

Bioinformatique pour le traitement des données de séquençage

Détection de variants

Maud Gautier, Annabelle Haudry, Thibault Latrille

17 Septembre 2019



Qu'est ce que la détection de variants ?



Détection de variants

Variant homozygote

Variant hétérozygote

AATTTATTATTAGGCGATACGGAGGCCGGAGCAGAGACAGC
ATTATTATTAGGCGATATGGAGGCAGAGCAGAGTCAGC

Individu diploïde

AATTTATTTTATAGGCGATACGGAGGCCAGAGCAGAGTCAGC

Genome de référence

Qu'est ce que la détection de variants ?

ACTGATCGATCGTACGTAGCTGA
TATGCTGCATGCATGCATG
CAGTCGATCGATCGT
CAGTCGATCGATCGTCAGTCGATCGAT
TATGCTGCATGCATGCA
GGGTATTATATATCTAGCT
GTATCTACGACTACTGCTACTGAC
TTATTACTGACTCGATGCA
GGGGATCTAAGCTGAGCC
GTATCTACGACTAGCT
GAGCTTTGAGTCG

Données de séquençage d'un individu diploïde (.fastq)

Alignement (matinée)

ATCTTATTATTAG TACGGAGGCGGAGCA C
ATTTTATTAT TACGGAGGCGGAGCA AGC
AATTTATT ATACGGATGCAGGGA AGC
AATT AGGCGATATGGAGG ACAGC
ATTATTAGGCGATACGGAGG CAGAGTCAGC
ATTATTAGGCGATACGGA CAGAGACAGC
TATTATTTGGCGATA AGCAGAGTCAGC
TATTATTAGGCGATA GGAGCAGAGACAG
TTTATTATTAGGCGAT AGGCGGAGCAGAGAC
AATTTATCATTAGG GGAGGCAGAGCAGAG
AATTTATTTTATAGGCGATACGGAGGCAGAGTCAGC

Données de séquençage alignées (.bam et .sam)

Genome de référence (.fasta)

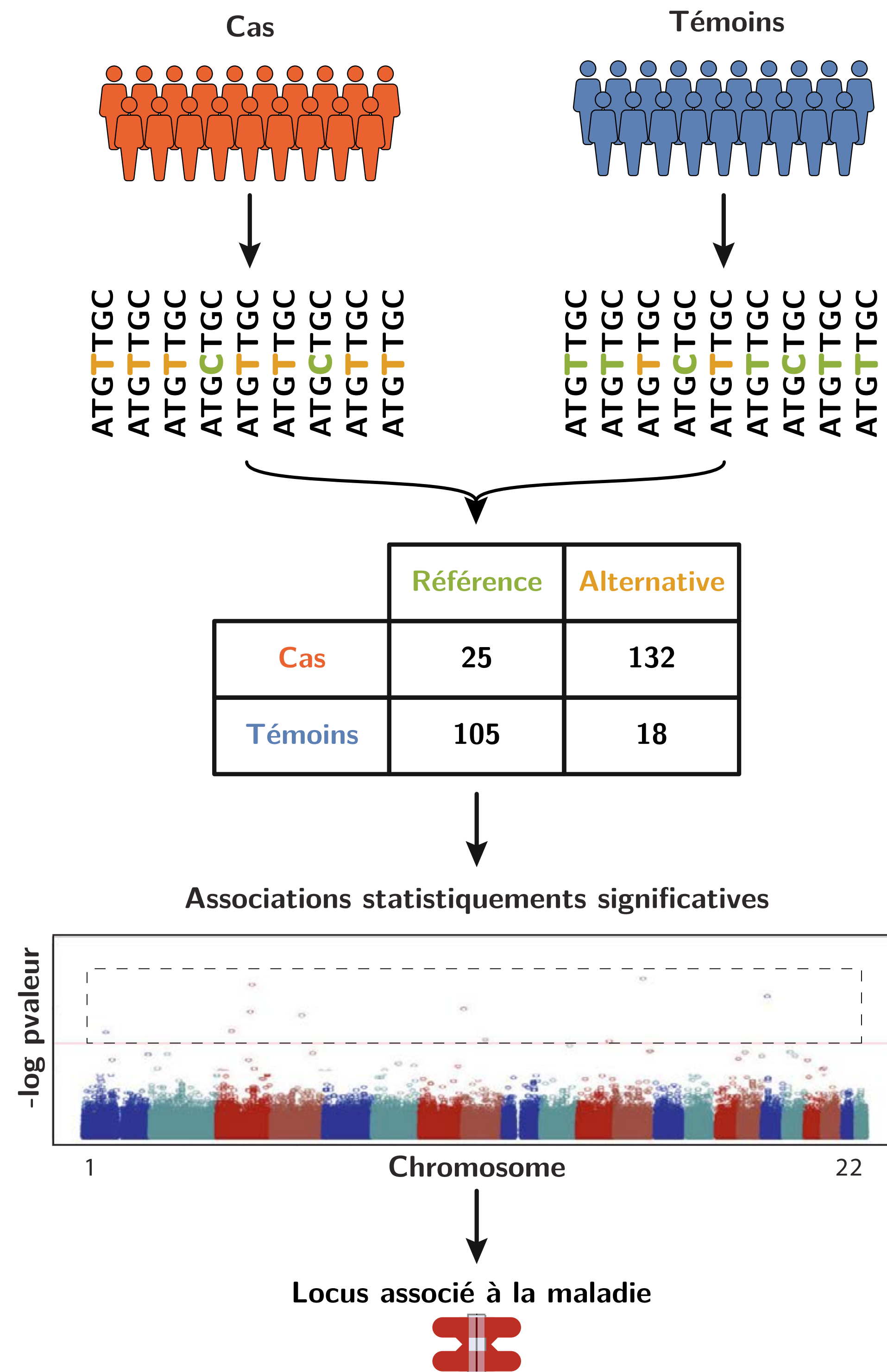
Détection de variants (après-midi)

AATTTATTATTAGGCGATACGGAGGCGGAGCAGAGACAGC
ATTTTATTATTAGGCGATATGGAGGCAGAGCAGAGTCAGC
AATTTATTTTATAGGCGATACGGAGGCAGAGTCAGC

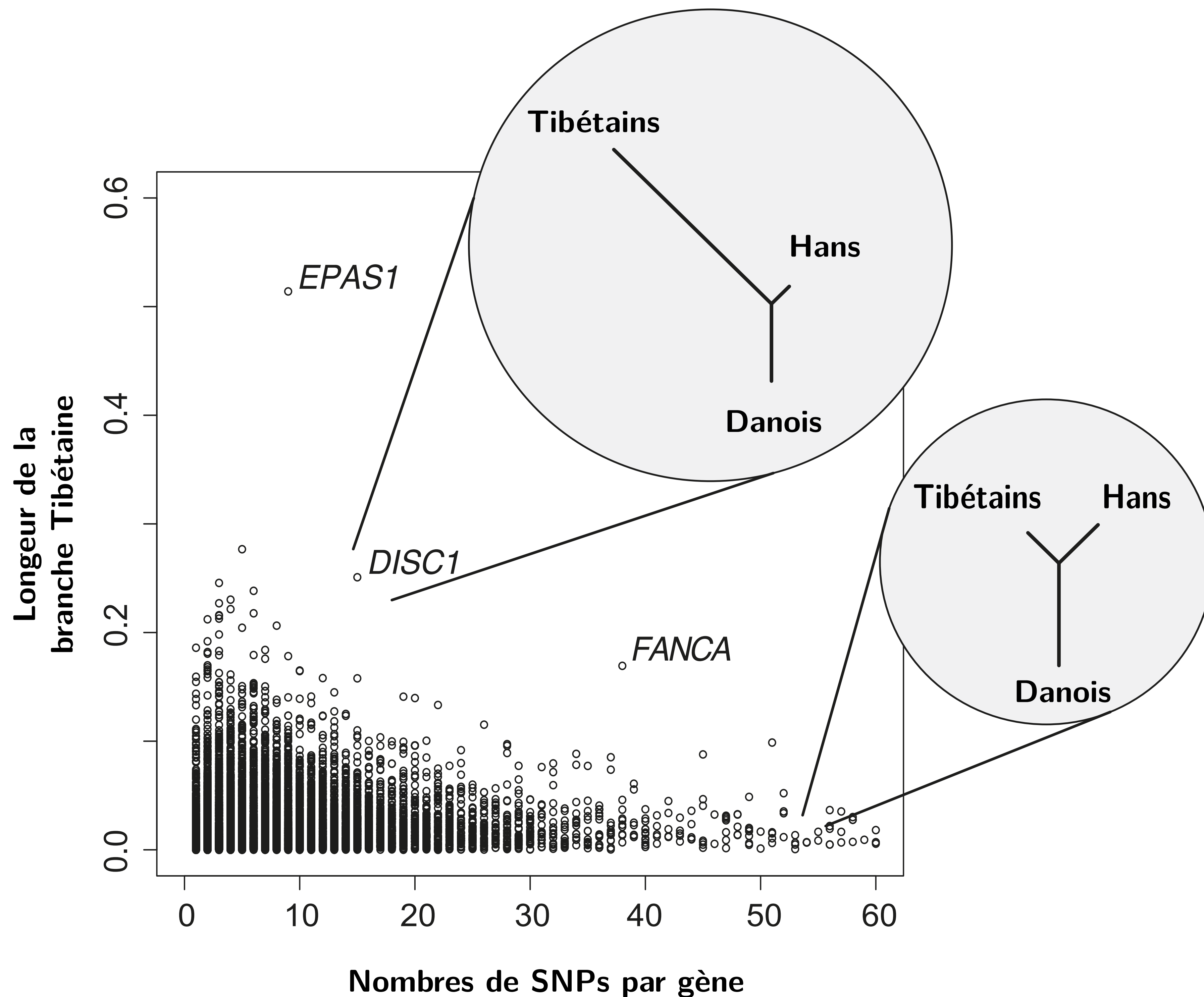
Individu diploïde (.vcf)

Genome de référence (.fasta)

Pourquoi cherche-t-on à détecter des variants ?



Quels gènes sont sélectionnés ?

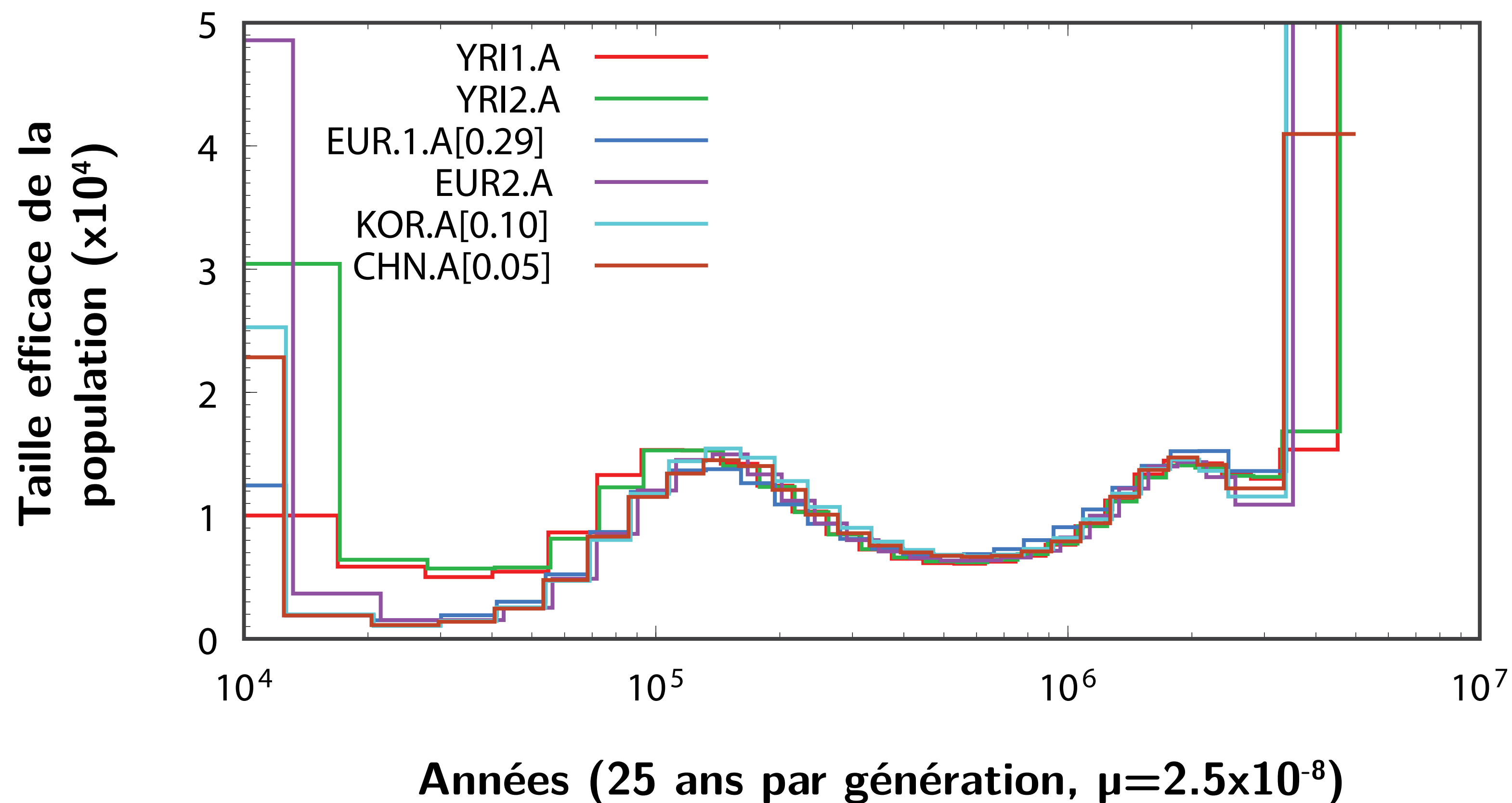


Yi *et al*, Science (2009)

Nos ancêtres étaient-ils nombreux ?

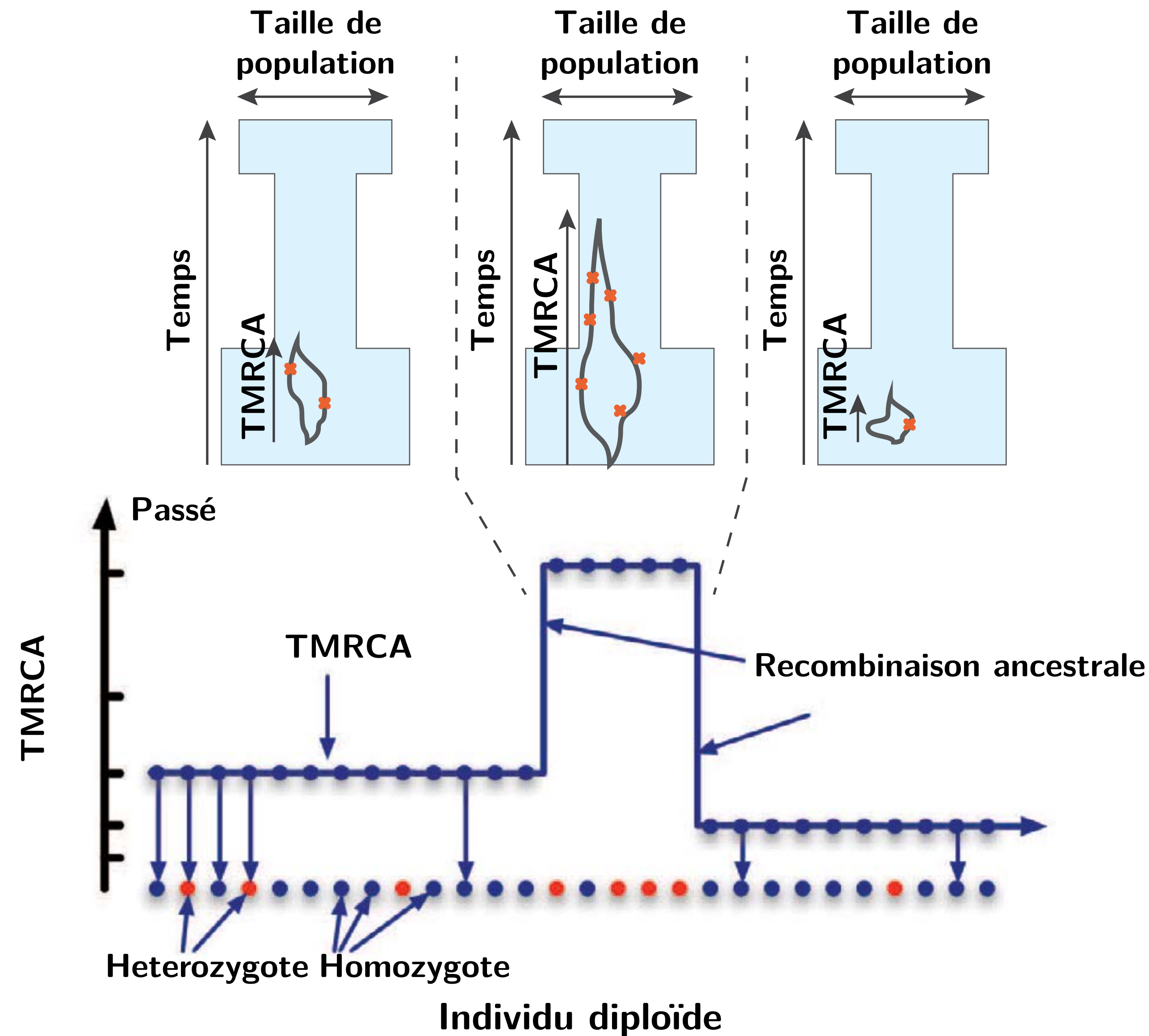
-----AATTATTATTAGGCGATACGGAGGC GGAGCAGAG ACAGC-----
-----A TTTTATTATTAGGCGATA TGGAGGCAGAGCAGAGTCAGC----- Individu diploïde

Inference of human population history
from individual whole-genome sequences,
H. Li & R. Durbin, *Nature* (2011)

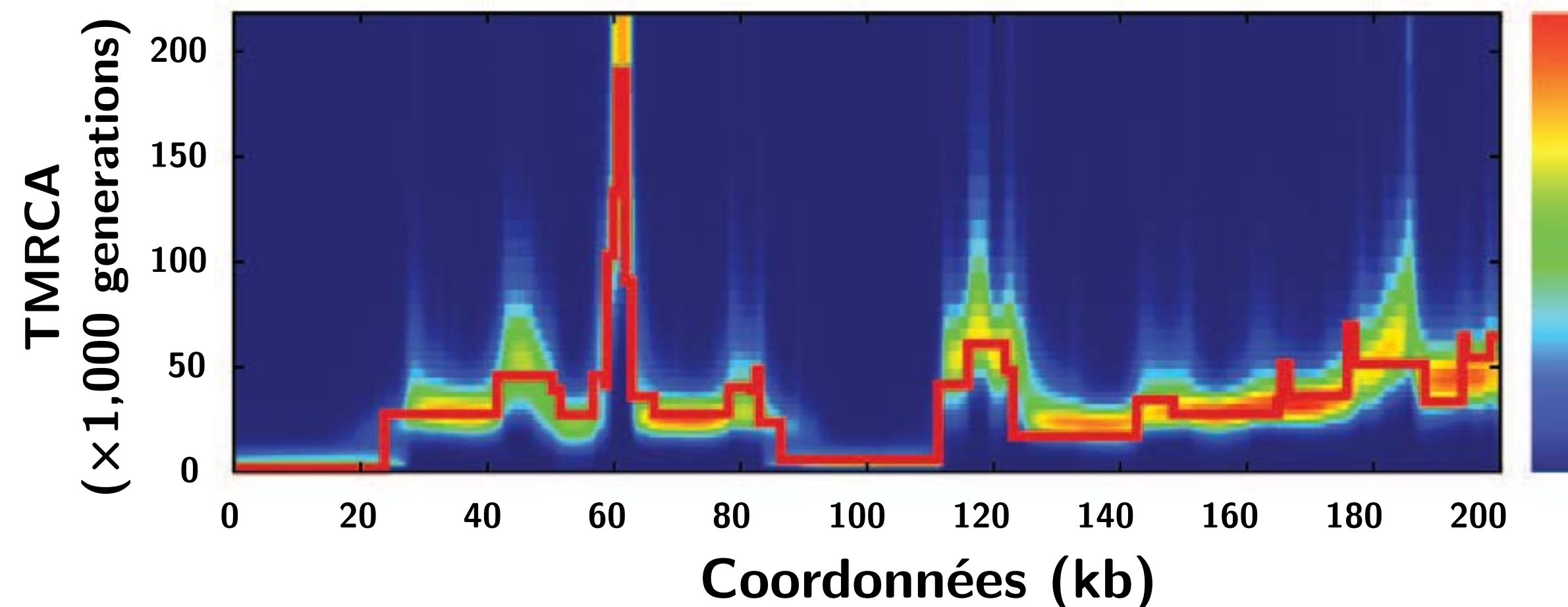


Li & Durbin, *Nature* (2011)

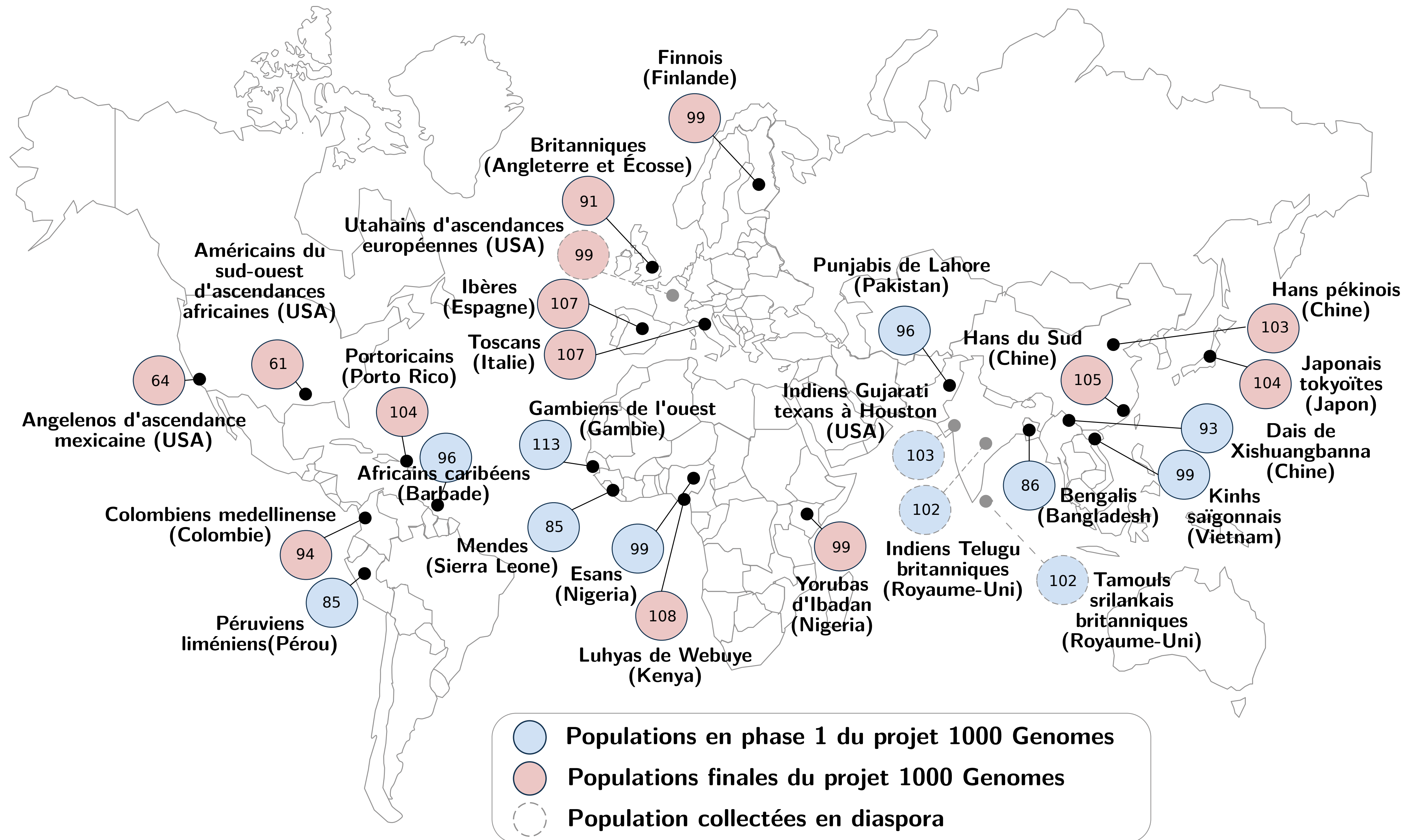
Nos ancêtres étaient-ils nombreux ?



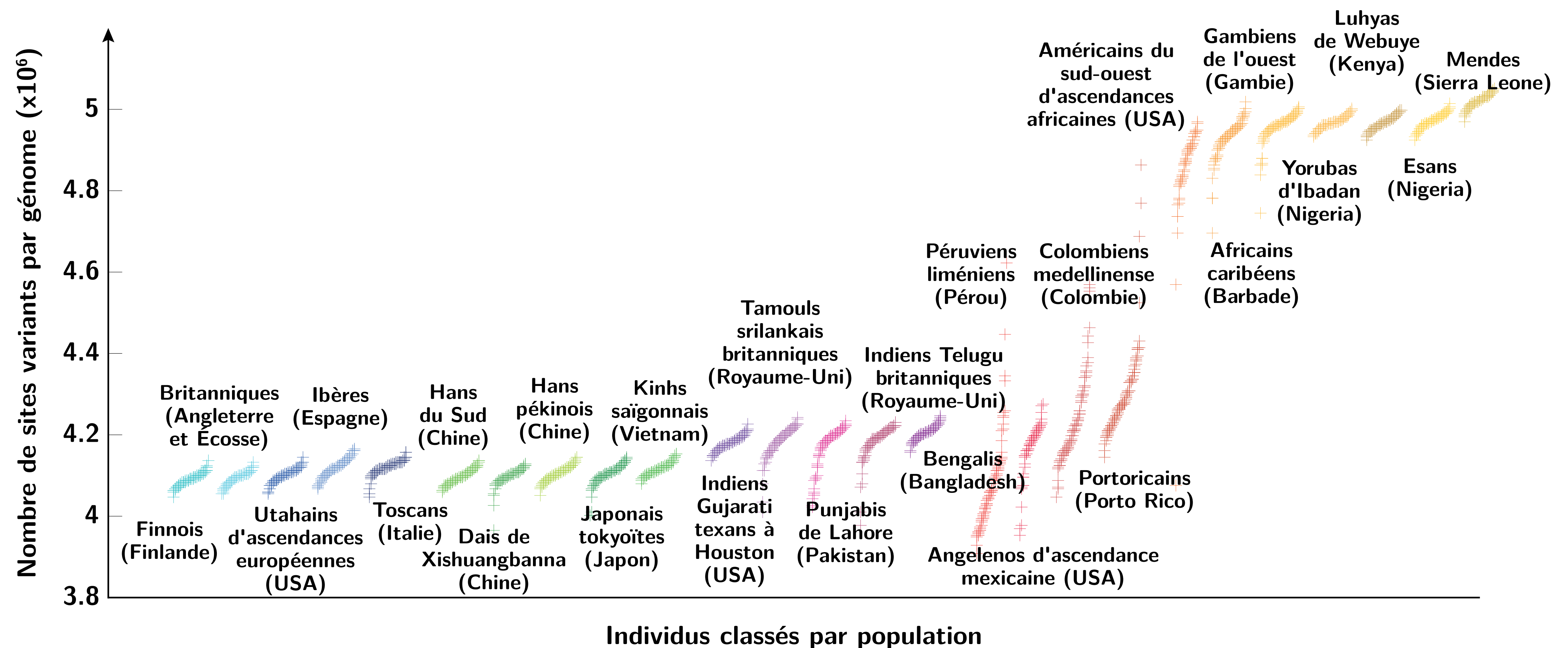
TMRCA:
Temps jusqu'à
l'ancêtre commun
le plus récent



Quelles données va-t-on utiliser ?

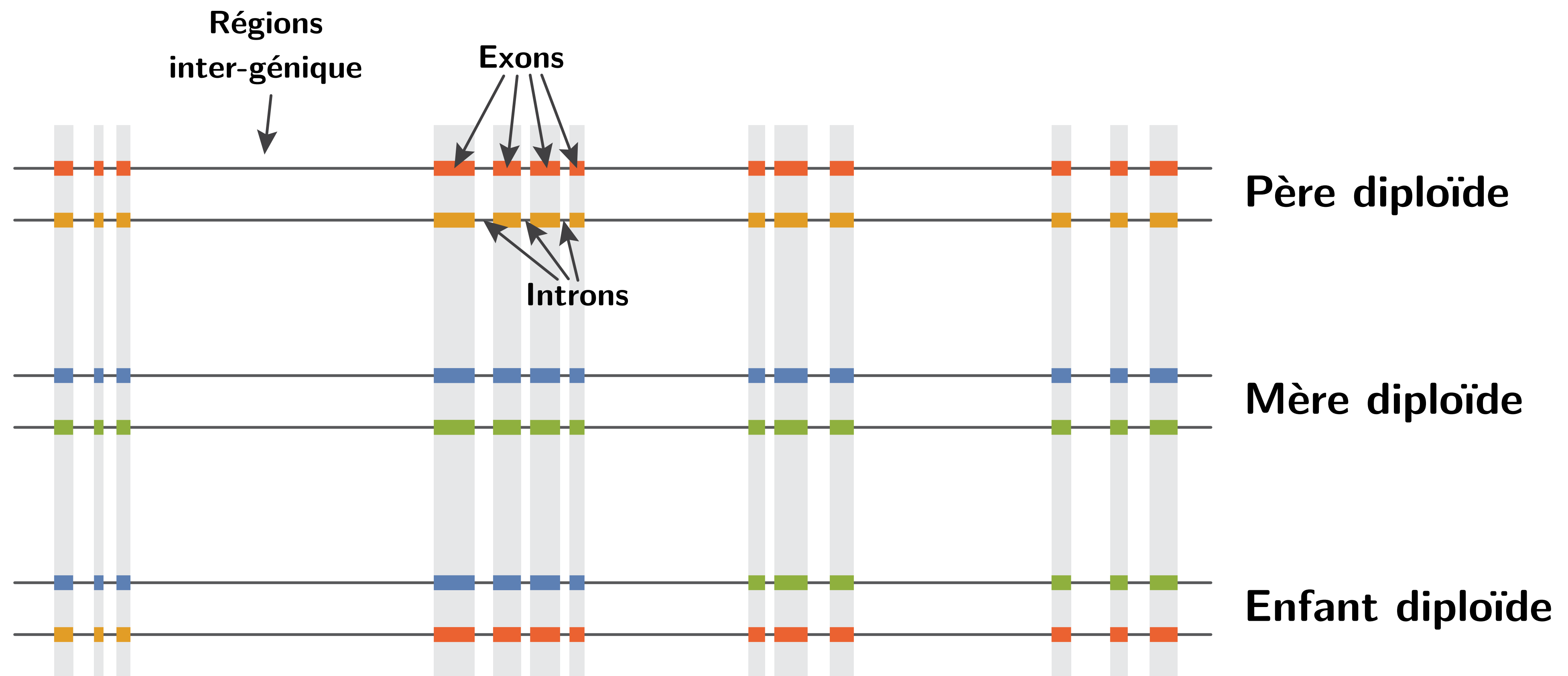


Quels sont les résultats du projet 1000 Génomes ?



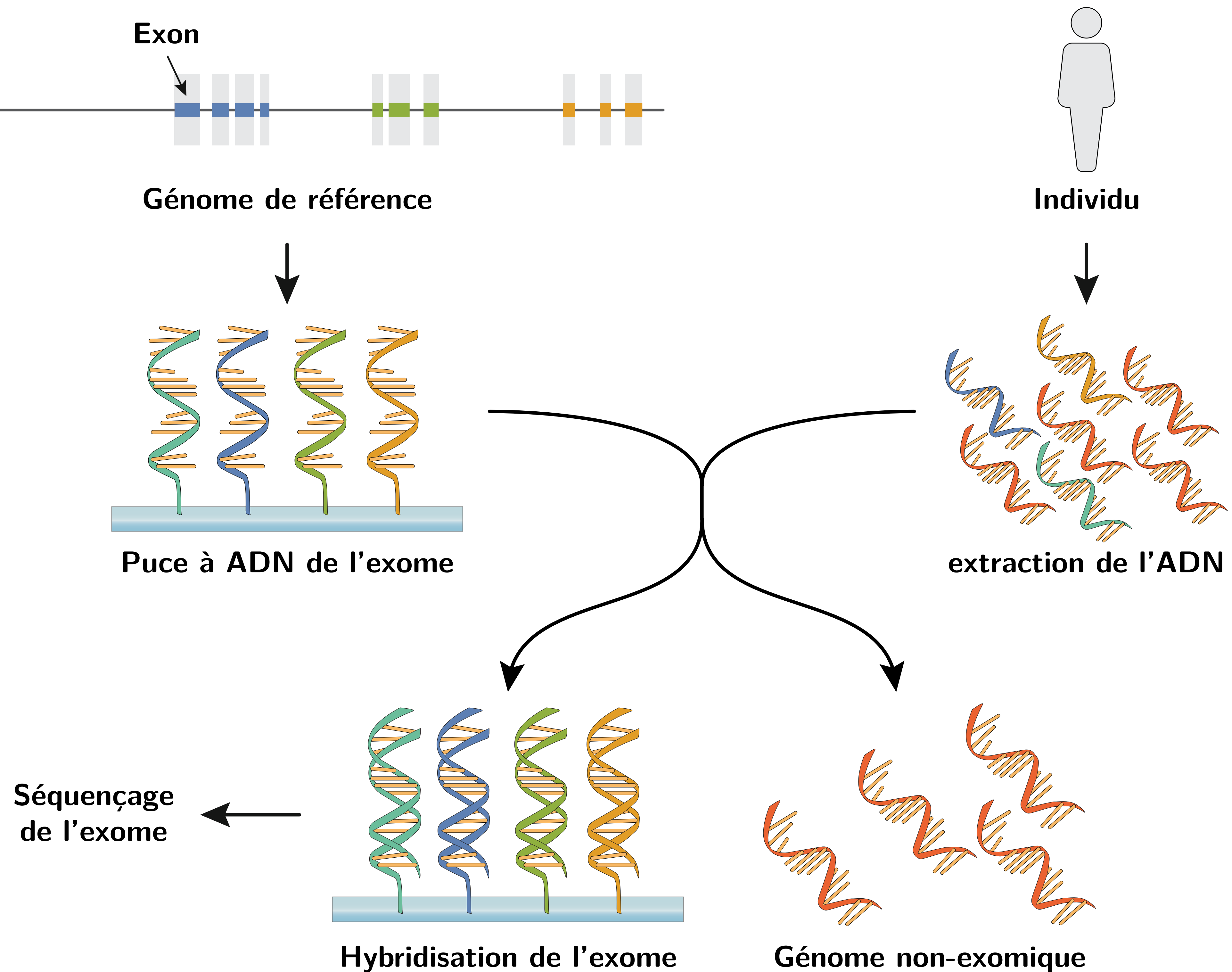
The 1000 Genomes Project Consortium, Nature (2015)

Va-t-on utiliser tout le jeu de données ?



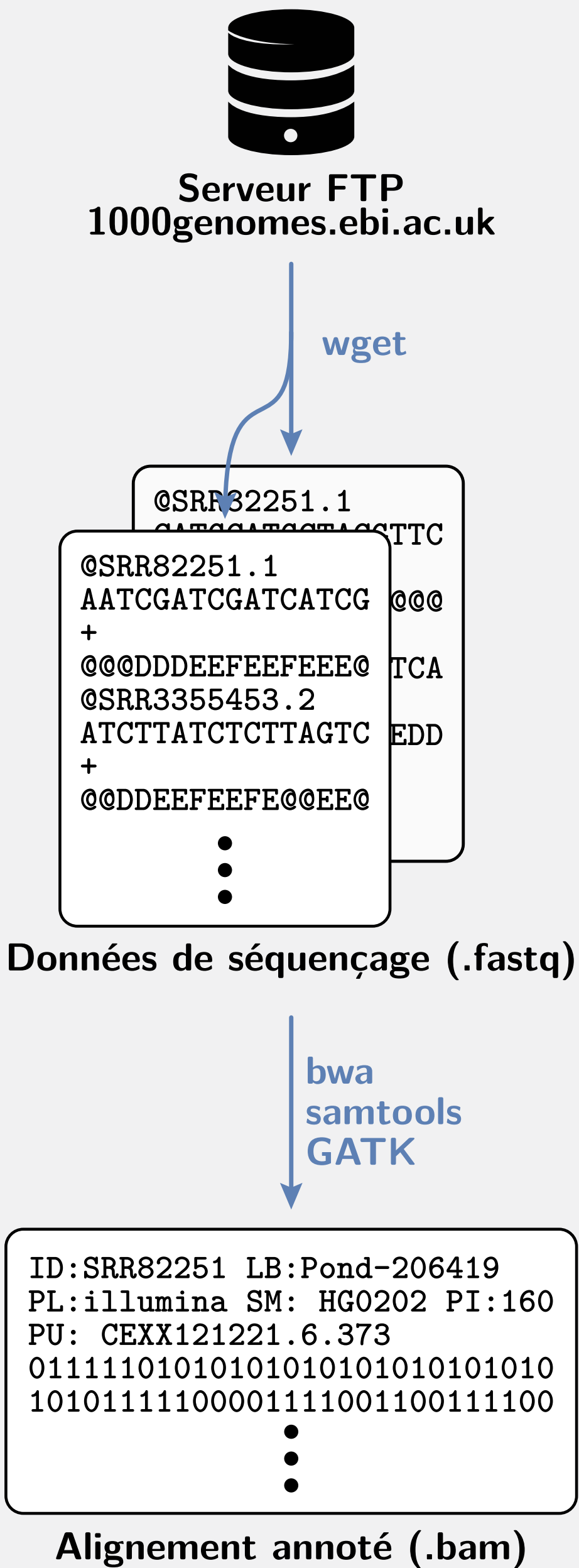
Séquençage de l'exome pour un
trio père-mère-enfant

Comment est obtenu l'exome ?

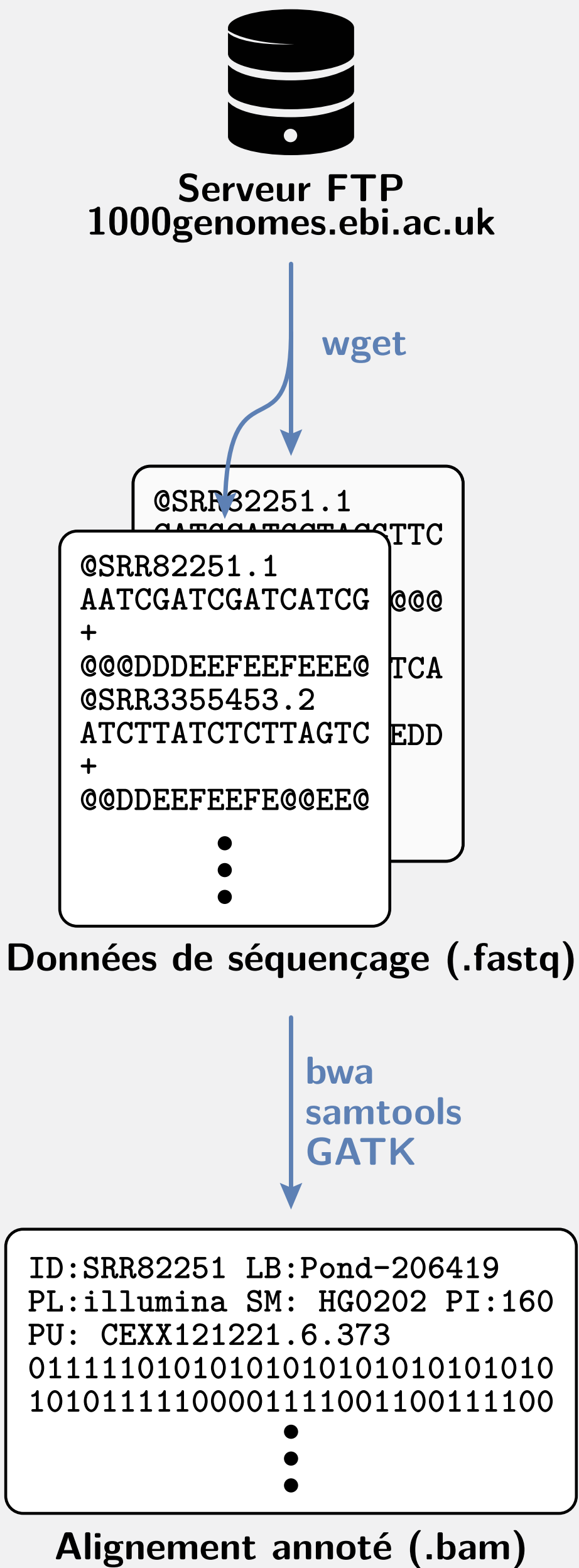


Donc on fait trois fois la même manipulation ?

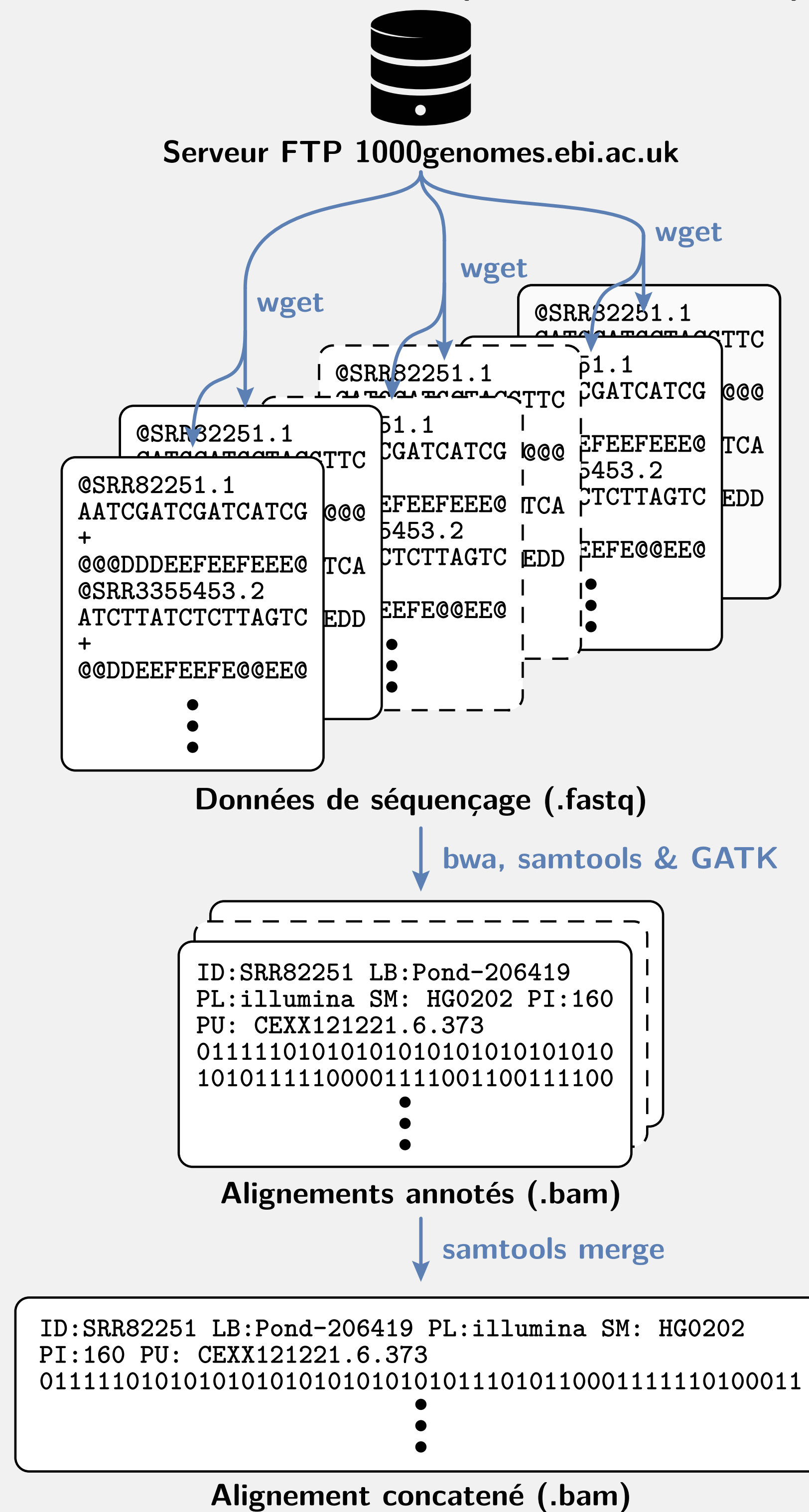
Niveau 1 - Enfant



Niveau 2 - Mère (avec variables)



Niveau 3 - Père (avec boucle)



Quels outils va-t-on utiliser ?

