While we are waiting:
- Connect to Tufts network (on campus or VPN)
- Chrome browser https://galaxy.cluster.tufts.edu
- Login with Tufts credentials
- Let me know if you have trouble logging in

# Intro to Next Generation Sequencing Data Analysis with Galaxy

Rebecca Batorsky

Pr Bioinformatics Scientist

Nov 2021

# Research Technology Team

**Delilah Maloney**
High Performance Computing Specialist

**Kyle Monahan**
Senior Data Science Specialist

**Shawn Doughty**
Manager, Research Computing

**Rebecca Batorsky**
Senior Bioinformatics Scientist

**Chris Barnett**
Senior Geospatial Analyst

**Tom Phimmasen**
Senior Data Consultant

**Patrick Florance**
Director, Academic Data Services

**Jake Perl**
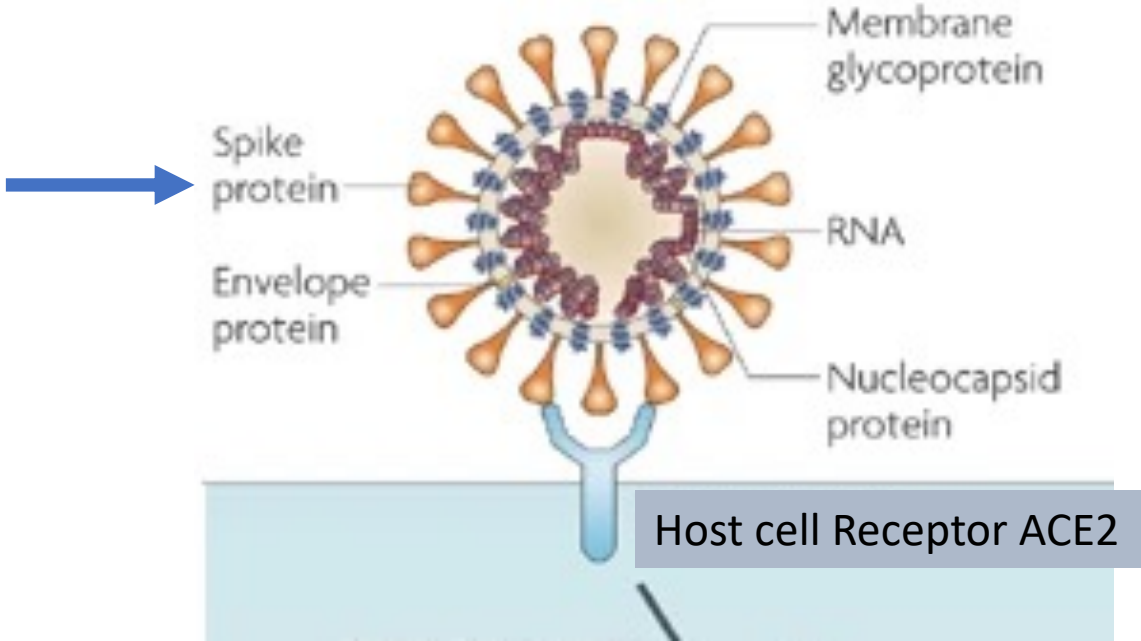Digital Humanities NLP Specialist

**Carolyn Talmadge**
Senior GIS Specialist

**Uku-Kaspar Uustalu**
Data Science Specialist

✓ Consultation on Projects and Grants
✓ High Performance Compute Cluster
✓ Workshops

https://it.tufts.edu/research-technology

# SARS-CoV2 Spike Protein



Host cell Receptor ACE2

# SARS-CoV2 Spike Protein VOCs

Spike protein



Major mutations of SARS-CoV-2
VOCs in Spike receptor binding domain (RBD)

Increase
- ACE2 affinity
- Transmissibility

# NCBI SARS-CoV-2 Resources

## Quick Navigation Guide

- Sequence Submission
- Literature
- Sequence-Related Resources
- Clinical Resources
- Other Websites

## SARS-CoV-2 Data

| | | |
|---|---|---|
| **1,457,511** | **2,276,801** | **3,215** |
| SRA runs | Nucleotide records | ClinicalTrials.gov |
| **198,246** | | **235,648** |
| PubMed | | PMC |

https://www.ncbi.nlm.nih.gov/sars-cov-2/

# Outline

# Viral Genome Next Generation Sequencing (NGS)

- Specimen Collected

- NGS reads

TAAGCGACGTA
Read1

TTACCAGATAGGTT
Read2

GGGCCAACTACC
Read3

# Viral Genome Next Generation Sequencing (NGS)

- Specimen Collected



- RNA extraction



- cDNA synthesis (using virus-specific primers)
- Amplification



- NGS library prep



- NGS sequencing



flowcell

- Alignment

# Short Read Alignment

# Paired end vs Single end reads



https://www.biostars.org/p/267167/

# Outline

# Log into Galaxy

- Connect to Tufts Network, either on campus or via VPN
- Visit https://galaxy.cluster.tufts.edu
- Log in with you Tufts credentials

# Galaxy on the Tufts High Performance Compute (HPC) Cluster

User laptop

## HPC Cluster

| node | node | node | node |
|------|------|------|------|
| node | node | node | node |

Proxy web-server (Apache)

Job scheduler (SLURM)

PAM Authentication

Galaxy web-server

PostgreSQL Database server

13

# User Interface

# User Interface

# Galaxy User Interface

To return to home screen

Minimize/Adjust toolbars

# History



Create New History

View all Histories

# History



Create New History

View all Histories

# Tools

# Tools



Click on the name of the tool to open it in the main panel

# Importing data

# Importing data

# Log into Galaxy and open course website

- Connect to Tufts Network, either on campus or via VPN
- Visit https://galaxy.cluster.tufts.edu
- Log in with you Tufts credentials
- Visit course website https://rbatorsky.github.io/intro-to-galaxy-ngs-sarscov2



Navigate to **Obtain_Data** section

# Outline

Introduction to the Galaxy Platform

Obtain data from NCBI

SARS-Cov-2 Alpha variant reference sequence
SARS-Cov-2 Delta variant NGS sample

Process Raw Reads (QC, adapter trimming)

Read Alignment

Visualization

# NGS details

# Next Generation Sequencing (NGS)

# Next Generation Sequencing (NGS)



**4. FRAGMENTS BECOME DOUBLE STRANDED**

Attached terminus  Free terminus  Attached terminus

The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

**5. DENATURE THE DOUBLE-STRANDED MOLECULES**

Attached  Attached

Denaturation leaves single-stranded templates anchored to the substrate.

**6. COMPLETE AMPLIFICATION**

Clusters

Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

https://sites.google.com/site/himbcorelab/illumina_sequencing

# Next Generation Sequencing (NGS)



https://sites.google.com/site/himbcorelab/illumina_sequencing

# Next Generation Sequencing (NGS)



This Illumina video is great for visualization!