

# Le séquençage nouvelle génération (NGS) en génétique humaine

Solena Le Scouarnec  
[solena.lescouarnec@univ-nantes.fr](mailto:solena.lescouarnec@univ-nantes.fr)

M2 Bioinformatique  
Module Bioinformatique appliquée

23 septembre 2013

# Objectifs

- **Partie 1**

Evolution des technologies de séquençage

Séquençage Illumina – obtention des données brutes

- **Partie 2 (Pierre Lindenbaum)**

Analyse bioinformatique

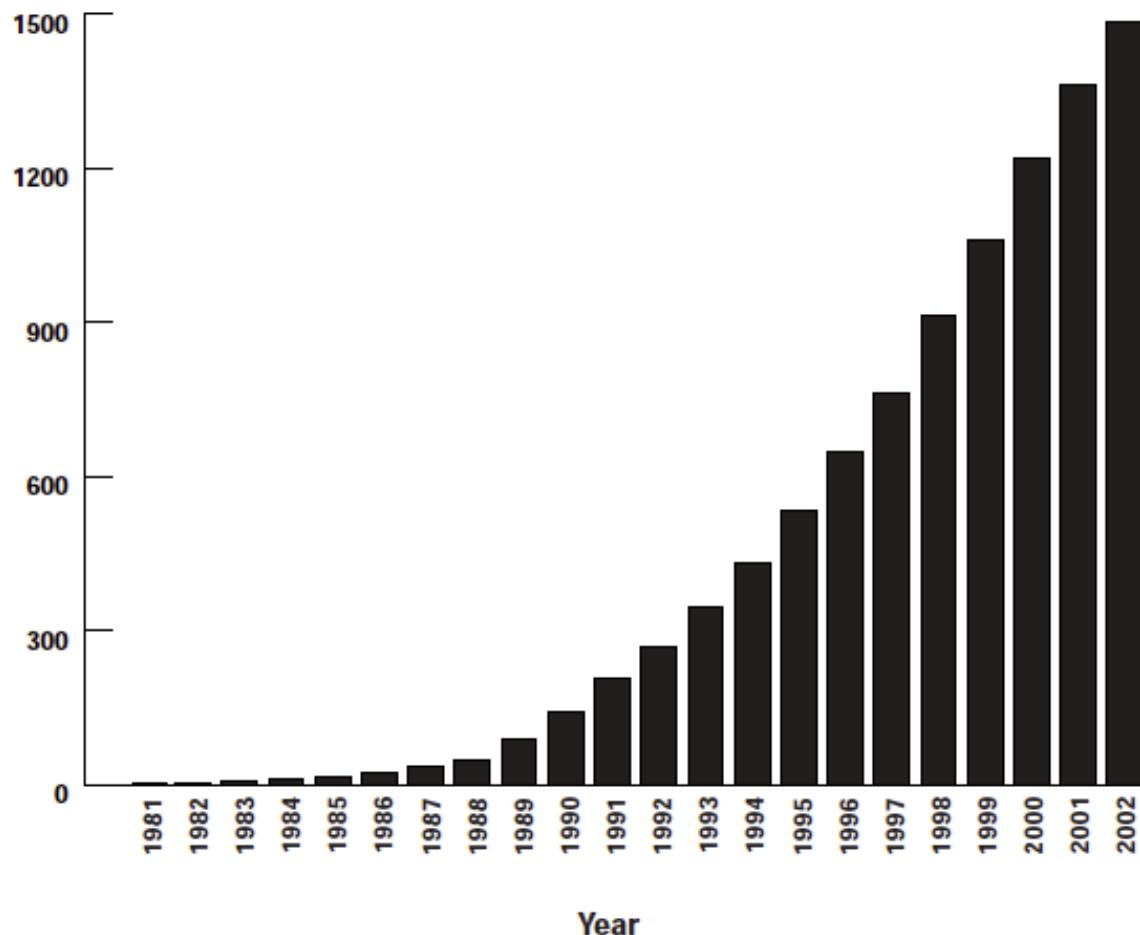
Format des fichiers

TP 9 & 16 octobre

# Plan

- Avant le NGS (méthode Sanger)
- La révolution NGS
- Depuis le NGS
- Séquençage avec la plateforme Illumina
- Exemple

# Maladies génétiques héréditaires



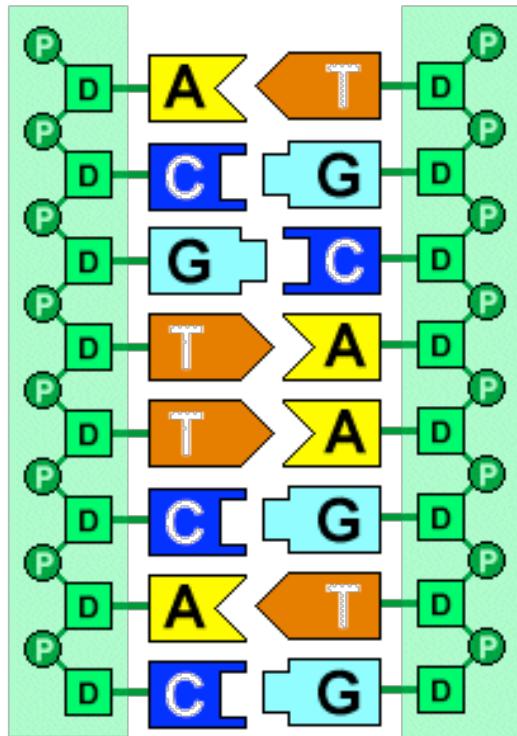
2002  
1485 gènes identifiés

2013  
3674 gènes identifiés

Base moléculaire inconnue pour >1000 maladies mendéliennes

Cumulative Pace of Disease Gene Discovery (1981-2002). The number of disease genes identified so far is 1,485. Data provided by Online Mendelian Inheritance in Man.

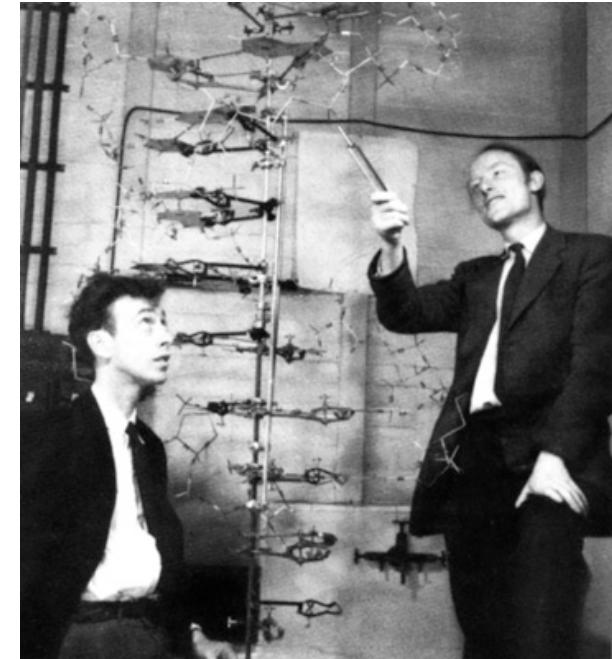
# L'ADN



A adénine  
T thymine  
C cytosine  
G guanine

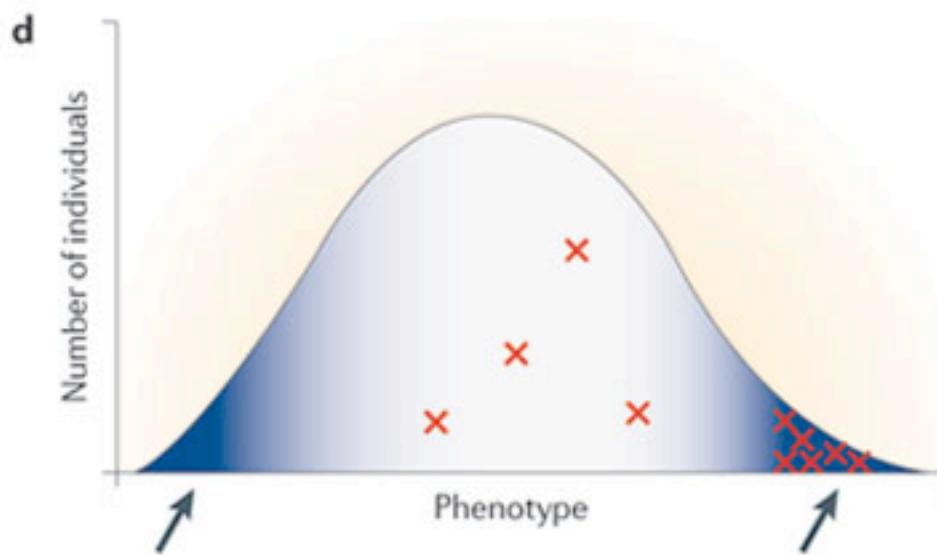
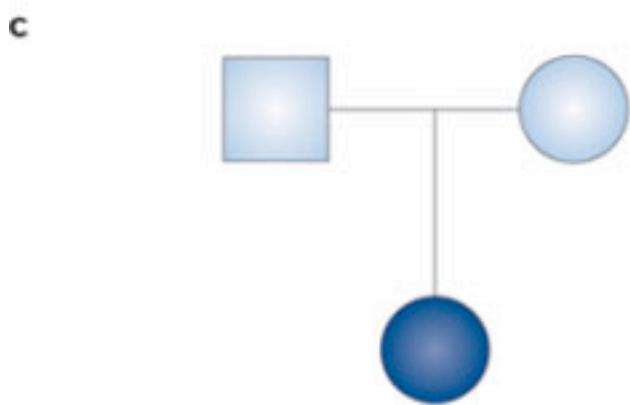
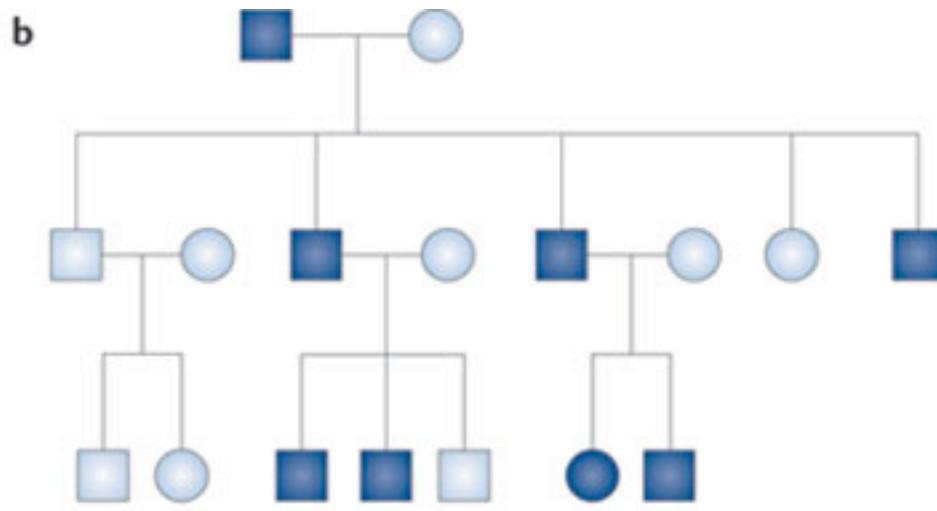
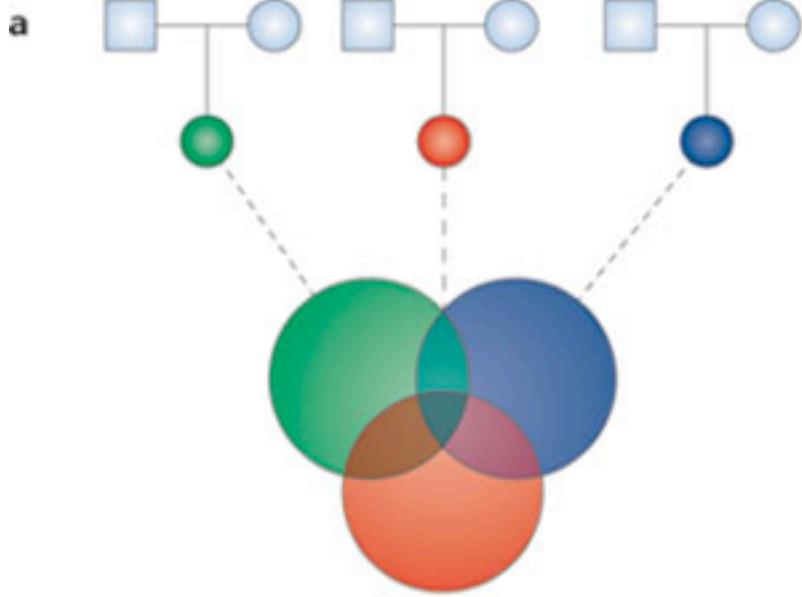
D P  
D désoxyribose  
acide phosphorique

P D A nucléotide



Séquençage d'ADN:

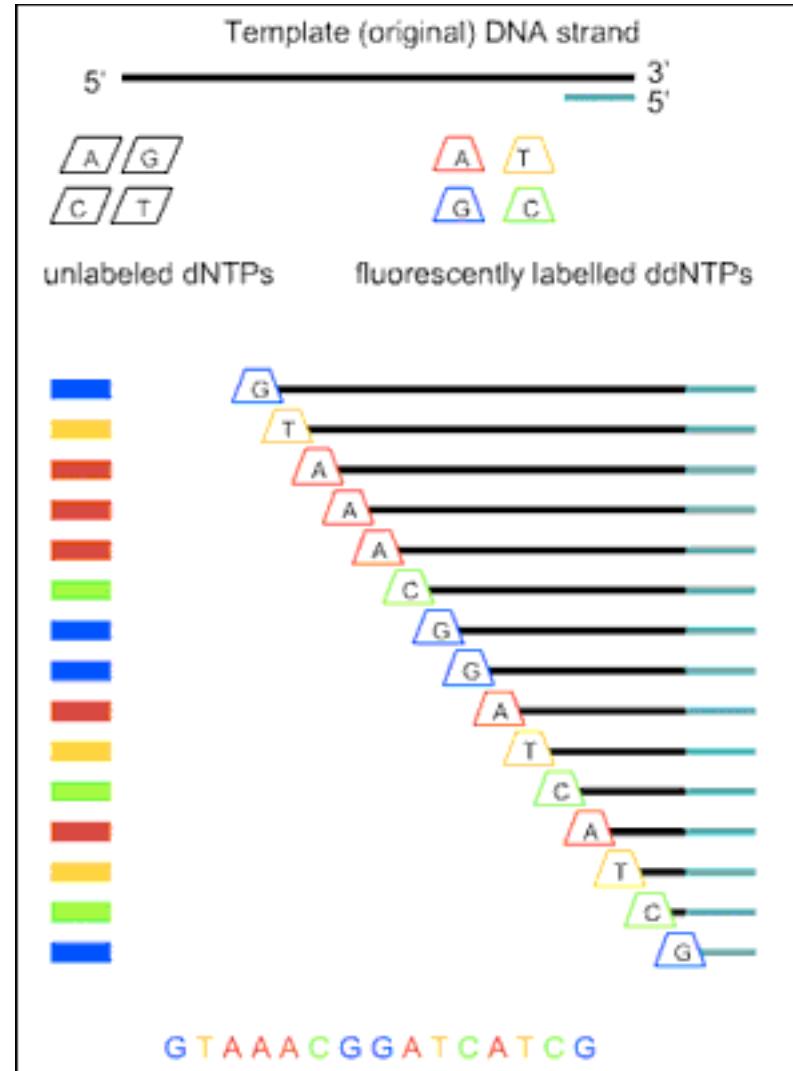
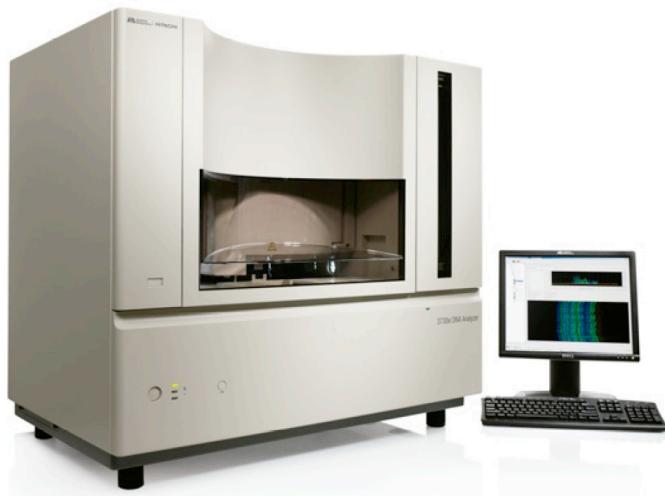
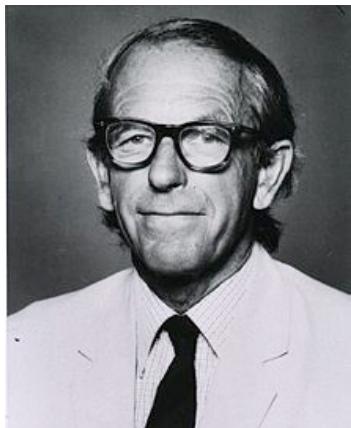
Déterminer l'ordre des nucléotides d'une molécule d'ADN



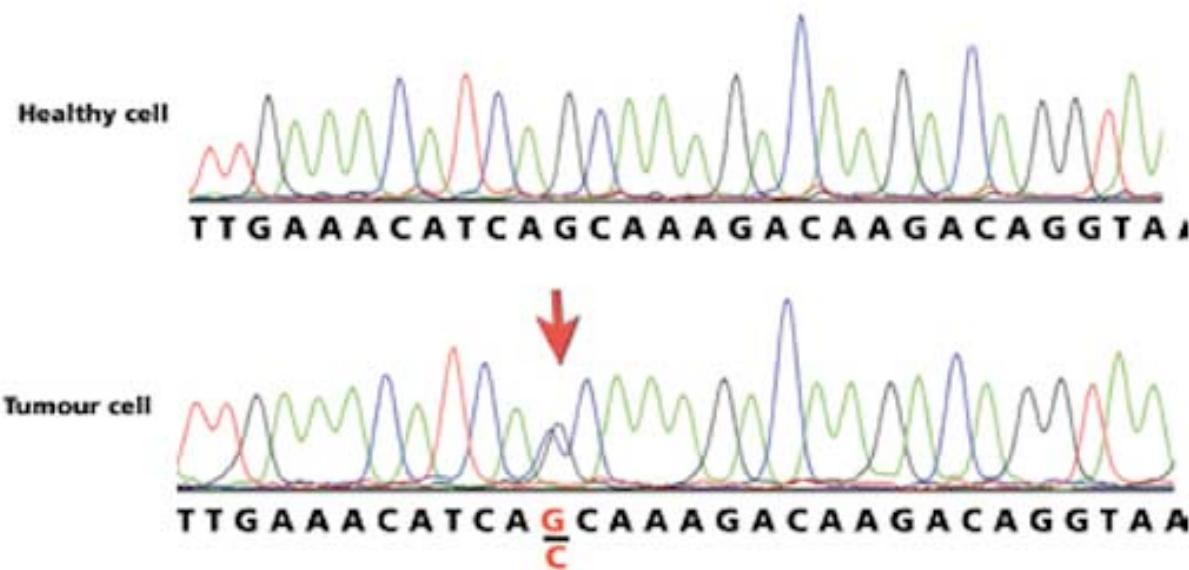
# Plan

- Avant le NGS (méthode Sanger)
- La révolution NGS
- Depuis le NGS
- Séquençage avec la plateforme Illumina
- Exemple

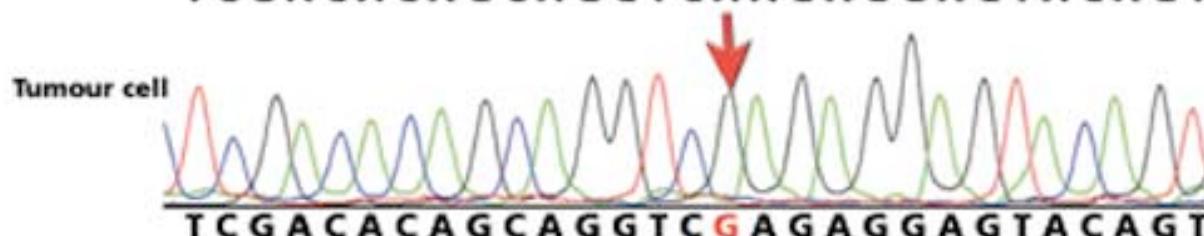
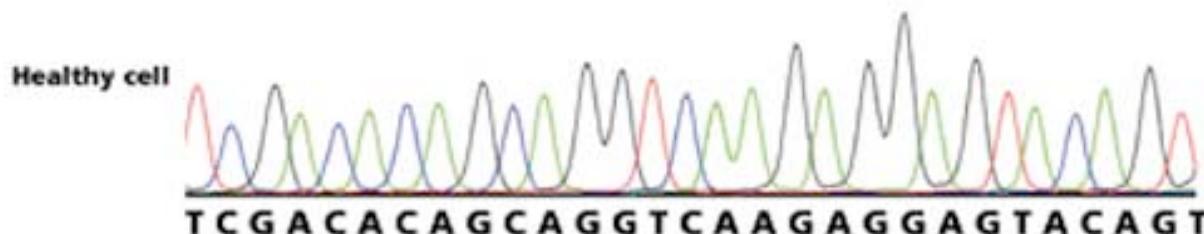
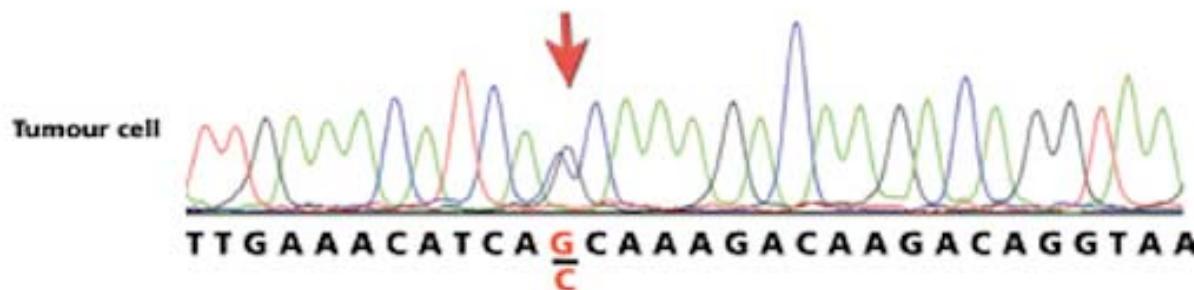
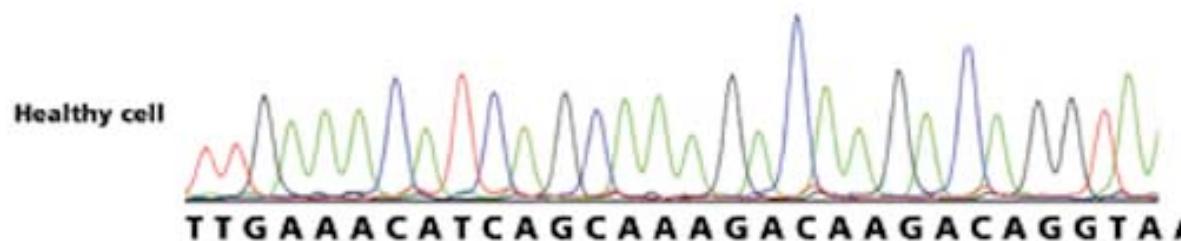
# Séquençage Sanger - Méthode -



# Séquençage Sanger - Analyse -

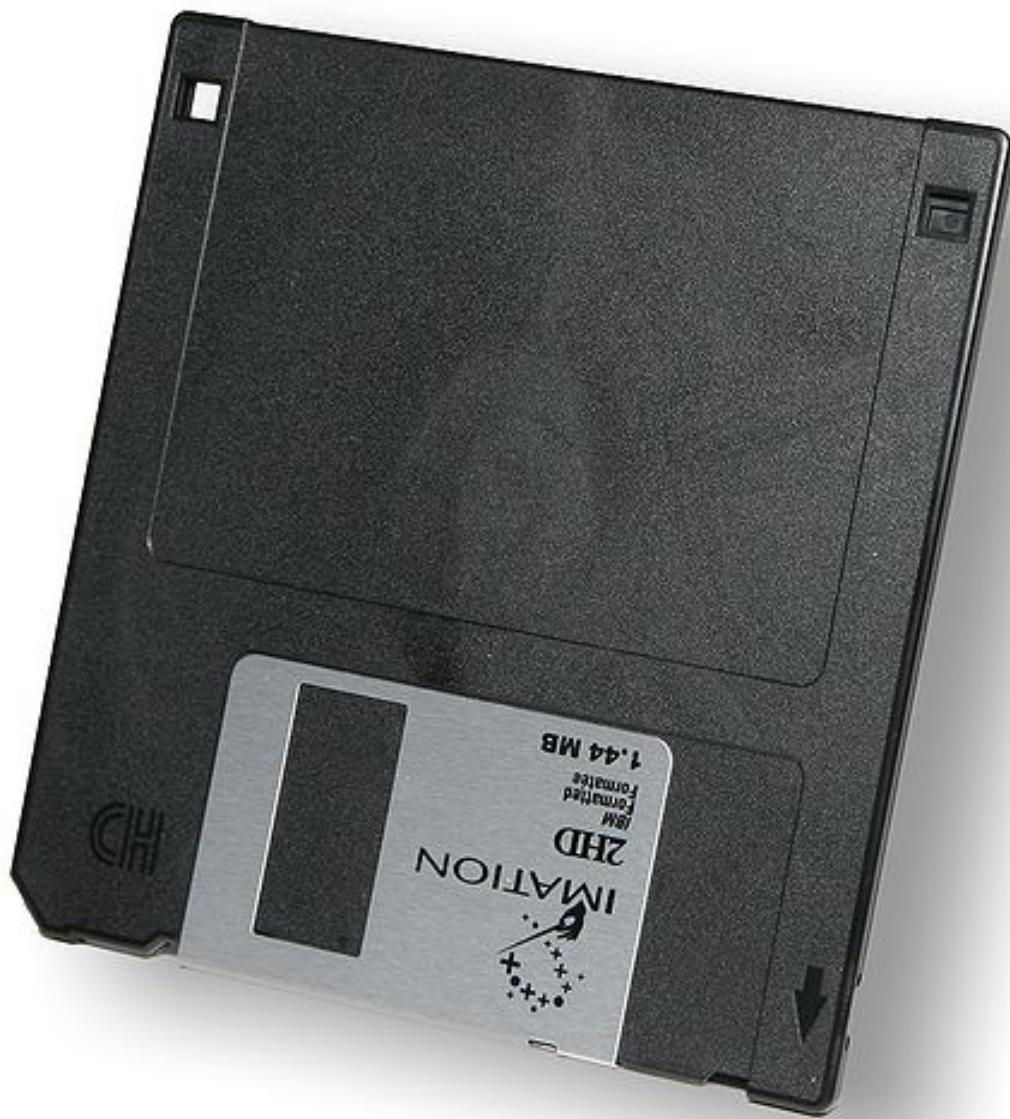


# Séquençage Sanger - Analyse -



Mutation  
hétérozygote

Mutation  
homozygote



# Séquençage Sanger

## Avantages

- Reads longs (1000 pb)
- Efficace pour petits projets
- Peu d'erreurs

## Limites

- Bas débit
- Coûteux

# Projet génome humain



Version la plus récente: GRCh37

# Plan

- Avant le NGS (méthode Sanger)
- La révolution NGS
- Depuis le NGS
- Séquençage avec la plateforme Illumina
- Exemple

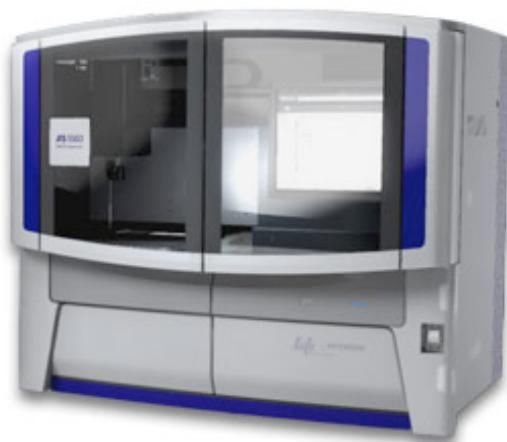
# Séquençage nouvelle génération (NGS)

2005



Roche GS FLX  
("pyroséquençage")

2006-2007



SOLiD 5500

2006-2007



HiSeq

# Séquençage Sanger / NGS

## Avantages

- Reads longs (1000 pb) / Reads plus courts
- Efficace pour petits projets / Petits & grands projets
- Peu d'erreurs / Erreurs de séquençage

## Limites

- Bas débit
- Coûteux

# Séquençage nouvelle génération (NGS)

Table 2 Next-generation DNA sequencing instruments

	Cost per base <sup>a</sup>	Read length (bp) <sup>b</sup>	Speed	
<b>Minimum cost per base</b>				
Complete Genomics	Low	Short	3 months	
HiSeq 2000 (Illumina)	Low	Mid	8 days	- coûteux
SOLID 5500xl (Life Technologies)	Low	Short	8 days	
<b>Maximum read length</b>				
454 GS FLX+ (Roche)	High	Long	1 day	séquences +
RS (Pacific Biosciences)	High	Very long	<1 day	longues





[http://www.flickr.com/photos/esquimo\\_2000/5241744434/](http://www.flickr.com/photos/esquimo_2000/5241744434/)

# Plan

- Avant le NGS (méthode Sanger)
- La révolution NGS
- Depuis le NGS
- Séquençage avec la plateforme Illumina
- Exemple

# Séquençage de génomes personnels

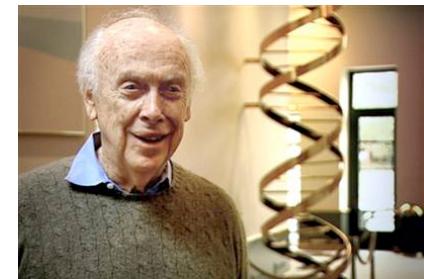
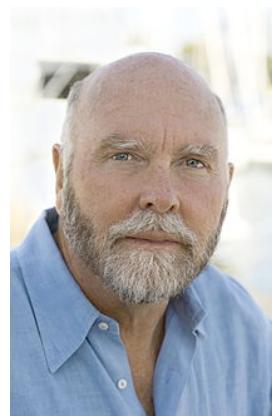
QUICKER, SMALLER, CHEAPER

Sanger

Sanger

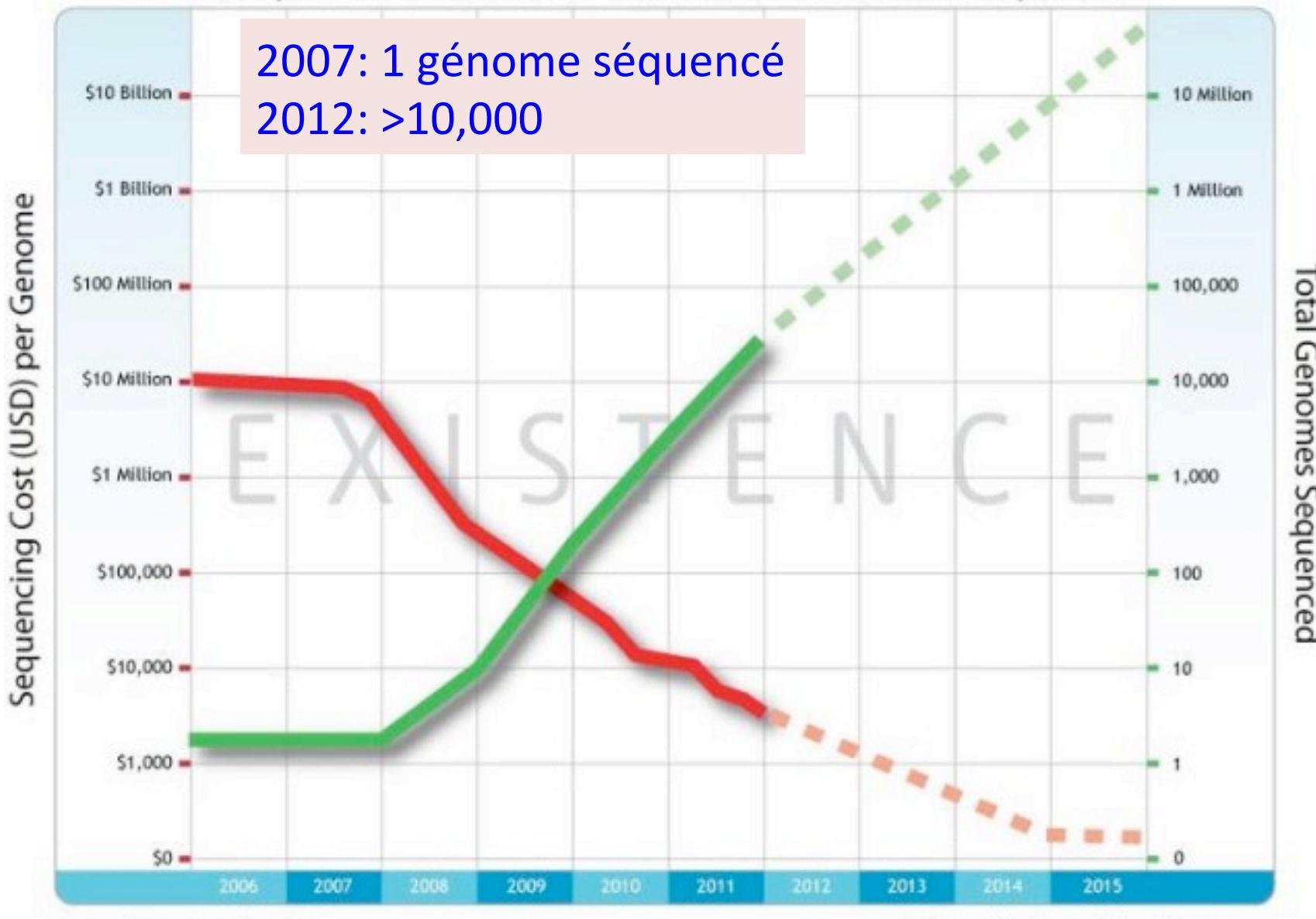
454

Genome sequenced (publication year)	HGP (2003)	Venter (2007)	Watson (2008)
Time taken (start to finish)	13 years	4 years	4.5 months
Number of scientists listed as authors	> 2,800	31	27
Cost of sequencing (start to finish)	\$2.7 billion	\$100 million	< \$1.5 million
Coverage	8-10 ×	7.5 ×	7.4 ×
Number of institutes involved	16	5	2
Number of countries involved	6	3	1

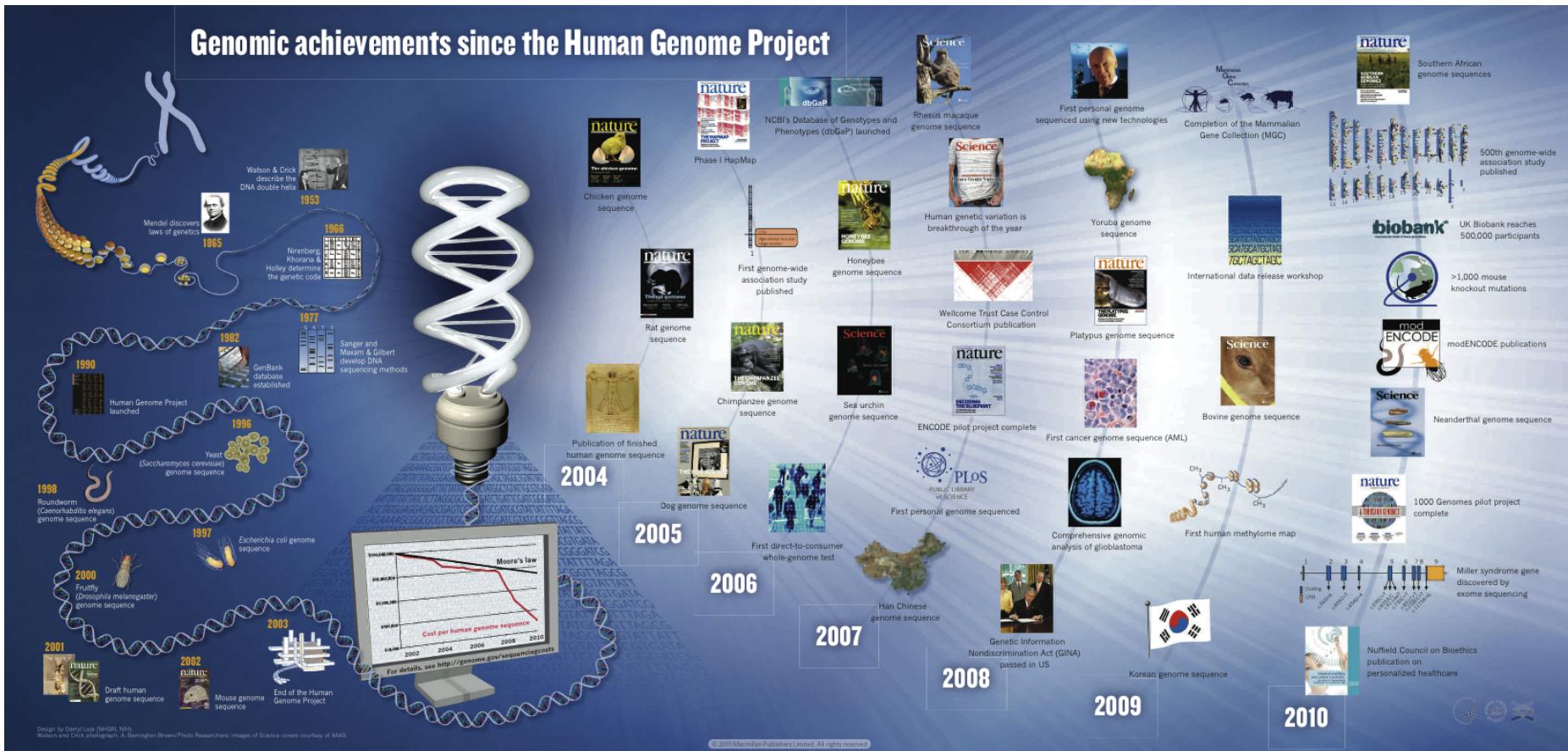


# Full Genome Sequencing & The Genetic Revolution

Cost per Human Genome vs Total Number of Genomes Sequenced



# Genomic achievements since the Human Genome Project



**nature**

Eric D. Green, Mark S. Guyer & National Human Genome Research Institute  
Nature 470, 204–213 (10 February 2011)

# Exemples

- Génome (ADN)

Séquençage *de novo*

Génome entier

Séquençage ciblé (exome, gènes...)

- Transcriptome (ARN)

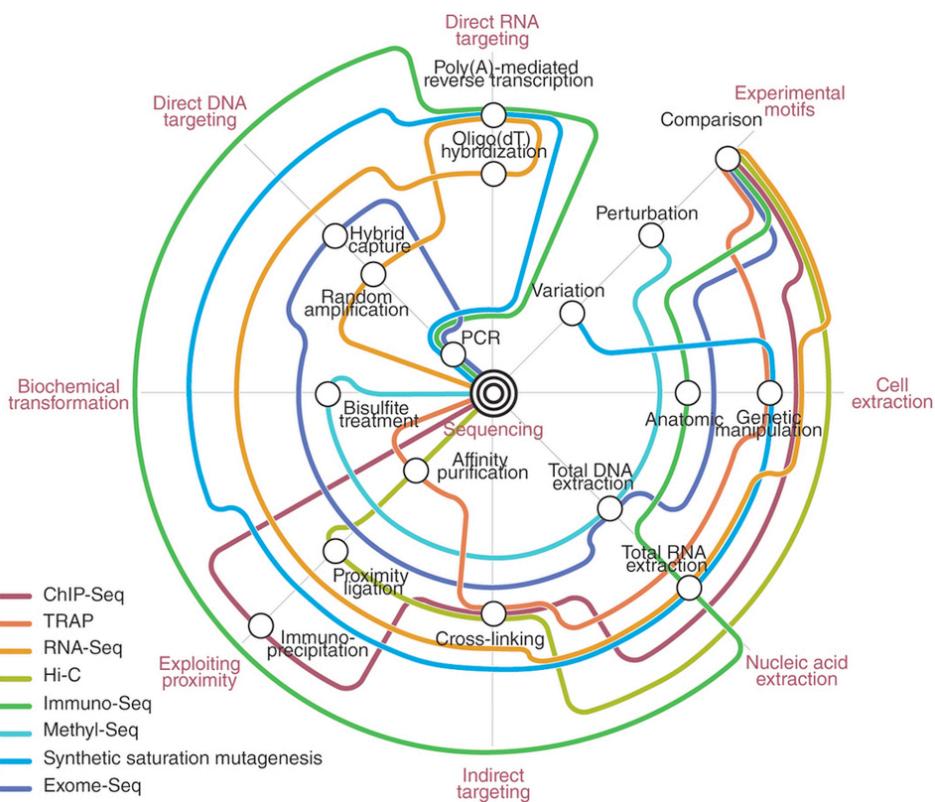
mRNA-Seq

Small RNA-Seq (miRNA)

- Autres

ChIP-Seq

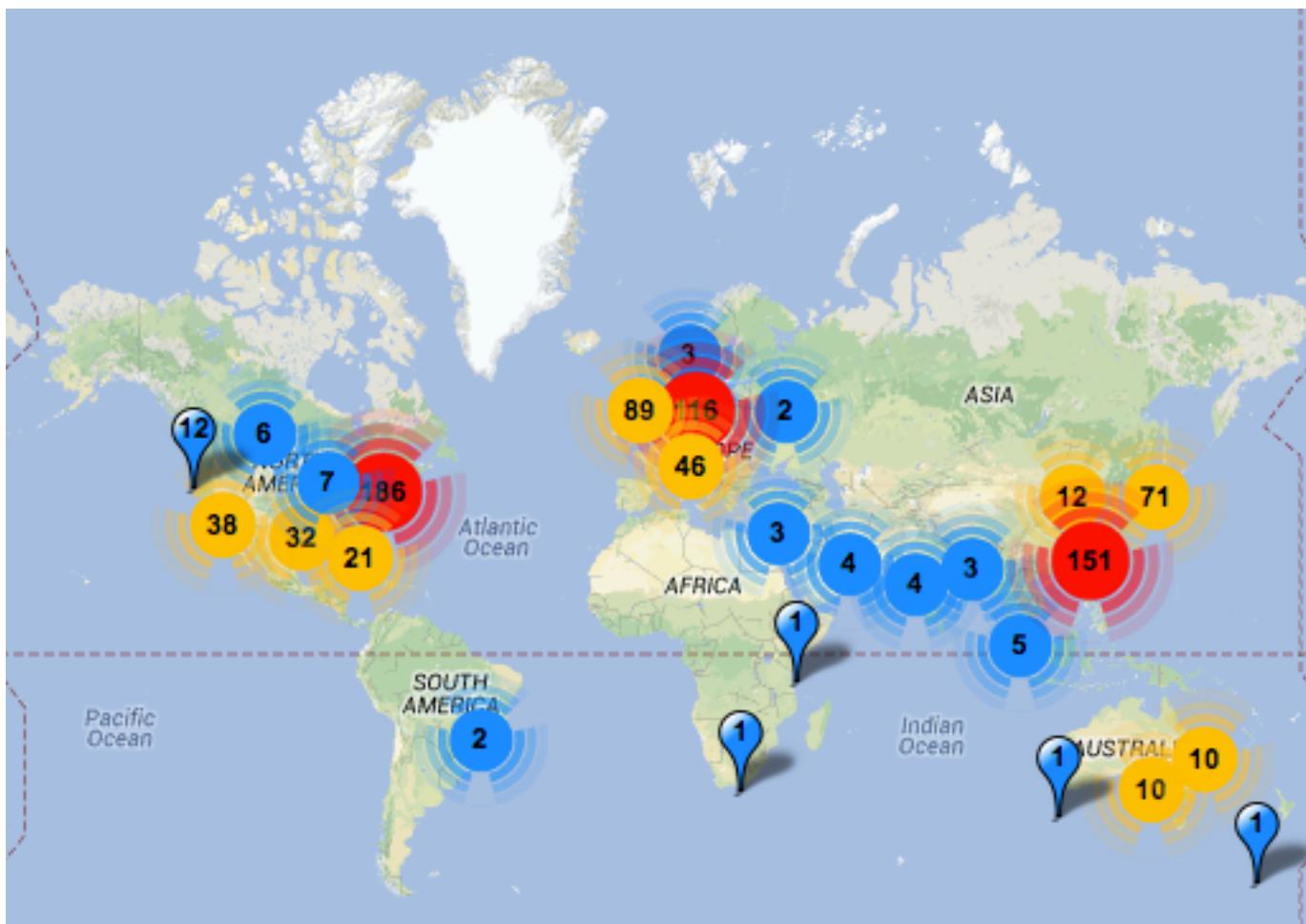
Methyl-Seq



# Plan

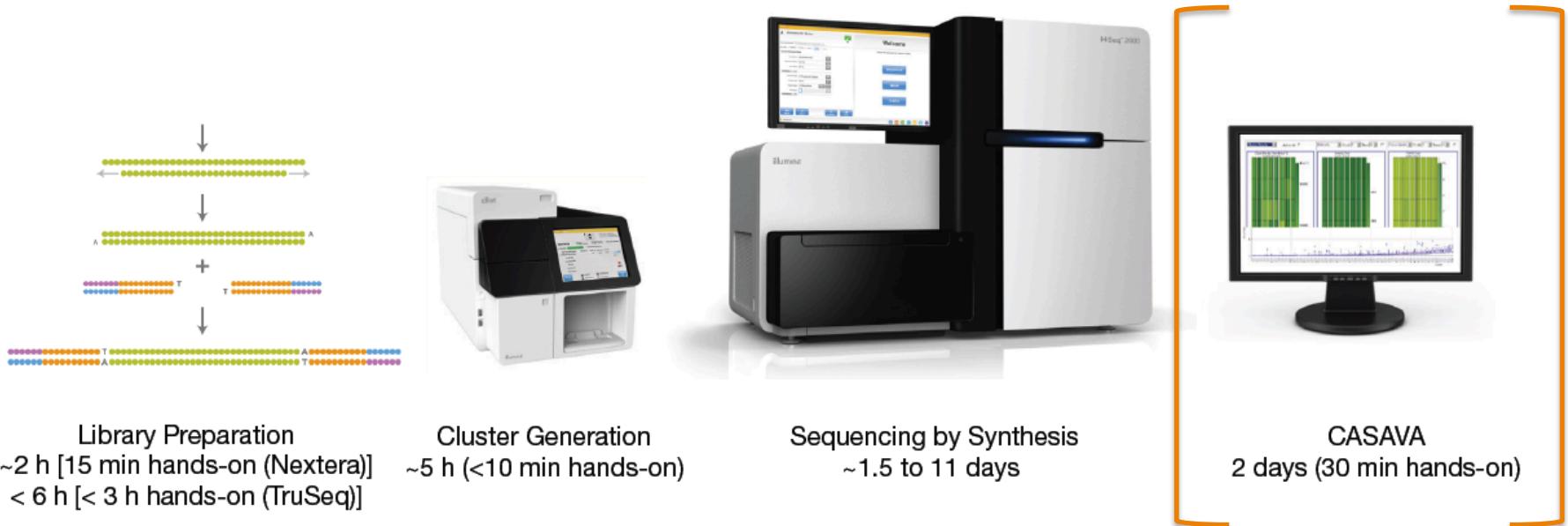
- Avant le NGS (méthode Sanger)
- La révolution NGS
- Depuis le NGS
- Séquençage avec la plateforme Illumina
- Exemple

# HiSeq (Illumina)



# Overview

Figure 3: Next-Generation Sequencing Simplified



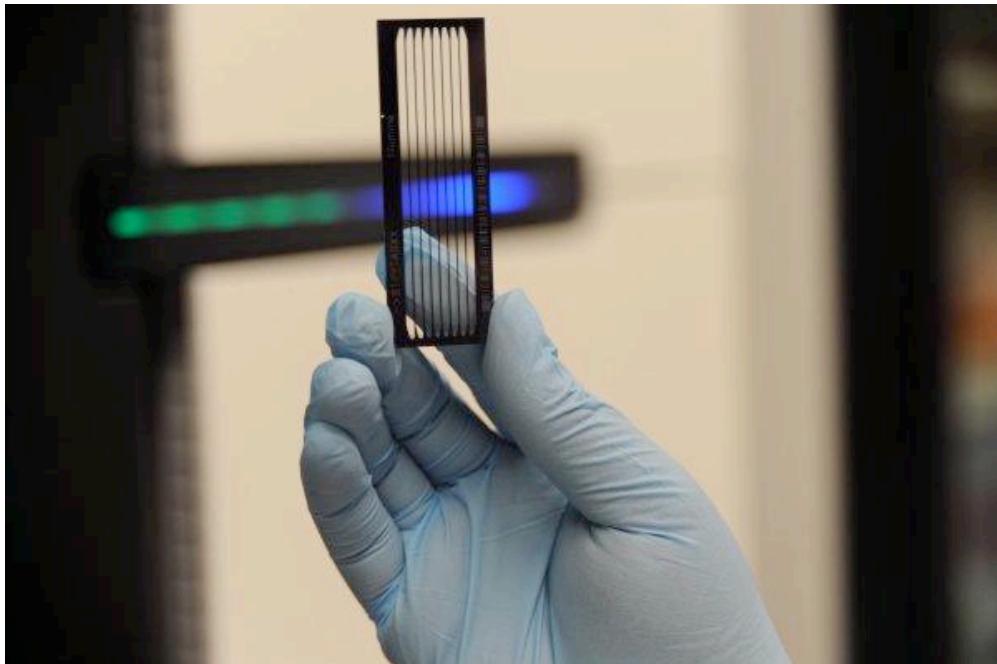
Préparation de  
la librairie

Génération  
de clusters

Séquençage  
=> reads

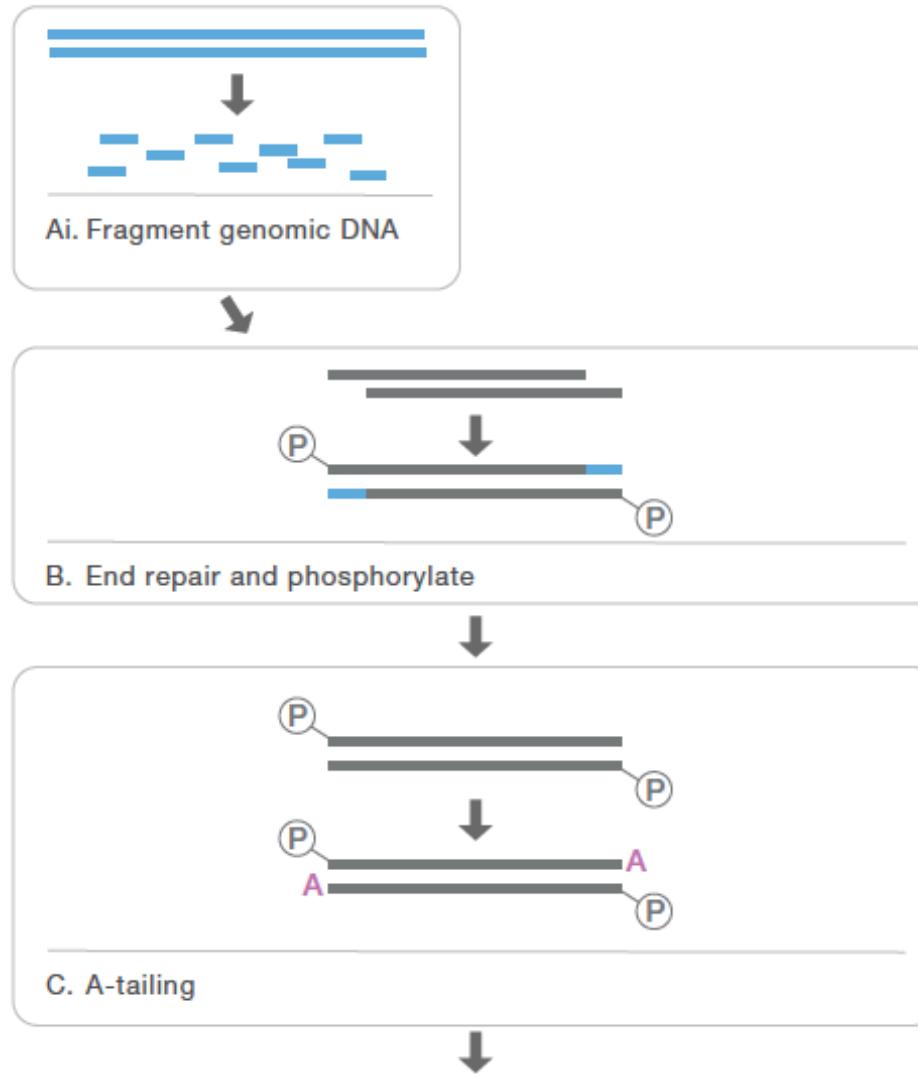
Bioinformatique  
+++

# Flow cell

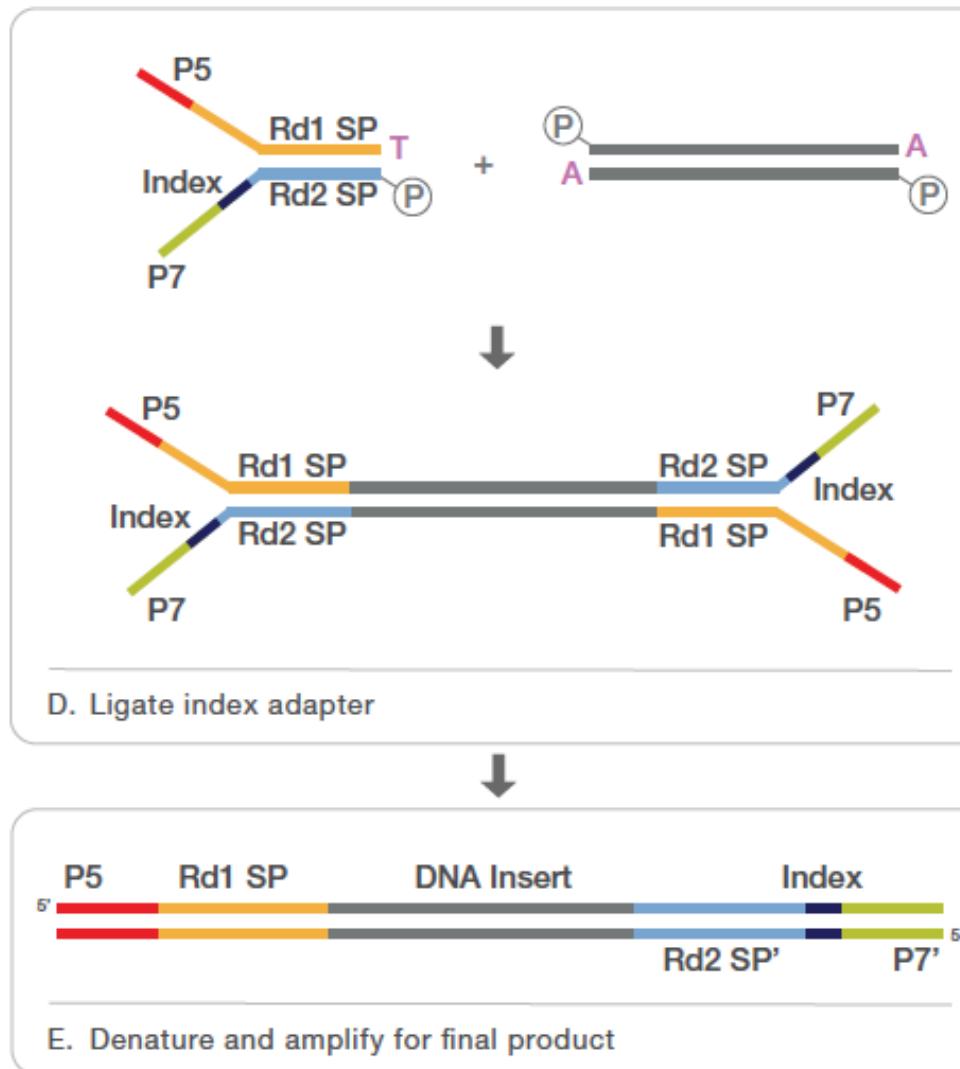


1 flow cell  
=8 lanes

# 1/ Préparation de librairie



# 1/ Préparation de librairie



# 2/ cBOT

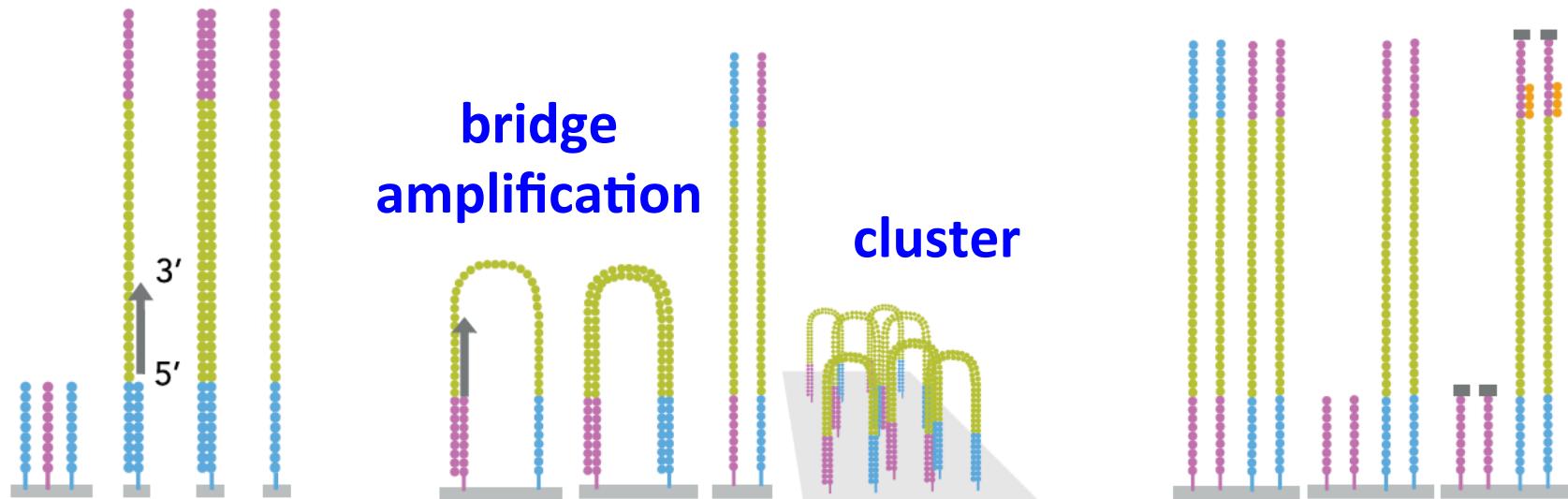
Figure 2: New cBot Features Enables Rapid and Streamlined Cluster Generation



The cBot cluster generation system is the next generation of workflow improvements for Illumina sequencing. Novel innovations include pre-packaged reagents, a single manifold, advanced fluidics and thermal stage features, integrated sensors, remote monitoring capabilities, and simplified data entry and tracking with the touch screen and barcode scanner.

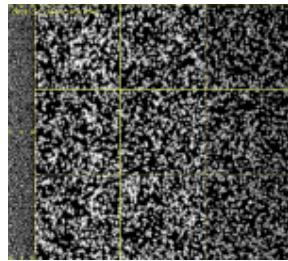
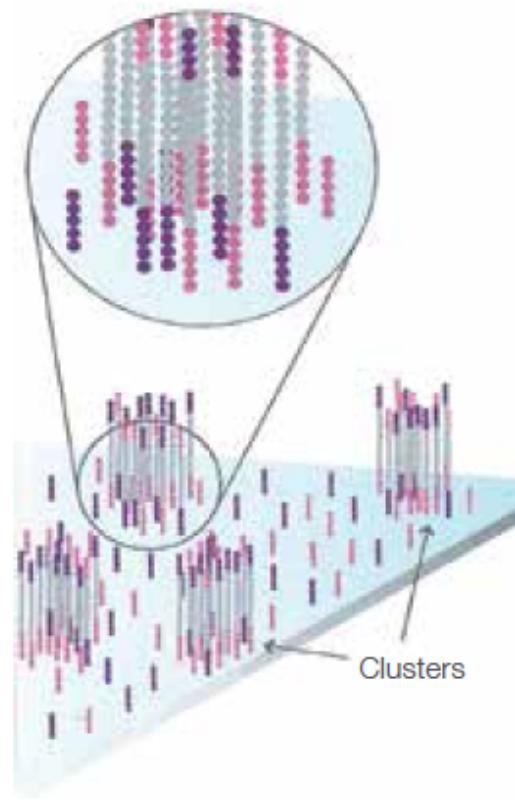
## 2/ Génération de clusters

Figure 3: Cluster Generation by Isothermal Bridge Amplification



Cluster generation from single-molecule DNA templates occurs within the sealed Illumina flow cell on the cBot instrument, and involves immobilization and 3' extension, bridge amplification, linearization, and hybridization.

## 2/ Génération de clusters



100-200 millions de clusters/lane

# 3/ HiSeq

Touch screen user interface facilitates step-by-step run setup. Simply enter read length, single- or paired-end read, and indexing information on-screen.



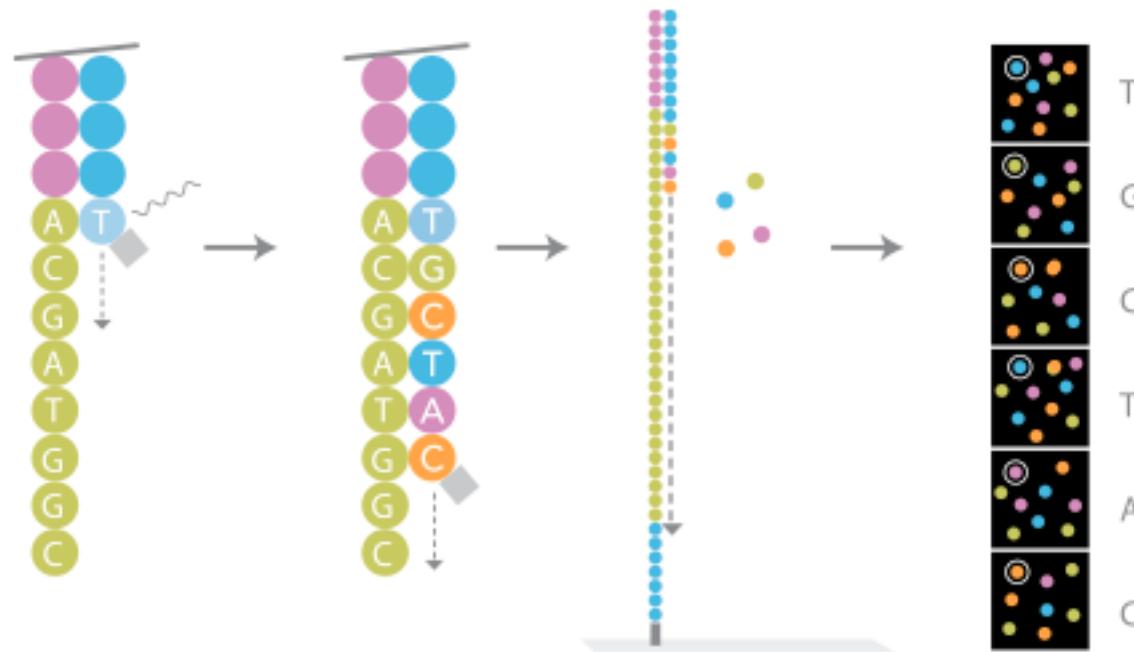
# 3/ Séquençage

Cycle 1



### 3/ Séquençage

Cycle 1   ...   Cycle 6   ...



→ reads

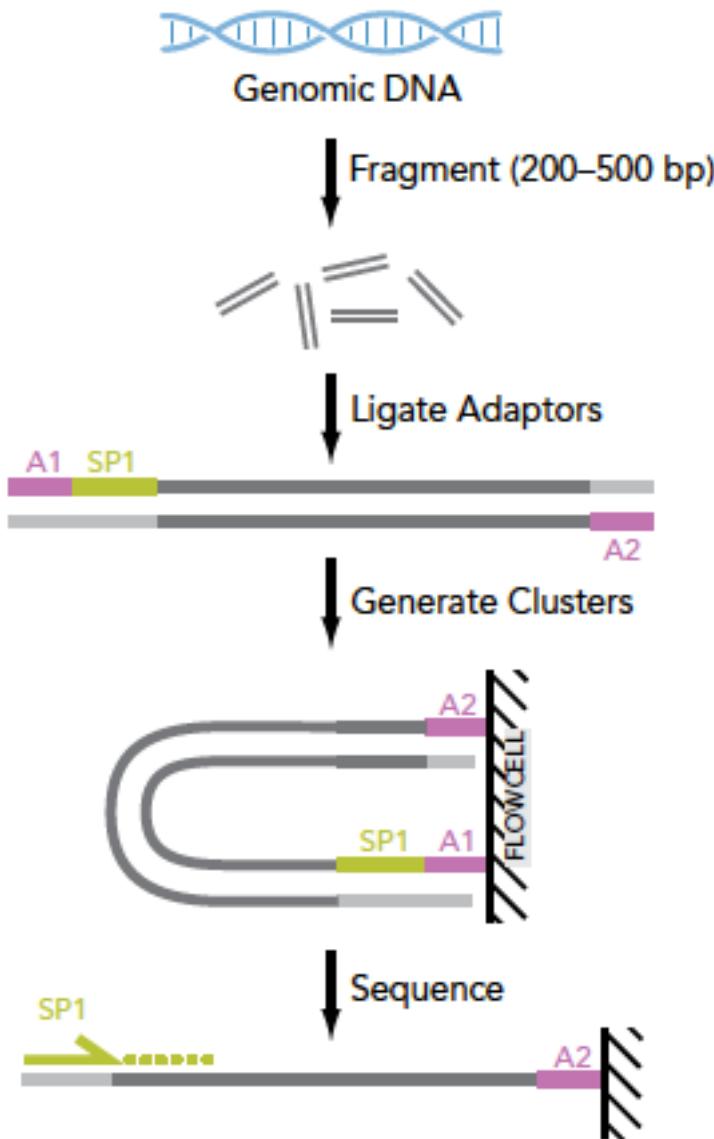
(ex. 50 pb, 75 pb ou 100 pb)

Chaque base d'un read est associée à une valeur de qualité

## Animation

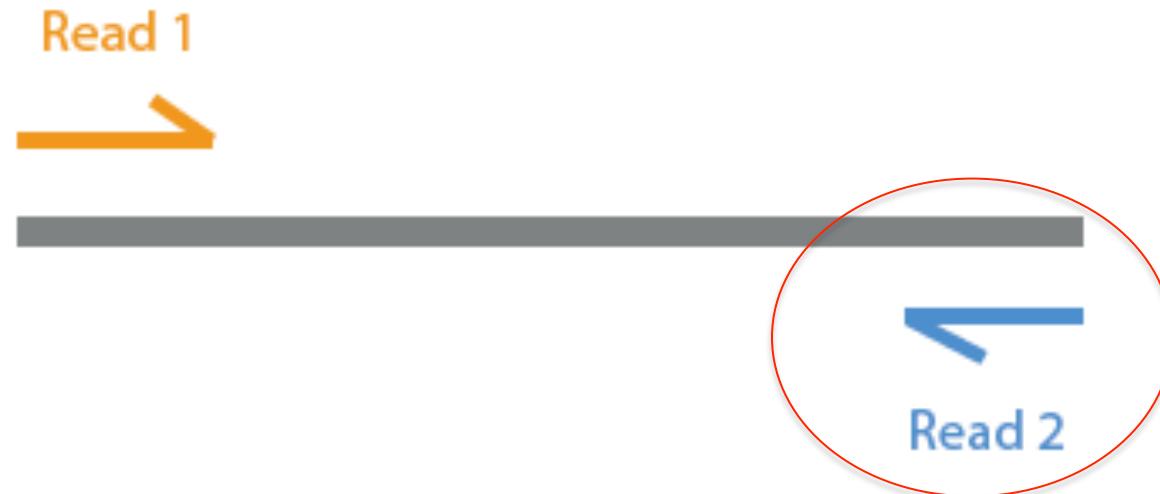
<http://www.yourgenome.org/teachers/speed.shtml>

# Single-end sequencing

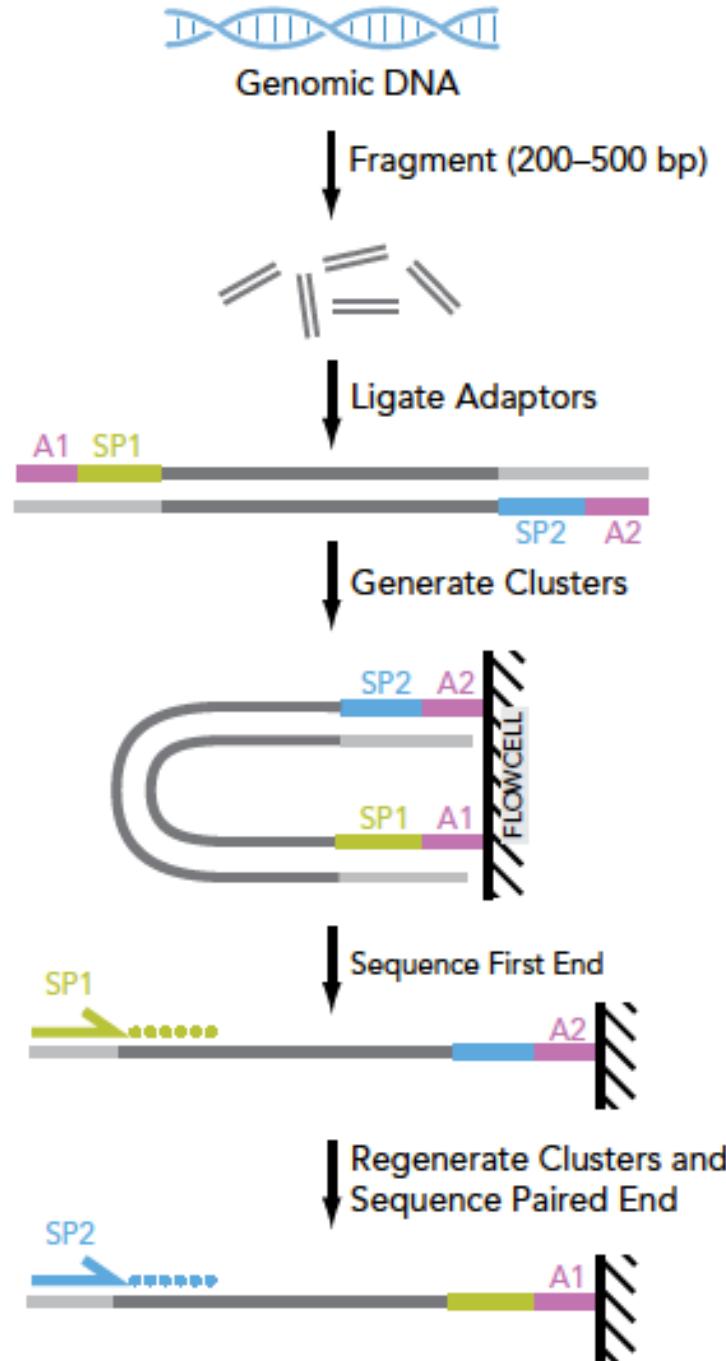


# Paired-end sequencing (1)

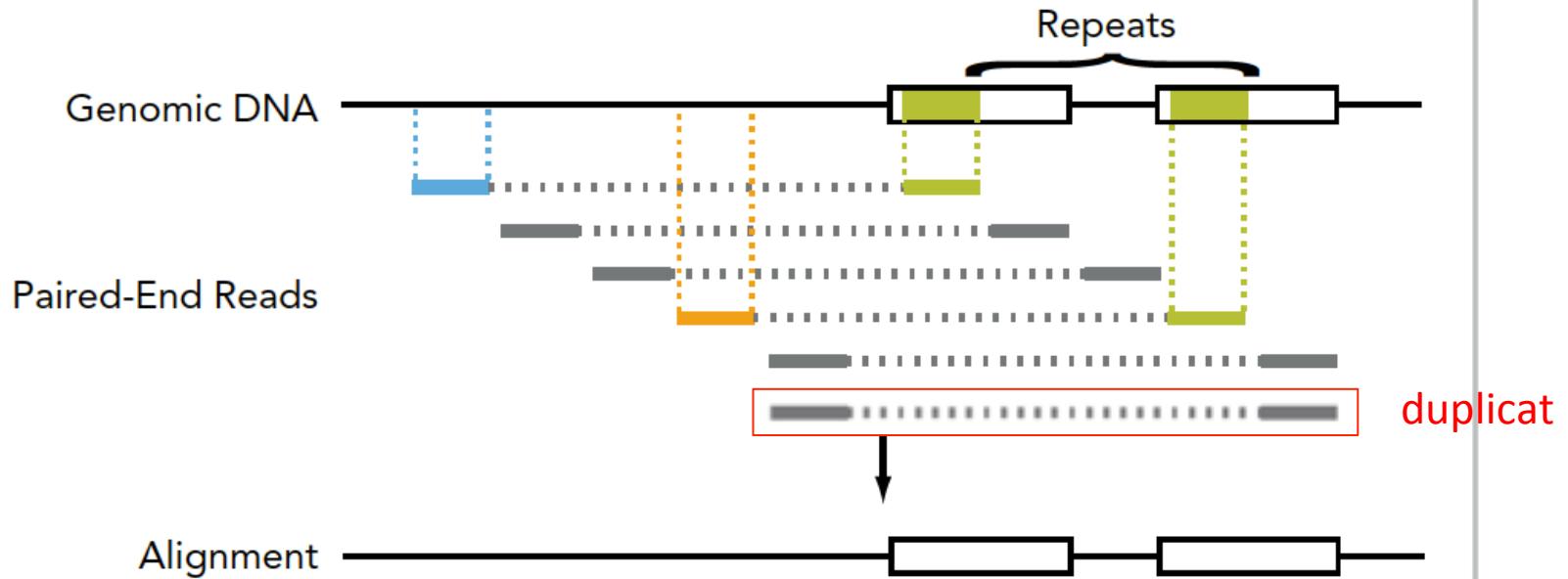
Paired-End Reads



# Paired-end sequencing (2)



## Figure 2: Unique Alignment Of Paired Reads In Repeats

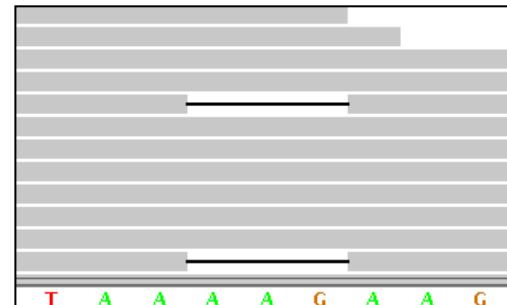
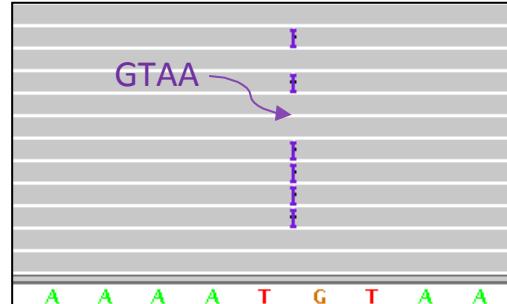
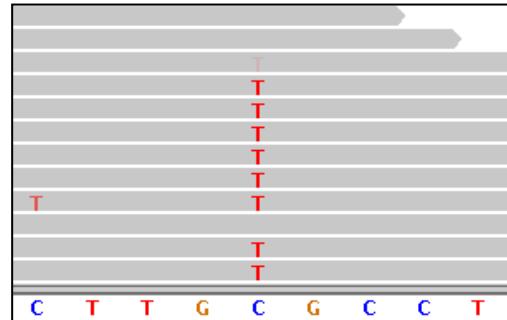


Reads in repeats (green) can be unambiguously aligned in complex genomes. Each read is associated with a paired read (blue or orange) and the separation between read pairs is known from the fragment size of the input DNA.

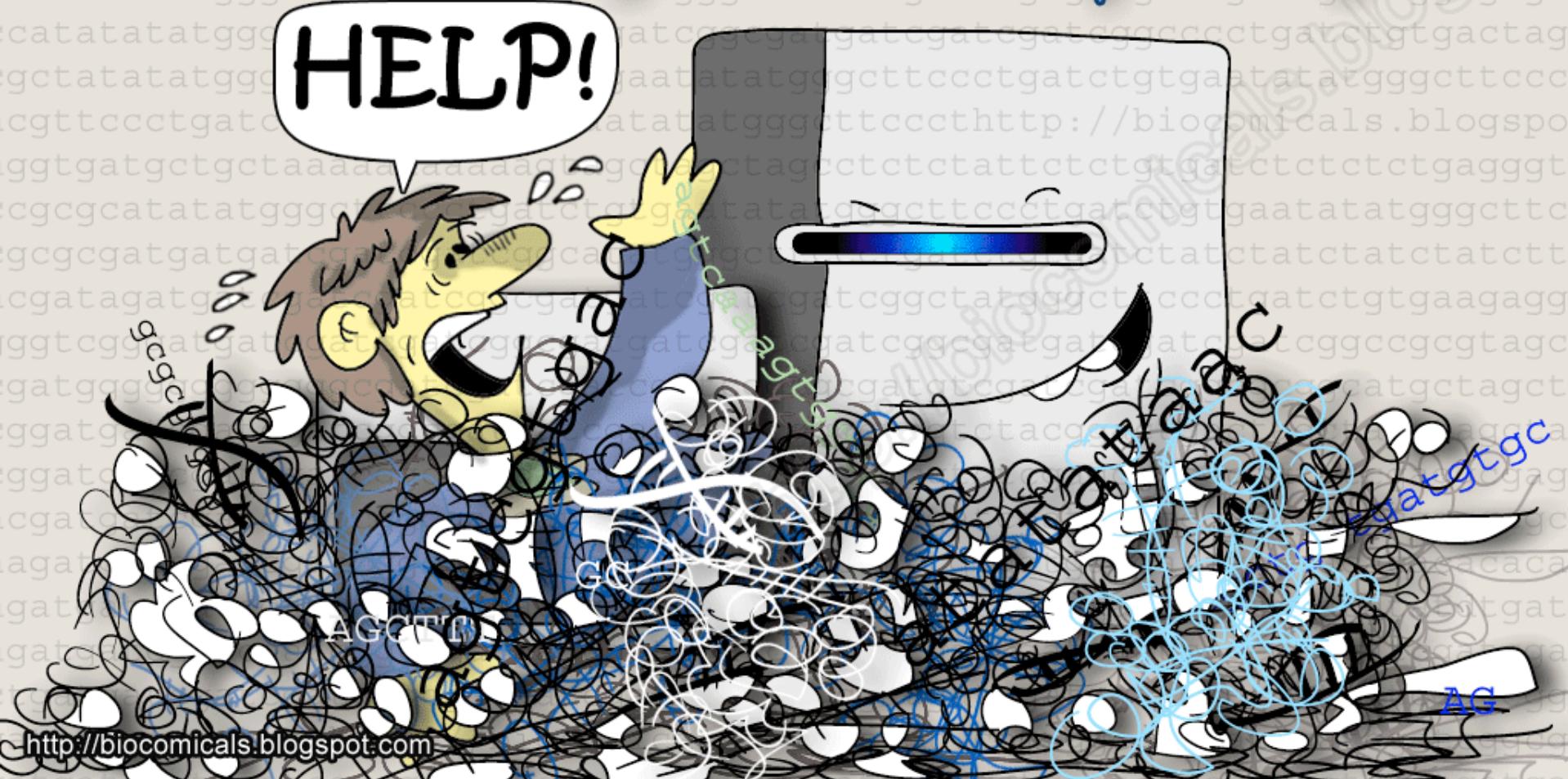
Table 1: Flexible Paired Sequencing Provides Optimal Detection Of Any Variant

Variant	Single Read	Short Insert Paired-Ends (200–500 bp)	Long Insert Mate Pairs (2–5 kb)	Paired-End And Mate Pair Combined
SNP	++	++++	++	++++
Small indels	++	++++	++	++++
Insertion	+	+++	+++	++++
Amplification	++	+++	+++	++++
Deletion	+	+++	++	++++
Inversion	+	+++	++	++++
Complex rearrangement	+	+++	++	++++
Large rearrangement	+	++	+++	++++

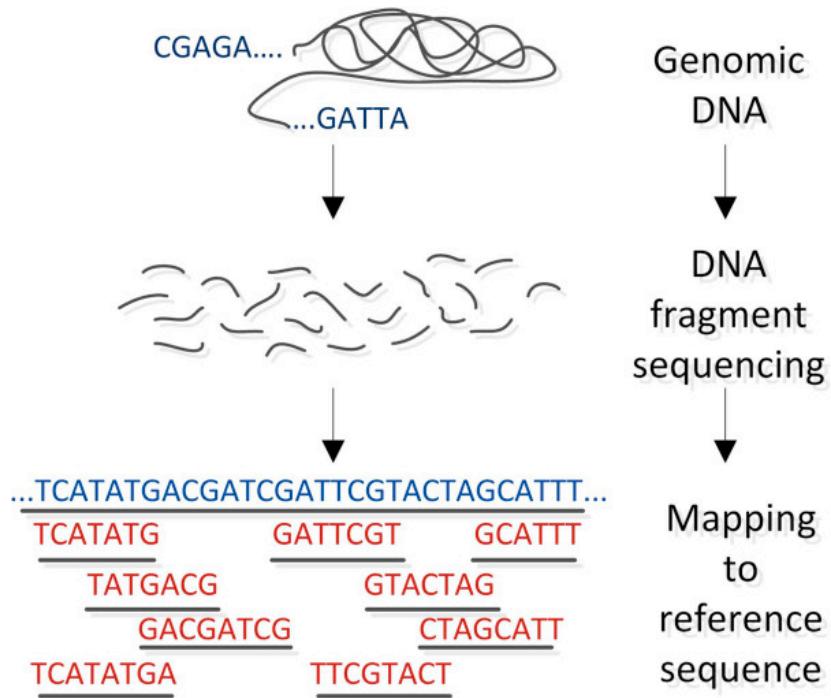
Only by combining short and long inserts can researchers be certain to find all different sizes and types of variants. In particular, short inserts are essential to identifying small indels and mate pairs are essential for identifying the largest rearrangements.



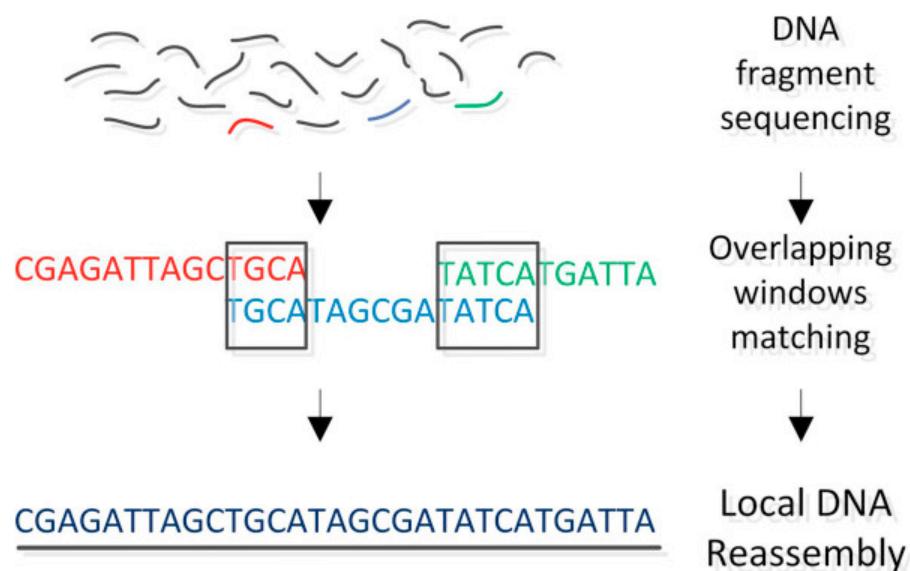
# Drowned in next generation sequencing data



## Génome de référence existant => alignement



## Assemblage *de novo*

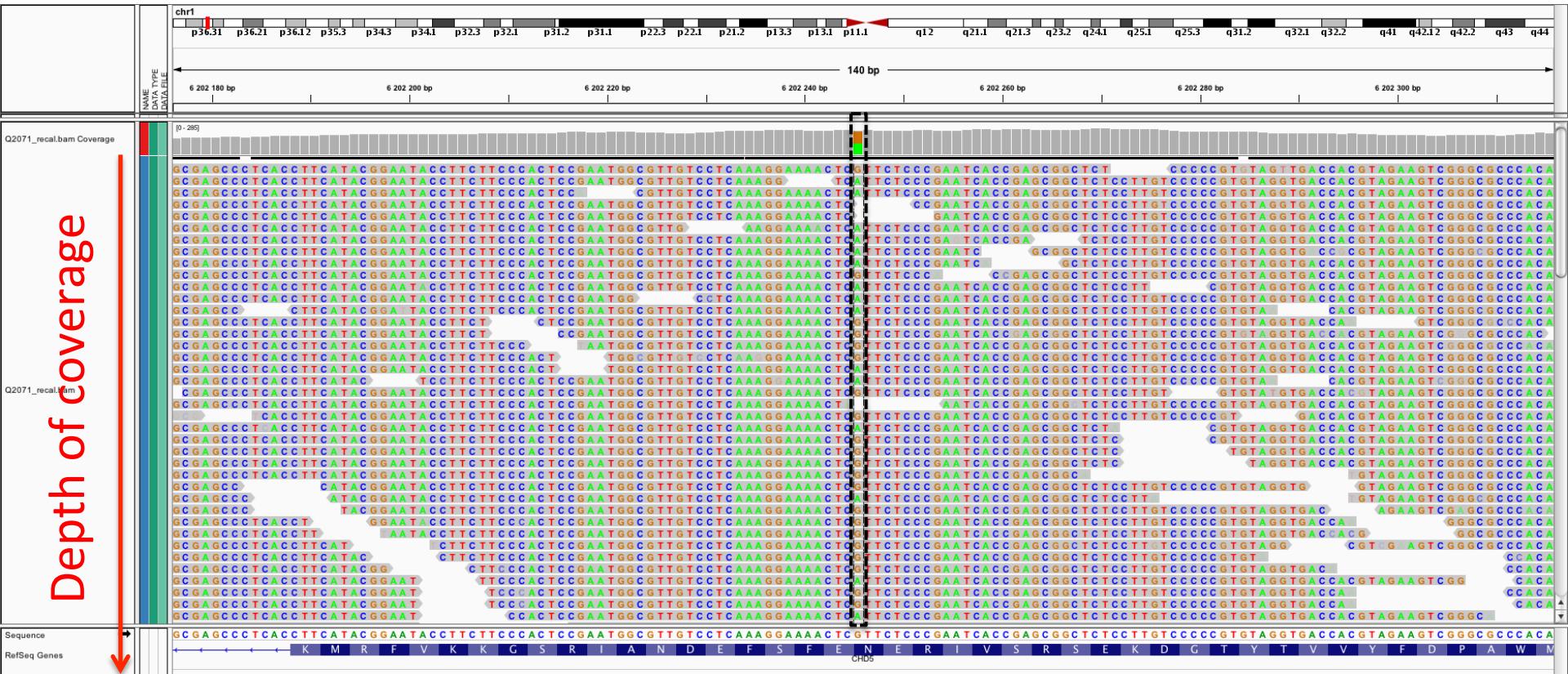


**Fichiers FASTQ  
(reads  
+qualités)**

**Fichiers BAM  
(reads alignés)**

**Fichiers BAM  
améliorés**

**Fichiers VCF  
(variants)**

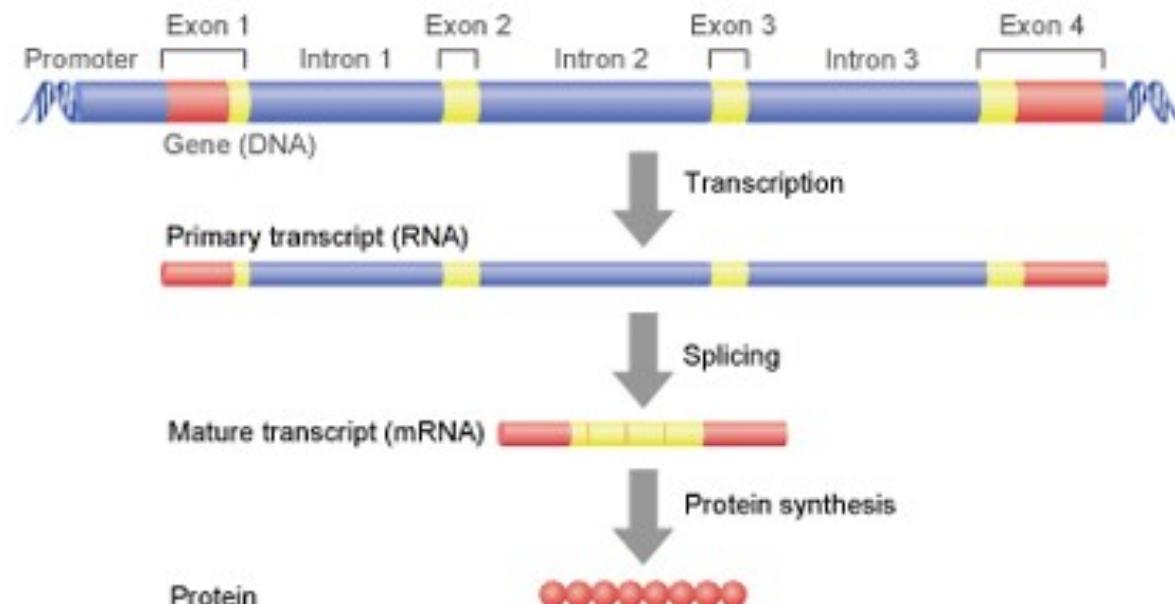


# Plan

- Avant le NGS (méthode Sanger)
- La révolution NGS
- Depuis le NGS
- Séquençage avec la plateforme Illumina
- Exemple

Identification de la mutation  
génétique responsable d'une  
maladie chez le patient X

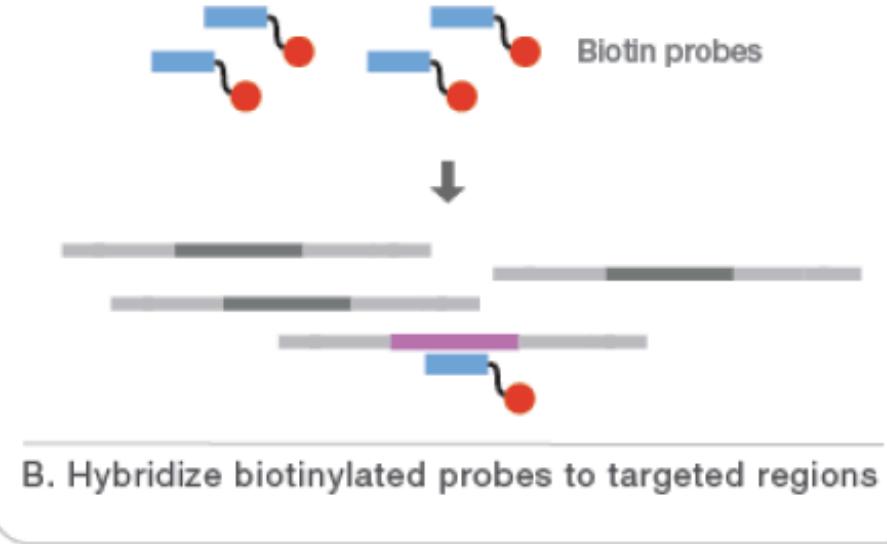
# Séquençage du génome entier ou de l'exome?



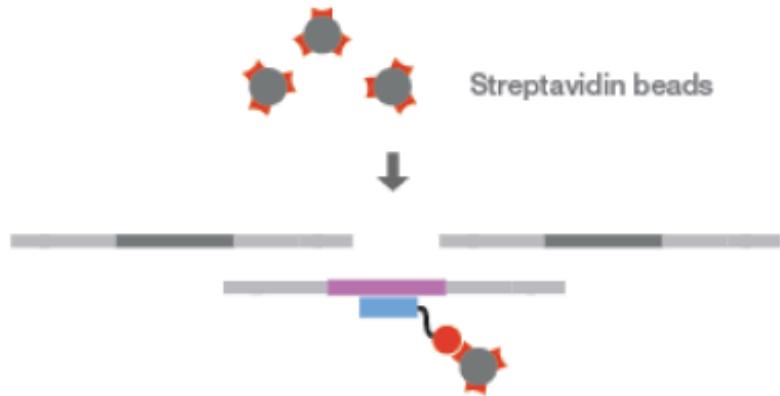
# Séquençage d'exome



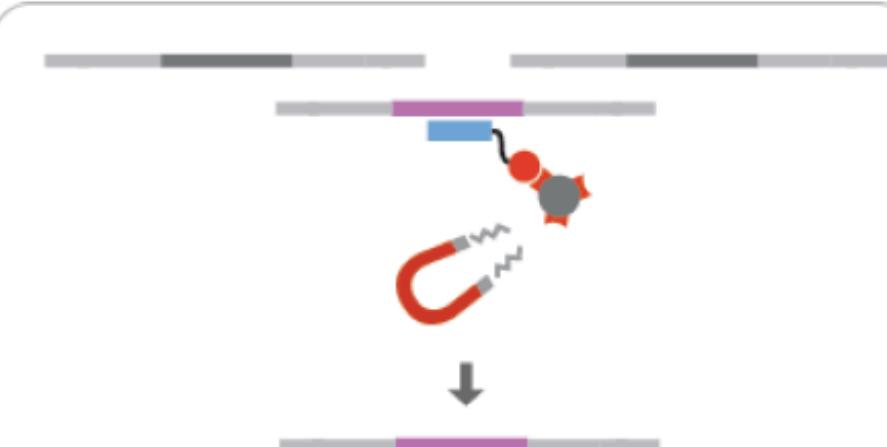
A. Denature double-stranded DNA library



B. Hybridize biotinylated probes to targeted regions

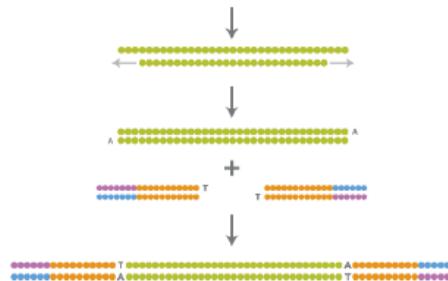


C. Enrichment using streptavidin beads



D. Elution from beads

Figure 3: Next-Generation Sequencing Simplified



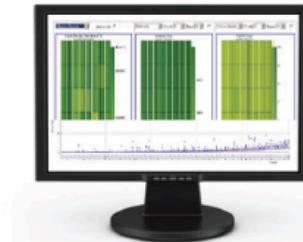
Library Preparation  
~2 h [15 min hands-on (Nextera)]  
< 6 h [< 3 h hands-on (TruSeq)]



Cluster Generation  
~5 h (<10 min hands-on)



Sequencing by Synthesis  
~1.5 to 11 days



CASAVA  
2 days (30 min hands-on)

From simplified sample preparation kits and automated cluster generation, to streamlined sequencing by synthesis and complete data analysis, Illumina HiSeq sequencing systems offer the industry's simplest next-generation sequencing workflow.

## Préparation de la librairie

## Génération de clusters

## Séquençage => reads

## Bioinformatique +++

Fichiers **FASTQ**  
(reads +qualités)

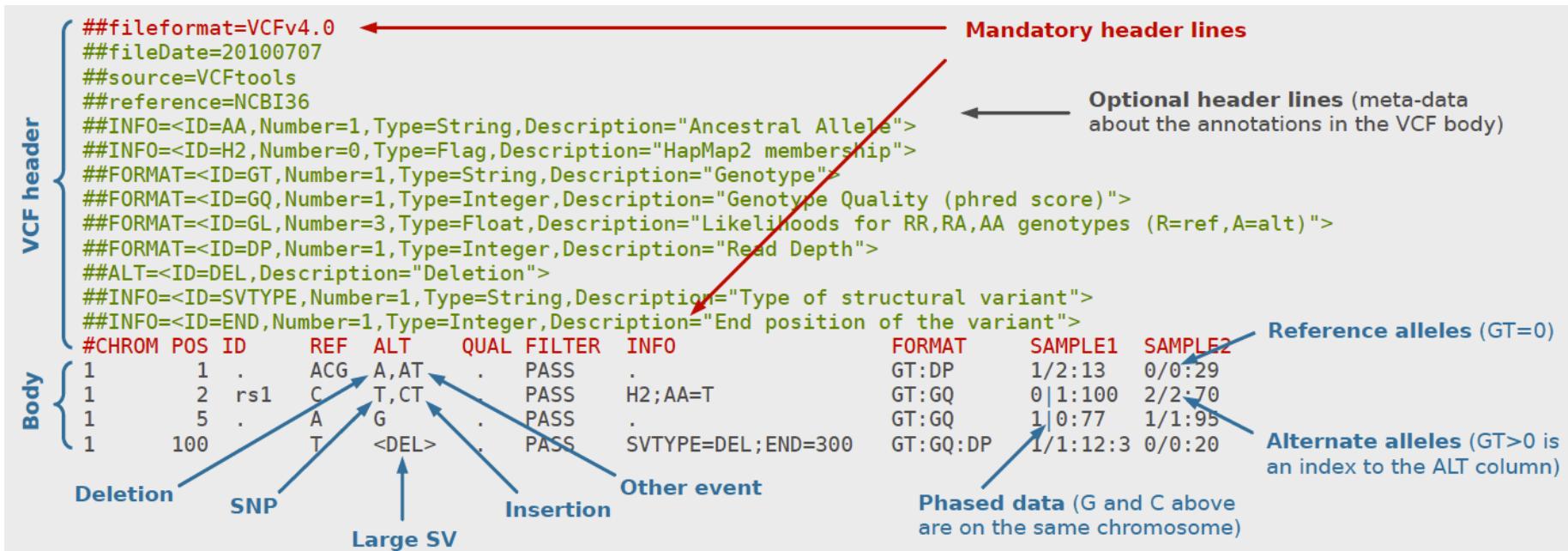
Fichiers **BAM**  
(reads alignés)

Fichiers **BAM**  
améliorés

Fichiers **VCF**  
(variants)

# Le format VCF

VCF = Variant Call Format



# Séquençage d'exome

Variant type	Mean number of variants ( $\pm$ sd) in African Americans	Mean number of variants ( $\pm$ sd) in European Americans
<b>Novel variants</b>		
Missense	303 ( $\pm$ 32)	192 ( $\pm$ 21)
Nonsense	5 ( $\pm$ 2)	5 ( $\pm$ 2)
Synonymous	209 ( $\pm$ 26)	109 ( $\pm$ 16)
Splice	2 ( $\pm$ 1)	2 ( $\pm$ 1)
Total	520 ( $\pm$ 53)	307 ( $\pm$ 33)
<b>Non-novel variants</b>		
Missense	10,828 ( $\pm$ 342)	9,319 ( $\pm$ 233)
Nonsense	98 ( $\pm$ 8)	89 ( $\pm$ 6)
Synonymous	12,567 ( $\pm$ 416)	10,536 ( $\pm$ 280)
Splice	36 ( $\pm$ 4)	32 ( $\pm$ 3)
Total	23,529 ( $\pm$ 751)	19,976 ( $\pm$ 505)
<b>Total variants</b>		
Missense	11,131 ( $\pm$ 364)	9,511 ( $\pm$ 244)
Nonsense	103 ( $\pm$ 8)	93 ( $\pm$ 6)
Synonymous	12,776 ( $\pm$ 434)	10,645 ( $\pm$ 286)
Splice	38 ( $\pm$ 5)	34 ( $\pm$ 4)
Total	24,049 ( $\pm$ 791)	20,283 ( $\pm$ 523)

*Sequence analysis*

Advance Access publication June 7, 2011

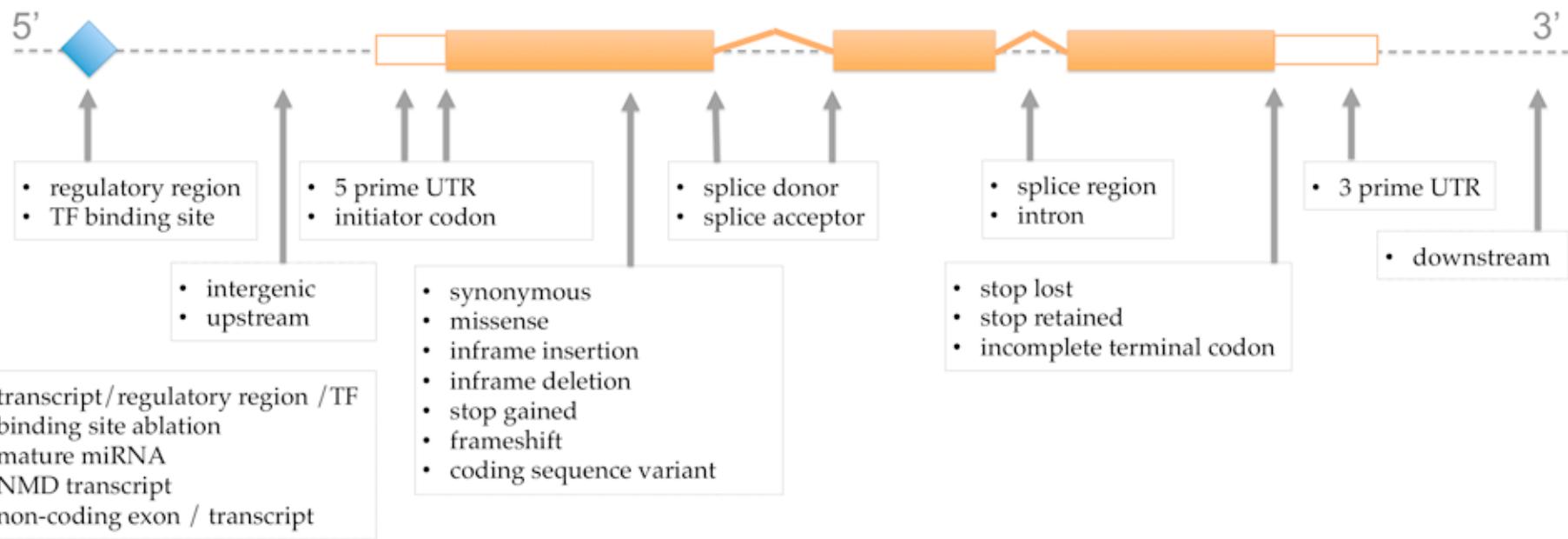
**The variant call format and VCFtools**

Petr Danecek<sup>1,†</sup>, Adam Auton<sup>2,†</sup>, Goncalo Abecasis<sup>3</sup>, Cornelis A. Albers<sup>1</sup>, Eric Banks<sup>4</sup>, Mark A. DePristo<sup>4</sup>, Robert E. Handsaker<sup>4</sup>, Gerton Lunter<sup>2</sup>, Gabor T. Marth<sup>5</sup>, Stephen T. Sherry<sup>6</sup>, Gilean McVean<sup>2,7</sup>, Richard Durbin<sup>1,\*</sup> and 1000 Genomes Project Analysis Group<sup>†</sup>

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge CB10 1SA, <sup>2</sup>Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK, <sup>3</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, <sup>4</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02141, <sup>5</sup>Department of Biology, Boston College, MA 02467, <sup>6</sup>National Institutes of Health National Center for Biotechnology Information, MD 20894, USA and <sup>7</sup>Department of Statistics, University of Oxford, Oxford OX1 3TG, UK

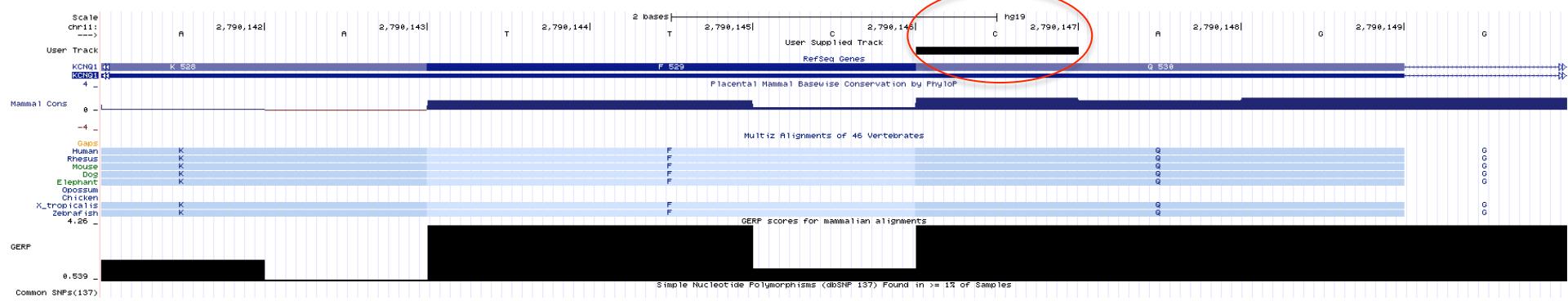
Associate Editor: John Quackenbush

# Annotation fonctionnelle (1)



# Annotation fonctionnelle (2)

CAG => AAG



#CHROM	POS	REF	ALT
Chr11	2790147	C	A

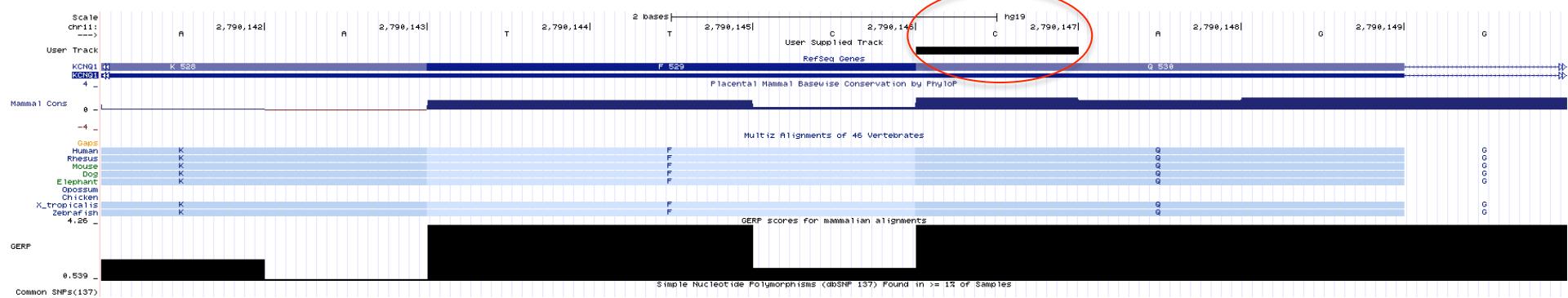
# le code génétique

ijk

	Deuxième lettre								
	U	C	A	G					
U	UUU	Phe	Ser	UAU	Tyr	UGU	Cys	U	
	UUC	Phe	Ser	UAC	Tyr	UGC	Cys	C	
	UUA	Leu	Ser	UAA	Stop	UGA	Stop	A	
	UUG	Leu	Ser	UAG	Stop	UGG	Trp	G	
C	CUU	Leu	Pro	CAU	His	CGU	Arg	U	Troisième lettre (côté 3')
	CUC	Leu	Pro	CAC	His	CGC	Arg	C	
	CUA	Leu	Pro	CAA	Gln	CGA	Arg	A	
	CUG	Leu	Pro	CAG	Gln	CGG	Arg	G	
A	AUU	Ile	Thr	AAU	Asn	AGU	Ser	U	
	AUC	Ile	Thr	AAC	Asn	AGC	Ser	C	
	AUA	Ile	Thr	AAA	Lys	AGA	Arg	A	
	AUG	Met	Thr	AAG	Lys	AGG	Arg	G	
G	GUU	val	Ala	GAU	Asp	GGU	Gly	U	
	GUC	val	Ala	GAC	Asp	GGC	Gly	C	
	GUA	val	Ala	GAA	Glu	GGA	Gly	A	
	GUG	val	Ala	GAG	Glu	GGG	Gly	G	
codon d'initiation					codon de terminaison				

# Annotation fonctionnelle (3)

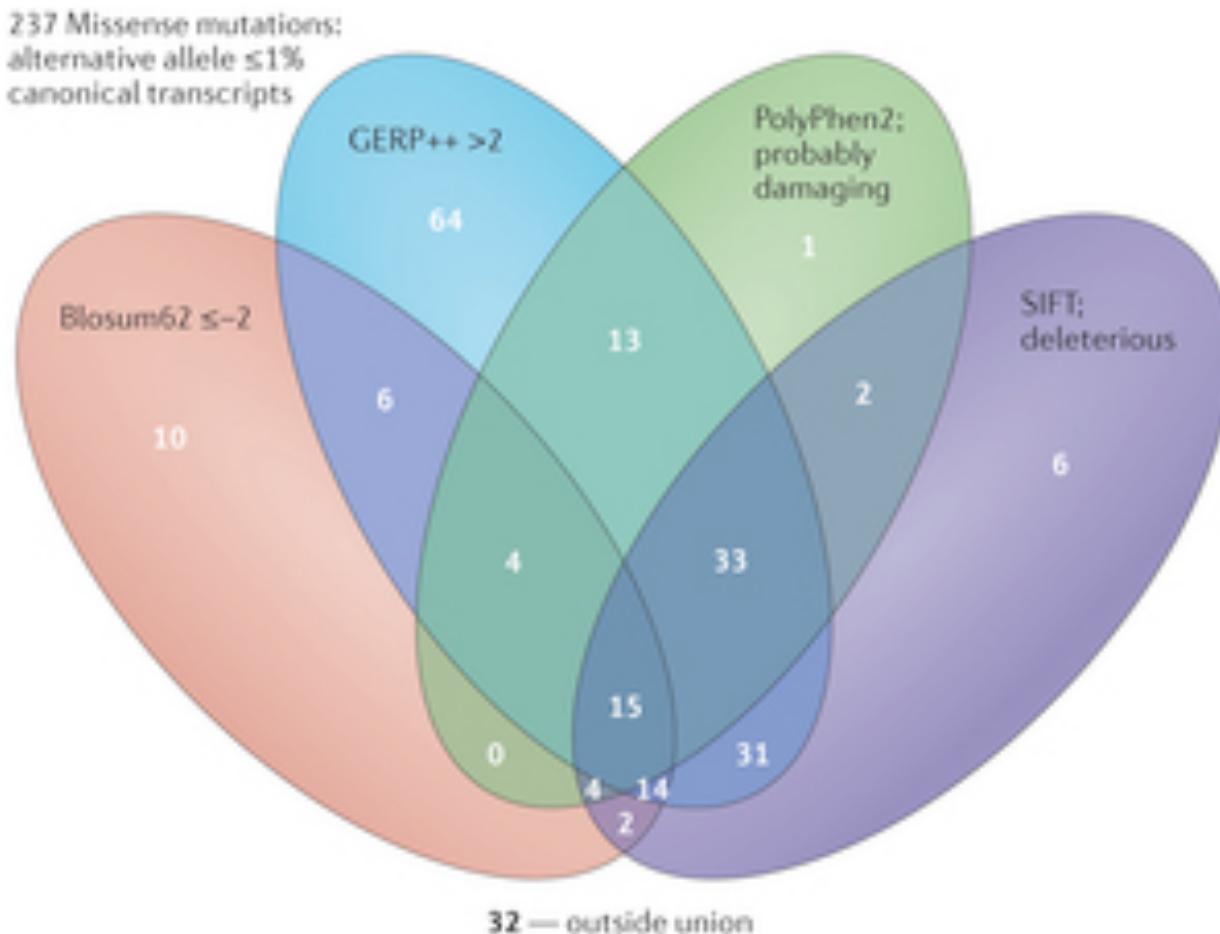
Q      K  
Gln      Lys  
CAG => AAG



#CHROM	POS	REF	ALT
Chr11	2790147	C	A

Mutation faux-sens Q530K

# Outils de prédiction



Nature Reviews | Genetics

# Prioritisation des variants

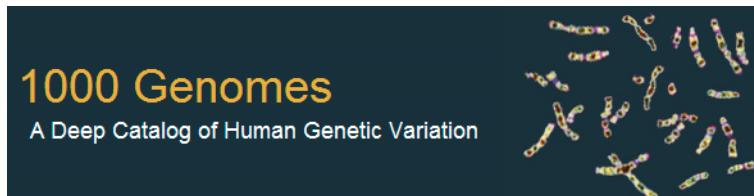


**INGENUITY<sup>®</sup>**  
S Y S T E M S



Johns  
Hopkins  
University

# Focus sur les variants rares?



Variant Pos	rs ID	Alleles	EA Allele #	AA Allele #	All Allele #	EA Genotype #	AA Genotype #	All Genotype #	MAF (%) (EA/AA/All)
3:38595746	<a href="#">rs45440596</a>	C>T	T=1/C=8327	T=1/C=3991	T=2/C=12318	TT=0/TC=1 /CC=4163	TT=0/TC=1 /CC=1995	TT=0/TC=2 /CC=6158	0.012/0.0251 /0.0162
3:38595747	<a href="#">rs369559182</a>	G>A	A=8/G=8318	A=0/G=4006	A=8/G=12324	AA=0/AG=8 /GG=4155	AA=0/AG=0 /GG=2003	AA=0/AG=8 /GG=6158	0.0961/0.0 /0.0649
3:38595795	<a href="#">rs141789366</a>	G>A	A=0/G=8444	A=5/G=4161	A=5/G=12605	AA=0/AG=0 /GG=4222	AA=0/AG=5 /GG=2078	AA=0/AG=5 /GG=6300	0.0/0.12/0.0397
3:38595795	<a href="#">rs141789366</a>	G>A	A=0/G=8444	A=5/G=4161	A=5/G=12605	AA=0/AG=0 /GG=4222	AA=0/AG=5 /GG=2078	AA=0/AG=5 /GG=6300	0.0/0.12/0.0397
3:38595801	<a href="#">rs376720757</a>	G>A	A=1/G=8453	A=1/G=4177	A=2/G=12630	AA=0/AG=1 /GG=4226	AA=0/AG=1 /GG=2088	AA=0/AG=2 /GG=6314	0.0118/0.0239 /0.0158
3:38595801	<a href="#">rs376720757</a>	G>A	A=1/G=8453	A=1/G=4177	A=2/G=12630	AA=0/AG=1 /GG=4226	AA=0/AG=1 /GG=2088	AA=0/AG=2 /GG=6314	0.0118/0.0239 /0.0158
3:38595989	<a href="#">rs199473618</a>	C>T	T=0/C=8396	T=7/C=4079	T=7/C=12475	TT=0/TC=0 /CC=4198	TT=0/TC=7 /CC=2036	TT=0/TC=7 /CC=6234	0.0/0.1713 /0.0561
3:38595990	<a href="#">rs372507927</a>	G>A	A=0/G=8398	A=1/G=4073	A=1/G=12471	AA=0/AG=0 /GG=4199	AA=0/AG=1 /GG=2036	AA=0/AG=1 /GG=6235	0.0/0.0245 /0.008
3:38596010	<a href="#">rs199473269</a>	C>T	T=1/C=8369	T=0/C=4026	T=1/C=12395	TT=0/TC=1 /CC=4184	TT=0/TC=0 /CC=2013	TT=0/TC=1 /CC=6197	0.0119/0.0 /0.0081

<http://www.1000genomes.org/>

<http://evs.gs.washington.edu/EVS/>



# Populations séquencées



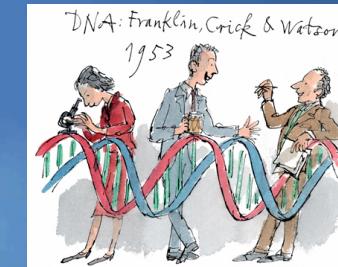
Tiny Faroe Islands to Begin Sequencing Genomes of All 50,000 Residents in Ambitious Effort to Advance Personalized Medicine

The SardiNIA Project

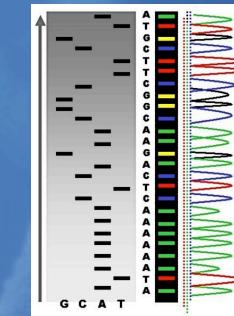


VACARME

# 1953: Structure en double hélice (Watson & Crick)



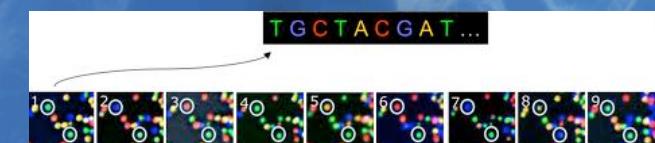
# 1977: Séquençage Sanger (``première génération``)



# 2001-2003: HGP



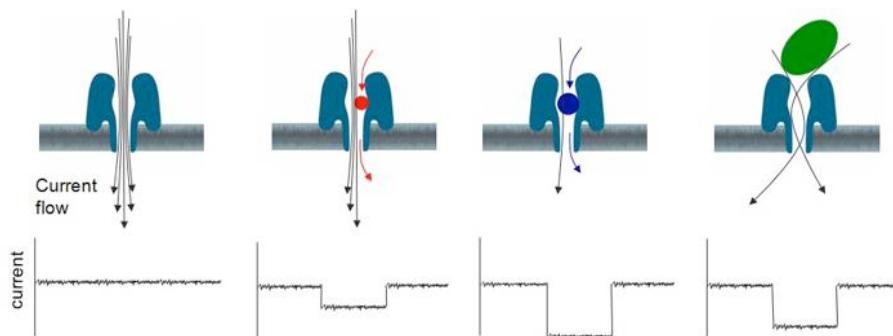
# 2005: Séquençage nouvelle génération (``seconde génération``)



# Next-next generation sequencing

## MinION: A complete DNA sequencer on a USB stick

By John Hewitt on March 29, 2013 at 10:07 am | 4 Comments



Helicos Bioxciences  
Pacific Biosciences  
Oxford Nanopore

...



« An exciting time for genomics »

