






SARS-CoV-2 Variants Associated with Vaccine Breakthrough in the Delaware Valley through Summer 2021

Andrew D. Marques,^a Scott Sherrill-Mix,^a John K. Everett,^a Shantanu Reddy,^a Pascha Hokama,^a Aoife M. Roche,^a Young Hwang,^a Abigail Glascock,^a Samantha A. Whiteside,^b Jevon Graham-Wooten,^b Layla A. Khatib,^b Ayannah S. Fitzgerald,^b Ahmed M. Moustafa,^{c,d} Colleen Bianco,^c Swetha Rajagopal,^c Jenna Helton,^e Regan Deming,^c Lidiya Denu,^c Azad Ahmed,^f Eimear Kitt,^{c,g} Susan E. Coffin,^{c,g} Claire Newbern,^e Josh Chang Mell,^f Paul J. Planet,^{c,g,h} Nitika Badjatia,ⁱ Bonnie Richards,^j Zi-Xuan Wang,^{i,k} Carolyn C. Cannuscio,^{l,m} Katherine M. Strelau,^{l,m} Anne Jaskowiak-Barr,ⁿ Leigh Cressman,ⁿ Sean Loughrey,ⁿ Arupa Ganguly,^o Michael D. Feldman,^p  Ronald G. Collman,^b Kyle G. Rodino,^p  Brendan J. Kelly,ⁿ  Frederic D. Bushman^a

^aDepartment of Microbiology, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

^bPulmonary, Allergy and Critical Care Division, Department of Medicine, University of Pennsylvania Perelman School of Medicine, Philadelphia, Pennsylvania, USA

^cDivision of Infectious Diseases, Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA

^dDivision of Gastroenterology, Hepatology & Nutrition, Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA

^eDivision of COVID-19 Containment, Philadelphia Department of Public Health, Philadelphia, Pennsylvania, USA

^fDepartment of Microbiology & Immunology, Center for Genomic Sciences, Drexel University College of Medicine, Philadelphia, Pennsylvania, USA

^gDepartment of Pediatrics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

^hSackler Institute for Comparative Genomics, American Museum of Natural History, New York, New York, USA

ⁱMolecular & Genomic Pathology Laboratory, Thomas Jefferson University Hospital, Philadelphia, Pennsylvania, USA

^jJefferson Occupational Health Network for Employees and Students (JOHN), Thomas Jefferson University, Philadelphia, Pennsylvania, USA

^kDepartment of Anatomy, Pathology, and Cell Biology, Thomas Jefferson University Hospital, Philadelphia, Pennsylvania, USA

^lLeonard Davis Institute of Health Economics, University of Pennsylvania, Philadelphia, Pennsylvania, USA

^mDepartment of Family Medicine and Community Health, University of Pennsylvania, Philadelphia, Pennsylvania, USA

ⁿDivision of Infectious Diseases, Department of Medicine & Department of Biostatistics, Epidemiology, and Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

^oDepartment of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

^pDepartment of Pathology and Laboratory Medicine, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

Andrew D. Marques, Scott Sherrill-Mix, and John K. Everett co-first authors. A.D.M. carried out both biochemical and bioinformatic analysis, S.S.-M. devised the statistical model, and J.K.E. carried out software engineering in support of the project.

ABSTRACT The severe acute respiratory coronavirus-2 (SARS-CoV-2) is the cause of the global outbreak of COVID-19. Evidence suggests that the virus is evolving to allow efficient spread through the human population, including vaccinated individuals. Here, we report a study of viral variants from surveillance of the Delaware Valley, including the city of Philadelphia, and variants infecting vaccinated subjects. We sequenced and analyzed complete viral genomes from 2621 surveillance samples from March 2020 to September 2021 and compared them to genome sequences from 159 vaccine breakthroughs. In the early spring of 2020, all detected variants were of the B.1 and closely related lineages. A mixture of lineages followed, notably including B.1.243 followed by B.1.1.7 (alpha), with other lineages present at lower levels. Later isolations were dominated by B.1.617.2 (delta) and other delta lineages; delta was the exclusive variant present by the last time sampled. To investigate whether any variants appeared preferentially in vaccine breakthroughs, we devised a model based on Bayesian autoregressive moving average logistic multinomial regression to allow rigorous comparison. This revealed that B.1.617.2 (delta) showed 3-fold enrichment in vaccine breakthrough cases (odds ratio of 3; 95% credible interval 0.89-11). Viral point substitutions could also be associated with vaccine breakthroughs, notably the N501Y substitution found in the alpha, beta and gamma variants (odds ratio 2.04; 95% credible interval of 1.25-3.18). This study thus overviews viral evolution and vaccine

Editor Stephen P. Goff, Columbia University/HHMI

Copyright © 2022 Marques et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Address correspondence to Frederic D. Bushman, bushman@penmedicine.upenn.edu, Brendan J. Kelly, brendank@penmedicine.upenn.edu, or Kyle G. Rodino, Kyle.Rodino@Penmedicine.upenn.edu.

The authors declare a conflict of interest. R.G.C. reports support to his lab for COVID-19 work unrelated to the current manuscript from OraSure, Inc, and ongoing collaborations with Resilient Biotics, Inc. All other authors have no competing interests.

This article is a direct contribution from Frederic D. Bushman, a Fellow of the American Academy of Microbiology, who arranged for and secured reviews by Angela Ciuffi, Institute of Microbiology, University Hospital Center and University of Lausanne, and Manuel Llano, University of Texas at El Paso.

Received 21 December 2021

Accepted 12 January 2022

Published 8 February 2022

breakthroughs in the Delaware Valley and introduces a rigorous statistical approach to interrogating enrichment of breakthrough variants against a changing background.

IMPORTANCE SARS-CoV-2 vaccination is highly effective at reducing viral infection, hospitalization and death. However, vaccine breakthrough infections have been widely observed, raising the question of whether particular viral variants or viral mutations are associated with breakthrough. Here, we report analysis of 2621 surveillance isolates from people diagnosed with COVID-19 in the Delaware Valley in southeastern Pennsylvania, allowing rigorous comparison to 159 vaccine breakthrough case specimens. Our best estimate is a 3-fold enrichment for some lineages of delta among breakthroughs, and enrichment of a notable spike substitution, N501Y. We introduce statistical methods that should be widely useful for evaluating vaccine breakthroughs and other viral phenotypes.

KEYWORDS SARS-CoV-2, COVID-19, coronavirus, genome sequencing, Philadelphia

The global COVID-19 pandemic is caused by infection with the virus SARS-CoV-2 (1). Analysis of whole-genome sequences from global viral samples shows ongoing changes in the composition of viral populations. RNA viruses have high mutation rates, so sequence change is expected in the viral genome due to random genetic drift (2). However, selection for efficient immune evasion and transmission between humans now seem likely to be major drivers of SARS-CoV-2 diversification (3–5).

Widespread vaccination against SARS-CoV-2 was introduced in the United States in the winter of 2020–2021. First to be implemented were vaccines based on modified mRNAs, followed by adenovirus vector delivery. Vaccines are highly protective against infection, severe disease, and death. However, infection of vaccinated individuals has been widely detected, albeit typically with much milder disease course compared to that experienced by unvaccinated individuals (6–8). Thus, interest turns to the question of which viral features are associated with vaccine breakthrough infections (7–9).

Several criteria can be applied to assessing whether sequence changes in a new variant have likely evolved to promote infection and vaccine breakthrough. Substitutions found in viral spike proteins can be tested in laboratory experiments to determine whether they promote more efficient replication in human cells or reduce binding of human antibodies (10–18). Other substitutions may alter epitopes targeted by the cellular immune system (6, 19, 20). Viral lineages with diverse combinations of these substitution have been identified and designated variants being monitored and variants of concern (VBM/VOC) by the Centers for Disease Control and Prevention (CDC). Some of the substitutions in these variants have been detected arising independently on multiple genetic backgrounds, such as the spike substitutions N501Y, E484K or the 69–70 deletion (4, 21–24), supporting a model of convergent evolution.

One indication of selection for increased transmission in humans is that several new variants have spread globally and rapidly displaced preexisting strains. This was first documented for the D614G substitution, which spread around the world in the Spring of 2020 and displaced most strains lacking this substitution (25–28). More recently, variants first identified in the UK (B.1.1.7 or alpha) (29, 30), South Africa (B.1.351 or beta), Brazil (P.1 or gamma), California (B.1.427 and B.1.429 or epsilon), New York (B.1.526 or iota) (31), and India (B.1.617.2 or delta) (8) have been suggested to be spreading at the expense of preexisting viral types. Against this background, intense interest focuses on whether particular variants are more efficient at infecting vaccinated individuals.

We have investigated viral genomic evolution in the Delaware Valley, which encompasses the city of Philadelphia, in a sample that includes vaccine breakthrough cases. Our initial report on the first wave of infection in this area (32) revealed that lineages in Philadelphia most closely matched sequences derived from New York City, approximately 100 miles away, which is larger than Philadelphia and had an earlier peak in infection. We also found that in some cases different viral sequence polymorphisms could be found in the same patients from different body sites or in longitudinal samples, suggestive of ongoing evolution within infected individuals (32).

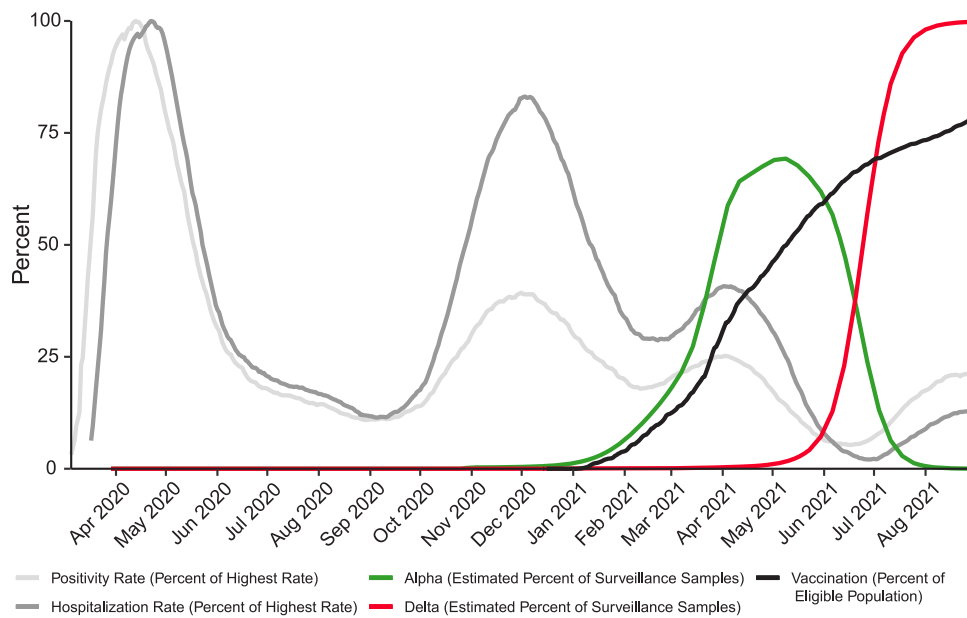


FIG 1 Longitudinal data from the COVID-19 pandemic in the city of Philadelphia. The y axis shows the daily test positivity rate (light gray) as a percentage of the highest value (26.57% positivity on 4/13/2020), the hospitalization rate (dark gray) as a percentage of the highest value (87 hospitalizations per day on 4/22/2020), the vaccination rate in adults 18 years old and older (black). The estimated percentage of surveillance samples classified as alpha (green) or delta (red) variant was estimated from the sequence data presented in this paper. Other data are from the city of Philadelphia “Testing Data: Programs and Initiatives.”

In this study, samples were collected from 2621 surveillance samples from infected individuals and 159 vaccine breakthrough cases through September 2021, and the representation of different variants compared. We found that several waves of variants rose and fell in prevalence over the course of our sampling period, by the end of sampling all genomes were identified as VOC delta lineages, and the delta lineage B.1.617.2 was potentially 3-fold enriched among vaccine breakthrough samples compared to surveillance samples. The amino acid substitution N501Y, found in the alpha, beta, and gamma variants, was also notably enriched in vaccine breakthrough samples. We introduce a rigorous statistical approach, Bayesian autoregressive moving average logistic multinomial regression, that should be widely useful for assessing enrichment of viral variants while controlling for the changing background of circulating strains.

RESULTS

The COVID-19 epidemic in the Delaware Valley. Sampling was carried out from March 2020 to September 2021. Over the course of the study, several waves of infection are evident as increased test positivity rates (Fig. 1, light gray curve) and COVID-19-attributed hospitalizations in the city of Philadelphia (Fig. 1, dark gray curve). Widespread vaccination was introduced in winter 2020–2021, reaching over 70% (in adults 18 years old and older) in Philadelphia by September 2021 (Fig. 1, black curve). As is described below, during this period, SARS-CoV-2 variants alpha and delta, designated VBM/VOC, became the overwhelming majority of genomes identified by our sequence surveys (Fig. 1, green and red curves).

Patient populations. Patient samples (nasal or nasopharyngeal swab and saliva) used for SARS-CoV-2 whole-genome sequence analysis are listed in Table S1. To preserve patient confidentiality, samples were deidentified except for collection date and rationale for collection. Surveillance samples ($n = 2621$) were defined as those acquired from clinical diagnostic laboratories across the Delaware Valley ($n = 2485$), from hospitalized subjects ($n = 116$), and asymptomatic subjects testing positive in a university screening program ($n = 20$) (33). Vaccine breakthrough cases ($n = 159$) were included only if detected at least 2 weeks after the final vaccine dose (second dose for Moderna and Pfizer-BioNTech mRNA vaccines or single dose for the Johnson & Johnson adenovirus vector platform) and the

subject tested positive by clinical laboratory assay. Data were not available on subject immune responses following vaccination.

As a positive control, *S* gene target failure samples ($n = 172$) were also compared. The TaqPath COVID-19 Combo kit by Thermo Fisher Scientific targets three regions of the SARS-CoV-2 genome for viral detection; the ORF1ab, nucleocapsid (N), and spike (S) genes. The region of the *S* gene interrogated by the assay overlies a characteristic deletion in the alpha variant (del 69–70), so samples containing alpha lineage virus are selectively negative for the spike amplicon, while the other two targets are detected. Samples with these characteristics were targeted for sequencing early during the wave of alpha infections to track the variant; here these spike target gene failures serve as positive controls for the statistical model.

Sequencing strategies. Several sequencing strategies were used to acquire whole-genome sequences. The POLAR protocol with ARTIC primers and Illumina sequencing was used for most samples (34). Smaller numbers of samples were acquired using the Paragon, Illumina RPIP, and Illumina CovidSeq methods. Viral genome sequences were judged to be high quality and included in the study if 95% of the viral genome was covered by at least 5 reads. In all, 2952 high-quality sequences were generated and analyzed. Viral variants were assigned using Pangolin lineage software (35, 36).

Variants detected in pooled surveillance data. Figure 2A shows the proportions of variants detected in surveillance samples over the course of the study from March 2020 to September 2021. Variants detected are summarized by the color code on the bar graph (variant designations used are summarized in Table S2). The numbers of genomes analyzed per week are indicated above each column. Numbers varied both as surveillance sequencing efforts accelerated and as the availability of samples varied.

Members of the B.1 lineage predominated March 2020 until fall 2020, at which time B.1.2 and B.1.243 became predominant. Later the B.1.526 and B.1.1.7 (alpha) variants emerged, with B.1.1.7 becoming predominant by spring 2021. Delta (B.1.617.2 and AY lineages) became detectable in late spring and predominant in early summer. By late summer, delta lineages were the only variants detected.

Assessing variant abundance in the surveillance population. Wide-spread vaccination was introduced in midwinter 2020–2021, but nevertheless breakthrough infections were detected in some vaccinated individuals, raising the question of whether specific viral features identifiable in sequence data might be associated with vaccine evasion. We sequenced 159 vaccine breakthrough samples collected between February 22 and September 3, 2021, and compared the distributions of lineages or genomic variations to those observed in surveillance samples from the same community ($n = 2621$).

Challenges in the analysis include the facts that: (i) the distribution of variants in the surveillance samples is changing over time; (ii) sampling is uneven over time; and (iii) sampling is subject to stochastic fluctuations. To address these issues, we developed a model based on Bayesian analysis that combines an autoregressive moving average model with multinomial logistic regression. The underlying variant proportions at each week of the study were estimated from the counts of SARS-CoV-2 variants assuming relatively smooth changes over time and accounting for stochastic fluctuations during sampling (Fig. 2B). These estimated surveillance proportions were compared to the variant counts observed in spike gene target failure samples (Fig. 2C) and vaccine breakthrough samples (Fig. 2D). Since the surveillance population proportions were estimated for each week, the weekly vaccine breakthrough or spike gene target failure variant counts could be compared to the corresponding time-matched surveillance estimates.

Examples of the temporal profiles of the modeled lineage succession are shown in more detail in Fig. 3. One benefit of the model is that uncertainty in the surveillance estimates can be included. For example, note that time points with fewer surveillance samples (lighter gray bars) have larger 95% credible intervals for the proportion estimate (colored shading). From this view, it is evident that several further lineages waxed and waned notably over the sampling period, including B.1.1.434 and B.1.526.

Analyzing lineages in spike gene target failures. As a control, we assessed the variants associated with spike (S) gene target failure samples (Fig. 2C and 4 and Table S3),

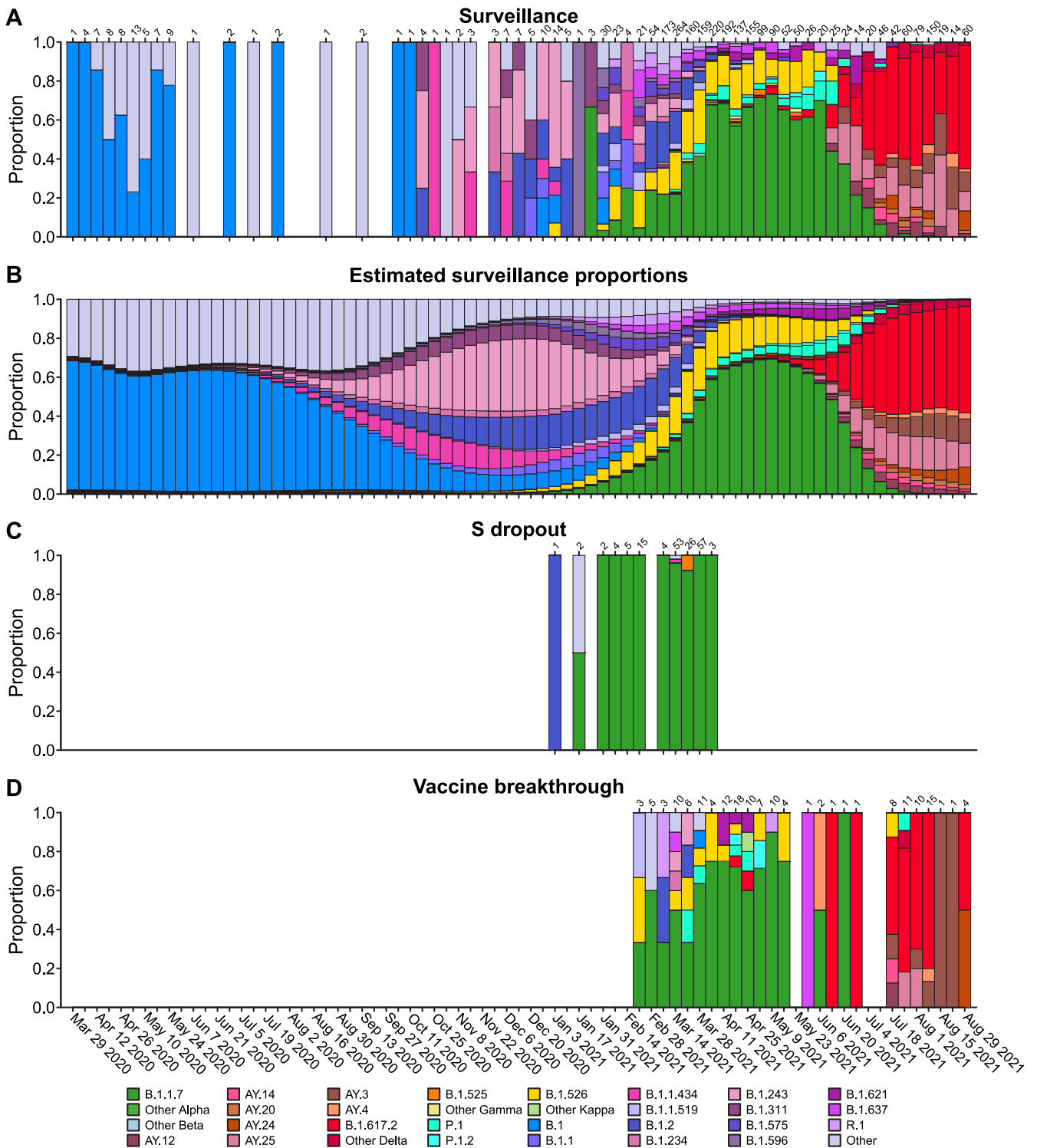


FIG 2 Comparison of viral genome sequence data from surveillance samples (A, B) to spike target gene failures (C) and vaccine breakthrough samples (D). (A) Longitudinal stacked bar graph depicting the SARS-CoV-2 variants present in surveillance samples from the Delaware Valley, shown as the proportion of genomes classified as each variant lineage within each week. The numbers of genomes sampled each week are shown above the graph. Variants are colored according to the key at the bottom of the figure. (B) Markings are the same as in (A), but showing the proportions of variants estimated from the count data in (A) using Bayesian autoregressive moving average multinomial logistic regression. (C) Markings as in (A), but showing counts of spike target gene failures samples. (D) Markings as in (A), but showing the counts of vaccine breakthrough samples. Designation of lineages as variants of concern and variants being monitored is presented in Table S2. For vaccine breakthroughs, the time window compared was from the introduction of widespread vaccination (March 1, 2021) to the end of our sampling period (September 3, 2021).

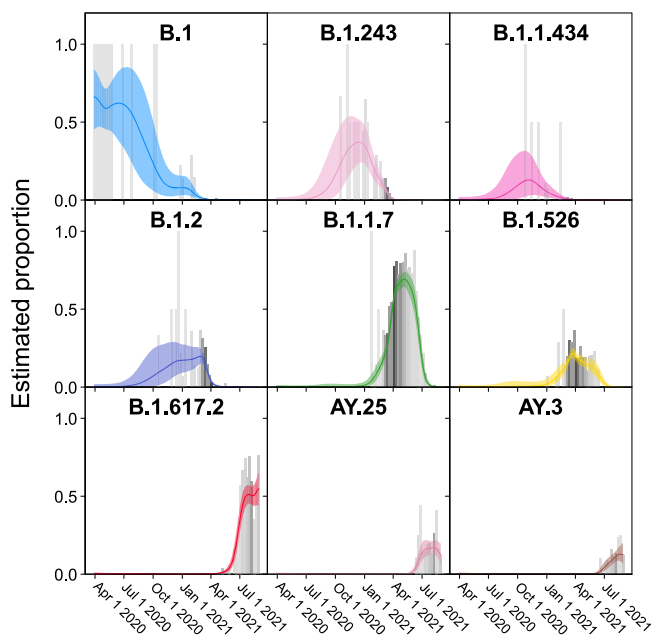


FIG 3 Frequencies of individual variants estimated using Bayesian autoregressive moving average multinomial logistic regression. Time is shown along the x axis and estimated proportions of the surveillance population along the y axis. The gray bars indicate raw proportions from the count data shaded by the number of observations observed in a given week (darker indicating more samples) while the colored lines indicate the proportion estimated by the Bayesian model. The light-colored envelopes around each line show the 95% credible intervals for the proportion. Only lineages achieving an estimated proportion of $>10\%$ in any given week are shown.

where an amplicon overlapping the 69–70 deletion in spike failed selectively. To track the spread and variation of B.1.1.7 during the alpha wave of infection, we sequenced 172 S gene target failure samples in the Delaware Valley from January to April 2021. Upon variant assignment, genomes from 96.1% of spike target gene failure samples corresponded to B.1.1.7. Rarely, other lineages were seen that share the spike deletion with B.1.1.7 (B.1.375 and B.1.525) or were likely stochastic failures of S amplification or genomes with novel combinations of mutations.

As expected, we found that the B.1.1.7 lineage was estimated to be highly enriched in the S gene target failure set over what would be expected from its frequency in the surveillance lineage counts (Fig. 4) (odds ratio of 120; 95% credible interval 38–470). The B.1.525 lineage, which also contains the S69-70 deletion, was estimated to be enriched as well (odds ratio 33; 95% credible interval 1.7–380). Note that this enrichment was detected based on counts of only two B.1.525 genomes in the S gene target failure set and 12 surveillance genomes, both emphasizing the sensitivity of the model and highlighting that the detection of enrichment can be easier in lineages rare in the surveillance population. Other lineages showed little association with spike target gene failures, as expected (note that only a single B.1.375 lineage sample was observed and thus it did not reach the threshold for lineage-specific analysis and was grouped in “Other”). Thus, the analysis of S target failures confirmed that the Bayesian autoregressive moving average categorical regression model was effective at identifying overrepresented lineages relative to the time-adjusted expectation from surveillance.

Enrichment of delta variant B.1.617.2 in vaccine breakthrough samples. We then applied the Bayesian model to analyzing vaccine breakthrough samples (Fig. 2D and 5 and Table S4). There was less clear enrichment in vaccine breakthroughs than in the S gene target failure set. The delta lineage B.1.617.2 did show signs of enrichment in breakthrough samples, with mean odds ratio of 3 (95% credible interval 0.89 to 11). The one-sided posterior probability of any enrichment for B.1.617.2 in vaccine breakthrough was estimated at 96% with a 73% probability of more than a 2-fold enrichment in odds of appearing in vaccine breakthrough. This enrichment was not observed for all lineages grouped within delta, with both a group-wise estimate for all delta lineages and the specific estimates for the

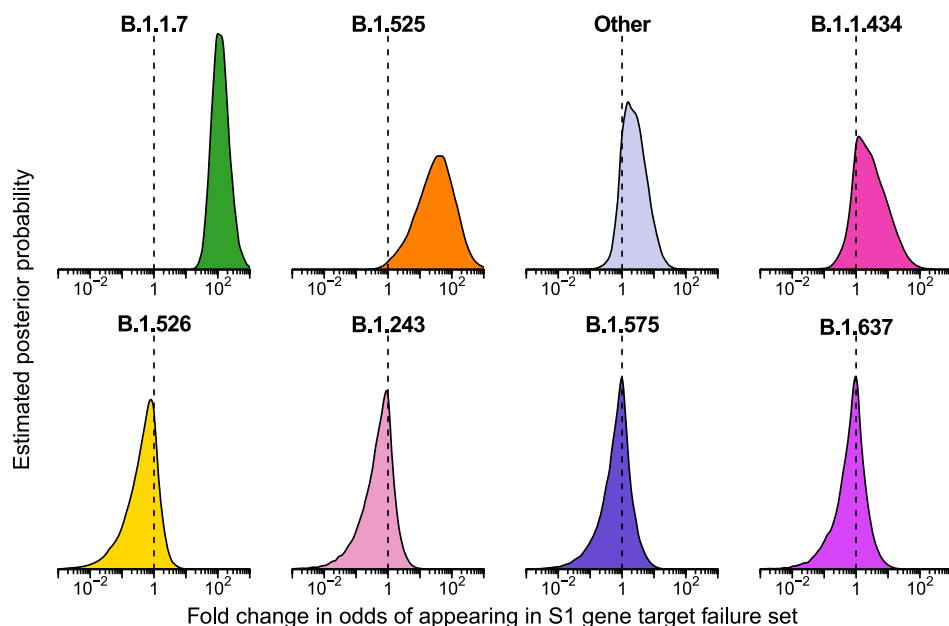


FIG 4 Estimated posterior probability densities for the enrichment of variants among spike gene target failures produced by Bayesian autoregressive moving average multinomial logistic regression. The x axis shows the fold enrichment/depletion in the odds of a variant (labeled above each plot) appearing in the spike gene target failure set relative to the proportions estimated in the surveillance population and the y axis shows the posterior probability. No enrichment (fold change of 1) is indicated by the dashed vertical line. Increased density to the right indicates greater likelihood of inclusion among spike gene target failures, increased density to the left indicates decreased likelihood. The four lineages with highest posterior probability of enrichment are shown on the top and four lineages with high posterior probability of depletion are shown on the bottom. Color coding indicates the variant queried with colors as in earlier figures.

other delta AY lineages, including AY.3, AY.4, AY.14, and AY.24, showing little enrichment over surveillance. These alternative delta lineages also had fewer samples observed, so a lack of power may be a partial explanation. Note that the power to detect enrichment of delta is limited by its rapid spread—observation of a delta vaccine breakthrough case at a

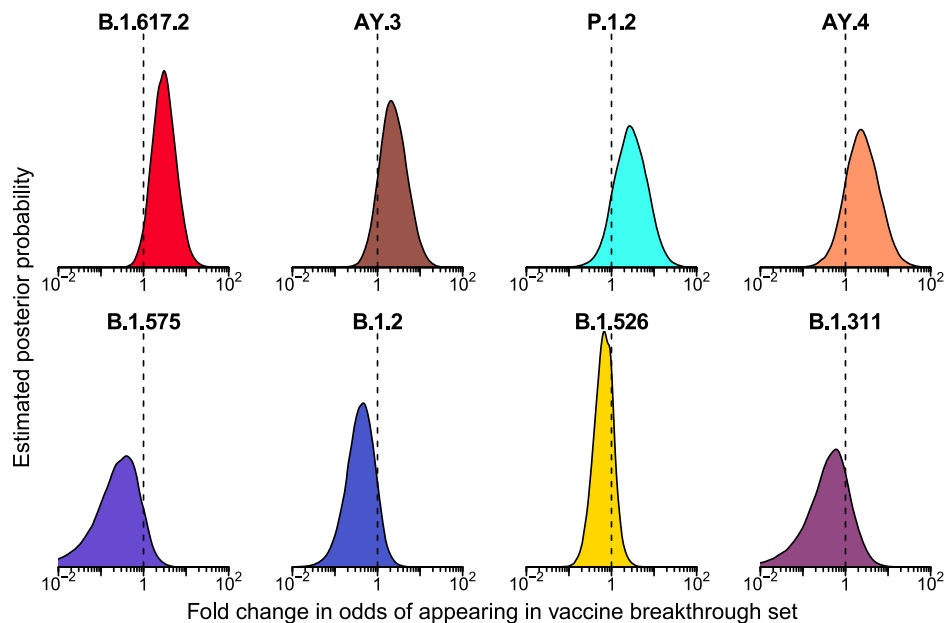


FIG 5 Posterior probability densities for enrichment of variants among vaccine breakthrough samples compared to the surveillance population as estimated by Bayesian autoregressive moving average multinomial logistic regression. Markings as in Fig. 4

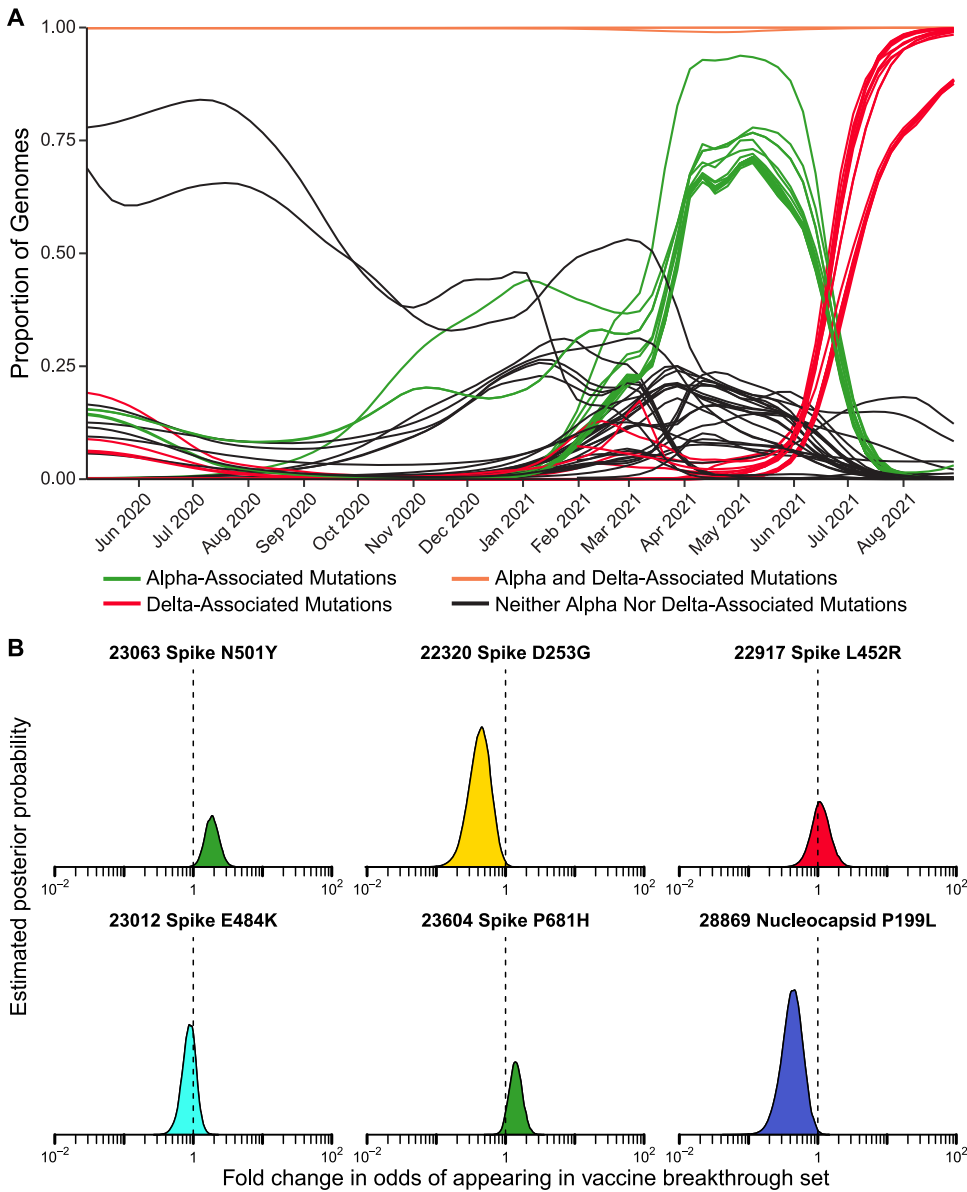


FIG 6 Assessment of enrichment of specific base substitutions and deletions in vaccine breakthrough samples. (A) Longitudinal frequencies of individual mutations estimated using Bayesian autoregressive moving average logistic regression models. Red indicates mutations commonly found in delta lineages. Green indicates mutations commonly found in alpha lineages. Orange indicates mutations shared by almost all lineages in the study. Black indicates mutations found in other subsets of lineages. (B) Estimated posterior probability densities summarizing the fold enrichment/depletion in odds of a mutation appearing in the vaccine breakthrough samples over its proportions estimated from surveillance samples. The mutations estimated as most enriched and most depleted are shown along with other mutations of interest. Markings as in Fig. 4 and 5.

time point when almost all surveillance population infections were also delta provides little information. Similarly, lineages that had already waned by the time of widespread vaccination are difficult to assess since there was little chance for them to appear in vaccine breakthrough cases.

Assessing possible enrichment of specific mutations in spike gene target failures and vaccine breakthroughs. We next investigated whether any individual base substitutions in the viral genome were selectively associated with spike gene target failure and vaccine breakthrough samples. A summary of mutations detected is presented in Table S5. We adapted the Bayesian autoregressive moving average logistic regression method to assess the behavior of single mutations and amino acid substitutions. Fig. 6 and Table S6

summarizes results. Many mutations varied in frequency over the course of the study often paralleling the profile of the most abundant lineages containing them (Fig. 6A).

As expected, the 69–70 deletion was found to be highly enriched in the S gene target failure set (odds ratio of 205, 95% credible interval of 65–581). Enrichment among S gene target failures was also seen for another 28 mutations (Table S6), where all the most strongly enriched substitutions were characteristic of the alpha variant and thus due to “hitchhiker” effects. Underrepresentation in spike target gene failures was also seen in mutations in multiple additional open reading frames, reflecting mutations found in variants lacking the 69–70 deletion. Each mutation was analyzed separately without accounting for genomic linkage so this is as expected.

For vaccine breakthroughs (Fig. 6B and Table S7), the most notably enriched substitution was N501Y, which showed an odds ratio of 2.04 (95% credible interval of 1.25–3.18). The N501Y substitution is found in multiple VBM, including alpha, beta, and gamma, and is reported to increase the affinity of spike protein binding to the ACE2 receptor (23, 37, 38) and to diminish binding of some human antibodies to spike (24). The spike substitution P681H was also slightly enriched; this substitution is near the furin cleavage site and may promote efficient proteolysis. Several additional substitutions in spike were potentially enriched (D614G, P681R, D950N), but the 95% credible interval included one, so the evidence for enrichment was weaker (Table S7). Two other closely studied spike substitutions, E484K and L452R, were not notably enriched among vaccine breakthrough cases, and the D253G substitution was modestly depleted. Among non-spike substitutions, the P199L substitution in the nucleocapsid showed notable depletion (Table S7). Stepping back, these findings together emphasize the potential importance of the N501Y substitution in vaccine breakthrough.

DISCUSSION

Vaccination against SARS-CoV-2 has been highly effective at preventing infection, hospitalization and death. However, infections occasionally take place despite vaccination. It has been suggested that certain lineages of SARS-CoV-2 are more prone to evade vaccination, however, interpretation of counts of vaccine breakthrough is difficult due to variation over time in circulating lineages, the numbers of vaccinated individuals, times since vaccination and the application of nonpharmaceutical interventions. Here, we introduce modeling based on Bayesian autoregressive moving average multinomial logistic regression, which estimates the proportions of lineages in the surveillance data over time, allowing comparison of lineage counts in populations of interest to time-matched surveillance estimates. Multiple studies have investigated whether certain lineages are more likely to appear in vaccine breakthrough infections (7, 8, 39–43). However, few studies have as large numbers of sequenced samples for both background surveillance and for vaccine breakthrough cases within a matched geographical region as reported here. Thus, our contribution provides targeted data on the nature of breakthroughs at nucleotide resolution.

Using this method, we found that there was evidence for enrichment in the odds of variant lineages appearing in the vaccine breakthrough population relative to the general surveillance population, with the delta variant B.1.617.2 showing the strongest signal with an estimated enrichment of 3-fold (95% credible interval 0.89 to 11). Several previous studies have also noted delta variants to be enriched among vaccine breakthroughs (7, 8); our data provides further support based on careful statistical analysis.

Using a similar Bayesian autoregressive moving average logistic regression model, we interrogated enrichment of point substitutions among vaccine breakthrough cases. The N501Y substitution stood out as particularly associated with vaccine breakthroughs. N501Y is found in several of the VBM, and this substitution is well known to increase affinity of binding for the viral spike protein to the ACE2 receptor (23, 37, 38), as well as reduce antibody binding to promote immune evasion (24). Recently the N501Y substitution was suggested to be a central feature in a meta-signature of up to 35 mutations which recur in alpha, beta, gamma and other lineages, and which mark a viral fitness peak reflecting optimization for spread in humans (44). In contrast, several substitutions were estimated to be less likely to appear in

vaccine breakthrough samples, including nucleocapsid mutation P199L (detected in B.1.2, B.1.526, and B.1.596; odds ratio of 0.5; 95% credible interval of 0.24–0.88), and spike mutation D253G (detected in B.1.526; odds ratio of 0.5; 95% credible interval of 0.23–0.89). The spike D253G substitution has been implicated in MAb evasion (45), and nucleocapsid P199L may alter the assembly of SARS-CoV-2 VLPs (46), functions possibly linked to their depletion in vaccine breakthrough.

This study has several limitations. For the vaccine breakthrough samples, we did not have paired immunological data, so it is unknown whether the vaccinations provoked protective immune responses. Our surveillance data were from a mixture of hospitalized patients, symptomatic subjects tested in clinical diagnostic laboratories, and asymptomatic subjects tested weekly at an academic institution, and so our sampling may not be entirely representative of viruses circulating in the community. Several different amplification and sequencing methods were used to acquire data. Documentation of vaccination status may be incomplete. For all the viral genome sequencing, only samples achieving a threshold level of viral RNA could be sequenced (roughly Ct <28 for swab-based testing, Ct <20 for saliva), so it is possible that other variants predominate in subjects with lower viral loads. Studies have been initiated to address some of these concerns.

In summary, similar to other regions in the US and around the world, VBM/VOC and other lineages waxed and waned in the Delaware Valley, with delta lineages ultimately comprising all recent samples in the fall of 2021. Widespread vaccination was introduced in the region in winter 2020–2021, and increasing numbers of breakthrough infections have since been detected. To compare the lineages found in breakthrough with those observed in general surveillance, we introduce analysis based on Bayesian autoregressive moving average multinomial logistic regression and show that delta variant B.1.617.2 showed 3-fold enrichment in vaccine breakthroughs. The N501Y substitution stood out among point substitutions for enrichment in vaccine breakthroughs. We expect that these modeling methods will be useful in monitoring the effectiveness of vaccination programs going forward as novel variants of SARS-CoV-2 continue to emerge.

MATERIALS AND METHODS

Human subjects. The University of Pennsylvania Institutional Review Board (IRB) reviewed the research protocol and deemed the limited data elements extracted with positive SARS-CoV-2 specimens to be exempt from human subject research per 45 CFR 46.104, category 4 (IRB #848605). For hospitalized subjects at the University of Pennsylvania, following informed consent (IRB protocol #823392), patients were sampled by collection of saliva, oropharyngeal and/or nasopharyngeal swabs, or endotracheal aspirates if intubated, as previously described (32). Clinical data were extracted from the electronic medical record. Further samples were collected from asymptomatic subjects detected in a screening program at the Perelman School of Medicine at the University of Pennsylvania and symptomatic subjects tested throughout the PennMedicine clinical network under IRB protocols #843565 and #848608. Human samples were collected at Children's Hospital of Philadelphia under protocol # 21-018726 approved by the Children's Hospital of Philadelphia IRB. Human samples were collected at Thomas Jefferson University under protocol number IRB Control # 21E.441 approved by the Thomas Jefferson University IRB. Vaccine breakthrough cases were identified by either report to HUP Infection and Control, Jefferson Infection and Control, CHOP Department of Infection Control and Prevention, the Philadelphia Department of Public Health, or, where applicable, chart review.

Sequencing methods. Several sequencing methods were used to acquire viral whole-genome sequences.

The POLAR protocol was used to acquire the majority of viral genome sequences (47). Illumina's NextSeq instrument was used to gather sequence data. In detail, 5 μ l of viral RNA, 0.5 μ l of 50 μ M Random Hexamers (Thermo Fisher, N8080127), 0.5 μ l of 10 mM dNTPs Mix (Thermo Fisher, 18427013), and 1 μ l nuclease-free water was incubated at for 5 min at 65°C proceeded by a 1 min incubation at 4°C. To perform reverse transcription, 6.5 μ l from the previous reaction, 0.5 μ l SuperScript III Reverse Transcriptase (Thermo Fisher, 18080085), 0.5 μ l of 0.1M DTT (Thermo Fisher, 18080085), 0.5 μ l of RNaseOut (Thermo Fisher, 18080051), and 2 μ l of 5X First-Strand Buffer (Thermo Fisher, 18080085) was incubated for 50 min at 42°C, followed by an incubation for 10 min at 70°C, and then held at 4°C. To amplify the cDNA, artic-ncov2019 version 3 primers were used (IDT). To perform PCR amplification of the viral cDNA, the following reagents were added to 2.5 μ l of the previous mixture: 0.25 μ l Q5 Hot Start DNA polymerase (NEB, M0493S), 5 μ l of 5X Q5 Reaction Buffer (NEB, M0493S), 0.5 μ l of 10 mM dNTPs Mix (NEB, N0447S), either 4.0 μ l of pooled primer set 1 or 3.98 μ l of pooled primer set 2, and nuclease-free water to bring to a final volume of 25 μ l. This PCR amplification of the viral cDNA used the following conditions: 98°C for 30 s for 1 cycle, 25 cycles at 98°C for 15 s and 65°C for 5 min, and then held at 4°C. Amplicons generated by the two primer sets from the same sample were pooled then diluted to a concentration of 0.25 ng/ μ l. The Nextera library was prepared using the Nextera XT Library Preparation kit (Illumina, FC-131-1096) and the IDT for Illumina DNA/RNA UD Indexes Set A and B (Illumina, 20027213, 20027214, 20027215, 20027216). The Quant-iT PicoGreen dsDNA quantitation assay kit was used to quantify the DNA of each sample

(Invitrogen, P7589). The samples were pooled in equal quantities, and the pooled library was quantified using the Qubit1X dsDNA HS assay kit (Invitrogen, Q33230). The library was sequenced on an Illumina NextSeq.

Another sequencing method, used at CHOP, was Paragon. Paragon sequencing was carried out as in (22). Briefly, RNA was extracted from nasopharyngeal swab samples using QIAamp Viral RNA Mini (Qiagen). Whole-genome sequencing was carried out by the Genomics Core Facility at Drexel University. Amplification was performed using Paragon Genomics CleanPlex SARS-CoV-2 Research and Surveillance NGS Panel 1 and 2. Libraries were quantified using the Qubit dsDNA HS (High Sensitivity) assay kit (Invitrogen) with the Qubit Fluorometer (Invitrogen). Library quality was assessed using Agilent High Sensitivity DNA kit and the 2100 Bioanalyzer instrument (Agilent). Libraries were normalized to 5 nM and pooled in equimolar concentrations. The resulting pool was quantified again using the Qubit dsDNA HS (High Sensitivity) assay kit (Invitrogen) and diluted to a final concentration of 4 nM; libraries were denatured and diluted according to Illumina protocols and loaded on the MiSeq at 10pM. Paired-end and dual-indexed 2x150bp sequencing was carried out using MiSeq Reagent Kits v3 (300 cycles).

The Thomas Jefferson University site carried out sequencing using the Illumina RPIIP and CovidSeq methods essentially as per the manufacturer's instructions. A nasopharyngeal swab specimen that was tested positive by PCR for SARS-CoV-2 was used for this study. Vaccine breakthrough specimens were identified by Jefferson Occupational Health Network for Employees and Students (JOHN) and were sequenced with a designation of VBT. Randomly selected residual positive specimens from the Molecular & Genomic Pathology Laboratory at Jefferson during the same period were sequenced for the purpose of epidemiology surveillance and labeled as SURV. RNA was extracted from 200 μ l of the specimen and eluted in 110 μ l using the bioMérieux EasyMag Extraction System, following EasyMag's Generic protocol. Whole-genome sequencing for SARS-CoV-2 was subsequently performed at the Molecular & Genomic Pathology Laboratory at Thomas Jefferson University Hospital. With the exception of a few samples, samples were sequenced using Illumina COVIDSeq protocol (Illumina) following manufacturer's guideline. Remaining specimens were sequenced using the Swift Normalase Amplicon Panels (SNAP) Core kit along with the SARS-CoV-2 Additional Genome Coverage Primer Panel following the manufacturer's protocol. For data analysis, Illumina's Local Run Manager's GenerateFASTQ module was used to generate the fastq files for all specimens. The fastq files were transferred to UPenn's secure data server for further processing to obtain QC metrics and lineage data.

Samples with VSP numbers lower than VSP00256 (Table S1) were previously described in Everett et al. (32).

Data analysis. To process sequence data, sequence reads are trimmed to remove low quality base calls ($< Q20$) and aligned to the original Wuhan reference sequence (NC_045512.2) with the BWA aligner tool (v0.7.17) (48), after which alignments are filtered with the Samtools package (v1.10) (49). Sequencing depth is determined for each position in the viral genome and genomes are accepted for analysis when $\geq 95\%$ of genome positions are covered with a read depth of ≥ 5 reads. Variant positions are called with the Bcftools package (v1.10.2-34) (50) requiring PHRED scores ≥ 20 and variant read frequencies $\geq 50\%$ of the total reads. The nature of substitutions is determined by retrieving reading frames from the reference GENBANK record, translating, and determining the native and mutant residues.

Variants were assigned using the Pangolin lineage software (Pangolin version 3.1.11 with the PangoleARN 2021-08-24 release). Note that these lineages are updated regularly and so are expected to change over time. Point mutations were assigned using a previously described bioinformatics pipeline (31).

Statistical analysis of variant enrichment and mutation enrichment in subsets of the data were separately assessed using a Bayesian model and Markov chain Monte Carlo sampling implemented in Stan (51). The model takes the vector of counts of the variants or mutations seen in the surveillance sampling for a given week, counts_w, and assumes they are multinomially distributed where:

$$\text{counts}_w \sim \text{Multinomial}(p_{*,w})$$

with probabilities modeled as multinomial logistic:

$$p_{i,w} = \frac{e^{x_{i,w}}}{\sum_{j=1}^n e^{x_{j,w}}}$$

where $p_{i,w}$ is the true proportion of variant or mutation i on week w . $p_{*,w}$ indicates the vector of all lineage probabilities for week w and n is the total number of variants or mutations observed. The underlying proportions for a given week are assumed to be centered around the proportions observed in the prior week (autoregressive) plus a change term that is itself correlated with the change observed in the prior week (moving average):

$$x_{i,w} = x_{i,w-1} + \text{change}_{i,w}$$

$$\text{change}_{i,w} \sim \text{Normal}(\text{change}_{i,w-1}, \sigma)$$

where $\text{change}_{i,w}$ is the change in log odds of variant or mutation i on week w . The initial starting proportions are given flat priors:

$$x_{i,1} = \begin{cases} 0, & i = 1 \\ \text{Normal}(0, 10), & i > 1 \end{cases}$$

The standard deviation for changes, σ , was given a $\text{Gamma}(1, 2)$ prior distribution.

To assess if there's enrichment of a particular variant or mutation in the vaccine breakthrough, counts of vaccine breakthrough samples for a given week, vaccineCounts_w , were modeled as:

$$\text{vaccineCounts}_w \sim \text{Multinomial}(v_{*,w})$$

where:

$$v_{i,w} = \frac{e^{y_{i,w}}}{\sum_{j=1}^n e^{y_{j,w}}}$$

and

$$y_{i,w} = x_{i,w} + \beta_i + \delta_{group_i}$$

where β_i measures the log odds enrichment/depletion of variant or mutation i in the vaccine breakthrough population over the surveillance population, $group_i$ indicates the WHO classification for Pango lineage i and δ_{group_i} indicates the enrichment or depletion for a CDC VBM/VOC classification containing more than one Pango lineage, e.g., Delta containing B.1.617.2 and other variants with an AY prefix. In these data, Delta, Gamma and Non-VBM/VOC groupings contained more than one lineage above the abundance threshold. For WHO classifications containing only a single lineage, δ_{group_i} was set to 0. The β and remaining δ were given DoubleExponential(0, 1) priors.

This was repeated equivalently in the samples collected as S gene target failures, with the counts, SFailCounts_w , modeled as:

$$\text{SFailCounts}_w \sim \text{Multinomial}(u_{*,w})$$

where:

$$u_{i,w} = \frac{e^{z_{i,w}}}{\sum_{j=1}^n e^{z_{j,w}}}$$

and

$$z_{i,w} = x_{i,w} + \alpha_i$$

The α were given DoubleExponential(0, 1) priors.

For the assessment of lineage enrichment, genomes from lineages that had more than 10 genomes assigned to them were included as their individual lineages, genomes from lineages with 10 or fewer assignments that were listed within VBM/VOC were included as miscellaneous classification for their particular WHO category, e.g., "Other delta" and all other genomes were grouped into an Other category. For the assessment of mutations, each mutation that was found in greater than 5% of genomes was independently modeled as above but replacing the Multinomial distributions with Binomial and removing the inapplicable lineage group δ terms.

Data availability. All viral genome sequences acquired in this study have been deposited in GISAID and at NCBI under accession numbers listed in Table S1. Analysis code is archived on Zenodo at <https://doi.org/10.5281/zenodo.5888338>. Sequence processing software and intermediate files used in this study are available at <https://doi.org/10.5281/zenodo.5559699>. A list of key reagents used in this study is in Table S8.

SUPPLEMENTAL MATERIAL

Supplemental material is available online only.

TABLE S1, PDF file, 0.3 MB.

TABLE S2, PDF file, 0.1 MB.

TABLE S3, PDF file, 0.03 MB.

TABLE S4, PDF file, 0.03 MB.

TABLE S5, PDF file, 1.9 MB.

TABLE S6, PDF file, 0.04 MB.

TABLE S7, PDF file, 0.04 MB.

TABLE S8, PDF file, 0.04 MB.

ACKNOWLEDGMENTS

We are grateful to individuals and their families who volunteered to provide specimens, members of the Bushman and Collman laboratories for help and suggestions, and to Laurie Zimmerman for artwork and help with manuscript preparation. We acknowledge help from all the staff of the Philadelphia Department of Public Health. Funding was provided by a

contract award from the Centers for Disease Control and Prevention (CDC BAA 200-2021-10986 and 75D30121C11102/000HCVL1-2021-55232), philanthropic donations to the Penn Center for Research on Coronaviruses and Other Emerging Pathogens, and in part by NIH grant R61/33-HL137063. B.J.K. is supported by NIH K23 AI 121485. Additional assistance was provided by the Penn Center for AIDS Research (P30-AI045008) and the Women's Committee of CHOP. This project has been funded in whole or in part with Federal funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services, under Contract No. 75N93021C00015.

A.D.M., J.E., S.S.-M., R.G.C., K.R., B.J.K., and F.D.B. designed the study; S.A.W., J.G.-W., L.A.K., A.S.F., A.D.M., C.B., S.R., J.H., R.D., L.D., A.A., E.K., S.C., C.N., J.C.M., P.P., C.C., K.S., A.G., M.F., R.G.C., K.R., A.J.-B., L.C., S.L., A.G., and B.J.K. managed acquisition of clinical specimens; A.D.M., C.B., S.R., L.D., A.A., J.C.M., P.P., N.B., B.R., Z.-X.W., P.H., A.M.R., A.G., A.M., and F.D.B. carried out sequencing; A.D.M., A.M.M., C.B., S.R., A.A., J.C.M., P.P., S.S.-M., B.J.K., J.E., and F.D.B. carried out bioinformatic and statistical analysis; A.D.M., S.S.-M., J.E., K.R., B.J.K., R.G.C., and F.D.B. wrote the article.

RGC reports support to his lab for COVID-19 work unrelated to the current manuscript from OraSure, Inc, and ongoing collaborations with Resilient Biotics, Inc.

We have no other conflicts of interest to report.

REFERENCES

- Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, Wang W, Song H, Huang B, Zhu N, Bi Y, Ma X, Zhan F, Wang L, Hu T, Zhou H, Hu Z, Zhou W, Zhao L, Chen J, Meng Y, Wang J, Lin Y, Yuan J, Xie Z, Ma J, Liu WJ, Wang D, Xu W, Holmes EC, Gao GF, Wu G, Chen W, Shi W, Tan W. 2020. Genomic characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor binding. *Lancet* 395:565–574. [https://doi.org/10.1016/S0140-6736\(20\)30251-8](https://doi.org/10.1016/S0140-6736(20)30251-8).
- Holmes EC. 2003. Error thresholds and the constraints to RNA virus evolution. *Trends Microbiol* 11:543–546. <https://doi.org/10.1016/j.tim.2003.10.006>.
- Holmes EC, Goldstein SA, Rasmussen AL, Robertson DL, Crits-Christoph A, Wertheim JO, Anthony SJ, Barclay WS, Boni MF, Doherty PC, Farrar J, Geoghegan JL, Jiang X, Leibowitz JL, Neil SJD, Skern T, Weiss SR, Worobey M, Andersen KG, Garry RF, Rambaut A. 2021. The origins of SARS-CoV-2: a critical review. *Cell* 184:4848–4856. <https://doi.org/10.1016/j.cell.2021.08.017>.
- Meng B, Kemp SA, Papa G, Dattir R, Ferreira IATM, Marelli S, Harvey WT, Lytras S, Mohamed A, Gallo G, Thakur N, Collier DA, Mlcochova P, Duncan LM, Carabelli AM, Kenyon JC, Lever AM, De Marco A, Saliba C, Culp K, Cameroni E, Matheson NJ, Piccoli L, Corti D, James LC, Robertson DL, Bailey D, Gupta RK, COVID-19 Genomics UK (COG-UK) Consortium. 2021. Recurrent emergence of SARS-CoV-2 spike deletion H69/V70 and its role in the Alpha variant B.1.1.7. *Cell Rep* 35:109292. <https://doi.org/10.1016/j.celrep.2021.109292>.
- Kistler KE, Huddleston J, Bedford T. 2021. Rapid and parallel adaptive mutations in spike S1 drive clade success in SARS-CoV-2. *bioRxiv* <https://doi.org/10.1101/2021.09.11.459844>.
- Goel RR, Painter MM, Apostolidis SA, Mathew D, Meng W, Rosenfeld AM, et al. 2021. mRNA Vaccination Induces Durable Immune Memory to SARS-CoV-2 with Continued Evolution to Variants of Concern. *bioRxiv* <https://doi.org/10.1101/2021.08.23.457229>.
- Scobie HM, Johnson AG, Suthar AB, Severson R, Alden NB, Balter S, Bertolino D, Blythe D, Brady S, Cadwell B, Cheng I, Davidson S, Delgado J, Devinney K, Duchin J, Duwell M, Fisher R, Fleischauer A, Grant A, Griffin J, Haddix M, Hand J, Hanson M, Hawkins E, Herlihy RK, Hicks L, Holtzman C, Hoskins M, Hyun J, Kaur R, Kay M, Kidrowski H, Kim C, Komatsu K, Kugeler K, Lewis M, Lyons BC, Lyons S, Lynfield R, McCaffrey K, McMullen C, Milroy L, Meyer S, Nolen L, Patel MR, Pogojans S, Reese HE, Saue A, Sell J, Sokol T, et al. 2021. Monitoring Incidence of COVID-19 Cases, Hospitalizations, and Deaths, by Vaccination Status - 13 U.S. Jurisdictions, April 4-July 17, 2021. *MMWR Morb Mortal Wkly Rep* 70:1284–1290. <https://doi.org/10.15585/mmwr.mm7037e1>.
- Kislaya I, Rodrigues EF, Borges V, Gomes JP, Sousa C, Almeida JP, et al. 2021. Delta variant and mRNA Covid-19 vaccines effectiveness: higher odds of vaccine breakthroughs. *medRxiv*.
- Saad-Roy CM, Morris SE, Metcalf CJE, Mina MJ, Baker RE, Farrar J, Holmes EC, Pybus OG, Graham AL, Levin SA, Grenfell BT, Wagner CE. 2021. Epidemiological and evolutionary considerations of SARS-CoV-2 vaccine dosing regimes. *Science* 372:363–370. <https://doi.org/10.1126/science.abg8663>.
- Fauver JR, Petrone ME, Hodcroft EB, Shioda K, Ehrlich HY, Watts AG, Vogels CBF, Brito AF, Alpert T, Muyombwe A, Razeq J, Downing R, Cheemarla NR, Wyllie AL, Kalinich CC, Ott IM, Quick J, Loman NJ, Neugebauer KM, Greninger AL, Jerome KR, Roychoudhury P, Xie H, Shrestha L, Huang M-L, Pitzer VE, Iwasaki A, Omer SB, Khan K, Bogoch II, Martinello RA, Foxman EF, Landry ML, Neher RA, Ko AI, Grubaugh ND. 2020. Coast-to-Coast Spread of SARS-CoV-2 during the Early Epidemic in the United States. *Cell* 181:990–996.e5. <https://doi.org/10.1016/j.cell.2020.04.021>.
- Oude Munnink BB, Nieuwenhuijse DF, Stein M, O'Toole A, Haverkate M, Mollers M, et al. 2020. Rapid SARS-CoV-2 whole-genome sequencing and analysis for informed public health decision-making in the Netherlands. *Nat Med*.
- Gonzalez-Reiche AS, Hernandez MM, Sullivan MJ, Ciferri B, Alshammari H, Obla A, Fabre S, Kleiner G, Polanco J, Khan Z, Albuquerque B, van de Guchte A, Dutta J, Francoeur N, Melo BS, Oussenko I, Deikus G, Soto J, Sridhar SH, Wang Y-C, Twyman K, Kasarskis A, Altman DR, Smith M, Sebra R, Aberg J, Krammer F, García-Sastre A, Luksza M, Patel G, Paniz-Mondolfi A, Gitman M, Sordillo EM, Simon V, van Bakel H. 2020. Introductions and early spread of SARS-CoV-2 in the New York City area. *Science* 369:297–301. <https://doi.org/10.1126/science.abc1917>.
- Moreno GK, Braun KM, Riemersma KK, Martin MA, Halfmann PJ, Crooks CM, et al. 2020. Distinct patterns of SARS-CoV-2 transmission in two nearby communities in Wisconsin, USA. *medRxiv* <https://doi.org/10.1101/2020.07.09.20149104>.
- Taboada B, Vazquez-Perez JA, Muñoz-Medina JE, Ramos-Cervantes P, Escalera-Zamudio M, Boukadida C, Sanchez-Flores A, Isa P, Mendieta-Condado E, Martínez-Orozco JA, Becerril-Vargas E, Salas-Hernández J, Grande R, González-Torres C, Gaytán-Cervantes FJ, Vazquez G, Pulido F, Araiza-Rodríguez A, Garcés-Ayala F, González-Bonilla CR, Grajales-Muñiz C, Borja-Aburto VH, Barrera-Badillo G, López S, Hernández-Rivas L, Perez-Padilla R, López-Martínez I, Ávila-Ríos S, Ruiz-Palacios G, Ramírez-González JE, Arias CF. 2020. Genomic Analysis of Early SARS-CoV-2 Variants Introduced in Mexico. *J Virol* 94 <https://doi.org/10.1128/JVI.01056-20>.
- Candido DS, Claro IM, de Jesus JG, Souza WM, Moreira FRR, Dellicour S, et al. 2020. Evolution and epidemic spread of SARS-CoV-2 in Brazil. *Science* <https://doi.org/10.1126/science.abd2161>.
- Jangra S, Ye C, Rathnasinghe R, Stadlbauer D, Krammer F, Simon V, et al. 2021. The E484K mutation in the SARS-CoV-2 spike protein reduces but does not abolish neutralizing activity of human convalescent and post-vaccination sera. *medRxiv*.
- Greaney AJ, Loes AN, Crawford KHD, Starr TN, Malone KD, Chu HY, Bloom JD. 2021. 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. <https://doi.org/10.1016/j.chom.2021.02.003>.
- Rockett RJ, Amott A, Lam C, Sadsad R, Timms V, Gray K-A, Eden J-S, Chang S, Gall M, Draper J, Sim EM, Bachmann NL, Carter I, Basile K, Byun R, O'Sullivan MV, Chen SC-A, Maddocks S, Sorrell TC, Dwyer DE, Holmes EC, Kok J, Prokopenko M, Sintchenko V. 2020. Revealing COVID-19 transmission in Australia by SARS-CoV-2 genome sequencing and agent-based modeling. *Nat Med* 26:1398–1404. <https://doi.org/10.1038/s41591-020-1000-7>.
- Grifoni A, Weiskopf D, Ramirez SI, Mateus J, Dan JM, Moderbacher CR, Rawlings SA, Sutherland A, Premkumar L, Jadi RS, Marrama D, de Silva

- AM, Frazier A, Carlin AF, Greenbaum JA, Peters B, Krammer F, Smith DM, Crotty S, Sette A. 2020. Targets of T Cell Responses to SARS-CoV-2 Coronavirus in Humans with COVID-19 Disease and Unexposed Individuals. *Cell* 181:1489–1501. <https://doi.org/10.1016/j.cell.2020.05.015>.
20. Agerer B, Koblishke M, Gudipati V, Montano-Gutierrez LF, Smyth M, Popa A, et al. 2021. SARS-CoV-2 mutations in MHC-I-restricted epitopes evade CD8(+) T cell responses. *Sci Immunol* 6. <https://doi.org/10.1126/sciimmunol.abg6461>.
 21. Skidmore PT, Kaelin EA, Holland LA, Maqsood R, Wu LI, Mellor NJ, Blain JM, Harris V, LaBaer J, Murugan V, Lim ES. 2021. Genomic Sequencing of SARS-CoV-2 E484K Variant B.1.243.1, Arizona, USA. *Emerg Infect Dis* 27: 2718–2720. <https://doi.org/10.3201/eid2710.211189>.
 22. Moustafa AM, Bianco C, Denu L, Ahmed A, Coffin SE, Neide B, Everett J, Reddy S, Rabut E, Deseignora J, Feldman MD, Rodino KG, Bushman F, Harris RM, Chang Mell J, Planet PJ. 2021. Comparative Analysis of Emerging B.1.1.7+E484K SARS-CoV-2 Isolates. *Open Forum Infect Dis* 8:ofab300. <https://doi.org/10.1093/ofid/ofab300>.
 23. Barton MI, MacGowan SA, Kutuzov MA, Dushek O, Barton GJ, van der Merwe PA. 2021. Effects of common mutations in the SARS-CoV-2 Spike RBD and its ligand, the human ACE2 receptor on binding affinity and kinetics. *Elife* 10. <https://doi.org/10.7554/eLife.70658>.
 24. Li Q, Nie J, Wu J, Zhang L, Ding R, Wang H, Zhang Y, Li T, Liu S, Zhang M, Zhao C, Liu H, Nie L, Qin H, Wang M, Lu Q, Li X, Liu J, Liang H, Shi Y, Shen Y, Xie L, Zhang L, Qu X, Xu W, Huang W, Wang Y. 2021. SARS-CoV-2 501Y.V2 variants lack higher infectivity but do have immune escape. *Cell* 184:2362–2371. <https://doi.org/10.1016/j.cell.2021.02.042>.
 25. Li X, Giorgi EE, Marichann MH, Foley B, Xiao C, Kong XP, et al. 2020. Emergence of SARS-CoV-2 through Recombination and Strong Purifying Selection. *bioRxiv*. <https://doi.org/10.1101/2020.03.20.000885>.
 26. Zhang L, Jackson CB, Mou H, Ojha A, Rangarajan ES, Izard T, et al. 2020. The D614G mutation in the SARS-CoV-2 spike protein reduces S1 shedding and increases infectivity. *bioRxiv* <https://doi.org/10.1101/2020.06.12.148726>.
 27. Yurkovetskiy L, Pascal KE, Tompkins-Tinch C, Nyalile T, Wang Y, Baum A, et al. 2020. SARS-CoV-2 Spike protein variant D614G increases infectivity and retains sensitivity to antibodies that target the receptor binding domain. *bioRxiv*.
 28. Daniloski Z, Guo X, Sanjana NE. 2020. The D614G mutation in SARS-CoV-2 Spike increases transduction of multiple human cell types. *bioRxiv* <https://doi.org/10.1101/2020.06.14.151357>.
 29. Kraemer MUG, Hill V, Ruis C, Dellicour S, Bajaj S, McCrone JT, Baele G, Parag KV, Battle AL, Gutierrez B, Jackson B, Colquhoun R, O'Toole Á, Klein B, Vespignani A, Volz E, Faria NR, Aanensen DM, Loman NJ, Du Plessis L, Cauchemez S, Rambaut A, Scarpino SV, Pybus OG, COVID-19 Genomics UK (COG-UK) Consortium. 2021. Spatiotemporal invasion dynamics of SARS-CoV-2 lineage B.1.1.7 emergence. *Science* 373:889–895. <https://doi.org/10.1126/science.abj0113>.
 30. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, Hinsley WR, Laydon DJ, Dabrera G, O'Toole Á, Amato R, Ragonnet-Cronin M, Harrison I, Jackson B, Ariani CV, Boyd O, Loman NJ, McCrone JT, Gonçalves S, Jorgensen D, Myers R, Hill V, Jackson DK, Gaythorpe K, Groves N, Sillitoe J, Kwiatkowski DP, Flaxman S, Ratmann O, Bhatt S, Hopkins S, Gandy A, Rambaut A, Ferguson NM, COVID-19 Genomics UK (COG-UK) consortium. 2021. Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7 in England. *Nature* 593:266–269. <https://doi.org/10.1038/s41586-021-03470-x>.
 31. Annavaiahala MK, Mohri H, Wang P, Nair M, Zucker JE, Sheng Z, et al. 2021. Emergence and expansion of SARS-CoV-2 B.1.526 after identification in New York. *Nature* <https://doi.org/10.1038/s41586-021-03908-2>.
 32. Everett J, Hokama P, Roche AM, Reddy S, Hwang Y, Kessler L, Glascock A, Li Y, Whelan JN, Weiss SR, Sherrill-Mix S, McCormick K, Whiteside SA, Graham-Wooten J, Khatib LA, Fitzgerald AS, Collman RG, Bushman F. 2021. SARS-CoV-2 Genomic Variation in Space and Time in Hospitalized Patients in Philadelphia. *mBio* 12. <https://doi.org/10.1128/mBio.03456-20>.
 33. Sherrill-Mix S, Hwang Y, Roche AM, Glascock A, Weiss SR, Li Y, Haddad L, Deraska P, Monahan C, Kromer A, Graham-Wooten J, Taylor LJ, Abella BS, Ganguly A, Collman RG, Van Duyn GD, Bushman FD. 2021. Detection of SARS-CoV-2 RNA using RT-LAMP and molecular beacons. *Genome Biol* 22:169. <https://doi.org/10.1186/s13059-021-02387-y>.
 34. Hilaire BGS, Durand NC, Mitra N, Pulido SG, Mahajan R, Blackburn A, et al. 2020. A rapid, low cost, and highly sensitive SARS-CoV-2 diagnostic based on whole genome sequencing. *Cochran Database Sys Rev* 8:1–129. <https://doi.org/10.1002/14651858.CD013705>.
 35. O'Toole Á, Scher E, Underwood A, Jackson B, Hill V, McCrone JT, Colquhoun R, Ruis C, Abu-Dahab K, Taylor B, Yeats C, Du Plessis L, Maloney D, Medd N, Attwood SW, Aanensen DM, Holmes EC, Pybus OG, Rambaut A. 2021. Assignment of epidemiological lineages in an emerging pandemic using the pangolin tool. *Virus Evol* 7:veab064. <https://doi.org/10.1093/ve/veab064>.
 36. Rambaut A, Holmes EC, O'Toole Á, Hill V, McCrone JT, Ruis C, Du Plessis L, Pybus OG. 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* 5:1403–1407. <https://doi.org/10.1038/s41564-020-0770-5>.
 37. Zahradník J, Marciano S, Shemesh M, Zoler E, Harari D, Chiaravalli J, Meyer B, Rudich Y, Li C, Marton I, Dym O, Elad N, Lewis MG, Andersen H, Gagne M, Seder RA, Douek DC, Schreiber G. 2021. SARS-CoV-2 variant prediction and antiviral drug design are enabled by RBD in vitro evolution. *Nat Microbiol* 6:1188–1198. <https://doi.org/10.1038/s41564-021-00954-4>.
 38. Zhu X, Mannar D, Srivastava SS, Berezuk AM, Demers J-P, Saville JW, Leopold K, Li W, Dimitrov DS, Tuttle KS, Zhou S, Chittori S, Subramaniam S. 2021. Cryo-electron microscopy structures of the N501Y SARS-CoV-2 spike protein in complex with ACE2 and 2 potent neutralizing antibodies. *PLoS Biol* 19:e3001237. <https://doi.org/10.1371/journal.pbio.3001237>.
 39. Farinholt T, Doddapaneni H, Qin X, Menon V, Meng Q, Metcalf G, et al. 2021. Transmission event of SARS-CoV-2 Delta variant reveals multiple vaccine breakthrough infections. *medRxiv*. <https://doi.org/10.1101/2021.06.28.21258780>.
 40. Kustin T, Harel N, Finkel U, Perchik S, Harari S, Tahor M, Caspi I, Levy R, Leshchinsky M, Ken Dror S, Bergeron G, Gadban H, Gadban F, Eliassian E, Shimron O, Saleh L, Ben-Zvi H, Keren Taraday E, Amichay D, Ben-Dor A, Sagas D, Strauss M, Shemer Avni Y, Huppert A, Kepten E, Balicer RD, Netzer D, Ben-Shachar S, Stern A. 2021. Evidence for increased breakthrough rates of SARS-CoV-2 variants of concern in BNT162b2-mRNA-vaccinated individuals. *Nat Med* 27:1379–1384. <https://doi.org/10.1038/s41591-021-01413-7>.
 41. Mlcochova P, Kemp S, Dhar MS, Papa G, Meng B, Ferreira I, et al. 2021. SARS-CoV-2 B.1.617.2 Delta variant replication and immune evasion. *Nature* <https://doi.org/10.1038/s41586-021-03944-y>.
 42. Lopez Bernal J, Andrews N, Gower C, Gallagher E, Simmons R, Thelwall S, Stowe J, Tessier E, Groves N, Dabrera G, Myers R, Campbell CNJ, Amirhalingam G, Edmunds M, Zambon M, Brown KE, Hopkins S, Chand M, Ramsay M. 2021. Effectiveness of Covid-19 Vaccines against the B.1.617.2 (Delta) Variant. *N Engl J Med* 385:585–594. <https://doi.org/10.1056/NEJMoa2108891>.
 43. McEwen AE, Cohen S, Bryson-Cahn C, Liu C, Pergam SA, Lynch J, et al. 2021. Variants of concern are overrepresented among post-vaccination breakthrough infections of SARS-CoV-2 in Washington State. *Clin Infect Dis* <https://doi.org/10.1093/cid/ciab581>.
 44. Martin DP, Weaver S, Tegally H, San JE, Shank SD, Wilkinson E, Lucaci AG, Gandhari J, Naidoo S, Pillay Y, Singh L, Lessells RJ, Gupta RK, Wertheim JO, Nekturenko A, Murrell B, Harkins GW, Lemey P, MacLean OA, Robertson DL, de Oliveira T, Kosakovsky Pond SL, COVID-19 Genomics UK (COG-UK). 2021. The emergence and ongoing convergent evolution of the SARS-CoV-2 N501Y lineages. *Cell* 184:5189–5200. <https://doi.org/10.1016/j.cell.2021.09.003>.
 45. McCallum M, De Marco A, Lempp FA, Tortorici MA, Pinto D, Walls AC, Beltramello M, Chen A, Liu Z, Zatta F, Zepeda S, di Iulio J, Bowen JE, Montiel-Ruiz M, Zhou J, Rosen LE, Bianchi S, Guarino B, Fregni CS, Abdelnabi R, Foo S-YC, Rothlauf PW, Bloyet L-M, Benigni F, Cameroni E, Neyts J, Riva A, Snell G, Telenti A, Whelan SPJ, Virgin HW, Corti D, Pizzuto MS, Veesler D. 2021. N-terminal domain antigenic mapping reveals a site of vulnerability for SARS-CoV-2. *Cell* 184:2332–2347. <https://doi.org/10.1016/j.cell.2021.03.028>.
 46. Syed AM, Taha TY, Khalid MM, Tabata T, Chen IP, Sreekumar B, et al. 2021. Rapid assessment of SARS-CoV-2 evolved variants using virus-like particles. *bioRxiv*.
 47. St Hilaire BG, Durand NC, Mitra N, Pulido SG, Mahajan R, Blackburn A, et al. 2020. A rapid, low cost, and highly sensitive SARS-CoV-2 diagnostic based on whole genome sequencing. *bioRxiv*.
 48. Li H, Durbin B. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25:1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
 49. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
 50. Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* 27:2987–2993. <https://doi.org/10.1093/bioinformatics/btr509>.
 51. Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M, et al. 2017. STAN: a Probabilistic Programming Language. *J Stat Software* 76:1–32.