

## Human genome integration of SARS-CoV-2 contradicted by long-read sequencing

Nathan Smits<sup>1,10</sup>, Jay Rasmussen<sup>2,10</sup>, Gabriela O. Bodea<sup>1,2,10</sup>, Alberto A. Amarilla<sup>3,10</sup>, Patricia Gerdes<sup>1</sup>, Francisco J. Sanchez-Luque<sup>4,5</sup>, Prabha Ajjikuttira<sup>2</sup>, Naphak Modhiran<sup>3</sup>, Benjamin Liang<sup>3</sup>, Jamila Faivre<sup>6</sup>, Ira W. Deveson<sup>7,8</sup>, Alexander A. Khromykh<sup>3,9</sup>, Daniel Watterson<sup>3,9</sup>, Adam D. Ewing<sup>1</sup>, Geoffrey J. Faulkner<sup>1,2\*</sup>

<sup>1</sup>Mater Research Institute - University of Queensland, TRI Building, Woolloongabba QLD 4102, Australia.

<sup>2</sup>Queensland Brain Institute, University of Queensland, Brisbane QLD 4072, Australia.

<sup>3</sup>School of Chemistry and Molecular Biosciences, University of Queensland, Brisbane QLD 4072, Australia.

<sup>4</sup>GENYO, Pfizer-University of Granada-Andalusian Government Centre for Genomics and Oncological Research, PTS Granada 18016, Spain.

<sup>5</sup>MRC Human Genetics Unit, Institute of Genetics and Cancer (IGC), University of Edinburgh, Western General Hospital, Edinburgh EH4 2XU, United Kingdom.

<sup>6</sup>INSERM, U1193, Paul-Brousse University Hospital, Hepatobiliary Centre, Villejuif 94800, France.

<sup>7</sup>Kinghorn Centre for Clinical Genomics, Garvan Institute of Medical Research, Sydney NSW 2010, Australia.

<sup>8</sup>St Vincent's Clinical School, Faculty of Medicine, University of New South Wales, Sydney NSW 2052, Australia.

<sup>9</sup>Australian Infectious Diseases Research Centre, Global Virus Network Centre of Excellence, Brisbane QLD 4072, Australia.

<sup>10</sup>These authors contributed equally.

\*Corresponding author: [faulknergj@gmail.com](mailto:faulknergj@gmail.com) (G.J.F)

1 **Abstract**

2 A recent study proposed severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)  
3 hijacks the LINE-1 (L1) retrotransposition machinery to integrate into the DNA of infected  
4 cells. If confirmed, this finding could have significant clinical implications. Here, we applied  
5 deep (>50×) long-read Oxford Nanopore Technologies (ONT) sequencing to HEK293T cells  
6 infected with SARS-CoV-2, and did not find any evidence of the virus existing as DNA. By  
7 examining ONT data from separate HEK293T cultivars, we resolved the complete sequences  
8 of 78 L1 insertions arising *in vitro* in the absence of L1 overexpression systems. ONT  
9 sequencing applied to hepatitis B virus (HBV) positive liver cancer tissues located a single  
10 HBV insertion. These experiments demonstrate reliable resolution of retrotransposon and  
11 exogenous virus insertions via ONT sequencing. That we found no evidence of SARS-CoV-2  
12 integration suggests such events *in vivo* are highly unlikely to drive later oncogenesis or  
13 explain post-recovery detection of the virus.

14 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is a single-stranded ~30kbp  
15 polyadenylated RNA betacoronavirus<sup>1</sup>. SARS-CoV-2 does not encode a reverse transcriptase  
16 (RT) and therefore is not expected to integrate into genomic DNA as part of its life cycle.  
17 This assumption is of fundamental importance to the accurate diagnosis and potential long-  
18 term clinical consequences of SARS-CoV-2 infection, as demonstrated by other viruses  
19 known to incorporate into genomic DNA, such as human immunodeficiency virus 1 (HIV-1)  
20 and hepatitis B virus (HBV)<sup>2-4</sup>.

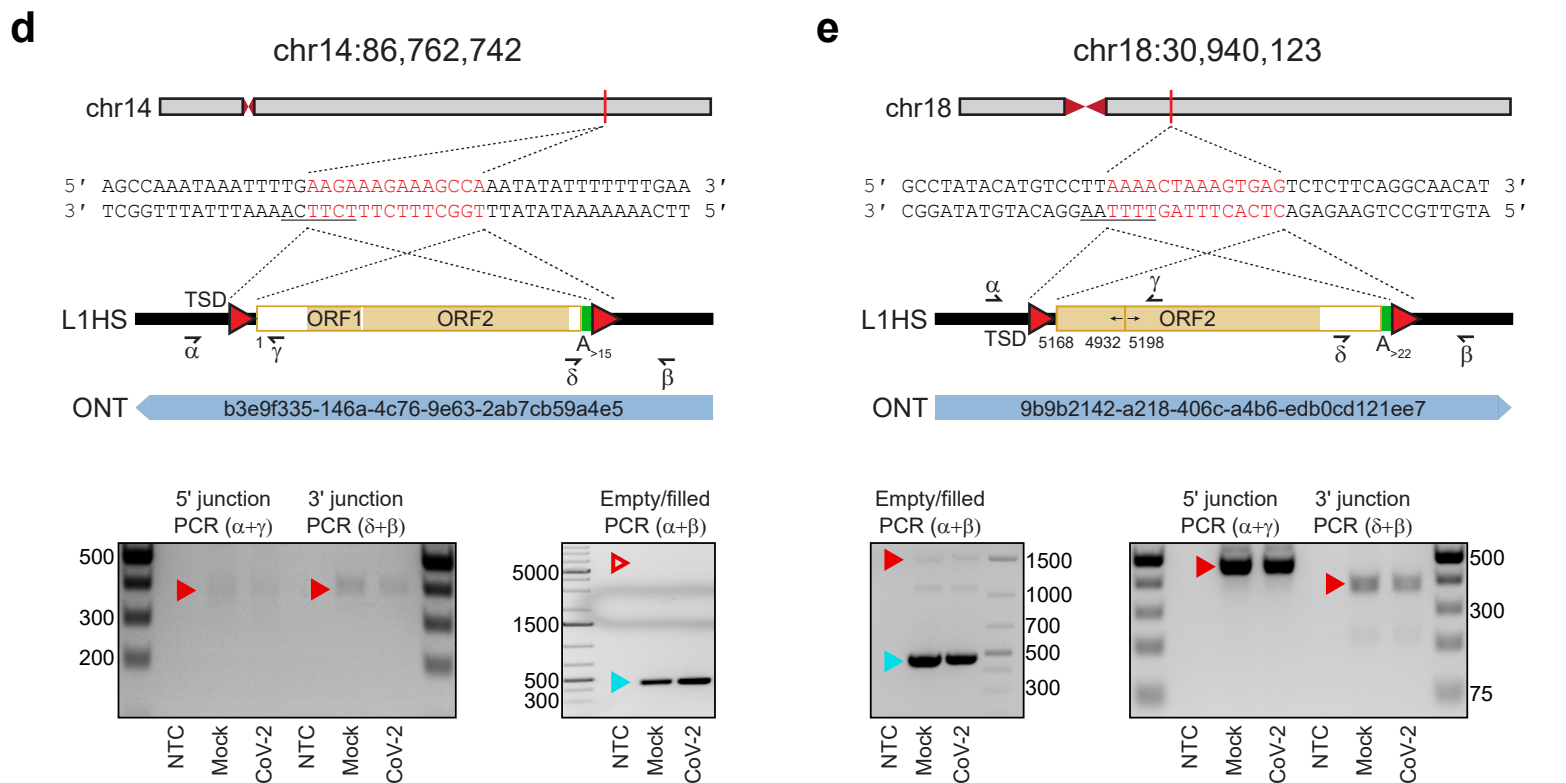
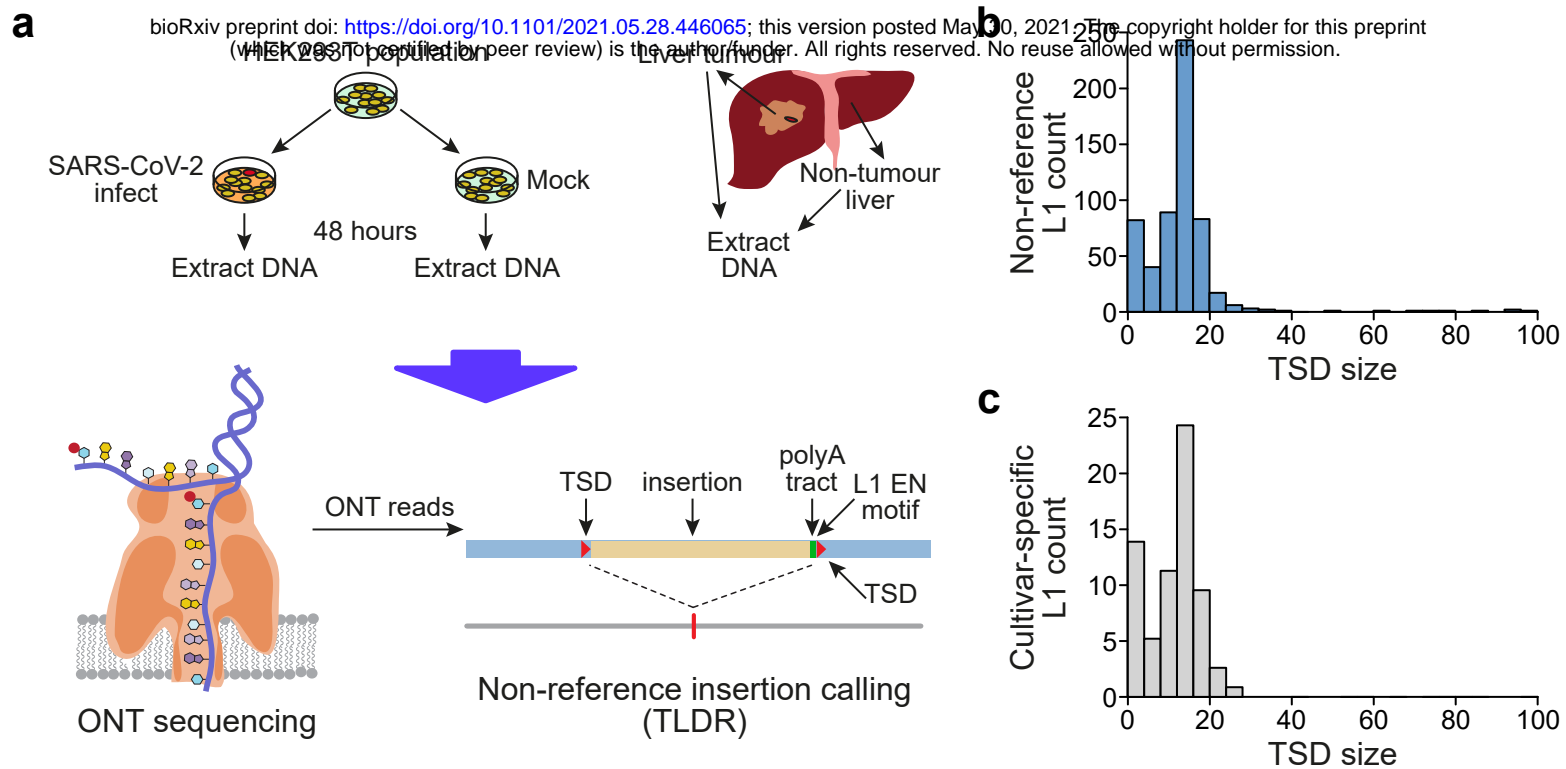
21 LINE-1 (L1) retrotransposons reside in all mammalian genomes<sup>5</sup>. In humans, L1  
22 transcribes a bicistronic mRNA encoding two proteins, ORF1p and ORF2p, essential to L1  
23 mobility<sup>6</sup>. ORF2p possesses endonuclease (EN) and RT activities, and exhibits strong *cis*  
24 preference for reverse transcription of L1 mRNA<sup>6-11</sup>. Nonetheless, the L1 protein machinery  
25 can *trans* mobilise polyadenylated cellular RNAs, including those produced by non-L1  
26 retrotransposons and protein-coding genes<sup>11-14</sup>. Somatic L1 *cis* mobilisation is observed in  
27 embryonic cells, the neuronal lineage, and various cancers<sup>15-20</sup>. By contrast, somatic L1-  
28 mediated *trans* mobilisation is very rare<sup>15,16,20</sup> and is likely repressed by various  
29 mechanisms<sup>8,16,21-23</sup>. Less than one cellular RNA *trans* insertion is expected for every 2000  
30 *cis* L1 insertions<sup>10</sup>. Both *cis* and *trans* L1-mediated insertions incorporate target site  
31 duplications (TSDs) and a 3' polyA tract, and integrate at the degenerate L1 EN motif 5'-  
32 TTTT/AA<sup>6,11-14,24-27</sup>. These sequence hallmarks can together discriminate artifacts from  
33 genuine insertions<sup>28</sup>.

34 In a recent study, Zhang *et al.* overexpressed L1 in HEK293T cells, infected these  
35 with SARS-CoV-2, and identified DNA fragments of the virus through PCR amplification<sup>29</sup>.  
36 These results, alongside other less direct<sup>30,31</sup> analyses, were interpreted as evidence of SARS-  
37 CoV-2 genomic integration<sup>29</sup>. Crucially, Zhang *et al.* then detected 63 putative SARS-CoV-2  
38 integrants by Oxford Nanopore Technologies (ONT) long-read sequencing. Of these, only a  
39 single integrant on chromosome X was spanned by an ONT read aligned to one locus, and  
40 was flanked by potential TSDs (**Extended Data Fig. 1**). However, this SARS-CoV-2  
41 integrant did not incorporate a 3' polyA tract, as is expected for an L1-mediated insertion, and  
42 involved an unusual 28kb internal deletion of the SARS-CoV-2 sequence. The SARS-CoV-2  
43 integrants reported by Zhang *et al.* were 26-fold enriched in exons, despite the L1 EN  
44 showing no preference for these regions<sup>32,33</sup>. Zhang *et al.* also used Illumina short-read  
45 sequencing to map putative SARS-CoV-2 integration junctions in HEK293T cells without L1  
46 overexpression. A lack of spanning reads and the tendency of Illumina library preparation to  
47 produce artefacts<sup>34</sup> leave this analysis open to interpretation.

48           The application of ONT sequencing to HEK293T cells nonetheless held conceptual  
49 merit. ONT reads can span germline and somatic retrotransposition events end-to-end, and  
50 resolve the sequence hallmarks of L1-mediated integration<sup>23,35</sup>. Through this approach, we  
51 previously found two somatic L1 insertions in the liver tumour sample of an individual  
52 positive for hepatitis C virus (HCV, a ~10kbp single-stranded non-polyadenylated RNA  
53 virus), including one PCR-validated L1 insertion spanned by a single ONT read<sup>23,36</sup>.  
54 HEK293T cells are arguably a favourable context to evaluate L1-mediated SARS-CoV-2  
55 genomic integration. They express L1 ORF1p<sup>37</sup>, readily accommodate engineered L1  
56 retrotransposition<sup>16,38,39</sup>, and support SARS-CoV-2 viral replication (**Extended Data Fig. 2**).  
57 Endogenous L1-mediated insertions can also be detected in cell culture by genomic analysis  
58 of separate cultivars derived from a common population<sup>40,41</sup>.

59           We therefore applied ONT sequencing (~54× genome-wide depth, read length N50 ~  
60 39kbp) to genomic DNA harvested from HEK293T cells infected with SARS-CoV-2 at a  
61 multiplicity of infection (MOI) of 1.0, as well as mock infected cells (~28× depth, N50 ~  
62 47kbp) (**Fig. 1a, Extended Data Fig. 2 and Supplementary Table 1**). As a positive control,  
63 we ONT sequenced the tumour and non-tumour liver tissue of a HBV-positive hepatocellular  
64 carcinoma patient<sup>36</sup>. To these data, we added those of Zhang *et al.*<sup>29</sup> and, as negative controls,  
65 the aforementioned HCV-positive hepatocellular carcinoma and normal liver samples<sup>23</sup>  
66 (**Supplementary Table 1**). We then used the Transposons from Long DNA Reads (TLDR)<sup>23</sup>  
67 software to call SARS-CoV-2, HBV, HCV and non-reference human-specific L1 (L1HS)  
68 insertions spanned by at least one uniquely aligned ONT read. TLDR detected no SARS-  
69 CoV-2, HBV or HCV insertions.

70           In total, TLDR identified 575 non-reference L1 insertions, which were typically  
71 flanked by TSDs with a median length of 14bp (**Fig. 1b and Supplementary Table 2**). No  
72 tumour-specific L1 insertions were found, apart from the two previously detected in the  
73 HCV-infected liver tumour<sup>23,36</sup>. Seventy-eight L1 insertions were found only in our SARS-  
74 CoV-2 infected HEK293T cells (66) or the mock infected control (12) and produced TSDs  
75 with a median length of 14bp (**Fig. 1c**). Of the 78 events, 69 (88.5%) were detected by a  
76 single spanning read and 13 carried a 3' transduction<sup>42,43</sup> (**Supplementary Table 2**). We  
77 chose at random 6/69 insertions detected by one spanning read for manual curation and PCR  
78 validation. All 6 insertions bore a TSD and a 3' polyA tract, and integrated at a degenerate L1  
79 EN motif (**Fig. 1d,e and Extended Data Fig. 3a-d**). Three were 5' inverted<sup>44,45</sup> (**Fig. 1e and**  
80 **Extended Data Fig. 3b,c**) and one carried a 3' transduction<sup>42</sup> traced to a mobile<sup>20</sup> full-length  
81 non-reference L1HS (**Extended Data Fig. 3b**). Two PCR amplified in the SARS-CoV-2 and



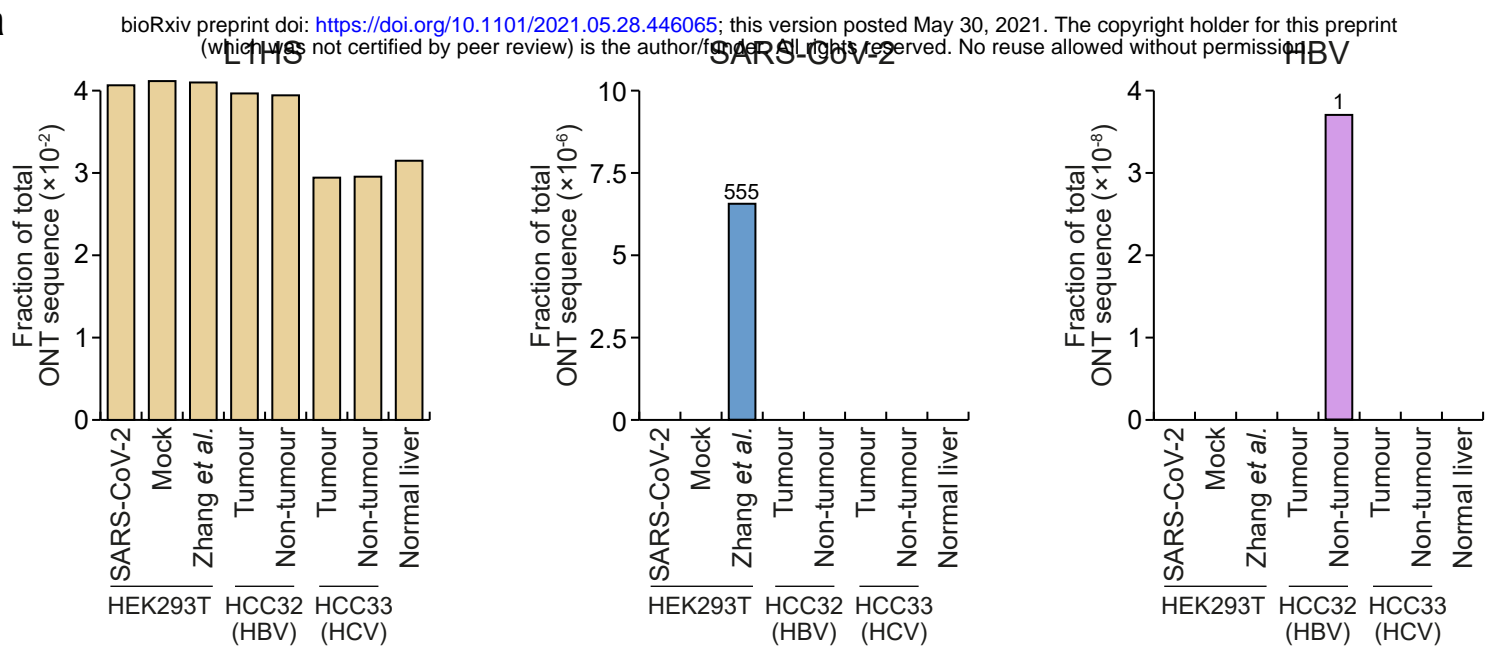
**Fig. 1: Detection of endogenous L1-mediated retrotransposition in human cells.** **a**, Experimental design. HEK-293T cells were divided into two populations (cultivars), which were then either SARS-CoV-2 infected or mock infected. DNA was extracted from each cultivar, as well as from hepatocellular carcinoma patient samples, and subjected to ONT sequencing. ONT reads were used to call non-reference L1 and virus insertions with TLDR, which also resolves TSDs and other retrotransposition hallmarks. TSDs: red triangles; polyA tract: green rectangle; ONT read: blue rectangle. Note: some illustrations are adapted from Ewing *et al.*<sup>23</sup>. **b**, TSD size distribution for non-reference L1 insertions, as annotated by TLDR. **c**, As for **b**, except showing data for L1 insertions found only in either our HEK293T cells infected with SARS-CoV-2 or our mock infected cells. **d**, Detailed characterisation of an L1 insertion detected in SARS-CoV-2 infected HEK293T cells by a single spanning ONT read aligned to chromosome 14. Nucleotides highlighted in red correspond to the integration site TSD. Underlined nucleotides correspond to the L1 EN motif. The cartoon indicates a full-length L1HS insertion flanked by TSDs (red triangles), and a 3' polyA tract (green). Numerals represent positions relative to the L1HS sequence L1.3<sup>46</sup>. The relevant spanning ONT read, with identifier, is positioned underneath the cartoon. Symbols ( $\alpha$ ,  $\beta$ ,  $\delta$ ,  $\gamma$ ) represent the approximate position of primers used for empty/filled site and L1-genome junction PCR validation reactions. Gel images display the results of these PCRs. Ladder band sizes are as indicated, NTC; non-template control. Red triangles indicate the expected size of L1 amplicons (empty triangle: no product observed; filled triangle: product observed). Blue triangles indicate expected empty site sizes. **e**, As for **d**, except for a 5' inverted/deleted L1HS located on chromosome 18.

82 mock infected samples (**Fig. 1d,e**) and four did not amplify in either sample (**Extended Data**  
83 **3a-d**). The 6 integration sites were on average spanned by 86 reads not containing the L1  
84 insertion (**Extended Data Fig. 3e**), a ratio (1:86) suggesting the L1s were absent from most  
85 cells. These and earlier<sup>23,35</sup> experiments show that lone spanning ONT reads can recover *bona*  
86 *fide* retrotransposition events, and highlight endogenous L1 activity in HEK293T cells  
87 lacking L1 overexpression systems.

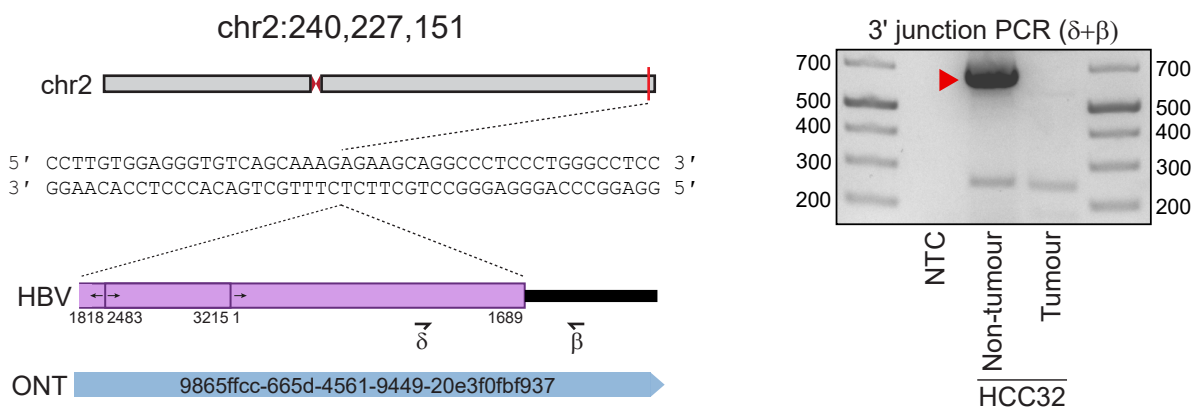
88 We next tested whether our computational analysis parameters excluded genuine  
89 HBV, HCV or SARS-CoV-2 insertions. We directly aligned our ONT reads to the genome of  
90 the SARS-CoV-2 isolate (QLD002, GISAID EPI\_ISL\_407896) used here, as well as to a  
91 geographically diverse set of HBV and HCV genomes (**Supplementary Table 1**), and a  
92 highly mobile L1HS sequence<sup>46</sup>. In total, 3.6% of our ONT sequence bases aligned to L1HS,  
93 whereas no alignments to the SARS-CoV-2 or HCV genomes were observed (**Fig. 2a**). One  
94 read from the HBV-infected non-tumour liver sample aligned to 2,770bp of a HBV genotype  
95 B isolate, and the remaining 2,901bp aligned to an intergenic region of chromosome 2 (**Fig.**  
96 **2b** and **Supplementary Table 2**). To validate this HBV insertion, we PCR amplified and  
97 capillary sequenced its 3' junction (**Fig. 2b**). The HBV sequence was linearised and  
98 rearranged (**Fig. 2b**) as per prior reports<sup>2-4</sup>. Direct inspection of ONT read alignments thus  
99 recovered a HBV integrant, which are found in ~1 per 10<sup>1</sup>-10<sup>4</sup> infected hepatocytes<sup>47-49</sup>, yet  
100 did not reveal reads alignable to the SARS-CoV-2 genome in our ONT datasets.

101 Reanalysing the ONT data generated by Zhang *et al.*, we found 555 reads (out of ~12  
102 million) that could be aligned to the SARS-CoV-2 genome (**Fig. 2a**), including one matching  
103 the aforementioned integrant on chromosome X that lacked a 3' polyA tract<sup>29</sup> (**Extended**  
104 **Data Fig. 1**). These reads (median length 924bp) were however 65.6% shorter than the  
105 overall dataset (2,686kbp) and were comprised of a much higher average proportion of  
106 SARS-CoV-2 sequence (52.3%) than the proportion of L1HS sequence found in reads  
107 aligned to L1HS (17.1%). An ONT read highlighted by Zhang *et al.* in support of a TSD-  
108 bearing SARS-CoV-2 insertion on chromosome 22 was also used to call a SARS-CoV-2  
109 insertion on chromosome 1 (**Extended Data Fig. 1**). Read alignment ambiguity to the  
110 genome resulted in TLDR calling neither the putative chromosome X or chromosome 22  
111 SARS-CoV-2 integrants. These analyses confirmed SARS-CoV-2 alignable reads were  
112 present in the Zhang *et al.* ONT dataset, yet these reads were unusually short and could  
113 include molecular artifacts interpreted by Zhang *et al.* as SARS-CoV-2 integrants.

114 In sum, we do not observe L1-mediated SARS-CoV-2 genomic integration in  
115 HEK293T cells, despite availability of the L1 machinery<sup>16,37-39</sup>. Our approach has several



**b**



**Fig. 2: ONT reads occasionally align to viral genome sequences.** **a**, Fractions of total ONT sequence alignable to L1HS (left), SARS-CoV-2 (middle) and HBV (right) isolate genomes. Read counts for SARS-CoV-2 and HBV are provided above histogram columns. No reads were aligned to the HCV isolate genomes. HEK293T data were generated here (SARS-CoV-2, mock) or by Zhang *et al.*<sup>29</sup>. HCC tumour/non-tumour liver pairs were sequenced here (HCC32; confirmed HBV-positive) or previously<sup>23</sup> (HCC33; HCV-positive). Normal liver ONT sequencing from our prior work<sup>23</sup> was included as an additional control. **b**, A HBV insertion detected in non-tumour liver. In this example, an ONT read from the non-tumour liver of HCC32 spanned the 3' junction of a HBV integrant located on chromosome 2. Of the HBV isolate genomes considered here (**Supplementary Table 1**), this read aligned best to a representative of genotype B (Genbank accession AB602818). The HBV sequence was rearranged consistent with its linearisation prior to integration<sup>2-4</sup>. Numerals indicate positions relative to AB602818. Symbols ( $\beta$ ,  $\delta$ ) represent the approximate position of primers used to PCR validate the HBV insertion. The gel image at right shows the PCR results. Ladder band sizes are as indicated. The red triangle indicates an on-target product.

116 notable differences and caveats when compared to that of Zhang *et al.*<sup>29</sup>. Each study used  
117 different SARS-CoV-2 isolates, and here the multiplicity of infection (MOI 1.0) was double  
118 that of Zhang *et al.* (MOI 0.5). The ONT library preparation kit and depth of sequencing  
119 applied to HEK293T cells by Zhang *et al.* (SQK-LSK109 kit, ~21× depth, N50 ~ 11kbp) and  
120 here (SQK-LSK110 kit, ~54× depth, N50 ~ 39kbp) differed. Zhang *et al.* applied ONT  
121 sequencing only to HEK293T cells transfected with an L1 expression plasmid, which human  
122 cells would not carry *in vivo*. We do not analyse SARS-CoV-2 patient samples although,  
123 arguably, HEK293T cells present an environment far more conducive to L1 activity than  
124 those cells accessed *in vivo* by SARS-CoV-2<sup>50,51</sup>. Widespread cell death post-infection also  
125 reduces the probability SARS-CoV-2 integrants would persist in the body<sup>52,53</sup>. Finally, the  
126 incredible enrichment reported by Zhang *et al.* for putative SARS-CoV-2 insertions in exons,  
127 which the L1 EN does not prefer<sup>32,33</sup>, contradicts the involvement of L1. We conclude L1 *cis*  
128 preference likely disfavours SARS-CoV-2 retrotransposition, making the phenomenon  
129 mechanistically plausible but likely very rare, as for other polyadenylated cellular RNAs<sup>6-11</sup>.

130

### 131 **Acknowledgements**

132 The authors thank the human subjects of this study who donated tissues to the Centre  
133 Hépatobiliaire, Paul-Brousse Hospital. We thank S. Richardson and R. Shukla for helpful  
134 discussions, K. Chappell and P. Young for project support, and Queensland Health for  
135 providing the SARS-CoV-2 virus isolate QLD02. This study was funded by an Australian  
136 Government Research Training Program (RTP) Scholarships (N.S.), an NHMRC-ARC  
137 Dementia Research Development Fellowship (GNT1108258, G.O.B.), seed funding provided  
138 by the Australian Infectious Disease Research Centre to establish SARS-CoV-2 research at  
139 the University of Queensland (A.A.K), an Australian Department of Health Medical Research  
140 Future Fund (MRFF) Novel Coronavirus Vaccine Development Grant (APP1202445-2020,  
141 D.W.), an MRFF Investigator Grant (MRF1175457, A.D.E.), an NHMRC Investigator Grant  
142 (GNT1173711, G.J.F.), a CSL Centenary Fellowship (G.J.F.), and the Mater Foundation.

143

### 144 **Author contributions**

145 N.S., J.R., G.O.B., A.A.A., P.G., F.J.S-L., P.A., N.M. and B.L. performed experiments and  
146 analysed data. A.D.E. and G.J.F. performed bioinformatic analysis. J.F., I.W.D., A.A.K.,  
147 D.W. and G.J.F. provided resources. G.J.F. designed the project and wrote the manuscript.

148

### 149 **Competing interests**



150 The authors declare no competing interests.

151

## 152 **References**

- 153 1. Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China.  
154 *Nature* **579**, 265–269 (2020).
- 155 2. Fujimoto, A. *et al.* Whole-genome sequencing of liver cancers identifies etiological  
156 influences on mutation patterns and recurrent mutations in chromatin regulators. *Nat.*  
157 *Genet.* **44**, 760–764 (2012).
- 158 3. Nagaya, T. *et al.* The mode of hepatitis B virus DNA integration in chromosomes of  
159 human hepatocellular carcinoma. *Genes Dev.* **1**, 773–782 (1987).
- 160 4. Jiang, Z. *et al.* The effects of hepatitis B virus integration into the genomes of  
161 hepatocellular carcinoma patients. *Genome Res.* **22**, 593–601 (2012).
- 162 5. Kazazian, H. H., Jr & Moran, J. V. Mobile DNA in Health and Disease. *N. Engl. J. Med.*  
163 **377**, 361–370 (2017).
- 164 6. Moran, J. V. *et al.* High frequency retrotransposition in cultured mammalian cells. *Cell*  
165 **87**, 917–927 (1996).
- 166 7. Kulpa, D. A. & Moran, J. V. Cis-preferential LINE-1 reverse transcriptase activity in  
167 ribonucleoprotein particles. *Nat. Struct. Mol. Biol.* **13**, 655–660 (2006).
- 168 8. Doucet, A. J., Wilusz, J. E., Miyoshi, T., Liu, Y. & Moran, J. V. A 3' Poly(A) Tract Is  
169 Required for LINE-1 Retrotransposition. *Mol. Cell* **60**, 728–741 (2015).
- 170 9. Monot, C. *et al.* The specificity and flexibility of 11 reverse transcription priming at  
171 imperfect T-tracts. *PLoS Genet.* **9**, e1003499 (2013).
- 172 10. Wei, W. *et al.* Human L1 retrotransposition: cispreference versus trans  
173 complementation. *Mol. Cell. Biol.* **21**, 1429–1439 (2001).
- 174 11. Garcia-Perez, J. L., Doucet, A. J., Bucheton, A., Moran, J. V. & Gilbert, N. Distinct  
175 mechanisms for trans-mediated mobilization of cellular RNAs by the LINE-1 reverse  
176 transcriptase. *Genome Res.* **17**, 602–611 (2007).
- 177 12. Dewannieux, M., Esnault, C. & Heidmann, T. LINE-mediated retrotransposition of  
178 marked Alu sequences. *Nat. Genet.* **35**, 41–48 (2003).
- 179 13. Esnault, C., Maestre, J. & Heidmann, T. Human LINE retrotransposons generate  
180 processed pseudogenes. *Nat. Genet.* **24**, 363–367 (2000).
- 181 14. Hancks, D. C., Ewing, A. D., Chen, J. E., Tokunaga, K. & Kazazian, H. H., Jr. Exon-  
182 trapping mediated by the human retrotransposon SVA. *Genome Res.* **19**, 1983–1991

- 183 (2009).
- 184 15. Evrony, G. D. *et al.* Cell lineage analysis in human brain using endogenous  
185 retroelements. *Neuron* **85**, 49–59 (2015).
- 186 16. Sanchez-Luque, F. J. *et al.* LINE-1 Evasion of Epigenetic Repression in Humans. *Mol.*  
187 *Cell* **75**, 590–604 (2019).
- 188 17. Feusier, J. *et al.* Pedigree-based estimation of human mobile element retrotransposition  
189 rates. *Genome Res.* **29**, 1567–1577 (2019).
- 190 18. Scott, E. C. *et al.* A hot L1 retrotransposon evades somatic repression and initiates  
191 human colorectal cancer. *Genome Res.* **26**, 745–755 (2016).
- 192 19. Schauer, S. N. *et al.* L1 retrotransposition is a common feature of mammalian  
193 hepatocarcinogenesis. *Genome Res.* **28**, 639–653 (2018).
- 194 20. Rodriguez-Martin, B. *et al.* Pan-cancer analysis of whole genomes identifies driver  
195 rearrangements promoted by LINE-1 retrotransposition. *Nat. Genet.* **52**, 306–319  
196 (2020).
- 197 21. Deniz, Ö., Frost, J. M. & Branco, M. R. Regulation of transposable elements by DNA  
198 modifications. *Nat. Rev. Genet.* **20**, 417–431 (2019).
- 199 22. Ahl, V., Keller, H., Schmidt, S. & Weichenrieder, O. Retrotransposition and Crystal  
200 Structure of an Alu RNP in the Ribosome-Stalling Conformation. *Mol. Cell* **60**, 715–727  
201 (2015).
- 202 23. Ewing, A. D. *et al.* Nanopore Sequencing Enables Comprehensive Transposable  
203 Element Epigenomic Profiling. *Mol. Cell* **80**, 915–928 (2020).
- 204 24. Jurka, J. Sequence patterns indicate an enzymatic involvement in integration of  
205 mammalian retroposons. *Proc. Natl. Acad. Sci. U. S. A.* **94**, 1872–1877 (1997).
- 206 25. Luan, D. D., Korman, M. H., Jakubczak, J. L. & Eickbush, T. H. Reverse transcription  
207 of R2Bm RNA is primed by a nick at the chromosomal target site: a mechanism for non-  
208 LTR retrotransposition. *Cell* **72**, 595–605 (1993).
- 209 26. Gilbert, N., Lutz, S., Morrish, T. A. & Moran, J. V. Multiple fates of L1  
210 retrotransposition intermediates in cultured human cells. *Mol. Cell. Biol.* **25**, 7780–7795  
211 (2005).
- 212 27. Raiz, J. *et al.* The non-autonomous retrotransposon SVA is trans-mobilized by the  
213 human LINE-1 protein machinery. *Nucleic Acids Res.* **40**, 1666–1683 (2012).
- 214 28. Faulkner, G. J. & Billon, V. L1 retrotransposition in the soma: a field jumping ahead.  
215 *Mob. DNA* **9**, 22 (2018).
- 216 29. Zhang, L. *et al.* Reverse-transcribed SARS-CoV-2 RNA can integrate into the genome

- 217 of cultured human cells and can be expressed in patient-derived tissues. *Proc. Natl.*  
218 *Acad. Sci. U. S. A.* **118**, (2021).
- 219 30. Kazachenka, A. & Kassiotis, G. SARS-CoV-2-host chimeric RNA-sequencing reads do  
220 not necessarily signify virus integration into the host DNA. *bioRxiv* 2021.03.05.434119  
221 (2021) doi:10.1101/2021.03.05.434119.
- 222 31. Yan, B. *et al.* Host-virus chimeric events in SARS-CoV2 infected cells are infrequent  
223 and artifactual. *J. Virol.* (2021) doi:10.1128/JVI.00294-21.
- 224 32. Sultana, T. *et al.* The Landscape of L1 Retrotransposons in the Human Genome Is  
225 Shaped by Pre-insertion Sequence Biases and Post-insertion Selection. *Mol. Cell* **74**,  
226 555–570 (2019).
- 227 33. Flasch, D. A. *et al.* Genome-wide de novo L1 Retrotransposition Connects  
228 Endonuclease Activity with Replication. *Cell* **177**, 837–851.e28 (2019).
- 229 34. Treiber, C. D. & Waddell, S. Resolving the prevalence of somatic transposition in  
230 *Drosophila*. *Elife* **6**, (2017).
- 231 35. Siudeja, K. *et al.* Unraveling the features of somatic transposition in the *Drosophila*  
232 intestine. *EMBO J.* **40**, e106388 (2021).
- 233 36. Shukla, R. *et al.* Endogenous retrotransposition activates oncogenic pathways in  
234 hepatocellular carcinoma. *Cell* **153**, 101–111 (2013).
- 235 37. Philippe, C. *et al.* Activation of individual L1 retrotransposon instances is restricted to  
236 cell-type dependent permissive loci. *Elife* **5**, (2016).
- 237 38. Kubo, S. *et al.* L1 retrotransposition in nondividing and primary human somatic cells.  
238 *Proc. Natl. Acad. Sci. U. S. A.* **103**, 8036–8041 (2006).
- 239 39. Niewiadomska, A. M. *et al.* Differential inhibition of long interspersed element 1 by  
240 APOBEC3 does not correlate with high-molecular-mass-complex formation or P-body  
241 association. *J. Virol.* **81**, 9577–9583 (2007).
- 242 40. Nguyen, T. H. M. *et al.* L1 Retrotransposon Heterogeneity in Ovarian Tumor Cell  
243 Evolution. *Cell Rep.* **23**, 3730–3740 (2018).
- 244 41. Klawitter, S. *et al.* Reprogramming triggers endogenous L1 and Alu retrotransposition in  
245 human induced pluripotent stem cells. *Nat. Commun.* **7**, 10286 (2016).
- 246 42. Holmes, S. E., Dombroski, B. A., Krebs, C. M., Boehm, C. D. & Kazazian, H. H., Jr. A  
247 new retrotransposable human L1 element from the LRE2 locus on chromosome 1q  
248 produces a chimaeric insertion. *Nat. Genet.* **7**, 143–148 (1994).
- 249 43. Moran, J. V., DeBerardinis, R. J. & Kazazian, H. H., Jr. Exon shuffling by L1  
250 retrotransposition. *Science* **283**, 1530–1534 (1999).

- 251 44. Ostertag, E. M. & Kazazian, H. H., Jr. Twin priming: a proposed mechanism for the  
252 creation of inversions in L1 retrotransposition. *Genome Res.* **11**, 2059–2065 (2001).
- 253 45. Kazazian, H. H., Jr *et al.* Haemophilia A resulting from de novo insertion of L1  
254 sequences represents a novel mechanism for mutation in man. *Nature* **332**, 164–166  
255 (1988).
- 256 46. Dombroski, B. A., Scott, A. F. & Kazazian, H. H., Jr. Two additional potential  
257 retrotransposons isolated from a human L1 subfamily that contains an active  
258 retrotransposable element. *Proc. Natl. Acad. Sci. U. S. A.* **90**, 6513–6517 (1993).
- 259 47. Mason, W. S. *et al.* HBV DNA Integration and Clonal Hepatocyte Expansion in Chronic  
260 Hepatitis B Patients Considered Immune Tolerant. *Gastroenterology* **151**, 986–998.e4  
261 (2016).
- 262 48. Rydell, G. E. *et al.* Abundance of non-circular intrahepatic hepatitis B virus DNA may  
263 reflect frequent integration into human DNA in chronically infected patients. *J. Infect.*  
264 *Dis.* (2020) doi:10.1093/infdis/jiaa572.
- 265 49. Tu, T., Budzinska, M. A., Vondran, F. W. R., Shackel, N. A. & Urban, S. Hepatitis B  
266 Virus DNA Integration Occurs Early in the Viral Life Cycle in an In Vitro Infection  
267 Model via Sodium Taurocholate Cotransporting Polypeptide-Dependent Uptake of  
268 Enveloped Virus Particles. *J. Virol.* **92**, (2018).
- 269 50. Sungnak, W. *et al.* SARS-CoV-2 entry factors are highly expressed in nasal epithelial  
270 cells together with innate immune genes. *Nat. Med.* **26**, 681–687 (2020).
- 271 51. Wiersinga, W. J., Rhodes, A., Cheng, A. C., Peacock, S. J. & Prescott, H. C.  
272 Pathophysiology, Transmission, Diagnosis, and Treatment of Coronavirus Disease 2019  
273 (COVID-19): A Review. *JAMA* **324**, 782–793 (2020).
- 274 52. Karki, R. *et al.* Synergism of TNF- $\alpha$  and IFN- $\gamma$  Triggers Inflammatory Cell Death,  
275 Tissue Damage, and Mortality in SARS-CoV-2 Infection and Cytokine Shock  
276 Syndromes. *Cell* **184**, 149–168.e17 (2021).
- 277 53. Varga, Z. *et al.* Endothelial cell infection and endotheliitis in COVID-19. *Lancet* **395**,  
278 1417–1418 (2020).
- 279 54. Amarilla, A. A. *et al.* An optimized high-throughput immuno-plaque assay for SARS-  
280 CoV-2. *Front. Microbiol.* **12**, 625136 (2021).
- 281 55. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,  
282 3094–3100 (2018).
- 283 56. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,  
284 2078–2079 (2009).

285 57. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).

286 58. Untergasser, A., Cutcutache, I. & Koressaar, T. Primer3—new capabilities and  
287 interfaces. *Nucleic acids* (2012).

## 288 **Methods**

### 289 ***SARS-CoV-2 infection of HEK293T cells***

290 HEK293T cells and African green monkey kidney cells (Vero E6) were maintained in standard  
291 Dulbecco's Modified Eagle Medium (DMEM). Culture media were supplemented with sodium  
292 pyruvate (11 mg/L), penicillin (100 U/mL), streptomycin (100 µg/mL) (P/S) and 10% foetal calf  
293 serum (FCS) (Bovogen, USA). Cells were maintained at 37 °C with 5% CO<sub>2</sub>.

294 An early Australian SARS-CoV-2 isolate (hCoV-19/Australia/QLD02/2020; GISAID  
295 Accession EPI\_ISL\_407896) was sampled from patient nasopharyngeal aspirates by  
296 Queensland Health Forensic and Scientific Services and used to inoculate Vero E6 African  
297 green monkey kidney cells (passage 2). A viral stock (passage 3) was then generated on Vero  
298 E6 cells and stored at -80°C. Viral titration was determined by immuno-plaque assay (iPA), as  
299 previously described<sup>54</sup>. To verify viral replication in HEK293T cells, a growth kinetic was  
300 assessed using a multiplicity of infection (MOI) of 1.0, and showed efficient SARS-CoV-2  
301 replication (**Extended Data Fig. 2**).

302 HEK293T viral infection was undertaken as follows: 3×10<sup>6</sup> HEK293T cells were  
303 seeded onto 6-well plates pre-coated with polylysine one day before infection. Cells were  
304 infected at MOI of 1 in 200 µL of DMEM (2% FCS and P/S) and incubated for 30 min at 37°C.  
305 Plates were rocked every 5 min to ensure the monolayer remained covered with inoculum. The  
306 inoculum was then removed, and the monolayer washed five times with 1 mL of additive-free  
307 DMEM. Finally, cells were maintained with 3 mL of DMEM (supplemented with 2% foetal  
308 bovine serum and P/S) and incubated at 37°C with 5% CO<sub>2</sub>. Cell supernatant was harvested 0,  
309 1, 2 and 3 days post-infection. The mock infected control differed only in that virus was not  
310 added to the inoculum media.

311 Genomic DNA was extracted from mock and SARS-CoV-2 infected (MOI 1.0)  
312 HEK293T cells sampled 2 days post-infection, using a Nanobind CBB Big DNA Kit  
313 (Circulomics) following the manufacturer's instructions for high molecular weight (HMW)  
314 DNA extraction. DNA was eluted in elution buffer (10 mM Tris-Cl, pH 8.5) and concentration  
315 measured by Qubit dsDNA High-Sensitivity Assay Kit on a Qubit Fluorometer (Life  
316 Technologies).

317

### 318 *Hepatocellular carcinoma patient samples*

319 Liver tumour and non-tumour tissue were previously obtained from a HBV-positive patient  
320 (HCC32, male, 73yrs) who underwent surgical resection at the Centre Hepatobiliaire, Paul-  
321 Brousse Hospital, and made available for research purposes with approval from the French  
322 Institute of Medical Research and Health (Reference: 11-047). Further ethics approvals were  
323 provided by the Mater Health Services Human Research Ethics Committee (Reference: HREC-  
324 15-MHS-52) and the University of Queensland Medical Research Review Committee  
325 (Reference: 2014000221). DNA was extracted from the HCC32 tissues in our earlier study<sup>36</sup>  
326 with a DNeasy Blood and Tissue Kit (QIAGEN, Germany) and stored at -80°C. To enrich for  
327 HMW DNA, 4.5µg of DNA from the patient HCC32 tumour and non-tumour liver samples  
328 was diluted to 75ng/µL in a 1.5mL Eppendorf DNA LoBind tube and processed with a Short  
329 Read Eliminator XS Kit (Circulomics) following the manufacturer's instructions.

330

### 331 *ONT sequencing*

332 DNA libraries were prepared at the Kinghorn Centre for Clinical Genomics (KCCG) using 3-  
333 4µg HMW input DNA, without shearing, and a SQK-LSK110 ligation sequencing kit. 350-  
334 500ng of each prepared library was sequenced separately on one PromethION (Oxford  
335 Nanopore Technologies) flow cell (FLO-PRO002, R9.4.1 chemistry) (**Supplementary Table**  
336 **1**). SARS-CoV-2 infected HEK293T DNA was sequenced on two flow cells. Flow cells were  
337 washed (nuclease flush) and reloaded at 24hr and 48hr with 350-500ng of additional library to  
338 maximise output. Bases were called with guppy 4.0.11 (Oxford Nanopore Technologies).  
339 Sequencing data were deposited in the Sequence Read Archive (SRA) under project  
340 PRJEB44816.

341

### 342 *ONT bioinformatic analyses*

343 To call non-reference insertions with TLDR<sup>23</sup>, ONT reads generated here, by Zhang *et al.*<sup>29</sup>,  
344 and by our previous ONT study of human tissues<sup>23</sup> (**Supplementary Table 1**) were aligned  
345 to the human reference genome build hg38 using minimap2<sup>55</sup> version 2.17 (index parameter:  
346 -x map-ont; alignment parameters: -ax map-ont -L -t 32) and samtools<sup>56</sup> version 1.12. BAM  
347 files were then processed as a group with TLDR<sup>23</sup> version 1.1 (parameters -e virus.fa -p 128 -  
348 m 1 --max\_te\_len 40000 --max\_cluster\_size 100 --min\_te\_len 100 --wiggle 100 --  
349 keep\_pickles -n nonref.collection.hg38.chr.bed.gz). The file virus.fa was composed of:  
350 representative HBV and HCV isolate genomes (**Supplementary Table 1**), the SARS-CoV-2  
351 isolate used here (GISAID Accession EPI\_ISL\_407896) and the mobile L1HS sequence

352 L1.3<sup>46</sup> (Genbank Accession L19088). The file nonref.collection.hg38.chr.bed.gz is a  
353 collection of known non-reference retrotransposon insertions available from  
354 github.com/adamewing/tldr/. The TLDR output table was further processed to remove calls  
355 not passing all TLDR filters, representing homopolymer insertions, where MedianMapQ < 50  
356 or family = “NA” or remappable = “FALSE” or UnmapCover < 0.75 or LengthIns < 100 or  
357 EndTE-StartTE < 100 or strand = “None” or SpanReads < 1 or L1HS insertions where  
358 EndTE < 6017. The filtered TLDR output table is provided as **Supplementary Table 2**.  
359 L1HS insertions detected in only our mock or SARS-CoV-2 infected HEK293T datasets, but  
360 not in both experiments, and not matching a known non-reference L1HS element, were  
361 designated as putative cultivar-specific insertions (**Supplementary Table 2**). Many if not  
362 most of these insertions were likely to have occurred in cell culture prior to the cultivars  
363 being separated.

364 To identify L1HS and viral sequences, we directly aligned all reads to the virus.fa file  
365 with minimap2 (index parameter: -x map-ont; alignment parameters: -ax map-ont -L -t 32).  
366 Reads containing alignments of  $\geq 100$ bp to a sequence present in virus.fa were counted with  
367 samtools idxstats. Alignments to HBV, HCV or SARS-CoV-2 were excluded if they  
368 overlapped with a genomic alignment of  $\geq 100$ bp. Read alignments were visualised with  
369 samtools view and the Integrative Genomics Viewer<sup>57</sup> version 2.8.6.

370

### 371 ***PCR validation***

372 We used Primer3<sup>58</sup> to design PCR primers for 6 L1 insertions found by a single spanning  
373 ONT read, using the reference genome and L1HS sequences as inputs (**Supplementary**  
374 **Table 2**). These validation experiments were conducted in three phases. Firstly, we  
375 performed an “empty/filled site” PCR using primers positioned on either side of the L1,  
376 where the filled site is the L1 allele, and the empty site is the remaining allele(s). Each  
377 empty/filled reaction was performed using a DNA Engine Tetrad 2 Thermal Cycler (Bio-  
378 Rad) and Expand Long Range Enzyme Mix, with 1X Expand Long Range Buffer with  
379 MgCl<sub>2</sub>, 50pmol of each primer, 0.5mM dNTPs, 5% DMSO, 100ng of template DNA and  
380 1.75U of enzyme, in a 25 $\mu$ L final volume. PCR cycling conditions were as follows: (92°C,  
381 3min) $\times$ 1; (92°C, 30sec; 54-57°C, 30sec; 68°C, 7min) $\times$ 10; (92°C, 30sec; 52-55°C, 30sec;  
382 68°C, 7min + 20sec/cycle) $\times$ 30; (68°C, 10min; 4°C, hold) $\times$ 1. Amplicons were visualised on a  
383 1% agarose gel stained with SYBR SAFE (Invitrogen). GeneRuler<sup>TM</sup> 1kb plus (Thermo  
384 Scientific) was used as the ladder. Secondly, we combined each empty/filled primer with a  
385 primer positioned within the L1 sequence, to amplify the 5' and 3' L1-genome junctions.

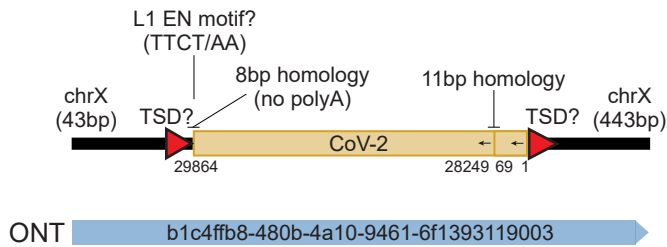
386 These reactions were conducted on a T100 Thermal Cycler (Bio-Rad), with MyTaq HS DNA  
387 polymerase, 1X MyTaq Reaction Buffer, 10pmol of each primer, 10ng of template DNA, and  
388 2.5U of enzyme, in a 25 $\mu$ L final volume. PCR cycling conditions were as follows: (95°C,  
389 1min) $\times$ 1; (95°C, 15sec; 53-55°C, 15sec; 72°C, 15sec) $\times$ 35; (72°C, 5min; 4°C, hold) $\times$ 1.  
390 Amplicons were visualised on a 1.5% agarose gel stained with SYBR SAFE (Invitrogen).  
391 Thirdly, we repeated the 5' L1-genome junction-specific PCR using 200ng template DNA.  
392 All PCRs were performed with non-template control, as well as DNA extracted from the  
393 same HEK293T cells (SARS-CoV-2 and mock) subjected to genomic analysis. Notably, L1  
394 insertions that did not amplify in either cultivar were still likely to be genuine events as they  
395 carried all of the relevant sequence hallmarks of L1-mediated retrotransposition.

396 PCR primers for the HBV insertion 3' junction (**Fig. 2b** and **Supplementary Table 2**)  
397 were designed with Primer3 using the reference genome and closest match HBV sequence  
398 (Genbank accession AB602818) as inputs. PCR amplification and capillary sequencing was  
399 conducted as per the L1 insertions, except using Expand Long Range polymerase (Roche)  
400 with 1X Expand Long Range buffer with MgCl<sub>2</sub>, 10pmol of each primer, 100ng of template  
401 DNA, 500 $\mu$ M of PCR Nucleotide Mix, and 3.5U of enzyme, in a 25 $\mu$ L final volume. PCR  
402 cycling conditions were as follows: (92°C, 2min) $\times$ 1; (92°C, 15sec; 65°C, 15sec; 68°C,  
403 7:30min) $\times$ 10; (92°C, 15sec; 65°C, 15sec; 68°C, 7min+ 20sec per cycle) $\times$ 35 (68°C, 10min;  
404 4°C, hold) $\times$ 1. Amplicons were visualized on a 1.2% agarose gel.

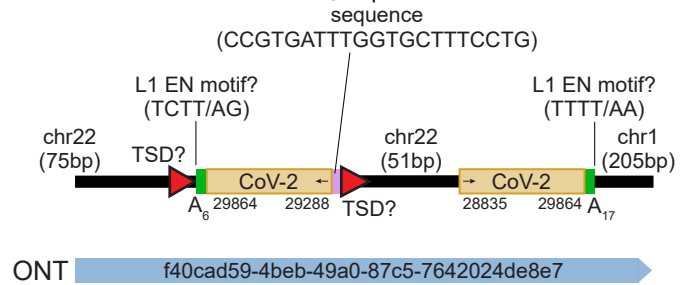
405 Amplicons in each experiment were visualised using a GelDoc (Bio-Rad) and, if of  
406 the correct size, gel extracted using a Qiagen MinElute Gel Extraction Kit and capillary  
407 sequenced by the Australian Genomics Research Facility (Brisbane).



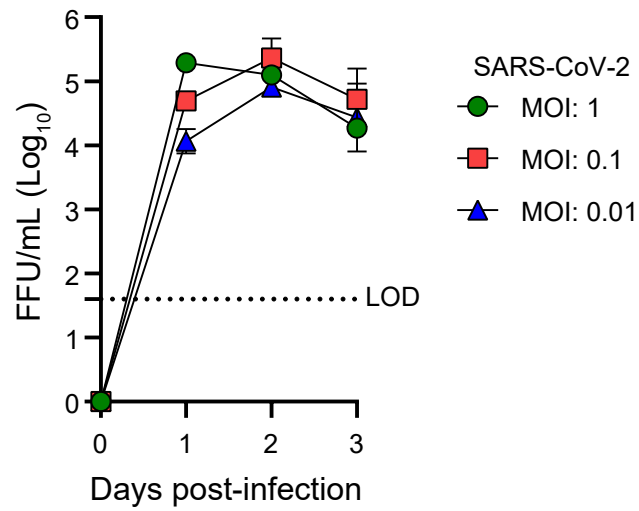
**a**



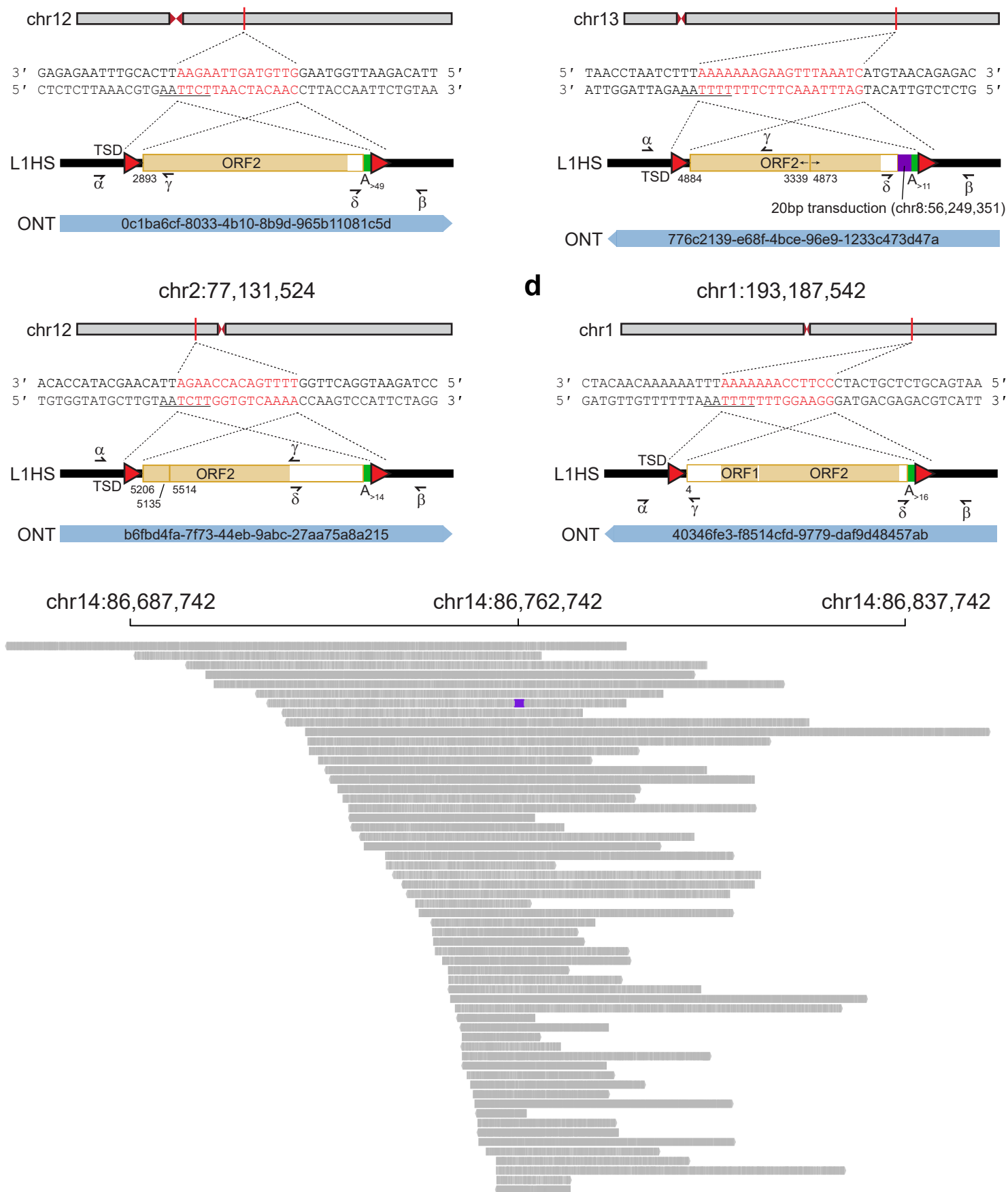
**b**



**Extended Data Fig. 1: Key SARS-CoV-2 insertions reported by Zhang *et al.*** **a**, A cartoon summarising the features of a putative SARS-CoV-2 integrant on chromosome X. Numerals underneath the SARS-CoV-2 sequence represent positions relative to the QLD02 virus isolate. Potential TSDs are shown as red triangles. No 3' polyA tract was found. Homologous regions at sequence junctions are marked. One spanning ONT read is positioned underneath the cartoon and its identifier is displayed. **b**, As for **a**, except showing an ONT read spanning two SARS-CoV-2 insertions, on chromosome 22 and chromosome 1. The alignments to chromosome 22 are flagged as supplementary by the minimap2 aligner. 3' polyA tracts are represented as green rectangles. Note: the chromosome 22 and chromosome X instances are the key examples reported by Zhang *et al.* in support of SARS-CoV-2 genomic integration. Neither example has a complete set of retrotransposition hallmarks (TSD, 3' polyA tract, L1 EN motif) *and* the support of a uniquely aligned ONT read.



**Extended Data Fig. 2: SARS-CoV-2 is replication competent in HEK293T cells.** HEK293T cells were infected with SARS-CoV-2 isolate QLD02 at an MOI of 0.01, 0.1 and 1.0. Inoculum was removed after infection and cells were washed before the addition of growth media. Supernatant was collected at the indicated time points and viral titres were quantified as focus-forming units (FFU) per mL by immuno-plaque assay (iPA)<sup>54</sup> with a limit of detection (LOD) as indicated.



**Extended Data Fig. 3: Additional HEK293T cell L1HS insertions.** **a**, A 5' truncated L1. **b**, A 5' inverted/deleted L1 carrying a 3' transduction (purple rectangle) traced to a known non-reference L1 present in HEK293T cells. **c**, A 5' inverted/deleted L1. **d**, A near full-length L1. Each panel shows the genomic coordinates of an L1 insertion, as well as the sequence at the insertion site. Nucleotides highlighted in red correspond to the integration site TSD. Underlined nucleotides correspond to the L1 EN motif. Cartoons summarise the features of each L1, with numerals representing positions relative to L1.3<sup>46</sup>, TSDs shown as red triangles, and 3' polyA tracts coloured as green rectangles. One spanning ONT read with its identifier is positioned underneath each cartoon. Symbols ( $\alpha$ ,  $\beta$ ,  $\delta$ ,  $\gamma$ ) represent the approximate position of primers used for empty/filled and L1-genome junction PCR validation reactions. No L1 amplicons were recovered by these assays. **e**. Integrative Genomics Viewer<sup>57</sup> visualisation of read alignments spanning the L1 integration site displayed in Fig. 1d. The L1 is coloured purple.