

1 **SARS-CoV-2 sequencing with cloud-based analysis illustrates expedient co-ordinated**
2 **surveillance of viral genomic epidemiology across six continents**

3
4 **Authors**

5 Daniel Amoako,¹ Nguyen To Anh,² Marc Brouard,³ Constanza Campano Romero,⁴ Andres
6 Castillo Ramirez,⁴ Bede Constantinides,³ Derrick Crook,^{3,5} Phan Manh Cuong,⁶ Moussa
7 Moise Diagne,⁷ Amadou Diallo,⁷ Nguyen Thanh Dung,⁸ Laura Dunn,⁵ Le Van Duyet,⁶ Josie
8 Everatt,¹ Katherine Fletcher,³ Philip Fowler,³ Mailie Gail,⁹ Hospital for Tropical Diseases
9 SARS-CoV-2 testing team,⁸ Nguyen Thi Thu Hong,² Martin Hunt,^{3,11} Zam Iqbal,^{3,11} Katie
10 Jeffery,⁵ Dikeledi Kekana,¹ Thomas Kesteman,² Jeff Knaggs,^{3,11} Marcela Lopes Alves,³ Dinh
11 Nguyen Huy Man,⁸ Amy J. Mathers,¹⁰ Nghiem My Ngoc,⁸ Sarah Oakley,⁵ Hardik Parikh,¹⁰
12 Tim Peto,^{3,5} Marcelo Rojas Herrera,⁴ Nicholas Sanderson,³ Vitali Sintchenko,⁹ Jeremy
13 Swann,³ Nguyen Thi Tam,² Le Van Tan,² Pham Ngoc Thach,⁶ Ndeye Marieme Top,⁷ Nguyen
14 Thu Trang,² Van Dinh Trang,⁶ H. Rogier Van Doorn,² Anne von Gottberg,^{1,12} Nicole
15 Wolter,^{1,12} Bernadette C Young^{3,5}

16
17 **Affiliations**

- 18 1. The National Institute for Communicable Diseases (NICD), a division of the National
19 Health Laboratory Service, Johannesburg., Republic of South Africa
20 2. Oxford University Clinical Research Unit, Ho Chi Minh City, Vietnam.
21 3. Nuffield Department of Medicine, University of Oxford, Oxford, United Kingdom
22 4. El Instituto de Salud Pública de Chile, Chile
23 5. Oxford University Hospitals NHS Trust, Oxford, United Kingdom
24 6. National Hospital for Tropical Diseases, Hanoi, Vietnam
25 7. The Institut Pasteur de Dakar, Senegal
26 8. Hospital for Tropical Diseases, Ho Chi Minh City, Vietnam
27 9. Microbial Genomics Reference Laboratory, Institute of Clinical Pathology and
28 Medical Research, New South Wales Health Pathology, Sydney, Australia

- 29 10. Division of Infectious Diseases, School of Medicine, University of Virginia,
30 Charlottesville, United States of America
- 31 11. European Bioinformatics Institute, Hinxton, UK
- 32 12. Faculty of Health Sciences, University of the Witwatersrand, Johannesburg, South
33 Africa
- 34
- 35 Corresponding author: Bernadette Young (Bernadette.young@ndm.ox.ac.uk)
- 36
- 37

38 **Abstract**

39 Viral sequencing has been critical in the COVID-19 pandemic response, but sequencing and
40 bioinformatics capacity remain inconsistent. To examine the utility of a cloud-based
41 sequencing analysis platform for SARS-CoV-2 sequencing, we conducted a cross-sectional
42 study incorporating seven countries in July 2022. Sites submitted sequential SARS-CoV-2
43 sequences over two weeks to the Global Pathogen Analysis Service (GPAS). The GPAS
44 bioinformatics cloud platform performs sequence assembly plus lineage and related sample
45 identification. Users can share information with collaborators while retaining data ownership.
46 Seven sites contributed sequencing reads from 5,346 clinical samples, of which 4,799/5,346
47 (89.8%) had a lineage identified. Omicron lineages dominated, with the vast majority being
48 BA.5, BA.4 and BA.2, commensurate with contemporary genomic epidemiological
49 observations. Phylogenetic analysis demonstrated low within-lineage diversity, and highly
50 similar sequences present in globally disparate sites. A cloud-based analysis platform like
51 GPAS addresses bioinformatics bottlenecks and facilitates collaboration in pathogen
52 surveillance, enhancing epidemic and pandemic preparedness.

53

54

55

56

57 **Introduction**

58 Viral genomic sequencing has been pivotal in understanding and responding to SARS-CoV-2
59 viral evolution during the COVID-19 pandemic. Global sharing of data and comparison of
60 samples has facilitated public health surveillance by tracking the emergence of variants
61 globally,[1] demonstrating the impact of travel restrictions on viral dynamics across
62 continents,[2] and of travel and quarantine on spread within countries,[3,4] as well as
63 demonstrating transmission routes within hospitals.[5] As successive waves driven by new
64 viral variants emerged, it became clear that rapidly detecting new lineages is critical both for
65 understanding and predicting disease epidemiology as well as for guiding both public health
66 measures and treatment of infection, including the use of neutralising antibodies.[6,7]
67 However there remains marked global variability in sequencing costs and capacity.[8] Much
68 SARS-CoV-2 sequencing has been undertaken at centralised reference laboratories, but the
69 speed and power of this tool is greatly enhanced when researchers and clinicians at any
70 location can contribute to the global view of viral evolution and surveillance.[9]

71

72 Several important challenges remain in expanding pathogen genomic surveillance, including
73 access to sequencing technology, the availability of bioinformatic resources and the global
74 interpretability of results generated using different bioinformatics approaches.[10,11,12] One
75 solution to the relative shortage of bioinformatic expertise in many settings is web-accessible
76 analysis infrastructure.[13] The Global Pathogen Analysis Service (GPAS) is a cloud-based
77 globally web-accessible platform that provides secure and rapid access to comprehensive
78 SARS-CoV-2 sequencing genome assembly and analysis software. GPAS supports multiple
79 sequencing platforms and tiling PCR primer schemes, and has been validated against widely
80 used methods. Users can opt to share their data on this platform, allowing those who opt in to
81 view their own results in a wider context, facilitating comparisons for rapid and
82 comprehensive assessments of viral data. The platform is designed for use by clinical and
83 public health laboratory scientists and has been established and used for over a year in a
84 tertiary hospital based clinical microbiology laboratory to ISO15189:2012 standards. Scaling

85 up this simple-to-use assembly, variant calling and analysis tool for SARS-CoV-2 empowers
86 laboratories to control the analysis of their sequence data. It provides proof of principle that
87 standard microbiology laboratories without bioinformatics expertise from across the world
88 can generate outputs for local surveillance and automated submission to public databases.
89 Such a service is an exemplar of how cloud-based services, accessible in all settings with
90 internet access, offer great promise for overcoming an important barrier to effective genomic
91 surveillance.

92

93 The primary aim of this project was to demonstrate that a globally broad snapshot of SARS-
94 CoV-2 viral diversity can be simultaneously gathered, processed and shared through a cloud-
95 based software, enabling aggregation for combined descriptive analysis.

96

97 **Methods**

98 *Sample frame and sequencing:* Sequencing centres participated in a two-week pilot in July to
99 August 2022. Centres were from seven countries: Senegal; Chile; South Africa; New South
100 Wales (NSW), Australia; Vietnam; United Kingdom; Virginia, United States (USA). Centres
101 were either accredited clinical microbiology or public health laboratories. All clinical
102 samples in which SARS-CoV-2 was detected were eligible for inclusion and underwent in-
103 house genomic sequencing at participating sites during the pilot were eligible for inclusion,
104 with no more than one submission from an individual patient.

105

106 Sequencing was performed using either Illumina or Oxford Nanopore Technology (ONT)
107 platforms. Multiple primer schemes were supported (ARCTIC V3, V4 or 4.1, and Midnight).
108 Sequencing platform and primer schemes were chosen by participating sites, who followed
109 their established sequencing protocols.

110

111 *Bioinformatic methods:* Unselected raw sequences were submitted to GPAS in FASTQ
112 format with sequencing information (instrument platform, primer scheme) and limited

113 associated metadata (sample name, whether sequence data was a control or study sample,
114 collection date, country of origin). At upload, a Universally Unique Identifier (UUID) was
115 randomly generated and assigned to each sample. The mapping between the centre's sample
116 identifier and the GPAS UUID was held only by the submitting organisation, ensuring that no
117 potentially personally identifiable information was transmitted to the GPAS platform. When
118 results were returned to the submitter after analysis, they could hence be re-linked to the
119 subject's record. At upload, only reads mapping to SARS-CoV-2 were kept,[14] guaranteeing
120 that no human reads were retained and pushed to the cloud. Tight access control to data
121 applied within GPAS. With the permission of all collaborators, data access controls were
122 configured such that all submitters could access each other's sequences, metadata and
123 analytical output data for the duration of the investigation.

124

125 Sequence data analysis was performed by the GPAS SARS-CoV-2 bioinformatics pipeline.
126 Genome assembly was performed by an amplicon-aware *de novo* genome assembly tool
127 Viridian v0.3.7.[15] Viridian applies quality control (QC) at both a read and amplicon level,
128 to produce a consensus assembly. Each genomic position is called A,C,G or T only where 10
129 or more reads passing amplicon-aware QC provide at least 70% support, otherwise it is called
130 ambiguously as N. After assembly, clade assignment and amino acid mutation calls were
131 made using Nextclade,[16] and Pango lineages assigned with Pangolin.[17] Samples were
132 further labelled according to the United Kingdom Health Security Agency (UKHSA)
133 nomenclature to define and track variants of concern using aln2type.[18] The GPAS QC suite
134 provided read-level QC output for each sample, reporting the number of mapped reads per
135 sample, median read depth, reference coverage at read depths of 10 and 20, and the number
136 of amplicons failing quality control. Consensus genome assemblies were aligned by Viridian
137 to the Wuhan-Hu-1 reference genome (MN908947.3).[19] Aligned sequences were compared
138 using findNeighbour4, which identifies sequences without conflicting variant calls (0 single
139 nucleotide variants (SNVs)), as well as those differing by one, two and three SNVs.[20,21]

140

141 Submitting sites verified the run quality reports of each sequencing batch, including review
142 of positive and negative controls, and outputs were passed or failed accordingly. All passed
143 samples were tagged for aggregate analysis, and released to the GPAS shared data pool,
144 which identifies other samples within three SNVs in the shared data pool. The GPAS
145 platform recorded the total number of sequencing runs and total number of samples. Analysis
146 duration was also monitored, along with descriptive reports of any unexpected issues with the
147 platform.

148

149 *Phylogenetic analysis:* In Oxford, genome assemblies in the shared data pool were
150 downloaded from the GPAS portal. A maximum likelihood phylogeny of aligned sequences
151 was constructed for the three largest Pango lineages in the sample set (BA.5, BA.4 and BA.2)
152 in using RAxML [22] (assuming a general time reversible (GTR) nucleotide substitution
153 model).

154

155

156 Results

157 All sites tagged their sequencing data and results granting access for aggregated analysis.
158 5,437 sequences were shared (Table 1, Table S1). As different sequencing centres worked
159 with samples collected at different times, the date of collection of submitted sequences varied
160 by centre, ranging from April to July 2022. Of submitted samples 91 (1.7%) were controls
161 and 5,346 (98.3%) were clinical samples. Among these, 4,797/5,346 (89.7%) were assembled
162 with at least 70% coverage (at a depth of 10 reads), and 4,799/5,346 (89.8%) were assembled
163 with sufficient coverage to be assigned a lineage (Table 2). 4,775/4,799 (99.5%) were
164 Omicron variants, with BA.5, BA.4 and BA.2 being the most common. A small number of
165 Delta variant sequences were identified (23/4,799 0.5%), all of which were collected prior to
166 July 2022.

167

168 Table 1: Samples submitted by study centre, with date of collection and genome coverage
169 (with minimum depth of 10 reads)

Centre	Submitted	Controls	<50% coverage	50-70% coverage	>70% coverage	Earliest sequence	Latest sequence
Senegal	197	0	14	10	173	13 th July 2022	20 th July 2022
Chile	1205	20	4	7	1174	16 th June 2022	21 st July 2022
South Africa	202	0	9	4	189	17 th May 2022	24 th June 2022
NSW, Australia	1194	31	34	23	1106	21 st June 2022	18 th July 2022
Vietnam	316	3	20	5	288	1 st April 2022	14 th July 2022
United Kingdom	1823	37	319	100	1367	12 th July 2022	26 th July 2022
Virginia, USA	500	0	0	0	500	7 th June 2022	29 th June 2022
Total	5437	91	400	149	4797	1 st April 2022	26 th July 2022

170

171

172

173

174

175

176 Table 2: SARS-CoV-2 lineage by study centre (where a lineage was assigned by Pangolin)

Centre	Omicron							Delta	Alpha	Other	Total
	BA.5	BA.4	BA.2	BE.1	BF.1	BA.1	BG.2	Delta	B.1.1.7	B.28	
Senegal	69	47	34	5	18	0	0	0	0	0	173
Chile	397	607	145	20	4	0	0	0	0	0	1,173
South Africa	89	86	3	6	0	5	0	0	0	0	189
NSW, Australia	736	130	184	54	3	0	0	0	0	0	1,107
Vietnam	8	0	256	0	0	4	0	22	0	0	290
United Kingdom	974	169	57	140	23	2	0	0	1	1	1,367
Virginia, USA	95	94	296	11	1	1	1	1	0	0	500
Total	2,368	1,133	975	236	49	12	2	23	1	1	4,799

177
 178 Comparing only clinical samples collected between 1st June 2022 and 31st July 2022, the
 179 proportion of samples assigned to the different Omicron sub-lineages varied substantially by
 180 study site (Figure 1). Each of BA.2, BA.4 or BA.5 dominated (constituted >50% of samples
 181 from) at least one site, while Senegal had no single dominant lineage. For each of the three
 182 most common lineages, a maximum likelihood phylogeny was constructed to compare
 183 diversity within and between sites (Figure 2).

184
 185 Phylogenetic analysis revealed global mixing of Omicron lineages. There were no examples
 186 of sub-lineages being completely geographically restricted within this study. One sub-lineage
 187 in BA.4 was predominantly found in Chile (Figure 2), however examples of this sub-lineage
 188 were also identified in the United Kingdom. Similarly, a sub-lineage of BA.2 was mainly
 189 found in Vietnam (Figure 2) but again an example of this genotype was found in the United
 190 Kingdom.

191
 192 Highly genetically similar samples were identified within lineages from globally distinct
 193 sites, with branch lengths corresponding to $\leq 0.01\%$ of the genome (approximately three
 194 nucleotides) separating many samples. This limited diversity was further explored by
 195 examining the relative likelihood of samples being from the same site with increasing
 196 genomic variation. By this measure, the greatest dispersal was seen within the BA.4 lineage

197 (Figure 3). In total 357/1,133 of BA.4 samples were genotypically identical to another
198 sequenced sample. Among these, 244/357 (68.4%) were only identical to other samples from
199 the same study centre. Allowing for just one single nucleotide variant (SNV) reduced the
200 proportion of closely related samples to the same study centre (and country) to 302/634
201 (47.6%). The proportions reduced to 279/819 (34.1%) and to 201/923 (21.8%) when allowing
202 for two and three SNVs respectively. Only 210/1,133 (18.5%) did not have a related sample
203 allowing up to three SNVs. The other major lineages BA.2 and BA.5 likewise showed
204 evidence of genotypically identical sequences identified at different study sites, with
205 increasing evidence of dispersal as SNV threshold increased (Figure 3).

206

207 Cloud processing performance was measured by time elapsed from upload until completion
208 of processing. The median sample processing time was 30.6 minutes (IQR 13.1-69.9
209 minutes). The distribution of processing times was much wider than expected, and processing
210 delays occurred for several batches from Australia, United Kingdom and United States
211 (Figure S1). Some of these coincided with data centre outages during an extreme heat event
212 in the United Kingdom, while other delays occurred due to errors in a software update for the
213 upload portal, which were resolved by improvements implemented in the user side command
214 line interface software. Other data upload issues were identified and addressed in real time,
215 including metadata file formatting errors, availability of upload clients for all required
216 operating systems, and user interface display errors preventing batch release.

217

218 **Discussion**

219 This pilot study reveals the potential for cloud-based synchronous global data processing and
220 analysis for public health. The GPAS platform provides an example of unified assembly,
221 variant calling and analysis of data from participating centres, which can facilitate both local
222 and globally aggregated analysis. In this cross-sectional observational study, the aggregate
223 analysis conforms to contemporary observations of SARS-CoV-2 genetic

224 epidemiology[23,24,25,26,27], though these observations remain limited by global disparities
225 in sequencing volume.

226

227 The opportunistic sampling frame used means this study is likely to be geographically
228 incomplete, but even with this limitation we observed cosmopolitan global mixing of
229 Omicron lineages between study sites. Even a BA.4 sub-lineage which demonstrated
230 extensive replacement in one geographical location (Chile) was not completely restricted to
231 that site and could be identified in a geographically distant site (United Kingdom). During the
232 period of the study, the United Kingdom site performed genomic sequencing on all clinical
233 isolates which had tested positive for SARS-CoV-2 on PCR, and despite being dominated by
234 BA.5, was also able to identify examples of a BA.2 sub-lineage that was otherwise restricted
235 to Vietnam. Such mixing is likely the result of movement for travel and migration, as SARS-
236 CoV-2 related travel restrictions had been relaxed in the countries with participating sites.

237

238 Given evidence of global mixing and limited diversity between sites, viral genome
239 sequencing is limited in its ability to inform studies of local viral transmission in this
240 Omicron dominated era, where both rapid global mixing and transmission drives spread.[26]
241 Allowing for up to three different single nucleotide variants, the majority of isolates from the
242 dominant lineages BA.2, BA.4 and BA.5 had at least one closely related sample within each
243 countries' contemporaneously sequenced samples. Depending on lineage, between 10.0%
244 and 31.5% of genotypically identical sequences in this study were identified in different
245 countries and continents. The use of a shared bioinformatics analysis pipeline, as described
246 here, highlights how lineages can be simultaneously and confidently mapped across distant
247 sites. This captures the evidence of well described rapid global dispersal of Omicron lineages
248 with successive selective sweeps of these new Omicron variants over short time scales.[27]

249

250 These observations highlight the potential for deriving insights from viral sequencing at a few
251 globally distributed sites. While this study pertains to a specific point in the pandemic, it

252 raises the question of how to better design genomic surveillance balancing the public health
253 need and available resources. It is recognised that a preponderance of sequence data came
254 from high-income countries, who sequenced a proportion of cases ten-fold higher than that
255 achieved in low- and middle-income countries during the first two years of the pandemic.[9]
256 Nor was the data collected evenly over the course of pandemic. In retrospect this was not
257 optimal. Better designed representative sampling would be immensely empowering and is
258 likely possible through well-designed and suitably powered sampling frames. Ideally, this
259 would be addressed through pandemic preparedness planning, enabling sequencing to be
260 focussed at selected sentinel sites globally cognisant of national jurisdictions and the benefit
261 of longitudinal data, as exemplified by a cohort design in the UK.[25]

262

263 Pathogen sequencing played a highly valuable role in the COVID-19 pandemic and remains a
264 critical global and national capability. This is highlighted by studies of viral transmission that
265 have yielded important insights into the modes of pathogen spread.[1-5] In addition,
266 identification of emerging variants-of-concern guided public health and policy responses,[28]
267 as well as informing assessments about the impact of emerging variants on the impact of
268 therapeutics (such as monoclonal antibodies) and vaccines.[7,29] The present cross-sectional
269 study has the limitation of lacking longitudinal data on how these findings changed over time.
270 Nevertheless, this study demonstrates how simply structured sampling across the world is
271 well placed to rapidly identify replacement by new lineages. The integration of genomic data
272 with subject disease manifestation and severity data would enable the ready investigation of
273 associations between genomic variation and severity of disease, as suggested by others [30]
274 and predictive modelling for anticipating the course of future epidemics as described so
275 effectively early in the SARS-CoV-2 pandemic [31].

276

277 The world is now better prepared to urgently initiate genomic pathogen surveillance. The
278 available sequencing platform capacity is much improved, and the achievements of global
279 data aggregation (illustrated by INSDC and GISAID databases) has alleviated a major

280 obstacle to effective genomic surveillance. A remaining obstacle is local turnkey
281 bioinformatics to perform sequence assembly, variant calling and analysis. These are needed
282 to deliver standardised and quality assured outputs – specifically the SARS-CoV-2 viral
283 sequence FASTA file – for phylogenetic analysis either for global analysis by an
284 international entity such as GISAID, or for local and national analyses. GPAS has developed
285 an exemplar cloud-based, web-enabled service design that is optimised for operation by a
286 laboratory scientist and removes dependence on in-laboratory bioinformatics expertise. Thus,
287 this work reveals the benefit of establishing software for simply and automatically ingesting
288 sequence data from a sequencer in a local lab and yielding data ready for local or
289 international analysis. A software infrastructure such as the GPAS platform fills a gap in
290 delivering global genomic surveillance. Flexible privacy protections allows users to control
291 sharing of genomic data. GPAS is freely available to low- and middle-income countries,
292 reducing some of the barriers to analysis services in low resource settings. Remaining
293 challenges include high speed computer network access and local epidemiological tools, and
294 expertise for detailed epidemiological analysis.

295

296 **Conclusions**

297 The COVID-19 pandemic has demonstrated the multiple ways in which viral sequencing can
298 inform pandemic responses, but resources and capabilities are not equitably distributed. A
299 service like GPAS democratises SARS-CoV-2 genomics, enabling researchers and
300 laboratories, regardless of their location or level of bioinformatics expertise, to participate
301 actively in global surveillance efforts. By providing accessible and user-friendly tools for
302 sequence assembly, lineage assignment, and data sharing, it fosters collaboration, harmonises
303 data, and ultimately enhances pandemic preparedness.

304

305

306 **Funding statement**

307 GPAS is free at the point of use for users in low- and middle-income countries. GPAS was
308 developed specifically to minimise the need for bioinformatics staff (a scarce and expensive
309 resource), to help bring genomic insights within reach for laboratories with limited resources.

310

311 GPAS operates as a non-profit organisation. All GPAS SARS-CoV-2 services were offered
312 free of charge to all sites for the purposes of this study. Staff costs for the 5C analysis team
313 were met by the University of Oxford.

314

315 The cloud infrastructure used in this study was donated by Oracle Corporation.

316

317 Sequencing activities for NICD are supported by a conditional grant from the South African
318 National Department of Health as part of the emergency COVID-19 response; a cooperative
319 agreement between the National Institute for Communicable Diseases of the National Health
320 Laboratory Service and the United States Centers for Disease Control and Prevention
321 (NU51IP000930); the South African Medical Research Council (SAMRC) with funds
322 received from the South African Department of Science and Innovation; the African Society
323 of Laboratory Medicine (ASLM) and Africa Centers for Disease Control and Prevention
324 through a sub-award from the Bill and Melinda Gates Foundation grant number INV-018978;
325 Africa PGI, the UK Foreign, Commonwealth and Development Office and Wellcome (Grant
326 no 221003/Z/20/Z); and the Department of Health and Social Care's Fleming Fund using UK
327 aid.. NICD sequencing was also supported by The Coronavirus Aid, Relief, and Economic
328 Security Act (CARES ACT) through the Centers for Disease Control and Prevention (CDC)
329 and the COVID International Task Force (ITF) funds through the CDC under the terms of a
330 subcontract with the African Field Epidemiology Network (AFENET) AF-NICD-001/2021.

331 Genomic surveillance conducted in Vietnam was supported by the Wellcome Trust
332 (222574/Z/21/Z). L.V.T. is supported by the Wellcome Trust of Great Britain
333 (204904/Z/16/Z and 226120/Z/22/Z).

334

335 **Ethical statement**

336 SARS-CoV-2 sequencing was performed for public health surveillance and ethical approval
337 for secondary analysis was not required. This determination was reviewed by the University
338 of Oxford Joint Research Office. No patient identifying information was shared as part of this
339 study. Sample identifiers remain with the submitter, and are never kept on GPAS or shared
340 with other users. GPAS generates identifiers for each sequence and writes a file linking
341 anonymised ID to the submitted identifiers. Only the user submitting sequence data has
342 access to these records. The GPAS upload client selects only SARS-CoV-2 sequence reads
343 and discards non SARS-CoV-2 reads (including all human sequences).

344

345 **Data availability**

346 Submitted sequencing reads (free from human reads) for all included samples are available
347 from European Nucleotide Archive study accession PRJEB70597.

348

349 **Conflicts of interest**

350 AvG and NW have received grant funding from Sanofi and The Bill and Melinda Gates
351 Foundation.

352

353 **Authors**

354 All authors have seen and approved the manuscript

355

356 **Acknowledgments**

- 357 • Microbial Genomics Reference Laboratory, New South Wales Health Pathology,
358 Sydney, Australia

- 359 • El Instituto de Salud Pública de Chile, Chile
- 360 • Centre for Respiratory Diseases and Meningitis (CRDM), National Institute for
- 361 Communicable Diseases (NICD), a division of the National Health Laboratory
- 362 Service, South Africa
- 363 • The Institut Pasteur de Dakar, Senegal
- 364 • Oxford University Hospitals NHS Trust, Oxford, United Kingdom
- 365 • University of Virginia, Charlottesville, United States of America
- 366 • Oxford University Clinical Research Unit, Vietnam
- 367 • Hospital for Tropical Diseases SARS-CoV-2 testing team, HTD Vietnam: Le
- 368 Manh Hung, Nguyen Le Nhu Tung, Nguyen Thanh Phong, Vo Minh Quang,
- 369 Pham Thi Ngoc Thoa, Nguyen Thanh Truong, Tran Nguyen Phuong Thao,
- 370 Dao Phuong Linh, Ngo Tan Tai, Ho The Bao, Vo Trong Vuong, Huynh Thi
- 371 Kim Nhung
- 372 • Oracle Global Health Business Unit

373

374 **References**

- 375 1. Li J, Lai S, Gao GF, Shi W. The emergence, genomic diversity and global spread of
- 376 SARS-CoV-2. *Nature*. 2021 Dec;600(7889):408-418. doi: 10.1038/s41586-021-
- 377 04188-6.
- 378 2. Wilkinson E, Giovanetti M, Tegally H, San JE, Lessells R, Cuadros D, et al. A year of
- 379 genomic surveillance reveals how the SARS-CoV-2 pandemic unfolded in Africa.
- 380 *Science*. 2021 Oct 22;374(6566):423-431. doi: 10.1126/science.abj4336.
- 381 3. McCrone JT, Hill V, Bajaj S, Pena RE, Lambert BC, Inward R, et al. Context-specific
- 382 emergence and growth of the SARS-CoV-2 Delta variant. *Nature*. 2022
- 383 Oct;610(7930):154-160. doi: 10.1038/s41586-022-05200-3.
- 384 4. Raghwani J, du Plessis L, McCrone JT, Hill SC, Parag KV, Thézé J, et al. Genomic
- 385 Epidemiology of Early SARS-CoV-2 Transmission Dynamics, Gujarat, India. *Emerg*
- 386 *Infect Dis*. 2022 Apr;28(4):751-758. doi: 10.3201/eid2804.212053.

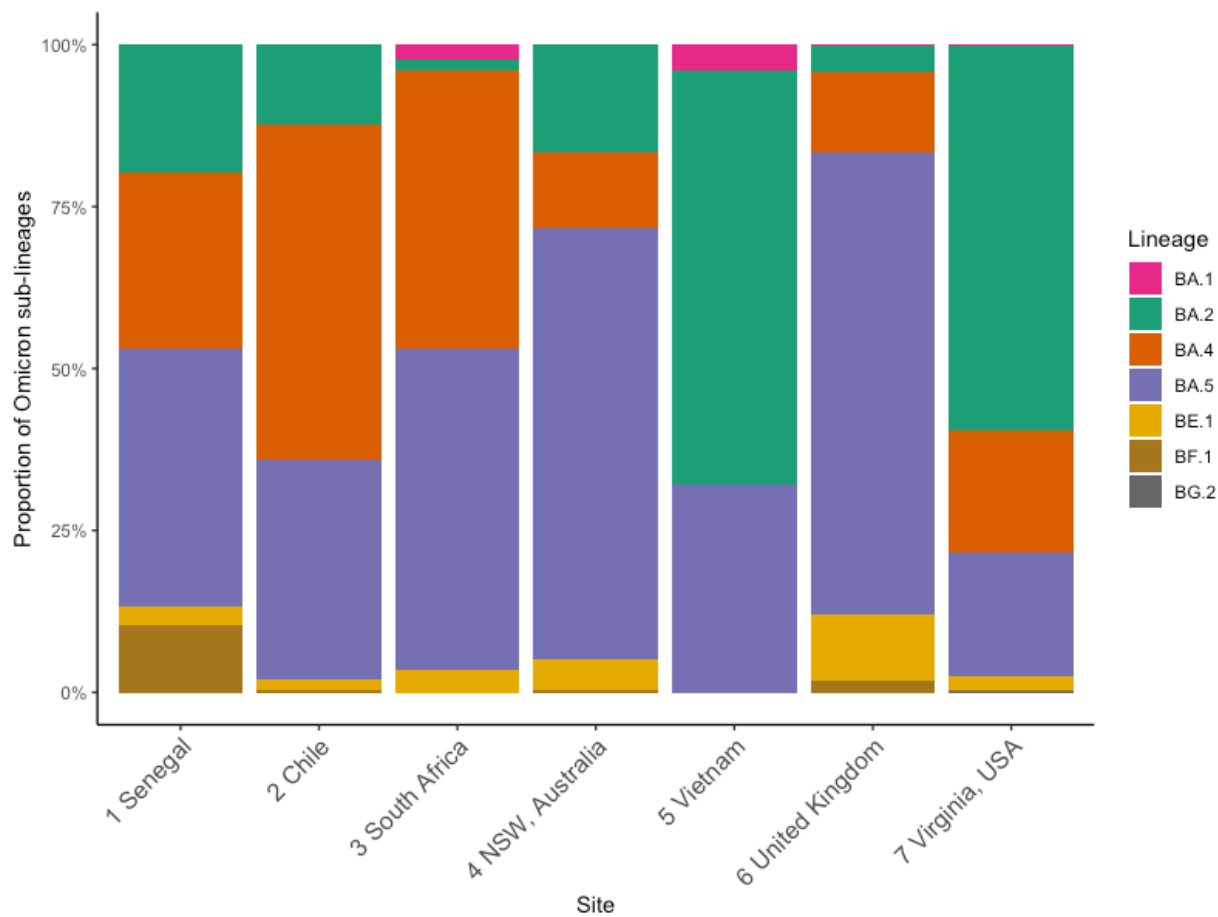
- 387 5. Snell LB, Fisher CL, Taj U, Stirrup O, Merrick B, Alcolea-Medina A, et al.
388 Combined epidemiological and genomic analysis of nosocomial SARS-CoV-2
389 infection early in the pandemic and the role of unidentified cases in transmission. *Clin*
390 *Microbiol Infect.* 2022 Jan;28(1):93-100. doi: 10.1016/j.cmi.2021.07.040
- 391 6. Tao K, Tzou PL, Nouhin J, Gupta RK, de Oliveira T, Kosakovsky Pond SL, et al. The
392 biological and clinical significance of emerging SARS-CoV-2 variants. *Nat Rev*
393 *Genet.* 2021 Dec;22(12):757-773. doi: 10.1038/s41576-021-00408-x.
- 394 7. Tuekprakhon A, Nutalai R, Dijokaite-Guraliuc A, Zhou D, Ginn HM, Selvaraj M, et
395 al. Antibody escape of SARS-CoV-2 Omicron BA.4 and BA.5 from vaccine and
396 BA.1 serum. *Cell.* 2022 Jul 7;185(14):2422-2433.e13. doi:
397 10.1016/j.cell.2022.06.005. Epub 2022 Jun 9. PMID: 35772405; PMCID:
398 PMC9181312.
- 399 8. Merhi G, Koweyes J, Salloum T, Khoury CA, Haidar S, Tokajian S. SARS-CoV-2
400 genomic epidemiology: data and sequencing infrastructure. *Future Microbiol.* 2022
401 Sep;17:1001-1007. doi: 10.2217/fmb-2021-0207.
- 402 9. Brito AF, Semenova E, Dudas G, Hassler GW, Kalinich CC, Kraemer MUG, et al.
403 Global disparities in SARS-CoV-2 genomic surveillance. *Nat Commun.* 2022 Nov
404 16;13(1):7003. doi: 10.1038/s41467-022-33713-y
- 405 10. Global genomic surveillance strategy for pathogens with pandemic and epidemic
406 potential, 2022–2032. Geneva: World Health Organization; 2022. Licence: CC BY-
407 NC-SA 3.0 IGO
- 408 11. Inzaule SC, Tessema SK, Kebede Y, Ogwell Ouma AE, Nkengasong JN. Genomic-
409 informed pathogen surveillance in Africa: opportunities and challenges. *Lancet Infect*
410 *Dis.* 2021 Sep;21(9):e281-e289. doi: 10.1016/S1473-3099(20)30939-7
- 411 12. Chen Z, Azman AS, Chen X, Zou J, Tian Y, Sun R, et al. Global landscape of SARS-
412 CoV-2 genomic surveillance and data sharing. *Nat Genet.* 2022 Apr;54(4):499-507.
413 doi: 10.1038/s41588-022-01033-y.

- 414 13. Ohlsen EC, Hawksworth AW, Wong K, Guagliardo SAJ, Fuller JA, Sloan ML, et al.
415 Determining Gaps in Publicly Shared SARS-CoV-2 Genomic Surveillance Data by
416 Analysis of Global Submissions. *Emerg Infect Dis.* 2022 Dec;28(13):S85-S92. doi:
417 10.3201/eid2813.220780.
- 418 14. Hunt M, Swann J, Constantinides B, Fowler PW, Iqbal Z. ReadItAndKeep: rapid
419 decontamination of SARS-CoV-2 sequencing reads. *Bioinformatics.* 2022 Jun
420 13;38(12):3291-3293. doi: 10.1093/bioinformatics/btac311.
- 421 15. https://github.com/iqbal-lab-org/viridian_workflow/wiki
- 422 16. Aksamentov I, Roemer C, Hodcroft EB, Neher RA. Nextclade: clade assignment,
423 mutation calling and quality control for viral genomes. *Journal of Open Source*
424 *Software*, 2021 6(67): 3773, <https://doi.org/10.21105/joss.03773>
- 425 17. O'Toole Á, Scher E, Underwood A, Jackson B, Hill V, McCrone JT, et al.
426 Assignment of epidemiological lineages in an emerging pandemic using the pangolin
427 tool. *Virus Evol.* 2021 Jul 30;7(2):veab064. doi: 10.1093/ve/veab064.
- 428 18. <https://github.com/connor-lab/aln2type>
- 429 19. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, et al. A new coronavirus
430 associated with human respiratory disease in China. *Nature.* 2020
431 Mar;579(7798):265-269. doi: 10.1038/s41586-020-2008-3
- 432 20. Mazariegos-Canellas O, Do T, Peto T, Eyre DW, Underwood A, Crook D, et al
433 BugMat and FindNeighbour: command line and server applications for investigating
434 bacterial relatedness. *BMC Bioinformatics.* 2017 Nov 13;18(1):477. doi:
435 10.1186/s12859-017-1907-2.
- 436 21. <https://github.com/davidhwyllie/findNeighbour4>
- 437 22. Stamatakis A. RAXML version 8: A tool for phylogenetic analysis and post-analysis
438 of large phylogenies *Bioinformatics.* 2014; 30(9),1312-3 doi:
439 10.1093/bioinformatics/btu033.
- 440 23. World Health Organisation Coronavirus (COVID-19) Dashboard,
441 <https://covid19.who.int> [Accessed 19th April 2023]

- 442 24. Tosta S, Moreno K, Schuab G, Fonseca V, Segovia FMC, Kashima S, et al. Global
443 SARS-CoV-2 genomic surveillance: What we have learned (so far). *Infect Genet*
444 *Evol.* 2023 Mar;108:105405. doi: 10.1016/j.meegid.2023.105405.
- 445 25. Foulkes S, Monk EJM, Sparkes D, Hettiarachchi N, Milligan ID, Munro K, et al.
446 Early Warning Surveillance for SARS-CoV-2 Omicron Variants, United Kingdom,
447 November 2021-September 2022. *Emerg Infect Dis.* 2023 Jan;29(1):184-188. doi:
448 10.3201/eid2901.221293.
- 449 26. Xu Y, Liu T, Li Y, Wei X, Wang Z, Fang M, et al. Transmission of SARS-CoV-2
450 Omicron Variant under a Dynamic Clearance Strategy in Shandong, China. *Microbiol*
451 *Spectr.* 2023 Mar 14;11(2):e0463222. doi: 10.1128/spectrum.04632-22.
- 452 27. Markov PV, Ghafari M, Beer M, Lythgoe K, Simmonds P, Stilianakis NI, et al. The
453 evolution of SARS-CoV-2. *Nat Rev Microbiol.* 2023 Apr 5. doi: 10.1038/s41579-
454 023-00878-2.
- 455 28. WHO Technical Advisory Group on SARS-CoV-2 Virus Evolution, Updated working
456 definitions and primary actions for SARS-Cov-2 variants, World Health Organisation,
457 March 2023. [https://www.who.int/publications/m/item/updated-working-definitions-](https://www.who.int/publications/m/item/updated-working-definitions-and-primary-actions-for--sars-cov-2-variants)
458 [and-primary-actions-for--sars-cov-2-variants](https://www.who.int/publications/m/item/updated-working-definitions-and-primary-actions-for--sars-cov-2-variants) [Accessed 19th April, 2023]
- 459 29. Aggarwal A, Akerman A, Milogiannakis V, Silva MR, Walker G, Stella AO, et al.
460 SARS-CoV-2 Omicron BA.5: Evolving tropism and evasion of potent humoral
461 responses and resistance to clinical immunotherapeutics relative to viral variants of
462 concern. *EBioMedicine.* 2022 Oct;84:104270. doi: 10.1016/j.ebiom.2022.104270.
- 463 30. Oude Munnink BB, Worp N, Nieuwenhuijse DF, Sikkema RS, Haagmans B,
464 Fouchier RAM, et al The next phase of SARS-CoV-2 surveillance: real-time
465 molecular epidemiology. *Nat Med.* 2021 Sep;27(9):1518-1524. doi: 10.1038/s41591-
466 021-01472-w.
- 467 31. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, et al. Assessing
468 transmissibility of SARS-CoV-2 lineage B.1.1.7 in England. *Nature.* 2021
469 May;593(7858):266-269. doi: 10.1038/s41586-021-03470-x.

470 **Figures**

471

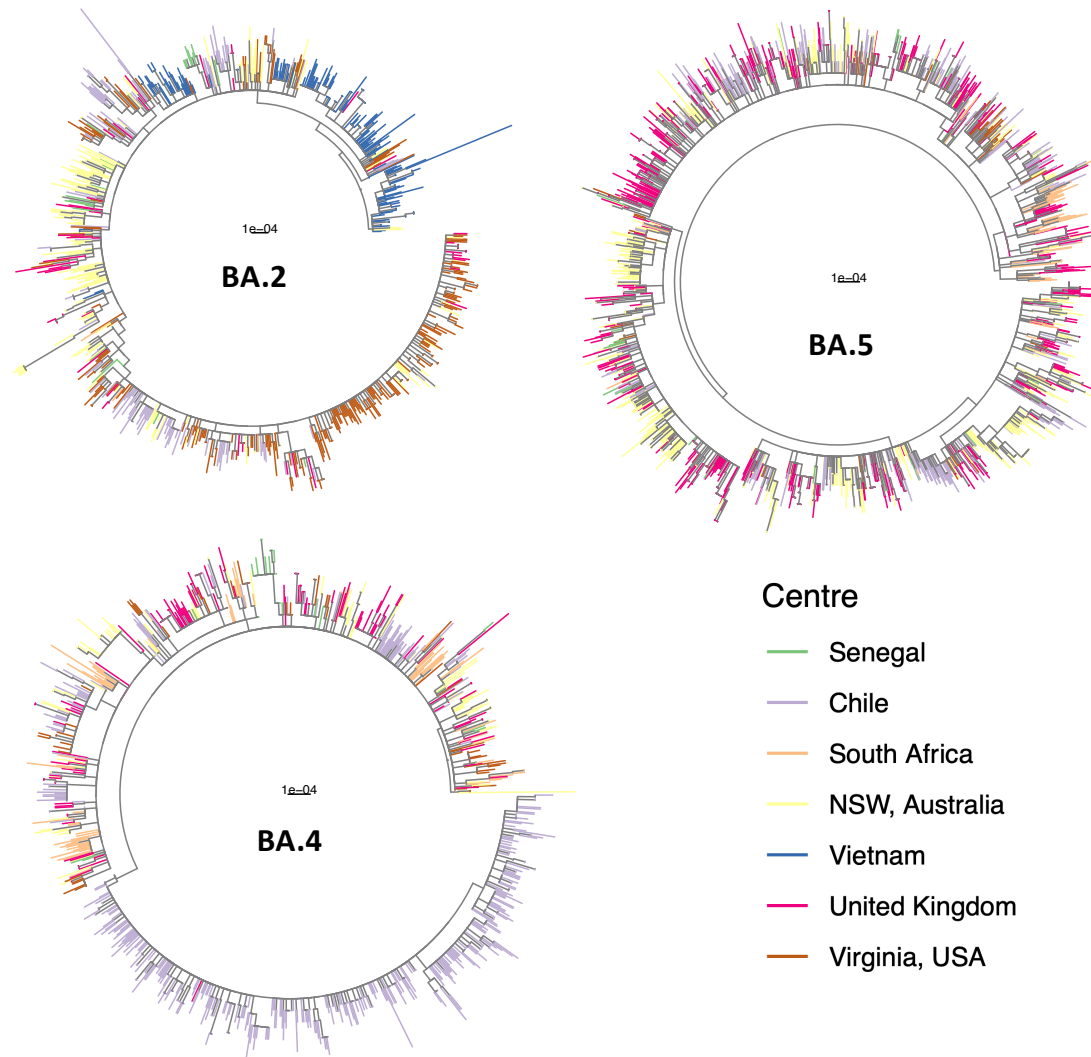


472

473

474 **Figure 1:** Omicron sub-lineages by study site. Proportion of samples assigned to Omicron
475 sub-lineage within each study site, where collection date was in June or July 2022.

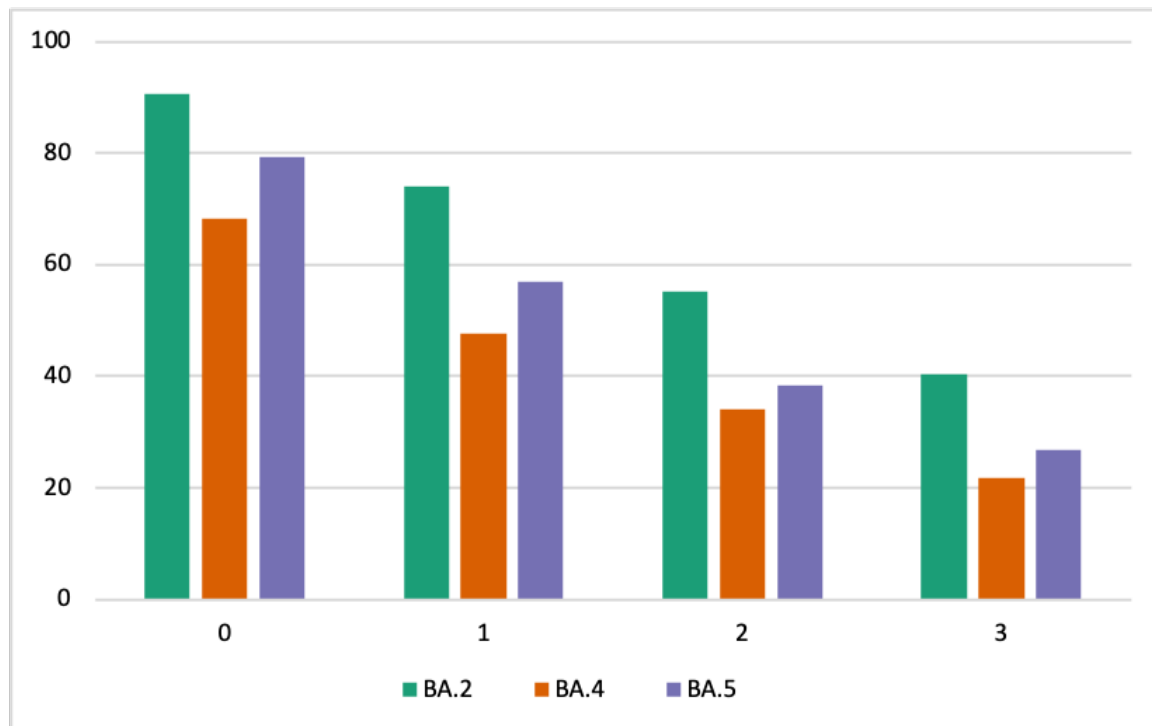
476



477

478 **Figure 2:** Maximum likelihood phylogeny of isolates from each major lineage (BA.2, BA.4,
479 BA.5) found in seven centres. Final branches are coloured according the centre in which they
480 were sequenced.

481



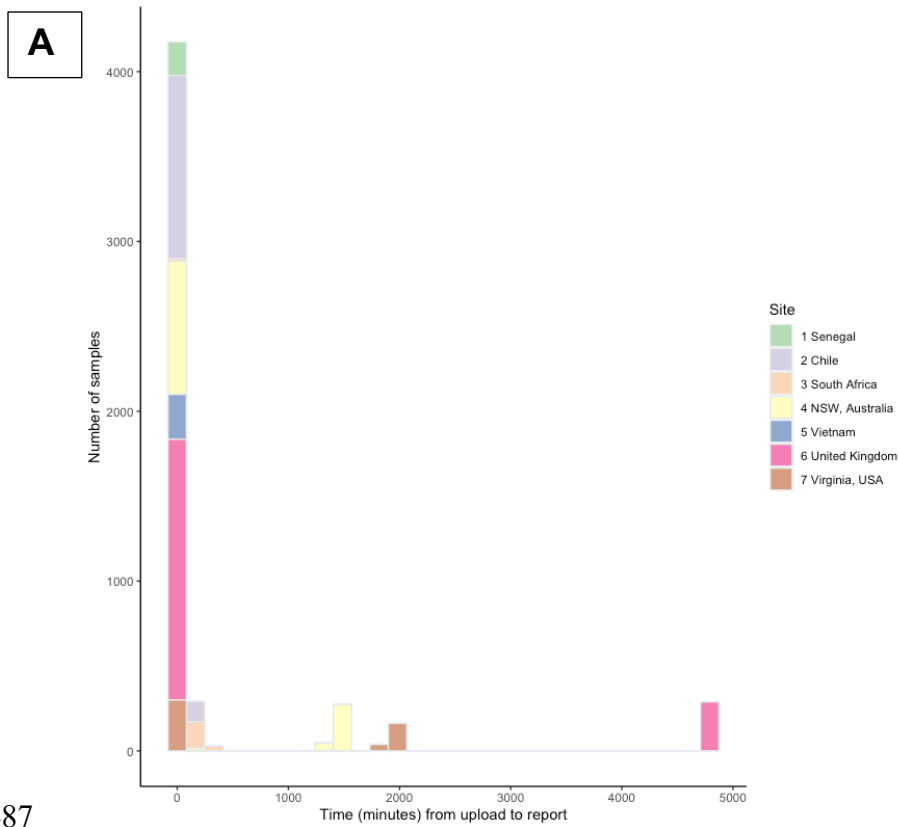
482

483 **Figure 3:** Percentage of related samples which are restricted to the same study site, allowing

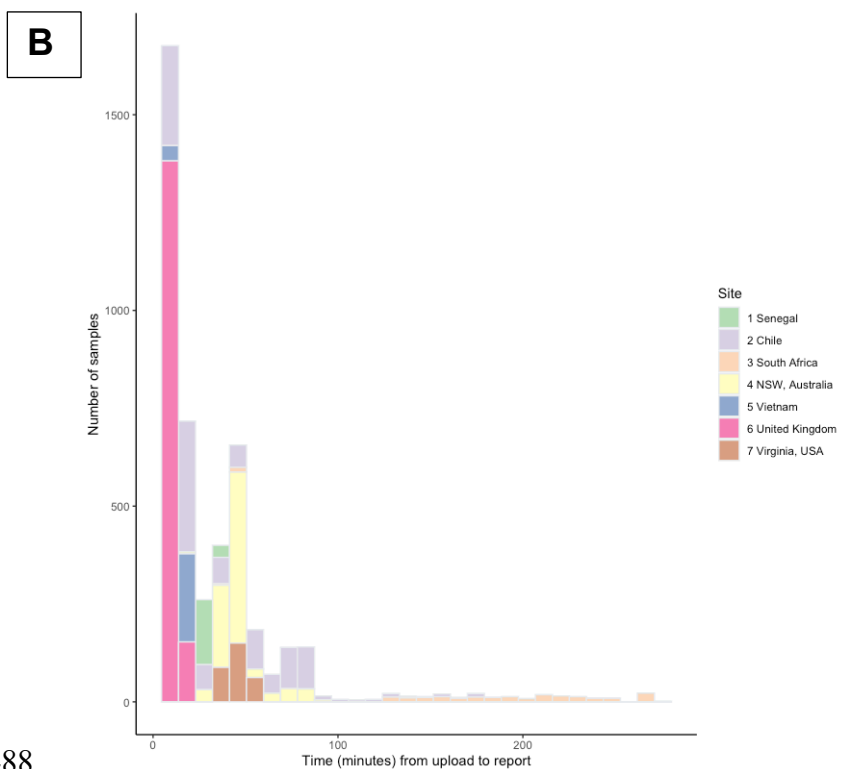
484 for 0, 1, 2 and 3 SNVs as the threshold for relatedness, for 3 major lineages (BA.2, BA.4 and

485 BA.5).

486



487



488

489 **Figure S1:** Distribution of processing time (in minutes) for each sample in GPAS pipeline
490 according to submitting site (a) all samples and (b) samples processed in under 1000 minutes

491