

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-21-00142R1	
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	
Funding Information:	U.S. Food and Drug Administration Medical Countermeasures (75F40120C00085)	Prof. Julian A. Hiscox
	MRC ((MR/W005611/1) G2P-UK)	Prof. Julian A. Hiscox
Abstract:	<p>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 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.</p>	
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 l'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:	<p>Reviewer #1: Comments: In this manuscript, the authors sequenced the SARS-CoV-2 transcriptomes of nasopharyngeal samples from 15 patients using both illumina sequencing and nanopore ARTIC primer3 aplicom sequencing, and developed a computational-pipeline called LeTRS to identify the junctions between the leader sequences in the 5' end of viral genome and the transcriptional regulatory sequence (TRS) within the viral genome (leader-TRS-junction). They first tested and applied their LeTRS tool in several published Nanopore RNA-sequencing data and their own sequencing data to analyses leader-TRS sequence information. They showed that the expression abundance and populations of viral subgenomic mRNA (sgmRNAs) with leader-TRS varies along the time points of post-infection. This study is important to understanding SARS-CoV-2 pathology. However, this article needs many improvements. My major suggestions are as follows:</p> <p>1. There are two types of leader sequences found in the SARS-CoV-2 sgmRNAs (Dongwan Kim et al., Cell 2020): leader with or without a TRS inside. In the current manuscript, the authors has used their LeTRS tool to identify the sgmRNAs with typical leader with TRS, but did not find the sgmRNAs with non-canonical leaders which do not include TRS inside (TRS-L-independent). I would suggest authors to further extend the studies to sgmRNAs with non-canonical leaders.</p> <p>Of note, the junctions in these noncanonical transcripts are not derived from a known TRS-B. Some junctions show short sequences (3–4 nt) common between the 50 and 30 sites, suggesting a partial complementarity-guided template switching (“polymerase jumping”). However, the majority do not have any obvious sequences. Thus, we cannot exclude a possibility that at least some of these transcripts are generated through a different mechanism(s).</p> <p>[Respond to comment 1: We have added a function in LeTRS to find sgmRNAs with non-conical leaders (TRS-L-independent) with the “-TRSLindependent” function. This function has been evaluated with the test sample (sequencing RNA from cells infected with SARS-CoV-2) as shown in Supplementary figure 2.]</p> <p>2. SARS-CoV2 genomic and subgenomic mRNAs has multiple types of RNA modifications, such as m6A, 5mC, etc. These modifications has been shown to be regulated and relevant to their polyA tail lengths in sgmRNAs (Kim et al., Cell 2020). I would suggest authors to address if and how RNA modifications levels or types will be</p>

dynamically relevant to sgRNA expression at different time points of post-infection. Also any preference of RNA modifications in certain types of sgRNAs (e.g. sgRNA: S which encodes spike-proteins).

[Respond to comment 2: We have directly sequenced the cell culture samples infected with SARS-CoV-2 at three time points for investigating the relationship between RNA modifications to sgRNA expression as shown in Supplementary Figures 8 and 9 and Supplementary Table 12. We specifically searched for two different types of methylation. We note that we can only sequence RNA from cell culture using direct RNA sequencing on the Nanopore. We have found that RNA concentration and quality in clinical samples was insufficient for direct RNA sequencing.]

3.I would suggest the authors to compare and evaluate the performance of their LeTRS tools with other similar tools, such as SuPER (Yang Y. et al., Mol. Biol. Evol. 2020), and SARS-CoV-2-leader (Alexandersen S. et al., Nature Communications 2020), to discuss the strength and weakness of their tool, though the authors have compared their LeTRS tool with another one (Periscope).

[Respond to comment 3: We have compared LeTRS with the tools listed by the reviewer using our test data (total RNA from cells infected with SARS-CoV-2) sequenced by three different approaches –ARTIC-Nanopore, ARTIC-Illumina and direct RNA sequencing. This data is presented in Table 1 and Supplementary Figure 3 A, B, C and D. We compare and contrast what the different tools have in common in terms of analysis function and what data types they can function with.]

4.I would suggest the authors to re-analyze the public patient's seq data (NCBI PRJNA636225) to examine if the same conclusion about the dysregulation of sgRNAs at later time points could be derived in different groups of patients.

[Respond to comment 4: We have reanalyzed sequencing data from a longitudinal study in two patients (NCBI PRJNA636225) using LeTRSs. The results also indicated a dysregulation of sgRNAs in late infection from the two patients (Supplementary Table 11). Apart from nuclease resistance and protection by cellular membranes, a phasic pattern of sgRNA synthesis may also contribute to the presence of sgRNAs at later time points.]

5.It would be nice to have a table to summarize the samples and individual information in this study, such as clinical symptoms of patients, gender and age group, and sample collection time point after infection.

[Respond to comment 5: Due to the different pathways clinical samples were obtained patient identifying information was not available. For example, samples sequenced using ARTIC-Nanopore were obtained via ISARIC-4C and some patient information was obtained (likely due to these being hospitalized cases – either for treatment or isolation). This is shown in Supplementary Table 10. Samples sequenced using ARTIC-Illumina were sequenced under the auspices of COG-UK and identifying information was not available.]

6.The dataset ID provided by this paper (NCBI PRJNA699398) could not be found in the NCBI database. Please the authors address this problem and make the dataset available for the public with a correct ID.

[Respond to comment 6: There is a link provided for reviewers:
<https://dataview.ncbi.nlm.nih.gov/object/PRJNA699398?reviewer=tro3da1gmlid1kk6mdjndh7pg0o>
We will release the data if the paper is accepted.]

7.The overall presentation, Figures, Tables and language of the paper could need some substantial improvement. The current manuscript includes many misused words, misused punctuation, grammatical errors, and mislabeling.

For examples:

(1) the title is too long. The author should conceive a title with concise but to the key-point.

[Respond to comment 7-1: We have shortened the title.]

(2) on page 4, the sentence "for SARS-CoV-2 the core motif is ACGAAC" could be revised as "The core motif of the TRS in SARS-CoV-2 is ACGAAC".

[Respond to comment 7-2: We have changed this.]

(3) on page 5, "cell infected in culture" is inaccurate. It could be expressed as "cultured cells with infection".

[Respond to comment 7-3: We have changed this.]

(4) on page 13, the word "commonality" might be replaced by "Common properties/features".

[Respond to comment 7-4: We have changed this.]

(5) the last sentence on page 13 also need language editing.

[Respond to comment 7-5: We have changed this.]

(6) on page 21, the subtitle "search leader-TRS" would be "searching leader-TRS". Pls keep the subtitle to be a short phrase, rather than beginning with a verb.

[Respond to comment 7-6: We have changed this.]

(7) pls keep the references in a consistent format. Pls correct the format of Ref. 26, 29 and 30 on page 25-26.

[Respond to comment 7-7: We have changed this.]

(8) The authors just need to acknowledge the COG-UK consortia and ISARIC4C consortia, rather than list names of all members in the consortia which occupy 8 pages' space.

[Respond to comment 7-8: We have removed these, apologies this was due to original rules around the consortium authorship statements/acknowledgements.]

(9) The x or y bar label and scales in most figures/suppl figures are too small to read.

[Respond to comment 7-9: We have increased the font on the labels.]

(10) The Figure legends of all figures are not clear enough and does not provide enough illustrations and explanations for the figures (e.g. Fig 1).

[Respond to comment 7-10: We have changed and expanded the Figure legends.]

(11) Supplemental Fig1 could be re-designed to be more clear. For instance, the authors can merge the same steps after the step of <SAM> or <BAM>, to avoid redundant information.

[Respond to comment 7-11: We have changed this.]

(12) The legend of table 8 seems exactly same as the legend of table 2. Pls check it.

[Respond to comment 7-12: We have changed this, Tables 1-8 have been moved to Supplementary Tables 1, 2 and 3.]

Reviewer #2: "Identification and quantification of SARS-CoV-2 leader subgenomic mRNA gene junctions in nasopharyngeal samples shows phasic transcription in animal models of COVID-19 and dysregulation at later time points that can also be identified in humans"

In this paper, Dong et al describe a new pipeline for identifying subgenomic mRNA from multiple types of sequence data, including amplicon (Illumina and Nanopore) as well as long read nanopore direct RNA or cDNA sequencing. It is useful to have a bioinformatics pipeline which can rapidly identify sgRNA in multiple types of sequence data and has the potential to open large amplicon datasets in particular for further analysis of sgRNA abundance. However, I believe that more validation of the accuracy of abundance estimates from amplicon data is required in order to give the research community more confidence in its use (and limitations).

Major comments:

1. More explanation/detail on methodology would be useful. The authors say that they find the most common peak for the break points of the disjunction site amongst all reads with a break point within a 20bp window of the expected breakpoint. Is there a threshold applied in terms of the difference between the most common and next-most-common breakpoint? Also for the novel sites, is there a clustering algorithm applied, or any site with more than 10 reads is reported?

[Respond to comment 1: We used the 20bp window (± 10 bp) of the true splicing sites (known) splicing sites for searching the known sgmRNAs. As noted in the manuscript although we refer to splicing – this is a fusion event. As the minimap2 paper indicated “When INDEL sequencing errors are frequent, it is difficult to find precise splicing sites in this case. If we allow up to 10bp distance from true splicing sites, 98.4% of aligned introns are approximately correct.” (<https://doi.org/10.1093/bioinformatics/bty191>). Because the known breakpoints are far from each other, the threshold was not defined between the most common and next-most-common breakpoint for the known breakpoints.

We used the coverage cut-off (>10 by default) for the novel sites because we found the novel sites usually have low sequence coverage and don't have a cluster like the known sgmRNAs. Alternatively, these novel sites could be due to RT and sequencing errors, and we note this in the manuscript. LeTRS reports these unknown sites as potential novel sites for future research as all other novel sgmRNAs in the research data.]

2. I would like a more direct comparison of sgRNA abundances estimated from amplicon based approach, vs using nanopore amplicon free approach? Its possible to do this only by comparing different tables. It would be easier to digest if there was a x-y plot comparing abundances from different approaches on the same sample. This would help give confidence that the amplicon based approach can provide good estimates. From looking at the tables 1 and 2, it seems that the amplicon approach estimates a lot less sgRNA than the amplicon free approach overall (in terms of normalized counts per million mapped reads). This is to be expected as most of the reads from the amplicon sequencing would be expected to come from the genome. It would be good to see which ORFs are under- and over- represented in the amplicon data, as I imagine this would also relate to which primer pairs are in the same amplicon pool in the arctic design.

Related to this, it would be good to have an analysis of how the primer design impacts detection of sgRNA. For example, I thought that only one of the primer pools includes a leader primer.

[Respond to comment 2 part 1: To address this question we infected cells in culture with SARS-CoV-2 and sequenced the viral RNA using three different approaches. Two were amplicon based – based around the ARTIC protocol (an amplicon based system) and also by direct RNA sequencing. This data is shown in Supplementary Tables 1, 2 and 3 to replace the old test data in the Tables 1-8.

With the Artic V3 pipeline, we used two primer pools for the PCR reactions in the whole virus amplification. Please find the primers used in the primer pool 1 and pool 2 at https://github.com/artic-network/artic-ncov2019/blob/master/primer_schemes/nCoV-2019/V3/nCoV-2019.tsv. For the Artic V3 pipeline, only the pool 1 includes a

5'(forward) primer located within the leader region (about < 80) on the genome (please find the position of Artic V3 primers on the virus genome at https://github.com/artic-network/artic-ncov2019/blob/master/primer_schemes/nCoV-2019/V3/nCoV-2019.primer.bed). The LeTRS (v2.0.1) has been modified to only identify the reads with primers in the pool 1, pool 2 or both pools. We compared the read counts evaluated by LeTRS in both ARTIC-Nanopore and ARTIC-Illumina test data for pool 1 and 2, and found only very few reads/read pairs contained the reverse primers with primer pool 2 (Supplementary Table 4 and 5), suggesting the primers in Artic pool 2 are almost not involved the sequencing of leader-TRS regions.

We have done the x-y plot as showed in Figure 3A and C for the reads with at least a primer sequence comparing abundances from different approaches on the same sample. The normalized counts showed a linear relationship between the amplicon based method to the direct sequencing method, while The Artic-Nanopore and Artic-Illumina showed same ratio of known sgmRNA as the nanopore direct RNA sequencing approach, except S and orf7a (Figure 3B and D for the reads with at least a primer sequence). This suggested an amplicon based approach can provide good estimates for most of the sgmRNAs, especially for N. This normalization method has been applied by <https://doi.org/10.1101/gr.268110.120> and <https://doi.org/10.1038/s41467-020-19883-7>.

PCR based approaches boosted value of denominator reduced the normalized count because a full length of mRNA is counted once with direct RNA sequencing approach will be counted many times with its the small amplicons. Artic illumina got even smaller normalized counts than Artic nanopore approach due to the probably the sequencing bias of illumina during bridge PCR (<https://doi.org/10.1186/gb-2013-14-5-r51>). Therefore, the normalized counts can only be used for the comparison of samples sequenced by same approach when that resulted same PCR and sequencing machine effects. The difference of normalized counts in the samples from amplicon based methods only indicate the relative difference.]

Further related to this, it would be good to have a plot which shows the proportion of read counts which are derived from left-primer only, right-primer only or both primers for each sgRNA, and how this compares to the overall ratio of left-only and right-only primers. It seemed odd to me at first glance that there are so many one-sided amplifications, but I imagine this is a small proportion overall, but a sizeable proportion of the reads which can identify sgRNA, due to the lack of primer pairs for many of the sgRNA. Based on this analysis, it would also be interesting to estimate what is the best depth of coverage of the amplicon panels to get reliable estimates of sgRNA abundance across the different ORFs.

[Respond to comment 2 part 2: We compared the ratio of reads with forward primers only and reverse primers only and both primers for each sgmRNAs to the overall ratios of reads with forward primers only and reverse primers only and both primers in all mapped reads of pool 1 and pool 2 and the mapped reads with any fusion sites in pool 1 and pool 2, found overall ratios showed abundant reads showed same pattern as the reads for sgmRNAs (Supplementary Figure 4). This suggested the mass of one side amplification is a nature of amplicon sequencing.]

3. It would be good to compare the novel breakpoints with those previously reported, e.g. in Taiorara et al, figure 2 and supplementary figure 6 (<https://doi.org/10.1101/2020.03.05.976167>). I can see that many of them line up with those you report in table 4, and I believe this sup

[Respond to comment 3: Taiorara et al didn't attach the exact breakpoints positions with their figure, but we generated a similar figure for comparison (Figure 7c). Figure 7c showed some similar breakpoints positions with Figure 2 of Taiorara et al's paper.]

4. Is there much overlap in the novel break points detected using nanopore amplicon ARTIC v3 vs nanopore dRNA? It would be good to have an extra column in Table 8 and table 4 indicating which of the breakpoints discovered in dRNA were also discovered in amplicon sequencing and vice versa. This will hopefully shed light on relative strengths of the two approaches. Similarly it would be useful to compare nanopore ARTIC and illumina ARTIC in this regard

[Respond to comment 4: As described above we have moved the new test data from a

unique cell culture sample to Supplementary Tables 1, 2 and 3 for Artic-Nanopore, Artic-Illumina and nanopore direct RNA sequencing. We didn't find any exactly the same novel fusion sites in these three approaches. To note in the publication describing minimap2 the paper details "In general, minimap2 is more consistent with existing annotations. It finds more junctions with a higher percentage being exactly or approximately correct" and "When INDEL sequencing errors are frequent, it is difficult to find precise splicing sites in this case. If we allow up to 10bp distance from true splicing sites, 98.4% of aligned introns are approximately correct." (<https://doi.org/10.1093/bioinformatics/bty191>). Therefore, it is very difficult to identify the exact novel fusion sites. Novel leader-TRS junctions were also known as leader dependent noncanonical fusions. LeTRS also has a function to identify leader independent long-distance (>5,000 nt) fusion and local joining yielding a deletion between proximal sites (20–5,000 nt distance) in the sequencing reads. If we look at the pattern of the fusion sites, some of the novel leader-TRS junctions (noncanonical fusions) and leader independent fusions in the test sample were supported by all three sequencing methods (Supplementary figure 2) with similar fusion sites.

The strength of LeTRS to identify the known breakpoints is much stronger than identifying novel sites, because LeTRS controls the aligner to search the known breakpoints with the guide of known annotations. As the paper said "In general, minimap2 is more consistent with existing annotations. It finds more junctions with a higher percentage being exactly or approximately correct" (<https://doi.org/10.1093/bioinformatics/bty191>).

5. Its hard to assess the evidence supporting the biphasic expression without having some idea of the error in the abundance estimates (also commented on this more below);

[Respond to comment 5: We have calculated the standard deviation of a binomial distribution as error bar. The data supports that biphasic expression/abundance of sgRNAs occurs.]

6. The conclusion of dysregulation in samples taken from patients many days into their infection is made only on a small number of samples. Also in Figure 4, the time post sample is not indicated. I presume the information is in one of the supplementary tables, but the submitted pdf has messed up these tables (its somewhere in the 729 page pdf). Nevertheless, it seems that the data supporting this conclusion is a bit thin, and I would be cautious in including that observation in the title of the paper.

[Respond to comment 6: We have changed the title to reflect this comment.]

Minor comments:

1. In figure legends (e.g. figure 1) you say the numbers in brackets are: reads with left primers, reads with right primers, reads with both primers. I can see from the numbers that these are not exclusive, but it might be easier to digest if you showed left-only, right-only and both

[Respond to minor comments 1: We modified the LeTRS to show forward-only, reverse-only and both primers.]

2. You make a point in the paper about whether the left break occurs at position 64 or 69. One thing I would worry about is that microhomology between TRS-L and TRS-R might make it difficult to be exactly sure of the breakpoint (because the sgRNA includes only one copy of the TRS, but its hard to know if it's the left or the right which is included, the aligner could equally well align to TRS-L and skip TRS-R or vice versa, and this would shift the coordinates slightly. Are there enough snp differences in TRS-L or TRS-R to be confident either way, and if so, does this have implications for whether TRS-L or TRS-R is retained in the sgRNA?

[Respond to minor comments 2: For the known sgRNA, we used the known annotation of breakpoints to guide the alignments and allowing a (± 10 bp) window of the true splicing/fusion sites for searching the breakpoints - if this would shift the coordinates slightly. Even if TRS-L or TRS-R is retained in the sgRNAs, the

implications will be random and equal to all samples with same sequencing approach and alignment tool. This should not affect the evaluation of the ratio of sgmRNAs and relative abundance across samples. We have also compared the number of reads for sgmRNAs with the other methods (tool called SARS-CoV-2-leader) that is to search a tag sequence within leader in reads but not the breakpoints of reads. SARS-CoV-2-leader produced a similar read count as LeTRS for the Artic-Nanopore (Supplementary Figure 3A) and Nanopore direct sequencing (Supplementary Figure 3C). SARS-CoV-2-leader produced more counts than LeTRS for Artic-Illumina, because LeTRS counts the read pairs but not reads (Supplementary Figure 3B). There are difficulties in searching for novel breakpoints, although we treat novel breakpoints as a potential sign of novel sgmRNAs for future research.]

3. Figure 1 panels B,C,D were a bit confusing. Why is the reference sequence in the middle. It would be good if the caption could be expanded to help the reader understand these panels in particular.

[Respond to minor comments 3: The figure legend has been changed but we would like to keep the reference sequence in the middle to show the forward and reverse amplification possibilities.]

4. The tables (table 1 to 8) and the figure 1A represent a lot of the same information, but the numbers don't line up exactly, because in the figures you only use counts which have both primers. It would be best to decide which to represent because it's confusing to have the same data presented twice essentially but in slightly different ways.

[Respond to minor comments 4: We have changed this and now consistently only used the reads containing at least one primer to plot data.]

5. In figure 1 you present the normalized abundance to 2 decimal places, but it's very unlikely that you have that level of precision. It would be good if you could add error bars to estimate the uncertainty in the abundance estimate (e.g. calculated using a binomial distribution).

[Respond to minor comments 5: We have calculated the standard deviation of a binomial distribution as an error bar.]

6. In figure 3, it's hard to know how much error there is in each of the measurements. By showing the normalized value, it's also hard to see what is the absolute change in the read counts. Ideally you would show either the read counts, or show error bars around the abundance estimates.

[Respond to minor comments 6: We now show error bars.]

7. Is there a mistake in the title of Table 8: "The LeTRS output table for novel sgmRNA in the tested Nanopore ARTIC v3 primers." Because the title of table 2 seems the same: "Table 2. The LeTRS output table for novel sgmRNA in the tested Nanopore ARTIC v3 primers". One of these approaches does not seem to find novel breakpoints, but the other does, presumably Table 8 should be Illumina based on the ordering?

[Respond to minor comments 7: We have changed this. Tables 1-8 have been moved to Supplementary Tables 1, 2 and 3.]

8. Error in caption of table 1: "Normalized count=(Read count-Total number of read mapped on reference genome)*1000000"

[Respond to minor comments 8: We have changed this. Tables 1-8 have been moved to Supplementary Tables 1, 2 and 3.]

9. In the supplementary figures, the captions you say: "Supplementary Figure 3. Raw (A and C) and normalised (B and D) expected (upper) and novel (lower) leader-TRS gene junctions count in the infecting SARS-CoV-2 inoculum source used for NHP study, sequenced by Illumina ARTIC method (Supplementary Table 8)."

	<p>I found the use of "expected" here confusing, because it implied to me that you had estimated expected counts. I would prefer the use of the term canonical, or something like that.</p> <p>[Respond to minor comments 9: We have changed "expected" to "canonical". Supplementary Figure 3 has become Supplementary Figure 5.]</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>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.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>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.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>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</p>	Yes

appropriate), referencing such data using 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 Dong¹, Rebekah Penrice-Randal¹, Hannah Goldswain¹, Tessa Prince¹, Nadine Randle¹,
5 I'ah Donovan-Banfield^{1,2}, Francisco J. Salguero³, Julia Tree³, Ecaterina Vamos¹, Charlotte Nelson¹,
6 Jordan Clark¹, Yan Ryan¹, James P. Stewart¹, Malcolm G. Semple^{1,2}, J. Kenneth Baillie⁴, Peter J. M.
7 Openshaw⁵, Lance Turtle^{1,2}, David A. Matthews⁶, Miles W. Carroll^{2,3}, Alistair C. Darby¹ and Julian
8 A. Hiscox^{1,2,7}.

9

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

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

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

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

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

15 ⁶University of Bristol, UK.

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

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

18

19 **Abstract**

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

41 **Importance**

42 When infecting cells, SARS-CoV-2 not only replicates its genome but also makes molecules called
43 subgenomic mRNAs (sgmRNAs) that are used as the template for many of the viral proteins,
44 including the spike glycoprotein. The sgmRNAs can only be found in infected cells, and therefore
45 their presence and ratio in a clinical sample is indicative that viral RNA synthesis has occurred,
46 and infected cells are present. The sgmRNAs are targets for diagnostic assays. We have developed
47 a rapid informatics methodology (LeTRS) to identify these unique molecules from multiple types
48 of sequencing data generated in response to the COVID-19 pandemic. We used this pipeline to
49 follow the pattern of sgmRNA abundance in nasopharyngeal samples taken from non-human
50 primate models and clinical samples from humans. We identified putative novel sgmRNAs that
51 may point to a potential new evolutionary mechanism in the virus. The data indicated that SARS-
52 CoV-2 RNA synthesis (and by inference infection) may occur in waves, and this has implications
53 for diagnostics and modelling of disease spread.

54

55 **Introduction**

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

75

76 The prevailing thought is that synthesis of sgmRNAs involves a discontinuous step during negative
77 strand synthesis[13, 14]. A natural consequence of this is recombination resulting in insertions
78 and deletions (indels) in the viral genome and the formation of defective viral RNAs. Thus, the
79 identification of the leader/sgmRNA complexes by sequencing provides information on the
80 abundance of the sgmRNAs and evidence that transcription has occurred in the tissue being
81 analysed. In terms of clinical samples, if infected cells are present, then leader/sgmRNA ‘fusion’
82 sequence can be identified, and inferences made about active viral RNA synthesis from the
83 relative abundance of the sgmRNAs. In the absence of published data from human challenge
84 models, the kinetics of virus infection are unknown, and most studies will begin with detectable
85 viral RNA on presentation of the patient with clinical symptoms. In general, models of infection
86 of humans with SARS-CoV-2 assume an exponential increase in viral RNA synthesis followed by a
87 decrease, as antibody levels increase[15].

88

89 To investigate the presence of SARS-CoV-2 sgmRNAs in clinical (and other) samples, a
90 bioinformatics tool (LeTRS), was developed to analyse sequencing data from SARS-CoV-2
91 infections by identifying the unique leader-TRS gene junction site for each sgmRNA. The utility of
92 this tool was demonstrated on cultured cells infected with SARS-CoV-2, nasopharyngeal samples
93 from humans with COVID-19 and longitudinal analysis of nasopharyngeal samples from two non-
94 human primate models infected with SARS-CoV-2. The tool is adaptable for other coronaviruses.
95 The results have implications for virus biology, diagnostics and disease modelling.

96

97 **Results**

98 A tool, LeTRS (named after the leader-TRS fusion site), was developed to detect and quantify
99 defined leader gene junctions of SARS-CoV-2 (and other coronaviruses) from multiple types of
100 sequencing data. This was used to investigate SARS-CoV-2 sgmRNA synthesis in humans and non-
101 human primate animal models. LeTRS was developed using the Perl programming language,
102 including a main program for the identification of sgmRNAs and a script for plotting graphs of the
103 results. The tool accepts FASTQ files derived from Illumina paired-end or Oxford Nanopore
104 sequencing (amplicon or direct RNA), or BAM files produced by a splicing alignment method with
105 a SARS-CoV-2 genome (Supplementary Figure 1). Note that SARS-CoV-2 sgmRNAs are not formed
106 by splicing, but this is the apparent observation from sequencing data because of the
107 discontinuous nature of transcription. By default, LeTRS analyses SARS-CoV-2 sequence data by
108 using 10 known leader-TRS junctions and an NCBI reference genome (NC_045512.2) to identify
109 leader dependent canonical sgmRNAs. However, given the potential heterogeneity in the leader-
110 TRS region and potential novel (leader dependent noncanonical) sgmRNAs the user can also
111 provide customised leader-TRS junctions and SARS-CoV-2 variants as a reference. As there is
112 some heterogeneity in the leader-TRS sites, LeTRS was also designed to search for multiple
113 features of sgmRNAs. This included the leader-TRS junction in a given interval, report on the 20
114 nucleotides at the 3' end of the leader sequence, the TRS, translate the first predicted orf of the
115 sgmRNA, and find the conserved ACGAAC sequences in the TRS. LeTRS can also be used to identify
116 the sequencing reads with leader independent fusion sites that has been suggested to probably
117 produce unknown ORFs yielding functional products [16]. The tool was designed to investigate
118 very large data sets that are produced during sequencing of multiple samples.

119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140

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[17], MIDNIGHT[18] 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, 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 Illumina read

141 length (usually 100-250 nts), both the forward and reverse primers would be detected at the
142 ends of each paired read (Figure 1B pink lines). If the amplicon was longer than the Illumina read
143 length, primer sequence would be only found at one end of each paired read (Figure 1B green
144 and brown lines), with the possibility of one of the paired reads having a fusion site. The extension
145 stage could also proceed with a single primer using cDNA derived from the sgmRNA as a template.
146 This type of PCR product has a very low amplification efficiency, but theoretically could also
147 generate the same Illumina paired-end read with a single primer sequence at one end (Figure
148 1C). These paired-end reads could include the fulllength of the leader sequence but might not
149 reach the 3' end of the sgmRNA, because of the limitation of Illumina sequencing length and
150 extension time (Figure 1C). Also, unless there are cryptic TRSs located towards the 3' end of the
151 genome, all sgmRNAs would be expected to be larger than the Illumina sequencing length.

152

153 In contrast, the different types of read alignment in the Nanopore based amplicon are simpler to
154 assign. The longer reads that tend to be generated by Nanopore sequencing (depending on
155 optimisation) enable the capture of full-length sequences of all amplicons. Provided the leader
156 sequence is included as a forward primer most of the reads spanning the leader-TRS junction
157 would contain the forward and reverse primer sequences at both ends (Figure 1D pink lines). If
158 the extension time allowed, single primer PCR amplification could take the Nanopore amplicon
159 sequencing reads to both the 3' and 5' ends of the sgmRNAs, and these types of reads would only
160 have a primer sequence at one end (Figure 1D brown lines). In the Nanopore direct RNA
161 sequencing (dRNAseq) approach, the full-length sgmRNA could be sequenced and mapped
162 entirely on the leader and TRS-orf regions (Figure 1E).

163

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

165 In order to assess the ability of LeTRS to identify the leader-TRS junctions from sequencing
166 information, a total RNA sample was prepared at 72 hours post-infection (hpi) from hACE2-A549
167 cells infected with SARS-CoV-2 (a lineage B isolate). This RNA was sequenced using an amplicon-
168 based approach (ARTIC) with either Nanopore (ARTIC-Nanopore) or Illumina (ARTIC-Illumina), or
169 alternatively by a Nanopore dRNAseq approach[16]. The ARTIC-Nanopore (Figure 2A,
170 Supplementary Table 1) and ARTIC-Illumina (Figure 2B, Supplementary Table 2) sequencing data
171 were evaluated with LeTRS by setting the analysis to both primers pools. For dRNAseq (Figure 2C,
172 Supplementary Table 3), data was evaluated with LeTRS using the default setting. All the major
173 known leader-TRS gene junctions were identified by these sequencing methods. Analysis
174 demonstrated an expected pattern of abundance of the leader-TRS gene junctions with the
175 leader-TRS nucleoprotein gene junction being most abundant (Figure 2A, B and C; Supplementary
176 Tables 1, 2 and 3). Novel low abundance leader-TRS gene junctions were also identified (Figure
177 2A, B and C; Supplementary Tables 1, 2 and 3). These known and novel leader-TRS junctions were
178 also known as leader dependent canonical and noncanonical fusions, respectively [2]. LeTRS also
179 has a function to identify leader independent long-distance fusion (>5,000 nt) and local joining
180 yielding a deletion between proximal sites (20-5,000 nt distance) in the sequencing reads. The
181 leader independent fusions (coverage ≥ 2) are shown in Supplementary Tables 1, 2 and 3. Indel
182 sequencing errors are frequent (defined as less than 20 nucleotides), especially in Nanopore
183 sequencing data, and therefore it is difficult to find precise fusion (apparent splicing) sites in this
184 case [19]. However, some of the novel leader-TRS junctions (noncanonical fusions) and leader

185 independent fusions in the test sample were supported by all three sequencing methods
186 (Supplementary Figure 2) with similar fusion sites. Many local fusions/deletions within the orf3,
187 E, M, orf6, orf7a, orf7b, orf8 and N genes were identified (Supplementary Figure 2 G, H and I)
188 confirmed previous findings [2, 20], and indicates these are common events. Some of the novel
189 leader-TRS junctions (noncanonical fusions) and leader independent fusions may be the result of
190 sequencing or reverse transcription errors, especially those with low abundance (Supplementary
191 Tables 1, 2 and 3; Supplementary Figure 2). The ARTIC-Illumina approach identified fewer novel
192 leader-TRS junctions (noncanonical fusions) and leader independent fusions than the other two
193 sequencing methodologies, probably due to lower sequencing coverage (Supplementary Tables
194 1, 2 and 3).

195
196 For ARTIC approaches, LeTRS was designed to analyse reads in the primers pool 1, pool 2 or both
197 pools. Only the ARTIC pool 1 included a forward primer that is located within the leader region
198 (< 80 nts) of the SARS-CoV-2 genome. The leader-TRS regions of sgRNAs can be PCR-amplified
199 by both forward and reverse primers in ARTIC pool 1, but only reverse primers in ARTIC pool 2.
200 The read counts evaluated by LeTRS in both ARTIC-Nanopore and ARTIC-Illumina were compared
201 in the test data for pool 1 and 2, and found only very few reads/read pairs contained the correct
202 primers (Supplementary Table 4 and 5), suggesting the primers in ARTIC pool 2 generally do not
203 contribute to sequencing of leader-TRS regions.

204

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

206 Other tools have been developed to identify sgmRNAs from ARTIC-Illumina and ARTIC-Nanopore
207 sequencing data, such as Periscope (v0.1.0) [21], SARS-CoV-2-leader
208 (<https://github.com/hyeshik/sars-cov-2-transcriptome>) [16] and SuPER
209 (<https://github.com/ncbi/SuPER>) [22]. These tools were compared with LeTRS as shown in Table
210 1. LeTRS and Periscope used the FASTQ files as input, while SARS-CoV-2-leader and SuPER
211 required SAM files from a user generated alignment. Searching fusion site and sequences tag in
212 the sequencing reads are two major methods used. LeTRS and SuPER analysed the fusion/splicing
213 information in sequence reads achieved by an alignment program and also take account of the
214 conserved ACGAAC sequences in the TRS. Periscope and SARS-CoV-2-leader are based on
215 searching for a short tag sequence in the leader from sequencing reads. However, searching for
216 a short tag sequence in the leader with the high error rate associated with Nanopore data can be
217 challenging. LeTRS and Periscope use primer information to differentiate reads mapping to
218 amplicons to reads mapping from original virus genomes. Besides Periscope, output from
219 dRNAseq is supported by the other available tools. Illumina sequencing reads are usually short (<
220 250 bases), paired and sequenced from both ends. If both reads in a single pair contain a fusion
221 site this will be counted twice by the other three tools (Figure 1B green and pink). However, if
222 only one of the reads in the pair contains a fusion site it will be counted once (Figure 1B brown).
223 This leads to biased counting. LeTRS takes this into account by treating each read pair as a single
224 event. LeTRS also has a unique function to analyse reads in the primers pool 1, pool 2 or both
225 pools from ARTIC based sequencing (Table 1).
226

227 To compare the performance to LeTRS, these three tools were evaluated using the hACE2-A549
228 cell culture sample sequenced by ARTIC-Nanopore, ARTIC-Illumina and Nanopore dRNAseq.
229 Using the ARTIC-Nanopore sequencing data, all the tools reported a similar number of read
230 counts for the 10 known sgmRNAs (Supplementary Figure 3A). LeTRS showed fewer counts for
231 the ARTIC-Illumina than the other three tools because of considering read pairs (Supplementary
232 Figure 3B). Interestingly, Periscope also identified fewer nucleoprotein sgmRNAs with the ARTIC-
233 Illumina sequencing data (Supplementary Figure 3B). As of writing, Periscope does not yet
234 support Nanopore dRNAseq data, therefore LeTRS, SARS-CoV-2-leader and SuPER were
235 compared. LeTRS and SARS-CoV-2-leader generally identified more dRNAseq reads than SuPER,
236 especially for the nucleoprotein sgmRNA (Supplementary Figure 3C). Finally, the ratio of read
237 counts with the 10 known sgmRNA (S:orf3:E:M:orf6:orf7a:orf7b:orf8:N:orf10) were compared,
238 and the three tools showed almost an identical ratio when analysing data from the same
239 sequencing methods (Supplementary Figure 3D). ARTIC-Nanopore and Nanopore dRNAseq
240 resulted in a higher ratio of read counts with M and orf7a respectively (Supplementary Figure
241 3D). The read counts ratio of sgmRNAs mapping to spike was much lower with dRNAseq
242 approaches (Supplementary Figure 3D).

243

244 **Normalisation of read counts for sgmRNA**

245 Normalisation of read counts has been widely used for RNAseq in the comparison of gene
246 expression level across samples [23]. The normalisation is generally based on the ratio of reads
247 mapped on the gene to the total number of reads in that sample. These tools use this algorithm
248 for the normalisation of read counts in searching for sgmRNA [21, 24]. LeTRS also incorporated a

249 method to differentiate the total reads mapped (i) or whether the reads have forward primer
250 only (ii), reverse primer only (iii), both primers (iv) or at least one primer (v) present. This is
251 achieved by (i) the total number of reads mapped on the SARS-CoV-2 genome for the number of
252 reads of leader-TRS fusion site as the numerator; (ii) the total number of reads with forward
253 primers only for the number of reads of leader-TRS fusion site with forward primers only as the
254 numerator; (iii) the total number of reads with reverse primers only for the number of reads of
255 leader-TRS fusion site with reverse primers only as the numerator; (iv) the total number of reads
256 with both primers for the number of reads of leader-TRS fusion site with both as the numerator
257 and (v) the total number of reads with at least one primer on one side for the number of reads
258 of leader-TRS fusion site with at least one primer on as the numerator (notes in Supplementary
259 Tables 1, 2 and 3).

260

261 Because LeTRS considers the primers; pool 1, pool 2 or both pools, normalisation could be
262 observed in ARTIC pool 1 only to minimise the effect from ARTIC pool 2 since primers in ARTIC
263 pool 2 are almost not involved the sequencing of leader-TRS regions (as described above). For
264 the same RNA derived from the hACE2-A549 cell culture sample sequenced by ARTIC-Nanopore,
265 ARTIC-Illumina or Nanopore dRNAseq approaches, the normalised counts for the known
266 sgmRNAs were much smaller with the pool 1 of PCR based amplicon methods (ARTIC-Nanopore
267 and ARTIC-Illumina) than the Nanopore dRNAseq approach (Figure 3A and C for the reads with
268 at least one primer sequence; Supplementary Tables 3, 4 and 5). However, the normalised counts
269 with ARTIC-Nanopore and ARTIC-Illumina showed the same ratio of known sgmRNA as the
270 Nanopore dRNAseq approach, except for sgmRNAs mapping to S and orf7a (Figure 3B and D for

271 the reads with at least a primer sequence). PCR based approaches increases the value of the
272 denominator and reduced the normalised count, because a full length of sgmRNA was counted
273 once with the dRNAseq approach compared to many times with the amplicon approaches. ARTIC-
274 Illumina had fewer normalised counts than ARTIC-Nanopore probably due to the sequencing bias
275 of Illumina during PCR [25]. Thus, if the samples were sequenced with the same methodology
276 they were comparable. With a PCR based method a normalised count should be used to show
277 the relative difference between samples.

278

279 LeTRS identified many reads with only one primer (one-sided amplification) with the PCR based
280 amplicon methods (Supplementary Tables 4 and 5). The ratio of reads with either forward and/or
281 reverse primers were compared for each sgmRNA to the overall ratios of reads, with forward
282 primers only or reverse primers only, both primers in all mapped reads of pool 1 and pool 2 and
283 the mapped reads with any fusion sites of pool 1 and pool 2. This indicated that abundant reads
284 were identified with a single pattern and these were similar to reads mapping to sgmRNAs,
285 suggesting a one sided amplification is associated with amplicon-based approaches
286 (Supplementary Figure 4).

287

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

291 Part of the difficulty of studying SARS-CoV-2 and the disease COVID-19 is establishing the
292 sequence of events from the start of infection. Most samples from humans are from

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

301

302 As expected, analysis of the sequence data using LeTRS from the inoculum used to infect the
303 NHPs indicated that leader gene junctions could be identified, but these did not follow the
304 pattern of abundance of leader TRS-gene junctions found in infected cells in culture, where the
305 leader TRS nucleoprotein gene junction was most abundant (Supplementary Figure 5). The
306 inoculum would be expected to contain mostly genomic RNA found in virions. In contrast,
307 analysis of the longitudinal sequencing data from nasopharyngeal aspirates from the NHP model
308 using LeTRS identified leader TRS-gene junctions associated with the major sgRNAs (Figure 4,
309 Supplementary Table 7) as well as novel leader-TRS gene junction sites (Supplementary Figures
310 6 and 7). Analysing the abundance of the leader-TRS-gene junctions for both model species over
311 the course of infection revealed a phasic nature of sgRNA synthesis in pool 1 to minimise the
312 effect from ARTIC pool 2 (Figure 4). The leader-TRS nucleoprotein gene junction was the most
313 abundant, and there was a phasic pattern of potential sgRNA abundance identified with the

314 ARTIC-Illumina method (Figure 4). For both species, viral load and hence sgRNA abundance had
315 decreased by 8 and 9 dpi.

316

317 **Analysis of leader-TRS-gene junction in human samples revealed expected and aberrant**
318 **abundances of sgRNAs**

319 To investigate the pattern of leader-TRS-gene junction abundance during infection of SARS-CoV-
320 2 in humans, nasopharyngeal swabs from patients with COVID-19 were sequenced by ARTIC-
321 Illumina (using samples from COG-UK) (N=15 patients) (Figure 5, Supplementary Table 8) or by
322 ARTIC-Nanopore (using samples from ISARIC-4C) (N=15 patients) (Figure 6, Supplementary Tables
323 9 and 10). In several samples, leader-TRS-gene junctions were identified and followed an
324 expected pattern, with the nucleoprotein gene junction being the most abundant (e.g., Sample
325 1 in Figures 5A and B, Patient 2 day1 in Figure 6A and B). However, in several of the samples there
326 was very large representation of single leader-TRS-gene junction (e.g., Sample 4 and 5 in Figures
327 5A and B). These tended to map to the nucleoprotein gene (Sample 5, 8 and 13 Figures 5A and
328 B). The heterogeneity in abundance of leader-TRS-gene junctions was reminiscent of that from
329 the NHP study with a defined and expected pattern near the start of infection but then becoming
330 phasic. The samples gathered under ISARIC-4C were from hospitalised patients and permitted
331 analysis in relation to reported date of symptom onset and sequential sampling. In general, the
332 data indicated that the first sample on admission to hospital contained an abundance of leader-
333 TRS-gene junctions which resembled the pattern seen in infected cells (Patient 6 day 1 and day 9
334 in Figures 6A and B). However, with further days post-sample, e.g. (Patient 7 day 7 Figures 6A
335 and B), the leader-TRS nucleoprotein gene junction was the most abundant and far exceeded any

336 other detectable species. The abundance of leader-TRS nucleoprotein gene junction in the
337 patients at a later stage of infection followed that observed in the NHP model (Figure 4).

338

339 **Analysis of sequencing data from a previously published study investigating SARS-CoV-2 RNA**
340 **in samples from patients**

341 Recent research detected sgmRNAs mapping to E, ORF7a and N in swabs up to 14 days in one
342 patient and ORF7a and N in another patient up to 17 days after first detection by using a high-
343 throughput amplicon sequencing method known as Ion AmpliSeq Coronavirus Research Panel on
344 an Ion S5 XL genetic sequencer. The authors concluded these sgmRNAs may be present for a
345 significant time after active infection due to nuclease resistance and protection by cellular
346 membranes [24]. The sequencing data from this study was reanalysed using LeTRS, and
347 confirmed the finding of sgmRNAs in late infection from the two patients (Supplementary Table
348 11). Apart from nuclease resistance and protection by cellular membranes, a phasic pattern of
349 sgmRNA synthesis may also contribute to the presence of sgmRNAs at later time points.

350

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

352 N6-methyladenosine (6mA) is a widely observed modification on cellular RNA, and 5-
353 methylcytosine methylation (5mC) has also been reported on viral RNAs [16]. Methylation of
354 SARS_CoV-2 RNA was examined using sequencing data from the Nanopore direct RNA seq
355 approach. Total RNA was purified at 6, 12 and 24 hpi from cells infected with SARS-CoV-2. The
356 total RNA was sequenced and reads mapping to sgmRNAs were extracted with LeTRS for 6mA
357 and 5mC examination. Almost all 10 observed sgmRNAs showed the same number of

358 modification sites of 6mA and 5mC at 6, 12 and 24 hpi. Modification with 5mC was more
359 abundant than 6mA in all 10 known sgmRNAs. There were differences in abundance of some
360 sgmRNAs especially the M and N subgenomic mRNAs (Supplementary Figures 8 and 9). However,
361 there did not appear to be a relationship between number of methylation sites and the
362 abundance of a particular sgmRNA (Supplementary Figures 8 and 9).

363

364 To further evaluate the relationship between time post-infection and modification by
365 methylation, a paired samples one-sided Wilcoxon test was used. This analysis suggested that
366 the 5mC modification fraction at 24 hpi was significantly less than compared to modification at 6
367 and 12 hpi (p -value < 0.05), except for ORF7b and S (Supplementary Table 12). Modification with
368 6mA at 24 hpi was also significantly less than at 6 hpi, but not at 12 hpi (p -value < 0.05) in S,
369 ORF3a, E, M, ORF6, ORF7a, ORF8 and N. The abundance of most sgmRNAs decreased with time
370 and both of these factors could account for the frequency of methylation.

371

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

373 The sequencing data from cells infected in culture (Supplementary Table 13), animal models and
374 clinical samples from humans indicated the presence of novel leader-TRS gene junctions. Their
375 detection generally increased with depth of coverage. Coronavirus replication and transcription
376 is promiscuous, and recombination is a natural result of this, resulting in indels and potential
377 gene rearrangements. Many of these novel leader-TRS junctions were centred around the known
378 gene orf but out of the search interval. These types of leader-TRS-gene junctions could be only
379 found with spike, membrane, ORF6, ORF7b and nucleocapsid orfs, in which the membrane orf

380 was the most common (Figure 7A). To define what might be genuine novel leader-TRS-gene
381 junctions, these were compared across the data in all ARTIC-Illumina data (Figure 7B,
382 Supplementary Table 14). Five novel leader-TRS-gene junctions were identified that were
383 common to all the data, and the majority of these were present immediately 5' of the membrane
384 orf). The novel leader-TRS-gene junctions from LeTRS (Figure 7C) showed a similar distribution as
385 a previous study, although this study did not detail the precise location [28].

386

387

388 **Discussion**

389 Coronavirus sgmRNAs are only synthesised during infection of cells and therefore their presence
390 in sequence data can be indicative of active viral RNA synthesis. The abundance of the sgmRNAs
391 in infected cells should follow a general pattern where the sgmRNA encoding the nucleoprotein
392 is the most abundant. Identification and quantification of the unique leader-TRS-gene junctions
393 for each sgmRNA can be used as a proxy for their abundance.

394

395 LeTRS was developed to interrogate sequencing datasets to identify the leader-TRS-gene
396 junctions present at the 5' end of the sgmRNAs. LeTRS was first evaluated and validated on cell
397 culture data from published datasets[2, 17] and from a cell culture experiment as part of this
398 study and then used in an analysis of nasopharyngeal samples from NHP and human clinical
399 samples. The results showed that the positions of the leader-TRS junction sites with peak read
400 counts were the same as the given reference positions. The exception was at the leader-TRS-
401 gene junction for orf7b in the Nanopore sequencing. The normalised count results confirmed the
402 reads spanning the junctions showed that the leader-TRS nucleoprotein gene junction was the
403 most abundant, and orf7b and orf10 were the most infrequent in line with other data[2, 24].
404 Several low abundant leader-TRS junctions were identified in all of the datasets (Supplementary
405 Figure 2) with the implication these were either from potential lower abundant novel sgmRNAs
406 or represented known sgmRNAs, but with different leader-TRS junctions. Likewise, at low
407 frequency these could represent an aberrant viral transcription, perhaps as a mechanism to
408 generate new orfs for selection or these could be artefacts of the different sequencing processes
409 (Figure 2). Traditionally, such sgmRNAs have been first identified in coronaviruses by either

410 northern blot and/or metabolic labelling [8] and sequencing approaches are likely to be more
411 sensitive giving the amplification steps involved. Several other groups have identified novel
412 leader-TRS-gene junctions and potential sgmRNAs for other coronaviruses, including avian
413 infectious bronchitis virus[29]. The best way of validating potential novel sgmRNAs would be
414 through matching proteomic data to confirm genuine ORFs [1]. Analysis of several published
415 sequencing datasets identified novel viral RNA molecules that the authors suggested were
416 sgmRNAs containing only the 5' region of orf1a [30]. Such species are likely to be defective RNAs,
417 that act as templates for replication, rather than sgmRNAs. Interestingly, at later time points
418 post-infection in cell culture, potential novel sgmRNAs were found to be generated non-
419 specifically [30]. This potentially ties in with a disconnect of leader-TRS-gene junctions observed
420 in our study both *in vivo* from the nasopharyngeal samples from latter time points in the NHP
421 models and in humans. This is also shown in published data from SARS-CoV-2 infections in cell
422 culture gathered at later time points compared to earlier time points [2, 17].

423

424 Advanced filtering can improve the confidence of the identified leader-TRS junction from
425 sequencing data. Amplicon sequencing provided a unique opportunity to filter the sequencing
426 reads. The reads spanning the junctions with the correct forward primer, reverse primer or both
427 primer sequences at the ends of reads proved the known/novel sgmRNA existing in tested ARTIC-
428 Illumina and ARTIC-Nanopore amplicon sequencing data (Supplementary Tables 1 and 2). For
429 Illumina sequencing, the same junction on paired reads with at least one primer provided extra
430 evidence for leader-TRS identification. Some reads were identified that did not have primer
431 sequences and these were likely to be erroneously mapped, from template sgmRNA or low-

432 quality sequence. These were present at very low abundance compared to authentically mapped
433 reads (Supplementary Tables 1 and 2). The Nanopore dRNAseq approach had the potential to
434 generate full-length mRNA sequences. The polyA sequences and leader-TRS junctions in the
435 reads can be good signals to prove the full-length sgmRNA in the test data (Supplementary Table
436 3). Crucially, LeTRS is the only tool to consider paired-end Illumina data and primer pools
437 currently, and therefore is suited for the paired-end Illumina data and provided the amplicon
438 sequencing information from either primer pools.

439

440 In terms of clinical samples (typically nasopharyngeal swabs), the presence of sgmRNAs will
441 generally be due to the presence of infected cells. This has been seen as indicative of active viral
442 RNA synthesis at the time of sampling[5, 31, 32], although these have also been postulated to be
443 present through resistant structures after infection has finished[33]. Analysis of inoculum
444 indicated that leader-TRS-gene junctions could be identified (Supplementary Figure 5) but that
445 these were not in the same ratio as found in cells infected in culture (e.g., Figure 2A, B and 2C).
446 Thus, if the abundance of leader-TRS-gene junctions follows an expected pattern of the leader-
447 TRS nucleoprotein gene junction being the most abundant followed by a general gradient in
448 sequence data from nasopharyngeal samples, then this may be indicative of an active infection
449 – and the presence of infected cells in a sample.

450

451 In the absence of a human challenge model, NHP models that closely resemble COVID-19 disease
452 in humans can be used to study SARS-CoV-2 infection from a very defined initial exposure. RNA
453 was sequenced from longitudinal nasopharyngeal samples from two NHP models, rhesus and

454 cynomolgus macaques[26]. LeTRS was used to identify the abundance of the leader-TRS-gene
455 junctions in this data. The analysis indicated a phasic pattern of sgmRNA synthesis with a large
456 drop off after 8 or 9 dpi in both NHP models. This phasic pattern may be explained by an initial
457 synchronous infection of respiratory epithelial cells followed by cell death. Released virus then
458 goes on to infect new epithelial cells, with virus infection increasing exponentially in waves but
459 becoming asynchronous. The decline in sgmRNA from 8 or 9 dpi overlaps with IgG seroconversion
460 and humoral immunity in both species[26], and follows similar kinetics to serology profiles
461 measured in patients with COVID-19.

462

463 The identification of sgmRNAs in nasopharyngeal samples and their kinetics has implications for
464 nucleic acid-based diagnostics (many of which have three targets, one in the orf1a/b region and
465 two which are shared between the genome and sgmRNAs – the nucleoprotein and the spike
466 genes). The phasic nature of leader-TRS-gene junctions in the longitudinal samples, and by
467 implication sgmRNAs, and overt abundance of the leader-TRS nucleoprotein gene junction found
468 in many of the human samples, suggests that it may not be possible to precisely identify where
469 in infection an individual is based on the abundance of sgmRNAs. Likewise, assuming equivalency
470 between the targets, if the nucleoprotein target is found to be more abundant than the spike
471 target than the genomic target, then this would suggest infected cells are present in the sample.
472 Decreases in Ct values associated with emerging variants could equally be explained by sloughed
473 cells being present in a nasopharyngeal sample as well as by increases in the amount of
474 virions/viral load. Therefore, we would caution that a decrease in Ct associated with RT-qPCR
475 based assays may not just be reflective of higher viral loads but also may be indicative of more

476 infected cells being present. These possibilities may be resolved by considering the relative ratios
477 of sgmRNAs identified.

478 **METHODS**

479 **Data input**

480 LeTRS was designed to analyse FASTQ files derived from Illumina paired-end or Nanopore
481 sequencing data derived from a SARS-CoV-2 amplicon protocol, or standard Nanopore SARS-CoV-
482 2 dRNAseq data (Figure 1). The Illumina/Nanopore FASTQ sequencing data were cleaned to
483 remove adapters and low-quality reads before input. Sequencing data derived from other
484 sequencing modes or platforms can also be analysed by LeTRS via input of a BAM file produced
485 by a custom splicing alignment method with a SARS-CoV-2 genome (NC_045512.2) as a reference
486 (Figure 1). This can also be rapidly adapted for other coronaviruses.

487

488 **Library preparations and sequencing**

489 We sequenced the 15 samples from human patients with Nanopore. Total RNA was isolated using
490 a QIAamp Viral RNA Mini Kit (Qiagen, Manchester, UK) by spin-column procedure according to
491 the manufacturer's instructions. Clinical samples were extracted with Trizol LS as described[4].
492 All RNA samples were treated with Turbo DNase (Invitrogen). SuperScript IV (Invitrogen) was
493 used to generate single-strand cDNA using random primer mix (NEB, Hitchin, UK). ARTIC V3 PCR
494 amplicons from the single-strand cDNA were generated following the Nanopore Protocol of PCR
495 tiling of SARS-CoV-2 virus (Version: PTC_9096_v109_revL_06Feb2020). Amplicons generated by
496 ARTIC PCR were purified and normalised to 200 fmol before DNA end preparation and barcode
497 and adapter ligation. Library was loaded onto a FLO-MIN106 flow cell and sequencing reads were
498 called with Guppy using the high-accuracy calling parameters.

499

500 The NHP samples and their inoculum, and our laboratory experiments conducted in cells were
501 sequenced with Illumina. The amplicons products for Illumina sequencing were prepared as per
502 the Nanopore sequencing above and then used in Illumina NEBNext Ultra II DNA Library
503 preparation. Following 4 cycles of amplification the library was purified using Ampure XP beads
504 and quantified using Qubit and the size distribution assessed using the Fragment analyzer. Finally,
505 the ARTIC library was sequenced on the Illumina® NovaSeq 6000 platform (Illumina®, San Diego,
506 USA) following the standard workflow. The generated raw FastQ files (2 x 250 bp) were trimmed
507 to remove Illumina adapter sequences using Cutadapt v1.2.1 [34]. The option “-O 3” was set, so
508 the that 3’ end of any reads which matched the adapter sequence with greater than 3 bp was
509 trimmed off. The reads were further trimmed to remove low quality bases, using Sickle v1.200
510 [35] with a minimum window quality score of 20. After trimming, reads shorter than 10 bp were
511 removed.

512

513 The LeTRS was also tested with a combined Nanopore-ARTIC v3 amplicon dataset of 7 published
514 viral cell culture samples (barcode01-barcode07) [17], and a dataset from a published direct RNA
515 Nanopore sequencing analysis Vero cells infected with SARS-CoV-2 or an uninfected negative
516 control [2].

517

518 **Sequencing data alignment and basic filtering**

519 LeTRS controlled Hisat2 v2.1.0 [36] to map the paired-end Illumina reads against the SARS-CoV-
520 2 reference genome (NC_045512.2) with the default setting, and Minimap2 v2.1 [19] to align the
521 Nanopore cDNA reads and direct RNA-seq reads on the viral genome using Minimap2 with “-ax

522 splice” and “-ax splice -uf -k14” parameters, respectively. LeTRS provided 10 known leader-TRS
523 junctions to improve alignment accuracy by using “--known-splicesite-infile” function in Hisat2
524 and “--junc-bed” function in Minimap2, but this application could be optionally switched off by
525 users. In order to remove low mapping quality and mis-mapped reads before searching the
526 leader-TRS junction sites, LeTRS used Samtools v1.9 [37] to have basic filtering for the reads in
527 the output Sam/Bam files according to their alignment states as shown (Table 9 - basic filtering).

528

529 **Searching the leader-TRS motifs**

530 After the mapping and basic filtering step, LeTRS searched aligned reads spanning the leader-TRS
531 junctions in the SARS-CoV-2 reference genome (Supplementary Figure 1). For the known leader-
532 TRS junctions, LeTRS searched the reads including the leader-TRS junctions within a given interval
533 around the known leader and TRS junctions sites. The leader break site interval is ± 10 nts, and
534 the TRS breaking sites interval is -20 nts to the 1 nt before the first known AUG in the default
535 setting (the intervals can be changed to custom values to investigate heterogeneity). LeTRS then
536 reported a peak count that was the number of reads carrying the most common leader-TRS
537 junctions within the given leader and TRS breaking sites intervals, and a cluster count that was
538 the number of all reads carrying leader-TRS junctions within the given leader and TRS breaking
539 sites intervals (Tables 1-6). LeTRS also searched the junctions out of the given intervals (the
540 genomic position of leader breaking site < 80) and reported the number of reads (>10 by default)
541 with novel leader-TRS junctions. These number of read counts were also reported by number of
542 reads in 1000000 as normalisation. The read including the known and novel leader-TRS junctions
543 could be optionally outputted in FastA format. Based on identified known and novel leader-TRS

544 junctions, LeTRS could report 20 nucleotides towards the 3' end of the leader sequence, the TRS
545 and translated the first orf of sgmRNAs sequence, and find the conserved ACGAAC sequences in
546 the TRS (Table S1-S6).

547

548 **Advance filtering**

549 Based on the alignment possibilities illustrated in Figure 2 and discussed, LeTRS further filters the
550 identified reads with known and novel leader-TRS junctions. This step is named as advance
551 filtering and can only applied when the input data is from Illumina paired-end reads, Nanopore
552 cDNA reads or Nanopore RNA reads (Table 2). If a BAM file is used as input data, the advanced
553 filtering step would be automatically skipped (Table 2). The number of reads including the known
554 and novel leader-TRS junctions, and the number of reads filtered with corresponding advance
555 filtering criteria were outputted into two tables in tab format (Tables 1-6).

556

557 **Leader-TRS junction plotting**

558 LeTRS-plot was developed as an automatic plotting tool that interfaces with the R package
559 ggplot2 v3.3.3 to view the leader-TRS junctions in the tables generated by LeTRS (Figure 3-5). The
560 plot shows peak count, filtered peak count, normalized peak count and normalized filtered peak
561 count for known leader-TRS junctions, and novel junction counts, filtered novel junction count,
562 normalized novel junction count and filtered normalized novel junction for novel leader-TRS
563 junctions.

564

565 **RNA modifications**

566 Total RNA extracted from cultured cells at 6, 12 and 24 hours were collected for Oxford
567 Nanopore direct RNA sequence. LeTRS was then run with a parameter of “extractfasta” to extract
568 subgenomic mRNAs reads in sequenced samples. The fast5 files that corresponds to the
569 extracted subgenomic mRNAs reads were withdrawn using fast5_subset in Oxford Nanopore
570 ont_fast5_api package (v0.3.2, https://github.com/nanoporetech/ont_fast5_api). The re-
571 squiggle algorithm in Tombo analysis pipelines (v1.5.1, <https://github.com/nanoporetech/tombo>)
572 defines a new assignment from raw signals to reference sequence with “--num-most-common-
573 errors 5” option. The resquiggled raw signals were further processed using “detect_modifications
574 alternative_model” functions in Tombo by setting “--rna and --alternate-bases 5mC” to identify
575 5-methylcytosine (5mC), and “predict_sites” in Nanom6A package (v2021_10_22) [38] with
576 default setting to identify N6-methyladenosine (6mA) in the subgenomic mRNAs reads.

577

578 **References**

- 579 1. Davidson, A.D., et al., *Characterisation of the transcriptome and proteome of SARS-CoV-2*
580 *reveals a cell passage induced in-frame deletion of the furin-like cleavage site from the*
581 *spike glycoprotein*. *Genome Med*, 2020. **12**(1): p. 68.
- 582 2. Kim, D., et al., *The Architecture of SARS-CoV-2 Transcriptome*. *Cell*, 2020. **181**(4): p. 914-
583 921 e10.
- 584 3. Nasir, J.A., et al., *A Comparison of Whole Genome Sequencing of SARS-CoV-2 Using*
585 *Amplicon-Based Sequencing, Random Hexamers, and Bait Capture*. *Viruses*, 2020. **12**(8).
- 586 4. Moore, S.C., et al., *Amplicon-Based Detection and Sequencing of SARS-CoV-2 in*
587 *Nasopharyngeal Swabs from Patients With COVID-19 and Identification of Deletions in the*

- 588 *Viral Genome That Encode Proteins Involved in Interferon Antagonism. Viruses, 2020.*
589 **12**(10).
- 590 5. Dorward, D.A., et al., *Tissue-Specific Immunopathology in Fatal COVID-19. Am J Respir Crit*
591 *Care Med, 2021. 203*(2): p. 192-201.
- 592 6. Graham, R.L., et al., *SARS coronavirus replicase proteins in pathogenesis. Virus Res, 2008.*
593 **133**(1): p. 88-100.
- 594 7. Pirc, K., et al., *Genome structure and transcriptional regulation of human coronavirus*
595 *NL63. Virol J, 2004. 1*: p. 7.
- 596 8. Hiscox, J.A., D. Cavanagh, and P. Britton, *Quantification of individual subgenomic mRNA*
597 *species during replication of the coronavirus transmissible gastroenteritis virus. Virus Res,*
598 *1995. 36*(2-3): p. 119-30.
- 599 9. Hiscox, J.A., et al., *Investigation of the control of coronavirus subgenomic mRNA*
600 *transcription by using T7-generated negative-sense RNA transcripts. J Virol, 1995. 69*(10):
601 p. 6219-27.
- 602 10. van Marle, G., et al., *Regulation of coronavirus mRNA transcription. J Virol, 1995. 69*(12):
603 p. 7851-6.
- 604 11. La Monica, N., K. Yokomori, and M.M. Lai, *Coronavirus mRNA synthesis: identification of*
605 *novel transcription initiation signals which are differentially regulated by different leader*
606 *sequences. Virology, 1992. 188*(1): p. 402-7.
- 607 12. Alonso, S., et al., *Transcription regulatory sequences and mRNA expression levels in the*
608 *coronavirus transmissible gastroenteritis virus. J Virol, 2002. 76*(3): p. 1293-308.

- 609 13. Sawicki, S.G., D.L. Sawicki, and S.G. Siddell, *A contemporary view of coronavirus*
610 *transcription*. J Virol, 2007. **81**(1): p. 20-9.
- 611 14. Jeong, Y.S. and S. Makino, *Evidence for coronavirus discontinuous transcription*. J Virol,
612 1994. **68**(4): p. 2615-23.
- 613 15. Cevik, M., et al., *Virology, transmission, and pathogenesis of SARS-CoV-2*. BMJ, 2020. **371**:
614 p. m3862.
- 615 16. !!! INVALID CITATION !!! [2].
- 616 17. Tyson, J.R., et al., *Improvements to the ARTIC multiplex PCR method for SARS-CoV-2*
617 *genome sequencing using nanopore*. bioRxiv, 2020.
- 618 18. Freed, N.E., et al., *Rapid and inexpensive whole-genome sequencing of SARS-CoV-2 using*
619 *1200 bp tiled amplicons and Oxford Nanopore Rapid Barcoding*. Biology Methods and
620 Protocols, 2020. **5**(1): p. bpaa014.
- 621 19. Li, H., *Minimap2: pairwise alignment for nucleotide sequences*. Bioinformatics, 2018.
622 **34**(18): p. 3094-3100.
- 623 20. Young, B.E., et al., *Effects of a major deletion in the SARS-CoV-2 genome on the severity*
624 *of infection and the inflammatory response: an observational cohort study*. The Lancet,
625 2020. **396**(10251): p. 603-611.
- 626 21. Parker, M.D., et al., *Subgenomic RNA identification in SARS-CoV-2 genomic sequencing*
627 *data*. Genome research, 2021. **31**(4): p. 645-658.
- 628 22. Yang, Y., et al., *Characterizing transcriptional regulatory sequences in coronaviruses and*
629 *their role in recombination*. Molecular Biology and Evolution, 2021. **38**(4): p. 1241-1248.

- 630 23. Anders, S., et al., *Count-based differential expression analysis of RNA sequencing data*
631 *using R and Bioconductor*. Nature protocols, 2013. **8**(9): p. 1765-1786.
- 632 24. Alexandersen, S., A. Chamings, and T.R. Bhatta, *SARS-CoV-2 genomic and subgenomic*
633 *RNAs in diagnostic samples are not an indicator of active replication*. Nature
634 communications, 2020. **11**(1): p. 1-13.
- 635 25. Ross, M.G., et al., *Characterizing and measuring bias in sequence data*. Genome biology,
636 2013. **14**(5): p. 1-20.
- 637 26. Salguero, F.J., et al., *Comparison of rhesus and cynomolgus macaques as an infection*
638 *model for COVID-19*. Nat Commun, 2021. **12**(1): p. 1260.
- 639 27. Ryan, K.A., et al., *Dose-dependent response to infection with SARS-CoV-2 in the ferret*
640 *model and evidence of protective immunity*. Nat Commun, 2021. **12**(1): p. 81.
- 641 28. Taiaroa, G., et al., *Direct RNA sequencing and early evolution of SARS-CoV-2*. BioRxiv, 2020.
- 642 29. Keep, S., et al., *Multiple novel non-canonically transcribed sub-genomic mRNAs produced*
643 *by avian coronavirus infectious bronchitis virus*. J Gen Virol, 2020. **101**(10): p. 1103-1118.
- 644 30. Nomburg, J., M. Meyerson, and J.A. DeCaprio, *Pervasive generation of non-canonical*
645 *subgenomic RNAs by SARS-CoV-2*. Genome Med, 2020. **12**(1): p. 108.
- 646 31. Corbett, K.S., et al., *Evaluation of the mRNA-1273 Vaccine against SARS-CoV-2 in*
647 *Nonhuman Primates*. N Engl J Med, 2020. **383**(16): p. 1544-1555.
- 648 32. Yu, J., et al., *DNA vaccine protection against SARS-CoV-2 in rhesus macaques*. Science,
649 2020. **369**(6505): p. 806-811.

- 650 33. Alexandersen, S., A. Chamings, and T.R. Bhatta, *SARS-CoV-2 genomic and subgenomic*
651 *RNAs in diagnostic samples are not an indicator of active replication*. Nat Commun, 2020.
652 **11**(1): p. 6059.
- 653 34. Martin, M., *Cutadapt removes adapter sequences from high-throughput sequencing reads*.
654 EMBnet.journal, 2011. **17**: p. <https://doi.org/10.14806/ej.17.1.200>.
- 655 35. Joshi, N.A. and J.N. Fass, *Sickle: A sliding-window, adaptive, quality-based trimming tool*
656 *for FastQ files*
657 *(Version 1.33)*. 2011: p. <https://github.com/najoshi/sickle>.
- 658 36. Kim, D., B. Langmead, and S.L. Salzberg, *HISAT: a fast spliced aligner with low memory*
659 *requirements*. Nat Methods, 2015. **12**(4): p. 357-60.
- 660 37. Li, H., et al., *The Sequence Alignment/Map format and SAMtools*. Bioinformatics, 2009.
661 **25**(16): p. 2078-9.
- 662 38. Gao, Y., et al., *Quantitative profiling of N 6-methyladenosine at single-base resolution in*
663 *stem-differentiating xylem of Populus trichocarpa using Nanopore direct RNA sequencing*.
664 Genome Biology, 2021. **22**(1): p. 1-17.
- 665
- 666

667 **Ethics approval and consent to participate**

668 All experimental work on NHPs was conducted under the authority of a UK Home Office approved
669 project license (PDC57C033) that had been subject to local ethical review at PHE Porton Down by
670 the Animal Welfare and Ethical Review Body (AWERB) and approved as required by the Home
671 Office Animals (Scientific Procedures) Act 1986 and the full ethics and NHP model are described.

672 **Consent for publication**

673 Not applicable

674 **Availability of data and materials**

675 LeTRS is available at <https://github.com/xiaofengdong83/LeTRS>.

676 Illumina and Nanopore test data sets are available under NCBI PRJNA699398.

677

678 **Competing interests**

679 The authors declare that they have no competing interests

680 **Funding**

681 This work was predominately funded by U.S. Food and Drug Administration Medical
682 Countermeasures Initiative contract (75F40120C00085) awarded to JAH. The article reflects the
683 views of the authors and does not represent the views or policies of the FDA. This work was also
684 supported by the MRC (MR/W005611/1) G2P-UK: A national virology consortium to address
685 phenotypic consequences of SARS-CoV-2 genomic variation (co-I JAH). JAH is also funded by the
686 Centre of Excellence in Infectious Diseases Research (CEIDR) and the Alder Hey Charity. The non-
687 human primate work was funded by the Coalition of Epidemic Preparedness Innovations (CEPI)
688 and the Medical Research Council Project CV220-060, Development of an NHP model of infection

689 and ADE with COVID-19 (SARS-CoV-2) both awarded to MWC. The ISARIC4C sample collection
690 and sequencing in this study was supported by a grants from the Medical Research Council (grant
691 MC_PC_19059), the National Institute for Health Research (NIHR; award CO-CIN-01), the Medical
692 Research Council (MRC; grant MC_PC_19059), and by the NIHR Health Protection Research Unit
693 (HPRU) in Emerging and Zoonotic Infections at University of Liverpool in partnership with Public
694 Health England (PHE), in collaboration with Liverpool School of Tropical Medicine and the
695 University of Oxford (award 200907), NIHR HPRU in Respiratory Infections at Imperial College
696 London with PHE (award 200927), Wellcome Trust and Department for International
697 Development (DID; 215091/Z/18/Z), the Bill and Melinda Gates Foundation (OPP1209135),
698 Liverpool Experimental Cancer Medicine Centre (grant reference C18616/A25153), NIHR
699 Biomedical Research Centre at Imperial College London (IS-BRC-1215-20013), PJMO is supported
700 by a NIHR senior investigator award (201385). The views expressed are those of the authors and
701 not necessarily those of the Department of Health and Social Care, DID, NIHR, MRC, Wellcome
702 Trust, or PHE. The funders had no role in the study design; in the collection, analysis, and
703 interpretation of data; in the writing of the report; or in the decision to submit the article for
704 publication.

705

706 **Authors' contributions**

707 X.D. developed the LeTRS software and performed the informatics analysis. X.D., A.D. and J.A.H.
708 analysed the data. J.S., J.T. and M.W.C. co-ordinated the NHP work and sample processing. R.P.-
709 R., J.P.S., H.G., T.P. and N.R. were involved in sequencing and informatics analysis of the NHP
710 samples with D.A.M. A.D. oversaw sequencing of the human clinical samples with E.V. and C.N

711 for the COG-UK data. R.P.-R. and J.A.H. oversaw sequencing of samples under the auspices of
712 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.
713 and M.W.C. initiated and led the study and wrote the manuscript with X.D., R.P.-R., A.D. with
714 other authors involved in editing the final version.

715 **Acknowledgments**

716 We would like to thank all members of the Hiscox Laboratory and the Centre for Genome
717 Research for supporting SARS-CoV-2/COVID-19 sequencing research. We would like to
718 acknowledge members of the COG-UK and ISARIC4C consortia for acquisition of the human
719 samples used in this study.

720

721

722

723

724

725

726

727

728

729

730

731

732

733 Table 1. Comparison of other Tools with LeTRS.

	LeTRS	Periscope	SARS-CoV-2-leader	SuPER
Input files	fastq	fastq	bam/sam	sam
Consider amplicon primer information used	yes	yes	no	no
Consider paired-end Illumina data	yes	no	no	no
Consider amplicon primer pool	yes	no	no	no
Consider the ACGAAC box	yes	no	no	yes
Support amplicon Illumina data	yes	yes	yes	yes
Support amplicon Nanopore data	yes	yes	yes	yes
Support Nanopore dRNAseq data	yes	no	yes	yes
Method	function searching	sequences tag searching	sequences tag searching	function searching

734

735

736

737

738

739

740

741

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

Output Filters	Illumina paired- end amplicon reads	Nanopore amplicon reads	Nanopore dRNAseq reads	Bam
MAPQ > 10	•	•	•	•
Read only one splicing junction	•	•	•	•
Basic Primary alignment only	•	•	•	•
filtering No supplementary alignment	•	•	•	•
Read mapped in pair	•			
No read reverse strand			•	
Read alignment 5' end includes forward primer	•	•		
Read alignment 3' end includes reverse primer	•	•		
Read alignment 5' end includes Advance forward primer and 3' end includes filtering reverse primer	•	•		
Paired read including at least one primer in each have same leader- TRS junction in alignments	•	•		
Read alignment 3' with > 1 ployA		•	•	
Read alignment 3' with > 5 ployA		•	•	

743

744 Figures

745 Figure 1. (A). Illustration of reads derived from sgmRNAs mapped onto the SARS-CoV-2 reference
746 genome with a splicing method. We note that splicing does not occur in coronaviruses but this is
747 the apparent observation of a fusion event between different parts of the genome. (B and C).
748 Illustration of the possible type of reads mapped on the SARS-CoV-2 reference genome for the
749 paired-end Illumina amplicon sequencing, where the lines with same colour implied paired reads,
750 (D) Nanopore amplicon sequencing and (E) Nanopore dRNAseq of the SARS-CoV-2 genome and
751 sgmRNAs. L and B in the boxes indicate the leader-TRS breaking sites on the leader side and TRS
752 side, respectively. Although we note these are where the apparent fusion site occurs. Yellow
753 colour indicates the leader region, black is the TRS and gene sequence, the red indicates a
754 sequence read that maps to SARS-CoV-2 sequence. Blue is a sequence that is present between
755 the leader sequence and the TRS. For (B) and (C) the same colour (brown, green and pink)
756 indicates that same paired read. For (B) the paired read contains both primers. For (C) the grey
757 and light blue colour is a paired read, but only contains one primer sequence at any end. The
758 vertical hash lines on (B, C, and D) indicates the position of a primer.

759

760 Figure 2. Analysis of reads mapping to the leader TRS-gene junctions with at least one primer
761 sequence at either end in sequencing data from hACE2-A549 cells infected with SARS-CoV-2 and
762 sequenced using either (A) an ARTIC-Nanopore approach, (B) an ARTIC-Illumina approach and (C)
763 a Nanopore dRNAseq approach. The data corresponds to that shown in detailed in
764 Supplementary Tables 1, 2 and 3. The standard deviation of a binomial distribution was calculated
765 to generate error bars. The data is presented as a histogram with a normalised count for each

766 sgmRNA starting at a particular position in the leader sequence as indicated in the line diagram
767 underneath. For each panel (A, B and C) the expected sgmRNA pattern is shown on the left and
768 novel sgmRNAs are shown on the right.

769

770 Figure 3. An X-Y/scatter plot using normalized counts of sgmRNAs (with greater than 5 A residues
771 at the 3' end – indicative of a polyA tail for the dRNAseq data). To generate the scatter plots
772 Nanopore dRNAseq data was plotted against the either the normalized count (at least one primer
773 sequence) of sgmRNAs with (A) ARTIC-Nanopore sequencing data and (C) ARTIC-Illumina
774 sequencing data or provided as ratio (B) and (D), respectively for
775 S:orf3:E:M:orf6:orf7a:orf7b:orf8:N:orf10 (using data from Supplementary Tables 3, 4 and 5).

776

777 Figure 4. Analysis of the abundance of reads mapping to the leader TRS-gene junctions that have
778 at least one primer sequence at either end in longitudinal nasopharyngeal samples taken from
779 two non-human primate models infected with SARS-CoV-2. The time post-infection in days is
780 indicated on the x-axis. The normalised count (read count/total number of reads mapped on the
781 reference genome)*1,000,000) of the leader TRS-gene junction abundance is shown on the left-
782 hand Y-axis with each unique leader TRS-gene junction colour coded. The right-hand Y axis is a
783 measure of the total depth of coverage for SARS-CoV-2 in that sample. Note the two scales are
784 different. SARS-CoV-2 was amplified and sequenced by ARTIC-Illumina. The data is organised into
785 groups of animals for the cynomolgus macaque groups 1 and 2 (A/E and B/F), and rhesus
786 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

787 for Day1 to Day9. The data corresponds to that shown in Supplementary Table 7. Standard
788 deviation of a binomial distribution was calculated to provide error bars.

789
790 Figure 5. Plots of normalised peak counts (A) and peak counts (B) of leader-TRS gene junctions of
791 reads with at least one primer sequences at either end derived from sequence data from 15
792 human patients sequenced with the ARTIC-Illumina approach and analysed by using sequence
793 derived from pool 1 primers. The data correspond to that shown in Supplementary Table 8.
794 Standard deviation of a binomial distribution was calculated to provide error bars.

795
796 Figure 6. Plots of normalised peak counts (A) and peak counts (B) of leader-TRS gene junctions of
797 reads with at least one primer sequence at either end derived from sequence data from 15
798 human patients sequenced with the ARTIC-Nanopore approach and analysed by using sequence
799 derived from pool 1 primers. The data correspond to that shown in Supplementary Table 9.
800 Standard deviation of a binomial distribution was calculated to provide error bars.

801
802 Figure 7. (A). Diagram of novel leader-TRS junctions centred around the known gene orf but out
803 of the search interval in the analysis of SARS-CoV-2 RNA from cell culture, non-human primate
804 and human sequencing data. Many novel junctions map to the leader-TRS membrane gene
805 junctions. (B). Venn diagram showing the overlap of novel leader-TRS gene junctions present in
806 SARS-CoV-2 infected cynomolgus and rhesus macaques, human patients, and Vero cells. Data
807 was obtained using the ATRIC-Illumina method (Supplementary Table 14). (C) Virus genome
808 position of the start of the fusion site (Y-axis) in the leader sequence plotted against the fusion

809 site present in the gene to show the potential positions of the novel leader-TRS junctions along
810 the SARS-CoV-2 genome (indicated above). A shown the colours present the novel leader-TRS
811 junctions identified in the different experimental condition (cynomolgus and rhesus macaques,
812 human patients, and Vero cells).

813 Supplementary Figures

814 Supplementary Figure 1. Bioinformatics pipeline for the identification of leader-TRS junctions in
815 sequencing data from SARS-CoV-2 infected material with LeTRS. This can be rapidly adapted for
816 other coronaviruses such as MERS-CoV and any newly emerged coronavirus. LeTRS can work
817 from Nanopore or Illumina amplicon data or more unbiased approaches such as direct RNA
818 sequencing, metagenomic or Illumina sequencing by using a BAM file.

819

820 Supplementary Figure 2. Novel (leader dependent noncanonical) fusions (count ≥ 2) found in the
821 cell culture test sample sequenced by (A) ARTIC-Nanopore, (B) ARTIC-Illumina and (C) Nanopore
822 dRNAseq approaches; leader independent long-distance ($>5,000$ nt) fusions (count ≥ 2) found in
823 the cell culture test sample sequenced by (D) ARTIC-Nanopore, (E) ARTIC-Illumina and (F)
824 Nanopore dRNAseq approaches; leader independent local joining yielding a deletion between
825 proximal sites (20–5,000 nt distance) fusions (count ≥ 2) found in the cell culture test sample
826 sequenced by (G) ARTIC-Nanopore, (H) ARTIC-Illumina and (I) Nanopore dRNAseq approaches.
827 The data correspond to that shown Supplementary Tables 1, 2 and 3.

828

829 Supplementary Figure 3. Comparison of different tools and LeTRS to evaluate sequencing data to
830 identify the unique sequencing features of SARS-CoV-2 sgmRNAs. Number of reads were
831 evaluated by LeTRS (all peak count), SARS-COV-2-leader, SuPER or periscope (High Quality count)
832 with the cell culture test sample sequenced by (A) ARTIC-Nanopore, (B) ARTIC-Illumina and (C)
833 Nanopore dRNAseq approaches; (D) Ratio of sgmRNAs (S:orf3:E:M:orf6:orf7a:orf7b:orf8:N:orf10)
834 identified by LeTRS (all peak count), SARS-COV-2-leader, SuPER or periscope (HQ count) with the

835 cell culture test sample sequenced by ARTIC-Nanopore, ARTIC-Illumina and Nanopore dRNAseq
836 approaches. The data are corresponded to that shown in Supplementary Tables 1, 2 and 3.

837

838 Supplementary Figure 4. Comparison of the ratio of reads in amplicon sequencing approaches
839 based on the ARTIC approach, with the forward primer only, reads with reverse primer only and
840 reads with both primers in sgRNAs to the overall ratio of reads with the forward primer only,
841 reads with reverse primer only and reads with both primers in all reads amplified by pool 1
842 primers, pool 2 primers and both pools of primers for the cell culture test sample sequenced by
843 (A) ARTIC-Nanopore and (B) ARTIC-Illumina approaches.

844

845 Supplementary Figure 5. Raw (A and C) and normalised (B and D) canonical (upper) and novel
846 (lower) leader-TRS gene junctions count in RNA purified from the inoculum of SARS-CoV-2 used
847 to infect either the cynomolgus or rhesus macaques. The RNA was sequenced by the ARTIC-
848 Illumina method (Supplementary Table 6). Standard deviation of a binomial distribution was
849 calculated to provide error bars.

850

851 Supplementary Figure 6. Novel leader-TRS gene junctions (count > 10) identified in RNA purified
852 from nasopharyngeal swabs taken daily from cynomolgus macaques infected with SARS-CoV-2
853 (Supplementary Table 7). The number before “-Day” indicated the group of cynomolgus
854 macaques. Standard deviation of a binomial distribution was calculated to provide error bars.

855

856

857 Supplementary Figure 7. Novel leader-TRS gene junctions (count > 10) identified in RNA purified
858 from nasopharyngeal swabs taken daily from from rhesus macaques (Supplementary Table 7).
859 The number before “-Day” indicated the group of cynomolgus macaques. Standard deviation of
860 a binomial distribution was calculated to provide error bars.

861

862 Supplementary Figure 8. Comparison of the fraction of 6mA modification (right-hand Y-axis) of
863 each site in sgmRNA at 6, 12 and 24 hours after post infection using direct RNA sequencing from
864 RNA purified from SARS-CoV-2 infected cells. The normalised count of the leader TRS-gene
865 junction abundance is shown on the left-hand Y-axis.

866

867 Supplementary Figure 9. Comparison of the fraction of 5mC modification (right-hand Y-axis) of
868 each site in sgmRNA at 6, 12 and 24 hours after post infection using direct RNA sequencing from
869 RNA purified from SARS-CoV-2 infected cells. The normalised count of the leader TRS-gene
870 junction abundance is shown on the left-hand Y-axis.

871

872

873

874

875

876

877

878

879 Supplementary Tables

880 Table S1. The LeTRS output tables for known sgmRNA, details of known sgmRNA, novel sgmRNA
881 (count ≥ 2), details of novel sgmRNA, and leader independent long-distance and local fusions
882 (count ≥ 2) evaluated in the cell culture test sample sequenced by the ARTIC-Nanopore approach.

883

884 Table S2. The LeTRS output tables for known sgmRNA, details of known sgmRNA, novel sgmRNA
885 (count ≥ 2), details of novel sgmRNA, and leader independent long-distance and local fusions
886 (count ≥ 2) evaluated in the cell culture test sample sequenced by the ARTIC-Illumina approach.

887

888 Table S3. The LeTRS output tables for known sgmRNA, details of known sgmRNA, novel sgmRNA
889 (count ≥ 2), details of novel sgmRNA, and leader independent long-distance and local fusions
890 (count ≥ 2) evaluated in the cell culture test sample sequenced by the Nanopore dRNAseq
891 approach.

892

893 Table S4. The LeTRS output table for known sgmRNA evaluated by primers of pool 1 and 2 in the
894 cell culture test sample sequenced by the ARTIC-Nanopore approach.

895

896 Table S5. The LeTRS output tables for known sgmRNA evaluated by primers of pool 1 and 2 in the
897 cell culture test sample sequenced by the ARTIC-Illumina approach.

898

899 Table S6. The LeTRS output tables for known sgmRNA and details of known sgmRNA with pool 1
900 primers, and novel sgmRNA (count > 10) and details of novel sgmRNA with both pools' primers

901 in the infecting SARS-CoV-2 inoculum source used for the NHP study, sequenced by the ARTIC-
902 Illumina method.

903

904 Table S7. The LeTRS output tables for known sgmRNA and details of known sgmRNA with pool 1
905 primers, and novel sgmRNA (count > 10) and details of novel sgmRNA with both pools' primers
906 in longitudinal nasopharyngeal samples taken from two non-human primate models (cynomolgus
907 and rhesus macaques) of SARS-CoV-2 in groups. SARS-CoV-2 was amplified using the ARTIC
908 approach and sequenced by Illumina. The data is organised into groups of animals for the
909 cynomolgus macaque groups 1 and 2 that were with "-1" and "-2" in the excel sheets.

910

911 Table S8. The LeTRS output tables for known sgmRNA and details of known sgmRNA in pool 1,
912 and novel sgmRNA (count > 10) and details of novel sgmRNA with both pools' primers from 15
913 human patients sequenced with ARTIC-Illumina.

914

915 Table S9. The LeTRS output tables for known sgmRNA and details of known sgmRNA in pool 1
916 from 15 human patients sequenced with ARTIC-Nanopore.

917

918 Table S10. The spreadsheet for the 15 human patients sequenced with the ARTIC-Nanopore
919 detailed in Table S9.

920

921 Table S11. Re-analysis of reads for known sgmRNAs in the (NCBI accession No. PRJNA636225)
922 [24].

923

924 Table S12. Evaluation of the difference of modification by the paired samples one-sided Wilcoxon
925 test to calculate p-value by treating the same nucleotides between any two time points as paired
926 data.

927

928 Table S13. The LeTRS output table for novel sgmRNA (count > 10) and details of novel sgmRNA
929 with both primer pools from VeroE6 cells infected in culture with SARS-CoV-2 (SCV2-006)
930 sequenced by ARTIC-Illumina primers. This sample is different from the one Table S2.

931

932 Table S14. Novel leader-TRS junctions centred around the known gene open reading frame but
933 out of the search interval in the analysis of cell culture, non-human primate and human
934 sequencing data.

935

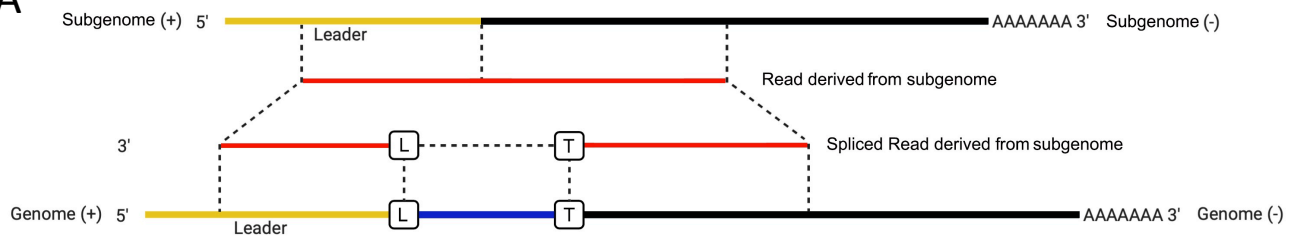
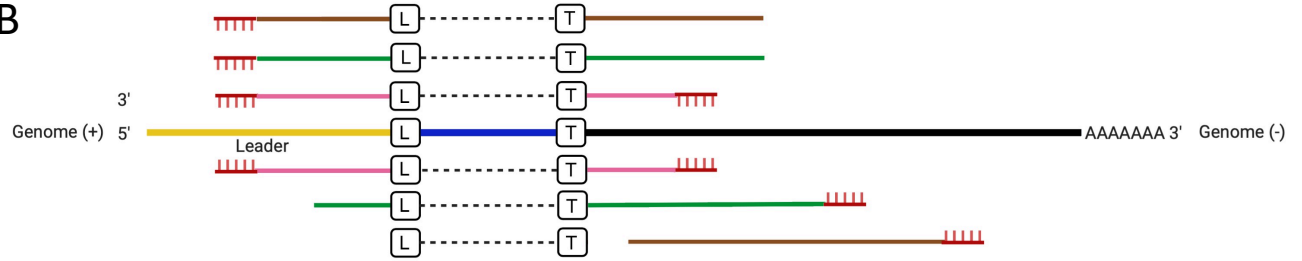
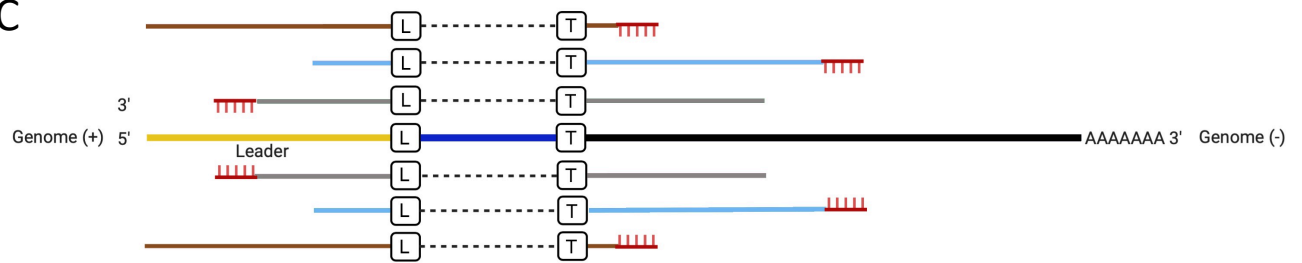
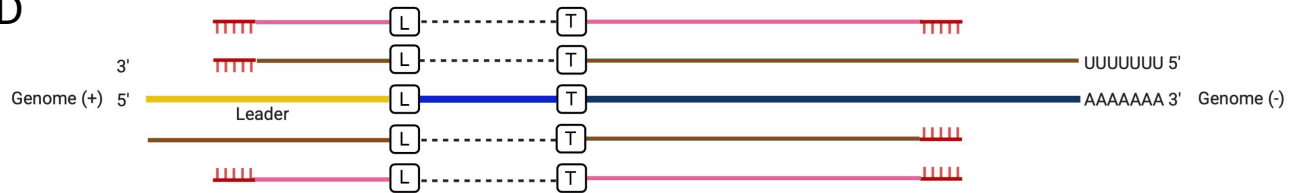
A**B****C****D****E**

Figure 1

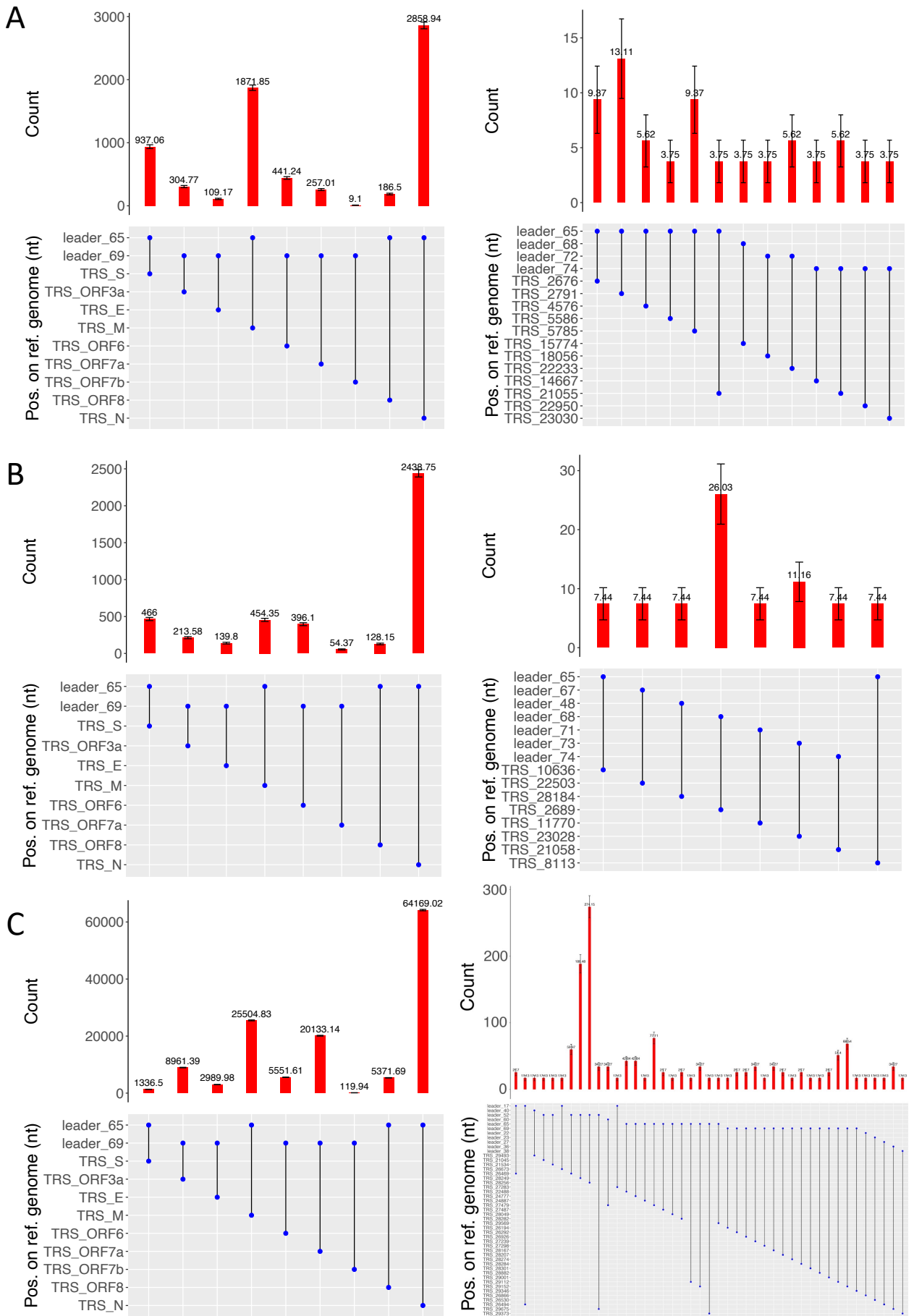


Figure 2

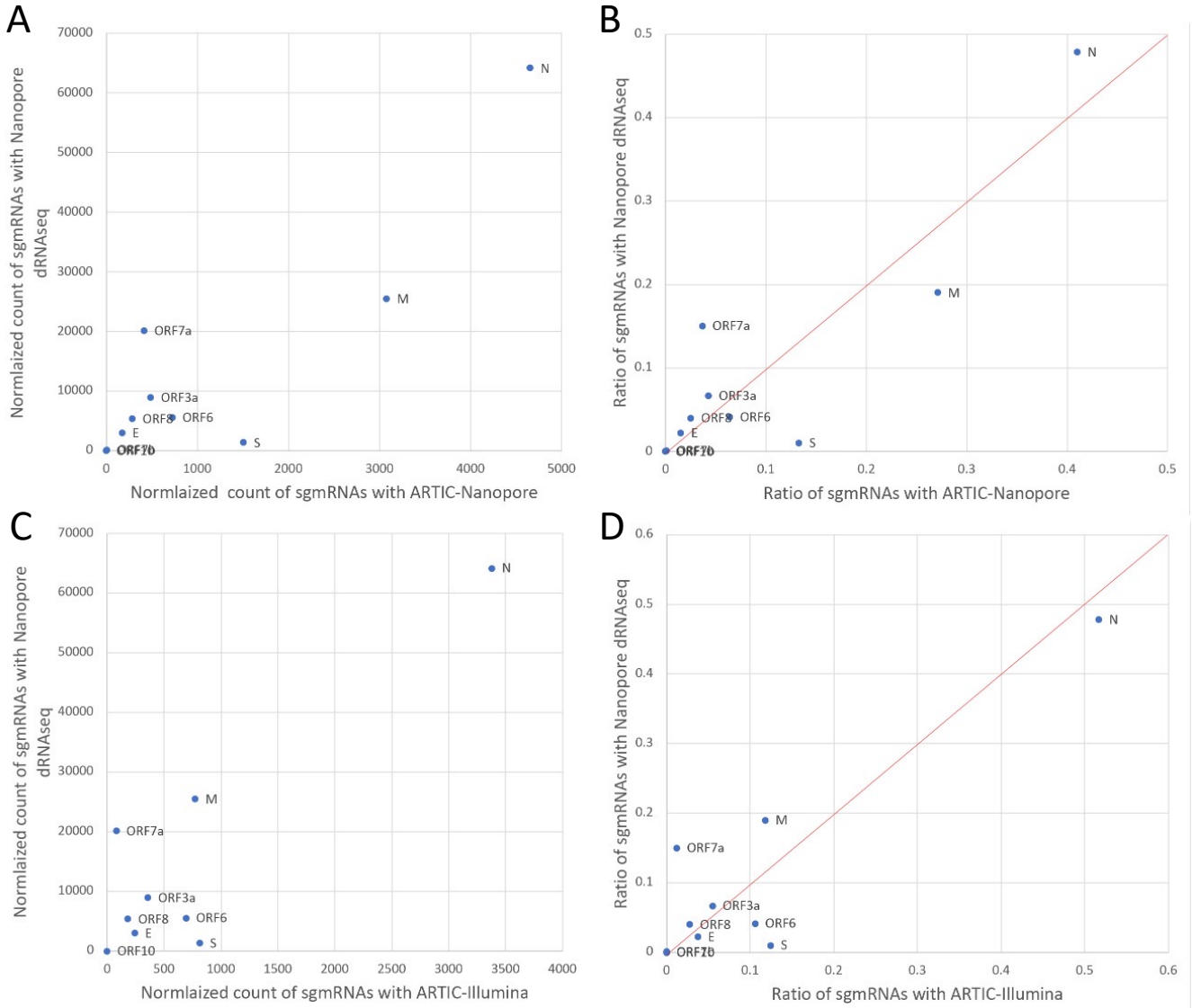


Figure 3

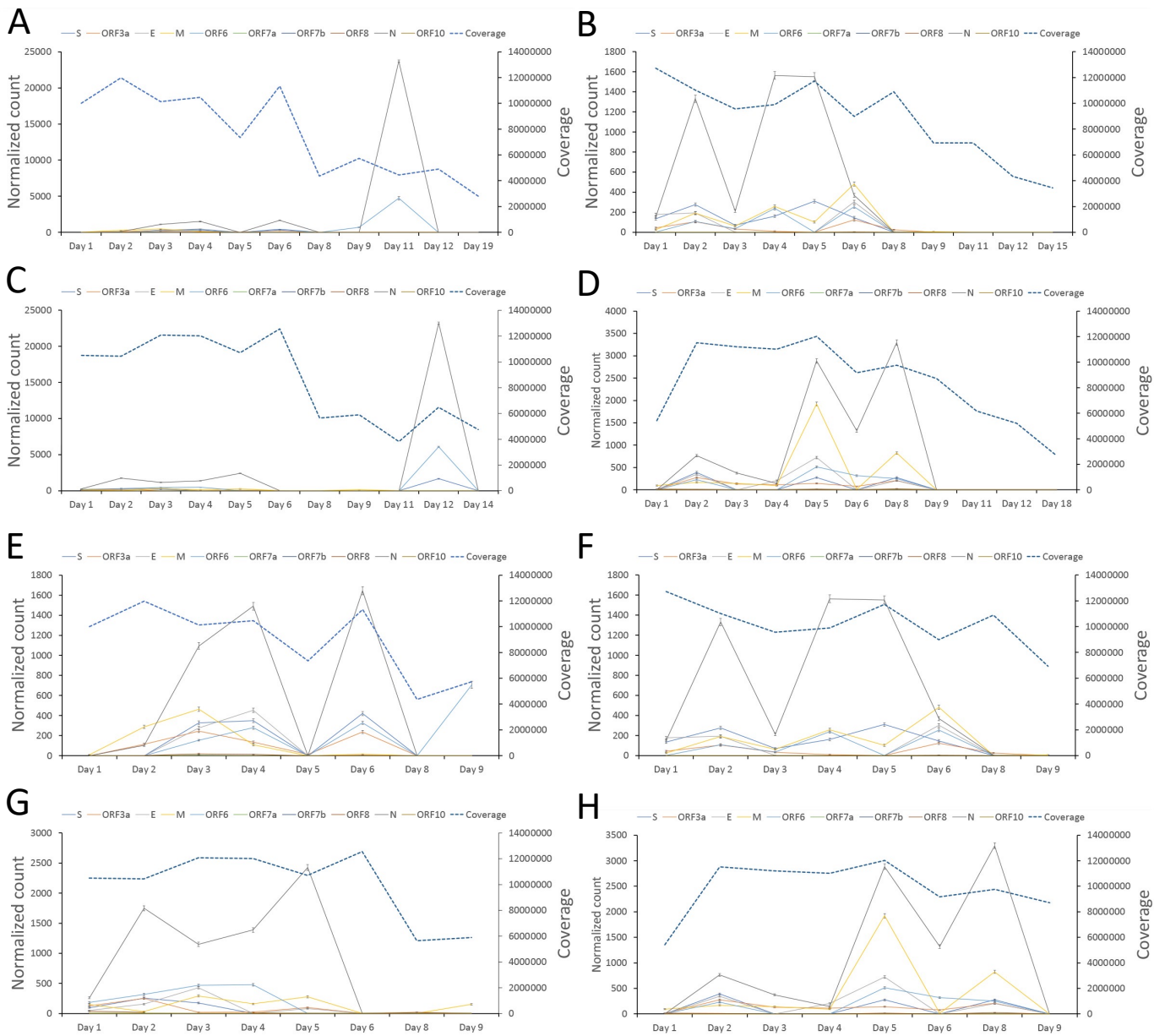


Figure 4

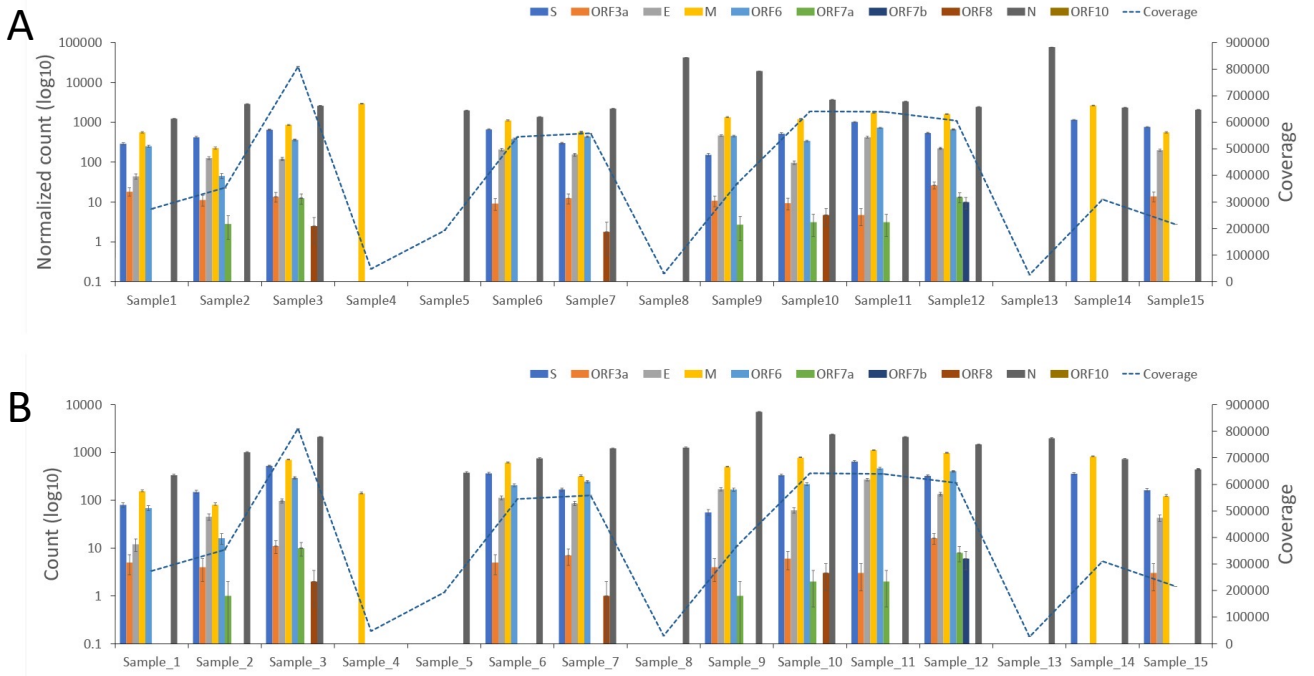


Figure 5

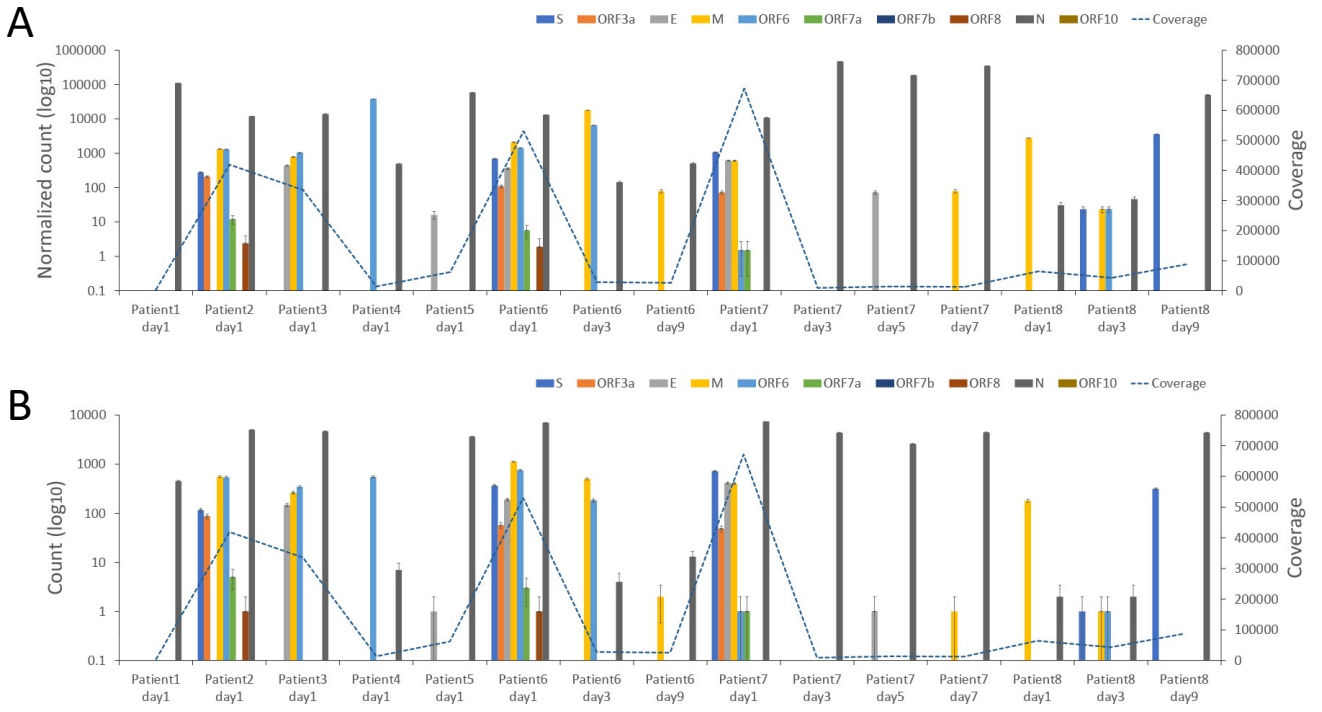
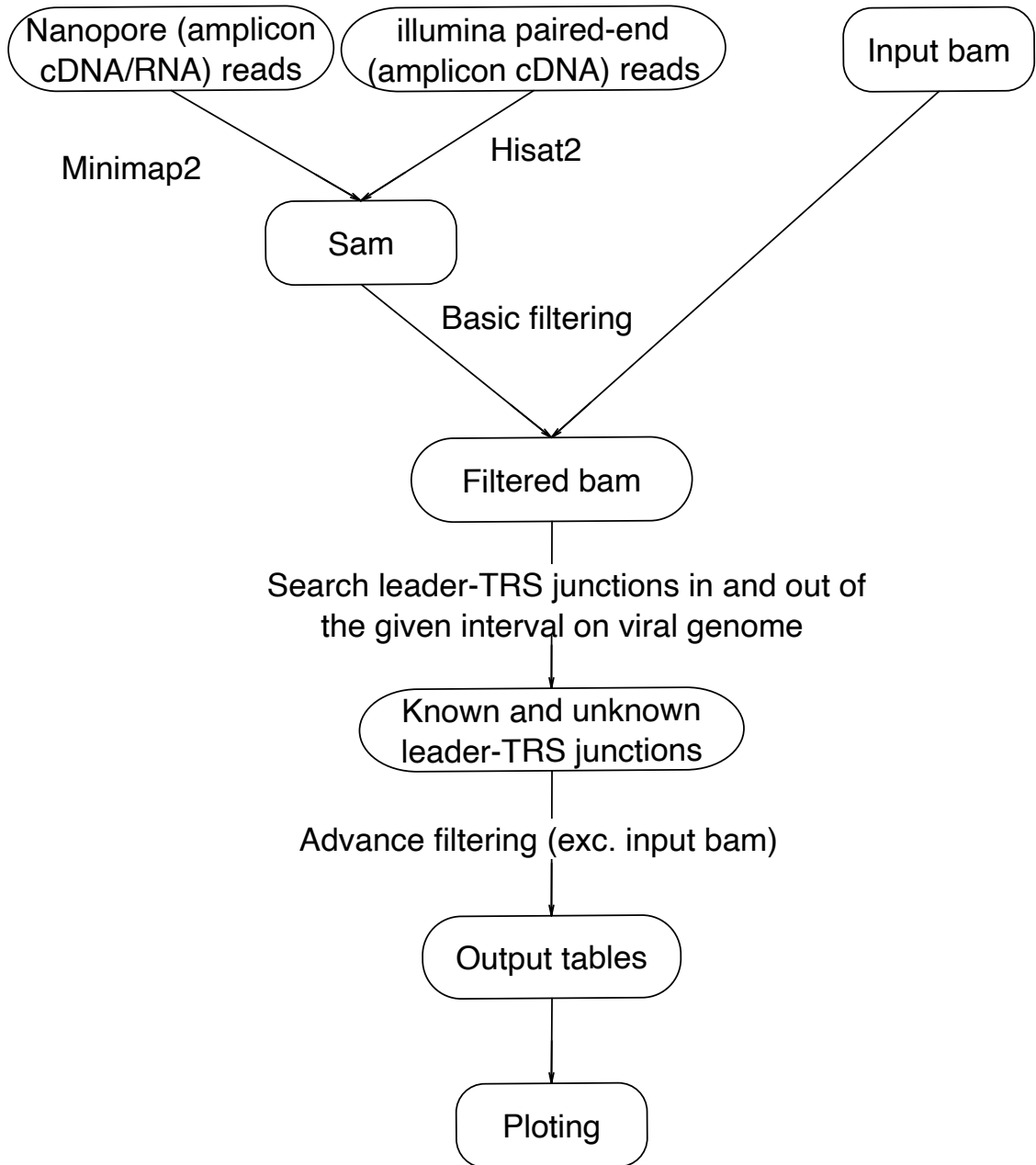
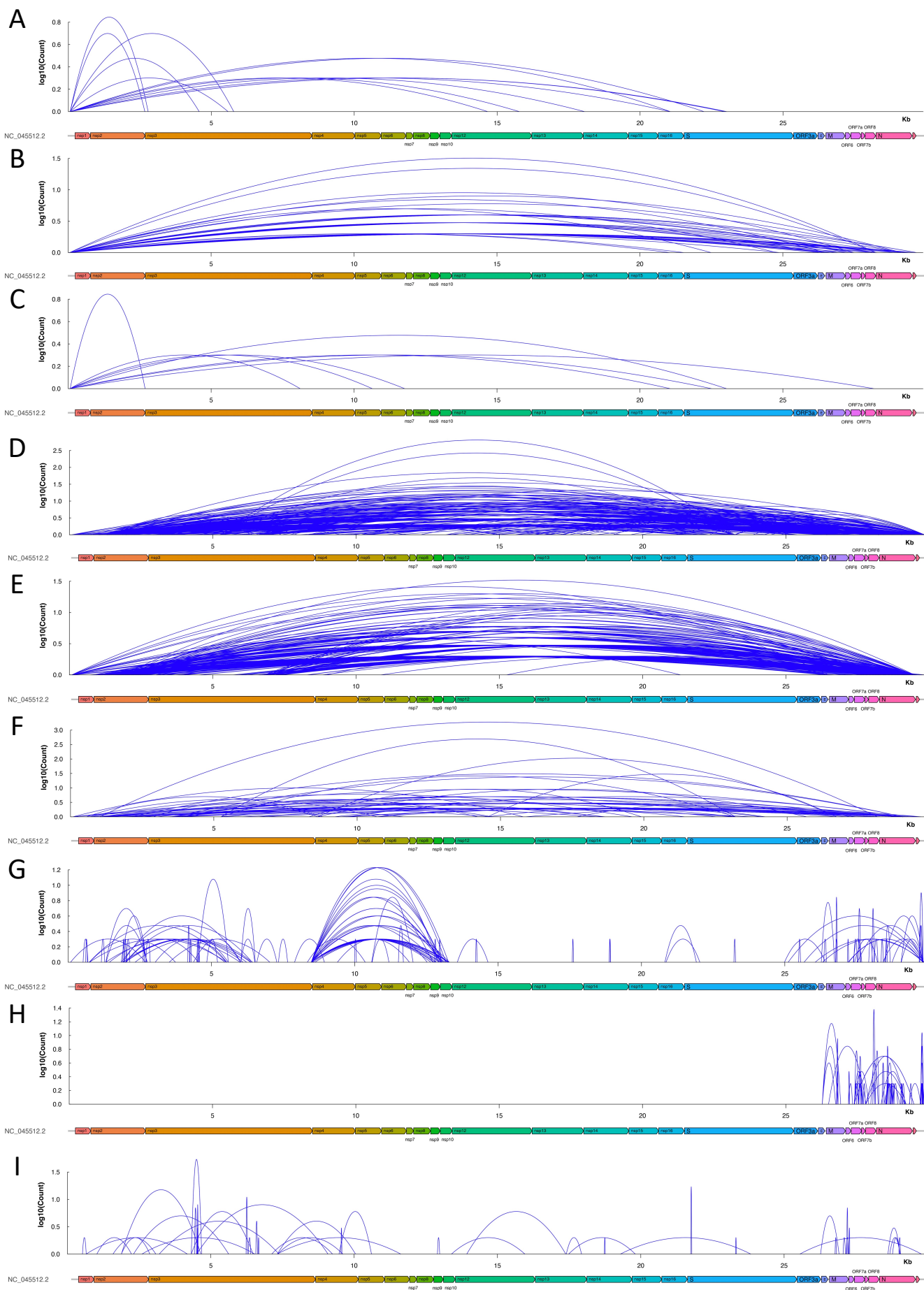


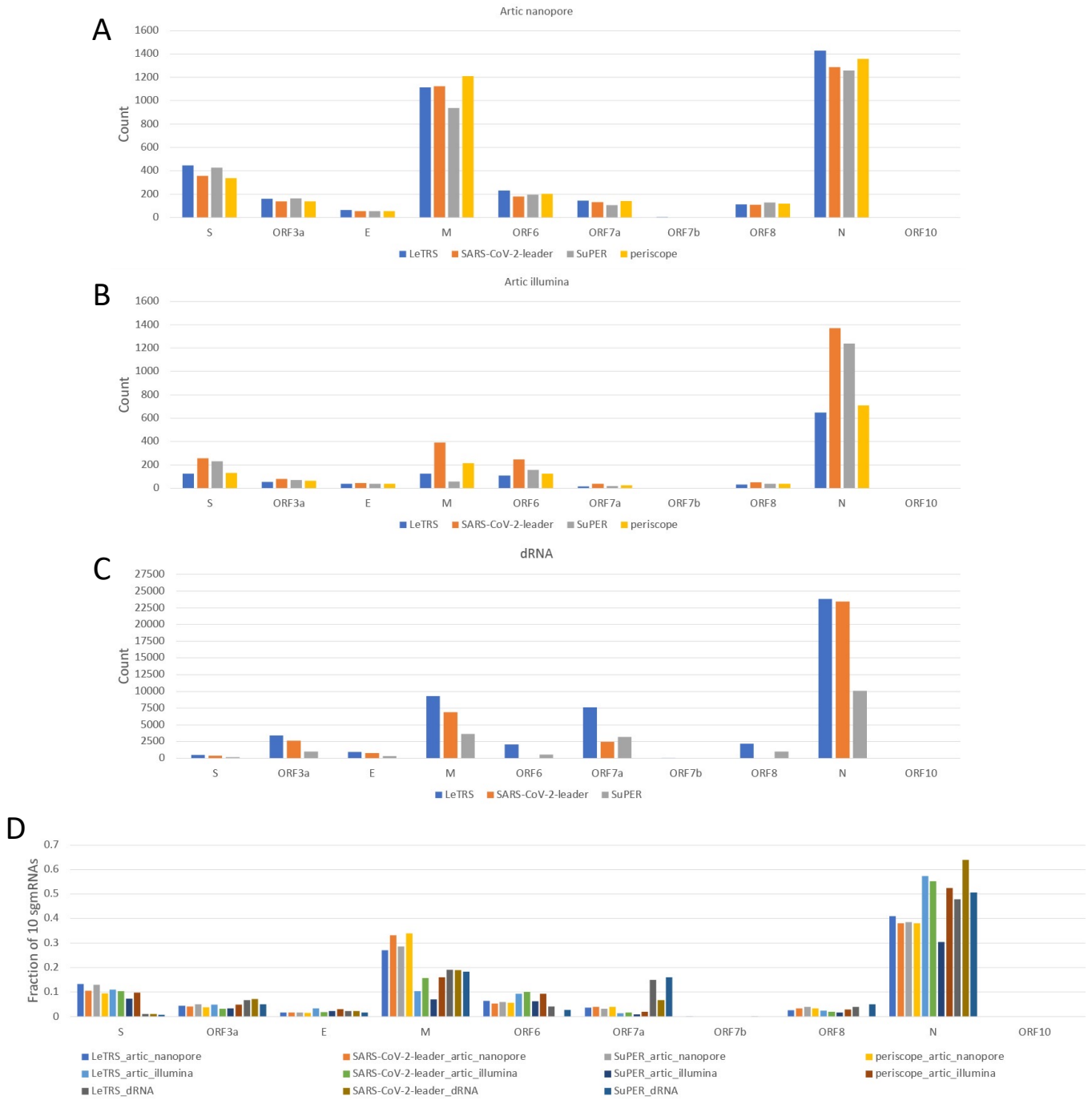
Figure 6



Supplementary Figure 1

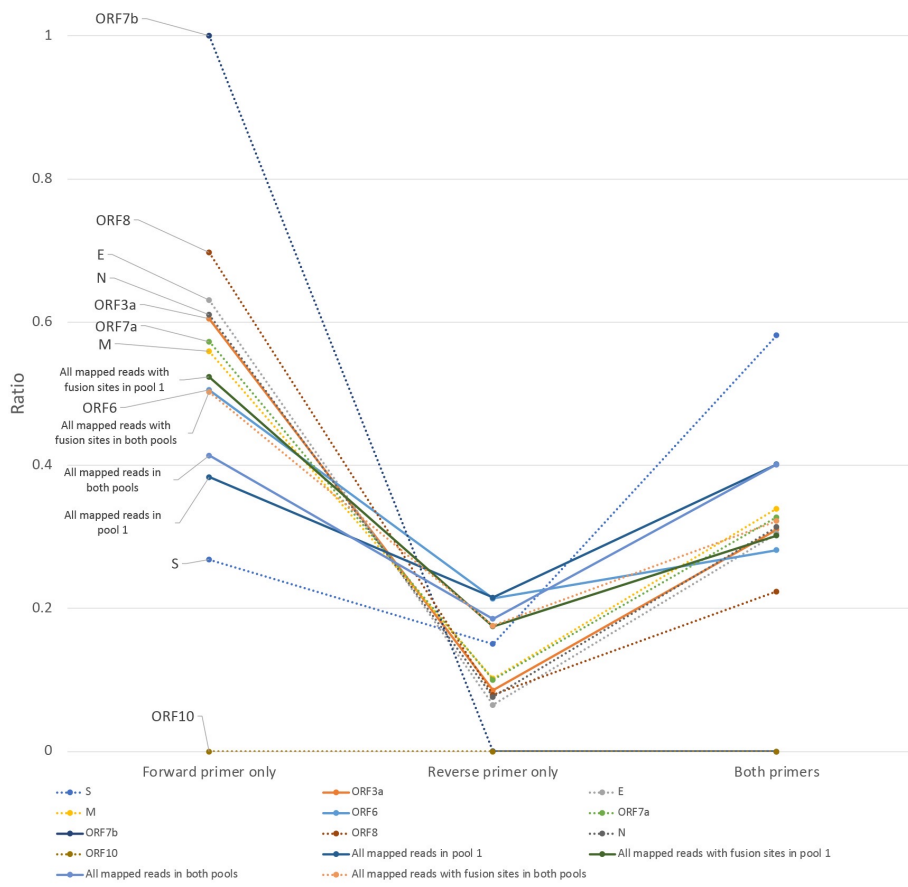


Supplementary Figure 2

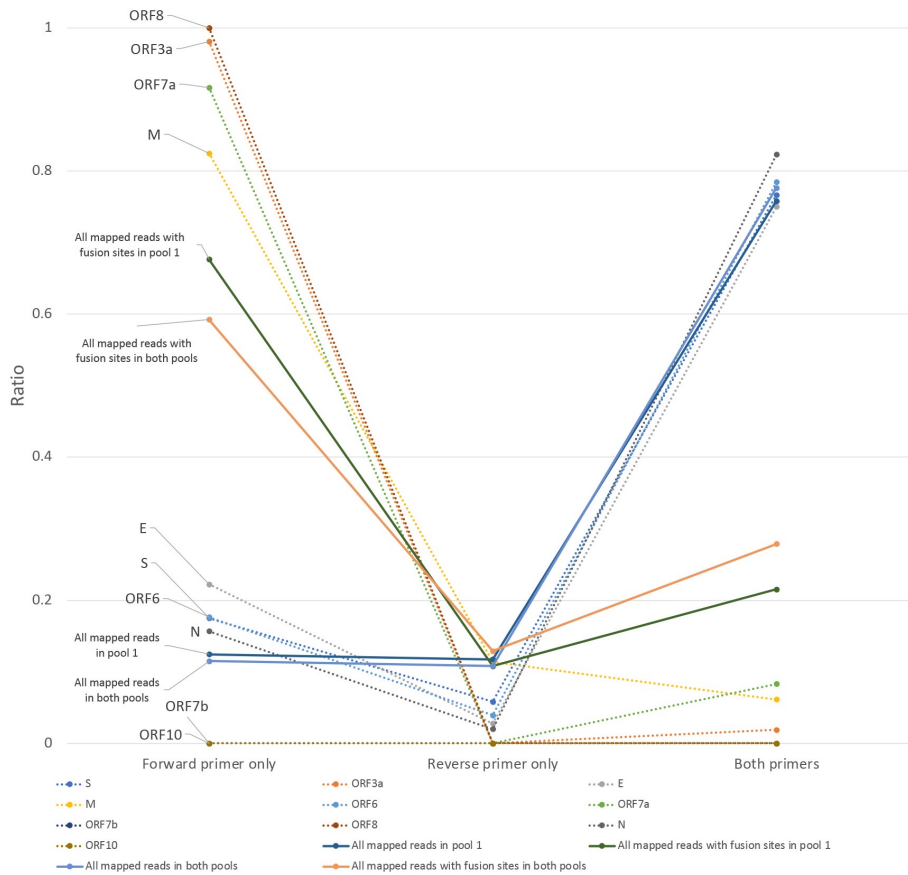


Supplementary Figure 3

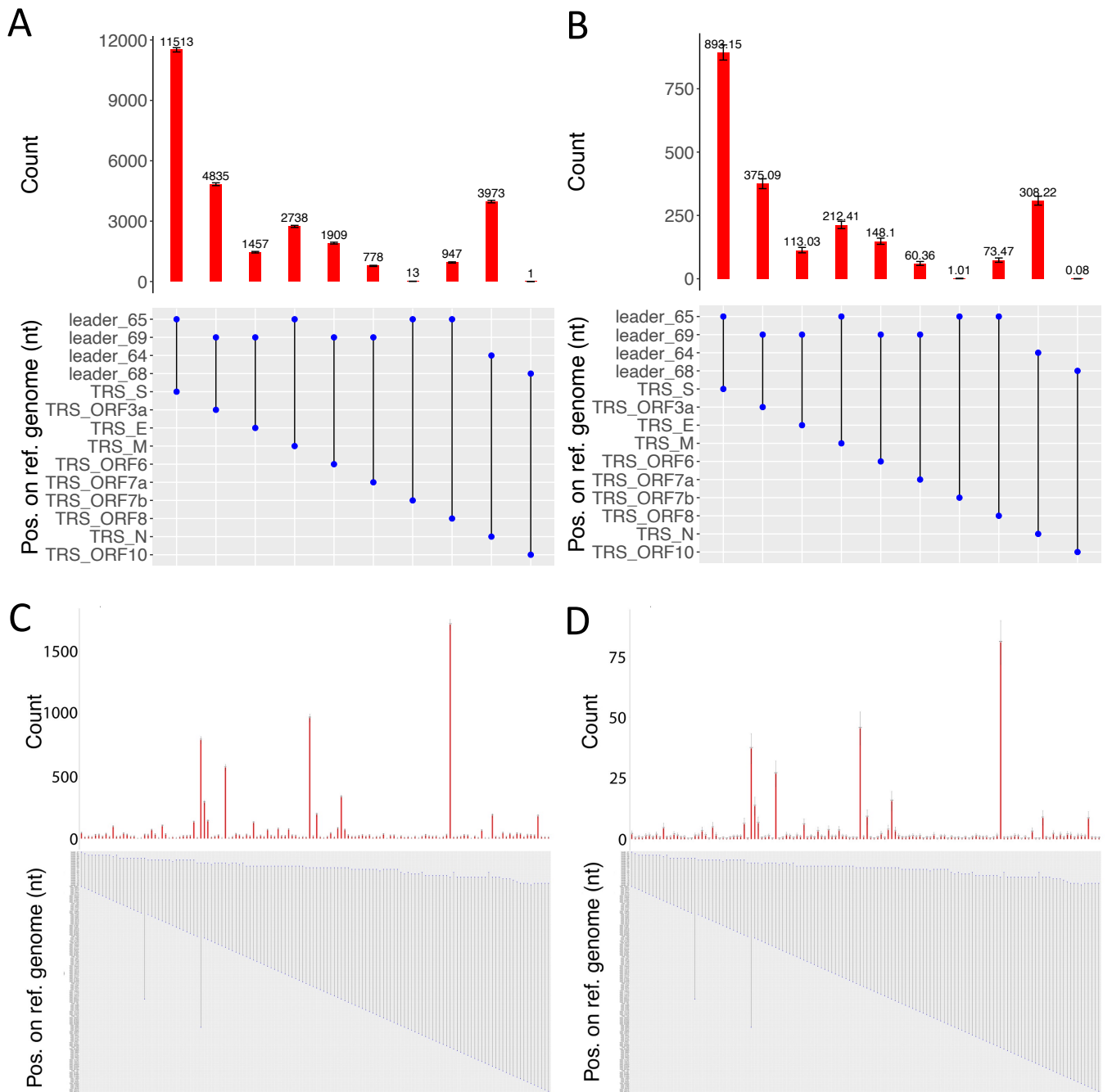
A



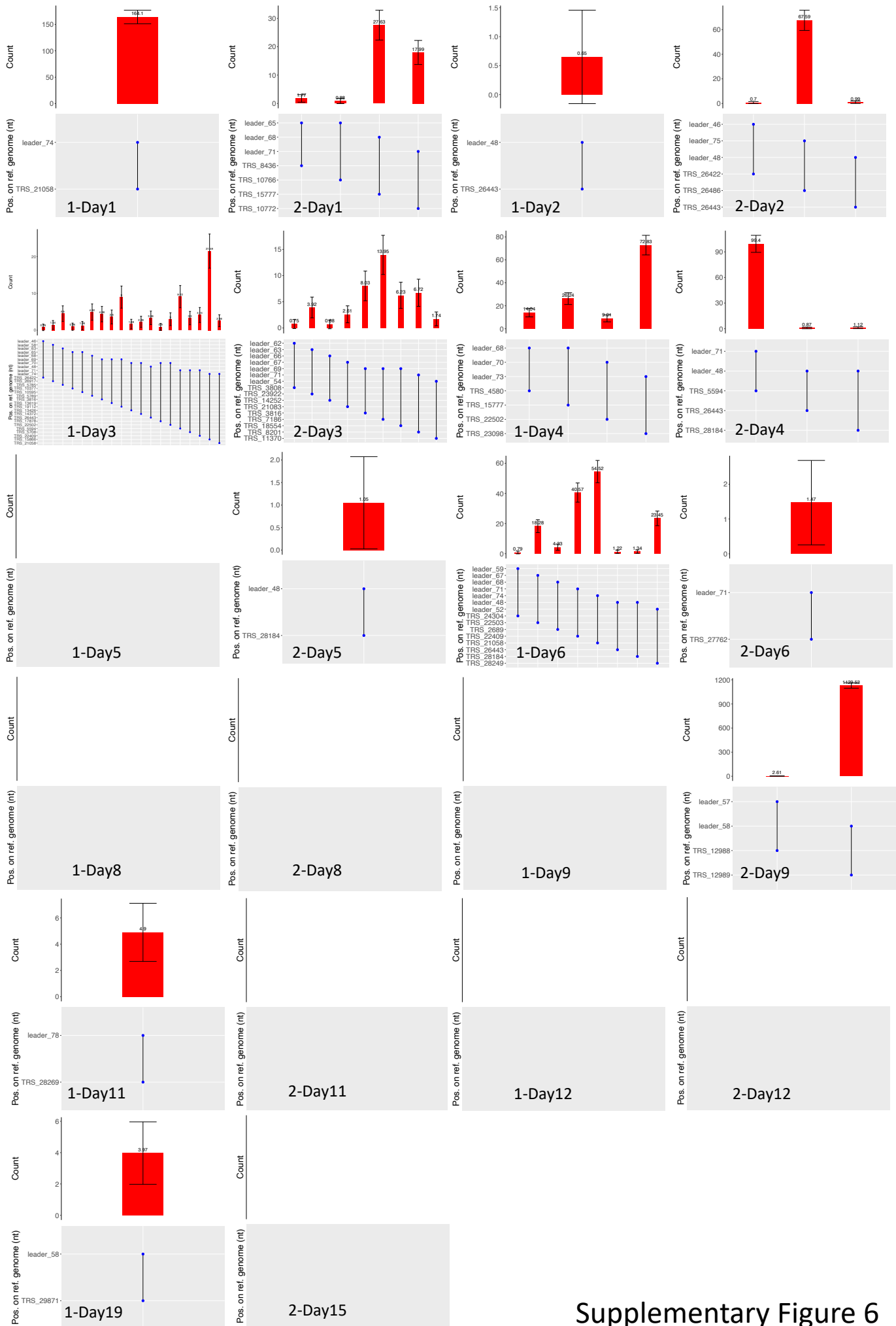
B



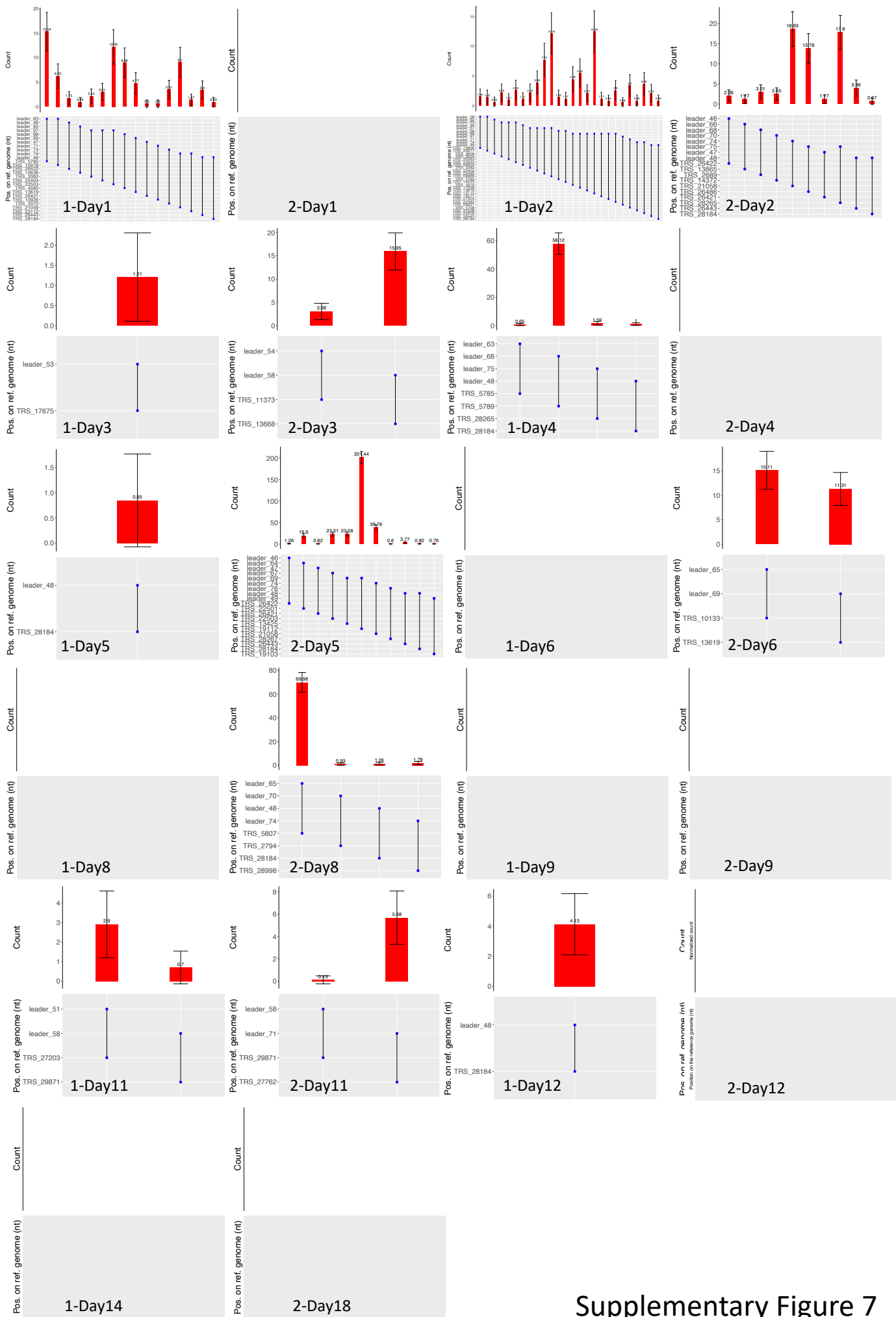
Supplementary Figure 4



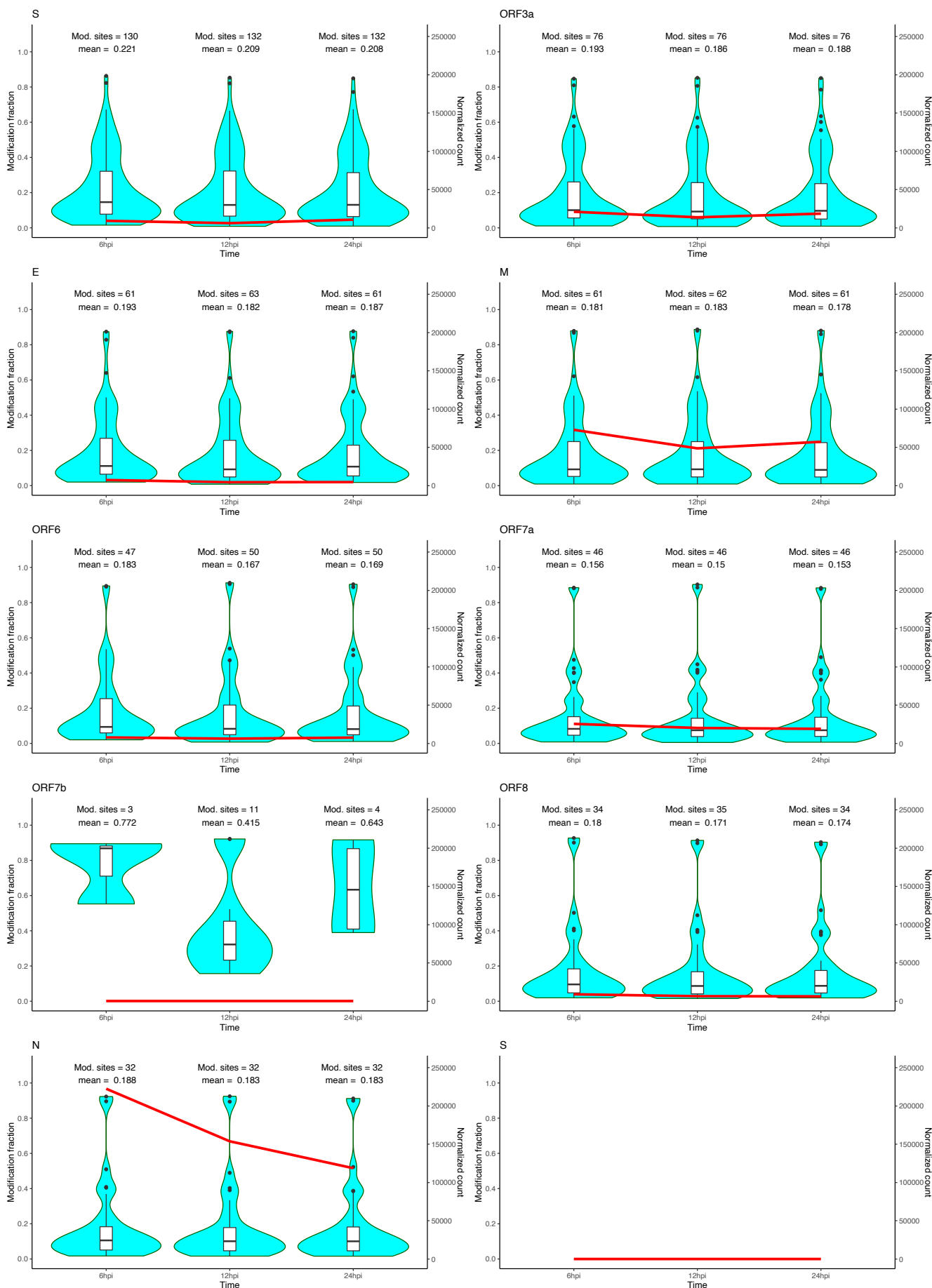
Supplementary Figure 5



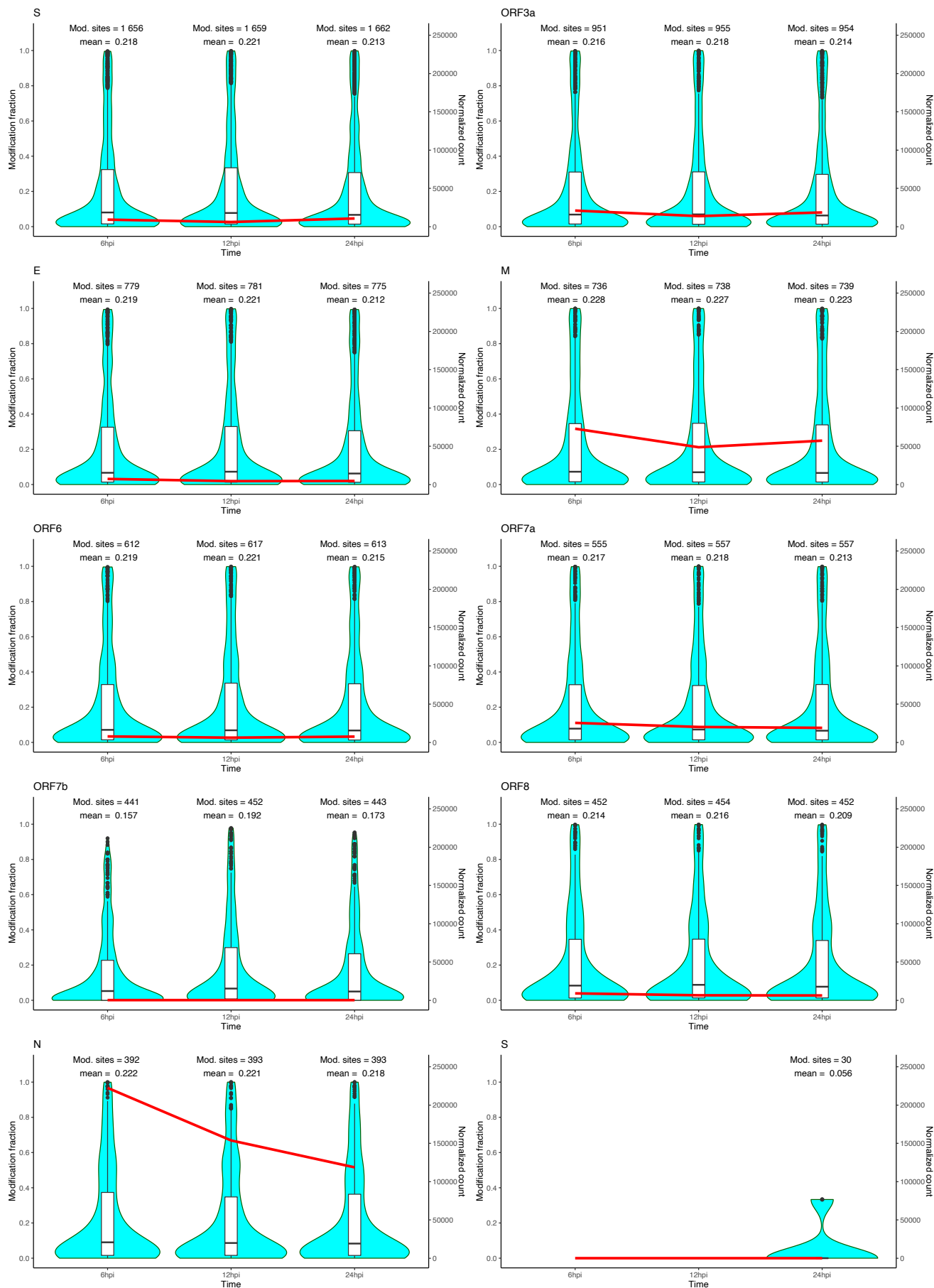
Supplementary Figure 6



Supplementary Figure 7



Supplementary Figure 8



Supplementary Figure 9





Click here to access/download
Supplementary Material
Supplementary_Table_2.xlsx









Click here to access/download
Supplementary Material
Supplementary_Table_5.xlsx

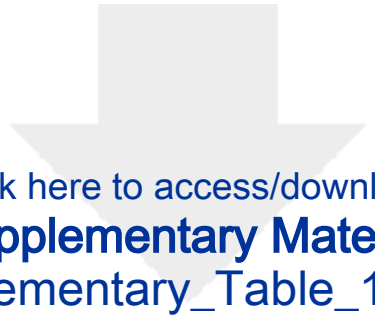












Click here to access/download
Supplementary Material
Supplementary_Table_10.xlsx





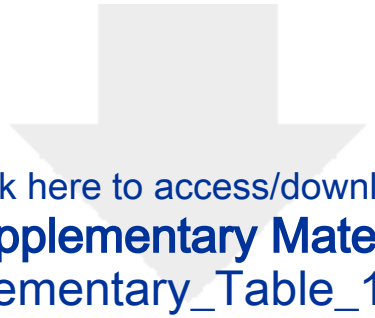
Click here to access/download
Supplementary Material
Supplementary_Table_11.xlsx





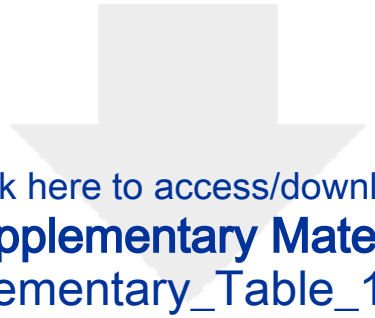
Click here to access/download
Supplementary Material
Supplementary_Table_12.xlsx





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 constructive 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 comments and conducted the additional experiments that the reviewers wanted. We provide a marked-up manuscript 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: In this manuscript, the authors sequenced the SARS-CoV-2 transcriptomes of nasopharyngeal samples from 15 patients using both illumina sequencing and nanopore ARTIC primer3 aplicom sequencing, and developed a computational-pipeline called LeTRS to identify the junctions between the leader sequences in the 5' end of viral genome and the transcriptional regulatory sequence (TRS) within the viral genome (leader-TRS-junction). They first tested and applied their LeTRS tool in several published Nanopore RNA-sequencing data and their own sequencing data to analyses leader-TRS sequence information. They showed that the expression abundance and populations of viral subgenomic mRNA (sgmRNAs) with leader-TRS varies along the time points of post-infection. This study is important to understanding SARS-CoV-2 pathology. However, this article needs many improvements. My major suggestions are as follows:

1. There are two types of leader sequences found in the SARS-CoV-2 sgmRNAs (Dongwan Kim et al., Cell 2020): leader with or without a TRS inside. In the current manuscript, the authors has used their LeTRS tool to identify the sgmRNAs with typical leader with TRS, but did not find the sgmRNAs with non-canonical leaders which do not include TRS inside (TRS-L-independent). I would suggest authors to further extend the studies to sgmRNAs with non-canonical leaders.

Of note, the junctions in these noncanonical transcripts are not derived from a known TRS-B. Some junctions show short sequences (3–4 nt) common between the 50 and 30 sites, suggesting a partial complementarity-guided template switching (“polymerase jumping”). However, the majority do not have any obvious sequences. Thus, we cannot exclude a possibility that at least some of these transcripts are generated through a different mechanism(s).

We have added a function in LeTRS to find sgmRNAs with non-conical leaders (TRS-L-independent) with the “-TRSLindependent” function. This function has been evaluated with the test sample (sequencing RNA from cells infected with SARS-CoV-2) as shown in Supplementary figure 2.

2. SARS-CoV2 genomic and subgenomic mRNAs has multiple types of RNA modifications, such as m6A, 5mC, etc. These modifications has been shown to be regulated and relevant to their polyA tail lengths in sgmRNAs (Kim et al., Cell 2020). I would suggest authors to address if and how RNA modifications levels or types will be dynamically relevant to sgmRNA expression at different time points of post-infection. Aso any preference of RNA modifications in certain types of sgmRNAs (e.g. sgmRNA: S which encodes spike-proteins).

We have direct RNA sequenced the cell cultural samples infected with SARS-CoV-2 at three time points for investigating the relationship between RNA modifications to sgmRNA expression as shown in Supplementary Figures 8 and 9 and Supplementary Table 12. We specifically searched for two different types of methylation. We note that we can only sequenced RNA from cell culture using direct RNA sequencing on the Nanopore. We have found that RNA concentration and quality in clinical samples was insufficient for direct RNA sequencing.

3. I would suggest the authors to compare and evaluate the performance of their LeTRS tools with other similar tools, such as SuPER (Yang Y. et al., Mol. Biol. Evol. 2020), and SARS-CoV-2-leader (Alexandersen S. et al., Nature

Communications 2020), to discuss the strength and weakness of their tool, though the authors has compared their LeTRS tool with another one (Periscope).

We have compared LeTRS with the tools listed by the reviewer using our test data (total RNA from cells infected with SARS-CoV-2) sequenced by three different approaches –ARTIC-Nanopore, ARTIC-Illumina and direct RNA sequencing. This data is presented in Table 1 and Supplementary Figure 3 A, B, C and D. We compare and contrast what the different tools have in common in terms of analysis function and what data types they can function with.

4. I would suggest the authors to re-analyze the public patient's seq data (NCBI PRJNA636225) to examine if the same conclusion about the dysregulation of sgmRNAs at later time points could be derived in different groups of patients.

We have reanalyzed sequencing data from a longitudinal study in two patients (NCBI PRJNA636225) using LeTRSs. The results also indicated a dysregulation of sgmRNAs in late infection from the two patients (Supplementary Table 11). Apart from nuclease resistance and protection by cellular membranes, a phasic pattern of sgmRNA synthesis may also contribute to the presence of sgmRNAs at later time points.

5. It would be nice to have a table to summary the samples and individual information in this study, such as clinical symptoms of patients, gender and age group, and sample collection time point after infection.

Due to the different pathways clinical samples were obtained patient identifying information was not available. For example, samples sequenced using ARTIC-Nanopore were obtained via ISARIC-4C and some patient information was obtained (likely due to these being hospitalized cases – either for treatment or isolation). This is shown in Supplementary Table 10. Samples sequencing using ARTIC-Illumina were sequenced under the auspices of COG-UK and identifying information was not available.

6. The dataset ID provided by this paper (NCBI PRJNA699398) could not be found in the NCBI database. Please the authors address this problem and make the dataset available for the public with a correct ID.

There is a link provided for reviewers:

<https://dataview.ncbi.nlm.nih.gov/object/PRJNA699398?reviewer=tro3da1gmlid1kk6mdjndh7pg0o>

We will release the data if the paper is accepted.

7. The overall presentation, Figures, Tables and language of the paper could need some substantial improvement. The current manuscript includes many misused words, misused punctuation, grammatical errors, and mislabeling.

For examples:

(1) the title is too long. The author should conceive a title with concise but to the key-point.

We have shortened the title.

(2) on page 4, the sentence "for SARS-CoV-2 the core motif is ACGAAC" could be revised as "The core motif of the TRS in SARS-CoV-2 is ACGAAC".

We have changed this.

(3) on page 5, "cell infected in culture" is inaccurate. It could be expressed as "cultured cells with infection".

We have changed this.

(4) on page 13, the word "commonality" might be replaced by "Common properties/features".

We have changed this.

(5) the last sentence on page 13 also need language editing.

We have changed this.

(6) on page 21, the subtitle "search leader-TRS" would be "searching leader-TRS". Pls keep the subtitle to be a short phrase, rather than beginning with a verb.

We have changed this.

(7) pls keep the references in a consistent format. Pls correct the format of Ref. 26, 29 and 30 on page 25-26.

We have changed this.

(8) The authors just need to acknowledge the COG-UK consortia and ISARIC4C consortia, rather than list names of all members in the consortia which occupy 8 pages' space.

We have removed these, apologies this was due to original rules around the consortium authorship statements/acknowledgements.

(9) The x or y bar label and scales in most figures/suppl figures are too small to read.

We have increased the font on the labels.

(10) The Figure legends of all figures are not clear enough and does not provide enough illustrations and explanations for the figures (e.g. Fig 1).

We have changed and expanded the Figure legends.

(11) Supplemental Fig1 could be re-designed to be more clear. For instance, the authors can merge the same steps after the step of <SAM> or <BAM>, to avoid redundant information.

We have changed this.

(12) The legend of table 8 seems exactly same as the legend of table 2. Pls check it.

We have changed this, Tables 1-8 have been moved to Supplementary Tables 1, 2 and 3.

Reviewer #2: "Identification and quantification of SARS-CoV-2 leader subgenomic mRNA gene junctions in nasopharyngeal samples shows phasic transcription in animal models of COVID-19 and dysregulation at later time points that can also be identified in humans"

In this paper, Dong et al describe a new pipeline for identifying subgenomic mRNA from multiple types of sequence data, including amplicon (Illumina and Nanopore) as well as long read nanopore direct RNA or cDNA sequencing. It is useful to have a bioinformatics pipeline which can rapidly identify sgRNA in multiple types of sequence data and has the potential to open large amplicon datasets in particular for further analysis of sgRNA abundance. However, I believe that more validation of the accuracy of abundance estimates from amplicon data is required in order to give the research community more confidence in its use (and limitations).

Major comments:

1. More explanation/detail on methodology would be useful. The authors say that they find the most common peak for the break points of the disjunction site amongst all reads with a break point within a 20bp window of the expected breakpoint. Is there a threshold applied in terms of the difference between the most common and next-most-common breakpoint? Also for the novel sites, is there a clustering algorithm applied, or any site with more than 10 reads is reported?

We used the 20bp window (± 10 bp) of the true splicing sites (known) splicing sites for searching the known sgmRNAs. As noted in the manuscript although we refer to splicing – this is a fusion event. As the minimap2 paper indicated “When INDEL sequencing errors are frequent, it is difficult to find precise splicing sites in this case. If we allow up to 10 bp distance from true splicing sites, 98.4% of aligned introns are approximately correct.” (<https://doi.org/10.1093/bioinformatics/bty191>). Because the known breakpoints are far from each other, the threshold was not defined between the most common and next-most-common breakpoint for the known breakpoints.

We used the coverage cut-off (>10 by default) for the novel sites because we found the novel sites usually have low sequence coverage and don't have a cluster like the known sgmRNAs. Alternatively, these novel sites could be due to RT and sequencing errors, and we note this in the manuscript. LeTRS reports these unknown sites as potential novel sites for future research as all other novel sgmRNAs in the research data.

2. I would like a more direct comparison of sgRNA abundances estimated from amplicon based approach, vs using nanopore amplicon free approach? Its possible to do this only by comparing different tables. It would be easier to digest if there was a x-y plot comparing abundances from different approaches on the same sample. This would help give confidence that the amplicon based approach can provide good estimates. From looking at the tables 1 and 2, it seems that the amplicon approach estimates a lot less sgRNA than the amplicon free approach overall (in terms of normalized counts per million mapped reads). This is to be expected as most of the reads from the amplicon sequencing would be expected to come from the genome. It would be good to see which ORFs are under- and over- represented in

the amplicon data, as I imagine this would also relate to which primer pairs are in the same amplicon pool in the arctic design.

Related to this, it would be good to have an analysis of how the primer design impacts detection of sgRNA. For example, I thought that only one of the primer pools includes a leader primer.

To address this question we infected cells in culture with SARS-CoV-2 and sequenced the viral RNA using three different approaches. Two were amplicon based – based around the ARTIC protocol (an amplicon based system) and also by direct RNA sequencing. This data is shown in Supplementary Tables 1, 2 and 3 to replace the old test data in the Tables 1-8.

With the Artic V3 pipeline, we used two primer pools for the PCR reactions in the whole virus amplification. Please find the primers used in the primer pool 1 and pool 2 at https://github.com/artic-network/artic-ncov2019/blob/master/primer_schemes/nCoV-2019/V3/nCoV-2019.tsv. For the Artic V3 pipeline, only the pool 1 includes a 5'(forward) primer located within the leader region (about < 80) on the genome (please find the position of Artic V3 primers on the virus genome at https://github.com/artic-network/artic-ncov2019/blob/master/primer_schemes/nCoV-2019/V3/nCoV-2019.primer.bed). The LeTRS (v2.0.1) has been modified to only identify the reads with primers in the pool 1, pool 2 or both pools. We compared the read counts evaluated by LeTRS in both ARTIC-Nanopore and ARTIC-Illumina test data for pool 1 and 2, and found only very few reads/read pairs contained the reverse primers with primer pool 2 (Supplementary Table 4 and 5), suggesting the primers in Artic pool 2 are almost not involved the sequencing of leader-TRS regions.

We have done the x-y plot as showed in Figure 3A and C for the reads with at least a primer sequence comparing abundances from different approaches on the same sample. The normalized counts showed a linear relationship between the amplicon based method to the direct sequencing method, while The Artic-Nanopore and Artic-Illumina showed same ratio of known sgRNA as the nanopore direct RNA sequencing approach, except S and orf7a (Figure 3B and D for the reads with at least a primer sequence). This suggested an amplicon based approach can provide good estimates for most of the sgRNAs, especially for N. This normalization method has been applied by <https://doi.org/10.1101/gr.268110.120> and <https://doi.org/10.1038/s41467-020-19883-7>.

PCR based approaches boosted value of denominator reduced the normalized count because a full length of mRNA is counted once with direct RNA sequencing approach will be counted many times with its the small amplicons. Artic illumina got even smaller normalized counts than Artic nanopore approach due to the probably the sequencing bias of illumina during bridge PCR (<https://doi.org/10.1186/gb-2013-14-5-r51>). Therefore, the normalized counts can only be used for the comparison of samples sequenced by same approach when that resulted same PCR and sequencing machine effects. The difference of normalized counts in the samples from amplicon based methods only indicate the relative difference.

Further related to this, it would be good to have a plot which shows the proportion of read counts which are derived from left-primer only, right-primer only or both primers for each sgRNA, and how this compares to the overall ratio of left-only and right-only primers. It seemed odd to me at first glance that there are so many one-sided amplifications, but I imagine this is a small proportion overall, but a sizeable proportion of the reads which can identify sgRNA, due to the lack of primer pairs for

many of the sgRNA. Based on this analysis, it would also be interesting to estimate what is the best depth of coverage of the amplicon panels to get reliable estimates of sgRNA abundance across the different ORFs.

We compared the ratio of reads with forward primers only and reverse primers only and both primers for each sgRNAs to the overall ratios of reads with forward primers only and reverse primers only and both primers in all mapped reads of pool 1 and pool 2 and the mapped reads with any fusion sites in pool 1 and pool 2, found overall ratios showed abundant reads showed same pattern as the reads for sgRNAs (Supplementary Figure 4). This suggested the mass of one side amplification is a nature of amplicon sequencing.

3. It would be good to compare the novel breakpoints with those previously reported, e.g. in Taiorara et al, figure 2 and supplementary figure 6 (<https://doi.org/10.1101/2020.03.05.976167>). I can see that many of them line up with those you report in table 4, and I believe this sup

Taiorara et al didn't attach the exact breakpoints positions with their figure, but we generated a similar figure for comparison (Figure 7c). Figure 7c showed some similar breakpoints positions with Figure 2 of Taiorara et al's paper.

4. Is there much overlap in the novel break points detected using nanopore amplicon ARTIC v3 vs nanopore dRNA? It would be good to have an extra column in Table 8 and table 4 indicating which of the breakpoints discovered in dRNA were also discovered in amplicon sequencing and vice versa. This will hopefully shed light on relative strengths of the two approaches. Similarly it would be useful to compare nanopore ARTIC and illumina ARTIC in this regard

As described above we have moved the new test data from a unique cell culture sample to Supplementary Tables 1, 2 and 3 for Artic-Nanopore, Artic-Illumina and nanopore direct RNA sequencing. We didn't find any exactly the same novel fusion sites in these three approaches. To note in the publication describing minimap2 the paper details "In general, minimap2 is more consistent with existing annotations. It finds more junctions with a higher percentage being exactly or approximately correct" and "When INDEL sequencing errors are frequent, it is difficult to find precise splicing sites in this case. If we allow up to 10 bp distance from true splicing sites, 98.4% of aligned introns are approximately correct." (<https://doi.org/10.1093/bioinformatics/bty191>). Therefore, it is very difficult to identify the exact novel fusion sites. Novel leader-TRS junctions were also known as leader dependent noncanonical fusions. LeTRS also has a function to identify leader independent long-distance (>5,000 nt) fusion and local joining yielding a deletion between proximal sites (20–5,000 nt distance) in the sequencing reads. If we look at the pattern of the fusion sites, some of the novel leader-TRS junctions (noncanonical fusions) and leader independent fusions in the test sample were supported by all three sequencing methods (Supplementary figure 2) with similar fusion sites.

The strength of LeTRS to identify the known breakpoints is much stronger than identifying novel sites, because LeTRS controls the aligner to search the known breakpoints with the guide of known annotations. As the paper said "In general, minimap2 is more consistent with existing annotations. It finds more junctions with a higher percentage being exactly or approximately correct" (<https://doi.org/10.1093/bioinformatics/bty191>).

5. Its hard to assess the evidence supporting the biphasic expression without having some idea of the error in the abundance estimates (also commented on this more below);

We have calculated the standard deviation of a binomial distribution as error bar. The data supports that biphasic expression/abundance of sgmRNAs occurs.

6. The conclusion of dysregulation in samples taken from patients many days into their infection is made only on a small number of samples. Also in Figure 4, the time post sample is not indicated. I presume the information is in one of the supplementary tables, but the submitted pdf has messed up these tables (its somewhere in the 729 page pdf) . Nevertheless, it seems that the data supporting this conclusion is a bit thin, and I would be cautious in including that observation in the title of the paper.

We have changed the title to reflect this comment.

Minor comments:

1. In figure legends (e.g. figure 1) you say the numbers in brackets are: reads with left primers, reads with right primers, reads with both primers. I can see from the numbers that these are not exclusive, but it might be easier to digest if you showed left-only, right-only and both

We modified the LeTRS to show forward-only, reverse-only and both primers

2. You make a point in the paper about whether the left break occurs at position 64 or 69. One thing I would worry about is that microhomology between TRS-L and TRS-R might make it difficult to be exactly sure of the breakpoint (because the sgRNA includes only one copy of the TRS, but its hard to know if it's the left or the right which is included, the aligner could equally well align to TRS-L and skip TRS-R or vice versa, and this would shift the coordinates slightly. Are the enough snp differences in TRS-L or TRS-R to be confident either way, and if so, does this have implications for whether TRS-L or TRS-R is retained in the sgRNA?

For the known sgmRNA, we used the known annotation of breakpoints to guide the alignments and allowing a (± 10 bp) window of the true splicing/fusion sites for searching the breakpoints - if this would shift the coordinates slightly. Even if TRS-L or TRS-R is retained in the sgmRNAs, the implications will be random and equal to all samples with same sequencing approach and alignment tool. This should not affect the evaluation of the ratio of sgmRNAs and relative abundance across samples. We have also compared the number of reads for sgmRNAs with the other methods (tool called SARS-CoV-2-leader) that is to search a tag sequence within leader in reads but not the breakpoints of reads. SARS-CoV-2-leader produced a similar read count as LeTRS for the Artic-Nanopore (Supplementary Figure 3A) and Nanopore direct sequencing (Supplementary Figure 3C). SARS-CoV-2-leader produced more counts than LeTRS for Artic-Illumina, because LeTRS counts the read pairs but not reads (Supplementary Figure 3B). There are difficulties in searching for novel breakpoints, although we treat novel breakpoints as a potential sign of novel sgmRNAs for future research.

3. Figure 1 panels B,C,D were a bit confusing. Why is the reference sequence in the middle. It would be good if the caption could be expanded to help the reader understand these panels in particular.

The figure legend has been changed but we would like to keep the reference sequence in the middle to show the forward and reverse amplification possibilities.

4. The tables (table 1 to 8) and the figure 1A represent a lot of the same information, but the numbers don't line up exactly, because in the figures you only use counts which have both primers. It would be best to decide which to represent because it's confusing to have the same data presented twice essentially but in slightly different ways.

We have changed this and now consistently only used the reads containing at least one primer to plot data.

5. In figure 1 you present the normalized abundance to 2 decimal places, but its very unlikely that you have that level of precision. It would be good if you could add error bars to estimate the uncertainty in the abundance estimate (e.g. calculated using a binomial distribution).

We have calculated the standard deviation of a binomial distribution as an error bar.

6. In figure 3, its hard to know how much error there is in each of the measurements. By showing the normalized value, its also hard to see what is the absolute change in the read counts. Ideally you would show either the read counts, or show error bars around the abundance estimates.

We now show error bars.

7. Is there a mistake in the title of Table 8: "The LeTRS output table for novel sgmRNA in the tested Nanopore ARTIC v3 primers." Because the title of table 2 seems the same: Table 2. The LeTRS output table for novel sgmRNA in the tested Nanopore ARTIC v3 primers" . One of these approaches does not seem to find novel breakpoints, but the other does, presumably Table 8 should be illumina based on the ordering?

We have changed this. Tables 1-8 have been moved to Supplementary Tables 1, 2 and 3.

8. Error in caption of table 1: " Normalized count=(Read count-Total number of read mapped on reference genome)*1000000"

We have changed this. Tables 1-8 have been moved to Supplementary Tables 1, 2 and 3.

9. In the supplementary figures, the captions you saay:" Supplementary Figure 3. Raw (A and C) and normalised (B and D) expected (upper) and novel (lower) leader-TRS gene junctions count in the infecting SARS-CoV-2 inoculum source used for NHP study, sequenced by Illumina ARTIC method (Supplementary Table 8)."

I found the use of "expected" here confusing, because it implied to me that you had estimated expected counts. I would prefer the use of the term canonical, or something like that.

We have changed "expected" to "canonical". Supplementary Figure 3 has become Supplementary Figure 5.