# Comprehensive in-vivo secondary structure of the SARS-CoV-2 genome reveals novel regulatory motifs and mechanisms

Nicholas C. Huston,[1]* Han Wan,[2]* Rafael de Cesaris Araujo Tavares,[3] Craig Wilen,[4,5] Anna Marie Pyle[2,3,6,]#

[1]Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, USA

[2]Department of Molecular, Cellular, and Developmental Biology, Yale University, New Haven, CT, USA

[3]Department of Chemistry, Yale University, New Haven, CT, USA

[4]Department of Laboratory Medicine, Yale School of Medicine, New Haven, CT, USA

[5]Department of Immunobiology, Yale School of Medicine, New Haven, CT, USA

[6]Howard Hughes Medical Institute, Chevy Chase, MD, USA

Running Head: Complete *in-vivo* SHAPE-MaP structure of the SARS-CoV-2 genome

#Address correspondence to Anna Marie Pyle, anna.pyle@yale.edu

*Nicholas C. Huston and Han Wan contributed equally to this work. Nicholas C. Huston is listed first because he performed the large majority of experimental work.

Abstract Word Count: 249

Text Word Count: 8,103

## Summary

25

26        SARS-CoV-2 is the positive-sense RNA virus that causes COVID-19, a disease that has

27    triggered a major human health and economic crisis. The genome of SARS-CoV-2 is unique among

28    viral RNAs in its vast potential to form stable RNA structures and yet, as much as 97% of its 30

29    kilobases have not been structurally explored in the context of a viral infection. Our limited knowledge

30    of SARS-CoV-2 genomic architecture is a fundamental limitation to both our mechanistic understanding

31    of coronavirus life cycle and the development of COVID-19 RNA-based therapeutics. Here, we apply a

32    novel long amplicon strategy to determine for the first time the secondary structure of the SARS-CoV-2

33    RNA genome probed in infected cells. In addition to the conserved structural motifs at the viral termini,

34    we report new structural features like a conformationally flexible programmed ribosomal frameshifting

35    pseudoknot, and a host of novel RNA structures, each of which highlights the importance of studying

36    viral structures in their native genomic context. Our in-depth structural analysis reveals extensive

37    networks of well-folded RNA structures throughout Orf1ab and reveals new aspects of SARS-CoV-2

38    genome architecture that distinguish it from other single-stranded, positive-sense RNA viruses.

39    Evolutionary analysis of RNA structures in SARS-CoV-2 shows that several features of its genomic

40    structure are conserved across beta coronaviruses and we pinpoint individual regions of well-folded

41    RNA structure that merit downstream functional analysis.  The native, complete secondary structure of

42    SAR-CoV-2 presented here is a roadmap that will facilitate focused studies on mechanisms of

43    replication, translation and packaging, and guide the identification of new RNA drug targets against

44    COVID-19.

45

46

47

48

49

## INTRODUCTION

Severe acute respiratory syndrome related coronavirus 2 (SARS-CoV2), which is responsible for the current global pandemic(Zhu et al., 2020), is a positive strand RNA virus in the genus *β-coronavirus*. To date, the outbreak of SARS-CoV2 has infected at least 12 million people globally, causing great economic loss and posing an ongoing public health threat(Dong et al., 2020). Included in the β-coronavirus genus are two related viruses, SARS-CoV and Middle East respiratory syndrome coronavirus (MERS-CoV), that caused global outbreaks in 2003 and 2012, respectively(de Wit et al., 2016). Despite the continued risk posed by β-coronaviruses, mechanistic studies of the family are limited and to date no effective antivirals or vaccines exist, highlighting the need for research that facilitates the development of therapeutics. With most research efforts focusing on viral proteins (Lan et al., 2020a, Yin et al., 2020, Wan et al., 2020), little is known about the viral RNA genome, especially its structural content. This is a major gap in our understanding because RNA structural elements in positive strand viruses play central roles in regulating all aspects of replication, translation, packaging and host defense (McMullan et al., 2007, Fricke et al., 2015, Pirakitikulr et al., 2016, Clyde and Harris, 2006, MacFadden et al., 2018), so an understanding of their location and function is critical for mechanistic understanding and strategies for viral control(Barrows et al., 2018, Adams et al., 2017). Given the success of antimicrobials targeted against conserved RNA structural elements in other pathogen genomes(Warner et al., 2018, Fedorova et al., 2018), there is an urgent, unmet need to elucidate the genome architecture of SARS-CoV-2.

Like other coronaviruses, the genome of SARS-CoV-2 is incredibly large (Maier et al., 2015, Zhu et al., 2020). Two open reading frames (ORF) for viral nonstructural proteins (Nsp) and 9 small ORFs that encode the structural proteins and a number of accessory genes comprise a ~30kb genome(Kim et al., 2020). These ORFs are flanked on either side by a 5' and 3'UTR that have been shown in other coronaviruses to possess conserved RNA structures with important functional roles in the viral life cycle(Yang and Leibowitz, 2015). Studies in Murine Hepatitis Virus (MHV) and Bovine Coronavirus (BCoV) suggest that the 5' viral termini folds in to 6 stems (SL1-SL6) that play roles in sgRNA synthesis

76    and viral replication(Madhugiri et al., 2018, Chen and Olsthoorn, 2010). In the 3'UTR, a pseudoknot

77    and a bulged stem loop (BSL) are essential for sgRNA synthesis in MHV(Zust et al., 2008).

78         One of the best-studied functional RNA elements in β-coronavirus genomes is the programmed

79    ribosomal frameshifting pseudoknot (PRF) that sits at the boundary between Orf1a and Orf1ab(Plant

80    and Dinman, 2008). The PRF, found in all coronaviruses, induces a -1 ribosomal frameshift that allows

81    for bypassing of the Orf1a stop codon and production of the orf1ab polyprotein, which includes the viral

82    replicase(Plant et al., 2005). Reporter assays using a truncated PRF construct showed that

83    programmed frameshifting occurs ~25% of the time in SARS-CoV(Kelly et al., 2020), and that it is

84    crucial for sgRNA synthesis(Plant et al., 2013).  Extensive mutational analysis has revealed a three-

85    stemmed pseudoknot structure for the SARS-CoV PRF(Plant et al., 2005). However, neither the

86    mechanism of frameshifting regulation nor the three-stem pseudoknot PRF conformation has been

87    validated in cells.

88         While recent computational studies suggest the 5'UTR, 3'UTR, and PRF functional elements

89    are conserved in the SARS-CoV-2 genome(Rangan et al., 2020, Andrews et al., 2020), these regions

90    account for a vanishingly small fraction of the total nucleotide content. Studies of other positive-sense

91    viral RNA genomes such as Hepatitis C virus (HCV) and Human Immunodeficiency Virus (HIV) have

92    revealed extensive networks of regulatory RNA structures contained within viral ORFs(Siegfried et al.,

93    2014, Pirakitikulr et al., 2016, Friebe and Bartenschlager, 2009, Li et al., 2018, You et al., 2004) which

94    direct critical aspects of viral function. It is therefore of crucial importance to assess and characterize

95    the structural features of the SARS-CoV-2 ORF, as elucidation of structural motifs will improve our

96    understanding of all coronaviruses and facilitate development of antiviral therapies for the entire family.

97         Recent advances in high-throughput structure probing methods (SHAPE-MaP, DMS-MaP) have

98    greatly facilitated the structural studies of long viral RNAs (Siegfried et al., 2014, Zubradt et al., 2017).

99    Recently, Manfredonia et al. performed full-length SHAPE-MaP analysis on *ex vivo* extracted and

100   refolded SARS-CoV-2 RNA(Manfredonia et al., 2020). However, structural studies on both viral and

101   messenger RNA have highlighted the importance of probing RNAs in their natural cellular

102   context(Simon et al., 2019, Rouskin et al., 2014). Lan et al performed full-length *in-vivo* DMS-MaPseq

103   on SARS-CoV2 infected cells(Lan et al., 2020b), but as DMS only reports on A and C nucleotides, the

104   data coverage is necessarily sparse. While both studies reveal important features of the structural

105   content in the SARS-CoV-2 genome and its evolutionary conservation, to date no work has been

106   published that captures information for every single nucleotide in an *in-vivo* context.

107        Here, we report for the first time the complete secondary structure of SARS-CoV-2 RNA

108   genome using in SHAPE-MaP data obtained in living cells. We deploy a novel long amplicon method

109   readily adapted to other long viral RNAs made possible by the highly processive reverse transcriptase

110   MarathonRT (Guo et al., 2020). The resulting genomic secondary structure map reveals functional

111   motifs at the viral termini that are structurally homologous to other coronaviruses, thereby fast-tracking

112   our understanding of the SARS-CoV-2 life cycle. We reveal conformational variability in the PRF,

113   highlighting the importance of studying viral structures in their native genomic context and underscoring

114   their dynamic nature. We also uncover elaborate networks of well-folded RNA structures dispersed

115   across Orf1ab, and we reveal features of the SARS-CoV-2 genome architecture that distinguish it from

116   other single-stranded, positive-sense RNA viruses. The analysis reveals large RNA structures within

117   the ORF that may ultimately prove to be as important for viral function as the PRF. Evolutionary

118   analysis of the full-length SARS-CoV-2 structure suggests that, not only do its architectural features

119   appear to be conserved across the β-coronavirus family, but individual regions of well-folded RNA may

120   be as well. Our work reveals the unique genomic architecture of SARS-CoV-2 in infected cells, points to

121   important viral strategies for infection and persistence, and identifies potential drug targets. The full-

122   length structure model we present here thus serves as an invaluable roadmap for future studies on

123   SARS-CoV-2 and other coronaviruses that emerge in the future.

124

125   **Keywords**

126   SARS-CoV-2, RNA genome architecture, *in-vivo* SHAPE-MaP, Functional RNA structure

127

128 **MATERIALS & METHODS**

129 **Cell Culture and SARS-CoV-2 Infection**

130       VeroE6 cells were cultured in Dulbecco's Modified Eagle Medium (DMEM) with 10% heat-

131 inactivated fetal bovine serum (FBS). Approximately $5\times10^6$ cells were plated in each of four T150 tissue

132 culture treated flasks. The following day media was removed and $10^5$ PFU in 4mL of media of SARS-

133 CoV-2 isolate USA-WA1/2020 (BEI Resources #NR-52281) was added to each flask. Virus was

134 adsorbed for 1 hour at 37°C and then 16mL of fresh media was added to each flask.

135

136 **RNA Probing and Purification**

137       Four days post-infection (dpi), the supernatant was aspirated from each flask, cells were

138 washed with 10mL of cold PBS-/- and then dislodged in 10ml PBS-/- with a cell scraper. The contents

139 were collected and centrifuged at 450g x 5 min at 4°C. The supernatant was removed and the cell

140 pellet was resuspended into 2ml of PBS-/- with 200µl DMSO or 2ml PBS with 200µl of 2M NAI (final

141 concentration = 200mM). Cells were incubated for 10 minutes at room temperature followed by

142 addition of 6mL of Trizol. RNA was extracted with the addition of 1.2mL of chloroform. The aqueous

143 phase was transferred to a new tube, followed by the addition of 12mL of 100% EtOH (70% final) and

144 precipitated overnight at -20°C. RNA was resuspended in 1xME buffer and purified using the Qiagen

145 RNeasy kit according to the manufacturer's protocol. RNA was eluted in 1xME buffer.

146

147 **Tiled-Amplicon Design**

148       Leveraging the extreme processivity of MarathonRT, a highly processive group II intron-

149 encoded RT(2), we designed fifteen 2000nt amplicon and a single 1300nt amplicons tiled across the

150 SARS-CoV-2 genome for full sequencing coverage. Adjacent amplicons were designed with a 100nt

151 overlap to ensure data is collected for regions otherwise masked by primer binding. Primers for reverse

152 transcription (RT) were designed using the OligoWalk tool(Lu and Mathews, 2008) to avoid highly-

153 structured primers and highly-structured regions of the SARS-CoV-2 genome. Forward and reverse

6

154    primer sets were designed for an optimal $T_m$ of 58°C. Reverse primers were inset 3nt from the 5'end of

155    the RT primer to enhance specificity of the PCR reaction.

156

157    **Reverse Transcription with MarathonRT**

158        MarathonRT purification was performed as described in (Guo et al., 2020). For each amplicon,

159    500ng of total cellular RNA was mixed with 1μL of the corresponding 1μM RT primer. Gene-specific

160    primers used for RT are listed in **Table S2**. Primers were annealed at 65°C for 5min then cooled to

161    room temperature, followed by addition of 8μL of 2.5x MarathonRT SHAPE-Map Buffer (125mM 1M

162    Tris-HCl pH 7.5, 500mM KCl, 12.5mM DTT, 1.25mM dNTPs, 2.5mM $Mn^{2+}$), 4μL of 100% glycerol, and

163    0.5μL of MarathonRT. RT reactions were incubated at 42°C for 3 hours. 1μL 3M NaOH was added to

164    each reaction and incubated at 95°C for 5min to degrade the RNA, followed by the addition of 1μL 3M

165    HCl to neutralize the reaction. cDNA was purified using AmpureXP beads (Cat**.** No. A63880) according

166    to manufacturer's protocol and a 1.8x bead-to-sample ratio. Purified cDNA was eluted in 10μL

167    nuclease-free water.

168

169    **SHAPE-MaP Library Construction**

170        Amplicons tiling the SARS-CoV-2 genome were generated using NEBNext UltraII Q5 MasterMix

171    (Cat. No. M0544L), gene-specific forward and reverse PCR primers, and 5μL of purified cDNA. Gene-

172    specific primers used for PCR are listed in **Table S3**. Touchdown cycling PCR conditions were used to

173    enhance PCR specificity (68-58°C annealing temperature gradient). PCR reaction products were

174    purified with Monarch DNA Clean-up Kits (NEB) with a binding buffer:sample ratio of 2:1 to remove

175    products smaller than 2kb. PCR products were visualized on 0.8% agarose gels to confirm production

176    of correctly sized amplicons. Amplicons were diluted to 0.2ng/uL and then pooled into two odd and two

177    even amplicon pools for downstream library preparation. Sequencing libraries were generated using a

178    NexteraXT DNA Library Preparation Kit (Illumina) according to manufacturer's protocol, but with 1/5[th]

179    the recommended volume. Libraries were quantified using a Qubit (Life Technologies) and a

180    BioAnalyzer (Agilent). Amplicon pools were recombined and sequenced on a NextSeq 500/550

181    platform using a 150 cycle mid-output kit.

182

183    **Structure Prediction**

184    All libraries were analyzed using ShapeMapper 2(Busan and Weeks, 2018), aligning reads to

185    SARS-CoV-2 genome (accession number: MN908947). Mutation rates between NAI-modified and

186    unmodified samples were tested for significance using the equal variance t-test. Using reactivities

187    output from ShapeMapper, ShapeKnots(Hajdin et al., 2013) was used to determine whether two

188    previously reported pseudoknots contained in the SARS-CoV-2 genome were predicted with

189    experimental SHAPE constraints. The two pseudoknots tested were the programmed ribosomal

190    frameshifting element that exists at the Orf1a/b boundary, and a pseudoknot in the 3'UTR that was

191    identified in the MHV and B-CoV genomes (references). We analyzed all 500nt windows separated by

192    a 100nt slide that contained each of the putative pseudoknots to determine if the pseudoknot was

193    successfully predicted.

194    SuperFold (Smola et al., 2015b) was used to generate a consensus structure prediction for the

195    entire SARS-CoV-2 genome using SHAPE reactivities obtained from biological replicate 1 as

196    constraints. We imposed a maximum pairing distance of 500nt. As our data only supported formation of

197    the pseudoknot contained in the programmed ribosomal frameshifting element, only this pseudoknot

198    was forced in this prediction. All structures output from the SuperFold prediction were visualized and

199    drawn using StructureEditor, a tool in the RNAStructure software suite(Reuter and Mathews, 2010).

200    Base-pairing distances were calculated from .ct structure files output from SuperFold full-length

201    SARS-CoV-2 consensus predictions, and compared to previously published, publically available full-

202    length genome structures for Dengue and Hepatitis C Virus generated with SHAPE constraints, a max-

203    pairing distance of 500nt, and the SuperFold pipeline (Mauger et al., 2015, Dethoff et al., 2018).

204

205    **Identification of Well-Folded Regions**

8

206    Two data signatures were used to identify well-folded regions: The first is the SHAPE reactivity

207    data generated with the SHAPE-MaP workflow and the ShapeMapper analysis tool(Busan and Weeks,

208    2018). The second is the Shannon entropy calculated from base-pairing probabilities determined during

209    the SuperFold partition function calculation(Smola et al., 2015b). Two replicate data sets were used,

210    including separate SuperFold predictions.

211    Local median SHAPE reactivity and Shannon Entropy were calculated in 55nt sliding windows.

212    The global median SHAPE reactivity or Shannon Entropy were subtracted from calculated values to aid

213    in data visualization. Regions with local SHAPE and Shannon Entropy signals 1) below the global

214    median 2) for stretches longer than 40 nucleotides 3) that appear in both replicate data sets were

215    considered well-folded. Disruptions, or regions where local SHAPE or Shannon Entropy rose above the

216    global median, are not considered to disqualify well-folded regions if they extended for less than 40

217    nucleotides. Arc plots generated from each replicate consensus structure predication were compared

218    for regions that meet sorting criteria described above in order to ensure agreement between secondary

219    structure models generated from each replicate SHAPE-MaP dataset.

220    Base-pairing distances of well-folded regions were calculated from .ct structure files output from

221    SuperFold consensus predictions, and compared to previously published, publicly available structures

222    for well-folded regions of the HIV genome generated with SHAPE constraints, a max-pairing distance of

223    500nt, and the SuperFold pipeline (Siegfried et al., 2014).

224

225    **Multiple sequence alignment**

226    To analyze evolutionary support for our consensus secondary structure prediction of the SARS-

227    CoV-2 genome, we generated two codon-based multiple sequence alignments (MSA) for Orf1a and

228    Orf1b constructed from genomes of closely related viral species (Douzery EJP,2018). All sequences

229    were chosen based on a phylogenetic study of SARS-CoV-2 (Ceraolo and Giorgi, 2020). All sequences

230    referenced below were downloaded from the NCBI Taxonomy browser(Benson et al., 2018).

231        A sarbecovirus MSA was generated using SARS-CoV-2 isolate Wuhan-Hu-1 (MN908947.3),

232    four bat coronaviruses (MG772934.1, JX993987.1, DQ022305.2, DQ648857.1), and five human SARS

233    coronaviruses (AY515512.1, AY274119.3, NC_004718.3, GU553363.1, DQ182595.1).

234        We also generated an "All β-coronavirus Alignment" using the sarbecovirus sequences

235    described above in addition to four MERS-CoV sequences (MK129253, KP209307, MF598594,

236    MG987420), one HKU-4 sequence (MH002337), three HKU-5 sequence (MH002342, NC009020,

237    MH002341), four HKU1 sequences (KY674942, KF686343, AY597011, DQ415903), three murine

238    hepatitis virus sequences (AY700211, AF208067, AB551247), three human coronavirus OC43

239    sequences (AY585229, NC006213, MN026164), two bovine coronavirus sequences (KU558922,

240    KU558923), and one camel coronavirus sequence (MN514966).

241        The orf1a and orf1b region were extracted from the full-length sequences based on the

242    GenBank annotation. Separate codon alignments for both Orf1a and orf1b were generated using

243    MACSE v2.0.3(Ranwez et al., 2018)  and default parameters (*-prog alignSequences*).

244

245    **Synonymous mutation rate analysis**

246        All codon alignments were visualized and edited using Jalview v 2.11.0(Waterhouse et al.,

247    2009). Synonymous mutation rates for each codon were estimated using the phylogenetic-based

248    parametric maximum likelihood (FUBAR) method(Murrell et al., 2013). Each codon was categorized as

249    base-paired or unpaired depending on strandedness of the nucleotide at the third position of each

250    codon in our SARS-CoV-2 consensus structure model(Dethoff et al., 2018).  The significance of

251    synonymous mutation rates between single- and double-stranded regions was determined using two-

252    tailed, equal variance *t*-test.

253

254    **Covariation analysis**

255        Covariation calculation and visualization was performed using R-chie(Lai et al., 2012). The

256    Sarbecovirus codon alignment described above was used for covariation analysis. Identification of

257    base-pairs with statistically significant evidence of covariation was performed on individual structures

258    using R-Scape (version 0.2.1)(Rivas et al., 2017) with the RAFSp statistics by using the "--RAFSp"

259    flag(default E-value:0.05 )(Tavares et al., 2019).

260

261    **Data Availability**

262    All ShapeMapper outputs, secondary structure files, and multiple sequence alignments use in this work

263    are available at the GitHub repository:   https://github.com/pylelab/SARS-CoV-

264    2_SHAPE_MaP_structure

265

266

267    **RESULTS**

268    ***In-vivo* SHAPE-MaP workflow yields high quality data suitable for structure prediction.**

269         To study the SARS-CoV-2 structure in the context of infected cells, the SARS-CoV-2 isolate

270    USA-WA1/2020, isolated from an oropharyngeal swab from a patient who had returned to the United

271    States from China and developed clinical disease, was used to infect VeroE6 cells (BEI Resources

272    #NR-52281). Infection was allowed to proceed for four days, at which point cells were collected and

273    treated with either NAI or DMSO. RNA was then extracted and purified. To generate sequencing

274    libraries, 2000 nucleotide (nt) overlapping amplicons were tiled across the entire SARS-CoV-2 genome

275    (**Fig. 1A**). Importantly, this approach is made possible by the utilization of the ultra-high processive

276    reverse transcriptase MarathonRT, which encodes NAI adducts as cDNA mutations during long-

277    amplicon SHAPE-MaP library construction(Guo et al., 2020). Specifically, gene-specific primers were

278    used to prime reverse transcription in the presence of manganese, followed by amplicon PCR with

279    gene-specific primers and cycling conditions designed to enhance specificity(Korbie and Mattick, 2008).

280    Gel electrophoresis confirmed successful amplification of all 16 amplicons (data not shown).

281    Sequencing of SHAPE-MaP libraries was performed using the Illumina NextSeq 500/550 platform.

282  After generating two independent biological replicates, the resulting sequencing data were

283  analyzed using the ShapeMapper pipeline(Smola et al., 2015b). Comprehensive datasets were

284  obtained, with median effective read depth > 70,000x and effective reactivity data for 99.7%

285  (29813/29903) of nucleotides in the SARS-CoV-2 genome in both replicate experiments.  To check the

286  SHAPE-MaP data quality, we analyzed the relative mutation rates of NAI-treated and DMSO-treated

287  RNA samples, revealing a significant elevation of mutation rates for NAI-treated samples (**Fig. 1B**, p-

288  value < 0.0001). This confirms that the full-length SARS-CoV-2 RNA was successfully modified *in-vivo*

289  and that these modifications were encoded as cDNA mutations.

290  To understand the relative SHAPE reactivity agreement within local regions of the genome, we

291  calculated Pearson correlation coefficients between two biological replicates. The Pearson's correlation

292  across the entire span of Orf1ab is 0.62, consistent with those previously reported for reactivities

293  calculated from *in-vivo* modified RNAs of this size(Smola et al., 2016). Across the sub genomic RNA

294  ORFs, the Pearson 's correlation is poor. We believe this reflects the fact that Amplicons 13, 14, 15,

295  and 16 will amplify both full-length *and* sub-genomic RNAs, and the difference in context will result in

296  different secondary structures(Tavares et al., 2020). For this reason, despite the fact all data have been

297  obtained globally, subsequent discrete structural analysis will focus on shared features of the viral

298  termini and the Orf1ab region.

299

300  **De novo structure prediction on full-length SARS-CoV-2 RNA identifies conserved functional**

301  **elements at the 5' and 3' genomic termini**

302  We performed secondary structure prediction with the SuperFold pipeline(Smola et al., 2015b),

303  using the *in-vivo* SHAPE reactivities to generate an experimentally constrained consensus secondary

304  structure prediction for the entire SARS-CoV-2 genome. As an extensive body of research has

305  elucidated structured RNA elements at the 5' and 3' viral termini as well as the Orf1ab boundary, with

306  conserved functions across β-coronaviruses, we first examined these regions from our consensus

12

307    prediction to determine whether they were stably folded and well-determined in the SARS-CoV-2

308    genome.

309        The 5' genomic terminus includes seven regions that have been identified and studied in other

310    coronaviruses (Reviewed in (Yang and Leibowitz, 2015)).  While sequence conservation suggested that

311    these elements might be conserved in SARS-CoV-2, our consensus structure prediction shows this to

312    be the case, and we derived a specific experimentally-determined structure for this section of the

313    genome.  The in-vivo SHAPE reactivity data correspond well with the resulting structural model (**Fig 2A,**

314    **inset**) and the low overall Shannon entropy values in this region (determined from base pair probability

315    calculation during the SuperFold prediction pipeline(Smola et al., 2015b)) support a well-determined

316    structure for the 5' genomic terminus (median$_{Nuc(1-400)}$ = $2.7x10^{-5}$ ; global median = 0.022).

317        Individual features that typify coronavirus structures are evident in the secondary structure of

318    the SARS-CoV-2 5'-UTR with good SHAPE reactivity agreement (**Fig 2A,** inset).  For example, we

319    observe a bipartite domain architecture for SL1, which was previously reported to play a role in

320    coronavirus replication(Li et al., 2008) (**Fig 2A**, labeled SL1). Similarity between SL1 structures

321    reported for other coronaviruses and the experimentally-determined structure reported here for SARS-

322    CoV-2 suggests that SL1 plays a similar role in SARS-CoV-2 life cycle.

323        Structural studies of SARS-CoV SL2 have shown that the SL2 pentaloop is stacked atop a 5-bp

324    stem. In addition, the pentaloop adopts a canonical CUYG fold in which the uracil is flipped out,

325    resulting in an architecture that is important for sgRNA synthesis(Lee et al., 2011). Our experimentally-

326    determined structure of SL2 from SARS-CoV-2 shows that it adopts exactly the same RNA fold (**Fig 2A**,

327    labeled SL2), again suggesting that it plays the same functional role in the SARS-CoV-2 life cycle.

328        The transcription regulatory sequence (TRS) is a conserved feature of β-coronaviruses and it is

329    required for sgRNA production(Yang and Leibowitz, 2015). The SARS-CoV leader TRS is predicted *in*

330    *silico* to be in stem loop 3 (SL3), with nucleotides exposed in its loop and base-paired in the stem(Chen

331    and Olsthoorn, 2010). The primary sequence of the SARS-CoV leader TRS is absolutely conserved

332    between SARS-CoV and SARS-CoV-2(Chen and Olsthoorn, 2010) ( 5'-ACGAAC-3'). Importantly, our

13

333 consensus prediction shows the SARS-CoV-2 leader TRS is also found in SL3, with a similar structural

334 organization as reported for other viruses (**Fig 2A**, labeled SL3, TRS indicated with solid black line)**.**

335       The SL4 of SARS-CoV-2 adopts a bipartite domain structure (**Fig 2A**, labeled SL4a, SL4b)

336 similar to that reported for MHV(Kang et al., 2006, Yang et al., 2015). Importantly, the AUG of the

337 predicted upstream ORF(uORF), which is phylogenetically conserved among β-coronaviruses(Raman

338 et al., 2003), is found in the top-most stem loop of SL4a, meaning it would be accessible for recognition

339 by a scanning ribosome (**Fig 2A**, uORF start codon indicated with solid grey line).

340       As predicted across coronaviruses(Chen and Olsthoorn, 2010), the trifurcated stem at the top of

341 SL5 is observed in the experimentally-determined structure of SARS-CoV-2 (**Fig 2A**, labeled SL5A-C).

342 This includes UUCGU pentaloop motifs in SL5A and SL5B, and a GNRA tetraloop in SLC. Previous

343 reports suggest this may represent a packaging signal for GroupIIB CoVs(Chen and Olsthoorn, 2010).

344       SL6 and SL7 are predicted in the SARS-CoV-2 structure, and the in-vivo SHAPE data agree

345 strongly support the existence of these stems (**Fig 2A**, labeled SL6 and SL7). However, functional

346 evidence for SL6 and SL7 is lacking in the literature for any coronavirus.

347       The 3' genomic terminus includes three well-studied stems, including the bulged-stem loop

348 (BSL), Stem Loop 1 (SLI), and a long-bulge stem that includes the hypervariable-region (HVR), the

349 S2M domain, the octanucleotide motif (ONM) subdomains, and a pseudoknot (Reviewed in (Yang and

350 Leibowitz, 2015)). The consensus structure recapitulates the secondary structure of all the three stems

351 with good SHAPE reactivity agreement (**Fig 2B, inset**) and overall low Shannon entropy

352 (median$_{Nuc(29,472-29,870)}$ = 0.016 ; global median = 0.022).  While the BSL is well determined in our

353 structure, the low reactivity for bulged nucleotides suggests the possibility of protein binding-partners

354 (**Fig 2B,** labeled BSL**)**.

355       A pseudoknot structure is proposed to exist between the base of the BSL stem loop and the

356 loop of SL1 in coronaviruses(Yang and Leibowitz, 2015). While pseudoknot formation is mutually

357 exclusive with the base of the BSL, studies in MHV have suggested that both structures contribute to

358 viral replication and the mutually exclusive structures are thought to function as a molecular switch in

14

359    different steps of RNA synthesis(Goebel et al., 2004). However, our *in-vivo* determined secondary

360    structure is inconsistent with formation of the pseudoknot (**Fig 2B**, putative base-pairing interactions

361    indicated by black lines). The low SHAPE reactivities for the nucleotides at the base of the BSL support

362    formation of the extended BSL stem, while high-reactivities of the nucleotides in the loop of SLI indicate

363    that it is highly accessible. Using the SHAPEKnots program for robust predication of pseudoknots

364    (implemented in RNA structure v5.8(Hajdin et al., 2013)), we found that a pseudoknot is never

365    predicted in three 500nt windows that cover the pseudoknotted region. Taken together, our data

366    strongly support the extended BSL conformation, indicating it is probably the dominant conformation *in-*

367    *vivo*.

368           The third stem in the 3' UTR includes three sub-domains. The HVR, so-named because it is

369    poorly conserved across group II coronaviruses(Goebel et al., 2007), is predicted to be mostly single-

370    stranded in our secondary structure, and the high reactivities across the span of this region lends

371    strong experimental support for an unstructured region (**Fig. 2B**, region indicated with solid black line

372    and labeled HVR**)**. The fact that this region is highly unstructured may also explain why it has been

373    experimentally demonstrated to tolerate numerous deletions, rearrangements, and point mutations in

374    MHV(Goebel et al., 2007).

375           The S2M region is contained within the apical part of the third stem. We observe that the first

376    three helices of S2M from SARS-CoV-2 exactly match the crystal structure determined for S2M from

377    SARS-CoV (Robertson et al., 2005). However, our in-vivo structure deviates significantly at the top of

378    the stem, with bases that are highly reactive (**Fig. 2B,** region indicated with solid black line and labeled

379    S2M**)**. It is possible that the SARS-CoV-2 S2M folds into a unique S2M conformation despite differing

380    by only a two bases (**Fig. 2B**, base-changes indicated by arrows; SARS-CoV base identity shown in

381    red). Indeed, as both single-nucleotide changes are transversions, any base-pairing interaction

382    involving these nucleotides in the SARS-CoV S2M structure could not be maintained in SARS-CoV-2.

383    Alternatively, this site could interact with factors *in-vivo* that are not captured in the crystallographic

384    study.

385    The ONM is predicted at the central bulge between the S2M region and the HVR region. The

386    sequence is absolutely conserved across β-coronavirus(Goebel et al., 2007), but no functional

387    significance has yet been shown. In our consensus structure, it is single-stranded (**Fig. 2B**, ONM

388    indicated with solid black line and labeled ONM).

389    Finally, we predict a completely different structure for the downstream terminal stem in the viral

390    3'UTR region (adjacent to the poly-A tail) than previously reported for other coronaviruses(Zust et al.,

391    2008). However, our structure prediction in this region is not highly accurate because of proximity of the

392    primer binding site.  That said, the putative stem is predicted to have high Shannon entropy

393    (median$_{Nuc(29472-29495,29853-29870)}$ = 0.2154 ; global median = 0.022), suggesting that it is not a well-ordered

394    structure in the cellular environment.

395

396    **Structure prediction of the programmed ribosomal frame-shifting element reveals**

397    **conformational flexibility**

398    One of the most well-studied RNA structures in the coronavirus coding region is the

399    programmed frame-shifting pseudoknot (PRF). It is located between orf1a and orf1b and plays an

400    important role in inducing a -1 frameshift in a translating ribosome, resulting in the synthesis of the

401    polyprotein ab, which includes the SARS-CoV-2 replicase (Plant and Dinman, 2008).

402    The PRF element previously characterized in SARS-CoV is proposed to contain three parts: an

403    attenuator stem loop, a conserved heptanucleotide "slippery" sequence, and a H-type pseudoknot

404    (Plant and Dinman, 2008). We performed SHAPEKnots(Hajdin et al., 2013) over four 500nt windows

405    that cover the pseudoknotted region in the SARS-CoV-2 genome to check if the PRF pseudoknot can

406    be discerned from our vivo SHAPE data. We found that the pseudoknot is successfully predicted in 3

407    out of 4 windows generated by ShapeKnots. Moreover, the nucleotides predicted to be involved in the

408    pseudoknotted helix have low SHAPE-reactivity (**Fig. 3A**, pseudoknot base-pairs indicated with red

409    lines). Our *in-vivo* SHAPE data therefore strongly support the formation of the pseudoknotted helix, and

16

410    the frame-shifting pseudoknot was thereafter included as a hard constraint during secondary structure

411    prediction.

412        The most probable, dominant structure of the PRF region, extracted from the full-length *in-vivo*

413    secondary structure, is shown in **Fig. 3A**. In our model, the SHAPE reactivity and Shannon entropy

414    calculation support a well-folded attenuator stem (AS) immediately upstream of the heptanucleotide

415    slippery sequence (HSS) (**Fig 3A**; AS and HSS indicated with labeled, solid black lines). The attenuator

416    stem has been demonstrated to be important for attenuating frameshifting in SARS-CoV(Cho et al.,

417    2013), and previous reports suggested that the attenuator stem structure is not well conserved between

418    SARS-CoV and SARS-CoV-2(Kelly et al., 2020). By contrast, our results suggest a SARS-CoV-2-

419    specific fold for the attenuator stem. The highly conserved heptanucleotide slippery sequence is

420    predicted to be single-stranded in our in-vivo structural model, which is consistent with studies on other

421    coronaviruses(Plant et al., 2005, Plant and Dinman, 2008).

422        Overall, the dominant structure prediction for the H-type pseudoknot in our structural model

423    differs from the one characterized in SARS-CoV. The H-type pseudokont in SARS-CoV is composed of

424    three coaxially stacked stems: SL1, SL2 and the pseudoknotted helix(Plant et al., 2005). The SL1 stem,

425    which contains the upstream pseudoknotted loop, is well folded in our consensus model as indicated by

426    SHAPE reactivity mapping (**Fig. 3A**; labeled SL1**)** and Shannon entropy (**Fig. 3C,** median$_{Nuc(13476-13503)}$

427    $= 1.9 \times 10^{-4}$; global median = 0.022). Importantly, the region reported to contain the SL2 stem(Rangan et

428    al., 2020, Plant et al., 2005) is predicted as single-stranded in our consensus structure, and

429    consequently does not include SL2  (**Fig 3A**; region indicated by dotted black line). Rather, the

430    dominant structure predicted for the PRF includes a different stem, SL3, that includes the downstream

431    pseudoknot arm (**Fig 3A**; labeled SL3). However, neither the single-stranded region nor SL3 are well-

432    determined in our structure as indicated by Shannon entropy mapping to the region (**Fig. 3C,**

433    median$_{Nuc(13503-13534)}$ = 0.24; global median = 0.022, labeled with a dotted black-line and SL3,

434    respectively).

435    As SuperFold calculates a partition function, lower probability base-pairing interactions are

436    captured during structure prediction steps. We therefore checked alternative, low probability base-pair

437    interactions captured for the PRF region. We found that the single-stranded region (**Fig. 3A**; indicated

438    by dotted black line) forms base-pairing interaction as many as 6 different regions in the SARS-CoV-2

439    genome (data not shown), We sought to determine if the previously reported SL2 conformation was

440    captured among them. Indeed, an alternate, lower probability structure containing an extended SL2 is

441    generated in the SuperFold prediction with the attenuator stem, heptanucleotide slippery sequence,

442    and SL1 intact. (**Fig. 3B**; alternate SL2 conformation labeled). In this variation, the SL2 stem is

443    predicted to fold with a median probability of 20% as determined from the probabilities of each

444    individual base-pair of the SL2 stem (**Fig. 3D;** individual base-pairs indicated with grey dots). In

445    contrast, the SL3 stem predicted in our dominant consensus structure has as much as 80% probability

446    of folding. The chemical probing data does not strongly support one structure over another (**Fig. 3A** and

447    **Fig. 3B**) and likely reflects structural flexibility and pairing promiscuity for the SL2 region. Taken

448    together, data determined *in-vivo* suggest that the frame-shifting pseudoknot in SARS-CoV-2 includes a

449    well-folded attenuator stem, SL1, and a pseudoknot, but that the region containing the putative SL2 is

450    conformationally flexible. Future studies are needed to explore if there is relationship between the

451    structural flexibility and the mechanistic role of the frameshifting pseudoknot.

452

453    **The secondary structure of SARS-CoV-2 Orf1ab reveals a network of unique RNA structural**

454    **elements**

455    While the successful identification of known, functional RNA structural elements lends strong

456    support for our methodology and for the overall secondary structural model, these known regions

457    account for only 3% of the total nucleotide content of the SARS-CoV-2 genome; little is known about

458    remaining 97%.

459    Here we report the first *in-vivo*-derived, SHAPE-constrained secondary structural model that

460    includes a description of the base-pairing interactions for all nucleotides within a coronavirus genome

18

461    (**Fig. 4A;** secondary structures described by arc plots underneath each of the three SHAPE/Shannon

462    plots). Representative secondary structural maps of small regions extracted from the consensus

463    prediction exemplify the types of substructures that are observed in protein-encoding regions of SARS-

464    CoV-2 (**Fig. 4B,** structures contained within Nsp3 and spanning Nsp6&7; nt4716-5682, nt11221-12043,

465    respectively). This resource is a valuable roadmap for ongoing studies, and to that end, a .ct file for the

466    full-length SARS-CoV-2 genome structure is freely available (see Data Availability).

467           To discover additional, well-folded RNA structures within the SARS-CoV-2 genome, we used a

468    sliding 55nt window to calculate the local median Shannon Entropy and we correlated these values with

469    experimentally-determined SHAPE reactivities (**Fig 4A**). Only regions with both median Shannon

470    entropy and SHAPE reactivity signals below the global median for stretches longer than 40nt, and

471    which appear in both replicate data sets, were considered well-determined and stable. In total, we

472    identify 40 such regions in Orf1ab (**Fig. 4B**, shaded). Hereafter, any structured region that meets these

473    above criteria will be referred to as "well-folded."

474           We also identified well-folded regions in the subgenomic RNA region (data not shown).

475    However, our previous correlation analysis suggests that the SHAPE signal from this region includes

476    reactivity signals from multiple RNA species, including genomic and subgenomic RNAs. Given the

477    method deployed to construct our SHAPE-MaP sequencing libraries, it is impossible to deconvolute

478    sgRNA data from genomic data, and new approaches will be required to separate genomic from

479    subgenomic structures.  While it will be interesting to explore this issue in subsequent studies, the

480    following analysis focuses on stable structures within the orf1ab region, which can be uniquely

481    determined.

482           To understand architectural organization of the overall "structuredness", or base-pair content

483    (BPC) within orf1ab, we calculated the double-strand content of individual protein domains within this

484    region of the genome (**Fig. 5A**, grey bars). We find that all protein domains have comparable BPC, with

485    an average of 56% (+/- 6.09%) of nucleotides involved in base-pairing interactions. However, the RNA

486    sequences within each protein domain are not equivalently well-folded (**Fig 5A,** black bars). For

487   example, we observe that ~50% of nucleotides within the 5'UTR, Nsp1, Nsp6, Nsp8, and Nsp12 are

488   concentrated in well-folded regions, suggesting these domains may be hubs for regulatory RNA

489   structures. By contrast, Nsp13, Nsp14, and Nsp16 have <15% of their nucleotide content lies in

490   discretely well-folded regions. At the most extreme end, Nsp10 contains no nucleotides in well-folded

491   regions. Considering that Nsp10 is located immediately upstream of the PRF, this lack of well-folded

492   structures may be important to direct proper frameshifting.

493        While analyzing the resulting secondary structural map, we noticed that the SARS-CoV-2

494   genome contains long-stretches of short, locally-folded stem loops (for example - **Fig. 4B**) with few

495   long-distance base-pairing interactions such as those indicated by large arcs in typical arc plots (**Fig.**

496   **4A**). Wondering if this was quantifiable feature unique to the SARS-CoV-2 genome, we calculated the

497   distance between base-paired nucleotides for every base-pairing interaction in our SARS-CoV-2

498   structural model. We compared these SARS-COV-2 base-pairing distances to those we calculated from

499   published full-length structural models for HCV(Mauger et al., 2015) and Dengue Virus(Dethoff et al.,

500   2018), where the data were prepared using the same structure prediction pipeline and constraints used

501   in our study. Interestingly, the median base-pairing distance in our SARS-CoV-2 consensus model is

502   25nt, and is significantly smaller than the median base-pairing distance in the HCV (median=40nt) and

503   Dengue Virus (median=33nt) consensus models (**Fig. 5B**). Even more, the upper bound of the

504   interquartile range (IQR) that describes the distribution of base-pairing distances in Dengue and HCV

505   genomes is much higher than the same bound in the SARS-CoV-2 genome (SARS-CoV-2 75[th]

506   percentile = 46nt; Dengue Virus 75[th] Percentile = 104nt; HCV 75[th] percentile = 101nt).  This suggests

507   SARS-CoV-2 has fewer long-distance base-paring interactions compared to Dengue and HCV genome.

508        We also calculated the median base-pairing distance for the well-folded regions of the SARS-

509   CoV-2 genome and compared the result to well-folded regions previously identified using the same Low

510   Shannon/Low SHAPE signatures in the HIV genome(Siegfried et al., 2014). We found that although

511   there is no significant difference in the size of well-folded regions in the SARS-CoV2 and HIV genomes

512   (data not shown), the median base-pairing distance in the well-folded regions of SARS-CoV-2 (median

513    = 26nt) is significantly lower than the base-pairing distance in well-folded regions of HIV (median = 34nt)

514    (**Fig. 5C).** Similarly, the upper bound of the IQR that describes the distribution of base-pairing distances

515    in well-folded regions of the HIV genome is much higher than the same bound in SARS-CoV-2 (SARS-

516    CoV-2 75$^{th}$ percentile = 44nt; HIV 75$^{th}$ percentile = 133nt).

517         Taken together, these results suggest that the SARS-CoV-2 genome folds into more local

518    secondary structures, such as the short stem-loops in **Fig. 4B**, and contains fewer long-range base-

519    pairing interactions than observed for other positive-sense RNA viruses.  Given the exceptional size of

520    the coronavirus genome (~30kb) relative to those of the positive-sense RNA viruses compared here

521    (~10kb), it is possible that the short base-pairing distance of SARS-CoV2 may carry functional

522    implications for maintaining genomic stability, preserving fidelity of translation, and evading innate

523    immune response.

524

525    **The overall structured-ness of the SARS-CoV-2 is conserved across all β-coronaviruses**

526         Synonymous mutations rates have been used previously to lend evolutionary support for well-

527    folded RNA secondary structures in other positive-sense RNA viruses(Dethoff et al., 2018, Tuplin et al.,

528    2002, Assis, 2014, Simmonds and Smith, 1999). This body of work has suggested lower synonymous

529    rates for double-stranded nucleotides when compared to single-stranded nucleotides in viral RNAs,

530    likely reflecting an evolutionary pressure to maintain base-pairing interactions of double-stranded

531    nucleotides. We therefore computed relative synonymous mutation rates to determine how evolutionary

532    pressure is applied to single- and double-stranded regions of the SARS-CoV2 genome.

533         To generate the codon-based alignment, 33 full-length genome sequences from the NCBI

534    Taxonomy database(Benson et al., 2018) were collected, including 12 SARS-CoV genomes, 8 MERS-

535    CoV genomes, and 13 more distantly related β-coronaviruses genomes. This alignment was then used

536    to calculate synonymous and non-synonymous mutation rates (dS and dN, respectively) for each codon

537    in the SARS-CoV2 orf1ab region using A Fast, Unconstrained Bayesian AppRoximation for inferring

538    selection (FUBAR) (Murrell et al., 2013). We then separated dS and dN into single- or double-stranded

21

539   bins as predicted in our consensus model. The strandedness of each codon was determined by the

540   strandedness at the third position of the codon(Dethoff et al., 2018).

541       In the "All β-Coronavirus" alignment, we observed a significantly lower synonymous mutation

542   rate (p<0.0001) for double-stranded codons (median = 3.765; IQR = 3.034-4.978) when compared to

543   single-stranded codons (median = 4.189; IQR = 3.232 - 5.562) in our consensus model (**Fig. 6A**). In

544   contrast, there was no significant difference (p = 0.86) observed for non-synonymous mutation rates

545   (dN) at single- (median = 0.4535; IQR = 0.103 – 0.685) or double-stranded codons (median = 0.453;

546   IQR = 0.139 – 0.675) (**Fig. 6B**) as dN reflects changes at the amino acid level. This suggests that

547   double-stranded regions of the SARS-CoV-2 genome experience stronger selective pressure against

548   synonymous mutations than single-stranded regions, which lends support to our consensus model and

549   suggests evolutionary maintenance of the observed secondary structure. Because an all β-coronavirus

550   alignment was used, our results indicate that the structural organization and overall base-pairing

551   content of Orf1ab is a conserved feature of the β-coronavirus family.

552       When analyzing relative synonymous mutation rates within individual protein domains, we

553   observed significantly decreased synonymous mutation rates for double-stranded codons in Nsp1,

554   Nsp2, Nsp3, Nsp4, Nsp6, Nsp8, Nsp12, Nsp13, and Nsp15 (**Fig. 6C**). Consistent with this, Nsp1, Nsp6,

555   Nsp8, and Nsp12 have >50% of their nucleotides localized within well-folded regions (**Fig. 5A,** black

556   bars). Taken together, this suggests that certain protein-coding domains contain regions of RNA

557   secondary structure that are conserved across β-Coronaviruses. For example, Nsp8, which is the most

558   well folded domain in SARS-CoV-2, is likely well-folded in other β-Coronaviruses.

559       By contrast, the base pairing content of Nsp5, Nsp7, Nsp9, Nsp10, Nsp14, and Nsp16 does not

560   appear to be conserved, as there is no significant difference in synonymous mutation rates of single-

561   and double-stranded codons (**Fig. 6C**). Consistent with this, Nsp14 and Nsp16 were shown to have

562   <15% of their nucleotides in well-folded regions, while Nsp10 does not contain any well-folded

563   nucleotides (**Fig. 5A**). Not only does this analysis support the observation that these regions of RNA

564    are not well-folded in SARS-CoV-2, our data suggest these regions may not be well folded in other β-

565    Coronaviruses.

566

567    **Evolutionary analysis for individual well-folded regions of the SARS-CoV-2 genome identifies**

568    **several conserved regions**

569            To further prioritize structural elements that may have conserved functional roles in the SARS-

570    CoV2 life cycle, we next applied our synonymous mutation rate analysis to each of the 40 discrete well-

571    folded domains identified by Low Shannon/Low SHAPE signatures (**Fig. 4B, Table S1**).  Four regions

572    (regions 23, 25, 34, and 36) that are well-determined based on their low median Shannon Entropy

573    values (**Table S1**) and SHAPE reactivity data showed significantly decreased synonymous mutation

574    rates at double-stranded codons when compared to single-stranded codons across the β-coronavirus

575    alignment (**Fig. 7A, 7B**).  Among those structures, region 25 and 34 are found at protein domain

576    boundaries. Region 25 ends exactly at the Nsp8/9 domain boundary, while Region 34 spans the

577    Nsp12/13 boundary. Region 23, 34, and 36 (**Fig. 7C, Fig. 7E, Fig. 7F**) contain a series of stem-loops

578    with small bulges. Region 25 contains a long-range duplex that closes a clover-leaf like structure with 8

579    stem-loops radiating from a central loop (**Fig. 7D**). This hub, or multi-helix junction might represent a

580    promising drug target, as multi-helix junctions often contain binding pockets with high binding affinity

581    and selectivity for small molecules(Warner et al., 2018).

582            Within the Sarbecovirus subgenus, we were able to identify four regions (regions 15, 22, 24, 27,

583    and 30) that are well-determined in our secondary structural model (based on low median Shannon

584    Entropy (**Table S1**) and SHAPE reactivity data) with significantly decreased synonymous mutation

585    rates in double-stranded relative to single-stranded codons (**Fig. 8A, 8B**).  Among these structures,

586    Region 24 contains two discrete multi-helix junctions, each with at least three stems radiating from

587    large central loops (**Fig. 8C**). Region 27, which contains a series of six stem-loops, is particularly

588    significant because it is only 100nt downstream of the PRF in Nsp12 (**Fig 8D**). Region 15, like Region

589    24, contains several well-determined long-range duplexes that segment the region into two discrete

23

590    multi-helix junctions (**Fig. 8E**). Region 22 contains a series of well-folded loops and it spans the Nsp5/6

591    boundary (**Fig. 8F**). Region 30 is a single stem-loop with bulges that divide the stem into distinct

592    duplexes (**Fig. 8G**)

593         To look for evolutionary evidence that directly supports conservation of specific base-pairing

594    interactions and secondary structures, we performed covariation analysis on the 5 structures that are

595    supported by Sarbecovirus-specific synonymous mutation rates. We visualized the base-pair

596    covariation levels using R-chie (Lai et al., 2012) and we used R-scape version 0.2.1(Rivas et al., 2017)

597    with the RAFSp statistics(Tavares et al., 2019) to test the statistical significance of putatively covarying

598    base-pairs. Specifically, we identified 3 regions (15, 22, and 30) that have covariation support **(Fig. 8E-**

599    **G)**. Region 15 has three significantly covarying base-pairs at the terminus of the downstream stem loop

600    (e-value < 0.05, **Fig. 8E**); Region 22 has 2 nucleotides with have one-sided variation in the most

601    upstream stem loop (e-value < 0.05, **Fig. 8F**); Region 30 has 3 covarying base-pairs at the very top of

602    the stem loop, and a single covarying pair at the bottom portion of the stem (e-value < 0.05, **Fig. 8G**).

603    Taken together, these results suggest the existence of stable, evolutionarily conserved structural

604    elements that merit subsequent functional analysis.

605

606    **Discussion**

607         Here we establish that the SARS-CoV-2 genomic RNA has an extraordinarily complex molecular

608    architecture, filled with elaborate secondary and tertiary structural features that persist in-vivo and

609    which are conserved through time, suggesting that this network of RNA secondary structural elements

610    plays a functional role in the virus lifecycle.  Indeed, the SARS-CoV-2 genome contains more well-

611    determined RNA structures than any virus studied to date, suggesting that its inherent "structuredness"

612    contributes in a unique way to viral fitness. This RNA secondary structural complexity is not just

613    confined to untranslated regions of the genome, as protein-coding sections of the SARS-CoV-2 open

614    reading frame are among the most well-structured regions.  Thus, as observed for HCV, coronavirus

615    reading frames experience evolutionary pressure that simultaneously shapes both protein sequence

616 and the surrounding RNA structures in which the proteins are encoded (a "code within the

617 code")(Pirakitikulr et al., 2016).  The secondary structure that we report is well-determined based on

618 available metrics in the field(Siegfried et al., 2014).  It is both a roadmap for navigating the vast RNA

619 landscape in coronaviruses, and a resource for orthogonal studies by others.  As such, the data

620 reported here are all publicly available for analysis and comparison by others

621 https://github.com/pylelab/SARS-CoV-2_SHAPE_MaP_structure.

622 　　Well-determined secondary structures of long RNA molecules are typically difficult to obtain *in-*

623 *vivo*(Mitchell et al., 2019, Leamy et al., 2016). They are usually derived from transcripts that have been

624 refolded and probed in-vitro, or from isolated cellular transcripts that have been stripped of cellular

625 components(Smola et al., 2015a, Siegfried et al., 2014). What is particularly surprising about this

626 SARS-CoV-2 study, and the high quality of the resulting secondary structure, is the fact that it was

627 entirely determined *in-vivo*, using infected cells that were treated directly with chemical probes. This

628 may be attributable to the fact that SARS-CoV-2 genomic RNA is so abundant in the infected cell,

629 ultimately becoming ~65% of the total cellular RNA(Kim et al., 2020).  With so much RNA material, it

630 becomes possible to maximize the signal to noise ratio in chemical probing experiments.  In addition,

631 the abundance of SARS-CoV-2 RNA may overwhelm the cell's ability to coat transcripts with

632 nonspecific RNA binding proteins, which can otherwise limit accessibility of chemical probes.  That said,

633 it will be interesting to compare the structure reported here with that obtained "*ex vivo*" (stripped of

634 protein), as that ∆SHAPE approach provides a useful way to flag possible protein binding sites(Smola

635 et al., 2015a).

636 　　The resulting experimental secondary structure provides new insights into known coronaviral

637 RNA motifs, and leads to the prediction of new ones that are likely to regulate viral function. The near

638 perfect structural homology of motifs at the 5' terminus for SARS-CoV-2 and other β-coronavirus

639 genomes suggests that the function of these upstream elements is conserved in coronaviruses

640 (reviewed in (Yang and Leibowitz, 2015)).  Furthermore, because our SARS-CoV-2 secondary structure

25

641  was determined *in-vivo*, our findings validate previous coronavirus structural models of 5'-elements, as

642  our data were obtained in a biologically relevant context.

643      Our SARS-CoV-2 secondary structure at the 3' viral terminus largely agrees with previous

644  studies on other β-coronavirus genomes (reviewed in(Yang and Leibowitz, 2015)). However, our model

645  of the 3' viral terminus deviates in one important way.  Neither the raw SHAPE reactivity data nor the

646  subsequent secondary structure prediction supports formation of a pseudoknot proposed between the

647  base of the BSL and SLI.  Indeed, the putative pseudoknot conformation is mutually exclusive with the

648  well-structured stem that we report at the base of the BSL. However, both conformations are proposed

649  to be essential in MHV(Goebel et al., 2004), so it is possible that the pseudoknot exists as a minority

650  conformation, or is transiently folded in SARS-CoV-2. Alternatively, because our structure represents

651  the first detailed description of a coronavirus 3'UTR structure *in-vivo*, it is possible this pseudoknot is

652  not present in other viruses.

653      Arguably the best-studied structural element in coronaviruses is the programmed ribosomal

654  frameshifting pseudoknot (PRF).  Required for proper replicase translation in all coronavirus family

655  members, the PRF adopts different conformations in the various coronaviruses, including three-

656  stemmed, two-stemmed, and kissing-loop pseudoknots (Baranov et al., 2005, Plant and Dinman, 2008).

657  The core of the SARS-CoV PRF, which shares an almost identical sequence with SARS-CoV-2, is

658  predicted to form a three-stem pseudoknot comprised of SLI, SL2, and a pseudoknot helix, with an

659  additional upstream attenuator stem that is poorly conserved in SARS-CoV-2(Kelly et al., 2020). Our

660  SHAPE reactivity and structure prediction are consistent with the existence of an attenuator stem, SL1,

661  and the pseudoknot. However, our data indicate that the region corresponding to SL2 is

662  conformationally flexible, adopting an SL2 stem with only a 20% probability. Consistent with our

663  reported distribution of structural isoforms, Kelly et al. use a reporter assay to suggest that the

664  frequency of successful frameshifting in SARS-CoV-2 is about 20%(Kelly et al., 2020), indicating that

665  the observed conformational variability of SL2 may be functional. Indeed, SL2 might function like a

666  switch: When SL2 is formed (~20% of the time), frameshifting occurs. When unfolded or forming base-

26

667    pairs with structures outside the PRF region, frameshifting would not occur. Further studies are

668    therefore required to explore the relationship between SL2 formation and SARS-CoV-2 frame-shifting

669    efficiency.

670          The study reported here provides a structure prediction for every single nucleotide in the SARS-

671    CoV-2 genome, enabling us to simultaneously interrogate both global and local features of genome

672    architecture. One can make two major observations about the global architecture the SARS-CoV-2

673    genome. First, this *in-vivo* derived, SHAPE-constrained model strongly agrees with the high double-

674    strand RNA content predicted from the entirely *in silico* model recently reported by our lab (Tavares et

675    al., 2020). Because the data herein were obtained *in-vivo,* this work confirms that the unusually high

676    double-strand content is maintained in a cellular context. Secondly, analysis of the experimental

677    secondary structure reveals that the SARS-CoV-2 genome has a shorter median base-pairing distance

678    when compared with other positive-sense RNA viral genomes, suggesting a role for extreme

679    compaction in the function of coronaviral genomes. Downstream analysis of synonymous mutation

680    rates suggests that global architectural features are conserved across β-coronaviruses. Considering

681    the exceptional size of these genomes, the high degree of dsRNA content may represent an

682    evolutionary strategy to enhance genome stability, as duplex RNA undergoes self-hydrolysis at a much

683    slower rate than single-stranded RNA and it is more resistant to cellular nucleases(Regulski and

684    Breaker, 2008, Wan et al., 2011). Interestingly, single-stranded regions in mRNA have been shown to

685    mediate phase separation at high cellular RNA concentrations(Van Treeck et al., 2018). Because

686    SARS-CoV-2 RNA is very abundant *in-vivo* (up to 65% of total cellular RNA content (Kim et al., 2020))

687    it is possible the high dsRNA content may provide a strategy to avoid phase separation during infection.

688    The preference for abundant locally folded, short stem-loop structures in β-coronavirus genomes may

689    also provide a conserved strategy for innate immune evasion.  Pattern recognition receptors such as

690    MDA5(Dias Junior et al., 2019) and ADAR modification(Nishikura, 2010) proteins recognize long RNA

691    duplexes as part of host defense processes, which could obviously be avoided by keeping duplex

692    lengths short.

27

693    Analysis of local features within the genome pinpoints 40 well-folded regions within the SARS-

694    CoV-2 orf1ab region. Of these 40 regions, at least five are conserved across all β-coronaviruses and

695    four are sarbecovirus specific. Four of the nine regions (Region 25, Region 34, Region 22, Region 24)

696    span boundaries between non-structural proteins, which may have relevance for polyprotein translation.

697    Previous studies have shown that RNA secondary structures can slow the rate of ribosome

698    translocation(Chen et al., 2013) and ribosome stalling is known to be important for proper folding of

699    nascent polypeptides(Collart and Weiss, 2020). Conserved, well-folded RNA structures at protein

700    domain boundaries may therefore slow or stall translocating ribosomes, thus allowing individual non-

701    structural proteins in the large Orf1a and Orf1ab poly-proteins to fold into their native conformations.

702    Intriguingly, three of the nine well-folded regions (Region 15, Region 24, Region 25) contain

703    complex, multi-helix junctions, or structural hubs. This is significant because multi-helix junctions often

704    comprise the core of RNA tertiary structures, like group II self-splicing introns, riboswitches and other

705    regulatory elements.  Because these elements are likely to contain well-defined pockets, they often

706    bind specifically to small molecules, and therefore serve as possible drug targets (Warner et al., 2018,

707    Hewitt et al., 2019, Fedorova et al., 2018).

708    One important cautionary observation from our work is the poor correlation of SHAPE

709    reactivities between two *in-vivo* biological replicates for regions encoding the subgenomic RNAs.

710    Previous *in silico* work from our lab has shown that individual subgenomic RNAs (sgRNAs), such as the

711    N sgRNA, fold differently than the corresponding regions in the genomic RNA due to differences in

712    upstream sequence context(Tavares et al., 2020). Though our tiled-amplicon design affords

713    sequencing coverage for the entire SARS-CoV-2 genome, it precludes deconvolution of reactivity

714    signals for regions shared between genomic- and subgenomic RNAs. This underscores the need for

715    methodological innovations that accurately assess the structural content specific to subgenomic RNA

716    molecules.  Absent such methodological advances, we caution others when interpreting reactivities

717    from the subgenomic region.

28

718     The *in-vivo*-determined SARS-CoV2 secondary structure present here provides a roadmap for

719     functional studies of the SARS-CoV2 genome and insights into mechanisms of the SARS-CoV-2 life

720     cycle.  Evolutionary support for consensus model across β-coronaviruses hints at conserved strategies

721     for genome stability, translation fidelity, and innate immune evasion. Finally, the identification of

722     individual well-folded regions conserved across β-coronaviruses, and within the sarbecovirus subgenus,

723     provide potential targets for the study of regulatory elements, and the search for much-needed

724     therapeutically active small molecules.

725

726     **Acknowledgments**

734

735     **Author Contributions**

736     N.C.H., H.W., and C.W. conducted experiments. N.C.H., H.W., C.W., and A.M.P. designed experiments.

737     N.C.H., H.W., R.C.A.T, A.M.P. wrote the paper.

738     **Declaration of Interests**

739

740     A patent application on MarathonRT has been filed by Yale University.

741

742     **Figure Legends**

29

743    **Figure 1**. Tiled-amplicon *in-vivo* SHAPE-MaP workflow yields high quality data for SARS-CoV-2 structure

744    prediction. **A)** Workflow of *in-vivo* SHAPE-MaP probing of full-length SARS-CoV-2 genomic RNA. **B)**

745    Mutation rates for two biological replicates confirm genomic RNA was successfully modified with NAI

746    electrophile. The boxes represent the interquartile range (IQR) of each data-set, with the median value

747    indicated by a line, average value indicated by a "x". Tukey-style whiskers extend 1.5 x IQR beyond each

748    box. Values outside this range are not shown. ****p<0.0001 by equal variance unpaired student t test.

749

750    **Figure 2.** *De novo* full-length structure prediction of SARS-CoV-2 genomic RNA identifies conserved

751    functional elements at the 5' and 3' viral termini. **A)** Consensus structure prediction for the 5' terminus

752    of SARS-CoV-2, colored by SHAPE Reactivity. Functional domains are labeled, including TRS sequence,

753    start codon of uORF, and start codon of Orf1a (indicated by black, grey, and green lines, respectively).

754    Inset – mapping of SHAPE reactivity data to single- and double-stranded regions, data are plotted with

755    a line indicating the median, and whiskers indicating the standard deviation **B)** Structure prediction for

756    the 3' terminus of SARS-CoV-2, colored by SHAPE reactivity. Functional domains are labeled. The

757    putative pseudoknot is indicated by solid black lines. Locations of the octanucleotide motif (ONM),

758    hypervariable region (HVR) and S2M are indicated by black lines.  Inset – mapping of SHAPE reactivity

759    to single- and double-stranded regions. Data are plotted with a line indicating the median, and

760    whiskers indicating the standard deviation ****p<0.0001 by equal variance unpaired student t test.

761

762    **Figure 3.** Structure prediction of the programmed ribosomal frame-shifting (PRF) element suggests

763    conformational variability of Stem Loop 2. **A)** Dominant PRF structural architecture, predicted by

764    SuperFold, colored by relative SHAPE Reactivity from this study. AS = Attenuator Stem; HSS =

765 Heptanucleotide Slippery Sequence; SL1 = Stem Loop 1; dotted line indicates region reported to form

766 stem loop 2 (SL2) or to form long-range interactions outside the PRF region in the SuperFold predition;

767 SL3 = Stem Loop 3; Red lines indicate pseudoknot interaction. **B)** Lower probability PRF conformation,

768 with fully-formed SL2, colored by relative SHAPE Reactivity **C)** Dominant PRF structure prediction

769 colored by relative Shannon entropy, labeled as in Panel A. **D)** Base-pairing probability for alternate SL2

770 conformation.  Each dot represents a base pair in SL2.  A base-pairing probability of 0.25 indicates a

771 25% probability of pairing for the indicated nucleotide.

772

773 **Figure 4.** Full-length genome structure prediction of SARS-CoV-2 Orf1ab reveals a network of well-

774 folded regions. **A)** Analysis of Shannon Entropy and SHAPE reactivities reveals 40 highly structured,

775 well-determined domains in Orf1ab. Nucleotide coordinates are indicated on the x-axis and numbered

776 in 1000 nucleotide intervals. Local median SHAPE reactivity and local median Shannon Entropy are

777 indicated by blue and orange lines, respectively. Well-folded regions are shaded with grey boxes. Arc

778 plots for all base-pairing interactions predicted by the structural model are shown beneath the local

779 SHAPE and Shannon entropy windows, corresponding to the genomic coordinates indicated on the x-

780 axis. The 5'UTR and non-structural protein (Nsp) domains are indicated by colored bars underneath arc

781 plot diagrams. **B)** Representative secondary structure predictions of two regions extracted from the

782 full-length consensus structure generated for the SARS-CoV-2 genome, with Nsp identity and genomic

783 position indicated.

784

785 **Figure 5.** Full-length genome structure prediction of SARS-CoV-2 Orf1ab reveals a unique genome

786 architecture.  **A)** Regions encoding individual non-structural protein (Nsp) domains have comparable

31

787    overall double-stranded RNA content (indicated by grey bars), but they do not adopt equally well-

788    folded substructures (indicated by black bars). A dotted line at 50% nucleotide content has been added

789    for clarity. **B)** SARS-CoV-2 has a shorter median base-pairing distance when compared to median base-

790    pairing distance in previously reported, full-length genome structures for two other positive-sense RNA

791    viruses (Mauger, *et. al.*, 2015; Dethoff, *et. al.,* 2018). Data are presented in Tukey-style box and

792    whiskers plot as described in Fig. 1B.  Asterisk definitions are below.   **C)** SARS-CoV-2 has a shorter

793    median base-pairing distance across well-folded regions of RNA when compared to those identified in

794    HIV (Siegfried, *et al.*, 2014). Data are presented as in B). *p<0.05, ****p<0.0001 by equal variance

795    unpaired student t test.

796

797    **Figure 6.** Structure-dependent variations in synonymous mutation rates suggest that all β-

798    coronaviruses have highly structured genomes (high BPC). **A)** Synonymous mutation rates calculated

799    across all β-coronaviruses for single- and double-stranded nucleotides of Orf1ab. Data are presented in

800    Tukey-style box and whiskers plot as described in Fig. 1B.  **B)** Non-synonymous mutation rates

801    calculated across all β-coronaviruses for single- and double-stranded nucleotides of Orf1ab. Data are

802    presented as in (A).  **C)** Comparison of synonymous mutation rates for single- and double-stranded

803    nucleotides within individual protein domains, calculated across all β-coronaviruses. Data are

804    presented as in (A). n.s. not significant,*p<0.05, ***p<0.001 ****p<0.0001 by equal variance unpaired

805    student t test.

806

807    **Figure 7.** Analysis of synonymous mutation rates within individual well-folded regions of the SARS-CoV-

808    2 genome identifies four regions that appear to be conserved across β-coronaviruses. **A)** Schematic of

809 well-folded regions in SARS-COV2 genome supported by Synonymous mutation rate analysis in β-

810 coronaviruses. **B)** Synonymous mutation rate separated by stranded-ness in four individual well-folded

811 regions. Data are plotted with a line indicating the median, and whiskers indicating the interquartile

812 range central. *p<0.05, **p<0.01 by equal variance unpaired student t test.  **C), D), E), F)** RNA

813 secondary structure diagrams of four well-folded regions supported by analysis of synonymous

814 mutation rates, colored by SHAPE reactivities, with genomic coordinates indicated below and in (A).

815

816 **Figure 8.** Analysis of synonymous mutation rates and covariation within individual regions of the SARS-

817 CoV-2 genome pinpoints five regions that are conserved only within the sarbecovirus

818 subgenus. **A)** Schematic of well-folded regions in the SARS-COV2 genome supported by Synonymous

819 mutation rate analysis in the sarbecovirus subgenus. **B)** Synonymous mutation rate separated by

820 stranded-ness in five individual well-folded regions. Data are plotted with a line indicating the median,

821 and whiskers indicating the interquartile range central . *p<0.05, **p<0.01  by equal variance unpaired

822 student t test.  **C),D)** RNA secondary structures of two well-folded regions colored by SHAPE

823 reactivity **E), F), G)**RNA secondary structure diagrams of three well-folded regions supported by both

824 synonymous mutation rate analysis and covariation in sarbecoviruses, colored by SHAPE reactivities.

825 Green boxes indicate significant covariation base pairs tested by Rscape-RAFSp(e-value<0.05).

826 Consensus nucleotides are colored by relative degree of sequence conservation within the alignment

827 (75% identity in gray, 90% identify in black, 97% identity in red). Individual nucleotides are represented

828 by circles according to their positional conservation and percentage occupancy thresholds (50%

829 occupancy in white, 75% occupancy in grey, 90% occupancy in black, 97% occupancy in red). Multiple

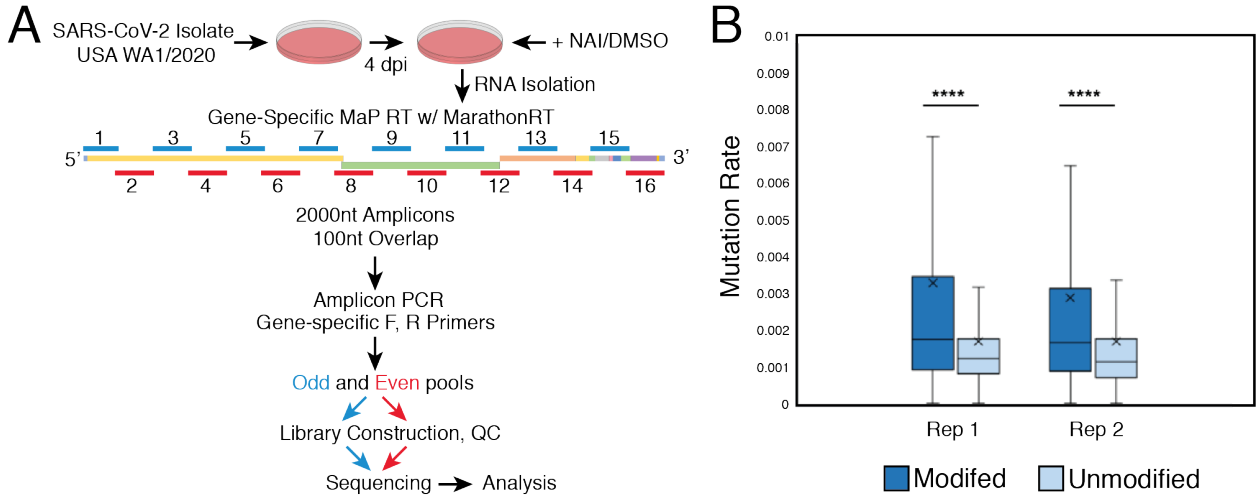830 sequence alignment files are provided in supplementary materials.

831 **Figure 1**
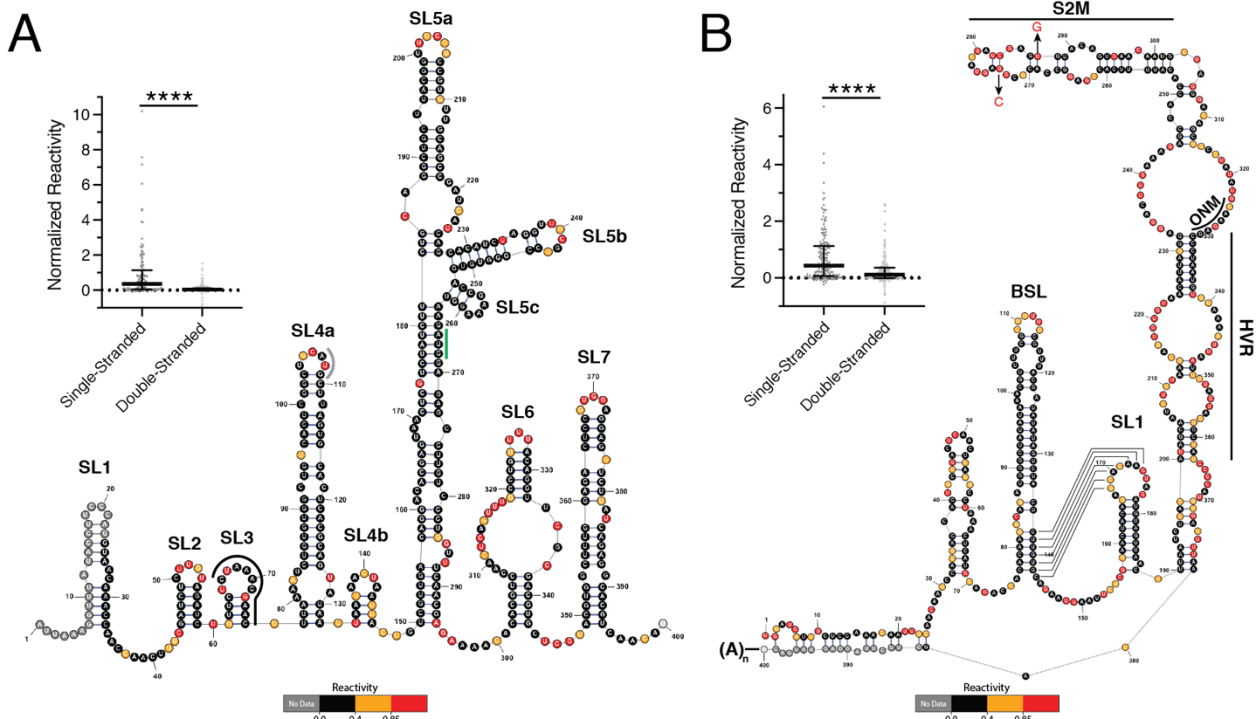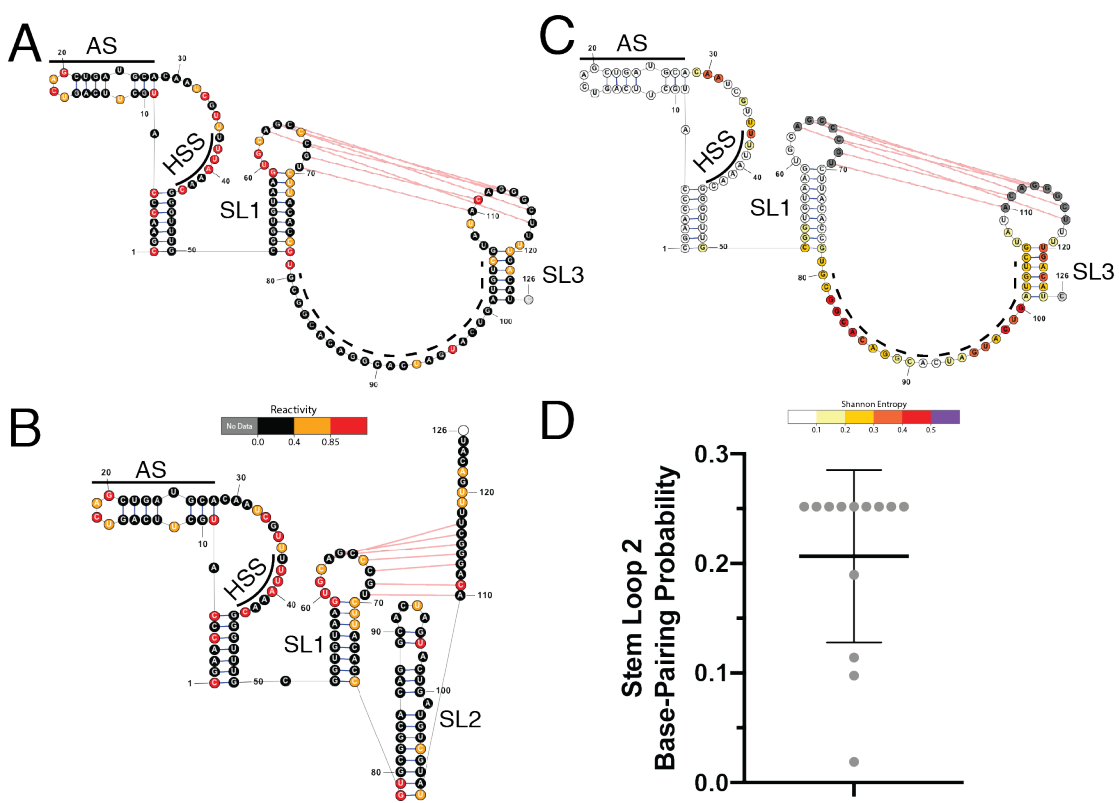


832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850 **Figure 2**

865    **Figure 3**



866

867

868

869

870

871

872

873

874

875

876

877

878

879 **Figure 4**



880

881

882

883

884    **Figure 5**



885

886

887

888

889

890

891

892

893

894

895

896

897 **Figure 6**

**Figure 7**



C Region 23
(11221nt-11470nt)

D Region 25
(12230nt-12686nt)

E Region 34
(16114nt-16260nt)

F Region 36
(17854nt-17938nt)

Reactivity

40

**Figure 8**



Region 24 (7717nt-8230nt)

Region 27 (10798nt-11039nt)

Region 15 (11552nt-11908nt)
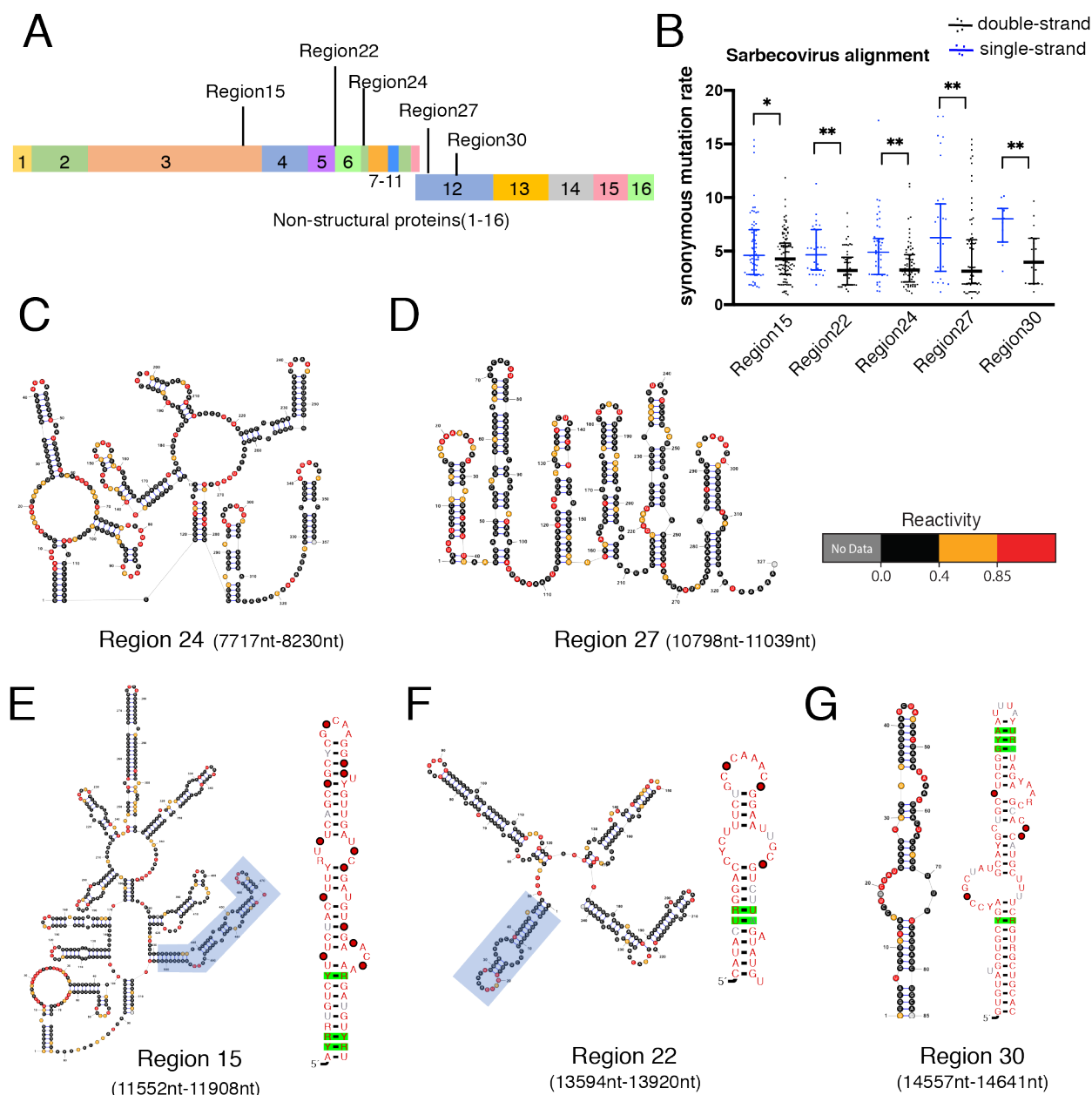
Region 22 (13594nt-13920nt)

Region 30 (14557nt-14641nt)

## References

922

923 ADAMS, R. L., PIRAKITIKULR, N. & PYLE, A. M. 2017. Functional RNA structures throughout the Hepatitis
924         C Virus genome. *Curr Opin Virol,* 24**,** 79-86.
925 ANDREWS, R. J., PETERSON, J. M., HANIFF, H. S., CHEN, J., WILLIAMS, C., GREFE, M., DISNEY, M. D. &
926         MOSS, W. N. 2020. An in silico map of the SARS-CoV-2 RNA Structurome. *bioRxiv.*
927 ASSIS, R. 2014. Strong epistatic selection on the RNA secondary structure of HIV. *PLoS Pathog,* 10**,**
928         e1004363.
929 BARANOV, P. V., HENDERSON, C. M., ANDERSON, C. B., GESTELAND, R. F., ATKINS, J. F. & HOWARD, M.
930         T. 2005. Programmed ribosomal frameshifting in decoding the SARS-CoV genome. *Virology,* 332**,**
931         498-510.
932 BARROWS, N. J., CAMPOS, R. K., LIAO, K. C., PRASANTH, K. R., SOTO-ACOSTA, R., YEH, S. C., SCHOTT-
933         LERNER, G., POMPON, J., SESSIONS, O. M., BRADRICK, S. S. & GARCIA-BLANCO, M. A. 2018.
934         Biochemistry and Molecular Biology of Flaviviruses. *Chem Rev,* 118**,** 4448-4482.
935 BENSON, D. A., CAVANAUGH, M., CLARK, K., KARSCH-MIZRACHI, I., OSTELL, J., PRUITT, K. D. & SAYERS,
936         E. W. 2018. GenBank. *Nucleic Acids Res,* 46**,** D41-D47.
937 BUSAN, S. & WEEKS, K. M. 2018. Accurate detection of chemical modifications in RNA by mutational
938         profiling (MaP) with ShapeMapper 2. *RNA,* 24**,** 143-148.
939 CERAOLO, C. & GIORGI, F. M. 2020. Genomic variance of the 2019-nCoV coronavirus. *J Med Virol,* 92**,**
940         522-528.
941 CHEN, C., ZHANG, H., BROITMAN, S. L., REICHE, M., FARRELL, I., COOPERMAN, B. S. & GOLDMAN, Y. E.
942         2013. Dynamics of translation by single ribosomes through mRNA secondary structures. *Nat*
943         *Struct Mol Biol,* 20**,** 582-8.
944 CHEN, S. C. & OLSTHOORN, R. C. 2010. Group-specific structural features of the 5'-proximal sequences
945         of coronavirus genomic RNAs. *Virology,* 401**,** 29-41.
946 CHO, C. P., LIN, S. C., CHOU, M. Y., HSU, H. T. & CHANG, K. Y. 2013. Regulation of programmed
947         ribosomal frameshifting by co-translational refolding RNA hairpins. *PLoS One,* 8**,** e62283.
948 CLYDE, K. & HARRIS, E. 2006. RNA secondary structure in the coding region of dengue virus type 2
949         directs translation start codon selection and is required for viral replication. *J Virol,* 80**,** 2170-82.
950 COLLART, M. A. & WEISS, B. 2020. Ribosome pausing, a dangerous necessity for co-translational events.
951         *Nucleic Acids Res,* 48**,** 1043-1055.
952 DE WIT, E., VAN DOREMALEN, N., FALZARANO, D. & MUNSTER, V. J. 2016. SARS and MERS: recent
953         insights into emerging coronaviruses. *Nat Rev Microbiol,* 14**,** 523-34.
954 DETHOFF, E. A., BOERNEKE, M. A., GOKHALE, N. S., MUHIRE, B. M., MARTIN, D. P., SACCO, M. T.,
955         MCFADDEN, M. J., WEINSTEIN, J. B., MESSER, W. B., HORNER, S. M. & WEEKS, K. M. 2018.
956         Pervasive tertiary structure in the dengue virus RNA genome. *Proc Natl Acad Sci U S A,* 115**,**
957         11513-11518.
958 DIAS JUNIOR, A. G., SAMPAIO, N. G. & REHWINKEL, J. 2019. A Balancing Act: MDA5 in Antiviral
959         Immunity and Autoinflammation. *Trends Microbiol,* 27**,** 75-85.
960 DONG, E., DU, H. & GARDNER, L. 2020. An interactive web-based dashboard to track COVID-19 in real
961         time. *Lancet Infect Dis,* 20**,** 533-534.
962 FEDOROVA, O., JAGDMANN, G. E., JR., ADAMS, R. L., YUAN, L., VAN ZANDT, M. C. & PYLE, A. M. 2018.
963         Small molecules that target group II introns are potent antifungal agents. *Nat Chem Biol,* 14**,**
964         1073-1078.

965 FRICKE, M., DUNNES, N., ZAYAS, M., BARTENSCHLAGER, R., NIEPMANN, M. & MARZ, M. 2015.
966     Conserved RNA secondary structures and long-range interactions in hepatitis C viruses. *RNA,* 21,
967     1219-32.
968 FRIEBE, P. & BARTENSCHLAGER, R. 2009. Role of RNA structures in genome terminal sequences of the
969     hepatitis C virus for replication and assembly. *J Virol,* 83, 11989-95.
970 GOEBEL, S. J., HSUE, B., DOMBROWSKI, T. F. & MASTERS, P. S. 2004. Characterization of the RNA
971     components of a putative molecular switch in the 3' untranslated region of the murine
972     coronavirus genome. *J Virol,* 78, 669-82.
973 GOEBEL, S. J., MILLER, T. B., BENNETT, C. J., BERNARD, K. A. & MASTERS, P. S. 2007. A hypervariable
974     region within the 3' cis-acting element of the murine coronavirus genome is nonessential for
975     RNA synthesis but affects pathogenesis. *J Virol,* 81, 1274-87.
976 GUO, L. T., ADAMS, R. L., WAN, H., HUSTON, N. C., POTAPOVA, O., OLSON, S., GALLARDO, C. M.,
977     GRAVELEY, B. R., TORBETT, B. E. & PYLE, A. M. 2020. Sequencing and Structure Probing of Long
978     RNAs Using MarathonRT: A Next-Generation Reverse Transcriptase. *J Mol Biol,* 432, 3338-3352.
979 HAJDIN, C. E., BELLAOUSOV, S., HUGGINS, W., LEONARD, C. W., MATHEWS, D. H. & WEEKS, K. M. 2013.
980     Accurate SHAPE-directed RNA secondary structure modeling, including pseudoknots. *Proc Natl*
981     *Acad Sci U S A,* 110, 5498-503.
982 HEWITT, W. M., CALABRESE, D. R. & SCHNEEKLOTH, J. S., JR. 2019. Evidence for ligandable sites in
983     structured RNA throughout the Protein Data Bank. *Bioorg Med Chem,* 27, 2253-2260.
984 KANG, H., FENG, M., SCHROEDER, M. E., GIEDROC, D. P. & LEIBOWITZ, J. L. 2006. Putative cis-acting
985     stem-loops in the 5' untranslated region of the severe acute respiratory syndrome coronavirus
986     can substitute for their mouse hepatitis virus counterparts. *J Virol,* 80, 10600-14.
987 KELLY, J. A., OLSON, A. N., NEUPANE, K., MUNSHI, S., SAN EMETERIO, J., POLLACK, L., WOODSIDE, M. T.
988     & DINMAN, J. D. 2020. Structural and functional conservation of the programmed -1 ribosomal
989     frameshift signal of SARS coronavirus 2 (SARS-CoV-2). *J Biol Chem*.
990 KIM, D., LEE, J. Y., YANG, J. S., KIM, J. W., KIM, V. N. & CHANG, H. 2020. The Architecture of SARS-CoV-2
991     Transcriptome. *Cell,* 181, 914-921 e10.
992 KORBIE, D. J. & MATTICK, J. S. 2008. Touchdown PCR for increased specificity and sensitivity in PCR
993     amplification. *Nat Protoc,* 3, 1452-6.
994 LAI, D., PROCTOR, J. R., ZHU, J. Y. & MEYER, I. M. 2012. R-CHIE: a web server and R package for
995     visualizing RNA secondary structures. *Nucleic Acids Res,* 40, e95.
996 LAN, J., GE, J., YU, J., SHAN, S., ZHOU, H., FAN, S., ZHANG, Q., SHI, X., WANG, Q., ZHANG, L. & WANG, X.
997     2020a. Structure of the SARS-CoV-2 spike receptor-binding domain bound to the ACE2 receptor.
998     *Nature,* 581, 215-220.
999 LAN, T. C. T., ALLAN, M., MALSICK, L., KHANDWALA, S., NYEO, S. Y., BATHE, M., GRIFFITHS, A. &
1000     ROUSKIN, S. 2020b. Structure of the full SARS-CoV-2 RNA genome in infected cells. *bioRxiv*.
1001 LEAMY, K. A., ASSMANN, S. M., MATHEWS, D. H. & BEVILACQUA, P. C. 2016. Bridging the gap between
1002     in vitro and in vivo RNA folding. *Q Rev Biophys,* 49, e10.
1003 LEE, C. W., LI, L. & GIEDROC, D. P. 2011. The solution structure of coronaviral stem-loop 2 (SL2) reveals
1004     a canonical CUYG tetraloop fold. *FEBS Lett,* 585, 1049-53.
1005 LI, L., KANG, H., LIU, P., MAKKINJE, N., WILLIAMSON, S. T., LEIBOWITZ, J. L. & GIEDROC, D. P. 2008.
1006     Structural lability in stem-loop 1 drives a 5' UTR-3' UTR interaction in coronavirus replication. *J*
1007     *Mol Biol,* 377, 790-803.

LI, P., WEI, Y., MEI, M., TANG, L., SUN, L., HUANG, W., ZHOU, J., ZOU, C., ZHANG, S., QIN, C. F., JIANG, T., DAI, J., TAN, X. & ZHANG, Q. C. 2018. Integrative Analysis of Zika Virus Genome RNA Structure Reveals Critical Determinants of Viral Infectivity. *Cell Host Microbe,* 24**,** 875-886 e5.

LU, Z. J. & MATHEWS, D. H. 2008. OligoWalk: an online siRNA design tool utilizing hybridization thermodynamics. *Nucleic Acids Res,* 36**,** W104-8.

MACFADDEN, A., O'DONOGHUE, Z., SILVA, P., CHAPMAN, E. G., OLSTHOORN, R. C., STERKEN, M. G., PIJLMAN, G. P., BREDENBEEK, P. J. & KIEFT, J. S. 2018. Mechanism and structural diversity of exoribonuclease-resistant RNA structures in flaviviral RNAs. *Nat Commun,* 9**,** 119.

MADHUGIRI, R., KARL, N., PETERSEN, D., LAMKIEWICZ, K., FRICKE, M., WEND, U., SCHEUER, R., MARZ, M. & ZIEBUHR, J. 2018. Structural and functional conservation of cis-acting RNA elements in coronavirus 5'-terminal genome regions. *Virology,* 517**,** 44-55.

MAIER, H. J., BICKERTON, E. & BRITTON, P. 2015. *Coronaviruses : methods and protocols,* New York, Humana Press ; Springer.

MANFREDONIA, I., NITHIN, C., PONCE-SALVATIERRA, A., GHOSH, P., WIRECKI, T., MARINUS, T., OGANDO, N. S., SNIDER, E. J., VAN HEMERT, M. J., BUJNICKI, J. M. & INCARNATO, D. 2020. Genome-wide mapping of therapeutically-relevant SARS-CoV-2 RNA structures. *bioRxiv*.

MAUGER, D. M., GOLDEN, M., YAMANE, D., WILLIFORD, S., LEMON, S. M., MARTIN, D. P. & WEEKS, K. M. 2015. Functionally conserved architecture of hepatitis C virus RNA genomes. *Proc Natl Acad Sci U S A,* 112**,** 3692-7.

MCMULLAN, L. K., GRAKOUI, A., EVANS, M. J., MIHALIK, K., PUIG, M., BRANCH, A. D., FEINSTONE, S. M. & RICE, C. M. 2007. Evidence for a functional RNA element in the hepatitis C virus core gene. *Proc Natl Acad Sci U S A,* 104**,** 2879-84.

MITCHELL, D., 3RD, ASSMANN, S. M. & BEVILACQUA, P. C. 2019. Probing RNA structure in vivo. *Curr Opin Struct Biol,* 59**,** 151-158.

MURRELL, B., MOOLA, S., MABONA, A., WEIGHILL, T., SHEWARD, D., KOSAKOVSKY POND, S. L. & SCHEFFLER, K. 2013. FUBAR: a fast, unconstrained bayesian approximation for inferring selection. *Mol Biol Evol,* 30**,** 1196-205.

NISHIKURA, K. 2010. Functions and regulation of RNA editing by ADAR deaminases. *Annu Rev Biochem,* 79**,** 321-49.

PIRAKITIKULR, N., KOHLWAY, A., LINDENBACH, B. D. & PYLE, A. M. 2016. The Coding Region of the HCV Genome Contains a Network of Regulatory RNA Structures. *Mol Cell,* 62**,** 111-20.

PLANT, E. P. & DINMAN, J. D. 2008. The role of programmed-1 ribosomal frameshifting in coronavirus propagation. *Front Biosci,* 13**,** 4873-81.

PLANT, E. P., PEREZ-ALVARADO, G. C., JACOBS, J. L., MUKHOPADHYAY, B., HENNIG, M. & DINMAN, J. D. 2005. A three-stemmed mRNA pseudoknot in the SARS coronavirus frameshift signal. *PLoS Biol,* 3**,** e172.

PLANT, E. P., SIMS, A. C., BARIC, R. S., DINMAN, J. D. & TAYLOR, D. R. 2013. Altering SARS coronavirus frameshift efficiency affects genomic and subgenomic RNA production. *Viruses,* 5**,** 279-94.

RAMAN, S., BOUMA, P., WILLIAMS, G. D. & BRIAN, D. A. 2003. Stem-loop III in the 5' untranslated region is a cis-acting element in bovine coronavirus defective interfering RNA replication. *J Virol,* 77**,** 6720-30.

RANGAN, R., ZHELUDEV, I. N. & DAS, R. 2020. RNA genome conservation and secondary structure in SARS-CoV-2 and SARS-related viruses: a first look. *RNA*.

051    RANWEZ, V., DOUZERY, E. J. P., CAMBON, C., CHANTRET, N. & DELSUC, F. 2018. MACSE v2: Toolkit for
052        the Alignment of Coding Sequences Accounting for Frameshifts and Stop Codons. *Mol Biol Evol,*
053        35, **2582-2584.**

054    REGULSKI, E. E. & BREAKER, R. R. 2008. In-line probing analysis of riboswitches. *Methods Mol Biol,* 419,
055        53-67.

056    REUTER, J. S. & MATHEWS, D. H. 2010. RNAstructure: software for RNA secondary structure prediction
057        and analysis. *BMC Bioinformatics,* 11, **129.**

058    RIVAS, E., CLEMENTS, J. & EDDY, S. R. 2017. A statistical test for conserved RNA structure shows lack of
059        evidence for structure in lncRNAs. *Nat Methods,* 14, **45-48.**

060    ROBERTSON, M. P., IGEL, H., BAERTSCH, R., HAUSSLER, D., ARES, M., JR. & SCOTT, W. G. 2005. The
061        structure of a rigorously conserved RNA element within the SARS virus genome. *PLoS Biol,* 3, e5.

062    ROUSKIN, S., ZUBRADT, M., WASHIETL, S., KELLIS, M. & WEISSMAN, J. S. 2014. Genome-wide probing of
063        RNA structure reveals active unfolding of mRNA structures in vivo. *Nature,* 505, **701-5.**

064    SIEGFRIED, N. A., BUSAN, S., RICE, G. M., NELSON, J. A. & WEEKS, K. M. 2014. RNA motif discovery by
065        SHAPE and mutational profiling (SHAPE-MaP). *Nat Methods,* 11, **959-65.**

066    SIMMONDS, P. & SMITH, D. B. 1999. Structural constraints on RNA virus evolution. *J Virol,* 73, **5787-94.**

067    SIMON, L. M., MORANDI, E., LUGANINI, A., GRIBAUDO, G., MARTINEZ-SOBRIDO, L., TURNER, D. H.,
068        OLIVIERO, S. & INCARNATO, D. 2019. In vivo analysis of influenza A mRNA secondary structures
069        identifies critical regulatory motifs. *Nucleic Acids Res,* 47, **7003-7017.**

070    SMOLA, M. J., CALABRESE, J. M. & WEEKS, K. M. 2015a. Detection of RNA-Protein Interactions in Living
071        Cells with SHAPE. *Biochemistry,* 54, **6867-75.**

072    SMOLA, M. J., CHRISTY, T. W., INOUE, K., NICHOLSON, C. O., FRIEDERSDORF, M., KEENE, J. D., LEE, D. M.,
073        CALABRESE, J. M. & WEEKS, K. M. 2016. SHAPE reveals transcript-wide interactions, complex
074        structural domains, and protein interactions across the Xist lncRNA in living cells. *Proc Natl Acad*
075        *Sci U S A,* 113, **10322-7.**

076    SMOLA, M. J., RICE, G. M., BUSAN, S., SIEGFRIED, N. A. & WEEKS, K. M. 2015b. Selective 2'-hydroxyl
077        acylation analyzed by primer extension and mutational profiling (SHAPE-MaP) for direct,
078        versatile and accurate RNA structure analysis. *Nat Protoc,* 10, **1643-69.**

079    TAVARES, R., MAHADESHWAR, G. & PYLE, A. M. 2020. The global and local distribution of RNA
080        structure throughout the SARS-CoV-2 genome. *bioRxiv.*

081    TAVARES, R. C. A., PYLE, A. M. & SOMAROWTHU, S. 2019. Phylogenetic Analysis with Improved
082        Parameters Reveals Conservation in lncRNA Structures. *J Mol Biol,* 431, **1592-1603.**

083    TUPLIN, A., WOOD, J., EVANS, D. J., PATEL, A. H. & SIMMONDS, P. 2002. Thermodynamic and
084        phylogenetic prediction of RNA secondary structures in the coding region of hepatitis C virus.
085        *RNA,* 8, **824-41.**

086    VAN TREECK, B., PROTTER, D. S. W., MATHENY, T., KHONG, A., LINK, C. D. & PARKER, R. 2018. RNA self-
087        assembly contributes to stress granule formation and defining the stress granule transcriptome.
088        *Proc Natl Acad Sci U S A,* 115, **2734-2739.**

089    WAN, Y., KERTESZ, M., SPITALE, R. C., SEGAL, E. & CHANG, H. Y. 2011. Understanding the transcriptome
090        through RNA structure. *Nat Rev Genet,* 12, **641-55.**

091    WAN, Y., SHANG, J., GRAHAM, R., BARIC, R. S. & LI, F. 2020. Receptor Recognition by the Novel
092        Coronavirus from Wuhan: an Analysis Based on Decade-Long Structural Studies of SARS
093        Coronavirus. *J Virol,* 94.

WARNER, K. D., HAJDIN, C. E. & WEEKS, K. M. 2018. Principles for targeting RNA with drug-like small molecules. *Nat Rev Drug Discov,* 17**,** 547-558.

WATERHOUSE, A. M., PROCTER, J. B., MARTIN, D. M., CLAMP, M. & BARTON, G. J. 2009. Jalview Version 2--a multiple sequence alignment editor and analysis workbench. *Bioinformatics,* 25**,** 1189-91.

YANG, D. & LEIBOWITZ, J. L. 2015. The structure and functions of coronavirus genomic 3' and 5' ends. *Virus Res,* 206**,** 120-33.

YANG, D., LIU, P., WUDECK, E. V., GIEDROC, D. P. & LEIBOWITZ, J. L. 2015. SHAPE analysis of the RNA secondary structure of the Mouse Hepatitis Virus 5' untranslated region and N-terminal nsp1 coding sequences. *Virology,* 475**,** 15-27.

YIN, W., MAO, C., LUAN, X., SHEN, D. D., SHEN, Q., SU, H., WANG, X., ZHOU, F., ZHAO, W., GAO, M., CHANG, S., XIE, Y. C., TIAN, G., JIANG, H. W., TAO, S. C., SHEN, J., JIANG, Y., JIANG, H., XU, Y., ZHANG, S., ZHANG, Y. & XU, H. E. 2020. Structural basis for inhibition of the RNA-dependent RNA polymerase from SARS-CoV-2 by remdesivir. *Science,* 368**,** 1499-1504.

YOU, S., STUMP, D. D., BRANCH, A. D. & RICE, C. M. 2004. A cis-acting replication element in the sequence encoding the NS5B RNA-dependent RNA polymerase is required for hepatitis C virus RNA replication. *J Virol,* 78**,** 1352-66.

ZHU, N., ZHANG, D., WANG, W., LI, X., YANG, B., SONG, J., ZHAO, X., HUANG, B., SHI, W., LU, R., NIU, P., ZHAN, F., MA, X., WANG, D., XU, W., WU, G., GAO, G. F., TAN, W., CHINA NOVEL CORONAVIRUS, I. & RESEARCH, T. 2020. A Novel Coronavirus from Patients with Pneumonia in China, 2019. *N Engl J Med,* 382**,** 727-733.

ZUBRADT, M., GUPTA, P., PERSAD, S., LAMBOWITZ, A. M., WEISSMAN, J. S. & ROUSKIN, S. 2017. DMS-MaPseq for genome-wide or targeted RNA structure probing in vivo. *Nat Methods,* 14**,** 75-82.

ZUST, R., MILLER, T. B., GOEBEL, S. J., THIEL, V. & MASTERS, P. S. 2008. Genetic interactions between an essential 3' cis-acting RNA pseudoknot, replicase gene products, and the extreme 3' end of the mouse coronavirus genome. *J Virol,* 82**,** 1214-28.