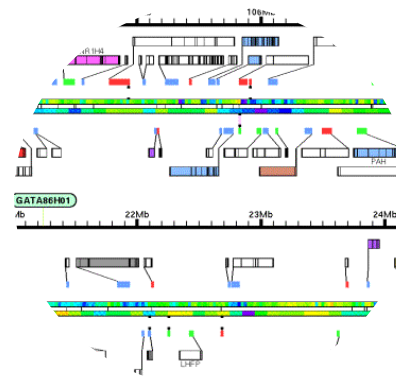
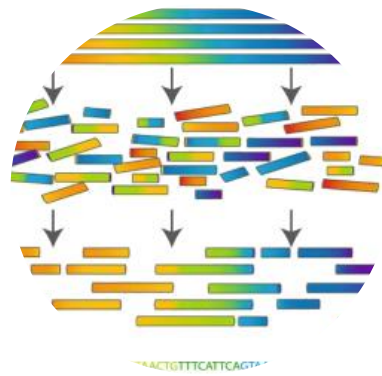
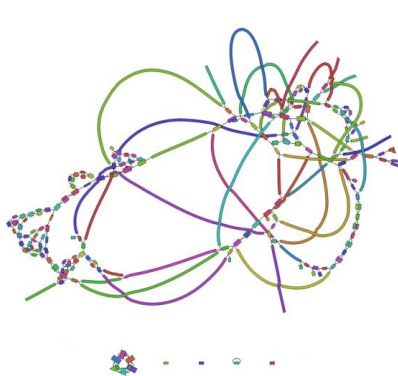


# **MASTER IN BIOPHARMACEUTICAL SCIENCES**

## **Bacterial Molecular Genetics**

**LISBON 2024**

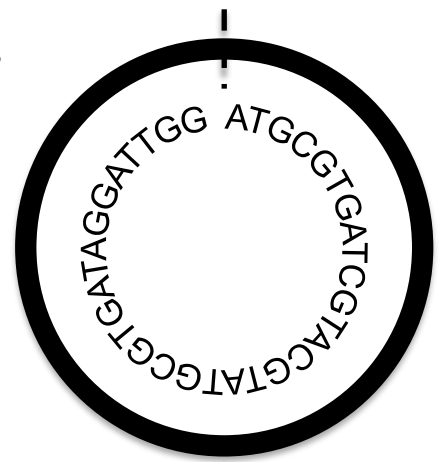
### *De novo Assembly*



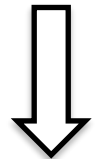
João Perdigão

# Whole Genome Sequencing... the Illumina way

Ideally...

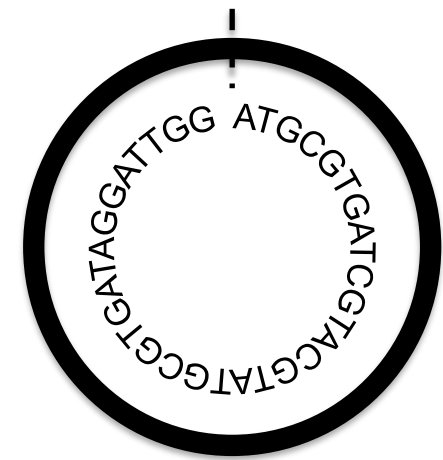


Sequencing

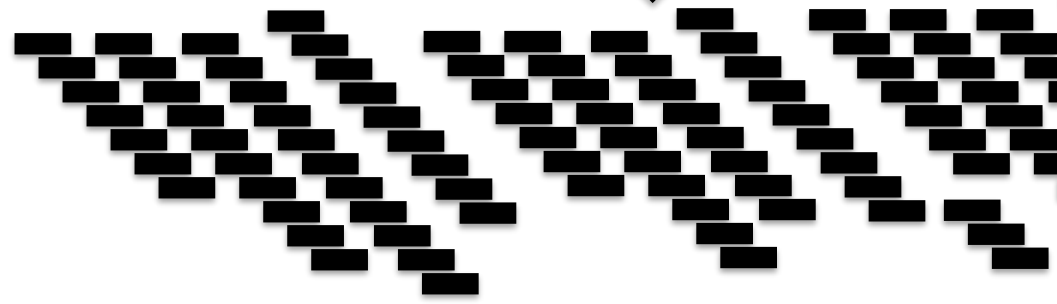


ATGCGTGATCGTACGTATGCGTGATAGGATTGG

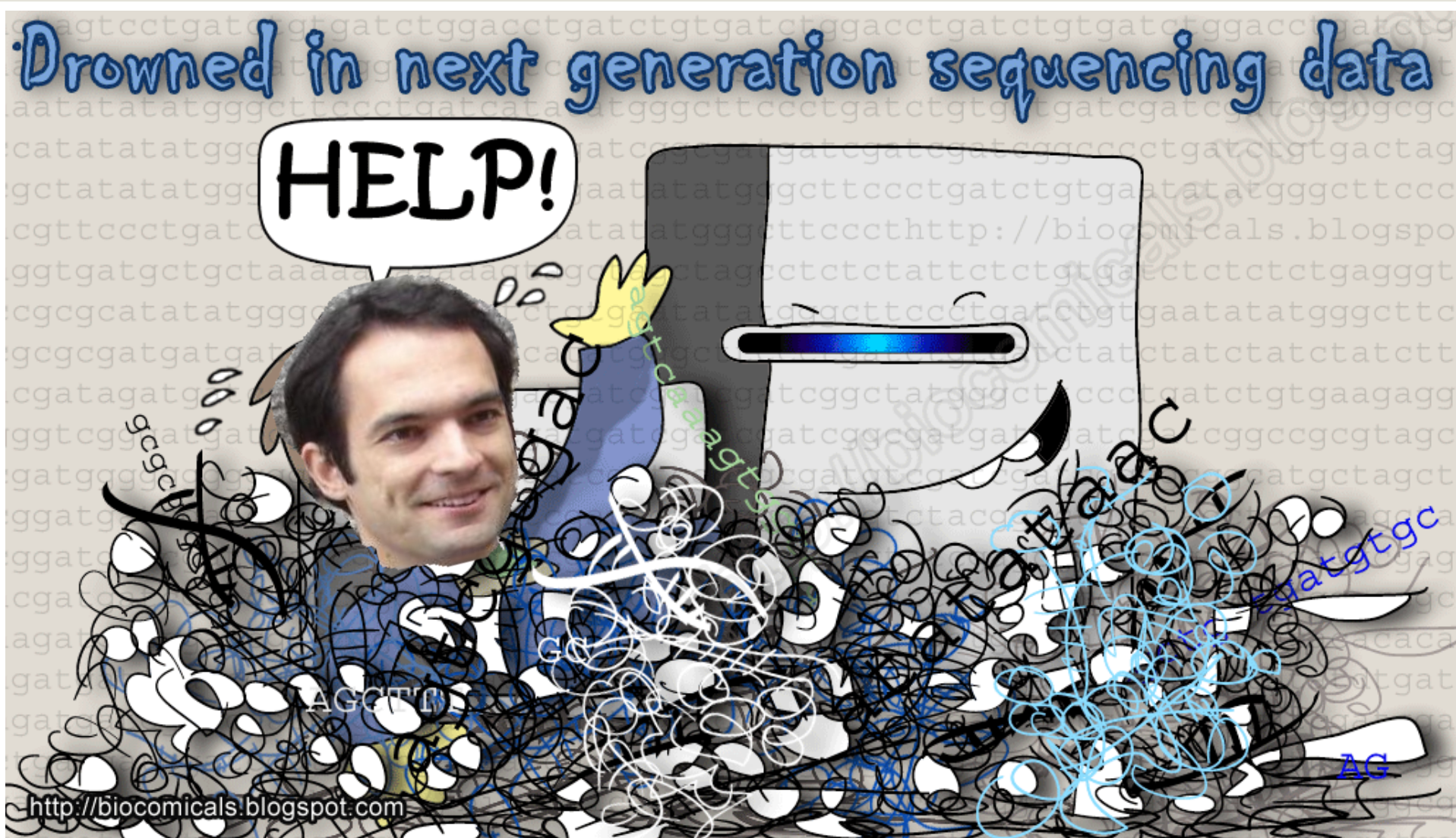
The harsh reality...



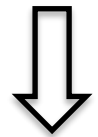
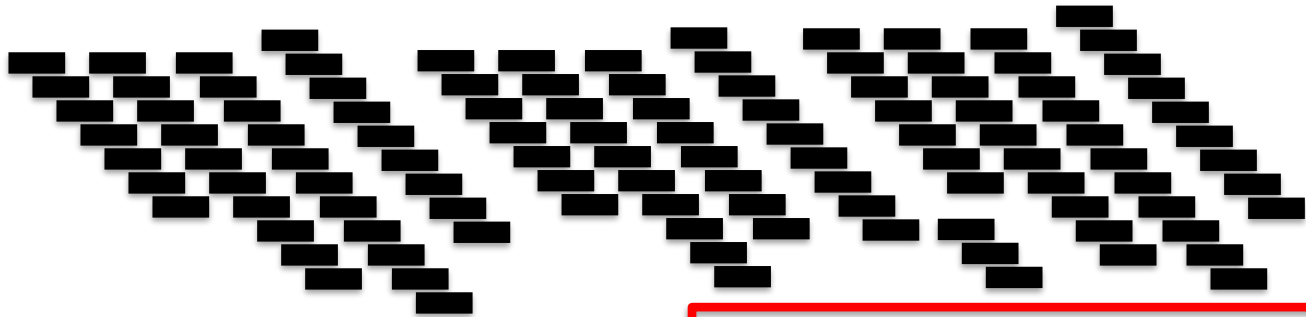
Sequencing



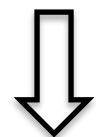
# Whole Genome Sequencing



# *Two main avenues for handling short-read data...*



**Mapping or Reference Assembly**

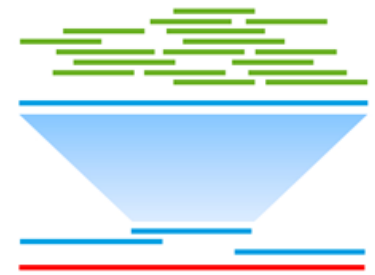


**De novo Assembly**

Short / Long reads

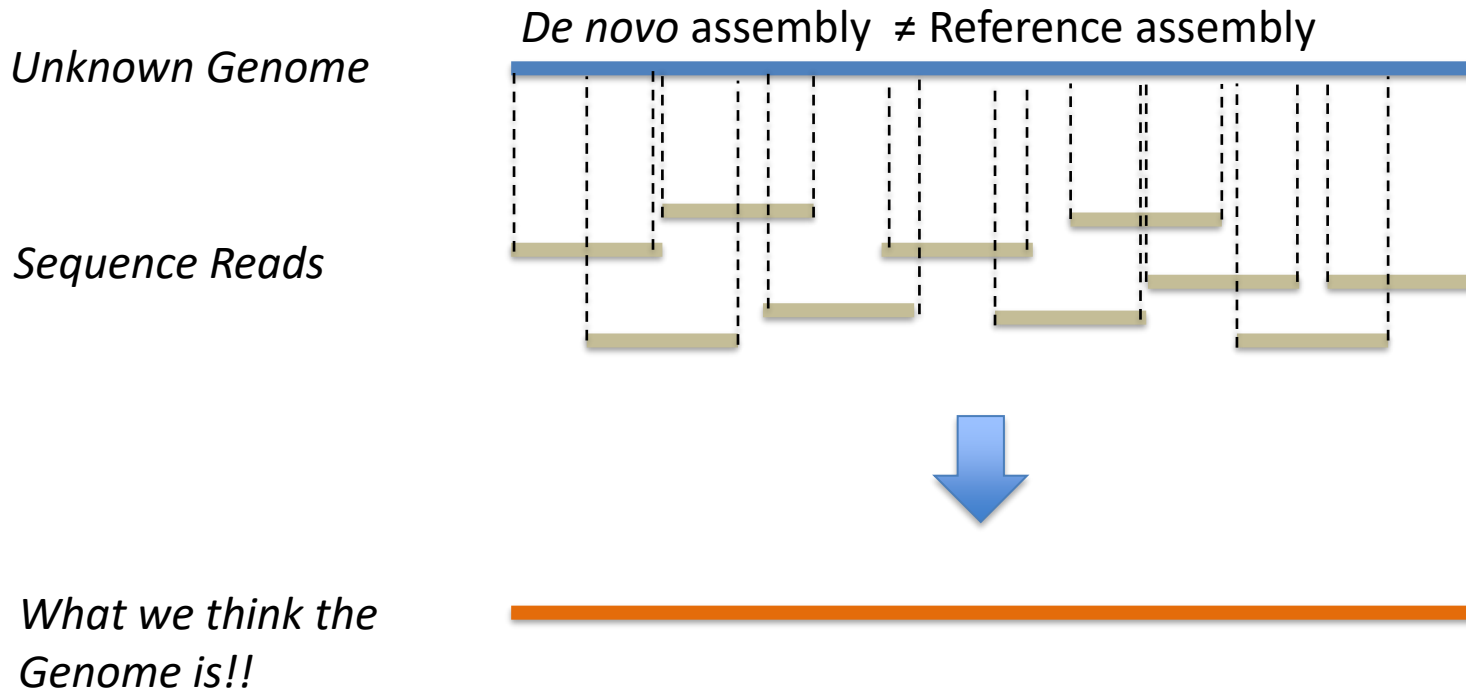
Contig assembly

Scaffold assembly  
Finished genome



# *De novo* Assembly

*De novo* assembly – genome assembly without a reference genome, i.e., starting *de novo* from sequence data



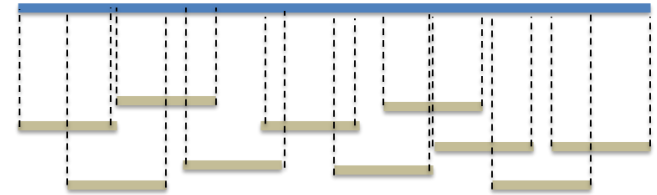
*A contig represents a stretch of DNA sequence assembled without any gaps, where the order and orientation of the nucleotides are known with high confidence.*

# *De novo* Assembly

*de novo* assembly attempts to reconstruct genomes by exploring read overlapping and contiguity

## Problems and Challenges:

- Large volume of sequence reads
- Sequencing errors
- Genomic repeat patterns/regions and homopolymers
- Uneven coverage/sequencing



**But, what is the definition of an assembly?**

*Best set of sequences that can approximate the sequenced genetic material*

***Implications?***



# Why to assemble?

## Objective/Purpose of the Assembly:

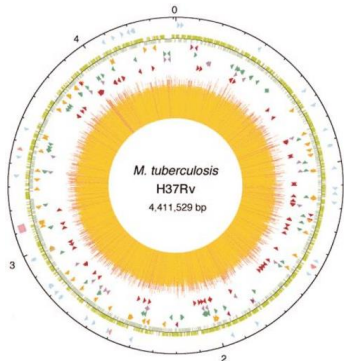
Obtain a reference genome



PacBio



Finished Genome Assembly



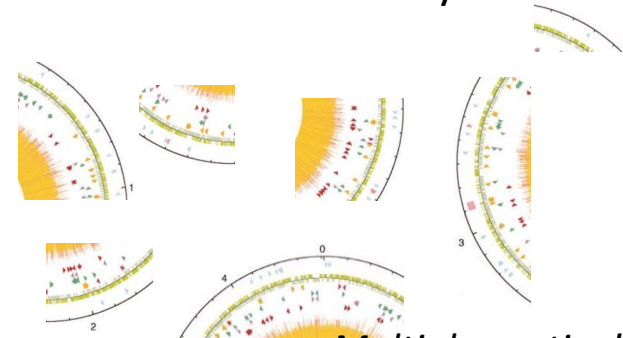
Gene content, insertions, deletions



Illumina



Draft Genome Assembly



Manual closure  
Mate-pair sequencing

*Multiple contigs!*

## Three major methods for assembly:

### i) overlap-and-extend

Finds read overlap and extends – suffix of a read is equal to the prefix of another read with a length that meets a defined threshold.

Software: SSAKE, VCAKE and SHARCGS

### ii) string graph

Construction of string graph where each read is a vertex with edges connecting overlapping nodes.

Software: Edena and BOA

Problems: high memory consumption and sequencing errors



## Three major methods for assembly:

### iii) de Bruijn Graph

Each vertex represents a substring of length  $k$  ( $k$ -mer) in a read. Edges connect vertexes if these are consecutive vertexes, i.e., the last  $k-1$  nucleotides in  $k$ -mer  $u$  are the same as the first  $k-1$  nucleotides of  $k$ -mer  $v$ .

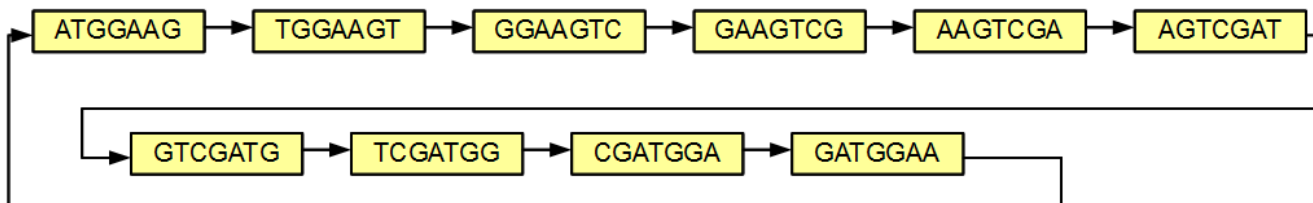
Software: Velvet, SOAPdenovo, SPAdes

Most widely used approaches.

Objective: represent every possible  $k$ -mer present in the genome!

ATGGAAGTCGATGGAAG

ATGGAAG  
TGAAGT  
GGAAGTC  
GAAGTCG  
AAGTCGA  
AGTCGAT  
GTCGATG  
TCGATGG  
CGATGGA  
GATGGAA  
ATGGAAG



## Three major methods for assembly:

### iii) de Bruijn Graph

Connecting the nodes:

- *Hamiltonian path*: visits **every vertex in the graph** (exactly once, because it is a path)
- *Eulerian trail*: visits **every edge in the graph** exactly once (because it is a trail, vertices may well be crossed more than once.)



William Hamilton



Leonhard Euler



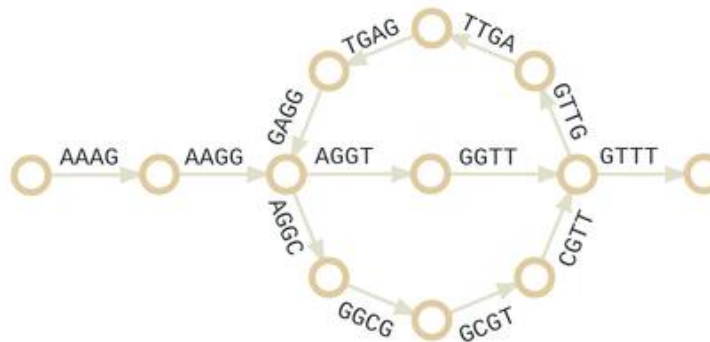
Nicolaas de Bruijn

#### A. Short read to $k$ -mers ( $k=4$ )

AAAGGCGTTGAGGTT

AAAG  
AAGG  
AGGC  
GGCG  
GCGT  
CGTT  
GTTG  
TTGA  
TGAG  
GAGG  
AGGT  
GGTT

#### B. Eulerian de Bruijn graph



#### C. Hamiltonian de Bruijn graph

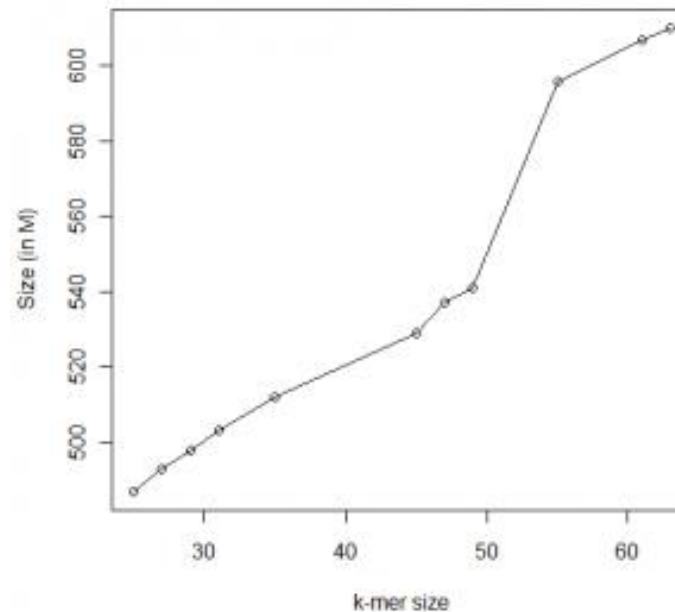
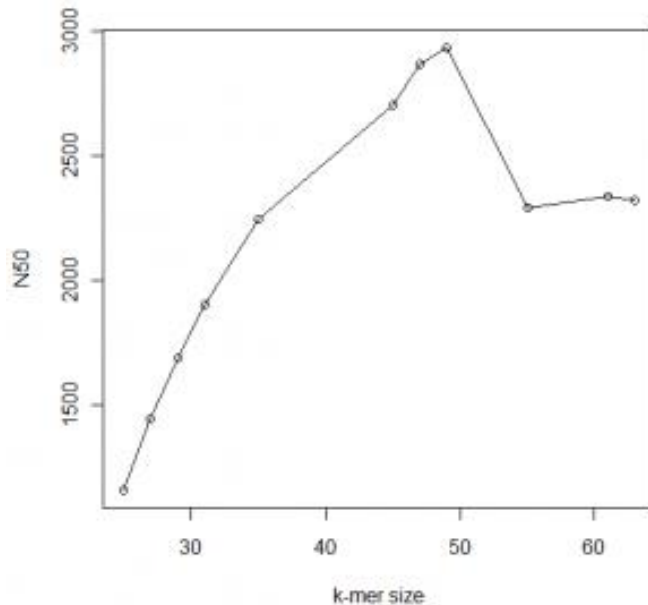


Three major methods for assembly:

## iii) de Bruijn Graph

Choice of the k-mer length: always below the read-length of your data!

The k-mer length should be an odd number: avoid palindromes – an odd-sized k-mer cannot form palindromes when reverse-complemented!



# De Bruijn Graphs - Exercise

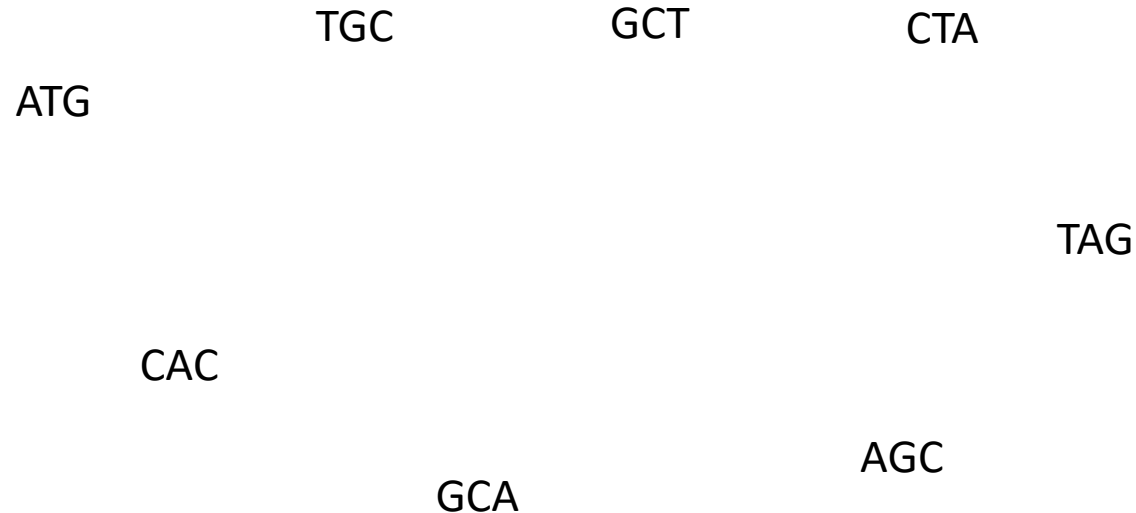
Let's consider the following reads:

ATGCTA

GCTAGC

TAGCAC

CTAGCA



1 List all 3bp k-mers

2 Establish links (edges) between k-mers differing by k-1 nucleotides

3 Visit all nodes and use the minimal path length

# De Bruijn Graphs - Exercise

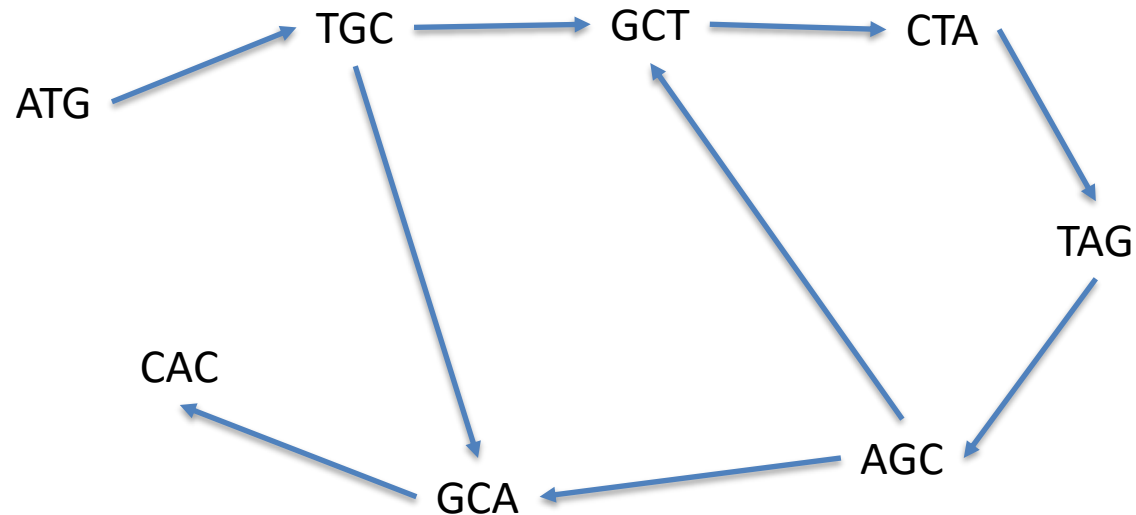
Let's consider the following reads:

ATGCTA

GCTAGC

TAGCAC

CTAGCA



1 List all 3bp k-mers

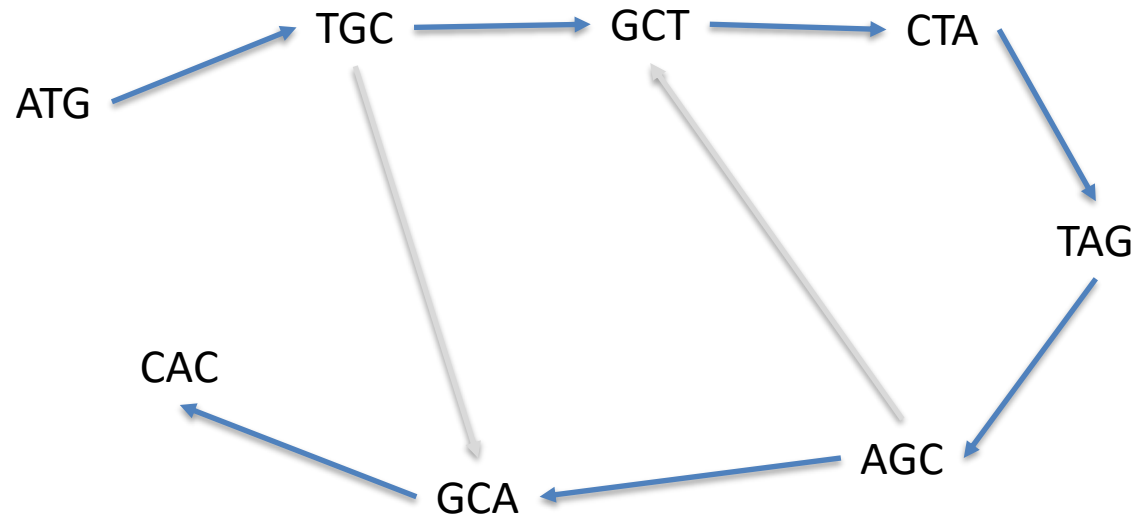
2 Establish links (edges) between k-mers differing by k-1 nucleotides

3 Visit all nodes only once

# De Bruijn Graphs - Exercise

Let's consider the following reads:

ATGCTA  
GCTAGC  
TAGCAC  
CTAGCA



ATGCTAGCAC

- 1 List all 3bp k-mers
- 2 Establish links (edges) between k-mers differing by k-1 nucleotides
- 3 Visit all nodes only once

# De Bruijn Graphs - Exercise

Let's consider the following reads:

ATGCTA  
GCTAGC  
TAGCAC  
CTAGCA

The Eulerian Path



- 1 List all 3bp k-mers
- 2 Define nodes between edges
- 3 Visit all edges only once

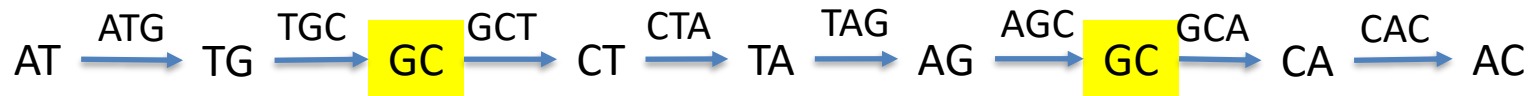


# De Bruijn Graphs - Exercise

Let's consider the following reads:

ATGCTA  
GCTAGC  
TAGCAC  
CTAGCA

The Eulerian Path



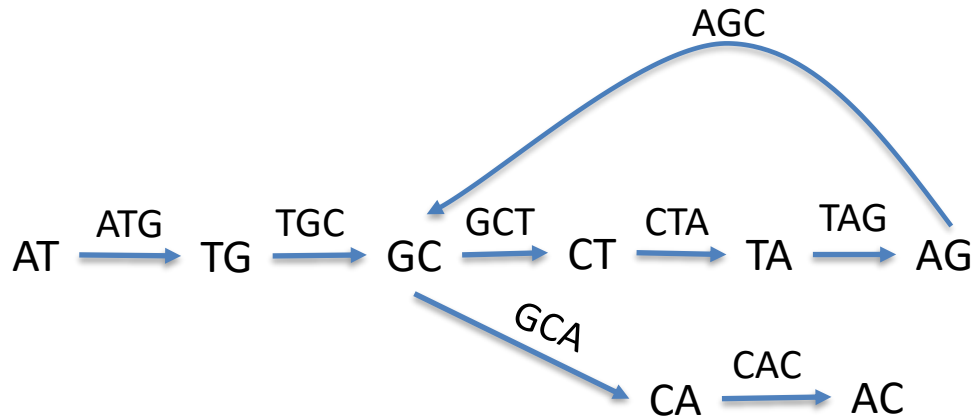
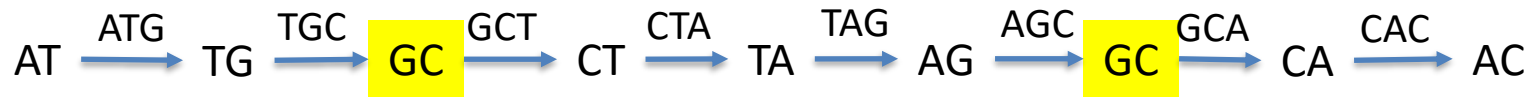
- 1 List all 3bp k-mers
- 2 Define nodes between edges
- 3 Visit all edges only once

# De Bruijn Graphs - Exercise

Let's consider the following reads:

ATGCTA  
GCTAGC  
TAGCAC  
CTAGCA

The Eulerian Path



ATGCTAGCAC

- 1 List all 3bp k-mers
- 2 Define nodes between edges
- 3 Visit all edges only once

# Evaluating and Comparing Assemblies

## Metrics to Evaluate and Compare Assemblies:

Evaluating an assembly can be reference-free or comparing to a reference genome!

### Purpose:

1. Assess the individual quality of an assembly;
2. Compare assemblers.

### Metrics commonly used:

- Number of contigs/scaffolds
- Total length of the assembly
- Length of the largest contig/scaffold
- Percentage of gaps in scaffolds ('N')
- N50/NG50 of contigs/scaffolds
- Number of predicted genes
- Number of core single-copy genes

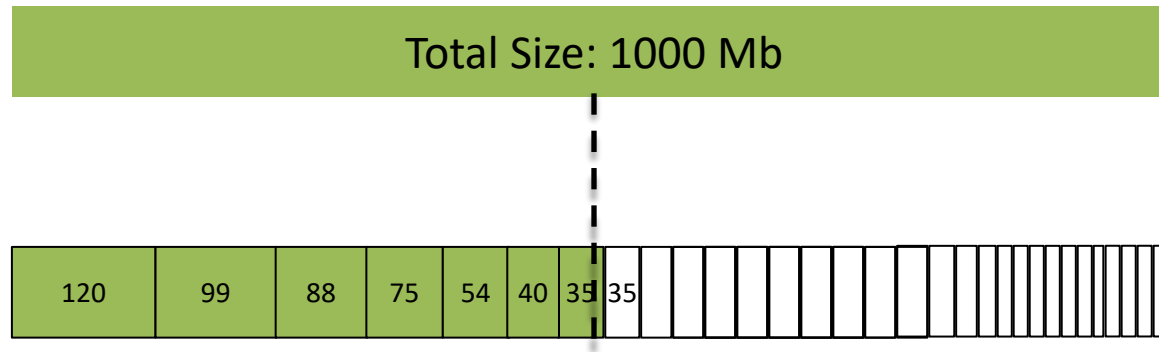
### Software:

**QUAST** – allows comparison against reference

**BUSCO** – evaluates the presence of core single-copy orthologous genes

**CheckM** - estimates of genome completeness and contamination by using collocated sets of genes that are ubiquitous and single-copy within a phylogenetic lineage

## N50 and L50 and other metrics



**N50** – shortest contig length spanning the midpoint of the assembly length (after sorting from largest to smallest contig);

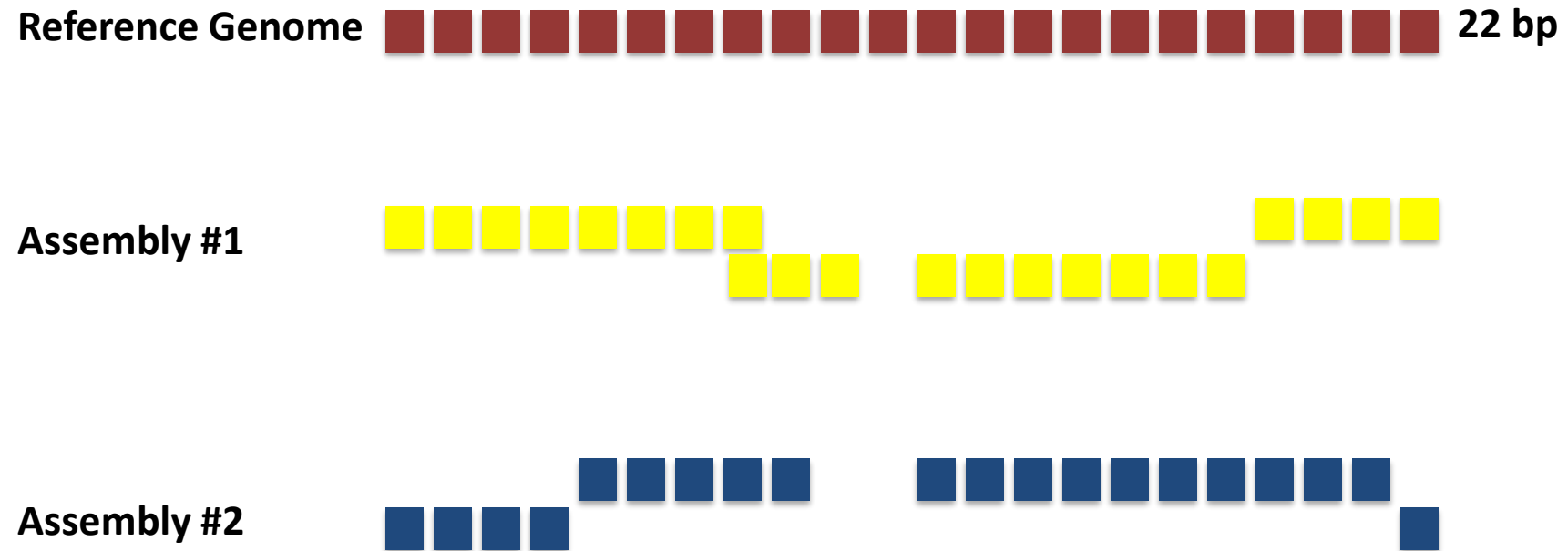
**NG50** – shortest contig length spanning the midpoint of the estimated genome size (after sorting from largest to smallest contig);

**L50** – number of contigs necessary to span the midpoint of the assembly length (after sorting from largest to smallest contig)

**LG50** – number of contigs necessary to span the midpoint of the estimated genome size (after sorting from largest to smallest contig)

*N50 and L50 in this exemple?*

# Evaluating and Comparing Assemblies

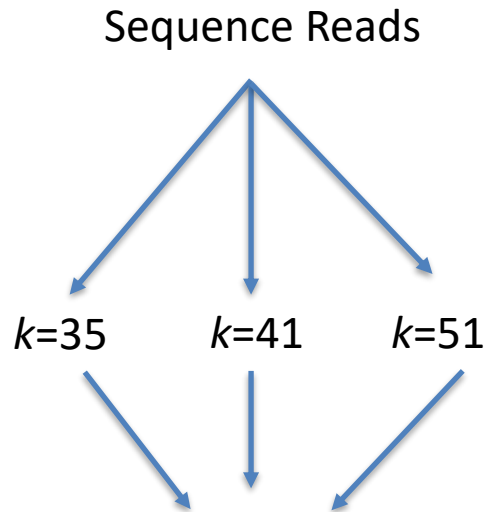


Calculate:	Assembly #1	Assembly #2
N50	7	10
NG50	7	5
Coverage	95.4%	90.9%
Total Length of Assembly	22	20

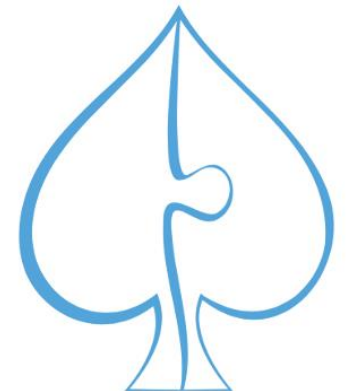
# Multi *k*-mer assembly – the way forward

## Multi k-mer assembly:

*Multi-k-mer strategies always provide better results than single k-mer*



Merge Assemblies -> Velvet+cd-hit+minimus2  
Iterative Assembly removing assembled reads – IDBA  
Multi-kmer de Bruijn graph - **SPAdes**



# Assembly pipelines



## As input, Unicycler takes one of the following:

- Illumina reads from a bacterial isolate (ideally paired-end, but unpaired works too)
- A set of long reads (either PacBio or Nanopore) from a bacterial isolate (uncorrected long reads are fine, though corrected long reads should work too)
- Illumina reads and long reads from the same isolate (best case)

## Reasons to use Unicycler:

- It circularises replicons without the need for a separate tool like [Circlator](#).
- It handles plasmid-rich genomes.
- It can use long reads of any depth and quality in hybrid assembly. 10x or more may be required to complete a genome, but Unicycler can make nearly-complete genomes with far fewer long reads.
- It produces an assembly *graph* in addition to a contigs FASTA file, viewable in [Bandage](#).
- It has very low misassembly rates.
- It can cope with very repetitive genomes, such as [Shigella](#).
- It's easy to use: runs with just one command and usually doesn't require tinkering with parameters.

## Reasons to not use Unicycler:

- You're assembling a eukaryotic genome or a metagenome (Unicycler is designed exclusively for bacterial isolates).
- Your Illumina reads and long reads are from different isolates (Unicycler struggles with sample heterogeneity).
- You're impatient (Unicycler is thorough but not especially fast).

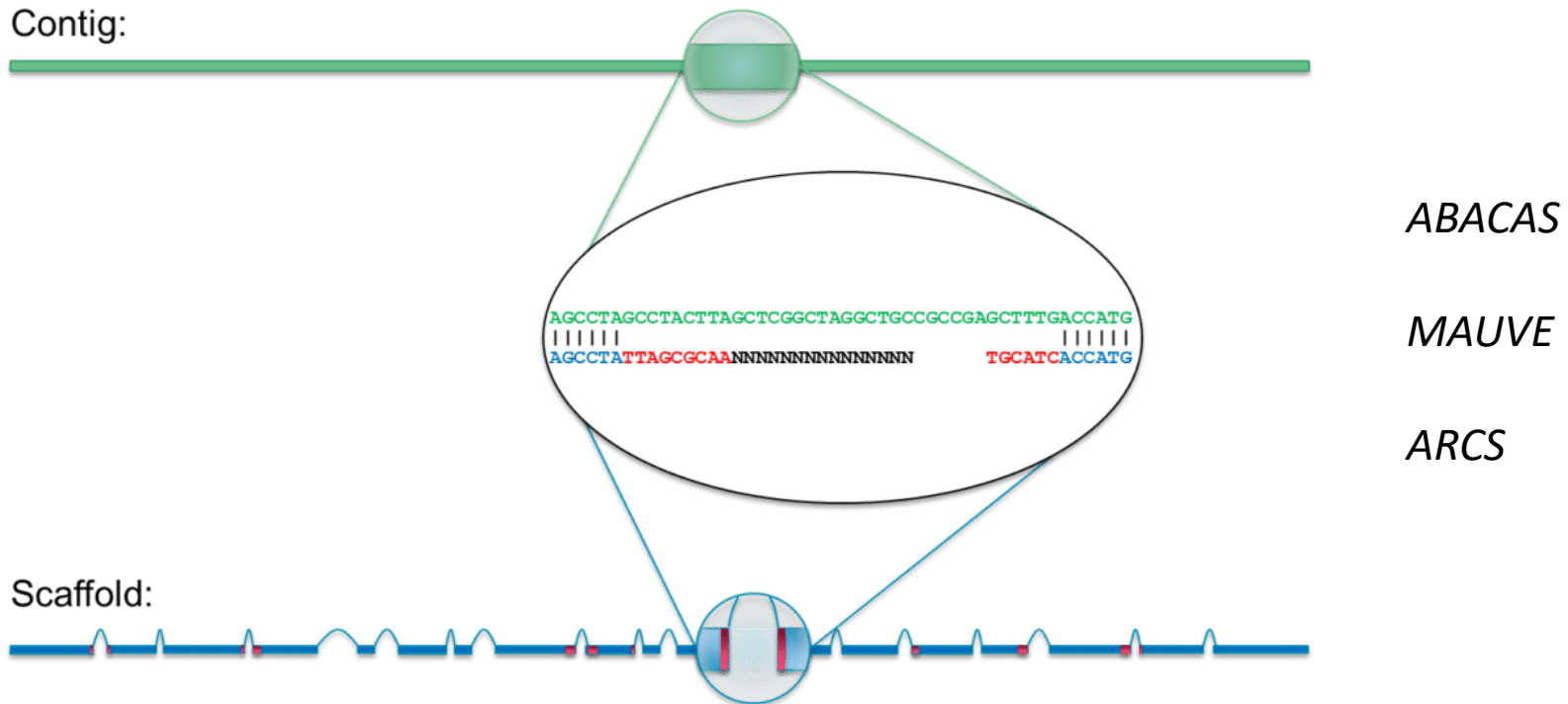
## Unicycler does:

- ***Short-read assembly***
- ***Long-read assembly***
- ***Hybrid assembly***



*Contigs are continuous stretches of sequence containing only A, C, G, or T bases without gaps.*

*Scaffolds are created by chaining contigs together **using additional information** about the relative position and orientation of the contigs in the genome*

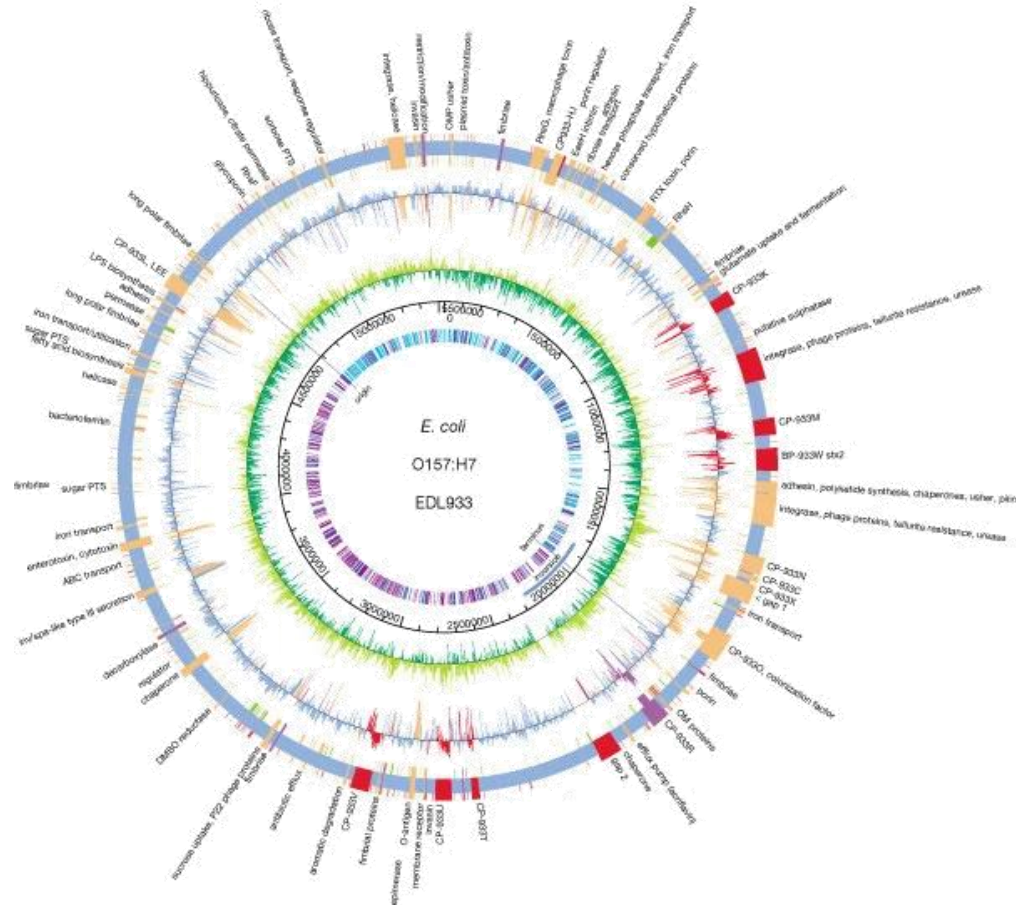


*Genome annotation refers to the identification of the location of genetic features...*

... most often protein-coding genes!

*This entails the exact pinpointing of feature coordinates and orientation in the assembled/finished genome.*

*Easier for bacterial genomes (high coding density ~90%, no introns) – less complex!*



*Two main approaches: intrinsic (ab initio) vs extrinsic (evidence based)*

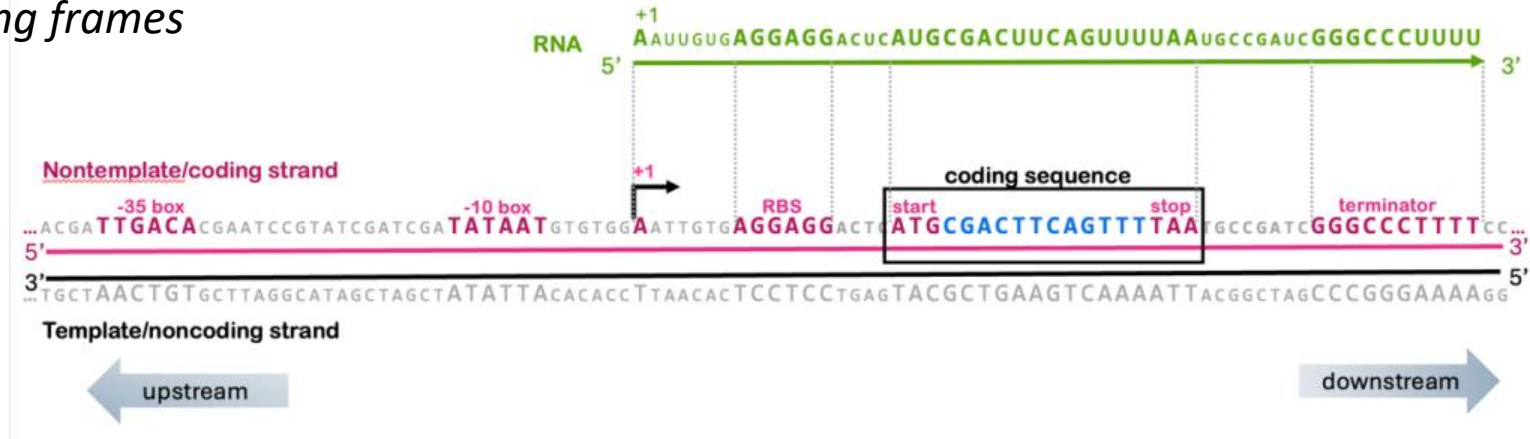
# Genome Annotation

## ***Intrinsic (ab initio) gene prediction***

- Solely based on DNA sequence without comparison to other sequences or databases
- Search for sequence signals associated with open reading frames

## **Software/Tools:**

GeneMark  
GLIMMER  
EasyGene  
GENSCAN

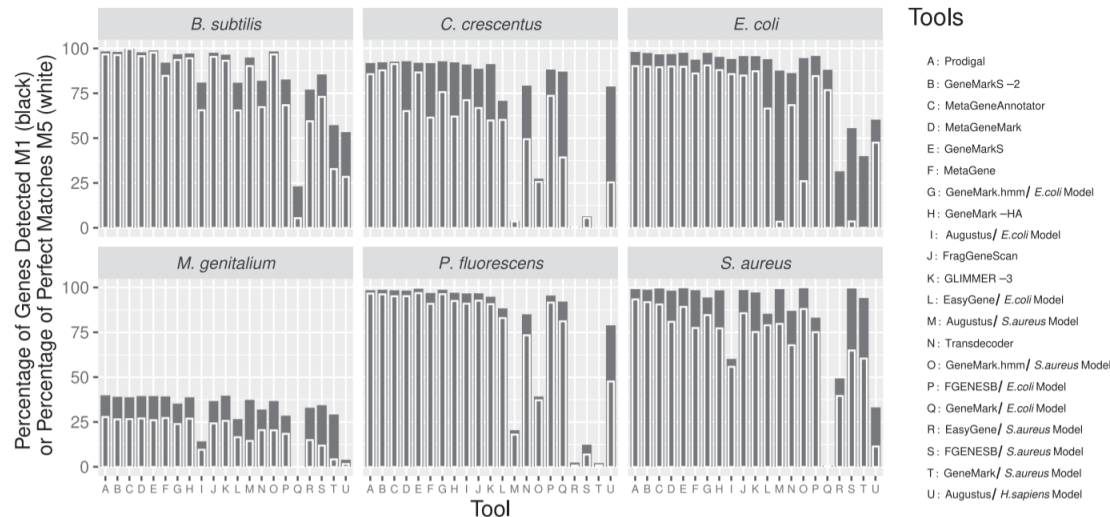
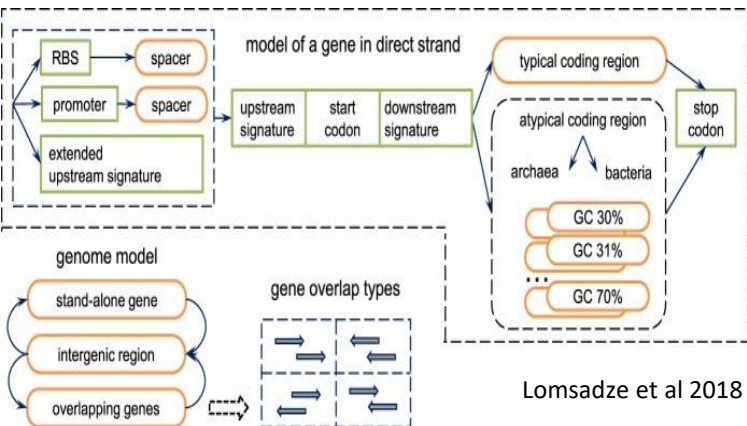


## ***Extrinsic – Evidence based gene prediction***

- Homology-based and comparative approach
- Comparison with databases of curated protein families or pre-existing genus/species-specific databases
- Can be supported via BLASTX!

# Genome Annotation

*There are multiple tools for gene prediction, often implementing complex statistical models (Hidden Markov Models) and Machine-Learning algorithms – its performance varies across organism*



Dimonaco et al 2022 Bioinformatics

## Some of the most widely known tools:

- Prodigal**: Very fast (*E. coli* K12 ~10s), unsupervised machine learning – no training, handles draft assemblies and metagenomes.
- GeneMarkS-2**: More advanced version of GeneMark, specifically for prokaryotes, this new version has the capability for self-training, used at NCBI, incorporates the detection of atypical coding regions;
- Glimmer3**: Interpolated Markov Models, useful for Bacteria, Archaea and viruses, developed and used at TIGR.

## RAST

*Web server for prokaryotic genome annotation – user friendly!*

*Initially calls tRNA and rRNA genes and subsequently calls protein-coding genes via implementation of GLIMMER2;*

*Able to establish phylogenetic context – via comparison of universal FIGfams and subsequent targeted search across closely related genomes.*

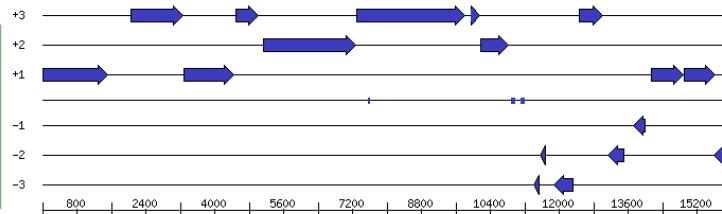
**The SEED Viewer** SEED Viewer version 2.0  
Welcome to the SEED Viewer - a read-only browser of the curated SEED data.  
For more information about The SEED please visit [theSEED.org](http://theSEED.org).

»Navigate »Organism »Comparative Tools »Help

**Browse Genome:** [Mycobacterium tuberculosis HPV115\\_08 \(1773.11223\)](#)

Location Focus Upload List

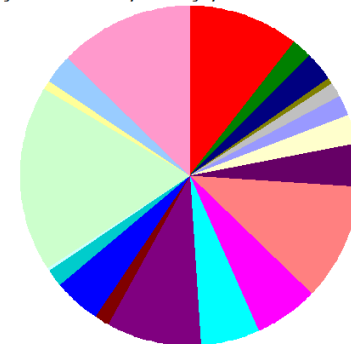
contig HPV115\_08 (4,361,434 bp) v  
start base 0  
window 16,000 bp v  
Color features by focus v  
<== draw ==>



Subsystem Coverage



Subsystem Category Distribution



Subsystem Feature Counts

- Cofactors, Vitamins, Prosthetic Groups, Pigments (176)
- Cell Wall and Capsule (36)
- Virulence, Disease and Defense (45)
- Potassium metabolism (7)
- Photosynthesis (0)
- Miscellaneous (23)
- Phages, Prophages, Transposable elements, Plasmids (2)
- Membrane Transport (35)
- Iron acquisition and metabolism (2)
- RNA Metabolism (44)
- Nucleosides and Nucleotides (71)
- Protein Metabolism (182)
- Cell Division and Cell Cycle (0)
- Mobility and Chemotaxis (0)
- Regulation and Cell signaling (104)
- Secondary Metabolism (1)
- DNA Metabolism (93)
- Fatty Acids, Lipids, and Isoprenoids (151)
- Nitrogen Metabolism (25)
- Dormancy and Sporulation (2)
- Respiration (74)
- Stress Response (26)
- Metabolism of Aromatic Compounds (6)
- Amino Acids and Derivatives (294)
- Sulfur Metabolism (12)
- Phosphorus Metabolism (47)
- Carbohydrates (204)

## Prokka

*“Prokka coordinates a suite of existing software tools to achieve a rich and reliable annotation of genomic bacterial sequences. Where possible, it will exploit multiple processing cores, and a typical bacterial genome can be annotated in ~10 min on a quad core desktop computer. It is well suited to iterative models of sequence analysis and integration into genomic software pipelines.”*

- (1) An optional user-provided set of annotated proteins, searched using BLAST+ blastp (Camacho et al., 2009).
- (2) All bacterial proteins in UniProt (Apweiler et al., 2004) that have real protein or transcript evidence and are not a fragment. This is 16 000 proteins, and typically covers 450% of the core genes in most genomes. BLAST+ is used for the search.
- (3) All proteins from finished bacterial genomes in RefSeq for a specified genus (optional). This captures domain-specific naming, and the databases vary in size and quality, depending on the popularity of the genus.
- (4) A series of hidden Markov model profile databases, including Pfam (Punta et al., 2012) and TIGRFAMs (Haft et al., 2013). This is performed using hmmscan from the HMMER 3.1 package (Eddy, 2011).
- (5) If no matches can be found, label as ‘hypothetical protein’.

## Prokka: rapid prokaryotic genome annotation

Torsten Seemann<sup>1,2</sup>

<sup>1</sup>Victorian Bioinformatics Consortium, Monash University, Clayton 3800 and <sup>2</sup>Life Sciences Computation Centre, Victorian Life Sciences Computation Initiative, Carlton 3053, Australia

Associate Editor: Alfonso Valencia

**Table 1.** Feature prediction tools used by Prokka

Tool (reference)	Features predicted
Prodigal (Hyatt 2010)	Coding sequence (CDS)
RNAmmer (Lagesen <i>et al.</i> , 2007)	Ribosomal RNA genes (rRNA)
Aragorn (Laslett and Canback, 2004)	Transfer RNA genes
SignalP (Petersen <i>et al.</i> , 2011)	Signal leader peptides
Infernal (Kolbe and Eddy, 2011)	Non-coding RNA



# Genome Annotation - Prokka

## Prokka

### Performance and output files

**Table 3.** Comparison of annotation of *E.coli* K-12 accession U00096.2

Feature	Reference	Prokka	RAST	xBase2
Total CDS	4321	<b>4305</b>	4512	4444
<i>Matching start</i>	—	<b>3828</b>	3571	3025
<i>Different start</i>	—	<b>318</b>	533	1052
<i>Missing CDS</i>	—	<b>172</b>	214	241
<i>Extra CDS</i>	—	<b>159</b>	405	367
<i>Hypothetical protein</i>	18	276	638	<b>156</b>
<i>With EC number</i>	1114	1050	<b>1118</b>	0
Total tRNA	89	<b>88</b>	86	<b>88</b>
Total rRNA	22	<b>22</b>	<b>22</b>	<b>22</b>

*Prokka produces a full set of output files that enable submission to the NCBI*

**Table 2.** Description of Prokka output files

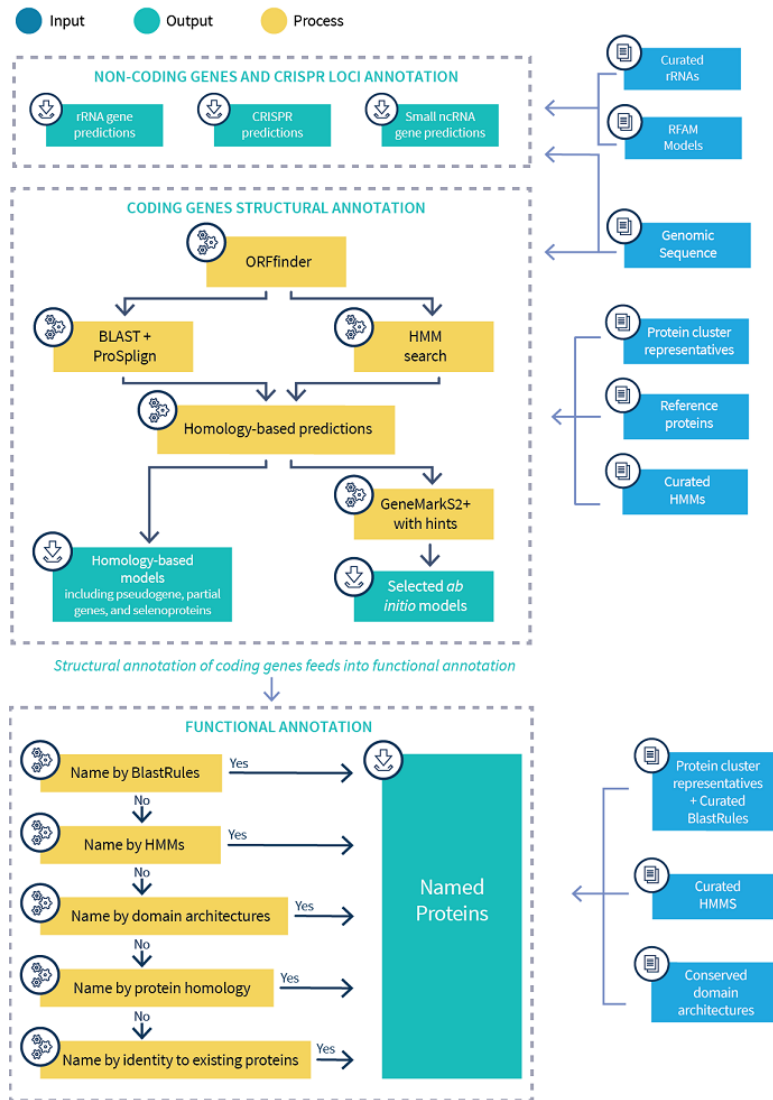
Suffix	Description of file contents
.fna	FASTA file of original input contigs (nucleotide)
.faa	FASTA file of translated coding genes (protein)
.ffn	FASTA file of all genomic features (nucleotide)
.fsa	Contig sequences for submission (nucleotide)
.tbl	Feature table for submission
.sqn	Sequin editable file for submission
.gbk	Genbank file containing sequences and annotations
.gff	GFF v3 file containing sequences and annotations
.log	Log file of Prokka processing output
.txt	Annotation summary statistics

*Performance for an E. coli genome (well studied) shows an overall better performance in comparison to RAST or xBase2.*



# Genome Annotation - PGAP

## NCBI – Prokaryotic Genome Annotation Pipeline



Combines *ab initio* gene prediction algorithms with homology based algorithms

Multi-level approach that enables detection of:

- Protein-coding genes
- structural RNAs
- tRNAs
- small RNAs
- Pseudogenes
- control regions
- direct and inverted repeats
- insertion sequences
- transposons and other mobile elements

Now available as a command-line utility (Docker Container) locally deployable!