

# SÉQUENÇAGE ET ANALYSE HAUT DÉBIT

# Génome humaine — 2001

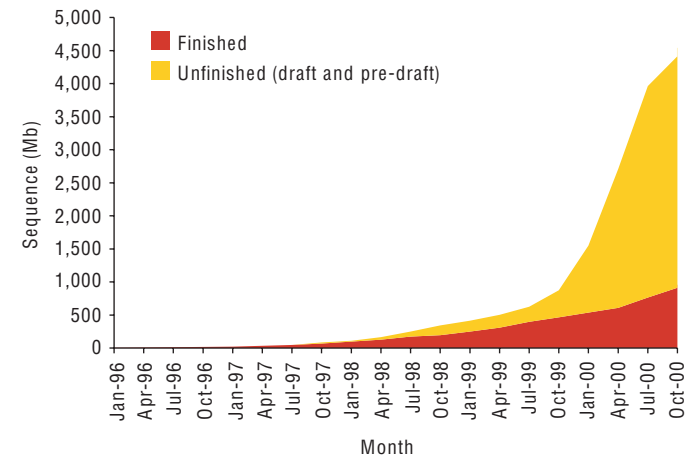
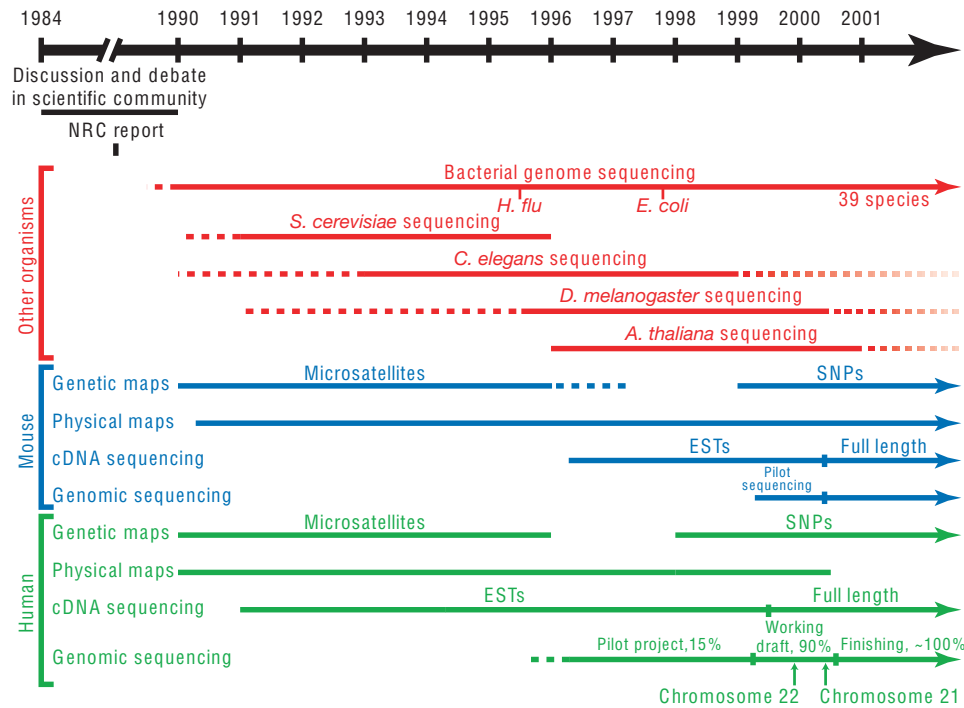
## Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium\*

860

© 2001 Macmillan Magazines Ltd

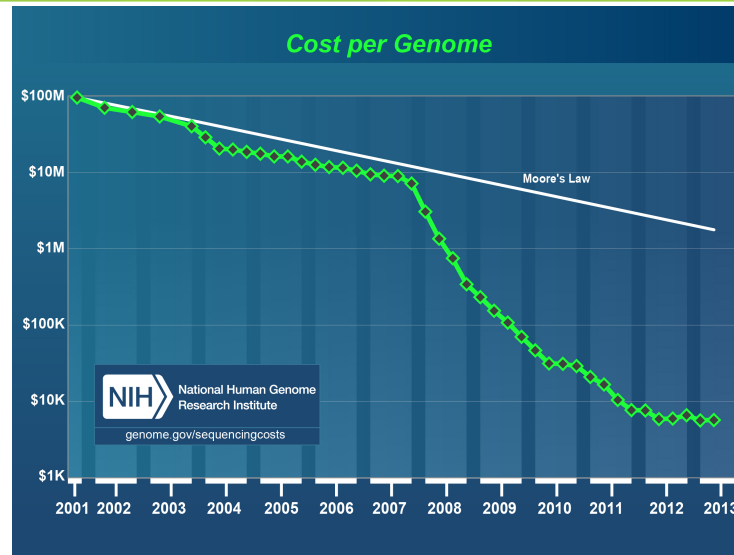
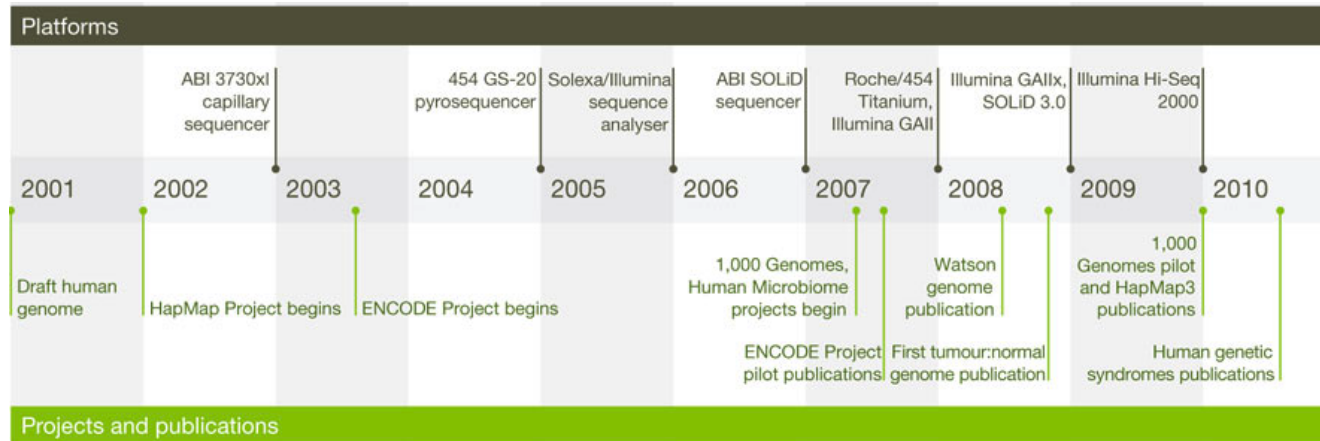
NATURE | VOL 409 | 15 FEBRUARY 2001 | www.nature.com



**Table 2 Total genome sequence from the collection of sequenced clones, by sequence status**

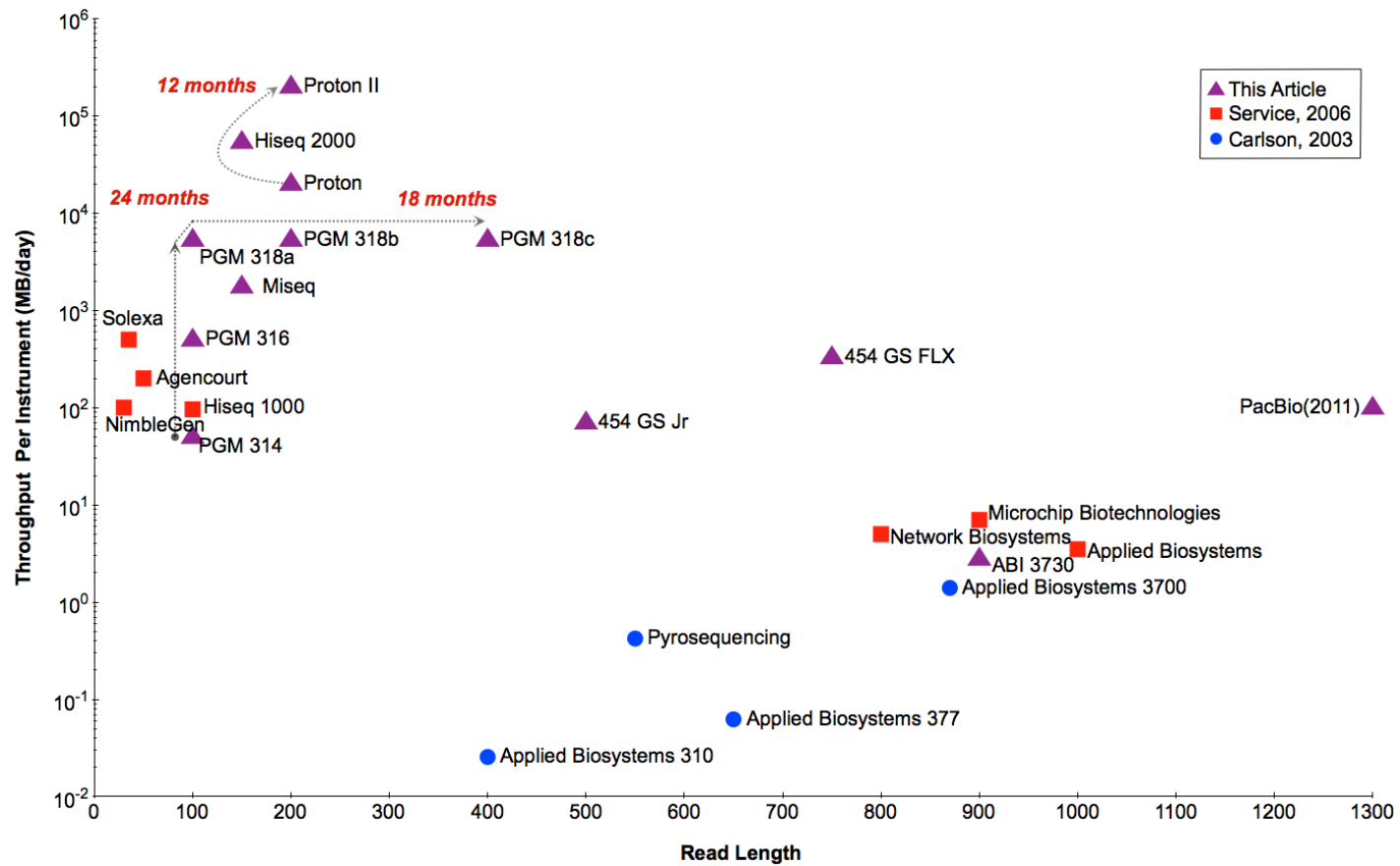
Sequence status	Number of clones	Total clone length (Mb)	Average number of sequence reads per kb*	Average sequence depth†	Total amount of raw sequence (Mb)
Finished	8,277	897	20–25	8–12	9,085
Draft	18,969	3,097	12	4.5	13,395
Predraft	2,052	267	6	2.5	667
<b>Total</b>					<b>23,147</b>

# Début de la siècle



Mardis *Nature* 470 :198 (2011)

# Nouvelles technologies



Rob Carlson, [synthesis.cc](http://synthesis.cc) (avril 2013)

%

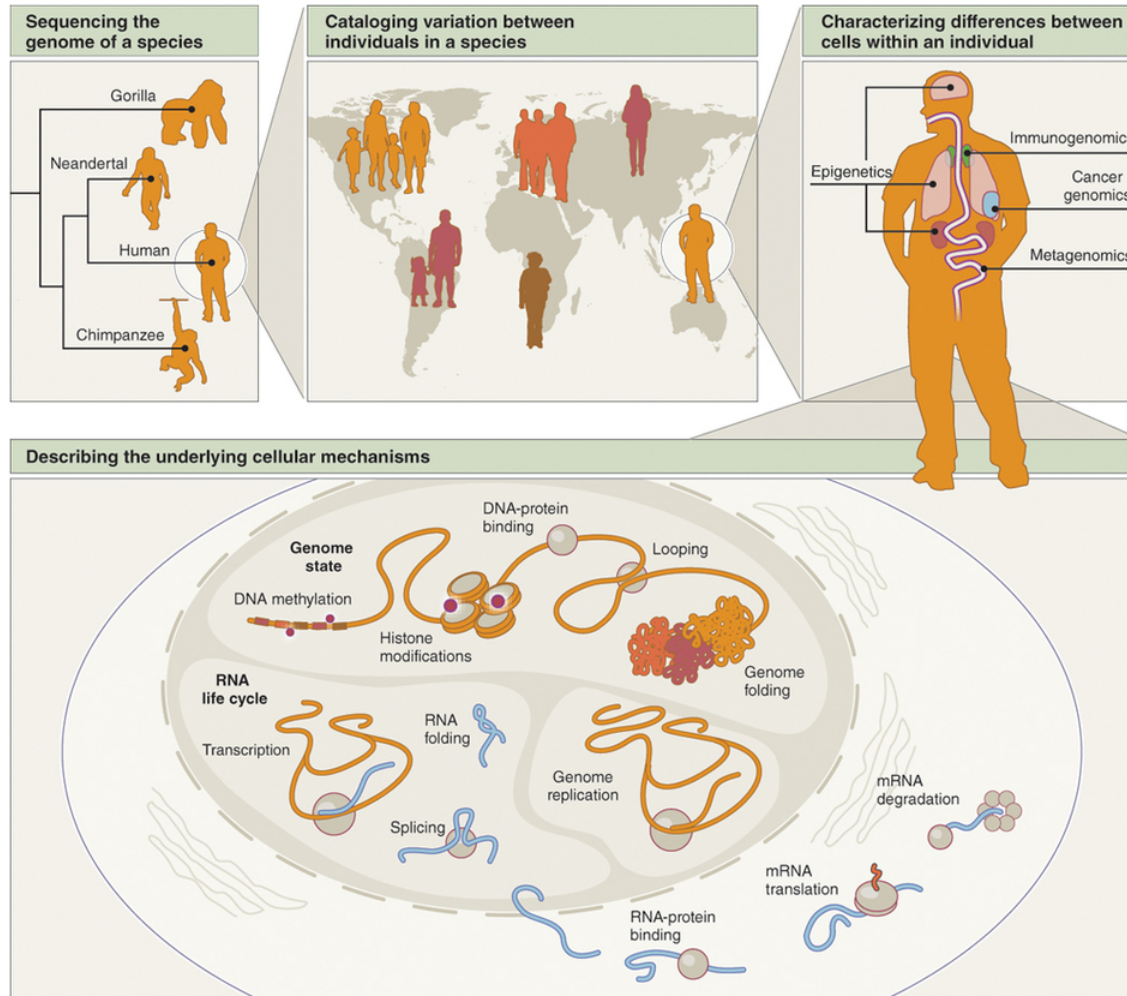
# Nouvelles technologies

**Table 2 Next-generation DNA sequencing instruments**

	Cost per base <sup>a</sup>	Read length (bp) <sup>b</sup>	Speed	Capital cost <sup>c</sup>
<b>Minimum cost per base</b>				
Complete Genomics	Low	Short	3 months	None (service)
HiSeq 2000 (Illumina)	Low	Mid	8 days	+++++++
SOLiD 5500xl (Life Technologies)	Low	Short	8 days	+++
<b>Maximum read length</b>				
454 GS FLX+ (Roche)	High	Long	1 day	+++++
RS (Pacific Biosciences)	High	Very long	<1 day	+++++++
<b>Maximum speed, minimum capital cost and minimum footprint</b>				
454 GS Junior (Roche)	High	Mid	<1 day	+
Ion Torrent PGM (Life Technologies)	Mid	Mid	<1 day	+
MiSeq (Illumina)	Mid	Long	1 day	+
<b>Combined prioritization of speed and throughput</b>				
Ion Torrent Proton (Life Technologies)	Low	Mid	<1 day	++
HiSeq 2500 (Illumina)	Low	Mid	2 days	+++++++

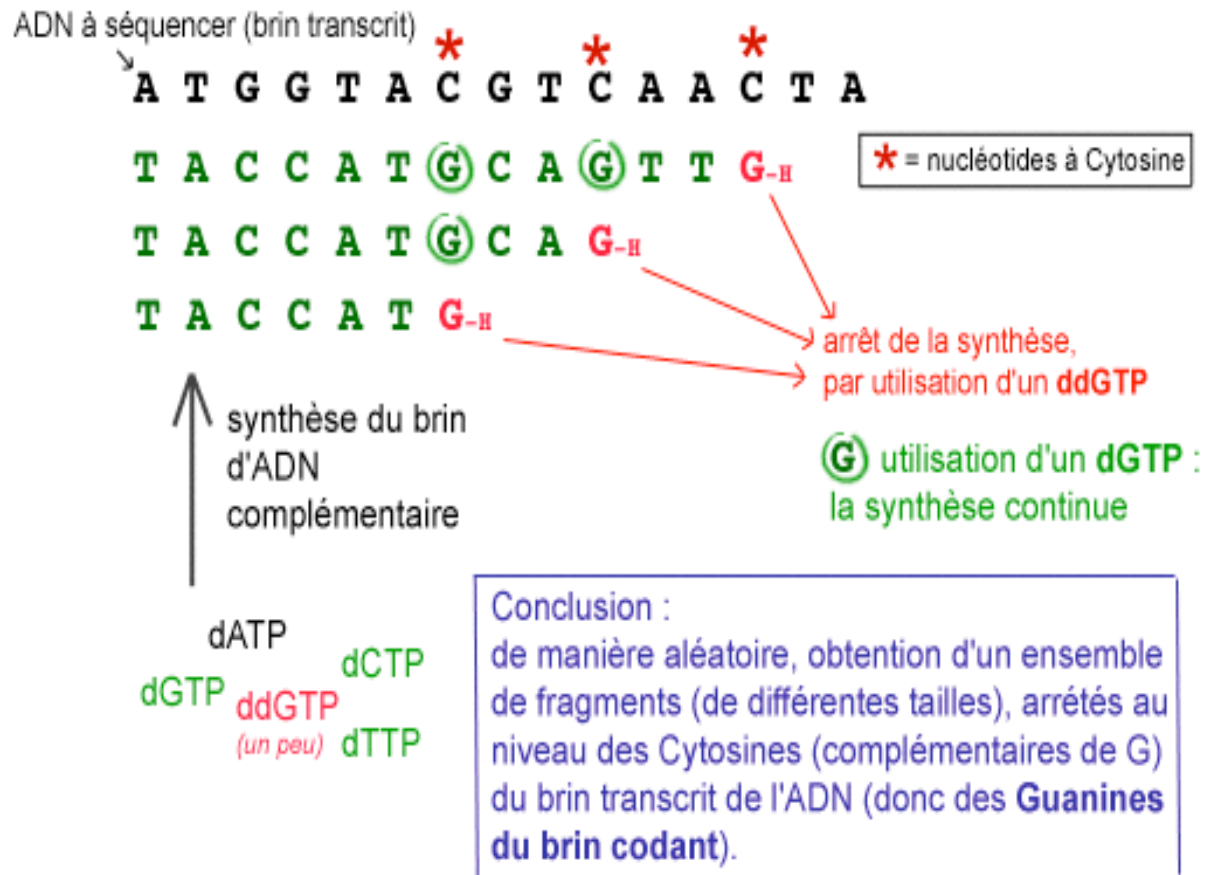
<sup>a</sup>'Low' is < \$0.10 per megabase, 'mid' is in-between and 'high' is > \$1 per megabase. <sup>b</sup>'Short' is < 200 bp, 'mid' is 200–400 bp, 'long' is > 400 bp and 'very long' is > 1,000 bp. <sup>c</sup>Each "+" corresponds to ~\$100,000. We list only commercialized instruments that can be purchased and for which performance data are publically available (as opposed to a comprehensive list of companies developing next-generation sequencing technologies). The categorizations refer to the aspect of sequencing performance to which the technology and/or its implementation in a specific instrument are primarily geared. These estimates were made at the time of publication, and the pace at which the field is moving makes it likely that they will be quickly outdated.

# Nouvelles applications



Shendure & Aiden *Nat Biotechnol* 30 :1084 (2012)

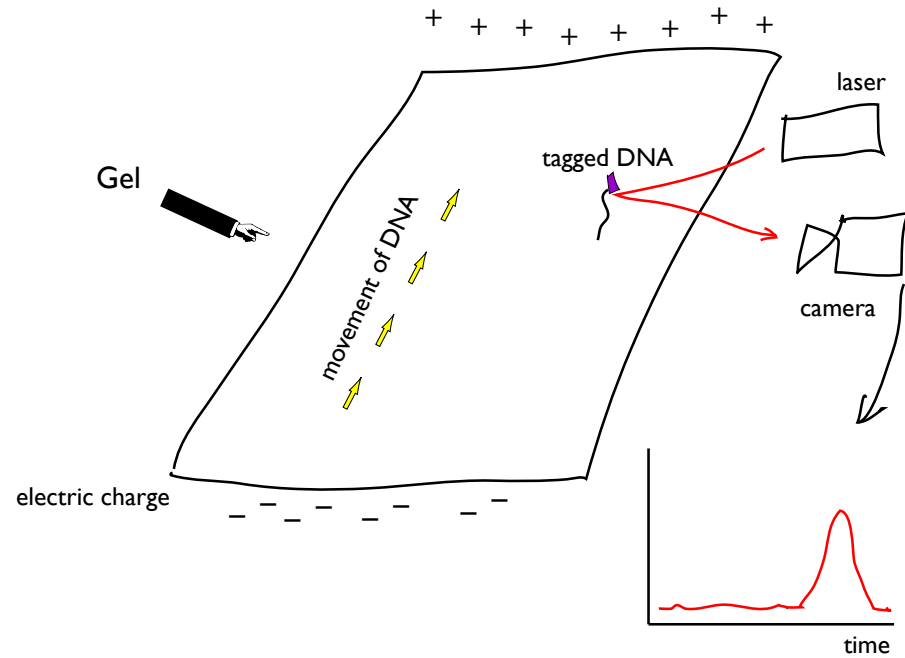
# Séquençage — méthode Sanger



Lodish et al. *Molecular Cell Biology*, 4th ed., W. H. Freeman 1999; Delarue et Furelaud, Jussieu

# Electrophorèse

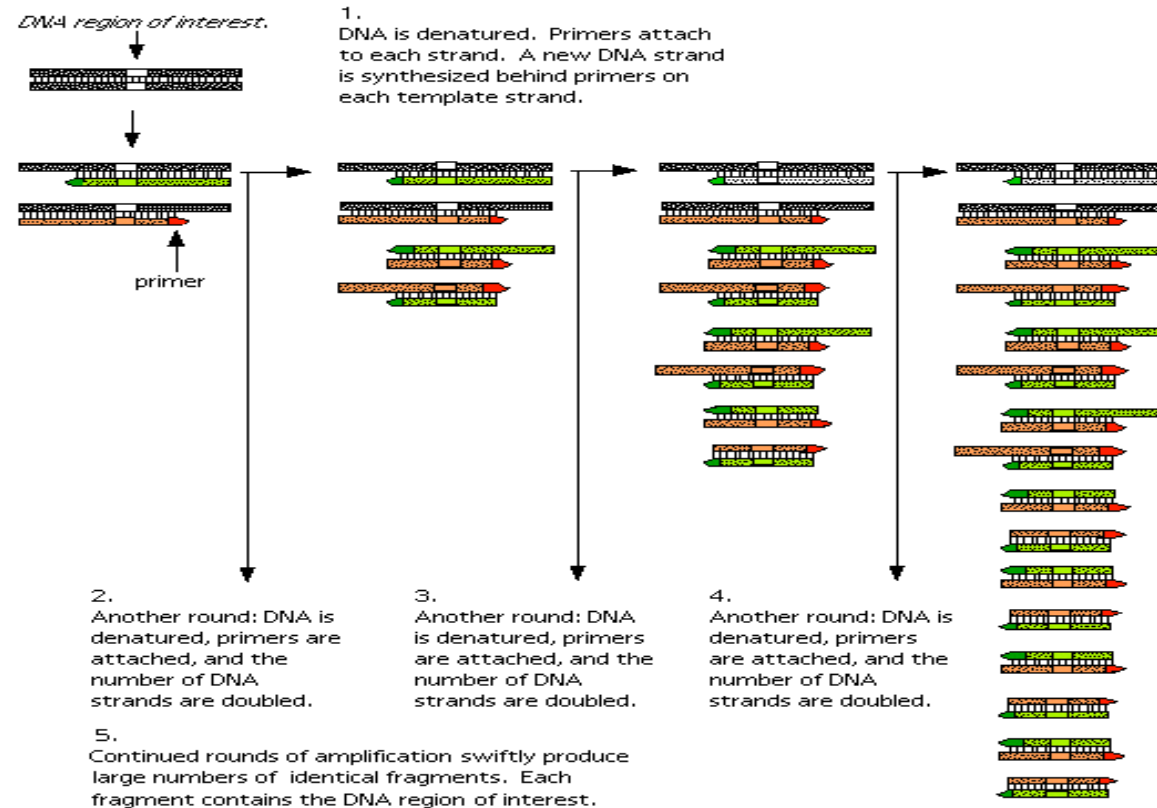
mesure la taille d'un fragment ADN





# Polymerase Chain Reaction

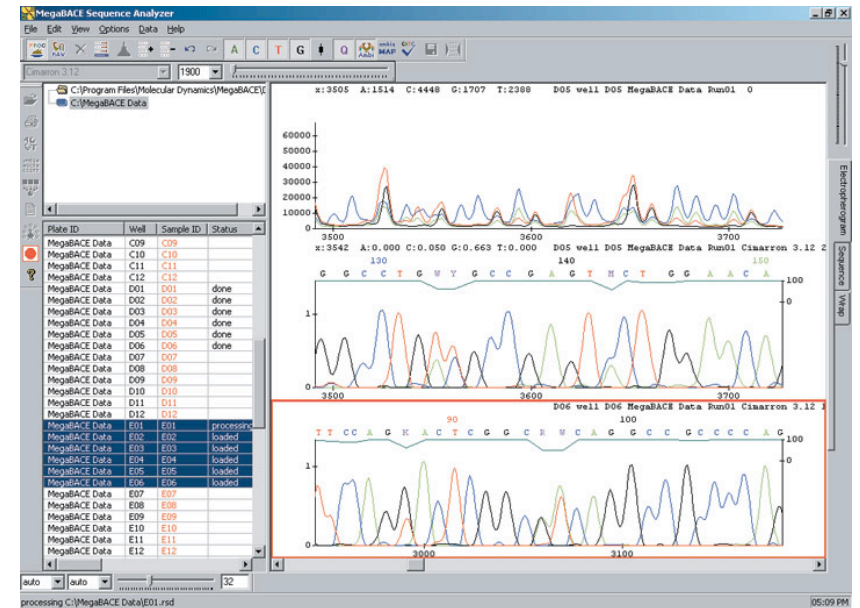
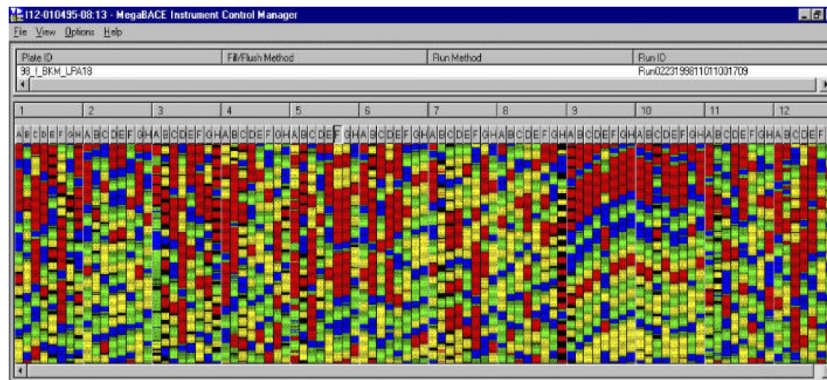
## POLYMERASE CHAIN REACTION



# Séquençage automatique

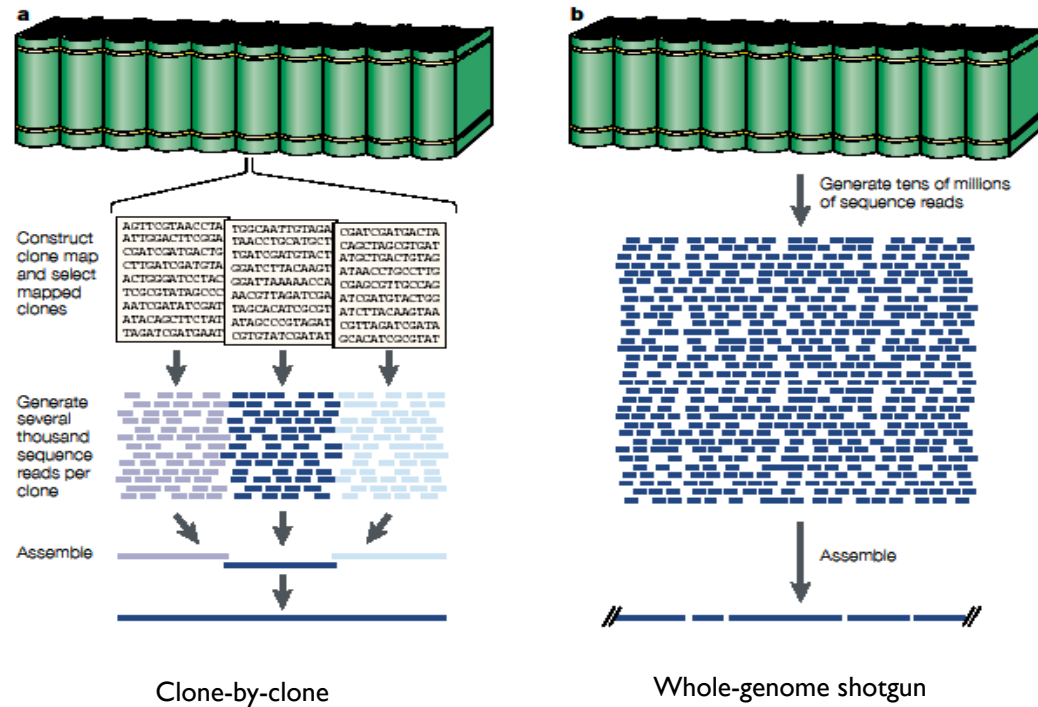
→ marqueurs fluorescents, multiplexage, capillaires, cycleurs thermiques

600–1000 pb, 96–384 capillaires, quelques heures



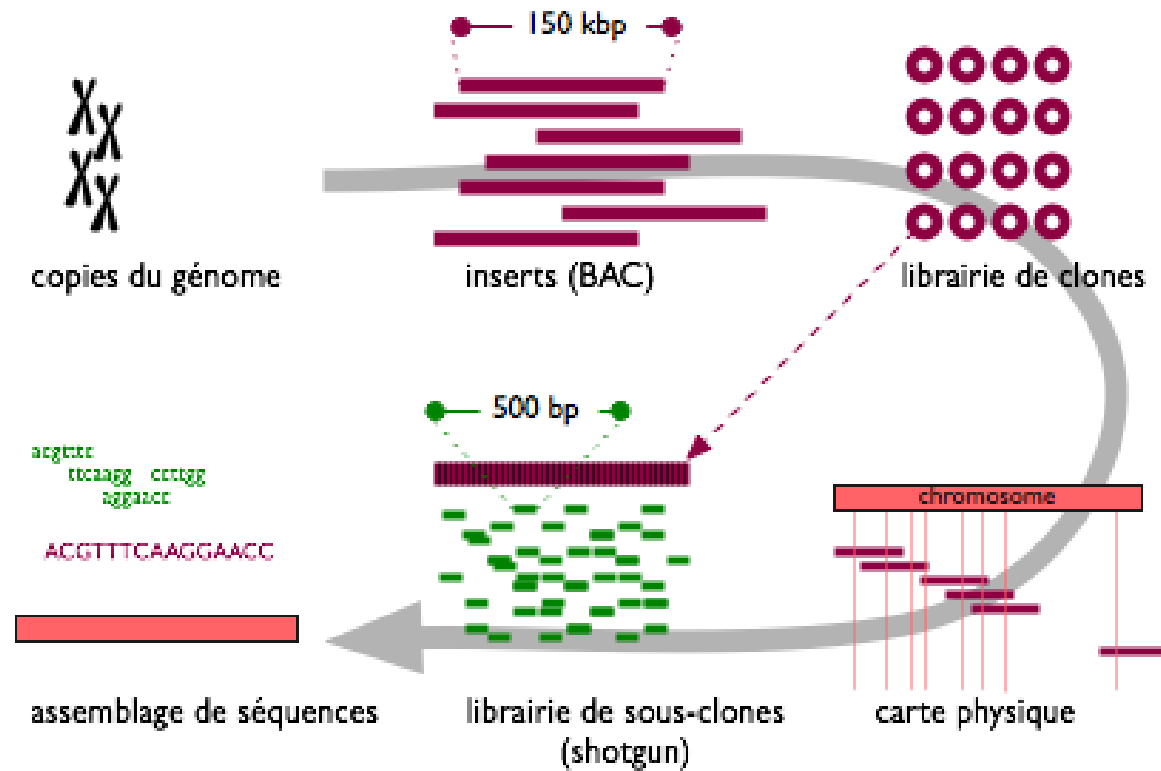
Molecular Dynamics → Amersham Biosciences → GE Life (Megabace)

# Séquençage du génome humaine



E. Green. Nature Reviews Genetics 2:573 (2001)

# Approche hiérarchique



# Phred

deux fichiers : séquences (fasta) et valeurs de qualité (.qual)  
nom du fichier de l'électrophérogram (image)

```
>SARS211.B21_A07_-_032.ab1 CHROMAT_FILE: SARS211.B21_A07_-_032.ab1 PHD_FILE:
SARS211.B21_A07_-_032.ab1.phd.1 CHEM: term DYE: big TIME: Thu May 15 13:44:06
2003 TEMPLATE: SARS211A07_ DIRECTION: fwd
GGAATAAATTCXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXXXXXCCCATTTGAACCAATTGNGNACTTTNAC
TAAAAGNACCAANTCTCCATTTAGAGCTTCACTACCTACAACATTGCTA
AAAATAGTGTAAAGAGTGTGCTAAATTATGTTGGATGCCGGCATTAAAT
TATGTGAAGTCACCCAAATTTCTAAATTGTTCACAATCGCTATGTGGCT

>SARS211.B21_A07_-_032.ab1 PHD_FILE: SARS211.B21_A07_-_032.ab1.phd.1
7 7 7 7 7 7 8 8 8 6 6 6 8 13 9 7 7 8 6 6 6 6 9 7
9 23 15 20 25 27 32 32 32 34 40 33 40 40 40 37 37 40
40 37 46 34 34 34 42 42 42 42 51 51 40 40 40 37 39
29 29 21 21 25 33 33 32 18 16 12 20 15 8 8 8 8 8
8 8 8 9 9 9 11 16 21 4 0 4 0 4 27 27 27 4 0 4 25 25
36 36 27 27 4 0 4 27 21 17 4 0 4 15 10 9 9 9 9 12 25
17 28 28 32 48 46 46 40 40 40 40 44 56 56 56 56 47
47 56 56 56 56 56 56 56 44 44 42 47 42 42 42 42 42
42 42 42 42 42 44 44 44 42 42 38 38 38 38 42 42
```

probabilité d'erreur pour qualité  $q$  :  $10^{-q/10}$

# FASTQ

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*(((((***+))%%%++) (%%%) .1***-+*''))**55CCF>>>>>CCCCCCC65
```

- ★ Ligne  $4k$  commence par '@' et introduit l'identificateur du morceau (*read*)
- ★ Ligne  $4k + 1$  donne la séquence
- ★ Ligne  $4k + 2$  commence par '+' et *peut* répéter l'identificateur de Ligne  $4k$  ou non
- ★ Ligne  $4k + 3$  donne les valeurs de qualité (encodage Sanger avec échelle Phred  $q + ' !' : 0 \rightarrow !, 1 \rightarrow ", \dots, 31 \rightarrow @, 32 \rightarrow A, \dots, 40 \rightarrow I$ )

# Méthodes de *next generation*

Principes :

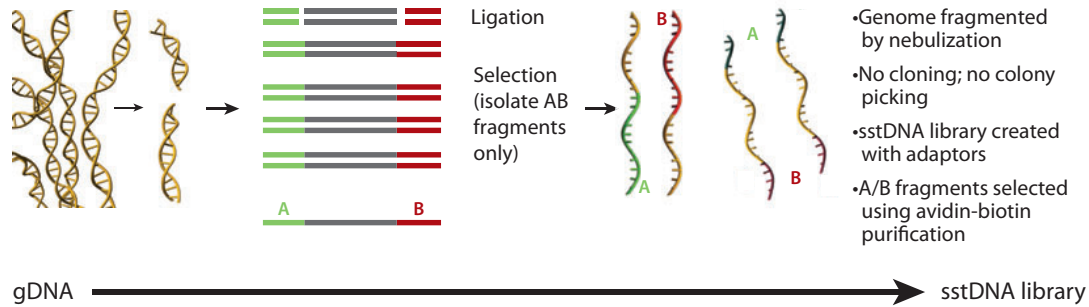
1. fixer fragments d'ADN ciblés sur une surface → arrangement en espace pour parallélisation («puce»)  
amplification (si nécessaire)
2. séquençage de fragments d'ADN + production d'image
3. analyse informatique de l'image : production de séquences + estimation d'erreurs de séquençage

# 454 : perles (*beads*)

**a**

## DNA library preparation

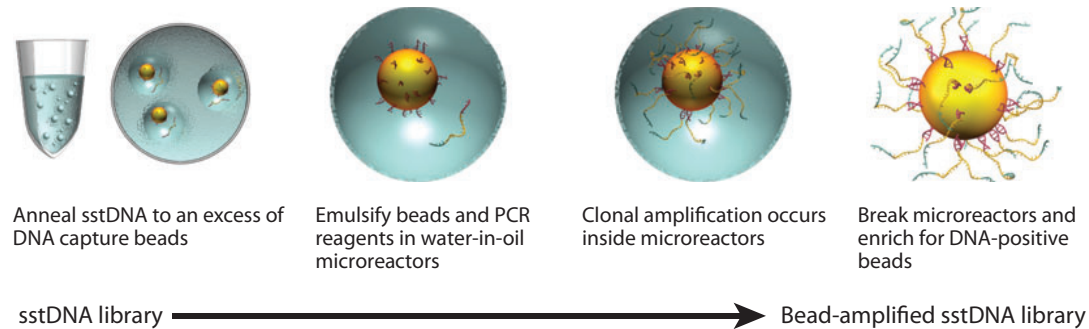
4.5 hours



**b**

## Emulsion PCR

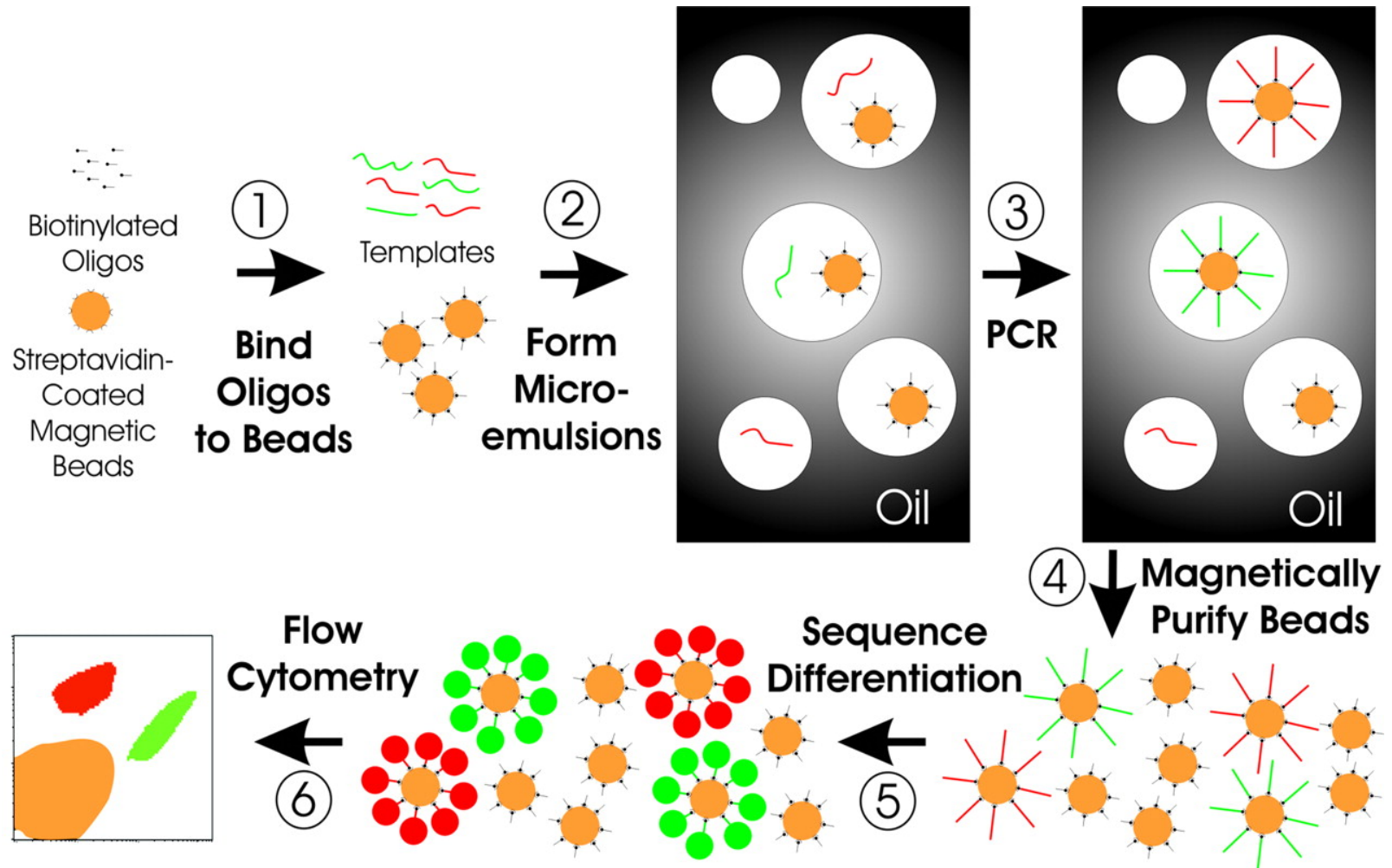
8 hours



Mardis *Annu Rev Genomics Hum Genet* 9 :387 (2008)

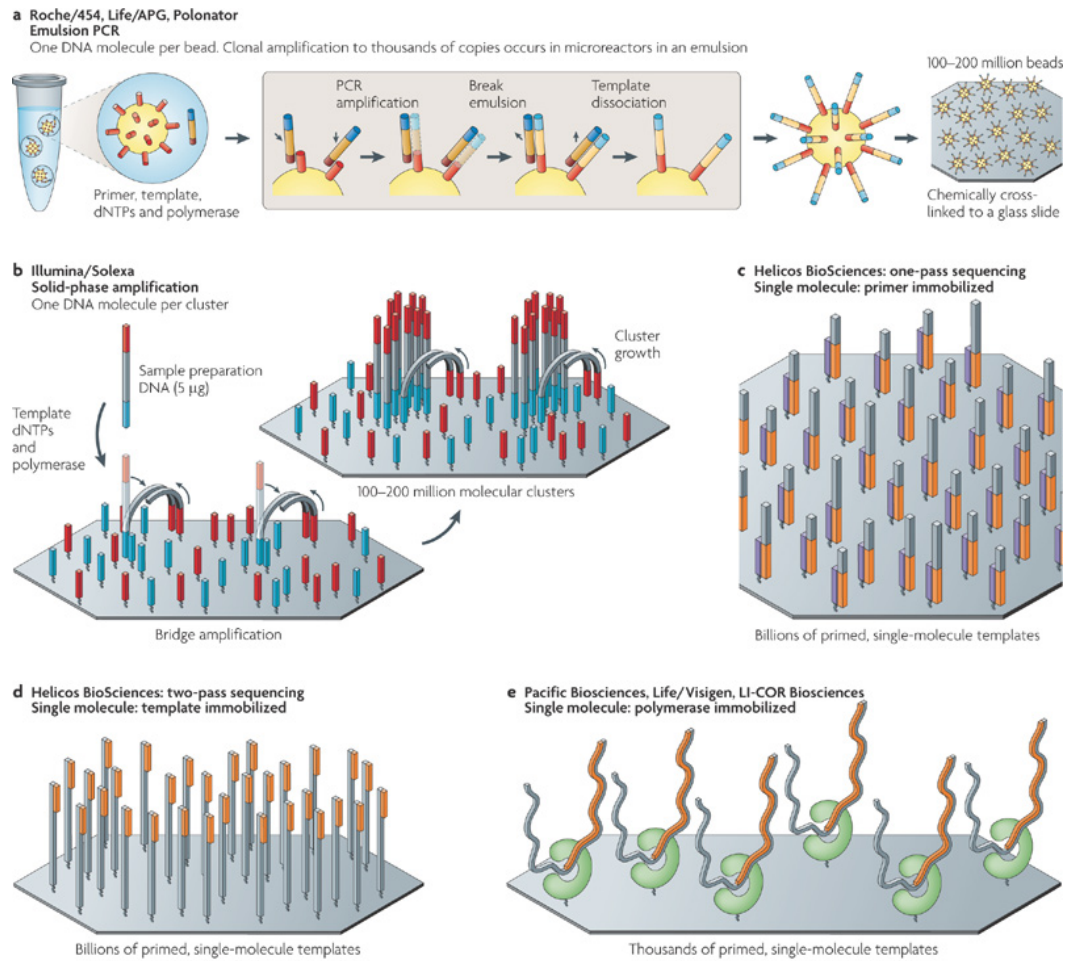


# PCR en emulsion



Dressman & al *PNAS* 100 :8817 (2003); →  

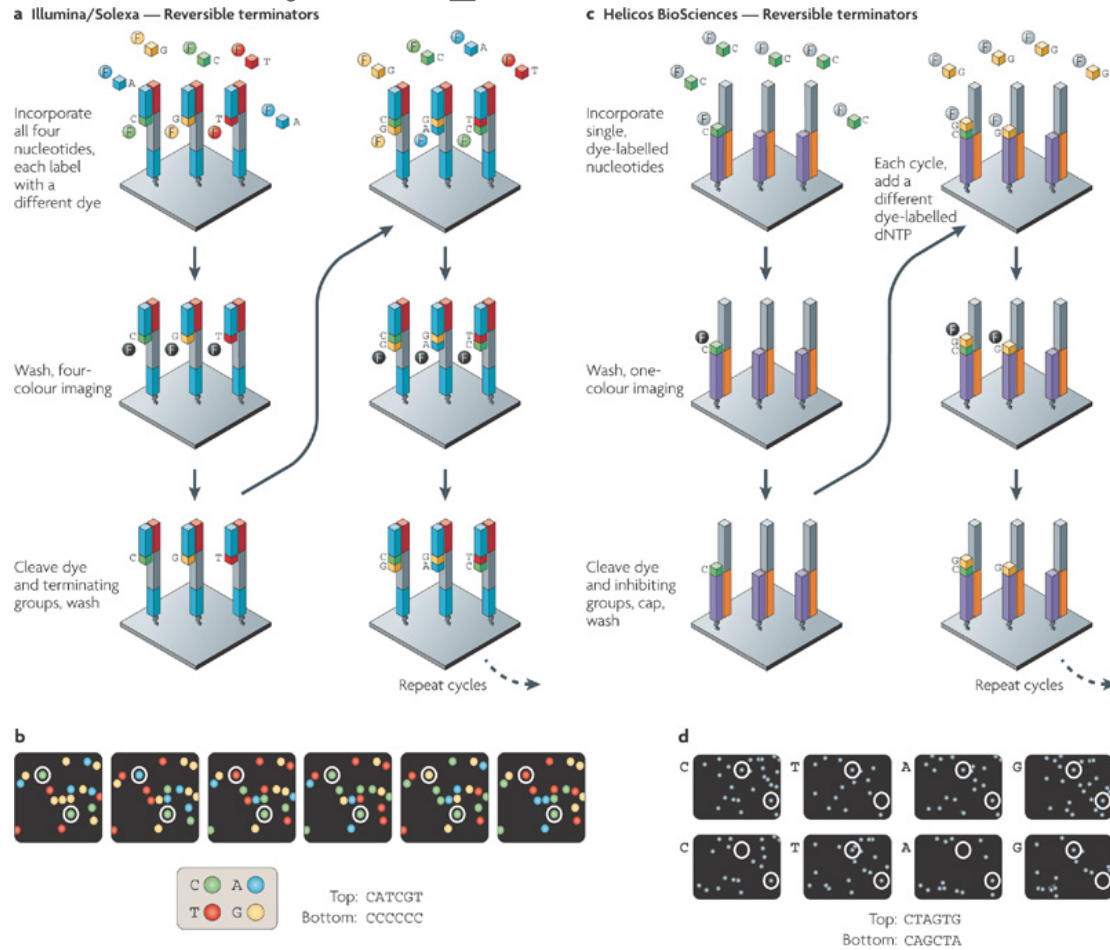
# Immobilisation du cible





Nature Reviews | Genetics

Metzker *Nat Rev Genet* 11 :31 (2010)

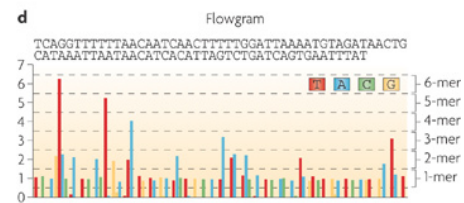
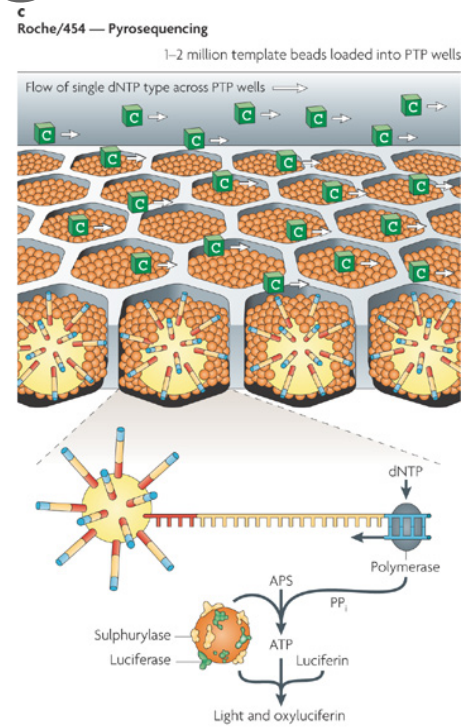
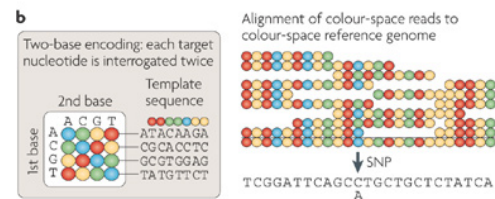
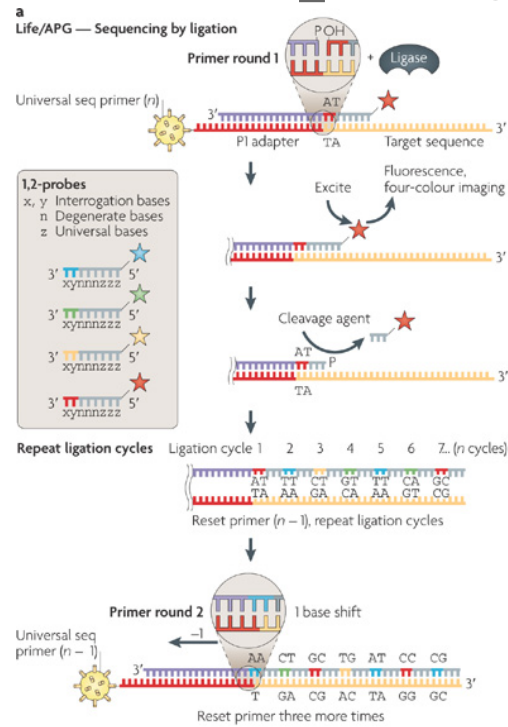
# Termination cyclique (Illumina)



Nature Reviews | Genetics

Metzker *Nat Rev Genet* 11 :31 (2010); →  

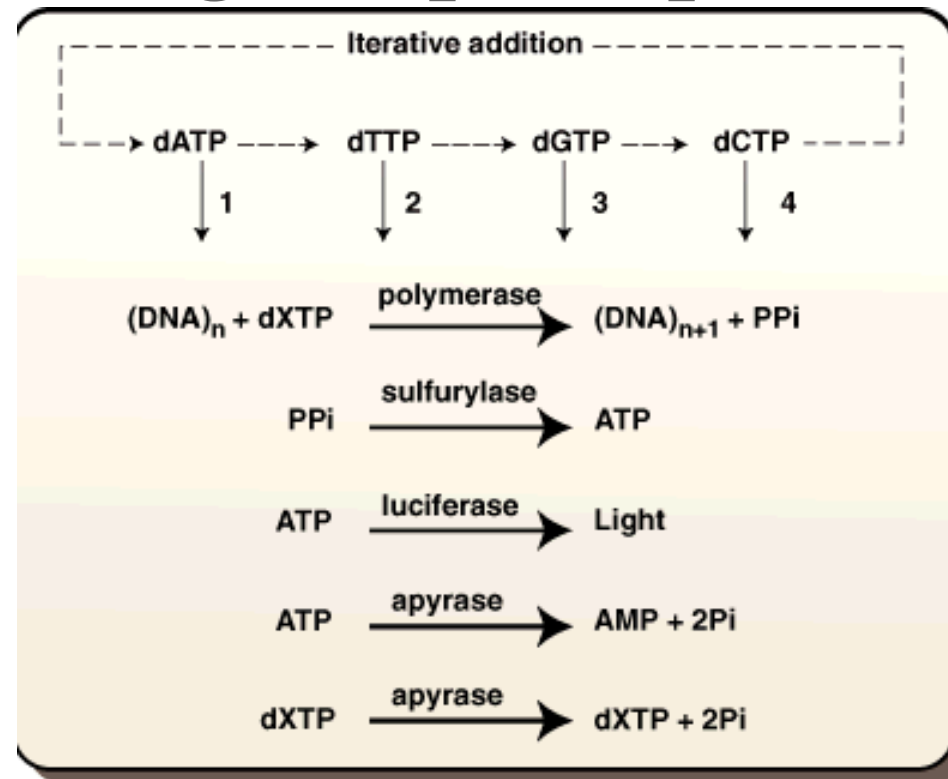
# Méthodes de séquençage



Nature Reviews | Genetics

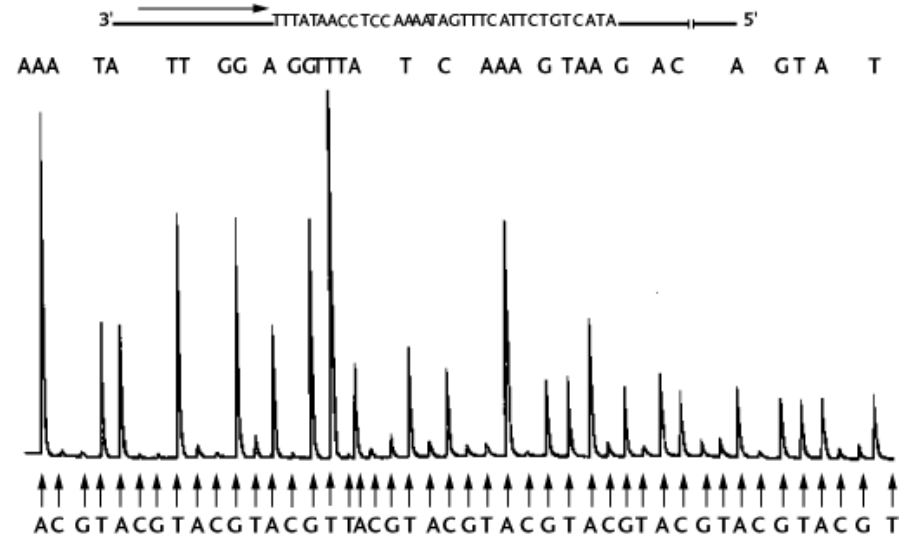
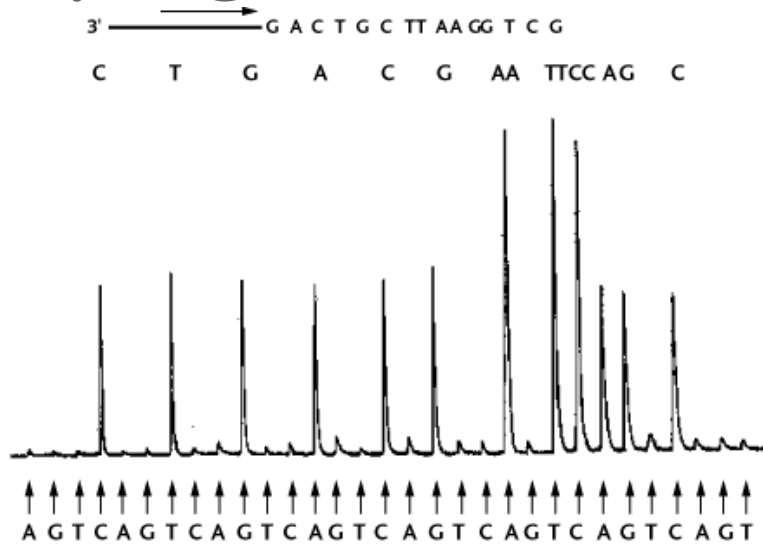
Metzker *Nat Rev Genet* 11 :31 (2010)

# Pyroséquençage — principe



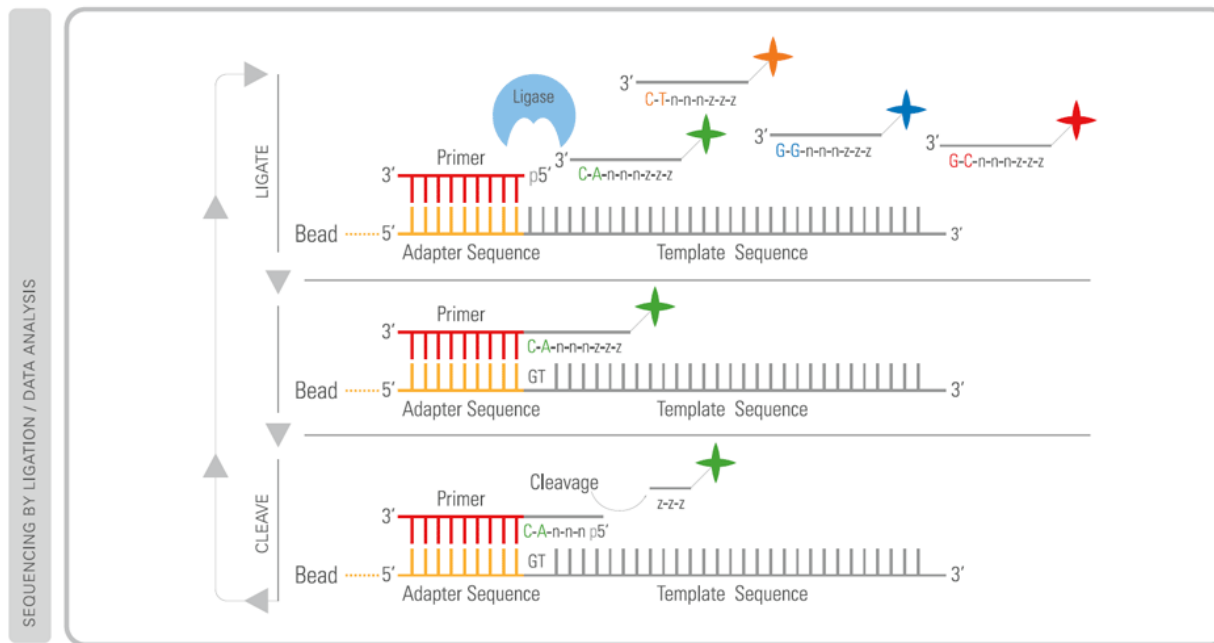
Ronaghi et al. *Science* 281 :363.

# Pyrogram



Ronaghi et al. *Science* 281 :363.

# Solid : sequencing by ligation



# Solid : encodage dans l'espace de couleurs

Séquence  $t[1..m + 1]$  encodée par  $t \cdot c[1..m]$ , où  $t = t[1]$ , et  $c[i] = t[i] \oplus t[i + 1]$ . L' $\oplus$  dénote l'OR exclusive dans l'encodage à 2 bits.

