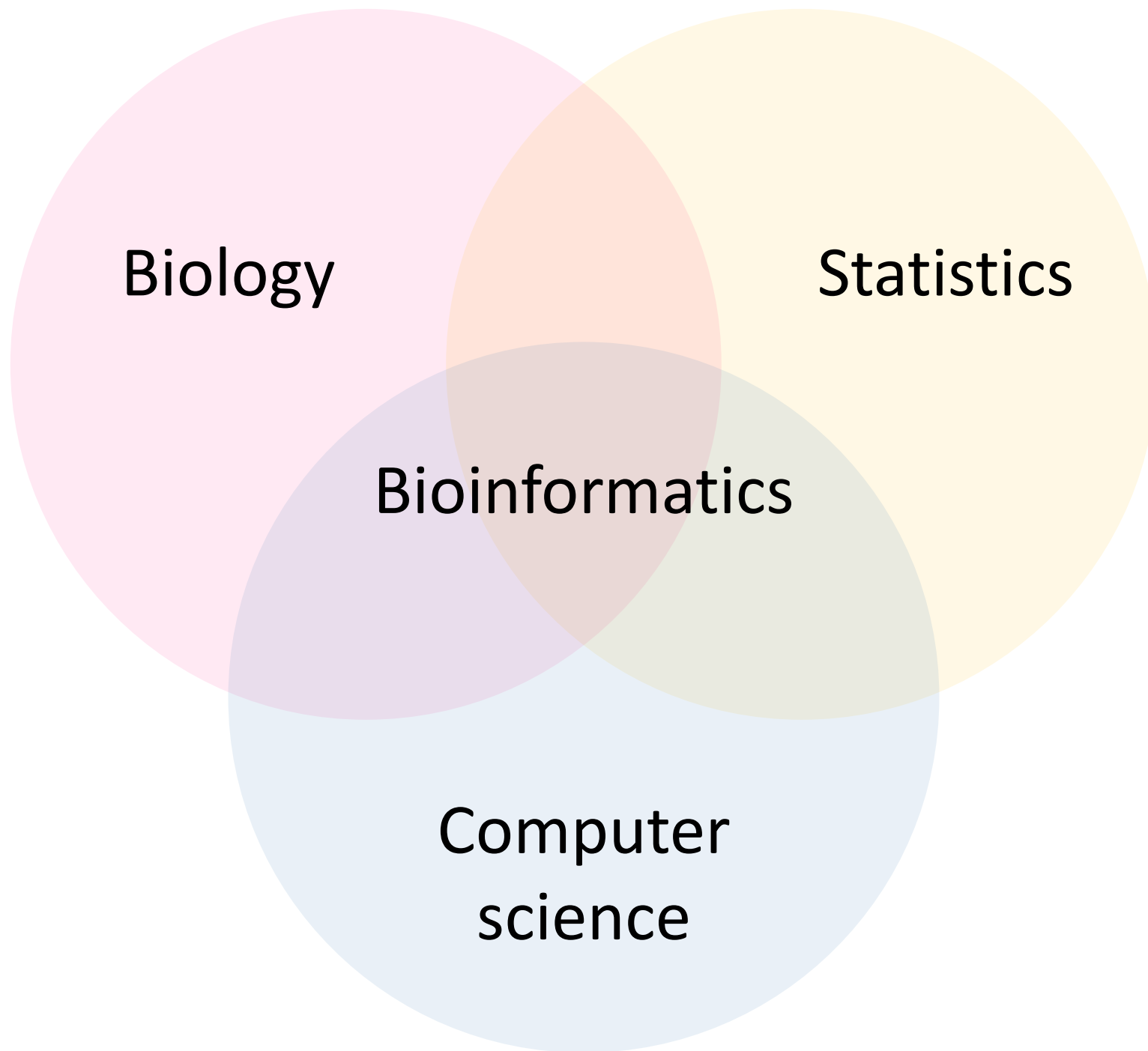


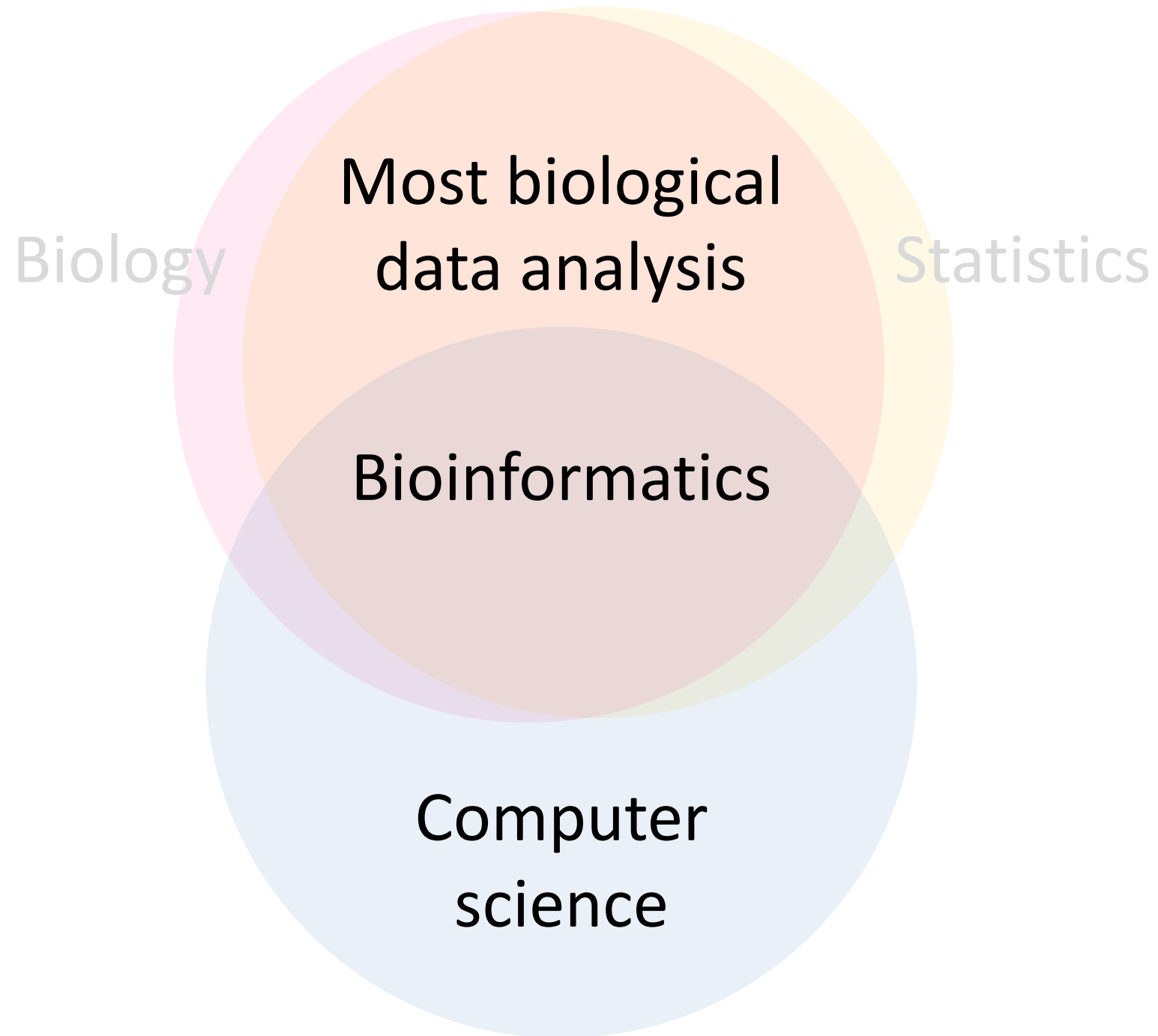
Bioinformatics for RPIP

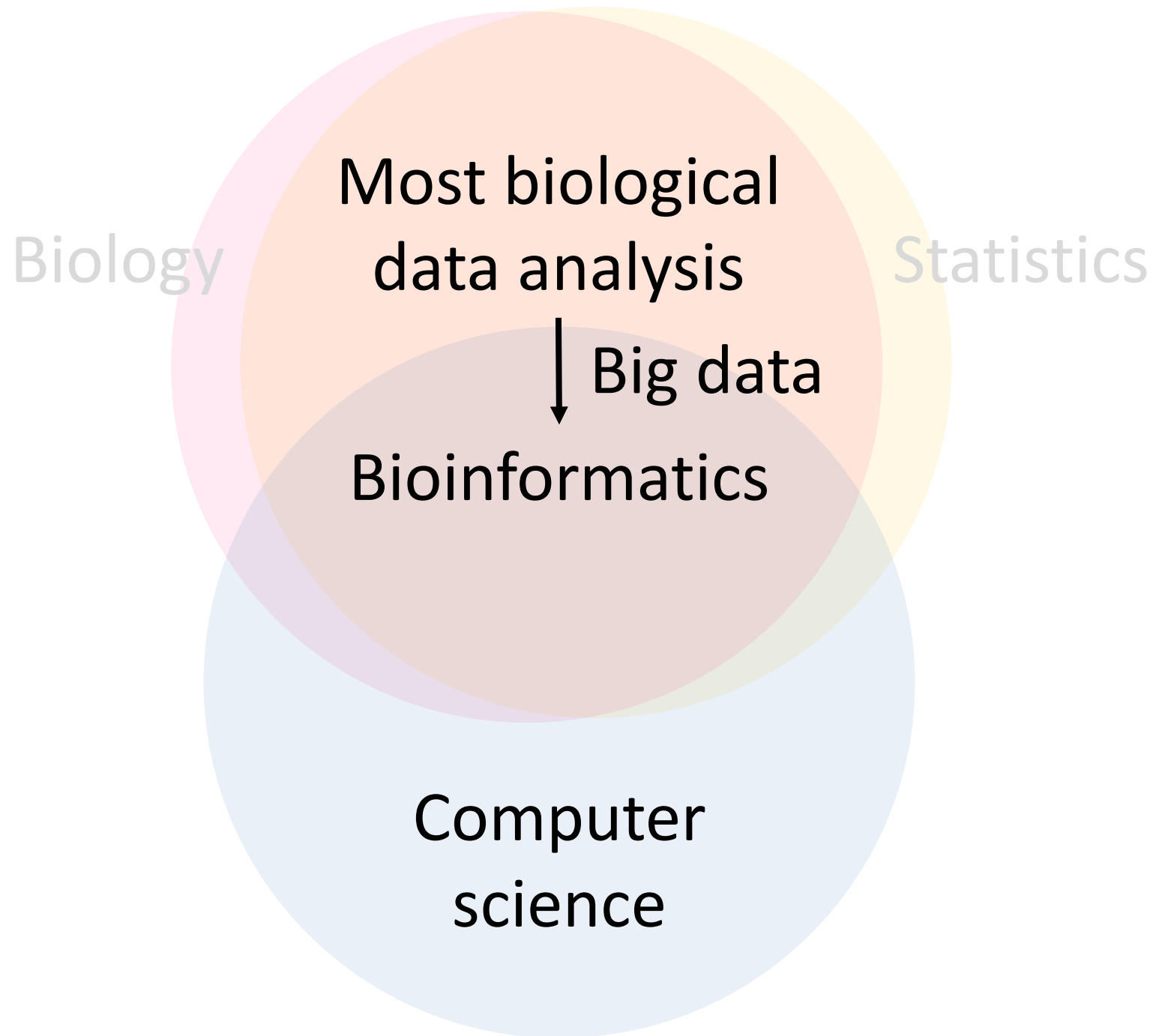
Kim Dill-McFarland

2021.05.18

kadm@uw.edu @kdillmcfarland







Big data are defined by

- Volume
 - 100s of samples
- Variety
 - 1000s of sequences per sample
- Velocity
 - 100s of new samples per week
+ a need for fast results
- Variability
 - Every sample (person) is unique

Goal

Build a bioinformatic pipeline to efficiently and reproducibly analyze Respiratory Pathogen ID/AMR Panel sequences

Specific aims

- Generate COVID19 consensus genomes for upload to public database
- Track variants of interest and variants of concern
- Track co-infections
- Create phylogenetic trees to
 - Track local variants
 - Identify local outbreaks
- Create reproducible, informative reports

Building a pipeline:
Consider the data

Respiratory Pathogen ID/AMR Panel (RPIP)

- Targeted sequencing of
 - All genes in SARS-CoV-2
 - All genes in influenza A/B
 - Identifying genes for 180 bacteria, 50 fungi, 40 viruses
 - Identifying genes for 1200 antimicrobial resistance markers
- Sequences from MiniSeq
 - 7 million (mid)
 - 20 million (rapid)
 - 22 million (high)

Raw data

.fastq

@SRR001666.1

GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCAC

+SRR001666.1

IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC

@Unique ID

Sequence

+Unique ID

Quality scores

.fasta

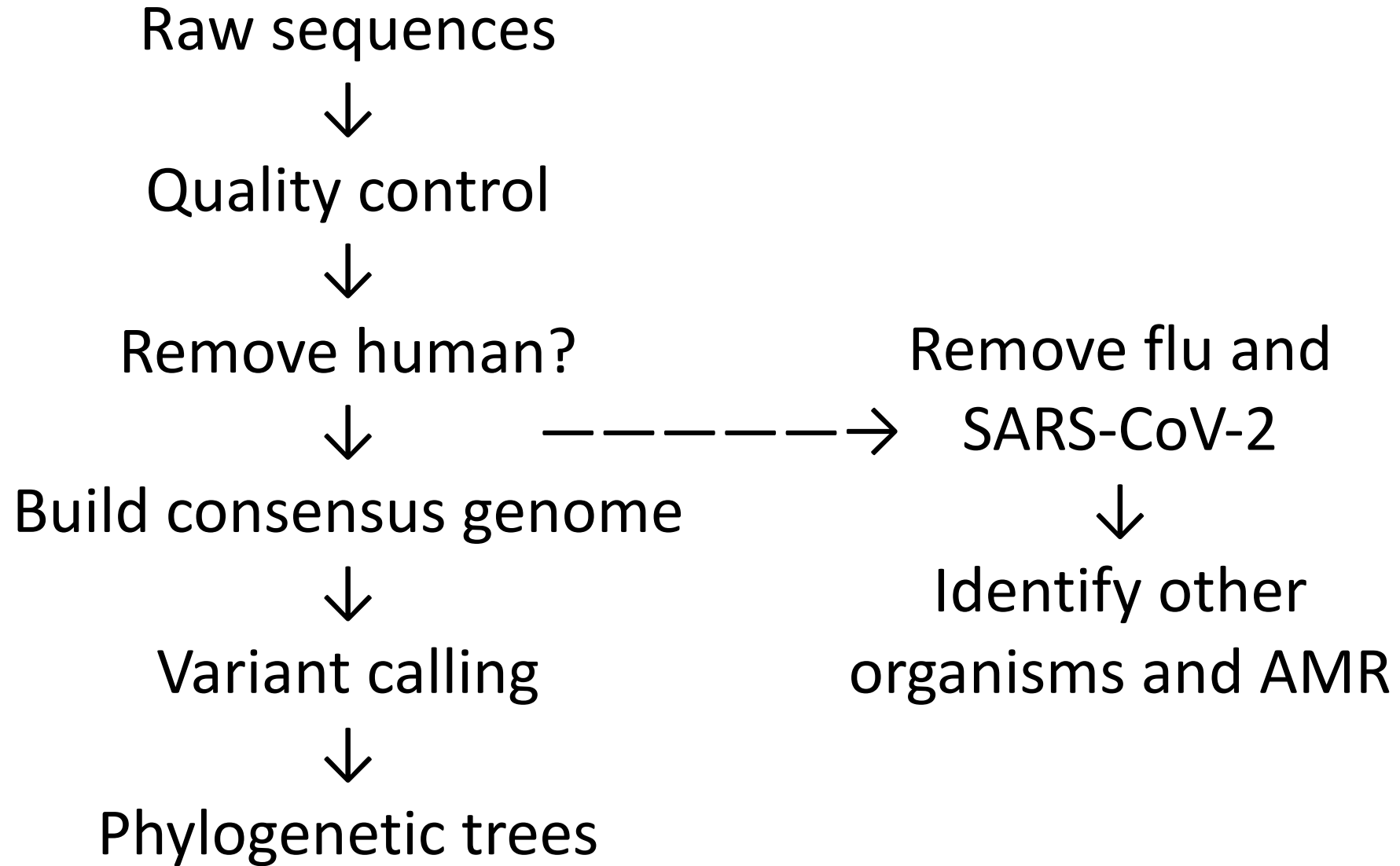
@SRR001666.1

GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCAC

@Unique ID

Sequence

Building a pipeline:
Outline a workflow



Building a pipeline:
Research software

Software considerations

0. Usefulness

- Does it do what you need it to do?

1. Documentation

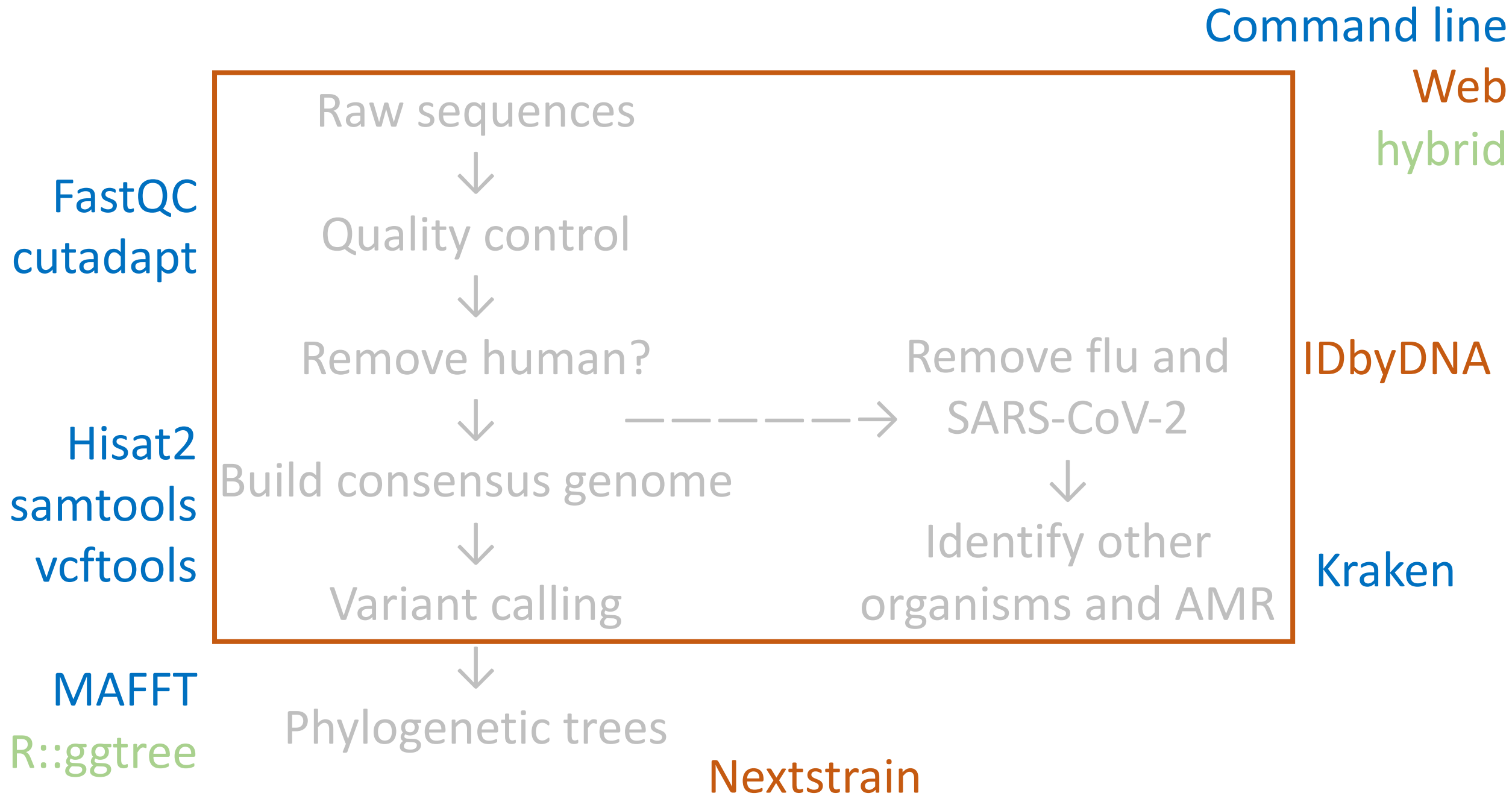
- Manuals, tutorials, GitHub
- Open-source or black box

2. Cost

- Ownership, annual membership, cost-per-use

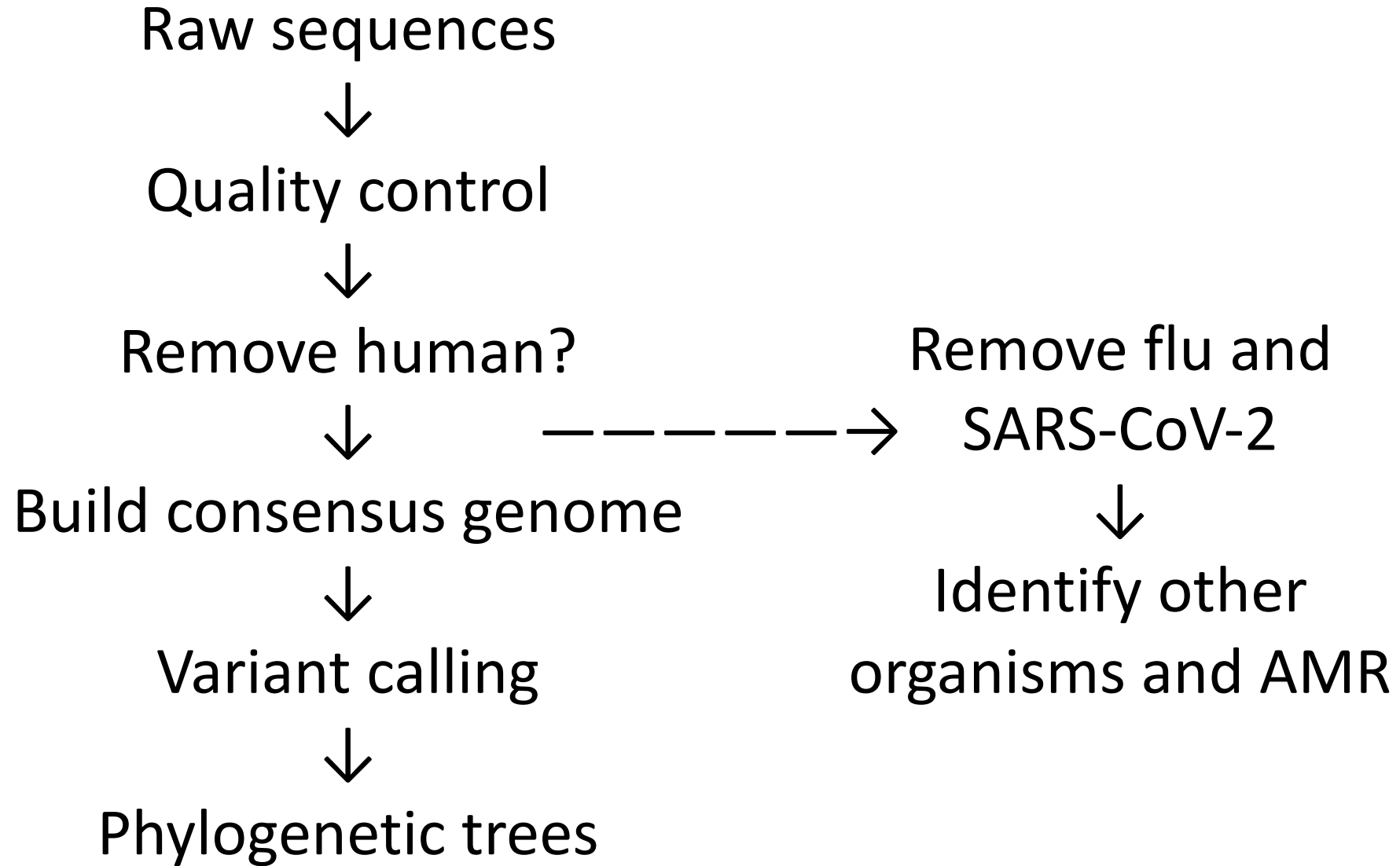
3. Platform

- Local (Windows, Mac, Linux), cloud, web
- Command line, graphical user interface (GUI), hybrid

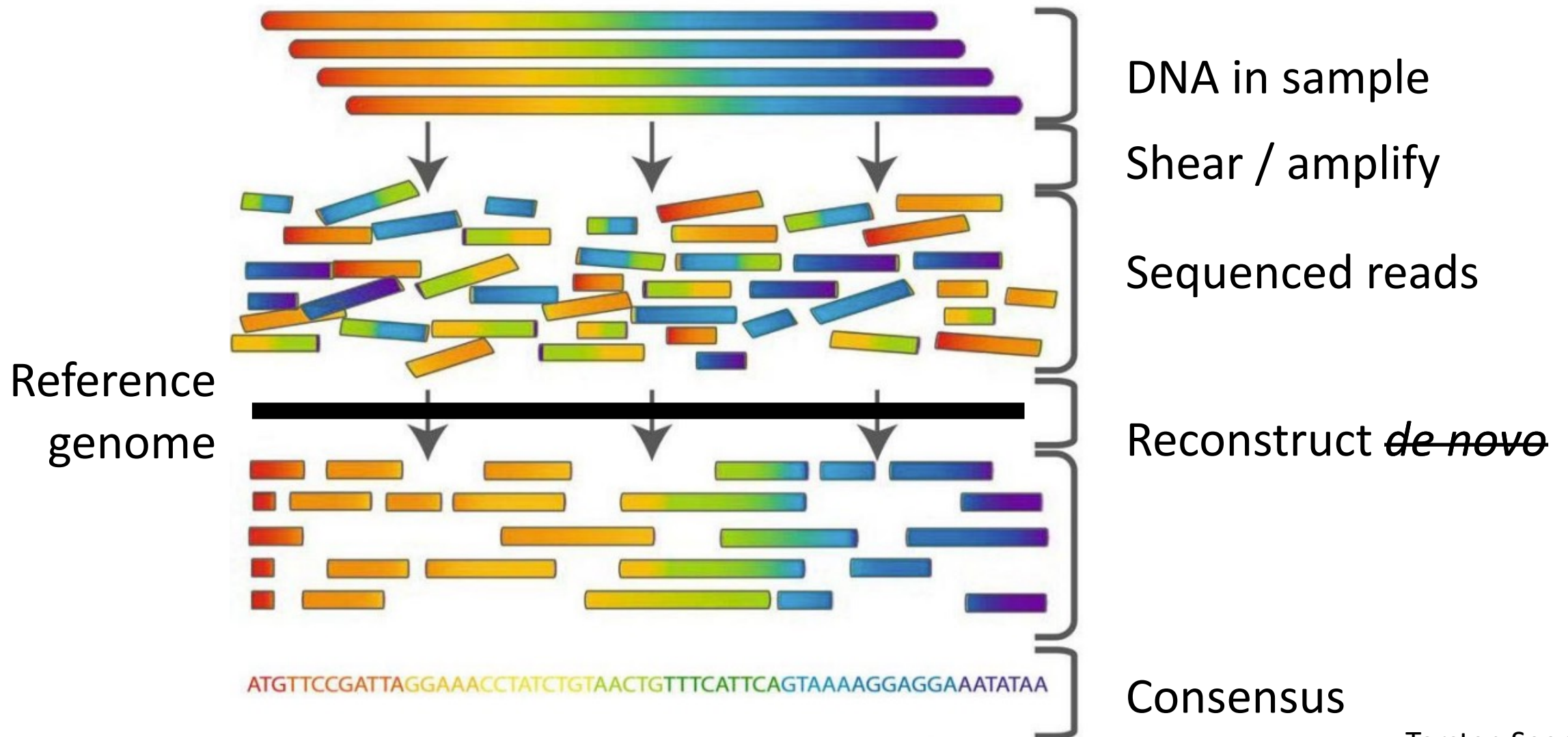


Specific aims

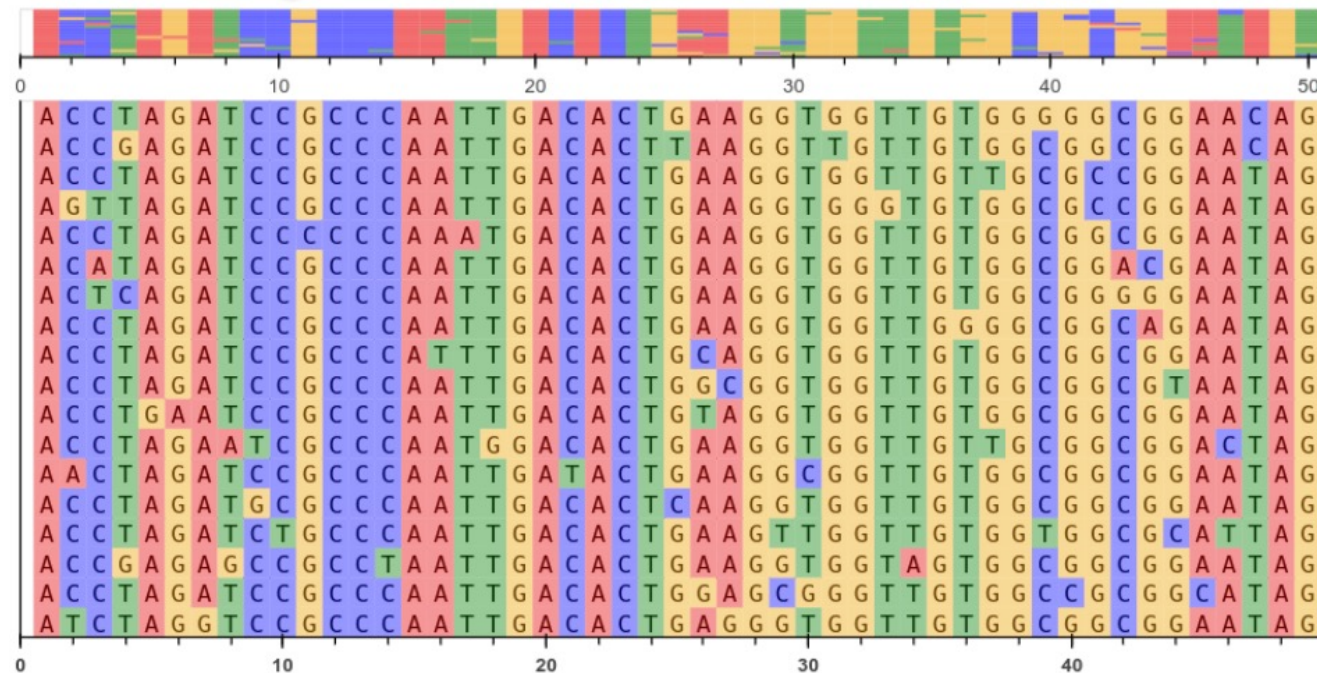
- Generate COVID19 consensus genomes for upload to public database
- Track variants of interest and variants of concern
 - How often does it update the global database?
- Track co-infections
- Create phylogenetic trees to
 - Track local variants
 - Identify local outbreaks
- Create reproducible, informative reports
 - How will software updates impact results?



Building a consensus genome



Errors in consensus



- Real or sequencing error? → PhiX control DNA
- If it's real...
 - IUPAC codes for uncertainty → R = A/G, N = ACTG, etc
 - Mixed population

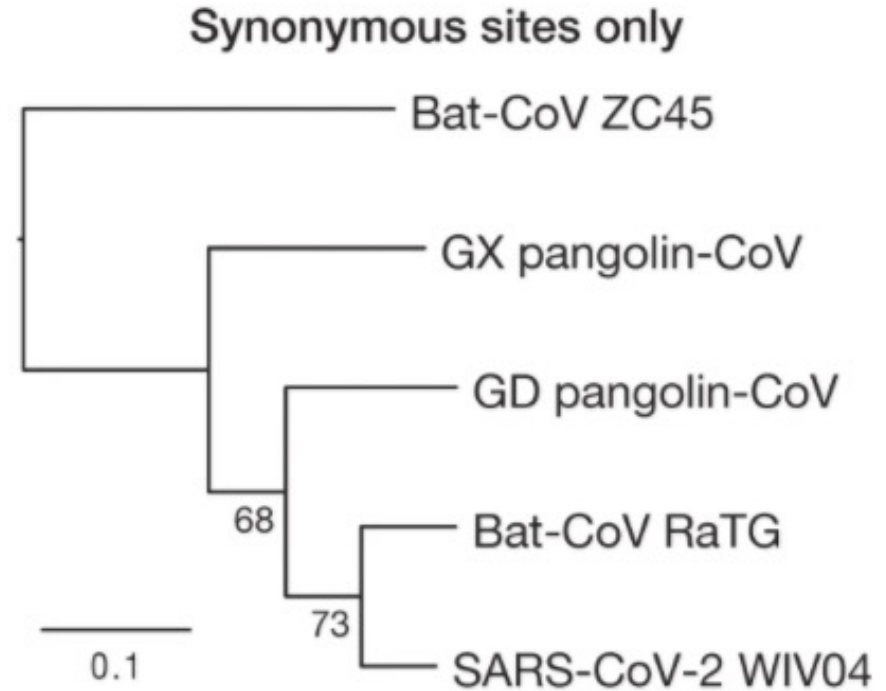
Variant calling

- Compare to database with variants of interest/concern
- Usually translated alignment

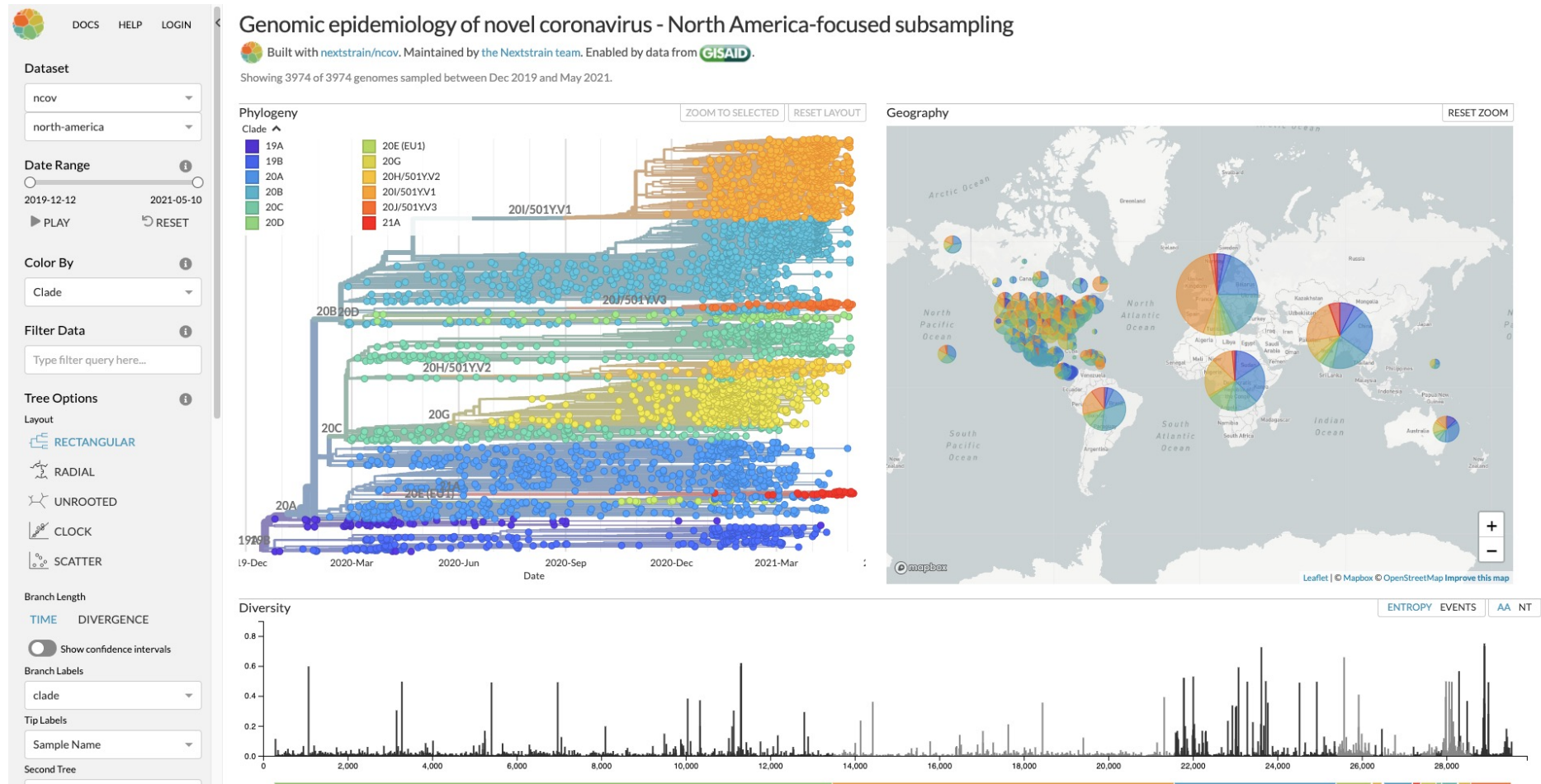
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|-----------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|-----|
| SARS-CoV | N | I | T | N | L | C | P | F | G | E | V | F | N | A | T | K | F | P | S | V | Y | A | W | E | R | K | K | I | S | N | C | V | A | D | Y | S | V | L | Y | N | S | T | F | F | S | T | F | K | C | Y | 367 |
| SARS-CoV-2 | N | I | T | N | L | C | P | F | G | E | V | F | N | A | T | R | F | A | S | V | Y | A | W | N | R | K | R | I | S | N | C | V | A | D | Y | S | V | L | Y | N | S | A | S | F | S | T | F | K | C | Y | 380 |
| GD pangolin-CoV | N | I | T | N | L | C | P | F | G | E | V | F | N | A | T | T | F | A | S | V | Y | A | W | N | R | K | R | I | S | N | C | V | A | D | Y | S | V | L | Y | N | S | T | S | F | S | T | F | K | C | Y | 380 |
| Bat-CoV RaTG | N | I | T | N | L | C | P | F | G | E | V | F | N | A | T | T | F | A | S | V | Y | A | W | N | R | K | R | I | S | N | C | V | A | D | Y | S | V | L | Y | N | S | T | S | F | S | T | F | K | C | Y | 380 |
| GX pangolin-CoV | N | I | T | N | L | C | P | F | G | E | V | F | N | A | S | K | F | A | S | V | Y | A | W | N | R | K | R | I | S | N | C | V | A | D | Y | S | V | L | Y | N | S | T | S | F | S | T | F | K | C | Y | 380 |
| Bat-CoV ZC45 | N | I | T | N | V | C | P | F | H | K | V | F | N | A | T | R | F | P | S | V | Y | A | W | E | R | T | K | I | S | D | C | I | A | D | Y | T | V | F | Y | N | S | T | S | F | S | T | F | K | C | Y | 376 |

Phylogenetic tree

- Based on translated alignment
- Groups more similar “species”



NextStrain



<https://nextstrain.org/ncov/north-america>

Building a pipeline:
Putting it all together

Pipeline wrappers

- Custom executable script
 - Short input in command line like `my_pipeline.sh data.fastq`
 - Download data + databases and run locally
- Web-based workflows (Terra Bio, Shiny apps)
 - Run executable script in the cloud
 - Data + databases can also be stored in the cloud
- Hybrid with IDbyDNA
 - Most steps on web
 - Download consensus sequence
 - Trees run locally, web, or cloud

Next steps

- Decide on desired pipeline format (platform, wrapper, specific software, etc)
- Build and test pipeline on pilot data

Addtl resources

- NextStrain <https://nextstrain.org/>
- Shiny apps <https://shiny.rstudio.com/gallery/>
- Terra Bio <https://terra.bio/>
 - Example workflow usage
https://www.youtube.com/watch?v=HObb_J9fPc0&t=604s