

# Mises en perspective

Analyse bioinformatique de séquences  
pour l'amélioration des plantes

# Mardi 11 février 2014

- **Concepts**

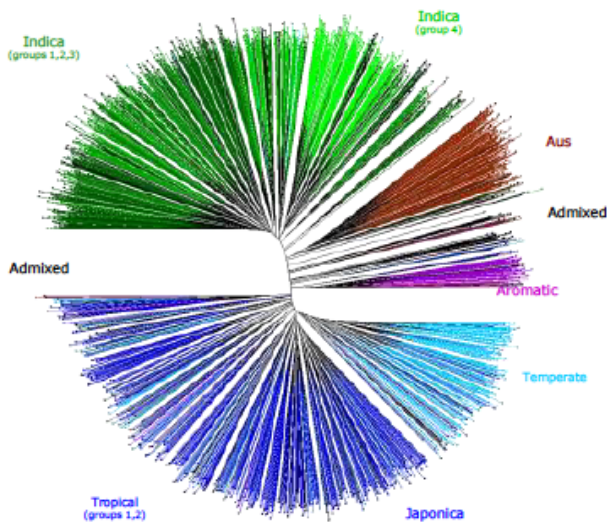
- Passage de données brutes à des données informatives
  - Alignement, élimination de duplicats techniques
  - Analyse de la profondeur de lecture
  - Identification des séquences présentes dans le jeu de données
  - Présence de polymorphismes (SNP, Indels), vrais ou faux...
  - Traitement des données: LD, design d'une puce SNP haut débit (Illumina Veracode)...
- Formats SAM (Sequence Alignment Map) / BAM (compressé), VCF (Variant Call Format)

- **Compétences**

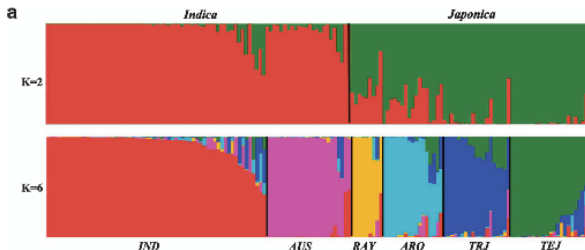
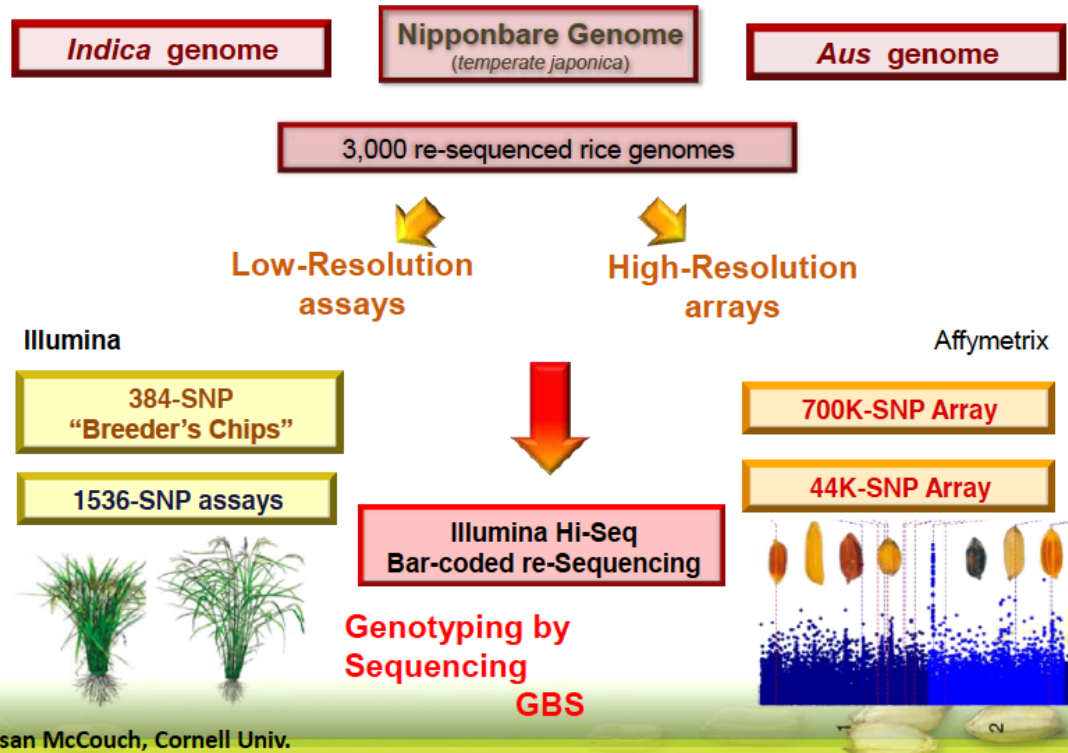
- Faire un assemblage de données NGS (MIRA)
- Positionner ces séquences sur une référence = mapping (BWA)
- Trier des fichiers
- Créer un workflow
- Détecter des SNP (GATK, SNIPlay)

# Mise en perspective: Polymorphisme et sélection

3000 diverse rice lines clustered by molecular markers

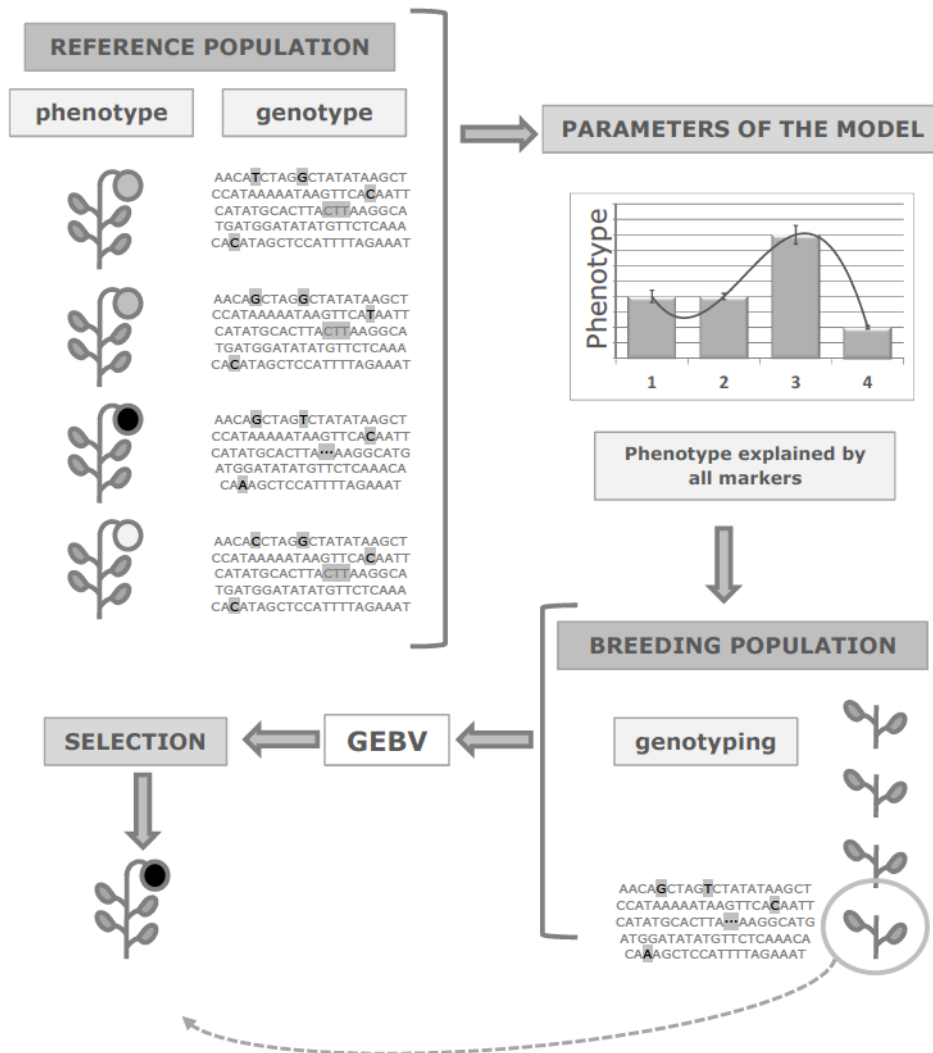


## SNP genotyping and analysis platforms for rice



Analyses de **diversité**: Stratégie de génotypage haute définition sur 3000 génomes de riz (Genotyping by sequencing)

# Mise en perspective: Polymorphisme et sélection



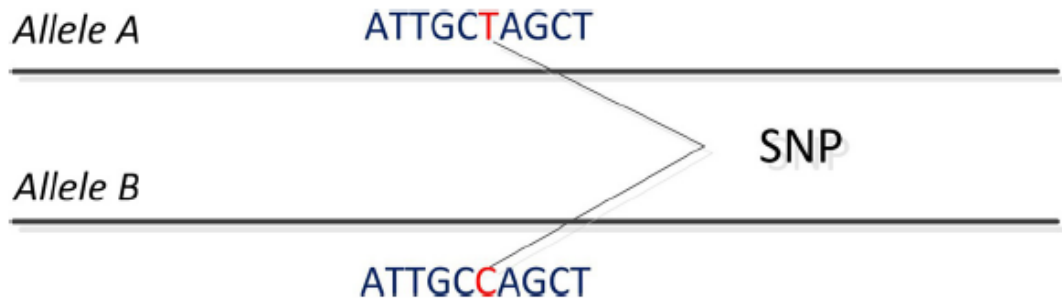
**Sélection génomique**  
(Perez de Castro et al 2012)

# Mise en perspective: Polymorphisme et sélection

Table 1 | A list of available non-commercial NGS genotype-calling software

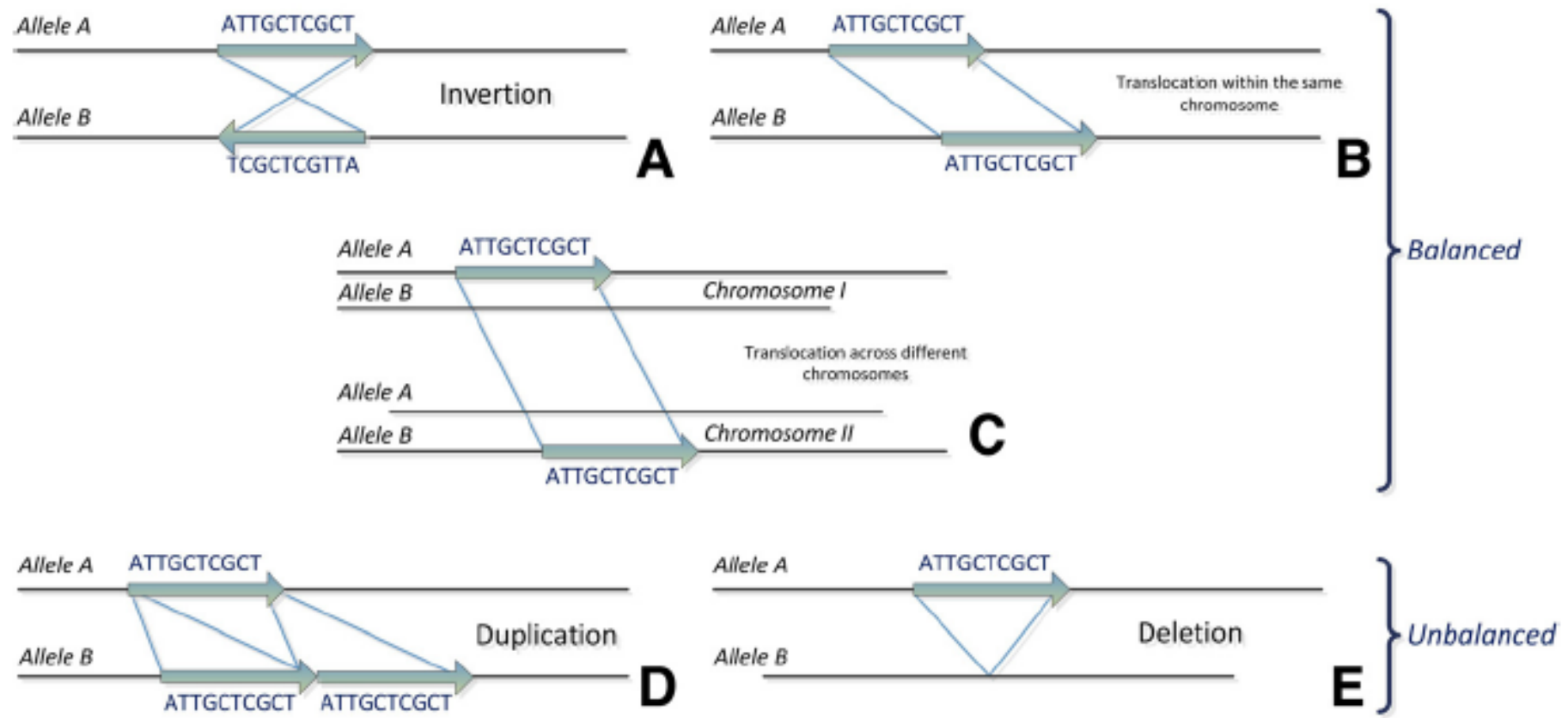
Software	Available from	Calling method	Prerequisites	Comments	Refs
SOAP2	<a href="http://soap.genomics.org.cn/index.html">http://soap.genomics.org.cn/index.html</a>	Single-sample	High-quality variant database (for example, dbSNP)	Package for NGS data analysis, which includes a single individual genotype caller (SOAPsnp)	15
realSFS	<a href="http://128.32.118.212/thorfinn/realSFS/">http://128.32.118.212/thorfinn/realSFS/</a>	Single-sample	Aligned reads	Software for SNP and genotype calling using single individuals and allele frequencies. Site frequency spectrum (SFS) estimation	-
Samtools	<a href="http://samtools.sourceforge.net/">http://samtools.sourceforge.net/</a>	Multi-sample	Aligned reads	Package for manipulation of NGS alignments, which includes a computation of genotype likelihoods (samtools) and SNP and genotype calling (bcftools)	53
GATK	<a href="http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit">http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit</a>	Multi-sample	Aligned reads	Package for aligned NGS data analysis, which includes a SNP and genotype caller (Unified Genotyper), SNP filtering (Variant Filtration) and SNP quality recalibration (Variant Recalibrator)	32,33
Beagle	<a href="http://faculty.washington.edu/browning/beagle/beagle.html">http://faculty.washington.edu/browning/beagle/beagle.html</a>	Multi-sample LD	Candidate SNPs, genotype likelihoods	Software for imputation, phasing and association that includes a mode for genotype calling	42
IMPUTE2	<a href="http://mathgen.stats.ox.ac.uk/impute/impute_v2.html">http://mathgen.stats.ox.ac.uk/impute/impute_v2.html</a>	Multi-sample LD	Candidate SNPs, genotype likelihoods	Software for imputation and phasing, including a mode for genotype calling. Requires fine-scale linkage map	44
QCall	<a href="ftp://ftp.sanger.ac.uk/pub/rd/OCALL">ftp://ftp.sanger.ac.uk/pub/rd/OCALL</a>	Multi-sample LD	'Feasible' genealogies at a dense set of loci, genotype likelihoods	Software for SNP and genotype calling, including a method for generating candidate SNPs without LD information (NLDA) and a method for incorporating LD information (LDA). The 'feasible' genealogies can be generated using Margarita ( <a href="http://www.sanger.ac.uk/resources/software/margarita">http://www.sanger.ac.uk/resources/software/margarita</a> )	54
MaCH	<a href="http://genome.sph.umich.edu/wiki/Thunder">http://genome.sph.umich.edu/wiki/Thunder</a>	Multi-sample LD	Genotype likelihoods	Software for SNP and genotype calling, including a method (GPT_Freq) for generating candidate SNPs without LD information and a method (thunder_glf_freq) for incorporating LD information	-

A more complete list is available from <http://seqanswers.com/wiki/Software/list>. LD, linkage disequilibrium; NGS, next-generation sequencing (Nielsen et al 2011)



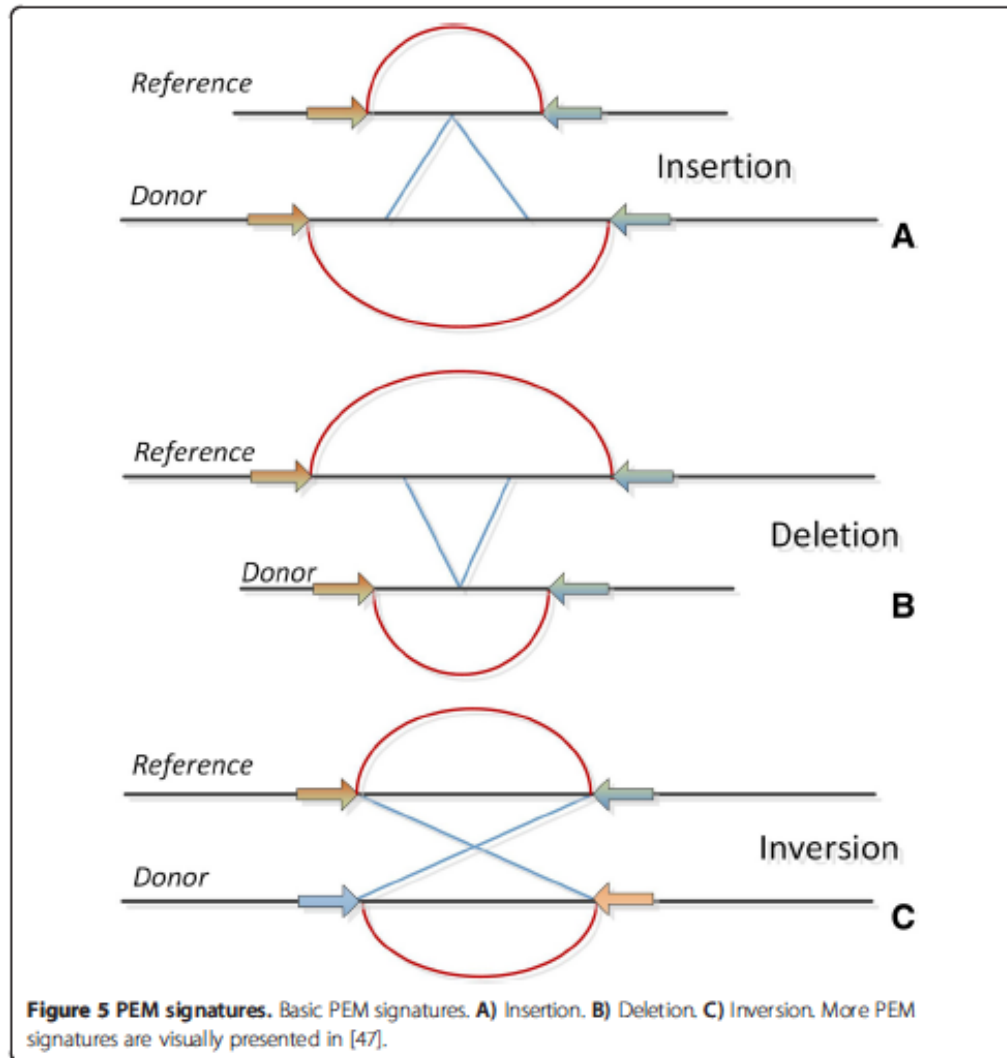
(Pavlopoulos et al 2013)

**Figure 3 SNP example.** A difference in a single nucleotide between two DNA fragments from different individuals. In this case we say that there are two alleles: C and T.



**Figure 4 Structural Variations.** This figure illustrates the basic structural variations. **A)** Inversion. **B)** Translocation within the same chromosome. **C)** Translocation across different chromosomes. **D)** Duplication. **E)** Deletion.

# Mise en perspective: Polymorphisme et sélection



(Pavlopoulos et al 2013)

Alignment	VCF representation	Variation
TGCACG TGTACG	POS 3 REF C ALT T	} SNP
TG_ACG TGTACG	POS 3 REF C ALT CT	
TGCACG TG__CG	POS 2 REF GCA ALT G	} Deletion
TGCGTG TG_TTG	POS 2 REF GCG ALT GT	

### Variation

POS	REF	ALT	INFO	
100	T	<DEL>	SVTYPE=DEL;END=300	} Large structural variants

**A**

```
##fileformat=VCFv4.0
##FORMAT=<ID=GQ,Number=1,Type=Float,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP Membership">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
```

} Header

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SAMPLE1	SAMPLE2	
1	3	rs2	ACG	A,AT	.	46.38	AN=2;DP=3;	GT:DP	1/2:8	0/0:10	} Body
1	2	.	C	T,CT	.	67.23	.	GT:GQ	0 1:60	2/2:30	
1	5	rs5	A	G	.	56.38	AC=2;AF=1	GT:GQ	1 0:63	1/1:85	
1	78	rs8	T	<DEL>	.	43.78	.	:DP	1/1:12	0/0:20	

Deletion      SNP      SV      Insertion

**B**

**Figure 9 VCF file.** This figure demonstrates an example of a CVF file. **A)** Different types of variations and polymorphisms that can be stored in CVF format. **B)** Example of a CVF format and its fields.



# Applications des polymorph

# Bibliographie

- G.A. Pavlopoulos et al. **Unravelling genomic variation from next generation sequencing data.** *Biodata Mining* (2013) 6:13
- A.M. Pérez-de-Castro et al. **Application of genomic tools in Plant breeding.** *Current genomics* (2012) 13, 179-195
- M.A. DePristo et al. **A framework for variation discovery and genotyping using next-generation DNA sequencing data.** *Nature Genetics* (2011) 43. 491-498
- R. Nielsen et al. **Genotype and SNP calling from next generation sequencing data.** *Nature Genetics* (2011) 12. 443-451