

## RESEARCH ARTICLE

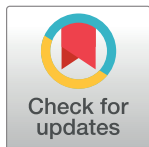
## Genome-wide association study between SARS-CoV-2 single nucleotide polymorphisms and virus copies during infections

Ke Li<sup>1,2\*</sup>, Chrispin Chaguza<sup>1,2</sup>, Julian Stamp<sup>3</sup>, Yi Ting Chew<sup>1,2</sup>, Nicholas F. G. Chen<sup>1</sup>, David Ferguson<sup>4,5</sup>, Sameer Pandya<sup>4,5</sup>, Nick Kerantzas<sup>4,5</sup>, Wade Schulz<sup>4,5</sup>, Yale SARS-CoV-2 Genomic Surveillance Initiative<sup>1</sup>, Anne M. Hahn<sup>1</sup>, C. Brandon Ogbunugafor<sup>2,6,7</sup>, Virginia E. Pitzer<sup>1,2</sup>, Lorin Crawford<sup>3,8,9</sup>, Daniel M. Weinberger<sup>1,2</sup>, Nathan D. Grubaugh<sup>1,2,6\*</sup>

**1** Department of Epidemiology of Microbial Diseases, Yale School of Public Health, New Haven, Connecticut, United States of America, **2** Public Health Modeling Unit, Yale School of Public Health, New Haven, Connecticut, United States of America, **3** Center for Computational Molecular Biology, Brown University, Providence, Rhode Island, United States of America, **4** Department of Laboratory Medicine, Yale School of Medicine, New Haven, Connecticut, United States of America, **5** Yale School of Medicine Biorepository, Yale University, New Haven, Connecticut, United States of America, **6** Department of Ecology and Evolutionary Biology, Yale University, New Haven, Connecticut, United States of America, **7** Santa Fe Institute, Santa Fe, New Mexico, United States of America, **8** Department of Biostatistics, Brown University, Providence, Rhode Island, United States of America, **9** Microsoft Research, Cambridge, Massachusetts, United States of America

**†** Authors of Yale SARS-CoV-2 Genomic Surveillance Initiative are listed in the Acknowledgments.

\* [ke.li.kl662@yale.edu](mailto:ke.li.kl662@yale.edu) (KL); [nathan.grubaugh@yale.edu](mailto:nathan.grubaugh@yale.edu) (NDG)



## OPEN ACCESS

**Citation:** Li K, Chaguza C, Stamp J, Chew YT, Chen NFG, Ferguson D, et al. (2024) Genome-wide association study between SARS-CoV-2 single nucleotide polymorphisms and virus copies during infections. *PLoS Comput Biol* 20(9): e1012469. <https://doi.org/10.1371/journal.pcbi.1012469>

**Editor:** Thomas Leitner, Los Alamos National Laboratory, UNITED STATES OF AMERICA

**Received:** March 3, 2024

**Accepted:** September 6, 2024

**Published:** September 17, 2024

**Copyright:** © 2024 Li et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** We used the R statistical software (v4.0.2) for all statistical analyses and visualization. Data and code used in this study are publicly available on Github: [https://github.com/grubaughlab/2024\\_paper\\_GWAS](https://github.com/grubaughlab/2024_paper_GWAS). All genome sequences used for the GWAS analysis and a subset of the associated metadata (accession number, virus name, collection date, originating lab and submitting lab, and the list of authors) in this dataset are published in GISAID's EpiCoV database: <https://doi.org/10.55876/gis8.240219fh>. The de-identified and coded clinical

## Abstract

Significant variations have been observed in viral copies generated during SARS-CoV-2 infections. However, the factors that impact viral copies and infection dynamics are not fully understood, and may be inherently dependent upon different viral and host factors. Here, we conducted virus whole genome sequencing and measured viral copies using RT-qPCR from 9,902 SARS-CoV-2 infections over a 2-year period to examine the impact of virus genetic variation on changes in viral copies adjusted for host age and vaccination status. Using a genome-wide association study (GWAS) approach, we identified multiple single-nucleotide polymorphisms (SNPs) corresponding to amino acid changes in the SARS-CoV-2 genome associated with variations in viral copies. We further applied a marginal epistasis test to detect interactions among SNPs and identified multiple pairs of substitutions located in the spike gene that have non-linear effects on viral copies. We also analyzed the temporal patterns and found that SNPs associated with increased viral copies were predominantly observed in Delta and Omicron BA.2/BA.4/BA.5/XBB infections, whereas those associated with decreased viral copies were only observed in infections with Omicron BA.1 variants. Our work showcases how GWAS can be a useful tool for probing phenotypes related to SNPs in viral genomes that are worth further exploration. We argue that this approach can be used more broadly across pathogens to characterize emerging variants and monitor therapeutic interventions.

metadata associated with the sequenced samples are available on Github: [https://github.com/grubaughlab/2024\\_paper\\_GWAS](https://github.com/grubaughlab/2024_paper_GWAS).

**Funding:** This project is supported by the Centers for Disease Control and Prevention (CDC) Broad Agency Announcement Contract 75D30122C14697 (to NDG). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** I have read the journal's policy and the authors of this manuscript have the following competing interests: NDG is a paid consultant for BioNTech, DMW has received consulting fees from Pfizer, Merck, and GSK, unrelated to this manuscript, and has been PI on research grants from Pfizer and Merck to Yale, unrelated to this manuscript.

## Author summary

Our study explores why viral load (copies measured by RT-qPCR) varies during SARS-CoV-2 infections by analyzing viral mutations and measuring viral copies in 9,902 individuals over two years. We aimed to understand how genetic differences in SARS-CoV-2 influence viral copies, considering host age and vaccination status. Using a genome-wide association study (GWAS), we identified several single-nucleotide polymorphisms (SNPs) in the virus linked to variations in viral levels. Notably, interactions between certain SNPs in the spike gene had non-linear effects on viral copies. Our analysis revealed that SNPs associated with higher viral copies were common in Delta and Omicron BA.2/BA.4/BA.5/XBB variants, while those linked to lower levels were mainly found in Omicron BA.1. This research highlights GWAS as a powerful tool for exploring virus genetics and suggests it can be broadly applied to monitor new variants of COVID-19 and other infectious diseases.

## Introduction

Continued SARS-CoV-2 transmission and evolution has propelled the COVID-19 pandemic. Peak viral replication in the upper respiratory tract occurs during the first few days of infection [1]. The viral load (or copies measured by RT-qPCR) in patient samples are valuable data to understand infection dynamics, such as inferring the likelihood of disease transmission [2]. However, it is challenging to use viral load data, and the challenge often arises from significant variations in viral load dynamics among sampled cases, which can be associated with 1) host heterogeneity, e.g., age [3] and vaccination status [4–6]; 2) distinct inherent properties of virus variants or sublineages [7], and 3) different sampling times [8]. For example, sampling during the early stages of infection may yield higher viral loads compared to later stages after viral replication has reached its peak. Nevertheless, the relative importance of these factors influencing viral load has not been completely explored [9,10].

Genome-wide association studies (GWAS) have emerged as a useful tool in the field of genetics, providing an approach to unraveling the complex interplay between genetic variations and observable traits, including diseases and drug resistance, as reviewed in [11]. Several studies have employed GWAS analysis to identify and investigate the association between human genetic variations across different individuals and the severity of COVID-19, shedding light on genetic variations that are related to severe infections [12–14]. However, few studies have utilized a GWAS method to study associations between the viral genome and viral traits [15–18]. The confluence of the extensive existing research on SARS-CoV-2 mutations and the millions of infections that have been sequenced provides us the opportunity to evaluate the application of GWAS for viral genomics. The hypothesis-free approach has the potential to enhance our understanding of genetic determinants influencing viral fitness and evolution and further inform effective public health strategies aimed at mitigating the spread and impact of SARS-CoV-2.

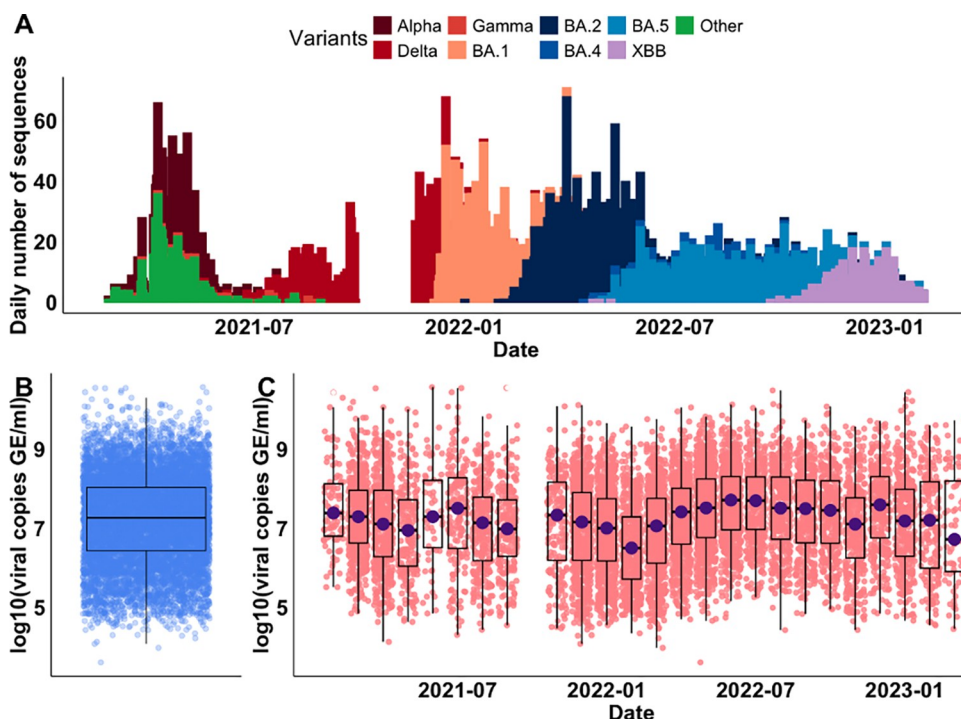
In this work, we aim to investigate the impact of intrinsic viral genetic substitutions (i.e., single nucleotide polymorphisms [SNPs]) on the changes in viral copies, adjusted for host age and vaccination status. For this, we apply a viral GWAS analysis to SARS-CoV-2 genomic sequencing and standardized RT-qPCR data collected from the Yale New Haven Hospital from February 2021 to March 2023. Using whole genome sequencing data on SARS-CoV-2 infections, along with relevant laboratory and patient metadata, we identify associations between viral SNPs and viral copies for different variants of concern (VOCs). We then

examine the temporal pattern of identified SNPs by constructing a phylogenetic tree, drawing upon subsamples, and analyzing the time series of the fraction of SNPs occurring in the sequences. This multifaceted analysis contributes to unraveling the complex dynamics of SARS-CoV-2 infections, providing valuable insights into the underlying viral SNPs that influence viral copies in different VOCs.

## Results

### Viral copies vary in SARS-CoV-2 infections

To better understand how SARS-CoV-2 viral load varies in infected individuals, we analyzed the viral copy data, along with associated host metadata (i.e., age and vaccination status), and genome sequencing data from a cohort of patients tested at the Yale New Haven Hospital (YNHH) located in Connecticut, US. We selected 9902 whole genome sequences with available viral copy data generated from remnant SARS-CoV-2 diagnostic samples over a 2-year period, from 03-Feb-2021 to 21-Mar-2023 (Fig 1A). The VOCs that we identified in our dataset during the sampling period included Alpha (n = 809), Delta (n = 1278), Gamma (n = 36), BA.1 (n = 1818), BA.2 (n = 2432), BA.4 (n = 293), BA.5 (n = 1992), XBB (n = 698), and the pre-VOC variant (named 'Other', n = 546). We conducted RT-qPCR using a standardized assay targeting the nucleocapsid (CDC 'N1' primers) for each sample to allow for cross-sample comparisons [19], except for a period during October 2021 when the PCR data were not generated. Across all samples, the viral copies, expressed as  $\log_{10}$ (viral copies per milliliter (Genome Equivalents/ml)), exhibited variations, ranging from 3.60 to 10.55, with a median value of 7.26



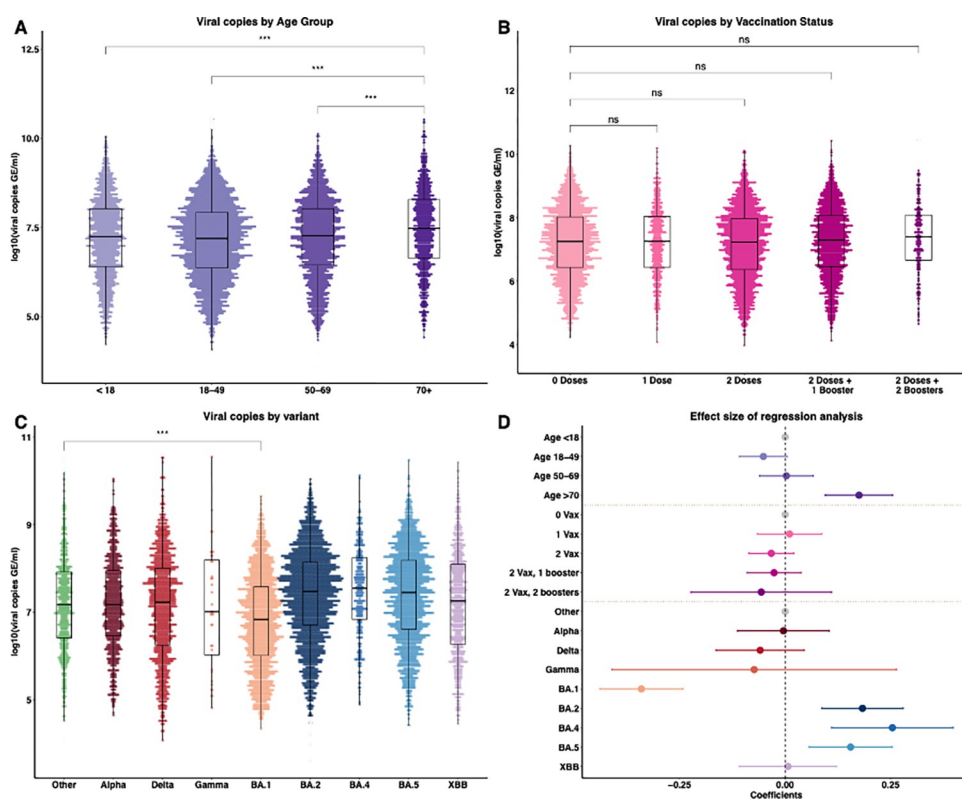
**Fig 1. Genomic sequences of SARS-CoV-2 infections and associated viral copies from cross-sectional samples collected in Connecticut, US.** (A) The daily number of genomic sequences of SARS-CoV-2 VOCs from February 2021 to March 2023. (B) The summary of viral copies of all samples, expressed as  $\log_{10}$ (viral copies per milliliter). (C) The summary of viral copies aggregated by month. The data gap in October 2021 is because we were unable to conduct PCR to obtain viral copies during this time.

<https://doi.org/10.1371/journal.pcbi.1012469.g001>

(Fig 1B). The variations in viral copies could be attributed either to the introduction and/or replacement of different VOCs, each with its own epidemic curve, or to the stochasticity from the sampling process. To reduce stochastic effects, we aggregated the viral copies by month and still observed large variations in the viral copies across the months, albeit with no consistent trend (Fig 1C). Notably, we observed the lowest median value of viral copies (median = 6.49) in February 2022, during which 96.3% of the sampled sequences tested positive for BA.1 infections. By contrast, we observed the highest median value of viral copies (median = 7.70) in June 2022, during which the sampled sequences tested positive for BA.2 (64.9%), BA.4 (6%), or BA.5 (29.1%) infections. Taken together, we showed a wide range of viral copies in the sampled SARS-CoV-2 infections with different VOCs, utilizing data from genomic surveillance and standardized RT-qPCR tests.

## Viral copies correlate with age and variants, but not with vaccination status

Having uncovered a large variability in the observed viral copies from the samples, we next assessed the factors associated with these changes. To do this, we first summarized and compared viral copies in various age groups (Fig 2A). A positive correlation has been previously reported between age and SARS-CoV-2 viral copies, showing that younger age groups had



**Fig 2. Viral copies by category and regression analysis results.** Comparison of viral copies stratified by (A) age groups, (B) vaccination statuses, (C) variant of concerns. (D) Association of age, vaccination status, and VOCs with viral copies, expressed as  $\log_{10}(\text{viral copies per milliliter (Genome Equivalents/ml)})$ . The reference groups (in gray) are Age <18 years old, 0 doses of vaccination, and the Other variant, respectively. The positive coefficients indicate the covariate is associated with higher viral copies value compared to the reference group, and vice versa. 0 Vax, 1 Vax, 2 Vax, 2 Vax 1 booster, and 2 vax 2 boosters denote vaccination statuses of 0 doses, 1 dose, 2 doses, 2 doses, and 1 booster, and 2 doses and 2 boosters, respectively, corresponding to the labels in (B). Results are shown as means with 95% confidence intervals. \*\*\*  $p < 0.001$ .

<https://doi.org/10.1371/journal.pcbi.1012469.g002>

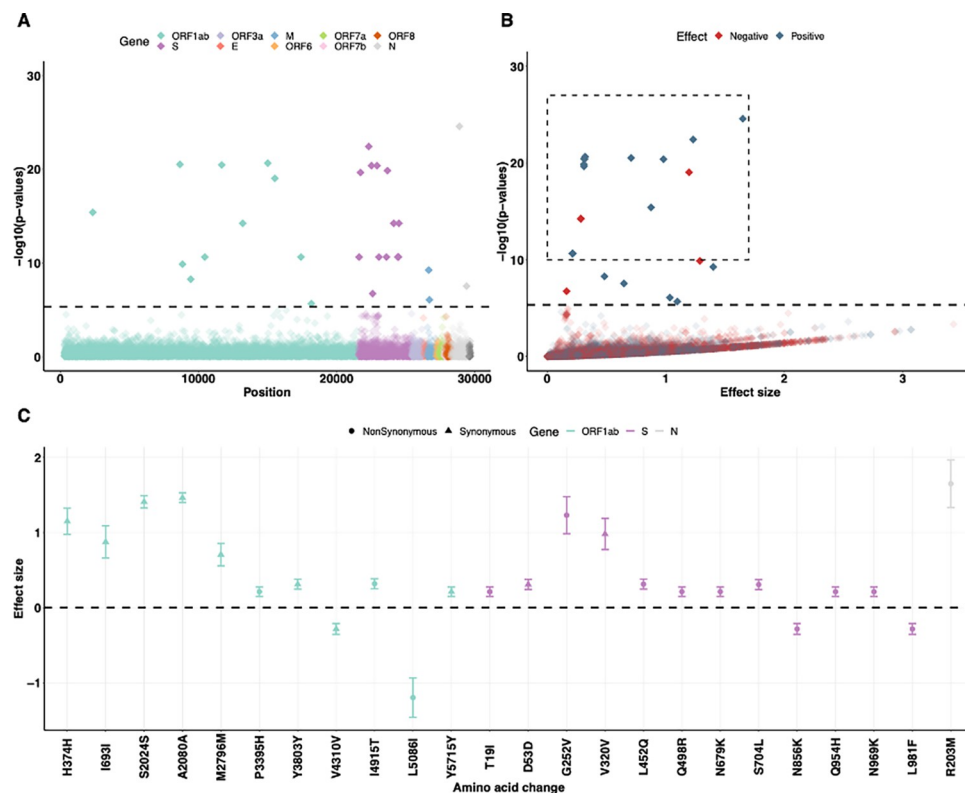
lower viral copies independent of gender and/or symptom duration [20]. We observed a similar result in our dataset and found that the oldest age group (i.e., >70 years old) had the highest viral copies compared with other age groups (mean = 7.47, 95% confidence interval (CI): [5.12, 9.49],  $p < 0.001$ , Wilcoxon signed-rank test). For the effect of vaccination on viral copies, some studies have demonstrated that although vaccination reduced the risk of infections with the Delta variant, no significant difference in peak viral copies was found between fully vaccinated and unvaccinated individuals [4,5,21]. In contrast, other studies have shown that vaccination reduced viral copies in BA.1 infections among boosted individuals compared to unvaccinated ones [6]. These results suggest the effect of vaccination on viral copies may depend on the characteristics of the infecting SARS-CoV-2 variant. We compared viral copies among groups with different vaccination statuses to assess the impact of vaccination on viral copies (Fig 2B), and no statistically significant differences were detected between the groups in our data ( $p > 0.05$ , Wilcoxon signed-rank test). Finally, we compared viral copies stratified by variant category (Fig 2C). Combining samples collected from all age and vaccination status groups for each variant, we found that the overall mean values of viral copies were lowest for infections with BA.1 (mean = 6.83, 95% CI: [4.87, 8.87],  $p < 0.001$ , Wilcoxon signed-rank test) compared to infections with other all non-BA.1 variants.

Since several factors may simultaneously impact the SARS-CoV-2 viral load, next, we sought to quantify the combined impact of age, vaccination status, and VOCs on the observed viral copies. To achieve this, we fitted a multivariate linear regression model, with viral copies as the outcome variable and age, vaccination, and VOCs as covariates (Fig 2D). We found that the older age group (i.e., age >70 years old) had a positive association with viral copies (mean = 0.17, 95% CI: [0.09, 0.25],  $p < 0.001$ ) compared with the reference group (i.e., age <18 years old). We also found that vaccination status was not associated with viral copies (i.e., 95% CIs of the vaccination coefficients span 0,  $p > 0.05$ ). Notably, we showed that infections with BA.1 were associated with reduced viral copies, with a mean effect size of -0.34 (95% CI: [-0.44, -0.24],  $p < 0.001$ ) in the same age group and vaccination status, compared to the Other variant. We also showed that infections with BA.2 (mean = 0.19, 95% CI: [0.09, 0.28],  $p < 0.001$ ), BA.4, or BA.5 (mean = 0.17, 95% CI: [0.07, 0.26],  $p < 0.001$ ) were associated with increased viral copies. Among them, infections with BA.4 were associated with the largest positive effect size (mean = 0.27, 95% CI: [0.12, 0.41],  $p < 0.001$ ). Our findings demonstrated that variations in viral copies were associated with infections caused by different SARS-CoV-2 variants and the older age group. This implies that intrinsic factors of the viruses, such as genetic mutations among distinct VOCs, are key determinants impacting viral copies.

### Viral GWAS reveals SARS-CoV-2 SNPs associated with viral copies

Having demonstrated that changes in SARS-CoV-2 viral copies are associated with infections caused by different viral variants or strains, especially Omicron BA.1/BA.2/BA.4/BA.5 variants (Fig 2D), we then sought to identify potential genetic mutations—specifically, SNPs—that contributed to these changes in viral copies. For this, we performed a GWAS analysis using high-quality genome sequences (i.e., genome coverage > 95%). We conducted whole-genome sequencing on the 9902 SARS-CoV-2 positive specimens collected from February 2021 to March 2023. Firstly, using Wuhan-Hu-1 (GenBank MN908937.3) as the reference genome, we identified 10,697 SNPs for further testing associated with viral copies as covariates. We then checked for the population structure of the 9902 genome sequences using a multidimensional scaling (MDS) method [22] (S1 Fig). We observed that Delta was an outgroup to other pre-Omicron variants (i.e., pre-VOC variant (Other), Alpha, and Gamma), and BA.1 was an outgroup to the BA.2/BA.4/BA.5/XBB cluster. In our model, we included the inferred four clusters





**Fig 3. GWAS analysis identifies several single nucleotide polymorphisms (SNPs) that are associated with the changes in viral copies.** (A) Genome-wide association results of the impact of identified SNPs on viral copies during SARS-CoV-2 infection. The dashed line indicates the permuted threshold for genome-wide significance  $p = 4.67 \times 10^{-6}$  (0.05/10697 SNPs). Significant SNPs are shown with solid colors. (B) SNPs (with  $p < 1 \times 10^{-10}$ ) that have positive (blue) or negative (red) effects on viral copies. (C) The corresponding synonymous (triangles) and non-synonymous (circles) amino acid changes that associate with increased or decreased viral copies. Data shown as means with 95% confidence intervals. The estimated effective sizes and associated standard deviations are given [S1 Table](#).

<https://doi.org/10.1371/journal.pcbi.1012469.g003>

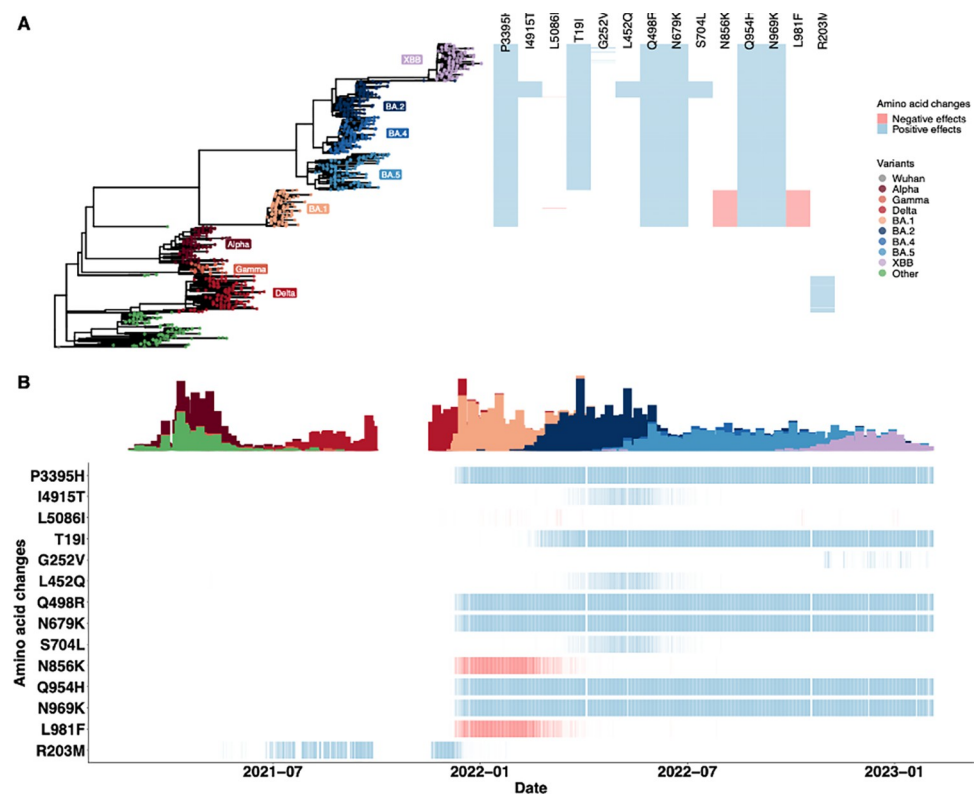
based on the MDS-computed distance to capture the viral population structure. Clusters were defined using a k-means clustering method ([S1 Fig](#)). The host ages and vaccination status were also included in the model as covariates.

Using the linear regression model on viral copies for each SNP, adjusted for viral population structure and host factors, we identified 31 SNPs exceeding the permuted threshold for genome-wide significance ( $p = 4.67 \times 10^{-6}$ , dashed line, [Fig 3A](#)). The threshold value was calculated as 0.05 divided by 10,697 SNPs [23]. We found that the observed distribution of  $p$ -values closely matches the expected distribution under the null hypothesis of no association ([S2A Fig](#)). To ascertain whether those SNPs have a negative or positive impact on viral copies and evaluate their effect size, we extracted the coefficients ( $\beta$ ) of the SNPs with  $p < 1 \times 10^{-10}$  and their standard deviations ( $\sigma$ ) from the regression model (dashed box, [Fig 3B](#)). We then annotated the SNPs to identify the associated amino acids, and among them, 14 SNPs were non-synonymous (i.e., changed the amino acid; [Fig 3C](#)). We found that a non-synonymous change N:R203M, located on the N gene, had the most significant association with increased viral copies ( $p = 2.68 \times 10^{-22}$ ,  $\beta = 1.65$ ,  $\sigma = 0.16$ ). By contrast, the amino acid change most strongly associated with a negative effect on viral copies was ORF1ab:L5086I ( $p = 9.20 \times 10^{-20}$ ,  $\beta = -1.20$ ,  $\sigma = 0.13$ ). We further conducted a marginal epistasis test [24–26] to detect the epistatic effects of SNPs on viral copies. We discovered multiple pairs of SNPs that exhibit positive epistatic effects on viral copies, with most interactions occurring in the S gene ([S3 Fig](#)).

To assess the impact of adjusting for the population structure of the SARS-CoV-2 strains using the MDS components on the regression results, we conducted a sensitivity analysis on the genome sequences using the inferred MDS components from the pairwise SNP distance matrix of SARS-CoV-2 sequences as covariates. By doing this, we identified 113 SNPs exceeding the permuted threshold (S4 Fig). The observed distribution of  $p$ -values also closely matched the expected distribution under the null hypothesis of no association (S2B Fig). The results may be more likely to reflect the SNPs that influence the viral copies dependent on lineage. We also examined the association between viral copies and SNPs after adjusting for the population structure based on the VOCs themselves, which broadly correspond to the identified sequence clusters. We showed that only a few SNPs were found (S5–S8 Figs), mostly within the Omicron BA.2/BA.4/BA.5/XBB cluster (S8 Fig).

### The impact of amino acid changes on viral copies is dependent on the variant

Having identified the 14 non-synonymous SNPs with statistically significant effects on viral copies in our primary analysis, we next sought to understand the temporal patterns of the emergence of these amino acid changes (Fig 4). To investigate the clustering of these SNPs, we



**Fig 4. The temporal dynamics of non-synonymous amino acid changes in the ORF1ab gene (P3395H, I4915T and L5086I), S gene (T19I, G252V, L452Q, Q498R, N679K, S704L, N856K, Q954H, N969K and L981F), and N gene (R203M) associated with changes in viral copies.** The results are based on the multivariate regression analysis using the sequence clusters (i.e., a categorical variable) inferred from the MDS components. (A) The phylogenetic tree estimated from a representative set of 996 genome sequences showing variant assignments and the locations of amino acid changes that increase (blue) or decrease (red) viral copies. (B) The temporal dynamics of the SNPs from February 2021 to March 2023. The transparency of the color corresponds to the mutation fraction in the daily sequence count: transparent color indicates low fractions, and opaque color indicates high fractions. The temporal dynamics of the SNPs, using MDS-inferred distance as a population control, are shown in S9–S11 Figs.

<https://doi.org/10.1371/journal.pcbi.1012469.g004>

randomly sampled approximately 120 genome sequences from each VOC category (only 36 sequences were available for Gamma in our dataset) and generated a phylogenetic tree drawing upon the subsamples (Fig 4A). We found a clear pattern in how these mutations emerged by VOC (Fig 4A heatmap). We found that all amino acid changes associated with a positive effect on viral copies were found in Delta and Omicron BA.2/BA.4/BA.5/XBB infections. Often, more than one amino acid change was observed in each sampled sequence, suggesting genetic linkage between these SNPs, as also shown in the epistasis test (S3 Fig), such as S:Q954H and N969K. In particular, we identified that the amino acid changes S:L452Q ( $p = 3.91 \times 10^{-25}$ ,  $\beta = 0.34$ ,  $\sigma = 0.03$ ) and S704L ( $p = 1.35 \times 10^{-24}$ ,  $\beta = 0.34$ ,  $\sigma = 0.03$ ) associated with a positive effect on viral copies were typically observed in combination with BA.2 infections—specifically, lineage BA.2.12.1. We also observed that the amino acid changes with negative effects on viral copies (ORF1ab:L5086I, S:N856K and L981F) were only associated with BA.1 infections.

To explore the temporal dynamics of these amino acid changes, we calculated the fraction of SNPs occurring in the sequences for each day, thereby accounting for the number of introductions to the population (Fig 4B). We observed most SNPs with a positive impact on viral copies emerging in sequences sampled from February 2022, when BA.2 was first detected in Connecticut. These SNPs were consistently observed in almost every sequence thereafter. By contrast, we found that the other two amino acid changes (S:L452Q and S704L) that had a positive effect on viral copies were only in the samples from BA.2 infections and did not arise again in sublineages of BA.4 or BA.5. S:G252V was associated with higher viral copies; however, we found that the SNP only appeared in a few sequences associated with XBB infections. Notably, the N:R203M mutation was only associated with Delta infections. For the SNPs (ORF1ab:L5086I, S:N856K and L981F) that had a negative association with viral copies, we observed that they were present in samples associated with BA.1 infections and did not persist when BA.1 was replaced by BA.2.

## Discussion

We conducted a GWAS analysis on 9,902 high-quality SARS-CoV-2 genome sequences generated from two years of genomic surveillance in Connecticut, US to identify and evaluate SNPs that were associated with variations in viral copies during infections. Using a GWAS approach, we were able to identify and examine virus-related factors that were associated with the observed variations in viral copies independent of host factors. This was achieved by combining data from a large cohort of individuals infected with different VOCs and employing a regression model for viral copies that accounted for virus-level factors (i.e., specific SNPs and genetic background), adjusted for individual factors (i.e., age and vaccination status). We identified several SNPs corresponding to non-synonymous amino acid changes in the SARS-CoV-2 genome that were individually or jointly associated with the variations in viral copies. In particular, temporal patterns of the SNPs revealed that SNPs associated with increased viral copies were predominantly observed in Delta and Omicron BA.2/BA.4/BA.5/XBB infections, whereas those associated with decreased viral copies were mostly observed in infections with Omicron BA.1 variants.

Using a GWAS approach, we successfully identified a subset of variant-defining amino acid changes in Delta and Omicron variants (S12 Fig). Note that we did not detect any substitutions in the Alpha and Gamma variants (likely due to the low sample size for Gamma). We also identified SNPs that did not define any major variant category, including S:L452Q and S704L that were specifically associated with BA.2.12.1, a sublineage of BA.2 that briefly dominated during the pandemic (i.e., dominated mainly in the US between March and May 2022). This highlights the application of GWAS for identifying SNPs associated with important



phenotypic effects without requiring a set of lineage-defining mutations to be defined a priori. Nevertheless, there are several reasons why we only detected a subset of the SNPs that defined different VOCs. Firstly, SNPs with small effect sizes may not be detected due to the stringent statistical significance thresholds applied in GWAS. Secondly, lineage-defining SNPs that are in low linkage disequilibrium with the causal mutations may not be detected [27], even if they may be functionally relevant. Our results showcase how GWAS can help to narrow the focus of SNPs associated with specific phenotypes, generating hypotheses for further investigation.

A key result from our analysis is that SNPs associated with viral copies did not exhibit the same temporal dynamics, even though they could have similar (either positive or negative) effects on viral copies, suggesting they may have independent effects on viral copies. Some amino acid changes, for example, ORF1ab:I4915T (positive effects), were only present in samples with BA.2 infections and disappeared when new Omicron variants emerged. Other SNPs (e.g., S:T19R), while also associated with higher viral copies, were observed and persisted in all BA.2/BA.4/BA.5/XBB infections. The distinct temporal pattern of SNPs, dependent on VOCs, may help explain the different fitness levels (e.g., intrinsic transmissibility or immune escape) of each variant [28,29]. Notably, we found three SNPs, ORF1ab:L5086I, S:N856K and L981F, were associated with decreased viral copies in BA.1 infections. The negative impact on viral copies should be interpreted with caution. Although the possibility that these SNPs have a direct impact on reducing viral copies cannot be ruled out, it is also likely that the estimated negative effects are due to a synthetic association with other SNPs. Further study may be required to disentangle the direct effects of these SNPs from the confounding influences of other genetic variations and to confirm their functional impact on viral copies.

In this study, we employed a series of single SNP regression models to identify the underlying SNPs associated with the changes in viral copies without accounting for potential interactions between SNPs. We noted that several synonymous SNPs located in the ORF1ab gene were identified to have an impact on viral copies. The synonymous SNPs were likely linked to non-synonymous SNPs that were under positive selection. In such cases, the synonymous SNPs can be carried along with the non-synonymous SNPs, resulting in their significance in the GWAS analysis, as shown in the subsequent epistasis test (S3 Fig). The epistasis test provided evidence that the synonymous mutations identified through GWAS analysis are likely the result of synthetic associations with other non-synonymous mutations. Nevertheless, the method provided an initial set of SNPs that are worth further exploration, pinpointing important mutations associated with viral copies and providing valuable insights into the overall genetic landscape of the viral population. The method, thus, represented an important first step towards understanding detailed epistatic effects among these mutations on viral copies. A paired or higher-order SNP regression study could be conducted as a subsequent step to test potential interactions or joint effects among different SNPs.

There are limitations to our study. First, we assumed that the distribution of times between infection and sample collection was similar through time and across variants as these data were not available. Given our samples were taken frequently over a 2-year period, we do not anticipate that this assumption will qualitatively impact our results. Second, our study primarily focuses on the genetic variants in VOCs, neglecting other factors such as host immune responses or environmental influences, partially captured by the host-associated covariates, including age and vaccination status in this study, that may also contribute to the changes in viral copies. Further study will be needed to address the impact of these factors on viral copies, for example, genome-to-genome analysis to reveal the impact of host-viral genetic interactions in SARS-CoV-2 infections [18,30]. Third, our data were obtained from a specific geographic region, whose population diversity may not necessarily be similar to other settings; therefore, extrapolating these findings to a broader population may require caution. Additionally,

focusing solely on consensus genomic changes in the analysis could overlook the genetic diversity within the sample, which may also influence variations in viral load. Despite these constraints, our study highlights the importance of sustained genomic surveillance and the need for comprehensive analyses to understand the nuanced impact of specific genetic variations on viral copies at the within-host level, and its implications for viral transmissibility and immune escape at the population level. Further work and collaborative efforts are essential to elucidate the complex interplay between viral genetics, host factors, and the dynamics of transmission associated with emerging variants. Such studies could inform predictive early warning public health systems regarding the emergence of potentially highly transmissible viral strains based on their constellation of mutations.

Recently, Duesterwald et al. [10] used genome sequence data and a machine-learning approach to predict cycle threshold (Ct) values of SARS-CoV-2 infections based on the  $k$ -mers. Similar to our findings, they suggested that S:L452 and P681 were hallmarks of VOCs, implying impacts on the observed Ct values in clinical samples. Although the machine-learning approach may capture broader patterns and interactions within the genome on Ct values, they lack interpretability compared to regression models. For example, regression-based models could offer insights into the direct association between specific genetic variants and viral copies. In addition, regression-based models may perform well even with limited sample sizes [16], provided that the assumptions of the model are met and the predictors are informative, whereas using machine-learning methods with small sample sizes can be challenging. However, the viral GWAS method may not be appropriate in situations where there is insufficient genetic diversity in the viral population under study, as this can limit the power to detect meaningful associations between mutations and viral traits. Additionally, it may not be suitable when the phenotypic traits of interest are not well-defined or accurately measured.

With the availability of high-quality whole-genome sequences for SARS-CoV-2, we demonstrated that GWAS analysis of the viral genome can identify SNPs that associate with positive or negative impacts on viral copies in VOCs, revealing important biological insights and enhancing our understanding of within-host viral dynamics. We argue that the application of GWAS analyses to study viral genomes provides a particularly tractable tool to identify potential SNPs of interest for further evaluation across different viral pathogens. It is particularly useful to understand the genetic basis of viral virulence, transmission, resistance to antiviral treatments, and host-virus interactions for several reasons. First, the small genome size of viruses and high evolutionary rates make it easier to perform comprehensive genome-wide scans for SNPs and to experimentally test the impacts of SNPs on specific traits. Second, significant phenotypic variations (e.g., viral loads and antibody responses) are often observed in viral infections, despite limited changes in the viral genome. GWAS can help to identify SNPs that correlate with these phenotypic variations, providing insights into the genetic basis of these traits. Third, the increasing accessibility to sequence viral genomes makes it possible to perform GWAS on rich datasets, enabling in-depth analysis of the temporal dynamics of viral evolution. Together, the applicability of GWAS analyses to study viral genomes can provide a new approach for exploring the intricate interplay between genetic mutations and phenotypes, informing strategies for managing and mitigating the impact of emerging viral variants, and contributing to the development of potential therapeutic interventions.

## Materials & methods

### Ethics statement

The Institutional Review Board from the Yale University Human Research Protection Program determined that the RT-qPCR testing and sequencing of de-identified remnant COVID-

19 clinical samples obtained from clinical partners conducted in this study is not research involving human subjects (IRB Protocol ID: 2000028599).

### Clinical sample collection and measurement of viral copies by RT-qPCR

SARS-CoV-2 positive samples (nasal swabs in viral transport media) were collected through the Yale New Haven Hospital (YNHH) System as a part of routine inpatient and outpatient testing and sent to the Yale SARS-CoV-2 Genomic Surveillance Initiative. Using the MagMAX viral/pathogen nucleic acid isolation kit, nucleic acid was extracted from 300 $\mu$ l of each clinical sample and eluted into 75 $\mu$ l of elution buffer. Extracted nucleic acid was then used as template for a “research use only” (RUO) RT-qPCR assay [19] to test for presence of SARS-CoV-2 RNA. Ct values from the nucleocapsid target (CDC-N1 primer-probe set [31]) were used to derive viral copy numbers using a previously determined standard curve for this primer set [32]. A positive RNA control with defined viral copy number (1000/ $\mu$ l) was used to standardize results across individual runs.

### Whole genome sequencing

Libraries were prepared for sequencing using the Illumina COVIDSeq Test (RUO version) and quantified using the Qubit High Sensitivity dsDNA kit. Negative controls were included for RNA extraction, cDNA synthesis, and amplicon generation. Prepared libraries were sequenced at the Yale Center for Genomic Analysis on the Illumina NovaSeq with a 2x150 approach and at least 1 million reads per sample.

Reads were then aligned to the Wuhan-Hu-1 reference genome (GenBank MN908937.3) using BWA-MEM v.0.7.15 [33]. Adaptor sequences were then trimmed, primer sequences masked, and consensus genomes called (simple majority >60% frequency) using iVar v1.3.133 [34] and SAMtools v1.11 [35]. When <20 reads were present at a site an ambiguous “N” was used, with negative controls consisting of  $\geq$ 99% Ns. The Pangolin lineage assignment tool [36] was used for assigning viral lineages.

### Clinical metadata

We obtained patient metadata and vaccination records from the YNHH system and the Center for Outcomes Research and Evaluation (CORE) and matched these records to sequencing data through unique sample identifiers. Duplicate patient records or those with missing or inconsistent metadata and vaccination date were removed from the GWAS analysis. We also removed patient records with persistent infections (>28 days since first positive test).

We determined vaccination status at time of infection by comparing the sample collection date to the patient’s vaccination record dates. We categorized vaccine statuses based on the number of vaccine doses received at least 14 days before the collection date. Patient vaccination statuses at the time of infection were categorized as follows: non-vaccine, one-dose vaccine, two-dose vaccine, two-dose vaccine with one booster, or two-dose vaccine with two boosters. We calculated the age of each patient as the difference between the date of birth and the sampling date.

### Single nucleotide polymorphisms

To identify single nucleotide polymorphisms (SNPs), we first aligned the 9902 genome sequences using *nextalign* (v3.2.1) [37] with the reference genome of the Wuhan-Hu-1 genome (GenBank accession: MN908937.3). Then, SNPs were identified using *snp-sites* (v2.4.1) [38], with the reference genome of the Wuhan-Hu-1 genome (GenBank accession:

MN908937.3). We also normalized the SNPs in the generated VCF file, such that multiallelic SNPs were separated into different rows. Normalizing the SNPs ensured that each SNP was one-hot encoded and analyzed separately. Note that we did not include ambiguous SNPs, deletions and insertions in our GWAS analysis. We used *vcf-annotator* (v0.7) to annotate SNPs to corresponding amino acid changes.

### Multidimensional scaling and population control

To reveal the underlying structure of the 9902 genome sequences. We first used *snp-dists* (v0.7.0) [39] to convert the aligned sequences (a FASTA alignment) to a SNP distance matrix. We then applied a multidimensional scaling (MDS) method [22] to transform the SNP distance matrix into a geometric configuration while preserving the original pairwise relationships. The scaling was conducted using *cmdscale* function in an R package *stats* (v3.6.2). We set the maximal dimensional parameter  $k = 2$ .

To measure the goodness of the transformation, we calculated the distance between the original genome sequencing data and compared it with the new distances determined by MDS. This involved arranging the two matrices of distances into two columns and computing the correlation coefficient (i.e.,  $r$ ) between them. Finally, we used  $r^2$  to measure the proportion of variance in the original distance matrix explained by the new computed distance matrix.

To determine the clusters (i.e., categorical variables) from MDS, we applied the k-means clustering method using the *kmeans* function implemented in R statistical software (v4.0.2). We set the number of centroids  $k = 4$ .

### Testing for associations between viral copies and SNPs

In this work, we conducted a series of single SNP regression analyses to test for associations between viral copies and SNPs, adjusted for host ages, vaccination status and viral population. The linear regression model is written as follows:

$$Y \sim \alpha W + \beta_i \text{SNP}_i + e, \quad (1)$$

where  $Y$  is a vector of normalized  $\log_{10}$ -transformed viral copies,  $W$  is a matrix of covariates, including age (a categorical variable with four age groups of “<18”, “18–49”, “50–69”, and “>70” years old), vaccination status (a categorical variable with vaccination statuses of “0 doses”, “1 dose”, “2 doses”, “2 doses and 1 booster”, “2 doses and 2 boosters”), a population control variable for different viral variants (a categorical variable with cluster numbers of “1”, “2”, “3” and “4”, see S1 Fig for detailed clusters), and an intercept, and  $\alpha$  is a vector that corresponds to coefficients of the covariates. In particular,  $\text{SNP}_i$  is a vector of genotype values for all samples at each SNP,  $i$ . It is a binary variable: 0 represents the SNP is not present in the genome sequence, whereas 1 represents its presence.  $\beta_i$  is the effective size of each identified SNP,  $i$ . We also conducted a sensitivity analysis including two terms  $\xi_1 d_1$ ,  $\xi_2 d_2$  as covariates in the model for population control. The vectors  $d_1$ ,  $d_2$  represent the two dimensions computed by MDS, and  $\xi_1$ ,  $\xi_2$  are the coefficients of the dimension covariates. The random effect of residual errors is presented here by  $e$ , which is assumed to follow a normal distribution with a mean of 0 and a standard deviation of  $\sigma_e$ .

### Marginal epistasis test

We applied the marginal epistasis test method to explore the interactions between SNPs on viral copies, using an R package *mvMAPIT* (v.2.0.3) [24–26]. This method maps SNPs with non-zero marginal epistatic effects—the combined pairwise interaction effects between a given

SNP and all other SNPs—identifying candidate variants involved in epistasis without needing to identify the exact partners with which the variants interact.

The method works by examining one variant at a time. For the  $j$ -th variant, the following linear model is applied,

$$\mathbf{y} = \mathbf{W}\boldsymbol{\gamma} + \mathbf{x}_j\beta_j + \mathbf{m}_j + \mathbf{g}_j + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim \text{MVN}(0, \tau^2 \mathbf{I}), \quad (2)$$

where  $\mathbf{y}$  is an  $n$ -vector of phenotypes (i.e., viral copies) for  $n$  individual samples;  $\mathbf{W}$  is a matrix of covariates including an intercept term with effects  $\boldsymbol{\gamma}$ ;  $\mathbf{x}_j$  is an  $n$ -vector for the  $j$ -th variant (i.e., SNP) that is focus of the model;  $\beta_j$  is the corresponding additive effect size for the  $j$ -th variant;  $\mathbf{m}_j = \sum_{l \neq j} \mathbf{x}_l \beta_l$  is the combined additive effects from all other variants, and effectively represents the additive effect of the  $j$ -th SNP under the polygenic background of all other SNPs; and  $\mathbf{g}_j = \sum_{l \neq j} (\mathbf{x}_j \circ \mathbf{x}_l) \alpha_l$  is the summation of all pairwise interaction effects between the  $j$ -th SNP and all other SNPs. Lastly,  $\boldsymbol{\varepsilon}$  is an  $n$ -dimensional vector of residual errors where MVN denotes a multivariate normal distribution,  $\tau^2$  is a variance term, and  $\mathbf{I}$  denotes an identity matrix.

To ensure model identifiability, MAPIT assumes that the additive and interaction effect sizes follow univariate normal distributions where  $\beta_l \sim N(0, \omega^2/(J-1))$  and  $\alpha_l \sim N(0, \sigma^2/(J-1))$  where  $J$  denotes the total number of variants in the dataset. This key assumption on the regression coefficients means that the two random effects can also be expressed probabilistically as: (i)  $\mathbf{m}_j \sim \text{MVN}(0, \omega^2 \mathbf{K}_j)$  where  $\mathbf{K}_j = \mathbf{X}_{-j} \mathbf{X}_{-j}^T / (J-1)$  is an additive genetic relatedness matrix that is computed using all genotypes other than the  $j$ -th SNP; and (ii)  $\mathbf{g}_j \sim \text{MVN}(0, \sigma^2 \mathbf{G}_j)$  where  $\mathbf{G}_j = \mathbf{D}_j \mathbf{K}_j \mathbf{D}_j$  is a non-additive relatedness matrix computed based on all pairwise interaction terms involving the  $j$ -th SNP and  $\mathbf{D}_j = \text{diag}(\mathbf{x}_j)$  denotes a diagonal matrix with the  $j$ -th genotype as its only nonzero elements.

The key takeaway from MAPIT is that the variance component  $\sigma^2$  represents a measure of the marginal epistatic effect for each SNP in the data. Therefore, to identify variants that have significant nonzero marginal epistatic effects, the model assesses the null hypothesis  $H_0: \sigma^2 = 0$  for each variant in the data set. The mvMAPIT software uses a method of moments algorithm to estimate model parameters and then uses a calibrated two-sided z-score (i.e., normal) test to derive  $p$ -values.

## Phylogenetic tree construction and comparison to variant-defining substitutions

We employed *iq-tree* (v2.2.2.6) [40] of a representative set using 996 of our 9902 genome sequences for tree construction, using Wuhan-Hu-1 (GenBank MN908937.3) as the reference genome. We specified the HKY substitution model and set the number of bootstrap replicates to 1,000. To visualize the phylogenetic tree, we used the *ggtree* (v1.4.11) implemented in the R statistical software (v4.0.2). The variant-defining amino acid changes were defined as those mutations with >75% prevalence in at least one lineage, as estimated on [outbreak.info](https://outbreak.info) website [41]. Note that we did not include deletions in variant-defining substitutions.

## Supporting information

**S1 Fig. Results of multidimensional scaling. The population structure of the 9902 genome sequences using a multidimensional scaling (MDS) method.** Clusters are defined using a  $k$ -means clustering method, as demonstrated on the bottom right corner.

(TIF)



**S2 Fig. Q-Q plots of GWAS p-values.** Q-Q plots (quantile-quantile plots) showing the p-values from GWAS analysis using (A) two MDS-computed components, or (B) MDS-inferred four clusters as covariates in the regression model.

(TIF)

**S3 Fig. Marginal epistasis tests identify single nucleotide polymorphisms (SNPs) that have epistatic interactions with others and are associated with the changes in viral copies.** (A) Marginal epistasis test results of the SNPs (annotated as amino acid changes) that have marginal epistatic effects on viral copies. The dashed line indicates the permuted threshold for genome-wide significance  $p = 0.05/171 = 2.74 \times 10^{-4}$ . Significant mutations are shown with solid colors. (B) The p-values and (C) the effect size of pairwise interaction tests among the significant mutations.

(TIF)

**S4 Fig. GWAS analysis identifies several single nucleotide polymorphisms (SNPs) that are associated with the changes in viral copies.** (A) Genome-wide association results of the impact of identified SNPs on viral copies during SARS-CoV-2 infection. The dashed line indicates the permuted threshold for genome-wide significance  $p = 4.67 \times 10^{-6}$ . Significant SNPs are shown with solid colors. (B) SNPs (with  $p < 1 \times 10^{-10}$ ) that have positive (blue) or negative (red) effects on viral copies. (C) The corresponding synonymous (triangles) and non-synonymous (circles) amino acid changes that associate with increased or decreased viral copies. Data is shown as means with 95% confidence intervals. The estimated effective sizes and associated standard deviations are given in [S1 Table](#). A Q-Q plot showing the observed distribution of p-value and the expected distribution is given in [S2 Fig](#).

(TIF)

**S5 Fig. GWAS analysis using only Cluster 1 data (shown in [S1 Fig](#)).** (A) Genome-wide association results of the impact of identified SNPs on viral copies during SARS-CoV-2 infection. The dashed line indicates the permuted threshold for genome-wide significance  $p = 4.03 \times 10^{-5}$  (0.05/1242 SNPs). Significant SNPs are shown with solid colors. (B) SNPs (with  $p < 1 \times 10^{-10}$ ) that have positive (blue) or negative (red) effects on viral copies. (C) The corresponding synonymous (triangles) amino acid changes that associate with increased or decreased viral copies. Data shown as means with 95% confidence intervals.

(TIF)

**S6 Fig. GWAS analysis using Cluster 2 data (shown in [S1 Fig](#)).** (A) Genome-wide association results of the impact of identified SNPs on viral copies during SARS-CoV-2 infection. The dashed line indicates the permuted threshold for genome-wide significance  $p = 3.68 \times 10^{-5}$  (0.05/1357 SNPs). Significant SNPs are shown with solid colors. (B) SNPs (with  $p < 1 \times 10^{-10}$ ) that have positive (blue) or negative (red) effects on viral copies. (C) The corresponding synonymous (triangles) amino acid changes that associate with increased or decreased viral copies. Data shown as means with 95% confidence intervals.

(TIF)

**S7 Fig. GWAS analysis using Cluster 3 data (shown in [S1 Fig](#)).** (A) Genome-wide association results of the impact of identified SNPs on viral copies during SARS-CoV-2 infection. The dashed line indicates the permuted threshold for genome-wide significance  $p = 2.80 \times 10^{-5}$  (0.05/1784 SNPs). Significant SNPs are shown with solid colors. (B) SNPs (with  $p < 1 \times 10^{-10}$ ) that have positive (blue) or negative (red) effects on viral copies. (C) The corresponding non-synonymous (circles) amino acid changes that associate with increased or decreased viral

copies. Data shown as means with 95% confidence intervals.  
(TIF)

**S8 Fig. GWAS analysis using Cluster 4 data (shown in S1 Fig).** (A) Genome-wide association results of the impact of identified SNPs on viral copies during SARS-CoV-2 infection. The dashed line indicates the permuted threshold for genome-wide significance  $p = 7.91 \times 10^{-6}$  (0.05/6314 SNPs). Significant SNPs are shown with solid colors. (B) SNPs (with  $p < 1 \times 10^{-10}$ ) that have positive (blue) or negative (red) effects on viral copies. (C) The corresponding synonymous (triangles) and non-synonymous (circles) amino acid changes that associate with increased or decreased viral copies. Data shown as means with 95% confidence intervals.  
(TIF)

**S9 Fig. The temporal dynamics of amino acid changes in the S gene associated with changes in viral copies.** The results are based on the multivariate regression analysis using the two MDS components as covariates. (A) The phylogenetic tree estimated from a representative set of 996 genome sequences showing variant assignments and the locations of amino acid changes that increase (blue) or decrease (red) viral copies. (B) The temporal dynamics of the SNPs from February 2021 to March 2023. The transparency of the color corresponds to the mutation fraction in the daily sequence count: transparent color indicates low fractions, and opaque color indicates high fractions.  
(TIF)

**S10 Fig. The temporal dynamics of amino acid changes in the ORF1ab gene associated with changes in viral copies.** The results are based on the multivariate regression analysis using the two MDS components as covariates. (A) The phylogenetic tree estimated from a representative set of 996 genome sequences showing variant assignments and the locations of amino acid changes that increase (blue) or decrease (red) viral copies. (B) The temporal dynamics of the SNPs from February 2021 to March 2023. The transparency of the color corresponds to the mutation fraction in the daily sequence count: transparent color indicates low fractions, and opaque color indicates high fractions.  
(TIF)

**S11 Fig. The temporal dynamics of amino acid changes in the ORF3a gene (S26L and T223I), M gene (D3G and I82T), ORF7b gene (T40I) and N gene (D63G and S413R) associated with changes in viral copies.** The results are based on the multivariate regression analysis using the two MDS components as covariates. (A) The phylogenetic tree estimated from a representative set of 996 genome sequences showing variant assignments and the locations of amino acid changes that increase (blue) or decrease (red) viral copies. (B) The temporal dynamics of the SNPs from February 2021 to March 2023. The transparency of the color corresponds to the mutation fraction in the daily sequence count: transparent color indicates low fractions, and opaque color indicates high fractions.  
(TIF)

**S12 Fig. Comparison of key variant-defining amino acid changes with GWAS-identified substitutions.** The comparison of the key amino acid changes (dark purple) in each variant, with GWAS-identified SNPs that were associated with negative (red) or positive (blue) effects on viral copies in the (A) S gene and (B) ORF1ab gene. The results of GWAS analysis using the two dimensions computed by MDS as covariates are shown as “GWAS 1”, and the results of the analysis using the categorical clusters as covariates are shown as “GWAS 2”. The effective sizes of identified SNPs using different population control methods are given in S1 and S2 Tables.  
(TIF)

**S1 Table. The identified amino acid changes associated estimated effective sizes and standard deviations using the multivariate linear regression model with categorical clusters as covariates.**

(DOCX)

**S2 Table. The identified amino acid changes associated estimated effective sizes and standard deviations using the multivariate linear regression model with MDS-computed dimensions as covariates.**

(DOCX)

## Acknowledgments

We would like to thank Verity Hill, Seth Redmond, Jiye Kwon, Rafael Lopes, Sophie Taylor, and Philip Jack for their helpful conversations and feedback on this work.

### Yale SARS-CoV-2 Genomic Surveillance Initiative:

Tara Alpert, Kaya Bilguvar, Kendall Billig, Mallery Breban, Anderson Brito, Christopher Castaldi, Rebecca Earnest, Bony De Kumar, Joseph Fauver, Chaney Kalinich, Tobias Koch, Marie Landry, Shrikant Mane, Isabel Ott, David Peaper, Mary Petrone, Kien Pham, Jessica Rothman, Irina Tikhonova, Chantal Vogels, Anne Watkins.

## Author Contributions

**Conceptualization:** Ke Li, Chrispin Chaguza.

**Data curation:** Ke Li, Nicholas F. G. Chen, David Ferguson, Sameer Pandya, Nick Kerantzas, Wade Schulz, Anne M. Hahn.

**Formal analysis:** Ke Li.

**Funding acquisition:** Nathan D. Grubaugh.

**Investigation:** David Ferguson, Sameer Pandya, Nick Kerantzas, Wade Schulz.

**Methodology:** Ke Li, Chrispin Chaguza, Julian Stamp, C. Brandon Ogbunugafor, Lorin Crawford.

**Software:** Julian Stamp, Lorin Crawford.

**Supervision:** Virginia E. Pitzer, Daniel M. Weinberger, Nathan D. Grubaugh.

**Visualization:** Ke Li, Yi Ting Chew.

**Writing – original draft:** Ke Li.

**Writing – review & editing:** Ke Li, Chrispin Chaguza, Julian Stamp, Yi Ting Chew, Nicholas F. G. Chen, Anne M. Hahn, C. Brandon Ogbunugafor, Virginia E. Pitzer, Lorin Crawford, Daniel M. Weinberger, Nathan D. Grubaugh.

## References

1. Killingley B, Mann AJ, Kalinova M, Boyers A, Goonawardane N, Zhou J, et al. Safety, tolerability and viral kinetics during SARS-CoV-2 human challenge in young adults. *Nat Med.* 2022; 28: 1031–1041. <https://doi.org/10.1038/s41591-022-01780-9> PMID: 35361992
2. Marks M, Millat-Martinez P, Ouchi D, Roberts CH, Alemany A, Corbacho-Monné M, et al. Transmission of COVID-19 in 282 clusters in Catalonia, Spain: a cohort study. *Lancet Infect Dis.* 2021; 21: 629–636. [https://doi.org/10.1016/S1473-3099\(20\)30985-3](https://doi.org/10.1016/S1473-3099(20)30985-3) PMID: 33545090

3. Jones TC, Biele G, Mühlemann B, Veith T, Schneider J, Beheim-Schwarzbach J, et al. Estimating infectiousness throughout SARS-CoV-2 infection course. *Science*. 2021;373. <https://doi.org/10.1126/science.abi5273> PMID: 34035154
4. Singanayagam A, Hakki S, Dunning J, Madon KJ, Crone MA, Koycheva A, et al. Community transmission and viral load kinetics of the SARS-CoV-2 delta (B.1.617.2) variant in vaccinated and unvaccinated individuals in the UK: a prospective, longitudinal, cohort study. *Lancet Infect Dis*. 2022; 22: 183–195. [https://doi.org/10.1016/S1473-3099\(21\)00648-4](https://doi.org/10.1016/S1473-3099(21)00648-4) PMID: 34756186
5. Kissler SM, Fauver JR, Mack C, Tai CG, Breban MI, Watkins AE, et al. Viral Dynamics of SARS-CoV-2 Variants in Vaccinated and Unvaccinated Persons. *N Engl J Med*. 2021; 385: 2489–2491. <https://doi.org/10.1056/NEJMc2102507> PMID: 34941024
6. Puhach O, Adea K, Hulo N, Sattonnet P, Genecand C, Iten A, et al. Infectious viral load in unvaccinated and vaccinated individuals infected with ancestral, Delta or Omicron SARS-CoV-2. *Nat Med*. 2022; 28: 1491–1500. <https://doi.org/10.1038/s41591-022-01816-0> PMID: 35395151
7. Boucau J, Marino C, Regan J, Uddin R, Choudhary MC, Flynn JP, et al. Duration of Shedding of Culturable Virus in SARS-CoV-2 Omicron (BA.1) Infection. *N Engl J Med*. 2022; 387: 275–277. <https://doi.org/10.1056/NEJMc2202092> PMID: 35767428
8. Hay JA, Kennedy-Shaffer L, Kanjilal S, Lennon NJ, Gabriel SB, Lipsitch M, et al. Estimating epidemiologic dynamics from cross-sectional viral load distributions. *Science*. 2021;373. <https://doi.org/10.1126/science.abh0635> PMID: 34083451
9. Fryer HR, Golubchik T, Hall M, Fraser C, Hinch R, Ferretti L, et al. Viral burden is associated with age, vaccination, and viral variant in a population-representative study of SARS-CoV-2 that accounts for time-since-infection-related sampling bias. *PLoS Pathog*. 2023; 19: e1011461. <https://doi.org/10.1371/journal.ppat.1011461> PMID: 37578971
10. Duesterwald L, Nguyen M, Christensen P, Wesley Long, Olsen R, Musser JM, et al. Using Genome Sequence Data to Predict SARS-CoV-2 Detection Cycle Threshold Values. <https://doi.org/10.1101/2022.11.14.22282297>
11. Uffelmann E, Huang QQ, Munung NS, de Vries J, Okada Y, Martin AR, et al. Genome-wide association studies. *Nature Reviews Methods Primers*. 2021; 1: 1–21. <https://doi.org/10.1038/s43586-021-00056-9>
12. Karim M, Dunham I, Ghousaini M. Mining a GWAS of Severe Covid-19. *The New England journal of medicine*. 2020. pp. 2588–2589. <https://doi.org/10.1056/NEJMc2025747> PMID: 33289971
13. Roberts GHL, Partha R, Rhead B, Knight SC, Park DS, Coignet MV, et al. Expanded COVID-19 phenotype definitions reveal distinct patterns of genetic association and protective effects. *Nat Genet*. 2022; 54: 374–381. <https://doi.org/10.1038/s41588-022-01042-x> PMID: 35410379
14. Genomewide Association Study of Severe Covid-19 with Respiratory Failure. *N Engl J Med*. 2020; 383: 1522–1534. <https://doi.org/10.1056/NEJMoa2020283> PMID: 32558485
15. Hahn G, Wu CM, Lee S, Lutz SM, Khurana S, Baden LR, et al. Genome-wide association analysis of COVID-19 mortality risk in SARS-CoV-2 genomes identifies mutation in the SARS-CoV-2 spike protein that colocalizes with P.1 of the Brazilian strain. *Genet Epidemiol*. 2021; 45: 685–693. <https://doi.org/10.1002/gepi.22421> PMID: 34159627
16. Power RA, Davaniah S, Derache A, Wilkinson E, Tanser F, Gupta RK, et al. Genome-Wide Association Study of HIV Whole Genome Sequences Validated using Drug Resistance. *PLoS One*. 2016; 11: e0163746. <https://doi.org/10.1371/journal.pone.0163746> PMID: 27677172
17. Ansari MA, Aranday-Cortes E, Ip CL, da Silva Filipe A, Lau SH, Bamford C, et al. Interferon lambda 4 impacts the genetic diversity of hepatitis C virus. *Elife*. 2019;8. <https://doi.org/10.7554/eLife.42463> PMID: 31478835
18. Ansari MA, Pedergrana V, L C Ip C, Magri A, Von Delft A, Bonsall D, et al. Genome-to-genome analysis highlights the effect of the human innate and adaptive immune systems on the hepatitis C virus. *Nat Genet*. 2017; 49: 666–673. <https://doi.org/10.1038/ng.3835> PMID: 28394351
19. Vogels CBF, Breban MI, Ott IM, Alpert T, Petrone ME, Watkins AE, et al. Multiplex qPCR discriminates variants of concern to enhance global surveillance of SARS-CoV-2. *PLoS Biol*. 2021; 19: e3001236. <https://doi.org/10.1371/journal.pbio.3001236> PMID: 33961632
20. Zhou C, Zhang T, Ren H, Sun S, Yu X, Sheng J, et al. Impact of age on duration of viral RNA shedding in patients with COVID-19. *Aging*. 2020; 12: 22399–22404. <https://doi.org/10.18632/aging.104114> PMID: 33223506
21. Acharya CB, Schrom J, Mitchell AM, Coil DA, Marquez C, Rojas S, et al. Viral Load Among Vaccinated and Unvaccinated, Asymptomatic and Symptomatic Persons Infected With the SARS-CoV-2 Delta Variant. *Open Forum Infect Dis*. 2022; 9: ofac135. <https://doi.org/10.1093/ofid/ofac135> PMID: 35479304
22. Torgerson WS. Multidimensional scaling: I. Theory and method. *Psychometrika*. 1952; 17: 401–419. <https://doi.org/10.1007/bf02288916>

23. VanderWeele TJ, Mathur MB. SOME DESIRABLE PROPERTIES OF THE BONFERRONI CORRECTION: IS THE BONFERRONI CORRECTION REALLY SO BAD? *Am J Epidemiol*. 2018; 188: 617–618. <https://doi.org/10.1093/aje/kwy250> PMID: 30452538
24. Crawford L, Zeng P, Mukherjee S, Zhou X. Detecting epistasis with the marginal epistasis test in genetic mapping studies of quantitative traits. *PLoS Genet*. 2017; 13: e1006869. <https://doi.org/10.1371/journal.pgen.1006869> PMID: 28746338
25. Stamp J, DenAdel A, Weinreich D, Crawford L. Leveraging the genetic correlation between traits improves the detection of epistasis in genome-wide association studies. *G3*. 2023;13. <https://doi.org/10.1093/g3journal/jkad118> PMID: 37243672
26. Zhou X. A UNIFIED FRAMEWORK FOR VARIANCE COMPONENT ESTIMATION WITH SUMMARY STATISTICS IN GENOME-WIDE ASSOCIATION STUDIES. *Ann Appl Stat*. 2017; 11: 2027–2051. <https://doi.org/10.1214/17-AOAS1052> PMID: 29515717
27. Tang X, Wu C, Li X, Song Y, Yao X, Wu X, et al. On the origin and continuing evolution of SARS-CoV-2. *Natl Sci Rev*. 2020; 7: 1012–1023. <https://doi.org/10.1093/nsr/nwaa036> PMID: 34676127
28. Carabelli AM, Peacock TP, Thorne LG, Harvey WT, Hughes J, COVID-19 Genomics UK Consortium, et al. SARS-CoV-2 variant biology: immune escape, transmission and fitness. *Nat Rev Microbiol*. 2023; 21: 162–177. <https://doi.org/10.1038/s41579-022-00841-7> PMID: 36653446
29. Souza PFN, Mesquita FP, Amaral JL, Landim PGC, Lima KRP, Costa MB, et al. The spike glycoprotein of SARS-CoV-2: A review of how mutations of spike glycoproteins have driven the emergence of variants with high transmissibility and immune escape. *Int J Biol Macromol*. 2022; 208: 105–125. <https://doi.org/10.1016/j.ijbiomac.2022.03.058> PMID: 35300999
30. Bartha I, Carlson JM, Brumme CJ, McLaren PJ, Brumme ZL, John M, et al. A genome-to-genome analysis of associations between human genetic variation, HIV-1 sequence diversity, and viral control. *Elife*. 2013; 2: e01123. <https://doi.org/10.7554/eLife.01123> PMID: 24171102
31. Lu X, Wang L, Sakthivel SK, Whitaker B, Murray J, Kamili S, et al. US CDC Real-Time Reverse Transcription PCR Panel for Detection of Severe Acute Respiratory Syndrome Coronavirus 2. *Emerg Infect Dis*. 2020; 26: 1654–1665. <https://doi.org/10.3201/eid2608.201246> PMID: 32396505
32. Vogels CBF, Brito AF, Wyllie AL, Fauver JR, Ott IM, Kalinich CC, et al. Analytical sensitivity and efficiency comparisons of SARS-CoV-2 RT-qPCR primer-probe sets. *Nature Microbiology*. 2020; 5: 1299–1305. <https://doi.org/10.1038/s41564-020-0761-6> PMID: 32651556
33. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv [q-bio. GN]*. 2013. Available: <http://arxiv.org/abs/1303.3997>.
34. Grubaugh ND, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, Main BJ, et al. An amplicon-based sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol*. 2019; 20: 8. <https://doi.org/10.1186/s13059-018-1618-7> PMID: 30621750
35. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of SAMtools and BCFtools. *Gigascience*. 2021;10. <https://doi.org/10.1093/gigascience/giab008> PMID: 33590861
36. O'Toole Á, Scher E, Underwood A, Jackson B, Hill V, McCrone JT, et al. Assignment of epidemiological lineages in an emerging pandemic using the pangolin tool. *Virus Evol*. 2021; 7: veab064. <https://doi.org/10.1093/ve/veab064> PMID: 34527285
37. Aksamentov I, Roemer C, Hodcroft E, Neher R. Nextclade: clade assignment, mutation calling and quality control for viral genomes. *J Open Source Softw*. 2021; 6: 3773. <https://doi.org/10.21105/joss.03773>
38. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, et al. SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microb Genom*. 2016; 2: e000056. <https://doi.org/10.1099/mgen.0.000056> PMID: 28348851
39. Creators Seemann, Torsten1 Show affiliations 1. The University of Melbourne. Source code for snp-dists software. <https://doi.org/10.5281/zenodo.1411986>
40. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, et al. Corrigendum to: IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Mol Biol Evol*. 2020; 37: 2461. <https://doi.org/10.1093/molbev/msaa131> PMID: 32556291
41. Gangavarapu K, Latif AA, Mullen JL, Alkuzweny M, Hufbauer E, Tsung G, et al. Outbreak.info genomic reports: scalable and dynamic surveillance of SARS-CoV-2 variants and mutations. *Nat Methods*. 2023; 20: 512–522. <https://doi.org/10.1038/s41592-023-01769-3> PMID: 36823332