

# Introduction to functional genomics

Pascal Martin



[Pascal.Martin@inra.fr](mailto:Pascal.Martin@inra.fr)



[@PgpMartin](https://twitter.com/PgpMartin)

- ✓ **Introduction**

- ✓ Genes, RNAs & transcriptomes
- ✓ Transcription and the epigenome
- ✓ Why measuring gene expression ?

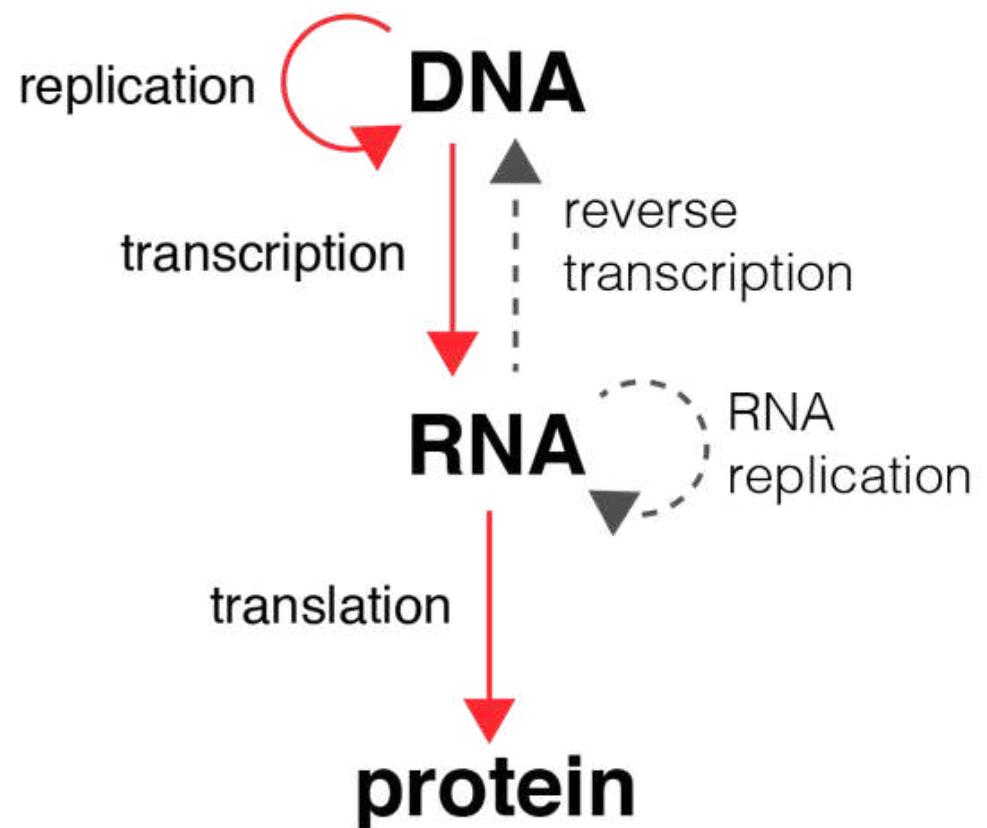
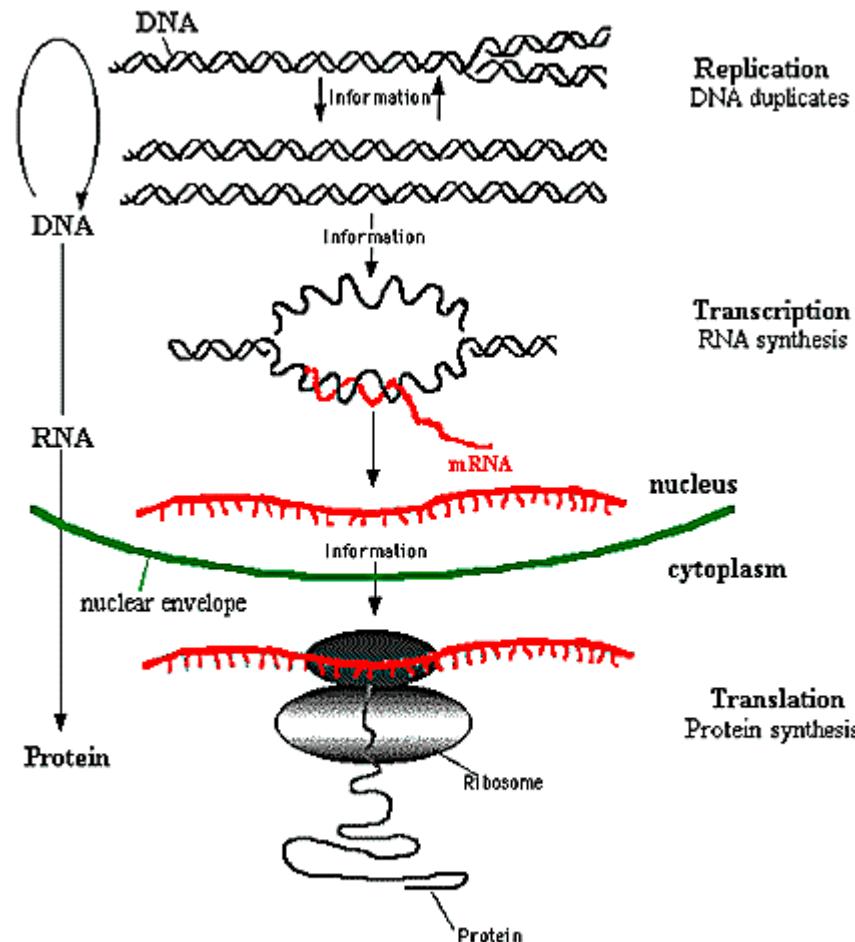
- ✓ **Transcriptomics**

- ✓ low throughput methods (Northen blot, RPA, qPCR)
- ✓ microarrays
- ✓ NGS / RNA-seq

- ✓ **Epigenomics**

- ✓ ChIP-seq
- ✓ Chromatin accessibility (MNase, DNase, FAIRE, ATAC)
- ✓ Chromatin conformation

# The Central Dogma of Molecular Biology



**The Central Dogma of Molecular Biology**

# What is a gene ?

Mendel, 1865  
Morgan, 1910  
Beadle & Tatum : "One gene, one enzyme",  
1941  
Avery, McLeod, McCarty, 1944  
Watson & Crick, Wilkins & Franklin, 1953  
Crick "Central dogma", 1958  
Jacob & Monod, lac operon, 1961  
Genetic code, 1965  
RNA splicing, 1967  
  
GENSCAN gene prediction, 1997  
Draft of the human genome, 2001  
ENCODE 1st phase, 2003-07

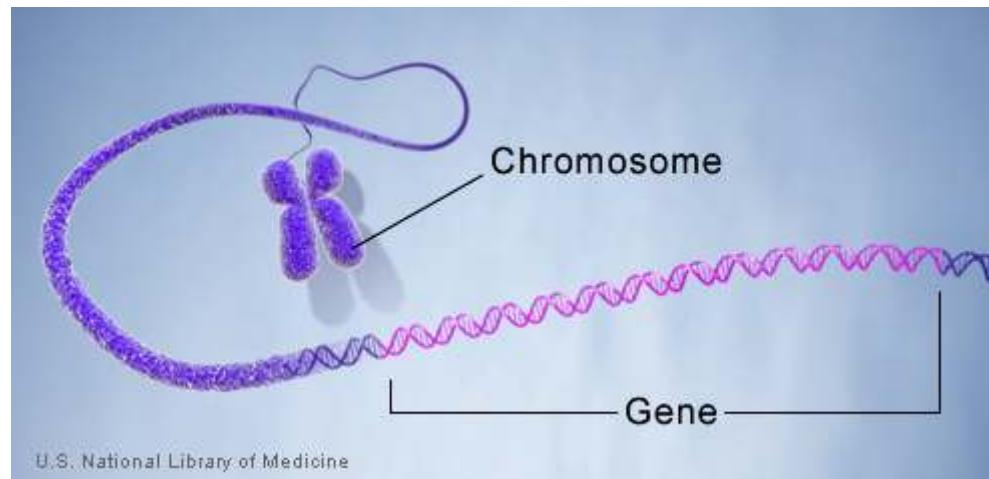
- Timeline
- Discrete Heredity Unit
  - Distinct Locus (linear organization, "beads on a string")
  - A protein blueprint
  - A physical molecule
  - A transcribed code
  - ORF sequence pattern
  - Annotated genomic entity

U	(RNA only)
C	T pyrimidines
G	A purines

## Post-ENCODE definition:

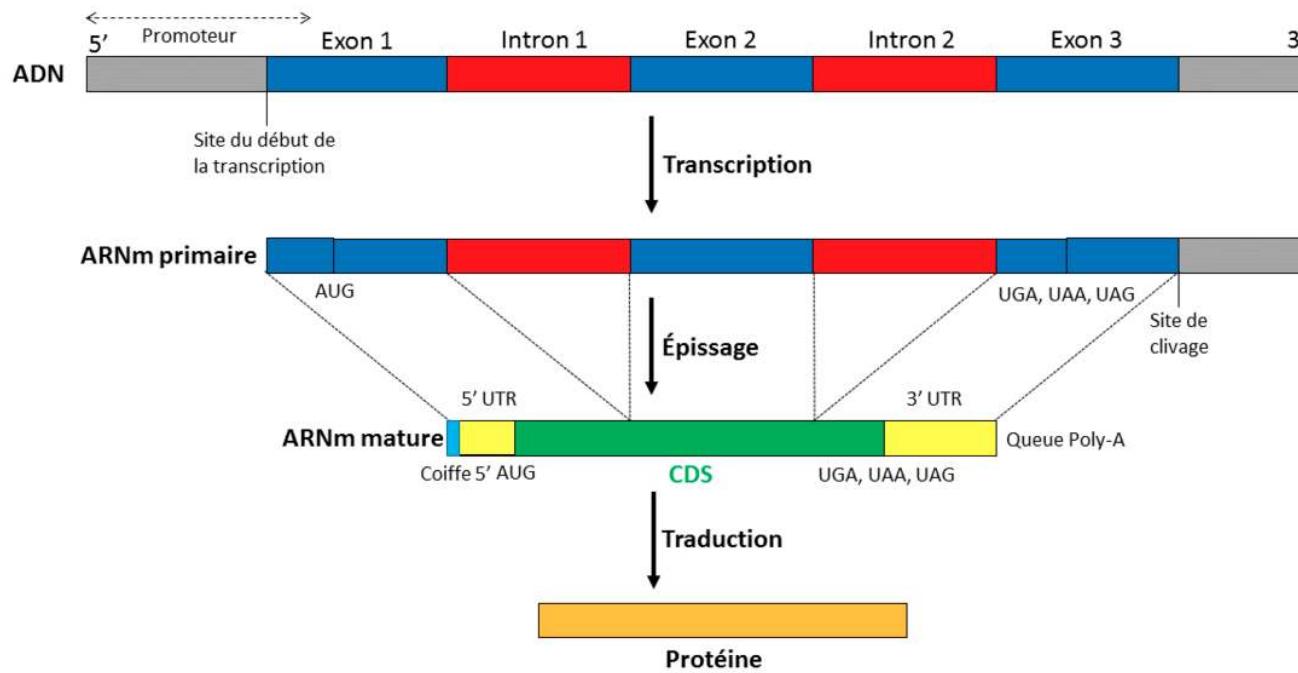
"A gene is a union of genomic sequences encoding a coherent set of potentially overlapping functional products"

Gerstein et al., Genome Res 2007

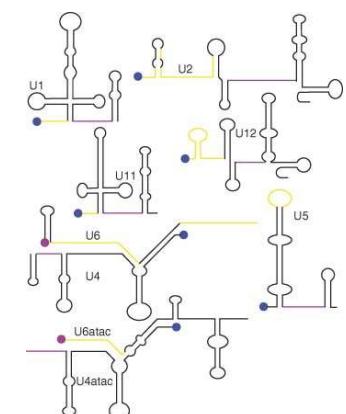
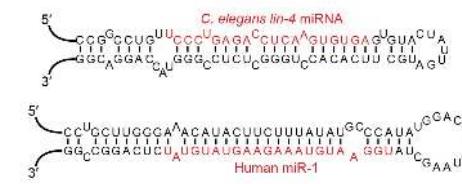
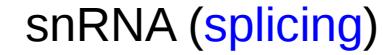
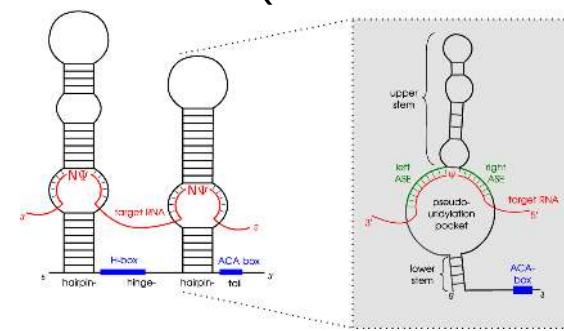
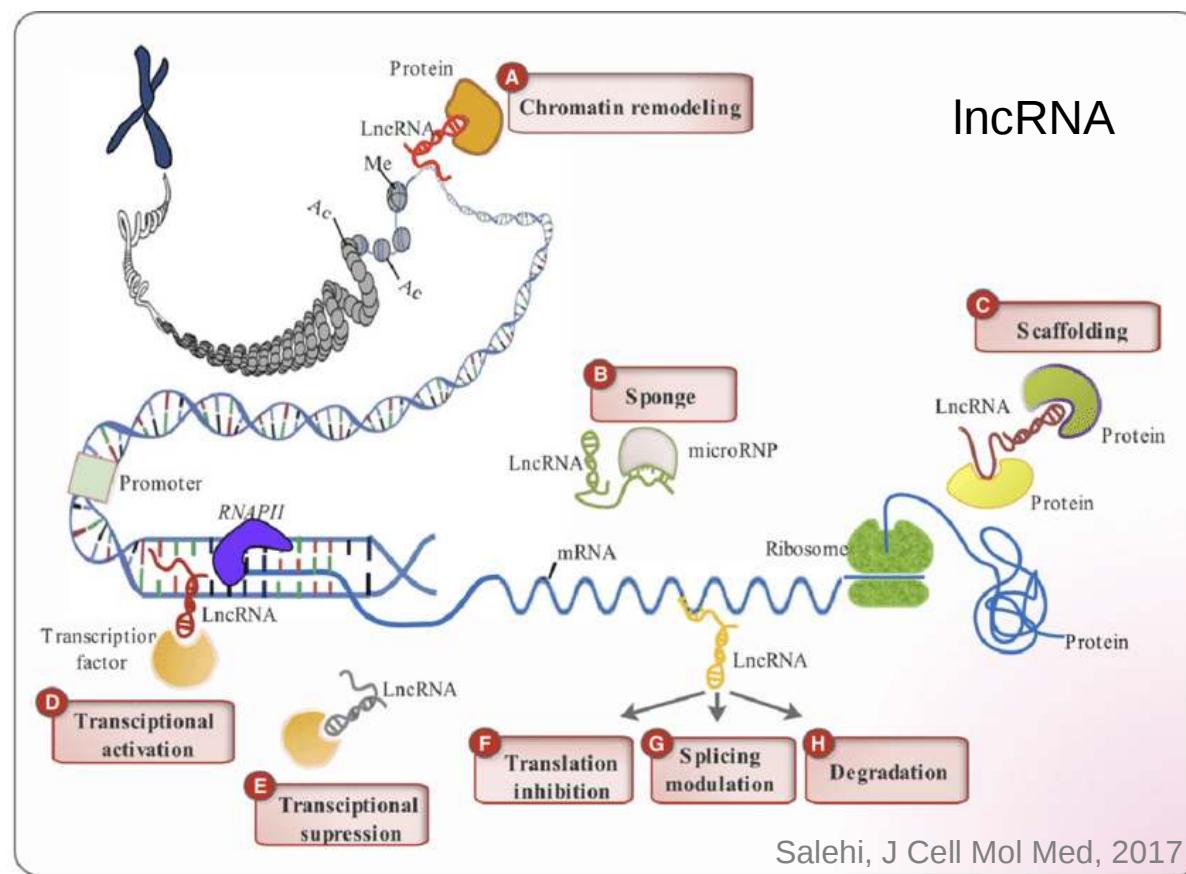
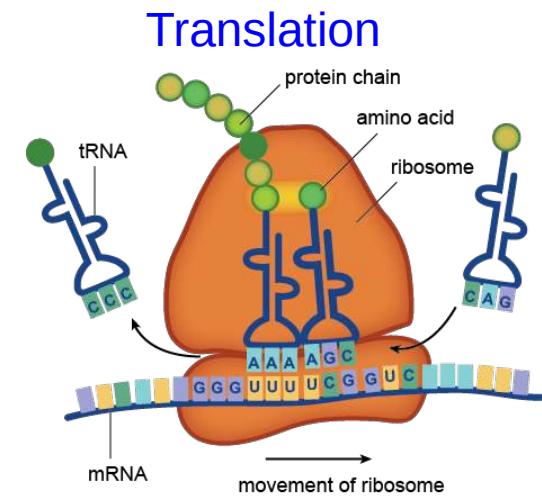
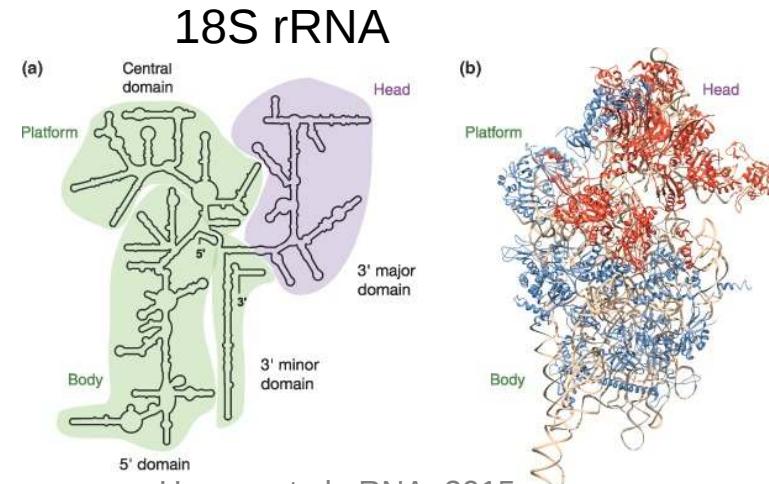
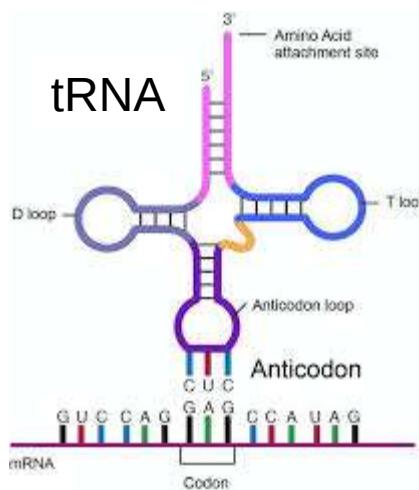


"A basic physical and functional unit of heredity"

# Many genes encode proteins

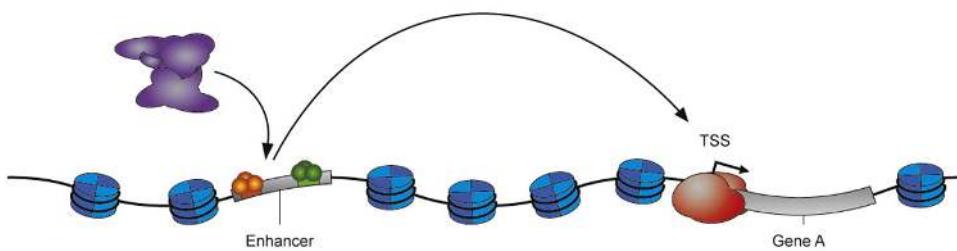


but many don't

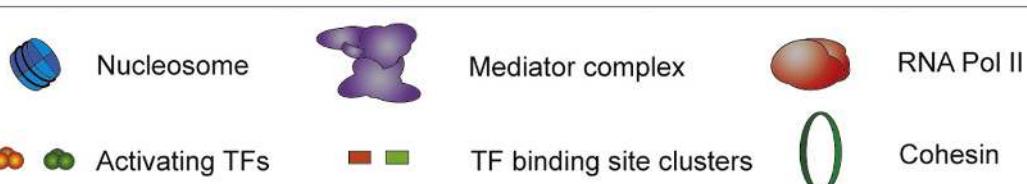
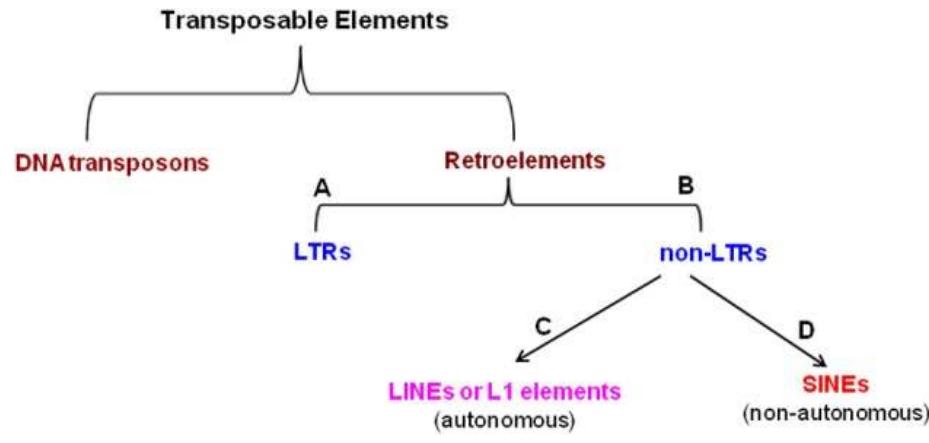
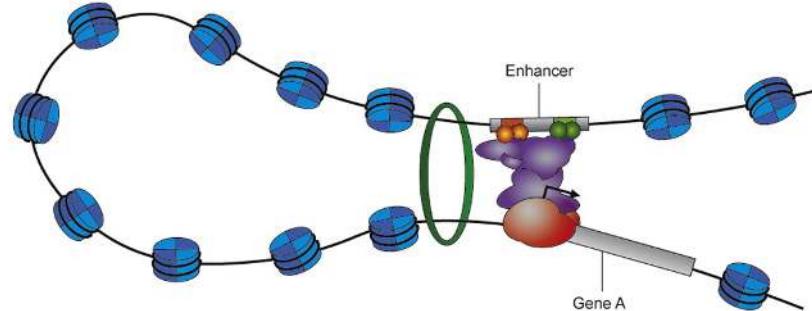
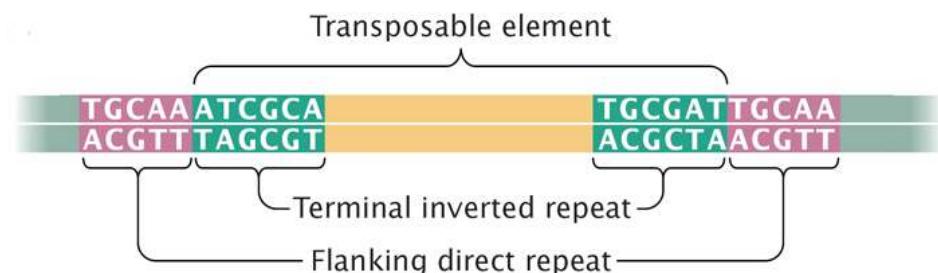


# Other transcribed regions

## Enhancers

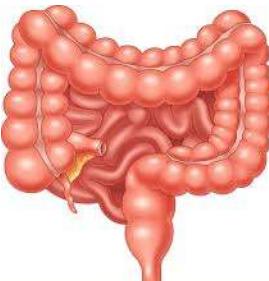
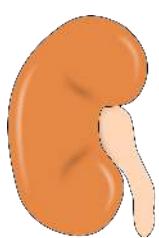
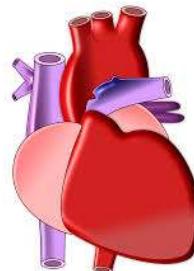
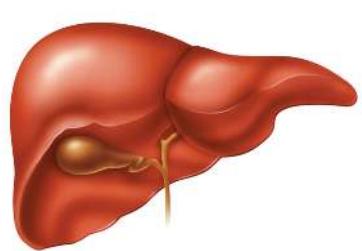


## Transposable elements / viruses

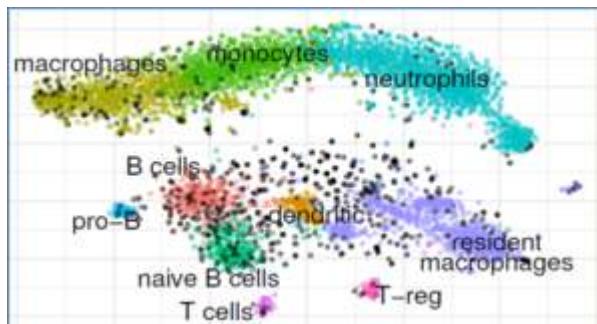


# Analyzing the transcriptomeS

Different tissues / conditions



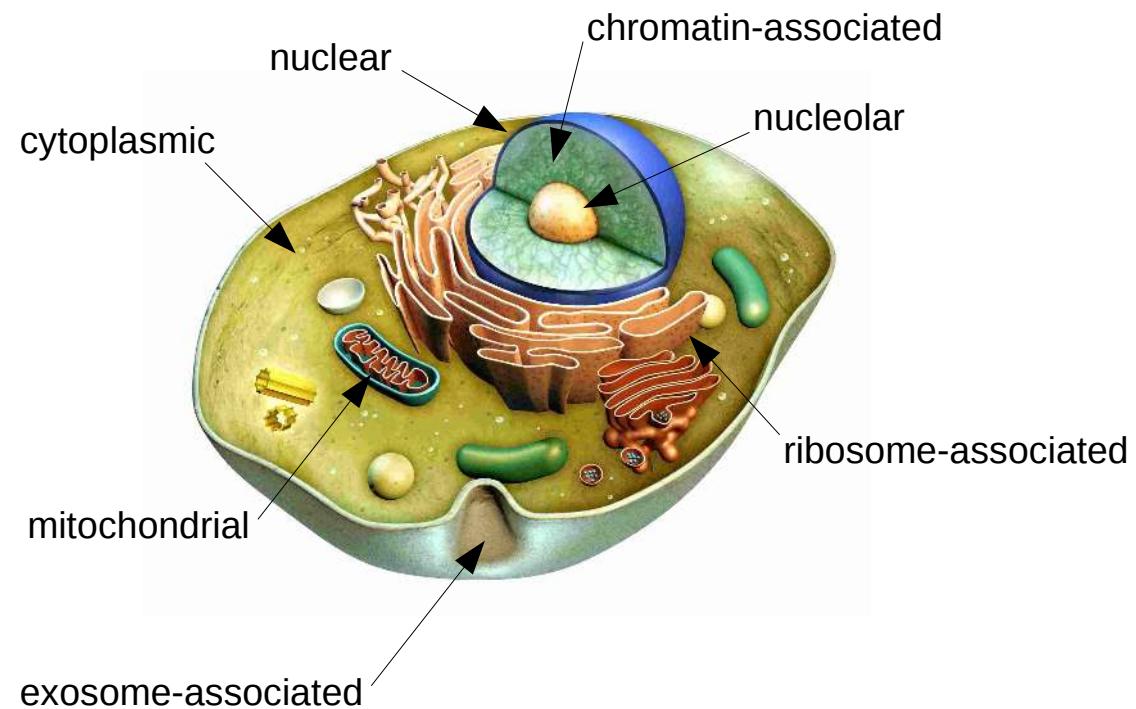
Single-cell transcriptomes



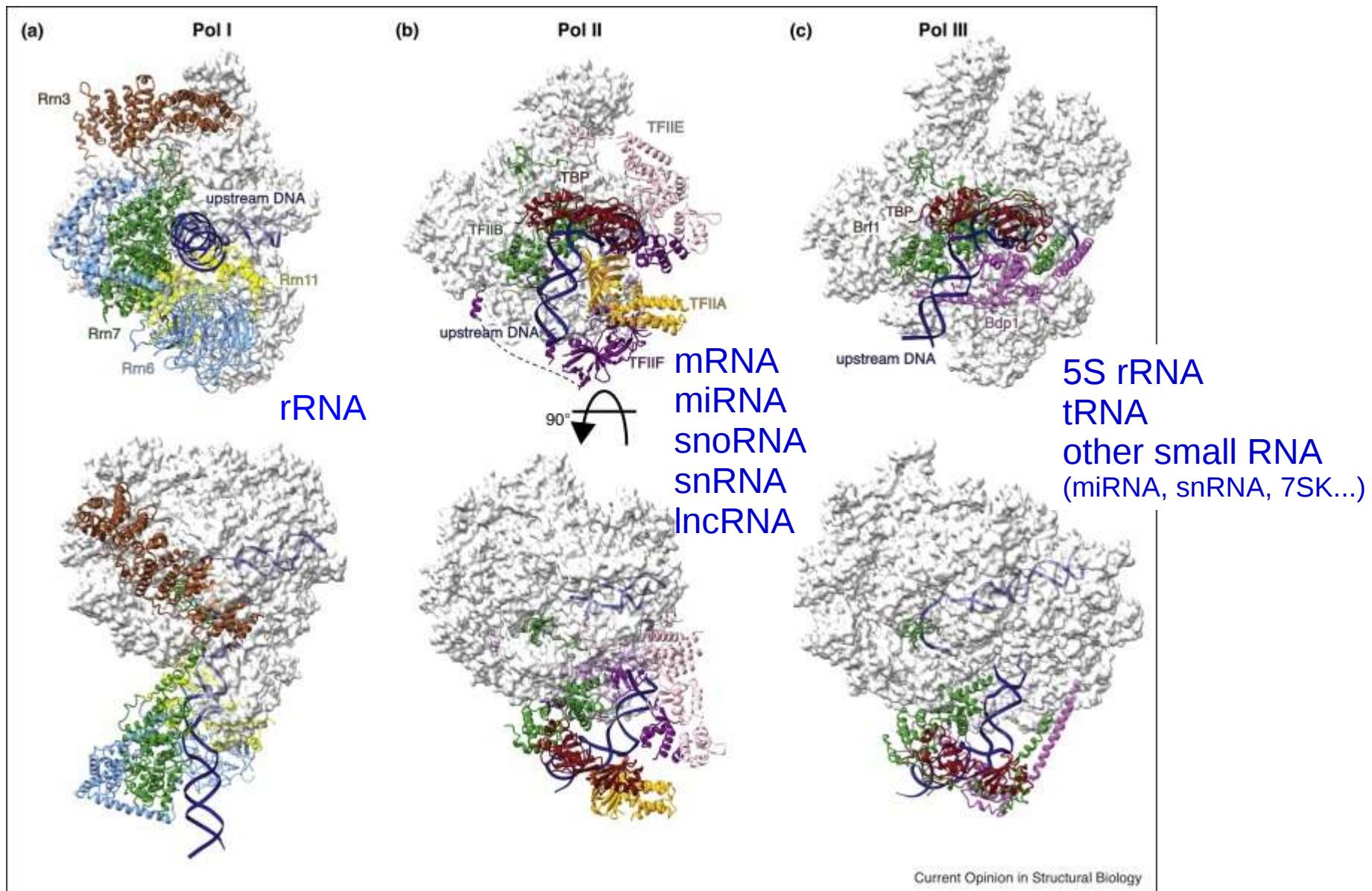
Different types of RNA

- PolyA+/-
- capped / uncapped
- short / long

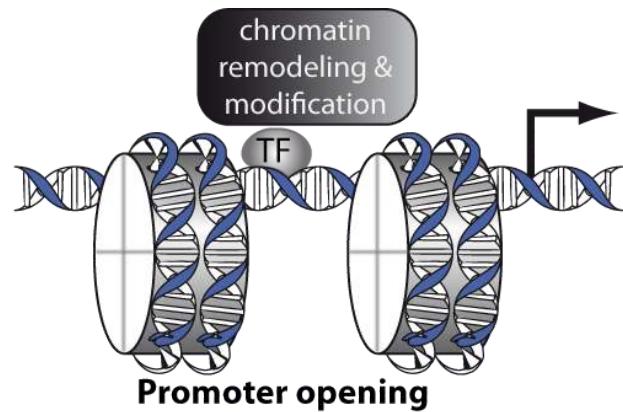
Subcellular fractions



# Transcribing DNA



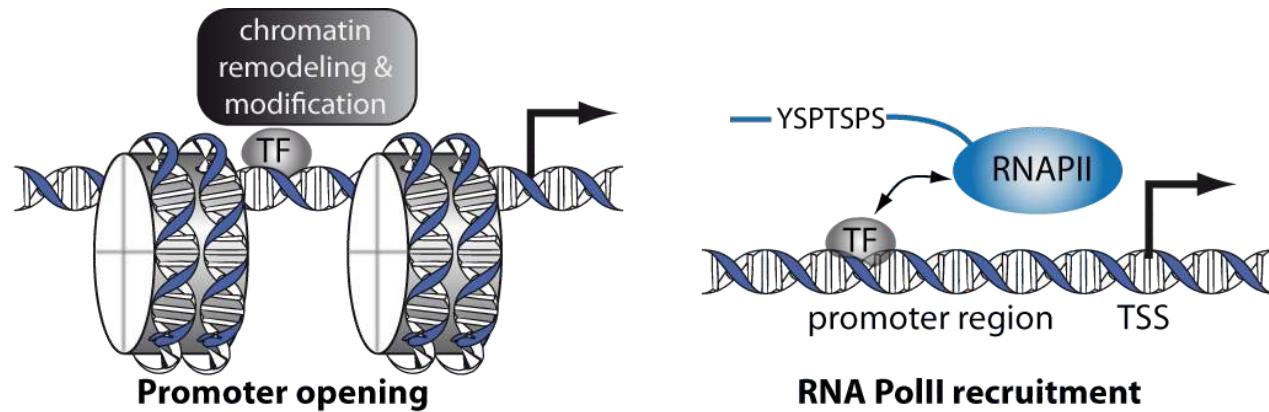
# RNA polymerase II transcription



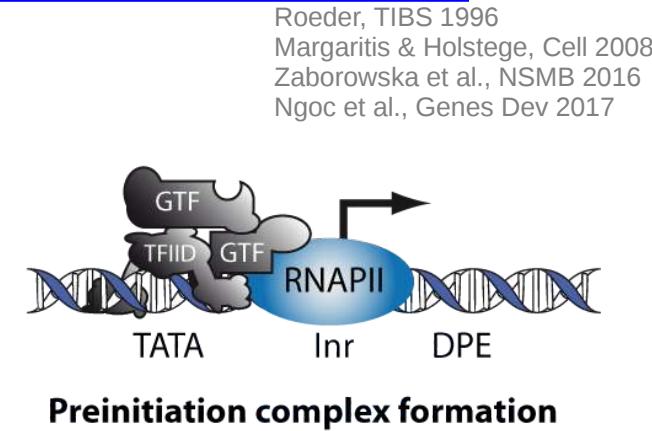
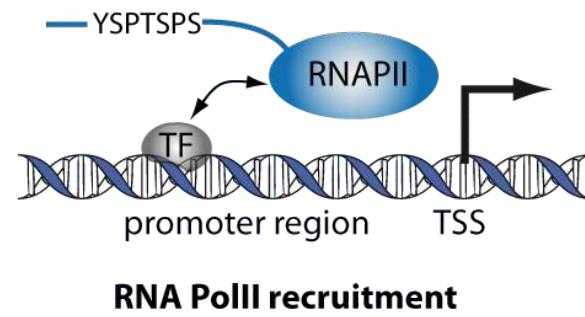
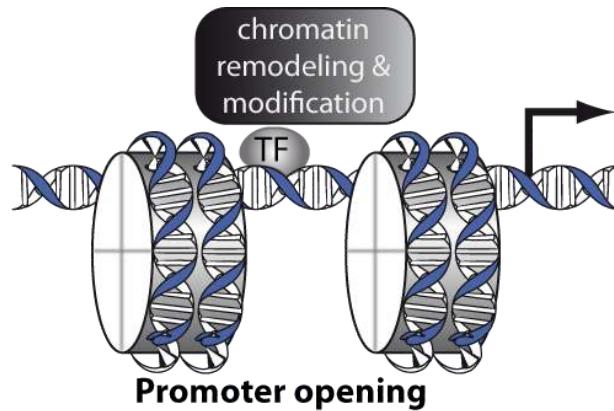
Roeder, TIBS 1996  
Margaritis & Holstege, Cell 2008

# RNA polymerase II transcription

Roeder, TIBS 1996  
Margaritis & Holstege, Cell 2008



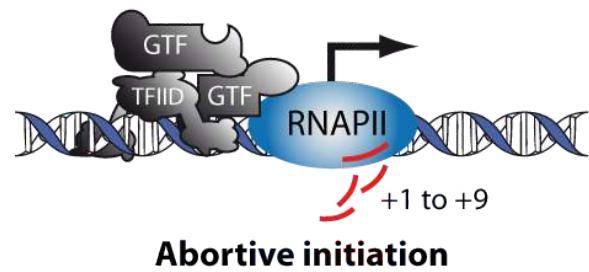
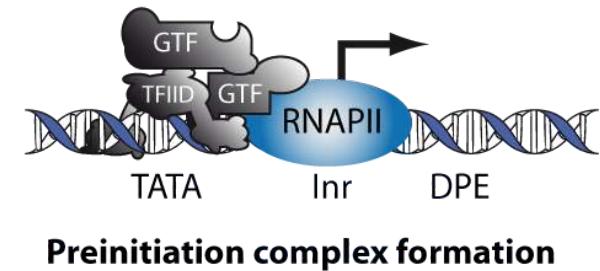
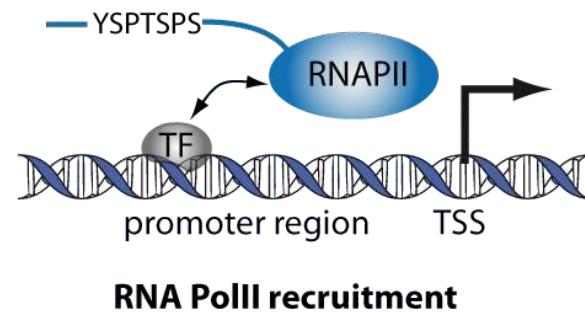
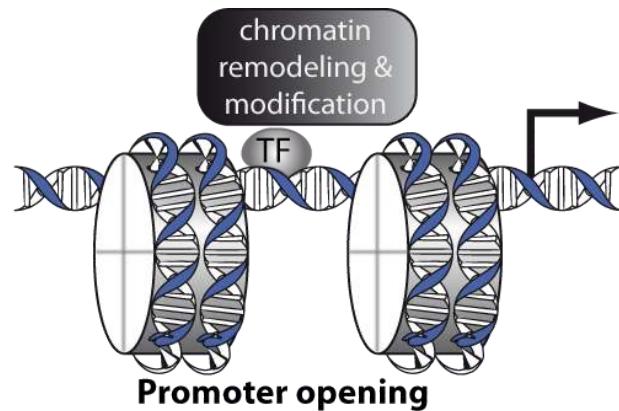
# RNA polymerase II transcription



Roeder, TIBS 1996  
Margaritis & Holstege, Cell 2008  
Zaborowska et al., NSMB 2016  
Ngoc et al., Genes Dev 2017

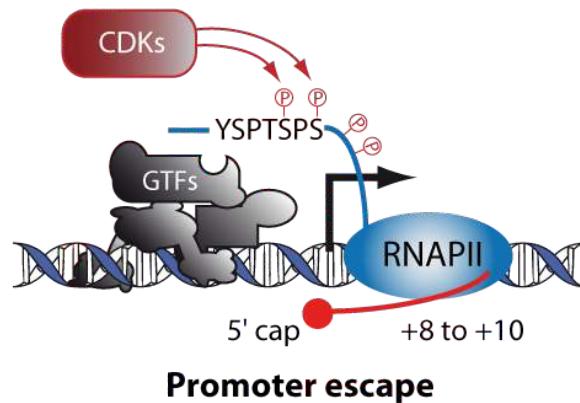
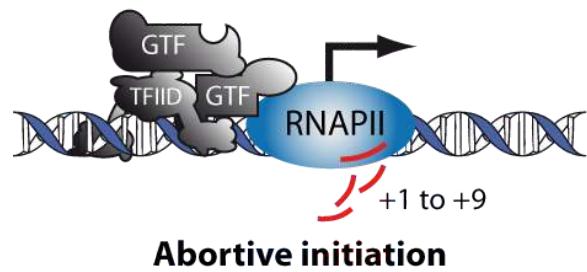
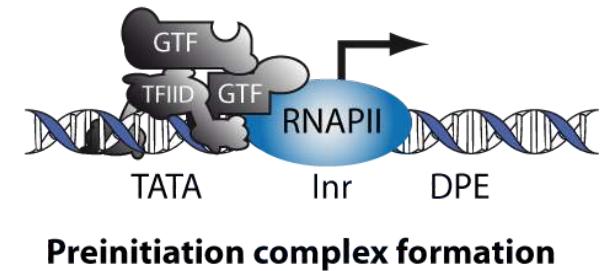
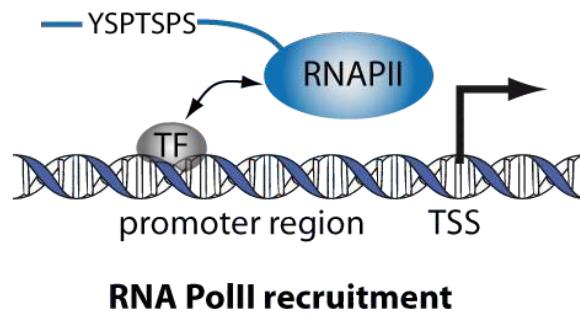
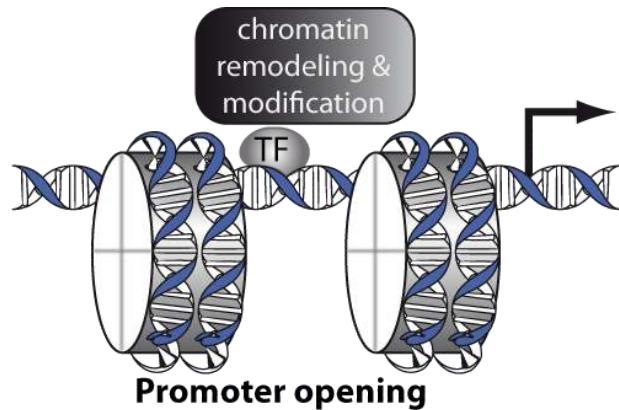
# RNA polymerase II transcription

Roeder, TIBS 1996  
Margaritis & Holstege, Cell 2008  
Zaborowska et al., NSMB 2016  
Ngoc et al., Genes Dev 2017



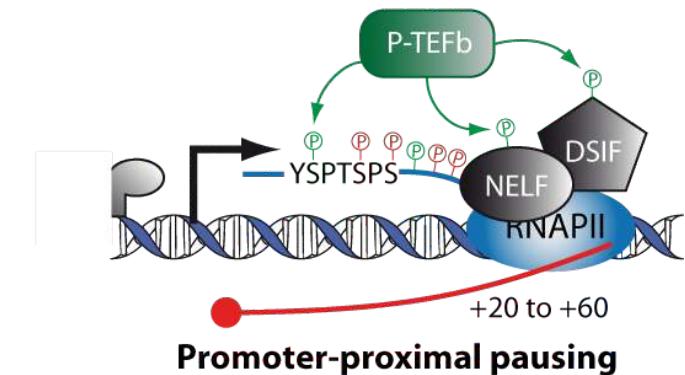
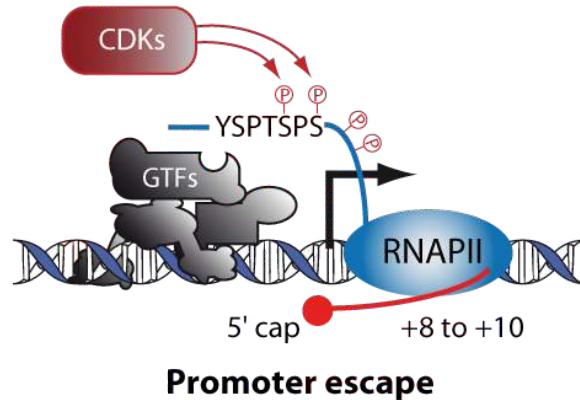
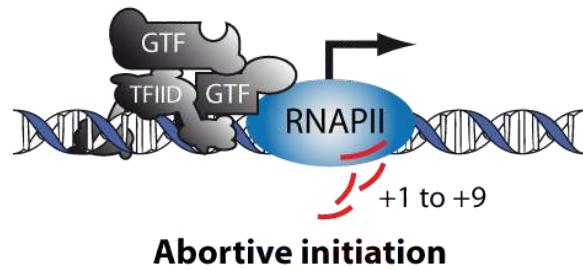
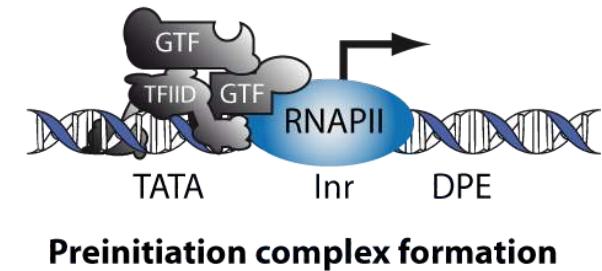
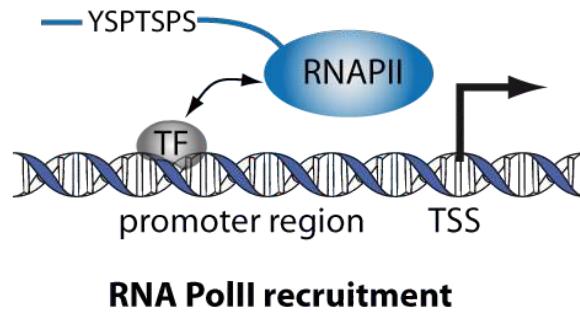
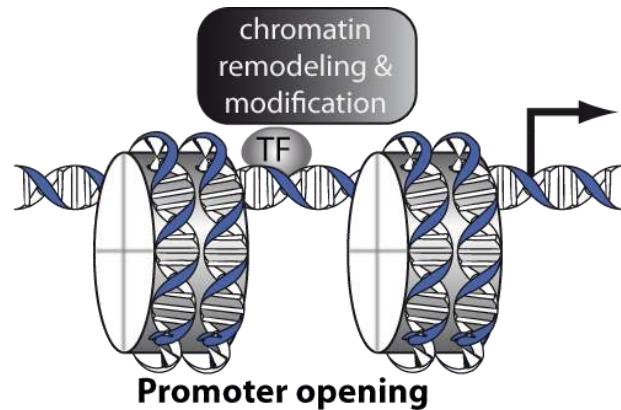
# RNA polymerase II transcription

Roeder, TIBS 1996  
Margaritis & Holstege, Cell 2008  
Zaborowska et al., NSMB 2016  
Ngoc et al., Genes Dev 2017



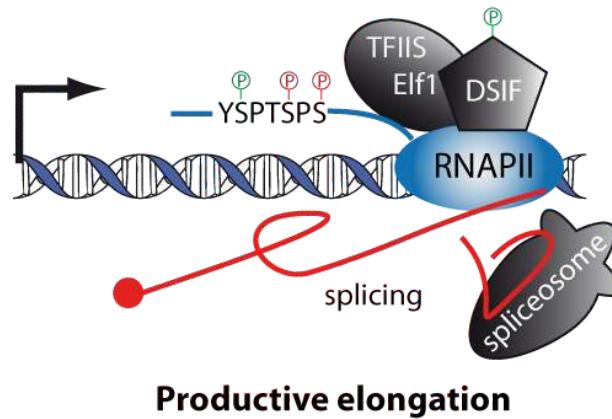
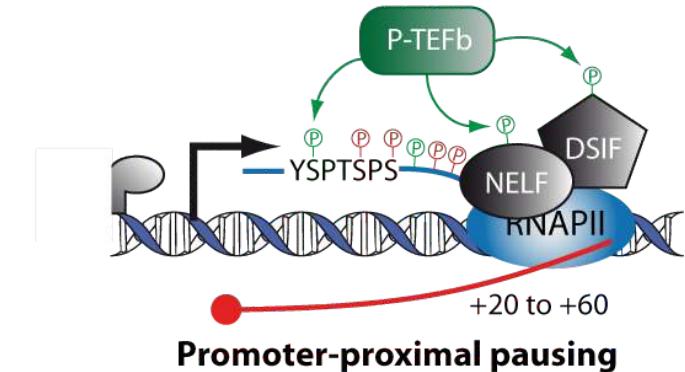
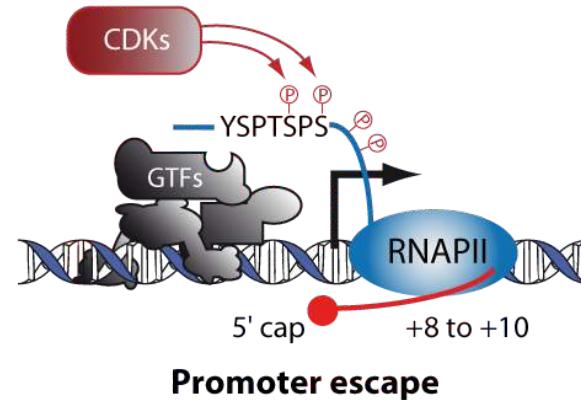
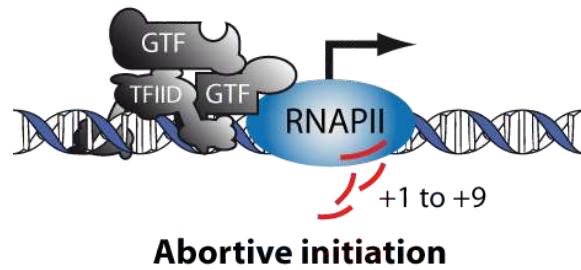
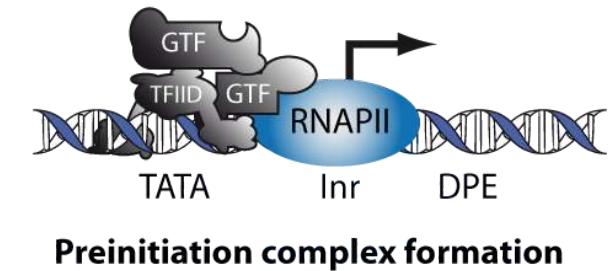
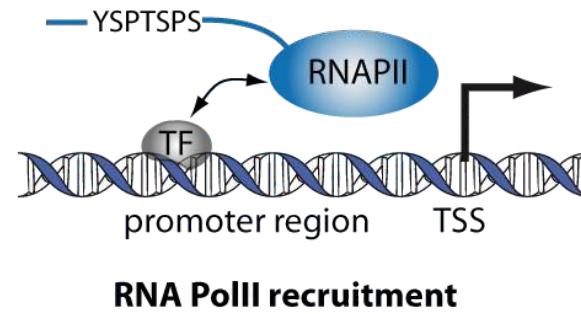
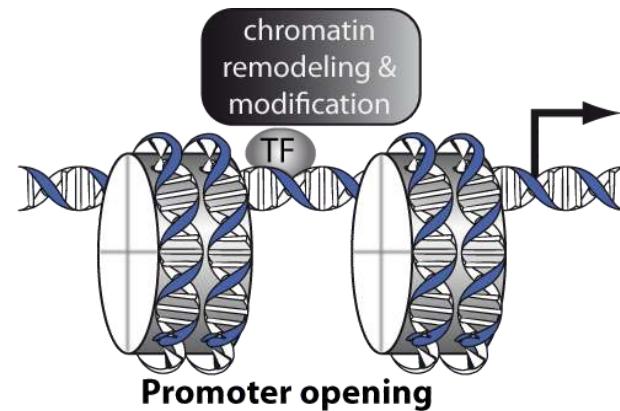
# RNA polymerase II transcription

Roeder, TIBS 1996  
Margaritis & Holstege, Cell 2008  
Zaborowska et al., NSMB 2016  
Ngoc et al., Genes Dev 2017



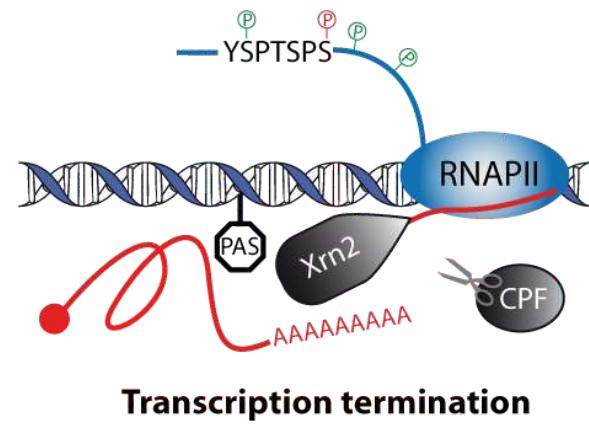
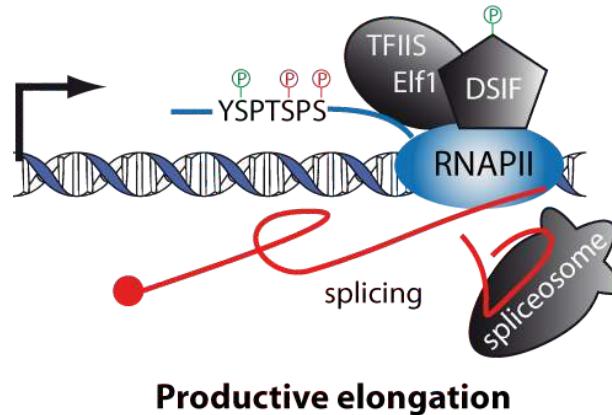
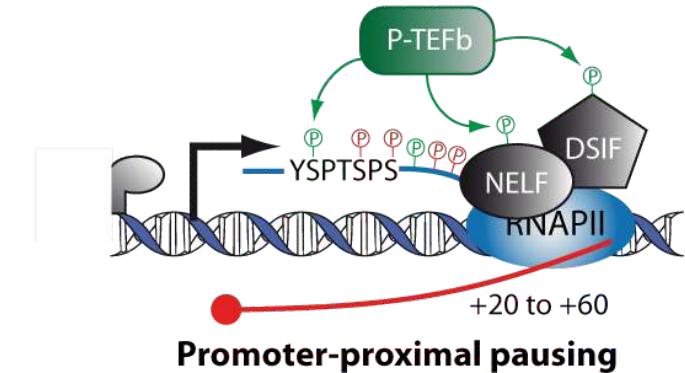
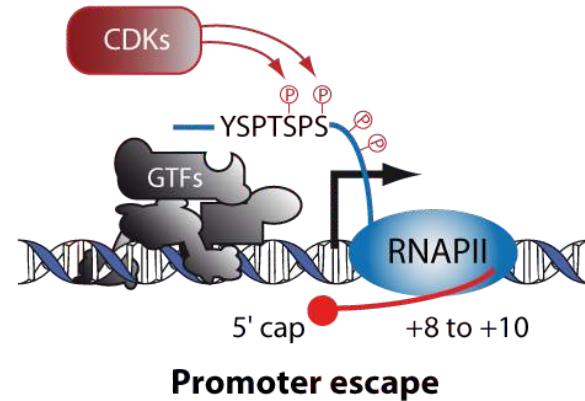
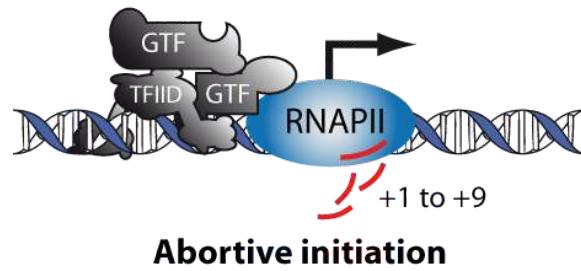
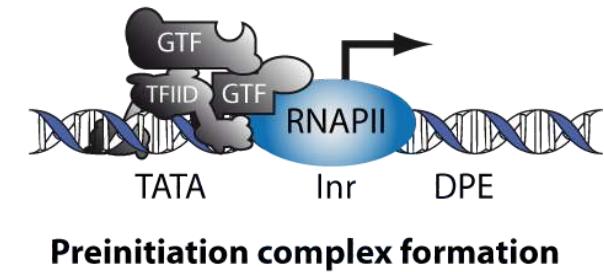
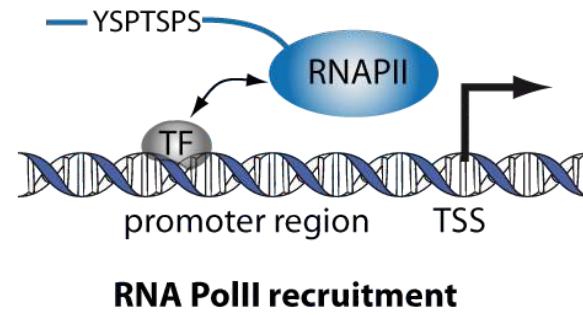
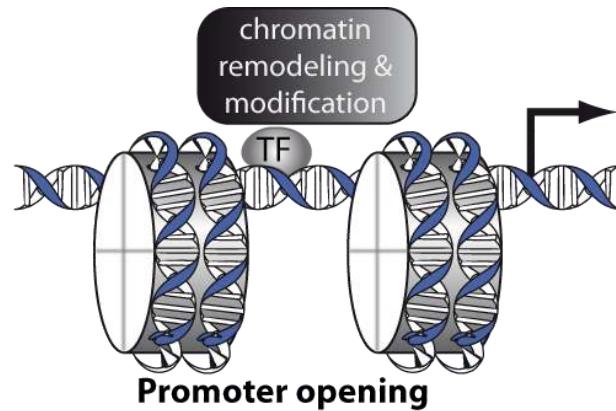
# RNA polymerase II transcription

Roeder, TIBS 1996  
 Margaritis & Holstege, Cell 2008  
 Zaborowska et al., NSMB 2016  
 Ngoc et al., Genes Dev 2017



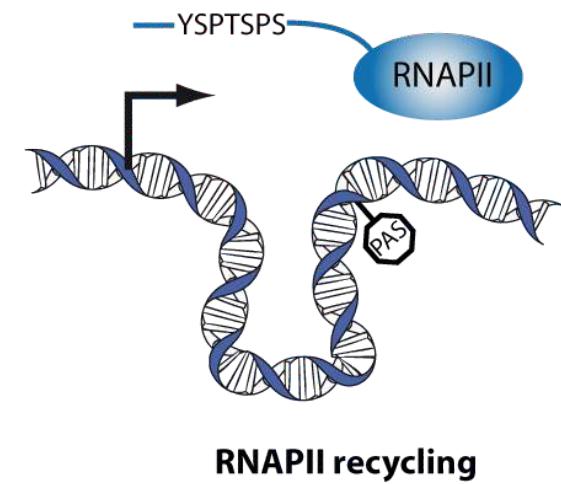
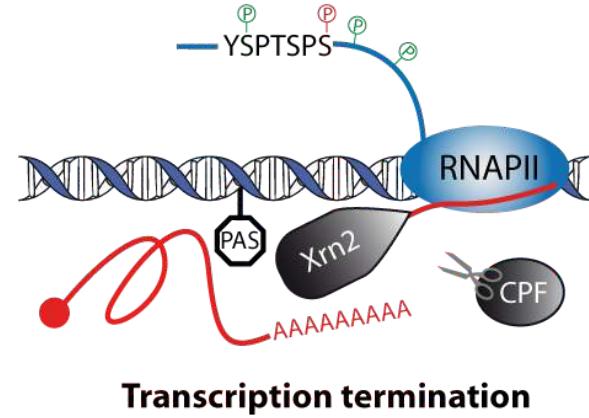
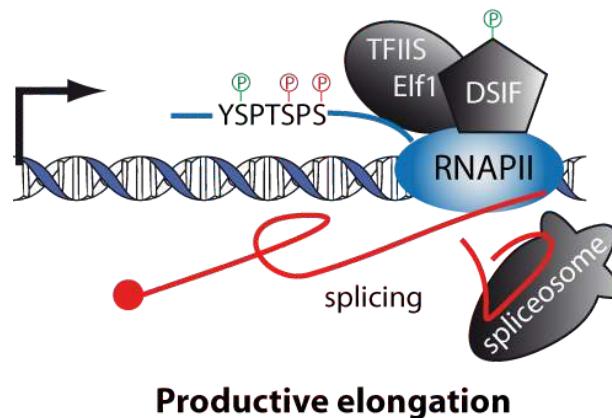
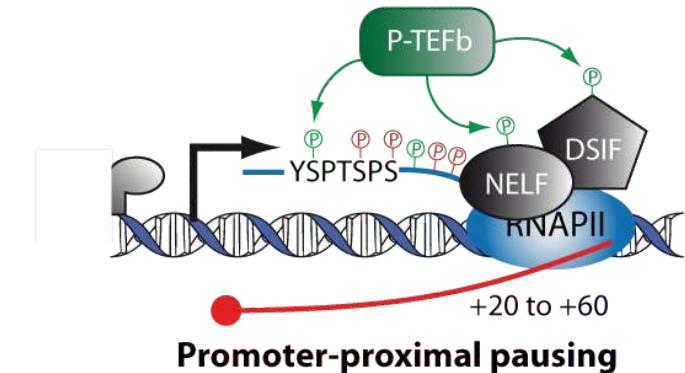
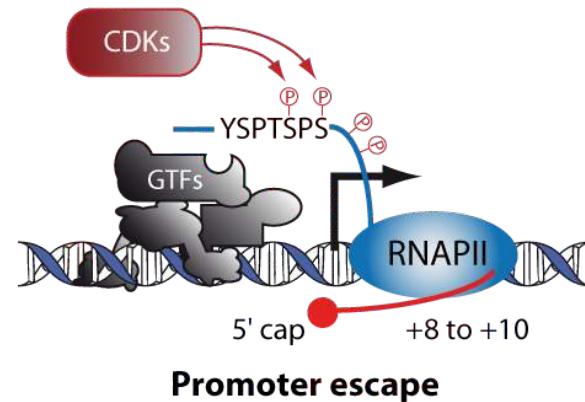
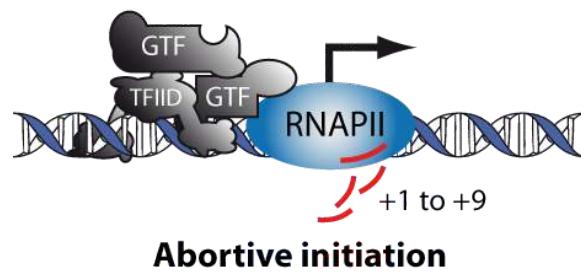
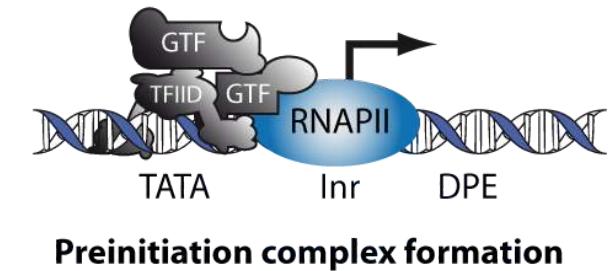
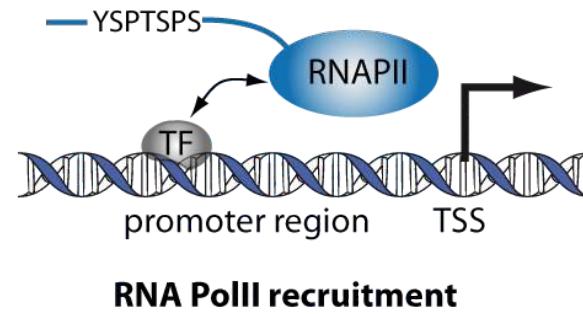
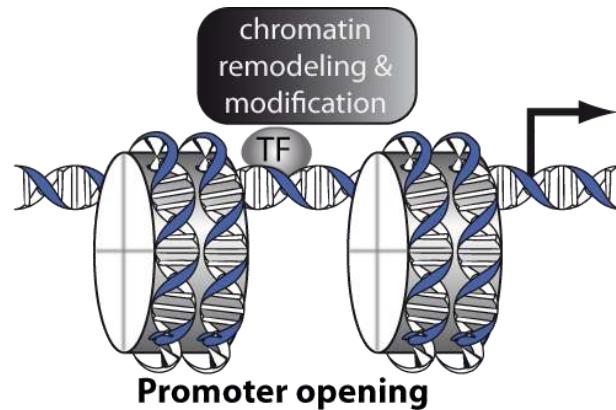
# RNA polymerase II transcription

Roeder, TIBS 1996  
 Margaritis & Holstege, Cell 2008  
 Zaborowska et al., NSMB 2016  
 Ngoc et al., Genes Dev 2017



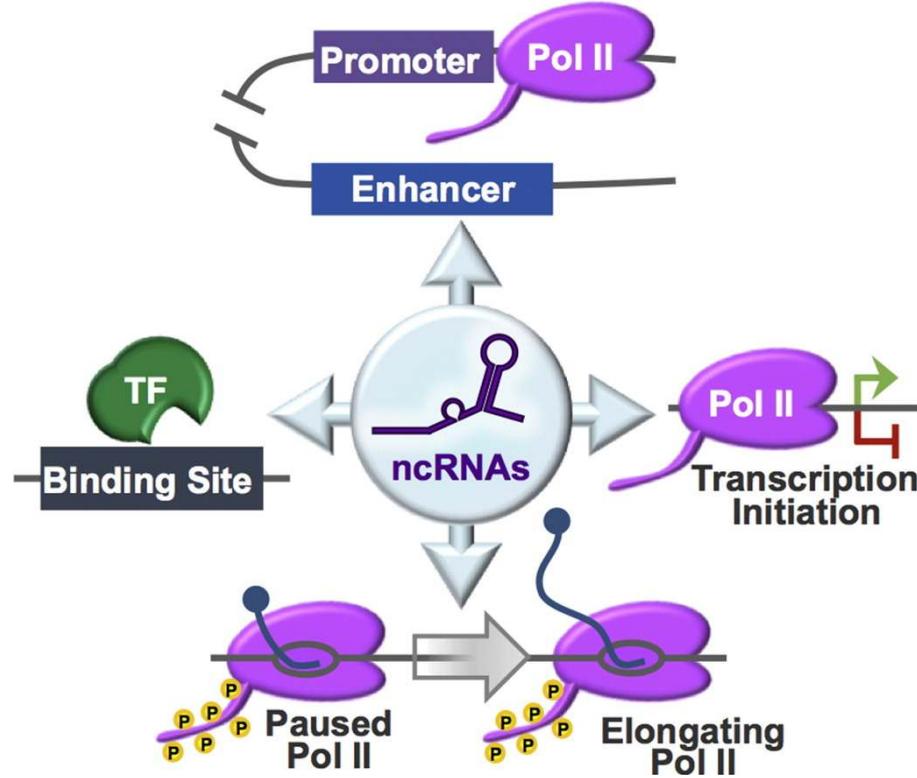
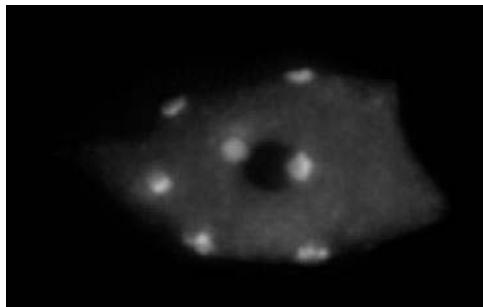
# RNA polymerase II transcription

Roeder, TIBS 1996  
 Margaritis & Holstege, Cell 2008  
 Zaborowska et al., NSMB 2016  
 Ngoc et al., Genes Dev 2017



# Non coding RNAs & DNA methylation

## Nucleus

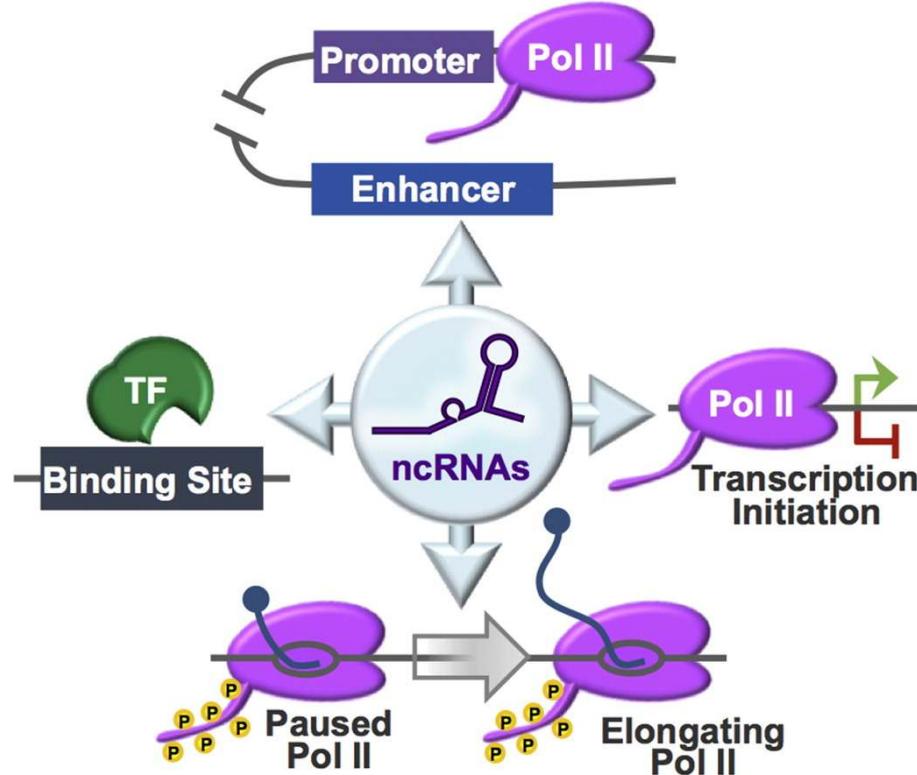
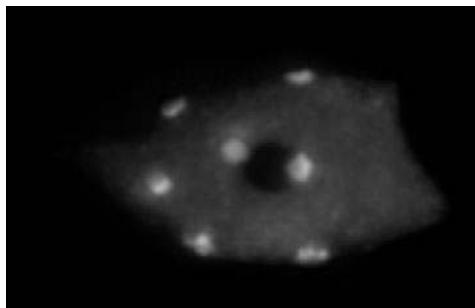


Sabin et al., Mol Cell, 2013  
Kornienko et al., BMC Biol, 2013  
Mercer & Mattick, NSMB, 2013

Eidem et al., J Mol Biol 2016

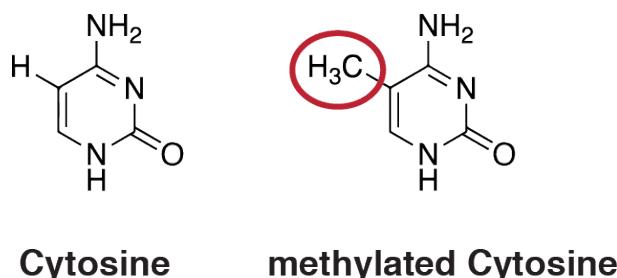
# Non coding RNAs & DNA methylation

## Nucleus

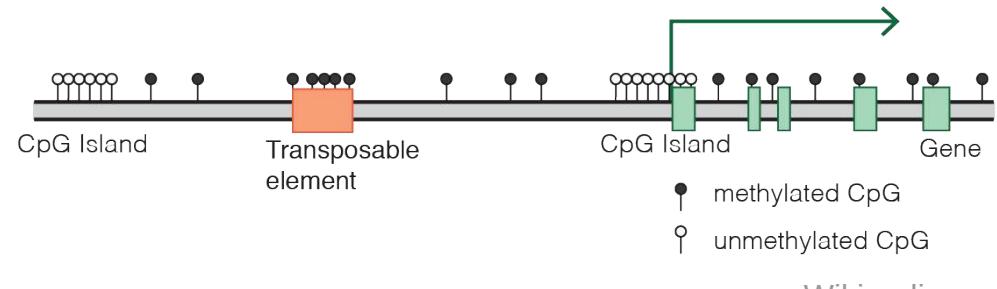


Sabin et al., Mol Cell, 2013  
Kornienko et al., BMC Biol, 2013  
Mercer & Mattick, NSMB, 2013

Eidem et al., J Mol Biol 2016



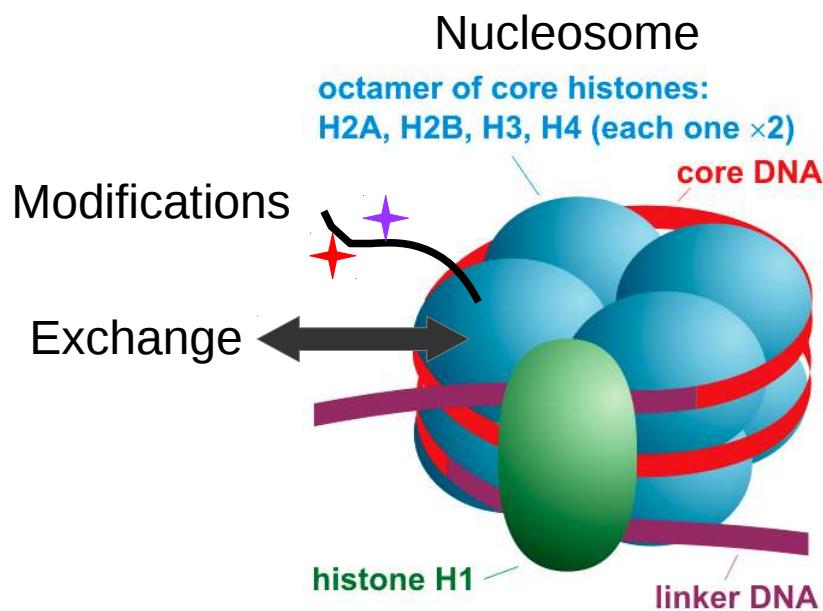
## Typical mammalian DNA methylation landscape



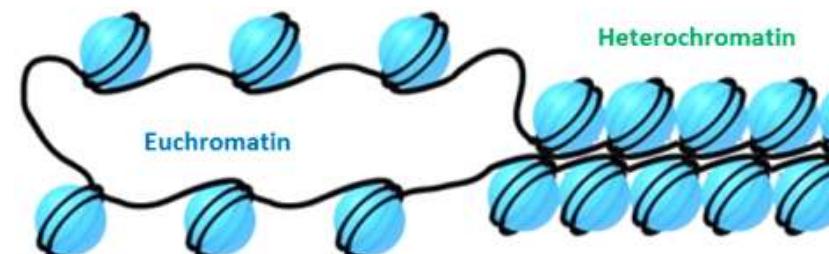
# Chromatin



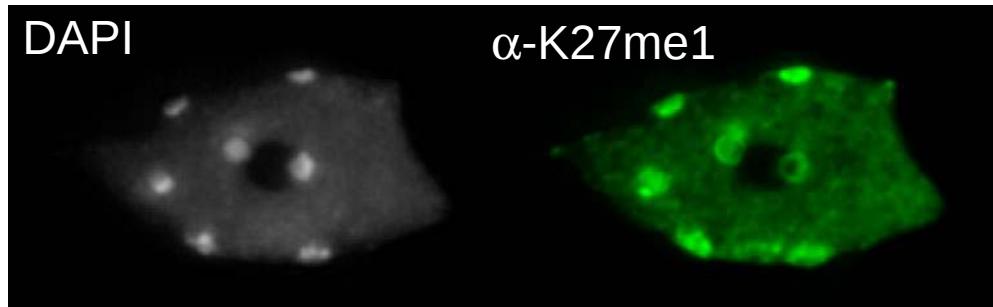
Chromatin = DNA + **proteins**



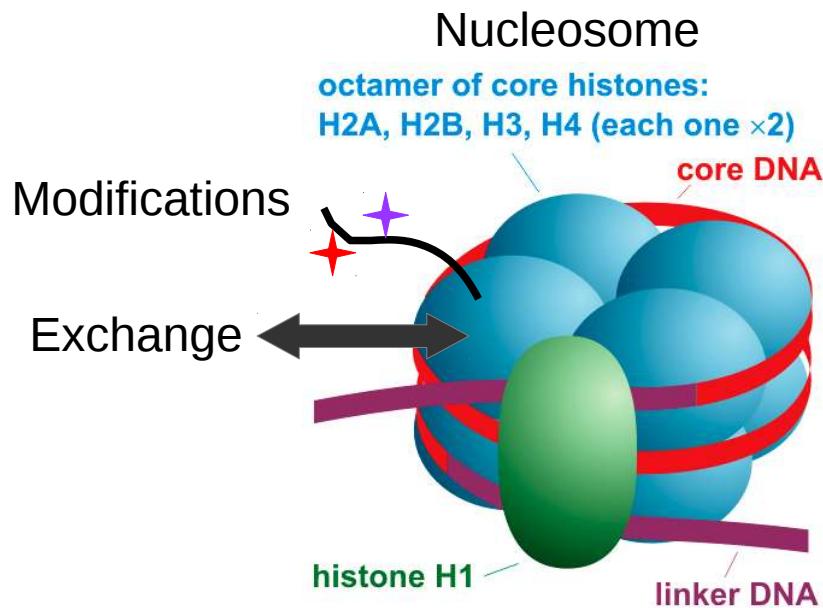
Lai & Pugh, Nat Rev MCB 2017  
Voong et al., Cell 2016  
Bednar et al., Mol Cell 2017



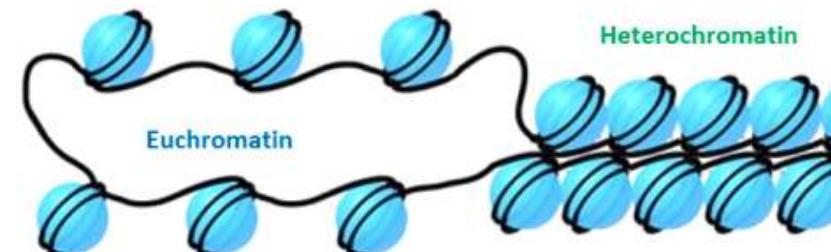
# Chromatin



Chromatin = DNA + **proteins**



Lai & Pugh, Nat Rev MCB 2017  
Voong et al., Cell 2016  
Bednar et al., Mol Cell 2017



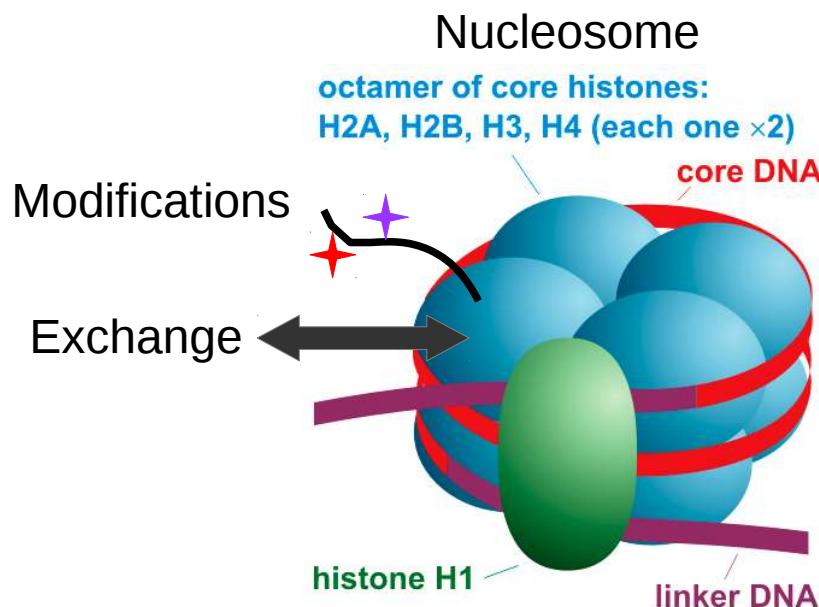
# Chromatin



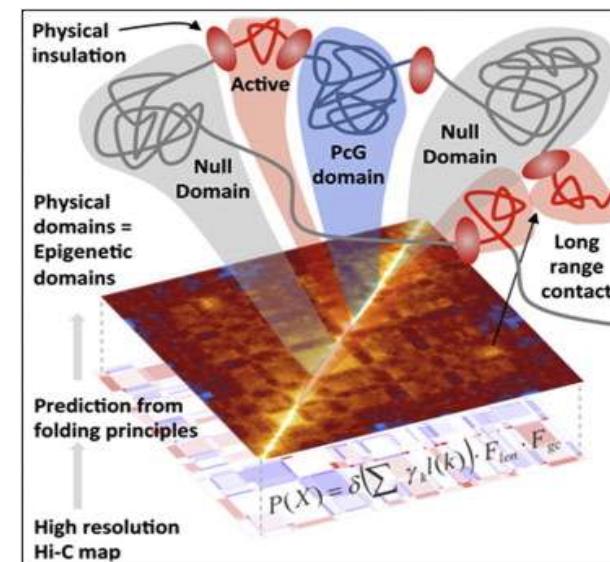
Chromatin = DNA + **proteins**



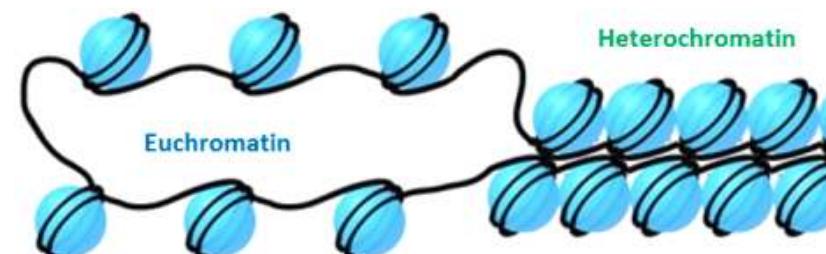
3D chromatin organization



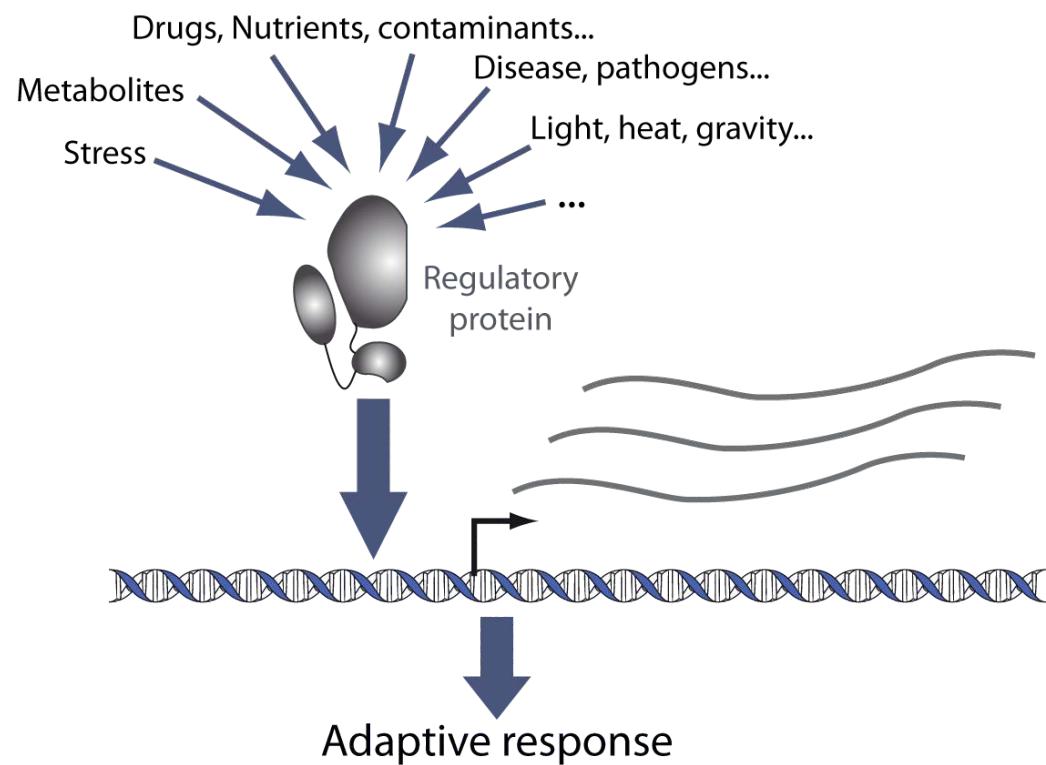
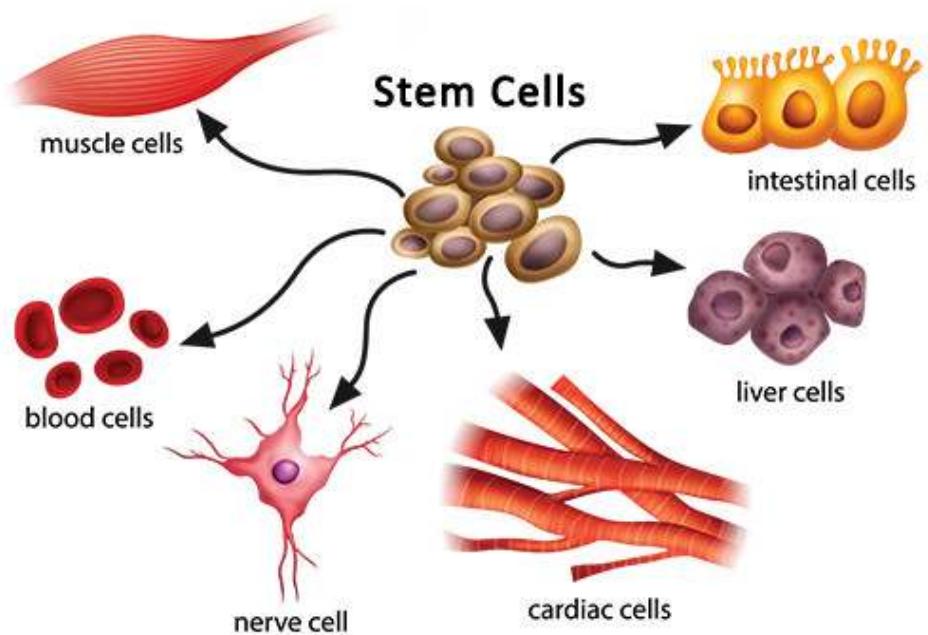
Lai & Pugh, Nat Rev MCB 2017  
Voong et al., Cell 2016  
Bednar et al., Mol Cell 2017



Sexton et al., Cell, 2012  
Sanyal et al., Nature 2012  
Nora et al., Nature 2012

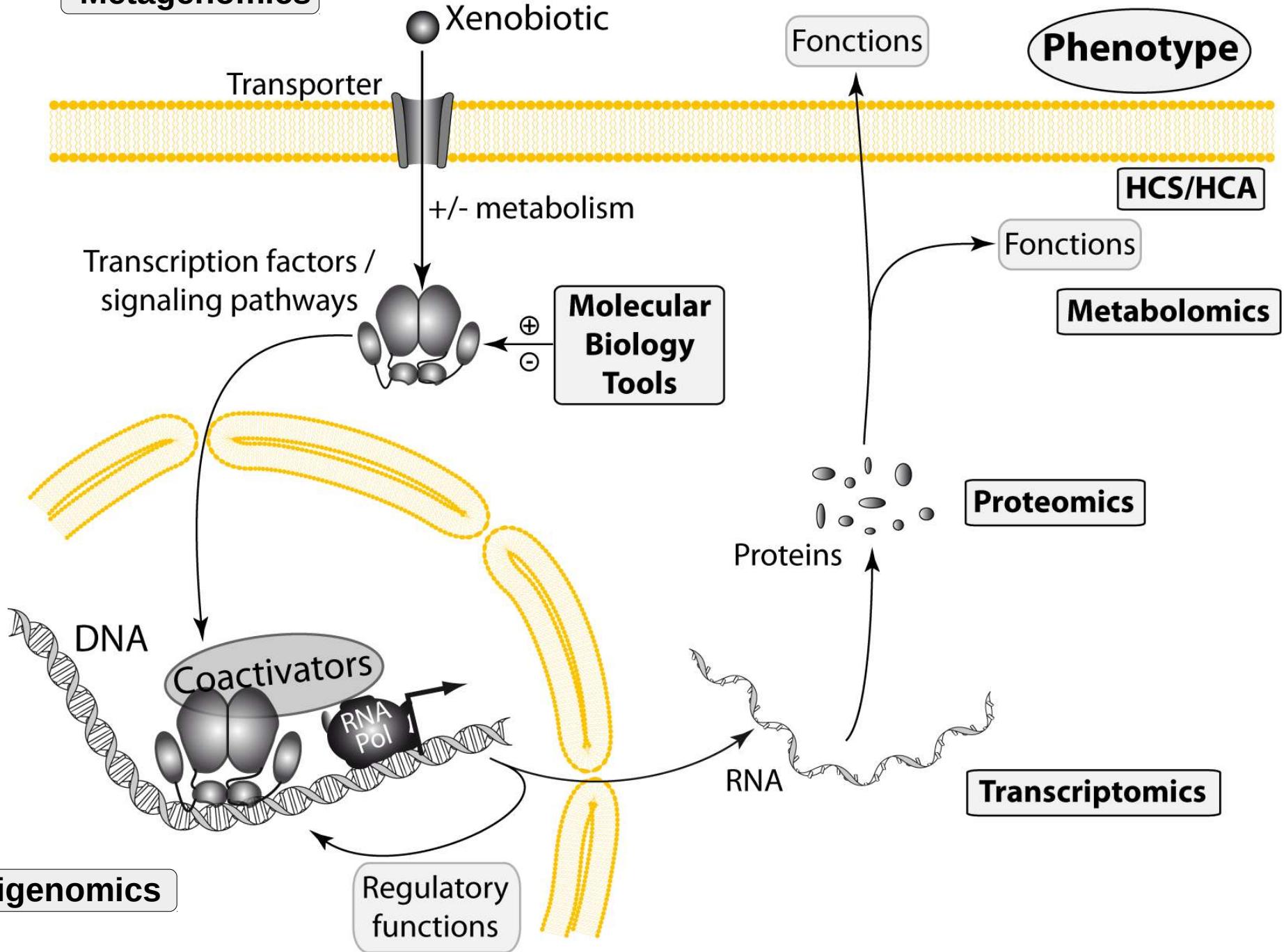


# Importance of gene expression regulations



# Toxicogenomics

Metagenomics



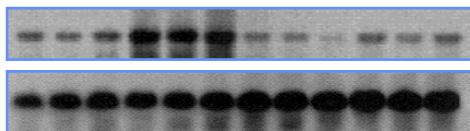
# Transcriptomics

# Quantify RNA abundance

Low throughput / targeted

Few genes

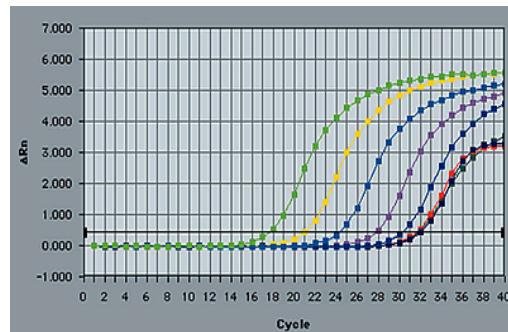
Northern blot / dot blot



Ribonuclease protection assay (RPA)



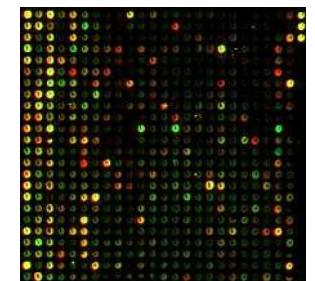
Real-time quantitative PCR (qPCR)



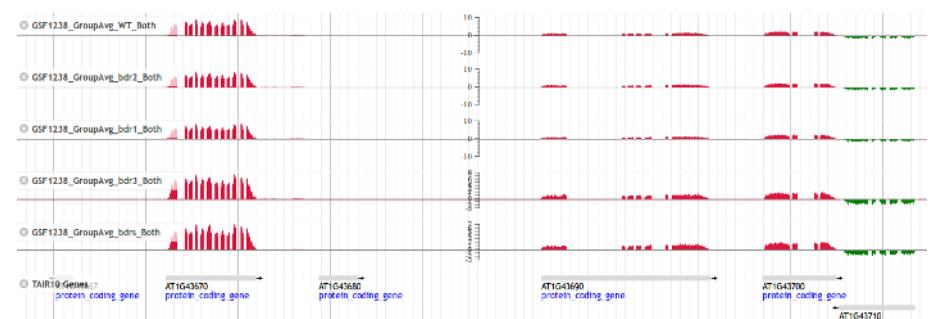
High throughput / genome wide

Many / all genes

Microarrays

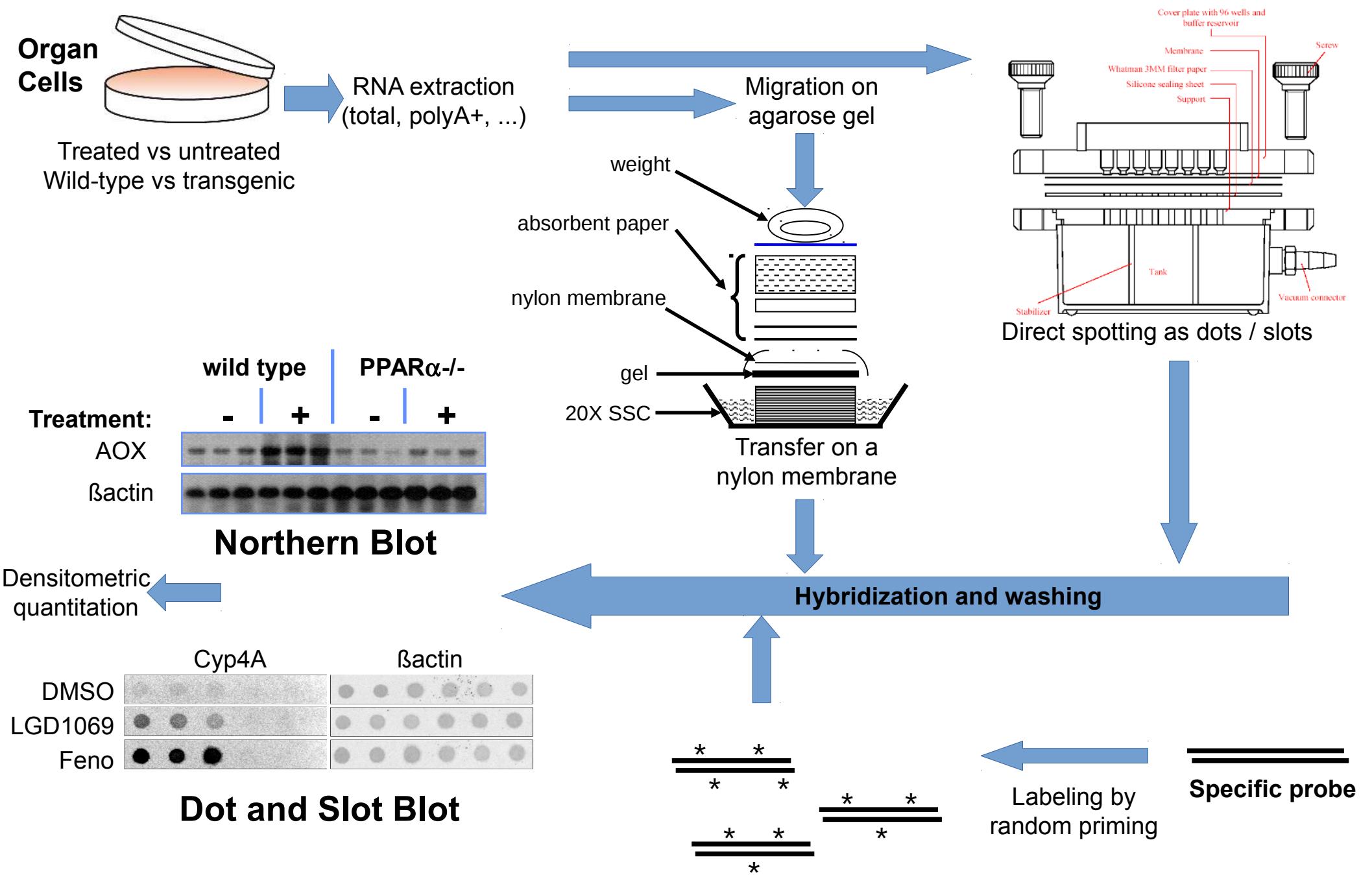


RNA-seq



# "Low" throughput methods

# Northern blot



# Northern blot

---

## Advantages:

- Detect alternative transcripts (if size are sufficiently different and probe is common to the different transcripts)
- Evaluate transcript(s) size
- Detect regions / exons shared among transcripts
- Low/medium sensitivity to RNA degradation + visualize / evaluate RNA degradation
- Low cost
- Quantitative and qualitative (see the "bands")

## Drawbacks / limits:

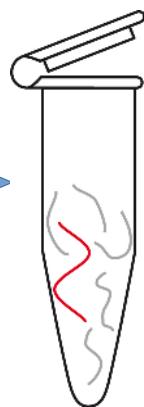
- Need to know the sequence (probe ~1Kb)
- Specificity may be hard to obtain (multigenic family)
- Low sensitivity
- Hard to set up (tedious, long, "dirty" : Ethidium Bromide, formaldehyde, formamide,  $^{32}\text{P}$ )

# RNase protection assay (RPA)

Organ Cells

Treated vs untreated  
Wild-type vs transgenic

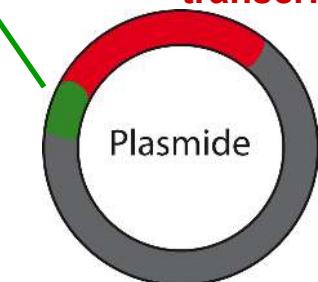
RNA extraction



Total RNA

Promoter of SP6 or  
T7 RNA polymerase

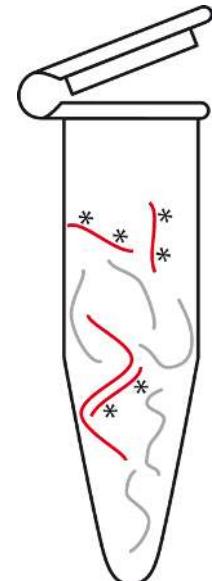
cDNA of studied  
transcript



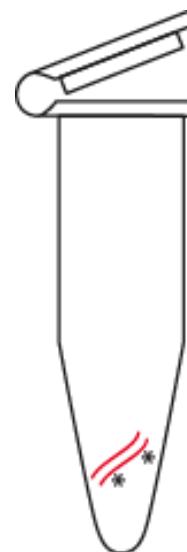
IVT + purification  
of labeled probe

antisense RNA probe  
(~100-500mers)

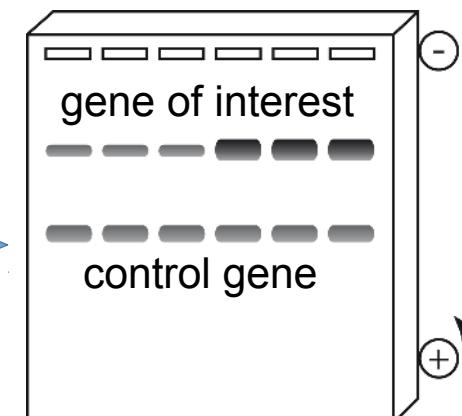
RNases digestion



Size purification



Electrophoresis



HYBRIDIZATION

Densitometric  
quantitation

## Advantages:

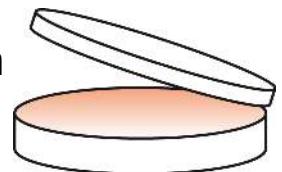
- Liquid phase hybridization => better sensitivity than NB (low abundance transcripts)
- Increase starting RNA (up to 100µg) => detect low abundance transcripts
- No expensive instrument
- Short probe => relatively low sensitivity to moderate RNA degradation
- Can study up to 10 different transcripts simultaneously (probes of different sizes)
- Easier to discriminate alternative transcripts of similar size
- Allows to study transcript ends and exon-exon junctions

## Drawbacks / limits:

- No info on transcript size
- Obtaining the probe can be tedious
- Manipulation of material that is sensitive to degradation (single stranded RNA probe)
- Sensitive to small sequence differences between the probe and the target

# Real-time quantitative PCR

Organ  
Cells



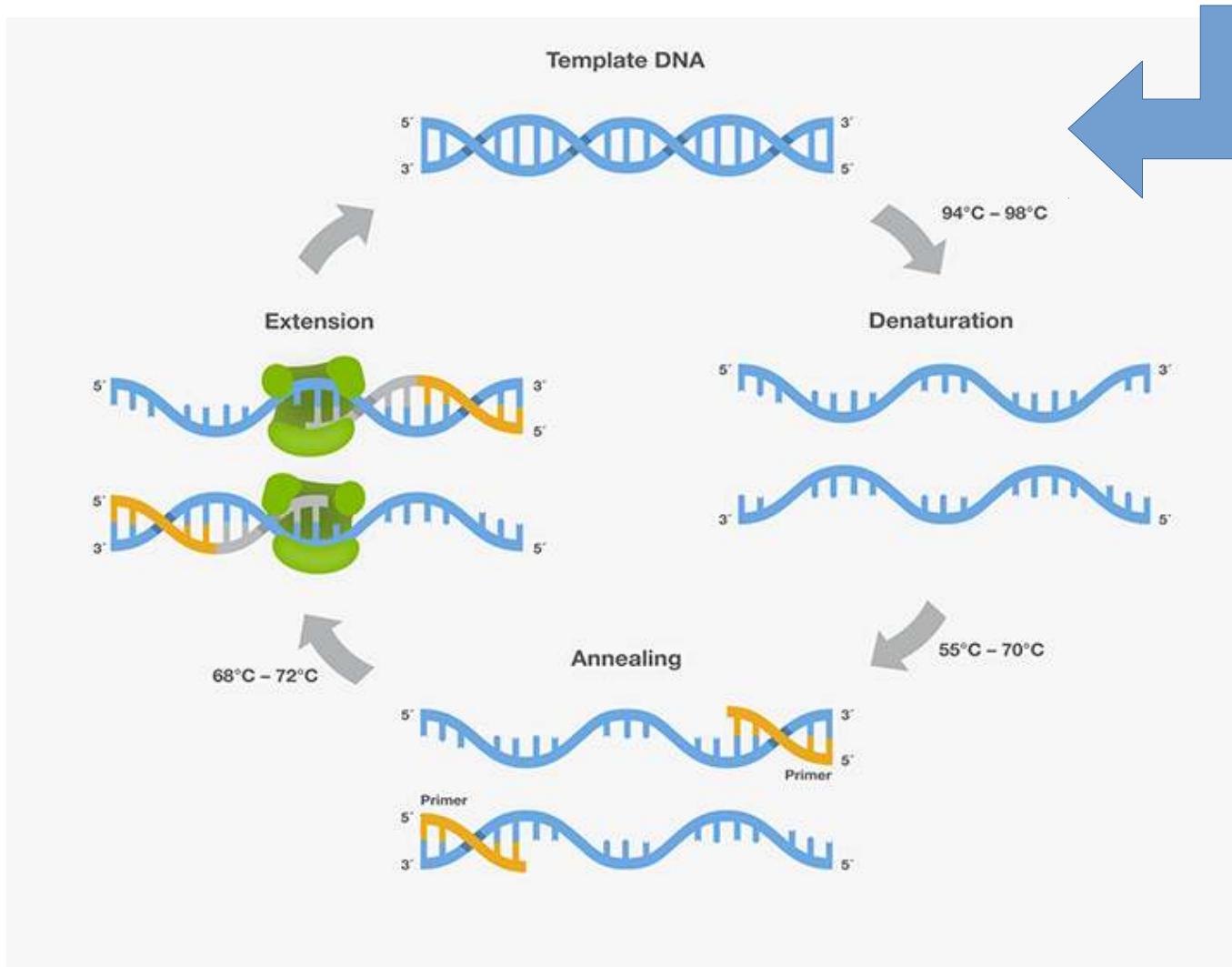
Treated vs untreated  
Wild-type vs transgenic

RNA extraction

Total RNA

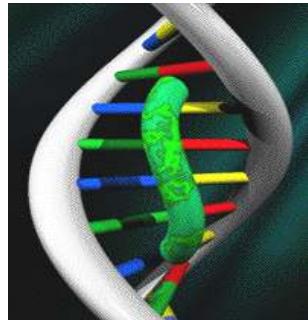
Reverse  
transcription

- oligodT
- random primers
- specific primers



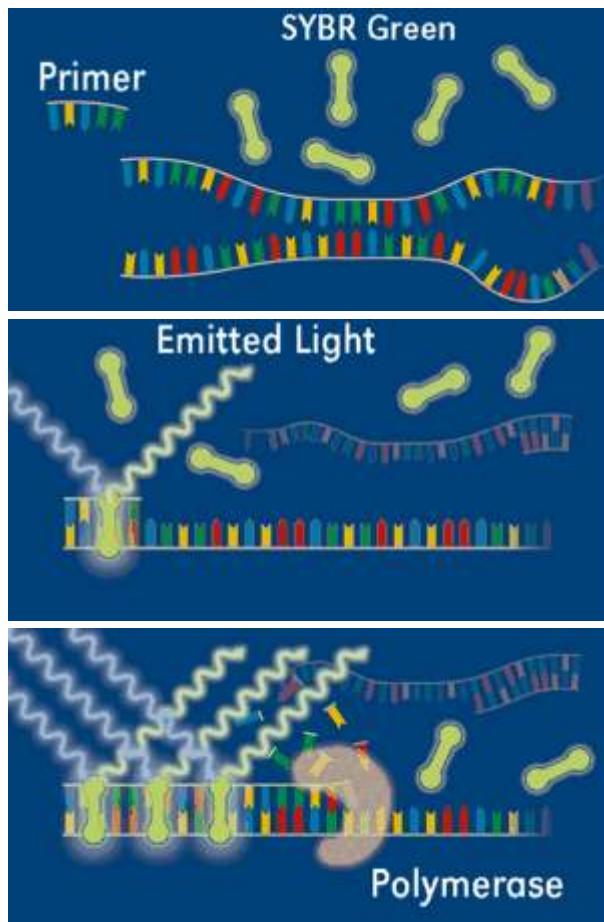
# qPCR chemistries

Follow in real-time the production of DNA amplicons



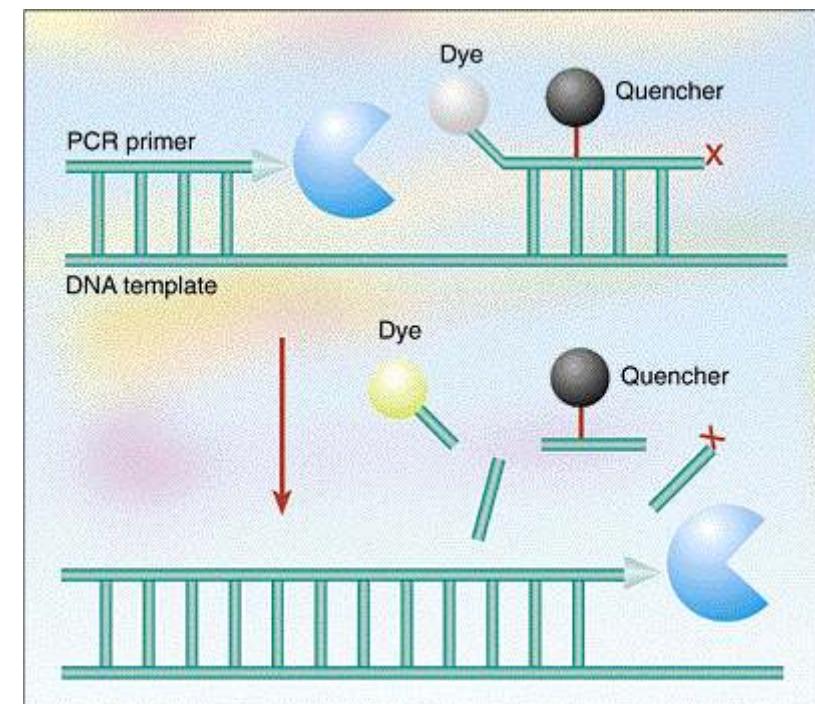
Intercalating agent  
(SYBR green, Eva Green...)

non sequence-specific

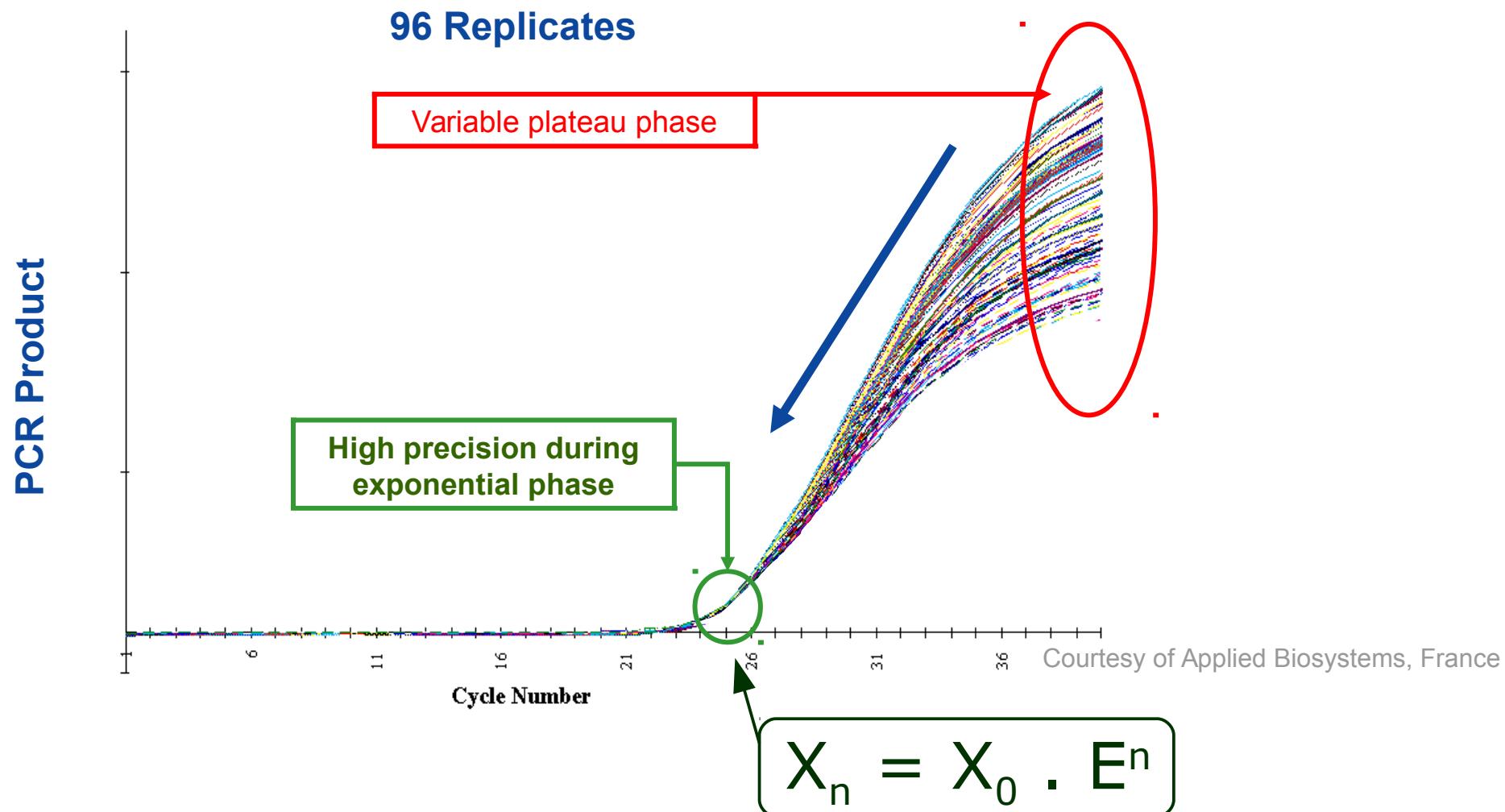


Probe-based chemistries  
(TaqMan, Lightcycler probes, ...)

sequence-specific => multiplexing



# qPCR reaction

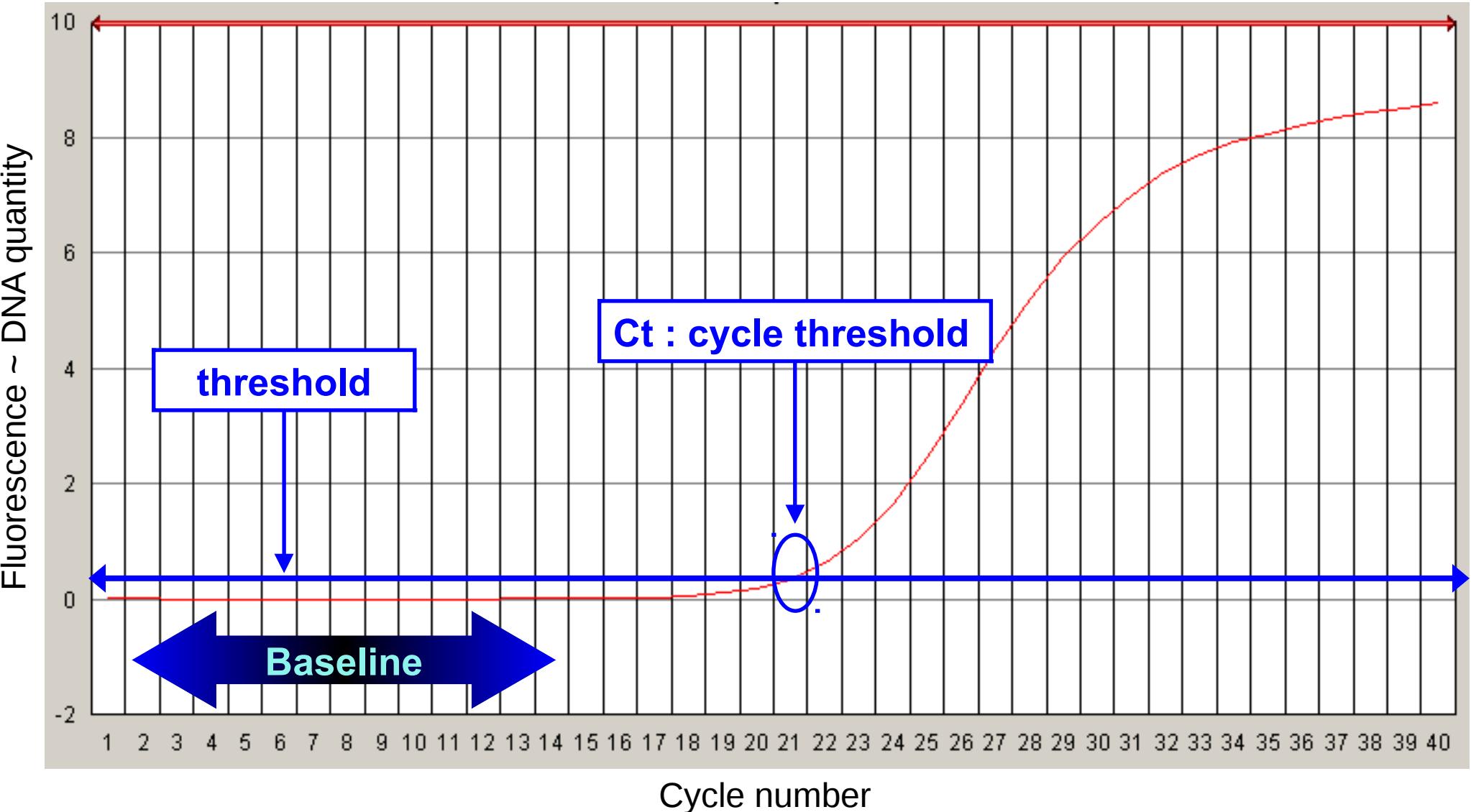


X<sub>n</sub> : # amplicons at cycle n

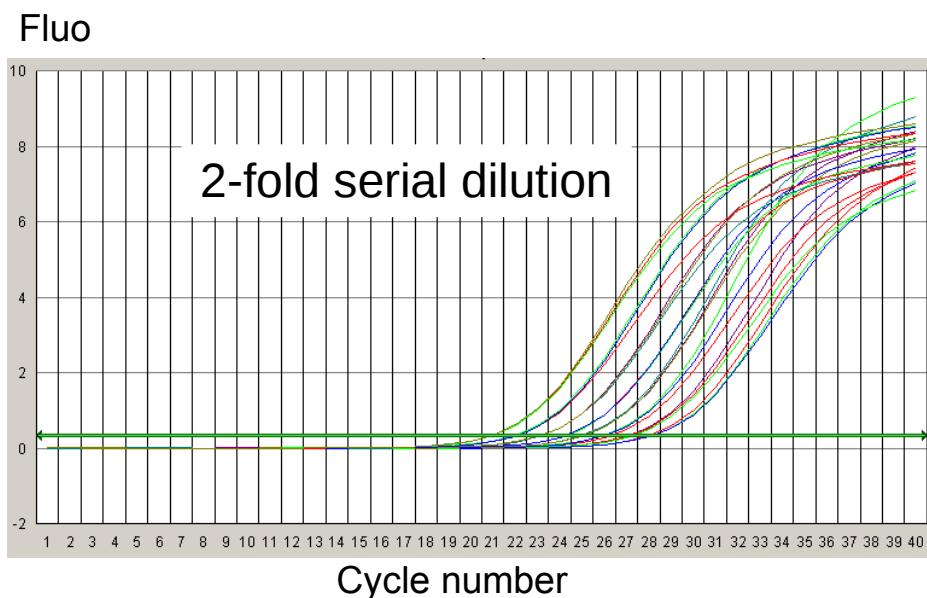
X<sub>0</sub> : initial # of transcripts

E = PCR efficiency in [1:2]  
(E=2  $\Leftrightarrow$  100% efficiency)

# qPCR reaction

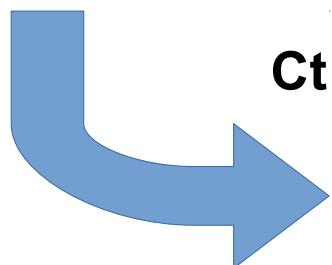
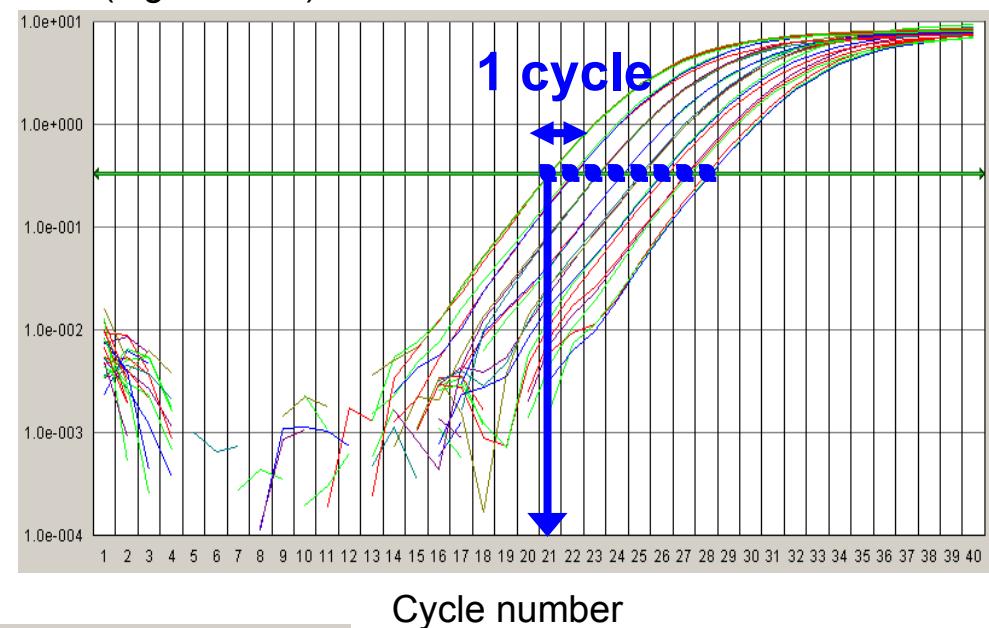


# qPCR reaction



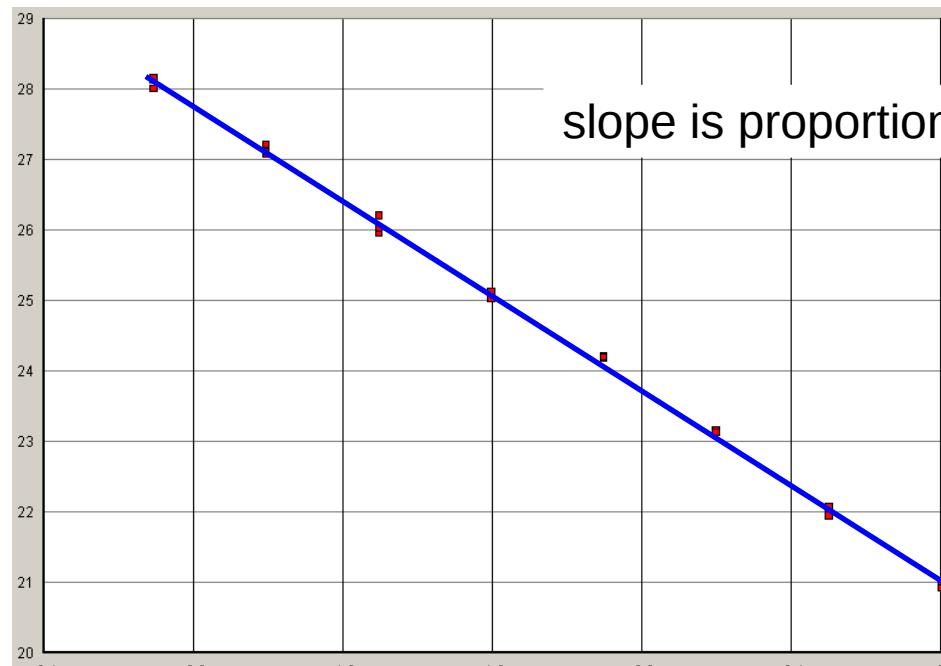
log

Fluo (log scale !!!)



Ct

slope is proportional to PCR efficiency



log(concentration)

### 3 fundamental principles:

- 1) Fluorescence is proportional to amplicon accumulation  
(specificity, primer/probe selection)
- 2) PCR efficiency is the same for all samples
- 3) Threshold must intersect all PCR curves at the exponential phase

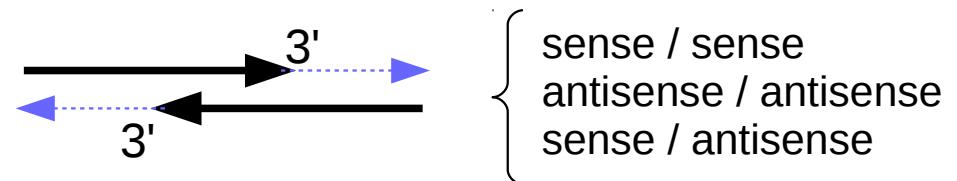
Peirson et al., NAR 2003

## qPCR primer design

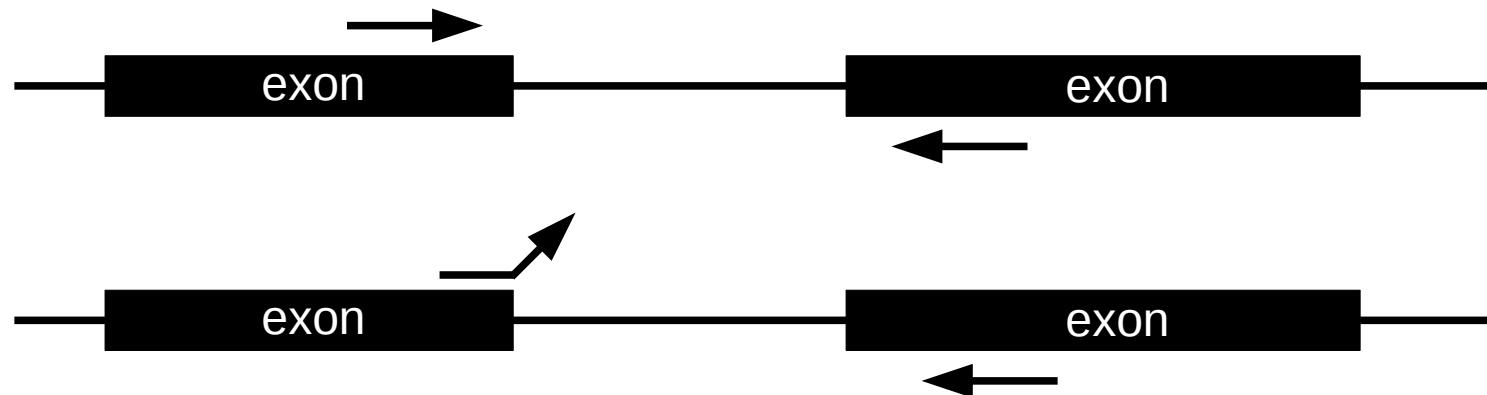
### 1) Fluorescence is proportional to amplicon accumulation (specificity, primer/probe selection)

- Design specific primers ("know your target")

- Avoid 3'-3' primer dimers



- Take advantage of DNA – RNA differences

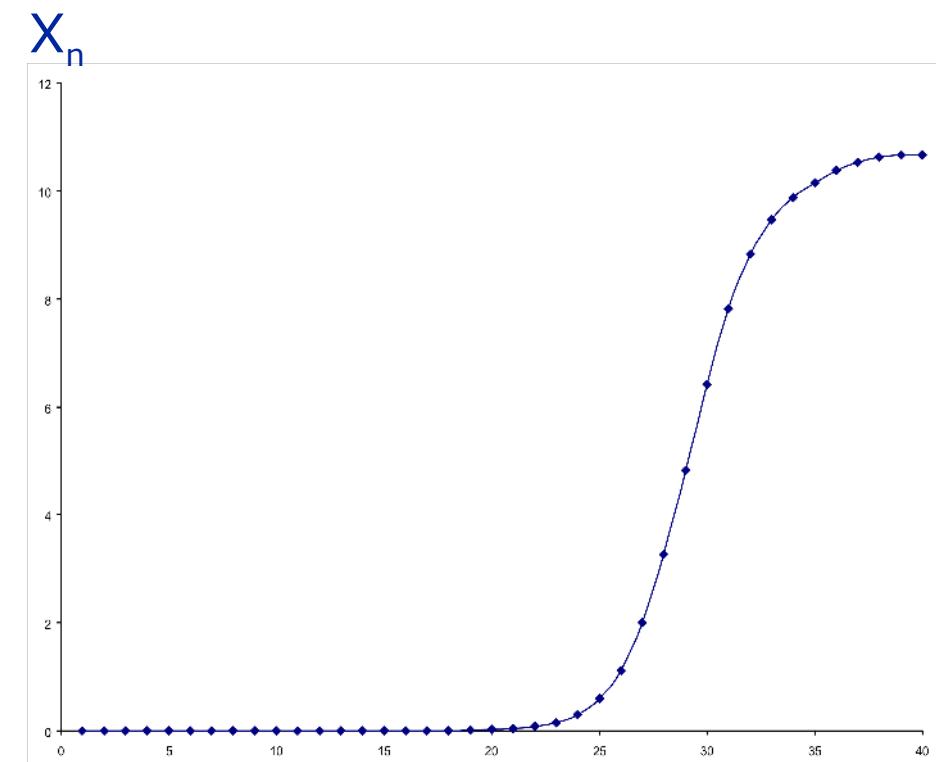
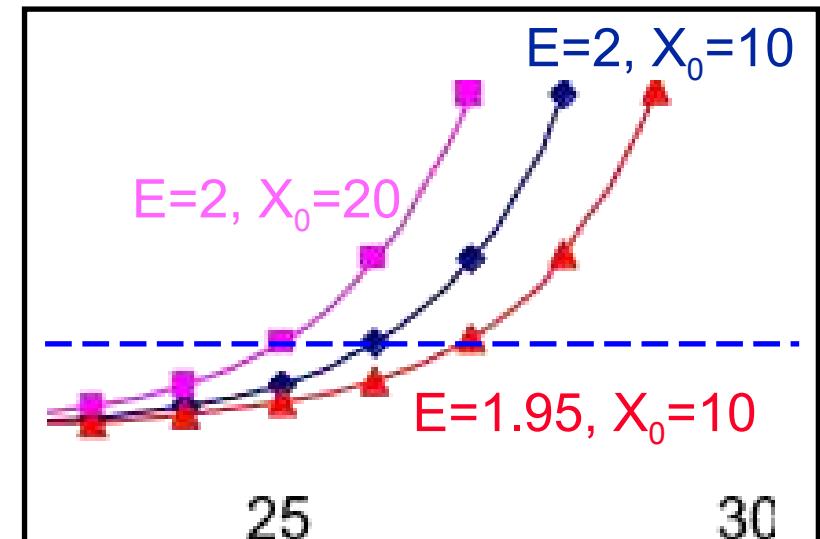


## qPCR efficiency

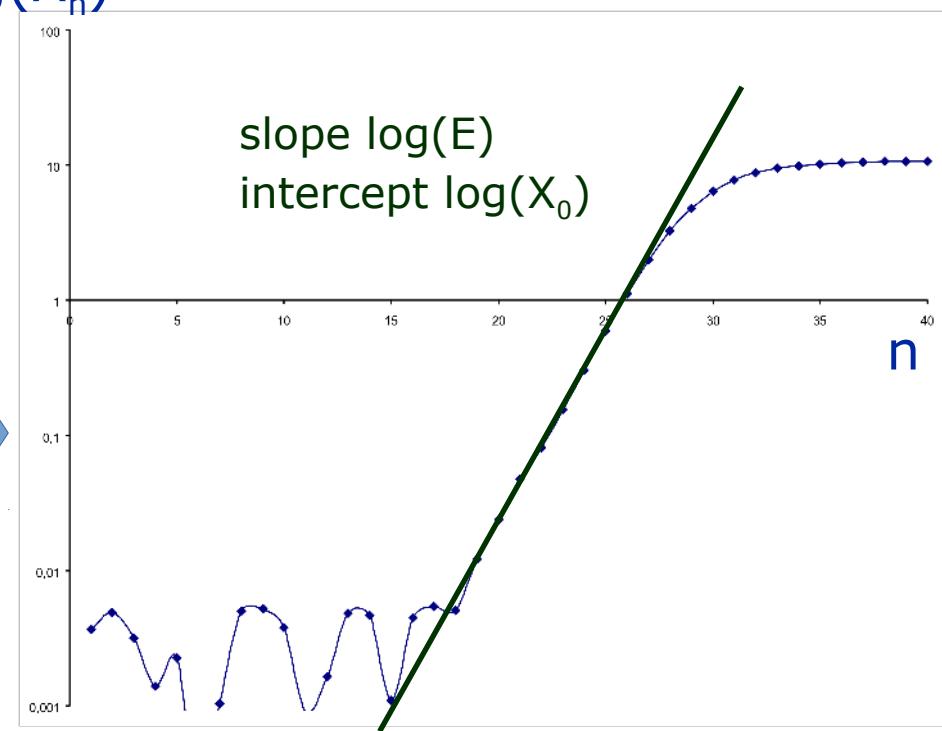
2) PCR efficiency is the same for all samples

$$X_n = X_0 \cdot E^n$$

log →  $\log(X_n) = \log(X_0) + n \cdot \log(E)$

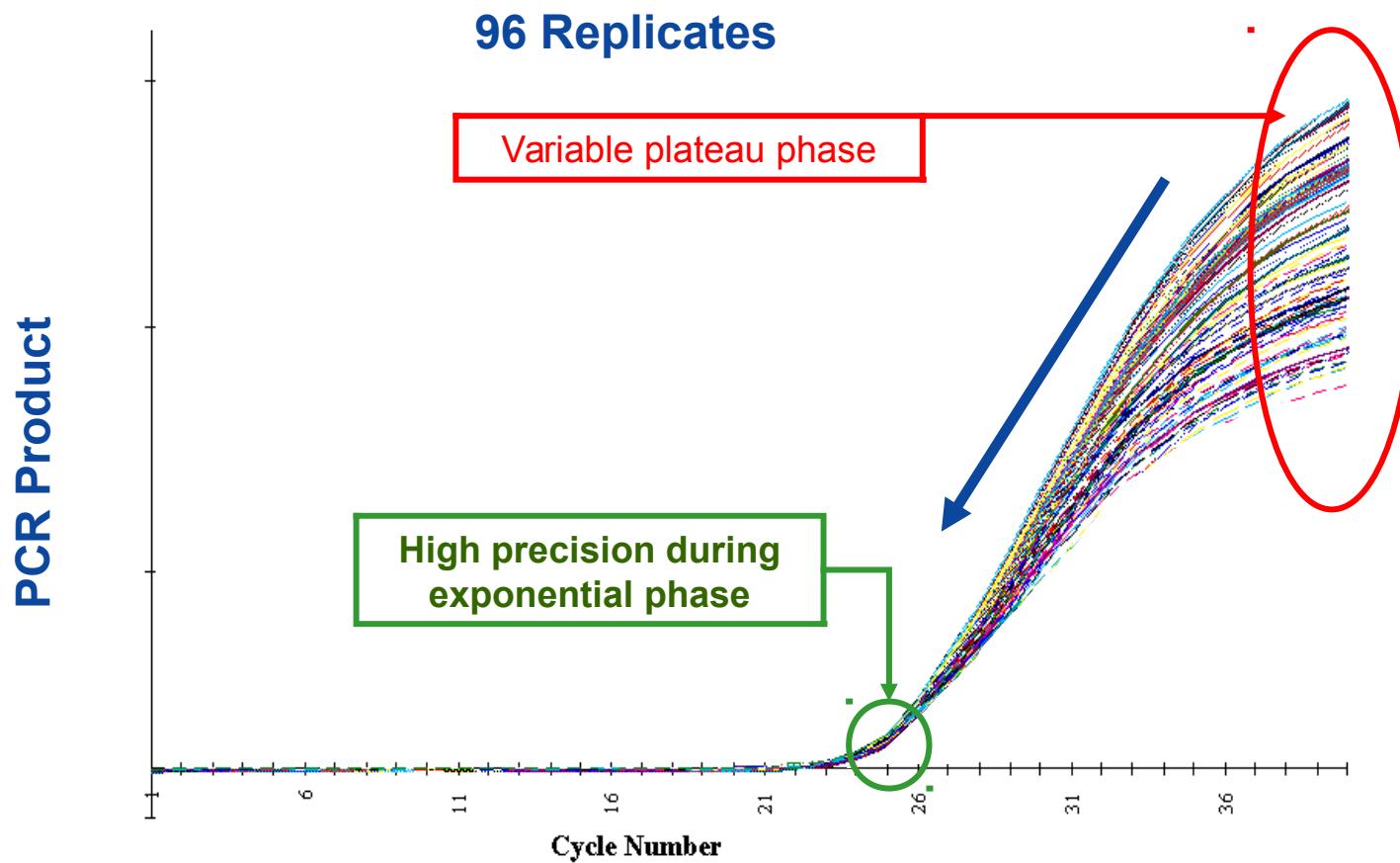


log →  $\log$



## qPCR primer design

3) Threshold must intersect all PCR curves at the exponential phase



→ Ct is linked to initial number of RNA molecules and not to different PCR kinetics

# qPCR

## Advantages:

- High sensitivity
- Easy and quick set-up
- Low sensitivity to RNA degradation (short amplicons)
- Highly specific (e.g. single base discrimination at primer 3' end)
- Multiple applications (gene expression, SNP genotyping, CNV, ChIP, etc.)
- Miniaturization



Roche Lightcycler 1536



1536 multiwell plates



Biomark HD (Fluidigm)



96 samples x 96 primers = 9216 reactions  
12-48 samples x ~770 primers : up to ~36K reactions

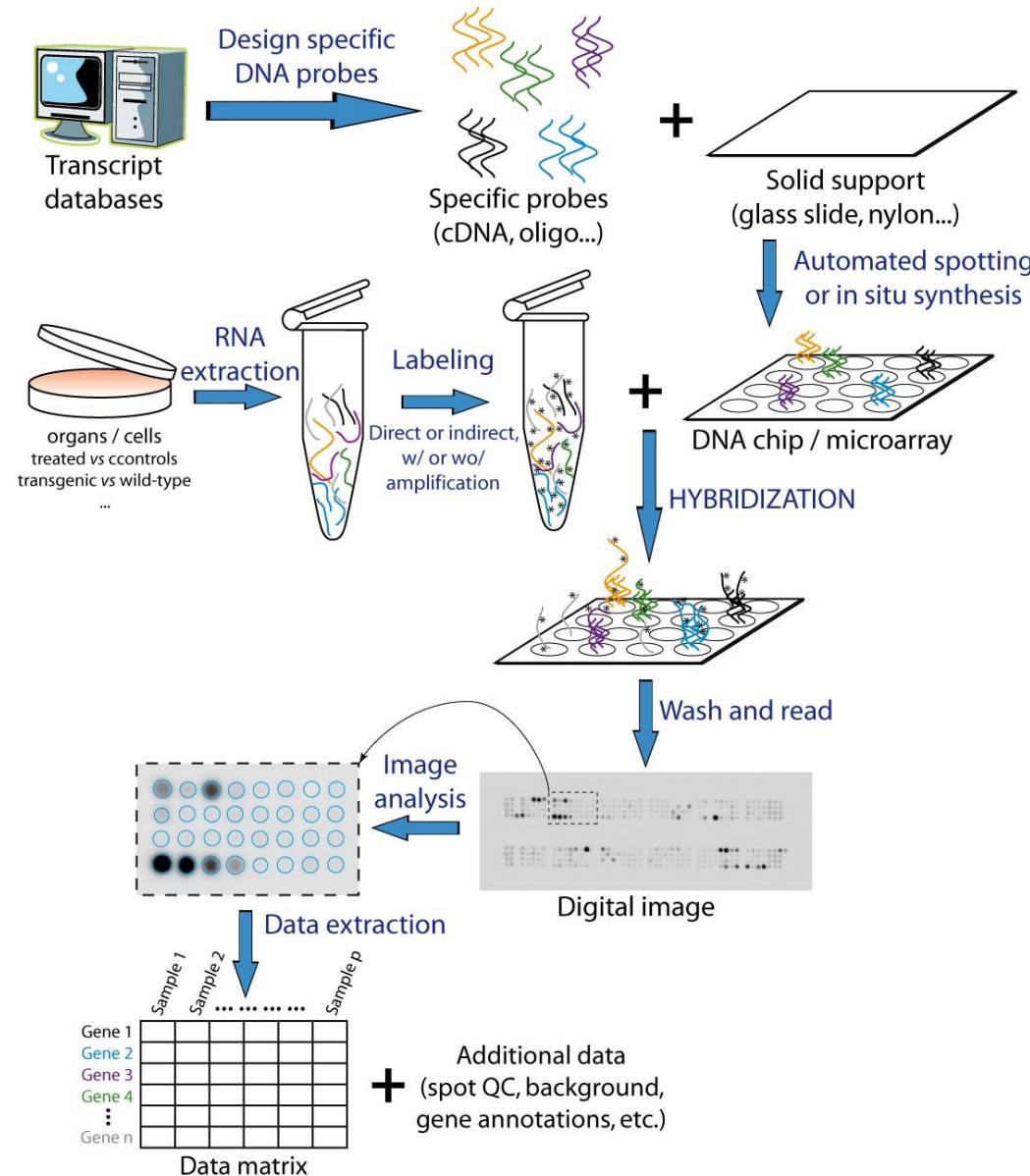
## Drawbacks / limits:

- Sensitive to PCR inhibitors (salts, ionic detergents, alcohols, ...)
- Need precise sequence information
- Data analysis a bit more complicated
- No info on transcript size or presence of alternative transcripts

# High throughput methods

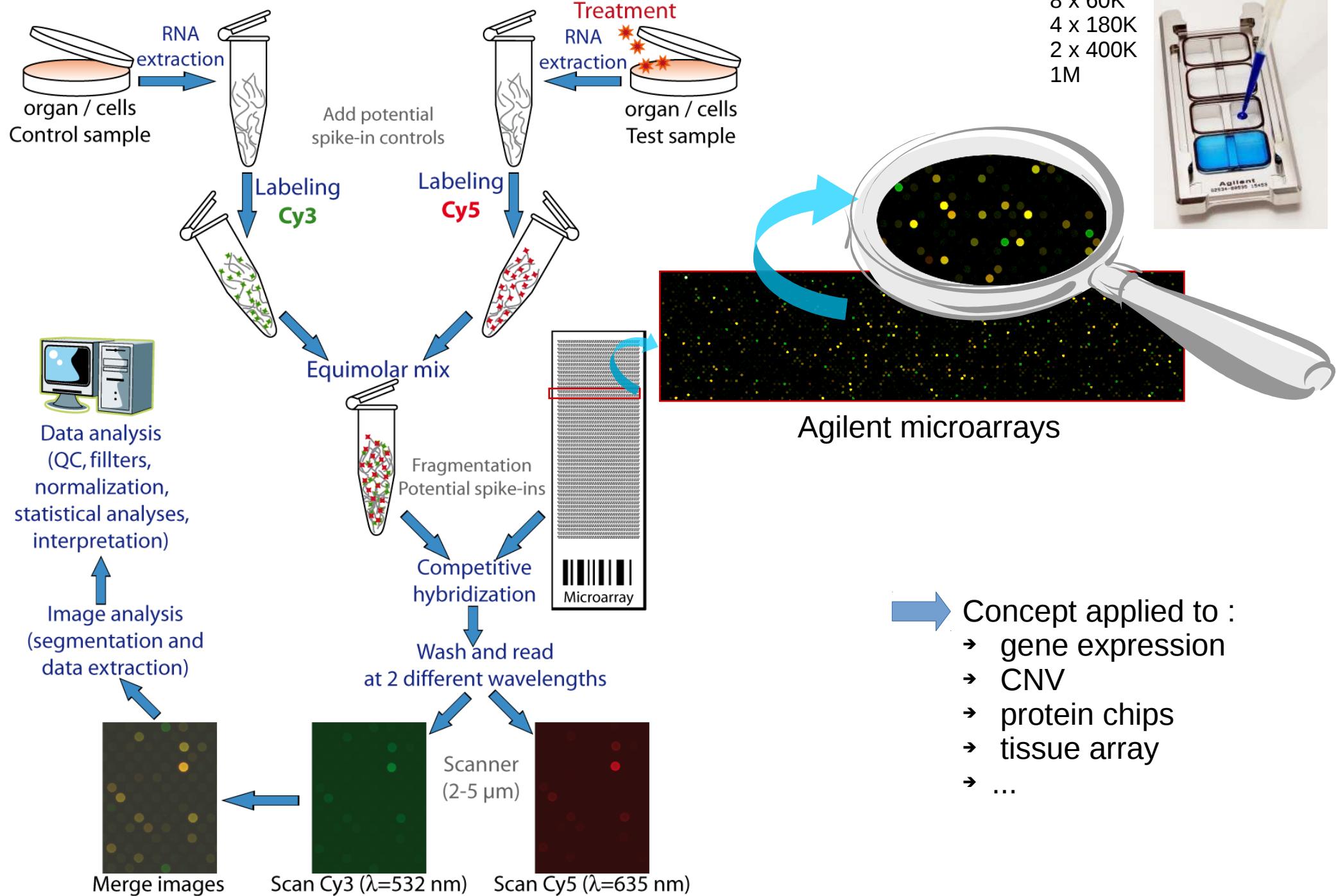
# Microarrays

Developed in the mid 90s thanks to improved knowledge of whole genome sequences



Affymetrix  
GeneChip

# Dual color microarrays



# Microarrays

---

## Advantages:

- Fast and easy to use
- Good reproducibility (better intra-platform, MAQC)
- Excellent transcriptome coverage (when sequence is known and annotated)
- Low cost (~150-200 € / sample)
- Measurements for different genes / transcripts are essentially independent from each other

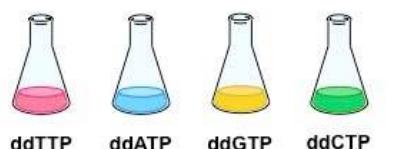
## Drawbacks / limits:

- Need a well annotated genome / transcriptome
- Cross hybridization can occur
- Expensive instruments (scanner) => facilities
- Data analysis is non trivial

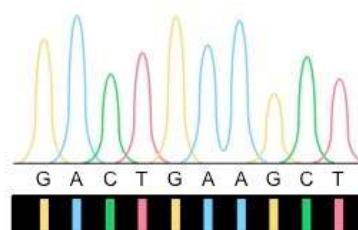
# NGS development

## Sanger sequencing (70s-80s)

4 x PCR (+ one dideoxynucleotide)



Use a sequencing machine

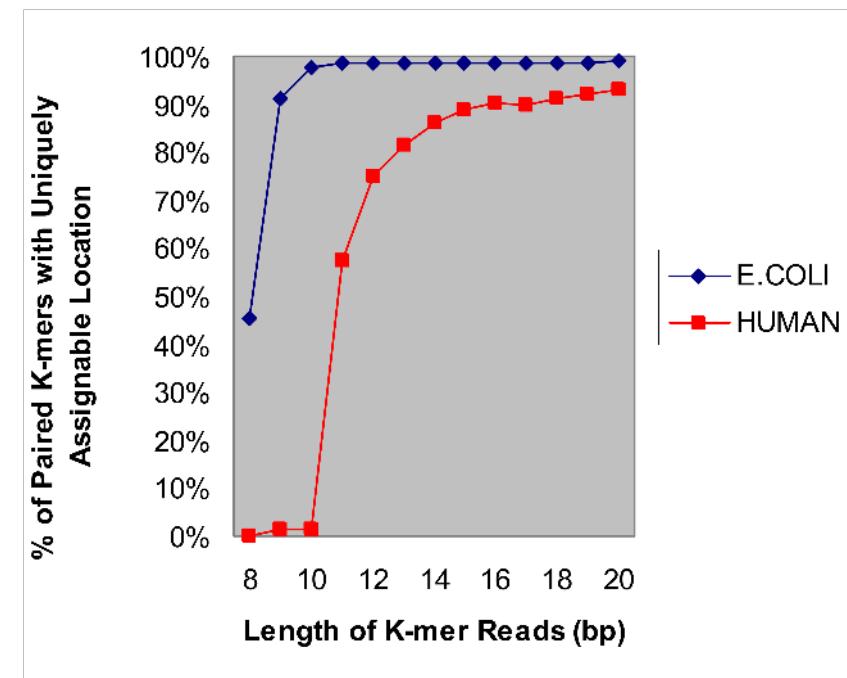


DNA sequence

Separate with a gel



For resequencing, short reads work



Jay Shendure

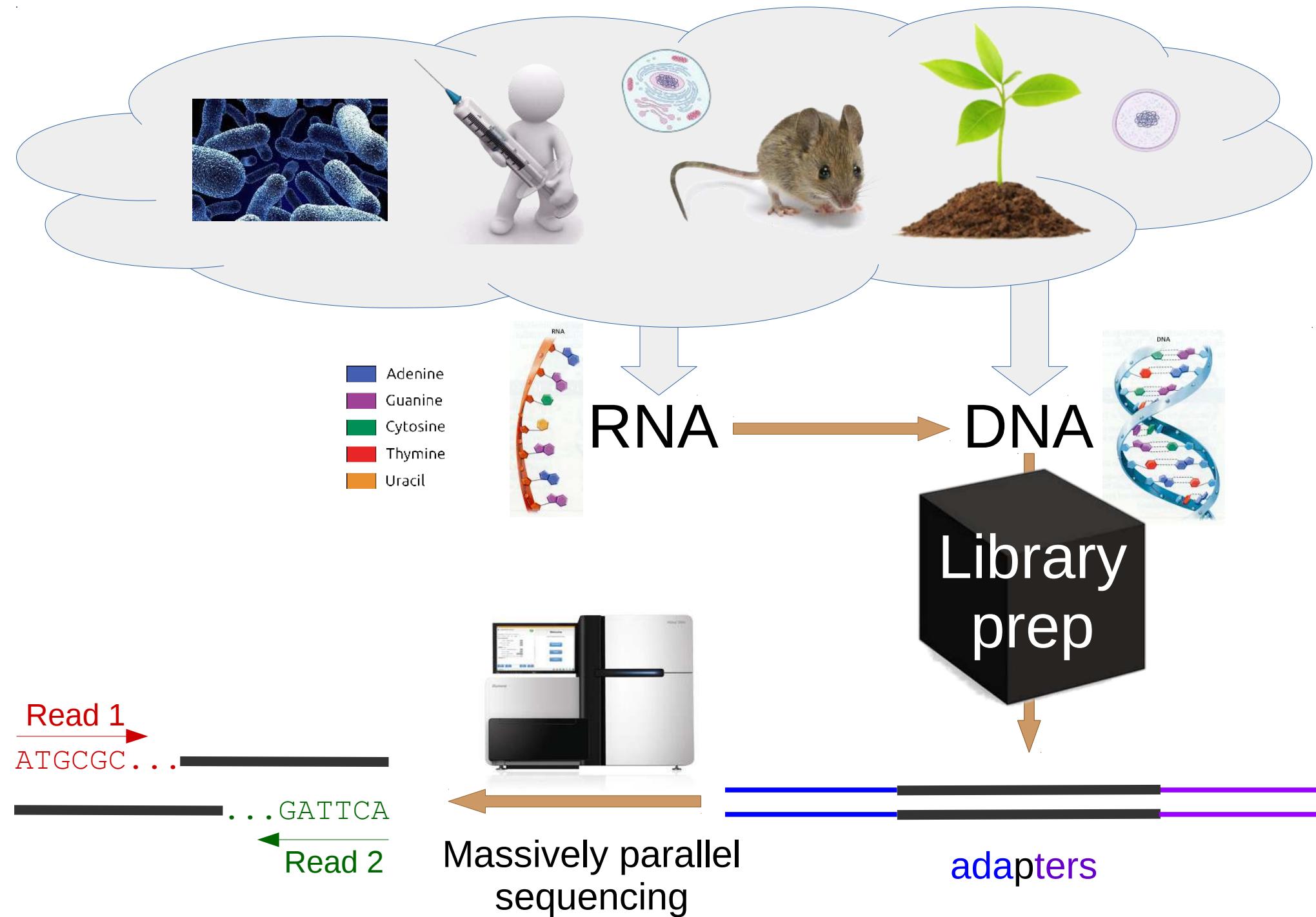
→ Human Genome Project (1990-2003)

→ Development of technologies for massively sequencing short tags (e.g. MPSS)

Brenner et al., Nat Biotech 2000

→ Development of currently used technologies (~2005)

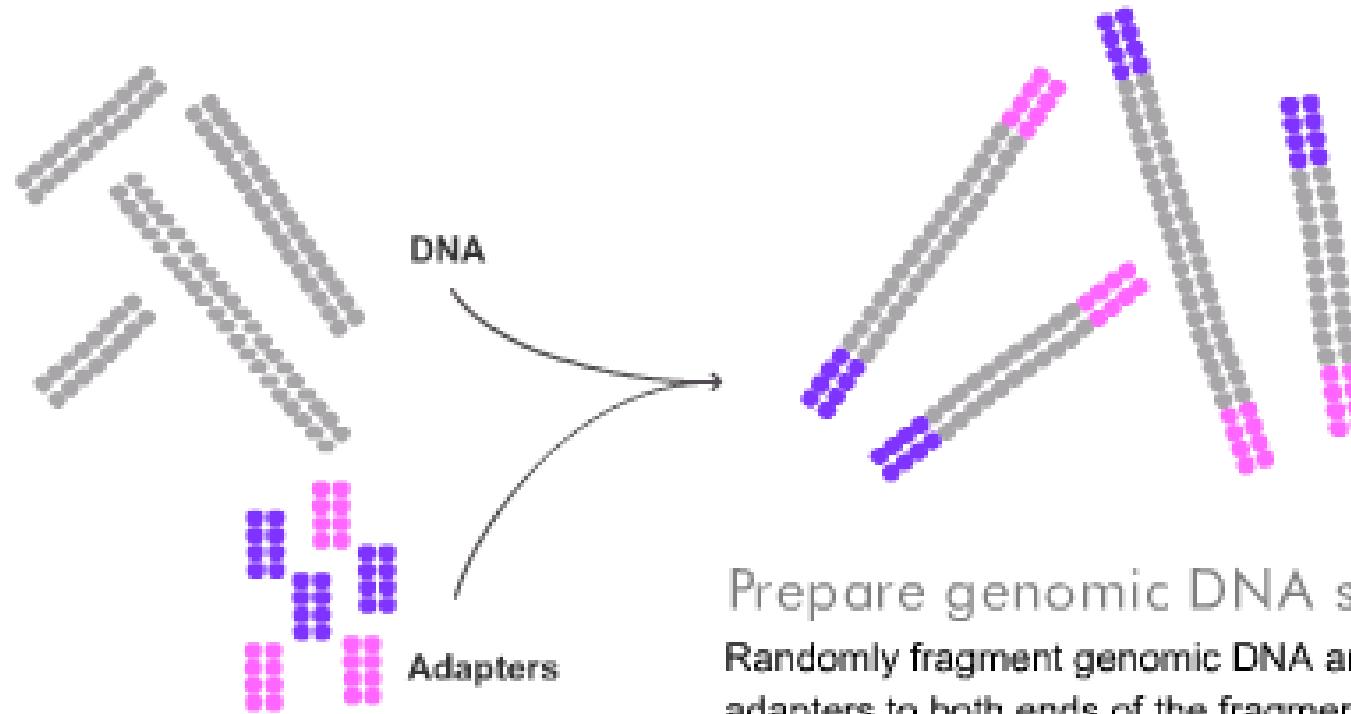
# What is NGS ?



# Illumina sequencing-by-synthesis

## Sequencing-By-Synthesis Demo

1



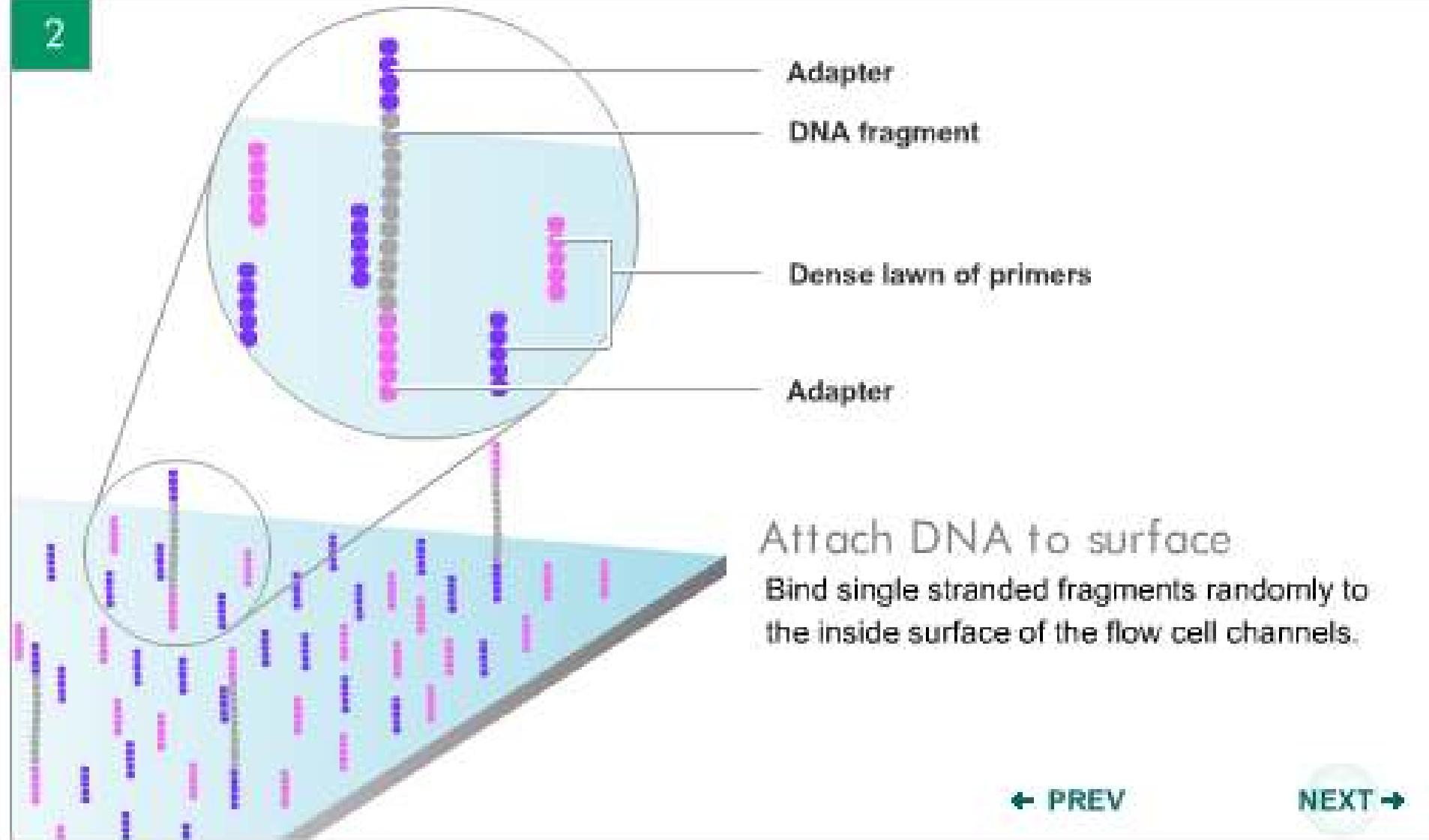
← PREV

NEXT →

# Illumina sequencing-by-synthesis

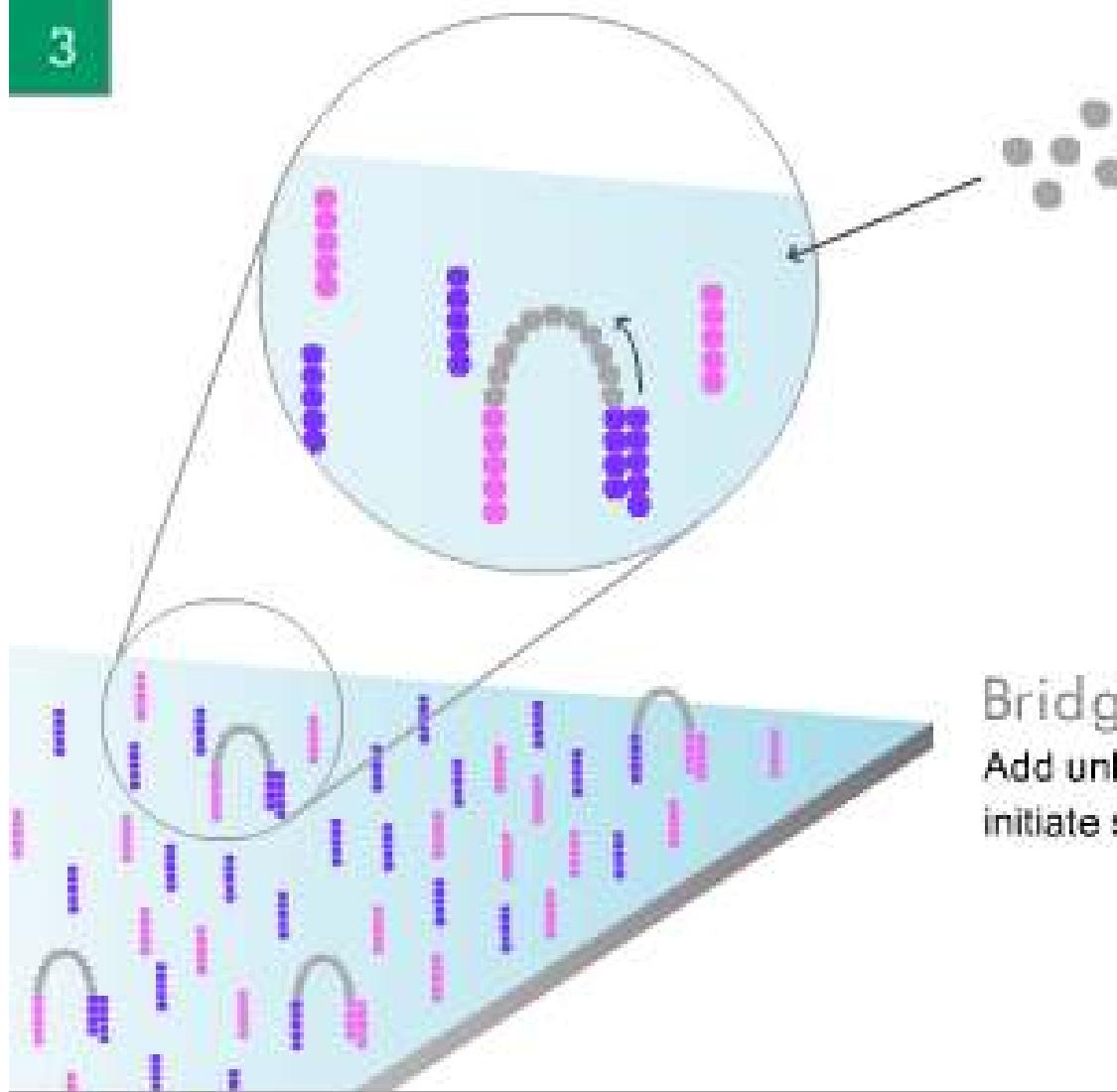
## Sequencing-By-Synthesis Demo

2



# Illumina sequencing-by-synthesis

3



## Bridge amplification

Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

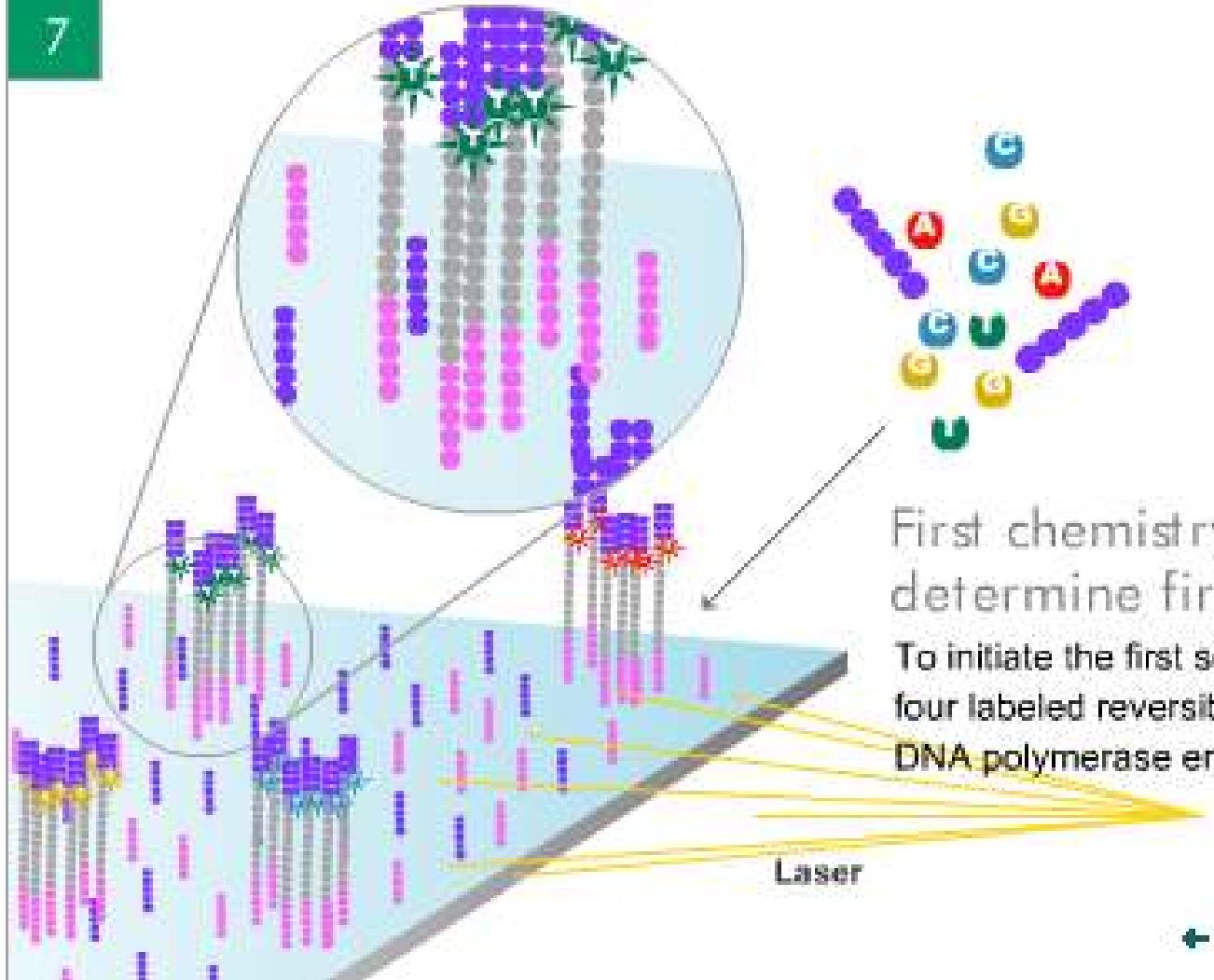
← PREV

NEXT →

# Illumina sequencing-by-synthesis

## Sequencing-By-Synthesis Demo

7



First chemistry cycle:  
determine first base

To initiate the first sequencing cycle, add all  
four labeled reversible terminators, primers and  
DNA polymerase enzyme to the flow cell.

← PREV

NEXT →

# Illumina sequencing-by-synthesis

10

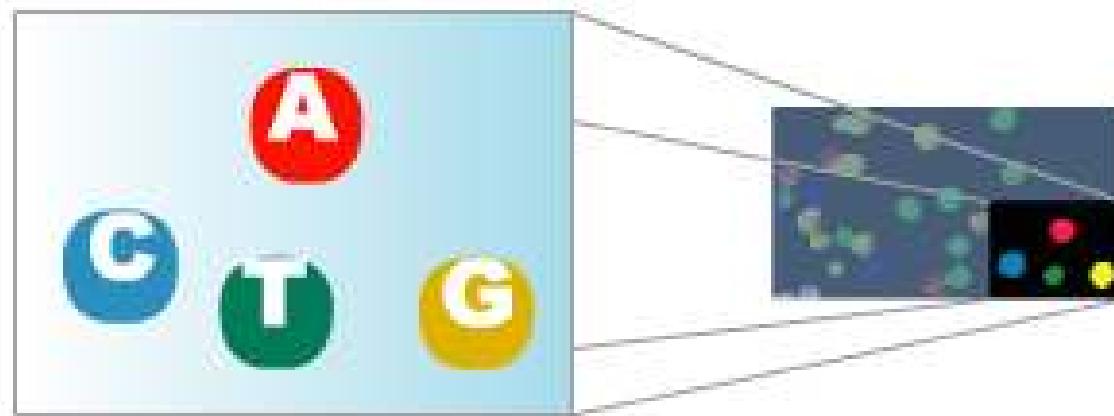


Image of second chemistry cycle  
is captured by the instrument

After laser excitation, collect the image data as  
before. Record the identity of the second base  
for each cluster.

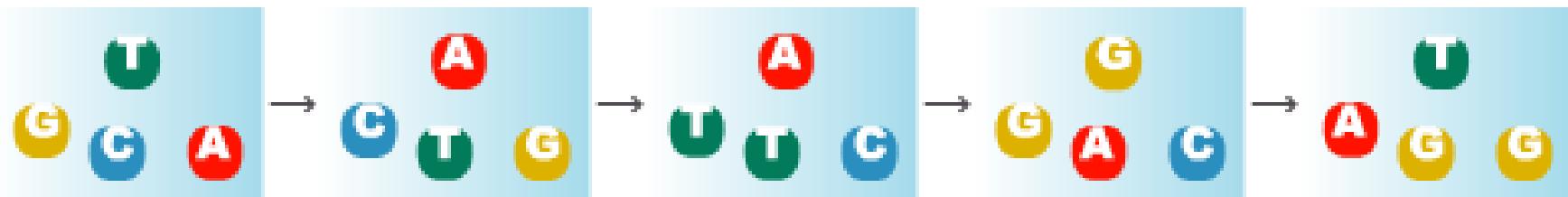
← PREV

NEXT →

# Illumina sequencing-by-synthesis

## Sequencing-By-Synthesis Demo

11



Cycle 1

Cycle 2

Cycle 3

Cycle 4

Cycle 5

**GCTGA....**

Sequence read over multiple chemistry cycles

Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

◀ PREV

NEXT ▶

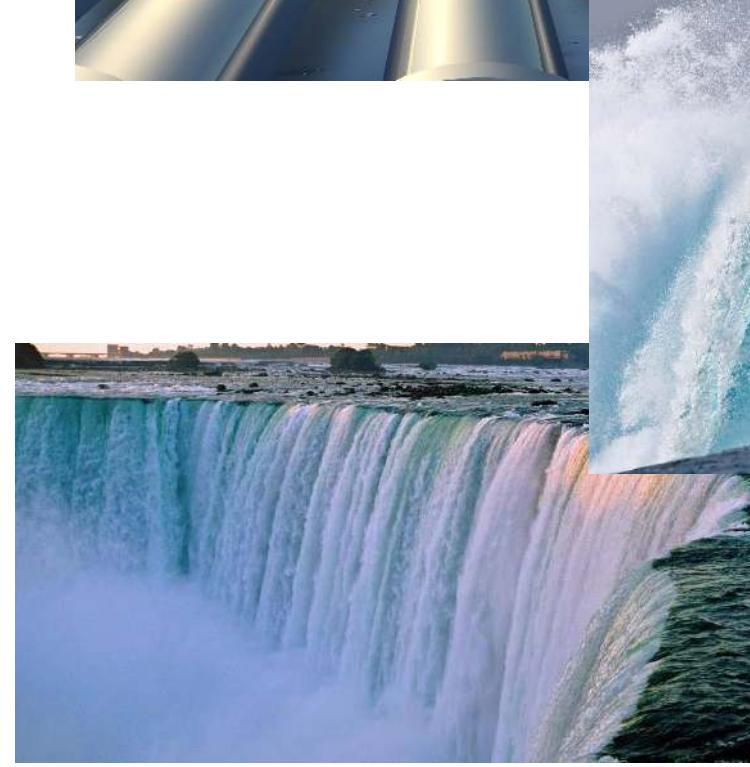
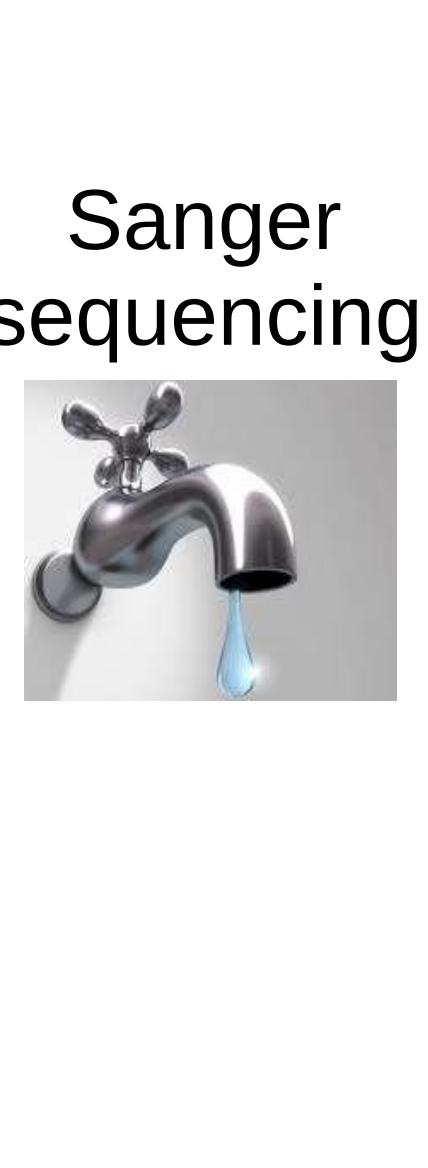
## What is NGS ?

---

= High-throughput sequencing

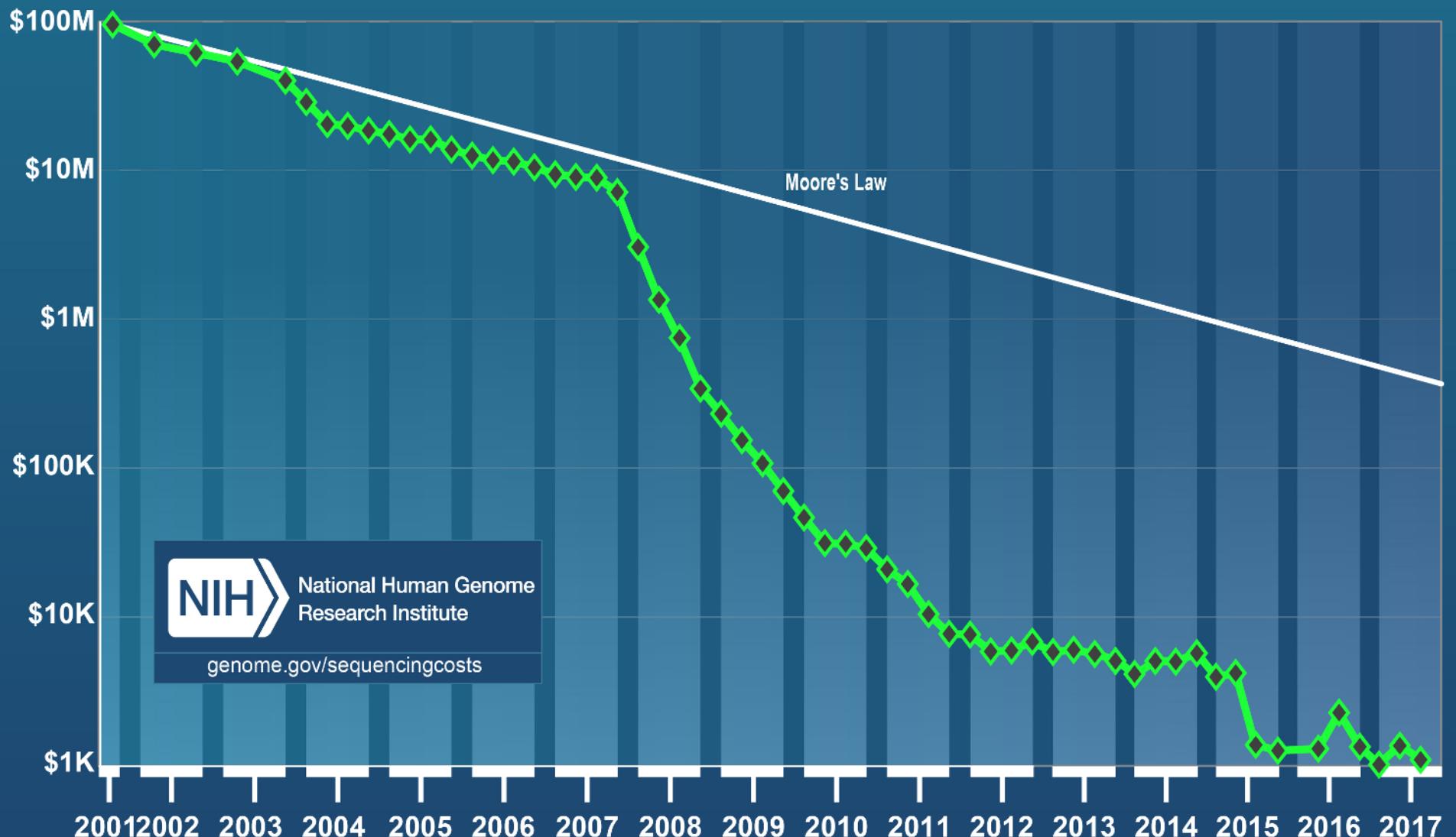


Sanger  
sequencing



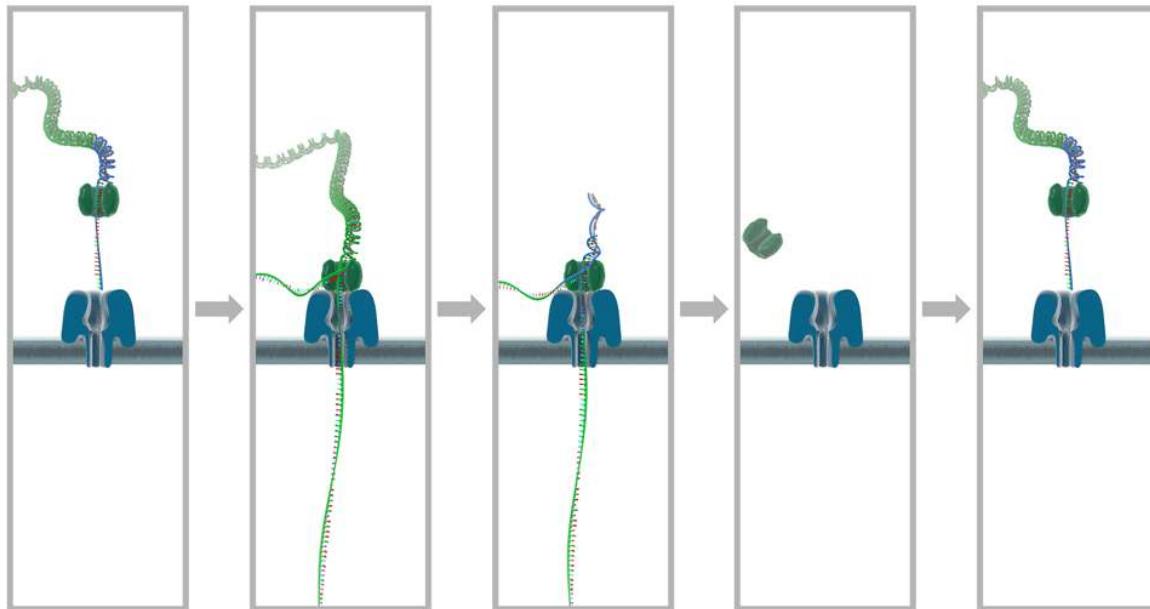
# Sequencing costs

## Cost per Genome

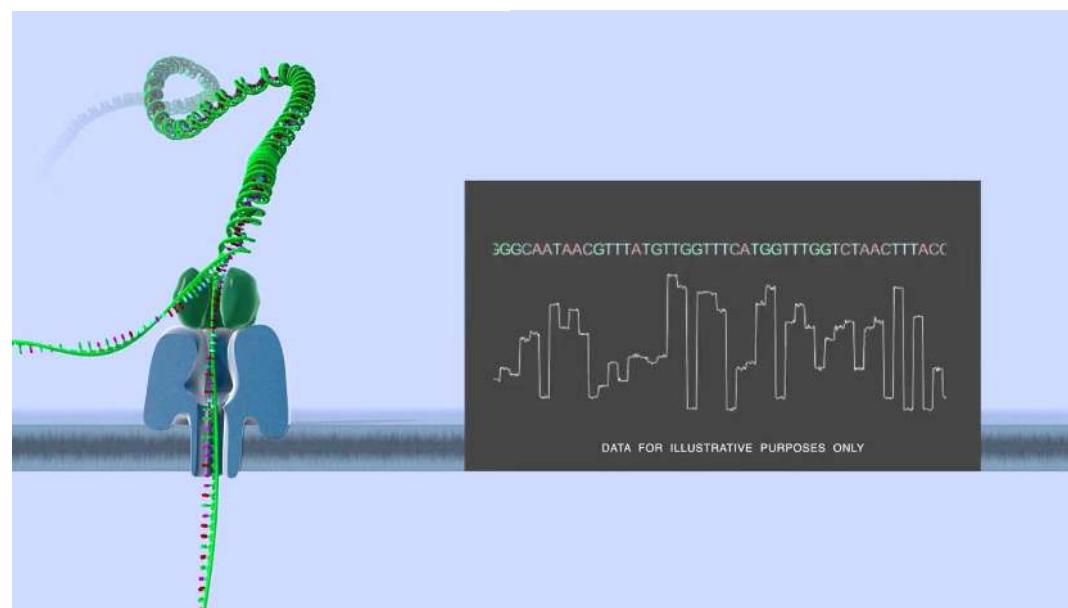


# Nanopore sequencing

→ Long reads are still important (e.g. assembly, phasing, splicing, etc.)



Pollard et al., Hum Mol Genet, 2018



# Sequencing technologies

Platform	Sample preparation	Technology / chemistry	Read length	Yield	Quality	Cost
Illumina	Bridge amplification	synthesis + fluorescence	50-250	++ +++++	+++++	++++
SoliD	emulsion PCR	ligation + fluorescence	50-75	++++	+++	+++
Ion Torrent	emulsion PCR	pyrosequencing + pH	200-400	++	+++	+++
Nanopore	Single molecule	synthesis + fluorescence	DNA-dependent up to 500Kb	+	+	++
Pacific Bioscience	Single molecule	synthesis + fluorescence	30Kb->100 Kb	+++	+	+++

Further reading :

[http://en.wikipedia.org/wiki/DNA\\_sequencing](http://en.wikipedia.org/wiki/DNA_sequencing)

Liu et al., J Biomed Biotechnol, 2012, ID 251364.

McGinn & Gut, N Biotechnol. 2013, 30(4):366-72.

Buermans & den Dunnen, Biochim Biophys Acta. 2014.

Shapiro E. et al. Nat Rev Genet. 2013, 14(9):618-30. Single-cell sequencing

Yang et al., J Nanosci Nanotechnol, 2013, 13(7):4521-38. Nanopore sequencing

# NGS applications

Soon et al., Mol Syst Biol, 2013

Population

Metagenomics

Protein-DNA / RNA

ChIP / RIP / CLIP

RNA

RNA / CAGE / GRO

chromatin accessibility

DNase / FAIRE / ATAC

chromatin structure

MNase / ChIP / Hi-C

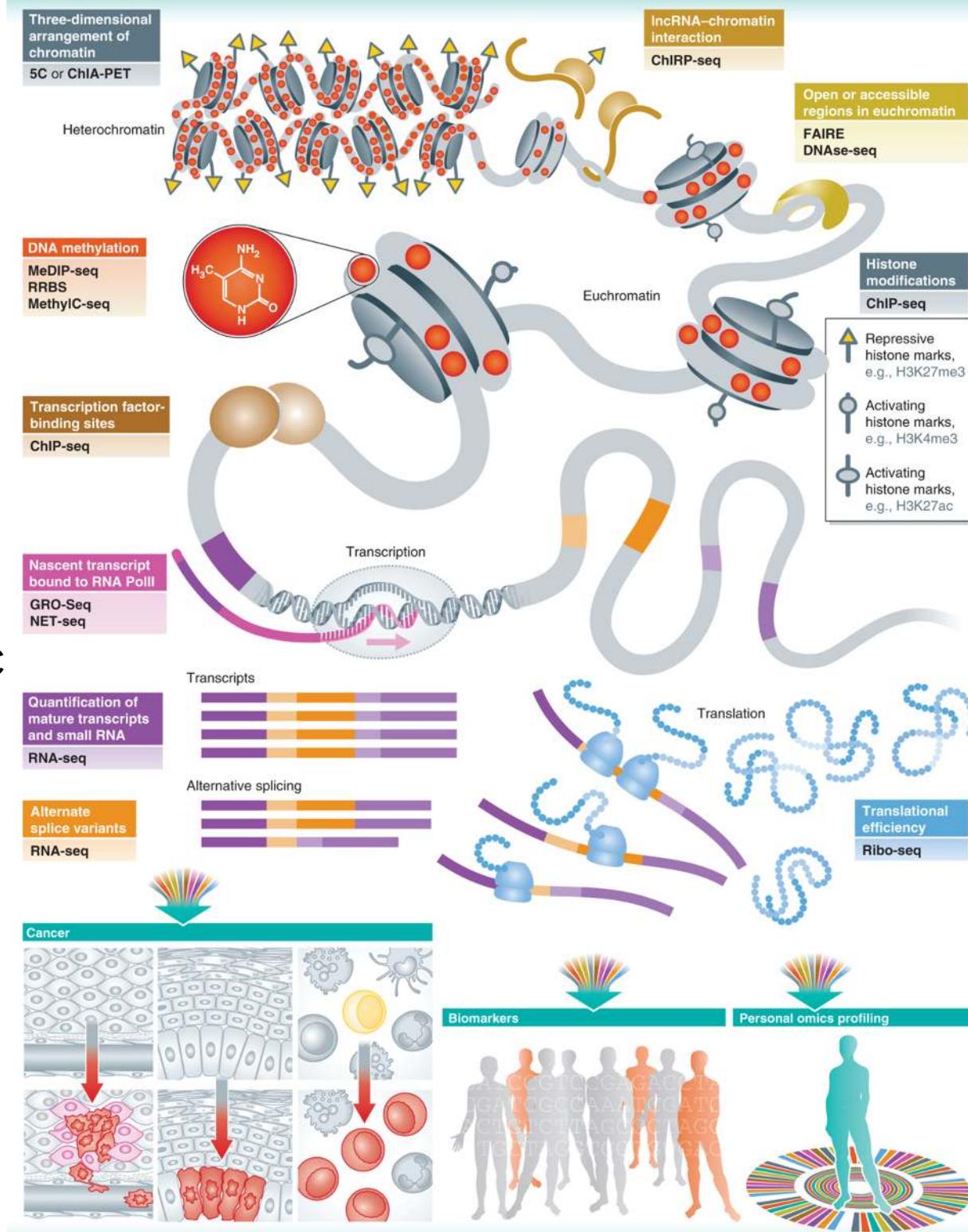
DNA methylation

BS / MeDIP / RRBS

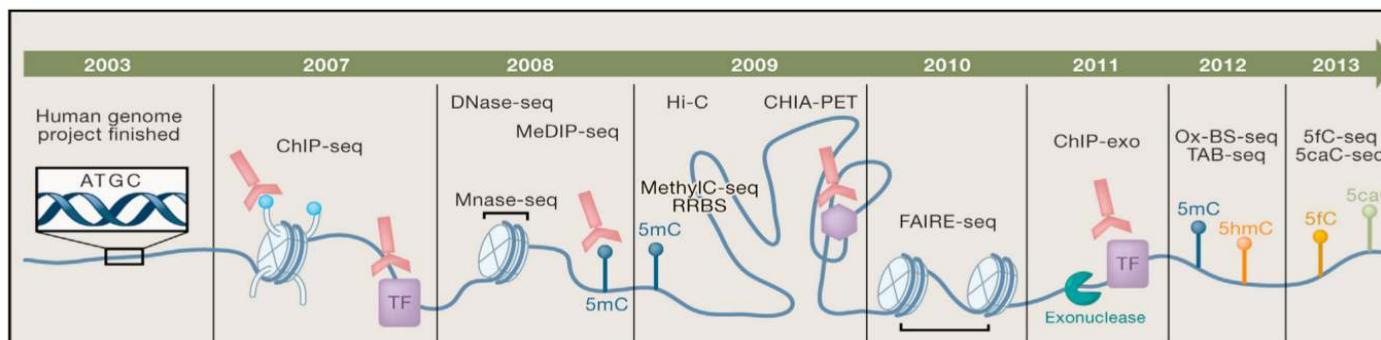
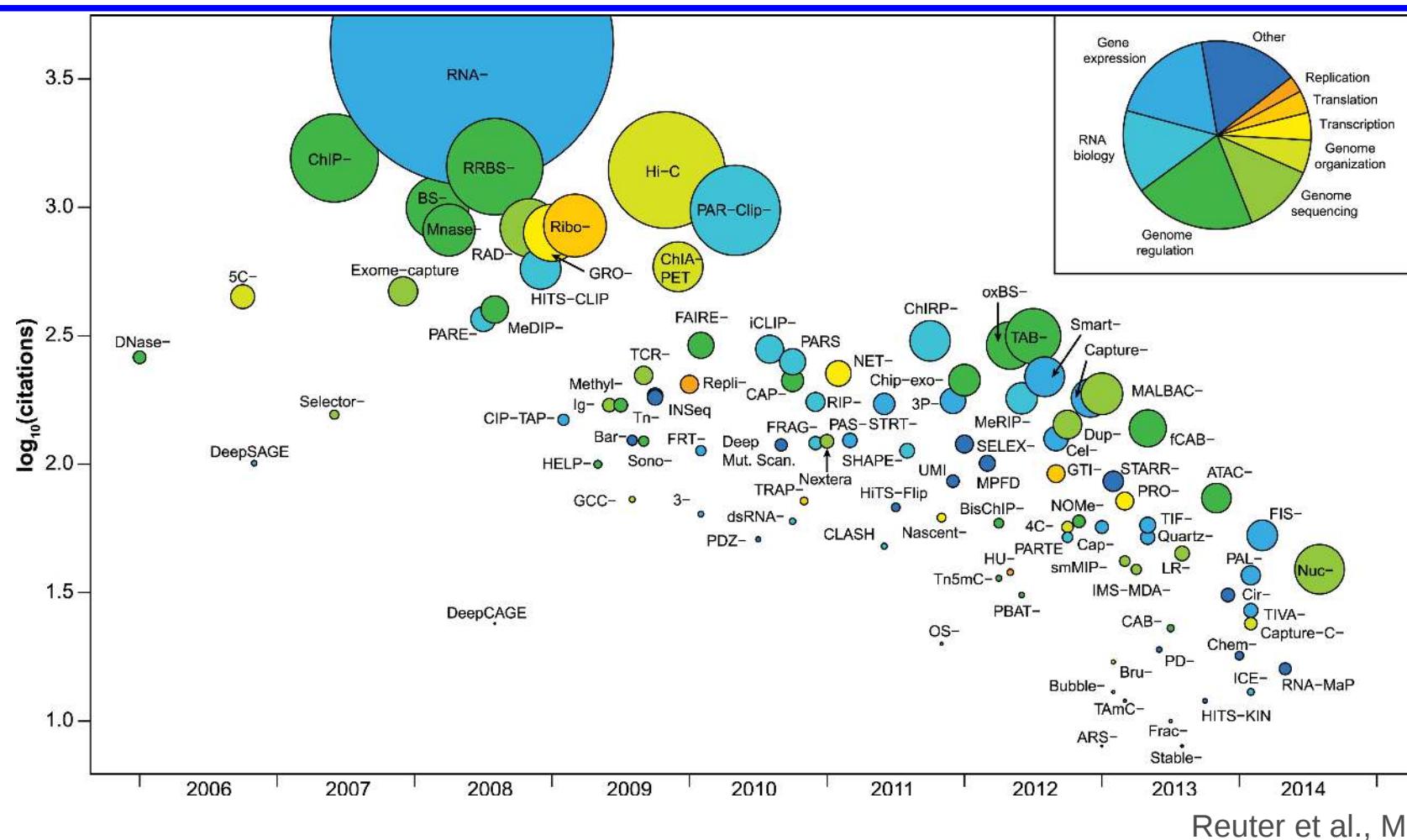
DNA

DNA / exomes

Many others ...



# NGS applications



Rivera & Ren, Cell, 2013

See also : Shendure & Lieberman Aiden, Nat Biotechnol, 2012

# Multiple objectives, often complementary

Analyze structures  
(interactions, secondary structures,...)

Identify variations  
(SNPs, indels, CNV...)

Identify peaks

Measure abundances

Count events / screen

Compare groups (KO/wt...)

Cluster / order  
(states, sequences, individuals,...)

Model  
(systems biology)

etc...

Functional Genomics

Genetics

NGS

Bioinformatics  
Biostatistics

Annotate

Search for motifs  
(TFs, ...)

Find associations  
(QTL, GWAS, ...)

Predict / integrate  
(clinical data, other omics, ...)

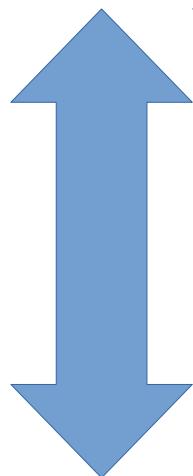
## What do the data look like ?

---

~50-200 bp



TCAGGTGCGCACTCTATGGC..... TAACGGGCTAAAATCAAGTCT  
GACTAACACATGGCGCGTCAG..... ATGCGAACTACGTGCGAGGC  
AGGACGACGGTCGGGAAC..... ATAGCGCGGGCCCATTGATT



Several million reads (large files!!)



TCAGGTGCGCACTCTATGGC..... GCCATGCGTAACATTAAGTCT  
TTAGACCGTAACGGCTATTG..... ATGCGAACTACGTGCGAGGC  
AGGACGACGGTCGGGAAC..... ATAGCGCGGGCCCATTGATT

# Fasta format

Start symbol  
Sequence identifier (no space)  
Description (spaces allowed)

>**NM\_028020 | Mus musculus integrator complex subunit 11, mRNA**

GCGCAGTGC GGAGGGACTGCCGCGGCCGGCGTCTGATGTGAGCGCGCGGGCTGCCTGCGGCTCT  
CCTGCGGCGGTGCACGCGAGCCGTGGAGGCCGTGGCTCGAGGCGACCTCCCTGACCATGCCCGAGATTAG  
GGTCACCTCCCTGGGGCTGCCAGGATGTAGGCCGAAGCTGCATCTGGTCTCCATTTCGGGAAGAAT  
GTCATGTTGGACTGTGGGATGCACATGGCTACAATGATGACAGGGCGTTCCGTGACTTTCTACATCA  
CCCAGAGTGGCCGCCTGACTGACTTCTGGACTGTGTGATCATCAGCCACCTCCACCTGGACCATTGTGG  
GGCACTCCCCTACTTCAGTGAAATGGTGGCTACGATGGACCCATCTATATGACCCATCCTACCCAGGCC  
ATCTGCCCATCCTGTTGAAAGACTACCGCAAGATTGCAGTGGACAAGAAGGGCGAGGCCAATTCTCA  
CTTCTCAGATGATCAAAGACTGTATGAAGAAGGTGGTGGCTGTTCACCTGCATCAGACAGTCAGGTGGA

Sequence

Extension	Alphabet	Utilisation
.fasta, .fa	any	Generic
.fna	DNA	Large genomic regions (contigs, chromosomes, ...)
.ffn	DNA	Coding sequence
.faa	Protein	Amino acid sequence
.frn	DNA	Non coding RNAs (e.g. tRNA, rRNA...)

# Multi-fasta

- Concatenation of fasta files
- The > sign separates sequences

>**Sequence1 tRNA A**

GCGCAGTGC GGAGGGACTGCCGCGGCCGGCGTCTGATGTGAGCGCGCGGGGTCTGCCTGCGGCTCT  
CCTGCGGGCGGTGCACGCGAGCCGTGGAGC

>**Sequence2 protein coding RNA B**

GGTCACTCCCTGGGGCTATGCAGGATGTAGGCCGAAGCTGCATCTGGTCTCCATTGGGAAGAAT  
GTCATGTTGGACTGTGGGATGCACATGGCTACAATGATGACAGGCGCTTCCTGACTTTCTTACATCA  
CCCAGAGTGGCCGCCTGACTGACTTTCTGGACTGTGTGATCATCAGCCACTCCACCTGGACCATTGTGG  
GGCACTCCCCTACTCAGTGAATGGTGGCTACGATGGACCCATCTATATGACCCATCCTACCCAGGCC  
ATCTGCCCATCCTGTTGGAAGACTACCGCAAGATTGCAGTGGACAAGAAGGGCGAGGCCATTCTTCA

>**Sequence3 non coding RNA C**

CTTCTCAGATGATCAAAGACTGTATGAAGAACGGTGGCTGTTCACCTGCATCAGACAGTCAGGTGGA  
TGATGAGCTGGAAATCAAGGCATACTATGCAGGCCACGTGCTGGGGCAGCCATGTTCCAGATTAAAGTG  
GGCTCGGAGTCTGTGGCTACACGGGTGACTATAACATGACCCCCAGACCGGCATTGGGGCTGCGTGG  
TTGACAAGTGTGCCCCAACTTGCTCATCACAGAACATCCACATATGCTACAACCATTGGAGATTCAAACG  
CTGCAGAGAGCGAGATTCTAAAGAAAGTTCATGAAACTGTGGAGCGTGGAGGAAAGGTGCTGATTCT  
GTGTTGCACTGGGCCGAGCACAAGAACTCTGCATCCTGCTGGAGACCTCTGGAGCGCATGAACCTGA  
AGGTGCCCATATACTCTACGGGCCTGACAGAGAAAGCCAACCAACTACTACAAGCTCTCATCACCTG  
GACCAACCAGAAGATCCGGAAGACATTGTCCAGAGGAACATGTTGAGTTAACGCACATCAAAGCCTT

# Phred quality score

- A quality score for each nucleotide
- Measures the reliability of the base call (physical process: chromatogram, image, pH...)
- Formalized in the Phred software developed for the Human Genome Project

Ewing et al., Genome Res. 1998. 8(3):175-85 & 186-94

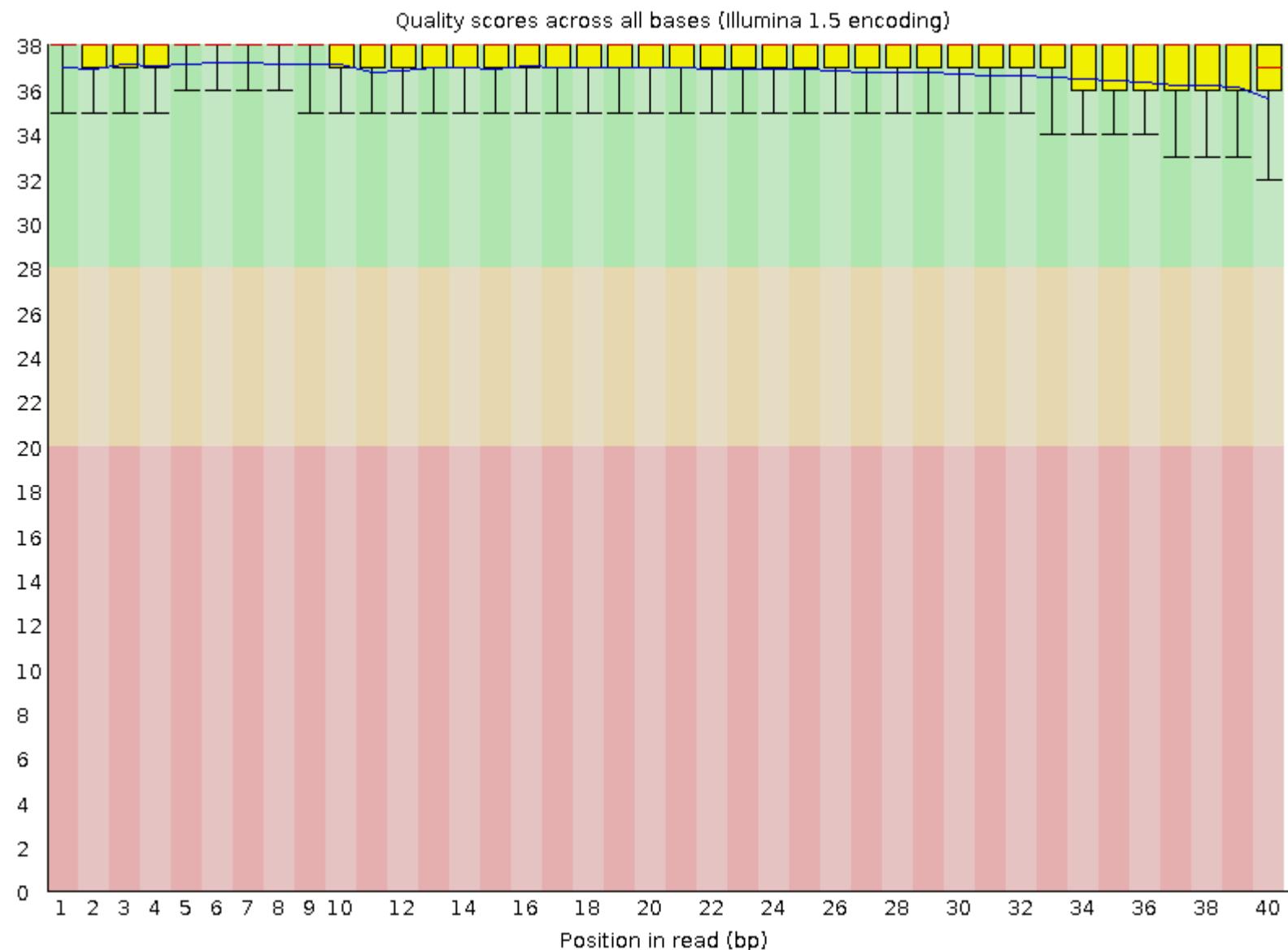
Phred quality score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1 000	99,9 %
40	1 in 10 000	99,99 %
50	1 in 100 000	99,999 %

$$Q = -10 \times \log_{10}(P) \leftrightarrow P = 10^{-Q/10}$$

- Q : Phred quality score
- P : Probability of base calling error

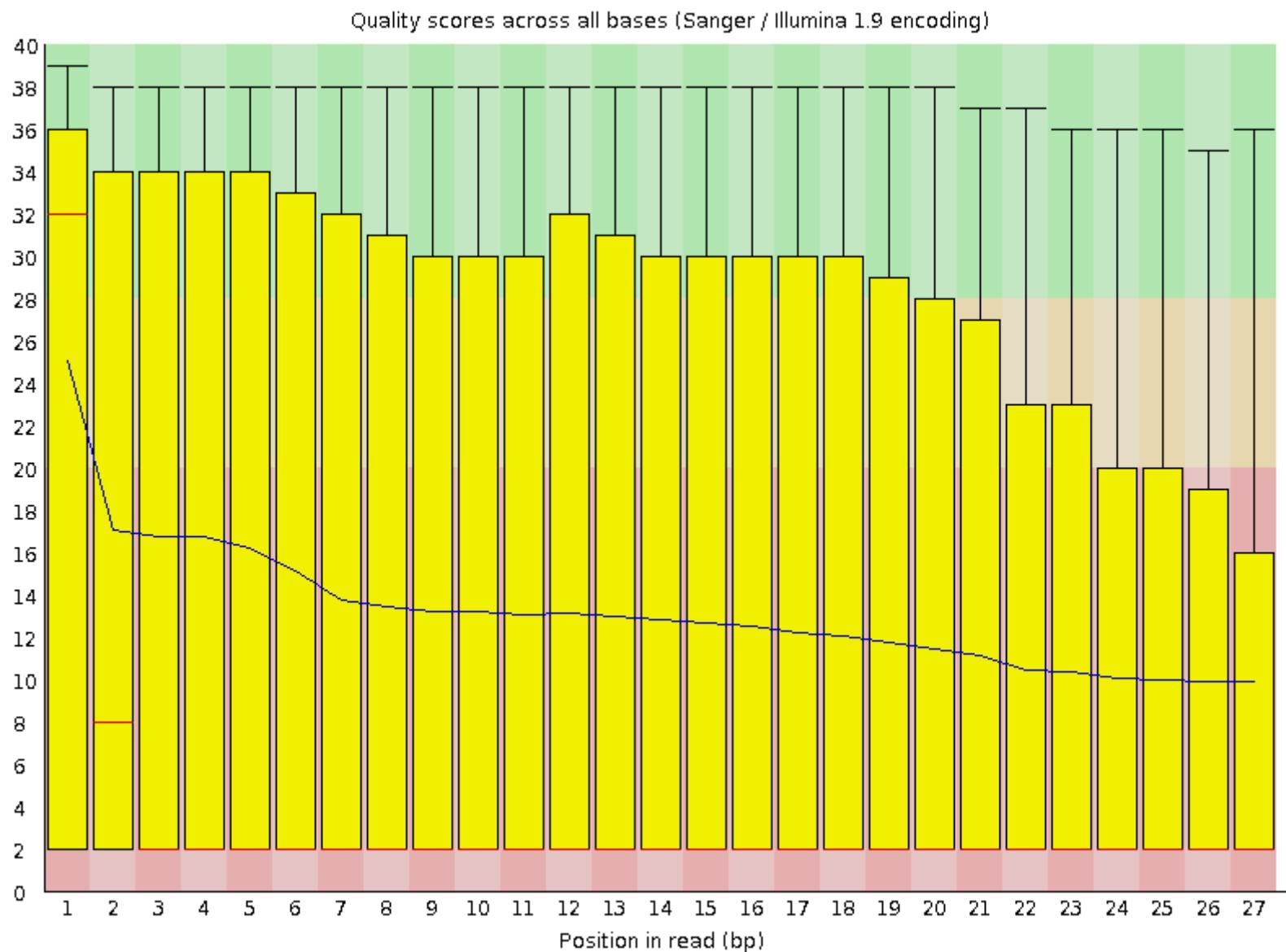
# Quality

➤ Exemple of high quality reads (Illumina)



# Quality

➤ Low quality reads (CAGE-seq, Illumina GAIIX)

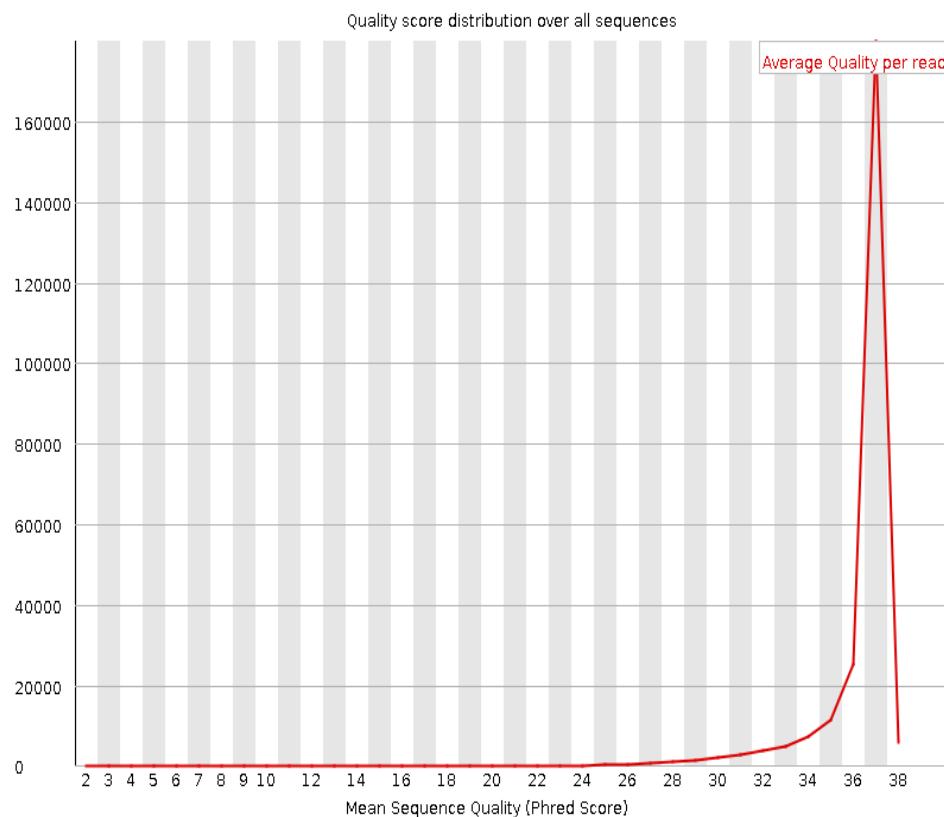


# Read quality

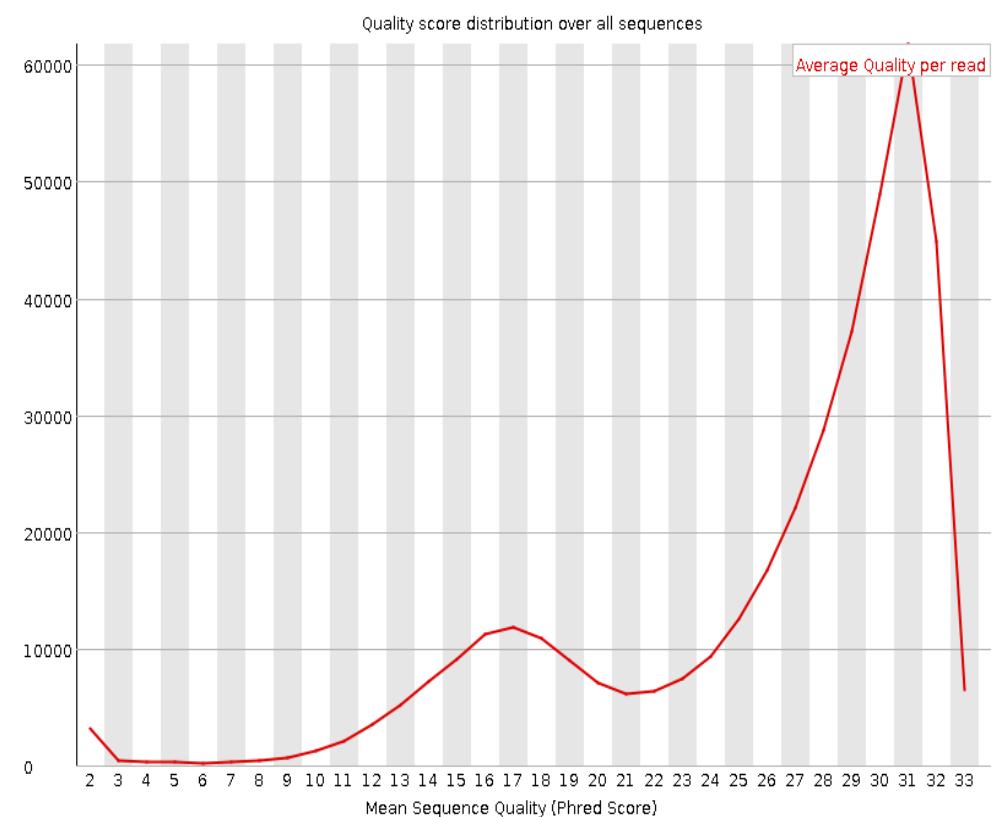
- Average phred score used to evaluate quality of entire read
- Different from alignment quality !!

Li et al., Genome Res. 2008. 18(11):1851-8. MAQ

Excellent quality (Illumina)



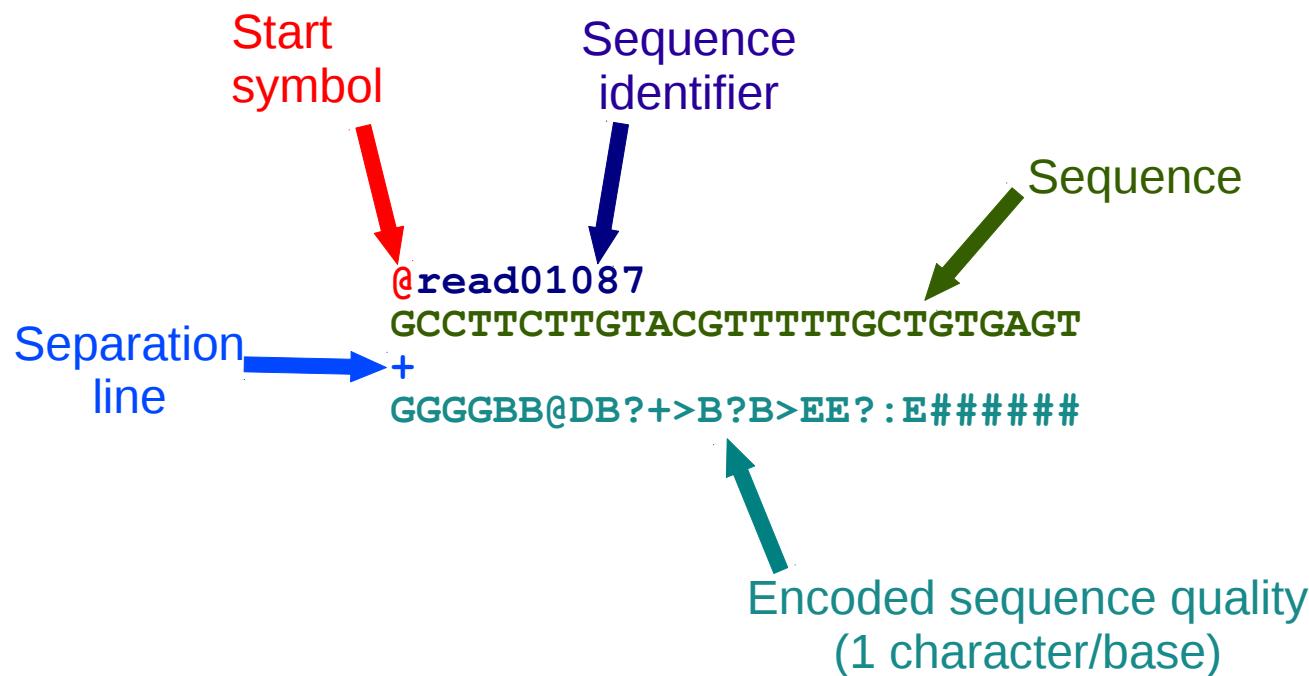
Low quality (Illumina)



# Fastq format

[http://en.wikipedia.org/wiki/FASTQ\\_format](http://en.wikipedia.org/wiki/FASTQ_format)

File extension: .fq or .fastq



## Sequence quality encoding (**ASCII**, Phred+33):

A horizontal scale representing a range from Q0 to Q40. The scale is marked with vertical tick marks at intervals of 10 units. Below the scale, five points are labeled with color-coded descriptions: 'bad' (red) at Q0, 'average' (orange) at Q10, 'ok' (yellow) at Q20, 'good' (green) at Q30, and 'excellent' (blue) at Q40.

# Multi-fastq

## ➤ Simple concatenation

```
@SRR488286.8911073
GCCTTCTTGTACGTTTGCTGTGAGT
+
GGGGBB@DB?+>B?B>EE?:E######
@SRR488286.10295868
GCTTCGTTCTCAAACCGTCGNCCAGA
+
G>GGGG:>?A;GGGGGB?;BB(BB??2
@SRR488286.2166439
GAGTTGGGCTTCCACCTCGACCGGGAA
+
C@<@B0:<'87;B0;B;;52;;06=:@
@SRR488286.8412402
GAGTAAAACATACAACTTAACGAGAA
+
A??;*A:>@.<@(@##########

```

- Often very large files (>10Go)
- Contain redundancy
- Compression (.gz, .bz2) => size often reduced to <20 % of orginal size
- Saves disk space and accelerate transfers but reading/writing is slower

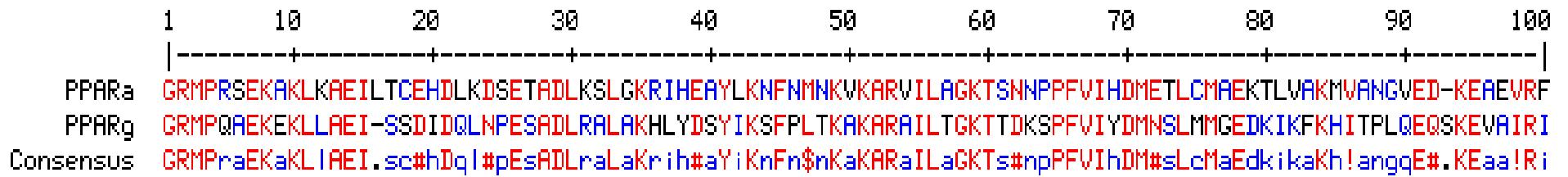
# Alignment

- We have :
  - Millions of reads in fastq file
  - A reference sequence in a fasta file

→ **Where do the reads come from?  
Alignment problem**

## What is an alignment ?

Sequence alignment on Wikipedia



- "Optimal" alignment / adjustment of 2 sequences.
- Match/mismatch ? gaps/extensions ? What criterion defines an optimal alignment ?

# Alignment : edit distance

The **edit distance** (Levenshtein, 1965) is the minimum number of characters that need to be replaced (mismatch), inserted or suppressed (gap/extension) to go from one character string to the other.

- Insertion (i)
- Deletion (d)
- Substitution (s)

TCGCGTA—CGTACG  
—CGC**C**TACCG**G**A—G  
d s i s d

Edit distance between query (q) and subject (s) :  $D(q,s)=5$

- Calculating the edit distance is not trivial.  
 **Dynamic programming**
- The algorithm provides the optimal alignment between 2 sequences

# Pairwise and multiple alignment

## ➤ Pairwise alignment (2 séquences) : Blast (n, p,...), blast2seq, Needle, Water...

```
query      1  GRMPRSEKAKLKAELITCEHDLKDSETADLKSLGKRIHEAYLKNFNMNKVKARVILAGKTSNNPPFVIHDMETLCMAEKTLVAKMVANGVED-KEAEVRF  99
subject    1  GRMPQAEKEKLLAEI-SSDIDQLNPESADLRALAKHLYDSYIKSFPLTKAKARAILTGKTTDKSPFVIYDMNSLMMGEDKIKFKHITPLQEQSKEVAIRI  99
Consensus   GRMP++EK KL AEI + + D + E+ADL++L K +++++Y+K+F + K KAR IL GKT++ PFVI+DM +L M E + K + E KE +R
```

- A **query** sequence vs a **subject** sequence
- Blast aligns 1 query vs thousands of subjects (database)
- provides an e-value (or E score) the number of hits one can "expect" to see by chance when searching a database of a particular size (background noise). See [BLAST tutorial](#).
- NGS aligner align millions of queries to a few references

## ➤ Multiple alignment (>2 sequences) : Multalin, ClustalW, Muscle, T-coffee ...

```
1      10     20     30     40     50     60     70     80     90     100
|-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
PPARa GRMPRSEKAKLKAELITCEHDLKDSETADLKSLGKRIHEAYLKNFNMNKVKARVILAGKTSNNPPFVIHDMETLCMAEKTLVAKMVANGVED-KEAEVRF
PPARb GRMPQAEKRKLVAGLTASEGCQHNQQLADLKAFSKHIYNAYLKNFMNTKKKARSILTGKSSHNAFPVIHDIEWLQAEKGLWKQLVNGLPPYNEISVHV
PPARg GRMPQAEKEKLLAEI-SSDIDQLNPESADLRALAKHLYDSYIKSFPLTKAKARAILTGKTTDKSPFVIYDMNSLMMGEDKIKFKHITPLQEQSKEVAIRI
Consensus GRMP.aEK.KL.Rei..s#.dq.#p#.ADLkal.Khiy#aYIKnFn$tK.KAR.IltGKts.n.PFVIhDm+tL.maEk. lv.K...ng.e..!r.
```

# Global vs Local alignment

➤ Global alignment:

FTFTALILLAVAV
F---TAL-LLA-AV

- Alignment on the whole length of the sequences
- Used for sequences of similar length
- Needleman-Wunsh algorithm

➤ Local alignment:

FT <b>FTALILL</b> -AVAV
<b>FTAL-LL</b>

- Only keeps the parts of the 2 sequences that match
- Well suited to search conserved protein domains
- Smith-Watermann algorithm

➤ Glocal or semi-local alignment:

FT <b>FTALILL-AV</b> AV	subject
<b>FTAL-LLAAV</b>	query

- global for the **query** and local for the **subject**

# Alignment : seed-and-extend

Reference sequence :

... ACTGGGTACATCGTACGATCGATCGATCGATCGATCGATCGGCTAGCTA ...

Short read:

GTCATCGTACGATCGATA**A**GATCGATCGATCGGCTA

 mismatch  
↓  
Slicing (11bp fragments)

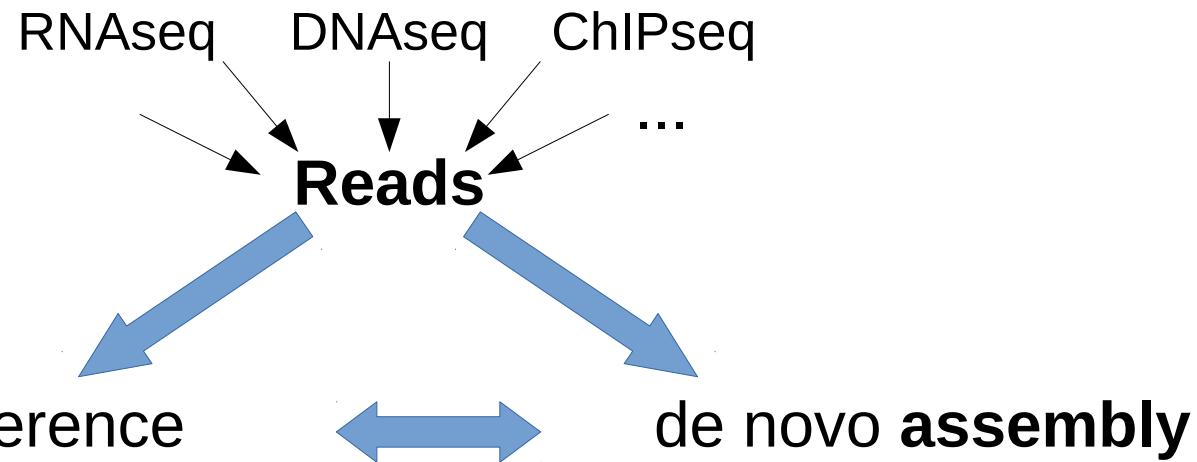
GTCATCGTACG      ATCGAT**A**GATCG      ATCGATCGGCTA  
Seed                  ↓  
Alignment of short fragments  
(search for 'anchor' or seed)

... ACTGG**GTCATCGTACG**ATCGATCGATCGATCGATCGGCTAGCTA ...  
↓  
Extension (Smith-Waterman)

  
GTCATCGTACGATCGAT**A**GATCGATCGATCGGCTA  
... ACTGG**GTCATCGTACG**ATCGATCGATCGATCGATCGGCTA**GCTAGCTA** ...

# NGS alignment

- Data **volume** imposes high-performance tools  
(computation, storage, algorithms, code)
- The reads have **subtle differences** compared to the reference genome  
(reference is assembled from multiple sequences, SNP, evolution, sequencing errors...)
  - ➡ Can't use exact string matching
- The genome contains a lot of **sequence repeats**
  - ➡ Some reads can align at multiple positions (up to >50 % sometimes...)
- Some reads do not align to the reference genome (primers, **contaminations**...)
- Often complementary strategies can be used :



# NGS alignment - tools

Li & Homer, Brief Bioinform, 2010 ; Fonseca et al., Bioinformatics, 2012

## ➤ Algorithms based hash tables

- ELAND
- BFAST
- MAQ (Li et al., Genome Res, 2008)
- SHRIMP, ...

→ An « index », obtained on reads and/or the reference

→ Mapping quality scores.

(ex : Ruffalo et al., Bioinformatics, 2012)

Phred scores  
Mismatchs/indels  
Repeats  
Paired reads

→ Use 'Seed-and-extend' concept

## ➤ Algorithms based on Burrows-Wheeler transform

→ reversible recoding of sequences that optimises their storage

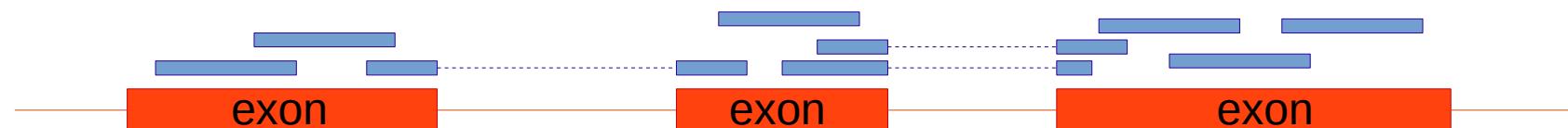
- Bowtie (Langmead et al., Genome Biol., 2009) / Bowtie2 (Langmead & Slazberg, Nat Methods, 2012)
- BWA (Li & Durbin, Bioinformatics, 2009 ; Li & Durbin, Bioinformatics, 2010)

## ➤ Spliced reads aligners (RNA-seq)

- HISAT2 (replaces TopHat2 Kim et al., Genome Biol, 2013, Pertea et al. Nat Protoc 2016)
- STAR (Dobin et al., 2013)

reads :

ADNg :



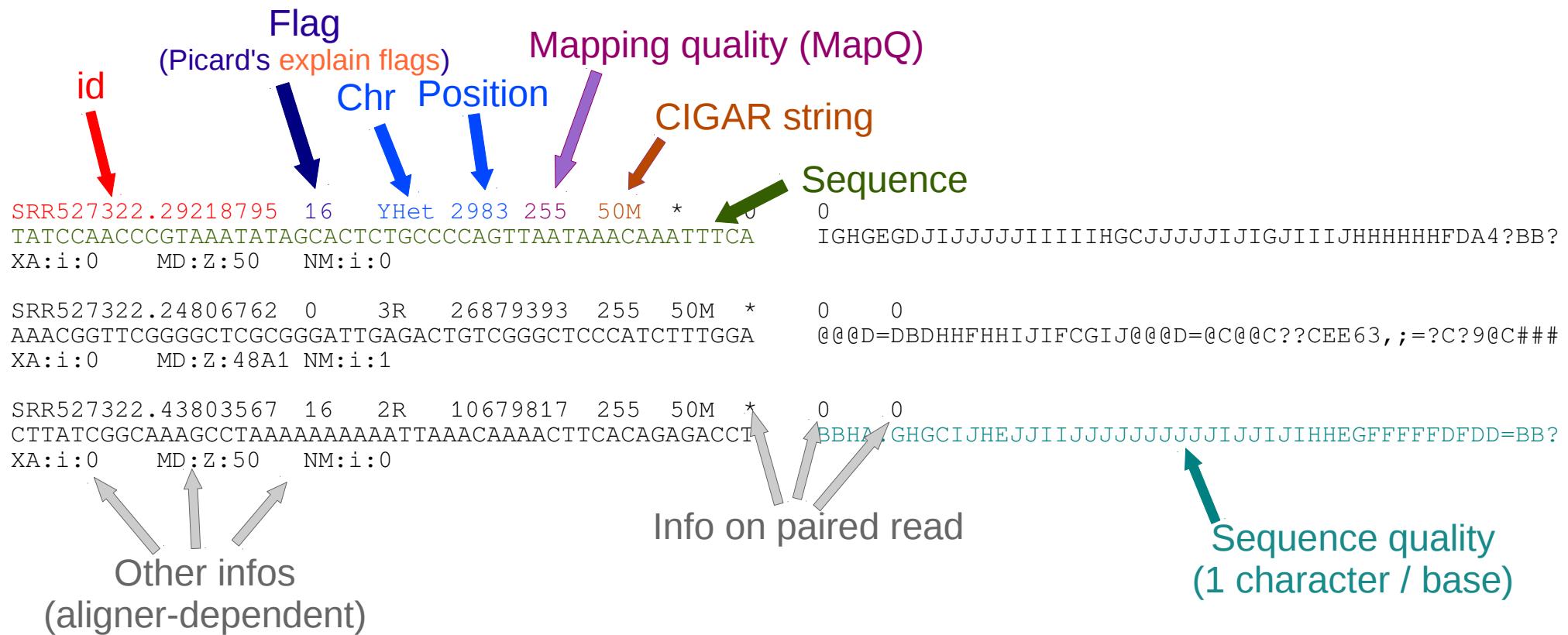
See also other alignment softwares on Wikipedia

# SAM / BAM file format

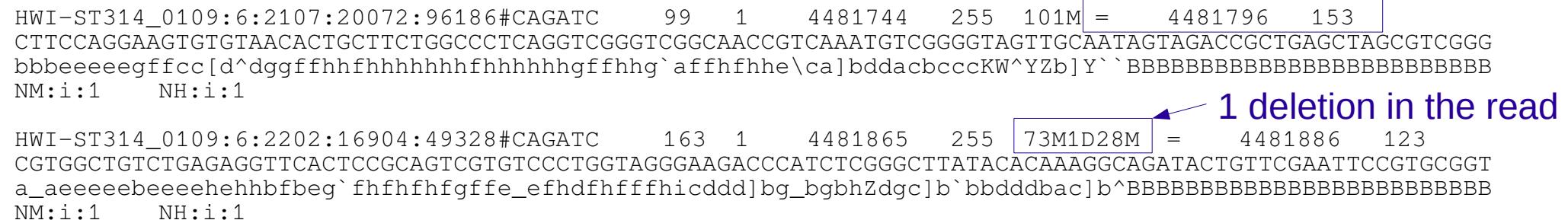
<http://samtools.github.io/>

Li et al., Bioinformatics, 2009

See also: SAMtools, Picard, genome.sph.umich.edu/wiki/SAM

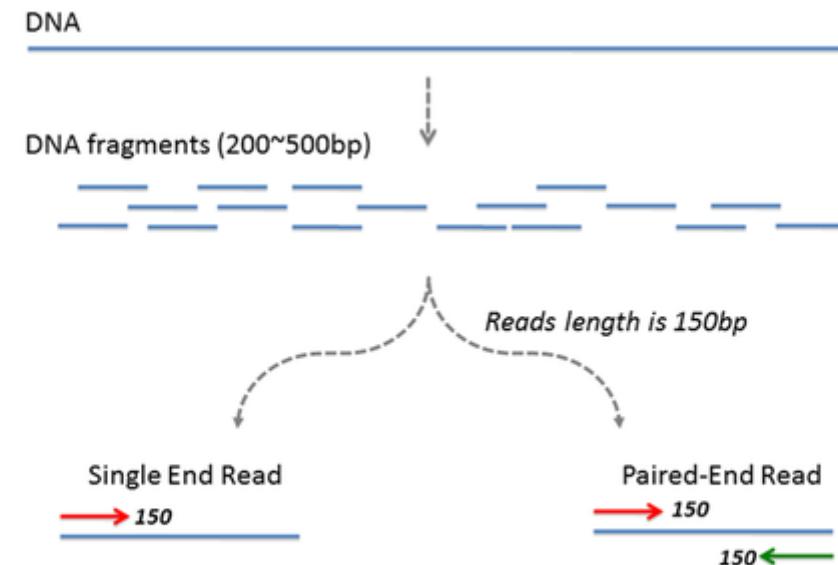


## ➤ Paired-end RNA-seq :



# SAM / BAM file format

## ➤ Paired-end reads (RNA-seq) :



**id**

**Flag**

**Chr**

**Position**

**MapQ**

**position paired read**

**fragment length**

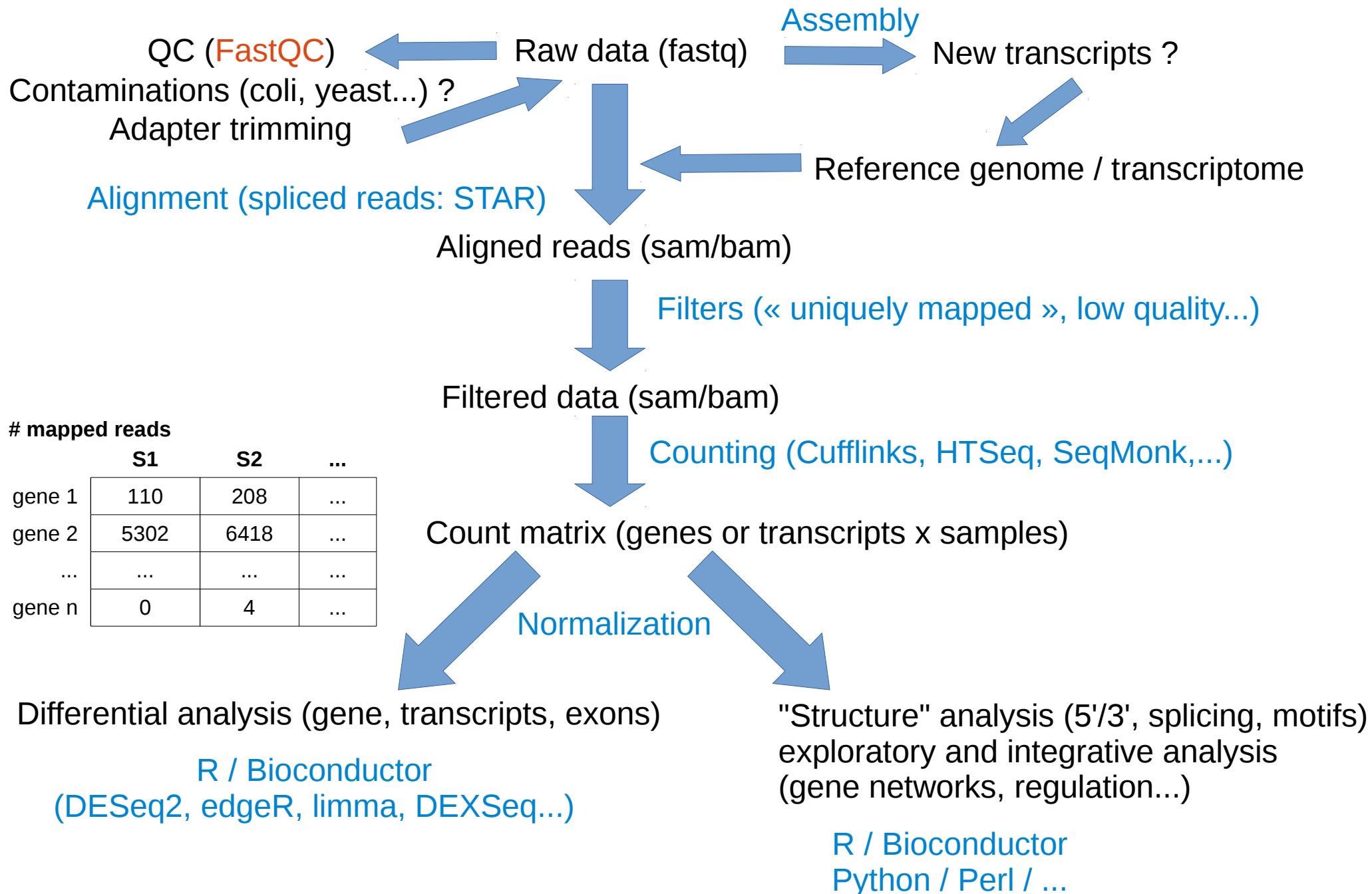
HWI-ST314\_0109:6:2107:20072:96186#CAGATC 99 1 4481744 255 101M = 4481796 153  
CTTCCAGGAAGTGTGTAACACTGCTTCTGGCCCTCAGGTGGGTGGCAACCGTCAAATGTCGGGGTAGTTGCAATAGTAGACCGCTGAGCTAGCGTCGGG  
bbbeeeeegffcc[d^dggffhhfhhhhhhfhhhhhgffhhg`afffhfhhe\ca]bddacbccckW^YZb]Y`BBBBBBBBBBBBBBBBBBBBBBBBBB  
NM:i:1 NH:i:1

HWI-ST314\_0109:6:2202:16904:49328#CAGATC 163 1 4481865 255 73M1D28M = 4481886 123  
CGTGGCTGTCTGAGAGGTTCACTCCGCAGTCGTGCCCTGGTAGGAAAGACCCATCTCGGGCTTATACACAAAGGCAGATACTGTTCGAATTCCGTGCGGT  
a\_aeeeeebeeeehehbfbeg`fhfhfhgf fe\_efhdfffffhicddd]bg\_bgbhZdgcb`bbdddbac]`^BBBBBBBBBBBBBBBBBBBBBBBBBB  
NM:i:1 NH:i:1

**1 deletion in the read**

**CIGAR string**

# Typical RNA-seq pipeline



## Advantages:

- Improved sensitivity and dynamic range
- Better reproducibility (**SEQC**)
- Not limited to prior knowledge of sequence and annotations
- Single base resolution of transcribed features (novel transcripts, splicing, allele-specific expression)
- Can detect quantitative as well as qualitative/structural changes to the transcriptome

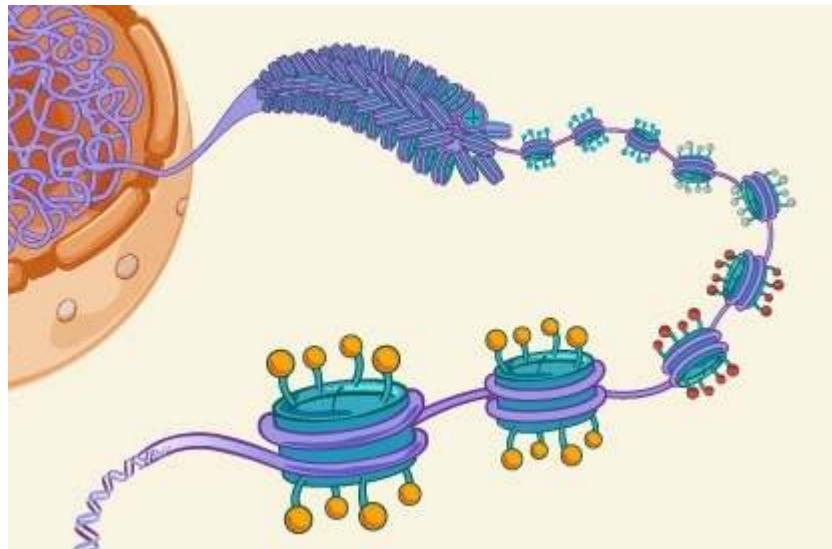
## Drawbacks / limits:

- Some GC content bias
- Mapping ambiguity for paralogous sequences
- Some protocols are highly sensitive to RNA degradation.
- Measurements for different genes are not independent from each other
- Higher cost per sample
- Data analysis is non trivial

# Epigenomics

# Epigenome

---



- DNA methylation
- Histone modifications
- Chromatin accessibility
- Chromatin conformation
- ncRNA

## International initiatives / consortium

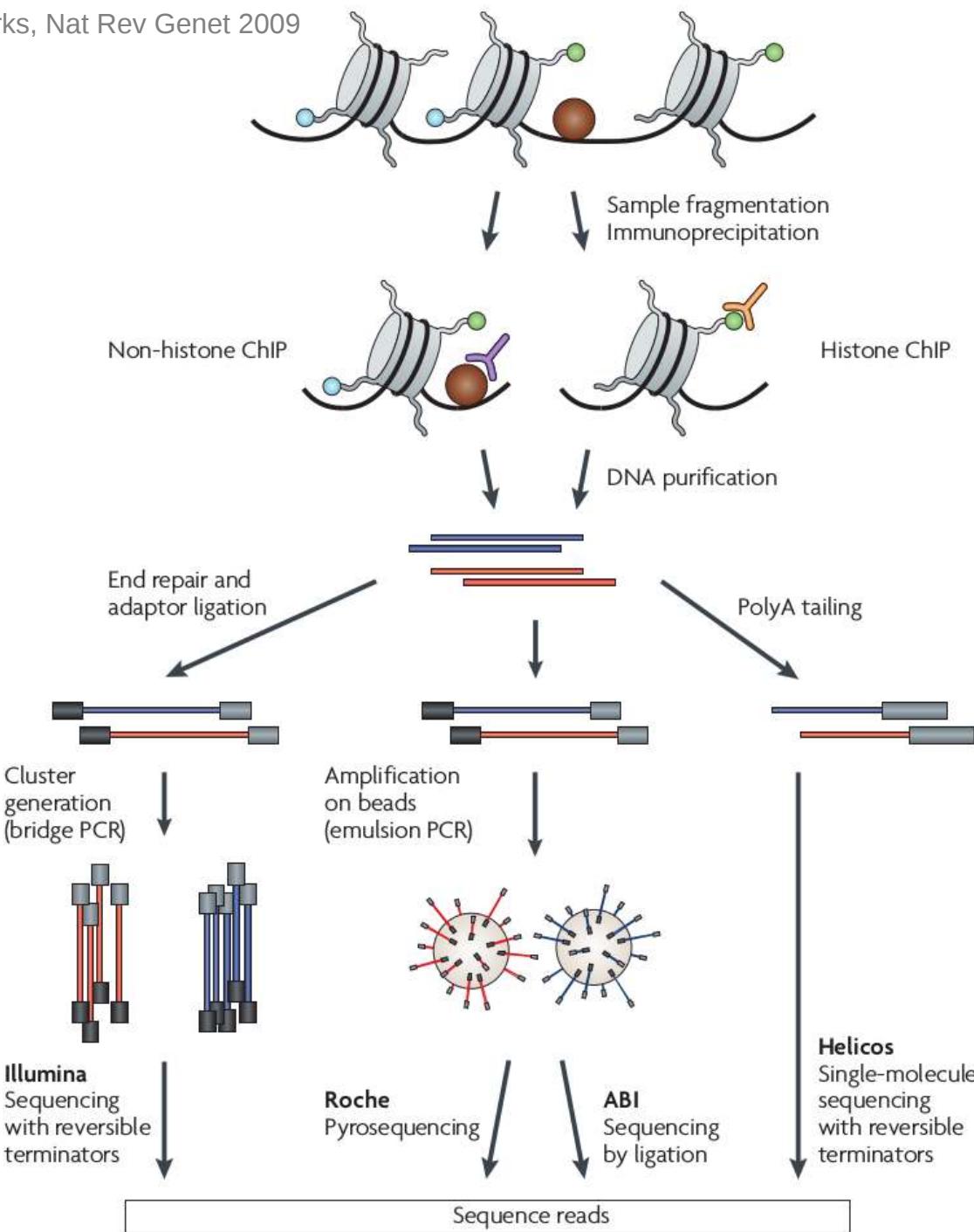
International Human Epigenome Consortium

- ENCODE and modENCODE
- NIH Roadmap Epigenomics
- Blueprint epigenome (Haematopoietic epigenomes)
- 4D Nucleome

...

# ChIP-seq

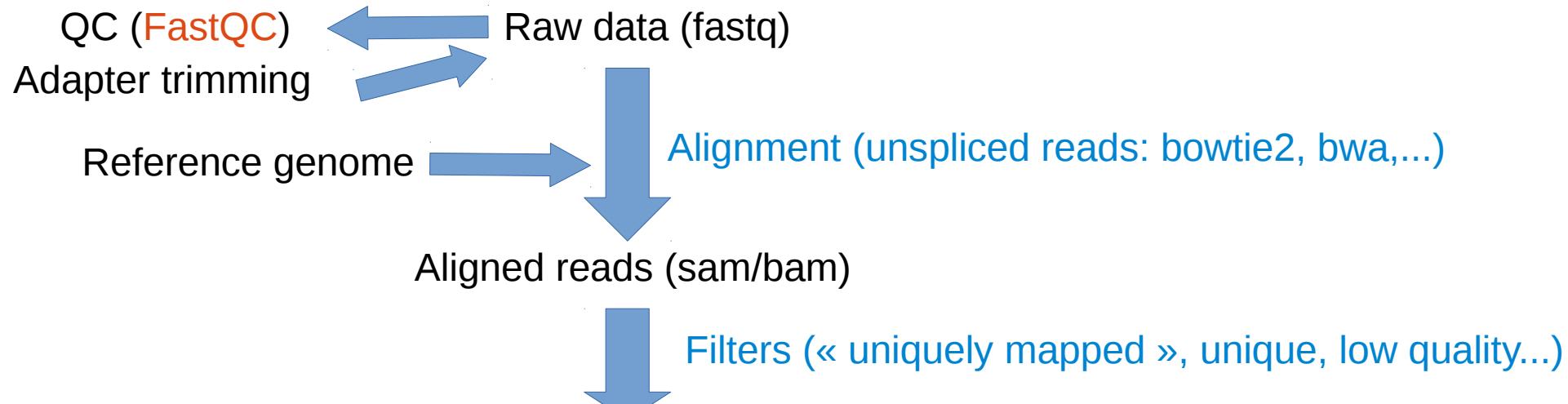
Parks, Nat Rev Genet 2009



- Mapping protein-DNA interactions genome-wide
- Transcription factors, RNA polymerase, histones, insulators, ...
- Numerous protocol variants
- Generally requires lots of starting material ( $>10^6$ - $10^7$  cells)

# Typical ChIP-seq pipeline

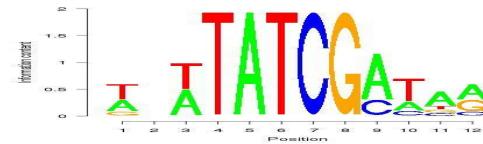
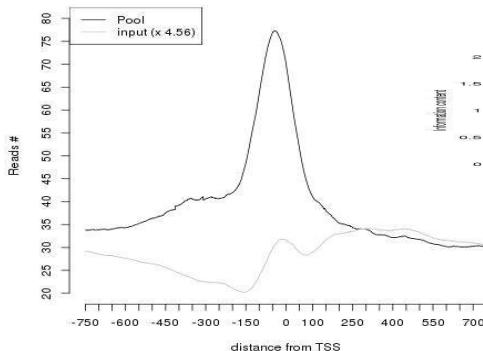
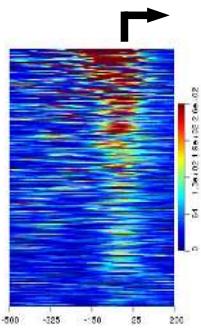
Voir aussi : Bailey et al., Plos Comp. Biol, 2013



Filtered reads (sam/bam)

coverage / profiles / heatmaps

R / Bioconductor  
Python / Perl / ...



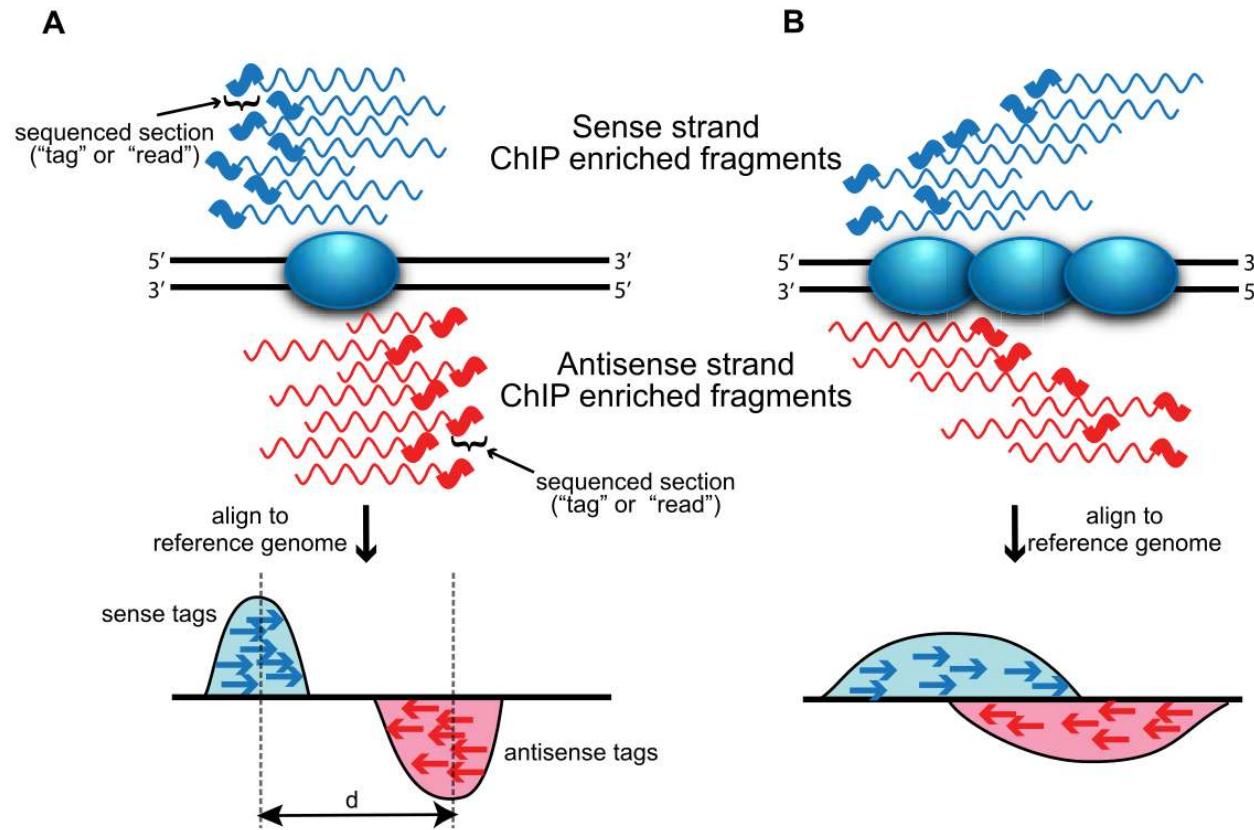
Peak calling → BED file  
MACS, MACS2, SPP  
Sicer (histone mod),...

Counting and differential analysis  
R / Bioconductor...

Motif search and analysis  
RSAT, XXmotif, MEME, R, ...

# Peak calling

The orientation of sequencing (5' to 3') and short read length create a bias in read distribution on +/- strands



A. Sequence-specific binding (TF...)

Wilbanks & Facciotti, PloS One, 2010

B. Broad binding (histone, RNA Pol,...)

## General principles for peak calling :

- Estimate fragment size ( $d$ ) then shift (or extend) the reads by  $d/2$
- Search for regions of enriched signal (# de reads) compared to background (local ideally)

Peak list : BED file

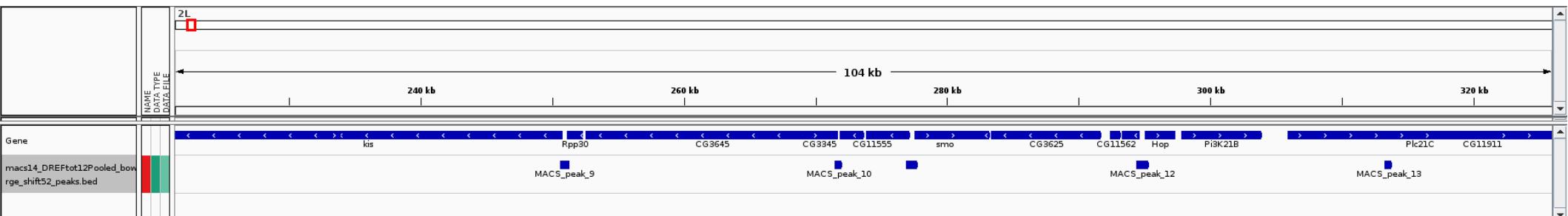
# BED format

[http://fr.wikipedia.org/wiki/BED\\_\(format\\_de\\_fichier\)](http://fr.wikipedia.org/wiki/BED_(format_de_fichier))  
<http://genome.ucsc.edu/FAQ/FAQformat.html>

Chromosome	Start	End	Other facultative info (9 columns)		
2L	47348	52731	MACS_peak_1	342.66	
2L	71901	73894	MACS_peak_2	175.37	
2L	347522	355682	MACS_peak_3	588.17	
2L	403198	404266	MACS_peak_4	75.21	
2L	489732	491282	MACS_peak_5	228.43	

Colonne	Contenu
1	chrom
2	chromStart
3	chromEnd
4	name
5	score
6	strand
7	thickStart
8	thickEnd
9	itemRgb
10	blockCount
11	blockSizes
12	blockStarts

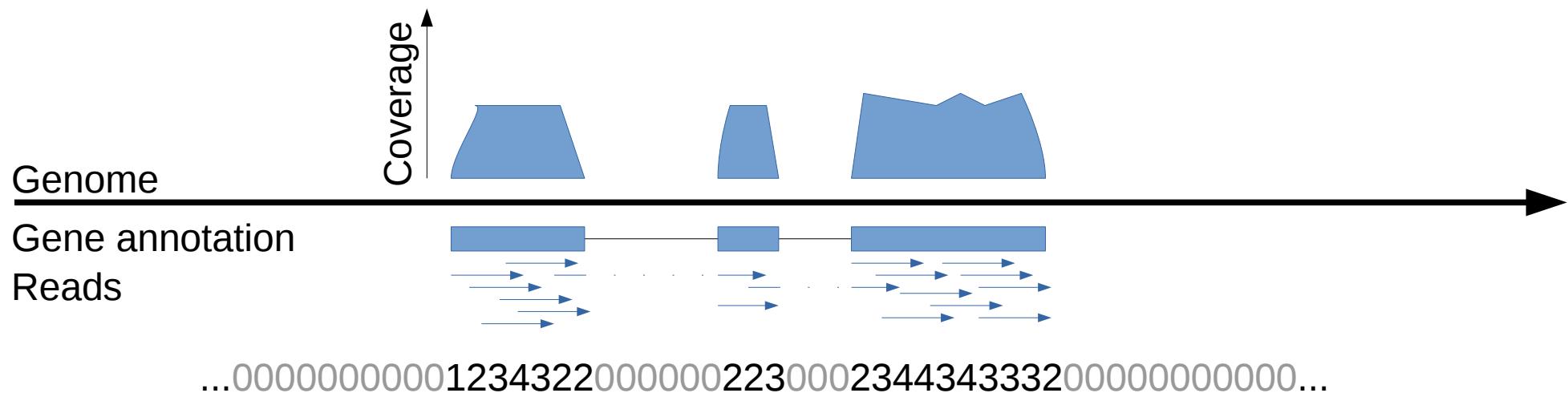
→ Storage of genomic intervals  
 Frequently used in genome browser to represent e.g. annotations or detected signal peaks



Other format for annotations : GFF, GTF (cf **UCSC**)

## Formats for data values along the genome

Wig/bigwig  
Bedgraph } Store a value (coverage) along the genome

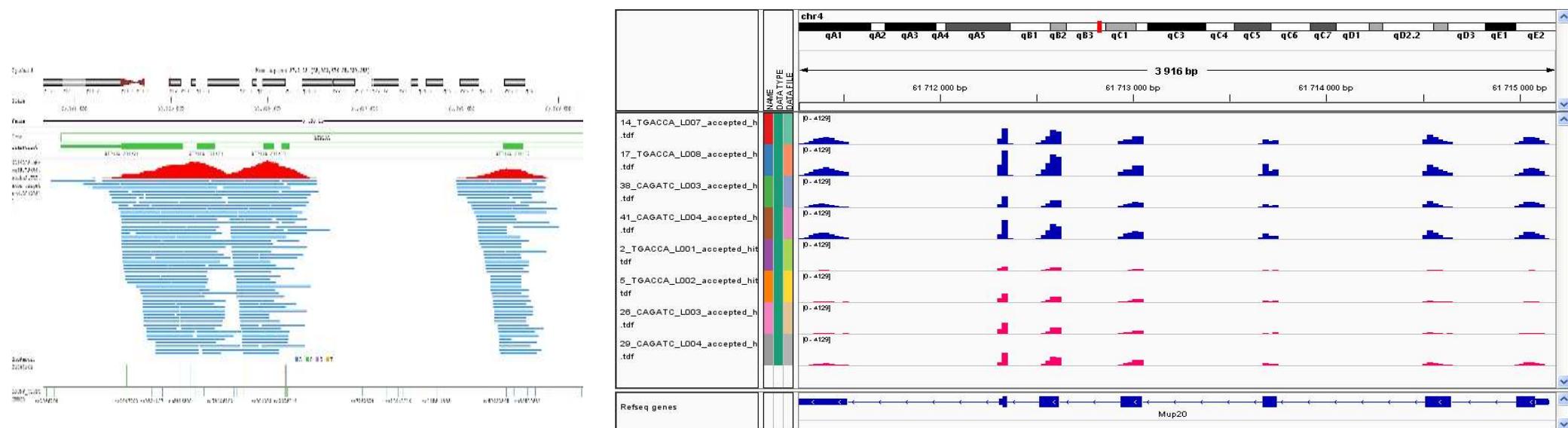
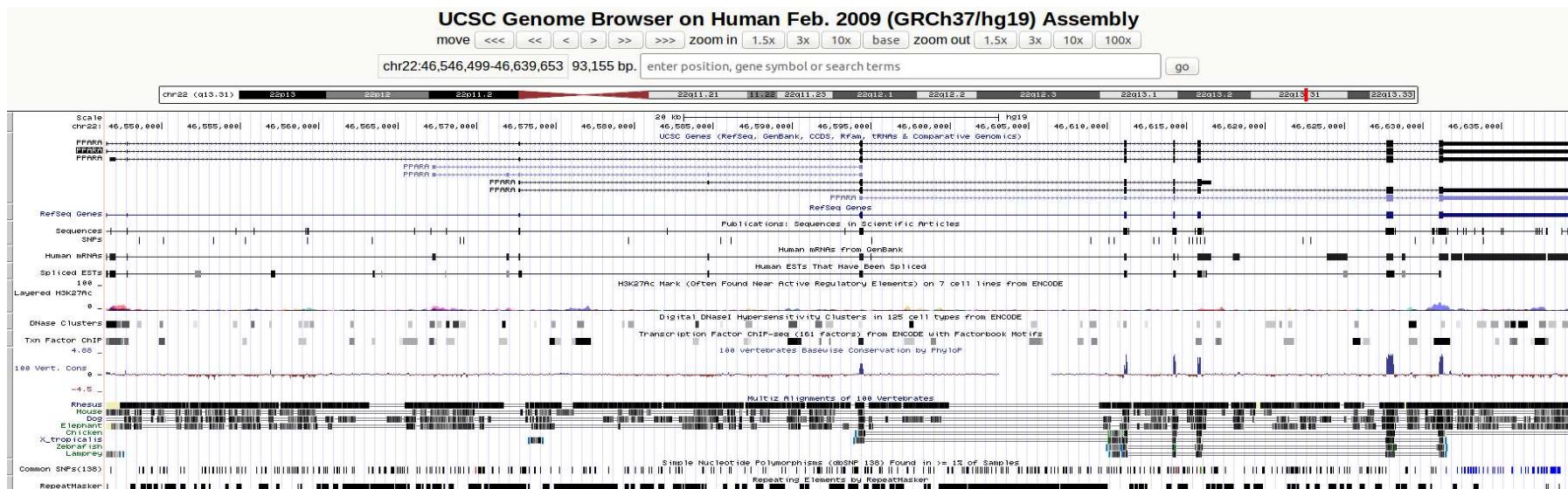


Chr	start	end	value
Chr1	1	125	0
Chr1	126	132	1
Chr1	133	134	2
...			

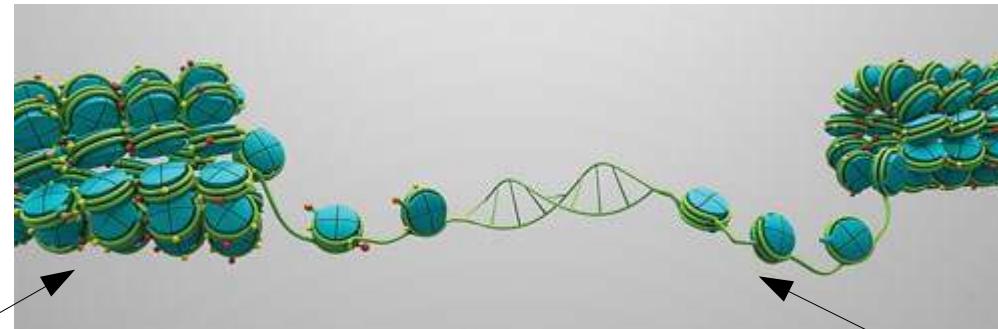
# Genome browser

A genome browser allows to navigate your data

Ex : UCSC GB, Integrative Genomics Viewer, Ensembl GB, jbrowse, ...

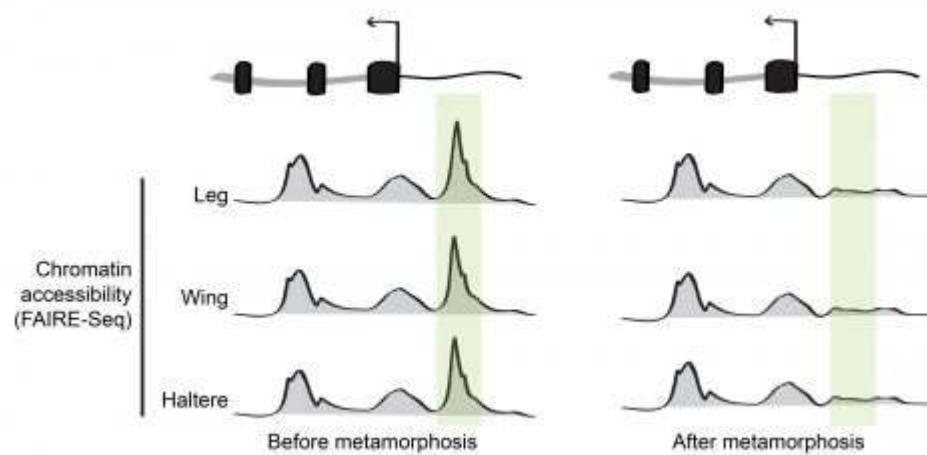


# Why mapping chromatin accessibility ?

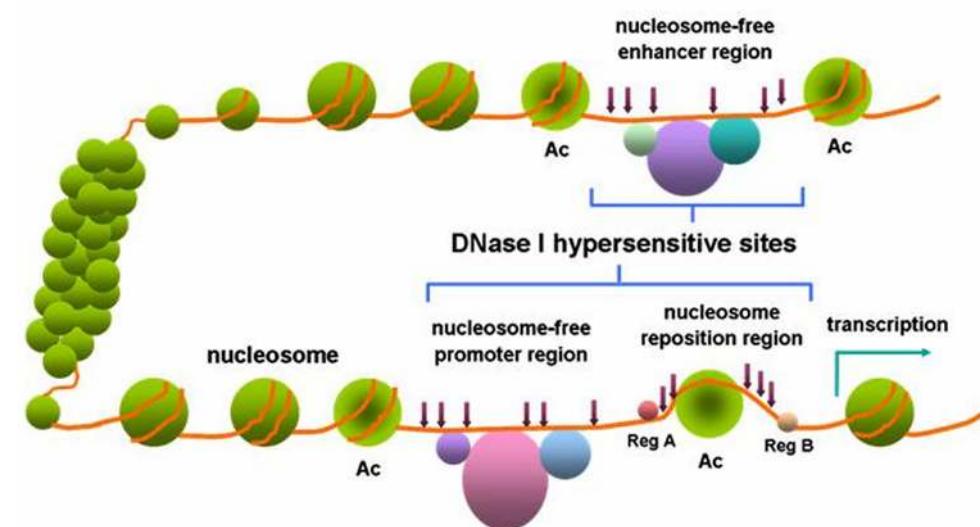


heterochromatin => inaccessible

euchromatin => more accessible



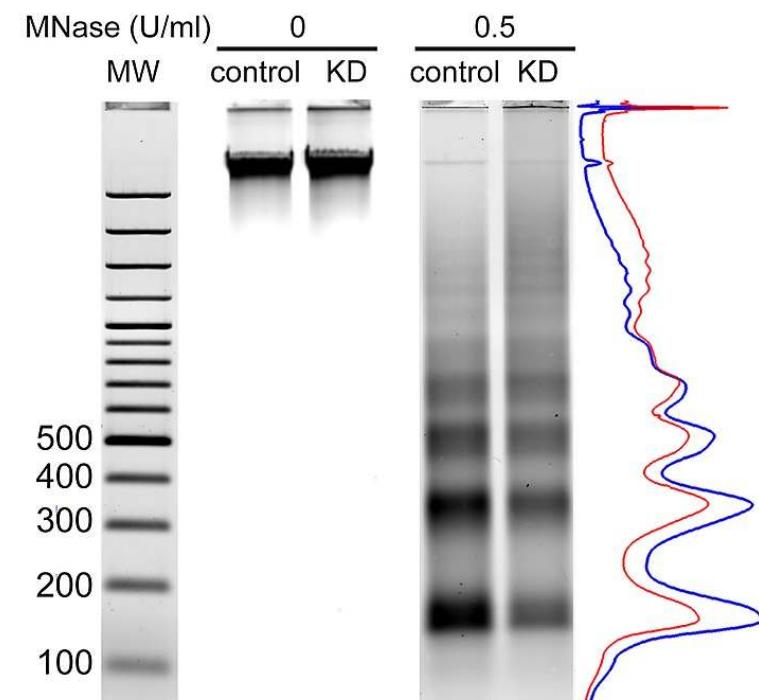
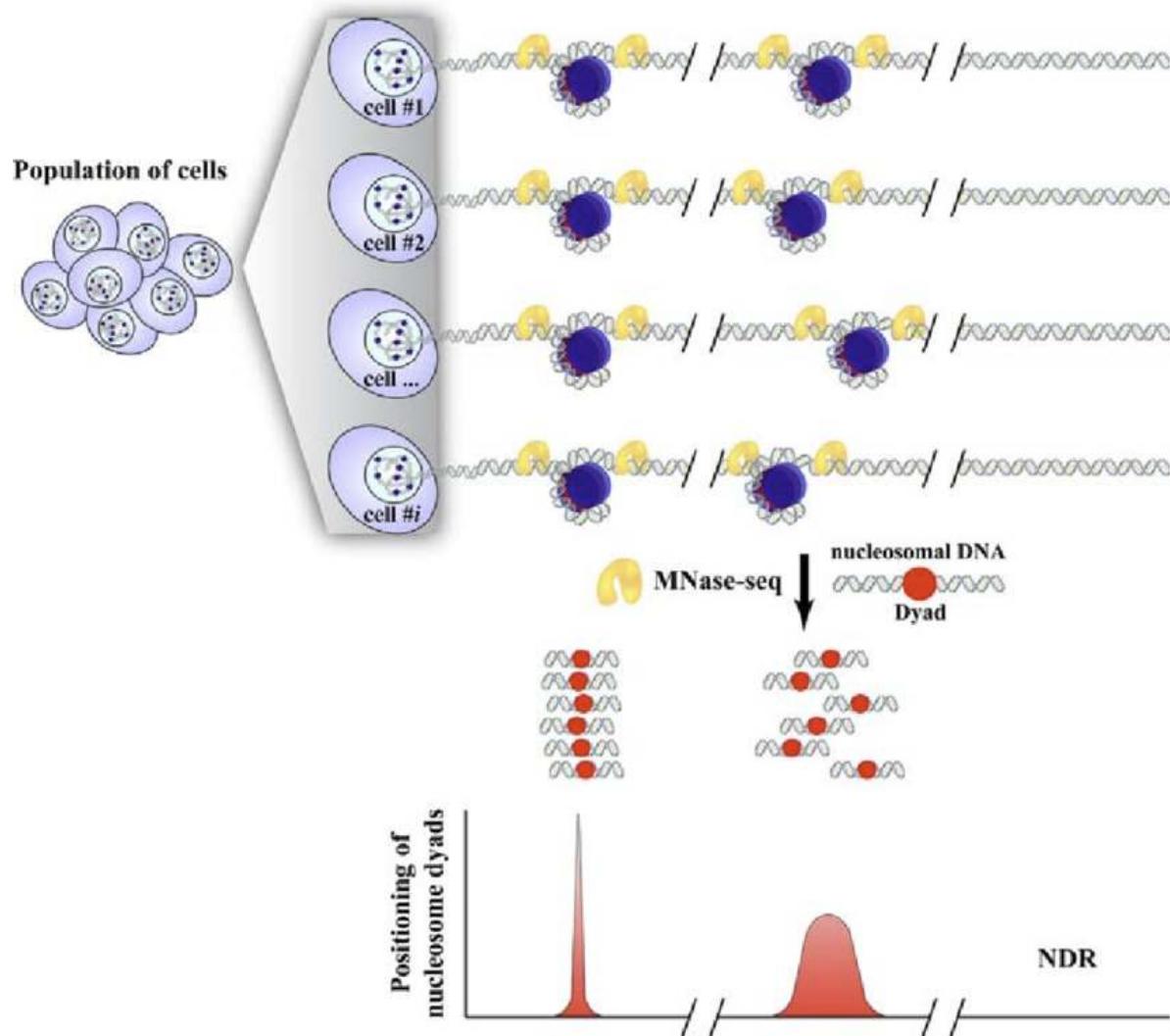
accessible <=> active



active promoters and enhancers  
are highly accessible regions

# MNase-seq

- Mapping nucleosome positions (MNase)



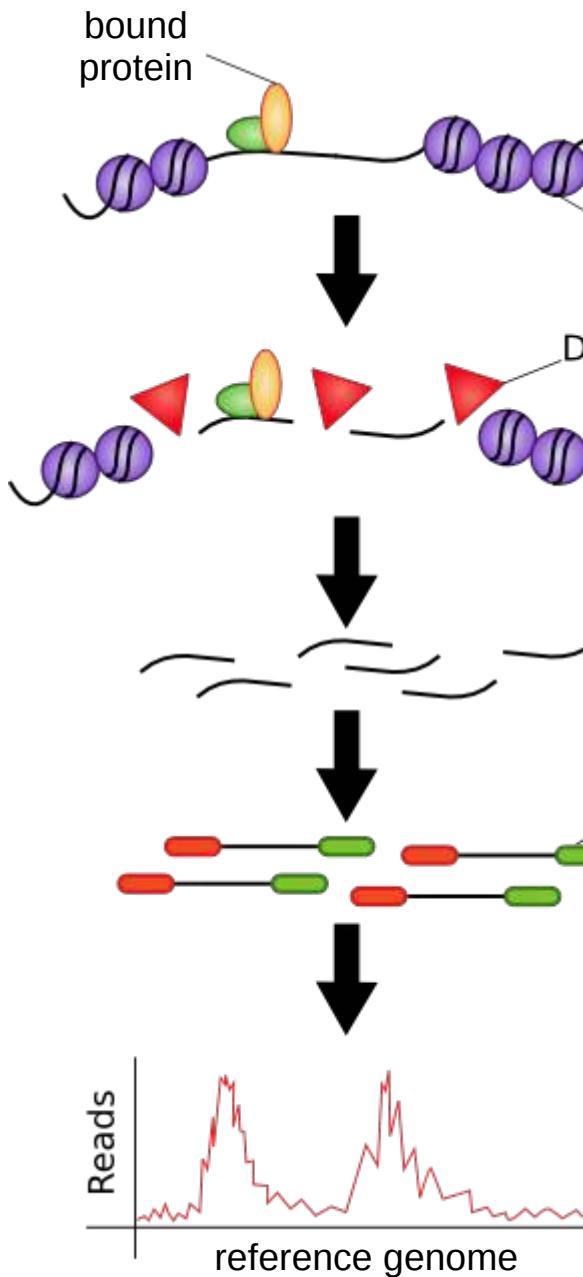
Celona et al., PloS Biol, 2011

Moshkin, AIMS Biophysics, 2015

- Can be combined with ChIP => Native ChIP

# DNase-seq

➤ Mapping accessible / regulatory regions (DNase)



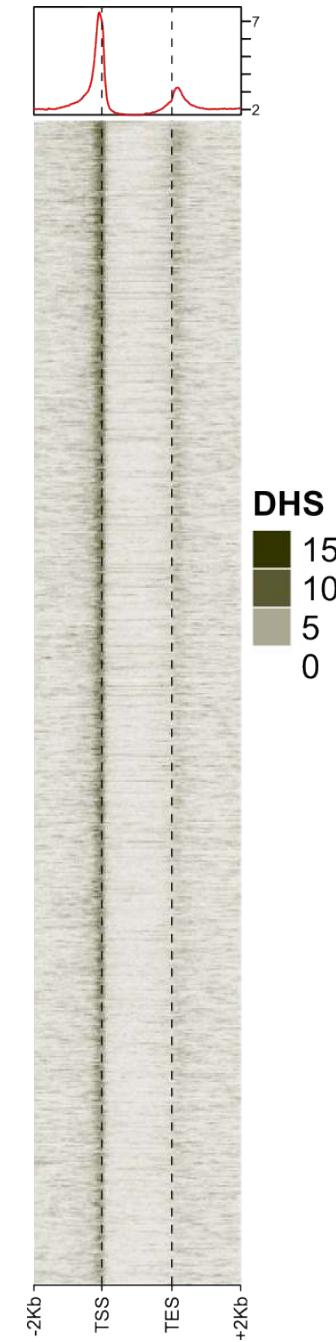
DNA extraction

DNase I digestion

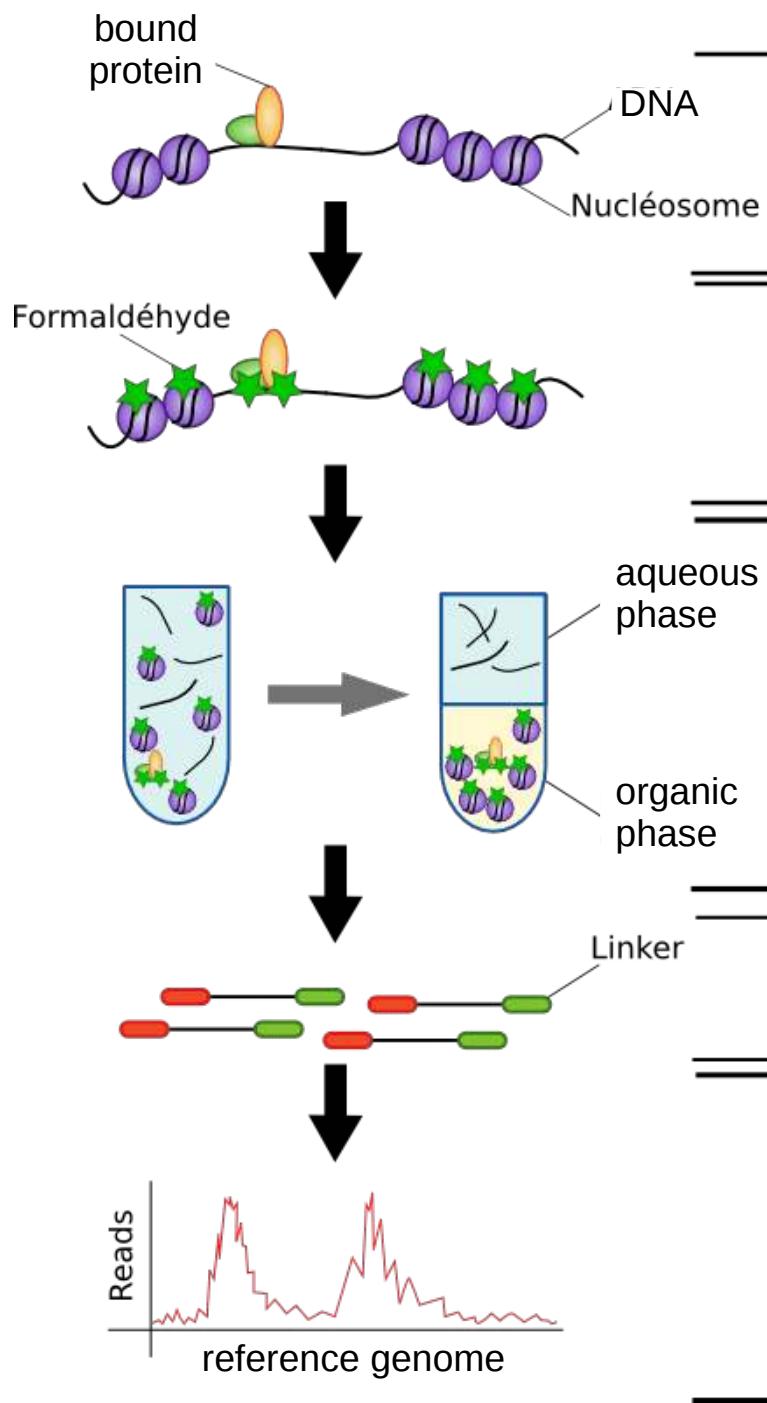
DNA purification

Library prep & amplification

Sequencing & bioinfo



# FAIRE-seq



➤ Mapping accessible / regulatory regions

DNA extraction

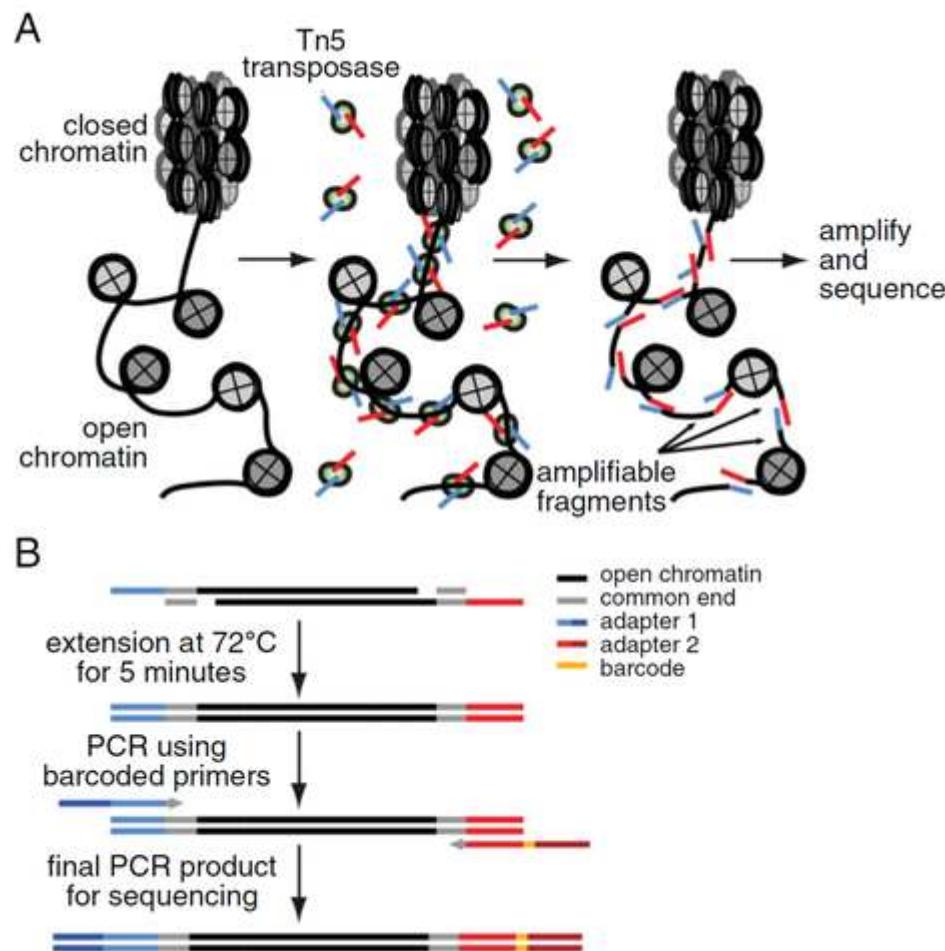
Crosslinking

Sonication and  
DNA purification

Library prep  
& amplification

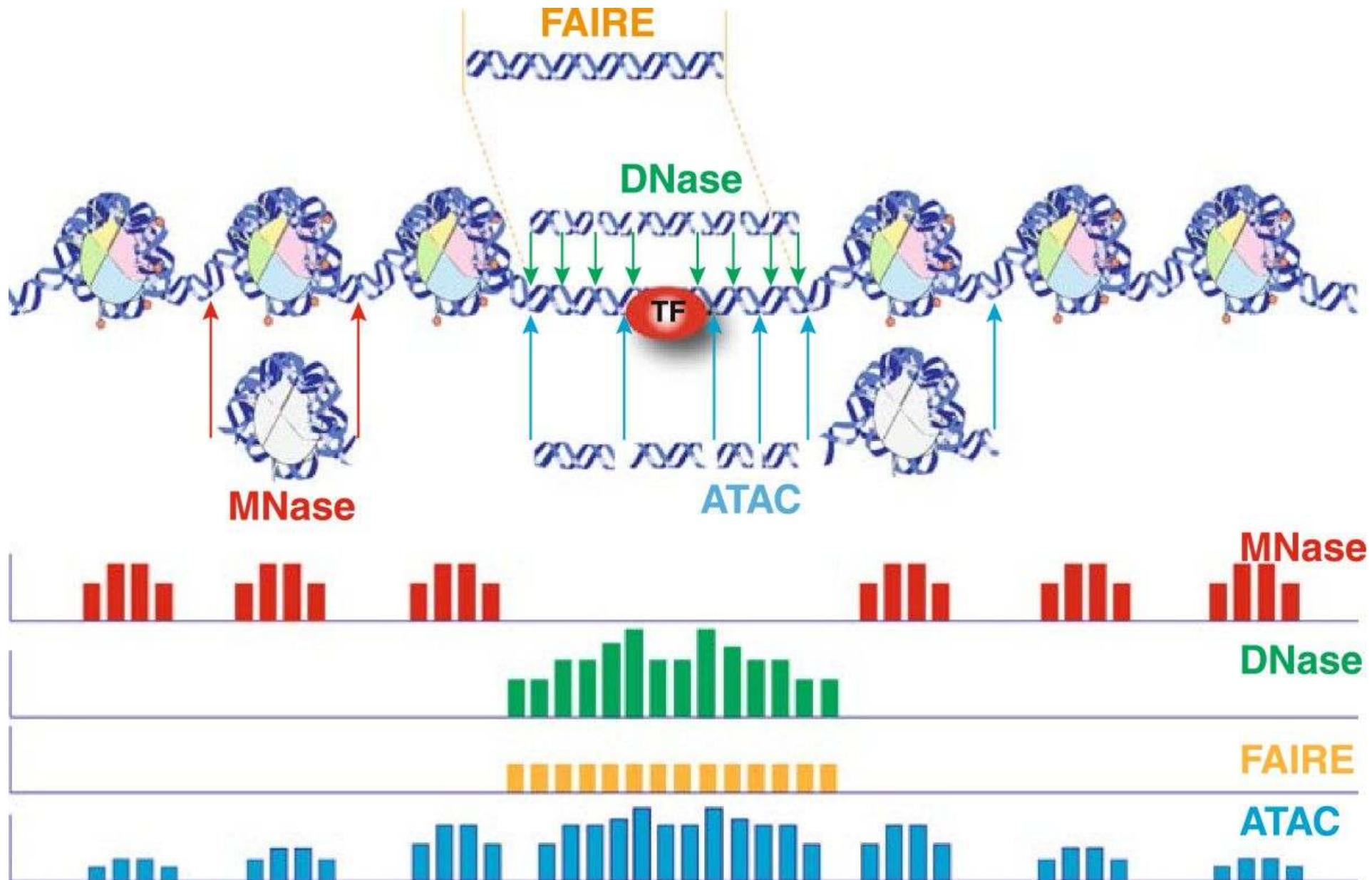
Sequencing &  
bioinfo

# ATAC-seq



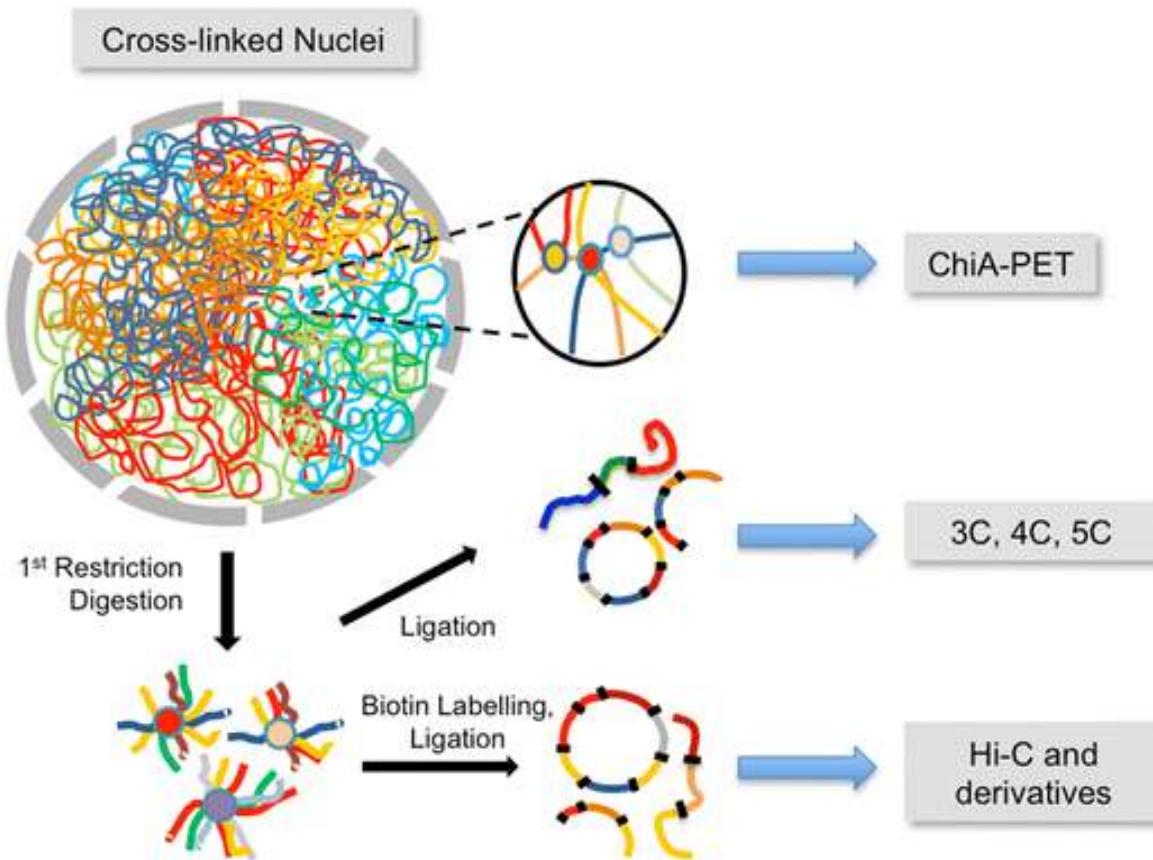
- Mapping of chromatin accessibility
- Fragment sizes carry important information (TF vs nucleosomal)
- Sequencing depth determines usable information
- Low input / single cell
- Often high mDNA contamination

# Chromatin accessibility methods



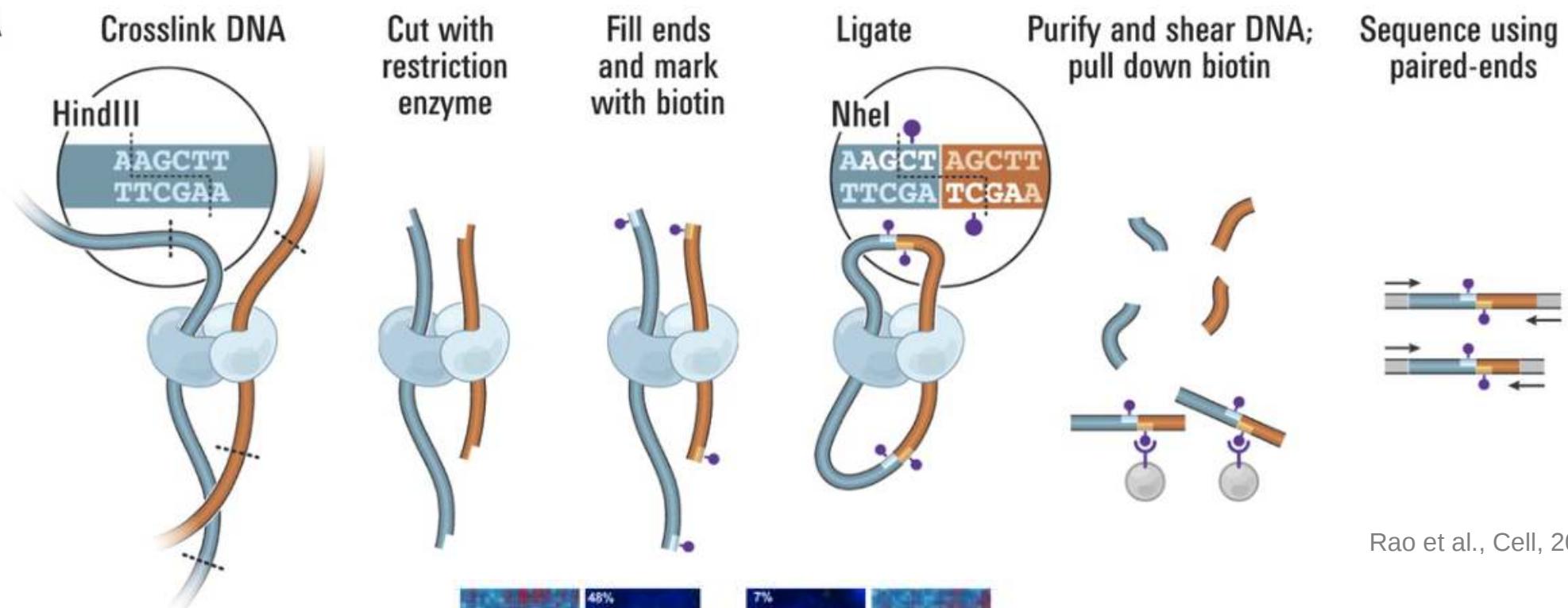
# Chromatin Conformation Capture (3C) and related methods

---

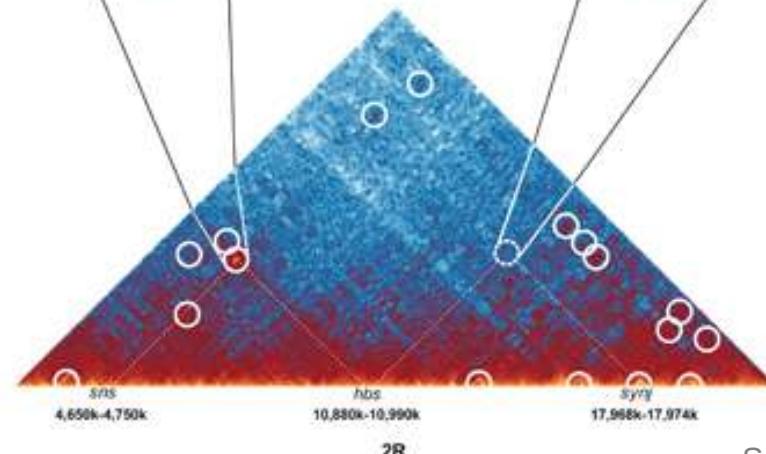


# Chromatin Conformation Capture (Hi-C)

A

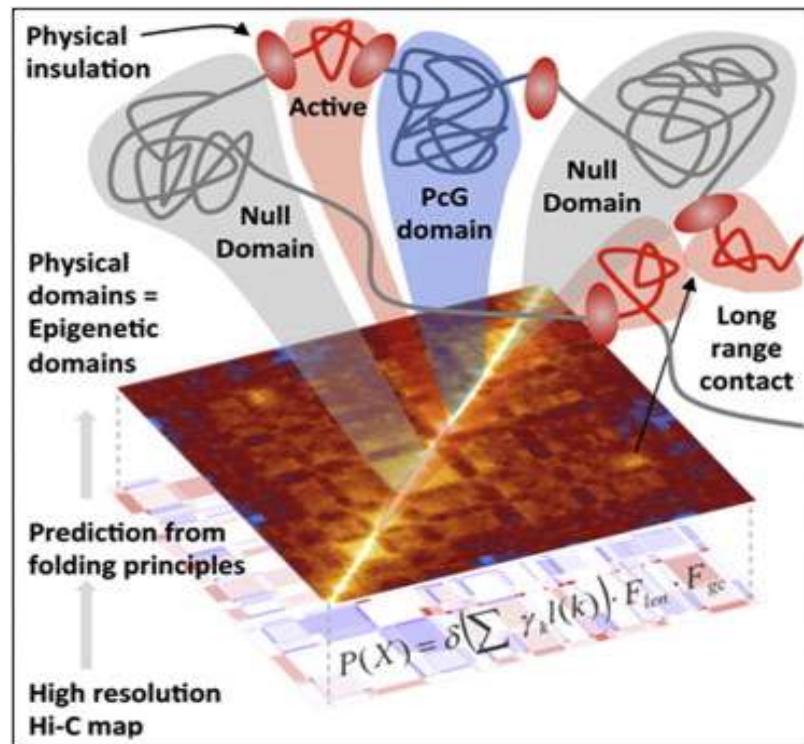
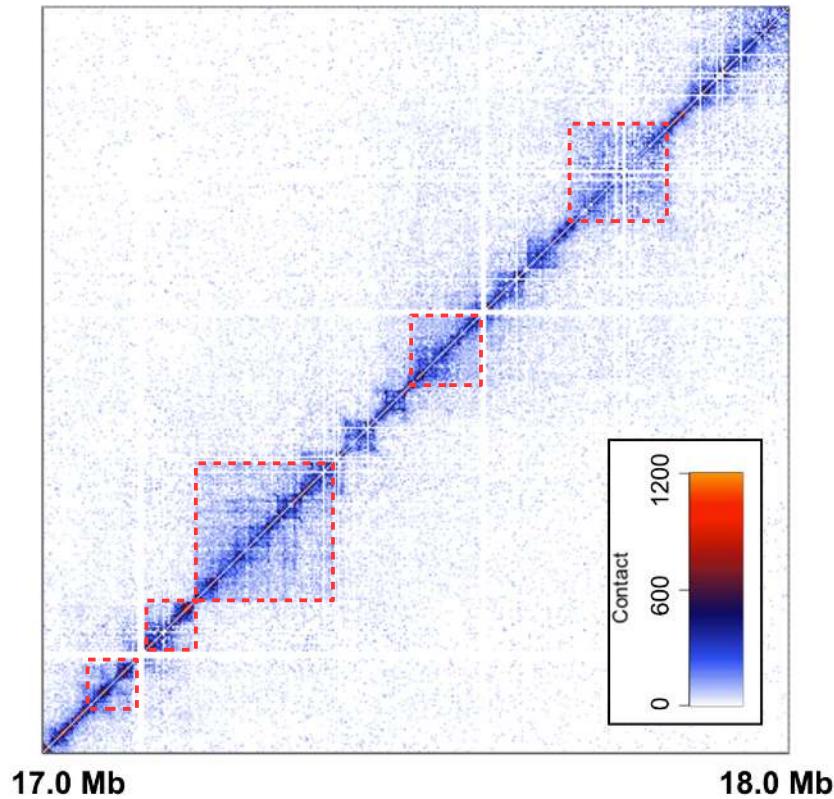


Rao et al., Cell, 2014



Sexton et al., Cell, 2012

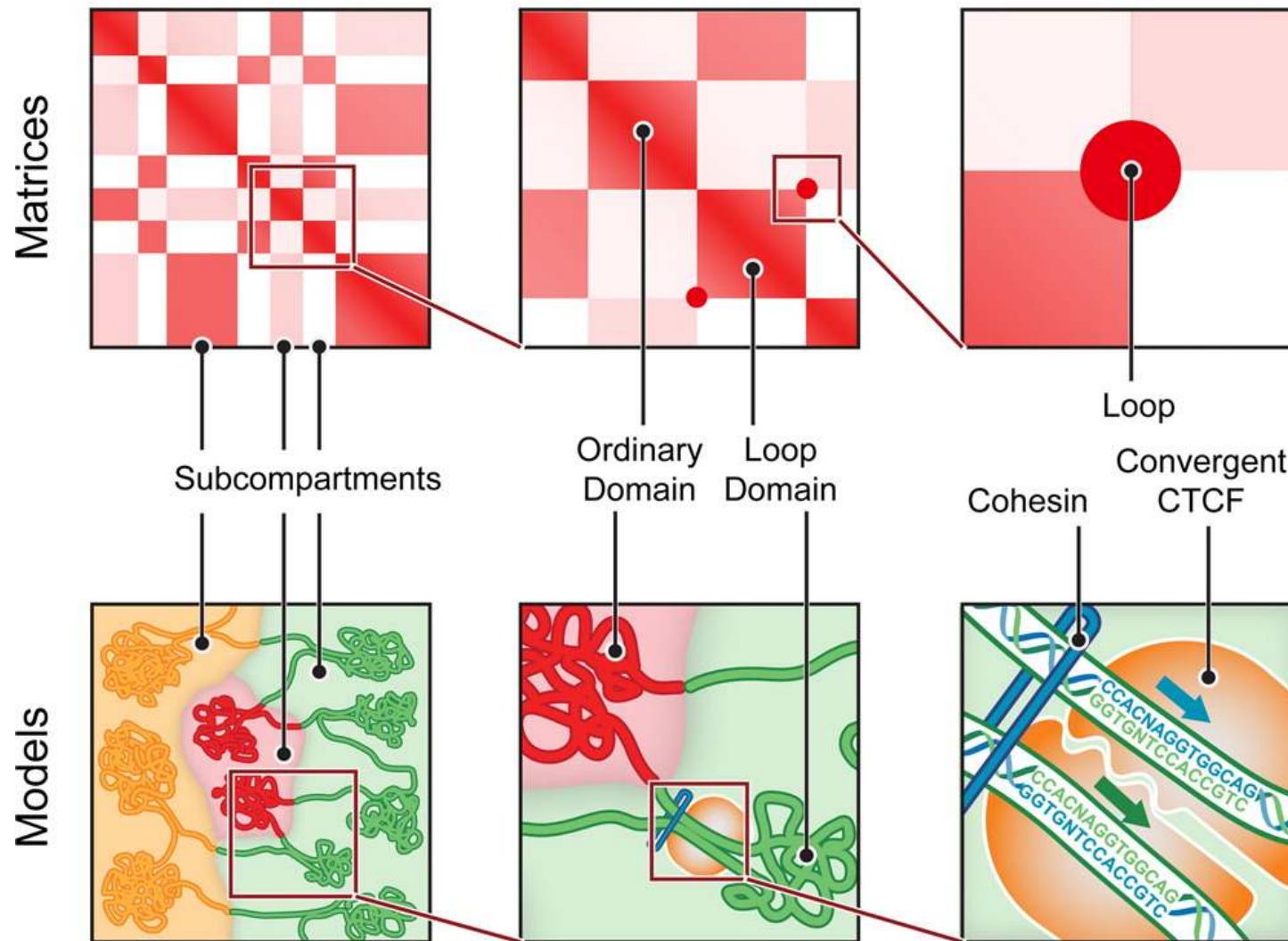
# Topologically associating domains



Sexton et al., Cell, 2012

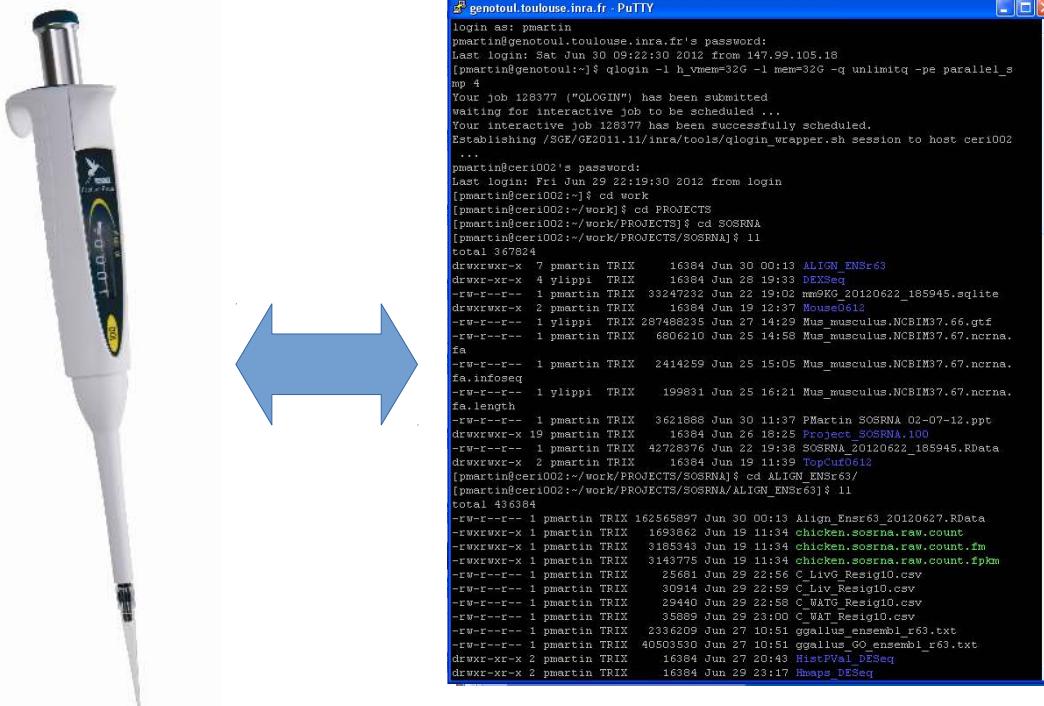
# Compartments and loops

## Hi-C Matrices and Models



# Conclusion

- NGS-based techniques provide an unprecedented view of genome function
- Improved scale, sensitivity and specificity
- New methods constantly developed (e.g. single cell)
- Need for standardization
- Data integration still a challenge



[github.com/pgpmartin](https://github.com/pgpmartin)

# Thank you for your attention



Pascal.Martin@inra.fr



@PgpMartin

