# Genome assembly and annotation

## Day 5: Genome annotation

**Antti Karkman**

Department of Microbiology – UH

antti.karkman@helsinki.fi

# Aims for this part of MMB-114

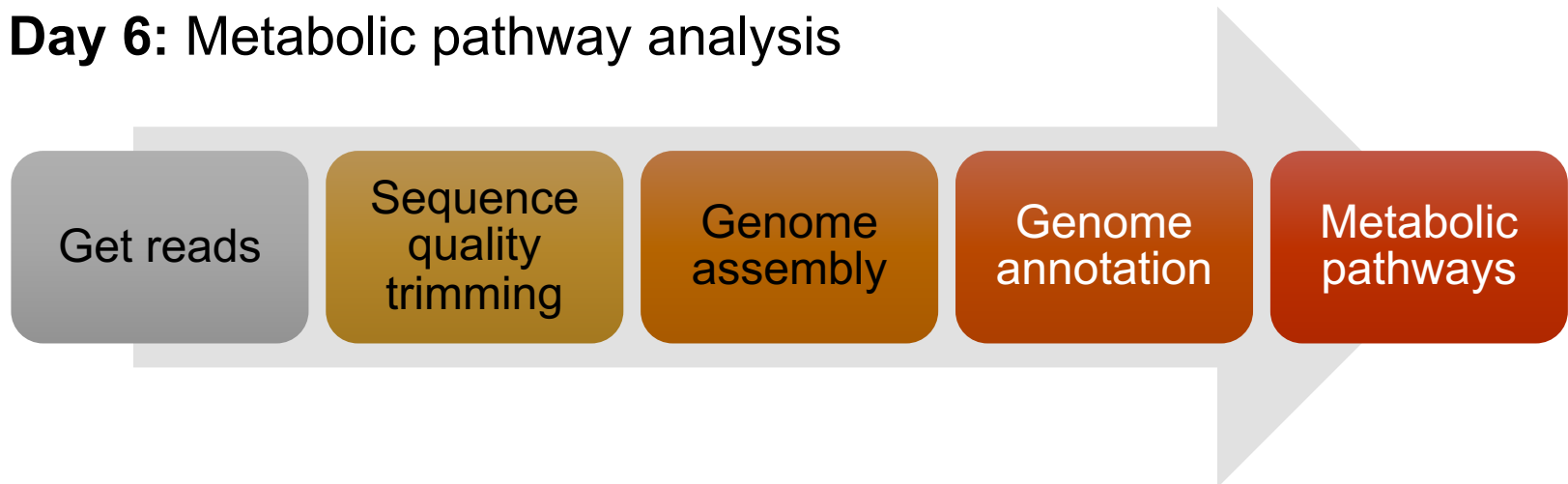**Day 1:** Basics of UNIX and working with the command line

**Day 2:** Handling of Illumina data

**Day 3:** Genome assembly

**Day 4:** Check-up and report

**Day 5:** Genome annotation
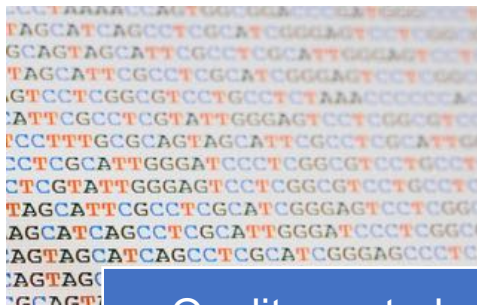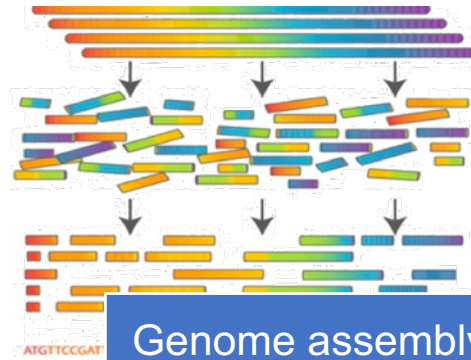
**Day 6:** Metabolic pathway analysis

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

# Recap from last week:

Isolation

DNA extraction

Sequencing

Quality control

Genome assembly

Genome annotation

# Annotation

Adding biological information to sequences (contigs)

Information that there is a gene x in contig y at location z

- Size of the gene
- Name of the gene
- Protein product
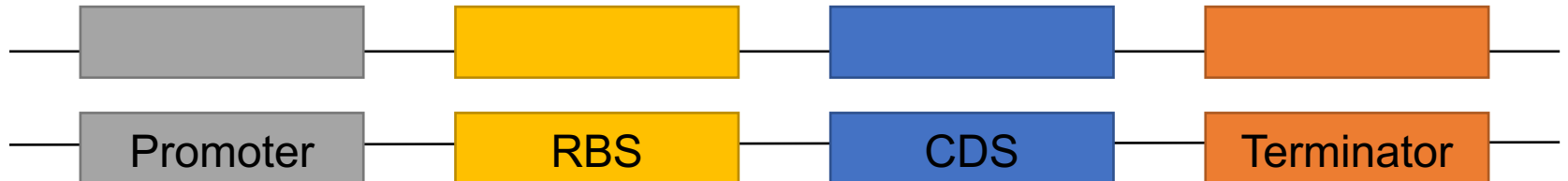


Contig y = 2,035 bp

# Bacterial genes

Promoter

Ribosome binding site (RBS)

Coding sequence (CDS)

Terminator

Also non-coding genes
- tRNA
- rRNA

# Two ways to identify protein-coding genes

**Sequence alignment (e.g. BLAST)**

Search contigs against a database

Computationally-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
```
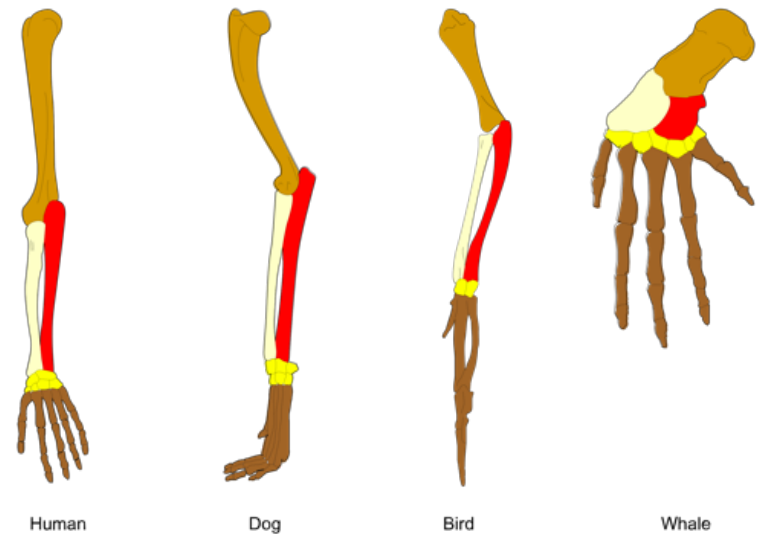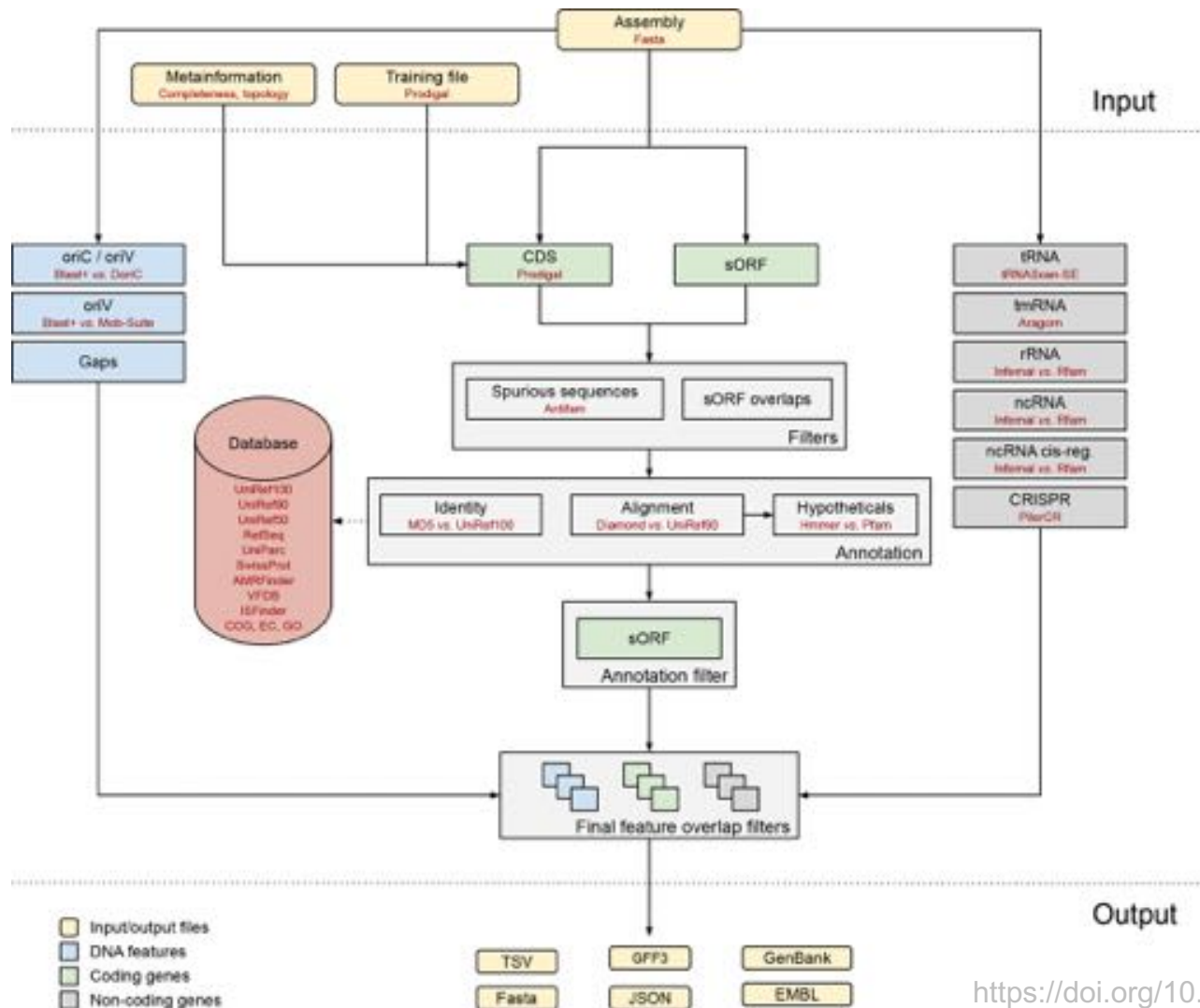
# Annotating genes

Homology

Statistical modelling of protein families/domains

Annotated 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

# Let's annotate our genome

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

(**Day 5:** Genome annotation)