

1 **Title:** Rapid, sensitive, full genome sequencing of Severe Acute Respiratory Syndrome Virus

2 Coronavirus 2 (SARS-CoV-2)

3 **Authors:** Clinton R Paden<sup>1</sup>, Ying Tao<sup>1</sup>, Krista Queen, Jing Zhang, Yan Li, Anna Uehara,

4 Suxiang Tong<sup>2</sup>

5 <sup>1</sup>These first authors contributed equally to this manuscript

6 <sup>2</sup>Corresponding author: sot1@cdc.gov

7 **Affiliations:**

8 Centers for Disease Control and Prevention, Atlanta, GA (all authors)

9 **Keywords:** COVID-19, SARS-CoV-2, coronavirus, genomics, High-Throughput Nucleotide  
10 Sequencing, Whole Genome Sequencing

11

## 12 **Abstract**

13           SARS-CoV-2 recently emerged, resulting a global pandemic. Rapid genomic information  
14 is critical to understanding transmission and pathogenesis. Here, we describe validated protocols  
15 for generating high-quality full-length genomes from primary samples. The first employs  
16 multiplex RT-PCR followed by MinION or MiSeq sequencing. The second uses singleplex,  
17 nested RT-PCR and Sanger sequencing.

18 In December 2019, SARS-CoV-2, the etiological agent of Coronavirus Disease 2019  
19 (COVID-19), emerged in Wuhan, China. Since then it has rapidly spread to the rest of the world  
20 (1-3). As of April 16, 2020, there have been 1,991,562 confirmed cases, including 130,885  
21 deaths, in 185 countries or regions (4).

22 Initial sequencing of SARS-CoV-2 showed limited genetic variation between cases, but  
23 did document specific changes that may be useful for understanding the source and transmission  
24 chains (5-8). Because SARS-CoV-2 has shown the capacity to spread rapidly and lead to a range  
25 of presentations in infected persons, from asymptomatic infection to mild, severe, or fatal  
26 disease, it is important to identify genetic variants in order to understand any changes in  
27 transmissibility, tropism, and pathogenicity. Sequence data can be used to inform decisions to  
28 better manage the spread of disease.

29 In this report, we describe the design and use of two PCR-based methods for sequencing  
30 SARS-CoV-2 clinical specimens. The first is a multiplex PCR panel followed by sequencing on  
31 either the Oxford Nanopore MinION or Illumina MiSeq. When coupled with MinION  
32 sequencing, the protocol can be implemented outside a traditional laboratory and can be  
33 completed in a single workday, similar to previous mobile genomic surveillance of Ebola and  
34 Zika virus outbreaks (9, 10). Additionally, we provide a complementary singleplex, nested PCR  
35 strategy, which improves sensitivity for samples with lower viral load and is compatible with  
36 Sanger sequencing.

## 37 **The Study**

38 On January 10, 2020, the first SARS-CoV-2 genome sequence was released online (11).  
39 That day, we designed two complementary panels of primers to amplify the virus genome for

40 sequencing. For one panel, we used the PRIMAL primer design tool (9) to design multiplex  
41 PCRs to amplify the genome in using only a few PCR reactions (Appendix). The final design  
42 consists of 6 pools of primers, targeting amplicon sizes of 550 base pairs (bp) with 100bp  
43 overlaps, to allow for sequencing on either the ONT MinION or Illumina MiSeq. For the second  
44 panel, we designed sets of primers to generate nested, tiling amplicons across the SARS-CoV-2  
45 genome (Appendix), for enhancing sensitivity in samples with lower viral loads. Each amplicon  
46 is 322-1030bp with an average overlap of 80bp. They are designed to be amplified and  
47 sequenced individually on Sanger instruments but may also be pooled for sequencing on next-  
48 generation sequencing platforms.

49 To determine the sensitivity of each sequencing strategy, we generated a set of six ten-  
50 fold serial dilutions of a SARS-CoV-2 isolate (12). Virus RNA was diluted into a constant  
51 background of A549 human cell line total nucleic acid (RNaseP  $C_T$  29). We quantitated each  
52 dilution using the CDC SARS-CoV-2 rRT-PCR for the N2 gene (13) (data not shown). The six  
53 dilutions spanned  $C_T$  values from 22-37, corresponding to ca. 2 to  $1.8 \times 10^5$  copies. We amplified  
54 triplicate samples at each dilution using the multiplex PCR pools. Next, we pooled, barcoded,  
55 and made libraries from each sample's amplicons using the ligation-based kit and PCR barcode  
56 expansion kit (methods in Appendix). MinION sequencing was performed on an R9.4.1 or R10.3  
57 flow cell until we obtained >1-2M raw reads. From those, 50-60% of reads could be  
58 demultiplexed (data not shown). Additionally, we sequenced these amplicons using the Illumina  
59 MiSeq for comparison (methods in Appendix).

60 For MinION sequencing, the reads were basecalled and analyzed using an in-house read  
61 mapping pipeline (detailed in Appendix). For samples with  $C_T \leq 29$ , we obtained >99% SARS-  
62 CoV-2 reads and >99% genome coverage at 20X depth, decreasing to an average of 93%

63 genome coverage at  $C_T$  33.2 and 48% at  $C_T$  35 (Figure 1A and 1B). Further, we were able to  
64 obtain full >20X genomes within the first 40-60 minutes of sequencing (Figure 1C).

65 Consensus accuracy, including SNPs and indels, is critical for determining coronavirus  
66 lineage and transmission networks. For high consensus level accuracy, we filtered reads based on  
67 length, mapped them to the reference sequence (RefSeq NC\_045512), trimmed primers based on  
68 position, and called variants with Medaka (<https://github.com/nanoporetech/medaka>) (details in  
69 Appendix). Each Medaka variant was filtered by coverage depth (>20X) and by the Medaka  
70 model-derived variant quality (>40). Here we used the variant quality score as a heuristic to filter  
71 remaining noise from the Medaka variants, compared to Sanger-derived sequences. After these  
72 steps, the data approaches 100% consensus accuracy (Table 1). Identical results were found  
73 using the R9.4.1 pore, through the  $C_T$  33.2 samples (data not shown). We noted larger deletions  
74 in some of the  $C_T$  33.2+ samples which likely reflect biases from limited copy numbers.

75 In the MiSeq data, we observed a similar trend in percent genome coverage at 100X  
76 depth, and a slightly lower percent mapped reads, compared to Nanopore data (Figure 1A and  
77 B). Increased read depth using the MiSeq potentially allows increased sample throughput,  
78 however the number of available dual unique barcodes limits actual throughput.

79 For the nested, singleplex PCR panel, we amplified the same serial dilutions with each  
80 nested primer set (methods in Appendix). The endpoint dilution for full genome coverage is  
81 approximately  $C_T$  35 (Figure 1B). At the  $C_T$  37 dilution, we observed significant amplicon  
82 dropout—at this dilution, there are <10 copies of the genome on average per reaction.

83 These protocols enabled rapid sequencing of the initial clinical cases of SARS-CoV-2 in  
84 the United States. For these cases, we amplified the virus genome using the singleplex PCR

85 amplicons, sequencing them with both MinION and Sanger instruments to validate MinION  
86 consensus accuracy. The MinION produced full-length genomes in <20 minutes of sequencing,  
87 while Sanger data was available the following day.

88 We used the multiplex PCR strategy in subsequent SARS-CoV-2 clinical cases (n=167),  
89 ranging in  $C_T$  values from 15.7 to 40 (mean 28.8, median 29.1). In cases below  $C_T$  33, we  
90 observed an average of 99.02% specific reads and 99.2% genome coverage at >20X depth  
91 (Figures 2A and 2B). Between  $C_T$  30-33, genome coverage varied by sample, and declined  
92 dramatically at higher  $C_T$  values, analogous to the isolate validation data. For these samples, we  
93 multiplexed 20-40 barcoded samples per flowcell. Enough data is obtained with 60 minutes of  
94 MinION sequencing for most samples, though for higher titer samples 10-20 minutes of  
95 sequencing is sufficient (Figure 2C).

96 Up-to-date primer sequences, protocols, and analysis scripts are found at  
97 [https://github.com/CDCgov/SARS-CoV-2\\_Sequencing/tree/master/protocols/CDC-](https://github.com/CDCgov/SARS-CoV-2_Sequencing/tree/master/protocols/CDC-Comprehensive)  
98 [Comprehensive](https://github.com/CDCgov/SARS-CoV-2_Sequencing/tree/master/protocols/CDC-Comprehensive). Data from this study is deposited in the NBCI SRA (BioProjects PRJNA622817  
99 and PRJNA610248).

## 100 **Conclusions**

101 Full genome sequencing is an indispensable tool in understanding emerging viruses. Here  
102 we present two validated PCR target-enrichment strategies that can be used with MinION,  
103 MiSeq, and Sanger platforms for sequencing SARS-CoV-2 clinical specimens. This ensures that  
104 most labs have access to one or more strategies.

105 The multiplex PCR strategy is effective at generating full genome sequences up to  $C_T$  33.  
106 The singleplex, nested PCR is effective up to  $C_T$  35, varying based on sample quality. The

107 turnaround time for the multiplex PCR MinION protocol is about 8 hours from nucleic acid to  
108 consensus sequence, compared to Sanger sequencing at about 14-18 hours (Table 2).  
109 Importantly, the multiplex PCR protocols offer an efficient, cost-effective, scalable system, and  
110 add little time and complexity as sample numbers increase (Table 2). The results from this study  
111 suggest multiplex PCR may be used effectively for routine sequencing, complemented by  
112 singleplex, nested PCR for low virus-titer samples and confirmation sequencing.

113

## 114 **Acknowledgments**

115 We would like to acknowledge the efforts of those in the Respiratory Viruses Branch at  
116 CDC who helped in organizing samples for this study, including Azaibi Tamin, Jennifer  
117 Harcourt, Natalie Thornburg, Shifaq Kamili, Xiaoyan Lu, and Stephen Lindstrom.

## 118 **Disclaimers**

119 The findings and conclusions in this report are those of the authors and do not necessarily  
120 represent the official position of the Centers for Disease Control and Prevention.

## 121 **Author Bio** (first author only, unless there are only 2 authors)

122 Clinton Paden is a virologist and bioinformatician in the CDC Pathogen Discovery Team,  
123 within the Division of Viral Diseases. His work is in identifying and characterizing novel and  
124 emerging pathogens.

125 **References**

- 126 1. Holshue ML, DeBolt C, Lindquist S, Lofy KH, Wiesman J, Bruce H, et al. First Case of 2019  
127 Novel Coronavirus in the United States. *The New England Journal of Medicine*.  
128 2020;382(10):929-36.
- 129 2. Patel A, Jernigan DB, Abdirizak F, Abedi G, Aggarwal S, Albina D, et al. Initial Public  
130 Health Response and Interim Clinical Guidance for the 2019 Novel Coronavirus Outbreak —  
131 United States, December 31, 2019–February 4, 2020. *Morbidity and Mortality Weekly*  
132 *Report*. 2020;69(5):140-6.
- 133 3. Wang C, Horby PW, Hayden FG, Gao GF. A novel coronavirus outbreak of global health  
134 concern. *Lancet (London, England)*. 2020;395(10223):470-3.
- 135 4. World Health Organization. Coronavirus disease 2019 (COVID-19) Situation Report 87.  
136 2020 16 April 2020 [cited; Available from: [https://www.who.int/emergencies/diseases/novel-](https://www.who.int/emergencies/diseases/novel-coronavirus-2019/situation-reports)  
137 [coronavirus-2019/situation-reports](https://www.who.int/emergencies/diseases/novel-coronavirus-2019/situation-reports)
- 138 5. Andersen K. Clock and TMRCA based on 27 genomes. 2020 25 January 2020 [cited;  
139 Available from: <http://virological.org/t/clock-and-tmrca-based-on-27-genomes/347>
- 140 6. Andersen KG, Rambaut A, Lipkin WI, Holmes EC, Garry RF. The proximal origin of SARS-  
141 CoV-2. *Nature Medicine*. 2020 2020/03/17.
- 142 7. Bedford T, Greninger AL, Roychoudhury P, Starita LM, Famulare M, Huang M-L, et al.  
143 Cryptic transmission of SARS-CoV-2 in Washington State. 2020:2020.04.02.20051417.
- 144 8. Deng X, Gu W, Federman S, Du Plessis L, Pybus O, Faria N, et al. A Genomic Survey of  
145 SARS-CoV-2 Reveals Multiple Introductions into Northern California without a  
146 Predominant Lineage. 2020:2020.03.27.20044925.



- 147 9. Quick J, Grubaugh ND, Pullan ST, Claro IM, Smith AD, Gangavarapu K, et al. Multiplex  
148 PCR method for MinION and Illumina sequencing of Zika and other virus genomes directly  
149 from clinical samples. *Nature Protocols*. 2017;12(6):1261-76.
- 150 10. Quick J, Loman NJ, Duraffour S, Simpson JT, Severi E, Cowley L, et al. Real-time, portable  
151 genome sequencing for Ebola surveillance. *Nature*. 2016;530(7589):228-32.
- 152 11. Holmes EC, Y Z. Novel 2019 coronavirus genome. 2020 [cited 2020 5 Apr]; Available  
153 from: <http://virological.org/t/novel-2019-coronavirus-genome/319>
- 154 12. Harcourt J, Tamin A, Lu X, Kamili S, Kumar Sakthivel S, Wang L, et al. Isolation and  
155 characterization of SARS-CoV-2 from the first US COVID-19 patient.  
156 2020:2020.03.02.972935.
- 157 13. Kujawski SA, Wong KK, Collins JP, Epstein L, Killerby ME, Midgley CM, et al. First 12  
158 patients with coronavirus disease 2019 (COVID-19) in the United States.  
159 2020:2020.03.09.20032896.

160

161 Address for correspondence: Suxiang Tong, Division of Viral Diseases, Centers for Diseases  
162 Control and Prevention, 1600 Clifton Rd NE, Mailstop H18-6, Atlanta, GA 30329; email:  
163 sot1@cdc.gov

164

165 Table 1. Genome consensus accuracy

<b>Virus tier (C<sub>T</sub>)</b>	<b>%Coverage (20X)<sup>a</sup></b>	<b>Indels</b>	<b>Indel bases</b>	<b>SNPs</b>	<b>%Identity<sup>b</sup></b>
22.3	99.659	0	0	0	100
	99.722	0	0	0	100
	99.635	0	0	0	100
25.7	99.635	0	0	0	100
	99.615	0	0	0	100
	99.642	0	0	0	100
29.2	99.508	0	0	0	100
	99.635	0	0	0	100
	99.615	0	0	0	100
33.2	93.024	1	1	0	100
	93.603	2	35	0	100
	87.894	0	0	0	100
35.6	41.653	1	1	0	100
	51.266	0	0	1	99.993
	50.821	1	15	2	99.987
37.6	14.634	0	0	1	99.977
	9.317	0	0	0	100
	12.363	0	0	0	100

166 <sup>a</sup> The 5' and 3' ends are primer sequence, so 100% coverage is not possible

167 <sup>b</sup> Percent of covered bases identical to reference sequence, excludes indels and low-coverage bases

168

169 Table 2. Comparison of input, time, and cost requirements for sequencing one or 96 specimens

Method	Input <sup>a</sup>	Single sample		96 Samples	
		Turnaround time	Approx. cost per sample <sup>c</sup>	Turnaround time	Approx. cost per sample <sup>c</sup>
<b>Multiplex/MinION</b>	10 uL	6-8 hours	\$528.70	8-10 hours	\$35.88
<b>Multiplex/MiSeq</b>	10 uL	30-68 hours <sup>b</sup>	\$1443.29	30-68 hours <sup>b</sup>	\$57.87
<b>Singleplex/Sanger</b>	190uL	16-18 hours	\$354.40	17-19 days	\$354.40

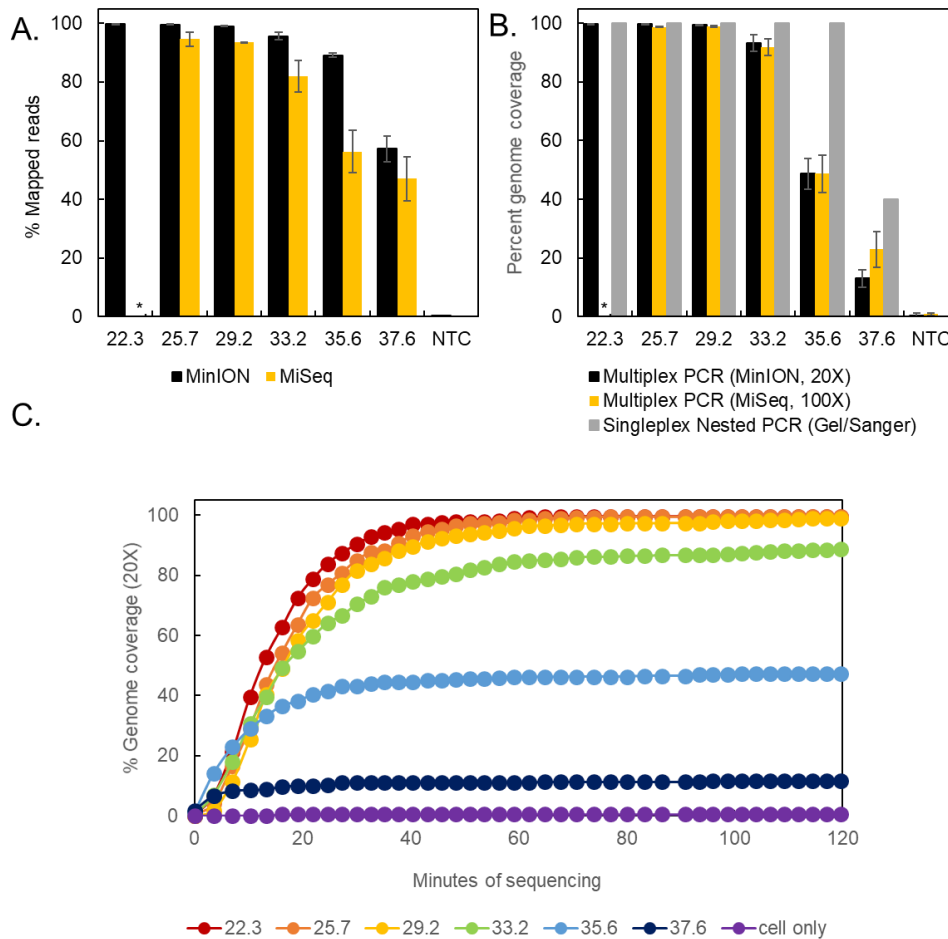
170 <sup>a</sup> Assumes a process of: 200uL resuspended respiratory specimen, extracted and eluted into 100uL

171 <sup>b</sup> Varies based on sequencing kit used

172 <sup>c</sup> Includes specific enzyme and reagent costs, excludes common laboratory supplies and labor costs

173

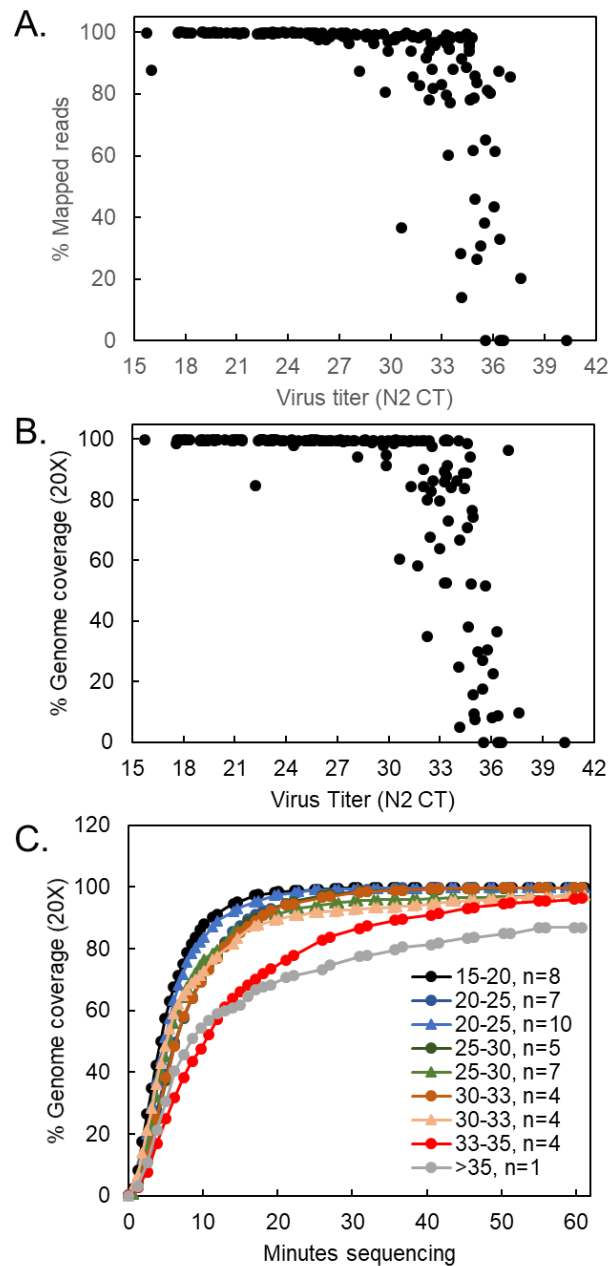
174 **Figure 1. Limits of detection.**



175

176 Triplicate serial dilutions of SARS-CoV-2 isolate A12 (12) amplified using the singleplex or  
177 multiplex primer set. The multiplex amplicons were barcoded, library-prepped, and sequenced  
178 on a MinION or MiSeq. (A) Percent of reads that map to the virus genome for each sample. (B)  
179 Percent of virus genome that is covered at >20X depth by the multiplex amplicons on the  
180 MinION (black) or >100X depth on the MiSeq (orange), or covered by the nested, singleplex  
181 amplicons (grey) (measured by presence or absence on a gel). (C) Real-time analysis of MinION  
182 sequencing data. Each data point represents the average 20X genome coverage of three  
183 replicates.

184 **Figure 2. Sequencing SARS-CoV-2 clinical samples.**



185

186 (A) Percent mapped and (B) percent genome coverage for 167 clinical SARS-CoV-2 samples,  
187 amplified with multiplex PCR strategy and sequenced on the MinION. (C) Time-lapse of 20X  
188 genome coverage obtained by MinION for clinical specimens at the indicated  $C_T$  values. Data  
189 points represent the average coverage for the indicated number of samples