

Course Manual - Phylogenetic Analysis

Scenario 3 – Retrieval of sequences and phylogenetic analysis

From the metagenomic analysis carried out you have determined the pathogen responsible for the outbreak above and now wish to identify a likely source of infection. Using phylogenetic analysis, what is the likely source of the infection?

Note: All analysis in this section will be carried out on the viral spike protein sequence – NCBI Protein id - **QHD43416.1** from MN908947.3 (Wuhan-1 strain)

Software for the session:

1. **Mafft** - Alignment tool
2. **MEGA** - Alignment Viewer and Editor
3. **Modeltest-ng** - Model Testing
4. **IQ-TREE** - Tree Building tool
5. **Figtree/MEGA** - Tree Viewer and Editor

Step 1 – Downloading related sequences

Website – GenBank (<https://www.ncbi.nlm.nih.gov/genbank/>)

An official website of the United States government [Here's how you know](#)

NIH National Library of Medicine
National Center for Biotechnology Information

sunandoroyud

GenBank Nucleotide MN908947 Search

GenBank Submit Genomes WGS Metagenomes TPA TSA INSDC Other

GenBank Overview

What is GenBank?

GenBank[®] is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences (*Nucleic Acids Research*, 2011 Jan 41(01):D38-42). GenBank is part of the [International Nucleotide Sequence Database Collaboration](#), which comprises the DNA Data Bank of Japan (DDBJ), the European Nucleotide Archive (ENA), and GenBank at NCBI. These three organizations exchange data on a daily basis.

A GenBank release occurs every two months and is available from the [ftp site](#). The [release notes](#) for the current version of GenBank provide detailed information about the release and notifications of upcoming changes to GenBank. Release notes for [previous GenBank releases](#) are also available. GenBank growth [statistics](#) for both the traditional GenBank divisions and the WGS division are available from each release.

An [annotated sample GenBank record](#) for a *Saccharomyces cerevisiae* gene demonstrates many of the features of the GenBank flat file format.

Access to GenBank

There are several ways to search and retrieve data from GenBank:

- Search GenBank for sequence identifiers and annotations with [Entrez Nucleotide](#).
- Search and align GenBank sequences to a query sequence using [BLAST](#) (Basic Local Alignment Search Tool). See [BLAST info](#) for more information about the numerous BLAST databases.
- Search, link, and download sequences programmatically using [NCBI utilities](#).
- The ASN.1 and flatfile formats are available at NCBI's anonymous FTP server: [ftp://ftp.ncbi.nlm.nih.gov/ncbi-asn1](#) and [ftp://ftp.ncbi.nlm.nih.gov/genbank](#).

GenBank Data Usage

The GenBank database is designed to provide and encourage access within the scientific community to the most up-to-date and comprehensive [PMI sequence information](#). Therefore, NCBI places no restriction on the reuse or distribution of the GenBank data. However,

GenBank Resources

- [GenBank Home](#)
- [Submission Types](#)
- [Submission Tools](#)
- [Search GenBank](#)
- [Update GenBank Records](#)

An official website of the United States government [Here's how you know](#)

NIH National Library of Medicine
National Center for Biotechnology Information

sunandoroyud

Nucleotide Nucleotide Search

Advanced Help

GenBank Send

Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1, complete genome

GenBank: MN908947.3
[FASTA](#) [Graphics](#)

Go to: [Annotations](#)

LOCUS MN908947 29983 bp ss-RNA linear VRL 18-MAR-2020
DEFINITION Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1, complete genome.
ACCESSION MN908947
VERSION MN908947.3
KEYWORDS
SOURCE Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)
ORGANISM *Severe acute respiratory syndrome coronavirus 2*
Viruses; Riboviria; Orthornavirae; Pisuviricota; Pisoniviricetes; Nidovirales; Coronavirinae; Coronaviridae; Orthocoronavirinae; Betacoronavirus; Sarbecovirus.
REFERENCE 1 (bases 1 to 29983)
AUTHORS Wu, F., Zhao, S., Yu, B., Chen, Y.M., Wang, W., Song, Z.G., Hu, Y., Tao, Z.W., Tian, Z.H., Pei, Y.Y., Yuan, M.L., Zhang, Y.L., Dai, F.H., Liu, Y., Wang, Q.M., Zheng, J.J., Xu, L., Holmes, E.C. and Zhang, Y.Z.
TITLE A new coronavirus associated with human respiratory disease in China
JOURNAL Nature 579 (7798), 265-269 (2020)
PUBMED 32015588
REFERENCE 2 (bases 1 to 29983)

Change region shown

Customize view

Analyze this sequence

- Run BLAST
- Pick Primers
- Highlight Sequence Features
- Find in this Sequence

NCBI Virus

Retrieve, view, and download SARS-CoV-2 coronavirus genomic and protein sequences.

Related information

- Assembly
- PubMed
- Taxonomy
- Full text in PMC
- Gene


gene

CDS

21563..25384
/gene="S"
21563..25384
/gene="S"
/note="structural protein"
/cdon_start=1
/product="surface glycoprotein"
/protein_id="ORF3a.1"

/translation="MEVFLVLLPLVSSQCVLLTTRTQLPPATNGFTRGVVPOKVR
SSVLSHQFLPFFSMTMPLHMSGTHGTRFDPVLPMDGVYFSTESKILR
QKGTGTLDSGQSLVNNATVYVCEQFENRPLGVVYNNMSSEFRTV
SSANKTFEYVSQFLNDLGGKQKQKRLRFVFNIDGVYKYSKHTPILNVDLPQ
GFSALPLVDLPILNITRFQTLALHRYLTPGSSSGVAGAAAYVGVLPRTFL
LKNMGITDAVDCALDPLSEKCTLKSFTVEKIVQTSNFRVQTESIVRFPILTH
LCPPGEVNAIRFASVYANRIRISNVADSVLYNSASFSTKCYGVSPKTLNDLCE
TNPVADSVIRGDEWQIAPQKTAIDVWYLPDFGCTVANSRHLDSVGNH
VLYRLFRSHLKPFRDSTETVQAGSTPCMGVEGFNCFPLQSYGFQPTNGVQPY
RVVLSFELLHAPATVCGPKSTNLVQKCNVFNGLTGTVLTESNKKLPPQFG
RDIADITDAVDPQTLLELIDTPCSGGVSVTPGTNTSNQAVLYQVNCIEVPVAI
HAQQLTPTRVYSTGSNFTVTRAGCLGAENHNSYECIDPEGALCASVQTNISPR
BARVASQSIATYHSLGAENSVSNISIAIPNFTISVTEILPVSHKTSVQDTH
YICGDSFCSHLLQVQSCFQLRHALTGLAEGENTQEVAGVQVYTPPTKQFG
GFNFSQLPDPSPKSRSEIIDLFNKVTADAGIKQVGDCLGTAANDLCAQKFN
GLTVLPPLTDEMLAQYTSALLAGTTSQNTFGAGAAQLPPANPQYFNGIVTQH
VLVYKQLTAKFISATGKIQSLSSITASALGKLDVNVNQAALNTLVKLSNFGA
TSSVLDLTLRLQKVEAFVQIDRLTGRQLSLQTYVQQLIRAEIRASANLAATHNS
EVLQGSRRDFCQKQYHNSPPGSPHGVYLVHTVYAGKRNFTAPATQDQKH
FPREGVVSAGTHAVTQRPVYEPQLETTONTIVSGKZNVGIVNTVYDPLQFELD
STKEELQKVFQHTSPVDLGGDSGNASVNIQKEIQRNEVKNLNSLDELQELG
KYEQVYKWPVYVNLGFAGLTAVVMTMLCONTSCCSCGKCCSGCKCFDEBDE
PVKGVLLHYT"

25393..26220
/gene="ORF3a"
25393..26220
/gene="ORF3a"


National Library of Medicine
 National Center for Biotechnology Information

[An official website of the United States government](#)
[Here's how you know](#)

[sunandoroyud](#)

BLAST® » blastn suite

[Home](#)
[Recent Results](#)
[Saved Strategies](#)
[Help](#)

blastn | blastp | blastx | tblastn | tblastx

Standard Nucleotide BLAST

BLASTN programs search nucleotide databases using a nucleotide query, more...

[Reset page](#)
[Bookmark](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

Query subrange [?](#)

From

To

Or, upload file [?](#)

Browse... No file selected.

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database

☒ Standard databases (nr etc.):
 ☐ rRNA/ITS databases
 ☐ Genomic + transcript databases
 ☐ Betacoronavirus

[?](#)

Organism

[exclude](#) [Add organism](#)

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown [?](#)

☐ Models (XM/XP) ☐ Uncultured/environmental sample sequences

☐ Sequences from type material

[YouTube](#) Create custom database

An official website of the United States government [Here's how you know](#)

NIH National Library of Medicine
National Center for Biotechnology Information

sunandoroyucl

BLAST® » blastn suite

Home Recent Results Saved Strategies Help

blastn blastp blastx tblastn tblastx

Standard Nucleotide BLAST

BLASTN programs search nucleotide databases using a nucleotide query. more...

Reset page Bookmark

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) Clear

MN908947.3

Query subrange [?](#)

From 21563 To 25384

Or, upload file [Browse...](#) No file selected. [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database ☒ Standard databases (nr etc.) ☐ rRNA/ITS databases ☐ Genomic + transcript databases ☐ Betacoronavirus [?](#)

Nucleotide collection (nr/nt) [?](#)

Organism [Optional](#)

Sars-Cov-2 [?](#) exclude [Add organism](#)

SARS-CoV-2 (taxid:2697049)

Reverse genetics vector pCCL-4K-SARS-CoV-2-ZsGreen (taxid:279...

Reverse genetics vector pCCL-4K-SARS-CoV-2-Wuhan-Hu-1 (taxid...

Reverse genetics vector pCCL-4K-SARS-CoV-2-mCherry (taxid:276...

Reverse genetics vector mC23-4K-SARS-CoV-2-NanoLuc (taxid:276...

Exclude [Optional](#)

Limit to [Optional](#)

Entrez Query [Optional](#)

Enter an Entrez query to limit search [?](#)

Program Selection

Optimize for ☐ Highly similar sequences (megablast) ☐ More dissimilar sequences (discontiguous megablast) ☒ Somewhat similar sequences (blastn) [?](#)

Choose a BLAST algorithm [?](#)

BLAST Search database Nucleotide collection (nr/nt) using Blastn (Optimize for somewhat similar sequences)

☐ Show results in a new window

+ Algorithm parameters

FOLLOW NCBI

An official website of the United States government [Here's how you know](#)

NIH National Library of Medicine
National Center for Biotechnology Information

sunandoroyucl

BLAST® » blastn suite » results for RID-9Z82W1A401R

Home Recent Results Saved Strategies Help

[< Edit Search](#) [Save Search](#) [Search Summary](#) [?](#) How to read this report? [BLAST Help Videos](#) [Back to Traditional Results Page](#)

i Your search is limited to records that exclude: SARS-CoV-2 (taxid:2697049)

Job Title gbj|MN908947.3|

RID 9Z82W1A401R [Search expires on 06-09 03:00 am](#) [Download All](#) [?](#)

Program BLASTN [Citation](#) [?](#)

Database nt [See details](#) [?](#)

Query ID MN908947.3

Description Severe acute respiratory syndrome coronavirus 2 isolate [?](#)

Molecule type nucleic acid

Query Length 3822

Other reports [Distance tree of results](#) [MSA viewer](#) [?](#)

Filter Results

Organism [only top 20 will appear](#) ☐ exclude

Type common name, binomial, taxid or group name

[+ Add organism](#)

Percent Identity to E value to Query Coverage to

[Filter](#) [Reset](#)

Descriptions [Graphic Summary](#) [Alignments](#) [Taxonomy](#)

Sequences producing significant alignments

Download [Select columns](#) [Show](#) 100 [?](#)

☒ select all 100 sequences selected

[GenBank](#) [Graphics](#) [Distance tree of results](#) [MSA Viewer](#)

Max Total Query E Per Acc.

[Feedback](#)

Sequences producing significant alignments Download Select columns Show 100

select all 100 sequences selected

GenBank Graphics Distance tree of results MSA Viewer

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/> Synthetic construct clone iSARS-CoV-2-nLuc-GFP-ORF1ab polyprotein (ORF1ab)_ORF1a polyprotein (ORF1a)_...	synthetic construct	6893	6893	100%	0.0	100.00%	30857	MT46167.1
<input checked="" type="checkbox"/> Synthetic construct clone iSARS-CoV-2-GFP-ORF1ab polyprotein (ORF1ab)_ORF1a polyprotein (ORF1a)_...	synthetic construct	6893	6893	100%	0.0	100.00%	30347	MT461670.1
<input checked="" type="checkbox"/> Synthetic construct clone iSARS-CoV-2-WT-ORF1ab polyprotein (ORF1ab)_ORF1a polyprotein (ORF1a)_...	synthetic construct	6893	6893	100%	0.0	100.00%	29903	MT461669.1
<input checked="" type="checkbox"/> Synthetic construct ORF1ab, spike, ORF3, E, M, ORF6, ORF8, and N genes, complete cds	synthetic construct	6893	6893	100%	0.0	100.00%	29891	MT108784.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_03 genome assembly, complete genom...	Severe acute re...	6889	6889	100%	0.0	99.97%	29903	HG994854.1
<input checked="" type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-ZsGreen, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	36033	MW289908.1
<input checked="" type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-NanoLuc, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	35853	MT926412.1
<input checked="" type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-mCherry, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	36048	MT926411.1
<input checked="" type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-Wuhan-Hu-1, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	35283	MT926410.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_07 genome assembly, complete genom...	Severe acute re...	6884	6884	100%	0.0	99.95%	29903	HG994856.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_10 genome assembly, complete genom...	Severe acute re...	6880	6880	100%	0.0	99.92%	29903	HG994859.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_08 genome assembly, complete genom...	Severe acute re...	6880	6880	100%	0.0	99.92%	29903	HG994857.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_09 genome assembly, complete genom...	Severe acute re...	6877	6877	100%	0.0	99.90%	29903	HG994858.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_06 genome assembly, complete genom...	Severe acute re...	6877	6877	100%	0.0	99.90%	29903	HG994855.1
<input checked="" type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_01 genome assembly, complete genom...	Severe acute re...	6866	6866	100%	0.0	99.84%	29900	HG994852.1
<input checked="" type="checkbox"/> Synthetic construct clone IDCCD-RG5, complete sequence	synthetic construct	6839	6839	100%	0.0	99.69%	30189	MZ128150.1
<input checked="" type="checkbox"/> Synthetic construct clone IDCCD-RG4, complete sequence	synthetic construct	6816	6816	100%	0.0	99.53%	30179	MZ128149.1
<input checked="" type="checkbox"/> Synthetic construct clone IDCCD-RG3, complete sequence	synthetic construct	6816	6816	100%	0.0	99.56%	30205	MZ128148.1
<input checked="" type="checkbox"/> Cloning vector p20020, complete sequence	Cloning vector p...	6796	6796	98%	0.0	99.84%	9004	MW045215.1
<input checked="" type="checkbox"/> Cloning vector p20020-BANCOVID-SARS-CoV-2 surface glycoprotein gene, complete cds	Cloning vector p...	6796	6796	98%	0.0	99.84%	4326	MW045214.1
<input type="checkbox"/> Synthetic construct clone IDCCD-RG3, complete sequence	synthetic construct	6816	6816	100%	0.0	99.56%	30205	MZ128148.1
<input type="checkbox"/> Cloning vector p20020, complete sequence	Cloning vector p...	6796	6796	98%	0.0	99.84%	9004	MW045215.1
<input type="checkbox"/> Cloning vector p20020-BANCOVID-SARS-CoV-2 surface glycoprotein gene, complete cds	Cloning vector p...	6796	6796	98%	0.0	99.84%	4326	MW045214.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_02 genome assembly, complete genom...	Severe acute re...	6861	6861	96%	0.0	99.92%	29901	HG994853.1
<input type="checkbox"/> Synthetic construct clone CSV41	synthetic construct	6587	6587	100%	0.0	98.22%	171907	MW036243.1
<input type="checkbox"/> Synthetic construct clone C46/53	synthetic construct	6582	6582	100%	0.0	98.19%	171918	MW036460.1
<input checked="" type="checkbox"/> Bat coronavirus isolate BANAL-20-52/Laos/2020, complete genome	Bat coronavirus	5971	5971	100%	0.0	94.81%	29638	MZ937006.1
<input checked="" type="checkbox"/> Bat coronavirus RaTG13, complete genome	Bat coronavirus	5671	5671	100%	0.0	92.89%	29655	MN096532.2
<input type="checkbox"/> Cloning vector pSF_Lentiv_SARS-CoV-2_partial-5'E/M/N, complete sequence	Cloning vector p...	5512	5512	80%	0.0	99.87%	13543	MT298905.1
<input type="checkbox"/> Synthetic construct clone NIBSC_20_138_4 sequence	synthetic construct	5512	5512	80%	0.0	99.87%	7558	MW059035.1
<input type="checkbox"/> Synthetic construct COM-VAC, complete sequence	synthetic construct	5470	5470	100%	0.0	91.58%	29655	MZ404503.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_11 genome assembly, complete genom...	Severe acute re...	5117	5492	94%	0.0	99.03%	29903	HG994860.1
<input checked="" type="checkbox"/> Bat coronavirus isolate BANAL-20-23/Laos/2020, complete genome	Bat coronavirus	4793	4793	98%	0.0	88.20%	29644	MZ937003.2
<input checked="" type="checkbox"/> Bat coronavirus isolate BANAL-20-103/Laos/2020, complete genome	Bat coronavirus	4763	4763	98%	0.0	88.01%	29632	MZ937001.1
<input checked="" type="checkbox"/> Pangolin coronavirus isolate cDNA31-S surface glycoprotein (S) gene, complete cds	Pangolin corona...	4159	4159	99%	0.0	84.25%	3798	MT799526.1
<input checked="" type="checkbox"/> Pangolin coronavirus isolate cDNA20-S surface glycoprotein (S) gene, complete cds	Pangolin corona...	4159	4159	99%	0.0	84.25%	3798	MT799525.1
<input checked="" type="checkbox"/> Pangolin coronavirus isolate cDNA18-S surface glycoprotein (S) gene, complete cds	Pangolin corona...	4159	4159	99%	0.0	84.25%	3798	MT799524.1
<input type="checkbox"/> Pangolin coronavirus isolate cDNA16-S surface glycoprotein (S) gene, complete cds	Pangolin corona...	4159	4159	99%	0.0	84.25%	3798	MT799523.1
<input type="checkbox"/> Pangolin coronavirus isolate cDNA8-S surface glycoprotein (S) gene, complete cds	Pangolin corona...	4159	4159	99%	0.0	84.25%	3798	MT799521.1
<input type="checkbox"/> Pangolin coronavirus isolate cDNA9-S surface glycoprotein (S) gene, complete cds	Pangolin corona...	4155	4155	99%	0.0	84.22%	3798	MT799522.1
<input type="checkbox"/> Pangolin coronavirus isolate MP789, complete genome	Pangolin corona...	4150	4150	99%	0.0	84.20%	29521	MT121216.1
<input type="checkbox"/> Pangolin coronavirus isolate PCoV_GX-P4L, complete genome	Pangolin corona...	4039	4039	100%	0.0	83.55%	29805	MT040333.1
<input type="checkbox"/> Pangolin coronavirus isolate PCoV_GX-P5L, complete genome	Pangolin corona...	4035	4035	100%	0.0	83.52%	29806	MT040335.1
<input type="checkbox"/> Pangolin coronavirus isolate PCoV_GX-P9E, complete genome	Pangolin corona...	4030	4030	100%	0.0	83.50%	29802	MT040336.1
<input type="checkbox"/> Pangolin coronavirus isolate CX_P2V, complete genome	Pangolin corona...	4028	4028	100%	0.0	83.48%	29729	MW032698.1

Other reports Distance tree of results MSA viewer

Descriptions Graphic Summary Alignments Taxonomy

Sequences producing significant alignments Download Select columns Show 100

FASTA (complete sequence) FASTA (aligned sequences) GenBank (complete sequence) Hit Table (text) Hit Table (CSV) Text Descriptions Table (CSV) XML ASN.1

Distance tree of results MSA Viewer

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input type="checkbox"/> Synthetic construct clone iSARS-CoV-2-nLuc-GFP-ORF1ab polyprotein (ORF1ab)_ORF1a polyprotein (ORF1a)_...	synthetic construct	6893	6893	100%	0.0	100.00%	30857	MT46167.1
<input type="checkbox"/> Synthetic construct clone iSARS-CoV-2-GFP-ORF1ab polyprotein (ORF1ab)_ORF1a polyprotein (ORF1a)_...	synthetic construct	6893	6893	100%	0.0	100.00%	30347	MT461670.1
<input type="checkbox"/> Synthetic construct clone iSARS-CoV-2-WT-ORF1ab polyprotein (ORF1ab)_ORF1a polyprotein (ORF1a)_...	synthetic construct	6893	6893	100%	0.0	100.00%	29903	MT461669.1
<input type="checkbox"/> Synthetic construct ORF1ab, spike, ORF3, E, M, ORF6, ORF8, and N genes, complete cds	synthetic construct	6893	6893	100%	0.0	100.00%	29891	MT108784.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_03 genome assembly, complete genom...	Severe acute re...	6889	6889	100%	0.0	99.97%	29903	HG994854.1
<input type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-ZsGreen, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	36033	MW289908.1
<input type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-NanoLuc, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	35853	MT926412.1
<input type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-mCherry, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	36048	MT926411.1
<input type="checkbox"/> Reverse genetics vector pCCI-4K-SARS-CoV-2-Wuhan-Hu-1, complete sequence	Reverse genetic...	6889	6889	100%	0.0	99.97%	35283	MT926410.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_07 genome assembly, complete genom...	Severe acute re...	6884	6884	100%	0.0	99.95%	29903	HG994856.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_10 genome assembly, complete genom...	Severe acute re...	6880	6880	100%	0.0	99.92%	29903	HG994859.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_08 genome assembly, complete genom...	Severe acute re...	6880	6880	100%	0.0	99.92%	29903	HG994857.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_09 genome assembly, complete genom...	Severe acute re...	6877	6877	100%	0.0	99.90%	29903	HG994858.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_06 genome assembly, complete genom...	Severe acute re...	6877	6877	100%	0.0	99.90%	29903	HG994855.1
<input type="checkbox"/> Severe acute respiratory syndrome-related coronavirus isolate H_SC_01 genome assembly, complete genom...	Severe acute re...	6866	6866	100%	0.0	99.84%	29900	HG994852.1
<input type="checkbox"/> Synthetic construct clone IDCCD-RG5, complete sequence	synthetic construct	6839	6839	100%	0.0	99.69%	30189	MZ128150.1
<input type="checkbox"/> Synthetic construct clone IDCCD-RG4, complete sequence	synthetic construct	6816	6816	100%	0.0	99.53%	30179	MZ128149.1

Feedback

Alternate Ways to Download from GenBank

Batch Entrez (<https://www.ncbi.nlm.nih.gov/sites/batchentrez>)

Batch Entrez

Given a file of Entrez accession numbers or other identifiers, Batch Entrez downloads the corresponding records.

Instructions

1. Start with a local file containing a list of accession numbers or identifiers
2. Select the database corresponding to the type of accession numbers or identifiers in your input file
3. Use the **Browse** or **Choose File...** button to select the input file
4. Press the **Retrieve** button to see a list of document summaries
5. Select a format in which to display the data for viewing, and/or saving
6. Select "Send to file" to save the file.

Tips

- To download entire genome records, check the NCBI FTP site, instead of using Batch Entrez.
- Some lists of record identifiers can be tens of thousands of lines long, so Batch Entrez may not retrieve all records from one list. Split the list of identifiers into smaller files using a file splitting software or a file split command at the command prompt in UNIX or LINUX systems.
- When loading large numbers of genome records, put several thousand record identifiers per file, one per line, left-adjusted.
- Please note that Batch Entrez will check for duplicate identifiers when reporting results from a list that you have imported.
- When retrieving a list of Nucleotide accessions, you must select the specific component database from which the accessions or GIs were saved. For Nucleotide, choose either the CoreNucleotide, the EST or the GSS selection from the database menu. If you have a mixed list of nucleotide accessions or UIDs, you will need to run the Batch Entrez search three times. (Select the database from the pull down menu)

Note: If you have Accession from different databases you will have to run Batch Entrez multiple times each for a unique sequence database

File to Use ~/Sunando/Betacoronavirus_Accession.txt

Entrez E-utilities <ftp://ftp.ncbi.nlm.nih.gov/entrez/entrezdirect/>

Manual - <https://www.ncbi.nlm.nih.gov/books/NBK179288/>

QuickStart - <http://bioinformatics.cvr.ac.uk/blog/ncbi-entrez-direct-unix-e-utilities/>

Browser

```
https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?
db=nucleotide&id=AY278488, AY304486, MN908947,
MT782115&rettype=fasta&retmode=text
```

Command Line

```
$ esearch -db "protein" -query "txid11270[Organism] AND L Protein Complete AND
refseq[filter]" | efetch -format fasta > outputfile.fasta

$ head outputfile.fasta
```

For the next step we will start with ~/Sunando/Spike.fas

Step 2 – Aligning Sequences.

Software Used

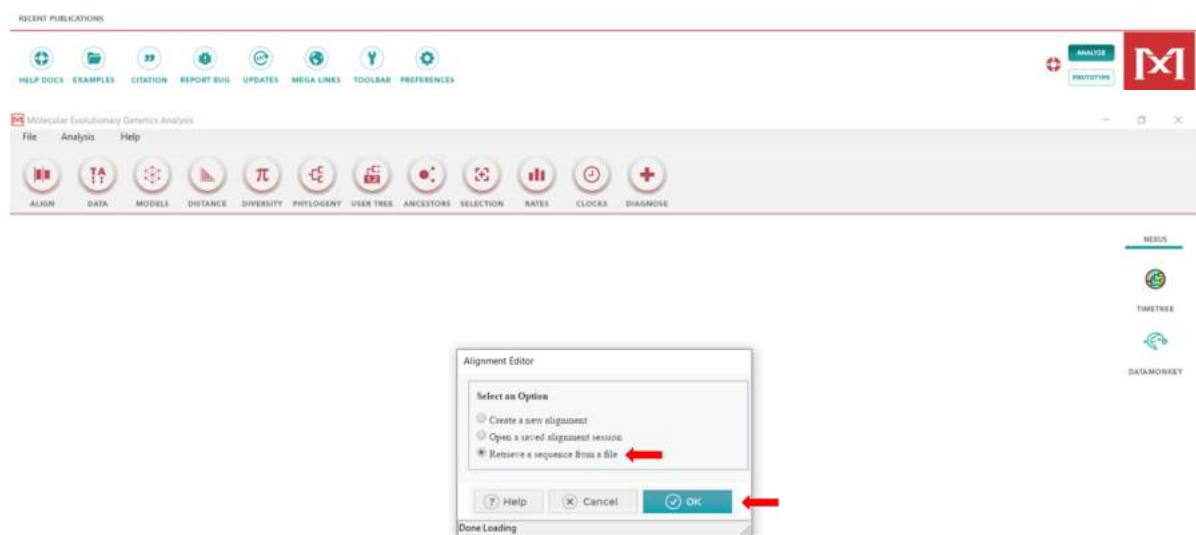
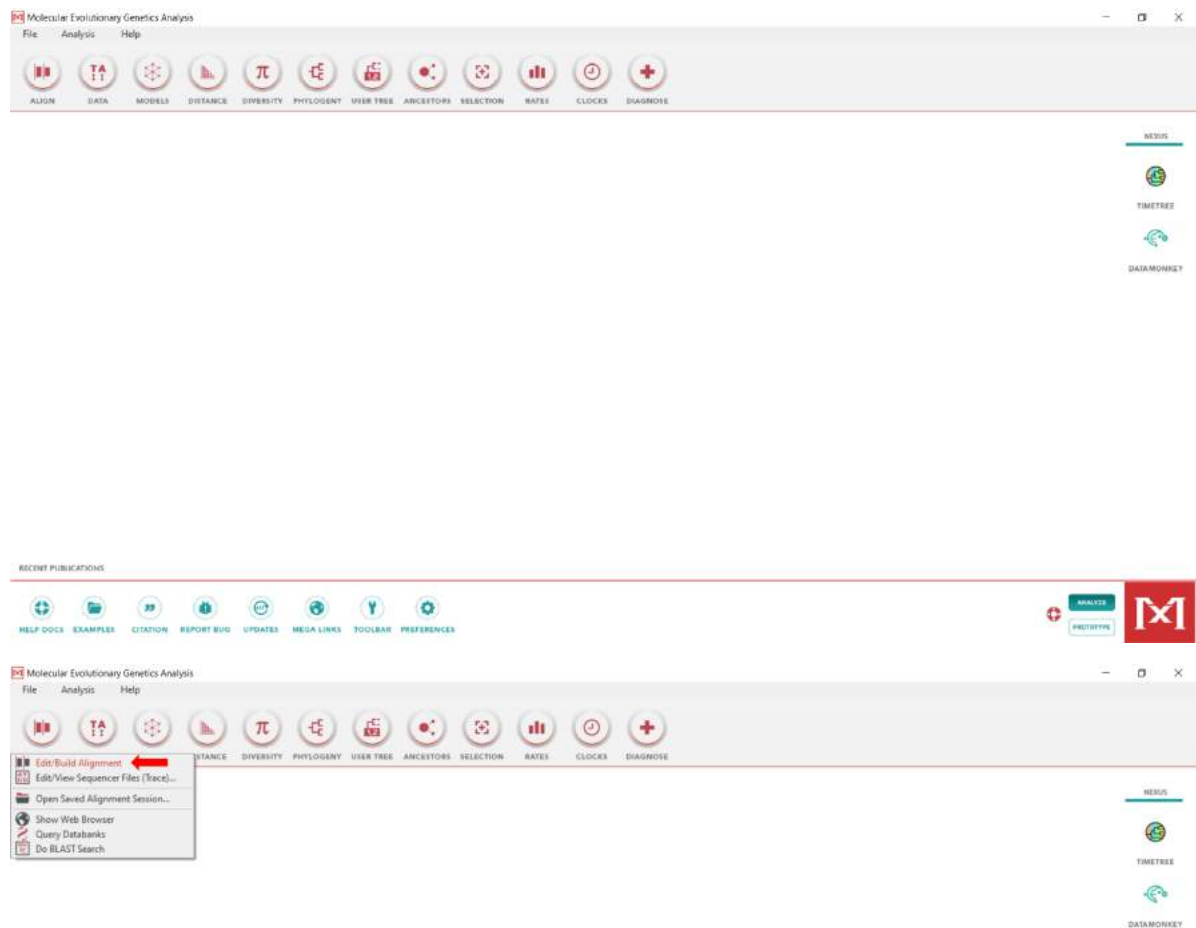
Mafft (<https://mafft.cbrc.jp/alignment/software/>)

MEGA (<http://droua.prabi.fr/software/seaview>)

Alternate Software – MUSCLE, CLUSTALW

To view/edit sequence files we will use MEGA. MEGA has a GUI and will launch as a standalone program

```
$ mega
```



Picture40


```
$ mafft ~/Sunando/Spike.fas > ~/Sunando/outputfile.fas
```

Uses the default models to align the sequences. For highly divergent sequences this may produce inaccurate alignments.

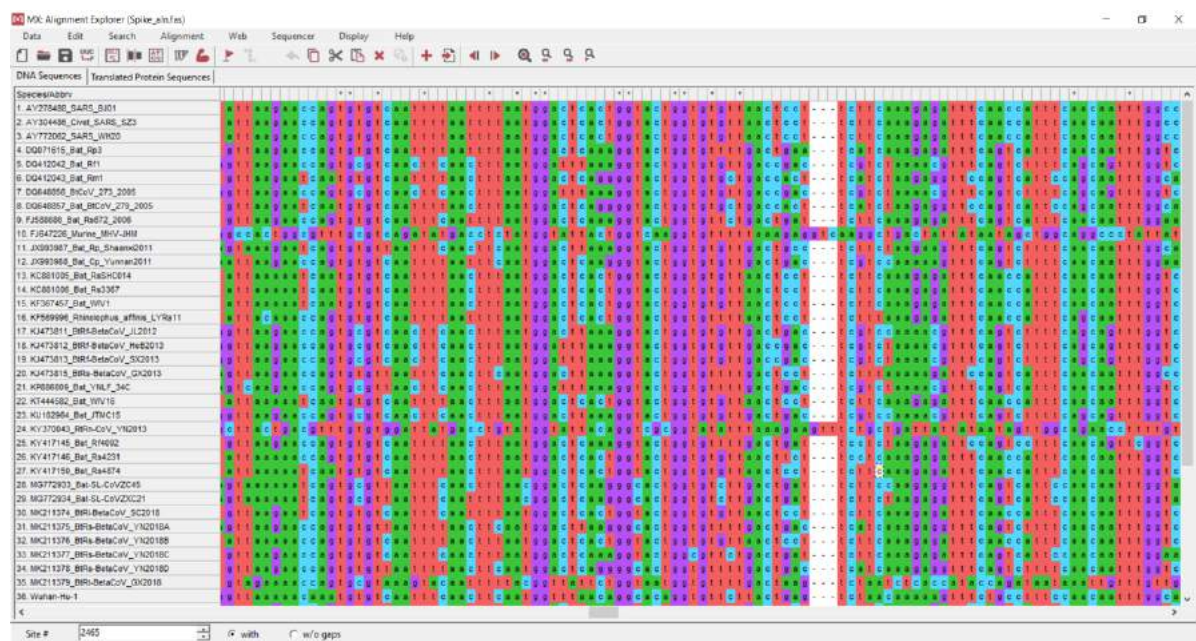
Alternatives are if you have a curated alignment `mafft --add` works to add new sequences to existing alignments which puts more weightage on the existing alignment

We can also use the L-INS-i algorithm in Mafft that aligns more divergent sequences using pairwise local alignments

```
$ mafft --maxiterate 1000 --localpair ~/Sunando/Spike.fas > outputfileinsi.fas
```

The final alignment file is ~/Sunando/Spike_aln.fas

You can also use the alignment file you have generated.



Step 3 – Constructing a Phylogeny.

Software Used

Modeltest-ng (<https://github.com/ddarriba/modeltest>)

IQ-TREE (<http://www.iqtree.org/>)

Alternate Software – PhyML, RAxML

For this session we will start with the aligned nucleotide sequences created in the last step.

Model Testing

To run model testing we will use Modeltest-ng

```
$ modeltest-ng -d aa -i ~/Sunando/Spike_aln.fas -o modeloutputfile -t ml -p 2
```

-i : Input file

-o : Output file

-t : Sets the starting tree topology

-p : Number of threads

We test three criteria to select the best fitting models BIC, AIC and AICc. The modeltest results for this alignment are in ~/Sunando/Spike_model.out and ~/Sunando/Spike_model.log

```
$ nano ~/Sunando/Spike_model.out
$ nano ~/Sunando/Spike_model.log
```

AICc	model	K	lnL	score	delta	weight
1	TVM+I+G4	9	-89441.2343	179094.4687	0.0000	0.5632
2	GTR+I+G4	10	-89440.4884	179094.9768	0.5081	0.4368
3	TVM+G4	8	-89467.5351	179145.0703	50.6016	0.0000
4	GTR+G4	9	-89467.0653	179146.1306	51.6620	0.0000
5	TPM3uf+I+G4	7	-89514.2112	179236.4224	141.9537	0.0000
6	TIM3+I+G4	8	-89514.1910	179238.3821	143.9134	0.0000
7	TPM2uf+I+G4	7	-89530.2483	179268.4965	174.0279	0.0000
8	TIM2+I+G4	8	-89530.0605	179270.1211	175.6524	0.0000
9	TPM3uf+G4	6	-89541.5959	179289.1919	194.7232	0.0000
10	TIM3+G4	7	-89541.6048	179291.2095	196.7409	0.0000

Best model according to AICc

Model: TVM+I+G4
lnL: -89441.2343
Frequencies: 0.2626 0.1931 0.1752 0.3691
Subst. Rates: 2.6387 5.3471 1.7471 1.9174 5.3471 1.0000
Inv. sites prop: 0.0673
Gamma shape: 0.9500
Score: 179094.4687
Weight: 0.5632

The model that we will use for tree building is TVM+I+G4

Tree building

To build a tree we are going to use IQ-TREE

```
$ iqtree -s ~/Sunando/Spike_aln.fas -bb 1000 -st DNA -nt 4 -alrt 1000 -pre  
treeoutfile
```

-s : Input File

-bb : ultrafast bootstrap

-st : data type

-nt : Number of threads

-alrt : SH-like approximate likelihood ratio test

-pre : Prefix for output file

IQ-TREE outputs multiple files. The final tree file we will use has an extension of .contree

The final output we will take forward to the next step while IQ-TREE completes running will be ~/Sunando/Spike_Tree.contree

Step 4 – Viewing and Modifying a Tree File.

Software Used

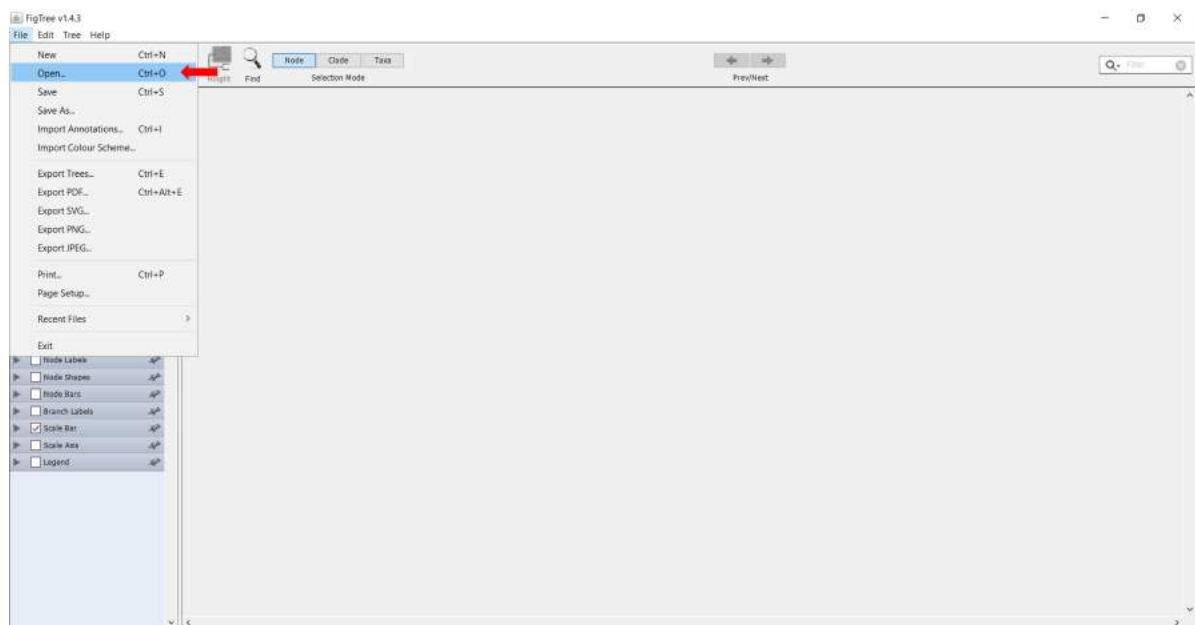
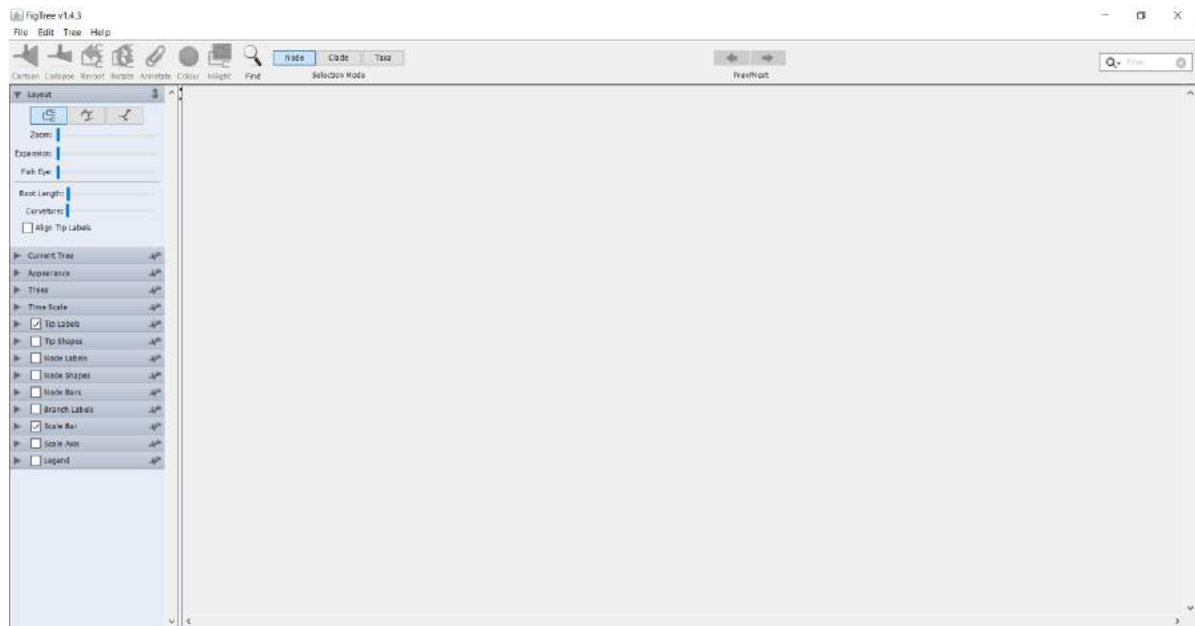
FigTree (<http://tree.bio.ed.ac.uk/software/figtree/>)

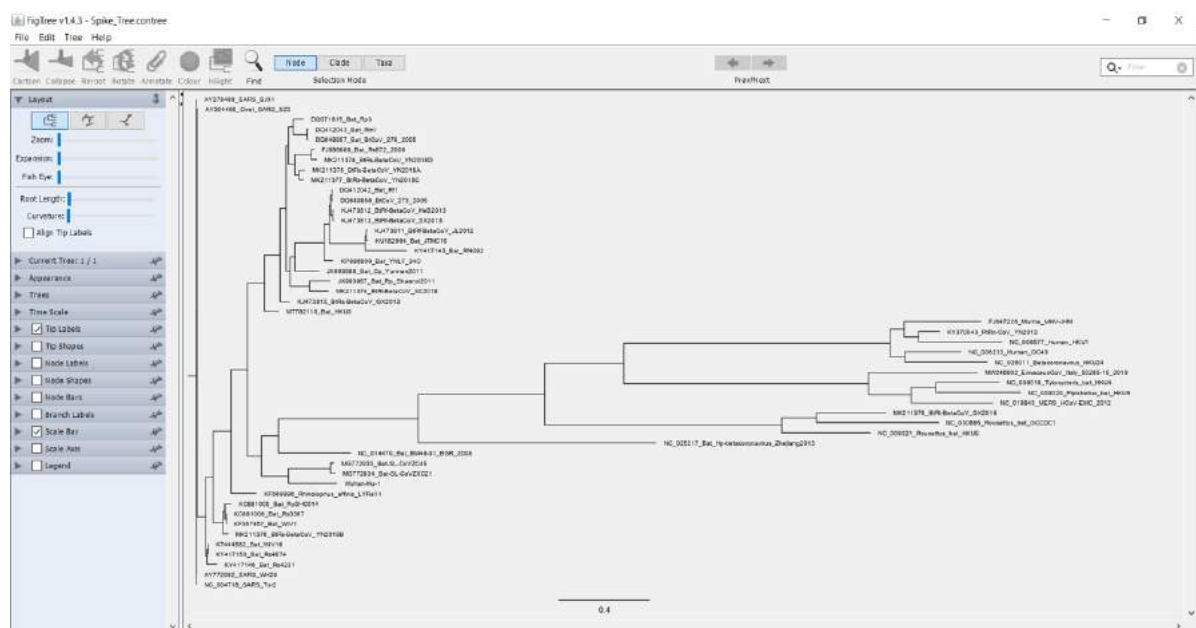
MEGA (<https://www.megasoftware.net/>)

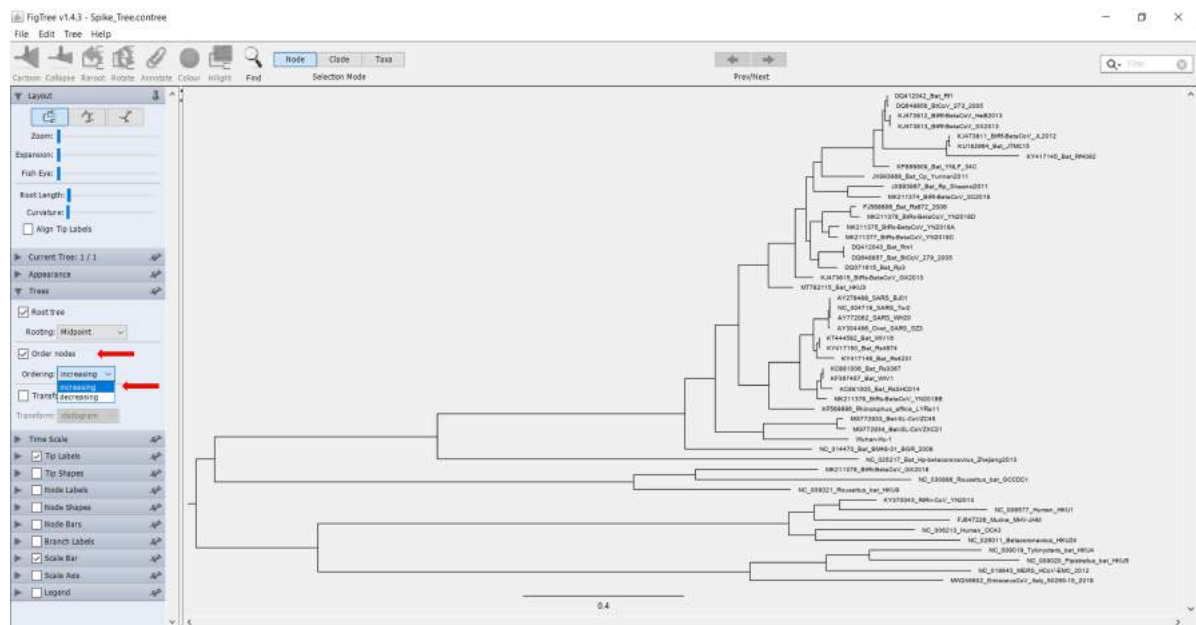
\$ figtree

Figtree is Java based and will launch a GUI

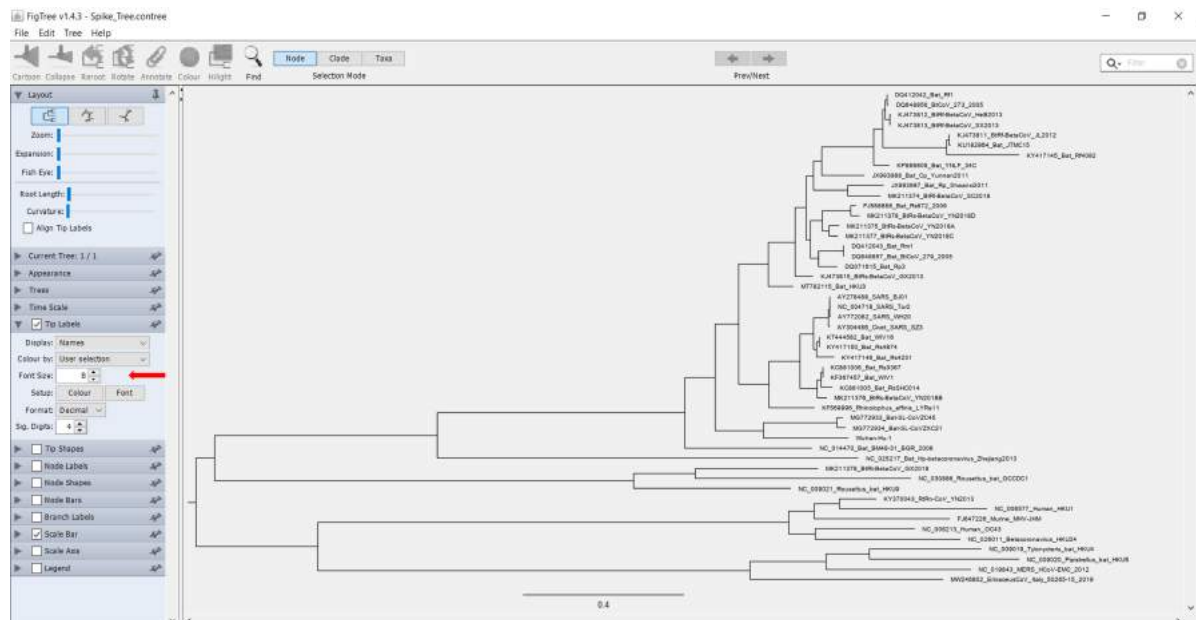
Open file - ~/Sunando/Spike_Tree.contree



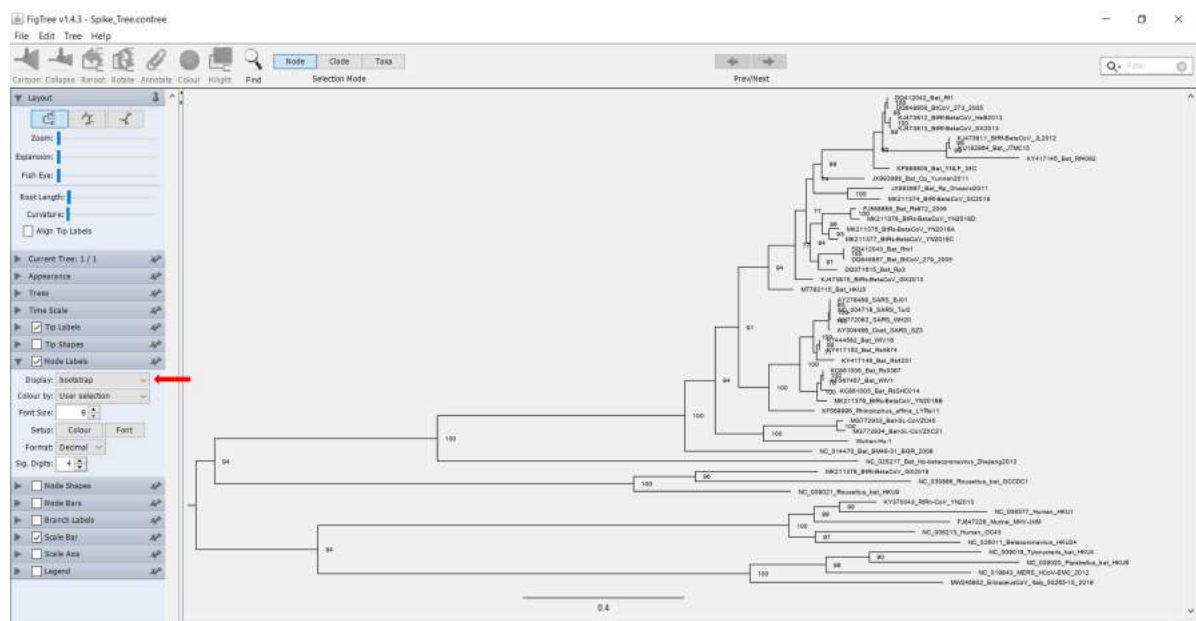
[illegible]

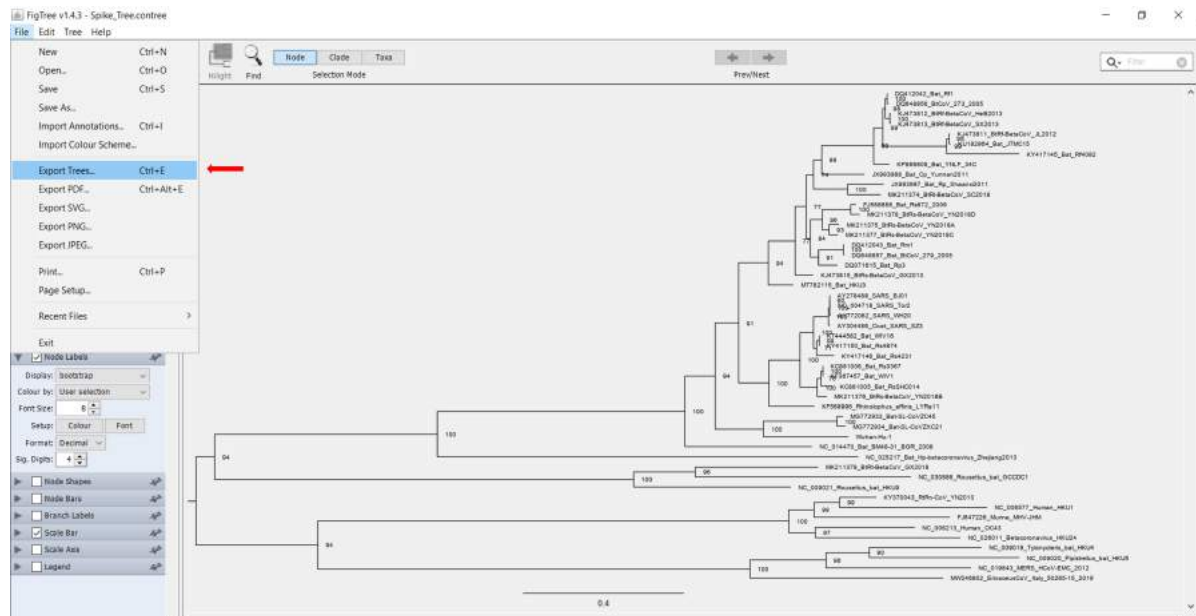


Modify Labels



Add bootstrap support

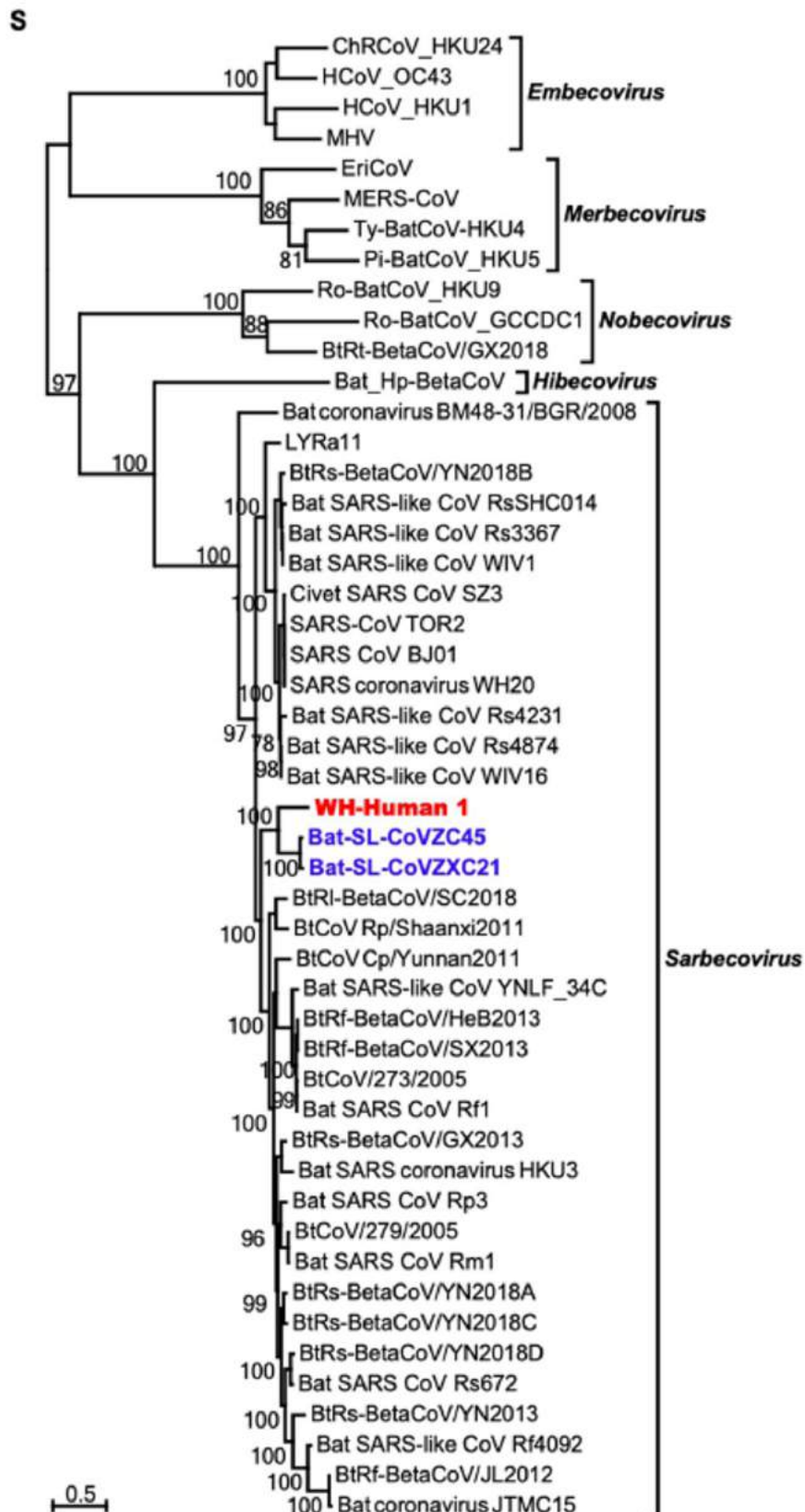




We can go through similar tree viewing in MEGA

Final Results

We further process the trees in **PowerPoint/ggtree** or other image editing tools to add metadata for each sequence in the alignment. In this case we have added virus groups for each virus in the alignment.



Based on this data one would infer that the isolated virus most likely lies within the Sarbecovirus group which includes the original SARS coronavirus but is quite distinct from it. The closest related species both are from bats which could suggest a potential origin.

