

*FC bio-informatique Fev 2015*

NGS : des reads aux SNPs  
(Vincent Ranwez – Montpellier SupAgro)

# *NGS : des reads aux SNPs*

- Introduction NGS
  - Exemples d'applications
  - Rapide comparatif des technologies
  - Les données reçues
  - Nettoyage des reads
- Assemblage de-novo
- Mapping sur un génome connu
- Détection de SNP
- Conclusions et discussions

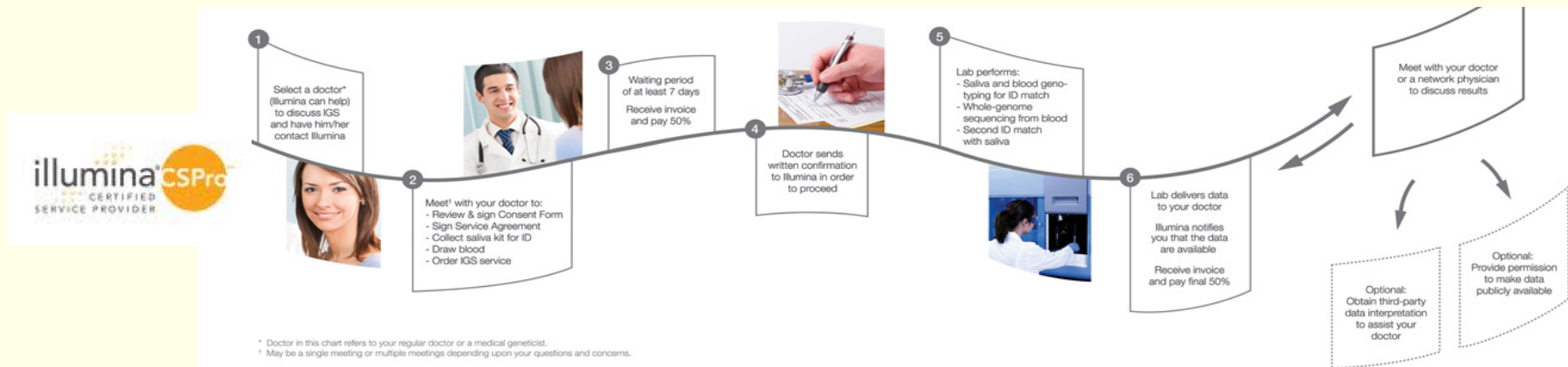
- Santé (humaine) : lier une maladie à une zone du génome
  - Reséquençage massif 1000 génomes, 2500 génomes etc.



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

## ■ Santé humaine :

- médecine personnalisée 😊
  - Assurance personnalisée ☹️
- ⇒ Nombreuses questions médicales, juridiques et éthiques



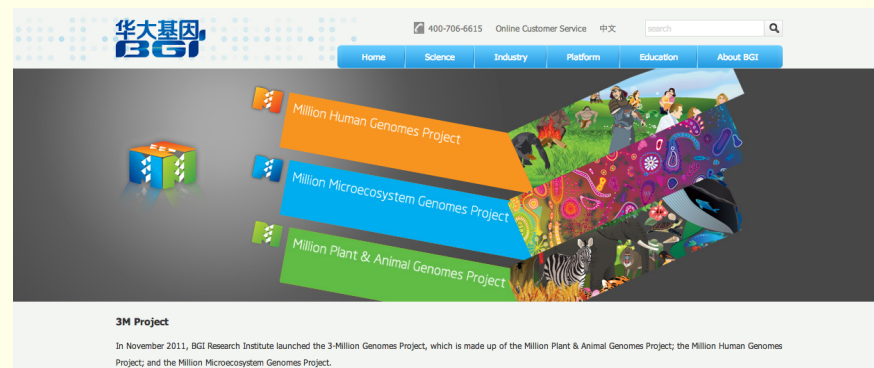
- Etude de la diversité inter et intra espèces
  - Phylogénie, flux de gènes, protection de la bio-diversité



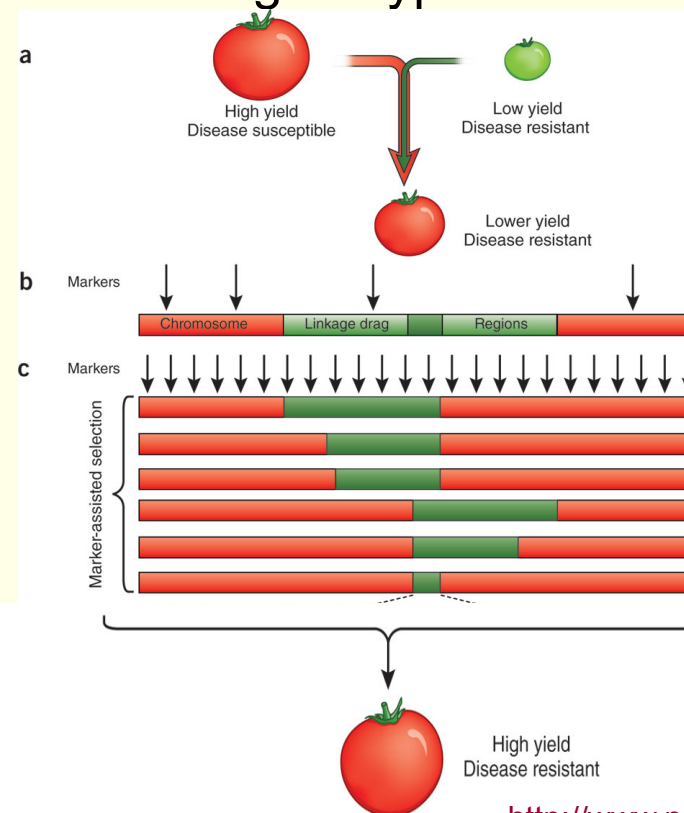
<http://www.1001genomes.org/>



<http://genome10k.soe.ucsc.edu/>



- Identification de marqueurs génomiques sur l'ensemble du génomes : SNP
  - ⇒ Sélection assistée par marqueurs
  - ⇒ lien génotype X phénotype
  - ⇒ interactions génotype X environnement



## Rapide comparatif des technologies

### NGS platforms overview



#### Technology

454

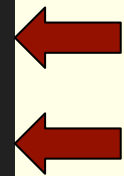
Solexa

SOLiD

Platform:	GS 20	FLX	TI	GA	GA II	1	2
Reads:	500 k	500 k	1 M	28 M	80 M		40 M 115 M

#### Fragment

Read length:	100	200	350	35	50	75	25 35
Run time:	6 hr	7 hr	9 hr	3 d	3 d	4 d	6 d 5 d
Yield:	50 Mb	100 Mb	400 Mb	1 Gb	4 Gb	6 Gb	1 Gb 4 Gb
Yield per hour:	8,3 Mb/hr	14,3 Mb/hr	44,4 Mb/hr	13,9 Mb/hr	55,6 Mb/hr	62,5 Mb/hr	6,9 Mb/hr 33,3 Mb/hr
Images:	11 GB	13 GB	27 GB	500 GB	1.1 TB	1.7 TB	1.8 TB 2.5 TB
PA Disk:	3 GB	3 GB	15 GB	175 GB	300 GB	350 GB	300 GB 750 GB
PA CPU:	10 hr	140 hr	220 hr	100 hr	70 hr	100 hr	NA NA
SRA:	500 MB	1 GB	4 GB	30 GB	50 GB	75 GB	100 GB 140 GB



## Rapide comparatif des technologies

www.biomedcentral.com - Table

www.biomedcentral.com/1471-2164/13/341/table/T1 *Quail et al. BMC Genomics 2012*

Cette page est en **anglais** Voulez-vous la traduire ?

**Table 1**  
**Technical specifications of Next Generation Sequencing platforms utilised in this study**

Platform	Illumina MiSeq	Ion Torrent PGM	PacBio RS	Illumina GAIIx	Illumina HiSeq 2000
Instrument Cost*	\$128 K	\$80 K**	\$695 K	\$256 K	\$654 K
Sequence yield per run	1.5-2Gb	20-50 Mb on 314 chip, 100-200 Mb on 316 chip, 1Gb on 318 chip	100 Mb	30Gb	600Gb
Sequencing cost per Gb*	\$502	\$1000 (318 chip)	\$2000	\$148	\$41
Run Time	27 hours***	2 hours	2 hours	10 days	11 days
Reported Accuracy	Mostly > Q30	Mostly Q20	<Q10	Mostly > Q30	Mostly > Q30
Observed Raw Error Rate	0.80 %	1.71 %	12.86 %	0.76 %	0.26 %
Read length	up to 150 bases	~200 bases	Average 1500 bases**** (C1 chemistry)	up to 150 bases	up to 150 bases
Paired reads	Yes	Yes	No	Yes	Yes
Insert size	up to 700 bases	up to 250 bases	up to 10 kb	up to 700 bases	up to 700 bases
Typical DNA requirements	50-1000 ng	100-1000 ng	~1 µg	50-1000 ng	50-1000 ng



## *Rapide comparatif des technologies*

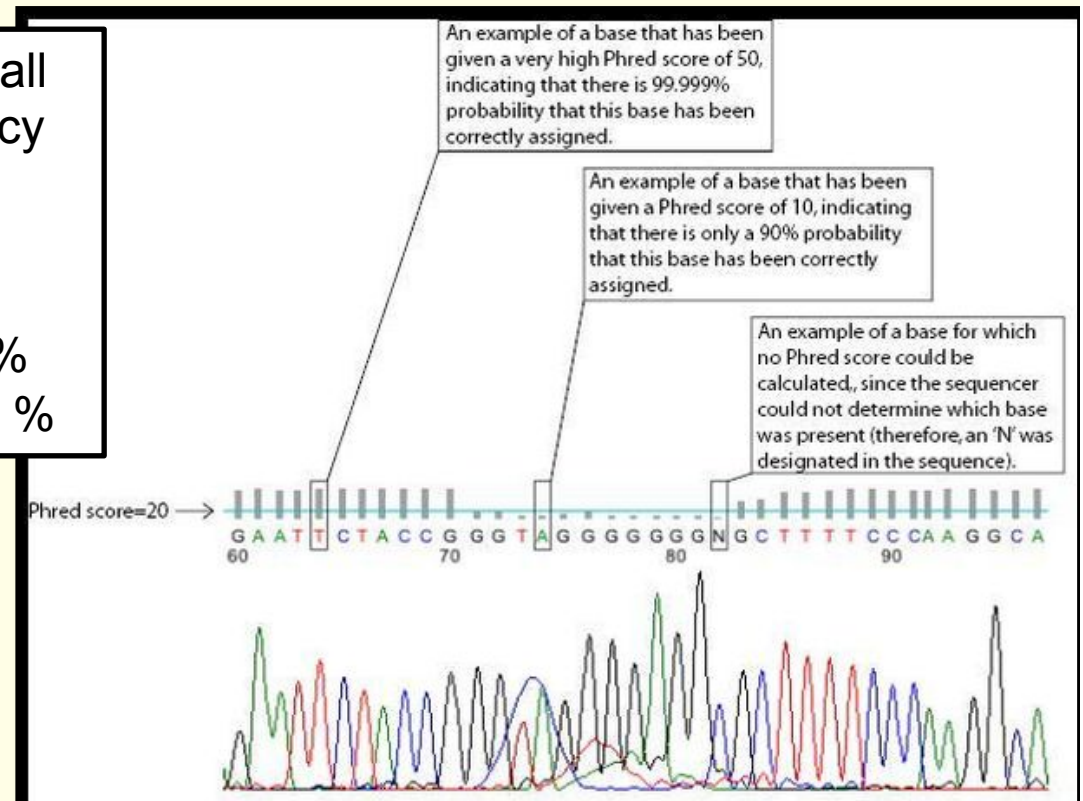
- Les comparatifs sont souvent obsolètes
- Chaque technologie a ses atouts
  - HiSeq, forte couverture reads fiables mais.... reads très courts
  - PacBio, reads très longs mais ... fort taux d'erreurs
- ⇒ Rien n'oblige à n'utiliser qu'une technologie !
- Comment choisir ?
  - En fonction de la question que vous souhaitez traiter !
  - Le prix par base ne fait pas tout
  - Mélange de science et de pression commerciale
  - Tous les prestataires ne sont pas égaux...

**Discuter avec des utilisateurs/clients  
avant de faire un choix**

## Ce que l'on reçoit

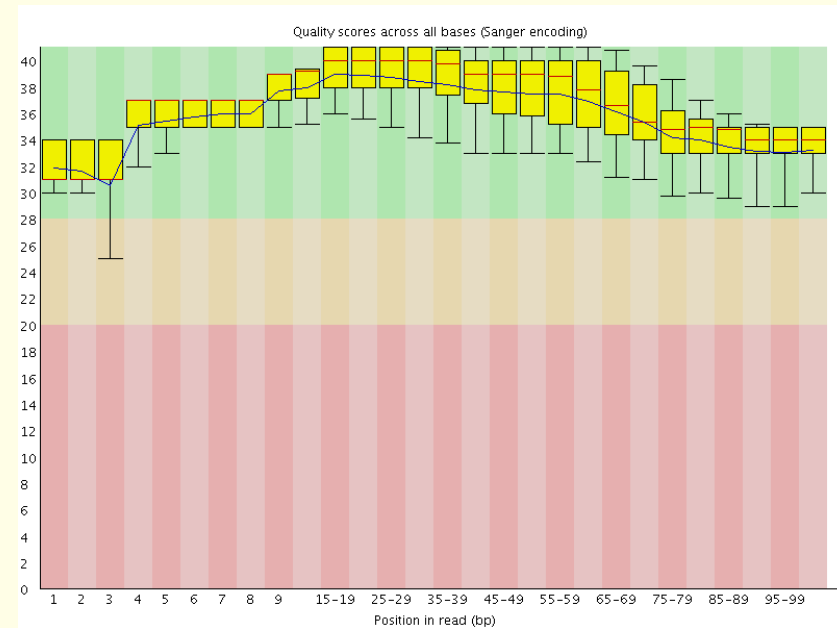
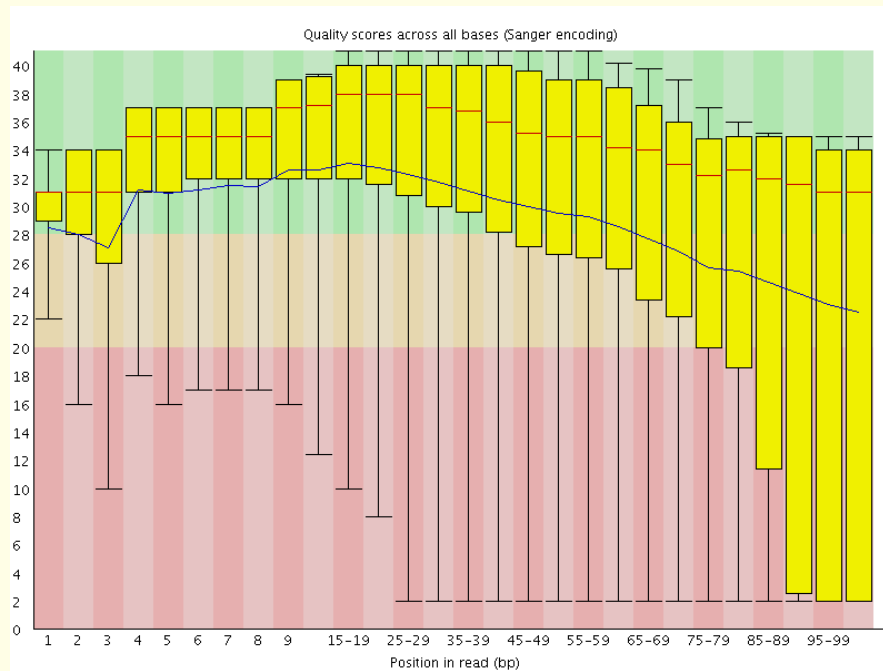
- Plein de séquences avec pour chaque nucléotide de chaque séquence un indice de qualité (Q)
- Q est obtenu à partir de la probabilité P que la base soit erronée :  **$Q = -10 \log(P)$**

Q	P	Base call accuracy
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1000	99.9 %
40	1 in 10000	99.99 %
50	1 in 100000	99.999 %



- **Adaptateurs**
  - ⇒ Normalement enlevés avant livraison
- **Les extrémités des reads sont souvent de moins bonne qualité**
  - ⇒ Elimination des bases extrêmes de faible qualité
  - ⇒ Masquages des bases internes de faible qualité
- **ACCTTCTTTTCCACGTCTT (seuil 5)**  
**2387878882888988433**
- **CTTCTTTNCCACGT**

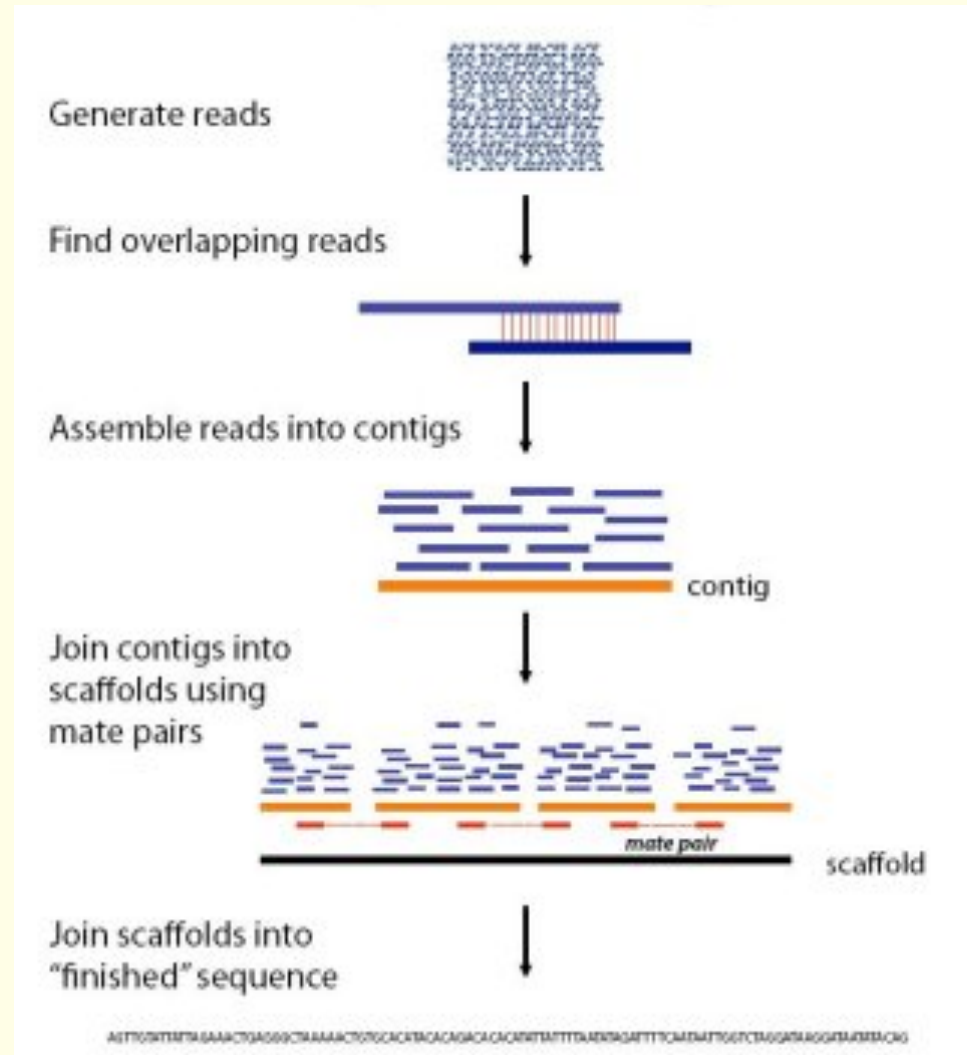
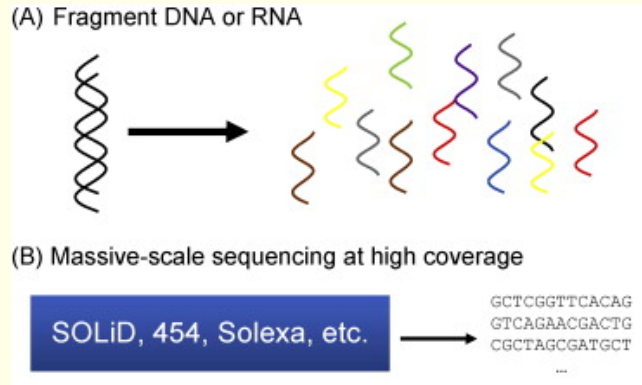
- S'assurer de la qualité avant toute analyse
  - ⇒ FastQC
- Exemple : Riz avant et après nettoyage



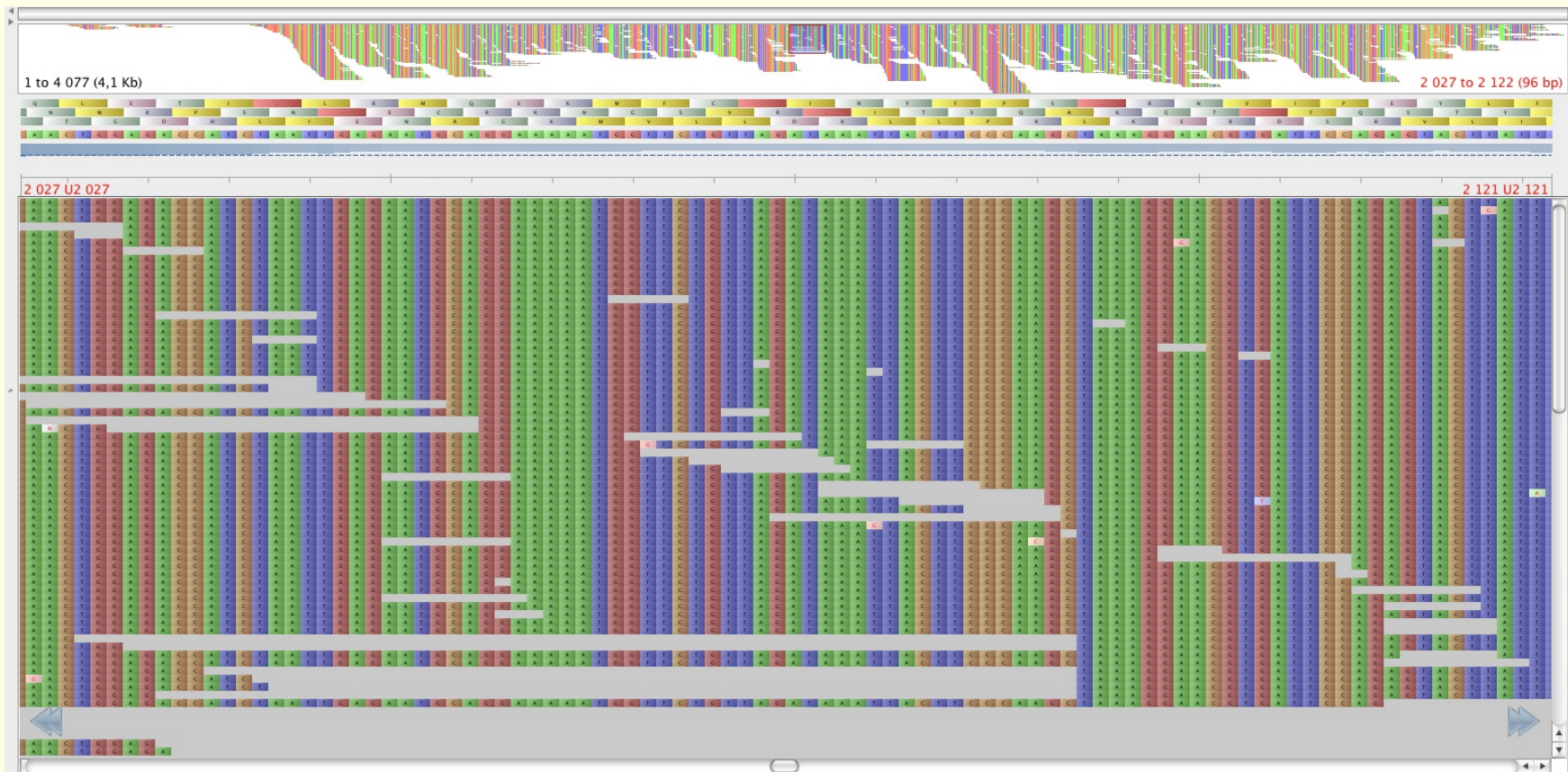
# *NGS : des reads aux SNPs*

- Introduction NGS
- **Assemblage de-novo**
  - Intuition, problématique
  - Premières approches « gloutonnes »
  - Approches utilisant des Graphes
  - Apports des paired-end
- Mapping sur un génome connu
- Détection de SNP
- Conclusions et discussions

■ Vision globale



- Vision globale



- Supposons que l'on séquence un tout petit (fragment de) génome

avec des reads de 10pbs on peut avoir :

CTCCCTGTCA

GTCATCTGTC

ACCCTCCCTG

- Comment retrouver le génome de départ ?
- CTCCCTGTCA

ACCCTCCCTG

GTCATCTGTC

CTCCCTGTCAACCCTCCCTGTGTCATCTGTC (génome 1)

- ACCCTCCCTGTGTCATCTGTC (génome 2)





*Premières approches gloutonnes*  
(e.g. SSAKE)

- Idées de base : construire le génome petit à petit en ajoutant un read après l'autre :

CTCCCTGTCA

GTCATCTGTC

ACCCTCCCTG

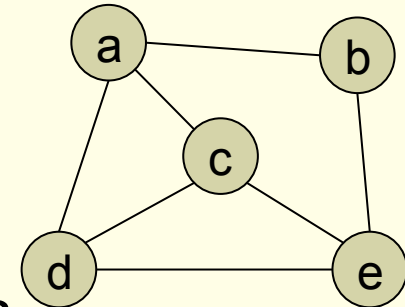
- Supposons que l'on parte de : CTCCCTGTCA (ajout?)  
CTCCCTGTCA



- Choisir le read (et la position) qui minimise le risque
  - ⇒ Plus grand chevauchement
  - ⇒ Construisez un exemple ou cela ne marche pas
- Choisir le read qui minimise le risque (alternatives)
  - ⇒ Plus petit ajout, plus grand rapport (chevauchement/ajout)

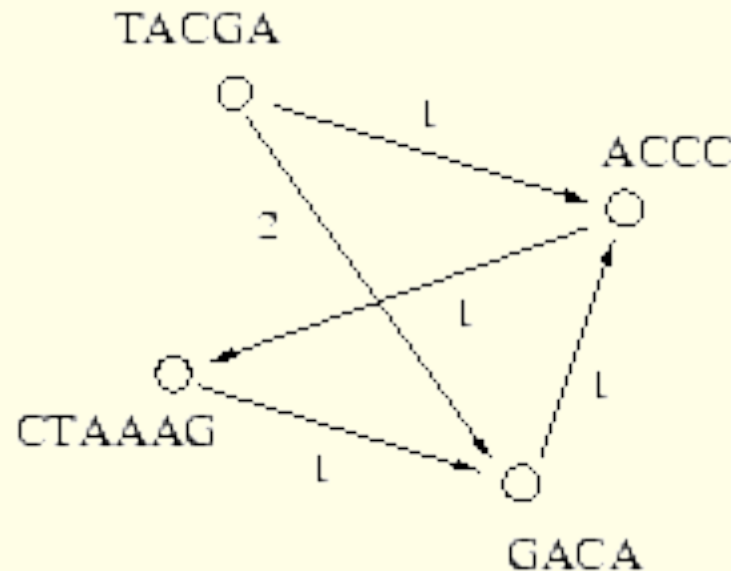
## Approche par graphe de chevauchement (e.g. CAP & Celera Assembler)

- Un graphe  $G = (V, E)$  est composé de :
  - $V$  : un ensemble de sommets/nœuds (Vertices)
  - $E$  : un ensemble d'arcs/arêtes (Edges)
- Un arc  $e = (u, v)$  est une paire de sommets
  - On peut associer des valeurs aux nœuds et aux arêtes
  - Les arêtes peuvent être ou non orientées
- Dans un graphe de chevauchement
  - Chaque séquence à assembler est un nœud
  - On met une arête entre deux nœuds si leurs séquences se chevauchent
  - Valeur d'une arête = taille du chevauchement
- Construisez le graphe de chevauchement pour  $E = \{TACGA, ACCC, CTAAAG, GACA\}$



*Approche par graphe de chevauchement  
(e.g. CAP & Celera Assembler)*

- Graphe pour  $E = \{TACGA,ACCC,CTAAAG,GACA\}$
- Comment se traduit le problème d'assemblage, que faut-il chercher dans le graphe de chevauchement ?

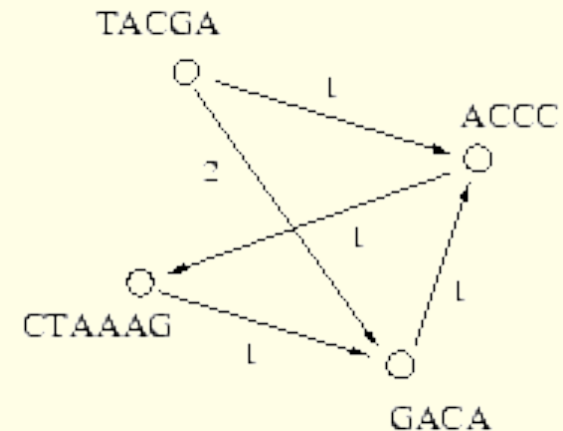


On cherche un chemin qui passe une et une seule fois par chaque sommet et qui soit de poids maximal

## Approche par graphe de chevauchement (e.g. CAP & Celera Assembler)

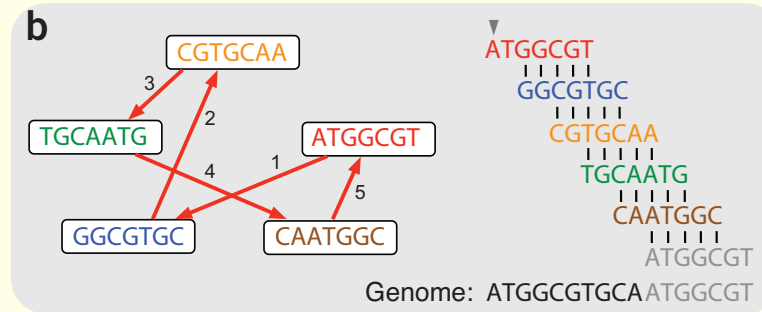
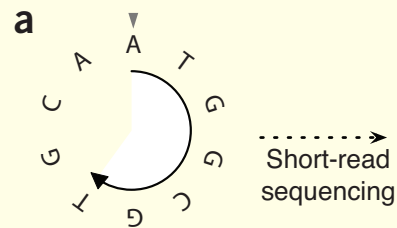
### ■ Limites de l'approche

- ⇒ La construction du graphe est coûteuse en temps de calculs  
combien de comparaisons pour 10 million de reads ?
- ⇒ La recherche d'un chemin passant une et une seule fois par  
chaque nœud est un problème difficile



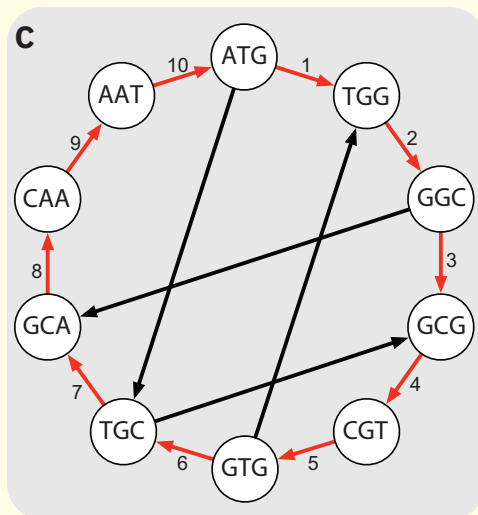
- Idée : décomposer les reads en fragments de taille k  
fixée à l'avance : k-mer

# Approche par graphe de « De Bruijn » (e.g. ABySS, trinity, velvet ...)

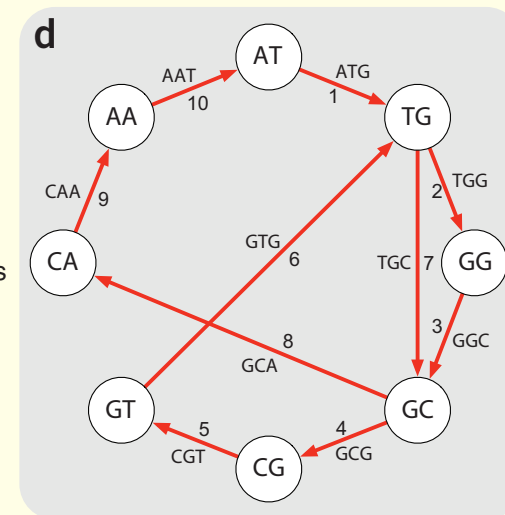
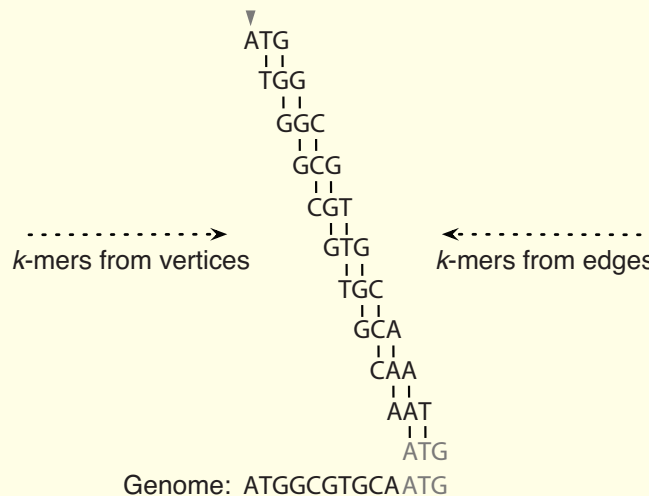


Vertices are  $k$ -mers  
Edges are pairwise alignments

Vertices are  $(k-1)$ -mers  
Edges are  $k$ -mers



**Hamiltonian cycle**  
Visit each vertex once  
(harder to solve)



**Eulerian cycle**  
Visit each edge once  
(easier to solve)

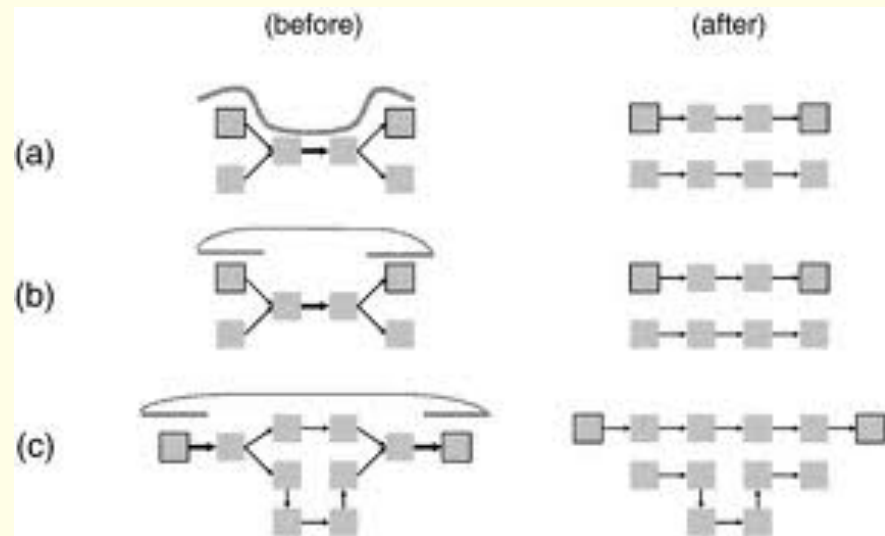
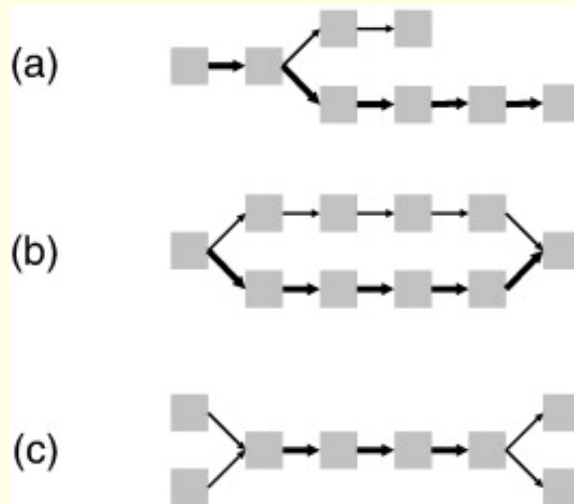
## Approche par graphe de « De Bruijn » (e.g. ABySS, trinity, velvet ...)

### ■ Résoudre les ambiguïtés

⇒ Profondeur,  
nombre de reads

⇒ Paired-end

⇒ Mated pairs

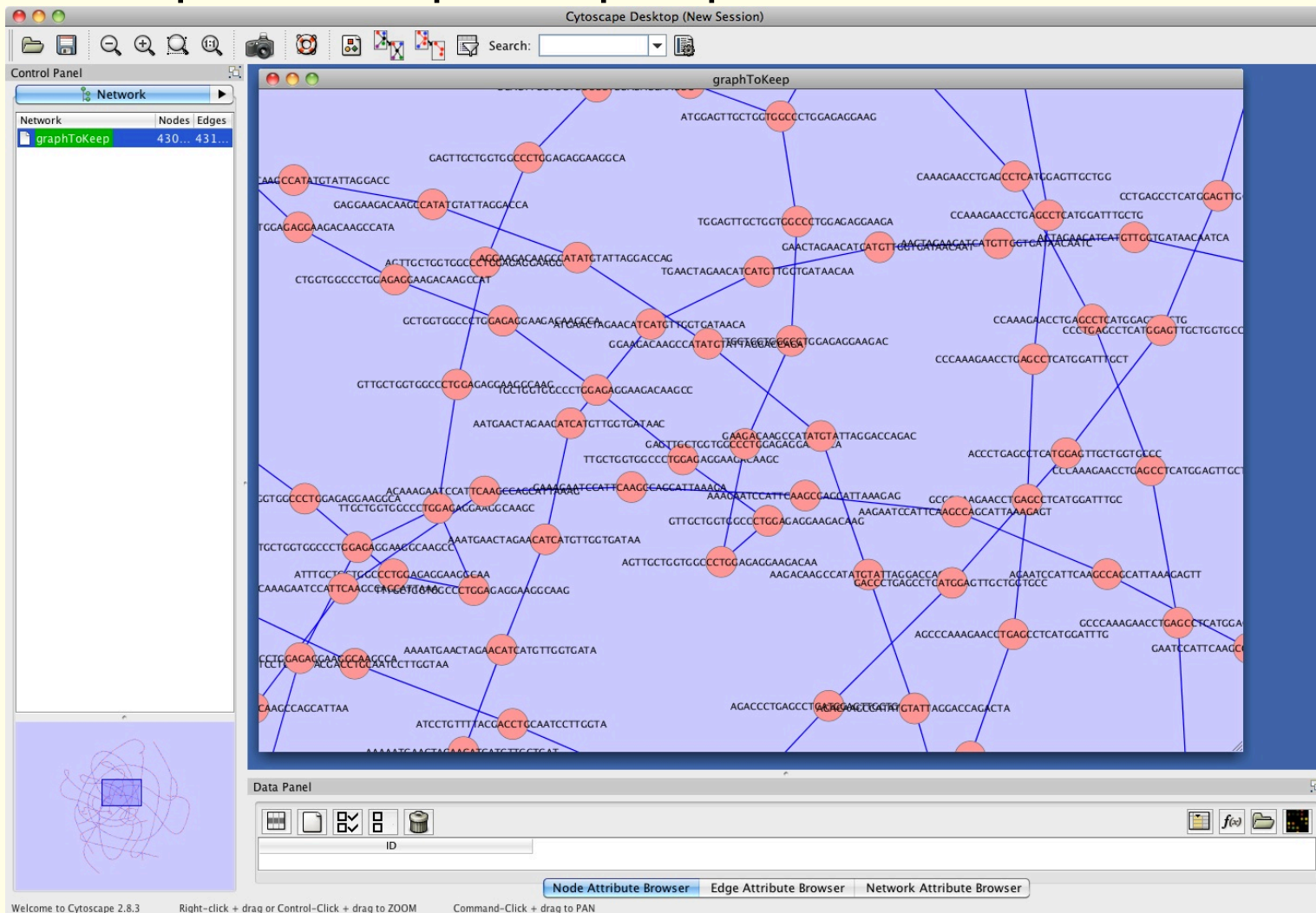


### ■ Les erreurs de séquençage augmentent la taille du graphe

⇒ Essayer de corriger avant d'assembler (e.g. ABySS)

# Approche par graphe de « De Bruijn » (e.g. ABySS, trinity, velvet ...)

- Tout petit exemple un peu plus réaliste



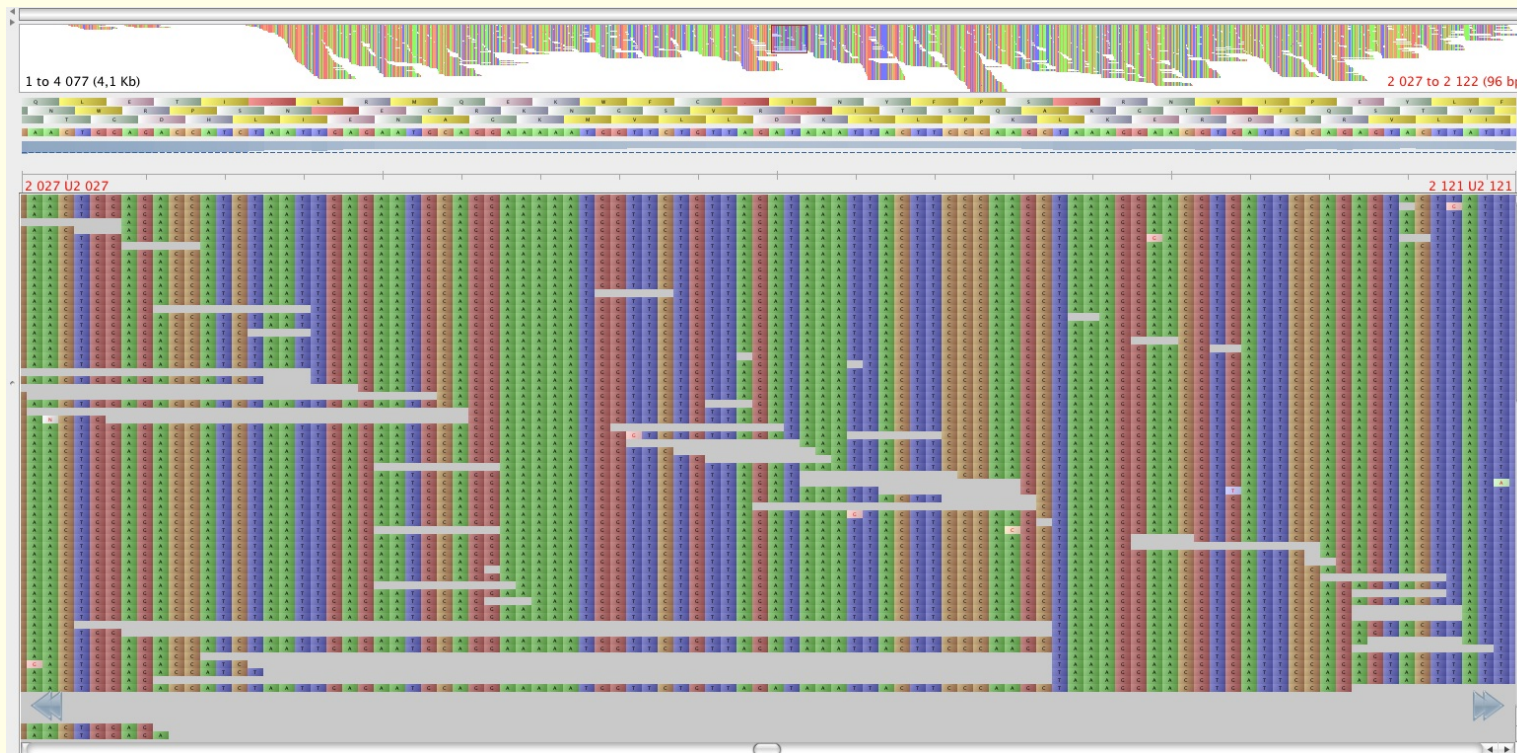
# *NGS : des reads aux SNPs*

- Introduction NGS
- Assemblage de-novo
- **Mapping sur un génome connu**
  - Approches par graines (et arbre des suffixes)
  - Approches par table de suffixe
- Détection de SNP
- Conclusions et discussions



## ■ Problématique

- On dispose d'un génome/transcriptome de référence
- On veut positionner les reads dessus



## Approches par graines et arbre de suffixe (e.g. SSAHA, SOAP, BLAT)

- Pour chaque read
  - Parcourir tout le génome et tester l'égalité read/portion génome
    - ⇒ Trop loooooong
- Chercher dans un livre de 500 pages où le mot «séquence » apparaît ☹

⇒ Avec un index ? 😊

- Pour chaque k-mer
  - Stocker ses position dans le génome
- Reads de 60 pb
  - Prendre des K-mer de 60 ?
    - ⇒ Lourd à stocker
  - Prendre des K-mer de 20 ?
    - ⇒ Tolérance aux erreurs
    - ⇒ Mieux, mais reste lourd à stocker

414 Index	Copyrighted Material
chimpanzee 36, 76, 105, 107, 112, 291, 303 chlorophyll 173 chloroplasts 3, 11, 48, 92-93 chromatin 14-15 remodelling 14-15 immunoprecipitation 384 chronic myeloid leukaemia 34 cis-regulatory region 396 cladistics 148, 150 classification 73 clone 55, 66, 202, 214 CLUSTAL-W 238 clustering 148-149 clustering coefficient 406 CODIS (Combined DNA Index System) 225 coevolution 138 co-expression patterns 385 codon usage 50 coiled-coils 354 Collins, F. 21 collision-induced dissociation 323 Combined DNA Index System (CODIS) 225 common ancestor 75, 145, 156 last universal (LUCA) 156 comparative genomics 107, 377 complementary base pairing 20 complexity 82, 368 computational 373, F and NP 374 dynamic 371 static 371 computer science 43-44, 233, 373 concanavalin A 316 conformation-sensitive gel electrophoresis 220, 258 conformational angles 310 conformational change 360 Commonwealth Scientific Industrial and Research Organization (CSIRO) 87, 285 conjugation, bacterial 92 constitutive mutant 396 contig 50, 214 contig map 38 core of protein family 326 cost of DNA sequencing 23, 208 Cox, T.M. 229 CpG islands 54 creatine kinase 356 Crick, F.H.C. 9, 14, 19, 21, 228, 317 Critical Assessment of Structure Prediction (CASP) 337-338 crossing-over 29 Cryo-electron microscopy 382 CSIRO (Commonwealth Scientific Industrial and Research Organization) 87, 285	cyanidin 285 cyanobacteria, photosynthesis 173 cystic fibrosis 38, 124 cytochrome c 114, 117  <b>D</b> DALI 253 Darwin, C.R. 3, 15-17, 73, 138, 142, 147, 346 databanks 247 Dayhoff, M.O. 45, 240, 251 Delbrück, M. 18 deletion 13, 32-33, 95 deletion loop 34 delphinidin 285 depression 111 deuterostome 76 developmental expression pattern 5 changes in <i>Drosophila melanogaster</i> 281 dexamethasone 272 diagnosis 269 diarray 274-275, 402 didoxynucleoside triphosphate 204-205, 207 dihydrokaempferol 285 diphosphoglycerate (DPG) 334 directed evolution 346-347 diseases, protein aggregation associated 350 disulphide bridges 309 dihydroretinol 318 DNA damage 402 DNA fingerprinting 37, 260, 221ff DNA packaging 389 DNA polymerase 204 DNA sequencing 202-203 automated 207 cost of 208 mass spectrometry 209-210 Maxam-Gilbert method 207 pyrosequencing 209 Sanger method 207 DNA structure 18-20 docking problem 352 dog genome 182 domestication 182 breeds 183 domain 101, 312 domain recombination networks 385 dosage compensation 70 dot plot 235-238, 259 double helix 3
	Down's syndrome 47 <i>Drosophila melanogaster</i> , genome 116 duplication 5, 9, 95 large-scale 102 whole-genome 103 dyneins 330 dyslexia 110  <b>E</b> ecdysosozoa 76 EcoCyc 375-376 effector 334 electrospray ionization (ESI) 321 EMBL data library 10 EMBOSS 238 Emden-Meyerhof pathway 274 ENCODE (Encyclopedia of DNA Elements) 71, 121, 123-124 Endangered Species Act 70 endoreplication 105 endosymbiont 3, 48, 71, 92-93 Ensembl 6 enthalpy 163 entropy 163, 368 enzyme design 348-349 ephrin 237 <i>Escherichia coli</i> gene lengths 89 <i>Escherichia coli</i> genome 89 erythropoietin 14 estrogen 58 ethical, legal, and social issues 3, 24, 40, 69, 87-88, 221ff euchromatin 116 eukaryotes 74 phylogenetic relationships 175 evolutionary significant unit (ESU) 69 exon 9, 50, 82 expressed sequence tag (EST) 49-50, 253-254, 272 expression chip 267 expression patterns, developmental changes in 5, 281 expression patterns, evolutionary changes in 291 extraction 138, 140, 200  <b>F</b> fX-174 204 feedback 360 feedback inhibition 56 Ferry, G. 215

## Approches par graines et arbre de suffixe (e.g. SSAHA, SOAP, BLAT)

### ■ Exemple

- Génome GCACAGCACA
- Indexation des 3-mers

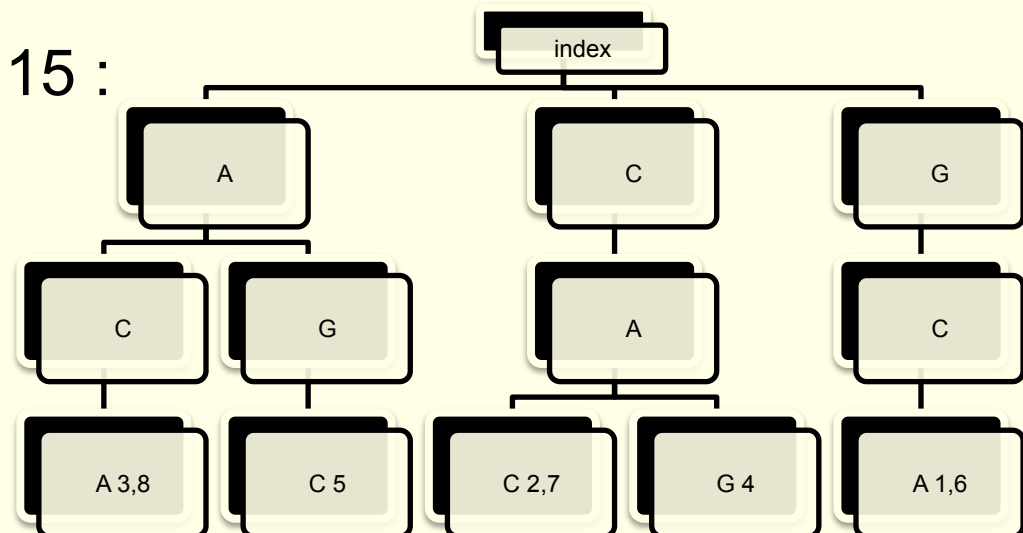
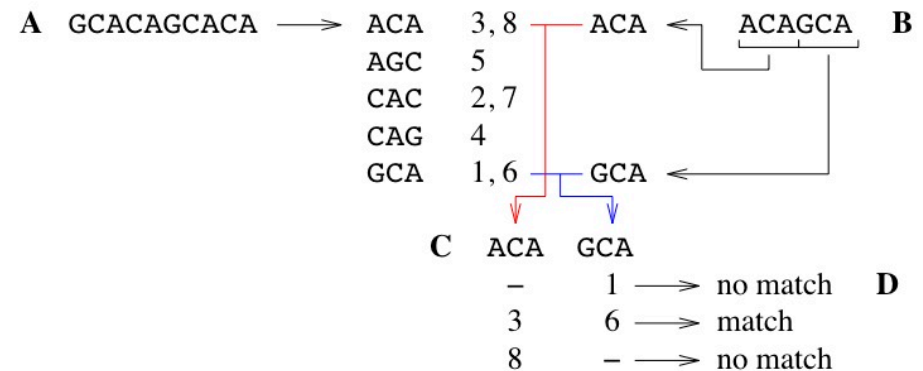
### ■ Recherche du read :

- Décomposition en 2 3-mers
- Recherche dans l'index de ces 3-mers
- Identification des zones avec les deux k-mers

### ■ Compression :

12 lettres au lieu de 15 :

⇒ 20% de gain

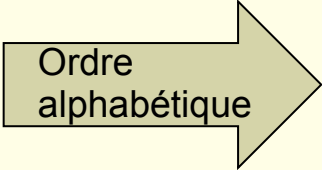


## Approches par table de suffixe

- Pour chaque read
  - Parcourir tout le génome et tester l'égalité read/portion génome
    - ⇒ Trop loooooong
- Chercher dans un livre de 500 pages où le mot «séquence » apparaît ☹
- Chercher la définition de « séquence » dans un dictionnaire 😊
- Pour cela on va construire une table des suffixes
  - Lister tous les suffixes du génome
  - Trier les suffixes par ordre alphabétique
  - Ex table des suffixes de **CATTATTAGGA** ?

## Approches par table de suffixe

### ■ Table des suffixes du génome CATTATTAGGA ?

1	CATTATTAGGA		1	11	A
2	ATTATTAGGA		2	8	AGGA
3	TTATTAGGA		3	6	ATTAGGA
4	TATTAGGA		4	2	ATTATTAGGA
5	ATTAGGA		5	1	CATTATTAGGA
6	TTAGGA		6	10	GA
7	TAGGA		7	9	GGA
8	AGGA		8	7	TAGGA
9	GGA		9	4	TATTAGGA
10	GA		10	6	TTAGGA
11	A		11	3	TTATTAGGA

- Trouver les positions des occurrences de TTA, de GGT ?
- En quoi le tri accélère t-il la recherche ?
  - ⇒ Est-ce un gain important ?
  - ⇒ Comment stocker la table en économisant l'espace ?

*Transformée de Burrow Wheeler*  
(e.g. *bwa*, *bowtie*)

- Reprend l'idée de la table de suffixe mais en beaucoup plus compact

## Format de sortie d'un mapping

### ■ Les fichiers « SAM » (Sequence Alignment/Map)

- ⇒ Fichier au format texte très lourd
- ⇒ Id du read, qualité du mapping, position du mapping, description « CIGAR » du mapping etc.

### ■ Format « CIGAR »

- ⇒ Description des événements associés à un mapping
- ⇒ M (match or mutation) I (insertion) D (Délétion)

```
ref ...CTTCATTACAG-TCTTTCG...  
read      ATC-CAGCT
```

**CIGAR : 3M1D3M1I1M** (on décrit ce qui se passe dans le read)

### ■ Les fichiers « BAM » versions compressées des SAM

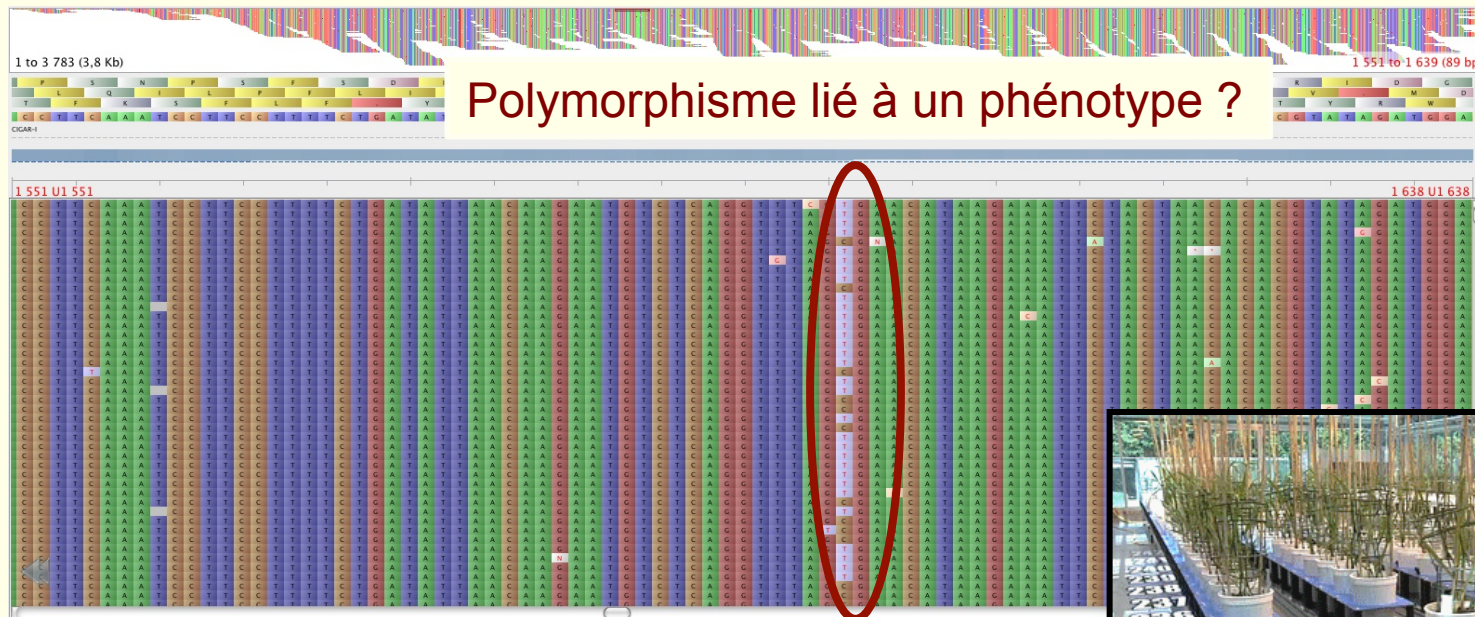
- ⇒ Plus rapide et moins lourd que les fichiers SAM
- ⇒ Fichier binaire non lisible par un humain

- Introduction NGS
- Assemblage de-novo
- Mapping sur un génome connu
- **Détection de SNP**
  - Approche par consensus
  - Approche probabiliste générale
  - Réglages et vérifications
- Conclusions et discussions



## ■ Problématique

- On dispose d'un génome/transcriptome de référence
- On à positionner les reads dessus
- On veut identifier les variations/polymorphismes



## Approche par consensus

(e.g. avec dnaSP)

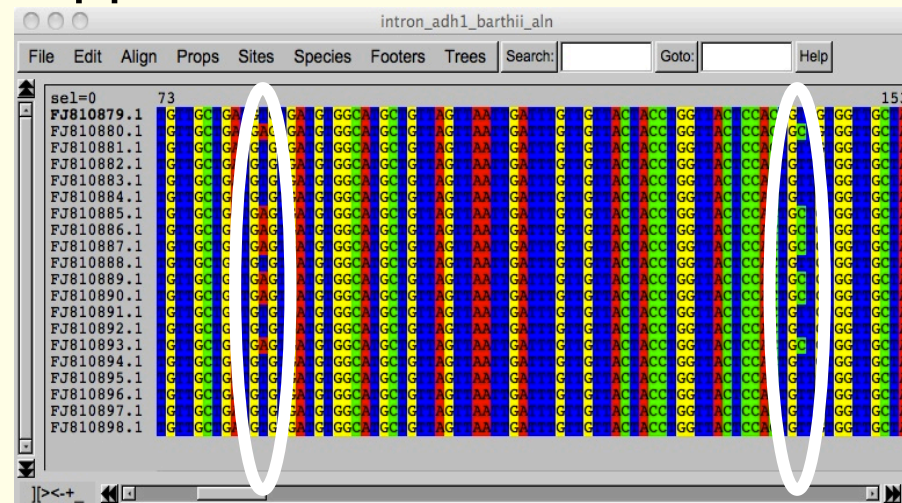
- Construire une séquence par individu
  - Assemblage par individu (pb de couverture ?)
  - Faire un consensus par individu après les avoir les mappés sur
    - Un génome de référence (l'idéal si possible)
    - Un assemblage fait en combinant les reads de tous les individus
  - Regarder pour chaque site si les individus diffèrent



# Approche par consensus

(e.g. avec dnaSP)

- Limite de l'approche ?



**Indiv 1**  
**TACAGTACGATC...**  
**TACAGT**  
**AGTACGATC**  
**ATCAC**  
**CAGCACG**  
**ACATCA**

**Indiv 2**  
**TACAGTACGATC...**  
**ACAGT**  
**AGTACGATC**  
**AGTACG**  
**ATCA**  
**TACAGCACG**  
**ACATCA**  
**TCAC**

<b>A</b>	<b>02050400100</b>
<b>C</b>	<b>00300030001</b>
<b>G</b>	<b>00003002000</b>
<b>T</b>	<b>10002000010</b>
<b>Indiv 1</b>	<b>: TACAGTACGATC...</b>
<b>Indiv 2</b>	<b>: TACAGTACGATC...</b>
<b>A</b>	<b>03060600100</b>
<b>C</b>	<b>00300040001</b>
<b>G</b>	<b>00004003000</b>
<b>T</b>	<b>10003000010</b>

**Approche probabiliste**  
(e.g. GATK, samtools, reads2SNP)

- $G_{a,b}$  : Génotype, ex  $G_{C,C}$  : homozygote C|C
- $D_{i,j}$  : Données observées, fréquences des nucléotides pour l'individu  $i$  au site  $j$

⇒  $D_{i,j} = \{f_A, f_C, f_G, f_T\}$

$$P(G_{ab} | D_{i,j}) = \frac{p(G_{ab}) p(D_{i,j} | G_{ab})}{p(D_{i,j})}$$

$$P(D_{i,j}) = \sum_a \sum_b p(G_{ab}) p(D_{i,j} | G_{ab}) \quad a,b = A,C,G \text{ ou } T$$

$$P(G_{ab} | D_{i,j}) = \frac{p(G_{ab}) p(D_{i,j} | G_{ab})}{\sum_a \sum_b p(G_{ab}) p(D_{i,j} | G_{ab})}$$

## Paramètres de la vraisemblance

### ■ Hardy-Weinberg

$$p(G_{ab}) = 2p_a p_b \text{ si } a \neq b$$

$$p(G_{ab}) = p_a^2 \text{ si } a = b$$

$p_a, p_b$  : fréquences au site  $j$  (tous individus confondus)

$$p(D_{i,j} | G_{ab}) = \begin{cases} P_{Multinomial} \{1 - 3\varepsilon, \varepsilon, \varepsilon, \varepsilon\} & \text{si } a = b = A \\ P_{Multinomial} \left\{ \frac{1}{2} - \varepsilon, \frac{1}{2} - \varepsilon, \varepsilon, \varepsilon \right\} & \text{si } a = A \text{ et } b = C \\ \text{etc...} & \end{cases}$$

**Approche probabiliste**  
(e.g. GATK, samtools, reads2SNP)

- $G_{a,b}$  : Génotype, ex  $G_{C,C}$  : homozygote C|C
- $D_{i,j}$  : Données observées, fréquences des nucléotides pour l'individu  $i$  au site  $j$   
 ⇒  $D_{i,j} = \{f_A, f_C, f_G, f_T\}$
- Probabilité postérieure et « SNP calling »

$$P(G_{ab} | D_{i,j}) = \frac{p(G_{ab}) p(D_{i,j} | G_{ab})}{\sum_a \sum_b p(G_{ab}) p(D_{i,j} | G_{ab})}$$

Attendu (stat, risque d'erreurs)  
Attendu (généet des pops)

## Réglages et vérifications

(e.g. GATK, samtools, reads2SNP)

- *Les fichiers vcf (variant call format)*
  - *Qualité du SNP, vraisemblance de chaque génotype etc.*
- *Fichiers bcf (binary call format)*
  - *Versions binaires compressées*
- *Qualité individuelle*
  - *Garder les SNP ayant une qualité > seuil (e.g. vcftools)*
  - *Garder les SNP isolés (3 SNP sur 10 sites consécutifs...???)*
  - *Etc.*
- *Qualité globale (comparaison aux attendus)*
  - *Ratio de SNP synonyme/non-synonyme ( $dN/dS \ll 1$ )*
  - *Ratio de SNP transition/transversion  $\sim 2$*
  - *Etc.*
- *Calibration : utilisation d'un jeu de SNP connus (GATK)*

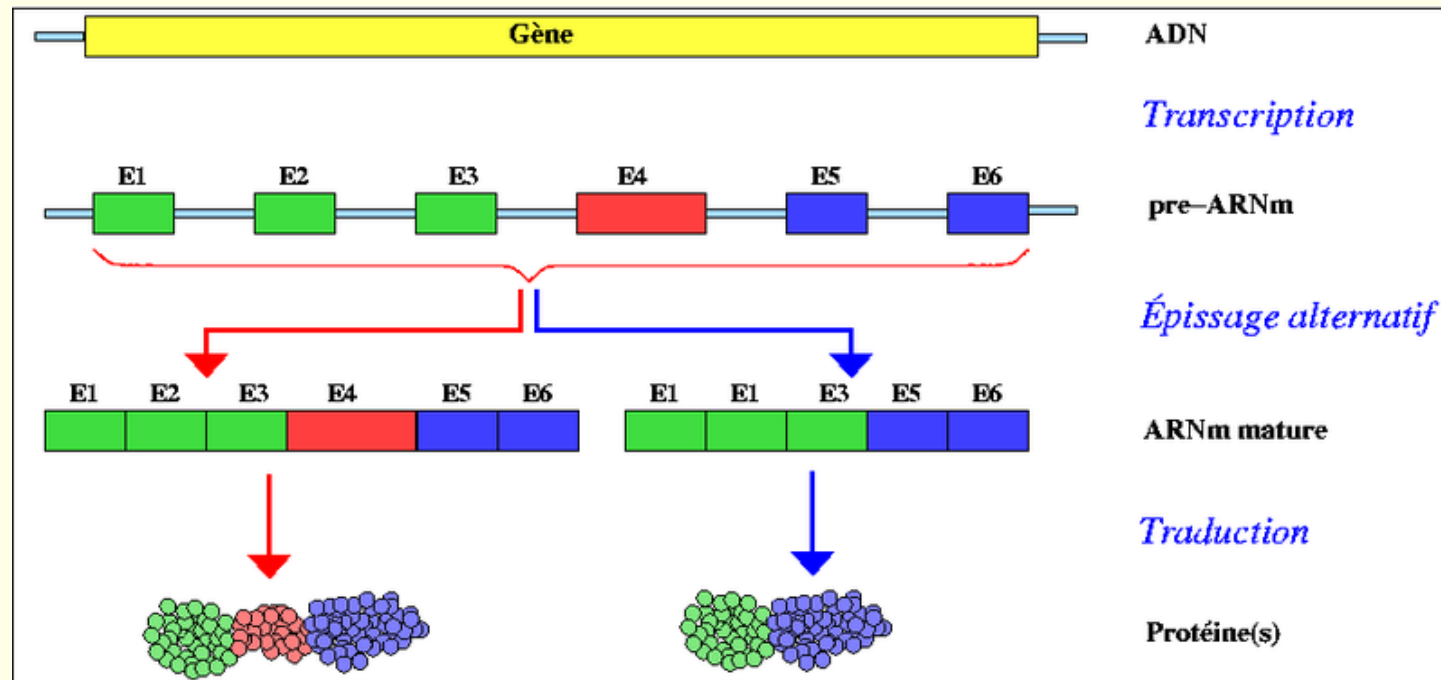
# *NGS : des reads aux SNPs*

- Introduction NGS
- Assemblage de-novo
- Mapping sur un génome connu
- Détection de SNP
- **Conclusions et discussions**

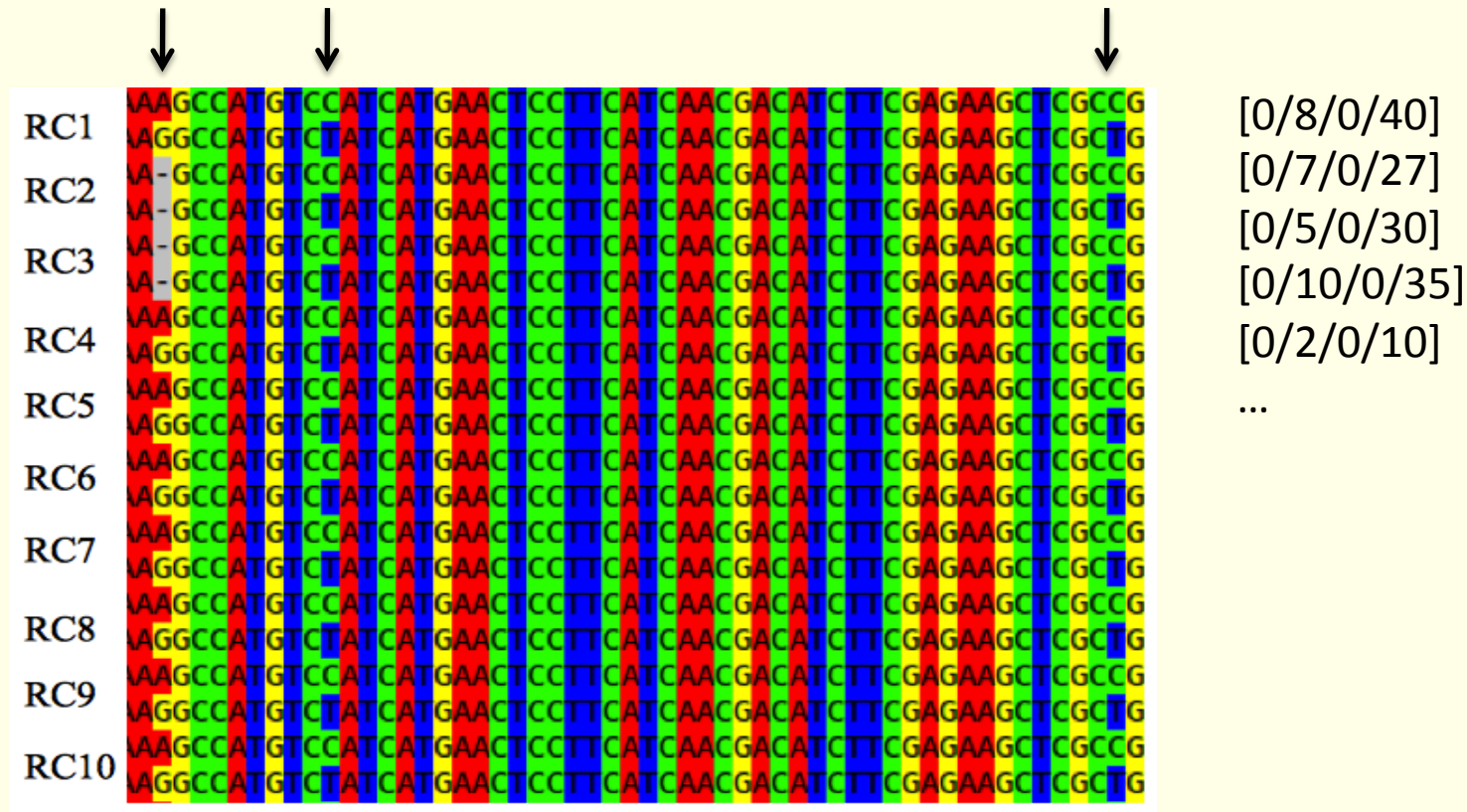


## Quelques problèmes ...

- Génomes de plantes
  - Beaucoup d'éléments répétés
  - ⇒ Complique l'assemblage et le mapping
- Transcriptomes plus simples ?



## Quelques problèmes ...



Exemple: Riz africain cultivé : très fortement autogame  
A quoi cela peut-il être dû ?

## Avant un projet NGS ?

- Séquencer pour quoi faire ?
  - Clarifier la question à laquelle on veut répondre
- Séquencer qui ?
  - Choix des individus, du type de tissu, de la date de prélèvement
  - Garder les individus vivants ou au moins de l'ADN (re-séquençage)
- Séquencer quoi ?
  - Génome ? Transcriptome ? Les deux ?
- Séquencer combien ?
  - Nombre d'individus => puissance statistique
  - Nombre de reads par individu => fiabilité, chevauchement
  - Choix de la technologie ...
- Evaluer les ressources nécessaires
  - CPU, stockage et ... temps humain (CDD, partenaires, sous-traiter)

Ne pas avoir peur de  
demander conseil,  
demander de l'aide...