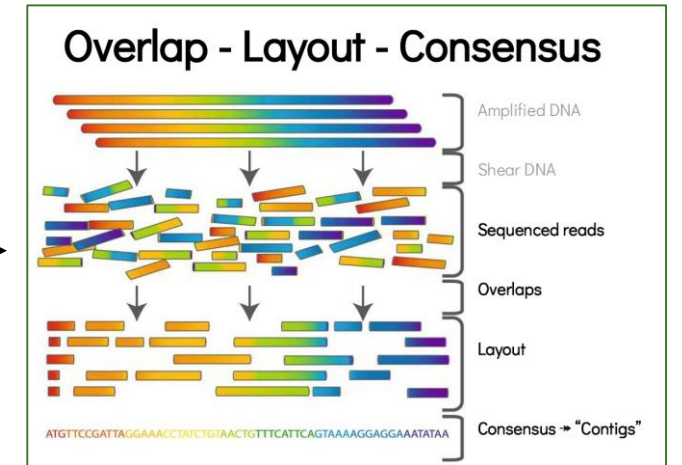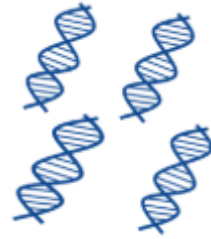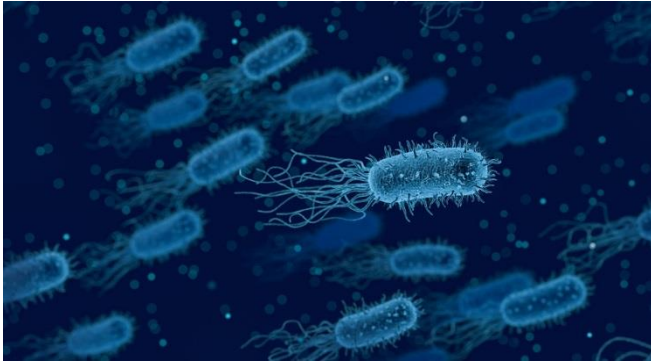# Alphabet Soup of WGS Analysis: MLST and SNP

Heather Blankenship, PhD Candidate Michigan State University,
Bioinformatics Intern Michigan Department of Health and Human Services
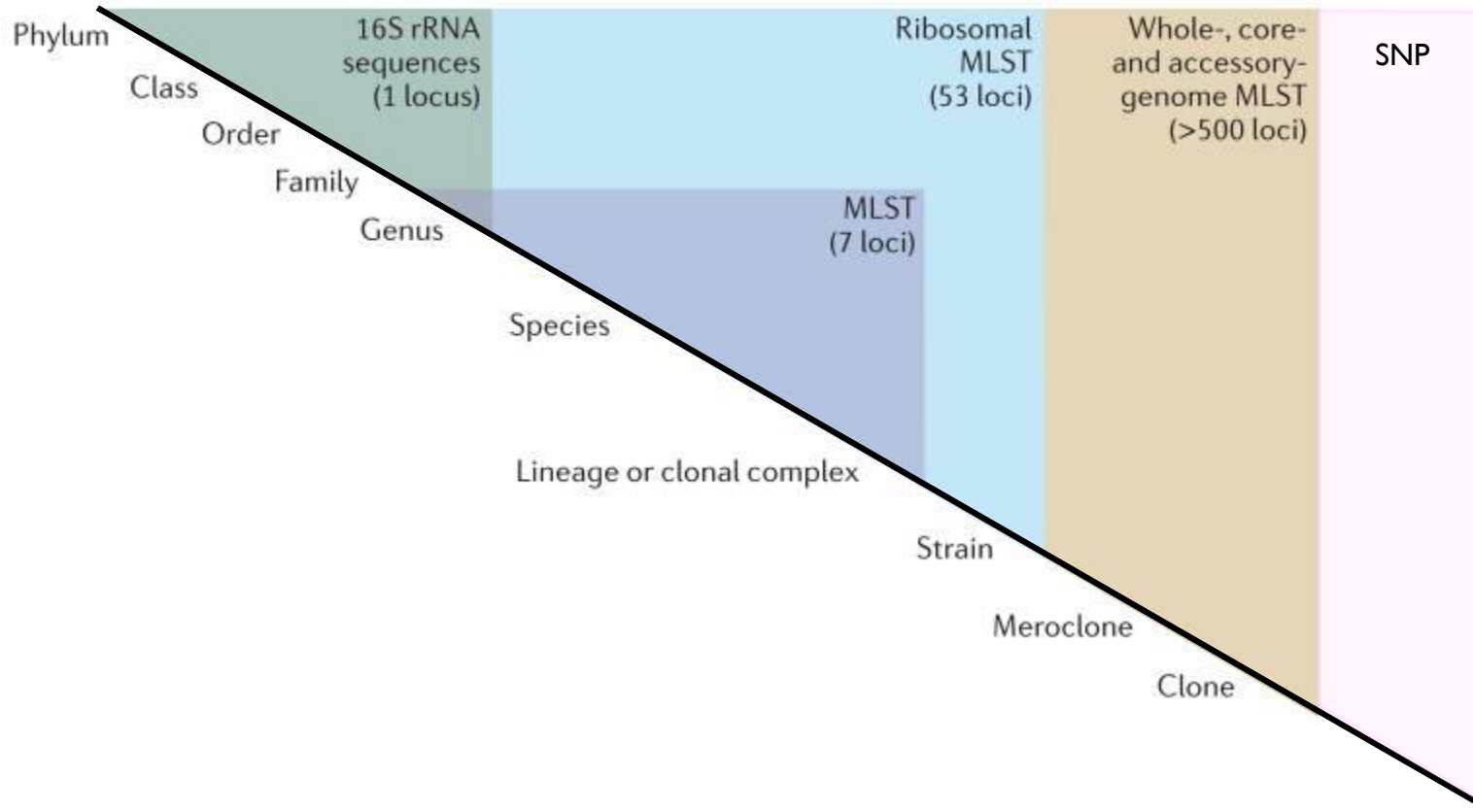
# WGS Pipeline



Overlap - Layout - Consensus

Amplified DNA

Shear DNA

Sequenced reads

Overlaps

Layout

Consensus → "Contigs"

ATGTTCCGATTAGGAAACCTATCTGTAACTGTTTCATTCAGTAAAAGGAGGAAATATAA

Torstern Seeman

Now what?

# Choosing a Method



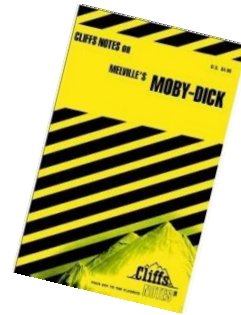Maiden MC *et al* (2013); Nat Rev Microbiol

- ▸ **MLST- Multi Locus Sequencing Typing**

- ▸ **SNP- Single Nucleotide Polymorphism**
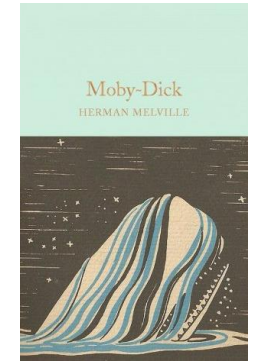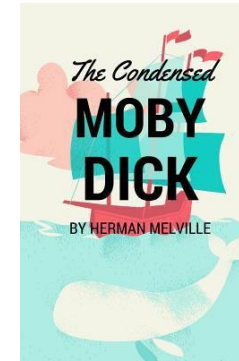
# WGS Methods

- **MLST**
  - CliffNotes version of a book

- **cgMLST**
  - Compare abridged versions of a book with other versions and only focus on the chapters that are the same in all of them

- **wgMLST analysis**
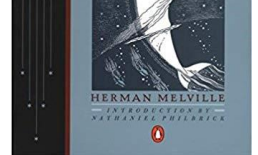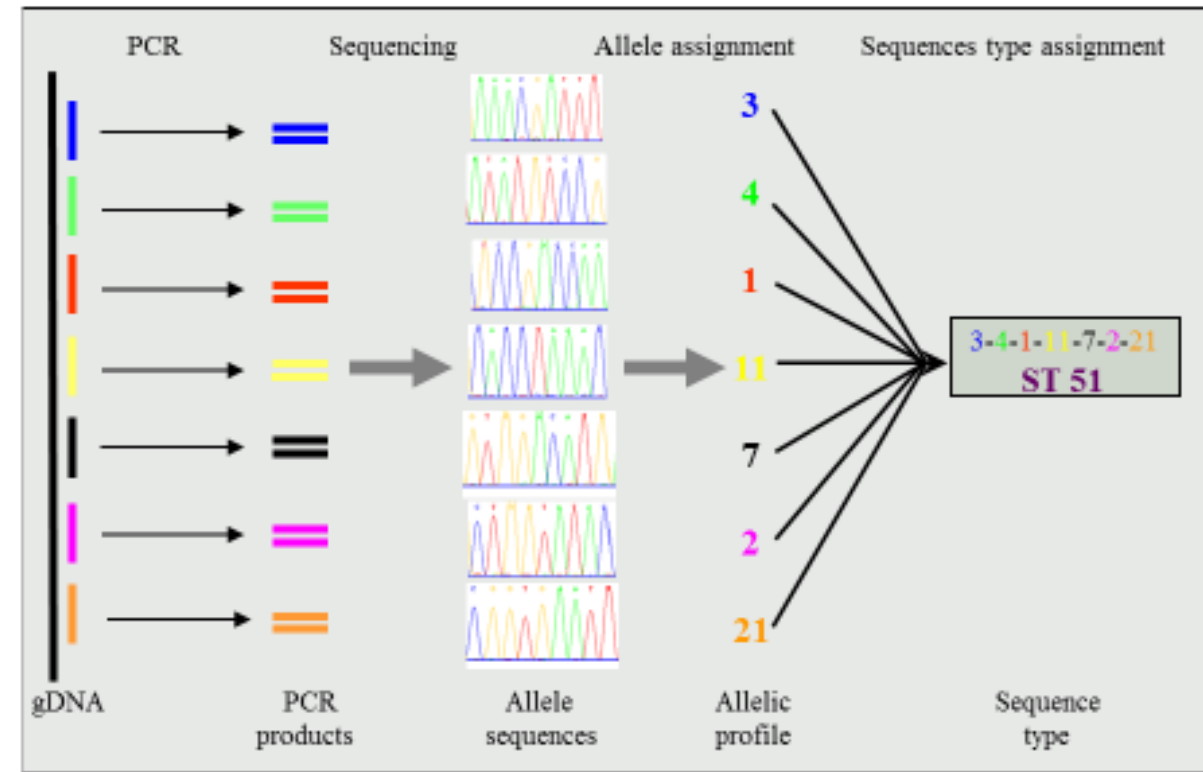  - Compare two books on a chapter by chapter basis

- **SNP analysis**
  - Compare two books character by character

The dog and cat watched the bird sing.
The **h**og and cat **m**atched the bird s**o**ng.

**Edition 1**    **Edition 8**

# MLST

- **Traditional MLST**
  - ~6-7 housekeeping genes
  - Allelic variation for each gene is identified
  - Sequence Type (ST) assigned based on allelic profile
  - Databases are internationally available
  - Suitable for examining differences at the population level

- **Allele**: one or more alternative forms of a gene
- **Locus**: gene or region of the gene that is being extracted and compared
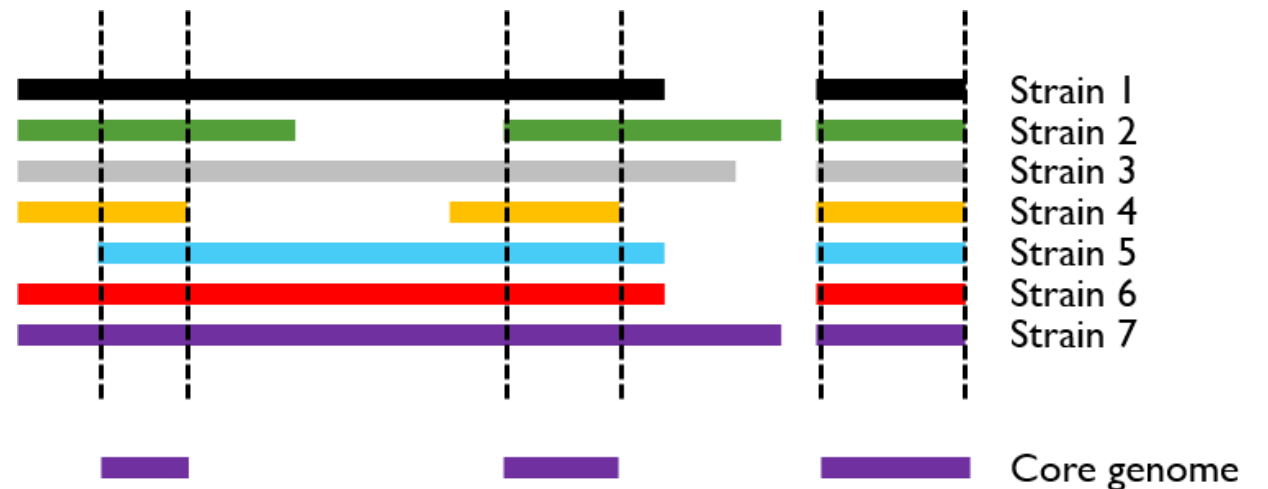


Ruppitsch W (2016); J Land Management, Food and Env

# wgMLST (whole genome)

- Compare genomes of interest gene by gene
- Database comprises genes from the genus/species of interest and represents a diverse genetic background within the genus/species
  - Does not have to be the full gene included in the database
  - A new allele is only added as it is encountered
- STANDARDIZED!

| Gene A | |
|---|---|
| Allele 1 | ATGTAGCGCTAGCC |
| Allele 2 | ATGTAGCCCTAGCC | SNP |
| Allele 3 | ATGTAGATGGCTAGCC | SNP + insertion |

# cgMLST (core genome)

▸ Similar to wgMLST, compares on a gene by gene basis

▸  What genes/portions of genes are common to all within the species

   ▸ Results in a smaller databases

▸ Still STANDARDIZED

▸ wgMLST or cgMLST

   ▸ wgMLST – higher resolution

   ▸ cgMLST – more stability



Strain 1
Strain 2
Strain 3
Strain 4
Strain 5
Strain 6
Strain 7

Core genome

# Allele codes/profiles

**MLST**
Sequence A – 4 . 4 . 5 . 6 . 10 . 1 . 5 – ST30
Sequence B – 4 . 6 . 5 . 6 . 10 . 1 . 5 – ST35
Sequence C – 1 . 4 . 5 . 6 . 10 . 1 . 1 – ST4

**wg/cgMLST**
Sequence A – 5 . 5 . 21 . 6
Sequence B – 5 . 5 . 21
Sequence C – 1 . 5

▸ The allele codes/profiles do not change regardless of the isolates that are added to an analysis

▸ Allows comparison with other labs/states/outbreaks etc

# Benefits/Disadvantages of MLST methods

▶ Disadvantage

  ▶ A lot of upfront development to develop and continually curate the databases

  ▶ SNPs, insertions, deletions are all treated the same, an allele may have multiple evolutionary events but it is not evident by looking at the allele number

  ▶ Comparison is based on the allele numbers and not the genetic sequences

  ▶ Requires genome assembly

  ▶ Does not include non-coding regions
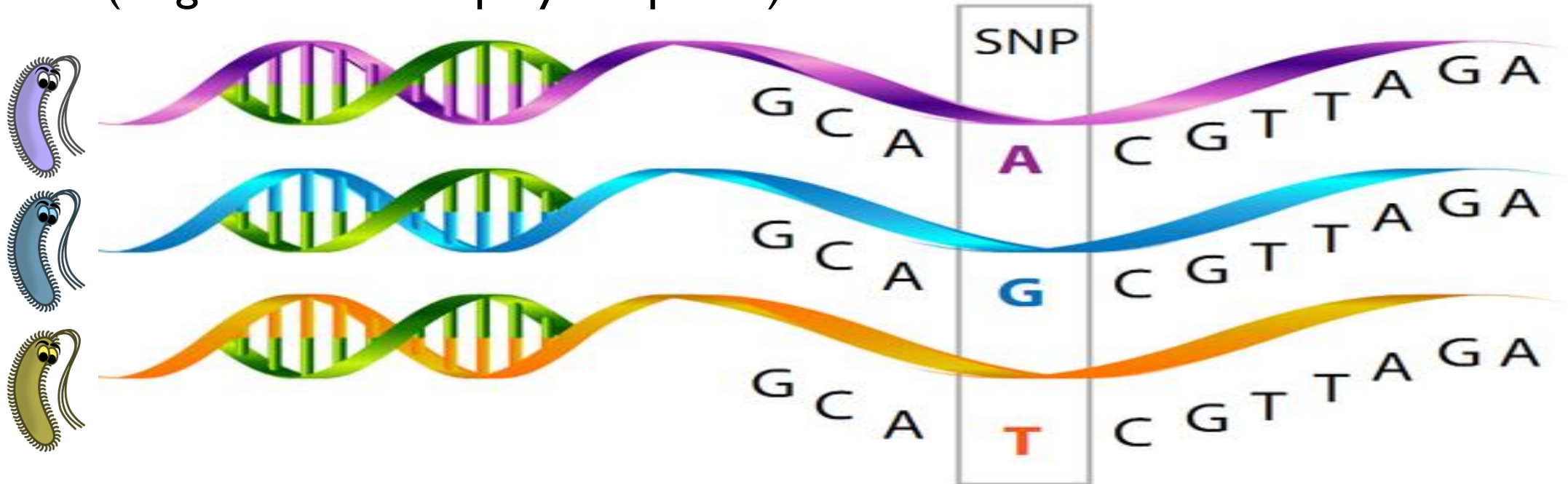
▶ Advantage

  ▶ Provides high enough differentiation of strains

  ▶ Allele call is stable

  ▶ Low computational power

  ▶ Standardized data

**Best for surveillance**

# SNP typing - Basics

- Map reads against a reference genome
- Identify SNPs between reference and reads
  - Quality, Coverage, Frequency
- Build phylogenetic tree based on concatenated SNPs
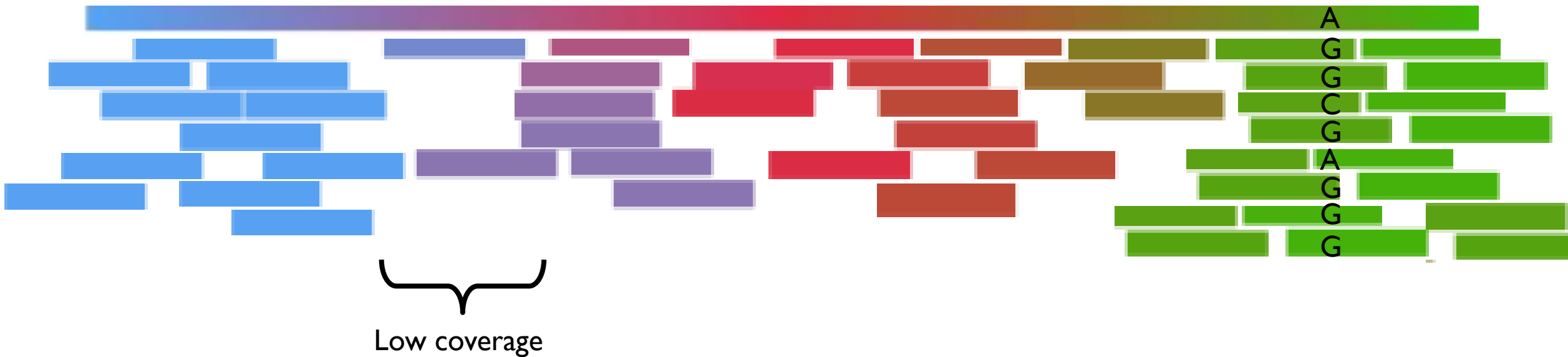- **SNP** (single nucleotide polymorphism)

# hqSNP (high quality)

▸ Takes into account the quality, coverage and frequency of the SNPs

▸ Quality: reads are filtered to ensure they pass a certain threshold



▸ Reads are cleaned to minimize the amount of erroneous data and improve average quality by removing read duplicates, reads with high frequency of ambiguous bases and adapter dimers
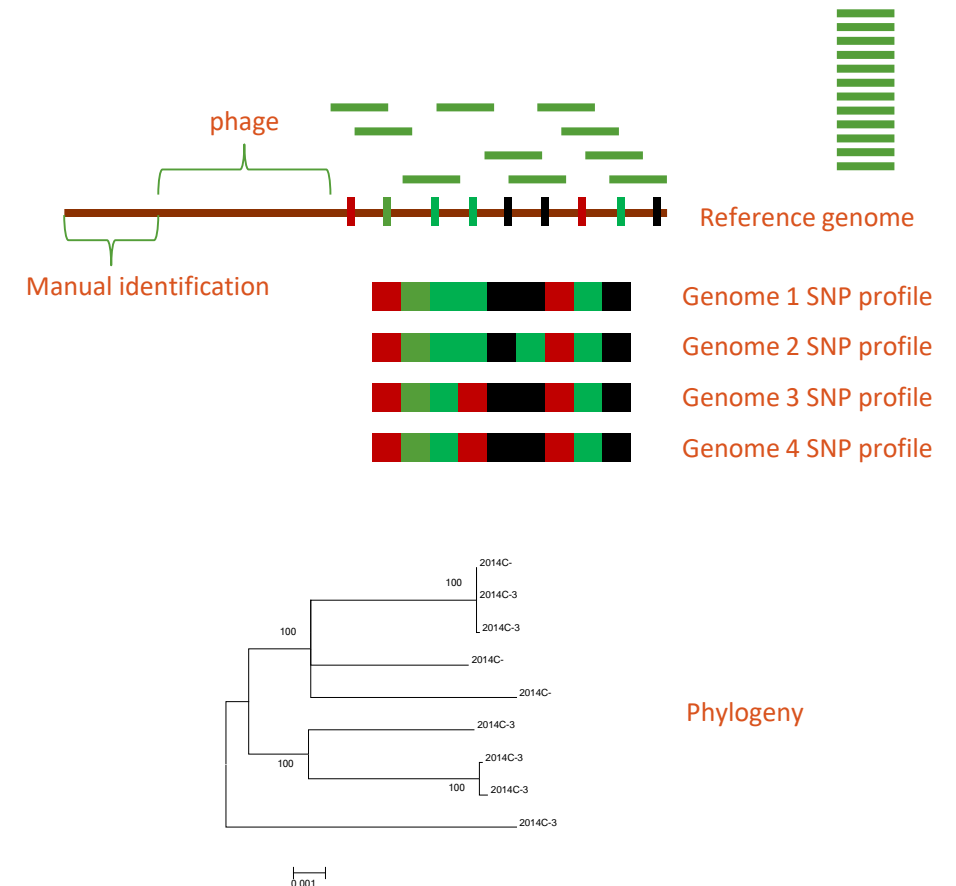
# hqSNP

- Coverage: how many reads do we have at a particular nucleotide

- Frequency: how many of the reads support the new SNP



Low coverage

# Lyve-Set (hqSNP) Process

- **Pre-Processing**
  - Phage discovery/masking
  - Manual identification of troublesome regions
  - Read cleaning
- **Mapping – SMALT**
  - 95% read identity
  - Unambiguous mapping
- **SNP calling – VarScan**
  - 75% consensus
  - 10X depth
- **Phylogeny inferring – RAxML**
  - Removal of clustered SNPs
  - Ascertainment bias model
  - Maximum likelihood
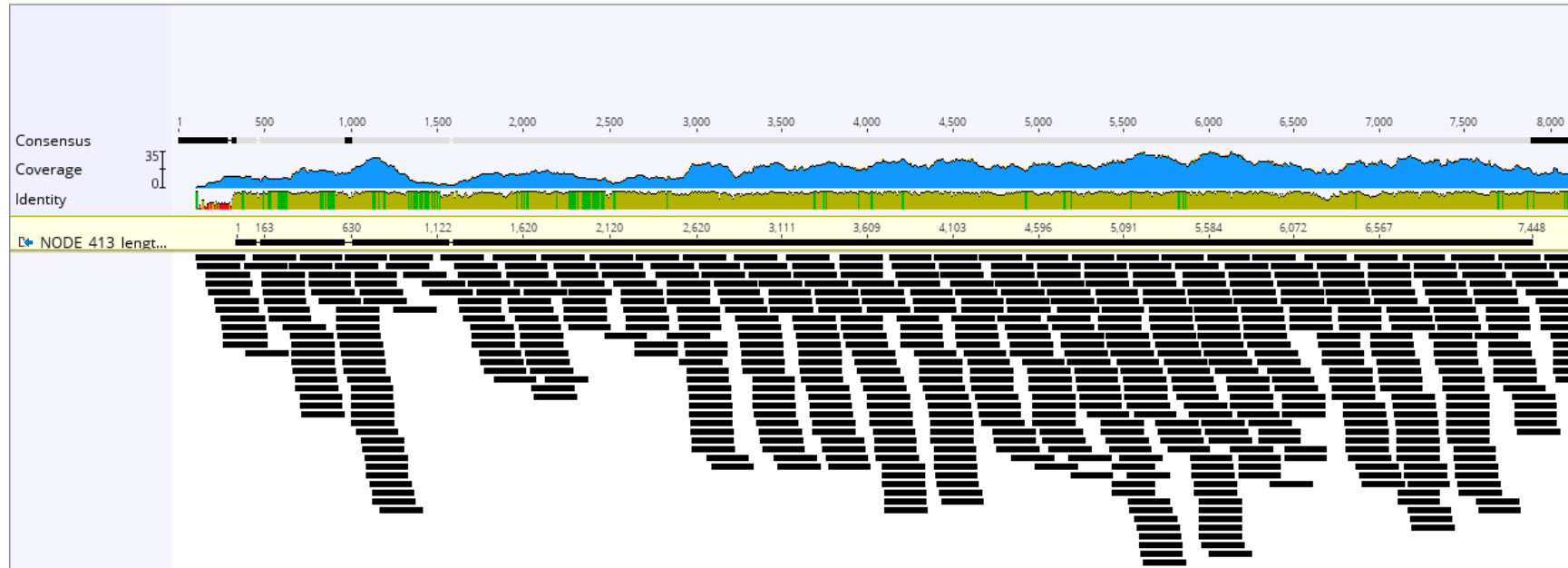


phage

Manual identification

Reference genome

Genome 1 SNP profile
Genome 2 SNP profile
Genome 3 SNP profile
Genome 4 SNP profile

2014C-
2014C-3
2014C-3
2014C-
2014C-
2014C-3
2014C-3
2014C-3
2014C-3

100
100
100
100
100

Phylogeny

0.001

Lee Katz

# SNP typing



▸ Allows for every nucleotide in the genome to be analyzed

▸ Very dependent upon the reference genome that is chosen

▸ All mobile elements (ie phages) are masked and not included in the analysis

# SNP visualization

# Benefits/Disadvantages of SNP methods

▶ Disadvantage

  ▶ Computationally intensive

  ▶ Reliant on the reference – not standardized

  ▶ Time and computationally intensive
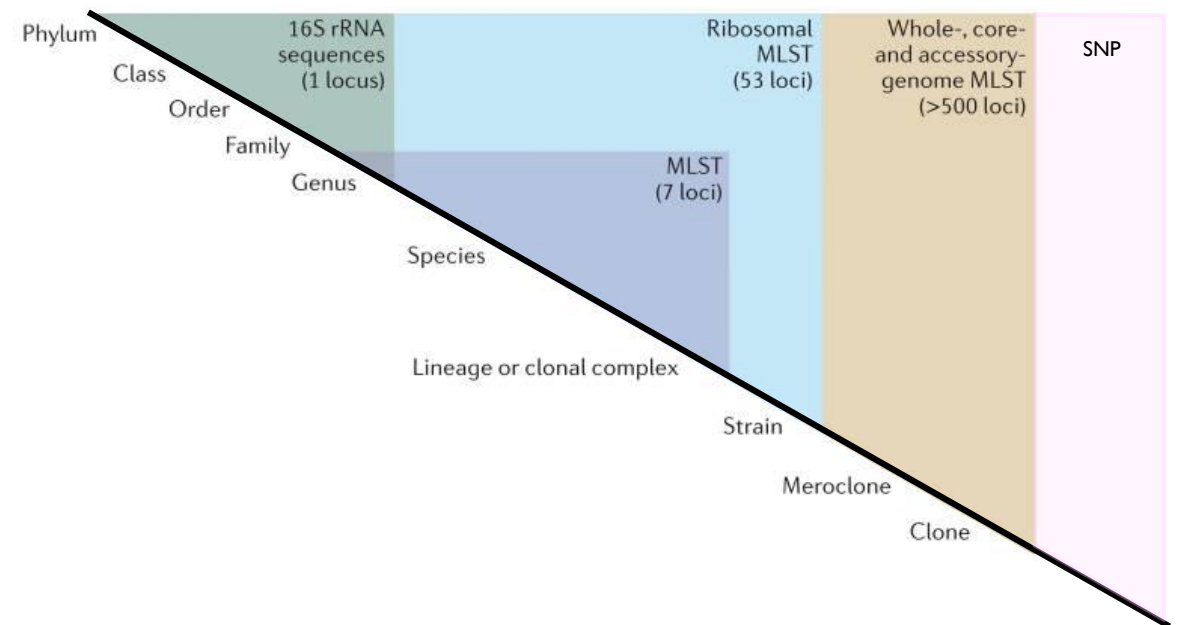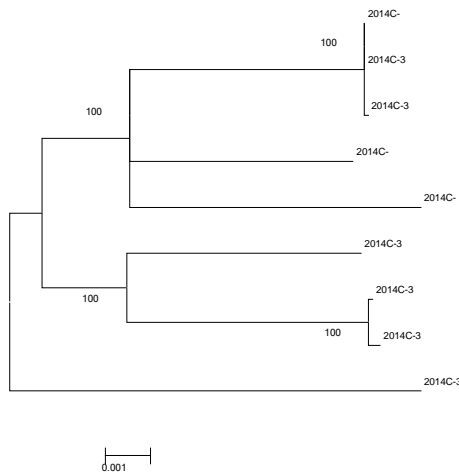
  ▶ Not stable and relies on the genomes present

▶ Advantage

  ▶ Does not require assembly of the genomes

  ▶ Very high discriminative power

  ▶ Each mutation is taken into account

  ▶ Non-coding regions are utilized

  ▶ Adaptable for all organisms

> **Best for high discriminatory power and outbreak investigation**

# Utilization of WGS Analysis

▶ Analysis will give us phylogenetic trees = <u>hypothesis</u>

▶ Strains that are genetically related based on analysis **MAY** share an epidemiological association

  ▶ What defines a cluster?

  ▶ Are cutoffs absolute?

  ▶ Does WGS "match" mean association?



Maiden MC *et al* (2013); Nat Rev Microbiol

# QUESTIONS?



▸ blankenshiph@michigan.gov

▸ selheime@msu.edu