

Refining SARS-CoV-2 Intra-host Variation by Leveraging Large-scale Sequencing Data

Fatima Mostefai^{1,2,3}, Jean-Christophe Grenier², Raphaël Poujol², Julie G. Hussin^{2,3,4*}

¹ Département de Biochimie et de Médecine Moléculaire, Université de Montréal, Québec, Canada

² Montreal Heart Institute, Québec, Canada

³ Mila - Quebec AI Institute, Université de Montréal, Québec, Canada

⁴ Département de Médecine, Université de Montréal, Québec, Canada

* Corresponding Author (julie.hussin@umontreal.ca)

Abstract

Understanding the evolution of viral genomes is essential for elucidating how viruses adapt and change over time. Analyzing intra-host single nucleotide variants (iSNVs) provides key insights into the mechanisms driving the emergence of new viral lineages, which are crucial for predicting and mitigating future viral threats. Despite the potential of next-generation sequencing (NGS) to capture these iSNVs, the process is fraught with challenges, particularly the risk of capturing sequencing artifacts that may result in false iSNVs. To tackle this issue, we developed a workflow designed to enhance the reliability of iSNV detection in large heterogeneous collections of NGS libraries. We use over 130,000 publicly available SARS-CoV-2 NGS libraries to show how our comprehensive workflow effectively distinguishes emerging viral mutations from sequencing errors. This approach incorporates rigorous bioinformatics protocols, stringent quality control metrics, and innovative usage of dimensionality reduction methods to generate representations of this high-dimensional dataset. We identified and mitigated batch effects linked to specific sequencing centers around the world and introduced quality control metrics that consider strand coverage imbalance, enhancing iSNV reliability. Additionally, we pioneer the application of the PHATE visualization approach to genomic data and introduce a methodology that quantifies how related groups of data points are within a

28 two-dimensional space, enhancing our ability to explain clustering patterns based on their
29 shared genetic characteristics. Our workflow sheds light on the complexities of viral ge-
30 nomic analysis with state-of-the-art sequencing technologies and advances the detection
31 of accurate intra-host mutations, opening the door for an enhanced understanding of viral
32 adaptation mechanisms.

33 1 Introduction

34 The advancements in high-throughput sequencing technologies have revolutionized the study
35 of viral genomes, particularly evident in the case of SARS-CoV-2 during the COVID-19 pan-
36 demic. The ability to track the virus's mutations and evolution during host infection is critical
37 in understanding the emergence of various variants of concern (VOCs). These VOCs, result-
38 ing from the accumulation of mutations, demonstrate the importance of selective pressures
39 both within an individual host (intra-host) and during transmission between hosts (inter-host)
40 (Lauring 2020). This complex interplay is key to the evolution of viral lineages, influenced
41 by factors like error-prone replications and host RNA-editing mechanisms (Di Giorgio et al.
42 2020). In the current literature, there are several hypotheses to explain the interplay between
43 intra-host and inter-host dynamics in the development of SARS-CoV-2 VOCs (Markov et al.
44 2023). These hypotheses include evolution within chronically infected individuals (Sonnleit-
45 ner et al. 2022; Quaranta et al. 2022; Hill et al. 2022; Ghafari et al. 2022; Oude Munnink
46 et al. 2021; Hale et al. 2022; Oreshkova et al. 2020; Bashor et al. 2021), spillovers from animal
47 populations (Washburne et al. 2022; Sacchetto et al. 2021; Robinson et al. 2023; Goldberg et
48 al. 2023; Rajendran et al. 2022), and emergence in regions with limited genomic surveillance.
49 Understanding these processes is vital to explain the rapid evolution of VOCs such as Delta
50 and Omicron, which have shown significant evolutionary leaps.

51 In response to the pandemic, a vast number of next-generation sequencing (NGS) libraries
52 for SARS-CoV-2 have been generated, primarily to construct consensus sequences for tracking
53 inter-host mutations and VOCs. However, they also provide valuable insights into intra-
54 host diversity, enabling the identification of intra-host single nucleotide variants (iSNVs) that
55 are key in exploring hypotheses around VOC emergence. Despite the considerable size and
56 breadth of available NGS libraries, the existing body of research on iSNV analysis remains
57 limited, with the majority of studies focusing on a relatively small number of NGS libraries
58 (Sun et al. 2023; Messali et al. 2023; Xi et al. 2023; Sun et al. 2023; Armero et al. 2021;
59 Wertheim et al. 2022; Y. Wang et al. 2021; Sonnleitner et al. 2022; Zhang et al. 2022; Quaranta
60 et al. 2022). This gap in research can also be attributed to challenges related to data quality,
61 such as the presence of sequencing artifacts that introduce errors and lead to false iSNVs.
62 To mitigate these challenges, current practices in intra-host viral analysis include the use of

63 technical replicates (Zhang et al. 2022), which, while effective, are resource-intensive, and the
64 application of hard filters on coverage and frequency, which lack uniformity across different
65 studies and often overlook noteworthy sequencing artifacts like strand bias (Roder et al. 2023;
66 Armero et al. 2021; Hedskog et al. 2010; Bull et al. 2011; Tonkin-Hill et al. 2021). This concern
67 underscores the need for more standardized methodologies in processing complex sequencing
68 data to ensure accurate and reliable iSNV analysis.

69 In computational biology, dimensionality reduction techniques are commonly used to sim-
70 plify the representation and analysis of complex datasets, like viral sequencing data, and
71 to uncover inherent data biases. These techniques, which have seen significant improve-
72 ments with the rise of high-dimensional data, include Principal Component Analysis (PCA)
73 (Novembre et al. 2008), often used for summarizing human genetic data, t-SNE (Platzer 2013;
74 Tamazian et al. 2022) for analyzing local structures, and PHATE (Moon et al. 2019), a novel
75 method that allows visualization of both global and local structures in high-dimensional data.
76 Despite the potential of these methods, their application to the extensive SARS-CoV-2 data
77 has been limited, often confined to analyzing consensus sequences (Hozumi et al. 2021; B.
78 Wang et al. 2021; Mostefai et al. 2022). This gap highlights an opportunity for a broader
79 application of the dimensionality reduction methods on viral genome data.

80 Here, we address this gap by using a comprehensive set of publicly available SARS-CoV-2
81 NGS libraries from the NCBI database, representing the pandemic’s initial years. We use a
82 combination of bioinformatics tools, stringent quality control measures, and dimensionality
83 reduction methods such as PHATE and t-SNE to identify intra-host mutations from sequenc-
84 ing artifacts. Our approach provides a workflow for analyzing SARS-CoV-2 sequencing data
85 and establishes adapted thresholds for the 2020 and 2021 datasets. In this study, we estab-
86 lish a framework for rapid and precise analysis of intra-host viral data, aiming to support
87 pandemic preparedness and response.

88 **2 Results**

89 **2.1 Curation Pipeline Overview**

90 While NGS data offers valuable insights into viral diversity and evolution, extracting meaning-
91 ful information demands rigorous bioinformatics and representation approaches. A systematic

92 methodology is crucial to process this data accurately, ensuring the reliability of identified
93 iSNVs. We, therefore, propose a comprehensive workflow to extract meaningful intra-host
94 mutations from NGS data. Our workflow is divided into two levels (Figure 1): the processing
95 and quality control of a set of libraries (Figure 1A) and the processing and quality control of
96 iSNVs within each library (Figure 1B).

97 To build a set of high-quality libraries, we meticulously processed a large set of Illumina
98 amplicon paired-end sequencing libraries, ensuring a representative sample across various time
99 points and locations. The data processing includes adapter and quality trimming, alignment
100 to the SARS-CoV-2 reference genome, primer trimming, and whole genome coverage quality
101 control (see Method section 4.1). Using the processed libraries, we performed iSNV calling and
102 computed key metrics such as Alternative Allele Frequency (*AAF*) and Strand bias likelihood
103 (*S*) (see Methods). These metrics help to accurately identify putative iSNVs while minimizing
104 artifacts. Next, dimensionality reduction methods, such as PHATE and t-SNE, are applied
105 to visualize and interpret the iSNV data through analyses of clustering structures. In this
106 process, we generate representations in two distinct ways: by the library, where each point in
107 the visualization represents a summary of the library, and by genomic position, where each
108 point corresponds to a specific genomic position summarizing its behaviour across libraries.
109 PHATE maintains meaningful distance between clusters (Moon et al. 2019), preserving hier-
110 archical relationships between sequencing libraries, we therefore use this technique to present
111 our main results. Similar findings are observed using t-SNE (see supplementary information
112 section 10.5). To differentiate between potential artifacts and biologically relevant patterns,
113 the clustering structures are measured using the Percentage of Nearest Neighbors (*PNN*)
114 presenting the same lineage label (as defined by the World Health Organization, WHO) or
115 sequencing center (SC), providing a robust metric to quantify clustering structures of different
116 sets of iSNVs. While the null hypothesis is for iSNVs to be randomly distributed, SC labels
117 are used to check for associations with sequencing centers as a proxy to identify potential arti-
118 facts. Conversely, WHO labels are expected to reflect biological relevance, with the limitation
119 that some sequencing centers favour the sequencing of some lineages over others.

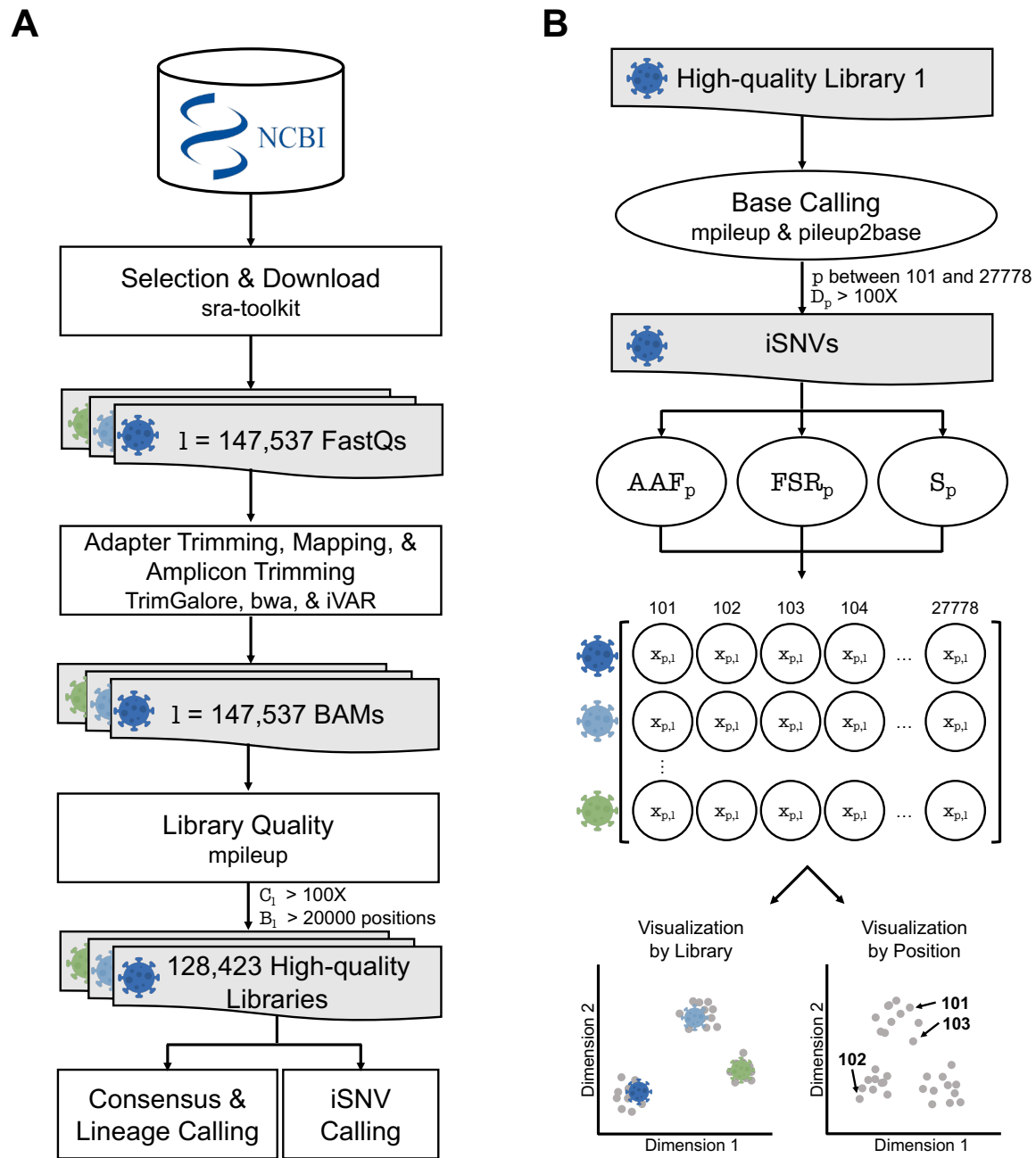


Figure 1: SARS-CoV-2 Sequencing Library Processing Workflows. **A**: Processing workflow of a set of SARS-CoV-2 sequencing libraries. The workflow starts with selecting and downloading 147,537 FastQ libraries from NCBI. Next, these libraries were trimmed for adapters, mapped to a reference, and trimmed again for primer targets. We set whole genome coverage filters of mean depth (C) $> 100X$ and breadth of coverage (B) > 20000 on each library l to keep only high-quality libraries, keeping 128,423 libraries for further analysis. **B** Processing workflow of a single RNA sequencing library. Within each high-quality sequencing library, base calling was done to extract iSNVs. During this process, the ends of the genome were removed, keeping genomic positions p between 101 and 27778 and the depth $D > 100X$ of p . Subsequently, for each iSNV, we computed the following quality metrics: Alternative Allele Frequency (AAF), Forward Strand Ratio (FSR), and Strand bias likelihood (S , equation 1). Thresholds for AAF and S metrics were established using dimensionality reduction visualization methods, reducing the data into two dimensions by either the libraries (left) or the genomic positions (right).

120 2.2 Extracting Emerging *de novo* iSNVs

121 We processed 128,423 high-quality SARS-CoV-2 sequencing libraries from the first two years
 122 of the COVID-19 pandemic, ensuring a representative sampling across time and geographic
 123 locations (Figure 2A). Genome data quality was assessed by coverage depth (Figure 2B,
 124 x-axis), breadth (Figure 2B, y-axis), and strand balance (Figure 2C). The distribution of
 125 depth and breadth of coverage reveals center-specific quality variability. In turn, the strand
 126 balance coverage shows unbalanced strand coverage across the genome with an oscillating
 127 pattern at the same genomic regions independent of the sequencing center (see supplementary
 128 information section 10.2).

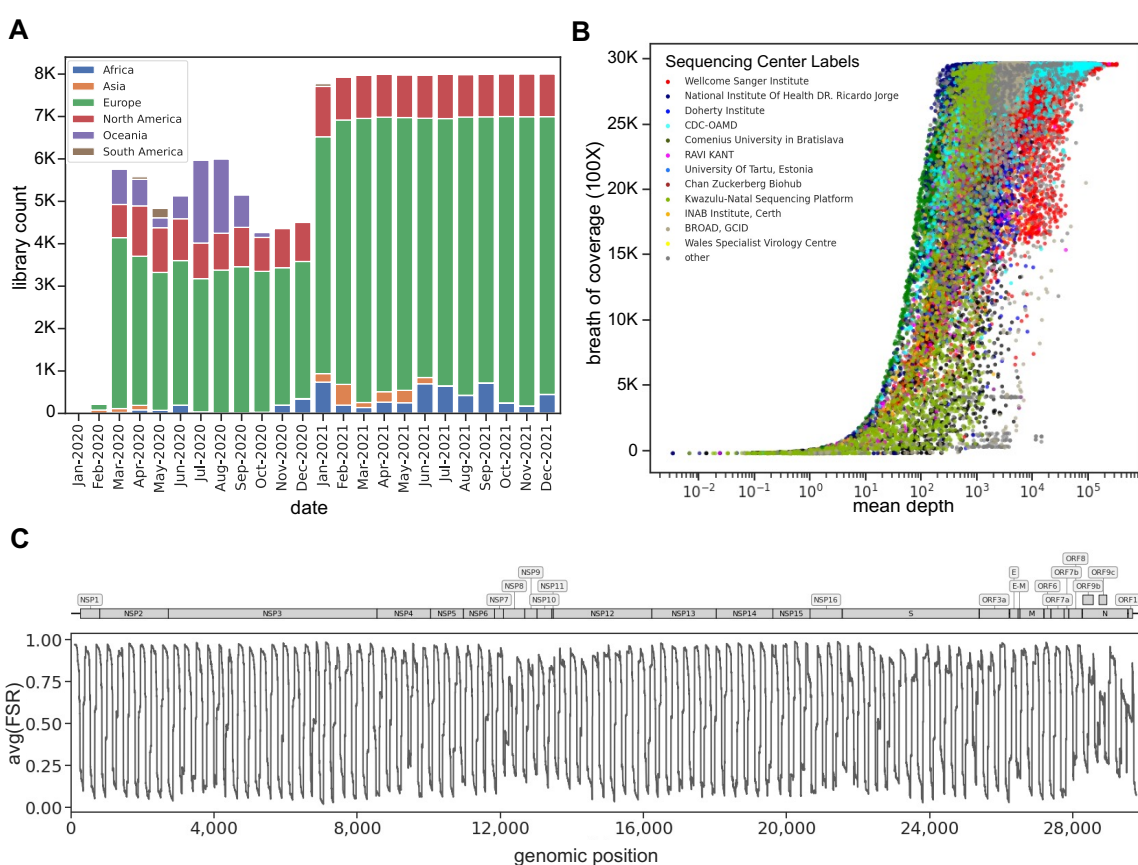


Figure 2: Data description and whole genome quality control. **A**: The 147,537 Illumina paired-end amplicon sequencing libraries were selected and downloaded from NCBI. Each bar in the graph represents the number of samples, categorized by their respective collection dates, with labels indicating their continents. **B**: The x-axis shows each library’s mean depth of coverage (log scale), while the y-axis shows the breadth of coverage. This breadth of coverage is the count of genomic positions covered by at least 100X of depth (D). Libraries with at least 1,000 libraries (93%) in our dataset are explicitly labelled with the sequencing centers, while the remaining libraries have been grouped under the "Others" label. **C**: Whole Genome forward strand ratio (FSR) averaged across libraries for each genomic position. The gene annotations are overlaid on the top panel.

129 We identified a total of 11,635,668 iSNVs in these libraries before any filtering steps, av-
130 eraging 91 iSNVs per library (see Methods section 4.3 and Table S1). PHATE representation
131 of the raw iSNV dataset distinctly discriminates libraries according to WHO lineage annota-
132 tions (Figure 3A). The mean Percentage of Nearest Neighbours from the same WHO lineage
133 (PNN_{WHO}) (Figure 3B) is at 98.39%, corroborating a strong lineage-specific signature in the
134 raw iSNVs (see Methods section 4.4).

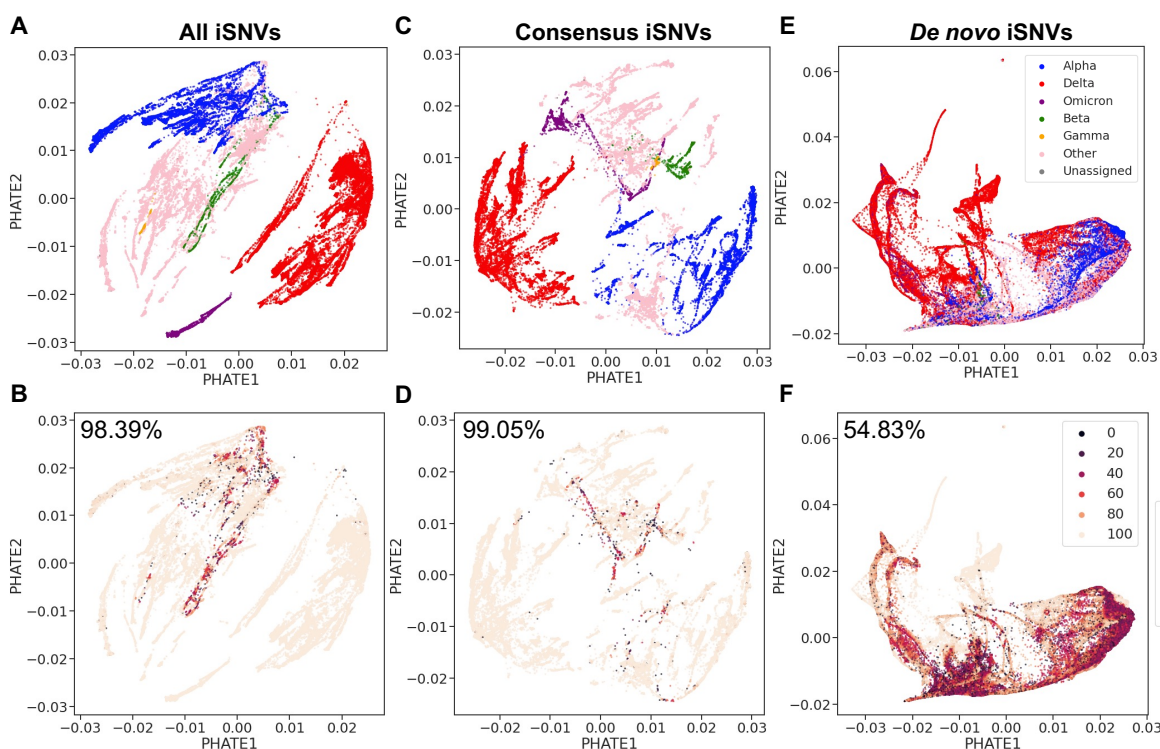


Figure 3: The PHATE representation organizes the unfiltered SARS-CoV-2 libraries according to WHO lineage annotation. **A**: PHATE visualization of the full dataset matrix with 11,635,668 iSNVs, using WHO lineage labels. **B**: The same dataset as **A**, labelled with the percentage of nearest neighbours that share the same WHO annotation as the library itself, and the total PNN_{WHO} value is displayed at the top left. Darker-coloured points signify a lower percentage of neighbouring points sharing the same label as the focal point. **C**: PHATE visualization of the consensus matrix, containing 3,634,563 iSNVs, with WHO lineage labels. **D**: Consensus matrix, similar to **C**, but with PNN_{WHO} labelling and the PNN_{WHO} value at the top left. **E**: PHATE visualization of the *de novo* matrix, including 8,000,668 iSNVs, labelled with WHO lineage annotations. **F**: *de novo* iSNV matrix, as in **E**, but with PNN_{WHO} labelling and the PNN_{WHO} value at the top left. Where a lineage lacks a WHO designation, it is labelled as "Other," and unassigned lineages are labelled as "Unassigned."

135 This result is likely driven by lineage-specific mutations, herein referred to as consensus
136 iSNVs, which are identified as having an Alternative Allele Frequency (AAF) over 75% and
137 are usually part of consensus sequences (Ferreira et al. 2021; Murall et al. 2021; Thielen

138 et al. 2021). These consensus iSNVs account for 3,634,563 iSNVs, averaging 28 per library
139 (Table S1), aligning with the SARS-CoV-2 mutation rate reported by NextStrain (Hadfield
140 et al. 2018) for this time period highlighting the reliability of these consensus iSNVs. PHATE
141 representation of consensus iSNVs only again shows strong alignment with WHO lineages
142 (Figure 3C), which is reflected in the high PNN_{WHO} values in PHATE of 99.05% (Figure
143 3D), confirming these iSNVs largely drive the lineage-specific clustering observed in our raw
144 iSNV dataset.

145 In contrast, we define putative *de novo* iSNVs ($AAF < 0.75$), representing emerging viral
146 mutations within the host, totalling 8,000,668 in the raw dataset, averaging 62 per library
147 (Table S1). The *de novo* iSNVs exhibit more heterogeneous clustering patterns (Fig. 3E)
148 with a lower PNN_{WHO} value of 54.83% (Figure 3F), suggesting a less pronounced lineage-
149 based structure. However, the clustering patterns of the *de novo* iSNVs show a stronger
150 alignment with lineage structure than expected by chance, with the baseline PNN_{WHO} from
151 random resampling at 32.82% for PHATE representation. This highlights the significance of
152 the observed PNN_{WHO} compared to the baseline value, suggesting a lineage-specific biolog-
153 ical relevance in the emerging mutations. The controlled sub-sampling experiments (Figure
154 S1, detailed in Methods section 4.5 and in supplementary information section 10.4) further
155 support these observations, underscoring the distinct clustering behaviours of consensus and
156 *de novo* iSNVs.

157 2.3 Resolving Artifacts in *de novo* iSNVs

158 Due to the geographic distribution of lineages, sequencing centers often sequence certain
159 lineages more frequently than others, potentially leading to technical artifacts that affect
160 lineage clustering in the *de novo* iSNV subset. This is confirmed by the clustering analysis of
161 the 8,000,668 *de novo* iSNVs (Figure 4), where the PHATE representation showed significant
162 sequencing center batch effects (Figure 4A), with a mean percentage of nearest neighbours
163 from the same sequencing center (PNN_{SC}) value of 62.31% (Figure 4B), greatly exceeding
164 the baseline value of 27.53%, expected by chance. This result indicates that our set of *de*
165 *novo* iSNVs likely contains sequencing artifacts. To filter out sequencing artifacts from the
166 set of *de novo* iSNVs and refine the dataset, we used the strand bias metric S (see Methods,
167 equation 1) and the Alternative Allele Frequency (AAF).

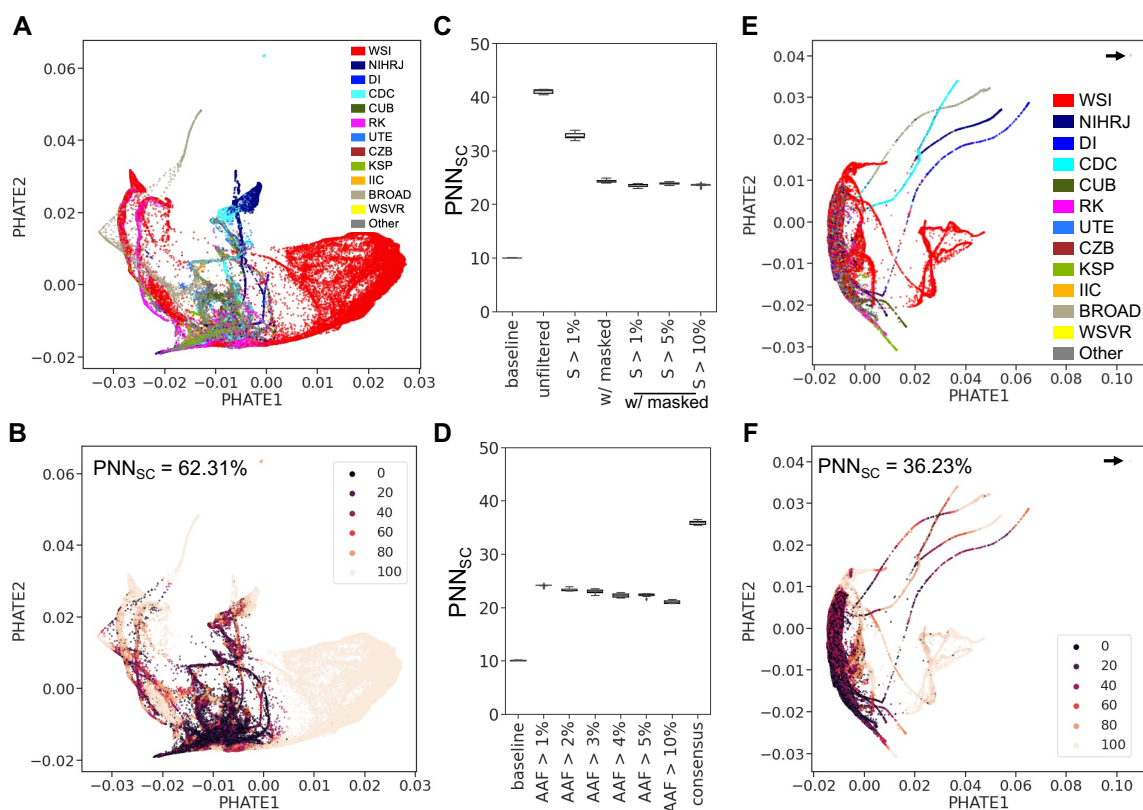


Figure 4: The Use of S and AAF Metrics Improves SARS-CoV-2 *de novo* iSNVs' PHATE Structure by Mitigating Sequencing Center Batch Effects and Artifacts. **A**: PHATE visualization of the unfiltered *de novo* matrix containing 8,000,668 iSNVs, labelled by the libraries' sequencing centers. **B**: PHATE visualization as in **A**, but with labels showing the percentage of $k=100$ nearest neighbours that share the same Sequencing Center (SC) annotation as the library itself and the total PNN_{SC} value is displayed at the top left. **C** and **D**: Boxplots displaying PNN_{SC} values for each PHATE visualization, derived from a sub-sampling controlled experiment across ten replicates (see method section 4.5). **C** shows PNN_{SC} values across various S metric thresholds, and **D** presents PNN_{SC} values across different AAF metric thresholds. **E**: PHATE visualization of the *de novo* matrix, filtered based on S and AAF thresholds, labelled by sequencing centers. **F**: PHATE visualization as in **E**, but with labels showing PNN_{SC} values, and the total PNN_{SC} value is displayed at the top left. In this representation sequencing centers with at least 1,000 libraries in our dataset are explicitly labelled with its sequencing center as follows: Welcome Sanger Institute (WSI), National Institute of Health DR. Ricardo Jorge (NIHRJ), Doherty Institute (DI), CDC-OAMD (CDC), Comenius University in Bratislava (CUB), Ravi Kant (RK), University of Tartu in Estonia (UTE), Chan Zuckerberg Biohub (CZB), Kwazulu-Natal Sequencing Platform (KSP), INAB Insitute in Certh (IIC), BROAD GCID (BROAD), Wales Specialist Virology Center (WSVR). While the remaining libraries were grouped under the "Other" label.

168 The PHATE visualization of the unfiltered *de novo* iSNVs prominently identifies the Well-
169 come Sanger Institute as a major cluster (Figure 4A) due to its significant representation of
170 75% in our library set. This underscores the potential impact of unbalanced sampling on
171 cluster formation and potentially PNN_{SC} values. To neutralize this imbalance, we designed
172 a controlled sub-sampling experiment, evenly selecting 10,000 libraries from each of the top
173 10 sequencing centers based on library counts (see Method section 4.5), aiming to reduce
174 the impact of sampling bias on the PNN_{SC} values. We thus assessed the impact of filtering
175 based on these two metrics, S and AAF , on the PHATE clustering structure measured with
176 PNN_{SC} using the controlled sub-sampling experiment to mitigate bias from uneven sampling
177 across sequencing centers (Figure 4C, D).

178 To address the observed strand coverage unbalanced in our dataset (Figure 2C), we used
179 the strand bias metric S , which assesses the likelihood of strand bias artifacts using the
180 alternative allele's strand coverage. Initially, filtering out iSNVs with $S < 1\%$ and 486 genomic
181 positions showing recurrent strand bias across libraries (see supplementary information section
182 10.3) significantly lowers sequencing center-specific artifacts. This was reflected in the reduced
183 PNN_{SC} values (Figure 4C) in the controlled sub-sampling experiments. However, PNN_{SC}
184 values remained stable when the S threshold was increased beyond 1%, suggesting no further
185 improvement based on this metric (Figure 4C).

186 Filtering based on allele frequency is a key metric in genomic studies. Some studies use a
187 low threshold, which may result in the inclusion of erroneous intra-host mutations (Y. Wang
188 et al. 2021; Armero et al. 2021; Popa et al. 2020; Tonkin-Hill et al. 2021; Lythgoe et al.
189 2021). In contrast, more stringent criteria could overlook the analysis of low-frequency, *de*
190 *nov*o intra-host mutations. Further refinement of our iSNV set based on the AAF metric
191 led to an additional decrease in PNN_{SC} values (Figure 4D), particularly when increasing
192 the AAF threshold to 5%. Despite testing additional combinations of thresholds, the final
193 PNN_{SC} metric did not reach the baseline value of 10%, suggesting that the optimal threshold
194 on the AAF metric is 5%.

195 Applying these optimal thresholds of 1% for S and 5% for AAF to filter out iSNVs, the
196 *de novo* iSNV count dropped from 8,000,668 to 468,651, averaging six iSNVs per library
197 (Table S1). This process notably decreased sequencing center batch effects in the PHATE
198 representation (Figure 4E-F), resulting in a PNN_{SC} value of 36.23%. While this value still

199 exceeds the baseline of 27.69%, the reduction marks an improvement in minimizing batch
200 effects. Additionally, the PNN_{SC} value for lineage-defining consensus iSNVs also does not
201 reach the baseline value (Figure 4D), implying that completely separating sequencing center
202 influences from lineage-specific signatures might represent an intractable challenge.

203 2.4 Identifying Outliers and Center-Specific Patterns

204 In our analysis of the 468,651 filtered *de novo* iSNVs, there remain outlier clusters showing
205 sequencing center homogeneity in the PHATE representation (Figure 4E-F). Notably, a small
206 but distinct set of libraries forms an outlier cluster, markedly separated from other libraries in
207 the PHATE representation (Figure 4E-F, indicated by an arrow). This observation suggests
208 that specific libraries from the same sequencing centers potentially have an excess of shared
209 iSNVs. We thus analyzed libraries' intra-host mutational load, defined as the number of
210 iSNVs in a library (see Method section 4.6). While most libraries in our dataset contain only
211 one or two iSNVs (Figure S2), some exhibit a high intra-host mutational load, with tens of
212 iSNVs per library.

213 To determine the optimal threshold for excess iSNVs in libraries, we computed the PNN_{SC}
214 value in PHATE representation after sequentially removing the top 1%, 5%, and 25% of
215 the most mutated libraries (Figure S2A-B). Removing the top 1% of outliers impacted the
216 PPN_{SC} value the most, decreasing it by 2%. Additional exclusions, even down to only
217 keeping libraries with one iSNV, did not further reduce the PNN_{SC} value (Figure S2B, 50th
218 percentile), underlining the impact of extreme outliers in the PHATE representation of the
219 full dataset.

220 To ensure biologically relevant libraries are not excluded, we explored in-depth the patterns
221 observed in the top 1% outlier libraries (1,270 outlier libraries) by computing the PHATE
222 representation only on these libraries (Figure 5A, S3A). They strongly cluster by sequencing
223 centers, indicating that their iSNVs are enriched for sequencing center-specific artifacts. In
224 this PHATE representation of the outlier libraries, we note four main clusters (Figure S3).
225 Cluster 1 is composed of 159 Doherty Institute libraries, corresponding to Australia's first
226 pandemic wave (March-August 2020) (Figure S3B). Cluster 2 comprises 104 libraries from
227 Scilifelab Stockholm, collected at the end of the second pandemic wave. Cluster 3 includes
228 109 libraries from the Kwazulu-Natal Sequencing Platform, with collections from January to

229 April 2021. Cluster 4 comprises 75 libraries from the Ravi Kant sequencing center, with a
 230 collection peak in May 2021.

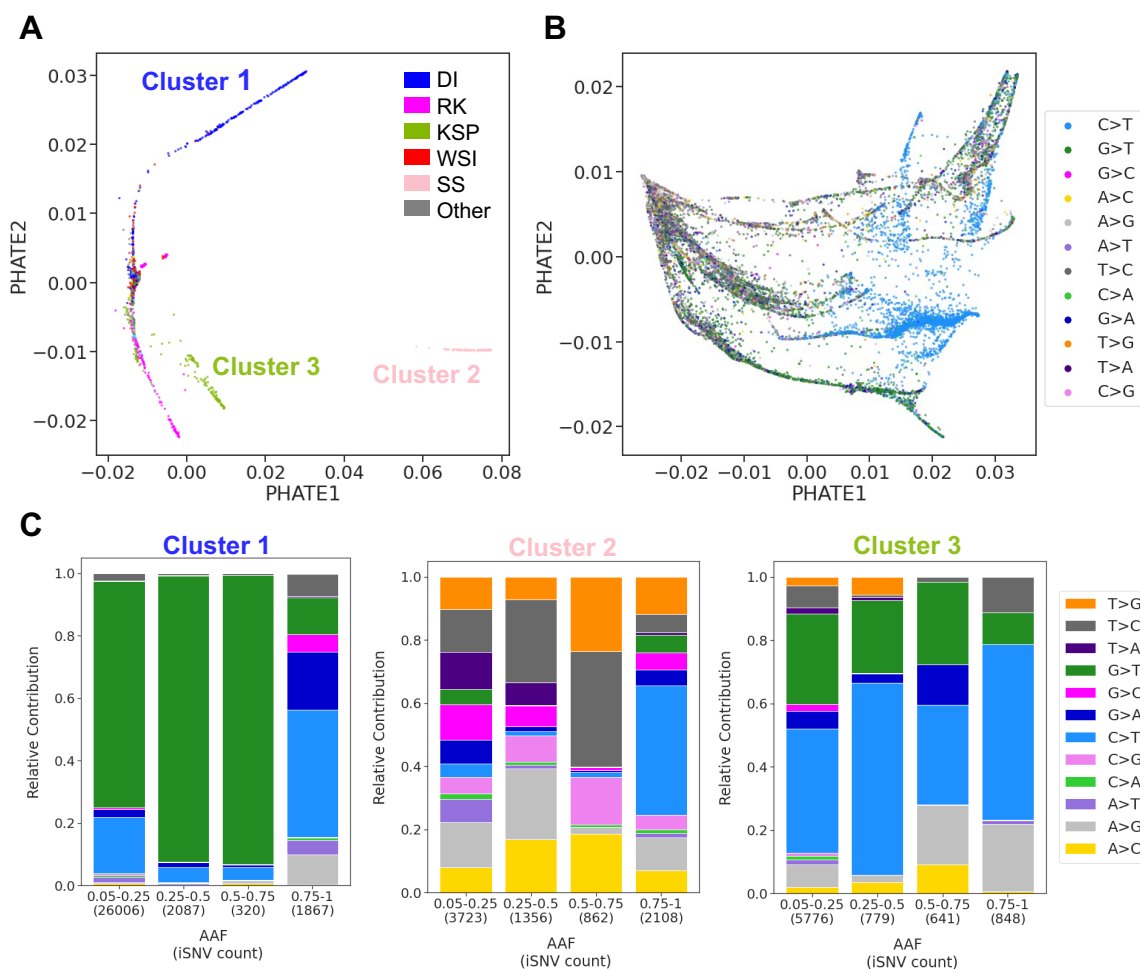


Figure 5: Unique Mutational Patterns in SARS-CoV-2 outlier libraries Tied to Sequencing Centers. **A**: PHATE visualization of outlier libraries, showcasing distinct clusters of SARS-CoV-2 libraries, each associated with specific sequencing centers. **B**: Displays the PHATE representation of genomic positions in outlier libraries, labelled with the most frequent substitution types observed across these libraries. **C**: Mutational patterns in iSNVs across the three distinct clusters in **A**, each associated with a specific sequencing center. **C** sequentially presents mutational patterns in iSNVs from Cluster 1 sequenced by the Doherty Institute, Cluster 2 associated with Scilifelab Stockholm, and Cluster 3 predominantly sequenced by the Kwazulu-Natal Sequencing Platform. The sequencing center's labels are as follows: Welcome Sanger Institute (WSI), Doherty Institute (DI), Ravi Kant (RK), Kwazulu-Natal Sequencing Platform (KSP), and Scilifelab Stockholm (SS).

231 To detect the mutational patterns responsible for these effects, we computed PHATE rep-
 232 resentation by genomic position on these outliers libraries (Figure 5B), showing clustering of
 233 C>T and G>T mutations. This contrasts with non-outlier libraries, which do not show clear
 234 clustering based on substitution patterns (Figure 6A). Quantifying the proportion of iSNVs

235 based on the substitution spectrum (see Methods section 4.7) revealed unique mutational
236 signatures within each of the three main clusters with the most libraries (Figure 5C). Each
237 cluster, associated with a specific sequencing center, exhibited mutational patterns distinct
238 from those in non-outlier libraries (Figure 6A). Cluster 1 displays a prominent G>T pattern
239 in *de novo* iSNVs, not seen in the consensus iSNVs from these same sequences. Interestingly,
240 we identified 40 genomic positions with a *de novo* iSNV in at least 80% of the libraries in
241 cluster 1. Cluster 2 libraries also displayed a unique mutational pattern in their iSNVs (Fig-
242 ure 5C, center), with T>G, T>C, A>G, and A>C as the predominant substitutions. These
243 also diverged from their respective consensus iSNVs except for T>G. A notable 30 genomic
244 positions have a *de novo* iSNV in at least 80% of the libraries in cluster 2. Lastly, cluster 3
245 libraries presented an excess of G>T and C>T that differed from their consensus iSNVs. In
246 this cluster 3, 14 genomic positions have a *de novo* iSNV in at least 80% of the libraries.

247 Overall, our outlier analysis revealed unique mutational patterns in *de novo* iSNVs across
248 different sequencing centers associated with an excess of iSNVs, showing the influence of
249 center-specific sequencing factors. These findings confirm the need to filter out the top 1%
250 outlier libraries with a mutational load above 44 iSNVs in our library set. Our results also
251 highlight the importance for sequencing centers to assess both the abundance of iSNVs and
252 the presence of unique mutational patterns as key indicators for evaluating their sequencing
253 processes.

254 2.5 Deriving a Final *de novo* iSNV Dataset

255 After our extensive curation, we kept 296,437 *de novo* iSNVs with $AAF > 5\%$ and $S > 1\%$,
256 from 72,470 non-outlier libraries with at least one iSNV, as our final curated dataset. The
257 PHATE visualization of the 296,437 retained *de novo* iSNVs by genomic position display
258 no clustering according to the mutational pattern, underscoring the optimal curation of the
259 dataset (Figure 6A). Additionally, the substitution spectrum of this curated set shows a
260 prevalence of C>T and G>T substitutions (Figure 6B), aligning with consensus iSNV patterns
261 and known SARS-CoV-2 mutational trends (Moshiri et al. 2023; Fumagalli et al. 2023; Bloom
262 et al. 2023; Saldivar-Espinoza et al. 2023).

263 The PHATE visualization by library (Figure 6C) shows greater sequencing center ho-
264 mogeneity compared to the initial representation of the raw iSNV data (Figure 4A). WHO

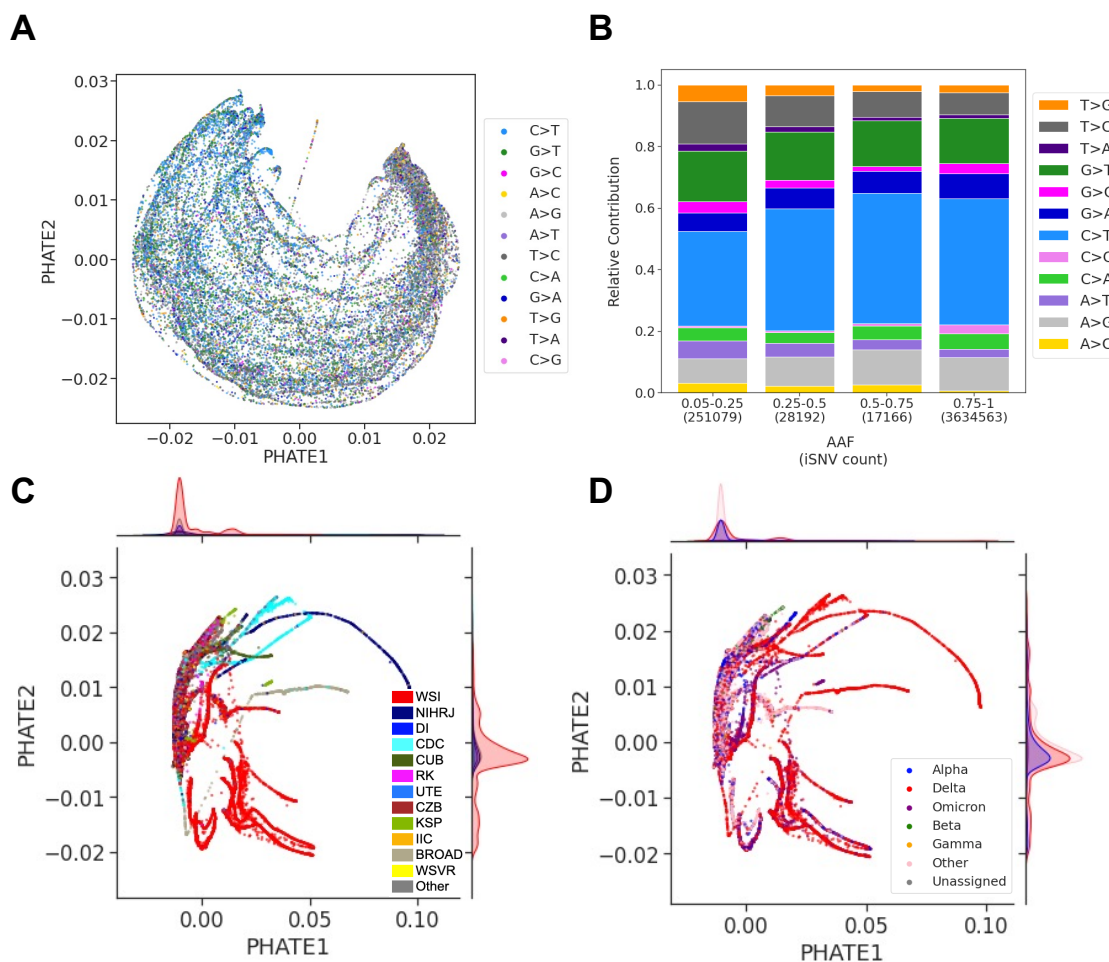


Figure 6: Attaining a Refined and Comprehensive Collection of SARS-CoV-2 Intra-host Sequencing Libraries and iSNVs via Meticulous Filtering. **A**: PHATE visualization of the refined library set, excluding outliers, with *de novo* iSNVs filtered based on *S* and *AAF* metrics. Each library is labelled by its sequencing center. **B**: Similar to **A**, but with labels showing the percentage of nearest neighbours (PNN_{SC}) for each sequencing center and the total PNN_{SC} value displayed at the top left. **C**: The total 296,437 *de novo* and consensus iSNVs, stratified by *AAF* and substitution types to reveal mutational biases. **D**: Presents a PHATE visualization of the transposed matrix for non-outlier libraries with filtered *de novo* iSNVs. Each point represents a genomic position of the SARS-CoV-2 genome, labelled by its most frequent substitution type across the libraries. Sequencing centers with at least 1,000 libraries are explicitly labelled with its sequencing center as follows: Welcome Sanger Institute (WSI), National Institute of Health DR. Ricardo Jorge (NIHRJ), Doherty Institute (DI), CDC-OAMD (CDC), Comenius University in Bratislava (CUB), Ravi Kant (RK), University of Tartu in Estonia (UTE), Chan Zuckerberg Biohub (CZB), Kwazulu-Natal Sequencing Platform (KSP), INAB Insitute in Certh (IIC), BROAD GCID (BROAD), Wales Specialist Virology Center (WSVR).

265 lineage annotations of the same PHATE representation show similar lineage homogeneity
266 (Figure 6D). Both sequencing center and WHO lineage annotations in the PHATE represen-
267 tation concentrate the majority of libraries into a single large cluster, as shown by the density
268 plots. Despite the presence of sequencing center-specific clusters (Figure 6C), lineage-specific
269 clustering is also noticeable (Figure 6D), suggesting that lineages from similar geographic
270 regions may share iSNV generation processes, meriting further investigation. Nevertheless,
271 the optimal refinement of the dataset is supported by a substantial decrease in the PNN_{SC}
272 value, from 62.31% to 33.26% (baseline value 26.29%).

273 3 Discussion

274 Emerging *de novo* mutations, or iSNVs, which occur during the intra-host phase of infection,
275 are critical for understanding viral diversity and evolution. These mutations can be detected
276 by analyzing sequencing libraries from infected hosts, although the sequencing process may
277 introduce artifacts, resulting in false iSNV calls. To address this challenge, we present a com-
278 prehensive two-step workflow tailored for intra-host viral NGS analysis, specifically focusing
279 on the SARS-CoV-2 RNA-seq libraries. It is specifically designed to robustly accommodate
280 and correct for artifacts arising from the diverse sources present in our heterogeneous dataset,
281 ensuring accurate detection of true iSNVs. First, we processed a large dataset of libraries
282 with stringent whole genome quality control. Subsequently, we use these libraries for iSNV
283 calling, employing specific quality metrics to differentiate putative iSNVs from artifacts. We
284 also implemented dimensionality reduction techniques like PHATE and t-SNE to visualize
285 and analyze library structures, enhancing our analysis with an explainability metric. Ap-
286 plying this workflow to a substantial SARS-CoV-2 dataset, we identified a set of emerging
287 (*de novo*) iSNVs for studying intra-host viral evolution, differentiating them from consensus
288 iSNVs using a 75% allele frequency threshold. This threshold is often used for its balance
289 between detecting true positives and minimizing false positives, at the expense of intra-host
290 diversity, by consensus callers (Ferreira et al. 2021; Murall et al. 2021; Thielen et al. 2021).
291 Additionally, we tackled the challenge of distinguishing *de novo* iSNVs from similar-frequency
292 artifacts using tailored quality metrics to establish appropriate thresholds for a given dataset,
293 ensuring our process is rigorous and non-arbitrary.

294 Sequencing accuracy is influenced by multiple factors, including sample preparation, PCR
295 amplification, and sequencing errors (Heguy et al. 2022; Zanini et al. 2017; McCrone et al.
296 2016; Grubaugh et al. 2019). This is especially the case when accurately detecting viral intra-
297 host diversity (McCrone et al. 2016; Illingworth et al. 2017; Zanini et al. 2017). Mutations
298 appearing on only one strand are likely due to amplification errors, as putative mutations
299 would be present on both strands. Known as strand bias artifacts, they have been overlooked
300 in the literature (Dinis et al. 2016; Illingworth et al. 2017; Zanini et al. 2017), but when
301 addressed in recent studies, it is typically through applying a stringent filter that counts
302 the appearances of an alternative allele on each strand (Sun et al. 2023; Xi et al. 2023;
303 N’Guessan et al. 2023). However, this common filtering approach fails to account for the
304 inherent imbalances in strand coverage frequently observed in targeted sequencing of SARS-
305 CoV-2. This oversight can significantly increase the risk of false negatives, with the rate of
306 missed variants varying unevenly across the genome. In response, our strand bias metric takes
307 a different approach by assessing the distribution of each iSNV’s alternative allele across both
308 strands, explicitly accounting for the imbalance in strand coverage observed in our SARS-
309 CoV-2 NGS libraries. This approach avoids the bias of traditional methods that only retain
310 genomic positions covered by both strands, a restriction that could impact about two-thirds
311 of the genome (Figure 2C). Additionally, our strand bias metric, while similar to a published
312 formula (McElroy et al. 2013), is tailored to a large viral NGS dataset. Interestingly, we
313 highlight a set of genomic positions frequently identified as strand bias artifacts supported
314 by our large and comprehensive dataset and see supplementary information section 10.3). By
315 masking these positions, we noted a significant reduction in sequencing center batch effects,
316 indicating that these positions may be specific to sequencing centers. Therefore, we highly
317 recommend masking these positions to mitigate sequencing errors and erroneous data analysis
318 and provide an efficient way to do so (Mostefai et al. 2024).

319 As intra-host viral genomic data grows in size and complexity (Chen et al. 2022; Smith
320 et al. 2023), the challenge of managing these datasets increases. Dimensionality reduction
321 methods are valuable for distilling this data into a more manageable form (Tapinos et al.
322 2019; Paradis 2022). However, interpreting these methods’ two-dimensional representations
323 can be challenging due to unclear biological significance (Karim et al. 2022). In our workflow,
324 we have incorporated PHATE and t-SNE alongside a metric that computes the percentage of

325 nearest neighbours sharing the same annotation (e.g. sequencing center, WHO variant). This
326 approach enhances the explainability of these techniques by highlighting relationships within
327 specific groups of libraries in the representation, establishing a novel approach to analyzing
328 high-dimensional viral sequence data. This methodology also facilitates the identification
329 of optimal iSNV filtering thresholds, a critical aspect of sequencing data quality control.
330 Implementing this approach allowed us to refine our quality metrics, resolve sequencing center
331 batch effects, and improve the reliability of our iSNV dataset. Moreover, we have pioneered
332 the use of PHATE in viral sequencing data analysis. We show that PHATE is especially
333 effective at handling libraries with varying iSNV counts, unlike t-SNE, which is impacted
334 by such libraries (see Figure S2 and supplementary information section 10.5). PHATE's
335 ability to accurately represent *de novo* mutations also demonstrates its potential for broader
336 applications in areas requiring *de novo* mutation analysis, including the study of cancer clonal
337 mutations (Muyas et al. 2023), evolutionary developmental biology (Short et al. 2018), and
338 metagenomics (Keegan et al. 2016).

339 Despite PHATE's ability to handle libraries with varying iSNV counts, outlier libraries
340 containing a large number of iSNVs significantly skewed the PHATE clustering structure,
341 highlighting a problematic aspect where a small subset disproportionately impacts the overall
342 analysis. The significant influence of these outlier libraries was apparent in the unique C>T
343 and G>T mutational patterns observed in PHATE's genomic position representation within
344 the outlier only libraries (Figure 5B), supporting the need to treat these libraries separately.
345 Additionally, the strong clustering by sequencing center of the top 1% outlier libraries suggests
346 that iSNVs within these libraries include sequencing artifacts specific to each center. This
347 was confirmed by the distinct mutational patterns and the recurrence of genomic positions
348 enriched for iSNVs within each outlier cluster, which, in turn, are associated with different
349 sequencing centers. The unique mutational signatures identified within the outlier clusters
350 also provide insight into the potential mechanisms of error introduction or bias in sequencing
351 workflows. For instance, the G T substitution pattern seen in the Peter Doherty Institute
352 libraries at the beginning of the pandemic (March to August 2020) may signal RNA degra-
353 dation. Following the adoption of improved sample storage protocols, this Institute noted
354 a reduction in the count of observed mutations (Peter Doherty Institute platform members,
355 personal communication). This change in sequencing libraries' quality emphasizes the neces-

356 sity of ongoing collaboration between sequencing centers and data analysts to adapt practices
357 and enhance sequencing data accuracy and reliability in real time.

358 Our approach, comprehensive as it is, faces some limitations. First, despite documented
359 instances of mixed infections (Vatteroni et al. 2022; Rockett et al. 2022) where lineage-defining
360 mutations appear at low frequencies, our current workflow is not designed to effectively cap-
361 ture these variations. In cases of mixed infections, iSNVs characterized as *de novo* under our
362 definition may actually stem from the co-presence of two (or more) different strains within a
363 host, as they would fall below our 75% threshold for emerging mutations. Therefore, partic-
364 ular care should be taken when analyzing datasets where a substantial proportion of samples
365 could be mixed infections. In our analysis, we found little evidence of mixed infections, but it
366 remains a possibility that some libraries—and consequently, iSNVs—could originate from such
367 infections, though they would likely have been excluded during our outlier analysis. Further-
368 more, our workflow is currently not specifically designed to address the complexity associated
369 with calling insertions and deletions (indels), which is an area for future development. In
370 particular, a benchmark of indel detection tools for intra-host data should be conducted to
371 enhance this aspect of viral genomic analysis. Bridging this gap poses a notable challenge
372 and offers a valuable opportunity for methodological innovation. This workflow is designed
373 specifically for Illumina sequencing data, favoured for its lower error rate (Fox et al. 2014), and
374 is less suitable for nanopore sequencing (Cook et al. 2024; Fournelle et al. 2024). The latter’s
375 high error rate of about 10% complicates the detection of low-frequency emerging mutations.
376 Our dataset, while diverse, primarily consists of libraries from European and North American
377 sources, mirroring the availability of publicly accessible sequencing data (Chen et al. 2022).
378 This situation underscores the need for improved sequence sharing and support for sequenc-
379 ing capabilities in underserved regions. Additionally, our reliance on publicly available single
380 instances of sequencing libraries leverages accessible data but complicates the confirmation of
381 variant calls due to the absence of multiple replicates, as done previously. We address this by
382 setting a minimum allele frequency threshold of 5%, higher than the typical Illumina error
383 rate of 1% (Fox et al. 2014), aligning with the literature that advocates stricter thresholds for
384 variant identification in the absence of replicates (Roder et al. 2023; Grubaugh et al. 2019).

385 Nonetheless, our workflow and dataset of high-quality intra-host iSNVs have proven in-
386 strumental in testing biological hypotheses and drawing conclusions on diverse areas of study.

387 Published applications include uncovering immune evasion mechanisms in SARS-CoV-2 through
388 sequence analysis and epitope mapping (N’Guessan et al. 2023), comparing intra-host vi-
389 ral evolution between immunosuppressed patients and the general population (Fournelle et
390 al. 2024), and investigating intra-host mutations that influence epitope binding predictions
391 (Caron et al. 2024). Additionally, this workflow and the identified set of *de novo* muta-
392 tions open up new avenues for exploring hypotheses concerning viral intra-host diversity and
393 evolution, providing a foundation for broader research initiatives in this field. For example,
394 we observed that intra-host library clustering based on WHO variants persisted above base-
395 line levels even after removing lineage-defining mutations. This leads us to hypothesize that
396 lineage-defining genetic factors may contribute to the intra-host mutational patterns, suggest-
397 ing a complex underlying mechanism of viral evolution within hosts. Our methodology has
398 proven robust in detecting these subtle lineage characteristics despite variations in sample
399 distribution, reinforcing the possibility of variant-specific effects on mutational events, a find-
400 ing supported by a recent study (Bradley et al. 2024). This intriguing result warrants further
401 investigation that could lead to the discovery of lineage dynamics and mutation impacts.

402 In conclusion, our robust viral intra-host processing and analysis workflow enhances the
403 use of existing cross-sectional sequencing libraries and improves the accuracy and depth of
404 viral genomic analyses. This advanced bioinformatics methodology is crucial for deepening
405 our understanding of intra-host diversity and strengthening preparedness strategies for fu-
406 ture pandemics, proving essential for responding effectively to other viruses in forthcoming
407 outbreaks.

408 4 Methods

409 4.1 Data Selection and Library Pre-processing

410 We downloaded a set of SARS-CoV-2 Illumina amplicon paired-end sequencing libraries
411 dataset from the first two years of the COVID-19 pandemic, ensuring a representative sam-
412 pling across time and geographic locations. For each month from January 2020 to December
413 2021, sequencing libraries were randomly chosen based on availability in the National Center
414 for Biotechnology Information (NCBI): up to 5,000 from the UK, up to 1,000 from the USA,
415 and up to 2,000 from other global regions, totalling a potential 8,000 libraries monthly (Figure

416 2A). This yielded a total of 147,537 downloaded libraries (supplementary information section
417 10.1).

418 For each library, Illumina sequencing adapters and bad-quality reads (Phred score <
419 20) were trimmed from the sequencing reads using TrimGalore V.0.6.0 (<https://github.com/FelixKrueger/TrimGalore>). The trimmed libraries were mapped to the SARS-CoV-
420 2 reference genome (NC045512.2) using BWA mem v.0.7.17-r1188 (Li et al. 2009), gen-
421 erating BAM files. Next, we used the iVar pipeline for primer trimming (Grubaugh et
422 al. 2019), using the ARTIC Network V3, V4, and V4.1 amplicon designs, as the sequenc-
423 ing centers in our dataset predominantly use these three kits during the sampling period
424 (<https://github.com/artic-network/primer-schemes>). We used the samtools mpileup
425 (with specific parameters `-Q 20 -q 0 -B -A -d 600000`) (Danecek et al. 2021) to gener-
426 ate pileup files containing read information for each BAM file. To parse the pileup files
427 and extract relevant data, we employed the publicly available script pileup2base (<https://github.com/riverlee/pileup2base>). We calculated the depth of coverage for each ge-
428 nomic position, which is the number of reads aligning to the position. The mean coverage
429 across all libraries is 10446X, so we labelled any position with depth below 100X (1% of the
430 mean) as low-quality. We calculated two metrics to evaluate each library's quality (Figure
431 2B): (1) C , the mean coverage (the mean number of reads per position) and (2) B , the
432 breadth of coverage (the number of genomic positions with a depth above 100X). We kept the
433 libraries with $C > 100X$ and $B > 20,000$ positions (representing two-thirds of the SARS-CoV-
434 2 genome), yielding a total of 128,423 high-quality libraries. (see supplementary information
435 section 10.1 for more details)

438 4.2 Consensus Sequences and Lineage Annotations

439 We obtained a consensus sequence for each of the 128,423 high-quality libraries using the
440 iVar pipeline consensus calling tool (`-q 20 -t 0.75 -m 20`) (Grubaugh et al. 2019). These
441 consensus sequences were annotated with Pango lineages using Pangolin 4.3 (Rambaut et al.
442 2020), which were next used to annotate with World Health Organization (WHO) lineages
443 (Alpha, Delta, Omicron, Delta, Gamma, and Others) using a custom script. Sequences with
444 no Pango lineage were annotated as 'Unassigned'.

445 4.3 iSNV Calling and Encoding

446 We called iSNVs present in the 128,423 high-quality sequencing libraries (Figure 1B) after
447 extracting genomic positions between positions 101 and 29,778, to exclude positions located
448 at both ends of the genome that are generally of lower quality. For each library, we used
449 pileup2base (Danecek et al. 2021) to obtain a base file, which contains, for the 29,678 positions,
450 the counts for each nucleotide (A, T, C, G) separated according to amplicon direction (forward
451 or reverse strand). Because we are focusing our analyses on single nucleotide substitutions,
452 we ignored the last two columns of the base file that report the number of reads with indels.
453 During this step, we kept only positions with a minimum coverage read depth of 100X.

454 We computed different iSNV metrics at the position level for each library using custom
455 scripts. We define the alternative allele (AA) as the most frequent allele at a given position
456 other than the reference allele. For each position and each library, we computed the Alter-
457 native Allele Frequency (AAF) as $AAF = (D_{AA})/D$, where D is the depth at the position
458 studied and D_{AA} is the depth for the alternative allele.

459 Due to the nature of targeted sequencing with amplicon design (Guo et al. 2012), it is
460 possible that a single position in the genome may not be sequenced in a balanced manner
461 between the forward and reverse directions. We thus compute the forward strand ratio as
462 $FSR = D^f/D$, with D^f the forward strand depth and D the total depth.

463 To evaluate if an alternative allele exhibits unbiased sampling across both strand direc-
464 tions, we used a binomial test. This test determines the probability of observing an allele
465 predominantly on one strand, indicating a higher artifact likelihood. For the forward strand,
466 let Y represent the expected count of reads bearing the alternative allele within the total
467 forward strand reads, D^f . Assuming Y follows a binomial distribution with a probability of
468 success given by the AAF , the probability of observing at least D_{AA}^f forward strand reads
469 with the alternative allele (AA) is calculated as:

$$S^f(Y \leq D_{AA}^f) = \sum_{y=0}^{D_{AA}^f} \binom{D^f}{y} (AAF)^y (1 - AAF)^{D^f - y} \quad (1)$$

470 This same approach is applied to the reverse strand reads, D^r , to calculate the likelihood
471 of observing at least D_{AA}^r reverse strand reads with the alternative allele. Finally, to ensure
472 a stringent evaluation, we take the minimum value of these calculated probabilities for both

473 the forward and reverse strands. This minimum value serves as the Strand Bias Likelihood
474 (S) metric for each iSNV, effectively quantifying the likelihood of no strand bias, and thus, a
475 low value reflects the potential for the presence of an artifact.

476 The resulting iSNVs for each sequencing library are represented by their AAF given
477 position p in a library l ($x_{p,l}$) (Figure 1B). This forms a matrix X , where the rows are our
478 128,423 high-quality sequencing libraries, and the columns are the genomic positions between
479 101 and 27778. We encode the initial unfiltered matrix X with $x_{p,l} = AAF_{p,l}$ for all iSNVs
480 from a given library l , and when an iSNV is filtered out based on thresholds for AAF and S ,
481 $x_{p,l}$ is set to 0.

482 4.4 Dimensionality Reduction and Clustering Evaluation

483 Given the high dimensionality of matrix X , we used dimensionality reduction methods to
484 explore the underlying structure within the high-quality libraries in two dimensions (2D). We
485 used incremental principal component analyses (PCA) for initialization and then obtained
486 2D representations of the PCA-transformed data using two different approaches: the widely-
487 used t-distributed Stochastic Neighbor Embedding (t-SNE) (Tamazian et al. 2022; Maaten
488 et al. 2008; Pedregosa et al. 2011) and the more recent heat diffusion for affinity-based tran-
489 sition embedding (PHATE) (Moon et al. 2019). We applied t-SNE with the Python library
490 `sklearn.manifold.T-SNE` and PHATE with the PHATE Python library (Moon et al. 2019).
491 The 2D embedding outputs from PHATE and t-SNE are visualized in scatter plots where
492 each library is coloured either by sequencing center or WHO lineage annotation.

493 To measure the impact of specific subgroups of iSNVs on clustering structures based on
494 either sequencing center (SC) or WHO lineage labels, we used a k-nearest neighbour (kNN)
495 approach, using $k=100$. This value of k is selected to simplify interpretation as a percentage
496 during neighbour selection and reflects the large number of libraries in our dataset. For
497 each library l within a representation, we identify the 100 nearest libraries $NN(l)$ using the
498 `sklearn.neighbors` Python package (Pedregosa et al. 2011). We then calculate $NN_z(l)$, the
499 count of nearest neighbours sharing the same z label as library l (where z is either WHO for
500 lineage or SC for sequencing center), and compute the percentage of nearest neighbours with
501 matching labels. For each representation, we derive a final PNN_z as the mean percentage
502 of nearest neighbours with matching labels across all libraries. A higher PNN_z indicates

503 that label z describes the data’s clustering structure. We also generate a baseline PNN_z ,
504 representing expected chance levels by randomly shuffling labels z before calculation. This
505 baseline acts as a standard for assessing the significance of observed patterns, emphasizing
506 the delta between observed and baseline PNN_z over the choice of k value.

507 **4.5 Experimental Design to Mitigate Sampling Biases**

508 We designed a controlled sub-sampling experiment by randomly sub-sampling libraries for
509 each WHO or sequencing center annotation to address the impact of biases stemming from
510 unbalanced sampling. To evaluate the influence of iSNVs on WHO patterns, we iteratively
511 sampled 1000 library rows for each of the Alpha, Beta, Delta, and Omicron variants from the
512 data matrix X ten times, generating replicates. This process resulted in ten matrices, each
513 comprising 4,000 rows. To investigate the effect of iSNVs on sequencing center patterns, we
514 used a similar approach, randomly selecting 1000 library rows. However, in this case, we ran-
515 domly sampled 1000 library rows from our dataset’s top 10 most frequent sequencing centers
516 (Table S3). This process resulted in ten matrices as replicates, each comprising 10,000 rows.
517 Within each matrix, $x_{p,l}$ values were set to 0 based on various AAF and S thresholds cut-
518 offs. After these two steps, we applied the same method as in Data Visualization to generate
519 PHATE and t-SNE visualization of the matrices. Subsequently, we quantified the clustering
520 structure of t-SNE and PHATE to derive a PNN_z value for each visualization. Specifically,
521 a high value of PNN_{SC} , indicating clustering primarily by sequencing center, would suggest
522 a dataset enriched for artifacts. Conversely, a high value of PNN_{WHO} , signifying clustering
523 primarily by WHO lineage annotations, would suggest a more biologically relevant dataset.

524 **4.6 Mutational Load**

525 The mutational load for each library was calculated as the total count of distinct iSNVs
526 identified regardless of allele frequencies. Per library, mutational loads were visualized using
527 histograms to illustrate the distribution of mutational loads across the dataset. We categorized
528 the libraries into different percentiles based on their mutational load, identifying those with
529 higher or lower numbers of iSNVs.

530 4.7 Substitution Spectrum Analysis

531 To assess the mutational landscape and identify specific patterns that may indicate underlying
532 mutational mechanisms or biases in the dataset, we looked at the substitution patterns within
533 iSNVs' different *AAF* frequencies. First, we categorized intra-host iSNVs into four *AAF* bins,
534 as follows: 5% to 25%, 25% to 50%, 50% to 75%, and 75% to 100%. This categorization was
535 based on the evidence for an alternative allele present in the iSNVs. Next, within each *AAF*
536 bin, we classified each iSNV in terms of its ancestral allele and alternative allele to obtain
537 12 categories of substitution types. These are A>G, A>C, A>T, C>A, C>G, C>T, G>A,
538 G>C, G>T, T>A, T>C, and T>G. This allowed us to analyze the relative contribution of
539 each substitution type within each *AAF* range.

540 5 Data and Code Access

541 The processing workflow's code can be found here: [https://github.com/HussinLab/IntraHost_](https://github.com/HussinLab/IntraHost_Covid_Pipeline.git)
542 [Covid_Pipeline.git](https://github.com/HussinLab/IntraHost_Covid_Pipeline.git). NCBI accession IDs utilized in this study and the high-quality iSNVs
543 identified within each sequencing library are accessible through a Mendeley data repository
544 (Mostefai et al. 2024, <https://doi.org/10.17632/8nvgtrkzdm.1>). We also provide the list
545 of recommended 477 genomic positions to mask in the same data repository (see supplemen-
546 tary information section 10.3).

547 6 Competing Interest Statement

548 This study was supported by funding from the Canada Foundation for Innovation (CFI)
549 (#40157), the IVADO COVID-19 Rapid Response grant (CVD19-030), the National Sciences
550 and Engineering Research Council (NSERC) (ALLRP-554923-20), the Canadian Institutes
551 of Health Research (CIHR) Project Grant (#174924), and the CIHR operating grant to the
552 Coronavirus Variants Rapid Response Network (CoVaRR-Net, ARR-175622). FM doctoral
553 studies are supported by the Hydro Quebec Scholarship. JGH is a Fonds de recherche du
554 Québec Santé (FRQS) Junior 2 Research Scholar.

555 **7 Acknowledgments**

556 The study was conducted in accordance with the Declaration of Helsinki and approved by
557 the Ethics Committee of the Montreal Heart Institute (protocol code 2021-2868 and date of
558 approval 23 July 2021). We express our gratitude to the members of Julie Hussin’s group for
559 their valuable discussions. The successful completion of this work was made possible by the
560 computational resources provided by the Digital Research Alliance of Canada clusters Narval
561 and Beluga. We also thank sequencing library submitters who made data available on NCBI,
562 who are listed in the available iSNV table (Mostefai et al. 2024). We especially extend our
563 appreciation to the Peter Doherty Institute for their insightful communication regarding the
564 distinctive mutational patterns observed within our dataset.

565 **8 Authors’ Contributions**

566 FM: conception, data acquisition, data pre-processing, data analysis, figure development,
567 and manuscript drafting. RP and JCG: data acquisition, data pre-processing, data analysis,
568 and manuscript critical revision and editing. JGH: funding, conception, supervision, and
569 co-drafting of the manuscript.

570 References

- 571 Armero A., Berthet N., and Avarre J.-C. 2021. Intra-Host Diversity of SARS-Cov-2 Should
572 Not Be Neglected: Case of the State of Victoria, Australia. *Viruses* **13**.
- 573 Bashor L., Gagne R. B., Bosco-Lauth A. M., Bowen R. A., Stenglein M., and VandeWoude S.
574 2021. SARS-CoV-2 evolution in animals suggests mechanisms for rapid variant selection.
575 *Proc. Natl. Acad. Sci. U. S. A.* **118**.
- 576 Bloom J. D., Beichman A. C., Neher R. A., and Harris K. 2023. Evolution of the SARS-CoV-2
577 Mutational Spectrum. *Mol. Biol. Evol.* **40**.
- 578 Bradley C. C. et al. 2024. Targeted accurate RNA consensus sequencing (tARC-seq) reveals
579 mechanisms of replication error affecting SARS-CoV-2 divergence. *Nat Microbiol.*
- 580 Bull R. A. et al. 2011. Sequential bottlenecks drive viral evolution in early acute hepatitis C
581 virus infection. *PLoS Pathog.* **7**: e1002243.
- 582 Caron E. et al. 2024. Integrating machine learning-enhanced immunopeptidomics and SARS-
583 CoV-2 population-scale analyses unveils novel antigenic features for Next-generation COVID-
584 19 vaccines. *Research Square*. DOI: [10.21203/rs.3.rs-3914861/v1](https://doi.org/10.21203/rs.3.rs-3914861/v1).
- 585 Chen Z. et al. 2022. Global landscape of SARS-CoV-2 genomic surveillance and data sharing.
586 *Nat. Genet.* **54**: 499–507.
- 587 Cook R. et al. 2024. The long and short of it: benchmarking viromics using Illumina, Nanopore
588 and PacBio sequencing technologies. *Microb Genom* **10**.
- 589 Danecek P. et al. 2021. Twelve years of SAMtools and BCFtools. *GigaScience* **10**. DOI: [10.1093/gigascience/giab008](https://doi.org/10.1093/gigascience/giab008).
- 590
- 591 Di Giorgio S., Martignano F., Torcia M. G., Mattiuz G., and Conticello S. G. 2020. Evi-
592 dence for host-dependent RNA editing in the transcriptome of SARS-CoV-2. *Sci Adv* **6**:
593 eabb5813.
- 594 Dinis J. M., Florek K. R., Fatola O. O., Moncla L. H., Mutschler J. P., Charlier O. K.,
595 Meece J. K., Belongia E. A., and Friedrich T. C. 2016. Deep Sequencing Reveals Potential
596 Antigenic Variants at Low Frequencies in Influenza A Virus-Infected Humans. *J. Virol.*
597 **90**: 3355–3365.

- 598 Ferreira V. H. et al. 2021. Prospective observational study and serosurvey of SARS-CoV-2
599 infection in asymptomatic healthcare workers at a Canadian tertiary care center. *PLoS*
600 *One* **16**: e0247258.
- 601 Fournelle D. et al. 2024. Intra-Host Evolution Analyses in an Immunosuppressed Patient
602 Supports SARS-CoV-2 Viral Reservoir Hypothesis. *Viruses* **16**.
- 603 Fox E. J., Reid-Bayliss K. S., Emond M. J., and Loeb L. A. 2014. Accuracy of Next Generation
604 Sequencing Platforms. *Next Gener Seq Appl* **1**.
- 605 Fumagalli S. E., Padhiar N. H., Meyer D., Katneni U., Bar H., DiCuccio M., Komar A. A.,
606 and Kimchi-Sarfaty C. 2023. Analysis of 3.5 million SARS-CoV-2 sequences reveals unique
607 mutational trends with consistent nucleotide and codon frequencies. *Virol. J.* **20**: 31.
- 608 Ghafari M., Liu Q., Dhillon A., Katzourakis A., and Weissman D. B. 2022. Investigating
609 the evolutionary origins of the first three SARS-CoV-2 variants of concern. *Frontiers in*
610 *Virology* **2**.
- 611 Goldberg A. R. et al. 2023. Wildlife exposure to SARS-CoV-2 across a human use gradient.
612 *bioRxiv*. DOI: [10.1101/2022.11.04.515237](https://doi.org/10.1101/2022.11.04.515237).
- 613 Grubaugh N. D. et al. 2019. An amplicon-based sequencing framework for accurately mea-
614 suring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol.* **20**: 8.
- 615 Guo Y., Long J., He J., Li C.-I., Cai Q., Shu X.-O., Zheng W., and Li C. 2012. Exome
616 sequencing generates high quality data in non-target regions. *BMC Genomics* **13**: 194.
- 617 Hadfield J., Megill C., Bell S. M., Huddleston J., Potter B., Callender C., Sagulenko P.,
618 Bedford T., and Neher R. A. 2018. Nextstrain: real-time tracking of pathogen evolution.
619 *Bioinformatics* **34**: 4121–4123. DOI: [10.1093/bioinformatics/bty407](https://doi.org/10.1093/bioinformatics/bty407).
- 620 Hale V. L. et al. 2022. SARS-CoV-2 infection in free-ranging white-tailed deer. *Nature* **602**:
621 481–486.
- 622 Hedskog C., Mild M., Jernberg J., Sherwood E., Bratt G., Leitner T., Lundeberg J., Ander-
623 sson B., and Albert J. 2010. Dynamics of HIV-1 quasispecies during antiviral treatment
624 dissected using ultra-deep pyrosequencing. *PLoS One* **5**: e11345.
- 625 Heguy A. et al. 2022. Amplification Artifact in SARS-CoV-2 Omicron Sequences Carrying
626 P681R Mutation, New York, USA. *Emerg. Infect. Dis.* **28**: 881–883.
- 627 Hill V. et al. 2022. The origins and molecular evolution of SARS-CoV-2 lineage B.1.1.7 in the
628 UK. *Virus Evol* **8**: veac080.

- 629 Hozumi Y., Wang R., Yin C., and Wei G.-W. 2021. UMAP-assisted K-means clustering of
630 large-scale SARS-CoV-2 mutation datasets. *Comput. Biol. Med.* **131**: 104264.
- 631 Illingworth C. J. R., Roy S., Beale M. A., Tutill H., Williams R., and Breuer J. 2017. On the
632 effective depth of viral sequence data. *Virus Evol* **3**: vex030.
- 633 Karim M. R., Islam T., Beyan O., Lange C., Cochez M., Rebholz-Schuhmann D., and Decker
634 S. 2022. Explainable AI for Bioinformatics: Methods, Tools, and Applications.
- 635 Keegan K. P., Glass E. M., and Meyer F. 2016. MG-RAST, a Metagenomics Service for
636 Analysis of Microbial Community Structure and Function. *Methods Mol Biol* **1399**: 207–
637 233.
- 638 Lauring A. S. 2020. Within-Host Viral Diversity: A Window into Viral Evolution. *Annu Rev*
639 *Viro* **7**: 63–81.
- 640 Li H. and Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler
641 transform. *Bioinformatics* **25**: 1754–1760. DOI: [10.1093/bioinformatics/btp324](https://doi.org/10.1093/bioinformatics/btp324).
- 642 Lythgoe K. A. et al. 2021. SARS-CoV-2 within-host diversity and transmission. *Science* **372**.
- 643 Maaten L. van der and Hinton G. 2008. Visualizing Data using t-SNE. *Journal of Machine*
644 *Learning Research* **9**: 2579–2605.
- 645 Markov P. V., Ghafari M., Beer M., Lythgoe K., Simmonds P., Stilianakis N. I., and Kat-
646 zourakis A. 2023. The evolution of SARS-CoV-2. *Nat. Rev. Microbiol.* **21**: 361–379.
- 647 McCrone J. T. and Lauring A. S. 2016. Measurements of Intrahost Viral Diversity Are Ex-
648 tremely Sensitive to Systematic Errors in Variant Calling. *J. Virol.* **90**: 6884–6895.
- 649 McElroy K., Zagordi O., Bull R., Luciani F., and Beerenwinkel N. 2013. Accurate single
650 nucleotide variant detection in viral populations by combining probabilistic clustering
651 with a statistical test of strand bias. *BMC Genomics* **14**: 501.
- 652 Messali S., Rondina A., Giovanetti M., Bonfanti C., Ciccozzi M., Caruso A., and Caccuri
653 F. 2023. Traceability of SARS-CoV-2 transmission through quasispecies analysis. *J. Med.*
654 *Viro* **95**: e28848.
- 655 Moon K. R. et al. 2019. Visualizing structure and transitions in high-dimensional biological
656 data. *Nat. Biotechnol.* **37**: 1482–1492.
- 657 Moshiri K., Mahmanzar M., Mahdavi B., Tokhanbigli S., Rahimian K., and Tavakolpour S.
658 2023. Mutation accumulation of SARS-CoV-2 genome in North America, South America,

- 659 and Oceania: Analysis of over 6.5 million sequences samples from Global Initiative on
660 Sharing Avian Influenza Data.
- 661 Mostefai F., Grenier J.-C., Poujol R., and Hussin J. 2024. SARS-CoV-2 Intra-host Mutational
662 Landscape: A Curated Dataset of iSNVs. *Mendeley Data*. DOI: [doi.org/10.17632/
663 8nvgtrkzdm.1](https://doi.org/10.17632/8nvgtrkzdm.1).
- 664 Mostefai F. et al. 2022. Population Genomics Approaches for Genetic Characterization of
665 SARS-CoV-2 Lineages. *Front. Med.*: 826746.
- 666 Murall C. L. et al. 2021. A small number of early introductions seeded widespread transmission
667 of SARS-CoV-2 in Québec, Canada. *Genome Med.* **13**: 169.
- 668 Muyas F., Sauer C. M., Valle-Inclán J. E., Li R., Rahbari R., Mitchell T. J., Hormoz S.,
669 and Cortés-Ciriano I. 2023. De novo detection of somatic mutations in high-throughput
670 single-cell profiling data sets. *Nat. Biotechnol.*
- 671 N’Guessan A. et al. 2023. Selection for immune evasion in SARS-CoV-2 revealed by high-
672 resolution epitope mapping and sequence analysis. *iScience* **26**: 107394.
- 673 Novembre J. et al. 2008. Genes mirror geography within Europe. *Nature* **456**: 98–101.
- 674 Oreshkova N. et al. 2020. SARS-CoV-2 infection in farmed minks, the Netherlands, April and
675 May 2020. *Euro Surveill.* **25**.
- 676 Oude Munnink B. B. et al. 2021. Transmission of SARS-CoV-2 on mink farms between humans
677 and mink and back to humans. *Science* **371**: 172–177.
- 678 Paradis E. 2022. Reduced multidimensional scaling. *Comput. Stat.* **37**: 91–105.
- 679 Pedregosa F. et al. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learn-
680 ing Research* **12**: 2825–2830.
- 681 Platzer A. 2013. Visualization of SNPs with t-SNE. *PLoS One* **8**: e56883.
- 682 Popa A. et al. 2020. Genomic epidemiology of superspreading events in Austria reveals mu-
683 tational dynamics and transmission properties of SARS-CoV-2. *Sci. Transl. Med.* **12**.
- 684 Quaranta E. G. et al. 2022. SARS-CoV-2 intra-host evolution during prolonged infection in
685 an immunocompromised patient. *Int. J. Infect. Dis.* **122**: 444–448.
- 686 Rajendran M. and Babbitt G. A. 2022. Persistent cross-species SARS-CoV-2 variant infectiv-
687 ity predicted via comparative molecular dynamics simulation. *R Soc Open Sci* **9**: 220600.
- 688 Rambaut A., Holmes E. C., O’Toole Á., Hill V., McCrone J. T., Ruis C., Plessis L. du, and
689 Pybus O. G. 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist

- 690 genomic epidemiology. *Nature Microbiology* **5**: 1403–1407. DOI: [10.1038/s41564-020-](https://doi.org/10.1038/s41564-020-0770-5)
691 [0770-5](https://doi.org/10.1038/s41564-020-0770-5).
- 692 Robinson S. J. et al. 2023. Surveillance for SARS-CoV-2 in Norway Rats (*Rattus norvegicus*)
693 from Southern Ontario. *Transbound. Emerg. Dis.* **2023**.
- 694 Rockett R. J. et al. 2022. Co-infection with SARS-CoV-2 Omicron and Delta variants revealed
695 by genomic surveillance. *Nat. Commun.* **13**: 2745.
- 696 Roder A. E. et al. 2023. Optimized quantification of intra-host viral diversity in SARS-CoV-2
697 and influenza virus sequence data. *MBio* **14**: e0104623.
- 698 Sacchetto L. et al. 2021. Lack of Evidence of Severe Acute Respiratory Syndrome Coron-
699 avirus 2 (SARS-CoV-2) Spillover in Free-Living Neotropical Non-Human Primates, Brazil.
700 *Viruses* **13**.
- 701 Saldivar-Espinoza B., Garcia-Segura P., Novau-Ferré N., Macip G., Martínez R., Puigbò P.,
702 Cereto-Massagué A., Pujadas G., and Garcia-Vallve S. 2023. The Mutational Landscape
703 of SARS-CoV-2. *Int. J. Mol. Sci.* **24**.
- 704 Short P. J. et al. 2018. De novo mutations in regulatory elements in neurodevelopmental
705 disorders. *Nature* **555**: 611–616.
- 706 Smith E. A. et al. 2023. Pathogen genomics in public health laboratories: successes, challenges,
707 and lessons learned from California’s SARS-CoV-2 Whole-Genome Sequencing Initiative,
708 California COVIDNet. *Microb Genom* **9**.
- 709 Sonnleitner S. T. et al. 2022. Cumulative SARS-CoV-2 mutations and corresponding changes
710 in immunity in an immunocompromised patient indicate viral evolution within the host.
711 *Nat. Commun.* **13**: 2560.
- 712 Sun B., Ni M., Liu H., and Liu D. 2023. Viral intra-host evolutionary dynamics revealed via
713 serial passage of Japanese encephalitis virus in vitro. *Virus Evol* **9**: veac103.
- 714 Tamazian G., Komissarov A. B., Kobak D., Polyakov D., Andronov E., Nechaev S., Kryzhe-
715 vich S., Porozov Y., and Stepanov E. 2022. t-SNE Highlights Phylogenetic and Temporal
716 Patterns of SARS-CoV-2 Spike and Nucleocapsid Protein Evolution: 255–262.
- 717 Tapinos A., Constantinides B., Phan M. V. T., Kouchaki S., Cotten M., and Robertson
718 D. L. 2019. The Utility of Data Transformation for Alignment, De Novo Assembly and
719 Classification of Short Read Virus Sequences. *Viruses* **11**.

- 720 Thielen P. M. et al. 2021. Genomic diversity of SARS-CoV-2 during early introduction into
721 the Baltimore-Washington metropolitan area. *JCI Insight* **6**.
- 722 Tonkin-Hill G. et al. 2021. Patterns of within-host genetic diversity in SARS-CoV-2. *Elife* **10**.
- 723 Vatteroni M. L., Capria A.-L., Spezia P. G., Frateschi S., and Pistello M. 2022. Co-infection
724 with SARS-CoV-2 omicron BA.1 and BA.2 subvariants in a non-vaccinated woman. *Lancet*
725 *Microbe* **3**: e478.
- 726 Wang B. and Jiang L. 2021. Principal Component Analysis Applications in COVID-19 Genome
727 Sequence Studies. *Cognit. Comput.*: 1–12.
- 728 Wang Y. et al. 2021. Intra-host variation and evolutionary dynamics of SARS-CoV-2 popu-
729 lations in COVID-19 patients. *Genome Med.* **13**: 30.
- 730 Washburne A., Jones A., Zhang D., Deigin Y., Quay S., and Massey S. E. 2022. Statistical
731 challenges for inferring multiple SARS-CoV-2 spillovers with early outbreak phylodynam-
732 ics. *bioRxiv*. DOI: [10.1101/2022.10.10.511625](https://doi.org/10.1101/2022.10.10.511625).
- 733 Wertheim J. O. et al. 2022. Detection of SARS-CoV-2 intra-host recombination during super-
734 infection with Alpha and Epsilon variants in New York City. *Nat. Commun.* **13**: 3645.
- 735 Xi B., Zeng X., Chen Z., Zeng J., Huang L., and Du H. 2023. SARS-CoV-2 within-host
736 diversity of human hosts and its implications for viral immune evasion. *MBio* **14**: e0067923.
- 737 Zanini F., Brodin J., Albert J., and Neher R. A. 2017. Error rates, PCR recombination, and
738 sampling depth in HIV-1 whole genome deep sequencing. *Virus Res.* **239**: 106–114.
- 739 Zhang Y. et al. 2022. SARS-CoV-2 intra-host single-nucleotide variants associated with disease
740 severity. *Virus Evol* **8**: veac106.

741 9 Supplemental Tables and Figures

Table S1: iSNV Count, Library Count, and Mutational Load

	iSNV Count	Library Count	Mutational Load
Total iSNVs	11,635,231	128,443	91
Consensus iSNVs	3,634,563	128,323	28
<i>De novo</i> iSNVs	8,000,668	128,352	62
S > 1% filtered <i>de novo</i> iSNVs	6,508,783	127,941	51
Masked <i>de novo</i> iSNVs	5,805,486	125,382	46
AAF > 5% filtered <i>de novo</i> iSNVs	468,651	73,729	6
Non-Outlier libraries <i>de novo</i> iSNVs	296,437	72,470	4

Table S2: Per Country Sequencing Libraries' Counts Before and After Coverage Filters.

	Countries	Before Filters (C_l)	After Filters (B_l)
Africa	Angola	519	266
	Cameroon	210	102
	Ethiopia	125	47
	Malawi	428	175
	Mozambique	287	152
	South Africa	3,662	2,527
	Zimbabwe	507	342
	Other ($n < 100$)	181	117
Asia	China	115	106
	India	437	350
	Israel	609	538
	Lebanon	367	275
	Pakistan	227	149
	Other ($n < 100$)	100	92
Europe	Austria	543	542
	Estonia	5,817	4,704
	Finland	5,331	4,394
	Greece	2,614	2,231
	Italy	452	323
	Norway	3,376	1,139
	Portugal	11,700	10,502
	Slovakia	5,982	5,553
	Switzerland	379	375
	United Kingdom	74,557	69,710
Other ($n < 100$)	129	101	
North America	Canada	625	588
	USA	20,880	16,300
	Other ($n < 100$)	91	82
Oceania	Australia	6,837	6,421
	Northern Mariana Islands	23	22
South America	Brazil	417	186
	Other ($n < 100$)	10	9
Total		147,537	128,420

Table S3: Per Sequencing Center Librerie's Counts Before and After Coverage Filters.

Sequencing Center	Before Filters (C_l)	After Filters (B_l)
Wellcome Sanger Institute	69,676	65,316
National Institute Of Health DR. Ricardo Jorge	11,700	10,502
Doherty Institute	6,699	6,306
CDC-OAMD	6,431	5,040
Ravi Kant	5,331	4,394
Kwazulu-Natal Sequencing Platform	4,711	2,827
Comenius University in Bratislava	4,648	4,458
University Of Tartu, Estonia	4,104	3,351
Norwegian Institute of Public Health (NIPH)	3,376	1,139
INAB Institute, Certh	2,614	2,231
BROAD, GCID	2,334	1,948
Chan Zuckerberg Biohub	1,971	1,850
Quadram Institute Bioscience	1,859	1,257
UPHL ID	1,853	1,062
Institute of Biomedicine and Translational Medicine	1,713	1,353
Wales Specialist Virology Centre	1,507	1,485
Chan Zuckerberg Biohub	1,364	1,229
Public Health Authority of the Slovak Republic	1,349	1,106
TX-SARS-COV-2	1,324	710
Public Health England (Colindale)	1,142	1,123
DCLS-NGS	1,062	492
California Department of Public Health	982	900
Liverpool Clinical Laboratories	825	741
CanCOGeN CPHLN	612	579
Tel Aviv University	609	538
NYC SARS-COV-2	562	493
CeMM	543	542
New Mexico Department of Health Scientific Laboratory	539	503
Colorado Department of Public Health and Environment	509	463
Gujarat Biotechnology Research Centre	436	349
West of Scotland Specialist Virology Centre, NHSGG	362	362
University Hospital of Basel	336	336
Delaware Public Health Lab	327	299
University of Kwazulu-Natal	277	180
Network for Genomic Surveillance in South Africa	275	274
CDC-PDD	253	222
Utah Public Health lab	233	169
Kwazulu-Natal Research and Sequencing Platform	229	137
UMIGS	223	209
University of Verona	220	188
SEARCH	211	191
Hospital Israelita Albert Einstein	208	0
SciLifeLab Stockholm	164	133
Institute of Clinical Pathology and Medical Research	138	135
LNCC	119	110
Ruijin Hospital, Shanghai Jiao Tong University of Medicine	112	103
Centers for Disease Control and Prevention	100	95

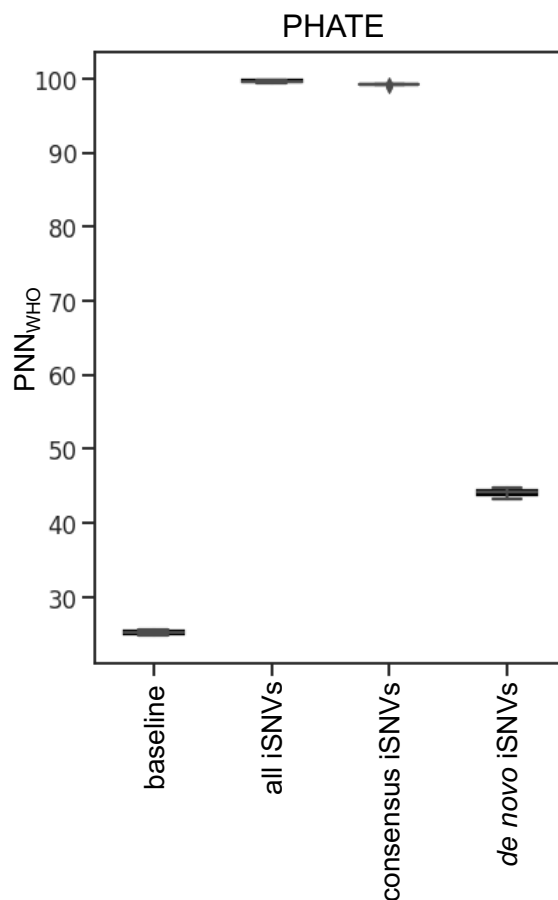


Figure S1: Unveiling WHO Lineage Patterns in SARS-CoV-2 iSNVs with PHATE Visualizations and PNN_{WHO} Metric. Boxplots show the distribution of the mean percentage of nearest neighbours (PNN_{WHO}) from the same WHO lineage annotation across libraries for each PHATE (**A**) visualization across the ten replicates from the sub-sampling controlled experiment (see method section 4.5). Before computing PNN_{WHO} , PHATE visualizations were generated on matrices containing a consistent sampling of 4,000 libraries from each of Alpha, Delta, Omicron, and Beta WHO annotated lineages. For PHATE, the first boxplot represents the expected PNN_{WHO} values by chance, followed by all iSNVs, consensus only iSNVs, and *de novo* only iSNVs. The number of nearest neighbours used in this experiment is $k=40$.

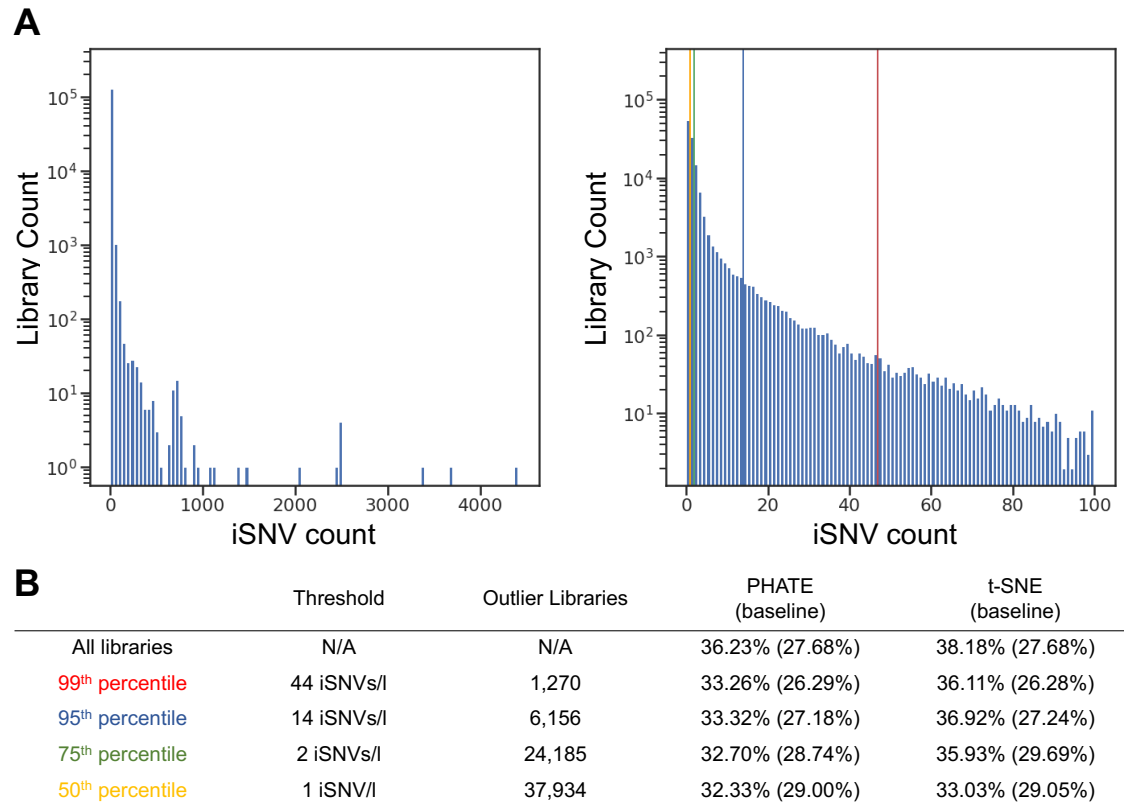


Figure S2: Analysis of Libraries' Mutational Load, which is the Number of iSNVs per Library. **A** The mutational load distribution across all libraries shows the variability in the number of iSNVs per library. **B** zoomed-in view of this distribution, focusing on libraries with up to 100 iSNVs. This view includes vertical lines to delineate various distribution percentiles. **C** A table summarizing the relationship between different outlier detection thresholds and their impact on library clustering structure on the PHATE visualizations. The table shows thresholds defined by the number of iSNVs per library, ranging from the 99th percentile (44 iSNVs/library) to the 50th percentile (1 iSNV/library). For each threshold, the table indicates the number of libraries classified as outliers and the corresponding percentage of nearest neighbours from the same WHO lineage (PNN_{WHO}), alongside the expected by chance PNN_{WHO} value.

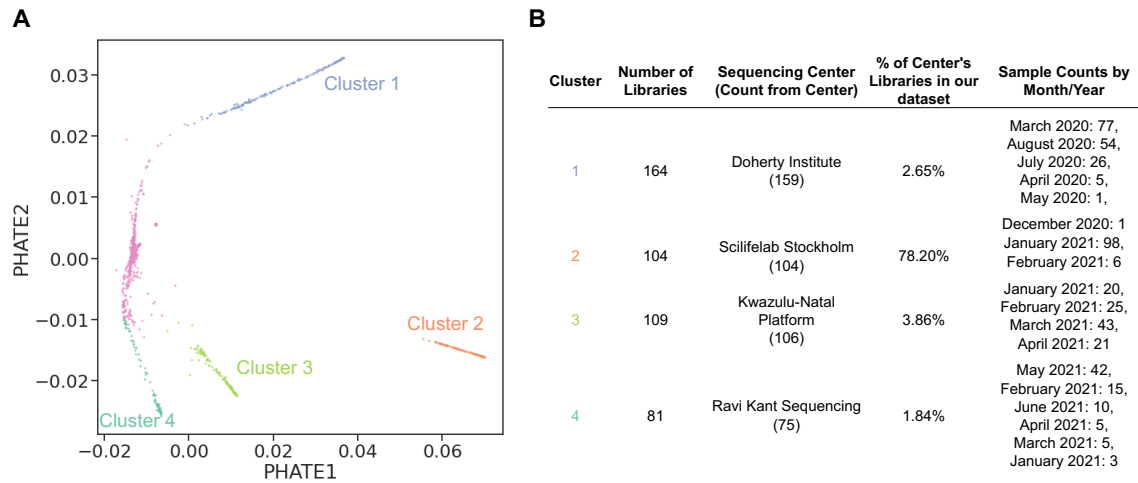


Figure S3: **A**: PHATE on the top 1% outlier libraries with the most iSNV count. The clusters on this PHATE representation were defined using K-means applied to the PHATE object. **B**: Table representing per cluster library information.

742 10 Supplemental Information

743 10.1 Details on Downloading SARS-CoV-2 Genomic Libraries from NCBI

744 A total of 147,537 SARS-CoV-2 Illumina amplicon paired-end sequencing reads were down-
745 loaded from NCBI, as follows: 51,837 Illumina SARS-CoV-2 sequencing libraries were down-
746 loaded from the NCBI database on June 4th, 2021, and another 95,700 Illumina sequencing
747 libraries were downloaded on February 12th, 2022. January and February 2020 were severely
748 underrepresented compared to the other months (Figure 2A). Most downloaded sequences
749 originated from Europe, constituting 75% of the dataset. Among the European sequences,
750 63% were obtained from the Wellcome Sanger Institute sequencing center, UK (Table S3), in-
751 dicating their significant contribution to the global sequencing efforts. Furthermore, a notable
752 number of downloaded libraries came from The Peter Doherty sequencing center, Australia,
753 between January and October 2020 (16% of the total libraries Table S3) as they led the
754 sequencing effort during that time in that region. Additionally, the dataset was enriched
755 with samples sequenced by North American sequencing centers, accounting for 15% of the
756 downloaded sequences (Tables S2 and S3). The underrepresentation of samples from January
757 and February 2020 reflects a limitation in the available data during the initial stages of the
758 pandemic. However, despite the initial disparities in data collection, which reflect the current
759 practical challenges faced by the scientific community (Chen et al. 2022), this dataset remains
760 highly informative, successfully capturing the global diversity of SARS-CoV-2 throughout the
761 later months of 2020 and extending into 2021.

762 Out of the total libraries downloaded, 134,879 had a mean coverage C above 100, and a
763 total of 138,723 libraries had a breadth of coverage B above 10,000, meaning that at least
764 10,000 genomic positions were covered at a depth of 100X or higher (Figure 2). The inter-
765 section of both filters allowed us to keep 128,423 high-quality sequencing libraries for further
766 analysis. The distributions of the breadth of coverage and mean depth show heterogeneity in
767 the coverage of the downloaded sequencing libraries. We also note the grouping of some se-
768 quencing centers (e.g. Wellcome Sanger Institute in red) and not others (e.g. the CDC's Office
769 of Advanced Molecular Detection - CDC-OAMD), displaying a heterogeneity across sequenc-
770 ing centers and within sequencing centers. Because we downloaded a representative sampling
771 of the available data on the NCBI database, this coverage distribution likely represents the

772 coverage heterogeneity of the available data on NCBI.

773 **10.2 Strand Coverage Across the Genome**

774 We evaluated the variation in strand coverage along the genome in our dataset using the
775 Forward Strand Ratio (FSR), which revealed a highly unbalanced distribution across the virus
776 sequence (Figure 2C). Only 31% of the viral genome in our dataset has a balanced coverage
777 from the forward and reverse read strands. Specifically, 40% of the genome is covered by the
778 plus strand, which is the number of genomic positions of the genome with an average forward
779 strand ratio above 90%. In contrast, 29% of the genome is covered by the minus strand,
780 with an average minus strand ratio above 90%. Thus, strand bias statistics in SARS-CoV-2
781 genomes need to consider strand coverage when evaluating if a *de novo* iSNVs is a stand bias
782 artifact, which motivates the development of our strand bias likelihood metric S .

783 **10.3 Recurrent Strand Bias Artifacts**

784 To better characterize strand bias artifacts, we analyzed a total of 1,491,885 intra-host single
785 nucleotide variants (iSNVs) identified as potential strand bias artifacts, with a likelihood of
786 no strand bias below 1% ($S;0.01$). We first examined their alternative allele frequency (AAF)
787 distribution. The AAF distribution of these excluded iSNVs does not differ significantly from
788 that of the other iSNVs, suggesting that strand bias artifacts can happen across a spectrum
789 of intra-host frequencies. This confirms that filtering based solely on AAF is insufficient to
790 eliminate strand bias artifacts.

791 Several genomic positions were found to be recurrent within these putative strand bias
792 artifacts. We computed the expected number of libraries with strand bias artifacts at a given
793 position, which has a mean of 4 and a 99th percentile of 68 libraries. We identified 486 genomic
794 positions that have a strand bias artifact reported in more than 68 libraries, labelling them as
795 recurrent strand bias artifacts, which we masked in our analyses across all libraries. To ensure
796 the robustness of iSNV analyses and to prevent the inclusion of recurrent spurious iSNVs, we
797 recommend evaluating and possibly masking these genomic positions in future SARS-CoV-2
798 intra-host studies.

799 10.4 Sub-sampling experiments to balance WHO variants

800 In our dataset, Alpha and Delta are overrepresented compared to other SARS-CoV-2 variants,
801 which may cause biases in the analysis results since unbalanced sampling can influence cluster
802 formation and PNN_{WHO} values (see, for example, Figure 3A, which distinctly marks Alpha
803 and Delta as dominant clusters). To address this, we conducted controlled sub-sampling
804 experiments, selecting 1,000 libraries each from the Alpha, Beta, Delta, and Omicron variants
805 (see Method section 4.5), aiming to mitigate variant sampling bias on PNN_{WHO} values in
806 the PHATE representation of iSNV subsets. We evaluated the clustering by WHO lineage
807 across three iSNV sets: unfiltered raw iSNVs, consensus iSNVs, and *de novo* iSNVs (Figure
808 S1). The raw and consensus iSNV datasets show high PNN_{WHO} values, indicating a strong
809 lineage-specific signature, primarily driven by frequent lineage-defining mutations, even when
810 samples per WHO variant are balanced. Conversely, *de novo* iSNVs exhibit lower PNN_{WHO}
811 values, indicating a subtler lineage-based structure but still above baseline, underscoring the
812 lineage-specific biological significance of emerging mutations. These controlled subsampling
813 experiments thus replicate our main findings with the full dataset (Figure 4). Therefore, the
814 lineage-specific signatures observed in our study are not a result of the uneven sampling of
815 WHO variants.

816 10.5 t-SNE Results Are Comparable to PHATE

817 In this section, we present results from t-SNE (t-Distributed Stochastic Neighbor Embedding)
818 analysis of SARS-CoV-2 genomic data, complementing the PHATE results found in the result
819 section (2). The method t-SNE is a machine learning algorithm used for dimensionality
820 reduction, offering an alternative approach to PHATE.

821 The t-SNE representation of the 128,423 high-quality sequencing libraries reveals dis-
822 tinct clusters by WHO lineage for both raw and consensus iSNV subsets, consistent with
823 PHATE's findings. For raw iSNVs, the Proportion of Nearest Neighbors (PNN_{WHO}) for
824 t-SNE is 99.43%, closely aligned with PHATE's 98.39%. Similarly, for consensus iSNVs, t-
825 SNE's PNN_{WHO} of 99.05% parallels PHATE's 99.37%, highlighting both methods' consistent
826 ability to identify lineage-specific mutations across the iSNV sets. Conversely, *de novo* iSNVs
827 (representing emerging mutations within the host) show less pronounced lineage-specific than
828 consensus iSNVs clustering in t-SNE representation, with a PNN_{WHO} value of 59.37%. This

829 suggests a deviation from the strong lineage alignment observed in raw and consensus iSNVs,
830 indicating that while *de novo* iSNVs still correlate with lineage structure more than baseline,
831 the association is less direct. The structure observed in *de novo* iSNVs through t-SNE com-
832 plements PHATE’s analysis, demonstrating consistent underlying data patterns regardless of
833 the representation method used.

834 Using the 8,000,668 unfiltered *de novo* iSNVs, both t-SNE and PHATE visualizations re-
835 vealed significant sequencing center batch effects, with t-SNE showing slightly higher PNN_{SC}
836 values (66.50%) compared to PHATE (62.31%). This indicates that both dimensionality re-
837 duction techniques captured the influence of sequencing center-specific artifacts within the *de*
838 *novo* iSNV dataset. Efforts to refine the dataset and mitigate these artifacts involved apply-
839 ing thresholds on the strand bias metric (S) and the Alternative Allele Frequency (AAF).
840 These measures effectively reduced sequencing center-specific artifacts, as evidenced by de-
841 creased PNN_{SC} values in both visualization methods after applying the filters, with the
842 t-SNE value (38.18%) slightly higher than PHATE (36.23%). Applying the filters effectively
843 reduced sequencing center-specific artifacts, as evidenced by decreased PNN_{SC} values in both
844 representation methods.

845 Similarly to PHATE, we also computed the PNN_{SC} values in t-SNE representation after
846 sequentially removing the top 1%, 5%, and 25% of the libraries with the most iSNV counts
847 (Figure S2B). As opposed to PHATE, the PNN_{SC} value of t-SNE did not drastically decrease
848 after the removal of the top 1% of our outliers. However, the PNN_{SC} values for both t-SNE
849 and PHATE only met after the exclusion of more libraries down to only keeping libraries with
850 one iSNV (Figure S2B, 50th percentile), underlining the stronger impact of outlier libraries
851 on t-SNE compared to PHATE.

852 Similar to the approach used with PHATE, we calculated the PNN_{SC} values for t-SNE
853 after removing the top 1%, 5%, and 25% of libraries based on iSNV counts (Figure S2B).
854 Unlike PHATE, the PNN_{SC} for t-SNE did not significantly decrease with the removal of the
855 top 1% of libraries. Both t-SNE and PHATE PNN_{SC} values converged after removing more
856 libraries, ultimately comparable for their PNN_{SC} values only when retaining those with a
857 single iSNV (Figure S2B, 50th percentile). This indicates that t-SNE is more susceptible to
858 bias from outlier libraries compared to PHATE.

859 This overall consistency between dimensionality reduction methods serves as compelling

860 evidence that the data's underlying structure is method-independent, suggesting that both
861 methods could be reliably applied to similar datasets to help inform future pre-processing
862 strategies in viral genomics. This alignment helps validate our pre-processing strategies in
863 viral genomics, demonstrating the robustness of our observations and the general applicability
864 of these techniques to analyze viral genomic data.