

Variants in SARS-CoV-2 Associated with Mild or Severe Outcome

Jameson D. Voss, MD, MPH, FACPM, FISO¹

Martin Skarzynski, PhD²

Erin M. McAuley, PhD²

Ezekiel J. Maier, PhD²

Thomas Gibbons, PhD³

Anthony C. Fries, PhD⁴

Richard R. Chapleau, PhD, MMOAS⁴

¹US Air Force Medical Readiness Agency, Falls Church, VA 22042

²Booz Allen Hamilton, Bethesda, MD 20814

³59th Medical Wing, Joint Base San Antonio, TX 78234

⁴US Air Force School of Aerospace Medicine, Wright Patterson AFB, OH 45433

Word count: Abstract – 224 words; Body – 2920 words; 2 Figures

Address correspondence and reprint requests to:

Jameson D. Voss, MD, MPH

7700 Arlington Blvd, Falls Church, VA, 22042

Phone: (719) 232-3509; E-mail: jameson.d.voss.mil@mail.mil

Conflict of interest statement: JDV, MS, EMM, EJM, TG, ACF, and RRC have no conflict of interest disclosures.

Disclaimer: The views expressed in this article are those of the authors and do not necessarily reflect the official policy or position of the Air Force, the Department of Defense, or the U.S. Government. Clearance PAIRS CASE #2020-0613.

Acknowledgements: The authors gratefully acknowledge the contributors, originating and submitting laboratories of the sequences from the GISAID EpiCoV Database (Elbe and Buckland-Merrett, 2017; Shu and McCauley, 2017), the basis of this research. A detailed list of contributing labs to GISAID is available in the Supplementary Information.

47 **Abstract:**

48 **Introduction:** The coronavirus disease 2019 (COVID-19) pandemic is a global public health
49 emergency causing a disparate burden of death and disability around the world. The molecular
50 characteristics of the virus that predict better or worse outcome are largely still being discovered.

51 **Methods:** We downloaded 155,958 severe acute respiratory syndrome coronavirus 2 (SARS-
52 CoV-2) genomes from GISAID and evaluated whether variants improved prediction of reported
53 severity beyond age and region. We also evaluated specific variants to determine the magnitude
54 of association with severity and the frequency of these variants among the genomes.

55 **Results:** Logistic regression models that included viral genomic variants outperformed other
56 models (AUC=0.91 as compared with 0.68 for age and gender alone; $p < 0.001$). Among
57 individual variants, we found 17 single nucleotide variants in SARS-CoV-2 have more than two-
58 fold greater odds of being associated with higher severity and 67 variants associated with ≤ 0.5
59 times the odds of severity. The median frequency of associated variants was 0.15% (interquartile
60 range 0.09%-0.45%). Altogether 85% of genomes had at least one variant associated with patient
61 outcome.

62 **Conclusion:** Numerous SARS-CoV-2 variants have two-fold or greater association with odds of
63 mild or severe outcome and collectively, these variants are common. In addition to
64 comprehensive mitigation efforts, public health measures should be prioritized to control the
65 more severe manifestations of COVID-19 and the transmission chains linked to these severe
66 cases.

67

68 **Introduction**

69 Since the coronavirus disease 2019 (COVID-19) pandemic emerged, humans have faced
70 unprecedented disruption from the newfound obligate parasite. Within the United States alone,
71 the unwelcome guest has already caused an estimated 2.5 million years of life lost.(1) Beyond
72 the United States, the global burden is substantial and growing, but it is not uniform; continents,
73 nations, communities, families, and patients are all affected differently. Understanding the basis
74 for this variability is an important global health priority.

75 One of the most common measures to describe the severity of COVID-19 is the infection fatality
76 ratio (IFR) or the number of deaths for every infection. There are widespread differences in IFR
77 between studies,(2-4) and this heterogeneity is not due to chance alone ($p < 0.001$). (2) One meta-
78 analysis of 26 studies estimated an IFR of 0.68% (0.53-0.82%) while cautioning this was likely
79 an “underestimate” and another meta-analysis with 61 studies estimated a median IFR of 0.26%,
80 outside the range of the first meta-analysis.(2, 3)

81 Other studies have suggested the variability in infection fatality ratio is more consistent when
82 infections are stratified by age, but is more likely to vary between studies when considering the
83 population above age 65.(4-6) In fact, one paper even estimated that “differences in age structure
84 of the population and the age-specific prevalence of COVID-19 explain nearly 90% of the
85 geographical variation in population IFR.”(4)

86 Although vulnerable age groups should be protected as recommended by the U.S. Centers for
87 Disease Control and Prevention (CDC), there are likely more factors at play in the IFR than age
88 alone. First, there is variability among some younger groups. In the most recent meta-analysis
89 among those age < 70 , IFR estimates varied from 0.00% to 0.31% with a median of 0.05%.(3)

90 Second, even within the same hospital systems, there are strong time trends in case fatality rate
91 after adjusting for baseline characteristics/risk (7) even at the national scale (8), which authors
92 attributed to better treatment. Unless baseline risk adjustments were inadequate, there are likely
93 time trends in IFR, which would also help reconcile the disparate estimates noted in the two
94 meta-analyses cited above.

95 Another possibility is severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) virulence
96 could differ across geographic locations because some local strains have higher or lower
97 virulence than others. Previous reports have argued that evolution of attenuated strains is
98 expected for RNA viruses and that “determination of naturally circulating attenuated SARS-
99 CoV-2 variants is an urgent matter.” (9) Previous publications on evolution of virulence
100 observed that pathogens transmitting with a respiratory route (unlike vector-borne diseases)
101 typically evolve toward lower virulence after emergence because healthier hosts tend to engage
102 in more social contact.(10)

103 Using publicly available data from GISAID (Global Initiative on Sharing Avian Influenza Data),
104 we interrogate the relationship between SARS-CoV-2 variants and associated patient outcomes
105 in the GISAID metadata. We differentiate between severe and mild patient outcomes and utilize
106 a logistic regression model in order to better understand how viral genomic SARS-CoV-2
107 variants are linked with COVID-19 patient outcomes.

108

109 **Methods:**

110 **Variant Alignment and Variant Calling**

111 SARS-CoV-2 genome sequences were obtained from GISAID (Global Initiative on Sharing All
112 Influenza Data) on October 21, 2020 (11, 12), (GISAID Acknowledgments Table,
113 Supplementary Table 1). GISAID sequences were filtered to include those of human origin.
114 FASTA sequences were aligned to the reference sequence, Wuhan-Hu-1 (NCBI: NC_045512.2;
115 GISAID: EPI_ISL_402125) using Minimap2 (version 2.17).(13) Resulting VCF (Variant Call
116 Format) files were annotated using SnpEff (version 5.0) and filtered using SnpSift.(14, 15) The
117 shell scripts used for variant alignment and variant calling, along with the Python scripts used to
118 perform the steps described below, are available on GitHub at <https://github.com/mskar/variants>.

119

120 **Metadata Preprocessing and Cohort Building**

121 Raw GISAID patient data was parsed from the JSON file using Python (version 3.8.2). Patient
122 outcomes were then aggregated into positive (“Mild”) outcomes or negative (“Severe”)
123 outcomes (detailed in Supplemental Figure 1). Briefly, “Mild” outcomes included: Outpatient,
124 Asymptomatic, Mild, Home/Isolated/Quarantined, and Not Hospitalized. “Severe” outcomes
125 included: Hospitalized (including severe, moderate, and stable) and Deceased (Death). Patient
126 outcomes that were unclear or empty were not included in our analyses.

127

128 **Variant and Metadata Modeling**

129 Annotated VCF files were parsed, pivoted to wide format, and joined with GISAID patient data
130 using Pandas (version 1.0.3).(16) Logistic regression models with the default L1 penalty (Lasso
131 regularization) were fit to the patient (rows) and variant (columns) data using Scikit-learn
132 (version 0.23.2).(17) Logistic regression model Area Under the Curve (AUC) and accuracy

133 values were calculated using Scikit-learn.(17) Models were persisted as pickle files using joblib
134 (version 0.14.1).

135

136 **Plotting and Statistical Analysis**

137 Scatter and bar plots were created using Pandas (version 1.0.3),(16) Matplotlib (version
138 3.2.1),(18) and Seaborn (version 0.10.1).(19) Logistic regression model AUC p-values and Chi-
139 square test p-values for association of variants with “Severe” outcomes were obtained using
140 Scipy (version 1.5.0).(20) Variant frequency was calculated using Pandas.(16) Genome position
141 tracks were added to scatterplots using DNA Features Viewer (version 3.0.3).(21) ROC curves
142 were plotted using Scikit-learn (version 0.23.2),(17) and Matplotlib.(18) Logistic regression
143 model Area Under the Curve (AUC) and accuracy values were calculated using Scikit-learn.(17)

144

145 **Results:**

146 *Sample population characteristics*

147 We collected a total of 155,958 viral genomes along with clinical metadata. The metadata
148 included numerous entries whereby the severity of the condition could not readily be discerned.
149 For example, “recovering”, “recovered and released”, and “mild symptoms inpatient for
150 observation” were found in the raw data and were not included. The full downloaded dataset
151 included 148,121 entries with empty or unknown clinical observations, and 4,200 entries for
152 which clinical severity could not be classified. From the remaining 3,637 sequences with clear
153 severity indications, we generated two classes (Supplemental Figure S1) by recoding the

154 observational metadata into consistent terminology and creating a “Severe” class of “deceased”,
155 “hospitalized”, “ICU”, and “pneumonia” (n=2,870); and a “Mild” class of “outpatient”, “mild”,
156 “epidemiology study”, “asymptomatic”, “screening”, and “stable in quarantine” (n=767). 85% of
157 these genomes had at least one variant associated with patient outcome. Viral sequences were
158 obtained from the six major geographical regions in GISAID between January and October 2020
159 (Supplemental Figure S2).

160 ***SARS-CoV-2 variants associated with “Severe” / “Mild” outcome categories***

161 The overwhelming majority of variants in the SARS-CoV-2 genomes assessed were rare, with
162 only 12 common variants with at least a 5% minor allele frequency (Figure 1C). Two of these
163 common variants (C26735T and C28311T) were associated with “Severe” or “Mild” outcomes,
164 as measured by having an odds ratio of greater than 2 or ≤ 0.5 , respectively. We also observed 84
165 of 157 rare variant associations with “Severe” or “Mild” outcomes. Collectively, 17 variants
166 were associated with “Severe” classification with at least an odds ratio of 2, while we found 67
167 associations with “Mild” classification ($OR \leq 0.5$). The variants associated with outcomes were
168 distributed across the genome, including the strongest “Severe” association within the C-terminal
169 end of the spike protein (Figure 1B). The majority of variants characterized here were transitions
170 (121, 71%), with 47% of those transitions (79) being C>T (Supplemental Figure S3). All
171 individual variant associations and frequencies are reported in Supplementary Table 2.

172

173 ***Predicting clinical outcomes for patients based upon clinical metadata and viral genomics***

174 Age and gender have been previously reported to be predictive of clinical outcomes.(22) Our
175 observations predicting “Severe” outcomes confirm these prior associations (Figure 2); however

176 we found that the c-statistic for predictions based on age or age and gender are only slightly
177 better than random chance at 0.677 (0.642-0.712) and 0.679 (0.644-0.714), respectively. We
178 hypothesized that viral genomic variants in the virus could also contribute to severity
179 classification. When accounting for the region of collection, we observed an increased c-statistic
180 to 0.817 (0.817-0.818). Previous observations show that regions tend to be dominated by
181 individual viral clades, and we found that adding clade to an age/gender/region model resulted in
182 moderate improvement of accuracy (81% vs. 86%, respectively) with a nearly identical AUC of
183 0.818 (0.817-0.818). While the difference between region and clade appears insignificant, adding
184 clade-level information increased the predictive ability of our model beyond age and gender
185 alone. We then considered whether variant-level information would further improve model
186 performance. We found that substituting clade with 4,499 genomic variants in the model
187 increased the c-statistic to 0.911 (0.910-0.911), significantly improving predictability for clinical
188 outcomes. In addition to the improvement in c-statistic, we compared the accuracy of our
189 predictions, which started at 81% for the age-only model and improved to 86% and 88% for each
190 additional step in the model building before finally reaching a maximum at 91% accuracy for the
191 age/gender/region/variant model.

192 Classifications based solely on age or on age and gender resulted in insignificant odds ratios in
193 our models. Both models had the same odds ratio (4.4), confidence interval (2.8-6.0), and p-
194 value (0.072). However, the other three models all had significant odds ratios with p-values less
195 than 0.0001. Consistent with the AUC results, the odds ratio was greatest for the full model
196 (age/gender/region/variant) followed by age/gender/region/clade and finally age/gender/region
197 (odds ratios: 12.3 (11.8-12.8), 8.4 (8.0-8.9), and 8.0 (7.6-8.4), respectively). Similarly, the
198 negative likelihood ratio for the full model displayed a large reduction in the likelihood of a

199 patient classified as “Mild” developing “Severe” symptoms (-LR=0.039) as compared to a
200 moderate reduction in the post-test outcomes for the age-only or age and gender models (-
201 LR=0.231 for both models).

202

203 **Discussion:**

204 We demonstrate that including genomic viral variants can substantially improve classification of
205 COVID-19 patient outcomes as compared with models using only age and region. Moreover, in
206 our models, we observe that some individual variants are particularly important with substantial
207 associations with severity, and that collectively these variants are not rare.

208 Associations between viral genomic variants and patient outcomes are not unexpected.

209 Consistent with known patterns in the evolution of virulence in RNA viruses (9, 10, 23), we
210 would expect many common strains have differing association with patient severity by this point
211 in the pandemic by chance alone and even as sampling is more likely to occur with severe
212 outcomes, variants correlated with mild outcomes are still being identified.

213 Though respiratory pathogens often evolve toward lower virulence, there have been historical
214 exceptions.(24) Modeling future fitness landscapes suggests that even partial isolation of
215 symptomatic cases can substantially reduce deaths with less transmission in the short term.

216 Importantly, this isolation can also potentially alter the evolutionary path by favoring less
217 virulent strains.(24) Alterations in virulence can happen with a small number of selections. For
218 example, *E. faecalis* evolved from a pathogen to a commensal strain in 15 passages in a worm
219 model, but most of the worm phenotype changed after just 5 rounds of bacterial selection.(25) A
220 mouse experiment showed higher virulence after 10 passages of a modified SARS-CoV-2 virus

221 with an additional ~1% drop in mouse body weight with each selection (26). Based upon these
222 findings and the large number of passages in the human outbreak, it can be expected that
223 significant evolution affecting virulence could occur in the SARS-CoV-2 genome.

224 Indeed, others have previously found and characterized individual variants with *in vitro* assays
225 (27) or provided correlates of severity with any change in protein coding,(28) or genomic
226 correlates of mortality.(29) We have taken a comprehensive approach to describe all variants
227 associated with mild or severe outcome regardless of whether it is synonymous. There are
228 challenges in identifying signatures of selection in non-coding regions, but “for RNA
229 viruses...critical aspects of the life cycle rely on molecular processes that are not reflected in
230 protein sequence.”(30) Empiric tests of selection and structural modeling of RNA and protein
231 interaction can identify regions under selection without using the ratio of nonsynonymous to
232 synonymous mutations.(30) Studies of selection in SARS-CoV-2 have also advised not to
233 underestimate the role of synonymous substitutions.(31) Beyond RNA molecule interactions,
234 there also appear to be selective pressures on which codon is used for an amino acid, which
235 could be attributed to tRNA abundance (32) or could be related to broader patterns of host RNA
236 editing involving deamination and similar mechanisms. (33-35) Alternatively, a variant
237 correlated with severity might represent an epiphenomenon that is linked to multiple variants that
238 each have a smaller association with severity. Regardless of the applicability of these
239 explanations for why synonymous changes could be important indicators of virulence, we
240 wanted to take an agnostic approach to characterize all variants correlated with severity so that
241 they could be further resolved with additional study and additional surveillance (particularly
242 among asymptomatic cases). These studies, in addition to comprehensive prognostic study, could
243 better clarify how unexpected a patient’s severity was as compared with additional risk factors.

244 For example, one mutation we identified (C13620T) is associated with 5.9 times the odds of
245 severe disease. Although it does not result in any change in an amino acid of the NSP12 (RNA-
246 dependent RNA polymerase), it could result in altered expression. Because NSP12 is required
247 for the transcription of all viral RNA in coronaviruses (36), increased replication could increase
248 virulence.

249 Other mutations were nonsynonymous. We identified two previously reported spike mutations,
250 V1176F (G25088T) and S477N (G22992A), as important indicators for COVID-19 disease
251 severity. Recent protein modeling studies have indicated that both mutations cause favorable
252 energetic changes that result in a more flexible Spike protein and can change RBD-ACE2
253 binding. Importantly, both mutations have been associated with higher mortality rates, and are
254 therefore expected to have significant impacts on public health.(37, 38) Other mutations were
255 observed less commonly, but could still be relevant for understanding higher viral pathogenesis.
256 The G26144T causes amino acid change G251V in Orf3a and was associated with 4.3 times the
257 odds of severe disease. Protein trafficking is a complex multifactorial process (39) and Orf3a
258 functions as a modulator of the trafficking properties of the spike protein of SARS-1 and is
259 dependent on the protein-protein interaction of Orf3a and S.(40) A structural analysis of the
260 G251V in Orf3a of SARS-CoV-2 results in significant changes in the overall protein structure
261 and weaker affinity for both the S and M proteins with Orf3a. One possible outcome of the
262 weaker Orf3a-S and Orf3a-M interaction could be an increased Orf3a-TRAF3 interaction
263 resulting in increased activation of the NLRP3 inflammasome by promoting TRAF3-dependent
264 ubiquitination of ASC.(41, 42) Thus, the altered protein-protein interaction of the G251V Orf3a
265 may impact the trafficking of Orf3a resulting in a higher propensity for inflammatory cytokine
266 activation.

267 We also identified that the C28311T mutation was associated with a lower odds ratio of 0.14.
268 This mutation lies within the probe of the N1 assay in the CDC's PCR assay (43), creating a
269 P13L change in the nucleocapsid protein. A previous study evaluated how this mutation may
270 alter protein-protein interaction and proposed it impacted virus stability, potentially contributing
271 to lower pathogenesis.(44) Additional follow-up studies will help to illuminate effects of these
272 variants on viral fitness, infectivity, host response, and evolutionary trajectory.

273 Identifying genetic variants associated with outcomes could provide mechanistic understanding
274 of the life cycle.(45) Efforts such as the CDC's annual influenza surveillance rely upon
275 understanding those key genetic variants to predict seasonal intensity and attempt to develop
276 effective countermeasures (vaccines) (<https://www.cdc.gov/flu/weekly/overview.htm>).

277 Although deep molecular insights are important, they are not necessary for public health
278 applications. The CDC offers symptom-based criteria for prioritizing testing and symptom-based
279 criteria are also included in recommendations for prioritizing contact tracing
280 ([https://www.cdc.gov/coronavirus/2019-ncov/php/contact-tracing/contact-tracing-plan/contact-](https://www.cdc.gov/coronavirus/2019-ncov/php/contact-tracing/contact-tracing-plan/contact-tracing.html#)
281 [tracing.html#](https://www.cdc.gov/coronavirus/2019-ncov/php/contact-tracing/contact-tracing-plan/contact-tracing.html#)). By prioritizing cases and contacts with symptoms (as compared with
282 asymptomatic cases and their asymptomatic contacts) in addition to the recommended global
283 mitigation efforts, there could be relative selective pressure against strains that are more likely to
284 cause symptoms. This could favor the emergence of attenuated strains over the long term.

285 The COVID-19 pandemic demonstrated the limitations of the global healthcare system in
286 intensive care units, mechanical ventilators, and emerging therapeutics and other medical
287 countermeasures.(46, 47) Early in the outbreak, cities such as New York became inundated with
288 infections and their ability to adequately sort and treat patients was quickly overwhelmed.(48)
289 The existence of a rapid and accurate tool that could help identify COVID-19 patients or clusters

290 that are more likely to experience severe symptoms or require intensive medical resources (e.g.,
291 inpatient hospitalization and ventilation) may be able to help healthcare systems allocate
292 resources to the regions with the most critical needs. Therefore, by providing a molecular risk
293 factor for more severe outcomes, these findings could help prioritize limited treatment supplies
294 to those at greatest risk, particularly as therapeutic interventions for infectious disease often need
295 to be given early in the disease course (e.g., empiric antivirals for influenza).

296 There are limitations with our analyses. First, the SARS-CoV-2 genomes uploaded to GISAID
297 are not necessarily representative of all circulating genomes, which can introduce a selection or
298 sampling bias into our analyses based on region, patient severity, or other unmeasured factors. In
299 Supplemental Figure S2 we show the sampling patterns over time by our patient severity
300 categorization. We sought to mitigate these limitations by eliminating the categories that had
301 ambiguous severity (e.g., “live” or “recovered”), and adjusting the associations for known
302 confounders. In addition, we do not seek to make causal claims about any specific viral genomic
303 variant. In aggregate these variants are predictive of outcome and the candidates we identify can
304 be further studied using molecular and other methods.

305 In summary, we have demonstrated that some SARS-CoV-2 genomic variants are strong
306 predictors of COVID-19 disease severity, and these variants appear to be commonly circulating.
307 This study provides a rationale for prioritizing control efforts for cases and populations
308 manifesting with unusually high severity, consistent with symptom-based criteria for testing used
309 by the CDC. Longitudinal monitoring of genomic variants within a novel pathogen such as
310 SARS-CoV-2 will be important for understanding drivers and effects of its evolution and
311 ultimately, its spread or control.

312

313 **References**

- 314 1. Elledge SJ. 2.5 Million Person-Years of Life Have Been Lost Due to COVID-19 in the United States.
315 medRxiv. 2020.
- 316 2. Meyerowitz-Katz G, Merone L. A systematic review and meta-analysis of published research
317 data on COVID-19 infection-fatality rates. medRxiv. 2020.
- 318 3. Ioannidis JP. The infection fatality rate of COVID-19 inferred from seroprevalence data. Bulletin
319 of the World Health Organization. 2020. Epub 14 Oct 2020.
- 320 4. Levin AT, Hanage WP, Owusu-Boaitey N, Cochran KB, Walsh SP, Meyerowitz-Katz G. ASSESSING
321 THE AGE SPECIFICITY OF INFECTION FATALITY RATES FOR COVID-19: SYSTEMATIC REVIEW, META-
322 ANALYSIS, AND PUBLIC POLICY IMPLICATIONS. medRxiv. 2020.
- 323 5. O'Driscoll M, Dos Santos GR, Wang L, Cummings DA, Azman AS, Paireau J, et al. Age-specific
324 mortality and immunity patterns of SARS-CoV-2 infection in 45 countries. medRxiv. 2020.
- 325 6. Onder G, Rezza G, Brusaferro S. Case-fatality rate and characteristics of patients dying in relation
326 to COVID-19 in Italy. *Jama*. 2020;323(18):1775-6.
- 327 7. Horwitz L, Jones SA, Cerfolio RJ, Francois F, Greco J, Rudy B, et al. Trends in Covid-19 risk-
328 adjusted mortality rates in a single health system. *Journal of Hospital Medicine*. 2020. Epub October 23,
329 2020. doi: 10.12788/jhm.3552.
- 330 8. Dennis J, McGovern A, Vollmer S, Mateen BA. Improving COVID-19 critical care mortality over
331 time in England: A national cohort study, March to June 2020. medRxiv. 2020.
- 332 9. Armengaud J, Delaunay-Moisan A, Thuret JY, Van Anken E, Acosta-Alvear D, Aragón T, et al. The
333 importance of naturally attenuated Sars-Cov-2 in the fight against Covid-19. *Environmental*
334 *Microbiology*. 2020;22(6):1997-2000.
- 335 10. Ewald PW. Evolution of virulence. *Infectious disease clinics of North America*. 2004;18(1):1.
- 336 11. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative contribution to
337 global health. *Global Challenges*. 2017;1(1):33-46.
- 338 12. Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data—from vision to reality.
339 *Eurosurveillance*. 2017;22(13):30494.
- 340 13. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;34(18):3094-
341 100.
- 342 14. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for annotating and
343 predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila*
344 *melanogaster* strain w1118; iso-2; iso-3. *Fly*. 2012;6(2):80-92.
- 345 15. Cingolani P, Patel VM, Coon M, Nguyen T, Land SJ, Ruden DM, et al. Using *Drosophila*
346 *melanogaster* as a Model for Genotoxic Chemical Mutational Studies with a New Program, SnpSift.
347 *Frontiers in Genetics*. 2012;3.
- 348 16. McKinney W, editor *Data structures for statistical computing in python*. Proceedings of the 9th
349 Python in Science Conference; 2010: Austin, TX.
- 350 17. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine
351 learning in Python. *the Journal of machine Learning research*. 2011;12:2825-30.
- 352 18. Hunter JD. Matplotlib: A 2D graphics environment. *Computing in science & engineering*.
353 2007;9(3):90-5.
- 354 19. Waskom M, Botvinnik O, Gelbart M, Ostblom J, Hobson P, Lukauskas S. mwaskom/seaborn: v0.
355 11.0 (September 2020); 2020. DOI; 2020.
- 356 20. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0:
357 fundamental algorithms for scientific computing in Python. *Nature methods*. 2020;17(3):261-72.
- 358 21. Zulkower V, Rosser S. DNA Features Viewer, a sequence annotations formatting and plotting
359 library for Python. bioRxiv. 2020.

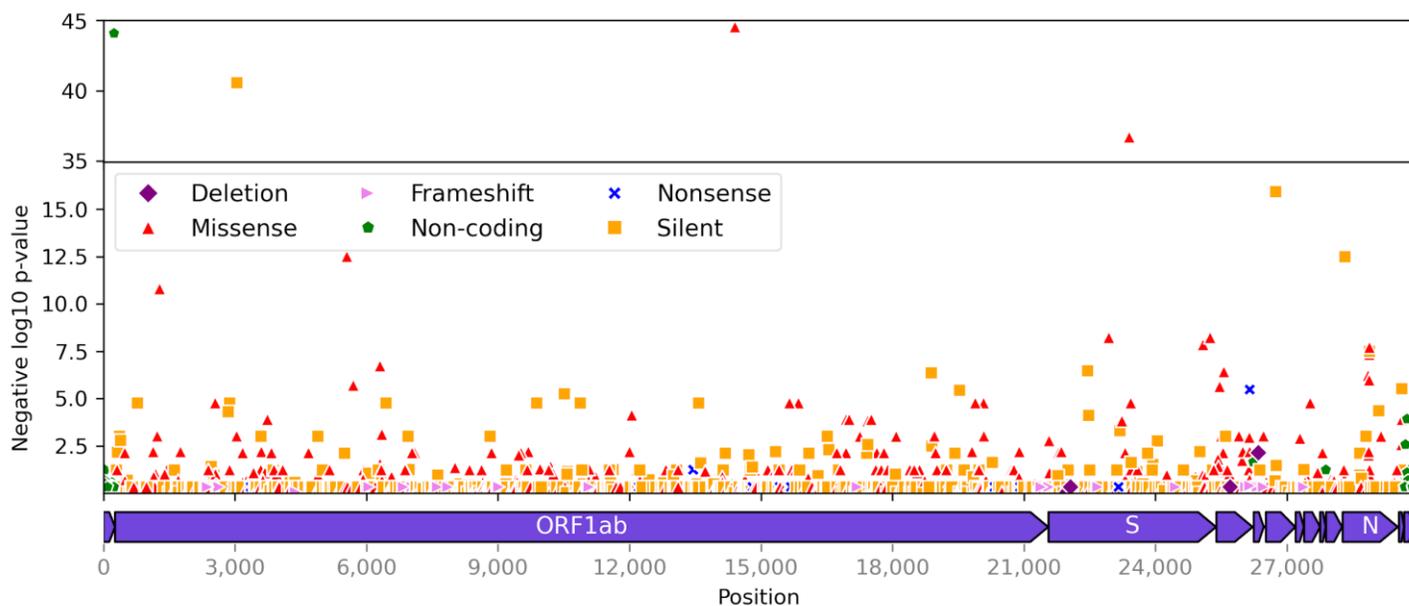
- 360 22. Matsushita K, Ding N, Kou M, Hu X, Chen M, Gao Y, et al. The relationship of COVID-19 severity
361 with cardiovascular disease and its traditional risk factors: A systematic review and meta-analysis. *Global*
362 *heart*. 2020;15(1).
- 363 23. Holmes EC. *The evolution and emergence of RNA viruses*: Oxford University Press; 2009.
- 364 24. Rochman ND, Wolf YI, Koonin EV. Evolution of Human Respiratory Virus Epidemics. medRxiv.
365 2020:2020.11.23.20237503. doi: 10.1101/2020.11.23.20237503.
- 366 25. King KC, Brockhurst MA, Vasieva O, Paterson S, Betts A, Ford SA, et al. Rapid evolution of
367 microbe-mediated protection against pathogens in a worm host. *The ISME journal*. 2016;10(8):1915-24.
- 368 26. Leist SR, Dinnon III KH, Schäfer A, Longping VT, Okuda K, Hou YJ, et al. A Mouse-Adapted SARS-
369 CoV-2 Induces Acute Lung Injury and Mortality in Standard Laboratory Mice. *Cell*. 2020.
- 370 27. Yao H-P, Lu X, Chen Q, Xu K, Chen Y, Cheng L, et al. Patient-derived mutations impact
371 pathogenicity of SARS-CoV-2. CELL-D-20-01124. 2020.
- 372 28. Nagy A, Pongor S, Gyorffy B. Different mutations in SARS-CoV-2 associate with severe and mild
373 outcome. medRxiv. 2020.
- 374 29. Hahn G, Wu CM, Lee S, Hecker J, Lutz SM, Haneuse S, et al. Mutations in SARS-CoV-2 spike
375 protein and RNA polymerase complex are associated with COVID-19 mortality risk. bioRxiv. 2020.
- 376 30. Berrio A, Gartner V, Wray GA. Positive selection within the genomes of SARS-CoV-2 and other
377 Coronaviruses independent of impact on protein function. *PeerJ*. 2020;8:e10234.
- 378 31. Velazquez-Salinas L, Zarate S, Eberl S, Gladue DP, Novella I, Borca MV. Positive Selection of
379 ORF1ab, ORF3a, and ORF8 Genes Drives the Early Evolutionary Trends of SARS-CoV-2 During the 2020
380 COVID-19 Pandemic. *Frontiers in Microbiology*. 2020;11:2592.
- 381 32. Dilucca M, Forcelloni S, Georgakilas AG, Giansanti A, Pavlopoulou A. Codon Usage and
382 Phenotypic Divergences of SARS-CoV-2 Genes. *Viruses*. 2020;12(5):498.
- 383 33. Simmonds P. Rampant C→U hypermutation in the genomes of SARS-CoV-2 and other
384 coronaviruses: Causes and consequences for their short-and long-term evolutionary trajectories.
385 *MSphere*. 2020;5(3).
- 386 34. Matyášek R, Kovařík A. Mutation patterns of human SARS-CoV-2 and Bat RaTG13 coronavirus
387 genomes are strongly biased towards C>U transitions, indicating rapid evolution in their hosts. *Genes*.
388 2020;11(7):761.
- 389 35. Di Giorgio S, Martignano F, Torcia MG, Mattiuz G, Conticello SG. Evidence for host-dependent
390 RNA editing in the transcriptome of SARS-CoV-2. *Science Advances*. 2020:eabb5813.
- 391 36. Subissi L, Posthuma CC, Collet A, Zevenhoven-Dobbe JC, Gorbalenya AE, Decroly E, et al. One
392 severe acute respiratory syndrome coronavirus protein complex integrates processive RNA polymerase
393 and exonuclease activities. *Proceedings of the National Academy of Sciences*. 2014;111(37):E3900-E9.
- 394 37. Farkas C, Mella A, Haigh JJ. Large-scale population analysis of SARS-CoV2 whole genome
395 sequences reveals host-mediated viral evolution with emergence of mutations in the viral Spike protein
396 associated with elevated mortality rates. medRxiv. 2020.
- 397 38. Turoňová B, Sikora M, Schürmann C, Hagen WJ, Welsch S, Blanc FE, et al. In situ structural
398 analysis of SARS-CoV-2 spike reveals flexibility mediated by three hinges. *Science*. 2020;370(6513):203-
399 8.
- 400 39. Gibbons TF, Storey SM, Williams CV, McIntosh A, Mitchel DM, Parr RD, et al. Rotavirus NSP4: Cell
401 type-dependent transport kinetics to the exofacial plasma membrane and release from intact infected
402 cells. *Virology journal*. 2011;8(1):1-20.
- 403 40. Tan Y-J. The Severe Acute Respiratory Syndrome (SARS)-coronavirus 3a protein may function as
404 a modulator of the trafficking properties of the spike protein. *Virology journal*. 2005;2(1):1-5.
- 405 41. Siu KL, Yuen KS, Castano-Rodriguez C, Ye ZW, Yeung ML, Fung SY, et al. Severe acute respiratory
406 syndrome Coronavirus ORF3a protein activates the NLRP3 inflammasome by promoting TRAF3-
407 dependent ubiquitination of ASC. *The FASEB Journal*. 2019;33(8):8865-77.

- 408 42. Issa E, Merhi G, Panossian B, Salloum T, Tokajian S. SARS-CoV-2 and ORF3a: Nonsynonymous
409 Mutations, Functional Domains, and Viral Pathogenesis. *Msystems*. 2020;5(3).
- 410 43. Lu X, Wang L, Sakthivel SK, Whitaker B, Murray J, Kamili S, et al. US CDC real-time reverse
411 transcription PCR panel for detection of severe acute respiratory syndrome coronavirus 2. *Emerging*
412 *infectious diseases*. 2020;26(8):1654.
- 413 44. Oulas A, Zanti M, Tomazou M, Zachariou M, Minadakis G, Bourdakou MM, et al. Generalized
414 linear models provide a measure of virulence for specific mutations in SARS-CoV-2 strains. *bioRxiv*. 2020.
- 415 45. Geoghegan JL, Holmes EC. The phylogenomics of evolving virus virulence. *Nature Reviews*
416 *Genetics*. 2018;19(12):756-69.
- 417 46. Grasselli G, Pesenti A, Cecconi M. Critical care utilization for the COVID-19 outbreak in
418 Lombardy, Italy: early experience and forecast during an emergency response. *Jama*.
419 2020;323(16):1545-6.
- 420 47. White DB, Lo B. A framework for rationing ventilators and critical care beds during the COVID-19
421 pandemic. *Jama*. 2020;323(18):1773-4.
- 422 48. Chin V, Samia NI, Marchant R, Rosen O, Ioannidis JP, Tanner MA, et al. A case study in model
423 failure? COVID-19 daily deaths and ICU bed utilisation predictions in New York State. *European Journal*
424 *of Epidemiology*. 2020;35(8):733-42.

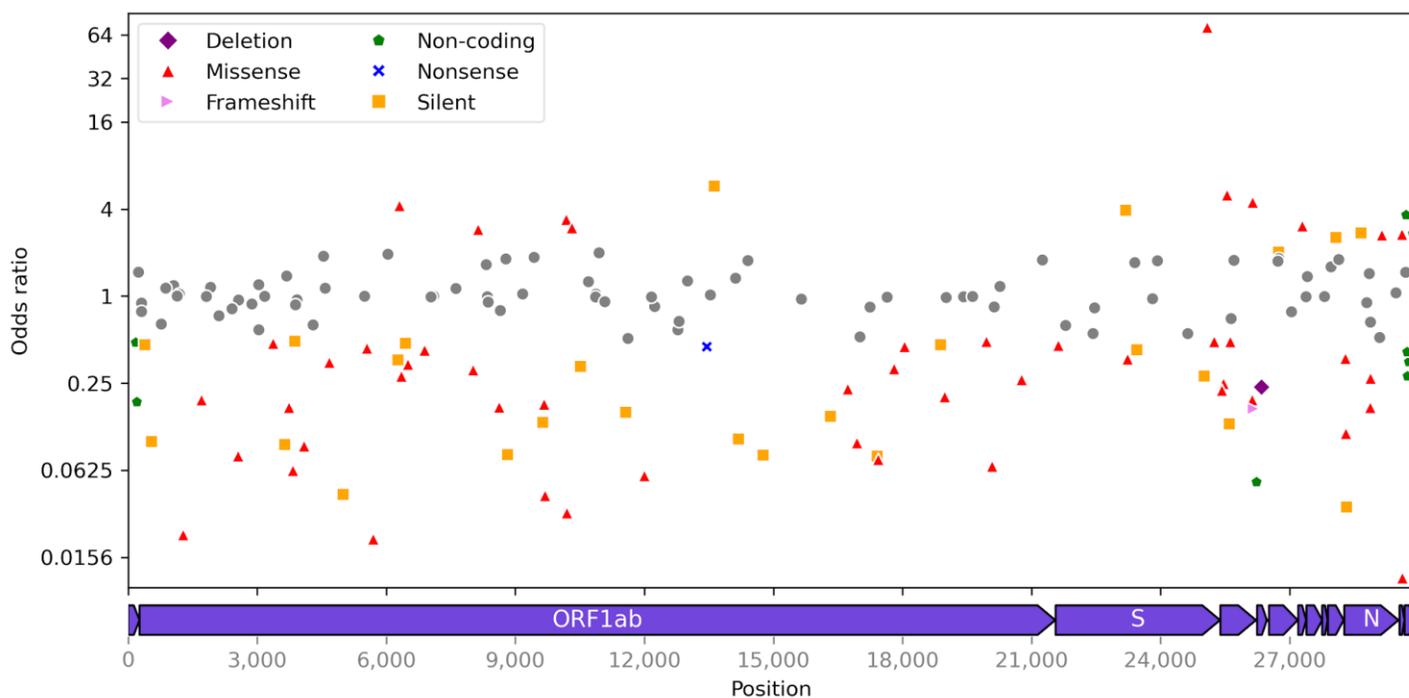
425

426 **Figures**

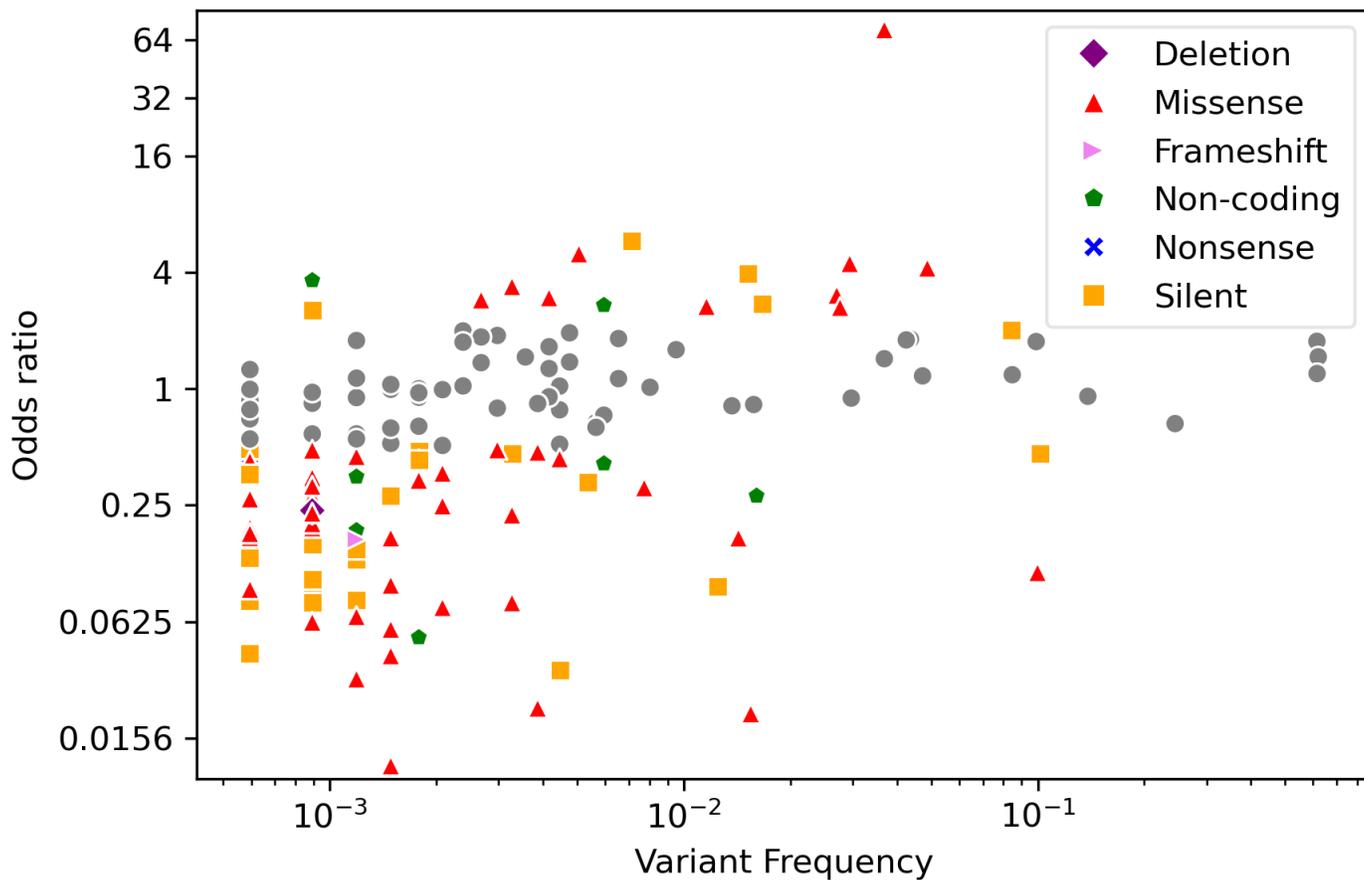
427 **Figure 1A:**



429 **Figure 1B:**

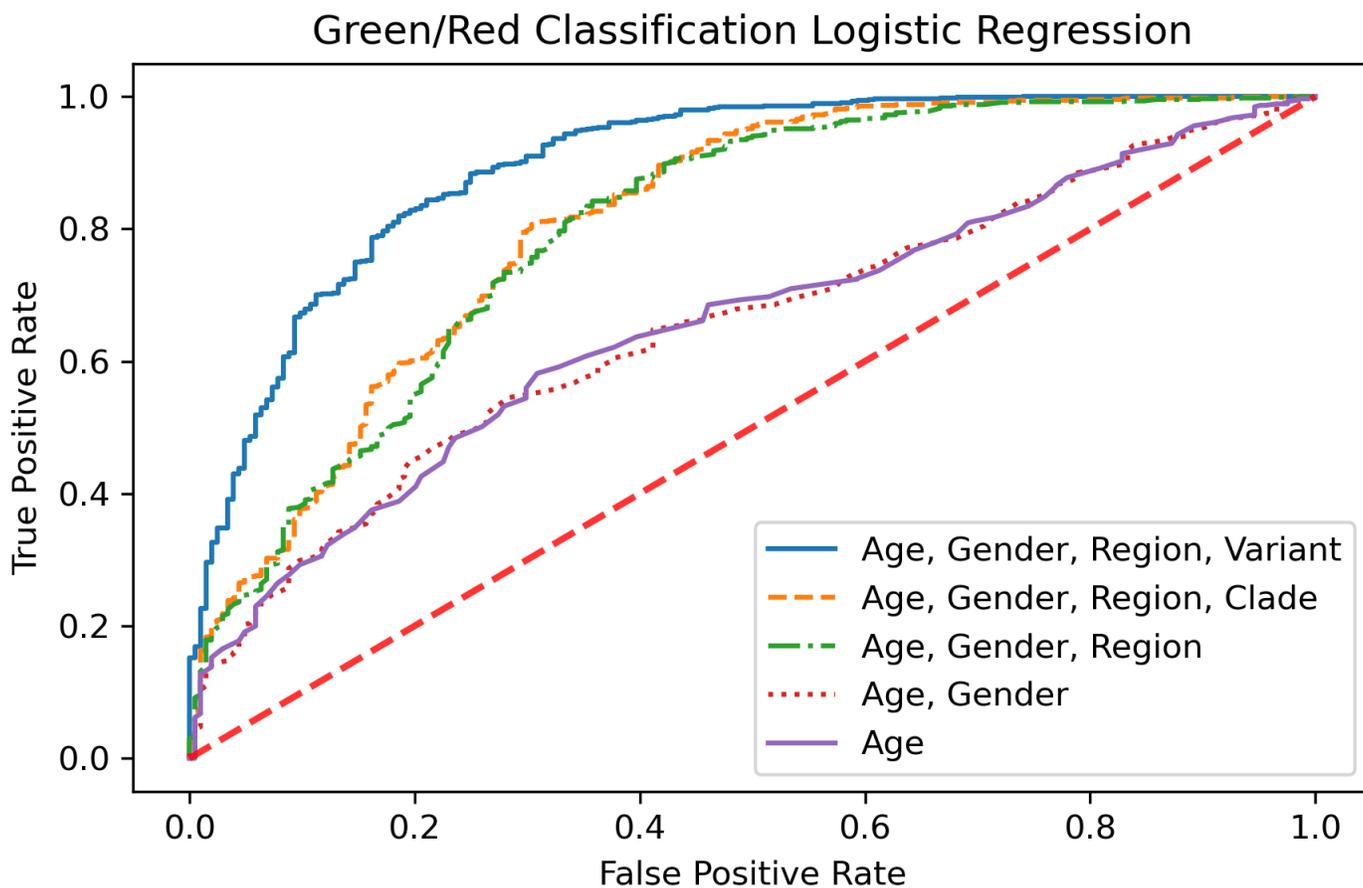


432 **Figure 1C:**



434

435 **Figure 2:**



437

438 **Figure Legends**

439 **Figure 1:** Overview of SARS-CoV-2 variants selected from GISAID data (n =
440 1595168).

441 A) Negative log₁₀ p-values of variant association (chi-square test) with “Severe”
442 outcome group (hospitalized, deceased, etc.) plotted against position of variants
443 (n = 4484) in the SARS-CoV-2 genome.

444 B) Odds ratios (log₂ scale) of “Severe” versus “Mild” (outpatient, asymptomatic,
445 etc.) outcome groups plotted against the positions of variants with odds ratios
446 not equal to one (n = 169) in the SARS-CoV-2 genome.

447 C) Odds ratios (log₂ scale) of “Severe” versus “Mild” outcome groups plotted
448 against log₁₀ frequency of variants (n = 169) in the patient subpopulation (n =
449 3363) without missing variables.

450 Points are labeled by mutation type (red: missense, green: non-coding, orange:
451 silent, yellow: nonsense, purple: deletion).

452

453 **Figure 2:** Comparison of nested logistic regression models.

454 Models are labeled based on the predictor variables (purple solid line: age; red
455 dotted line: [age, gender], green dash-dotted line: [age, gender, region], orange
456 dashed line: [age, gender, region, clade], blue solid line: [age, gender, region,
457 variant]) used to predict whether SARS-CoV-2 patients (n = 3363) belong to
458 “Severe” (hospitalized, deceased) or “Mild” (outpatient, asymptomatic, etc.)
459 outcome groups.

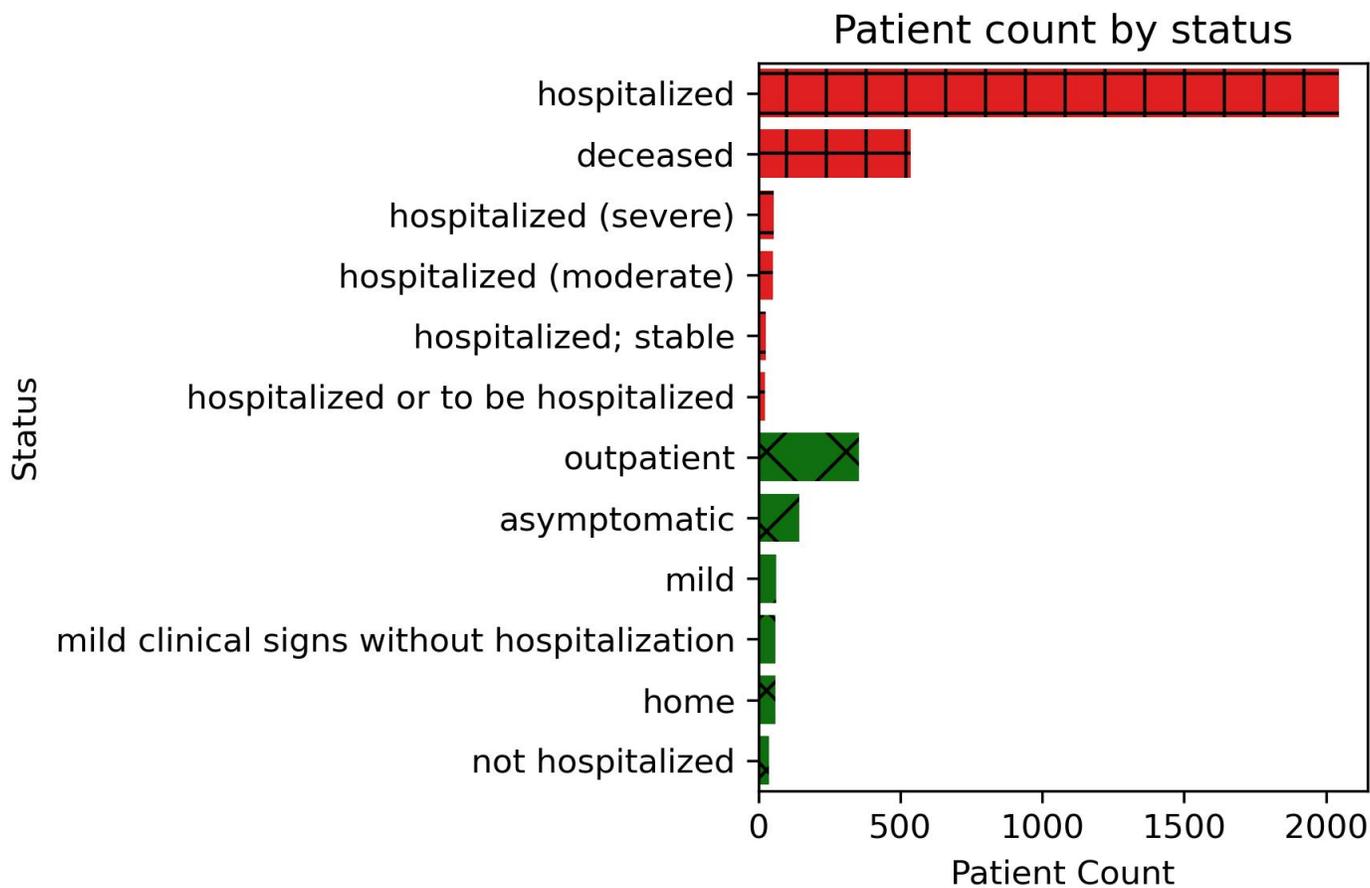
460

461

462 **Supplemental Figures**

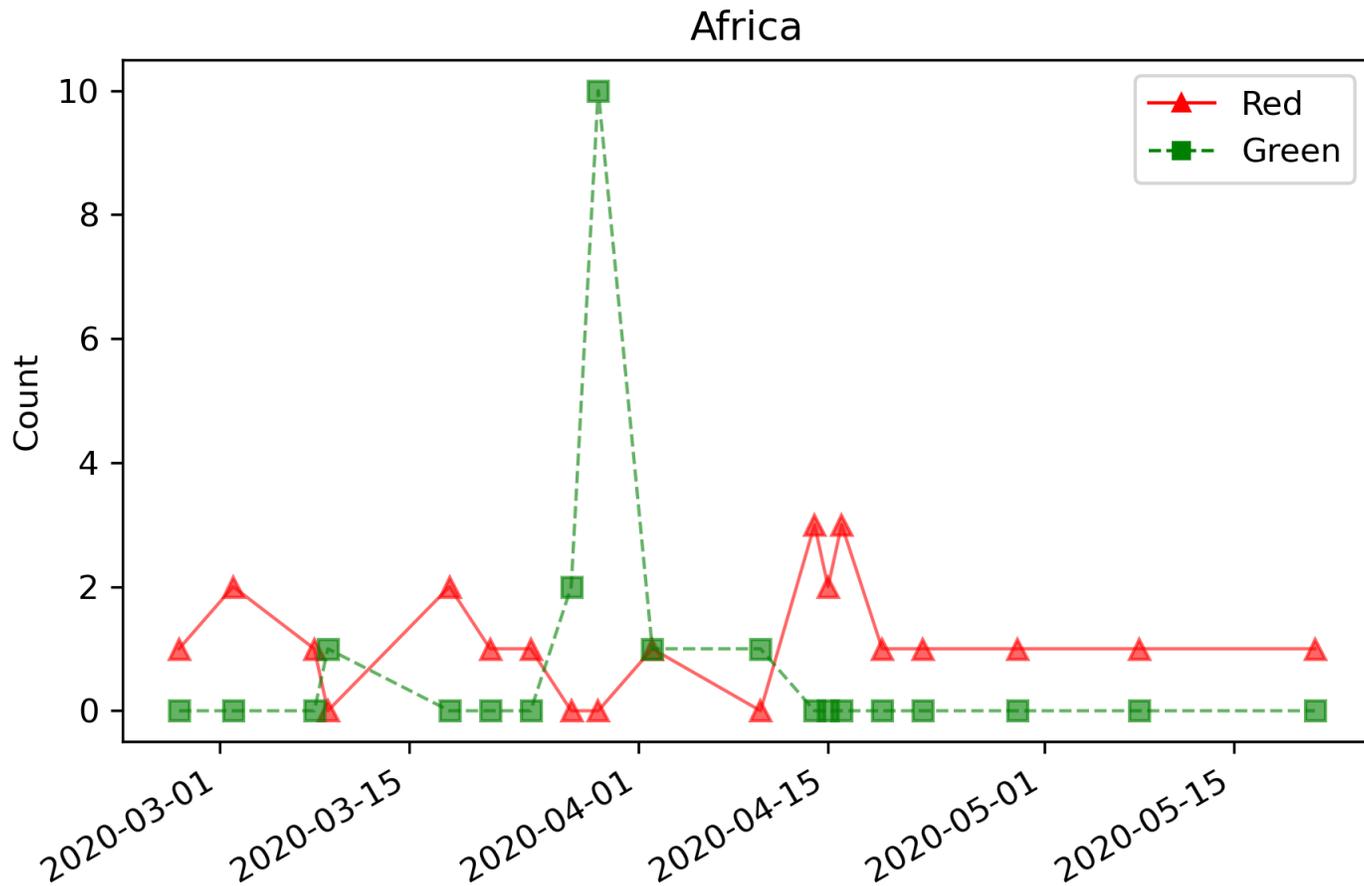
463

464 **Figure S1:**

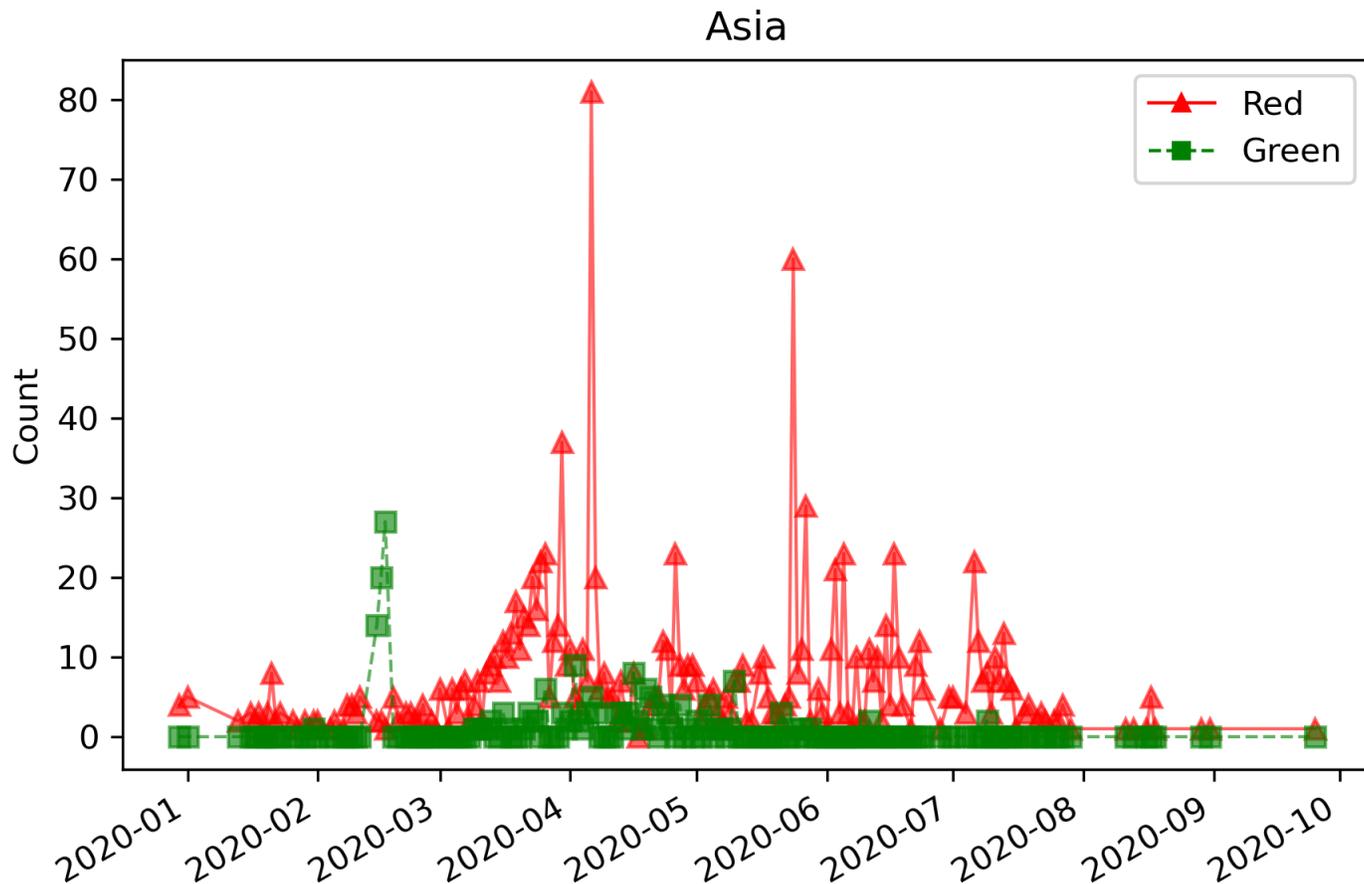


466

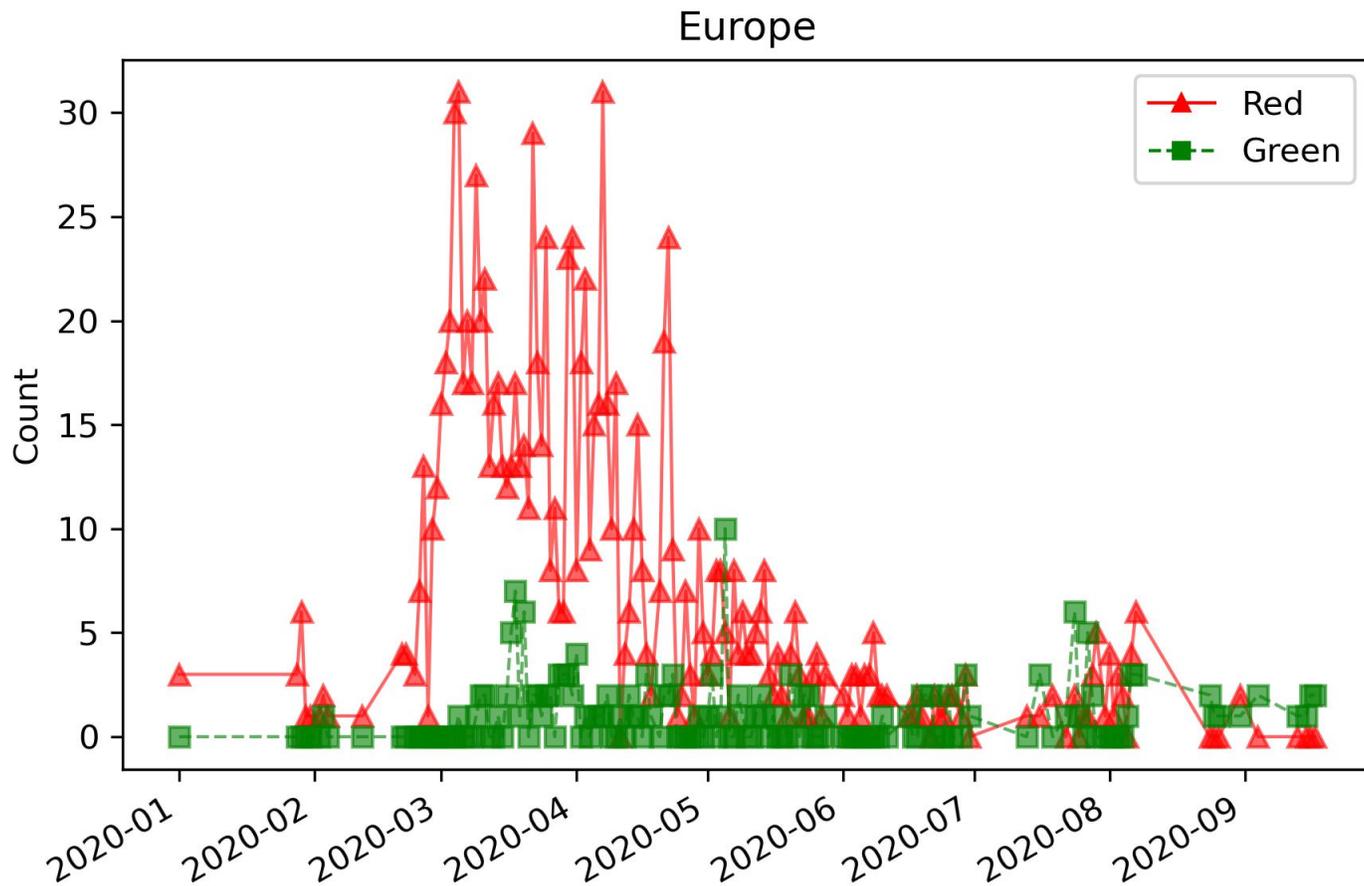
467 **Figure S2A:**



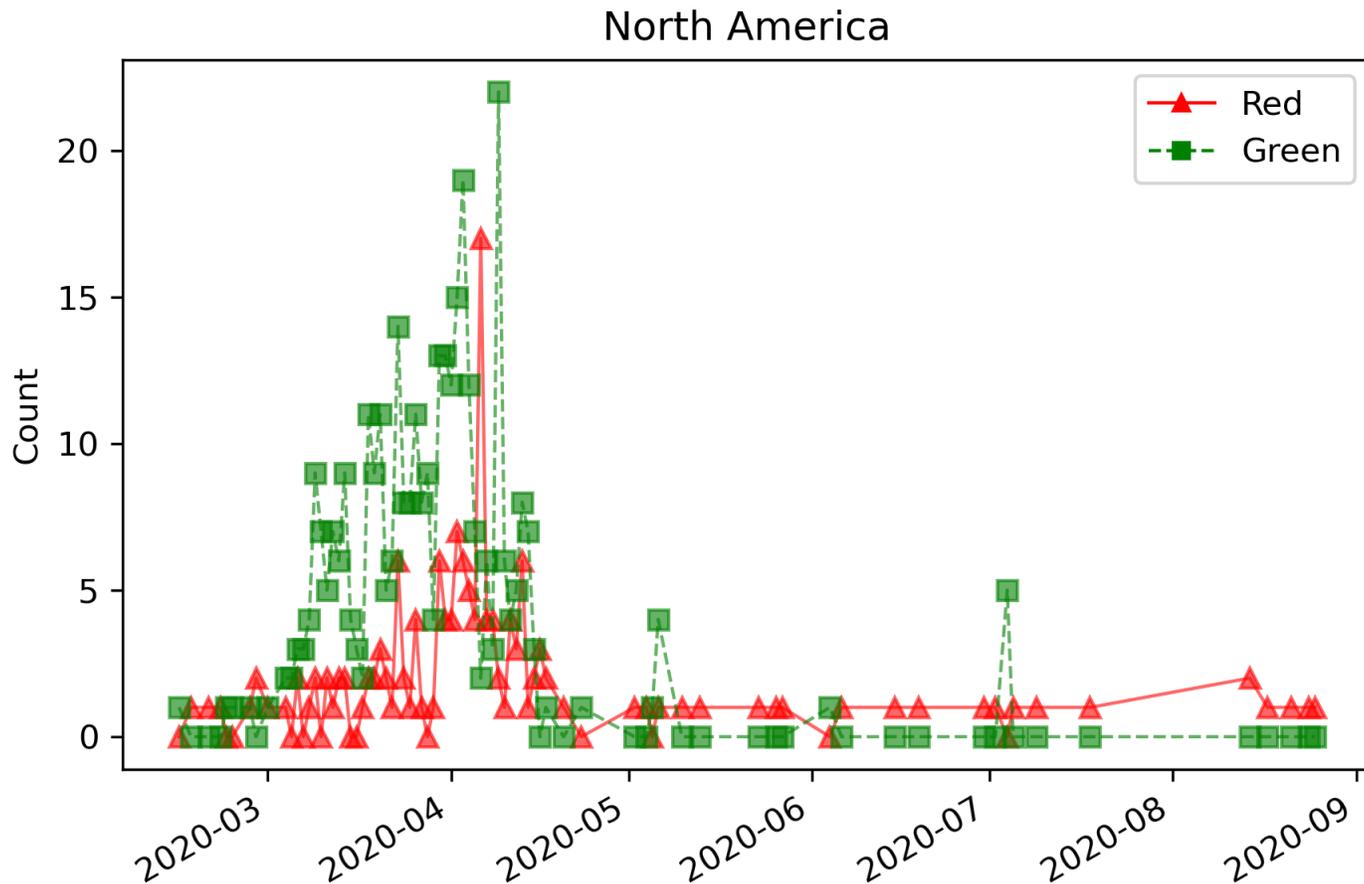
470 **Figure S2B:**



472 **Figure S2C:**

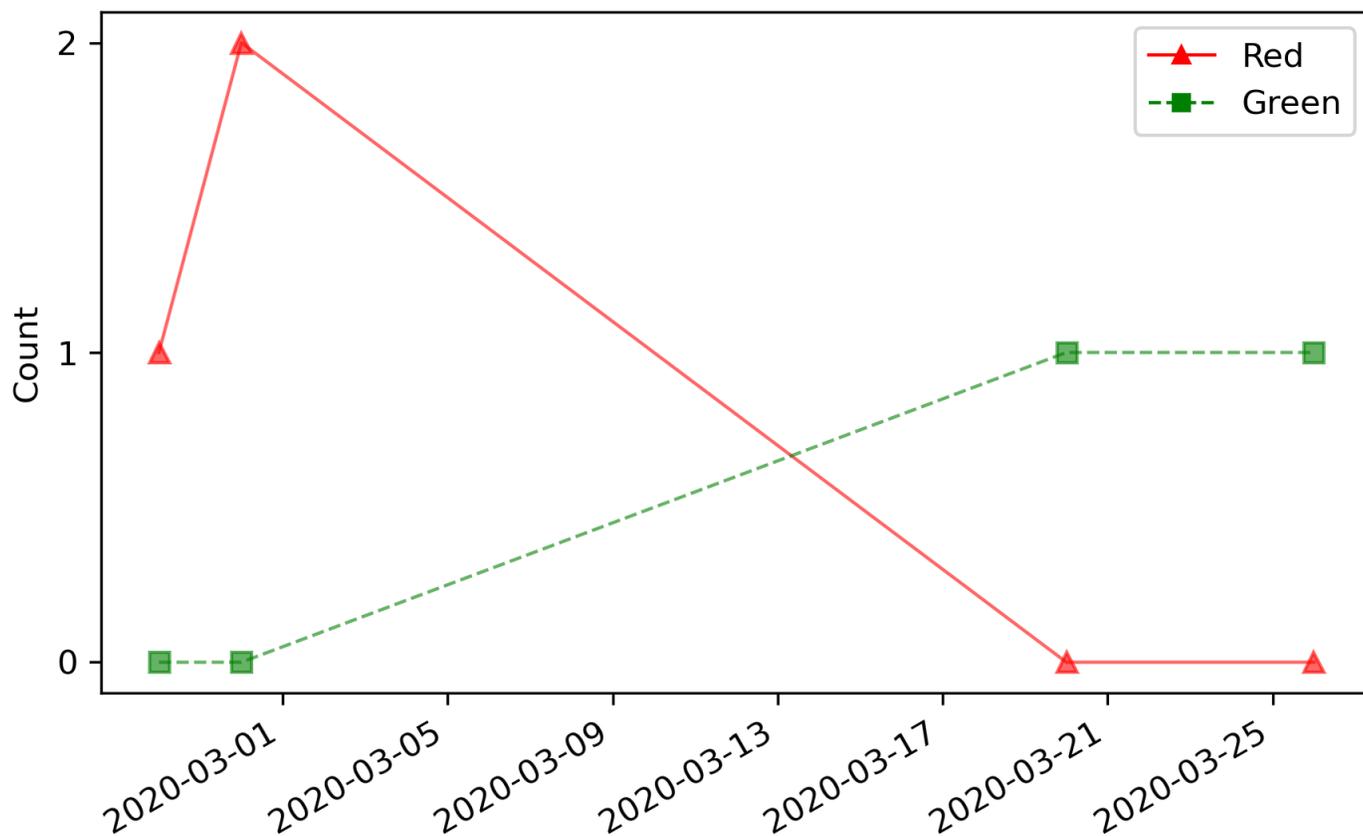


475 **Figure S2D:**



477 **Figure S2E:**

Oceania

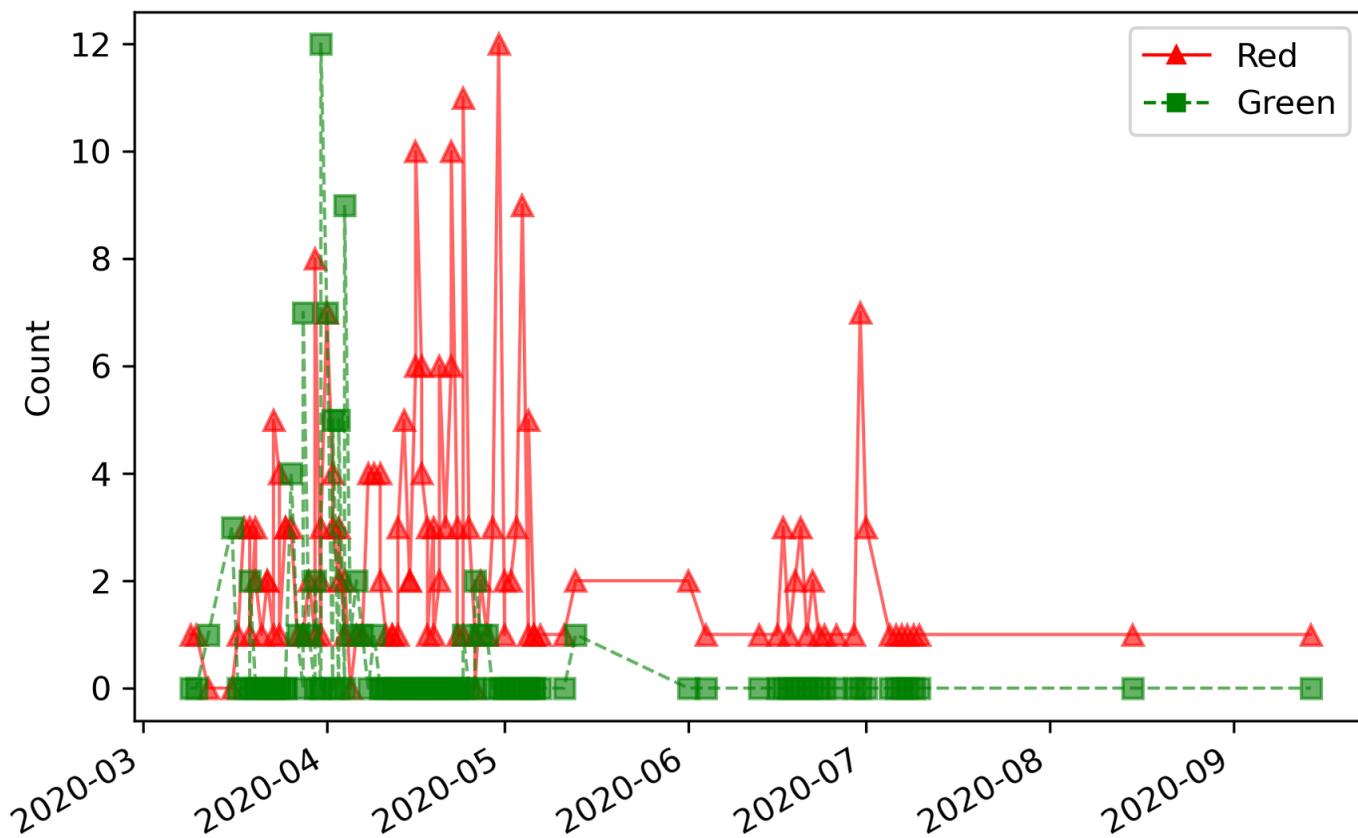


479

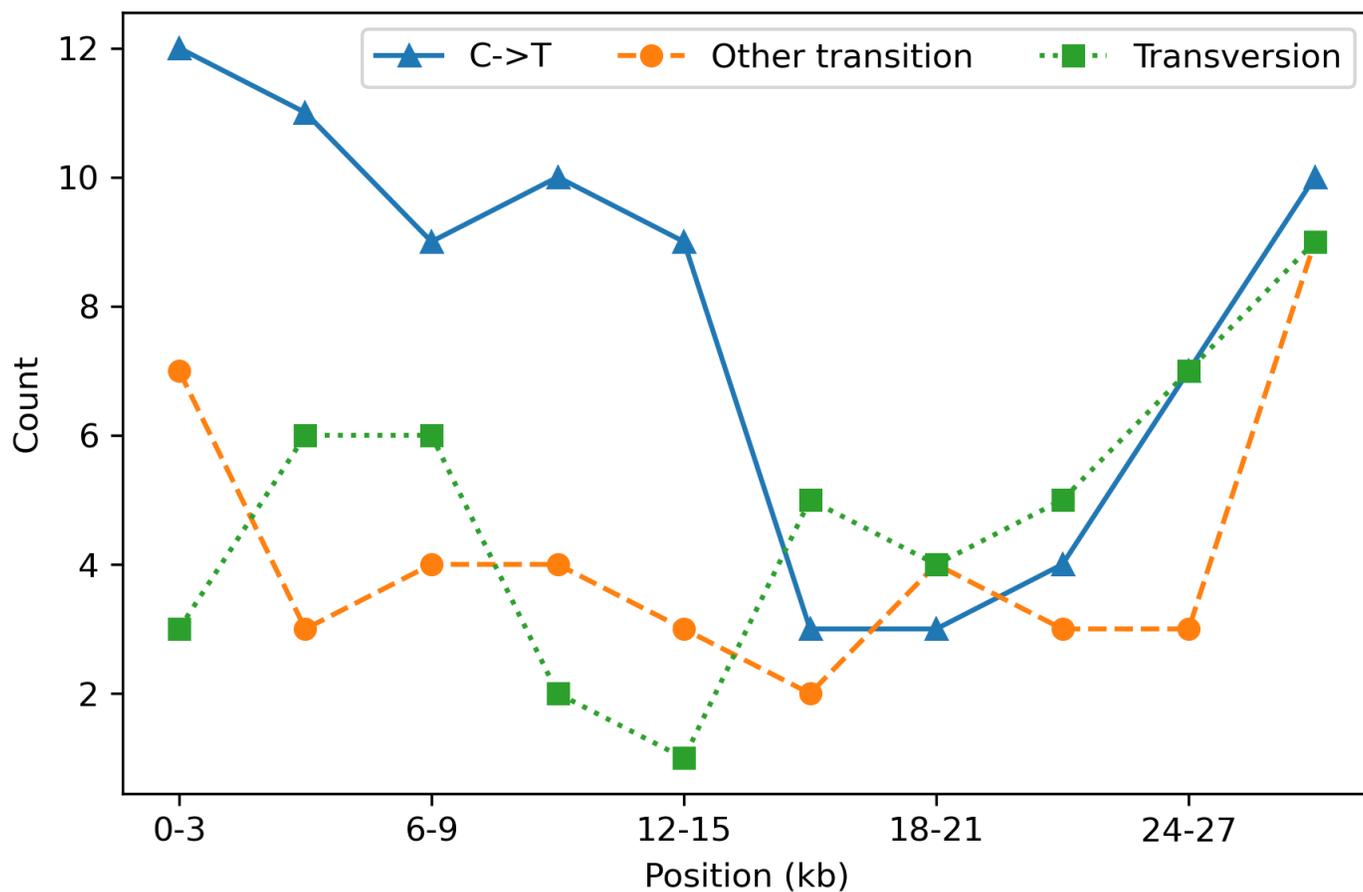
480

Figure S2F:

South America

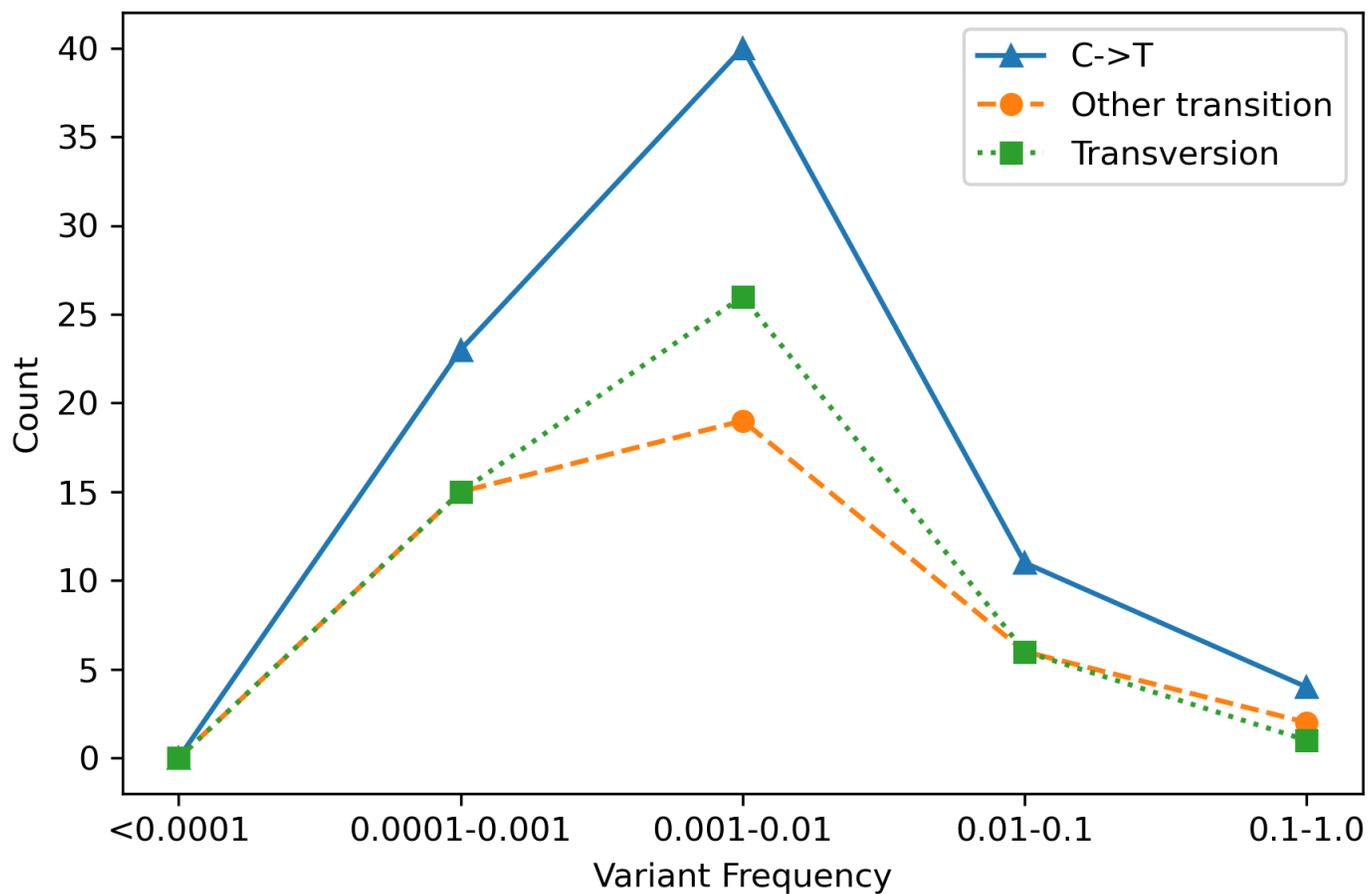


483 **Figure S3A:**



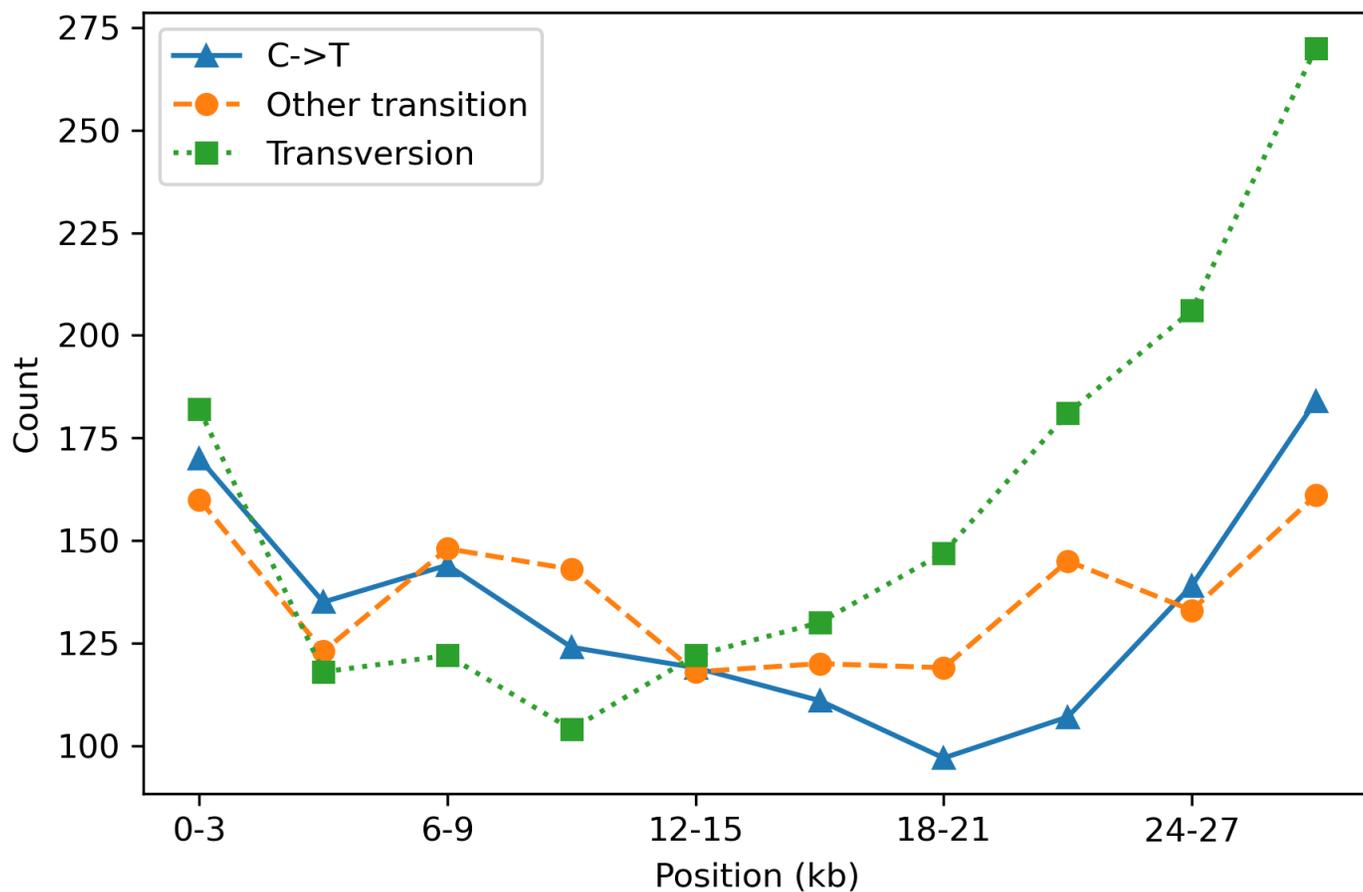
485

486 **Figure S3B:**



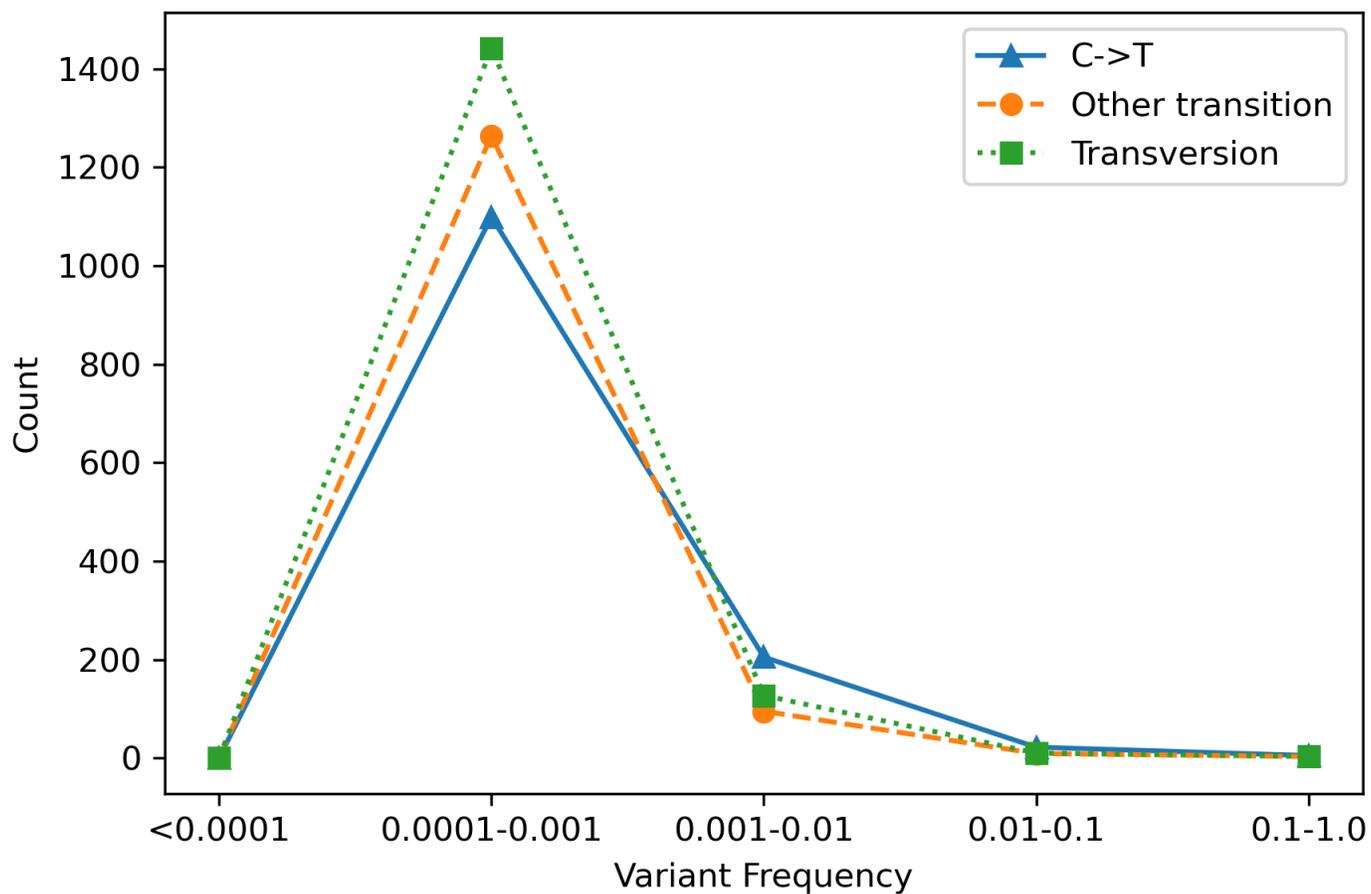
488

489 **Figure S3C:**



491

492 **Figure S3D:**



494

495 **Supplemental Figure Legends**

496 **Figure S1:** SARS-CoV-2 patient (n = 3611) composition of “Severe” and “Mild”
497 outcome groups.

498 Bars are labeled by group membership (red vertical-horizontal hatch: [severe,
499 symptomatic, deceased], green diagonal hatch: [live, released, mild, recovered,
500 asymptomatic]).

501 **Figure S2:** “Severe” and “Mild” outcome counts (n = 3611) over time for all
502 GISAID regions.

503 A) Africa (n = 37)

504 B) Asia (n = 1454)

505 C) Europe (n = 1265)

506 D) North America (n = 499)

507 E) Oceania (n = 5)

508 F) South America (n = 351)

509 Curves are labeled by outcome group (red solid line: “Severe” outcome, green
510 dotted line: “Mild” outcome).

511 Dates (x-axis) are shown in YYYY-MM or YYYY-MM-DD format.

512 **Figure S3:** Number of C to T transitions (C->T) compared to other mutation
513 types.

514 A) The counts of variants with odds ratios not equal to one (n=170) plotted
515 against log₁₀ variant frequency half-open intervals, e.g. [0.01, 0.1).

516 B) The counts of variants with odds ratios not equal to one (n=170) plotted
517 against half-open intervals of 3 kilobases (kb), e.g. [0, 3000) in the SARS-CoV-2
518 genome.

519 C) The counts of variants used for logistic regression modeling (n=4484) plotted
520 against log₁₀ variant frequency half-open intervals, e.g. [0.01, 0.1).

521 D) The counts of variants used for logistic regression modeling (n=4484) plotted
522 against half-open intervals of 3 kilobases (kb), e.g. [0, 3000) in the SARS-CoV-2
523 genome.

524 Curves are labeled based on mutation type (blue solid line: C->T transition, orange dashed line:
525 other transition [C->T, T->C, A->G, G->A], green dotted line: transversion [C->A, A->C, T->G, G-
526 >T, C->G, G->C, A->T, T->A]).