1    Rapid selection of P323L in the SARS-CoV-2 polymerase (NSP12) in humans and non-human

2    primate models and confers a large plaque phenotype

3    Xiaofeng Dong[1†], Hannah Goldswain[1†], Rebekah Penrice-Randal[1†], Ghada T. Shawli[1†], Tessa

4    Prince[1,2†], Maia Kavanagh Williamson[3†,], Nadine Randle[1], Benjamin Jones[1], Francisco J

5    Salguero[4], Julia A. Tree[4], Yper Hall[4], Catherine Hartley[1], Maximilian Erdmann[3], James Bazire[3],

6    Tuksin Jearanaiwitayakul[3,5] ISARIC4C investigators, Malcolm G. Semple[1,2,6], Peter J. M.

7    Openshaw[7], J. Kenneth Baille[8], Stevan R. Emmett[9,10], Paul Digard[8], David A. Matthews[3], Lance

8    Turtle[1,2], Alistair Darby[1], Andrew D. Davidson[3], Miles W. Carroll[2,4,11] and Julian A. Hiscox[1,2,12*].

9    [1]Institute of Infection, Veterinary and Ecological Sciences, University of Liverpool, UK.

10   [2]NIHR Health Protection Unit in Emerging and Zoonotic Infections, Liverpool, UK.

11   [3]School of Cellular and Molecular Medicine, University of Bristol, UK.

12   [4]UK Health Security Agency, Porton Down, UK.

13   [5]Department of Microbiology, Mahidol University, Thailand.

14   [6]Department of Respiratory Medicine, Alder Hey Children's Hospital, Liverpool, UK.

15   [7]National Heart and Lung Institute, Imperial College London, UK.

16   [8]The Roslin Institute, University of Edinburgh, UK.

17   [9]Royal United Hospitals Bath NHS Foundation Trust, UK.

18   [10]Bristol Medical School University of Bristol, UK.

19   [11]Nuffield Department of Medicine, University of Oxford, UK.

20    [12]A*STAR Infectious Diseases Laboratories (A*STAR ID Labs), Agency for Science, Technology

21    and Research (A*STAR), Singapore.

22    [†]These authors contributed equally.

23    Correspondence: julian.hiscox@liverpool.ac.uk

24　**Abstract**

25　The mutational landscape of SARS-CoV-2 varies at both the dominant viral genome sequence

26　and minor genomic variant population. An early change associated with transmissibility was the

27　D614G substitution in the spike protein. This appeared to be accompanied by a P323L

28　substitution in the viral polymerase (NSP12), but this latter change was not under strong

29　selective pressure. Investigation of P323L/D614G changes in the human population showed

30　rapid emergence during the containment phase and early surge phase of wave 1 in the UK. This

31　rapid substitution was from minor genomic variants to become part of the dominant viral

32　genome sequence. A rapid emergence of 323L but not 614G was observed in a non-human

33　primate model of COVID-19 using a starting virus with P323 and D614 in the dominant genome

34　sequence and 323L and 614G in the minor variant population. In cell culture, a recombinant

35　virus with 323L in NSP12 had a larger plaque size than the same recombinant virus with P323.

36　These data suggest that it may be possible to predict the emergence of a new variant based on

37　tracking the distribution and frequency of minor variant genomes at a population level, rather

38　than just focusing on providing information on the dominant viral genome sequence e.g.,

39　consensus level reporting. The ability to predict an emerging variant of SARS-CoV-2 in the global

40　landscape may aid in the evaluation of medical countermeasures and non-pharmaceutical

41　interventions.

## Introduction

There are many distinct lineages of SARS-CoV-2 currently circulating worldwide and some that have become extinct [1]. Sequence data show that that the genome of SARS-CoV-2 is changing as the pandemic continues. Replication and transcription of the SARS-CoV-2 genome directly drives three types of genetic change in the virus. The first is recombination, and this is a natural consequence of the way in which the virus synthesizes its subgenomic messenger RNAs (sgmRNAs). This may account for insertions and deletions, for example observed in and around the furin cleavage site in the spike glycoprotein [2] and other genes [3]. The second driver of genetic change is the continual accruing of point mutations. These changes may confer advantages in transmission, such as the A23402G, encoding the D614G substitution in the spike protein [4], which has come to predominate in global SARS-CoV-2 sequences since the start of the outbreak [5]. Such point mutations may be driven by errors during RNA synthesis by the viral encoded RNA dependent RNA polymerase (NSP12) and larger replication complex and/or by host mediated processes [6,7]. The third mechanism is the potential generation and selection of new transcription regulatory signals (TRSs) and the synthesis of new viral sgmRNAs and proteins [8]. Promiscuous recombination and mutation in coronaviruses may allow these viruses to overcome selection pressures, transit population bottlenecks and result in the emergence of new variants [9,10].

This variation exists in individual humans/animals infected with SARS-CoV-2, where there will be a dominant viral genome sequence(s) with minor genomic variants [10]. These latter genomes will have both synonymous (non-coding) and non-synonymous (coding) variations (changes) around the dominant viral genome sequence. These variations may be selected for and become

64 the dominant viral genome sequence when the virus enters a new host, as has been

65 demonstrated with the adaptation of Ebola virus in a guinea pig model of infection [11].

66 Alternatively, the variation may exist at a minor variant level but nevertheless impact upon

67 virus biology, for example with the Ebola virus RNA dependent RNA polymerase (L protein) and

68 the relationship with overall viral load in patients with Ebola virus disease [12].

69 Since the start of the COVID-19 pandemic different dominant viral genome sequences and non-

70 synonymous changes appear to rise and fall in the SARS-CoV-2 global sequences [1]. The D614G

71 spike protein variant of SARS-CoV-2 was first observed in February 2020 and by May 2020

72 approximately 80% of viruses sequenced contained this substitution. The major clade

73 containing D614G (Pango lineages B.1 and sub-lineages) contained potentially linked

74 substitutions, including C14407U in NSP12 that confers a P323L substitution. However, some

75 lineages, such as A.19 and A.2.4, gained D614G in the spike protein but not P323L in NSP12 [13].
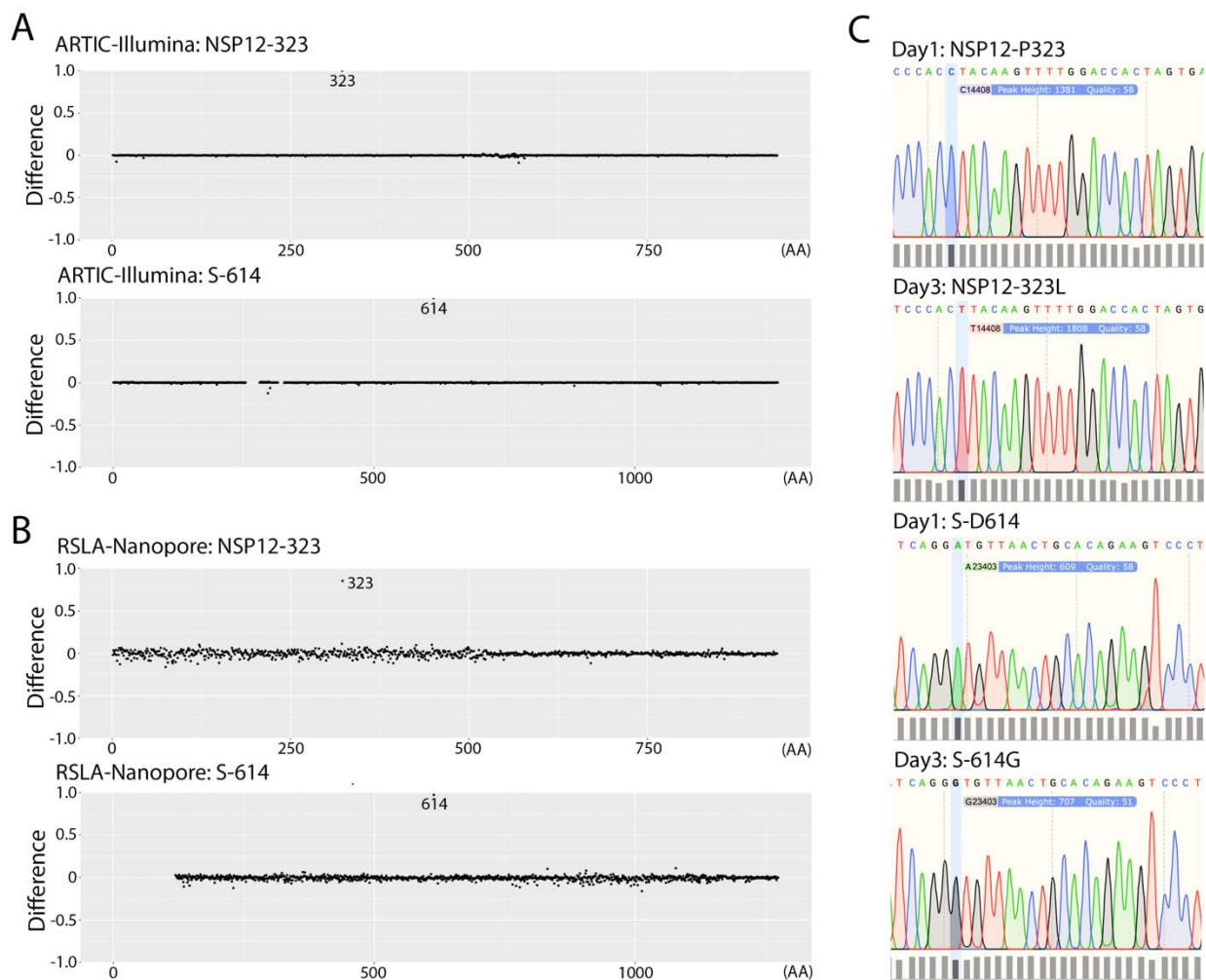
76 Therefore, whether P323L in NSP12 conferred a fitness advantage and was subject to selection

77 pressure is unknown. To investigate the within host selection pressure for the P323L variant,

78 sequential samples from patients with COVID-19 prior to and during the D614G/P323L change

79 in the UK were sequenced to study both the dominant viral genome sequence and minor

80 variant genomes. Additionally, a lineage B SARS-CoV-2 with 323L and 614G in the minor variant

81 population was used to infect two non-human primate models [13], cynomolgus (*Macaca*

82 *fascicularis*) and rhesus (*Macaca mulatta*) macaques. Longitudinal sampling indicated that 323L

83 became part of the dominant viral genome sequence, but not 614G. Reverse genetics analysis

84 of P323L in the background of a 614G virus indicated that the 323L variant grew with a larger

85 plaque phenotype. Overall, this change provided an additive advantage to D614G in the spike

86    protein. In the wider context the work indicated that an emerging dominant sequence could be

87    predicted by analysis of minor variant genomes.

88   **RESULTS**

89   **Identification of a P323L substitution in NSP12 in the same human patient.** To identify

90   whether the P323L substitution occurred rapidly in NSP12, nasopharyngeal swabs were

91   identified in the ISARIC-4C biobank that were obtained from patients infected with lineage B

92   SARS-CoV-2 prior to the major shift from P323 to 323L and D614 to 614G. Samples were further

93   down selected based on clinical information providing a dates of symptom onset, first sample

94   and subsequent longitudinal samples. This provided samples from a total of 472

95   nasopharyngeal swabs. RNA was isolated from the swabs and used as templates for the

96   amplification of SARS-CoV-2 genome and sgmRNAs using both short (ARTIC-Illumina) and

97   longer-read length (Rapid Sequencing Long Amplicons-Nanopore, RSLA-Nanopore) [14,15].

98   Longitudinal samples from 12 patients had sufficient read depth to call a consensus for the

99   dominant viral genome sequence in each sample and to derive information on the frequency of

100  minor genomic variants, focusing on codon 323 in NSP12 and 614 in the spike protein. In one

101  patient, who was admitted to the intensive care unit at the Royal Liverpool Hospital, both

102  sequencing approaches indicated that the P323L and D614G substitution occurred in the SARS-

103  CoV-2 genome between the 1st sample and 2nd samples taken two days apart (Figure 1A and 1B,

104  respectively). To independently confirm this observation, the source RNA was Sanger

105  sequenced with primers to generate longer amplicons around the potential substitution sites.

106  The data validated that for NSP12 the codon encoding the amino acid at position 323 changed

107  from CCU (encoding P) to CUU (encoding L) (Figure 1C). For the spike protein, the codon

108  encoding the amino acid at position 614 changed from GAU (encoding D) to GGU (encoding G)

109  (Figure 1C). Therefore, the data suggested that both P323L and D614G were rapidly selected in

110     the patient over a two-to-three-day period. Another possibility is that the patient was infected

111     with a P323/D614 variant and subsequently became infected with a 323L/614G variant through

112     nosocomial infection in the hospital setting. However, we consider this possibility unlikely; as

113     this patient was one of the first cases admitted to the intensive care unit of Liverpool University

114     Hospitals, when there were relatively few other patients present in the hospital at that period

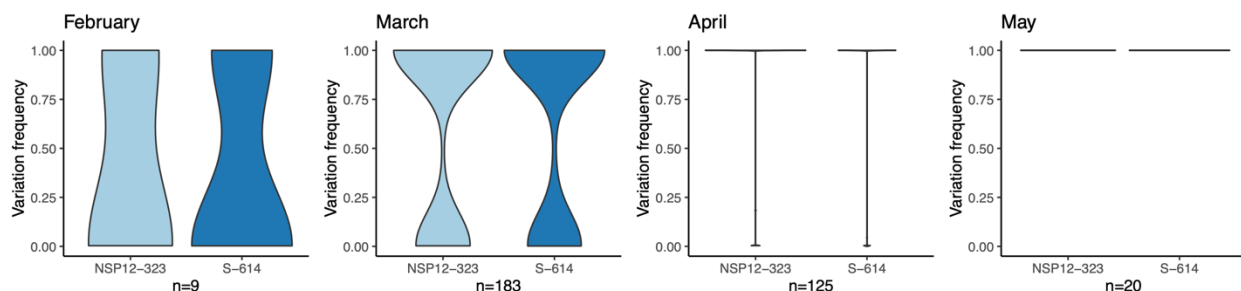115     of the containment phase.

116



117     *Figure 1. Sequence analysis and amino acid substitution in NSP12 (P323L) and the spike protein*

118     *(D614G) between an initial sample and one taken two days later in a single patient. Three*

119 *different sequencing approaches were used: (A) an ARTIC -Illumina approach and (B) an RSLA-*

120 *Nanopore approach. Individual dots represent a codon position on either NSP12 or the spike*

121 *protein compared to the Wuhan reference sequence. The difference between the sampling days*

122 *is indicated by a positive difference indicating divergence of the day 3 sequence away from the*

123 *Wuhan-Hu-1 complete genome reference sequence (NC_045512), and a negative difference*

124 *indicating divergence of the first sample taken towards the Wuhan reference sequence. In both*

125 *cases considering the ratio of a particular position for dominant viral genome sequence versus*

126 *minor variant. (C) Sanger sequence analysis of the amplicons used to investigate the dominant*

127 *viral genome sequence around the sites within NSP12 (codon 323) and spike protein (codon 614)*

128 *that changed between the first and third days of sampling in a patient hospitalized with COVID-*

129 *19.*

130

131 The distribution of P323L and D614G at the minor genomic variant level was evaluated in the

132 human population between January 2020 and June 2020, when these substitutions became

133 part of the dominant viral genome sequence. SARS-CoV-2 was sequenced from nasopharyngeal

134 swabs sampled from 522 patients over that time and usable data obtained from 377 (Figure 2).



135

136 *Figure 2. Analysis of the ratio of P323L (light blue) and D614G (blue) at a dominant viral genome*

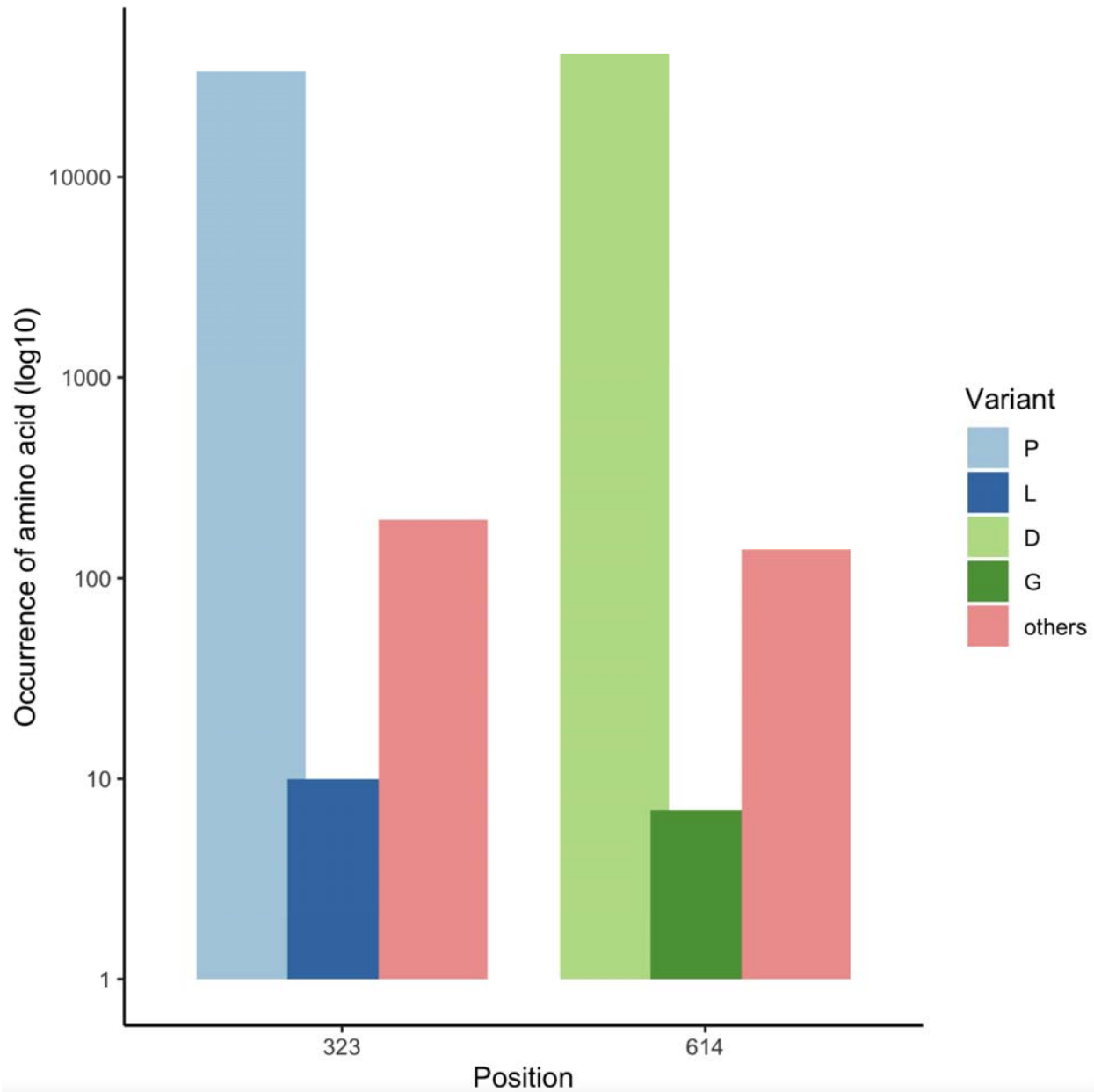137 *sequence and minor variant genomes in 377 patients between February 2020 and May 2020 in*

9

138 *the UK. SARS-CoV-2 sequence was obtained from nasopharyngeal swabs from 377 hospitalized*

139 *patients. The width of the violin plot indicates the number of samples/patients with the*

140 *frequency on the y-axis. The data shows the transition from P323L and D614G over time in the*

141 *minor variant genomes, such that by April 2020 in the UK, the 323L and 614G substitutions were*

142 *part of the dominant viral genome sequence and by May 2020, there was no evidence of P323*

143 *and D614.*

144

145 The data (Figure 2) indicated that there was increasing prevalance from P323 to 323L and D614

146 to 614G in the February to March sampling period. For both February and March 2020, patients

147 had mixed populations of P323L and D614G. However, for the patients sampled in April and

148 May 2020 the dominant viral genome sequence in each patient had 323L and 614G, suggesting

149 either strong selection pressure and/or multiple founder effects.

150 **Longitudinal analysis of variation in non-human primates and cell culture**

151 To investigate whether the P323L substitution was driven by strong selection pressure,

152 nasopharyngeal swabs were taken longitudinally from cynomolgus and rhesus macaques (12

153 animals of each species, a mix of males and females) that had been infected with an isolate of

154 SARS-CoV-2 prior to the P323L and D614G changes; SARS-CoV-2 Victoria/01/202040, that had

155 been sampled on the 24[th] January 2020 [16]. The isolate had been passaged three times in cell

156 culture to generate stock virus prior to infection of the cynomolgus and rhesus macaques.

157 Sequencing of the stock virus indicated a very low proportion of NSP12 323L and spike 614G
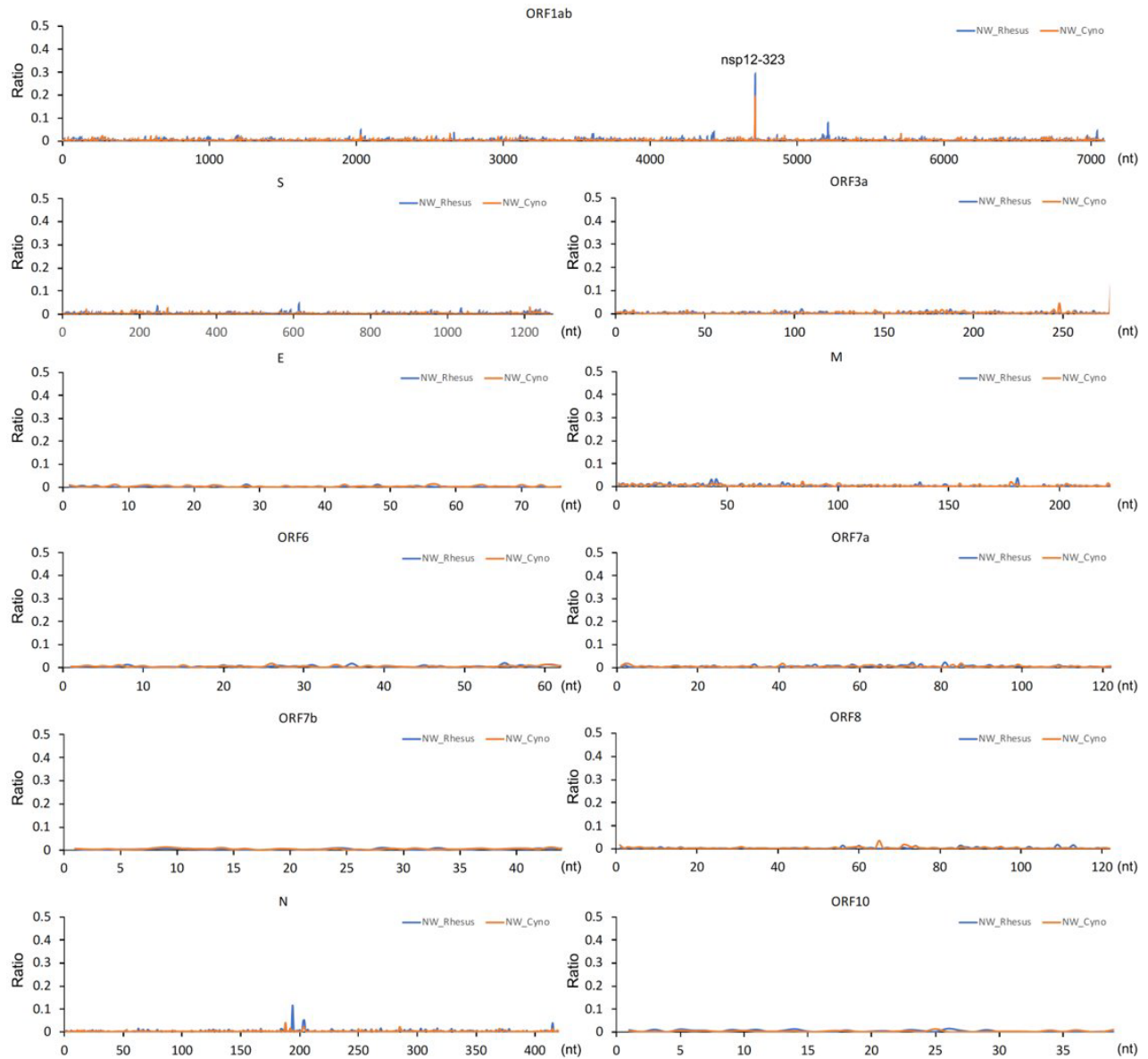
158 (Figure 3).

159

160

*Figure 3. Histogram showing the amino acid coverage at position 323 in NSP12 and 614 in the spike protein in the SARS-CoV-2 Victoria/01/202040 stock as determined by ARTIC-Illumina sequencing. Site coverage is shown on the y-axis. The proportion of amino acids mapped are shown, light blue or light green is the P323 or D614 at the 323 positions in NSP12 and 614 in the spike protein, respectively. The proportion of the L or G in NSP12 and the spike protein,*

12

166    *respectively, is indicated in dark blue and dark green, respectively. The frequency of other amino*

167    *acids at those positions is indicated in pink. We note that data were obtained through an ARTIC-*

168    *Illumina based approach and as such PCR duplicates could not be removed. This may impact on*

169    *the reported ratios.*

170

171    Nasal washes were taken daily from each animal during infection [13]. RNA was purified and

172    sequenced using two independent approaches, shotgun sequencing on an Illumina platform

173    and via ARTIC-Illumina with the latter for specifically sequencing SARS-CoV-2 RNA. Dominant

174    viral genome sequence and minor genomic variants were determined for SARS-CoV-2 for each

175    sample in which genome coverage could be obtained. To obtain a global overview and identify

176    whether there were any hot spots for minor genomic variants, these were plotted as an

177    average over the course of the infections in the non-human primates (NHPs) (Figure 4). The

178    data indicated that minor genomic variants occurred throughout the genome, but the greatest

179    variation occurred at position 14,408 in the orf1ab region, which resulted in a C to U change.

180    This resulted in a non-synonymous change in NSP12 with the substitution of P323L (amino acid

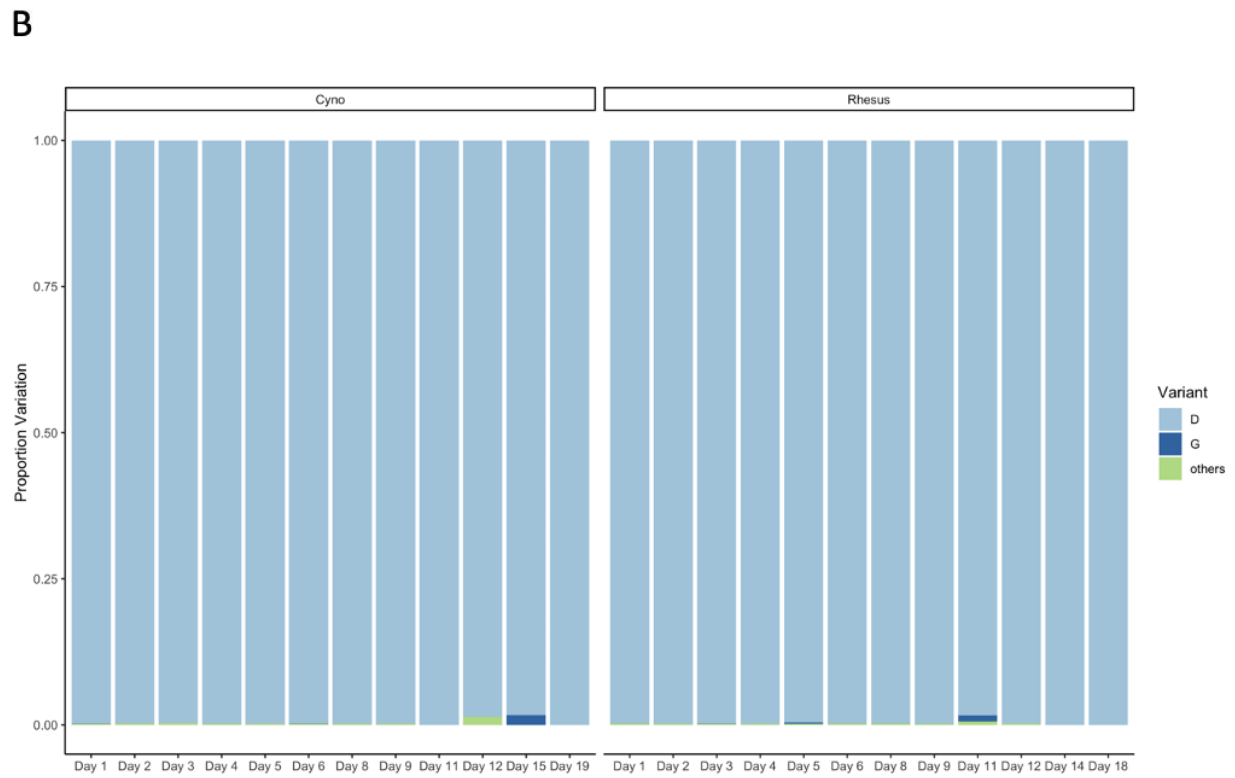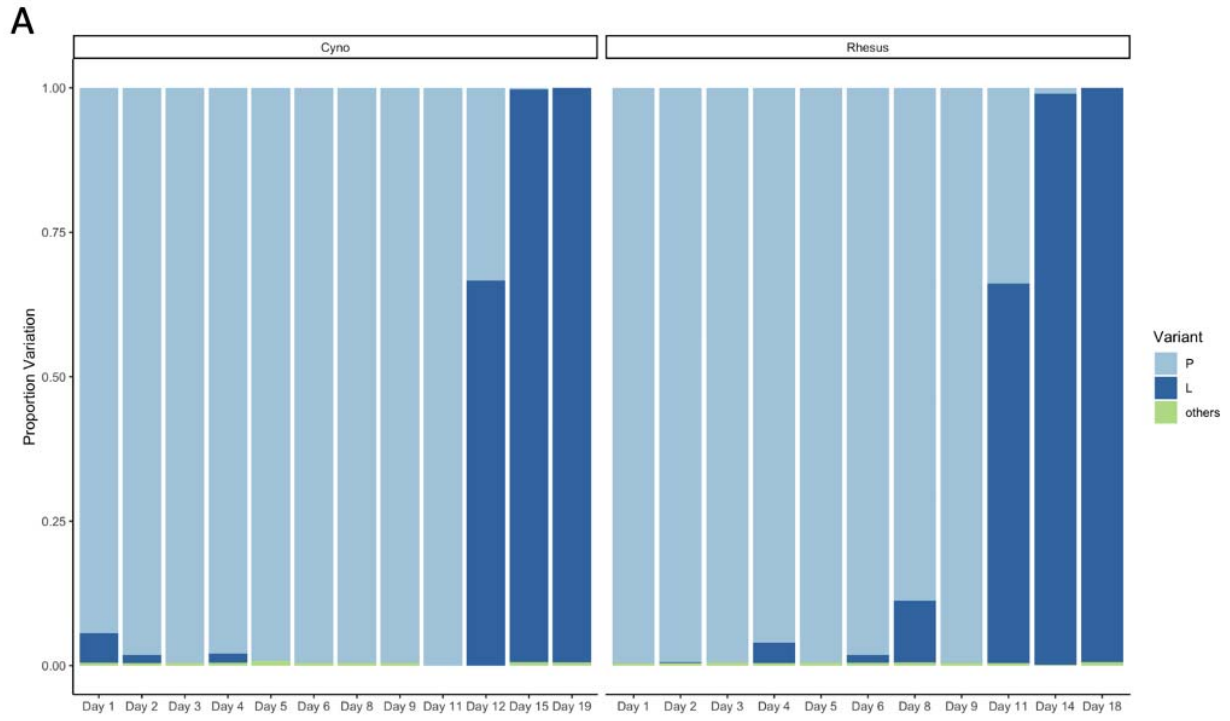181    position 4715 with respect to the ORF1AB polyprotein).

*Figure 4. Analysis of minor variant genomes in cynomolgus and rhesus macaques infected with the SARS-CoV-2 Victoria/01/202040 isolate using data from shotgun Illumina RNA sequencing of nasal washes (NW). Data presented as a global average over the course of the infection from sequencing SARS-CoV-2 from longitudinal samples. Each SARS-CoV-2 open reading frame is indicated above the appropriate panel. The major difference was at position 323 in NSP12.*

189    To determine how rapidly these mutations were selected in the individual animals, sequences

190    from longitudinal samples were analyzed (Figure 5, showing ARTIC-Illumina data)

191    (Supplementary Figure 1, showing both ARTIC-Illumina and ARTIC-Nanopore approaches and

192    coverage). The sequencing data, using the two different approaches, showed that the P323L

193    mutation was already present as a minor genomic variant (at higher levels than the inoculum)

194    by Day 1 in some animals, as well as the presence of other minor genomic variants at this

195    position. However, as infection progressed the frequency of the 323L minor genomic variant

196    increased and became part of the dominant viral genome sequence by the end point of

197    infection. This was the general pattern for all individual animals whether cynomolgus or rhesus
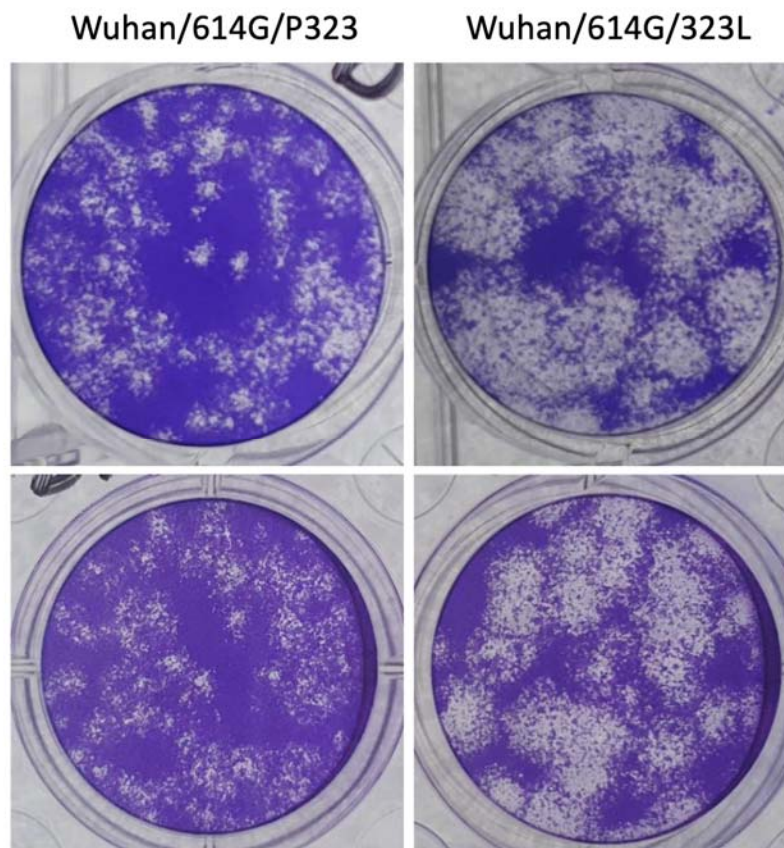
198    macaque.

199

200

201 *Figure 5. Analysis of NSP12 position 323 (A) and the spike protein position 614 (B) in SARS-CoV-2*

202 *from nasopharyngeal swabs taken longitudinally from infected cynomolgus and rhesus*

203 *macaques. Data in this figure is from the ARTIC-Illumina approach to specifically amplify SARS-*

204 *CoV-2 RNA. The day post infection is shown for the animals. In some cases, where there was*

205 *more than one animal for each day, or usable sequence was obtained, the average value was*

206 *calculated. For each position of interest either the P (for position 323 in NSP12) or D (for position*

207 *614 in the spike protein) is shown in light blue, and the substitution of L or G, shown in dark*

208 *blue, respectively. Green indicates other substitutions at that position. The left-hand y-axis*

209 *indicates the % variation at the indicated position). The % variation was only shown for these*

210 *sites with coverage > 5.*

211

212 **The P323L substitution in NSP12 confers a growth advantage in the context of a recombinant**

213 **virus with 614G in the spike protein**

214 Previous data indicated that Victoria/01/202040 grew with a small plaque phenotype and lower

215 titer compared to more contemporary variants including Variants of Concern (VOCs), that grew

216 to higher titres with larger or mixed plaque morphologies [17]. The later virus isolates contained

217 the P323L and D614G substitutions in NSP12 and the spike protein, respectively, as the

218 dominant viral genome sequence, as well as other changes. To investigate whether the 323L

219 substitution conferred an advantage over and above the 614G change in the spike protein, two

220 recombinant viruses were created that were based on the 614G background, one with P323

221 (Wuhan/614G/P323) and the other with 323L (Wuhan/614G/323L) in NSP12. Growth of these
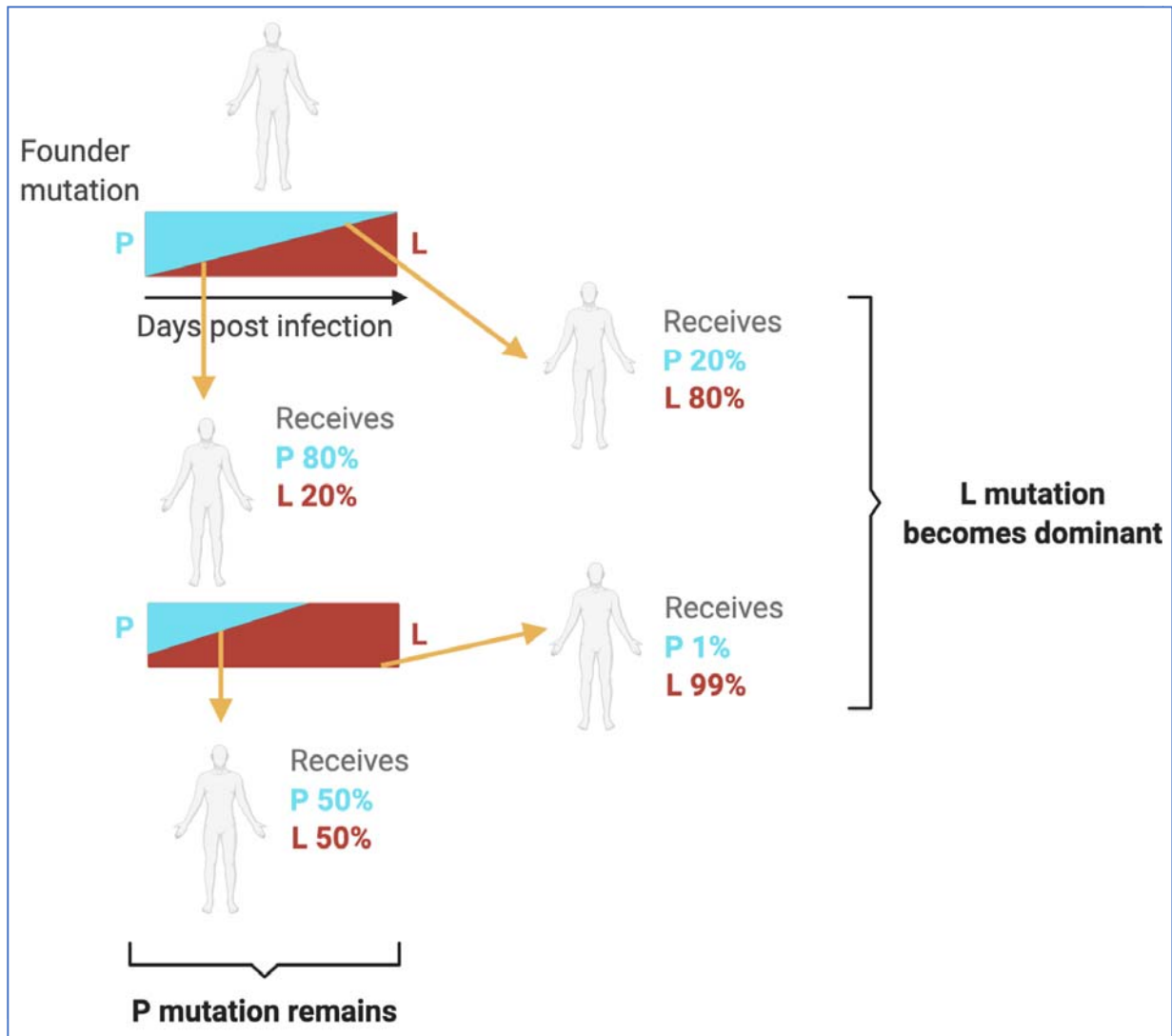
17

222    two recombinant viruses were compared in cell culture by examining plaque morphology. The

223    data indicated that Wuhan/614G/323L had a large plaque phenotype whereas

224    Wuhan/614G/P323 had a small plaque phenotype (Figure 6), suggesting that the 323L

225    substitution conferred a growth advantage.

226



227

*Figure 6. Representative images of plaques formed by two recombinant viruses that have the*

*Wuhan-Hu-1 background (NC_045512) and an engineered D614G substitution in the spike*

*protein and differed at position 323 in NSP12 with either a P or L, these were termed*

*Wuhan/614G/P323 and Wuhan/614G/323L, respectively.*

232    **Maintenance of variation at position 323 in NSP12 in the population**

233    Based on the experimental data presented in this study, we propose a model where the

234    emergence and distribution of minor variant genomes and dominant viral genome sequence for

235    SARS-CoV-2 is dependent on selection pressure and time post-infection at which a virus

236    population is transmitted onwards to another individual (Figure 7).
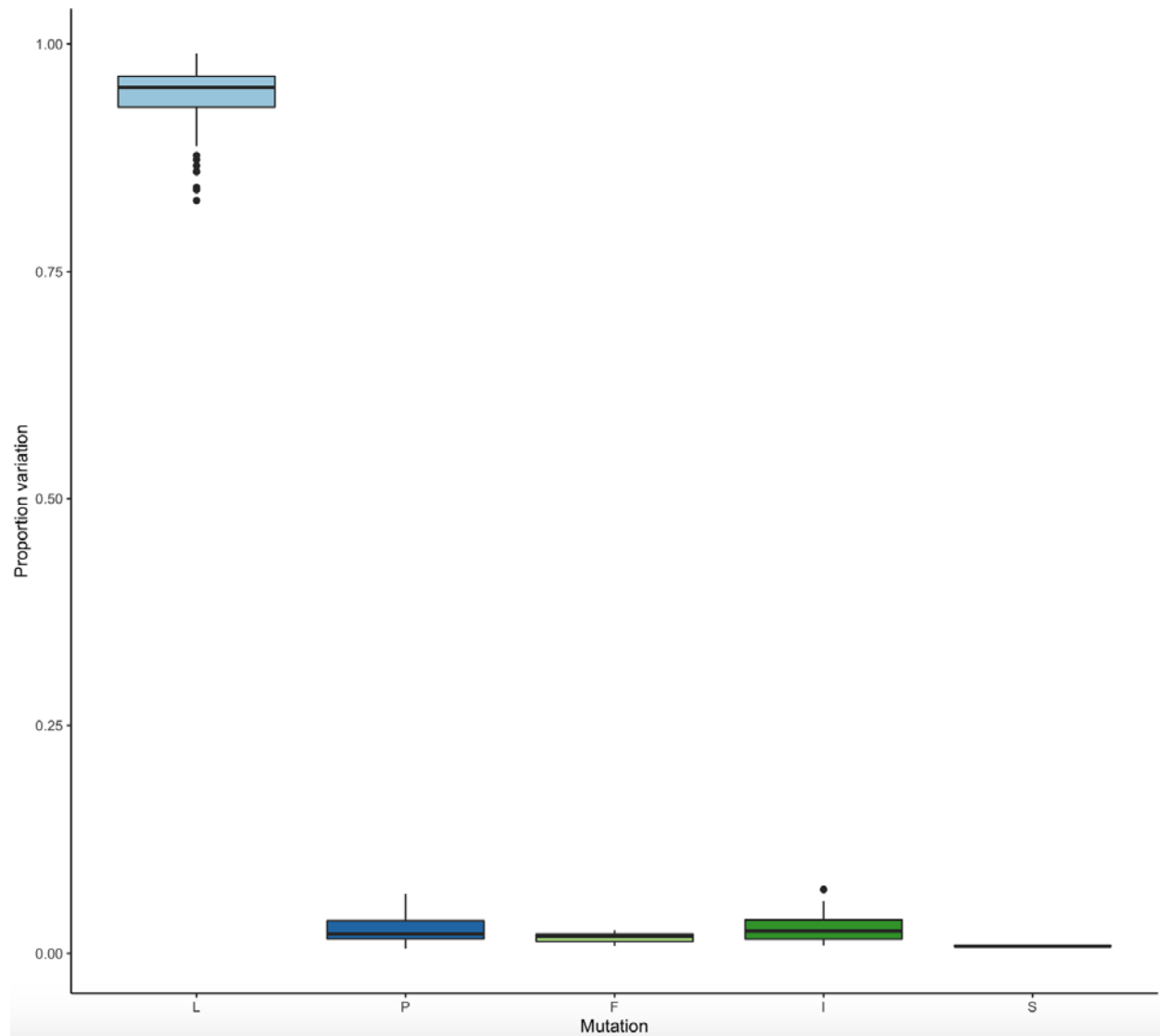


237

238    *Figure 7. Model for the transmission of variant genomes which encodes amino acids under*

239    *strong selection pressure showing the potential options for growth and transmission of viral*

19

240   *populations with either consensus viral genomes with P323 (cyan) and 323L (red) present in*

241   *minor variant genomes or in equilibrium or where 323L is in dominant viral genome sequence*

242   *and P323 is present in the minor variant genomes. Given the potential strong selection pressure*

243   *on this position the time post-infection transmission occurs is crucial in determining which*

244   *variant becomes dominant viral genome sequence.*

245

246   One of the predictions of this model is that whilst 323L in NSP12 might now be part of the

247   dominant viral genome sequence, other variants at this position will be present and persist (e.g.

248   P323) at this position. To test this contemporary sequence data (post the P323L and D614G

249   substitutions) that had been deposited between July and September 2021 on the Short Read

250   Archive was examined for variation at position 323 in NSP12 (Figure 8). The data indicated that

251   323L is the dominant variant, but P323 and other substitutions such as 323F are present as

252   minor genomic variants.

253

Figure 8. Amino acid mutations at site 323 in NSP12 in samples sequenced using the ARTIC-Nanopore approach (n=101) from July-September 2021 obtained from the Short Read Archive. The bioinformatics tool DiversiTools was used to generate proportions of the counts of amino acids at site 323 and showed that L is dominant in viral sequences from mid-late 2021, with P remaining a small proportion of the population alongside amino acids F, S and I.

21

259    **Discussion**

260    Several variants have come to dominate the global landscape of SARS-CoV-2 infections,

261    including ones with the initial D614G and P323L polymorphisms in the spike protein and NSP12

262    respectively (B.1), followed by Alpha (B.1.1.7), Delta (B.1.617.2) and Omicron (B.1.1529). These

263    have occurred in waves and are likely linked to increases in transmissibility [4], coupled with spike

264    variation-mediated immune escape [18,19], founder effects [20-22], behaviour patterns of hosts and

265    population density [23,24] and non-pharmaceutical interventions [25]. Whilst VoCs have differed in

266    terms of transmissibility, in general there has been no marked change in inherent morbidity

267    and mortality, although an early variant with a deletion in ORF8 was associated with a less

268    severe inflammatory response and better patient outcome [3].

269    Among the first major changes in the dominant viral genome sequence of SARS-CoV-2 were the

270    P323L and the D614G substitutions in NSP12 and the spike protein, respectively. Focus has

271    been placed on spike D614G and its association with increased infectivity [26]. We wanted to

272    investigate the selection pressure at these two sites by analysing the virus population in

273    humans over the period when the two substitutions became part of the dominant viral genome

274    sequence, as well as studying this in two non-human primate animal models. The first analysis

275    suggested rapid selection of P323L in NSP12 and D614G in the spike protein within humans.

276    This was reflected in the substitutions 323L and 614G polymorphisms in the minor genomic

277    variant population becoming the dominant viral genome sequence and replacing P323 and

278    D614 within a few days of within host selection (Figure 2). At the population level, data

279    suggested this selection was established over a two-month period in the UK (February and

280    March 2020). We note that although samples used in this study were collected early in the

281    pandemic in the UK, during the containment phase and in the early surge phase of Wave 1,

282    there was no evidence that the change from P323L in NSP12 and D614G in the spike protein

283    resulted in an increase in disease severity.

284    The selection pressure at these two positions (within an isolate close to the original Wuhan

285    outbreak) was evaluated in two non-human primate models for COVID-19 that recapitulate the

286    mild disease observed in most humans [13]. Here, the SARS-CoV-2 variant used for infection had

287    P323 in NSP12 and D614 in the spike protein in the dominant consensus sequence. At the

288    minor variant genome level, 323L in NSP12 was present with a frequency of 0.03% and 614G in

289    the spike protein at 0.02%. The sequence analysis indicated that for those animals where later

290    time points returned usable viral genomic information, the dominant viral genome sequence

291    now contained 323L in NSP12, but not necessarily 614G in the spike protein (Figure 5).

292    Recombinant viruses that differed at codon 323 in NSP12 in the context of a background with

293    D614G in the spike protein and showed that the P323 virus grew with a smaller plaque

294    morphology than a version with 323L. There are several different determinants of plaque size

295    including those related to *in vitro* growth rate, evasion of antiviral responses and cell to cell

296    fusion [27,28]. NSP12 has been shown to attenuate type I interferon production [29], and this may be

297    variant dependent. The mechanism behind the selection pressure acting on the P323L

298    substitution in both humans and non-human primate animal models is unknown. However,

299    NSP12 is the RNA dependent RNA polymerase, and such polymerase complexes can be

300    composed of both viral and host cell proteins[30,31]. We speculate that the P323L substitution

301    may alter the composition of the replication complex by altering interactions with the host cell

302    proteome and thereby facilitating virus replication. Therefore, it is tempting to speculate that

23

303    growth of viruses in cell lines from the original host species might drive the selection back. This

304    might provide a mechanism to narrow down candidates for the original zoonotic event(s).

305    In our model (Figure 7), an individual with the substitution present in a minor variant genome

306    with a selective advantage will see an increase in the proportion of this genome as infection

307    progresses. Under this pressure the minor variant genome will become the dominant viral

308    genome sequence. If transmission occurs early in infection, then the variant will be maintained

309    at a minor genomic variant level. If selective pressure is strong then the viral population that is

310    being transmitted will have the substitution as part of the dominant viral genome sequence –

311    and this will persist during further infections. Another consequence is that the sudden

312    emergence of a substitution as part of the dominant genome sequence may be due to founder

313    effect. For example, 323F in NSP12 that was identified in a cluster of cases in Norther Nevada

314    and in Nigeria (B.1.525). However, this substitution has not become part of the global dominant

315    viral genome sequence, despite that 323F was identified in samples from early 2020.

316    The data in this study indicates that in some cases it may be possible to predict the emergence

317    of a new dominant viral genome sequence and hence new variant. This would be based on

318    tracking the distribution and frequency of minor variant genomes at a population level, rather

319    than just focusing on providing information on the dominant viral genome sequence e.g.,

320    consensus level reporting. Whilst computationally more intensive and perhaps requiring higher

321    quality samples and sequencing data, the ability to earlier predict a newly emerging variant of

322    SARS-CoV-2 in the global landscape may aid in the evaluation of medical countermeasures and

323    non-pharmaceutical interventions.

324 **Materials and methods**

325 **Illumina for NHP NW samples**

326 Total RNA in each sample was extracted with QIAmp viral RNA extraction kit and eluted in pure

327 water. Following the manufacturer's protocols, total RNA was used as input material in to the

328 QIAseq FastSelect –rRNA HMR (Qiagen) protocol to remove cytoplasmic and mitochondrial

329 rRNA with a fragmentation time of 7 or 15 minutes. Subsequently, the NEBNext® Ultra™ II

330 Directional RNA Library Prep Kit for Illumina® (New England Biolabs) was used to generate the

331 RNA libraries, followed by 11 cycles of amplification and purification using AMPure XP beads.

332 Each library was quantified using Qubit and the size distribution assessed using the Agilent 2100

333 Bioanalyser, and the final libraries were pooled in equimolar ratios. The raw FASTQ files (2 x

334 150 bp) generated by an Illumina® NovaSeq 6000 (Illumina®, San Diego, USA) were trimmed to

335 remove Illumina adapter sequences using Cutadapt v1.2.1 [32]. The option "–O 3" was set, so the

336 that 3' end of any reads which matched the adapter sequence with greater than 3 bp was

337 trimmed off. The reads were further trimmed to remove low quality bases, using Sickle v1.200

338 [33] with a minimum window quality score of 20. After trimming, reads shorter than 10 bp were

339 removed.

340 The minor variations of amino acid in the genes of virus were called as our previous description

341 [34]. Hisat2 v2.1.0 [35] was used to map the trimmed reads on the cynomolgus (*M. fascicularis*) and

342 rhesus (*M. mulatta*) reference genome assemblies (release-94) downloaded from the Ensembl

343 FTP site. The unmapped reads were extracted by bam2fastq (v1.1.0) and then mapped on the

344 inoculum SARS-CoV-2 genome (GenBank sequence accession: NC_045512.2) using Bowtie2

345  v2.3.5.1 [35] by setting the options to parameters "--local -X 2000 --no-mixed", followed by SAM

346  file to BAM file conversion, sorting, and removal of the reads with a mapping quality score

347  below 11 using SAMtools v1.9 [36]. After that, the PCR and optical duplicate reads in the BAM

348  files were discarded using the MarkDuplicates in the Picard toolkit v2.18.25

349  (http://broadinstitute.github.io/picard/) with the option of "REMOVE_DUPLICATES=true". This

350  BAM file was then processed by the diversiutils script in DiversiTools

351  (http://josephhughes.github.io/btctools/) with the "-orfs" function to generate the number of

352  amino acid changes caused by the nucleotide deviation at each site in the protein. In order to

353  distinguish low frequency variants from Illumina sequence errors, the diversiutils script used

354  the calling algorithms based on the Illumina quality scores to calculate a P-value for each

355  variant at each nucleotide site [37]. The amino acid change was then filtered based on the P-value

356  (<0.05) to remove the low frequency variants from Illumina sequence errors.

357

358  **ARTIC Illumina for longitudinal swab samples and NHP NW samples**

359  Samples from clinical specimens were processed at CL3 at the University of Liverpool as part of

360  the study described in this chapter. Nasopharyngeal swabs were collected in viral transport

361  media. Swabs were left to defrost in a Tripass I cabinet in CL3. The swab was removed from the

362  tube and dipped in virkon before disposal to reduce dripping and aerosol generation. 250ml of

363  viral transport media was removed from the swab sample and added to 750ml of Trizol LS

364  (Invitrogen (10296028)) and mixed well. Remaining extraction was continued under CL2

365  conditions. All RNA samples were then treated with Turbo DNase (Invitrogen). SuperScript IV

366  (Invitrogen) was used to generate single-strand cDNA using random primer mix (NEB, Hitchin,

26

367  UK). ARTIC V3 PCR amplicons from the single-strand cDNA were generated following the

368  Nanopore      Protocol     of     PCR     tiling     of     SARS-CoV-2     virus     (Version:

369  PTC_9096_v109_revL_06Feb2020). The amplicons products were then used in Illumina

370  NEBNext Ultra II DNA Library preparation. Following 4 cycles of amplification the library was

371  purified using Ampure XP beads and quantified using Qubit and the size distribution assessed

372  using the Fragment Analyzer. Finally, the ARTIC library was sequenced on the Illumina®

373  NovaSeq 6000 platform (Illumina®, San Diego, USA) following the standard workflow. The

374  generated raw FASTQ files (2 x 250 bp) were trimmed to remove Illumina adapter sequences

375  using Cutadapt v1.2.1 26. The option "–O 3" was set, so the that 3' end of any reads which

376  matched the adapter sequence with greater than 3 bp was trimmed off. The reads were further

377  trimmed to remove low quality bases, using Sickle v1.200 27 with a minimum window quality

378  score of 20. After trimming, reads shorter than 10 bp were removed. The NHP NW total RNA

379  have been extracted and sequenced in our previous paper [38].

380  The variations of amino acid in the genes of the virus were called as our previous description [34].

381  Hisat2 v2.1.0 [35] was used to map the trimmed reads onto the human reference genome

382  assembly GRCh38 (release-91) downloaded from the Ensembl FTP site. The unmapped reads

383  were extracted by bam2fastq (v1.1.0) and then mapped on a known SARS-CoV-2 genome

384  (GenBank sequence accession: NC_045512.2) using Bowtie2 v2.3.5.1 [35] by setting the options to

385  parameters "--local -X 500 --no-mixed", followed by SAM file to BAM file conversion, sorting,

386  and removal of the reads with a mapping quality score below 11, not in pair, and not primary

387  and supplementary alignment using SAMtools v1.9 [36]. Bamclipper (v 1.0.0) [39] was used to trim

388  the ARTIC primer sequences on the mapped reads within the BAM files. The reads without

389    ARTIC primer sequences were also excluded in the further analysis. This trimmed BAM file was

390    then        processed        by        the        diversiutils        script        in        DiversiTools

391    (http://josephhughes.github.io/DiversiTools/) with the "-orfs" function to generate the number

392    of amino acid changes caused by the nucleotide deviation at each site in the protein in

393    comparison to the reference SARS-CoV-2 genome (NC_045512.2). In order to distinguish low

394    frequency variants from Illumina sequence errors, the diversiutils script used the calling

395    algorithms based on the Illumina quality scores to calculate a P-value for each variant at each

396    nucleotide site [37].

397    **Rapid Sequencing Long Amplicons (RSLA) nanopore for longitudinal swab samples**

398    Total RNA of longitudinal swab samples were extracted as described above. Sequencing

399    libraries for amplicons generated by RSLA [14] were prepared following the 'PCR tiling of SARS-

400    CoV-2 virus with Native Barcoding' protocol provided by Oxford Nanopore Technologies using

401    LSK109 and EXP-NBD104/114. The artic-ncov2019 pipeline v1.2.1 (https://artic.network/ncov-

402    2019/ncov2019-bioinformatics-sop.html) was used to filter the passed FASTQ files produced by

403    Nanopore sequencing with lengths between 800 and 1600. This pipeline was then used to map

404    the filtered reads on the reference SARS-CoV-2 genome (NC_045512.2) by minimap2 and

405    assigned each read alignment to a derived amplicon and excluded primer sequences based on

406    the RSLA primer schemes in the BAM files. These BAM files were further analysed using

407    DiversiTools (http://josephhughes.github.io/btctools/) with the "-orfs" function to generate the

408    ratio of amino acid change in the reads and coverage at each site of the protein in comparison

409    to the reference SARS-CoV-2 genome (NC_045512.2). The amino acids with highest ratio and

410    coverage > 10 were used to assemble the consensus protein sequences.

411

412 **Sanger sequencing**

413 cDNA template was amplified using Q5 High-Fidelity DNA Polymerase following the PCR

414 conditions: denaturation at 98$^{\circ}$C for 30 sec followed by 39 cycles of 10 sec denaturation at

415 98$^{\circ}$C, 30 sec annealing at 66$^{\circ}$C, and then 50 sec of extension at 72$^{\circ}$C. A final extension step was

416 done for 2 min at 72$^{\circ}$C. The primer sets used for amplification were (SARS-CoV-

417 2_15_LEFT=ATACGCCAACTTAGGTGAACG, SARS-CoV-2_15_RIGHT= AACATGTTG-TGCCAACCACC)

418 to detect the P323L mutation or (SARS-CoV-2_24_LEFT= TTGAACTTCTACATGCACCAGC, SARS-

419 CoV-2_RIGHT=CCAGAAGTGATTGTACCCGC) to detect the D614G mutation. PCR products were

420 purified using AMPure XP beads (Beckman Coulter) and quantified using the Qubit High

421 Sensitivity 1X dsDNA kit (Invitrogen). To visualise band quality, PCR products were run on a

422 1.5% agarose gel. 10 ng of each amplified product was sent for sanger sequencing (Source

423 Bioscience, UK).

424

425 **Cells**

426 African green monkey kidney C1008 (Vero E6) cells (Public Health England, PHE) were cultured

427 in Dulbecco's minimal essential medium (DMEM) (Sigma) with 10% foetal bovine serum (FBS)

428 (Sigma) and 0.05mg/ml gentamicin at 37°C/5% CO2. Vero/hSLAM cells (PHE) were grown in

429 DMEM with 10% FBS and 0.05mg/ml gentamicin (Merck) with the addition of 0.4mg/ml

430 Geneticin (G418; Thermofisher) at 37°C/5% CO2. Human ACE2-A549 (hACE2-A549), a lung

431 epithelial cell line which overexpresses the ACE-2 receptor [40], were cultured in DMEM with 10%

432    FBS and 0.05mg/ml gentamicin with the addition of 10µg/ml Blasticidin (Invitrogen). Only

433    passage 3-10 cultures were used for experiments.

434

435    **Generation and culture of recombinant viruses**

436    Recombinant SARS-CoV-2 viruses were generated by reverse genetics using the

437    "transformation-associated recombination" in yeast approach [41]. 11 cDNA fragments with 70 bp

438    end-terminal overlaps which spanned the entire SARS-CoV-2 isolate Wuhan-Hu-1 genome

439    (GenBank accession: NC_045512) were produced by GeneArt™ synthesis (Invitrogen™,

440    ThermoFisher) as inserts in sequence verified, stable plasmid clones. The 5′ terminal cDNA

441    fragment was modified to contain a T7 RNA polymerase promoter and an extra "G" nucleotide

442    immediately upstream of the SARS-CoV-2 5′ sequence, whilst the 3′ terminal cDNA fragment

443    was modified such that the 3' end of the SARS-CoV-2 genome was followed by a stretch of 33

444    "A"s followed by the unique restriction enzyme site Asc I. The inserts were amplified by PCR

445    using a Platinum SuperFi II mastermix (ThermoFisher) and assembled into full length SARS-CoV-

446    2 cDNA clones in the YAC vector pYESL1 using a GeneArt™ High-Order Genetic Assembly System

447    (A13285, Invitrogen™, ThermoFisher) according to the manufacturer's instructions. RNA

448    transcripts produced from the YAC clones by transcription with T7 polymerase were used to

449    recover infectious virus. Two viruses were produced on the Wuhan-Hu-1 background and had a

450    D614G substitution in the spike protein and differed at amino acid position 323 in NSP12 with

451    either a P or L, these were termed Wuhan/614G/P323 and Wuhan/614G/323L, respectively.

452    Whole genome sequencing confirmed the presence of these changes. Stocks of the viruses

453    were cultured in Vero E6 cells in DMEM containing 2% FBS, 0.05mg/ml gentamicin and

454    harvested 72 hours post inoculation. Virus stocks were aliquoted and stored at -80°C. All stocks

455    were titred by plaque assay on Vero E6 cells and pictures of the resulting plaques recorded.

456

457    **Serial passage of SARS-CoV-2 Victoria/01/2020**

458    SARS-CoV-2 Victoria/01/2020 was passaged three times in Vero/hSLAM cells prior to receiving

459    it. hACE2-A549 cells were then infected at an MOI of 0.01 and incubated for 72 hours (Passage

460    4). Following this, 100μl was passaged to fresh cells and incubated at 37C for 1 hour. After the

461    incubation, media was topped up with DMEM containing 2% FBS, 0.05mg/ml gentamicin and

462    incubated for 72 hours (Passage 5). This process was repeated until Passage 13 (a total of ten

463    passages through hACE2-A549 cells).

464

465    **Analysis of global sequences from July-September 2021**

466    Sequences were obtained from the Short Read Archive (SRA) under accession numbers: ERR6343731,

467    ERR6343734, ERR6343745, ERR6343747, ERR6343749, ERR6344225, ERR6346453, ERR6346456,

468    ERR6346459, ERR6758978, ERR6758981, ERR6759296, ERR6761288, ERR6761458, ERR6761562,

469    ERR6761570, ERR6761711, ERR6761986, ERR6762387, ERR6762545, ERR6762546, ERR6825821,

470    ERR6878898, ERR6879599, ERR6879604, ERR6887797, ERR6887811, ERR6887812, ERR6887820,

471    ERR6888048, ERR6888063, ERR6888078, ERR6888265, ERR6888283, SRR16376487, SRR16376490,

472    SRR16376491, SRR16376494, SRR16376495, SRR16376496, SRR16376497, SRR16376501, SRR16376502,

473    SRR16376505, SRR16376510, SRR16376515, SRR16376516, SRR16376522, SRR16376523, SRR16376524,

474    SRR16376526, SRR16376529, SRR16376530, SRR16376531, SRR16376536, SRR16376540, SRR16376543,

475    SRR16376544, SRR16376547, SRR16376551, SRR16376552, SRR16376554, SRR16376557, SRR16376559,

476    SRR16376573, SRR16376580, SRR16376589, SRR16376599, SRR16376608, SRR16376613, SRR16376614,

477    SRR16376648, SRR16376678, SRR16376782, SRR16376802, SRR16376804, SRR16376807, SRR16376810,

478    SRR16376884, SRR16376904, SRR16376907, SRR16376912, SRR16376913, SRR16376914, SRR16376916,

479    SRR16376921, SRR16376922, SRR16376925, SRR16376927, SRR16376928, SRR16376929, SRR16376932,

480    SRR16376935, SRR16376939, SRR16376940, SRR16376941, SRR16376943, SRR16376944, SRR16376946,

481    SRR16376949, SRR16376951. All sequences were ARTIC-Nanopore sequenced using the V3 primer

482    scheme and downloaded as SRA files. The SRA files were converted to FASTQ files using the SRA Toolkit

483    v2.11.3 (https://github.com/ncbi/sra-tools) command fastq-dump. The FASTQ files were processed

484    through the artic-ncov2019 v1.2.1 pipeline (https://artic.network/ncov-2019/ncov2019-bioinformatics-

485    sop.html) and the DiversiTools tool (https://github.com/josephhughes/DiversiTools) as described above.

486 **References**

487 1 Worobey, M. *et al.* The emergence of SARS-CoV-2 in Europe and North America. *Science*

488 **370**, 564-570, doi:10.1126/science.abc8169 (2020).

489 2 Davidson, A. D. *et al.* Characterisation of the transcriptome and proteome of SARS-CoV-

490 2 reveals a cell passage induced in-frame deletion of the furin-like cleavage site from the

491 spike glycoprotein. *Genome Med* **12**, 68, doi:10.1186/s13073-020-00763-0 (2020).

492 3 Young, B. E. *et al.* Effects of a major deletion in the SARS-CoV-2 genome on the severity

493 of infection and the inflammatory response: an observational cohort study. *Lancet* **396**,

494 603-611, doi:10.1016/S0140-6736(20)31757-8 (2020).

495 4 Hou, Y. J. *et al.* SARS-CoV-2 D614G variant exhibits efficient replication ex vivo and

496 transmission in vivo. *Science*, doi:10.1126/science.abe8499 (2020).

497 5 Yang, H. C. *et al.* Analysis of genomic distributions of SARS-CoV-2 reveals a dominant

498 strain type with strong allelic associations. *Proc Natl Acad Sci U S A*,

499 doi:10.1073/pnas.2007840117 (2020).

500 6 Simmonds, P. Rampant C-->U Hypermutation in the Genomes of SARS-CoV-2 and Other

501 Coronaviruses: Causes and Consequences for Their Short- and Long-Term Evolutionary

502 Trajectories. *mSphere* **5**, doi:10.1128/mSphere.00408-20 (2020).

503   7      Ratcliff, J. & Simmonds, P. Potential APOBEC-mediated RNA editing of the genomes of

504          SARS-CoV-2 and other coronaviruses and its impact on their longer term evolution.

505          *Virology* **556**, 62-72, doi:10.1016/j.virol.2020.12.018 (2021).

506   8      Dong, X. *et al.* Identification and quantification of SARS-CoV-2 leader subgenomic mRNA

507          gene junctions in nasopharyngeal samples shows phasic transcription in animal models

508          of COVID-19 and dysregulation at later time points that can also be identified in

509          humans. *bioRxiv*, 2021.2003.2003.433753, doi:10.1101/2021.03.03.433753 (2021).

510   9      Peacock, T. P., Penrice-Randal, R., Hiscox, J. A. & Barclay, W. S. SARS-CoV-2 one year on:

511          evidence for ongoing viral adaptation. *J Gen Virol* **102**, doi:10.1099/jgv.0.001584 (2021).

512   10     Lythgoe, K. A. *et al.* SARS-CoV-2 within-host diversity and transmission. *Science* **372**,

513          doi:10.1126/science.abg0821 (2021).

514   11     Dowall, S. D. *et al.* Elucidating variations in the nucleotide sequence of Ebola virus

515          associated with increasing pathogenicity. *Genome Biol* **15**, 540, doi:10.1186/PREACCEPT-

516          1724277741482641 (2014).

517   12     Dong, X. *et al.* Variation around the dominant viral genome sequence contributes to

518          viral load and outcome in patients with Ebola virus disease. *Genome Biol* **21**, 238,

519          doi:10.1186/s13059-020-02148-3 (2020).

520     13      Salguero, F. J. *et al.* Comparison of rhesus and cynomolgus macaques as an infection

521             model for COVID-19. *Nat Commun* **12**, 1260, doi:10.1038/s41467-021-21389-9 (2021).


522     14      Moore, S. C. *et al.* Amplicon-Based Detection and Sequencing of SARS-CoV-2 in

523             Nasopharyngeal Swabs from Patients With COVID-19 and Identification of Deletions in

524             the Viral Genome That Encode Proteins Involved in Interferon Antagonism. *Viruses* **12**,

525             doi:10.3390/v12101164 (2020).


526     15      Nasir, J. A. *et al.* A Comparison of Whole Genome Sequencing of SARS-CoV-2 Using

527             Amplicon-Based Sequencing, Random Hexamers, and Bait Capture. *Viruses* **12**,

528             doi:10.3390/v12080895 (2020).


529     16      Caly, L. *et al.* Isolation and rapid sharing of the 2019 novel coronavirus (SARS-CoV-2)

530             from the first patient diagnosed with COVID-19 in Australia. *Med J Aust* **212**, 459-462,

531             doi:10.5694/mja2.50569 (2020).


532     17      Prince, T. *et al.* Sequence analysis of SARS-CoV-2 in nasopharyngeal samples from

533             patients with COVID-19 illustrates population variation and diverse phenotypes, placing

534             the in vitro growth properties of B.1.1.7 and B.1.351 lineage viruses in context. *bioRxiv*,

535             2021.2003.2030.437704, doi:10.1101/2021.03.30.437704 (2021).

536     18     Wang, B. *et al.* Resistance of SARS-CoV-2 variants to neutralization by convalescent

537             plasma from early COVID-19 outbreak in Singapore. *NPJ Vaccines* **6**, 125,

538             doi:10.1038/s41541-021-00389-2 (2021).

539     19     Saad-Roy, C. M. *et al.* Epidemiological and evolutionary considerations of SARS-CoV-2

540             vaccine dosing regimes. *Science* **372**, 363-370, doi:10.1126/science.abg8663 (2021).

541     20     Gomez-Carballa, A., Bello, X., Pardo-Seco, J., Martinon-Torres, F. & Salas, A. Mapping

542             genome variation of SARS-CoV-2 worldwide highlights the impact of COVID-19 super-

543             spreaders. *Genome Res* **30**, 1434-1448, doi:10.1101/gr.266221.120 (2020).

544     21     Diez-Fuertes, F. *et al.* A Founder Effect Led Early SARS-CoV-2 Transmission in Spain. *J*

545             *Virol* **95**, doi:10.1128/JVI.01583-20 (2021).

546     22     Tasakis, R. N. *et al.* SARS-CoV-2 variant evolution in the United States: High

547             accumulation of viral mutations over time likely through serial Founder Events and

548             mutational bursts. *PLoS One* **16**, e0255169, doi:10.1371/journal.pone.0255169 (2021).

549     23     Ward, T. *et al.* Growth, reproduction numbers and factors affecting the spread of SARS-

550             CoV-2 novel variants of concern in the UK from October 2020 to July 2021: a modelling

551             analysis. *BMJ Open* **11**, e056636, doi:10.1136/bmjopen-2021-056636 (2021).

552    24    Rader, B. *et al.* Crowding and the shape of COVID-19 epidemics. *Nat Med* **26**, 1829-1834,

553            doi:10.1038/s41591-020-1104-0 (2020).

554    25    Kraemer, M. U. G. *et al.* Spatiotemporal invasion dynamics of SARS-CoV-2 lineage B.1.1.7

555            emergence. *Science* **373**, 889-895, doi:10.1126/science.abj0113 (2021).

556    26    Korber, B. *et al.* Tracking Changes in SARS-CoV-2 Spike: Evidence that D614G Increases

557            Infectivity of the COVID-19 Virus. *Cell* **182**, 812-827 e819, doi:10.1016/j.cell.2020.06.043

558            (2020).

559    27    Goh, K. C. *et al.* Molecular determinants of plaque size as an indicator of dengue virus

560            attenuation. *Sci Rep* **6**, 26100, doi:10.1038/srep26100 (2016).

561    28    Kato, F. *et al.* Characterization of large and small-plaque variants in the Zika virus clinical

562            isolate ZIKV/Hu/S36/Chiba/2016. *Sci Rep* **7**, 16160, doi:10.1038/s41598-017-16475-2

563            (2017).

564    29    Wang, W. *et al.* SARS-CoV-2 nsp12 attenuates type I interferon production by inhibiting

565            IRF3 nuclear translocation. *Cell Mol Immunol* **18**, 945-953, doi:10.1038/s41423-020-

566            00619-y (2021).

567    30    Munday, D. C. *et al.* Interactome analysis of the human respiratory syncytial virus RNA

568            polymerase complex identifies protein chaperones as important cofactors that promote

569     L-protein stability and RNA synthesis. *J Virol* **89**, 917-930, doi:10.1128/JVI.01783-14

570     (2015).

571   31   Noton, S. L., Aljabr, W., Hiscox, J. A., Matthews, D. A. & Fearns, R. Factors affecting de

572     novo RNA synthesis and back-priming by the respiratory syncytial virus polymerase.

573     *Virology* **462-463**, 318-327, doi:10.1016/j.virol.2014.05.032 (2014).

574   32   Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing

575     reads. *EMBnet. journal* **17**, 10-12 (2011).

576   33   Joshi, N. & Fass, J.   (2011).

577   34   Dong, X. *et al.* Variation around the dominant viral genome sequence contributes to

578     viral load and outcome in patients with Ebola virus disease. *Genome biology* **21**, 1-20

579     (2020).

580   35   Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory

581     requirements. *Nature methods* **12**, 357 (2015).

582   36   Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-

583     2079 (2009).

584    37    Morelli, M. J. *et al.* Evolution of foot-and-mouth disease virus intra-sample sequence

585          diversity during serial transmission in bovine hosts. *Veterinary research* **44**, 12 (2013).

586    38    Dong, X. *et al.* Identification and quantification of SARS-CoV-2 leader subgenomic mRNA

587          gene junctions in nasopharyngeal samples shows phasic transcription in animal models

588          of COVID-19 and aberrant pattens in humans. *bioRxiv* (2021).

589    39    Au, C. H., Ho, D. N., Kwong, A., Chan, T. L. & Ma, E. S. BAMClipper: removing primers

590          from alignments to minimize false-negative mutations in amplicon next-generation

591          sequencing. *Scientific reports* **7**, 1-7 (2017).

592    40    Buchrieser, J. *et al.* Syncytia formation by SARS-CoV-2-infected cells. *EMBO J* **39**,

593          e106267, doi:10.15252/embj.2020106267 (2020).

594    41    Thi Nhu Thao, T. *et al.* Rapid reconstruction of SARS-CoV-2 using a synthetic genomics

595          platform. *Nature* **582**, 561-565, doi:10.1038/s41586-020-2294-9 (2020).

596

**Acknowledgments**

40

623

**Author contributions**

625    Conceptualization: DAM, AD, MWC and JAH. Data curation: HG, XD, RP-R, DAM, AD and JAH.

626    Formal analysis: HG, XD, NR, RP-R, PD, ADD, TP, AD and JAH. Funding acquisition: MGS, PJMO,

627    JKB, DAM, LT, AD, ADD, MWC and JAH. Investigation: HG, XD, NR, RP-R, ADD, GTS, BJ, MKW,

628    ME, JB, TJ, FJS, SRE, JT, CH and JAH. Methodology: HG, XD, GTS, RP-R, ADD, MKW, FJS and JT.

629    Project administration: JAH. Resources: MWC, FJS, JT, YH, MGS, PJMO, ADD, JKB, LT. Software:

630    HG, XD, RP-R, DAM and AD. Supervision: MWC, AD, ADD, SRE and JAH. Validation: HG, XD, RP-R,

631    NR, CH, GTS, DAM, AD and JAH. Visualisation: HG, XD and RP-R. Writing – original draft: HG, XD,

632    RP-R and JAH. Writing – reviewing and editing. HG, XD, RP-R, DAM, AD, LT, ADD, MWC and JAH.

633

**Availability of data and materials**

635    All viral sequence data used in this analysis were deposited with the National Center for

636    Biotechnology Information under the project accession number PRJNA789459 and can be

637    accessed via https://www.ncbi.nlm.nih.gov/bioproject/PRJNA789459.

638
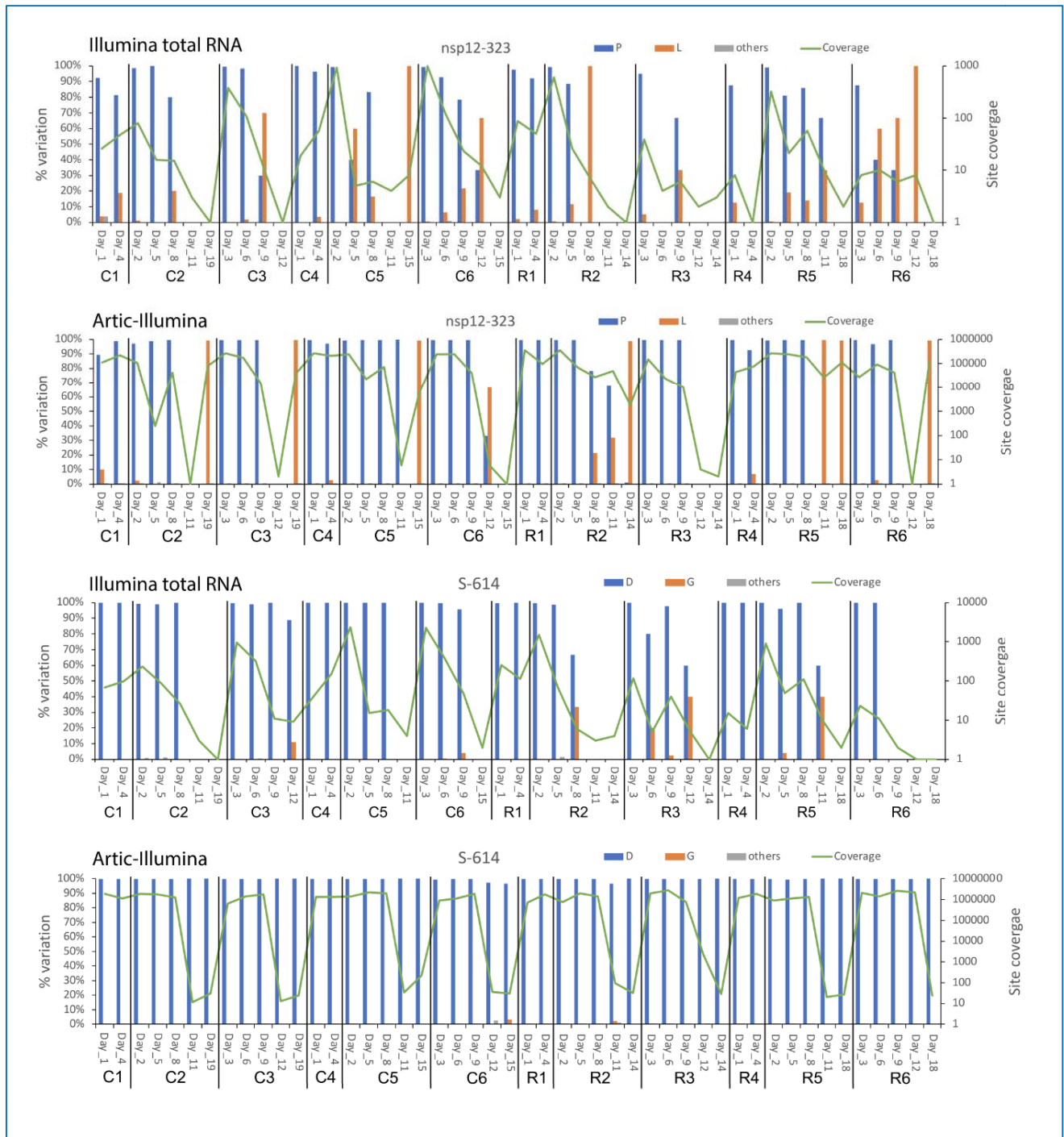
**Competing interests**

640    The authors declare that they have no competing interests.

41

641

642 **Ethics approval and consent to participate**

643 Patients were recruited under the International Severe Acute Respiratory and emerging

644 Infection Consortium (ISARIC) Clinical Characterisation Protocol CCP (https://isaric.net/ccp) by

645 giving informed consent. ISARIC CCP was reviewed and approved by the national research

646 ethics service, Oxford (13/SC/0149). All experimental work on non-human primates was

647 conducted under the authority of a UK Home Office approved project license (PDC57C033) that

648 had been subject to local ethical review at PHE Porton Down by the Animal Welfare and Ethical

649 Review Body (AWERB) and approved as required by the Home Office Animals (Scientific

650 Procedures) Act 1986. None of the animals had been used previously for experimental

651 procedures.

652 **Supplementary information**

653



654  *Supplementary Figure 1. Analysis of NSP12 position 323 and the spike protein position 614 in*

655  *SARS-CoV-2 from nasopharyngeal swabs taken longitudinally from infected cynomolgus (C) and*

656     *rhesus (R) macaques. Sequencing was performed using both an Illumina shot gun sequencing*

657     *approach (Illumina total RNA) or using an ARTIC-Illumina approach to specifically amplify SARS-*

658     *CoV-2 RNA. The day post infection is shown for each individual animal (number after the C or R)*

659     *(x-axis). For each position of interest either the P (for position 323 in NSP12) or D (for position*

660     *614 in the spike protein) is shown in blue, and the substitution of L or G, shown in orange,*

661     *respectively. Grey indicates other substitutions at that position. The left-hand y-axis indicates*

662     *the % variation at the indicated position and the right-hand x-axis shows amino acid site*

663     *coverage for each position (green line). The % variation was only shown for these sites with*

664     *coverage > 5.*

665

666     Supplementary information

667     ISARIC4C Investigators

668     Consortium Lead Investigator: J Kenneth Baillie.

669

670     Chief Investigator: Malcolm G Semple.

671

672     Co-Lead Investigator: Peter JM Openshaw.

673

674     ISARIC Clinical Coordinator: Gail Carson.

675

676     Co-Investigator: Beatrice Alex, Petros Andrikopoulos, Benjamin Bach, Wendy S Barclay, Debby

677     Bogaert, Meera Chand, Kanta Chechi, Graham S Cooke, Ana da Silva Filipe, Thushan de Silva,

678     Annemarie B Docherty, Gonçalo dos Santos Correia, Marc-Emmanuel Dumas, Jake Dunning,

679     Tom Fletcher, Christoper A Green, William Greenhalf, Julian L Griffin, Rishi K Gupta, Ewen M

680     Harrison, Julian A Hiscox, Antonia Ying Wai Ho, Peter W Horby, Samreen Ijaz, Saye Khoo, Paul

681     Klenerman, Andrew Law, Matthew R Lewis, Sonia Liggi, Wei Shen Lim, Lynn Maslen, Alexander J

682     Mentzer, Laura Merson, Alison M Meynert, Shona C Moore, Mahdad Noursadeghi, Michael

683    Olanipekun, Anthonia Osagie, Massimo Palmarini, Carlo Palmieri, William A Paxton, Georgios

684    Pollakis, Nicholas Price, Andrew Rambaut, David L Robertson, Clark D Russell, Vanessa Sancho-

685    Shimizu, Caroline J Sands, Janet T Scott, Louise Sigfrid, Tom Solomon, Shiranee Sriskandan,

686    David Stuart, Charlotte Summers, Olivia V Swann, Zoltan Takats, Panteleimon Takis, Richard S

687    Tedder, AA Roger Thompson, Emma C Thomson, Ryan S Thwaites, Lance CW Turtle, Maria

688    Zambon.

689

690    Project Manager: Hayley Hardwick, Chloe Donohue, Fiona Griffiths, Wilna Oosthuyzen.

691

692    Project Administrator: Cara Donegan, Rebecca G. Spencer.

693

694    Data Analyst: Lisa Norman , Riinu Pius, Thomas M Drake, Cameron J Fairfield, Stephen R Knight,

695    Kenneth A Mclean, Derek Murphy, Catherine A Shaw.

696

697    Data and Information System Manager: Jo Dalton, Michelle Girvan, Egle Saviciute, Stephanie

698    Roberts, Janet Harrison, Laura Marsh, Marie Connor, Sophie Halpin, Clare Jackson, Carrol

699    Gamble, Daniel Plotkin, James Lee.

700

46

701 Data Integration and Presentation: Gary Leeming, Andrew Law, Murray Wham, Sara Clohisey,

702 Ross Hendry, James Scott-Brown.

703

704 Material Management: Victoria Shaw, Sarah E McDonald.

705

706 Patient Engagement: Seán Keating.

707

708 Outbreak Laboratory Staff and Volunteers: Katie A. Ahmed, Jane A Armstrong, Milton

709 Ashworth, Innocent G Asiimwe, Siddharth Bakshi, Samantha L Barlow, Laura Booth, Benjamin

710 Brennan, Katie Bullock, Benjamin WA Catterall, Jordan J Clark, Emily A Clarke, Sarah Cole, Louise

711 Cooper, Helen Cox, Christopher Davis, Oslem Dincarslan, Chris Dunn, Philip Dyer, Angela Elliott,

712 Anthony Evans, Lorna Finch, Lewis WS Fisher, Terry Foster, Isabel Garcia-Dorival, Philip

713 Gunning, Catherine Hartley, Rebecca L Jensen, Christopher B Jones, Trevor R Jones, Shadia

714 Khandaker, Katharine King, Robyn T. Kiy, Chrysa Koukorava, Annette Lake, Suzannah Lant, Diane

715 Latawiec, Lara Lavelle-Langham, Daniella Lefteri, Lauren Lett, Lucia A Livoti, Maria Mancini,

716 Sarah McDonald, Laurence McEvoy, John McLauchlan, Soeren Metelmann, Nahida S Miah,

717 Joanna Middleton, Joyce Mitchell, Shona C Moore, Ellen G Murphy, Rebekah Penrice-Randal,

718 Jack Pilgrim, Tessa Prince, Will Reynolds, P. Matthew Ridley, Debby Sales, Victoria E Shaw,

719 Rebecca K Shears, Benjamin Small, Krishanthi S Subramaniam, Agnieska Szemiel, Aislynn

47

720     Taggart, Jolanta Tanianis-Hughes, Jordan Thomas, Erwan Trochu, Libby van Tonder, Eve

721     Wilcock, J. Eunice Zhang, Lisa Flaherty, Nicole Maziere, Emily Cass, Alejandra Doce Carracedo,

722     Nicola Carlucci , Anthony Holmes, Hannah Massey.

723

724     Edinburgh Laboratory Staff and Volunteers: Lee Murphy, Sarah McCafferty, Richard Clark, Angie

725     Fawkes, Kirstie Morrice, Alan Maclean, Nicola Wrobel, Lorna Donnelly, Audrey Coutts,

726     Katarzyna Hafezi, Louise MacGillivray, Tammy Gilchrist.

727

728     Local Principal Investigators: Kayode Adeniji, Daniel Agranoff, Ken Agwuh, Dhiraj Ail, Erin L.

729     Aldera, Ana Alegria, Sam Allen, Brian Angus, Abdul Ashish, Dougal Atkinson, Shahedal Bari,

730     Gavin Barlow, Stella Barnass, Nicholas Barrett, Christopher Bassford, Sneha Basude, David

731     Baxter, Michael Beadsworth, Jolanta Bernatoniene, John Berridge, Colin Berry, Nicola Best,

732     Pieter Bothma, David Chadwick, Robin Brittain-Long, Naomi Bulteel, Tom Burden, Andrew

733     Burtenshaw, Vikki Caruth, David Chadwick, Duncan Chambler, Nigel Chee, Jenny Child, Srikanth

734     Chukkambotla, Tom Clark, Paul Collini, Catherine Cosgrove, Jason Cupitt, Maria-Teresa Cutino-

735     Moguel, Paul Dark, Chris Dawson, Samir Dervisevic, Phil Donnison, Sam Douthwaite, Andrew

736     Drummond, Ingrid DuRand, Ahilanadan Dushianthan, Tristan Dyer, Cariad Evans, Chi Eziefula,

737     Chrisopher Fegan, Adam Finn, Duncan Fullerton, Sanjeev Garg, Sanjeev Garg, Atul Garg,

738     Effrossyni Gkrania-Klotsas, Jo Godden, Arthur Goldsmith, Clive Graham, Elaine Hardy, Stuart

739     Hartshorn, Daniel Harvey, Peter Havalda, Daniel B Hawcutt, Maria Hobrok, Luke Hodgson, Anil

740    Hormis, Michael Jacobs, Susan Jain, Paul Jennings, Agilan Kaliappan, Vidya Kasipandian,

741    Stephen Kegg, Michael Kelsey, Jason Kendall, Caroline Kerrison, Ian Kerslake, Oliver Koch, Gouri

742    Koduri, George Koshy, Shondipon Laha, Steven Laird, Susan Larkin, Tamas Leiner, Patrick Lillie,

743    James Limb, Vanessa Linnett, Jeff Little, Mark Lyttle, Michael MacMahon, Emily MacNaughton,

744    Ravish Mankregod, Huw Masson, Elijah Matovu, Katherine McCullough, Ruth McEwen, Manjula

745    Meda, Gary Mills, Jane Minton, Mariyam Mirfenderesky, Kavya Mohandas, Quen Mok, James

746    Moon, Elinoor Moore, Patrick Morgan, Craig Morris, Katherine Mortimore, Samuel Moses,

747    Mbiye Mpenge, Rohinton Mulla, Michael Murphy, Megan Nagel, Thapas Nagarajan, Mark

748    Nelson, Lillian Norris, Matthew K. O'Shea, Igor Otahal, Marlies Ostermann, Mark Pais, Carlo

749    Palmieri, Selva Panchatsharam, Danai Papakonstantinou, Hassan Paraiso, Brij Patel, Natalie

750    Pattison, Justin Pepperell, Mark Peters, Mandeep Phull, Stefania Pintus, Jagtur Singh Pooni, Tim

751    Planche, Frank Post, David Price, Rachel Prout, Nikolas Rae, Henrik Reschreiter, Tim Reynolds,

752    Neil Richardson, Mark Roberts, Devender Roberts, Alistair Rose, Guy Rousseau, Bobby Ruge,

753    Brendan Ryan, Taranprit Saluja, Matthias L Schmid, Aarti Shah, Prad Shanmuga, Anil Sharma,

754    Anna Shawcross, Jeremy Sizer, Manu Shankar-Hari, Richard Smith, Catherine Snelson, Nick

755    Spittle, Nikki Staines, Tom Stambach, Richard Stewart, Pradeep Subudhi, Tamas Szakmany, Kate

756    Tatham, Jo Thomas, Chris Thompson, Robert Thompson, Ascanio Tridente, Darell Tupper-Carey,

757    Mary Twagira, Nick Vallotton, Rama Vancheeswaran, Lisa Vincent-Smith, Shico Visuvanathan,

758    Alan Vuylsteke, Sam Waddy, Rachel Wake, Andrew Walden, Ingeborg Welters, Tony

759    Whitehouse, Paul Whittaker, Ashley Whittington, Padmasayee Papineni, Meme Wijesinghe,

760    Martin Williams, Lawrence Wilson, Sarah Cole, Stephen Winchester, Martin Wiselka, Adam

761    Wolverson, Daniel G Wootton, Andrew Workman, Bryan Yates, Peter Young.