

1 **Molecular Architecture of Early Dissemination**
2 **and Massive Second Wave of the SARS-CoV-2 Virus in a**
3 **Major Metropolitan Area**

4

5 **Running Title: Two waves of COVID-19 disease in Houston, Texas**

6

7 **S. Wesley Long,^{a,b,1} Randall J. Olsen,^{a,b,1} Paul A. Christensen,^{a,1} David W.**
8 **Bernard,^{a,b} James J. Davis,^{c,d} Maulik Shukla,^{c,d} Marcus Nguyen,^{c,d} Matthew**
9 **Ojeda Saavedra,^a Prasanti Yerramilli,^a Layne Pruitt,^a Sishir Subedi,^a Hung-**
10 **Che Kuo,^e Heather Hendrickson,^a Ghazaleh Eskandari,^a Hoang A. T.**
11 **Nguyen,^a J. Hunter Long,^a Muthiah Kumaraswami,^a Jule Goike,^e Daniel**
12 **Boutz,^f Jimmy Gollihar,^{a,f} Jason S. McLellan,^e Chia-Wei Chou,^e Kamyab**
13 **Javanmardi,^e Ilya J. Finkelstein,^{e,g}, and James M. Musser^{a,b#}**

14

15 ^aCenter for Molecular and Translational Human Infectious Diseases Research,
16 Department of Pathology and Genomic Medicine, Houston Methodist Research
17 Institute and Houston Methodist Hospital, 6565 Fannin Street, Houston, Texas
18 77030

19 ^bDepartments of Pathology and Laboratory Medicine, and Microbiology and
20 Immunology, Weill Cornell Medical College, 1300 York Avenue, New York, New
21 York 10065

22 ^cConsortium for Advanced Science and Engineering, University of Chicago, 5801

23 South Ellis Avenue, Chicago, Illinois, 60637

24 ^dComputing, Environment and Life Sciences, Argonne National Laboratory, 9700

25 South Cass Avenue, Lemont, Illinois 60439

26 ^eDepartment of Molecular Biosciences and Institute of Molecular Biosciences,

27 The University of Texas at Austin, Austin, Texas 78712

28 ^fCCDC Army Research Laboratory-South, University of Texas, Austin, Texas 78712

29 ^gCenter for Systems and Synthetic Biology, University of Texas at Austin, Austin,
30 Texas 78712

31

32 ¹S.W.L., R.J.O., and P.A.C. contributed equally to this article. The order of co-first
33 authors was determined by discussion and mutual agreement between the three
34 co-first authors.

35

36 #Address correspondence to: James M. Musser, M.D., Ph.D., Department of
37 Pathology and Genomic Medicine, Houston Methodist Research Institute, 6565
38 Fannin Street, Suite B490, Houston, Texas 77030. Tel: 713.441.5890, E-mail:
39 jmmusser@houstonmethodist.org.

40

41 This article is a direct contribution from James M. Musser, a Fellow of the
42 American Academy of Microbiology, who arranged for and secured reviews by
43 Barry N. Kreiswirth, Center for Discovery and Innovation, Hackensack Meridian

44 Health, New Jersey; and David M. Morens, National Institute of Allergy and
45 Infectious Diseases, National Institutes of Health, Maryland.

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66 **ABSTRACT** We sequenced the genomes of 5,085 SARS-CoV-2 strains
67 causing two COVID-19 disease waves in metropolitan Houston, Texas, an
68 ethnically diverse region with seven million residents. The genomes were
69 from viruses recovered in the earliest recognized phase of the pandemic in
70 Houston, and an ongoing massive second wave of infections. The virus
71 was originally introduced into Houston many times independently. Virtually
72 all strains in the second wave have a Gly614 amino acid replacement in the
73 spike protein, a polymorphism that has been linked to increased
74 transmission and infectivity. Patients infected with the Gly614 variant
75 strains had significantly higher virus loads in the nasopharynx on initial
76 diagnosis. We found little evidence of a significant relationship between
77 virus genotypes and altered virulence, stressing the linkage between
78 disease severity, underlying medical conditions, and host genetics. Some
79 regions of the spike protein - the primary target of global vaccine efforts -
80 are replete with amino acid replacements, perhaps indicating the action of
81 selection. We exploited the genomic data to generate defined single amino
82 acid replacements in the receptor binding domain of spike protein that,
83 importantly, produced decreased recognition by the neutralizing
84 monoclonal antibody CR30022. Our study is the first analysis of the
85 molecular architecture of SARS-CoV-2 in two infection waves in a major
86 metropolitan region. The findings will help us to understand the origin,
87 composition, and trajectory of future infection waves, and the potential

88 **effect of the host immune response and therapeutic maneuvers on SARS-**
89 **CoV-2 evolution.**

90

91 **IMPORTANCE** There is concern about second and subsequent waves of
92 COVID-19 caused by the SARS-CoV-2 coronavirus occurring in
93 communities globally that had an initial disease wave. Metropolitan
94 Houston, Texas, with a population of 7 million, is experiencing a massive
95 second disease wave that began in late May 2020. To understand SARS-
96 CoV-2 molecular population genomic architecture, evolution, and
97 relationship between virus genotypes and patient features, we sequenced
98 the genomes of 5,085 SARS-CoV-2 strains from these two waves. Our study
99 provides the first molecular characterization of SARS-CoV-2 strains
100 causing two distinct COVID-19 disease waves.

101

102 **KEYWORDS:** SARS-CoV-2, COVID-19 disease, genome sequencing, molecular
103 population genomics, evolution

104

105 [Introduction]

106 **P**andemic disease caused by the severe acute respiratory syndrome
107 coronavirus 2 (SARS-CoV-2) virus is now responsible for massive human
108 morbidity and mortality worldwide (1-5). The virus was first documented to cause
109 severe respiratory infections in Wuhan, China, beginning in late December 2019
110 (6-9). Global dissemination occurred extremely rapidly and has affected major
111 population centers on most continents (10, 11). In the United States, the Seattle
112 and the New York City (NYC) regions have been especially important centers of
113 COVID-19 disease caused by SARS-CoV-2. For example, as of August 19,
114 2020, there were 227,419 confirmed SARS-CoV-2 cases in NYC, causing 56,831
115 hospitalizations and 19,005 confirmed fatalities and 4,638 probable fatalities (12).
116 Similarly, in Seattle and King County, 17,989 positive patients and 696 deaths
117 have been reported as of August 18, 2020 (13).

118 The Houston metropolitan area is the fourth largest and most ethnically
119 diverse city in the United States, with a population of approximately 7 million
120 (14, 15). The 2,400-bed Houston Methodist health system has seven hospitals
121 and serves a large, multiethnic, and socioeconomically diverse patient population
122 throughout greater Houston (13, 14). The first COVID-19 case in metropolitan
123 Houston was reported on March 5, 2020 with community spread occurring one
124 week later (16). Many of the first cases in our region were associated with
125 national or international travel in areas known to have SARS-CoV-2 virus
126 outbreaks (16). A central molecular diagnostic laboratory serving all Houston
127 Methodist hospitals and our very early adoption of a molecular test for the SARS-

128 CoV-2 virus permitted us to rapidly identify positive patients and interrogate
129 genomic variation among strains causing early infections in the greater Houston
130 area. Our analysis of SARS-CoV-2 genomes causing disease in Houston has
131 continued unabated since early March and is ongoing. Genome sequencing and
132 related efforts were expanded extensively in late May as we recognized that a
133 prominent second wave was underway (**Figure 1**).

134 Here, we report that SARS-CoV-2 was introduced to the Houston area
135 many times, independently, from diverse geographic regions, with virus
136 genotypes representing genetic clades causing disease in Europe, Asia, South
137 America and elsewhere in the United States. There was widespread community
138 dissemination soon after COVID-19 cases were reported in Houston. Strains with
139 a Gly614 amino acid replacement in the spike protein, a polymorphism that has
140 been linked to increased transmission and *in vitro* cell infectivity, increased
141 significantly over time and caused virtually all COVID-19 cases in the massive
142 second disease wave. Patients infected with strains with the Gly614 variant had
143 significantly higher virus loads in the nasopharynx on initial diagnosis. Some
144 naturally occurring single amino acid replacements in the receptor binding
145 domain of spike protein resulted in decreased reactivity with a neutralizing
146 monoclonal antibody, consistent with the idea that some virus variants arise due
147 to host immune pressure.

148

149 **RESULTS**

150 **Description of metropolitan Houston.** Houston, Texas, is located in the
151 southwestern United States, 50 miles inland from the Gulf of Mexico. It is the
152 most ethnically diverse city in the United States (14). Metropolitan Houston is
153 comprised predominantly of Harris County plus parts of eight contiguous
154 surrounding counties. In the aggregate, the metropolitan area includes 9,444
155 square miles. The estimated population size of metropolitan Houston is 7 million
156 (<https://www.houston.org/houston-data>).

157 **Epidemic curve characteristics over two disease waves.** The first
158 confirmed case of COVID-19 in the Houston metropolitan region was reported on
159 March 5, 2020 (16), and the first confirmed case diagnosed in Houston Methodist
160 hospitals was reported on March 6, 2020. The epidemic curve indicated a first
161 wave of COVID-19 cases that peaked around April 11-15, followed by a decline
162 in cases until May 11. Soon thereafter, the slope of the case curve increased with
163 a very sharp uptick in confirmed cases beginning on June 12 (**Figure 1B**). We
164 consider May 11 as the transition between waves, as this date is the inflection
165 point of the cumulative new cases curve and had the absolute lowest number of
166 new cases in the mid-May time period. Thus, for the data presented herein, wave
167 1 is defined as March 5 through May 11, 2020, and wave 2 is defined as May 12
168 through July 7, 2020. Epidemiologic trends within the Houston Methodist Hospital
169 population were mirrored by data from Harris County and the greater
170 metropolitan Houston region (**Figure 1A**). Through the 7th of July, 25,366
171 COVID-19 cases were reported in Houston, 37,776 cases in Harris County, and
172 53,330 in metropolitan Houston, including 9,823 cases in Houston Methodist

173 facilities (inpatients and outpatients) (<https://www.tmc.edu/coronavirus-updates/infection-rate-in-the-greater-houston-area/> and
174 <https://harriscounty.maps.arcgis.com/apps/opsdashboard/index.html#/c0de71f8ea484b85bb5efcb7c07c6914>).
175
176

177 During the first wave (early March through May 11), 11,476 COVID-19
178 cases were reported in Houston, including 1,729 cases in the Houston Methodist
179 Hospital system. Early in the first wave (from March 5 through March 30, 2020),
180 we tested 3,080 patient specimens. Of these, 406 (13.2%) samples were positive
181 for SARS-CoV-2, representing 40% (358/898) of all confirmed cases in
182 metropolitan Houston during that time period. As our laboratory was the first
183 hospital-based facility to have molecular testing capacity for SARS-CoV-2
184 available on site, our strain samples are likely representative of COVID-19
185 infections during the first wave.

186 For the entire study period (March 5 through July 7, 2020), we tested
187 68,418 specimens from 55,800 patients. Of these, 9,121 patients (16.4%) had a
188 positive test result, representing 17.1% (9,121/53,300) of all confirmed cases in
189 metropolitan Houston. Thus, our strain samples are also representative of those
190 responsible for COVID-19 infections in the massive second wave.

191 To test the hypothesis that, on average, the two waves affected different
192 groups of patients, we analyzed individual patient characteristics (hospitalized
193 and non-hospitalized) in each wave. Consistent with this hypothesis, we found
194 significant differences in the COVID-19 patients in each wave (**Table S1**). For
195 example, patients in the second wave were significantly younger, had fewer

196 comorbidities, were more likely to be Hispanic/Latino (by self-report), and lived in
197 zip codes with lower median incomes (**Table S1**). A detailed analysis of the
198 characteristics of patients hospitalized in Houston Methodist facilities in the two
199 waves has recently been published (17).

200 **SARS-CoV-2 genome sequencing and phylogenetic analysis.** To
201 investigate the genomic architecture of the virus across the two waves, we
202 sequenced the genomes of 5,085 SARS-CoV-2 strains dating to the earliest time
203 of confirmed COVID-19 cases in Houston. Analysis of SARS-CoV-2 strains
204 causing disease in the first wave (March 5 through May 11) identified the
205 presence of many diverse virus genomes that, in the aggregate, represent the
206 major clades identified globally to date (**Figure 1B**). Clades G, GH, GR, and S
207 were the four most abundantly represented phylogenetic groups (**Figure 1B**).
208 Strains with the Gly614 amino acid variant in spike protein represented 82% of
209 the SARS-CoV-2 strains in wave 1, and 99.9% in wave 2 ($p<0.0001$; Fisher's
210 exact test) (**Figure 1B**). This spike protein variant is characteristic of clades G,
211 GH, and GR. Importantly, strains with the Gly614 variant represented only 71%
212 of the specimens sequenced in March, the early part of wave 1 (**Figure 1B**). We
213 attribute the decrease in strains with this variant observed in the first two weeks
214 of March (**Figure 1B**) to fluctuation caused by the relatively fewer COVID-19
215 cases occurring during this period.

216 **Relating spatiotemporal genome analysis with virus genotypes over**
217 **two disease waves.** We examined the spatial and temporal mapping of genomic
218 data to investigate community spread during wave 1 (**Figure 2**). Rapid and

219 widespread community dissemination occurred soon after the initial COVID-19
220 cases were reported in Houston. The heterogenous virus genotypes present very
221 early in wave 1 indicate that multiple strains independently entered metropolitan
222 Houston, rather than introduction and spread of a single strain. An important
223 observation was that strains of most of the individual subclades were distributed
224 over broad geographic areas (**Figure S1**). These findings are consistent with the
225 known ability of SARS-CoV-2 to spread very rapidly from person to person.

226 **Relationship between virus clades, clinical characteristics of infected**
227 **patients, and additional metadata.** It is possible that SARS-CoV-2 genome
228 subtypes have different clinical characteristics, analogous to what is believed to
229 have occurred with Ebola virus (18-20) and known to occur for other pathogenic
230 microbes (21). As an initial examination of this issue in SARS-CoV-2, we tested
231 the hypothesis that patients with disease severe enough to warrant
232 hospitalization were infected with a non-random subset of virus genotypes. We
233 also examined the association between virus clades and disease severity based
234 on overall mortality, highest level of required care (intensive care unit,
235 intermediate care unit, inpatient or outpatient), need for mechanical ventilation,
236 and length of stay. There was no simple relationship between virus clades and
237 disease severity using these four indicators. Similarly, there was no simple
238 relationship between virus clades and other metadata, such as sex, age, or
239 ethnicity (**Figure S2**).

240 **Machine learning analysis.** Machine learning models can be used to
241 identify complex relationships not revealed by statistical analyses. We built

242 machine learning models to test the hypothesis that virus genome sequence can
243 predict patient outcomes including mortality, length of stay, level of care, ICU
244 admission, supplemental oxygen use, and mechanical ventilation. Models to
245 predict outcomes based on virus genome sequence alone resulted in low F1
246 scores less than 50% (0.41 – 0.49) and regression models showed similarly low
247 R² values (-0.01 – -0.20) (**Table S2**). F1 scores near 50% are indicative of
248 classifiers that are performing similarly to random chance. The use of patient
249 metadata alone to predict patient outcome improved the model's F1 scores by 5-
250 10% (0.51 – 0.56) overall. The inclusion of patient metadata with virus genome
251 sequence data improved most predictions of outcomes, compared to genome
252 sequence alone, to 50% to 55% F1 overall (0.42 – 0.55) in the models (**Table**
253 **S2**). The findings are indicative of two possibilities that are not mutually
254 exclusive. First, patient metadata, such as age and sex, may provide more signal
255 for the model to use and thus result in better accuracies. Second, the model's
256 use of single nucleotide polymorphisms (SNPs) may have resulted in overfitting.
257 Most importantly, no SNP predicted a significant difference in outcome. A table of
258 classifier accuracy scores and performance information is provided in **Table S2**.

259 **Patient outcome and metadata correlations.** Overall, very few metadata
260 categories correlated with patient outcomes (**Table S3**). Mortality was
261 independently correlated with both Rh factor-positive blood type and increasing
262 age, with Pearson correlation coefficients (PCC) equal to 0.22 and 0.27,
263 respectively. This means that 22-27% of the variation in mortality can be
264 predicted from Rh factor-positive blood type and patient age, respectively. Length

265 of stay correlated independently with increasing age (PCC=0.20). Positive Rh
266 factor also correlated with ICU length of stay (PCC=0.20). All other patient
267 metadata correlations to outcomes had PCC less than 0.20 (**Table S3**).

268 We further analyzed outcomes correlated to isolates from wave 1 and 2,
269 and the presence of the Gly614 variant in spike protein. Being in wave 1 was
270 independently correlated with mechanical ventilation days, overall length of stay,
271 and ICU length of stay, with PCC equal to 0.20, 0.18, and 0.14, respectively.

272 Importantly, the presence of the Gly614 variant did not correlate with patient
273 outcomes (**Table S3**).

274 **Analysis of the *nsp12* polymerase gene.** The SARS-CoV-2 genome
275 encodes an RNA-dependent RNA polymerase (RdRp, also referred to as Nsp12)
276 used in virus replication (22-25). Two amino acid substitutions (Phe479Leu and
277 Val556Leu) in RdRp each confer significant resistance *in vitro* to remdesivir, an
278 adenosine analog (26). Remdesivir is inserted into RNA chains by RdRp during
279 replication, resulting in premature termination of RNA synthesis and inhibition of
280 virus replication. This compound has shown prophylactic and therapeutic benefit
281 against MERS-CoV and SARS-CoV-2 experimental infection in rhesus
282 macaques (27, 28). Recent reports indicate that remdesivir has therapeutic
283 benefit in some COVID-19 hospitalized patients (29-33), leading it to be now
284 widely used in patients worldwide. Thus, it may be important to understand
285 variation in RdRp in large strain samples.

286 To acquire data about allelic variation in the *nsp12* gene, we analyzed our
287 5,085 virus genomes. The analysis identified 265 SNPs, including 140

288 nonsynonymous (amino acid-altering) SNPs, resulting in amino acid
289 replacements throughout the protein (**Table 1**, **Figure 3**, **Figure 4**, **Figure S3**,
290 and **Figure S4**). The most common amino acid change was Pro322Leu,
291 identified in 4,893 of the 5,085 (96%) patient isolates. This amino acid
292 replacement is common in genomes from clades G, GH, and GR, which are
293 distinguished from other SARS-CoV-2 clades by the presence of the Gly614
294 amino acid change in the spike protein. Most of the other amino acid changes in
295 RdRp were present in relatively small numbers of strains, and some have been
296 identified in other isolates in a publicly available database (34). Five prominent
297 exceptions included amino acid replacements: Ala15Val in 138 strains, Met462Ile
298 in 59 strains, Met600Ile in 75 strains, Thr907Ile in 45 strains, and Pro917Ser in
299 80 strains. All 75 Met600Ile strains were phylogenetically closely related
300 members of clade G, and also had the Pro322Leu amino acid replacement
301 characteristic of this clade (**Figure S3**). These data indicate that the Met600Ile
302 change is likely the evolved state, derived from a precursor strain with the
303 Pro322Leu replacement. Similarly, we investigated phylogenetic relationships
304 among strains with the other four amino acid changes noted above. In all cases,
305 the vast majority of strains with each amino acid replacement were found among
306 individual subclades of strains (**Figure S3**).

307 Importantly, none of the observed amino acid polymorphisms in RdRp
308 were located precisely at two sites known to cause *in vitro* resistance to
309 remdesivir (26). Most of the amino acid changes are located distantly from the
310 RNA-binding and catalytic sites (**Figure S4** and **Table 1**). However,

311 replacements at six amino acid residues (Ala442Val, Ala448Val, Ala553Pro/Val,
312 Gly682Arg, Ser758Pro, and Cys812Phe) may potentially interfere with either
313 remdesivir binding or RNA synthesis. Four (Ala442Val, Ala448Val,
314 Ala553Pro/Val, and Gly682Arg) of the six substitution sites are located
315 immediately above the nucleotide-binding site, that is comprised of Lys544,
316 Arg552, and Arg554 residues as shown by structural studies (**Figure 4**). The
317 positions of these four variant amino acid sites are comparable to Val556 (**Figure**
318 **4**), for which a Val556Leu mutation in SARS-CoV was identified to confer
319 resistance to remdesivir *in vitro* (26). The other two substitutions (Ser758Pro and
320 Cys812Phe) are inferred to be located either at, or in the immediate proximity of,
321 the catalytic active site, that is comprised of three contiguous residues (Ser758,
322 Asp759, and Asp760). A proline substitution we identified at Ser758 (Ser758Pro)
323 is likely to negatively impact RNA synthesis. Although Cys812 is not directly
324 involved in the catalysis of RNA synthesis, it is only 3.5 Å away from Asp760.
325 The introduction of the bulkier phenylalanine substitution at Cys812 (Cys812Phe)
326 may impair RNA synthesis. Consequently, these two substitutions are expected
327 to detrimentally affect virus replication or fitness.

328 **Analysis of the gene encoding the spike protein.** The densely glycosylated
329 spike protein of SARS-CoV-2 and its close coronavirus relatives binds directly to
330 host-cell angiotensin-converting enzyme 2 (ACE2) receptors to enter host cells
331 (35-37). Thus, the spike protein is a major translational research target, including
332 intensive vaccine and therapeutic antibody (35-64). Analysis of the gene
333 encoding the spike protein identified 470 SNPs, including 285 that produce

334 amino acid changes (**Table 2, Figure 5**). Forty-nine of these replacements
335 (V11A, T51A, W64C, I119T, E156Q, S205A, D228G, L229W, P230T, N234D,
336 I235T, T274A, A288V, E324Q, E324V, S325P, S349F, S371P, S373P, T385I,
337 A419V, C480F, Y495S, L517F, K528R, Q628E, T632I, S708P, T719I, P728L,
338 S746P, E748K, G757V, V772A, K814R, D843N, S884A, M902I, I909V, E918Q,
339 S982L, M1029I, Q1142K, K1157M, Q1180R, D1199A, C1241F, C1247G, and
340 V1268A) are not represented in a publicly available database (34) as of August
341 19, 2020. Interestingly, 25 amino acid sites have three distinct variants (that is,
342 the reference amino acid plus two additional variant amino acids), and five amino
343 acid sites (amino acid positions 21, 27, 228, 936, and 1050) have four distinct
344 variants represented in our sample of 5,085 genomes (**Table 2, Figure 5**).

345 We mapped the location of amino acid replacements onto a model of the
346 full-length spike protein (35, 65) and observed that the substitutions are found in
347 each subunit and domain of the spike (**Figure 6**). However, the distribution of
348 amino acid changes is not uniform throughout the protein regions. For example,
349 compared to some other regions of the spike protein, the receptor binding
350 domain (RBD) has relatively few amino acid changes, and the frequency of
351 strains with these substitutions is low, each occurring in fewer than 10 isolates.
352 This finding is consistent with the functional constraints on RBD to mediate
353 interaction with ACE2. In contrast, the periphery of the S1 subunit NTD contains
354 a dense cluster of substituted residues, with some single amino acid
355 replacements found in 10–20 isolates (**Table 2, Figure 5, Figure 6**). Clustering
356 of amino acid changes in a distinct region of the spike protein may be a signal of

357 positive selection. Inasmuch as infected patients make antibodies against the
358 NTD, we favor the idea that host immune selection is one force contributing to
359 some of the amino acid variation in this region. One NTD substitution, H49Y, was
360 found in 142 isolates. This position is not well exposed on the surface of the NTD
361 and is likely not a result of immune pressure. The same is true for another highly
362 represented substitution, F1052L. This substitution was observed in 167 isolates,
363 and F1052 is buried within the core of the S2 subunit. The substitution observed
364 most frequently in the spike protein in our sample is D614G, a change observed
365 in 4,895 of the isolates. As noted above, strains with the Gly614 variant
366 significantly increased in wave 2 compared to wave 1.

367 As observed with RdRp, the majority of strains with each single amino
368 acid change in the spike protein were found on a distinct phylogenetic lineage
369 (**Figure S5**), indicating identity by descent. A prominent exception is the
370 Leu5Phe replacement that is present in all major clades, suggesting that this
371 amino acid change arose multiple times independently or very early in the course
372 of SARS-CoV-2 evolution. Finally, we note that examination of the phylogenetic
373 distribution of strains with multiple distinct amino acid replacements at the same
374 site (e.g., Arg21Ile/Lys/Thr, Ala27Ser/Thr/Val, etc.) revealed that they were
375 commonly found in different genetic branches, consistent with independent origin
376 (**Figure S5**).

377 **Cycle threshold (Ct) comparison of SARS-CoV-2 strains with either**
378 **the Asp614 or Gly614 amino acid replacements in spike protein.** It has been
379 reported that patients infected with strains having spike protein Gly614 variant

380 have, on average, higher virus loads on initial diagnosis (66-70). To determine if
381 this is the case in Houston strains, we examined the cycle threshold (Ct) for
382 every sequenced strain that was detected from a patient specimen using the
383 SARS-CoV-2 Assay done by the Hologic Panther instrument. We identified a
384 significant difference ($p<0.0001$) between the mean Ct value for strains with an
385 Asp614 ($n=102$) or Gly614 ($n=812$) variant of the spike protein (**Figure 7**).
386 Strains with Gly614 had a Ct value significantly lower than strains with the
387 Asp614 variant, indicating that patients infected with the Gly614 strains had, on
388 average, higher virus loads on initial diagnosis than patients infected by strains
389 with the Asp614 variant (**Figure 7**). This observation is consistent with the
390 conjecture that, on average, strains with the Gly614 variant are better able to
391 disseminate.

392 **Characterization of recombinant proteins with single amino acid**
393 **replacements in the receptor binding domain region of spike protein.** The
394 RBD of spike protein binds the ACE2 surface receptor and is also targeted by
395 neutralizing (36, 37, 41, 43-46, 48-62, 71). Thus, single amino acid replacements
396 in this domain may have functional consequences that enhance virus fitness. To
397 begin to test this idea, we expressed spike variants with the Asp614Gly
398 replacement and 13 clinical RBD variants identified in our genome sequencing
399 studies (**Figure 8, Table S4, S5**). All RBD variants were cloned into an
400 engineered spike protein construct that stabilizes the perfusion state and
401 increases overall expression yield (spike-6P, here referred to as spike) (64).

402 We first assessed the biophysical properties of spike-Asp614Gly, an
403 amino acid polymorphism that is common globally and increased significantly in
404 our wave 2 strain isolates. Pseudotyped viruses expressing spike-Gly614 have
405 higher infectivity for host cells *in vitro* than spike-Asp614 (66, 67, 69, 72, 73). The
406 higher infectivity of spike-Gly614 is correlated with increased stability and
407 incorporation of the spike protein into the pseudovirion (73). We observed a
408 higher expression level (**Figure 8A, B**) and increased thermostability for the
409 spike protein construct containing this variant (**Figure 8C, D**). The size exclusion
410 chromatography (SEC) elution profile of spike-Asp614 was indistinguishable from
411 spike-Gly614, consistent with a trimeric conformation (**Figure 8A**). These results
412 are broadly consistent with higher-resolution structural analyses of both spike
413 variants.

414 Next, we purified and biophysically characterized 13 RBD mutants that
415 each contain Gly614 and one additional single amino acid replacement we
416 identified by genome sequencing our clinical samples (**Table S6**). All variants
417 eluted as trimers, indicating the global structure, remained intact (**Figure 8** and
418 **Figure S6**). However, several variants had reduced expression levels and
419 virtually all had decreased thermostability relative to the variant that had only a
420 D614G single amino acid replacement (**Figure 8D**). The A419V and A522V
421 mutations were especially deleterious, reducing yield and precluding further
422 downstream analysis (**Figure 8B**). We next assayed the affinity of the 11 highest-
423 expressing spike variants for ACE2 and the neutralizing monoclonal antibody
424 CR3022 via enzyme-linked immunosorbent assays (ELISAs) (**Figure 8E-G** and

425 **Table S6).** Most variants retained high affinity for the ACE2 surface receptor.
426 However, importantly, three RBD variants (F338L, S373P, and R408T) had
427 substantially reduced affinity for CR3022, a monoclonal antibody that disrupts the
428 spike protein homotrimerization interface (63, 74). Notably, the S373P mutation
429 is one amino acid away from the epitope recognized by CR3022. These results
430 are consistent with the interpretation that some RBD mutants arising in COVID-
431 19 patients may have increased ability to escape humoral immune pressure, but
432 otherwise retain strong ACE2 binding affinity.

433

434 **DISCUSSION**

435 In this work we analyzed the molecular population genomics, sociodemographic,
436 and medical features of two waves of COVID-19 disease occurring in
437 metropolitan Houston, Texas, between early March and early July 2020. We also
438 studied the biophysical and immunologic properties of some naturally occurring
439 single amino acid changes in the spike protein RBD identified by sequencing the
440 5,085 genomes. We discovered that the first COVID-19 wave was caused by a
441 heterogenous array of virus genotypes assigned to several different clades. The
442 majority of cases in the first wave are related to strains that caused widespread
443 disease in European and Asian countries, as well as other localities. We
444 conclude that the SARS-CoV-2 virus was introduced into Houston many times
445 independently, likely by individuals who had traveled to or from different parts of
446 the world, including other communities in the United States. In support of this
447 conclusion, the first cases in metropolitan Houston were associated with a travel

448 history to a known COVID-19 region (16). The data are consistent with the fact
449 that Houston is a large international city characterized by a multi-ethnic
450 population and is a prominent transport hub with direct flights to major cities
451 globally.

452 The second wave of COVID-19 cases also is characterized by SARS-
453 CoV-2 strains with diverse genotypes. Virtually all cases in the second and
454 ongoing disease wave were caused by strains with the Gly614 variant of spike
455 protein (**Figure 1B**). Our data unambiguously demonstrate that strains with the
456 Gly614 variant increased significantly in frequency in wave 2 relative to wave 1 in
457 the Houston metropolitan region. This shift occurred very rapidly in a matter of
458 just a few months. Amino acid residue Asp614 is located in subdomain 2 (SD-2)
459 of the spike protein and forms a hydrogen bond and electrostatic interaction with
460 two residues in the S2 subunit of a neighboring protomer. Replacement of
461 aspartate with glycine would eliminate both interactions, thereby substantively
462 weakening the contact between the S1 and S2 subunits. We previously
463 speculated (75) that this weakening produces a more fusogenic spike protein, as
464 S1 must first dissociate from S2 before S2 can refold and mediate fusion of virus
465 and cell membranes. Stated another way, virus strains with the Gly614 variant
466 may be better able to enter host cells, potentially resulting in enhanced spread.
467 Consistent with this idea, Korber et al. (66) showed that the Gly614 variant grows
468 to higher titer as pseudotyped virions. On initial diagnosis infected individuals had
469 lower RT-PCR cycle thresholds suggesting higher upper respiratory tract viral
470 loads. Our data (**Figure 7**) are fully consistent with that finding Zhang et al. (73)

471 reported that pseudovirus with the 614Gly variant infected ACE2-expressing cells
472 more efficiently than the 614Asp. Similar results have been described by Hu et
473 al. (67) and Lorenzo-Redondo et al. (68). Plante et al. (76) recently studied
474 isogenic mutant SARS-CoV-2 strains with either the 614Asp or 614Gly variant
475 and found that the 614Gly variant virus had significantly increased replication in
476 human lung epithelial cells *in vitro* and increased infectious titers in nasal and
477 trachea washes obtained from experimentally infected hamsters. These results
478 are consistent with the idea that the 614Gly variant bestows increased virus
479 fitness in the upper respiratory tract (76).

480 Additional work is needed to investigate the potential biomedical relevance
481 and public health importance of the Asp614Gly polymorphism, including but not
482 limited to virus dissemination, overall fitness, impact on clinical course and
483 virulence, and development of vaccines and therapeutics. Although it is possible
484 that stochastic processes alone may account for the rapid increase in COVID-19
485 disease frequency caused by viruses containing the Gly614 variant, we do not
486 favor that interpretation in part because of the cumulative weight of the
487 epidemiologic, human RT-PCR diagnostics data, *in vitro* experimental findings,
488 and animal infection studies using isogenic mutant virus strains. In addition, if
489 stochastic processes solely are responsible, we believe it is difficult to explain
490 essentially simultaneous increase in frequency of the Gly614 variant in
491 genetically diverse viruses in three distinct clades (G, GH, and GR) in a
492 geographically large metropolitan area with 7 million ethnically diverse people.
493 Regardless, more research on this important topic is warranted.

494 Several groups have analyzed the relationship between human genetics,
495 blood types, and patient outcomes (33, 34, 77-80). In this regard, we identified a
496 correlation between Rh factor-positive blood type and increased mortality and
497 length of ICU stay. (Supplemental Table 3). Thus, our data are consistent with
498 studies suggesting that there are host genetic factors that contribute to disease
499 severity and outcome. More research is required to address the important
500 possibility.

501 The diversity present in our 1,026 virus genomes from the first disease
502 wave contrasts somewhat with data reported by Gonzalez-Reiche et al., who
503 studied 84 SARS-CoV-2 isolates causing disease in patients in the New York
504 City region (11). Those investigators concluded that the vast majority of disease
505 was caused by progeny of strains imported from Europe. Similarly, Bedford et al.
506 (10) reported that much of the COVID-19 disease in the Seattle, Washington
507 area was caused by strains that are progeny of a virus strain recently introduced
508 from China. Some aspects of our findings are similar to those reported recently
509 by Lemieux et al. based on analysis of strains causing disease in the Boston
510 area (81). Our findings, like theirs, highlight the importance of multiple
511 importation events of genetically diverse strains in the epidemiology of COVID-19
512 disease in this pandemic. Similarly, Icelandic and Brazilian investigators
513 documented that SARS-CoV-2 was imported by individuals traveling to or from
514 many European and other countries (82, 83).

515 The virus genome diversity and large sample size in our study permitted
516 us to test the hypothesis that distinct virus clades were nonrandomly associated

517 with hospitalized COVID-19 patients or disease severity. We did not find
518 evidence to support this hypothesis, but our continuing study of COVID-19 cases
519 accruing in the second wave will further improve statistical stratification.

520 We used machine learning classifiers to identify if any SNPs contribute to
521 increased infection severity or otherwise affect virus-host outcome. The models
522 could not be trained to accurately predict these outcomes from the available virus
523 genome sequence data. This may be due to sample size or class imbalance.

524 However, we do not favor this interpretation. Rather, we think that the inability to
525 identify particular virus SNPs predictive of disease severity or infection outcome
526 likely reflects the substantial heterogeneity in underlying medical conditions and
527 treatment regimens among COVID-19 patients studied herein. An alternative but
528 not mutually exclusive hypothesis is that patient genotypes play an important role
529 in determining virus-human interactions and resulting pathology. Although some
530 evidence has been presented in support of this idea (84, 85), available data
531 suggest that in the aggregate, host genetics does not play an overwhelming role
532 in determining outcome in the great majority of adult patients, once virus infection
533 is established.

534 Remdesivir is a nucleoside analog reported to have activity against
535 MERS-CoV, a coronavirus related to SARS-CoV-2. Recently, several studies
536 have reported that remdesivir shows promise in treating COVID-19 patients (29-
537 33), leading the FDA to issue an emergency use authorization. Because *in vitro*
538 resistance of SARS-CoV to remdesivir has been reported to be caused by either
539 of two amino acid replacements in RdRp (Phe479Leu or Val556Leu), we

540 interrogated our data for polymorphisms in the *nsp12* gene. Although we
541 identified 140 different inferred amino acid replacements in RdRp in the 5,085
542 genomes analyzed, none of these were located precisely at the two positions
543 associated with *in vitro* resistance to remdesivir. Inasmuch as remdesivir is now
544 being deployed widely to treat COVID-19 patients in Houston and elsewhere, our
545 findings suggest that the majority of SARS-CoV-2 strains currently circulating in
546 our region should be susceptible to this drug.

547 The amino acid replacements Ala442Val, Ala448Val, Ala553Pro/Val, and
548 Gly682Arg that we identified occur at sites that, intriguingly, are located directly
549 above the nucleotide substrate entry channel and nucleotide binding residues
550 Lys544, Arg552, and Arg554 (22, 23) (**Figure 4**). One possibility is that
551 substitution of the smaller alanine or glycine residues with the bulkier side chains
552 of Val/Pro/Arg may impose structural constraints for the modified nucleotide
553 analog to bind, and thereby disfavor remdesivir binding. This, in turn, may lead to
554 reduced incorporation of remdesivir into the nascent RNA, increased fidelity of
555 RNA synthesis, and ultimately drug resistance. A similar mechanism has been
556 proposed for a Val556Leu change (23).

557 We also identified one strain with a Lys477Asn replacement in RdRp. This
558 substitution is located close to a Phe479Leu replacement reported to produce
559 partial resistance to remdesivir *in vitro* in SARS-CoV patients from 2004,
560 although the amino acid positions are numbered differently in SARS-CoV and
561 SARS-CoV-2. Structural studies have suggested that this amino acid is surface-
562 exposed, and distant from known key functional elements. Our observed

563 Lys477Asn change is also located in a conserved motif described as a finger
564 domain of RdRp (**Figure 3 and 4**). One speculative possibility is that Lys477 is
565 involved in binding a yet unidentified cofactor such as Nsp7 or Nsp8, an
566 interaction that could modify nucleotide binding and/or fidelity at a distance.
567 These data warrant additional study in larger patient cohorts, especially in
568 individuals treated with remdesivir.

569 Analysis of the gene encoding the spike protein identified 285 polymorphic
570 amino acid sites relative to the reference genome, including 49 inferred amino
571 acid replacements not present in available databases as of August 19, 2020.
572 Importantly, 30 amino acid sites in the spike protein had two or three distinct
573 replacements relative to the reference strain. The occurrence of multiple variants
574 at the same amino acid site is one characteristic that may suggest functional
575 consequences. These data, coupled with structural information available for
576 spike protein, raise the possibility that some of the amino acid variants have
577 functional consequences, for example including altered serologic reactivity and
578 shown here. These data permit generation of many biomedically relevant
579 hypotheses now under study.

580 A recent study reported that RBD amino acid changes could be selected
581 *in vitro* using a pseudovirus neutralization assay and sera obtained from
582 convalescent plasma or monoclonal antibodies (86). The amino acid sites
583 included positions V445 and E484 in the RBD. Important to note, variants G446V
584 and E484Q were present in our patient samples. However, these mutations
585 retain high affinity to CR3022 (**Figure 8F, G**). The high-resolution structure of the

586 RBD/CR3022 complex shows that CR3022 makes contacts to residues 369-386,
587 380-392, and 427-430 of RBD (74). Although there is no overlap between
588 CR3022 and ACE2 epitopes, CR3022 is able to neutralize the virus through an
589 allosteric effect. We found that the Ser373Pro change, which is located within the
590 CR3022 epitope, has reduced affinity to CR3022 (**Figure 8F, G**). The F338L and
591 R408T mutations, although not found directly within the interacting epitope, also
592 display reduced binding to CR3022. Other investigators (86) using *in vitro*
593 antibody selection identified a change at amino acid site S151 in the N-terminal
594 domain, and we found mutations S151N and S151I in our patient samples. We
595 also note that two variant amino acids (Gly446Val and Phe456Leu) we identified
596 are located in a linear epitope found to be critical for a neutralizing monoclonal
597 antibody described recently by Li et al. (87).

598 In the aggregate, these findings suggest that mutations emerging within
599 the spike protein at positions within and proximal to known neutralization
600 epitopes may result in escape from antibodies and other therapeutics currently
601 under development. Importantly, our study did not reveal that these mutant
602 strains had disproportionately increased over time. The findings may also bear
603 on the occurrence of multiple amino acid substitutions at the same amino acid
604 site that we identified in this study, commonly a signal of selection. In the
605 aggregate, the data support a multifaceted approach to serological monitoring
606 and biologics development, including the use of monoclonal antibody cocktails
607 (46, 47, 88).

608

609 CONCLUDING STATEMENT

610 Our work represents analysis of the largest sample to date of SARS-CoV-2
611 genome sequences from patients in one metropolitan region in the United States.
612 The investigation was facilitated by the fact that we had rapidly assessed a
613 SARS-CoV-2 molecular diagnostic test in January 2020, more than a month
614 before the first COVID-19 patient was diagnosed in Houston. In addition, our
615 large healthcare system has seven hospitals and many facilities (e.g., outpatient
616 care centers, emergency departments) located in geographically diverse areas of
617 the city. We also provide reference laboratory services for other healthcare
618 entities in the Houston area. Together, our facilities serve patients of diverse
619 ethnicities and socioeconomic status. Thus, the data presented here likely reflect
620 a broad overview of virus diversity causing COVID-19 infections throughout
621 metropolitan Houston. We previously exploited these features to study influenza
622 and *Klebsiella pneumoniae* dissemination in metropolitan Houston (89, 90). We
623 acknowledge that every “twig” of the SARS-CoV-2 evolutionary tree in Houston is
624 not represented in these data. The samples studied are not comprehensive for
625 the entire metropolitan region. For example, it is possible that our strain samples
626 are not fully representative of individuals who are indigent, homeless, or of very
627 low socioeconomic groups. In addition, although the strain sample size is
628 relatively large compared to other studies, the sample represents only about 10%
629 of all COVID-19 cases in metropolitan Houston documented in the study period.
630 In addition, some patient samples contain relatively small amounts of virus
631 nucleic acid and do not yield adequate sequence data for high-quality genome

632 analysis. Thus, our data likely underestimate the extent of genome diversity
633 present among SARS-CoV-2 causing COVID-19 and will not identify all amino
634 acid replacements in the virus in this geographic region. It will be important to
635 sequence and analyze the genomes of additional SARS-CoV-2 strains causing
636 COVID-19 cases in the ongoing second massive disease wave in metropolitan
637 Houston, and these studies are underway. Data of this type will be especially
638 important to have if a third and subsequent waves were to occur in metropolitan
639 Houston, as it could provide insight into molecular and epidemiologic events
640 contributing to them.

641 The genomes reported here are an important data resource that will
642 underpin our ongoing study of SARS-CoV-2 molecular evolution, dissemination,
643 and medical features of COVID-19 in Houston. As of August 19, 2020, there
644 were 135,866 reported cases of COVID-19 in metropolitan Houston, and the
645 number of cases is increasing daily. Although the full array of factors contributing
646 to the massive second wave in Houston is not known, it is possible that the
647 potential for increased transmissibility of SARS-CoV-2 with the Gly614 may have
648 played a role, as well as changes in behavior associated with the Memorial Day
649 and July 4th holidays, and relaxation of some of the social constraints imposed
650 during the first wave. The availability of extensive virus genome data dating from
651 the earliest reported cases of COVID-19 in metropolitan Houston, coupled with
652 the database we have now constructed, may provide critical insights into the
653 origin of new infection spikes and waves occurring as public health constraints
654 are further relaxed, schools and colleges re-open, holidays occur, commercial air

655 travel increases, and individuals change their behavior because of COVID-19
656 “fatigue.” The genome data will also be useful in assessing ongoing molecular
657 evolution in spike and other proteins as baseline herd immunity is generated,
658 either by natural exposure to SARS-CoV-2 or by vaccination. The signal of
659 potential selection contributing to some spike protein diversity and identification
660 of naturally occurring mutant RBD variants with altered serologic recognition
661 warrant close attention and expanded study.

662

663 MATERIALS AND METHODS

664 **Patient specimens.** All specimens were obtained from individuals who
665 were registered patients at Houston Methodist hospitals, associated facilities
666 (e.g., urgent care centers), or institutions in the greater Houston metropolitan
667 region that use our laboratory services. Virtually all individuals met the criteria
668 specified by the Centers for Disease Control and Prevention to be classified as a
669 person under investigation.

670

671 **SARS-CoV-2 molecular diagnostic testing.** Specimens obtained from
672 symptomatic patients with a high degree of suspicion for COVID-19 disease were
673 tested in the Molecular Diagnostics Laboratory at Houston Methodist Hospital
674 using an assay granted Emergency Use Authorization (EUA) from the FDA
675 (<https://www.fda.gov/medical-devices/emergency-situations-medical-devices/faqs-diagnostic-testing-sars-cov-2#offeringtests>). Multiple testing
676 platforms were used, including an assay that follows the protocol published by
677

678 the WHO (<https://www.who.int/docs/default-source/coronavirus/protocol-v2-1.pdf>) using the EZ1 virus extraction kit and EZ1 Advanced XL instrument or
679 QIASymphony DSP Virus kit and QIASymphony instrument for nucleic acid
680 extraction and ABI 7500 Fast Dx instrument with 7500 SDS software for reverse
681 transcription RT-PCR, the COVID-19 test using BioFire Film Array 2.0
682 instruments, the Xpert Xpress SARS-CoV-2 test using Cepheid GeneXpert
683 Infinity or Cepheid GeneXpert Xpress IV instruments, the SARS-CoV-2 Assay
684 using the Hologic Panther instrument, and the Aptima SARS-CoV-2 Assay using
685 the Hologic Panther Fusion system. All assays were performed according to the
686 manufacturer's instructions. Testing was performed on material obtained from
687 nasopharyngeal or oropharyngeal swabs immersed in universal transport media
688 (UTM), bronchoalveolar lavage fluid, or sputum treated with dithiothreitol (DTT).
689 To standardize specimen collection, an instructional video was created for
690 Houston Methodist healthcare workers
691 (<https://vimeo.com/396996468/2228335d56>).
692

693
694 **Epidemiologic curve.** The number of confirmed COVID-19 positive cases
695 was obtained from USAFacts.org (<https://usafacts.org/visualizations/coronavirus-covid-19-spread-map/>) for Austin, Brazoria, Chambers, Fort Bend, Galveston,
696 Harris, Liberty, Montgomery, and Waller counties. Positive cases for Houston
697 Methodist Hospital patients were obtained from our Laboratory Information
698 System and plotted using the documented collection time.
699

700

701 **SARS-CoV-2 genome sequencing.** Libraries for whole virus genome
702 sequencing were prepared according to version 1 or version 3 of the ARTIC
703 nCoV-2019 sequencing protocol (<https://artic.network/ncov-2019>). Long reads
704 were generated with the LSK-109 sequencing kit, 24 native barcodes (NBD104
705 and NBD114 kits), and a GridION instrument (Oxford Nanopore). Short reads
706 were generated with the NexteraXT kit and NextSeq 550 instrument (Illumina).

707

708 **SARS-CoV-2 genome sequence analysis.** Consensus virus genome
709 sequences from the Houston area isolates were generated using the ARTIC
710 nCoV-2019 bioinformatics pipeline. Publicly available genomes and metadata
711 were acquired through GISAID on August 19, 2020. GISAID sequences
712 containing greater than 1% N characters, and Houston sequences with greater
713 than 5% N characters were removed from consideration. Identical GISAID
714 sequences originating from the same geographic location with the same
715 collection date were also removed from consideration to reduce redundancy.
716 Nucleotide sequence alignments for the combined Houston and GISAID strains
717 were generated using MAFFT version 7.130b with default parameters (91).
718 Sequences were manually curated in JalView (92) to trim the ends and to
719 remove sequences containing spurious inserts. Phylogenetic trees were
720 generated using FastTree with the generalized time-reversible model for
721 nucleotide sequences (93). CLC Genomics Workbench (QIAGEN) was used to
722 generate the phylogenetic tree figures.

723

724 **Geospatial mapping.** The home address zip code for all SARS-CoV-2
725 positive patients was used to generate the geospatial maps. To examine
726 geographic relatedness among genetically similar isolates, geospatial maps were
727 filtered to isolates containing specific amino acid changes.

728

729 **Time series.** Geospatial data were filtered into wave 1 (3/5/2020-
730 5/11/2020) and wave 2 (5/12/2020-7/7/2020) time intervals to illustrate the
731 spread of confirmed SARS-CoV-2 positive patients identified over time.

732

733 **Machine learning.** Virus genome alignments and patient metadata were
734 used to build models to predict patient metadata and outcomes using both
735 classification models and regression. Metadata considered for prediction in the
736 classification models included age, ABO and Rh blood type, ethnic group,
737 ethnicity, sex, ICU admission, IMU admission, supplemental oxygen use, and
738 ventilator use. Metadata considered for prediction in regression analysis included
739 ICU length of stay, IMU length of stay, total length of stay, supplemental oxygen
740 use, and ventilator use. Because sex, blood type, Rh factor, age, age decade,
741 ethnicity, and ethnic group are features in the patient features and combined
742 feature sets, models were not trained for these labels using patient and
743 combined feature sets. Additionally, age, length of stay, IMU length of stay, ICU
744 length of stay, mechanical ventilation days, and supplemental oxygen days were
745 treated as regression problems and XGBoost regressors were built while the rest
746 were treated as classification problems and XGBoost classifiers were built.

747 Three types of features were considered for training the XGBoost
748 classifiers: alignment features, patient features, and the combination of alignment
749 and patient features. Alignment features were generated from the consensus
750 genome alignment such that columns containing ambiguous nucleotide bases
751 were removed to ensure the models did not learn patterns from areas of low
752 coverage. These alignments were then one-hot encoded to form the alignment
753 features. Patient metadata values were one-hot encoded with the exception of
754 age, which remained as a raw integer value, to create the patient features.
755 These metadata values consisted of age, ABO, Rh blood type, ethnic group,
756 ethnicity, and sex. All three types of feature sets were used to train models that
757 predict ICU length of stay, IMU length of stay, overall length of stay, days of
758 supplemental oxygen therapy, and days of ventilator usage while only alignment
759 features were used to train models that predict age, ABO, Rh blood type, ethnic
760 group, ethnicity, and sex.

761 A ten-fold cross validation was used to train XGBoost models (94) as
762 described previously (95, 96). Depths of 4, 8, 16, 32, and 64 were used to tune
763 the models, but accuracies plateaued after a depth of 16. SciKit-Learn's (97)
764 classification report and r2 score were then used to access overall accuracy of
765 the classification and regression models, respectively.

766

767 **Patient metadata correlations.** We encoded values into multiple columns
768 for each metadata field. For example, the ABO column was divided into four
769 columns for A, B, AB, and O blood type. Those columns were encoded with a 1

770 for the patients' ABO type, with all other columns encoded with 0. This was
771 repeated for all non-outcome metadata fields. Age, however, was not re-
772 encoded, as the raw integer values were used. Each column was then correlated
773 to the various outcome values for each patient (deceased, ICU length, IMU
774 length, length of stay, supplemental oxygen length, and ventilator length) to
775 obtain a Pearson coefficient correlation value for each metadata label and
776 outcome.

777

778 **Analysis of the *nsp12* polymerase and S protein genes.** The *nsp12*
779 virus polymerase and S protein genes were analyzed by plotting SNP density in
780 the consensus alignment using Python (Python v3.4.3, Biopython Package
781 v1.72). The frequency of SNPs in the Houston isolates was assessed, along with
782 amino acid changes for nonsynonymous SNPs.

783

784 **Cycle threshold (Ct) comparison of SARS-CoV-2 strains with either**
785 **Asp614 or Gly614 amino acid replacements in the spike protein.** The cycle
786 threshold (Ct) for every sequenced strain that was detected from a patient
787 specimen using the SARS-CoV-2 Assay on the Hologic Panther instrument was
788 retrieved from the Houston Methodist Hospital Laboratory Information System.
789 Statistical significance between the mean Ct value for strains with an aspartate
790 ($n=102$) or glycine ($n=812$) amino acid at position 614 of the spike protein was
791 determined with the Mann-Whitney test (GraphPad PRISM 8).

792

793 **Creation and characterization of spike protein RBD variants.** Spike
794 RBD variants were cloned into the spike-6P (HexaPro; F817P, A892P, A899P,
795 A942P, K986P, V987P) base construct that also includes the D614G substitution
796 (pIF638). Briefly, a segment of the gene encoding the RBD was excised with
797 EcoRI and NheI, mutagenized by PCR, and assembled with a HiFi DNA
798 Assembly Cloning Kit (NEB).

799 FreeStyle 293-F cells (Thermo Fisher Scientific) were cultured and
800 maintained in a humidified atmosphere of 37°C and 8% CO₂ while shaking at
801 110-125rpm. Cells were transfected with plasmids encoding spike protein
802 variants using polyethylenimine. Three hours post-transfection, 5µM kifunensine
803 was added to each culture. Cells were harvested four days after transfection and
804 the protein containing supernatant was separated from the cells by two
805 centrifugation steps: 10 min at 500rcf and 20 min at 10,000rcf. Supernatants
806 were kept at 4°C throughout. Clarified supernatant was loaded on a Poly-Prep
807 chromatography column (Bio-Rad) containing Strep-Tactin Superflow resin (IBA),
808 washed with five column volumes (CV) of wash buffer (100mM Tris-HCl pH 8.0,
809 150mM NaCl; 1mM EDTA), and eluted with four CV of elution buffer (100mM
810 Tris-HCl pH 8.0, 150mM NaCl, 1mM EDTA, 2.5mM d-Desthiobiotin). The eluate
811 was spin-concentrated (Amicon Ultra-15) to 600µL and further purified via size-
812 exclusion chromatography (SEC) using a Superose 6 Increase 10/300 column
813 (G.E.) in SEC buffer (2mM Tris pH 8.0, 200mM NaCl and 0.02% NaN₃). Proteins
814 were concentrated to 300µL and stored in SEC buffer.

815 The RBD spike mutants chosen for analysis were all RBD amino acid
816 mutants identified by our genome sequencing study as of June 15, 2020. We
817 note that the exact boundaries of the RBD domain varies depending on the paper
818 used as reference. We used the boundaries demarcated in Figure 1A of Cai et al.
819 Science paper 21 July) (98) that have K528R located at the RBD-CTD1 interface.
820

821 **Differential scanning fluorimetry.** Recombinant spike proteins were
822 diluted to a final concentration of 0.05mg/mL with 5X SYPRO orange (Sigma) in
823 a 96-well qPCR plate. Continuous fluorescence measurements ($\lambda_{\text{ex}}=465\text{nm}$,
824 $\lambda_{\text{em}}=580\text{nm}$) were collected with a Roche LightCycler 480 II. The temperature
825 was increased from 22°C to 95°C at a rate of 4.4°C/min. We report the first
826 melting transition.

827

828 **Enzyme-linked immunosorbent assays.** ELISAs were performed to
829 characterize binding of S6P, S6P D614G, and S6P D614G-RBD variants to
830 human ACE2 and the RBD-binding monoclonal antibody CR3022. The ACE2-
831 hFc chimera was obtained from GenScript (Z03484), and the CR3022 antibody
832 was purchased from Abcam (Ab273073). Corning 96-well high-binding plates
833 (CLS9018BC) were coated with spike variants at 2 $\mu\text{g}/\text{mL}$ overnight at 4°C. After
834 washing four times with phosphate buffered saline + 0.1% Tween20 (PBST;
835 300 $\mu\text{L}/\text{well}$), plates were blocked with PBS+2% milk (PBSM) for 2 h at room
836 temperature and again washed four times with PBST. These were serially diluted
837 in PBSM 1:3 seven times in triplicate. After 1 h incubation at room temperature,

838 plates were washed four times in PBST, labeled with 50 μ L mouse anti-human
839 IgG1 Fc-HRP (SouthernBlots, 9054-05) for 45 min in PBSM, and washed again
840 in PBST before addition of 50 μ L 1-step Ultra TMB-ELISA substrate (Thermo
841 Scientific, 34028). Reactions were developed for 15 min and stopped by addition
842 of 50 μ L 4M H₂SO₄. Absorbance intensity (450nm) was normalized within a plate
843 and EC₅₀ values were calculated through 4-parameter logistic curve (4PL)
844 analysis using GraphPad PRISM 8.4.3.

845

846 **ACKNOWLEDGMENTS**

847 We thank Dr. Steven Hinrichs and colleagues at the Nebraska Public Health
848 Laboratory, and Dr. David Persse and colleagues at the Houston Health
849 Department for providing samples used to validate our initial SARS-CoV-2
850 molecular assay. We thank Drs. Jessica Thomas and Zejuan Li, Erika Walker,
851 Concepcion C. Cantu, the very talented and dedicated molecular technologists,
852 and the many labor pool volunteers in the Molecular Diagnostics Laboratory for
853 their dedication to patient care. We also thank Brandi Robinson, Harrold Cano,
854 Cory Romero, Brooke Burns, and Hayder Mahmood for technical assistance. We
855 are indebted to Drs. Marc Boom and Dirk Sostman for their support, and to many
856 very generous Houston philanthropists for their tremendous support of this
857 ongoing project, including but not limited to anonymous, Ann and John Bookout
858 III, Carolyn and John Bookout, Ting Tsung and Wei Fong Chao Foundation, Ann
859 and Leslie Doggett, Freeport LNG, the Hearst Foundations, Jerold B. Katz
860 Foundation, C. James and Carole Walter Looke, Diane and David Modesett, the

861 Sherman Foundation, and Paula and Joseph C. "Rusty" Walter III. We gratefully
862 acknowledge the originating and submitting laboratories of the SARS-CoV-2
863 genome sequences from GISAID's EpiFluTM Database used in some of the work
864 presented here. We also thank many colleagues for critical reading of the
865 manuscript and suggesting improvements, and Sasha Pejerrey, Adrienne
866 Winston, Heather McConnell, and Kathryn Stockbauer for editorial contributions.
867 We are especially indebted to Drs. Nancy Jenkins and Neal Copeland for their
868 scholarly suggestions to improve an early version of the manuscript.

869

870 **Author contributions:** J.M.M. conceptualized and designed the project; S.W.L,
871 R.J.O., P.A.C., D.W.B., J.J.D., M.S., M.N., M.O.S., C.C.C., P.Y., L.P., S.S., H.-C.
872 K., H.H., G.E., H.A.T.N., J.H.L., M.K., J.G., D.B., J.G., J.S.M., C.-W.C., K.J., and
873 I.F. performed research. All authors contributed to writing the manuscript.

874 Data and material availability: The spike-6P ("HexaPro") plasmid is available from
875 Addgene (ID: 154754) or from I.J.F. under a material transfer agreement with
876 The University of Texas at Austin. Additional plasmids are available upon request
877 from I.J.F.

878

879 This study was supported by the Fondren Foundation, Houston Methodist
880 Hospital and Research Institute (to J.M.M.), NIH grant AI127521 (to J.S.M.), NIH
881 grants GM120554 and GM124141 to I.J.F., the Welch Foundation (F-1808 to
882 I.J.F.), and the National Science Foundation (1453358 to I.J.F.). I.J.F. is a CPRIT
883 Scholar in Cancer Research. J.J.D., M.S., and M.N. are supported by the NIAID

It is made available under a CC-BY-NC-ND 4.0 International license .

884 Bacterial and Viral Bioinformatics resource center award (contract number
885 75N93019C00076). J.J.D. and M.N. are also funded by the United States
886 Defense Advanced Research Projects Agency Award iSENTRY Friend or Foe
887 program award (contract number HR0011937807).

888 REFERENCES

- 889 1. 2020. World Health Organization Coronavirus Disease 2019 (COVID-19) Situation Report.
890 https://www.who.int/docs/default-source/coronavirus/situation-reports/20200420-sitrep-91-covid-19.pdf?sfvrsn=fclf0670b_4. Accessed April 21.
- 892 2. Gorbalenya AE, Baker SC, Baric RS, de Groot RJ, Drosten C, Gulyaeva AA, Haagmans BL,
893 Lauber C, Leontovich AM, Neuman BW, Penzar D, Perlman S, Poon LLM, Samborskiy DV,
894 Sidorov IA, Sola I, Ziebuhr J, Coronaviridae Study Group of the International Committee on
895 Taxonomy of V. 2020. The species Severe acute respiratory syndrome-related coronavirus:
896 classifying 2019-nCoV and naming it SARS-CoV-2. *Nature Microbiology* 5:536-544.
- 897 3. Wang C, Horby PW, Hayden FG, Gao GF. 2020. A novel coronavirus outbreak of global health
898 concern. *Lancet* 395:470-473.
- 899 4. Perlman S. 2020. Another Decade, Another Coronavirus. *New England Journal of Medicine*
900 382:760-762.
- 901 5. Allel K, Tapia-Muñoz T, Morris W. 2020. Country-level factors associated with the early spread
902 of COVID-19 cases at 5, 10 and 15 days since the onset. *Glob Public Health*
903 doi:10.1080/17441692.2020.1814835:1-14.
- 904 6. Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, Zhang L, Fan G, Xu J, Gu X, Cheng Z, Yu T, Xia
905 J, Wei Y, Wu W, Xie X, Yin W, Li H, Liu M, Xiao Y, Gao H, Guo L, Xie J, Wang G, Jiang R,
906 Gao Z, Jin Q, Wang J, Cao B. 2020. Clinical features of patients infected with 2019 novel
907 coronavirus in Wuhan, China. *Lancet* 395:497-506.
- 908 7. Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, Zhao X, Huang B, Shi W, Lu R, Niu P, Zhan F,
909 Ma X, Wang D, Xu W, Wu G, Gao GF, Tan W. 2020. A Novel Coronavirus from Patients with
910 Pneumonia in China, 2019. *New England Journal of Medicine* 382:727-733.
- 911 8. Chan JF, Yuan S, Kok KH, To KK, Chu H, Yang J, Xing F, Liu J, Yip CC, Poon RW, Tsui HW,
912 Lo SK, Chan KH, Poon VK, Chan WM, Ip JD, Cai JP, Cheng VC, Chen H, Hui CK, Yuen KY.
913 2020. A familial cluster of pneumonia associated with the 2019 novel coronavirus indicating
914 person-to-person transmission: a study of a family cluster. *Lancet* 395:514-523.
- 915 9. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH, Pei YY, Yuan ML,
916 Zhang YL, Dai FH, Liu Y, Wang QM, Zheng JJ, Xu L, Holmes EC, Zhang YZ. 2020. A new
917 coronavirus associated with human respiratory disease in China. *Nature* 579:265-269.
- 918 10. Bedford T, Greninger AL, Roychoudhury P, Starita LM, Famulare M, Huang M-L, Nalla A,
919 Pepper G, Reinhardt A, Xie H, Shrestha L, Nguyen TN, Adler A, Brandstetter E, Cho S, Giroux
920 D, Han PD, Fay K, Frazer CD, Ilcisin M, Lacombe K, Lee J, Kiavand A, Richardson M, Sibley
921 TR, Truong M, Wolf CR, Nickerson DA, Rieder MJ, Englund JA, Hadfield J, Hodcroft EB,
922 Huddleston J, Moncla LH, Müller NF, Neher RA, Deng X, Gu W, Federman S, Chiu C, Duchin J,
923 Gautam R, Melly G, Hiatt B, Dykema P, Lindquist S, Queen K, Tao Y, Uehara A, Tong S, et al.
924 2020. Cryptic transmission of SARS-CoV-2 in Washington State. *medRxiv*
925 doi:10.1101/2020.04.02.20051417:2020.04.02.20051417.
- 926 11. Gonzalez-Reiche AS, Hernandez MM, Sullivan MJ, Ciferri B, Alshammary H, Obla A, Fabre S,
927 Kleiner G, Polanco J, Khan Z, Alburquerque B, van de Guchte A, Dutta J, Francoeur N, Melo BS,
928 Oussenko I, Deikus G, Soto J, Sridhar SH, Wang Y-C, Twyman K, Kasarskis A, Altman DR,
929 Smith M, Sebra R, Aberg J, Krammer F, García-Sastre A, Lukiszka M, Patel G, Paniz-Mondolfi A,
930 Gitman M, Sordillo EM, Simon V, van Bakel H. 2020. Introductions and early spread of SARS-
931 CoV-2 in the New York City area. *Science* 369:297-301.
- 932 12. Health N. 2020. COVID-19 Data. <https://www1.nyc.gov/site/doh/covid/covid-19-data.page>.
933 Accessed August 19.
- 934 13. County K. 2020. Daily COVID-19 outbreak summary.
935 <https://www.kingcounty.gov/depts/health/covid-19/data/daily-summary.aspx>. Accessed August
936 18.
- 937 14. Cline M, Emerson M, Bratter J, Howell J, Jeanty P. 2012. Houston Region Grows More
938 Racially/Ethnically Diverse, With Small Declines in Segregation. A Joint Report Analyzing Census
939 Data from 1990, 2000, and 2010.
- 940 15. Emerson M, Bratter J, Howell J, Jeanty P, Cline M. 2012. Houston Region Grows More
941 Racially/Ethnically Diverse, With Small Declines in Segregation. A Joint Report Analyzing

- 942 Census Data from 1990, 2000, and 2010. Kinder Institute for Urban Research & the Hobby
943 Center for the Study of Texas,
944 16. Services THaH. 2020. Texas Health and Human Services. <https://hhs.texas.gov/>. Accessed
945 August 18.
946 17. Vahidy FS, Drews AL, Masud FN, Schwartz RL, Askary BB, Boom ML, Phillips RA. 2020.
947 Characteristics and Outcomes of COVID-19 Patients During Initial Peak and Resurgence in the
948 Houston Metropolitan Area. *Jama* doi:10.1001/jama.2020.15301.
949 18. Diehl WE, Lin AE, Grubaugh ND, Carvalho LM, Kim K, Kyaw PP, McCauley SM, Donnard E,
950 Kucukural A, McDonel P, Schaffner SF, Garber M, Rambaut A, Andersen KG, Sabeti PC, Luban
951 J. 2016. Ebola Virus Glycoprotein with Increased Infectivity Dominated the 2013-2016 Epidemic.
952 *Cell* 167:1088-1098.e6.
953 19. Urbanowicz RA, McClure CP, Sakuntabhai A, Sall AA, Kobinger G, Müller MA, Holmes EC,
954 Rey FA, Simon-Loriere E, Ball JK. 2016. Human Adaptation of Ebola Virus during the West
955 African Outbreak. *Cell* 167:1079-1087.e5.
956 20. Dietzel E, Schudt G, Krähling V, Matrosovich M, Becker S. 2017. Functional Characterization of
957 Adaptive Mutations during the West African Ebola Virus Outbreak. *J Virol* 91.
958 21. Kachroo P, Eraso JM, Beres SB, Olsen RJ, Zhu L, Nasser W, Bernard PE, Cantu CC, Saavedra
959 MO, Arredondo MJ, Strope B, Do H, Kumaraswami M, Vuopio J, Grondahl-Yli-Hannuksela K,
960 Kristinsson KG, Gottfredsson M, Pesonen M, Pensar J, Davenport ER, Clark AG, Corander J,
961 Caugant DA, Gaini S, Magnussen MD, Kubik SL, Nguyen HAT, Long SW, Porter AR, DeLeo
962 FR, Musser JM. 2019. Integrated analysis of population genomics, transcriptomics and virulence
963 provides novel insights into *Streptococcus pyogenes* pathogenesis. *Nat Genet* 51:548-559.
964 22. Gao Y, Yan L, Huang Y, Liu F, Zhao Y, Cao L, Wang T, Sun Q, Ming Z, Zhang L, Ge J, Zheng
965 L, Zhang Y, Wang H, Zhu Y, Zhu C, Hu T, Hua T, Zhang B, Yang X, Li J, Yang H, Liu Z, Xu W,
966 Guddat LW, Wang Q, Lou Z, Rao Z. 2020. Structure of the RNA-dependent RNA polymerase
967 from COVID-19 virus. *Science* doi:10.1126/science.abb7498:eabb7498.
968 23. Yin W, Mao C, Luan X, Shen D-D, Shen Q, Su H, Wang X, Zhou F, Zhao W, Gao M, Chang S,
969 Xie Y-C, Tian G, Jiang H-W, Tao S-C, Shen J, Jiang Y, Jiang H, Xu Y, Zhang S, Zhang Y, Xu
970 HE. 2020. Structural basis for inhibition of the RNA-dependent RNA polymerase from SARS-
971 CoV-2 by remdesivir. *Science* 368:1499-1504.
972 24. Shannon A, Le NT, Selisko B, Eydoux C, Alvarez K, Guillemot JC, Decroly E, Peersen O, Ferron
973 F, Canard B. 2020. Remdesivir and SARS-CoV-2: Structural requirements at both nsp12 RdRp
974 and nsp14 Exonuclease active-sites. *Antiviral Res* 178:104793.
975 25. Gordon CJ, Tchesnokov EP, Woolner E, Perry JK, Feng JY, Porter DP, Gotte M. 2020.
976 Remdesivir is a direct-acting antiviral that inhibits RNA-dependent RNA polymerase from severe
977 acute respiratory syndrome coronavirus 2 with high potency. *J Biol Chem*
978 doi:10.1074/jbc.RA120.013679.
979 26. Agostini ML, Andres EL, Sims AC, Graham RL, Sheahan TP, Lu X, Smith EC, Case JB, Feng
980 JY, Jordan R, Ray AS, Cihlar T, Siegel D, Mackman RL, Clarke MO, Baric RS, Denison MR.
981 2018. Coronavirus Susceptibility to the Antiviral Remdesivir (GS-5734) Is Mediated by the Viral
982 Polymerase and the Proofreading Exoribonuclease. *mBio* 9.
983 27. de Wit E, Feldmann F, Cronin J, Jordan R, Okumura A, Thomas T, Scott D, Cihlar T, Feldmann
984 H. 2020. Prophylactic and therapeutic remdesivir (GS-5734) treatment in the rhesus macaque
985 model of MERS-CoV infection. *Proc Natl Acad Sci U S A* 117:6771-6776.
986 28. Williamson BN, Feldmann F, Schwarz B, Meade-White K, Porter DP, Schulz J, van Doremalen
987 N, Leighton I, Yinda CK, Pérez-Pérez L, Okumura A, Lovaglio J, Hanley PW, Saturday G, Bosio
988 CM, Anzick S, Barbian K, Cihlar T, Martens C, Scott DP, Munster VJ, de Wit E. 2020. Clinical
989 benefit of remdesivir in rhesus macaques infected with SARS-CoV-2. *Nature* 585:273-276.
990 29. Grein J, Ohmagari N, Shin D, Diaz G, Asperges E, Castagna A, Feldt T, Green G, Green ML,
991 Lescure FX, Nicastri E, Oda R, Yo K, Quiros-Roldan E, Studemeister A, Redinski J, Ahmed S,
992 Bennett J, Chelliah D, Chen D, Chihara S, Cohen SH, Cunningham J, D'Arminio Monforte A,
993 Ismail S, Kato H, Lapadula G, L'Her E, Maeno T, Majumder S, Massari M, Mora-Rillo M, Mutoh
994 Y, Nguyen D, Verweij E, Zoufaly A, Osinusi AO, DeZure A, Zhao Y, Zhong L, Chokkalingam A,
995 Elboudwarej E, Telep L, Timbs L, Henne I, Sellers S, Cao H, Tan SK, Winterbourne L, Desai P,
996 et al. 2020. Compassionate Use of Remdesivir for Patients with Severe Covid-19. *N Engl J Med*
997 doi:10.1056/NEJMoa2007016.

- 998 30. Goldman JD, Lye DCB, Hui DS, Marks KM, Bruno R, Montejano R, Spinner CD, Galli M, Ahn
999 MY, Nahass RG, Chen YS, SenGupta D, Hyland RH, Osinusi AO, Cao H, Blair C, Wei X,
1000 Gaggar A, Brainard DM, Towner WJ, Muñoz J, Mullane KM, Marty FM, Tashima KT, Diaz G,
1001 Subramanian A. 2020. Remdesivir for 5 or 10 Days in Patients with Severe Covid-19. N Engl J
1002 Med doi:10.1056/NEJMoa2015301.
1003 31. Beigel JH, Tomashek KM, Dodd LE, Mehta AK, Zingman BS, Kalil AC, Hohmann E, Chu HY,
1004 Luetkemeyer A, Kline S, Lopez de Castilla D, Finberg RW, Dierberg K, Tapson V, Hsieh L,
1005 Patterson TF, Paredes R, Sweeney DA, Short WR, Touloumi G, Lye DC, Ohmagari N, Oh MD,
1006 Ruiz-Palacios GM, Benfield T, Fätkenheuer G, Kortepeter MG, Atmar RL, Creech CB, Lundgren
1007 J, Babiker AG, Pett S, Neaton JD, Burgess TH, Bonnett T, Green M, Makowski M, Osinusi A,
1008 Nayak S, Lane HC. 2020. Remdesivir for the Treatment of Covid-19 - Preliminary Report. N Engl
1009 J Med doi:10.1056/NEJMoa2007764.
1010 32. Spinner CD, Gottlieb RL, Criner GJ, Arribas López JR, Cattelan AM, Soriano Viladomiu A,
1011 Ogbuagu O, Malhotra P, Mullane KM, Castagna A, Chai LYA, Roestenberg M, Tsang OTY,
1012 Bernasconi E, Le Turnier P, Chang SC, SenGupta D, Hyland RH, Osinusi AO, Cao H, Blair C,
1013 Wang H, Gaggar A, Brainard DM, McPhail MJ, Bhagani S, Ahn MY, Sanyal AJ, Huhn G, Marty
1014 FM. 2020. Effect of Remdesivir vs Standard Care on Clinical Status at 11 Days in Patients With
1015 Moderate COVID-19: A Randomized Clinical Trial. Jama doi:10.1001/jama.2020.16349.
1016 33. Olander SA, Perez KK, Go AS, Balani B, Price-Haywood EG, Shah NS, Wang S, Walunas TL,
1017 Swaminathan S, Slim J, Chin B, De Wit S, Ali SM, Soriano Viladomiu A, Robinson P, Gottlieb
1018 RL, Tsang TYO, Lee IH, Haubrich RH, Chokkalingam AP, Lin L, Zhong L, Bekele BN, Mera-
1019 Giler R, Gallant J, Smith LE, Osinusi AO, Brainard DM, Hu H, Phulpin C, Edgar H, Diaz-Cuervo
1020 H, Bernardino JI. 2020. Remdesivir for Severe COVID-19 versus a Cohort Receiving Standard of
1021 Care. Clin Infect Dis doi:10.1093/cid/ciaa1041.
1022 34. (CNCB) CNCfB. 2020. 2019 Novel Coronavirus Resource (2019nCoVR).
1023 <https://bigd.big.ac.cn/ncov/about?lang=en>. Accessed August 19.
1024 35. Wrapp D, Wang N, Corbett KS, Goldsmith JA, Hsieh C-L, Abiona O, Graham BS, McLellan JS.
1025 2020. Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. Science
1026 367:1260-1263.
1027 36. Walls AC, Park YJ, Tortorici MA, Wall A, McGuire AT, Veesler D. 2020. Structure, Function,
1028 and Antigenicity of the SARS-CoV-2 Spike Glycoprotein. Cell 181:281-292.e6.
1029 37. Wang Q, Zhang Y, Wu L, Niu S, Song C, Zhang Z, Lu G, Qiao C, Hu Y, Yuen KY, Wang Q,
1030 Zhou H, Yan J, Qi J. 2020. Structural and Functional Basis of SARS-CoV-2 Entry by Using
1031 Human ACE2. Cell doi:10.1016/j.cell.2020.03.045.
1032 38. Jackson LA, Anderson EJ, Roushpal NG, Roberts PC, Makhene M, Coler RN, McCullough MP,
1033 Chappell JD, Denison MR, Stevens LJ, Pruijssers AJ, McDermott A, Flach B, Doria-Rose NA,
1034 Corbett KS, Morabito KM, O'Dell S, Schmidt SD, Swanson PA, 2nd, Padilla M, Mascola JR,
1035 Neuzil KM, Bennett H, Sun W, Peters E, Makowski M, Albert J, Cross K, Buchanan W, Pikaart-
1036 Tautges R, Ledgerwood JE, Graham BS, Beigel JH. 2020. An mRNA Vaccine against SARS-
1037 CoV-2 - Preliminary Report. N Engl J Med doi:10.1056/NEJMoa2022483.
1038 39. Folegatti PM, Ewer KJ, Aley PK, Angus B, Becker S, Belij-Rammerstorfer S, Bellamy D, Bibi S,
1039 Bittaye M, Clutterbuck EA, Dold C, Faust SN, Finn A, Flaxman AL, Hallis B, Heath P, Jenkin D,
1040 Lazarus R, Makinson R, Minassian AM, Pollock KM, Ramasamy M, Robinson H, Snape M,
1041 Tarrant R, Voysey M, Green C, Douglas AD, Hill AVS, Lambe T, Gilbert SC, Pollard AJ. 2020.
1042 Safety and immunogenicity of the ChAdOx1 nCoV-19 vaccine against SARS-CoV-2: a
1043 preliminary report of a phase 1/2, single-blind, randomised controlled trial. Lancet 396:467-478.
1044 40. Zhu FC, Guan XH, Li YH, Huang JY, Jiang T, Hou LH, Li JX, Yang BF, Wang L, Wang WJ, Wu
1045 SP, Wang Z, Wu XH, Xu JJ, Zhang Z, Jia SY, Wang BS, Hu Y, Liu JJ, Zhang J, Qian XA, Li Q,
1046 Pan HX, Jiang HD, Deng P, Gou JB, Wang XW, Wang XH, Chen W. 2020. Immunogenicity and
1047 safety of a recombinant adenovirus type-5-vectorized COVID-19 vaccine in healthy adults aged 18
1048 years or older: a randomised, double-blind, placebo-controlled, phase 2 trial. Lancet 396:479-488.
1049 41. Brouwer PJM, Caniels TG, van der Straten K, Snitselaar JL, Aldon Y, Bangaru S, Torres JL, Okba
1050 NMA, Claireaux M, Kerster G, Bentlage AEH, van Haaren MM, Guerra D, Burger JA, Schermer
1051 EE, Verheul KD, van der Velde N, van der Kooi A, van Schooten J, van Breemen MJ, Bijl TPL,
1052 Sliepen K, Aartse A, Derkking R, Bontjer I, Kootstra NA, Wiersinga WJ, Vidarsson G, Haagmans

- 1053 BL, Ward AB, de Bree GJ, Sanders RW, van Gils MJ. 2020. Potent neutralizing antibodies from
1054 COVID-19 patients define multiple targets of vulnerability. *Science* 369:643-650.
1055 42. Chi X, Yan R, Zhang J, Zhang G, Zhang Y, Hao M, Zhang Z, Fan P, Dong Y, Yang Y, Chen Z,
1056 Guo Y, Zhang J, Li Y, Song X, Chen Y, Xia L, Fu L, Hou L, Xu J, Yu C, Li J, Zhou Q, Chen W.
1057 2020. A neutralizing human antibody binds to the N-terminal domain of the Spike protein of
1058 SARS-CoV-2. *Science* 369:650-655.
1059 43. Wec AZ, Wrapp D, Herbert AS, Maurer DP, Haslwanter D, Sakharkar M, Jangra RK, Dieterle
1060 ME, Lilov A, Huang D, Tse LV, Johnson NV, Hsieh C-L, Wang N, Nett JH, Champney E,
1061 Burnina I, Brown M, Lin S, Sinclair M, Johnson C, Pudi S, Bortz R, Wirchnianski AS,
1062 Laudermilch E, Florez C, Fels JM, O'Brien CM, Graham BS, Nemazee D, Burton DR, Baric RS,
1063 Voss JE, Chandran K, Dye JM, McLellan JS, Walker LM. 2020. Broad neutralization of SARS-
1064 related viruses by human monoclonal antibodies. *Science* 369:731-736.
1065 44. Zost SJ, Gilchuk P, Case JB, Binshtain E, Chen RE, Nkolola JP, Schäfer A, Reidy JX, Trivette A,
1066 Nargi RS, Sutton RE, Suryadevara N, Martinez DR, Williamson LE, Chen EC, Jones T, Day S,
1067 Myers L, Hassan AO, Kafai NM, Winkler ES, Fox JM, Shrihari S, Mueller BK, Meiler J,
1068 Chandrashekhar A, Mercado NB, Steinhardt JJ, Ren K, Loo YM, Kallewaard NL, McCune BT,
1069 Keeler SP, Holtzman MJ, Barouch DH, Gralinski LE, Baric RS, Thackray LB, Diamond MS,
1070 Carnahan RH, Crowe JE, Jr. 2020. Potently neutralizing and protective human antibodies against
1071 SARS-CoV-2. *Nature* 584:443-449.
1072 45. Greaney AJ, Starr TN, Gilchuk P, Zost SJ, Binshtain E, Loes AN, Hilton SK, Huddleston J, Eguia
1073 R, Crawford KH, Dingens AS, Nargi RS, Sutton RE, Suryadevara N, Rothlauf PW, Liu Z, Whelan
1074 SP, Carnahan RH, Crowe JE, Bloom JD. 2020. Complete mapping of mutations to the SARS-
1075 CoV-2 spike receptor-binding domain that escape antibody recognition. *bioRxiv*
1076 doi:10.1101/2020.09.10.292078:2020.09.10.292078.
1077 46. Baum A, Copin R, Ajithdoss D, Zhou A, Lanza K, Negron N, Ni M, Wei Y, Atwal GS, Oyejide
1078 A, Goez-Gazi Y, Dutton J, Clemons E, Staples HM, Bartley C, Klaffke B, Alfson K, Gazi M,
1079 Gonzales O, Dick E, Carrion R, Pessant L, Porto M, Cook A, Brown R, Ali V, Greenhouse J,
1080 Taylor T, Andersen H, Lewis MG, Stahl N, Murphy AJ, Yancopoulos GD, Kyratsous CA. 2020.
1081 REGN-COV2 antibody cocktail prevents and treats SARS-CoV-2 infection in rhesus macaques
1082 and hamsters. *bioRxiv* doi:10.1101/2020.08.02.233320:2020.08.02.233320.
1083 47. Baum A, Fulton BO, Wloga E, Copin R, Pascal KE, Russo V, Giordano S, Lanza K, Negron N, Ni
1084 M, Wei Y, Atwal GS, Murphy AJ, Stahl N, Yancopoulos GD, Kyratsous CA. 2020. Antibody
1085 cocktail to SARS-CoV-2 spike protein prevents rapid mutational escape seen with individual
1086 antibodies. *Science* 369:1014-1018.
1087 48. Barnes CO, West AP, Jr., Huey-Tubman KE, Hoffmann MAG, Sharaf NG, Hoffman PR, Koranda
1088 N, Gristick HB, Gaebler C, Muecksch F, Lorenzi JCC, Finkin S, Hägglöf T, Hurley A, Millard
1089 KG, Weisblum Y, Schmidt F, Hatzioannou T, Bieniasz PD, Caskey M, Robbiani DF,
1090 Nussenzweig MC, Bjorkman PJ. 2020. Structures of Human Antibodies Bound to SARS-CoV-2
1091 Spike Reveal Common Epitopes and Recurrent Features of Antibodies. *Cell* 182:828-842.e16.
1092 49. Alsoussi WB, Turner JS, Case JB, Zhao H, Schmitz AJ, Zhou JQ, Chen RE, Lei T, Rizk AA,
1093 McIntire KM, Winkler ES, Fox JM, Kafai NM, Thackray LB, Hassan AO, Amanat F, Krammer F,
1094 Watson CT, Kleinstein SH, Fremont DH, Diamond MS, Ellebedy AH. 2020. A Potently
1095 Neutralizing Antibody Protects Mice against SARS-CoV-2 Infection. *J Immunol* 205:915-922.
1096 50. Salazar E, Kuchipudi SV, Christensen PA, Eagar T, Yi X, Zhao P, Jin Z, Long SW, Olsen RJ,
1097 Chen J, Castillo B, Leveque C, Towers D, Lavinder JJ, Gollihar J, Cardona JA, Ippolito GC,
1098 Nissly RH, Bird I, Greenawalt D, Rossi RM, Gontu A, Srinivasan S, Poojary I, Cattadori IM,
1099 Hudson P, Josley NM, Prugar L, Huie KE, Herbert AS, Bernard DW, Dye JM, Kapur V, Musser
1100 JM. 2020. Convalescent plasma anti-SARS-CoV-2 spike protein ectodomain and receptor binding
1101 domain IgG correlate with virus neutralization. *The Journal of Clinical Investigation*
1102 doi:10.1172/JCI141206.
1103 51. Salazar E, Christensen PA, Graviss EA, Nguyen DT, Castillo B, Chen J, Lopez BV, Eagar TN, Yi
1104 X, Zhao P, Rogers J, Shehabeldin A, Joseph D, Leveque C, Olsen RJ, Bernard DW, Gollihar J,
1105 Musser JM. 2020. Treatment of COVID-19 Patients with Convalescent Plasma Reveals a Signal
1106 of Significantly Decreased Mortality. *Am J Pathol* doi:10.1016/j.ajpath.2020.08.001.
1107 52. Starr TN, Greaney AJ, Hilton SK, Ellis D, Crawford KHD, Dingens AS, Navarro MJ, Bowen JE,
1108 Tortorici MA, Walls AC, King NP, Veesler D, Bloom JD. 2020. Deep Mutational Scanning of

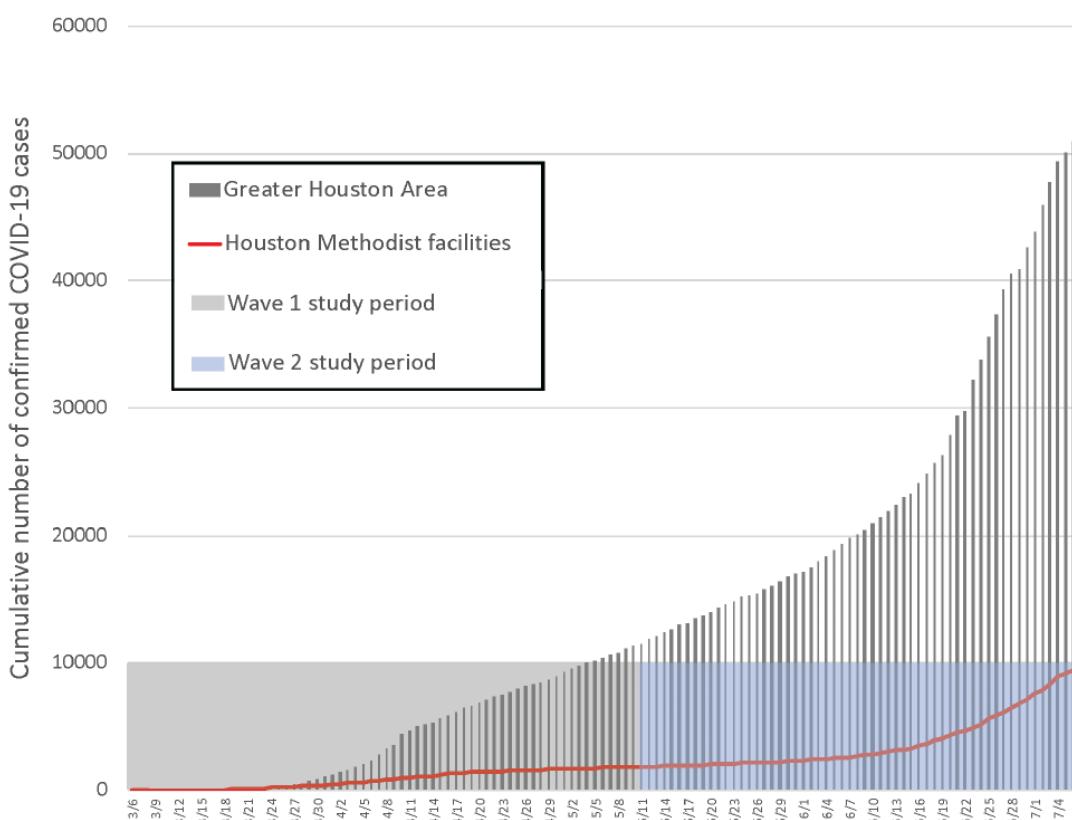
- 1109 SARS-CoV-2 Receptor Binding Domain Reveals Constraints on Folding and ACE2 Binding. *Cell*
1110 doi:10.1016/j.cell.2020.08.012.
1111 53. Steffen TL, Stone ET, Hassert M, Gearling E, Grimberg BT, Espino AM, Pantoja P, Climent C,
1112 Hoft DF, George SL, Sariol CA, Pinto AK, Brien JD. 2020. The receptor binding domain of
1113 SARS-CoV-2 spike is the key target of neutralizing antibody in human polyclonal sera. *bioRxiv*
1114 doi:10.1101/2020.08.21.261727:2020.08.21.261727.
1115 54. Corbett KS, Flynn B, Foulds KE, Francica JR, Boyoglu-Barnum S, Werner AP, Flach B,
1116 O'Connell S, Bock KW, Minai M, Nagata BM, Andersen H, Martinez DR, Noe AT, Douek N,
1117 Donaldson MM, Nji NN, Alvarado GS, Edwards DK, Flebbe DR, Lamb E, Doria-Rose NA, Lin
1118 BC, Louder MK, O'Dell S, Schmidt SD, Phung E, Chang LA, Yap C, Todd J-PM, Pessant L, Van
1119 Ry A, Browne S, Greenhouse J, Putman-Taylor T, Strasbaugh A, Campbell T-A, Cook A, Dodson
1120 A, Steingrebe K, Shi W, Zhang Y, Abiona OM, Wang L, Pegu A, Yang ES, Leung K, Zhou T,
1121 Teng I-T, Widge A, et al. 2020. Evaluation of the mRNA-1273 Vaccine against SARS-CoV-2 in
1122 Nonhuman Primates. *New England Journal of Medicine* doi:10.1056/NEJMoa2024671.
1123 55. van Doremalen N, Lambe T, Spencer A, Belij-Rammerstorfer S, Purushotham JN, Port JR,
1124 Avanzato VA, Bushmaker T, Flaxman A, Ulaszewska M, Feldmann F, Allen ER, Sharpe H,
1125 Schulz J, Holbrook M, Okumura A, Meade-White K, Pérez-Pérez L, Edwards NJ, Wright D,
1126 Bissett C, Gilbride C, Williamson BN, Rosenke R, Long D, Ishwarbhai A, Kailath R, Rose L,
1127 Morris S, Powers C, Lovaglio J, Hanley PW, Scott D, Saturday G, de Wit E, Gilbert SC, Munster
1128 VJ. 2020. ChAdOx1 nCoV-19 vaccine prevents SARS-CoV-2 pneumonia in rhesus macaques.
1129 Nature doi:10.1038/s41586-020-2608-y.
1130 56. Wang C, Li W, Drabek D, Okba NMA, van Haperen R, Osterhaus A, van Kuppeveld FJM,
1131 Haagmans BL, Grosveld F, Bosch BJ. 2020. A human monoclonal antibody blocking SARS-CoV-
1132 2 infection. *Nat Commun* 11:2251.
1133 57. Ju B, Zhang Q, Ge J, Wang R, Sun J, Ge X, Yu J, Shan S, Zhou B, Song S, Tang X, Yu J, Lan J,
1134 Yuan J, Wang H, Zhao J, Zhang S, Wang Y, Shi X, Liu L, Zhao J, Wang X, Zhang Z, Zhang L.
1135 2020. Human neutralizing antibodies elicited by SARS-CoV-2 infection. *Nature* 584:115-119.
1136 58. Liu L, Wang P, Nair MS, Yu J, Rapp M, Wang Q, Luo Y, Chan JF, Sahi V, Figueroa A, Guo XV,
1137 Cerutti G, Bimela J, Gorman J, Zhou T, Chen Z, Yuen KY, Kwong PD, Sodroski JG, Yin MT,
1138 Sheng Z, Huang Y, Shapiro L, Ho DD. 2020. Potent neutralizing antibodies against multiple
1139 epitopes on SARS-CoV-2 spike. *Nature* 584:450-456.
1140 59. Rogers TF, Zhao F, Huang D, Beutler N, Burns A, He W-t, Limbo O, Smith C, Song G, Woehl J,
1141 Yang L, Abbott RK, Callaghan S, Garcia E, Hurtado J, Parren M, Peng L, Ramirez S, Ricketts J,
1142 Ricciardi MJ, Rawlings SA, Wu NC, Yuan M, Smith DM, Nemazee D, Teijaro JR, Voss JE,
1143 Wilson IA, Andrabi R, Briney B, Landais E, Sok D, Jardine JG, Burton DR. 2020. Isolation of
1144 potent SARS-CoV-2 neutralizing antibodies and protection from disease in a small animal model.
1145 Science 369:956-963.
1146 60. Hassan AO, Case JB, Winkler ES, Thackray LB, Kafai NM, Bailey AL, McCune BT, Fox JM,
1147 Chen RE, Alsoussi WB, Turner JS, Schmitz AJ, Lei T, Shrihari S, Keeler SP, Fremont DH, Greco
1148 S, McCray PB, Jr., Perlman S, Holtzman MJ, Ellebedy AH, Diamond MS. 2020. A SARS-CoV-2
1149 Infection Model in Mice Demonstrates Protection by Neutralizing Antibodies. *Cell* 182:744-
1150 753.e4.
1151 61. Chandrashekhar A, Liu J, Martinot AJ, McMahan K, Mercado NB, Peter L, Tostanoski LH, Yu J,
1152 Maliga Z, Nekorchuk M, Busman-Sahay K, Terry M, Wrijil LM, Ducat S, Martinez DR, Atyeo C,
1153 Fischinger S, Burke JS, Stein MD, Pessant L, Van Ry A, Greenhouse J, Taylor T, Blade K, Cook
1154 A, Finneyfrock B, Brown R, Teow E, Velasco J, Zahn R, Wegmann F, Abbink P, Bondzie EA,
1155 Dagotto G, Gebre MS, He X, Jacob-Dolan C, Kordana N, Li Z, Lifton MA, Mahrokhan SH,
1156 Maxfield LF, Nityanandam R, Nkolola JP, Schmidt AG, Miller AD, Baric RS, Alter G, Sorger
1157 PK, Estes JD, et al. 2020. SARS-CoV-2 infection protects against rechallenge in rhesus macaques.
1158 Science 369:812-817.
1159 62. Mercado NB, Zahn R, Wegmann F, Loos C, Chandrashekhar A, Yu J, Liu J, Peter L, McMahan K,
1160 Tostanoski LH, He X, Martinez DR, Rutten L, Bos R, van Manen D, Vellinga J, Custers J,
1161 Langedijk JP, Kwaks T, Bakkers MJG, Zuidgeest D, Huber SKR, Atyeo C, Fischinger S, Burke
1162 JS, Feldman J, Hauser BM, Caradonna TM, Bondzie EA, Dagotto G, Gebre MS, Hoffman E,
1163 Jacob-Dolan C, Kirilova M, Li Z, Lin Z, Mahrokhan SH, Maxfield LF, Nampanya F,
1164 Nityanandam R, Nkolola JP, Patel S, Ventura JD, Verrington K, Wan H, Pessant L, Ry AV,

- 1165 Blade K, Strasbaugh A, Cabus M, et al. 2020. Single-shot Ad26 vaccine protects against SARS-CoV-2 in rhesus macaques. *Nature* doi:10.1038/s41586-020-2607-z.
1166
1167 63. Yuan M, Wu NC, Zhu X, Lee CD, So RTY, Lv H, Mok CKP, Wilson IA. 2020. A highly
1168 conserved cryptic epitope in the receptor binding domains of SARS-CoV-2 and SARS-CoV.
1169 *Science* 368:630-633.
1170 64. Hsieh C-L, Goldsmith JA, Schaub JM, DiVenere AM, Kuo H-C, Javanmardi K, Le KC, Wrapp D,
1171 Lee AG, Liu Y, Chou C-W, Byrne PO, Hjorth CK, Johnson NV, Ludes-Meyers J, Nguyen AW,
1172 Park J, Wang N, Amengor D, Lavinder JJ, Ippolito GC, Maynard JA, Finkelstein IJ, McLellan JS.
1173 2020. Structure-based design of prefusion-stabilized SARS-CoV-2 spikes. *Science* 369:1501-
1174 1505.
1175 65. Woo H, Park SJ, Choi YK, Park T, Tanveer M, Cao Y, Kern NR, Lee J, Yeom MS, Croll TI, Seok
1176 C, Im W. 2020. Developing a Fully Glycosylated Full-Length SARS-CoV-2 Spike Protein Model
1177 in a Viral Membrane. *J Phys Chem B* 124:7128-7137.
1178 66. Korber B, Fischer WM, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Hengartner N, Giorgi
1179 EE, Bhattacharya T, Foley B, Hastie KM, Parker MD, Partridge DG, Evans CM, Freeman TM, de
1180 Silva TI, McDanal C, Perez LG, Tang H, Moon-Walker A, Whelan SP, LaBranche CC, Saphire
1181 EO, Montefiori DC. 2020. Tracking Changes in SARS-CoV-2 Spike: Evidence that D614G
1182 Increases Infectivity of the COVID-19 Virus. *Cell* 182:812-827.e19.
1183 67. Hu J, He C-L, Gao Q-Z, Zhang G-J, Cao X-X, Long Q-X, Deng H-J, Huang L-Y, Chen J, Wang
1184 K, Tang N, Huang A-L. 2020. D614G mutation of SARS-CoV-2 spike protein enhances viral
1185 infectivity. *bioRxiv* doi:10.1101/2020.06.20.161323:2020.06.20.161323.
1186 68. Lorenzo-Redondo R, Nam HH, Roberts SC, Simons LM, Jennings LJ, Qi C, Achenbach CJ,
1187 Hauser AR, Ison MG, Hultquist JF, Ozer EA. 2020. A Unique Clade of SARS-CoV-2 Viruses is
1188 Associated with Lower Viral Loads in Patient Upper Airways. *medRxiv*
1189 doi:10.1101/2020.05.19.20107144:2020.05.19.20107144.
1190 69. Cassia Wagner PR, Chris D. Frazar, Jover Lee, Nicola F. Müller, Louise H. Moncla, James
1191 Hadfield, Emma B. Hodcroft, Benjamin Pelle, Matthew Richardson, Caitlin Behrens, Meei-Li
1192 Huang, Patrick Mathias, Gregory Pepper, Lasata Shrestha, Hong Xie, Amin Addetia, Truong
1193 Nguyen, Victoria M Rachleff, Romesh Gautam, Geoff Melly, Brian Hiatt, Philip Dykema,
1194 Amanda Adler, Elisabeth Brandstetter, Peter D. Han, Kairsten Fay, Misja Ilcisin, Kirsten
1195 Lacombe, Thomas R. Sibley, Melissa Truong, Caitlin R. Wolf, Karen Cowgill, Stephanie Schrag,
1196 Jeff Duchin, Michael Boeckh, Janet A. Englund, Michael Famulari, Barry R. Lutz, Mark J.
1197 Rieder, Matthew Thompson, Richard A. Neher, Geoffrey S. Baird, Lea M. Starita, Helen Y. Chu,
1198 Jay Shendure, Scott Lindquist, Deborah A. Nickerson, Alexander L. Greninger, Keith R. Jerome,
1199 Trevor Bedford. 2020. Comparing viral load and clinical outcomes in Washington State across
1200 D614G substitution in spike protein of SARS-CoV-2. <https://github.com/blab/ncov-wa-d614g>.
1201 Accessed September 8.
1202 70. Volz EM, Hill V, McCrone JT, Price A, Jorgensen D, O'Toole A, Southgate JA, Johnson R,
1203 Jackson B, Nascimento FF, Rey SM, Nicholls SM, Colquhoun RM, da Silva Filipe A, Shepherd
1204 JG, Pascall DJ, Shah R, Jesudason N, Li K, Jarrett R, Pacchiarini N, Bull M, Geidelberg L,
1205 Siveroni I, Goodfellow IG, Loman NJ, Pybus O, Robertson DL, Thomson EC, Rambaut A,
1206 Connor TR. 2020. Evaluating the effects of SARS-CoV-2 Spike mutation D614G on
1207 transmissibility and pathogenicity. *medRxiv*
1208 doi:10.1101/2020.07.31.20166082:2020.07.31.20166082.
1209 71. Lv Z, Deng Y-Q, Ye Q, Cao L, Sun C-Y, Fan C, Huang W, Sun S, Sun Y, Zhu L, Chen Q, Wang
1210 N, Nie J, Cui Z, Zhu D, Shaw N, Li X-F, Li Q, Xie L, Wang Y, Rao Z, Qin C-F, Wang X. 2020.
1211 Structural basis for neutralization of SARS-CoV-2 and SARS-CoV by a potent therapeutic
1212 antibody. *Science* 369:1505-1509.
1213 72. Yurkovetskiy L, Wang X, Pascal KE, Tomkins-Tinch C, Nyalile T, Wang Y, Baum A, Diehl WE,
1214 Dauphin A, Carbone C, Veinotte K, Egri SB, Schaffner SF, Lemieux JE, Munro J, Rafique A,
1215 Barve A, Sabeti PC, Kyratsous CA, Dudkina N, Shen K, Luban J. 2020. Structural and Functional
1216 Analysis of the D614G SARS-CoV-2 Spike Protein Variant. *bioRxiv*
1217 doi:10.1101/2020.07.04.187757:2020.07.04.187757.
1218 73. Zhang L, Jackson CB, Mou H, Ojha A, Rangarajan ES, Izard T, Farzan M, Choe H. 2020. The
1219 D614G mutation in the SARS-CoV-2 spike protein reduces S1 shedding and increases infectivity.
1220 *bioRxiv* doi:10.1101/2020.06.12.148726:2020.06.12.148726.

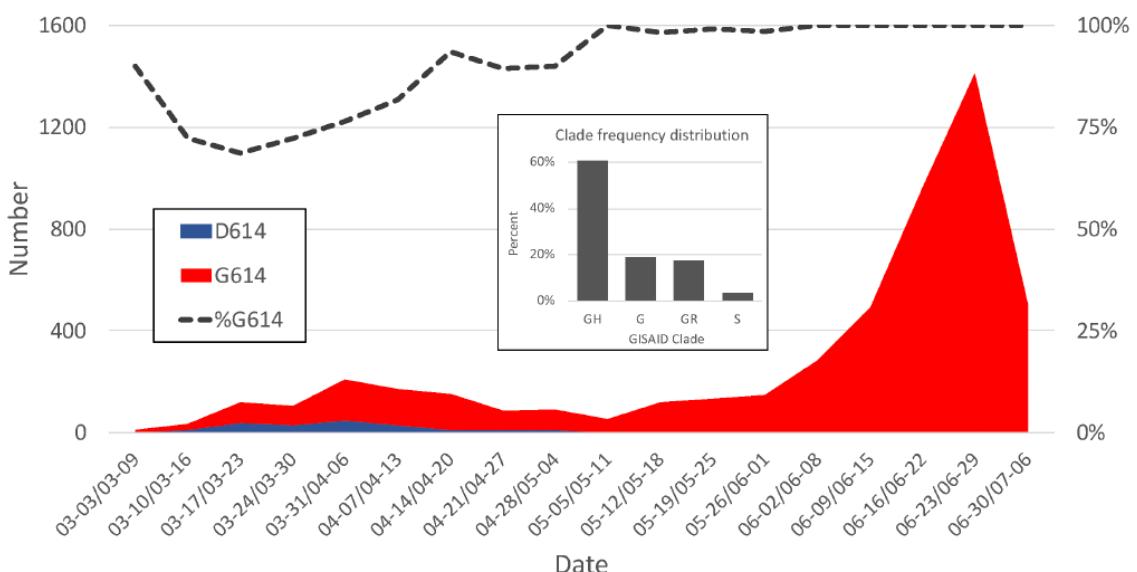
- 1221 74. Huo J, Zhao Y, Ren J, Zhou D, Duyvesteyn HME, Ginn HM, Carrique L, Malinauskas T, Ruza
1222 RR, Shah PNM, Tan TK, Rijal P, Coombes N, Bewley KR, Tree JA, Radecke J, Paterson NG,
1223 Supasa P, Mongkolsapaya J, Scream GR, Carroll M, Townsend A, Fry EE, Owens RJ, Stuart DI.
1224 2020. Neutralization of SARS-CoV-2 by Destruction of the Prefusion Spike. *Cell Host Microbe*
1225 doi:10.1016/j.chom.2020.06.010.
- 1226 75. Long SW, Olsen RJ, Christensen PA, Bernard DW, Davis JR, Shukla M, Nguyen M, Ojeda
1227 Saavedra M, Cantu CC, Yerramilli P, Pruitt L, Subedi S, Hendrickson H, Eskandari G,
1228 Kumaraswami M, McLellan JS, Musser JM. 2020. Molecular Architecture of Early Dissemination
1229 and Evolution of the SARS-CoV-2 Virus in Metropolitan Houston, Texas. *bioRxiv*
1230 doi:10.1101/2020.05.01.072652:2020.05.01.072652.
- 1231 76. Plante JA, Liu Y, Liu J, Xia H, Johnson BA, Lokugamage KG, Zhang X, Muruato AE, Zou J,
1232 Fontes-Garfias CR, Mirchandani D, Scharton D, Bilello JP, Ku Z, An Z, Kalveram B, Freiberg
1233 AN, Menachery VD, Xie X, Plante KS, Weaver SC, Shi P-Y. 2020. Spike mutation D614G alters
1234 SARS-CoV-2 fitness and neutralization susceptibility. *bioRxiv*
1235 doi:10.1101/2020.09.01.278689:2020.09.01.278689.
- 1236 77. Latz CA, DeCarlo C, Boitano L, Png CYM, Patell R, Conrad MF, Eagleton M, Dua A. 2020.
1237 Blood type and outcomes in patients with COVID-19. *Ann Hematol* 99:2113-2118.
- 1238 78. Wu BB, Gu DZ, Yu JN, Yang J, Shen WQ. 2020. Association between ABO blood groups and
1239 COVID-19 infection, severity and demise: A systematic review and meta-analysis. *Infect Genet
1240 Evol* 84:104485.
- 1241 79. Zhao J, Yang Y, Huang H, Li D, Gu D, Lu X, Zhang Z, Liu L, Liu T, Liu Y, He Y, Sun B, Wei M,
1242 Yang G, Wang X, Zhang L, Zhou X, Xing M, Wang PG. 2020. Relationship between the ABO
1243 Blood Group and the COVID-19 Susceptibility. *Clin Infect Dis* doi:10.1093/cid/ciaa1150.
- 1244 80. Zietz M, Tatonetti NP. 2020. Testing the association between blood type and COVID-19 infection,
1245 intubation, and death. *medRxiv* doi:10.1101/2020.04.08.20058073.
- 1246 81. Lemieux J, Siddle KJ, Shaw BM, Loreth C, Schaffner S, Gladden-Young A, Adams G, Fink T,
1247 Tomkins-Tinch CH, Krasilnikova LA, Deruff KC, Rudy M, Bauer MR, Lagerborg KA,
1248 Normandin E, Chapman SB, Reilly SK, Anahtar MN, Lin AE, Carter A, Myhrvold C, Kemball M,
1249 Chaluvadi S, Cusick C, Flowers K, Neumann A, Cerrato F, Farhat M, Slater D, Harris JB, Branda
1250 J, Hooper D, Gaeta JM, Baggett TP, O'Connell J, Gnirke A, Lieberman TD, Philippakis A, Burns
1251 M, Brown C, Luban J, Ryan ET, Turbett SE, LaRocque RC, Hanage WP, Gallagher G, Madoff
1252 LC, Smole S, Pierce VM, Rosenberg ES, et al. 2020. Phylogenetic analysis of SARS-CoV-2 in the
1253 Boston area highlights the role of recurrent importation and superspreading events. *medRxiv*
1254 doi:10.1101/2020.08.23.20178236:2020.08.23.20178236.
- 1255 82. Gudbjartsson DF, Helgason A, Jonsson H, Magnusson OT, Melsted P, Norddahl GL,
1256 Saemundsdottir J, Sigurdsson A, Sulem P, Agustsdottir AB, Eiriksdottir B, Fridriksdottir R,
1257 Gardarsdottir EE, Georgsson G, Gretarsdottir OS, Gudmundsson KR, Gunnarsdottir TR, Gylfason
1258 A, Holm H, Jensson BO, Jonasdottir A, Jonsson F, Josefsdottir KS, Kristjansson T, Magnusdottir
1259 DN, le Roux L, Sigmundsdottir G, Sveinbjornsson G, Sveinsdottir KE, Sveinsdottir M,
1260 Thorarensen EA, Thorgjornsson B, Löve A, Masson G, Jonsdottir I, Möller AD, Gudnason T,
1261 Kristinsson KG, Thorsteinsdottir U, Stefansson K. 2020. Spread of SARS-CoV-2 in the Icelandic
1262 Population. *N Engl J Med* doi:10.1056/NEJMoa2006100.
- 1263 83. Candido DS, Claro IM, de Jesus JG, Souza WM, Moreira FRR, Dellicour S, Mellan TA, du
1264 Plessis L, Pereira RHM, Sales FCS, Manuli ER, Thézé J, Almeida L, Menezes MT, Voloch CM,
1265 Fumagalli MJ, Coletti TM, da Silva CAM, Ramundo MS, Amorim MR, Hoeltgebaum HH, Mishra
1266 S, Gill MS, Carvalho LM, Buss LF, Prete CA, Ashworth J, Nakaya HI, Peixoto PS, Brady OJ,
1267 Nicholls SM, Tanuri A, Rossi ÁD, Braga CKV, Gerber AL, de C. Guimarães AP, Gaburo N,
1268 Alencar CS, Ferreira ACS, Lima CX, Levi JE, Granato C, Ferreira GM, Francisco RS, Granja F,
1269 Garcia MT, Moretti ML, Perroud MW, Castilheiras TMPP, Lazari CS, et al. 2020. Evolution and
1270 epidemic spread of SARS-CoV-2 in Brazil. *Science* 369:1255-1260.
- 1271 84. Ellinghaus D, Degenhardt F, Bujanda L, Buti M, Albillos A, Invernizzi P, Fernández J, Prati D,
1272 Baselli G, Assetta R, Grimsrud MM, Milani C, Aziz F, Kässens J, May S, Wendorff M,
1273 Wienbrandt L, Uellendahl-Werth F, Zheng T, Yi X, de Pablo R, Chercoles AG, Palom A, Garcia-
1274 Fernandez AE, Rodriguez-Frias F, Zanella A, Bandera A, Protti A, Aghemo A, Lleo A, Biondi A,
1275 Caballero-Garralda A, Gori A, Tanck A, Carreras Nolla A, Latiano A, Fracanzani AL, Peschuck
1276 A, Julià A, Pesenti A, Voza A, Jiménez D, Mateos B, Nafria Jimenez B, Quereda C, Paccapelo C,

- 1277 Gassner C, Angelini C, Cea C, Solier A, et al. 2020. Genomewide Association Study of Severe
1278 Covid-19 with Respiratory Failure. *N Engl J Med* doi:10.1056/NEJMoa2020283.
- 1279 85. 2020. The COVID-19 Host Genetics Initiative, a global initiative to elucidate the role of host
1280 genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic. *Eur J Hum
1281 Genet* 28:715-718.
- 1282 86. Weisblum Y, Schmidt F, Zhang F, DaSilva J, Poston D, Lorenzi JCC, Muecksch F, Rutkowska M,
1283 Hoffmann H-H, Michailidis E, Gaebler C, Agudelo M, Cho A, Wang Z, Gazumyan A, Cipolla M,
1284 Luchsinger L, Hillyer CD, Caskey M, Robbiani DF, Rice CM, Nussenzweig MC, Hatziloannou T,
1285 Bieniasz PD. 2020. Escape from neutralizing antibodies by SARS-CoV-2 spike protein variants.
1286 *bioRxiv* doi:10.1101/2020.07.21.214759:2020.07.21.214759.
- 1287 87. Li T, Han X, Wang Y, Gu C, Wang J, Hu C, Li S, Wang K, Luo F, Huang J, Long Y, Song S,
1288 Wang W, Hu J, Wu R, Mu S, Hao Y, Chen Q, Gao F, Shen M, Long S, Gong F, Li L, Wu Y, Xu
1289 W, Cai X, Qu D, Yuan Z, Gao Q, Zhang G, He C, Nai Y, Deng K, Du L, Tang N, Xie Y, Huang
1290 A, Jin A. 2020. A key linear epitope for a potent neutralizing antibody to SARS-CoV-2 S-RBD.
1291 *bioRxiv* doi:10.1101/2020.09.11.292631:2020.09.11.292631.
- 1292 88. Hansen J, Baum A, Pascal KE, Russo V, Giordano S, Wloga E, Fulton BO, Yan Y, Koon K, Patel
1293 K, Chung KM, Hermann A, Ullman E, Cruz J, Rafique A, Huang T, Fairhurst J, Libertiny C,
1294 Malbec M, Lee W-y, Welsh R, Farr G, Pennington S, Deshpande D, Cheng J, Watty A, Bouffard
1295 P, Babb R, Levenkova N, Chen C, Zhang B, Romero Hernandez A, Saotome K, Zhou Y, Franklin
1296 M, Sivapalasingam S, Lye DC, Weston S, Logue J, Haupt R, Frieman M, Chen G, Olson W,
1297 Murphy AJ, Stahl N, Yancopoulos GD, Kyratsous CA. 2020. Studies in humanized mice and
1298 convalescent humans yield a SARS-CoV-2 antibody cocktail. *Science* 369:1010-1014.
- 1299 89. Long SW, Olsen RJ, Eagar TN, Beres SB, Zhao P, Davis JJ, Brettin T, Xia F, Musser JM. 2017.
1300 Population Genomic Analysis of 1,777 Extended-Spectrum Beta-Lactamase-Producing
1301 Klebsiella pneumoniae Isolates, Houston, Texas: Unexpected Abundance of Clonal
1302 Group 307. *mBio* 8:e00489-17.
- 1303 90. Stucker KM, Schobel SA, Olsen RJ, Hodges HL, Lin X, Halpin RA, Fedorova N, Stockwell TB,
1304 Tovchigrechko A, Das SR, Wentworth DE, Musser JM. 2015. Haemagglutinin mutations and
1305 glycosylation changes shaped the 2012/13 influenza A(H3N2) epidemic, Houston, Texas. *Euro
1306 Surveill* 20.
- 1307 91. Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
1308 improvements in performance and usability. *Mol Biol Evol* 30:772-80.
- 1309 92. Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. 2009. Jalview Version 2--a
1310 multiple sequence alignment editor and analysis workbench. *Bioinformatics* 25:1189-91.
- 1311 93. Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood trees for
1312 large alignments. *PLoS One* 5:e9490.
- 1313 94. Chen T, Guestrin C. 2016. XGBoost: A Scalable Tree Boosting System, abstr Proceedings of the
1314 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San
1315 Francisco, California, USA, Association for Computing Machinery,
- 1316 95. Nguyen M, Brettin T, Long SW, Musser JM, Olsen RJ, Olson R, Shukla M, Stevens RL, Xia F,
1317 Yoo H. 2018. Developing an in silico minimum inhibitory concentration panel test for *Klebsiella
1318 pneumoniae*. *Scientific reports* 8:421.
- 1319 96. Nguyen M, Long SW, McDermott PF, Olsen RJ, Olson R, Stevens RL, Tyson GH, Zhao S, Davis
1320 JJ. 2019. Using machine learning to predict antimicrobial MICs and associated genomic features
1321 for nontyphoidal Salmonella. *Journal of Clinical Microbiology* 57:e01260-18.
- 1322 97. Pedregosa F, Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
1323 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M.,
1324 Perrot, M., and Duchesnay, E. 2011. Scikit-learn: Machine Learning in Python. *Journal of
1325 Machine Learning Research* 12 (2011) 2825-2830.
- 1326 98. Cai Y, Zhang J, Xiao T, Peng H, Sterling SM, Walsh RM, Jr., Rawson S, Rits-Volloch S, Chen B.
1327 2020. Distinct conformational states of SARS-CoV-2 spike protein. *Science*
1328 doi:10.1126/science.abd4251.
- 1329

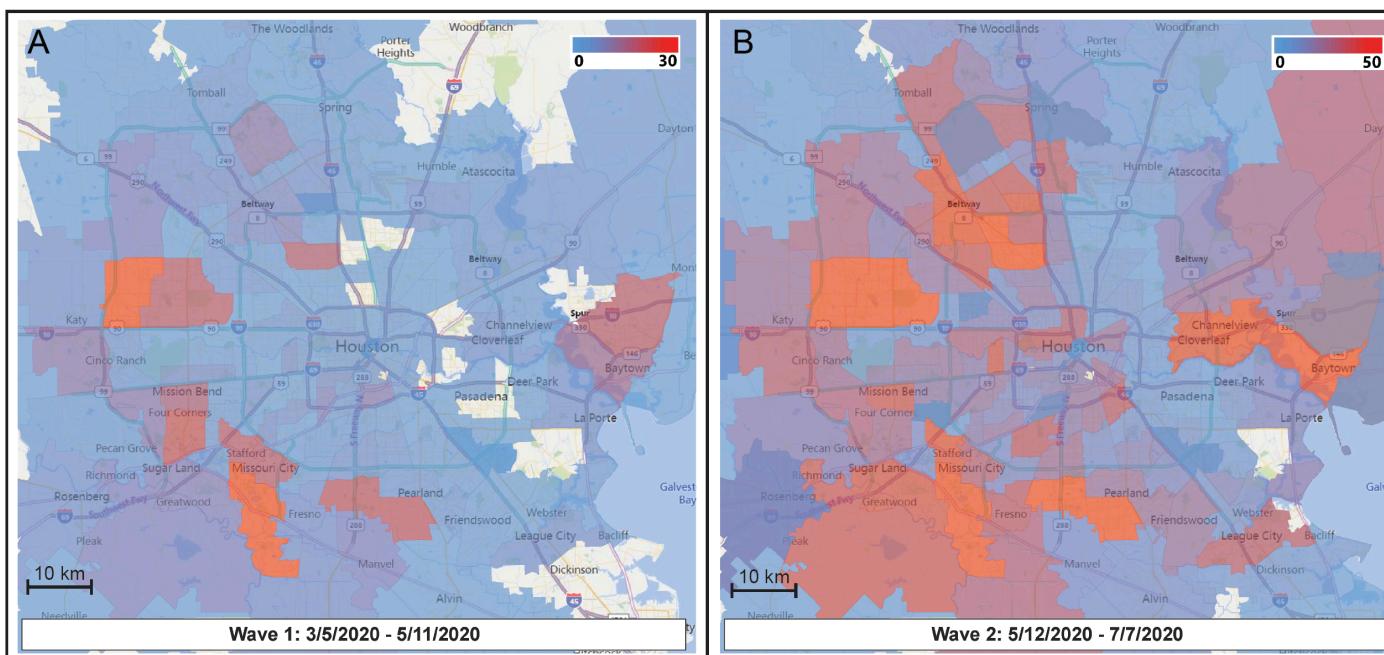
A



B



1331 **FIG 1** (A) Confirmed COVID-19 cases in the Greater Houston Metropolitan region.
1332 Cumulative number of COVID-19 patients over time through July 7, 2020. Counties
1333 include Austin, Brazoria, Chambers, Fort Bend, Galveston, Harris, Liberty, Montgomery,
1334 and Waller. The shaded area represents the time period during which virus genomes
1335 characterized in this study were recovered from COVID-19 patients. The red line
1336 represents the number of COVID-19 patients diagnosed in the Houston Methodist
1337 Hospital Molecular Diagnostic Laboratory. (B) Distribution of strains with either the
1338 Asp614 or Gly614 amino acid variant in spike protein among the two waves of COVID-
1339 19 patients diagnosed in the Houston Methodist Hospital Molecular Diagnostic
1340 Laboratory. The large inset shows major clade frequency for the time frame studied.
1341

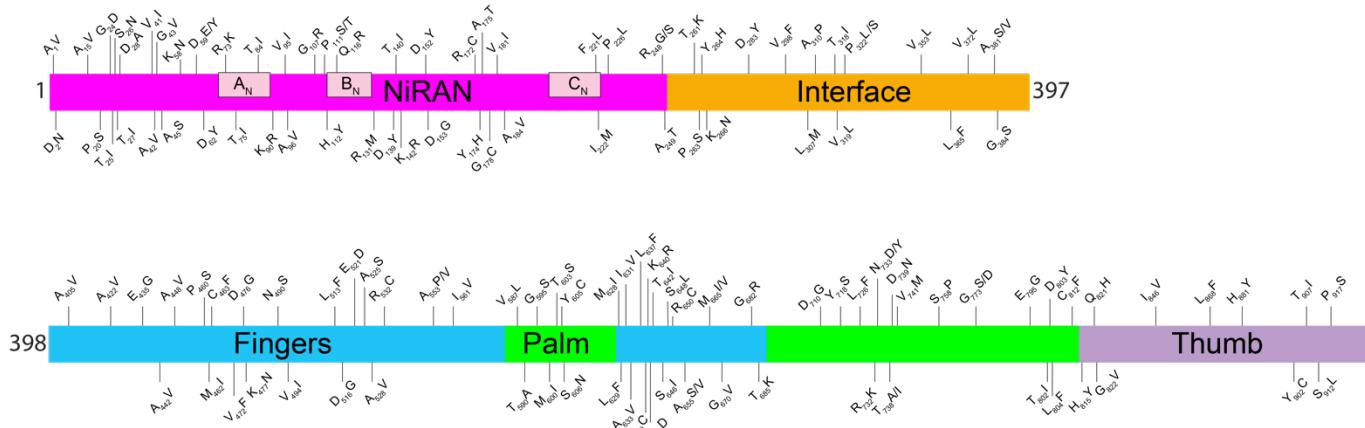


1342

1343 **FIG 2** Sequential time-series heatmaps for all COVID-19 Houston Methodist patients
1344 during the study period. Geospatial distribution of COVID-19 patients is based on zip
1345 code. Panel A (left) shows geospatial distribution of sequenced SARS-CoV-2 strains in
1346 wave 1 and panel B (right) shows wave 2 distribution. The collection dates are shown at
1347 the bottom of each panel. The insets refer to numbers of strains in the color spectrum
1348 used. Note difference in numbers of strains used in panel A and panel B insets.

1349

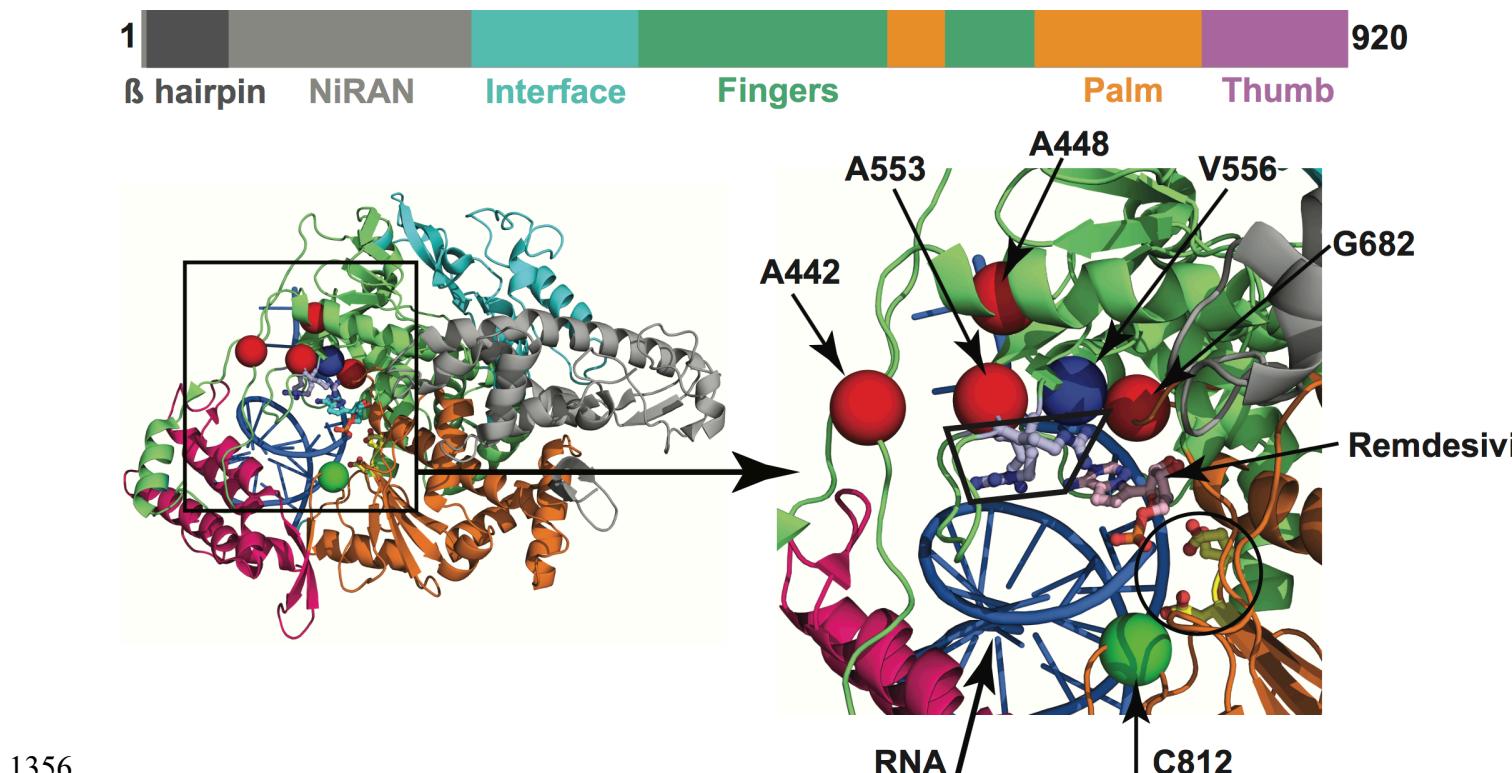
RNA-dependent RNA polymerase (RdRp)



1350

1351 **FIG 3** Location of amino acid replacements in RNA-dependent RNA polymerase
1352 (RdRp/Nsp12) among the 5,085 genomes of SARS-CoV-2 sequenced. The various
1353 RdRp domains are color-coded. The numbers refer to amino acid site. Note that several
1354 amino acid sites have multiple variants identified.

1355

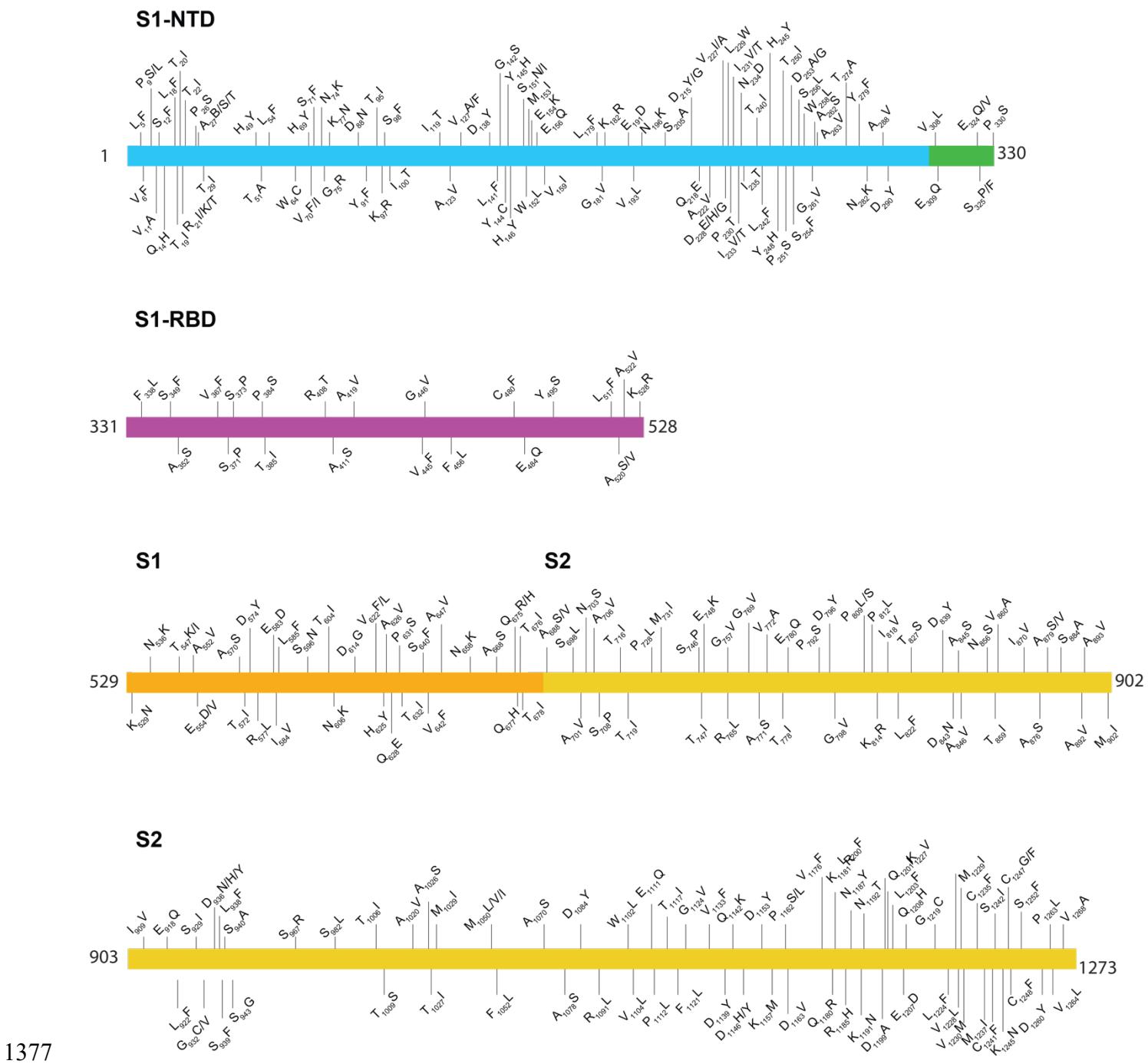


1356

1357 **FIG 4** Amino acid changes identified in Nsp12 (RdRp) in this study that may influence
1358 interaction with remdesivir. The schematic at the top shows the domain architecture of
1359 Nsp12. (Left) Ribbon representation of the crystal structure of Nsp12-remdesivir
1360 monophosphate-RNA complex (PDB code: 7BV2). The structure in the right panel
1361 shows a magnified view of the boxed area in the left panel. The Nsp12 domains are
1362 colored as in the schematic at the top. The catalytic site in Nsp12 is marked by a black
1363 circle in the right panel. The side chains of amino acids comprising the catalytic site of
1364 RdRp (Ser758, Asp759, and Asp760) are shown as balls and stick and colored yellow.
1365 The nucleotide binding site is boxed in the right panel. The side chains of amino acids
1366 participating in nucleotide binding (Lys544, Arg552, and Arg554) are shown as balls and
1367 sticks and colored light blue. Remdesivir molecule incorporated into the nascent RNA is
1368 shown as balls and sticks and colored light pink. The RNA is shown as a blue cartoon
1369 and bases are shown as sticks. The positions of C_α atoms of amino acids identified in

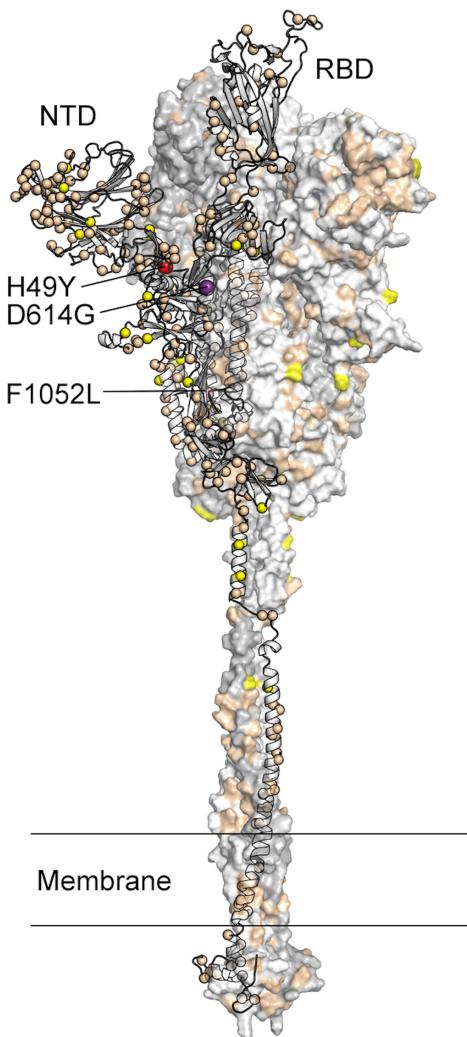
1370 this study are shown as red and green spheres and labeled. The amino acids that are
1371 shown as red spheres are located above the nucleotide binding site, whereas Cys812
1372 located at the catalytic site is shown as a green sphere. The side chain of active site
1373 residue Ser758 is shown as ball and sticks and colored yellow. The location of C α atoms
1374 of remdesivir resistance conferring amino acid Val556 is shown as blue sphere and
1375 labeled.

1376



1378 **FIG 5** Location of amino acid replacements in spike protein among the 5,085 genomes
 1379 of SARS-CoV-2 sequenced. The various spike protein domains are color-coded. The
 1380 numbers refer to amino acid site. Note that many amino acid sites have multiple variants
 1381 identified.

1382

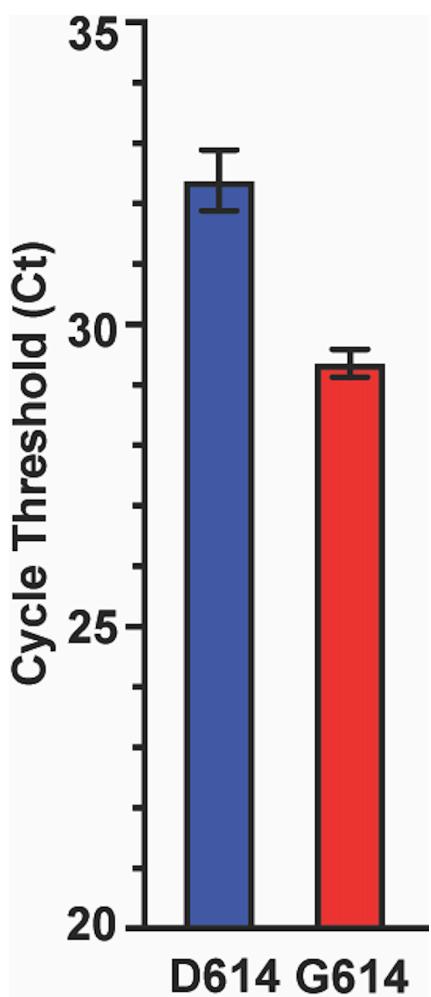


1383

1384 **FIG 6** Location of amino acid substitutions mapped on the SARS-CoV-2 spike protein.
1385 Model of the SARS-CoV-2 spike protein with one protomer shown as ribbons and the
1386 other two protomers shown as a molecular surface. The Ca atom of residues found to be
1387 substituted in one or more virus isolates identified in this study is shown as a sphere on
1388 the ribbon representation. Residues found to be substituted in 1–9 isolates are colored
1389 tan, 10–99 isolates yellow, 100–999 isolates colored red (H49Y and F1052L), and >1000
1390 isolates purple (D614G). The surface of the aminoterminal domain (NTD) that is distal to

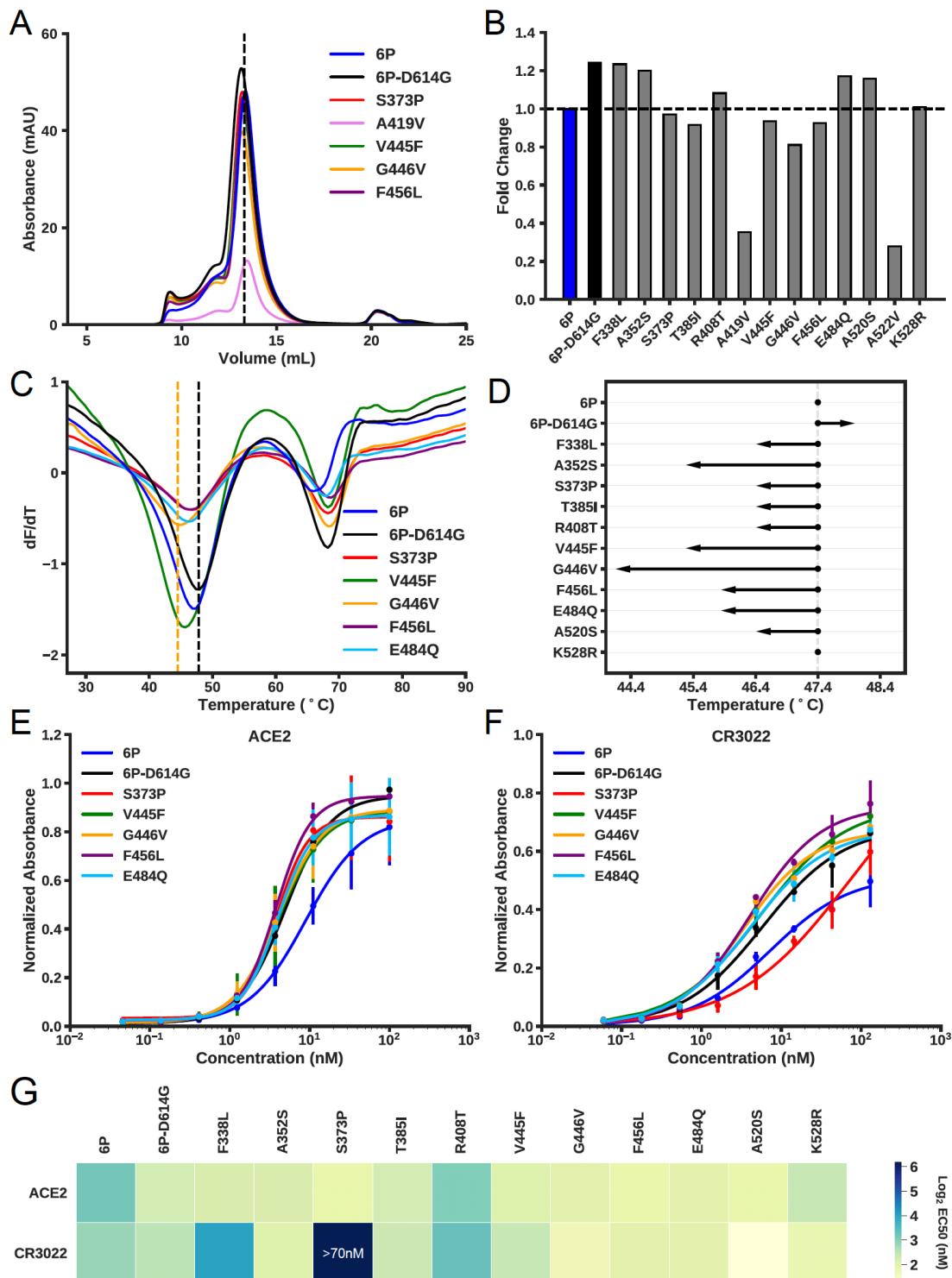
1391 the trimeric axis has a high density of substituted residues. RBD, receptor binding

1392 domain.



1393

1394 **FIG 7** Cycle threshold (Ct) for every SARS-CoV-2 patient sample tested using the
1395 Hologic Panther assay. Data are presented as mean +/- standard error of the mean for
1396 strains with an aspartate (D614, n=102 strains, blue) or glycine (G614, n=812 strains,
1397 red) at amino acid 614 of the spike protein. Mann-Whitney test, *P<0.0001.



1398

1399 **FIG 8** Biochemical characterization of spike RBD variants. (A) Size-exclusion
 1400 chromatography (SEC) traces of the indicated spike-RBD variants. Dashed line indicates
 1401 the elution peak of spike-6P. (B) The relative expression of all RBD variants as

1402 determined by the area under the SEC traces. All expression levels are normalized
1403 relative to spike-6P. (C) Thermostability analysis of RBD variants by differential scanning
1404 fluorimetry. Each sample had three replicates and only mean values were plotted. Black
1405 vertical dashed line indicates the first melting temperature of 6P-D614G and orange
1406 vertical dashed line indicates the first melting temperature of the least stable variant
1407 (spike-G446V). (D) First apparent melting temperature of all RBD variants. (E) ELISA-
1408 based binding affinities for ACE2 and (F) the neutralizing antibody CR3022 to the
1409 indicated RBD variants. (G) Summary of EC50s for all measured RBD variants.
1410
1411

1412 **Table 1.** Nonsynonymous SNPs of SARS-CoV-2 *nsp12*.

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
13446	C3T	A1V	N-terminus		2	2
13448	G5A	D2N	N-terminus	1		1
13487	C44T	A15V	N-terminus		138	138
13501	C58T	P20S	N-terminus		1	1
13514	G71A	G24D	N-terminus		3	3
13517	C74T	T25I	N-terminus		4	4
13520	G77A	S26N	N-terminus		1	1
13523	C80T	T27I	N-terminus		1	1
13526	A83C	D28A	N-terminus		1	1
13564	G121A	V41I	B hairpin		1	1
13568	C125T	A42V	B hairpin	1		1
13571	G128T	G43V	B hairpin	1		1
13576	G133T	A45S	B hairpin		12	12
13617	G174T	K58N	NiRAN		1	1
13618	G175T	D59Y	NiRAN		24	24
13620	C177G	D59E	NiRAN		1	1
13627	G184T	D62Y	NiRAN	1		1
13661	G218A	R73K	NiRAN		1	1
13667	C224T	T75I	NiRAN		2	2
13694	C251T	T84I	NiRAN		1	1
13712	A269G	K90R	NiRAN		1	1
13726	G283A	V95I	NiRAN		1	1
13730	C287T	A96V	NiRAN	2	2	4
13762	G319C	G107R	NiRAN		1	1
13774	C331A	P111T	NiRAN		1	1
13774	C331T	P111S	NiRAN		15	15
13777	C334T	H112Y	NiRAN		1	1
13790	A347G	Q116R	NiRAN		2	2
13835	G392T	R131M	NiRAN		1	1
13858	G415T	D139Y	NiRAN		3	3
13862	C419T	T140I	NiRAN	1	5	6
13868	A425G	K142R	NiRAN	1		1
13897	G454T	D152Y	NiRAN		4	4
13901	A458G	D153G	NiRAN		2	2
13957	C514T	R172C	NiRAN		2	2
13963	T520C	Y174H	NiRAN		1	1
13966	G523A	A175T	NiRAN		1	1
13975	G532T	G178C	NiRAN		4	4

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
13984	G541A	V181I	NiRAN		1	1
13994	C551T	A184V	NiRAN		8	8
14104	T661C	F221L	NiRAN	2		2
14109	A666G	I222M	NiRAN	1		1
14120	C677T	P226L	NiRAN		2	2
14185	A742G	R248G	NiRAN		1	1
14187	G744T	R248S	NiRAN		1	1
14188	G745A	A249T	NiRAN		1	1
14225	C782A	T261K	Interface		4	4
14230	C787T	P263S	Interface		1	1
14233	T790C	Y264H	Interface		1	1
14241	G798T	K266N	Interface		1	1
14290	G847T	D283Y	Interface		1	1
14335	G892T	V298F	Interface		8	8
14362	C919A	L307M	Interface		2	2
14371	G928C	A310P	Interface	1		1
14396	C953T	T318I	Interface		1	1
14398	G955T	V319L	Interface		1	1
14407	C964T	P322S	Interface		2	2
14408	C965T	P322L	Interface	843	4050	4893
14500	G1057T	V353L	Interface		5	5
14536	C1093T	L365F	Interface		1	1
14557	G1114T	V372L	Fingers		4	4
14584	G1141T	A381S			1	1
14585	C1142T	A381V	Fingers		10	10
14593	G1150A	G384S	Fingers	1		1
14657	C1214T	A405V	Fingers		1	1
14708	C1265T	A422V			1	1
14747	A1304G	E435G	Fingers		2	2
14768	C1325T	A442V	Fingers		21	21
14786	C1343T	A448V	Fingers	3	6	9
14821	C1378T	P460S			1	1
14829	G1386T	M462I	Fingers		59	59
14831	G1388T	C463F	Fingers		3	3
14857	G1414T	V472F			1	1
14870	A1427G	D476G	Fingers		5	5
14874	G1431T	K477N	Fingers	1		1
14912	A1469G	N490S	Fingers	1	1	2
14923	G1480A	V494I	Fingers		2	2

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
14980	C1537T	L513F	Fingers	1	1	2
14990	A1547G	D516G			1	1
15006	G1563C	E521D	Fingers	2	3	5
15016	G1573T	A525S	Fingers		3	3
15026	C1583T	A528V	Fingers	5	1	6
15037	C1594T	R532C	Fingers		1	1
15100	G1657C	A553P	Fingers		1	1
15101	C1658T	A553V	Fingers		1	1
15124	A1681G	I561V	Fingers		2	2
15202	G1759C	V587L	Palm		7	7
15211	A1768G	T590A			1	1
15226	G1783A	G595S	Palm		1	1
15243	G1800T	M600I	Palm	71	4	75
15251	C1808G	T603S	Palm		1	1
15257	A1814G	Y605C			1	1
15260	G1817A	S606N	Palm		1	1
15327	G1884T	M628I	Fingers	3	1	4
15328	C1885T	L629F	Fingers	1		1
15334	A1891G	I631V	Fingers		1	1
15341	C1898T	A633V	Fingers		1	1
15352	C1909T	L637F	Fingers		1	1
15358	C1915T	R639C	Fingers		1	1
15362	A1919G	K640R	Fingers		1	1
15364	C1921G	H641D	Fingers	1		1
15368	C1925T	T642I	Fingers	1		1
15380	G1937T	S646I	Fingers	1		1
15386	C1943T	S648L	Fingers		2	2
15391	C1948T	R650C	Fingers		1	1
15406	G1963T	A655S	Fingers		3	3
15407	C1964T	A655V	Fingers	1		1
15436	A1993G	M665V	Fingers		2	2
15438	G1995T	M665I	Fingers		24	24
15452	G2009T	G670V	Fingers		28	28
15487	G2044C	G682R	Palm		1	1
15497	C2054A	T685K	Palm		1	1
15572	A2129G	D710G	Palm		1	1
15596	A2153G	Y718S	Palm	2		2
15619	C2176T	L726F	Palm	1		1
15638	G2195A	R732K	Palm	1		1

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
15640	A2197G	N733D	Palm	1		1
15640	A2197T	N733Y	Palm	1		1
15655	A2212G	T738A	Palm		2	2
15656	C2213T	T738I	Palm		2	2
15658	G2215A	D739N	Palm		2	2
15664	G2221A	V741M	Palm		1	1
15715	T2272C	S758P	Palm		1	1
15760	G2317A	G773S	Palm	1		1
15761	G2318A	G773D	Palm		1	1
15827	A2384G	E795G	Palm	1		1
15848	C2405T	T802I	Palm		1	1
15850	G2407T	D803Y	Palm		1	1
15853	C2410T	L804F	Palm		2	2
15878	G2435T	C812F	Palm		1	1
15886	C2443T	H815Y	Palm		1	1
15906	G2463T	Q821H	Thumb	1	1	2
15908	G2465T	G822V	Thumb		1	1
15979	A2536G	I846V	Thumb	4		4
16045	C2602T	L868F	Thumb		1	1
16084	C2641T	H881Y	Thumb		1	1
16148	A2705G	Y902C	Thumb		1	1
16163	C2720T	T907I	Thumb		45	45
16178	C2735T	S912L	Thumb		2	2
16192	C2749T	P917S	Thumb		80	80

1413

1414

14**Table 2.** Nonsynonymous SNPs in SARS-CoV-2 spike protein.

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
21575	C13T	L5F	S1	11	25	36
21578	G16T	V6F	S1		1	1
21587	C25T	P9S	S1	2		2
21588	C26T	P9L	S1	1	1	2
21594	T32C	V11A	S1		1	1
21597	C35T	S12F	S1		6	6
21604	G42T	Q14H	S1		1	1
21614	C52T	L18F	S1 - NTD	1	11	12
21618	C56T	T19I	S1 - NTD	1	1	2
21621	C59T	T20I	S1 - NTD		1	1
21624	G62T	R21I	S1 - NTD		6	6
21624	G62A	R21K	S1 - NTD		1	1
21624	G62C	R21T	S1 - NTD		3	3
21627	C65T	T22I	S1 - NTD	2	4	6
21638	C76T	P26S	S1 - NTD		17	17
21641	G79T	A27S	S1 - NTD	1	1	2
21641	G79A	A27T	S1 - NTD	1		1
21642	C80T	A27V	S1 - NTD		1	1
21648	C86T	T29I	S1 - NTD	1	4	5
21707	C145T	H49Y	S1 - NTD		142	142
21713	A151G	T51A	S1 - NTD		1	1
21724	G162T	L54F	S1 - NTD		11	11
21754	G192T	W64C	S1 - NTD		1	1
21767	C205T	H69Y	S1 - NTD	1	7	8
21770	G208A	V70I	S1 - NTD		1	1
21770	G208T	V70F	S1 - NTD	1		1
21774	C212T	S71F	S1 - NTD		1	1
21784	T222A	N74K	S1 - NTD	1		1
21785	G223C	G75R	S1 - NTD		1	1
21793	G231T	K77N	S1 - NTD		1	1
21824	G262A	D88N	S1 - NTD		1	1
21834	A272T	Y91F	S1 - NTD		1	1
21846	C284T	T95I	S1 - NTD	1	10	11
21852	A290G	K97R	S1 - NTD		1	1
21855	C293T	S98F	S1 - NTD	1	2	3
21861	T299C	I100T	S1 - NTD		2	2

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
21918	T356C	I119T	S1 - NTD	1		1
21930	C368T	A123V	S1 - NTD		1	1
21941	G379T	V127F	S1 - NTD		1	1
21942	T380C	V127A	S1 - NTD		4	4
21974	G412T	D138Y	S1 - NTD	2		2
21985	G423T	L141F	S1 - NTD		1	1
21986	G424A	G142S	S1 - NTD		2	2
21993	A431G	Y144C	S1 - NTD	1		1
21995	T433C	Y145H	S1 - NTD	2		2
21998	C436T	H146Y	S1 - NTD	1	2	3
22014	G452A	S151N	S1 - NTD		1	1
22014	G452T	S151I	S1 - NTD		2	2
22017	G455T	W152L	S1 - NTD	1	1	2
22021	G459T	M153I	S1 - NTD		1	1
22021	G459A	M153I	S1 - NTD		1	1
22022	G460A	E154K	S1 - NTD		1	1
22028	G466C	E156Q	S1 - NTD	2		2
22037	G475A	V159I	S1 - NTD	1		1
22097	C535T	L179F	S1 - NTD		1	1
22104	G542T	G181V	S1 - NTD		1	1
22107	A545G	K182R	S1 - NTD		1	1
22135	A573T	E191D	S1 - NTD		1	1
22139	G577T	V193L	S1 - NTD		1	1
22150	T588G	N196K	S1 - NTD	1		1
22175	T613G	S205A	S1 - NTD		1	1
22205	G643T	D215Y	S1 - NTD		1	1
22206	A644G	D215G	S1 - NTD		2	2
22214	C652G	Q218E	S1 - NTD		1	1
22227	C665T	A222V	S1 - NTD		1	1
22241	G679A	V227I	S1 - NTD		2	2
22242	T680C	V227A	S1 - NTD	1		1
22244	G682C	D228H	S1 - NTD		2	2
22245	A683G	D228G	S1 - NTD	1		1
22246	T684G	D228E	S1 - NTD	2		2
22248	T686G	L229W	S1 - NTD	1		1
22250	C688A	P230T	S1 - NTD	1		1
22253	A691G	I231V	S1 - NTD	1		1

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
22254	T692C	I231T	S1 - NTD	1		1
22259	A697G	I233V	S1 - NTD	1		1
22260	T698C	I233T	S1 - NTD	1		1
22262	A700G	N234D	S1 - NTD	1		1
22266	T704C	I235T	S1 - NTD	1		1
22281	C719T	T240I	S1 - NTD		5	5
22286	C724T	L242F	S1 - NTD		1	1
22295	C733T	H245Y	S1 - NTD		2	2
22304	T742C	Y248H	S1 - NTD		3	3
22311	C749T	T250I	S1 - NTD	1	4	5
22313	C751T	P251S	S1 - NTD		2	2
22320	A758G	D253G	S1 - NTD		2	2
22320	A758C	D253A	S1 - NTD	1		1
22323	C761T	S254F	S1 - NTD		3	3
22329	C767T	S256L	S1 - NTD	1		1
22335	G773T	W258L	S1 - NTD	1		1
22344	G782T	G261V	S1 - NTD	3		3
22346	G784T	A262S	S1 - NTD		4	4
22350	C788T	A263V	S1 - NTD	1		1
22382	A820G	T274A	S1 - NTD		1	1
22398	A836T	Y279F	S1 - NTD	1		1
22408	T846G	N282K	S1 - NTD		1	1
22425	C863T	A288V	S1 - NTD		1	1
22430	G868T	D290Y	S1 - NTD	1		1
22484	G922T	V308L	S1		3	3
22487	G925C	E309Q	S1		1	1
22532	G970C	E324Q	S1		1	1
22533	A971T	E324V	S1		1	1
22535	T973C	S325P	S1		1	1
22536	C974T	S325F	S1		1	1
22550	C988T	P330S	S1 - RBD		2	2
22574	T1012C	F338L	S1 - RBD	1		1
22608	C1046T	S349F	S1 - RBD		1	1
22616	G1054T	A352S	S1 - RBD		7	7
22661	G1099T	V367F	S1 - RBD		1	1
22673	T1111C	S371P	S1 - RBD		3	3
22679	T1117C	S373P	S1 - RBD		1	1

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
22712	C1150T	P384S	S1 - RBD		1	1
22716	C1154T	T385I	S1 - RBD	3		3
22785	G1223C	R408T	S1 - RBD		1	1
22793	G1231T	A411S	S1 - RBD		1	1
22818	C1256T	A419V	S1 - RBD	1		1
22895	G1333T	V445F	S1 - RBD		1	1
22899	G1337T	G446V	S1 - RBD	2		2
22928	T1366C	F456L	S1 - RBD	1		1
23001	G1439T	C480F	S1 - RBD		1	1
23012	G1450C	E484Q	S1 - RBD	1		1
23046	A1484C	Y495S	S1 - RBD		1	1
23111	C1549T	L517F	S1 - RBD		1	1
23120	G1558T	A520S	S1 - RBD	1	6	7
23121	C1559T	A520V	S1 - RBD		1	1
23127	C1565T	A522V	S1 - RBD	1	1	2
23145	A1583G	K528R	S1 - RBD		2	2
23149	G1587T	K529N	S1		1	1
23170	C1608A	N536K	S1		1	1
23202	C1640A	T547K	S1		2	2
23202	C1640T	T547I	S1		1	1
23223	A1661T	E554V	S1		2	2
23224	G1662T	E554D	S1	4	31	35
23270	G1708T	A570S	S1		3	3
23277	C1715T	T572I	S1	5	5	10
23282	G1720T	D574Y	S1		1	1
23292	G1730T	R577L	S1	1		1
23311	G1749T	E583D	S1		6	6
23312	A1750G	I584V	S1		1	1
23315	C1753T	L585F	S1	1	7	8
23349	G1787A	S596N	S1		1	1
23373	C1811T	T604I	S1		2	2
23380	C1818A	N606K	S1		2	2
23403	A1841G	D614G	S1	841	4054	4895
23426	G1864T	V622F	S1		2	2
23426	G1864C	V622L	S1		2	2
23435	C1873T	H625Y	S1		1	1
23439	C1877T	A626V	S1		1	1

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
23444	C1882G	Q628E	S1		7	7
23453	C1891T	P631S	S1	1		1
23457	C1895T	T632I	S1		1	1
23481	C1919T	S640F	S1	1	42	43
23486	G1924T	V642F	S1		1	1
23502	C1940T	A647V	S1		1	1
23536	C1974A	N658K	S1		4	4
23564	G2002T	A668S	S1		1	1
23586	A2024G	Q675R	S1		14	14
23587	G2025C	Q675H	S1		1	1
23587	G2025T	Q675H	S1		4	4
23589	C2027T	T676I	S1	1	2	3
23593	G2031T	Q677H	S1	1	1	2
23595	C2033T	T678I	S1	1		1
23624	G2062T	A688S	S2		4	4
23625	C2063T	A688V	S2		16	16
23655	C2093T	S698L	S2		1	1
23664	C2102T	A701V	S2		21	21
23670	A2108G	N703S	S2		1	1
23679	C2117T	A706V	S2		1	1
23684	T2122C	S708P	S2		1	1
23709	C2147T	T716I	S2		1	1
23718	C2156T	T719I	S2		1	1
23745	C2183T	P728L	S2	1		1
23755	G2193T	M731I	S2	3	1	4
23798	T2236C	S746P	S2		1	1
23802	C2240T	T747I	S2		1	1
23804	G2242A	E748K	S2		1	1
23832	G2270T	G757V	S2		1	1
23856	G2294T	R765L	S2		1	1
23868	G2306T	G769V	S2		3	3
23873	G2311T	A771S	S2		8	8
23877	T2315C	V772A	S2		1	1
23895	C2333T	T778I	S2		1	1
23900	G2338C	E780Q	S2		1	1
23936	C2374T	P792S	S2		1	1

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
23948	G2386T	D796Y	S2		2	2
23955	G2393T	G798V	S2	1		1
23987	C2425T	P809S	S2		2	2
23988	C2426T	P809L	S2		1	1
23997	C2435T	P812L	S2		1	1
24003	A2441G	K814R	S2		1	1
24014	A2452G	I818V	S2 - FP		5	5
24026	C2464T	L822F	S2 - FP		97	97
24041	A2479T	T827S	S2 - FP		4	4
24077	G2515T	D839Y	S2	2		2
24089	G2527A	D843N	S2	1	1	2
24095	G2533T	A845S	S2		5	5
24099	C2537T	A846V	S2		1	1
24129	A2567G	N856S	S2		7	7
24138	C2576T	T859I	S2		5	5
24141	T2579C	V860A	S2		1	1
24170	A2608G	I870V	S2		3	3
24188	G2626T	A876S	S2		1	1
24197	G2635T	A879S	S2		31	31
24198	C2636T	A879V	S2		1	1
24212	T2650G	S884A	S2		11	11
24237	C2675T	A892V	S2		1	1
24240	C2678T	A893V	S2	1		1
24268	G2706T	M902I	S2		1	1
24287	A2725G	I909V	S2 - HR1		2	2
24314	G2752C	E918Q	S2 - HR1	1		1
24328	G2766C	L922F	S2 - HR1		2	2
24348	G2786T	S929I	S2 - HR1		1	1
24356	G2794T	G932C	S2 - HR1		1	1
24357	G2795T	G932V	S2 - HR1		1	1
24368	G2806A	D936N	S2 - HR1		3	3
24368	G2806C	D936H	S2 - HR1		1	1
24368	G2806T	D936Y	S2 - HR1	3	4	7
24374	C2812T	L938F	S2 - HR1		3	3
24378	C2816T	S939F	S2 - HR1		4	4
24380	T2818G	S940A	S2 - HR1		5	5
24389	A2827G	S943G	S2 - HR1		6	6

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
24463	C2901A	S967R	S2 - HR1	2		2
24507	C2945T	S982L	S2 - HR1		1	1
24579	C3017T	T1006I	S2 - CH		1	1
24588	C3026G	T1009S	S2 - CH		1	1
24621	C3059T	A1020V	S2 - CH	1		1
24638	G3076T	A1026S	S2 - CH		2	2
24642	C3080T	T1027I	S2 - CH	5		5
24649	G3087T	M1029I	S2 - CH		1	1
24710	A3148T	M1050L	S2		1	1
24710	A3148G	M1050V	S2	1	1	2
24712	G3150T	M1050I	S2		2	2
24718	C3156A	F1052L	S2	1	166	167
24770	G3208T	A1070S	S2		2	2
24794	G3232T	A1078S	S2 - CD	3	2	5
24812	G3250T	D1084Y	S2 - CD	1	29	30
24834	G3272T	R1091L	S2 - CD	1		1
24867	G3305T	W1102L	S2 - CD		1	1
24872	G3310T	V1104L	S2 - CD		1	1
24893	G3331C	E1111Q	S2 - CD		2	2
24897	C3335T	P1112L	S2 - CD	2	2	4
24912	C3350T	T1117I	S2 - CD		1	1
24923	T3361C	F1121L	S2 - CD		2	2
24933	G3371T	G1124V	S2 - CD	1	2	3
24959	G3397T	V1133F	S2 - CD		1	1
24977	G3415T	D1139Y	S2 - CD		1	1
24986	C3424A	Q1142K	S2	1		1
24998	G3436T	D1146Y	S2		4	4
24998	G3436C	D1146H	S2		13	13
25019	G3457T	D1153Y	S2		11	11
25032	A3470T	K1157M	S2	1		1
25046	C3484T	P1162S	S2		5	5
25047	C3485T	P1162L	S2		3	3
25050	A3488T	D1163V	S2		2	2
25088	G3526T	V1176F	S2		18	18
25101	A3539G	Q1180R	S2		1	1
25104	A3542G	K1181R	S2		4	4
25116	G3554A	R1185H	S2		2	2

Genomic Locus	Gene Locus	Amino Acid Change	Domain	Wave 1 (n=1026)	Wave 2 (n=4059)	Total (n=5085)
25121	A3559T	N1187Y	S2		1	1
25135	G3573T	K1191N	S2		1	1
25137	A3575C	N1192T	S2	1		1
25158	A3596C	D1199A	S2		1	1
25160	C3598T	L1200F	S2		1	1
25163	C3601A	Q1201K	S2		1	1
25169	C3607T	L1203F	S2	1		1
25183	G3621T	E1207D	S2		1	1
25186	G3624T	Q1208H	S2	1		1
25217	G3655T	G1219C	S2	1	3	4
25234	G3672T	L1224F	S2		1	1
25241	A3679G	I1227V	S2	1		1
25244	G3682T	V1228L	S2		2	2
25249	G3687T	M1229I	S2		1	1
25249	G3687C	M1229I	S2		2	2
25250	G3688A	V1230M	S2		1	1
25266	G3704T	C1235F	S2		4	4
25273	G3711T	M1237I	S2		2	2
25284	G3722T	C1241F	S2		1	1
25287	G3725T	S1242I	S2		4	4
25297	G3735T	K1245N	S2		1	1
25301	T3739G	C1247G	S2		1	1
25302	G3740T	C1247F	S2		4	4
25305	G3743T	C1248F	S2		2	2
25317	C3755T	S1252F	S2		1	1
25340	G3778T	D1260Y	S2		2	2
25350	C3788T	P1263L	S2	1	2	3
25352	G3790T	V1264L	S2		1	1
25365	T3803C	V1268A	S2		1	1

1416

1417 The domain region of RBD is based on structural information found in Cai et al.
1418 2020 (98).

1419

1420

1421 Forty-nine of these amino acid replacements (V11A, T51A, W64C, I119T,
1422 E156Q, S205A, D228G, L229W, P230T, N234D, I235T, T274A, A288V, E324Q,
1423 E324V, S325P, S349F, S371P, S373P, T385I, A419V, C480F, Y495S, L517F,
1424 K528R, Q628E, T632I, S708P, T719I, P728L, S746P, E748K, G757V, V772A,

It is made available under a [CC-BY-NC-ND 4.0 International license](#).

1425 K814R, D843N, S884A, M902I, I909V, E918Q, S982L, M1029I, Q1142K,
1426 K1157M, Q1180R, D1199A, C1241F, C1247G, and V1268A) were not
1427 represented in a publicly available database (34) as of August 19, 2020

It is made available under a [CC-BY-NC-ND 4.0 International license](#).

1428