# Direct RNA sequencing and early evolution of SARS-CoV-2

George Taiaroa[1,2], Daniel Rawlinson[1], Leo Featherstone[1], Miranda Pitt[1], Leon Caly[2], Julian Druce[2], Damian Purcell[1], Leigh Harty[1], Thomas Tran[2], Jason Roberts[2], Nichollas Scott[1], Mike Catton[2], Deborah Williamson[1,3], Lachlan Coin[1·], Sebastian Duchene[1·]

1. Department of Microbiology and Immunology, University of Melbourne at The Peter Doherty Institute for Infection and Immunity, Melbourne, Australia
2. Victorian Infectious Diseases Reference Laboratory, Royal Melbourne Hospital, at the Peter Doherty Institute for Infection and Immunity, Victoria, Australia.
3. Department of Microbiology, Royal Melbourne Hospital, Victoria, Australia
· These authors contributed equally to the work

Corresponding author:

George Taiaroa, The Peter Doherty Institute for Infection and Immunity, University of Melbourne and Victorian Infectious Diseases Reference Laboratory, Melbourne, Australia Tel: +61 (0)3 8344 5466. Email: george.taiaroa@unimelb.edu.au

Abstract - Fundamental aspects of SARS-CoV-2 biology remain to be described, having the potential to provide insight to the response effort for this high-priority pathogen. Here we describe the first native RNA sequence of SARS-CoV-2, detailing the coronaviral transcriptome and epitranscriptome, and share these data publicly. A data-driven inference of viral genetic features and evolutionary rate is also made. The rapid sharing of sequence information throughout the SARS-CoV-2 pandemic represents an inflection point for public health and genomic epidemiology, providing early insights into the biology and evolution of this emerging pathogen.

The pandemic of severe acute respiratory syndrome 2 (SARS-CoV-2), causing the disease COVID-19 and originating in Wuhan, China, has spread to more than 200 countries and territories, and has caused more than 1,000,000 cases globally [1-4]. SARS-CoV-2 is a positive-sense single-stranded RNA ((+)ssRNA) virus, belonging to the *Coronaviridae* family and *betacoronavirus* genus [5]. Related betacoronaviruses are capable of infection and

35 ongoing transmission in mammalian and avian hosts, resulting in illness in humans such as

36 Middle East respiratory syndrome (MERS) and the original severe acute respiratory

37 syndrome (SARS) as examples [6-7]. Based on the limited sampling of potential reservoir

38 species, SARS-CoV-2 has been found to be most similar to bat betacoronaviruses on a

39 genomic level, potentially indicating that bats are a natural reservoir [5,8].

40

41 The genome sequence of SARS-CoV-2 was rapidly determined and shared on January 5th

42 of 2020, being 29,903 nucleotides in length, and annotated based on sequence similarity to

43 other coronaviruses (GenBank: MN908947.3). As the emergence of SARS-CoV-2 has

44 escalated, genomic analyses have played a key role in public health responses, including in

45 the design of appropriate molecular diagnostics and supporting epidemiological efforts to

46 track and contain the outbreak [9,10]. Taken together, publicly available sequence data

47 suggest a recently occurring, point-source outbreak, as described in online sources [10-12].

48

49 Aspects of the response assume that the biology of SARS-CoV-2 is comparable with

50 previously characterised coronaviruses, including the annotation of genes and the estimation

51 of molecular evolutionary rates [11-12]. It remains highly relevant to determine these

52 features experimentally with SARS-CoV-2-specific data, potentially revealing other insights

53 into the biology of this emergent pathogen. To address this, here we describe (i) the first

54 native RNA sequence of SARS-CoV-2, detailing the coronaviral transcriptome and

55 epitranscriptome, and (ii) estimates of coronaviral evolutionary rates and related timescales,

56 based on data available at this stage of the outbreak.

57

58 Characterised coronaviruses have some of the largest genomes among RNA viruses, and

59 express their genetic content as a nested set of polyadenylated mRNA transcripts (Figure 1),

60 with lengths corresponding to each encoded open reading frame (ORF). These include two

61 large ORFs, ORF1a and ORF1ab, encoded by the complete viral genome and expressed

62 upon cell entry. Other subgenomic mRNAs are generated through a mechanism termed

63 discontinuous extension of minus strands, encoding structural proteins (spike protein (S),

64 envelope protein (E), membrane protein (M) and nucleocapsid protein (N)) and accessory

65 proteins (3a, 6, 7a, 7b, 8 and 10). The subgenomic mRNAs have a common 5′ leader

66 sequence, near-identical to that located in the 5′-UTR of the viral genome; the transcription

67 mechanism repositions the 5′ leader sequence upstream of ORFs, with each translation

68 start site being located at the primary position for ribosome scanning. These ORFs, although

69 annotated, are yet to be shown as expressed experimentally. Standard sequencing

70   technologies are unable to produce reads representing (i) complete RNA viral genomes or

71   (ii) subgenomic mRNAs needed to verify annotated ORFs, as these methods generate short

72   reads and have a reliance on amplification to generate complementary DNA (cDNA)

73   sequences.

74

75   To define the architecture of the coronaviral transcriptome, a recently established direct RNA

76   sequencing approach was applied, using a highly parallel array of nanopores [16]. In brief,

77   nucleic acids were prepared from cell culture material with high levels of SARS-CoV-2

78   growth, this being expected to include examples of both genomic mRNA and transcripts

79   corresponding to each ORF. These were sequenced with Oxford Nanopore Technologies,

80   including poly(T) adaptors and an R9.4 flowcell on a GridION platform (Oxford Nanopore

81   Technologies). Through this approach, the electronic current is measured as individual

82   strands of RNA translocate through a nanopore, with derived signal-space data basecalled

83   to infer the corresponding nucleobases. As a comparator, virion material of SARS-CoV-2

84   was also prepared and sequenced through this approach, with complete viral genome

85   sequences expected to predominate rather than subgenomic transcripts.

86

87   The cellular-derived material was used to generate 680,347 reads, comprising 860Mb of

88   sequence information (BioProject PRJNA608224). Aligning to the genome of the cultured

89   SARS-CoV-2 isolate (MT007544.1), a subset of reads were attributed to coronavirus

90   sequences (28.9%), comprising 367Mb of sequence distributed across the 29,893 base

91   genome. Of these, a number had lengths >20,000 bases, capturing the majority of the

92   SARS-CoV-2 genome on a single molecule. This direct RNA sequencing approach

93   generated an average 12,230 fold coverage of the coronaviral genome, biased towards

94   sequences proximal to the polyadenylated 3' tail; coverage ranged from 34 fold to >160,000

95   fold (Figure 1B), reflecting the higher abundance of subgenomic mRNAs carrying these

96   sequences, as well as the directional sequencing from the polyadenylated 3' tail. The virion

97   material generated fewer reads, and included a calibration standard added during library

98   preparation (430,923 reads, BioProject PRJNA608224).

99

100  Many features of SARS-CoV-2 biology are captured in these direct RNA sequence data,

101  including the transcriptome, as well as RNA base modifications or 'epitranscriptome'. To

102  define the transcriptome, the shared 5' leader sequence was used as a marker to identify

103  intact transcripts, these corresponding to subgenomic mRNAs and having a low abundance

104  in the virion-derived data (Supplementary Figures 1 and 2). In SARS-CoV-2, we identify

105  eight major viral mRNAs in addition to the viral genome; each annotated gene was observed

106  as a distinct subgenomic mRNA, outside of ORF7b and ORF10 (Figure 1C, Supplementary

107  Table 1). In SARS, ORF7a and ORF7b are encoded on a shared subgenomic mRNA, with

108  translation of ORF7b being achieved through ribosome leaky scanning, explaining the

109  absence of a dedicated ORF7b-encoding subgenomic mRNA [17]. There is however no

110  satisfactory explanation for the absence of an ORF10-encoding subgenomic mRNA.

111

112  ORF10 is the last predicted coding sequence upstream of the poly-A sequence, and the

113  shortest of the predicted coding sequences at 117 bases in length. ORF10 also has no

114  annotated function, and the putative encoded peptide does not appear in SARS-CoV-2

115  proteomes [18,19] or have a homolog in the SARS-CoV-1 proteome (Proteome ID:

116  UP000000354 [20]). Subgenomic mRNAs corresponding to ORF10 are not identifiable in our

117  reads (Supplementary Figure 3). These data suggest that the sequence currently annotated

118  as ORF10 does not have a protein coding function in SARS-CoV-2. Ongoing molecular

119  evolution at this locus should be considered in light of this finding.

120

121  Instead of encoding a protein sequence, the locus annotated as ORF10 immediately

122  upstream of the 3' UTR may act itself or as a precursor of other RNAs in the regulation of

123  gene expression, replication or modulating translation efficiency or cellular antiviral

124  pathways; the 3' UTR of coronaviruses contains domains critical for regulating viral RNA

125  synthesis and other aspects of viral biology [21]. An initial region of the 3' UTR appears

126  essential for viral replication, and an area further 3' includes the stem-loop II-like motif (s2m)

127  a feature conserved in SARS-CoV-2 and other coronaviruses [22,23], the s2m having a

128  proposed role in recruiting host translational machinery [24]. A small number of cell culture-

129  derived SARS-CoV-2 genomes carry a shared deletion at an area of the 3' UTR including an

130  aspect of the s2m (Supplementary Figure 4), this parallel molecular evolution further

131  suggesting the region may have functional roles *in vivo*.

132

133  An analysis of transcript breakpoints further illustrates the potential for 5' UTR positions

134  outside of the canonical leader sequence to enable transcript production, with low-frequency

135  non-canonical variations in mRNA splice co-ordinates (Figure 2, Supplementary Figures 5

136  and 6). Low frequency variants may be generated during the preparation of nucleic acids for

137  sequencing, with the rate of chimeric read formation being unknown; this could be explored

138  through analysis of *in vitro* transcribed RNA control material.

139

140    In addition to RNA modifications such as the methylation of the 5' cap structure and

141    polyadenylation of the 3' terminus needed for efficient translation of coding sequences, other

142    RNA modifications may have functional roles in SARS-CoV-2 [25]. A range of modifications

143    are identifiable using direct RNA sequence data [16,25]; our available SARS-CoV-2 direct

144    RNA sequence data providing adequate coverage to confidently call specific modifications.

145    Through analysis of signal-space data, we identified 42 positions with predicted 5-

146    methylcytosine modifications, appearing at consistent positions between subgenomic

147    mRNAs (Supplementary Figure 7, and Supplementary Table 2). In other positive ssRNA

148    viruses, RNA methylation can change dynamically during the course of infection [26],

149    influencing host-pathogen interaction and viral replication. Other modifications may become

150    apparent once training datasets are available for direct RNA sequence data, with little known

151    of the epitranscriptomic landscape of coronaviruses [25,27].

152

153    As well as investigating the above assumed features of SARS-CoV-2 genetics, sequence

154    data also enable an estimate of the molecular evolutionary rate, with globally sourced

155    genome sequences being shared and publicly available. Evolutionary rate estimates from

156    other coronaviruses such as Middle East Respiratory Syndrome (MERS) are not necessarily

157    applicable here, particularly because MERS had multiple independent introductions into

158    humans [28-30]. To estimate the evolutionary rate and time of origin of the SARS-CoV-2

159    outbreak, we carried out Bayesian phylogenetic analyses using a curated set of 122 high

160    quality publicly available SARS-CoV-2 genome sequences, each having a known collection

161    date (Figure 3, Supplementary Table 3). The sampling times were sufficient to calibrate a

162    molecular clock and infer the evolutionary rate and timescale of the outbreak using a

163    Bayesian approach; the evolutionary rate of SARS-CoV-2 was estimated to be $1.20 \times 10^{-3}$

164    substitutions/site/year (95% HPD $8.91\times10^{-4}$ - $1.52\times10^{-3}$), and the of time of origin to be late

165    November 2019 (95% HPD August 2019, December 2019), which is in agreement with

166    epidemiological evidence and other recent analyses (Figure 3A) [1-4, 31, 32].

167

168    A further set of 66 high quality genomes collected earlier in the outbreak (Supplementary

169    Table 4), and maximal diversity data set from all data available in GISAID to March 28th to

170    show the utility of capturing varying degrees of genetic diversity (Supplementary Table 5).

171    This Bayesian approach demonstrated improved precision in estimates of evolutionary rates

172    using our dataset with highest genetic diversity. This may be explained by the stochastic

173    variation typical in data from early in an outbreak having a smaller impact as the virus

174    accumulates genetic variation. These results are also supported by root-to-tip regression, a

175    visual assessment of the degree of clocklike evolution in the data (Supplementary Figures 8

176    and 9). The evolutionary rate generated the high diversity set of 100 genomes was 1.56 ×

177    $10^{-3}$ substitutions/site/year (95% HPD $1.09 \times 10^{-3}$ - $2.05 \times 10^{-3}$), whereas that based on 66

178    genomes was 1.16 × $10^{-3}$ substitutions/site/year (95% HPD $6.32 \times 10^{-4}$ - $1.69 \times 10^{-3}$). Our

179    estimate of the evolutionary rate of SARS-CoV-2 is in line with those of other coronaviruses

180    (Figure 3B), and the low genomic diversity and recent timescale of the outbreak support a

181    recently occurring, point-source transfer to humans.

182

183    Other phylodynamic inferences may soon become possible for SARS-CoV-2, as further

184    genomic data becomes available and the sampling rate becomes more consistent. The

185    current distribution of sampling times (Supplementary Figure 8) appears to be prohibitive to

186    phylodynamic inference of the SARS-CoV-2 effective population size ($N_e$, not included here).

187    Although a required threshold of genomes to allow such phylodynamic investigation may

188    have been crossed, the temporal spread of these isolates may differ too much to satisfy

189    constant sampling assumptions underlying many phylodynamic skyline approaches inferring

190    $N_e$ over time. Again, as sampling continues a more consistent rate of sampling is likely to

191    emerge, allowing such analyses.

192

193    Insights are provided on the molecular biology of SARS-CoV-2, revealed through the use of

194    direct RNA sequence and publicly available data. The rapid sharing of these and other

195    genetic data support the global response effort and represents an inflection point for

196    communicable diseases and genomic epidemiology, with complete data shared openly and

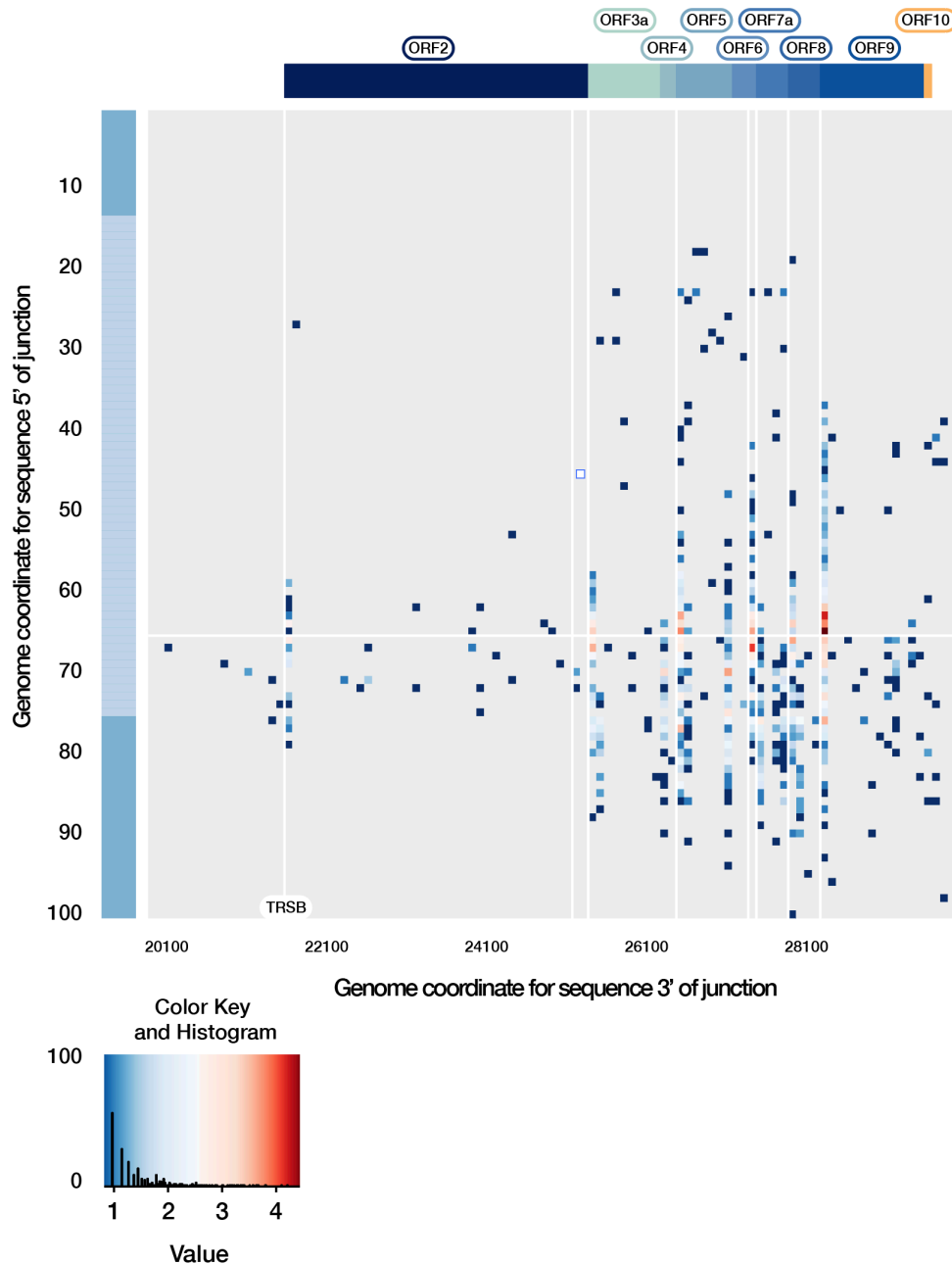197    rapidly between academic and public health groups.

198



199

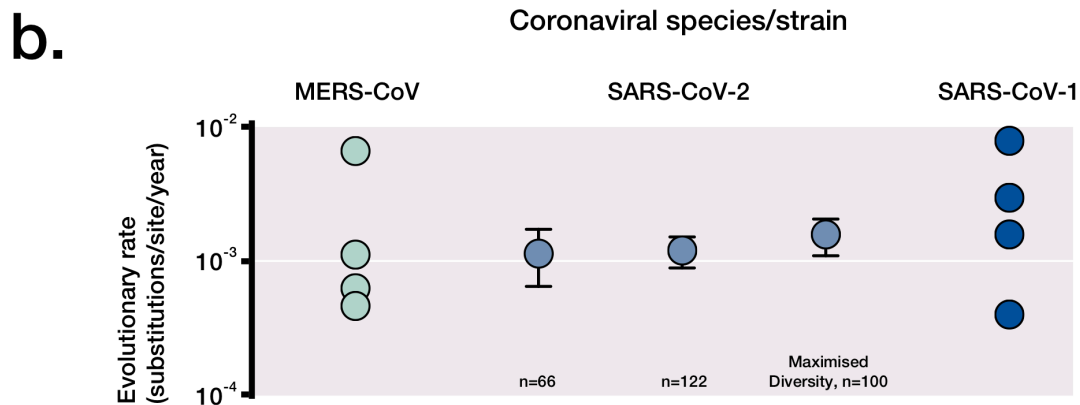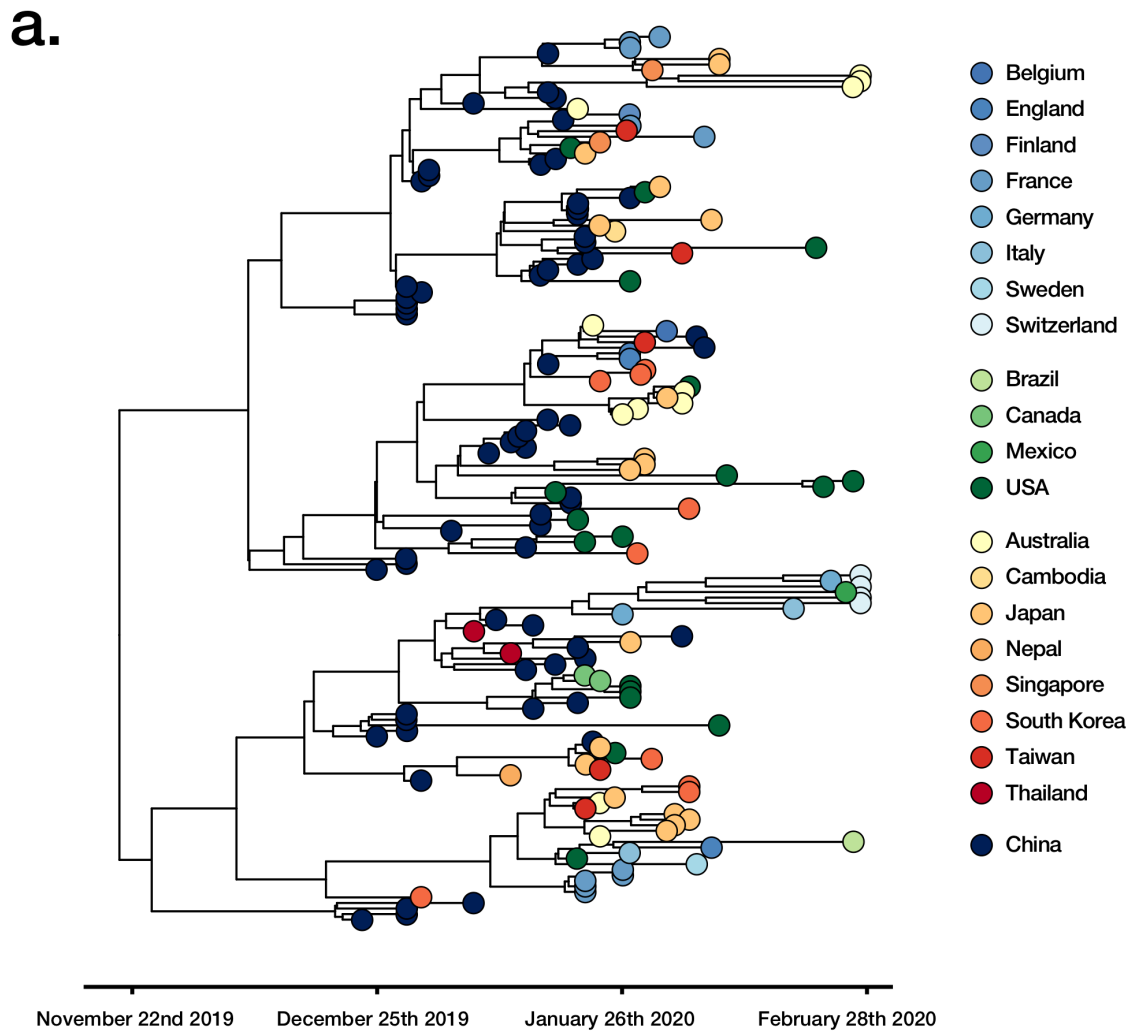200 **Figure 1. SARS-CoV-2 genetics and transcriptome architecture.**

201 A) Schematic of the early stages of SARS-CoV-2 cell entry and transcript production,

202 including *in vivo* synthesis of positive sense genome-length RNA molecules and subgenomic

203 mRNAs. B) Read coverage of direct RNA reads from cell-culture material, aligned to the

204 local SARS-CoV-2 genome (29,893 bases), showing a bias towards the 3' polyadenylated

205 end. C) Read length histogram, showing subgenomic mRNAs attributed to coding

206 sequences.

**Figure 2. Breakpoint analysis of the SARS-CoV-2 transcriptome.**

Direct RNA reads carrying a breakpoint relative to the 5' leader sequence are shown, representing potentially viable transcripts. These breakpoints are localised at the same position on the leader sequence (positions 62-68), and on the 3' to predicted transcription regulating sequences in the body of the genome (TRS-Bs, highlighted by vertical weight lines), generating common subgenomic mRNAs. Of note, many low frequency breakpoints are detected, although few near the sequence currently annotated as ORF10. The key shows the distribution of transcript breakpoints. Colour is matched to a 'value' measuring the number of reads with break points at that position, log10-scaled. The histogram component illustrates the number of transcripts with a given abundance value.

**Figure 3. Assessment of viral evolutionary rate and outbreak timing with SARS-CoV-2-specific data.** A) A timed highest clade-credibility phylogenetic tree of curated SARS-CoV-2 genomes as inferred in BEAST. B) Comparison of SARS-CoV-2 rate estimates with varying datasets, and previously published estimates of other coronaviruses.

223   References

224

225   [1]   World Health Organization. Pneumonia of unknown cause — China. 2020

226         (https://www.who.int/csr/don/05-january-2020-pneumonia-of-unkown-cause-

227         china/en/).

228

229   [2]   United Nations Development Programme – March 2020. UNDP support for

230         coronavirus-affected countries goes beyond health.

231         https://www.undp.org/content/undp/en/home/blog/2020/undp-support-for-

232         coronavirus-affected-countries-goes-beyond-heal.html

233

234   [3]   Dong, Ensheng, Hongru Du, and Lauren Gardner. "An interactive web-based

235         dashboard to track COVID-19 in real time." The Lancet Infectious Diseases (2020).

236         DOI: 10.1016/S1473-3099(20)30120-1

237

238   [4]   World Health Organization. Statement on the second meeting of the International

239         Health Regulations (2005) Emergency Committee regarding the outbreak of novel

240         coronavirus (2019-nCoV). January 30, 2020 (https://www.who.int/news-

241         room/detail/30-01-2020-statement-on-the-second-meeting-of-the-international-

242         health-regulations-(2005)-emergency-committee-regarding-the-outbreak-of-novel-

243         coronavirus-(2019-ncov).

244

245   [5]   Lu, Roujian, et al. "Genomic characterisation and epidemiology of 2019 novel

246         coronavirus: implications for virus origins and receptor binding." The Lancet (2020).

247         DOI: 10.1016/S0140-6736(20)30251-8

248

249   [6]   Peiris JS, Yuen KY, Osterhaus AD, Stöhr K. The severe acute respiratory syndrome.

250         New England Journal of Medicine. 2003 Dec 18;349(25):2431-41. DOI:

251         10.1056/NEJMra032498

252

253   [7]   Assiri A, McGeer A, Perl TM, Price CS, Al Rabeeah AA, Cummings DA, Alabdullatif

254         ZN, Assad M, Almulhim A, Makhdoom H, Madani H. Hospital outbreak of Middle East

255         respiratory syndrome coronavirus. New England Journal of Medicine. 2013 Aug

256         1;369(5):407-16. DOI: 10.1056/NEJMoa1306742

257

258  [8]   Perlman, S. Another decade, another coronavirus. New England Journal of Medicine.
259        2020 February 20; 382:760-762. DOI: 10.1056/NEJMe2001126

260

261  [9]   Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data–from
262        vision to reality. Eurosurveillance. 2017 Mar 30;22(13). https://www.gisaid.org/

263

264  [10]  Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P,
265        Bedford T, Neher RA. Nextstrain: real-time tracking of pathogen evolution.
266        Bioinformatics. 2018 Dec 1;34(23):4121-3. https://nextstrain.org/ncov

267

268  [11]  Andersen KG, Rambaut A, Lipkin WI, Holmes EC, Garry RF. The Proximal Origin of
269        SARS-CoV-2. Virological, accessed on 27/02/2020. http://virological.org/t/the-
270        proximal-origin-of-sars-cov-2/398

271

272  [12]  Rambaut A. Phylodynamic Analysis | 129 genomes | 24 Feb 2020. Virological,
273        accessed 27/02/2020. http://virological.org/t/phylodynamic-analysis-129-genomes-
274        24-feb-2020/356

275

276  [13]  Yount B, Curtis KM, Fritz EA, Hensley LE, Jahrling PB, Prentice E, Denison MR,
277        Geisbert TW, Baric RS. Reverse genetics with a full-length infectious cDNA of severe
278        acute respiratory syndrome coronavirus. Proceedings of the National Academy of
279        Sciences. 2003 Oct 28;100(22):12995-3000. DOI: 10.1073/pnas.1735582100

280

281  [14]  Brian DA, Baric RS. Coronavirus genome structure and replication. InCoronavirus
282        replication and reverse genetics 2005 (pp. 1-30). Springer, Berlin, Heidelberg.

283

284  [15]  Chen Y, Cai H, Xiang N, Tien P, Ahola T, Guo D. Functional screen reveals SARS
285        coronavirus nonstructural protein nsp14 as a novel cap N7 methyltransferase.
286        Proceedings of the National Academy of Sciences. 2009 Mar 3;106(9):3484-9. DOI:
287        10.1073/pnas.0808790106

288

289  [16]  Garalde DR, Snell EA, Jachimowicz D, Sipos B, Lloyd JH, Bruce M, Pantic N,
290        Admassu T, James P, Warland A, Jordan M. Highly parallel direct RNA sequencing
291        on an array of nanopores. Nature methods. 2018 Mar;15(3):201. DOI:
292        10.1038/nmeth.4577

293

294   [17]   Schaecher SR, Mackenzie JM, Pekosz A (2007). The ORF7b Protein of Severe
295          Acute Respiratory Syndrome Coronavirus (SARS-CoV) Is Expressed in Virus-
296          Infected Cells and Incorporated into SARS-CoV Particles. Journal of Virology, 81(2),
297          718–731. DOI: 10.1128/JVI.01691-06

298

299   [18]   Bojkova D, Klann K, Koch B, Widera M, Krause D, Ciesek S, Cinatl J, Münch C
300          (2020) SARS-CoV-2 infected host cell proteomics reveal potential therapy targets.
301          DOI: 10.21203/rs.3.rs-17218/v1

302

303   [19]   Davidson AD, Williamson MK, Lewis S, Shoemark D, Carroll MW, Heesom K,
304          Zambon M, Ellis J, Lewis PA, Hiscox JA, Matthews DA (2020) Characterisation of the
305          transcriptome and proteome of SARS-CoV-2 using direct RNA sequencing and
306          tandem mass spectrometry reveals evidence for a cell passage induced in-frame
307          deletion in the spike glycoprotein that removes the furin-like cleavage site. DOI:
308          10.1101/2020.03.22.002204

309

310   [20]   He R, Dobie F, Ballantine M, Leeson A, Li Y, Bastien N, Cutts T, Andonov A, Cao J,
311          Booth TF, Plummer FA, Tyler S, Baker L, Li X (2004) . Analysis of multimerization of
312          the SARS coronavirus nucleocapsid protein. Biochemical and Biophysical Research
313          Communications, 316(2), 476–483. DOI: 10.1016/j.bbrc.2004.02.074

314

315   [21]   Goebel SJ, Hsue B, Dombrowski TF, Masters PS. Characterization of the RNA
316          components of a putative molecular switch in the 3' untranslated region of the murine
317          coronavirus genome. J Virol. 2004 Jan;78(2):669-82. DOI: 10.1128/jvi.78.2.669-
318          682.2004

319

320   [22]   Tengs T, Kristoffersen AB, Bachvaroff TR, Jonassen CM. A mobile genetic element
321          with unknown function found in distantly related viruses. Virol J. 2013 Apr 25;10:132.
322          DOI: 10.1186/1743-422X-10-132

323

324   [23]   Rangan R, Zheludev IN, Das R. RNA genome conservation and secondary structure
325          in SARS-CoV-2 and SARS-related viruses. bioRxiv preprint 2020 DOI:
326          10.1101/2020.03.27.012906.

327

328   [24]   Robertson MP, Igel H, Baertsch R, Haussler D, Ares M, Scott WG. The Structure of a
329          Rigorously Conserved RNA Element within the SARS Virus Genome. PLoS Biol.
330          2005 Jan; 3(1): e5. DOI: 10.1371/journal.pbio.0030005
331

332   [25]   Viehweger A, Krautwurst S, Lamkiewicz K, Madhugiri R, Ziebuhr J, Hölzer M, Marz
333          M. Direct RNA nanopore sequencing of full-length coronavirus genomes provides
334          novel insights into structural variants and enables modification analysis. Genome
335          research. 2019 Sep 1;29(9):1545-54. DOI: 10.1101/gr.247064.118
336

337   [26]   Lichinchi G, Zhao BS, Wu Y, Lu Z, Qin Y, He C, Rana TM. Dynamics of human and
338          viral RNA methylation during Zika virus infection. Cell host & microbe. 2016 Nov
339          9;20(5):666-73. DOI: 10.1016/j.chom.2016.10.002
340

341   [27]   Denison MR, Graham RL, Donaldson EF, Eckerle LD, Baric RS. Coronaviruses: an
342          RNA proofreading machine regulates replication fidelity and diversity. RNA biology.
343          2011 Mar 1;8(2):270-9. DOI: 10.4161/rna.8.2.15013
344

345   [28]   Chinese SARS Molecular Epidemiology Consortium. Molecular evolution of the
346          SARS coronavirus during the course of the SARS epidemic in China. Science. 2004
347          Mar 12;303(5664):1666-9. DOI: 10.1126/science.1092002
348

349   [29]   Dudas G, Carvalho LM, Rambaut A, Bedford T. MERS-CoV spillover at the camel-
350          human interface. Elife. 2018 Jan 16;7:e31257. DOI: 10.7554/eLife.31257
351

352   [30]   Cotten M, Watson SJ, Kellam P, Al-Rabeeah AA, Makhdoom HQ, Assiri A, Al-Tawfiq
353          JA, Alhakeem RF, Madani H, AlRabiah FA, Al Hajjar S. Transmission and evolution
354          of the Middle East respiratory syndrome coronavirus in Saudi Arabia: a descriptive
355          genomic study. The Lancet. 2013 Dec 14;382(9909):1993-2002. DOI:
356          10.1016/S0140-6736(13)61887-5
357

358   [31]   Andersen K. Clock and TMRCA based on 27 genomes. Virological, accessed on
359          27/02/2020. http://virological.org/t/clock-and-tmrca-based-on-27-genomes/347
360

361   [32]   Bedford, T. Phylodynamic estimation of incidence and prevalence of novel
362          coronavirus (nCoV) infections through time. Virological, accessed on 27/02/2020.

363        http://virological.org/t/phylodynamic-estimation-of-incidence-and-prevalence-of-novel-

364        coronavirus-ncov-infections-through-time/391
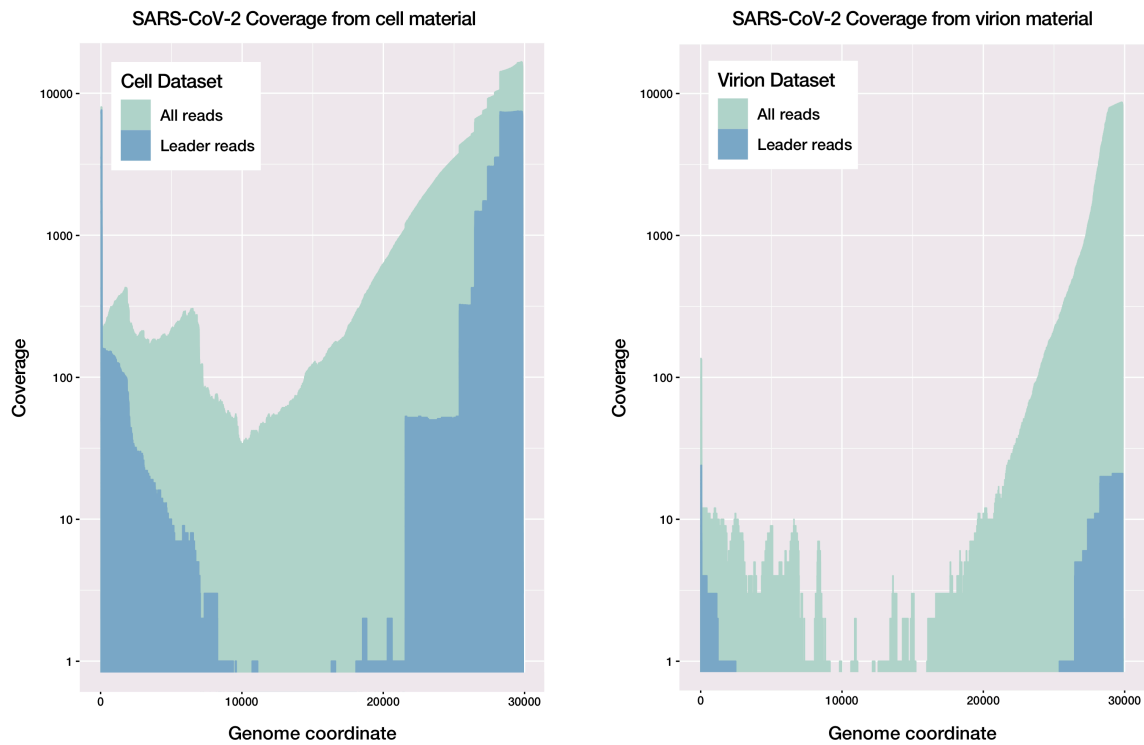
365

366

367    Acknowledgements

368

380

381

382    Supplementary Figure 1

383    Native RNA sequence coverage of the SARS-CoV-2 genome for cell-culture and virion-

384    derived material. A) Coverage of the SARS-CoV-2 genome for the cell-culture dataset, for all

385    reads and for those predicted to be intact mRNA transcripts or 'leader reads', showing an

386    abundance of such transcripts. B) Coverage of the SARS-CoV-2 genome for the virion-

387    derived dataset, showing a relative paucity of intact mRNA transcripts.

388

389      Supplementary Figure 2

390      Distribution of native RNA reads between intact transcripts ('leader reads') and other partial

391      transcripts and genomic sequences. Intact transcripts include the leader sequence at the 5'

392      and a polyadenylated 3' end (A and E, for cell-culture material and virion material

393      respectively), while other partial transcripts or genomes either contain a leader sequence

394      and lack an appropriate 3' sequence (B and F), vice versa (C and D), or lack both a leader

395      sequence and an expected 3' end.

396

397

398 Supplementary Figure 3

399 Absence of observed coding potential for ORF10 in SARS-CoV-2. A) Read length

400 histogram, showing subgenomic mRNAs attributed to coding sequences, with the area

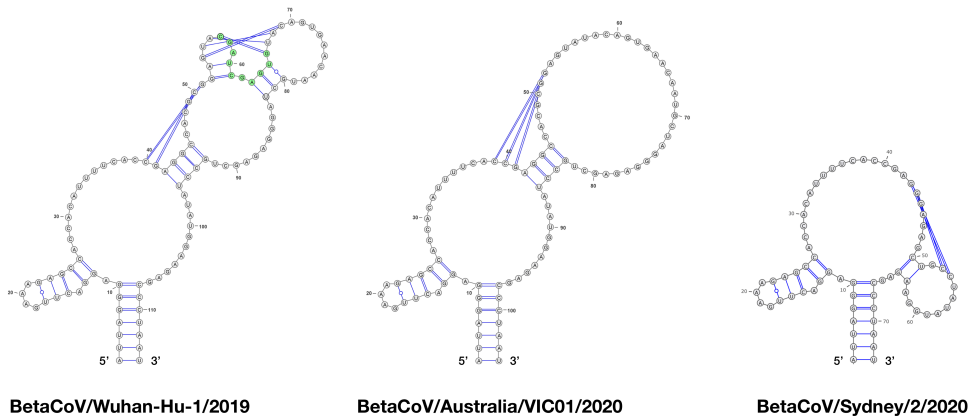401 highlighted shown in detail in a second panel. B) Read length histogram, showing read

402 counts of lengths corresponding to those of the ORF10 subgenomic mRNA (~360 bases), if

403 present in the dataset. Of the <500 base reads shown, none align to ORF10.

**a.** Alignment of the SARS-CoV-2 3' UTR for selected isolates



**b.** Schematic of predicted pseudoknots in the SARS-CoV-2 3' UTR affected by culture-derived deletions



BetaCoV/Wuhan-Hu-1/2019     BetaCoV/Australia/VIC01/2020     BetaCoV/Sydney/2/2020

404

405    Supplementary Figure 4

406    Structured RNAs in the SARS-CoV-2 3' UTR. A) An alignment of SARS-CoV-2 3' UTR

407    sequences, including the original Wuhan-Hu-1 sourced from Wuhan, China and considered

408    the reference genome for the outbreak, and two examples of cultured SARS-CoV-2 isolates

409    exhibiting deletions in a shared 3' UTR region predicted to form a pseudoknot structure. B)

410    Predicted pseudoknot structure of the SARS-CoV-2 3'UTR affected by the above culture-

411    derived deletions.

412

413    Supplementary Figure 5

414    Extended breakpoint analysis of the SARS-CoV-2 transcriptome. The genome coordinates

415    3' of the breakpoint are extended to include potential 3' sequences positioned between 1-

416    10,001 of the genome. This highlights low frequency breakpoints, increasing in frequency

417    near the sequence annotated as ORF10 and the 3' end of the genome. The key shows the

418    distribution of transcript breakpoints. Colour is matched to a 'value' measuring the number of

419    reads with break points at that position, log10-scaled. The histogram component illustrates

420    the number of transcripts with a given abundance value.

421

Supplementary Figure 6

ORF-specific breakpoint analyses. The corresponding breakpoints for each currently annotated ORF in the SARS-CoV-2 genome are shown (A-I), highlighting a canonical breakpoint for ORFs with a corresponding subgenome mRNA, and a low frequency of non-canonical splice sites often centred on a canonical site. Of note, low frequency splice sites can be seen for an area between ORF 7a and ORF8, likely corresponding to ORF7b (G). There is an absence of splice sites for ORF in this dataset (I). The key shows the distribution of transcript breakpoints. Colour is matched to a 'value' measuring the number of reads with break points at that position, log10-scaled. The histogram component illustrates the number of transcripts with a given abundance value.

432

433    Supplementary Figure 7

434    Subgenomic mRNA abundance and predicted sites of modification. A) Coverage of relevant

435    coding sequences achieved by alignment of subgenomic mRNAs to the SARS-CoV-2

436    genome (log scale). Red lines indicate the first base of each coding sequence from ORF2-

437    10. B) Schematic of relevant annotated coding sequences. C) Position of predicted m5C

438    positions in subgenomic mRNAs. Dark blue lines indicate positions predicted to have >90%

439    base modification; light blue lines indicate positions predicted to have between 50% and

440    90% base modification.

441

442

443    Supplementary Figure 8

444    Assessment of SARS-CoV-2 phylogenetics and viral evolutionary rate based on 66 early

445    genomes made publicly available. A) A timed highest clade-credibility phylogenetic tree of

446    curated SARS-CoV-2 66 genomes as inferred in BEAST. B) Comparison of the SARS-CoV-

447    2 rate estimate for the n=66 set and previously published estimates of other coronaviruses.

448

449    Supplementary Figure 9

450    Time-to-most-common-recent-ancestor (TMRCA)and root-to-tip regression of both early and

451    maximally diverse SARS-CoV-2 genome datasets. A) TMRCA and root-to-tip regression of

452    122 high quality complete SARS-CoV-2 genomes made available early in the pandemic.

453    B) TMRCA and root-to-tip regression of 100 maximally diverse SARS-CoV-2 genomes,

454    selected from the first 700 genomes made publicly available.

Distribution of SARS-CoV-2 sampling times

455

456 Supplementary Figure 10

457 Distribution of SARS-CoV-2 sampling times used to generate publicly available genomes.

458 The distribution has notable deviations from an expected exponential growth in the number

459 of genomes available, such as in mid-February, with constant sampling being an underlying

460 assumption for many phylodynamic skyline approaches inferring effective population size.

461    Methods

462    Samples for direct RNA sequencing

463    The SARS-CoV-2 material was prepared from the first Australian case of COVID-2019

464    (Australia/VIC01/2020), maintained in cell culture. In brief, African green monkey kidney

465    cells expressing the human signalling lymphocytic activation molecule (SLAM; termed

466    Vero/hSLAM cells accordingly) with associated SARS-CoV-2 infection were grown at 37°C

467    at 5% $CO_2$ in media consisting of 10 mL Earle's minimum essential medium, 7% FBS

468    (Bovogen Biologicals, Keilor East, Aus), 2 mM L-Glutamine, 1 mM Sodium pyruvate, 1500

469    mg/L sodium bicarbonate, 15 mM HEPES and 0.4 mg/ml geneticin in 25cm² flasks. This

470    isolate is to the best of our knowledge typical for SARS-CoV-2 isolates, with the genome of

471    the cultured isolate (MT007544.1) having three single nucleotide variants (T19065C,

472    T22303G, G26144T) relative to the SARS-CoV-2 Wuhan-Hu-1 reference genome

473    (MN908947.3), and a 10 base deletion in the 3' UTR. Both the T22303G and 3' UTR

474    variants have been confirmed as culture-derived through Sanger sequencing of clinical and

475    culture material, and do not appear in the earlier virion-derived data.

476

477    Nucleic acids were prepared from clarified cell-free supernatant (reflecting virion material)

478    and infected cell culture material (representing actively transcribed and translated viral

479    material), following inactivation with linear acrylamide and ethanol. RNA was extracted from

480    100µl of supernatant and a modest pellet for the cell-culture material (~200mg) respectively,

481    using manually prepared wide-bore pipette tips and minimal steps to maintain RNA length

482    for long read sequencing, and a QIAamp Viral RNA Mini Kit (Qiagen, Hilden, Germany).

483    Carrier RNA was not added to Buffer AVL, with 1% linear acrylamide (Life Technologies,

484    Carlsbad, CA, USA) added instead.  Wash buffer AW1 was omitted from the purification

485    stage, with RNA eluted in 50 µl of nuclease free water, followed by DNase treatment with

486    Turbo DNase (Thermo Fisher Scientific, Waltham, MA, USA) 37°C for 30 min.  RNA was

487    cleaned and concentrated to 10 µl using the RNA Clean & Concentrator-5 kit (Zymo

488    Research, Irvine, CA, USA), as per manufacturer's instructions.

489

490    Nanopore sequencing of direct RNA

491    Prepared RNA (~1µg) was carried into a direct RNA sequence library preparation with the

492    Oxford Nanopore DRS protocol (SQK- RNA002, Oxford Nanopore Technologies) following

493    the manufacturer's specifications, with addition of the control RNA in the virion sample.

494    Libraries were loaded on R9.4 flow cells and sequenced on a GridION device for the cell-

495    derived material and a MinION device for the virion-derived material (Oxford Nanopore

496 Technologies), and sequenced for 40 hours. Signal-space data was used to generate

497 nucleobase sequences ('basecalled') using Guppy, either as a standalone program or as

498 ont-guppy-for-gridion 3.0.6. Both signal-space and basecalled read data are available at

499 BioProject PRJNA608224. It should be noted that non-polyadenylated RNAs are not

500 expected to be detected with this approach.

501

502 Characterisation of SARS-CoV-2 transcriptome architecture

503 Direct RNA reads passing the above given quality thresholds were aligned to the genome of

504 the cultured Australian SARS-COV-2 isolate (MT007544.1), with parallel and concordant

505 analyses in Geneious Prime (2019.2.1, [M1]) and minimap2 v 2.11 using the "spliced" preset

506 [M2]. Coverage statistics were determined from the resulting read alignments. To identify

507 complete subgenomic mRNAs, reads were aligned to a 62 base SARS-COV-2 leader

508 sequence (5'ACCUUCCCAGGUAACAAACCAACCAACUUUCGAUCUCUUGUAGAU

509 CUGUUCUCUAAACGAAC), with reads aligning to the leader sequence being pooled and

510 visualized in a length histogram. Significant peaks were identified visually and confirmed

511 with a smoothed z-score algorithm. Reads captured in this binning-by-length strategy were

512 re-aligned to the reference genome using the above methods and visualized in Tablet [M3].

513 Subgenome bins were refined to remove reads which did not originate at the 3' poly-A tail as

514 expected for intact subgenomic mRNAs, or which had leader sequences at least 10bp

515 longer than expected. Subgenome bins were re-aligned, with coverage calculated in

516 SAMtools [M4], and plotted using ggplot2 [M5] in R [M6]. Breakpoints in mRNAs were

517 determined with CIGAR string manipulation; any given spliced region longer than 100bp

518 (represented by Ns in the CIGAR string after aligning with minimap2) was regarded as a

519 spliced transcript and the 5' and 3' genome co-ordinates of the breakpoint were recorded for

520 analysis. The IPKnot webserver [M7] was used to predict the RNA secondary structures,

521 and the VARNA visualization applet [M8] to produce schematics.

522

523 As an alternate method of defining the SARS-CoV-2 transcriptome, reads carrying a

524 breakpoint relative to the 5' leader sequence are shown, representing potentially viable

525 transcripts. This was determined through CIGAR string manipulation. Any spliced region

526 longer than 100bp (represented by Ns in the CIGAR string after aligning with minimap2) was

527 regarded as a spliced transcript and the 5' and 3' genome co-ordinates of the breakpoint

528 were recorded for analysis. Locations of Transcription Regulating Sequence in the body of

529 the genome (TRS-B) were determined with a Position Weight Matrix (PWM) search. Portions

530 of reads aligning to the conserved TRS in the leader sequence (TRS-L) were transformed

531     into a count matrix, which was then passed into the FIMO program version 5.5.1 for motif

532     detection with NRDB as the background distribution [M9]. Detected TRS-B sites are plotted

533     alongside breakpoint heatmaps.

534

535     Data availability

536     All signal-space (fast5) and basecalled data (fastq) generated in this work are publicly

537     available on the sequence read archive (SRA), as part of the BioProject PRJNA608224 (See

538     Supplementary Table 6 for relevant accession numbers).

539

540     Identification of 5mC methylation

541     Nanopore sequencing preserves *in vivo* base modifications and enables their detection from

542     raw voltage signal information. In brief, the signal-space fast5 files corresponding to

543     identified subgenomic mRNAs were assessed to identify signal changes corresponding to

544     5mC methylation. These were first retrieved using the fast5_fetcher_multi function in

545     SquiggleKit [M10]. Reads were processed to align raw signal with basecalled sequence data

546     using Tombo v1.5 [https://github.com/nanoporetech/tombo]. Canonical reference sequences

547     were made for each subgenomic mRNAs, with the binned fast5 files input into the

548     detect_modifications function, with 5mC as the alternate-model parameter. Outputs were

549     converted to dampened_fraction wiggle files and exported for visualization and analysis.

550

551     Assessment of publicly available proteomes

552     Proteomic datasets were downloaded from the PRIDE proteomic database [M11] (Pride

553     accession: PXD017710) and processes using Maxquant (1.6.3.4 [M12]) allowing semi-

554     specific free-N-terminus tryptic as a protease specificity. Quantitation was set to TMT11 plex

555     labelling and human proteome (Uniprot: UP000005640) and SARS-CoV2 (build in house)

556     databases were used. Additional searches were made of the SARS-CoV reference

557     proteome (Proteome ID: UP000000354), and recently available SARS-CoV-2 proteomic

558     data [19].

559

560     SARS-CoV-2 Phylogenetics

561     In order to estimate the evolutionary rate and time of origin of SARS-CoV-2, we carried out

562     phylogenetic analyses in BEAST v1.101 [M13 on three datasets. The first dataset included

563     66 high quality genomes available up to February 10th 2020, the second consisted of 122

564     available up to February 24th 2020, both from GISAID and GenBank (Supplementary Table

565     3). A third maximal diversity dataset (n=100) was included to demonstrate the utility of

566 capturing varying degrees of genetic diversity, these genomes being selected from the first

567 700 genomes available on GISAID and maximised phylogenetic diversity achieved using

568 Treemer [M14] (Supplementary Table 4). Temporal signal was assessed using BETS [M15].

569 Initially we determined whether the evolutionary signal and time over which the genome data

570 were collected was sufficient to calibrate the molecular clock, allowing for the evolutionary

571 rate and timescale of the outbreak to be inferred. The model selection approach from BETS

572 supported a strict molecular clock model with genome sampling times for calibration and a

573 coalescent exponential tree prior, which posits that the number of infected individuals grows

574 exponentially over time. We used the HKY+$\Gamma$ substitution model, and set the following priors

575 for key parameters:

- A continuous time Markov chain for the evolutionary rate

- A Laplace distribution with mean of 0 and scale of 100 for the growth rate

- An exponential distribution with mean of 1 for the effective population size.

579 A Markov chain Monte Carlo of length $10^7$ was set, sampling every $10^3$ steps, and assessed

580 sufficient sampling by verifying that the effective sample size for all parameters was at least

581 200 as determined in Tracer [M16], automatically discarding 10% of the burn in. We

582 summarised the posterior distribution of phylogenetic trees by selecting the highest clade

583 credibility tree alongside calculating posterior node probabilities and the distribution of node

584 ages. Comparison to other coronaviral evolutionary rates included studies [M17-24].

585

586 A root-to-tip regression usually produces lower evolutionary rate estimates than explicit

587 phylogenetic methods [M25], although is commonly used to inspect temporal signal in the

588 data. In our analyses, the data set that maximised phylogenetic diversity had a higher $R^2$

589 than that with 122 samples collected earlier, and an evolutionary rate that was more similar

590 to that obtained in BEAST. Although, both data sets had temporal signal according to BETS,

591 the root-to-tip regressions demonstrate that including more genetic diversity can produce

592 improved estimates, probably because stochastic variation has a stronger impact in smaller

593 data sets that are collected early in the outbreak.

594

595 [M1]   Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S,

596         Cooper A, Markowitz S, Duran C, Thierer T. Geneious Basic: an integrated and

597         extendable desktop software platform for the organization and analysis of sequence

598         data. Bioinformatics. 2012 Jun 15;28(12):1647-9.

599

600   [M2]   Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics. 2018
601         Sep 15;34(18):3094-100.
602

603   [M3]   Milne I, Bayer M, Stephen G, Cardle L, Marshall D. Tablet: visualizing next-
604         generation sequence assemblies and mappings. Bioinformatics. 2016 (pp. 253-268).
605         Humana Press, New York, NY.
606

607   [M4]   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
608         Durbin R. The sequence alignment/map format and SAMtools. Bioinformatics. 2009
609         Aug 15;25(16):2078-9.
610

611   [M5]   Wickham H. ggplot2: elegant graphics for data analysis. Springer; 2016 Jun 8
612

613   [M6]   R Core Team (2018). R: A language and environment for statistical computing. R
614         Foundation for Statistical Computing, Vienna, Austria.
615

616   [M7]   Sato, Kengo, et al. "IPknot: fast and accurate prediction of RNA secondary structures
617         with pseudoknots using integer programming." Bioinformatics. 2011 27.13: i85-i93.
618

619   [M8]   Darty, Kévin, Alain Denise, and Yann Ponty. "VARNA: Interactive drawing and
620         editing of the RNA secondary structure." Bioinformatics 2009 25.15: 1974.
621

622   [M9]   Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif.
623         Bioinformatics 2011 27: 1017-1018
624

625   [M10]  Vizcaino JA, Csordas A, del-Toro N, Dianes JA, Griss J, Lavidas I, et al. Update of
626         the PRIDE database and its related tools. Nucleic Acids Res. 2016;44(D1):D447-56.
627

628   [M11]  Cox J, Mann M. MaxQuant enables high peptide identification rates, individualized
629         p.p.b.-range mass accuracies and proteome-wide protein quantification. Nat
630         Biotechnol. 2008;26(12):1367-72. DOI: 10.1038/nbt.1511.
631

632   [M12]  Ferguson JM, Smith MA. SquiggleKit: A toolkit for manipulating nanopore signal
633         data. Bioinformatics. 2019 Dec 15;35(24):5372-3.
634

635  [M13]  Drummond AJ, Rambaut A. BEAST: Bayesian evolutionary analysis by sampling
636        trees. BMC evolutionary biology. 2007 Dec;7(1):214.
637

638  [M14]  Menardo F, Loiseau C, Brites D, Coscolla M, Gygli SM, Rutaihwa LK, Trauner A,
639        Beisel C, Borrell S, Gagneux S. Treemmer: a tool to reduce large phylogenetic
640        datasets with minimal loss of diversity. BMC Bioinformatics 2018 19(1), 164. DOI:
641        10.1186/s12859-018-2164-8
642

643  [M15]  Duchene S, Stadler T, Ho SY, Duchene DA, Dhanasekaran V, Baele G. Bayesian
644        Evaluation of Temporal Signal in Measurably Evolving Populations. bioRxiv. 2019
645        Jan 1:810697.
646

647  [M16]  Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior summarization in
648        Bayesian phylogenetics using Tracer 1.7. Systematic biology. 2018 Sep;67(5):901.
649

650  [M17]  Cotten M, Watson SJ, Kellam P, Al-Rabeeah AA, Makhdoom HQ, Assiri A, Al-Tawfiq
651        JA, Alhakeem RF, Madani H, AlRabiah FA, Al Hajjar S. Transmission and evolution
652        of the Middle East respiratory syndrome coronavirus in Saudi Arabia: a descriptive
653        genomic study. The Lancet. 2013 Dec 14;382(9909):1993-2002.
654

655  [M18]  Cotten M, Watson SJ, Zumla AI, Makhdoom HQ, Palser AL, Ong SH, Al Rabeeah
656        AA, Alhakeem RF, Assiri A, Al-Tawfiq JA, Albarrak A. Spread, circulation, and
657        evolution of the Middle East respiratory syndrome coronavirus. MBio. 2014 Feb
658        28;5(1):e01062-13.
659

660  [M19]  Dudas G, Carvalho LM, Rambaut A, Bedford T. MERS-CoV spillover at the camel-
661        human interface. Elife. 2018 Jan 16;7:e31257.
662

663  [M20]  Zhao Z, Li H, Wu X, Zhong Y, Zhang K, Zhang YP, Boerwinkle E, Fu YX. Moderate
664        mutation rate in the SARS coronavirus genome and its implications. BMC
665        evolutionary biology. 2004 Dec 1;4(1):21.
666

667  [M21]  Salemi M, Fitch WM, Ciccozzi M, Ruiz-Alvarez MJ, Rezza G, Lewis MJ. Severe
668        acute respiratory syndrome coronavirus sequence characteristics and evolutionary

669        rate estimate from maximum likelihood analysis. Journal of virology. 2004 Feb

670        1;78(3):1602-3.

671

672   [M22]  Wu SF, Du CJ, Wan P, Chen TG, Li JQ, Li D, Zeng YJ, Zhu YP, He FC. The genome

673        comparison of SARS-CoV and other coronaviruses. Yi chuan= Hereditas. 2003

674        Jul;25(4):373-82.

675

676   [M23]  Chinese SARS Molecular Epidemiology Consortium. Molecular evolution of the

677        SARS coronavirus during the course of the SARS epidemic in China. Science. 2004

678        Mar 12;303(5664):1666-9.

679

680   [M24]  Sanjuán R. From molecular genetics to phylodynamics: evolutionary relevance of

681        mutation rates across viruses. PLoS pathogens. 2012 May;8(5).

682

683   [M25]  Duchêne S, Geoghegan JL, Holmes EC, Ho SY. Estimating evolutionary rates using

684        time-structured data: a general comparison of phylogenetic methods. Bioinformatics.

685        2016 32(22), 3375-3379.