

1 **METHODS**

2 **Cell culture and virus infection**

3 HEK293T (Cell Resource Center, Institute of Basic Medicine Chinese Academy of
4 Medical Sciences, 3111C0001CCC000091), Huh-7 (Cell Resource Center, Institute of
5 Basic Medicine Chinese Academy of Medical Sciences, 3111C0001CCC000679),
6 Calu-3 (Cell Resource Center, Institute of Basic Medicine Chinese Academy of
7 Medical Sciences, 3111C0001CCC000032), Cos-7 (Kunming Cell Bank, Chinese
8 Academy of Sciences, 3153C0001000000037) and MA-104 (Chinese Typical Culture
9 Preservation Center Cell Bank, 3142C0001000000041) cells were cultured in DMEM
10 (ThermoFisher,0030034DJ) supplemented with 10% fetal bovine serum (Gibco,
11 10099-141C) and penicillin-streptomycin (ThermoFisher,15140122) at 37°C, 5% CO₂.
12 Cells were infected with SARS-CoV-2 (CDC of Guangdong province, GD108#), at a
13 multiplicity of infection (MOI) of 0.1, and were collected for RNA isolation 24 hours
14 post virus infection. All experiments with the SARS-CoV-2 virus were performed in
15 the BSL-4 laboratory.

16

17 **RNA purification, library construction and sequencing**

18 Total RNA from SARS-CoV-2 infected cell samples was extracted with Trizol reagent
19 (Invitrogen, 15596026), and then subjected to the rRNA depletion with the Ribo-off
20 rRNA Depletion Kit (Human/Mouse/Rat)(Vazyme, N406-02) following the
21 manufacturer's instructions. Libraries were constructed using the KAPA RNA
22 HyperPrep Kit (KAPA Biosystems, KK8541) following the manufacturer's instructions.
23 Sequencing was performed on Illumina NovaSeq 6000 system with paired end 150 bp
24 read length. Meanwhile, mixed sample of SARS-CoV-2 infected Huh-7 cells with
25 zebrafish embryonic RNA was also used to prepare RNA-seq libraries, serving as the
26 artificial chimeric RNA-seq reads control. For each infected sample, two replicates
27 were performed, while three replicates were performed for the mixed sample.

28

29 **Genomic DNA isolation, library construction and whole genome sequencing**

30 DNA samples were isolated by Universal Genomic DNA Kit (CW BIO, CW2298)
31 according to the manufacturer's instructions. The optical density values at 260/280
32 were approximately 1.6~1.8. Genomic DNA of the same SARS-CoV-2 infected cell
33 lines was prepared and whole genome shotgun libraries were constructed by using

34 TruePrep DNA library Prep kit V2 for Illumina (Vazyme, Cat No. TD501-TD503),
35 followed by sequencing on Illumina NovaSeq 6000 with paired end 150 bp read length.
36 For each infected sample, two replicates were performed. About 100 Gb data were
37 obtained for each replicate (Table S1).

38

39 **Bioinformatics analysis of RNA-seq data**

40 Raw FASTQ reads was processed to filter low quality bases and cut adapter sequence
41 by fastp (version 0.20.1)(Chen et al., 2018). The clean RNA-seq reads from SARS-
42 CoV-2 infected cell were aligned to host genome appending with SARS-CoV-2 genome
43 (NC_045512) by using STAR (version 2.7.7a, parameters as followed: --chimOutType
44 Junctions SeparateSAMold WithinBAM HardClip --chimSegmentMin 50 --
45 chimScoreJunctionNonGTAG 0 --alignSJstitchMismatchNmax -1 -1 -1 -1 --
46 chimJunctionOverhangMin 50 --outSAMtype BAM SortedByCoordinate --quantMode
47 TranscriptomeSAM GeneCounts)(Dobin et al., 2013). The versions of genomes include
48 GRCh38 for human (annotation release 102), Vero_WHO_p1.0 for green monkey
49 (NCBI Chlorocebus sabaues annotation release 102) and Mmul_10 for rhesus monkey
50 (NCBI Macaca mulatta annotation release 103), and GRCz11 for zebrafish (NCBI:
51 GCA_000002035.4). The duplicates were discarded. For RNA-seq reads from the
52 mixed libraries of SARS-CoV-2 infected cells and uninfected zebrafish embryos were
53 also aligned to human genome, zebrafish genome and SARS-CoV-2 genome with the
54 same pipeline. The chimeric reads were also gained between each two genomes.

55

56 **Bioinformatics analysis of whole genome sequencing data**

57 Reads from whole genome sequencing data were mapped to human genome plus
58 SARS-CoV-2 genome by using BWA (version 0.7.17)(Li and Durbin, 2009), and
59 genome coverage statistics was obtained by samtools (version 1.12)(Li et al., 2009).

60

61 **Gene expression analysis**

62 Gene expression level were determined by RPKM (Reads Per Kilobase per Million
63 mapped reads) and reads counts of genes were calculated by cufflinks (Trapnell et al.,
64 2012). The expressed genes were then filtered by $RPKM \geq 1$. To note, in Fig. 4c, the
65 expression levels were defined as CPM without normalization of transcript length to
66 represent the fragment depth of each transcript, directly. The p value for difference of

67 expressions between chimeric genes and non-chimeric genes were calculated by Mann-
68 Whitney U test (Table S2, S4 and S6).

69

70 **Chimeric gene analysis**

71 Chimeric reads between host and SARS-CoV-2 were obtained via the above STAR
72 analysis pipeline. For these reads, the fragments from host genome were annotated by
73 corresponding genome annotation reference by bedtools intersectBed (Quinlan and
74 Hall, 2010). The chimeric level was then defined as CPM (Counts Per Million reads),
75 which were calculated by all chimeric reads from same gene and normalized by viral
76 reads depth but not sequencing depth, except for results in Extended Data Fig. 1. The
77 chimeric genes were then filtered by $CPM \geq 1$ (Table S3, S5 and S6).

78 The chimeric events were defined by loci of junctions between host and SARS-CoV-2.
79 Chimeric events at different resolutions for (10 nt, 50 nt and 100 nt) were re-defined by
80 corresponding steps, which were obtained by sliding windows at various resolution
81 along host and viral genome.

82

83 **Analysis for accumulated expression and chimeric levels**

84 The accumulated expression and chimeric levels were defined by sum of genes in
85 related bins. Genes were first sorted in order of loci along each chromosome and then
86 separated into one bin for each 100 genes. Those genes in the terminal of chromosome
87 were defined as one bin though the number of genes might be less than 100. The
88 accumulated expression level and chimeric level of each bin were calculated by sums
89 of expression level and chimeric level from genes in the bin, respectively.

90

91 **Statistics and reproducibility**

92 Non-parametric Mann-Whitney U-test (Wilcoxon rank-sum test, two-sided) is applied
93 for calculating p value to assess the statistical significance of differences between two
94 groups, which has also been mentioned in the related figure legends. Linear regression
95 model is performed to evaluate the correlation between chimeric level and expression
96 level (both RPKM and CPM).

97

98 **Data Availability**

99 The RNA-seq and whole genome sequencing data supporting the conclusions of this

100 article has been deposited in the Genome Sequence Archive under accession number
101 CRA004187 linked to the project PRJCA004871.

102

103 **REFERENCES**

104 Chen, S., Zhou, Y., Chen, Y., and Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ
105 preprocessor. *Bioinformatics* 34, i884-i890.

106 Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,
107 Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner.
108 *Bioinformatics* 29, 15-21.

109 Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-
110 Wheeler transform. *Bioinformatics* 25, 1754-1760.

111 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,
112 Abecasis, G., Durbin, R., and Genome Project Data Processing, S. (2009). The
113 Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078-2079.

114 Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for
115 comparing genomic features. *Bioinformatics* 26, 841-842.

116 Trapnell, C., Roberts, A., Goff, L., Pertea, G., Kim, D., Kelley, D.R., Pimentel, H.,
117 Salzberg, S.L., Rinn, J.L., and Pachter, L. (2012). Differential gene and transcript
118 expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* 7,
119 562-578.

120 **Figure S1. Viral chimeric gene transcripts were integrated with both mRNAs and**
121 **non-coding RNAs from host.**

122 **(A)** Heatmap displaying the chimeric level for filtered chimeric genes in each sample.
123 The chimeric level was defined as CPM of chimeric reads and normalized by
124 sequencing depth.

125 **(B)** Venn diagram showing the number of shared and specific chimeric genes from **(A)**
126 among 3 human cell lines. The overlapped genes between Huh-7 and Calu-3 cells
127 mainly encode mRNAs, lncRNAs and miscRNAs.

128 **(C)** Pie charts showing the proportions of reads aligned to human (blue) and SARS-
129 CoV-2 (pink) genomes in infected 293T, Huh-7 and Calu-3 cells.

130 **(D)** Line chart displaying the proportion of chimeric genes supported by various
131 numbers of chimeric reads in 293T (yellow), Huh-7 (blue) and Calu-3 (pink) cells.

132 **(E)** Barplot displaying the proportions of RNA types for chimeric genes in 3 human
133 cell lines.

134 **(F)** IGV tracks displaying the junction loci in both human RNA and SARS-CoV-2 RNA
135 for common chimeric genes *ATP5F1A* (top), *EEF2* (middle) and *RPS19* (bottom). The
136 lines with different colors indicate the sources of chimeric reads. Yellow, blue and pink
137 represents 293T, Huh-7 and Calu-3 cells, respectively. Blocks with different colors
138 represent 5'UTR, ORF1ab, S, ORF3a, E, M, ORF6, ORF7ab, ORF8, N and 3'UTR
139 along the SARS-CoV-2 genome from left to right.

140 **(G)** Line chart displaying the proportions of chimeric events supported by various
141 numbers of chimeric reads in 293T (yellow), Huh-7 (blue) and Calu-3 (pink) cells.

142 **(H)** Pie charts showing the proportions of reads aligned to green monkey (green) and
143 SARS-CoV-2 (pink) in infected Cos-7 cells (left), and proportions of reads aligned to
144 rhesus monkey (yellow) and SARS-CoV-2 (pink) in infected MA-104 cells (right).

145 **(I)** Line chart displaying the proportion of chimeric genes represented by various
146 numbers of chimeric reads in Cos-7 (green) and MA-104 (yellow) cells.

147 **(J)** Barplot displaying the proportion of RNA types for chimeric genes in Cos-7 and
148 MA-104 cells.

149 **(K)** Line chart displaying the proportions of chimeric events represented by various
150 numbers of chimeric reads in Cos-7 (left) and MA-104 (right) cells.

151

152

153 **Figure S2. The chimeric events were not conserved between replicates.**

154 (A) Venn diagrams showing the number of shared chimeric genes for two replicates
155 from same cell line (top) and conserved chimeric events on shared chimeric genes at
156 different resolutions (1 nt, 10 nt, 50 nt and 100 nt). For venn diagrams of chimeric
157 events, the values in brackets represents the numbers of chimeric genes containing the
158 shared chimeric events.

159 (B) IGV tracks displaying the junction loci in both human RNA and SARS-CoV-2 RNA
160 for chimeric genes in each replicate (top). The shared chimeric events were displayed
161 at different nucleotide resolutions (1 nt, 10 nt, 50 nt and 100 nt). Blocks with different
162 colors represent 5'UTR, ORF1ab, S, ORF3a, E, M, ORF6, ORF7ab, ORF8, N and
163 3'UTR along the SARS-CoV-2 genome.

164

165 **Figure S3. Chimeric genes were enriched in highly expressed genes for human and**
166 **SARS-CoV-2.**

167 (A) Cumulative distributions of expression levels for chimeric genes (dark color) and
168 non-chimeric genes (light color) in 293T, Huh-7 and Calu-3 cells. *P* values were
169 calculated by Mann-Whitney U test.

170 (B) Number of genes within different expression levels (grey) and proportion of
171 chimeric genes in each pool within various expression levels. Yellow, blue and pink
172 represent 293T, Huh-7 and Calu-3 cells, respectively.

173 (C) IGV tracks displaying the intensity of chimeric reads along SARS-CoV-2 genome
174 in infected 293T (top), Huh-7 (middle) and Calu-3 (bottom) cells.

175 (D) IGV tracks displaying the depth of perfect matched reads along SARS-CoV-2
176 genome in two replicates of each cell line.

177

178 **Figure S4. Both chimeric genes and events in Cos-7 (green monkey) and MA-104**
179 **(rehesus monkey) cells showed similar pattern to human.**

180 (A) Scatter plot displaying the chimeric levels for chimeric genes in two replicates. The
181 sizes of scatters represent number of chimeric genes with the same distribution of
182 chimeric levels in two replicates.

183 (B) Venn diagrams showing the number of shared chimeric genes for two replicates of
184 each cell line (top) and conserved chimeric events on shared chimeric genes at different
185 nucleotide resolutions (1 nt, 10 nt, 50 nt and 100 nt). For venn diagrams of chimeric

186 events, the values in brackets represent the numbers of chimeric genes containing the
187 shared chimeric events.

188 **(C)** IGV tracks displaying the junction loci in both monkey RNA and SARS-CoV-2
189 RNA for chimeric genes in two replicates (top). The shared chimeric events were
190 displayed at different resolutions (1nt, 10nt, 50nt and 100nt). Blocks with different
191 colors represent 5'UTR, ORF1ab, S, ORF3a, E, M, ORF6, ORF7ab, ORF8, N and
192 3'UTR along the SARS-CoV-2 genome.

193 **(D)** Boxplot displaying the distributions of expression level for chimeric genes (dark
194 color) and non-chimeric genes (light color) in infected Cos-7 and MA-104 cells.

195 **(E)** Cumulative distributions of expression levels for chimeric genes (dark color) and
196 non-chimeric genes (light color) in Cos-7 and MA-104 cells. *p* values were calculated
197 by Mann-Whitney U test.

198 **(F)** Number of genes within different expression levels (grey) and proportion of
199 chimeric genes in each pool within varicose expression levels.

200 **(G)** IGV tracks displaying the intensity of chimeric reads along SARS-CoV-2 genome
201 in two replicates of infected Cos-7 (top) and MA-104 (bottom) cells.

202 **(H)** IGV tracks displaying the depth of perfect matched reads along SARS-CoV-2
203 genome.

204

205 **Figure S5. Information of chimeric genes along chromosomes and coverage of**
206 **reads from whole genome sequencing.**

207 **(A)** Barplot displaying the numbers of chimeric genes along different chromosomes in
208 infected 293T, Huh-7 and Calu-3 cells.

209 **(B)** Barplot displaying the numbers of all annotated genes along different chromosomes.

210 **(C)** Chimeric levels and corresponding expression levels of genes along the bin with
211 the highest accumulated chimeric levels (chr14: 39,385,404-49,852,821) were
212 displayed by scatters and lines, respectively. The arrows point out the chimeric levels
213 and gene loci of chimeric genes.

214 **(D)** The frequency and accumulated expression levels of chimeric genes in
215 corresponding bins were displayed by scatters and lines, respectively.

216 **(E)** Chimeric levels and corresponding expression levels of genes along the bins with
217 the most chimeric genes were displayed by scatters and lines, respectively. The arrows
218 indicate the loci of chimeric genes.

219 (F) IGV tracks showing the coverage of whole genome sequencing along genome in
220 each sample.

221

222 **Figure S6. Chimeric events might come from library constructions but not natural**
223 **events.**

224 (A) Pie chart showing the proportions of reads aligned to human (blue), SARS-CoV-2
225 (pink) and zebrafish (green) in the mixed library.

226 (B) Line chart displaying the proportion of chimeric events represented by various
227 numbers of chimeric reads in mixed library. Blue and green lines represent viral
228 chimeric genes with human or zebrafish RNA, respectively.

229 (C) Cumulative distributions of expression levels for chimeric genes (dark color) and
230 non-chimeric genes (light color) for human (left) and zebrafish (right) RNAs.

231 (D) Violin plot displaying the distributions of expression level for chimeric genes (dark
232 color) and non-chimeric genes (light color) for human RNA (left) and zebrafish RNA
233 (right). e, Number of genes within different expression levels (grey) and proportion of
234 chimeric genes in each pool within various expression levels. The left and right panels
235 showing the human RNA and zebrafish RNA from the same mixed library, respectively.

236 (E) Number of genes within different expression levels (grey) and proportion of
237 chimeric genes in each pool within various expression levels. The left and right panels
238 showing the human RNA and zebrafish RNA from the same mixed library, respectively.

239 (F) Distribution of chimeric (dark green) and non-chimeric reads (light green) across
240 the length of chimeric zebrafish mRNAs. 5'UTRs, CDSs, and 3'UTRs of human
241 mRNAs were individually binned into regions spanning 1% of their total length, and
242 the percentages of chimeric and non-chimeric reads that fall within each bin were
243 determined, respectively.

244 (G) IGV tracks displaying the depth of viral RNA along SARS-CoV-2 genome (top)
245 and intensity of chimeric reads with human and zebrafish, individually (bottom).

246 (H) Scatter plot displaying the correlation between chimeric level and expression level
247 (RPKM) of chimeric genes for human (left) and zebrafish (right) from the same mixed
248 library.

249 (I) Pie chart showing the proportions of reads identified as perfect matched reads
250 (including reads aligned to human, zebrafish and SARS-CoV-2), and human-zebrafish
251 chimeric reads (red).











