



ORIGINAL RESEARCH

Genomic Perspectives on the Emerging SARS-CoV-2 Omicron Variant



Wentai Ma^{1,2,#}, Jing Yang^{1,2,#}, Haoyi Fu^{1,2}, Chao Su³, Caixia Yu⁴,
 Qihui Wang³, Ana Tereza Ribeiro de Vasconcelos⁵, Georgii A. Bazykin^{6,7},
 Yiming Bao^{2,4}, Mingkun Li^{1,2,8,*}

¹ CAS Key Laboratory of Genomic and Precision Medicine, Beijing Institute of Genomics, Chinese Academy of Sciences and China National Center for Bioinformation, Beijing 100101, China

² University of Chinese Academy of Sciences, Beijing 100049, China

³ CAS Key Laboratory of Pathogenic Microbiology and Immunology, Institute of Microbiology, Chinese Academy of Sciences, Beijing 100101, China

⁴ National Genomics Data Center, Beijing Institute of Genomics, Chinese Academy of Sciences and China National Center for Bioinformation, Beijing 100101, China

⁵ Laboratório de Bioinformática, Laboratório Nacional de Computação Científica, Petrópolis 25651-075, Brazil

⁶ Skolkovo Institute of Science and Technology, Moscow 121205, Russia

⁷ Kharkevich Institute for Information Transmission Problems of the Russian Academy of Sciences, Moscow 127051, Russia

⁸ Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650201, China

Received 27 December 2021; revised 6 January 2022; accepted 7 January 2022

Available online 13 January 2022

Handled by Fangqing Zhao

KEYWORDS

Omicron;
 Genomics;
 Mutation;
 Variant of concern;
 SARS-CoV-2

Abstract A new variant of concern for SARS-CoV-2, Omicron (B.1.1.529), was designated by the World Health Organization on November 26, 2021. This study analyzed the viral genome sequencing data of 108 samples collected from patients infected with Omicron. First, we found that the enrichment efficiency of viral nucleic acids was reduced due to mutations in the region where the primers anneal to. Second, the Omicron variant possesses an excessive number of mutations compared to other variants circulating at the same time (median: 62 vs. 45), especially in the *Spike* gene. Mutations in the *Spike* gene confer alterations in 32 amino acid residues, more than those observed in other SARS-CoV-2 variants. Moreover, a large number of nonsynonymous mutations occur in the codons for the amino acid residues located on the surface of the Spike protein, which could potentially affect the replication, infectivity, and antigenicity of SARS-CoV-2. Third, there are 53 mutations between the Omicron variant and its closest sequences available in public

* Corresponding author.

E-mail: limk@big.ac.cn (Li M).

Equal contribution.

Peer review under responsibility of Beijing Institute of Genomics, Chinese Academy of Sciences / China National Center for Bioinformation and Genetics Society of China.

<https://doi.org/10.1016/j.gpb.2022.01.001>

1672-0229 © 2022 The Authors. Published by Elsevier B.V. and Science Press on behalf of Beijing Institute of Genomics, Chinese Academy of Sciences / China National Center for Bioinformation and Genetics Society of China.

This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

databases. Many of these mutations were rarely observed in public databases and had a low mutation rate. In addition, the linkage disequilibrium between these mutations was low, with a limited number of mutations concurrently observed in the same genome, suggesting that the Omicron variant would be in a different evolutionary branch from the currently prevalent variants. To improve our ability to detect and track the source of new variants rapidly, it is imperative to further strengthen genomic surveillance and data sharing globally in a timely manner.

Introduction

On November 22, 2021, the first genome sequence of a new variant of concern (VOC), Omicron (also known as B.1.1.529), was released in Global Initiative on Sharing All Influenza Data (GISAID: EPI_ISL_6590782) [1]. The sample was obtained from a patient who arrived in Hong Kong, China on November 11 from South Africa via Doha in Qatar (<https://news.sky.com/story/covid-19-how-the-spread-of-omicron-went-from-patient-zero-to-all-around-the-globe-12482183>). To date, the first known Omicron variant sample was collected on November 5, 2021 in South Africa (GISAID: EPI_ISL_7456440). Until December 12, 2021, there were over 2000 Omicron sequences submitted to the GISAID from South Africa, Botswana, Ghana, the United Kingdom, and many other countries. The emergence of this variant has attracted much attention due to the sheer number of mutations in the *Spike* gene, which may affect the viral transmissibility, replication, and binding of antibodies, and its dramatic increase in South Africa [2]. Preliminary studies have shown that the new variant could substantially evade immunity from prior infection and vaccination [3,4]. Meanwhile, a study has proposed that the emergence of the Omicron variant is associated with an increased risk of SARS-CoV-2 reinfection [5]. However, it is still unclear where the new variant came.

In this study, we characterized the genomic features of the Omicron variant using data from 108 patients infected with the Omicron variant, which were generated by the Network for Genomic Surveillance in South Africa (NGS-SA) [2,6], and we speculate that the new variant is unlikely derived from recently discovered variants through either mutation or recombination.

Results

Reduced enrichment efficiency of the PCR-tiling amplicon protocols on the Omicron variant

Among 207 Omicron samples sequenced and shared by NGS-SA, 158 samples had more than 90% of the viral genome covered by at least 5-fold, which were used in the subsequent analysis. Notably, two sequencing protocols were implemented. The first was to enrich the viral genome with the Midnight V6 primer sets followed by sequencing on the GridION platform (hereinafter referred to as Midnight, [dx.doi.org/10.17504/protocols.io.bwypfvn](https://doi.org/10.17504/protocols.io.bwypfvn)). The second protocol involved enrichment by the Artic V4 primer sets, and the amplicons were sequenced on the Illumina MiSeq platform (hereinafter referred to as Artic, [dx.doi.org/10.17504/protocols.io.bdp7i5rn](https://doi.org/10.17504/protocols.io.bdp7i5rn)). Fifty samples were sequenced using both protocols, and we found a high consistency in the major allele frequency between the two protocols (Figure S1). Artic data were preferred due to higher sequencing depth (median: 191 for Midnight vs. 250 for

Artic; $P < 0.01$, Mann–Whitney U test). Finally, 49 samples sequenced by the Midnight protocol and 59 samples sequenced by the Artic protocol were included in the study.

Both protocols enabled efficient enrichment of viral nucleic acids from total RNA, and the fractions of SARS-CoV-2 reads in the sequencing data were 84% and 94% for the Midnight and Artic protocols, respectively. Although the Artic protocol had a relatively higher in-target percentage ($P < 0.001$, Mann–Whitney U test), the evenness of the sequencing depth of the SARS-CoV-2 was higher for the Midnight protocol (variance of the sequencing depth, 0.121 for Midnight vs. 0.159 for Artic; $P < 0.001$, Mann–Whitney U test). The sequencing depth profiles of the SARS-CoV-2 genome were similar among samples sequenced by the same protocol but differed markedly between the two protocols (Figure 1A). The sequencing depths varied among different genomic regions, reflecting the differential enrichment efficiency of the primers used for amplification. Moreover, we found that the large number of mutations possessed by the Omicron variant had a significant impact on the enrichment efficiency of the primers. In particular, the enrichment efficiency of seven primers in the Artic protocol and three primers in the Midnight protocol was affected by at least one mutation (Figure 1A). The worst coverages of the three regions for Primers 76, 79, and 90 using the Artic protocol were all associated with the presence of mutations in these regions where these primers annealed to, whose sequencing depths were reduced by 2586-fold, 246-fold, and 234-fold, respectively, compared to the expected depth (Figure 1B). Strikingly, five mutations were located in the region where the 5' end of the least efficient Primer 76 annealed to. The enrichment efficiency of another four primers in the Artic protocol (Primers 10, 27, 88, and 89) was less affected by the mutations, which showed 1.3-fold, 1.4-fold, 3.4-fold, and 1.9-fold reductions, respectively. Thus, the results suggest that the Omicron mutations can decrease the enrichment efficiency by PCR amplification, and there is an urgent need to update the Arctic V4 primers. We noted that the developer of the Artic protocol had already proposed a solution to this, and all seven affected primers had been updated (<https://community.artic.network/t/sars-cov-2-v4-1-update-for-omicron-variant/342>). In contrast, the efficiency of Midnight primers was less influenced by mutations in the Omicron variant. The three affected primers, Primers 10, 24, and 28, showed no reduction, 2-fold reduction, and 28-fold reduction, respectively, in sequencing depth compared to the expected depth.

An extraordinary number of mutations in the *Spike* gene of the Omicron variant

The number of mutations (with major allele frequency $\geq 70\%$) of the Omicron variant varied from 61 to 64, and 61 of them were identified in more than 90% of the samples, which included 54 SNPs, 6 deletions, and 1 insertion. All these

mutations were fixed at the individual level (Figure 2A). The total number of mutations of the Omicron variant was significantly higher than that of other variants detected in South Africa in November (median: 62 vs. 45; $P < 0.001$, Mann–Whitney U test). Strikingly, over half of these mutations (34, 55.7%) were located in the *Spike* gene, whose length was 12.8% of the whole genome. Moreover, 32 of these mutations were nonsynonymous mutations. Such proportion was significantly higher than that observed in the same region in other variants (94% vs. 67%; $P < 0.001$, Fisher's exact test, Ka/Ks [7] = 8.65), suggesting positive selection on this gene.

The Omicron variant showed a greater number of mutations than other VOCs (Figure 2B). The difference was more marked in the *Spike* gene. As a result, the Omicron variant possessed 1.6–2.7 times more amino acid changes in the Spike protein, and 4–14 times more amino acid changes in the receptor-binding domain (RBD) region of the Spike protein than other VOCs collected simultaneously (Figure 2C and D). Strikingly, the divergence in the amino acid sequence between the Omicron variant and the early SARS-CoV-2 sequence (Wuhan-Hu-1) in the Spike protein and its RBD region was greater than or equivalent to that between SARS-like coronavirus (Pangolin MP789, Bat BANAL-20-52, and Bat RaTG13) and Wuhan-Hu-1 [8–12]. The dramatic changes in the Spike protein and its RBD region may substantially change the antigenicity and susceptibility to pre-existing antibodies.

Potential risks associated with Omicron mutations

Most mutations occurred on the surface of the trimeric Spike protein, especially in the RBD region (Figure S2). Eight of the 15 amino acid mutations in the RBD region (K417N, G446S, E484A, Q493R, G496S, Q498R, N501Y, and Y505H) were located at the positions that were proposed to be critical for viral binding to the host receptor angiotensin-converting enzyme 2 (ACE2) [13]. Among them, the K417N and N501Y mutations, which were also identified in the Beta variant, have been reported to influence binding to human ACE2 [14]; N501Y confers a higher affinity of the viral Spike protein to ACE2 [15]. How other mutations affect the affinity to ACE2 of humans and other animal hosts is still unknown.

Moreover, some other amino acid changes in the Spike protein are known to be associated with changes in replication and infectivity of the virus. For example, $\Delta 69-70$ could enhance infectivity associated with increased cleaved Spike incorporation [16]; P681H could potentially confer replication advantage through increased cleavage efficacy by furin and adaptation to resist innate immunity [3,17]; H655Y was suspected to be an adaptive mutation that could increase the infectivity of the virus in both human and animal models [16]. In addition, amino acid mutations in other proteins, such as R203K and G204R in the Nucleoprotein protein, could also potentially increase the infectivity, fitness, and virulence of the virus [18]. Of note, the function of these mutations was investigated because they were present in other VOCs. The effect of other less frequent mutations and the combination of the aforementioned mutations on the biology of the virus warrants further investigation.

Mutations in the RBD region of the Spike protein, which is the target of many antibodies, may compromise the neutralization of existing antibodies induced by vaccination or natural

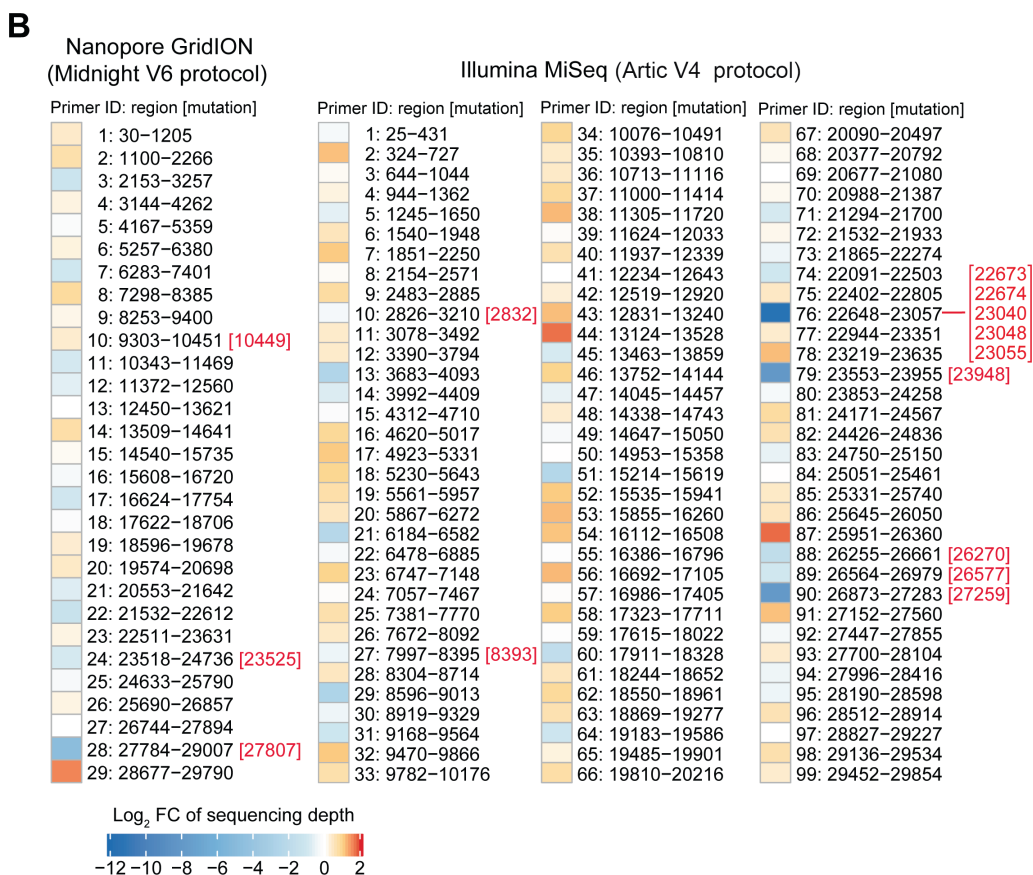
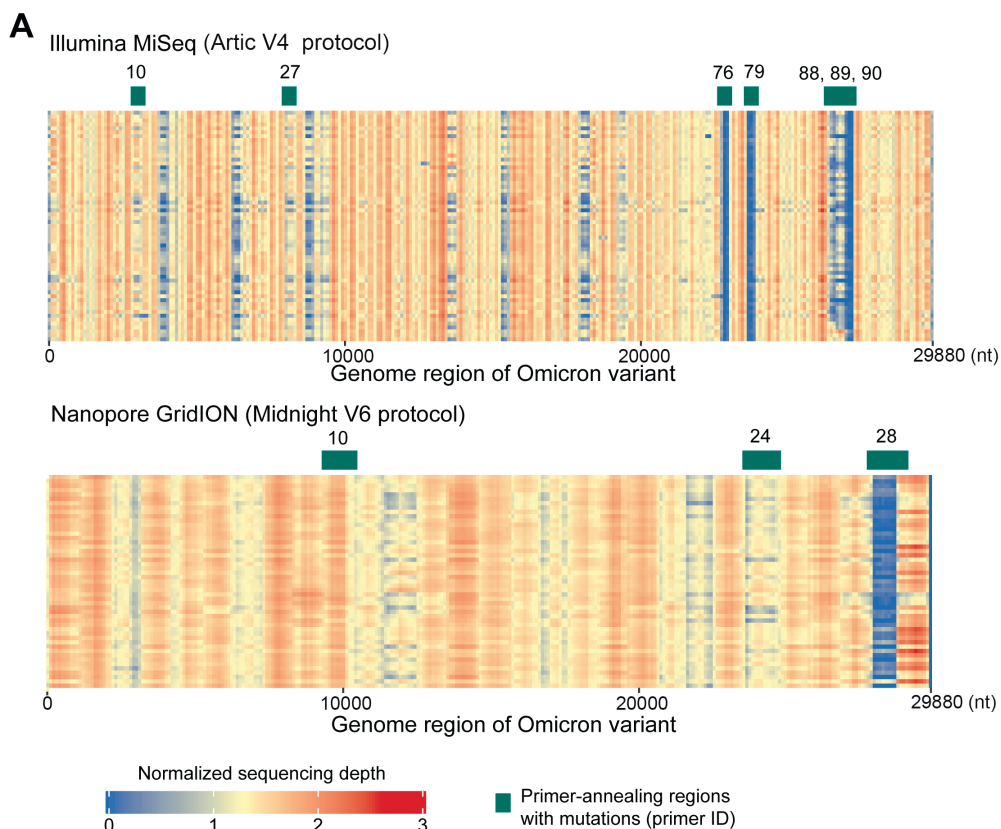
infection [19]. Recent studies have shown severely reduced neutralization of the Omicron variant by monoclonal antibodies and vaccine sera [4,20,21]. Meanwhile, preliminary studies suggested that the Omicron variant caused two times more reinfection than previous strains, further supporting the speculation that the new variant can evade immunity from prior infection and vaccination [5]. However, the escape from pre-existing immunity is incomplete, and a vaccine booster shot is likely to provide a high level of protection against the Omicron variant [4]. Here, we analyzed the epitope regions of 182 protein complex structures of antibodies that bind to SARS-CoV-2 Spike [including the RBD, N-terminal domain (NTD), and other regions] from the Protein Data Bank. We found that mutations in the Omicron variant were enriched in the epitope region of the Spike protein (Figure 3A). The median number of antibodies bound to the Omicron mutation sites was 53, which was significantly higher than those bound to other positions (median = 3; $P < 0.001$, Mann–Whitney U test). Moreover, by analyzing the deep mutational scanning data [22], we found that these mutations could potentially impact the binding of different classes of antibodies (Figure 3B), which was classified by the location and conformation of antibody binding [23], suggesting that the therapeutic strategy of antibody cocktails may also be affected.

Obscure evolutionary trajectory of the Omicron variant

In addition to the 61 shared mutations, some specific mutations were identified in different individuals, ranging from one to three, indicating relatively low population diversity at the time of sampling (Figure 4A). Meanwhile, no obvious clusters were found in the phylogenetic tree, suggesting that the Omicron variant was still in the early transmission stage during sampling. The time to the most recent common ancestor (TMRCA) was estimated to be in the middle of October 2021 (95% highest density interval: October 7 to October 20).

To screen for the possible predecessors of the Omicron variant, the 108 Omicron sequences were used as queries to look for the closest sequences in public databases, which included more than 5 million sequences released before November 1, 2021. We found three closest sequences to the queries, which differed by 53–56 nucleotides from the Omicron genomes. The three sequences were from lineage B.1.1 and collected between March and June 2020. They all had eight mutations relative to Wuhan-Hu-1, and seven of the mutations were shared among them (Figure 4A). The presence of a large number of differences suggests that the Omicron lineage was separated from other lineages a long time ago and has never been sequenced since then. This is an uncommon situation considering more than 5 million genomes have been sequenced in over 180 countries and regions. The distribution of the number of differences between all haplotype sequences in public databases and their closest sequences showed that 53 is approximately 1-fold higher than the maximum number of differences observed in public databases (20 when at least three sequences were required to eliminate the influence of sequencing or assembly errors, Figure 4B), again emphasizing the distinctiveness of the Omicron variant.

Most Omicron lineage-specific mutations (52/54) were identified in public databases (Figure 4C). However, they were unlikely to be present in one sequence by chance. First, over



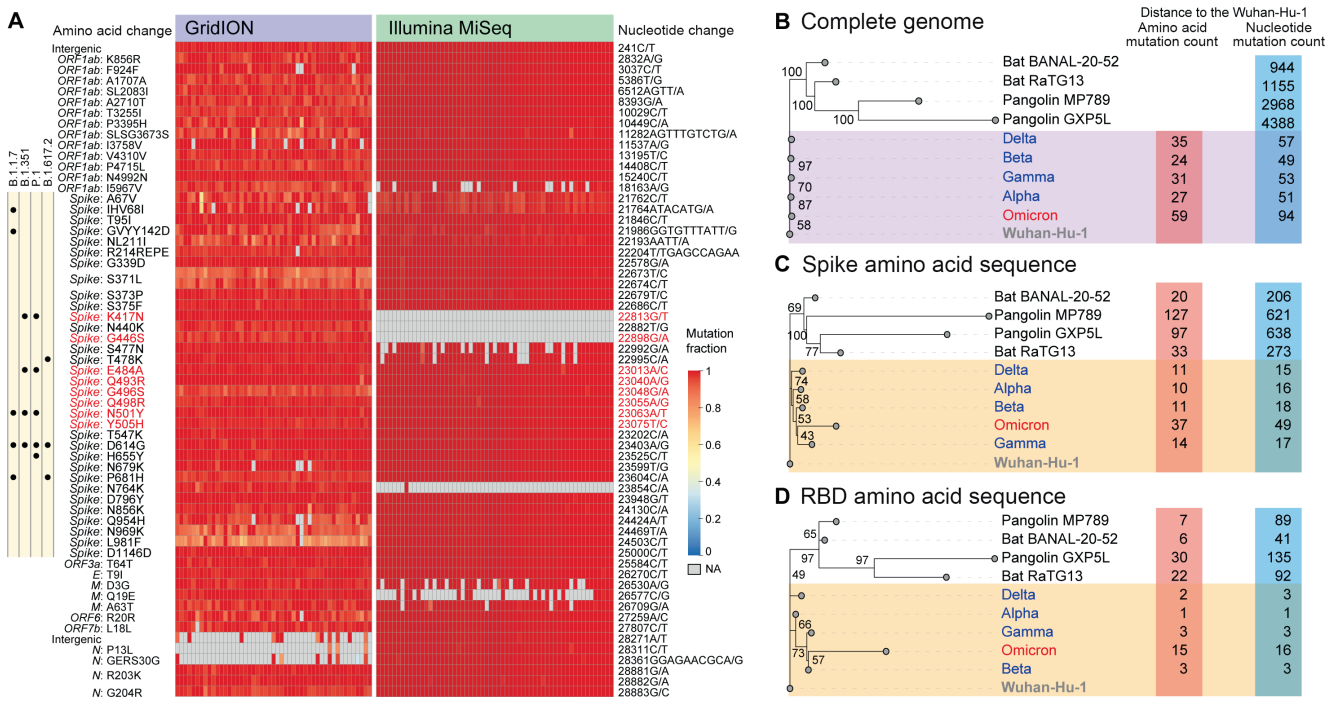


Figure 2 Mutations in the Omicron genome and its evolutionary relationship with other variants and SARS-like coronaviruses

A. Summary of mutations in the Omicron genome. Each row represents a mutation, and changes in nucleotides and amino acids are marked on two sides of the heatmap. Mutations located in the sites critical for viral binding to the human receptor ACE2 are marked in red [13]. Mutations observed in the *Spike* gene of other VOCs are listed on the left of the heatmap. **B.** Phylogenetic tree of five VOCs and SARS-like coronaviruses based on the nucleotide sequences. **C.** Phylogenetic tree of five VOCs and SARS-like coronaviruses based on the amino acid sequences of the Spike protein. **D.** Phylogenetic tree of five VOCs and SARS-like coronaviruses based on the amino acid sequences of the RBD region of the Spike protein. Two bat coronaviruses (Bat BANAL-20–52 and Bat RaTG3) whose genomes are most similar to SARS-CoV-2 [8,9], two pangolin coronaviruses (Pangolin MP789 and Pangolin GXP5L) [10,11], and sequences of four recently collected VOCs [Alpha variant (GISAID: EPI_ISL_6141707), Beta variant (GISAID: EPI_ISL_6774033), Gamma variant (GISAID: EPI_ISL_6898988), Delta variant (GISAID: EPI_ISL_6585201)] were included in the analysis of the phylogenetic tree. All sequences of the Alpha, Beta, Gamma, and Delta variants were collected in November 2021, and those collected in South Africa were preferred. The Wuhan-Hu-1 sequence is shown as the outgroup of the tree for better visualization [12]. The number of mutations relative to Wuhan-Hu-1 is listed on the right of the tree. Insertion of multiple bases is considered as a single mutation, while deletion of multiple bases is considered as multiple single-base deletions. ACE2, angiotensin-converting enzyme 2; VOC, variant of concern; RBD, receptor-binding domain; GISAID, Global Initiative on Sharing All Influenza Data; NA, not available.

half of the mutations were rarely detected in the populations, *i.e.*, 33 mutations were detected in less than 1000 samples out of five million sequences (16 mutations were detected in less than 100 samples). Second, the mutation rate (represented by the occurrence number of mutations on the phylogenetic tree) was extremely low for 11 of the mutations (occurring only once in the evolution of SARS-CoV-2, mutation rate = 1).

Third, the linkage disequilibrium between these mutations was low, and only four mutation pairs had r^2 greater than 0.8. Moreover, we further examined whether any combination of these mutations appeared in public databases and found that the maximum number of mutations in the same genome was six. Therefore, the evolutionary trajectory of the Omicron lineage cannot be resolved by the current genome data.

Figure 1 Sequence enrichment efficiency of the Omicron variant using different protocols

A. Distribution of the sequencing depth of the Omicron variant. The average sequencing depth is shown for each non-overlapping window of 100 bp after normalization by the total number of reads in the sample. The primers affected by the mutations in Omicron are labeled on top of each heatmap. **B.** The efficiency of each primer in amplifying the nucleic acids of the Omicron variant. The color represents the fold change of enrichment efficiency, calculated by the sum of the depths of all samples in this region divided by the expected value (assuming no differences among regions). The overlapping region of adjacent primers was excluded from the analysis. The Omicron mutations located in the region where primers anneal to are labeled on the right of the primer ID.

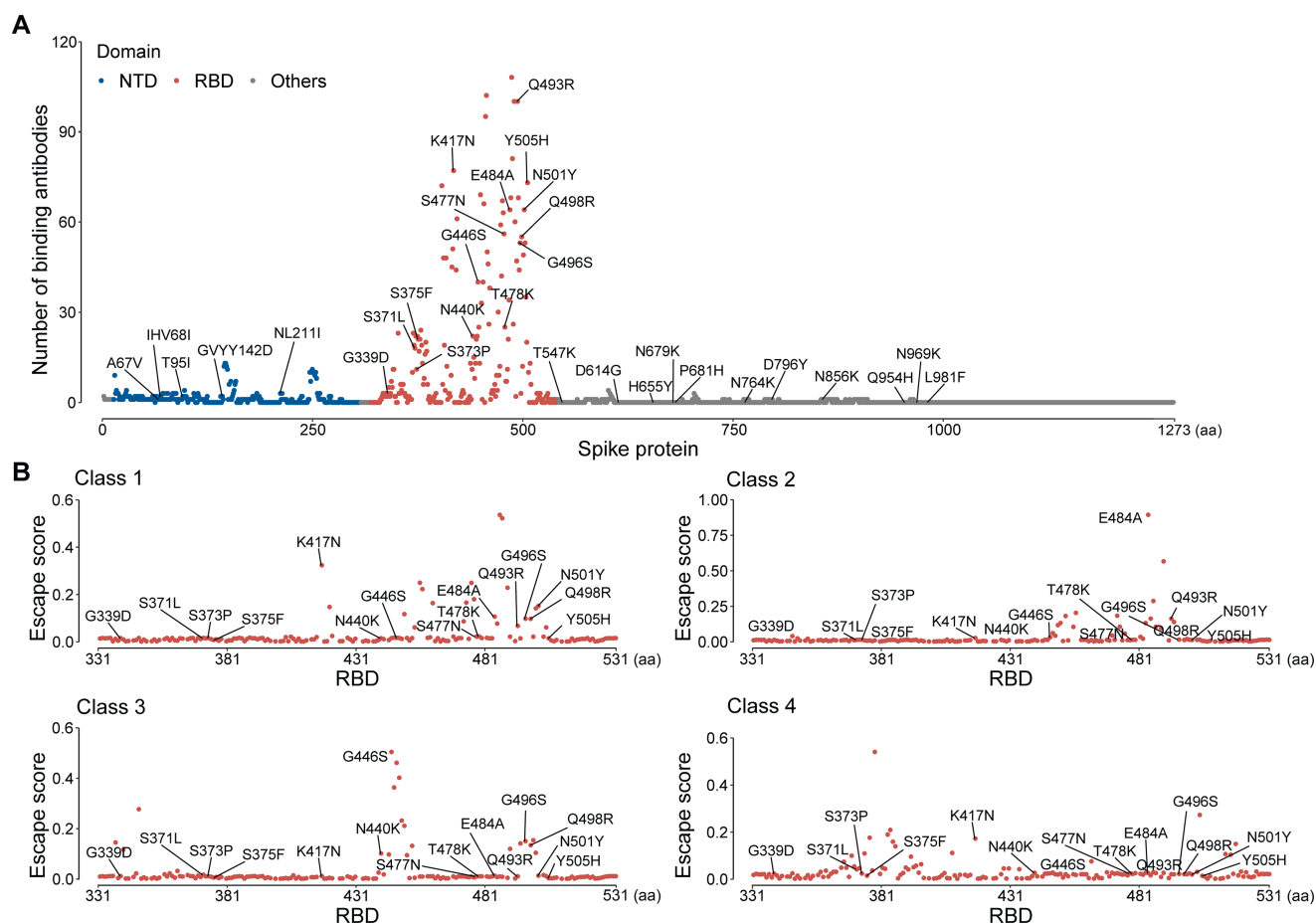


Figure 3 Distribution of the Omicron mutations at the antibody binding positions

A. The number of binding antibodies at the Omicron mutation sites in the Spike protein. All Omicron mutations in the Spike protein except R214REPE were labeled in this panel. **B.** The escape score of the Omicron mutations estimated from deep mutational scanning. The escape score for each position was calculated as the mean of the scores of all antibodies belonging to the same class. NTD, N-terminal domain.

Discussion

The unique genome features of the Omicron variant make it the most special SARS-CoV-2 variant to date. The excess number of nonsynonymous mutations in the *Spike* gene implies that the Omicron variant might evolve under selection pressure, which may come from antibodies or adaptation to new hosts. It is speculated that it may have been incubated in a patient chronically infected with SARS-CoV-2, e.g., HIV patients with immunocompromising conditions. This hypothesis has been supported by the accelerated viral evolution observed in immunocompromised patients and has been previously proposed to explain how the Alpha variant was generated [24,25] (<https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563>). If this hypothesis is true for the Omicron variant, we suspect that the original virus that infected the patient might still be missing in public databases because the current closest sequences were circulating in population one and a half years ago; the time was too long, even for a chronic infection. Another hypothesis involves a spillover from humans to animals followed by a spillover back from animals to humans; such

a process has been proposed to be possible in mink [26]. Interestingly, a recent study has proposed that the progenitor of the Omicron variant seems to have evolved in mice for some time before jumping back into humans [27]. The binding affinity test between the Omicron RBD and animal ACE2 may help to test this hypothesis. A third hypothesis is that the virus split with other variants a long time ago and was transmitted cryptically in the population. Since viral genome surveillance is poor in many countries, it is difficult to reject this hypothesis, which again underscores the importance of strengthening viral surveillance on a global scale. Moreover, a hypothesis of acquisition by recombination between different variants is unlikely since the components that make up the Omicron genome could not be found in the current SARS-CoV-2 databases, and of course, we cannot reject the possibility that the Omicron genome consists of a combination of components that have not been sequenced. More discussion of the possible origin of the Omicron variant can be found in other studies [28].

Benefiting from the establishment of the viral genome surveillance network and extensive research on the function of viral mutations, it took less than a week to designate the new VOC Omicron since the first identification of its genome,

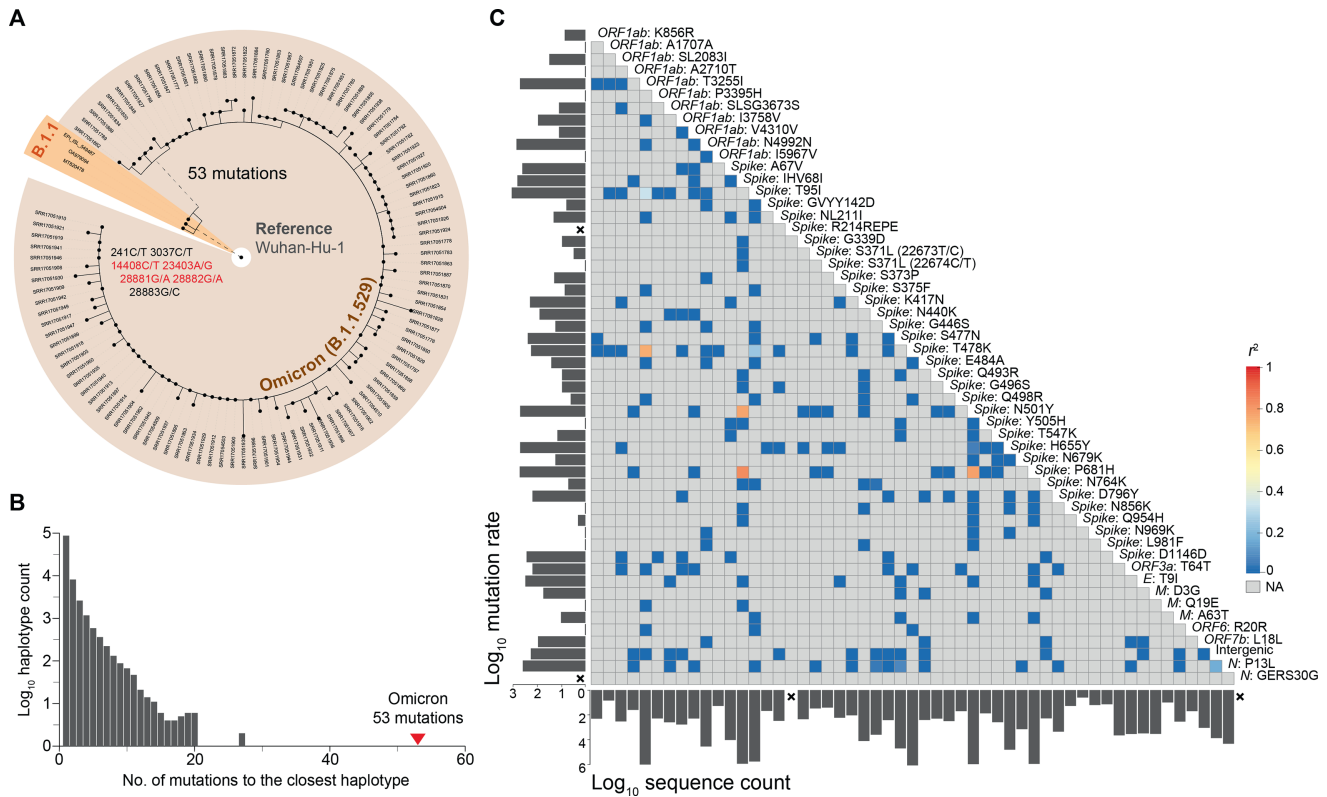


Figure 4 Evolutionary features of the Omicron variant

A. The phylogenetic tree of 108 Omicron sequences and their closest sequences in public databases. Wuhan-Hu-1 is shown as the outgroup of the tree. The three closest sequences belonging to lineage B.1.1 are highlighted in orange. Nonsynonymous mutations are marked in red. **B.** The distribution of the number of differences between all haplotypes (nonredundant sequences) in public databases and their closest sequences. The minimum number of sequences required for a valid haplotype was set to 3. **C.** Correlation between different Omicron mutations. Only 54 Omicron lineage-specific mutations were included in the analysis. The color in the heatmap represents the linkage disequilibrium coefficient (r^2) between mutations. The mutation rate and the number of sequences in public databases that possess the same mutation are labeled on the left and bottom of the heatmap, respectively. A cross is labeled if the mutation was not observed in public databases.

which is much faster than the designation of previous VOCs. However, it will still take several months to verify the risk of the new VOC. There have been over 200,000 new infections per day in the past year. Undoubtedly, we will face more mutant variants in the future, which may result in significant changes in transmissibility, infectivity, and pathogenicity. Unfortunately, it is still impossible to predict the evolutionary direction of the viral genome; hence, we have no hint at what the next VOC will be. To enhance the ability to rapidly respond to the emergence of new VOCs, we should further strengthen genome surveillance on a worldwide scale and develop experimental and computational methods for rapid and high-throughput resolution of mutational functions.

Materials and methods

Data collection

The sequencing data were retrieved from the Sequence Read Archive (SRA) database in NCBI (BioProject: PRJNA784038), which were generated by the NGS-SA [2,6]. In total 211 samples were downloaded on November 30, 2021 (Table S1). The virus lineage was assigned by Pango

[29], and 4 samples that cannot be assigned to the Omicron lineage were discarded. All the remaining 207 samples were assigned to Omicron BA.1.

Quality control and mutation detection

Quality control and adaptor trimming were performed by FASTP [30]. The resultant reads were mapped to Wuhan-Hu-1 (GenBank: NC_045512.2) using minimap2 (-ax sr) [31]. Primer alignment and trimming were performed by the align_trim function from Artic (https://artic-tools.readthedocs.io/en/latest/commands/#align_trim). The mpileup file and the read count file were generated by SAMtools [32] and Varscan2, respectively [33]. The consensus sequence was obtained using the following criteria: 1) depth ≥ 5 -fold; and 2) frequency of the major allele $\geq 70\%$.

Sequence depth analysis

The sequencing depth was calculated for each non-overlapping window with a size of 100 bp, except for the last window, which ranged from 29801 nt to 29880 nt. The fold change of each primer region was calculated by the sum of the depth

of all samples in this region divided by the expected value (assuming no differences among regions).

Identification of epitope regions on the Spike protein

We downloaded the structures of 182 protein complexes of antibodies that bind to the SARS-CoV-2 Spike or its RBD or NTD from the Protein Data Bank (PDB; all structures available before August 8, 2021, <https://www.rcsb.org/>). The residues in the Spike protein involved in binding to antibodies were identified by a distance of less than 4.5 Å between two counterparts in which van der Waals interactions occur. Deep mutational scanning results were obtained from https://jbloombio.github.io/SARS2_RBD_Ab_escape_maps/, which includes information on sites in the SARS-CoV-2 Spike RBD where mutations reduce binding by antibodies/sera [22]. The escape score at each position was calculated as the mean of the scores of all antibodies belonging to the same class.

Display of the Omicron mutations on the structure of the Spike protein

We downloaded the cryogenic electron microscopy structure of SARS-CoV-2 Spike extracellular domain (PDB: 6VYB) and the crystal structure of RBD-hACE2 complex (PDB: 6LZG) from the PDB (<https://www.rcsb.org/>). The structure of the RBD region was extracted from the RBD-hACE2 complex. All structures were visualized by PyMOL software (<https://pymol.org/2/>). Omicron mutations relative to Wuhan-Hu-1 are labeled on the structure except for those invisible in the structure.

Construction of the phylogenetic tree

The amino acid sequences were converted from nucleotide sequences using MEGA-X (10.1.8) [34]. Phylogenetic construction was performed by IQ-TREE (1.6.12) [35]. The GTR+I model was used for nucleotide sequences, while the Blosum62 model was used for amino acid sequences.

TMRCAs estimation

The estimation of TMRCA and mutation rate was performed by BEAST (v2.6.4) [36] using 108 sequences collected between November 13, 2021 and November 23, 2021. The HKY85 nucleotide substitution model and strict molecular clock were used.

Search for the closest sequences in public databases

The distance of two SARS-CoV-2 sequences was represented by the mutation difference, which was calculated by an online tool at National Genomics Data Center, China National Center for Bioinformation (<https://ngdc.cncb.ac.cn/ncov/online/tool/genome-tracing/?lang=en>). Publicly available SARS-CoV-2 sequences were downloaded from the GISAID, NCBI, and RCoV19 databases (November 1, 2021); only

high-quality and complete sequences were included in the analysis [1,37].

Calculation of linkage disequilibrium

The r^2 statistic was used to measure the strength of the linkage disequilibrium between each pair of mutations [38]. The calculation of linkage disequilibrium was based on all unique haplotypes from public databases.

CRedit author statement

Wentai Ma: Methodology, Formal analysis, Writing - original draft. **Jing Yang:** Methodology, Formal analysis, Writing - original draft. **Haoyi Fu:** Methodology, Formal analysis. **Chao Su:** Methodology, Formal analysis. **Caixia Yu:** Resources. **Qihui Wang:** Methodology. **Ana Tereza Ribeiro de Vasconcelos:** Resources, Writing - review & editing. **Georgii A. Bazykin:** Resources, Writing - review & editing. **Yiming Bao:** Methodology, Writing - review & editing. **Mingkun Li:** Conceptualization, Methodology, Supervision, Writing - original draft, Writing - review & editing. All authors have read and approved the final manuscript.

Competing interests

The authors declare that they have no competing interests.

Acknowledgments

We thank Dr. Jennifer Giandhari, Dr. Eduan Wilkinson, and Dr. Tulio de Oliveira from Centre for Epidemic Response and Innovation (CERI), Stellenbosch University, South Africa for sharing the data and workflow, GISAID and associated laboratories and researchers for the shared sequence information, Dr. Aiping Wu from Suzhou Institute of Systems Medicine, Chinese Academy of Medical Sciences for help in the phylogenetic tree analysis, Mikhail Moldovan from Skolkovo Institute of Science and Technology, Russia for help in the recombination analysis. This study was funded by the National Natural Science Foundation of China (Grant No. 82161148009), the Strategic Priority Research Program of Chinese Academy of Sciences (Grant No. XDB38030400), the Capital Health Development and Research Special Programme (Grant No. 2021-1G-3012), the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) - NGS-BRICS - n°: 440931/2020-7, and the Russian Foundation for Basic Research (RFBR) (Grant No. 20-54-80014).

Supplementary material

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.gpb.2022.01.001>.

ORCID

ORCID 0000-0003-1931-8687 (Wentai Ma)
 ORCID 0000-0002-3934-7883 (Jing Yang)
 ORCID 0000-0001-9696-5445 (Haoyi Fu)
 ORCID 0000-0002-5824-7968 (Chao Su)
 ORCID 0000-0002-3882-9979 (Caixia Yu)
 ORCID 0000-0003-3768-0401 (Qihui Wang)
 ORCID 0000-0002-4632-2086
 (Ana Tereza Ribeiro de Vasconcelos)
 ORCID 0000-0003-2334-2751 (Georgii A. Bazykin)
 ORCID 0000-0002-9922-9723 (Yiming Bao)
 ORCID 0000-0003-1041-1172 (Mingkun Li)

References

- [1] Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Glob Chall* 2017;1:33–46.
- [2] Viana R, Moyo S, Amoako DG, Tegally H, Scheepers C, Althaus CL, et al. Rapid epidemic expansion of the SARS-CoV-2 Omicron variant in southern Africa. *Nature* 2022;603:679–86.
- [3] Lista MJ, Winstone H, Wilson HD, Dyer A, Galao RP, Lorenzo GD, et al. The P681H mutation in the Spike glycoprotein confers Type I interferon resistance in the SARS-CoV-2 Alpha (B.1.1.7) variant. *bioRxiv* 2021;467693.
- [4] Cele S, Jackson L, Khoury DS, Khan K, Moyo-Gwete T, Tegally H, et al. Omicron extensively but incompletely escapes Pfizer BNT162b2 neutralization. *Nature* 2022;602:654–6.
- [5] Pulliam JRC, van Schalkwyk C, Govender N, von Gottberg A, Cohen C, Groome MJ, et al. Increased risk of SARS-CoV-2 reinfection associated with emergence of the Omicron variant in South Africa. *Science* 2022;376:eabn4947.
- [6] Wilkinson E, Giovanetti M, Tegally H, San JE, Lessells R, Cuadros D, et al. A year of genomic surveillance reveals how the SARS-CoV-2 pandemic unfolded in Africa. *Science* 2021;374:423–31.
- [7] Zhang Z. KaKs_Calculator 3.0: calculating selective pressure on coding and non-coding sequences. *Genomics Proteomics Bioinformatics* 2022. <https://doi.org/10.1016/j.gpb.2021.12.002>.
- [8] Temmam S, Vongphayloth K, Salazar EB, Munier S, Bonomi M, Regnault B, et al. Coronaviruses with a SARS-CoV-2-like receptor-binding domain allowing ACE2-mediated entry into human cells isolated from bats of Indochinese peninsula. *Research Square* 2021. <https://doi.org/10.21203/rs.3.rs-871965/v1>.
- [9] Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, et al. A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* 2020;579:270–3.
- [10] Liu P, Jiang JZ, Wan XF, Hua Y, Li L, Zhou J, et al. Are pangolins the intermediate host of the 2019 novel coronavirus (SARS-CoV-2)? *PLoS Pathog* 2020;16:e1008421.
- [11] Lam TY, Jia N, Zhang YW, Shum MH, Jiang JF, Zhu HC, et al. Identifying SARS-CoV-2-related coronaviruses in Malayan pangolins. *Nature* 2020;583:282–5.
- [12] Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, et al. A new coronavirus associated with human respiratory disease in China. *Nature* 2020;579:265–9.
- [13] Wang Q, Zhang Y, Wu L, Niu S, Song C, Zhang Z, et al. Structural and functional basis of SARS-CoV-2 entry by using human ACE2. *Cell* 2020;181:894–904.e9.
- [14] Laffey C, de Koning K, Kanaar R, Lebbink JHG. Experimental evidence for enhanced receptor binding by rapidly spreading SARS-CoV-2 variants. *J Mol Biol* 2021;433:167058.
- [15] Liu Y, Liu J, Plante KS, Plante JA, Xie X, Zhang X, et al. The N501Y spike substitution enhances SARS-CoV-2 infection and transmission. *Nature* 2022;602:294–9.
- [16] Meng B, Kemp SA, Papa G, Datir R, Ferreira I, Marelli S, et al. Recurrent emergence of SARS-CoV-2 spike deletion H69/V70 and its role in the Alpha variant B.1.1.7. *Cell Rep* 2021;35:109292.
- [17] Lubinski B, Fernandes MHV, Frazier L, Tang T, Daniel S, Diel DG, et al. Functional evaluation of the P681H mutation on the proteolytic activation of the SARS-CoV-2 variant B.1.1.7 (Alpha) spike. *iScience* 2022;25:103589.
- [18] Wu H, Xing N, Meng K, Fu B, Xue W, Dong P, et al. Nucleocapsid mutations R203K/G204R increase the infectivity, fitness, and virulence of SARS-CoV-2. *Cell Host Microbe* 2021;29:1788–801.e6.
- [19] Hastie KM, Li H, Bedinger D, Schendel SL, Dennison SM, Li K, et al. Defining variant-resistant epitopes targeted by SARS-CoV-2 antibodies: a global consortium study. *Science* 2021;374:472–8.
- [20] Wilhelm A, Widera M, Grikscheit K, Toptan T, Schenk B, Pallas C, et al. Reduced neutralization of SARS-CoV-2 Omicron variant by vaccine sera and monoclonal antibodies. *medRxiv* 2021;21267432.
- [21] Cao Y, Wang J, Jian F, Xiao T, Song W, Yisimayi A, et al. Omicron escapes the majority of existing SARS-CoV-2 neutralizing antibodies. *Nature* 2022;602:657–63.
- [22] Greaney AJ, Starr TN, Barnes CO, Weisblum Y, Schmidt F, Caskey M, et al. Mapping mutations to the SARS-CoV-2 RBD that escape binding by different classes of antibodies. *Nat Commun* 2021;12:4196.
- [23] Barnes CO, Jette CA, Abernathy ME, Dam KM, Esswein SR, Gristick HB, et al. SARS-CoV-2 neutralizing antibody structures inform therapeutic strategies. *Nature* 2020;588:682–7.
- [24] Choi B, Choudhary MC, Regan J, Sparks JA, Padera RF, Qiu X, et al. Persistence and evolution of SARS-CoV-2 in an immunocompromised host. *N Engl J Med* 2020;383:2291–3.
- [25] Kemp SA, Collier DA, Datir RP, Ferreira I, Gayed S, Jahun A, et al. SARS-CoV-2 evolution during treatment of chronic infection. *Nature* 2021;592:277–82.
- [26] Oude Munnink BB, Sikkema RS, Nieuwenhuijse DF, Molenaar RJ, Munger E, Molenkamp R, et al. Transmission of SARS-CoV-2 on mink farms between humans and mink and back to humans. *Science* 2021;371:172–7.
- [27] Wei C, Shan KJ, Wang W, Zhang S, Huan Q, Qian W, et al. Evidence for a mouse origin of the SARS-CoV-2 Omicron variant. *J Genet Genomics* 2021;48:1111–21.
- [28] Kupferschmidt K. Where did “weird” Omicron come from? *Science* 2021;374:1179.
- [29] Rambaut A, Holmes EC, O’Toole Á, Hill V, McCrone JT, Ruis C, et al. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* 2020;5:1403–7.
- [30] Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 2018;34:i884–90.
- [31] Li H, Birol I. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 2018;34:3094–100.
- [32] Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 2009;25:2078–9.
- [33] Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, et al. VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res* 2012;22:568–76.
- [34] Kumar S, Stecher G, Li M, Knyaz C, Tamura K, Battistuzzi FU. MEGA X: molecular evolutionary genetics analysis across computing platforms. *Mol Biol Evol* 2018;35:1547–9.

-
- [35] Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, et al. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol* 2020;37:1530–4.
- [36] Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A, et al. BEAST 2.5: an advanced software platform for Bayesian evolutionary analysis. *PLoS Comput Biol* 2019;15:e1006650.
- [37] Song S, Ma L, Zou D, Tian D, Li C, Zhu J, et al. The global landscape of SARS-CoV-2 genomes, variants, and haplotypes in 2019nCoV-R. *Genomics Proteomics Bioinformatics* 2020;18:749–59.
- [38] Slatkin M. Linkage disequilibrium—understanding the evolutionary past and mapping the medical future. *Nat Rev Genet* 2008;9:477–85.