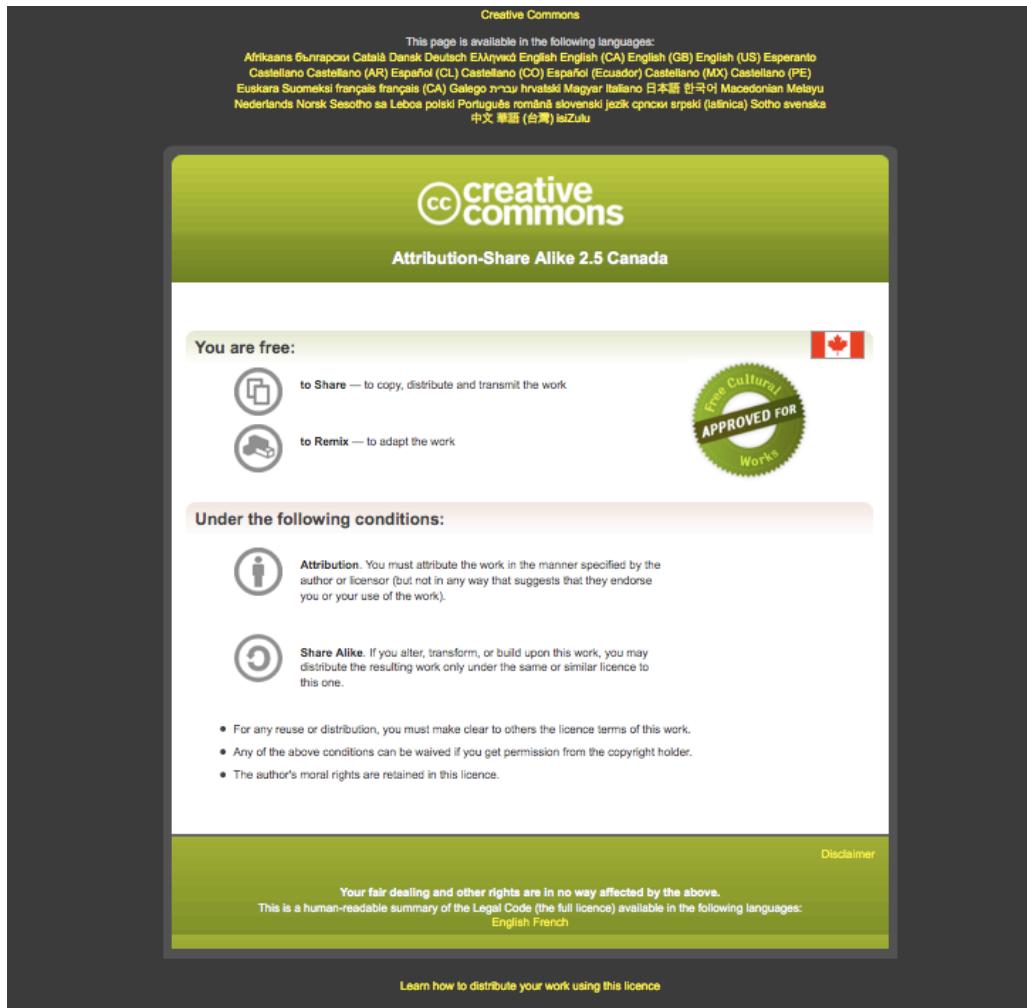




Canadian Bioinformatics Workshops

www.bioinformatics.ca

bioinformaticsdotca.github.io



Contributions

- This slide deck includes material from John Tyson (BCCDC)

Module 4: Viral Pathogen Genomic Analysis



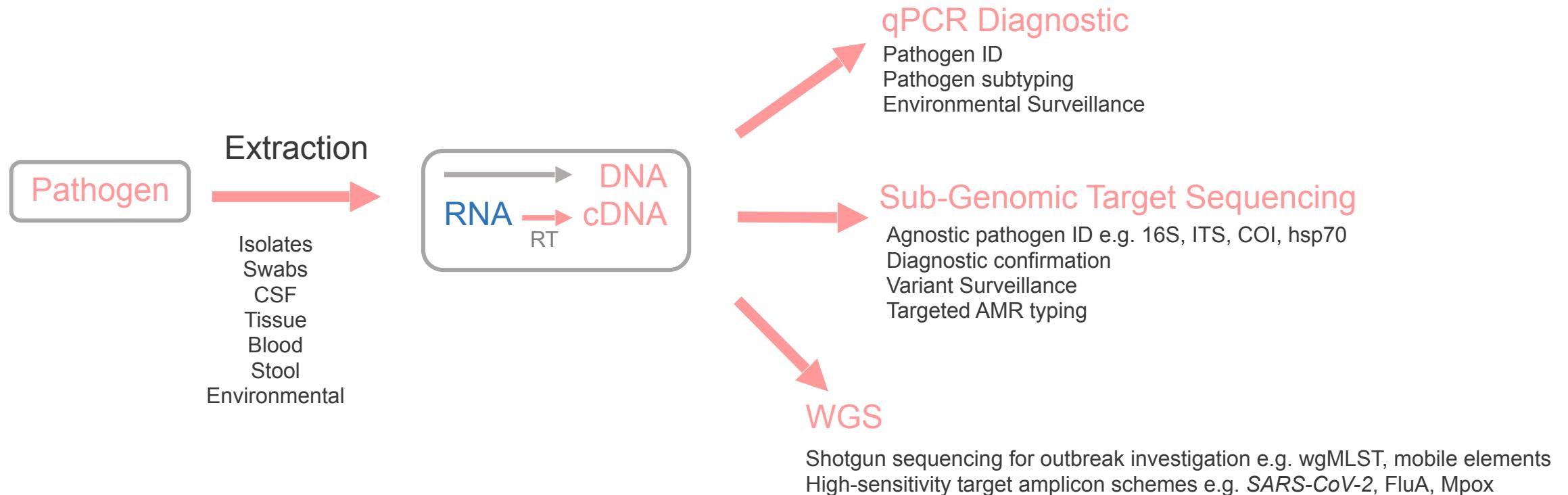
Jared Simpson
Infectious Disease Genomic Epidemiology
April 18-21, 2023



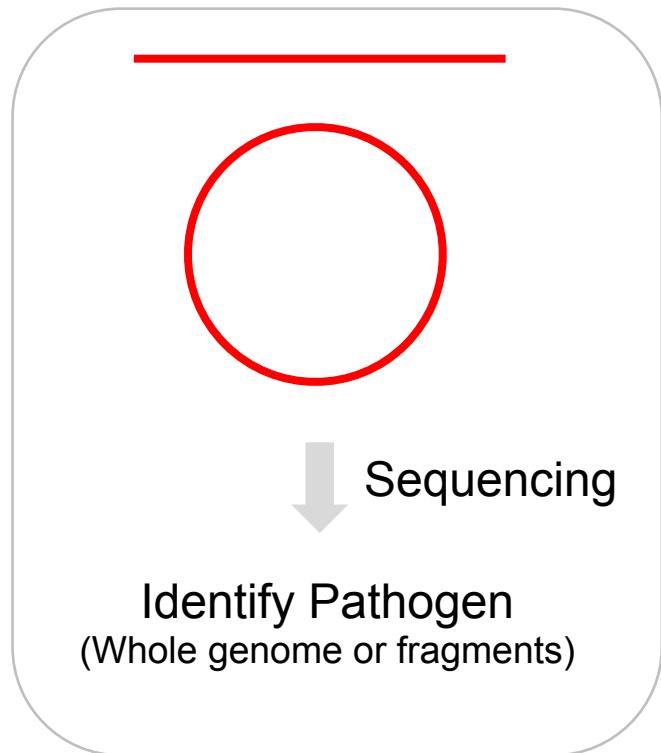
Learning Objectives

- By the end of this lecture, you will:
 - Understand the different approaches to sequence a viral genome
 - Know the key steps for analyzing amplicon data
 - Interpret the results of an amplicon-based analysis pipelines
 - Perform quality control on your data

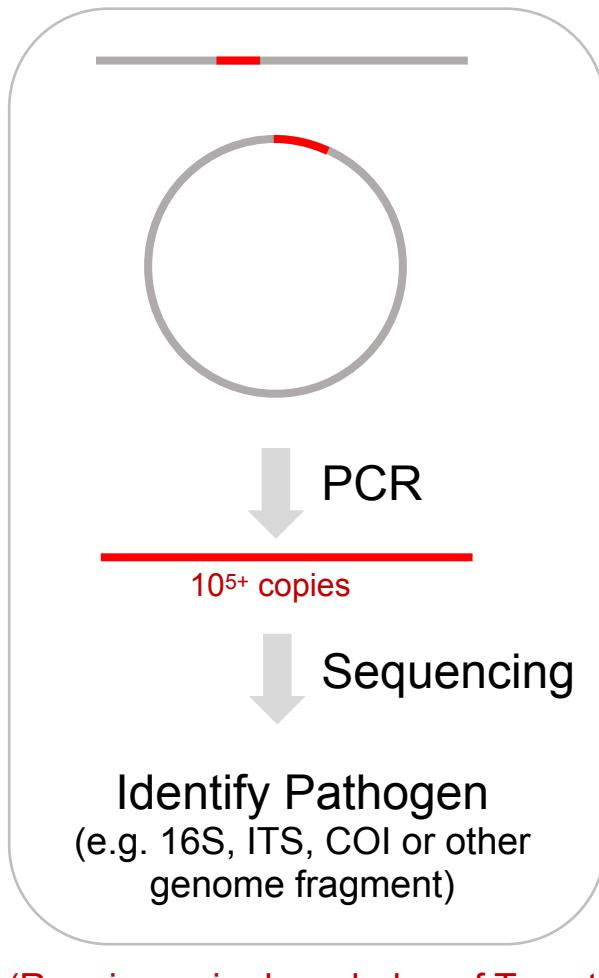
Pathogen Genetic testing & typing



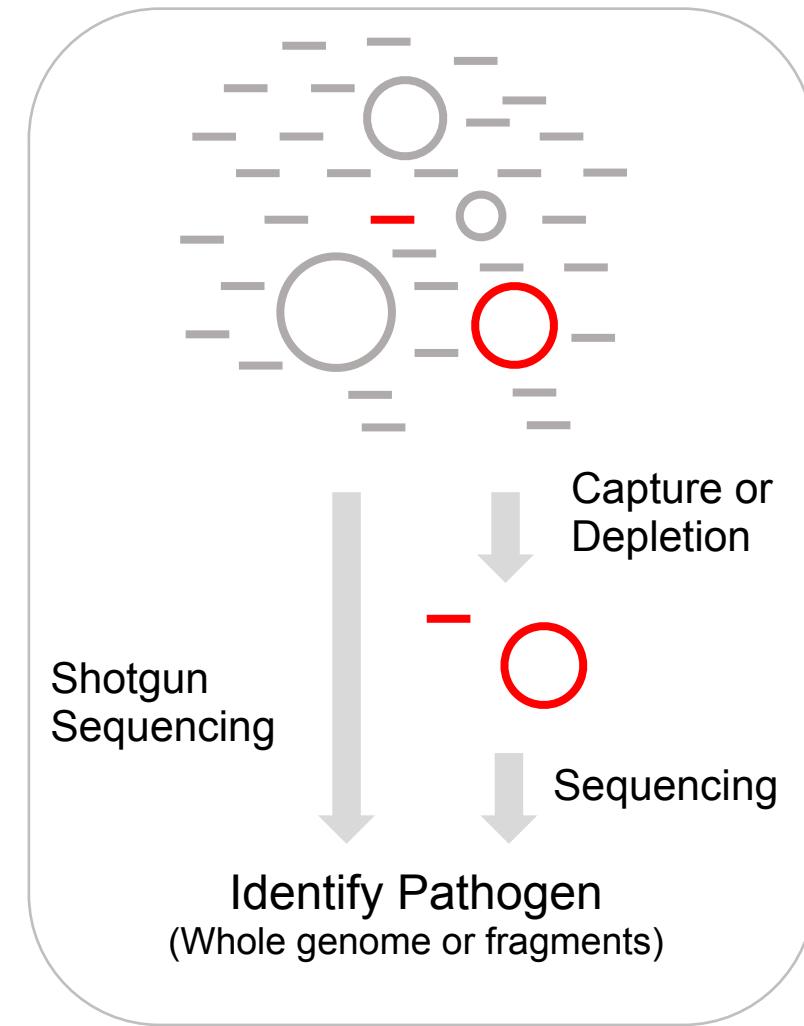
Isolate Sequencing

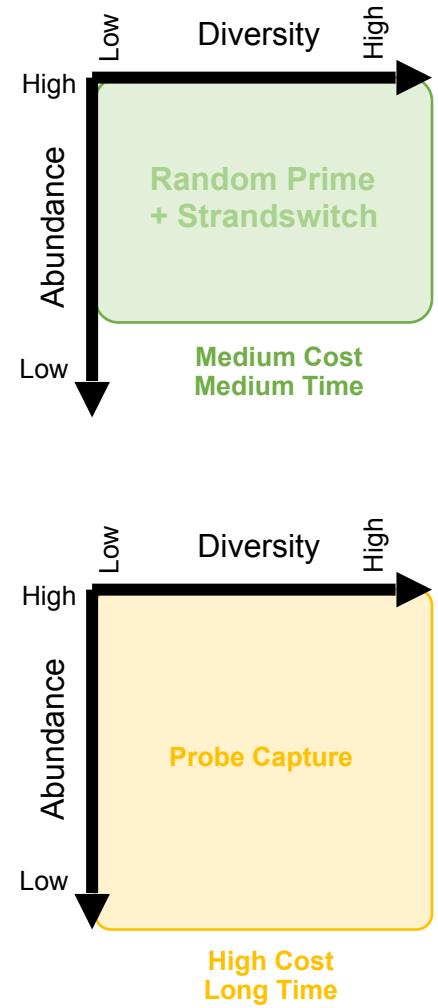
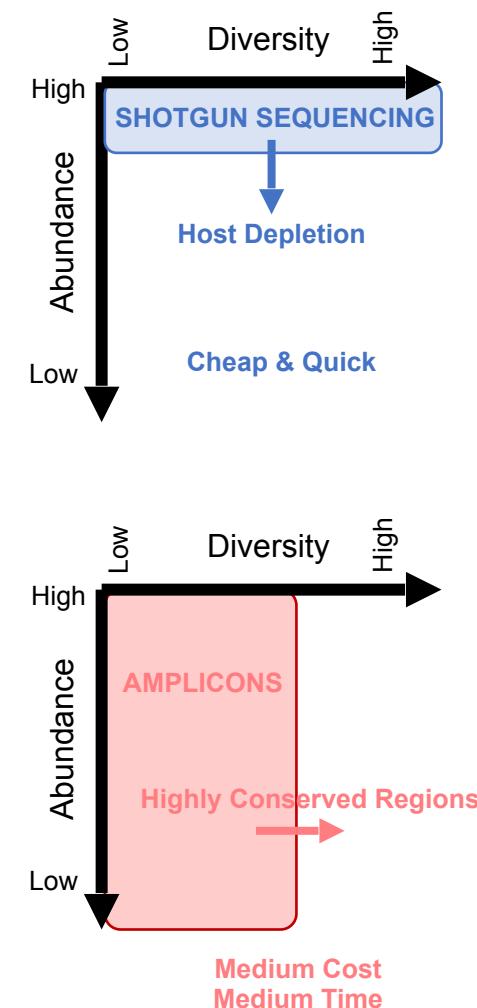
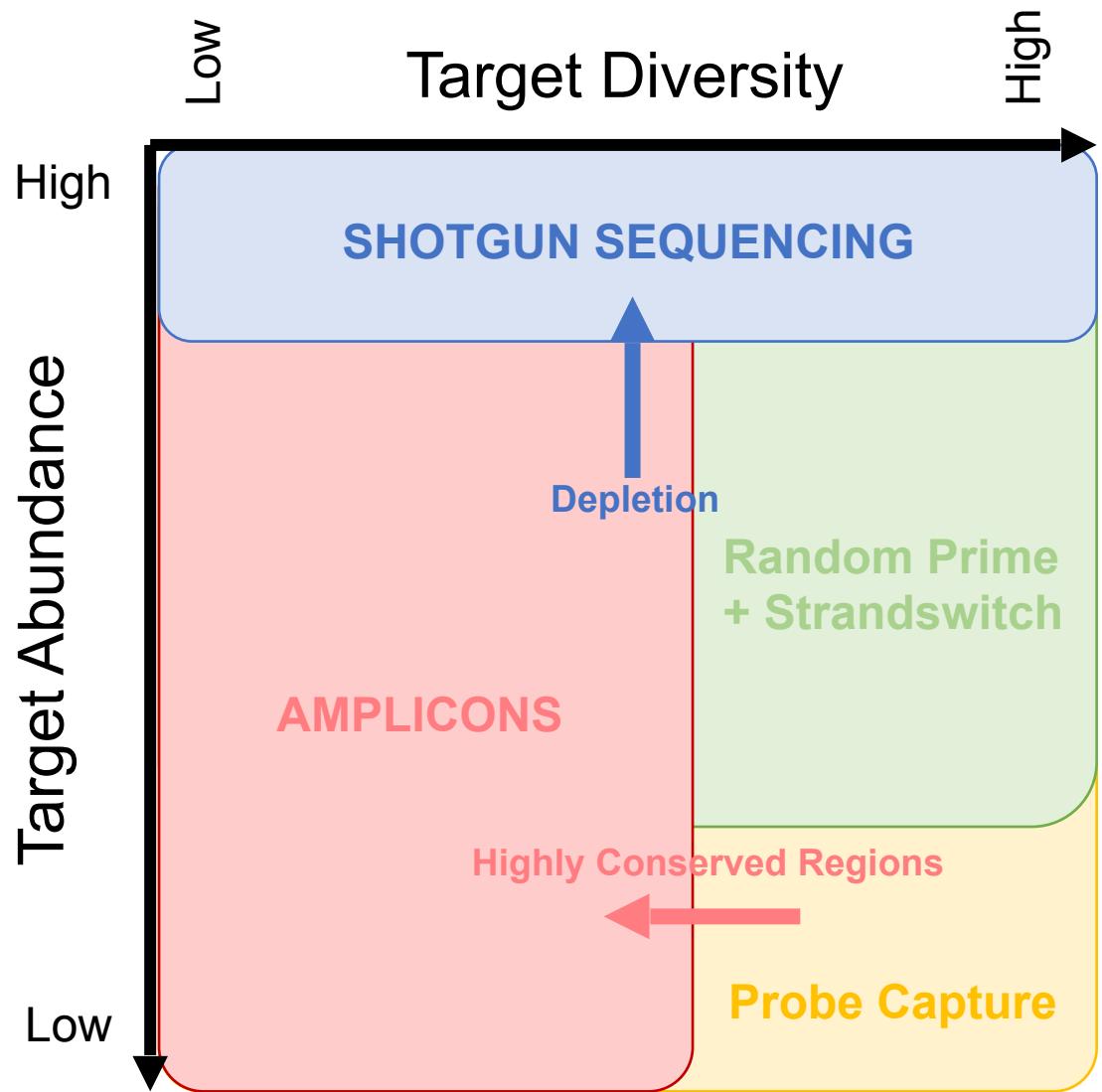


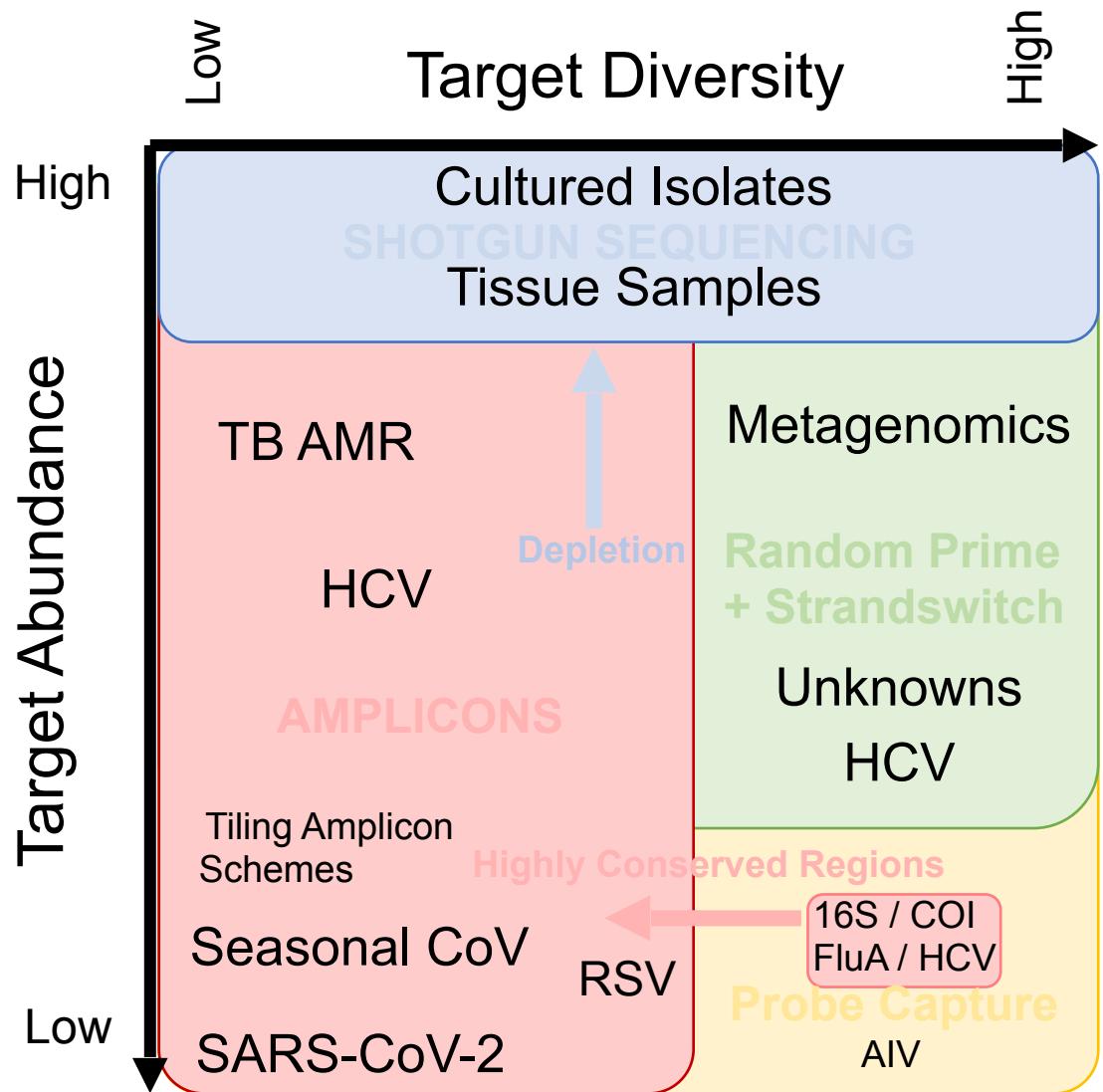
Targeted Sequencing



Complex Mixture Sequencing



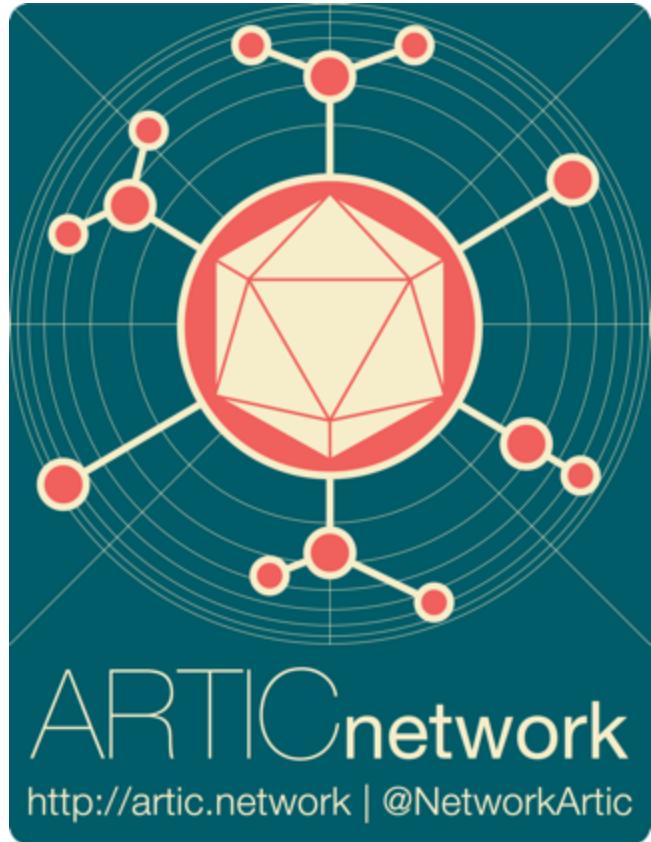
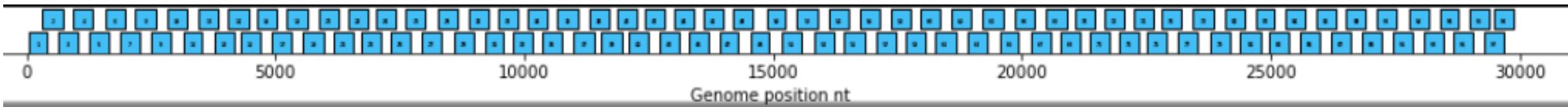




The Different Sequencing Technologies

- There are a number of “Next Generation” sequencing technologies available:
 - Illumina short read length sequencers :- small to high capacity instruments, short read lengths (<300b)
 - PacBio long read length sequencers :- Medium capacity instruments, long read lengths (<50Kb CCS, ~20Kb HiSeq consensus)
 - Oxford Nanopore long read length sequencers :- Small to large capacity, Ultra long read lengths (250 to 2Mb+)

Tiled Amplicon Sequencing



Artic Protocol (V3)

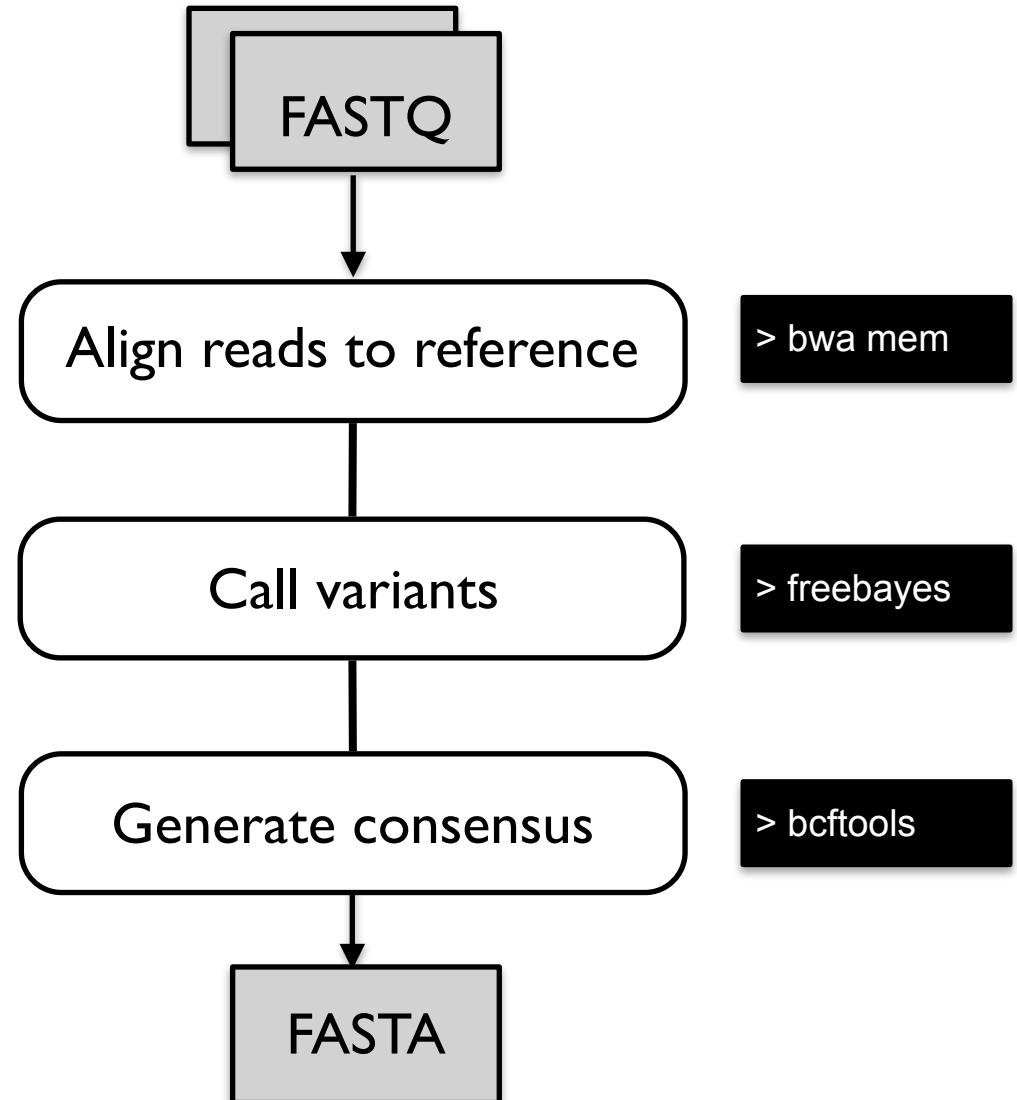
- Highly multiplexed PCR reaction
- 98 primer pairs over two primer pools
- 400 bp amplicons
- Balancing amplicons → sequencing extra deep to accommodate

Improvements to the ARTIC multiplex PCR method for SARS-CoV-2 genome sequencing using nanopore

[John R Tyson](#),¹ [Phillip James](#),² [David Stoddart](#),² [Natalie Sparks](#),³ [Arthur Wickenhagen](#),⁴ [Grant Hall](#),⁵ [Ji Hyun Choi](#),⁶ [Hope Lapointe](#),⁶ [Kimia Kamelian](#),⁷ [Andrew D Smith](#), [Natalie Prystajecky](#),^{7,8} [Ian Goodfellow](#),⁵ [Sam J Wilson](#),⁴ [Richard Harrigan](#),⁶ [Terrance P Snutch](#),¹ [Nicholas J Loman](#),³ and [Joshua Quick](#)³

Analysis Pipelines

WGS



(This pipeline is intentionally simplified)

Read Mapping and Alignment

Typically reads are much shorter (100-200bp) than the reference genome (SARS-CoV-2: ~30,000bp)

We need to find where the read *maps* to in the reference genome and how the bases in the read *align* to the reference

reference	GAGATACATGAGAGAGTATCTCGACTCTAGGCCGATACCATTGTA AGTATCTTGACTCTA
read	

Key definitions:

mapping: the region in the reference that is most similar to the read

alignment: how the read lines up to the reference base-by-base

Consensus sequences

- Compare reads to the SARS-CoV-2 reference genome

reference	GATCCATGTAGTACCATTAGTACAGTACCATATAT
	GATCCATGTAGTACCATTAGTACAGTACCATA
	ATCCATGTAGTACCATTAGTACAGTACCATATAT
	CATGTAGTACCATTAGTACAGTACCATATAT
	GTAGTACCATTAGTACAGTACCATATAT
	GTAGTACCATTAGTACAGTACCATATAT
reads	TAGTACCATTAGTACAGTACCATATAT
	TAGTACCATTAGTACAGTACCATATAT
	AGTACCATTAGTACAGTACCATATAT
	AGTACCATTAGTACAGTACCATATAT
	GTACCATTAGTACAGTACCATATAT
consensus	GATCCATGTAGTACCATTAGTACAGTACCATATAT

Consensus sequences

- *variants* or *mutations* are detected when the aligned reads differ from the reference

reference

GATCCATGTAGTACCAT**T**AGTACAGTACCATATAT

GATCCATGTAGTACCAT**C**AGTACAGTACCATATA

ATCCATGTAGTACCAT**C**AGTACAGTACCATATAT

CATGTAGTACCAT**C**AGTACAGTACCATATAT

GTAGTACCAT**C**AGTACAGTACCATATAT

GTAGTACCAT**C**AGTACAGTACCATATAT

TAGTACCAT**C**AGTACAGTACCATATAT

TAGTACCAT**C**AGTACAGTACCATATAT

AGTACCAT**C**AGTACAGTACCATATAT

AGTACCAT**C**AGTACAGTACCATATAT

GTACCAT**C**AGTACAGTACCATATAT

GATCCATGTAGTACCATC**AGTACAGTACCATATAT**

reads

consensus

Consensus sequences

- *ambiguous or mixed bases* are detected when the aligned reads have evidence for more than one type of base reference

reads

GATCCATGTAGTACCAT**T**AGTACAGTACCATAT
 GATCCATGTAGTACCAT**T**AGTACAGTACCATA
 ATCCATGTAGTACCAT**C**AGTACAGTACCATAT
 CATGTAGTACCAT**C**AGTACAGTACCATAT
 GTAGTACCAT**C**AGTACAGTACCATAT
 GTAGTACCAT**T**AGTACAGTACCATAT
 TAGTACCAT**C**AGTACAGTACCATAT
 TAGTACCAT**C**AGTACAGTACCATAT
 AGTACCAT**T**AGTACAGTACCATAT
 AGTACCAT**T**AGTACAGTACCATAT
 GTACCAT**C**AGTACAGTACCATAT

consensus

GATCCATGTAGTACCATY**AGTACAGTACCATAT**



IUPAC ambiguity code

Consensus sequences

- Can you think of explanations for *why* this can happen?

reference

GATCCATGTAGTACCAT**T**AGTACAGTACCATATAT

GATCCATGTAGTACCAT**T**AGTACAGTACCATATA

ATCCATGTAGTACCAT**C**AGTACAGTACCATATAT

CATGTAGTACCAT**C**AGTACAGTACCATATAT

GTAGTACCAT**C**AGTACAGTACCATATAT

GTAGTACCAT**T**AGTACAGTACCATATAT

TAGTACCAT**C**AGTACAGTACCATATAT

TAGTACCAT**C**AGTACAGTACCATATAT

AGTACCAT**T**AGTACAGTACCATATAT

AGTACCAT**T**AGTACAGTACCATATAT

GTACCAT**C**AGTACAGTACCATATAT

GATCCATGTAGTACCATY**AGTACAGTACCATATAT**

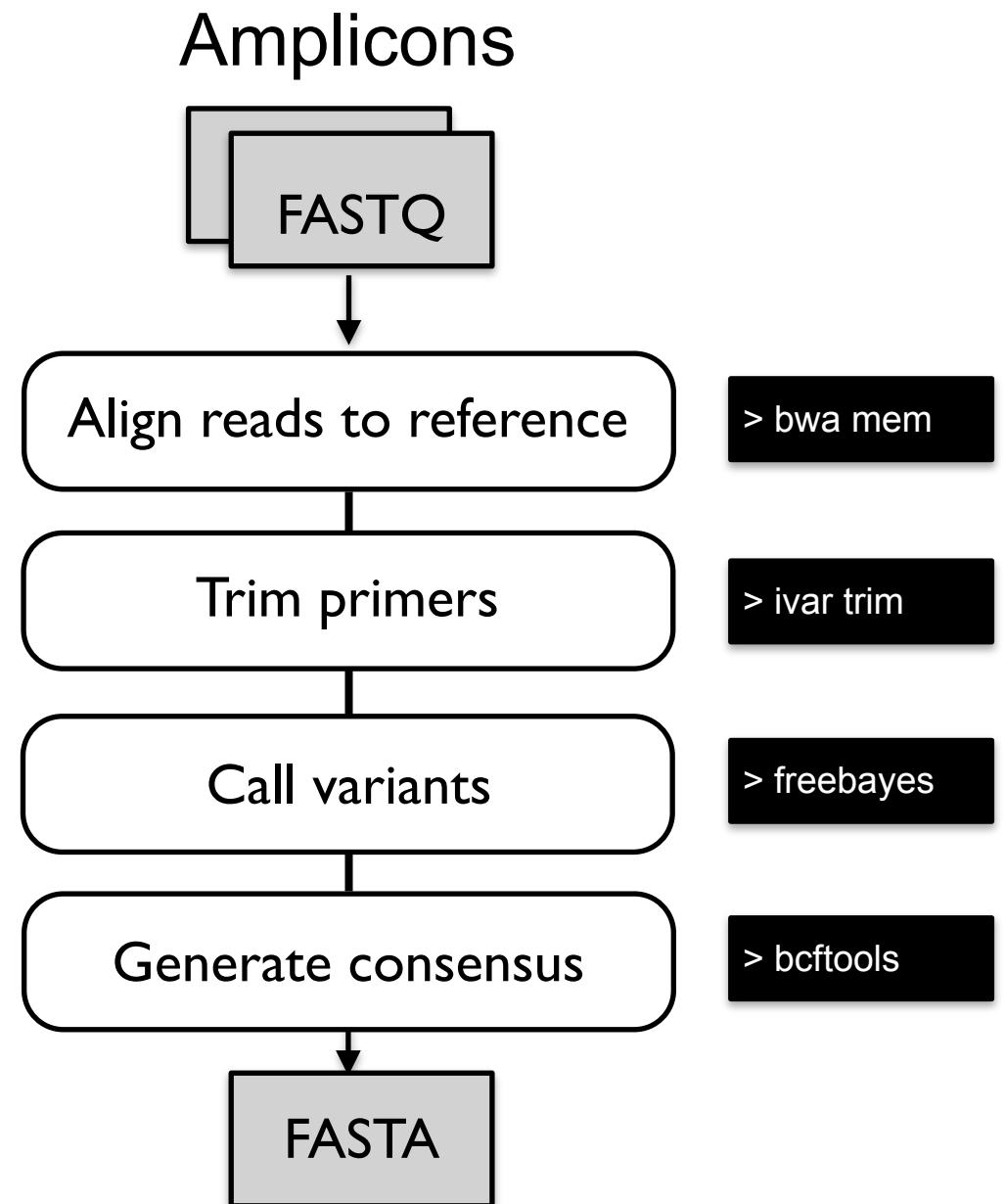
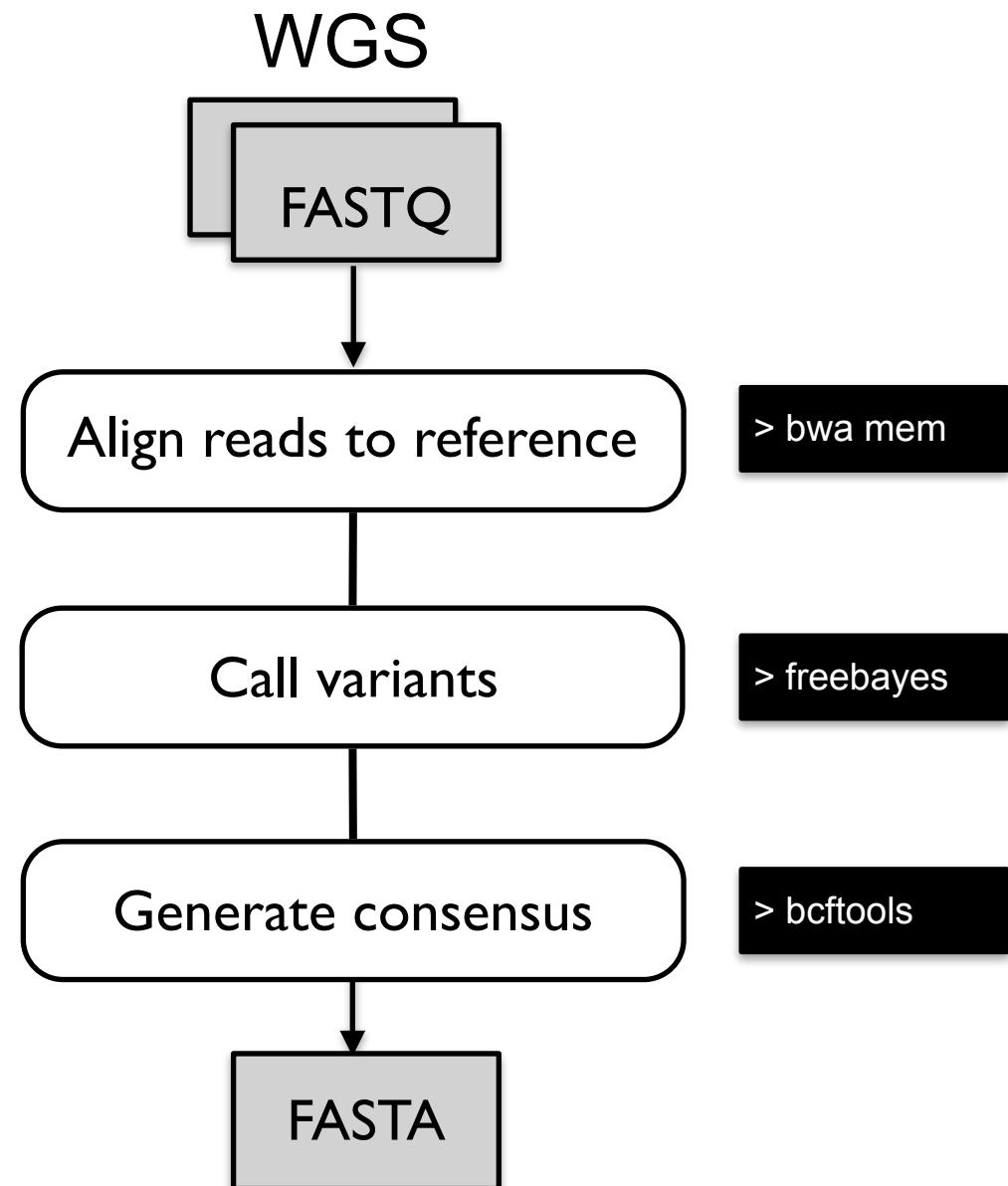
reads

consensus



IUPAC ambiguity code

Analysis Pipelines



Why do we need to trim primers?

Example Illumina read:

```
ERR5338522.5
ACCAACCAACTTCGATCTCTGTAGATCTGTTCTAAACGAACTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCACTCACGCAGTATAATTAACTA
+
FFFFF:FFF:FFFFF:FFFFFFFFFFFFFFFF,FFF,,FFFFFFFFFF:FFFF:FFFF:FFFFFFFFFF:FFFFFFFFFF:FFFF:FFFF
```

Why do we need to trim primers?

Example Illumina read:

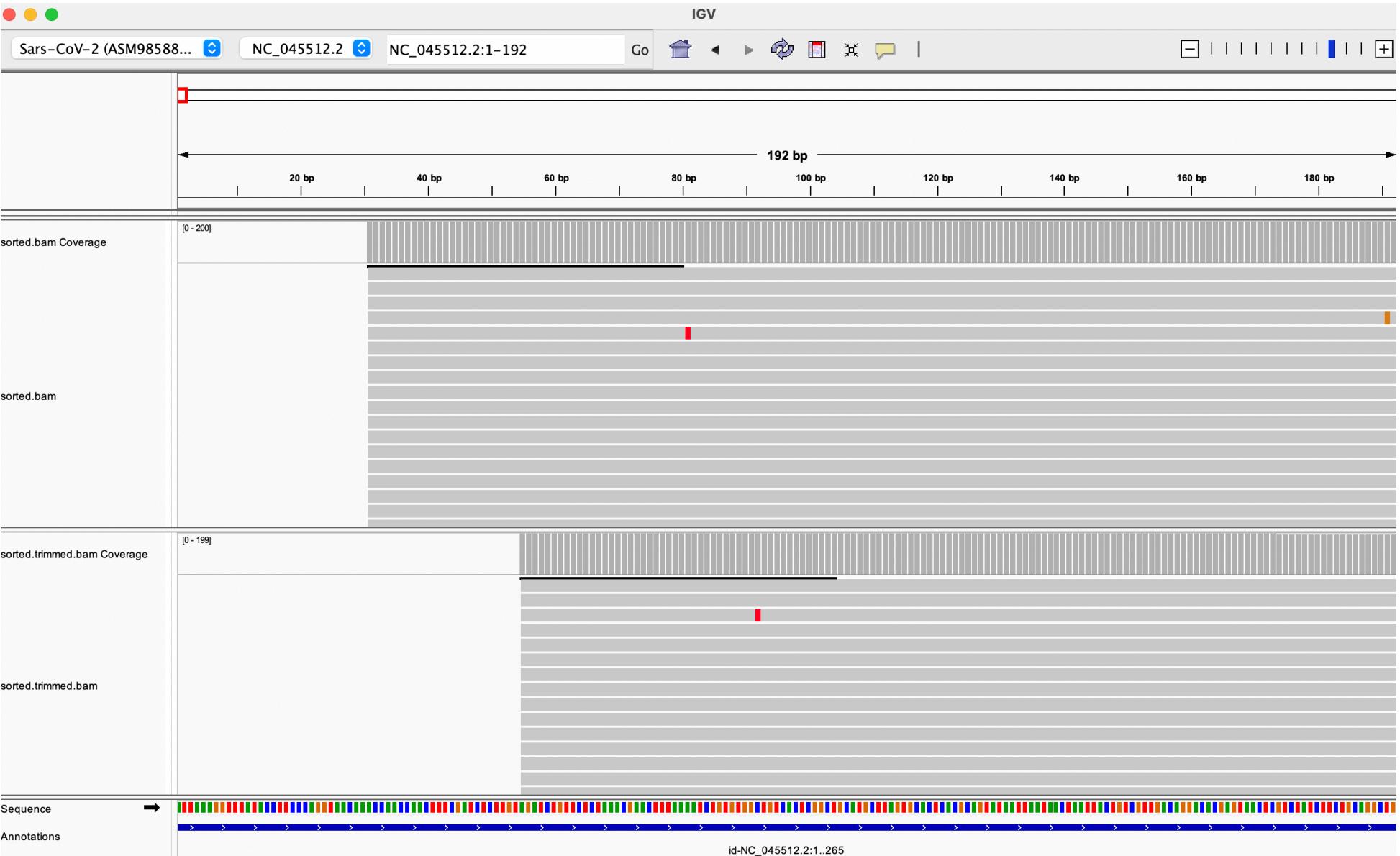
```
ERR5338522.5
ACCAACCAACTTCGATCTCTGTAGATCTGTTCTAAACGAACTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCACTCACGCAGTATAATTAACTA
+
FFFFF:FFF:FFFFF:FFFFFFFFFFFFFFFF, FFF,, FFFFFFFFFFFF:FFFFFFF:FFFFFFFFFFFF:FFFFFFF:FFFF:FFFF:FFFF:FFFF:FFFF
```

ARTIC V3 primer highlighted in red

Primer sequence is always identical to reference and cannot be used to detect mutations

Unless removed from the reads we will over-call reference alleles

Why do we need to trim primers?



How do we quality check our results?

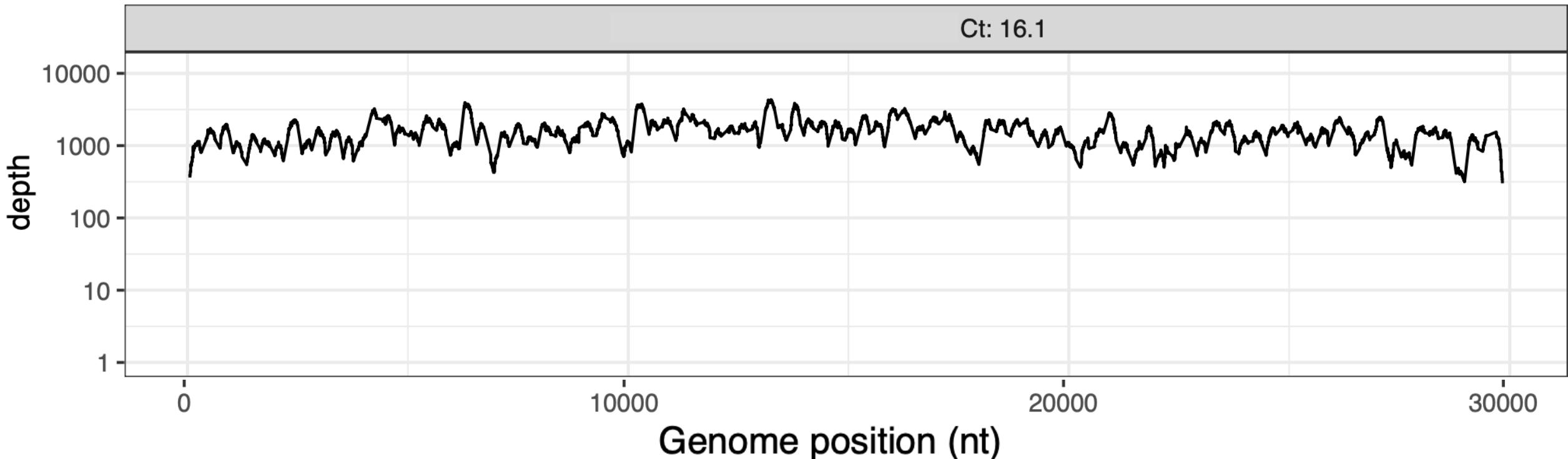
- Was the genome successfully sequenced?
 - Virus abundance varies sample-to-sample, storage conditions and extraction methods may differ in ability to recover complete genomes
- Is the genome sequence accurate?
 - Low abundance samples may have amplification artifacts, low coverage may lead to consensus errors
- Is the sequencing run contaminated?
 - Sequencing protocols are highly multiplexed, many amplification cycles increases risk of cross-contamination on the plate

What do we mean by quality?

- **Was the genome successfully sequenced?**
 - Virus abundance varies sample-to-sample, storage conditions and extraction methods may differ in ability to recover complete genomes
- Is the genome sequence accurate?
 - Low abundance samples may have amplification artifacts, low coverage may lead to consensus errors
- Is the sequencing run contaminated?
 - Sequencing protocols are highly multiplexed, many amplification cycles increases risk of cross-contamination on the plate

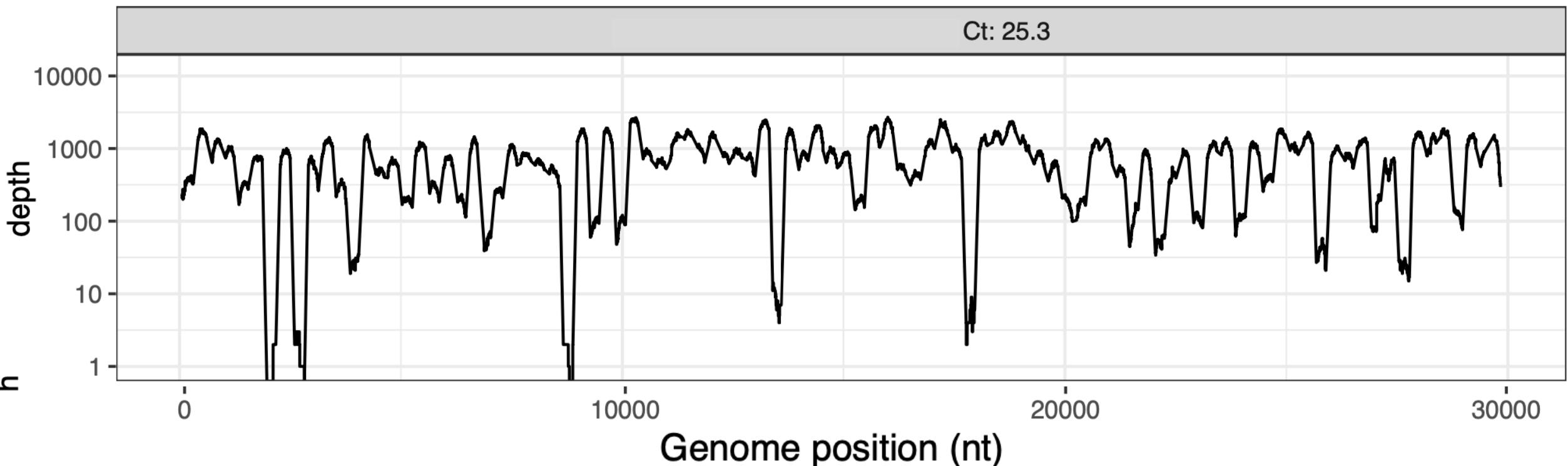
Genome Coverage

Low Ct sample, ideal coverage profile



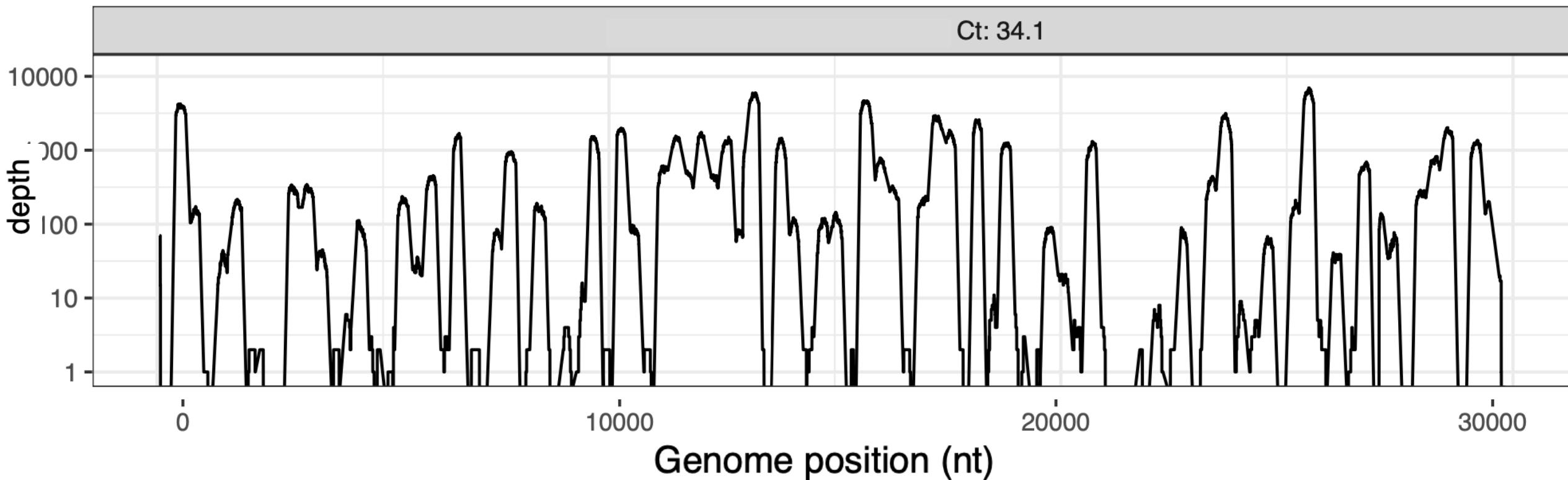
Genome Coverage

Moderate Ct sample, some dropouts



Genome Coverage

High Ct sample, many dropouts



Genome Completeness

- The bioinformatics pipelines take the mapped reads and output a consensus sequence
 - Positions with insufficient coverage are masked with “N” bases
- Genome completeness is defined as the proportion of non-N bases

```
>sampleA  
ACGGNNNNACA
```

Genome completeness = 7/10 = 0.7

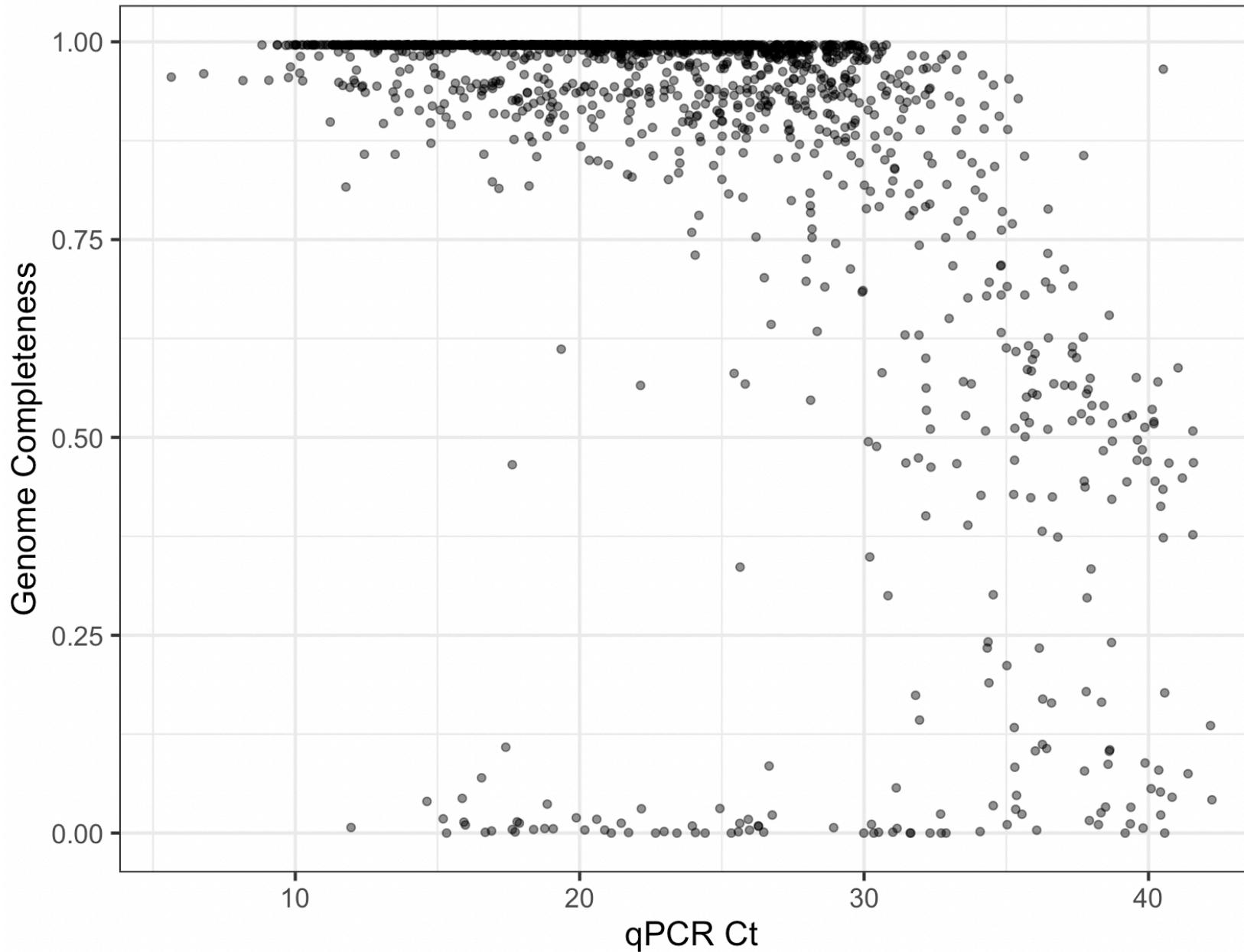
```
>sampleB  
ACGGANTACA
```

Genome completeness = 9/10 = 0.9

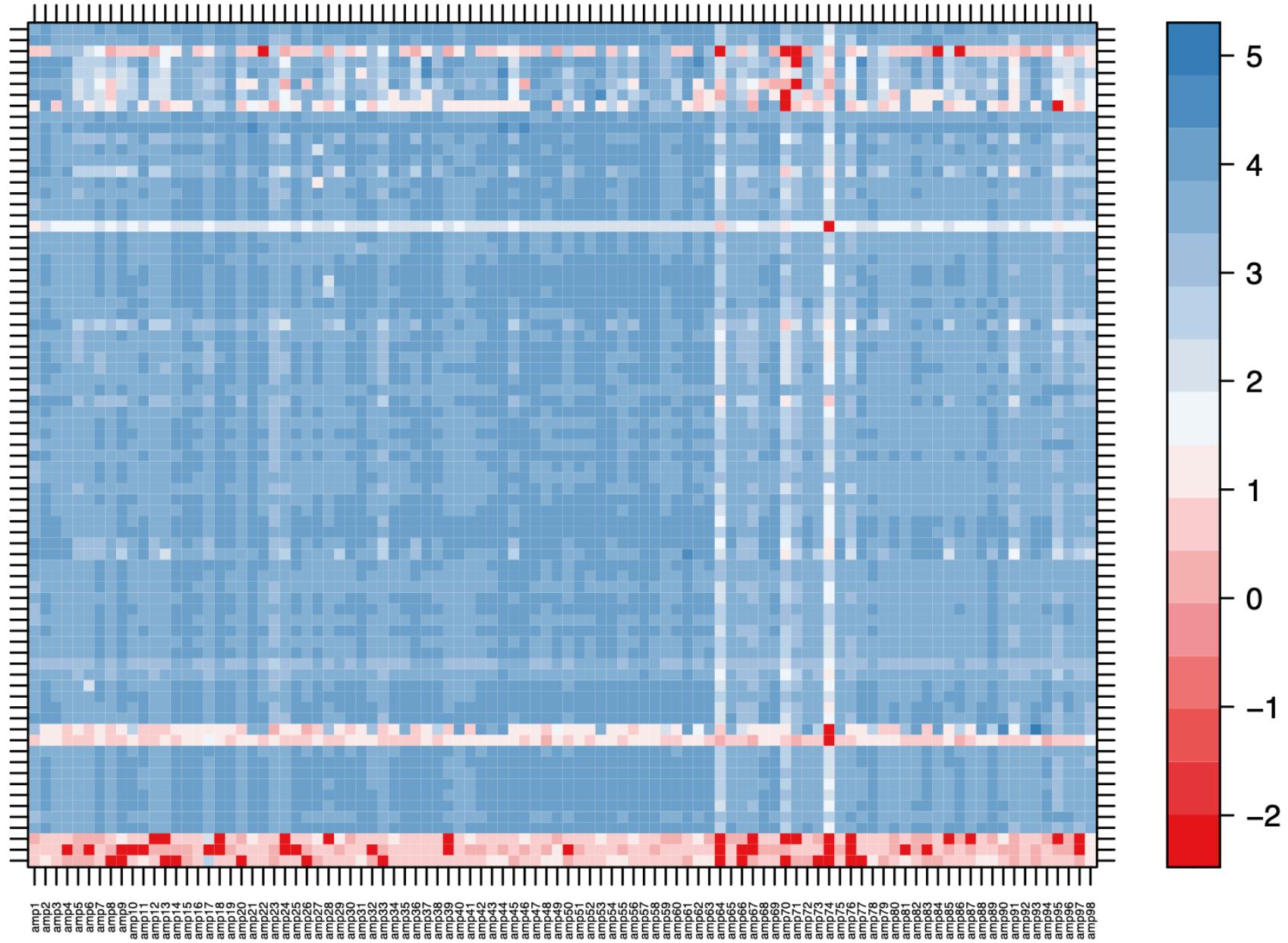
Genome Completeness

- Incomplete genomes are harder to analyze (e.g. may cause problems for multiple sequence alignment or building a tree)
- Many projects use genome completeness as a primary criteria for submission to public repositories
 - For CanCOGeN genomes with completeness >90% are submitted to GISAID

Relationship between Ct and completeness



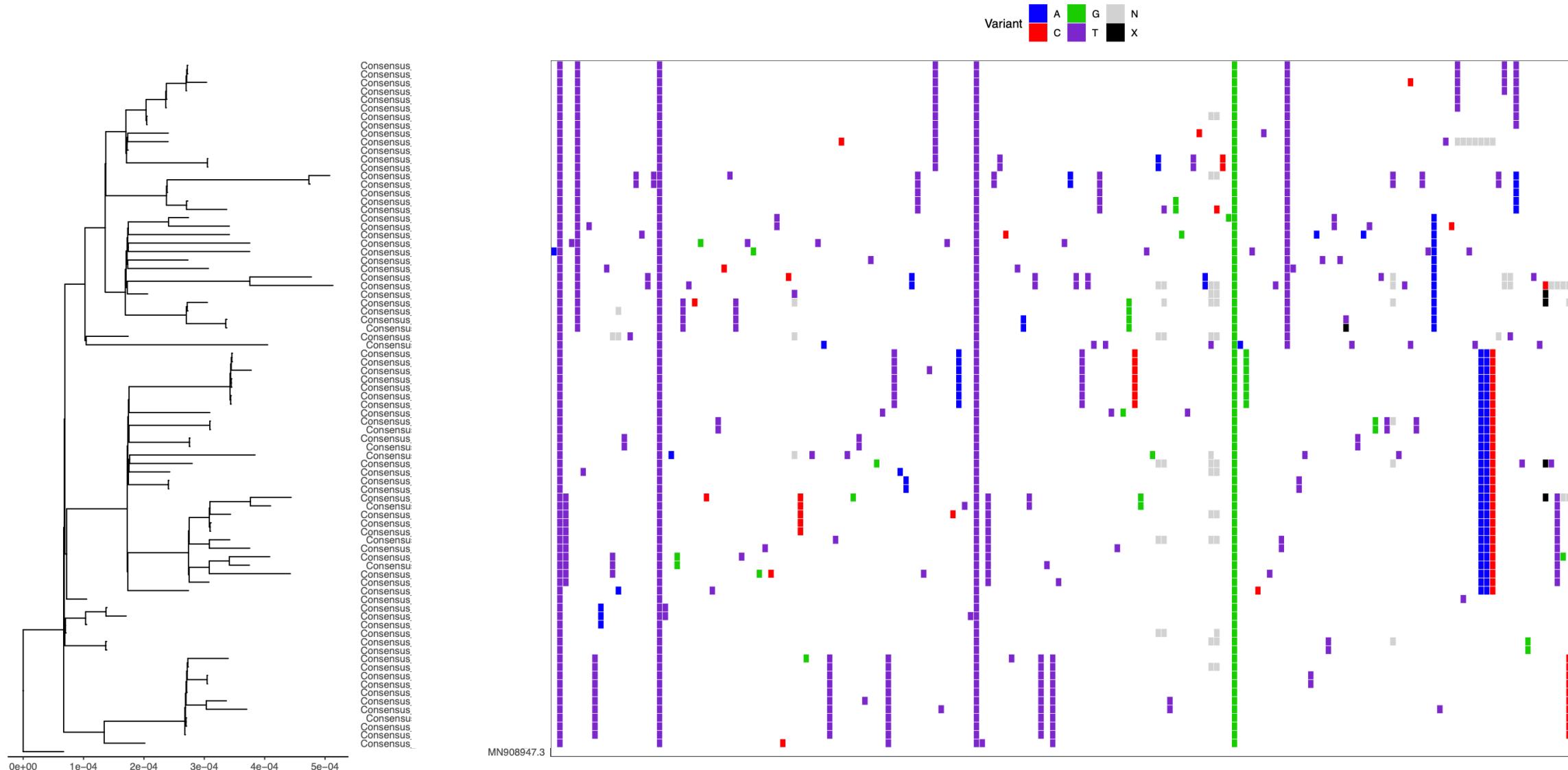
Genome Coverage QC



What do we mean by quality?

- Was the genome successfully sequenced?
 - Virus abundance varies sample-to-sample, storage conditions and extraction methods may differ in ability to recover complete genomes
- **Is the genome sequence accurate?**
 - Low abundance samples may have amplification artifacts, low coverage may lead to consensus errors
- Is the sequencing run contaminated?
 - Sequencing protocols are highly multiplexed, many amplification cycles increases risk of cross-contamination on the plate

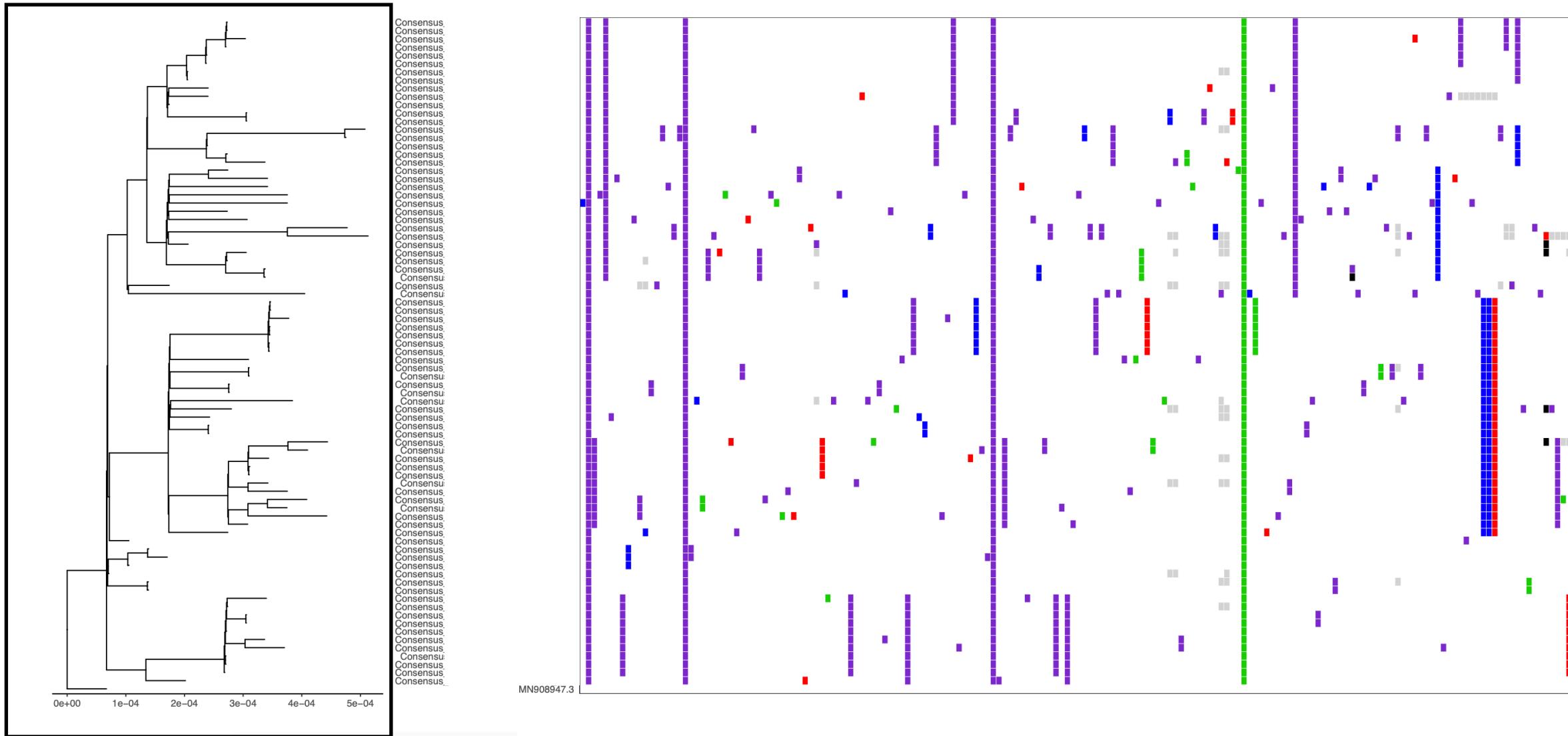
Assessing variants - ncov-tools tree/snps



This plot was inspired by work from Mads Albertsen

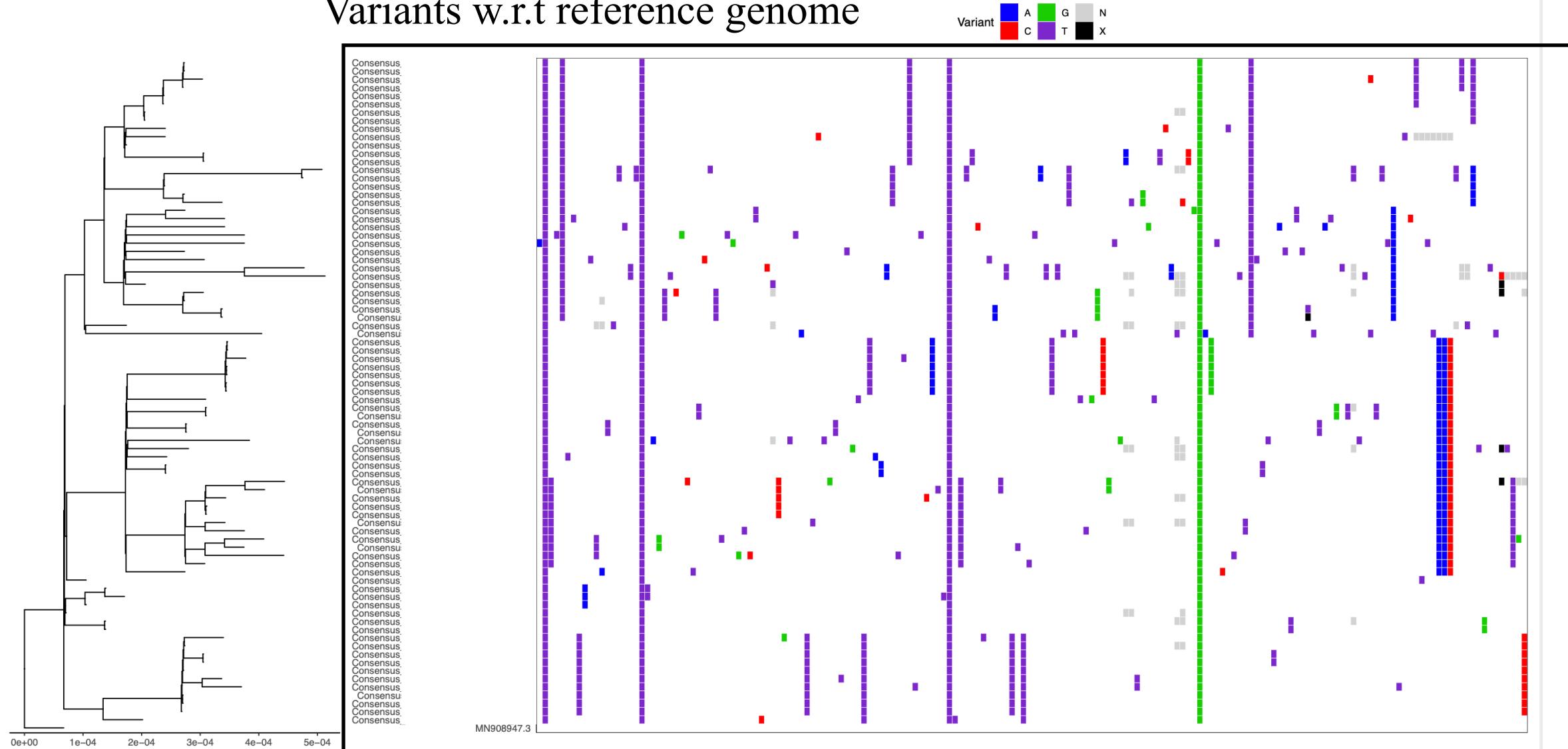
Assessing variants - ncov-tools tree/snps

tree built by MAFFT+iqtree

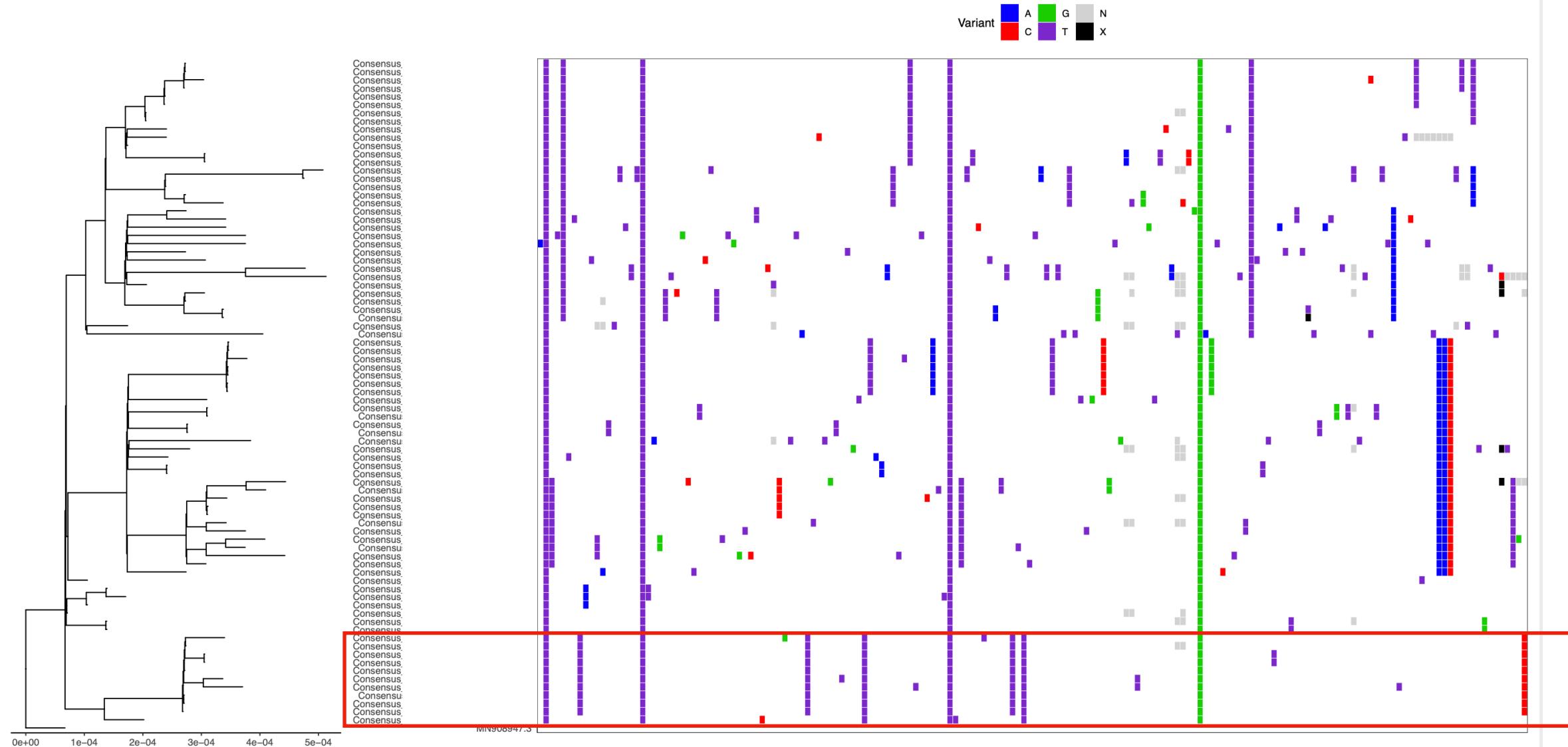


Assessing variants - ncov-tools tree/snps

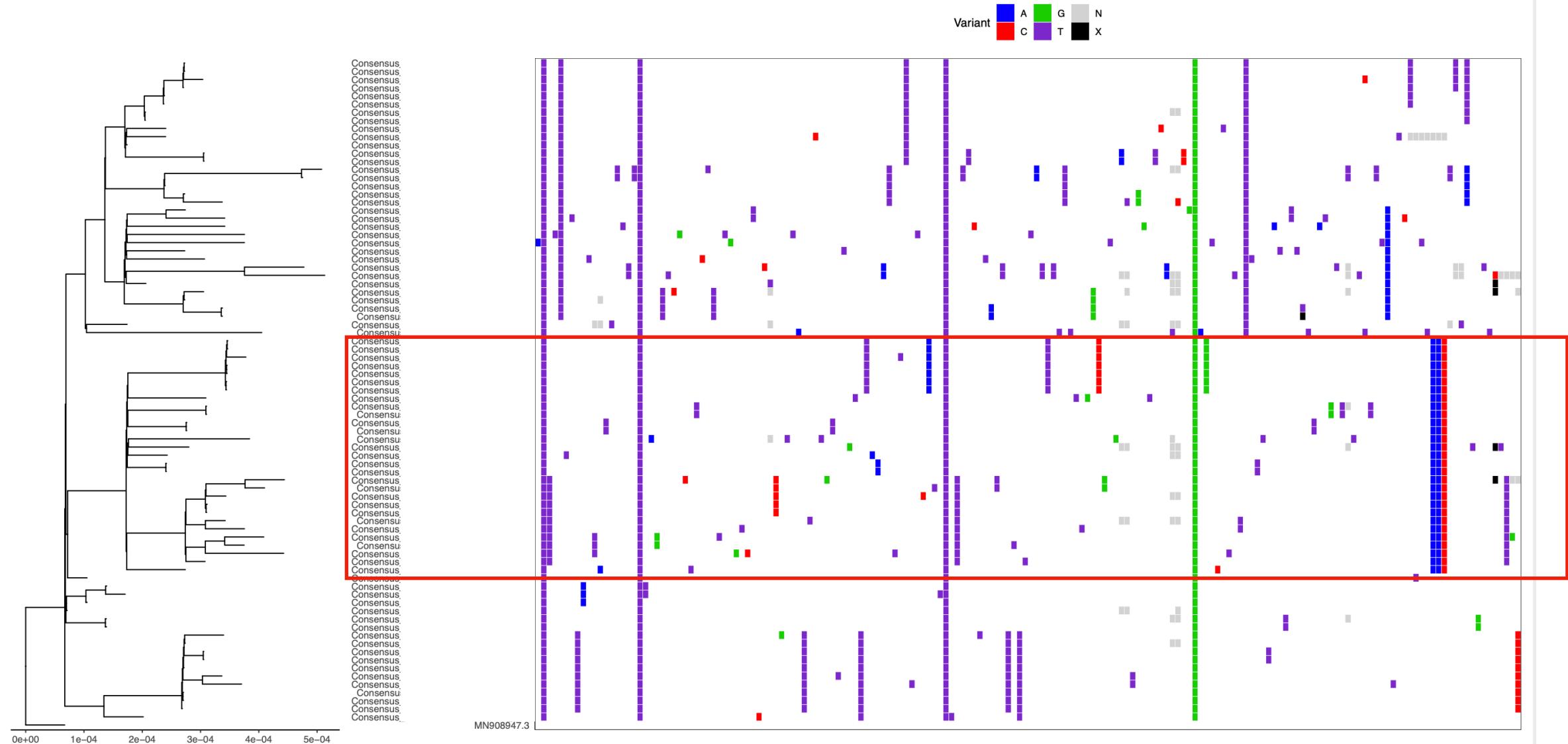
Variants w.r.t reference genome



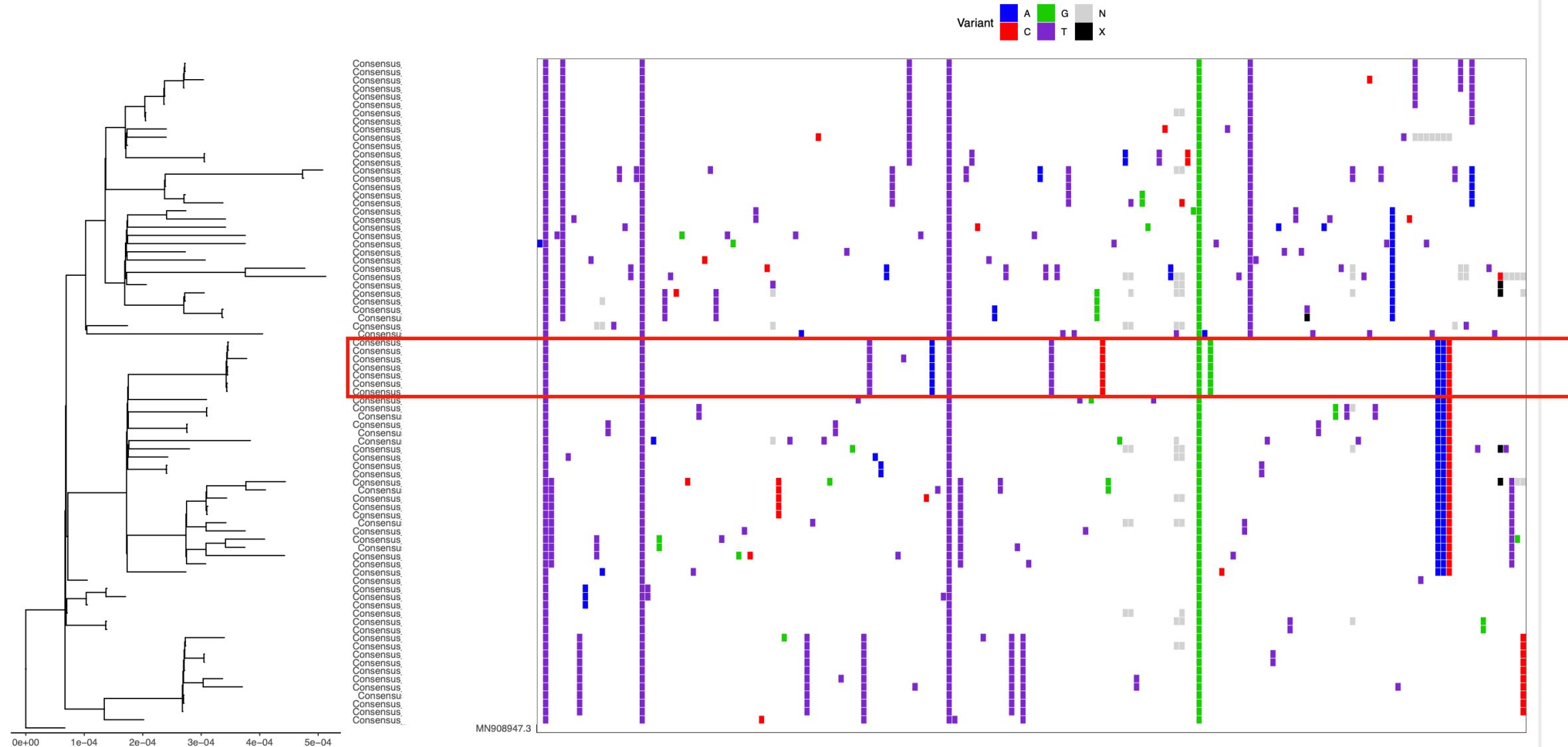
Assessing variants - ncov-tools tree/snps



Assessing variants - ncov-tools tree/snps



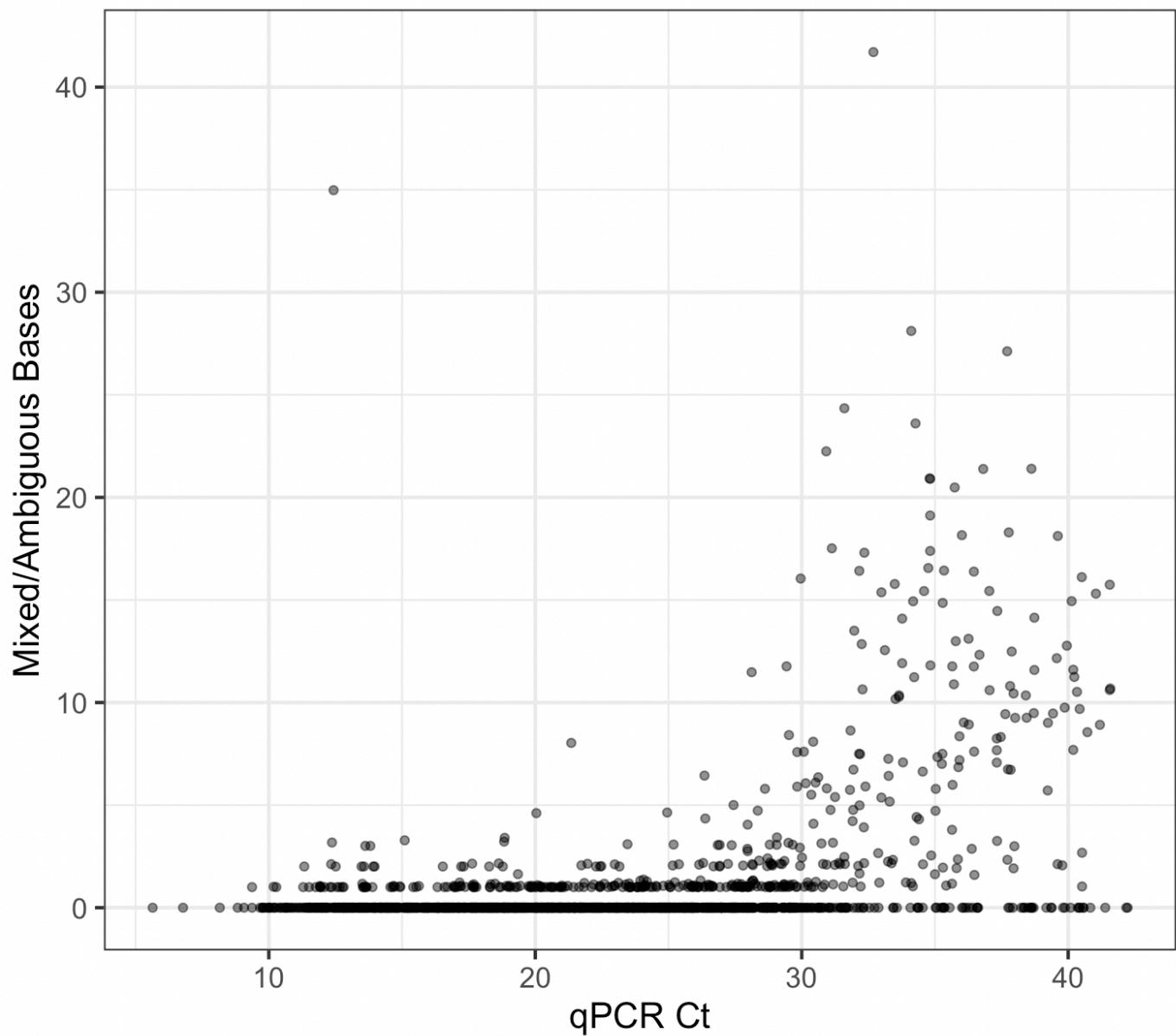
Assessing variants - ncov-tools tree/snps



Variant QC - Mixed/Ambiguous Bases

Number of observed ambiguous bases increases for high Ct samples

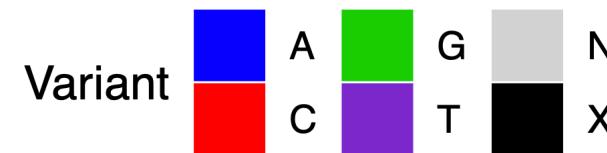
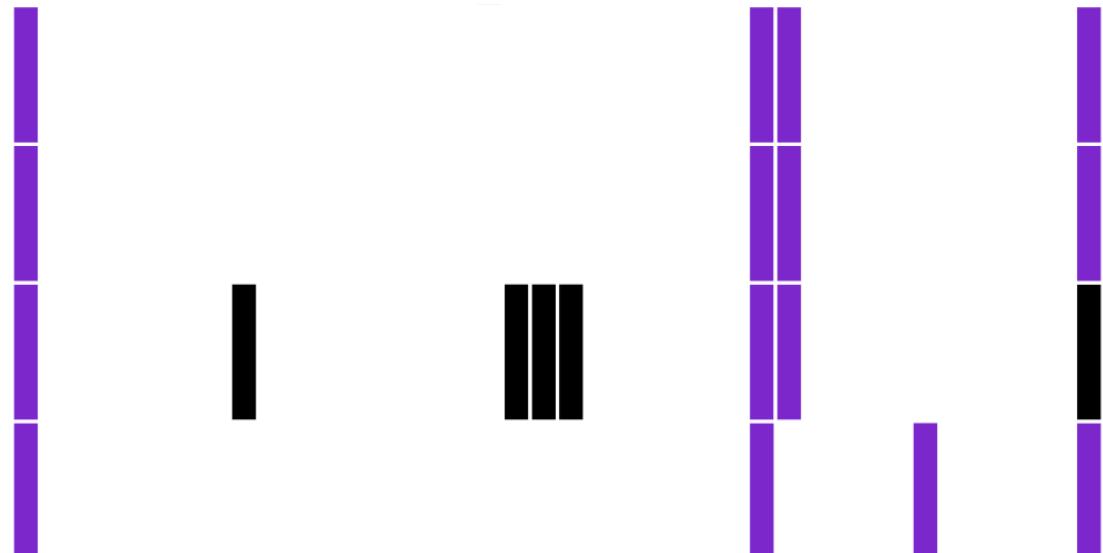
Possible RNA edits, RT or PCR errors from low-template samples



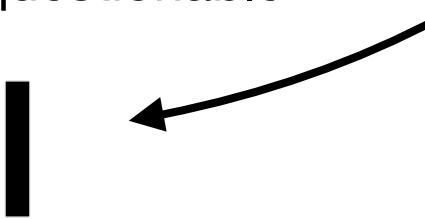
Variant QC - Mixed/Ambiguous Bases

Ambiguous bases (IUPAC symbols) are displayed using black rectangles

We QC-fail samples that have 5 or more ambiguous positions



Samples with many mixed positions are questionable

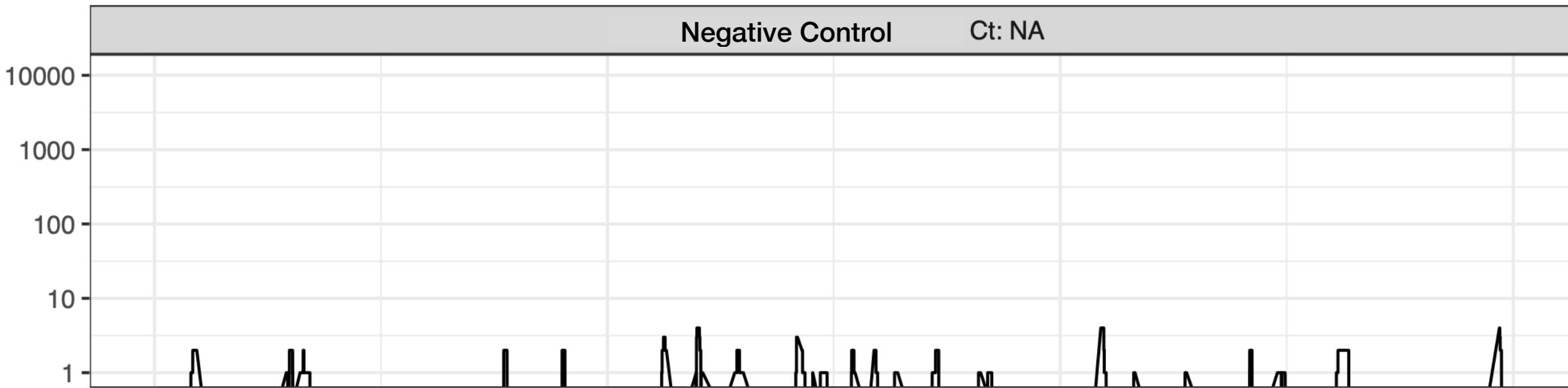


What do we mean by quality?

- Was the genome successfully sequenced?
 - Virus abundance varies sample-to-sample, storage conditions and extraction methods may differ in ability to recover complete genomes
- Is the genome sequence accurate?
 - Low abundance samples may have amplification artifacts, low coverage may lead to consensus errors
- **Is the sequencing run contaminated?**
 - Sequencing protocols are highly multiplexed, many amplification cycles increases risk of cross-contamination on the plate

Contamination - Assess by Negative Control

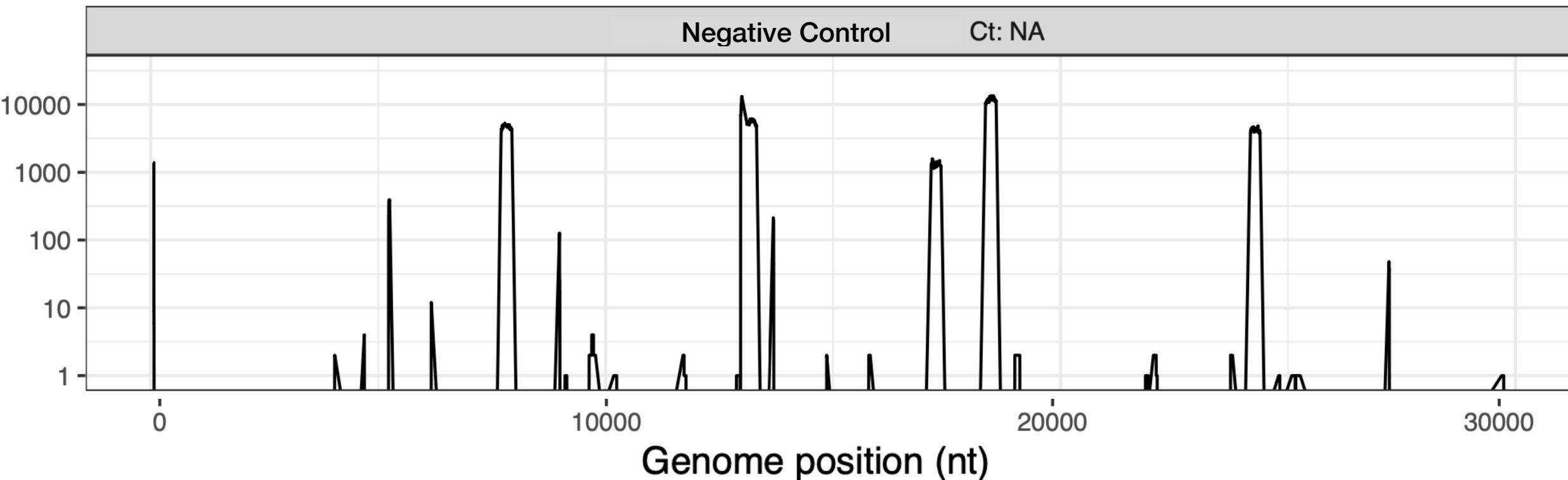
Ideal: very few reads mapped to genome



Contamination - Assess by Negative Control

Entire amplicons present in negative control

Recommendation: discard run, or mask these regions

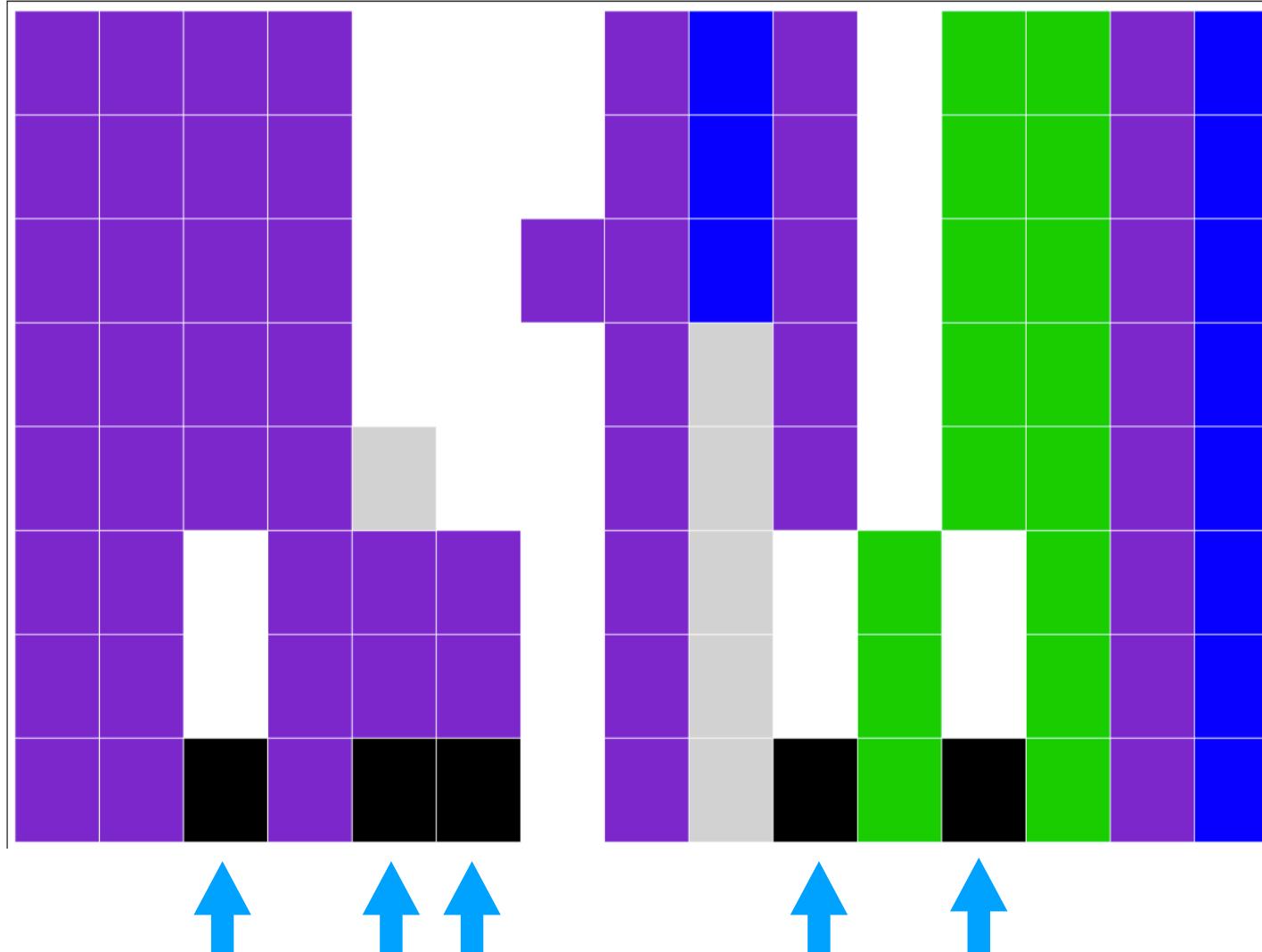


Sample contamination

SampleC has evidence of alleles from both sampleA and sampleB

Recommendation: discard sample after investigation

sampleA



sampleB

sampleC

We are on a Coffee Break & Networking Session

Workshop Sponsors:



Canadian Centre for
Computational
Genomics



HPC4Health

