

Supplementary Material for Predicting the Emergence of SARS-CoV-2 Clades

Siddharth Jain¹, Xiongye Xiao², Paul Bogdan^{2,*}, and Jehoshua Bruck^{1,*}

¹California Institute of Technology, Electrical Engineering, Pasadena, 91125, USA

²University of Southern California, Electrical and Computer Engineering, Los Angeles, 90007, USA

*pbogdan@usc.edu, bruck@caltech.edu

ABSTRACT

Evolution is a process of change where mutations in the viral RNA are selected based on their fitness for replication and survival. Given that current phylogenetic analysis of SARS-CoV-2 identifies new viral clades after they exhibit evolutionary selections, one wonders whether we can identify the viral selection and predict the emergence of new viral clades? Inspired by the Kolmogorov complexity concept, we propose a generative complexity (algorithmic) framework capable to analyze the viral RNA sequences by mapping the multiscale nucleotide dependencies onto a state machine, where states represent subsequences of nucleotides and state-transitions probabilities encode the higher order interactions between these states. We apply computational learning and classification techniques to identify the active state-transitions and use those as features in clade classifiers to decipher the transient mutations (still evolving within a clade) and stable mutations (typical to a clade). As opposed to current analysis tools that rely on the edit distance between sequences and require sequence alignment, our method is computationally local, does not require sequence alignment and is robust to random errors (substitution, insertions and deletions). Relying on GISAID viral sequence database, we demonstrate that our method can predict clade emergence, potentially aiding with the design of medications and vaccines.

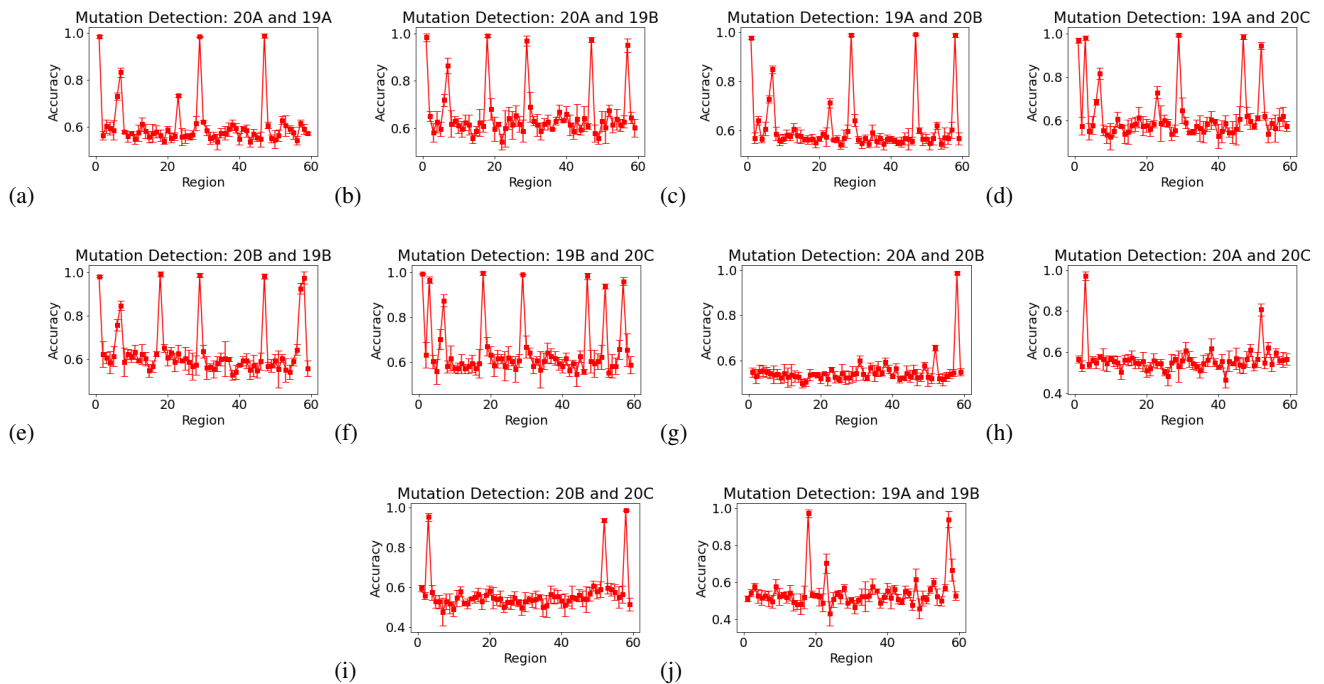


Figure S1. Mutation detection with region length 500. SARS-CoV-2 RNA is divided into continuous non-overlapping regions of length 500. Spikes represent the regions where mutations differentiating clade pairs 20A and 19A (a), 20A and 19B (b), 19A and 20B (c), 19A and 20C (d), 20B and 19B (e), 19B and 20C (f), 20A and 20B (g), 20A and 20C (h), 20B and 20C (i), 19A and 19B (j) are detected.

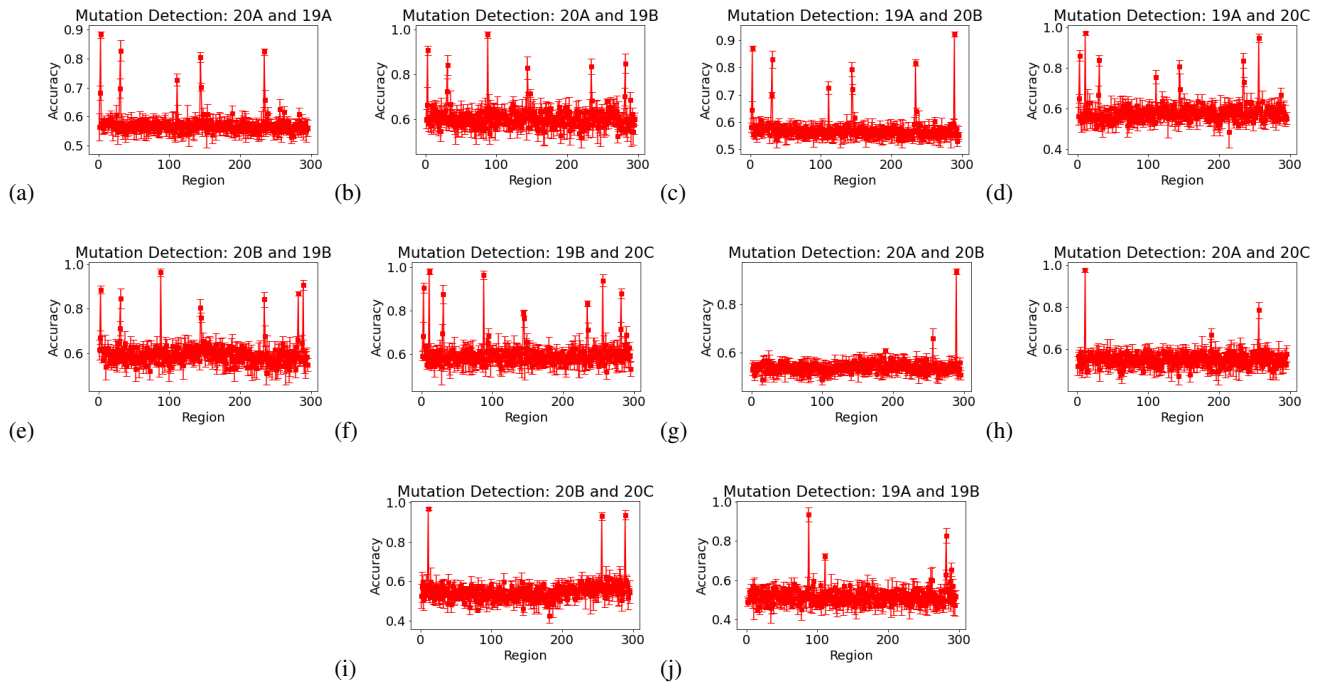


Figure S2. Mutation detection with region length 100. SARS-CoV-2 RNA is divided into continuous non-overlapping regions of length 100. Spikes represent the regions where mutations differentiating clade pairs 20A and 19A (a), 20A and 19B (b), 19A and 20B (c), 19A and 20C (d), 20B and 19B (e), 19B and 20C (f), 20A and 20B (g), 20A and 20C (h), 20B and 20C (i), 19A and 19B (j) are detected.

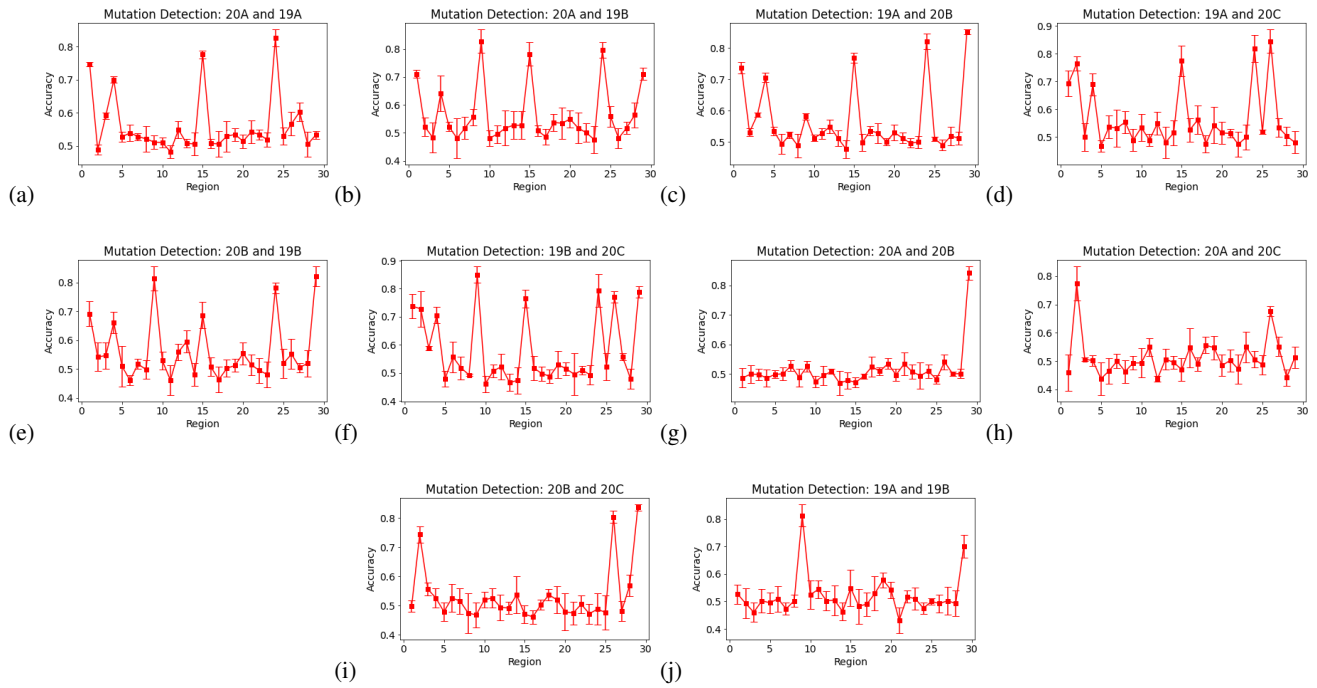


Figure S3. Mutation detection in the presence of synthetically added noise. We synthetically added 10% noise in the SARS-CoV-2 viral RNA sequences by adding substitution, insertion and deletion mutations. We show here that the mutations differentiating clades can still be detected, meaning that the detection using state transition probability matrix is robust to noise.

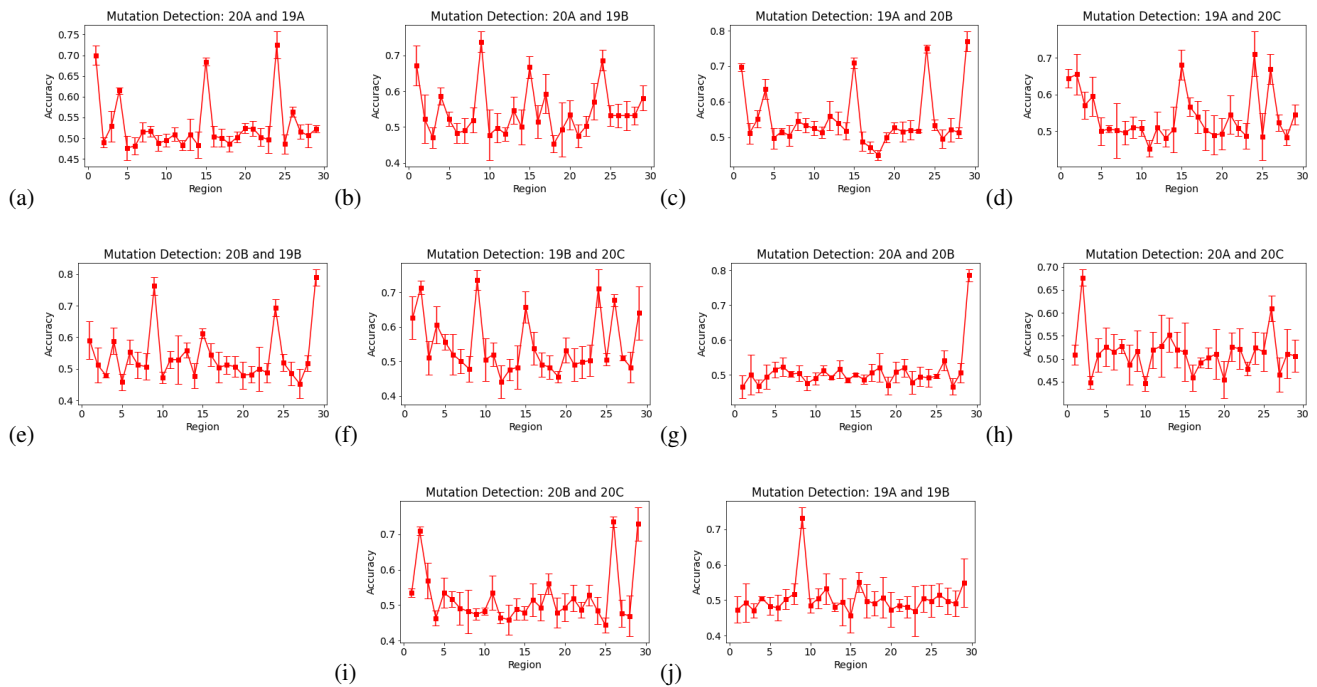


Figure S4. Mutation detection in the presence of synthetically added noise. We synthetically added 15% noise in the SARS-CoV-2 viral RNA sequences by adding substitution, insertion and deletion mutations. We show here that the mutations differentiating clades can still be detected, meaning that the detection using state transition probability matrix is robust to noise.

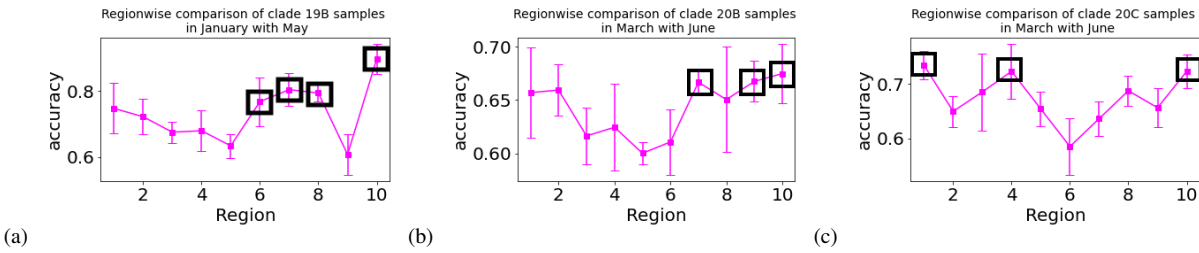


Figure S5. Evolution activity within clades 19B, 20B, 20C. (a) Regions 6 (15000-18000), 7 (18000-21000), 8 (21000-24000) and 10 (27000-30000) are actively mutating within clade 19B and could correspond to the location of mutations in the new clade that may emerge from 19B in the future. (b) Regions 7 (18000-21000), 9 (24000-27000) and 10 (27000-30000) contributes to maximum classification accuracy and hence are the most actively mutating within clade 20B which indicates the possibility of the emergence of a new clade from 20B with mutations in region 7,9 and 10. (c) Regions 1 (0-3000), 4 (9000-12000) and 10 (27000-30000) are the most actively mutating regions within clade 20C which indicates the possibility of emergence of a new clade from 20C with mutations in these regions.

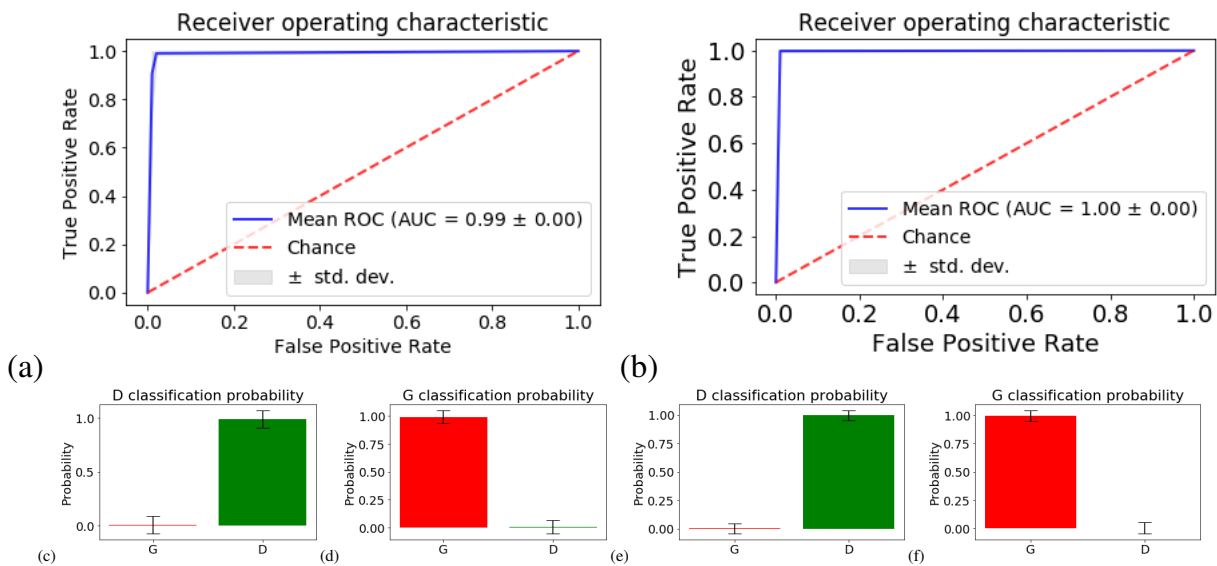


Figure S6. The state machine model can distinguish spike protein mutation sequences from others. We demonstrate here that the sequences with spike protein mutation (D614G) can be distinguished from others using the state machine transformation. Our state machine model accurately captures the higher order differences and classifies sequences with spike protein mutation from the others. (a), (b) ROC curves for the classifier built using the state transition probability as the feature vector to differentiate between sequences with spike mutation from the others. $k = 4, b = 1$ and $k = 5, b = 1$ for (a) and (b) respectively. (c)-(f) show the classification probability on the validation set for sequences with spike mutation and the others using classifier built with state transition probabilities as features.

Table S1. Highly Likely Transitions

probability ≥ 0.5
<i>aaagc</i> → <i>t</i> , <i>aaatc</i> → <i>a</i> , <i>aacga</i> → <i>a</i> , <i>aaactc</i> → <i>a</i> , <i>acctc</i> → <i>a</i> , <i>actat</i> → <i>t</i> , <i>actca</i> → <i>a</i> , <i>agcct</i> → <i>t</i> , <i>ataac</i> → <i>a</i> , <i>atacg</i> → <i>t</i> , <i>atcga</i> → <i>t</i> , <i>attcg</i> → <i>t</i> , <i>atttc</i> → <i>a</i> , <i>caacg</i> → <i>t</i> , <i>caagc</i> → <i>t</i> , <i>cacag</i> → <i>a</i> , <i>cagac</i> → <i>a</i> , <i>cagecg</i> → <i>t</i> , <i>catac</i> → <i>a</i> , <i>ccaag</i> → <i>a</i> , <i>cccat</i> → <i>t</i> , <i>cccca</i> → <i>a</i> , <i>cccgc</i> → <i>a</i> , <i>ccctc</i> → <i>a</i> , <i>ccgag</i> → <i>g</i> , <i>ccgcg</i> → <i>a</i> , <i>ccggg</i> → <i>t</i> , <i>ccggt</i> → <i>a</i> , <i>cctcc</i> → <i>a</i> , <i>cctgc</i> → <i>t</i> , <i>cctgg</i> → <i>t</i> , <i>cgaaa</i> → <i>t</i> , <i>cgact</i> → <i>a</i> , <i>cgatt</i> → <i>t</i> , <i>cgagc</i> → <i>a</i> , <i>cggat</i> → <i>g</i> , <i>cgtag</i> → <i>t</i> , <i>cgtat</i> → <i>a</i> , <i>cgtgg</i> → <i>t</i> , <i>cgtgt</i> → <i>t</i> , <i>ctaac</i> → <i>a</i> , <i>ctacc</i> → <i>a</i> , <i>ctagg</i> → <i>t</i> , <i>ctcgc</i> → <i>a</i> , <i>ctctc</i> → <i>a</i> , <i>gaagc</i> → <i>t</i> , <i>gaagg</i> → <i>t</i> , <i>gaccc</i> → <i>t</i> , <i>gacgg</i> → <i>t</i> , <i>gactc</i> → <i>a</i> , <i>gaggc</i> → <i>t</i> , <i>gagtg</i> → <i>t</i> , <i>gatgc</i> → <i>t</i> , <i>gatgg</i> → <i>t</i> , <i>gatta</i> → <i>t</i> , <i>gcccc</i> → <i>t</i> , <i>gcctc</i> → <i>a</i> , <i>gcgag</i> → <i>a</i> , <i>gcggc</i> → <i>a</i> , <i>gcgtc</i> → <i>a</i> , <i>gctac</i> → <i>t</i> , <i>gctcc</i> → <i>a</i> , <i>gctgc</i> → <i>t</i> , <i>gctgg</i> → <i>t</i> , <i>ggacc</i> → <i>t</i> , <i>ggagc</i> → <i>t</i> , <i>ggatc</i> → <i>a</i> , <i>ggcag</i> → <i>a</i> , <i>gggtg</i> → <i>t</i> , <i>gggtt</i> → <i>t</i> , <i>ggtga</i> → <i>t</i> , <i>ggtgc</i> → <i>t</i> , <i>ggtgg</i> → <i>t</i> , <i>ggtgt</i> → <i>t</i> , <i>gtacc</i> → <i>a</i> , <i>gtatc</i> → <i>t</i> , <i>gtccc</i> → <i>t</i> , <i>gtctc</i> → <i>t</i> , <i>gtcgc</i> → <i>t</i> , <i>gtggc</i> → <i>t</i> , <i>gttag</i> → <i>a</i> , <i>gttct</i> → <i>t</i> , <i>tactc</i> → <i>a</i> , <i>tagtc</i> → <i>t</i> , <i>tatcc</i> → <i>t</i> , <i>tatcg</i> → <i>t</i> , <i>tattc</i> → <i>t</i> , <i>tcacc</i> → <i>t</i> , <i>tcgag</i> → <i>g</i> , <i>tcggt</i> → <i>a</i> , <i>tctcg</i> → <i>t</i> , <i>tcgag</i> → <i>a</i> , <i>tgcga</i> → <i>a</i> , <i>tggac</i> → <i>a</i> , <i>tggtc</i> → <i>a</i> , <i>ttagc</i> → <i>t</i> , <i>ttata</i> → <i>a</i> , <i>ttatc</i> → <i>t</i> , <i>ttccc</i> → <i>t</i> , <i>ttcga</i> → <i>t</i> , <i>ttcgc</i> → <i>a</i> , <i>ttgac</i> → <i>a</i> , <i>ttggt</i> → <i>g</i> .
probability ≥ 0.55
<i>aacga</i> → <i>a</i> , <i>aaactc</i> → <i>a</i> , <i>agcct</i> → <i>t</i> , <i>ataac</i> → <i>a</i> , <i>atacg</i> → <i>t</i> , <i>attcg</i> → <i>t</i> , <i>caagc</i> → <i>t</i> , <i>cagecg</i> → <i>t</i> , <i>ccaag</i> → <i>a</i> , <i>cccat</i> → <i>t</i> , <i>cccca</i> → <i>a</i> , <i>cccgc</i> → <i>a</i> , <i>ccgag</i> → <i>g</i> , <i>ccgcg</i> → <i>a</i> , <i>ccggg</i> → <i>t</i> , <i>ccggt</i> → <i>a</i> , <i>cctcc</i> → <i>a</i> , <i>cctgc</i> → <i>t</i> , <i>cctgg</i> → <i>t</i> , <i>cgaaa</i> → <i>t</i> , <i>cgact</i> → <i>a</i> , <i>cgagc</i> → <i>a</i> , <i>cggat</i> → <i>g</i> , <i>cgtag</i> → <i>t</i> , <i>ctaac</i> → <i>a</i> , <i>ctacc</i> → <i>a</i> , <i>ctcgc</i> → <i>a</i> , <i>gaagc</i> → <i>t</i> , <i>gaagg</i> → <i>t</i> , <i>gaccc</i> → <i>t</i> , <i>gaggc</i> → <i>t</i> , <i>gcccc</i> → <i>t</i> , <i>gcgag</i> → <i>a</i> , <i>gcggc</i> → <i>a</i> , <i>gcgtc</i> → <i>a</i> , <i>gctac</i> → <i>t</i> , <i>gctcc</i> → <i>a</i> , <i>gctgc</i> → <i>t</i> , <i>gctgg</i> → <i>t</i> , <i>ggagc</i> → <i>t</i> , <i>ggatc</i> → <i>a</i> , <i>ggcag</i> → <i>a</i> , <i>gggtg</i> → <i>t</i> , <i>ggtgg</i> → <i>t</i> , <i>ggtgt</i> → <i>t</i> , <i>gtacc</i> → <i>a</i> , <i>gtatc</i> → <i>t</i> , <i>gtccc</i> → <i>t</i> , <i>gtctc</i> → <i>t</i> , <i>tatcg</i> → <i>t</i> , <i>tattc</i> → <i>t</i> , <i>tcggt</i> → <i>a</i> , <i>tctcg</i> → <i>t</i> , <i>tgcga</i> → <i>a</i> , <i>tggac</i> → <i>a</i> , <i>ttagc</i> → <i>t</i> , <i>ttcga</i> → <i>t</i> , <i>ttgac</i> → <i>a</i> .
probability ≥ 0.6
<i>atacg</i> → <i>t</i> , <i>caagc</i> → <i>t</i> , <i>ccgag</i> → <i>g</i> , <i>ccgcg</i> → <i>a</i> , <i>ccggg</i> → <i>t</i> , <i>cctcc</i> → <i>a</i> , <i>cggac</i> → <i>a</i> , <i>cggat</i> → <i>g</i> , <i>cgtag</i> → <i>t</i> , <i>ctaac</i> → <i>a</i> , <i>ctacc</i> → <i>a</i> , <i>gaagg</i> → <i>t</i> , <i>gaccc</i> → <i>t</i> , <i>gaggc</i> → <i>t</i> , <i>gcccc</i> → <i>t</i> , <i>gcgag</i> → <i>a</i> , <i>gctcc</i> → <i>a</i> , <i>gctgg</i> → <i>t</i> , <i>ggtgt</i> → <i>t</i> , <i>gtccc</i> → <i>t</i> , <i>gtctc</i> → <i>t</i> , <i>tatcg</i> → <i>t</i> , <i>ttagc</i> → <i>t</i> .
probability ≥ 0.65
<i>ccggg</i> → <i>t</i> , <i>cggat</i> → <i>g</i> , <i>gctcc</i> → <i>a</i> , <i>gtccc</i> → <i>t</i> , <i>gtctc</i> → <i>t</i> .
probability ≥ 0.7
<i>gtccc</i> → <i>t</i> .

Table S2. Mutations differentiating different clades (obtained from nextstrain.org)

Clades	Mutations
19A	Root clade
19B	C8782T, T28144C
20A	C14408T, A23403G
20B	C14408T, A23403G, G28881A, G28882A, G28883C
20C	C14408T, A23403G, C1059T, G25563T