

Rapid epidemic expansion of the SARS-CoV-2 Omicron variant in southern Africa

Raquel Viana^{1*}, Sikhulile Moyo^{2,3,4*}, Daniel G Amoako^{5*}, Houriiyah Tegally^{6*}, Cathrine Scheepers^{5,7*}, Christian L Althaus⁸, Ugochukwu J Anyaneji⁶, Phillip A Bester^{9,10}, Maciej F Boni¹¹, Mohammed Chand¹², Wonderful T Choga³, Rachel Colquhoun¹³, Michaela Davids¹⁴, Koen Deforche¹⁵, Deelan Doolabh¹⁶, Susan Engelbrecht¹⁷, Josie Everatt⁵, Jennifer Giandhari⁶, Marta Giovanetti^{18,19}, Diana Hardie^{16,20}, Verity Hill¹³, Nei-Yuan Hsiao^{16,20,21}, Arash Iranzadeh²², Arshad Ismail⁵, Charity Joseph¹², Rageema Joseph¹⁶, Legodile Koopile², Sergei L Kosakovsky Pond²³, Moritz UG Kraemer²⁴, Lesego Kuate-Lere²⁵, Oluwakemi Laguda-Akingba^{26,27}, Onalethatha Lesetedi-Mafoko²⁸, Richard J Lessells⁶, Shahin Lockman^{2,29}, Alexander G Lucaci²³, Arisha Maharaj⁶, Boitshoko Mahlangu⁵, Tongai Maponga¹⁷, Kamela Mhlakwane^{17,30}, Zinhle Makatini³¹, Gert Marais^{16,20}, Dorcas Maruapula², Kereng Masupu⁴, Mogomotsi Matshaba^{4,32,33}, Simnikiwe Mayaphi³⁴, Nokuzola Mbhele¹⁶, Mpaphi B Mbulawa³⁵, Adriano Mendes¹⁴, Koleka Mlisana^{36,37}, Anele Mnguni⁵, Thabo Mohale⁵, Monika Moir³⁸, Kgomotso Moruisi²⁵, Mosepele Mosepele^{4,39}, Gerald Motsatsi⁵, Modisa S Motswaledi^{4,40}, Thongbotho Mphoyakgosi³⁵, Nokukhanya Msomi⁴¹, Peter N Mwangi^{10,42}, Yeshnee Naidoo⁶, Noxolo Ntuli⁵, Martin Nyaga^{10,42}, Lucier Olubayo^{21,22}, Sureshnee Pillay⁶, Botshelo Radibe², Yajna Ramphal⁶, Upasana Ramphal⁶, James E San⁶, Lesley Scott⁴³, Roger Shapiro^{2,29}, Lavanya Singh⁶, Pamela Smith-Lawrence²⁵, Wendy Stevens⁴³, Amy Strydom¹⁴, Kathleen Subramoney³¹, Naume Tebeila⁵, Derek Tshiabuila⁶, Joseph Tsui²⁴, Stephanie van Wyk³⁸, Steven Weaver²³, Constantinos K Wibmer⁵, Eduan Wilkinson³⁸, Nicole Wolter^{5,44}, Alexander E Zarebski²⁴, Boitumelo Zuze², Dominique Goedhals^{10,45}, Wolfgang Preiser^{17,30}, Florette Treurnicht³¹, Marietje Venter¹⁴, Carolyn Williamson^{16,20,21,46}, Oliver G Pybus²⁴, Jinal Bhiman^{5,7}, Allison Glass^{1,47}, Darren P Martin^{21,46}, Andrew Rambaut¹³, Simani Gaseitsiwe^{2,3**}, Anne von Gottberg^{5,44**}, Tulio de Oliveira^{6,38,48**} ✉

¹Lancet Laboratories, Johannesburg, South Africa

²Botswana Harvard AIDS Institute Partnership, Botswana Harvard HIV Reference Laboratory, Gaborone, Botswana

³Harvard T.H. Chan School of Public Health, Boston, Massachusetts

⁴Botswana Presidential COVID-19 Taskforce, Gaborone, Botswana

⁵National Institute for Communicable Diseases (NICD) of the National Health Laboratory Service (NHLS), Johannesburg, South Africa

⁶KwaZulu-Natal Research Innovation and Sequencing Platform (KRISP), Nelson R Mandela School of Medicine, University of KwaZulu-Natal, Durban, South Africa

⁷South African Medical Research Council Antibody Immunity Research Unit, School of Pathology, Faculty of Health Sciences, University of the Witwatersrand, Johannesburg, South Africa

⁸Institute of Social and Preventive Medicine, University of Bern, Bern, Switzerland

⁹Division of Virology, National Health Laboratory Service, Bloemfontein, South Africa

¹⁰Division of Virology, University of the Free State, Bloemfontein, South Africa

¹¹Center for Infectious Disease Dynamics, Department of Biology, Pennsylvania State University, University Park, PA, USA

¹²Diagnofirm Medical Laboratories, Gaborone, Botswana

¹³Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK

¹⁴Zoonotic Arbo and Respiratory Virus Program, Centre for Viral Zoonoses, Department of Medical Virology, University of Pretoria, Pretoria, South Africa

¹⁵Emweb bv, Herent, Belgium

¹⁶Division of Medical Virology, Faculty of Health Sciences, University of Cape Town, Cape Town, South Africa

¹⁷Division of Medical Virology, Faculty of Medicine and Health Sciences, Stellenbosch University, Tygerberg, Cape Town, South Africa

¹⁸Laboratorio de Flavivirus, Fundacao Oswaldo Cruz, Rio de Janeiro, Brazil

¹⁹Laboratório de Genética Celular e Molecular, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

²⁰Division of Virology, NHLS Groote Schuur Laboratory, Cape Town, South Africa

²¹Wellcome Centre for Infectious Diseases Research in Africa (CIDRI-Africa)

²²Division of Computational Biology, Faculty of Health Sciences, University of Cape Town

²³Institute for Genomics and Evolutionary Medicine, Department of Biology, Temple University, Pennsylvania, USA

²⁴Department of Zoology, University of Oxford, Oxford, UK

²⁵Health Services Management, Ministry of Health and Wellness, Gaborone, Botswana

²⁶NHLS Port Elizabeth Laboratory, Port Elizabeth, South Africa

²⁷Faculty of Health Sciences, Walter Sisulu University, Eastern Cape, South Africa

²⁸Public Health Department, Integrated Disease Surveillance and Response, Ministry of Health and Wellness, Gaborone, Botswana

²⁹Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public Health, Boston, Massachusetts, USA

³⁰NHLS Tygerberg Laboratory, Tygerberg, Cape Town, South Africa

³¹Department of Virology, Charlotte Maxeke Johannesburg Academic Hospital, Johannesburg, South Africa

³²Botswana-Baylor Children's Clinical Centre of Excellence

³³Baylor College of Medicine, Houston, Texas, USA

³⁴Department of Medical Virology, University of Pretoria, Pretoria, South Africa

³⁵National Health Laboratory, Health Services Management, Ministry of Health and Wellness, Gaborone, Botswana

³⁶National Health Laboratory Service (NHLS), Johannesburg, South Africa

³⁷Centre for the AIDS Programme of Research in South Africa (CAPRISA), Durban, South Africa

³⁸Centre for Epidemic Response and Innovation (CERI), School of Data Science and Computational Thinking, Stellenbosch University, Stellenbosch, South Africa

³⁹Department of Medicine, Faculty of Medicine, University of Botswana, Gaborone, Botswana

⁴⁰Department of Medical Laboratory Sciences, School of Allied Health Professions, Faculty of Health Sciences, University of Botswana, Gaborone, Botswana

⁴¹Discipline of Virology, School of Laboratory Medicine and Medical Sciences and National Health Laboratory Service (NHLS), University of KwaZulu–Natal, Durban, South Africa

⁴²Next Generation Sequencing Unit, Division of Virology, Faculty of Health Sciences, University of the Free State, Bloemfontein, South Africa

⁴³Department of Molecular Medicine and Haematology, University of the Witwatersrand, Johannesburg, South Africa

⁴⁴School of Pathology, Faculty of Health Sciences, University of the Witwatersrand, Johannesburg, South Africa

⁴⁵PathCare Vermaak, Pretoria, South Africa

⁴⁶Institute of Infectious Disease and Molecular Medicine, University of Cape Town, Cape Town, South Africa

⁴⁷Department of Molecular Pathology, School of Pathology, Faculty of Health Sciences, University of the Witwatersrand, Johannesburg, South Africa

⁴⁸Department of Global Health, University of Washington, Seattle, WA, USA

*These authors contributed equally: Raquel Viana, Sikhulile Moyo, Daniel G Amoako, Houriiyah Tegally, Cathrine Scheepers

**These authors jointly supervised the work: Simani Gaseitsiwe, Anne von Gottberg, Tulio de Oliveira

Summary

The severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) epidemic in southern Africa has been characterised by three distinct waves. The first was associated with a mix of SARS-CoV-2 lineages, whilst the second and third waves were driven by the Beta and Delta variants respectively¹⁻³. In November 2021, genomic surveillance teams in South Africa and Botswana detected a new SARS-CoV-2 variant associated with a rapid resurgence of infections in Gauteng Province, South Africa. Within three days of the first genome being uploaded, it was designated a variant of concern (Omicron) by the World Health Organization and, within three weeks, had been identified in 87 countries. The Omicron variant is exceptional for carrying over 30 mutations in the spike glycoprotein, predicted to influence antibody neutralization and spike function⁴. Here, we describe the genomic profile and early transmission dynamics of Omicron, highlighting the rapid spread in regions with high levels of population immunity.

Introduction

Since the onset of the COVID-19 pandemic in December 2019, variants of SARS-CoV-2 have emerged repeatedly. Some variants have spread worldwide and made major contributions to the cyclical infection waves that occur asynchronously in different regions. Between October and December 2020, the world witnessed the emergence of the first variants of concern (VOC). These variants exhibited increased transmissibility and/or immune evasion properties that threatened global efforts to control the pandemic. Although the Alpha, Beta and Gamma VOCs^{2,5} that emerged during this time disseminated globally and drove epidemic resurgences in many different countries, it was the highly transmissible Delta variant that subsequently displaced all other VOC in most regions of the world⁶. During its spread, the Delta variant evolved into multiple sub-lineages⁷, some of which demonstrated signs of having a growth advantage in certain locations⁸, prompting speculation that the next VOC to drive a resurgence of infections would be likely derived from Delta. However, in October 2021, while Delta was continuing to exhibit high levels of transmission in the Northern hemisphere, a large Delta wave was subsiding in southern Africa. The culmination of this wave coincided with the emergence of a novel SARS-CoV-2 variant that, within days of its near-simultaneous discovery in four individuals in Botswana, a traveler from South Africa in Hong Kong, and 54 individuals in South Africa, was designated by the World Health Organization as Omicron: the fifth VOC of SARS-CoV-2.

Results

Epidemic dynamics and detection of Omicron

The three distinct epidemic waves of SARS-CoV-2 experienced by southern African countries were each driven by different variants: the first by descendants of the B.1 lineage¹, the second by the Beta VOC^{2,9}, and the third by the Delta VOC³, with an estimated 2-5% of third wave cases in South Africa attributed to the C.1.2 lineage¹⁰ (**Fig. 1A**). Serosurveys conducted before

the Delta wave suggested high levels of exposure to SARS-CoV-2 (40-60%) in South Africa^{11,12}, Malawi¹³, and Zimbabwe¹⁴, and modelled estimates suggested seroprevalence of 70-80% across South Africa by October 2021¹⁵. Accordingly, the weeks following the third wave in South Africa, between 10 October and 15 November 2021, were marked by a period of lower-level transmission as indicated by a low incidence of reported COVID-19 cases (100-200 new cases per day) and low (<2%) test positivity rates (**Fig. 1A-1C**).

A rapid increase in COVID-19 cases was observed in mid-November 2021 in Gauteng province, the economic hub of South Africa containing the cities of Tshwane (Pretoria) and Johannesburg. Specifically, rising case numbers and test positivity rates were first noticed in Tshwane, initially associated with outbreaks in higher education settings. This resurgence of cases was accompanied by an increasing frequency of S-gene target failure (SGTF) during TaqPath-based (Thermo Fisher Scientific) diagnostic PCR testing: a phenomenon previously observed with the Alpha variant due to a deletion at amino acid positions 69 and 70 ($\Delta 69-70$) in the SARS-CoV-2 spike protein¹⁶. Given the low prevalence of Alpha in South Africa (**Fig. 1A**), targeted whole-genome sequencing of these specimens was prioritized.

On 19 November 2021, sequencing results of an initial batch of 8 SGTF samples collected between 14-16 November 2021 indicated that all were a new and genetically-distinct lineage of SARS-CoV-2. Further rapid sequencing identified the same variant in 29 of 32 routine diagnostic samples from multiple locations in Gauteng Province, indicating widespread circulation of this new variant by the second week of November. Crucially, this rise immediately preceded a sharp increase in reported case numbers (**Fig. 1C, Extended Data Fig. 1**). In the

following four days this lineage was confirmed by sequencing in another two provinces: KwaZulu-Natal (KZN) and the Western Cape (**Fig. 1B**).

Concurrently in Gaborone, Botswana (~360km from Tshwane), four genomes generated from samples collected on 11 November 2021, and sequenced on 17-18 November 2021 as part of weekly surveillance, displayed an unusual set of mutations. These were reported to the Botswana Ministry of Health and Wellness on 22 November 2021, as “unusual sequences” that were linked to a group of visitors (non-residents) on a diplomatic mission. The sequences were uploaded to GISAID^{17,18} on 23 November 2021, and it became apparent that they belonged to a new lineage. A further 15 genomically confirmed cases (not epidemiologically linked to the first four) were identified within the same week from various other locations in Botswana. All of these either had travel links from South Africa, or were contacts of someone with travel links.

On 24 November 2021, these SARS-CoV-2 genomes from both South Africa and Botswana were designated as belonging to a new PANGO lineage (B.1.1.529)¹⁹, later divided into sub-lineages alias BA.1 (the main clade), BA.2 and BA.3. On 26 November 2021, the lineage was designated a VOC and named Omicron by the WHO on the recommendation of the Technical Advisory Group on SARS-CoV-2 Virus Evolution²⁰. By the first week of December 2021, Omicron was causing a rapid and sustained increase in cases in South Africa and Botswana (**Fig. 1C, Extended Data Fig. 2** for Botswana). In Gauteng, weekly test positivity rates increased from <1% in the week beginning 31 October, to 16% in the week beginning 21 November 2021, and to 35% in the week beginning 28 November, concurrently with an exponential rise in COVID-19 incidence (**Fig. 1C, Extended Data Fig. 1**). Nationally, daily case numbers exceeded 22 000 (84% of the peak of the previous wave of infections) by 9 December 2021. At the same time, the proportion of TaqPath PCR tests with SGTF increased rapidly in all

provinces of South Africa reaching ~90% nationally by the week beginning 21 November 2021, strongly indicating that the fourth wave was being driven by Omicron: an indication that has now been confirmed by virus genome sequencing in all provinces (**Fig. 1C**). Similarly, Botswana experienced a sharp increase in cases, doubling every 2-3 days late November to early December 2021, transitioning from a 7-day moving average of <10 cases/100 000 to above 25 cases/100,000 in less than 10 days (**Extended Data Fig. 2**).

By 16 December 2021, Omicron had been detected in 87 countries, both in samples from travelers returning from southern Africa, and in samples from routine community testing (**Extended Data Fig. 3**).

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

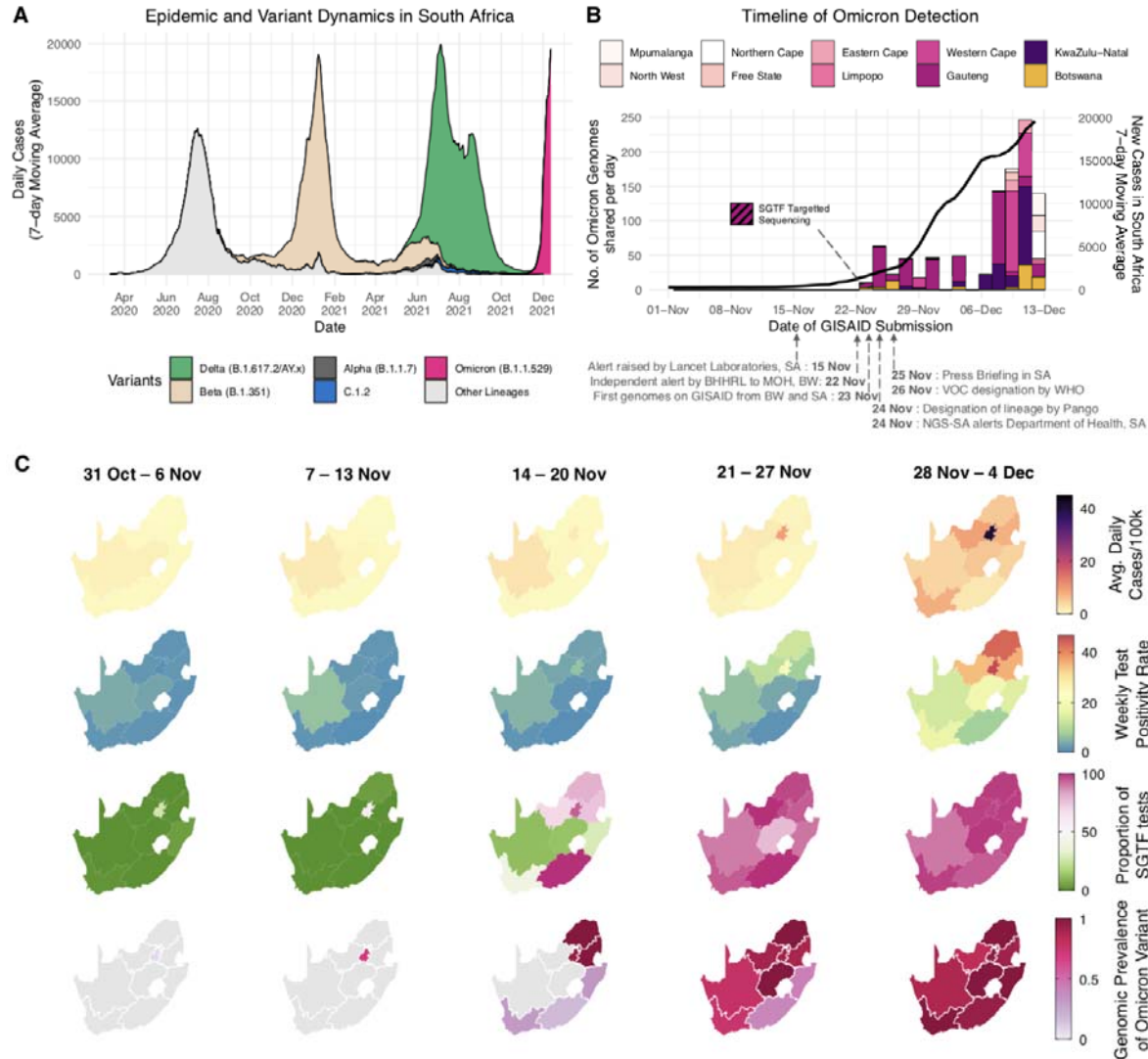


Figure 1: Detection of Omicron variant A) The progression of daily reported cases in South Africa from March 2020 to December 2021. The 7-day rolling average of daily case numbers is coloured by the inferred proportion of variants responsible for the infections, as calculated by genomic surveillance data on GISAID. B) Timeline of Omicron detection in Botswana and South Africa. Bars represent the number of Omicron genomes shared per day, according to the date they were uploaded to GISAID, while the line represents the 7-day moving average of daily new cases in South Africa. C) Weekly progression of average daily cases per 100,000, test positivity rates, proportion of SGTF tests (on the TaqPath COVID-19 PCR assay) and genomic prevalence of Omicron in nine provinces of South Africa for five weeks from 31 October to 4

December 2021. Note that because of heterogeneous use of the TaqPath PCR assay across provinces, the proportion of SGTF tests illustrated for the Eastern Cape Province in weeks of 14 - 20 Nov and 21 - 27 Nov are based on only 2 and 4 data points respectively.

The evolutionary origins of Omicron

To determine when and where Omicron likely originated, we analyzed all 686 available Omicron genomes (including 248 from southern Africa and 438 from elsewhere in the world) retrieved from GISAID (date of access 7 December 2021)^{17,18}, in the context of a global reference set of representative SARS-CoV-2 genomes (n=12 609) collected between December 2019 and November 2021. Preliminary maximum-likelihood phylogenies identified the BA.1/Omicron sequences as a monophyletic clade rooted within the B.1.1 lineage (Nextstrain clade 20B), with no clear basal progenitor (**Fig. 2A**). Importantly, the BA.1/Omicron cluster is highly phylogenetically distinct from any known VOC or variants of interest (VOI) and from any other lineages known to be circulating in southern Africa (e.g. C.1.2) (**Fig. 2A**). More recently, two related lineages have emerged (BA.2 and BA.3), both sharing many, but not all of the characteristic mutations of BA.1/Omicron and both having many unique mutations of their own. We primarily focus here on the BA.1 lineage which is rapidly spreading in multiple countries around the world and is the lineage first officially designated as the Omicron VOC.

Time-calibrated Bayesian phylogenetic analysis of all BA.1 assigned genomes from southern Africa (as of 11 December 2021, n=553) estimated the time when the most recent common ancestor of the analysed BA.1 lineage sequences existed to be 9 October 2021 (95% credible intervals 30 September - 20 October) with a per-day exponential growth rate of 0.136 (95% confidence interval (CI) 0.100 – 0.173) reflecting a doubling time of 5.1 days (95% CI 4.0 – 6.9) (**Fig 2B**). These estimates are robust to whether the evolutionary rate is estimated from the data or fixed to previously estimated values (**Extended Data Table S1**). Limiting the analysis to a

subset of genomes from Gauteng Province only (279 genomes) yields a faster growth rate estimate with a doubling time of 1.8 days (95% CI 1.4 – 3.0) (**Extended Data S1**). Using a phylodynamic model that accounts for variable genome sampling through time (birth-death skyline model) yields doubling times of Omicron in South Africa between 3.88 and 3.99 days and mean effective reproduction number estimates (R_e) of 2.74 to 2.79 for the period of early November to early December (Extended Data Table X). Spatiotemporal phylogeographic analysis indicates that the BA.1/Omicron variant spread from the Gauteng province of South Africa to seven of the eight other provinces and to two regions of Botswana from late October to late November 2021, and shows evidence of more recent transmission within and between other South African provinces (Fig 2C).

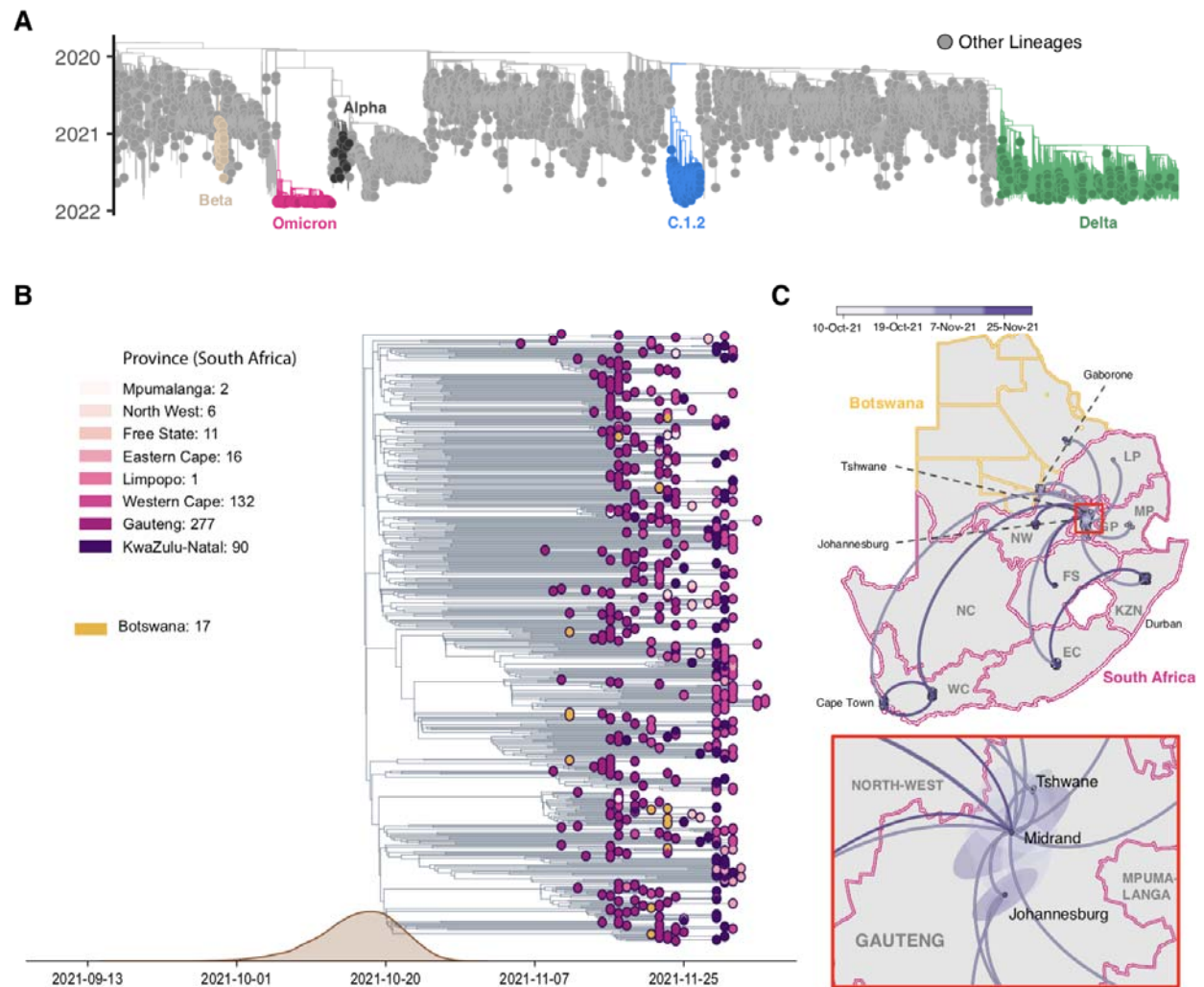


Figure 2: Evolution of Omicron. A) Time-resolved maximum likelihood phylogeny of 13,295 SARS-CoV-2 sequences; 9,944 of these are from Africa (denoted with tip point circle shapes). Alpha, Beta and Delta VOCs and the C.1.2 lineage, recently circulating in South Africa, are denoted in black, brown, green and blue respectively. The newly identified SARS-CoV-2 Omicron variant is shown in pink. Genomes of other lineages are shown in grey. B) Time-resolved maximum clade credibility phylogeny of the Omicron cluster of southern African genomes ($n = 553$), with locations indicated. The distribution of estimated time of origin is also shown. C) Spatiotemporal reconstruction of the spread of the Omicron variant in Southern Africa

with an inset of Gauteng province. Circles represent nodes of the maximum clade credibility phylogeny, coloured according to their inferred time of occurrence (scale in top panel). Shaded areas represent the 80% highest posterior density interval and depict the uncertainty of the phylogeographic estimates for each node. Solid curved lines denote the links between nodes and the directionality of movement is anticlockwise along the curve.

Molecular profile of Omicron

Compared to Wuhan-Hu-1, Omicron carries 15 mutations in the spike receptor-binding domain (RBD) (**Fig. 3**), five of which (G339D, N440K, S477N, T478K, N501Y) have been shown individually to enhance hACE2 binding²¹. Seven of the RBD mutations (K417N, G446S, E484A, Q493R, G496S, Q498R and N501Y) are expected to have moderate to strong impacts on binding of at least three of the four major classes of RBD-targeted neutralizing antibodies (NAbs)²²⁻²⁴. These RBD mutations coupled with four amino acid substitutions (A67V, T95I, G142D, and L212I), three deletions (69-70, 143-145 and 211) and an insertion (EPE between 214 and 215) in the N-terminal domain (NTD)²⁵, are predicted to underlie the substantially reduced sensitivity of Omicron to neutralization by anti-SARS-CoV-2 antibodies induced by either infection or vaccination^{26,27}. These mutations also involve key structural epitopes targeted by some of the currently authorized monoclonal antibodies, particularly bamlanivimab + etesevimab and casirivimab + imdevimab²⁸⁻³⁰. Preliminary analysis suggests that although the spike mutations involve a number of T cell and B cell epitopes, the majority of epitopes (>70%) remain unaffected³¹.

Omicron also has a cluster of three mutations (H655Y, N679K and P681H) adjacent to the S1/S2 furin cleavage site (FCS) which are likely to enhance spike protein cleavage and fusion with host cells^{32,33} and which could also contribute to enhanced transmissibility³⁴ (**Extended**

Data Fig. 4).

Outside of the spike protein, a deletion in nsp6 (105-107del), in the same region as deletions seen in Alpha, Beta, Gamma and Lambda, may have a role in evasion of innate immunity³⁵; and the double mutation in nucleocapsid (R203K, G204R), also present in Alpha, Gamma and C.1.2, has been associated with enhanced infectivity in human lung cells³⁶.

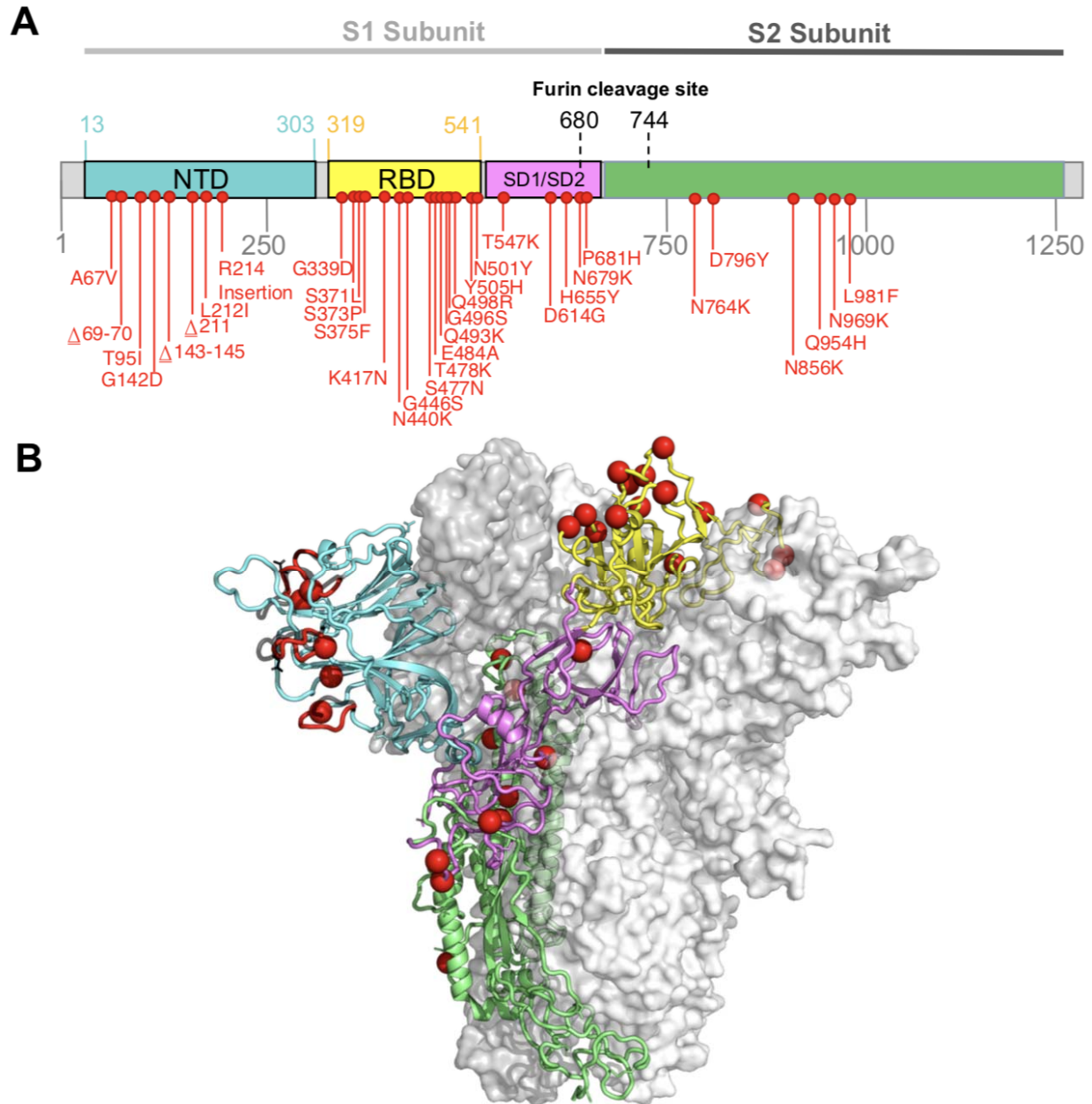


Figure 3. Molecular profile of Omicron A) Amino-acid mutations on the spike gene of the BA.1/Omicron variant. B) Structure of the SARS-CoV-2 Spike trimer, showing a single spike protomer in cartoon view. The N terminal domain, receptor binding domain, subdomains 1 and 2, and the S2 protein are shown in cyan, yellow, pink, and green respectively. Red spheres indicate the alpha carbon positions for each omicron variant residue. NTD-specific loop

insertions/deletions are shown in red, with the original loop shown in transparent black.

Omicron is not obviously recombinant

Given the large number of mutations differentiating BA.1/Omicron and BA.2 from other known SARS-CoV-2 lineages it was considered plausible that either (i) both of these lineages might have descended from a common recombinant ancestor, or (ii) that one of the BA lineages might have originated via recombination between a virus in the other BA lineage and a virus in a non-BA lineage. We tested these hypotheses using a variety of recombination detection approaches (implemented in the programs GARD³⁷; 3SEQ³⁸; and RDP5³⁹) to identify potential signals of recombination in sequence datasets containing BA.1 and BA.2 sequences together with sequences representative of global SARS-CoV-2 genomic diversity.

Potential evidence of recombination was identified by 3SEQ, GARD and RDP5 in these datasets. The most likely recombination breakpoint locations were located between nucleotide positions 20520 and 21619 (near the start of the S-gene; supported by GARD, RDP5 and 3SEQ) and between 23609 and 23614 (near the middle of the S-gene; supported by GARD and 3SEQ). 3SEQ identified an additional breakpoint at nucleotide position 24513 (toward the end of the S-gene). Phylogenetic analysis of the genome regions bounded by these breakpoints (genome coordinates 1450-20520, 21619-23609 and 23614-24513) revealed no support for a recombinant origin for either the BA.1 or BA.2 lineages (**Extended Data Fig. 5**). Although one BA.1 isolate (Botswana/R43B66) displayed evidence of having potentially inherited nucleotides 23614-24513 from a Delta virus by recombination, there was no strong phylogenetic support for the clustering of this sequence with Delta viruses. Further, read coverage in this region of the Botswana/R43B66 sequence was so low that we were unable to exclude the possibility that the

apparent recombination signal was attributable to a combination of miscalled/uncalled nucleotides and alignment uncertainty.

Although we found no convincing phylogenetic or statistical evidence of either the most recent common ancestor of BA.1 and BA.2 being recombinant, or of the most recent common ancestors of either the BA.1 or BA.2 lineages having been derived through recombination, it should be noted that recombination tests in general will not have sufficient statistical power to reliably identify evidence of individual recombination events that result in transfers of less than ~5 contiguous polymorphic nucleotide sites between genomes^{37,40,41}. Further, if BA.1 and/or BA.2 are the products of a series of multiple partially overlapping recombination events occurring across multiple temporally clustered replication cycles, the complex patterns of nucleotide variation that might result could be extremely difficult to interpret as recombination using the methods applied here⁴².

Signs of strong selection and epistasis during the origin and ongoing evolution of the Omicron lineage

We applied a selection analysis pipeline to all available sequences designated as BA.1 in GISAID as of 8 December 2021. The analysis followed the procedure described previously³⁵, and downsampled alignments of individual protein encoding regions to obtain a median of 25 unique Omicron haplotype sequences and 107 unique haplotype sequences for each gene/ORF from a representative selection of other SARS-CoV-2 lineages (used as background sequences to contextualize evolution within the Omicron sub-clade).

We detected evidence of gene-wide positive selection (using the BUSTED method⁴³) acting on six genes/ORFs since the ancestral BA.1/Omicron and BA.2 lineage split from the B.1.1 lineage: S-gene ($p < 0.0001$), exonuclease ($p < 0.0001$), nsp6 ($p = 0.001$), M-gene ($p = 0.002$), N-gene

($p = 0.006$), and E-gene ($p = 0.05$). In all six genes, this selection was strong ($dN/dS > 10$), and occurred in bursts ($\leq 6\%$ of branch/site combinations selected). The branch separating BA.1/Omicron from its most recent B.1.1 ancestor had the most prominent selection signal (which was strongest in the S-gene; BUSTED p -value < 0.0001 with $dN/dS > 100$ at $\sim 0.5\%$ of S-gene codon sites⁴⁴), strongly supporting the hypothesis that adaptive evolution played a significant role in the mutational divergence of Omicron from other B.1.1 SARS-CoV-2 lineages. Relative to the intensity of selection evident within the background B.1.1 lineages, selection in three genes was likely significantly intensified in the ancestral Omicron lineage: S-gene (intensification factor $K = 2.0$; $p < 0.0001$ ⁴⁵), exonuclease ($K = 4.0$; $p < 0.0001$), and nsp6 ($K = 5.1$; $p = 0.02$).

Among 294 codon sites that are polymorphic among the BA.1/Omicron sequences analysed, 32 were found to have experienced episodic positive selection since BA.1 split from the B.1.1 lineage (MEME $p \leq 0.01$, **Extended Data Table S2**⁴⁶). Sixteen (50%) of these codon sites are in the S-gene, 13 of which contain BA.1 lineage-defining mutations (i.e. these selection signals reflect mutations that occurred within the ancestral Omicron lineage). The three positively selected codon sites that did not correspond to sites of lineage-defining mutations (S/346, S/452, and S/701) are particularly notable as these are attributable to mutations that have occurred since the MRCA of the analysed BA.1 sequences. The mutations driving the positive selection signals at these three sites in the Omicron S-gene converge on mutations seen in other VOCs or VOIs (R346K in Mu, L452R in Delta, and A701V in Beta and Iota). The A701V mutation, the precise impact of which is currently unknown, is one of 19 in a proposed “501Y lineage Spike meta-signature” comprising the set of mutations that were most adaptive during

the evolution of the Alpha, Beta and Gamma VOC lineages³⁵. Further, both R346K and L452R are known to impact antibody binding²³ and both of the codon sites where these mutations occur display evidence for directional selection (using the FADE method⁴⁷). These selective patterns suggest that, during its current explosive spread, Omicron may be undergoing additional evolution to modify its neutralization profile.

Potential for increased transmissibility and immune evasion

We estimated that Omicron had a growth advantage of 0.24 (95% CI: 0.16-0.33) per day over Delta in Gauteng, South Africa (**Fig. 4A**). This corresponds to a 5.4-fold (95% CI: 3.1-10.1) weekly increase in cases compared to Delta. The growth advantage of Omicron is likely to be mediated by (i) an increase relative to other variants of its intrinsic transmissibility, (ii) an increase relative to other variants in its capacity to infect, and be transmitted from, previously infected and vaccinated individuals; or (iii) both.

The predicted combination of transmissibility and immune evasion for Omicron strongly depends on the assumed level of current population immunity against infection by, and transmission of, the competing variant Delta that is afforded by prior-infections with wild-type Wuhan, Beta, Delta, and other strains, and/or vaccination (**Fig. 4B**). For moderate levels of population immunity against Delta ($\Omega = 0.4$), immune evasion alone cannot explain the observed growth advantage of Omicron (**Fig. 4C**). For medium levels of immunity against Delta ($\Omega = 0.6$), very high levels of immune evasion could explain the observed growth advantage without an additional increase in transmissibility (**Fig. 4D**). For high levels of population immunity against Delta ($\Omega = 0.8$), even moderate levels of immune evasion (~25-50%) can explain the observed growth advantage without an additional increase in transmissibility (**Fig. 4E**). The results of seroprevalence studies and vaccination coverage (~40% of the adult

population) in South Africa suggest that the proportion of the population with potential immunity against Delta and earlier variants is likely to be above 60%^{11,12}. We thus argue that the population level of protective immunity against Delta is high, and that partial immune evasion is a major driver for the observed dynamics of Omicron in South Africa. This notion is supported by recent findings that show an increased risk of SARS-CoV-2 reinfection associated with the emergence of Omicron in South Africa⁴⁸ and the initial results from neutralization assays^{26,27}. In addition to immune evasion, an increase, or decrease, in the transmissibility of Omicron compared to Delta cannot, however, be ruled out.

There are a number of limitations to this analysis. First, we estimated the growth advantage of Omicron based on early sequence data only. These data could be biased due to targeted sequencing of SGTF samples and stochastic effects (e.g., superspreading) in a low incidence setting, which can lead to overestimates of the growth advantage, and consequently of the increased transmissibility and immune evasion. Second, without reliable estimates of the level of protective immunity against Delta in South Africa, we cannot obtain precise estimates of transmissibility or immune evasion of Omicron.

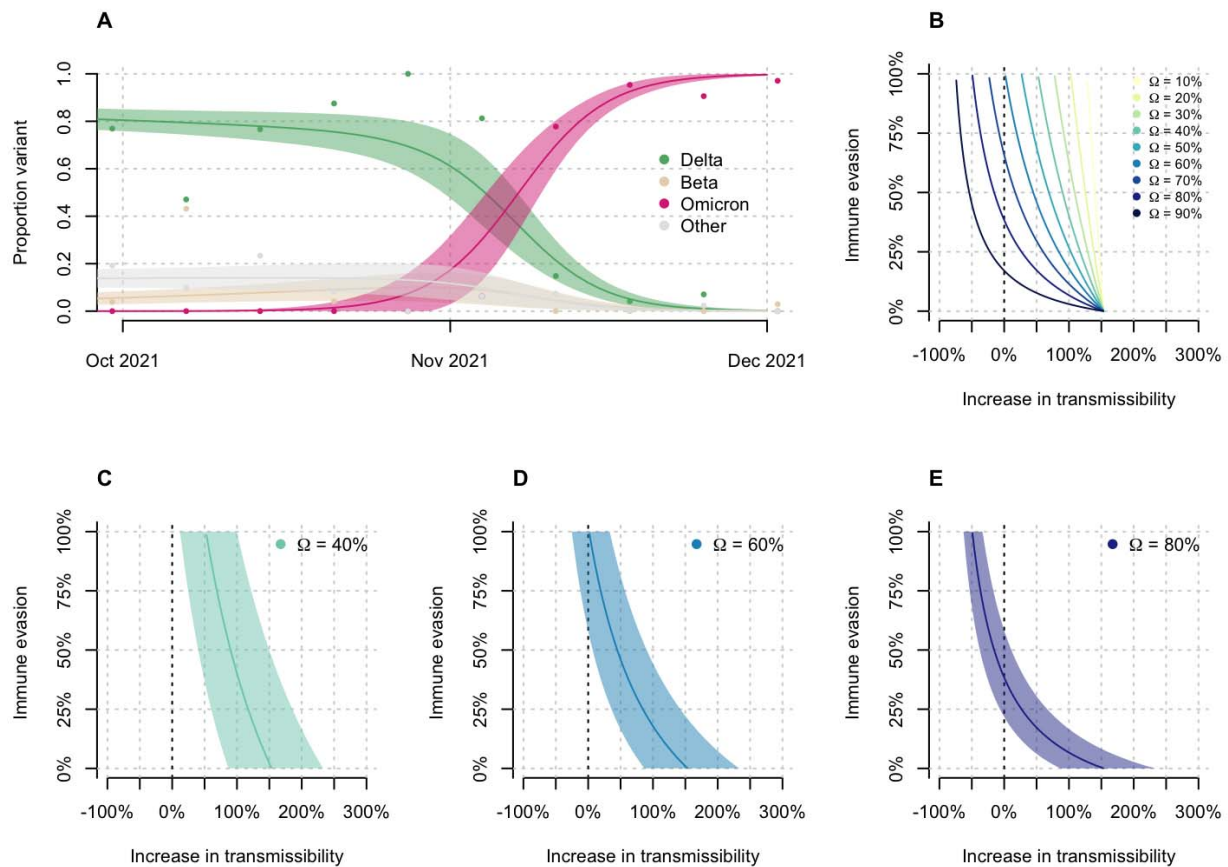


Figure 4: Growth of Omicron in Gauteng, South Africa, and relationship between potential increase in transmissibility and immune evasion. (A) Omicron rapidly outcompeted Delta in November 2021. Model fits are based on a multinomial logistic regression. Dots represent the weekly proportions of variants. (B) The relationship between the potential increase in transmissibility and immune evasion strongly depends on the assumed level of current population immunity against Delta (Ω). (C-E) Relationship for a population immunity of 40%, 60%, and 80% against infection and transmission with Delta. The dark vertical dashed line indicates equal transmissibility of Omicron compared to Delta. Shaded areas correspond to the 95% CIs of the model estimates.

Conclusion

Strong genomic surveillance systems in South Africa and Botswana enabled the identification of Omicron within a week of observing a resurgence in cases in Gauteng Province. Immediate notification of the WHO and early designation as a VOC has stimulated global scientific efforts and has given other countries time to prepare their response. Omicron is now driving a fourth wave of the SARS-CoV-2 epidemic in southern Africa, and is spreading rapidly in several other countries. Genotypic and phenotypic data suggest that Omicron has the capacity for substantial evasion of neutralizing antibody responses, and modelling suggests that immune evasion could be a major driver of the observed transmission dynamics. Close monitoring of the spread of Omicron in countries outside southern Africa will be necessary to better understand its transmissibility and the capacity of this variant to evade post-infection and vaccine-elicited immunity. Neutralizing antibodies are only one component of the immune protection from vaccines and prior infection, and the cellular immune response is predicted to be less affected by the mutations in Omicron. Vaccination therefore remains critical to protect those at highest risk of severe disease and death. The emergence and rapid spread of Omicron poses a threat to the world and a particular threat in Africa, where fewer than one in ten people is fully vaccinated.

Main references

1. Tegally, H. *et al.* Sixteen novel lineages of SARS-CoV-2 in South Africa. *Nat. Med.* **27**, 440–446 (2021).
2. Tegally, H. *et al.* Detection of a SARS-CoV-2 variant of concern in South Africa. *Nature* **592**, 438–443 (2021).
3. Tegally, H. *et al.* Rapid replacement of the Beta variant by the Delta variant in South Africa. *medRxiv* (2021) doi:10.1101/2021.09.23.21264018.
4. Martin, D. P. *et al.* Selection analysis identifies significant mutational changes in Omicron that are likely to influence both antibody neutralization and Spike function (part 1 of 2). <https://virological.org/t/selection-analysis-identifies-significant-mutational-changes-in-omicron-that-are-likely-to-influence-both-antibody-neutralization-and-spike-function-part-1-of-2/771> (2021).
5. Faria, N. R. *et al.* Genomics and epidemiology of the P.1 SARS-CoV-2 lineage in Manaus, Brazil. *Science* **372**, 815–821 (2021).
6. Dhar, M. S. *et al.* Genomic characterization and epidemiology of an emerging SARS-CoV-2 variant in Delhi, India. *Science* **374**, 995–999 (2021).
7. New AY lineages – Pango Network. <https://www.pango.network/new-ay-lineages/>.
8. Eales, O. *et al.* SARS-CoV-2 lineage dynamics in England from September to November 2021: high diversity of Delta sub-lineages and increased transmissibility of AY.4.2. *medRxiv* (2021).
9. Wilkinson, E. *et al.* A year of genomic surveillance reveals how the SARS-CoV-2 pandemic unfolded in Africa. *Science* **374**, 423–431 (2021).
10. Scheepers, C. *et al.* The continuous evolution of SARS-CoV-2 in South Africa: a new lineage with rapid accumulation of mutations of concern and global detection. *medRxiv* (2021) doi:10.1101/2021.08.20.21262342.
11. Kleynhans, J. *et al.* SARS-CoV-2 Seroprevalence in a Rural and Urban Household Cohort

during First and Second Waves of Infections, South Africa, July 2020-March 2021.

Emerging Infect. Dis. **27**, 3020–3029 (2021).

12. Vermeulen, M. *et al.* Prevalence of anti-SARS-CoV-2 antibodies among blood donors in South Africa during the period January-May 2021. *Res. Sq.* (2021) doi:10.21203/rs.3.rs-690372/v1.
13. Mandolo, J. *et al.* Dynamics of SARS-CoV-2 exposure in Malawian blood donors: a retrospective seroprevalence analysis between January 2020 and February 2021. *medRxiv* (2021) doi:10.1101/2021.08.18.21262207.
14. Fryatt, A. *et al.* Community SARS-CoV-2 seroprevalence before and after the second wave of SARS-CoV-2 infection in Harare, Zimbabwe. *EClinicalMedicine* **41**, 101172 (2021).
15. South African COVID-19 Modelling Consortium. *COVID-19 modelling update: Considerations for a potential fourth wave*. 20 <https://www.nicd.ac.za/wp-content/uploads/2021/11/SACMC-Fourth-wave-report-17112021-final.pdf> (2021).
16. Volz, E. *et al.* Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7 in England. *Nature* (2021) doi:10.1038/s41586-021-03470-x.
17. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision to reality. *Euro Surveill.* **22**, 30494 (2017).
18. Elbe, S. & Buckland-Merrett, G. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Global Challenges* **1**, 33–46 (2017).
19. Rambaut, A. *et al.* A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat. Microbiol.* **5**, 1403–1407 (2020).
20. World Health Organization. Classification of Omicron (B.1.1.529): SARS-CoV-2 Variant of Concern. [https://www.who.int/news-room/statements/26-11-2021-classification-of-omicron-\(b.1.1.529\)-sars-cov-2-variant-of-concern](https://www.who.int/news-room/statements/26-11-2021-classification-of-omicron-(b.1.1.529)-sars-cov-2-variant-of-concern) (2021).
21. Starr, T. N. *et al.* Deep Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints on Folding and ACE2 Binding. *Cell* **182**, 1295-1310.e20 (2020).

22. Greaney, A. J. *et al.* Mapping mutations to the SARS-CoV-2 RBD that escape binding by different classes of antibodies. *Nat. Commun.* **12**, 4196 (2021).
23. Greaney, A. J. *et al.* Complete Mapping of Mutations to the SARS-CoV-2 Spike Receptor-Binding Domain that Escape Antibody Recognition. *Cell Host Microbe* **29**, 44-57.e9 (2021).
24. Greaney, A. J. *et al.* Comprehensive mapping of mutations in the SARS-CoV-2 receptor-binding domain that affect recognition by polyclonal human plasma antibodies. *Cell Host Microbe* **29**, 463-476.e6 (2021).
25. McCallum, M. *et al.* N-terminal domain antigenic mapping reveals a site of vulnerability for SARS-CoV-2. *Cell* **184**, 2332-2347.e16 (2021).
26. Cele, S. *et al.* SARS-CoV-2 Omicron has extensive but incomplete escape of Pfizer BNT162b2 elicited neutralization and requires ACE2 for infection. *medRxiv* (2021) doi:10.1101/2021.12.08.21267417.
27. Rössler, A., Riepler, L., Bante, D., Laer, D. von & Kimpel, J. SARS-CoV-2 B.1.1.529 variant (Omicron) evades neutralization by sera from vaccinated and convalescent individuals. *medRxiv* (2021) doi:10.1101/2021.12.08.21267491.
28. Starr, T. N., Greaney, A. J., Dingens, A. S. & Bloom, J. D. Complete map of SARS-CoV-2 RBD mutations that escape the monoclonal antibody LY-CoV555 and its cocktail with LY-CoV016. *Cell Rep. Med.* **2**, 100255 (2021).
29. Starr, T. N. *et al.* Prospective mapping of viral mutations that escape antibodies used to treat COVID-19. *Science* **371**, 850–854 (2021).
30. Cao, Y. R. *et al.* B.1.1.529 escapes the majority of SARS-CoV-2 neutralizing antibodies of diverse epitopes. *BioRxiv* (2021) doi:10.1101/2021.12.07.470392.
31. Bernasconi, A. *et al.* Report on Omicron Spike mutations on epitopes and immunological/epidemiological/kinetics effects from literature - SARS-CoV-2 coronavirus / nCoV-2019 Genomic Epidemiology - Virological. <https://virological.org/t/report-on-omicron-spike-mutations-on-epitopes-and-immunological-epidemiological-kinetics-effects-from->

- literature/770 (2021).
32. Brown, J. C. *et al.* Increased transmission of SARS-CoV-2 lineage B.1.1.7 (VOC 2020212/01) is not accounted for by a replicative advantage in primary airway cells or antibody escape. *BioRxiv* (2021) doi:10.1101/2021.02.24.432576.
 33. Saito, A. *et al.* SARS-CoV-2 spike P681R mutation enhances and accelerates viral fusion. *BioRxiv* (2021) doi:10.1101/2021.06.17.448820.
 34. Mlcochova, P. *et al.* SARS-CoV-2 B.1.617.2 Delta variant replication and immune evasion. *Nature* **599**, 114–119 (2021).
 35. Martin, D. P. *et al.* The emergence and ongoing convergent evolution of the SARS-CoV-2 N501Y lineages. *Cell* **184**, 5189-5200.e7 (2021).
 36. Wu, H. *et al.* Nucleocapsid mutations R203K/G204R increase the infectivity, fitness, and virulence of SARS-CoV-2. *Cell Host Microbe* (2021) doi:10.1016/j.chom.2021.11.005.
 37. Kosakovsky Pond, S. L., Posada, D., Gravenor, M. B., Woelk, C. H. & Frost, S. D. W. GARD: a genetic algorithm for recombination detection. *Bioinformatics* **22**, 3096–3098 (2006).
 38. Lam, H. M., Ratmann, O. & Boni, M. F. Improved algorithmic complexity for the 3SEQ recombination detection algorithm. *Mol. Biol. Evol.* **35**, 247–251 (2018).
 39. Martin, D. P. *et al.* RDP5: a computer program for analyzing recombination in, and removing signals of recombination from, nucleotide sequence datasets. *Virus Evol.* **7**, veaa087 (2021).
 40. Boni, M. F., Posada, D. & Feldman, M. W. An exact nonparametric method for inferring mosaic structure in sequence triplets. *Genetics* **176**, 1035–1047 (2007).
 41. Posada, D. & Crandall, K. A. Evaluation of methods for detecting recombination from DNA sequences: computer simulations. *Proc Natl Acad Sci USA* **98**, 13757–13762 (2001).
 42. van der Walt, E. *et al.* Rapid host adaptation by extensive recombination. *J. Gen. Virol.* **90**, 734–746 (2009).

43. Wisotsky, S. R., Kosakovsky Pond, S. L., Shank, S. D. & Muse, S. V. Synonymous Site-to-Site Substitution Rate Variation Dramatically Inflates False Positive Rates of Selection Analyses: Ignore at Your Own Peril. *Mol. Biol. Evol.* **37**, 2430–2439 (2020).
44. Smith, M. D. *et al.* Less is more: an adaptive branch-site random effects model for efficient detection of episodic diversifying selection. *Mol. Biol. Evol.* **32**, 1342–1353 (2015).
45. Wertheim, J. O., Murrell, B., Smith, M. D., Kosakovsky Pond, S. L. & Scheffler, K. RELAX: detecting relaxed selection in a phylogenetic framework. *Mol. Biol. Evol.* **32**, 820–832 (2015).
46. Murrell, B. *et al.* Detecting individual sites subject to episodic diversifying selection. *PLoS Genet.* **8**, e1002764 (2012).
47. Kosakovsky Pond, S. L., Poon, A. F. Y., Leigh Brown, A. J. & Frost, S. D. W. A maximum likelihood method for detecting directional evolution in protein sequences and its application to influenza A virus. *Mol. Biol. Evol.* **25**, 1809–1824 (2008).
48. Pulliam, J. R. C. *et al.* SARS-CoV-2 reinfection trends in South Africa: analysis of routine surveillance data. *medRxiv* (2021) doi:10.1101/2021.11.11.21266068.

Figure legends

Figure 1: Detection of Omicron variant A) The progression of daily reported cases in South Africa from March 2020 to December 2021. The 7-day rolling average of daily case numbers is coloured by the inferred proportion of variants responsible for the infections, as calculated by genomic surveillance data on GISAID. B) Timeline of Omicron detection in Botswana and South Africa. Bars represent the number of Omicron genomes shared per day, according to the date they were uploaded to GISAID, while the line represents the 7-day moving average of daily new cases in South Africa. C) Weekly progression of average daily cases per 100,000, test positivity rates, proportion of SGTF tests (on the TaqPath COVID-19 PCR assay) and genomic prevalence of Omicron in nine provinces of South Africa for five weeks from 31 October to 4 December 2021. Note that because of heterogeneous use of the TaqPath PCR assay across provinces, the proportion of SGTF tests illustrated for the Eastern Cape Province in weeks of 14 - 20 Nov and 21 - 27 Nov are based on only 2 and 4 data points respectively.

Figure 2: Evolution of Omicron A) Time-resolved maximum likelihood phylogeny of 13,295 SARS-CoV-2 sequences; 9,944 of these are from Africa (denoted with tip point circle shapes). Alpha, Beta and Delta VOCs and the C.1.2 lineage, recently circulating in South Africa, are denoted in black, brown, green and blue respectively. The newly identified SARS-CoV-2 Omicron variant is shown in pink. Genomes of other lineages are shown in grey. B) Time-resolved maximum clade credibility phylogeny of the Omicron cluster of southern African genomes ($n = 553$), with locations indicated. The distribution of estimated time of origin is also shown. C) Spatiotemporal reconstruction of the spread of the Omicron variant in Southern Africa with an inset of Gauteng province. Circles represent nodes of the maximum clade credibility phylogeny, coloured according to their inferred time of occurrence (scale in top panel). Shaded areas represent the 80% highest posterior density interval and depict the uncertainty of the

phylogeographic estimates for each node. Solid curved lines denote the links between nodes and the directionality of movement is anticlockwise along the curve.

Figure 3. Molecular profile of Omicron A) Amino-acid mutations on the spike gene of the BA.1/Omicron variant. B) Structure of the SARS-CoV-2 Spike trimer, showing a single spike protomer in cartoon view. The N terminal domain, receptor binding domain, subdomains 1 and 2, and the S2 protein are shown in cyan, yellow, pink, and green respectively. Red spheres indicate the alpha carbon positions for each omicron variant residue. NTD-specific loop insertions/deletions are shown in red, with the original loop shown in transparent black.

Figure 4: Growth of Omicron in Gauteng, South Africa, and relationship between potential increase in transmissibility and immune evasion. (A) Omicron rapidly outcompeted Delta in November 2021. Model fits are based on a multinomial logistic regression. Dots represent the weekly proportions of variants. (B) The relationship between the potential increase in transmissibility and immune evasion strongly depends on the assumed level of current population immunity against Delta (Ω). (C-E) Relationship for a population immunity of 40%, 60%, and 80% against infection and transmission with Delta. The dark vertical dashed line indicates equal transmissibility of Omicron compared to Delta. Shaded areas correspond to the 95% CIs of the model estimates.

Methods

Epidemiological dynamics

We analyzed daily cases of SARS-CoV-2 in South Africa up to 14 December 2021 from publicly released data provided by the National Department of Health and the National Institute for Communicable Diseases. This was accessible through the repository of the Data Science for Social Impact Research Group at the University of Pretoria (<https://github.com/dsfsi/covid19za>)^{49,50}. The National Department of Health releases daily updates on the number of confirmed new cases, deaths and recoveries, with a breakdown by province. Daily case numbers for Botswana were obtained via Our World in Data (OWID) COVID-19 data repository (<https://github.com/owid/covid-19-data>). We consulted estimates for the R_e of SARS-CoV-2 in South Africa and Botswana from the 'covid-19-Re' data repository (<https://github.com/covid-19-Re/dailyRe-Data>)⁵¹. We obtained test positivity data from weekly reports from the National Institute for Communicable Diseases (NICD)⁵². Data to calculate the proportion of positive Thermo Fisher TaqPath COVID-19 PCR tests with SGTF in South Africa was obtained from the National Health Laboratory Service and Lancet Laboratories. Test positivity data for Botswana was obtained from the National Health Laboratory through 6 December 2021. All data visualization was generated through the ggplot package in R⁵³.

SARS-CoV-2 sampling

As part of the NGS-SA, seven sequencing hubs in South Africa receive randomly selected samples for sequencing every week according to approved protocols at each site⁵⁴. These samples include remnant nucleic acid extracts or remnant nasopharyngeal and oropharyngeal swab samples from routine diagnostic SARS-CoV-2 PCR testing from public and private laboratories in South Africa. In response to a focal resurgence of COVID-19 in the City of Tshwane Metropolitan Municipality in Gauteng Province in November, we enriched our routine sampling with additional samples from the affected area, including initial targeted sequencing of

SGTF samples. In Botswana, all public and private laboratories submit randomly selected residual nasopharyngeal and oropharyngeal PCR positive samples weekly to the National Health Laboratory (NHL) and the Botswana Harvard HIV Reference Laboratory (BHHRL) for sequencing.

Ethical statement

The genomic surveillance in South Africa was approved by the University of KwaZulu–Natal Biomedical Research Ethics Committee (BREC/00001510/2020), the University of the Witwatersrand Human Research Ethics Committee (HREC) (M180832), Stellenbosch University HREC (N20/04/008_COVID-19), University of Cape Town HREC (383/2020), University of Pretoria HREC (H101/17) and the University of the Free State Health Sciences Research Ethics Committee (UFS-HSD2020/1860/2710). The genomic sequencing in Botswana was conducted as part of the national vaccine roll-out plan and was approved by the Health Research and Development Committee (Health Research Ethics body, HRDC#00948 and HRDC#00904). Individual participant consent was not required for the genomic surveillance. This requirement was waived by the Research Ethics Committees.

Ion Torrent Genexus Integrated Sequencer methodology for rapid whole genome sequencing of SARS-CoV-2

Viral RNA was extracted using the MagNA Pure 96 DNA and Viral Nucleic Acid kit on the automated MagNA Pure 96 system (Roche Diagnostics, USA) as per the manufacturer's instructions. Extracts were then screened by qPCR to acquire the mean cycle threshold (Ct) values for the SARS-CoV-2 N-gene and ORF1ab-gene using the TaqMan 2019-nCoV assay kit v1 (ThermoFisher Scientific, USA) on the ViiA7 Real-time PCR system (ThermoFisher Scientific, USA) as per the manufacturer's instructions. Extracts were sorted into batches of N=8 within a Ct range difference of 5 for a maximum of two batches per run. Extracts with <200

copies were sequenced using the low viral titer protocol. Next-generation sequencing was performed using the Ion AmpliSeq SARS-CoV-2 Research Panel on the Ion Torrent Genexus Integrated Sequencer (ThermoFisher Scientific, USA) which combines automated cDNA synthesis, library preparation, templating preparation and sequencing within 24 hours. The Ion AmpliSeq SARS-CoV-2 Research Panel consists of 2 primer pools targeting 237 amplicons tiled across the SARS-CoV-2 genome providing >99% coverage of the SARS-CoV-2 genome (~30 kb) and an additional 5 primer pairs targeting human expression controls. The SARS-CoV-2 amplicons range from 125 to 275 bp in length. TRINITY was utilised for de novo assembly and the Iterative Refinement Meta-Assembler (IRMA) for genome assisted assembly as well as FastQC for quality checks.

Whole-genome sequencing and genome assembly

RNA was extracted on an automated Chemagic 360 instrument, using the CMG-1049 kit (Perkin Elmer, Hamburg, Germany). The RNA was stored at -80°C prior to use. Libraries for whole genome sequencing were prepared using either the Oxford Nanopore Midnight protocol with Rapid Barcoding or the Illumina COVIDseq Assay.

Illumina Miseq/NextSeq

For the Illumina COVIDseq assay, the libraries were prepared according to the manufacturer's protocol. Briefly, amplicons were tagmented, followed by indexing using the Nextera UD Indexes Set A. Sequencing libraries were pooled, normalized to 4 nM and denatured with 0.2 N sodium acetate. A 8 pM sample library was spiked with 1% PhiX (PhiX Control v3 adaptor-ligated library used as a control). We sequenced libraries on a 500-cycle v2 MiSeq Reagent Kit on the Illumina MiSeq instrument (Illumina). On the Illumina NextSeq 550 instrument,

sequencing was performed using the Illumina COVIDSeq protocol (Illumina Inc, USA), an amplicon-based next-generation sequencing approach. The first strand synthesis was carried using random hexamers primers from Illumina and the synthesized cDNA underwent two separate multiplex PCR reactions. The pooled PCR amplified products were processed for fragmentation and adapter ligation using IDT for Illumina Nextera UD Indexes. Further enrichment and cleanup was performed as per protocols provided by the manufacturer (Illumina Inc). Pooled samples were quantified using Qubit 3.0 or 4.0 fluorometer (Invitrogen Inc.) using the Qubit dsDNA High Sensitivity assay according to manufacturer's instructions. The fragment sizes were analyzed using TapeStation 4200 (Invitrogen). The pooled libraries were further normalized to 4nM concentration and 25 µl of each normalized pool containing unique index adapter sets were combined in a new tube. The final library pool was denatured and neutralized with 0.2N sodium hydroxide and 200 mM Tris-HCL (pH7), respectively. 1.5 pM sample library was spiked with 2% PhiX. Libraries were loaded onto a 300-cycle NextSeq 500/550 HighOutput Kit v2 and run on the Illumina NextSeq 550 instrument (Illumina, San Diego, CA, USA).

Midnight Protocol

For Oxford Nanopore sequencing, the Midnight primer kit was used as described by Freed and Silander⁵⁵. cDNA synthesis was performed on the extracted RNA using LunaScript RT mastermix (New England BioLabs) followed by gene-specific multiplex PCR using the Midnight Primer pools which produce 1200bp amplicons which overlap to cover the 30-kb SARS-CoV-2 genome. Amplicons from each pool were pooled and used neat for barcoding with the Oxford Nanopore Rapid Barcoding kit as per the manufacturer's protocol. Barcoded samples were pooled and bead-purified. After the bead clean-up, the library was loaded on a prepared R9.4.1 flow-cell. A GridION X5 or MinION sequencing run was initiated using MinKNOW software with the base-call setting switched off.

Genome assembly

We assembled paired-end and nanopore .fastq reads using Genome Detective 1.132 (<https://www.genomedetective.com>) which was updated for the accurate assembly and variant calling of tiled primer amplicon Illumina or Oxford Nanopore reads, and the Coronavirus Typing Tool⁵⁶. For Illumina assembly, GATK HaploTypeCaller --min-pruning 0 argument was added to increase mutation calling sensitivity near sequencing gaps. For Nanopore, low coverage regions with poor alignment quality (<85% variant homogeneity) near sequencing/amplicon ends were masked to be robust against primer drop-out experienced in the Spike gene, and the sensitivity for detecting short inserts using a region-local global alignment of reads, was increased. In addition, we also used the wf_artic (ARTIC SARS-CoV-2) pipeline as built using the nextflow workflow framework⁵⁷. In some instances, mutations were confirmed visually with .bam files using Geneious software V2020.1.2 (Biomatters). The reference genome used throughout the assembly process was NC_045512.2 (numbering equivalent to MN908947.3).

Raw reads from the Illumina COVIDSeq protocol were assembled using the Exatype NGS SARS-CoV-2 pipeline v1.6.1, (<https://sars-cov-2.exatype.com/>). This pipeline performs quality control on reads and then maps the reads to a reference using Examap. The reference genome used throughout the assembly process was NC_045512.2 (Accession number: MN908947.3).

Several of the initial Ion Torrent genomes contained a number of frameshifts, which caused unknown variant calls. Manual inspection revealed that these were likely to be sequencing errors resulting in mis-assembled regions (likely due to the known error profile of Ion Torrent sequencers)⁵⁸. To resolve this, the raw reads from the IonTorrent platform were assembled using the SARSCoV2 RECoVERY (REconstruction of COronaVirus gEnomes & Rapid anaLYsis) pipeline implemented in the Galaxy instance ARIES (<https://aries.iss.it>). This pipeline fixed the observed frameshifts, confirming that they were artefacts of mis-assembly; this subsequently

resolved the variant calls. The Exatype and RECoVERY pipelines each produce a consensus sequence for each sample. These consensus sequences were manually inspected and polished using Aliview v1.27 (<http://ormbunkar.se/aliview/>).

All of the sequences were deposited in GISAID (<https://www.gisaid.org/>)^{17,18}, and the GISAID accession identifiers are included as part of **Supplementary Table S3**. Raw reads for our sequences have also been deposited at the NCBI Sequence Read Archive (BioProject accession PRJNA784038).

The number and position of the Omicron mutations has affected a number of primers and caused primer drop-outs across a range of sequencing protocols, especially within the RBD (<https://primer-monitor.neb.com/lineages>). These primer drop-outs have resulted in a number of genomes missing stretches of the RBD, and can affect estimates of mutation prevalence and the determination of the true set of lineage-defining mutations. Given this, bam files of all initial genomes were inspected with IG Viewer to confirm mutation calls where reference calls were suspected to be from low coverage at primer dropout sites⁵⁹.

Lineage classification

We used the widespread dynamic lineage classification method from the 'Phylogenetic Assignment of Named Global Outbreak Lineages' (PANGOLIN) software suite (<https://github.com/hCoV-2019/pangolin>)¹⁹. This is aimed at identifying the most epidemiologically important lineages of SARS-CoV-2 at the time of analysis, enabling researchers to monitor the epidemic in a particular geographic region. For the Omicron variant described in this study, the corresponding PANGO lineage designation is BA.1 (lineages v1.2.106). When first characterized the lineage was designated as B.1.1.529 but the emergence of three sibling lineages to Omicron resulted in the split into sub-lineages (B.1.1.529.1,

B.1.1.529.2 and B.1.1.529.3, aliased as BA.1, BA.2 and BA.3). BA.1 contains all the genomes with the original mutational constellation that was designated as Omicron and, at time of writing, is the dominant sub-lineage.

Recombination testing

To test for the possibility that the Omicron lineage is a recombinant of other SARS-CoV-2 lineages, we used a global subsample of sequences spanning January 2021 to August 2021. Using the NCBI SARS-CoV-2 Data hub^{60,61}, we constructed a dataset containing 221 sequences by randomly sampling five sequences from each month for each continent. No Oceania samples were available from July or August, and no South American sequences were available from July 2021⁶². These sequences were aligned together with a set of five high quality BA.1 and seven BA.2 sequences (representing the known diversity of these clades on 5 December 2021) using MAFFT⁶³ with default settings. Whereas 3SEQ³⁸, and RDP5³⁹ were used to analyse this dataset, a subsample of the 39 most divergent sequences from the dataset was analysed using the GARD recombination detection method³⁷. Default program settings were used throughout for recombination analyses, with the exception of RDP5 analysis, in which sequences were treated as linear and the window sizes for the SiScan and BootScan methods (two of the seven recombination detection methods applied in RDP5) were changed to 2000 nucleotides.

Selection analyses

We investigated the nature and extent of selective forces acting on BA.1 genes encoding individual protein products (a median of 25 unique BA.1 sequences per protein product encoding genome region). A subset of publicly available sequences (from the Virus Pathogen Database and Analysis Resource (ViPR) (<https://www.viprbrc.org/>)) were included as

background sequences to contextualize selection signals detectable within the BA.1 lineage at the levels of complete protein product encoding regions, and individual codons (a median of 106 sequences per protein coding region). Sequences were selected quality checked, aligned and subjected to BUSTED, RELAX, MEME, FADE, FEL, and BGM selection analyses (all implemented in HyPhy v2.5.31⁶⁴) using the automated RASCL pipeline as outlined previously^{2,9,35}.

Structure modeling

We modelled the spike protein on the basis of the Protein Data Bank coordinate set 7A94, showing the first step of the spike protein trimer activation with one RBD domain in the up position, bound to the human ACE2 receptor⁶⁵. We used the Pymol program (The PyMOL Molecular Graphics System, version 2.2.0) for visualization.

Phylogenetic analysis

All sequences on GISAID^{17,18} designated Omicron (n=686; date of access: 7 December 2021) were analyzed against a globally representative reference set of SARS-CoV-2 genotypes (n=12 609) spanning the entire genetic diversity observed since the start of the pandemic. In short, the reference set included: 1. All genomes from Africa assigned to PANGO lineage B.1.1 or any of its descendents, excluding those belonging to a VOC clade; 2. A representative subsampling of global data from the publicly maintained global build of Nexstrain (<https://nextstrain.org/ncov/gisaid/global>); 3. The top thirty BLAST hits when querying GISAID BLAST for BA.1 and BA.2 sequences. This sampling scheme ensures that we analyze Omicron against the closest variants of the virus. Omicron and reference sequences were aligned with Nextalign⁶⁶. A maximum-likelihood (ML) tree topology was inferred in FastTree⁶⁷ under the following parameters: a General Time Reversible (GTR) model of nucleotide substitution and a total of 100 bootstrap replicates⁶⁸. The resulting ML-tree topology was transformed into a time-

calibrated phylogeny where branches along the tree are scaled in calendar time using TreeTime⁶⁹. The resulting tree was then visualized and annotated in ggtree in R⁷⁰.

Time-calibrated BEAST analysis

To estimate a time-scale and growth rate from the genome sequence data, BEAST v1.10.4^{71,72} was used to sample phylogenetic trees under an exponential growth coalescent model using a strict molecular clock. All BA.1 assigned genomes from South Africa and Botswana (as of 11 December 2021) were included with some lower coverage genomes removed leaving a total of 553 genomes. The single South African BA.2 (CERI-KRISP-K032307, EPI_ISL_6795834) was included to help stabilize the root of the BA.1 clade but the exponential growth coalescent model was only applied to BA.1 (a constant population size coalescent was used for the rest of the tree). The rate of molecular evolution was estimated from the data. Two runs of 100 million iterations were compared to assess convergence and then post-burnin samples pooled to summarize parameter estimates.

Birth-death phylogenetic analysis

We analysed the full South Africa & Botswana dataset (n = 552) and the reduced dataset containing only Gauteng Province genomes (n = 277) using the serially sampled birth-death skyline (BDSKY) model⁷³, implemented in BEAST2 v2.5.2⁷⁴. To allow for changes in genomic sampling intensity shortly after the discovery of the new lineage, we allowed the sampling proportion to vary with time while keeping all other models parameters constant over the study period. The choice of prior distributions for the model parameters is summarised in Extended Data Table 3.

For each analysis, we used a strict clock model with a fixed clock rate of 7.5×10^{-4} substitution/site/year and a HKY substitution model. The mean duration of infectiousness was

fixed at 10 days^{75,76}. The effective reproductive number, R_e , was assumed to be constant with time. The sampling proportion, s , was assumed to be 0 before the collection time of the first sample (2021-11-04) and allowed to change at fixed times that were approximately equidistantly spaced between the first sample and the most recent sample (2021-12-05). The maximum clade credibility (MCC) tree generated from the analysis of the full South Africa and Botswana dataset with a Skygrid coalescent tree prior was used as the starting tree. We kept the subsequent tree topologies fixed such that the resulting MCMC chain only sampled internal node heights.

To assess the robustness of our estimates of R_e under different assumptions of temporal variations in the sampling proportion, we repeated the analyses with different numbers of equidistant change-time points (3 or 4). All other model parameters and priors were kept the same.

For each analysis, we ran two independent chains of 1×10^8 MCMC steps and sampled parameters every 10,000 steps. We used Tracer v1.7⁷⁷ to evaluate MCMC convergence for each of the individual chains ($ESS > 200$) which were then combined using LogCombiner to obtain the final posterior distribution after removing 10% of each chain as burn-in.

The resulting estimates for the doubling time and time of origin are summarised in Extended Data Table 4. The 3-epoch and 4-epoch BDSKY models resulted in similar estimates of the effective reproductive number, R_e , for both datasets: 2.74 (95% CI: 2.56-2.92) compared to 2.79 (95% CI: 2.60-2.99) for the South Africa & Botswana dataset, and 4.17 (95% CI: 3.76-4.60) compared to 3.85 (95% CI: 3.50-4.22) for the Gauteng Province dataset.

Phylogeographic analysis

Markov Chain Monte Carlo (MCMC) analyses were run in duplicate in BEAST v1.10.4^{71,72} for a total of 100 million iterations sampling every 10,000 steps in the chain. Convergence of runs was assessed in Tracer v1.7.1⁷⁷ based on high effective sample sizes (>200) and good mixing in the chains. Maximum clade credibility trees for each run were summarized in TreeAnnotator after discarding the first 10% of the chain as burn in. Finally, the spatiotemporal dispersal of Omicron was mapped using the R package “seraphim”⁷⁸.

Estimating transmission advantage

We analyzed 805 SARS-CoV-2 sequences from Gauteng, South Africa, that were uploaded to GISAID with sample collection dates from 1 September - 1 December 2021¹⁷. We used a multinomial logistic regression model to estimate the growth advantage of Omicron compared to Delta at the time point where the proportion of Omicron reached 50%^{79,80}. We fitted the model using the *multinom* function of the *nnet* package and estimated the growth advantage using the package *emmeans* in R.

The difference in the net growth rates (i.e., the growth advantage) between a variant (Omicron) and the wild-type (Delta) can be expressed as follows⁸¹:

$$\rho = (1 + \tau)\beta(S + \square(1 - S)) - \beta S,$$

where τ is the increase of the intrinsic transmissibility, \square is the level of immune evasion, β is the transmission rate of the wild-type, and S is the proportion of the population that is susceptible to the wild-type. This relation can be algebraically solved for τ and \square . We further define $R_w = \beta SD$ as the effective reproduction number of the wild-type with D being the generation time. $\Omega = 1 - S$ corresponds to the proportion of the population with protective immunity against infection and subsequent transmission with the wild-type.

We estimated \square for different levels of τ and Ω . To propagate the uncertainty, we constructed 95% credible intervals (CIs) of the estimates from 10,000 parameter samples of ρ , D , and R_w . We assumed D to be normally distributed with a mean of 5.2 days and a standard deviation of 0.8 days⁸². We sampled from publicly available estimates of the daily R_w based on confirmed cases during the early growth phase of Omicron in South Africa (1 October - 31 October 2021; range: 0.78-0.85 (<https://github.com/covid-19-Re>)⁵¹).

Data availability

All SARS-CoV-2 whole genome sequences produced by NGS-SA are deposited in the GISAID sequence database and are publicly available subject to the terms and conditions of the GISAID database. The GISAID accession numbers of sequences used in the phylogenetic analysis, including Omicron and global references, are provided in the **Supplementary Table S1**.

Code availability

All input files (e.g. alignments or XML files), along with all resulting output files and scripts used in the present study will be made available upon request and publicly shared on GitHub at final publication.

Methods references

49. Marivate, V. & Combrink, H. M. Use of Available Data To Inform The COVID-19 Outbreak in South Africa: A Case Study. *Data Sci. J.* **19**, (2020).
50. Marivate, V. *et al.* Coronavirus disease (COVID-19) case data - South Africa. *Zenodo* (2020) doi:10.5281/zenodo.3819126.
51. Huisman, J. S. *et al.* Estimation and worldwide monitoring of the effective reproductive number of SARS-CoV-2. *medRxiv* (2020) doi:10.1101/2020.11.26.20239368.

52. National Institute for Communicable Diseases. WEEKLY TESTING SUMMARY - NICD.
<https://www.nicd.ac.za/diseases-a-z-index/disease-index-covid-19/surveillance-reports/weekly-testing-summary/>.
53. Wickham, H. ggplot2. *WIREs Comp Stat* **3**, 180–185 (2011).
54. Msomi, N., Mlisana, K., de Oliveira, T. & Network for Genomic Surveillance in South Africa writing group. A genomics network established to respond rapidly to public health threats in South Africa. *Lancet Microbe* **1**, e229–e230 (2020).
55. SARS-CoV2 genome sequencing protocol (1200bp amplicon “midnight” primer set, using Nanopore Rapid kit). <https://dx.doi.org/10.17504/protocols.io.bwyppfvn>.
56. Cleemput, S. *et al.* Genome Detective Coronavirus Typing Tool for rapid identification and characterization of novel coronavirus genomes. *Bioinformatics* **36**, 3552–3555 (2020).
57. GitHub - epi2me-labs/wf-artic: ARTIC SARS-CoV-2 workflow and reporting.
<https://github.com/epi2me-labs/wf-artic#readme>.
58. Bragg, L. M., Stone, G., Butler, M. K., Hugenholtz, P. & Tyson, G. W. Shining a light on dark sequencing: characterising errors in Ion Torrent PGM data. *PLoS Comput. Biol.* **9**, e1003031 (2013).
59. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
60. Hatcher, E. L. *et al.* Virus Variation Resource - improved response to emergent viral outbreaks. *Nucleic Acids Res.* **45**, D482–D490 (2017).
61. National Library of Medicine. NCBI Virus: SARS-CoV-2 Data Hub.
https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/virus?SeqType_s=Nucleotide&VirusLineage_ss=SARS-CoV-2,%20taxid:2697049.
62. covid19-omicron-origins-recombination/aligned_234.shortnames.afa at main · bonilab/covid19-omicron-origins-recombination · GitHub.
<https://github.com/bonilab/covid19-omicron-origins-recombination/blob/main/4%20GS5%20plus%20Canada%20Outlier%20Lineage/4.2%20ali>

gned_mafft_addfrag_wref/aligned_234.shortnames.afa.

63. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002).
64. Kosakovsky Pond, S. L. *et al.* HyPhy 2.5-A Customizable Platform for Evolutionary Hypothesis Testing Using Phylogenies. *Mol. Biol. Evol.* **37**, 295–299 (2020).
65. Benton, D. J. *et al.* Receptor binding and priming of the spike protein of SARS-CoV-2 for membrane fusion. *Nature* **588**, 327–330 (2020).
66. Aksamentov, I., Roemer, C., Hodcroft, E. & Neher, R. Nextclade: clade assignment, mutation calling and quality control for viral genomes. *JOSS* **6**, 3773 (2021).
67. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 — approximately maximum-likelihood trees for large alignments. *PLoS ONE* **5**, e9490 (2010).
68. Felsenstein, J. Confidence limits on phylogenies: an approach using the bootstrap. *Evolution* **39**, 783–791 (1985).
69. Sagulenko, P., Puller, V. & Neher, R. A. TreeTime: Maximum-likelihood phylodynamic analysis. *Virus Evol.* **4**, vex042 (2018).
70. Yu, G. Using ggtree to Visualize Data on Tree-Like Structures. *Curr. Protoc. Bioinformatics* **69**, e96 (2020).
71. Suchard, M. A. *et al.* Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol.* **4**, vey016 (2018).
72. Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969–1973 (2012).
73. Stadler, T., Kühnert, D., Bonhoeffer, S. & Drummond, A. J. Birth-death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc Natl Acad Sci USA* **110**, 228–233 (2013).
74. Bouckaert, R. *et al.* BEAST 2.5: An advanced software platform for Bayesian evolutionary

- analysis. *PLoS Comput. Biol.* **15**, e1006650 (2019).
75. Benvenuto, D. *et al.* The global spread of 2019-nCoV: a molecular evolutionary analysis. *Pathog. Glob. Health* **114**, 64–67 (2020).
 76. Byrne, A. W. *et al.* Inferred duration of infectious period of SARS-CoV-2: rapid scoping review and analysis of available evidence for asymptomatic and symptomatic COVID-19 cases. *BMJ Open* **10**, e039856 (2020).
 77. Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior summarization in bayesian phylogenetics using tracer 1.7. *Syst. Biol.* **67**, 901–904 (2018).
 78. Dellicour, S., Rose, R., Faria, N. R., Lemey, P. & Pybus, O. G. SERAPHIM: studying environmental rasters and phylogenetically informed movements. *Bioinformatics* **32**, 3204–3206 (2016).
 79. Davies, N. G. *et al.* Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England. *Science* **372**, (2021).
 80. Campbell, F. *et al.* Increased transmissibility and global spread of SARS-CoV-2 variants of concern as at June 2021. *Euro Surveill.* **26**, (2021).
 81. Althaus, C. L. *et al.* A tale of two variants: Spread of SARS-CoV-2 variants Alpha in Geneva, Switzerland, and Beta in South Africa. *medRxiv* (2021)
doi:10.1101/2021.06.10.21258468.
 82. Ganyani, T. *et al.* Estimating the generation interval for coronavirus disease (COVID-19) based on symptom onset data, March 2020. *Euro Surveill.* **25**, (2020).
 83. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
 84. Boni, M. F. *et al.* Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the COVID-19 pandemic. *Nat. Microbiol.* **5**, 1408–1417 (2020).

Acknowledgements

We thank Linda de Gouveia, Amelia Buys, Cardia Fourie, Noluthando Duma, Malusi Ndlovu and other members of the NICD Centre for Respiratory Diseases and Meningitis and Sequencing Core Facility. We thank Nevashan Govender, Genevieve Ntshoe, Andronica Moipone Shonhiwa, Darren Muganhiri, Itumeleng Matiea, Eva Mathatha, Fhatuwani Gavhi, Teresa Mashudu Lamola, Matimba Makhubele, Mmaborwa Matjokotja, Simbulele Mdleleni, Masingita Makhubela from the national SARS-CoV-2 NICD surveillance team for NMCSS case data, and Fazil Mckenna, Trevor Graham Bell, Ndivhuwo Munava, Stanford Kwenda, Muzammil Raza Bano and Jimmy Khosa from NICD IT for NMCSS case and test data (in particular, SGTF data). We also thank the following people from the diagnostic laboratories for their assistance: Kubendran Reddy, Lilishia Gounder and Cherise Naicker from NHLS Inkosi Albert Luthuli Central Hospital Laboratory; Stephen Korsman from NHLS Groote Schuur Laboratory; and Annabel Enoch at NHLS Green Point Laboratory. Equally, we thank the global laboratories that generated and made public the SARS-CoV-2 sequences (through GISAID) used as reference dataset in this study (a complete list of individual contributors of sequences is provided in Supplementary Table S3).

The research reported in this publication was supported by the Strategic Health Innovation Partnerships Unit of the South African Medical Research Council, with funds received from the South African Department of Science and Innovation. CA received funding from the European Union's Horizon 2020 research and innovation programme - project EpiPose (No 101003688). DPM was funded by the Wellcome Trust (222574/Z/21/Z). RC & AR acknowledge support from the Wellcome Trust (Collaborators Award 206298/Z/17/Z - ARTIC network) and AR from the European Research Council (grant agreement number 725422 – ReservoirDOCS). VH was supported by the Biotechnology and Biological Sciences Research Council (BBSRC) (grant number BB/M010996/1). AEZ, JT, MUGK, OGP acknowledge support from the Oxford Martin

School. MUGK acknowledges support from the Rockefeller Foundation, Google.org, and the European Horizon 2020 programme MOOD (#874850). MV and the ZARV members, UP was funded through the ANDEMIA G7 Global Health Concept: contributions to improvement of International Health, COVID19 funds through the Robert Koch Institute. The genomic sequencing at UCT/NHLS is funded from the South African Medical Research Council and Department of Science and Innovation; and by the Wellcome Centre for Infectious Diseases Research in Africa (CIDRI-Africa) which is supported by core funding from the Wellcome Trust [203135/Z/16/Z and 222754]. CW and JNB are funded by the EDCTP (RADIATES Consortium; RIA2020EF-3030). Sequencing activities at the NICD were supported by: a conditional grant from the South African National Department of Health as part of the emergency COVID-19 response; a cooperative agreement between the National Institute for Communicable Diseases of the National Health Laboratory Service and the United States Centers for Disease Control and Prevention (grant number 5 U01IP001048-05-00); the African Society of Laboratory Medicine (ASLM) and Africa Centers for Disease Control and Prevention through a sub-award from the Bill and Melinda Gates Foundation grant number INV-018978; the UK Foreign, Commonwealth and Development Office and Wellcome (Grant no 221003/Z/20/Z); the South African Medical Research Council (Reference number SHIPNCD 76756); the UK Department of Health and Social Care, managed by the Fleming Fund and performed under the auspices of the SEQAFRICA project. The genomic sequencing in Botswana was supported by the Foundation for Innovative New Diagnostics and Fogarty International Center (5D43TW009610), NIH (5K24AI131924-04; 5K24AI131928-05), as well in kind support from the Botswana government through the Ministry of Health & Wellness and Presidential COVID-19 Task Force. SM was supported in part by the Bill & Melinda Gates Foundation [036530]. Under the grant conditions of the Foundation, a Creative Commons Attribution 4.0 Generic License has already been assigned to the Author Accepted Manuscript version that might arise from this submission

Author Contributions

Genomic data generation: RV, SM, DGA, HT, CS, JG, JE, SG, WTC, DM, BZ, BR, LK, RS, SL, MBM, PS, MM, MM, KM, AM, AI, BM, MSM, JES, NN, GM, SP, TM, UR, YN, CW, SE, TM, WP, LS, UJA, MM, SvW, DT, KD, DH, KM, DD, RJ, AI, DG, PAB, MMN, PNM, JNB;

Sample collection and metadata curation: RV, SM, DGA, AM, AS, MD, SM, WTC, DM, PS, MC, CJ, LK, OL, KM, NT, NH, NM, KM, AS, AM, MD, ZM, OL, YR, AM, KS, DG, PAB, FT, MV

Data analysis: HT, CS, RJL, NW, JE, AR, CA, EW, CKW, DPM, VH, RC, JES, MG, SP, AGL, SW, MFB, AEZ, JT, MUGK, OGP

Study design and data interpretation: RV, SM, DGA, RJL, AR, CA, SG, MM, MM, KM, LK, OL, MSM, KM, CW, OGP, AG, FT, MV, JNB, AvG, TdO

Manuscript writing: SM, HT, RJL, JG, JE, AR, CA, EW, DPM, JNB, AvG, TdO

All authors reviewed the manuscript

Competing interests statement

The authors declare no competing interests

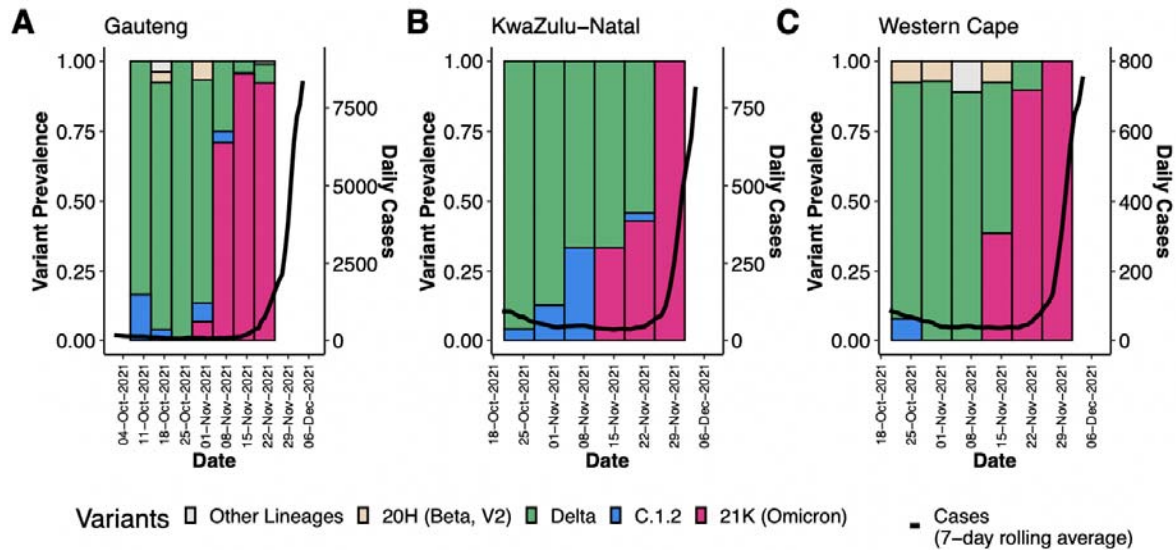
Additional information

Supplementary Information is available for this paper

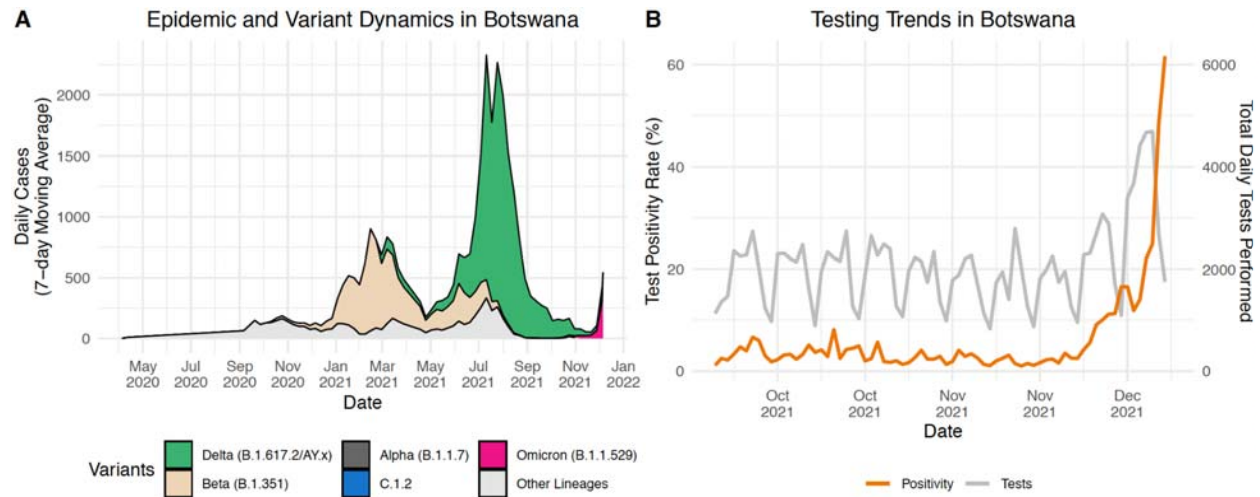
Correspondence and requests for materials should be addressed to Professor Tulio de Oliveira, Centre for Epidemic Response and Innovation (CERI), School of Data Science and Computational Thinking, Stellenbosch University, Stellenbosch, South Africa, tulio@sun.ac.za

Reprints and permissions information is available at www.nature.com/reprints

Extended Data Figures



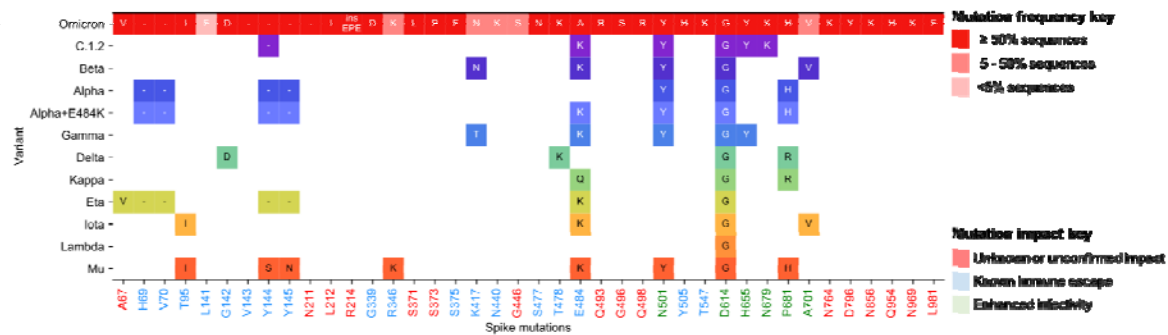
Extended Data Figure 1: Progression of daily recorded cases and variant proportions in Gauteng (A), KwaZulu-Natal (B) and Western Cape (C) provinces between October and December 2021. A sharp increase in 7-day rolling average of the number of cases is observed in all three of the biggest provinces in South Africa at the emergence of the Omicron variant.



Extended Data Figure 2: Epidemic Progression in Botswana. A) Epidemic and variant dynamics in Botswana from May 2020 to December 2021, with 7-day rolling average of the number of recorded cases coloured by the proportion of variants as inferred by genomic surveillance data available on GISAID. At the end of November 2021, a big Delta-driven wave was coming to the end and an Omicron wave was starting at the end of November 2021. B) Trends in testing numbers and positivity rates in Botswana between October and December 2021, showing a sharp increase in positivity rate mid-November 2021.

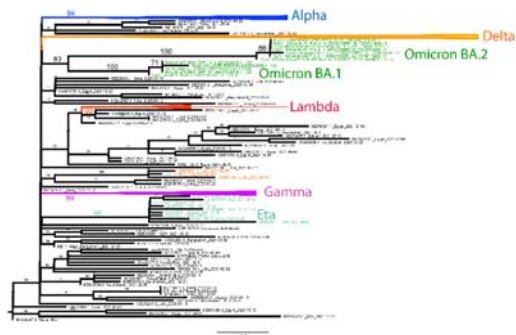
tracing or other reasons), routine surveillance or unknown if no information has been provided.

Countries are ordered by number of sequences available on GISAID as of December 16, 2021.

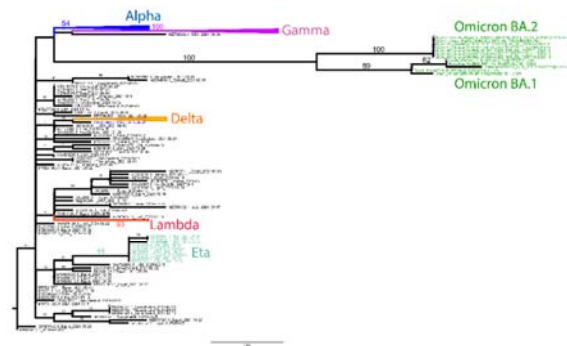


Extended Data Figure 4: Omicron/BA.1 spike mutations shared with other VOC/VOIs. All spike mutations seen in Omicron/BA.1 are listed at the top in red and coloured according to prevalence. Prevalence was calculated by number of mutation detections / total number of sequences. However, primer drop-outs have affected the RBD region spanning K417N, N440K and G446S, and so it is likely that these mutations may actually be more prevalent than indicated here. For the VOC/VOIs only mutations that are shared with Omicron and seen in $\geq 50\%$ of the respective VOC/VOI sequences are shown and are coloured according to Nextstrain clade. The mutations listed at the bottom are shaded according to known immune escape (blue), enhanced infectivity (green) or for unknown/unconfirmed impact (red).

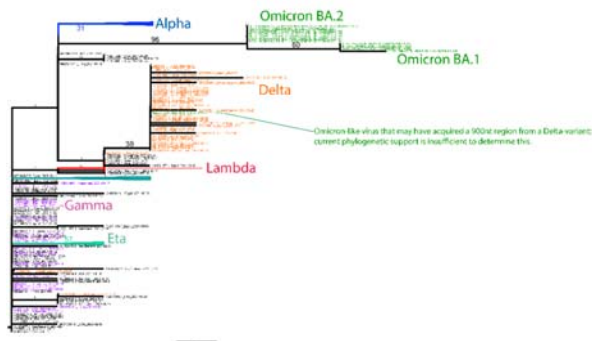
Breakpoint-free region, positions 1450 - 20520



Breakpoint-free region, positions 21619 - 23609



Breakpoint-free region, positions 23614 - 24513



Extended Data Figure 5. Maximum-likelihood trees (inferred with RAxML v8.2.12⁸³) for three key breakpoint-free regions (BFR)⁸⁴ of a SARS-CoV-2 genome alignment including global sequences from 2021 (N=221) and 12 sequences of the BA.1/Omicron and BA.2 lineages. All three regions show Omicron as an independently-evolved lineage (bootstrap ≥ 83) with no indication of any regions showing closer/further clustering with non-Omicron lineages. A 900nt segment (in the S2 region) of Botswana/R43B66_BHP_521004487/2021 may have been inserted into this virus through a recombination with a Delta-variant virus, however the clustering of Botswana/R43B66 with the Delta lineage is supported by a low bootstrap value of 38 due to

the short length of this region. Two other candidate BFRs -- one in S2 (positions 26159-27269) and one at the 5' end of the genome (positions 1-1060) -- showed little genetic diversity and no phylogenetic evidence of recombination. Regions shorter than 900nt had too little genetic diversity to show evidence of phylogenetic clustering. Bootstrap values are shown on branches with relevant values magnified for readability. All trees were rooted on the Wuhan-Hu-1 sequence.

Extended Data Table 1. Parameter estimates from BEAST for the full South Africa & Botswana data set and the reduced data set of only Gauteng Province genomes. 95% Highest Posterior Density (HPD) intervals in parentheses.

| Data set | Evolutionary rate $\times 10^{-3}$ changes/site/year | BA.1 Time of most recent common ancestor (TMRCA) | Exponential growth rate /day | Doubling time days |
|---|--|---|------------------------------------|-----------------------|
| South Africa + Botswana 553 Genomes | 1.20 (0.92, 1.49) | 9 Oct 2021 (30 Sep, 20 Oct) | 0.137 (0.100, 0.174) | 5.1 (3.8, 7.0) |
| South Africa + Botswana 553 Genomes | 0.75 fixed | 1 Oct 2021 (21 Sep, 13 Oct) | 0.139 (0.099, 0.183) | 5.0 (4.0, 7.0) |
| Gauteng Province, South Africa only 279 genomes | 0.30 (0.08, 0.5) | 16 Sep 2021 (15 Jul, 21 Oct) | 0.222 (0.048, 0.396) | 3.1 (1.7, 14.5) |
| Gauteng Province, South Africa only 279 genomes | 1.1 fixed | 21 Oct 2021 (16 Oct, 28 Oct) | 0.384 (0.217, 0.562) | 1.8 (1.2, 3.2) |

Extended Data Table 2. Sites in the Omicron sequences that have been subject to episodic diversifying selection

| Coordinate (SARS-CoV-2) | Gene/ORF | Codon (in gene/ORF) | # of selected branches | AA composition | p-value | Notes |
|----------------------------|----------|------------------------|---------------------------|------------------------|---------|-------------------|
| 1459 | ORF1a | 399 | 1 | K/33, N/1 | 0.0003 | |
| 2092 | ORF1a | 610 | 1 | S/33, L/1 | 0.0013 | |
| 2674 | ORF1a | 804 | 1 | P/32, L/1, -/1 | 0.0003 | |
| 8833 | ORF1a | 2857 | 1 | V/22, -/1, A/1 | 0.0085 | |
| 10447 | ORF1a | 3395 | 1 | H/15, P/2, Y/1 | 0.0009 | |
| 10636 | ORF1a | 3458 | 1 | G/15, -/1, F/1, V/1 | 0.0012 | |
| 17673 | ORF1b | 1402 | 1 | I/29, -/1, V/1 | 0.0060 | |
| 18264 | ORF1b | 1599 | 1 | R/20, -/1, Q/1, E/1 | 0.0001 | |
| 18267 | ORF1b | 1600 | 1 | E/20, -/2, T/1 | 0.0005 | |
| 18270 | ORF1b | 1601 | 1 | E/20, -/2, C/1 | 0.0001 | |
| 18273 | ORF1b | 1602 | 1 | A/20, -/2, C/1 | 0.0000 | |
| 21033 | ORF1b | 2522 | 1 | L/10, F/2 | 0.0001 | |
| 21844 | S | 95 | 2 | I/82, T/3, -/3 | 0.0028 | Clade defining |
| 22576 | S | 339 | 5 | D/77, G/7, -/4 | 0.0000 | Clade |

| | | | | | | |
|-------|---|-----|---|---------------------|--------|---------------------------|
| | | | | | | defining |
| 22597 | S | 346 | 2 | R/75, K/11, - /2 | 0.0062 | Affects Ab binding |
| 22672 | S | 371 | 2 | L/75, -/7, S/6 | 0.0000 | Clade defining |
| 22684 | S | 375 | 2 | F/75, S/8, -/5 | 0.0061 | Clade defining |
| 22810 | S | 417 | 2 | N/58, -/18, K/12 | 0.0037 | Clade defining |
| 22879 | S | 440 | 3 | K/54, -/27, N/7 | 0.0006 | Clade defining |
| 22897 | S | 446 | 2 | S/56, -/24, G/8 | 0.0038 | Clade defining |
| 22915 | S | 452 | 1 | L/59, -/25, R/4 | 0.0003 | Affects Ab binding |
| 23662 | S | 701 | 2 | A/81, V/6, -/1 | 0.0097 | 501Y Metasignatur e |
| 23851 | S | 764 | 2 | K/75, -/7, N/6 | 0.0083 | Clade defining |
| 23947 | S | 796 | 2 | Y/82, D/5, -/1 | 0.0044 | Clade defining |
| 24127 | S | 856 | 2 | K/82, N/5, -/1 | 0.0006 | Clade defining |
| 24421 | S | 954 | 3 | H/81, Q/7 | 0.0058 | Clade |

| | | | | | | |
|-------|------|-----|---|----------------------|--------|-------------------|
| | | | | | | defining |
| 24466 | S | 969 | 2 | K/81, N/6, -/1 | 0.0045 | Clade defining |
| 24502 | S | 981 | 2 | F/79, L/8, -/1 | 0.0048 | Clade defining |
| 28250 | ORF8 | 120 | 2 | L/45, F/35, - /12 | 0.0011 | |
| 28369 | N | 33 | 0 | -/36, S/6, G/1 | 0.0058 | |
| 28459 | N | 63 | 2 | D/31, G/6, -/6 | 0.0000 | |
| 29299 | N | 343 | 1 | D/39, G/3, -/1 | 0.0059 | |

Extended Data Table 3 Prior distributions used for the BDSKY analyses. The becoming non-infectious rate was fixed to 36.5/day which corresponds to an infectious period of 10 days.

| Parameter | Prior distribution |
|--|---|
| clock rate | Fixed at 0.00075 substitution/site/year |
| kappa | lognormal(M = 1, s = 1.25) |
| gamma shape | Exponential(m = 1) |
| effective reproductive number, R_e | lognormal(M=0.8, s=0.5) |
| becoming non-infectious rate, δ | Fixed at 36.5/day |
| sampling proportion, s | beta(alpha = 2, beta = 1000) |
| time of origin, t_{Origin} | lognormal(M = -2, s = 0.2) |

Extended Data Table 4 Doubling time and time of origin estimates for the full South Africa & Botswana dataset and the reduced dataset of only Gauteng Province genomes under two different assumptions of temporal variation in sampling proportion.

| | doubling time, \square_d | | time of origin, t_{Origin} | |
|---|----------------------------|----------------------|--|--|
| | 3-epoch | 4-epoch | 3-epoch | 4-epoch |
| South Africa + Botswana BA.1 (n=552) | 3.99 [3.61, 4.45] | 3.88 [3.47, 4.30] | 2021-10-08 [2021-09-27, 2021-10-16] | 2021-10-07 [2021-09-26, 2021-10-15] |
| Gauteng Province only BA.1 (n=279) | 2.19 [1.90, 2.48] | 2.43 [2.15, 2.77] | 2021-10-12 [2021-10-04, 2021-10-18] | 2021-10-10 [2021-10-02, 2021-10-17] |