



Putative Host-Derived Insertions in the Genomes of Circulating SARS-CoV-2 Variants

Yiyang Yang,^a Keith Dufault-Thompson,^a Rafaela Salgado Fontenele,^a Xiaofang Jiang^a

^aNational Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA

ABSTRACT Insertions in the SARS-CoV-2 genome have the potential to drive viral evolution, but the source of the insertions is often unknown. Recent proposals have suggested that human RNAs could be a source of some insertions, but the small size of many insertions makes this difficult to confirm. Through an analysis of available direct RNA sequencing data from SARS-CoV-2-infected cells, we show that viral-host chimeric RNAs are formed through what are likely stochastic RNA-dependent RNA polymerase template-switching events. Through an analysis of the publicly available GISAID SARS-CoV-2 genome collection, we identified two genomic insertions in circulating SARS-CoV-2 variants that are identical to regions of the human 18S and 28S rRNAs. These results provide direct evidence of the formation of viral-host chimeric sequences and the integration of host genetic material into the SARS-CoV-2 genome, highlighting the potential importance of host-derived insertions in viral evolution.

IMPORTANCE Throughout the COVID-19 pandemic, the sequencing of SARS-CoV-2 genomes has revealed the presence of insertions in multiple globally circulating lineages of SARS-CoV-2, including the Omicron variant. The human genome has been suggested to be the source of some of the larger insertions, but evidence for this kind of event occurring is still lacking. Here, we leverage direct RNA sequencing data and SARS-CoV-2 genomes to show that host-viral chimeric RNAs are generated in infected cells and two large genomic insertions have likely been formed through the incorporation of host rRNA fragments into the SARS-CoV-2 genome. These host-derived insertions may increase the genetic diversity of SARS-CoV-2 and expand its strategies to acquire genetic material, potentially enhancing its adaptability, virulence, and spread.

KEYWORDS SARS-CoV-2, insertion, host-viral chimeric reads, ribosomal RNA

During the COVID-19 pandemic, insertions have been frequently acquired in SARS-CoV-2 lineages (1–4). Insertions have been associated with several globally circulating lineages, including the insertion of one amino acid at position 146 of the S protein (ins146N) of the variant of interest Mu (B.1.621) (4), insertions at the recurrent insertion site 214 of the N-terminal domain (NTD) region on the S protein that occurred in the lineages B.1.214.2 (ins214TDR) and A.2.5 (ins214AAG) (1), and the insertion ins214EPE in the recently emerged variant of concern, Omicron (5). Although there is insufficient evidence to show the direct impact these insertions have on viral spread and interference with immune responses, the fact that variants carrying those insertions have circulated for long periods suggests that they might be advantageous or neutral for transmission. Results from a long-term *in vitro* experiment where SARS-CoV-2 was coinoculated with highly neutralizing antibodies have also shown that an 11-amino-acid insertion (ins248KTRNKSTSRRE) at the NTD N5 loop of the S protein was able to drive antibody escape, suggesting a potential role of insertions in enhancing infectivity and virulence (6). Taken together, insertions have the potential to increase genetic diversity in SARS-CoV-2 and contribute to the continued evolution of the virus.

Previous research has shown that most small insertions in the SARS-CoV-2 genome likely originated from template sliding, local duplication, or template switching

Editor Sergio Baranzini, University of California, San Francisco

Ad Hoc Peer Reviewer Jan Postberg, HELIOS Medical Center Wuppertal, Witten/Herdecke University

This is a work of the U.S. Government and is not subject to copyright protection in the United States. Foreign copyrights may apply.

Address correspondence to Xiaofang Jiang, xiaofang.jiang@nih.gov.

The authors declare no conflict of interest.

Received 23 February 2022

Accepted 2 May 2022

Published 18 May 2022

between viruses (2). Longer insertions (equal to or larger than nine nucleotides) have been detected in multiple coronavirus genomes, including in variants of concern like the Omicron variant, but their origin remains unknown. Host genetic material has been suggested as a possible source for these insertions (5, 7). Venkatakrishnan et al. suggested that the unique insertion (ins214EPE) in the Omicron variant could have originated from the human common cold virus HCoV-229E or the human genome based on BLAST search (5), and the human genome has been speculated to be the source of multiple other small insertions (7). However, given that these insertion sequences are typically short, sequence comparisons tend to be less informative, and false-positive matches have a high chance of occurring. Additionally, coronavirus replication occurs in modified endoplasmic reticulum-derived double-membrane vesicles, providing a physical barrier between viral and host genetic material (8), and coronavirus replication complexes are known to contain enzymes with proofreading activity (9), both of which likely play roles in limiting the formation of host-viral chimeric sequences.

Human-derived insertions in the SARS-CoV-2 genome would likely be generated through RNA-dependent RNA polymerase (RdRp)-driven template-switching events between SARS-CoV-2 and host RNAs. While template-switching events between coronaviruses are common (10–13) and likely contribute to the emergence of SARS-CoV-2 lineages, including the Deltacron variant (14), template-switching events between coronaviruses and host RNAs are rarely documented (15, 16). Chimeric reads between SARS-CoV-2 RNA and human RNA have been detected but were interpreted as a signal of SARS-CoV-2 integration into the human genome in a previous controversial study (17). Others have suggested that the chimeric reads were likely to be template-switching artifacts mediated by reverse transcriptase or PCR during library preparation (18–21). One possible explanation that was largely omitted in these studies is that the SARS-CoV-2-host chimeric RNA could be generated by RdRp-driven template switching.

Here, to investigate the possible existence of SARS-CoV-2-host chimeric RNA, we take advantage of the publicly available Nanopore direct RNA sequencing data of SARS-CoV-2. Direct transcriptome sequencing (RNA-seq) sequences the individual polyadenylated RNAs, directly mitigating the possible formation of chimeric reads during library preparation or amplification. We first identified SARS-CoV-2-host chimeric RNA from direct RNA-seq data and showed that RdRp-driven template switching between SARS-CoV-2 and host mRNA occurs, but it is infrequent and stochastic. We also found that highly expressed host genes and structural RNA genes have a higher chance to be observed in chimeric RNA reads. We then systematically analyzed the SARS-CoV-2 genomes deposited in the GISAID (Global initiative on sharing all influenza data) database (22), resulting in the identification of two insertions in functional SARS-CoV-2 genomes that likely originated from the host 18S and 28S rRNAs.

RESULTS

Host-viral mRNA chimeras are rare but do exist. We first analyzed direct RNA-seq data from SARS-CoV-2-infected cell lines to identify sequences formed from chimeric host-viral RNAs. The direct RNA-seq data were quality filtered and mapped to both the host and SARS-CoV-2 transcriptomes to identify potential chimeric sequences. Out of the 30 samples that were analyzed, host-viral chimeric reads were detected in 16 of the samples, with an average of 0.029% (standard deviation, 0.048%) of the reads mapped to SARS-CoV-2 being chimeric (see Table S1 in the supplemental material). Chimeric reads were typically rare, making up 0.206% of one sample, but less than 0.06% of the other 15 samples, and these rates may be an overestimation due to the cell lines used compared to what would be observed in *in vivo* conditions. Additionally, chimeric reads detected in five samples were further investigated using paired-end sequencing short reads from the same samples (Table S2). Approximately 1.4% (5 out of 357) of chimeric reads were supported by at least five read pairs spanning the junctions. This finding implies that a small fraction of the host-viral chimeric mRNA molecules could function as the templates for RNA replication.

We then analyzed the chimeric reads to identify trends in how the viral and host RNA sequences were joined. All of the viral-derived sequences in chimeric reads were

annotated as positive-sense RNA, and a majority (92.49%) of the reads contained host-derived positive-sense sequences. Upon further examination, the few host reads that were identified as being negative sense were largely long noncoding RNAs that were present in the raw reads as the negative-sense sequences, making it likely that they were misannotated rather than actually being derived from negative-sense RNA. These results suggest that the host-viral chimeric sequences are not the result of the integration of the viral genetic material into the host genome, which would have resulted in a nearly equal mix of positive- and negative-sense viral sequences (17). Most likely, these host-viral chimeric sequences were created from positive- to positive-strand template-switching events (23, 24).

Viral-host chimeric read formation is likely a stochastic process. The chimeric reads were then analyzed to determine if there were any patterns in the composition of the sequences and in which positions relative to the references they were formed. Both viral-to-host and host-to-viral chimeric sequences were detected in the direct RNA-seq data, but the chimeric reads did not show a preference for either organization (Table S1). Both types of sequences were seen in approximately the same frequency, with viral-to-host reads making up 55% of the chimeric sequences and host-to-viral reads making up 45%. This lack of strong preference may indicate that host RNA can be readily recognized by viral RdRp, but other factors, like the exclusion of host RNA by the formation of the double-membrane vesicles, might prevent the formation of chimeric RNAs. When examining the positions of the junctions on the viral RNA sequences, we found there was a bias toward the junction sites being located in the dense coding region near the three-prime end of the sequence, with fewer junctions being identified in the ORF1ab genes, the largest region of the genome (Fig. 1). This is likely due to the ORF1ab region not being retained in the canonical SARS-CoV-2 subgenomic RNAs, resulting in fewer viral RNAs being synthesized with these regions that could form chimeric RNAs (25). It suggests that the process by which chimeric sequences are formed is likely stochastic, depending on the availability of template RNA molecules.

Previous studies have also found that indel formation and template-switching events preferentially occur in the loops and stems formed in the RNA secondary structure (2, 3). First, a permutation test was used to investigate if junction sites were commonly located in stems (positions that form base pairs) or nonstem regions (non-base-paired positions) in the viral RNA. The results of this test showed a significant ($P < 0.01$) preference for the formation of junctions in non-base-paired regions of the RNA secondary structure (Fig. 1). One-sided Fisher's exact tests were performed to explore if junction sites were enriched in specific types of RNA structures. Consistent with the results of the permutation test, stems were underrepresented at the junction sites (Table S3). We speculate that the non-base-paired regions of the SARS-CoV-2 RNA may be more susceptible to stochastic template-switching events due to their more "open" configurations, where the viral RdRp could easily attach or detach as it moves along the RNA.

An examination of the types of human gene sequences found in the chimeric sequences revealed an enrichment of noncoding RNAs and highly expressed genes. We found that a disproportionate number of noncoding RNAs, mainly long noncoding RNAs (lncRNAs), were forming parts of the chimeric reads compared to their abundance in the human genomes. These noncoding RNA chimeric sequences made up 8.8% and 10.5% of the chimeric reads detected in the Caco and Calu cell lines, respectively, while noncoding sequences made up only 4% of the genes annotated in the human genome. This enrichment of noncoding RNA chimeric sequences was tested using Fisher's exact test confirming that the trend was significant (Caco cells, odds ratio, 2.2, $P = 0.043$; Calu cells, odds ratio, 2.8, $P = 0.001$). When analyzed in the context of the expression level of the host genes in each sample, we also observed an enrichment for highly expressed genes forming parts of the chimeric sequences (Fig. 2). This enrichment was confirmed through the Mann-Whitney U tests showing that the trend was significant in the two human cell lines ($P < 2.2e-16$ for both) and the *Chlorocebus sabaeus* (green monkey) cell line ($P < 2.2e-16$). These results appear to highlight two groups of sequences that are forming chimeric RNAs, structural RNAs like lncRNAs,

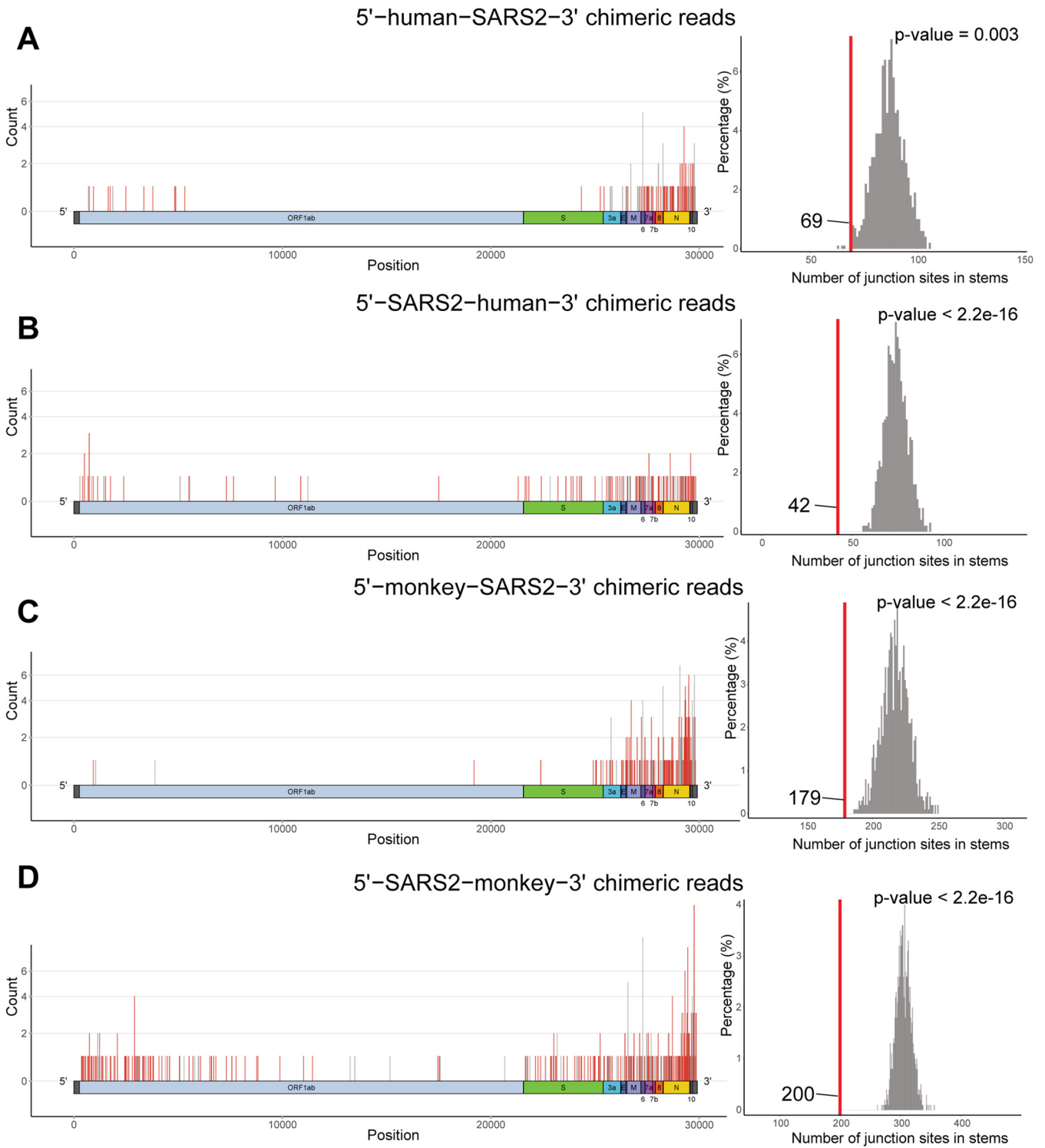


FIG 1 Locations of the chimeric read junction sites and permutation tests for the number of junction sites in stems. Diagrams show how frequently junction sites occur at each position on the SARS-CoV-2 genome for 5'-human-SARS2-3' (A), 5'-SARS2-human-3' (B), 5'-monkey-SARS2-3' (C), and 5'-SARS2-monkey-3' (D) chimeric reads. Positions are colored based on the secondary structure of the SARS-CoV-2 RNA, with red lines indicating that the position is in the nonstem region, while gray indicates that the position is located in the stem region. Histograms following each diagram show the corresponding results of permutation tests used to test if the junction sites of chimeric reads are within base-paired regions of the viral RNA. Each test consists of 1,000 permutations, and the actual frequency of junction sites occurring in the stem regions is marked with a vertical red line.

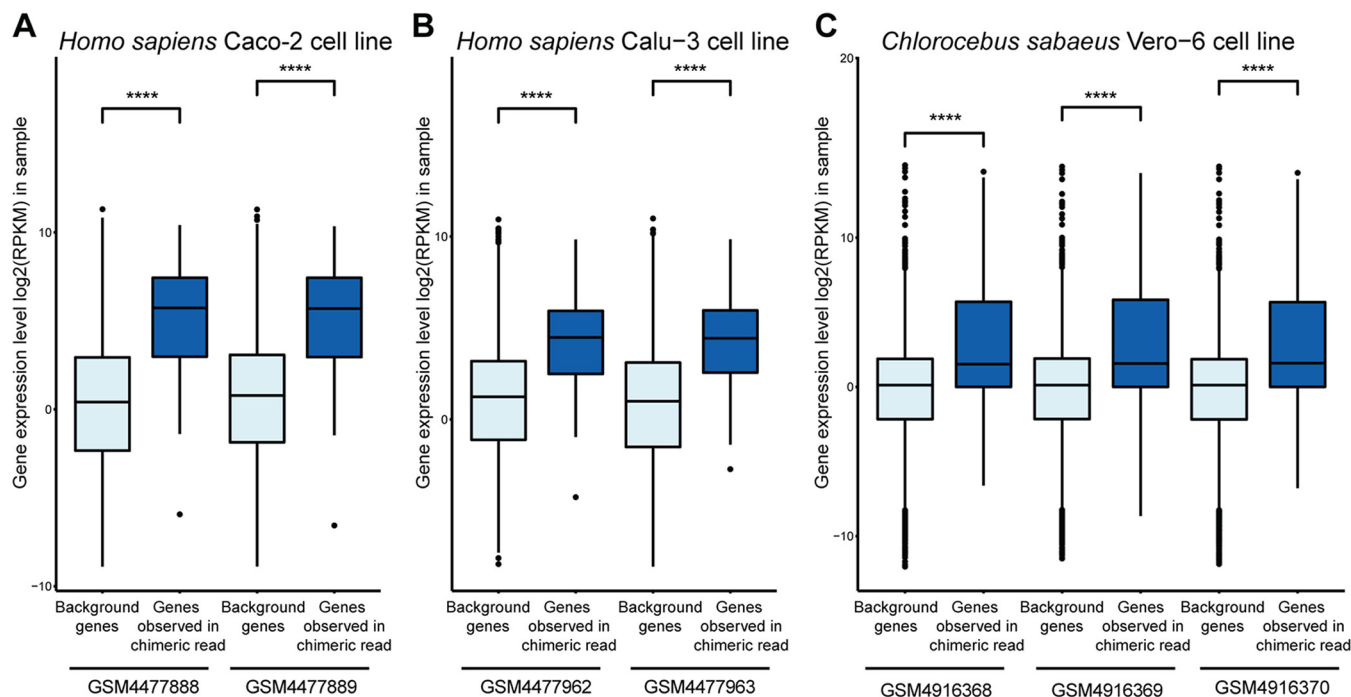


FIG 2 Expression levels of host genes observed in chimeric reads. The expression level of host protein-coding genes observed in chimeric reads is significantly higher than the background protein-coding gene expression level based on studies on *Homo sapiens* Caco-2 cell line (A), *Homo sapiens* Calu-3 cell line (B), and *Chlorocebus sabaesus* Vero-6 cell line (C).

which may be susceptible due to their secondary structures, and highly expressed genes, which would have more RNA molecules present for template-switching events to occur with. This suggests that the formation of chimeras is largely stochastic, with factors like the abundance of RNAs playing a large role, but that certain RNA molecules may be more susceptible to these events due to their structure.

A systematic search for host-derived insertions in SARS-CoV-2 genomes. We performed a survey of the GISAID SARS-CoV-2 genomes to identify insertions with potential host origins. Insertions were detected based on alignments and comparison to the Wuhan-Hu-1/2019 reference genome. Only insertions greater than or equal to 21 nucleotides long that were found outside the 5' and 3' untranslated regions were considered in subsequent analyses (Table S4). Of the 36 insertions that were found, 17 of them were found in multiple SARS-CoV-2 genomes but were not monophyletic. Upon further examination, the genomes containing these insertions tended to be sequenced by the same labs at around the same time, making it likely that these detected insertions are due to library preparation or sequencing errors rather than the result of multiple independent insertion events in different viral lineages. Of the 19 other insertions, 16 of them were only detected in a single genome, and while many of these had plausible hits to human genes, it is difficult to assess if these are true insertions or library preparation or sequencing artifacts due to their limited presence.

The three remaining insertions were from monophyletic virus variants and were further examined to determine if they had plausible homologous sequences in the human genome. Two of the insertions were found to be identical to conserved segments of the 28S and 18S rRNAs and were analyzed further. The remaining insertion was 21 nucleotides long and was found in six SARS-CoV-2 genomes of the Alpha B.1.1.7 lineage. These genomes were collected in early March of 2021 from England, United Kingdom, by two laboratories and sequenced at the same location using the same sequencing platform. The raw reads were available for two of the genomes, namely, England/ALDP-13C8C28/2021 (EPI_ISL_1331302) and England/QEUH-13C1955/2021 (EPI_ISL_1332461), and were examined directly, providing confirmation that the insertion was present and likely not an artifact. Unfortunately, no plausible source for this insertion

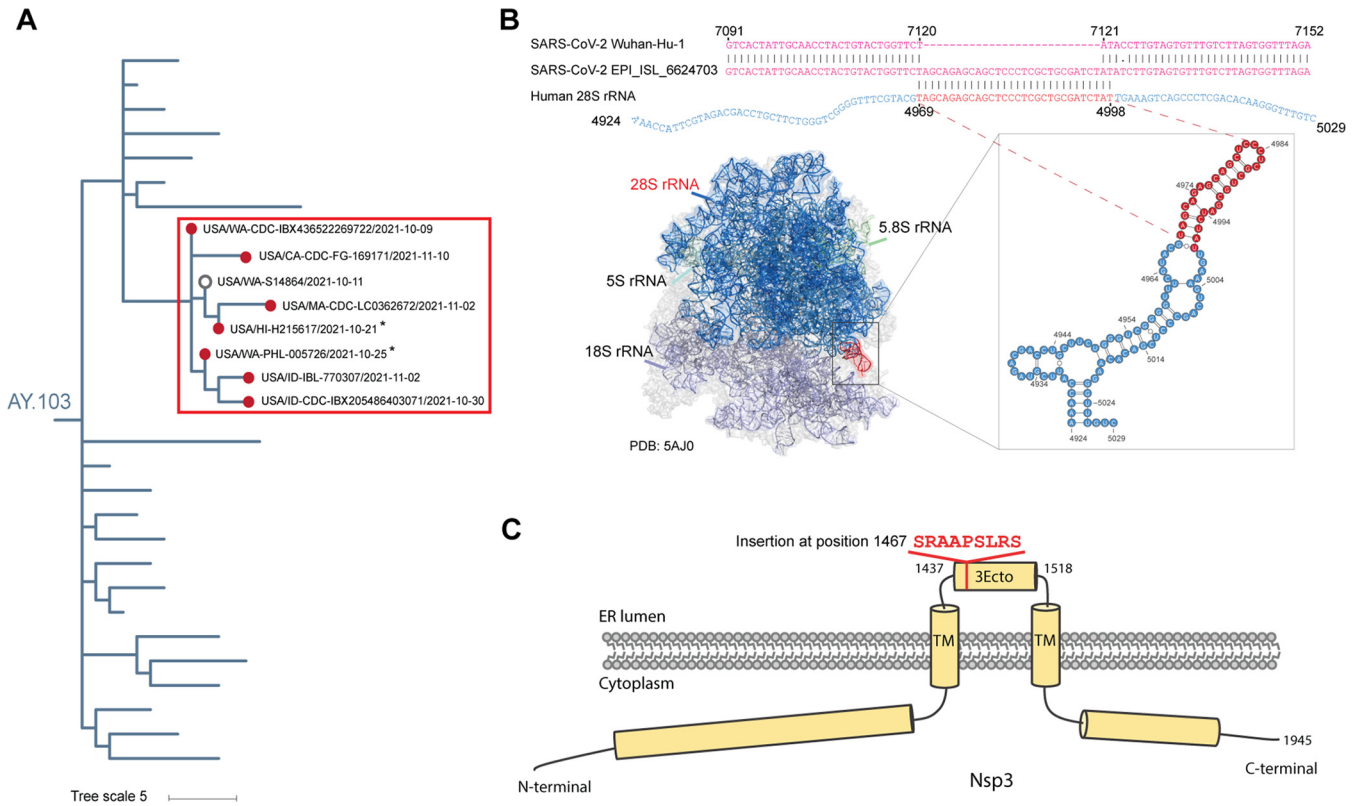


FIG 3 28S rRNA-derived insertion in SARS-CoV-2 genomes. (A) Phylogeny tree showing the genomes containing the human 28S-derived insertion. The clade where the insertion was detected is highlighted with a red box, and the genomes with the insertion are marked with red circles at the tips. Asterisk indicates that the insertion should be present in the variant based on raw sequencing data. (B) The insertion in SARS-CoV-2 genomes potentially originates from the host 28S rRNA, shown by the sequence alignment of SARS-CoV-2 reference genome (GenBank accession no. [NC_045512](#), GISAID accession no. China/Wuhan-Hu-1/2019) (pink), USA/CA-CDC-FG-169171/2021 (GenBank accession no. [OL591909](#), GISAID accession no. EPI_ISL_6624703) (pink), and human 28S rRNA (chain A2 of PDB ID [5AJ0](#)) (blue). There are five possible alignments for mapping this insertion to the reference. Only the alignment with the sequence inserted after the 3rd position of 2,285th codon in ORF1ab is shown. The putative insertion origin is colored in red. The numbers listed above and below the alignment indicate the positions of aligned bases in the original sequences. The insertion sequence (red) was mapped to the 28s rRNA (blue) in a human polysome three-dimensional structure (PDB ID [5AJ0](#)). A zoomed-in view of the RNA secondary structure shows that the insertion is located on the no. 94 stem of domain 7 (positions 4969 to 4998) 28S rRNA region (highlighted in red). (C) Diagram shows the position of the human 28S rRNA-derived insertion in the ectodomain (3Ecto) of the Nsp3 protein.

was able to be identified using a BLAST search in the NCBI nonredundant nucleotide database and a collection of coronavirus genomes with a cutoff E value of $1e-2$, and it was not analyzed further.

28S rRNA-derived insertion in SARS-CoV-2 genomes. We detected a 27-nucleotide-long insertion in five SARS-CoV-2 genomes (Table S5 and Fig. 3A) at position 7120 of the reference genome (China/Wuhan-Hu-1/2019). The five genomes containing the 28S rRNA-derived insertions were collected by different laboratories and were sequenced on different sequencing platforms, making it extremely unlikely that laboratory error is responsible for the presence of the insertions. The five genomes belong to a monophyletic group. In this clade, there are three other variants whose assembled genomes do not contain the insertion. We were able to obtain access to the raw genome sequencing data of two of the three variants, USA/WA-PHL-005726/2021 (EPI_ISL_6259191) and USA/HI-H215617/2021 (EPI_ISL_6540096). We then did further analysis on the raw sequencing to check if the insertion was indeed missing. First, we generated consensus genome sequences based on the alignment of sequencing reads to the SARS-CoV-2 reference genome and found the consensus sequences did not contain the insertion. Next, we manually added the 28S rRNA-derived insertion at position 7120 of the consensus genome and compared the reads exclusively aligned to the consensus genome with the insertion and the reads exclusively aligned to the consensus genome without the insertion. We found that 99.76% (8,700/8,721 for EPI_ISL_6259191) and 99.93% (1,502/1,503 for EPI_ISL_6540096)

of the exclusively mapped reads support the presence of the insertion in the genomes. The reason that the insertion is missing in the submitted genomes (EPI_ISL_6259191 and EPI_ISL_6540096) is likely that the assembly was generated using an insertion-unaware approach, such as reference-based consensus calling. For the only one variant that was missing the insertion in the genomes, we are not able to assess if it is due to failure to identify the insertion based on the consensus caller or the subsequent loss of the inserted sequence.

By performing the BLAST search for this insertion against the human transcripts (NCBI Homo sapiens Annotation Release 109 RNAs), an exact match (E value, $2e-06$) of this insertion was found in the nucleotide sequences of 28S rRNA (Fig. 3B). We observed an extra three overlapping bases in the pairwise alignment of SARS-CoV-2 variants containing the insertion and the human 28S rRNA sequence, extending the length of identity nucleotide bases from 27 nucleotides to 30 nucleotides. The identical region was located at positions 4969 to 4998 of the human 28S rRNA (based on the structure of PDB ID [5AJ0](#), chain A2) and makes up part of the highly conserved loop 94 stem of domain 7 of the rRNA molecule according to the Gorski et al.'s segmentation of human 28S rRNA (26) (Fig. 3B).

Due to the high level of sequence conservation of 28S rRNA, asserting the origin of the insertion-related 30 nucleotide sequences is impossible based on sequence identity alone. In the human genome (GRCh38 release 105), three 28S rRNA gene copies in chromosome 21 and one copy in chromosome 12 contain the exact 30 nucleotide sequences. When we searched the 30 nucleotide sequences in the LSU rRNA database downloaded from SILVA (27), 98 organisms were found to contain the sequences. The last common ancestor of these 98 organisms is *Euteleostomi* (bony vertebrates). Given the fact that the insertion emerged from the SARS-CoV-2 variant circulating in humans, the originating organism of the 28S rRNA-derived insertion is most likely humans.

The nine-amino-acid insertion is located at position 1467 of the ectodomain (3Ecto) in the Nsp3 protein, the only domain of this protein located on the luminal side of the endoplasmic reticulum (Fig. 3C). Nsp3, along with Nsp4 and Nsp6, has been shown to be involved in the formation of double-membrane vesicles in coronavirus-infected cells (28, 29). The 3Ecto domain is specifically involved in the recruitment of Nsp4 and has been shown to be an essential component of Nsp3 for correct double-membrane vesicle formation (28). At this point, it is unclear if this insertion would have had an effect on viral fitness, but given its location in the 3Ecto domain, it is possible that the insertion could have an effect on the interactions between Nsp3 and other proteins and on the membrane rearrangement process.

The monophyletic group with the 28S rRNA-derived insertion belonged to the AY.103 group of the Delta lineage (30) (Fig. 3A). The AY.103 variant was first detected worldwide on 1 January 2021 and in the United States on 2 January 2021. The clade containing the 28S rRNA-derived insertion is defined by five nucleotide mutations (T7900C, A10420T, C18646T, C25721T, and C29668T). By September 2021, AY.103 had become the most common Delta lineage in the United States and has continued to be responsible for a significant fraction of cases until the recent emergence of the Omicron variant (31). The five genomes containing the 28S rRNA-derived insertion were collected between 9 October and 10 November 2021 from the states of Washington, Idaho, Massachusetts, and California, indicating that these variants were likely being transmitted over this time frame. The extent to which it was being spread seems to be low, as Idaho was the only state where multiple genomes were collected from, and no genomes containing the insertion have been reported since. Based on the limited spread of the viruses containing the 28S rRNA-derived insertion, it is likely that the insertion might not confer phenotypic advantages or is possibly disadvantageous to the virus. Nonetheless, our data show that AY.103 lineages containing this insertion were viable and were transmitted for a short period of time.

18S rRNA-derived insertion in SARS-CoV-2 genomes. A 24-nucleotide insertion was detected in two genomes at position 27492 in the genome of the reference genome (China/Wuhan-Hu-1/2019) (Table S5). A sequence search against human transcripts (NCBI Homo sapiens Annotation Release 109 RNAs) was performed using BLAST

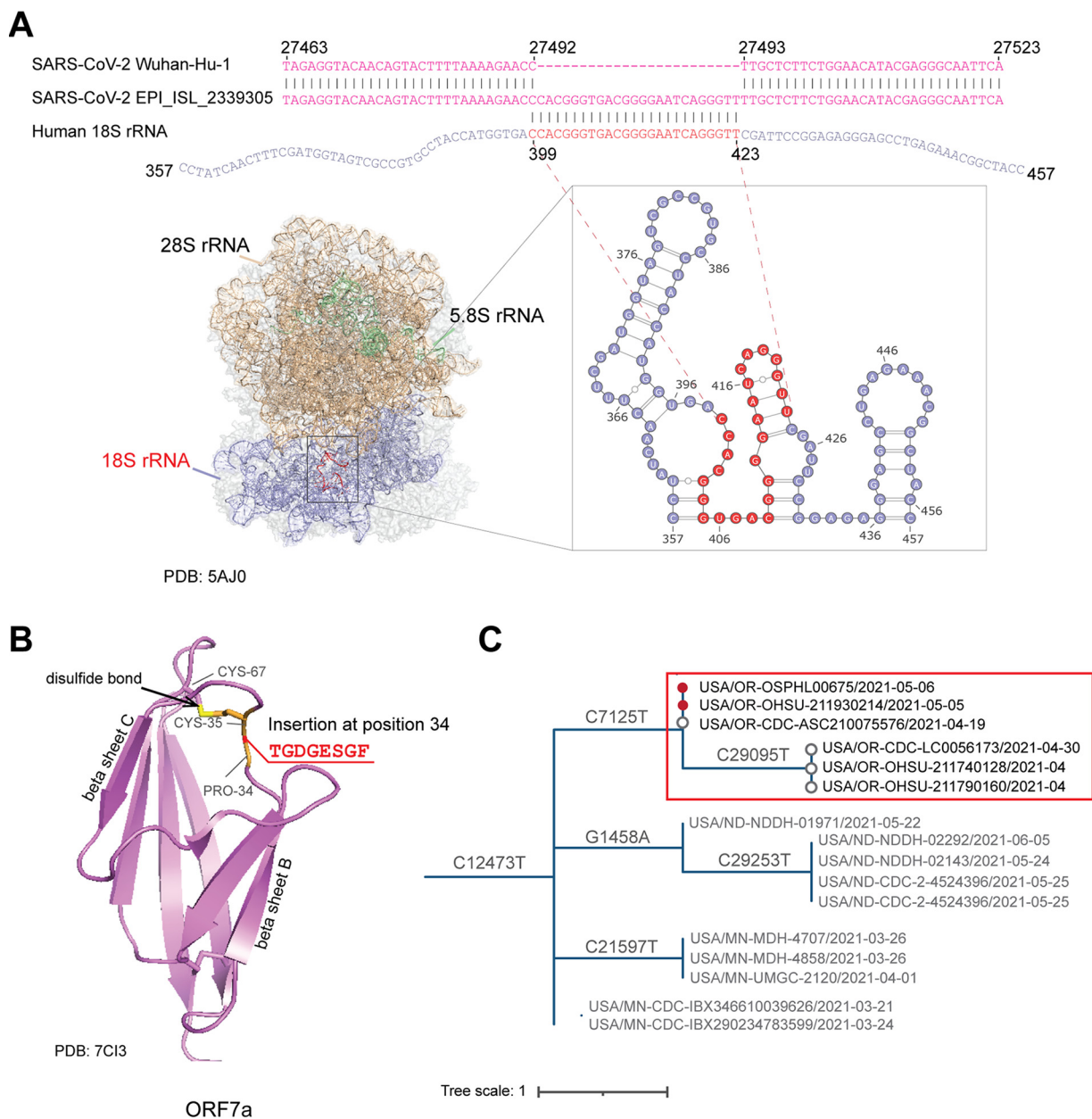


FIG 4 18S rRNA-derived insertion in SARS-CoV-2 genomes. (A) The insertion in SARS-CoV-2 genomes potentially originates from the host 18S rRNA, shown by the sequence alignment of SARS-CoV-2 reference genome (GenBank accession no. [NC_045512](#), GISAID accession no. China/Wuhan-Hu-1/2019) (pink), USA/OR-OSPHL00675/2021 (GISAID accession no. EPI_ISL_2339305) (pink), and human 18S rRNA (purple). The putative insertion origin is colored in red. The numbers listed above and below the alignment indicate the positions of aligned bases in the original sequences. The insertion sequence (red) was mapped to the 18S rRNA (purple) in a human polysome three-dimensional structure (PDB ID [5AJ0](#)). A zoomed-in view of the RNA secondary structure shows that the insertion covers parts of helices 11 and 12 of the 5' domain of the 18S rRNA. The location of the putative insertion sequence is highlighted red. (B) Diagram showing the position of the human 18S-derived insertion on the structure of the SARS-CoV-2 ORF7a protein (PDB ID [7CI3](#)). (C) Phylogeny tree showing the genomes containing the human 18S-derived insertion. The clade where the insertion was detected is highlighted with a red box, and the genomes with the insertion are marked with red circles.

(32), resulting in the identification of an exact match to a 24-nucleotide stretch (E value, $2e-5$) of the 18S rRNA sequence. When aligned to the full 18S rRNA sequence, it was found that the identical region extended one additional nucleotide outside the insertion region, bringing the identical stretch to 25 nucleotides (Fig. 4A). The insertion was identical to a highly conserved region of the 18S rRNA (at positions 399 to 423 in 18S rRNA), consisting of a portion of the helix 12 of the 5' domain (33, 34). In the human genome alone, there are five copies of the 18S rRNA gene on chromosome 21

that contain identical matches for this 25-nucleotide sequence. Compared to the SSU rRNA SILVA database (27), identical sequences were found in the 18S sequences of 2,289 organisms, which had a common ancestor of Opisthokonta (Fungi/Metazoa group). Considering that the viral samples were circulating in human populations, it is highly likely that the insertion was derived from human 18S rRNA.

The insertion is in the SARS-CoV-2 ORF7a protein, encoding an eight-amino-acid sequence that is located between the proline and cysteine at positions 34 and 35 in the reference protein sequence (Fig. 4B). The cysteine at position 35 is known to form a disulfide bond with a cysteine at position 67 and is thought to help stabilize the beta-sheet structure (35, 36), and the possible function of the proline at position 34 is not known. The ORF7a protein has been shown to contain an immunoglobulin-like ectodomain between residues 16 and 96 on the protein, which is thought to have a role in binding to human immune cells and modulating immune response (35–37). Given the proximity of the insert to the disulfide bond-forming cysteine at position 34 and the size of the insert, it is possible that this insert would have an effect on the overall structure and immunoregulatory functions of ORF7a, but without additional evidence, the effect of this insertion on the fitness of the virus remains unknown.

The two genomes containing the 18S rRNA insertion were from the same clade in the Alpha B.1.1.7 SARS-CoV-2 lineage, which was first identified in England, United Kingdom, in mid-December of 2020 (Fig. 4C). This variant was designated a variant of concern due to its transmissibility and a large number of mutations and quickly became the dominant variant in England while spreading to other countries (38). The genomes containing the 18S rRNA-derived insertion, along with the other four genomes in the same clade, were collected in April and May of 2021 in Oregon, United States. The genomes from the variants containing the insertion were collected and sequenced by different labs using different sequencing platforms, making it unlikely that the insertion was a sequencing or library preparation artifact. We did not detect the insertion in any of the other four genomes from this clade, indicating that either they do not have the insertion, they have it but it was not detected, or that the insertion was only acquired in a subclade within this group. After May of 2021, no new genomes containing this insertion were collected, indicating that the period during which these lineages were circulating may have been brief. While these viral variants seem to be viable and transmitted for a short period of time, the insertion likely does not confer a significant advantage or may be disadvantageous for the virus, resulting in its limited spread.

DISCUSSION

Insertions in the SARS-CoV-2 genome can be introduced through multiple mechanisms and have the potential to give rise to new variants with enhanced infectivity, pathogenicity, and antibody escape (2, 6), but the source of these insertions is often difficult to determine and has been hotly debated (5, 7). Leveraging available direct RNA sequencing data and an analysis of SARS-CoV-2 genomes, we have found evidence of the formation of viral-host chimeric RNA sequences and described two novel human-derived genomic insertions present in circulating variants of SARS-CoV-2.

Through our screening of direct RNA-seq data from SARS-CoV-2-infected cell lines, we found that viral-host chimeric RNAs were rare but were present in approximately half of the samples analyzed. The chimeric reads all contained positive-sense viral RNA sequences, indicating that these chimeric sequences are not the result of the integration of the viral genetic material into the host genome, which would have resulted in a nearly equal mix of positive- and negative-sense viral sequences (17). This process does appear to be stochastic in nature, though, with no preference for starting with host or viral sequences during chimera formation and a higher frequency of chimeras being formed with highly expressed genes in the cells. The regions in the RNA where these template-switching events occur appear to be influenced by the secondary structure of the viral RNA, possibly due to certain structures being more susceptible to template-switching events, similar to what has been reported in previous studies (2, 3).

The accurate determination of the exact junction boundaries and potential base pairings was hindered by the high error rate of 14% in direct RNA sequencing data and the limited number of host-viral chimeras detected in this study. The exact molecular basis for the viral-host chimera remains unclear, and future investigation with larger sets of error-corrected direct RNA-seq data of SARS-CoV-2 could be beneficial to address this question.

The formation of host-viral chimeric mRNAs or subgenomic RNAs could mostly be transient events, not having a long-term impact on viral fitness, but the possibility of human-derived insertions in the coronavirus genomes could have significant implications considering the role that genomic insertions seem to have in the evolution of new SARS-CoV-2 variants (5, 6). The putative 18S- and 28S-derived insertions were identified in circulating variants of the SARS-CoV-2, and while these particular variants did not seem to spread widely, they do provide evidence that human genetic material can be a source of genomic insertions in SARS-CoV-2. Interestingly, rRNAs have been established to be a source of insertions in influenza genomes, in some cases resulting in significantly more pathogenic viral variants (39, 40). It has been speculated that these recombination events often occur with host rRNAs due to their abundance in the cells, the presence of recombination hot spots on rRNA molecules, and the utilization of host rRNAs during viral replication (39). Similar factors may play a role in the formation of these rRNA-derived insertions in SARS-CoV-2, but the formation of double-membrane vesicles during SARS-CoV-2 would seemingly complicate this process. There may be accidental capture of host RNAs inside the double-membrane vesicles during their formation or some crossover of host RNA from the cytosol, but evidence of this is lacking and warrants further investigation.

Conclusions. Overall, our results suggest that viral-host chimeric sequences can be formed, likely through stochastic RdRp template-switching events. Furthermore, we have identified two long insertions in SARS-CoV-2 genomes in previously circulating variants which are likely derived from human ribosomal RNAs. While the source of smaller insertions that are present in many SARS-CoV-2 genomes are still difficult to identify due to their short lengths, these results provide evidence that bolsters the hypothesis that some of them are derived from human genetic material. The mechanisms at work in the formation of these chimeric RNAs and genomic insertions are still unclear but warrant further study, considering the potential importance of these processes in viral evolution and the emergence of new variants.

MATERIALS AND METHODS

Identification of host-viral chimeric reads in SARS-CoV-2 direct RNA-seq data. The Nanopore direct RNA-seq data from SARS-CoV-2-infected cell lines were downloaded from the NCBI SRA database (see Table S1 in the supplemental material). All reads were quality trimmed using NanoFilt v2.8.0 (41) to remove the first 50 nucleotides of each read and require an average quality score of at least 10 over the length of the read. The trimmed reads were then mapped using minimap2 v2.23 (42) to the SARS-CoV-2 reference genome (NCBI GenBank accession no. [NC_045512.2](https://www.ncbi.nlm.nih.gov/nuccore/NC_045512.2)) (43) and either a reference *Chlorocebus sabaues* transcriptome (ftp://ftp.ensembl.org/pub/release-105/fasta/chlorocebus_sabaeus/) or human transcriptome (ftp://ftp.ensembl.org/pub/release-105/fasta/homo_sapiens/). The mapping files were converted to the Pairwise mApping Format (PAF) using the `paftools` script that is part of minimap2 (42). Reads that mapped to both the host and SARS-CoV-2 transcriptomes were extracted for analysis as potential chimeric sequences. To avoid including chimeric reads that resulted from technical artifacts such as those caused by misinterpretation of open-pore states by base-calling software (19), additional quality filtering was applied to the chimeric reads. The distance between the mapped regions of the virus and the host sequence on the chimeric reads was required to be less than 15 nucleotides, the junction was required to be formed in the middle of the genes (not within the last 50 nucleotides of the first gene sequence nor the first 50 nucleotides of the second gene sequence), and the quality score within 20 bp of either side of the junction was required to be higher than the 20th percentile quality score for that read.

Mapping short reads to direct RNA-seq chimeric reads. We collected paired-end sequencing data on five samples with corresponding direct RNA sequencing data. The short reads were first preprocessed with `fastp` v0.23.1 (44) and then mapped to the chimeric reads from the same samples by using `minimap2` v2.23 (42) with options “-ax sr -w 5” to tolerate the high error rate of the Nanopore direct RNA sequencing reads (45). Read pairs spanning the junctions were detected and counted with a custom script. The numbers of read pairs supporting the chimeric reads are provided in Table S2.

Analysis of junction positions in relation to viral RNA secondary structure. The RNA secondary structure of the SARS-CoV-2 reference genome was obtained from previous studies (46, 47), and bpRNA

(48) was used to assign each residue to secondary structure elements. A junction site was considered in the stem if the two flanking nucleotides were in the same stem. To investigate if junctions tend to happen in nonstem regions, the number of junctions occurring in base-paired positions was calculated and compared with a background distribution for the numbers of junctions located in stems derived from a 1,000-time random sampling of the same number of sites along the viral RNA strand. To further examine which types of structural elements are over- or underrepresented at junction sites in viral-host chimeric reads and in host-viral chimeric reads, a one-sided Fisher's exact test was performed.

Analysis of the expression level of host genes observed in chimeric reads. Gene expression profiles for two SARS-CoV-2-infected Caco-2 cell line samples (GEO accession nos. [GSM4477888](#) and [GSM4477889](#)), two SARS-CoV-2-infected Calu-3 cell line samples (GEO accession nos. [GSM4477962](#) and [GSM4477963](#)), and three SARS-CoV-2-infected Vero-6 cell line samples (GEO accession nos. [GSM4916368](#), [GSM4916369](#), and [GSM4916370](#)) were downloaded from the GEO database. The read counts of each gene were normalized by the total number of reads in each sample and by the gene length (reads per kilobase per million [RPKM]) to represent the gene expression level. The background gene set was composed of all expressed protein-coding genes in the cell line. To evaluate whether the expression level of the host protein-coding genes in chimeric reads is significantly greater than the expression level of the background gene set, a one-sided Mann-Whitney U test was performed for each sample.

Identification of insertions in SARS-CoV-2 genomes. The SARS-CoV-2 genomes available at GISAID (<https://www.gisaid.org/>) on 17 December 2021 were downloaded for analysis ($n = 6,163,073$). The sequences were then processed by Nextclade CLI v1.7.0 (49), which generated a multiple-sequence alignment against the reference genome (Wuhan-Hu-1/2019) and provided a list of single nucleotide polymorphisms, insertions, and deletions associated with each genome sequence. Only sequences that passed all quality controls and were assessed as "good" applied by Nextclade were used for further analysis ($n = 5,226,229$). Insertions greater than or equal to 21 nucleotides long and found outside the 5' and 3' untranslated regions of the viral genomes were kept. They were searched in the NCBI nonredundant nucleotide database and a collection of coronavirus genomes with BLASTN (E value $\leq 1e-2$) (32) to explore their possible origins.

Monophyletic test. To check if the insertions of interest formed a monophyletic group, all genomes that contained the same insertion were analyzed using Ultrafast Sample placement on Existing tRee v0.5.1 (USHER) (50) against a phylogenetic tree with available genomes ($n = 6,257,569$) from GISAID, GenBank, COG-UK, and CNCB (the China National Center for Bioinformatics) generated by sarscov2phylo pipeline v13-11-20 (51). The sequences were placed within an updated global subsampled SARS-CoV-2 phylogenetic tree, and local subtrees were computed to show more sequences with the same context as the ones being analyzed.

Verification of insertions with raw sequencing data. The raw genome sequencing data of USA/WA-PHL-005726/2021 (EPI_ISL_6259191) and USA/HI-H215617/2021 (EPI_ISL_6540096) were analyzed to check if the insertion was indeed missing. The raw sequencing reads were processed for quality control using fastp v0.23.1 (44) with default parameters and mapped to the SARS-CoV-2 reference genome using BWA-MEM v0.7.17 (52). Primer sequences in reads of EPI_ISL_6259191 were soft clipped using iVar Trim (parameters, `-m 1 -q 0 -s 4 -e`), and reads in amplicons with variants in primer binding sites were removed by iVar removereads v1.3.1 (53). The sequencing data of EPI_ISL_6540096 were preprocessed by the providing laboratory, and the primers were removed. Consensus genome sequences were generated based on the alignments, and it was found that the consensus sequences did not contain the insertion. The 28S rRNA-derived insertion was manually added at position 7120 of the consensus genomes to generate consensus genomes with the insertion. The alignment files were converted to FASTQ format using SAMtools fastq command v1.14 (54) and realigned to the consensus genomes with or without the insertion using bowtie v2.4.4 (45) (parameter, `-xeq`). Reads exclusively aligned to the consensus genome with the insertion and exclusively aligned to the consensus genome without the insertion were identified with a custom script (https://github.com/ncbi/SARS2_host_derived_insertions/blob/main/verify_insertion_insertion_match_reads.py).

This work utilized the computational resources of the NIH HPC Biowulf cluster (<http://hpc.nih.gov>).

Data availability. The data sets generated in this study and scripts are available in the GitHub repository at https://github.com/ncbi/SARS2_host_derived_insertions.

SUPPLEMENTAL MATERIAL

Supplemental material is available online only.

TABLE S1, XLSX file, 0.01 MB.

TABLE S2, XLSX file, 0.04 MB.

TABLE S3, XLSX file, 0.1 MB.

TABLE S4, XLSX file, 0.01 MB.

TABLE S5, XLSX file, 0.1 MB.

TABLE S6, XLSX file, 0.01 MB.

ACKNOWLEDGMENTS

We thank Eugene V. Koonin and Sofya K. Garushyants for their thoughtful comments on our manuscript and their code on how to perform permutation tests and plot the results provided at https://github.com/garushyants/covid_insertions_paper. We gratefully

acknowledge the researchers from the originating laboratories responsible for obtaining the specimens and the submitting laboratories where genetic sequence data were generated and shared via the GISAID Initiative on which this research is based (see Table S6 in the supplemental material). We want to particularly thank the Washington State Public Health Laboratories and the State of Hawaii Laboratories Division for sharing the raw sequencing data for the genomes USA/WA-PHL-005726/2021 and USA/HI-H215617/2021.

We declare that we have no competing interests.

All authors are supported by the Intramural Research Program of the NIH, National Library of Medicine.

Y.Y. was involved in the execution of the analyses, interpretation of the results, and writing and revision of the manuscript. K.D.-T. was involved in the interpretation of the results and writing and revision of the manuscript. R.S.F. was involved in the execution of the analyses and writing of the manuscript. X.J. was involved in the conceptualization, planning, interpretation of the results, and revision of the manuscript. All authors read and approved the final manuscript.

REFERENCES

1. Gerdol M, Dishnica K, Giorgetti A. 2022. Emergence of a recurrent insertion in the N-terminal domain of the SARS-CoV-2 spike glycoprotein. *Virus Res* 310:198674. <https://doi.org/10.1016/j.virusres.2022.198674>.
2. Garushyants SK, Rogozin IB, Koonin EV. 2021. Template switching and duplications in SARS-CoV-2 genomes give rise to insertion variants that merit monitoring. *Commun Biol* 4:9. <https://doi.org/10.1038/s42003-021-02858-9>.
3. Chrisman BS, Paskov K, Stockham N, Tabatabaei K, Jung J-Y, Washington P, Varma M, Sun MW, Maleki S, Wall DP. 2021. Indels in SARS-CoV-2 occur at template-switching hotspots. *BioData Min* 14:20. <https://doi.org/10.1186/s13040-021-00251-0>.
4. Laiton-Donato K, Franco-Muñoz C, Álvarez-Díaz DA, Ruiz-Moreno HA, Usme-Ciro JA, Prada DA, Reales-González J, Corchuelo S, Herrera-Sepúlveda MT, Naizaque J, Santamaría G, Rivera J, Rojas P, Ortiz JH, Cardona A, Malo D, Prieto-Alvarado F, Gómez FR, Wiesner M, Martínez MLO, Mercado-Reyes M. 2021. Characterization of the emerging B. 1.621 variant of interest of SARS-CoV-2. *Infect Genet Evol* 95:105038. <https://doi.org/10.1016/j.meegid.2021.105038>.
5. Venkatakrishnan A, Anand P, Lenehan PJ, Suratekar R, Raghunathan B, Niesen MJ, Soundararajan V. 2021. Omicron variant of SARS-CoV-2 harbors a unique insertion mutation of putative viral or human genomic origin. *OSF Preprints* <https://doi.org/10.31219/osf.io/f7bxy>.
6. Andreano E, Piccini G, Licastro D, Casalino L, Johnson NV, Paciello I, Dal Monego S, Pantano E, Manganaro N, Manenti A, Manna R, Casa E, Hyseni I, Benincasa L, Montomoli E, Amaro RE, McLellan JS, Rappuoli R. 2021. SARS-CoV-2 escape from a highly neutralizing COVID-19 convalescent plasma. *Proc Natl Acad Sci U S A* 118:e2103154118. <https://doi.org/10.1073/pnas.2103154118>.
7. Peacock TP, Bauer DL, Barclay WS. 11 October 2021, posting date. Putative host origins of RNA insertions in SARS-CoV-2 genomes. <https://virological.org/t/putative-host-origins-of-rna-insertions-in-sars-cov-2-genomes/761>.
8. Knoops K, Kikkert M, Worm SHvd, Zevenhoven-Dobbe JC, Van Der Meer Y, Koster AJ, Mommaas AM, Snijder EJ. 2008. SARS-coronavirus replication is supported by a reticulovesicular network of modified endoplasmic reticulum. *PLoS Biol* 6:e226. <https://doi.org/10.1371/journal.pbio.0060226>.
9. Robson F, Khan KS, Le TK, Paris C, Demirbag S, Barfuss P, Rocchi P, Ng W-L. 2020. Coronavirus RNA proofreading: molecular basis and therapeutic targeting. *Mol Cell* 79:710–727. <https://doi.org/10.1016/j.molcel.2020.07.027>.
10. Banner LR, Lai MM. 1991. Random nature of coronavirus RNA recombination in the absence of selection pressure. *Virology* 185:441–445. [https://doi.org/10.1016/0042-6822\(91\)90795-D](https://doi.org/10.1016/0042-6822(91)90795-D).
11. Liao C, Lai M. 1992. RNA recombination in a coronavirus: recombination between viral genomic RNA and transfected RNA fragments. *J Virol* 66: 6117–6124. <https://doi.org/10.1128/JVI.66.10.6117-6124.1992>.
12. Yang Y, Yan W, Hall AB, Jiang X. 2021. Characterizing transcriptional regulatory sequences in coronaviruses and their role in recombination. *Mol Biol Evol* 38:1241–1248. <https://doi.org/10.1093/molbev/msaa281>.
13. Simon-Loriere E, Holmes EC. 2011. Why do RNA viruses recombine? *Nat Rev Microbiol* 9:617–626. <https://doi.org/10.1038/nrmicro2614>.
14. Bolze A, White S, Basler T, Dei Rossi A, Roychoudhury P, Greninger AL, Hayashibara K, Wyman D, Kil E, Dai H. 2022. Evidence for SARS-CoV-2 Delta and Omicron co-infections and recombination. *medRxiv* <https://doi.org/10.1101/2022.03.09.22272113>.
15. Forni D, Cagliani R, Clerici M, Sironi M. 2017. Molecular evolution of human coronavirus genomes. *Trends Microbiol* 25:35–48. <https://doi.org/10.1016/j.tim.2016.09.001>.
16. Yan B, Chakravorty S, Mirabelli C, Wang L, Trujillo-Ochoa JL, Chauss D, Kumar D, Lionakis MS, Olson MR, Wobus CE, Afzali B, Kazemian M. 2021. Host-virus chimeric events in SARS-CoV-2-infected cells are infrequent and artifactual. *J Virol* 95:e00294-21. <https://doi.org/10.1128/JVI.00294-21>.
17. Zhang L, Richards A, Barrasa MI, Hughes SH, Young RA, Jaenisch R. 2021. Reverse-transcribed SARS-CoV-2 RNA can integrate into the genome of cultured human cells and can be expressed in patient-derived tissues. *Proc Natl Acad Sci U S A* 118:e2105968118. <https://doi.org/10.1073/pnas.2105968118>.
18. Briggs E, Ward W, Rey S, Law D, Nelson K, Bois M, Ostrov N, Lee HH, Laurent JM, Mita P. 2021. Assessment of potential SARS-CoV-2 virus integration into human genome reveals no significant impact on RT-qPCR COVID-19 testing. *Proc Natl Acad Sci U S A* 118:e2113065118. <https://doi.org/10.1073/pnas.2113065118>.
19. Parry R, Gifford RJ, Lytras S, Ray SC, Coin LJ. 2021. No evidence of SARS-CoV-2 reverse transcription and integration as the origin of chimeric transcripts in patient tissues. *Proc Natl Acad Sci U S A* 118:e2109066118. <https://doi.org/10.1073/pnas.2109066118>.
20. Smits N, Rasmussen J, Bodea GO, Amarilla AA, Gerdes P, Sanchez-Luque FJ, Ajjikuttira P, Modhiran N, Liang B, Faivre J, Deveson IW, Khromykh AA, Watterson D, Ewing AD, Faulkner GJ. 2021. No evidence of human genome integration of SARS-CoV-2 found by long-read DNA sequencing. *Cell Rep* 36:109530. <https://doi.org/10.1016/j.celrep.2021.109530>.
21. Zhang L, Richards A, Barrasa MI, Hughes SH, Young RA, Jaenisch R. 2021. Response to Parry et al.: strong evidence for genomic integration of SARS-CoV-2 sequences and expression in patient tissues. *Proc Natl Acad Sci U S A* 118:e2109497118. <https://doi.org/10.1073/pnas.2109497118>.
22. Shu Y, McCauley J. 2017. GISAID: Global Initiative on Sharing All Influenza Data—from vision to reality. *Eurosurveillance* 22:30494. <https://doi.org/10.2807/1560-7917.ES.2017.22.13.30494>.
23. Wang D, Jiang A, Feng J, Li G, Guo D, Sajid M, Wu K, Zhang Q, Ponty Y, Will S, Liu F, Yu X, Li S, Liu Q, Yang X-L, Guo M, Li X, Chen M, Shi Z-L, Lan K, Chen Y, Zhou Y. 2021. The SARS-CoV-2 subgenome landscape and its novel regulatory features. *Mol Cell* 81:2135–2147.e5. <https://doi.org/10.1016/j.molcel.2021.02.036>.
24. Wu H-Y, Brian DA. 2007. 5'-Proximal hot spot for an inducible positive-to-negative-strand template switch by coronavirus RNA-dependent RNA polymerase. *J Virol* 81:3206–3215. <https://doi.org/10.1128/JVI.01817-06>.
25. Kim D, Lee J-Y, Yang J-S, Kim JW, Kim VN, Chang H. 2020. The architecture of SARS-CoV-2 transcriptome. *Cell* 181:914–921.e10. <https://doi.org/10.1016/j.cell.2020.04.011>.

26. Gorski JL, Gonzalez IL, Schmickel RD. 1987. The secondary structure of human 28S rRNA: the structure and evolution of a mosaic rRNA gene. *J Mol Evol* 24:236–251. <https://doi.org/10.1007/BF02111237>.
27. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. 2013. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res* 41: D590–D596. <https://doi.org/10.1093/nar/gks1219>.
28. Hagemeyer MC, Monastyrska I, Griffith J, van der Sluijs P, Voortman J, En Henegouwen PMvB, Vonk AM, Rottier PJ, Reggiori F, De Haan CA. 2014. Membrane rearrangements mediated by coronavirus nonstructural proteins 3 and 4. *Virology* 458–459:125–135. <https://doi.org/10.1016/j.virol.2014.04.027>.
29. Lei J, Kusov Y, Hilgenfeld R. 2018. Nsp3 of coronaviruses: structures and functions of a large multi-domain protein. *Antiviral Res* 149:58–74. <https://doi.org/10.1016/j.antiviral.2017.11.001>.
30. Khare S, Gurry C, Freitas L, Schultz MB, Bach G, Diallo A, Akite N, Ho J, Lee RT, Yeo W, Curation Team GC, Maurer-Stroh S. 2021. GISAID's role in pandemic response. *China CDC Wkly* 3:1049–1051. <https://doi.org/10.46234/ccdcw2021.255>.
31. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P, Bedford T, Neher RA. 2018. Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* 34:4121–4123. <https://doi.org/10.1093/bioinformatics/bty407>.
32. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10:9. <https://doi.org/10.1186/1471-2105-10-421>.
33. Gopanenko AV, Malygin AA, Karpova GG. 2015. Exploring human 40S ribosomal proteins binding to the 18S rRNA fragment containing major 3'-terminal domain. *Biochim Biophys Acta* 1854:101–109. <https://doi.org/10.1016/j.bbapap.2014.11.001>.
34. Granneman S, Petfalski E, Swiatkowska A, Tollervey D. 2010. Cracking pre-40S ribosomal subunit structure by systematic analyses of RNA-protein cross-linking. *EMBO J* 29:2026–2036. <https://doi.org/10.1038/emboj.2010.86>.
35. Cao Z, Xia H, Rajsbaum R, Xia X, Wang H, Shi P-Y. 2021. Ubiquitination of SARS-CoV-2 ORF7a promotes antagonism of interferon response. *Cell Mol Immunol* 18:746–748. <https://doi.org/10.1038/s41423-020-00603-6>.
36. Zhou Z, Huang C, Zhou Z, Huang Z, Su L, Kang S, Chen X, Chen Q, He S, Rong X, Xiao F, Chen J, Chen S. 2021. Structural insight reveals SARS-CoV-2 ORF7a as an immunomodulating factor for human CD14+ monocytes. *IScience* 24:102187. <https://doi.org/10.1016/j.isci.2021.102187>.
37. Su C-M, Wang L, Yoo D. 2021. Activation of NF- κ B and induction of proinflammatory cytokine expressions mediated by ORF7a protein of SARS-CoV-2. *Sci Rep* 11:1–12. <https://doi.org/10.1038/s41598-021-92941-2>.
38. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, Hinsley WR, Laydon DJ, Dabrera G, O'Toole Á, Amato R, Ragonnet-Cronin M, Harrison I, Jackson B, Ariani CV, Boyd O, Loman NJ, McCrone JT, Gonçalves S, Jorgensen D, Myers R, Hill V, Jackson DK, Gaythorpe K, Groves N, Sillitoe J, Kwiatkowski DP, Flaxman S, Ratmann O, Bhatt S, Hopkins S, Gandy A, Rambaut A, Ferguson NM, COVID-19 Genomics UK (COG-UK) Consortium. 2021. Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7 in England. *Nature* 593:266–269. <https://doi.org/10.1038/s41586-021-03470-x>.
39. Gulyaev AP, Spronken MI, Funk M, Fouchier RA, Richard M. 2021. Insertions of codons encoding basic amino acids in H7 hemagglutinins of influenza A viruses occur by recombination with RNA at hotspots near snoRNA binding sites. *RNA* 27:123–132. <https://doi.org/10.1261/rna.077495.120>.
40. Khatchikian D, Orlich M, Rott R. 1989. Increased viral pathogenicity after insertion of a 28S ribosomal RNA sequence into the haemagglutinin gene of an influenza virus. *Nature* 340:156–157. <https://doi.org/10.1038/340156a0>.
41. De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. 2018. NanoPack: visualizing and processing long-read sequencing data. *Bioinformatics* 34:2666–2669. <https://doi.org/10.1093/bioinformatics/bty149>.
42. Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34:3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.
43. Wu F, Zhao S, Yu B, Chen Y-M, Wang W, Song Z-G, Hu Y, Tao Z-W, Tian J-H, Pei Y-Y, Yuan M-L, Zhang Y-L, Dai F-H, Liu Y, Wang Q-M, Zheng J-J, Xu L, Holmes EC, Zhang Y-Z. 2020. A new coronavirus associated with human respiratory disease in China. *Nature* 579:265–269. <https://doi.org/10.1038/s41586-020-2008-3>.
44. Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 34:i884–i890. <https://doi.org/10.1093/bioinformatics/bty560>.
45. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9:357–359. <https://doi.org/10.1038/nmeth.1923>.
46. Cao C, Cai Z, Xiao X, Rao J, Chen J, Hu N, Yang M, Xing X, Wang Y, Li M, Zhou B, Wang X, Wang J, Xue Y. 2021. The architecture of the SARS-CoV-2 RNA genome inside virion. *Nat Commun* 12:14. <https://doi.org/10.1038/s41467-021-22785-x>.
47. Huston NC, Wan H, Strine MS, Tavares RCA, Wilen CB, Pyle AM. 2021. Comprehensive in vivo secondary structure of the SARS-CoV-2 genome reveals novel regulatory motifs and mechanisms. *Mol Cell* 81:584–598. e5. <https://doi.org/10.1016/j.molcel.2020.12.041>.
48. Danaee P, Rouches M, Wiley M, Deng D, Huang L, Hendrix D. 2018. bpRNA: large-scale automated annotation and analysis of RNA secondary structure. *Nucleic Acids Res* 46:5381–5394. <https://doi.org/10.1093/nar/gky285>.
49. Aksamentov I, Roemer C, Hodcroft EB, Neher RA. 2021. Nextclade: clade assignment, mutation calling and quality control for viral genomes. *Joss* 6:3773. <https://doi.org/10.21105/joss.03773>.
50. Turakhia Y, Thornlow B, Hinrichs AS, De Maio N, Gozashiti L, Lanfear R, Haussler D, Corbett-Detig R. 2021. Ultrafast Sample placement on Existing tRees (USHER) enables real-time phylogenetics for the SARS-CoV-2 pandemic. *Nat Genet* 53:809–816. <https://doi.org/10.1038/s41588-021-00862-7>.
51. Lanfear R, Mansfield R. 2020. A global phylogeny of SARS-CoV-2 sequences from GISAID. *Zenodo* 10. <https://doi.org/10.5281/zenodo.3958883>.
52. Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25:1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
53. Grubaugh ND, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, Main BJ, Tan AL, Paul LM, Brackney DE, Grewal S, Gurfield N, Van Rompay KKA, Isern S, Michael SF, Coffey LL, Loman NJ, Andersen KG. 2019. An amplicon-based sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol* 20:8–19. <https://doi.org/10.1186/s13059-018-1618-7>.
54. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.