

Drivers of COVID-19 variant wave dynamics: inferring oncoming wave size using global data with genomics

S. Molan¹, N.K. Smith^{2,3}, V. Gandhi¹, M. Li^{4,5}, C. Colijn¹, C.L. Murall³, and J.E.
Stockdale^{1,*}

¹Department of Mathematics, Simon Fraser University

²Department of Community Health and Epidemiology, Dalhousie University

³National Microbiology Laboratory Branch, Public Health Agency of Canada

⁴Department of Mathematics and Statistics, McMaster University

⁵Science Policy Integration Branch, Public Health Agency of Canada

*Corresponding author: 8888 University Drive, Burnaby, BC, V5A 1S6,

jessica_stockdale@sfu.ca

Abstract

The continued evolution of the SARS-CoV-2 virus drove waves of infection worldwide throughout the pandemic. These evolutionary dynamics posed significant challenges for public health forecasting and, specifically, for predicting the size of COVID-19 waves. In this work we leverage a range of global public data, with a focus on features derived from pathogen genomic sequences, to model and predict the relative size of COVID-19 waves (as compared to the previous wave) across countries. Focusing on Omicron BA.1 and BA.2, we develop statistical models to assess the predictive power of these data in forecasting future variant-driven wave peaks. We find that, while forecasting wave size is a challenging task, variables such as genomic variant characteristics, prior wave dynamics, and demographic features e.g. life expectancy were informative, whereas seasonality was not. Our results show that the importance of features changed markedly between Omicron waves, reflecting the evolving epidemiological and genomic landscape. This work provides insights into improving predictive models for future outbreaks and pandemics, and prioritizing data collection efforts to enhance forecasting accuracy.

Keywords: COVID-19; Prediction; Omicron variants; Genomic data; Random forest

1 Introduction

Since the appearance of SARS-CoV-2 in late 2019, heterogeneity in public health response, vaccine uptake, health care system capacity, and the continued evolution of more transmissible and immune-evasive variants created variable transmission dynamics, with multiple waves of COVID-19 cases per year in most countries. Many respiratory viruses have clear seasonal wave dynamics, with increased transmission and activity during the colder months, but the wave dynamics of COVID-19 appear to have had faster persistent oscillations. COVID-19 waves have been less predictable in timing or magnitude, posing a challenge for public health preparation and response.

25 During the pandemic, surveillance and data collection efforts made unprecedented amounts of
26 global public data available, including epidemiological, genomic, clinical, serological, mobility,
27 and demographic data. There were many efforts to collect and standardize population-level data
28 for public use. For example, the COVID-19 Data Repository by the Center for Systems Sci-
29 ence and Engineering (CSSE) at Johns Hopkins University collected basic epidemiological data
30 such as reported cases, hospitalizations and mortality from all countries around the world¹; sim-
31 ilarly the World Health Organization (WHO) collated global epidemiological data²; Apple and
32 Google mobility released population-level aggregated movement data publicly to support COVID-
33 19 research^{3,4}; GISAID hosted SARS-CoV-2 genomic sequences⁵; SeroTracker collated the re-
34 sults of sero studies globally⁶; the Oxford COVID-19 Government Response Tracker reported on
35 the changing levels of government pandemic policy and interventions⁷; and Our World in Data
36 (OWID)⁸ collated and shared data from many of the above sources. Compared to any historical
37 disease outbreak, there is a very large volume of publicly available data for COVID-19.

38 Despite the dissimilarities of wave dynamics between countries, the emergence of the Omicron
39 variant (B.1.1.529) in late 2021 was relatively consistent around the world. In a span of approxi-
40 mately three months, the Omicron wave overwhelmed the testing and epidemiological data collec-
41 tion capacity of many jurisdictions. Due to its significantly higher rate of transmission and strong
42 immune evasive properties, Omicron and its descendants completely replaced the previous variants
43 that were more genetically similar to the wildtype. Omicron continued to evolve immune evasive
44 variants with very successful lineages that circulated year-round. With hundreds of circulating
45 lineages at any given time, monitoring and identifying the next fittest variant that would drive an
46 uptick in cases or evade currently used drugs or vaccines became a collective non-stop effort by
47 scientists around the world⁹⁻¹¹. National and local public health authorities monitored the activity
48 of COVID-19 in their jurisdictions and compared to other countries (e.g. neighbouring countries
49 or countries that had a head start in the emergence of the latest variant of interest) to evaluate the
50 magnitude of the impact of the next wave. These wave dynamic comparisons were particularly im-
51 portant at the height of the pandemic, as decisions on whether to implement or release public health

52 measures were responsive and critical. However, despite the large volumes of real-time globally
53 shared data, this predictive work was done largely using human intuition that tried to synthesize
54 a myriad of likely predictive factors (current season, differences in population vaccination status,
55 timing of previous wave, etc.) along with the various variant characteristics and current epidemi-
56 ological trends and thus, understandably, fell short of reliably predicting important indicators of
57 impact, such as the relative size of the next local wave.

58 In this work, we collate a wide range of publicly available data, including epidemiological, ge-
59 nomic, clinical, serological, mobility, and demographic data, to explore the wave dynamics of
60 Omicron BA.1 and BA.2 in a variety of countries. In particular, we focus on the inclusion of fea-
61 tures derived from global SARS-CoV-2 genomic sequences. We seek to determine whether the
62 relative size of an upcoming case peak (BA.1 or BA.2, compared to a recent peak) could be pre-
63 dicted using a statistical model and data up until the initial growth phase of the respective variant.
64 We also seek to identify which of the features were the most informative of the relative wave peak
65 sizes, with a focus on the genomic features derived from global sequence data. Alongside this
66 manuscript, we make available our curated and cleaned dataset which collates this wide range of
67 features, from December 2019 until January 2023, across 248 countries and jurisdictions.

68 **2 Materials and Methods**

69 **2.1 Prediction task**

70 We focus on predicting the peak size of a new variant-driven wave relative to the peak size of
71 the previous wave in a country. We study two waves: Omicron BA.1 relative to Delta and Omi-
72 cron BA.2 relative to BA.1, which we refer to as BA.1/Delta and BA.2/BA.1, respectively. These
73 periods corresponded to a time of major change in the pandemic, with pronounced increases in
74 transmission. Epidemiological and genomic data sources were still regularly collected at this time,

75 along with data on COVID-19 testing and vaccination/booster campaigns. In the prediction task,
76 our response variable is the relative size of the maximum case counts attributed to the emerging
77 variant of concern (VOC) vs. the maximum VOC-attributed case counts in the previous wave (a
78 different dominant variant). For example, the relative size of BA.2 is the ratio of the peak cases
79 of BA.2 and of BA.1. Compared to the task of predicting absolute peak sizes, this minimizes the
80 impact of variation in the levels of case detection in different countries. Although case count time
81 series are impacted by changes in at-home test reporting, access to testing and so on, it remains
82 the case that reported infections were one of the primary ways countries assessed their COVID-19
83 burden, and of high interest for public health decision makers. We seek to assess how routinely
84 collected data, particularly genomic sequence data, could serve as an early indicator of the relative
85 wave size, for which case counts are the earliest measure.

86 In the sections below, we first describe the full collated dataset, for all available countries from the
87 start of the pandemic until the end of January 2023. We then describe the subset of the data that is
88 used in our predictive model.

89 **2.2 Data**

90 We collated publicly available data on variables that could be expected to affect the size of COVID-
91 19 waves across countries. Data were obtained from online repositories including Our World in
92 Data (OWID)⁸, which collected information from sources such as the Johns Hopkins University
93 COVID-19 data repository¹, the WHO database², and the Oxford COVID-19 Government Re-
94 sponse Tracker⁷. Additional data came from SeroTracker⁶, Google Mobility³, and GISAID⁵. Some
95 variables were taken directly from these sources (e.g., case counts), and some we derived from the
96 available data, as will be described in the data processing section. Where available, we collected
97 a daily time series of each variable for each country (a total of 248 countries and jurisdictions)
98 throughout a study window of 2019-12-29 to 2023-01-29.

99 We sourced global viral genomic data of human origin (not animal or environmental) from GI-
100 SAID⁵ on May 4, 2023. Sequences were filtered to include only those that had complete collection
101 dates and were submitted on or before February 1, 2023. The variants were designated using
102 Pangolin version 4.2¹², PUSHER-v1.18.1.1. We then classified Pango lineages into variants of
103 concern (VOC) and variants of interest (VOI) groups that align with the Nextclade clade classifi-
104 cation system (a higher-order grouping of successful lineages^{13,14}). Full details on the VOC and
105 VOI classifier function are available in the GitHub repository associated with this work, as is a
106 supplemental table describing the 1,048,516 sequences used, available on GISAID and accessible
107 at doi.org/10.55876/gis8.241023pk. To infer the dates of global emergence of each
108 clade, a method that uses the time to the most recent common ancestor (tMRCA)¹⁵ was applied to
109 each Nextstrain clade. A curated date for Omicron BA.1 was taken from Viana et al.¹⁶. To avoid
110 sequences with data entry errors, we removed any sequences dated prior to the inferred date of
111 global emergence of that particular lineage's clade. We collapsed the genomic metadata dataset
112 into a weekly format from Monday to Sunday, with the week's total counts being summed on
113 Sunday. To obtain mutation and amino acid substitution calls for each sequence, the global viral
114 sequence dataset was run through Nextclade CLI^{17,18} v. 3.4.0 that aligns the genomes relative to
115 the wildtype (Wuhan-Hu-1/2019 (MN908947)), assigns clades, runs quality checks and gives mu-
116 tational information relative to the original Wuhan reference. To minimize noise in the data for
117 the predictive modelling task, upon combining genomic data, we converted all data to a weekly
118 timescale, in which weeks run from Monday to Sunday.

119 In total, we collated data on 113 variables, which we group into 11 broad categories: COVID-19
120 cases and reproduction rates, deaths, excess mortality, hospitalizations, policy responses, testing,
121 vaccination, serology studies, genomic variables, demographic, and miscellaneous. The final cat-
122 egory includes variables such as mobility data, season, and latitude/longitude. The full list of
123 variables is shown in Figure ???. The full dataset, along with variable descriptions and sources is
124 available in the GitHub repository associated with this work.

125 These data draw from many different sources and are subject to many limitations, including miss-
126 ing data (almost certainly not at random), changes in data collection policy over time, differences
127 in definitions across countries, and more. Despite this, these data represent perhaps our most
128 comprehensive international understanding of the burden of COVID-19, the immunity against it,
129 and the viral evolution of SARS-CoV-2 that is publicly available, so, noting these limitations, we
130 proceed.

131 **2.3 Data processing**

132 To summarize the genomic data in a variant-specific way, we derived additional variables, includ-
133 ing VOC-specific case counts, a measure of sequencing effort, genetic diversity and distance, and
134 growth advantage.

135 We estimated the case counts for each VOC by multiplying the total weekly case counts (from
136 OWID⁸) by the proportion of the total VOC sequences (from GISAID⁵). This generated a time
137 series of weekly case counts per VOC per country. During the emergence of each new variant, we
138 calculated the time (number of days) since the previous variant's peak. We created a sequencing
139 effort variable to capture the variation in case testing and genomic sequencing between and within
140 countries throughout the pandemic. We defined sequencing effort as the total weekly number of
141 sequences divided by the total weekly number of cases in each country.

For each country, we calculated several diversity metrics. The simplest was richness, defined as the total number of unique Pango lineages detected per week. We also included the number of distinct VOCs present as a separate variable. To account for the relative abundance of lineages, we calculated Hill diversity of order 1 (D_1), equivalent to the exponentiated Shannon index (H), which incorporates both lineage richness and evenness, providing a more complete picture of diversity

than richness.¹⁹ The Shannon index was calculated using the *vegan* package²⁰ with the formula:

$$H(t) = - \sum_i p_i(t) \log(p_i(t))$$

142 where $p_i(t)$ is the proportion of sequences belonging to the Pango lineage i at time t . The Hill
143 diversity was then calculated as $D_1(t) = \exp(H(t))$.

144 To quantify the genetic distance between an emerging variant and concurrently circulating lineages,
145 we calculated the mean number of unique amino acid substitutions in the circulating lineages rel-
146 ative to the emerging variant (Supplementary A.1). For each emerging variant, a high-quality
147 and representative early reference sequence was manually selected. For example, the BA1 refer-
148 ence sequence was from Botswana on 2021-12-20 (GISAID Accession ID: EPI_ISL_9002788),
149 and the BA.2 reference sequence was from Denmark on 2021-12-20 (GISAID Accession ID:
150 EPI_ISL_8347448). Within each country's respective fitting window (see below), we computed
151 the genetic distance between the emerging variant's reference and every sequence in the dataset.
152 The mean of these distances was then calculated for both the S-gene and the total genome, which
153 were used as features in the model. Further details can be found in the Supplementary Materials.

To include the relative fitness of the variants, we estimated the selection coefficient (s_i) for each emerging variant. We fitted a multinomial logistic growth model to variant proportions in each country during defined fitting windows. These selection coefficients have become the gold standard for estimating the relative growth rate of emerging variants (e.g. *CoV-Spectrum.org*²¹) and are used to quantify the speed of variant replacement; faster-growing variants are considered more likely to drive subsequent waves of infection. Here, the frequency $p_i(t)$ of variant i at time t is expected to increase relative to all other circulating variants by

$$p_i(t) = p_i(0) \exp(s_i t) / \sum_j p_j(0) \exp(s_j t),$$

154 where s_i is the selection coefficient of variant i per day, assuming constant selection for the fitting

155 period. We followed the fitting procedure and the code base developed by members of CAMEO
156 (CoVaRR-Net’s Computational Analysis, Modelling and Evolutionary Outcomes team²²), which
157 is used in the Canadian real-time variant tracking site *Duotang*²³ and is described by Gill et al.²⁴.

158 Serology data was available as a single numeric value per serology study, rather than a time series.
159 To handle this, we used Last Observation Carried Forward (LOCF) imputation to generate a piece-
160 wise constant time series. This carries forward the percent seroprevalence identified in each study
161 from the start of that study’s time window, through to the start of the next seroprevalence study
162 (or until February 1 2023, whichever is earlier). We selected LOCF for its simplicity, as it avoids
163 making assumptions such as linear growth in seroprevalence in the unobserved periods between
164 studies. We performed a similar LOCF imputation for the vaccine manufacturer data, which shows
165 how many of a given manufacturer’s vaccines were used in a given country. No further adjustments
166 or corrections were made to the raw data at this stage.

167 **2.4 Modelling approach**

168 While our master dataset covers the entire pandemic up to February 2023, for the task of predicting
169 BA.1/Delta and BA.2/BA.1, we focused on the subset of the data corresponding to the time period
170 of emergence of the respective new variant. This was designed to model the time in which the
171 question of emerging variant wave size would be of most interest to local authorities. We mask
172 from the predictive model the knowledge of any data after the end of a fixed fitting window.

173 We defined the fitting window as the week in which the emerging variant passed 10% of weekly
174 sequences in the country, and the four weeks prior (35 days total). For example, if the emerging
175 variant passed 10% in frequency on 2022-02-13, then the fitting window ranges from 2022-01-10
176 to 2022-02-13. This means that each country has its own fitting window for each wave, relative
177 to the variant’s growth in that country. We label t_0 as the first Sunday on which the emerging
178 variant’s sequences passed 10% of the total. The previous 4 Sundays are labelled t_{-1} , t_{-2} , t_{-3} ,

179 and t_{-4} , resulting in a fitting window of $[t_{-4}, t_0]$. We trained our model on slices of the master
180 data set within the fitting window (that is, all the features from dates t_{-4} , t_{-3} , t_{-2} , t_{-1} , and t_0).
181 An imbalance in the total number of variables (113) compared to the number of countries with
182 sufficient data for prediction (below) motivated the use of multiple data points per country. To
183 further minimize this imbalance, we discarded some redundant variables, prioritizing (i) summary
184 variables and (ii) those with least missingness. For example, we discarded daily new tests per
185 100,000 population and individual policy measures such as school closures, but retained daily new
186 tests and overall stringency index.

187 Many countries had significant amounts of missing epidemiological data or relatively limited ge-
188 nomic data. To ensure our analysis focused on countries with robust genomic surveillance, we
189 selected a subset of countries based on their sequencing effort. We identified 43 countries (Table
190 ??) that maintained a minimum of 10 sequences per week between 2021-06-01 and 2022-06-01.
191 These countries were also found to have the most robust data across other variables. The sub-
192 set of 25 countries belonging to the Organization for Economic Co-operation and Development
193 (OECD)²⁵ (Table ??) had further decreased missingness: the results of a model trained on these
194 OECD countries are included in the Supplementary Materials. In the rest of the main text, we
195 focus on a model trained on 46 variables (Table ??) and 215 (43×5) observations for 43 countries.

196 We used multiple imputation to handle the remaining missing data, using the *mice*²⁶ package in R
197 with Predictive Mean Matching (PMM), excluding the response variable of relative peak size. We
198 observed, across all countries, an all-or-nothing pattern of variable reporting: for each country and
199 variable, data were either fully reported across all five time points or entirely missing. Therefore,
200 we considered the “country” variable as a blocking factor for the imputation. We performed 30
201 separate imputations for each missing value, generating a set of plausible datasets to take into
202 account imputation uncertainty.

203 The Random Forest algorithm was chosen for its ability to model complex interactions among
204 variables. For each imputation of the dataset, the *randomForest*²⁷ package in R was used to train a

205 random forest model on a training set. We implemented a leave-one-out cross-validation (LOOCV)
206 procedure. For each iteration, all five data points from one country were held out as the test set,
207 while the remaining countries were used for training. We used grid search to identify the optimal
208 value for the number of features tried at each split, and each random forest model was configured
209 to grow 100,000 trees. Convergence was monitored through the Out-Of-Bag (OOB) error rate,
210 ensuring that the model's error rate stabilized as the number of trees increased.

211 Predictions were made for the held-out country using each of the 30 imputed dataset models.
212 These were averaged to obtain a single prediction per observation. The average predictions were
213 compared against the ground truth relative wave sizes to compute the Root Mean Squared Error
214 (RMSE) and the Coefficient of Determination R^2 . Additionally, the model's ability to correctly
215 predict whether the subsequent wave is larger or smaller than the last one was assessed by classi-
216 fying predictions into four quadrants based on the true and predicted ratios, with correct classifi-
217 cations falling into the first and third quadrants. Lastly, feature importance was calculated for each
218 imputed dataset and aggregated to understand the overall importance of each feature in predicting
219 the target variable.

220 **3 Results**

221 Figure 1 shows the overall VOC dynamics from 2020 to February 2023, derived from total global
222 case counts (top) and the VOC proportions from global whole genome sequencing data (top and
223 middle). The bottom panel shows the VOC replacement wave dynamics, illustrated by the VOC-
224 specific case waves. Note that, with the exception of Alpha and BQ, all VOCs had reached $> 85\%$
225 of the total global weekly cases (middle panel, Figure 1) by the time case counts peaked globally
226 (top panel), demonstrating a tight temporal coupling between the VOC replacement dynamics and
227 the case wave dynamics. In the first year of the pandemic, wildtype variants dominated and peaked
228 in late 2020. The emergence of the VOCs in early 2021 led to more oscillatory dynamics, with Al-

229 pha and Delta dominating during 2021. However, the emergence of Omicron was a massive global
230 contagion event that led to the largest infectious disease wave ever recorded. Following the promi-
231 nent BA.1 wave was a global smaller BA.2 wave, though case counting had been overwhelmed
232 in many countries after BA.1, and thus all Omicron waves should be considered under counted,
233 particularly beyond 2023 as public health policies shifted away from COVID-19 case counting.

234 While the aggregated global dynamics show a large BA.1 wave followed by a shorter and wider
235 BA.2 wave, the country-level reality varied significantly. In fact, the BA.1 and BA.2 Omicron lin-
236 eages emerged very close together in time, and whether BA.2 would drive a wave locally was hard
237 to interpret at the time. Figure 2 illustrates this variability among countries. The BA.1 and BA.2
238 wave dynamics can be divided into four categories: larger BA.2 than BA.1 (e.g. Thailand), equal
239 BA.2 (e.g. Sweden), smaller BA.2 (e.g. Canada), and no BA.2 wave (e.g. Peru). Delta was con-
240 sistently before the emergence of BA.1, but relative sizes at the country level also varied, despite
241 the picture of a smaller Delta wave depicted by global dynamics in Figure 1. See Supplementary
242 Figure ?? for the full Delta/BA.1 wave dynamics.

243 In countries with a larger BA.2 wave (defined as relative size > 1.1), t_0 occurred at the same time
244 or before the BA.1 peak (Austria, Denmark, Germany, Malaysia, South Korea, India, Singapore,
245 and Thailand, as shown in Figure 2). In countries where BA.2 and BA.1 had comparable sizes
246 (relative size ≥ 0.9 and ≤ 1.1), t_0 occurred before or after the peak of BA.1 (Sweden is before,
247 Finland is after). Of the 32 countries with smaller BA.2, 26 (81%) had t_0 after the BA.1 peak
248 (relative size ≥ 0.05 and < 0.9 ; Argentina, Australia, Brazil, Bulgaria, Canada, Chile, Colombia,
249 Costa Rica, Croatia, France, Ireland, Israel, Italy, Japan, Luxembourg, Mexico, Panama, Poland,
250 Portugal, Slovakia, Slovenia, South Africa, Spain, Switzerland, USA, United Kingdom). In one
251 country t_0 occurred at the same time as the BA.1 peak (Belgium), and in the remaining five (16%)
252 it occurred before the BA.1 peak (Indonesia, Netherlands, Norway, Romania, Russia). In the only
253 country with no BA.2 wave (relative size < 0.05), t_0 occurred after the BA.1 peak (Peru). Overall,
254 we see that our fitting windows generally end well before the peak of the wave to be predicted, and

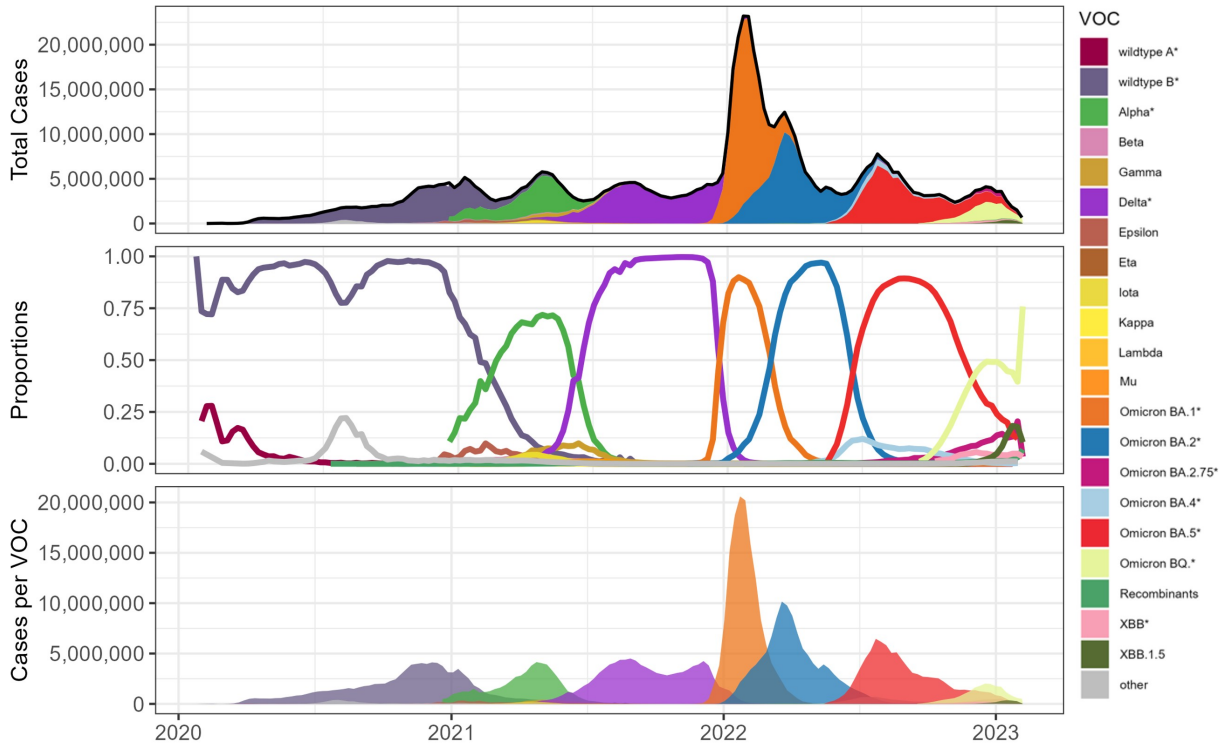


Figure 1: World VOC dynamics. Top panel: Total global reported cases filled in with stacked VOC proportions (middle panel). Middle panel: global VOC proportion plot from sequencing data. Bottom panel: VOC-specific cases over time (determined by multiplying the total cases by the VOC proportion). VOCs with * includes all sublineages.

255 indeed often before or during the peak of the prior wave. See Supplementary Figure ?? for the full
256 BA.1/BA.2 wave dynamics.

257 Selecting a subset of countries to use in the predictive model significantly reduces missing data
258 (Figure 3). The 43 chosen countries had at most 35/113 (31%) entirely missing variables, before
259 data sub-setting and imputation. We were also able to maintain reasonable global representation
260 among our chosen countries. Further visualizations and examples of data and missingness are
261 available in the Supplementary Materials, Figures ?? through ??.

262 Panel A of Figure 4 shows the predictions of the random forest model, with the response variable
263 modelling the relative size of BA.1 to the Delta wave. The RMSE was 4.09 and R^2 was 0.16. To
264 put this into perspective, the range of observed ratios of BA.1 to Delta was from 0.12 to 34.10.

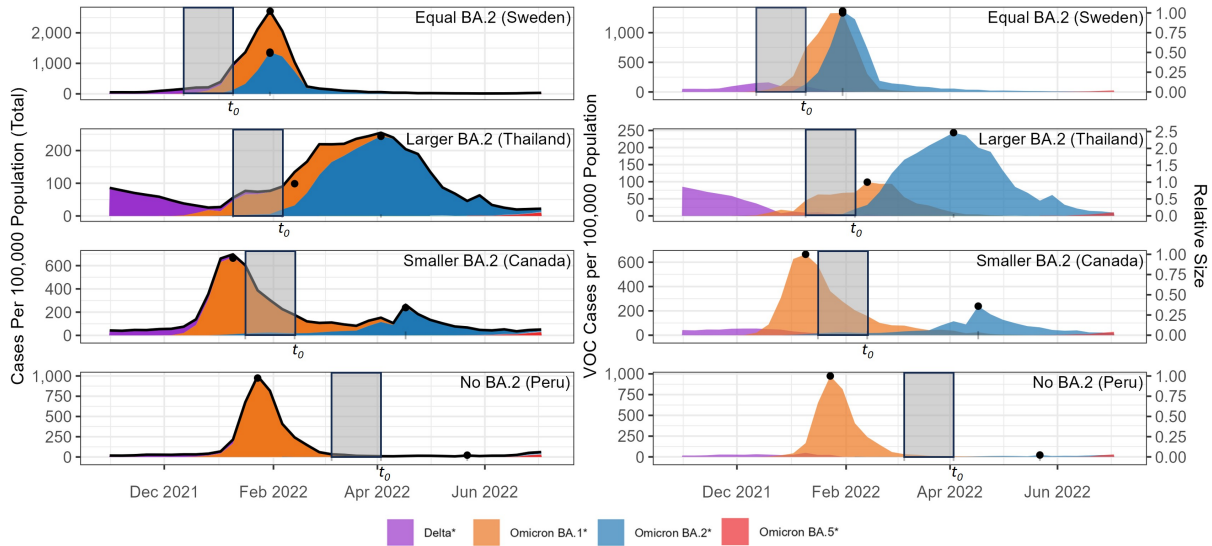


Figure 2: Example categories of country-specific wave dynamics of BA.1 and BA.2. The left column shows the total cases per country filled in by VOC proportion. The right column shows case waves attributed to specific VOCs in a single country over time. The top row shows an example (Sweden) with approximately equal BA.1 and BA.2 waves. The second row shows an example (Thailand) of a larger BA.2 wave. The third row shows an example (Canada) of a smaller BA.2 wave. The bottom row shows the only example (Peru) among the 43 countries with no BA.2 wave. The shaded area represents the five-week fitting window used in our analysis, indexed by the week in which BA.2 exceeds 10% of local sequences (t_0). Black dots represent the BA.1 and BA.2 peaks determined from the VOC-specific cases. In all countries, BA.1 always comes before BA.2, except in Sweden where they peak at the same time. The secondary Y-axis in the right column shows the response in our predictive model, i.e. the relative size of BA.2/BA.1 per country.

265 The prediction for South Korea has a large uncertainty interval. Australia is an outlier. Prior to
 266 the BA.1 wave, Australia had an incredibly low number of COVID-19 cases, under their pursued
 267 ‘COVID-zero’ strategy. This strategy was ended in late 2021, close to the time of BA.1 emergence
 268 globally. Although our model predicted a large growth of cases in Australia (predicted ratio 9.33),
 269 it was unable to predict the true scale of such a change (actual ratio 34.10).

270 Panel B of Figure 4 shows the predictions for the relative wave size of BA.2 to BA.1. The RMSE
 271 was 0.57 and R^2 was 0.39. The range of the observed ratios of BA.2 to BA.1 was 0.02 to 6.38.
 272 The RMSE of the BA.2/BA.1 model was much smaller than that of the BA.1/Delta model, which is
 273 expected given the smaller range of observed ratios for BA.2 to BA.1. India is the outlier here, with
 274 a much larger BA.2 wave than predicted. Interestingly, India, along with Denmark and Sweden,

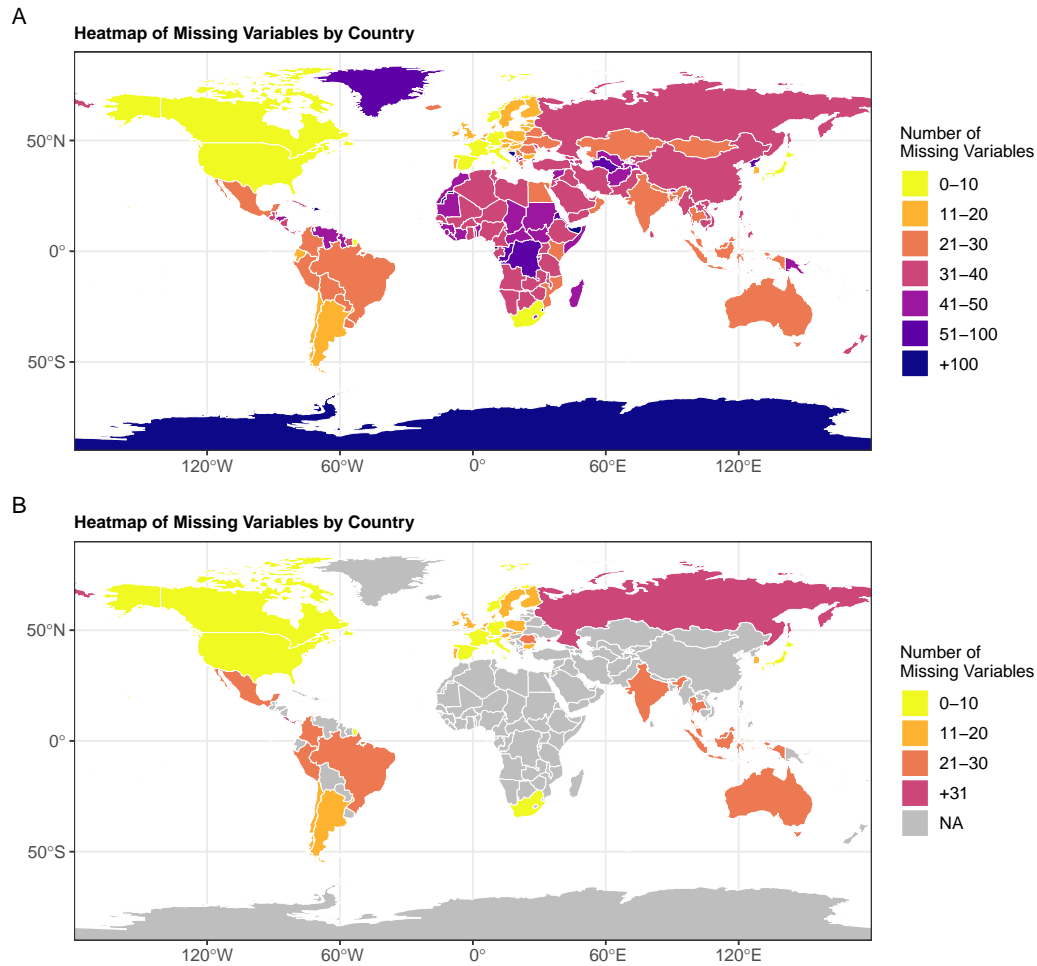


Figure 3: Number of missing variables (out of a total 113) per country, (A) among all countries for which data was collected, (B) among the 43 countries selected for prediction. We define here missing as entirely missing, for all weeks throughout the entire period of study 2019-12-29 to 2023-01-29.

275 had the highest mean genetic distance from the emerging variant BA.2, while also being one of the
276 first countries to detect this lineage (Figure ?? and Figure ??).

277 We also calculated how many times the random forest models predicted the general trend of the
278 subsequent wave, i.e., is the next wave larger, smaller, or equal in size to the last wave. This
279 measure of accuracy is perhaps more relevant in practice if the question of public health interest
280 is whether the next wave will be better or worse than the previous one. Our models correctly
281 predicted whether the subsequent wave was larger, smaller, or equal size in 91% and 80% of the

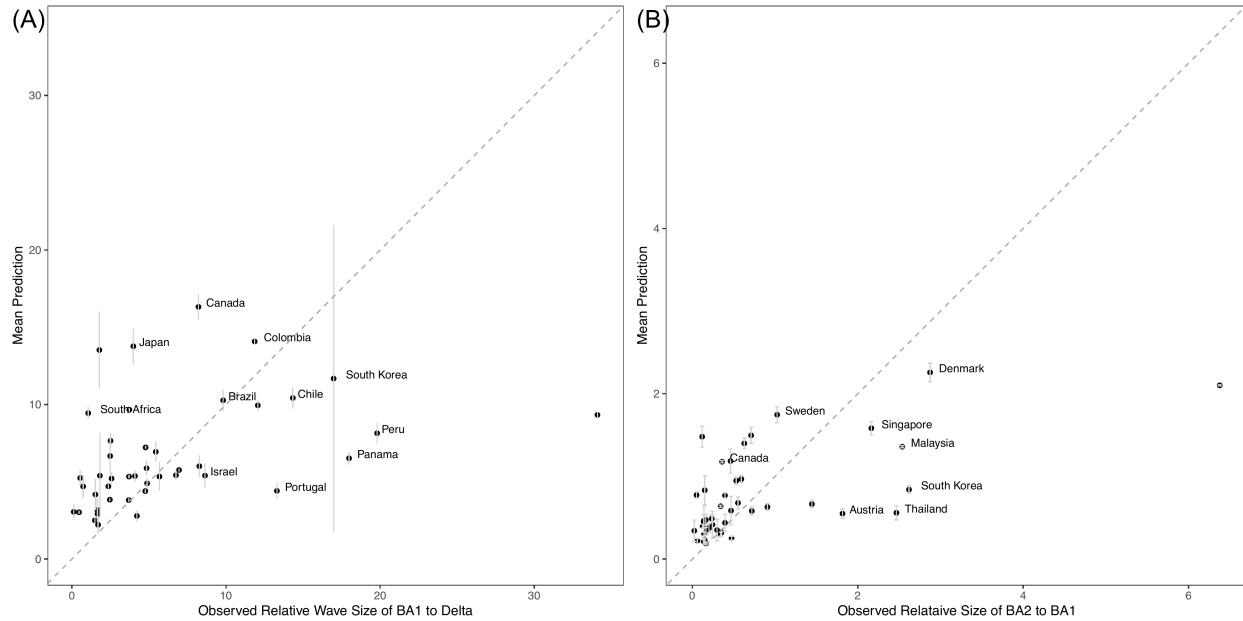


Figure 4: Mean predictions of the wave size of (A) BA.1 relative to Delta, and (B) BA.2 relative to BA.1, across 43 countries using a Leave-One-Country-Out cross-validation approach with a Random Forest model. The x-axis represents the observed relative wave sizes. The y-axis shows the corresponding mean predicted values. The dashed line represents the 1:1 line, indicating perfect agreement between observed and predicted values. Error bars represent the standard deviation of the predictions across multiple imputed datasets.

282 countries for BA.1/Delta and BA.2/BA.1, respectively (Figure 5).

283 Figure 6 shows the feature importance plots for the two models, visualizing the importance of each
284 predictor variable in contributing to the model's accuracy. The top panel ranks the features ac-
285 cording to the importance score, calculated as the mean increase in the mean squared error (MSE)
286 for the BA.2/BA.1 model when that feature is removed from the model. The bottom panel shows
287 the equivalent scores for the BA.1/Delta model. Genomically derived variables are among the fea-
288 tures of high importance. For the transition from BA.1 to BA.2, the genomic distance measure
289 of all genes was the most important feature, although it was only moderately important for the
290 transition from Delta to Omicron. The genomic distance of the S-gene and the selection coeffi-
291 cient were also moderately to highly informative in both models, as was the number of days since
292 the previous VOC peak. We identified several demographic and health variables of importance in
293 each wave, for example, cardiovascular death rates, life expectancy, and longitude were features

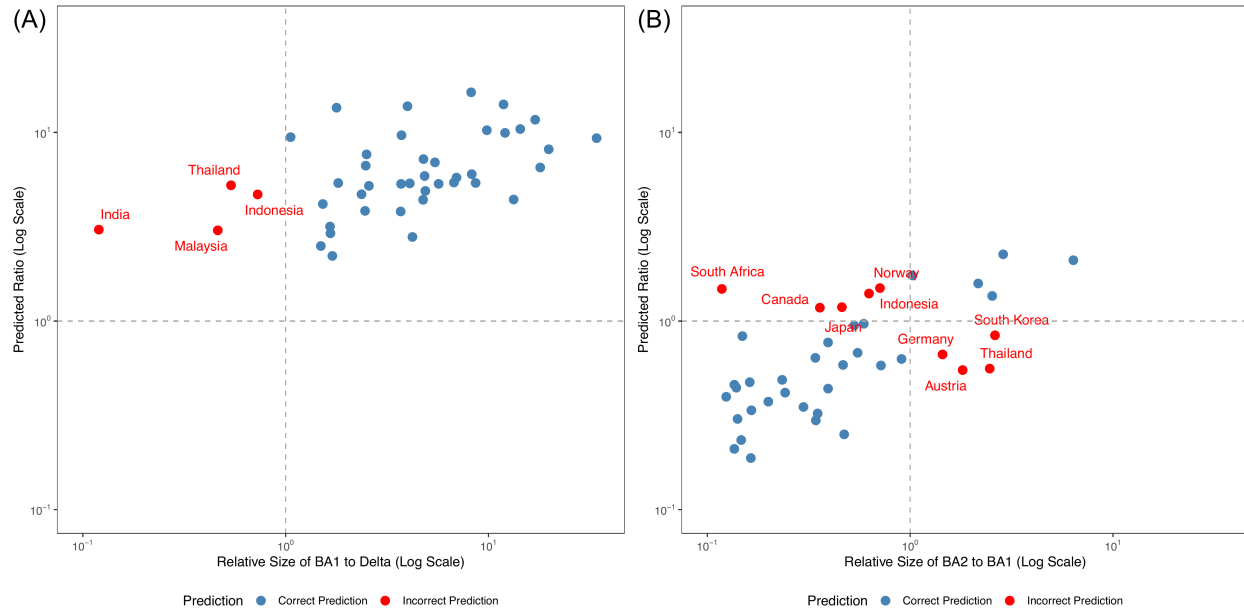


Figure 5: Prediction accuracy plots for two models: (A) BA.1 to Delta and (B) BA.2 to BA.1 variants, both on a log scale. Points in the first (top-right) and third (bottom-left) quadrants represent correct predictions where the model accurately identifies if the ratio of the subsequent waves is larger or smaller than 1. Points in the second and fourth quadrants represent incorrect predictions, where the model incorrectly identifies the ratio of the subsequent waves.

294 of highest importance for BA.1/Delta and somewhat lower importance for BA.2/BA.1. Diabetes
295 prevalence was more important in the transition from BA.1 to BA.2. Interestingly, we see consid-
296 erable changes in feature importance between waves, although several features (e.g., total deaths
297 per million, population density, GDP per capita) are informative across both waves.

298 The stringency index was not found to be important in predicting the relative wave size for either
299 wave, nor were counts of new deaths or cases during the fitting window (although total cases were
300 informative), nor mobility features (parks, residential, etc.) or serology features. However, we note
301 that the missingness in the serology data was high. Of note, all the features pertaining to seasons
302 of the year fell last or nearly last in importance, suggestive that the season a country is in at the
303 time of emergence may not have played a significant role in relative Omicron wave size (a marked
304 departure from what is expected of respiratory viruses).

305 The results of the 25 OECD countries model are available in the Supplementary Materials. Despite

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

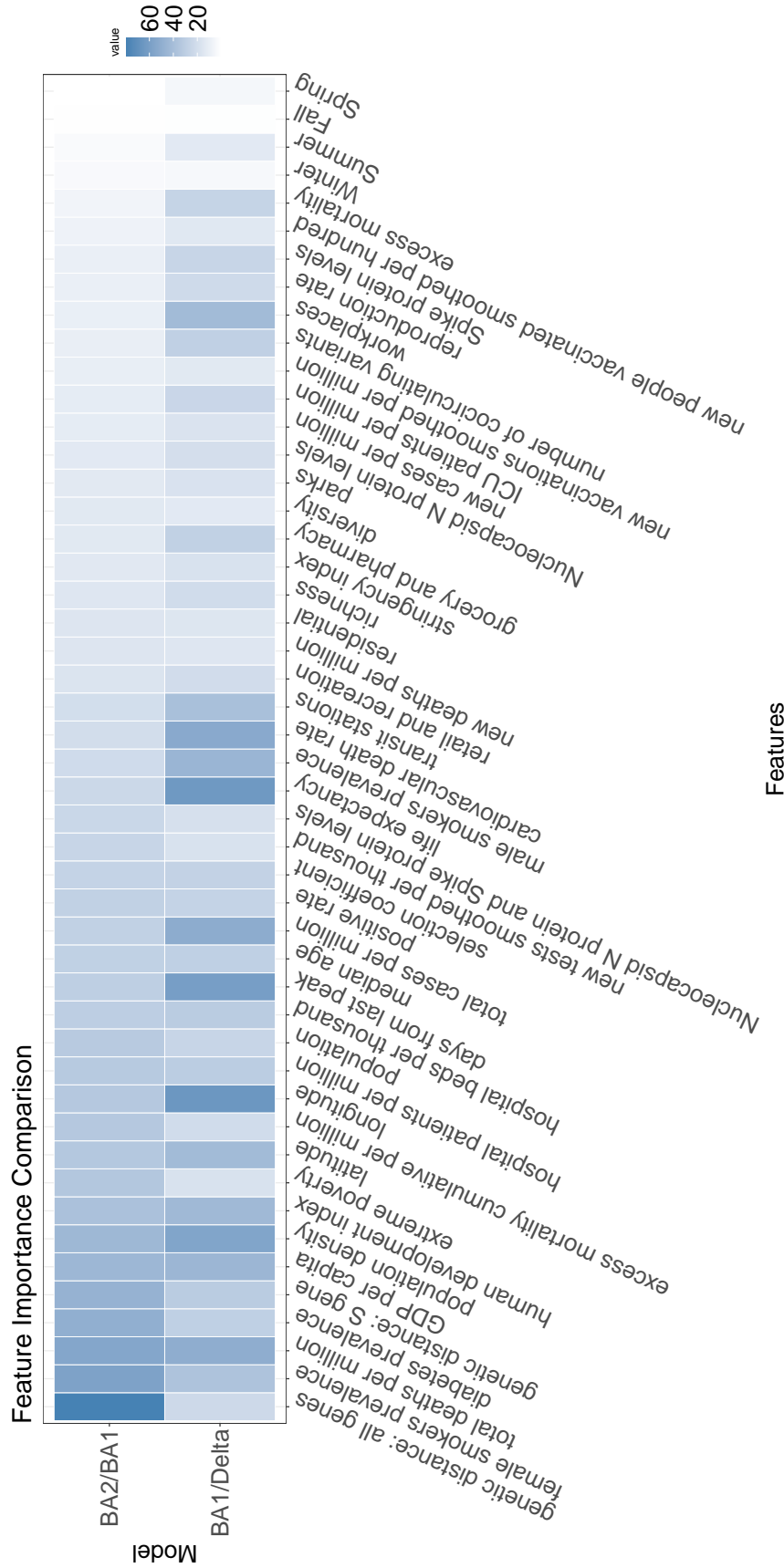


Figure 6: Feature importance comparison for Random Forest models predicting the relative wave size of BA.1 to Delta and BA.2 to BA.1 among 43 countries. The heatmap shows the Mean Decrease in Accuracy (MDA) of various features for each model, with darker shades indicating higher importance.

306 more consistent data quality in the OECD countries, the results were similar. The model trained
307 on 25 OECD countries has slightly higher accuracy in predicting if the next wave will be larger or
308 smaller (100% for BA.1/Delta and 88% for BA.2/BA.1 compared to 91% and 80%, respectively),
309 but RMSE and R^2 are very similar. Notably, all countries had a larger BA.1 wave than a Delta
310 wave. This is a promising result: suggesting that including countries with less consistent data is
311 not detrimental to model performance.

312 **4 Discussion**

313 Using a large well-curated global dataset that combines various data sources (epidemiological,
314 genomic, and other data), we find that random forest models can reasonably predict the relative
315 size of a wave (compared to the previous wave). In practice, the most pressing question upon the
316 emergence of a new variant is whether it is expected to have a stronger or weaker impact than the
317 previous variant, and the relative case wave size is seen as an important indicator of local impact:
318 in answering this, our model has very strong performance (Figure 5).

319 We find that the most important features for the transition between Delta and BA.1 included ge-
320 ographic factors (longitude), population health metrics (life expectancy and cardiovascular death
321 rate), and the number of days since the peak of the Delta wave. This suggests that demographic
322 conditions were the most critical in explaining the shift from Delta to BA.1. During the Delta pe-
323 riod, higher-income nations were vaccinated at very high rates, while many lower-income countries
324 struggled to secure vaccines²⁸: it is possible that demographic features, such as life expectancy,
325 may capture a broader economic effect on relative wave size through impacts on population im-
326 munity.

327 For the transition between BA.1 and BA.2, the most important features included genetic distance
328 measures (total number of amino acid substitutions and the number of S-gene amino acid substitu-

329 tions) between the emerging variant and circulating variants, and some population health metrics
330 (total female smokers, total deaths per million, and diabetes prevalence). Note that the total genetic
331 distance was more informative than the VOC's genetic distance in the S-gene alone. This finding
332 underscores that viral fitness is a polygenic trait and that mutations outside the Spike protein,
333 perhaps influencing replication efficiency or innate immune antagonism, were also critical deter-
334 minants of Omicron's success. The association with diabetes prevalence, a known risk factor for
335 severe COVID-19²⁹, may indicate that in a population with high levels of baseline immunity from
336 recent infection, the underlying comorbidities become more significant predictors of a variant's
337 success by marking the remaining pockets of vulnerability. This suggests that both the inherent
338 characteristics of the variants themselves and the demographics were strongly associated with this
339 transition.

340 We propose that the absence of shared top features between the two waves likely reflects the rapidly
341 evolving host-pathogen landscape. The transition from Delta to BA.1 occurred when global vac-
342 cine coverage was highly unequal, making broad demographic and economic factors and the result-
343 ing differences in population immunity the dominant predictors. In contrast, the subsequent tran-
344 sition from BA.1 to BA.2 occurred after the massive BA.1 wave had generated more widespread
345 (although short-lived) immunity, likely increasing the relative importance of the subtler genetic
346 differences between these two closely related sublineages. This finding suggests a challenge for
347 pandemic forecasting more broadly: if driving factors are markedly changing dynamically, we
348 cannot naively rely on models trained on previous waves to predict future waves.

349 Although seasonality was considered in the analysis, it did not show a strong influence on variant
350 dynamics in either wave studied. This finding challenges the assumption that COVID-19 would
351 evolve into a purely seasonal virus. The absence of seasonality as an important feature has implica-
352 tions for public health, suggesting that COVID-19 waves could occur outside of typical respiratory
353 virus seasons, affecting clinical capacities and complicating vaccine roll-out plans, as seen with
354 recent summer waves. The rate of new vaccinations was also not a top predictor for either wave.

355 This is likely due to two factors. First, by this stage of the pandemic, many of the countries in
356 our dataset had already achieved high coverage for the primary vaccine series, so the rate of new
357 vaccinations was low. Second, the protective effects of vaccination are not immediate.

358 Along with our findings from the subset of OECD countries (that sequenced larger amounts of
359 genomic data), our results suggest that viral genomics can offer helpful insights in forecasting
360 efforts as to the likely impact of new variants of concern. Furthermore, we found that measures
361 of genetic distance to the previously circulating variants in each country, overall or in the S-gene,
362 were more or similarly informative as the VOC's selection coefficient. That is, beyond just the
363 competitive advantage of a new VOC, the local genomic landscape and the VOC's ability to evade
364 immunity are influential in determining a VOC's success. This differs from the transition from
365 Delta to BA.1, where fewer people had previously been infected, making high transmissibility
366 alone sufficient for the new variant to spread.

367 As noted, many countries had missing data, and it was challenging to obtain timely and consistent
368 measures of serological features that reflect immunity. This may have contributed to the low impor-
369 tance of serological features in our predictive model. While serology was a key surveillance tool,
370 as the pandemic progressed and population exposure to SARS-CoV-2 increased, serology became
371 a less reliable indicator of recent infection. When most of the population had been exposed at least
372 once, the prevalence of *N* or *S* antibodies rose, making it difficult to detect reinfections through
373 serological data. Furthermore, numbers of reported cases depend on testing policy and reporting
374 inclinations of different jurisdictions, and testing may have at times focused on those most at risk
375 of a severe outcome, travellers, the general public, or on those with other risks. Polymerase Chain
376 Reaction (PCR) testing is expensive, and particularly during (and after) the first omicron wave,
377 many jurisdictions ceased formal testing for COVID-19 unless the result of the test would impact
378 medical treatment for the individual. Even hospitalization and death data suffer from incompara-
379 bility across jurisdictions, as protocols, reporting requirements and definitions varied. Yet, these
380 data provide our most comprehensive understanding of the pandemic and are the result of many

381 dedicated collection efforts.

382 Statistical models like the ones we have presented here have limitations. Indeed, there were many
383 limitations to COVID-19 forecasting efforts in general,³⁰ including challenges of data quality and
384 highly stochastic underlying transmission processes. Machine learning methods are an attractive
385 solution for modelling these processes without requiring a mechanistic representation of the deter-
386 minants of wave size, particularly as we can still incorporate our mechanistic understanding into
387 engineering features such as our variables from sequence data. However, these methods are data-
388 intensive. We bolstered our model with multiple time points per country to counteract the challenge
389 of having many more features than countries with sufficient data for prediction. An alternative ap-
390 proach might instead train the model on a single time point from several prior waves. However, as
391 our results revealed, the relative importance of features changed hugely between waves, potentially
392 reflecting large changes in local epidemiological context. It is important to acknowledge that our
393 analysis was conducted retrospectively, and thus is not presented as a real-time forecasting tool.
394 Although we held out all information from the model after the end of our fitting windows, we did
395 not take into account the possibility of reporting delays in the genomic data or other data streams.
396 Rather, we present this work as an exploratory analysis to identify which factors were predictive
397 during the Omicron BA.1 and BA.2 waves. We note, however, that while many predictive efforts
398 are hampered by reporting delays, the genomic features central to our model may be more robust
399 to this challenge and could be valuable in future forecasting frameworks, as they rely on VOC
400 proportions and not total counts.

401 Determining the relative size of the next variant-driven wave remains an important challenge, both
402 for SARS-CoV-2 and future pandemics. Here we presented a systematic procedure for modelling
403 two VOC waves, providing insights into the driving factors of Omicron BA.1 and BA.2 using the
404 same approach, and yet there are differences in feature importance reflecting the evolving nature
405 of reality and the value of information. As a future VOC emerges, it is unclear to what extent vac-
406 cination and previous infection confer immunity, and it is unclear how populations, governments,

407 and public health institutions will respond (and how data will be reported). We sought to determine
408 how well the relative wave sizes of BA.1 and BA.2 could be predicted, and the results were rea-
409 sonably good, particularly for the question of whether the new wave will be larger or smaller than
410 the most recent one. However, it remains unclear to what extent past waves will be good models
411 for future waves. To date, Omicron variants and recombinants have not evolved to become more
412 severe (as Delta was), though in evolutionary terms there is little to prevent this from occurring.
413 Although there is considerable knowledge of the mutations that govern increased transmissibility
414 and immune escape, with the data currently available, it is challenging to characterize the portion
415 of a given population that will be protected from a new variant. If it is deemed worthwhile to pre-
416 dict the burden and timing of a new variant, additional data on its infectivity and immune evasion
417 in the relevant population will likely be required.

418 **5 Data and code availability**

419 All data and code from this study are accessible on GitHub at [https://github.com/ShabMolan/](https://github.com/ShabMolan/VOC)
420 VOC. This repository contains accession numbers for all GISAID sequences used in this study, also
421 accessible at doi.org/10.55876/gis8.241023pk.

422 **6 Acknowledgements**

423 We gratefully acknowledge all data contributors, including the Authors and their Originating lab-
424 oratories responsible for obtaining the specimens, and their submitting laboratories for generating
425 the genetic sequence and metadata and sharing via the GISAID Initiative, on which this research
426 is based. We thank the countries who collected and shared their data with the various sources used
427 in this work (OWID, Johns Hopkins University COVID-19 Data Repository, Oxford COVID-19
428 Government Response Tracker, Google Community Mobility Reports, SeroTracker) as well as the

429 teams that have worked to maintain these websites.

430 We thank Dr Finlay Maguire, Dalhousie University, for support and technical discussions in the
431 early stages of the project.

432 **7 Funding**

433 The authors received funding from the Natural Sciences and Engineering Research Council of
434 Canada https://www.nserc-crsng.gc.ca/index_eng.asp through Discovery Grants
435 RGPIN-2023-4509 (JES) and RGPIN-2019-06624 (CC), and the Canadian Institutes of Health Re-
436 search <https://cihr-irsc.gc.ca/e/193.html> through the CGS-M program (SM). The
437 authors were supported by the Federal Government of Canada's Canada 150 Research Chair pro-
438 gramme [https://www.canada150.chairs-chaires.gc.ca/home-accueil-eng.](https://www.canada150.chairs-chaires.gc.ca/home-accueil-eng.aspx)
439 [aspx](https://www.canada150.chairs-chaires.gc.ca/home-accueil-eng.aspx) (CC), the Canadian Network for Modelling Infectious Diseases (CANMOD) [https://](https://canmod.net/)
440 canmod.net/ (CC), the Public Health Agency of Canada [https://www.canada.ca/en/](https://www.canada.ca/en/public-health.html)
441 [public-health.html](https://www.canada.ca/en/public-health.html) (ML, VG, CLM), and the Genomics Research and Development Ini-
442 tiative (GRDI) by the Government of Canada <https://grdi.canada.ca/en> (NKS). The
443 funders had no role in study design, data collection and analysis, decision to publish, or prepara-
444 tion of the manuscript.

445 **References**

446 ¹ Ensheng Dong, Hongru Du, and Lauren Gardner. An interactive web-based dashboard to track
447 COVID-19 in real time. *The Lancet infectious diseases*, 20(5):533–534, 2020.

448 ² World Health Organization. WHO Coronavirus (COVID-19) dashboard. [https://data.](https://data.who.int/dashboards/covid19/about)
449 [who.int/dashboards/covid19/about](https://data.who.int/dashboards/covid19/about), 2023.

- 450 ³ Google LLC. Google COVID-19 Community Mobility Reports. <https://www.google.com/covid19/mobility/>, 2022.
- 451
- 452 ⁴ Sean Jewell, Joseph Futoma, Lauren Hannah, Andrew C. Miller, Nicholas J. Foti, and Emily B. Fox. It's Complicated: Characterizing The Time-varying Relationship Between Cell Phone Mobility and COVID-19 Spread in the US. In *Nature Digital Medicine*, 2021.
- 453
- 454
- 455 ⁵ Shruti Khare, Céline Gurry, Lucas Freitas, Mark B Schultz, Gunter Bach, Amadou Diallo, Nancy Akite, Joses Ho, Raphael TC Lee, Winston Yeo, et al. GISAID's role in pandemic response. *China CDC weekly*, 3(49):1049, 2021.
- 456
- 457
- 458 ⁶ Rahul K Arora, Abel Joseph, Jordan Van Wyk, Simona Rocco, Austin Atmaja, Ewan May, Tingting Yan, Niklas Bobrovitz, Jonathan Chevrier, Matthew P Cheng, et al. SeroTracker: a global SARS-CoV-2 seroprevalence dashboard. *The Lancet Infectious Diseases*, 21(4):e75–e76, 2021.
- 459
- 460
- 461
- 462 ⁷ Thomas Hale, Noam Angrist, Rafael Goldszmidt, Beatriz Kira, Anna Petherick, Toby Phillips, Samuel Webster, Emily Cameron-Blake, Laura Hallas, Saptarshi Majumdar, et al. A global panel database of pandemic policies (Oxford COVID-19 Government Response Tracker). *Nature human behaviour*, 5(4):529–538, 2021.
- 463
- 464
- 465
- 466 ⁸ Our World in Data. COVID-19 Dataset by Our World in Data. <https://covid.ourworldindata.org/>, 2023.
- 467
- 468 ⁹ Pango Lineages. Pango lineage designation issues. <https://github.com/cov-lineages/pango-designation/issues>. Accessed: 2024-09-26.
- 469
- 470 ¹⁰ ETH Zurich. CoV-Spectrum: SARS-CoV-2 variant tracking tool. <https://cov-spectrum.org>. Accessed: 2024-09-26.
- 471
- 472 ¹¹ World Health Organization. Tracking sars-cov-2 variants. <https://www.who.int/activities/tracking-SARS-CoV-2-variants>. Accessed: 2024-09-26.
- 473

- 474 ¹² Áine O’Toole, Emily Scher, Anthony Underwood, Ben Jackson, Verity Hill, John T McCrone,
475 Rachel Colquhoun, Chris Ruis, Khalil Abu-Dahab, Ben Taylor, et al. Assignment of epidemio-
476 logical lineages in an emerging pandemic using the pangolin tool. *Virus evolution*, 7(2):veab064,
477 2021.
- 478 ¹³ James Hadfield, Colin Megill, Sidney M Bell, John Huddleston, Barney Potter, Charlton Callen-
479 der, Pavel Sagulenko, Trevor Bedford, and Richard A Neher. Nextstrain: real-time tracking of
480 pathogen evolution. *Bioinformatics*, 34(23):4121–4123, December 2018.
- 481 ¹⁴ Nextstrain. Clade definitions. [https://github.com/nextstrain/ncov/blob/](https://github.com/nextstrain/ncov/blob/master/defaults/clades.tsv)
482 [master/defaults/clades.tsv](https://github.com/nextstrain/ncov/blob/master/defaults/clades.tsv). Accessed: 2024-09-26.
- 483 ¹⁵ Pavel Sagulenko, Vadim Puller, and Richard A Neher. TreeTime: Maximum-likelihood phylo-
484 dynamic analysis. *Virus evolution*, 4(1):vex042, 2018.
- 485 ¹⁶ Raquel Viana, Sikhulile Moyo, Daniel G Amoako, Houriiyah Tegally, Cathrine Scheepers,
486 Christian L Althaus, Ugochukwu J Anyaneji, Phillip A Bester, Maciej F Boni, Mohammed
487 Chand, et al. Rapid epidemic expansion of the SARS-CoV-2 Omicron variant in southern Africa.
488 *Nature*, 603(7902):679–686, 2022.
- 489 ¹⁷ Ivan Aksamentov, Cornelius Roemer, Emma B Hodcroft, and Richard A Neher. Nextclade:
490 clade assignment, mutation calling and quality control for viral genomes. *Journal of open source*
491 *software*, 6(67):3773, 2021.
- 492 ¹⁸ Nextclade CLI Documentation. Nextclade CLI. [https://docs.nextstrain.org/](https://docs.nextstrain.org/projects/nextclade/en/stable/)
493 [projects/nextclade/en/stable/](https://docs.nextstrain.org/projects/nextclade/en/stable/). Accessed: 2024-09-26.
- 494 ¹⁹ Mark O Hill. Diversity and evenness: a unifying notation and its consequences. *Ecology*,
495 54(2):427–432, 1973.
- 496 ²⁰ Jari Oksanen, Gavin L. Simpson, F. Guillaume Blanchet, Roeland Kindt, Pierre Legendre, Pe-
497 ter R. Minchin, R.B. O’Hara, Peter Solymos, M. Henry H. Stevens, Eduard Szoecs, Helene Wag-

598 ner, Matt Barbour, Michael Bedward, Ben Bolker, Daniel Borcard, Gustavo Carvalho, Michael
599 Chirico, Miquel De Caceres, Sebastien Durand, Heloisa Beatriz Antoniazi Evangelista, Rich
600 FitzJohn, Michael Friendly, Brendan Furneaux, Geoffrey Hannigan, Mark O. Hill, Leo Lahti,
601 Dan McGlinn, Marie-Helene Ouellette, Eduardo Ribeiro Cunha, Tyler Smith, Adrian Stier,
602 Cajo J.F. Ter Braak, and James Weedon. *vegan: Community Ecology Package*, 2022. R package
603 version 2.6-4.

604 ²¹ C. Chen, S. Nadeau, M. Yared, P. Voinov, N. Xie, C. Roemer, and T. Stadler. CoV-Spectrum:
605 analysis of globally shared SARS-CoV-2 data to identify and characterize new variants. *Bioin-*
606 *formatics*, 38(6):1735–1737, mar 2022.

607 ²² CoVaRR-Net. Computational analysis, modelling and evo-
608 lutionary outcomes (cameo). [https://covarrnet.ca/
609 computational-analysis-modelling-and-evolutionary-outcomes-cameo/](https://covarrnet.ca/computational-analysis-modelling-and-evolutionary-outcomes-cameo/).
610 Accessed: 2024-09-26.

611 ²³ CoVaRR-Net SARS-CoV-2 Duotang. [https://covarr-net.github.io/duotang/
612 duotang.html](https://covarr-net.github.io/duotang/duotang.html). Accessed: 2024-09-27.

613 ²⁴ Erin E Gill, Baofeng Jia, Carmen Lia Murall, Raphaël Poujol, Muhammad Zohaib An-
614 war, Nithu Sara John, Justin Richardsson, Ashley Hobb, Abayomi S Olabode, Alexan-
615 dru Lepsa, et al. The Canadian VirusSeq Data Portal & Duotang: open resources for
616 SARS-CoV-2 viral sequences and genomic epidemiology. *ArXiv*, 2024. arXiv preprint:
617 <https://arxiv.org/pdf/2405.04734>.

618 ²⁵ Organisation for Economic Co-operation and Development. <https://www.oecd.org/>.
619 Accessed: 2024-06-17.

620 ²⁶ Stef van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained
621 equations in R. *Journal of Statistical Software*, 45(3):1–67, 2011.

522 ²⁷ Andy Liaw and Matthew Wiener. Classification and Regression by randomForest. *R News*,
523 2(3):18–22, 2002.

524 ²⁸ Peter B McIntyre, Rakesh Aggarwal, Ilesh Jani, Jaleela Jawad, Sonali Kochhar, Noni MacDon-
525 ald, Shabir A Madhi, Ezzeddine Mohsni, Kim Mulholland, Kathleen M Neuzil, et al. COVID-19
526 vaccine strategies must focus on severe disease and global equity. *The Lancet*, 399(10322):406–
527 410, 2022.

528 ²⁹ U.S. Centers for Disease Control and Prevention. Underlying Medical Conditions Associ-
529 ated with Higher Risk for Severe COVID-19. [https://www.cdc.gov/coronavirus/
530 2019-ncov/your-health/underlying-medical-conditions.html](https://www.cdc.gov/coronavirus/2019-ncov/your-health/underlying-medical-conditions.html), 2025. Ac-
531 cessed: 2025-07-27.

532 ³⁰ John PA Ioannidis, Sally Cripps, and Martin A Tanner. Forecasting for COVID-19 has failed.
533 *International journal of forecasting*, 38(2):423–438, 2022.