# Supplementary Materials for

Introductions and early spread of SARS-CoV-2 in France, 24 January to 23 march 2020

**This file includes:**

**Materials and Methods**

Sample collection

After reports of severe pneumonia in late December 2019, enhanced surveillance was implemented in France to detect suspected infections. For each suspected case, respiratory samples from the upper respiratory tract (nasopharyngeal swabs or aspirates) and when possible from the lower respiratory tract, were sent to the NRC, to perform SARS-CoV-2-specific real-time RT-qPCR. Demographic information, date of illness onset, and travel history were obtained when possible. A subset of samples were selected according to the viral load and their sampling location in order to have a broad representation across different regions of France.

Molecular test

RNA extraction was performed with the NucleoSpin Dx Virus Extraction kit (Macherey Nagel). RNA was extracted from 100 µl of specimen, eluted in 100 µl of water and used as a template for RT-qPCR. Samples were tested with a one-step RT-qPCR using three sets of primers as described on the WHO website (https://www.who.int/docs/default-source/coronaviruse/real-time-rt-pcr-assays-for-the-detection-of-sars-cov-2-institut-pasteur-paris.pdf?sfvrsn=3662fcb6_2).

Virus Sequencing

Viral genome sequences were generated by two different approaches. The first consisted in direct metagenomic sequencing, which resulted in complete or near complete genome sequences for samples with viral load higher than $1.45 \times 10^4$ viral genome copies/µl, which corresponds to a Ct value of 25.6 with the IP4 primer set (Fig. 1B, Data S1, Table S3).

2

42  Briefly, extracted RNA was first treated with Turbo DNase (Ambion) followed by

43  purification using SPRI beads Agencourt RNA clean XP (Beckman Coulter). RNA was

44  converted to double stranded cDNA. Libraries were then prepared using the Nextera XT

45  DNA Library Prep Kit (Illumina) and sequenced on an Illumina NextSeq500 (2×150

46  cycles) on the Mutualized Platform for Microbiology (P2M) at Institut Pasteur.

47  For samples with lower viral load, we implemented a highly multiplexed PCR amplicon

48  approach [1] using the ARTIC Network multiplex PCR primers set v1

49  (https://artic.network/ncov-2019), with modification as suggested in [2]. Synthesized

50  cDNA was used as template and amplicons were generated using two pooled primer

51  mixtures for 35 rounds of amplification. We prepared sequencing libraries using the

52  NEBNext Ultra II DNA Library Prep Kit for Illumina (New England Biolabs) and barcoded

53  with NEBNext Multiplex Oligos for Illumina (Dual Index Primers Set 1) (New England

54  Biolabs). We sequenced prepared libraries on an Illumina MiSeq using MiSeq Reagent

55  Kit v3 (2×300 cycles) at the biomics platform of Institut Pasteur.

56  <u>Genome assembly</u>

57  Raw reads were trimmed using Trimmomatic v0.36 [3] to remove Illumina adaptors and

58  low quality reads, as well as primer sequences for samples sequenced with the amplicon-

59  based approach. We assembled reads from all sequencing methods into genomes using

60  Megahit, and also performed direct mapping against reference genome Wuhan/Hu-

61  1/2019 (NCBI Nucleotide – NC_045512, GenBank – MN908947) using the CLC

62  Genomics Suite v5.1.0 (QIAGEN). We then used SAMtools v1.3 to sort the aligned bam

63  files and generate alignment statistics [4]. Aligned reads were manually inspected using

64  Geneious prime v2020.1.2 (2020) (https://www.geneious.com/), and consensus
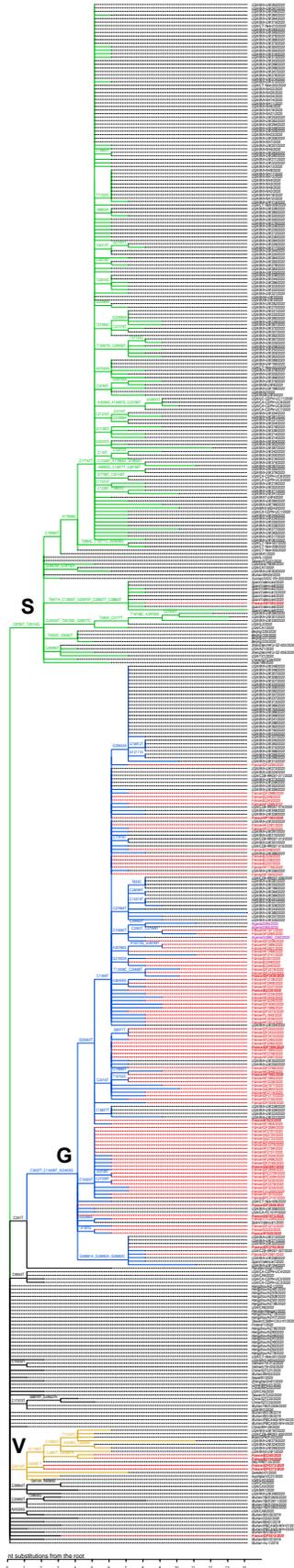
65   sequences were generated using a minimum of 3X read-depth coverage to make a base

66   call. No genomic deletions were detected in the genomes analyzed.

67   <u>Phylogenetic analysis</u>

68   A set of 100 SARS-CoV-2 sequences generated in this study (97 from France, 3 from

69   Algeria) was complemented with 338 genomes published or freely available sequences

70   on GenBank or the GISAID database. From the latter, only published sequences were

71   chosen [5, 6] (Data S1, Table S2). A total of 438 full genome sequences were analyzed

72   with augur and auspice as implemented in the Nextstrain pipeline [7] version from March

73   20, 2020 (https://github.com/nextstrain/ncov). Within the pipeline, sequences were

74   aligned to the reference Wuhan/Hu-1/2020 strain of SARS-CoV-2 (GenBank accession

75   MN908947). The alignment was visually inspected and sequences from France were

76   subset to analyze shared SNPs. No evidence of recombination was detected using RDP

77   v4.97 [8]. A maximum likelihood phylogenetic tree was built based on the GTR model,

78   after masking 130 and 50 nucleotides from the 5' and 3' ends of the alignment,

79   respectively, as well as single nucleotides at positions 18529, 29849, 29851, 29853 to

80   reduce the possibility of including variants due to assembly artifacts as performed in

81   Fauver *et al.*, 2020 [5] and following the Nextstrain implementation for SARS-CoV-2. We

82   checked for temporal signal using Tempest v1.5.3 [9]. The temporal phylogenetic

83   analyses were performed with augur and TreeTime [10], assuming clock rate of

84   0.0008±0.0004 (SD) substitutions/site/year [11], coalescent skyline population growth

85   model and the root set on the branch leading to the Wuhan/Hu-1/2020 sequence. The

86   time and divergence trees were visualized with FigTree v1.4.4

87   (http://tree.bio.ed.ac.uk/software/figtree/). Nucleotide substitutions from the reference

88  sequence that define internal nodes of the tree were extracted from the final Nextstrain

89  build file and annotated on the tree using a custom R script (https://www.R-project.org.

90  Adobe Illustrator 2020 was used to prepare final tree figures. Sequence metadata (Data

91  S1, Table S2), Nextstrain build, and R script is available at https://github.com/Simon-

92  LoriereLab/SARS-CoV-2-France. The phylogeny can be visualized interactively at

93  https://nextstrain.org/community/Simon-LoriereLab/SARS-CoV-2-France. In this study,

94  we used the proposed nomenclature from GISAID to annotate three major clades V, G,

95  and S according to specific single-nucleotide polymorphisms that are shared by all

96  sequences in the clade. Clade defining variants according to GISAID nomenclature are
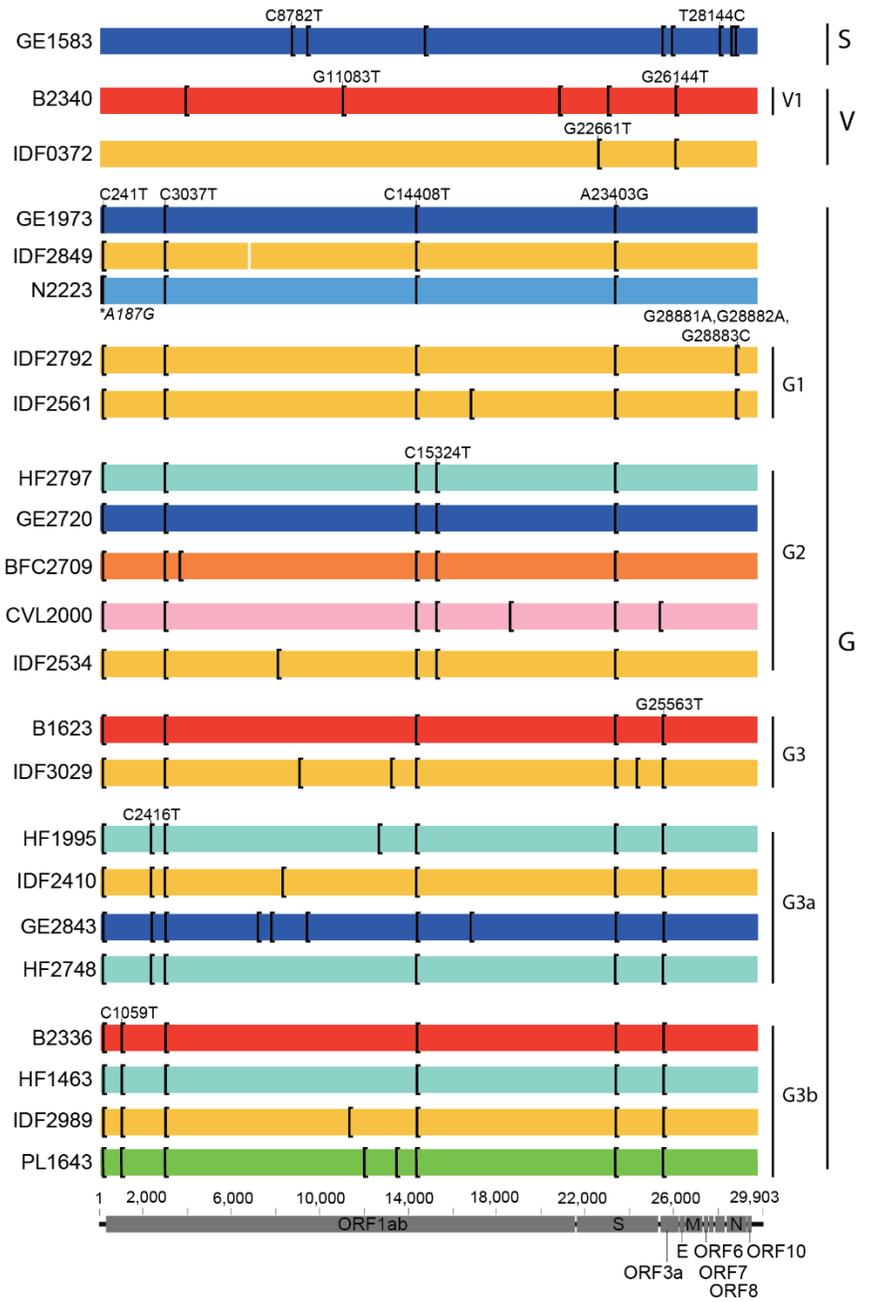
97  included in Data S1, Table S1.

98

nt substitutions from the root

100   **Fig. S1**. Phylogenetic divergence tree of all labeled SARS-CoV-2 sequences used in this

101   study. Maximum-likelihood tree including all sequences from Northern France, Algerian

102   sequences and publicly available global SARS-CoV-2 sequences, corresponding to the

103   collapsed tree shown in Fig. 3 (same ordering as in Fig. 2 and Fig. 3). GISAID clades are

104   indicated next to the corresponding nodes and branches are colored distinctly. Tips

105   indicate strain names, colored in red for sequences from France, and are noted in bold if

106   discussed in this study.

107

109 **Fig. S2**. Single-nucleotide polymorphisms representing the diversity among sequences

110 across the regions of Northern France. Multiple sequence alignment of all SARS-CoV-2

111 genomes sampled across the northern part of France from different clades and lineages.

112 Single nucleotide variants with respect to the reference (MN908947) are shown as black

113 vertical bars and shared substitutions among the sequences of each clade or lineage are

114 annotated. A substitution only found in sequences from Normandie is noted in italic.

115

116 **Data S1.** (Separate file)

117 A single Excel document with multiple sheets, representing tables below.

118 Table S1. Clade defining SNPs according to GISAID nomenclature.

119 Table S2. Metadata associated with the sequences used in this study.

120 Table S3. Viral RNA load and genome recovery data.

121 Table S4. List of collaborators in the RENAL network in the north of France.

122

123

**References:**

1. Quick J, Grubaugh ND, Pullan ST, Claro IM, Smith AD, Gangavarapu K, et al. Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus genomes directly from clinical samples. Nat Protoc. 2017;12(6):1261-76.

2. Kentaro I, Tsuyoshi S, Masanori H, Rina T, Makoto K. Aproposal ofan alternative primer fortheARTIC Network's multiplex PCRto improve coverageof SARS-CoV-2 genomesequencing. bioRxiv 2020.

3. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics. 2014;30(15):2114-20.

4. WysokerA F, RuanJ H, MarthG A. DurbinR. 2009a. The sequence alignment/map format and SAMtools. Bioinformatics.25(16):2078-9.

5. Fauver JR, Petrone ME, Hodcroft EB, Shioda K, Ehrlich HY, Watts AG, et al. Coast-to-Coast Spread of SARS-CoV-2 during the Early Epidemic in the United States. Cell. 2020;181(5):990-6 e5.

6. Deng X, Gu W, Federman S, du Plessis L, Pybus OG, Faria N, et al. Genomic surveillance reveals multiple introductions of SARS-CoV-2 into Northern California. Science. 2020.

7. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. Nextstrain: real-time tracking of pathogen evolution. Bioinformatics. 2018;34(23):4121-3.

8. Martin DP, Murrell B, Golden M, Khoosal A, Muhire B. RDP4: Detection and analysis of recombination patterns in virus genomes. Virus evolution. 2015;1(1).

9. Rambaut A, Lam TT, Max Carvalho L, Pybus OG. Exploring the temporal structure of heterochronous sequences using TempEst (formerly Path-O-Gen). Virus evolution. 2016;2(1):vew007.

10. Sagulenko P, Puller V, Neher RA. TreeTime: Maximum-likelihood phylodynamic analysis. Virus evolution. 2018;4(1):vex042.

148    11.    Rambaut A. Phylogenetic analysis of nCoV-2019 genomes http://virological.org/t/phylodynamic-

149    analysis-176-genomes-6-mar-2020/3562020

150

151