1 Contrasting epidemiology and population genetics of COVID-19 infections defined with

- 2 74 polymorphic loci in SARS-CoV-2 genomes sampled globally.
- 3 Felicia Chan^{1,2}, Ricardo Ataide^{3,4}, Jack S. Richards^{2,3,5} and Charles A. Narh^{2,3,5*}

- ⁵ ¹Central Clinical School, Monash University, Melbourne, Australia.
- ⁶ ²Burnet Institute for Medical Research, VIC 3004, Melbourne, Australia.
- ⁷ ³Department of Medicine, University of Melbourne, Australia.
- ⁴Walter and Eliza Hall Institute, Melbourne, Australia
- ⁵Departmet of Infectious Diseases, Monash University, Melbourne, Australia.
- 10
- 11 *Correspondence: Charles A. Narh
- 12 charles.narh@burnet.edu.au
- 13 Short running title: Genetic tool for monitoring SARS-CoV-2 infections.
- 14 Keywords: COVID-19, SARS-CoV-2, epidemiology, genetics, multilocus-genotypes,
- 15 evolution, linkage, mutation and transmission.

16 Abstract

SARS-CoV-2, the coronavirus causing COVID-19, has infected and killed several millions of 17 people worldwide. Since the first COVID-19 outbreak in December 2019, SARS-CoV-2 has 18 19 evolved with a few genetic variants associated with higher infectivity. We aimed to identify polymorphic loci in SARS-CoV-2 that can be used to define and monitor the viral 20 epidemiology and population genetics in different geographical regions. Between December 21 2019 and September 2020, we sampled 5,959 SARS-CoV-2 genomes. More than 80% of the 22 genomes sampled in Africa, Asia, Europe, North America, Oceania and South America were 23 24 reportedly isolated from clinical infections in older patients, > 20 years. We used the first indexed genome (NC 045512.2) as a reference and constructed multilocus genotypes (MLGs) 25 for each sampled genome based on amino acids detected at 74 polymorphic loci located in 26 ORF1ab, ORF3a, ORF8, matrix (M), nucleocapsid (N) and spike (S) genes. Eight of the 74 27 28 loci were informative in estimating the risk of carrying infections with mutant alleles among different age groups, gender and geographical regions. Four mutant alleles - ORF1ab L₄₇₁₅, S 29 G₆₁₄, and N K₂₀₃ and R₂₀₄ reached 90% prevalence globally, coinciding with peaks in 30 transmission but not COVID-19 severity, from March to August 2020. During this period, the 31 32 MLG genetic diversity was moderate in Asia, Oceania and North America; in contrast to Africa, Europe and South America, where lower genetic diversity and absence of linkage 33 disequilibrium indicated clonal SARS-CoV-2 transmission. Despite close relatedness to Asian 34 MLGs, MLGs in the global population were genetically differentiated by geographic region, 35 suggesting structure in SARS-CoV-2 populations. Our findings demonstrate the utility of the 36 74 loci as a genetic tool to study and monitor SARS-CoV-2 transmission dynamics and 37 evolution, which can inform future control interventions. 38

⁴

39 Introduction

Knowledge of the epidemiology and transmission dynamics of SARS-CoV-2, the causative 40 agent of the COVID-19 pandemic, is seminal to public health control efforts. SARS-CoV-2 41 has more than 110 million confirmed cases – which resulted in 2.5 million deaths – as of March 42 43 2021. Since the first outbreaks in China in December 2019, SARS-CoV-2 transmission hotspots have shifted spatio-temporally from Asia to Europe, followed by North America and 44 South America [1, 2]. Testing and isolation of infected individuals have been integral, as public 45 health interventions, to control the virus' spread. Global reports have shown significant 46 differences in the prevalence, distribution and demographics of COVID-19 cases; however, 47 little is known about how these epidemiological differences relate to infection dynamics. 48

49

50 SARS-CoV-2 is a beta-coronavirus with a positive single-stranded RNA genome of \sim 30 kb. It

shares 80% of its genome with SARS-CoV, which caused the 2003-2004 SARS outbreak [3].

52 A prominent structural feature is the spike glycoprotein (S), which facilitates cell entry and is

53 a target of host immune responses. Other structural proteins include the nucleocapsid (N),

54 envelope (E) and matrix/membrane (M) proteins, which are involved in viral assembly and

55 priming of host immune responses [4]. Nearly half of the genome comprises an opening reading

frame 1ab (*orf1ab*), which encodes 16 non-structural proteins (NSP 1-16) that constitute the

replicase machinery. Notable NSPs include NSP1 (suppresses host immune responses), NSP5
(encodes viral 3C-like protease) and NSP12 (encodes the RNA-dependent RNA polymerase,

i.e., RdRP). Other ORFs, including ORF3a (induces apoptosis in host cells), ORF6, ORF7a,

60 ORF8 (ORF8 mediates immune suppression and evasion) and ORF10 encode accessory

61 proteins that are involved in viral replication and host immune dysregulation [4].

62

Since first being identified, the virus has evolved, with numerous genetic variants being 63 64 associated with higher infectivity. The geographical distribution and probable risk factors (e.g., demographics and clinical factors) for infection with mutant genotypes remain unknown. 65 Comparative genomic analysis of SARS-CoV-2 infections collected globally suggests that the 66 virus is adapting to its human host. A few genetic variants harbouring E484K, N501Y and 67 D614G mutations in the S protein have been associated with higher infectivity than the wild-68 69 type variant, Wuhan NC 045512.2 [5, 6]. Variants with these mutations rose to predominance in many parts of the United Kingdom and South Africa [7, 8]. Other mutations, including 70 orflab P4715L, Orf3a G251V and orf8 L84S, have been associated with higher infectivity and 71 viral density, respectively [9]. Whether these and other unreported mutations are linked, under 72 73 selection and can be used to source-track infections within and between different geographical regions has not been investigated. 74

75

More polymorphic and informative loci, representative of the global SARS-CoV-2 genetic
diversity, are needed to accurately differentiate closely related variants and interrogate the virus
population genetics in different geographical regions [9, 10]. Comparison of SARS-CoV-2
whole genomes identified phylogenetic clusters, defined by Single Nucleotide Polymorphisms
(SNPs) in < 10 codons/loci, that differentiated European and Asian infections [11, 12];

81 however, these loci lacked the needed resolution to differentiate variants circulating globally.

Multilocus genotyping using amino acid changes in SARS-CoV-2 can reduce the complexity in the genomic data and provide informative and virologically relevant data that can provide insights into the transmission dynamics and evolution of variants causing COVID-19. This approach on multiple polymorphic loci can estimate and monitor the genetic diversity of SARS-CoV-2 populations spatiotemporally and in response to control interventions.

87

This study evaluated the epidemiology and population genetics of 5,959 SARS-CoV-2 88 genomes sampled globally to identify risk factors associated with infection with mutant 89 variants and gain insights into how the viral population had evolved geographically eight 90 months into the pandemic. Briefly, we identified 74 polymorphic loci, of which eight loci 91 located in orflab, orf3a, orf8, N and S genes, were considered informative in explaining the 92 risk of infection with mutant variants among different demographics and COVID-19 disease 93 94 phenotypes. Multilocus genotyping at the 74 loci allowed us to genetically differentiate closely related variants circulating globally and gain insights into the viral population genetics in 95

96 different geographical regions.

98 Material and Methods

99 **Data curation and study variables.** The current study sought to investigate the epidemiology and population genetics of SARS-CoV-2 genetic variants causing the COVID-19 pandemic. It 100 was conducted retrospectively by analyzing SARS-CoV-2 whole genomes of ~30 kb. These 101 genomes were isolated from human infections - asymptomatic and symptomatic. From 102 103 December 2019 to September 2020, a total of 5,959 complete genomes with their associated clinical and patient data were retrieved from the Global Initiative on Sharing Avian Influenza 104 Data (GISAID) database [13]. The demographic data included age (grouped into four 105 categories: 0 - 19, 20 - 39, 40 - 59 and ≥ 60 years), gender and the geographical region 106 107 (continent) where the infection was diagnosed. The associated metadata included clinical outcomes - asymptomatic (no symptoms) or symptomatic (mild or severe/critical). Other 108 metadata included specimen type - upper respiratory tract (URT) or lower respiratory tract 109 (LRT) and the technology/chemistry used to sequence the viral genome. 110

Sequence alignments and selection of polymorphic loci. The genomes were aligned to the 111 Wuhan reference strain (NC 045512.2) using minimap version 2.17, and the SNPs, including 112 the corresponding amino acid changes, were called using the Geneious Prime SNP caller [14, 113 15]. Amino acids identical to the reference strain at the investigated loci were considered wild-114 type; else, they were considered mutants. Only polymorphic (≥ 2 alleles, i.e. amino acids, 115 including the wild-type) loci with a minor allele frequency (MAF) of 0.01 were retained and 116 analyzed in this study. These criteria were implemented to ensure unbiased construction of 117 haplotypes (within a gene) and multilocus genotypes, i.e. MLGs (across \geq two genes) [16]. An 118 allele was designated by the amino acid followed by the codon (referred hereafter as locus) 119 120 number. E.g. S D₆₁₄ and S G₆₁₄ indicate glutamine (wild-type allele) and glycine (mutant allele) at locus/codon 614 of the S protein. 121

Population genetics. Genetic diversity indices – the number of haplotypes or MLGs, eMLGs 122 (normalized MLG based on smallest sample) and expected heterozygosity (H_e) were estimated 123 using poppr V2.8.5 [16]. The eMLG was then plotted using the R package Vegan [17], as a 124 125 rarefaction curve to estimate the depth/richness in sampling. The evenness (E5) statistic was used to evaluate whether the haplotypes or eMLGs found within the population were evenly 126 distributed. Its score ranges from 0 (presence of predominant haplotypes or MLGs) to 1 127 (haplotypes or MLGs are evenly distributed). The H_e is a measure of genetic diversity, scoring 128 from 0 (no genetic diversity, i.e., genomes carry the same haplotype or MLG) to 1 (complete 129 diversity, i.e., genomes carry unique haplotypes or MLGs). 130

131

132 To determine linkage disequilibrium (LD), i.e., non-random association of alleles at two or

more loci, the standardized index of association $(\bar{r}d)$ was estimated using *poppr*. The presence

134 of genomes with identical MLGs, i.e. clones within a population, can overestimate the LD [18].

135 To account for this, the LD was clone-corrected using the dataset consisting of unique MLGs.

- 136 To determine whether the LD was 'structured' between specific gene pairs, a pairwise LD was
- 137 performed as described elsewhere [19]. The \bar{r} d score ranges from 0 to 1, with 0 indicating no

LD and 1 indicating complete LD. The statistical significance of the score was supported by a
P-value < 0.05.

140

Genetic differentiation (Nei's GST) among MLG populations within and between continents 141 was estimated using mmod [20]. The GST score ranges from 0 (no genetic differentiation, i.e., 142 populations are similar or have identical MLGs) to 1 (i.e., complete genetic differentiation, i.e., 143 populations are dissimilar or have unique MLGs); values ranging from 0 to 0.09, 0.1 to 0.19 144 and ≥ 0.2 indicate little, moderate and great genetic differentiation, respectively [21]. We then 145 performed a Discriminant Analysis of Principal Components (DAPC) - a multivariate 146 method for identifying genetic clusters of closely related MLGs [22]. Briefly, the global MLG 147 dataset was trained on a K-means algorithm, implemented in *adegenet* [22] to identify the 148 optimum number of genetic clusters within the global population. The DAPC was then 149 performed on the genetic clusters retained during a PCA by maximizing the genetic variance 150 between populations while minimizing the variance within populations [23]. The DAPC 151 assigned population membership probability to each MLG, which was plotted using ggplot2 152 [24]. To visualize the genetic relationships and clonal complexes among the MLGs, the 153 goeBURST FULL MST algorithm implemented in *Phyloviz* V2 was used to construct networks 154 of minimum-spanning trees [25]. 155

156

157 **Statistical analysis**. Statistical analysis was performed in R v3.5.2 [26] and STATA v16 [27]. 158 Proportions were compared using the chi-square or Fisher's exact test. Multiple testing was 159 adjusted using the Holm-Bonferroni method. Logistic regression was performed to determine 160 the association between the study variables and the odds of harbouring a mutant allele. Age, 161 geographical region, and gender were considered possible confounders and adjusted for in the 162 final model. The adjusted odds ratio (OR) was considered statistically significant for all 163 analysis where the P-value was < 0.05.

164

166 **Results and discussion**

Demographics of the study population. The majority of the SARS-CoV-2 whole genomes 167 we sampled (N=5,959) from GISAID between December 2019 and September 2020 were 168 reportedly from Asia (26.5%), Europe (19.9%) and Oceania (27.6%), with $\leq 10\%$ each being 169 from Africa, North America and South America (Table 1). Despite significant differences (P-170 value ≤ 0.005) in the proportion of genomes sampled for the study variables, > 35.0% of the 171 genomes sampled in Africa and Oceania were reportedly from the 20 -39 years age group. In 172 contrast, the majority of genomes (> 31.0%) from the rest of the world were reportedly from 173 the 40 - 59 years age group (Table 1, S1 and S2). These data are consistent with the age 174 175 disparities in COVID-19 cases [28, 29]. Except in Africa, where a significantly higher proportion of genomes were isolated from females (56.2%, P-value = 0.002), most were 176 reportedly isolated from males. 177

		Global	Africa	Asia	Europe	N.America	Oceania	S.America
Factor	Category	5959	601	1579	1188	597	1646	348
Age group [#] (years)	0-19	323 (5.4)	68 (11.3)	123 (7.8)	43 (3.6)	20 (3.6)	58 (3.5)	11 (3.2)
	20 - 39	2054 (34.5)	276 (45.9)	559 (35.4)	273 (22.9)	196 (32.8)	634 (38.5)	116 (33.3)
	40 - 59	1996 (33.5)	185 (30.8)	561 (35.5)	389 (32.7)	217 (36.4)	515 (31.3)	129 (37.1)
	60+	1586 (26.6)	72 (11.9)	336 (21.3)	483 (40.7)	164 (27.5)	439 (26.7)	92 (26.4)
Gender [#]	Females	2664 (44.7)	338 (56.2)	581 (36.8)	581 (48.9)	250 (41.9)	754 (45.8)	160 (45.9)
	Males	3295 (55.3)	263 (43.8)	998 (63.2)	607 (51.1)	347 (58.1)	892 (54.2)	188 (54.0)
Clinical [#] Severity	Asymptomatic	60 (1.5)	0 (0.0)	41 (2.7)	18 (1.9)	0 (0.0)	0 (0.0)	1 (0.4)
	Mild	557 (14.1)	15 (2.5)	134 (8.8)	64 (6.7)	282 (47.2)	2 (12.5)	60 (24.9)
	Severe	3321 (84.33)	586 (97.5)	1356(88.6)	870 (91.4)	315 (52.8)	14 (87.5)	180 (74.7)
	Missing data*	2021	0	48	236	0	1630	107

178 Table 1. Demographics of the study population

denotes number of genomes and the percentage, N (%). * denotes the number of genomes
with missing clinical data and this was not included in the percentage calculations. N. America

- (North America) and S. America (South America).
- 182

The majority of the SARS-CoV-2 genomes were reportedly isolated from throat swabs 183 and were sequenced using Illumina. Half of the genomes we analyzed had the associated 184 data on the specimen type collected for diagnosis or isolating the virus genome. URT 185 specimens constituted 92.0%. Of these, throat swabs were the majority (65.1%) (Figure S2). 186 This was observed for all the study variables (Figure S3-S5) except in South America, where 187 more than 60.3% of the genomes were isolated from nose swabs (Figure S4). Globally, more 188 than 68.0% of the genomes we sampled were reportedly sequenced using Illumina except in 189 Asia and South America, where a higher proportion of genomes were reportedly sequenced 190 using Ion Torrent (39.53%) and Nanopore (43.10%), respectively (Figure S6). 191

192

193 Seventy-four polymorphic loci were selected for multilocus genotyping of SARS-CoV-2

genomes. The majority of mutations in SARS-CoV-2 variants have been considered neutral,

i.e. associated with demographic processes [30, 31]. A few others, considered homoplasic

(recur independently) and adaptive (associated with viral transmissibility and/or 196 pathogenicity), have been detected in clinical infections circulating worldwide [9, 31, 32]. 197 Based on these reports and our filtering criteria of a MAF ≥ 0.01 , 74 polymorphic loci (≥ 2 198 alleles) were used to construct the MLGs (Figure S7). These loci are located within ORF1ab -199 NSP1, NSP2, NSP3, NSP4, NSP5, NSP8, NSP12, NSP13 and NSP14, two accessory proteins 200 - ORF3a and ORF8, and three structural proteins - M, N and S (Table S3). The moderate genetic 201 differentiation (Gst = ~ 0.10) observed for these loci demonstrate their utility as markers for 202 differentiating SARS-CoV-2 variants (Table S4). The orflab gene was the most polymorphic 203 204 loci (Table S3-S4), indicating a mutational hotspot in SARS-CoV-2 [33]. However, its low H_e , ≤ 0.44 , compared to the moderate-to-high $H_e \geq 0.5$, in the orf3a, S and N genes is consistent 205 with previous reports using nucleotide data [34]. Furthermore, our data also suggest that most 206 mutations in the orflab were not under strong selection compared to those in the accessory and 207 structural genes [30, 31, 35]. Indeed, B cell epitopes on the N and S proteins were shown to be 208 highly diversified, allowing the virus to evade host immune responses [36]. 209

210

The spatiotemporal selection of variants carrying mutant alleles - ORF1ab L₄₇₁₅, S G₆₁₄, 211 and N K₂₀₃ and R₂₀₄ was associated with spikes in COVID-19 cases. Of the 149 mutant 212 alleles detected among the 74 loci (Figure S8-S10 and Table S4), eight alleles (Figure 1) were 213 considered putatively adaptive, having been previously associated with infectivity, 214 pathogenesis and host immune dysregulation [4]. The first four - ORF1ab L₄₇₁₅ (located in 215 216 NSP12/RdRP), S G₆₁₄, and N K₂₀₃ and R₂₀₄ rose sharply in frequency from < 6% in February to > 90% in August. This rise was associated with significant peaks in global COVID-19 cases 217 between March and September (Figure 1) [37]. The remaining four alleles, including ORF1ab 218 I₂₆₅ (NSP1), ORF1ab F₃₆₀₆ (NSP5), ORF8 S₈₄ and ORF3a H₅₇, were detected 'transiently' at a 219 lower frequency, ranging from 4 to 35% (Figure 1 and Figure S8-S10). In contrast, most of the 220 remaining 141 alleles circulated at a lower frequency, < 25%, and were not detected throughout 221 the study period and in all continents (Figure S8-S10 and Table S4). These alleles, particularly 222 those in the ORF1ab, have been considered neutral [31]. It is worth noting that the N501Y 223 mutation in the S protein, associated with higher infectivity among UK and South African 224 variants [5, 6], was carried by 0.8% of the genomes we sampled from Australia, Oceania. These 225 genomes were reportedly isolated in June, suggesting that the S Y₅₀₁ allele emerged earlier than 226 previously reported [38]. 227

228

229

The risk of harbouring a mutant allele varied with age, gender, geographical region and 230 COVID-19 phenotype. A critical gap in the surveillance of SARS-CoV-2 infections and 231 understanding COVID-19 clinical severity is the lack of data on the relationship between 232 SARS-CoV-2 variants causing clinical infections and the clinical and demographic factors of 233 infected individuals. To investigate this, we estimated the risk of harbouring a wild-type versus 234 mutant allele at eight informative loci among our study variables. These informative loci -235 ORF1ab (265, 3606 and 4715), ORF3a 57, ORF8 84, N (203 and 204) and S 614 had alleles 236 that were evenly distributed (E.5 score ≥ 0.6) in the global population and recurred throughout 237 the study period (Figure 1 and Table S3). 238

239

Briefly, the risk of harbouring a mutant allele was associated with age, gender and COVID-19 240 phenotype in all the six continents (Figure 2). Compared to patients below 20 years, patients 241 aged 20+ years had a higher risk of harbouring mutant alleles at four loci - N 203 and 204 242 (Asian males), ORF3a 57 (Europe) and ORF8 84 (North America) (Figure 2). This risk among 243 the 20+ years cohort was maintained at the S 614 locus for patients in Oceania but not in North 244 America. Host immune responses to the accessory and structural proteins have been implicated 245 as drivers of SARS-CoV-2 evolution, with immune responses of older patients associated with 246 infections carrying mutant alleles [39]. Interestingly, compared to asymptomatic cases, severe 247 cases in Africa, Asia and North America were more likely to harbour the S D₆₁₄ allele (Figure 248 2). Although the S G_{614} has been associated with enhanced viral transmission, our data indicates 249 that it was not associated with severe COVID-19 as reported elsewhere [40]. 250

251

252 The majority of SARS-CoV-2 MLGs within each continent were detected except in Asia

and Oceania. A major hurdle to tracking the spread and understanding the evolution of SARS-CoV-2 is the complexity and close relatedness of infecting genomes within and between different geographical regions. Therefore, we utilized the 74 polymorphic loci to differentiate the 5,959 genomes by constructing MLGs. We then used the MLG data to measure the richness in our sampling and to interrogate the population genetics of SARS-CoV-2 at the global and continental levels.

The MLGs were moderately distributed, as indicated by the E.5 score ranging from 0.41 for 259 Europe to 0.60 for South America. Asia and Oceania had the highest number of observed (\geq 260 181) and expected (\geq 77) unique MLGs (Table S5). South America had the lowest number, 40 261 (Table 3 and S6). The plateauing of the rarefaction curve for Africa, Europe and the Americas 262 indicated that no new MLGs were detected with further sampling in these populations (Figure 263 3A). In contrast, the steep rise in Asia and Oceania's curve suggested that we had under-264 sampled and therefore did not capture most of the MLGs within these populations. As indicated 265 by the positive correlation (r = 0.94, P-value = 0.017) between sample size and MLG 266 abundance, deep sampling is critical for detecting most SARS-CoV-2 variants causing 267 COVID-19 in affected communities. However, detecting most variants may be challenging in 268 situations where logistics for COVID-19 testing are inadequate. Indeed, < 2% of the genomes 269 we sampled were reportedly isolated from patients with asymptomatic infections. Considering 270 that asymptomatic infections constitute ~80% of all COVID-19 cases and refuel and sustain 271 the virus's transmission worldwide [41], they must be included during surveillance. 272

Linkage structures in SARS-CoV-2 populations vary among geographical regions. We 273 investigated the possibility that specific alleles could be linked and contribute to SARS-CoV-274 2 fitness, including transmissibility and pathogenicity. We first quantified the genetic diversity 275 of SARS-CoV-2 populations in the different geographical regions. The multilocus genetic 276 diversity was lowest in South America ($H_e = 0.15$) and highest in Asia, North America and 277 Oceania ($H_e \ge 0.26$) (Table 2). Our data is consistent with previous reports indicating that 278 SARS-CoV-2 phylogenies in the latter two regions depicted the diversity that existed 279 worldwide as of July 2020 [37, 42, 43]. There was significant genome-wide LD (non-random 280

association among alleles) both at the global and continental levels ($\bar{r}d \ge 0.034$, P-value < 0.001) (Table 2). However, the decay of this LD ($\bar{r}d \le 0.007$, P-value ≥ 0.166) in Africa, Europe

- and South America after repeating the analysis with the unique MLGs was indicative of clonal
- SARS-CoV-2 transmission in these geographical regions. Indeed, multiple outbreaks in Europe
 and South America were due to local transmissions and were largely associated with clusters
- 286 of closely related infections [43, 44].

Interestingly, we observed that the genome-wide LD was driven by specific gene pairs, i.e., 287 'structure'. We detected the strongest LD signal ($\bar{r}d \ge 0.3$, P-value < 0.001) between NSP12 288 and S, NSP1 and ORF3a, NSP4 and ORF8, and NSP4 and N (Figure 3B and Figure S12), 289 consistent with previous reports using nucleotide data [39]. While the former three LD 290 291 structures either decayed or were maintained at the continental level, the NSP12 and S LD structures were consistently prominent in all the geographical regions (Figure S11). A probable 292 driver of the NSP12-S LD structure could be the strong co-selection of the orflab L4715 and S 293 G₆₁₄ alleles, which putatively enhance viral replication and infectivity, respectively [45]. Our 294 findings imply that the evolutionary pressures shaping the virus vary in different geographical 295 regions - with public health interventions, demographic factors, and host immune responses 296 being the most likely key drivers [39, 46, 47]. Nonetheless, it is worth noting that these LD 297 structures may change as the virus continues to evolve, underscoring the need for continuous 298 genomic surveillance. 299

Continent	Ν	MLGs	eMLGs (SE)	E.5	He	<i>r</i> d (P-value)	<i>r</i> d-cc (P-value)
Africa	601	74	56.4 (3.1)	0.57	0.19	0.034 (0.001)	0.007 (0.230)
Asia	1579	185	78.1(5.2)	0.48	0.26	0.080 (0.001)	0.022 (0.001)
Europe	1188	97	53.2 (3.8)	0.41	0.18	0.082 (0.001)	0.009 (0.166)
N.America	597	69	48.6 (3.3)	0.47	0.30	0.219 (0.001)	0.065 (0.001)
Oceania	1646	181	77.5 (5.0)	0.45	0.30	0.082 (0.001)	0.020 (0.001)
S.America	348	40	40.0 (0.0)	0.60	0.15	0.080 (0.001)	0.007 (0.278)
Total	5959	472	95.3 (5.8)	0.34	0.26	0.079 (0.001)	0.019 (0.001)

300 Table 2: Genetic diversity estimates for SAR-COV-2 populations circulating globally.

The number of observed MLGs was normalized by the smallest sample size to obtained the expected MLGs (eMLGs) with the standard error (SE). The standardized index of association $(\bar{r}d)$ was clone-corrected ($\bar{r}d$ -cc) using the unique MLGs dataset. N. America (North America) and S. America (South America).

305

306 SARS-CoV-2 MLGs in Asia and Oceania were representative of the global MLG 307 population. The minimum spanning tree (Figure 4A) was drawn to visualize the network 308 relationships among the 472 unique MLGs. Here, at least 11 major clusters, including eight 309 global clusters (GC1 to 8), were identified. Nearly all clusters comprised a considerable number 310 of Asian MLGs (Figure 4A), indicative of admixture populations, as shown in Figure 4B. The

majority of MLGs in each continent were predicted to have 20-50% geographical assignment 311 to Asia and Oceania (Figure 4B). This data suggests that the majority of SARS-CoV-2 variants 312 in the world are of Asian descent. It also supports contact tracing data that showed that most 313 SARS-CoV-2 cases during the early (March to May) stages of the pandemic were linked to 314 imported cases from Asia [48]. However, a few MLGs were unique to each continent, as 315 indicated by the 40-60% within-continent membership assignments. Prominent among 316 continental clusters were the AC1 and OC1 clusters, detected in Asia and Oceania, respectively 317 (Figure 4A). The ~60% membership assignment of African MLGs to Europe (Figure 4B) 318 suggested that most African infections were more likely to have been imported from Europe, 319 contradicting previous reports of an American source based solely on travel data [49, 50]. This 320 underscores the need to build strong surveillance systems utilizing both travel and genomic 321 data. 322

323

SARS-CoV-2 MLGs were genetically differentiated within and between continents. We 324 detected many mutant alleles restricted to each continent (Table S6), suggesting geographical 325 structuring in SARS-CoV-2 populations. Asia had the highest number of unique (?) alleles, 19 326 in total, including ORF1ab D₁₈₁₂, P₆₇₆, G₂₅₈₆ and I₆₂₉₇, being detected in 7.4, 4.5, 3.1 and 2.8%, 327 respectively of genomes sampled. Among the 16 private alleles detected in Oceania, ORF1ab 328 329 S₂₇₁ (4.4%) was predominant. Europe had 12 private alleles with S G₃₂₀ (10.7%) being predominant, while ORF3a S₂₅₁ (5.4%) was predominant among eight private alleles detected 330 in North America. Nine private alleles, including ORF1ab H₄₀₈₀ (15.6%) and N L₁₈₇ (8.5%), 331 were detected in Africa. Only two private alleles – ORF3a A_{165} (0.6%) and ORF8 A_{62} (0.6%) 332 were detected in South America. 333

Further support for the geographical structuring in the global SARS-CoV-2 populations was 334 obtained from the DAPC analysis in which we identified four main genetic clusters (Figure 335 4C). One cluster consisted of a subset of MLGs from each continent, consistent with the 336 previously described admixture populations in Figure 4B. In contrast, the other three clusters 337 consisted of MLGs, predominantly from Africa, Asia and North America (Figure 4C). A minor 338 proportion of MLGs from Oceania showed clinal differentiation into North America (Figure 339 4B and C), which likely represent closely related SARS-CoV-2 variants that spread between 340 the two regions [51]. A minor proportion of European and South American MLGs showed 341 clinal differentiation into the African cluster (Figure 4C). The Gst estimates also supported 342 evidence of geographical structuring. There was moderate-to-high genetic differentiation (Gst 343 \geq 0.204) among the continents except between Oceania, separately with Asia and Europe, and 344 between Europe and South America, where there was little genetic differentiation (Gst \leq 0.079) 345 (Figure 4D). 346

Spatial connectivity, including cross-border migrations and international travel, is a major conduit for spreading the virus (i.e., resulting in gene flow) among countries. Hence, we expected to see little genetic differentiation among 'regional blocks' within a continent. This hypothesis was valid for Europe, where there was little to moderate genetic differentiation (Gst ≤ 0.181) among regional blocks except between Eastern and Northern Europe (Gst = 0.248) (Figure S12). Interestingly, we detected moderate-to-high genetic differentiation (Gst ≥ 0.111)

among regional blocks in Africa and Asia (Figure S12). This may be a reflection of the fast and strict travel bans that were put into effect early during the spread of the virus in both regions.

356

Conclusion. The disproportionate distribution of SARS-CoV-2 genomes among the young and 357 older age groups in this study was representative of the age distribution of COVID-19 cases 358 reported globally. Throat swabs were the preferred specimen for COVID-19 diagnosis, but the 359 invasive nature involved in sampling may limit its utility for surveillance. In building a robust 360 and efficient surveillance system, access to affordable sequencing and rapid analysis of 361 complex genomic data will be seminal to inform control efforts. The utility of the 74 362 polymorphic loci as markers for studying the epidemiology and population genetics of SARS-363 CoV-2 infections represents significant progress in developing molecular tools for SARS-CoV-364 2 Surveillance. In particular, the eight informatic loci revealed contrasting epidemiology and 365 transmission dynamics of the virus among different demographics, geographical regions and 366 COVID-19 phenotypes. The selection of alleles at these loci and the maintenance of key LD 367 structures in the infecting genomes indicate that the virus is evolving and adapting resulting in 368 enhanced transmissibility to humans. An effect of this evolution was the structuring we 369 observed in the viral population, which allowed us to differentiate closely related variants 370 between different geographical regions genetically. Future studies can include additional loci 371 372 to increase the differentiation among variants within the same geographical region.

373

Limitations. The data presented in this study needs further investigations to draw definite conclusions on the association between age, gender and geographical region and the SARS-CoV-2 variant causing COVID-19. Nearly all the genomes we sampled did not have metadata on where the patient got infected besides the country where the infection was diagnosed and/or genome isolated. Thus, it was difficult to infer a source of infection based solely on the MLG data.

References 380

- Lin, Y.-C., et al., The spatiotemporal estimation of the risk and the international transmission of COVID-19: a global 1. perspective. Scientific Reports, 2020. 10(1): p. 20021.
 - 2. Romero-Severson, E.O., et al., Change in global transmission rates of COVID-19 through May 6 2020. PLOS ONE, 2020. 15(8): p. e0236776.
- 3. Narh, C.A., Genomic Cues From Beta-Coronaviruses and Mammalian Hosts Sheds Light on Probable Origins and Infectivity of SARS-CoV-2 Causing COVID-19. Frontiers in Genetics, 2020. 11(902).
- 4. Zhu, G., et al., Minireview of progress in the structural study of SARS-CoV-2 proteins. Current Research in Microbial Sciences, 2020. 1: p. 53-61.
- 5 Wise, J., Covid-19: New coronavirus variant is identified in UK. BMJ, 2020. 371: p. m4857.
- 6. Zhang, L., et al., SARS-CoV-2 spike-protein D614G mutation increases virion spike density and infectivity. Nature Communications, 2020. 11(1): p. 6013.
- 7. Hodcroft, E.B., et al., Emergence and spread of a SARS-CoV-2 variant through Europe in the summer of 2020. medRxiv, 2020: p. 2020.10.25.20219063.
- 8. CDC. Emerging SARS-CoV-2 Variants. 2021 [cited 2021 01 March]; Available from: https://www.cdc.gov/coronavirus/2019ncov/more/science-and-research/scientific-brief-emerging-variants.html.
- 9. Yao, H., et al., Patient-derived SARS-CoV-2 mutations impact viral replication dynamics and infectivity in vitro and with clinical implications in vivo. Cell Discovery, 2020. 6(1): p. 76.
- 10. Moya, A., E.C. Holmes, and F. González-Candelas, The population genetics and evolutionary epidemiology of RNA viruses. Nature Reviews Microbiology, 2004. 2(4): p. 279-288.
- Forster, P., et al., Phylogenetic network analysis of SARS-CoV-2 genomes. Proc Natl Acad Sci U S A, 2020. 117(17): p. 9241-11. 9243
- 12. Mercatelli, D. and F.M. Giorgi, Geographic and Genomic Distribution of SARS-CoV-2 Mutations. Frontiers in microbiology, 2020. 11: p. 1800-1800.
- 13. Elbe, S. and G. Buckland-Merrett, Data, disease and diplomacy: GISAID's innovative contribution to global health. Global challenges (Hoboken, NJ), 2017. 1(1): p. 33-46.
- 14. Kearse, M., et al., Geneious Basic: An integrated and extendable desktop software platform for the organization and analysis of sequence data. Bioinformatics, 2012. 28(12): p. 1647-1649.
- 15. Li, H., Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics, 2018. 34(18): p. 3094-3100.
- Kamvar, Z.N., J.F. Tabima, and N.J. Grünwald, Poppr: an R package for genetic analysis of populations with clonal, partially 16. clonal, and/or sexual reproduction. PeerJ, 2014. 2: p. e281.
- 17. Oksanen, J., et al., vegan: Community Ecology Package. R package version 2.4-3. Vienna: R Foundation for Statistical Computing.[Google Scholar], 2016.
- 18. Mangin, B., et al., Novel measures of linkage disequilibrium that correct the bias due to population structure and relatedness. Heredity, 2012. 108(3): p. 285-291.
- 19. Garcia, V., et al., Clonal interference can cause wavelet-like oscillations of multilocus linkage disequilibrium. Journal of The Royal Society Interface, 2018. 15(140): p. 20170921.
- 20. Winter, D., MMOD: An R library for the calculation of population differentiation statistics. Molecular Ecology Resources, 2012. 12.
- 21. Hedrick, P.W., A Standardized Genetic Differentiation Measure. Evolution, 2005. 59(8): p. 1633-1638.
- 22. Jombart, T., S. Devillard, and F. Balloux, Discriminant analysis of principal components: a new method for the analysis of genetically structured populations. BMC genetics, 2010. 11(1): p. 94.
- 23. Jombart, T., adegenet: a R package for the multivariate analysis of genetic markers. Bioinformatics, 2008. 24(11): p. 1403-5. Wickham, H., ggplot2: elegant graphics for data analysis. J Stat Softw, 2010. 35(1): p. 65-88. 24.
- 25. Francisco, A.P., et al., PHYLOViZ: phylogenetic inference and data visualization for sequence based typing methods. BMC Bioinformatics, 2012. 13(1): p. 87.
- 26. RCore, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Austria, 2015. 2018. 27. StataCorp, Stata Statistical Software:. 2019.
- Ausubel, J., Populations skew older in some of the countries hit hard by COVID-19. Pew Research Center. https://www. 28.
 - pewresearch. org/fact-tank/2020/04/22/populations-skew-older-in-some-of-the-countries-hit-hard-by-covid-19, 2020.
- 29. Davies, N.G., et al., Age-dependent effects in the transmission and control of COVID-19 epidemics. Nature Medicine, 2020. **26**(8): p. 1205-1211.
- Leung, D.T., et al., Antibody response of patients with severe acute respiratory syndrome (SARS) targets the viral nucleocapsid. 30. J Infect Dis, 2004. 190(2): p. 379-86.
- 31. van Dorp, L., et al., Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. Infection, Genetics and Evolution, 2020. 83: p. 104351.
- Banerjee, S., et al., Mutational spectra of SARS-CoV-2 orflab polyprotein and signature mutations in the United States of 32. America. Journal of Medical Virology. n/a(n/a).
 - 33. Pachetti, M., et al., Emerging SARS-CoV-2 mutation hot spots include a novel RNA-dependent-RNA polymerase variant. Journal of Translational Medicine, 2020. 18(1): p. 179.
- 440 441 442 Morais, I.J., et al., The global population of SARS-CoV-2 is composed of six major subtypes. Scientific Reports, 2020. 10(1): p. 34. 18289.
- 35. van Dorp, L., et al., No evidence for increased transmissibility from recurrent mutations in SARS-CoV-2. Nature Communications, 2020. 11(1): p. 5986.
- Forni, D., et al., Antigenic variation of SARS-CoV-2 in response to immune pressure. Molecular Ecology, 2020. n/a(n/a). 36.
- 37. Russell, T.W., et al., Effect of internationally imported cases on internal spread of COVID-19: a mathematical modelling study. The Lancet Public Health, 2021. 6(1): p. e12-e20.
- 443 444 445 446 447 448 449 Leung, K., et al., Early transmissibility assessment of the N501Y mutant strains of SARS-CoV-2 in the United Kingdom, October 38. to November 2020. Euro surveillance : bulletin Europeen sur les maladies transmissibles = European communicable disease bulletin, 2021. 26(1): p. 2002106.
- 450 451 Zeng, H.-L., et al., Global analysis of more than 50,000 SARS-CoV-2 genomes reveals epistasis between eight viral genes. 39. Proceedings of the National Academy of Sciences, 2020. 117(49): p. 31519-31526.

- 40. Korber, B., et al., Tracking Changes in SARS-CoV-2 Spike: Evidence that D614G Increases Infectivity of the COVID-19 Virus. Cell, 2020. 182(4): p. 812-827.e19.
- 41. McArthur, L., et al., Review of Burden, Clinical Definitions, and Management of COVID-19 Cases. The American journal of tropical medicine and hygiene, 2020. 103(2): p. 625-638.

42. Geoghegan, J.L., et al., Genomic epidemiology reveals transmission patterns and dynamics of SARS-CoV-2 in Aotearoa New Zealand. Nature Communications, 2020. 11(1): p. 6351.

Zhao, Z., et al., Genetic grouping of SARS-CoV-2 coronavirus sequences using informative subtype markers for pandemic spread 43. visualization. PLOS Computational Biology, 2020. 16(9): p. e1008269.

44. Poterico, J.A. and O. Mestanza, Genetic variants and source of introduction of SARS-CoV-2 in South America. J Med Virol, 2020. 92(10): p. 2139-2145.

45. Caccuri, F., et al., A persistently replicating SARS-CoV-2 variant derived from an asymptomatic individual. Journal of Translational Medicine, 2020. 18(1): p. 362.

46. Islam, M.R., et al., Genome-wide analysis of SARS-CoV-2 virus strains circulating worldwide implicates heterogeneity. Scientific Reports, 2020. 10(1): p. 14004.

- 47. Wang, R., et al., Host Immune Response Driving SARS-CoV-2 Evolution. Viruses, 2020. 12(10).
- du Plessis, L., et al., Establishment and lineage dynamics of the SARS-CoV-2 epidemic in the UK. Science, 2021: p. eabf2946. 48. 49. Sun, H., et al., Importations of COVID-19 into African countries and risk of onward spread. BMC Infectious Diseases, 2020. 20(1): p. 598.
 - 50. Tegally, H., et al., Sixteen novel lineages of SARS-CoV-2 in South Africa. Nature Medicine, 2021.
- 51. Seemann, T., et al., Tracking the COVID-19 pandemic in Australia using genomics. Nature Communications, 2020. 11(1): p. 472 4376.
- 473

474

476 Conflict of Interest

477 None to declare.

478 Author Contributions

479 CAN, RA and JSR conceived and designed the study. CAN and FC acquired the data and performed
480 the analysis. CAN and FC drafted the manuscript. RA and JSR critically revised the manuscript. All
481 authors read and approved the manuscript for publication.

482 Funding

This work was partly supported by the National Health and Medical Research Council (NHMRC) of
Australia [APP1161076 to JSR]. Burnet Institute received funding from the NHMRC Independent
Research Institutes Infrastructure Support Scheme, and the Victorian State Government Operational
Infrastructure Support Scheme. The funders had no role in study design, data collection and analysis,
decision to publish, or preparation of the manuscript.

- 488 Acknowledgements. We are grateful to the authors from the laboratories responsible for obtaining the
- 489 specimens and the submitting laboratories where genetic sequence data were generated and shared via
- 490 the GISAID Initiative, on which this research is based.

1 Figure legends

- 2
- 3



4

5 Figure 1. Spatiotemporal prevalence of eight SARS-CoV-2 mutant alleles and COVID-19 new

6 cases. The prevalence data for the alleles and newly confirmed cases (WHO report 2020) are

7 reported for December 2019 to September 2020. Four mutant alleles - ORF1ab L₄₇₁₅, S G₆₁₄,

- 8 and N K_{203} and R_{204} were associated with spikes in COVID-19 cases.
- 9



10

11 Figure 2. Eight informatic loci in SARS-CoV-2 associated with carriage of mutant alleles

12 among different age groups, gender, geographical regions and COVID-19 clinical phenotype.

13 The adjusted odds ratio (OR) with the P-value are shown with orange, blue and grey matrices

indicating OR > 1 and P-value < 0.05, OR < 1 and P-value < 0.05 and 1 < OR > 1 and P-value

15 > 0.05, respectively. The OR was not estimated for study variables with < 5 samples.



Figure 3. Richness in MLG sampling and the standardized index of association ($\overline{r}d$) among 18 genes in the global SAR-CoV-2 population sampled from December 2019 to September 2020. 19 A. Rarefaction curve of the MLGs sampled in each geographical region. The plateau in the 20 21 curves indicated no new MLGs were detected with further sampling in Africa, Europe, the American SARS-CoV-2 populations. B. Pairwise LD estimates among genes in the 5,959 22 genomes sampled globally. The $\bar{r}d$ ranges from 0 (no LD) to 1 (complete LD). The values in 23 the coloured heatmap indicate the P-value associated with the pairwise $\bar{r}d$ estimates. The 24 25 strongest LD signal ($\bar{r}d > 0.3$, P-value < 0.001) was detected between NSP12 and S, NSP1 and ORF3a, NSP4 and ORF8, and NSP4 and N. 26

27







Figure 4. Genetic relatedness and differentiation among SAR-CoV-2 MLGs from different 31 32 geographical regions. A. Relatedness among the 472 unique MLGs in the global population. Eleven clusters including eight global clusters (GC1-8) and three continental clusters – Asia 33 (AC1-2) and Oceania (OC1) were detected. Nearly all clusters contained a considerable 34 number of Asian MLGs. B. Population membership assignment of each MLG. Admixture 35 populations were prominent in all geographical regions. C. DAPC analysis identified one 36 global cluster (MLGs from all regions, central axis of PCA plot) and three continental clusters 37 - Africa, Asia and North America. D. Moderate to high genetic differentiation (Nei's Gst) was 38 observed among MLGs from different continents. 39

- 40
- 41