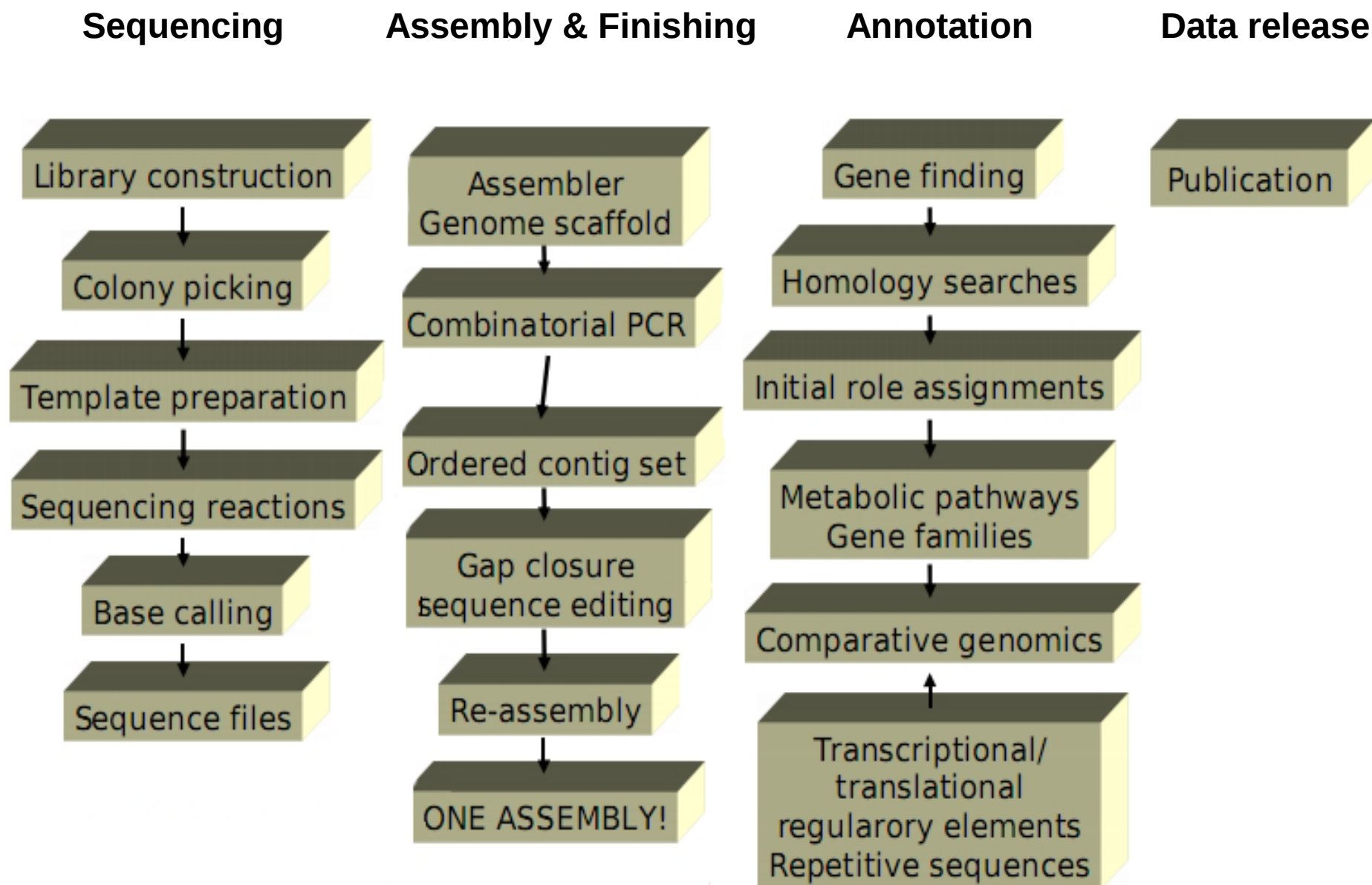


## Genomes assemblies and finishing

*Bioinformatics teachings*

*<http://bioinfomed.fr> - Olivier Croce -*

# Summary



# 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/>

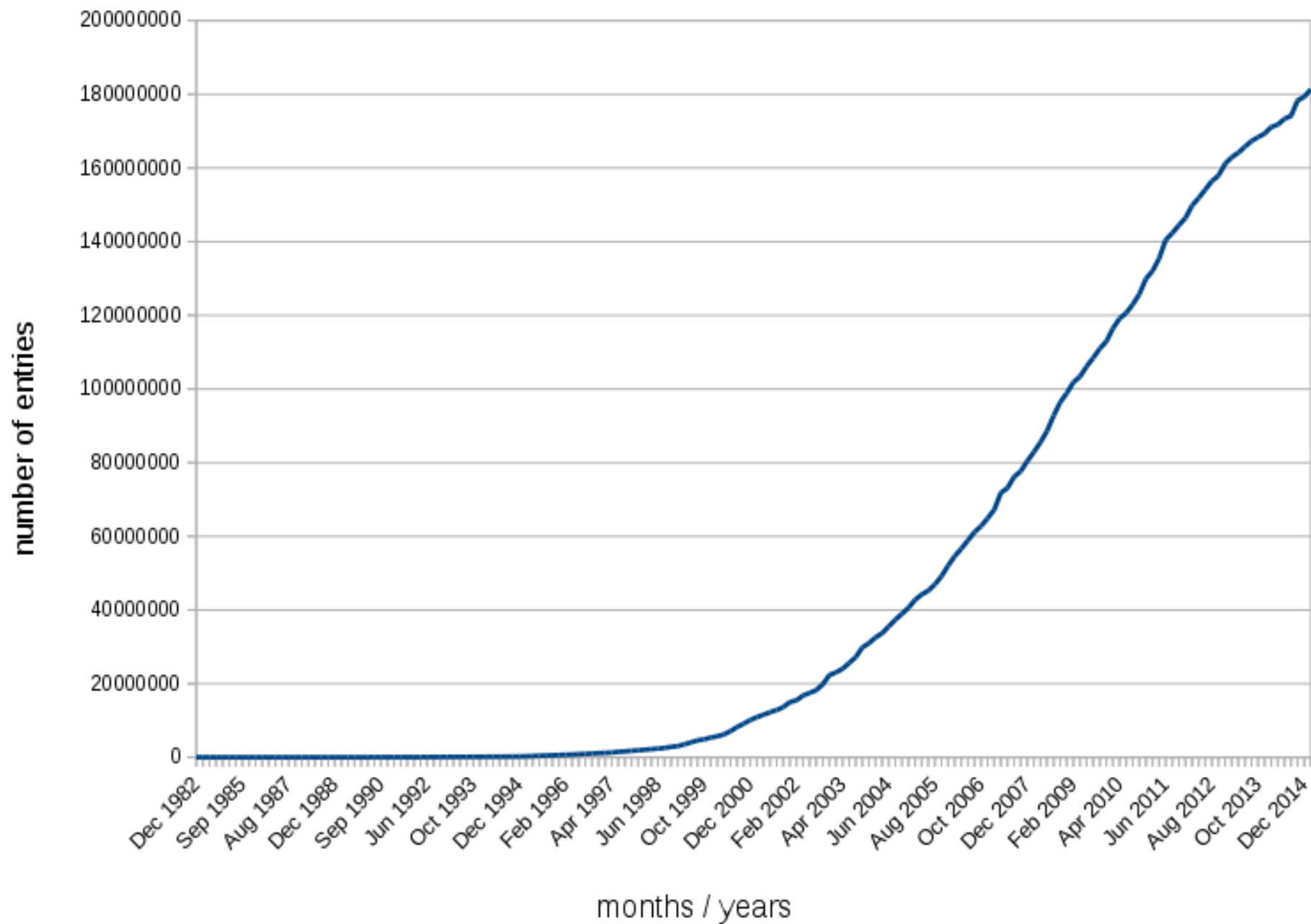
They are associated (International Nucleotide Sequence Database Collaboration) and exchange the same data which is periodically duplicated together

## Embl, Ebi, Ddbj contain:

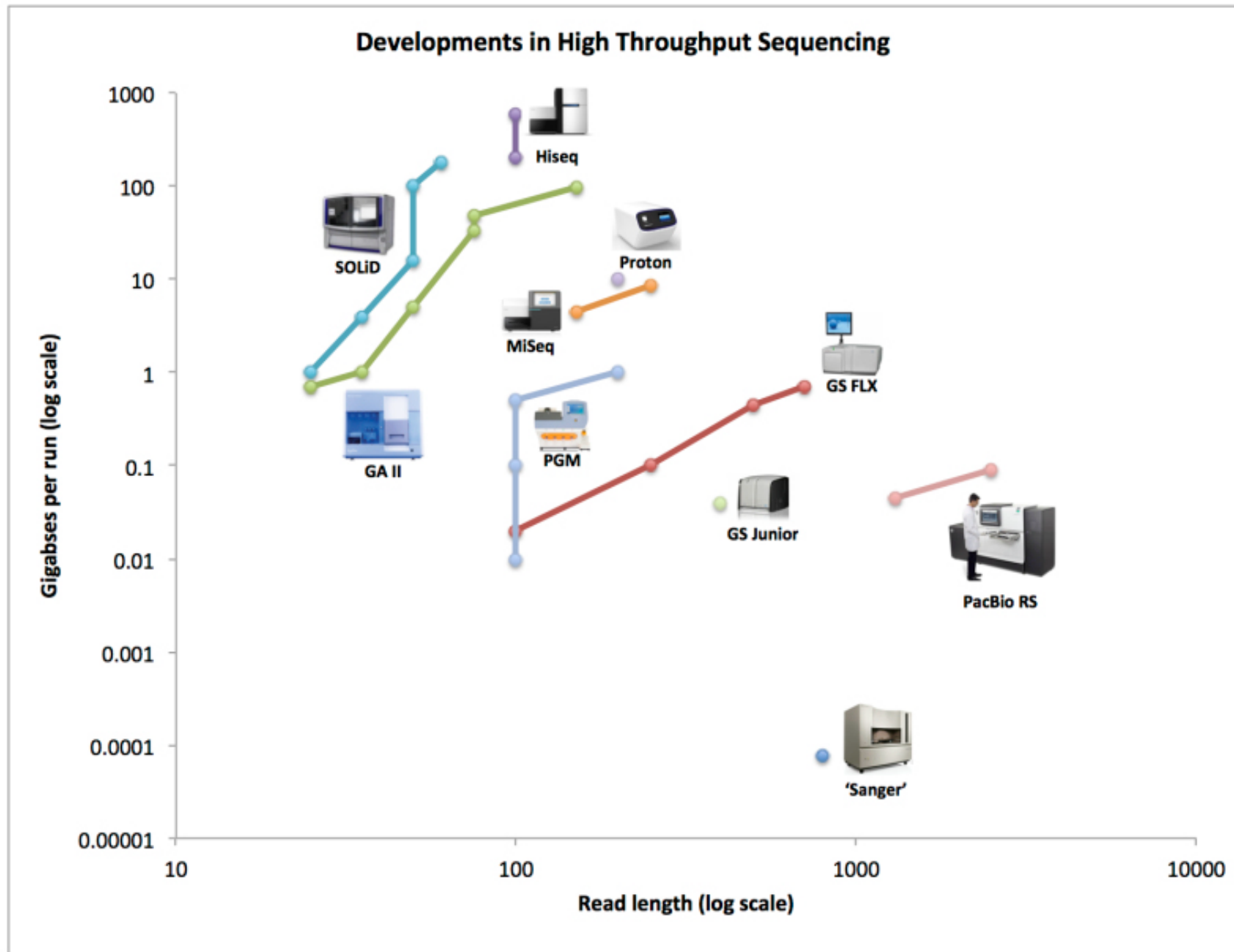
- Sequences of DNA or RNA from various sequencers technologies and from many labs
  - \* Some genome fragments : one or more genes, intergenic sequences, parts of a genome
  - \* Completed genomes
  - \* mRNA, tRNA, rRNA (ie. 16s)
- Annotations

# Data release

GenBank Size (GenBank.txt/gbrel.txt)



# Sequencing



## 2 conceptions of the sequences finishing:

- **The genome sequence must be completed and with a high quality 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.

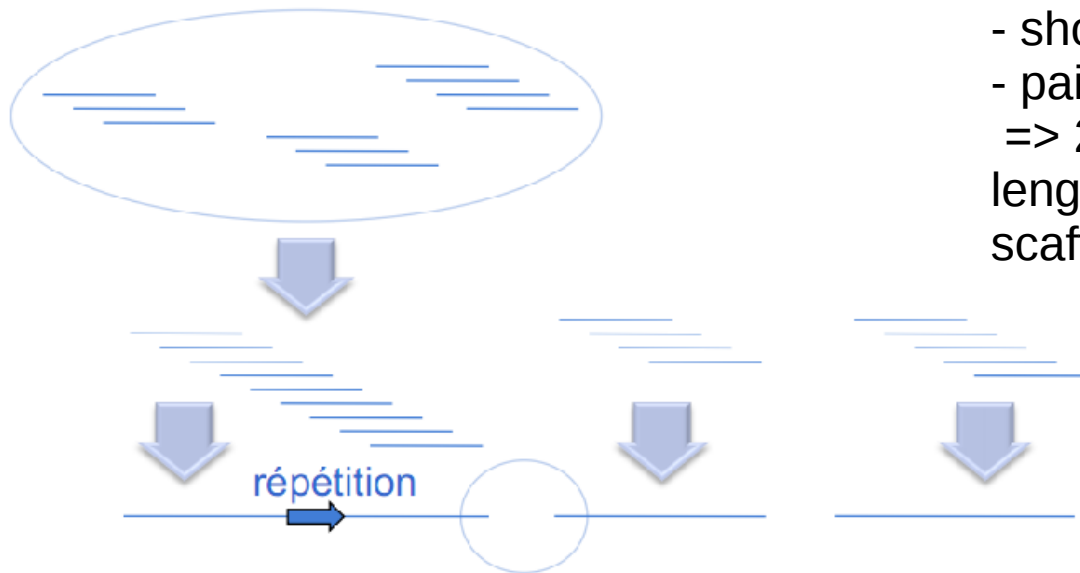
- **The sequence should be uncompleted with a draft quality, whether most of the genes are sequenced and identified.**

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

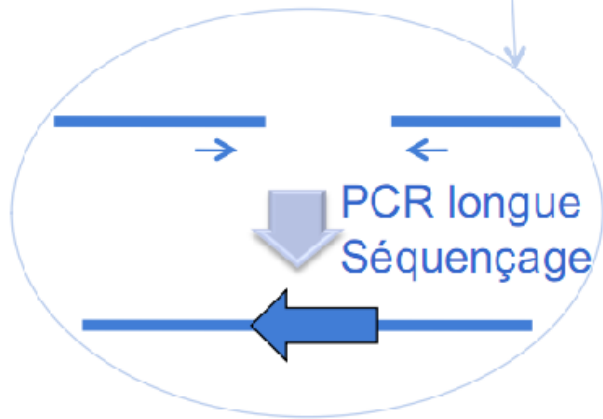
=> In general, fundamental research usually performs high quality genomes and applicative research (industry, our lab) usually performs draft genomes

Construction of a library of genomes fragments:

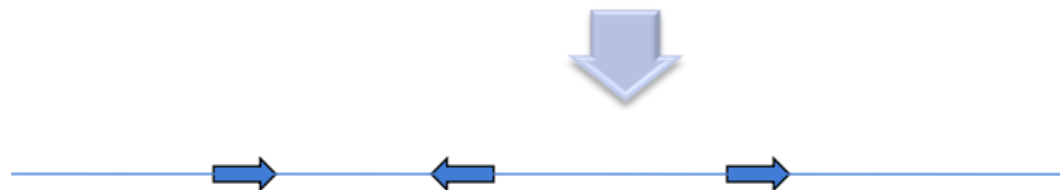
- shotgun = single fragments
- paired-end (or mate-pair)  
=> 2 fragments linked by an insert of a known length (~5 kb for 454 or Illumina), needed for scaffolding



Assembly process => construction of contigs (and scaffolds) from reads



*In silico* finishing + PCR to fill gaps



Checking with a closed reference or using annotations

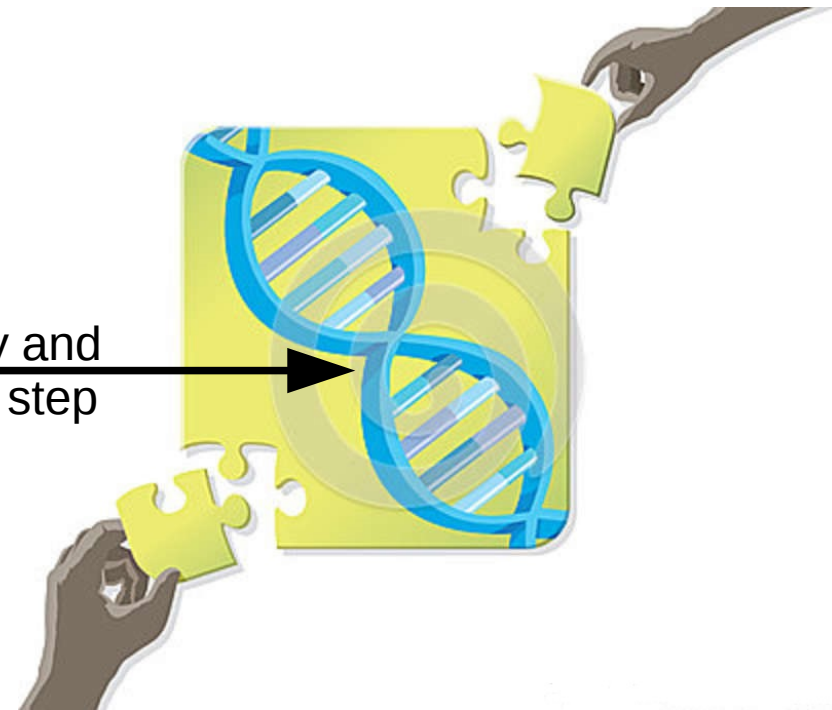
# Sequencing




Sequencing step: reads have heterogeneous distribution



Assembly and finishing step






  
 ATCGATGCGTAGCAGACTACCGTTACGATGCCTT...
   
 TAGCTACGCATCGTCTGATGGCAATGCTACGGAA...

Fragmentation + sequencing  
=> sets of reads


  
 TAGCTACGCATCGT
   
 ATCGATGCGTAGC
   
 TAGCAGACTACCGTT
   
 GTTACGATGCCTT

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

Build of contigs with overlapping regions

CGATGCGTAGCA
   
 ATCGATGCGTAGC
   
 TAGCAGACTACCGTT
   
 GTTACGATGCCTT

Assembly :  
=> alignements of reads + consensus

↓
   
 .....ATCGATGCGTAGCAGACTACCGTTACGATGCCTT.....

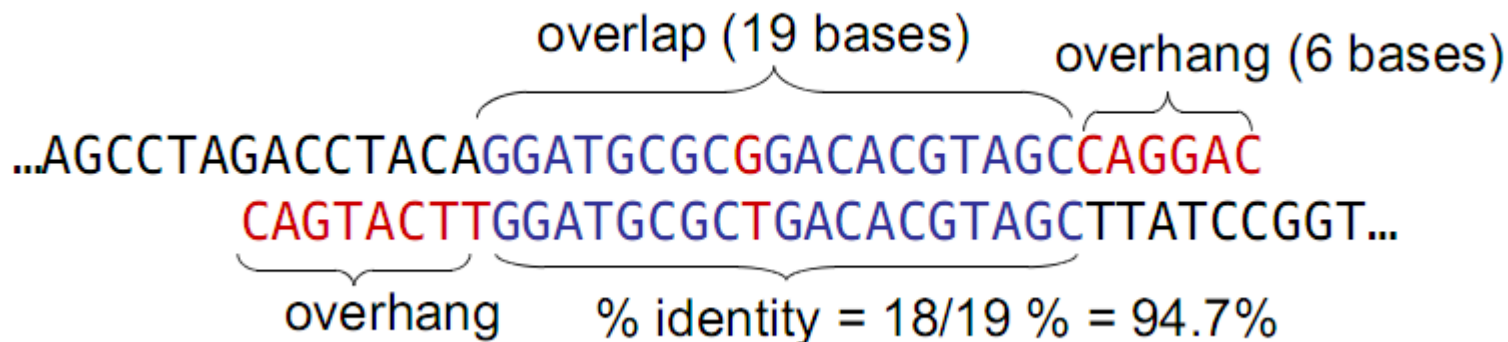
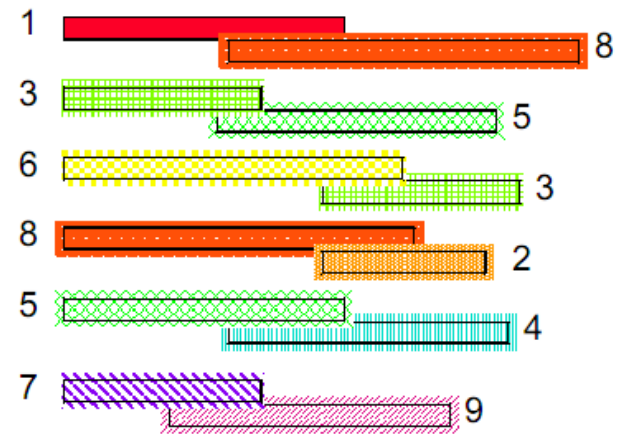
## Search for best pairings

- Compare each sequence (and its reverse complement) against every others sequence to find the best overlapping

=> list of best candidates with similarities criteria

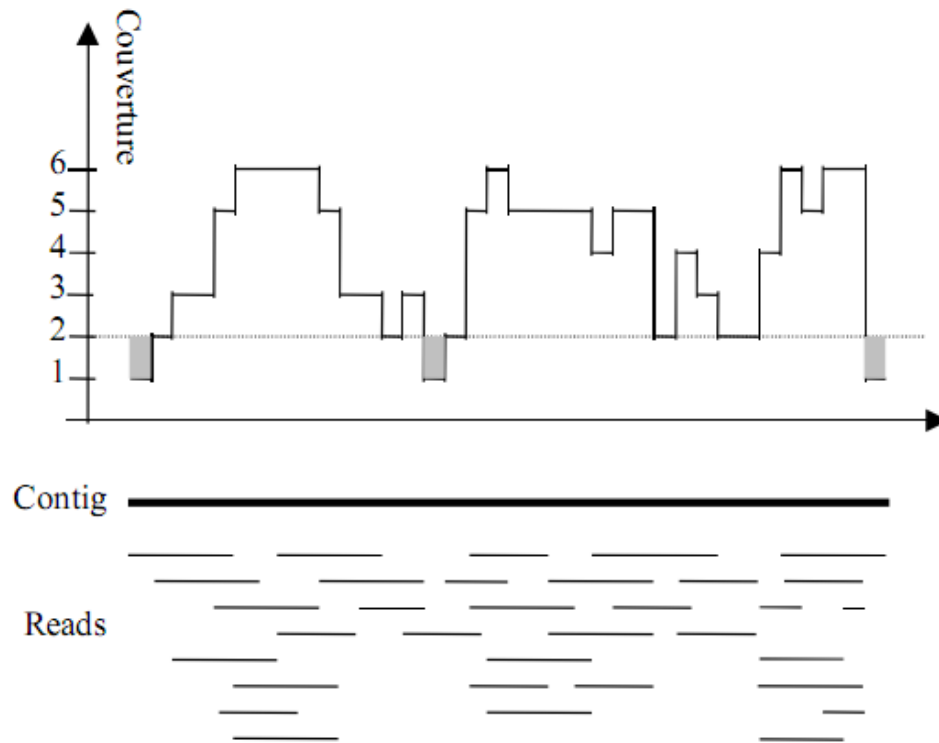
Best candidate is a compromise between :

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



## Main remaining problems:

- Bad assembly of reads
- Low coverage of reads



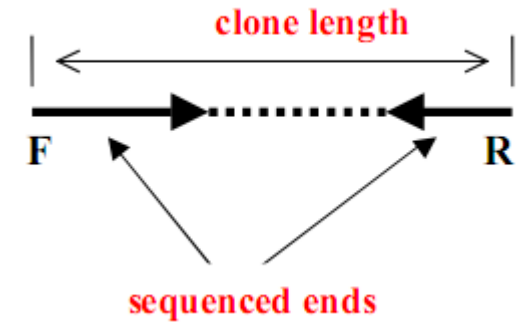
- Bad insert size estimation
- Different orientation of contigs
- Error of sequencing
- Repeat sequence ambiguities

Remap reads on a reference genome or assembled genome itself  
Highlighting errors, ie. sequencing error or SNP, show coverage

The screenshot shows the EagleView software interface. At the top, the title bar reads "EagleView". Below it is a menu bar with "File", "Configure", "Preferences", and "Help". A navigation bar contains navigation buttons (left, right, home, end) and a position field showing "5025". To the right, the reference sequence is identified as "MTDNA-C6 5187". Below the navigation bar, a coverage bar is visible with a tooltip that reads "=> P=33, Q=0". The main area displays a sequence alignment with multiple rows of reads. The reference sequence is shown at the top, and the reads are aligned below it. A mouse cursor is pointing to a specific position in the alignment, and a tooltip displays the coverage information "=> P=33, Q=0". The alignment shows several reads with some mismatches highlighted in red, indicating sequencing errors or SNPs. The reads are aligned to the reference sequence, and the coverage bar above them shows the number of reads at each position. The reference sequence is shown in a light blue font, and the reads are shown in a dark blue font. The alignment is shown in a grid format, with the reference sequence on the left and the reads on the right. The coverage bar is shown above the reads, and the tooltip is shown below the mouse cursor.

## Reads paired-end (similar to mate-pair)

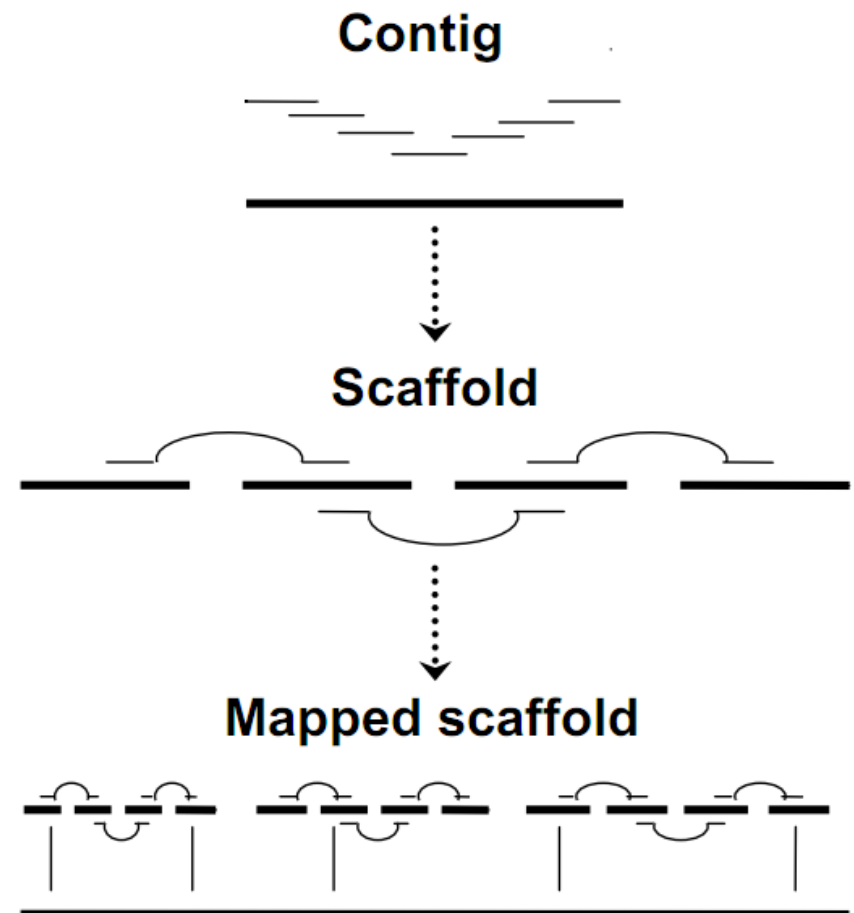
- The distance between the reads is known (length of the insert), with some experimental uncertainty
- Distance of insert depends of technology (454 or Illumina => 3-8 kb)

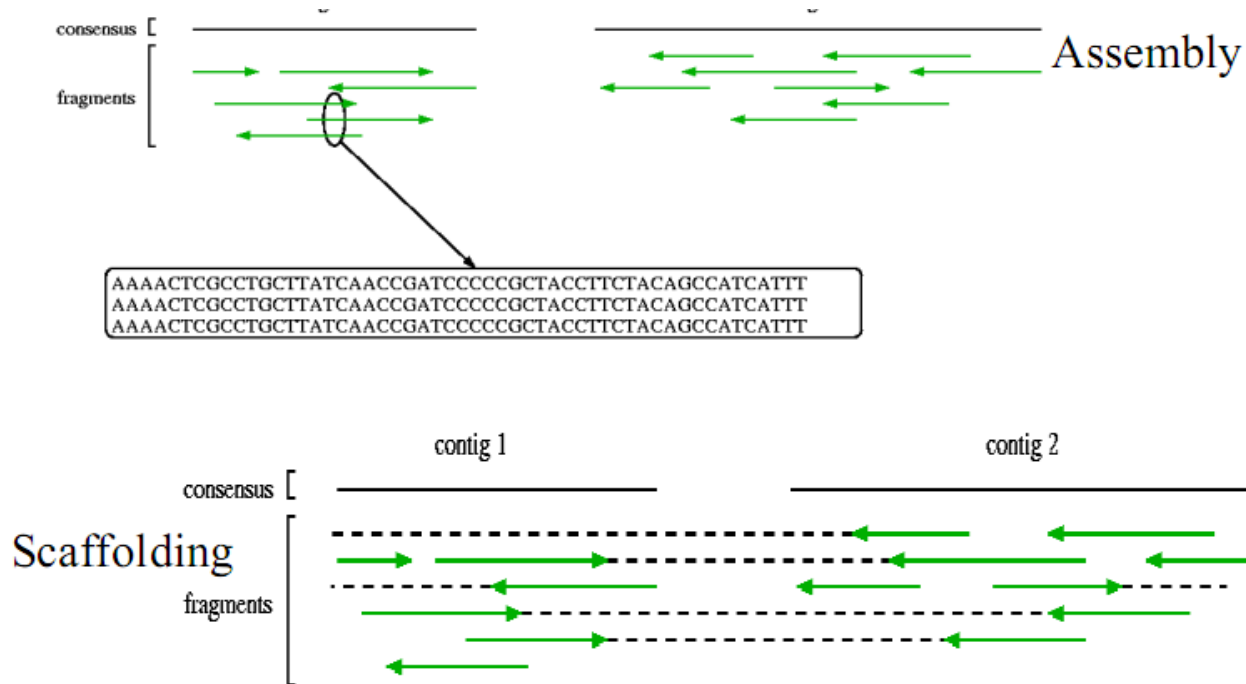


**Contigs** : group of overlapping reads, without gap

**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 !





## Finishing :

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)
- Fill the gaps by extending the boundaries of each gap using ends of mapping reads
- Order (or reorder) contigs
- Desassemble some areas if they seem to be false