# GigaScience

## Analysis of SARS-CoV-2 known and novel subgenomic mRNAs in cell culture, animal model and clinical samples using LeTRS, a bioinformatic tool to identify unique sequence identifiers.

### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-21-00142R2 |
| Full Title: | Analysis of SARS-CoV-2 known and novel subgenomic mRNAs in cell culture, animal model and clinical samples using LeTRS, a bioinformatic tool to identify unique sequence identifiers. |
| Article Type: | Research |

| Abstract: | SARS-CoV-2 has a complex strategy for the transcription of viral subgenomic mRNAs (sgmRNAs), which are targets for nucleic acid diagnostics. Each of these sgmRNAs has a unique 5' sequence, the leader-transcriptional regulatory sequence gene junction (leader-TRS-junction), that can be identified using sequencing. High resolution sequencing has been used to investigate the biology of SARS-CoV-2 and the host response in cell culture and animal models and from clinical samples. LeTRS, a bioinformatics tool, was developed to identify leader-TRS-junctions and be used as a proxy to quantify sgmRNAs for understanding virus biology. LeTRS is readily adaptable for other coronaviruses such as Middle East respiratory syndrome coronavirus (MERS-CoV) or a future newly discovered coronavirus. LeTRS was tested on published datasets and novel clinical samples from patients and longitudinal samples from animal models with COVID-19. LeTRS identified known leader-TRS-junctions and identified putative novel sgmRNAs that were common across different mammalian species. This may be indicative of an evolutionary mechanism where plasticity in transcription generates novel open reading frames, that can then subject to selection pressure. The data indicated multi-phasic abundance of sgmRNAs in two different animal models. This recapitulates the relative sgmRNA abundance observed in cells at early points in infection, but not at late points. This pattern is reflected in some human nasopharyngeal samples, and therefore has implications for transmission models and nucleic acid-based diagnostics. LeTRS provides a quantitative measure of sgmRNA abundance from sequencing data. This can be used to assess the biology of SARS-CoV-2 (or other coronaviruses) in clinical and non-clinical samples, especially to evaluate different variants and medical countermeasures that may influence viral RNA synthesis. |
|---|---|

| Corresponding Author: | Julian Hiscox University of Liverpool Liverpool, UNITED KINGDOM |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | University of Liverpool |
| Corresponding Author's Secondary Institution: | |
| First Author: | Xiaofeng Dong |
| First Author Secondary Information: | |
| Order of Authors: | Xiaofeng Dong |
| | Rebekah Penrice-Randal |
| | Hannah Goldswain |
| | Tessa Prince |

| | |
|---|---|
| | Nadine Randle |
| | Donavan-Banfield I'ah |
| | Francisco J Salguero |
| | Julia Tree |
| | Ecaterina Vamos |
| | Charlotte Nelson |
| | Jordan Clark |
| | Yan Ryan |
| | James P. Stewart |
| | Malcolm G. Semple |
| | John Kenneth Baillie |
| | Peter J. Openshaw |
| | Lance Turtle |
| | David A. Matthews |
| | Miles W. Carroll |
| | Alistair C. Darby |
| | Julian A. Hiscox |
| **Order of Authors Secondary Information:** | |
| **Response to Reviewers:** | Reviewer reports:<br>Reviewer #1: Comments: It is an important study. Except for a few minor points, the authors have addressed most of the reviewers' suggestions. This manuscript will be considered for acceptance after addressing the following minor suggestions:<br><br>1.The authors have compared the algorithm design, input, and output, and the counts of predicted sgmRNA across four tools. However, it would be nice if the authors could compare these tools' performances regarding prediction accuracy, F-measure, sensitivity, and specific scores. These will let the readers and potential users have a better sense of choosing a different tool for different purposes.<br><br>[We have added the prediction accuracy, F-measure, sensitivity, and specific scores, calculated based on simulated Illumina and Nanopore reads, in the Table 1.]<br><br>2.It is unclear what the red line means in Supplemental Figure 8-9.<br><br>[The red lines in Supplemental Figure 8 and 9 are for the normalized count of sgmRNA identified by LeTRS. We have moved this to Supplementary Table 12.]<br><br>3.On page 18, lines 364-370. The analysis and significance that the authors stated in that paragraph do not show the apparent trends in Supplemental Figure 9. Would the authors update the figure types to reflect the results of their statistical tests?<br><br>[We have updated the boxplots in Supplemental Figures 8 and 9. We used a paired samples one-sided Wilcoxon test that takes account the difference at each modification site of two compared sgmRNAs in different time points. A large amount of modification sites with differences resulted a low p-value even the trends in boxplots are not very large.]<br><br>4.On page 18, line 370. The author mentioned that "The abundance of most sgmRNAs decreased with time, and both of these factors could account for the frequency of methylation." Based on the context, it seems that the conclusion could not be derived. Because the methylation frequency is a ratio, then it may not correlate with the abundance of the sgmRNAs. |

| | [We have removed this sentence to reflect the reviewer's content.] |
|---|---|
| | Reviewer #2: Happy with revisions, no further comments |
| **Additional Information:** | |
| **Question** | **Response** |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using | Yes |

a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.

Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?

1    Analysis of SARS-CoV-2 known and novel subgenomic mRNAs in cell culture, animal model and

2    clinical samples using LeTRS, a bioinformatic tool to identify unique sequence identifiers.

3

4    Xiaofeng Dong[1], Rebekah Penrice-Randal[1], Hannah Goldswain[1], Tessa Prince[1], Nadine Randle[1],

5    I'ah Donovan-Banfield[1,2], Francisco J. Salguero[3], Julia Tree[3], Ecaterina Vamos[1], Charlotte Nelson[1],

6    Jordan Clark[1], Yan Ryan[1], James P. Stewart[1], Malcolm G. Semple[1,2], J. Kenneth Baillie[4], Peter J. M.

7    Openshaw[5], Lance Turtle[1,2], David A. Matthews[6], Miles W. Carroll[2,3], Alistair C. Darby[1] and Julian

8    A. Hiscox[1,2,7].

9

10    [1]Institute of Infection, Veterinary and Ecological Sciences, University of Liverpool, UK.

11    [2]NIHR Health Protection Research Unit in Emerging and Zoonotic Infections, Liverpool, UK.

12    [3]UK-Health Security Agency, Salisbury, UK.

13    [4]The Roslin Institute, University of Edinburgh, UK.

14    [5]National Heart and Lung Institute, Imperial College London, UK.

15    [6]University of Bristol, UK.

16    [7]Infectious Diseases Horizontal Technology Centre (ID HTC), A*STAR, Singapore.

17    Corresponding author: julian.hiscox@liverpool.ac.uk

18

19    **ORCID iDs:**

20    Julian A Hiscox [0000-0002-6582-0275]; Rebekah Penrice-Randal [0000-0002-0653-2097];

21    Hannah Goldswain [0000-0003-4194-8714]; Tessa Prince [0000-0002-8796-2629]; Nadine

22    Randle [0000-0002-3775-9585]; Francisco J Salguero [0000-0002-5315-3882]; Julia Tree [0000-

23    0002-1720-6764]; James P Stewart [0000-0002-8928-2037]; Malcolm G Semple [0000-0001-

24    9700-0418]; John Kenneth Baillie [0000-0001-5258-793X]; Peter J Openshaw [0000-0002-7220-

25    2555]; Lance Turtle [0000-0002-0778-1693]; David A Matthews [0000-0003-4611-8795]; Miles

26    W Carroll [0000-0002-7026-7187]; Alistair C Darby [0000-0002-3786-6209]

27    **Abstract**

28    SARS-CoV-2 has a complex strategy for the transcription of viral subgenomic mRNAs (sgmRNAs),

29    which are targets for nucleic acid diagnostics. Each of these sgmRNAs has a unique 5' sequence,

30    the leader-transcriptional regulatory sequence gene junction (leader-TRS-junction), that can be

31    identified using sequencing. High resolution sequencing has been used to investigate the biology

32    of SARS-CoV-2 and the host response in cell culture and animal models and from clinical samples.

33    LeTRS, a bioinformatics tool, was developed to identify leader-TRS-junctions and be used as a

34    proxy to quantify sgmRNAs for understanding virus biology. LeTRS is readily adaptable for other

35    coronaviruses such as Middle East respiratory syndrome coronavirus (MERS-CoV) or a future

36    newly discovered coronavirus. LeTRS was tested on published datasets and novel clinical samples

37    from patients and longitudinal samples from animal models with COVID-19. LeTRS identified

38    known leader-TRS-junctions and identified putative novel sgmRNAs that were common across

39    different mammalian species. This may be indicative of an evolutionary mechanism where

40    plasticity in transcription generates novel open reading frames, that can then subject to selection

41    pressure. The data indicated multi-phasic abundance of sgmRNAs in two different animal models.

42    This recapitulates the relative sgmRNA abundance observed in cells at early points in infection,

43    but not at late points. This pattern is reflected in some human nasopharyngeal samples, and

44    therefore has implications for transmission models and nucleic acid-based diagnostics. LeTRS

45    provides a quantitative measure of sgmRNA abundance from sequencing data. This can be used

46    to assess the biology of SARS-CoV-2 (or other coronaviruses) in clinical and non-clinical samples,

47    especially to evaluate different variants and medical countermeasures that may influence viral

48    RNA synthesis.

49  **Importance**

50  When infecting cells, SARS-CoV-2 not only replicates its genome but also makes molecules called

51  subgenomic mRNAs (sgmRNAs) that are used as the template for many of the viral proteins,

52  including the spike glycoprotein. The sgmRNAs can only be found in infected cells, and therefore

53  their presence and ratio in a clinical sample is indicative that viral RNA synthesis has occurred,

54  and infected cells are present. The sgmRNAs are targets for diagnostic assays. We have developed

55  a rapid informatics methodology (LeTRS) to identify these unique molecules from multiple types

56  of sequencing data generated in response to the COVID-19 pandemic. We used this pipeline to

57  follow the pattern of sgmRNA abundance in nasopharyngeal samples taken from non-human

58  primate models and clinical samples from humans. We identified putative novel sgmRNAs that

59  may point to a potential new evolutionary mechanism in the virus. The data indicated that SARS-

60  CoV-2 RNA synthesis (and by inference infection) may occur in waves, and this has implications

61  for diagnostics and modelling of disease spread.

62

**Introduction**

Various sequencing approaches are used to characterise SARS-CoV-2 RNA synthesis in cell culture [1, 2], ex vivo models [3] and clinical samples. This can include nasopharyngeal swabs from patients with COVID-19 [4] to post-mortem samples from patients who died of severe disease [5]. Bioinformatic interrogation of this data can provide critical information on the biology of the virus. SARS-CoV-2 genomes are message sense, and the 5' two thirds of the genome is translated and proteolytically cleaved into a variety of functional subunits, many of which are involved in the synthesis of viral RNA [6]. The remaining one third of the genome is expressed through a nested set of subgenomic mRNAs (sgmRNAs). These have common 5' and 3' ends with the coronavirus genome, including a leader sequence, and are thus co-terminal. Many studies have shown that the sgmRNA located towards the 3' end of the genome, which encodes the nucleoprotein, generally has a higher abundance than those located immediately after the 1a/b region and the genome itself in infected cells [7, 8]. However, there is not necessarily a precise transcription gradient of the sgmRNAs. The 5' leader sequence on the sgmRNAs is immediately abutted to a short sequence called a transcriptional regulatory sequence (TRS) that is involved in the control of sgmRNA synthesis [9, 10]. These TRSs are located along the genome and are proximal to the start codons of the open reading frames [11]. In the negative sense the TRSs are complementary to a short portion of the genomic leader sequence. The TRS is composed of a short core motif that is conserved and flanking sequences [9, 10, 12]. The core motif of the TRS in SARS-CoV-2 is ACGAAC.

84    The prevailing thought is that synthesis of sgmRNAs involves a discontinuous step during negative

85    strand synthesis [13, 14]. A natural consequence of this is recombination resulting in insertions

86    and deletions (indels) in the viral genome and the formation of defective viral RNAs. Thus, the

87    identification of the leader/sgmRNA complexes by sequencing provides information on the

88    abundance of the sgmRNAs and evidence that transcription has occurred in the tissue being

89    analysed. In terms of clinical samples, if infected cells are present, then leader/sgmRNA 'fusion'

90    sequence can be identified, and inferences made about active viral RNA synthesis from the

91    relative abundance of the sgmRNAs. In the absence of published data from human challenge

92    models, the kinetics of virus infection are unknown, and most studies will begin with detectable

93    viral RNA on presentation of the patient with clinical symptoms. In general, models of infection

94    of humans with SARS-CoV-2 assume an exponential increase in viral RNA synthesis followed by a

95    decrease, as antibody levels increase [15].

96

97    To investigate the presence of SARS-CoV-2 sgmRNAs in clinical (and other) samples, a

98    bioinformatics tool (LeTRS), was developed to analyse sequencing data from SARS-CoV-2

99    infections by identifying the unique leader-TRS gene junction site for each sgmRNA. The utility of

100   this tool was demonstrated on cultured cells infected with SARS-CoV-2, nasopharyngeal samples

101   from humans with COVID-19 and longitudinal analysis of nasopharyngeal samples from two non-

102   human primate models infected with SARS-CoV-2. The tool is adaptable for other coronaviruses.

103   The results have implications for virus biology, diagnostics and disease modelling.

104

**Results**

A tool, LeTRS (named after the leader-TRS fusion site), was developed to detect and quantify defined leader gene junctions of SARS-CoV-2 (and other coronaviruses) from multiple types of sequencing data. This was used to investigate SARS-CoV-2 sgmRNA synthesis in humans and non-human primate animal models. LeTRS was developed using the Perl programming language, including a main program for the identification of sgmRNAs and a script for plotting graphs of the results. The tool accepts FASTQ files derived from Illumina paired-end or Oxford Nanopore sequencing (amplicon or direct RNA), or BAM files produced by a splicing alignment method with a SARS-CoV-2 genome (Supplementary Figure 1). Note that SARS-CoV-2 sgmRNAs are not formed by splicing, but this is the apparent observation from sequencing data because of the discontinuous nature of transcription. By default, LeTRS analyses SARS-CoV-2 sequence data by using 10 known leader-TRS junctions and an NCBI reference genome (NC_045512.2) to identify leader dependent canonical sgmRNAs. However, given the potential heterogeneity in the leader-TRS region and potential novel (leader dependent noncanonical) sgmRNAs the user can also provide customised leader-TRS junctions and SARS-CoV-2 variants as a reference. As there is some heterogeneity in the leader-TRS sites, LeTRS was also designed to search for multiple features of sgmRNAs. This included the leader-TRS junction in a given interval, report on the 20 nucleotides at the 3' end of the leader sequence, the TRS, translate the first predicted orf of the sgmRNA, and find the conserved ACGAAC sequences in the TRS. LeTRS can also be used to identify the sequencing reads with leader independent fusion sites that has been suggested to probably produce unknown ORFs yielding functional products [2]. The tool was designed to investigate very large data sets that are produced during sequencing of multiple samples.

127

**Combinations of read alignments with the leader-TRS junction that are considered for**

**identifying leader-TRS junction sites**

Various approaches have been used to sequence the SARS-CoV-2 genome and in most cases, this

would also include any sgmRNAs as they are 3' co-terminal and share common sequence

extending from the 3' end. Methods such as ARTIC[16], MIDNIGHT[17] and RSLA[4] use primer

sets to generate overlapping amplicons that span the entire genome, and also amplify sgmRNA.

Included is a primer to the leader sequence, so that the unique 5' end of these moieties are

sequenced. Primer sets of ARTIC, MIDNIGHT and RSLA are generally formed of 2 pools. For the

ARTIC method, at the time of the study, only the pool 1 included a forward primer located within

the leader region (< 80 nts) of the SARS-CoV-2 genome (https://github.com/artic-network/artic-

ncov2019/blob/master/primer_schemes/nCoV-2019/V3/nCoV-2019.primer.bed).     Therefore,

LeTRS was designed with a function to analyse reads in the primer pool 1, pool 2 or both pools.

Unbiased sequencing can also be used in methodologies to identify SARS-CoV-2 sequence. Data

in the GISAID database have been generated by Oxford Nanopore (minority) or Illumina

(majority) based approaches. These can give different types of sequencing reads derived from

the sgmRNAs that can be mapped back on the reference SARS-CoV-2 genome by splicing

alignment (Figure 1A). For example, there are several different types of reads that can be derived

from mapping Illumina-based amplicon sequencing onto the reference viral genome (Figure 1B

and 1C). During the PCR stage, the extension time allows the leader-TRS region on the sgmRNAs

to be PCR-amplified by the forward primer and the reverse primer before and after leader-TRS

junction in different primer sets, respectively. If the amplicon had a length shorter than the

8

149    Illumina read length (usually 100-250 nts), both the forward and reverse primers would be

150    detected at the ends of each paired read (Figure 1B pink lines). If the amplicon was longer than

151    the Illumina read length, primer sequence would be only found at one end of each paired read

152    (Figure 1B green and brown lines), with the possibility of one of the paired reads having a fusion

153    site. The extension stage could also proceed with a single primer using cDNA derived from the

154    sgmRNA as a template. This type of PCR product has a very low amplification efficiency, but

155    theoretically could also generate the same Illumina paired-end read with a single primer

156    sequence at one end (Figure 1C). These paired-end reads could include the full length of the

157    leader sequence but might not reach the 3' end of the sgmRNA, because of the limitation of

158    Illumina sequencing length and extension time (Figure 1C). Also, unless there are cryptic TRSs

159    located towards the 3' end of the genome, all sgmRNAs would be expected to be larger than the

160    Illumina sequencing length.

161

162    In contrast, the different types of read alignment in the Nanopore based amplicon are simpler to

163    assign. The longer reads that tend to be generated by Nanopore sequencing (depending on

164    optimisation) enable the capture of full-length sequences of all amplicons. Provided the leader

165    sequence is included as a forward primer most of the reads spanning the leader-TRS junction

166    would contain the forward and reverse primer sequences at both ends (Figure 1D pink lines). If

167    the extension time allowed, single primer PCR amplification could take the Nanopore amplicon

168    sequencing reads to both the 3' and 5' ends of the sgmRNAs, and these types of reads would only

169    have a primer sequence at one end (Figure 1D brown lines). In the Nanopore direct RNA

170    sequencing (dRNAseq) approach, the full-length sgmRNA could be sequenced and mapped

171    entirely on the leader and TRS-orf regions (Figure 1E).

172

173    **Evaluation of LeTRS on SARS-CoV-2 infection in cell culture.**

174    In order to assess the ability of LeTRS to identify the leader-TRS junctions from sequencing

175    information, a total RNA sample was prepared at 72 hours post-infection (hpi) from hACE2-A549

176    cells infected with SARS-CoV-2 (a lineage B isolate). This RNA was sequenced using an amplicon-

177    based approach (ARTIC) with either Nanopore (ARTIC-Nanopore) or Illumina (ARTIC-Illumina), or

178    alternatively by a Nanopore dRNAseq approach [18]. The ARTIC-Nanopore (Figure 2A,

179    Supplementary Table 1) and ARTIC-Illumina (Figure 2B, Supplementary Table 2) sequencing data

180    were evaluated with LeTRS by setting the analysis to both primers pools. For dRNAseq (Figure 2C,

181    Supplementary Table 3), data was evaluated with LeTRS using the default setting. All the major

182    known leader-TRS gene junctions were identified by these sequencing methods. Analysis

183    demonstrated an expected pattern of abundance of the leader-TRS gene junctions with the

184    leader-TRS nucleoprotein gene junction being most abundant (Figure 2A, B and C; Supplementary

185    Tables 1, 2 and 3). Novel low abundance leader-TRS gene junctions were also identified (Figure

186    2A, B and C; Supplementary Tables 1, 2 and 3). These known and novel leader-TRS junctions were

187    also known as leader dependent canonical and noncanonical fusions, respectively [2]. LeTRS also

188    has a function to identify leader independent long-distance fusion (>5,000 nt) and local joining

189    yielding a deletion between proximal sites (20-5,000 nt distance) in the sequencing reads. The

190    leader independent fusions (coverage >= 2) are shown in Supplementary Tables 1, 2 and 3. Indel

191    sequencing errors are frequent (defined as less than 20 nucleotides), especially in Nanopore

10

192   sequencing data, and therefore it is difficult to find precise fusion (apparent splicing) sites in this

193   case [19]. However, some of the novel leader-TRS junctions (noncanonical fusions) and leader

194   independent fusions in the test sample were supported by all three sequencing methods

195   (Supplementary Figure 2) with similar fusion sites. Many local fusions/deletions within the orf3,

196   E, M, orf6, orf7a, orf7b, orf8 and N genes were identified (Supplementary Figure 2 G, H and I)

197   confirmed previous findings [2, 20], and indicates these are common events. Some of the novel

198   leader-TRS junctions (noncanonical fusions) and leader independent fusions may be the result of

199   sequencing or reverse transcription errors, especially those with low abundance (Supplementary

200   Tables 1, 2 and 3; Supplementary Figure 2). The ARTIC-Illumina approach identified fewer novel

201   leader-TRS junctions (noncanonical fusions) and leader independent fusions than the other two

202   sequencing methodologies, probably due to lower sequencing coverage (Supplementary Tables

203   1, 2 and 3).

204

205   For ARTIC approaches, LeTRS was designed to analyse reads in the primers pool 1, pool 2 or both

206   pools. Only the ARTIC pool 1 included a forward primer that is located within the leader region

207   (< 80 nts) of the SARS-CoV-2 genome. The leader-TRS regions of sgmRNAs can be PCR-amplified

208   by both forward and reverse primers in ARTIC pool 1, but only reverse primers in ARTIC pool 2.

209   The read counts evaluated by LeTRS in both ARTIC-Nanopore and ARTIC-Illumina were compared

210   in the test data for pool 1 and 2, and found only very few reads/read pairs contained the correct

211   primers (Supplementary Table 4 and 5), suggesting the primers in ARTIC pool 2 generally do not

212   contribute to sequencing of leader-TRS regions.

213

214 **Comparison with other informatic tools that can identify leader TRS gene junctions.**

215 Other tools have been developed to identify sgmRNAs from ARTIC-Illumina and ARTIC-Nanopore

216 sequencing data, such as Periscope (v0.1.0) [21], SARS-CoV-2-leader

217 (https://github.com/hyeshik/sars-cov-2-transcriptome) [18] and SuPER

218 (https://github.com/ncbi/SuPER) [22]. These tools were compared with LeTRS as shown in Table

219 1. LeTRS and Periscope used the FASTQ files as input, while SARS-CoV-2-leader and SuPER

220 required SAM files from a user generated alignment. Searching fusion site and sequences tag in

221 the sequencing reads are two major methods used. LeTRS and SuPER analysed the fusion/splicing

222 information in sequence reads achieved by an alignment program and also take account of the

223 conserved ACGAAC sequences in the TRS. Periscope and SARS-CoV-2-leader are based on

224 searching for a short tag sequence in the leader from sequencing reads. However, searching for

225 a short tag sequence in the leader with the high error rate associated with Nanopore data can be

226 challenging. LeTRS and Periscope use primer information to differentiate reads mapping to

227 amplicons to reads mapping from original virus genomes. Besides Periscope, output from

228 dRNAseq is supported by the other available tools. Illumina sequencing reads are usually short (<

229 250 bases), paired and sequenced from both ends. If both reads in a single pair contain a fusion

230 site this will be counted twice by the other three tools (Figure 1B green and pink). However, if

231 only one of the reads in the pair contains a fusion site it will be counted once (Figure 1B brown).

232 This leads to biased counting. LeTRS takes this into account by treating each read pair as a single

233 event. LeTRS also has a unique function to analyse reads in the primers pool 1, pool 2 or both

234 pools from ARTIC based sequencing (Table 1). Accuracy, sensitivity, specificity and the F-measure

235 score were calculated with simulated Illumina and Nanopore sequencing reads. All of these tools

236    performed better for analysing the simulated Illumina sequencing reads compared to the

237    simulated Nanopore sequencing reads (Table 1). LeTRS showed greater sensitivity and F-measure

238    score than the other tools for processing the simulated Nanopore sequencing reads (Table 1).

239

240    To compare the performance to LeTRS, these three tools were evaluated using the hACE2-A549

241    cell culture sample sequenced by ARTIC-Nanopore, ARTIC-Illumina and Nanopore dRNAseq.

242    Using the ARTIC-Nanopore sequencing data, all the tools reported a similar number of read

243    counts for the 10 known sgmRNAs (Supplementary Figure 3A). LeTRS showed fewer counts for

244    the ARTIC-Illumina than the other three tools because of considering read pairs (Supplementary

245    Figure 3B). Interestingly, Periscope also identified fewer nucleoprotein sgmRNAs with the ARTIC-

246    Illumina sequencing data (Supplementary Figure 3B). As of writing, Periscope does not yet

247    support Nanopore dRNAseq data, therefore LeTRS, SARS-CoV-2-leader and SuPER were

248    compared. LeTRS and SARS-CoV-2-leader generally identified more dRNAseq reads than SuPER,

249    especially for the nucleoprotein sgmRNA (Supplementary Figure 3C). Finally, the ratio of read

250    counts with the 10 known sgmRNA (S:orf3:E:M:orf6:orf7a:orf7b:orf8:N:orf10) were compared,

251    and the three tools showed almost an identical ratio when analysing data from the same

252    sequencing methods (Supplementary Figure 3D). ARTIC-Nanopore and Nanopore dRNAseq

253    resulted in a higher ratio of read counts with M and orf7a respectively (Supplementary Figure

254    3D). The read counts ratio of sgmRNAs mapping to spike was much lower with dRNAseq

255    approaches (Supplementary Figure 3D).

256

257    **Normalisation of read counts for sgmRNA**

258    Normalisation of read counts has been widely used for RNAseq in the comparison of gene

259    expression level across samples [23]. The normalisation is generally based on the ratio of reads

260    mapped on the gene to the total number of reads in that sample. These tools use this algorithm

261    for the normalisation of read counts in searching for sgmRNA [21, 24]. LeTRS also incorporated a

262    method to differentiate the total reads mapped (i) or whether the reads have forward primer

263    only (ii), reverse primer only (iii), both primers (iv) or at least one primer (v) present. This is

264    achieved by (i) the total number of reads mapped on the SARS-CoV-2 genome for the number of

265    reads of leader-TRS fusion site as the numerator; (ii) the total number of reads with forward

266    primers only for the number of reads of leader-TRS fusion site with forward primers only as the

267    numerator; (iii) the total number of reads with reverse primers only for the number of reads of

268    leader-TRS fusion site with reverse primers only as the numerator; (iv) the total number of reads

269    with both primers for the number of reads of leader-TRS fusion site with both as the numerator

270    and (v) the total number of reads with at least one primer on one side for the number of reads

271    of leader-TRS fusion site with at least one primer on as the numerator (notes in Supplementary

272    Tables 1, 2 and 3).

273

274    Because LeTRS considers the primers; pool 1, pool 2 or both pools, normalisation could be

275    observed in ARTIC pool 1 only to minimise the effect from ARTIC pool 2 since primers in ARTIC

276    pool 2 are almost not involved the sequencing of leader-TRS regions (as described above). For

277    the same RNA derived from the hACE2-A549 cell culture sample sequenced by ARTIC-Nanopore,

278    ARTIC-Illumina or Nanopore dRNAseq approaches, the normalised counts for the known

279    sgmRNAs were much smaller with the pool 1 of PCR based amplicon methods (ARTIC-Nanopore

280    and ARTIC-Illumina) than the Nanopore dRNAseq approach (Figure 3A and C for the reads with

281    at least one primer sequence; Supplementary Tables 3, 4 and 5). However, the normalised counts

282    with ARTIC-Nanopore and ARTIC-Illumina showed the same ratio of known sgmRNA as the

283    Nanopore dRNAseq approach, except for sgmRNAs mapping to S and orf7a (Figure 3B and D for

284    the reads with at least a primer sequence). PCR based approaches increases the value of the

285    denominator and reduced the normalised count, because a full length of sgmRNA was counted

286    once with the dRNAseq approach compared to many times with the amplicon approaches. ARTIC-

287    Illumina had fewer normalised counts than ARTIC-Nanopore probably due to the sequencing bias

288    of Illumina during PCR [25]. Thus, if the samples were sequenced with the same methodology

289    they were comparable. With a PCR based method a normalised count should be used to show

290    the relative difference between samples.

291

292    LeTRS identified many reads with only one primer (one-sided amplification) with the PCR based

293    amplicon methods (Supplementary Tables 4 and 5). The ratio of reads with either forward and/or

294    reverse primers were compared for each sgmRNA to the overall ratios of reads, with forward

295    primers only or reverse primers only, both primers in all mapped reads of pool 1 and pool 2 and

296    the mapped reads with any fusion sites of pool 1 and pool 2. This indicated that abundant reads

297    were identified with a single pattern and these were similar to reads mapping to sgmRNAs,

298    suggesting a one sided amplification is associated with amplicon-based approaches

299    (Supplementary Figure 4).

300

**Analysis of sequencing data from longitudinal nasopharyngeal samples taken from two non-human primate models of COVID-19 indicated multi-phasic sgmRNA synthesis and novel sgmRNAs.**

Part of the difficulty of studying SARS-CoV-2 and the disease COVID-19 is establishing the sequence of events from the start of infection. Most samples from humans are from nasopharyngeal aspirates taken when clinical symptoms develop. This tends to be 5 to 6 days post-exposure. In the absence of a human challenge model, animal models can be used to study the kinetics of SARS-CoV-2[26, 27]. Two separate non-human primate (NHP) models, cynomolgus and rhesus macaques, were established for the study of SARS-CoV-2 that mirrored disease in most humans[26]. To study the pattern of sgmRNA synthesis over the course of infection, nasopharyngeal samples were sequentially gathered daily from 1 dpi up to 18 dpi from the two NHP models. RNA was purified from these longitudinal samples as well as the inoculum virus and viral RNA sequenced using ARTIC-Illumina.

As expected, analysis of the sequence data using LeTRS from the inoculum used to infect the NHPs indicated that leader gene junctions could be identified, but these did not follow the pattern of abundance of leader TRS-gene junctions found in infected cells in culture, where the leader TRS nucleoprotein gene junction was most abundant (Supplementary Figure 5). The inoculum would be expected to contain mostly genomic RNA found in virions. In contrast, analysis of the longitudinal sequencing data from nasopharyngeal aspirates from the NHP model using LeTRS identified leader TRS-gene junctions associated with the major sgmRNAs (Figure 4, Supplementary Table 7) as well as novel leader-TRS gene junction sites (Supplementary Figures

323    6 and 7). Analysing the abundance of the leader-TRS-gene junctions for both model species over

324    the course of infection revealed a phasic nature of sgmRNA synthesis in pool 1 to minimise the

325    effect from ARTIC pool 2 (Figure 4). The leader-TRS nucleoprotein gene junction was the most

326    abundant, and there was a phasic pattern of potential sgmRNA abundance identified with the

327    ARTIC-Illumina method (Figure 4). For both species, viral load and hence sgmRNA abundance had

328    decreased by 8 and 9 dpi.

329

330    **Analysis of leader-TRS-gene junction in human samples revealed expected and aberrant**

331    **abundances of sgmRNAs**

332    To investigate the pattern of leader-TRS-gene junction abundance during infection of SARS-CoV-

333    2 in humans, nasopharyngeal swabs from patients with COVID-19 were sequenced by ARTIC-

334    Illumina (using samples from COG-UK) (N=15 patients) (Figure 5, Supplementary Table 8) or by

335    ARTIC-Nanopore (using samples from ISARIC-4C) (N=15 patients) (Figure 6, Supplementary Tables

336    9 and 10). In several samples, leader-TRS-gene junctions were identified and followed an

337    expected pattern, with the nucleoprotein gene junction being the most abundant (e.g., Sample

338    1 in Figures 5A and B, Patient 2 day1 in Figure 6A and B). However, in several of the samples there

339    was very large representation of single leader-TRS-gene junction (e.g., Sample 4 and 5 in Figures

340    5A and B). These tended to map to the nucleoprotein gene (Sample 5, 8 and 13 Figures 5A and

341    B). The heterogeneity in abundance of leader-TRS-gene junctions was reminiscent of that from

342    the NHP study with a defined and expected pattern near the start of infection but then becoming

343    phasic. The samples gathered under ISARIC-4C were from hospitalised patients and permitted

344    analysis in relation to reported date of symptom onset and sequential sampling. In general, the

345    data indicated that the first sample on admission to hospital contained an abundance of leader-

346    TRS-gene junctions which resembled the pattern seen in infected cells (Patient 6 day 1 and day 9

347    in Figures 6A and B). However, with further days post-sample, e.g. (Patient 7 day 7 Figures 6A

348    and B), the leader-TRS nucleoprotein gene junction was the most abundant and far exceeded any

349    other detectable species. The abundance of leader-TRS nucleoprotein gene junction in the

350    patients at a later stage of infection followed that observed in the NHP model (Figure 4).

351

352    **Analysis of sequencing data from a previously published study investigating SARS-CoV-2 RNA**

353    **in samples from patients**

354    Recent research detected sgmRNAs mapping to E, ORF7a and N in swabs up to 14 days in one

355    patient and ORF7a and N in another patient up to 17 days after first detection by using a high-

356    throughput amplicon sequencing method known as Ion AmpliSeq Coronavirus Research Panel on

357    an Ion S5 XL genetic sequencer. The authors concluded these sgmRNAs may be present for a

358    significant time after active infection due to nuclease resistance and protection by cellular

359    membranes [24]. The sequencing data from this study was reanalysed using LeTRS, and

360    confirmed the finding of sgmRNAs in late infection from the two patients (Supplementary Table

361    11). Apart from nuclease resistance and protection by cellular membranes, a phasic pattern of

362    sgmRNA synthesis may also contribute to the presence of sgmRNAs at later time points.

363

364    **Analysis of sgmRNA modification in longitudinal samples in cell culture.**

365    N6-methyladenosine (6mA) is a widely observed modification on cellular RNA, and 5-

366    methylcytosine methylation (5mC) has also been reported on viral RNAs [18]. Methylation of

367    SARS_CoV-2 RNA was examined using sequencing data from the Nanopore direct RNA seq

368    approach. Total RNA was purified at 6, 12 and 24 hpi from cells infected with SARS-CoV-2. The

369    total RNA was sequenced and reads mapping to sgmRNAs were extracted with LeTRS for 6mA

370    and 5mC examination. Almost all 10 observed sgmRNAs showed the same number of

371    modification sites of 6mA and 5mC at 6, 12 and 24 hpi (Supplementary Table 12). Modification

372    with 5mC was more abundant than 6mA in all 10 known sgmRNAs (Supplementary Table 12).

373    There were differences in abundance of some sgmRNAs especially the M and N subgenomic

374    mRNAs (Supplementary Table 12). However, there did not appear to be a relationship between

375    number of methylation sites and the abundance of a particular sgmRNA (Supplementary Table

376    12).

377    To further evaluate the relationship between time post-infection and modification by

378    methylation, a paired samples one-sided Wilcoxon test was used. This analysis suggested that

379    the 5mC modification fraction at 24 hpi was significantly less than compared to modification at 6

380    and 12 hpi (p-value < 0.05), except for ORF7b and ORF10 (Supplementary Figures 8 and 9;

381    Supplementary Table 13). Modification with 6mA at 24 hpi was also significantly less than at 6

382    hpi, but not at 12 hpi (p-value < 0.05) in S, ORF3a, E, M, ORF6, ORF7a, ORF8 and N.

383

384    **Common properties/features of novel leader-TRS gene junctions and sgmRNAs**

385    The sequencing data from cells infected in culture (Supplementary Table 14), animal models and

386    clinical samples from humans indicated the presence of novel leader-TRS gene junctions. Their

387    detection generally increased with depth of coverage. Coronavirus replication and transcription

388    is promiscuous, and recombination is a natural result of this, resulting in indels and potential

389  gene rearrangements. Many of these novel leader-TRS junctions were centred around the known

390  gene orf but out of the search interval.  These types of leader-TRS-gene junctions could be only

391  found with spike, membrane, ORF6, ORF7b and nucleocapsid orfs, in which the membrane orf

392  was the most common (Figure 7A). To define what might be genuine novel leader-TRS-gene

393  junctions, these were compared across the data in all ARTIC-Illumina data (Figure 7B,

394  Supplementary Table 15). Five novel leader-TRS-gene junctions were identified that were

395  common to all the data, and the majority of these were present immediately 5' of the membrane

396  orf). The novel leader-TRS-gene junctions from LeTRS (Figure 7C) showed a similar distribution as

397  a previous study, although this study did not detail the precise location [28].

398

399

**400**   **Discussion**

**401**   Coronavirus sgmRNAs are only synthesised during infection of cells and therefore their presence

**402**   in sequence data can be indicative of active viral RNA synthesis. The abundance of the sgmRNAs

**403**   in infected cells should follow a general pattern where the sgmRNA encoding the nucleoprotein

**404**   is the most abundant. Identification and quantification of the unique leader-TRS-gene junctions

**405**   for each sgmRNA can be used as a proxy for their abundance.

**406**

**407**   LeTRS was developed to interrogate sequencing datasets to identify the leader-TRS-gene

**408**   junctions present at the 5' end of the sgmRNAs. LeTRS was first evaluated and validated on cell

**409**   culture data from published datasets [2, 16] and from a cell culture experiment as part of this

**410**   study and then used in an analysis of nasopharyngeal samples from NHP and human clinical

**411**   samples. The results showed that the positions of the leader-TRS junction sites with peak read

**412**   counts were the same as the given reference positions. The exception was at the leader-TRS-

**413**   gene junction for orf7b in the Nanopore sequencing. The normalised count results confirmed the

**414**   reads spanning the junctions showed that the leader-TRS nucleoprotein gene junction was the

**415**   most abundant, and orf7b and orf10 were the most infrequent in line with other data [2, 24].

**416**   Several low abundant leader-TRS junctions were identified in all of the datasets (Supplementary

**417**   Figure 2) with the implication these were either from potential lower abundant novel sgmRNAs

**418**   or represented known sgmRNAs, but with different leader-TRS junctions. Likewise, at low

**419**   frequency these could represent an aberrant viral transcription, perhaps as a mechanism to

**420**   generate new orfs for selection or these could be artefacts of the different sequencing processes

**421**   (Figure 2). Traditionally, such sgmRNAs have been first identified in coronaviruses by either

422    northern blot and/or metabolic labelling [8] and sequencing approaches are likely to be more

423    sensitive giving the amplification steps involved. Several other groups have identified novel

424    leader-TRS-gene junctions and potential sgmRNAs for other coronaviruses, including avian

425    infectious bronchitis virus [29]. The best way of validating potential novel sgmRNAs would be

426    through matching proteomic data to confirm genuine ORFs [1]. Analysis of several published

427    sequencing datasets identified novel viral RNA molecules that the authors suggested were

428    sgmRNAs containing only the 5' region of orf1a [30]. Such species are likely to be defective RNAs,

429    that act as templates for replication, rather than sgmRNAs. Interestingly, at later time points

430    post-infection in cell culture, potential novel sgmRNAs were found to be generated non-

431    specifically [30]. This potentially ties in with a disconnect of leader-TRS-gene junctions observed

432    in our study both *in vivo* from the nasopharyngeal samples from latter time points in the NHP

433    models and in humans. This is also shown in published data from SARS-CoV-2 infections in cell

434    culture gathered at later time points compared to earlier time points [2, 16].

435

436    Advanced filtering can improve the confidence of the identified leader-TRS junction from

437    sequencing data. Amplicon sequencing provided a unique opportunity to filter the sequencing

438    reads. The reads spanning the junctions with the correct forward primer, reverse primer or both

439    primer sequences at the ends of reads proved the known/novel sgmRNA existing in tested ARTIC-

440    Illumina and ARTIC-Nanopore amplicon sequencing data (Supplementary Tables 1 and 2). For

441    Illumina sequencing, the same junction on paired reads with at least one primer provided extra

442    evidence for leader-TRS identification. Some reads were identified that did not have primer

443    sequences and these were likely to be erroneously mapped, from template sgmRNA or low-

444 quality sequence. These were present at very low abundance compared to authentically mapped

445 reads (Supplementary Tables 1 and 2). The Nanopore dRNAseq approach had the potential to

446 generate full-length mRNA sequences. The polyA sequences and leader-TRS junctions in the

447 reads can be good signals to prove the full-length sgmRNA in the test data (Supplementary Table

448 3). Currently, LeTRS is the only tool to consider paired-end Illumina data and primer pools, and

449 therefore is suited for interrogating paired-end Illumina data and providing data from amplicon

450 sequencing information from either primer pools.

451

452 In terms of clinical samples (typically nasopharyngeal swabs), the presence of sgmRNAs will

453 generally be due to the presence of infected cells. This has been seen as indicative of active viral

454 RNA synthesis at the time of sampling[5, 31, 32], although these have also been postulated to be

455 present through resistant structures after infection has finished [33]. Analysis of inoculum

456 indicated that leader-TRS-gene junctions could be identified (Supplementary Figure 5) but that

457 these were not in the same ratio as found in cells infected in culture (e.g., Figure 2A, B and 2C).

458 Thus, if the abundance of leader-TRS-gene junctions follows an expected pattern of the leader-

459 TRS nucleoprotein gene junction being the most abundant followed by a general gradient in

460 sequence data from nasopharyngeal samples, then this may be indicative of an active infection

461 – and the presence of infected cells in a sample.

462

463 In the absence of a human challenge model, NHP models that closely resemble COVID-19 disease

464 in humans can be used to study SARS-CoV-2 infection from a very defined initial exposure. RNA

465 was sequenced from longitudinal nasopharyngeal samples from two NHP models, rhesus and

466    cynomolgus macaques [26]. LeTRS was used to identify the abundance of the leader-TRS-gene

467    junctions in this data. The analysis indicated a phasic pattern of sgmRNA synthesis with a large

468    drop off after 8 or 9 dpi in both NHP models. This phasic pattern may be explained by an initial

469    synchronous infection of respiratory epithelial cells followed by cell death. Released virus then

470    goes on to infect new epithelial cells, with virus infection increasing exponentially in waves but

471    becoming asynchronous. The decline in sgmRNA from 8 or 9 dpi overlaps with IgG seroconversion

472    and humoral immunity in both species [26], and follows similar kinetics to serology profiles

473    measured in patients with COVID-19.

474

475    The identification of sgmRNAs in nasopharyngeal samples and their kinetics has implications for

476    nucleic acid-based diagnostics (many of which have three targets, one in the orf1a/b region and

477    two which are shared between the genome and sgmRNAs – the nucleoprotein and the spike

478    genes). The phasic nature of leader-TRS-gene junctions in the longitudinal samples, and by

479    implication sgmRNAs, and overt abundance of the leader-TRS nucleoprotein gene junction found

480    in many of the human samples, suggests that it may not be possible to precisely identify where

481    in infection an individual is based on the abundance of sgmRNAs. Likewise, assuming equivalency

482    between the targets, if the nucleoprotein target is found to be more abundant than the spike

483    target than the genomic target, then this would suggest infected cells are present in the sample.

484    Decreases in Ct values associated with emerging variants could equally be explained by sloughed

485    cells being present in a nasopharyngeal sample as well as by increases in the amount of

486    virions/viral load. Therefore, we would caution that a decrease in Ct associated with RT-qPCR

487    based assays may not just be reflective of higher viral loads but also may be indicative of more

488    infected cells being present. These possibilities may be resolved by considering the relative ratios

489    of sgmRNAs identified.

490 **METHODS**

491 **Data input**

492 LeTRS was designed to analyse FASTQ files derived from Illumina paired-end or Nanopore

493 sequencing data derived from a SARS-CoV-2 amplicon protocol, or standard Nanopore SARS-CoV-

494 2 dRNAseq data (Figure 1). The Illumina/Nanopore FASTQ sequencing data were cleaned to

495 remove adapters and low-quality reads before input. Sequencing data derived from other

496 sequencing modes or platforms can also be analysed by LeTRS via input of a BAM file produced

497 by a custom splicing alignment method with a SARS-CoV-2 genome (NC_045512.2) as a reference

498 (Figure 1). This can also be rapidly adapted for other coronaviruses.

499

500 **Library preparations and sequencing**

501 We sequenced the 15 samples from human patients with Nanopore. Total RNA was isolated using

502 a QIAamp Viral RNA Mini Kit (Qiagen, Manchester, UK) by spin-column procedure according to

503 the manufacturer's instructions. Clinical samples were extracted with Trizol LS as described[4].

504 All RNA samples were treated with Turbo DNase (Invitrogen). SuperScript IV (Invitrogen) was

505 used to generate single-strand cDNA using random primer mix (NEB, Hitchin, UK). ARTIC V3 PCR

506 amplicons from the single-strand cDNA were generated following the Nanopore Protocol of PCR

507 tiling of SARS-CoV-2 virus (Version: PTC_9096_v109_revL_06Feb2020). Amplicons generated by

508 ARTIC PCR were purified and normalised to 200 fmol before DNA end preparation and barcode

509 and adapter ligation. Library was loaded onto a FLO-MIN106 flow cell and sequencing reads were

510 called with Guppy using the high-accuracy calling parameters.

511

512    The NHP samples and their inoculum, and our laboratory experiments conducted in cells were

513    sequenced with Illumina. The amplicons products for Illumina sequencing were prepared as per

514    the Nanopore sequencing above and then used in Illumina NEBNext Ultra II DNA Library

515    preparation.  Following 4 cycles of amplification the library was purified using Ampure XP beads

516    and quantified using Qubit and the size distribution assessed using the Fragment analyzer. Finally,

517    the ARTIC library was sequenced on the Illumina® NovaSeq 6000 platform (Illumina®, San Diego,

518    USA, RRID:SCR_016387) following the standard workflow. The generated raw FastQ files (2 x 250

519    bp) were trimmed to remove Illumina adapter sequences using Cutadapt v1.2.1

520    (RRID:SCR_011841)[34]. The option "−O 3" was set, so the that 3' end of any reads which

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

522    trimmed to remove low quality bases, using Sickle v1.200 [35] with a minimum window quality

523    score of 20. After trimming, reads shorter than 10 bp were removed.

524

525    The LeTRS was also tested with a combined Nanopore-ARTIC v3 amplicon dataset of 7 published

526    viral cell culture samples (barcode01-barcode07) [16], and a dataset from a published direct RNA

527    Nanopore sequencing analysis Vero cells infected with SARS-CoV-2 or an uninfected negative

528    control [2].

529

530    **Sequencing data alignment and basic filtering**

531    LeTRS controlled Hisat2 v2.1.0 (RRID:SCR_015530)[36] to map the paired-end Illumina reads

532    against the SARS-CoV-2 reference genome (NC_045512.2) with the default setting, and

533    Minimap2 v2.1 [19] to align the Nanopore cDNA reads and direct RNA-seq reads on the viral

534  genome using Minimap2 with "–ax splice" and "-ax splice -uf -k14" parameters, respectively.

535  LeTRS provided 10 known leader-TRS junctions to improve alignment accuracy by using "--

536  known-splicesite-infile" function in Hisat2 and "--junc-bed" function in Minimap2, but this

537  application could be optionally switched off by users. In order to remove low mapping quality

538  and mis-mapped reads before searching the leader-TRS junction sites, LeTRS used Samtools v1.9

539  (RRID:SCR_002105)[37] to have basic filtering for the reads in the output Sam/Bam files according

540  to their alignment states as shown (Table 9 - basic filtering).

541

542  **Searching the leader-TRS motifs**

543  After the mapping and basic filtering step, LeTRS searched aligned reads spanning the leader-TRS

544  junctions in the SARS-CoV-2 reference genome (Supplementary Figure 1). For the known leader-

545  TRS junctions, LeTRS searched the reads including the leader-TRS junctions within a given interval

546  around the known leader and TRS junctions sites. The leader break site interval is ±10 nts, and

547  the TRS breaking sites interval is -20 nts to the 1 nt before the first known AUG in the default

548  setting (the intervals can be changed to custom values to investigate heterogeneity). LeTRS then

549  reported a peak count that was the number of reads carrying the most common leader-TRS

550  junctions within the given leader and TRS breaking sites intervals, and a cluster count that was

551  the number of all reads carrying leader-TRS junctions within the given leader and TRS breaking

552  sites intervals (Tables 1-6). LeTRS also searched the junctions out of the given intervals (the

553  genomic position of leader breaking site < 80) and reported the number of reads (>10 by default)

554  with novel leader-TRS junctions. These number of read counts were also reported by number of

555  reads in 1000000 as normalisation. The read including the known and novel leader-TRS junctions

556 could be optionally outputted in FastA format. Based on identified known and novel leader-TRS

557 junctions, LeTRS could report 20 nucleotides towards the 3' end of the leader sequence, the TRS

558 and translated the first orf of sgmRNAs sequence, and find the conserved ACGAAC sequences in

559 the TRS (Table S1-S6).

560

561 **Advance filtering**

562 Based on the alignment possibilities illustrated in Figure 2 and discussed, LeTRS further filters the

563 identified reads with known and novel leader-TRS junctions. This step is named as advance

564 filtering and can only applied when the input data is from Illumina paired-end reads, Nanopore

565 cDNA reads or Nanopore RNA reads (Table 2). If a BAM file is used as input data, the advanced

566 filtering step would be automatically skipped (Table 2). The number of reads including the known

567 and novel leader-TRS junctions, and the number of reads filtered with corresponding advance

568 filtering criteria were outputted into two tables in tab format (Tables 1-6).

569

570 **Leader-TRS junction plotting**

571 LeTRS-plot was developed as an automatic plotting tool that interfaces with the R package

572 ggplot2 v3.3.3 to view the leader-TRS junctions in the tables generated by LeTRS (Figure 3-5). The

573 plot shows peak count, filtered peak count, normalized peak count and normalized filtered peak

574 count for known leader-TRS junctions, and novel junction counts, filtered novel junction count,

575 normalized novel junction count and filtered normalized novel junction for novel leader-TRS

576 junctions.

577

**Simulation of Illumina and Nanopore reads**

578

579 To assess the performance of LeTRS and other tools, simulated Illumina reads were generated

580 using ART (v2.5.8) [38] and Nanopore reads were generated using NanoSim (v2.6.0,

581 RRID:SCR_018243) [39]. The real reads generated by the ARTIC-Nanopore approach, ARTIC-

582 Illumina approach and Nanopore dRNAseq approach for the hACE2-A549 cells infected with

583 SARS-CoV-2 were used to create custom Illumina and Nanopore read quality/error profiles with

584 ART and NanoSim. Illumina paired reads (2x250 bp) and Nanopore cDNA-1D read for both ARTIC

585 and sgmRNA amplicons were simulated at 50000 × coverage for each amplicon and 2,000,000

586 reads in total, respectively. Nanopore dRNAseq reads (2,000,000) of the sgmRNA and viral

587 genome were generated using transcriptome mode.

588

**RNA modifications**

589

590 Total RNA extracted from cultured cells at 6, 12 and 24 hours were collected for Oxford

591 Nanopore direct RNA sequence. LeTRS was then run with a parameter of "extractfasta" to extract

592 subgenomic mRNAs reads in sequenced samples. The fast5 files that corresponds to the

593 extracted subgenomic mRNAs reads were withdrawn using fast5_subset in Oxford Nanopore

594 ont_fast5_api package (v0.3.2, https://github.com/nanoporetech/ont_fast5_api). The re-

595 squiggle algorithm in Tombo analysis pipelines (v1.5.1,

596 https://github.com/nanoporetech/tombo) defines a new assignment from raw signals to

597 reference sequence with "--num-most-common-errors 5" option. The resquiggled raw signals

598 were further processed using "detect_modifications alternative_model" functions in Tombo by

599 setting "--rna and --alternate-bases 5mC" to identify 5-methylcytosine (5mC), and "predict_sites"

600 in Nanom6A package (v2021_10_22) [40] with default setting to identify N6-methyladenosine

601 (6mA) in the subgenomic mRNAs reads.

602

603 **Ethics approval and consent to participate**

604 All experimental work on NHPs was conducted under the authority of a UK Home Office approved

605 project license (PDC57C033) that had been subject to local ethical review at PHE Porton Down by

606 the Animal Welfare and Ethical Review Body (AWERB) and approved as required by the Home

607 Office Animals (Scientific Procedures) Act 1986 and the full ethics and NHP model are described.

608 **Consent for publication**

609 Not applicable

610

611 **Data Availability**

612 Illumina and Nanopore test data sets are available under NCBI PRJNA699398. Snapshots of the

613 code are available in the *GigaScience* GigaDB repository[41].

614 **Availability and requirements**

615    • Project name: LeTRS
616    • Project home page: https://github.com/Hiscox-lab/LeTRS
617    • Operating system(s): Platform independent
618    • Programming language: Perl
619    • Other requirements: samtools(>=1.11), hisat2(=2.1.0), minimap2(=2.17),
620      portcullis(>=1.1.2)
621    • License: Apache 2.0
622    • RRID:SCR_022138
623

624 **Competing interests**

625 The authors declare that they have no competing interests

648    Trust, or PHE. The funders had no role in the study design; in the collection, analysis, and

649    interpretation of data; in the writing of the report; or in the decision to submit the article for

650    publication.

651

652    **Authors' contributions**

653    X.D. developed the LeTRS software and performed the informatics analysis. X.D., A.D. and J.A.H.

654    analysed the data. J.S., J.T. and M.W.C. co-ordinated the NHP work and sample processing. R.P.-

655    R., J.P.S., H.G., T.P. and N.R. were involved in sequencing and informatics analysis of the NHP

656    samples with D.A.M. A.D. oversaw sequencing of the human clinical samples with E.V. and C.N

657    for the COG-UK data. R.P.-R. and J.A.H. oversaw sequencing of samples under the auspices of

658    ISARIC-4C with clinical samples collected and managed by J.K.B, L.T., M.G.S. and P.J.M.O. J.A.H.

659    and M.W.C. initiated and led the study and wrote the manuscript with X.D., R.P.-R., A.D. with

660    other authors involved in editing the final version.

661    **Acknowledgments**

662    We would like to thank all members of the Hiscox Laboratory and the Centre for Genome

663    Research for supporting SARS-CoV-2/COVID-19 sequencing research. We would like to

664    acknowledge members of the COG-UK and ISARIC4C consortia for acquisition of the human

665    samples used in this study.

666

667

668

669

670 **References**

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

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

673 *spike glycoprotein.* Genome Med, 2020. **12**(1): p. 68.

674 2. Kim, D., et al., *The Architecture of SARS-CoV-2 Transcriptome.* Cell, 2020. **181**(4): p. 914-

675 921 e10.

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

677 *Amplicon-Based Sequencing, Random Hexamers, and Bait Capture.* Viruses, 2020. **12**(8).

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

679 *Nasopharyngeal Swabs from Patients With COVID-19 and Identification of Deletions in the*

680 *Viral Genome That Encode Proteins Involved in Interferon Antagonism.* Viruses, 2020.

681 **12**(10).

682 5. Dorward, D.A., et al., *Tissue-Specific Immunopathology in Fatal COVID-19.* Am J Respir Crit

683 Care Med, 2021. **203**(2): p. 192-201.

684 6. Graham, R.L., et al., *SARS coronavirus replicase proteins in pathogenesis.* Virus Res, 2008.

685 **133**(1): p. 88-100.

686 7. Pyrc, K., et al., *Genome structure and transcriptional regulation of human coronavirus*

687 *NL63.* Virol J, 2004. **1**: p. 7.

688 8. Hiscox, J.A., D. Cavanagh, and P. Britton, *Quantification of individual subgenomic mRNA*

689 *species during replication of the coronavirus transmissible gastroenteritis virus.* Virus Res,

690 1995. **36**(2-3): p. 119-30.

691     9.      Hiscox, J.A., et al., *Investigation of the control of coronavirus subgenomic mRNA*

692             *transcription by using T7-generated negative-sense RNA transcripts.* J Virol, 1995. **69**(10):

693             p. 6219-27.

694     10.     van Marle, G., et al., *Regulation of coronavirus mRNA transcription.* J Virol, 1995. **69**(12):

695             p. 7851-6.

696     11.     La Monica, N., K. Yokomori, and M.M. Lai, *Coronavirus mRNA synthesis: identification of*

697             *novel transcription initiation signals which are differentially regulated by different leader*

698             *sequences.* Virology, 1992. **188**(1): p. 402-7.

699     12.     Alonso, S., et al., *Transcription regulatory sequences and mRNA expression levels in the*

700             *coronavirus transmissible gastroenteritis virus.* J Virol, 2002. **76**(3): p. 1293-308.

701     13.     Sawicki, S.G., D.L. Sawicki, and S.G. Siddell, *A contemporary view of coronavirus*

702             *transcription.* J Virol, 2007. **81**(1): p. 20-9.

703     14.     Jeong, Y.S. and S. Makino, *Evidence for coronavirus discontinuous transcription.* J Virol,

704             1994. **68**(4): p. 2615-23.

705     15.     Cevik, M., et al., *Virology, transmission, and pathogenesis of SARS-CoV-2.* BMJ, 2020. **371**:

706             p. m3862.

707     16.     Tyson, J.R., et al., *Improvements to the ARTIC multiplex PCR method for SARS-CoV-2*

708             *genome sequencing using nanopore.* bioRxiv, 2020.

709             https://doi.org/10.1101/2020.09.04.283077

710     17.     Freed, N.E., et al., *Rapid and inexpensive whole-genome sequencing of SARS-CoV-2 using*

711             *1200 bp tiled amplicons and Oxford Nanopore Rapid Barcoding.* Biology Methods and

712             Protocols, 2020. **5**(1): p. bpaa014.

713    18.    !!! INVALID CITATION !!! [2].

714    19.    Li, H., *Minimap2: pairwise alignment for nucleotide sequences.* Bioinformatics, 2018.

715           **34**(18): p. 3094-3100.

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

717           *of infection and the inflammatory response: an observational cohort study.* The Lancet,

718           2020. **396**(10251): p. 603-611.

719    21.    Parker, M.D., et al., *Subgenomic RNA identification in SARS-CoV-2 genomic sequencing*

720           *data.* Genome research, 2021. **31**(4): p. 645-658.

721    22.    Yang, Y., et al., *Characterizing transcriptional regulatory sequences in coronaviruses and*

722           *their role in recombination.* Molecular Biology and Evolution, 2021. **38**(4): p. 1241-1248.

723    23.    Anders, S., et al., *Count-based differential expression analysis of RNA sequencing data*

724           *using R and Bioconductor.* Nature protocols, 2013. **8**(9): p. 1765-1786.

725    24.    Alexandersen, S., A. Chamings, and T.R. Bhatta, *SARS-CoV-2 genomic and subgenomic*

726           *RNAs in diagnostic samples are not an indicator of active replication.* Nature

727           communications, 2020. **11**(1): p. 1-13.

728    25.    Ross, M.G., et al., *Characterizing and measuring bias in sequence data.* Genome biology,

729           2013. **14**(5): p. 1-20.

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

731           *model for COVID-19.* Nat Commun, 2021. **12**(1): p. 1260.

732    27.    Ryan, K.A., et al., *Dose-dependent response to infection with SARS-CoV-2 in the ferret*

733           *model and evidence of protective immunity.* Nat Commun, 2021. **12**(1): p. 81.

734    28.    Taiaroa, G., et al., *Direct RNA sequencing and early evolution of SARS-CoV-2.* BioRxiv,

735          2020.

736    29.    Keep, S., et al., *Multiple novel non-canonically transcribed sub-genomic mRNAs produced*

737          *by avian coronavirus infectious bronchitis virus.* J Gen Virol, 2020. **101**(10): p. 1103-1118.

738    30.    Nomburg, J., M. Meyerson, and J.A. DeCaprio, *Pervasive generation of non-canonical*

739          *subgenomic RNAs by SARS-CoV-2.* Genome Med, 2020. **12**(1): p. 108.

740    31.    Corbett, K.S., et al., *Evaluation of the mRNA-1273 Vaccine against SARS-CoV-2 in*

741          *Nonhuman Primates.* N Engl J Med, 2020. **383**(16): p. 1544-1555.

742    32.    Yu, J., et al., *DNA vaccine protection against SARS-CoV-2 in rhesus macaques.* Science,

743          2020. **369**(6505): p. 806-811.

744    33.    Alexandersen, S., A. Chamings, and T.R. Bhatta, *SARS-CoV-2 genomic and subgenomic*

745          *RNAs in diagnostic samples are not an indicator of active replication.* Nat Commun, 2020.

746          **11**(1): p. 6059.

747    34.    Martin, M., *Cutadapt removes adapter sequences from high-throughput sequencing*

748          *reads.* EMBnet.journal, 2011. **17**: p. https://doi.org/10.14806/ej.17.1.200.

749    35.    Joshi, N.A. and J.N. Fass, *Sickle: A sliding-window, adaptive, quality-based trimming tool*

750          *for FastQ files*

751    *(Version 1.33).* 2011: p. https://github.com/najoshi/sickle.

752    36.    Kim, D., B. Langmead, and S.L. Salzberg, *HISAT: a fast spliced aligner with low memory*

753          *requirements.* Nat Methods, 2015. **12**(4): p. 357-60.

754    37.    Li, H., et al., *The Sequence Alignment/Map format and SAMtools.* Bioinformatics, 2009.

755          **25**(16): p. 2078-9.

756    38.    Huang, W., et al., *ART: a next-generation sequencing read simulator.* Bioinformatics,

757          2012. **28**(4): p. 593-594.

758    39.    Yang, C., et al., *NanoSim: nanopore sequence read simulator based on statistical*

759          *characterization.* GigaScience, 2017. **6**(4): p. gix010.

760    40.    Gao, Y., et al., *Quantitative profiling of N 6-methyladenosine at single-base resolution in*

761          *stem-differentiating xylem of Populus trichocarpa using Nanopore direct RNA sequencing.*

762          Genome Biology, 2021. **22**(1): p. 1-17.

763    41.    Dong X; Penrice-Randal R; Goldswain H; Prince T; Randle N; Donavan-Banfield I; Salguero

764          FJ; Tree J; Vamos E; Nelson C; Clark J; Ryan Y; Stewart JP; Semple MG; Baillie JK; M

765          Openshaw PJ; Turtle L; Matthews DA; Carroll MW; Darby AC; Hiscox JA (2022): Supporting

766          data for"Identification and quantification of SARS-CoV-2 leader subgenomic mRNA gene

767          junctions in nasopharyngeal samples shows phasic transcription in animal models of

768          COVID-19 and dysregulation at later time points that can also be identified in humans"

769          GigaScience Database. http://dx.doi.org/10.5524/102209

770

771

772 **Ethics approval and consent to participate**

773 All experimental work on NHPs was conducted under the authority of a UK Home Office approved

774 project license (PDC57C033) that had been subject to local ethical review at PHE Porton Down by

775 the Animal Welfare and Ethical Review Body (AWERB) and approved as required by the Home

776 Office Animals (Scientific Procedures) Act 1986 and the full ethics and NHP model are described.

777 **Consent for publication**

778 Not applicable

779

780 **Data Availability**

781 Illumina and Nanopore test data sets are available under NCBI PRJNA699398. Snapshots of the

782 code are available in the *GigaScience* GigaDB repository[41].

783

784 LeTRS is available at https://github.com/xiaofengdong83/LeTRS.

785

786 **Competing interests**

787 The authors declare that they have no competing interests

788 **Funding**

789 This work was predominately funded by U.S. Food and Drug Administration Medical

790 Countermeasures Initiative contract (75F40120C00085) awarded to JAH.  The article reflects the

791 views of the authors and does not represent the views or policies of the FDA. This work was also

792 supported by the MRC (MR/W005611/1) G2P-UK: A national virology consortium to address

793 phenotypic consequences of SARS-CoV-2 genomic variation (co-I JAH). JAH is also funded by the

813

814    **Authors' contributions**

815  X.D. developed the LeTRS software and performed the informatics analysis. X.D., A.D. and J.A.H.

816  analysed the data. J.S., J.T. and M.W.C. co-ordinated the NHP work and sample processing. R.P.-

817  R., J.P.S., H.G., T.P. and N.R. were involved in sequencing and informatics analysis of the NHP

818  samples with D.A.M. A.D. oversaw sequencing of the human clinical samples with E.V. and C.N

819  for the COG-UK data. R.P.-R. and J.A.H. oversaw sequencing of samples under the auspices of

820  ISARIC-4C with clinical samples collected and managed by J.K.B, L.T., M.G.S. and P.J.M.O. J.A.H.

821  and M.W.C. initiated and led the study and wrote the manuscript with X.D., R.P.-R., A.D. with

822  other authors involved in editing the final version.

828

829

830

831

832

833

834

835

836

837    Table 1. Comparison of other Tools with LeTRS.

| | | LeTRS | Periscope | SARS-CoV-2-leader | SuPER |
|---|---|---|---|---|---|
| Input files | | fastq | fastq | bam/sam | sam |
| Consideration of amplicon primer information used | | yes | yes | no | no |
| Consideration of paired-end Illumina data | | yes | no | no | no |
| Consideration of amplicon primer pool | | yes | no | no | no |
| Consideration of the ACGAAC box | | yes | no | no | yes |
| Support of amplicon Illumina data | | yes | yes | yes | yes |
| Support of amplicon Nanopore data | | yes | yes | yes | yes |
| Support of Nanopore dRNAseq data | | yes | no | yes | yes |
| Method | | Fusion site searching | Sequences tag searching | Sequences tag searching | Fusion site searching |
| Accuracy | ARTIC-Illumina | 1.0000 | 0.9998 | 0.9998 | 0.9996 |
| | ARTIC-Nanopore | 0.9985 | 0.9981 | 0.9980 | 0.9979 |
| | Nanopore dRNAseq | 0.9982 | - | 0.9948 | 0.9937 |
| Sensitivity | ARTIC-Illumina | 0.9997 | 0.9498 | 0.9644 | 0.9230 |
| | ARTIC-Nanopore | 0.6294 | 0.5326 | 0.5154 | 0.4843 |
| | Nanopore dRNAseq | 0.8448 | - | 0.5949 | 0.4817 |
| Specificity | ARTIC-Illumina | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| | ARTIC-Nanopore | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| | Nanopore dRNAseq | 1.0000 | - | 1.0000 | 1.0000 |
| F-measure | ARTIC-Illumina | 0.9998 | 0.9499 | 0.9655 | 0.9243 |
| | ARTIC-Nanopore | 0.7621 | 0.6699 | 0.6611 | 0.6215 |
| | Nanopore dRNAseq | 0.9157 | - | 0.7140 | 0.5934 |

838    Accuracy, sensitivity, specificity and F-measure score were calculated with simulated Illumina and

839    Nanopore sequencing reads for the known sgmRNAs.

840

841

842    Table 2. The criteria of basic and advanced filtering for four different types of input data for

843    LeTRS.

| | Output Filters | Illumina paired-end amplicon reads | Nanopore amplicon reads | Nanopore dRNAseq reads | Bam |
|---|---|---|---|---|---|
| **Basic filtering** | MAPQ > 10 | ● | ● | ● | ● |
| | Read only one splicing junction | ● | ● | ● | ● |
| | Primary alignment only | ● | ● | ● | ● |
| | No supplementary alignment | ● | ● | ● | ● |
| | Read mapped in pair | ● | | | |
| | No read reverse strand | | | ● | |
| **Advance filtering** | Read aligment5' end includes forward primer | ● | ● | | |
| | Read aligment3' end includes reverse primer | ● | ● | | |
| | Read aligment5' end includes forward primer and 3' end includes reverse primer | ● | ● | | |
| | Paired read including at least one primer in each have same leader-TRS junction in alignments | ● | ● | | |
| | Read aligment3' with > 1ployA | | ● | ● | |

43

Read aligment3' with > 5ployA                    ●              ●

844

845     Figures

846     Figure 1. (A). Illustration of reads derived from sgmRNAs mapped onto the SARS-CoV-2 reference

847     genome with a splicing method. We note that splicing does not occur in coronaviruses but this is

848     the apparent observation of a fusion event between different parts of the genome. (B and C).

849     Illustration of the possible type of reads mapped on the SARS-CoV-2 reference genome for the

850     paired-end Illumina amplicon sequencing, where the lines with same colour implied paired reads,

851     (D) Nanopore amplicon sequencing and (E) Nanopore dRNAseq of the SARS-CoV-2 genome and

852     sgmRNAs. L and B in the boxes indicate the leader-TRS breaking sites on the leader side and TRS

853     side, respectively. Although we note these are where the apparent fusion site occurs. Yellow

854     colour indicates the leader region, black is the TRS and gene sequence, the red indicates a

855     sequence read that maps to SARS-CoV-2 sequence. Blue is a sequence that is present between

856     the leader sequence and the TRS. For (B) and (C) the same colour (brown, green and pink)

857     indicates that same paired read. For (B) the paired read contains both primers. For (C) the grey

858     and light blue colour is a paired read, but only contains one primer sequence at any end. The

859     vertical hash lines on (B, C, and D) indicates the position of a primer.

860

861     Figure 2. Analysis of reads mapping to the leader TRS-gene junctions with at least one primer

862     sequence at either end in sequencing data from hACE2-A549 cells infected with SARS-CoV-2 and

863     sequenced using either (A) an ARTIC-Nanopore approach, (B) an ARTIC-Illumina approach and (C)

864     a Nanopore dRNAseq approach. The data corresponds to that shown in detailed in

865     Supplementary Tables 1, 2 and 3. The standard deviation of a binomial distribution was calculated

866     to generate error bars. The data is presented as a histogram with a normalised count for each

867    sgmRNA starting at a particular position in the leader sequence as indicated in the line diagram

868    underneath. For each panel (A, B and C) the expected sgmRNA pattern is shown on the left and

869    novel sgmRNAs are shown on the right.

870

871    Figure 3. An X-Y/scatter plot using normalized counts of sgmRNAs (with greater than 5 A residues

872    at the 3' end – indicative of a polyA tail for the dRNAseq data). To generate the scatter plots

873    Nanopore dRNAseq data was plotted against the either the normalized count (at least one primer

874    sequence) of sgmRNAs with (A) ARTIC-Nanopore sequencing data and (C) ARTIC-Illumina

875    sequencing    data    or    provided    as    ratio    (B)    and    (D),    respectively    for

876    S:orf3:E:M:orf6:orf7a:orf7b:orf8:N:orf10 (using data from Supplementary Tables 3, 4 and 5).

877

878    Figure 4. Analysis of the abundance of reads mapping to the leader TRS-gene junctions that have

879    at least one primer sequence at either end in longitudinal nasopharyngeal samples taken from

880    two non-human primate models infected with SARS-CoV-2. The time post-infection in days is

881    indicated on the x-axis. The normalised count (read count/total number of reads mapped on the

882    reference genome)*1,000,000) of the leader TRS-gene junction abundance is shown on the left-

883    hand Y-axis with each unique leader TRS-gene junction colour coded. The right-hand Y axis is a

884    measure of the total depth of coverage for SARS-CoV-2 in that sample. Note the two scales are

885    different. SARS-CoV-2 was amplified and sequenced by ARTIC-Illumina. The data is organised into

886    groups of animals for the cynomolgus macaque groups 1 and 2 (A/E and B/F), and rhesus

887    macaque groups 1 and 2 (C/G and D/H). E, F, G and H zoom in to see the details of A, B, C and D

888    for Day1 to Day9. The data corresponds to that shown in Supplementary Table 7. Standard

889    deviation of a binomial distribution was calculated to provide error bars.

890

891    Figure 5. Plots of normalised peak counts (A) and peak counts (B) of leader-TRS gene junctions of

892    reads with at least one primer sequences at either end derived from sequence data from 15

893    human patients sequenced with the ARTIC-Illumina approach and analysed by using sequence

894    derived from pool 1 primers. The data correspond to that shown in Supplementary Table 8.

895    Standard deviation of a binomial distribution was calculated to provide error bars.

896

897    Figure 6. Plots of normalised peak counts (A) and peak counts (B) of leader-TRS gene junctions of

898    reads with at least one primer sequence at either end derived from sequence data from 15

899    human patients sequenced with the ARTIC-Nanopore approach and analysed by using sequence

900    derived from pool 1 primers. The data correspond to that shown in Supplementary Table 9.

901    Standard deviation of a binomial distribution was calculated to provide error bars.

902

903    Figure 7. (A). Diagram of novel leader-TRS junctions centred around the known gene orf but out

904    of the search interval in the analysis of SARS-CoV-2 RNA from cell culture, non-human primate

905    and human sequencing data. Many novel junctions map to the leader-TRS membrane gene

906    junctions. (B). Venn diagram showing the overlap of novel leader-TRS gene junctions present in

907    SARS-CoV-2 infected cynomolgus and rhesus macaques, human patients, and Vero cells. Data

908    was obtained using the ATRIC-Illumina method (Supplementary Table 15). (C) Virus genome

909    position of the start of the fusion site (Y-axis) in the leader sequence plotted against the fusion

910    site present in the gene to show the potential positions of the novel leader-TRS junctions along

911    the SARS-CoV-2 genome (indicated above). A shown the colours present the novel leader-TRS

912    junctions identified in the different experimental condition (cynomolgus and rhesus macaques,

913    human patients, and Vero cells).

914 Supplementary Figures

915 Supplementary Figure 1. Bioinformatics pipeline for the identification of leader-TRS junctions in

916 sequencing data from SARS-CoV-2 infected material with LeTRS. This can be rapidly adapted for

917 other coronaviruses such as MERS-CoV and any newly emerged coronavirus. LeTRS can work

918 from Nanopore or Illumina amplicon data or more unbiased approaches such as direct RNA

919 sequencing, metagenomic or Illumina sequencing by using a BAM file.

920

921 Supplementary Figure 2. Novel (leader dependent noncanonical) fusions (count >=2) found in the

922 cell culture test sample sequenced by (A) ARTIC-Nanopore, (B) ARTIC-Illumina and (C) Nanopore

923 dRNAseq approaches; leader independent long-distance (>5,000 nt) fusions (count >=2) found in

924 the cell culture test sample sequenced by (D) ARTIC-Nanopore, (E) ARTIC-Illumina and (F)

925 Nanopore dRNAseq approaches; leader independent local joining yielding a deletion between

926 proximal sites (20–5,000 nt distance) fusions (count >=2) found in the cell culture test sample

927 sequenced by (G) ARTIC-Nanopore, (H) ARTIC-Illumina and (I) Nanopore dRNAseq approaches.

928 The data correspond to that shown Supplementary Tables 1, 2 and 3.

929

930 Supplementary Figure 3. Comparison of different tools and LeTRS to evaluate sequencing data to

931 identify the unique sequencing features of SARS-CoV-2 sgmRNAs. Number of reads were

932 evaluated by LeTRS (all peak count), SARS-COV-2-leader, SuPER or periscope (High Quality count)

933 with the cell culture test sample sequenced by (A) ARTIC-Nanopore, (B) ARTIC-Illumina and (C)

934 Nanopore dRNAseq approaches; (D) Ratio of sgmRNAs

935 (S:orf3:E:M:orf6:orf7a:orf7b:orf8:N:orf10) identified by LeTRS (all peak count), SARS-COV-2-

936     leader, SuPER or periscope (HQ count) with the cell culture test sample sequenced by ARTIC-

937     Nanopore, ARTIC-Illumina and Nanopore dRNAseq approaches. The data are corresponded to

938     that shown in Supplementary Tables 1, 2 and 3.

939

940     Supplementary Figure 4. Comparison of the ratio of reads in amplicon sequencing approaches

941     based on the ARTIC approach, with the forward primer only, reads with reverse primer only and

942     reads with both primers in sgmRNAs to the overall ratio of reads with the forward primer only,

943     reads with reverse primer only and reads with both primers in all reads amplified by pool 1

944     primers, pool 2 primers and both pools of primers for the cell culture test sample sequenced by

945     (A) ARTIC-Nanopore and (B) ARTIC-Illumina approaches.

946

947     Supplementary Figure 5. Raw (A and C) and normalised (B and D) canonical (upper) and novel

948     (lower) leader-TRS gene junctions count in RNA purified from the inoculum of SARS-CoV-2 used

949     to infect either the cynomolgus or rhesus macaques. The RNA was sequenced by the ARTIC-

950     Illumina method (Supplementary Table 6). Standard deviation of a binomial distribution was

951     calculated to provide error bars.

952

953     Supplementary Figure 6. Novel leader-TRS gene junctions (count > 10) identified in RNA purified

954     from nasopharyngeal swabs taken daily from cynomolgus macaques infected with SARS-CoV-2

955     (Supplementary Table 7). The number before "-Day" indicated the group of cynomolgus

956     macaques. Standard deviation of a binomial distribution was calculated to provide error bars.

957

958

959     Supplementary Figure 7. Novel leader-TRS gene junctions (count > 10) identified in RNA purified

960     from nasopharyngeal swabs taken daily from from rhesus macaques (Supplementary Table 7).

961     The number before "-Day" indicated the group of cynomolgus macaques. Standard deviation of

962     a binomial distribution was calculated to provide error bars.

963

964     Supplementary Figure 8. Comparison of the fraction of 6mA modification (right-hand Y-axis) of

965     each site in sgmRNA at 6, 12 and 24 hours after post infection using direct RNA sequencing from

966     RNA purified from SARS-CoV-2 infected cells. Only the sites with modification in at least one of

967     the 6hpi, 12hpi and 24hpi were analysed.

968

969     Supplementary Figure 9. Comparison of the fraction of 5mC modification (right-hand Y-axis) of

970     each site in sgmRNA at 6, 12 and 24 hours after post infection using direct RNA sequencing from

971     RNA purified from SARS-CoV-2 infected cells. Only the sites with modification in at least one of

972     the 6hpi, 12hpi and 24hpi were analysed.

973

974

975

976

977

978

979

980

981    Supplementary Tables

982    Table S1. The LeTRS output tables for known sgmRNA, details of known sgmRNA, novel sgmRNA

983    (count >=2), details of novel sgmRNA, and leader independent long-distance and local fusions

984    (count >=2) evaluated in the cell culture test sample sequenced by the ARTIC-Nanopore

985    approach.

986

987    Table S2. The LeTRS output tables for known sgmRNA, details of known sgmRNA, novel sgmRNA

988    (count >=2), details of novel sgmRNA, and leader independent long-distance and local fusions

989    (count >=2) evaluated in the cell culture test sample sequenced by the ARTIC-Illumina approach.

990

991    Table S3. The LeTRS output tables for known sgmRNA, details of known sgmRNA, novel sgmRNA

992    (count >=2), details of novel sgmRNA, and leader independent long-distance and local fusions

993    (count >=2) evaluated in the cell culture test sample sequenced by the Nanopore dRNAseq

994    approach.

995

996    Table S4. The LeTRS output table for known sgmRNA evaluated by primers of pool 1 and 2 in the

997    cell culture test sample sequenced by the ARTIC-Nanopore approach.

998

999    Table S5. The LeTRS output tables for known sgmRNA evaluated by primers of pool 1 and 2 in the

1000   cell culture test sample sequenced by the ARTIC-Illumina approach.

1001

1002    Table S6. The LeTRS output tables for known sgmRNA and details of known sgmRNA with pool 1

1003    primers, and novel sgmRNA (count > 10) and details of novel sgmRNA with both pools' primers

1004    in the infecting SARS-CoV-2 inoculum source used for the NHP study, sequenced by the ARTIC-

1005    Illumina method.

1006

1007    Table S7. The LeTRS output tables for known sgmRNA and details of known sgmRNA with pool 1

1008    primers, and novel sgmRNA (count > 10) and details of novel sgmRNA with both pools' primers

1009    in longitudinal nasopharyngeal samples taken from two non-human primate models (cynomolgus

1010    and rhesus macaques) of SARS-CoV-2 in groups. SARS-CoV-2 was amplified using the ARTIC

1011    approach and sequenced by Illumina. The data is organised into groups of animals for the

1012    cynomolgus macaque groups 1 and 2 that were with "-1" and "-2" in the excel sheets.

1013

1014    Table S8. The LeTRS output tables for known sgmRNA and details of known sgmRNA in pool 1,

1015    and novel sgmRNA (count > 10) and details of novel sgmRNA with both pools' primers from 15

1016    human patients sequenced with ARTIC-Illumina.

1017

1018    Table S9. The LeTRS output tables for known sgmRNA and details of known sgmRNA in pool 1

1019    from 15 human patients sequenced with ARTIC-Nanopore.

1020

1021    Table S10. The spreadsheet for the 15 human patients sequenced with the ARTIC-Nanopore

1022    detailed in Table S9.

1023

1024   Table S11. Re-analysis of reads for known sgmRNAs in the (NCBI assession No. PRJNA636225)

1025   [24].

1026

1027   Table S12. Summary of normalized count, number of modification sites and average modification

1028   fraction in of each gmRNA at 6hpi, 12hpi and 24hpi.

1029

1030   Table S13. Evaluation of the difference of modification by the paired samples one-sided Wilcoxon

1031   test to calculate p-value by treating the same nucleotides between any two time points as paired

1032   data.

1033

1034   Table S14. The LeTRS output table for novel sgmRNA (count > 10) and details of novel sgmRNA

1035   with both primer pools from VeroE6 cells infected in culture with SARS-CoV-2 (SCV2-006)

1036   sequenced by ARTIC-Illumina primers. This sample is different from the one Table S2.

1037

1038   Table S15. Novel leader-TRS junctions centred around the known gene open reading frame but

1039   out of the search interval in the analysis of cell culture, non-human primate and human

1040   sequencing data.

1041

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5

Figure 6

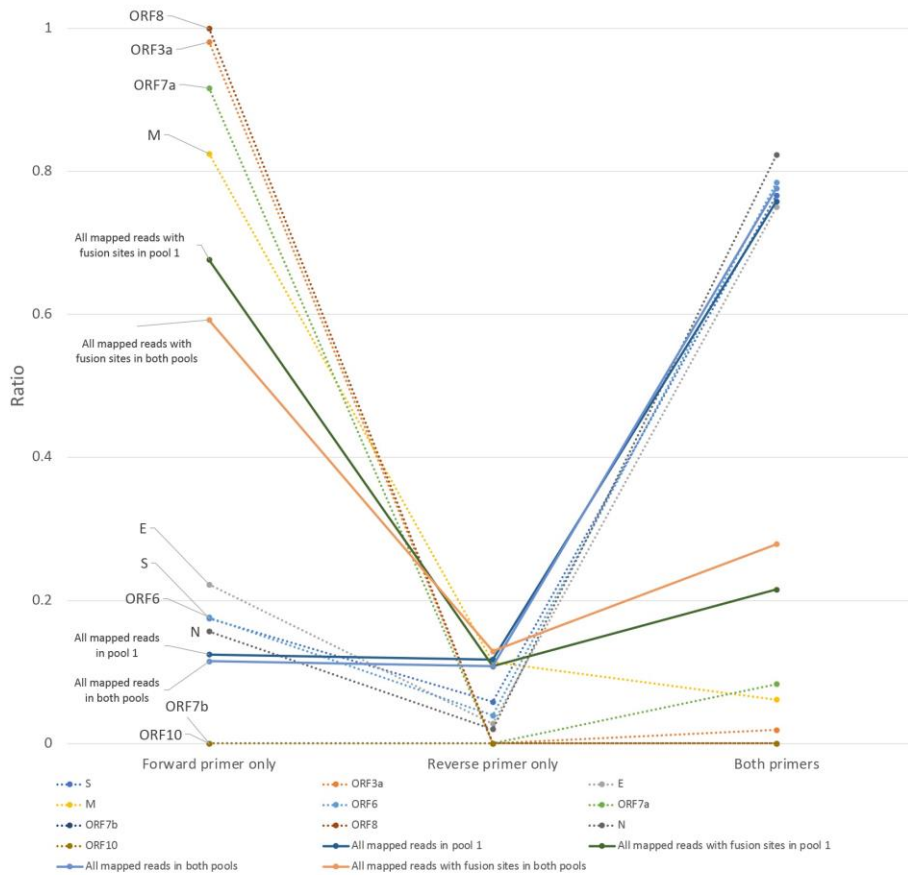Figure 7

Supplementary Figure 1

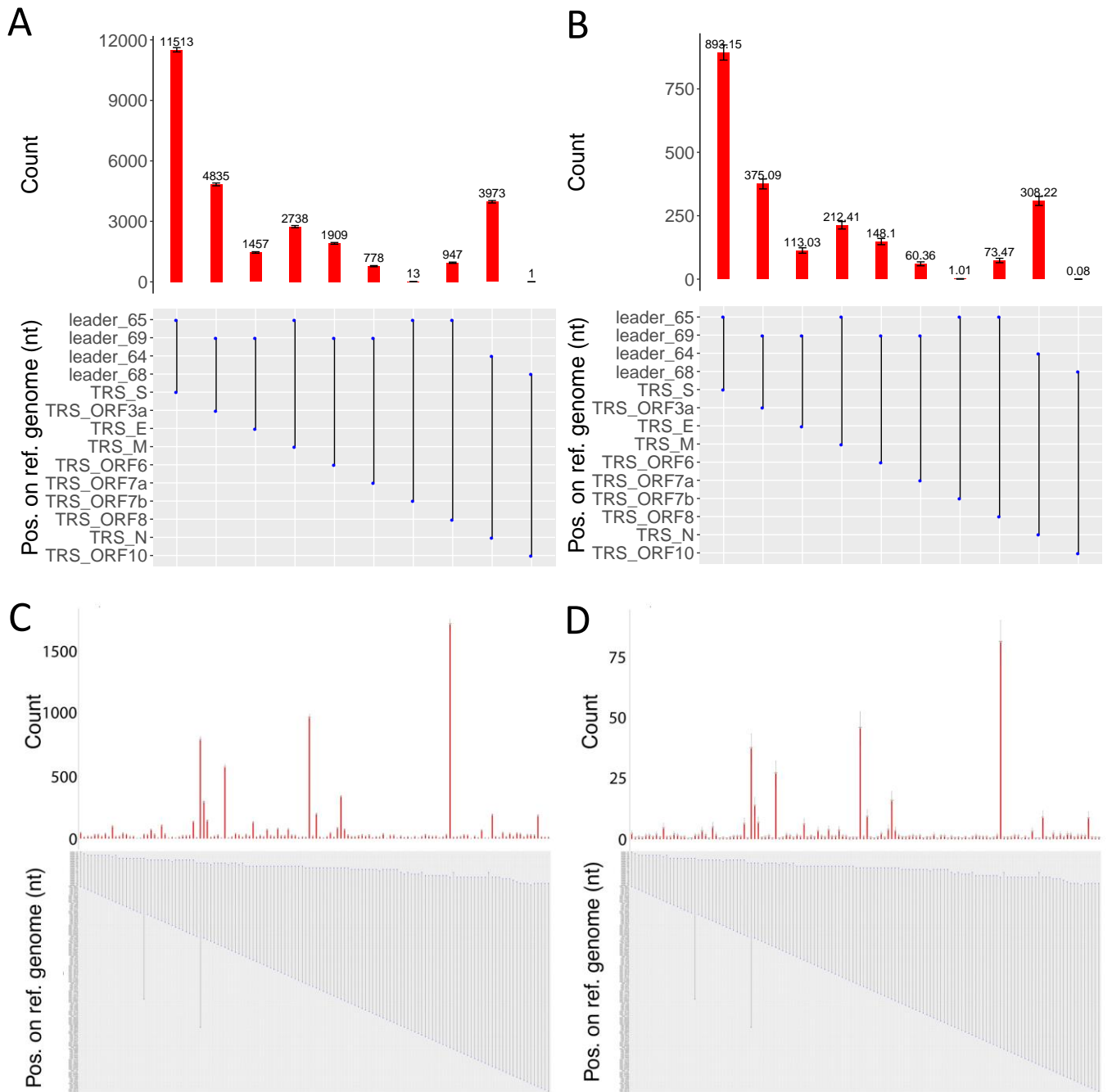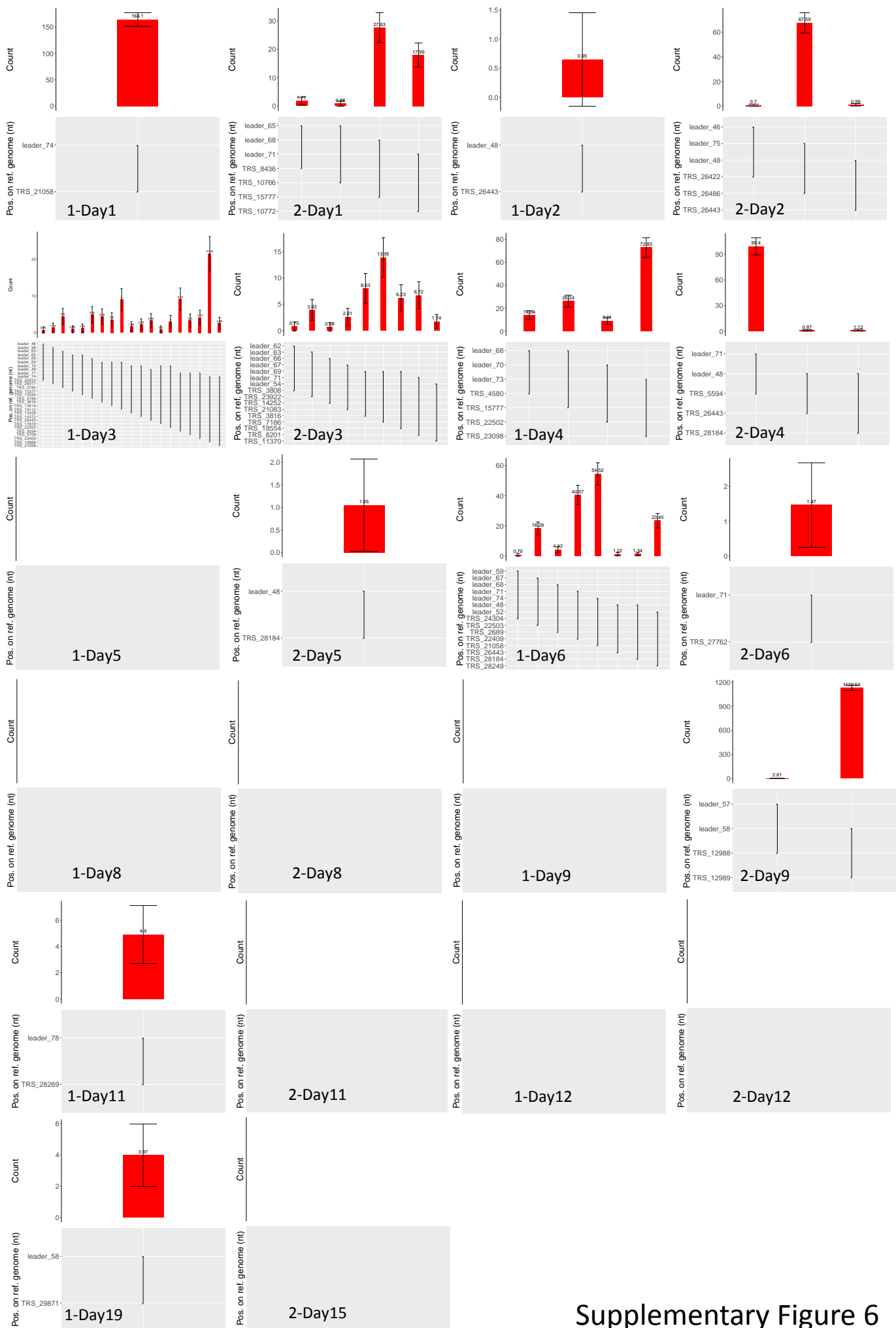Supplementary Figure 2

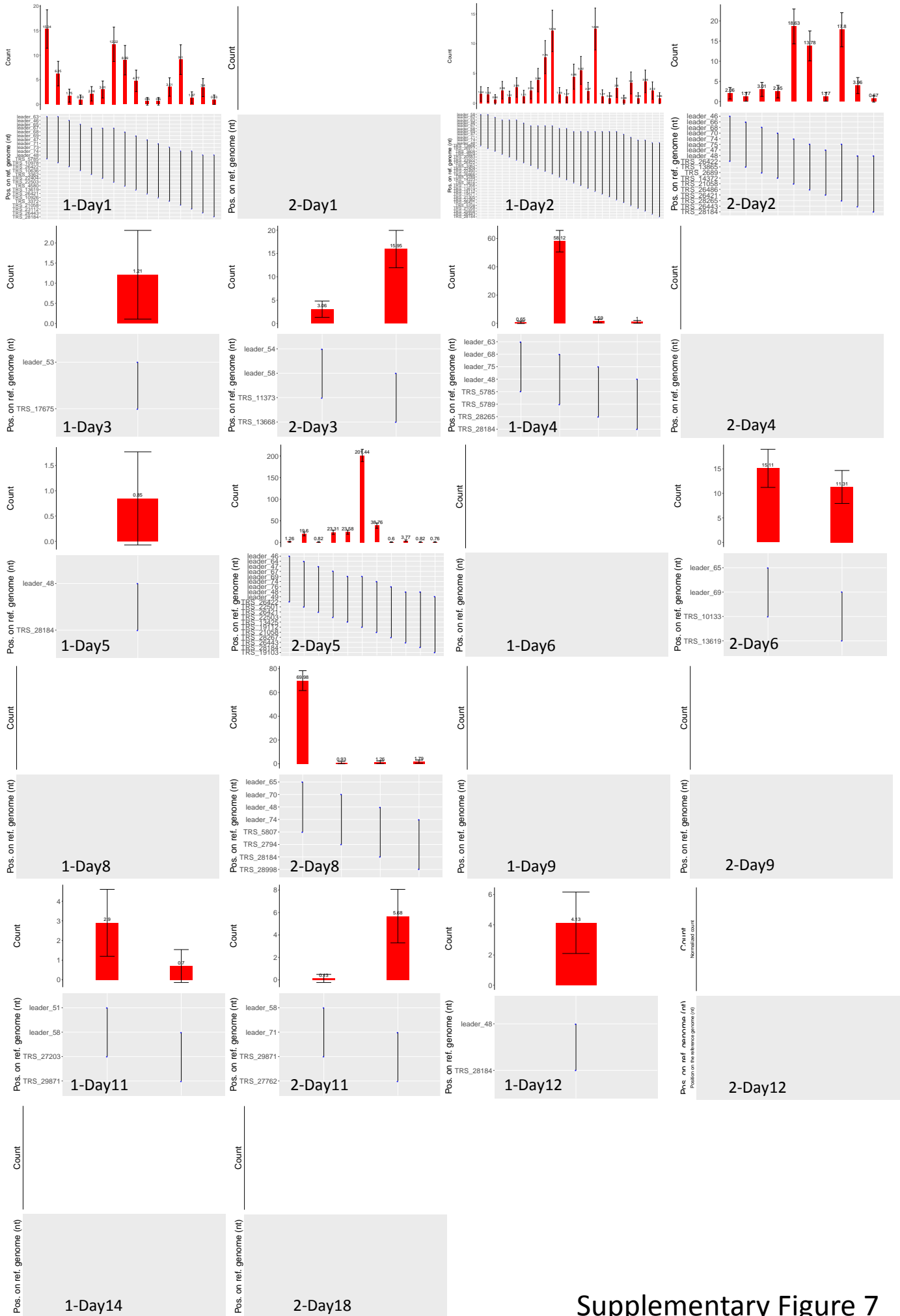Supplementary Figure 3

A



B

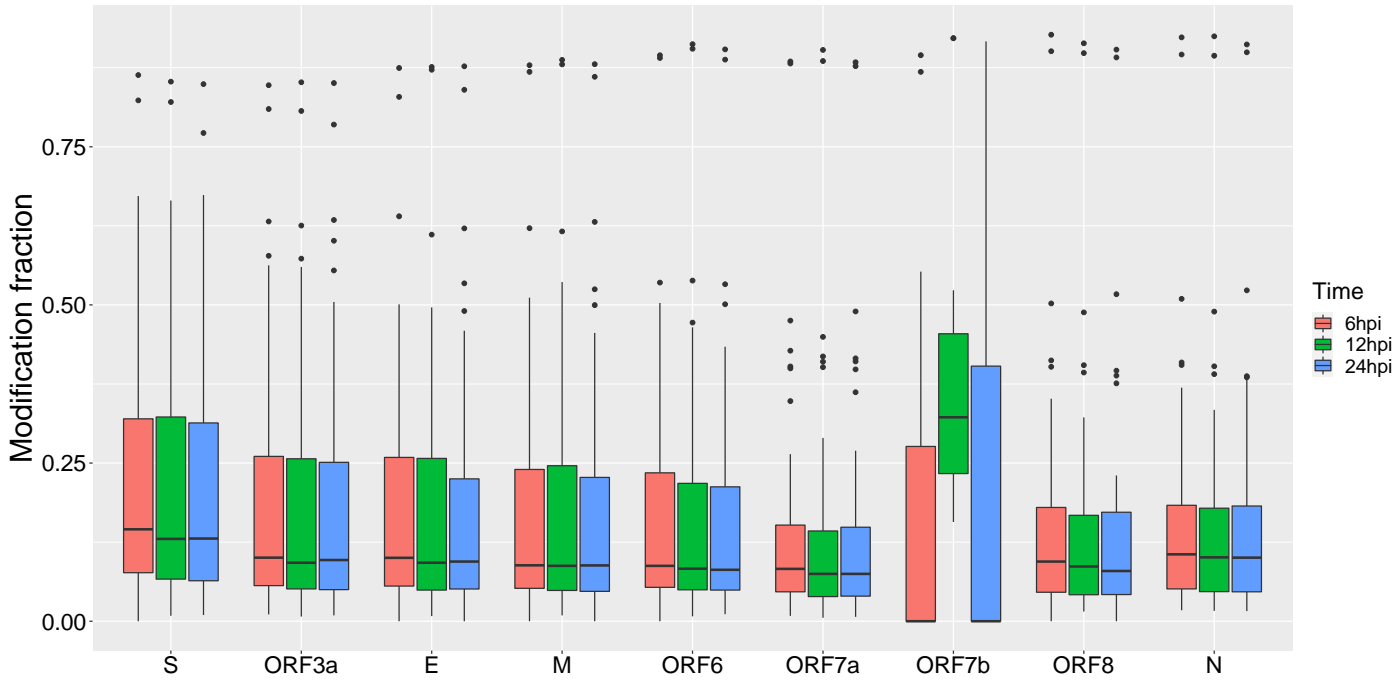Supplementary Figure 4
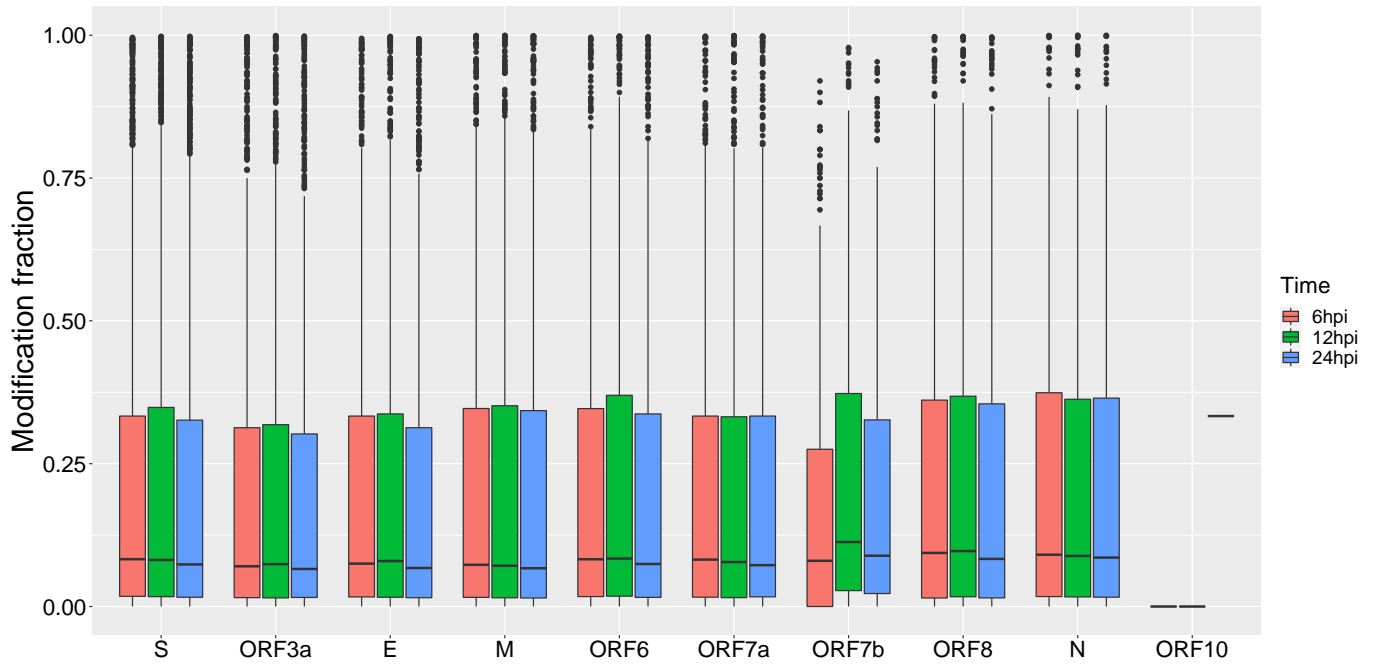
Supplementary Figure 5

Supplementary Figure 6

Supplementary Figure 7

Supplementary Figure 8

Supplementary Figure 9

Click here to access/download
**Supplementary Material**
Supplementary_Table_1.xlsx

Click here to access/download
**Supplementary Material**
Supplementary_Table_2.xlsx

Click here to access/download
**Supplementary Material**
Supplementary_Table_3.xlsx

Supplementary Table

Click here to access/download
**Supplementary Material**
Supplementary_Table_4.xlsx

Click here to access/download
**Supplementary Material**
Supplementary_Table_5.xlsx

Supplementary Table

Click here to access/download
**Supplementary Material**
Supplementary_Table_6.xlsx

Click here to access/download

**Supplementary Material**
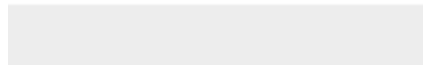
Supplementary_Table_7.xlsx

Supplementary Table

Click here to access/download
**Supplementary Material**
Supplementary_Table_9.xlsx

Click here to access/download
**Supplementary Material**
Supplementary_Table_10.xlsx

Supplementary Table

Click here to access/download
**Supplementary Material**
Supplementary_Table_11.xlsx

Supplementary Table

Click here to access/download
**Supplementary Material**
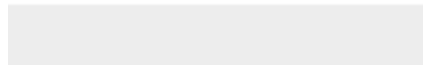Supplementary_Table_12.xlsx

Supplementary Table

Click here to access/download

**Supplementary Material**

Supplementary_Table_13.xlsx

Click here to access/download

**Supplementary Material**

Supplementary_Table_14.xlsx

**Prof. Julian A. Hiscox**

**Chair in Infection and Global Health**

**Deputy Associate Pro-Vice Chancellor
Research and Impact (FHLS)**

The University of Liverpool
Department of Infection Biology
Institute of Infection, Veterinary and Ecological
Sciences
Liverpool Science Park IC2
146 Brownlow Hill
Liverpool
L3 5RF

**Tel:** +44 (0)7812238359.
**Email:** julian.hiscox@liverpool.ac.uk

Dear GigaScience

Many thanks for reviewing our manuscript describing a bioinformatic tool we developed to study coronavirus biology, specifically demonstrated on clinical and model samples infected with SARS-CoV-2. We very much appreciate the very constructive and detailed reviews. Below we detail our point-by-point responses (in red) to the thoughts and suggestions of the reviewers (in black). We have acted on all these new comments and conducted the additional experiments that the reviewers wanted. We provide a marked-up manuscript of the initial revised version showing alterations from the original submitted version and a clean version with all changes etc accepted.

Yours sincerely,

Prof. Julian A. Hiscox.

Reviewer reports:

Reviewer #1: Comments: It is an important study. Except for a few minor points, the authors have addressed most of the reviewers' suggestions. This manuscript will be considered for acceptance after addressing the following minor suggestions:

1.  The authors have compared the algorithm design, input, and output, and the counts of predicted sgmRNA across four tools. However, it would be nice if the authors could compare these tools' performances regarding prediction accuracy, F-measure, sensitivity, and specific scores. These will let the readers and potential users have a better sense of choosing a different tool for different purposes.
<span style="color:red">We have added the prediction accuracy, F-measure, sensitivity, and specific scores, calculated based on simulated Illumina and nanopore reads, in the Table 1.</span>

2.  It is unclear what the red line means in Supplemental Figure 8-9.
<span style="color:red">The red lines in Supplemental Figure 8 and 9 are for the normalized count of sgmRNA identified by LeTRS. We have moved this to Supplementary Table 12.</span>

3.  On page 18, lines 364-370. The analysis and significance that the authors stated in that paragraph do not show the apparent trends in Supplemental Figure 9. Would the authors update the figure types to reflect the results of their statistical tests?
<span style="color:red">We have updated the boxplots in Supplemental Figures 8 and 9. We used a paired samples one-sided Wilcoxon test that takes account the difference at each modification site of two compared sgmRNAs in different time points. A large amount of modification sites with differences resulted a low p-value even the trends in boxplots are not very large.</span>

4.  On page 18, line 370. The author mentioned that "The abundance of most sgmRNAs decreased with time, and both of these factors could account for the frequency of methylation." Based on the context, it seems that the conclusion could not be derived. Because the methylation frequency is a ratio, then it may not correlate with the abundance of the sgmRNAs.
<span style="color:red">We have removed this sentence to reflect the reviewer's content.</span>


Reviewer #2: Happy with revisions, no further comments