

<https://ngstc.iutms.umontpellier.fr/formations/m2bs/TP-rnaseq/>

# L'ère NGS....

# NGS et Transcriptomics

- Introduction: l'ère NGS et état des lieux, principaux séquenceurs (vu en intro)
- Les enjeux de l'analyse, qqls applications en diagnostic (ADN/ARN)
- Transcriptome et RNAseq, les questions biologiques
  - ✓ Les différentes stratégies d'analyse
  - ✓ Le mapping, exemple et limites
  - ✓ Annotations génomiques, transcriptome de référence
  - ✓ Comment passer à l'échelle

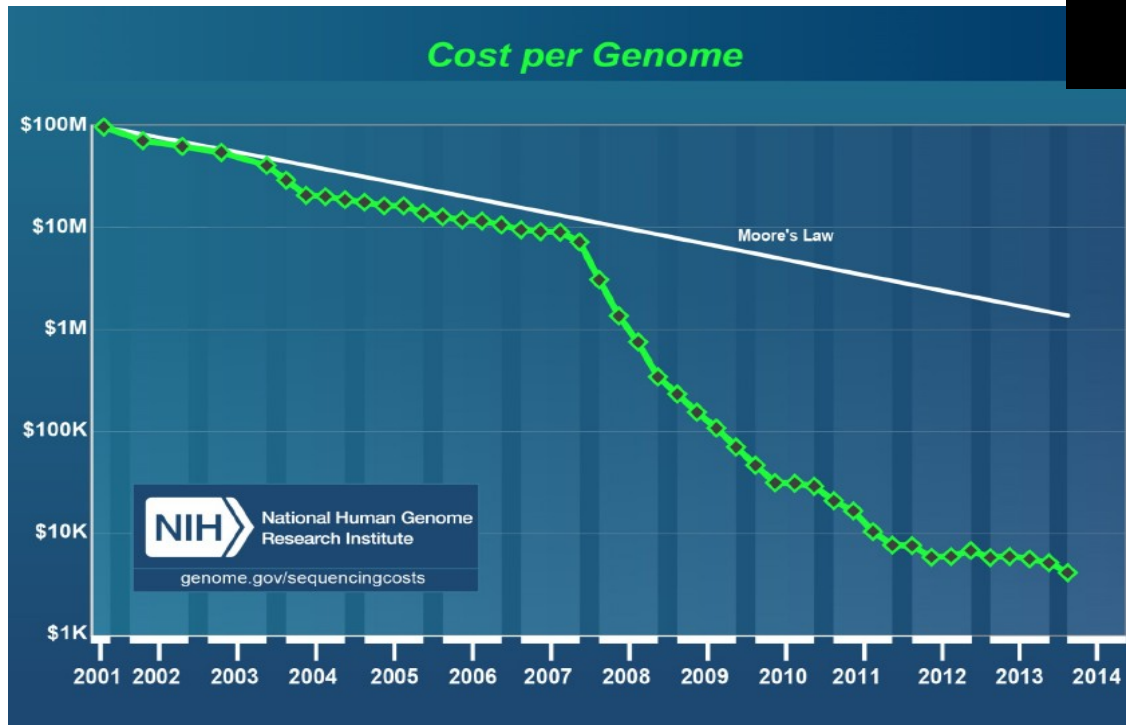
# NGS: une nouvelle ère

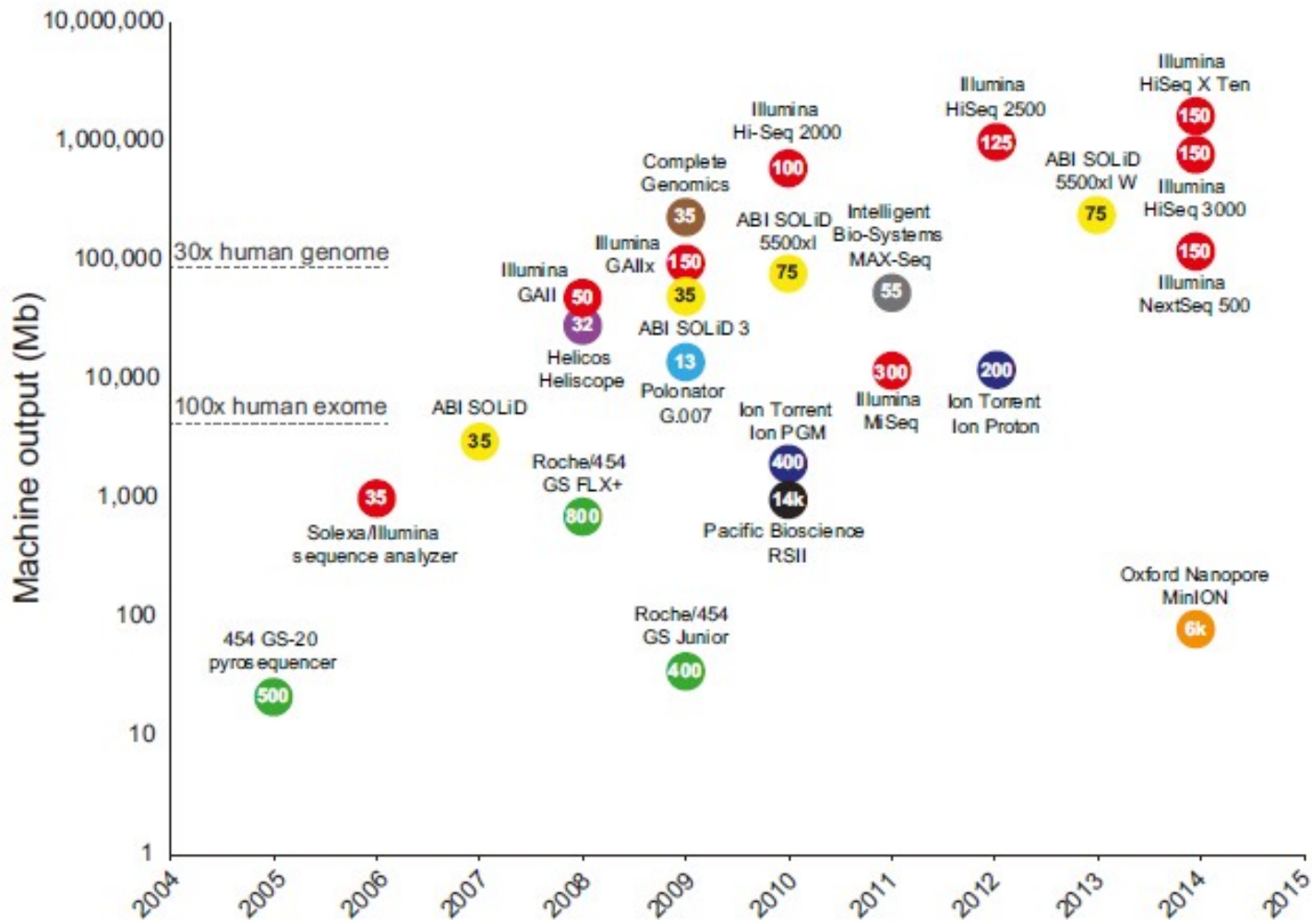
- Coût (qqls milliers d'euros)
- Temps (Un génome humain en quelques jours)
- Puissance (Milliards bases en quelques heures)

June 26th 2000: official announcement of the completion of the draft of the human genome sequence (truly finished in 2004)

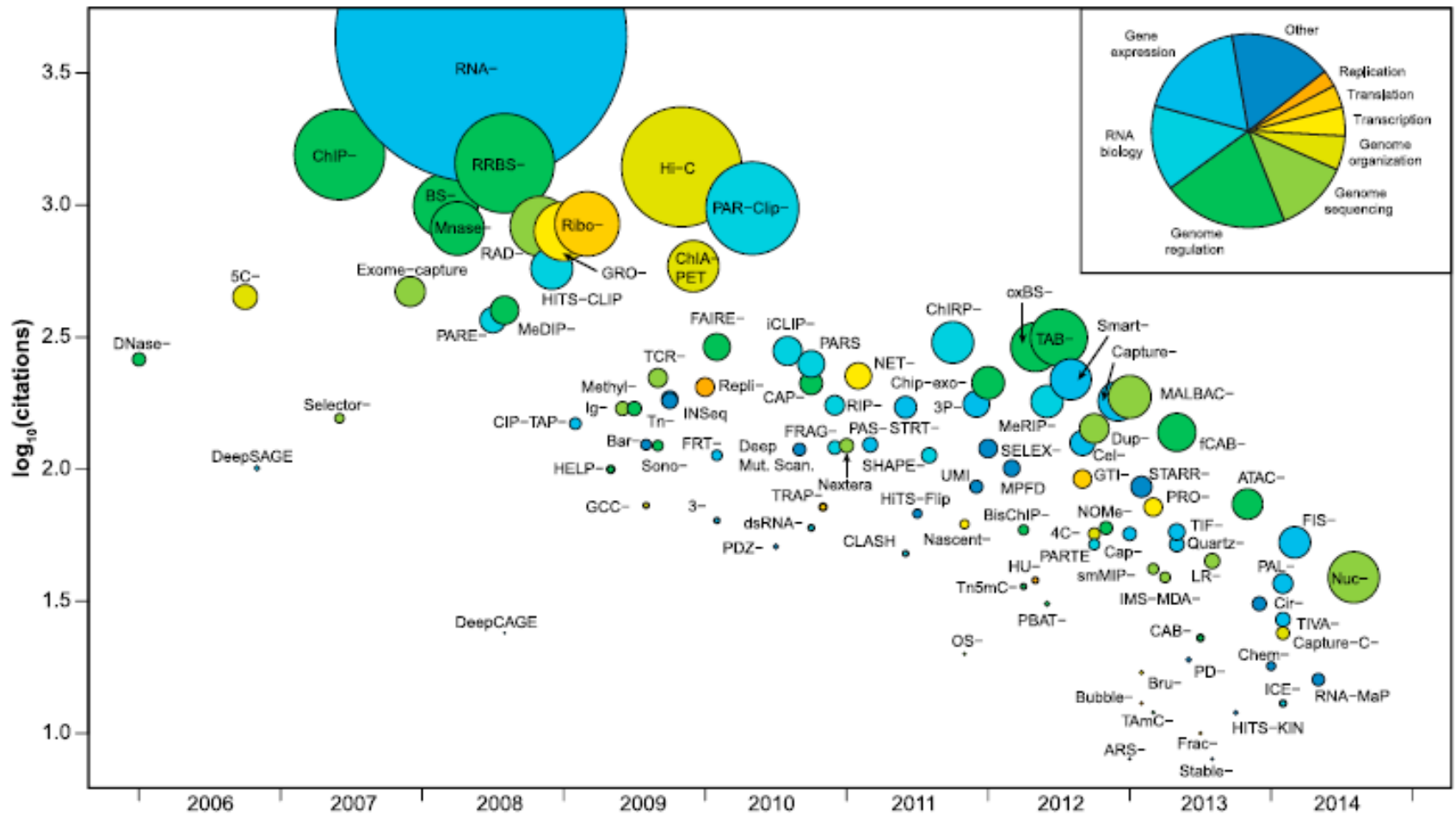


Costs:	
HGP:	Celera:
3 billion \$	200 million \$
15 years	2 years

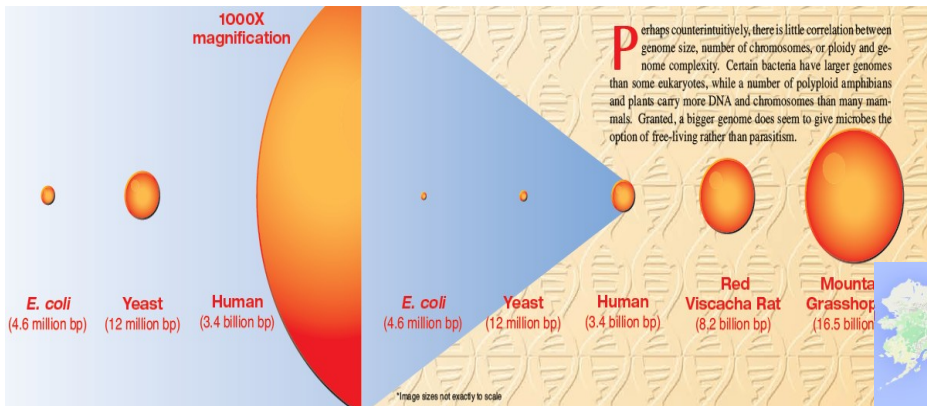




# 1-Explosion des séquenceurs NGS



## 2-Explosion des applications -Seq



## IGSR and the 1000 Genomes Project



## International network of cancer genome projects

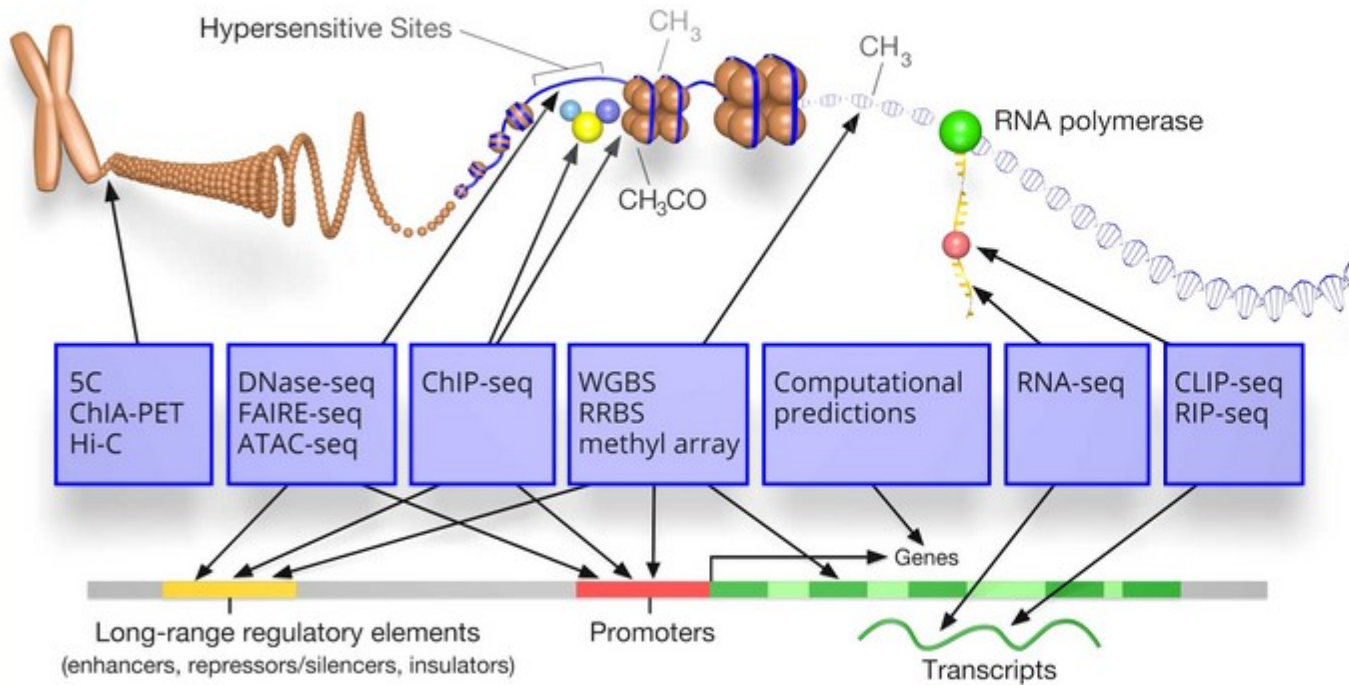
The International Cancer Genome Consortium\*

The International Cancer Genome Consortium (ICGC) was launched to coordinate large-scale cancer genome studies in tumours from 50 different cancer types and/or subtypes that are of clinical and societal importance across the globe. Systematic studies of more than 25,000 cancer genomes at the genomic, epigenomic and transcriptomic levels will reveal the repertoire of oncogenic mutations, uncover traces of the mutagenic influences, define clinically relevant subtypes for prognosis and therapeutic management, and enable the development of new cancer therapies.

## 3-Explosion des grands programmes de séquençage

# ENCODE: Encyclopedia of DNA Elements

T'



Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

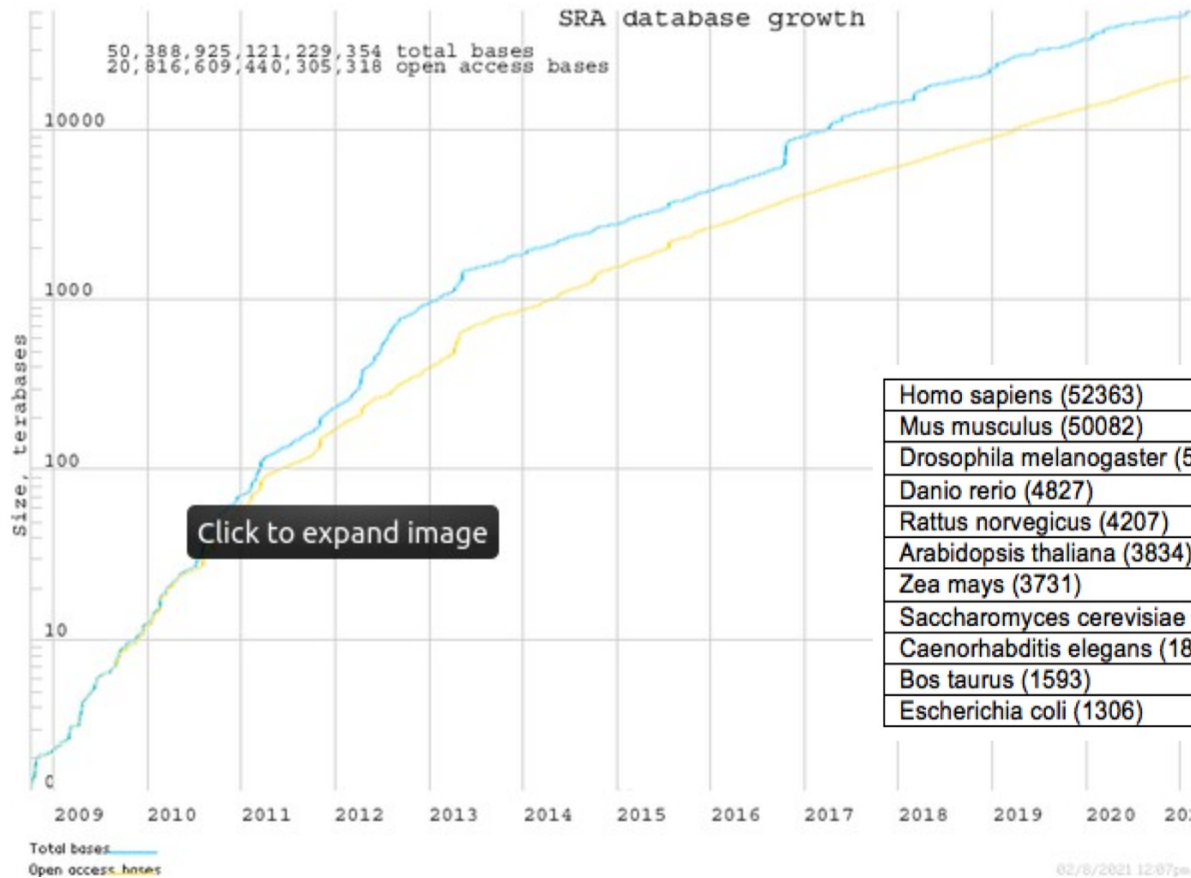
HUMAN

MOUSE

WORM

FLY

## 3-Explosion des grands programmes de séquençage

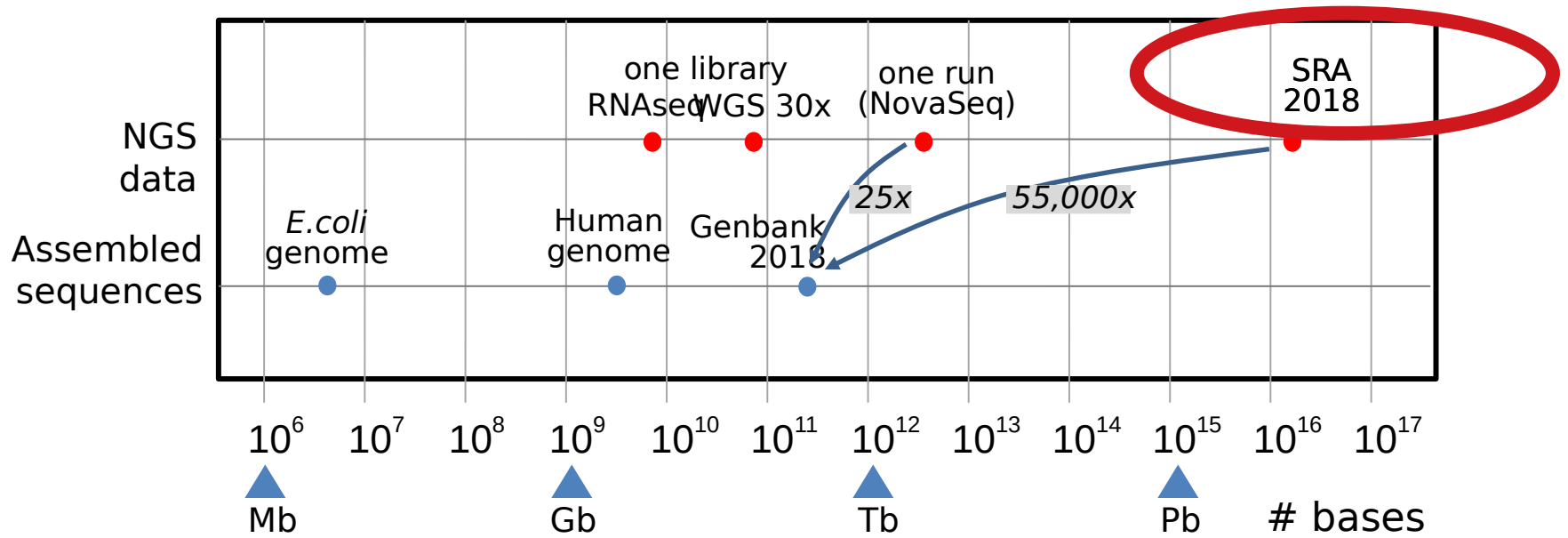


Homo sapiens (52363)	Oryza sativa (1210)
Mus musculus (50082)	Macaca mulatta (942)
Drosophila melanogaster (5256)	Gallus gallus (904)
Danio rerio (4827)	Zaire ebolavirus (867)
Rattus norvegicus (4207)	human metagenome (861)
Arabidopsis thaliana (3834)	Glycine max (811)
Zea mays (3731)	Mycobacterium tuberculosis (780)
Saccharomyces cerevisiae (3321)	Solanum lycopersicum (761)
Caenorhabditis elegans (1814)	Equus caballus (759)
Bos taurus (1593)	All other taxa (53604)
Escherichia coli (1306)	

Top20 des données RNA-seq dans  
SRA par espèces

## 4-Explosion des données publiques

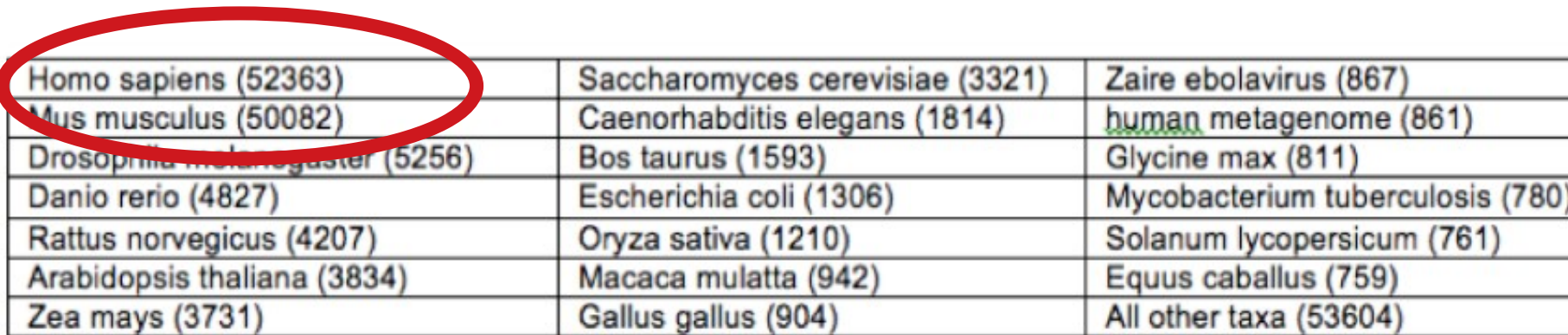




## 4-Explosion des données publiques

# Transcriptome and RNAseq

- Public datasets in SRA



Homo sapiens (52363)	Saccharomyces cerevisiae (3321)	Zaire ebolavirus (867)
Mus musculus (50082)	Caenorhabditis elegans (1814)	human metagenome (861)
Drosophila melanogaster (5256)	Bos taurus (1593)	Glycine max (811)
Danio rerio (4827)	Escherichia coli (1306)	Mycobacterium tuberculosis (780)
Rattus norvegicus (4207)	Oryza sativa (1210)	Solanum lycopersicum (761)
Arabidopsis thaliana (3834)	Macaca mulatta (942)	Equus caballus (759)
Zea mays (3731)	Gallus gallus (904)	All other taxa (53604)

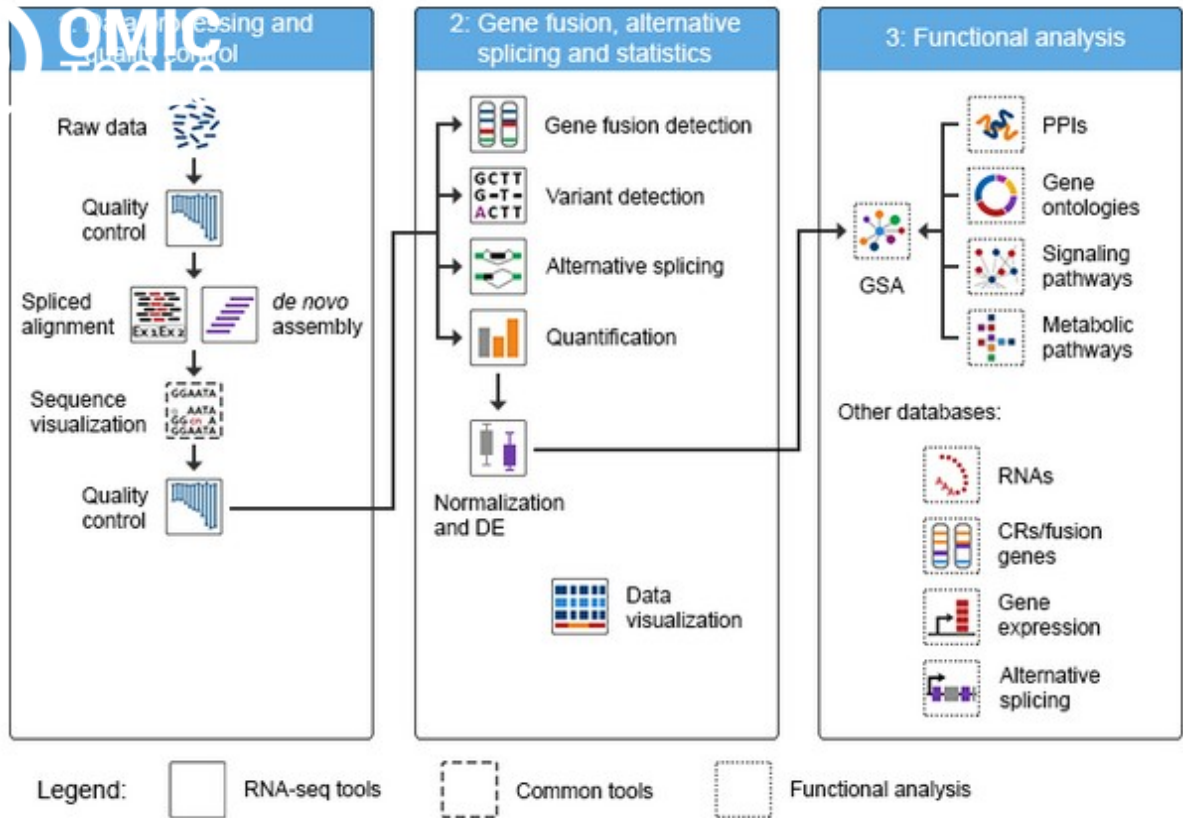
*Figure 1: The top 20 species ranked by number of RNA-seq libraries (parenthesis) available in the SRA database*

TGCA | 10696

GTEX | 8555

ENCODE | 657

GEUVADIS | 465



<https://omictools.com/>

## 5-Profusion d'Outils Bioinformatiques

## RNA-SEQ ANALYSIS APPLICATIONS



### Experimental design

9 tools



### Spliced read alignment

52 tools



### Assembly evaluation

2 tools



### Read count

5 tools



### Transposable element expression

3 tools



### Alternative splicing

60 tools



### Driver gene fusion prediction

4 tools



### Sex-linked gene detection

2 tools



### Variant detection

16 tools



### Circular RNA detection

15 tools



### HLA genotyping

7 tools



### Co-expression network analysis

2 tools



### Quality control

31 tools



### De novo transcriptome assembly

32 tools



### Read realignment

3 tools



### Transcript quantification

82 tools



### Normalization/differential expression

89 tools



### Gene fusion detection

36 tools



### Gene prediction

6 tools



### Transcriptome annotation

9 tools



### Allele-specific expression

17 tools



### RNA editing

10 tools



### Alternative polyadenylation

9 tools

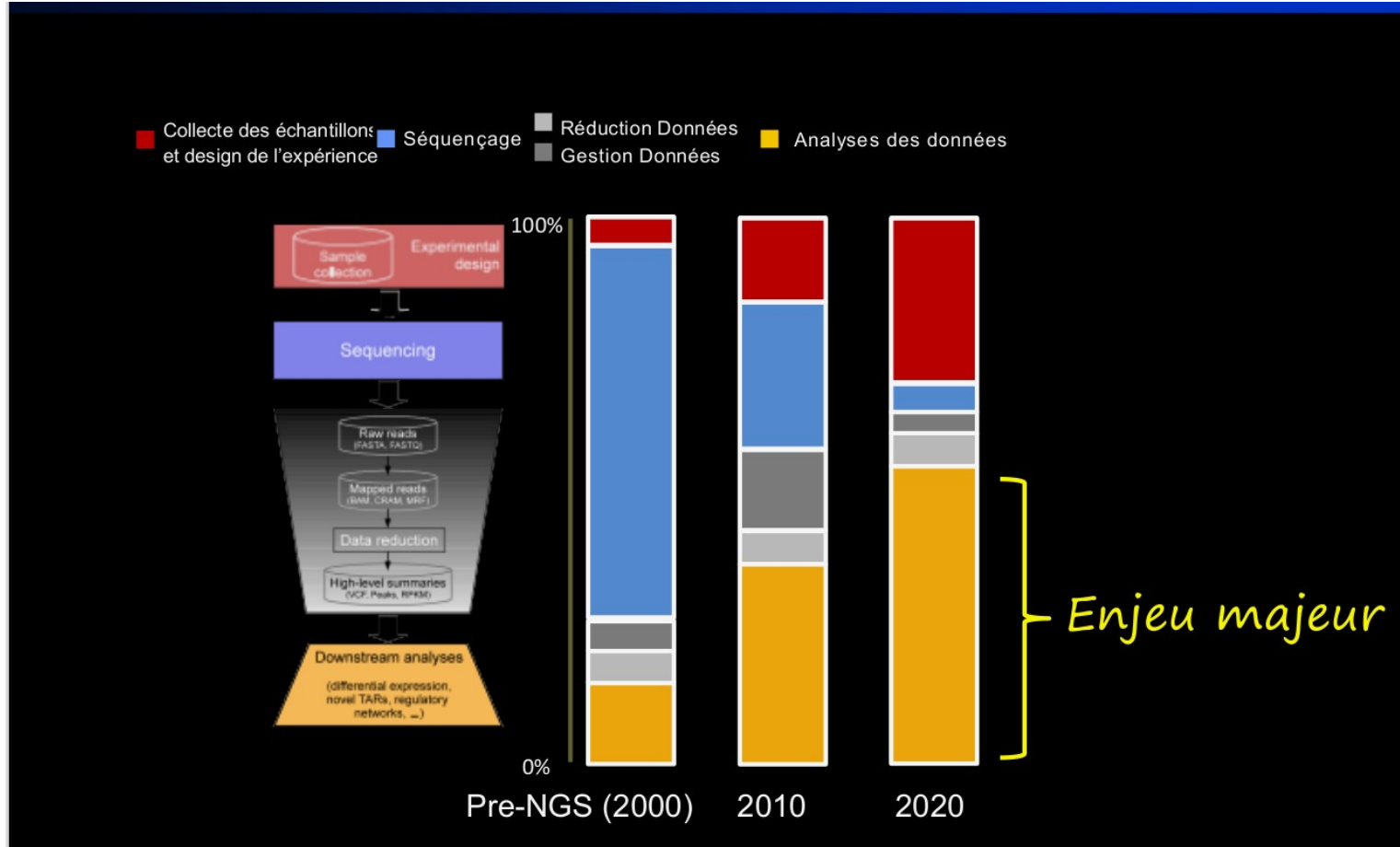


### Gene expression clustering

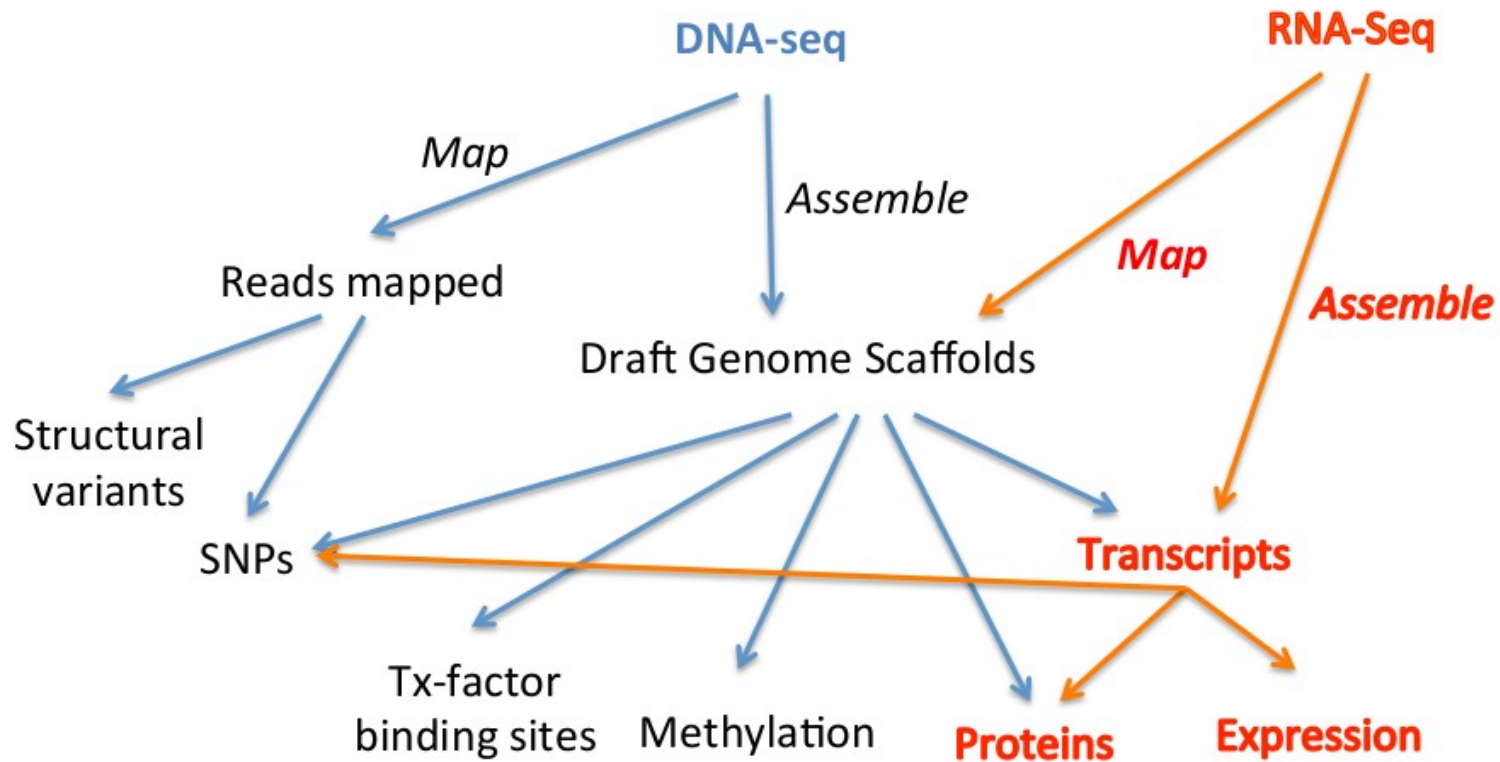
7 tools

<https://omictools.com/whole-exome-sequencing-category>

# Quelques enjeux de l'analyse des données



# Applications



Genome et Transcriptome / Whole DNAseq and RNAseq

Les points communs de l'analyse

Les spécificités

Les champs d'application exemples en diagnostic humain

Quelques limitations

# Principales Applications en diagnostic moléculaire

## **ADN (Panel de gènes ciblés, Exomes, Whole genome ...)**

Mutations (Substitutions, Indels, Tandem repeats..)

Anomalies chromosomiques (translocations, Inversions, gènes de fusion)

Variants structuraux à large échelle (CNV)

ADN circulant

## **ARN (RNA-seq)**

Expression des gènes codants

Long non-codingRNA (lncRNA), microRNA (miRNA), ARN circulaires (circRNA)...

Variants d'épissage ...

ARN de fusion...

Indels, mutations.....

## **ADN&Epigénome, Méthylation, Hi-C (prochaine étape)**

Métagenomique (Microbiologie, Virologie, Microbiome..)

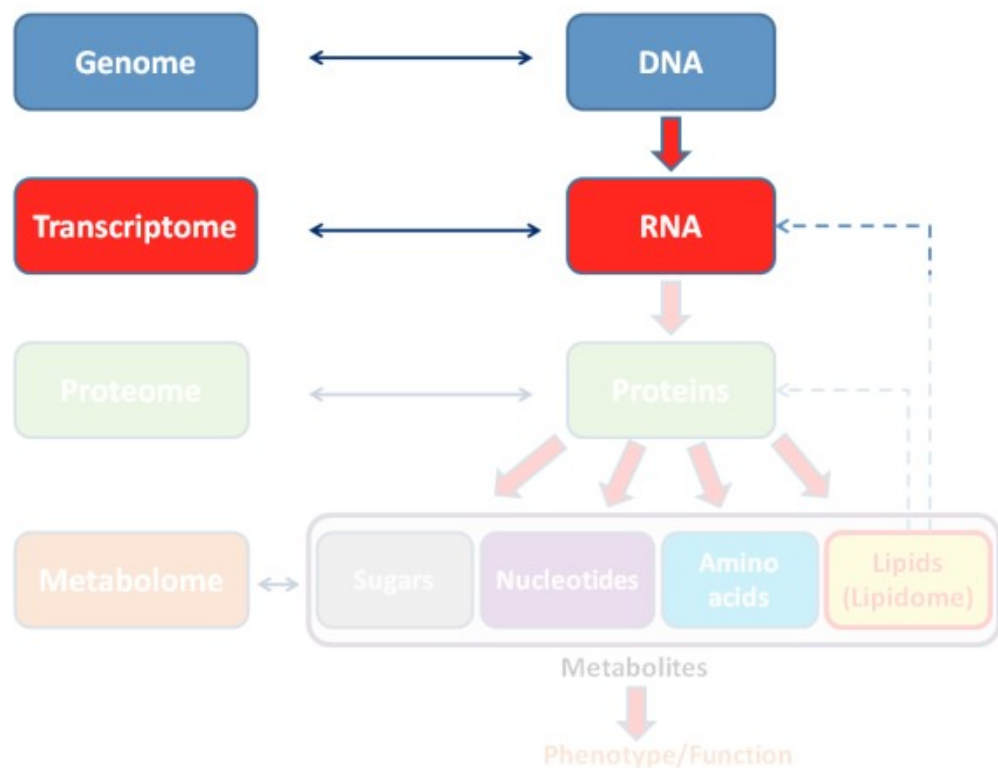


# DNA-RNA-Applications

- |  |  |   |
|--|--|---|
| <ul style="list-style-type: none"><li>• Panel</li></ul>    | <ul style="list-style-type: none"><li>• Exome</li></ul>  | <ul style="list-style-type: none"><li>• Whole genome</li></ul>      |
| ADN/ARN  | <ul style="list-style-type: none"><li>• RNAseq</li></ul> |   |
| <ul style="list-style-type: none"><li>• 50-500Mo</li></ul> | <ul style="list-style-type: none"><li>• 1-4Go</li></ul>  | <ul style="list-style-type: none"><li>• 40-400Go</li></ul>          |
| <ul style="list-style-type: none"><li>• labTop</li></ul>   | <ul style="list-style-type: none"><li>• server</li></ul> | <ul style="list-style-type: none"><li>• server/<br/>cloud</li></ul> |

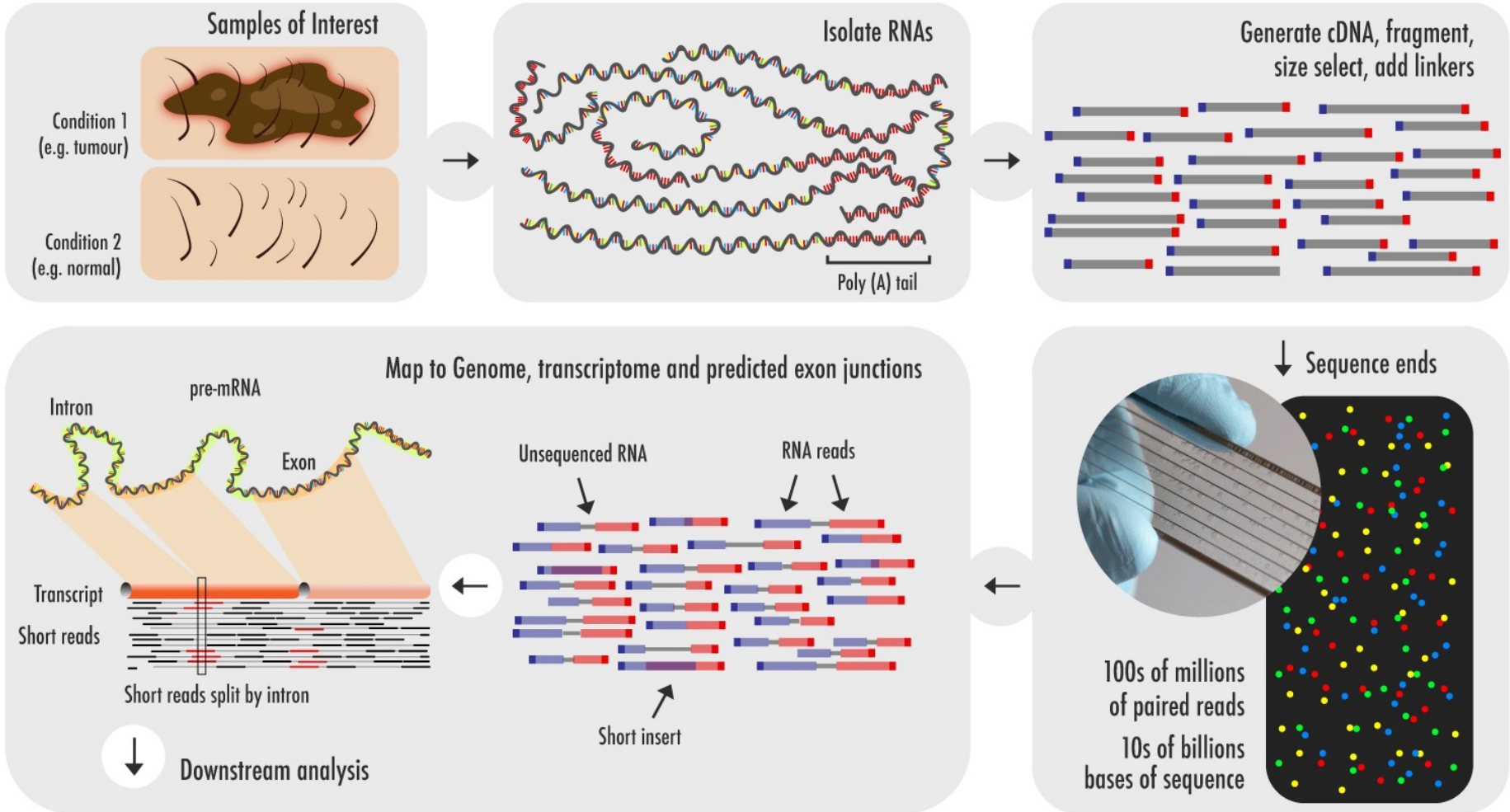
# Transcriptomics

**Transcriptome** = complete set of all RNA molecules ("transcripts") produced from a genome **OR** specific subset of transcripts present in a particular cell type or under specific growth conditions



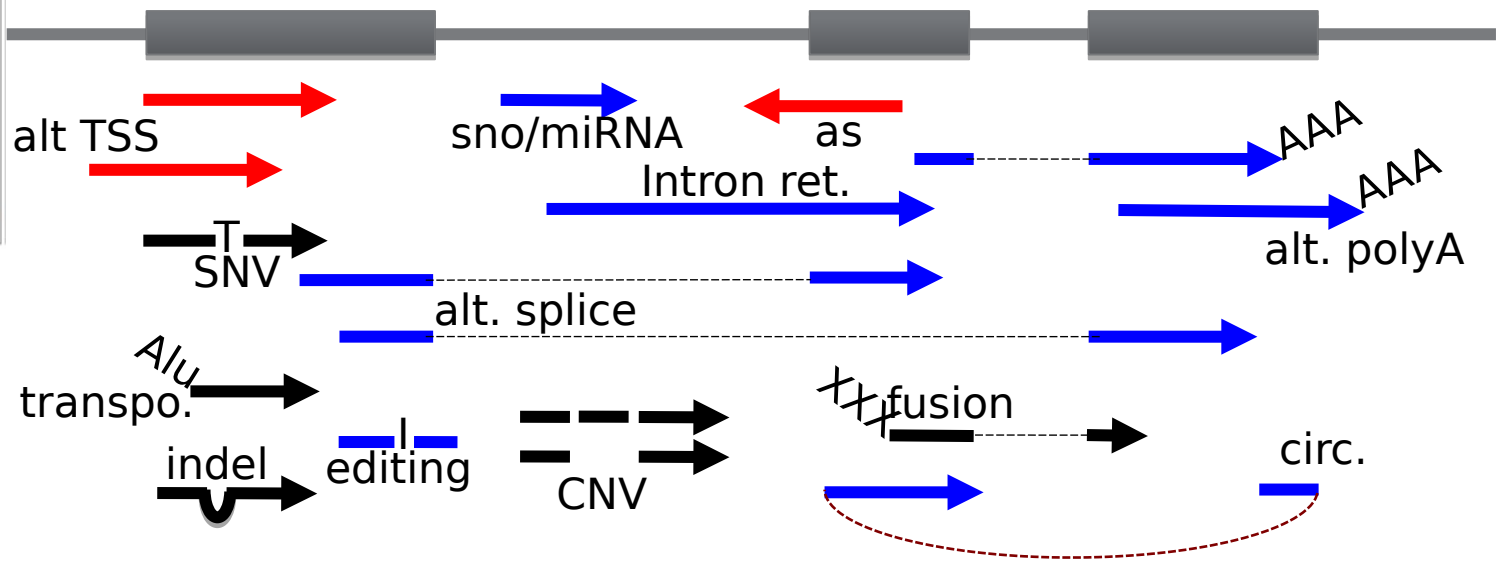
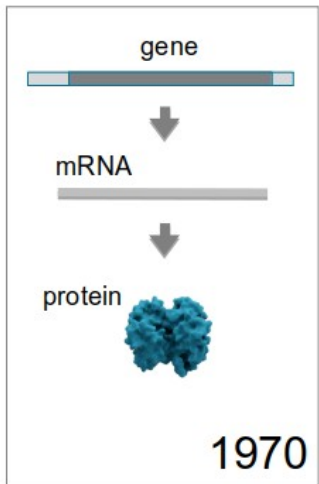
**Transcriptomics** involves large-scale analysis of RNAs to follow when, where, and under what conditions genes are expressed.

# Transcriptome and RNAseq



A simple method

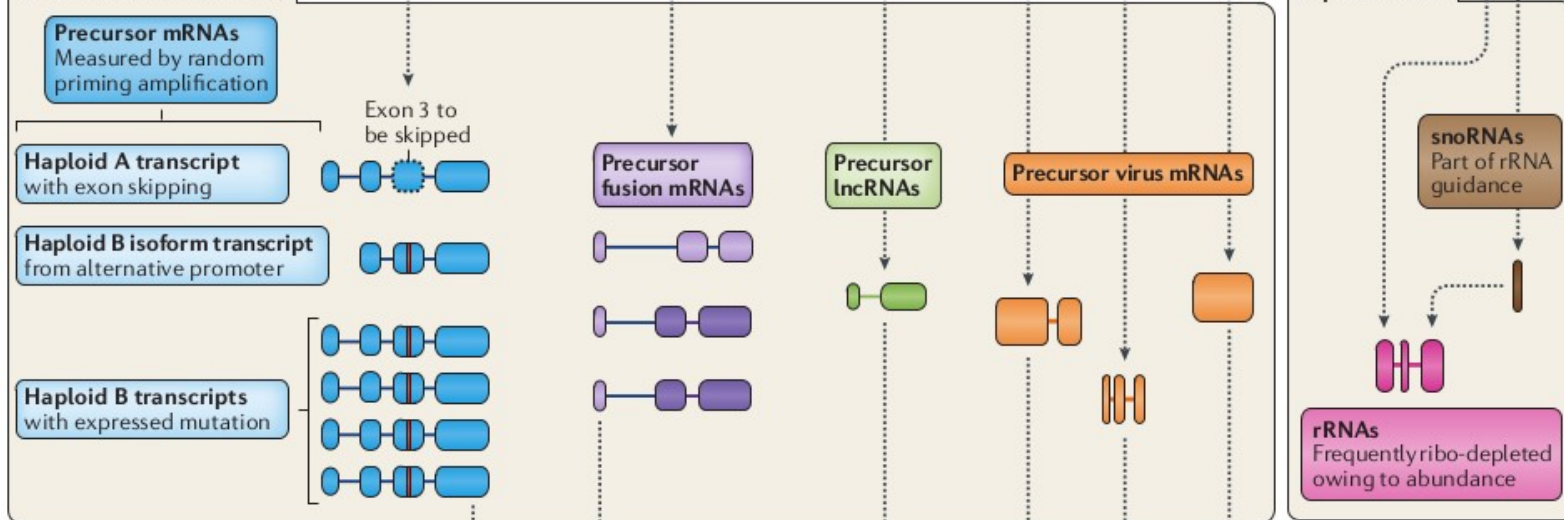
# The Real Transcriptome is combinatorial



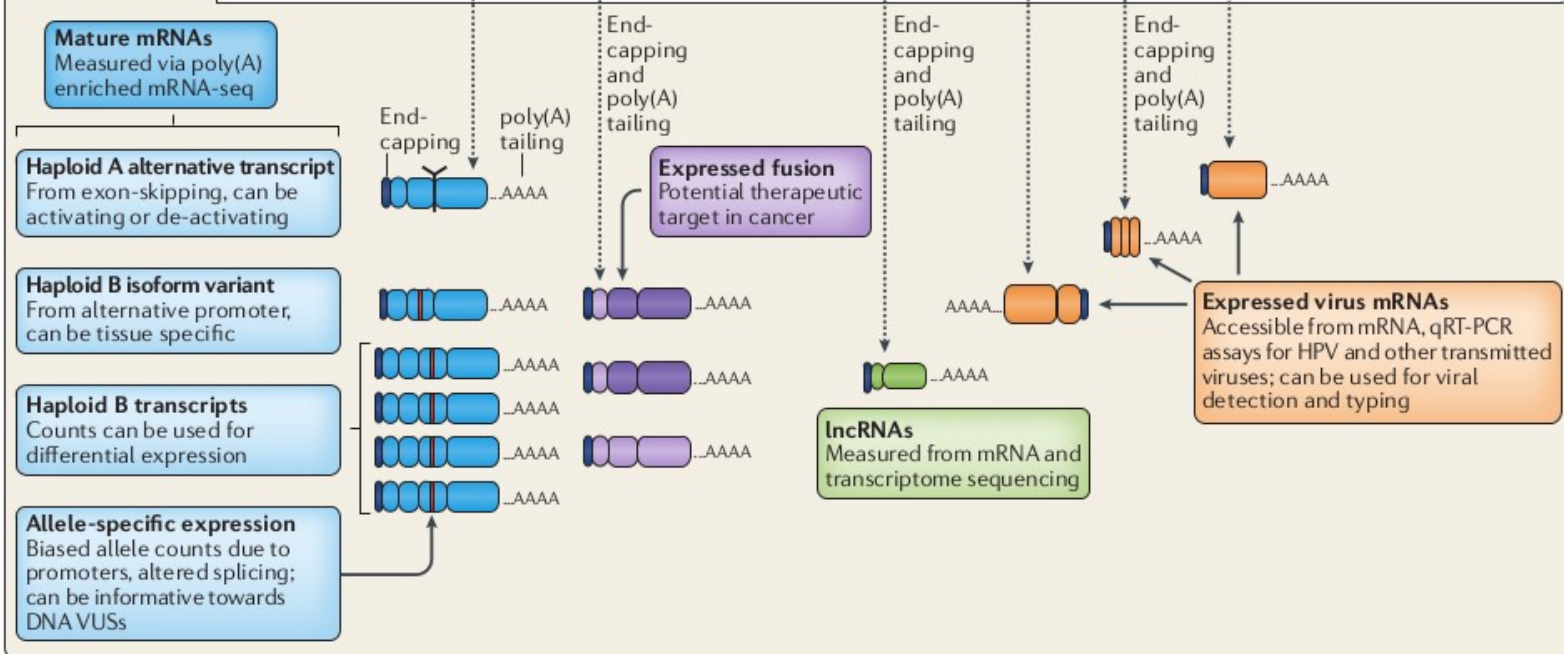
**Genetic** X **Transcriptional** X **Post-transcriptional variation**

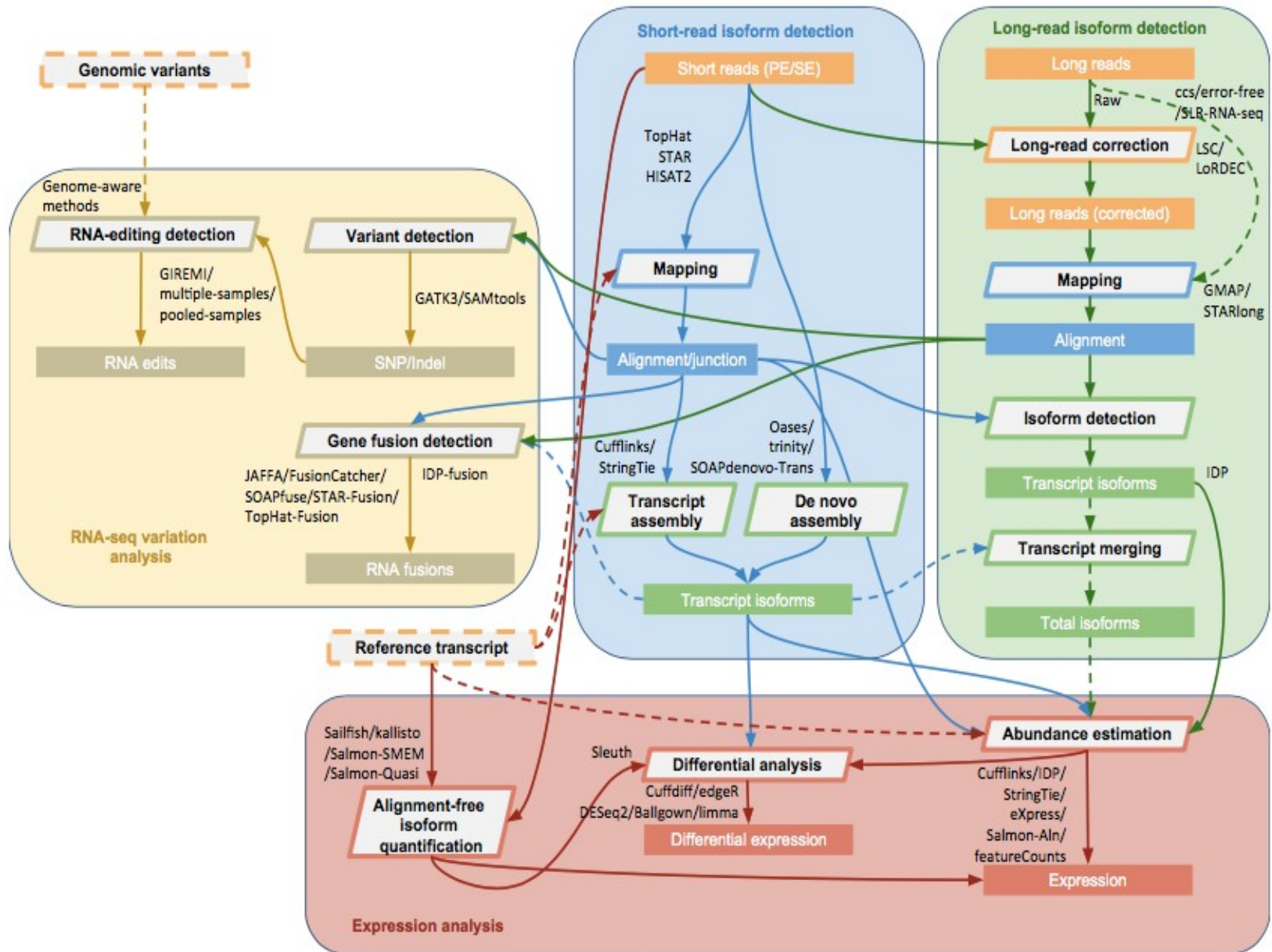
DNA → RNA

### Transcriptome sequencing



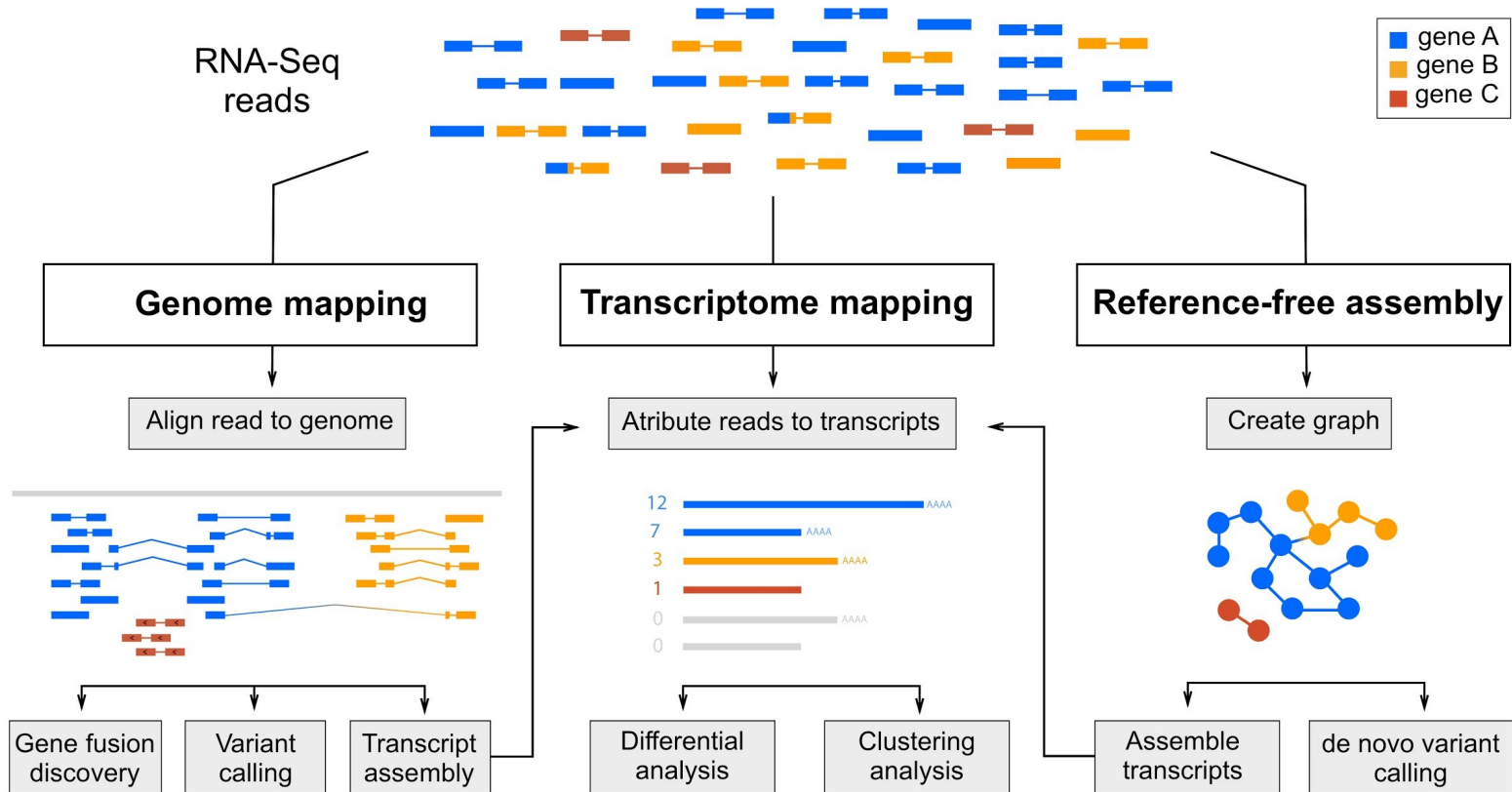
### mRNA sequencing





.RNAseq analysis (nat genetics)

# Transcriptome and RNA-seq



A complex analysis





# Enjeux de l'analyse

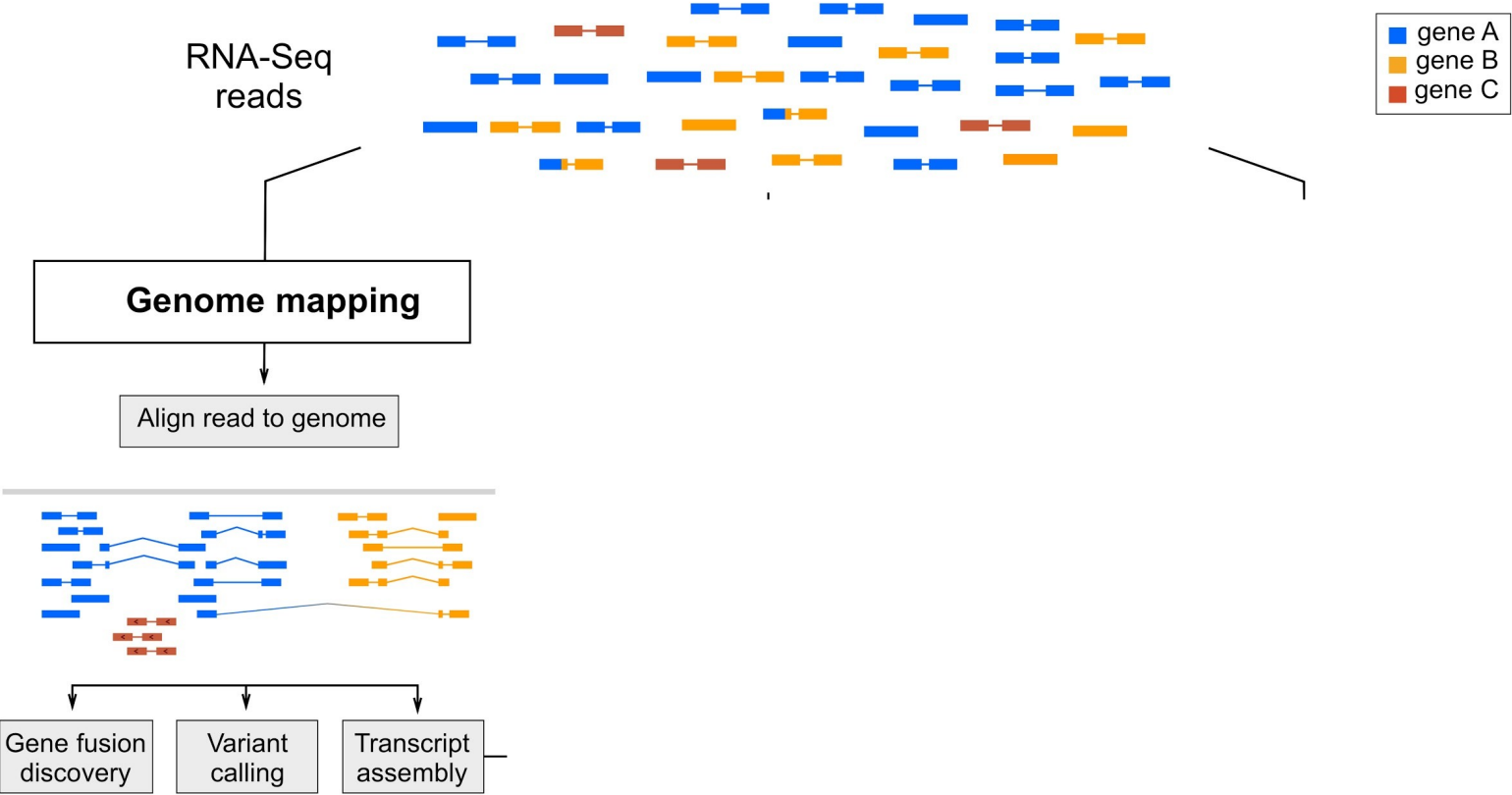
## **De nombreux softwares :**

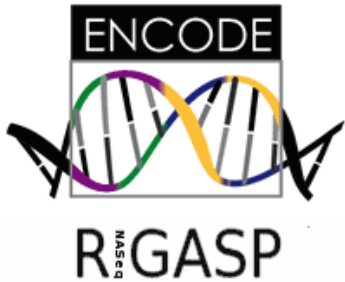
Comment choisir le meilleur et le plus approprié à une question biologique donnée?

Du principe algorithmique à la biologie (et inversement)

Exemple de mapping, limites..

# RNA-seq and genome mapping





## RGASP Round 3: RNA-seq Read Alignment Assessment

One of the lessons learned from rounds 1 & 2 of the project was that the initial step of aligning the reads has a major influence on the quality of gene predictions produced. Therefore, a third round of RGASP was conducted to focus primarily on read mapping to the genome.

The project was related to the "Sequence Mapping and Assembly Assessment Project (SMAAP)", a collaborative effort to compare and evaluate methods and strategies for de novo genome assembly (dnGASP) and RNA-seq read alignment (RGASP3) using data from second generation sequencing platforms.

RGASP3 is organised by [Paul Bertone \(EBI\)](#) with input from the Wellcome Trust Sanger Institute and the CRG. [[Contact](#)]

### Goal of RGASP 3

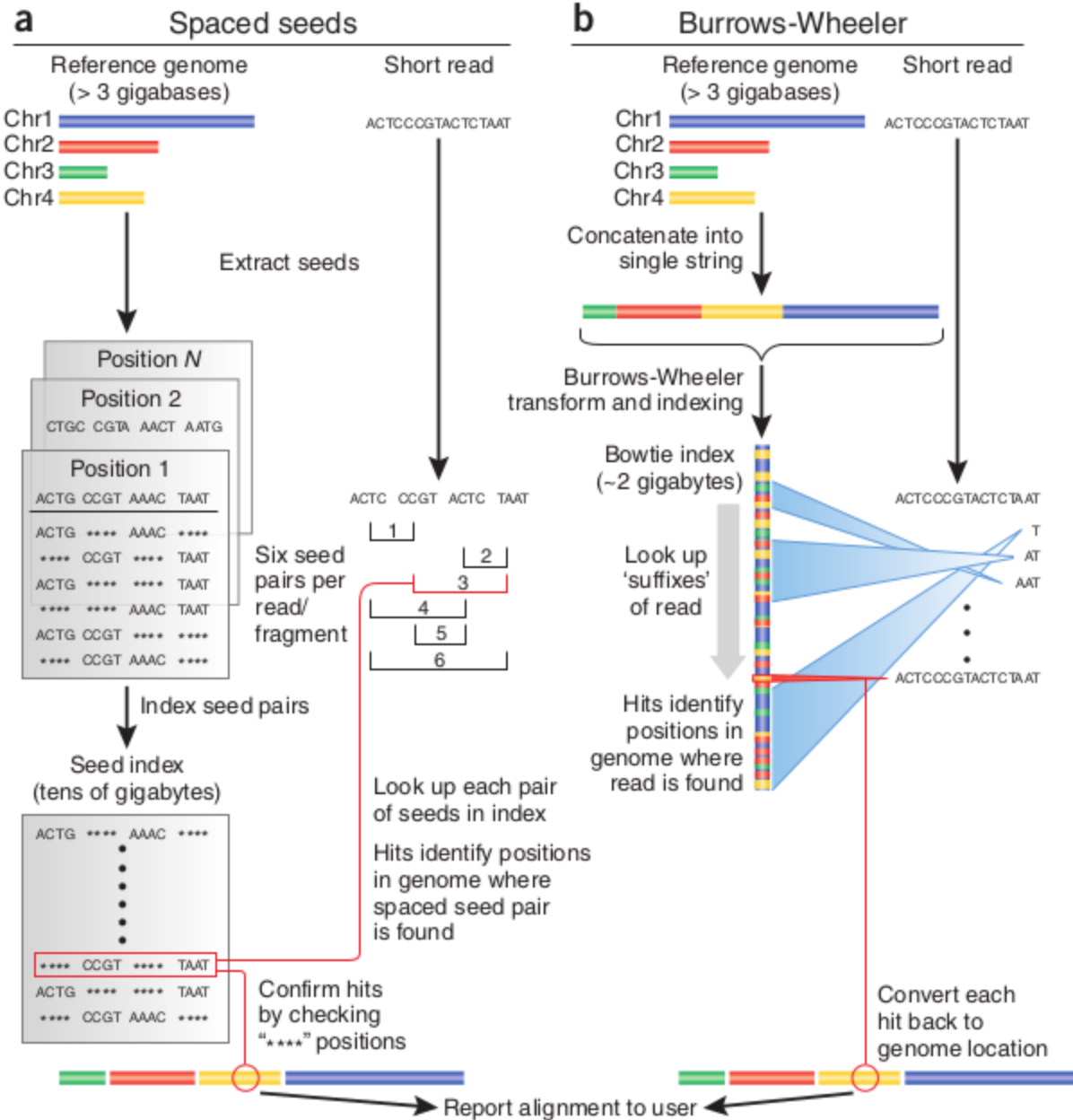
The principal aim of the RGASP3 project is to allow an unbiased evaluation of different analysis methods within the community generating **high-quality RNA-seq read alignments** that can be used for efficient transcriptome characterization (transcript discovery and quantitation). A total of 26 spliced alignment protocols based on 11 programs and pipelines were evaluated based on alignment yield, basewise accuracy, mismatch and gap placement, exon junction discovery and suitability of alignments for transcript reconstruction. These results will be published in a forthcoming paper and additional data will be posted here.

### Source input data

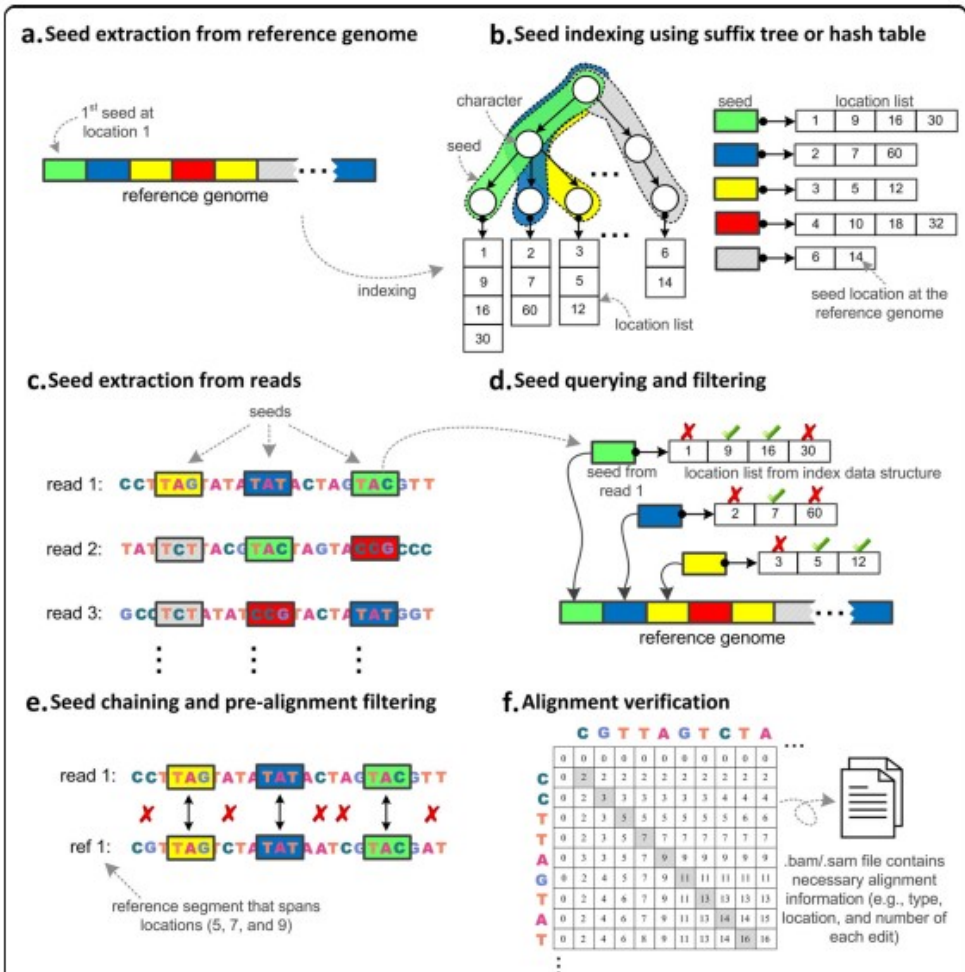
1. Mouse whole brain RNA-seq data (David Adams lab, WTSI/UK)  
Paired-end Illumina 76bp reads, insert sizes 175-225 bp
2. K562 cell line (human chronic myelogenous leukemia) RNASeq data (Tom Gingeras lab,

Project
Phase 2 GENCODE Goals
Data
Statistics - Human
Statistics - Mouse
BioDalliance
Participants
Publications
lncRNA microarray
RGASP 1/2
RGASP 3
Blog
GENCODE workshops
Contact us

L'étape de mapping n'est pas toujours simple.....



**Technology dictates algorithms:**  
*recent developments in read alignment*  
 Alser et al. *Genome Biology* (2021)  
 22:249



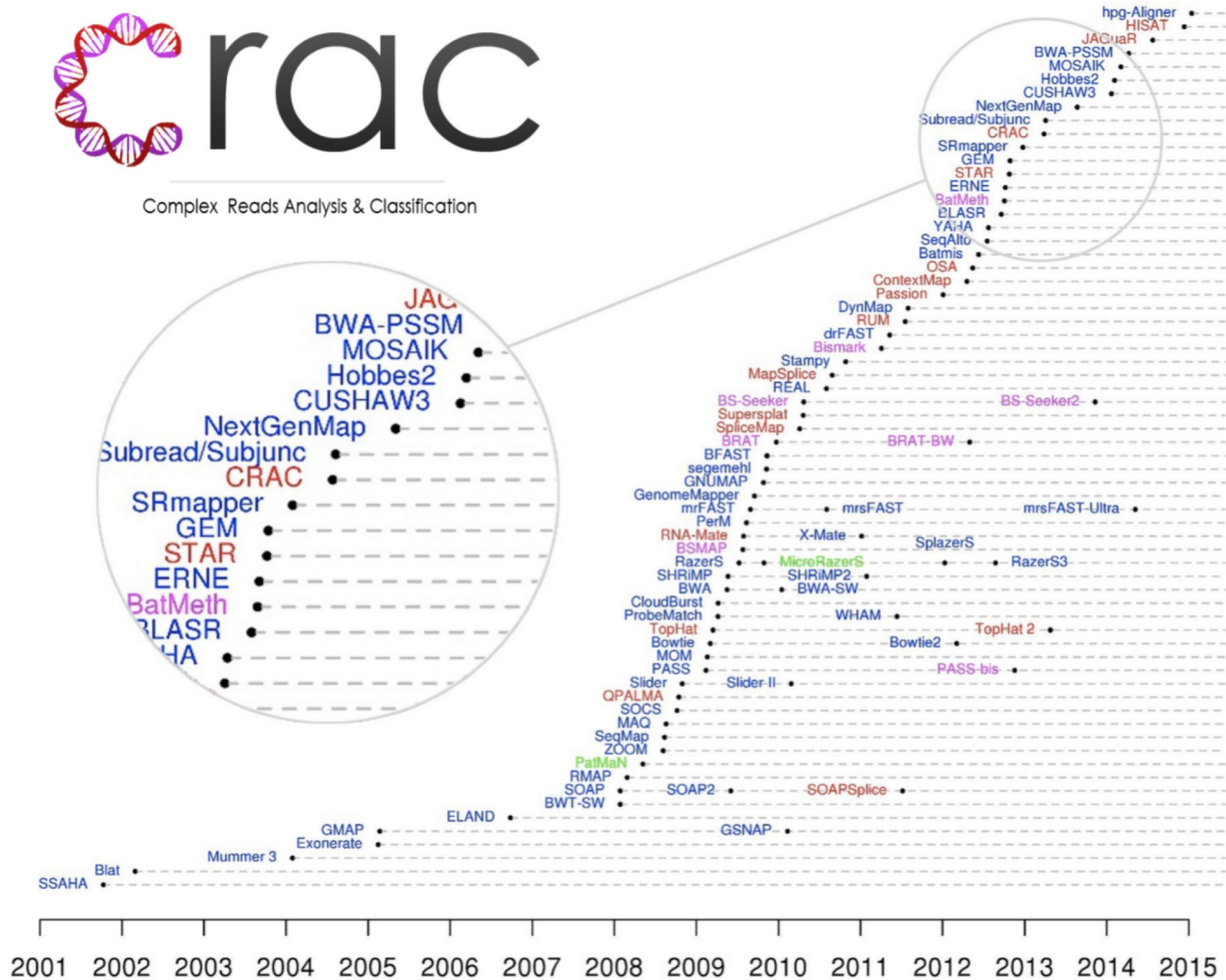
**Fig. 1** Overview of a read alignment algorithm. **a** The seeds from the reference genome sequence are extracted. **b** Each extracted seed and all its occurrence locations in the reference genome are stored using the data structure of choice (suffix tree and hash table are presented as an example). Common prefixes of the seeds are stored once in the branches of the suffix tree, while the hash table stores each seed individually. **c** The seeds from each read sequence are extracted. **d** The occurrences of each extracted seed in the reference genome are determined by querying the index database. In this example, the three seeds from the first read appear adjacent at locations 5, 7, and 9 in the reference genome. Two of the same seeds appear also adjacent at another two locations (12 and 16). Other non-adjacent locations are filtered out (marked with X) as they may not span a good match with the first read. **e** The adjacent seeds are linked together to form a longer chain of seeds by examining the mismatches between the gaps. Pre-alignment filters can also be applied to quickly decide whether or not the computationally expensive DP calculation is needed. **f** Once the pre-alignment filter accepts the alignment between a read and a region in the reference genome, then DP-based (or non-DP-based) verification algorithms are used to generate the alignment file (in BAM or SAM formats), which contains alignment information such as the exact number of differences, location of each difference, and their type.

**problème = # logiciels × # paramètres × # d'applications**



# racc

Complex Reads Analysis & Classification

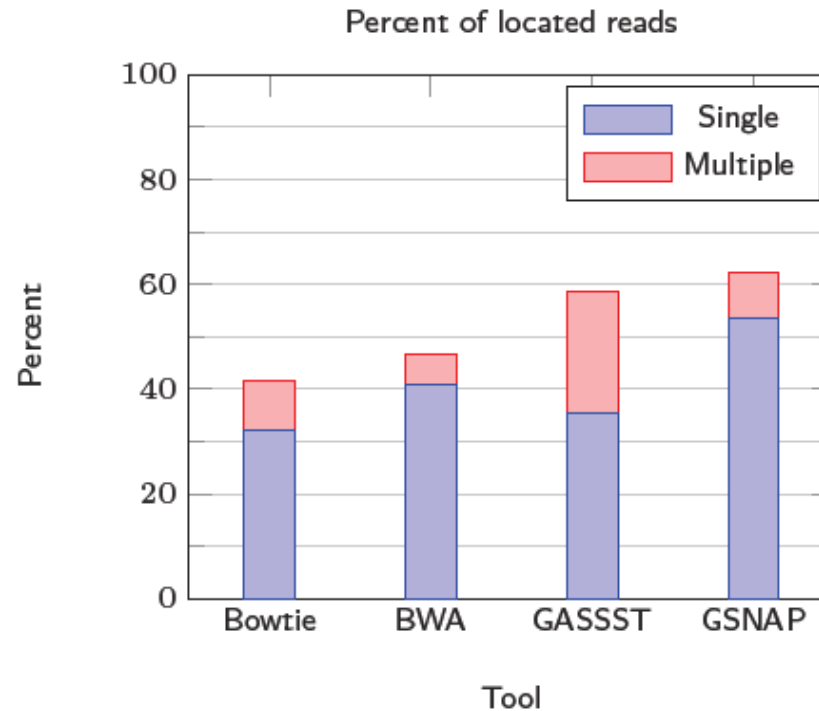


Nuno A. Fonseca, Johan Rung, Alvis Brazma, John C. Marioni, **Tools for mapping high-throughput sequencing data**. *Bioinformatics* (2012) 28 (24): 3169-3177.

# Exemple de problèmes bioinformatiques: 1.le mapping

## An example

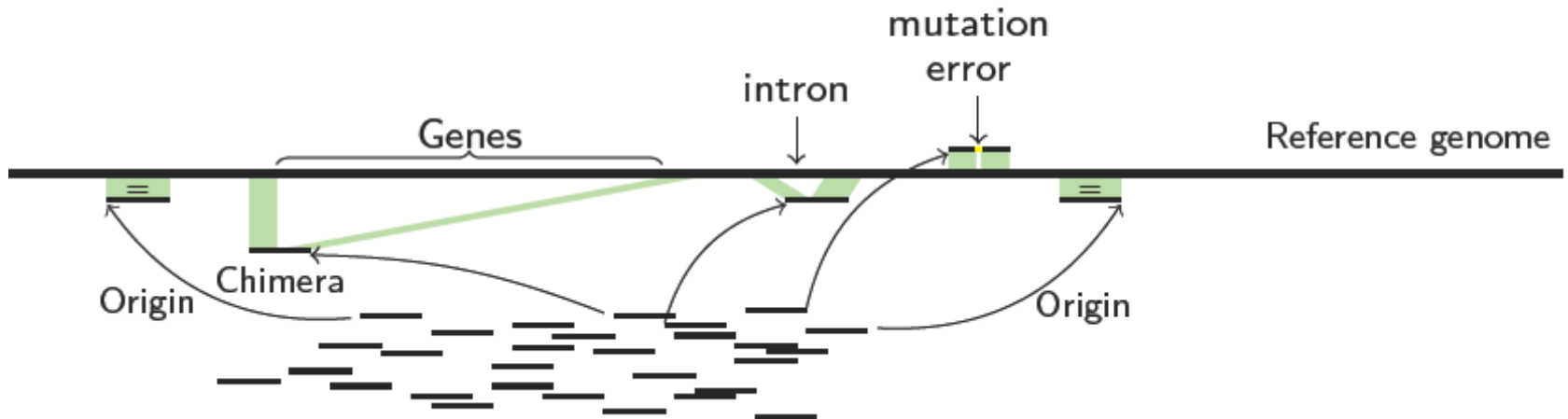
- RNA-Seq data from a K562 library (human cancer cells)
- ~ 70 millions of reads of 75 bp



# Inputs

RNA-seq: collection de millions de reads (~100 bp)

Reference genome:  $3 \times 10^9$  pb (genome humain)



« *Etapas de Mapping complexes* »



---

# CRAC: an integrated approach for RNA-seq analysis

For each read  
k factors and their  
genomic mapping

*CRAC: an approach to classify reads*

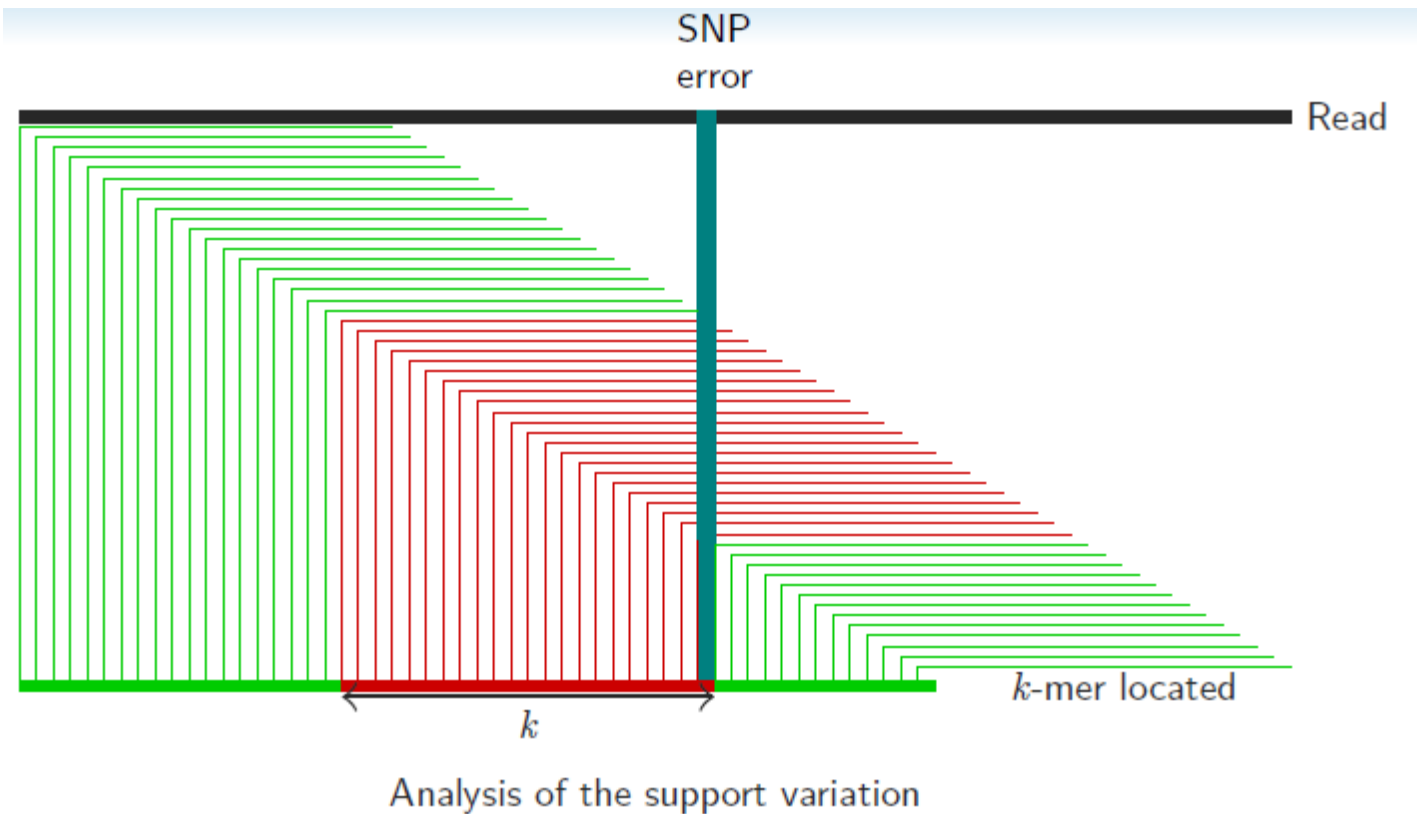
Mapping information  
used for SNP  
Error and splice...



*Based on the Gk arrays structure to index large sets of reads*

*Based on a k-mer approach*

*Length of k ~22 nt for th human genome*

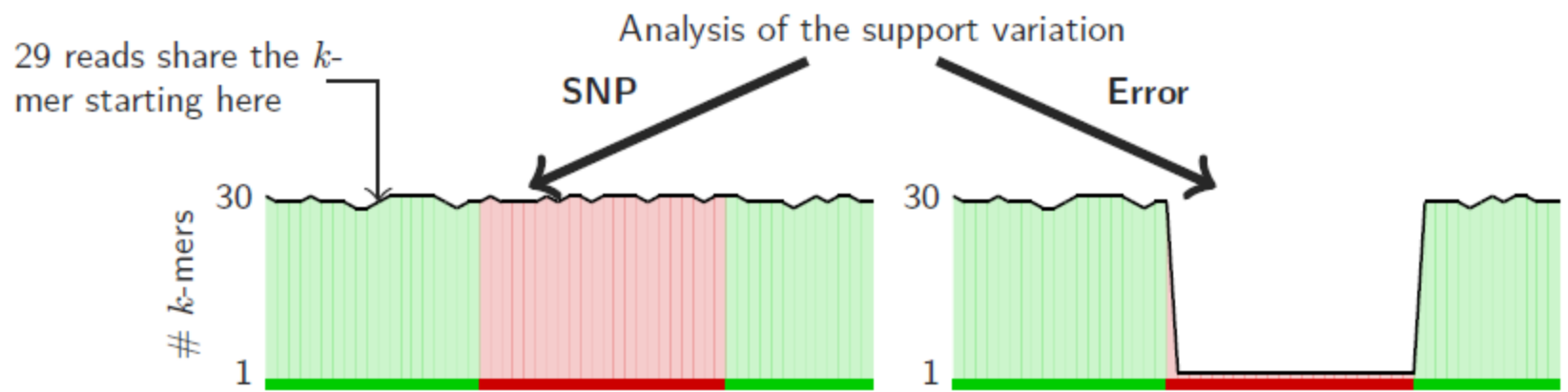
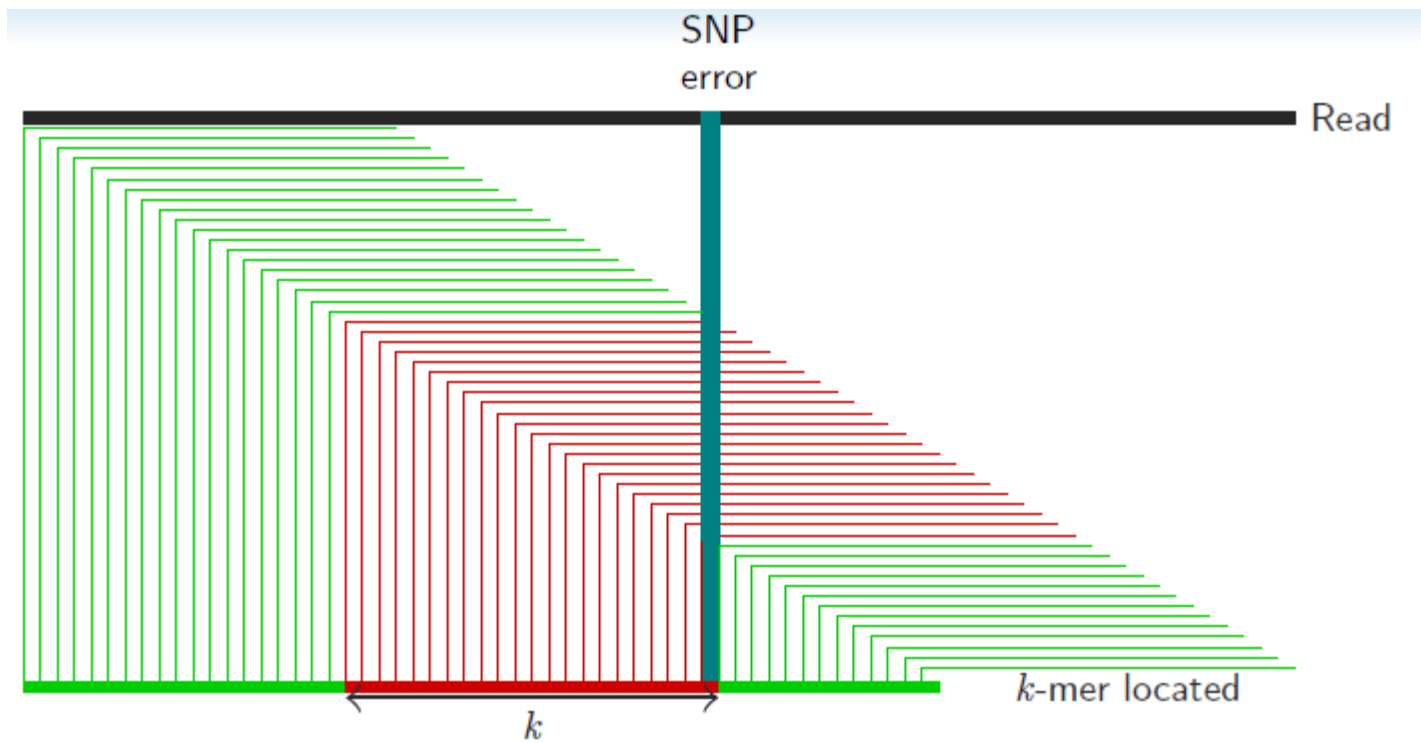


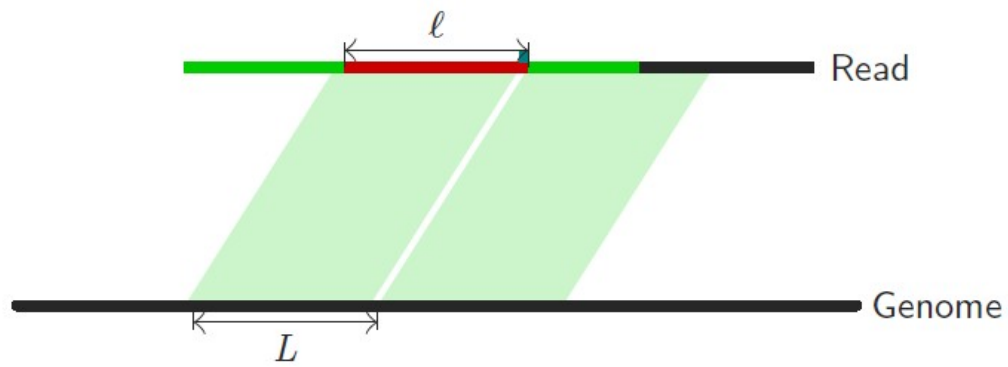
**break**

An interruption of the  $k$ -mers location on the genome

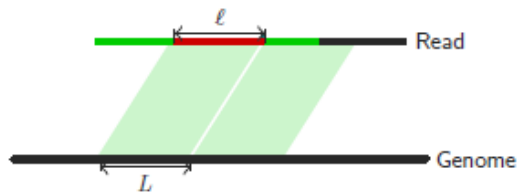
**support**

A  $k$ -mer support is the number of reads that contain it (local coverage)

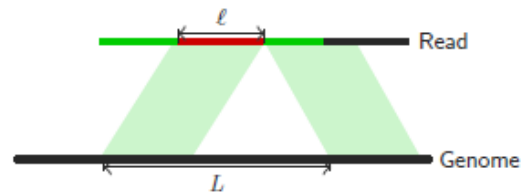




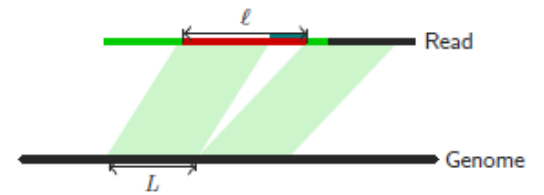
Substitution



Insertion

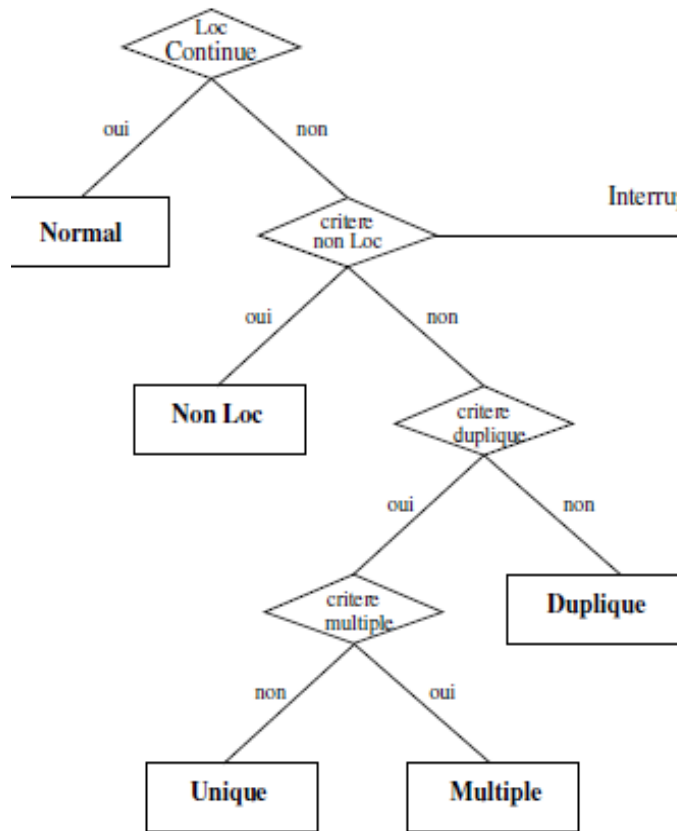


Deletion

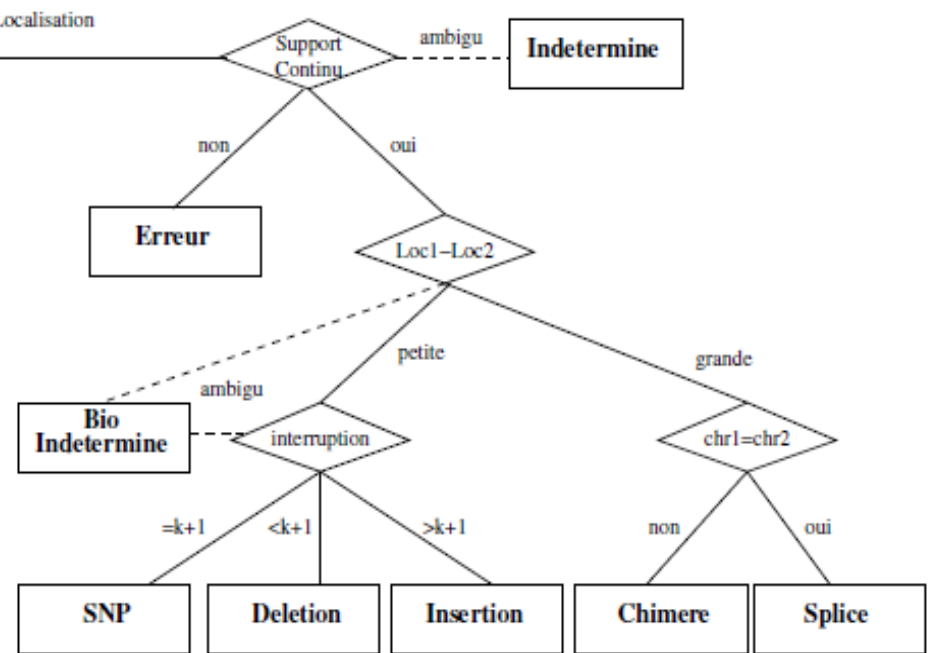


Splicing junction

## Localisation

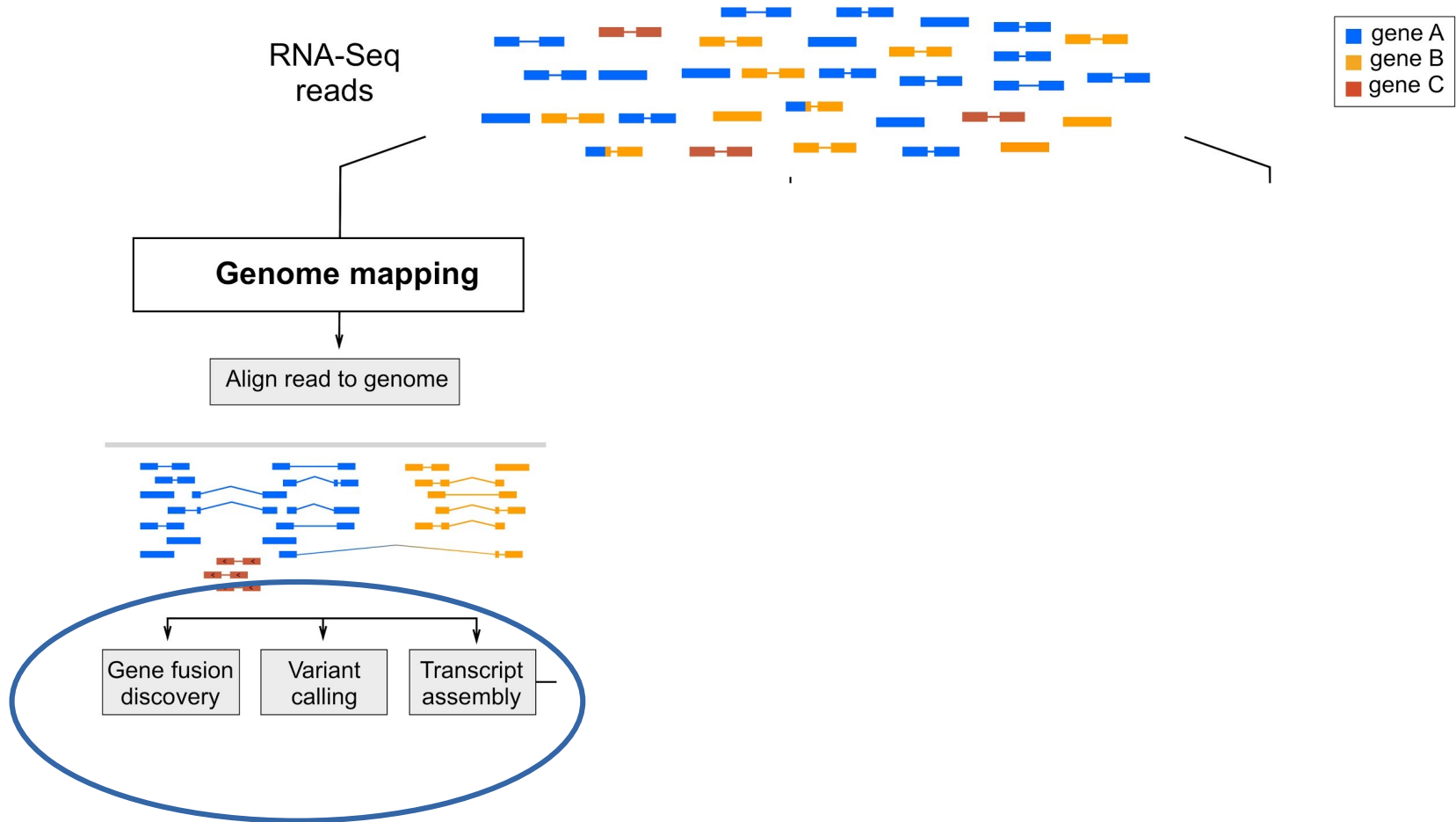


## Classification



Interruption de Localisation

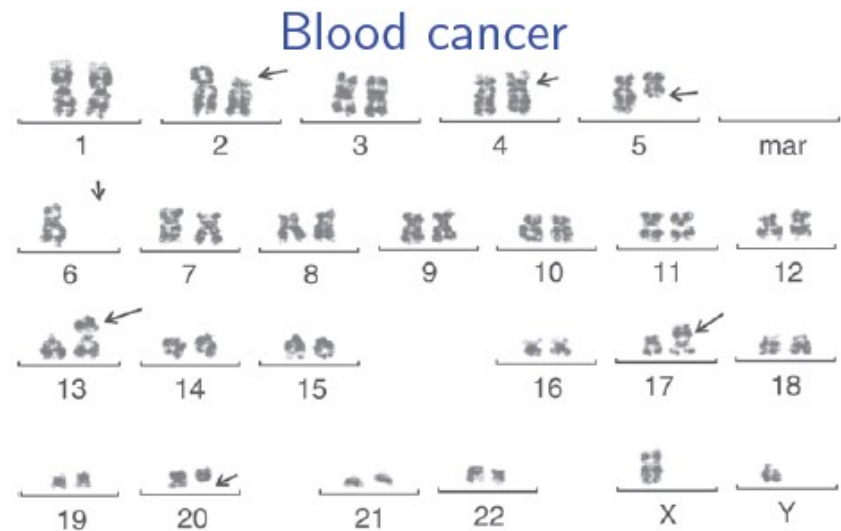
# RNA-seq and genome mapping



# Trouver des événements rares

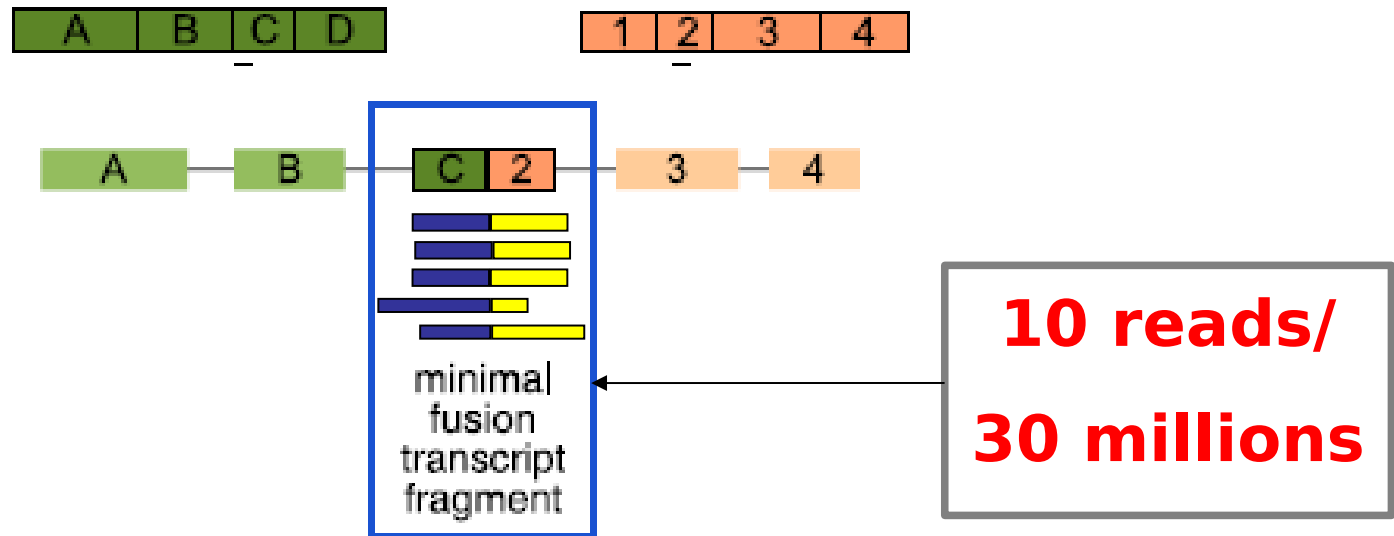
## Abnormal chromosome pool in cancer

- Diagnosis of chronic myelogenous leukemia (CML)
- Prognosis in myelodysplastic syndrome



# Trouver des événements rares (ARNde fusion)

1/ trouver des séquences rares et spécifiques dans une collection de read ( $10 / 30 \cdot 10^6$ )



2/ Distinguer les différents événements biologiques et les différencier des artefacts



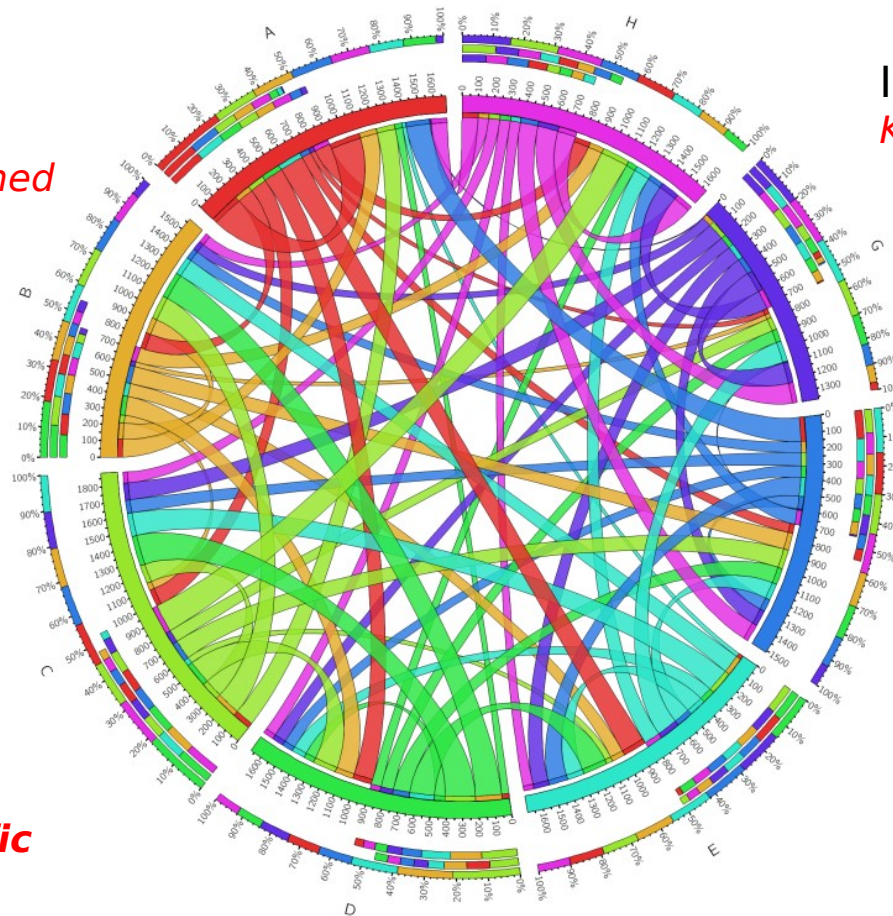
# The **chimeric** transcriptome in AML

Translocation  
*New translocation*  
*Variant of well-defined translocation*

Inversion  
*Known Inversion*

Tandem repeat  
*Inverted & direct*

Read-through  
*Tumor and non-tumor specific*



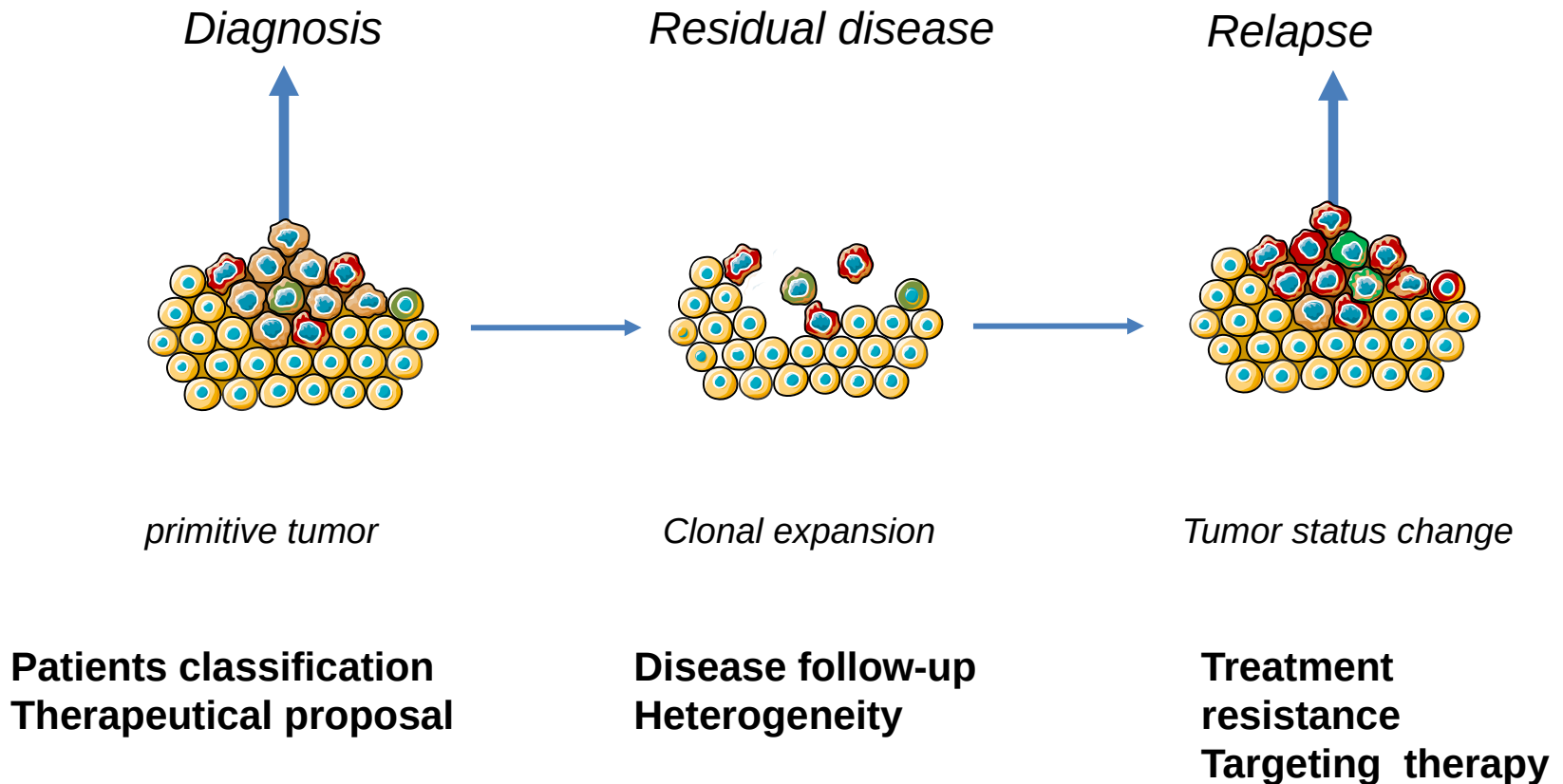
*New classes 2,3, 4*  
 Trans-splicing ?  
 CNV ?  
 Other

**New Chimeric RNA in all classes**

17 biological Validations

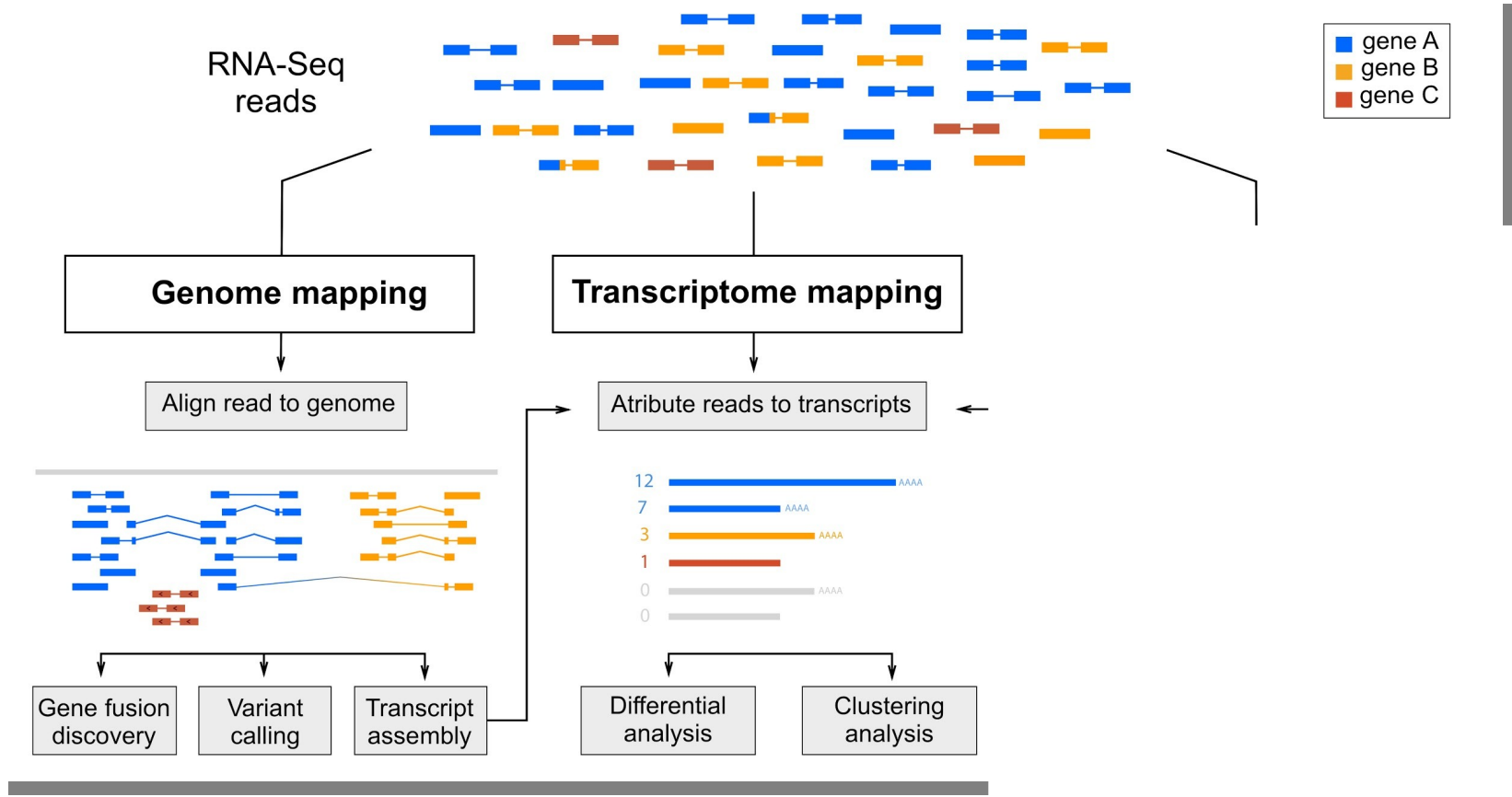
# Biomarqueurs et médecine de précision

*Diagnostic, suivi du patient en cancerologie*

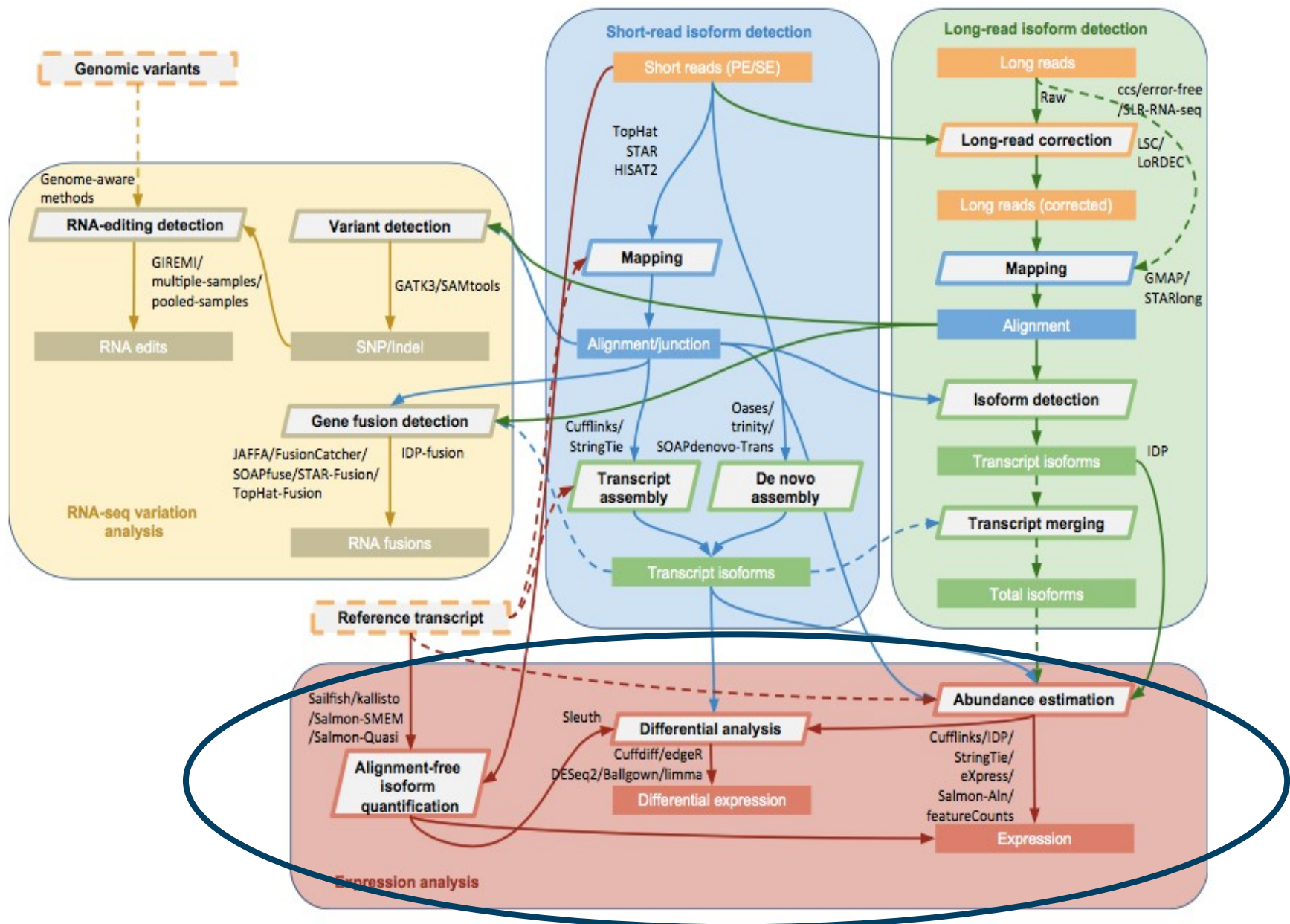


**La solution : une combinaison de Biomarqueurs**  
(Mutations, Fusion genes, Gene expression, Splicing events, lncRNA, miRNA)

# Transcriptome and RNA-seq



Transcriptome mapping for DGE



Differentiel gene expression (DGE) ...Application la plus fréquente...

# NGS et Transcriptomics

- Introduction: l'ère NGS et état des lieux, principaux séquenceurs (vu en intro)
- Les enjeux de l'analyse, principales applications en diagnostic (ADN/ARN)
- Transcriptome et RNAseq, les questions biologiques
  - ✓ Les différentes stratégies d'analyse
  - ✓ Le mapping, exemple et limites
  - ✓ **Annotations génomiques, transcriptome de reference**
  - ✓ Comment passer à l'échelle

## Statistics about the current Human GENCODE Release (version 25)

\* The statistics derive from the [gtf file](#) <sup>Ⓜ</sup> that contains only the annotation of the main chromosomes.

For details about the calculation of these statistics please see the [README\\_stats.txt](#) <sup>Ⓜ</sup> file.



[Compare with the previous release \(GENCODE 24\)](#) »

### Version 25 (March 2016 freeze, GRCh38) - Ensembl 86

#### General stats

Total No of Genes	58037	Total No of Transcripts	198093
Protein-coding genes	19950	Protein-coding transcripts	80087
Long non-coding RNA genes	15767	- full length protein-coding:	54755
Small non-coding RNA genes	7258	- partial length protein-coding:	25332
Pseudogenes	14650	Nonsense mediated decay transcripts	13769
- processed pseudogenes:	10725	Long non-coding RNA loci transcripts	27692
- unprocessed pseudogenes:	3400		
- unitary pseudogenes:	214		
- polymorphic pseudogenes:	51		
- pseudogenes:	21	Total No of distinct translations	60033
Immunoglobulin/T-cell receptor gene segments		Genes that have more than one distinct translations	13536
- protein coding segments:	411		
- pseudogenes:	239		



## Human

### Statistics about the current GENCODE Release (version 38)

The statistics derive from the [gtf file](#) that contains only the annotation of the main chromosomes.

For details about the calculation of these statistics please see the [README\\_stats.txt file](#).

#### General stats

Total No of Genes	60649	Total No of Transcripts	237012
Protein-coding genes	19955	Protein-coding transcripts	86757
Long non-coding RNA genes	17944	- full length protein-coding	61015
Small non-coding RNA genes	7567	- partial length protein-coding	25742
Pseudogenes	14773	Nonsense mediated decay transcripts	18881
- processed pseudogenes	10667	Long non-coding RNA loci transcripts	48752
- unprocessed pseudogenes	3565		
- unitary pseudogenes	241		
- polymorphic pseudogenes	49		
- pseudogenes	15	Total No of distinct translations	63968
Immunoglobulin/T-cell receptor gene segments		Genes that have more than one distinct tra	13689
- protein coding segments	409		
- pseudogenes	236		