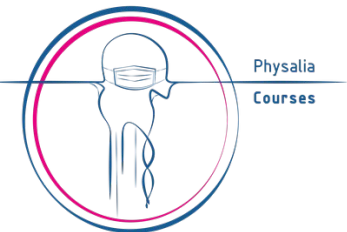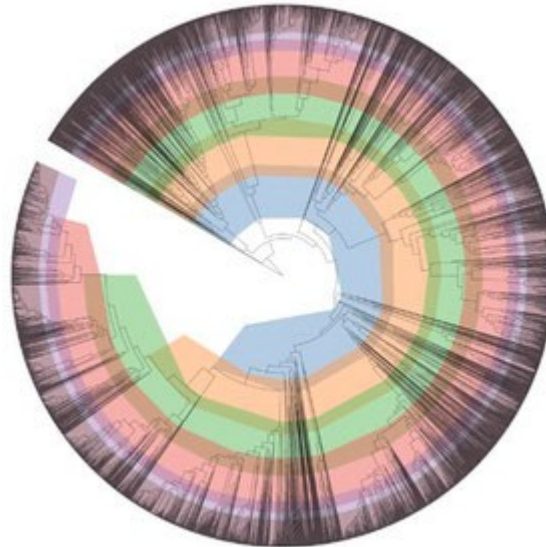# ENVIRONMENTAL METAGENOMICS

Physalia course, online, 13-17 October 2025

## Metagenome de-novo assembly and quality control

Nikolay Oskolkov, Group Leader of Metabolic Research Group at LIOS, Riga, Latvia
Samuel Aroney, Postdoctoral Research Fellow, Queensland University of Technology



Physalia
Courses

# Typical analysis methods used in metagenomics

## 1) Alignment:



## 2) Classification:



Centrifuge

MetaPhlan

Clark

Reference based:
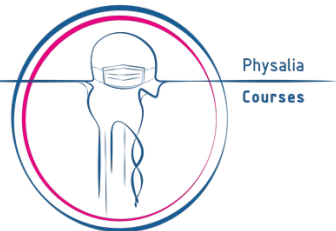assume similarity to reference

## 3) De-novo assembly:



```
>seq1
GCCGTAGTCC...
>seq2
...
```

*Binning-based analysis*

genome_A

genome_B

genome_C

Assembly     gene prediction/ annotation

*Assembly-based analysis*

Phylogenetic binning

Reference free:
unbiased but challenging

# Strenghts and weaknesses of read-based metagenomics

**Comprehensiveness**
Can provide an aggregate picture of community function or structure, but is based only on the fraction of reads that map effectively to reference dbs

**Community complexity**
Can deal with communities of arbitrary complexity given sufficient sequencing depth and satisfactory reference database coverage

**Novelty**
Cannot resolve organisms for which genomes of close relatives are unknown

**Computational burden**
Can be performed efficiently, enabling large meta-analyses

**Genome-resolved metabolism**
Can typically resolve only the aggregate metabolism of the community, and links with phylogeny are only possible in the context of known reference genomes

**Expert manual supervision**
Usually does not require manual curation, but selection of reference genomes to use could involve human supervision

**Integration with microbial genomics**
Obtained profiles cannot be directly put into the context of genomes derived from pure cultured isolates

Physalia
Courses

# Summary metagenome analysis strategies

**Read-based analysis**

Functional/taxonomic annotation directly on reads (often partial genes)

```
>seq1
GCCGTAGTCC...
>seq2
...
```

Assembly

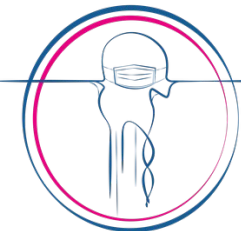gene prediction/ annotation

**Binning-based analysis**

genome$_A$
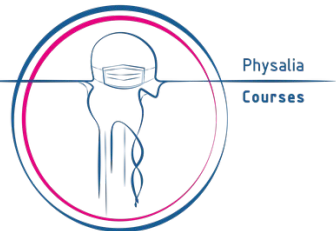
genome$_B$

genome$_C$

**Assembly-based analysis**

Phylogenetic binning

Physalia Courses

# Strengths and weaknesses of assembly-based metagenomics

**Comprehensiveness**

Can construct multiple whole genomes, but only for organisms with enough coverage to be assembled and binned

**Community complexity**

In complex communities, only a fraction of the genomes can be resolved by assembly

**Novelty**

Can resolve genomes of entirely novel organisms with no sequenced relatives

**Computational burden**

Requires computationally costly assembly, mapping and binning

**Genome-resolved metabolism**

Can link metabolism to phylogeny through completely assembled genomes, even for novel diversity

**Expert manual supervision**

Manual curation required for accurate binning and scaffolding and for misassembly detection

**Integration with microbial genomics**

Assemblies can be fed into microbial genomic pipelines designed for analysis of genomes from pure cultured isolates
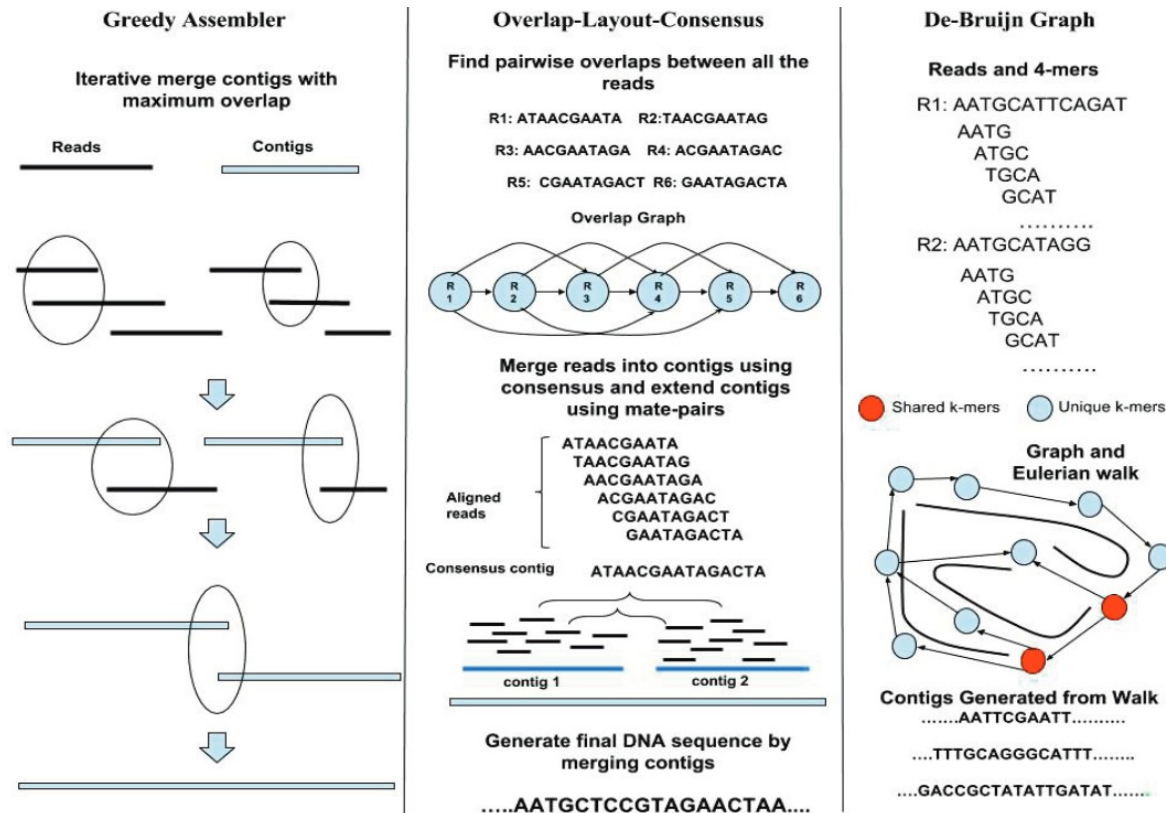
Physalia
Courses

# De novo assembly

Assemble short nucleotide sequences into longer sequences by finding their overlap / consensus without reference genome

- No reference available
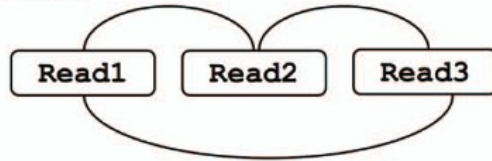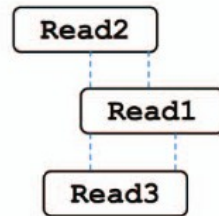- Uneven and complex communities



Assembly

Metagenomic binning

Metagenome assembled genomes

Physalia
Courses

**Greedy Assembler**

Iterative merge contigs with maximum overlap

Reads          Contigs

**Overlap-Layout-Consensus**

Find pairwise overlaps between all the reads

R1: ATAACGAATA    R2: TAACGAATAG

R3: AACGAATAGA    R4: ACGAATAGAC

R5: CGAATAGACT    R6: GAATAGACTA

Overlap Graph

Merge reads into contigs using consensus and extend contigs using mate-pairs

Aligned reads
ATAACGAATA
TAACGAATAG
AACGAATAGA
ACGAATAGAC
CGAATAGACT
GAATAGACTA

Consensus contig    ATAACGAATAGACTA

contig 1        contig 2

Generate final DNA sequence by merging contigs

.....AATGCTCCGTAGAACTAA....

**De-Bruijn Graph**

Reads and 4-mers

R1: AATGCATTCAGAT
AATG
 ATGC
  TGCA
   GCAT
..........
R2: AATGCATAGG
AATG
 ATGC
  TGCA
   GCAT
..........

● Shared k-mers    ○ Unique k-mers

Graph and Eulerian walk

Contigs Generated from Walk
.......AATTCGAATT..........

....TTTGCAGGGCATTT........

....GACCGCTATATTGATAT.......

**Figure 1: Overview of different de novo assembly paradigms.** Schematic representation of the three main paradigms for genome assembly – Greedy, Overlap-Layout-Consensus, and de Bruijn. In Greedy assembler, reads with maximum overlaps are iteratively merged into contigs. In Overlap-Layout-Consensus approach, a graph is constructed by finding overlaps between all pairs of reads. This graph is further simplified and contigs are constructed by finding branch-less paths in the graph, and taking the consensus sequence of the overlapping reads implied by the corresponding paths. Contigs are further organized and extended using mate pair information. In de Bruijn graph assemblers, reads are chopped into short overlapping segments (k-mers) which are organized in a de Bruijn graph structure based on their co-occurrence across reads. The graph is simplified to remove artifacts due to sequencing errors, and branch-less paths are reported as contigs.

Ghurye et al., Yale Journal of Biology and Medicine, 2016

Physalia
Courses

**(a) Overlap, Layout, Consensus assembly**

**(i) Find overlaps**

Read1 Read2 Read3

**(ii) Layout reads**

Read2
Read1
Read3

**(iii) Build consensus**

CGATTCTA
TTCTAAGT
GATT**G**TAA
―――――――――
CGATTCTAAGT

**(b) De Bruijn graph assembly**

**(i) Make kmers**

| Read1: TTCTAAGT | Read2: CGATTCTA | Read3: GATT**G**TAA |
|---|---|---|
| Kmers: TTC | Kmers: CGA | Kmers: GAT |
| TCT | GAT | ATT |
| CTA | ATT | **TTG** |
| TAA | TTC | **TGT** |
| AAG | TCT | **GTA** |
| AGT | CTA | TAA |

**(ii) Build graph**

CGA →T GAT →T ATT →C TTC →T TCT →A CTA →A TAA →G AAG →T AGT
ATT →G TTG →T TGT →A GTA →A TAA

**(iii) Walk graph and output contigs**

CGA →T GAT →T ATT →C TTC →T TCT →A CTA →A TAA →G AAG →T AGT

CGATTCTAAGT

**Figure 1.** Two different approaches to genome assembly: **(a)** in Overlap, Layout, Consensus assembly, (i) overlaps are found between reads and an overlap graph constructed (edges indicate overlapping reads). (ii) Reads are laid out into contigs based on the overlaps (dashed lines indicate overlapping portions). (iii) The most likely sequence is chosen to construct consensus sequence. **(b)** In dBg assembly, (i) reads are decomposed into kmers by sliding a window of size *k* across the reads. (ii) The kmers become vertices in the dBg, with edges connecting overlapping kmers. Polymorphisms (red) form branches in the graph. A count is kept of how many times a kmer is seen, shown here as numbers above kmers. (iii) Contigs are built by walking the graph from edge nodes. A variety of heuristics handle branches in the graphs—for example, low coverage paths, as shown here, may be ignored.

SOAPdenovo
Velvet
Spades

# Popular de Bruijn graph *de-novo* metagenomic assemblers for short Illumina reads
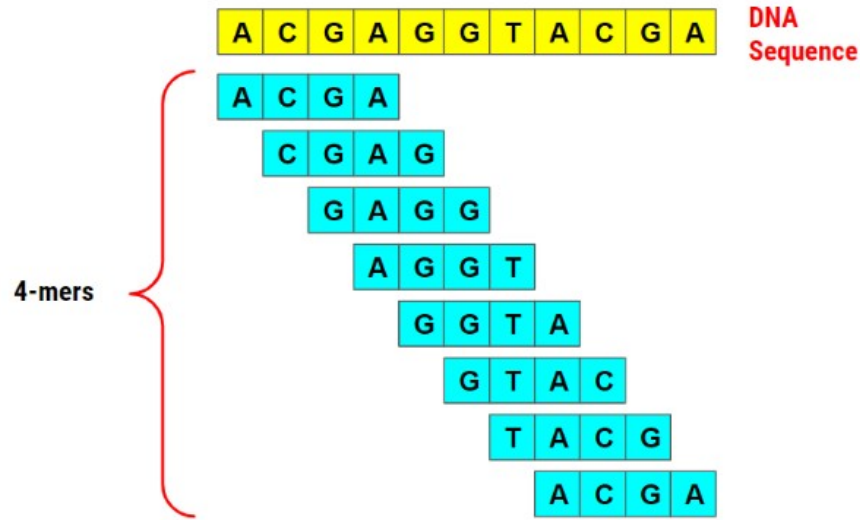


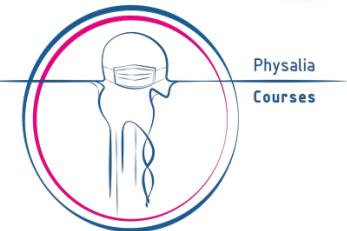We are going to use Megahit in the exercises

# *De novo* assembly using MEGAHIT

**MEGAHIT**: de Bruijn-graph assembler using a distribution of different k-mer lengths inferred from the length of the sequencing data

| | | | | | | | | | | | | **DNA** |
|A|C|G|A|G|G|T|A|C|G|A| **Sequence** |

4-mers

A C G A
C G A G
G A G G
A G G T
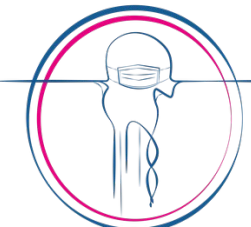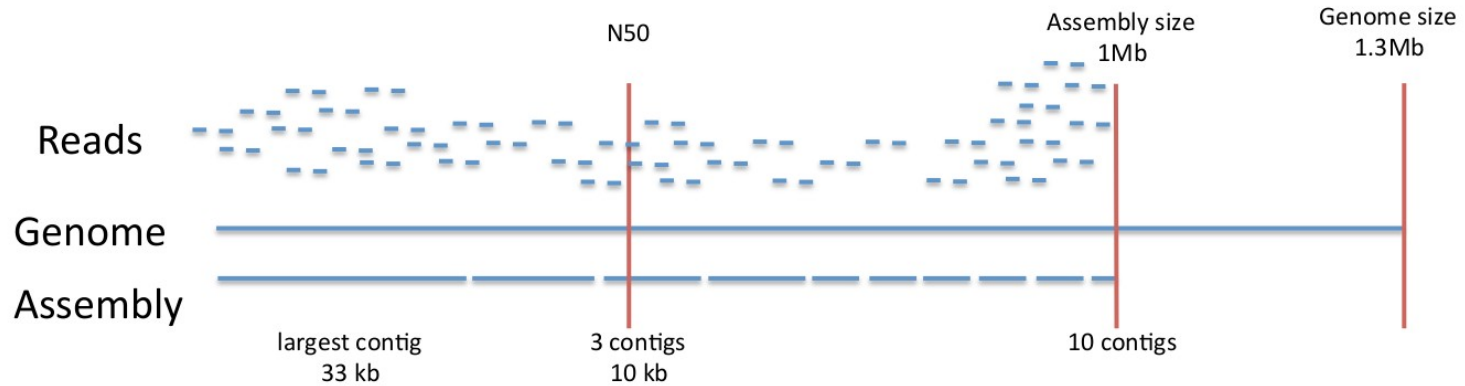G G T A
G T A C
T A C G
A C G A

reasons for using MEGAHIT:
- **low-memory** footprint
- has little issues with the **presence of ancient DNA damage**
- works with **single-end data**

BUT: lower assembly quality than other assemblers for modern sequencing data (see CAMI II challenge; DOI: 10.1038/s41592-022-01431-4)

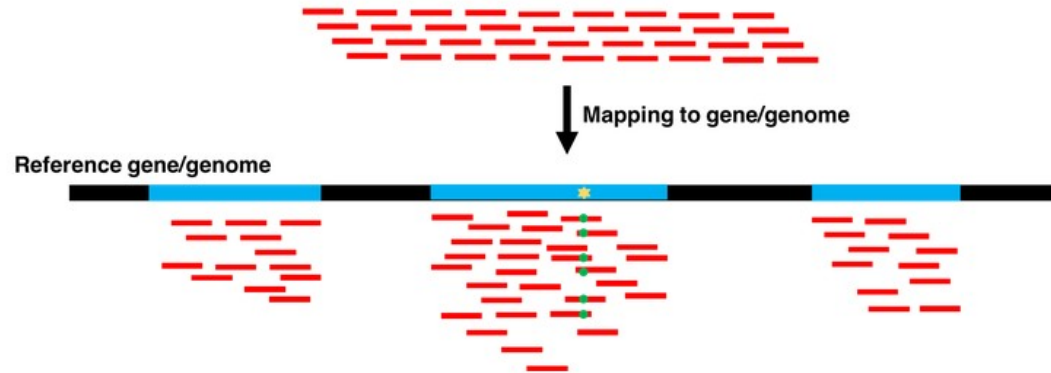Physalia
**Courses**

# Assembly metrics

- assembly size
- number of contigs, largest contig
- N50

# Alignment against the contigs

Many of the following steps require the alignment of the short-read data against the de novo assembled contigs, e.g.

- correction of the contig sequences
- binning of the contigs into MAGs (coverage along the contigs)
- quantification of the presence of ancient DNA damage



Mapping to gene/genome

Reference gene/genome

Physalia
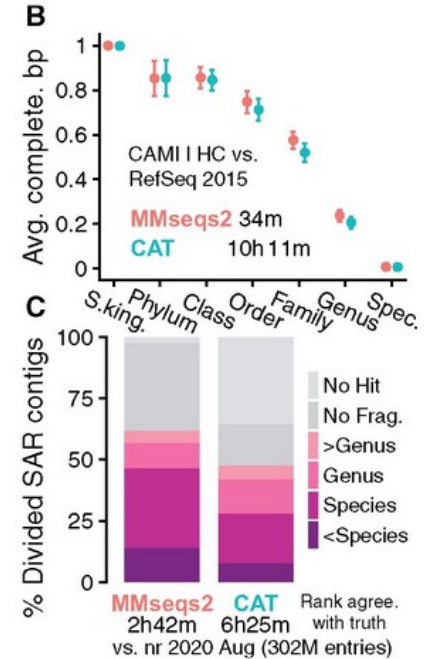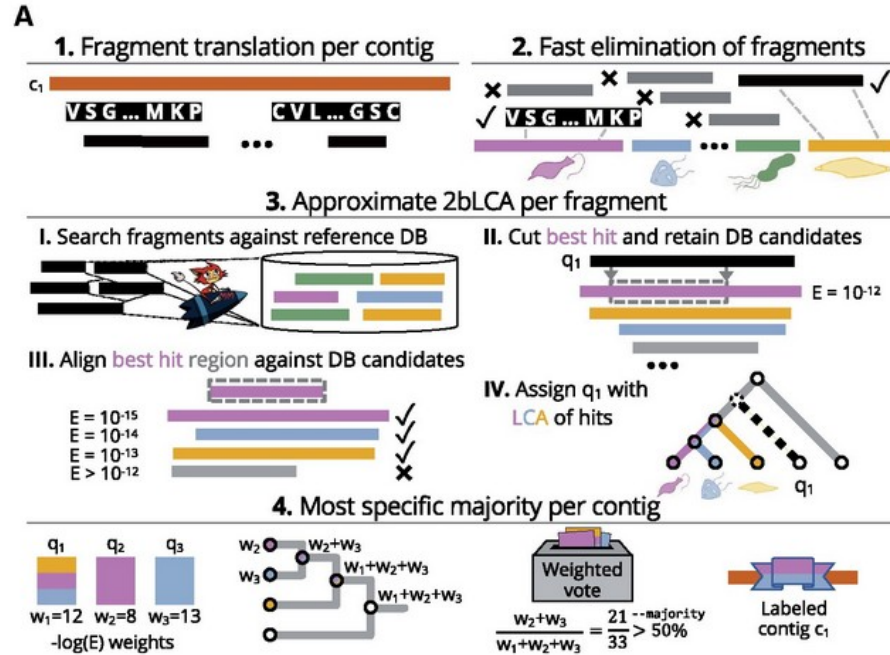Courses

# Taxonomic classification - on contig level

The likely taxonomic origin of contigs can be determined by aligning them against a reference database.

**available aligners:**
- BLAST/DIAMOND
- Kraken2
- Centrifuge
- **MMSeqs2**

**available databases:**
- NCBI NT/RefSeq
- **GTDB**



Mirdita *et al.* (Bioinformatics, 2021; doi: bioinformatics/btab184) Fig. 1