# Day 5: Genome annotation

## MMB-114

# Schedule

**Day 1:** Basics of UNIX and working with the command line
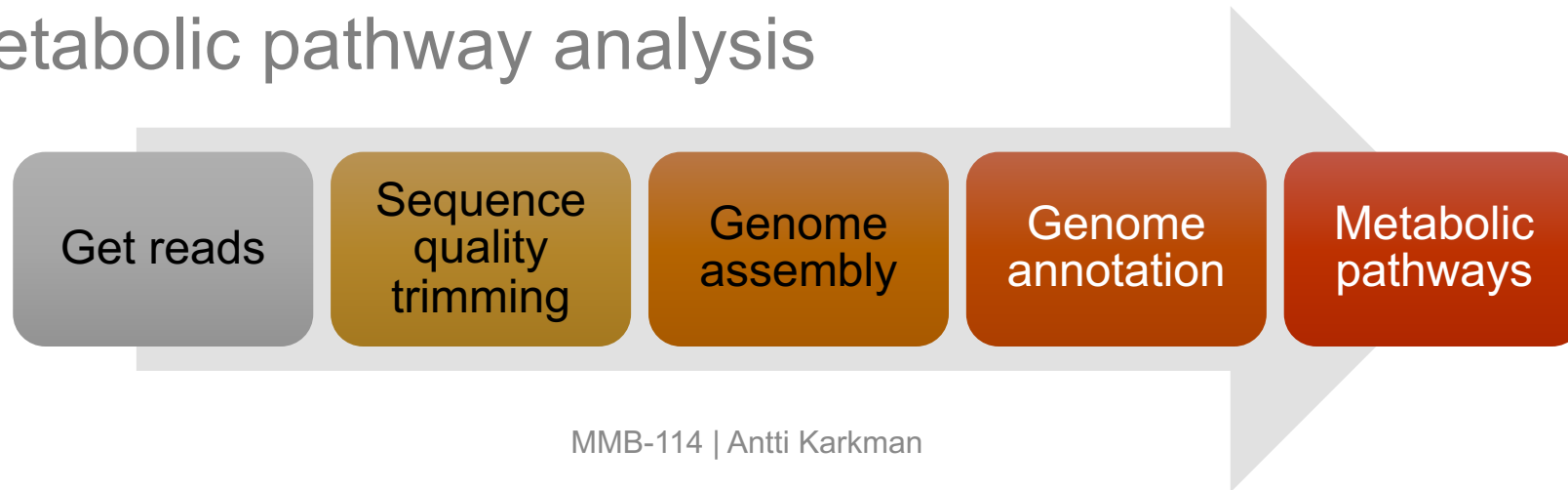
**Day 2:** Handling of Nanopore/Illumina data

**Day 3:** Check-up

**Day 4:** Genome assembly

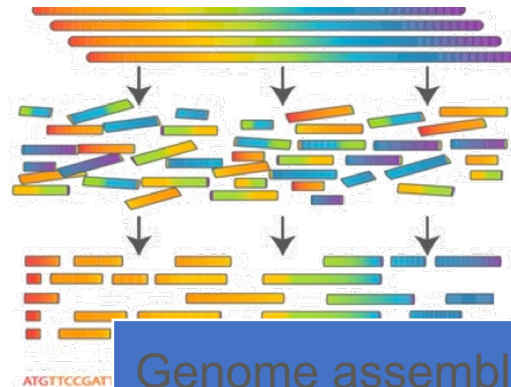**Day 5:** Genome annotation

**Day 6:** Metabolic pathway analysis

Get reads → Sequence quality trimming → Genome assembly → Genome annotation → Metabolic pathways

# Quick recap


Isolation


DNA extraction


Sequencing


Quality control


Genome assembly


Genome annotation

# Gene annotation

- Adding biological information to sequences
- There is a gene X in contig Y on location Z
  - Size of the gene
  - Name of the gene
  - Function of the gene (protein / RNA gene)

| | *nifH* | | *nifD* | | *nifK* | | *nifE* | | *nifN* | | *nifB* | |

Contig Y = 20 035 bp

# Ways to identify protein coding genes

- Sequence alignments
    - E.g. BLAST
    - Search contigs against a database
    - Computationally (and manually) intensive

- Gene finding
    - Start codon (ATG)
    - Open reading frame (ORF)
    - Stop codon (TAA, TAG, TGA)

```
1. ATG CAA TGG GGA AAT GTT ACC AGG TCC GAA CTT ATT GAG GTA AGA CAG ATT TAA
2. A TGC AAT GGG GAA ATG TTA CCA GGT CCG AAC TTA TTG AGG TAA GAC AGA TTT AA
3. AT GCA ATG GGG AAA TGT TAC CAG GTC CGA ACT TAT TGA GGT AAG ACA GAT TTA A
```

# Functions to genes

- Homology
- Statictical modelling of protein families/domains
- Annotation databases
  - NCBI
  - KEGG
  - COG
  - SEED
  - GO
  - UNIPROT
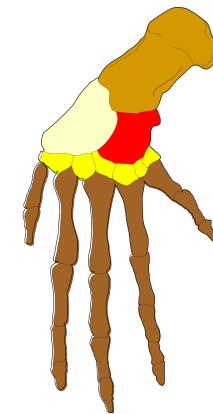  - INTERPRO
  - PFAM
  - TIGR
  - …



Human          Dog          Bird          Whale

# BAKTA

- rapid and standardized annotation of bacterial genomes via alignment-free sequence identification

# Taxonomy and completeness of your genome

**CheckM2**

- **Predicts genome completeness and contamination based on ML model**

- **Designed for metagenome-assembled genomes (MAGs)**

**GTDB-Tk**

- **The Genome Taxonomy Database Toolkit**

- **Taxonomic assignment based on GTDB**

- **Domain-specific concatenated protein reference trees**

# Material and methods examples

**Nanopore:**

- https://www.nature.com/articles/s41587-020-0422-6#Sec2
- https://journals.asm.org/doi/10.1128/msystems.00491-22

**Illumina:**

- https://journals.asm.org/doi/10.1128/msphere.00538-22#sec-4
- https://link.springer.com/article/10.1186/s40793-022-00424-2#Sec2

**Sanger:**

"Sanger sequencing was done on the purified products using BigDye v3.1 Chemistry and primers XXX and XXX (see Supplement Table 2) and analyzed on an ABI3130xl Capillary Sequencer (Thermo, Life Technologies). The obtained sequences were edited…" https://doi.org/10.1002/jmv.27418

MMB-114 | Antti Karkman

# Let's annotate your genome

Go to Github and follow the instructions:

[https://github.com/karkman/MMB-114_Genomics](https://github.com/karkman/MMB-114_Genomics)

(**Day 5:** Genome annotation)