# More data, more problems?
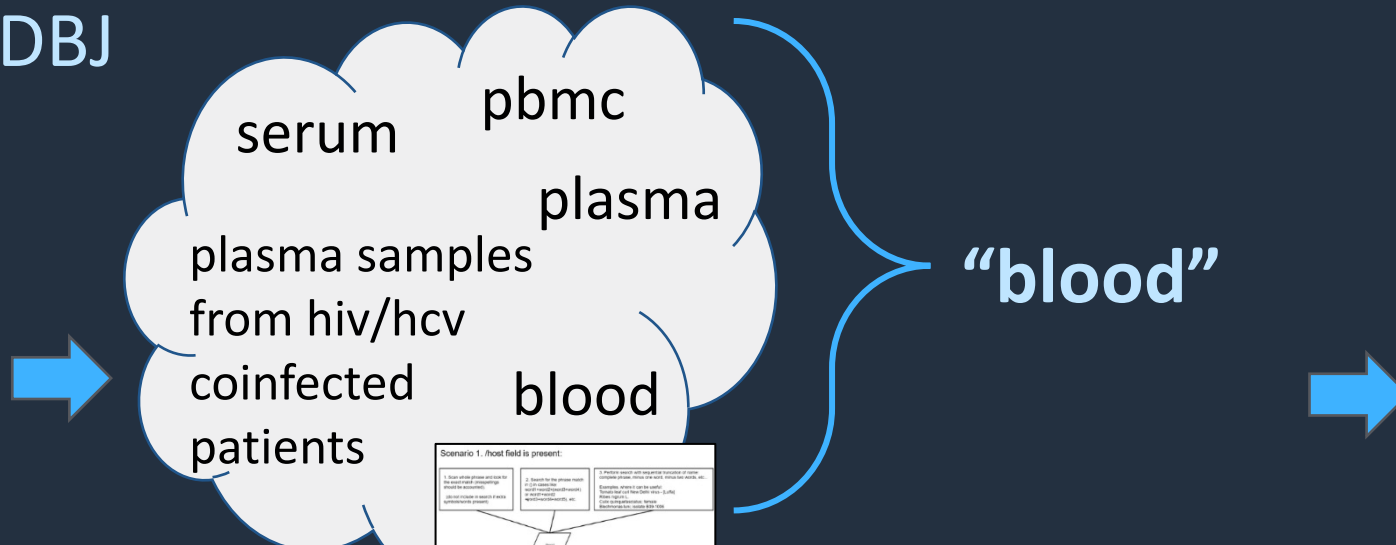
- How do you find data for viruses that have certain characteristics?

- From a geographic location, or collected recently?

- Which name do you search with?
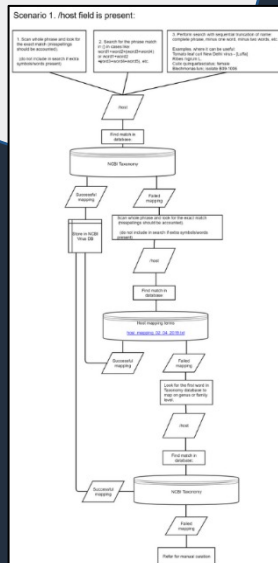
- Part of a lineage or subtype?
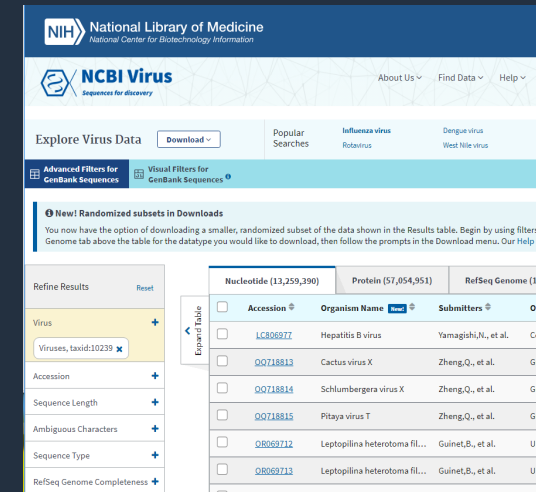
# NCBI Virus

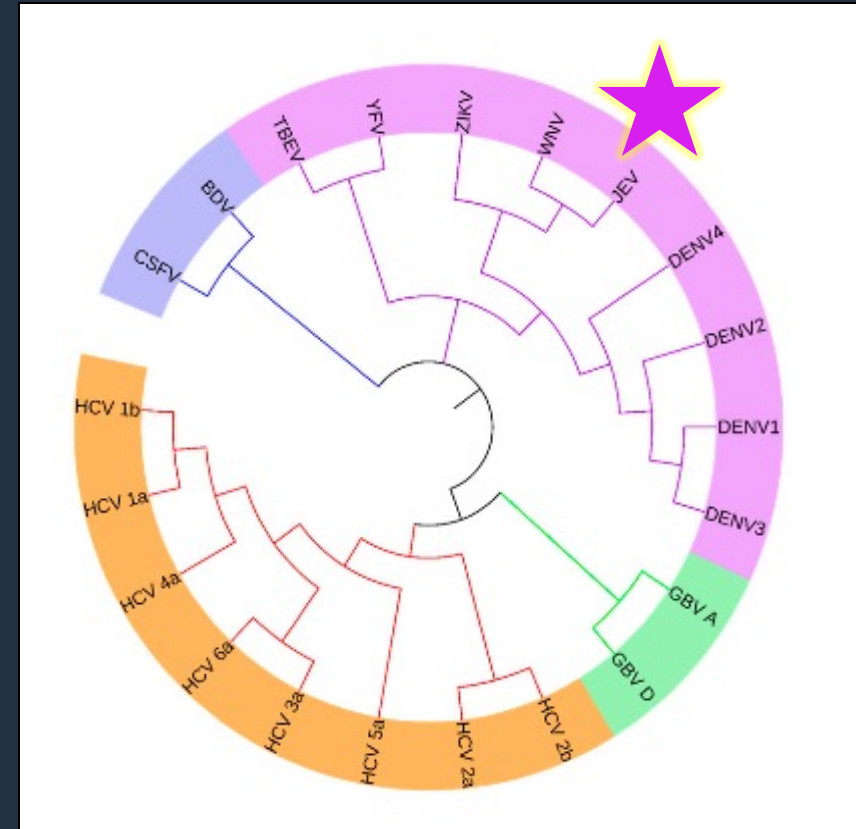# NCBI Virus Curation



Submission to GenBank (or DDBJ or ENA)

serum

pbmc

plasma

plasma samples from hiv/hcv coinfected patients

blood

"blood"

Host
Sample location in host
Collection date
Geographic regions
Lab or vaccine strain
Environmental strain

# RefSeq records

- Representatives of a species or unclassified group

- Based on GenBank record

- Usually, 1 per species

- ICTV exemplars and "new" viruses

- Standard gene & protein annotation for well characterized viruses

# Role of RefSeqs in Public Health Response

- GenBank team + NCBI Virus team

- 24 hrs - RefSeq with improved annotation

- Standardized annotation – protein coding regions & names

- Demonstrated SARS-CoV-2 taxonomy

- RefSeq was used to create VADR models
  - Submissions are normally manually reviewed - Automated submission of sequences

- Provides widely available sequence coordinates,  S:D614G

# NCBI Virus Results Table

- Accession: GenBank, SRA, BioProject, BioSample

- Geographic region (North America; USA; GA)

- Submitting author, organization, location
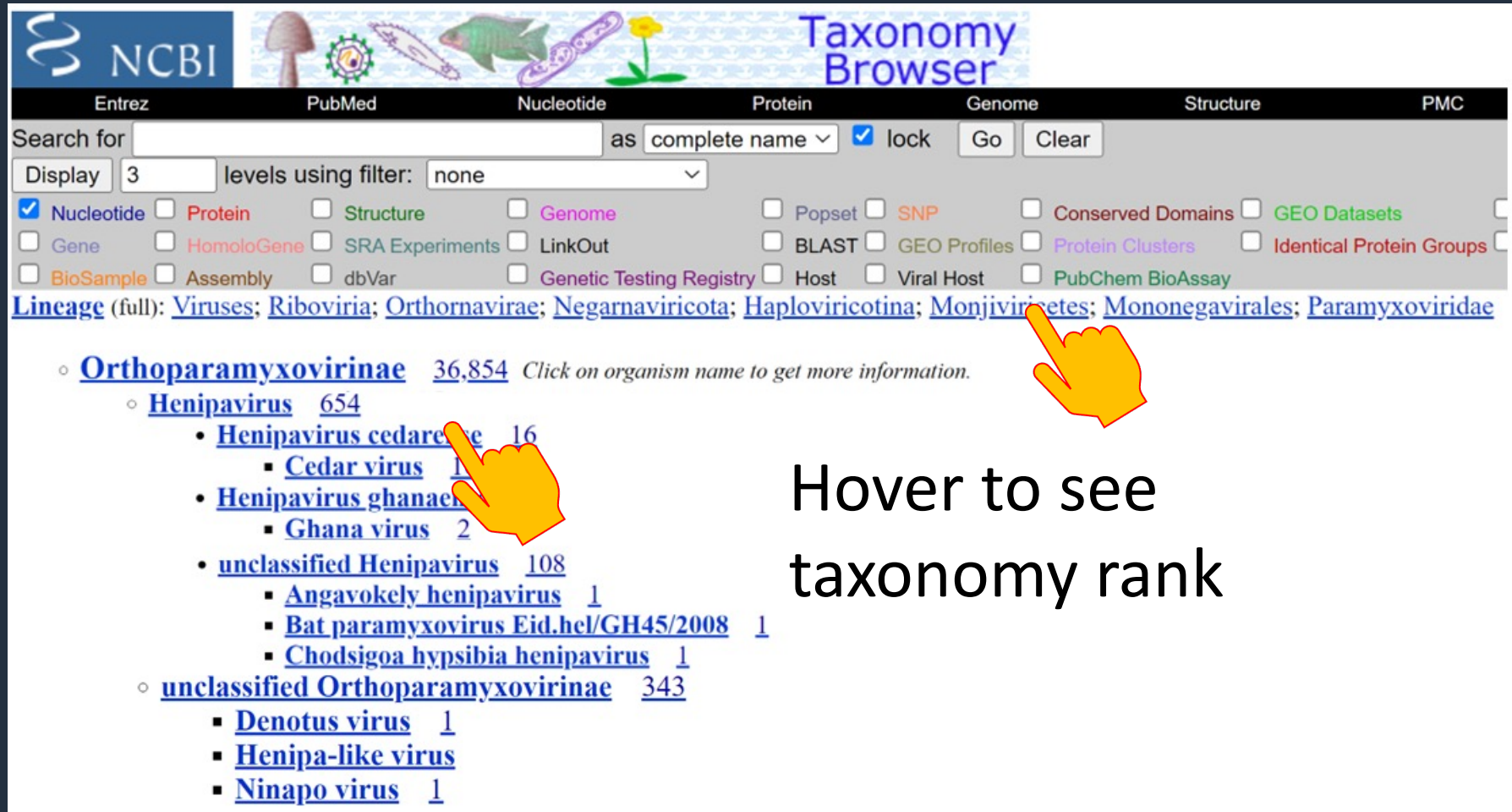
- Several more

# Pango Lineages in NCBI Virus

- Automated tool ☺

- Each evening: Check for updates to the tools, Run on all SARS-CoV-2 sequences

- Filter by lineage

- Include lineage in Results Table download, or defline in downloads

  - Versions are also included in Results Table download:

  - 4.3.1/1.26/v0.1.12/0.3.19/0.6.2

  - pangolin 4.3.1 / pangolin-data 1.24 / constellations v0.1.12 / scorpio 0.3.19 / UShER 0.6.2

# NCBI Taxonomy

- Virus species: ICTV

  - Gives us a <u>formalized</u> way to talk about virus groups

- Virus name, common names: Nipah virus, SARS-CoV-2

- When no official taxonomy is available yet, NCBI places sequence records into "unclassified" bins

- NCBI taxonomy also includes synonyms, virus names, common names, etc. to make data searches easier

  - Nipah-virus, 2019-nCoV, Newcastle disease virus

- Includes nodes below the species level

# NCBI Taxonomy Bowser

# NCBI Taxonomy Bowser



Unclassified bins for viruses without ICTV taxonomy

# NCBI Taxonomy Bowser – New!



Search "NCBI Datasets" → Taxonomy

Genomes ≠ Nucleotide records

Please share your feedback

# NCBI Virus = NCBI Taxonomy



Search with <u>any</u> level of taxonomy, common name, many abbreviations, etc

Interactive dashboard for exploring sequence data, understanding biases, and refining selections

# Flexible download formats

FASTA: nucleotide, coding regions, or proteins
- Customize the defline/header

Accessions lists

Results table
- Choose which fields to include in the table

Randomized subsets, stratified by country, collection year, release year, or host

# NCBI Datasets – Command line or API

- Programmatically access virus data
- Nucleotide & Protein sequences, annotation, and metadata report
- Same normalized metadata as NCBI Virus
- Great for large data downloads
- https://www.ncbi.nlm.nih.gov/datasets/docs/how-tos/virus/

# Help us help you.

- Lineages can change over time – but data provided during submission is archival
- NCBI Virus can associate sequence records with up-to-date lineages
- & help distribute lineage references
- We need:
  - Community-accepted classification schemes
  - Automated classification tools with long-term support

# Acknowledgements

**NCBI Virus team**

Olga Blinkova

Brett Spurrier

David Kristensen

Yuri Ostapchuk

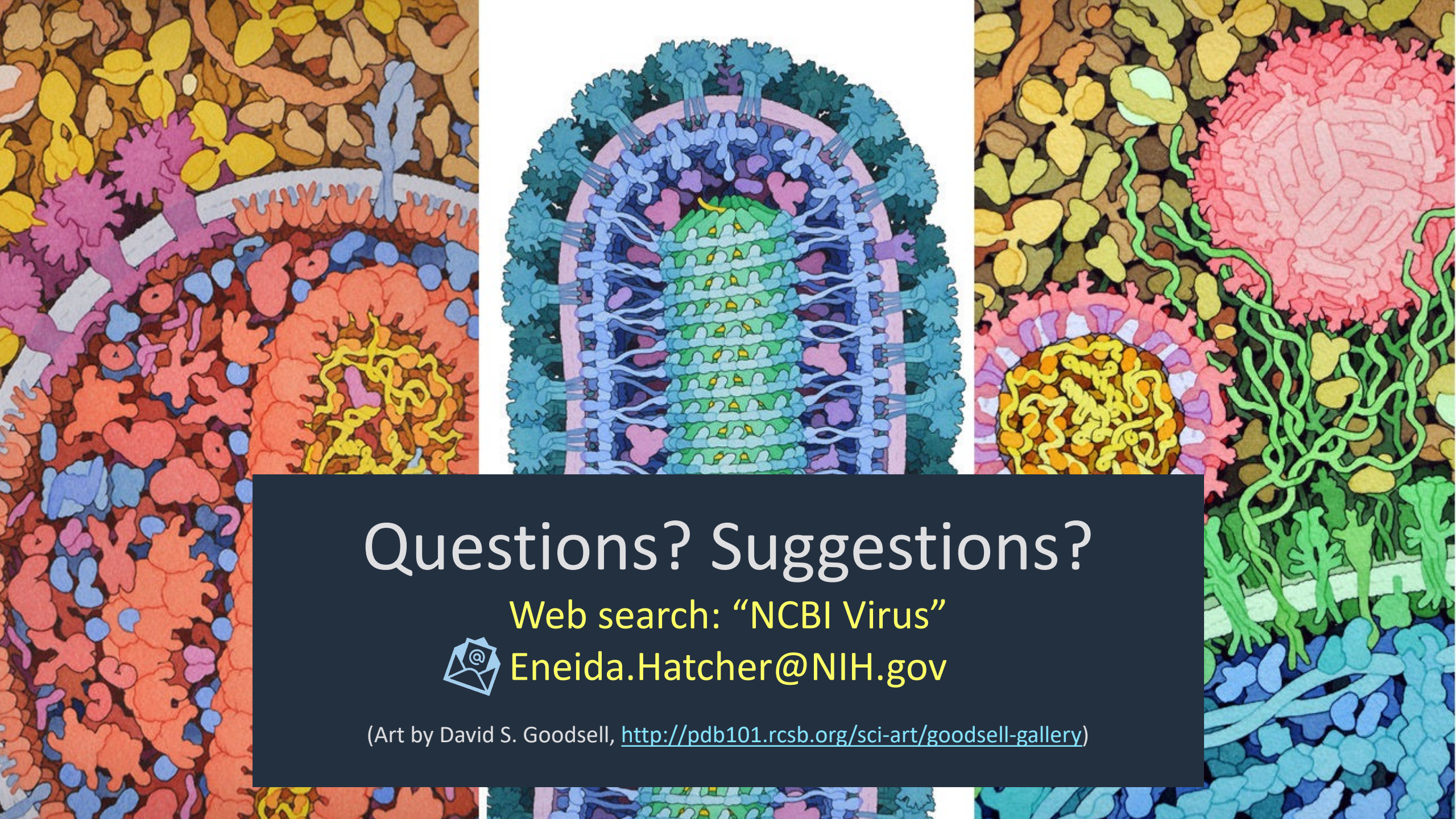Sergey Resenchuk

Igor Tolstoy

Anna Glodek

Ravinder Eskandary

Ryan Connor – former

J. Rodney Brister - former

**NCBI & NLM colleagues**

Everyone who has contributed to the shared public data

NIH U.S. National Library of Medicine
National Center for Biotechnology Information

NCBI

# Questions? Suggestions?

Web search: "NCBI Virus"

✉ Eneida.Hatcher@NIH.gov

(Art by David S. Goodsell, http://pdb101.rcsb.org/sci-art/goodsell-gallery)