# Analyse de variants génétiques (SNPs, indels)

# I- SNP calling and genotype assignation
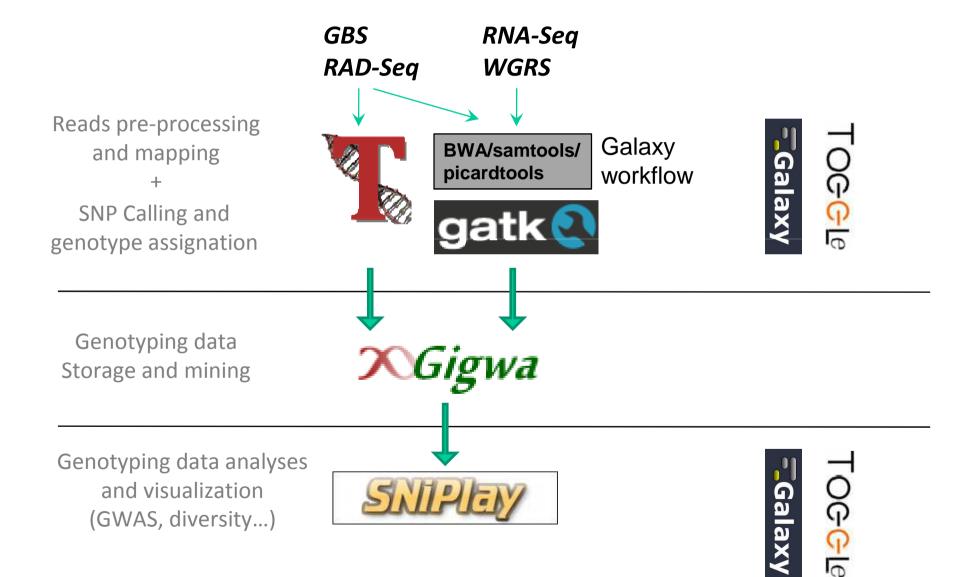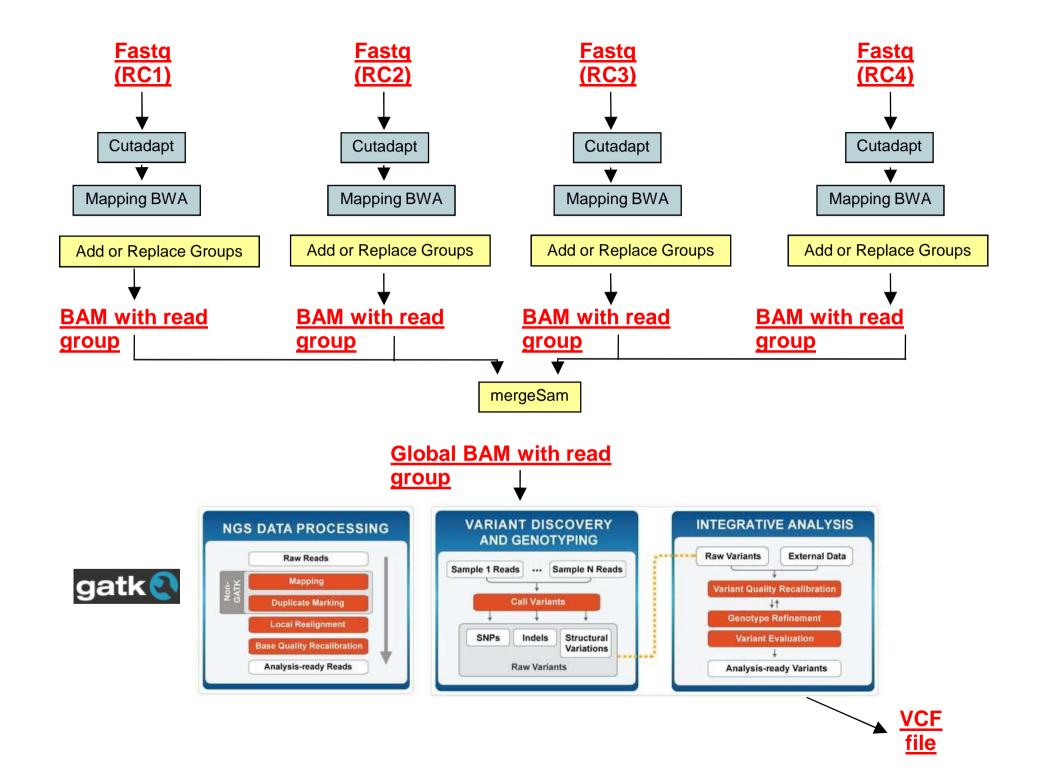
# II- SNP data analyses and visualization

# I- SNP calling and genotype assignation

# Format SAM/BAM

A header listing at least the different sequences in the reference and their length.
At least one line for each read. The different fields are separated by tabulation.

```
ILLUMINA-199BA2:90:FC:3:73:6197:17241 16    chr11      66288622 37   72M *    0    0
TGGGTGTAAGTGAGGAGATTGGGACTTTAGACCAGGCACTCACTCAGAGGAGTTACTGACAGGGCA
GGGAGG ;BHHHIFHFIHIHIIHIHIIIIGDIIIIGIEIIHIIIIHIIIFIHIIIIIIHIIIHIIHIIIIIIHIFIIIIII    XT:A:U    NM:i:0
     X0:i:1      X1:i:0     XM:i:0     XO:i:0     XG:i:0     MD:Z:72
```

Read name
Bitwise flag
Reference sequence
Start position of the alignment on the reference
Mapping quality (255 means unknown quality)
CIGAR String
Reference sequence and alignment start of the mate
Observed fragment length
Nucleotidic and quality sequence of the read
And optional fields

# Format Pileup

- Another format for variant calling (generated by samtools)
- Describe alignment row by row (not line by line like in SAM format)
- Used by softwares such as **Varscan** (varscan pileup2snp)
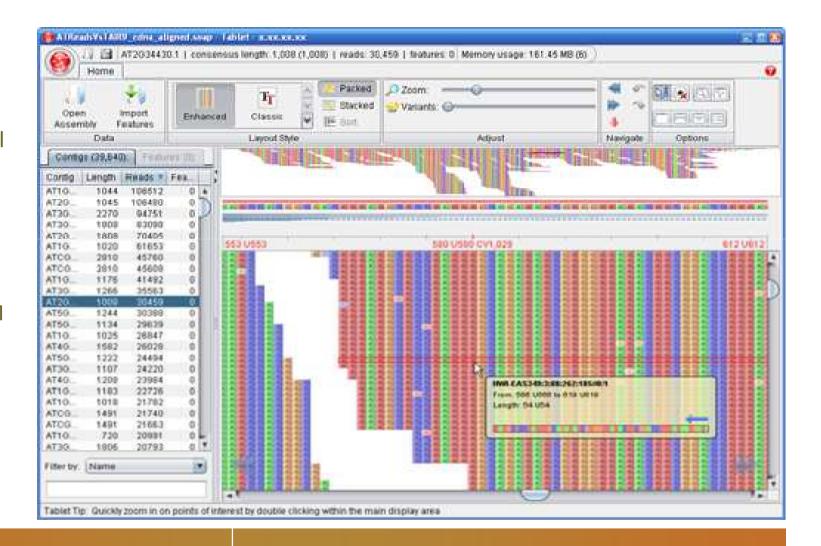- Frequently used for rare variants, with a low frequency (e.g. viral pop)

```
seq1 272 T 24  ,.$.....,,,,,....,,,,,..^+. <<<+;<<<<<<<<<<<<=<;<;7<&
seq1 273 T 23  ,.....,,,.,,,....,,,,..A <<<;<<<<<<<<<3<=<<<;<<+
seq1 274 T 23  ,.$....,,,.,,,....,,,,.    7<7;<;<<<<<<<<<=<;<;<<6
seq1 275 A 23  ,$....,,,.,,,....,,,,..^l. <+;9*<<<<<<<<<<=<<:;<<<<
seq1 276 G 22  ...T,,,.,,,....,,,,.... 33;+<<7=7<<7<&<<1;<<6<
seq1 277 T 22  ....,,,.,,,.C.,,,,,...G. +7<;<<<<<<<&<=<<:;<<&<
seq1 278 G 23  ....,,,.,,,....,,,,....^k. %38*<<;<7<<7<=<<<;<<<<<
seq1 279 C 23  A..T,,,.,,,....,,,,..... ;75&<<<<<<<<<<=<<<9<<:<<
```

Alexis Dereeper – Christine Tranchant

# Tablet

• Graphical tool to visualize assemblies

• Accept many formats
ACE, SAM, BAM

# GATK (Genome Analysis ToolKit)

- Software package to analyse NGS data.

- Implemented to analyse human resequencing data, for medical purpose (1000 genomes, The Cancer Genome Atlas)

- Includes depth analyses, quality score recalibration, SNP/InDel detection

- Complementary with other packages: SamTools, PicardTools, VCFtools, BEDtools

PREPROCESS:

  * Index human genome (Picard), we used HG18 from UCSC.
  * Convert Illumina reads to Fastq format
  * Convert Illumina 1.6 read quality scores to standard Sanger scores

FOR EACH SAMPLE:

  1. Align samples to genome (BWA), generates SAI files.
  2. Convert SAI to SAM (BWA)
  3. Convert SAM to BAM binary format (SAM Tools)
  4. Sort BAM (SAM Tools)
  5. Index BAM (SAM Tools)
  6. Identify target regions for realignment (Genome Analysis Toolkit)
  7. Realign BAM to get better Indel calling (Genome Analysis Toolkit)
  8. Reindex the realigned BAM (SAM Tools)
  9. Call Indels (Genome Analysis Toolkit)
  10. Call SNPs (Genome Analysis Toolkit)
  11. View aligned reads in BAM/BAI (Integrated Genome Viewer)

# GATK (Genome Analysis ToolKit)

- IndelRealigner module: realigns around indels in order to avoid false positive SNPs



HiSeq data, raw BWA alignments

HiSeq data, after MSA

# Format VCF (Variant Call Format)

→ Advantages:

Variation description for each position + genotype assignations

Indexed flat files.

Binary files also exist: BCF format



```
##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS     ID       REF  ALT      QUAL FILTER INFO                          FORMAT       Sample1         Sample2        Sample3
2      4370    rs6057   G    A        29   .      NS=2;DP=13;AF=0.5;DB;H2       GT:GQ:DP:HQ  0|0:48:1:52,51  1|0:48:8:51,51 1/1:43:5:.,.
2      7330    .        T    A        3    q10    NS=5;DP=12;AF=0.017           GT:GQ:DP:HQ  0|0:46:3:58,50  0|1:3:5:65,3   0/0:41:3
2      110696  rs6055   A    G,T      67   PASS   NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
2      130237  .        T    .        47   .      NS=2;DP=16;AA=T               GT:GQ:DP:HQ  0|0:54:7:56,60  0|0:48:4:56,51 0/0:61:2
2      134567  microsat1 GTCT G,GTACT 50   PASS   NS=2;DP=9;AA=G                GT:GQ:DP     0/1:35:4        0/2:17:2       1/1:40:3
```

# Other GATK functionalities

• Module DepthOfCoverage:

Allows to get sequencing depth for each gene, each position and each individual

• Module ReadBackedPhasing:

Allows to set, if possible, associations between alleles (phase and haplotypes) when we are in an heterozygote situation.

```
#CHROM  POS   ID    REF   ALT   QUAL   FILTER   INFO   FORMAT   SAMP1        SAMP2
chr1    1     .     A     G     99     PASS     .      GT:GL:GQ      0/1:-100,0,-100:99        0/1:-100,0,-100:99
chr1    2     .     A     G     99     PASS     .      GT:GL:GQ:PQ   1|1:-100,0,-100:99:60     0|1:-100,0,-100:99:50
chr1    3     .     A     G     99     PASS     .      GT:GL:GQ:PQ   0|1:-100,0,-100:99:60     0|0:-100,0,-100:99:60
chr1    4     .     A     G     99     FAIL     .      GT:GL:GQ      0/1:-100,0,-100:99        0/1:-100,0,-100:99
chr1    5     .     A     G     99     PASS     .      GT:GL:GQ:PQ   0|1:-100,0,-100:99:70     1|0:-100,0,-100:99:60
chr1    6     .     A     G     99     PASS     .      GT:GL:GQ:PQ   0/1:-100,0,-100:99        1|1:-100,0,-100:99:70
chr1    7     .     A     G     99     PASS     .      GT:GL:GQ:PQ   0|1:-100,0,-100:99:80     0|1:-100,0,-100:99:70
chr1    8     .     A     G     99     PASS     .      GT:GL:GQ:PQ   0|1:-100,0,-100:99:90     0|1:-100,0,-100:99:80
```

The proper interpretation of these records is that SAMP1 has the following haplotypes at positions 1-5 of chromosome 1:

1. AGAAA
2. GGGAG

And two haplotypes at positions 6-8:

1. AAA
2. GGG

Et non
AGG
GGA

# Haplotypes and phasing

- **Haplotype**: Specific groups of genes or alleles that progeny inherited from one parent

- **Phasing**: Determination of haplotype phase.
Process of statistical estimation of haplotypes from genotype data.

- Can be infered by statistics methods using non-ambigous haplotypes present in the dataset (Gevalt, ShapeIT, Phase)

- Can be resolved using physical association of alleles within the reads
(GATK ReadBackedPhasing, GATK HaplotypeCaller)

For GBS data

Tassel pipeline
Version 5

# II- SNP data analyses and visualization

# Projet Gigwa, pour la gestion des données massives de variants (GBS, RADSeq, WGRS)

« With NGS arise serious computational challenges in terms of storage, search, sharing, analysis, and data visualization, that redefine some practices in data management. »

- Based on NoSQL technology 

- Handles VCF files (Variant Call Format) and annotations

- Supports multiple variant types: SNPs, InDels, SSRs, SV

- Powerful genotyping queries

- Easily scalable with MongoDB sharding

- Transparent access

- Takes phasing information into account when importing/exporting in VCF format

http://gigwa.southgreen.fr/gigwa/

# SNP annotation using SnpEff



- It annotates and predicts the effects of variants on genes (amino acid changes…)
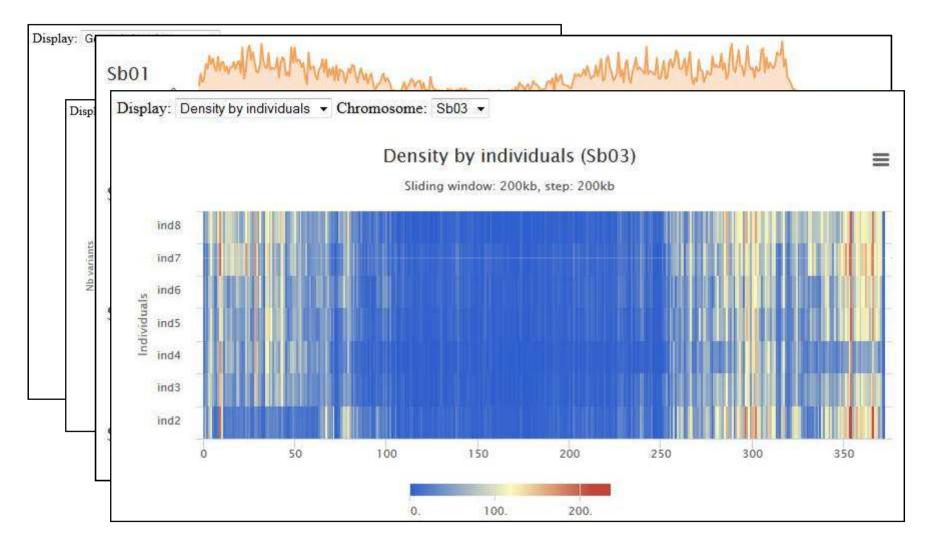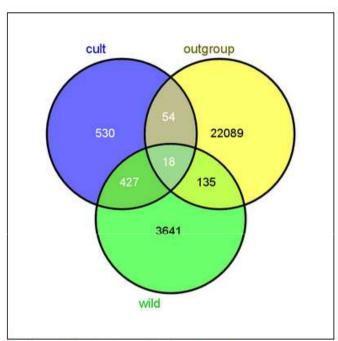- Uses as input GFF annoation file and VCF

Specific and shared polymorphisms between groups
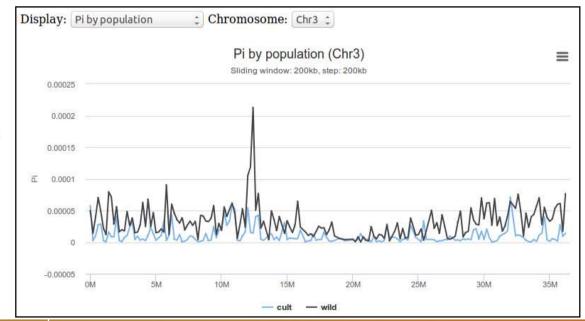
## Comparison between individuals

Fst: Fixation index: measure of population differentiation due to genetic structure.

Pi: Nucleotide diversity: Average number of nucleotide differences per site between any two DNA sequences chosen randomly from the sample population
Used to measure the degree of polymorphism within a population

+ 2186 polymorphisms inter-group

## Diversity analysis

SNP density by individuals can allow the detection of introgression event.

**Introgression** = Movement of a exogene region (gene flow) from one species into the gene pool of another by the repeated backcrossing of an interspecific hybrid with one of its parent species
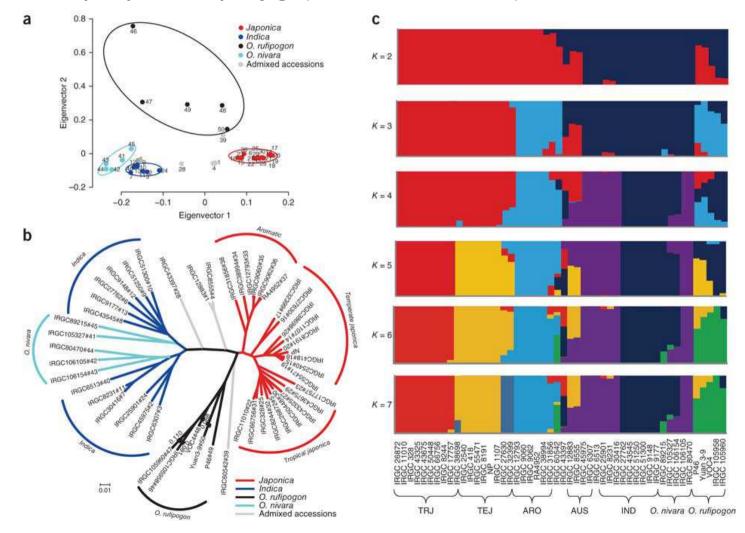
Widely used in agronomy obtained but can occurs naturally

# Population structure

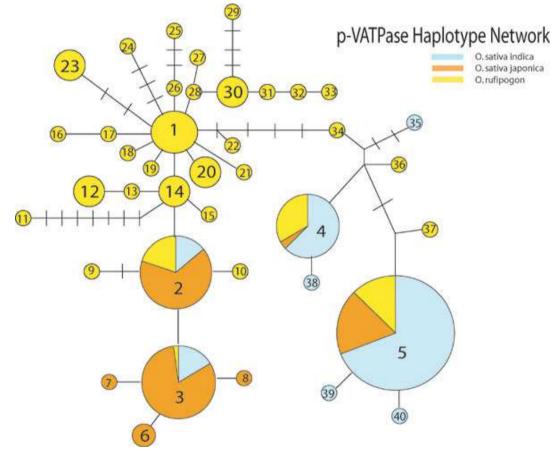**Ex: Riz asiatique après re-séquençage (Xun et al, Nature, 2011)**

# Haplotype network

**Exemple d'une région génomique chez le Riz**



p-VATPase Haplotype Network
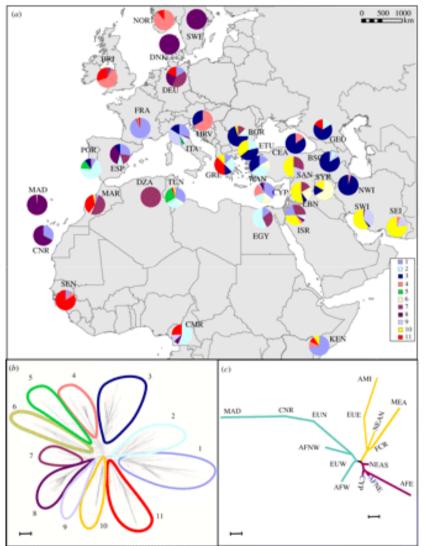
O. sativa indica
O. sativa japonica
O. rufipogon

# Haplotype and geographical distribution

Différenciation génétique de la souris domestique (Bonhomme et al, 2010)

## GWAS (Genome-Wide Association Studies)
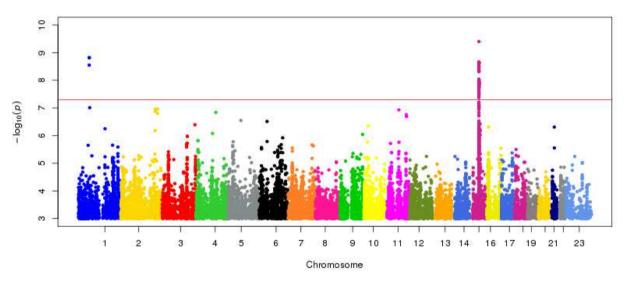
- Estimate association between a marker and a phenotypic character

- Manhattan plots: displays GWAS statistical tests (-log10 pvalue) along chromosomes

- TASSEL, MLMM sofwares

- False positives because of the studied structuration panel
=> correction using structure population et and kinship

# GWAS issues

• **Choice of genotypic panel**: phenotypic diversity for target traits must be sufficient (core-collection, MAGIC lines, NAM…)

• **Population structure** induces high rates of false associations (false positives)
• Correction using structure population et and kinship. Mixed models:
  o Q
  o K (widely used)
  o Q+K (widely used)

• **Density of markers** must be enough to provide a good genome cover. Density can be also highly variable.

• **Linkage disequilibrium (LD) landscape**: level of intra- and inter-chromosomal LD (number of loci in LD with loci from other chromosomes). Ideally, LD profile must be flat to avoid distorsion in association patterns.
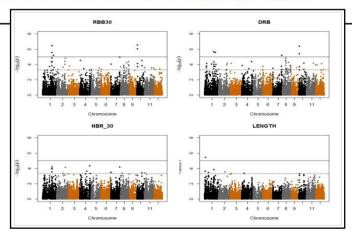
# Study of root characters using GWAS in Oryza sativa japonica. Influence of a correction using structure and kinship
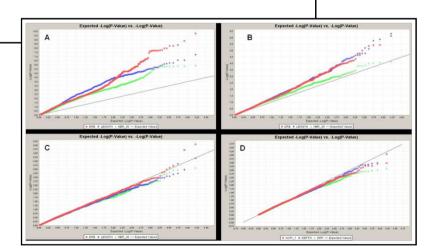
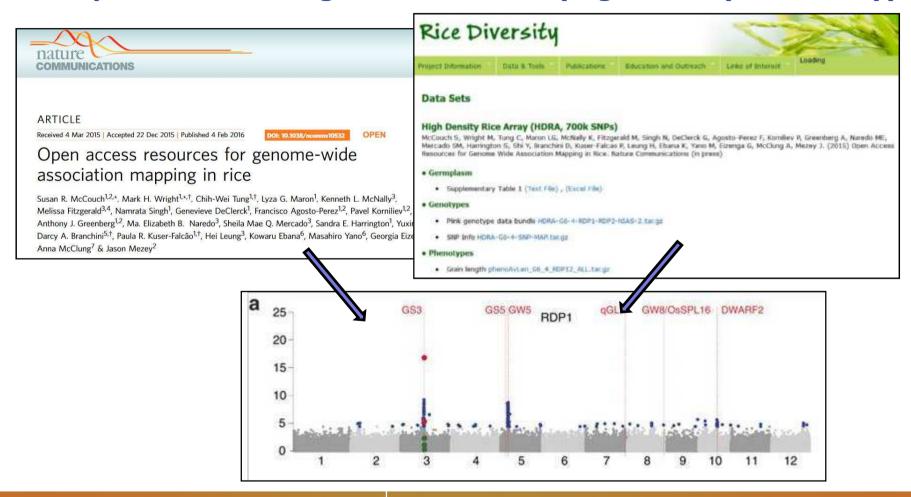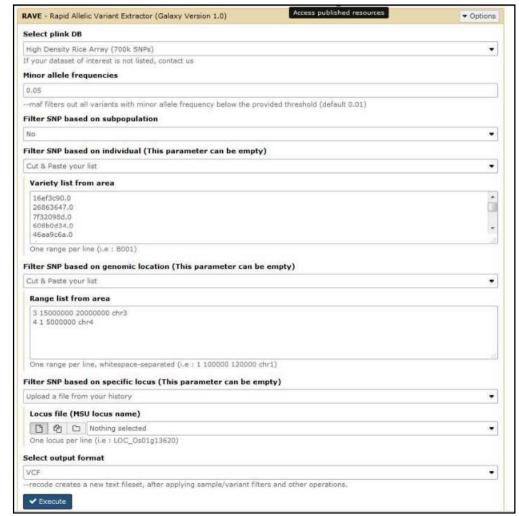# Exemple du Riz: 3000 genomes+ HDRA (High density Rice Array)

# Exemple du Riz: 3000 genomes+ HDRA (High density Rice Array)

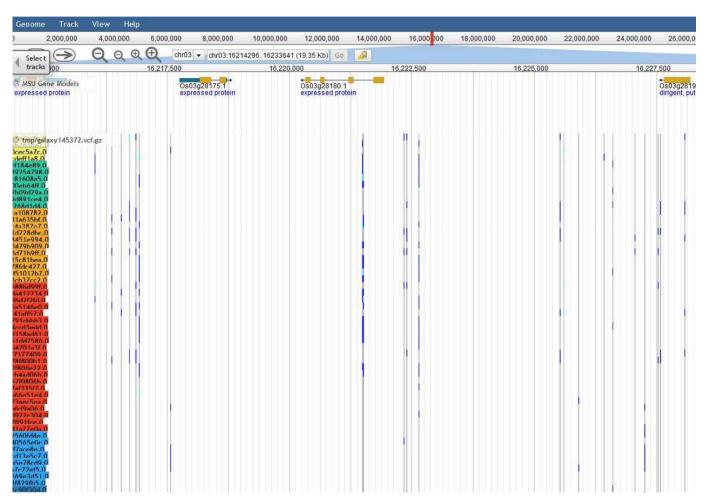Extraction rapide des variants après sélection d'une region / population donnée.

# Exemple du Riz: 3000 genomes+ HDRA (High density Rice Array)

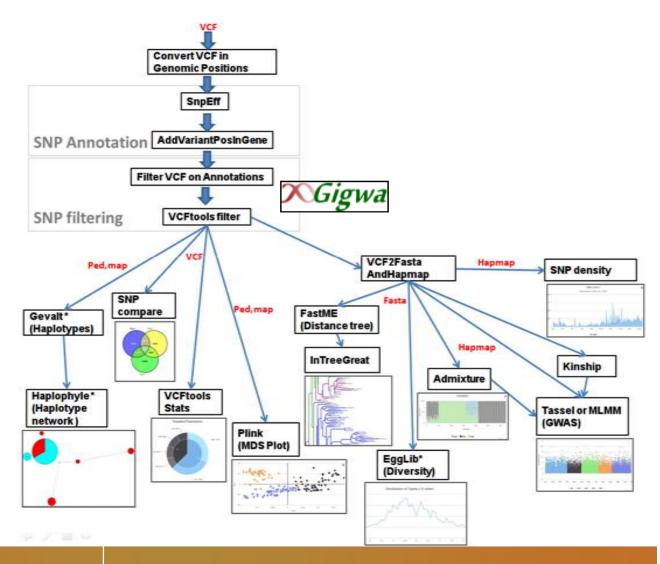Visualisation du contexte génomique dans un génome browser (plugin Jbrowse)

**SNiPlay: Web application for polymorphism analyses**

http://sniplay.southgreen.fr

# SNiPlay Site web



**http://sniplay.southgreen.fr**

## "Galaxy4Sniplay" : SNiPlay sous Galaxy