

**March 1st, 2022**

# **Workshop 1**

**Downloading SARS-COV2 genomic sequences from NCB**

**Natya Hans March 1st, 2022**





## Click here

## **COVID-19 Information**

[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#) | [Español](#)

**BLAST®**

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

# Basic Local Alignment Search Tool

**BLAST** finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

[Learn more](#)

11

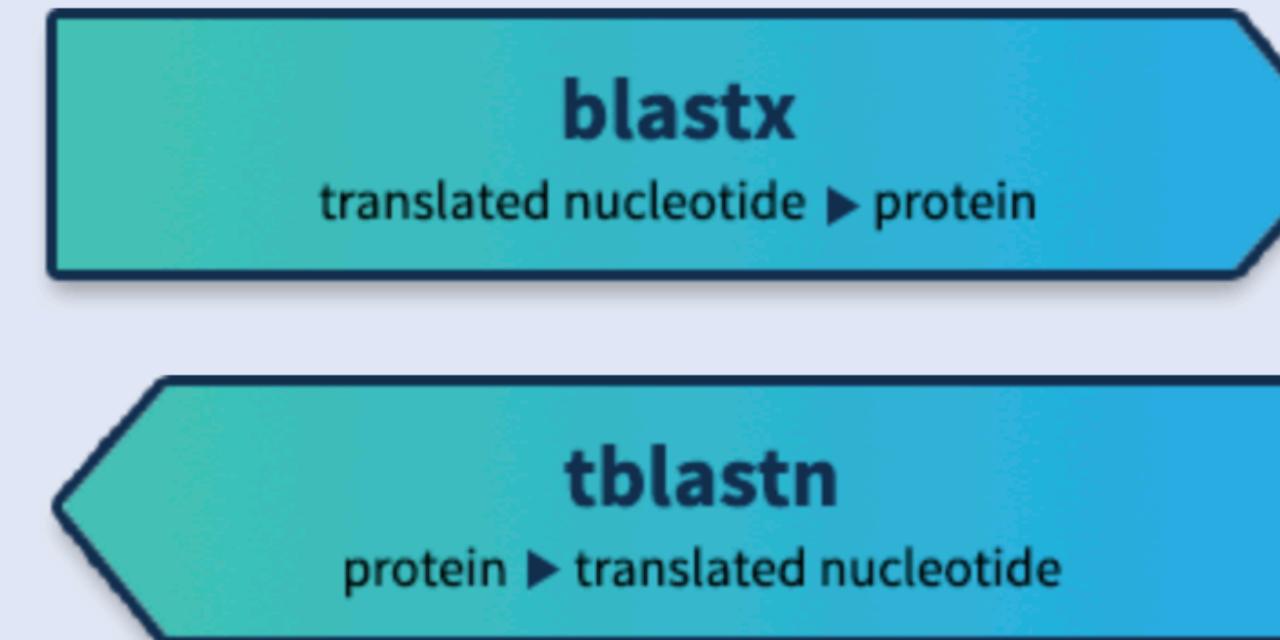
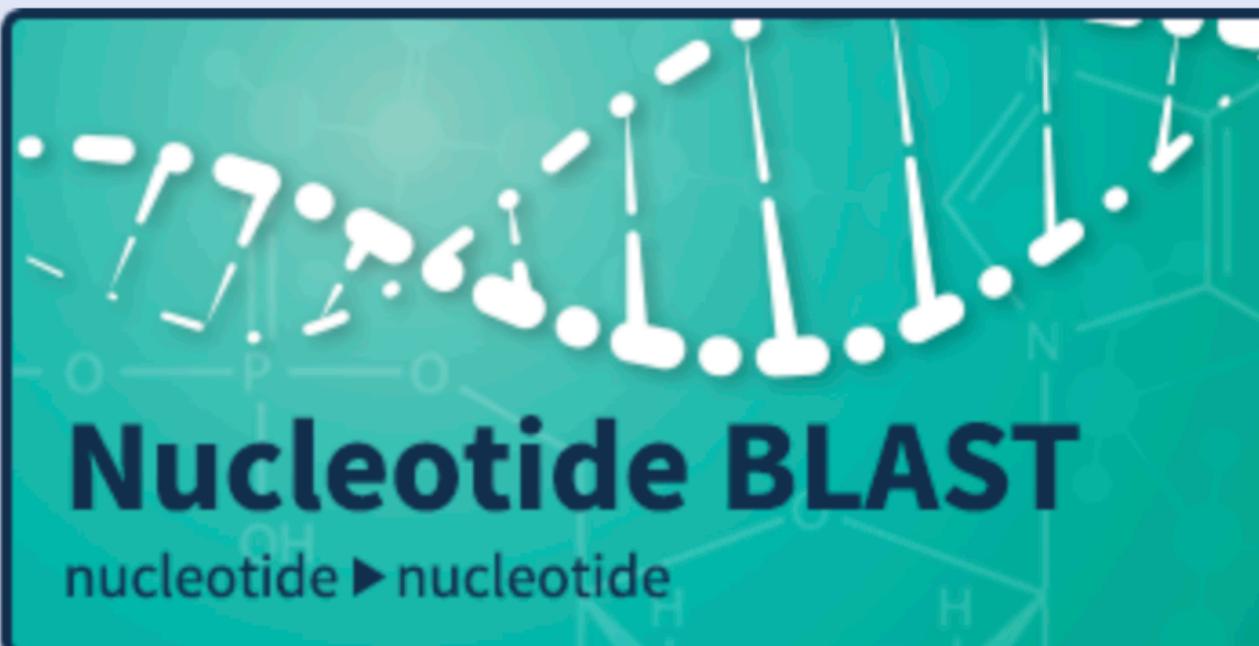
**ElasticBLAST is here**

ElasticBLAST is a new cloud based tool to run your BLAST searches faster and make you more effective.

Mon, 07 Feb 2022 12:00:00 EST

 More BLAST news...

**Web BLAST**





## COVID-19 Information

[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#) | [Español](#)



# NCBI SARS-CoV-2 Resources

Click here

## SARS-CoV-2 Data

3,099,725

SRA runs

4,070,899

Nucleotide records

3,215

ClinicalTrials.gov

231,749

PubMed

288,388

PMC

**Quick Navigation Guide**

- [Sequence Submission](#)
- [Literature](#)
- [Sequence-Related Resources](#)
- [Clinical Resources](#)
- [Other Websites](#)

## SARS-CoV-2 Data Hub

**Click here**

[Download](#)
[Quick Links](#)
[Betacoronavirus BLAST](#)
[CDC Outbreak Information](#)
[SARS-CoV-2 Articles in](#)
[PubMed](#)
[SRA Data](#)
[NCBI SARS-CoV-2 Resources](#)
[Datasets command line](#)
[Tabular View](#)
[Dashboard Visualizations](#)
[Mutations in SRA](#)
[Complete Tree](#)

Selected Results: 0

[Align](#)
[Build Phylogenetic Tree](#)

**Refine Results**

[Reset](#)

Virus +

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049 X

Accession +

Sequence Length +

**Nucleotide (4,083,187)** **Protein (24,012,480)** **RefSeq Genome (1)** [Select Columns](#)

**Accession** **Submitters** **Release Date** **Pangolin** **Isolate** **Species**

[Expand Table](#)

<input type="checkbox"/>	<a href="#">NC_045512</a> <small>RefSeq</small>	Wu,F., et al.	2020-01-13	B	Wuhan-Hu-1	Severe acute respirat
<input type="checkbox"/>	<a href="#">OM840138</a>	Andrews,K....	2022-02-27	B.1.2	ID-U1-IIDS-U0847	Severe acute respirat
<input type="checkbox"/>	<a href="#">OM840139</a>	Andrews,K....	2022-02-27	B.1.2	ID-U1-IIDS-U0850	Severe acute respirat
<input type="checkbox"/>	<a href="#">OM840140</a>	Andrews,K....	2022-02-27	B.1	ID-U1-IIDS-U0852	Severe acute respirat

[Tabular View](#)[Dashboard Visualizations](#)[Mutations in SRA](#)[Complete Tree](#)[Click here](#)**Statistics****1**

RefSeq Genomes

**23,865,769**

All Proteins

**4,070,899**

All Nucleotides

**38**

RefSeq Proteins

**892,702**

Complete Nucleotides

**Geographic and Time Distribution**

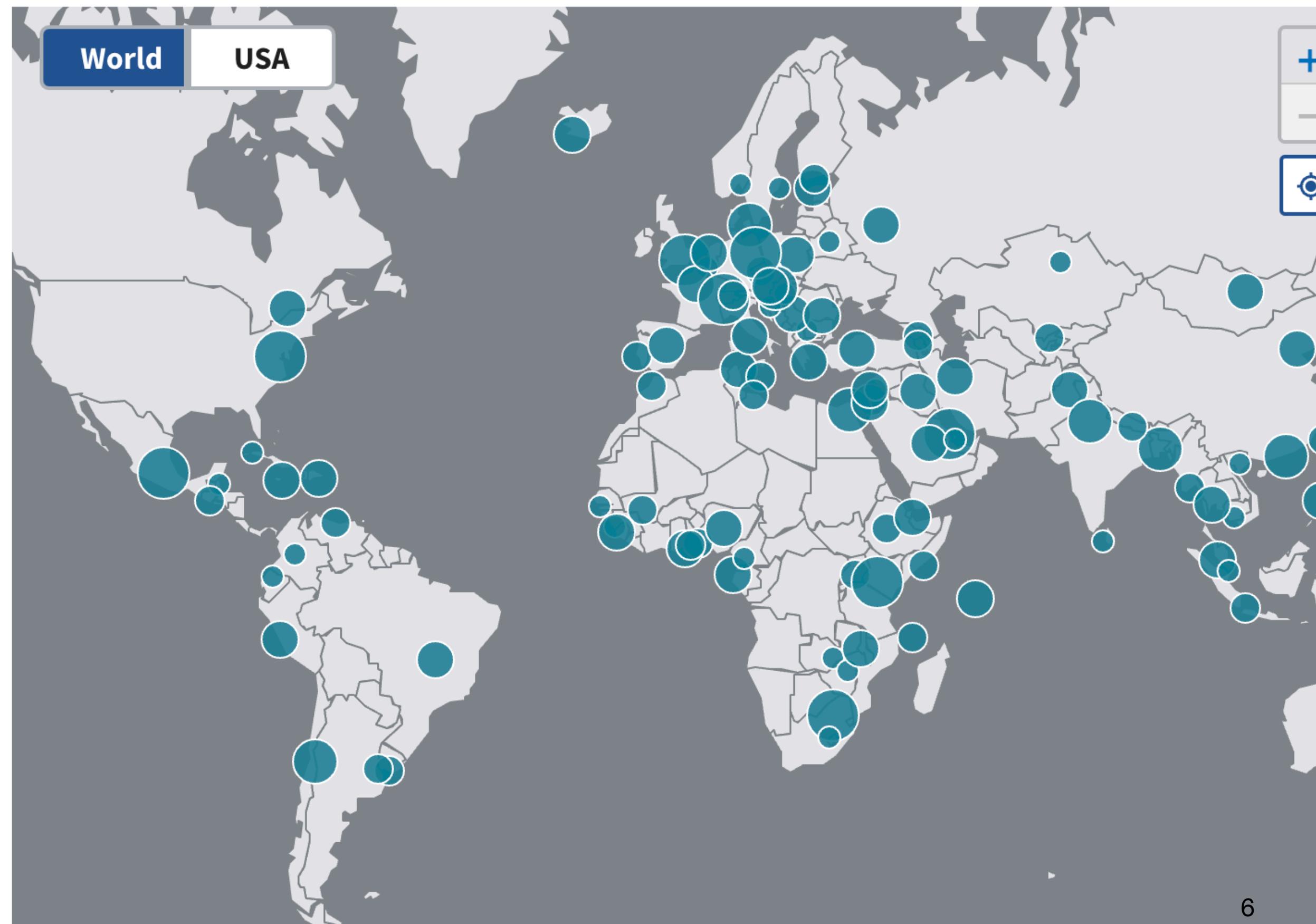
Choose locations to select SARS-CoV-2 sequence records with that collection location. Use the sliders or click date columns to select SARS-CoV-2 records by their sample collection date and/or their GenBank release date.

**Geographic Distribution**

Search for a country

**Collection Date**

Weekly

**Collection Date**

Weekly

1/1/2020 - 1/7/2020

2/23/2022 - 3/1/2022

**Release Date**

Weekly

1/1/2020 - 1/7/2020

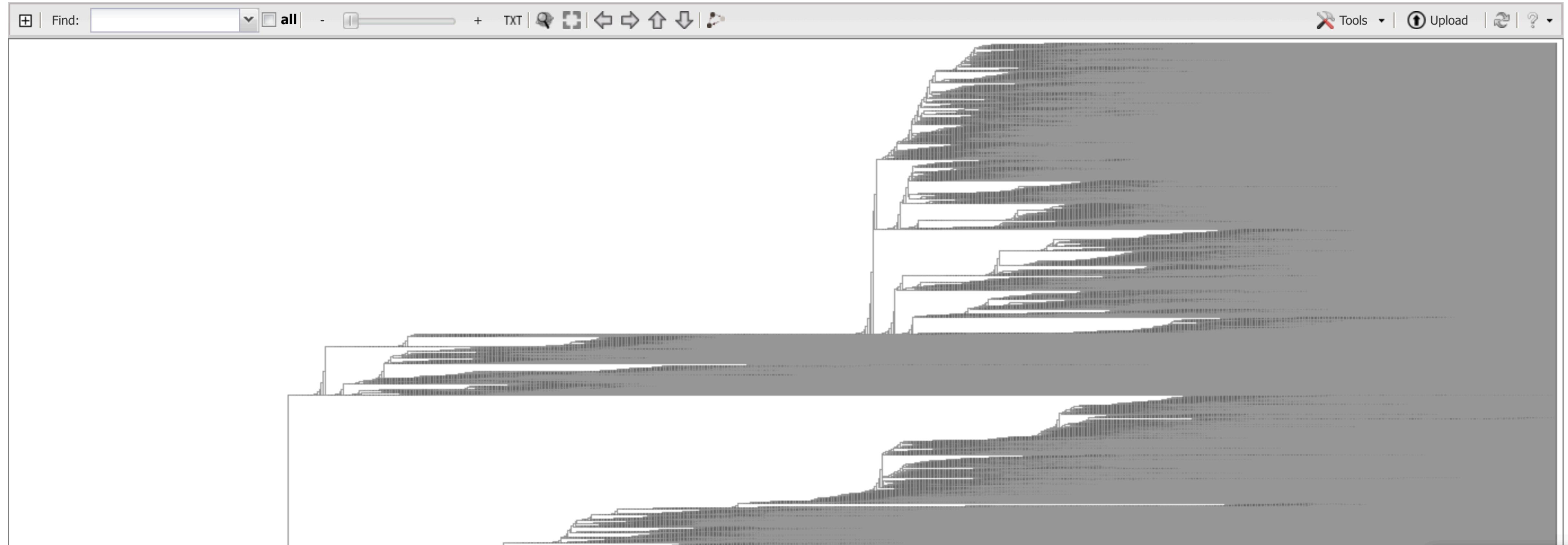
2/23/2022 - 3/1/2022

**Click here**[About Us](#) | [Find Data](#) | [Help](#) | [How to Participate](#) | [Submit Sequences](#)[Contact Us](#)

## Tree of complete SARS-CoV-2 Sequences

[Algorithm and parameters](#)

Updated: Fri Feb 25 2022



# Tree of complete SARS-CoV-2 Sequences

## Algorithm and parameters ↴

Our goal is to show the biological diversity of complete and nearly-complete natural SARS-CoV-2 sequences. All SARS-CoV-2 sequences which meet the following requirements are included:

- Sequence length between 29600 and 31000 bp after trimming ends which do not align with RefSeq
- Must have collection date (year + month)
- The number of ambiguous nucleotides must be < 1% in the trimmed sequences

Each sequence is aligned to the RefSeq ([NC\\_045512](#)), and mutations such as SNPs, indels, and long substitutions are calculated.

Distances are computed between pairs of sequences as the number of different mutations, where non-ambiguous different mutations are counted as 1, and ambiguous vs. non-ambiguous different mutations are counted 0.1 (one sequence has a non-ambiguous mutation and another sequence has an ambiguous mutation in the same place).

Any term in the leaf labels can be searched using the "Find" option.

The leaf labels are in this format: **accession.version|host|isolation source|geographic location|collection date|GenBank release date**,  
for example "**MT370839.1|Homo sapiens|oronasopharynx|USA|2020-03-14|2020-04-23**".

Updated: Fri Feb 25 2022

[Tabular View](#)[Dashboard Visualizations](#)[Mutations in SRA](#)[Complete Tree](#)[Click here](#)**Statistics****1**

RefSeq Genomes

**23,865,769**

All Proteins

**4,070,899**

All Nucleotides

**38**

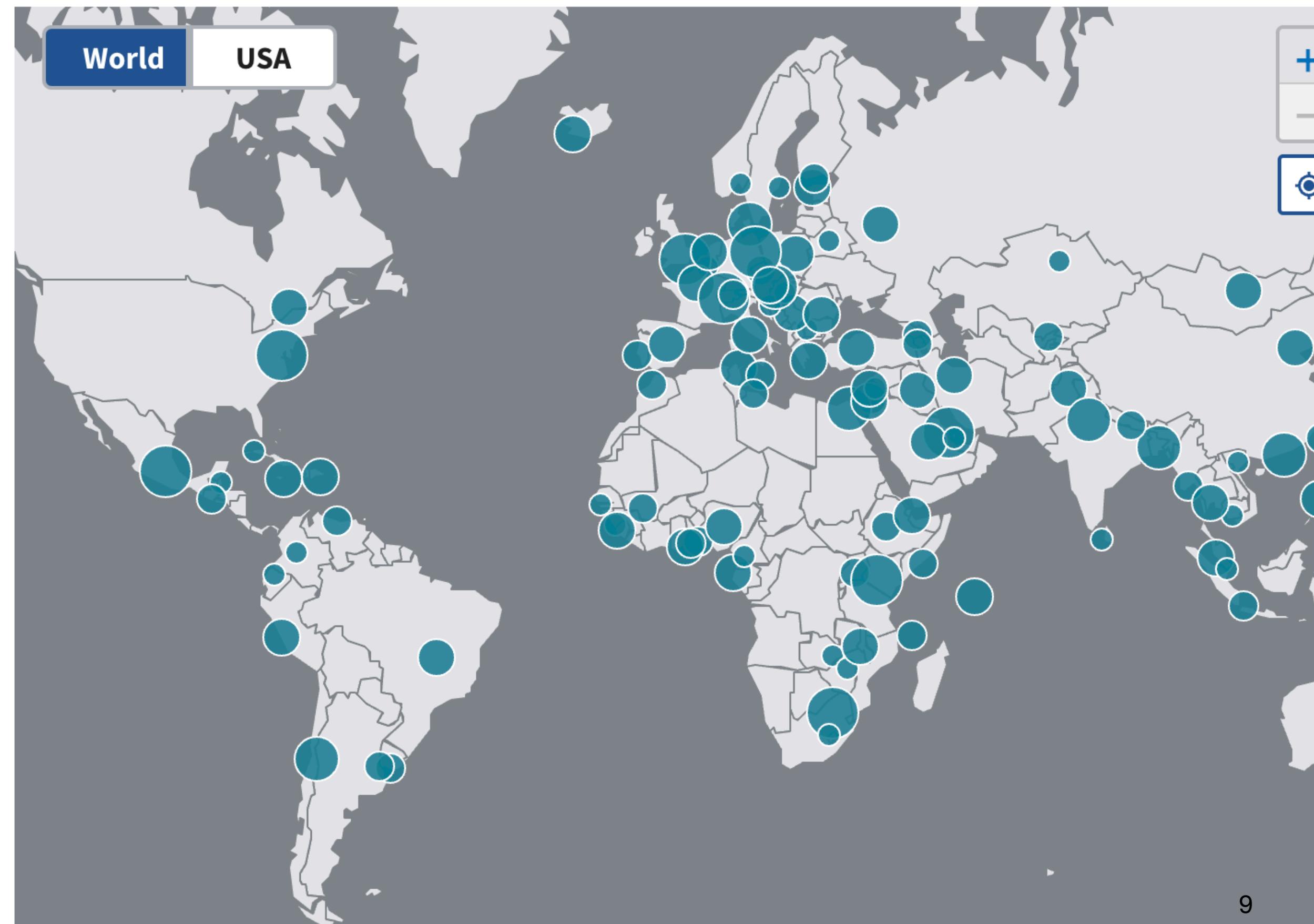
RefSeq Proteins

**892,702**

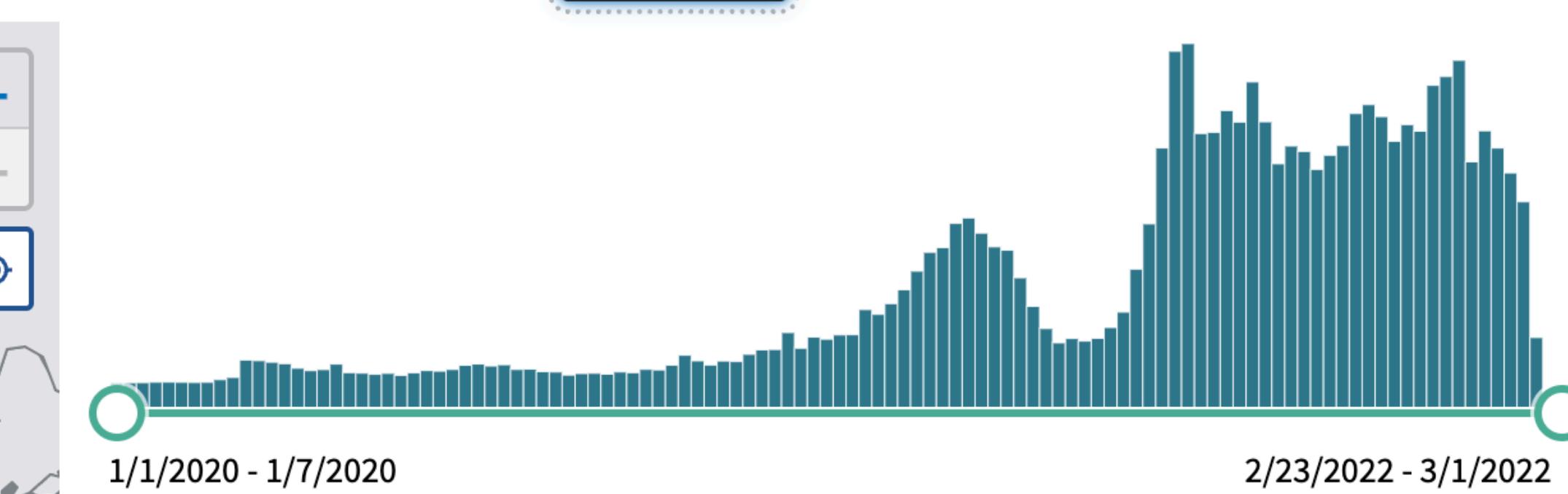
Complete Nucleotides

**Geographic and Time Distribution**

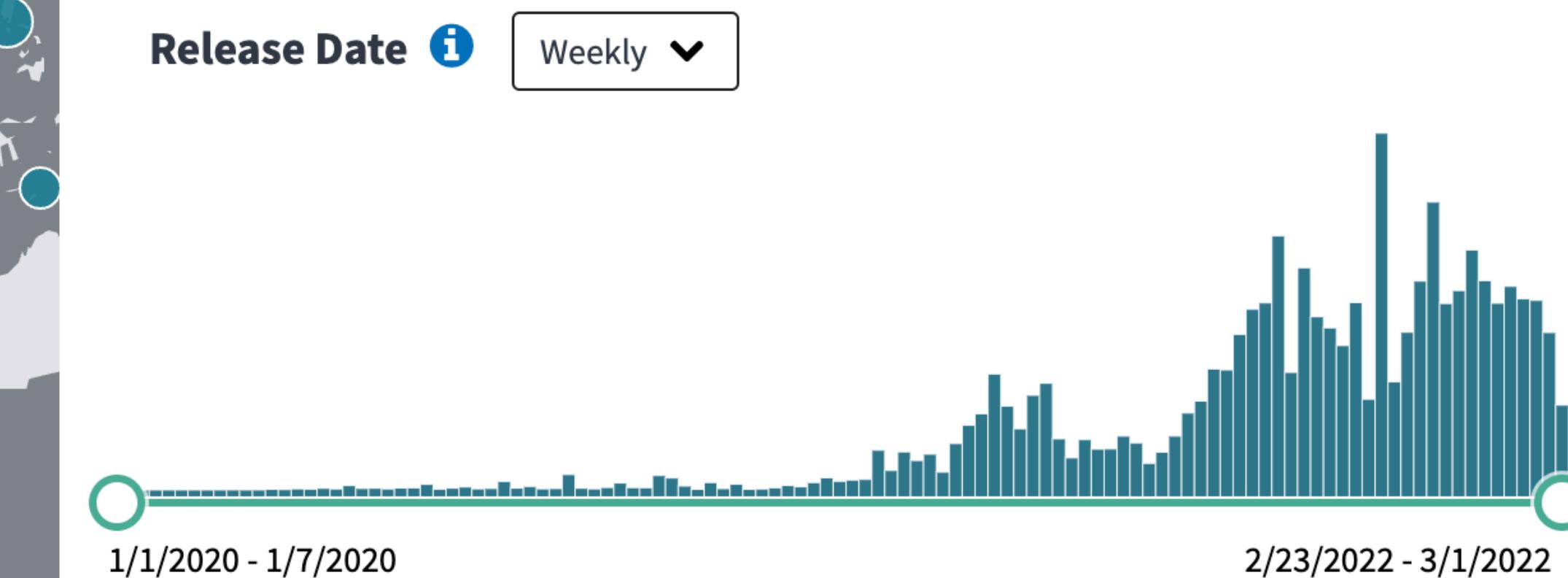
Choose locations to select SARS-CoV-2 sequence records with that collection location. Use the sliders or click date columns to select SARS-CoV-2 records by their sample collection date and/or their GenBank release date.

**Geographic Distribution** [i](#)[Search for a country](#)**Collection Date** [i](#)

Weekly

**Release Date** [i](#)

Weekly



**Then click here**

# SARS-CoV-2 Data Hub

Download ▾

Quick Links

Betacoronavirus BLAST  
CDC Outbreak Information

SARS-CoV-2 Articles in PubMed  
SRA Data

SARS-CoV-2 Resources  
Datasets  
Download line

Tabular View

Dashboard Visualizations

Mutations in SRA i Complete Tree i

Selected Results: 38

Align

Build Phylogenetic Tree

**Click here first**

Refine Results

Reset

Virus

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049

Accession

Sequence Length

Ambiguous Characters

Sequence Type

RefSeq x

RefSeq Genome Completeness +

Nucleotide Completeness +

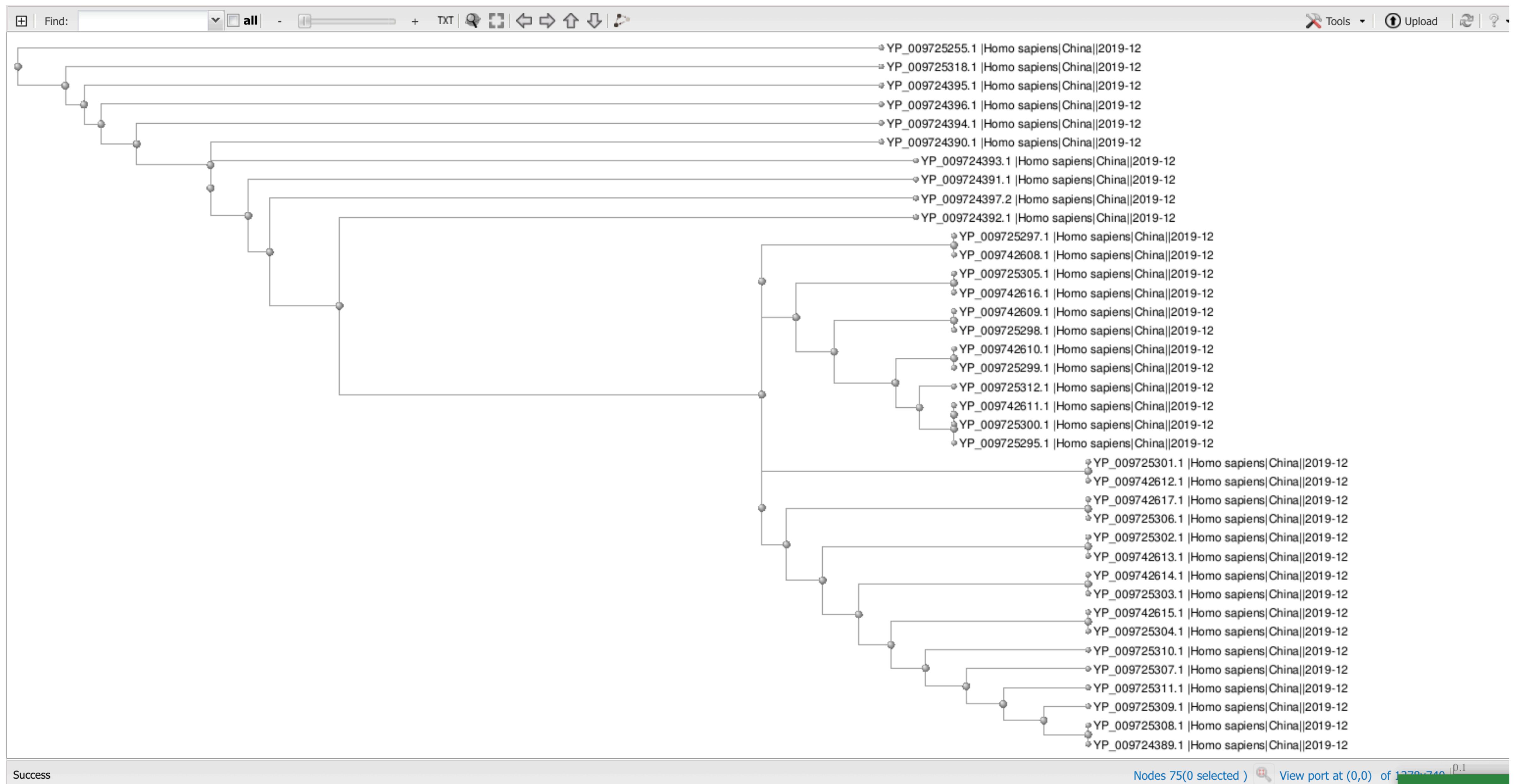
Pango lineage +

Random Sampling New! +

Expand Table

	Accession	Release Date	Pangolin	Isolate	Species	Length	Protein	Geo Location
	<input checked="" type="checkbox"/> <a href="#">YP_009742608</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	180	leader protein	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742609</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	638	nsp2	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742610</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	1945	nsp3	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742611</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	500	nsp4	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742612</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	306	3C-like proteinase	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742613</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	290	nsp6	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742614</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	83	nsp7	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742615</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	198	nsp8	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742616</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	113	nsp9	China
	<input checked="" type="checkbox"/> <a href="#">YP_009742617</a> RefSeq	2020-03-30	B	Wuhan-Hu-1	Severe acute respiratory s...	139	nsp10	China

# Phylogenetic Tree



## SARS-CoV-2 Data Hub

[Download](#) ▾

[Quick Links](#)
[Betacoronavirus BLAST](#)
[CDC Outbreak Information](#)
[SARS-CoV-2 Articles in](#)
[PubMed](#)
[SRA Data](#)
[NCBI SARS-CoV-2 Resources](#)
[Datasets command line](#)
[Tabular View](#)
[Dashboard Visualizations](#)
[Mutations in SRA](#) ⓘ

[Complete Tree](#) ⓘ

Selected Results: 0

[Align](#)
[Build Phylogenetic Tree](#)
**Click here**

[Refine Results](#)
[Reset](#)
[Virus](#)


Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049


[Accession](#)

[Sequence Length](#)


Nucleotide (4,083,187)
Protein (24,012,480)
RefSeq Genome (1)
[Select Columns](#)

<input type="checkbox"/>	Accession	Submitters	Release Date	Pangolin	Isolate	Species
<input type="checkbox"/>	<a href="#">NC_045512</a> <small>RefSeq</small>	Wu,F., et al.	2020-01-13	B	Wuhan-Hu-1	Severe acute respirat
<input type="checkbox"/>	<a href="#">OM840138</a>	Andrews,K....	2022-02-27	B.1.2	ID-U1-IIDS-U0847	Severe acute respirat
<input type="checkbox"/>	<a href="#">OM840139</a>	Andrews,K....	2022-02-27	B.1.2	ID-U1-IIDS-U0850	Severe acute respirat
<input type="checkbox"/>	<a href="#">OM840140</a>	Andrews,K....	2022-02-27	B.1	ID-U1-IIDS-U0852	Severe acute respirat

Expand Table

## Refine Results

Reset

Virus

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049

Accession

Sequence Length

Ambiguous Characters

Sequence Type

RefSeq Genome Completeness

Nucleotide Completeness

Pango lineage

Random Sampling New!

Isolate

Proteins

Provirus

Geographic Region

Host

Submitters

Isolation Source

Collection Date

Release Date

Environmental Source

Lab Host

Vaccine Strain

## Nucleotide Completeness

 complete (899,433)

 partial (3,183,754)
Random Sampling New!

Filters sequences that were collected for a purpose of baseline surveillance.

 Include (899,433)

 Exclude (403,540)

 Only (495,893)

## Isolation Source

 feces (26)

 lung (74)

 lung, oronasopharynx (103)

 oronasopharynx (45,206)

 oronasopharynx, oronasopharynx (10)

 placenta (3)

 saliva, oronasopharynx (2,613)

**Nucleotide (74)**      [Protein \(862\)](#)      [RefSeq Genome \(0\)](#)

[Select Columns](#)

**Show Filters**

<input type="checkbox"/> Accession	Release Date	Pangolin	Isolate	Species	Length	Nuc Completeness
<input type="checkbox"/> <a href="#">OL913104</a>	2021-12-20	None	C57MA14	Severe acute respiratory s...	29903	complete
<input type="checkbox"/> <a href="#">MT582498</a>	2020-06-09	B.3	NRW-02.1	Severe acute respiratory s...	29782	complete
<input type="checkbox"/> <a href="#">MW306667</a>	2020-11-30	B.1.9	ETLKVET2	Severe acute respiratory s...	29867	complete
<input type="checkbox"/> <a href="#">MT457400</a>	2020-05-12	B.1.8	NB03_index	Severe acute respiratory s...	29890	complete
<input type="checkbox"/> <a href="#">MT457401</a>	2020-05-12	B.1.8	NB04_index	Severe acute respiratory s...	29891	complete
<input type="checkbox"/> <a href="#">MT396266</a>	2020-04-28	B.1.8	1	Severe acute respiratory s...	29880	complete
<input type="checkbox"/> <a href="#">MZ710319</a>	2021-08-05	B.1.369	ICU-TA-113-496	Severe acute respiratory s...	29877	complete
<input type="checkbox"/> <a href="#">MZ710321</a>	2021-08-05	B.1.369	ICU-TA-523-496	Severe acute respiratory s...	29835	complete
<input type="checkbox"/> <a href="#">MZ710322</a>	2021-08-05	B.1.369	ICU-TA-113-496-Ampliseq	Severe acute respiratory s...	29878	com

[Feedback](#)

## Select Columns

### Add or Remove Columns From the Results Table

#### Columns to add:

+ [SRA Accession](#)

+ [Submitters](#)

+ [Genus](#)

+ [Family](#)

+ [Molecule type](#)

+ [Sequence Type](#)

+ [Genotype](#)

+ [Segment](#)

+ [Publications](#)

+ [Country](#) New!

+ [BioSample](#)

#### Displayed columns

(click to remove from the view):

✗ [Accession](#)

✗ [Release Date](#)

✗ [Pangolin](#)

✗ [Isolate](#) New!

✗ [Species](#)

✗ [Length](#)

✗ [Nuc Completeness](#)

✗ [Geo Location](#)

✗ [USA](#)

✗ [Host](#)

✗ [Isolation Source](#)

Then click here



SARS-CoV-2 Data Hub

Download ▾

Quick Links

Betacoronavirus BLAST  
CDC Outbreak Information

SARS-CoV-2 Articles in  
PubMed  
SRA Data

NCBI SARS-CoV-2 Resources  
Datasets command line

Tabular View

Dashboard Visualizations

Mutations in SRA ⓘ Complete Tree ⓘ

Selected Results: 74

Align

Build Phylogenetic Tree

Click here first

Refine Results

Reset

Virus



Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049



Accession



Sequence Length



Ambiguous Characters



Nucleotide (74) Protein (862) RefSeq Genome (0)					
<input checked="" type="checkbox"/>	Accession	Release Date	Pangolin	Isolate	Species
<input checked="" type="checkbox"/>	<a href="#">OL913104</a>	2021-12-20	None	C57MA14	Severe acute respiratory s...
<input checked="" type="checkbox"/>	<a href="#">MT582498</a>	2020-06-09	B.3	NRW-02.1	Severe acute respiratory s...
<input checked="" type="checkbox"/>	<a href="#">MW306667</a>	2020-11-30	B.1.9	ETLKVET2	Severe acute respiratory s...
<input checked="" type="checkbox"/>	<a href="#">MT457400</a>	2020-05-12	B.1.8	NB03_index	Severe acute respiratory s...
<input checked="" type="checkbox"/>	<a href="#">MT457401</a>	2020-05-12	B.1.8	NB04_index	Severe acute re...

Select Columns

Feedback

## Download Results

X

### Step 1 of 3: Select Data Type

Sequence data  
(FASTA Format)

Nucleotide

Coding Region

Protein

Accession List

Nucleotide

Protein

Assembly

Current table view result

CSV format

XML format

**Next**

## Download Results

X

### Step 2 of 3: Select Records

Download Selected Records

Download All Records

**Back**

**Next**

# Download Results

X

## Step 3 of 3: Select FASTA definition line

Use default: Accession GenBank Title

Build custom: Accession GenBank Title

Assembly	
SRA Accession	<span>Add ➤</span> <span>◀ Remove</span>
Submitters	
Release Date	
Danogolin	

Accession

GenBank Title

Back

Download

# Download Results

X

## Step 3 of 3: Select FASTA definition line

Use default: Accession GenBank Title

Build custom: Accession Genus Species Geo Location

Segment
Publications
Country
USA
Host

Add ➤  
◀ Remove

Accession
Genus
Species
Geo Location

Back

Download



Downloads — vi sequences-2.fasta — 119×31

>OM062573.1 |Betacoronavirus|Severe acute respiratory syndrome-related coronavirus|China: Beijing  
ATCAAAGGTTATACCTTCCCAGGTAACAAACCAACCAACTTCGATCTCTGTAGATCT  
GTTCTCTAAACGAACCTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCAC  
CACGCAGTATAATTAAATAACTAATTACTGTCGTTGACAGGACACGAGTAACCTTCTATC  
TTCTGCAGGCTGCTTACGGTTCGTCCGTGTTGCAGCCGATCATCAGCACATCTAGGTTT  
TGTCCGGGTGTGACCGAAAGGTAAGATGGAGAGCCTTGTCCCTGGTTCAACGAGAAAAC  
ACACGTCCAACTCAGTTGCCTGTTTACAGGTTCGCGACGTGCTCGTACGTGGCTTGG  
AGACTCCGTGGAGGAGGTCTTATCAGAGGCACGTCAACATCTAAAGATGGCACTTGTGG  
CTTAGTAGAAGTTGAAAAAGGCCTTGCCTCAACTTGAACAGCCCTATGTGTTCATCAA  
ACGTTCGGATGCTCGAACACTGCACCTCATGGTCATGTTATGGTGAGCTGGTAGCAGAACT  
CGAAGGCATTCAAGTACGGTCGTAGTGGTGAGACACTTGGTGTCCCTGTCCCTCATGTGGG  
CGAAATACCAGTGGCTTACCGCAAGGTTCTTCGTAAGAACGGTAATAAAGGAGCTGG