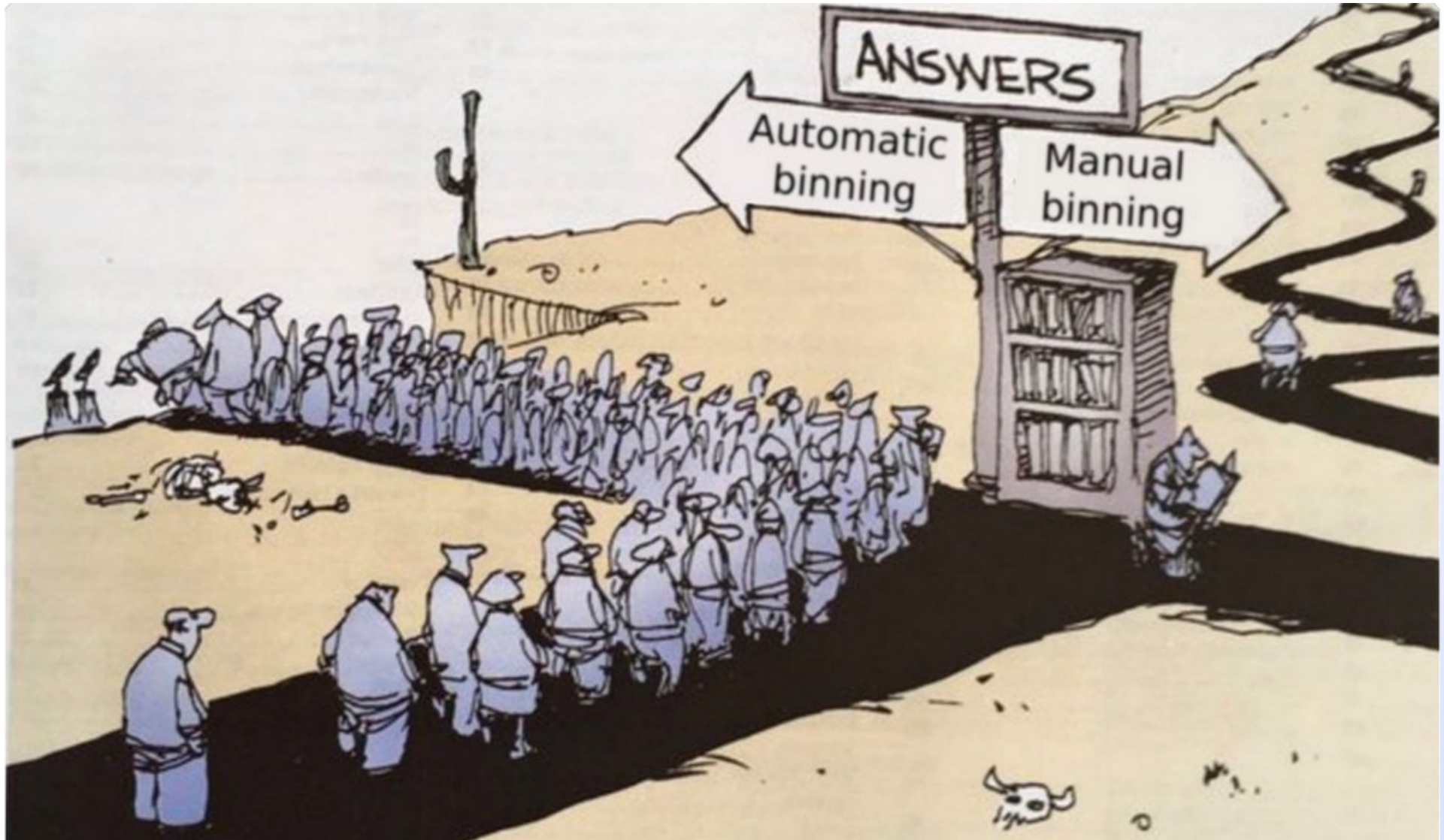# Module V – Metagenomic binning and MAGs
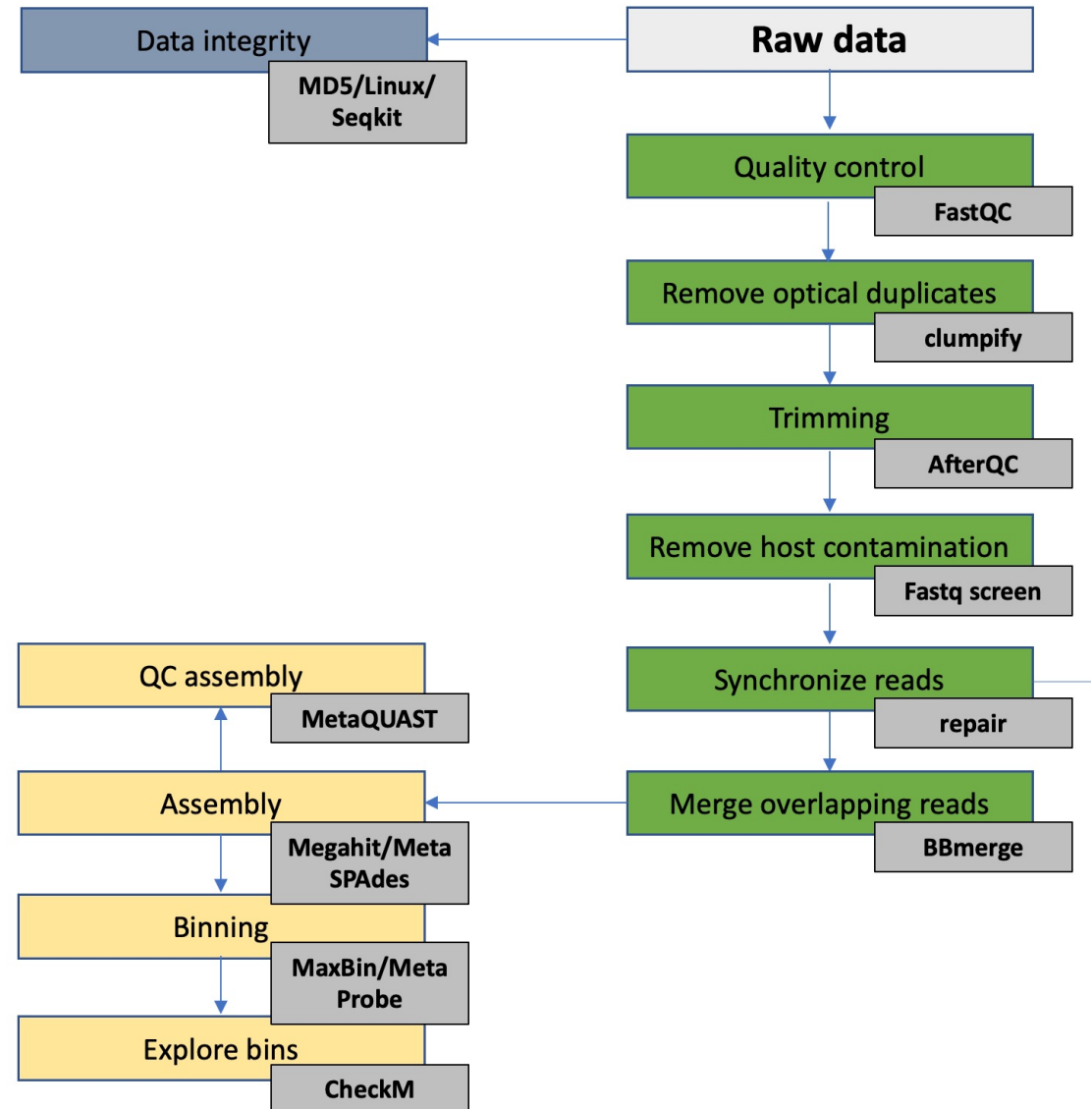


Tom Delmont

# Overview

Binning of contigs

Binning of reads

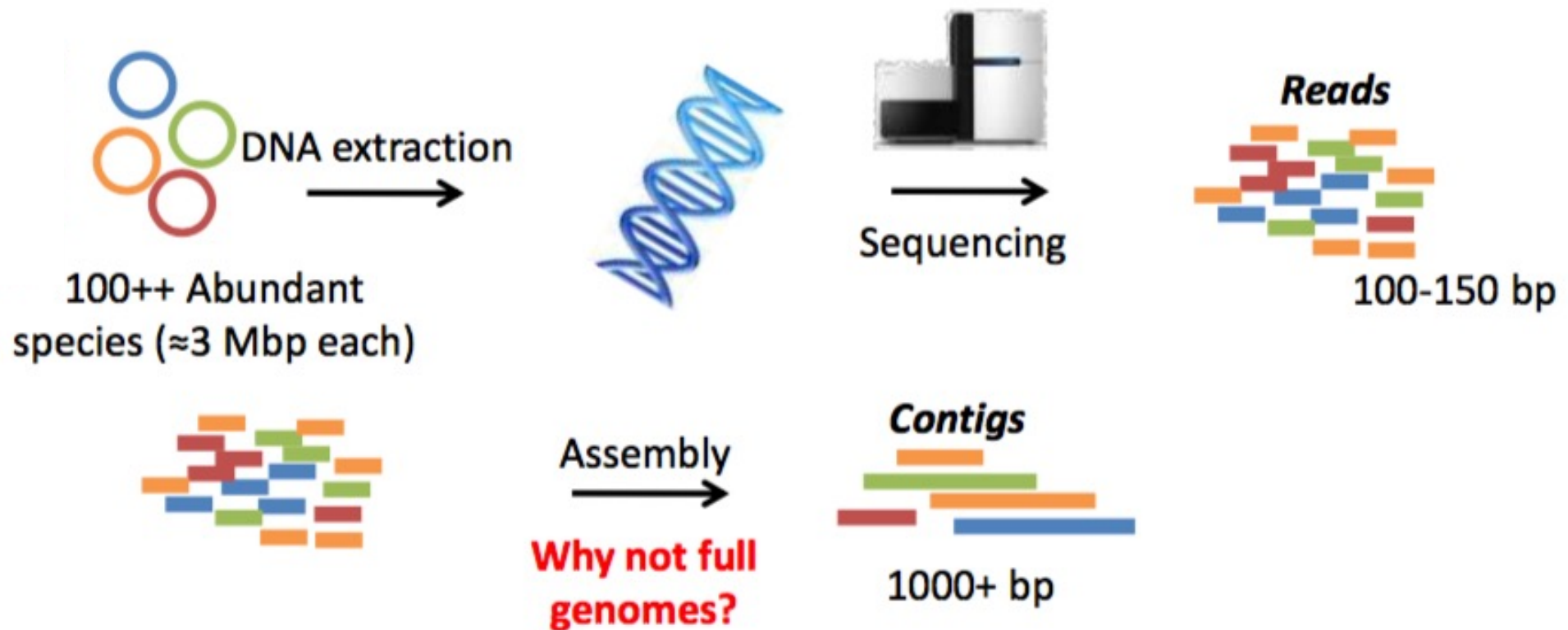Evaluate metagenomic bins

Taxonomic classification of bins

# Recap: obtaining a genome sequence from a metagenomes

Metagenomic Assembled Genomes (MAGs)

Trying to reconstruct the individual genomes of a mixture of DNA from an entire population

Metagenomic assemblies will still be highly fragmented - Binning

# Binning

Method to sort data values into a smaller groups or "bins"

For example to group animals into more taxon-specific bins
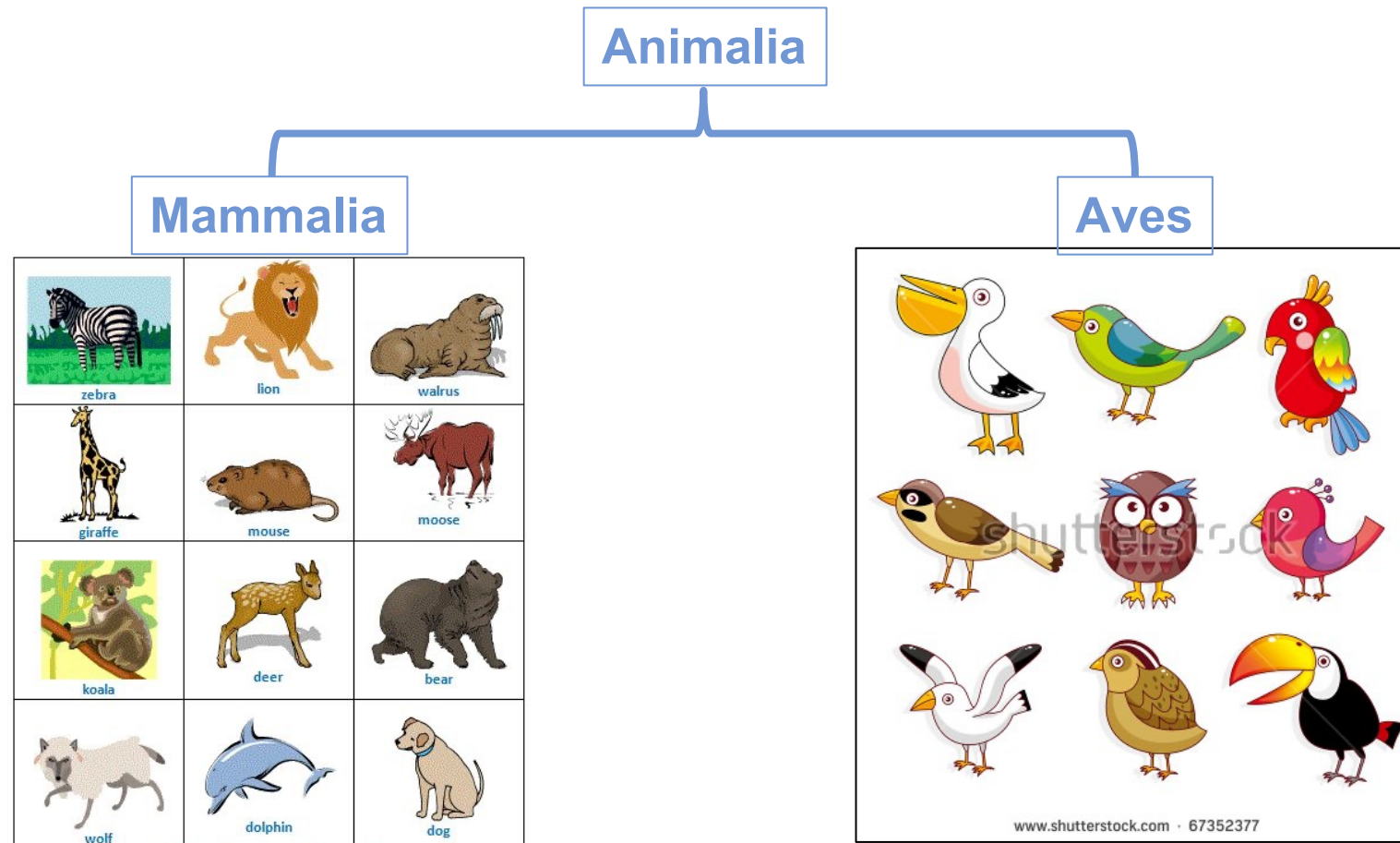


Jungle Animals Real Life Wallpapers Hd

# Binning

Method to sort data values into a smaller groups or "bins"

For example to group animals into more taxon-specific bins

Various taxonomic levels: All belong to Kingdom = *Animalia*, but Class = *Aves* <u>AND</u> *Mammalia*

**Animalia**

**Mammalia**

**Aves**

# Binning - metagenomics

Group contigs or reads belonging to the same specie

Group nucleotide sequences based on **composition**

Group nucleotide sequences based on **abundance**



SlideShare / CoMet: Coverage and Composition based binning of Metagenomes

# Two types of binning strategies

Taxonomy dependent and taxonomy independent strategies

# Binning – Taxonomy independent methods

Referred to as un-supervised

Enables the discovery of new microbial of new organisms

Two types of features used for classification

**Sequence composition based**

- Assumption that the genome composition is unique for each taxon

- DNA fragments from the same genome are more similar than those from different genomes

- Cluster formation being defined by k-mer composition

**Abundance based**

- Coverage reflecting abundance of given tax

- Cluster formation being defined by k-mer abundance
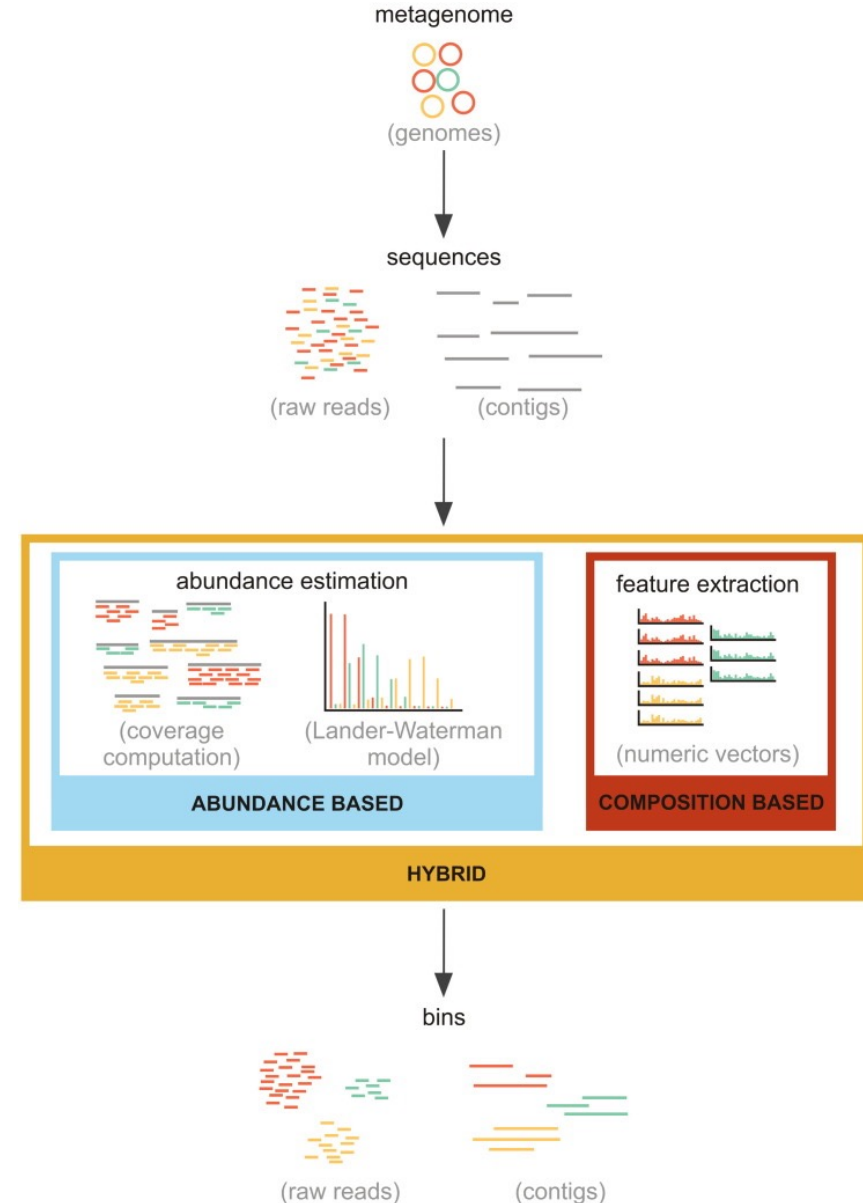
# Binning – Taxonomy independent methods

## Hybrid binning

Combine abundance and composition

Give more accurate binning results

## Can be performed on either sequence reads or assembled contigs

## Potentially separate subspecies into individual bins



K. Sedlar et al. / Computational and Structural Biotechnology Journal 15 (2017)
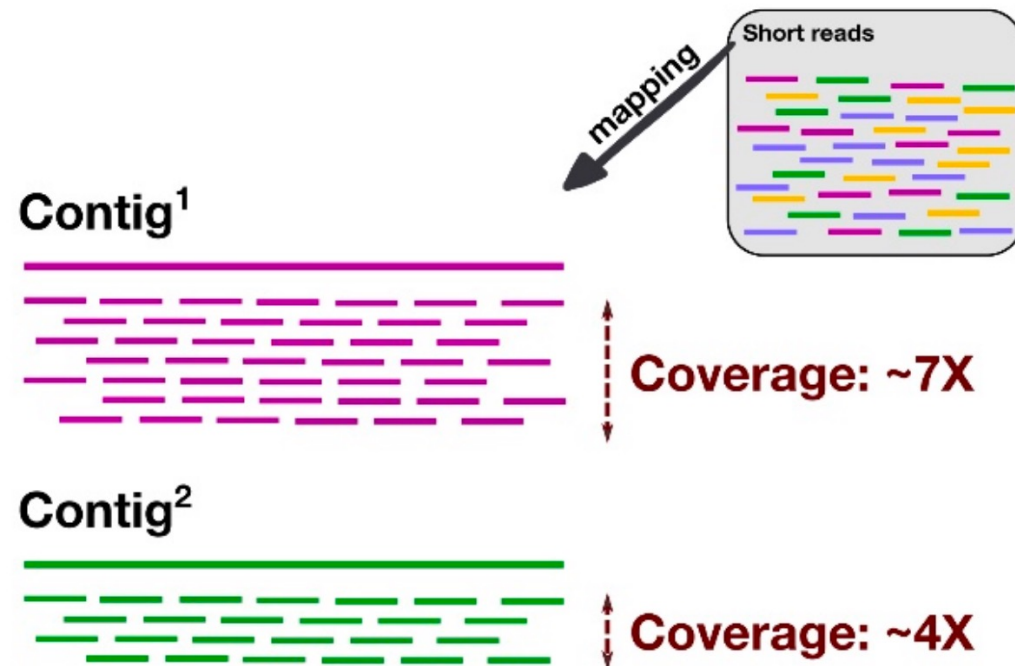
# How are nucleotide sequences binned?

Abundance based binning

Also called coverage based binning

Sequences originating from the same specie will have similar abundance in the sample



SlideShare / A. Murat Eren, Intro to metagenomic binning

# How are nucleotide sequences binned?

## Composition based binning

Genomic signatures have been shown to display a species-specific pattern

GC content is simple and commonly used genomic signature

More widely used genomic signature is tetranucleotide frequencies



SlideShare / A. Murat Eren, Intro to metagenomic binning

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG
       AATT
        ATTC
         TTCC
          TCCG
           CCGG

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG
AATT
  ATTC
    TTCC
      TCCG
        CCGG

|      | AATT | ATTC | TTCC | TCCG | CCGG |
|------|------|------|------|------|------|
| Seq1 |      |      |      |      |      |

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG
  AATT
   ATTC
    TTCC
     TCCG
      CCGG

| | AATT | ATTC | TTCC | TCCG | CCGG |
|---|---|---|---|---|---|
| Seq1 | 1 | 1 | 1 | 1 | 1 |

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG          Seq2: AATTAAGG
         AATT                          AATT
         ATTC                          ATTA
         TTCC                          TTAA
         TCCG                          TAAG
         CCGG                          AAGG

| | AATT | ATTC | TTCC | TCCG | CCGG | ATTA | TTAA | TAAG | AAGG |
|---|---|---|---|---|---|---|---|---|---|
| Seq1 | 1 | 1 | 1 | 1 | 1 | | | | |
| Seq2 | 1 | | | | | 1 | 1 | 1 | 1 |

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG   Seq2: AATTAAGG   Seq3: AAGGAAGG   Seq4: AATTAATT   Seq5: GGAAGGAA

| Seq1 | Seq2 | Seq3 | Seq4 | Seq5 |
|------|------|------|------|------|
| AATT | AATT | AAGG | AATT | GGAA |
| ATTC | ATTA | AGGA | ATTA | GAAG |
| TTCC | TTAA | GGAA | TTAA | AAGG |
| TCCG | TAAG | GAAG | TAAT | AGGA |
| CCGG | AAGG | AAGG | AATT | GGAA |

| | AATT | ATTC | TTCC | TCCG | CCGG | ATTA | TTAA | TAAG | AAGG | TCCC | AGGA | GGAA | GAAG | TAAT |
|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| Seq1 | 1 | 1 | 1 | 1 | 1 | | | | | | | | | |
| Seq2 | 1 | | | | | 1 | 1 | 1 | 1 | | | | | |
| Seq3 | | | | | | | | | 2 | | 1 | 1 | 1 | |
| Seq4 | 2 | | | | | 1 | 1 | | | | | | | 1 |
| Seq5 | | | | | | | | | 1 | | 1 | 2 | 1 | |

# Composition based binning

## Computation of tetranucleotide frequencies (k-mer =4)

| Seq1: AATTCCGG | Seq2: AATTAAGG | Seq3: AAGGAAGG | Seq4: AATTAATT | Seq5: GGAAGGAA |
|---|---|---|---|---|
| AATT | AATT | AAGG | AATT | GGAA |
| ATTC | ATTA | AGGA | ATTA | GAAG |
| TTCC | TTAA | GGAA | TTAA | AAGG |
| TCCG | TAAG | GAAG | TAAT | AGGA |
| CCGG | AAGG | AAGG | AATT | GGAA |

|  | AATT | ATTC | TTCC | TCCG | CCGG | ATTA | TTAA | TAAG | AAGG | TCCC | AGGA | GGAA | GAAG | TAAT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Seq1 | 1 | 1 | 1 | 1 | 1 |  |  |  |  |  |  |  |  |  |
| Seq2 | 1 |  |  |  |  | 1 | 1 | 1 | 1 |  |  |  |  |  |
| Seq3 |  |  |  |  |  |  |  |  | 2 |  | 1 | 1 | 1 |  |
| Seq4 | 2 |  |  |  |  | 1 | 1 |  |  |  |  |  |  | 1 |
| Seq5 |  |  |  |  |  |  |  |  | 1 |  | 1 | 2 | 1 |  |

# Composition based binning

## Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG
AATT
ATTC
TTCC
TCCG
CCGG

Seq2: AATTAAGG
AATT
ATTA
TTAA
TAAG
AAGG

Seq3: AAGGAAGG
AAGG
AGGA
GGAA
GAAG
AAGG

Seq4: AATTAATT
AATT
ATTA
TTAA
TAAT
AATT

Seq5: GGAAGGAA
GGAA
GAAG
AAGG
AGGA
GGAA

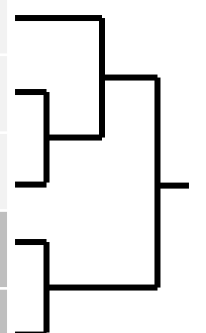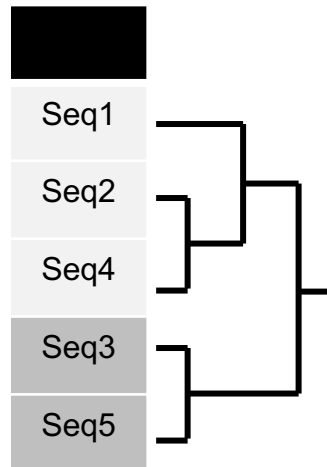|      | AATT | ATTC | TTCC | TCCG | CCGG | ATTA | TTAA | TAAG | AAGG | TCCC | AGGA | GGAA | GAAG | TAAT |
|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| Seq1 | 1    | 1    | 1    | 1    | 1    |      |      |      |      |      |      |      |      |      |
| Seq2 | 1    |      |      |      |      | 1    | 1    | 1    | 1    |      |      |      |      |      |
| Seq4 | 2    |      |      |      |      | 1    | 1    |      |      |      |      |      |      | 1    |
| Seq3 |      |      |      |      |      |      |      |      | 2    |      | 1    | 1    | 1    |      |
| Seq5 |      |      |      |      |      |      |      |      | 1    |      | 1    | 2    | 1    |      |

# Composition based binning

## Computation of tetranucleotide frequencies (k-mer =4)

Seq1: AATTCCGG    Seq2: AATTAAGG    Seq3: AAGGAAGG    Seq4: AATTAATT    Seq5: GGAAGGAA

| Seq1 | Seq2 | Seq3 | Seq4 | Seq5 |
|------|------|------|------|------|
| AATT | AATT | AAGG | AATT | GGAA |
| ATTC | ATTA | AGGA | ATTA | GAAG |
| TTCC | TTAA | GGAA | TTAA | AAGG |
| TCCG | TAAG | GAAG | TAAT | AGGA |
| CCGG | AAGG | AAGG | AATT | GGAA |

|       | AATT | ATTC | TTCC | TCCG | CCGG | ATTA | TTAA | TAAG | AAGG | TCCC | AGGA | GGAA | GAAG | TAAT |
|-------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| Seq1  | 1    | 1    | 1    | 1    | 1    |      |      |      |      |      |      |      |      |      |
| Seq2  | 1    |      |      |      |      | 1    | 1    | 1    | 1    |      |      |      |      |      |
| Seq4  | 2    |      |      |      |      | 1    | 1    |      |      |      |      |      |      | 1    |
| Seq3  |      |      |      |      |      |      |      |      | 2    |      | 1    | 1    | 1    |      |
| Seq5  |      |      |      |      |      |      |      |      | 1    |      | 1    | 2    | 1    |      |

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

| Seq1: AATTCCGG | Seq2: AATTAAGG | Seq3: AAGGAAGG | Seq4: AATTAATT | Seq5: GGAAGGAA |
|---|---|---|---|---|
| AATT | AATT | AAGG | AATT | GGAA |
| ATTC | ATTA | AGGA | ATTA | GAAG |
| TTCC | TTAA | GGAA | TTAA | AAGG |
| TCCG | TAAG | GAAG | TAAT | AGGA |
| CCGG | AAGG | AAGG | AATT | GGAA |

# Composition based binning

Computation of tetranucleotide frequencies (k-mer =4)

| Seq1: AATTCCGG | Seq2: AATTAAGG | Seq3: AAGGAAGG | Seq4: AATTAATT | Seq5: GGAAGGAA |
|---|---|---|---|---|
| AATT | AATT | AAGG | AATT | GGAA |
| ATTC | ATTA | AGGA | ATTA | GAAG |
| TTCC | TTAA | GGAA | TTAA | AAGG |
| TCCG | TAAG | GAAG | TAAT | AGGA |
| CCGG | AAGG | AAGG | AATT | GGAA |

# Composition based binning

Computation of tetranucleotide frequencies (k-mer = 4)

Seq1: AATTCCGG    Seq2: AATTAAGG    Seq3: AAGGAAGG    Seq4: AATTAATT    Seq5: GGAAGGAA

| Seq1 | Seq2 | Seq3 | Seq4 | Seq5 |
|------|------|------|------|------|
| AATT | AATT | AAGG | AATT | GGAA |
| ATTC | ATTA | AGGA | ATTA | GAAG |
| TTCC | TTAA | GGAA | TTAA | AAGG |
| TCCG | TAAG | GAAG | TAAT | AGGA |
| CCGG | AAGG | AAGG | AATT | GGAA |

# Taxonomy independent binning of contigs - MaxBin

Binning of assembled contigs using an expectation-maximization algorithm

Bins are predicted from initial identification of marker genes in assembled sequences

Tetranucleotide frequencies and scaffold coverages are combined to organize metagenomic sequences into individual bins

Estimation of genome completeness – 107 marker genes



Wu et al., Microbiome 2014 2:26

# MaxBin - preformance

Sequencing depths highly affect the results

10-genome simulated datasets - 20X versus 80X coverage

# Taxonomy independent binning of sequence reads- MetaProb

Assembly-assisted tool for binning of reads

Phase 1 groups overlapping reads into groups

Phase 2 builds the probabilistic sequence signatures of independent reads and merges the groups into clusters

# Binning results for the CAMI data sets

Investigated performance when recovering individual genome bins

Large variation:

Average genome completeness (34% to 80%)

Average purity (70% to 97%)



Critical Assessment of Metagenome Interpretation—a benchmark of metagenomics software

Alexander Sczyrba ✉, Peter Hofmann […] Alice C McHardy ✉

# Selecting a binning method

Highly dependent on the sample and the aim of the project and available resources

The length of the metagenomic sequences – key factor

- Ultra-short sequences (75 bp) - assembly step becomes a necessity
- Short length sequences (200-400 bp)- alignment-based or hybrid binning methods
- Long length sequences - alignment-based as well as composition-based binning methods

Are you aiming to identify novel un-culturable species?

Human microbiota

- Most species are known
- Presence or absence of one or several species?

Environmental samples

- Most species are unknown

# Refinement of bins - RefineM

Methods for outlier filtering reduces the total number of contigs being binned
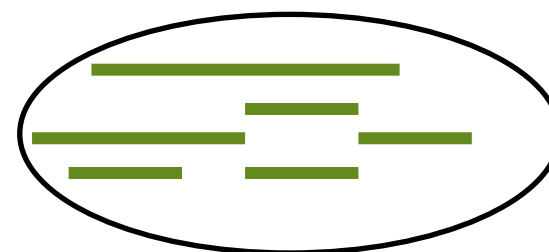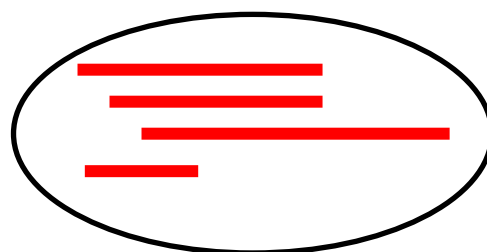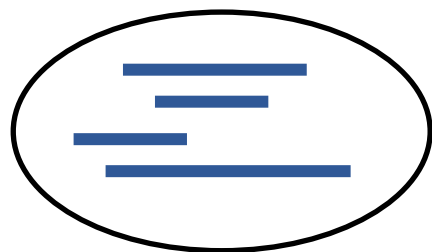
Deviating GC

Deviating tetranucleotide composition

Deviating coverage depth

Identify contigs with a coding density suggestive of a Eukaryotic origin



Refinement of bins

# Estimate completeness and contamination of MAGs

Assembly statistics

Total size of MAG (sum of contigs in bin)

Contig size (N50 value)

Presence and absence of lineage-specific genes
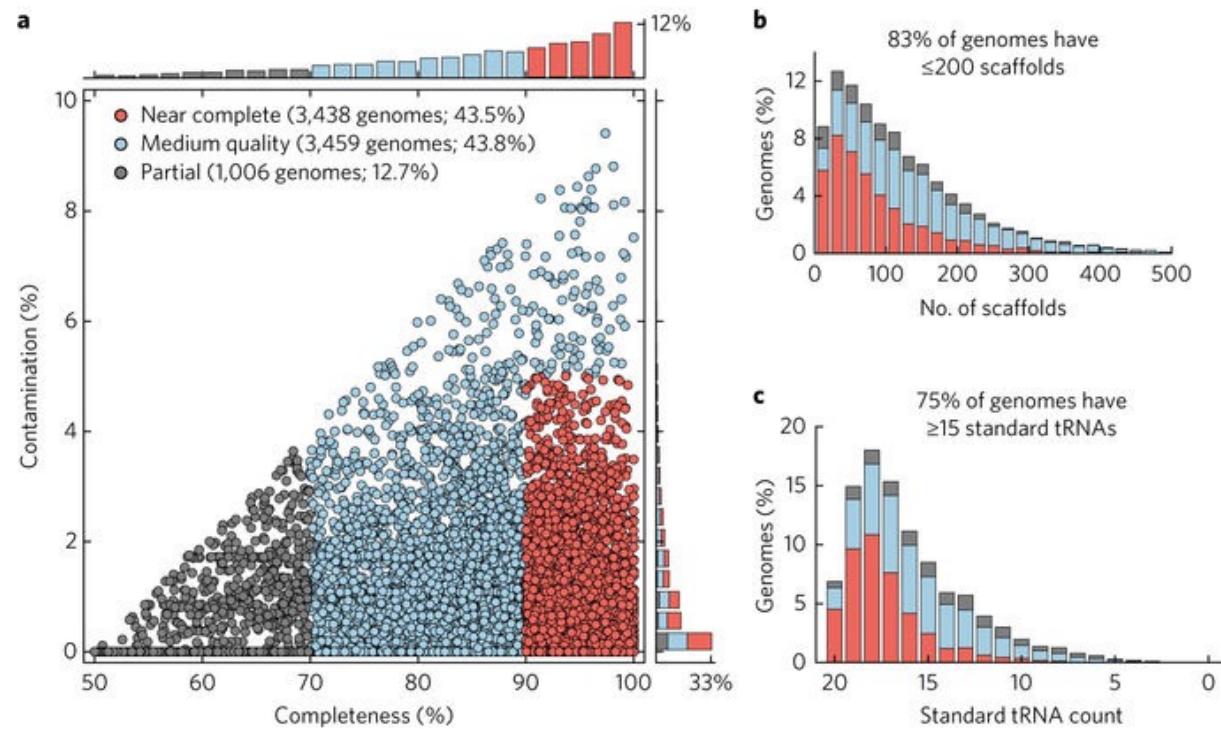
Presence of 20 standard tRNAs

# Estimate completeness and contamination of MAGs

CheckM assess the quality of genomes recovered from metagenomes

Estimate genome completeness and contamination

Using collocated single-copy marker genes within a phylogenetic lineage
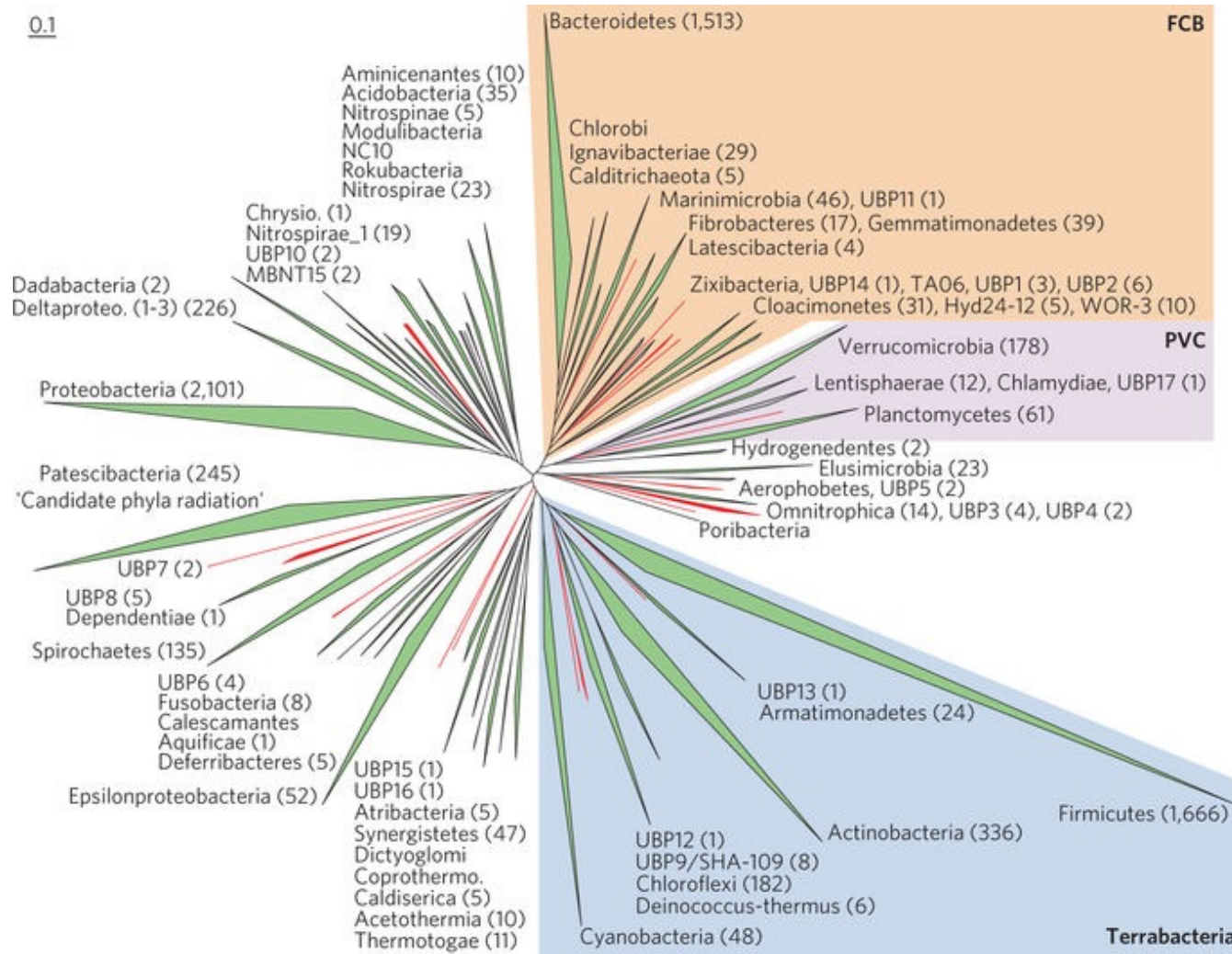
- Bacteria: 104 markers organized into 58 sets

- Archaea: 150 markers organized into 108 sets



Parks et al. / *Nat Microbiol.*
*2017 Nov;2(11):1533-1542.*

# Example of important outcome from MAGs

Analysed 1500 metagenomics datasets

First genomes from 17 bacterial phyla and 3 archaeal phyla

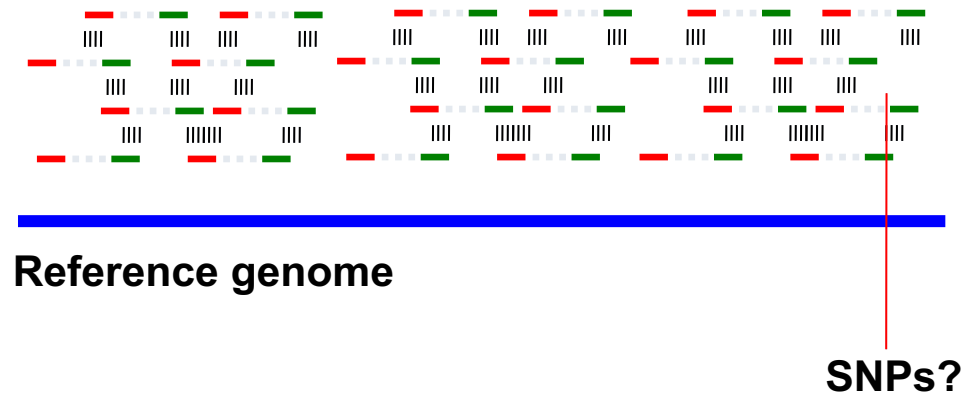**Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life**

Donovan H. Parks, Christian Rinke, Maria Chuvochina, Pierre-Alain Chaumeil, Ben J. Woodcroft, Paul N. Evans, Philip Hugenholtz* and Gene W. Tyson*
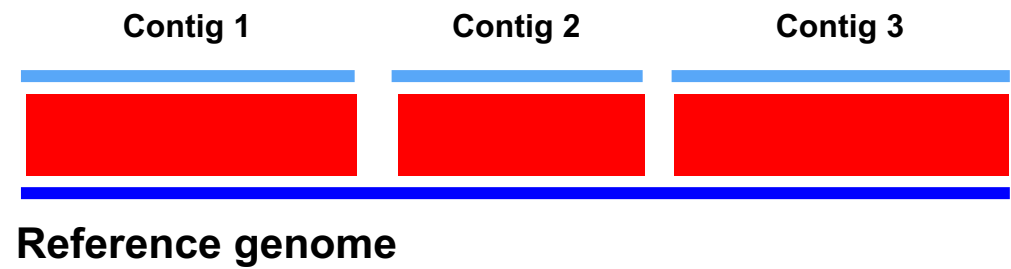
# Two slides on mapping against a reference sequence

Two methods for mapping against a reference sequence

Aligning reads against reference

Aligning contigs against reference

**Reference genome**

**SNPs?**

**Contig 1**  **Contig 2**  **Contig 3**

**Reference genome**

# Overall read mapping process: