



Atelier Variant niveau 2

Cathy PHILLIPE - CEA

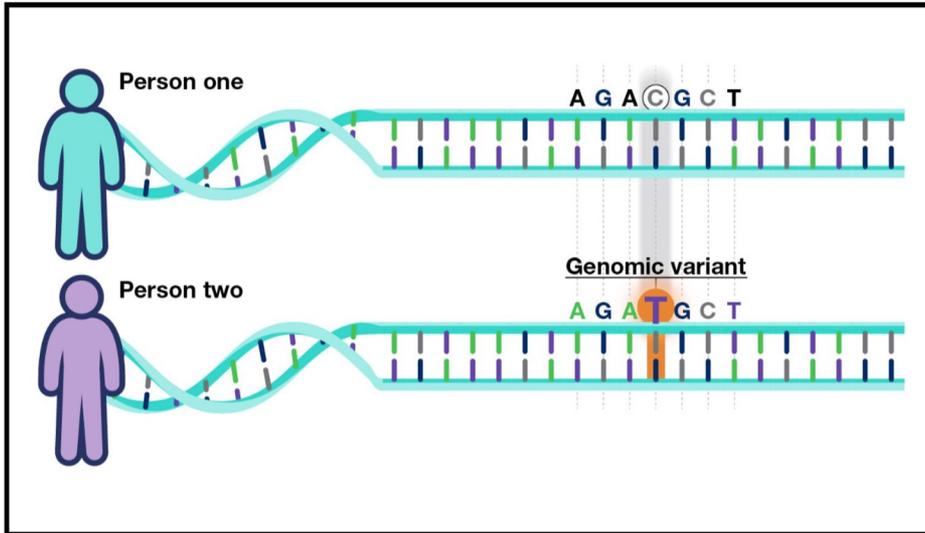
Nadia BESSOLTANE - INRAe

Pauline FRANCOIS - ANSES

Définition

1. Qu'est ce que c'est une variation génomique ?

Une variation génomique est un changement, d'une ou plusieurs bases nucléotides, dans une séquence d'ADN particulière en comparaison avec une séquence d'ADN (un génome) de référence (1). Les variations génomiques se distinguent en deux catégories : [polymorphismes](#) et [mutations](#).



Il existe différents types de variations :

- **SNV** : Single Nucleotide Variant
- **INDEL** : INsertion ou DELEtion
- **SV** (Structural Variant)
- **CNV** (Copy Number Variation)

du Fastq aux VCFs (ebain1)



Reads (Fastq)

Fastq Quality Control
-- *FastQC* --

Mapping
-- *Bwa* --

Reference genome (Fasta)

Processing Post Alignment
-- *GATK/Picard* --

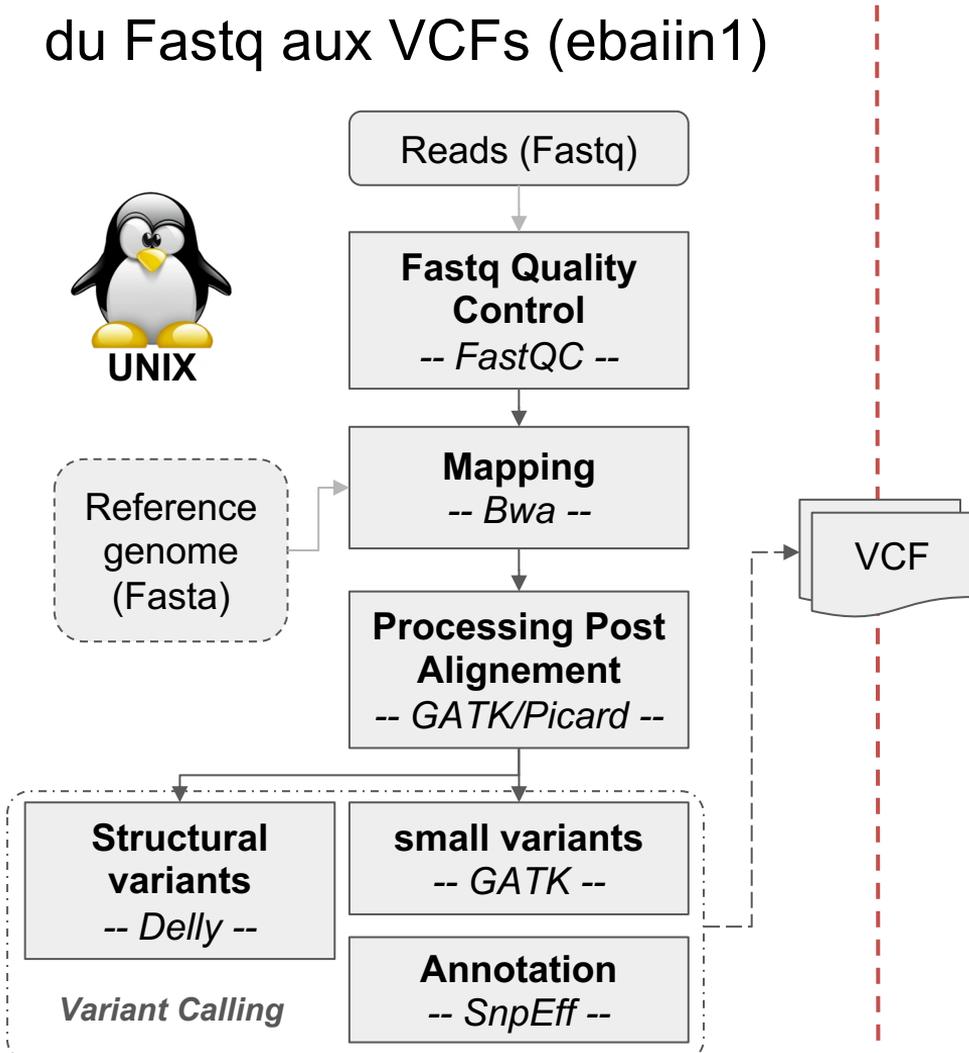
Structural variants
-- *Delly* --

small variants
-- *GATK* --

Annotation
-- *Snpeff* --

Variant Calling

VCF



VCF (variant call format)

VCF header

```
##fileformat=VCFv4.0
##fileDate=20100707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
```

Mandatory header lines (points to ##fileformat=VCFv4.0)

Optional header lines (meta-data about the annotations in the VCF body) (points to ##INFO=AA, ##INFO=H2, ##FORMAT=GL, ##INFO=SVTYPE, ##INFO=END)

Body

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SAMPLE1	SAMPLE2
1	1	.	ACG	A,AT	.	PASS	.	GT:DP	1/2:13	0/0:29
1	2	rs1	C	T,CT	.	PASS	H2;AA=T	GT:GQ	0 1:100	2/2:70
1	5	.	A	G	.	PASS	.	GT:GQ	1 0:77	1/1:95
1	100	.	T		.	PASS	SVTYPE=DEL;END=300	GT:GQ:DP	1/1:12:3	0/0:20

Reference alleles (GT=0) (points to 0/0:29)

Alternate alleles (GT>0 is an index to the ALT column) (points to 1|0:77)

Phased data (G and C above are on the same chromosome) (points to 0|1:100)

Deletion (points to in ALT)

SNP (points to A,AT in ALT)

Large SV (points to in ALT)

Insertion (points to T,CT in ALT)

Other event (points to H2;AA=T in INFO)

VCF → 3 parties principales

METADATA

```
##fileformat=VCFv4.2
##FILTER=<ID=LowQual,Description="Low quality">
##FORMAT=<ID=AD,Number=R,Type=Integer,Description="Allelic depths for the ref and alt alleles in the order listed">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Approximate read depth (reads with MQ=255 or with bad mates are filtered)">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=PL,Number=G,Type=Integer,Description="Normalized, Phred-scaled likelihoods for genotypes as defined in the VCF specification">
##GATKCommandLine=<ID=HaplotypeCaller,CommandLine="HaplotypeCaller --min-base-quality-score 18 --emit-ref-confidence NONE">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency, for each ALT allele, in the same order as listed">
##INFO=<ID=AN,Number=1,Type=Integer,Description="Total number of alleles in called genotypes">
##INFO=<ID=BaseQRankSum,Number=1,Type=Float,Description="Z-score from Wilcoxon rank sum test of Alt Vs. Ref base qualities">
...
##INFO=<ID=ReadPosRankSum,Number=1,Type=Float,Description="Z-score from Wilcoxon rank sum test of Alt vs. Ref read position">
##INFO=<ID=SOR,Number=1,Type=Float,Description="Symmetric Odds Ratio of 2x2 contingency table to detect strand bias">
##contig=<ID=6,length=119458736>
##source=HaplotypeCaller
```

VCF header

Body

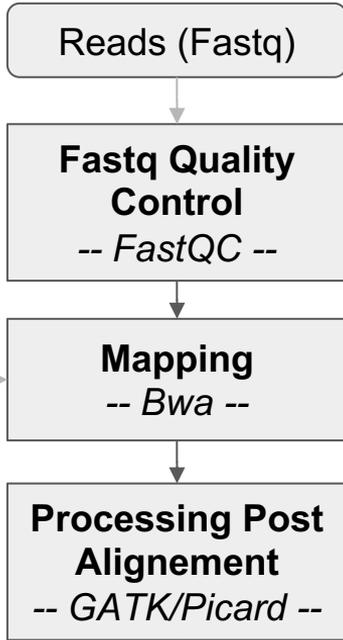
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
6	2	.	T	A	67.64	.	AC=1;AF=0.500;...
6	4	.	GT	G	58.60	.	AC=1;AF=0.500;...
6	9	.	C	CA	55.60	.	AC=1;AF=0.500;...

FORMAT	SRR1262731	SRR1262732
GT:AD:DP:GQ:PL	0/1:3,2:5:75:75,0,105	0/1:3,2:5:75:75,
GT:AD:DP:GQ:PL	0/1:1,2:3:28:66,0,28	0/1:1,2:3:28:66,
GT:AD:DP:GQ:PL	0/1:7,2:9:63:63,0,279	0/1:7,2:9:63:63,

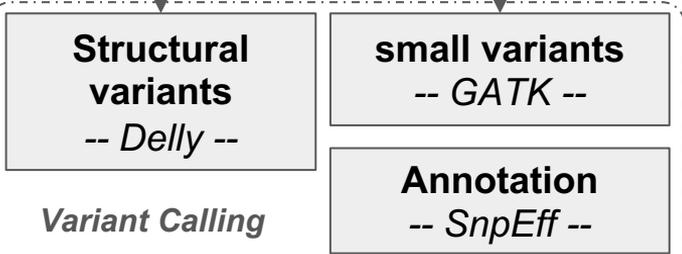
INFO

GENOTYPE

du Fastq aux VCFs (ebain1)



Reference genome (Fasta)



du VCFs aux marqueurs (ebain2)



Rappel : Filtre des variants

- De nombreux filtres peuvent être appliqués sur le VCF
 - type de variants à garder (SNVs seulement, Indels...)
 - région d'intérêt
 - filtres sur la qualité (seuils arbitraires : profondeur, génotype (0/1, 1/1), ratio allélique...)

- Filtres difficilement transposables entre analyse :
 - dépendent de la **question biologique**
 - dépendent des outils utilisés

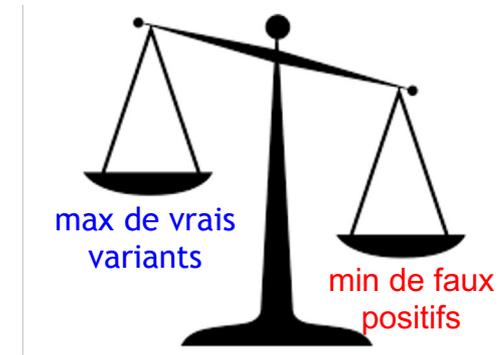
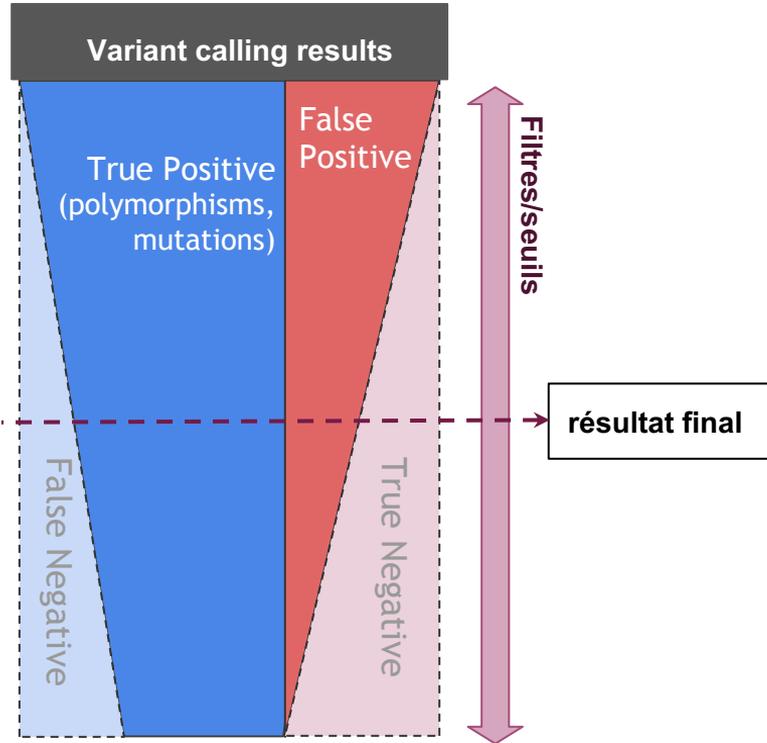
→ d'où l'intérêt de faire ses propres scripts

Rappel : Faux positif

- De nombreux variants **Faux Positifs** peuvent survenir des étapes précédentes :
 - Artéfacts issus des **cycle PCR** pendant la préparation des échantillons
 - Artéfacts liés à la **technologie de séquençage** (PacBio, HiSeq, NextSeq, ...)
 - Difficultés d'**alignement** (régions d'ADN répétées)
 - **Erreurs de lecture** lors du “BaseCalling”

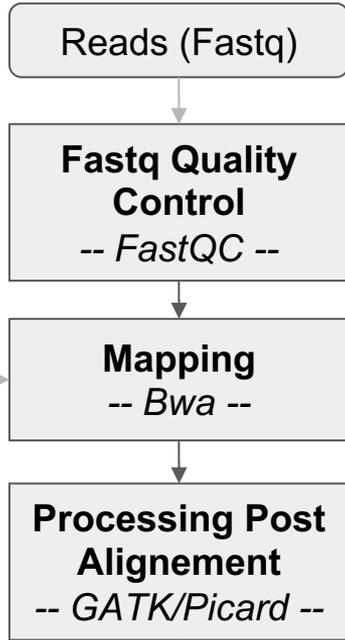
Rappel :

Faux positifs vs Filtres qualité

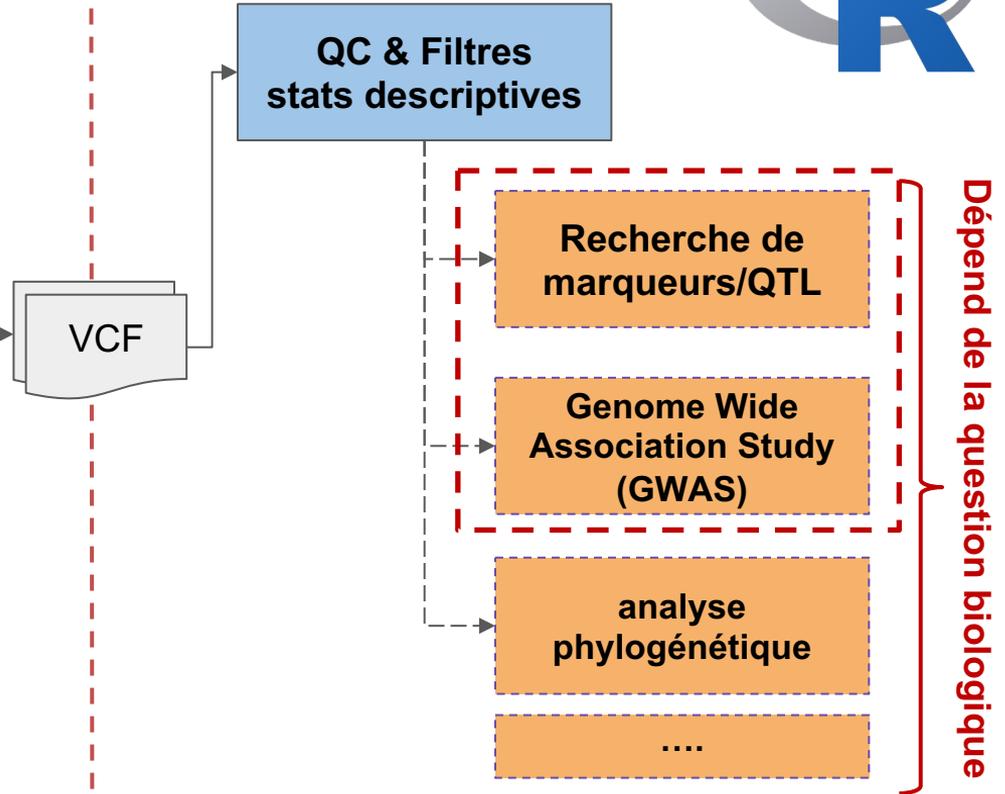


Plus on est stringent plus on va éliminer les **faux positifs** mais avec le risque de perdre de **vrais variants**

du Fastq aux VCFs (ebain1)



du VCFs aux marqueurs (ebain2)



Objectifs de l'atelier variants

- I. Analyse post-VCF des petits variants (SNV & InDel) → 2h (Nadia)
- II. Etudes d'association génétique à grande échelle (GWAS) → 4h (Cathy)

Objectifs de l'atelier variants

I. Analyse post-VCF des petits variants (SNV & InDel) → 2h (Nadia)

1- Se connecter sur Rstudio via OnDemand
(1 cpu; 8G ram)

2- Aller sur votre espace projet

```
> setwd("/shared/projects/olala_olala")
```

3- Copier le matériel de TP dans le dir TP_variants

```
> file.copy(from = "/shared/projects/2514_ebiii_n2/dnaseq/TP_small_variants",  
           to   = "./", recursive=TRUE)
```

4- Positionner l'espace du travail

```
> setwd("TP_small_variants/")
```

Objectifs de l'atelier variants

- I. Analyse post-VCF des petits variants (SNV & InDel) → 2h (Nadia)
- II. Etudes d'association génétique à grande échelle (GWAS) → 4h (Cathy)

1- Se connecter sur Rstudio via OnDemand
(4 cpu; 16G ram)

2- Aller sur votre espace projet

```
> setwd("/shared/projects/olala_olala")
```

3- Copier le matériel de TP dans le dir TP_variants

```
> file.copy(from = "/shared/projects/2514_ebiii_n2/dnaseq/GWAS",  
           to   = "./", recursive=TRUE)
```

4- Positionner l'espace du travail

```
> setwd("GWAS")
```