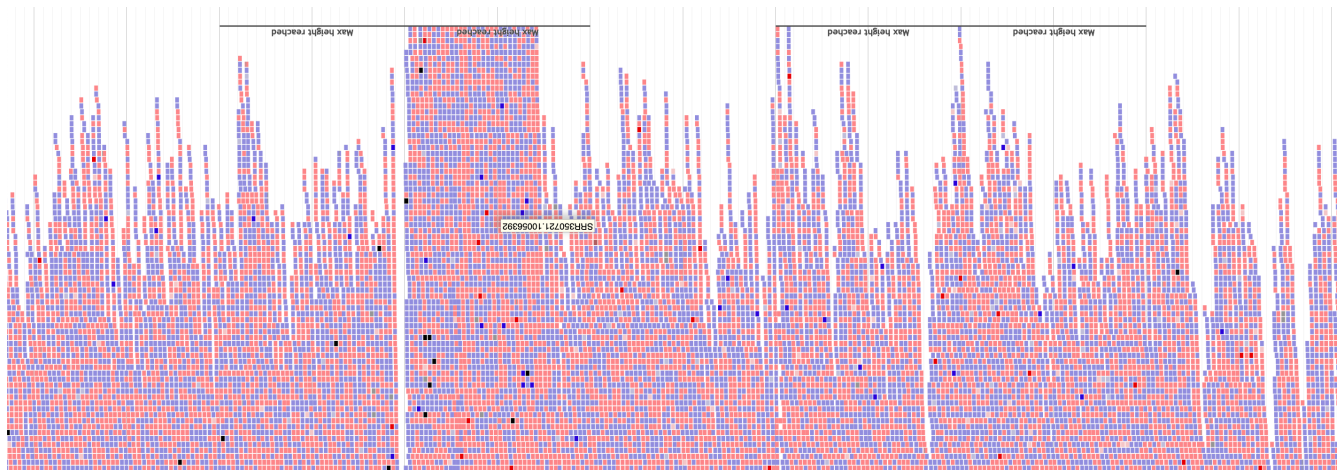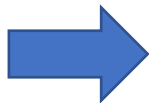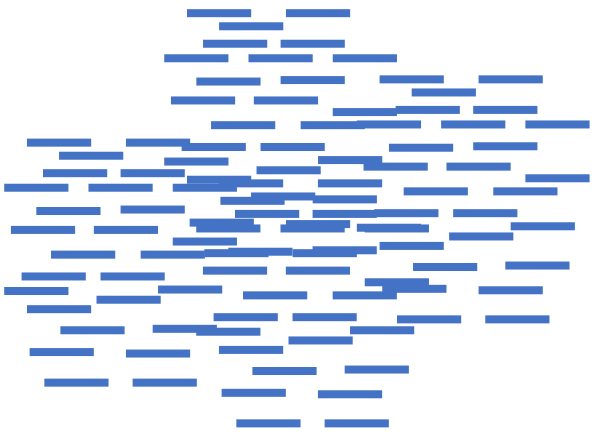# Variation data from high throughput sequencing

Omar S. Harb, PhD
University of Pennsylvania
VEuPathDB

Experimental Sequencing Reads

Aligned Reads
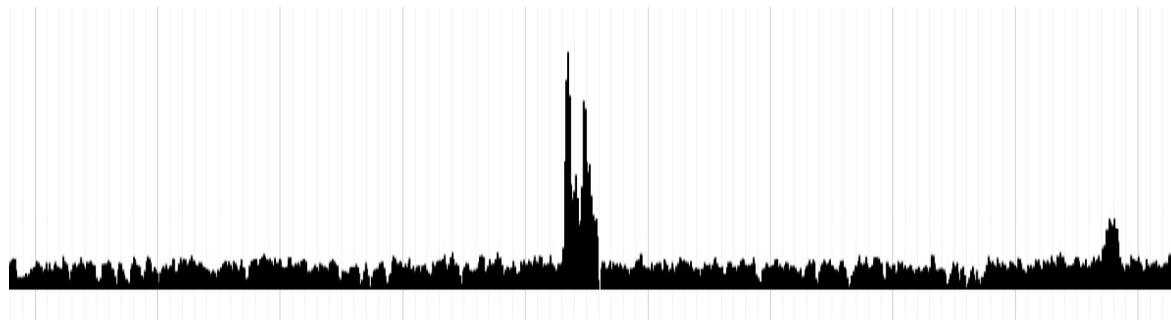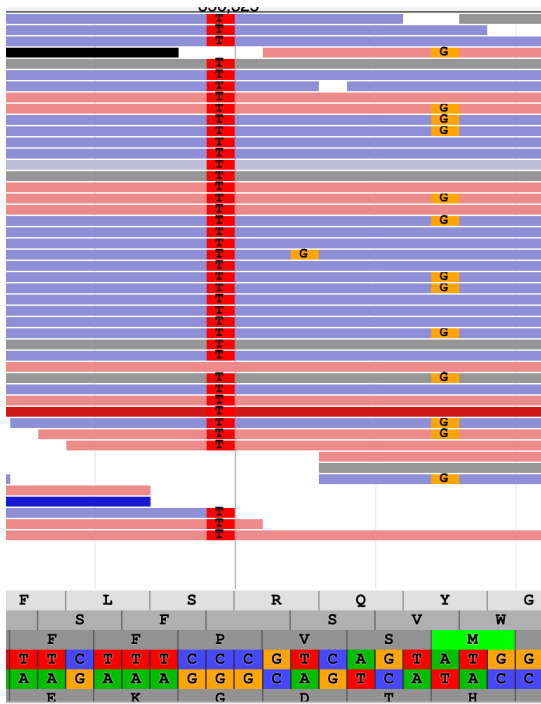
Reference Genome

Variants

CNV

# Single Nucleotide Polymorphisms (SNPs)
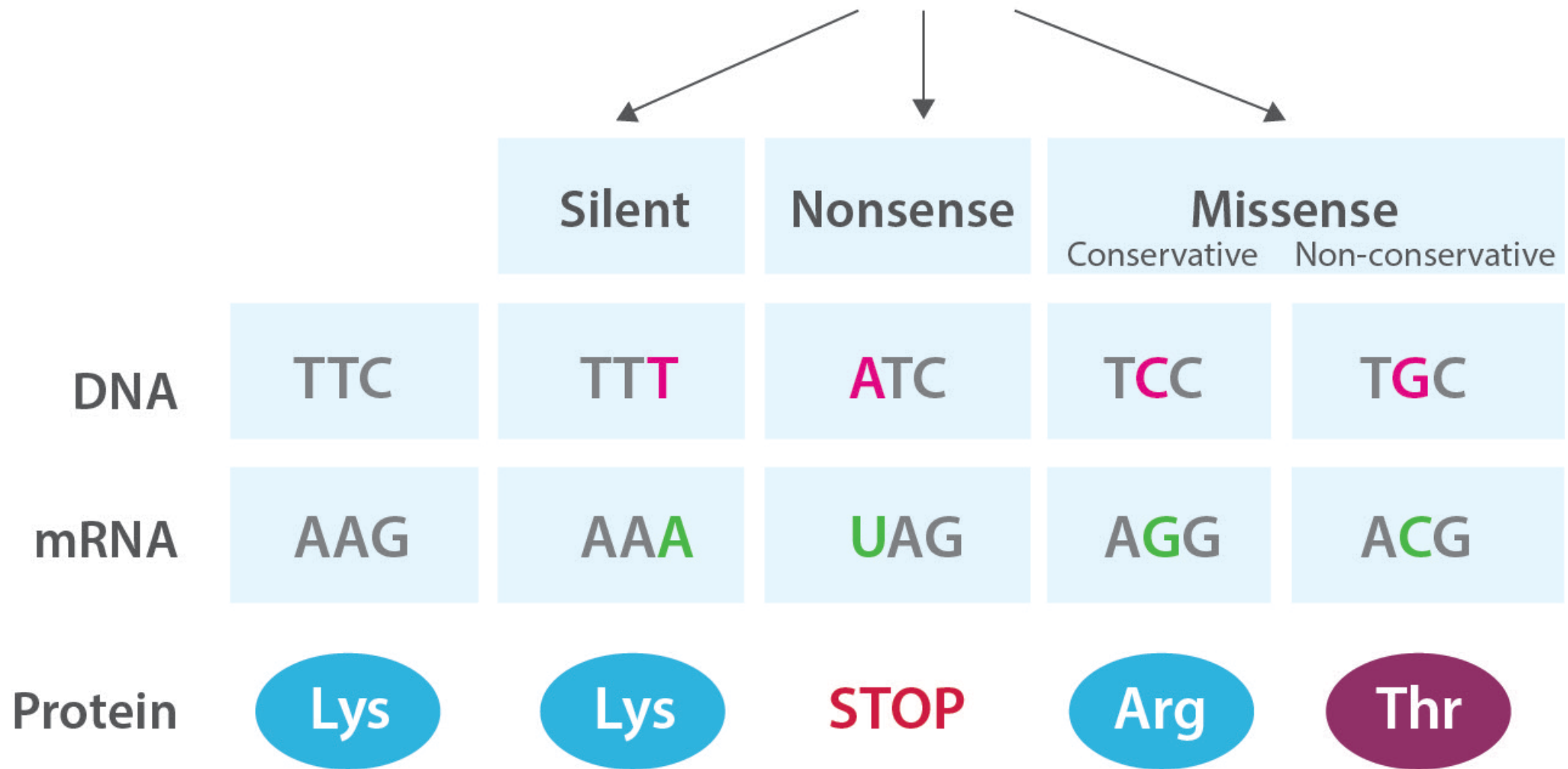
# What are SNPs

SNPs are genetic variations. Each SNP represents a single change in one nucleotide.

There are different types of SNPs but broadly can be divided into synonymous, nonsynonymous non-coding.

**_Synonymous_**: a  nucleotide change that does not change the amino acid sequence – also called silent mutation.

**_Nonsynonymous_**: a nucleotide change that results an amino acid change.  These can be missense (results in a different amino acid) or nonsense (results in a stop codon) mutations.

# Point Mutations

|  | | Silent | Nonsense | Missense | |
|---|---|---|---|---|---|
|  | | | | Conservative | Non-conservative |
| **DNA** | TTC | TTT | ATC | TCC | TGC |
| **mRNA** | AAG | AAA | UAG | AGG | ACG |
| **Protein** | Lys | Lys | STOP | Arg | Thr |

Important factors to think about when considering SNPs:

- Can a SNP be called based on the data?
- How many of your sample had suffient reads to call a SNP at that position?
- What is the frequency of the SNP?  What could affect read frequency?
- How many of your samples have that SNP? Minor/major allele frequency.

# Isolate X aligned sequencing reads



Nucleotide position = 1567
Aligned reads = 10
Reference = A
Isolate reads = G (10)
Read Frequency = 100%

Nucleotide position = 1583
Aligned reads = 4
No Call: <5 reads

Nucleotide position = 1600
Aligned reads = 10
Reference = G
Isolate reads = A (4) and G (6)
Read Frequency G = 60%
Read Frequency A = 40%
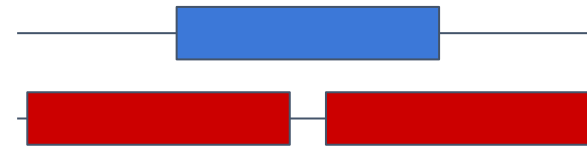
# Copy Number Variations

# What are Copy Number Variations (CNVs)?

A phenomenon in which a region of the genome appears a different number of times in different individuals

No CNV

Duplication

Deletion

# Types of CNV

- **Polyploidy**
  - Duplication of the whole genome
  - Common in plants
  - Happened twice in the vertebrate lineage leading to humans
- **Aneuploidy**
  - Duplication of some but not all chromosomes
  - Common in *Leishmania* and *Trypanosoma cruzi*
- **Segmental duplications**
  - Duplication of part of a chromosome
  - We'll look at an example in *Candida*
- **Gene duplications**
  - Duplication of a single gene
  - *ROP5* in *Toxoplasma gondii*
- **Short CNVs**
  - Duplications of short regions, e.g., CAG repeats in Huntingdon's disease or BRC repeats in *BRCA2*

# Types of CNV

- **Polyploidy**
  - Duplication of the whole genome
  - Common in plants
  - Happened twice in the vertebrate lineage leading to humans
- **Aneuploidy**
  - Duplication of some but not all chromosomes
  - Common in *Leishmania* and *Trypanosoma cruzi*
- **Segmental duplications**
  - Duplication of part of a chromosome
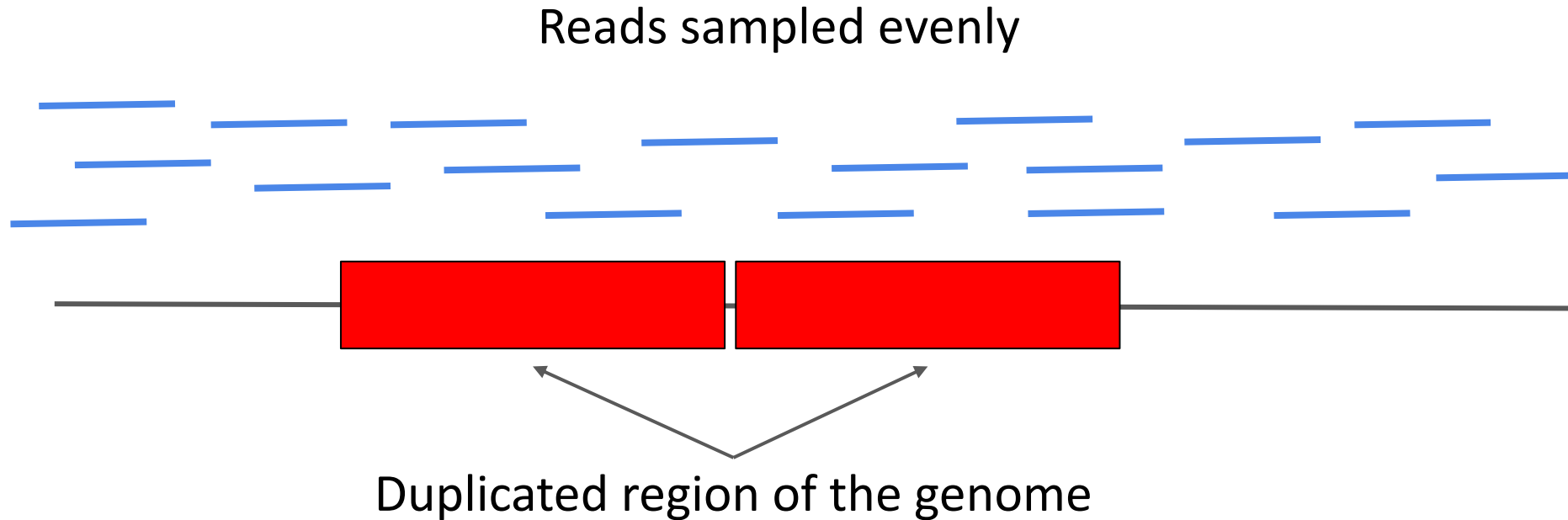  - We'll look at an example in *Candida*
- **Gene duplications**
  - Duplication of a single gene
  - *ROP5* in *Toxoplasma gondii*
- **Short CNVs**
  - Duplications of short regions, e.g., CAG repeats in Huntingdon's disease or BRC repeats in *BRCA2*
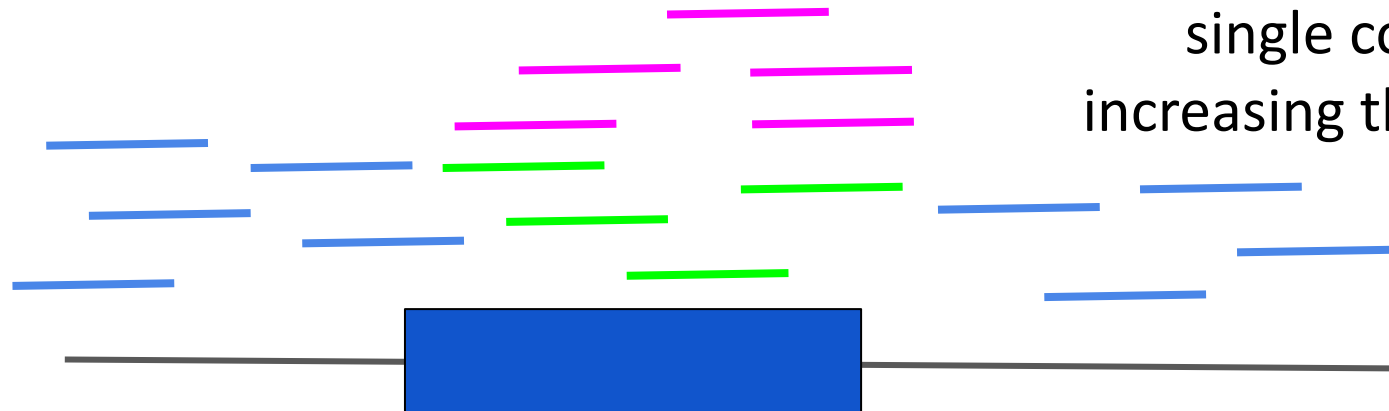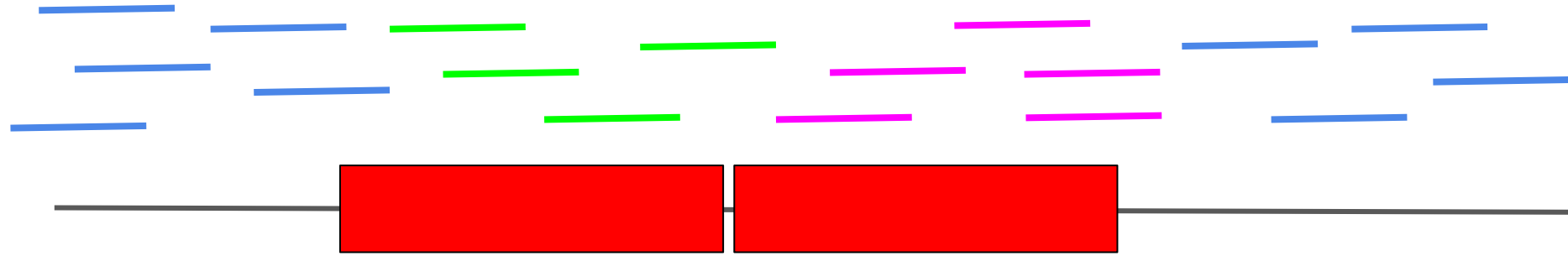
# CNVs and Coverage

In whole genome sequencing (WGS) one would expect to sample evenly along the genome resulting in even coverage after alignment

Reads sampled evenly

Duplicated region of the genome

This is subject to minor variation associated with GC content and stochastic variation
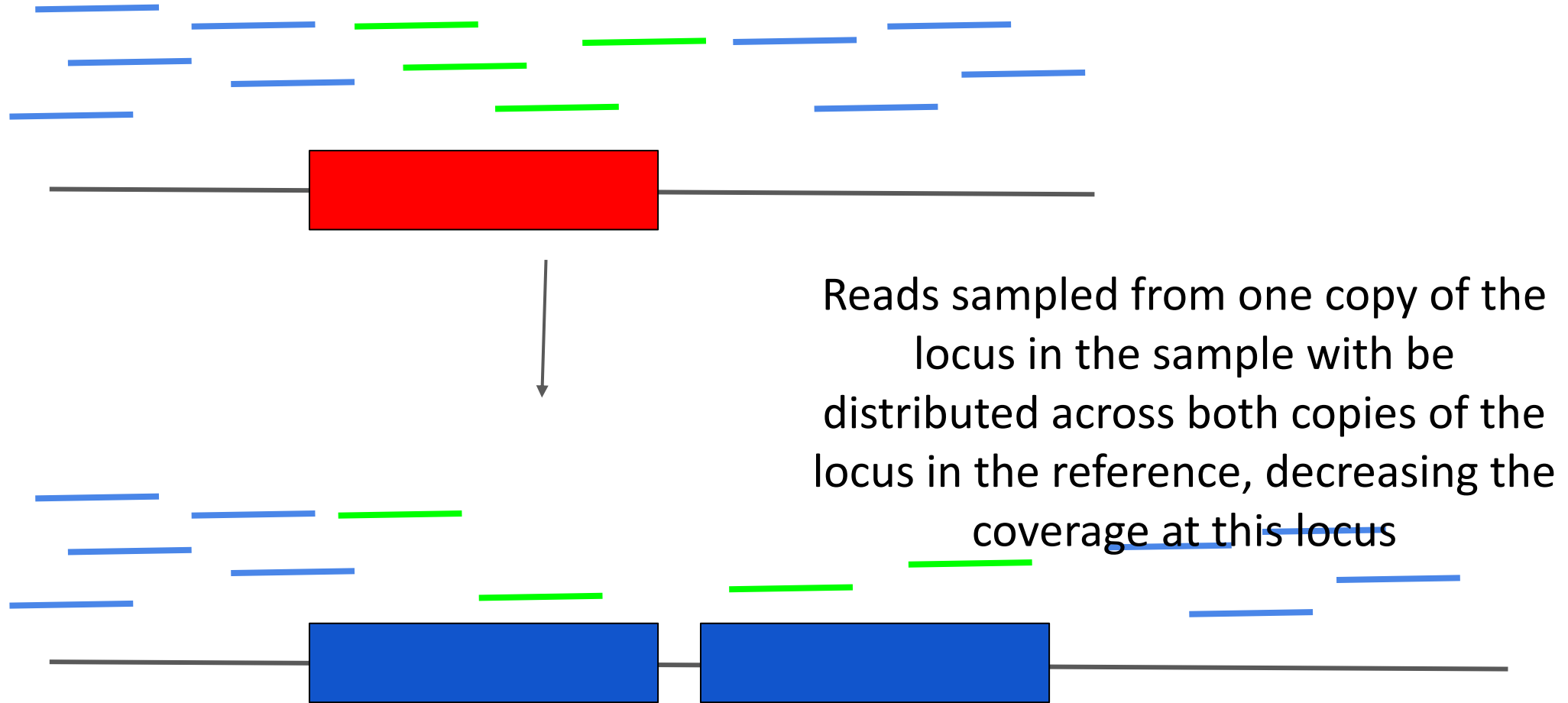
# Duplication results in increased coverage



Reads sampled from both copies of the duplication will all map to the single copy in the reference, increasing the coverage at this locus

Reference genome without duplication

# Deletion results in decreased coverage



Reads sampled from one copy of the locus in the sample with be distributed across both copies of the locus in the reference, decreasing the coverage at this locus

# Duplications Associated with Drug Resistance



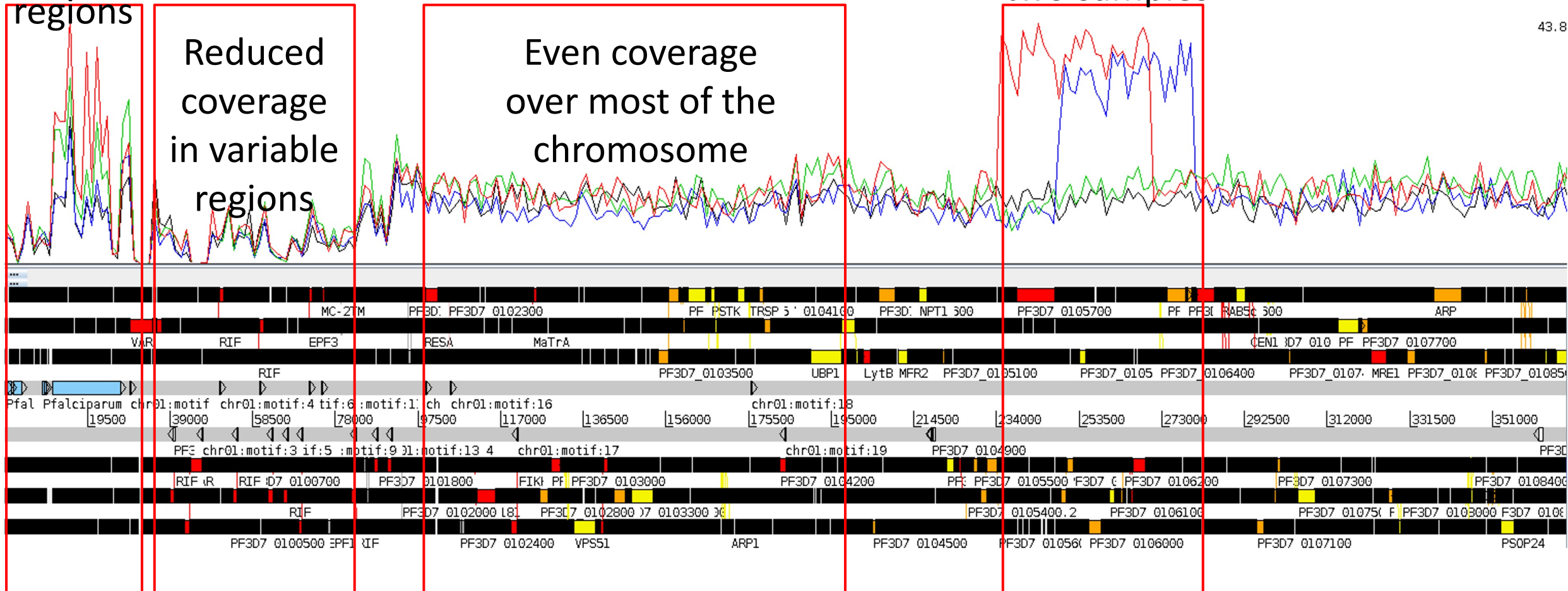Variable coverage in repetitive regions

Reduced coverage in variable regions

Even coverage over most of the chromosome

Duplications in two samples

# In VEuPathDB

- Search for supernumerary chromosomes in different isolates
- Search for genes with increased or decreased copy number in different isolates
- Explore coverage in JBrowse