# The UCSC SARS-CoV-2 Genome Browser: One-stop Shopping for the Latest Molecular Details of SARS-CoV-2

**David Haussler**
**UC Santa Cruz Genomics Institute**
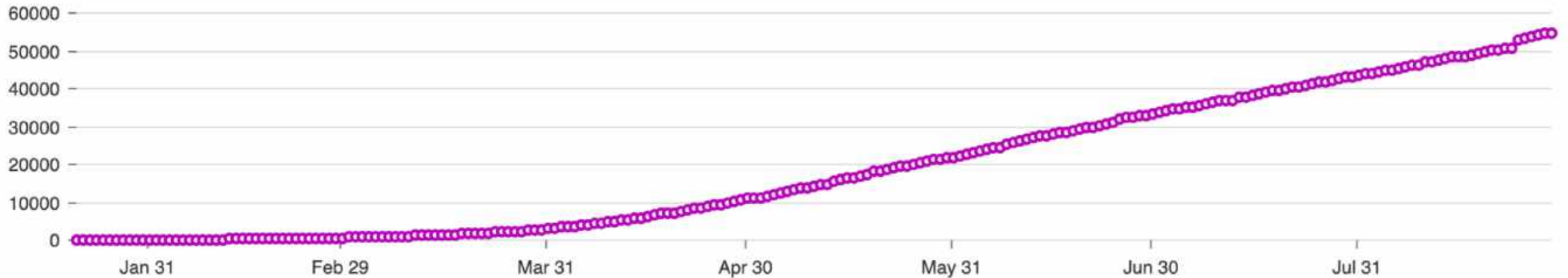
https://genome.ucsc.edu/covid19.html

https://genome.ucsc.edu/cgi-bin/hgTracks?db=wuhCor1

# SARS-CoV-2 Research is Generating Data at an Astonishing Pace

Jan 21, 2020 - Aug 30, 2020

**4121** New in the Past 7 Days | **54917** Cumulative Papers



**In April, SARS CoV-2 papers had a doubling time of ~14.5 days.**

**(The virus doubling time in April was ~7 days)**

*Source: primer.ai*

# Genomic Data has also grown at an exponential rate



**First virus genome released on Jan 10, 2020**

# Genomic Data has also grown at an exponential rate



**More than 100,000 genomes now sequenced!**

*Source: GISAID*

# How do we make use of all this genomic and molecular data for analysis?



Experiments — Predictions/Annotations

Standard Data Formats

Scale
NC_045512v2:

10 kb | wuhCor1
5,000 | 10,000 | 15,000 | 20,000 | 25,000

NCBI Genes2

ORF1a
ORF1ab
S
ORF3a
E
M
ORF6
ORF7a
ORF7b
ORF8
N
ORF10

**Visualize in Genome Browser**

**Compare Datasets**

UNIVERSITY OF CALIFORNIA
SANTA CRUZ Genomics Institute

# The Genome Browser annotates nucleotides with information stored in tracks



http://genome.ucsc.edu/s/SARS_CoV2/Figure1

# Users add annotations via "crowd-sourced" annotations

## Users add annotations to spreadsheet at: http://bit.ly/cov2annots



# Insert your annotations below. Mouse-over the headers to see instructions. Contact maxh@ucsc.edu if you have questions on this form or suggestions.

# Note that the annotations do not go immediately to the public site, they are only made public once per day. To show the current version of the annotations track, click this link: https://genome-test.gi.ucsc.edu/

| Start | End | Label | Category | Long descriptive text | URL to website or paper wit | Your email |
|---|---|---|---|---|---|---|
| 1 | 450 | 5UTR | genes | 5' UTR structured RNA | https://www.biorxiv.org/conten | jferna10@ucsc.edu |
| 23605 | 23617 | furin_cleavage | proteins | Novel polybasic protein cleavage site (aa seq RRAR) that can be processed by | https://www.biorxiv.org/conten | haussler@ucsc.edu |
| 22871 | 23086 | ACE2_receptor | proteins | receptor biding site motif in the virus S protein for the human ACE2 protein | https://www.nature.com/article | haussler@ucsc.edu |
| 23923 | 23980 | fusion_peptide | proteins | fusion peptide in the viral S protein facilitating fusion of viral membrane with hos | https://www.nature.com/article | haussler@ucsc.edu |

## Annotation appears on genome.ucsc.edu after approval:



| Scale | | 50 bases | | | | | | | | | wuhCor1 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| NC_045512v2: | 23,890 | 23,900 | 23,910 | 23,920 | 23,930 | 23,940 | 23,950 | 23,960 | 23,970 | 23,980 | |

UniProt Protein Products (Polypeptide Chains)

S glycoprotein  Q D K N T Q E V F A Q V K Q I Y K T P P I K D F G G F N F S Q I L

Spike protein S2 Q D K N T Q E V F A Q V K Q I Y K T P P I K D F G G F N F S Q I L

Crowd-sourced data: annotations contributed via bit.ly/cov2annots

fusion_peptide >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>

## Click element for information

| Category | proteins |
|---|---|
| Long descriptive text | fusion peptide in the viral S protein facilitating fusion of viral membrane with host cell membrane |
| URL to website or paper with further info | https://www.nature.com/articles/s41422-020-0305-x |


UNIVERSITY OF CALIFORNIA SANTA CRUZ Genomics Institute

# Sequencing artifacts can impact phylogenetic inferences

**Mutations can trace transmission BUT artifacts can confound analysis**

**Some "mutations" that influence tree topology are lab-specific artifacts**



**List of "problematic sites" to mask from analysis available on genome.ucsc.edu**
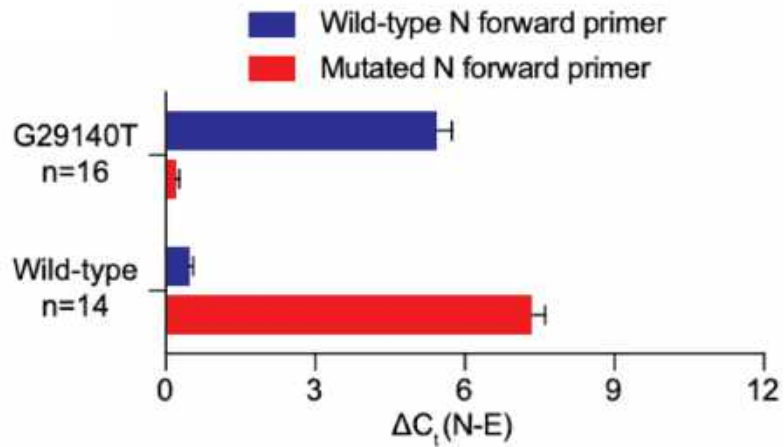
# Variation can affect the ability to accurately detect virus

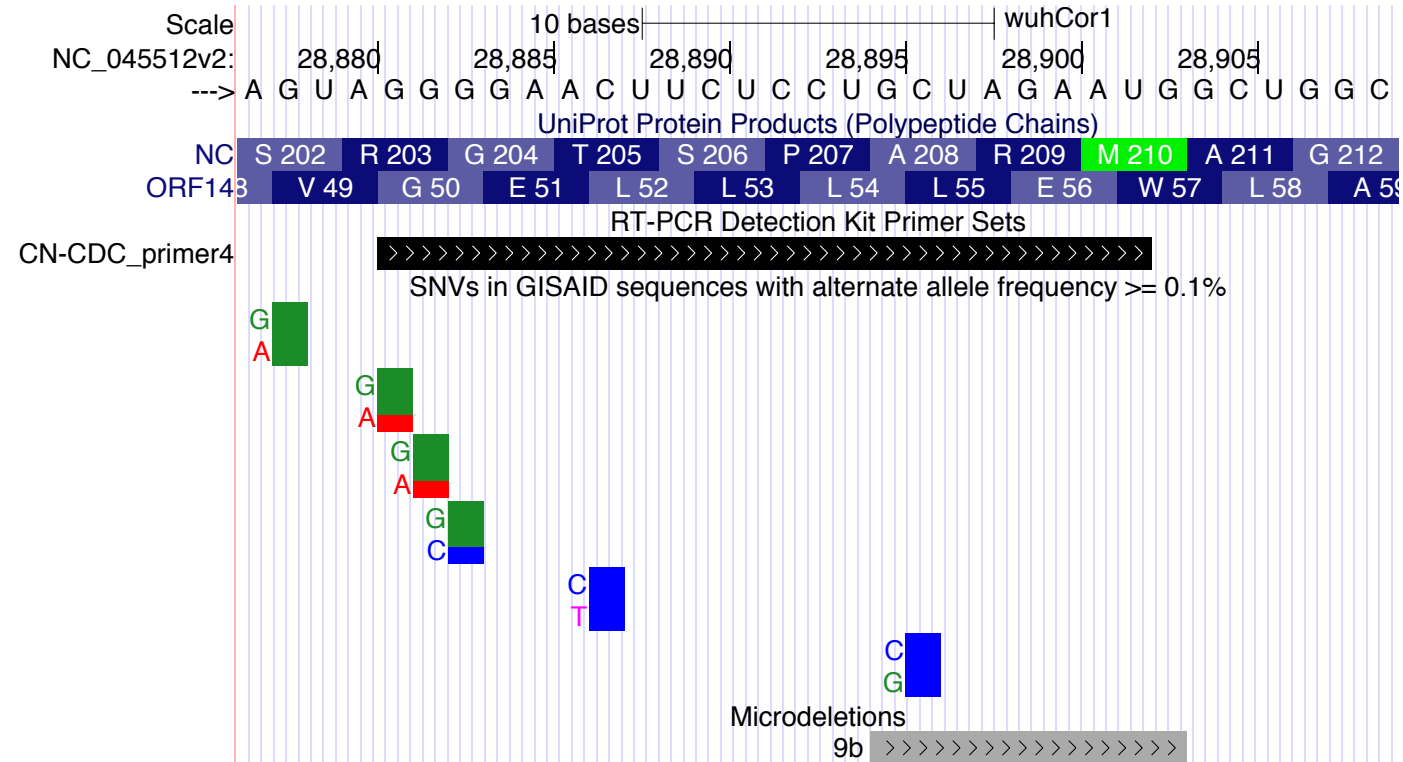## CZI Biohub recently reported variant that affected detection



**Identification of a polymorphism in the N gene of SARS-CoV-2 that adversely impacts detection by a widely-used RT-PCR assay**

Manu Vanaerschot, Sabrina A. Mann, James T. Webber, Jack Kamm, Sidney M. Bell, John Bell, Si Noon Hong, Minh Phuong Nguyen, Lienna Y. Chan, Karan D. Bhatt, Michelle Tan, Angela M. Detweiler, Alex Espinosa, Wesley Wu, Joshua Batson, David Dynerman, CLIAHUB Consortium, Debra A. Wadford, Andreas S. Puschnik, Norma Neff, Vida Ahyong, Steve Miller, Patrick Ayscue, Cristina M. Tato, Simon Paul, Amy Kistler, Joseph L. DeRisi, Emily D. Crawford

**doi:** https://doi.org/10.1101/2020.08.25.265074

## Genome Browser overlays standard detection primers with emerging variants



## Deletions and variants alter primer choice as pandemic progresses.

**Genome browser view of variants and predicted antibody-spike contacts**

**Interactive viewer on click**



**Neutralizing antibodies could contact mutable residues.
However, no contact sites for S antibodies have variants >1% frequency.
Nothing to worry about so far!**

PDB from Wu et al., *Science*, 2020

# Virus and host receptor interfaces rapidly evolve

**Alignment of Coronaviruses (SARS-CoV-2 Browser)**

**Alignment of 100 vertebrates (Human Genome Browser: hg38)**



Green = synonymous mutations   Red = non-synonymous mutations   Yellow = alignment gap

## Residues in S-ACE2 interface are rapidly evolving in both virus and host.

## SARS-CoV-2 is evolved to be successful in humans.

ACE2 residues that contact S

ACE2

S



UNIVERSITY OF CALIFORNIA SANTA CRUZ Genomics Institute

**Workflow:**

1. Sequence & Assemble Genome

2. Upload sequences to database.

3. Place new sequences in context of existing global phylogenetic tree.

4. Trace spread via genomic epidemiology.

Although vastly better than previous efforts, each step is not truly "real-time". Specifically, current phylogenetics software is not built to scale to 100,000 genomes!



*Nexstrain.org*

# Ultrafast Sample placement on Existing tRees (UShER) is a step towards real-time viral genomics

**Workflow:**

1. Sequence & Assemble Genome



User Sample
Existing Sample 1
Existing Sample 2

2. Upload sequences to database.

3. **Place new sequences in context of existing global phylogenetic tree.**

4. Trace spread via genomics.

Although vastly better than previous efforts, each step is not truly "real-time". Current methods not built to scale for 100,000 genomes!

| Method | Time to Place 1000 Sequences |
|---|---|
| PAGAN2 | 24+ Hours |
| IQ-TREE2 | 24+ Hours |
| TreeBeST | 24+ Hours |
| RAxML epa | 24+ Hours |
| **UShER** | **43.2 SECONDS** |

UsHER uses parsimony annotations of tree branches & an optimized binary file

UNIVERSITY OF CALIFORNIA SANTA CRUZ | Genomics Institute

# UShER is now available and integrated into the SARS-CoV-2 Genome Browser



**New samples added in blue in interactive environment with alignment and mutation calls**

**Link & Demonstration:** https://genome.ucsc.edu/cgi-bin/hgPhyloPlace
www.github.com/russcd/USHER_DEMO/

UNIVERSITY OF CALIFORNIA SANTA CRUZ | Genomics Institute

# An idealized roadmap for how future outbreaks might be traced in true real-time

1. Sequence & assemble genome using nanopore and laptop in the field

2. Upload sequences to database automatically.

3. Place new sequences in context of existing global phylogenetic tree and get analysis immediately.
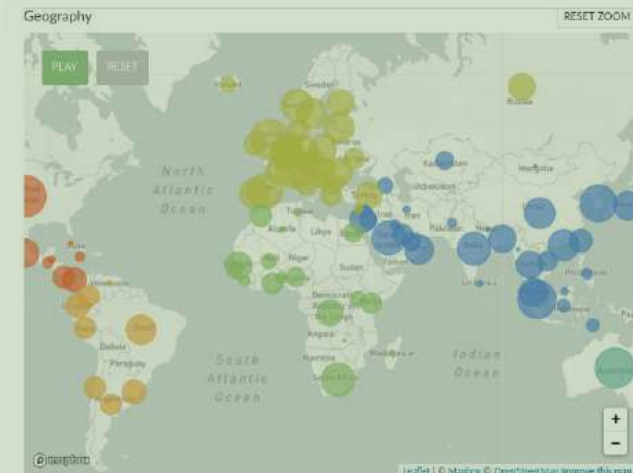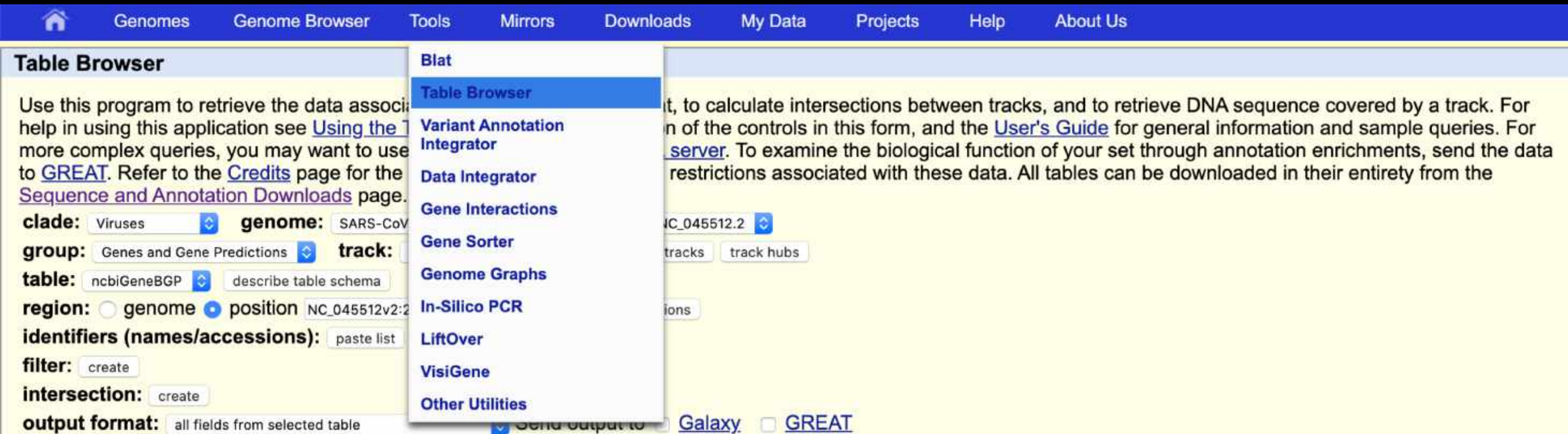
**UShER**

4. True real time genomic contact tracing!

UNIVERSITY OF CALIFORNIA
SANTA CRUZ | Genomics Institute

# All data is easily accessible via the SARS-CoV-2 Browser

| | Genomes | Genome Browser | Tools | Mirrors | Downloads | My Data | Projects | Help | About Us |
|---|---|---|---|---|---|---|---|---|---|

## Table Browser

Use this program to retrieve the data associa̶̶̶̶̶̶̶̶̶̶t, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see Using the T̶̶̶̶̶̶̶̶̶̶n of the controls in this form, and the User's Guide for general information and sample queries. For more complex queries, you may want to use̶̶̶̶̶̶̶̶̶̶ server. To examine the biological function of your set through annotation enrichments, send the data to GREAT. Refer to the Credits page for the̶̶̶̶̶̶̶̶̶̶ restrictions associated with these data. All tables can be downloaded in their entirety from the Sequence and Annotation Downloads page.

**Blat**
**Table Browser**
**Variant Annotation Integrator**
**Data Integrator**
**Gene Interactions**
**Gene Sorter**
**Genome Graphs**
**In-Silico PCR**
**LiftOver**
**VisiGene**
**Other Utilities**

**clade:** Viruses  **genome:** SARS-CoV  NC_045512.2

**group:** Genes and Gene Predictions  **track:**  tracks  track hubs

**table:** ncbiGeneBGP  describe table schema

**region:** ○ genome ● position  NC_045512v2:2  ions

**identifiers (names/accessions):** paste list

**filter:** create

**intersection:** create

**output format:** all fields from selected table  Galaxy  GREAT

## Downloading Data using MariaDB (MySQL)

The UCSC Genome Browser uses MariaDB as the backend database server. MariaDB is a community-developed, commercially supported fork of the MySQL relational database management system, intended to remain free and open-source software under the GNU General Public License.

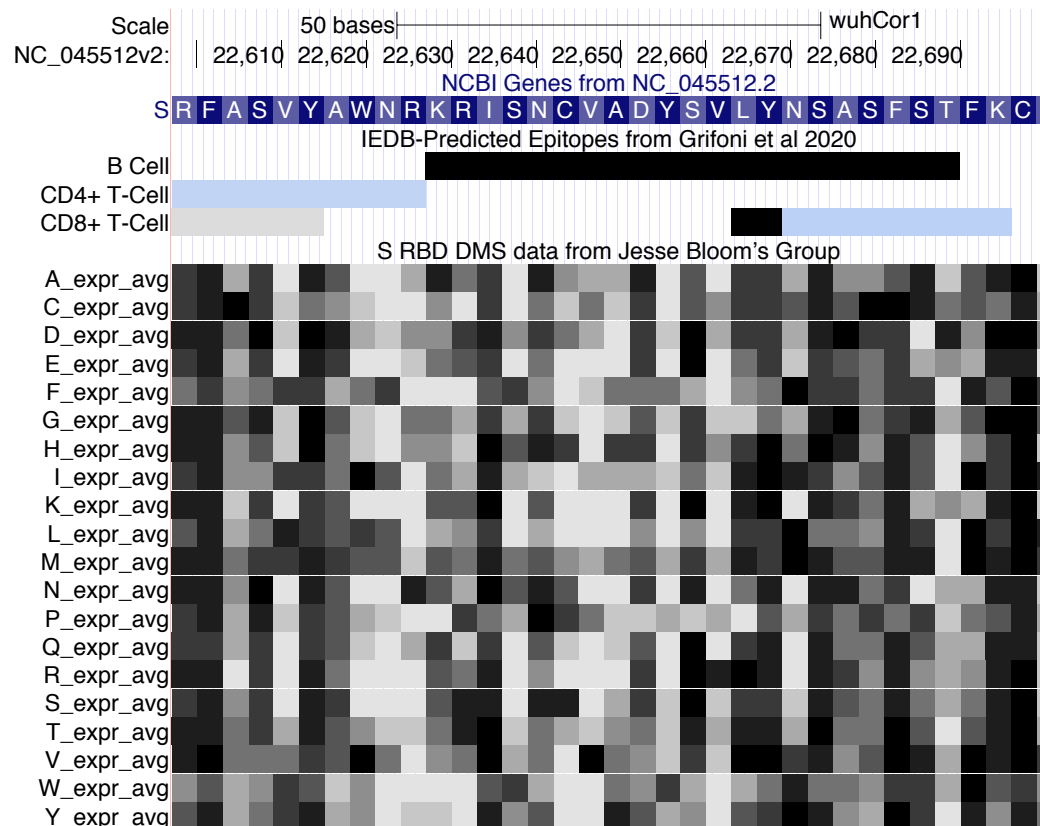We have two MariaDB databases for public access:

- **genome-mysql.soe.ucsc.edu** (located on the US west coast)

- **genome-euro-mysql.soe.ucsc.edu** (located in Europe)

These servers allow MySQL access to the same set of data currently available on our public Genome Browser site. The data are synchronized weekly with the main databases on our public site. During synchronization, the MariaDB server can be intermittently out of sync with the main website for a short period of time. The weekly synchronization takes place on Monday mornings from 4:00 am to 9:00 am Pacific Time (GMT -7:00 in summer, GMT -8:00 in winter).
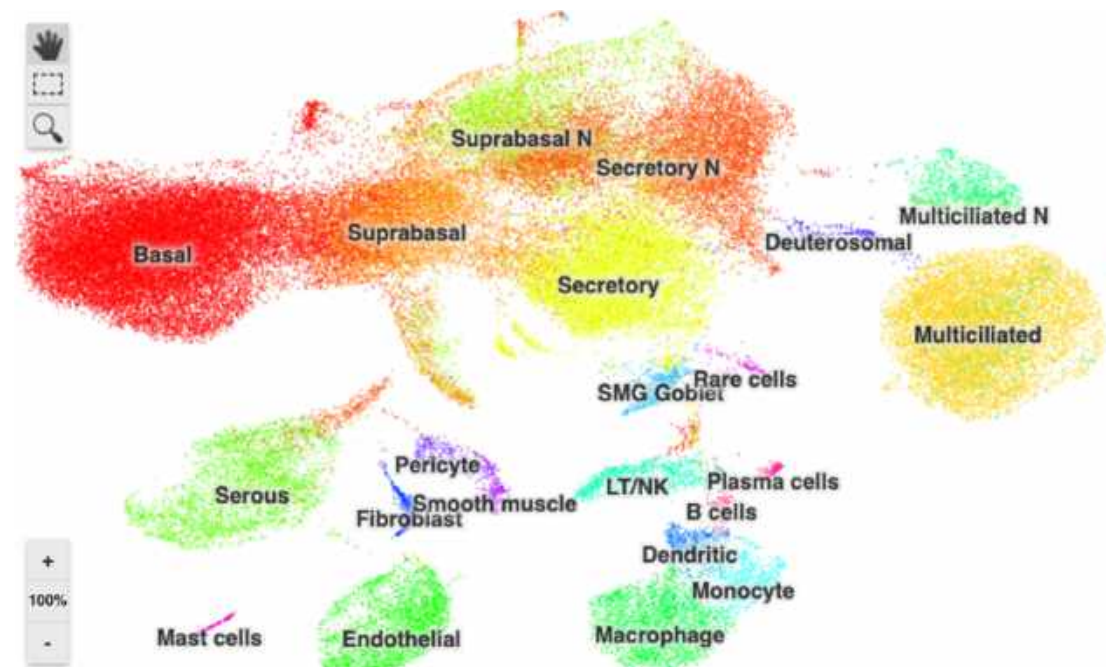
# Please consider adding your genomic data!

## Data from many HHMI colleagues already online:

**S protein Deep Mutational Scanning Data** *(Jesse Bloom Lab)* available at genome.ucsc.edu:

**COVID19 Cell Atlas** *(Mark Krasnow Lab)* available in interactive scRNA-seq browser at **cells.ucsc.edu:**

# Acknowledgements

## UCSC SARS-CoV-2 Browser

Jason Fernandes
Hiram Clawson
Angie Hinrichs
Jairo Navarro Gonzalez
Brian T . Lee
Luis R. Nassar
Brian J. Raney
Kate R. Rosenbloom
Santrupti Nerli
Arjun A. Rao
Daniel Schmelter
Alastair Fyfe
Nathan Maulding
Ann S. Zweig
Todd M. Lowe
Manuel Ares Jr
Jim Kent
Max Haeussler

## Recurrent Errors & UShER

Russ Corbett Lab (UCSC)
      Bryan Thornlow
      Landen Gozashti

Yatish Turakhia (UCSC now, starting lab at UCSD, 2021)

Rob Lanfear (Australian National Univ)

Nick Goldman Lab (EBI)
      Nicola De Maio
      Conor R. Walker
      Lukas Weilguny

Rui Borges (Institut für Populationsgenetik)
Greg Slodkowicz (MRC)