

---

# ENHANCING GLOBAL PREPAREDNESS DURING AN ONGOING PANDEMIC FROM PARTIAL AND NOISY DATA

---

Pascal Klamser<sup>1,2,†</sup>, Valeria d'Andrea<sup>3,†</sup>, Francesco Di Lauro<sup>4</sup>, Adrian Zachariae<sup>1,2</sup>, Sebastiano Bontorin<sup>3,5</sup>, Antonello di Nardo<sup>6</sup>, Matthew Hall<sup>4</sup>, Benjamin F. Maier<sup>1,2</sup>, Luca Ferretti<sup>4</sup>, Dirk Brockmann<sup>1,2</sup>, Manlio De Domenico<sup>7,8,\*</sup>

<sup>1</sup>Robert Koch-Institute, Nordufer 20, 13353 Berlin, Germany

<sup>2</sup>Institute for Theoretical Biology, Humboldt-University of Berlin, Philippstr. 13, D-10115 Berlin, Germany

<sup>3</sup>Fondazione Bruno Kessler, Via Sommarive 18, 38123 Povo (TN), Italy

<sup>4</sup>Big Data Institute, University of Oxford, Li Ka Shing Centre for Health Information and Discovery, Oxford, UK

<sup>5</sup>Department of Physics, University of Trento, Via Sommarive 14, 38123 Povo (TN), Italy

<sup>6</sup>The Pirbright Institute, Ash Road, Pirbright, Woking, Surrey, United Kingdom

<sup>7</sup>Department of Physics and Astronomy, G. Galilei, University of Padua, Via Francesco Marzolo 8, 35131, Padua, Italy

<sup>8</sup>Padua Center for Network Medicine, University of Padua, Via Francesco Marzolo 8, 35131, Padua, Italy

<sup>†</sup>Contributed equally to this work.

\*Corresponding author. E-mail: [manlio.dedomenico@unipd.it](mailto:manlio.dedomenico@unipd.it)

August 19, 2022

## ABSTRACT

As the coronavirus disease 2019 (COVID-19) spread globally, emerging variants such as B.1.1.529 quickly became dominant worldwide. Sustained community transmission favors the proliferation of mutated sub-lineages with pandemic potential, due to cross-national mobility flows, which are responsible for consecutive cases surge worldwide. We show that, in the early stages of an emerging variant, integrating data from national genomic surveillance and global human mobility with large-scale epidemic modeling allows to quantify its pandemic potential, providing quantifiable indicators for pro-active policy interventions. We validate our framework on worldwide spreading variants and gain insights about the pandemic potential of BA.5 and BA.2.75 sub-lineages. Country-level epidemic intelligence is not enough to contrast the pandemic of respiratory pathogens such as SARS-CoV-2 and a scalable integrated approach, i.e. pandemic intelligence, is required to enhance global preparedness.

## Introduction

The coronavirus disease (COVID-19) outbreak, caused by the SARS-CoV-2 virus and first detected in China in early 2020, likely originated from the Huanan seafood wholesale market in Wuhan [1] and continues to spread worldwide. It has forced national governments to pursue country-level elimination strategies [2, 3, 4] or mitigation policies relying on both non-pharmaceutical interventions (NPI) – e.g., physical distancing, wearing masks, hand hygiene, limit large gathering of people, curfews and, in the worst cases, lockdowns [5] – and pharmaceutical ones, such as massive vaccination campaigns and antiviral therapies [6, 7, 8]. Early strict interventions have been shown to be more effective

than longer moderate ones in containing national outbreaks in curbing epidemic growth [9], for similar intermediate distress and infringement on individual freedom [10].

In contrast to policy during the early stages of the pandemic, when pharmaceutical interventions were not yet available, most current national efforts to control the virus rely on reactive strategies which alternate enhancement and lifting of NPIs, with the ultimate goal of prevention, or reduction, of pressure on national health systems. To achieve successful containment, such reactive strategies require high capacity for testing and sequencing to continuously monitor the potential emergence of novel viral strains of SARS-CoV-2, whose mutations might be responsible for more severe and/or more transmissible variants with pandemic potential. We define pandemic potential as the ability of a variant to escape population immunity acquired by vaccination or previous infections and to quickly spread worldwide.

Although the emergence of within-host variants with immune escape is likely to be relatively rare [11], sustained community transmission might favor it. When a new variant emerges, it is crucial for policy and decision-making to characterize novel mutations [12, 13, 14], estimate the growth advantage of the new variant with respect to the existing ones [15] and quantify the effectiveness of currently available vaccines [16, 17]. Consequently, any delay in identifying an emerging variant and in determining its key epidemiological parameters introduces uncertainties in the timeline of community transmissions and imported cases which limit, if not completely prevent, effective mitigation responses to take place, similarly to the cryptic transmission of the wild type SARS-CoV-2 which led to the first COVID-19 wave [18]. Combined with limited testing capacity, porous travel screening [19] – at national and, overall, cross-national levels, where international travel play a significant role to amplify the pandemic potential [20, 21, 18] – and lifting of national NPI, the same delays might seriously hinder the timely detection of an emerging variant. The COVID-19 pandemic has been characterized by the regular emergence of such variants [22]. Three important questions arise during the early stages of such a variant, at which point data is missing and noisy: i) can we reconstruct its geographical origin? ii) can we estimate how long it has been spreading undetected in that location? iii) can we quantify the risk of importation to other locations?

In this work, we devise a protocol to quantitatively answer these questions. We show that, by integrating phylogenetic, epidemiological and behavioral analyses within a framework for data-driven and model-informed pandemic intelligence, it is possible to quantify the pandemic potential of an emerging variant and predict the dynamics of subsequent national outbreaks with satisfactory precision.

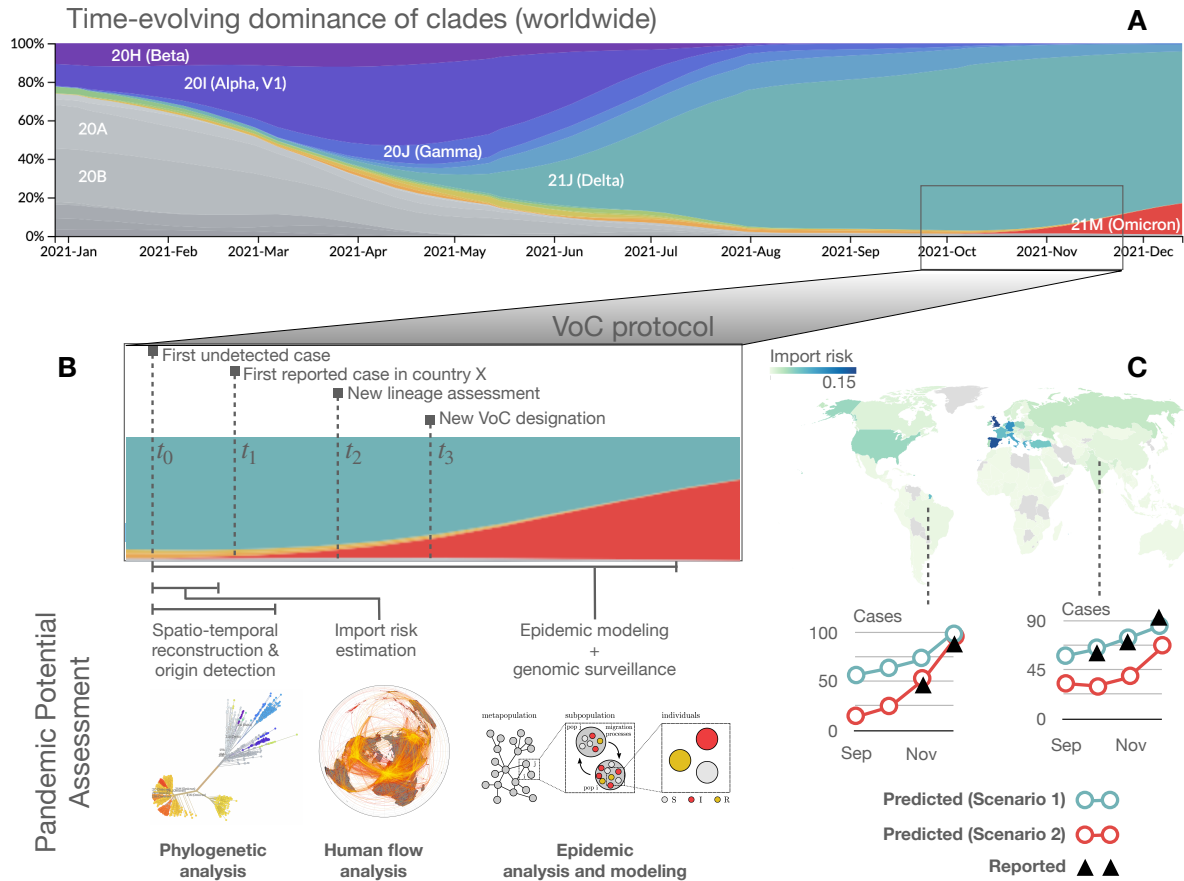
## Results

**Blueprint for a pandemic intelligence framework.** Reliably quantifying the pandemic potential of an emerging variant requires data, and acquiring data requires time. Between the time  $t_0$  of the first undetected case and the time  $t_1$  of the first reported case and its subsequent lineage designation at time  $t_2$ , an emerging variant can silently spread within its country of origin and beyond. For example, let us consider the B.1.1.529 (in the following referred to BA.1) lineage of the Omicron variant (also known as BA.1). This was first reported by genomic surveillance teams in South Africa and Botswana on November 25th 2021. Priority actions have been established by the World Health Organization (WHO) for member states on November 26th, with designation as a variant of concern (VOC) [23] required to raise the level of international alert ( $t_3$ ). By December 16th 2021, there were several reports of an estimated reduction in both vaccine effectiveness against infection and severe disease [24, 25, 26, 17, 27], together with characterisations of the epidemiology of the variant in South Africa [28], Denmark [29] and Norway [30]. Early phylogenetic analysis placed  $t_0$  during the third week of October 2021, about one month before  $t_1$ . Three weeks later it had been identified in 87 countries [28].

Figure 1 summarizes this timeline for BA.1, while highlighting the main analytical steps required to define a self-consistent protocol to characterize the pandemic potential of an emerging variant (see Supplementary Fig. S1 for more mechanistic scheme). Figure 1B illustrates how genomic surveillance data and epidemic modeling can be used to infer the spatio-temporal coordinates of the variant's origin, thus providing information on  $t_0$ . This information is used to estimate the importation risk for all countries in the world due to cross-national human flows. Finally, imported cases are used as seeds for community transmission leading to country-level outbreaks, while accounting for the epidemiological parameters characterizing the new variant. Unavoidable uncertainties about  $t_0$  and epidemiological parameters are propagated through the workflow. Plausible scenarios are presented, accounting for distinct levels of case under-reporting in each destination country (Fig. 1C).

In the following, we describe each step of the procedure, detailing our pandemic intelligence framework and the underlying modeling assumptions.

**Reconstructing the origin of an emerging variant in space and time.** For all SARS-CoV-2 sequences belonging to the B.1.1.7 (Alpha), B.1.617.2 (Delta), BA.1, BA.2, BA.5, and BA.2.75 (Omicron) lineages from GISAID [31, 32, 33],



**Figure 1: Schematic illustration of our pandemic intelligence workflow.** (A) Evolutionary dynamics of SARS-CoV-2 variants, coded by colour. The panel is obtained from [nextstrain.org](https://nextstrain.org), based on GISAID data. (B) For the B.1.1.529 lineage (or BA.1, *Omicron*, according to the WHO nomenclature), we identify four distinct time points in the process of characterising the variant, from the time of the first undetected case to the designation as Variant of Concern. This illustrates how genomic surveillance data is used in combination with global human movement data and epidemic modeling to: i) perform a spatiotemporal reconstruction of the patient zero to identify the country of origin of an emerging variant and estimate its epidemiological parameters and ii) calculate the importation risk for all other countries worldwide. (C) For a subset of about 50 countries worldwide (depending on sequencing data availability), we forecast the increase in the number of cases due to the consequent community transmission according to what-if scenarios, accounting for distinct levels of under-reporting. For a more mechanistic workflow scheme see Supplementary Fig. S1.

we retained only those generated from cases reported during the early stage of the corresponding wave from the country of evolutionary origin, from 20 up to a total of 100 sequences per lineage. Where there were multiple candidate countries of origin, we estimated the outbreak country by a simple trait model. We then generated 3 alignments, comprised of respectively 20%, 50% and 100% of the sequence set. These were subsequently cleaned by trimming the 5' and 3' untranslated regions and gap-only sites. Bayesian evolutionary reconstruction of the dated phylogenetic history [34] was used to obtain posterior distributions of the growth rate  $t$ , the parameters of the molecular clock, and the time of the most recent common ancestor (tMRCA). See Materials and Methods for details.

In this way obtain an estimate of  $t_0$ , the time of the first unreported case, as well as of other epidemic parameters such as the growth rate. From these we estimated the effective reproduction number and generation interval. Indicating the number of infected individuals and number of deaths at time  $t$  by  $I(t)$  and  $D(t)$  respectively, we consider the time period during which there is co-circulation of an existing variant  $v$  and an emerging one  $\omega$ . We approximate the epidemic evolution by

$$\begin{aligned}
 I(t_0 + \Delta t) &= I_v(t_0 + \Delta t) + I_\omega(t_0 + \Delta t) = \\
 &= I_v(t_0)R_v(t_0)^{\Delta t/GI_v} + I_\omega(t_0)R_\omega(t_0)^{\Delta t/GI_\omega},
 \end{aligned}
 \tag{1}$$

where  $I_x(t)$  is the number of infections due to variant  $x$  at time  $t$ ,  $R_x(t_0)$  is the effective reproduction number at time  $t_0$ , and  $GI_x$  is the generation interval. Similarly, the deaths due to the co-circulating variants are approximated by  $D(t) = D_v(t) + D_\omega(t)$ , with

$$D_x(t_0 + \Delta t + \tau_x) = I_x(t_0 + \Delta t) \times IFR_x, \quad x = v, \omega \quad (2)$$

where  $IFR_x$  denotes the infection fatality rate of variant  $x$  and  $\tau_x$  is the time lag between infection and death. To fit the unknown epidemiological parameters, i.e. the ones related to variant  $\omega$  for which we obtain a joint probability distribution, we use an optimization procedure (see Materials and Methods)

In the case of BA.1, we obtain  $t_0 = 28$  October 2021 (95% HPD: 20 October–5 November) and a daily growth rate estimate of 0.582 (95% HPD: 0.117–1.035) from the phylogenetic analysis and  $t_0 = 19$  October 2021 (95% CL: 15 October–23 October) from epidemic modelling, with  $R_t = 2.56$  (95% CL: 2.16–3.19) and  $GI = 7.36$  (95% CL: 6.12–9.17). Our results are in good agreement with the literature, reporting  $t_0 = 9$  October 2021 (95% HPD: 30 September–20 October), exponential growth rate of 0.137 (95% HPD: 0.099–0.175) per day [28] and  $GI = 6.84$  days (95% credible intervals: 5.72–8.60) [35].

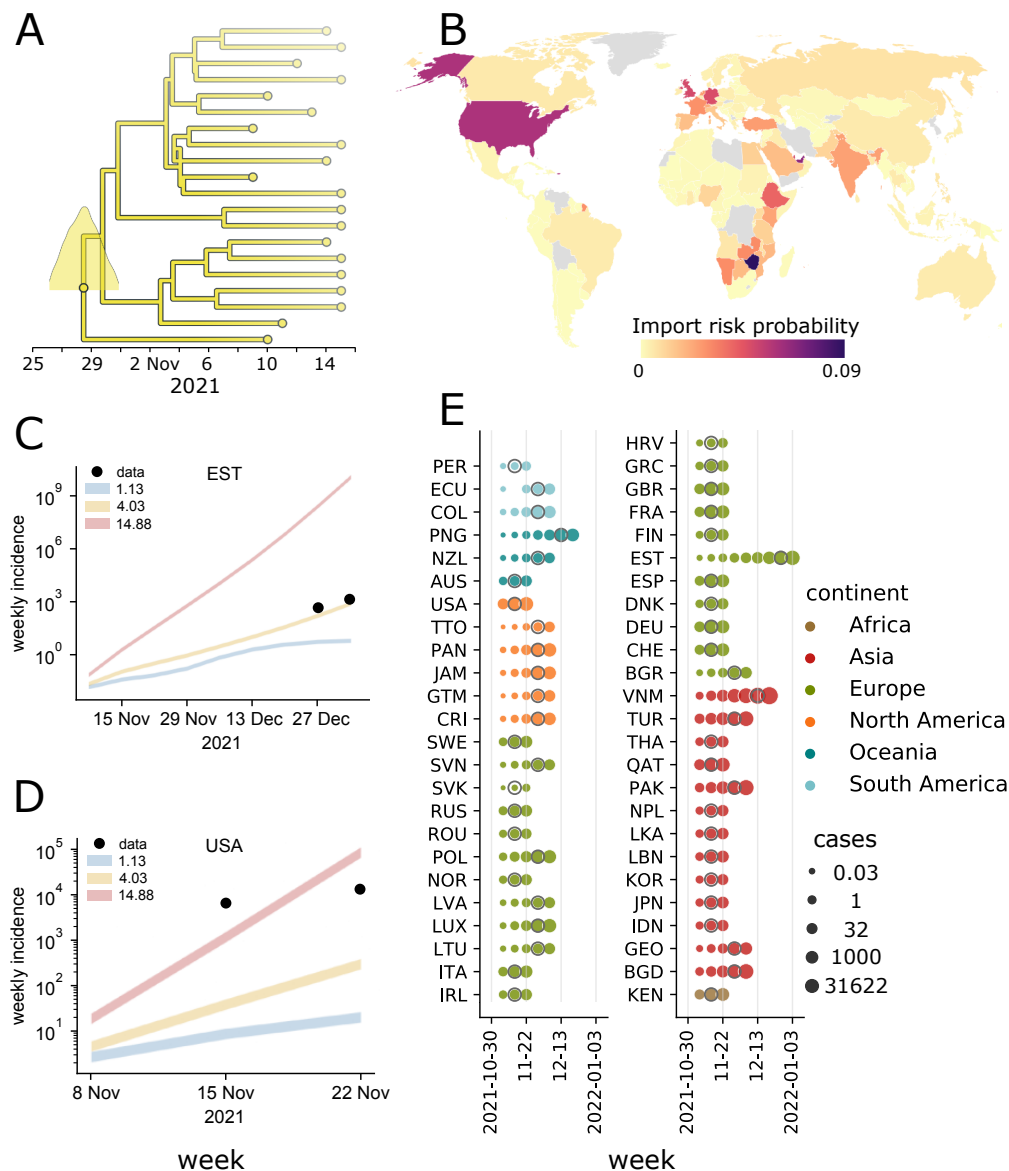
For further details, refer to Materials and Methods and Supplementary Figs. S2–S3.

**Estimating the importation risk of an emerging variant by country.** We use monthly seat capacities of flights between airports from the Official Airline Guide [36], encoding how many people could have travelled if all seats were occupied on flights from airport A to B in the month of the estimated  $t_0$ . We indicate the corresponding flow matrix by  $\mathbf{F}$ , where entry  $F_{ij}$  describes the maximal passenger flow to  $i$  from  $j$ . The travelling population in the catchment area of an airport is obtained by  $N_i = F_i$ , with  $F_i = \sum_j F_{ji}$ , i.e., we assume that the population in the catchment area of the airport is equal to the airports outflow. For each emerging variant, the resulting large-scale network of international travels corresponding to the month of  $t_0$  is used. The import risk is calculated as in [37]: based on a reconstructed effective distance graph [21] and a random walk with an exit option we estimate how likely it is that one infected individual reaches any airport worldwide given they set off from an airport in the country where the variant emerged. To work at country level, we aggregate the import risk of all airports of the outbreak country by computing the weighted mean with weights

$$w_n = \frac{N_n}{\sum_{m \in C(n)} N_m} \quad (3)$$

where  $C(n)$  the set of airports that belong to the same country as airport  $n$  does. We have performed an extensive analysis to validate the estimated import risk against available data, such as the official arrival times as obtained from the WHO, for each emerging variant. We find considerable correlation between arrival time and import risk distance for different variants (Alpha, Beta, Delta, Gamma) with a median of  $r = 0.55$  (range  $r \in [0.41, 0.56]$ ). This median is the largest we found, when compared with several alternative distance measures (see Supplementary Figs. S2, S3, S4). The stochastic nature of the reported variant arrival times (based on the by-country rate of genome sequencing, the probabilistic distribution of infected individuals among passengers, etc.) is one possible reason for the imperfect correlation, but another possibility is that the assumed outbreak location is incorrect. To test this, we attempted to identify outbreak locations by recomputing the correlation for all countries (similarly to [21]). For Beta, Gamma, and BA.1 the country declared by the WHO as the outbreak source had the greatest degree of correlation. For Delta and Alpha the WHO candidate had the 2nd and 5th best correlation respectively (see Supplementary Figs. S4, S5). We extended the analysis to sublineages of Omicron and previously circulating variants of interest (VOIs) by estimating arrival times and outbreak countries from GISAID data (see Material and Methods). For 13 of 17 variants the suspected outbreak location from GISAID had at least the 3rd-largest correlation coefficient (of 183), and for all variants the GISAID candidate was at least on the 12th rank (see Supplementary Figs. S7, S8, S9).

**Modeling country-level epidemic spread of an emerging variant under distinct scenarios.** We use results from the previous step of the pipeline as inputs for an epidemic model in order to forecast the potential surge in cases due to an emerging variant in a target country. First, we estimate the daily number of infected people (seeds) traveling to the target country from the country where the VoC emerged (source country), based upon four elements: 1) results of our phylogenetic analysis, which inform both the growth rate and the time of emergence of the variant of concern, 2) genomic surveillance in the source country, 3) estimates of prevalence in the source country (incorporating under-reporting), and 4) the import risk score of the target based on estimates from our analysis. Then, we produce short term estimates of the daily incidence of the VoC in the target country by means of a Renewal process [38, 39, 40], in which we take into account both the introductions of seeds from the source country and the local epidemic dynamics caused by secondary cases. The renewal equation approach comes with three main advantages with respect to other models,



**Figure 2: Quantifying the pandemic potential of the BA.1 lineage.** (A) Phylogenetic reconstruction and estimation of the most recent common ancestor (MRCA), identified South Africa on 28 October 2021 (95% HPD: 20 October–5 November) as the most likely MRCA. (B) Import risk map: countries are colored by their probability to import infectious individuals carrying the BA.1 lineage. (C, D) Projected weekly incidence in Estonia and the U.S. obtained from epidemic modeling, under different  $R_t$  scenarios indicated by coloured lines. Line thickness represents the range between the minimum and maximum assumed values of underreporting in the source country (here South Africa). Points represent the observed incidence. (E) Case counts simulated using the  $R_t$  scenario that corresponds to the mean growth rate from the phylogenetic analysis. For each country, the date of the first reported case is indicated with a grey circle.

such as SIR [41]. In fact, 1) it does not require to include in the dynamics the immunological status of the population in the target country; 2) the VoC dynamics can be considered as independent from the ones of the co-circulating VoCs, thus avoiding the need of estimating additional parameters for concurring spreading processes; 3) the model explicitly includes the most relevant epidemiological observables, such as  $R_t$ , the serial interval distribution [42], and the immune escape of the VoC. For further details we refer to Materials and Methods and Supplementary Figs. S11–S12.

**Assessing the pandemic potential of emerging variants.** In Fig. 2 we show the result of each step described above in determining the genomic and epidemiological parameters of the BA.1 lineage and, accordingly, quantify its

pandemic potential. We refer the reader to Supplementary Figs. S13-S14 for a more detailed analysis of errors in these estimates. Figure 2A displays a time-resolved maximum clade credibility phylogeny of the lineage. Panel B is the map of import risk across the world. Panels C and D show, for two example countries, the simulated epidemic projections, plotted as weekly incidence. For each reproduction number, the shaded area represents the interval between the estimates derived using the minimum and maximum values of underreporting in the source country. Panel E provides model estimates of case counts in all considered countries.

Figure 3 shows the results obtained for the SARS-CoV-2 lineages B.1.1.7 (Alpha), B.1.617.2 (Delta), BA.1, BA.2 and BA.5 (Omicron). The date of the most recent common ancestors and the growth rate are shown, together with the temporal evolution of the number of expected cases around 50 countries (varies depending on available sequencing data; Alpha: 59, Delta: 55, BA.2: 51, BA.5: 49 countries). Point estimates of the mean and 95% HPD regions are further provided in Table 1.

To assess the prediction error of our workflow we computed the normalized root-mean-square error (nRMSE) between prediction scenarios and observations. The nRMSE is zero, if the observation lies in between the simulation scenarios. Otherwise, the nRMSE is the RMSE between observation and the closest prediction scenario, normalized to the range that is spanned by the observations in the respective target country (for details see Material and Methods). Figure 4 captures the absolute and relative frequency of countries according to their nRMSE. Our predictions are in very good agreement (nRMSE = 0) for Alpha in 81.4%, BA.1 in 53.1%, BA.2 in 52.9%, BA.5 in 49% and for Delta in 12.7% of all considered countries. Note that even though Delta has the smallest amount of countries with incidences falling within scenarios prediction, more than 75% of the countries have a nRMSE  $\leq 2.5$ .

SARS-CoV-2 Lineage	tMRCA [95% HPD]	Growth Rate [95% HPD]
B.1.1.7 (Alpha)	9 Sep 2020 [28 Aug 2020–19 Sep 2020]	0.097 [0.008–0.202]
B.1.617.2 (Delta)	24 Oct 2020 [5 Jul 2020–10 Oct 2020]	0.020 [0.008–0.033]
BA.1 (Omicron)	28 Oct 2020 [20 Oct 2020–5 Nov 2020]	0.582 [90.117–1.035]
BA.2 (Omicron)	23 Oct 2020 [4 Oct 2020–9 Nov 2020]	0.144 [0.046–0.262]
BA.5 (Omicron)	9 Jan 2021 [19 Dec 2020–29 Jan 2021]	0.113 [0.051–0.177]
BA.2.75 (Omicron)	4 Apr 2022 [9 Mar 2022–27 Apr 2022]	0.096 [0.037–0.162]

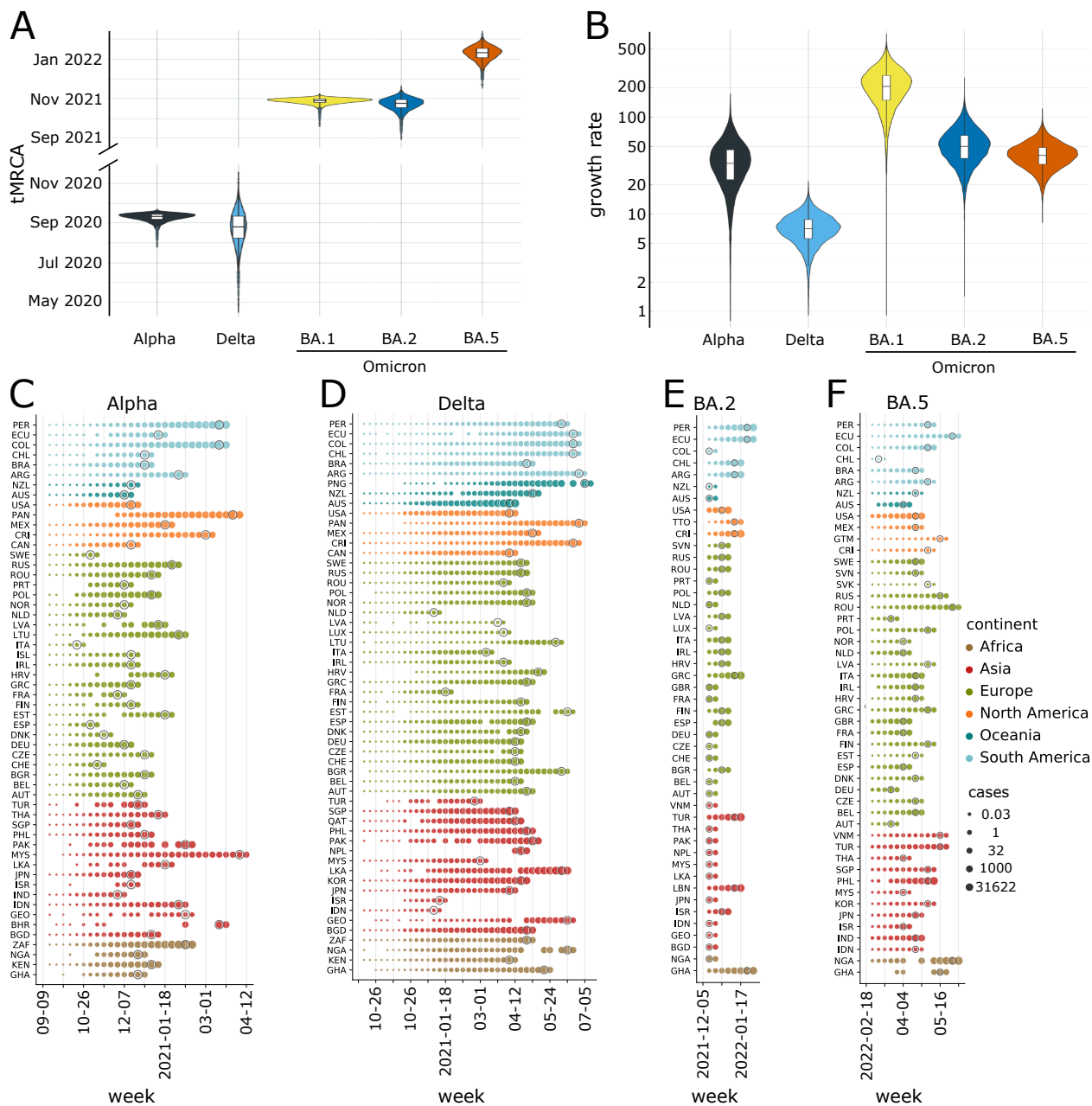
Table 1: Phylogenetic estimates of the time of most recent common ancestor (tMRCA) and daily growth rate for SARS-CoV-2 B.1.1.7 (Alpha), B.1.617.2 (Delta), BA.1, BA.2, BA.5 and BA.2.75 (Omicron) lineages.

## Discussion and conclusions

Here we have presented an integrated framework that combines phylogenetic analysis of genomic surveillance data with international human mobility data and large-scale epidemic modeling, in order to characterize in nearly real time the pandemic potential of an emerging variant. This concept is intended to provide quantitative indicators about the ability of a variant to escape population immunity acquired previous infections and/or vaccination, and quickly spread at a global level through human activities.

Our framework naturally deals with missing and noisy information to infer, through a Bayesian approach, the most likely origin – in space (on the country level) and time – of an emerging variant and its growth rate. Spatial and temporal coordinates are used to feed an analytical technique to estimate the probability that a given number of infectious individuals, departing from the country where the variant first appeared, travel to other countries with no exposure to it. This crucial step is based on international travel data, providing information about human movements between countries. Note that our approach is more powerful than naive estimates based only origin-destination pairs: in fact, we make full use of the knowledge we have about the underlying international travel network and its latent geometry [21, 43, 44], known to play a crucial role to amplify the spread of an emerging pathogen [18]. The last stage of our framework is to use importation risk to quantify the number of imported infectious cases to each country and, accordingly, estimate the consequent unfolding of the epidemic due to the emerging variant. The epidemic model is intended to quantify undetected infections that occur well before the first genomic sequence is isolated from a case in a country. Note that it is also possible to estimate the time at which the emerging variant will become dominant in the destination country, as shown in Supplementary Figs. S15-S16. The estimation is based on a logistic growth equation for the relative fraction of a new variant. These predictions will be less accurate if growth advantages in different countries are heterogeneous, for example due to immune escape. In this scenario, they would require calibration for each variant.

AUGUST 19, 2022



**Figure 3: Pan-viral pandemic potential: comparing multiple lineages.** (A–B) tMRCA and growth rate estimates for Alpha, Delta, BA.1, BA.2 and BA.5 from phylogenetic analysis. (C–F) Estimates of case numbers in all the considered countries for the same variants. For each lineage and country, the epidemic simulation starts at the time of infection  $t_0$  of the first undetected case as identified using the phylogenetic analysis. The simulation stops at the third date at which sequences belonging to the considered lineages are greater than zero. Results are provided in logarithmic scale and times at which the first case is reported are marked with grey circles.

Only the early phase of spread of a new lineage is estimated and the proposed model can safely take advantages of assumptions like a homogeneous mixing and the lack of feedback loops in the epidemic dynamics.

We have validated our integrated framework on existing variants, including B.1.1.7 (Alpha), B.1.617.2 (Delta), BA.1, BA.2 and BA.5 (Omicron), finding excellent agreement with independent estimates of the relevant phylogenetic and epidemiological parameters. By accounting for different scenarios in the progress of the epidemic in each country, we provide quantifiable indicators to inform decision makers and support pro-active policy interventions to mitigate

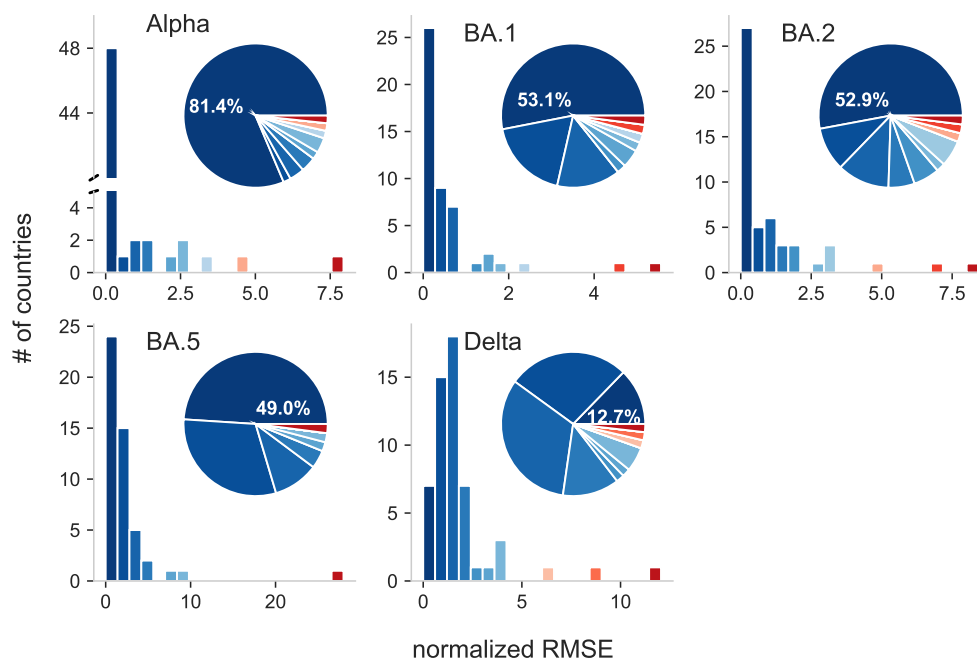


Figure 4: **Pandemic intelligence workflow error estimation.** Absolute (bars) and relative frequency (segments) of countries according to their normalized root-mean-square error (nRMSE) for the Alpha (B.1.1.7), Delta (B.1.617.2), BA.1, BA.2, BA.5 (Omicron) lineages. The normalized RMSE is zero if the number of infected people evaluated from data is inside the range spanned by the epidemic scenarios. Otherwise, it is the RMSE between observed incidence and the incidence of the closest epidemic scenario, normalized to the range spanned by observed incidences in the respective country. The order and color of the bars and segments is identical, i.e. the bars serves as color legend for the segments. For example, dark blue always represents the number or percentage of countries with the smallest nRMSE.

the potentially harmful effect of an emerging variant. For the current variant of most concern, BA.5, we estimate that its most recent common ancestor existed in early January 2022 (9 January 2022, 95% HPD: 19 December 2021–29 January 2022), with a daily growth rate of 0.113 (95% HPD: 0.051–0.177).

Overall, our findings show that it is possible to aim at pandemic intelligence, even with partial and noisy data. We must caution that the estimates of the pandemic potential of an emerging SARS-CoV-2 variant are largely driven by the uncertainty in the spatio-temporal coordinates of its origin. Our most unreliable estimates are obtained for countries where genomic surveillance is poor, propagating uncertainties to short-term projections which, in turn, exhibit larger variability. Importantly, note that only the validation of our scenario predictions relies on large enough sequencing rates in the target country, but not its application. That means our framework is perfectly suitable for low- and middle-income countries with little to no genomic surveillance.

Failures in international cooperation with a view to finding global solutions have undoubtedly shaped the COVID-19 pandemic. We have provided robust evidence that epidemic intelligence at country level is not enough, alone, to contrast the pandemic of respiratory pathogens such as SARS-CoV-2, in the absence of well-coordinated genomic surveillance – especially in low-income and middle-income countries currently lacking and adequate response capacity [45] – and global projections of variant’s pandemic potential. Our approach is inherently integrated and scalable, adding to ongoing modeling efforts and pan-viral analyses [46, 47, 48, 49, 22] and responding to global calls for coordinated action [45, 50, 51]. The data-driven approach provides a vital step in the path towards pandemic intelligence – where the interconnected and interdependent nature of human activities [21, 18, 52] is naturally accounted for at a global level – as well a means of enhancing global preparedness against future emerging variants.

## References

- [1] Michael Worobey, Joshua I Levy, Lorena Malpica Serrano, Alexander Crits-Christoph, Jonathan E Pekar, Stephen A Goldstein, Angela L Rasmussen, Moritz UG Kraemer, Chris Newman, Marion PG Koopmans, et al.



- The huanan seafood wholesale market in wuhan was the early epicenter of the covid-19 pandemic. *Science*, 2022.
- [2] Juanjuan Zhang, Maria Litvinova, Yuxia Liang, Yan Wang, Wei Wang, Shanlu Zhao, Qianhui Wu, Stefano Merler, Cécile Viboud, Alessandro Vespignani, et al. Changes in contact patterns shape the dynamics of the covid-19 outbreak in china. *Science*, 368(6498):1481–1486, 2020.
- [3] Matteo Chinazzi, Jessica T Davis, Marco Ajelli, Corrado Gioannini, Maria Litvinova, Stefano Merler, Ana Pastore y Piontti, Kunpeng Mu, Luca Rossi, Kaiyuan Sun, et al. The effect of travel restrictions on the spread of the 2019 novel coronavirus (covid-19) outbreak. *Science*, 368(6489):395–400, 2020.
- [4] Moritz UG Kraemer, Chia-Hung Yang, Bernardo Gutierrez, Chieh-Hsi Wu, Brennan Klein, David M Pigott, Open COVID-19 Data Working Group?, Louis du Plessis, Nuno R Faria, Ruoran Li, et al. The effect of human mobility and control measures on the covid-19 epidemic in china. *Science*, 368(6490):493–497, 2020.
- [5] Nicola Perra. Non-pharmaceutical interventions during the COVID-19 pandemic: A review. *Physics Reports*, 913:1–52, may 2021.
- [6] John S Tregoning, Katie E Flight, Sophie L Higham, Ziyin Wang, and Benjamin F Pierce. Progress of the covid-19 vaccine effort: viruses, vaccines and variants versus efficacy, effectiveness and escape. *Nature Reviews Immunology*, 21(10):626–636, 2021.
- [7] Rishi R. Goel, Mark M. Painter, Sokratis A. Apostolidis, Divij Mathew, Wenzhao Meng, Aaron M. Rosenfeld, Kendall A. Lundgreen, Arnold Reynaldi, David S. Khoury, Ajinkya Pattekar, Sigrid Gouma, Leticia Kuri-Cervantes, Philip Hicks, Sarah Dysinger, Amanda Hicks, Harsh Sharma, Sarah Herring, Scott Korte, Amy E. Baxter, Derek A. Oldridge, Josephine R. Giles, Madison E. Weirick, Christopher M. McAllister, Moses Awofolaju, Nicole Tanenbaum, Elizabeth M. Drapeau, Jeanette Dougherty, Sherea Long, Kurt D’Andrea, Jacob T. Hamilton, Maura McLaughlin, Justine C. Williams, Sharon Adamski, Oliva Kuthuru, Ian Frank, Michael R. Betts, Laura A. Vella, Alba Grifoni, Daniela Weiskopf, Alessandro Sette, Scott E. Hensley, Miles P. Davenport, Paul Bates, Eline T. Luning Prak, Allison R. Greenplate, and E. John Wherry. mRNA vaccines induce durable immune memory to SARS-CoV-2 and variants of concern. *Science*, 374(6572), dec 2021.
- [8] David S Khoury, Deborah Cromer, Arnold Reynaldi, Timothy E Schlub, Adam K Wheatley, Jennifer A Juno, Kanta Subbarao, Stephen J Kent, James A Triccas, and Miles P Davenport. Neutralizing antibody levels are highly predictive of immune protection from symptomatic sars-cov-2 infection. *Nature medicine*, pages 1–7, 2021.
- [9] Benjamin F Maier and Dirk Brockmann. Effective containment explains subexponential growth in recent confirmed covid-19 cases in china. *Science*, 368(6492):742–746, 2020.
- [10] Laura Di Domenico, Chiara E. Sabbatini, Pierre-Yves Boëlle, Chiara Poletto, Pascal Crépey, Juliette Paireau, Simon Cauchemez, François Beck, Harold Noel, Daniel Lévy-Bruhl, and Vittoria Colizza. Adherence and sustainability of interventions informing optimal control against the COVID-19 pandemic. *Communications Medicine*, 1(1), December 2021.
- [11] Katrina A Lythgoe, Matthew Hall, Luca Ferretti, Mariateresa de Cesare, George MacIntyre-Cockett, Amy Trebes, Monique Andersson, Newton Otecko, Emma L Wise, Nathan Moore, et al. Sars-cov-2 within-host diversity and transmission. *Science*, 372(6539), 2021.
- [12] Tyler N Starr, Allison J Greaney, Sarah K Hilton, Daniel Ellis, Katharine HD Crawford, Adam S Dingens, Mary Jane Navarro, John E Bowen, M Alejandra Tortorici, Alexandra C Walls, et al. Deep mutational scanning of sars-cov-2 receptor binding domain reveals constraints on folding and ace2 binding. *Cell*, 182(5):1295–1310, 2020.
- [13] Allison J Greaney, Tyler N Starr, Pavlo Gilchuk, Seth J Zost, Elad Binshtein, Andrea N Loes, Sarah K Hilton, John Huddleston, Rachel Eguia, Katharine HD Crawford, et al. Complete mapping of mutations to the sars-cov-2 spike receptor-binding domain that escape antibody recognition. *Cell host & microbe*, 29(1):44–57, 2021.
- [14] Tyler N Starr, Allison J Greaney, Amin Addetia, William W Hannon, Manish C Choudhary, Adam S Dingens, Jonathan Z Li, and Jesse D Bloom. Prospective mapping of viral mutations that escape antibodies used to treat covid-19. *Science*, 371(6531):850–854, 2021.
- [15] Nicholas G Davies, Sam Abbott, Rosanna C Barnard, Christopher I Jarvis, Adam J Kucharski, James D Munday, Carl AB Pearson, Timothy W Russell, Damien C Tully, Alex D Washburne, et al. Estimated transmissibility and impact of SARS-CoV-2 lineage B. 1.1. 7 in England. *Science*, 372(6538), 2021.
- [16] Laith J Abu-Raddad, Hiam Chemaitelly, and Adeel A Butt. Effectiveness of the BNT162b2 Covid-19 Vaccine against the B. 1.1. 7 and B. 1.351 Variants. *New England Journal of Medicine*, 385(2):187–189, 2021.

AUGUST 19, 2022

- [17] Nick Andrews, Julia Stowe, Freja Kirsebom, Samuel Toffa, Tim Rickeard, Eileen Gallagher, Charlotte Gower, Meaghan Kall, Natalie Groves, Anne-Marie O'Connell, et al. Covid-19 vaccine effectiveness against the Omicron (B. 1.1. 529) variant. *New England Journal of Medicine*, 386(16):1532–1546, 2022.
- [18] Jessica T. Davis, Matteo Chinazzi, Nicola Perrà, Kunpeng Mu, Ana Pastore y Piontti, Marco Ajelli, Natalie E. Dean, Corrado Gioannini, Maria Litvinova, Stefano Merler, Luca Rossi, Kaiyuan Sun, Xinyue Xiong, Ira M. Longini, M. Elizabeth Halloran, Cécile Viboud, and Alessandro Vespignani. Cryptic transmission of SARS-CoV-2 and the first COVID-19 wave. *Nature*, oct 2021.
- [19] Katelyn Gostic, Ana CR Gomez, Riley O Mummah, Adam J Kucharski, and James O Lloyd-Smith. Estimated effectiveness of symptom and risk screening to prevent the spread of COVID-19. *Elife*, 9, 2020.
- [20] Duygu Balcan, Vittoria Colizza, Bruno Gonçalves, Hao Hu, José J Ramasco, and Alessandro Vespignani. Multiscale mobility networks and the spatial spreading of infectious diseases. *Proceedings of the National Academy of Sciences*, 106(51):21484–21489, 2009.
- [21] Dirk Brockmann and Dirk Helbing. The Hidden Geometry of Complex, Network-Driven Contagion Phenomena. *Science*, 342(6164):1337–1342, dec 2013.
- [22] O. Eales, L. de Oliveira Martins, A.J. Page, et al. Dynamics of competing sars-cov-2 variants during the omicron epidemic in england. *Nature Communications*, 13:4375, 2022.
- [23] Enhancing response to Omicron SARS-CoV-2 variant — who.int. [https://www.who.int/publications/m/item/enhancing-readiness-for-omicron-\(b.1.1.529\)-technical-brief-and-priority-actions-for-member-states](https://www.who.int/publications/m/item/enhancing-readiness-for-omicron-(b.1.1.529)-technical-brief-and-priority-actions-for-member-states). [Accessed 30-Jul-2022].
- [24] Alexander Wilhelm, Marek Widera, Katharina Grikscheit, Tuna Toptan, Barbara Schenk, Christiane Pallas, Melinda Metzler, Niko Kohmer, Sebastian Hoehl, Fabian A Helfritz, et al. Reduced neutralization of sars-cov-2 omicron variant by vaccine sera and monoclonal antibodies. *MedRxiv*, 2021.
- [25] Report 48 - The value of vaccine booster doses to mitigate the global impact of the Omicron SARS-CoV-2 variant — imperial.ac.uk. <https://www.imperial.ac.uk/mrc-global-infectious-disease-analysis/covid-19/report-48-global-omicron/>. [Accessed 30-Jul-2022].
- [26] Annika Rössler, Lydia Riepler, David Bante, Dorothee von Laer, and Janine Kimpel. Sars-cov-2 omicron variant neutralization in serum from vaccinated and convalescent persons. *New England Journal of Medicine*, 386(7):698–700, 2022.
- [27] Henning Gruell, Kanika Vanshylla, Pinkus Tober-Lau, David Hillus, Philipp Schommers, Clara Lehmann, Florian Kurth, Leif E Sander, and Florian Klein. mRNA booster immunization elicits potent neutralizing serum activity against the SARS-CoV-2 Omicron variant. *Nature medicine*, 28(3):477–480, 2022.
- [28] Raquel Viana, Sikhulile Moyo, Daniel G Amoako, Houriiyah Tegally, Cathrine Scheepers, Christian L Althaus, Ugochukwu J Anyaneji, Phillip A Bester, Maciej F Boni, Mohammed Chand, et al. Rapid epidemic expansion of the sars-cov-2 omicron variant in southern africa. *Nature*, 603(7902):679–686, 2022.
- [29] Laura Espenhain, Tjede Funk, Maria Overvad, Sofie Marie Edslev, Jannik Fonager, Anna Cäcilia Ingham, Morten Rasmussen, Sarah Leth Madsen, Caroline Hjorth Espersen, Raphael N Sieber, et al. Epidemiological characterisation of the first 785 sars-cov-2 omicron variant cases in denmark, december 2021. *Eurosurveillance*, 26(50):2101146, 2021.
- [30] Lin T Brandal, Emily MacDonald, Lamprini Veneti, Tine Ravlo, Heidi Lange, Umaer Naseer, Siri Feruglio, Karoline Bragstad, Olav Hungnes, Liz E Ødeskaug, et al. Outbreak caused by the SARS-CoV-2 Omicron variant in Norway, November to December 2021. *Eurosurveillance*, 26(50):2101147, 2021.
- [31] Stefan Elbe and Gemma Buckland-Merrett. Data, disease and diplomacy: Gisaïd's innovative contribution to global health. *Global challenges*, 1(1):33–46, 2017.
- [32] Yuelong Shu and John McCauley. Gisaïd: Global initiative on sharing all influenza data—from vision to reality. *Eurosurveillance*, 22(13):30494, 2017.
- [33] Shruti Khare, Céline Gurry, Lucas Freitas, Mark B Schultz, Gunter Bach, Amadou Diallo, Nancy Akite, Jose Ho, Raphael TC Lee, Winston Yeo, et al. Gisaïd's role in pandemic response. *China CDC Weekly*, 3(49):1049, 2021.
- [34] Stephen W Attwood, Sarah C Hill, David M Aanensen, Thomas R Connor, and Oliver G Pybus. Phylogenetic and phylodynamic approaches to understanding and combating the early SARS-CoV-2 pandemic. *Nature Reviews Genetics*, pages 1–16, 2022.

- [35] Mattia Manica, Alfredo De Bellis, Giorgio Guzzetta, Pamela Mancuso, Massimo Vicentini, Francesco Venturelli, Alessandro Zerbini, Eufemia Bisaccia, Maria Litvinova, Francesco Menegale, et al. Intrinsic generation time of the sars-cov-2 omicron variant: an observational study of household transmission. *The Lancet Regional Health-Europe*, 19:100446, 2022.
- [36] Official Airline Guide. Official Airline Guide: global airline schedules data, 2022.
- [37] Pascal P. Klamser, Adrian Zachariae, Benjamin F. Maier, Olga Baranov, Clara Jongen, Frank Schlosser, and Dirk Brockmann. Inferring country specific import risk of diseases from the World-Aviation-Network. *In preparation*, 2022.
- [38] Christophe Fraser. Estimating Individual and Household Reproduction Numbers in an Emerging Epidemic. *PLOS ONE*, 2(8):e758, 2007.
- [39] Anne Cori, Neil M. Ferguson, Christophe Fraser, and Simon Cauchemez. A New Framework and Software to Estimate Time-Varying Reproduction Numbers During Epidemics. *American Journal of Epidemiology*, 178(9):1505–1512, November 2013.
- [40] William D. Green, Neil M. Ferguson, and Anne Cori. Inferring the reproduction number using the renewal equation in heterogeneous epidemics. *Journal of The Royal Society Interface*, 19(188):20210429, 2022.
- [41] Matt J Keeling and Ken T.D Eames. Networks and epidemic models. *Journal of The Royal Society Interface*, 2(4):295–307, sep 2005.
- [42] Luca Ferretti, Alice Ledda, Chris Wymant, Lele Zhao, Virginia Ledda, Lucie Abeler, Michelle Kendall, Anel Nurtay, Hao-Yuan Cheng, Ta-Chou Ng, Hsien-Ho Lin, Rob Hinch, Joanna Masel, A Marm Kilpatrick, and Christophe Fraser. The timing of COVID-19 transmission. *medRxiv*, page 16, 2020.
- [43] Chittaranjan Hens, Uzi Harush, Simi Haber, Reuven Cohen, and Baruch Barzel. Spatiotemporal signal propagation in complex networks. *Nature Physics*, 15(4):403–412, 2019.
- [44] Marian Boguna, Ivan Bonamassa, Manlio De Domenico, Shlomo Havlin, Dmitri Krioukov, and M Ángeles Serrano. Network geometry. *Nature Reviews Physics*, 3(2):114–135, 2021.
- [45] Nelson Aghoghov Evaborhene. A strong and independent Africa CDC would benefit the world. *The Lancet*, jul 2022.
- [46] Mary Bushman, Rebecca Kahn, Bradford P Taylor, Marc Lipsitch, and William P Hanage. Population impact of SARS-CoV-2 variants with enhanced transmissibility and/or partial immune escape. *Cell*, 184(26):6229–6242, 2021.
- [47] Finlay Campbell, Brett Archer, Henry Laurenson-Schafer, Yuka Jinnai, Franck Konings, Neale Batra, Boris Pavlin, Katelijn Vandemaale, Maria D Van Kerkhove, Thibaut Jombart, et al. Increased transmissibility and global spread of SARS-CoV-2 variants of concern as at June 2021. *Eurosurveillance*, 26(24):2100509, 2021.
- [48] Jalen Singh, Pranav Pandit, Andrew G McArthur, Arinjay Banerjee, and Karen Mossman. Evolutionary trajectory of SARS-CoV-2 and emerging variants. *Virology journal*, 18(1):1–21, 2021.
- [49] Christiaan H van Dorp, Emma E Goldberg, Nick Hengartner, Ruian Ke, and Ethan O Romero-Severson. Estimating the strength of selection for new SARS-CoV-2 variants. *Nature communications*, 12(1):1–13, 2021.
- [50] Lorenzo Subissi, Anne von Gottberg, Lipi Thukral, Nathalie Worp, Bas B Oude Munnink, Surabhi Rathore, Laith J Abu-Raddad, Ximena Aguilera, Erik Alm, Brett N Archer, et al. An early warning system for emerging SARS-CoV-2 variants. *Nature medicine*, pages 1–6, 2022.
- [51] World Health Organization et al. Strengthening pandemic preparedness planning for respiratory pathogens: policy brief, 27 april 2022. In *Strengthening pandemic preparedness planning for respiratory pathogens: policy brief, 27 April 2022*. World Health Organization, 2022.
- [52] Colin J. Carlson, Gregory F. Albery, Cory Merow, Christopher H. Trisos, Casey M. Zipfel, Evan A. Eskew, Kevin J. Olival, Noam Ross, and Shweta Bansal. Climate change increases cross-species viral transmission risk. *Nature*, 607(7919):555–562, jul 2022.

## Acknowledgments

The authors acknowledge the GISAID initiative and all the authors from the originating laboratories where genetic sequence data were generated for sharing such data through GISAID, which has made this work possible.

AUGUST 19, 2022

## **Funding**

A.D.N. acknowledges support from the Department for Environment, Food and Rural Affairs (Defra), United Kingdom [research grant: SE2945], and the Biotechnology and Biological Research Council (BBSRC), United Kingdom [project: BBS/E/I/0007035].

## **Author contributions**

M.D.D. designed the study; P.K., V.D., F.D.L., A.Z., S.B., A.D.N., M.H. and L.F. performed the numerical experiments; all authors analyzed the data and wrote the manuscript.

## **Competing interests**

The authors declare no competing interests.

## **Data and materials availability**

Both the data and analysis material will be available online at Zenodo upon publication. This work is licensed under a Creative Commons Attribution 4.0 International (CC BY 4.0) license, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. To view a copy of this license, visit <https://creativecommons.org/licenses/by/4.0/>. This license does not apply to figures/photos/artwork or other content included in the article that is credited to a third party; obtain authorization from the rights holder before using such material.

# SI Appendix

## Material and Methods and Supplementary Information

### “Enhancing global preparedness during an ongoing pandemic from partial and noisy data”

Pascal Klamsner<sup>1,2,†</sup>, Valeria d’Andrea<sup>3,†</sup>, Francesco Di Lauro<sup>4</sup>, Adrian Zachariae<sup>1,2</sup>, Sebastiano Bontorin<sup>3,5</sup>, Antonello di Nardo<sup>6</sup>, Matthew Hall<sup>4</sup>, Benjamin F. Maier<sup>1,2</sup>, Luca Ferretti<sup>4</sup>, Dirk Brockmann<sup>1,2</sup>, Manlio De Domenico<sup>7,8,\*</sup>

<sup>1</sup>Robert Koch-Institute, Nordufer 20, 13353 Berlin, Germany

<sup>2</sup>Institute for Theoretical Biology, Humboldt-University of Berlin, Philippstr. 13, D-10115 Berlin, Germany

<sup>3</sup>Fondazione Bruno Kessler, Via Sommarive 18, 38123 Povo (TN), Italy

<sup>4</sup>Big Data Institute, University of Oxford, Li Ka Shing Centre for Health Information and Discovery, Oxford, UK

<sup>5</sup>Department of Physics, University of Trento, Via Sommarive 14, 38123 Povo (TN), Italy

<sup>6</sup>The Pirbright Institute, Ash Road, Pirbright, Woking, Surrey, United Kingdom

<sup>7</sup>Department of Physics and Astronomy, G. Galilei, University of Padua, Via Francesco Marzolo 8, 35131, Padua, Italy

<sup>8</sup>Padua Center for Network Medicine, University of Padua, Via Francesco Marzolo 8, 35131, Padua, Italy

<sup>†</sup>Contributed equally to this work.

<sup>\*</sup>Corresponding author. E-mail: [manlio.dedomenico@unipd.it](mailto:manlio.dedomenico@unipd.it)

## Contents

<b>I Phylogenetic Reconstruction</b>	<b>2</b>
I.1 Genomic dataset compilation . . . . .	2
I.2 Phylogenetic estimates of epidemiological parameters . . . . .	2
I.3 Estimates based on epidemic modeling . . . . .	3
<b>II Import Risk estimation</b>	<b>5</b>
II.1 International travel dataset compilation . . . . .	5
II.2 Quantifying the Import Risk . . . . .	5
II.3 Relation to distance and arrival time . . . . .	8
II.4 Alternative distance measures . . . . .	8
II.5 Data for arrival time and outbreak region . . . . .	9
II.6 Outbreak detection based on 1st count GISAID data . . . . .	10
<b>III Epidemic Scenarios</b>	<b>13</b>
III.1 Renewal equation . . . . .	13
III.2 A fully worked out example: the Alpha variant . . . . .	14
III.3 Prediction error . . . . .	16
<b>IV Variant Dominance</b>	<b>19</b>
<b>V Information Distance</b>	<b>22</b>

## Materials and Methods

### I Phylogenetic Reconstruction

#### I.1 Genomic dataset compilation

We retrieved all SARS-CoV-2 sequences belonging to the Alpha B.1.1.7, Delta B.1.617.2, Omicron BA.1, BA.2, BA.5, and BA.2.75 lineages from GISAID. Each genomic dataset was filtered by only retaining those sequences that were generated from cases reported during the initial wave and from the country of evolutionary origin, up to a total of 100 sequences per lineage. We then generated 3 alignments using MAFFT 7.505 [1], each comprised of 20%, 50% and 100% of the total number of sequences, which were subsequently cleaned by trimming the 5' and 3' untranslated regions and gap-only sites.

#### I.2 Phylogenetic estimates of epidemiological parameters

We performed a common Bayesian evolutionary reconstruction of timed phylogenetic history using BEAST 1.10.5 [2] that was source compiled from its github repository (<https://github.com/beast-dev/beast-mcmc>). We modelled the nucleotide substitution process according to an  $HKY85 + \Gamma$  parameterisation, setting a strict molecular clock and an exponential growth model as coalescent prior. We used a  $Lognormal(\mu = 9 \times 10^{-4}, \sigma^2 = 1 \times 10^{-5})$  prior for the molecular rate of evolution, a  $Laplace(\mu = 0, b = 100)$  prior for the rate of exponential growth and a  $Lognormal(\mu = 5.7, \sigma^2 = 2.3)$  prior for the exponentially growing viral population size. We further set an initial calibration for the time of the most recent common ancestor (tMRCA) at an age of  $\sim 6$  months before the most recent sample included in the alignment. All the remaining priors were left at their default values.

Bayesian inference through Markov chain Monte Carlo (MCMC) was performed for  $2 \times 10^8$  generations, sampling every 20,000 generations and using the BEAGLE 3.1.2 library to increase computational performance [3]. MCMC

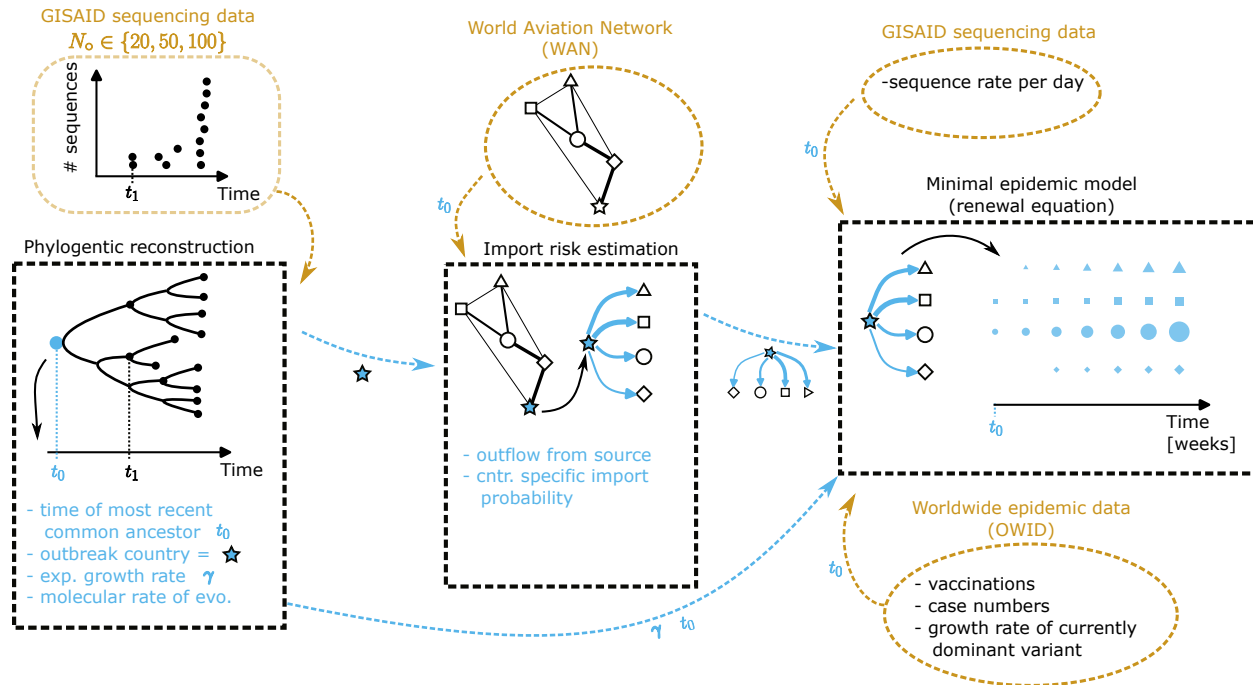


Figure S1: **Schematic mechanistic pipeline workflow.** The three pieces of our pipeline (black dashed boxes) are illustrated and what input they get from external data sources (orange colored) or from the output of earlier parts of the pipeline (blue arrows). The  $t_0$  at the orange arrows means that data from the external sources is used at or prior to  $t_0$ , which is the estimated time of the most recent ancestor (the output of the first part of the pipeline, the phylogenetic reconstruction).

convergence and mixing properties were inspected using Tracer 1.7.2 [4] to ensure that effective sample size (ESS) values associated with estimated parameters were all  $>200$ . After discarding 10% of sampled trees as burn-in, estimates of the growth rate, molecular clock and tMRCAs were extracted along with their posterior distributions (Figure S2).

### I.3 Estimates based on epidemic modeling

We obtain an independent estimate for  $t_0$ , the time of the first unreported case, and for other epidemic parameters, such as the effective reproduction number and the generation interval. By indicating with  $I(t)$  the number of infected individuals at time  $t$  and with  $D(t)$  the number of deaths, we consider the stage with the co-circulation of an existing variant  $v$  and the emerging one  $\omega$ . Since we consider the final stage of the contagions due to  $v$  and the early stage of the contagions due to  $\omega$ , we approximate the epidemic evolution by

$$\begin{aligned} I(t_0 + \Delta t) &= I_v(t_0 + \Delta t) + I_\omega(t_0 + \Delta t) = \\ &= I_v(t_0)R_v(t_0)^{\Delta t/GI_v} + I_\omega(t_0)R_\omega(t_0)^{\Delta t/GI_\omega}, \end{aligned} \quad (S1)$$

where  $I_x(t)$  is the number of infections due to variant  $x$  at time  $t$ ,  $R_x(t_0)$  is the effective reproduction number at time  $t_0$  and  $GI_x$  is the generation interval. Similarly, the deaths due to the co-circulating variants are approximated by

$$D_v(t_0 + \Delta t + \tau_v) = I_v(t_0 + \Delta t) \times \text{IFR}_v, \quad (S2)$$

$$D_\omega(t_0 + \Delta t + \tau_\omega) = I_\omega(t_0 + \Delta t) \times \text{IFR}_\omega, \quad (S3)$$

$$D(t) = D_v(t) + D_\omega(t) \quad (S4)$$

where  $\text{IFR}_x$  denotes the infection fatality rate of variant  $x$  and  $\tau_x$  is the lag between infection and death. To fit the unknown parameters, i.e. the ones related to variant  $\omega$ , we use particle swarm optimization [5] to minimize the loss function

$$\phi(\theta) = \frac{1}{2} \frac{\sqrt{\text{Var}[\log(1 + I(t)) - \log(1 + I_{obs}(t))]} + \frac{1}{2} \frac{\sqrt{\text{Var}[D(t) - D_{obs}(t)]}}{\sqrt{\text{Var}[D_{obs}(t)]}}, \quad (S5)$$

where  $I_{obs}(t)$  and  $D_{obs}(t)$  are the number of infected individuals and deaths from empirical data [6],  $\text{Var}$  indicates the variance in time and  $\theta = \{t_0; R_\omega(t_0); GI_\omega; \text{IFR}_\omega; \tau_\omega\}$  is the vector of the epidemiological parameters characterizing the emerging variant, for which we obtain a joint probability distribution.

AUGUST 19, 2022

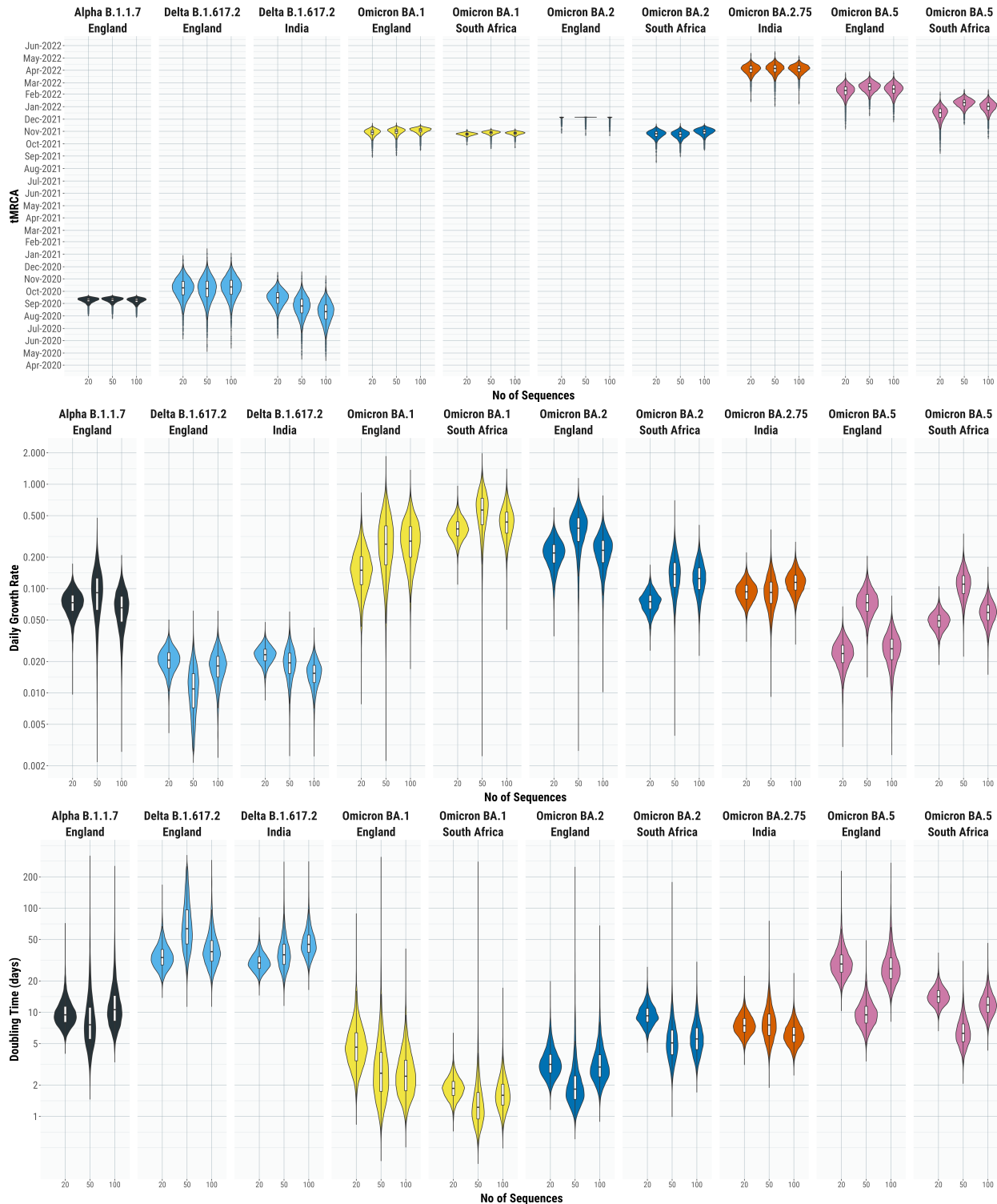


Figure S2: **Pan-variant phylogenetic analysis.** Posterior distributions of the time of the most recent common ancestor (tMRCA), daily growth rate and doubling time estimated for each of the Alpha B.1.1.7, Delta B.1.617.2, Omicron BA.1, BA.2, BA.5, and BA.2.75 SARS-Cov-2 lineages using alignments of 20, 50 and 100 sequences.



## II Import Risk estimation

### II.1 International travel dataset compilation

We retrieve the monthly seat capacities between airports from the OAG (Official Airline Guide). Note, that it does not represent the actual passengers that flew from airport A to B in one month, but the maximal capacity, i.e. how many could have travelled if all seats were occupied. It is therefore an upper limit for the passenger flux and we refer to it as the flow matrix  $\mathbf{F}$ , where  $F_{ij}$  describes the maximal passenger flow to  $i$  from  $j$ . We estimate the travelling population in the catchment area of an airport by  $N_i = F_i$ , with  $F_i = \sum_j F_{ji}$ , i.e. we assume that the population is proportional to the outflux of the airport. For each variant we use the WAN at the month of the outbreak day of the respective variant.

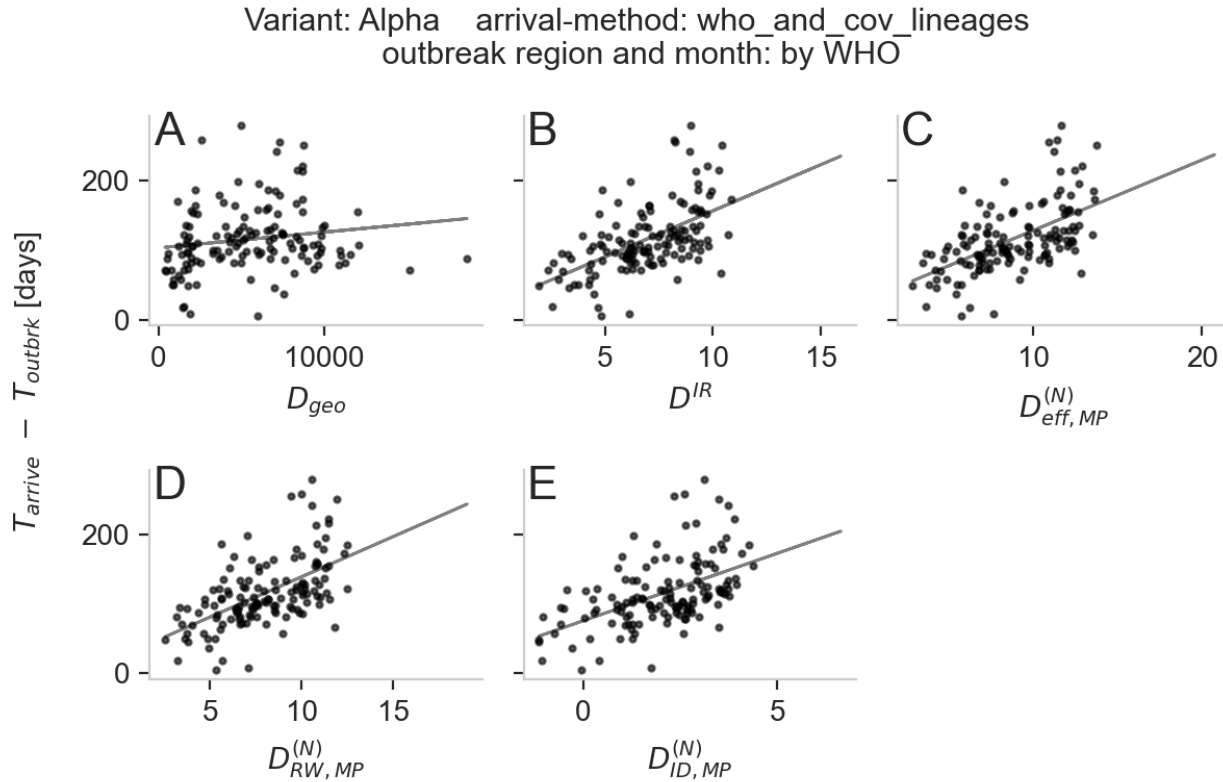


Figure S3: **Distance measures vs. arrivals for Alpha variant.** The distance measures are the geographic distance  $D_{geo}$  (A), the import risk distance  $D^{IR}$  (B), the effective distance  $D_{eff,MP}^{(N)}$  (C), the random walk distance  $D_{RW,MP}^{(N)}$  (D) and the information diffusion distance  $D_{ID,MP}^{(N)}$  (E) whereby the latter three (C, D, E) are generalized to weighted multiple paths.

### II.2 Quantifying the Import Risk

The import risk method is introduced in a separate study [9] where it is compared to another data-driven estimate. Here we present a short outline of the method. To know how many passengers leave at node  $j$  given they started at node  $i$ , we introduce the shortest path exit probability  $q(j|i)$  (SPEX). It is based on the shortest path tree of the effective distance [10], and combines the exit probability with all possible paths that end in  $j$ . The resulting import risk is therefore an extension of the SPEX.

In order to compute the SPEX we first define, with the flow matrix (maximal passenger flux)  $\mathbf{F}$  and the travelling population of the catchment area  $N_i$ , the transition matrix  $\mathbf{P}$ , where the element  $P_{ij} = F_{ij} / \sum_i F_{ij} = F_{ij} / F_j$  is the probability to transition to  $i$  from  $j$ . Now, the effective distance graph [10] is  $D_{ij} = d_0 - \log(P_{ij})$ , with  $d_0$  as the distance offset which we set to  $d_0 = 1$  (the larger  $d_0$  the more  $D_{ij}$  increases with increasing hop-distance). Let  $\mathbf{T}(n_0)$

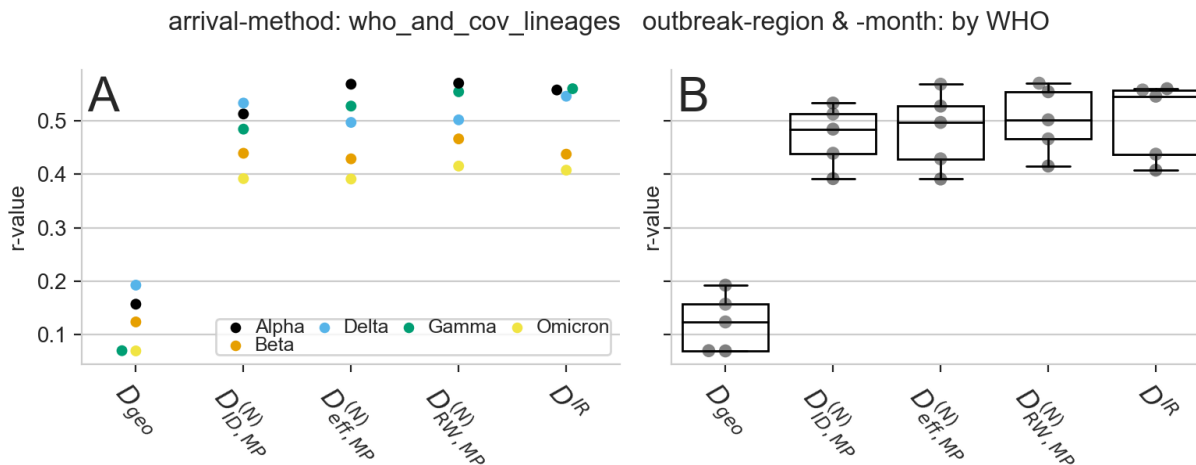


Figure S4: **Correlation comparison between different distance measures.** The distance measures are the geographic distance  $D_{geo}$ , the import risk distance  $D_{IR}$ , the effective distance  $D_{eff,MP}^{(N)}$ , the random walk distance  $D_{RW,MP}^{(N)}$  and the information diffusion distance  $D_{ID,MP}^{(N)}$  whereby the latter three are generalized to weighted multiple paths.

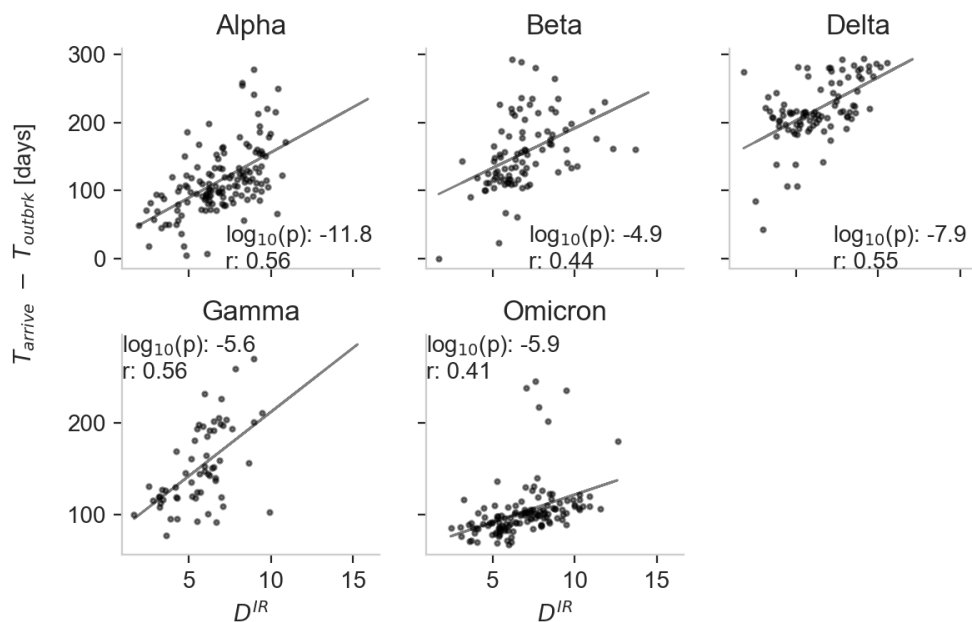


Figure S5: **Correlation of arrival times of variants with the import risk distance  $D^{IR}$ .** For the import risk distance  $D^{IR}(m|n_0) = -\log(p_\infty(m|n_0))$  the WAN of the WHO outbreak month is used and the WHO outbreak location as source country. The arrival time are taken from the "cov-lineages.org" [7, 8] project.

be the shortest path tree on  $\mathbf{D}$  for the point of origin  $n_0$ . With respect to node  $n$  the downstream nodes  $\Omega(n|n_0)$  are those nodes that can be reached from the source  $n_0$  through node  $n$  on  $\mathbf{T}(n_0)$ .

Now we compute the SPEX  $q(i|n_0)$  by assuming that all passenger that start at  $n_0$ , travel along the shortest path tree  $\mathbf{T}(n_0)$  and distribute to other airports according to their respective populations  $N_n$ . We assume that the exit probability at  $i$  is proportional to the ratio of the population at  $i$  (i.e.  $N_i$ ) to the population of all of  $i$ 's downstream nodes  $\sum_{n \in \Omega(i|n_0)} N_n$  plus  $N_i$ :

AUGUST 19, 2022

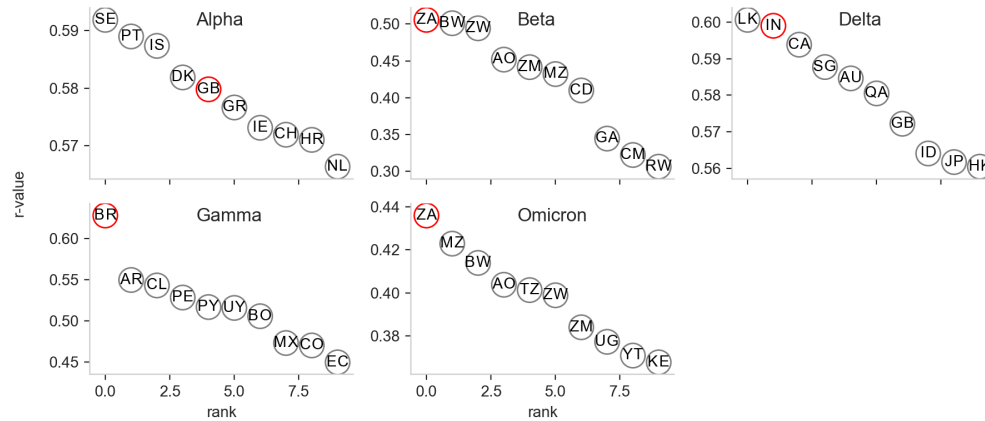


Figure S6: **Arrival prediction (r-value) for the 10 best outbreak candidate.** The r-value between the import risk distance  $D^{IR}(m|n_0) = -\log(p_\infty(m|n_0))$  and the arrival time for the 10 best ranked outbreak countries ( $n_0$ ). The 2 Letters in the circles are the countries ISO alpha-2 codes. The red circle marks the country declared as outbreak country by the WHO.

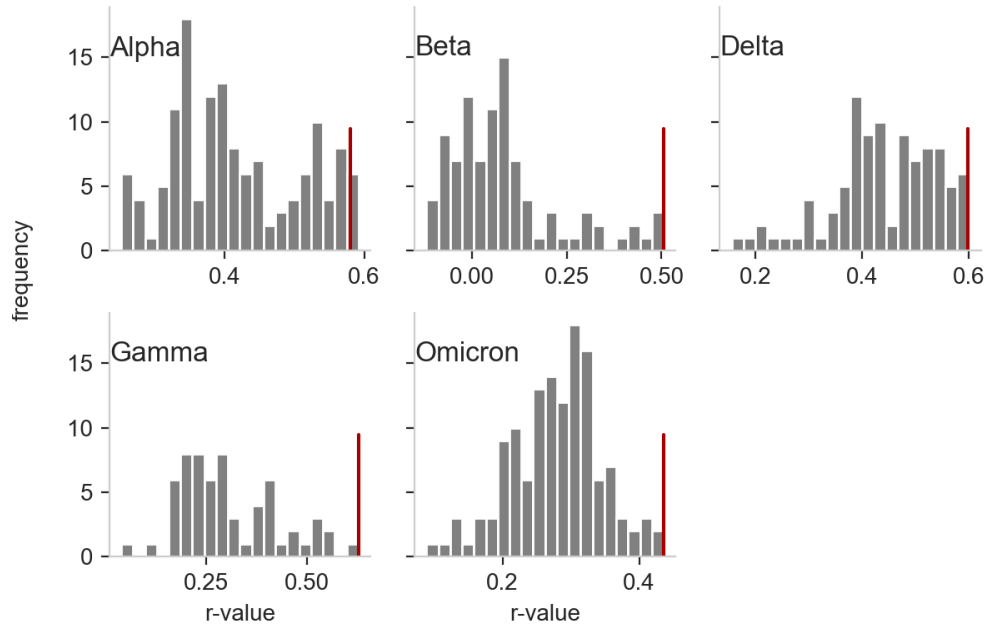


Figure S7: **Arrival prediction performane (r-value) for the outbreak country candidates.** The frequency of the r-value between the import risk distance  $D^{IR}(m|n_0) = -\log(p_\infty(m|n_0))$  and the arrival time for all possible outbreak countries. The red vertical line marks the r-value using the country declared as outbreak country by the WHO.

$$q(i|n_0) = \frac{N_i}{N_i + \sum_{n \in \Omega(i|n_0)} N_n} . \quad (S6)$$

Now, we use the SPEx on a random walk that starts at  $n_0$  and the walker exits at node  $i$  with probability  $q(i|n_0)$  or continues its walk with probability  $1 - q(i|n_0)$ . Thus, the probability to be at node  $m$  if the walker was before at node  $m - 1$  is

$$S(m, m - 1|n_0) = P_{m, m-1}(1 - q(m - 1|n_0)) . \quad (S7)$$

Consequently, the probability to take a path  $\Gamma$  starting at  $n_0$  and exiting at  $m$  is

$$p(\Gamma) = q(m|n_0) \prod_{(i,j) \in \Gamma} S(i,j|n_0). \quad (\text{S8})$$

The probability to exit at node  $m$  from all possible paths (of all possible lengths) is

$$p_\infty(m|n_0) = q(m|n_0) \left[ \sum_{k=1}^{\infty} \mathbf{S}^k(n_0) \right]_{m,n_0} \quad (\text{S9})$$

$$= q(m|n_0) \left[ (\mathbf{1} - \mathbf{S}(n_0))^{-1} - \mathbf{1} \right]_{m,n_0}. \quad (\text{S10})$$

Note that  $\mathbf{S}^k(n_0)_{m,n_0}$  is the probability sum of all paths that started in  $n_0$  and end after  $k$  steps in  $m$ . We aggregate all airports of the same country by computing the weighted mean with weights

$$w_n = \frac{N_n}{\sum_{m \in C(i)} N_m} \quad (\text{S11})$$

with  $C(n)$  as the set of airports that belong the same country as node  $n$  does.

### II.3 Relation to distance and arrival time

In order to assess the quality of the import risk, we compare it with the arrival time of past variants. Clearly, the higher the import risk to a country, the earlier it is to arrive and the direct relation between the probability of travel to a city  $m$  from a city  $n_0$  and the mean first arrival time  $t_1$  is

$$t_1(m|n_0) = d_0 - c \log(P(m|n_0)) \quad (\text{S12})$$

which is the effective distance [11, 12]. Thus, we define the import risk distance as

$$D^{IR}(m|n_0) = -\log(p_\infty(m|n_0)) \quad (\text{S13})$$

which is proportional to the mean first arrival time.

### II.4 Alternative distance measures

There are alternative measures to estimate the arrival time [10, 13, 14], and we want to compare our import risk distance to these established measures. However, please note that the alternative measures have a clear qualitative relation to the arrival time, but it is not possible to directly infer the number of passengers that travel between airports from them (what the import risk is especially designed for). The already introduced alternative measure is the effective distance [11] that uses the flow between airports to estimate the probability to travel from airport  $n$  to  $m$

$$d_{eff}(m,n) = d_0 - \log(P_{m,n}). \quad (\text{S14})$$

Now, the distance along a specific path  $\Gamma$  that connects  $m$  and  $n_0$  is the sum of the path elements distances

$$d_{eff}(\Gamma) = \sum_{(m,n) \in \Gamma} d_{eff}(m,n). \quad (\text{S15})$$

Finally the effective distance from airport  $n_0$  to  $m$ , also not directly connected airports, is the minimal effective distance of all possible paths  $\Omega(m,n_0)$  they are connected through

$$D_{eff}(m|n_0) = \min_{\Gamma \in \Omega(m,n_0)} (d_{eff}(\Gamma)). \quad (\text{S16})$$

An extension to the effective distance is the random-walk effective distance [14] that considers all possible paths connecting two airports  $\Omega(m,n_0)$  instead of only taking the dominant path with the shortest distance:

$$D_{RW}(m|n_0) = -\ln \left( \sum_{\Gamma \in \Omega(m,n_0)} e^{-d_{eff}(\Gamma)} \right). \quad (\text{S17})$$

Note that the sum of path distances via their exponential is due to the linkage to the arrival time as explained in [14].

We also add a comparison with a metric derived from Diffusion Distance [13] which exploits the definition of a random walk Laplacian on top of the WAN. We further explain this Information Distance  $D^{ID}$  in the dedicated section V.

## Country-Level aggregation.

The country-level aggregation of the import risk distance  $D^{IR}$  is done by first aggregating the import risk on country-level (as described in Sect. II.2) and then applying Eq. S13.

To aggregate the other distances ( $D_{eff}$ ,  $D_{RW}$ ) we could either take (along the line of  $D_{eff}$ ) the minimal distance between two countries (of all relevant airport pairs), or use a weighted multipath approach as used in the derivation of  $D_{RW}$ . We will highlight the latter in the following; however, we also computed the minimal measure and found that it is outperformed by the multipath distance (not shown, but it is the basic finding in [14]).

As shown in [12], the effective distance of two paths combined is

$$e^{-D_{eff}(\{\Gamma_a, \Gamma_b\})} = e^{-d_{eff}(\Gamma_a)} + e^{-d_{eff}(\Gamma_b)}. \quad (S18)$$

Thus, the multipath (MP) effective distance that considers all shortest paths from country  $S$  to  $M$  is:

$$D_{eff,MP}(M|S) = -\ln \left( \sum_{m \in M, s \in S} e^{-D_{eff}(m|s)} \right) \quad (S19)$$

with  $M$  as the set of all target airports in country  $M$  and  $S$  all source airports of country  $S$ .

Since the distance of source airports with a larger population are more important, we additionally weight the source airport with  $w_i = F_i / \sum_{s \in S} F_s$ , which represents the probability of an infected to start in location  $n$ . Now, we compute the population weighted multipath effective distance by

$$D_{eff,MP}^{(N)}(M|S) = -\ln \left( \sum_{m \in M, s \in S} w_s e^{-D_{eff}(m|s)} \right). \quad (S20)$$

Note that the weighting for the effective distance can be reformulated to

$$D_{eff,MP}^{(N)}(M|S) = -\ln \left( \sum_{m \in M, s \in S} w_s \prod_{k,l \in \Gamma_{m,s}} e^{-d_0} P_{k,l} \right) \quad (S21)$$

which corresponds to multiplying the probability to start at the source airport  $s$  to the first step of each path. Analogously the

## II.5 Data for arrival time and outbreak region

We compare the import risk to measured arrival times of different variants. Therefore, we need to define the outbreak-country and -month and the arrival times. We defined these variables in different ways.

**(I) external sources** Here we rely on peer reviewed [8] or official [15] sources. The outbreak country and the outbreak month are taken from the website of the World Health Organization (WHO) "Tracking SARS-CoV-2 variants" [15] and the arrival times of the variants Alpha, Beta, Delta, Gamma and Omicron were externally computed with "grinch" [8] and taken from their project website [7]. If arrival times are before the official outbreak they are removed from the analysis (for Delta=1, Gamma=1 and Omicron=19 countries are removed).

**(II) GISAID data** To also use the other variants to validate our import risk method we design a simple arrival time algorithm. First, we need to define the outbreak day. Instead of relying on an official definition from the WHO, we use GISAID data. The outbreak time  $T_{X,out}$  of variant  $X$  is defined by

$$T_{X,out} = T(F_X(t) \geq g \cdot \max(F_X)) - 30 \text{days} \quad (S22)$$

with  $F_X(t)$  being the fraction of variant  $X$  to all sequenced probes at time  $t$  and  $T(F_X(t) \geq g \cdot \max(F_X))$  the time when  $F_X(t)$  crosses the first time the threshold  $g \cdot \max(F_X)$  where  $g \in ]0, 1[$  and we set  $g = 0.025$ . In other words, the outbreak is defined by 30 days before the variant reached 2.5% of its world wide peak. We estimate the arrival time of variant  $X$  in an country by the most simple way: the first time the variant is detected (according to GISAID data). In Fig. S8 the estimated outbreak time, official WHO and arrival times of each country are shown. Since for some variants (Alpha, Delta, BA.2) many arrival times fall clearly before our estimated and even the official outbreak date, we recomputed for these countries the arrival time to the first GISAID-detection after the outbreak date. We argue that either (i) the sequencing of the variant in these countries was error-prone (1. count is very sensitive to any wrong predetection) or (ii) the spreading was slow and the variant did not dominate the local epidemic until it reached a susceptible country (low NPIs) from where it did spread more easily (probably the case for Delta).

AUGUST 19, 2022

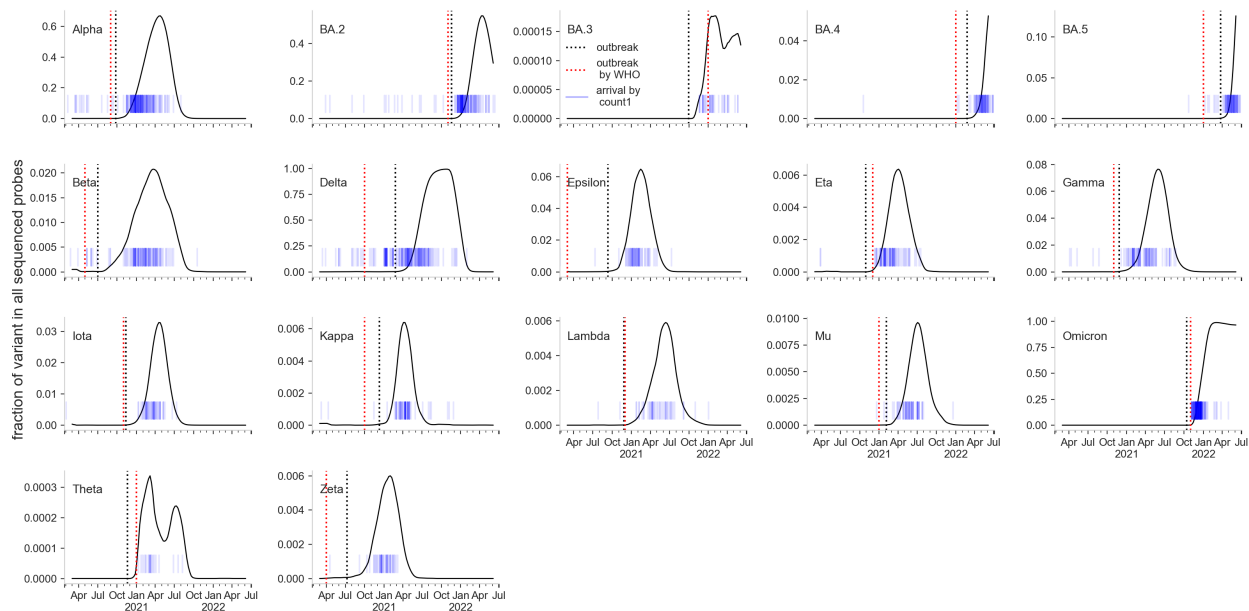


Figure S8: **Outbreak defined by fraction of all sequenced probes.** The outbreak date (black dashed vertical line) of a variant can be defined by the first time the fraction of a variant  $X$  of all sequenced probes reaches 2.5% of its current world wide peak. To exclude maldetections of 1st. arrival times in countries, we exclude all arrival times (blue short vertical lines) that are before the outbreak date and set the arrival time as the first detection in the respective country after the outbreak date. The official outbreak date by WHO is marked by a red dashed vertical line.

## II.6 Outbreak detection based on 1st count GISAID data

If we repeat the outbreak detection method using all variants and the arrival times estimated via GISAID data (arrival by first detection, Fig. S8), we see that the outbreak detection via the best correlation between import risk distance  $D^{IR}$  and arrival times  $T_{arrival}$  in general confirms the outbreak regions declared by the WHO (see Figs. S10, S9). There is a discrepancy for Delta. While using WHO and "cov-lineages.org" data, the official outbreak country India (IN) was second best, it is only on rank 12 if our GISAID estimates are used. A possible explanation is, that our outbreak date estimation is 5 months after the WHO date. In order to not lose the countries with arrivals before the outbreak date, we recompute the arrivals by the first count after the estimated outbreak date. One can argue that Delta did locally spread much stronger in South Africa (ZA, the top ranked country), and therefore is ZA for the worldwide distribution of larger importance than India. An alternative explanation is that the passenger flow in the WAN was too low and when it increased ZA had a more active Delta epidemic.

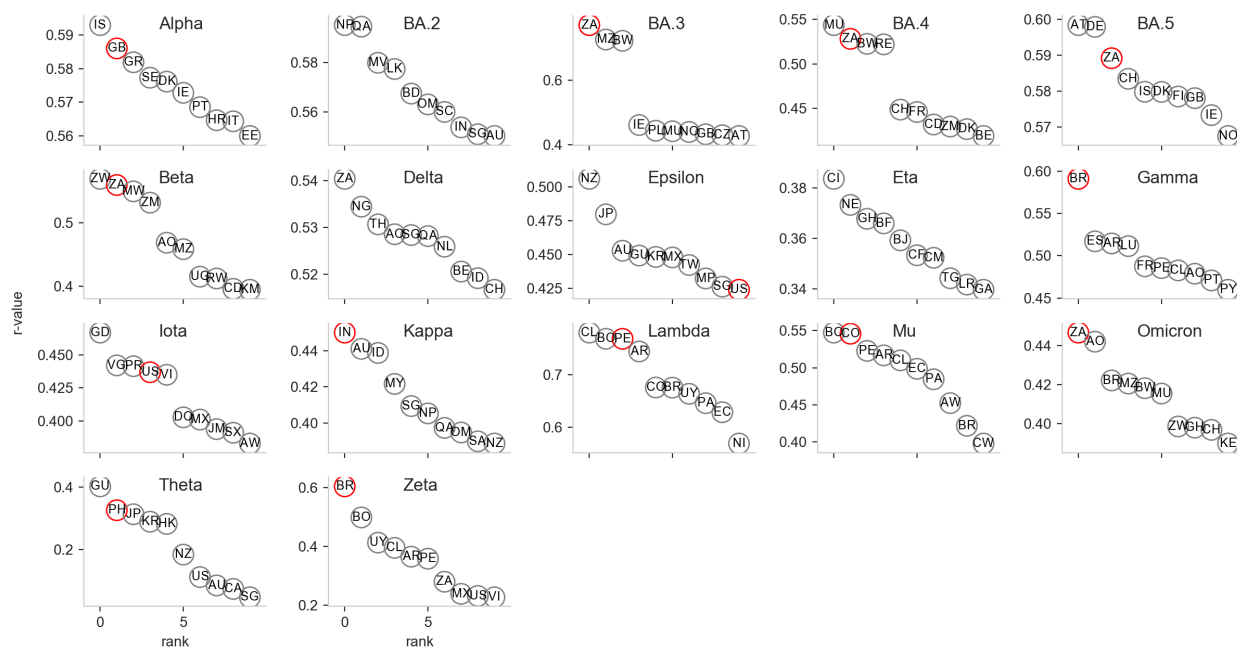


Figure S9: **Arrival prediction (r-value) for the 10 best outbreak candidate.** The r-value between the import risk distance  $d_{\infty}(m|n_0) = -\log(p_{\infty}(m|n_0))$  and the arrival time for the 10 best ranked outbreak countries ( $n_0$ ). The 2 Letters in the circles are the countries ISO alpha-2 codes. The red circle marks the country estimated as outbreak country based on GISAID arrival times. In contrast to Fig. S6: the arrival times and outbreak dates are estimated via GISAID data (arrival by first count, outbreak date by reaching the first time 2.5% of world wide peak of the respective variant).

AUGUST 19, 2022

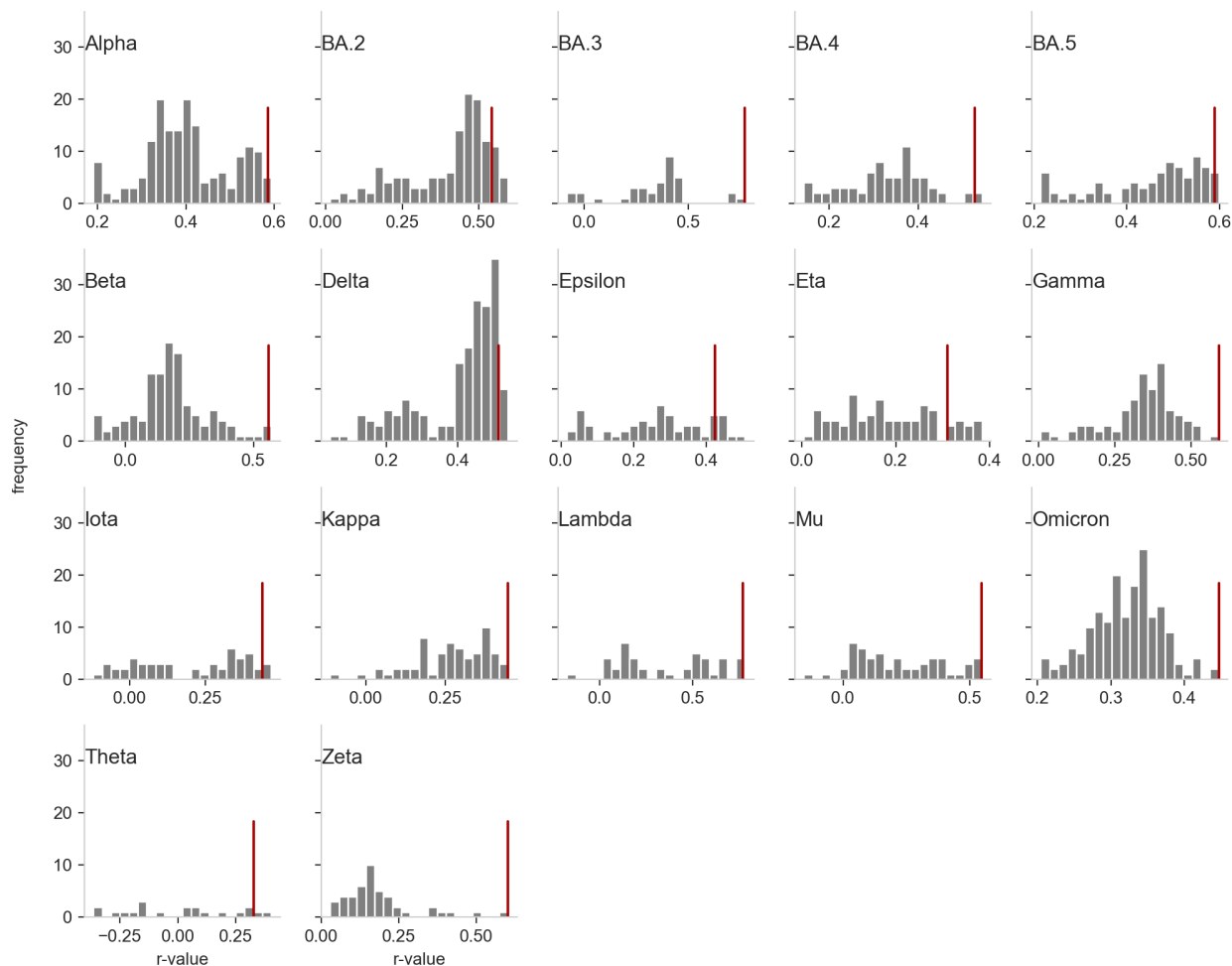


Figure S10: **Arrival prediction performance (r-value) for the outbreak country candidates.** The frequency of the r-value between the import risk distance  $D^{IR}(m|n_0) = -\log(p_\infty(m|n_0))$  and the arrival time for all possible outbreak countries. The red vertical line marks the r-value using the country estimated as outbreak country based on GISAID arrival times. In contrast to Fig. S7: the arrival times and outbreak dates are estimated via GISAID data (arrival by first count, outbreak date by reaching the first time 2.5% of world wide peak of the respective variant).



### III Epidemic Scenarios

We consider two distinct models to project the number of daily new infected people, namely, a renewal equation based model and a multi-strain SIR-like model. The first one is actually part of the pipeline, while the second one is used as validation.

#### III.1 Renewal equation

The renewal equation approach is a well-known technique, widely used in epidemiology [16, 17, 18]. The reason why renewal equations are such strong candidates for early projection of new cases, is the fact that informing them requires only the reproduction number of the new variant of concern, its generation interval distribution, and the number of people infected by the new variant who travel into the target country from the source country. This allows easily to explore scenarios with different values of epidemiological quantities of interest, such as the effective reproduction number of a new variant as it spreads from the source country to others through travelers.

For now, we assume that the susceptible population is much larger than the number of active cases, and that the mixing between the infected and the susceptible is homogeneous. This allows to exclude feedback loops in the dynamics, e.g. the fact that immunity to the new variant builds up through infection, which would modify the dynamics itself. Such strong assumptions are acceptable as long as we restrict our projections to the very first few weeks from the introduction of the new variant in the target country.

The model assumes that the number of newly infected people at day  $t$ ,  $I(t)$ , is given by two distinct processes: a) the arrival of infected individuals from the source country ( $I_{out}(t)$ ) and b) the daily new infections ( $I_{in}(t)$ ) happening in the target country due to the endogenous spreading. The former is estimated from section II, while the latter can be estimated through the renewal equation

$$I_{in}(t) = \sum_{s=t_0}^t \Gamma_s \mathcal{R}_s I(s), \quad (\text{S23})$$

where  $t_0$  is the day the first infected cases arrived in the target country,  $\mathcal{R}_s$  is the daily reproductive number on day  $s$ , and  $\Gamma_s$  is the generation time distribution, i.e. the fraction of transmissions that would occur on day  $s$  after infection. Finally  $I(t) = I_{out}(t) + I_{in}(t)$ . This is the simplest renewal process, which does not include the fact that the target population might have an inhomogeneous immunological landscape, due to previous infections or vaccination. To model this phenomenon, we reinterpret the term on the right side of equation (S23) as the number of inoculations spreading from currently infecting people, which will turn into infections depending on the susceptibility of the recipients. If we assume that previous infections (with other variants) protect against reinfection with an efficacy of  $n_e$ , and, analogously, vaccination has an effectiveness of  $\nu_e$ , then we can explicitly account for removals by modifying equation (S23) into

$$I_{in}(t) = \sum_{s=t_0}^t \Gamma_s \mathcal{R}_s I(s) \left(1 - n_e \frac{R^{(old)}(t)}{N}\right) \left(1 - \nu_e \frac{V(t)}{N}\right), \quad (\text{S24})$$

where  $R^{(old)}(t)$  is the number of recovered people from previous variants that still have some protection against infections, and  $V(t)$  is the total number of vaccinated people. This assumes that the number of recovered or vaccinated people is uniformly distributed across the population, and that the events 'being vaccinated' and 'having been infected' are independent. This also assumes no gradual waning of protection against infection. However, we can consider as recovered or vaccinated only people who were infected or vaccinated recently, rather than from the beginning of the pandemic. For instance, considering only people who got either infected or their second dose up to six months prior to  $t$  is equivalent to assuming that there is an abrupt waning of efficacy against protection six months after getting infected or vaccinated.

Although these hypothesis might seem unrealistic, the lack of readily available data about waning and immunological landscapes of various countries, and the fact that this should be used only for short-term scenario explorations, allow us to avoid introducing further complexity into the model.

The cumulative number of cases and amount of fully vaccinated individuals at each day are the ones reported in the public repository at <sup>1</sup>. We select the values for vaccine efficacy and protection from previous infection from available works. In particular we set the vaccine efficacy  $\nu_e$  to 0 for Alpha, 0.5 for Delta, BA1 and BA2 and to 0.12 for BA.5

<sup>1</sup><https://ourworldindata.org/>

( [19, 20, 21, 22]). The selected protection against reinfection  $n_e$  is 1 for Delta, 0.56 for BA.1 and BA.2 Omicron lineages and 0.13 for BA.5 ( [23, 24, 22]).

The second model is a multi-strain SIR inspired by [25]. This is a two-strain model in which people who recover after being infected with the former variant are not completely immune to infection from the latter variant. The equations governing this system are

$$\begin{cases} \frac{dS}{dt} &= -(\lambda_0(t) + \lambda_1(t))S(t) \\ \frac{dI^{(0)}}{dt} &= \lambda_0 S(t) - \gamma I^{(0)}(t) \\ \frac{dI^{(1)}}{dt} &= \lambda_1 S(t) + (1 - n_e \alpha) \lambda_1 R^{(1)}(t) - \gamma I^{(1)}(t) \\ \frac{dR^{(0)}}{dt} &= \gamma I^{(0)}(t) - (1 - n_e \alpha) \lambda_1 R^{(0)}(t) \\ \frac{dR^{(1)}}{dt} &= \gamma I^{(1)}(t) \end{cases} \quad (\text{S25})$$

where  $\lambda_i(t) = \beta_i \frac{I^{(i)}(t)}{N}$ ,  $\beta_i$  being the transmission rate of the variant  $i$ , and  $\gamma$  being the recovery rate. The initial condition  $S(t_0), I_0(t_0), I_1(t_0), R_0(t_0), R_1(t) = \{S_0, I_0^{(0)}, I_{out}(t_0) + I_0^{(1)}, R_0^{(0)}, R_0^{(1)}\}$ . Note that, since  $I^{out}(t)$  represents the arrivals from the source country at the beginning of each day, the system is not closed. This is not a problem because we are considering countries, so  $\frac{I^{out}(t)}{N} \ll 1$ . Since the dynamics does not include, per se, the fact that the initial condition changes every day due to arrivals, we can solve this system on a daily basis, updating the initial condition and restarting the system accordingly. The advantage of this system is that it includes feedback phenomena, which is good when validating the model, as it may need to run for more than a few weeks. The drawbacks are that informing the model requires good point estimates of the various compartments, and the interpretation of the transmissibility coefficient related to the measured  $\mathcal{R}_t$ , which may not be straight-forward. For such reasons, this model is used to validate the renewal equation approach, in particular for countries where no new cases were observed after a few weeks from their emergence (not shown). Projections errors valuated with the SIR model relative to Alpha lineage are shown in

### III.2 A fully worked out example: the Alpha variant

We apply our pipeline to a real case, the Alpha variant of concern (VOC), that was identified in the UK on 20 September 2020 [8]. We assume that the UK is the source country and we demonstrate how the pipeline works. In the following, we consider as the generation time interval distribution the one inferred from the literature [26].

Starting from the phylogenetic part of our pipeline, we take the time of emergence estimated when  $n = 20$  sequences were collected, to simulate a realistic scenario where only few information is available. This gives a central estimate for the time of emergence of the Alpha variant around the 9<sup>th</sup> of November 2020. The daily growth rate estimated is  $r = 0.097$  (95% HPD: 0.008–0.202). To translate this into  $R_t$  in the source country, we assume that all the growth rate advantage of Alpha relative to the previous circulating variants is given only by transmission advantage (limited capacity of reinfections with Alpha). Further, typical generation time distributions are Gammas, as in [26]. This allows us to estimate the  $R_t$  using formula 2.2 in [16]:

$$R = \frac{(r + b)^a}{b^a}, \quad (\text{S26})$$

where  $b$  and  $a$  are the shape and rate of the Gamma distribution generation time. In our case,  $a = 5.9, b = 1.13$ , therefore  $R_{t(\alpha)} = 1.62(1.04, 2.63)$ .

For any target country, the projection of the number of cases infected with Alpha in the next weeks is performed in two steps: first, we estimate the number of infected travellers (referred to as seeds) who arrive in the target country from the source country, then we use the renewal equation (S24) on each possible scenario, to account for endogenous transmission of the secondary cases in the target country. The first step consists in using the import risk estimates described in section II to compute the number of daily travellers from source country to other target countries. We use import risk probability from source to target times the average daily outflow of passengers from source country using WAN data. We then determined the number of travellers infected with Alpha. This is done by considering the proportion of sequenced cases that are Alpha times the 7 – day moving average of daily incidence of new cases, assuming that sequences are taken randomly from the infected population. This estimate does not include undercounting in the source country, which we can estimate as follows.

For a given country, we use the daily new estimated COVID-19 infections from the IHME model which is a hybrid with two main components: a statistical “death model” component produces death estimates that are used to fit an

SEIR model component<sup>2</sup>. For a complete overview of this model and a comparison with other estimates we refer to OWID<sup>3</sup>. The data we used for our estimation are publicly available<sup>4</sup>. In a given temporal window, we integrate over time the confirmed number of cases (7d moving average) and the estimated true number of cases, as well as the estimates for its lower and upper bounds defining the 95% uncertainty interval. The mean undercounting factor is estimated by the ratio between the integrated estimate of the true cases and the confirmed ones in the temporal window, and similarly we estimate the corresponding uncertainty interval. We show in Fig. SS11 the undercounting factor obtained for all countries for which the data is available, whereas Fig. SS12 shows the evolution of this factor along periods of 6 months for some representative countries.

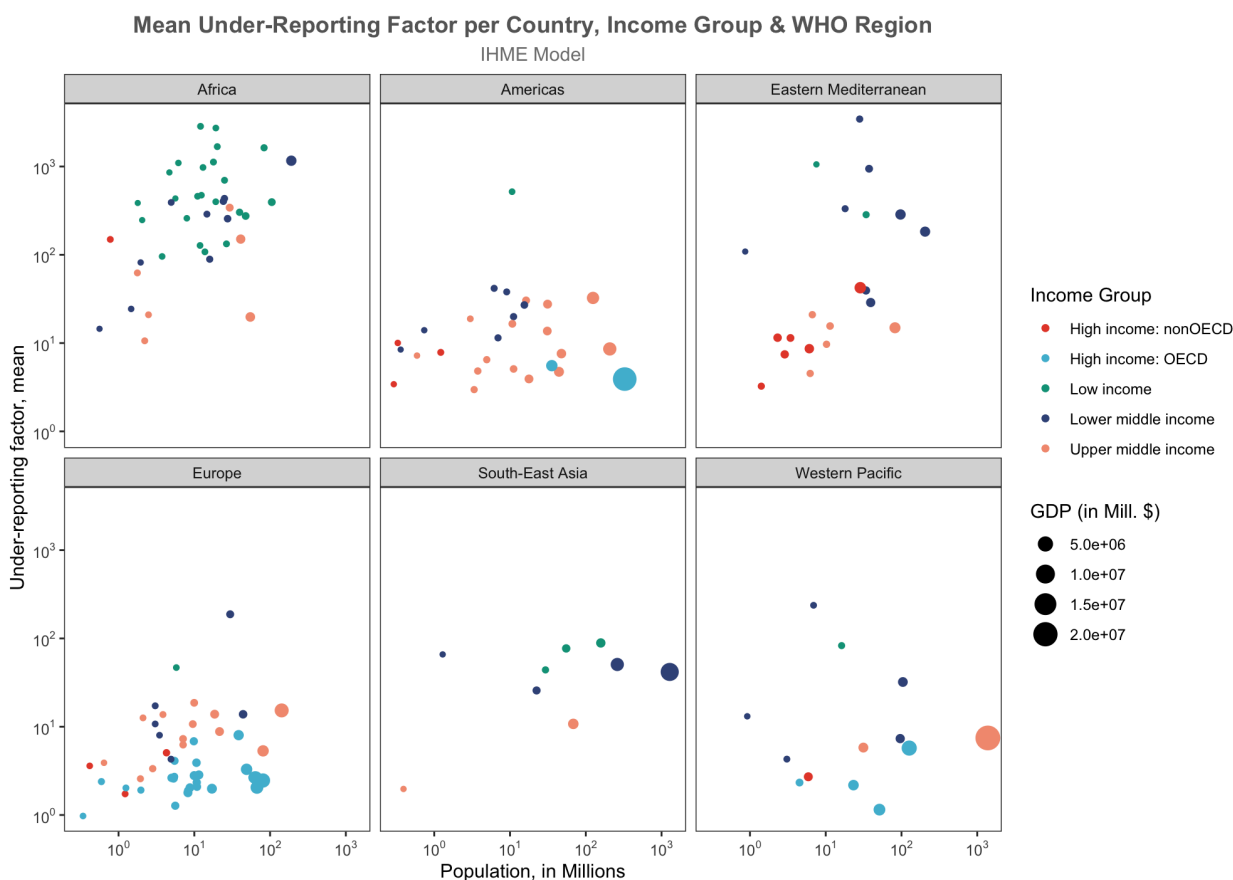


Figure S11: **Undercounting factors by WHO region and income group.** Estimates of the factor accounting for missing confirmed cases: values larger than 1 indicate that a country is counting and confirming less COVID-19 cases than the real number. The reference period is the first semester of 2022. See the text for further details.

To allow for variability in undercounting, we consider two extreme scenarios: the best one, where undercounting is assumed to be 2.27, and the worst one, where undercounting is assumed to be 2.97. The number of infected travelers from the source country to the target country is then computed by multiplying the number of travellers into the target country by the proportion of infected people in the source country. This is often not a natural number. This is not a problem, as the renewal equation does not need to use integer number of infected people, and we interpret this as the results of the various averaging performed through all the steps. The model produces the total number of infected people in the target country given the seeds and the  $\mathcal{R}_t$  by day of infection. To validate the model, we need to estimate how many people infected with the VOC were present in the target country during the considered period. We do so in the same way we estimate prevalence in the source country: by multiplying the proportion of sequenced cases that turned out to be Alpha times the daily incidence in the target country, scaled by the estimated undercounting factor.

<sup>2</sup><https://covid19.healthdata.org/>

<sup>3</sup><https://ourworldindata.org/covid-models>.

<sup>4</sup><https://ourworldindata.org/grapher/daily-new-estimated-covid-19-infections-ihme-model>.

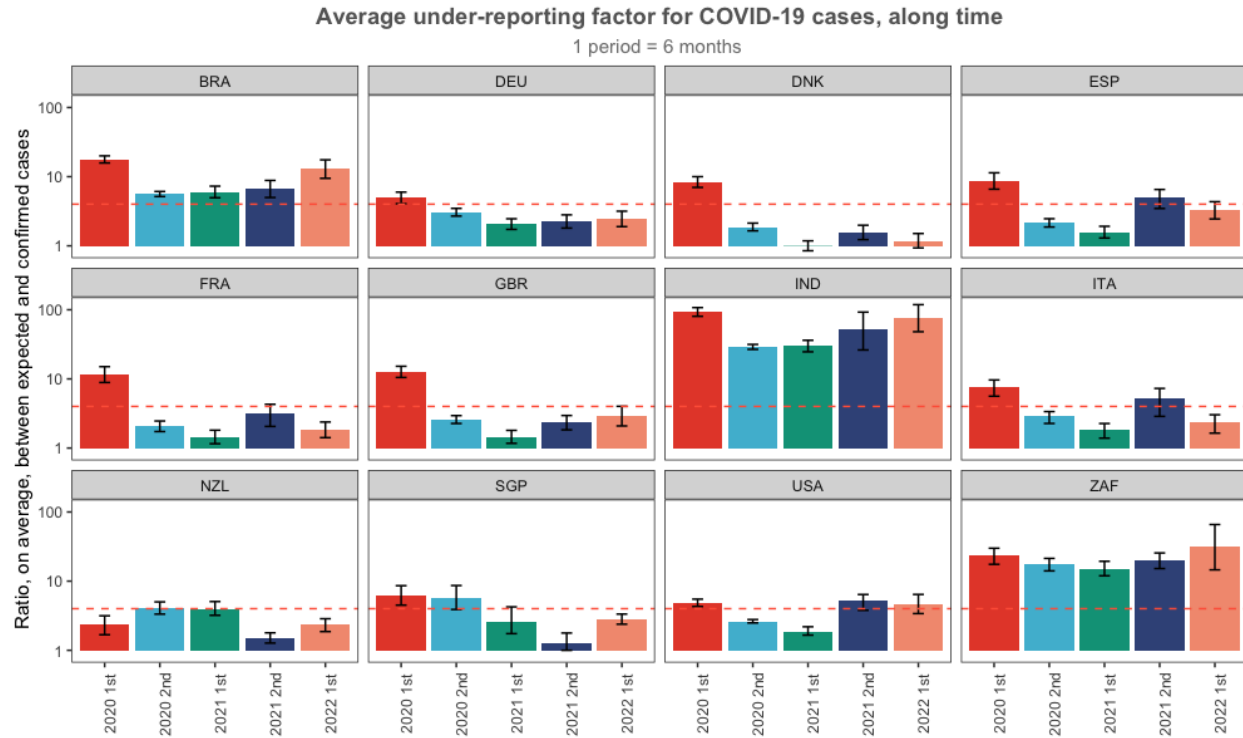


Figure S12: **Undercounting factors over time.** Estimates of the factor accounting for missing confirmed cases as in Fig. S11, where each panel describes the evolution along periods of 6 months for some representative countries. The dashed line indicates the value 4. See the text for further details.

The total number of different scenarios computed is, in this case  $2 \times 2 \times 3$ : undercounting in both the source and the target countries, and the different reproduction number of the VOC. Results are shown in Figure 3C and in Figure S13A.

### III.3 Prediction error

For each lineage we evaluate different scenarios with a) low and high values of under reporting in both source and target country b) three different basic reproduction numbers  $R_t$  that correspond to the range of growth rate values estimated from the phylogenetic reconstruction.

We infer from data the number of infected individuals with the emerging lineage in the target country  $m$  and we evaluate the prediction error as zero if this estimated number is included in the range identified by different epidemic scenarios. If the number of infected people evaluated from data is out of the range spanned by the epidemic curves, then the prediction error is evaluated as the root-mean-square error, normalized to the range of the data observed in the target country  $m$ , between observed and the closest simulated epidemic curve:

$$nRMSE(m) = \frac{1}{\max_t (I_m^{(data)}) - \min_t (I_m^{(data)})} \sqrt{\frac{1}{n_t} \sum_{t=1}^{n_t} [I_m^{(data)}(t) - I_m^{(model)}(t)]^2} \quad (S27)$$

where  $n_t$  is the number of weeks with number of sequences greater than zero for the selected lineage in the considered country  $m$ , that is  $n_t$  is the number of available data points with not null infected people. Since the scenario simulations stop at the 3 week after sequencing was reported in country  $m$ ,  $n_t$  is always  $n_t = 2$ . The idea behind the normalization by the data range is that it reflects the noise of reported sequences, i.e. if the sequencing rate is low, we expect a large variation and the sequencing data is less reliable. Prediction errors evaluated for all the considered lineages are shown in Figure 3 of the main document. All the panels report the nRMSE in each country as a function of both the number of daily passengers normalized to the total country population (x-axis, values for 100000 individuals) and the number of total collected daily sequences normalized to the total number of confirmed cases (y-axis, values for 100000 cases).

AUGUST 19, 2022

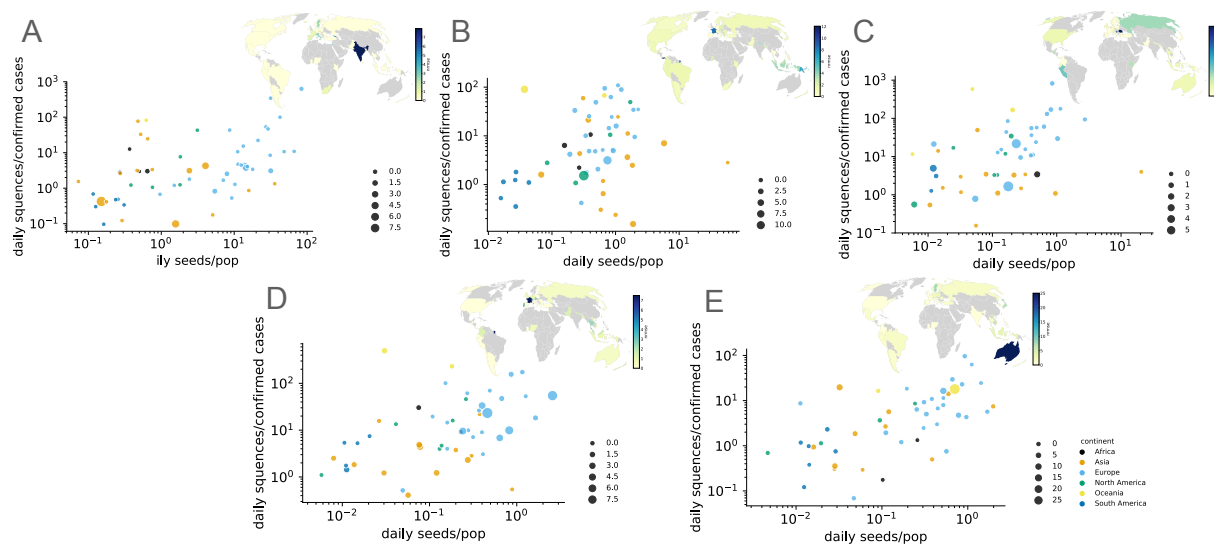


Figure S13: **Epidemic prediction errors.** Estimated errors between the number of individuals infected with an emerging lineage and the epidemic curves simulated in the considered scenarios. X-axis show the number of daily passengers normalized to the population in each country (for 100,000 individuals), y-axis report the number of collected daily sequences, without any classification per lineage, normalized to the total number of confirmed cases (for 100,000 cases). Inset panels show the map of prediction errors in each country. Panels A-E refer to, respectively, Alpha, Delta, BA.1, BA.2 and BA.5 lineages.

Insets show the evaluated error in each country. Results assess that, in most of the country, the simulated scenarios encompass the data and the prediction error is evaluated as zero. Moreover, error values greater than zero can be found for countries with higher passenger flows.

AUGUST 19, 2022

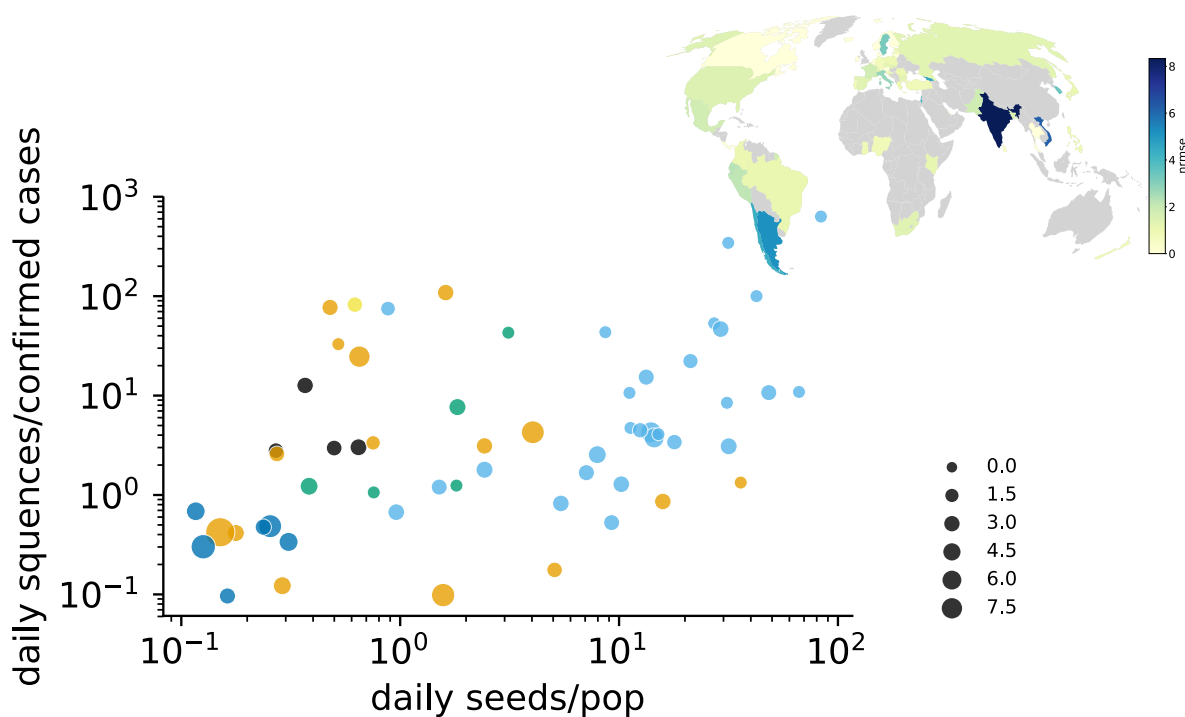


Figure S14: **Epidemic prediction errors with SIR model, Alpha lineage.** Estimated errors between the number of individuals infected with an emerging lineage and the epidemic curves simulated in the considered scenarios. X-axis show the number of daily passengers normalized to the population in each country (for 100,000 individuals), y-axis report the number of collected daily sequences, without any classification per lineage, normalized to the total number of confirmed cases (for 100,000 cases). Inset panels show the map of prediction errors in each country.

## IV Variant Dominance

As already described for the Alpha variant by Fort [27], the relative fraction of a new variant can be described by a simple and well known logistic growth equation. We use an equivalent formula with a slightly different convention, using the growth rate  $g$  of the variants as the fitness

$$\frac{dx}{dt} = x(g_v - \hat{g}) \quad (\text{S28})$$

This equation is solved by the logistic growth function

$$x = \frac{1}{1 + \frac{1-x_0}{x_0} e^{\hat{g}t(1-g_v)}}$$

Assuming the initial import is small, we can simplify this to

$$x \approx \frac{1}{1 + \frac{1}{x_0} e^{\hat{g}t(1-g_v)}}$$

Solving this equation for the inflection point of the sigmoid curve, the point where the new variant becomes the dominant variant in the system (where  $x(t = t_{inf}) \approx 0.5$ ), gives us:

$$t_{inf} = \frac{\log\left(\frac{x_0}{1-x_0}\right)}{-g_v \hat{g}} \approx -\frac{\log(x_0)}{g_v \hat{g}} \quad (\text{S29})$$

For our evaluation we used the GISAID variant data aggregated to weekly values. We show the validity of the inflection time approximation in Eq. S29 by estimating the ratio  $x_0$  by the ratio between imported cases by domestic cases  $x_0 \propto \text{cases}/\text{import}$  in Fig. S15.

This function can be directly fitted to available sequence data using logistic regression methods. This will generally perform well in the absence of sampling biases. However, in the presence of sampling bias often present in the first data points of a new outbreak, we found robust regression methods to be more reliable. Transformation using the logit function results in a linear relationship between time and the sample fraction of the new variant:

$$\text{logit}(x(t)) = \log\left(\frac{x_0}{1-x_0}\right) + g_v \hat{g}t \approx \log(x_0) + g_v \hat{g}t$$

Then we can perform a robust regression using the RANSAC algorithm. We used the regression method of the python sklearn library [28] and extracted the fitted values for  $t_{inf}$ .

AUGUST 19, 2022

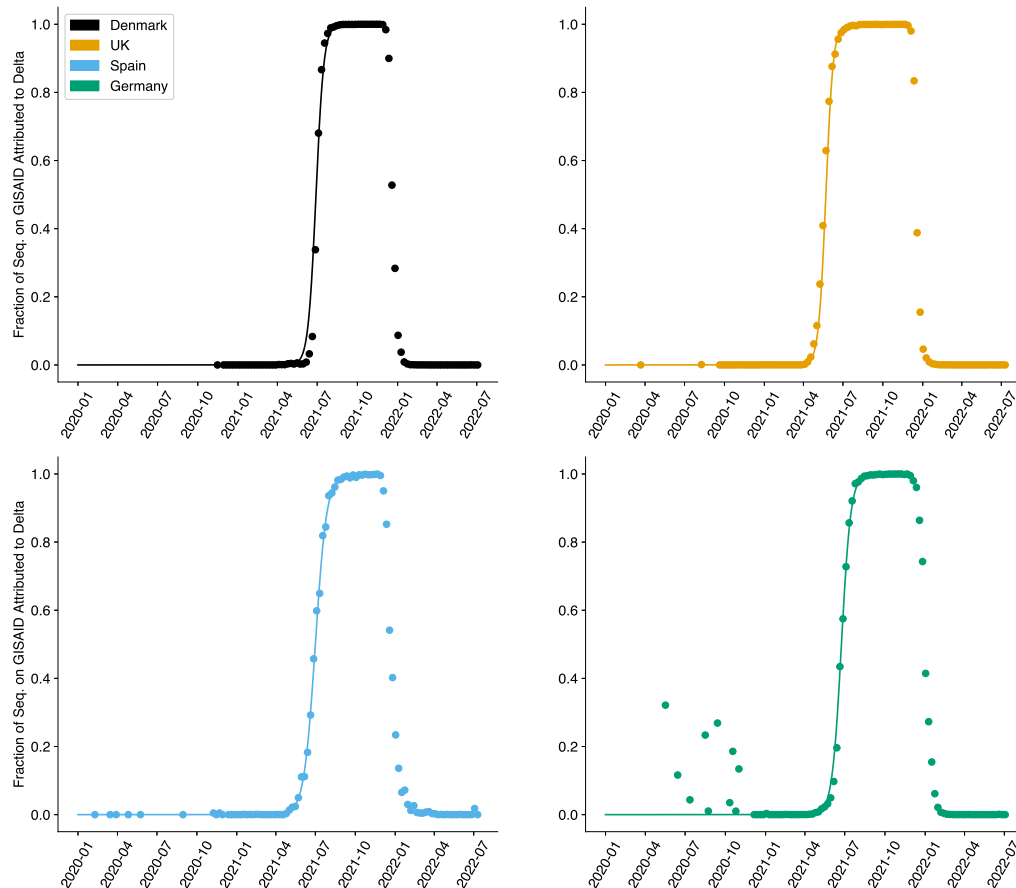


Figure S15: The fraction of seq. on GISAID attributed to the Delta variant for four example countries. As described for the Alpha variant by Fort [27], the relative fraction of a new variant can be accurately described by a simple logistic growth equation (Eq. S28).



AUGUST 19, 2022

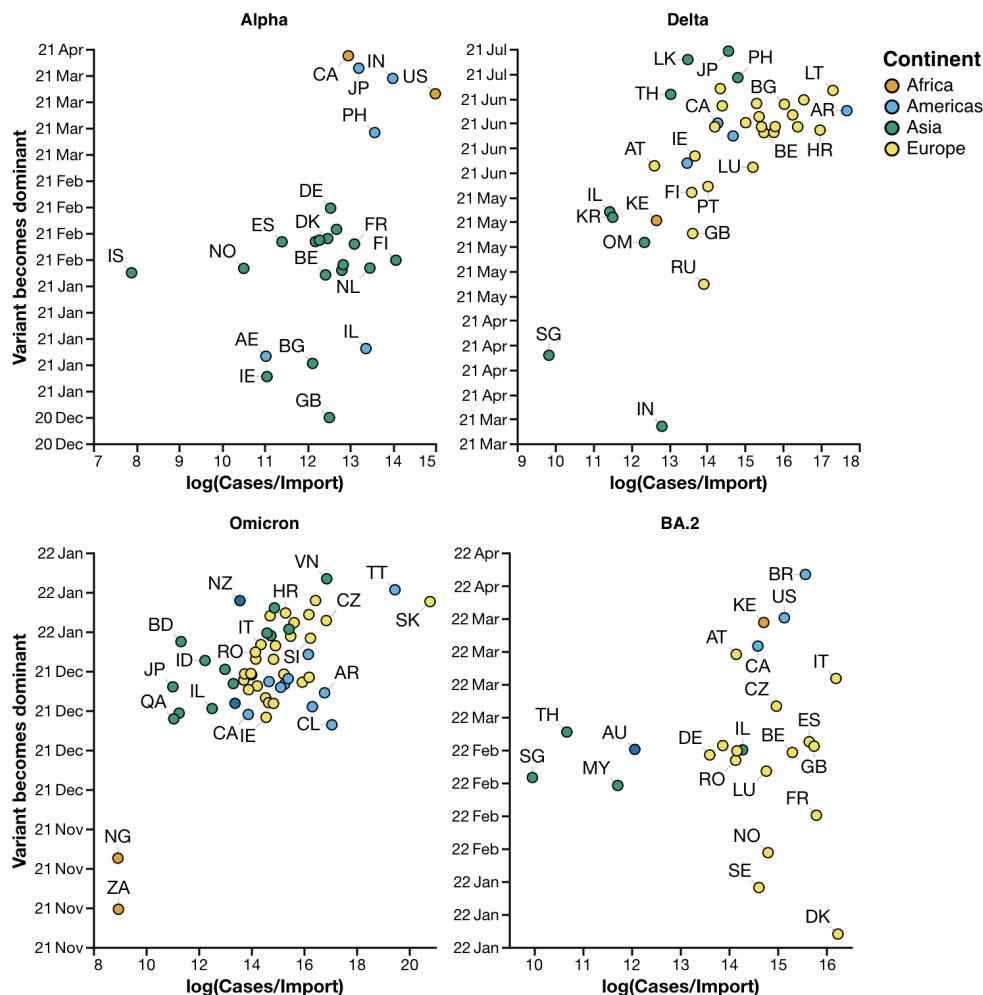


Figure S16: The date at which four example variants became dominant in different target countries, and an easy to obtain ratio of target country cases and target country's import risk. Assuming a constant growth advantage  $f$ , the logistic growth equation can be solved for the inflection time  $t_{inf}$  (new variant constitutes 50% of new cases), describing a linear relationship between the infection time and the logarithm of the fraction  $x_0$  of imported cases to domestic cases ( $t_{inf} \propto -\frac{\log(x_0)}{f-1}$ ), in the limit of small  $x_0$ . The fraction  $x_0$  can be approximated using the Import Risk and case numbers. For the delta variant, the relationship is linear and highly correlated (Pearson's  $r = 0.81$ ,  $p < 10^{-10}$ ), for the other variants the relationship is less clear. The time a variant becomes dominant is based on weekly aggregates of GISAID data, transformed with the logit function to obtain a linear relationship between time and the fraction of variant sequences and fitted using a robust regression (RANSAC algorithm, using the scikit-learn python package [28]). For the BA.2 variant, the absolute number of cases is multiplied by the fraction of Omicron cases, both for the fitted data and the cases in the ratio.

## V Information Distance

We also devise an alternative definition of distance on top of a network which embeds information from multiple-pathways diffusion as an additional comparison to the import risk measure. Distances based on the diffusive properties of the system have been of interest in the recent years [10, 14]. Another key example is the Diffusion Distance [13] which estimates a metric distance between nodes based on how similarly the random walkers explore the network by using those nodes as sources, under the assumption that a mesoscale structure is recovered during the time scales in which the random walker explores its functional community.

Starting from Diffusion Distance definition, we propose an educated rewrite of the measure that fits the problem under study to predict arrival times of a random walker on the network, such as an infectious traveller from a source country. The probability  $\mathbf{p}(t | i)$  of a walker to be in any point in the network at time  $t$ , starting from node  $i$ , embeds information of multiple paths via successive applications of the Laplacian operator. We introduce a new measure that merges this concept from Diffusion Distance and also embeds information from Effective Distance [10], namely, the idea that low probabilities  $p_k(t | i)$  are associated with large distances. This can be embedded by taking the negative of the logarithm of the probability, in analogy with Shannon's entropy. We now introduce this candidate measure for diffusive dynamics which we define Information Distance:

$$D_{(s \rightarrow k)}^{ID}(t) = -\log_{10} p_k(t | s) \quad (\text{S30})$$

in which  $p_k(t | s)$  represents the  $k$ -th entry associated with node  $k$  of the probability state  $\mathbf{p}(t | s) = \mathbf{v}_s \cdot e^{-tL^{RW}}$ . Here  $\mathbf{v}_s$  is the initial condition probability for the walker starting from node  $s$ , the canonical vector with  $s$ -th component equal to 1. The random walk normalized Laplacian ( $L^{RW}$ ) [29] term encodes the probability to move from node  $i$  to node  $j$  in its matrix elements. Its off-diagonal terms can be computed as the negative value of  $P_{ij}$ , which is directly estimated from the WAN weighted links as stated in subsection II.2. Given the multiple timescales involved in this definition, we evaluate the metric at different scales  $t$  to find the timescale at which  $D^{ID}(t)$  performs better.

## References

- [1] K. Katoh and D. M. Standley, “MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability,” *Molecular Biology and Evolution*, vol. 30, pp. 772–780, 01 2013.
- [2] M. A. Suchard, P. Lemey, G. Baele, D. L. Ayres, A. J. Drummond, and A. Rambaut, “Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10,” *Virus Evolution*, vol. 4, 06 2018. vey016.
- [3] D. L. Ayres, M. P. Cummings, G. Baele, A. E. Darling, P. O. Lewis, D. L. Swofford, J. P. Huelsenbeck, P. Lemey, A. Rambaut, and M. A. Suchard, “BEAGLE 3: Improved Performance, Scaling, and Usability for a High-Performance Computing Library for Statistical Phylogenetics,” *Systematic Biology*, vol. 68, pp. 1052–1061, 04 2019.
- [4] A. Rambaut, A. J. Drummond, D. Xie, G. Baele, and M. A. Suchard, “Posterior Summarization in Bayesian Phylogenetics Using Tracer 1.7,” *Systematic Biology*, vol. 67, pp. 901–904, 04 2018.
- [5] J. Kennedy and R. Eberhart, “Particle swarm optimization,” in *Proceedings of ICNN’95-international conference on neural networks*, vol. 4, pp. 1942–1948, IEEE, 1995.
- [6] H. Ritchie, E. Mathieu, L. Rodés-Guirao, C. Appel, C. Giattino, E. Ortiz-Ospina, J. Hasell, B. Macdonald, D. Beltekian, and M. Roser, “Coronavirus pandemic (covid-19),” *Our World in Data*, 2020. <https://ourworldindata.org/coronavirus>.
- [7] Cov-Lineages.org, “Global Lineage Reports,” 2022.
- [8] Á. O’Toole, V. Hill, O. G. Pybus, A. Watts, I. I. Bogoch, K. Khan, J. P. Messina, H. Tegally, R. R. Lessells, J. Giandhari, S. Pillay, K. A. Tumedi, G. Nyepetsi, M. Keabonye, M. Matsheka, M. Mine, S. Tokajian, H. Hassan, T. Salloom, G. Merhi, J. Koweyes, J. L. Geoghegan, J. de Ligt, X. Ren, M. Storey, N. E. Freed, C. Pattabiraman, P. Prasad, A. S. Desai, R. Vasanthapuram, T. F. Schulz, L. Steinbrück, T. Stadler, A. Parisi, A. Bianco, D. García de Viedma, S. Buenestado-Serrano, V. Borges, J. Isidro, S. Duarte, J. P. Gomes, N. S. Zuckerman, M. Mandelboim, O. Mor, T. Seemann, A. Arnott, J. Draper, M. Gall, W. Rawlinson, I. Deveson, S. Schlegelbusch, J. McMahon, L. Leong, C. K. Lim, M. Chironna, D. Loconsole, A. Bal, L. Josset, E. Holmes, K. St. George, E. Lasek-Nesselquist, R. S. Sikkema, B. Oude Munnink, M. Koopmans, M. Brytting, V. Sudha rani, S. Pavani, T. Smura, A. Heim, S. Kurkela, M. Umair, M. Salman, B. Bartolini, M. Rueca, C. Drosten, T. Wolff, O. Silander, D. Eggink, C. Reusken, H. Vennema, A. Park, C. Carrington, N. Sahadeo, M. Carr, G. Gonzalez, T. de Oliveira, N. Faria, A. Rambaut, and M. U. G. Kraemer, “Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 with grinch,” *Wellcome Open Research*, vol. 6, p. 121, sep 2021.
- [9] P. P. Klamser, A. Zachariae, B. F. Maier, O. Baranov, C. Jongen, F. Schlosser, and D. Brockmann, “Inferring country specific import risk of diseases from the World-Aviation-Network,” *In preparation*, 2022.
- [10] D. Brockmann and D. Helbing, “The Hidden Geometry of Complex, Network-Driven Contagion Phenomena,” *Science*, vol. 342, pp. 1337–1342, dec 2013.
- [11] D. Brockmann and D. Helbing, “The hidden geometry of complex, network-driven contagion phenomena,” *science*, vol. 342, no. 6164, pp. 1337–1342, 2013.
- [12] A. Gautreau, A. Barrat, and M. Barthélemy, “Global disease spread: Statistics and estimation of arrival times,” *Journal of Theoretical Biology*, vol. 251, pp. 509–522, apr 2008.
- [13] M. De Domenico, “Diffusion Geometry Unravels the Emergence of Functional Clusters in Collective Phenomena,” *Physical Review Letters*, vol. 118, p. 168301, apr 2017.
- [14] F. Iannelli, A. Koher, D. Brockmann, P. Hövel, and I. M. Sokolov, “Effective distances for epidemics spreading on complex networks,” *Physical Review E*, vol. 95, p. 012313, jan 2017.
- [15] WHO, “WHO - Tracking SARS-CoV-2 variants,” 2022.
- [16] A. Cori, N. M. Ferguson, C. Fraser, and S. Cauchemez, “A New Framework and Software to Estimate Time-Varying Reproduction Numbers During Epidemics,” *American Journal of Epidemiology*, vol. 178, pp. 1505–1512, Nov. 2013.
- [17] L. Ferretti, C. Wymant, M. Kendall, L. Zhao, A. Nurtay, L. Abeler-Dörner, M. Parker, D. Bonsall, and C. Fraser, “Quantifying SARS-CoV-2 transmission suggests epidemic control with digital contact tracing,” *Science (80-. )*, vol. 368, p. eabb6936, may 2020.
- [18] I. Dorigatti and E. Al., “Report 4: Severity of 2019-Novel Coronavirus (nCoV) (10 February 2020); [www.imperial.ac.uk/media/imperial-college/medicine/sph/ide/gida-fellowships/Imperial-College-COVID19-severity-10-02-2020.pdf](http://www.imperial.ac.uk/media/imperial-college/medicine/sph/ide/gida-fellowships/Imperial-College-COVID19-severity-10-02-2020.pdf),” no. February, pp. 1–14, 2020.

- [19] M. Fabiani, M. Puopolo, C. Morciano, M. Spuri, S. Spila Alegiani, A. Filia, F. D'Ancona, M. Del Manso, F. Riccardo, M. Tallon, V. Proietti, C. Sacco, M. Massari, R. Da Cas, A. Mateo-Urdiales, A. Siddu, S. Battilomo, A. Bella, A. T. Palamara, P. Popoli, S. Brusaferrero, G. Rezza, F. Menniti Ippolito, and P. Pezzotti, "Effectiveness of mRNA vaccines and waning of protection against SARS-CoV-2 infection and severe covid-19 during predominant circulation of the delta variant in Italy: retrospective cohort study," *BMJ*, vol. 376, 2022.
- [20] F. C. M. Kirsebom, N. Andrews, J. Stowe, S. Toffa, R. Sachdeva, E. Gallagher, N. Groves, A.-M. O'Connell, M. Chand, M. Ramsay, and J. L. Bernal, "COVID-19 vaccine effectiveness against the omicron (BA.2) variant in England," *Lancet Infect. Dis.*, vol. 22, pp. 931–933, jul 2022.
- [21] N. Andrews, J. Stowe, F. Kirsebom, S. Toffa, T. Rickeard, E. Gallagher, C. Gower, M. Kall, N. Groves, A.-M. O'Connell, D. Simons, P. B. Blomquist, A. Zaidi, S. Nash, N. Iwani Binti Abdul Aziz, S. Thelwall, G. Dabrera, R. Myers, G. Amirthalangam, S. Gharbia, J. C. Barrett, R. Elson, S. N. Ladhani, N. Ferguson, M. Zambon, C. N. J. Campbell, K. Brown, S. Hopkins, M. Chand, M. Ramsay, and J. Lopez Bernal, "Covid-19 Vaccine Effectiveness against the Omicron (B.1.1.529) Variant," *N. Engl. J. Med.*, vol. 386, no. 16, pp. 1532–1546, 2022.
- [22] Q. Wang, Y. Guo, S. Iketani, M. S. Nair, Z. Li, H. Mohri, M. Wang, J. Yu, A. D. Bowen, J. Y. Chang, J. G. Shah, N. Nguyen, Z. Chen, K. Meyers, M. T. Yin, M. E. Sobieszczyk, Z. Sheng, Y. Huang, L. Liu, and D. D. Ho, "Antibody evasion by SARS-CoV-2 Omicron subvariants BA.2.12.1, BA.4, & BA.5," *Nature*, 2022.
- [23] V. J. Hall, S. Foulkes, A. Charlett, A. Atti, E. J. M. Monk, R. Simmons, E. Wellington, M. J. Cole, A. Saei, B. Oguti, K. Munro, S. Wallace, P. D. Kirwan, M. Shrotri, A. Vusirikala, S. Rokadiya, M. Kall, M. Zambon, M. Ramsay, T. Brooks, C. S. Brown, M. A. Chand, S. Hopkins, N. Andrews, A. Atti, H. Aziz, T. Brooks, C. S. Brown, D. Camero, C. Carr, M. A. Chand, A. Charlett, H. Crawford, M. Cole, J. Conneely, S. D'Arcangelo, J. Ellis, S. Evans, S. Foulkes, N. Gillson, R. Gopal, L. Hall, V. J. Hall, P. Harrington, S. Hopkins, J. Hewson, K. Hoschler, D. Ironmonger, J. Islam, M. Kall, I. Karagiannis, O. Kay, J. Khawam, E. King, P. Kirwan, R. Kyffin, A. Lackenby, M. Lattimore, E. Linley, J. Lopez-Bernal, L. Mabey, R. McGregor, S. Miah, E. J. M. Monk, K. Munro, Z. Naheed, A. Nissr, A. M. O'Connell, B. Oguti, H. Okafor, S. Organ, J. Osbourne, A. Otter, M. Patel, S. Platt, D. Pople, K. Potts, M. Ramsay, J. Robotham, S. Rokadiya, C. Rowe, A. Saei, G. Sebbage, A. Semper, M. Shrotri, R. Simmons, A. Soriano, P. Staves, S. Taylor, A. Taylor, A. Tengbe, S. Tonge, A. Vusirikala, S. Wallace, E. Wellington, M. Zambon, D. Corrigan, M. Sartaj, L. Cromey, S. Campbell, K. Braithwaite, L. Price, L. Haahr, S. Stewart, E. D. Lacey, L. Partridge, G. Stevens, Y. Ellis, H. Hodgson, C. Norman, B. Larru, S. Mcwilliam, S. Winchester, P. Ciecwiwa, A. Pai, C. Loughrey, A. Watt, F. Adair, A. Hawkins, A. Grant, R. Temple-Purcell, J. Howard, N. Slawson, C. Subudhi, S. Davies, A. Bexley, R. Penn, N. Wong, G. Boyd, A. Rajgopal, A. Arenas-Pinto, R. Matthews, A. Whileman, R. Laugharne, J. Ledger, T. Barnes, C. Jones, D. Botes, N. Chitalia, S. Akhtar, G. Harrison, S. Horne, N. Walker, K. Agwuh, V. Maxwell, J. Graves, S. Williams, A. O'Kelly, P. Ridley, A. Cowley, H. Johnstone, P. Swift, J. Democratis, M. Meda, C. Callens, S. Beazer, S. Hams, V. Irvine, B. Chandrasekaran, C. Forsyth, J. Radmore, C. Thomas, K. Brown, S. Roberts, P. Burns, K. Gajee, T. M. Byrne, F. Sanderson, S. Knight, E. Macnaughton, B. J. L. Burton, H. Smith, R. Chaudhuri, K. Hollinshead, R. J. Shorten, A. Swan, R. J. Shorten, C. Favager, J. Murira, S. Baillon, S. Hamer, K. Gantert, J. Russell, D. Brennan, A. Dave, A. Chawla, F. Westell, D. Adeboyeke, P. Papineni, C. Pegg, M. Williams, S. Ahmad, S. Ingram, C. Gabriel, K. Pagget, P. Ciecwiwa, G. Maloney, J. Ashcroft, I. Del Rosario, R. Crosby-Nwaobi, C. Reeks, S. Fowler, L. Prentice, M. Spears, G. McKerron, K. McLelland-Brooks, J. Anderson, S. Donaldson, K. Templeton, L. Coke, N. Elumogo, J. Elliott, D. Padgett, M. Mirfenderesky, A. Cross, J. Price, S. Joyce, I. Sinanovic, M. Howard, T. Lewis, P. Cowling, D. Potoczna, S. Brand, L. Sheridan, B. Wadams, A. Lloyd, J. Mouland, J. Giles, G. Pottinger, H. Coles, M. Joseph, M. Lee, S. Orr, H. Chenoweth, C. Auckland, R. Lear, T. Mahungu, A. Rodger, K. Penny-Thomas, S. Pai, J. Zamikula, E. Smith, S. Stone, E. Boldock, D. Howcroft, C. Thompson, M. Aga, P. Domingos, S. Gormley, C. Kerrison, L. Marsh, S. Tazzyman, L. Allsop, S. Ambalkar, M. Beekes, S. Jose, J. Tomlinson, A. Jones, C. Price, J. Pepperell, M. Schultz, J. Day, A. Boulos, E. Defever, D. McCracken, K. Brown, K. Gray, A. Houston, T. Planche, R. Pritchard Jones, D. Wycherley, S. Bennett, J. Marrs, K. Nimako, B. Stewart, N. Kalakonda, S. Khanduri, A. Ashby, M. Holden, N. Mahabir, J. Harwood, B. Payne, K. Court, N. Staines, R. Longfellow, M. E. Green, L. E. Hughes, M. Halkes, P. Mercer, A. Roebuck, E. Wilson-Davies, L. Gallego, R. Lazarus, N. Aldridge, L. Berry, F. Game, T. Reynolds, C. Holmes, M. Wiselka, A. Higham, M. Booth, C. Duff, J. Alderton, H. Jory, E. Virgilio, T. Chin, M. Z. Qazzafi, A. M. Moody, R. Tilley, T. Donaghy, K. Shipman, R. Sierra, N. Jones, G. Mills, D. Harvey, Y. W. J. Huang, J. Birch, L. Robinson, S. Board, A. Broadley, C. Laven, N. Todd, D. W. Eyre, K. Jeffery, S. Dunachie, C. Duncan, P. Klennerman, L. Turtle, T. De Silva, H. Baxendale, and J. L. Heeney, "SARS-CoV-2 infection rates of antibody-positive compared with antibody-negative health-care workers in England: a large, multicentre, prospective cohort study (SIREN)," *Lancet*, vol. 397, pp. 1459–1469, apr 2021.
- [24] H. N. Altarawneh, H. Chemaitelly, M. R. Hasan, H. H. Ayoub, S. Qassim, S. AlMukdad, P. Coyle, H. M. Yassine, H. A. Al-Khatib, F. M. Benslimane, Z. Al-Kanaani, E. Al-Kuwari, A. Jeremijenko, A. H. Kaleeckal, A. N. Latif,

AUGUST 19, 2022

- R. M. Shaik, H. F. Abdul-Rahim, G. K. Nasrallah, M. G. Al-Kuwari, A. A. Butt, H. E. Al-Romaihi, M. H. Al-Thani, A. Al-Khal, R. Bertollini, P. Tang, and L. J. Abu-Raddad, "Protection against the Omicron Variant from Previous SARS-CoV-2 Infection," *N. Engl. J. Med.*, vol. 386, no. 13, pp. 1288–1290, 2022.
- [25] P. Stefanelli, F. Trentini, G. Guzzetta, V. Marziano, A. Mammone, M. Sane Schepisi, P. Poletti, C. Molina Grané, M. Manica, M. del Manso, X. Andrianou, M. Ajelli, G. Rezza, S. Brusaferrero, S. Merler, and C.-. N. M. S. S. Group, "Co-circulation of SARS-CoV-2 Alpha and Gamma variants in Italy, February and March 2021," *Eurosurveillance*, vol. 27, no. 5, 2022.
- [26] L. Ferretti, A. Ledda, C. Wymant, L. Zhao, V. Ledda, L. Abeler, M. Kendall, A. Nurtay, H.-Y. Cheng, T.-C. Ng, H.-H. Lin, R. Hinch, J. Masel, A. M. Kilpatrick, and C. Fraser, "The timing of COVID-19 transmission," p. 16.
- [27] H. Fort, "A very simple model to account for the rapid rise of the alpha variant of sars-cov-2 in several countries and the world," *Virus Research*, vol. 304, p. 198531, 2021.
- [28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [29] M. Newman, "Networks, 2nd edn oxford," *UK: Oxford University Press.[Google Scholar]*, 2018.