# Targeted Amplification and Genetic Sequencing of the Severe Acute Respiratory Syndrome Coronavirus 2 Surface Glycoprotein

**Running Title**: Targeted SARS-CoV-2 S-Gene Sequencing

**Authors**:

Matthew W. Keller[1][#], Lisa M. Keong[1], Benjamin L. Rambo-Martin[1], Norman Hassell[1], Kristine Lacek[1], Malania M. Wilson[1], Marie K. Kirby[1], Jimma Liddell[1], D. Collins Owuor[1], Mili Sheth[2], Joseph Madden[2], Justin S. Lee[2], Rebecca J. Kondor[1], David E. Wentworth[1], and John R. Barnes[1][#]

**Affiliations**:

1) Influenza Division, National Center for Immunization and Respiratory Diseases (NCIRD), Centers for Disease Control and Prevention (CDC), Atlanta, Georgia, USA

2) Biotechnology Core Facility Branch, Division of Scientific Resources, National Center for Emerging and Zoonotic Infectious Diseases, Centers for Disease Control and Prevention, Atlanta, Georgia, USA

# Corresponding authors:

Matthew W. Keller, PhD: nqp3@cdc.gov

John R. Barnes, PhD: fzq9@cdc.gov

## Abstract

The SARS-CoV-2 spike protein is a highly immunogenic and mutable protein that is the target of vaccine prevention and antibody therapeutics. This makes the encoding S-gene an important sequencing target. The SARS-CoV-2 sequencing community overwhelmingly adopted tiling amplicon-based strategies for sequencing the entire genome. As the virus evolved, primer mismatches inevitably led to amplicon drop-out. Given the exposure of the spike protein to host antibodies, mutation occurred here most rapidly, leading to amplicon failure over the most insightful region of the genome. To mitigate this, we developed SpikeSeq, a targeted method to amplify and sequence the S-gene. We evaluated 20 distinct primer designs through iterative *in silico* and *in vitro* testing to select the optimal primer pairs and run conditions. Once selected, periodic *in silico* analysis monitor primer conservation as SARS-CoV-2 evolves. Despite being designed during the Beta wave, the selected primers remain > 99% conserved through Omicron as of 2023-04-14. To validate the final design, we compared SpikeSeq data and National SARS-CoV-2 Strain Surveillance whole-genome data for 321 matching samples. Consensus sequences for the two methods were highly identical (99.998%) across the S-gene. SpikeSeq can serve as a complement to whole-genome surveillance or be leveraged where only S-gene sequencing is of interest. While SpikeSeq is adaptable to other sequencing platforms, the Nanopore platform validated here is compatible with low to moderate throughputs, and its simplicity better enables users to achieve accurate results, even in low resource settings.

## Introduction

In December 2019 an outbreak of pneumonia of unknown cause began in Wuhan, China (1). This illness (COVID-19) was found to be caused by a novel betacoronavirus, severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2, SC2) (2). The virus quickly spread around the world, and in March 2020, the World Health Organization officially declared COVID-19 a pandemic (3). As of May 2023, COVID-19 has caused roughly 677 million infections and 6.9 million deaths (4).

As the COVID-19 pandemic progressed, waves of new variants spread (5), and mutations within the surface glycoprotein (spike) accumulated (6, 7). The spike protein is key to the viral replication cycle as its binding to the human angiotensin-converting enzyme 2 (ACE2) receptor initiates cellular entry of the virus (8). It also bears clinical significance as it is the target of vaccine prevention (9) and antibody therapeutics (10). The continual evolution of SARS-CoV-2 to evade immune pressures has led to a plethora of spike mutations that have been deleterious to vaccine effectiveness (11) and antibody neutralization (12-14). Importantly, many of these mutations are located within the receptor-binding domain (RBD) (15) where 90% of neutralizing antibodies target SARS-CoV-2 (16). As such, the spike protein encoding S-gene is an important sequencing target, and complete and accurate data for the S-gene is paramount for high quality surveillance information.

54 Genomic tools, such as the widely used Artic SARS-CoV-2 primer set, have required numerous

55 updates to remain effective against new variants (17-20) (https://github.com/artic-network/artic-

56 ncov2019/tree/master/primer_schemes/nCoV-2019). It can be challenging for surveillance labs, which

57 are likely operating at surge capacity, to examine available alternative methods and validate revisions

58 to the method for use. When overlooked, these limitations can lead to overt sequencing gaps or areas

59 of low coverage, usually within the S-gene (21). This issue has occurred multiple times during the

60 COVID-19 pandemic with major variant transitions (Origin strain to Alpha, Alpha to Delta, and Delta to

61 Omicron). But this problem is not limited to major variant shifts. Variability in sequencing protocols and

62 the design of sequencing primers in highly mutable regions of the SARS-CoV-2 spike protein causes

63 intermittent sequencing dropouts, even with moderate amounts of variation. This is complicated further

64 by the variability of organization for countrywide sequencing. Some countries have centralized health

65 systems with more direct control of sequencing protocols and communication. Whereas other countries

66 contract out sequencing to private labs, which can yield higher data volume but have more variability in

67 sequencing methods and directness of communication. This can lead to protocol issues that are very

68 slow to address, causing blind spots to critical regions of the spike protein as evolution occurs. As an

69 illustration, we examined global surveillance data across the SARS-CoV-2 RBD over different time

70 periods (**Figure S01**). During the transition from Delta to Omicron, this critical region was missing a

71 significant amount of data. Moreover, the shape of that missing data resembles the amplicon 76

72 dropout known to affect the Artic SARS-CoV-2 primer set during the emergence of Omicron (19). The

73 coverage across this region has since improved, but this does illustrate the issues of using mutation-

74 sensitive amplification methods through a highly mutable region of a highly mutable virus.

75 Efforts have been made to focus surveillance to the S-gene; however, these methods have serious

76 limitations. One such effort is to use eight overlapping amplicons (342-979 bp) and sanger sequencing

77 to bring SARS-CoV-2 surveillance to low resource areas (22). Unfortunately, the need for eight RT-

78 PCRs per sample and the use of sanger sequencing is costly, labor intensive, and seriously limits

79 throughput potential. In an effort to improve the throughput of S-gene only sequencing, a modified

80 version of Artic V3 SARS-CoV-2 primer set, HiSpike, was developed (23). HiSpike retains many of the

81 limitations of the Artic SARS-CoV-2 protocol, most notably, the use of small amplicons (~400 bp) and

82 the need for many primers to bind within the spike coding region. Using many primers to generate

83 many small overlapping amplicons is not ideally suited to the surveillance of a rapidly evolving RNA

84 virus and will likely again lead to sequencing dropouts due to primer mismatches. Indeed, multiple

85 primers from both studies have conservation issues.

86 Because of these challenges, it was critical to develop a robust method for obtaining rapid sequence

87 information, specifically for the S-gene. For this purpose, we developed SpikeSeq, a targeted method to

3

88  amplify and sequence the S-gene (**Figure 1**). SpikeSeq uses four carefully selected and highly

89  conserved primers (**Table 1**) to produce two overlapping amplicons that yield full coverage of the S-

90  gene (**Figure 2**).

## Results

### Primer Selection and Validation

93  We used the conservation of all available SARS-CoV-2 sequences to identify (**Table S01**) and evaluate

94  (**Table S02**) candidate primers. We identified three candidates for each of the 4 needed primers (S1F,

95  S1R, S2F, and S2R) with additional candidates for S2R where SARS-CoV-2 (Wuhan-hu-1,

96  NC_045512.2) and SARS-CoV-1 (NC_004718.3) shared identity. We eliminated those with < 95%

97  conservation for all available SARS-CoV-2 sequences. By testing candidate primer combinations

98  across an annealing temperature gradient (**Table S03; Figures S02-S03**), we were able to

99  simultaneously eliminate possible combinations with poor performance and select 60°C as the

100 annealing temperature. Finally, a limit of detection assay (**Tables S03-S04; Figure S04**) was used to

101 select S1F_21358, S1R_23813, S2F_23288, and S2R_25460 as the final primers (**Table 1; Seq S01**).

102 We periodically monitor the conservation of these primers, and as of 2023-07-28, the selected primers

103 remain highly conserved against SARS-CoV-2 using three months of US data, three months of global

104 data, and all global data (**Figures S05-S16**). SpikeSeq primers also show some conservation against

105 related coronaviruses (**Seq S02; Figure S17**). If a particular subvariant is of concern, we can perform a

106 more focused conservation analysis. Such was the case with Omicron XBB, XBB.1.5, and derivatives.

107 Analysis against those subvariants, as of 2023-01-18, demonstrated that our primers remained

108 conserved (**Table S05**).

109 Our design results in an amplification strategy where two overlapping amplicons, each in their own RT-

110 PCR reaction, span the entire gene. The four selected primers, which will generally be known as S1F,

111 S1R, S2F, and S2R, avoid mutations and regions of high diversity (**Figure 2**).

### SpikeSeq Runs

113 We performed 14 Nanopore sequencing runs to validate and characterize SpikeSeq. A summary of

114 these runs is available in the supplemental materials (**Table S06**).

### Sensitivity and Specificity

116 The limit of detection (LOD) via MinION flow cell sequencing was ~ 100 copies/µL and Ct 30 (**Table**

117 **S07; Figure S18**). Via Flongle flow cell sequencing, the LOD was ~ 100 copies/µL and Ct 27 (**Table**

118 **S08; Figure S19**).

119 No reads from the 84 NTCs mapped to SARS-CoV-2 (**Tables S09-S10; Figure S20-S21**).

4

**SpikeSeq Validation**

120

121   We tested 377 samples via SpikeSeq for a pairwise comparison to National SARS-CoV-2 Strain

122   Surveillance (NS3) whole-genome data. Of those, 321 samples passed SpikeSeq (**Seq S04**) and

123   whole-genome sequencing (**Seq S05**) to be carried forward for further analysis. The S-gene consensus

124   sequences were highly identical with 1,225,156 identities out of 1,225,185 positions (99.998%

125   identical). Analyzing SpikeSeq data via Nextclade or Pangolin is limited by some group defining

126   mutation residing outside of the S-gene. Still, with the widespread use of these analytical tools, we

127   wanted to characterize the Nextclade results of SpikeSeq derived S-gene sequences in comparison to

128   whole-genome sequences. Of the 281 samples that had a variant assignment (e.g., Delta or Omicron),

129   Nextclade assignment of these variants was 100% concordant between SpikeSeq and whole-genome

130   data. These assignments included the variants Alpha, Beta, Gamma, Delta, Epsilon, Eta, Iota, Lambda,

131   Mu, and Omicron (24, 25). Identical clades were assigned for 91% of samples. As expected, clades

132   with identical S-gene sequences, such as clades 21A (Delta) and 21J (Delta), were often conflated.

133   Clades 21K (Omicron) and 21L (Omicron) were accurately assigned due to the S-gene diversity

134   between those clades. Identical Nextclade_pango lineages were assigned for 72% of samples. Similar

135   to clade identification, the resolution of SpikeSeq lineage identification is limited by a great number of

136   named lineages and their identification being based off mutations outside the S-gene. (**Table S11**)

137   Importantly, SpikeSeq identified 4,428 mutations which includes all 4,422 spike protein mutations

138   identified by whole-genome sequencing. For six samples, SpikeSeq identified one additional mutation

139   each (**Table S11**). Further investigation of raw read data confirmed that these additional mutations

140   were due to minor subpopulations at > 20% frequency amplified at variable proportions due to separate

141   rounds of PCR between SpikeSeq and NS3. In any case, correctly identifying all 4,422 presumably true

142   spike protein mutations reflects a high degree of accuracy that is more than sufficient for surveillance

143   purposes.

144   A subset of 277 clinical specimens from the NS3 project were used for additional characterization of

145   SpikeSeq. By comparing Ct values to SpikeSeq coverage results (**Table S12**), we found that 98-99% of

146   samples with a Ct value less than 25 (n = 217) passed the coverage threshold of requiring ≥ 50x

147   coverage at every position. For samples with Ct values between 25 and 30 (n = 44), 89% of samples

148   passed. And for samples with Ct values over 30 (n = 16), 81% of samples passed (**Figure S22**).

149   For this subset of 277 clinical specimens, we split the three Nanopore libraries for loading on standard

150   MinION flow cells (FLO-MIN106) for 72 hours and disposable Flongle flow cells (FLO-FLG001) for 24

151   hours. These flow cell types are known to have disparate sequencing yields, and indeed the Flongle

152   flow cells produced just ~1% of the average coverage compared to the MinION flow cells. However, the

153    coverage thresholds for SpikeSeq only requires a full assembly be made and ≥ 50x coverage at every

154    position. Using those requirements, MinION flow cell sequencing passed 267/277 samples (96%), and

155    Flongle flow cell sequencing passed 241/277 samples (87%). In other words, ~ 1% of the average

156    coverage from Flongle flow cells passed 90% (241/267) as many samples with respect to MinION flow

157    cells (**Table S13; Figure S23**).

158    For this same subset of 277 clinical specimens, we diverted a portion of the spike amplicons to Illumina

159    sequencing, and 251 samples passed both techniques. We compared consensus level identity between

160    spike amplicons sequenced via Nanopore (MIN) to those same amplicons sequenced via Illumina (ILL;

161    **Seq S06**). For 251 samples, consensus sequences were highly identical with 958,236 identities out of

162    958,239 positions (99.9997% identical). This was expanded to a three-way comparison that includes

163    the corresponding NS3 generated S-gene sequences (NS3; **Table S14**). The MINvNS3 consensus

164    sequences were 99.9972% (958,212/958,239) identical, and the ILLvNS3 consensus sequences were

165    99.9969% (958,210/958,240) identical. All 251 samples had 100% identity between at least two of the

166    three methods, and 20 samples had discrepant results. Because at least two of the methods always

167    agreed, the discrepant results always appeared in pairs, were of identical magnitude, and shared a

168    common method. For example, the three-way blast results of sample 3002648260 for MINvILL is 100%

169    identical whereas MINvNS3 and ILLvNS3 are both 99.974% identical. This indicates the discrepancy

170    lies with the NS3 derived S-gene consensus sample 3002648260. Of the 20 samples with discrepant

171    results, 18 are due to discrepancies with the NS3 derived S-gene consensus, and 2 are due to

172    discrepancies with the Illumina sequenced spike amplicons. This distribution of discrepancies is

173    expected as the NS3 samples were independently amplified, processed, and analyzed. Ultimately

174    though, these discrepancies are very minor and more than acceptable for surveillance purposes.

175    **Phylogenetics**

176    We visualized the nextclade results in auspice to generate a tanglegram (**Figure S24**) of matching

177    samples (n = 321) that passed both SpikeSeq and whole-genome sequencing. As detailed in **Table**

178    **S11**, variant assignment was 100% concordant and clade assignment was highly concordant with

179    ambiguities appearing for different clades with identical S-gene sequences.

180    **Discussion**

181    We have developed and validated a robust method for amplifying and sequencing the SARS-CoV-2

182    surface glycoprotein. The length of the S-gene necessitated internal primers and separate overlapping

183    RT-PCRs. Still, we were able to limit the total number of primers to four and the number of primers

184    within the coding region of the S-gene to two. With so few primers required, we were able to evaluate

185    many candidates for conservation and efficacy. We also ensured the primer binding sites avoided any

186  major structural/functional elements. This design and process of primer selection gives SpikeSeq its

187  best odds at avoiding mutations that might affect primer annealing. Indeed, despite being originally

188  designed during the beta wave, the selected primers have remained highly conserved through Omicron

189  as of 2023-04-14.

190  We validated this method against hundreds of clinical specimens collected for genomic surveillance by

191  the NS3 project. We compared SpikeSeq data to the available whole-genome data to confirm that

192  SpikeSeq accurately represents the S-gene amino acid mutations. Having only one amplicon in each

193  reaction and two amplicons total, QC via electrophoresis can immediately reveal dropouts. This,

194  combined with strict coverage requirements, ensures only complete and high quality S-gene

195  sequencing data is reported.

196  As the pandemic progressed, the methods used by NS3 required several revisions, and occasionally,

197  relied upon SpikeSeq for sequence completion (**Figure S25**) or confirmation of recombination (26). The

198  speed of SpikeSeq has also proven useful during the floods of high priority samples associated with the

199  appearance of new variants. In one case, SpikeSeq was used to confirm the presence of Omicron in

200  the US Virgin Islands which allowed the territory to acquire antibody therapeutics best suited at the time

201  for infection with the Omicron variant.

202  SpikeSeq can serve as a complement to whole-genome sequencing data with S-gene coverage gaps,

203  be leveraged as a tool for projects in which only S-gene sequencing is of interest, and stand alone as a

204  means of surveillance. SpikeSeq was evaluated and approved as a research use only method by the

205  CDC Infectious Disease Test Review Board, and is currently being deployed to partner laboratories. In

206  collaboration with The World Health Organization, we are hosting intensive week-long regional trainings

207  where country representatives will gain hands-on experience and receive reagents, consumables, and

208  a Mk1C sequencing device sufficient to perform a year of SARS-CoV-2 S-gene (and influenza A virus)

209  surveillance. These regional trainings dedicate a great deal of time on foundational knowledge about

210  the viruses themselves, the fundamentals of a quality surveillance effort, the importance of each step of

211  data analysis and curation, and critically evaluating all step of the surveillance pipeline to ensure only

212  quality data is submitted to public databases. While it is tempting to simply ship out point-and-click

213  solutions, we want to develop a strong foundational knowledge about surveillance and data curation

214  when training and equipping labs and countries new to next-generation sequencing (NGS) surveillance.

215  This not only ensures the best use of resources, but gives those labs and countries the best chance at

216  being successful in generating quality data and participating in this global surveillance effort. Moreover,

217  because SpikeSeq targets a portion of the genome, a given amount of surveillance capacity could

218  cover several times more samples compared to whole-genome sequencing on the same platform. The

219  Nanopore platform used by SpikeSeq is compatible with low to moderate throughputs, and its simplicity

220 better enables users to achieve accurate results, even in low resource settings. Finally, the relatively

221 low capital expenditure makes this strategy an ideal starting point for public health laboratories new to

222 NGS surveillance. As of April 2023, public health representatives from 59 countries have received this

223 training with 23 more scheduled by the end of August 2023.

224 Whole-genome sequencing by a variety of methods will remain an integral part of SARS-CoV-2

225 surveillance, and we are not intending SpikeSeq to simply be a replacement. Whole-genome

226 sequencing is the only way to properly assign phylogenetic relationships or monitor for amino acid

227 mutations outside of the S-gene that can, for example, affect viral replication and pathogenesis (27).

228 Moreover, quality whole-genome data is necessary to monitor primer conservation for any targeted

229 amplification strategy.

230 SpikeSeq represents a refocusing on essential information needed from surveillance data. Whole-

231 genome surveillance of SARS-CoV-2 has occasionally, and unfortunately, prioritized getting any result

232 at the expense of sequence completeness and quality. As an example, eagerness to define new

233 clades/lineages based on trivial differences has convoluted the classification of SARS-CoV-2 viruses

234 and obscured the relationships between similar or disparate S-gene mutations that carry clinical

235 significances. By focusing on the S-gene, imposing strict coverage and quality metrics, and applying

236 lessons learned through surveillance of the diverse RNA influenza viruses, we hope to supplement

237 SARS-CoV-2 surveillance with complete and quality reporting on the rapidly mutating S-gene.

238 **Materials and Methods**

239 **Molecular Workflow**

240 To amplify the S-gene, we produced overlapping amplicons (S1 and S2) via separate SuperScript™ IV

241 One-Step RT-PCR System (Thermo Fischer Scientific, USA) reactions. The RT-PCR mixture

242 contained: 4.25 µL nuclease-free water, 12.5 µL SSIV 2X reaction mix, 0.25 µL SSIV RT Mix, 5 µL S1

243 or S2 primer pairs, and 3 µL of RNA. The RT-PCR conditions are as follows: 10 minutes at 50°C; 2

244 minutes at 98°C; 40 cycles of 10 seconds at 98°C, 10 seconds at 60°C, and 1 minute 15 seconds at

245 72°C; a final elongation of 5 minutes at 72°C; and a hold at 4°C. Electrophoresis quality control was

246 performed on individual RT-PCRs. After QC, corresponding S1 and S2 amplicons were combined,

247 cleaned via SPRI beads (1x) with ethanol washes, and eluted into 15 µL of nuclease-free water.

248 Nanopore libraries were prepared using SQK-LSK109 and EXP-NBD196 and sequenced on GridION

249 (Oxford Nanopore Technologies, UK) using FLO-MIN106 or FLO-FLG001 flow cells.

250 Laboratory procedures for RT-PCR and library preparation are available in the supplemental (**Text S01**

251 and **Text S02**).

252   For Illumina sequencing, a portion of the cleaned amplicons were taken and prepared using the

253   Nextera XT sample preparation kit. Since the SARS-CoV-2 S-gene amplicons are of a similar size to

254   the influenza virus amplicons, they were processed via the standard influenza surveillance pipeline

255   used by the CDC Genomics and Diagnostics Team (28, 29).

256   **Sequencing Data Analysis**

257   During the sequencing run, we used the GridION MinKNOW to perform super-accuracy basecalling live

258   (ont-guppy-for-gridion 5.0.17 or 5.1.13), to trim the barcodes, and to filter the reads. We trimmed

259   primers using BBDuk (30), restricted the trimming using restrictleft=50 and restrictright=50, and referred

260   to the primer sequences (**Seq S01**). We assembled reads using IRMA

261   (https://wonder.cdc.gov/amd/flu/irma/irma.html) with the CoV-s-gene module (IRMA v.1.0.3

262   https://wonder.cdc.gov/amd/flu/irma/release_notes.html) and mapped to the S-gene reference (28). For

263   a sample to pass SpikeSeq, it must meet coverage and quality metrics. Specifically, it must have a

264   complete S-gene assembly, have at least 50x coverage at every position, and be free of frameshift

265   mutations. Mutations were identified using Nextclade Web version 2.6.1 (https://clades.nextstrain.org;

266   accessed September 30, 2022) SARS-CoV-2 without recombinants (24).

267   Analysis tools are available online https://cdcgov.github.io/MIRA (31).

268   **Primer Selection and Validation**

269   We selected four primer target regions where S1F and S2R would lie outside of the S-gene coding

270   region, and S1R and S2F would be on opposite sides of the S1/S2 cleavage site and avoid major

271   structural elements. We identified multiple sets of candidate primers for each S1F, S1R, S2F, and S2R.

272   For S2R, we also evaluated an area where SARS-CoV-2 (Wuhan-hu-1, NC_045512.2) and SARS-

273   CoV-1 (NC_004718.3) shared identity (**Table S01**). During the Beta wave (March 2021), we evaluate

274   the conservation of primer candidates against 476,466 SARS-CoV-2 genomes (**Table S02**). Twenty

275   primer combinations were tested (**Table S03**). We initially screened the candidate primer pairs across a

276   temperature gradient using RNA from B.1.351 (Beta) with a ct 25 as determined by the Flu SC2

277   multiplex assay (32). We used an LOD of B.1.351 (Beta) from ct 14-30 (846k to 16 copies/μL; **Table

278   S04**) to finalize the primer selection. The presence of amplicons was determined using a QIAxcel HT

279   fragment analyzer.

280   We monitored the conservation of the primers by downloading data from GISAID. Downloaded genomic

281   data was aligned to the Wuhan-Hu-1 reference (NCBI accession MN908947.3) genome using SSW

282   (33). Aligned genome primer regions were regularly compared for mismatches against each individual

283   primer sequence. This information was used to highlight potential assay issues with new emerging

284  variants. We downloaded diversity (entropy) data from Nextstrain

285  (https://nextstrain.org/ncov/gisaid/global/6m; accessed March 6, 2023) (24).

**Sensitivity and Specificity**

287  To measure the absolute limit of detection, we used a custom synthetic RNA fragment from Twist

288  Bioscience (CA, USA) based on the Delta lineage virus hCoV-19/USA/CO-CDC-MMB09467199/2021.

289  The sequence for this fragment (TwistDeltaFragment_4276451.fasta) is available in the supplemental

290  materials (**Seq S03**). The 4,626-nucleotide fragment spans the S-gene and extends into neighboring

291  genes. The synthetic fragment aliquots were delivered at 629,000 copies/µL as determined via

292  manufacturer ddPCR. To measure the viral limit of detection, we used a propagated isolate of Delta

293  SARS-CoV-2 and measured the Ct value of the serial dilutions using the Flu SC2 multiplex assay (32).

294  The limit of detection was determined by the most dilute sample to pass coverage and quality

295  thresholds for all the replicates.

296  We prepared RT-PCR master mixes for triplicate limit of detection assays using both synthetic and viral

297  material. The dilution series were a 5-fold serial dilution through 7 steps with a water NTC as the 8th

298  step. LOD amplicons were split at the end-prep stages for sequencing on both MinION and flongle flow

299  cells. For sequencing on MinION flow cells, we included 48 additional water NTCs. For sequencing on

300  Flongle flow cells, we included 24 additional A549 RNA (Rp Ct 22) NTCs.

**SpikeSeq Validation**

302  To validate this method, we tested a total of 377 specimens from the NS3 project. We started with a

303  retrospective analysis of 277 clinical specimens that were collected from March to August 2021 and

304  that capture the diversity of SARS-CoV-2 into the Delta wave. During the omicron wave, we continued

305  SpikeSeq validation concurrently with NS3. These additional 100 samples were collected from

306  November 2021 to January 2022. Of these 377 samples, 321 passed SpikeSeq (**Seq S04**) and whole-

307  genome sequencing (**Seq S05**) to be carried forward for further analysis.

308  We compared matching samples (n = 321) that passed both SpikeSeq and whole-genome sequencing

309  using ncbi-blast+/2.9.0 (34) and Nextclade Web version 2.6.1 (https://clades.nextstrain.org; accessed

310  September 30, 2022) SARS-CoV-2 without recombinants (24). Using the output of Nextclade, we

311  evaluated the concordance of variant, clade, and lineage assignment. We also compared the reported

312  S-gene amino acid mutations for complete matches of corresponding samples and by counting

313  individual mutations for corresponding samples (**Table S11**).

314  A subset of 277 samples through Delta  was used to compare Ct values to coverage, Nanopore

315  sequencing yield on two flow cell types (FLO-MIN106 versus FLO-FLG001), and Nanopore sequencing

316  accuracy to Illumina sequencing. Each time the RNA was thawed, we tested it with the Flu SC2

10

317   multiplex assay (32) to determine the Ct value and amplified the S-gene using the methods presented

318   here. For samples with an undefined Ct value (n =2), a Ct value of 40 was assigned. We then split the

319   spike amplicons to both Illumina and Nanopore sequencing methods. For Nanopore sequencing, we

320   prepared libraries using the methods described here and loaded both standard MinION flow cells (FLO-

321   MIN106) for 72 hours and disposable Flongle flow cells (FLO-FLG001) for 24 hours.

322   All 277 samples from this subset (pass or fail) were used to assess the relative pass rates of standard

323   MinION flow cells (FLO-MIN106) versus Ct value (**Table S12 and Figure S22**) and versus disposable

324   Flongle flow cells (FLO-FLG001; **Table S13 and Figure S23**).

325   From that subset of 277 samples, 251 samples passed both Nanopore (FLO-MIN106) and Illumina

326   sequencing of the SpikeSeq (**Seq S06**). For each of these 251 samples, we used ncbi-blast+/2.9.0 (34)

327   to generate a three-way comparison between: SpikeSeq amplification and Nanopore sequencing (MIN),

328   Illumina sequencing of those same amplicons (ILL), and NS3 surveillance results for the S-gene (NS3;

329   **Table S14**).

330   **Phylogenetics**

331   We compared matching samples (n = 321) that passed both SpikeSeq and whole-genome sequencing

332   using Nextclade Web version 2.6.1 (https://clades.nextstrain.org; accessed September 30, 2022)

333   SARS-CoV-2 without recombinants (24). From this analysis, we exported the phylogenetics and

334   visualized them with Auspice (https://auspice.us; accessed October 6, 2022). We added a metadata

335   sheet to label and highlight added sequences above the backbone sequences.

336   **Primer Kit Manufacturing**

337   CDC Division of Scientific Resources manufactured primers for use in this study and distribution to

338   public health laboratories. The Oligo Synthesis Laboratory synthesized the primers, purified via HPLC,

339   and verified by mass spectrophotometry. Following initial synthesis and purification, we received three

340   QC aliquots for limit of detection analysis and excess material for use in this study. The remaining

341   material (5 mmol each primer) was then transferred to the Diagnostic Manufacturing Laboratory for

342   stochiometric mixing of forward and reverse primers, dispensing, drying, and kit assembly. We received

343   three aliquots for QC testing.

344   **Supplemental Material**

345   Supplemental material for this article may be found at

346   https://figshare.com/articles/dataset/Supplemental_Material/22762076.

347   Supplemental legends are available in the supplemental (**Text S03**).

## Acknowledgements

348

349 We thank the CDC/NCEZID/Division of Scientific Resources/Biotechnology Core Facility Branch/Oligo

350 Synthesis Laboratory for synthesizing the primers used in this study.

351 We thank the CDC/NCEZID/Division of Scientific Resources/Reagent and Diagnostic Services

352 Branch/Diagnostic Manufacturing Laboratory for manufacturing primer kits.

## Competing interests

353

354 We declare no competing interests.

## Data availability

355

356 Corresponding SpikeSeq (Nanopore sequencing) S-gene consensus sequences and NS3 whole-

357 genome consensus sequences are available in the supplemental materials (n = 321 each; **Seq S04-**

358 **S05**). SpikeSeq amplification and Illumina sequencing derived S-gene consensus sequences (n = 251)

359 are available in the supplemental materials (**Seq S06**).

360 https://figshare.com/articles/dataset/Supplemental_Material/22762076

361 FASTQ reads (that BLAT matched to IRMA reference) are available online at NCBI under BioProject:

362 PRJNA999712. The BioSamples (n=810) include the 321 primary validation samples (320 FLO-MIN106

363 and 1 FLO-FLG001), the 238 flongle yield replicates that passed, and 251 Illumina accuracy replicates

364 that passed. https://www.ncbi.nlm.nih.gov/bioproject/PRJNA999712

## References

365

366 1. WHO. 2020. Pneumonia of unknown cause – China. https://www.who.int/emergencies/disease-
367 outbreak-news/item/2020-DON229. Accessed October 4, 2022.
368 2. 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
369 D, Xu W, Wu G, Gao GF, Tan W, China Novel Coronavirus I, Research T. 2020. A novel coronavirus from
370 patients with pneumonia in China, 2019. N Engl J Med 382:727-733.
371 3. Cucinotta D, Vanelli M. 2020. WHO Declares COVID-19 a Pandemic. Acta Biomed 91:157-160.
372 4. Dong E, Du H, Gardner L. 2020. An interactive web-based dashboard to track COVID-19 in real time.
373 Lancet Infect Dis 20:533-534.
374 5. Lambrou AS SP, Steele MK, et al. 2022. Genomic Surveillance for SARS-CoV-2 Variants: Predominance of
375 the Delta (B.1.617.2) and Omicron (B.1.1.529) Variants — United States, June 2021–January 2022.
376 MMWR Morb Mortal Wkly Rep 71:206-211.
377 6. Frampton D, Rampling T, Cross A, Bailey H, Heaney J, Byott M, Scott R, Sconza R, Price J, Margaritis M.
378 2021. Genomic characteristics and clinical effect of the emergent SARS-CoV-2 B. 1.1. 7 lineage in
379 London, UK: a whole-genome sequencing and hospital-based cohort study. Lancet Infect Dis 21:1246-
380 1256.
381 7. Dhar MS, Marwal R, VS R, Ponnusamy K, Jolly B, Bhoyar RC, Sardana V, Naushin S, Rophina M, Mellan TA,
382 Mishra S, Whittaker C, Fatihi S, Datta M, Singh P, Sharma U, Ujjainiya R, Bhatheja N, Divakar MK, Singh
383 MK, Imran M, Senthivel V, Maurya R, Jha N, Mehta P, A V, Sharma P, VR A, Chaudhary U, Soni N, Thukral
384 L, Flaxman S, Bhatt S, Pandey R, Dash D, Faruq M, Lall H, Gogia H, Madan P, Kulkarni S, Chauhan H,
385 Sengupta S, Kabra S, Gupta RK, Singh SK, Agrawal A, Rakshit P, Nandicoori V, Tallapaka KB, Sowpati

386  Divya T, et al. 2021. Genomic characterization and epidemiology of an emerging SARS-CoV-2 variant in
387  Delhi, India. Science 374:995-999.
388  8.  Yan R, Zhang Y, Li Y, Xia L, Guo Y, Zhou Q. 2020. Structural basis for the recognition of SARS-CoV-2 by
389      full-length human ACE2. Science 367:1444-1448.
390  9.  Thompson MG, Burgess JL, Naleway AL, Tyner H, Yoon SK, Meece J, Olsho LEW, Caban-Martinez AJ,
391      Fowlkes AL, Lutrick K, Groom HC, Dunnigan K, Odean MJ, Hegmann K, Stefanski E, Edwards LJ, Schaefer-
392      Solle N, Grant L, Ellingson K, Kuntz JL, Zunie T, Thiese MS, Ivacic L, Wesley MG, Mayo Lamberte J, Sun X,
393      Smith ME, Phillips AL, Groover KD, Yoo YM, Gerald J, Brown RT, Herring MK, Joseph G, Beitel S, Morrill
394      TC, Mak J, Rivers P, Poe BP, Lynch B, Zhou Y, Zhang J, Kelleher A, Li Y, Dickerson M, Hanson E, Guenther
395      K, Tong S, Bateman A, Reisdorf E, et al. 2021. Prevention and Attenuation of Covid-19 with the
396      BNT162b2 and mRNA-1273 Vaccines. New Eng J Med 385:320-329.
397  10. Chen RE, Winkler ES, Case JB, Aziati ID, Bricker TL, Joshi A, Darling TL, Ying B, Errico JM, Shrihari S,
398      VanBlargan LA, Xie X, Gilchuk P, Zost SJ, Droit L, Liu Z, Stumpf S, Wang D, Handley SA, Stine WB, Shi P-Y,
399      Davis-Gardner ME, Suthar MS, Knight MG, Andino R, Chiu CY, Ellebedy AH, Fremont DH, Whelan SPJ,
400      Crowe JE, Purcell L, Corti D, Boon ACM, Diamond MS. 2021. In vivo monoclonal antibody efficacy against
401      SARS-CoV-2 variant strains. Nature 596:103-108.
402  11. Harvey WT, Carabelli AM, Jackson B, Gupta RK, Thomson EC, Harrison EM, Ludden C, Reeve R, Rambaut
403      A, Peacock SJ, Robertson DL, Consortium C-GU. 2021. SARS-CoV-2 variants, spike mutations and immune
404      escape. Nat Rev Microbiol 19:409-424.
405  12. Rees-Spear C, Muir L, Griffith SA, Heaney J, Aldon Y, Snitselaar JL, Thomas P, Graham C, Seow J, Lee N,
406      Rosa A, Roustan C, Houlihan CF, Sanders RW, Gupta RK, Cherepanov P, Stauss HJ, Nastouli E, Doores KJ,
407      van Gils MJ, McCoy LE. 2021. The effect of spike mutations on SARS-CoV-2 neutralization. Cell Rep
408      34:108890.
409  13. VanBlargan LA, Errico JM, Halfmann PJ, Zost SJ, Crowe JE, Purcell LA, Kawaoka Y, Corti D, Fremont DH,
410      Diamond MS. 2022. An infectious SARS-CoV-2 B.1.1.529 Omicron virus escapes neutralization by
411      therapeutic monoclonal antibodies. Nat Med 28:490-495.
412  14. Greaney AJ, Starr TN, Gilchuk P, Zost SJ, Binshtein E, Loes AN, Hilton SK, Huddleston J, Eguia R, Crawford
413      KHD, Dingens AS, Nargi RS, Sutton RE, Suryadevara N, Rothlauf PW, Liu Z, Whelan SPJ, Carnahan RH,
414      Crowe JE, Bloom JD. 2021. Complete Mapping of Mutations to the SARS-CoV-2 Spike Receptor-Binding
415      Domain that Escape Antibody Recognition. Cell Host Microbe 29:44-57.e9.
416  15. Cao Y, Wang J, Jian F, Xiao T, Song W, Yisimayi A, Huang W, Li Q, Wang P, An R, Wang J, Wang Y, Niu X,
417      Yang S, Liang H, Sun H, Li T, Yu Y, Cui Q, Liu S, Yang X, Du S, Zhang Z, Hao X, Shao F, Jin R, Wang X, Xiao J,
418      Wang Y, Xie XS. 2022. Omicron escapes the majority of existing SARS-CoV-2 neutralizing antibodies.
419      Nature 602:657-663.
420  16. Piccoli L, Park Y-J, Tortorici MA, Czudnochowski N, Walls AC, Beltramello M, Silacci-Fregni C, Pinto D,
421      Rosen LE, Bowen JE. 2020. Mapping neutralizing and immunodominant sites on the SARS-CoV-2 spike
422      receptor-binding domain by structure-guided high-resolution serology. Cell 183:1024-1042. e21.
423  17. Itokawa K, Sekizuka T, Hashino M, Tanaka R, Kuroda M. 2020. Disentangling primer interactions
424      improves SARS-CoV-2 genome sequencing by multiplex tiling PCR. PLOS ONE 15:e0239403.
425  18. Davis JJ, Long SW, Christensen PA, Olsen RJ, Olson R, Shukla M, Subedi S, Stevens R, Musser JM, Pride
426      DT. 2021. Analysis of the ARTIC Version 3 and Version 4 SARS-CoV-2 Primers and Their Impact on the
427      Detection of the G142D Amino Acid Substitution in the Spike Protein. Microbiology Spectrum 9:e01803-
428      21.
429  19. Artic Network. 2021.  SARS-CoV-2 V4.1 update for Omicron variant.
430      https://community.artic.network/t/sars-cov-2-v4-1-update-for-omicron-variant/342. Accessed 03 March
431      2022.
432  20. Tyson JR, James P, Stoddart D, Sparks N, Wickenhagen A, Hall G, Choi JH, Lapointe H, Kamelian K, Smith
433      AD, Prystajecky N, Goodfellow I, Wilson SJ, Harrigan R, Snutch TP, Loman NJ, Quick J. 2020.

434       Improvements to the ARTIC multiplex PCR method for SARS-CoV-2 genome sequencing using nanopore.
435       bioRxiv doi:10.1101/2020.09.04.283077:2020.09.04.283077.
436  21.  Sanderson T, Barrett JC. 2021. Variation at Spike position 142 in SARS-CoV-2 Delta genomes is a
437       technical artifact caused by dropout of a sequencing amplicon. Wellcome Open Res 6:305.
438  22.  Salles TS, Cavalcanti AC, da Costa FB, Dias VZ, de Souza LM, de Meneses MDF, da Silva JAS, Amaral CD,
439       Felix JR, Pereira DA, Boatto S, Guimarães MAAM, Ferreira DF, Azevedo RC. 2022. Genomic surveillance
440       of SARS-CoV-2 Spike gene by sanger sequencing. PLOS ONE 17:e0262170.
441  23.  Fass E, Zizelski Valenci G, Rubinstein M, Freidlin PJ, Rosencwaig S, Kutikov I, Werner R, Ben-Tovim N,
442       Bucris E, Erster O, Zuckerman NS, Mor O, Mendelson E, Dveyrin Z, Rorman E, Nissan I. 2022. HiSpike
443       Method for High-Throughput Cost Effective Sequencing of the SARS-CoV-2 Spike Gene. Front Med 8.
444  24.  Aksamentov I, Roemer C, Hodcroft EB, Neher RA. 2021. Nextclade: clade assignment, mutation calling
445       and quality control for viral genomes. Journal of Open Source Software 6:3773.
446  25.  Rambaut A, Holmes EC, O'Toole Á, Hill V, McCrone JT, Ruis C, du Plessis L, Pybus OG. 2020. A dynamic
447       nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. Nat Microbiol 5:1403-
448       1407.
449  26.  Lacek KA, Rambo-Martin B, Batra D, Zheng X-y, Keller MW, Wilson M, Sheth M, Davis M, Burroughs M,
450       Gerhart J, Hassell N, Lee J, Shepard SS, Cook PW, Wentworth DE, Barnes JR, Kondor R, Paden CR,
451       Peacock TPR, Sakaguchi H. 2022. Identification of a Novel SARS-CoV-2 Delta-Omicron Recombinant Virus
452       in the United States. bioRxiv doi:10.1101/2022.03.19.484981:2022.03.19.484981.
453  27.  Johnson BA, Zhou Y, Lokugamage KG, Vu MN, Bopp N, Crocquet-Valdes PA, Kalveram B, Schindewolf C,
454       Liu Y, Scharton D, Plante JA, Xie X, Aguilar P, Weaver SC, Shi PY, Walker DH, Routh AL, Plante KS,
455       Menachery VD. 2022. Nucleocapsid mutations in SARS-CoV-2 augment replication and pathogenesis.
456       PLoS Pathog 18:e1010627.
457  28.  Shepard SS, Meno S, Bahl J, Wilson MM, Barnes J, Neuhaus E. 2016. Viral deep sequencing needs an
458       adaptive approach: IRMA, the iterative refinement meta-assembler. BMC Genomics 17:708.
459  29.  Rambo-Martin BL, Keller MW, Wilson MM, Nolting JM, Anderson TK, Vincent AL, Bagal UR, Jang Y,
460       Neuhaus EB, Davis CT, Bowman AS, Wentworth DE, Barnes JR. 2020. Influenza A virus field surveillance
461       at a swine-human interface. mSphere 5.
462  30.  Bushnell B. 2014. BBMap: a fast, accurate, splice-aware aligner.  Lawrence Berkeley National Lab.(LBNL),
463       Berkeley, CA (United States),
464  31.  Rambo-Martin BL, Lacek KA, Chau R. 2023. MIRA: An Interactive Dashboard for Influenza Genome and
465       SARS-CoV-2 Spike-Gene Assembly and Curation, https://cdcgov.github.io/MIRA/index.html.
466  32.  Shu B, Kirby MK, Davis WG, Warnes C, Liddell J, Liu J, Wu K-H, Hassell N, Benitez AJ, Wilson MM, Keller
467       MW, Rambo-Martin BL, Camara A, Winter J, Kondor RJ, Zhou B, Spies S, Rose LE, Winchell JM, Limbago
468       BM, Wentworth DE, Barnes JR. 2021. Multiplex real-time reverse transcription PCR for influenza A virus,
469       influenza B virus, and severe acute respiratory syndrome coronavirus 2. Emerg Infect Dis 27:1821-1830.
470  33.  Zhao M, Lee W-P, Garrison EP, Marth GT. 2013. SSW Library: An SIMD Smith-Waterman C/C++ Library
471       for Use in Genomic Applications. PLOS ONE 8:e82138.
472  34.  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TLJBB. 2009. BLAST+:
473       architecture and applications.  10:421.

474

475 **Tables and Figures**

476 **Table 1: Primers**

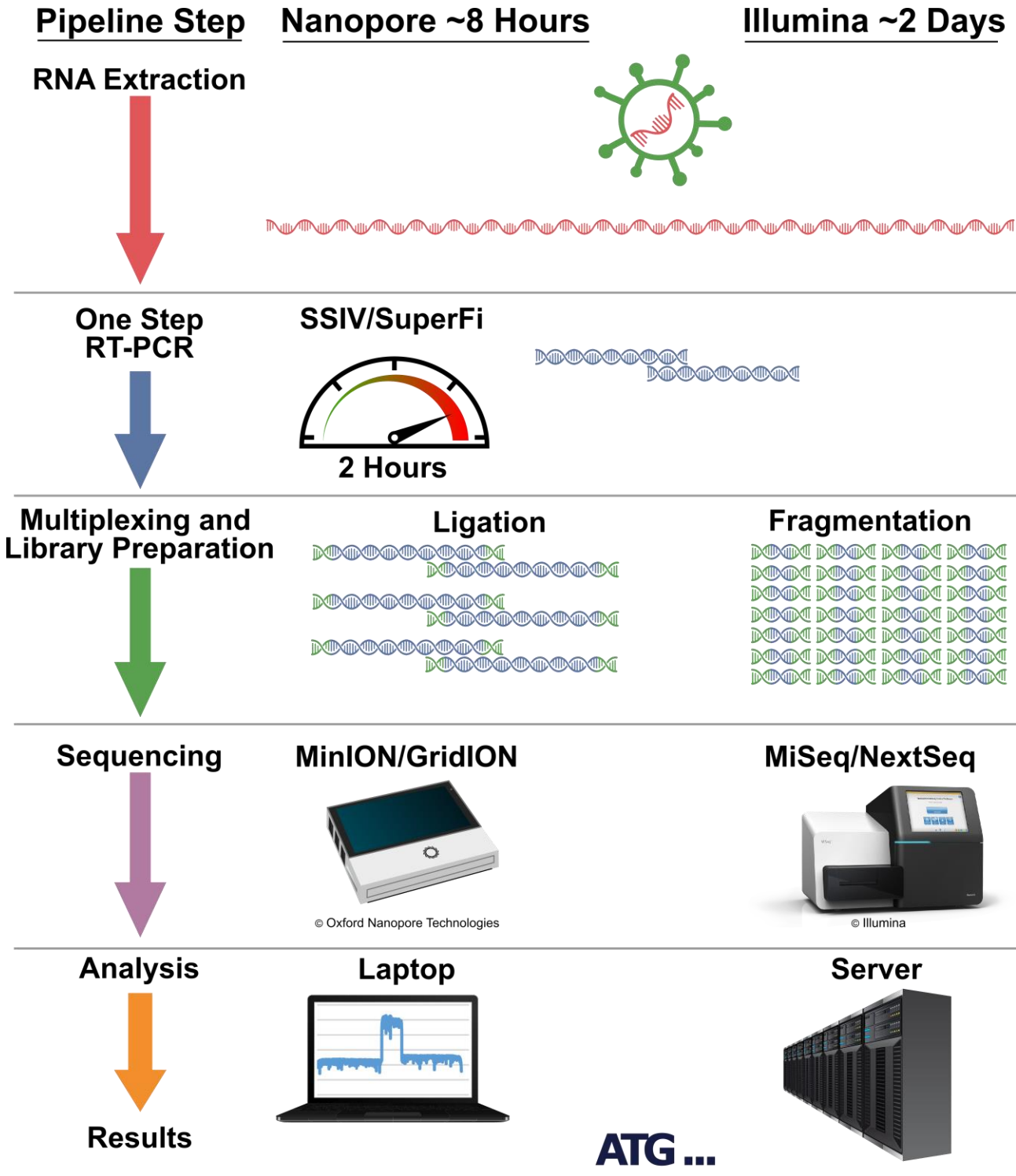477 SpikeSeq primer sequences and working stock concentration.

| S1 primer pool | | |
| --- | --- | --- |
| **Oligo** | **Sequence 5'-3'** | **µM in pool** |
| S1F_21358 | ACAAATCCAATTCAGTTGTCTTCCTATTC | 5 |
| S1R_23813 | TGCTGCATTCAGTTGAATCACC | 5 |
| **S2 primer pool** | | |
| **Oligo** | **Sequence 5'-3'** | **µM in pool** |
| S2F_23288 | GTCCGTGATCCACAGACACTT | 5 |
| S2R_25460 | GCATCCTTGATTTCACCTTGCTTC | 5 |

478

479

480 **Figure 1: SpikeSeq Workflow**



481

482 SpikeSeq is presented here as RT-PCR amplification and Nanopore sequencing. The workflow is

483 designed to be flexible as the amplicons can be diverted to other sequencing platforms.

484   **Figure 2: SpikeSeq Amplification Strategy**



485

486   SpikeSeq amplicons are green with primer locations highlited in red and traced down to the ORFs.

487   Diversity (entropy) across the region is plotted with small circles. Detected amino acid mutations are

488   maked with Xs. SARS-CoV-2 ORFs are black with the Surface Glycoprotein ORF highlighted blue.
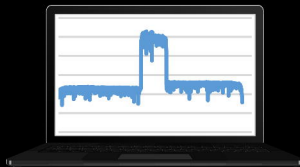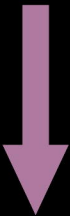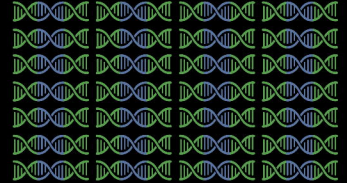
489   Separate one-step RT-PCRs generate overlapping S1 and S2 amplicons that are 2.2 kb and 2.5 kb

490   respectively. These amplicons extend beyond the coding region and overlap across the S1-S2 subunit

491   cleavage site.

# SpikeSeq Amplification Strategy

ATG ...