

## NGS and Genomes assemblies

*Bioinformatics teachings*

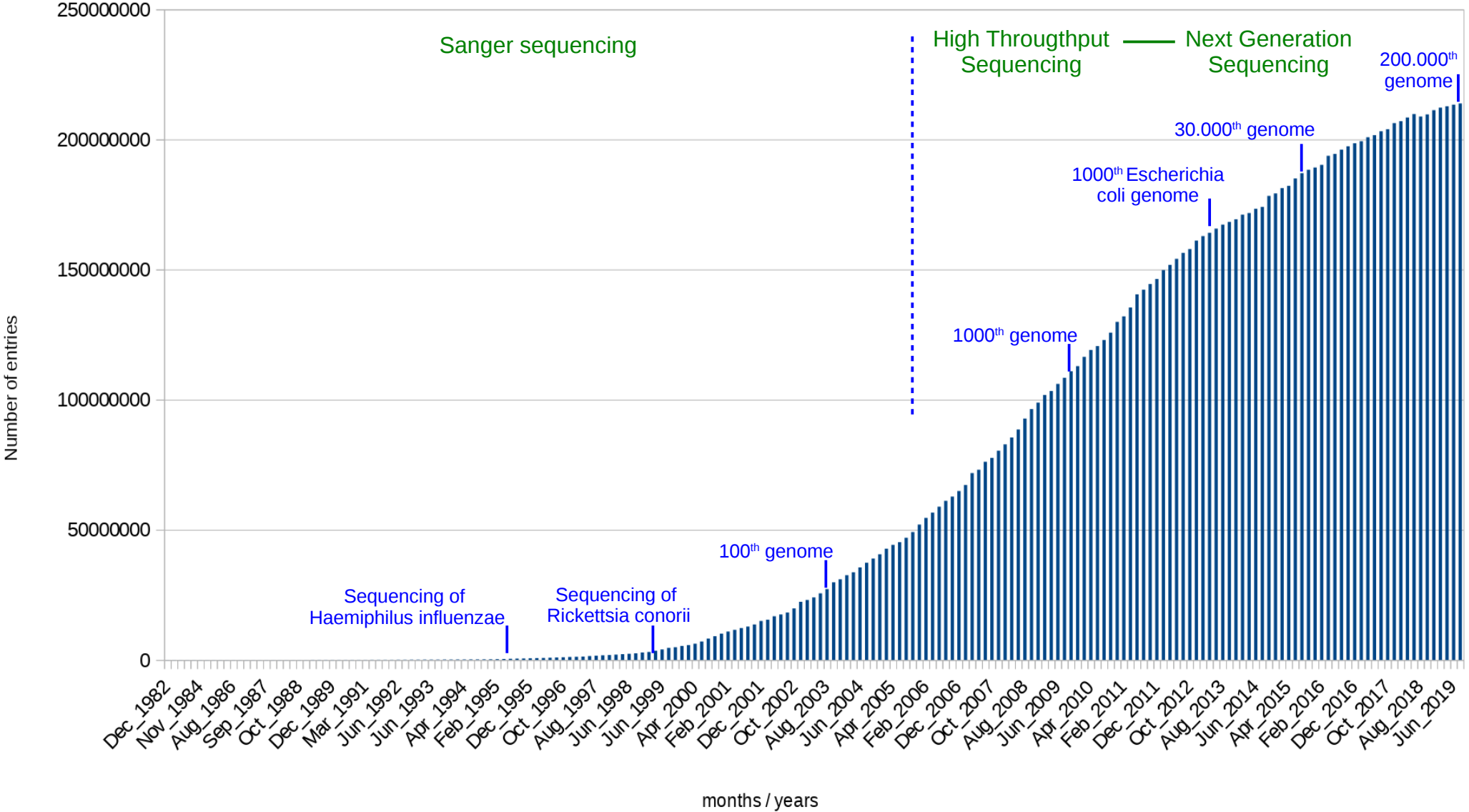
*<http://bioinfomed.fr> - Olivier Croce ([croce@unice.fr](mailto:croce@unice.fr))*

# Data release

## GenBank size

(GenBank.txt/gbrel.txt)

(in blue) : history of bacterial genomes sequencing



# Data release

- Submission of the sequence on public databases
- Not always => publication

## **3 main public databases:**

- EMBL-EBI - ENA (European Nucleotide Archive)  
<http://www.ebi.ac.uk/embl/>
- GenBank (USA) – NCBI  
<http://www.ncbi.nlm.nih.gov/Genbank/>
- DDBJ (DNA DataBank of Japon) – CIB  
<http://www.ddbj.nig.ac.jp/>



They are associated (International Nucleotide Sequence Database Collaboration) and exchange the same data which is periodically duplicated together

## **Contain:**

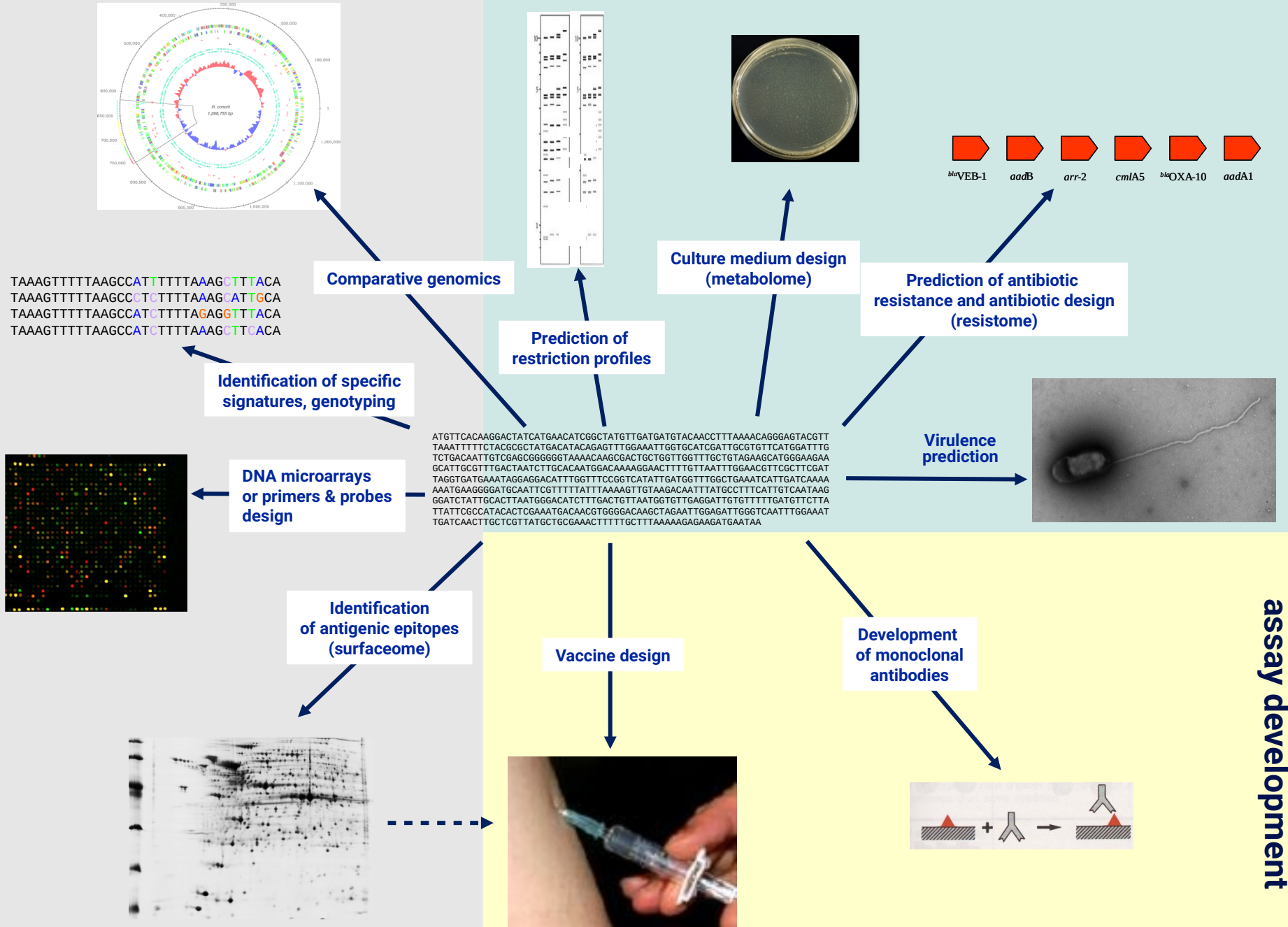
- Sequences of DNA or RNA from various sequencers technologies and from many labs
  - \* Some genome fragments : one or more genes, intergenic sequences, parts of a genome
  - \* Completed genomes
  - \* mRNA, tRNA, rRNA (ie. 16s)
- Annotations

# Aims

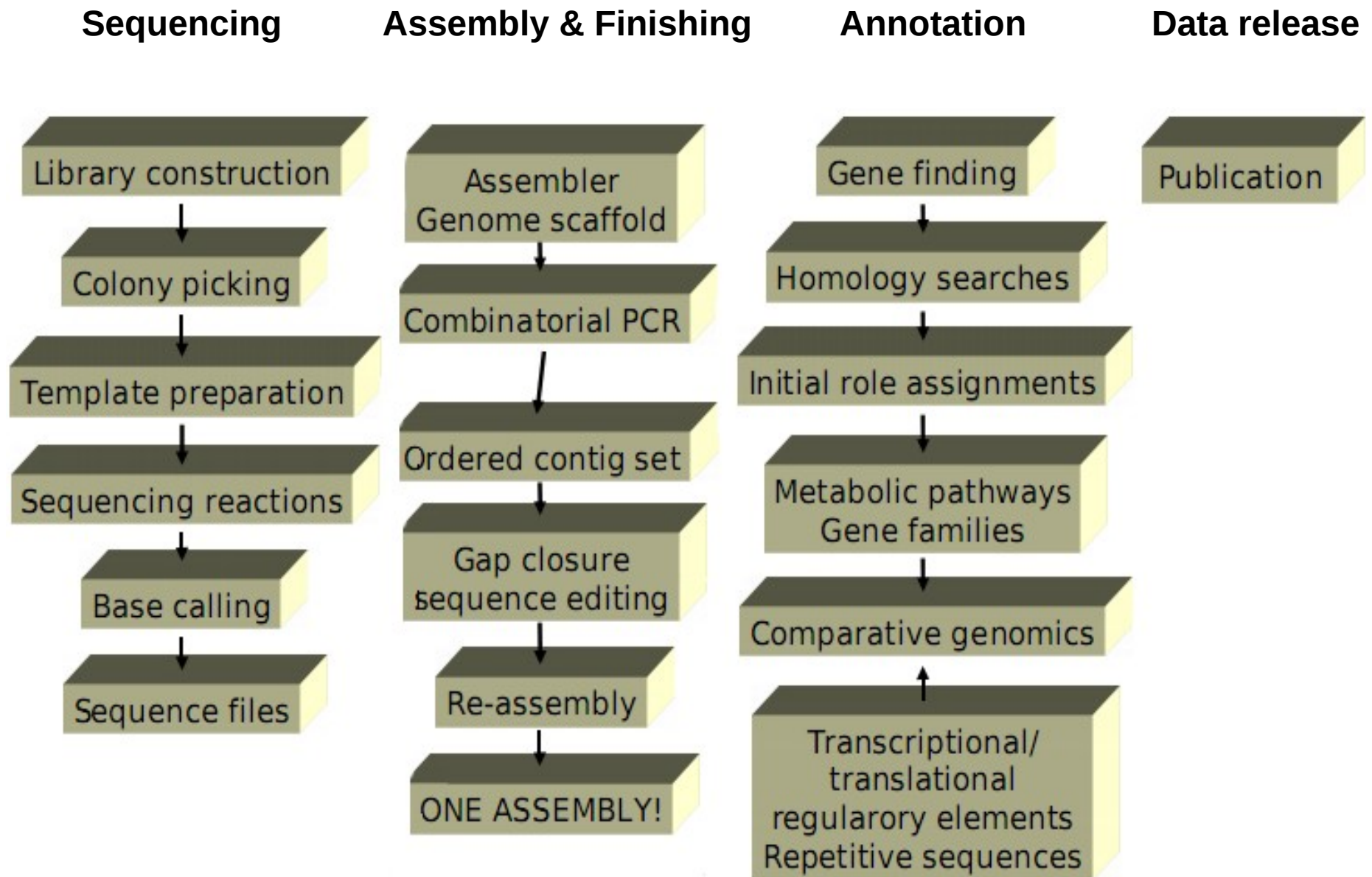
## Molecular detection and identification

## Phenotype prediction

## Vaccine & serological assay development



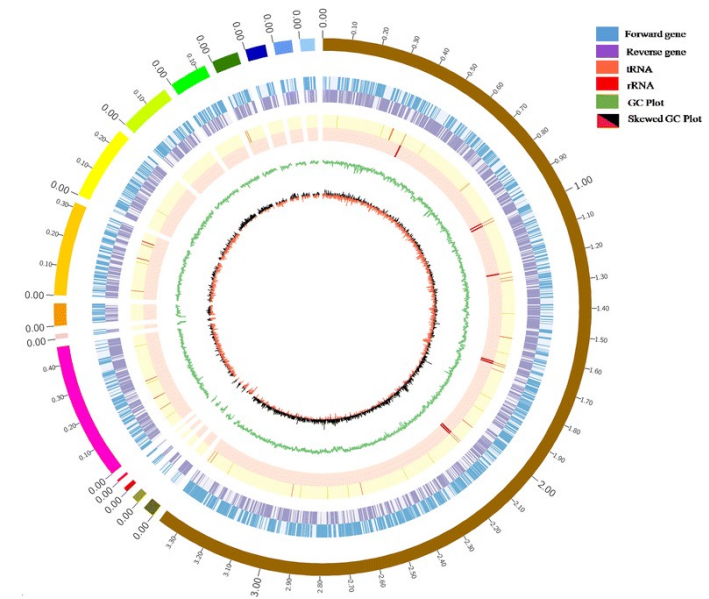
# From the bench to the publication



# Quality level of genomes

- **The genome sequence must be completed with a high quality and annotated before the release**

Of course the best, but very time consuming.  
Actually, 90-95 % of a microorganism genome could be easily covered without finishing, but the 5-10 % remained can take many weeks or months to be ended  
=> now easier using long reads sequencing



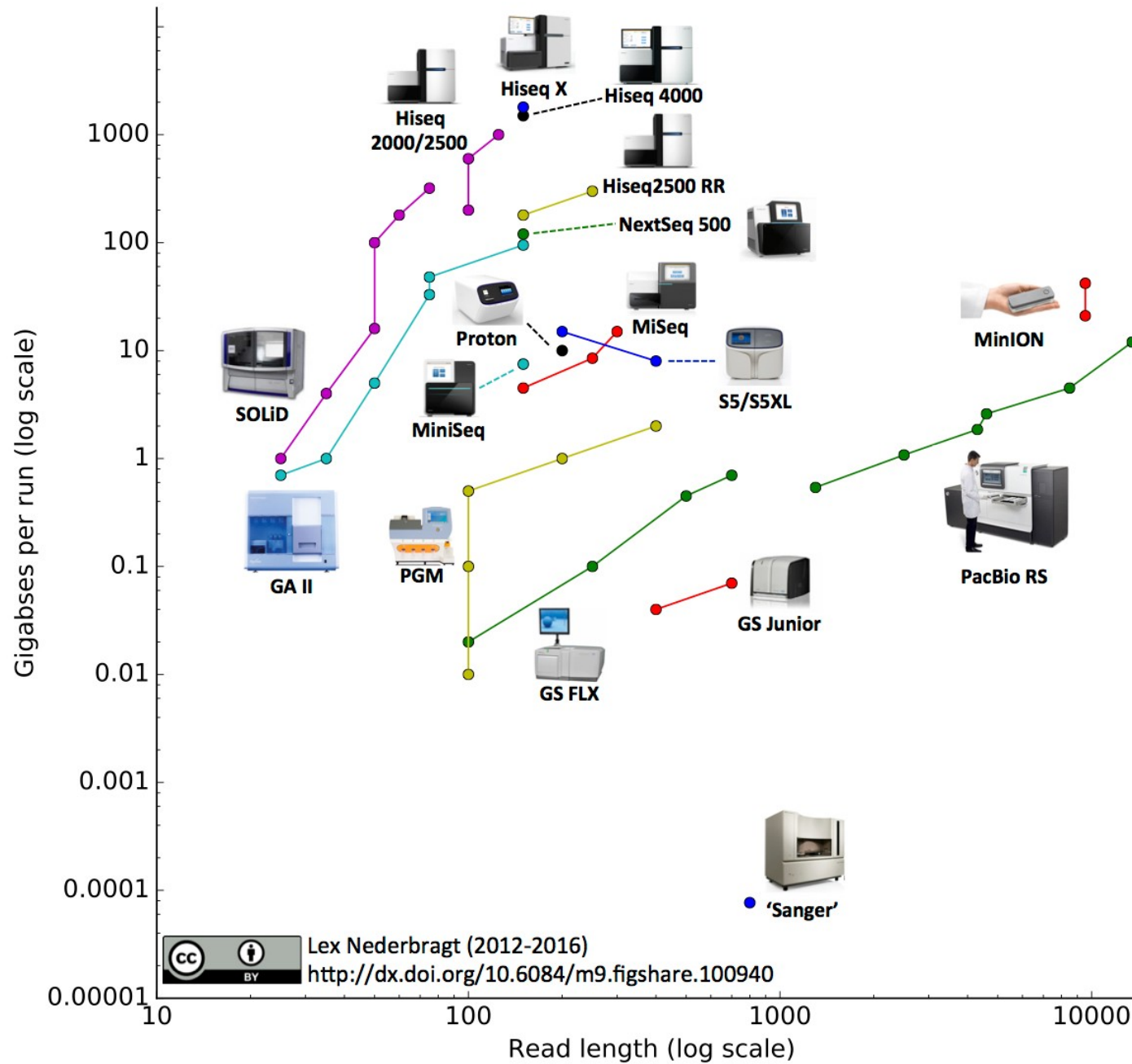
- **The sequence should be uncompleted with a draft quality, whether we suppose most of the genes are sequenced (and identified)**

Many eukaryote genomes are only draft genomes, because of the complexity of finishing

=> In general, fundamental research usually performs high quality genomes and applicative research (industry, part of clinical) usually performs draft genomes

=> depending of the project : time and experience (bioinformatician), money (coverage of NGS), organisms, the question to answer

# Sequencing technologies



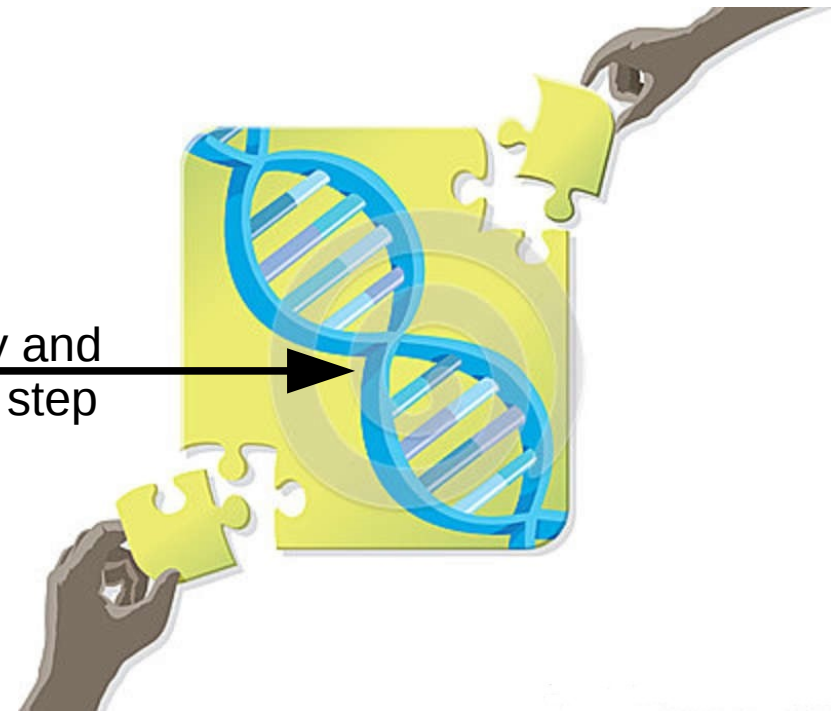
# Principle of sequencing and assembly



Sequencing step: reads have heterogeneous distribution



Assembly and finishing step





# Principle of sequencing and assembly

ATCGATGCGTAGCAGACTACCGTTACGATGCCTT...  
TAGCTACGCATCGTCTGATGGCAATGCTACGGAA...

Fragmentation + sequencing  
=> sets of reads

TAGCTACGCATCGT  
ATCGATGCGTAGC  
TAGCAGACTACCGTT  
GTTACGATGCCTT

ATCGATGCGTAGC  
TAGCAGACTACCGTT  
GTTACGATGCCTT  
TGCTACGCATCG → CGATGCGTAGCA  
(sequence inv-compl)

CGATGCGTAGCA  
ATCGATGCGTAGC  
TAGCAGACTACCGTT  
GTTACGATGCCTT

Build of contigs with overlapping regions

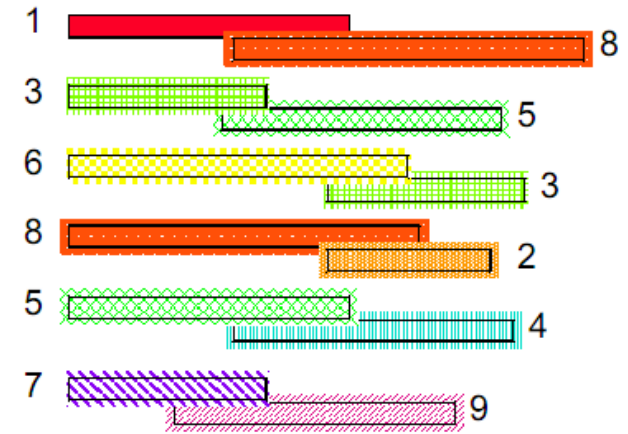
Assembly :  
=> alignements of reads + consensus

.....ATCGATGCGTAGCAGACTACCGTTACGATGCCTT.....

# Principle of sequencing and assembly

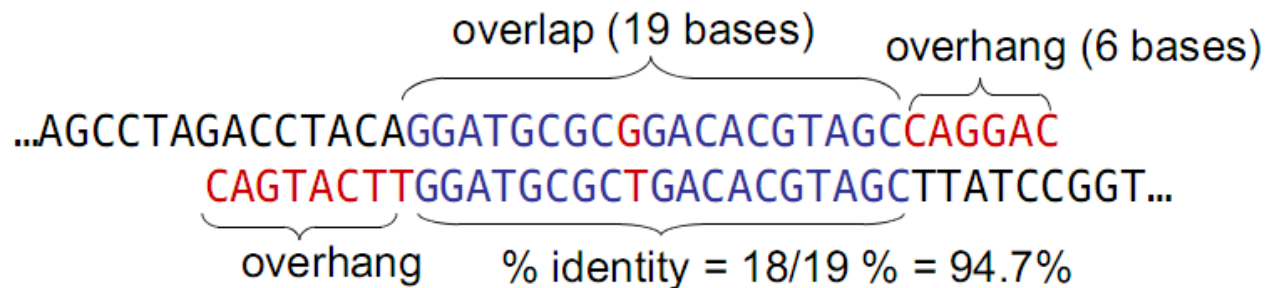
Search for best pairings by comparing each sequence (and its reverse complement) against every others sequences to find the best overlapping

=> list of best candidates with similarities criteria



Best candidate is a compromise between :

- maximum overlap length - region of similarity between regions
- minimum overhang length - unaligned ends of the sequences
- maximum % identity in overlap region
- minimum repeat length



=> Many **assemblers** tools existed (depending of sequencings technologies, libraries, genomes size, etc..)

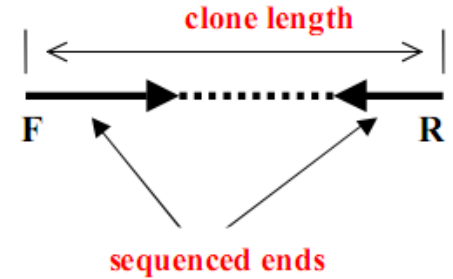
# Constructions of library from genomes fragments

\* **single - end (= shotgun)** : reads sequenced independently

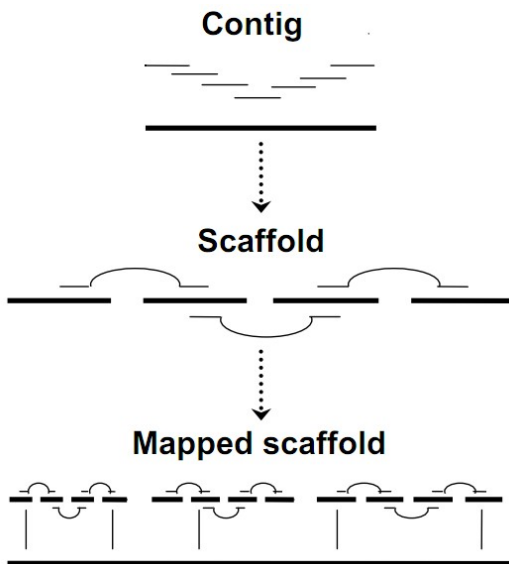
\* **paired-end** (similar to **mate-pair**) : reads are sequenced by pairs

- The distance between the reads is known (length of the insert), with some experimental uncertainty

- Distance of insert depends of technology (ie. Illumina ~150 nt for paired-end, ~1-5 kb mate-paired)



**Why using PE/MP ?** length of reads is limited => assemble repetitive regions by using reads as “anchors”

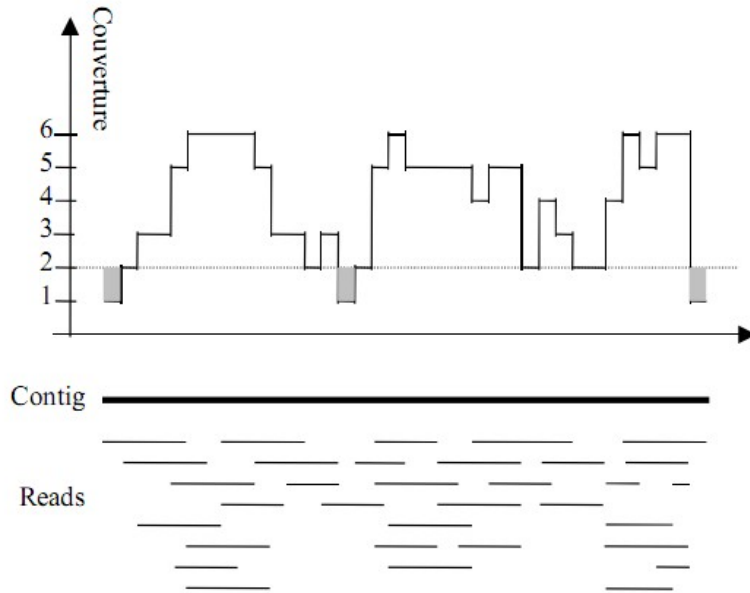


**Contigs** : group of overlapping reads, without gap

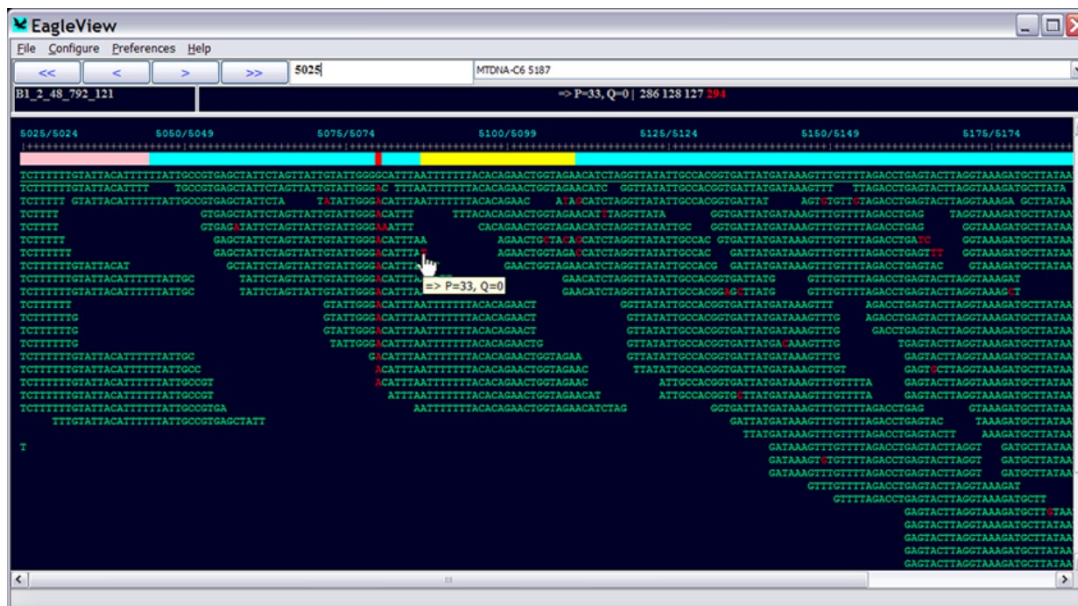
**Scaffold** : group of contigs order and in the same sens. Gap ("NNN") could existed and their length are known. Scaffolds exists only if a paired-ends (or mate pairs) sequencing was performed !

**Mapped scaffolds** : scaffolds mapped along a reference. Order, orientation and length of gaps are estimated, but not sure !

# Main remaining problems



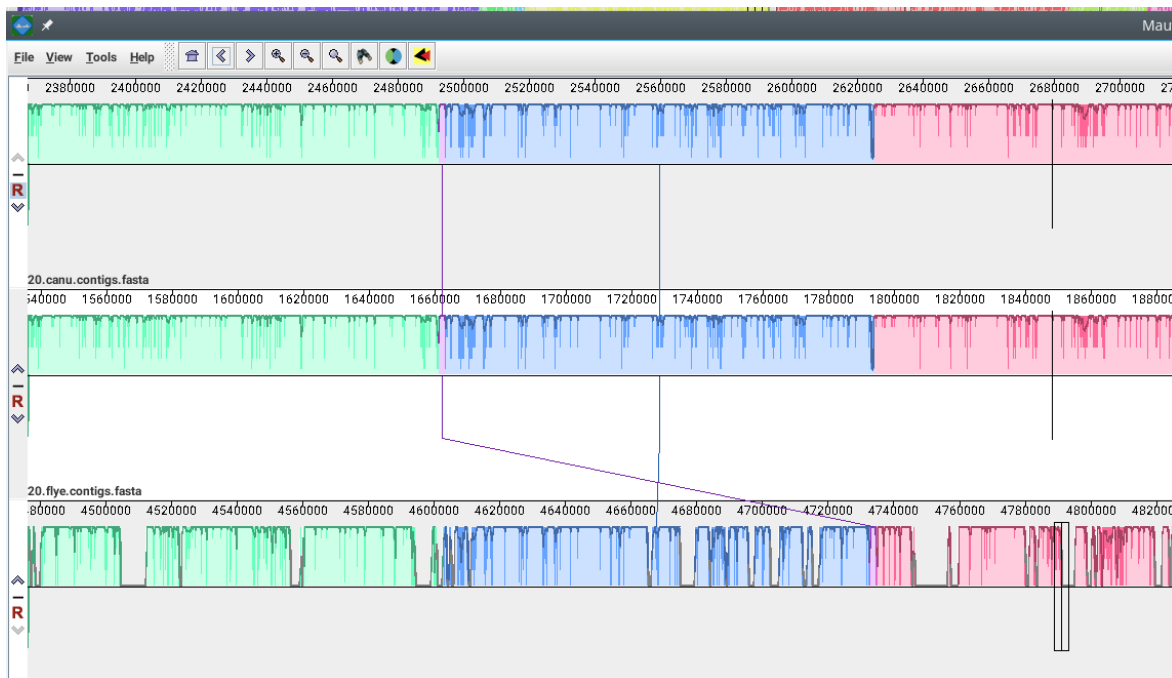
- Bad assembly of reads
- Low coverage of reads
- Bad insert size estimation
- Different orientation of contigs
- Error of sequencing
- Repeat sequence ambiguities



Re-Mapping, to see coverage, SNP and potential errors

# Finishing

- (re)Mapping of reads along the assembled genome (or/and a reference)
- help to correct the low quality/coverage areas
- Check the order of contigs
- Check the redundancy of contigs (false contigs or true repeat contigs like rRNA operons)
- Compare synteny between multiple assemblers (global alignment)
- Fill the gaps by extending the boundaries of each gap using ends of mapping reads (or use PCR)
- Order (or reorder) contigs
- Desassemble some areas if they seem to be false



*Bacillus cereus* assemblies using 3 assemblers tools. 2 first genomes are very similar, the third show many differences

=> High improvement with new long-reads technology (Minlon Nanopore, PacBio)