the importance of considering mutations in viral genomes not only from the perspective of their
42 impact on protein structure, but also how they may impact other molecular processes critical to
43 the viral life cycle.

44

45 **Introduction**

46 An important challenge in understanding zoonotic events is identifying the genetic changes that
47 allow a pathogen to infect a new host. Such information can highlight molecular processes in
48 both the pathogen and host that have practical value. The recent outbreak of SARS-CoV-2, a
49 novel coronavirus, provides both a challenge and an opportunity to learn more about the specific
50 adaptations that enable the virus to thrive in human hosts and that endow it with traits distinct
51 from previously described coronaviruses. Formal tests for natural selection are a powerful tool in
52 this endeavor because they can be applied in an unbiased manner throughout the viral genome:
53 evidence of negative selection can reveal regions of the genome that are broadly constrained
54 functionally and thus unlikely to contribute to species-specific traits, while evidence of branch-
55 specific positive selection can identify candidate regions of the genome where molecular
56 processes may have diverged from that of other species.

57 Several recent studies have tested for natural selection in the SARS-CoV-2 genome based on the
58 ratio of synonymous to non-synonymous (dN/dS) substitutions relative to other coronaviruses
59 (Tang et al., 2020; Chaw et al., 2020; Li et al., 2020a). The most prominent signal to emerge
60 from these studies is a mix of positive and purifying selection within the gene encoding the Spike
61 glycoprotein, which mediates invasion of host cells by binding to the angiotensin-converting

62 enzyme 2 (ACE2) receptor in host cells (Gallagher & Buchmeier, 2001; Tortorici & Veesler,
63 2019). This finding makes good biological sense, because structural changes in the spike protein
64 are common and are known to influence the ability of the virus to infect new hosts and jump
65 between species (Hulswit, de Haan & Bosch, 2016). A single nucleotide polymorphism (SNP)
66 that results in an amino acid substitution in Spike protein has increased in frequency during the
67 global pandemic more rapidly than other SNPs (Korber et al., 2020), leading to speculation that
68 it is an adaptation that alters the interaction between Spike and ACE2, FURIN and
69 TMPRSS2 (Eaaswarkhanth, Al Madhoun & Al-Mulla, 2020).

70 Beyond mutations that alter Spike protein, however, there exists little understanding of positive
71 selection within the SARS-CoV-2 genome and how this may have shaped viral traits. Few
72 convincing signals of positive selection exist for any of the other viral proteins (Cagliani et al.,
73 2020; Velazquez-Salinas et al., 2020; Chaw et al., 2020). For RNA viruses, however, critical
74 aspects of the life cycle rely on molecular processes that are not reflected in protein sequence. In
75 particular, in positive-strand RNA viruses such as coronaviruses, the single RNA molecule that
76 constitutes the genome is first transcribed and translated to produce the replicase polyprotein 1a
77 and 1ab that is cleaved into multiple non-structural proteins, some of which participate in the
78 assembly of a cellular structure known as the replicase-transcriptase complex (RTC), where the
79 proper environment for viral replication and transcription is created. Then, the RNA-dependent-
80 RNA-polymerase (RdRp or Nsp12) produces negative sense genomic and subgenomic RNAs
81 that are used as template strands that are then transcribed in the opposite direction to make more
82 positive-sense viral genomes and a variety of RNA molecules that are translated into structural
83 proteins for packaging (Fehr & Perlman, 2015; Kim et al., 2020). Although the viral proteins that
84 help mediate these processes are visible to tests for selection that rely on dN/dS ratios, the RNA
85 molecules with which they interact are not. This leaves the operation of natural selection on
86 important molecular functions within the viral life cycle largely unexamined.

87 Here we utilize a test for positive selection that identifies an excess of nucleotide substitutions
88 within a defined window in the genome relative to neutral expectation using a likelihood ratio
89 framework (Wong & Nielsen, 2004; Haygood et al., 2007). We implemented this test using
90 *adaptiPhy* (Berrio, Haygood & Wray, 2020) to infer regions of the genome that were likely
91 targets of branch-specific positive selection in several *Sarbecovirus* species from bat, pangolin,
92 and human hosts. Our results recapitulate results from dN/dS-based tests that highlight *S*, the

93 gene encoding Spike protein, as a prominent target of natural selection within the SARS-CoV-2
94 genome (Cagliani et al., 2020; Chaw et al., 2020; Li et al., 2020a). Importantly, we also identify
95 genomic regions not previously reported to be targets of positive selection. Based on structural
96 modeling of RNA and protein, we argue that these newly identified regions of positive selection
97 likely affect species-specific RNA, rather than protein, function. These genomic regions are
98 candidates for understanding the molecular mechanisms that endow SARS-CoV-2 with some of
99 its unique biological properties.

100 **Materials & Methods**

101 **Sequence Alignment**

102 To identify branch specific positive selection, it is necessary to obtain a query and a reference
103 alignment. We downloaded six high quality reference genomes from the subgenus Sarbecovirus
104 (Table 1). Next, we used MAFFT (Katoh & Standley, 2013) plugin in Geneious Prime v.2.1
105 (Kearse et al., 2012) with default settings to build a sequence alignment. Next, we refined the
106 alignment using a gene by gene procedure.

107

108 **Testing for Positive Selection**

109 Despite *adaptiPhy* was originally designed to investigate regions of complex genomes under
110 positive selection, it can be used to identify regions of a viral sequence alignment where the
111 foreground branch is evolving at faster rates than the expectation from the background species.
112 We performed a selection analysis on sliding windows of 300 bp with a step of 150 bp along a
113 sequence alignment of 5 reference genome sequences of coronaviruses of the subgenus
114 *Sarbecovirus* and two sequences of Pangolin Coronavirus recently published (Liu, Chen & Chen,
115 2019; Lam et al., 2020). This procedure generates partitions where a tree topology can be fitted.
116 To investigate the extent of positive selection or branches with long substitution rates along the
117 SARS-CoV-2 genome, we used a branch-specific method known as *adaptiPhy* that was initially
118 developed in 2007 (Haygood et al., 2007) and recently improved (Berrio, Haygood & Wray,
119 2020). This computational methodology makes use of a likelihood ratio test based on the
120 maximum likelihood estimates obtained from HyPhy v2.5 (Pond, Frost & Muse, 2005; Pond et
121 al., 2020). The branch of interest (e.g., SARS-CoV-2 branch) is used as the foreground and the
122 rest of the alignment is used as the background. To obtain data from nucleotide substitutions
123 alone, we used *msa_split* from PHAST (Hubisz, Pollard & Siepel, 2011) to remove insertions

124 and any sequence gaps that were present in the genomes of the background virus species relative
125 to the SARS-CoV-2 genome. The assumption for the background species is the same for both the
126 null and alternative models; specifically, only neutral evolution and negative (purifying)
127 selection are permitted. While in the foreground, the assumptions are the same as for the
128 background in the null model. In the alternative model, all three types of evolution are permitted
129 (neutral evolution, negative selection, and positive selection) in the foreground of the following
130 topology: (((((SARS_CoV_2, Bat_CoV_RaTG13), Pa_CoV_Guangdong),
131 Pa_CoV_Guangxi_P4L), (Bat_CoV_LYRa11, SARS_CoV)), Bat_CoV_BM48). This method is
132 highly sensitive and specific and can differentiate between positive selection and relaxation of
133 constraint. *AdaptiPhy* requires at least 3 kb reference alignment for each species that is used as a
134 putatively neutral proxy for computing substitution rates. Viruses' genomes lack non-functional
135 regions, therefore, the most reasonable proxy for neutral evolution has to be found in the regions
136 outside the query window. To do this, we concatenated twenty regions of 300 bp of the viral
137 genome alignment that were drawn randomly with replacement from the entire genome
138 alignment. Then, for each query alignment, we built a reference alignment of 6 kb as it produces
139 a stable evolutionary standard of recombination rates. To control for the stochasticity of the
140 evolutionary process, we run each query against ten bootstrapped samples of reference
141 alignments. Finally, we used a custom R script to compute the likelihood ratio, which was used
142 as a test statistic for a chi-squared test with one degree of freedom to calculate a *P*-value for each
143 query. Then, we corrected the distribution of all *P*-values per query region using the *p.adjust()* R
144 function with the *fdr* method. Next, we classified a query window to be under positive selection
145 if the *P*-adjusted value was < 0.05 . We were unable to successfully run *adaptiPhy* on two
146 windows because the outgroup species (Bat_CoV_BM48) contained a deletion of 406 bp relative
147 to SARS-CoV-2, which spans the entire ORF8.

148 Next, we calculated the distribution of substitution rates for each branch and nodes in each query
149 and reference sequence using *phyloFit* (Hubisz, Pollard & Siepel, 2011). To visualize the
150 strength of selection comprehensively, we computed the statistic ζ (zeta), representing the
151 evolutionary rate. To calculate this rate, we compared the substitution rate in the query with their
152 respective reference alignments. This parameter, “ ζ ”, is analogous to ω (omega), the ratio of
153 dN/dS, where a value of $\omega < 1$ indicates constraint or negative selection; a value of $\omega = 1$ indicates
154 neutrality; and a value of $\omega > 1$ indicates positive selection (Wong & Nielsen, 2004).

155 **Testing for Conservation**

156 To test for conservation, we used the *phastCons* computational method from PHAST (Siepel et
157 al., 2005; Hubisz, Pollard & Siepel, 2011). To run this tool, we used the models obtained with
158 *phyloFit* for the reference alignments and then, we generated an average estimate of the
159 conserved and non-conserved states of the models with *phyloBoot* (Hubisz, Pollard & Siepel,
160 2011). Finally, we run the final analysis using *phastCons* on the query alignments using the
161 previous models to generate *phastCons* values for each base-pair along the sequence. To plot
162 these we took the average from each alignment and plot it using the library Gviz and
163 Bioconductor (Hahne & Ivanek, 2016) in R.

164 **Testing for Recombination**

165 Inference of branch specific selection can be confounded by recombination given that a single
166 phylogenetic tree may not explain the evolution of viruses. Recombination is common in
167 coronaviruses (Hon et al., 2008; Graham & Baric, 2010; Lau et al., 2015; Hu et al., 2017; Li et
168 al., 2020b; Lam et al., 2020) and it should be accounted for as an alternative explanation of
169 selection at the nucleotide level. Here, we screened for evidence of recombination in two ways,
170 one, by estimating phylogenetic trees in sliding windows of 500 bp and a step of 150 along
171 coronavirus alignment using RaXML-NG v0.9 (Kozlov et al., 2019).

172 **Evaluating polymorphic diversity in the pandemics of 2020**

173 We downloaded complete sequences of SARS-CoV-2 genomes from the NCBI Virus database
174 (https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/virus?SeqType_s=Nucleotide). As of June 26,
175 2020, we obtained and aligned 5597 SARS-CoV-2 genomes sequenced worldwide. To align
176 these sequences, we used MAFFT (Kato & Standley, 2013) plugin from Geneious Prime 2.1
177 (Kearse et al., 2012), eliminating 597 sequences with the highest number of differences and
178 ambiguities relative to the reference sequence (RefSeq: NC_045512.2), for a total of 5,000
179 sequences. Next, we estimated the frequency of SNP variants using the Find Variations/SNPs
180 tool with a minimum coverage of 4,900 sequences and a minimum frequency of 0.01, to identify
181 nucleotide variants among a subset of high quality sequenced genomes in order to evaluate
182 ongoing evolution in the regions under positive selection.

183 **Analysis of RNA and Protein structures**

184 To investigate potential structural changes in Nsp4 and Nsp16 at both the RNA and protein level,
185 we performed minimum free energy (MFE) prediction analysis using the RNAfold WebServer

186 (Gruber et al., 2008; Lorenz et al., 2011) and consensus homology modeling using PHYRE2's
187 intensive mode (Kelley et al., 2016). These analyses were performed for both Nsp4 and Nsp16
188 sequences for SARS-CoV-2, Bat-CoV_RaTG13, Pan-CoV-Guangdong, and SARS-CoV.
189 RNAfold uses a loop-based energy model and a dynamic programming algorithm to predict the
190 structure of the sequence such that the free energy of the structure is minimized. The RNAfold
191 WebServer generates graphical outputs for both the MFE and Centroid structures, which display
192 the base pairing probabilities by color (blue = 0, red = 1). These two MFE structures correspond
193 to the MFE and the Centroid traces in the mountain plot, which is a positional representation of
194 the secondary structure. In our figures, we show the MFE structure prediction.
195 PHYRE2 aligns input protein sequences using Position-Specific Iterated BLAST (PSI-BLAST)
196 against sequences of experimentally resolved protein structures. A 3D model of the input
197 sequence is then constructed based on homology-matched templates, optimizing for greatest
198 sequence coverage and highest confidence. Regions of the input sequence without a matching
199 template sequence are modeled *ab initio* and with Poing, a multi-template modelling tool.
200 Pairwise comparisons of predicted protein structures were visualized using PyMOL software
201 (DeLano, 2002). Alignment and structural comparisons performed by FATCAT (Ye & Godzik,
202 2004).

203 **Results**

204 **Positive and negative selection are highly localized within coronavirus genomes**

205 We tested for branch-specific selection on nucleotide sequences in coronavirus genomes,
206 focusing on six species from the *Sarbecovirus* Subgenus (Coronaviridae family) and Bat-CoV-
207 BM48-31/BGR/2008 as an outgroup. Using 300 bp windows with a step size of 150 bp, we
208 scanned the genome alignment for concentrations of fixed mutations that exceed the neutral
209 expectation based on the genome as whole relative to that particular window's evolutionary
210 history among the seven species. This test identifies regions of the genome showing the most
211 extreme divergence in nucleotide sequence on a particular branch relative to its specific
212 background rate of evolution across the entire phylogeny and without reference to the genetic
213 code. Fig 1A shows windows of inferred positive selection (red dots) on the branch leading to
214 each species. The results reveal several signals of positive selection that are unique to a single
215 species and others that are recapitulated in multiple species. The latter finding suggests that some
216 segments of the viral genome have repeatedly experienced adaptive modification. In general, the

217 distribution of positive selection is more similar in closely related species than in divergent ones
218 (Figs 1A and 1B), suggesting that some molecular functions have been altered over an interval
219 that extends beyond the origin of a single species but not across the entire Sarbecovirus
220 radiation.

221 Next, we identified regions of the genome that are highly conserved across the sarbecovirus
222 genomes examined in this study using PhastCons (Siepel et al., 2005) (Fig 2A). As with positive
223 selection, conservation is highly localized (Figs 2A and 2B). Based on a criterion of PhastCons >
224 0.9, we found high levels of conservation in regions encoding seven proteins: 3CL-Pro, Nsp6,
225 Nsp8, Nsp9, Nsp10, Nsp11, RNA-dependent RNA polymerase (RdRp), Protein 3a, Nucleocapsid
226 phosphoprotein, and Envelope (Figs 2A-D). These loci of exceptional sequence conservation
227 highlight critical molecular features: Nucleocapsid and Envelope are essential structural proteins
228 of the coronavirus capsid, while the other proteins regulate a variety of molecular process during
229 viral replication (Tan et al., 2005; Lu et al., 2006; Minakshi et al., 2009; Freundt et al., 2010;
230 Fuchs, 2012; Yue et al., 2018).

231 **The gene encoding Spike protein is under persistent positive selection**

232 In all ingroup species examined we detected signals of positive selection within the S
233 gene, which encodes the Spike protein, and in five of the six species this was the most prominent
234 signal in the entire genome (Fig 1A). This finding confirms previous studies that used the dN/dS
235 ratio to test for selection on protein function (Tang et al., 2020; Chaw et al., 2020; Li et al.,
236 2020a). Interestingly, we observed that the specific regions showing signatures of positive
237 selection differed between species (Figs 1B and 1C). In SARS-CoV-2, we detected signals of
238 positive selection in four segments of the S gene. First, the region encoding the entire receptor
239 binding domain (RBD) shows an extended signal (Figs 1B-C and 3A); as others have noted,
240 structural changes in this region may improve binding to human ACE2 (Wang et al., 2020;
241 Wrapp et al., 2020; Wang, Liu & Gao, 2020). The second segment encodes another externally
242 facing region, the S1 subunit N-terminal (NTD) domain, which includes the first disulfide bond
243 (amino acids 13 - 113) and several glycosylation sites. The third signal of positive selection
244 within S is located around the derived furin cleavage site (amino acids 664 - 812) that has been
245 found to be essential for infection of lung cells (Hoffmann, Kleine-Weber & Pöhlmann, 2020).
246 The fourth signal is located in a segment encoding the S2 and S2' subunits that includes the
247 Heptad repeat 2 (amino acids 1114 - 1213). These heptad repeats were previously associated

248 with episodes of selection for amino acids that increase the stability of the six-helix bundle
249 formed by both heptad repeats in MERS and other coronaviruses (Forni et al., 2015); they are
250 also thought to determine host expansions and therefore, facilitate virus cross-species
251 transmission (Graham & Baric, 2010).

252 The distribution of inferred positive selection in the S gene of SARS-CoV differed from that of
253 SARS-CoV-2 described above. Notably, there was no signal in the ACE2 binding domain (Figs
254 1B and 1C). Moreover, a signal was present throughout the N-terminal domain and in the
255 boundary region between the S1 and the S2 subunits (Fig 1), a region that includes the
256 proteolytic cleavage (M de Haan et al., 2004). Interestingly, this region evolved a novel furin
257 cleavage site in SARS-CoV-2 that may increase the cleavage efficiency and cell-cell fusion
258 activity and changes in the virulence of the virion as seen in mutant studies of SARS-CoV and
259 SARS-CoV-2 (Follis, York & Nunberg, 2006; Hoffmann, Kleine-Weber & Pöhlmann, 2020).

260 **Genes encoding Nsp4 and Nsp16 contain branch-specific signals of positive selection**

261 We also detected two shorter signals of positive selection within the SARS-CoV-2
262 genome that are located outside of the S gene, in ORF1a and ORF1b (Fig 1A). Interestingly,
263 both encode small proteins that contribute to viral replication. The first is Nsp4, which encodes a
264 membrane-bound protein with a cytoplasmic C-terminal domain; it is thought to anchor the
265 Viral-Replication-Transcription Complex (RTC) to the modified endoplasmic reticulum
266 membranes in the host cell (Oostra et al., 2008; Hagemeyer et al., 2011, 2014; Snijder, Decroly
267 & Ziebuhr, 2016). The SARS-CoV-2 Nsp4 protein differs from that of closely related
268 sarbecoviruses by two nearly adjacent amino acids: V380A and V382I. Although this region of
269 the genome as a whole is not highly conserved (Fig 2), both of these positions are V residues in
270 all of the in-group species we examined except SARS-CoV-2 (Fig 3B and 4A). This signal is too
271 weak to be scored as positive selection using dN/dS-based tests (Sharp, 1997; Nielsen, 1997;
272 Yang & Nielsen, 2008) and indeed may not affect protein function given the biochemically
273 similar side-chains of the amino acids involved.

274 The second signal of positive selection outside of the S gene lies within Nsp16. This gene
275 encodes a 2'-O-methyltransferase that modifies the 5'-cap of viral mRNAs (Decroly et al., 2008;
276 Bouvet et al., 2010) and assists in evasion of the innate immune system of host cells (Züst et al.,
277 2011; Menachery, Debbink & Baric, 2014; Nelemans & Kikkert, 2019). Of note, this is the only
278 signal of positive selection within the SARS-CoV-2 genome that lacks any nonsynonymous

279 substitutions (Fig 2). All of the nucleotide substitutions in Nsp16 during the origin of SARS-
280 CoV-2 are synonymous, while the Nsp16 genes of SARS-CoV-2, Bat-Cov-RaTG13, and Pan-
281 CoV-GD (Guangdong) all encode identical proteins (Figs 3C and 5A). This suggests a complex
282 mechanism of selection in the form of purifying selection at the protein level and positive
283 selection at the nucleotide level. Ancestral state reconstruction of Nsp16 indicates that 20
284 synonymous substitutions likely occurred in the lineage leading to SARS-CoV-2 after the split
285 from the common ancestor with BatCoV-RaTG13, while 19 substitutions are synonymous
286 substitutions that occurred in the lineage leading to Bat-CoV-RaTG13 (Supplementary data).
287 Eleven of these twenty substitutions are concentrated within the region scoring high for positive
288 selection in SARS-CoV-2 and twelve within the positively selected region in Bat-CoV-RaTG13.
289 Surprised by these findings, we hypothesized that the Nsp16 RNA secondary structure may
290 differ among species in ways that affect molecular functions mediated directly (although not
291 solely) by RNA, such as replication, transcription, translation, or evasion of the host immune
292 system. To investigate this possibility, we first compared the secondary structure and minimum
293 free energy (MFE) of RNA in the vicinity of Nsp4 and Nsp16 among the genomes of SARS-
294 CoV-2, Bat-CoV-RaTG13, Pan-CoV-GD, and SARS-CoV using RNAfold (Gruber et al., 2008).
295 Both the predicted secondary structures and mountain plots, which show the free energy
296 predictions along the length of the sequence by position, reveal differences in RNA folding
297 dynamics across the four species (Figs 4B and 5B). Analysis of the reconstructed sequence of the
298 SARS-CoV-2 + Bat-CoV-RaTG13 ancestor reveal that most of these differences evolved during
299 the origin of SARS-CoV-2 (S1 Fig). These differences among species in predicted secondary
300 structures within Nsp4 and Nsp16 stand in contrast to the 5' UTR, which is thought to fold into a
301 stable secondary structure that is markedly conserved among Sarbecovirus species (S3 Fig).
302 Together, these observations suggest that the signal of positive selection within Nsp16 in the
303 SARS-CoV-2 genome may reflect changes in RNA, rather than protein, function that are unique
304 to this species of coronavirus.

305 While the focus here is on SARS-CoV-2, it is worth noting that we also detected signals of
306 positive selection outside of the S gene in the other sarbecovirus genomes examined here. The
307 distribution of positive selection in the genome of SARS-CoV, for instance, shows some
308 similarities to, but also notable differences from, that of SARS-Cov-2 (Fig 1). In both species, S
309 and Nsp16 contain signals of positive selection, although in distinct regions of the two genes (Fig

310 1). In addition, the genome of SARS-CoV contains signals of positive selection in Nsp2, Nsp3,
311 and ORF3a, none of which shows elevated rates of substitution in SARS-CoV-2. The first two
312 genes encode proteins with important roles in viral replication: Nsp2 may disrupt intracellular
313 signaling in the host cell (Cornillez-Ty et al., 2009) while Nsp3 cleaves itself, Nsp1, and Nsp2
314 from the replicase polyproteins (Báez-Santos, St. John & Mesecar, 2015), assists in the assembly
315 of the double membrane vesicles of the RTC system (Hagemeijer et al., 2014), and antagonizes
316 the host innate immune response (Tsuchida, Kawai & Akira, 2009; Frieman et al., 2009;
317 Matthews et al., 2014).

318 **Recombination does not account for most signals of positive selection**

319 Recombination from another species can be a confounding factor in the inference of positive
320 selection using the framework employed here, because the inserted genomic segment may be
321 more divergent than the rest of the foreground genome is from nearby background species.
322 Several instances of recombination have been reported in coronaviruses, including SARS-CoV-2
323 (Hon et al., 2008; Lam et al., 2020; Boni et al., 2020; Li et al., 2020a), making it important to
324 distinguish regions of recombination from positive selection. The two processes produce distinct
325 genetic signatures, with recombination the result of a single event (possibly later further
326 recombined) and positive selection as detected here the result of multiple independent mutations
327 that were fixed over an extended interval and are spatially concentrated. In order to test for
328 regions of the SARS-CoV-2 that contain recombined segments from other species, we estimated
329 the phylogenetic history of 500 bp segments of the genome with a step size of 150 bp among the
330 aligned genomes of the seven species examined in this study. We used RaXML-NG v0.9
331 (Kozlov et al., 2019) to reconstruct topology for each segment independently and searched for
332 cases where the topology differed from the expected topology based on the entire genome: (Bat-
333 CoV-BM48, ((Bat-CoV-LYRa11, SARS-CoV), (Pa-CoV-GX, (Pa-CoV-GD, (SARS-CoV-2, Ba-
334 CoV-RaTG13))))). Recombination from a divergent species should produce an incongruent
335 topology in one or more adjacent windows, revealing a recombined region and its approximate
336 breakpoints. We identified 12 regions where the topology differed from the expected (Fig 6). Of
337 note, these regions are somewhat more concentrated in the part of the genome that encodes
338 structural proteins. Consistent with a previous report (Li et al., 2020a), we observed overlap
339 between regions scoring high for positive selection and recombination in S, the gene encoding
340 Spike protein (Fig 6M), specifically the region that encodes for the ACE2 binding domain and a

341 region that includes the furin-cleavage site (Figs 6F-G). Importantly, however, none of the
342 putatively recombined regions overlap with the windows scoring high for positive selection
343 within the genes encoding Nsp4 and Nsp16 proteins in SARS-CoV-2.

344 **Recent changes in allele frequency may result from positive selection and hitch-hiking**

345 To gain insight into the evolutionary mechanisms that have shaped genetic variation more
346 recently within the SARS-CoV-2 genome, we compiled a list of known mutations, based on
347 5,000 accessions sequenced since the beginning of the current pandemic (see Methods). As
348 expected, the vast majority of variants are singletons, representing either mutations that are not
349 segregating or sequencing errors. The density distribution of polymorphisms (regardless of
350 frequency) is elevated within 2-3 kb at both ends the SARS-CoV-2 genome (Fig 2D). The site-
351 frequency spectrum is strongly left-skewed (S2 Fig). Given that the effective population size of
352 SARS-CoV-2 is likely very large, this distribution suggests that most SNPs are not subject to
353 positive selection and that negative selection prevents most new mutations from rising in
354 frequency due to drift. However, we did observe four SNPs that are present at high alternative
355 allele frequency (AAF > 0.6) (Fig 2C), a situation that can reflect positive selection, drift, or
356 hitch-hiking. Interestingly, all four of these SNPs are in tight LD (Toyoshima et al., 2020), which
357 suggests that positive selection on one of them may have driven the other three to high frequency
358 due to hitch-hiking.

359 We next investigated the likely consequences for altered molecular function due to each of these
360 four high-frequency derived SNPs. Two are located within regions of the genome that are highly
361 conserved among sarbecovirus species (Figs 2B-C). The first is a C>U substitution at position
362 241 in the 5'UTR, a region of the genome where RNA secondary structure is highly conserved
363 across Coronavirus species (Madhugiri et al., 2016; Rangan et al., 2020; Alhatlani, 2020). Using
364 RNAfold (Gruber et al., 2008) we found that this C>U transition had no impact on the stem-loop
365 structure established for SARS-CoV (S3 Fig). The other mutation in a conserved region of the
366 genome is a nonsynonymous substitution in the RdRp gene (14,408; P323L) at the interface
367 domain, which is thought to mediate protein-protein interactions (Pachetti et al., 2020; Hillen et
368 al., 2020). Because proline residues can influence secondary structure, we used PHYRE2 to
369 predict the impact of the P232L mutation on protein structure. Comparison of the two predicted
370 structures using FATCAT shows they are nearly identical (Table S1). The other two high-
371 frequency derived SNPs are located in regions that are neither highly conserved nor highly

372 divergent. One is a synonymous SNP in Nsp3 (3,037) and the other a nonsynonymous SNP in S
373 (23,403; D614G). This last SNP effectively removes a charged side-chain between the receptor
374 binding domain and the furin cleavage site of S, a region of recurrent positive selection among
375 the Sarbecovirus species we examined. Thus, of the four high-frequency derived SNPs, the
376 nonsynonymous substitution in S the most plausible candidate for altering molecular function
377 and thus becoming a target of natural selection.

378

379 **Discussion**

380 A crucial feature contributing to the global spread of Covid19 is that viral shedding starts before
381 the onset of symptoms (He et al., 2020); in contrast, shedding began two to ten days after the
382 onset of symptoms during the SARS epidemic of 2003 (Peiris et al., 2003; Pitzer, Leung &
383 Lipsitch, 2007). This striking difference suggests that one or more molecular mechanisms during
384 host cell invasion, virus replication, or immune avoidance may have changed during the origin of
385 SARS-CoV-2. Mutations contributing to viral transmission would likely be favored by natural
386 selection, making tests for positive selection a useful tool for identifying candidate genetic
387 changes responsible for the unique properties of SARS-CoV-2. Here, we searched for regions of
388 possible positive selection within the genomes of six coronavirus species, including SARS-CoV
389 and SARS-CoV-2. The method we used tests for an excess of branch-specific nucleotide
390 substitutions within a defined window relative to a neutral expectation for divergence in that
391 window and without regard to the genetic code (Wong & Nielsen, 2004; Haygood et al., 2007;
392 Berrio, Haygood & Wray, 2020).

393 Several prior studies have identified *S*, the gene encoding the Spike glycoprotein, as a target of
394 recurrent positive selection in coronavirus genomes, including SARS-CoV-2, based on ω , the
395 ratio of synonymous to nonsynonymous substitutions (Andersen et al., 2020; Cagliani et al.,
396 2020; Tang et al., 2020; Armijos-Jaramillo et al., 2020; Lam et al., 2020; Li et al., 2020a). *S*
397 thus serves as a positive control for our ability to detect signals of positive selection using a
398 different approach, which uses a likelihood ratio framework to identify regions of elevated,
399 branch-specific nucleotide substitution rates relative to a model that allows only drift (Wong &
400 Nielsen, 2004; Haygood et al., 2007; Berrio, Haygood & Wray, 2020). Consistent with this
401 expectation, we found that portions of the gene encoding Spike showed a striking elevation of
402 sequence divergence relative to the rest of the genome on the branches leading to all six species

403 examined. The specific regions of S containing high divergence differs markedly, however,
404 among species (Fig 1B). In SARS-CoV and Bat-CoV-LYRa11, these regions include the N-
405 terminal region, which contains glycosylation sites important for viral camouflage (Watanabe et
406 al., 2019; Yang et al., 2020) and a site of proteolytic cleavage that allows entry into the host cell
407 (Belouzard, Chu & Whittaker, 2009) (Fig 1C and 3A). In contrast, signals of positive selection in
408 SARS-CoV-2 and Bat-CoV-RaTG13 are concentrated in the domain that mediates binding to the
409 host receptor ACE2 (Fig 1C and 3A). These distinct distributions suggest that modifications in
410 different aspects of Spike function took place as various coronaviruses adapted to novel hosts. In
411 particular, the concentration of derived amino acid substitutions in the receptor binding domain
412 of Spike (Figs 1B and 1C) in SARS-CoV-2 and Bat-CoV-RaTG13 may reflect selection for
413 amino acid substitutions that result in higher affinity for ACE2 protein in different host species.
414 Importantly, we also detected signals of positive selection in two additional regions of the
415 SARS-CoV-2 genome, specifically within the genes encoding Nsp4 and Nsp16 (Fig 1A). Of
416 note, the Nsp16 region also shows a parallel signal of positive selection on the branch leading to
417 SARS-CoV. To our knowledge, this is the first report of possible adaptive change in molecular
418 function during the evolutionary origin of SARS-CoV-2 outside of the gene encoding Spike
419 protein. Prior scans for positive selection within the SARS-CoV-2 genome used elevated ω as
420 the signal of positive selection, which restricts attention to positive selection based on changes in
421 protein function. For coronaviruses this is a notable limitation, given that many aspects of the
422 lifecycle involve RNA function (Madhugiri et al., 2016; Ziv et al., 2020; Alhatlani, 2020). In
423 addition, the secondary structure of some segments within the RNA genome is well conserved
424 among coronavirus species, which implies a functional role (Rangan et al., 2020; Sanders et al.,
425 2020; Huston et al., 2020a). Indeed, the SARS-CoV-2 genome is reported to contain more well-
426 structured regions than any other known virus, including both coding and noncoding regions of
427 the genome (Huston et al., 2020a). We therefore examined nucleotide substitutions within
428 regions of putative positive selection in Nsp4 and Nsp16 for their likely impact on both protein
429 and RNA structure (Fig 4 and 5).

430 In the case of Nsp4 protein, two nearly adjacent nonsynonymous substitutions at residues 380
431 and 382 occurred on the branch leading to SARS-CoV-2 (Fig 3B). These both involve changing
432 side chains with similar biochemical properties, respectively valine to alanine and valine to
433 isoleucine. Homology-directed modeling of protein structure suggests that these two amino acid

434 substitutions have very little impact on either secondary or tertiary structure when comparing the
435 SARS-CoV-2 protein orthologue to those of the other species examined (Fig 4A). In the case of
436 Nsp16 protein, no nonsynonymous substitutions evolved on the branch leading to SARS-CoV-2.
437 Thus, the signal of positive selection within Nsp4 is unlikely to reflect changes in protein
438 structure or function, while the signal within Nsp16 cannot affect either because the encoded
439 polypeptide is identical.

440 With highly similar and identical protein structures predicted for Nsp4 and Nsp16, respectively,
441 we considered the possibility that the signals of positive selection instead reflect changes in RNA
442 structure and function. Previous studies found that neither the Nsp4 nor Nsp16 regions stand out
443 as particularly well folded regions of the genome, although Nsp16 does contain a single well-
444 folded region and Nsp4 two moderately well folded regions (Rangan et al., 2020; Huston et al.,
445 2020b). Further, both genes show significantly decreased sequence divergence among
446 coronavirus species within predicted double-stranded regions] (Rangan et al., 2020; Sanders et
447 al., 2020; Huston et al., 2020a). Indeed, the well-folded region within Nsp16 is the only such
448 region in the SARS-CoV-2 genome that is also well conserved among related coronaviruses
449 (Sanders et al., 2020). These published observations suggest possible functional roles for folded
450 structures within Nsp4 and Nsp16. Our minimum free energy (MFE) predictions reveal that the
451 likely secondary structure of the RNA genome in the region of the Nsp4 and Nsp16 genes likely
452 differs among the six coronavirus species we examined (Fig 4B and 5B, top rows). The MFE
453 predictions also indicate differences among species in entropy across the regions containing the
454 signals of positive selection, indicating likely differences in the stability of the folded molecule
455 (Fig 4B and 5B, bottom rows). Together, these new results indicate that the folded regions of
456 Nsp4 and Nsp16 in the SARS-Cov-2 genome may differ in shape from those of related
457 coronaviruses.

458 Unfortunately, little is currently known about the molecular functions of secondary structures in
459 coronavirus genomes. Most of the attention has been directed towards the 5' UTR, 3' UTR, and
460 frameshift element at the junction between Orf1a and Orf1b, which together contain the most
461 well-folded regions in the SARS-CoV-2 genome (Andrews et al., 2020; Sanders et al., 2020;
462 Huston et al., 2020a). Thus, it is not possible at this time to link structural and thermodynamic
463 features within Nsp4 and Nsp16 that are unique to SARS-CoV-2 to specific molecular functions.
464 As discussed above, however, published evidence suggests that RNA secondary structures within

465 these regions of the genome may be functional (Rangan et al., 2020; Sanders et al., 2020; Huston
466 et al., 2020a). These functions could, in principle, affect genome or transcript function, or both.
467 Plausible possibilities include secondary structures that recruit specific RNA-binding proteins to
468 mediate transcriptional regulation or transcript processing (Pirakitikulr et al., 2016; Pan et al.,
469 2020), that mediate looping for other reasons (Gebhard, Filomatori & Gamarnik, 2011; Ziv et al.,
470 2020), or that simply facilitate or impede processivity of the replication or translation machinery
471 (MacFadden et al., 2018).

472 **Conclusions**

473 Scans for positive selection typically focus on changes in protein function and far less often
474 consider the possibility of adaptive change in RNA function. By shining a light on regions of the
475 SARS-CoV-2 genome that appear to be under positive selection yet are unlikely to alter protein
476 function, our results illustrate the value of evaluating the potential for adaptive changes in
477 secondary structures within the genomes of RNA viruses. In particular, we identify Nsp4 and
478 Nsp16 as regions of the SARS-Cov-2 genome that may contain mutations that contribute to the
479 unique biological and epidemiological features of this recently emerged pathogen.
480 While it is tempting to speculate about the possible adaptive role of changes in RNA structure
481 within these accelerated regions, we suggest that this is best done in the context of relevant
482 experimental results. For example, it might be informative to modify the primary sequence of the
483 genome so as to encode the same protein sequence while altering or disrupting secondary
484 structure within Nsp4 or Nsp16, then assay the consequences for viral replication and for specific
485 molecular functions. We hope that our results inspire these or other experiments aimed at better
486 understand the evolving functions of RNA secondary structure within the SARS-CoV-2 genome.

488 **Acknowledgements**

489 The authors would like to thank all the members of the lab groups of David McClay and
490 Greg Wray for comments and the Compact for Open-Access Publishing Equity (COPE) program
491 at Duke University for supporting publication costs.

492 **References**

493 Alhatlani BY. 2020. In silico identification of conserved cis -acting RNA elements in the SARS-
494 CoV-2 genome. *Future Virology*:fv1-2020-0163. DOI: 10.2217/fv1-2020-0163.

- 495 Andersen KG, Rambaut A, Lipkin WI, Holmes EC, Garry RF. 2020. The proximal origin of
496 SARS-CoV-2. *Nature Medicine* 26:450–452. DOI: 10.1038/s41591-020-0820-9.
- 497 Andrews RJ, Peterson JM, Haniff HS, Chen J, Williams C, Grefe M, Disney MD, Moss WN.
498 2020. An in silico map of the SARS-CoV-2 RNA Structurome. *bioRxiv*: the preprint
499 server for biology:2020.04.17.045161. DOI: 10.1101/2020.04.17.045161.
- 500 Armijos Jaramillo V, Yeager J, Muslin C, Perez Castillo Y. 2020. SARS-CoV-2, an
501 evolutionary perspective of interaction with human ACE2 reveals undiscovered amino acids
502 necessary for complex stability. *Evolutionary Applications*:eva.12980. DOI:
503 10.1111/eva.12980.
- 504 Báez-Santos YM, St. John SE, Mesecar AD. 2015. The SARS-coronavirus papain-like protease:
505 Structure, function and inhibition by designed antiviral compounds. *Antiviral Research*
506 115:21–38. DOI: 10.1016/j.antiviral.2014.12.015.
- 507 Belouzard S, Chu VC, Whittaker GR. 2009. Activation of the SARS coronavirus spike protein
508 via sequential proteolytic cleavage at two distinct sites. *Proceedings of the National*
509 *Academy of Sciences of the United States of America* 106:5871–5876. DOI:
510 10.1073/pnas.0809524106.
- 511 Berrio A, Haygood R, Wray GA. 2020. Identifying branch-specific positive selection throughout
512 the regulatory genome using an appropriate proxy neutral. *BMC Genomics* 21:359. DOI:
513 10.1186/s12864-020-6752-4.
- 514 Boni MF, Lemey P, Jiang X, Lam TT-Y, Perry BW, Castoe TA, Rambaut A, Robertson DL.
515 2020. Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the
516 COVID-19 pandemic. *Nature Microbiology*:2020.03.30.015008. DOI: 10.1038/s41564-
517 020-0771-4.
- 518 Bouvet M, Debarnot C, Imbert I, Selisko B, Snijder EJ. 2010. In Vitro Reconstitution of SARS-
519 Coronavirus mRNA Cap Methylation. *PLoS Pathog* 6:1000863. DOI:
520 10.1371/journal.ppat.1000863.
- 521 Cagliani R, Forni D, Clerici M, Sironi M. 2020. Computational Inference of Selection
522 Underlying the Evolution of the Novel Coronavirus, Severe Acute Respiratory Syndrome
523 Coronavirus 2 Downloaded from. *JVI.asm.org* 1 *Journal of Virology* 94:411–431. DOI:
524 10.1128/JVI.00411-20.
- 525 Chaw S-M, Tai J-H, Chen S-L, Hsieh C-H, Chang S-Y, Yeh S-H, Yang W-S, Chen P-J, Wang

- 526 H-Y. 2020. The origin and underlying driving forces of the SARS-CoV-2 outbreak. *Journal*
527 *of biomedical science* 27:73. DOI: 10.1186/s12929-020-00665-8.
- 528 Cornillez-Ty CT, Liao L, Yates Iii JR, Kuhn P, Buchmeier MJ. 2009. Severe Acute Respiratory
529 Syndrome Coronavirus Nonstructural Protein 2 Interacts with a Host Protein Complex
530 Involved in Mitochondrial Biogenesis and Intracellular Signaling. *JOURNAL OF*
531 *VIROLOGY* 83:10314–10318. DOI: 10.1128/JVI.00842-09.
- 532 Decroly E, Imbert I, Coutard B, Bouvet M, Selisko B, Alvarez K, Gorbalenya AE, Snijder EJ,
533 Canard B. 2008. Coronavirus Nonstructural Protein 16 Is a Cap-0 Binding Enzyme
534 Possessing (Nucleoside-2'O)-Methyltransferase Activity. *Journal of Virology* 82:8071–
535 8084. DOI: 10.1128/jvi.00407-08.
- 536 DeLano WL. 2002. Pymol: An open-source molecular graphics tool. :40, 82–92.
- 537 Eaaswarkhanth M, Al Madhoun A, Al-Mulla F. 2020. Could the D614G substitution in the
538 SARS-CoV-2 spike (S) protein be associated with higher COVID-19 mortality?
539 *International journal of infectious diseases* □: *IJID* □: *official publication of the*
540 *International Society for Infectious Diseases* 96:459–460. DOI: 10.1016/j.ijid.2020.05.071.
- 541 Fehr AR, Perlman S. 2015. Coronaviruses: An overview of their replication and pathogenesis.
542 In: *Coronaviruses: Methods and Protocols*. Springer New York, 1–23. DOI: 10.1007/978-
543 1-4939-2438-7_1.
- 544 Follis KE, York J, Nunberg JH. 2006. Furin cleavage of the SARS coronavirus spike
545 glycoprotein enhances cell-cell fusion but does not affect virion entry. *Virology* 350:358–
546 369. DOI: 10.1016/j.virol.2006.02.003.
- 547 Forni D, Filippi G, Cagliani R, De Gioia L, Pozzoli U, Al-Daghri N, Clerici M, Sironi M. 2015.
548 The heptad repeat region is a major selection target in MERS-CoV and related
549 coronaviruses. *Scientific Reports* 5:1–10. DOI: 10.1038/srep14480.
- 550 Freundt EC, Yu L, Goldsmith CS, Welsh S, Cheng A, Yount B, Liu W, Frieman MB, Buchholz
551 UJ, Sreaton GR, Lippincott-Schwartz J, Zaki SR, Xu X-N, Baric RS, Subbarao K, Lenardo
552 MJ. 2010. The Open Reading Frame 3a Protein of Severe Acute Respiratory Syndrome-
553 Associated Coronavirus Promotes Membrane Rearrangement and Cell Death. *Journal of*
554 *Virology* 84:1097–1109. DOI: 10.1128/jvi.01662-09.
- 555 Frieman M, Ratia K, Johnston RE, Mesecar AD, Baric RS. 2009. Severe acute respiratory
556 syndrome coronavirus papain-like protease ubiquitin-like domain and catalytic domain

- 557 regulate antagonism of IRF3 and NF-kappaB signaling. *Journal of virology* 83:6689–705.
558 DOI: 10.1128/JVI.02220-08.
- 559 Fuchs SY. 2012. Ubiquitination-mediated regulation of interferon responses. *Growth Factors*
560 30:141–148. DOI: 10.3109/08977194.2012.669382.
- 561 Gallagher TM, Buchmeier MJ. 2001. Coronavirus Spike Proteins in Viral Entry and
562 Pathogenesis. *Virology* 279:371–374. DOI: 10.1006/viro.2000.0757.
- 563 Gebhard LG, Filomatori C V., Gamarnik A V. 2011. Functional RNA elements in the dengue
564 virus genome. *Viruses* 3:1739–1756. DOI: 10.3390/v3091739.
- 565 Graham RL, Baric RS. 2010. Recombination, reservoirs, and the modular spike: mechanisms of
566 coronavirus cross-species transmission. *Journal of virology* 84:3134–46. DOI:
567 10.1128/JVI.01394-09.
- 568 Gruber AR, Lorenz R, Bernhart SH, Neuböck R, Hofacker IL. 2008. The Vienna RNA websuite.
569 *Nucleic acids research* 36:70–74. DOI: 10.1093/nar/gkn188.
- 570 Hagemeyer MC, Monastyrska I, Griffith J, van der Sluijs P, Voortman J, van Bergen en
571 Henegouwen PM, Vonk AM, Rottier PJM, Reggiori F, de Haan CAM. 2014. Membrane
572 rearrangements mediated by coronavirus nonstructural proteins 3 and 4. *Virology* 458–
573 459:125–135. DOI: 10.1016/j.virol.2014.04.027.
- 574 Hagemeyer MC, Ulasli M, Vonk AM, Reggiori F, Rottier PJM, de Haan CAM. 2011. Mobility
575 and Interactions of Coronavirus Nonstructural Protein 4. *Journal of Virology* 85:4572–4577.
576 DOI: 10.1128/JVI.00042-11.
- 577 Hahne F, Ivanek R. 2016. Visualizing genomic data using Gviz and bioconductor. In: *Methods in*
578 *Molecular Biology*. Humana Press Inc., 335–351. DOI: 10.1007/978-1-4939-3578-9_16.
- 579 Haygood R, Fedrigo O, Hanson B, Yokoyama K-D, Wray GA. 2007. Promoter regions of many
580 neural- and nutrition-related genes have experienced positive selection during human
581 evolution. *Nature genetics* 39:1140–4. DOI: 10.1038/ng2104.
- 582 He X, Lau EHY, Wu P, Deng X, Wang J, Hao X, Lau YC, Wong JY, Guan Y, Tan X, Mo X,
583 Chen Y, Liao B, Chen W, Hu F, Zhang Q, Zhong M, Wu Y, Zhao L, Zhang F, Cowling BJ,
584 Li F, Leung GM. 2020. Temporal dynamics in viral shedding and transmissibility of
585 COVID-19. *Nature Medicine* 26:672–675. DOI: 10.1038/s41591-020-0869-5.
- 586 Hillen HS, Kokic G, Farnung L, Dienemann C, Tegunov D, Cramer P. 2020. Structure of
587 replicating SARS-CoV-2 polymerase. *Nature*. DOI: 10.1038/s41586-020-2368-8.

- 588 Hoffmann M, Kleine-Weber H, Pöhlmann S. 2020. A Multibasic Cleavage Site in the Spike
589 Protein of SARS-CoV-2 Is Essential for Infection of Human Lung Cells. *Molecular Cell*
590 78:779-784.e5. DOI: 10.1016/j.molcel.2020.04.022.
- 591 Hon C-C, Lam T-Y, Shi Z-L, Drummond AJ, Yip C-W, Zeng F, Lam P-Y, Chi F, Leung C.
592 2008. Evidence of the Recombinant Origin of a Bat Severe Acute Respiratory Syndrome
593 (SARS)-Like Coronavirus and Its Implications on the Direct Ancestor of SARS
594 Coronavirus. *JOURNAL OF VIROLOGY* 82:1819–1826. DOI: 10.1128/JVI.01926-07.
- 595 Hu B, Zeng L-P, Yang X-L, Ge X-Y, Zhang W, Li B, Xie J-Z, Shen X-R, Zhang Y-Z, Wang N,
596 Luo D-S, Zheng X-S, Wang M-N, Daszak P, Wang L-F, Cui J, Shi Z-L. 2017. Discovery of
597 a rich gene pool of bat SARS-related coronaviruses provides new insights into the origin of
598 SARS coronavirus. *PLOS Pathogens* 13:e1006698. DOI: 10.1371/journal.ppat.1006698.
- 599 Hubisz MJ, Pollard KS, Siepel A. 2011. PHAST and RPHAST: phylogenetic analysis with
600 space/time models. *Briefings in bioinformatics* 12:41–51. DOI: 10.1093/bib/bbq072.
- 601 Hulswit RJG, de Haan CAM, Bosch BJ. 2016. Coronavirus Spike Protein and Tropism Changes.
602 In: *Advances in Virus Research*. Academic Press Inc., 29–57. DOI:
603 10.1016/bs.aivir.2016.08.004.
- 604 Huston NC, Wan H, Araujo Tavares R de C, Wilen C, Pyle AM. 2020a. Comprehensive in-vivo
605 secondary structure of the SARS-CoV-2 genome reveals novel regulatory motifs and
606 mechanisms. *bioRxiv*: the preprint server for biology:1–46. DOI:
607 10.1101/2020.07.10.197079.
- 608 Huston NC, Wan H, Araujo Tavares R de C, Wilen C, Pyle AM. 2020b. Comprehensive in-vivo
609 secondary structure of the SARS-CoV-2 genome reveals novel regulatory motifs and
610 mechanisms. *bioRxiv*: the preprint server for biology:1–46. DOI:
611 10.1101/2020.07.10.197079.
- 612 Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version 7:
613 Improvements in Performance and Usability. *Molecular Biology and Evolution* 30:772–780.
614 DOI: 10.1093/molbev/mst010.
- 615 Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A,
616 Markowitz S, Duran C, Thierer T, Ashton B, Meintjes P, Drummond A. 2012. Geneious
617 Basic: an integrated and extendable desktop software platform for the organization and
618 analysis of sequence data. *Bioinformatics (Oxford, England)* 28:1647–9. DOI:

- 619 10.1093/bioinformatics/bts199.
- 620 Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE. 2016. The Phyre2 web portal for
621 protein modeling, prediction and analysis. *Nature Protocols* 10:845–858. DOI:
622 10.1038/nprot.2015.053.
- 623 Kim D, Lee J-Y, Yang J-S, Kim JW, Kim VN, Chang H. 2020. The Architecture of SARS-CoV-
624 2 Transcriptome. *Cell* 181:914-921.e10. DOI: 10.1016/j.cell.2020.04.011.
- 625 Korber B, Fischer WM, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Hengartner N, Giorgi
626 EE, Bhattacharya T, Foley B, Hastie KM, Parker MD, Partridge DG, Evans CM, Freeman
627 TM, de Silva TI, Sheffield COVID-19 Genomics Group, McDanal C, Perez LG, Tang H,
628 Moon-Walker A, Whelan SP, LaBranche CC, Sapphire EO, Montefiori DC. 2020. Tracking
629 Changes in SARS-CoV-2 Spike: Evidence that D614G Increases Infectivity of the COVID-
630 19 Virus. *Cell*. DOI: 10.1016/j.cell.2020.06.043.
- 631 Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A, Wren J. 2019. RAxML-NG: A fast,
632 scalable and user-friendly tool for maximum likelihood phylogenetic inference.
633 *Bioinformatics* 35:4453–4455. DOI: 10.1093/bioinformatics/btz305.
- 634 Lam TT-Y, Jia N, Zhang Y-W, Shum MH-H, Jiang J-F, Zhu H-C, Tong Y-G, Shi Y-X, Ni X-B,
635 Liao Y-S, Li W-J, Jiang B-G, Wei W, Yuan T-T, Zheng K, Cui X-M, Li J, Pei G-Q, Qiang
636 X, Cheung WY-M, Li L-F, Sun F-F, Qin S, Huang J-C, Leung GM, Holmes EC, Hu Y-L,
637 Guan Y, Cao W-C. 2020. Identifying SARS-CoV-2-related coronaviruses in Malayan
638 pangolins. *Nature* 583:282–285. DOI: 10.1038/s41586-020-2169-0.
- 639 Lau SKP, Feng Y, Chen H, Luk HKH, Yang W-H, Li KSM, Zhang Y-Z, Huang Y, Song Z-Z,
640 Chow W-N, Fan RYY, Ahmed SS, Yeung HC, Lam CSF, Cai J-P, Wong SSY, Chan JFW,
641 Yuen K-Y, Zhang H-L, Woo PCY. 2015. Severe Acute Respiratory Syndrome (SARS)
642 Coronavirus ORF8 Protein Is Acquired from SARS-Related Coronavirus from Greater
643 Horseshoe Bats through Recombination. *Journal of Virology* 89:10532–10547. DOI:
644 10.1128/jvi.01048-15.
- 645 Li X, Giorgi EE, Marichannelgowda MH, Foley B, Xiao C, Kong X-P, Chen Y, Gnanakaran S,
646 Korber B, Gao F. 2020a. Emergence of SARS-CoV-2 through recombination and strong
647 purifying selection. *Science Advances* 6:eabb9153. DOI: 10.1126/sciadv.abb9153.
- 648 Li Y, Yang X, Wang N, Wang H, Yin B, Yang X, Jiang W. 2020b. The divergence between
649 SARS-CoV-2 and RaTG13 might be overestimated due to the extensive RNA modification.

- 650 *Future Virology*:fv1-2020-0066. DOI: 10.2217/fv1-2020-0066.
- 651 Liu P, Chen W, Chen J-P. 2019. Viral Metagenomics Revealed Sendai Virus and Coronavirus
652 Infection of Malayan Pangolins (*Manis javanica*). *Viruses* 11:979. DOI:
653 10.3390/v11110979.
- 654 Lorenz R, Bernhart SH, Höner zu Siederdisen C, Tafer H, Flamm C, Stadler PF, Hofacker IL.
655 2011. ViennaRNA Package 2.0. *Algorithms for Molecular Biology* 6:26. DOI:
656 10.1186/1748-7188-6-26.
- 657 Lu W, Zheng BJ, Xu K, Schwarz W, Du L, Wong CKL, Chen J, Duan S, Deubel V, Sun B.
658 2006. Severe acute respiratory syndrome-associated coronavirus 3a protein forms an ion
659 channel and modulates virus release. *Proceedings of the National Academy of Sciences of*
660 *the United States of America* 103:12540–12545. DOI: 10.1073/pnas.0605402103.
- 661 M de Haan CA, Stadler K, Godeke G-J, Jan Bosch B, M Rottier PJ. 2004. Cleavage Inhibition of
662 the Murine Coronavirus Spike Protein by a Furin-Like Enzyme Affects Cell-Cell but Not
663 Virus-Cell Fusion Downloaded from. *JOURNAL OF VIROLOGY* 78:6048–6054. DOI:
664 10.1128/JVI.78.11.6048-6054.2004.
- 665 MacFadden A, Odonoghue Z, Silva PAGC, Chapman EG, Olsthoorn RC, Sterken MG, Pijlman
666 GP, Bredenbeek PJ, Kieft JS. 2018. Mechanism and structural diversity of exoribonuclease-
667 resistant RNA structures in flaviviral RNAs. *Nature Communications* 9:1–11. DOI:
668 10.1038/s41467-017-02604-y.
- 669 Madhugiri R, Fricke M, Marz M, Ziebuhr J. 2016. Coronavirus cis-Acting RNA Elements. In:
670 *Advances in Virus Research*. Academic Press Inc., 127–163. DOI:
671 10.1016/bs.aivir.2016.08.007.
- 672 Matthews K, Schäfer A, Pham A, Frieman M. 2014. The SARS coronavirus papain like protease
673 can inhibit IRF3 at a post activation step that requires deubiquitination activity. *Virology*
674 *Journal* 11:209. DOI: 10.1186/s12985-014-0209-9.
- 675 Menachery VD, Debbink K, Baric RS. 2014. Coronavirus non-structural protein 16: Evasion,
676 attenuation, and possible treatments. *Virus Research* 194:191–199. DOI:
677 10.1016/j.virusres.2014.09.009.
- 678 Minakshi R, Padhan K, Rani M, Khan N, Ahmad F, Jameel S. 2009. The SARS Coronavirus 3a
679 protein causes endoplasmic reticulum stress and induces ligand-independent
680 downregulation of the type 1 interferon receptor. *PloS one* 4:e8342. DOI:

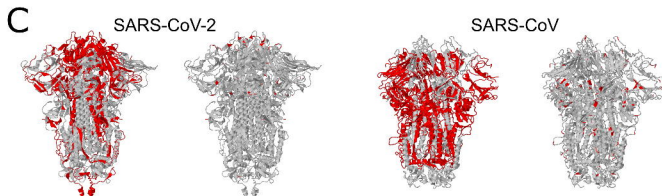
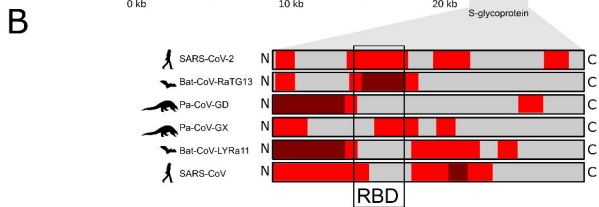
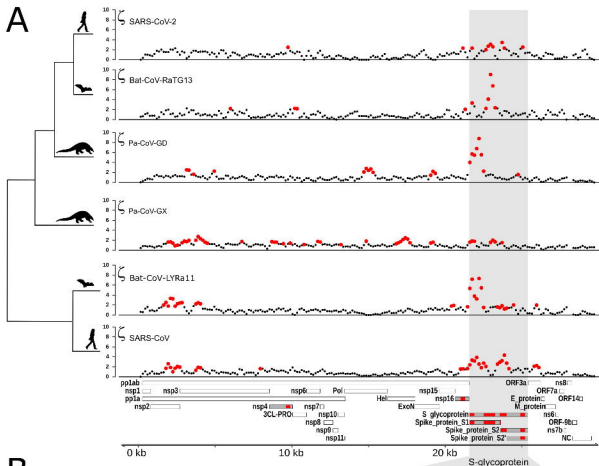
- 681 10.1371/journal.pone.0008342.
- 682 Nelemans T, Kikkert M. 2019. Viral innate immune evasion and the pathogenesis of emerging
683 RNA virus infections. *Viruses* 11. DOI: 10.3390/v11100961.
- 684 Nielsen R. 1997. The ratio of replacement to silent divergence and tests of neutrality. *Journal of*
685 *Evolutionary Biology* 10:217–231. DOI: 10.1046/j.1420-9101.1997.10020217.x.
- 686 Oostra M, Hagemeyer MC, van Gent M, Bekker CPJ, te Lintelo EG, Rottier PJM, de Haan
687 CAM. 2008. Topology and Membrane Anchoring of the Coronavirus Replication Complex:
688 Not All Hydrophobic Domains of nsp3 and nsp6 Are Membrane Spanning. *Journal of*
689 *Virology* 82:12392–12405. DOI: 10.1128/jvi.01219-08.
- 690 Pachetti M, Marini B, Benedetti F, Giudici F, Mauro E, Storici P, Masciovecchio C, Angeletti S,
691 Ciccozzi M, Gallo RC, Zella D, Ippodrino R. 2020. Emerging SARS-CoV-2 mutation hot
692 spots include a novel RNA-dependent-RNA polymerase variant. *Journal of Translational*
693 *Medicine* 18:1–9. DOI: 10.1186/s12967-020-02344-6.
- 694 Pan J, Qian X, Lattmann S, El Sahili A, Yeo TH, Jia H, Cressey T, Ludeke B, Noton S, Kalocsay
695 M, Fearn R, Lescar J. 2020. Structure of the human metapneumovirus polymerase
696 phosphoprotein complex. *Nature* 577:275–279. DOI: 10.1038/s41586-019-1759-1.
- 697 Peiris JSM, Chu CM, Cheng VCC, Chan KS, Hung IFN, Poon LLM, Law KI, Tang BSF, Hon
698 TYW, Chan CS, Chan KH, Ng JSC, Zheng BJ, Ng WL, Lai RWM, Guan Y, Yuen KY.
699 2003. Clinical progression and viral load in a community outbreak of coronavirus-
700 associated SARS pneumonia: A prospective study. *Lancet* 361:1767–1772. DOI:
701 10.1016/S0140-6736(03)13412-5.
- 702 Pirakitikulr N, Kohlway A, Lindenbach BD, Pyle AM. 2016. The Coding Region of the HCV
703 Genome Contains a Network of Regulatory RNA Structures. *Molecular cell* 62:111–20.
704 DOI: 10.1016/j.molcel.2016.01.024.
- 705 Pitzer VE, Leung GM, Lipsitch M. 2007. Estimating variability in the transmission of severe
706 acute respiratory syndrome to household contacts in Hong Kong, China. *American Journal*
707 *of Epidemiology* 166:355–363. DOI: 10.1093/aje/kwm082.
- 708 Pond SLK, Frost SDW, Muse S V. 2005. HyPhy: hypothesis testing using phylogenies.
709 *Bioinformatics (Oxford, England)* 21:676–9. DOI: 10.1093/bioinformatics/bti079.
- 710 Pond SLK, Poon AFY, Velazquez R, Weaver S, Hepler NL, Murrell B, Shank SD, Magalis BR,
711 Bouvier D, Nekrutenko A, Wisotsky S, Spielman SJ, Frost SDW, Muse S V. 2020. HyPhy

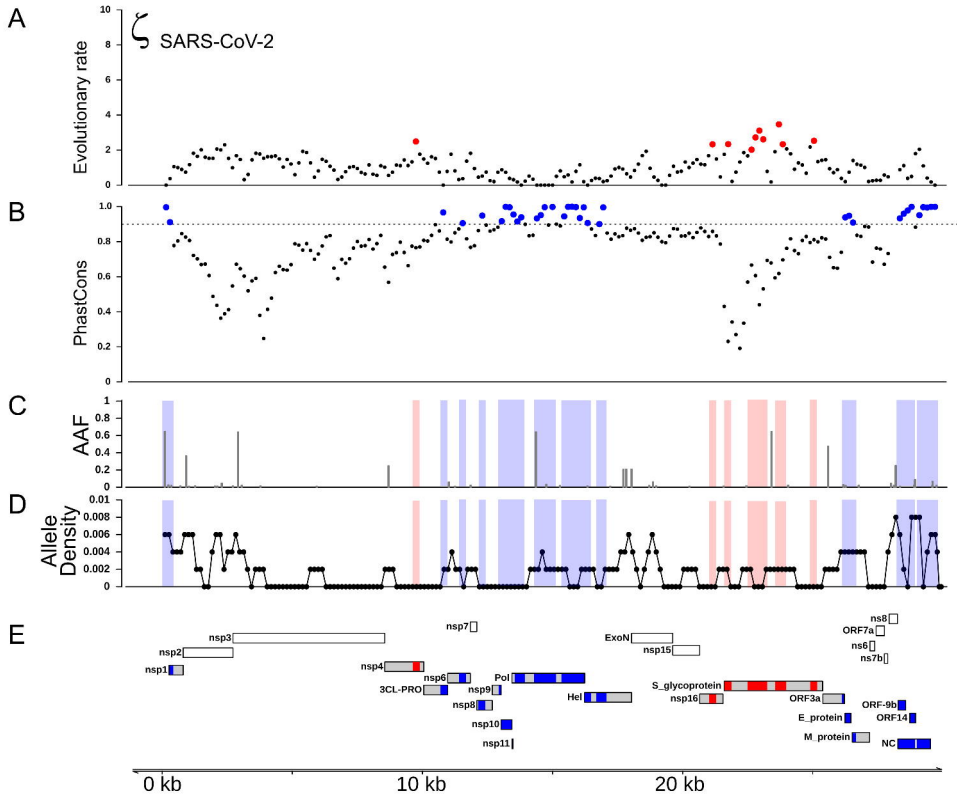
- 712 2.5-A Customizable Platform for Evolutionary Hypothesis Testing Using Phylogenies.
713 *Molecular biology and evolution* 37:295–299. DOI: 10.1093/molbev/msz197.
- 714 Rangan R, Zheludev IN, Hagey RJ, Pham EA, Wayment-Steele HK, Glenn JS, Das R. 2020.
715 RNA genome conservation and secondary structure in SARS-CoV-2 and SARS-related
716 viruses: a first look. *RNA (New York, N.Y.)* 26:937–959. DOI: 10.1261/rna.076141.120.
- 717 Sanders W, Fritch EJ, Madden EA, Graham RL, Vincent HA, Heise MT, Baric RS, Moorman
718 NJ. 2020. Comparative analysis of coronavirus genomic RNA structure reveals
719 conservation in SARS-like coronaviruses. *bioRxiv*: the preprint server for biology. DOI:
720 10.1101/2020.06.15.153197.
- 721 Sharp PM. 1997. In search of molecular darwinism. *Nature* 385:111–112. DOI:
722 10.1038/385111a0.
- 723 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J,
724 Hillier LDW, Richards S, Weinstock GM, Wilson RK, Gibbs RA, Kent WJ, Miller W,
725 Haussler D. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
726 genomes. *Genome Research* 15:1034–1050. DOI: 10.1101/gr.3715005.
- 727 Snijder EJ, Decroly E, Ziebuhr J. 2016. The Nonstructural Proteins Directing Coronavirus RNA
728 Synthesis and Processing. In: *Advances in Virus Research*. Academic Press Inc., 59–126.
729 DOI: 10.1016/bs.aivir.2016.08.008.
- 730 Tan Y-J, Tham P-Y, Chan DZL, Chou C-F, Shen S, Fielding BC, Tan THP, Lim SG, Hong W.
731 2005. The Severe Acute Respiratory Syndrome Coronavirus 3a Protein Up-Regulates
732 Expression of Fibrinogen in Lung Epithelial Cells. *Journal of Virology* 79:10083–10087.
733 DOI: 10.1128/jvi.79.15.10083-10087.2005.
- 734 Tang X, Wu C, Li X, Song Y, Yao X, Wu X, Duan Y, Zhang H, Wang Y, Qian Z, Cui J, Lu J.
735 2020. On the origin and continuing evolution of SARS-CoV-2. *National Science Review*
736 7:1012–1023. DOI: 10.1093/nsr/nwaa036.
- 737 Tortorici MA, Veessler D. 2019. Structural insights into coronavirus entry. *Advances in virus*
738 *research* 105:93–116. DOI: 10.1016/bs.aivir.2019.08.002.
- 739 Toyoshima Y, Nemoto K, Matsumoto S, Nakamura Y, Kiyotani K. 2020. SARS-CoV-2 genomic
740 variations associated with mortality rate of COVID-19. *Journal of Human Genetics*:1–8.
741 DOI: 10.1038/s10038-020-0808-9.
- 742 Tsuchida T, Kawai T, Akira S. 2009. Inhibition of IRF3-dependent antiviral responses by

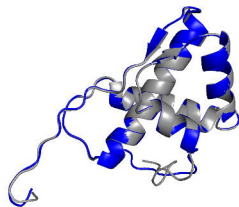
- 743 cellular and viral proteins. *Cell Research* 19:3–4. DOI: 10.1038/cr.2009.1.
- 744 Velazquez-Salinas L, Zarate S, Eberl S, Gladue DP, Novella I, Borca M V. 2020. Positive
745 selection of ORF3a and ORF8 genes drives the evolution of SARS-CoV-2 during the 2020
746 COVID-19 pandemic. *bioRxiv*:2020.04.10.035964. DOI: 10.1101/2020.04.10.035964.
- 747 Wang Y, Liu M, Gao J. 2020. Enhanced receptor binding of SARS-CoV-2 through networks of
748 hydrogen-bonding and hydrophobic interactions. *Proceedings of the National Academy of
749 Sciences of the United States of America* 117:13967–13974. DOI:
750 10.1073/pnas.2008209117.
- 751 Wang Q, Zhang Y, Wu L, Niu S, Song C, Zhang Z, Lu G, Qiao C, Hu Y, Yuen KY, Wang Q,
752 Zhou H, Yan J, Qi J. 2020. Structural and Functional Basis of SARS-CoV-2 Entry by Using
753 Human ACE2. *Cell* 181:894-904.e9. DOI: 10.1016/j.cell.2020.03.045.
- 754 Watanabe Y, Bowden TA, Wilson IA, Crispin M. 2019. Exploitation of glycosylation in
755 enveloped virus pathobiology. *Biochimica et biophysica acta. General subjects* 1863:1480–
756 1497. DOI: 10.1016/j.bbagen.2019.05.012.
- 757 Wong WSW, Nielsen R. 2004. Detecting selection in noncoding regions of nucleotide
758 sequences. *Genetics* 167:949–58. DOI: 10.1534/genetics.102.010959.
- 759 Wrapp D, Wang N, Corbett KS, Goldsmith JA, Hsieh C-L, Abiona O, Graham BS, McLellan JS.
760 2020. Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science*
761 367:1260–1263. DOI: 10.1126/science.abb2507.
- 762 Yang T-J, Chang Y-C, Ko T-P, Draczkowski P, Chien Y-C, Chang Y-C, Wu K-P, Khoo K-H,
763 Chang H-W, Hsu S-TD. 2020. Cryo-EM analysis of a feline coronavirus spike protein
764 reveals a unique structure and camouflaging glycans. *Proceedings of the National Academy
765 of Sciences* 117:1438–1446. DOI: 10.1073/pnas.1908898117.
- 766 Yang Z, Nielsen R. 2008. Mutation-selection models of codon substitution and their use to
767 estimate selective strengths on codon usage. *Molecular Biology and Evolution* 25:568–579.
768 DOI: 10.1093/molbev/msm284.
- 769 Ye Y, Godzik A. 2004. FATCAT: a web server for flexible structure comparison and structure
770 similarity searching. *Nucleic Acids Research* 32:W582–W585. DOI: 10.1093/nar/gkh430.
- 771 Yue Y, Nabar NR, Shi CS, Kamenyeva O, Xiao X, Hwang IY, Wang M, Kehrl JH. 2018. SARS-
772 Coronavirus Open Reading Frame-3a drives multimodal necrotic cell death. *Cell Death and
773 Disease* 9:1–15. DOI: 10.1038/s41419-018-0917-y.

774 Ziv O, Price J, Shalamova L, Kamenova T, Goodfellow I, Weber F, Miska EA. 2020. The short-
775 and long-range RNA-RNA Interactome of SARS-CoV-2 Co-first authors.
776 *bioRxiv*:2020.07.19.211110. DOI: 10.1101/2020.07.19.211110.

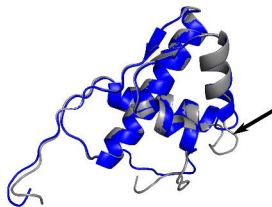
777 Züst R, Cervantes-Barragan L, Habjan M, Maier R, Neuman BW, Ziebuhr J, Szretter KJ, Baker
778 SC, Barchet W, Diamond MS, Siddell SG, Ludewig B, Thiel V. 2011. Ribose 2'-O-
779 methylation provides a molecular signature for the distinction of self and non-self mRNA
780 dependent on the RNA sensor Mda5. *Nature immunology* 12:137–43. DOI:
781 10.1038/ni.1979.
782



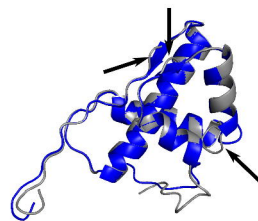


A**SARS-CoV-2**

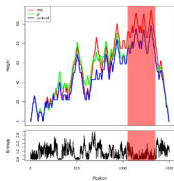
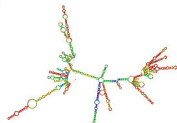
Bat-CoV-RaTG13



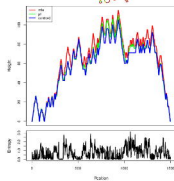
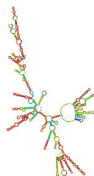
Pa-CoV-GD



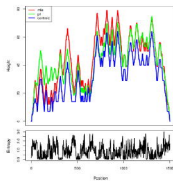
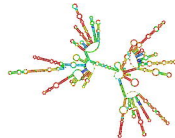
SARS-CoV

B

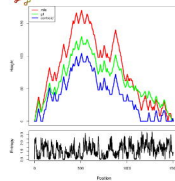
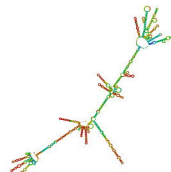
SARS-CoV-2



Bat-CoV-RaTG13



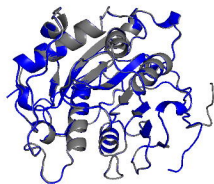
Pa-CoV-GD



SARS-CoV

A

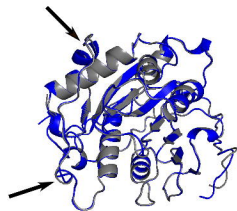
SARS-CoV-2



Bat-CoV-RaTG13

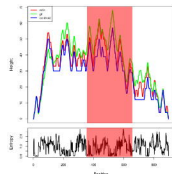
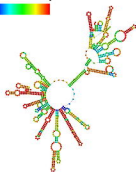


Pa-CoV-GD

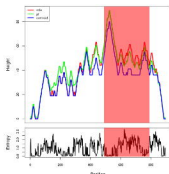
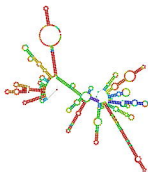


SARS-CoV

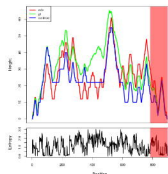
B



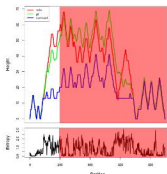
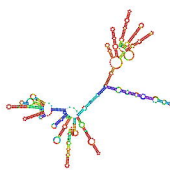
SARS-CoV-2



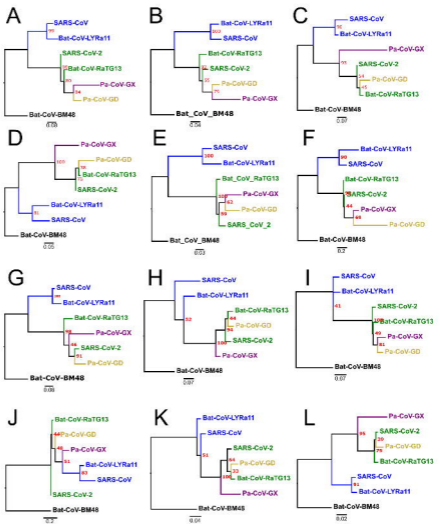
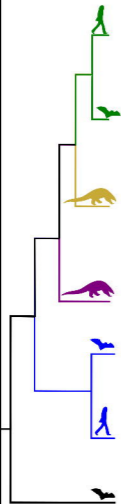
Bat-CoV-RaTG13



Pa-CoV-GD



SARS-CoV



M

