

# Exploring the coronavirus epidemic using the new WashU Virus Genome Browser

Jennifer A. Flynn<sup>1\*</sup>, Deepak Purushotham<sup>1\*</sup>, Mayank NK Choudhary<sup>1\*</sup>, Xiaoyu Zhuo<sup>1\*</sup>, Changxu Fan<sup>1\*</sup>, Gavriel Matt<sup>1\*</sup>, Daofeng Li<sup>1†</sup> and Ting Wang<sup>1,2†</sup>

\* These authors contributed equally to this work.

† These authors jointly supervised this work. Co-corresponding author emails: [dli23@wustl.edu](mailto:dli23@wustl.edu) and [twang@genetics.wustl.edu](mailto:twang@genetics.wustl.edu)

<sup>1</sup>The Edison Family Center for Genome Sciences & Systems Biology, Department of Genetics, Washington University, 4515 McKinley Avenue, Campus Box 8510, St. Louis, MO 63110, USA

<sup>2</sup>McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO 63108, USA

## Abstract

Since its debut in mid-December, 2019, the novel coronavirus (2019-nCoV) has rapidly spread from its origin in Wuhan, China, to several countries across the globe, leading to a global health crisis. As of February 7, 2020, 44 strains of the virus have been sequenced and uploaded to NCBI's GenBank [1], providing insight into the virus's evolutionary history and pathogenesis. Here, we present the WashU Virus Genome Browser, a web-based portal for viewing virus genomic data. The browser is home to 16 complete 2019-nCoV genome sequences, together with hundreds of related viral sequences including severe acute respiratory syndrome coronavirus (SARS-CoV), Middle East respiratory syndrome coronavirus (MERS-CoV), and Ebola virus. In addition, the browser features unique customizability, supporting user-provided upload of novel viral sequences in various formats. Sequences can be viewed in both a track-based representation as well as a phylogenetic tree-based view, allowing the user to easily compare sequence features across multiple strains. The WashU Virus Genome Browser inherited many features and track types from the WashU Epigenome Browser, and additionally incorporated a new type of SNV track to address the specific needs of viral research. Our Virus Browser portal can be accessed at <https://virusgateway.wustl.edu>, and documentation is available at <https://virusgateway.readthedocs.io/>.

## Introduction

On December 12, 2019, the first case of a novel coronavirus (2019-nCoV) was reported in Wuhan, China, and by February 6, 2020, the virus spread to 24 additional countries, infecting more than 27,000 individuals and resulting in 565 fatalities, according to the World Health Organization (WHO) [2]. The 2019-nCoV is a member of the *Betacoronavirus* genus, which is one of four genera of coronaviruses of the subfamily Orthocoronavirinae in the family Coronaviridae, of the order Nidovirales [3, 4]. The species in this genus are enveloped, contain a positive single-stranded RNA genome, and are of zoonotic, likely bat, origins [5]. 2019-nCoV is one of the largest RNA virus genomes varying from 27kb to 32kb in size, with this particular strain ringing in at

41 29,903 bps long [6]. The virus is one of 7 coronaviruses known to infect humans, and along with  
42 the severe acute respiratory syndrome coronavirus (SARS-CoV) and the Middle East respiratory  
43 syndrome coronavirus (MERS-CoV), 2019-nCoV is one of the species responsible for severe  
44 respiratory distress in humans as well as other animals [4]. In an effort to better understand the  
45 pathogenesis of this family of viruses, several groups have sequenced individual strains, providing  
46 a powerful resource hosted by NCBI.

47  
48 The WashU Epigenome Browser is a powerful tool for visualizing multiple functional genomic  
49 datasets and data types simultaneously [5-8]. The general layout of the Epigenome Browser  
50 displays the genome on the x-axis, and individual tracks encompassing many different varieties  
51 can be loaded and viewed in the context of the genome and accompanying metadata. Recent  
52 updates to the browser have incorporated new functionality, including live browsing, greatly  
53 enhancing its functionality [5]. With this powerful tool in-hand, we sought to adapt the browser for  
54 use of visualizing viral genomes, to support more efficient research and more rapid knowledge  
55 dissemination in response to the recent 2019-nCoV outbreak. To accomplish this, we created the  
56 WashU Virus Genome Browser, adapted from the WashU Epigenome Browser. The Virus  
57 Genome Browser houses reference genomes for 2019-nCoV, MERS, SARS, and Ebola virus,  
58 along with several annotation tracks including gene annotation, putative antibody-binding  
59 epitopes, CG density, and sequence diversity. Complete genomes of individual strains of each  
60 virus species (16, 551, 332, and 1574, respectively as of February 7, 2020, and periodically  
61 updated) are available as a database for instant viewing on the Virus Browser via multiple track  
62 types designed to display pairwise comparison to the references. Additionally, we aligned the  
63 genomes of all available strains in the database and generated a phylogenetic tree for each virus  
64 species that allows the user to directly select strains from the tree and view as tracks in the  
65 genomic display. In addition to all track types supported by the Epigenome Browser, we designed  
66 a new SNV track type to display sequence variation. Users can upload their own alignment results  
67 from any aligner and display them as SNV tracks on the browser.

68  
69 The functionality of the Virus Browser is not limited to the 4 species currently housed. Users can  
70 upload their own reference genome in FASTA format and display tracks in the context of the user-  
71 specified reference. While maintaining the same functionality as that of the Epigenome Browser  
72 and providing novel functionality to aid specifically in viral genome research, we hope that the  
73 Virus Browser may facilitate research against new epidemic viruses.

74

## 75 Materials and Methods

### 76 Reference sequences, additional strains, and gene annotations:

77 Genomic sequences of all viral strains were downloaded as FASTA files from NCBI  
78 [Supplementary Table 1]. All available sequences as of January 31, 2020, for 2019-nCoV, MERS,  
79 SARS, and Ebola were downloaded (n=16, 551, 332, and 1574, respectively). The reference  
80 genomic sequence of the selected virus (2019-nCoV: NC\_045512.2; MERS: NC\_019843.3;  
81 SARS: NC\_004718.3, Ebola: KM034562.1) is automatically displayed as a color coded track  
82 when opening the genomic track browser viewing format. Genic annotations of reference  
83 genomes were downloaded as GFF3 files from NCBI and converted to refBed format for viewing  
84 on the browser.

## 85 Sequence alignment and tree generation:

86 The genomes of all individual strains of each virus were aligned to the reference genome using  
87 the pairwise alignment tool stretcher [9] with parameters “-gapopen 16 -gapextend 4”. To generate  
88 the phylogenetic trees, we used the MAFFT program, employing the fast option to align individual  
89 strains of each viral genome to its reference [10, 11]. Phylogenetic trees were built using FastTree  
90 with the GTR model [12, 13].

## 91 Data Tracks:

### 92 Genome Comparison Track:

93 We adopted the genome comparison tracks from the WashU Epigenome Browser. Any pairwise  
94 alignment results in markx3 or FASTA format can be converted with our publicly accessible script  
95 “aligned\_fa\_2\_genomealign.py” [14] and directly displayed as genome comparison tracks on the  
96 Virus Browser.

### 97 SNV Track:

98 We developed the SNV track type to display sequence variation of individual strains relative to  
99 their reference. Variations from the reference genome, including mismatches and deletions, are  
100 displayed with customizable colors. Insertions compared to the reference genome can be  
101 expanded upon selecting to show the nucleotides inserted. When viewing large regions, such as  
102 the whole genome, it is not possible to display all individual variation events. Therefore, the  
103 frequency of variation events is also displayed in a “density mode” where a high value over a  
104 region signifies multiple sequence variation events within the region.

## 105 Congeneric (or Closely-related) Immune Epitope Locations:

106 We wrote a text processing utility to import antibody-binding epitopes curated by the Immune  
107 Epitope Database and Analysis Resource (IEDB) for MERS-CoV and SARS-CoV [15].  
108 Subsequently, we used tblastn to align linear epitopes to the Wuhan seafood market pneumonia  
109 virus isolate Wuhan-Hu-1 (Taxonomy ID: 2697049; NCBI:txid2697049). We found 955 out of  
110 2,817 linear epitopes identified in SARS had at least 1 “hit” in the 2019-nCoV genome  
111 [Supplementary Data 1]. Three epitopes have 2 “hits” each. However, the secondary hit is on the  
112 negative strand with very low percent identity (37.5% to 53.8%) to the 2019-nCoV genome and  
113 are hence filtered out as 2019-nCoV is a (+) ssRNA virus. Similarly, we found 1 hit out of 38 linear  
114 epitopes identified in MERS. We also provide scripts [14] that can be used to obtain a quick  
115 overview of the similarity of linear epitopes identified in other viruses in databases like IEDB.  
116 These tracks can provide researchers preliminary data to support exploratory analyses pertaining  
117 to the immunogenicity of 2019-nCoV—an actively explored vertical of 2019-nCoV research.

### 118 GC Density Track:

119 GC density tracks were created for each reference genome, displaying the percentage of G  
120 (guanine) and C (cytosine) bases in 5-bp windows.

## 121 Sequence Diversity Track and Shannon Track:

122 In order to display a measure of sequence conservation across the genome, we calculated the  
123 percentage of each of the 4 nucleotides at each position in the genome across all strains for a  
124 given virus species. The resulting bed tracks display the percentages each nucleotide comprises  
125 across all strain for each genomic position. We also calculated Shannon entropy for each position  
126 along the genome using the percentages of each of the 4 nucleotides. A high Shannon entropy  
127 at a position signifies that the 4 possible nucleotides are equally likely across all strains of this  
128 virus, and thus the position is likely divergent. A low Shannon entropy at a position means that  
129 the identity of the nucleotide at this position is highly conserved across all strains. The entropy()  
130 function of the R package “entropy” was used for calculations.

## 131 Resources for User-Defined Bed and Categorical Tracks:

132 In addition to our housed data tracks, we also offer scripts (“publicParseAlignment.py”,  
133 “publicAlignment.py”, and “publicConvertMarkx3.py”) to convert any markx3 or FASTA-formatted  
134 alignment into displayable bed and categorical formats, and a script (“publicJsonGen.py”) to  
135 generate a json file for uploading multiple data files together for display  
136 [<https://github.com/debugpoint136/WashU-Virus-Genome-Browser>]. A default color code for  
137 sequence variation is also included in the script.

# 138 Results

## 139 Organization of the Virus Genome Browser

140 The WashU Virus Genome Browser houses consensus reference genomic sequences for 4  
141 different pathogenic virus species: 2019-nCoV, MERS, SARS, and Ebola, as well as a  
142 comprehensive set of genome assemblies for the individual strains of each virus (16, 551, 332,  
143 and 1574, respectively). When users first navigate to the WashU Virus Browser and select  
144 “Browse Data”, they are directed to a page with several customizable options, including a drop-  
145 down menu from which they may choose a reference genome [Figure 1]. Corresponding with the  
146 reference genome selected, a metadata table is displayed containing sortable features such as  
147 species, strain, isolate, isolation source, host, country, and collection date, to allow for quick and  
148 easy sorting of individual strains. The user may select viral isolates from the metadata table to be  
149 visualized in one of our two displayable platforms: the track view (green arrow, Figures 2 and 3)  
150 or the phylogenetic tree view (orange arrow, Figures 4 and 5).

## 151 The Track View

152 The track view option has a standard genome browser layout similar to that of the WashU  
153 Epigenome Browser, in which a reference genome sequence is visualized as a sliding window.  
154 Various annotation data tracks are hosted on the browser and can be loaded for visualization in  
155 a genomic context. For each virus, we downloaded publicly available annotations of the reference  
156 genome and converted these annotations into refBed tracks that can be visualized in the genome  
157 browser. Likewise, immune epitopes identified in SARS were aligned to the 2019-nCoV reference  
158 [Materials and Methods], and a track displaying their coordinates in 2019-nCoV is provided. GC-  
159 density tracks were also created for each reference genome, and display the percentage of Gs  
160 (Guanines) and Cs (Cytosines) per 5bp window. An entropy track [Materials and Methods]  
161 showing the degree of sequence diversity at each position and a diversity track [Materials and  
162 Methods] showing the percentage of each of the 4 nucleotides at each position across all strains  
163 of the given virus species are also included in the database. In addition to hosting 4 virus species

164 reference genomes, The Virus Genome Browser also supports displaying user-specified  
165 genomes provided in FASTA format, as shown in the top left part of Figure 2A, under the browser  
166 logo.

167  
168 The WashU Virus Browser supports a “zoomed-out” view of the entire viral genome. The zoomed-  
169 out view can help the user quickly determine the regions of interest that have high frequencies of  
170 variation from the reference (SNV track), and also the regions with high nucleotide diversity  
171 among all strains (Shannon tracks) [Figure 2A]. Figure 2A illustrates a genome-level browser view  
172 of the 2019-nCoV reference genome and 2 SARS strains, each aligned to the SARS reference  
173 genome (AY278488.2 = BJ01, DQ071615.1 = Bat rp3, NC\_045512.2 = 2019-nCoV). Sequence  
174 variation displayed in density mode [Materials and Methods] shows that the divergence between  
175 the 2019-nCoV reference genome (red) and the SARS reference genome is higher than the  
176 divergence between the two additional SARS strains (green) and the SARS reference genome.  
177 For AY278488.2, the variation from reference is mainly confined to the beginning of the genome,  
178 while the remainder of the genome is relatively consistent with the reference. However, for  
179 DQ071615.1 (bat-derived), the 5' end of gene S displays high variation from the reference  
180 genome. Likewise, the SARS Shannon track shows that the SARS genome is highly diverse  
181 across different strains at gene S.

182  
183 Once a region of interest is identified, the standard magnification tool of the browser can be used  
184 to quickly zoom into the region [Figure 2A]. Upon zooming in, a genome comparison track can be  
185 used to inspect variations from the reference genome, particularly useful for comparing cross-  
186 species alignments and viewing structural variations [Figure 2B]. The genome comparison track  
187 is adopted from the Epigenome Browser. The top navy-colored horizontal bar represents the  
188 reference genome loaded (SARS in the case of Figure 2B) and the bottom purple-colored  
189 horizontal bar represents the sequence being aligned to the reference (the 2019-nCoV reference  
190 sequence, NC\_045512.2, in this case). Insertions and deletions are represented as gaps in either  
191 the reference or the query. Matches are represented by black lines linking the 2 genomes while  
192 mismatches are distinguished by omission of the black bar. When the user hovers over a specific  
193 nucleotide, the alignment details around that specific nucleotide are shown.

194  
195 Upon further magnification, regions can be inspected on a nucleotide level. Mismatches,  
196 insertions, and deletions are color-coded in the SNV tracks and stretches of grey signify positions  
197 matching the reference [Figure 2C]. Detailed information, such as inserted nucleotides, is  
198 displayed upon clicking. When zoomed into individual nucleotides, as shown in Figure 2C, The  
199 diversity bed track shows the percentage of each nucleotide across all strains of SARS at the  
200 specific position.

201  
202 The versatility of the WashU browser framework makes it possible to adapt the browser to address  
203 various questions of interest. Figure 3 demonstrates the utility of using the browser for immune  
204 epitope conservation discovery. We recapitulated Zhou et al.'s [16] alignment results of two SARS  
205 strains to the reference 2019-nCoV nucleocapsid protein sequence [Figure 3A, 3B]. Upon  
206 inspection of the region, we could directly observe that many immune epitopes are conserved  
207 between SARS and 2019-nCoV [Figure 3C]. The user can identify the amino acid sequence of an  
208 epitope by simply clicking the track.

209  
210 Encouraged by the high sequence similarity between SARS-CoV and the 2019-nCoV reference  
211 strain (NCBI:txid2697049), we mined the list of experimentally identified linear epitopes from T-  
212 cell, B-cell and MHC-ligand assays from IEDB [15]. We identified a list of 320 high-confidence  
213 linear epitopes [Supplementary Table 2] whose amino acids are identical to predicted translated

214 products from the 2019-nCoV reference strain. These provide a catalogue of epitopes for  
215 researchers testing immune targets that can potentially elicit T-cell, B-cell and antibody response  
216 to 2019-nCoV.

217  
218 We also provide these as an annotated bed track to the reference 2019-nCoV genome. Along  
219 with the individual strains' SNV tracks, the epitope tracks can provide a quick, intuitive and visual  
220 resource to guide prioritization of experimental resources towards developing diagnostics and  
221 therapeutics against 2019-nCoV. The value of our novel SNV tracks will only increase as  
222 additional strains are sequenced, helping us better understand the evolving 2019-nCoV genome  
223 and prioritize epitopes.

224

## 225 The Phylogenetic Tree View

226 The second viewing option offered by the WashU Virus Genome Browser is a "tree" format, in  
227 which the evolutionary relationships of different viral isolates can be visualized as a phylogenetic  
228 tree [17]. When the user navigates to the data page of the browser, and selects "Tree View"  
229 [Figure 1], all viral genomes hosted on the browser for the selected virus species are displayed in  
230 the form of a right-aligned phylogenetic tree, where solid lines indicate branch lengths [Figure 4].  
231 To the right of the tree is a metadata heatmap displaying strain-specific details such as isolate,  
232 isolation source, host, country, and collection date. Additionally, if the user added any individual  
233 tracks to their cart from the main page, those selected will display a checkmark to the right,  
234 allowing the user to easily see where their strains of interest lie among all other strains.

235

236 In addition to the right-aligned tree view, the browser also supports a more traditional left-aligned  
237 linear tree view and a radial view. The left-aligned tree view displays branch lengths indicating  
238 relatedness of isolates [Figure 5A]. We noticed that in each virus type, several individual strains  
239 maintained high sequence similarity, resulting in several short branch lengths and a long vertical  
240 tree. In order to improve visualization, we also created a radial tree view [Figure 5B].

241

## 242 Discussion

243 Maps help us understand the world around us and navigate it. Moreover, they play a critical role  
244 in disaster management during disease outbreaks. Herein, we describe the first genetic mapping,  
245 exploration, and visualization tool from the WashU Epigenome Browser team that is specifically  
246 dedicated to viral genomes. We provide reference genome maps and genomic datasets related  
247 to 4 viral disease outbreaks: SARS (2002-03), MERS (2012), Ebola (2014-16) and the latest nCoV  
248 (2019-20). More importantly, we not only present publicly available information in the format of  
249 easily accessible data tracks, but also offer a platform with high customizability and flexibility  
250 where individual investigators and teams can upload and visualize their own genomic datasets in  
251 a plethora of formats. In this report, we have demonstrated using the Virus Browser to 1) quickly  
252 and intuitively compare multiple viral genomes and study the viral genome at multiple levels  
253 [Figure 2, Figure 4, Figure 5]; and 2) combine viral genome information with other functional  
254 genomic information (amino acid sequence and putative immune epitope locations, as shown  
255 Figure 3) through multiple track types the browser supports, and identify potential therapeutic  
256 targets.

257

258 We expect that the WashU Virus Browser can support research related to the latest novel  
259 Coronavirus outbreak of 2019-20, and hope that this tool helps accelerate research to further our  
260 understanding of 2019-nCoV and aid in the development of therapeutics. In addition, our platform  
261 supports the study of any user-specified viral genome, and can be expanded to other viral  
262 research.

263  
264 To aid in the battle against this crisis, we are releasing the browser at first moment. The browser  
265 is still under active construction and is constantly being updated. General feedback, suggestions  
266 for additional tracks, and bug reports may be sent to the WashU Virus Genome Browser team  
267 by opening an issue request at [https://github.com/debugpoint136/WashU-Virus-Genome-  
268 Browser/issues](https://github.com/debugpoint136/WashU-Virus-Genome-Browser/issues).

269  
270  
271

## 272 References

- 273  
274  
275 1. **NCBI GeneBank** [<https://www.ncbi.nlm.nih.gov/genbank/2019-ncov-seqs/>]  
276 2. **Novel Coronavirus (2019-nCoV) Situation Report - 17**  
277 [[https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200206-  
278 sitrep-17-ncov.pdf?sfvrsn=17f0dca\\_4](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200206-sitrep-17-ncov.pdf?sfvrsn=17f0dca_4)]  
279 3. **International Committee on Taxonomy of Viruses (ICTV)**  
280 [<https://talk.ictvonline.org/taxonomy/>]  
281 4. Cui J, Li F, Shi ZL: **Origin and evolution of pathogenic coronaviruses.** *Nat Rev Microbiol*  
282 2019, **17**(3):181-192.  
283 5. Li D, Hsu S, Purushotham D, Sears RL, Wang T: **WashU Epigenome Browser update**  
284 **2019.** *Nucleic Acids Res* 2019, **47**(W1):W158-W165.  
285 6. Zhou X, Li D, Zhang B, Lowdon RF, Rockweiler NB, Sears RL, Madden PA, Smirnov I,  
286 Costello JF, Wang T: **Epigenomic annotation of genetic variants using the Roadmap**  
287 **Epigenome Browser.** *Nature biotechnology* 2015, **33**(4):345-346.  
288 7. Zhou X, Lowdon RF, Li D, Lawson HA, Madden PA, Costello JF, Wang T: **Exploring long-**  
289 **range genome interactions using the WashU Epigenome Browser.** *Nat Methods* 2013,  
290 **10**(5):375-376.  
291 8. Zhou X, Maricque B, Xie M, Li D, Sundaram V, Martin EA, Koebbe BC, Nielsen C, Hirst M,  
292 Farnham P *et al*: **The Human Epigenome Browser at Washington University.** *Nat*  
293 *Methods* 2011, **8**(12):989-990.  
294 9. Myers EW, Miller W: **Optimal alignments in linear space.** *Comput Appl Biosci* 1988,  
295 **4**(1):11-17.  
296 10. Katoh K, Misawa K, Kuma K, Miyata T: **MAFFT: a novel method for rapid multiple**  
297 **sequence alignment based on fast Fourier transform.** *Nucleic Acids Res* 2002,  
298 **30**(14):3059-3066.  
299 11. Katoh K, Standley DM: **MAFFT multiple sequence alignment software version 7:**  
300 **improvements in performance and usability.** *Mol Biol Evol* 2013, **30**(4):772-780.

- 301 12. Price MN, Dehal PS, Arkin AP: **FastTree: computing large minimum evolution trees with**  
302 **profiles instead of a distance matrix.** *Mol Biol Evol* 2009, **26**(7):1641-1650.
- 303 13. Price MN, Dehal PS, Arkin AP: **FastTree 2--approximately maximum-likelihood trees for**  
304 **large alignments.** *PLoS One* 2010, **5**(3):e9490.
- 305 14. **Virus Browser Source Code** [[https://github.com/debugpoint136/WashU-Virus-Genome-](https://github.com/debugpoint136/WashU-Virus-Genome-Browser)  
306 [Browser](https://github.com/debugpoint136/WashU-Virus-Genome-Browser)]
- 307 15. Vita R, Mahajan S, Overton JA, Dhanda SK, Martini S, Cantrell JR, Wheeler DK, Sette A,  
308 Peters B: **The Immune Epitope Database (IEDB): 2018 update.** *Nucleic Acids Res* 2019,  
309 **47**(D1):D339-D343.
- 310 16. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL *et al*: **A**  
311 **pneumonia outbreak associated with a new coronavirus of probable bat origin.** *Nature*  
312 2020.
- 313 17. Shank SD, Weaver S, Kosakovsky Pond SL: **phylotree.js - a JavaScript library for**  
314 **application development and interactive data visualization in phylogenetics.** *BMC*  
315 *Bioinformatics* 2018, **19**(1):276.  
316

## 317 Acknowledgements

318 We thank doctors, nurses, investigators, and all other people fighting on the front line against  
319 this viral outbreak, and we sincerely hope that this tool will aid in this battle.

## 320 Author Contribution:

321 Conceptualization, T.W. Web development, D.L and D.P. SNV track development, J.F. and C.F.  
322 Immune epitope analysis, M.C. Data download, metadata generation and annotation, G.M.  
323 Sequence alignments and tree generation, X.Z. Manuscript preparation, J.F, C.F, M.C, G.M,  
324 T.W.  
325

## 326 Author Support:

327 J.F. is supported in part by the Siteman Cancer Center Precision Medicine Pathway.  
328 X.Z. is supported in part by 5R25DA027995.  
329 TW is supported by NIH grants R01HG007175, U24ES026699, U01CA200060,  
330 U01HG009391, and U41HG010972, and by the American Cancer Society Research Scholar  
331 grant RSG-14-049-01-DMC.  
332

## 333 Figure Captions

334 **Figure 1:** Screenshot of the WashU Virus Genome Browser data page. This view demonstrates  
335 several customizable features of the browser, including which genome reference to use, which  
336 data tracks to select based on several metadata features, and which browser view to use:  
337 “genomic” view (green arrow) or phylogenetic tree view (orange arrow).



338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366

**Figure 2:** Illustration of genomic-level and nucleotide-level track views. A: “zoomed out” track view of the entire genome. 2019-nCoV reference genome (shown in red, NC045512.2) and 2 SARS strains (shown in green, DQ071615.1 and AY278488.2) are aligned to the SARS reference genome (NC\_004718.3). The box in the top left corner allows users to upload and use any sequence in FASTA format as the reference genome. The shaded vertical bar demonstrates the user’s ability to select a region by mouse for further magnification. B: “Zoomed in” view of the sequence flanking the 5’ end of the S protein. C: A further “zoomed in” view to the level of individual nucleotides. Stretches of grey indicate matching while variations are color coded.

**Figure 3:** Alignment of the genomic region encoding the nucleocapsid protein. A: 2 SARS strains (DQ071615.1 and AY278488.2) and 5 2019-nCoV strains (MN938384.1, MN975262.1, MN985325.1, MN988668.1, and MN988669.1) are aligned to the 2019-nCoV reference. The region encoding the nucleocapsid protein is shown. Putative SARS immune epitopes [Materials and Methods] are displayed in “density mode”. B: A zoomed-in view of A (orange box), displaying the first 9 amino acids of the reference. Results show a “TCA” insertion in the AY278488.2 alignment between positions 28294 and 28295 of the 2019-nCoV reference sequence, which is not present in DQ071615.1. These results are consistent with the results reported in Extended Data Figure 5 of Zhou et al. [16]. C: A zoomed-in view of A (purple box), displaying a region conserved between SARS and 2019-nCoV, overlapping several putative immune epitopes.

**Figure 4:** Screenshot of a linear, right-aligned tree view displaying all housed 2019-nCoV sequences with accompanying metadata. Solid lines signify distance.

**Figure 5:** A: Screenshot of a linear, left-aligned phylogenetic tree view, displaying all 2019-nCoV strains hosted by the browser. B: Screenshot of a radial tree view for all 2019-nCoV strains.

# Figure 1

WashU Virus Genome Browser

nCov Reference

TREE VIEW

DATA

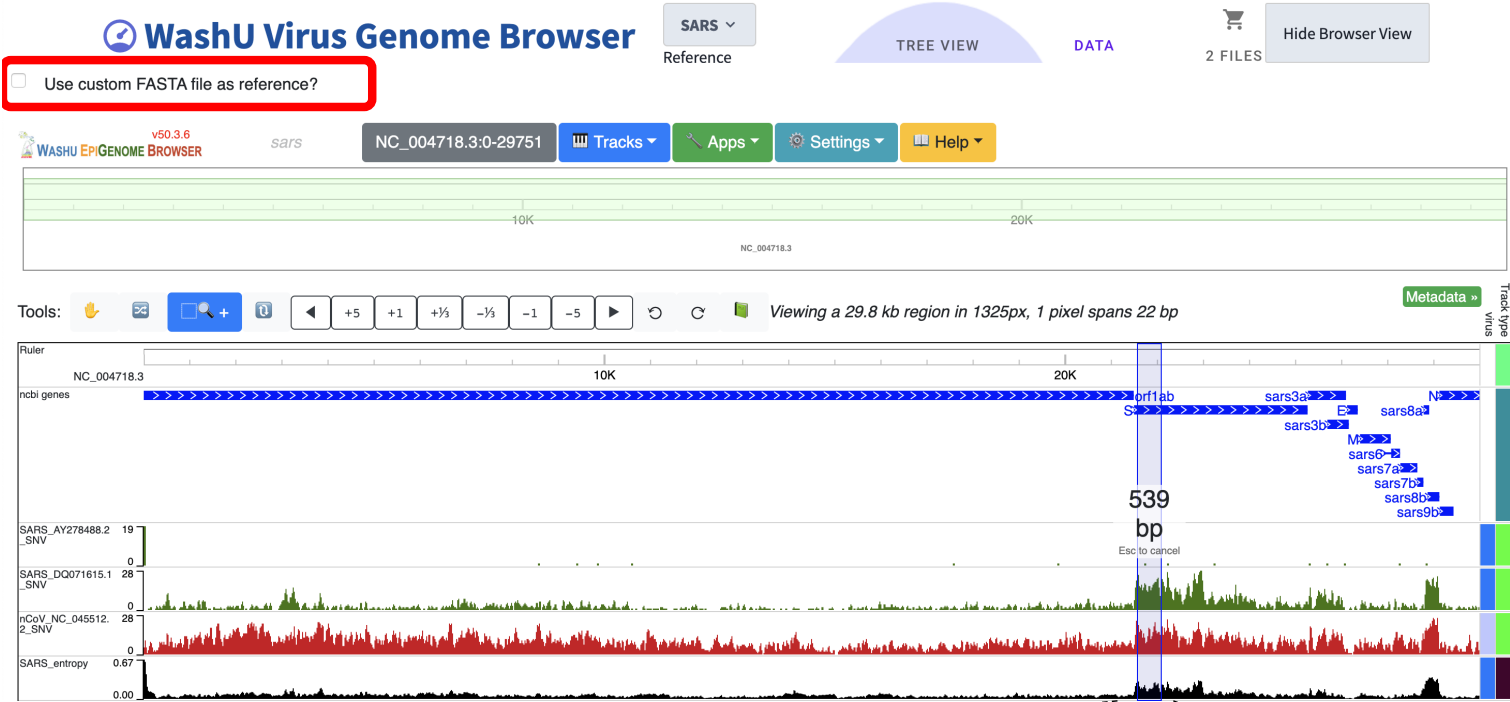
0 FILES

Show Browser View

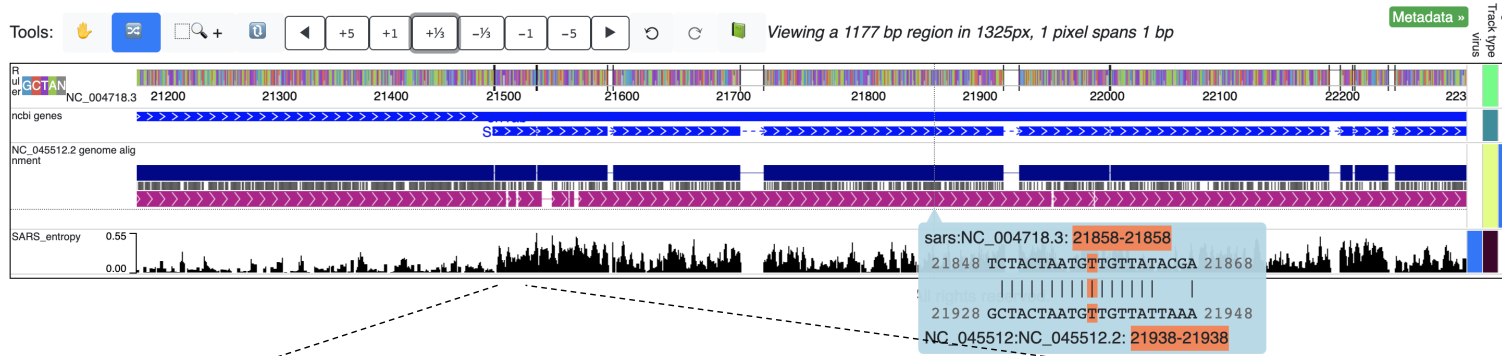
ID	Accession	Isolate	Molecule Type	Country	Collection Date
1	NC_045512.2	Wuhan-Hu-1	genomic RNA	China	Dec-2019
2	MN938384.1	2019-nCoV_HKU-SZ-002a_2020	genomic RNA	China: Shenzhen	Jan-2020
3	MN975262.1	2019-nCoV_HKU-SZ-005b_2020	genomic RNA	China	Jan-2020
4	MN985325.1	2019-nCoV/USA-WA1/2020	genomic RNA	USA	19-Jan-2020
5	MN988713.1	2019-nCoV/USA-IL1/2020	genomic RNA	USA: Illinois	21-Jan-2020
6	MN994467.1	2019-nCoV/USA-CA1/2020	genomic RNA	USA: CA	23-Jan-2020
7	MN994468.1	2019-nCoV/USA-CA2/2020	genomic RNA	USA: CA	22-Jan-2020
8	MN997409.1	2019-nCoV/USA-AZ1/2020	genomic RNA	USA: AZ	22-Jan-2020
9	MN988668.1	2019-nCoV WHU01	genomic RNA	China	02-Jan-2020
10	MN988669.1	2019-nCoV WHU02	genomic RNA	China	02-Jan-2020
11	MN996527.1	WIV02	genomic RNA	China: Wuhan	30-Dec-2019
12	MN996528.1	WIV04	genomic RNA	China: Wuhan	30-Dec-2019
13	MN996529.1	WIV05	genomic RNA	China: Wuhan	30-Dec-2019
14	MN996530.1	WIV06	genomic RNA	China: Wuhan	30-Dec-2019
15	MN996531.1	WIV07	genomic RNA	China: Wuhan	30-Dec-2019
16	MT007544.1	Australia/VIC01/2020	genomic RNA	Australia: Victoria	25-Jan-2020

# Figure 2

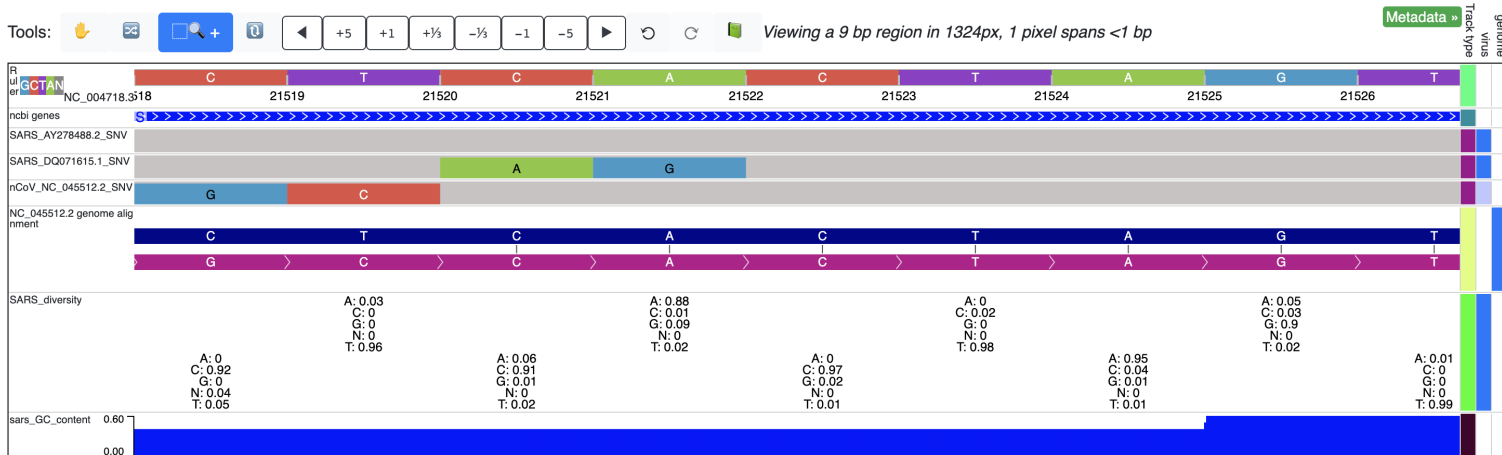
## A



## B

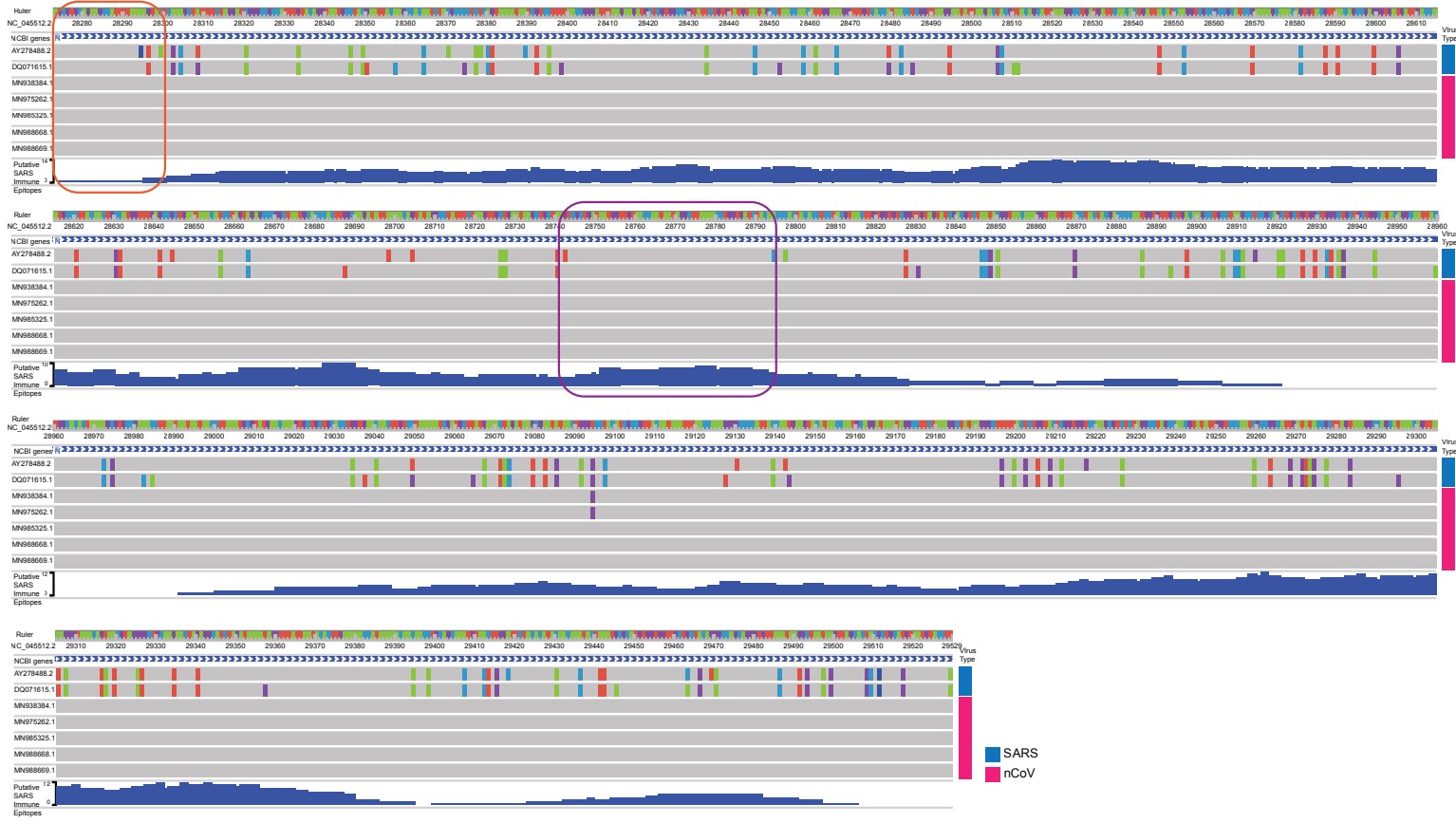


## C

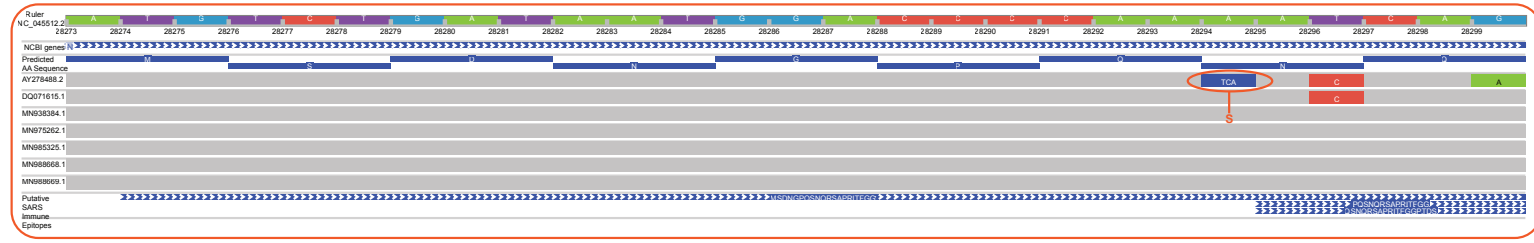


# Figure 3

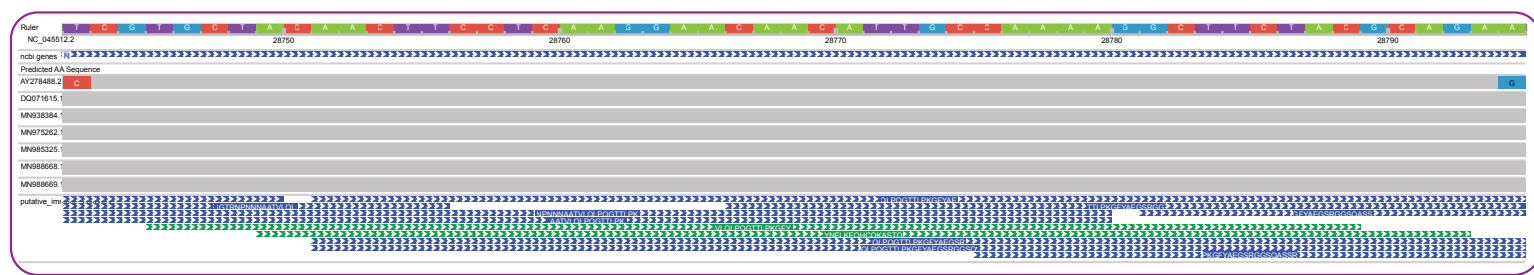
## A



## B



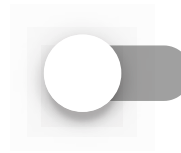
## C



# Figure 4

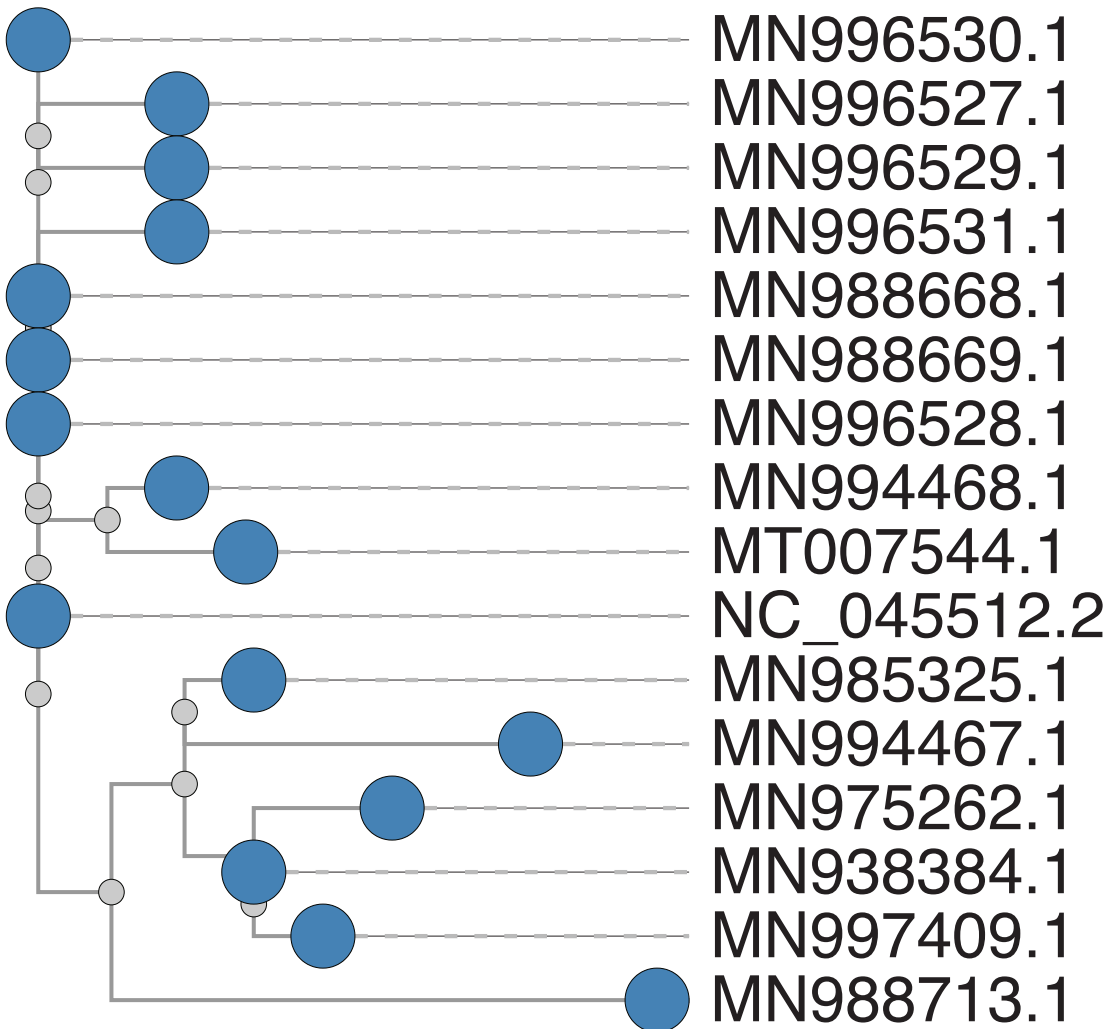


Linear

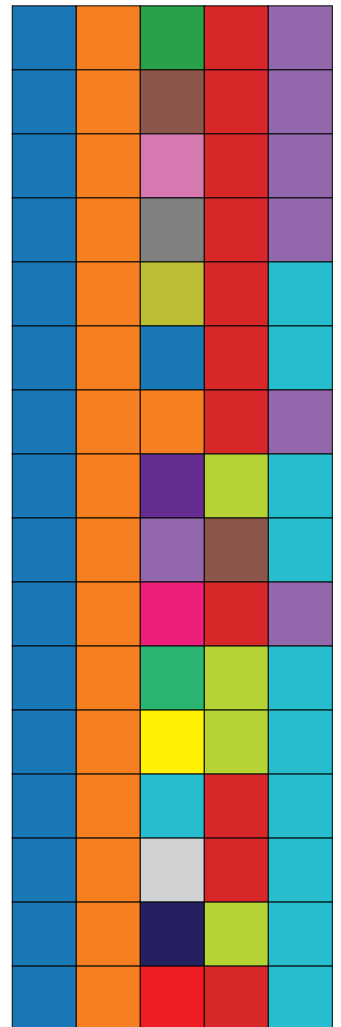


Radial

0.00020

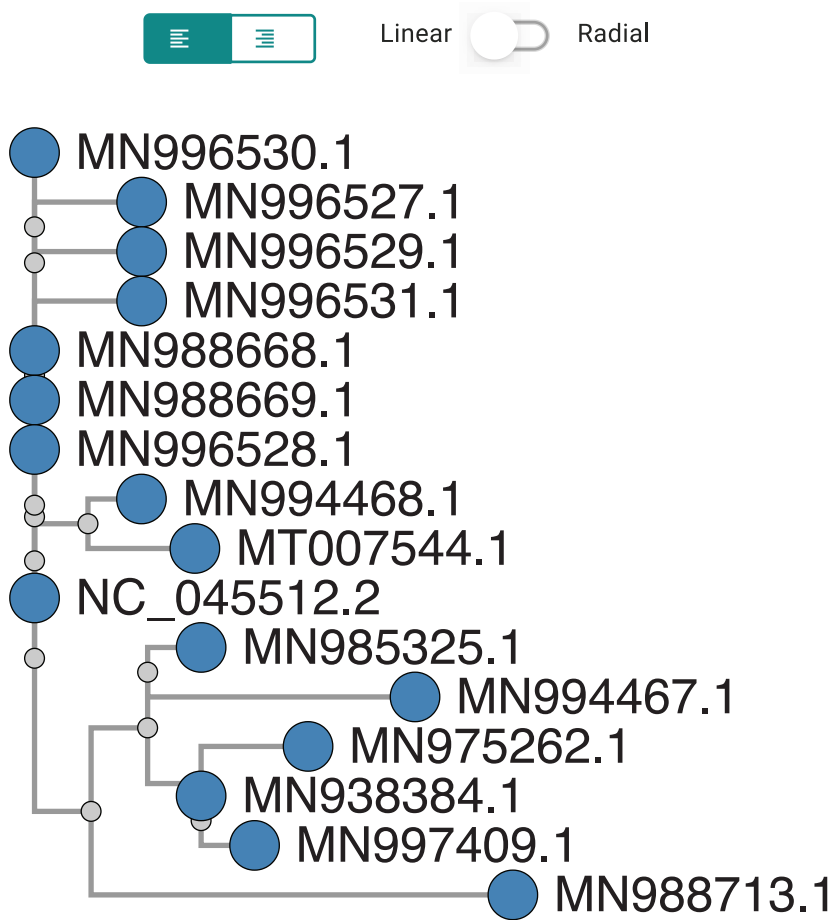


Virus Species  
 Molecule Type  
 Isolate  
 Country  
 Year



# Figure 5

## A



## B

