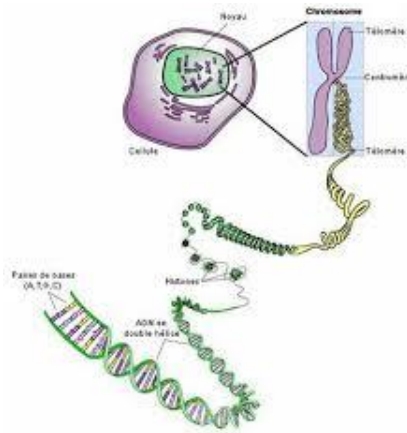


# Les marqueurs moléculaires

**Philippe Barre**

INRA Centre Poitou-Charentes

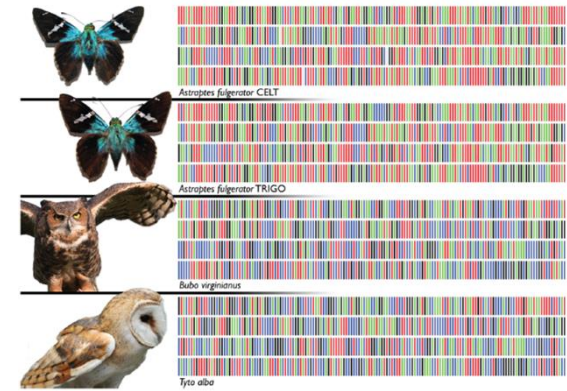
UR4 : Unité de Recherche Pluridisciplinaire, Prairies et Plantes Fourragères  
(URP<sup>3</sup>F)



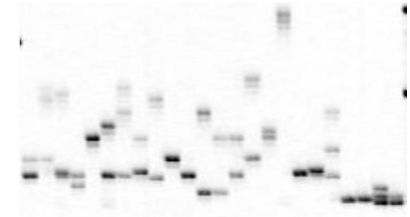
Poitiers M1 Biologie Ecologie 10 Février 2016

# Plan

É Introduction

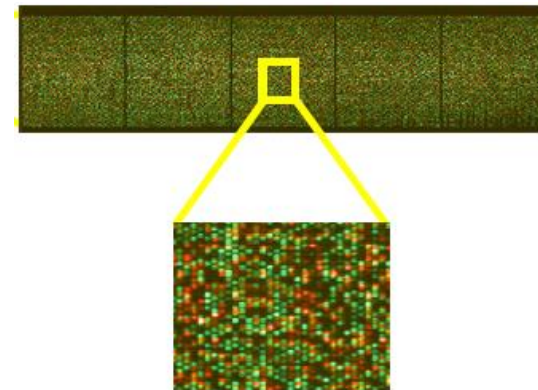


É Les marqueurs « bas- moyen débit »



É Les marqueurs « haut-très haut débit »

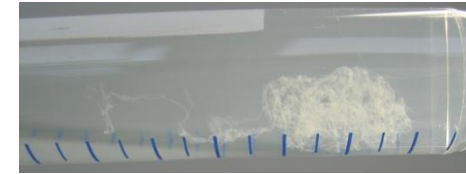
É Conclusion





# INTRODUCTION

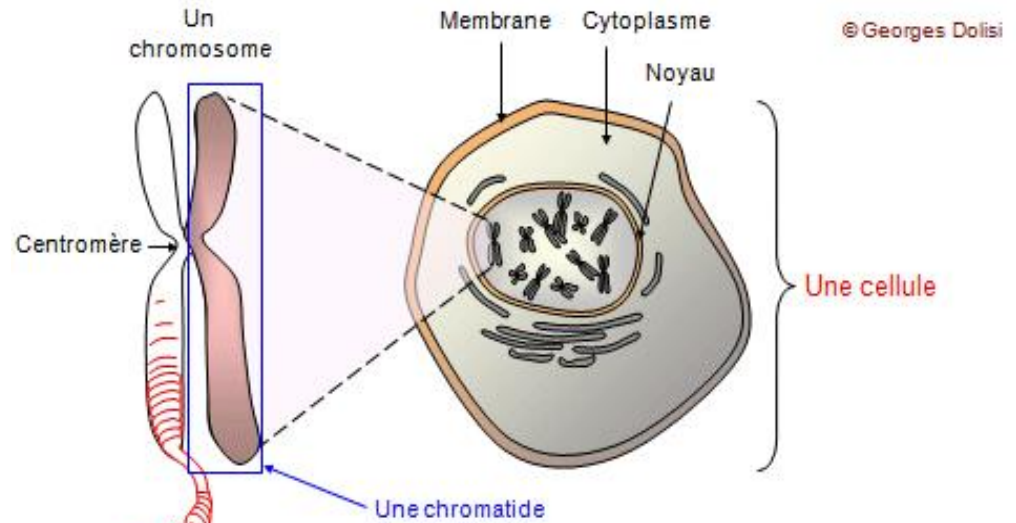
# L'ADN ?



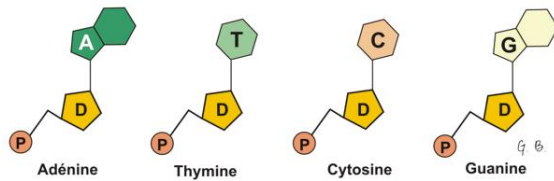
## Acide DésoxyriboNucléique

C'est le support de l'information  
génétique

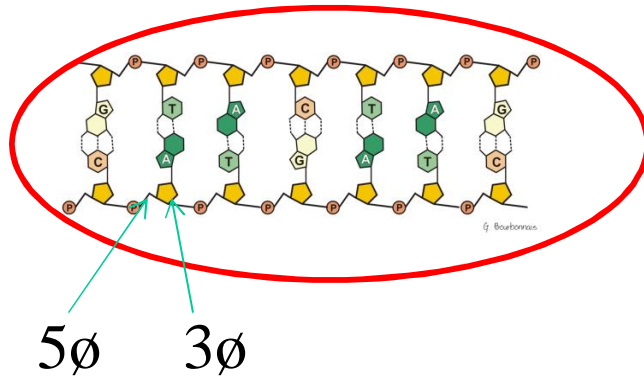
**Chromosome**  
(métaphase) :  $\rightarrow$  Forme condensée de  
l'ADN (= forme visible)



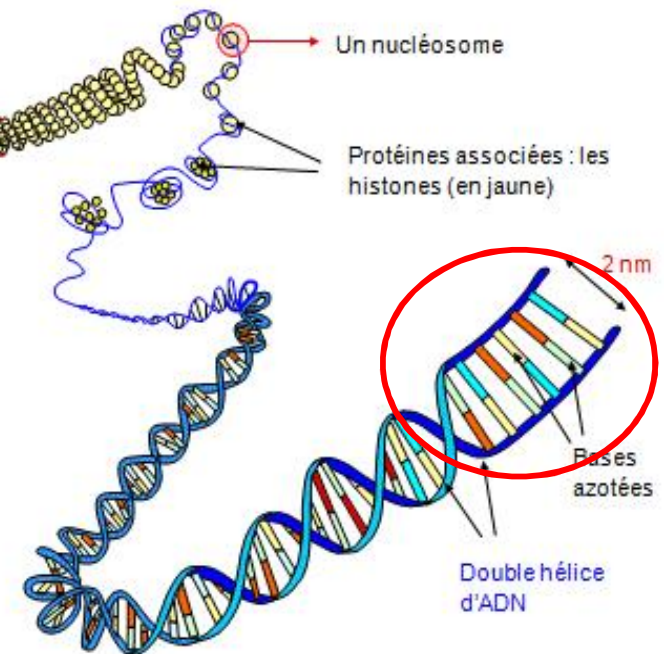
Éléments constitutifs : **4 nucléotides**



Organisation en double brin :

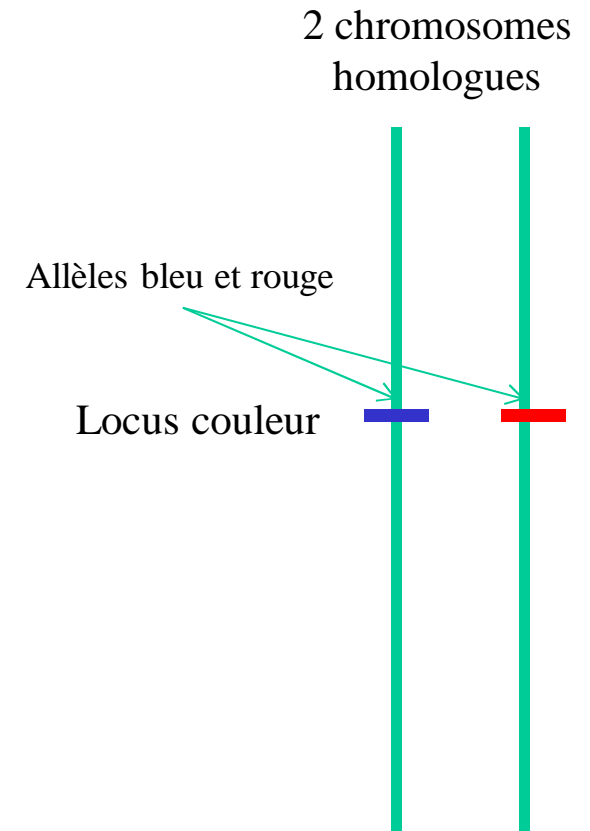


- Cytosine
- Guanine
- Adénine
- Thymine



# Définition des marqueurs moléculaires

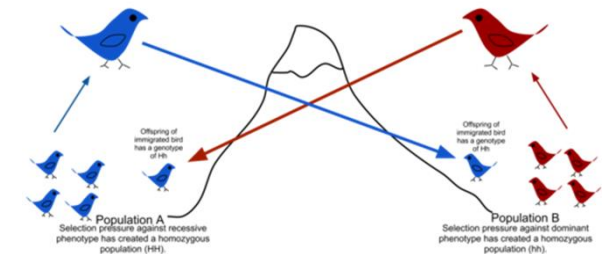
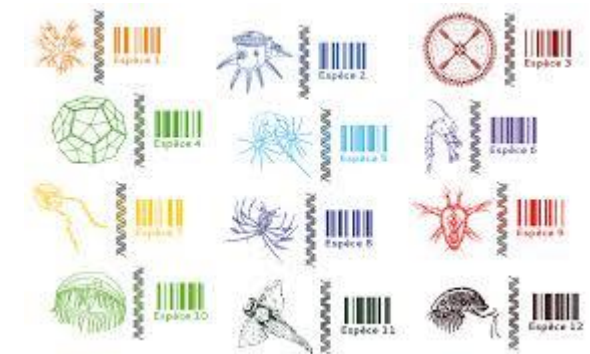
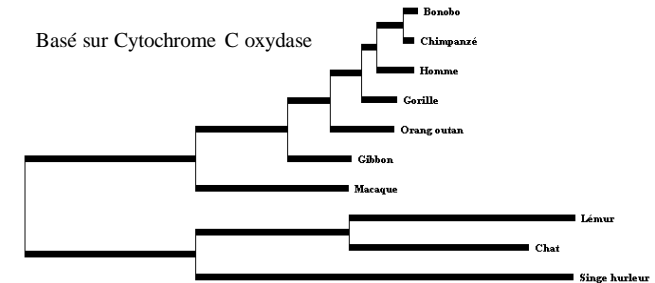
- “ Un endroit sur le génome : **locus**
- “ Présentant du polymorphisme entre et au sein des individus : **allèles**
- “ Les différents allèles sont identifiables en totalité (**co-dominant**) ou en partie (**dominant**) pour chaque individu



# A quoi servent les marqueurs ?

“ Etude de la diversité entre individus de la même espèce ou de espèces différentes

- Evolution des espèces
- Identification des espèces (barcoding)
- Identification des individus (criminalité)
- Structuration de la diversité au sein des espèces (flux de gènes, sélection, dérive, mutation)



# A quoi servent les marqueurs ?

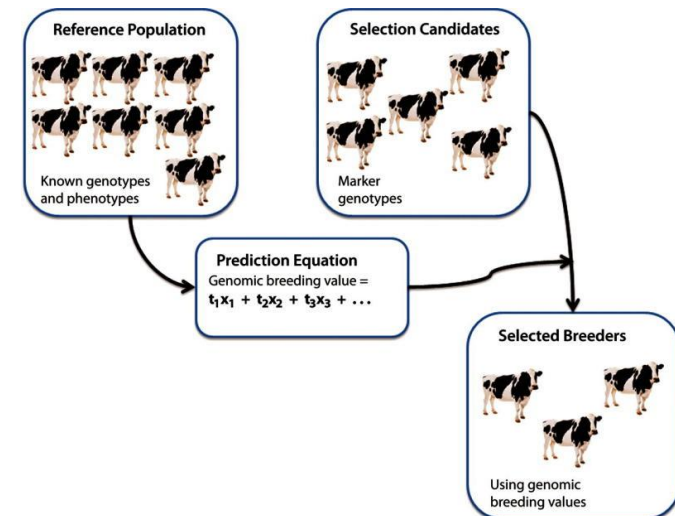
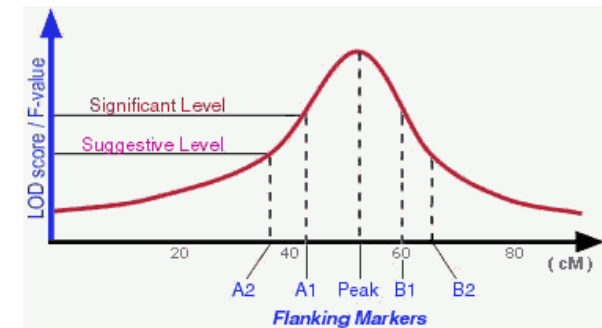
” Relation entre polymorphisme moléculaire et polymorphisme phénotypique

- Identification de régions du génome ou de gènes impliqués dans la variation d'un caractère quantitatif (QTL)

” Construction d'un individu sur la base des allèles aux QTL

- Prédiction de la valeur génétique d'un individu à l'aide de son génome

” Sélection génomique



# Une évolution fulgurante!!!

É Découverte de la double hélice  
de l'ADN 1953 (James Watson  
and Francis Cricks)

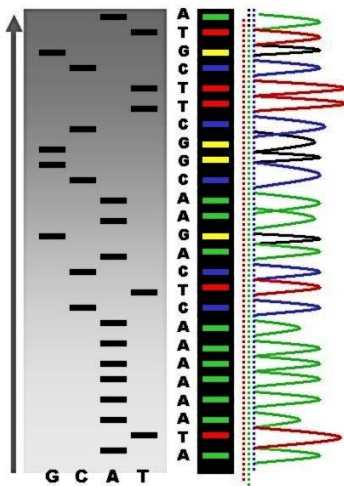
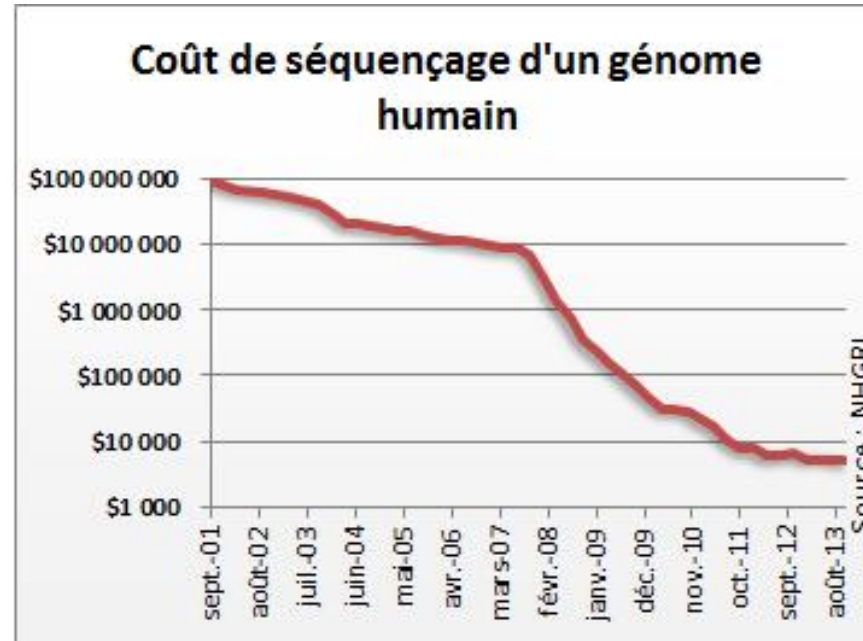
É Techniques de hybridation  
(années 1970)

É Séquençage de l'ADN 1955  
(Fred Sanger)

É Invention de la PCR 1993 (Kary  
Bank Mullis)

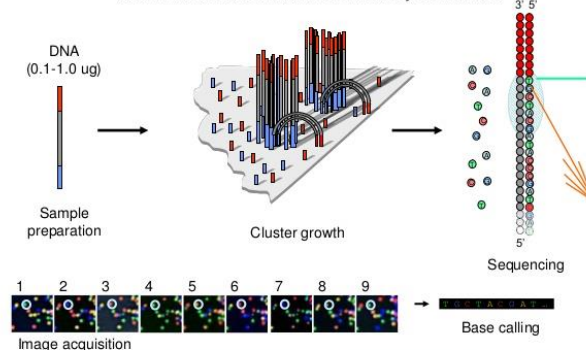


# Une évolution fulgurante!!!



## ILLUMINA Sequencing Technology

Robust Reversible Terminator Chemistry Foundation



[http://www.ncbi.nlm.nih.gov/genome/pr  
obe/doc/Technologies.shtml](http://www.ncbi.nlm.nih.gov/genome/pr<br/>obe/doc/Technologies.shtml)

<http://www2.warwick.ac.uk/fac/sci/lifesci/research/vegin/geneticimprovement/geneticmarker/>

<http://get.genotoul.fr/>

<http://gentyane.clermont.inra.fr/>

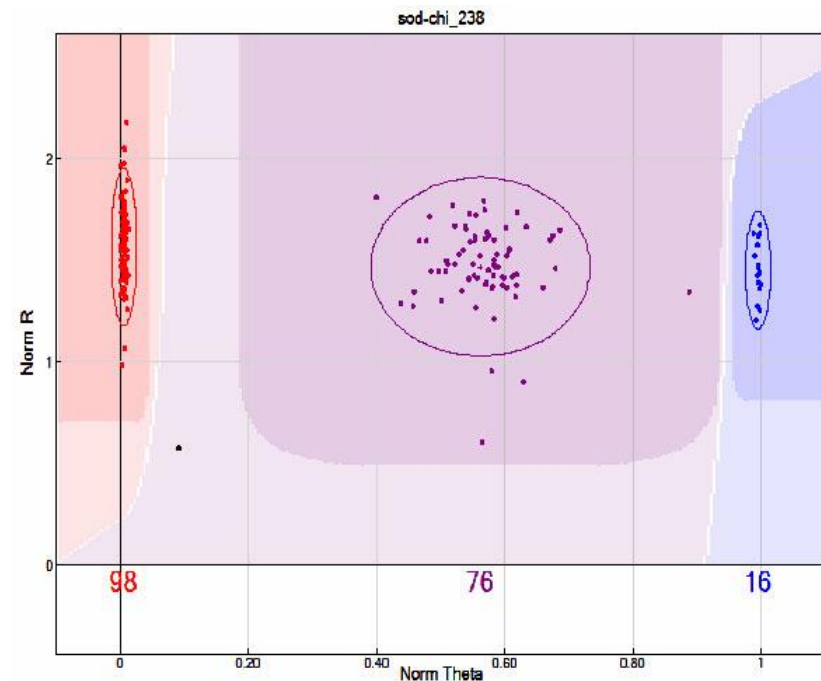
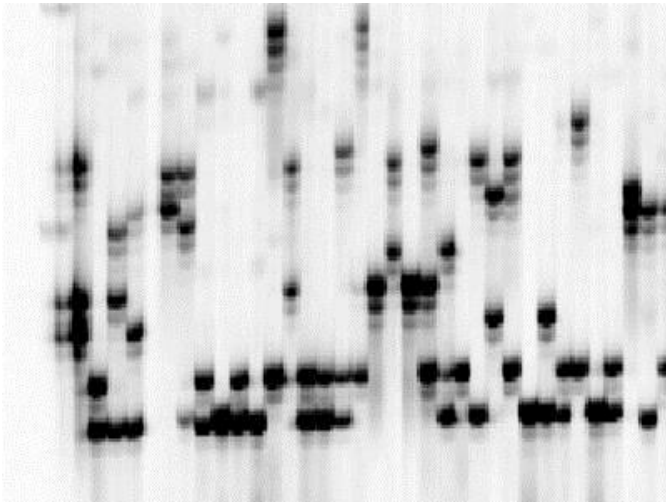
# Une évolution fulgurante!!!

Coût du génotypage dépend beaucoup du nombre de marqueurs et du nombre d'individus

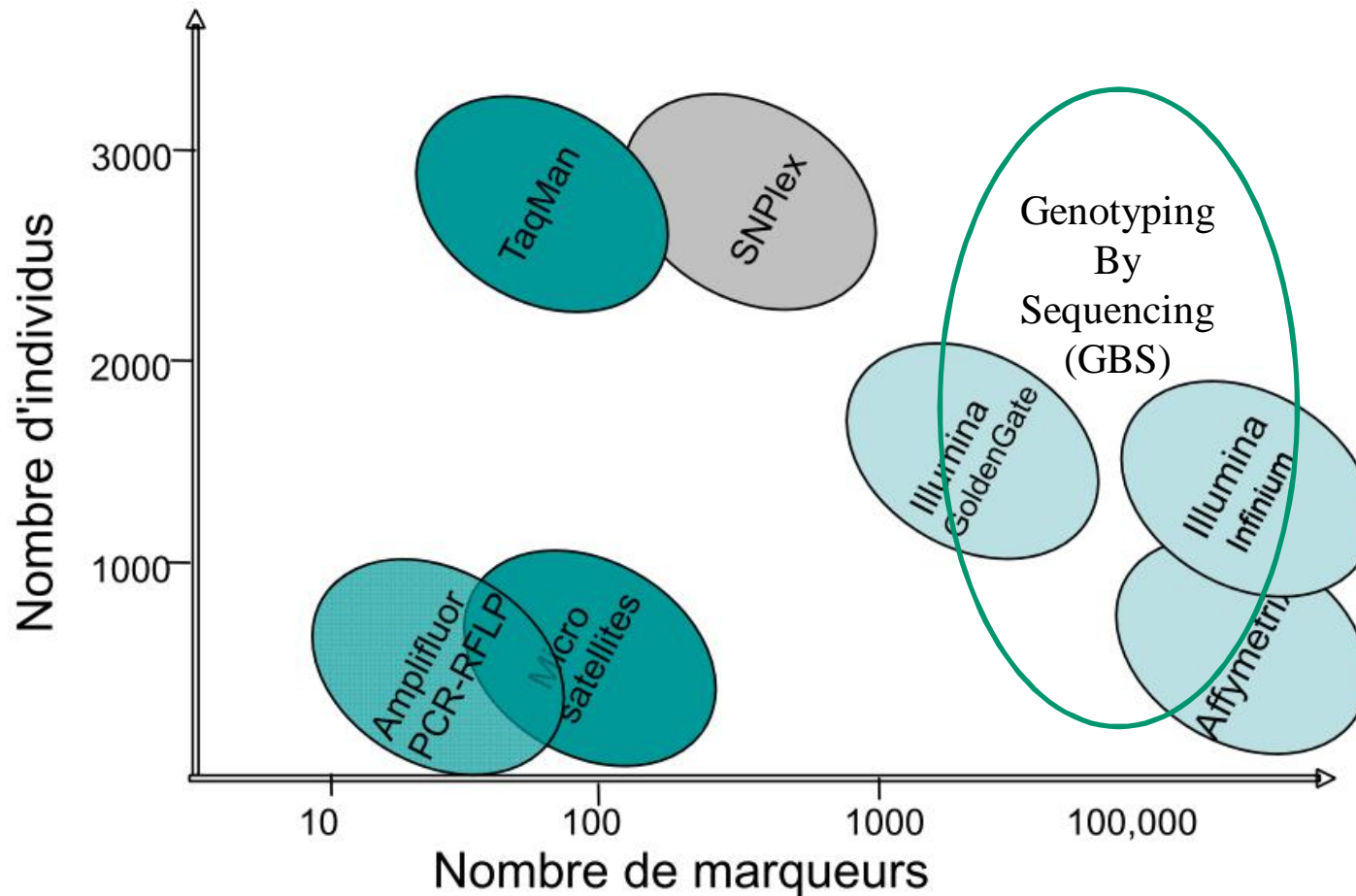
Exemples:

- 0,3 " /point pour un marqueur sur un génotype
- 0,05" /point pour 384 marqueurs sur 1536 individus
- 0,003" /point pour 20000 marqueurs sur 500 individus

Moins le coût du point est cher plus il faut génotyper de marqueurs et d'individus



# De multiples méthodes de génotypage



A définir selon les besoins !!!

# Quelles régions du génome?

É Répartis sur l'ensemble du génome

É Ciblés sur des régions particulières:

- ó Gènes candidats

- ó Ensemble des gènes

É Connaissance préalable du génome

- ó Séquençage complet / partiel du génome

- ó Séquençage du transcriptome



# **LES MARQUEURS « BAS- MOYEN DÉBIT »**

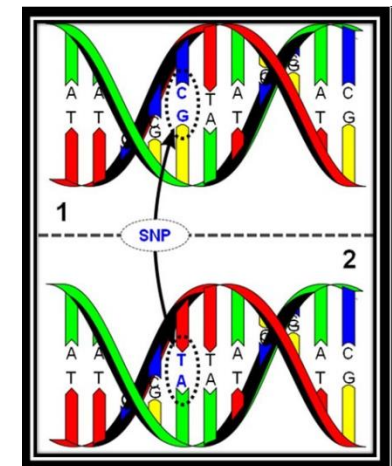
# De multiples méthodes

## “ Polymorphisme de longueur

- . par hybridation: Restriction Fragments Length Polymorphism (pour mémoire)
- . par PCR: Sequence Tagged Site (STS) including Simple Sequence Repeat (SSR), Cleaved Amplified Polymorphic Sequences (CAP), Amplified Fragment Length Polymorphism (AFLP).õ

## “ Polymorphisme de séquence : Single Nucleotide Polymorphism (SNP)

- . Taqman
- . Kompetitive Allelic Specific PCR (KASP)õ



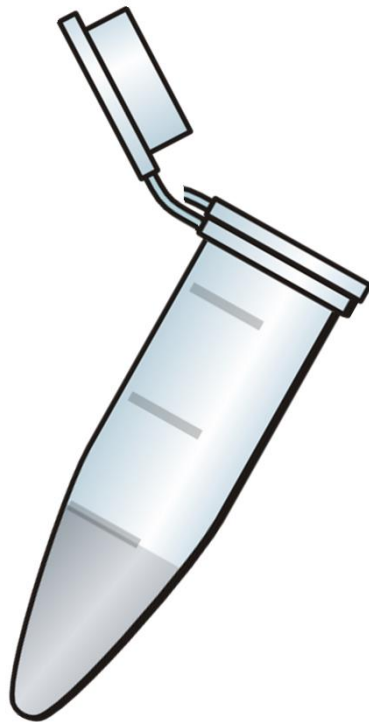
# La PCR

Multiplier un fragment d'ADN un grand nombre de fois afin de le repérer.

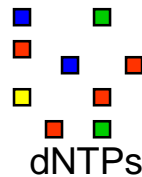


Technique de la PCR : **P**olymerase **C**hain **R**eaction

## Réactifs nécessaires à la PCR



tampon PCR 10X



MgCl<sub>2</sub>



amorces



Taq polymérase



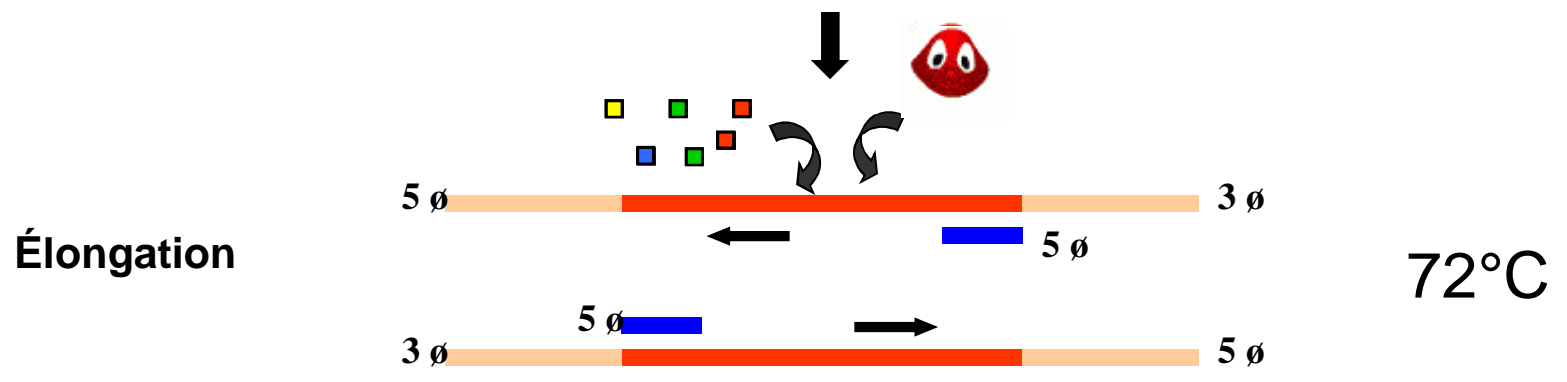
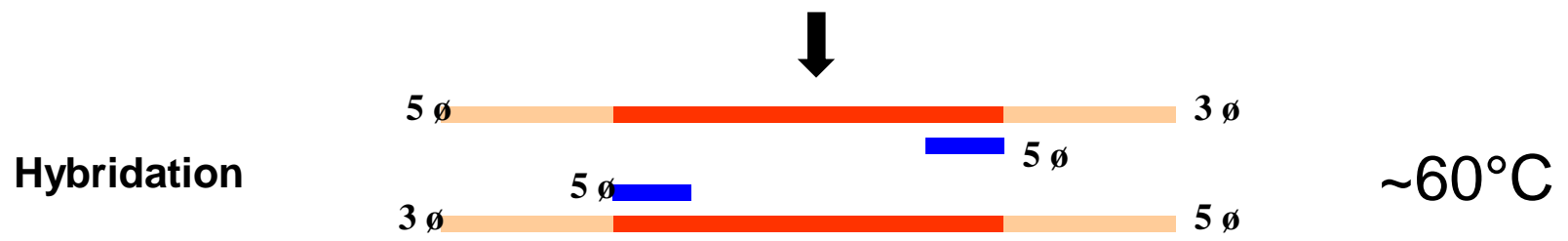
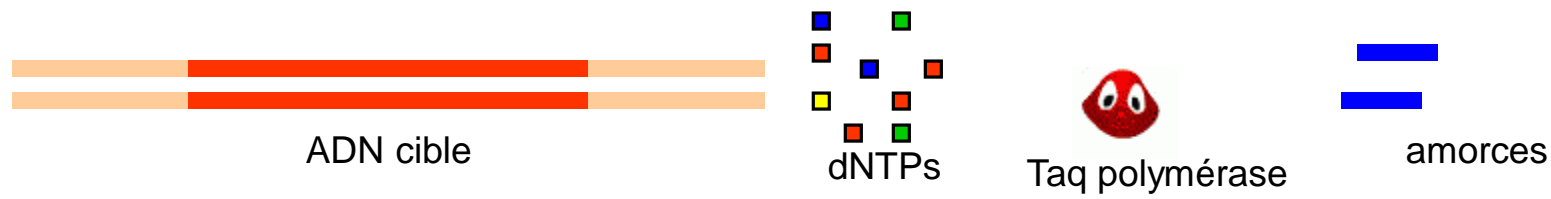
ADN



eau

## Thermocycleur





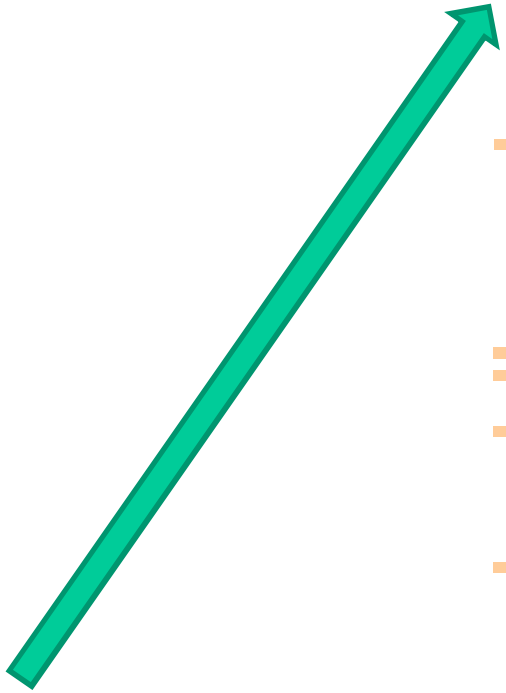
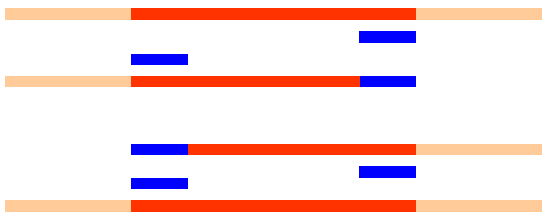




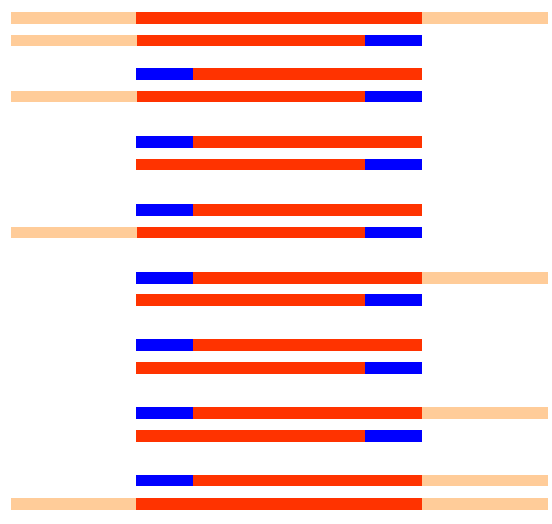
1 copie



2 copie



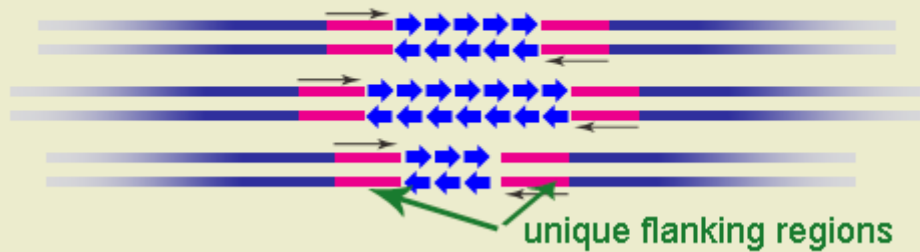
4 copie



8 copie

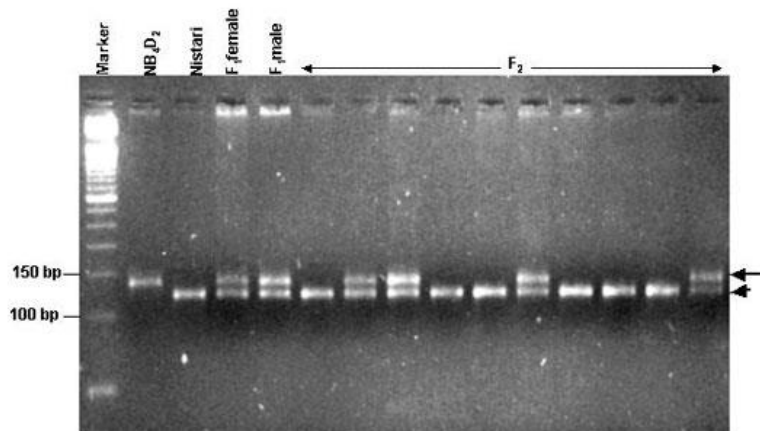
# Exemple 1: SSR, microsatellite

The number of SSRs is highly variable among individuals



Nécessité de connaître la séquence pour définir des amorces :

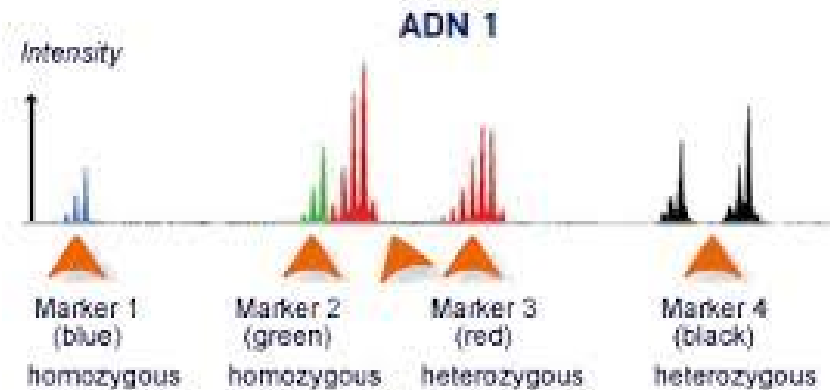
- Banques enrichies en SSR
- Résidus du séquençage de transcriptomes



Inheritance and segregation of SSR loci in the silkworm. SSR profiles obtained using Bmsat070 on two parental strains, Nistari and NB<sub>1</sub>D<sub>2</sub> and their F<sub>1</sub> and F<sub>2</sub> offspring (10Nos). The segregation of NB<sub>1</sub>D<sub>2</sub> (Female) and Nistari (Male) specific polymorphic loci is shown by arrow and arrow heads respectively.

Electrophorèse sur gel d'agarose

SSR (Single Sequence Repeat)



Electrophorèse sur séquenceur capillaire

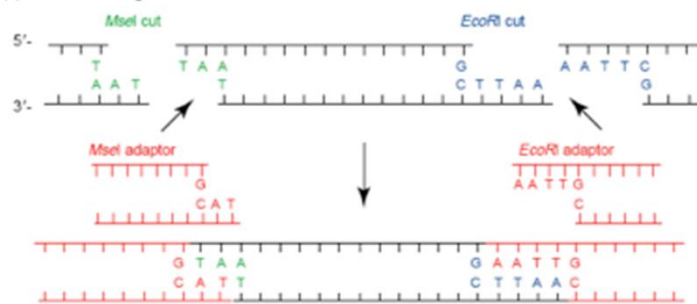
# Exemple 2: Amplified Fragment Length Polymorphism (AFLP)

## AFLP® assays

(a) AFLP template preparation  
Whole genomic DNA

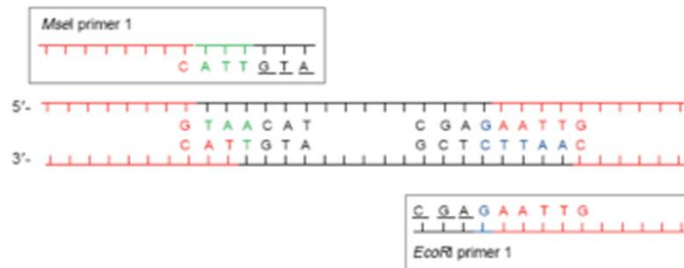


(b) Restriction and ligation

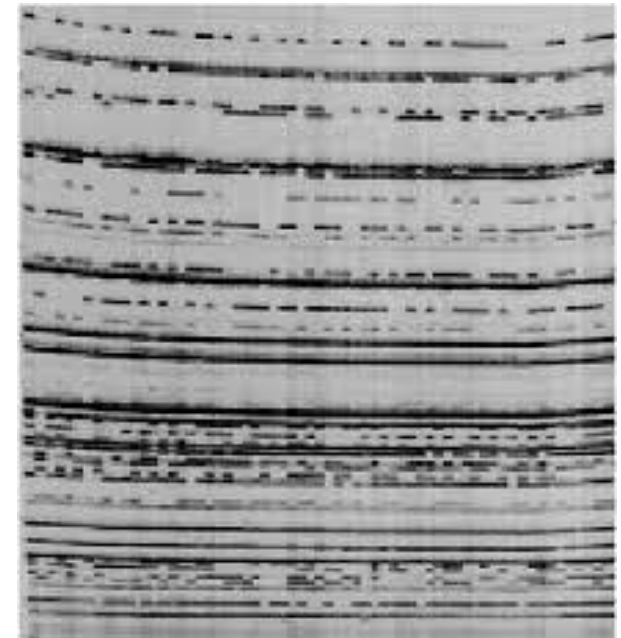


Pre-amplify with one selective nucleotide to give 4 sub-sets of "Preamp" reactions

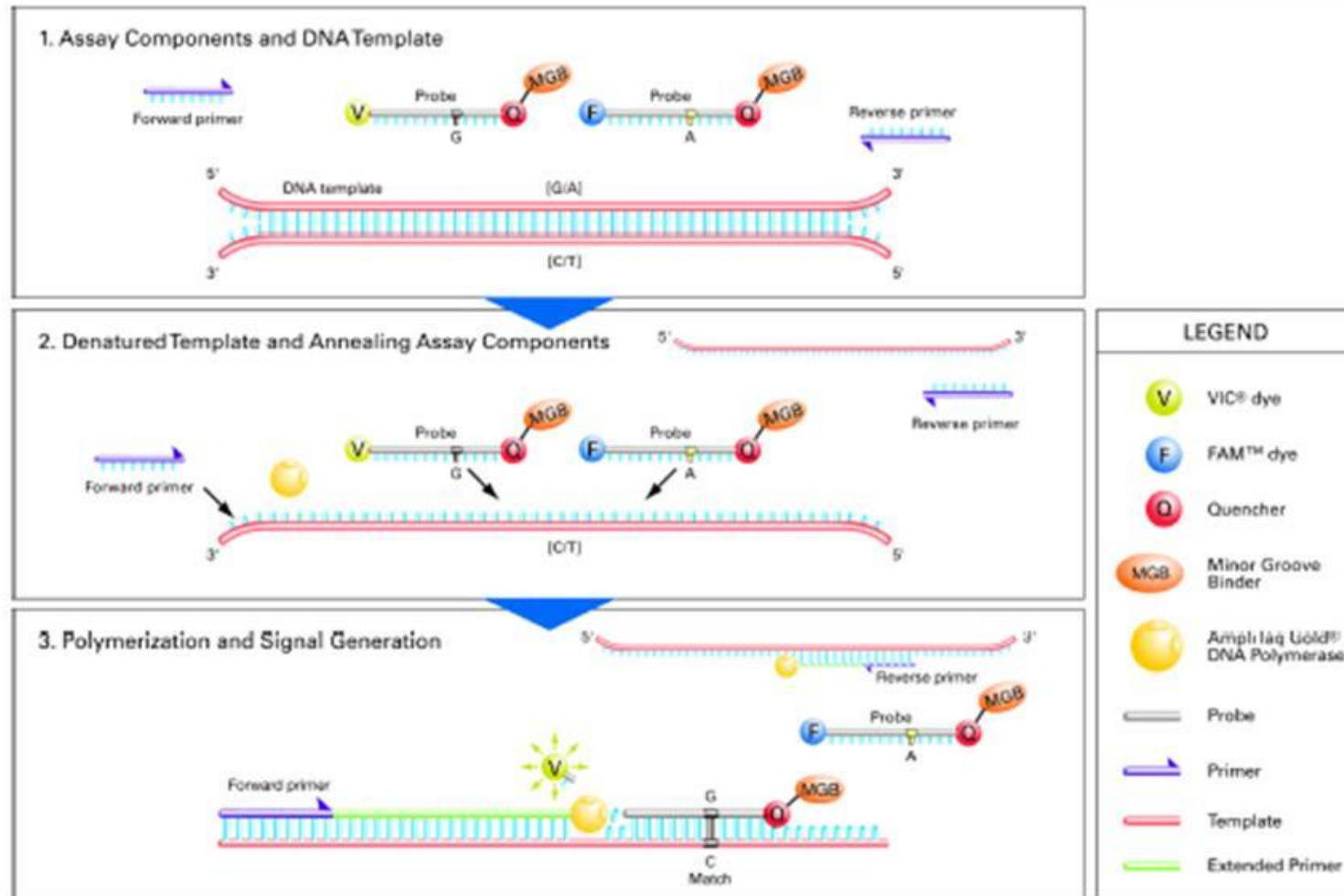
(c) Selective amplification (one of many primer combinations shown)



AFLP technology developed by KeyGene in early 1990's  
See Vos et al. (1995).



# Exemple 3: Système Taqman

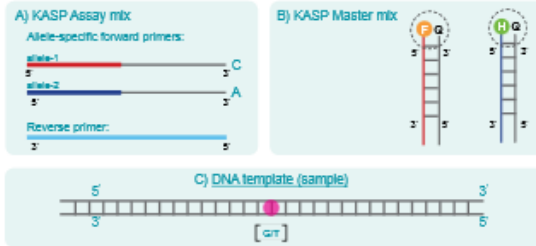


Cher à développer mais une fois développer très peu cher mais pas « multiplexable » !

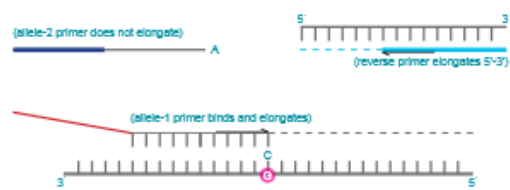
# Exemple 4: Système KASP KBioscience

## 1) Assay components:

KASP uses three components: test DNA with the SNP of interest; KASP Assay Mtx containing two different, allele-specific, competing forward primers with unique tail sequences and one reverse primer; the KASP Master mix containing FRET cassette plus Taq polymerase in an optimised buffer solution.

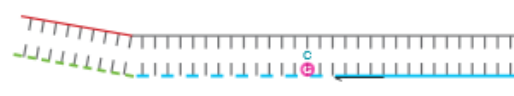


## 2) Denatured template and annealing components – PCR round 1:



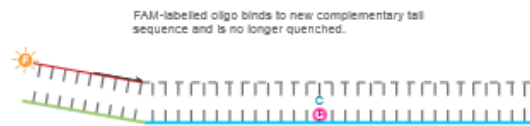
In the first round of PCR, one of the allele-specific primers matches the target SNP and, with the common reverse primer, amplifies the target region.

## 3) Complement of allele-specific tail sequence generated – PCR round 2:

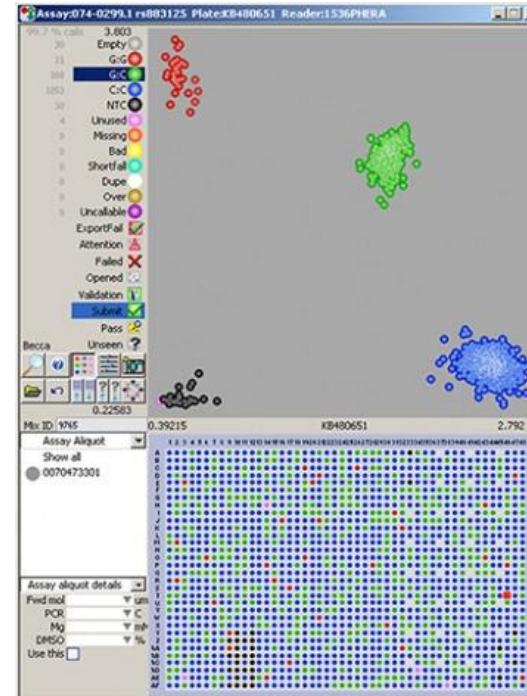
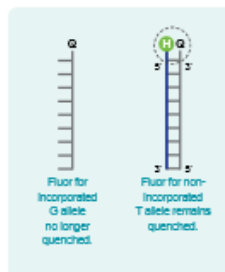


(Reverse primer binds, elongates and makes a complementary copy of the allele-1 tail.)

## 4) Signal generation – PCR round 3:



In further rounds of PCR, levels of allele-specific tail increase. The fluor labelled part of the FRET cassette is complementary to new tail sequences and binds, releasing the fluor from the quencher to generate a fluorescent signal.

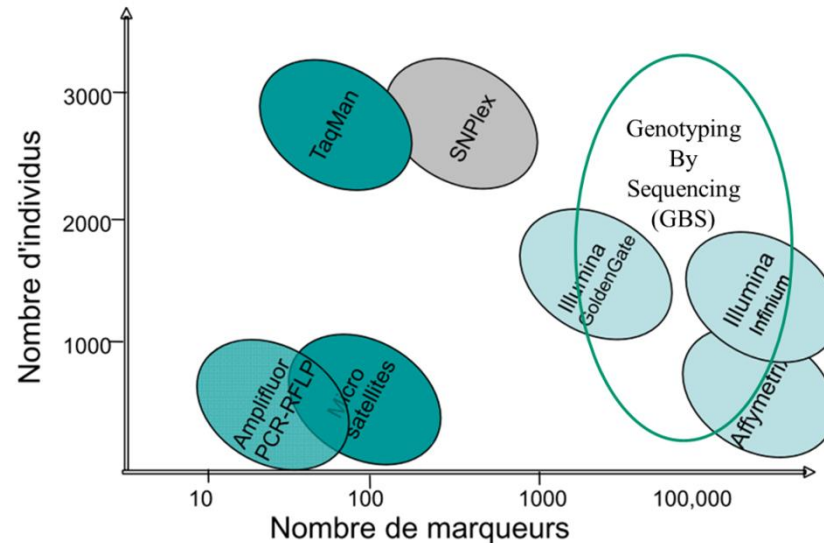


- Utilisation de la PCR + lecteur fluorescence
- Kit pour 96, 384, 1536 individus



**LES MARQUEURS « HAUT-  
TRES HAUT DÉBIT »**

# Haut óTrès haut débit ?



<http://crgs.genopole-toulouse.prd.fr>



É A la fois beaucoup de marqueurs (>1000) et beaucoup d'individus (>1000)

É Nécessité de robot de pipetage

É Nécessité de mémoire informatique et puissance de calcul

# Deux grands types

É Génotypage par codage

ó Illumina golden gate

ó Illumina Infinium

ó Affymetrixí

É Génotypage par séquençage : GBS

ó Différentes régions du génome  
séquencées

ó De multiples techniques de  
séquençage

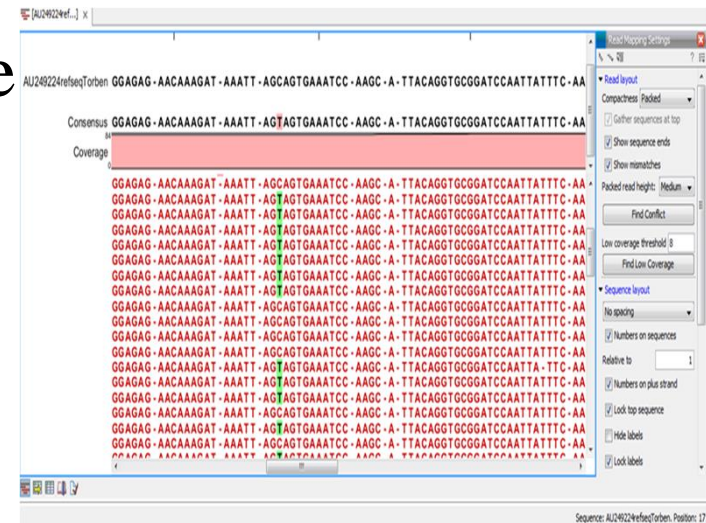
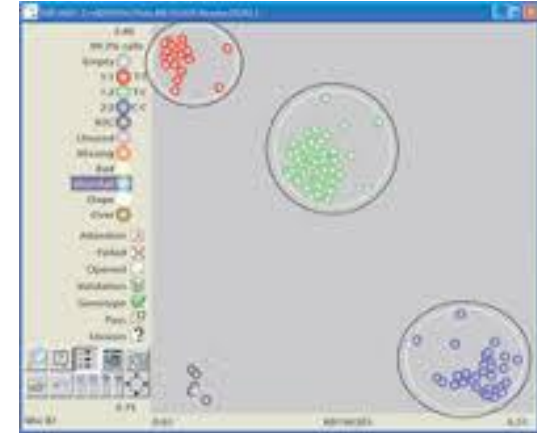
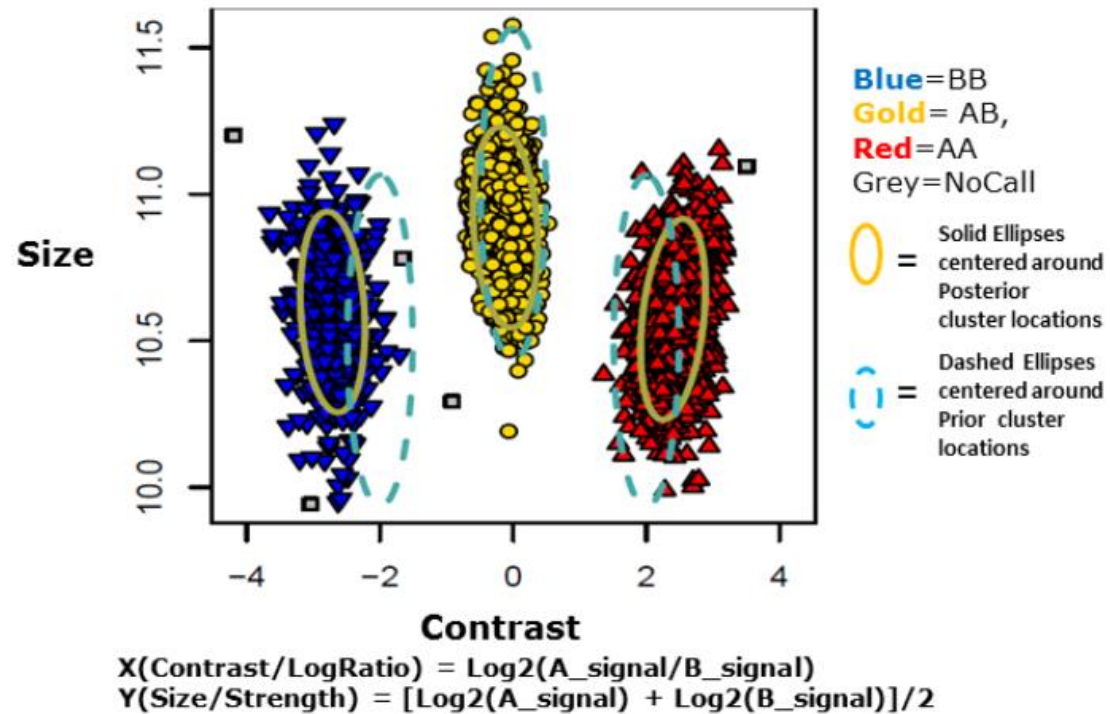




Figure 2.2 SNP Cluster Plot produced by the SNPolisher package, via the *Ps\_Visualization* function.



# GÉNOTYPAGE PAR CODAGE

# Exemple 1: KASP

|                           | <b>LC480®</b>                             | <b>Biomark®</b>   | <b>Tecan Infinite</b>                      |
|---------------------------|---|---|--|
| <b>Support</b>            | <b>Plaque 384</b>                         | <b>Puces : - 48*48<br/>- 96*96</b>                            | <b>Plaque 96, 394 ou 1536</b>              |
| <b>Type de projet</b>     | <b>Quelques SNPs</b>                      | <b>Au moins 48 SNPs</b>                                       | <b>De 1 à 1000 SNPs</b>                    |
| <b>Volume réactionnel</b> | <b>µL</b>                                 | <b>nL</b>   | <b>µL</b>                                  |
| <b>Débit</b>              | <b>20 plaques/jour :<br/>7680 données</b> | <b>2 à 3 puces/jour :<br/>18240 données<br/>(puces 96*96)</b> | <b>60 plaques/jour :<br/>23040 données</b> |
| <b>Coût / data</b>        | <b>0,16 € HT</b>                          | <b>0,1 € HT</b>   | <b>0,08 € HT</b>                           |
|                           | <b>Moyen débit</b>                        | <b>Haut ótrès haut débit</b>                                  |  |

## Exemple 2: Illumina Golden Gate Veracode (BeadXpress) <http://uthscsa.edu/csb/CSBPDFfiles/genomics-GoldenGateGenotypingOverview.pdf>

É Plusieurs degrés de multiplexage disponibles :

ó 48plex

ó 96plex

ó 144plex

ó 192plex

ó 384plex

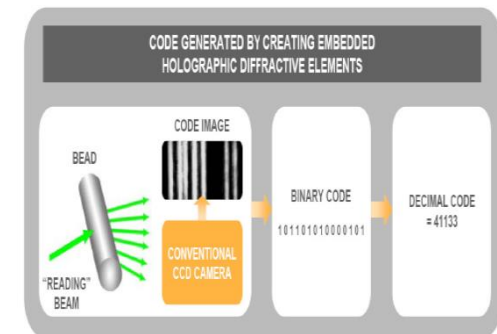
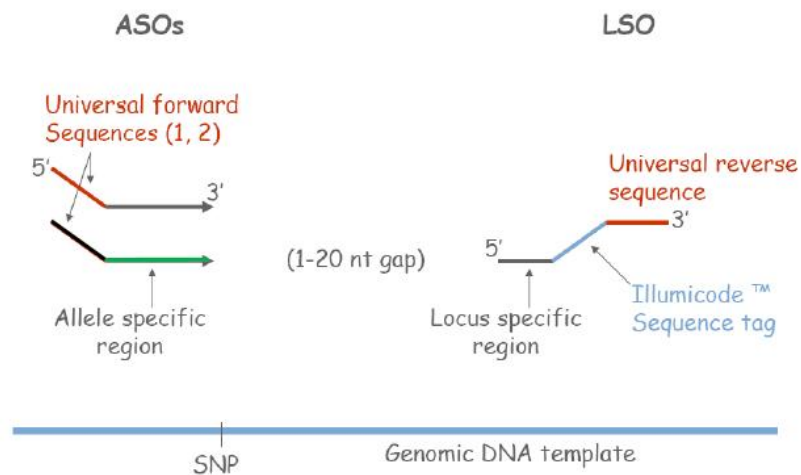
É Nombre minimum d'individus génotypés : 480 ( $5 \times 96$ ) puis  
par multiple de 96

É Coût à la data : de 0,045p HT à 0,22p HT selon  
configuration (nbre ech/nbre de SNP)

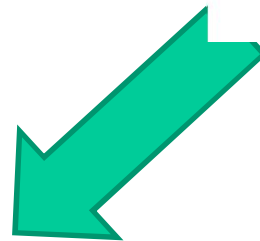
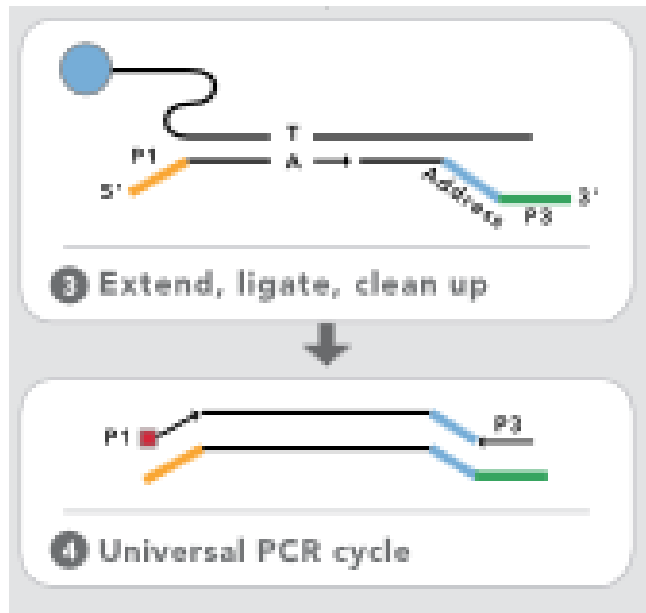
Plateforme GENTYANE INRA Clermont-Ferrand (C. Poncet)

# Choix des SNP

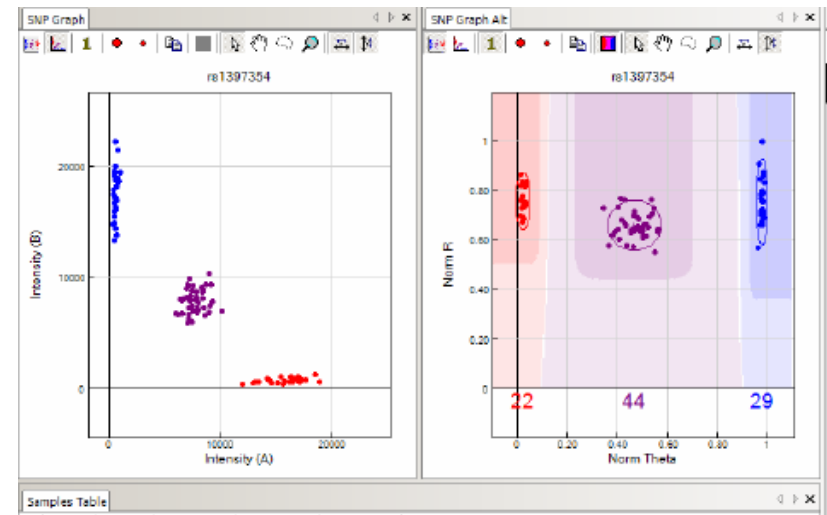
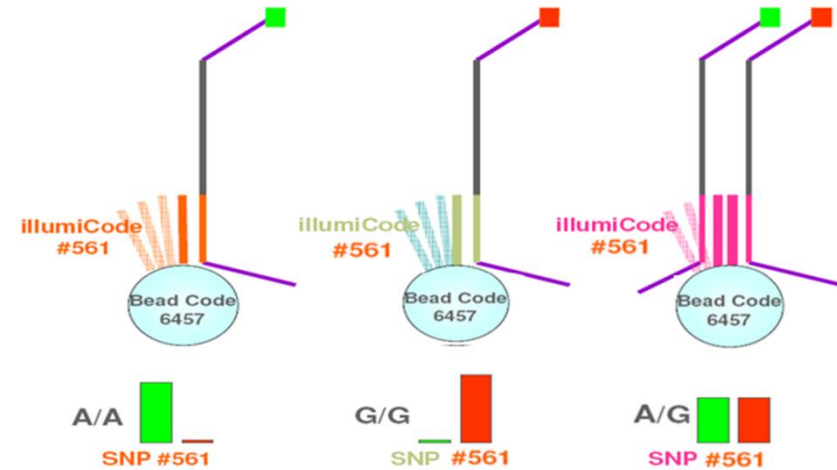
- É Soumission à Illumina d'un fichier en ligne
- É Besoin de 50bp autour du SNP
- É Obtention d'un score de qualité par SNP (Illumina, résultat instantané)
- É Choix définitif des SNP intégrant score et position
- É Commande de l'OPA et design des amorces par ILLUMINA (2 mois)
- É Réception OPA et génotypage (1 semaine pour 480 individus)



# Présentation de la Chimie GoldenGate



Hybridization to VeraCode Beads – one SNP per bead type



# Exemple 3: Axiom Affymetrix Genetitan

## Solution automatisée de génotypage SNP par hybridation



### Genetitan Affymetrix

Hybridation / lavage / lecture  
des plaques (PEG) 96 puces

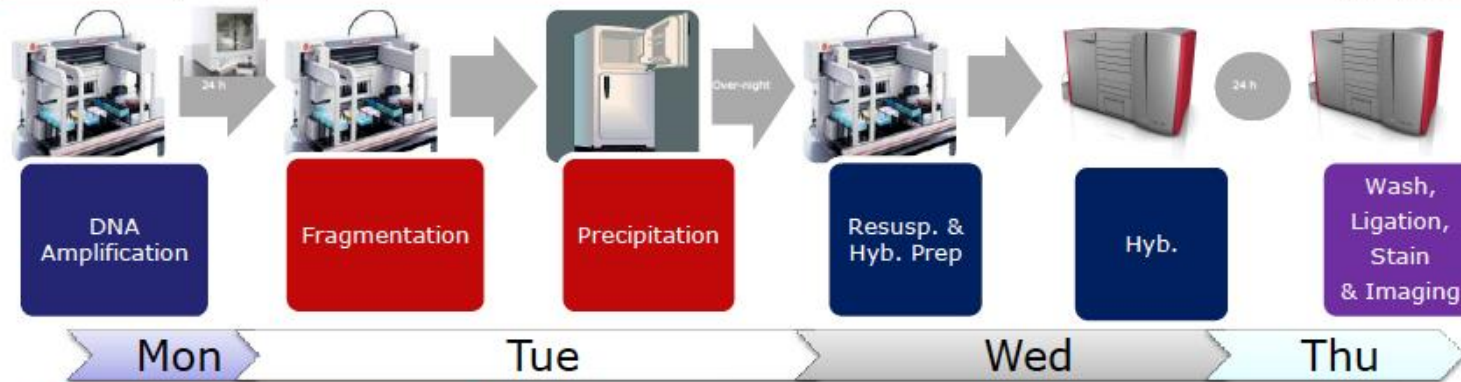
### Biomeck FxP Beckman

Préparation des échantillons avant  
passage sur GeneTitan  
Permet un débit de 8 plaques  
(PEG)/semaine

[http://media.affymetrix.com/support/downloads/manuals/axiom\\_genotyping\\_solution\\_analysis\\_guide.pdf](http://media.affymetrix.com/support/downloads/manuals/axiom_genotyping_solution_analysis_guide.pdf)

Plateforme GENTYANE INRA Clermont-Ferrand (C. Poncet)

# Axiom™ Genotyping Solution



- Throughput**
  - Highest throughput/operator & Fewest interventions
  - **768 samples per week**
- Hands-on-time**
  - Lowest in industry
  - **Less than 2.5 hr/plate**
- Ease of use**
  - Color coded reagent tubes, templates
  - **Reagents, buffers, GeneTitan consumables**
- Support**
  - Partnership with Beckman
  - **Methods supported by Affymetrix FAS**

# Schematic of the Axiom Assay



Target prep

Hybridization

Ligation

Signal amplification

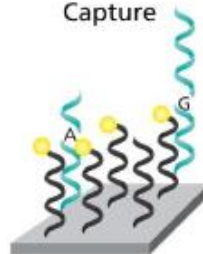
Amplify



Fragment



Capture



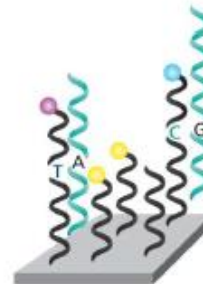
+

Label

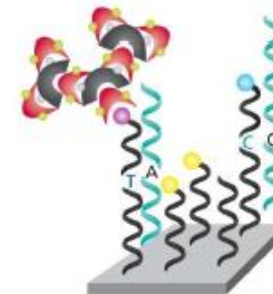


Labeled solution probe

Differentiate



Stain and image

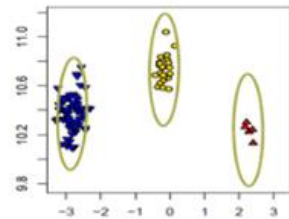




# Le codage en génotype

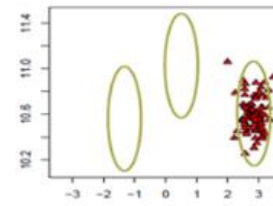
Figure 3.3 Cluster Plot examples and descriptions of the seven SNP classification categories. OTV SNPs are discussed further in [Adjust Genotype Calls for OTV SNPs on page 30](#).

## Poly High Resolution



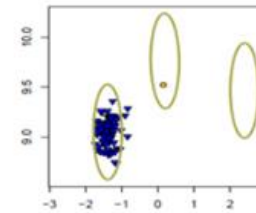
- Good cluster resolution
- At least 2 examples of minor allele

## Mono High Resolution



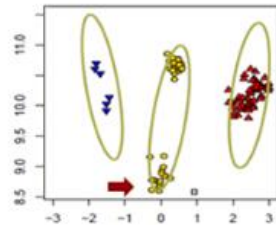
- High absolute Contrast value
- All genotyped samples are monomorphic

## No Minor Homozygote



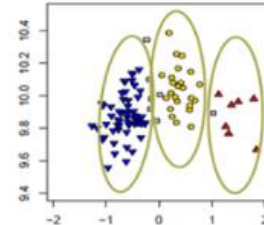
- Good cluster resolution
- No minor homozygous examples

## Off-Target Variant



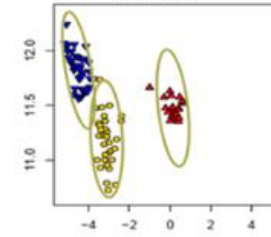
- Has an *Off-Target Variant* cluster (arrow).

## Call Rate Below Threshold



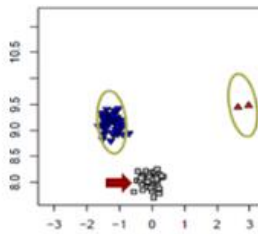
- Call Rate is below threshold, but all other cluster properties are good

## Other

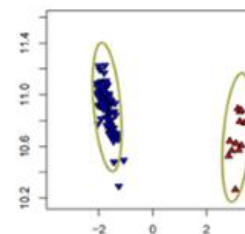


- One or more cluster properties are below threshold values

## Hemizygous



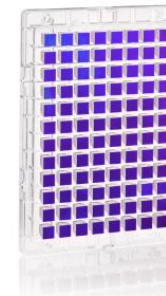
Y SNP  
females are set to  
No Call (arrow)



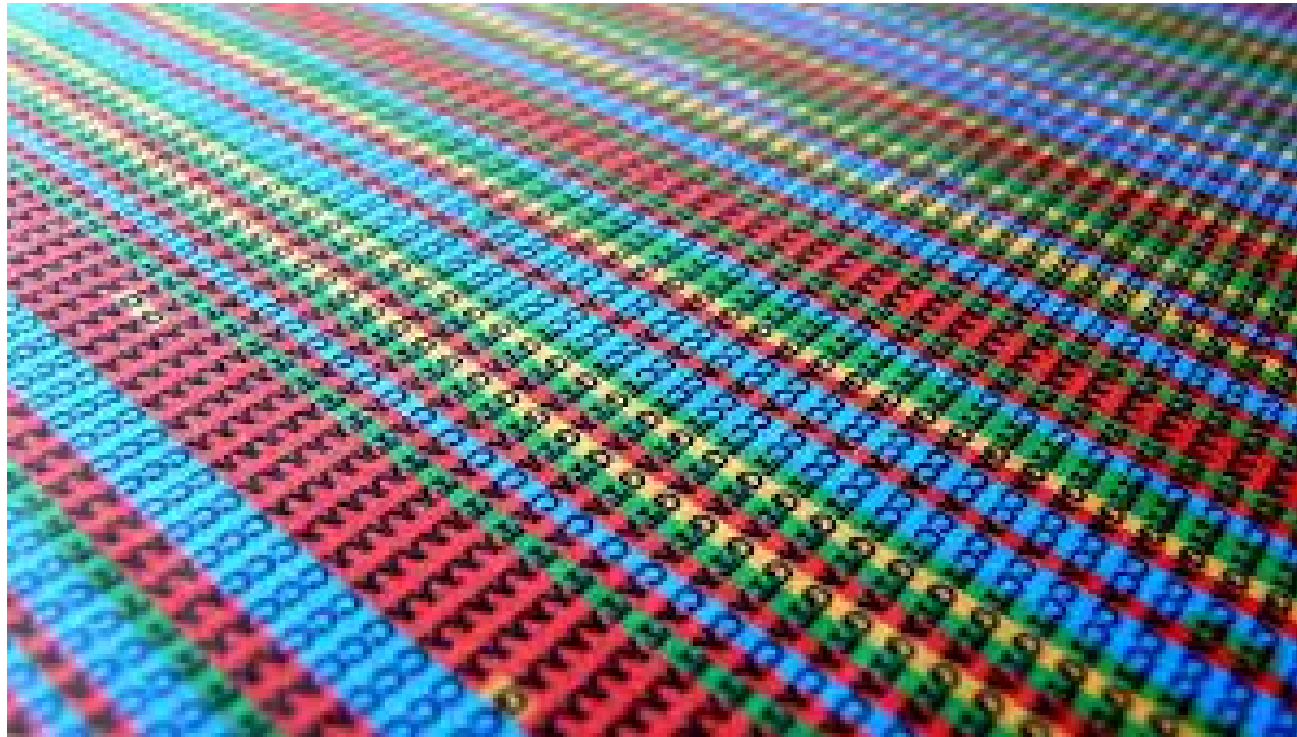
Mito SNP

## Avantages et Inconvénients de cette technologie

- + Débit très important
- + Besoin humain limité / quantité de données produites
- + Analyse rapide sous R
- + Coût à la data très faible de 0,08€ HT à 0,00006€ HT selon config
- + Souplesse nombre de SNP de 1500 à 675000 par puce
- Format 96 individus imposé
- Minimum de commande pour 480 individus
- Ticket d'entrée minimum de 60k€ HT



**Technologie proposée sur la plateforme  
Gentyane depuis Septembre 2013**



# **GÉNOTYPAGE PAR SÉQUENÇAGE**

# Génotypage par séquençage

É Séquencer tout le génome : oui mais cher !!

É Séquencer une partie du génome :

- ó Les gènes : transcriptome

- ó Des régions réparties sur le génome mais non  
« choisies » : GBS avec enzyme de restriction

- ó Des régions choisies : capture

# Génotypage par séquençage

É Trois étapes:

- ó Fabrication de la « banque » : mélange de fragments d'ADN : peut être réalisée dans un labo de biologie moléculaire classique + compétence BM
- ó Séquençage : nécessité des machines très chères et qui évoluent très vite
- ó Analyse bio-informatique : nécessité des ordinateurs très puissants + compétence informatique

# Les banques

É Banques d'ADNc : attention aux tissus choisis selon la question et si besoin standardisation

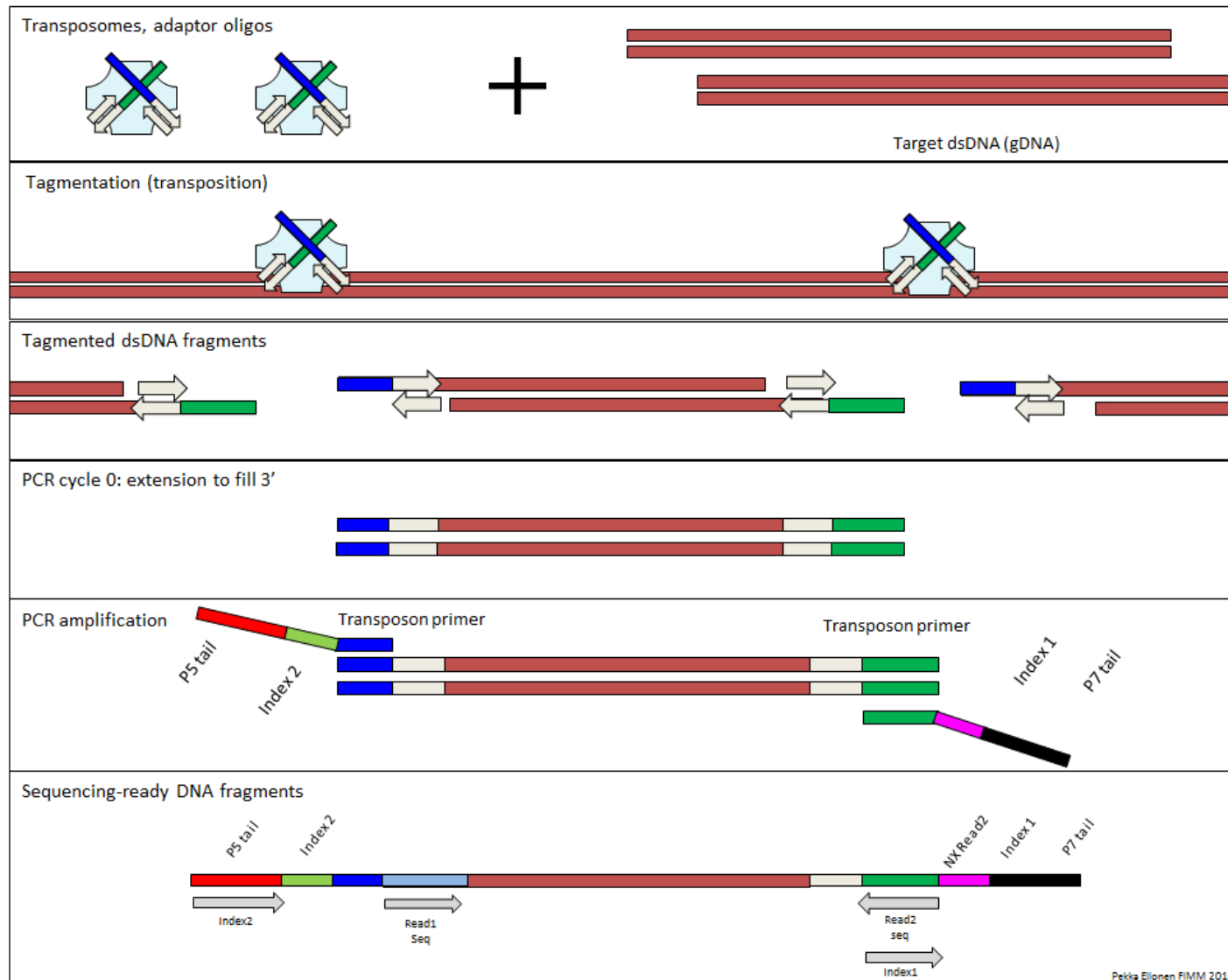
É Banque d'ADN génomique :

- ó ADN total fragmenté de manière aléatoire

- ó ADN issu de l'amplification d'une partie du génome par PCR: séquences cibles (ex. barcoding) ou enzymes de restriction

- ó ADN issu de l'hybridation sur des séquences connues : capture

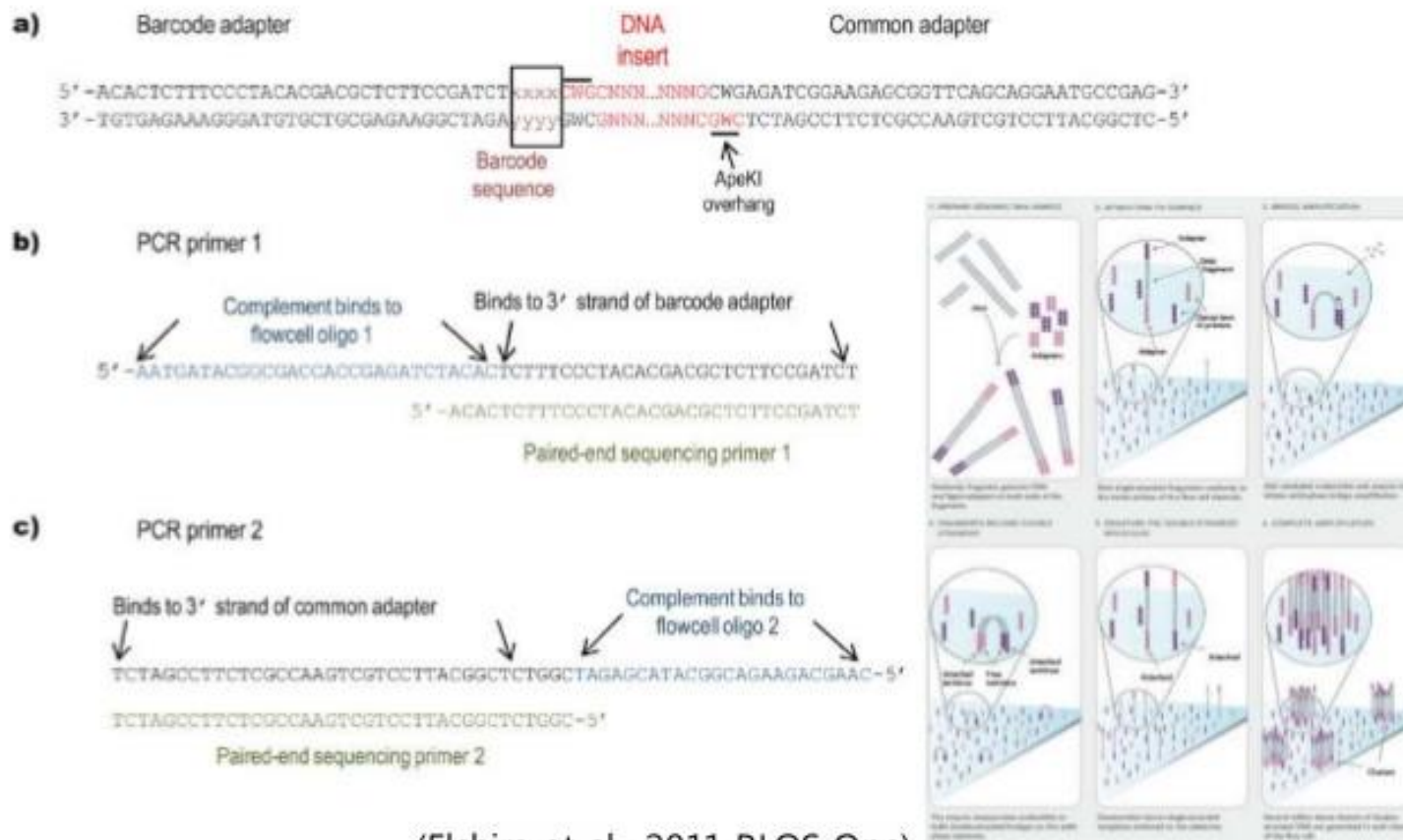
# Les banques



Système  
NEXTERA  
de Illumina

# Exemple : GBS Elshire et al.2011 PLOS One

Digestion de l'ADN génomique par une enzyme de restriction (ici ApeKI)

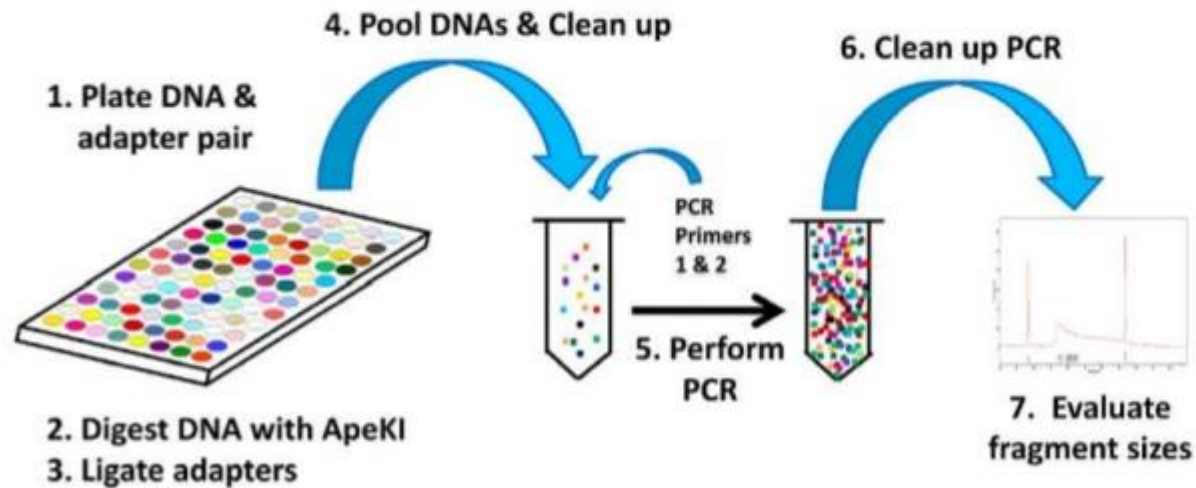


(Elshire et al., 2011 PLOS One)



# Exemple : GBS Elshire et al.2011 PLOS One











## GBS library construction



(Elshire et al., 2011 PLOS One)

Il existe beaucoup de variantes RADseqí .

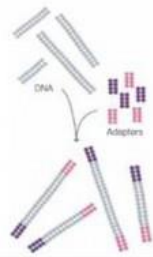
# Séquençage: plusieurs plateformes selon les besoins

| Séquenceurs 2 <sup>ème</sup> génération |   |   |  |   |   |   |   |   |   |   |   |
|---|---|---|--|---|---|---|---|---|---|---|---|
| Société                                 | Roche   |   |  | Illumina  |   |   | Life Technologies   |   |   |   |   |
| Plateforme                              |  |  |  |  |  |  |  |  |  |  |  |
| Technologie                             | GS Junior   | 454   |  | MiSeq   | HiSeq 1000  | HiSeq 2000  | Genome Analyzer Ix  | Ion Torrent PGM   | SOLiD 4   | SOLiD 5500  | SOLiD 5500xl  |
| Technologie                             | Titanium  | FLX Titanium<br>FLX +   |  |   |   |   |   | Chip 314<br>Chip 316<br>Chip 318  |   |   |   |
| Acides nucléiques (matrice)             |   |   |  |   |   |   |   |   |   |   |   |
| Ligation adaptateurs                    |   |   |  |   |   |   |   |   |   |   |   |
| Méthode d'amplification                 | PCR en émulsion   |   |  | « Bridge PCR »  |   |   | PCR en émulsion   |   |   |   |   |
| Méthode de séquençage                   | Synthèse (Pyroséquençage)   |   |  | Synthèse  |   |   | Ligation  |   |   |   |   |
| Durée de séquençage/run                 | 10h   | 10h 20h   |  | 26h   | 8jrs  | 8jrs  | 14jrs   | 2h  | 12jrs   | 8jrs  | 8jrs  |
| Capacité (Mb) séquençage/run            | 50  | 500 900   |  | 1500  | 100000  | 200000  | 95000   | >10 >100 >1000  | 70000   | 80000   | 150000  |
| Taille moyenne des reads                | 400   | 400 700   |  | 150+150   | 100+100   | 100+100   | 150+150   | 100 >100 >100   | 50+35   | 75+35   | 75+35   |
| Coût (\$) /run                          | 1100  | 6200  |  | 750   | 10000   | 20000   | 11500   | 500 750 950   | 8150  | 6100  | 10500   |
| Coût machine + annexes ((K\$))          | 110+25  | 500+30  |  | 125   | 560   | 690   | 250   | 50+20   | 480+55  | 350+55  | 600+55  |
| Exactitude de séquençage (%)            | 99  | 99  |  | 99,9  | 99,9  | 99,9  | 99,9  | 99  | 99,95   | 99,95   | 99,99   |

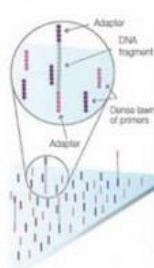
# Séquenceur Illumina : évolution rapide

|                              | <br>MiniSeq System                   | <br>MiSeq Series  | <br>NextSeq Series   | <br>HiSeq Series  | <br>HiSeq X Series*  |
|------------------------------|---|---|---|--|---|
| <b>Key Methods</b>           | Amplicon, targeted RNA, small RNA, and targeted gene panel sequencing.  | Small genome, amplicon, and targeted gene panel sequencing.   | Everyday exome, transcriptome, and targeted resequencing.   | Production-scale genome, exome, transcriptome sequencing, and more.  | Population- and production-scale whole-genome sequencing.   |
| <b>Maximum Output</b>        | 7.5 Gb  | 15 Gb   | 120 Gb  | 1500 Gb  | 1800 Gb   |
| <b>Maximum Reads per Run</b> | 25 million  | 25 million <sup>†</sup>   | 400 million   | 5 billion  | 6 billion   |
| <b>Maximum Read Length</b>   | 2 × 150 bp  | 2 × 300 bp  | 2 × 150 bp  | 2 × 150 bp   | 2 × 150 bp  |
| <b>Run Time</b>              | 4–24 hours  | 4–55 hours  | 12–30 hours   | <1–3.5 days (HiSeq 3000/HiSeq 4000)<br>7 hours–6 days (HiSeq 2500)   | <3 days   |
| <b>Benchtop Sequencer</b>    | Yes   | Yes   | Yes   | No   | No  |
| <b>System Versions</b>       | <ul style="list-style-type: none"> <li>• MiniSeq System for low-throughput targeted DNA and RNA sequencing</li> </ul> | <ul style="list-style-type: none"> <li>• MiSeq System for targeted and small genome sequencing</li> <li>• MiSeq FGx System for forensic genomics</li> <li>• MiSeqDx System for molecular diagnostics</li> </ul> | <ul style="list-style-type: none"> <li>• NextSeq 500 System for everyday genomics</li> <li>• NextSeq 550 System for both sequencing and cytogenomic arrays</li> </ul> | <ul style="list-style-type: none"> <li>• HiSeq 3000/HiSeq 4000 Systems for production-scale genomics</li> <li>• HiSeq 2500 Systems for large-scale genomics</li> </ul> | <ul style="list-style-type: none"> <li>• HiSeq X Five System for production-scale whole-genome sequencing</li> <li>• HiSeq X Ten System for population-scale whole-genome sequencing</li> </ul> |

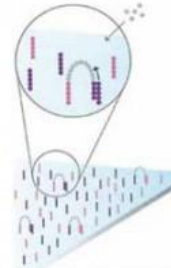
# Exemple séquençage Hiseq



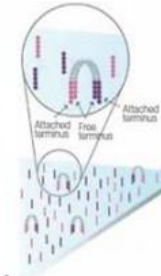
**Figure 1**  
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.



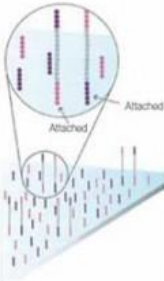
**Figure 2**  
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.



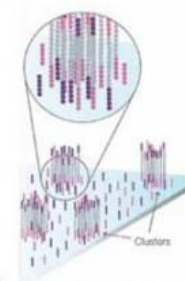
**Figure 3**  
Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.



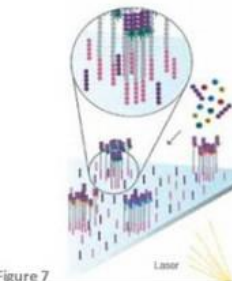
**Figure 4**  
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.



**Figure 5**  
Denaturation leaves single-stranded templates anchored to the substrate.



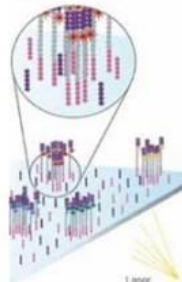
**Figure 6**  
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.



**Figure 7**  
The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.



**Figure 8**  
After laser excitation, the emitted fluorescence from each cluster is captured and the first base is identified.



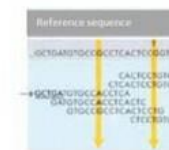
**Figure 9**  
The next cycle repeats the incorporation of four labeled reversible terminators, primers, and DNA polymerase.



**Figure 10**  
After laser excitation, the image is captured as before, and the identity of the second base is recorded.



**Figure 11**  
The sequencing cycles are repeated to determine the sequence of bases in a fragment, one base at a time.



**Figure 12**  
The data are aligned and compared to reference, and sequencing differences are identified.

# Bio-informatique

É Enormément de données à traiter : penser à la mémoire !!! : machines dédiées

É Logiciels :

ó Quelques uns sous Windows dont CLC genomics: convivial, payant, peu modulable

ó Beaucoup d'application sous Linux : gratuit, très modulable, moyennement convivial

ó Un système intermédiaire : Galaxy  
<https://usegalaxy.org/>

# Bio-informatique : plusieurs étapes

É Démultiplexage : séparation des séquences par génotypes

É Nettoyage : enlever les adaptateurs, bonne qualité, longueur minimale

É Alignement des séquences:

ó Sur une séquence de référence

ó Sur un assemblage de novo

É Détection des SNP

ó Profondeur de séquence

ó Fréquence minimale d'un allèle





**CONCLUSIONS**

# Conclusions

É Une multitude de méthodes avec des avantages et inconvénients

É Adapter la méthode à la question posée et aux connaissances de l'espèce d'intérêt

É Génotypage de plus en plus spécialisé à sous traiter tout ou partie (Macrogen, Keygene, LGC genomicsí )





# Conclusions

- É Dans l'avenir séquençage complet systématique??
- É Outils de bio-informatique cruciaux, problème de stockage
- É Changement des outils d'analyse des données de génotypage et de phénotypage maintenant en très grand nombre
- É Des outils impressionnants pour mieux comprendre la diversité du monde vivants !!!!

