

Détection de variants

Thomas Faraut

Toulouse, 13 mars 2023

Introduction
oooooooooooo

Detection
oooooooooooo

Genotype likelihoods
oooooooooooo

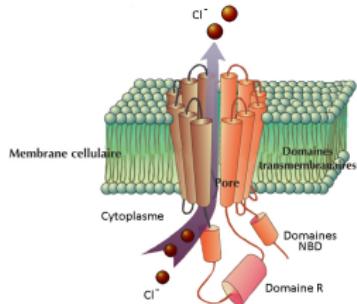
Complications
oooooooooooo

Outline

- What is a genomic variant ?
- Why should we care about them ?
- How can we detect variant ?

Genetic disorders

Cystic fibrosis is caused by mutations in the CFTR gene

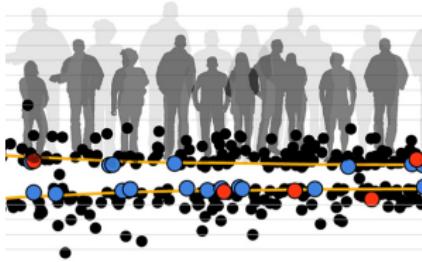
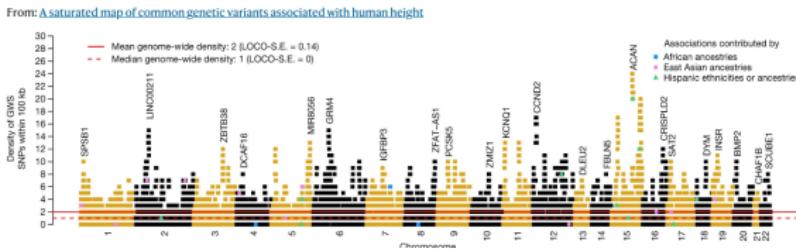


CFTR sauvage	Classe 1	Classe 2	Classe 3	Classe 4	Classe 5	Classe 6
Type de défauts	Ø protéine	Ø traffic	Ø fonction	Diminution de la fonction	Moins de protéines	Moins stable
Exemples de mutations	p.Gly542*	p.Gly556lu p.Ile57del	p.Val520Phe p.Ser549Arg	p.Arg117His p.Arg334Trp p.Ser1235Arg	p.Arg556Gu c.1689-188A-G c.2651-56A	Résiduel p.Phe508del p.Gln1412*
Approches requises	Restaurer la synthèse protéique	Corriger le transportement de la protéine	Restaurer la conformation du canal	Restaurer la conformation du canal	Maturat°/ Restaurer la conformat° de l'épissage	Favoriser la conformat° protéique

La mucoviscidose. Du gène à la thérapeutique Medecine et Sciences (2021)

Genomic variants and human height

10% to 40% of all variation in height can be explained by the genome

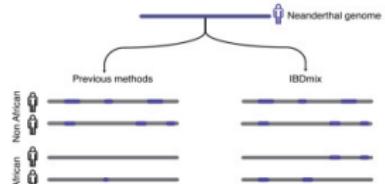


A saturated map of common genetic variants associated with human height. Nature (2022)

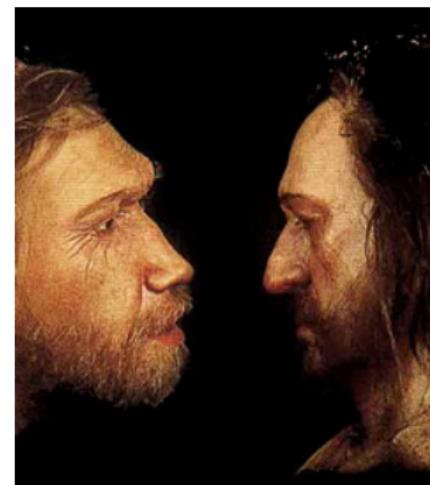
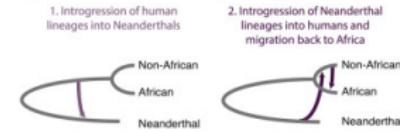
Neanderthal our cousin

"East Asian populations have approximately 20% more Neanderthal ancestry compared to Europeans"

Higher signal of Neanderthal ancestry in African individuals than previously thought



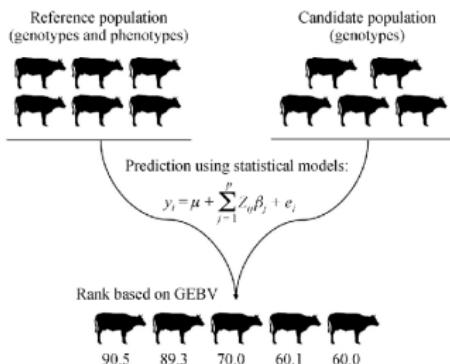
Signal of Neanderthal ancestry in Africa due to two events



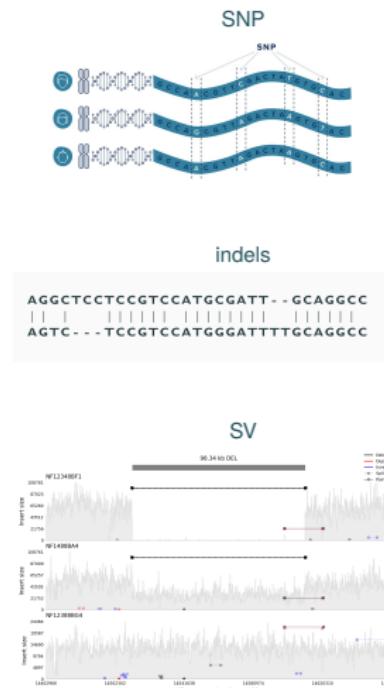
Identifying and Interpreting Apparent Neanderthal Ancestry in African Individuals. Cell (2020)

Genomic selection

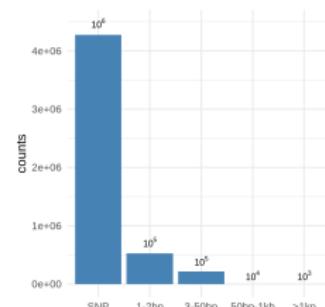
L'estimation d'une relation entre variations et phénotype permet de prédire le potentiel génétique d'un individu.



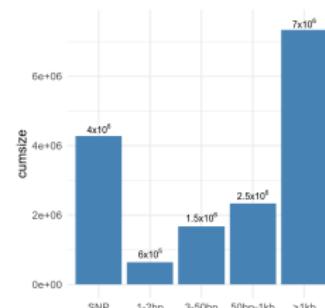
Variations



Number of variants

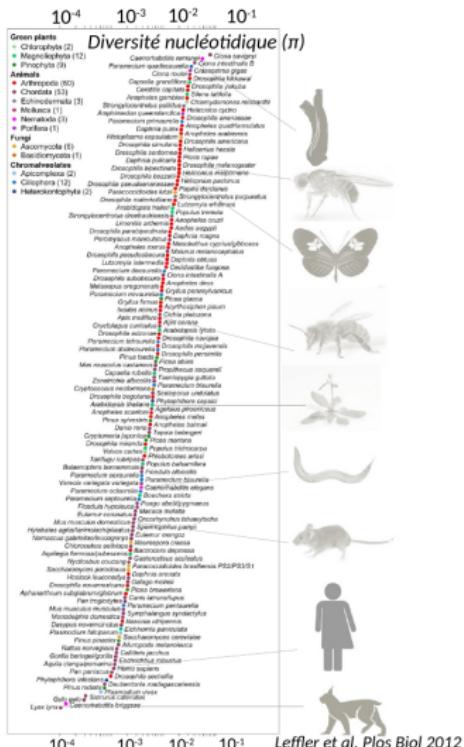


Cumulative size



Nucleotide diversity

The level of polymorphism, the nucleotide diversity, differs between species



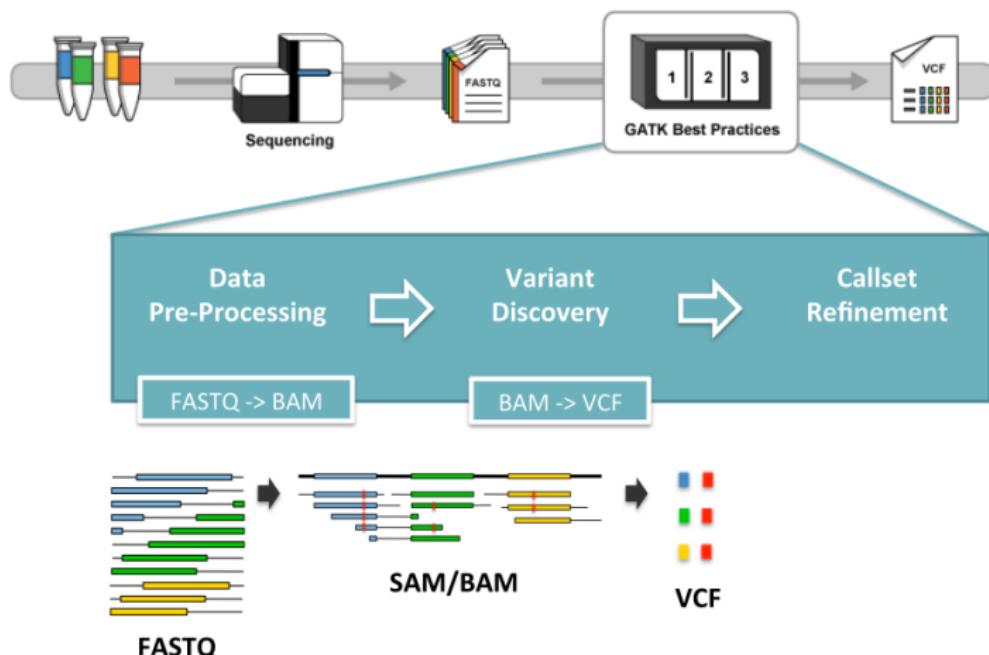
Introduction
oooooooo●ooooo

Detection
oooooooooooo

Genotype likelihoods
oooooooooooo

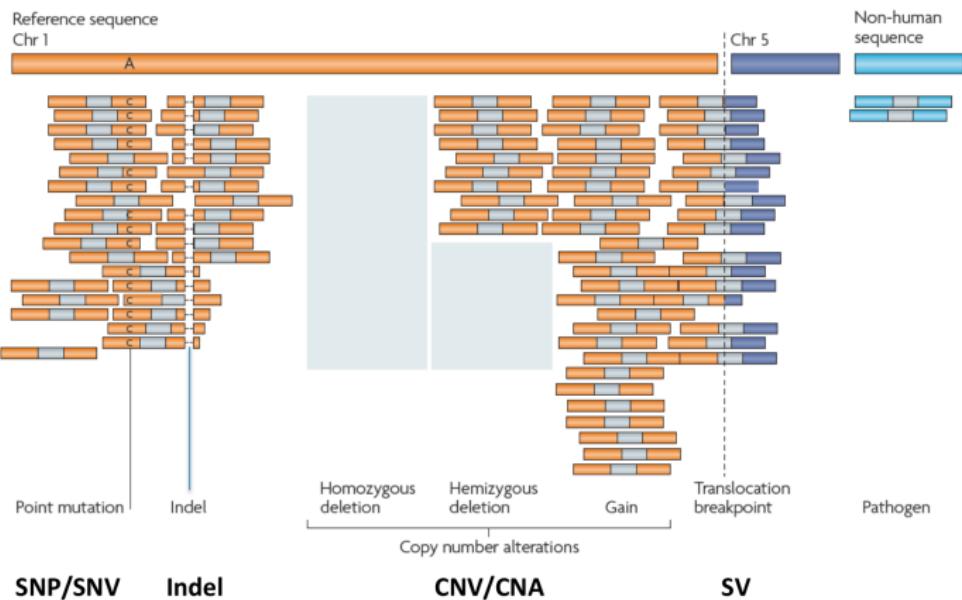
Complications
oooooooo

NGS workflow



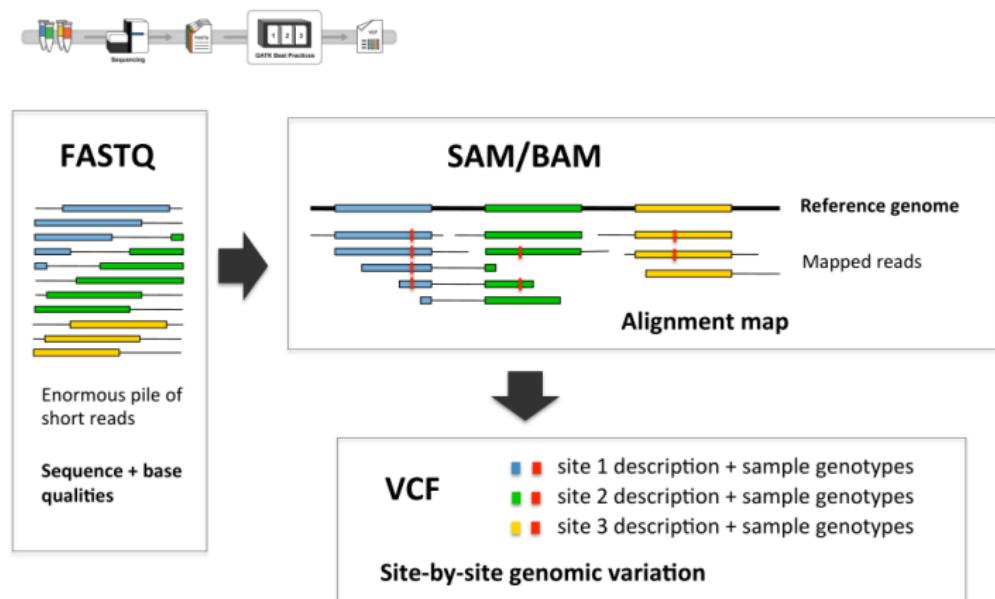
GATK a primer

Detecting variants with reads



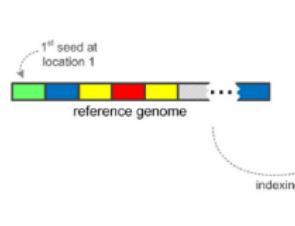
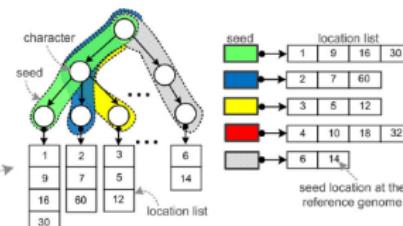
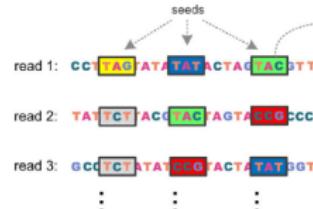
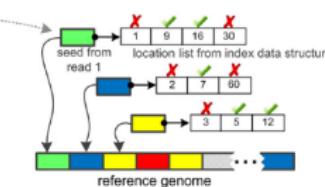
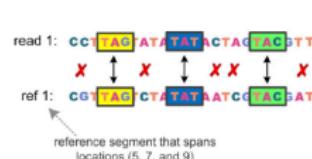
GATK a primer

Variant discovery pipeline



GATK a primer

Map to the reference

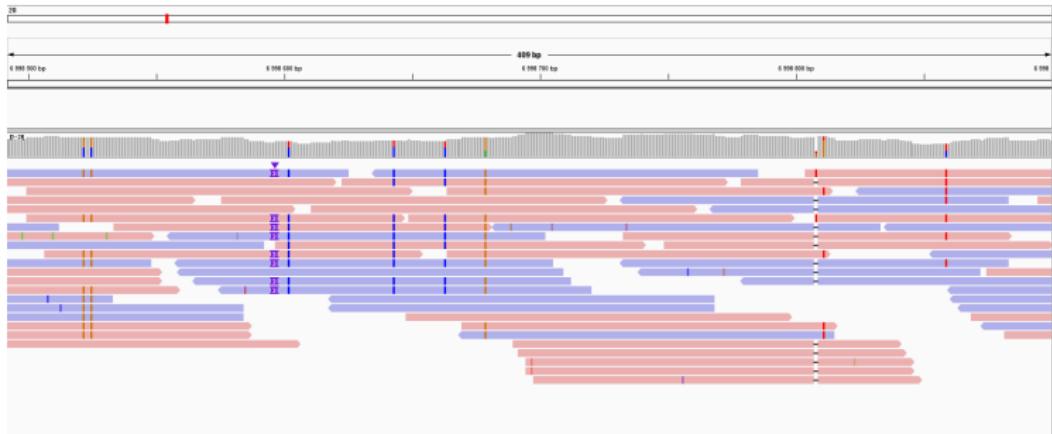
a. Seed extraction from reference genome**b. Seed indexing using suffix tree or hash table****c. Seed extraction from reads****d. Seed querying and filtering****e. Seed chaining and pre-alignment filtering****f. Alignment verification**

C	G	T	T	A	T	G	T	C	T	A	...
0	0	0	0	0	0	0	0	0	0	0	
0	2	2	2	2	2	2	2	2	2	2	
C	0	2	3	3	3	3	3	4	4	4	
T	0	2	3	3	7	5	2	2	6	6	
T	0	2	3	5	2	7	7	7	7	7	
A	0	3	3	2	7	9	9	9	9	9	
G	0	2	4	5	7	9	11	11	11	11	
T	0	2	4	6	7	9	11	11	13	13	
A	0	2	4	6	7	9	11	12	14	14	
T	0	2	4	6	8	9	11	13	14	16	

.bam/.sam file contains necessary alignment information (e.g., type, location, and number of each edit)

Technology dictates algorithms: recent developments in read alignment

Starting from Aligned reads



- Reads exhibit discrepancies with the reference assembly
- How can we decide if it is an SNP or not?

Empirical approach

		SNP ?	
		↓	
ref	gagatctgGagcatgCaggtggatgtca	T	40
read2	gagatctgGGagcatgCaggtggatgtca	C	40
read3	gagatctgGGagcatgCaggtggatgtca	C	40
read4	gagatctgGGagcatgCaggtggatgtca	C	40
read5	gagatctgGGagcatgCaggtggatgtca	C	40
read6	gagatctgGGagcatgAaggtggatgtca	T	40
read7	gagatctgGGagcatgAaggtggatgtca	T	40
read8	gagatctgGGagcatgAaggtggatgtca	T	40
read9	gagatctg..agcatgAaggtggatgtca	T	38
read10	gagatctgGGagcatgCaggtggatgtca	C	40
read11	gagatctgGGagcatgCaggtggatgtca	G	40

- Count the reads associated with each allele : 4 A and 7 C
- Take a decision according to the proportion of reads for each allele \Rightarrow heterozygote

Empirical approach

Caveats

- What is a reasonable choice for the minimum required number of alternate alleles ?
- Cutoff for deciding when an individual is heterozygote ?
- How can we handle the uncertainty inherent of low-coverage data ?
- How to deal with the sequencing errors ?
- No measure of uncertainty in the genotype inference

Probabilistic approach

SNP ?
↓

		10	20	30	40	
ref	gagatctgGGagcatgCaggtggatgtcaTcagagaTacc					40
read2	gagatctgGGagcatgCaggtggatgtcaCcagagaTacc					40
read3	gagatctgGGagcatgCaggtggatgtcaCcagagaTacc					40
read4	gagatctgGGagcatgCaggtggatgtcaCcagagaTacc					40
read5	gagatctgGGagcatgCaggtggatgtcaCcagagaTacc					40
read6	gagatctgGGagcatgAaggtggatgtcaTcagagaTacc					40
read7	gagatctgGGagcatgAaggtggatgtcaTcagagaGacc					40
read8	gagatctgGGagcatgAaggtggatgtcaTcagagaTacc					40
read9	gagatctg..agcatgAaggtggatgtcaTcagagaTacc					38
read10	gagatctgGGagcatgCaggtggatgtcaCcagagaTacc					40
read11	gagatctgGGagcatgCaggtggatgtcaCcagagaGacc					40

- What is the probability that a column signs an SNP ?
- What is the most likely genotype for that SNP ?

Introduction
oooooooooooo

Detection
ooo●oooooooo

Genotype likelihoods
oooooooooooo

Complications
oooooooo

Probabilistic approach

A single individual

Probabilistic approach

$$P(G_i|D) = \frac{P(G_i)P(D|G_i)}{P(D)}$$

- G_i : one of the possible genotypes
- D : data, i.e. a column of the multiple alignment of reads
- Suppose for the moment a single diploïd individual, hence 3 possible genotypes {0/0, 0/1, 1/1}

Bayes theorem

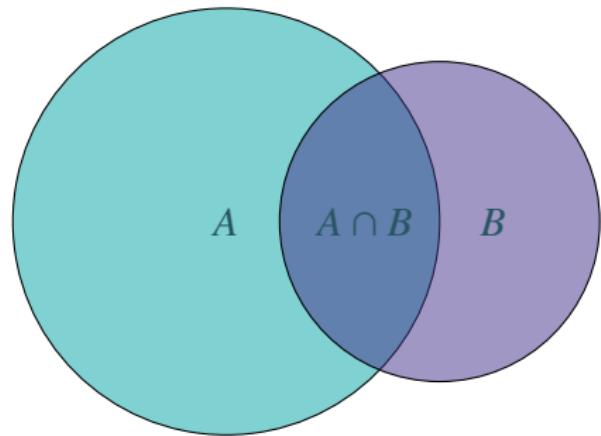
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

and

$$P(A \cap B) = P(A)P(B|A)$$

hence

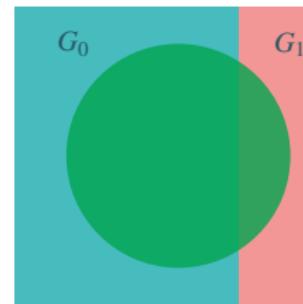
$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}$$



Bayes theorem

Consider a haploid genome, hence two genotypes 0 and 1

$$P(G_i|D) = \frac{P(G_i)P(D|G_i)}{P(D)}$$



$$P(D) = P(D \cap G_0) + P(D \cap G_1)$$

hence

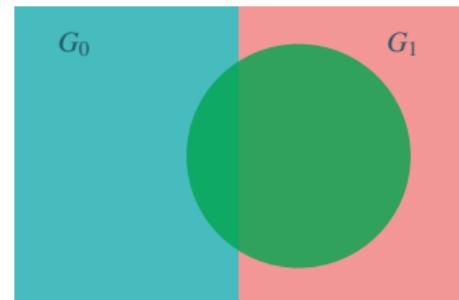
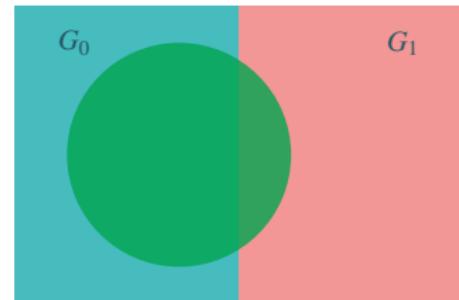
$$P(G_i|D) = \frac{P(G_i)P(D|G_i)}{P(G_0)P(D|G_0) + P(G_1)P(D|G_1)}$$

Bayes theorem

Consider a haploid genome, hence two genotypes 0 and 1

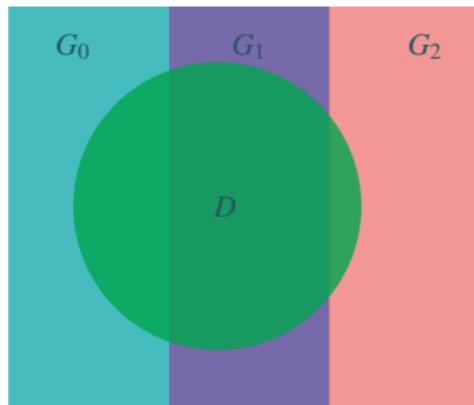
$$P(G_i|D) = \frac{P(G_i)P(D|G_i)}{P(D)}$$

$$\propto P(G_i)P(D|G_i)$$



Bayes theorem

- Consider a diploid genome with 3 genotypes, $G_0 = 0/0$, $G_1 = 0/1$ and $G_2 = 1/1$
- Reads data D gives us some clues about the most likely genotype



Probabilistic approach

$$P(G_i|D) \propto P(G_i)P(D|G_i)$$

Probability of data given the genotype

- $P(D|G_i)$ is the probability to observe the reads for a given genotype
- In the absence of sequencing errors it is straightforward

Probabilistic approach

	SNP ?	
	↓	
ref	gagatctgG G agcatgCaggtggatgtca T cagagaTacc	40
read2	gagatctgGGagcatgCaggtggatgtcaC c agagaTacc	40
read3	gagatctgGGagcatgCaggtggatgtcaC c agagaTacc	40
read4	gagatctgGGagcatgCaggtggatgtcaC c agagaTacc	40
read5	gagatctgGGagcatgCaggtggatgtcaC c agagaTacc	40
read6	gagatctgGGagcatgAaggtggatgtcaTcagagaTacc	40
read7	gagatctgGGagcatgAaggtggatgtcaTcagagaGacc	40
read8	gagatctgGGagcatgAaggtggatgtcaTcagagaTacc	40
read9	gagatctg..agcatgAaggtggatgtcaTcagagaTacc	38
read10	gagatctgGGagcatgCaggtggatgtcaC c agagaTacc	40
read11	gagatctgGGagcatgCaggtggatgtcaC c agagaGacc	40

Figure – A candidate SNP with 10 reads, 6 C and 4 A

- Suppose a candidate SNP with two alleles (0 and 1)

$$P(D|G_i)$$

- Let ε be the error rate, and
- D be a set of n reads, k reads with 0 and $(n - k)$ reads with 1

$$P(D|00) = \binom{n}{k} (1 - \varepsilon)^k \varepsilon^{n-k}$$

$$P(D|11) = \binom{n}{k} \varepsilon^k (1 - \varepsilon)^{n-k}$$

$$P(D|01) = \binom{n}{k} p^k (1 - p)^{n-k}$$

with

$$p = \frac{1}{2}\varepsilon + \frac{1}{2}(1 - \varepsilon) = \frac{1}{2}$$

- The right hand sides are called the genotype likelihoods

$$P(G|D) \propto L(G|D) = P(D|G)$$

$L(G_i) = P(D|G_i)$ the genotype likelihood

We are (only) interested in finding the most likely genotype

$$P(D|00) \propto (1 - \varepsilon)^k \varepsilon^{n-k}$$

$$P(D|11) \propto \varepsilon^k (1 - \varepsilon)^{n-k}$$

$$P(D|01) \propto \left(\frac{1}{2}\right)^n$$

- The genotype caller provides the genotype likelihoods coded using a Phred score.

$$PL = -10 \lfloor \log_{10}(P(D|G)) \rfloor$$

where $\lfloor x \rfloor$ means "the largest integer smaller than x"

- The **called** genotype is the one with the higher likelihood

Introduction
○○○○○○○○○○

Detection
○○○○○○○○○○○○

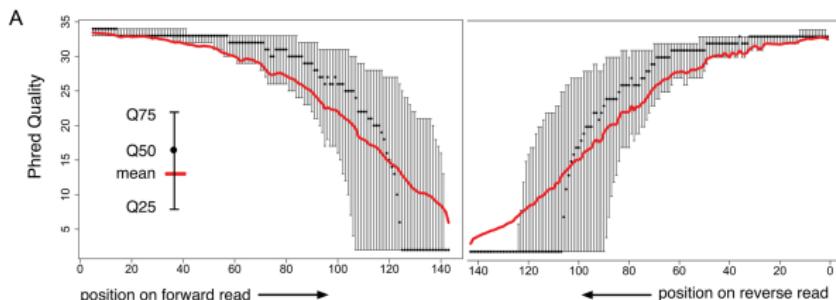
Genotype likelihoods
○○●○○○○○○○○○○

Complications
○○○○○○○○

Probabilistic approach

BaseQuality Score

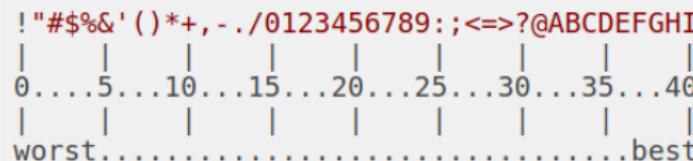
Sequencing errors



- DNA sequencing is not exempt of errors
- Manufacturers, like Illumina, provide Base Quality Score coded in a phred score

Sequence quality

FASTQ quality values



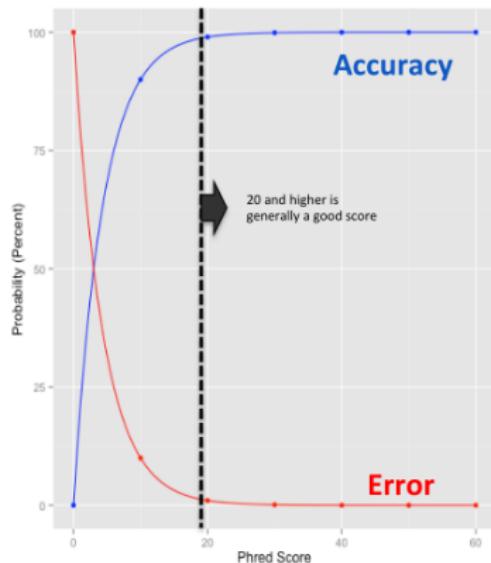
- ASCII coding of the quality (*Phred* quality score)

$$Q = -10 \lfloor \log_{10}(p) \rfloor$$

Q	p	erreur
10	10^{-1}	10%
20	10^{-2}	1%
30	10^{-3}	1‰

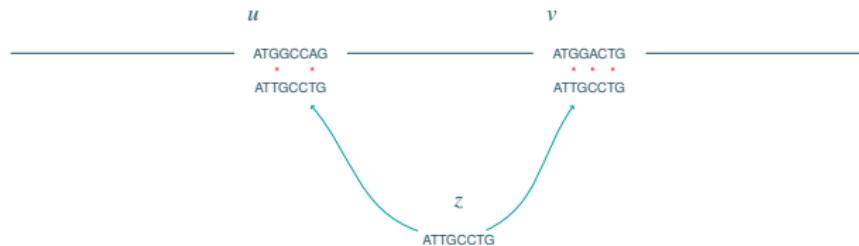
Sequencing errors

- $\text{Phred value} = -10 * \log_{10}(\epsilon)$
- Examples:
 - 90% confidence (10% error rate) = Q10
 - 99% confidence (1% error rate) = Q20
 - 99.9% confidence (.1% error rate) = Q30
- SAM encoding adds 33 to the value
(because ASCII 33 is the first visible character)



GATK a primer

Mapping quality



Let $p(z|u)$, the probability of z aligned at u , here for example with two errors

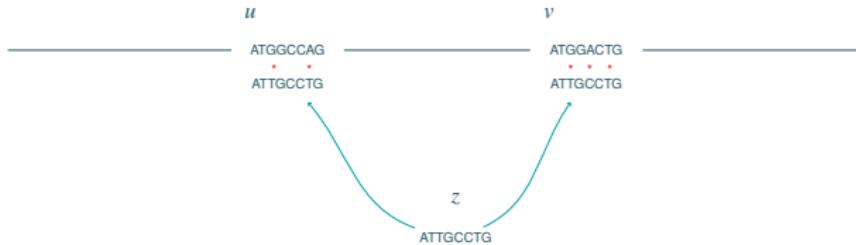
$$p(z|u) = \varepsilon_3 * \varepsilon_7$$

the probability that the read R comes from position u can be computed using

$$p(u|z) = \frac{p(u)p(z|u)}{\sum_v p(v)p(z|v)} = \frac{p(z|u)}{\sum_v p(z|v)}$$

where the sum \sum_v is over all the alternative positions.

Mapping quality



The mapping quality is defined by

$$\begin{aligned} mapQ &= -10 \log_{10} [\text{read is wrongly mapped}] \\ &= -10 \left[\log_{10} (1 - p(u|z)) \right] \end{aligned}$$

$mapqQ = 0 \implies$ many alternate positions are equally likely.

$L(G_i)$: genotype likelihood

- The base quality provided in the (fastq) sam files enables to compute a position-specific, read-specific error probability
- If the three genotypes are coded 0, 1, 2, the likelihood takes the form

$$\mathcal{L}(g) = \left(\frac{1}{2}\right)^n \prod_{j=1}^k [(2-g)\varepsilon_j + g(1-\varepsilon_j)] \prod_{j=k+1}^n [(2-g)(1-\varepsilon_j) + g\varepsilon_j]$$

Easy to compute !

Probabilistic approach

Prior on genotype

$$P(G_i|D) = \frac{P(G_i)P(D|G_i)}{P(D)}$$

How do we compute the prior on genotypes ?

$P(G_i)$: prior on genotypes

- Single sample case and diploid. The current practice is a uniform prior on the 3 genotypes {0/0, 0/1, 1/1}
- Multiple samples case : all the different genotype configurations are not equally likely, see below

Genotype likelihood : synthèse

- La génotype d'un individu est celui qui maximise la vraisemblance
- C'est à dire $P(D|G)$ la probabilité des données sachant le génotype
- Une mesure d'incertitude est fournie pour chaque génotype sous la forme d'une différence entre les log-vraisemblance des différents génotypes

Introduction
oooooooooooo

Detection
oooooooooooo

Genotype likelihoods
oooooooooooo●

Complications
oooooooooooo

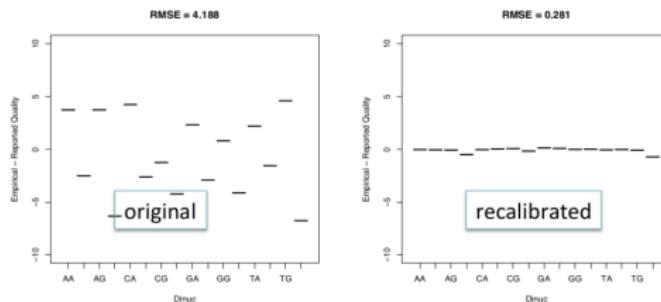
Probabilistic approach

Complications

Complications

- Base quality score are not correctly calibrated \Rightarrow BQSR
- Alignments around indels are problematic \Rightarrow Indel Realignment
- Considering a single column is problematic \Rightarrow use haplotypes

BQSR : Base Quality Score Recalibration

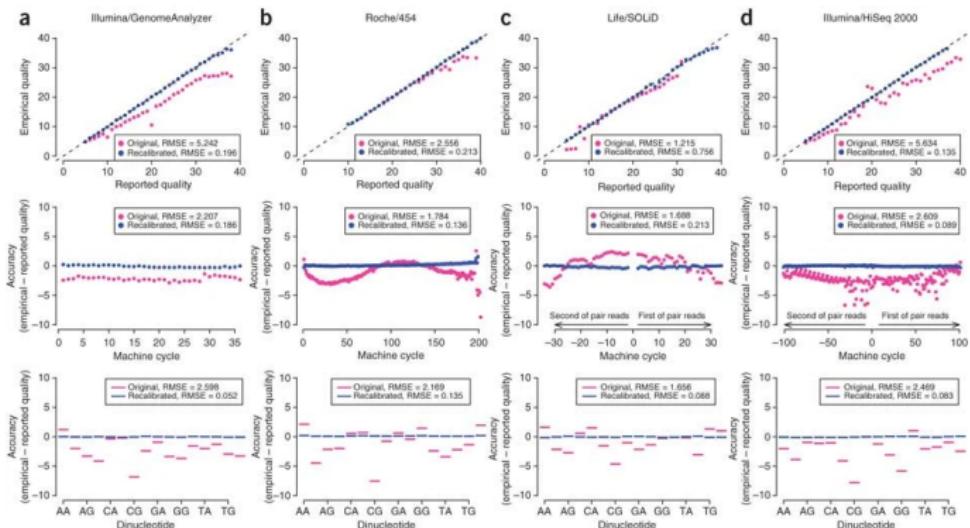


- Error rates can be corrected from the set of reads
- From all the reads, tabulate in bins (original quality score, position in read, dinucleotide context)
- In each bin, count the number of bases and the number of mismatches \Rightarrow new error estimates

GATK a primer

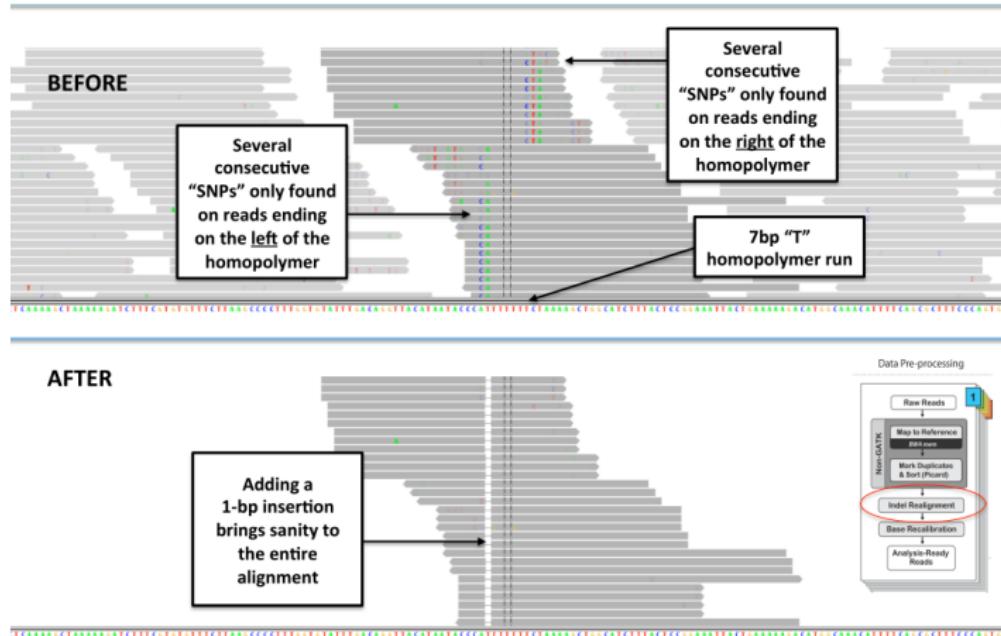


BQSR : Base Quality Score Recalibration



A framework for variation discovery and genotyping using next-generation DNA sequencing data

Indel realignment



GATK a primer



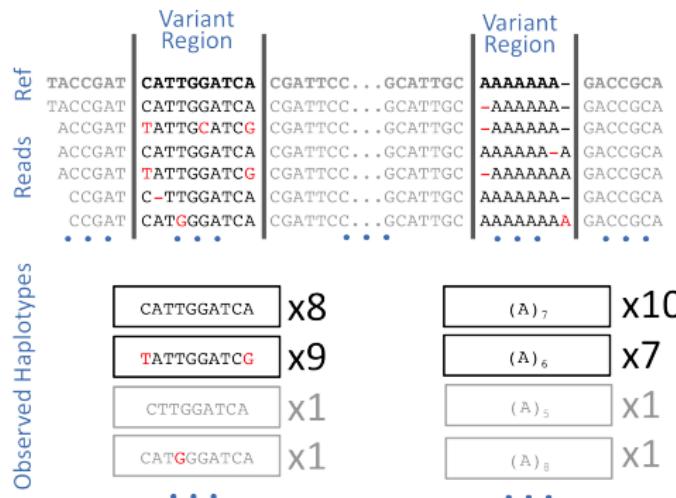
Indel realignment



GATK a primer

Haplotypes with Freebayes

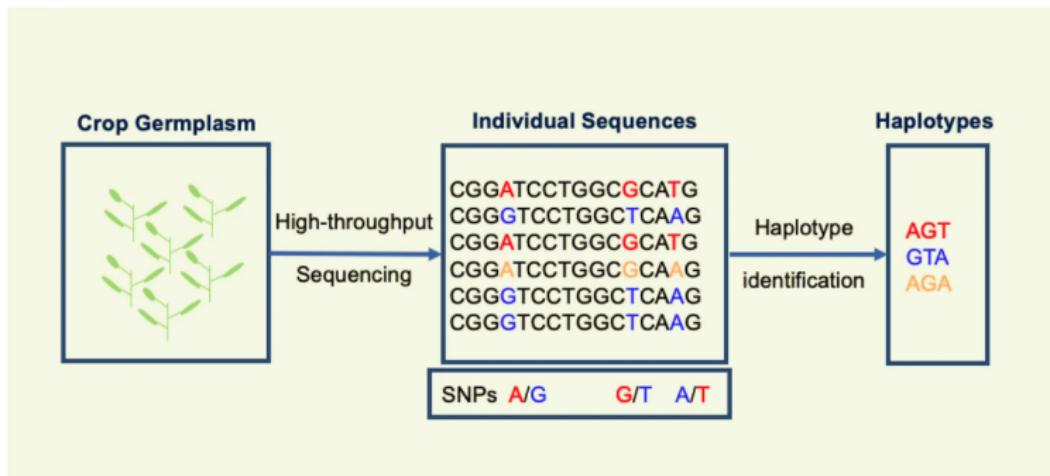
- Indels, and more generally alignments are sometime problematic \Rightarrow seeking to free itself from alignments



Freebayes

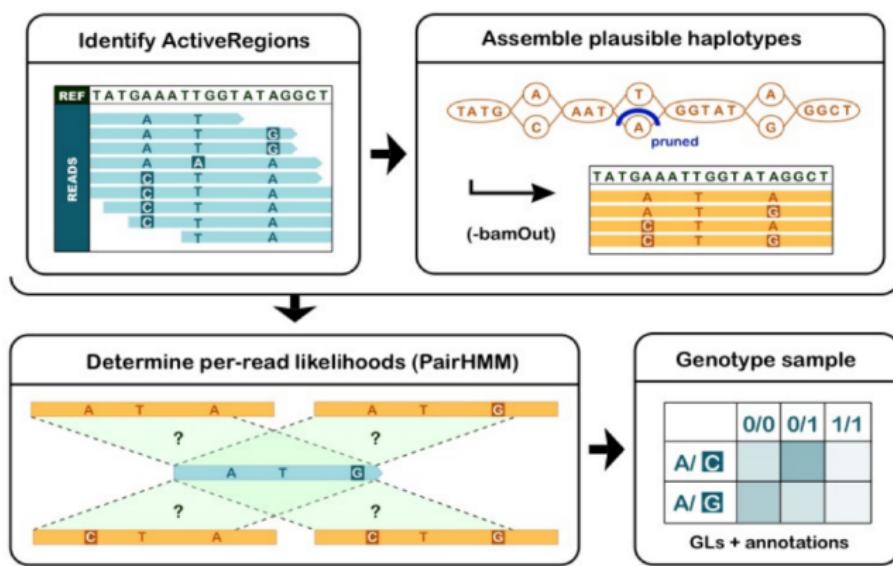
Haplotype : recall

- Haplotypes describe the succession of alleles along a chromosome (segment)



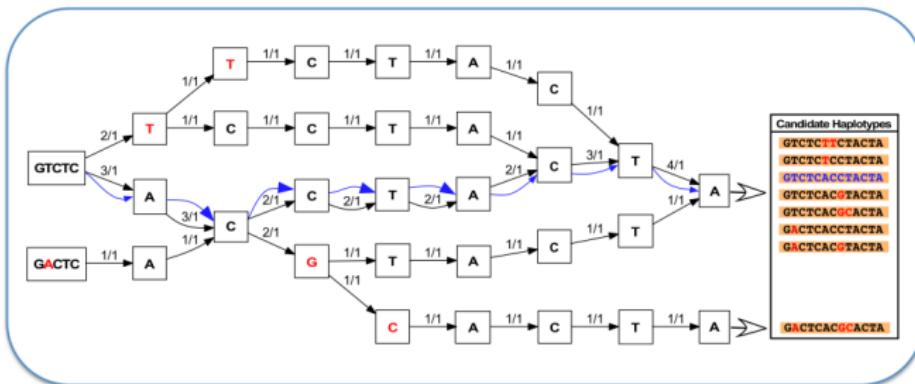
Features and applications of haplotypes in crop breeding

Haplotypes with HaplotypeCaller



A framework for variation discovery and genotyping using next-generation DNA sequencing data

Haplotypes with HaplotypeCaller



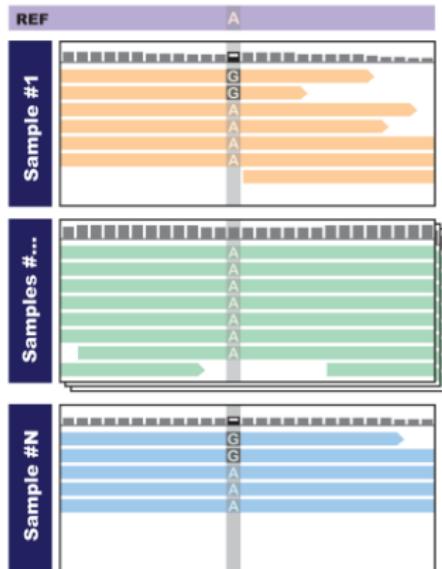
- Previous alignments are ignored
 - K-mers consist of every possible sequence combination based on the reads
 - Most likely paths through the graph are scored

GATKr12-5-Variant calling joint genotyping

Probabilistic approach

Multiple individuals

Multiple samples

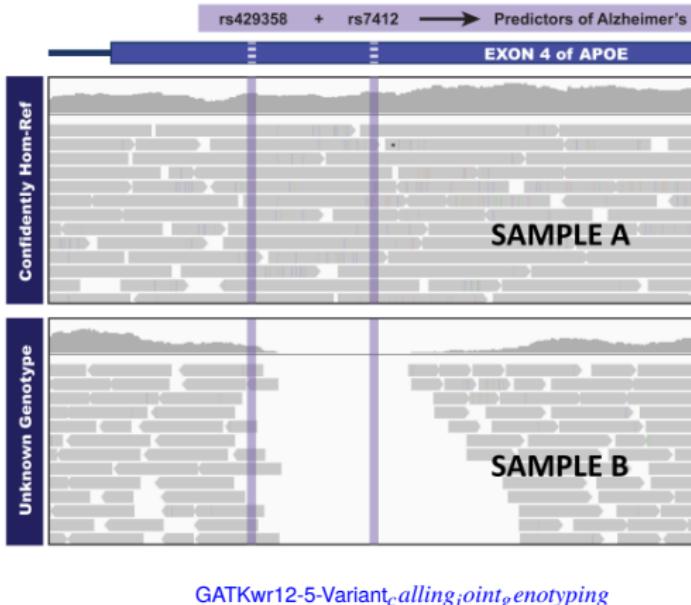


- Sample #1 or Sample #N alone:
 - weak evidence for variant
 - may miss calling the variant
- Both samples seen together:
 - unlikely to be artifact
 - call the variant more confidently

GATKr12-5-Variant calling joint genotyping

Multiple samples

- Analyzed individually:**
 - No call for either sample
 - Very different reasons!
- In joint analysis with other samples:**
 - Hom-ref call and no-call genotypes emitted



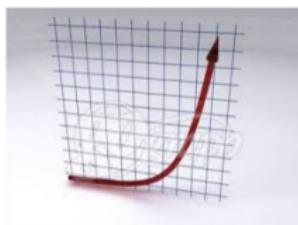
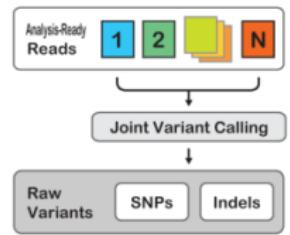
GATKr12-5-Variant_calling_joint_genotyping

Probabilistic approach

$$P(G_1, \dots, G_n | R_1, \dots, R_n) = \frac{P(R_1, \dots, R_n) P(R_1, \dots, R_n | G_1, \dots, G_n)}{P(R_1, \dots, R_n)}$$

- G_1, \dots, G_n : the possible genotypes for the n individuals
- R_1, \dots, R_n : data, i.e. a column of the multiple alignment of reads for all individuals
- The major difference is that the prior probability on genotypes, $P(R_1, \dots, R_n)$ are generally not taken as uniform
- All the different genotype configurations are not equally likely : for example say 10 0/0, 20 0/1 and 150 1/1) (think of Hardy Weinberg) \implies population priors (freebayes, bcftools call, GLnexus)

Multiple samples

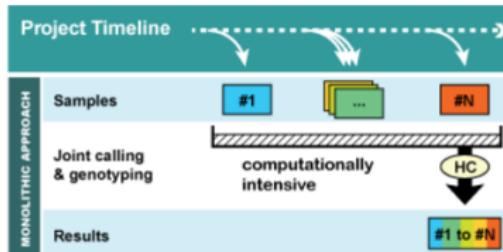


gg4366215 www.gograph.com

Compute requirements scale exponentially with number of samples

(combinatorial problem)

It gives us the right answers, but...



Want to add new samples?

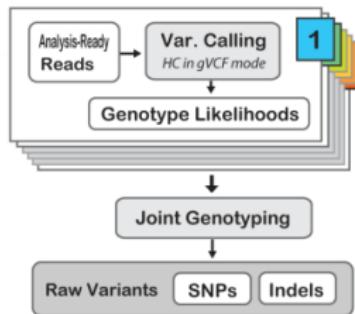
Got to re-run pipeline from scratch!

GATKwr12-5-Variant_calling joint_genotyping

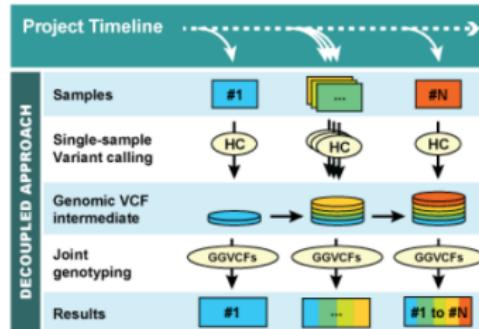
Multiple samples : gvcf files



Scalable over sample size



+ Incremental over time



GATKr12-5-Variant calling joint genotyping

Variant quality score

Each variant is associated with an uncertainty measure (field QUAL in vcf file), coded as a Phred score

$$QUAL = -10 \lfloor \log_{10}(P[\text{not a variant site}]) \rfloor$$

$$= -10 \lfloor \log_{10}(P(G = 00|D)) \rfloor$$

In the case of multiple individuals it measures the probability of having at least one alternate allele.

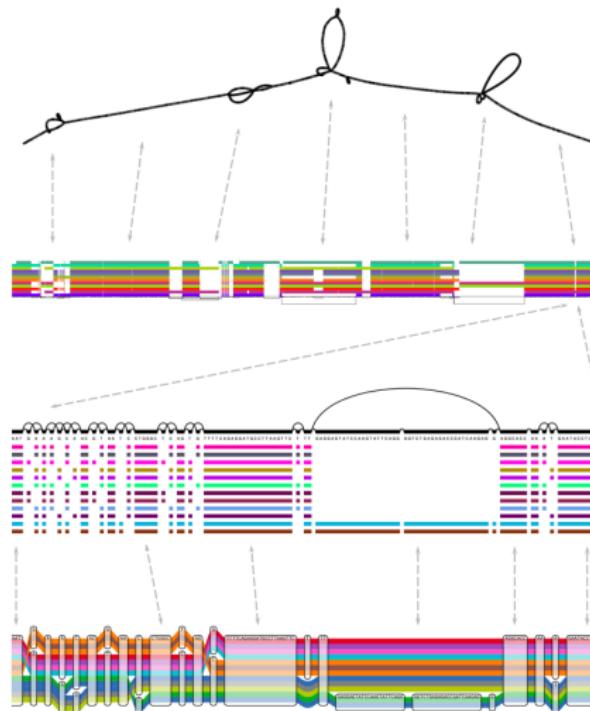
Variant detection today

- Structural variant detection and variation graphs
 - ▶ In order to bypass the bias induced by the use of a single reference assembly, current approaches aim at combining multiple genomes in a single graph structure and use this structure as a reference on which the reads are mapped
- Variant detection using deep learning
 - ▶ Recent developments in variant discovery use deep learning approach to distinguish real variants from sequencing errors or alignment artifacts

Multiple samples
○○○○○○

Variation detection today
●●○○

Variation graphs



Sequence tube maps: making graph genomes intuitive to commuters

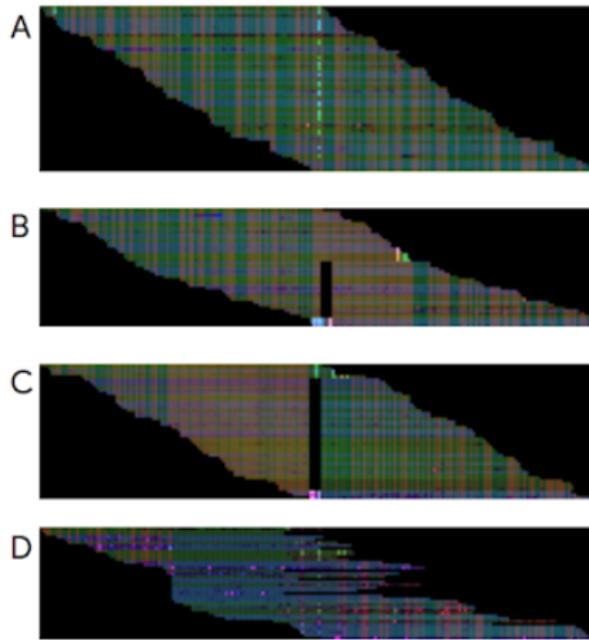
Détection de variants



Thomas Faraut

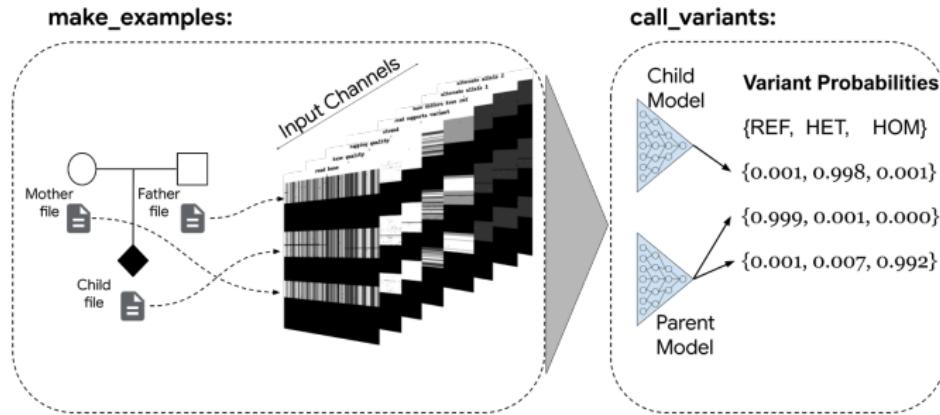
54/57

Deep learning for Variant detection



A universal SNP and small-indel variant caller using deep neural networks

Deep learning for Variant detection



A universal SNP and small-indel variant caller using deep neural networks

Deep learning in bioinformatics



Johann Gregor Mendel: the victory of statistics over human imagination