# JCAD 2018

**Utilisation des ressources CIMENT
dans le cadre du projet epimed:**

**Les besoins spécifiques de la bioinformatique pour
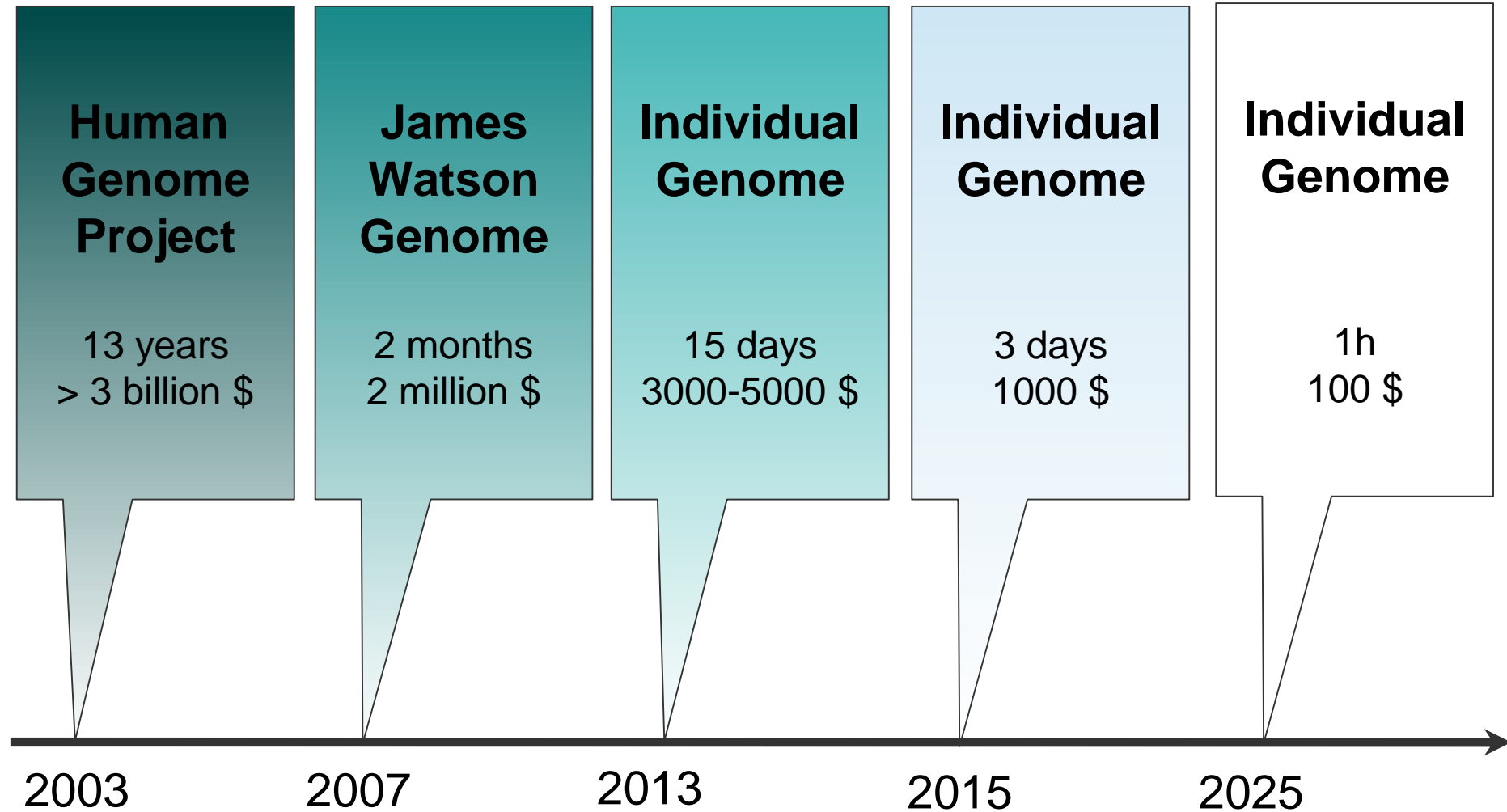l'analyse des données d'exome.**

**Quentin Testard**, Julien Thevenon, Laure Raymond, Jean-François Taly
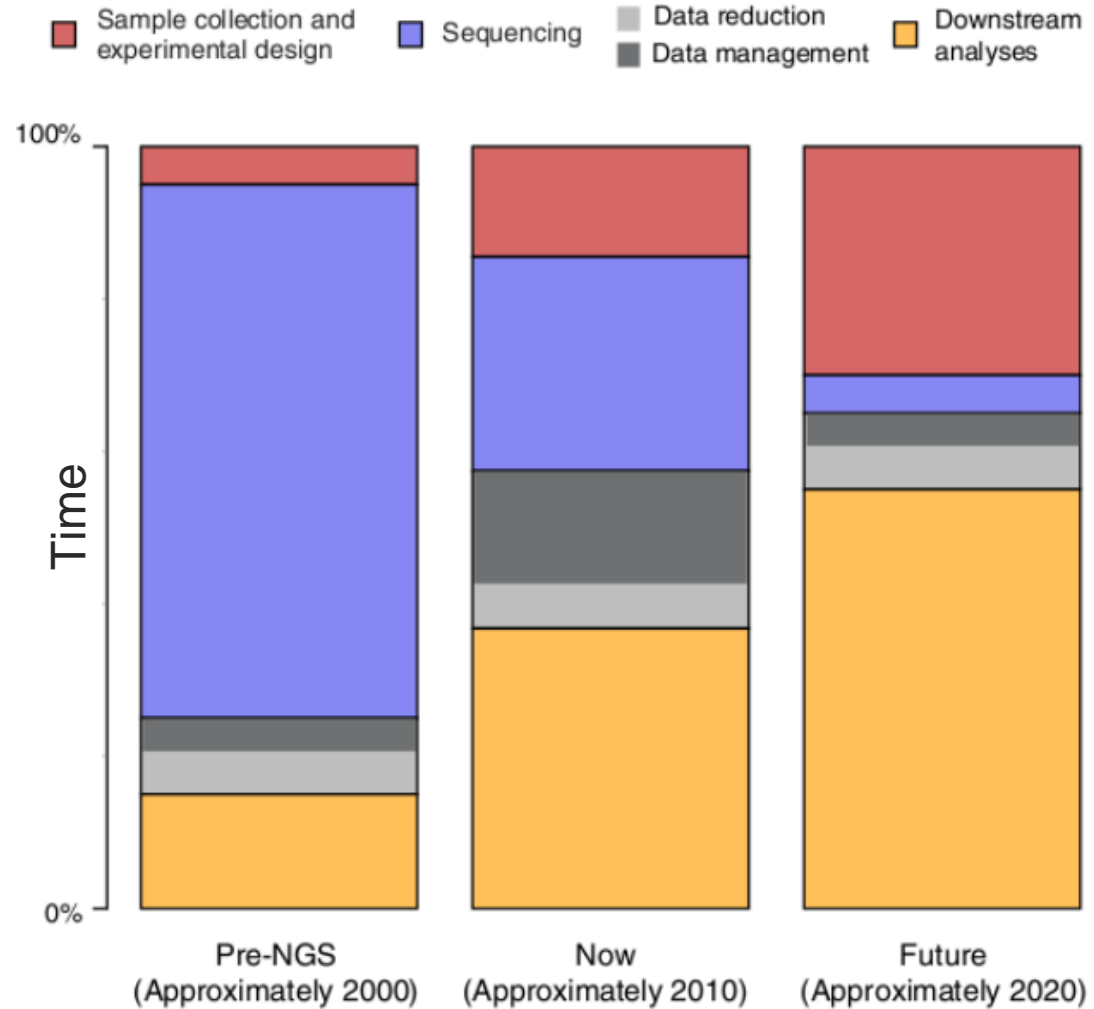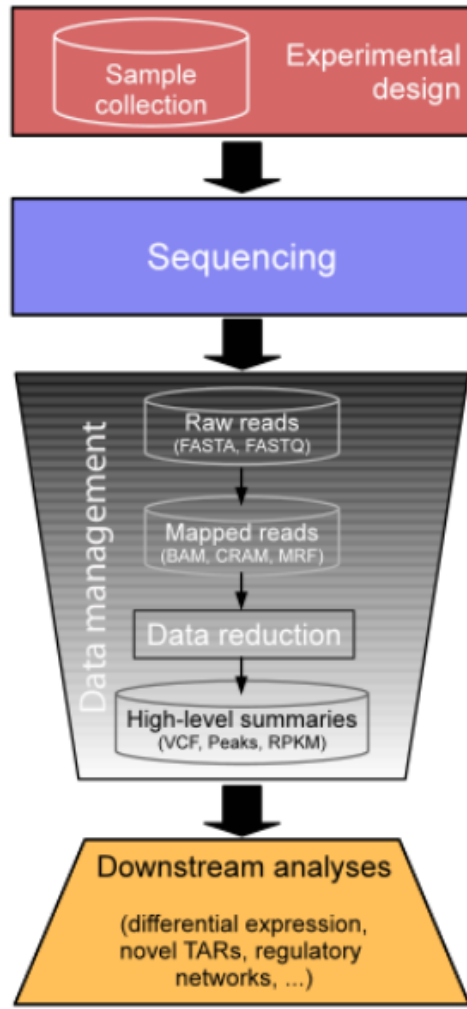
*25/10/18*

# My thesis objectives

- **Industrialized** processes for exome analysis $\Rightarrow$ **Workflow**

- **Speed up** / use **novels approaches** for exome analysis

$\Rightarrow$ **Containers**

- **Master** pipeline execution

- **Control** pipeline possible failures $\Rightarrow$ **Pipeline certification**

- **Will be used in routine health diagnosis activity**

# DNA sequencing

| Human Genome Project | James Watson Genome | Individual Genome | Individual Genome | Individual Genome |
|---|---|---|---|---|
| 13 years > 3 billion $ | 2 months 2 million $ | 15 days 3000-5000 $ | 3 days 1000 $ | 1h 100 $ |
| 2003 | 2007 | 2013 | 2015 | 2025 |

# Bioinformatics analysis

# The Epimed initiative

*Medical Epigenetics and Bioinformatics*

- **Objectives :**

**Facilitate a translational research** in epigenetics, between the fundamental research teams and the medical teams.

**Analyze large-scale whole-genome data** that is essential to understand the epigenome.

Organize an interactive **database** associating **"omics" data** with **biological and clinical data.**

**Community**

**Collaborative network** of about **50 people**:
- IAB, CHU, TIMC-IMAG, LJK, CEA
- International collaborations

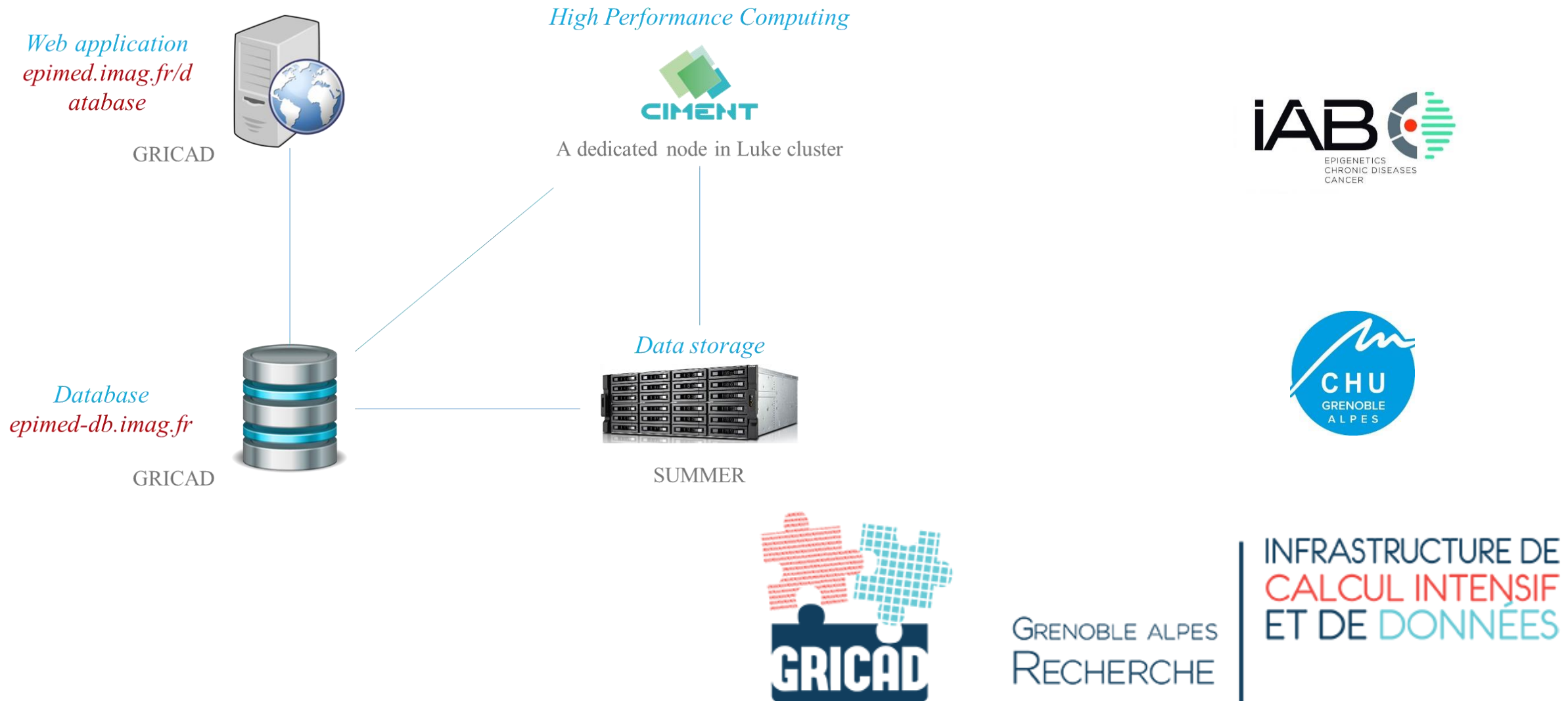Sophie Rousseaux (DR)
epigenetics,
medical research

Anne-Laure Vitte (IE)
molecular biology

Florent Chuffart (IR)
statistics, computing

Ekaterina Flin (IR)
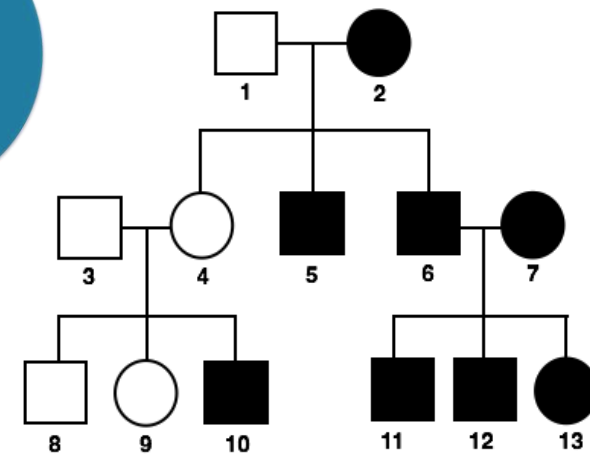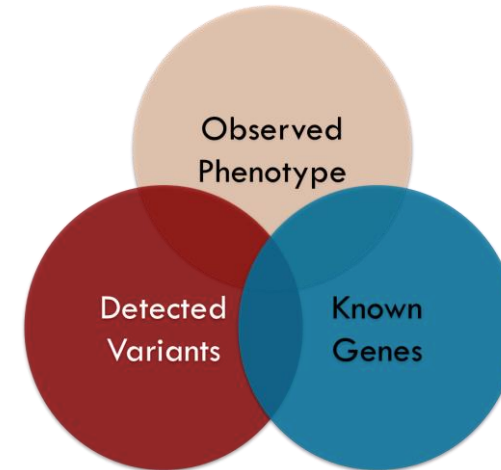databases, web systems

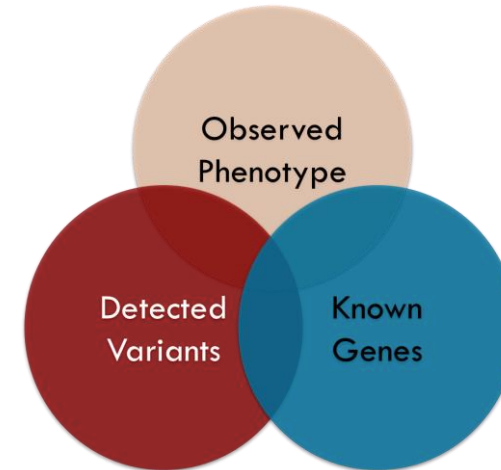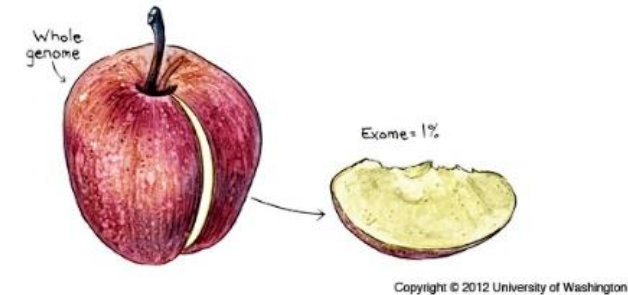# The Epimed infrastructure

*High Performance Computing*

*Web application*
*epimed.imag.fr/d atabase*

GRICAD

CIMENT

A dedicated node in Luke cluster

iAB
EPIGENETICS
CHRONIC DISEASES
CANCER

*Data storage*

*Database*
*epimed-db.imag.fr*

GRICAD

SUMMER

CHU
GRENOBLE
ALPES

GRICAD

GRENOBLE ALPES
RECHERCHE

INFRASTRUCTURE DE
CALCUL INTENSIF
ET DE DONNÉES
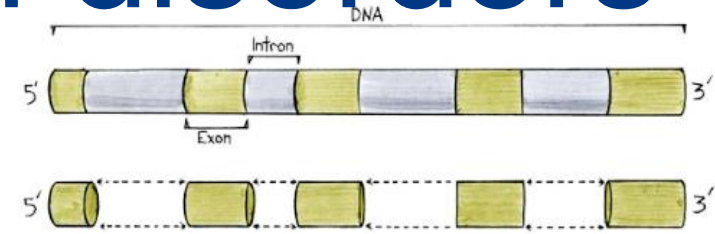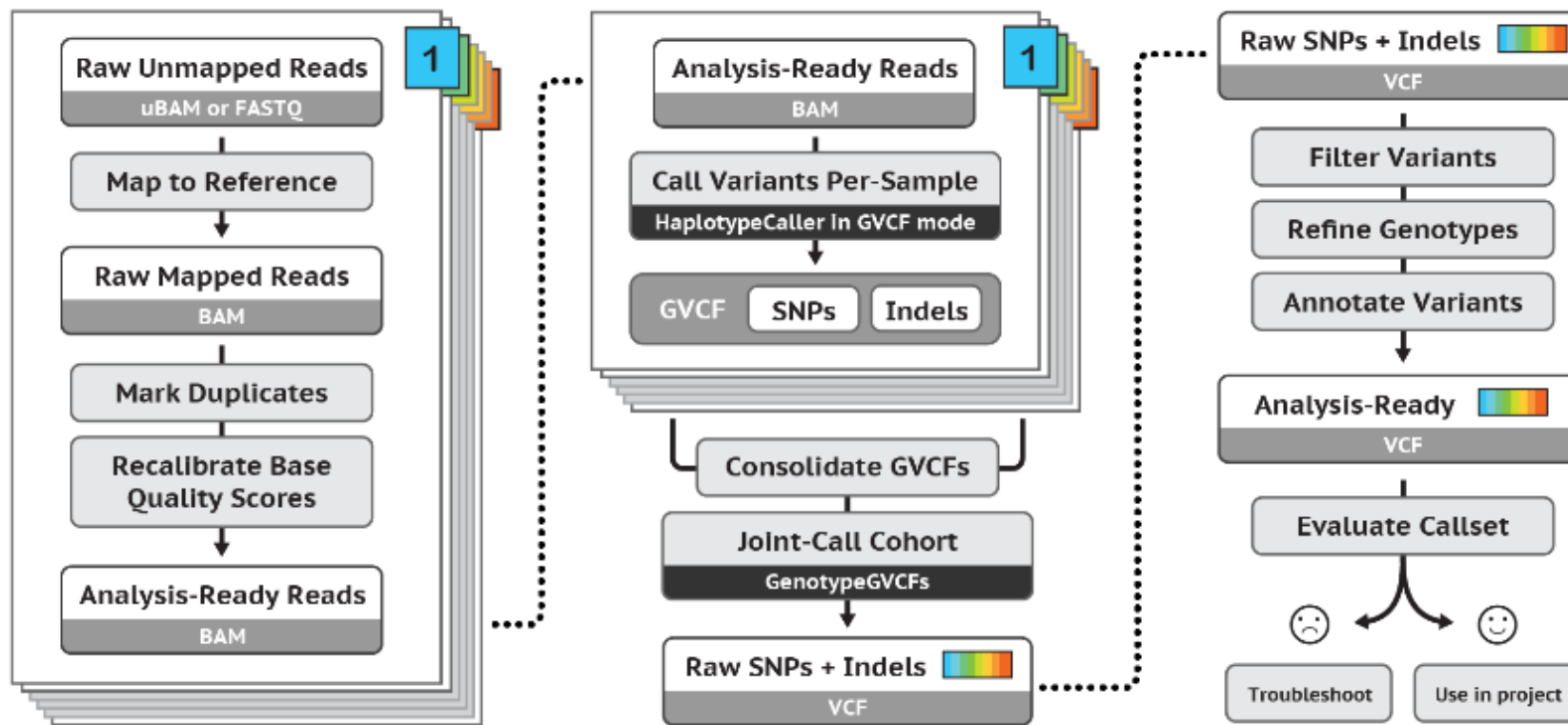
# Exome and rare Mendelian disorders

- Mostly **genetic related**

- Mostly **monogenic/mendelian diseases**

- Quite **difficult** to **diagnose** :
  - Often **long** to **diagnose** => 25 % between 5 to 25 years with **traditional methods**

- **Severe :**
  - Child **mortality rates** around **30%** under **5 y/o**

- **Numerous :**
  - Around **7000 different diseases**
  - **7-8% of total worlwide population (around 30M people in Europe)**

- **Exome analysis :**
  - **Compare sample DNA material to a human « reference »**
  - Find **discordance** between both
  - Relate difference to **patient clinical features** (**phenotype**)
  - Resolve 40% of undiagnosed case

# Bioinformatics analysis

**From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline.**

Van der Auwera GA[1], Carneiro MO[1], Hartl C[1], Poplin R[1], Del Angel G[1], Levy-Moonshine A[1], Jordan T[1], Shakir K[1], Roazen D[1], Thibault J[1], Banks E[1], Garimella KV[2], Altshuler D[1], Gabriel S[1], DePristo MA[1].

https://software.broadinstitute.org/gatk/

# « Required » ressources

- **Cluster Dahu CIMENT :**
  - 72 Dell C6420 (2304 cores)
  - 2x Intel Skylake Gold 6130
  - 192 GB RAM
  - SSD 240GB + SSD 446GB + HDD 4TB

| Analysis step | Input | Time | RAM | CPU | Size output | Read disk |
|---|---|---|---|---|---|---|
| Mapping | ≈ 2 * 4-5 Gb | ≈ 50 – 60 min | ≈ 24 Gb | ≈ 16000% | ≈ 8 - 10 Gb | ≈ 30 - 35 Gb |
| Mark Duplicates | ≈ 8 - 10 Gb | ≈ 20 - 25 min | ≈ 36 - 38 Gb | ≈ 10000% | ≈ 10 - 12 Gb | ≈ 20 - 21 Gb |
| Base Recalibration 1 | ≈ 10 - 12 Gb | ≈ 55 - 60 min | ≈ 36 - 38 Gb | ≈ 10000% | ≈ 18 - 20 Gb | ≈ 45 - 50 Gb |
| Base Recalibration 2 | ≈ 18 - 20 Gb | ≈ 60 - 70 min | ≈ 36 - 38 Gb | ≈ 10000% | ≈ 18 - 20 Gb | ≈ 55 - 65 Gb |
| Haplotypecaller | ≈ 18 - 20 Gb | ≈ 30 - 40 min | ≈ 36 - 38 Gb | ≈ 10000% | ≈ 50 - 90 Mb | ≈ 30 - 40 Gb |
| CombineGVCFs | X * 50 - 90 Mb + 0,5 - 3 Gb | ≈ 10 - 40 min | ≈ 36 - 38 Gb | ≈ 10000% | 24 * 0,5 - 3 Gb | ≈ 10 - 40 Gb |
| GenotypeGVCFs | 24 * 0,5 - 3 Gb | ≈ 1 - 30 min | ≈ 36 - 38 Gb | ≈ 10000% | X * 4 - 5 Mb | ≈ 1 - 12 Gb |
| Annotation | X * 4 - 5 Mb | ≈ 18 - 20 min | ≈ 48 - 50 Gb | ≈ 8000% | X * 30 - 40 Mb | ≈ 14 - 20 Gb |

# Used tools

Nat Biotechnol. 2017 Apr 11;35(4):316-319. doi: 10.1038/nbt.3820.

## Nextflow enables reproducible computational workflows.

Di Tommaso P[1], Chatzou M[1,2], Floden EW[1,2], Barja PP[1,2], Palumbo E[1], Notredame C[1].

Genome Res. 2010 Sep;20(9):1297-303. doi: 10.1101/gr.107524.110. Epub 2010 Jul 19.

## The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.

McKenna A[1], Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA.
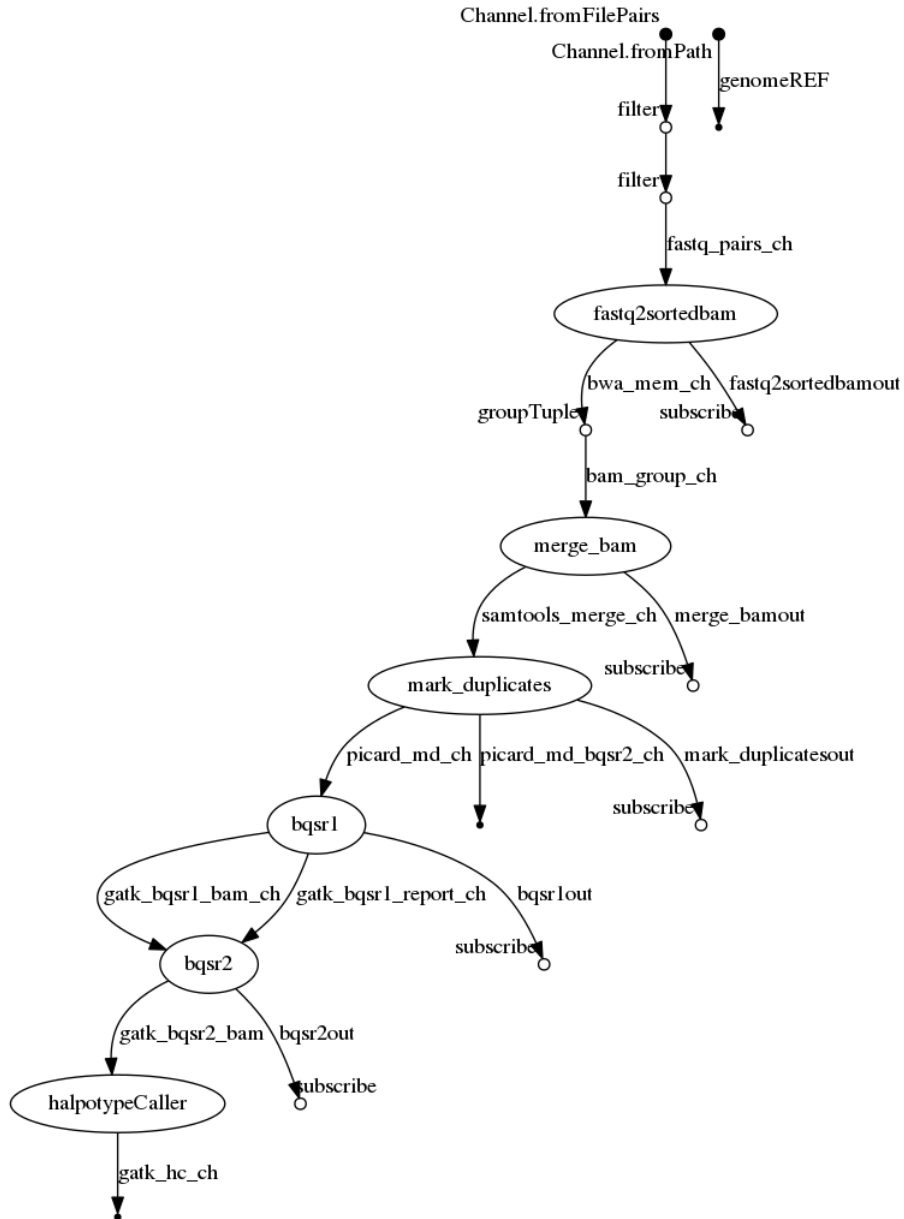
Singularity

docker

# nextflow

- https://www.nextflow.io/

- **Dataflow manager**

- « **Crash** » **management** and « **resume** »fonction

- **Automatic pipeline launching**

- « **Compatible** » with **OAR**, CIMENT batch scheduler

- **Summary reports**



process

```
process fastqc {
  publishDir "${params.resultDir}/fastqc", mode: 'copy'

  input:
  file fastq from fastq_list

  output:
  file '*.zip' into fastq_out_zip
  file '*.html' into fastq_out_html
  val "${params.resultDir}" into fastqc_rep

  """
  fastqc $fastq
  """
}

$fastqc {container = "${params.singularityDir}/fastqc-0.11.15.img"}
```
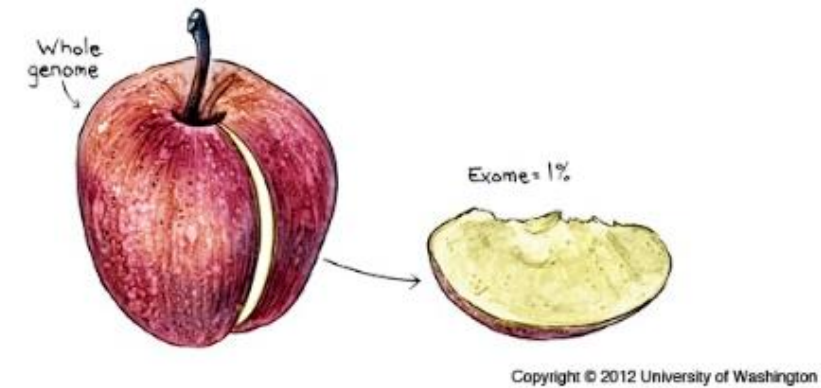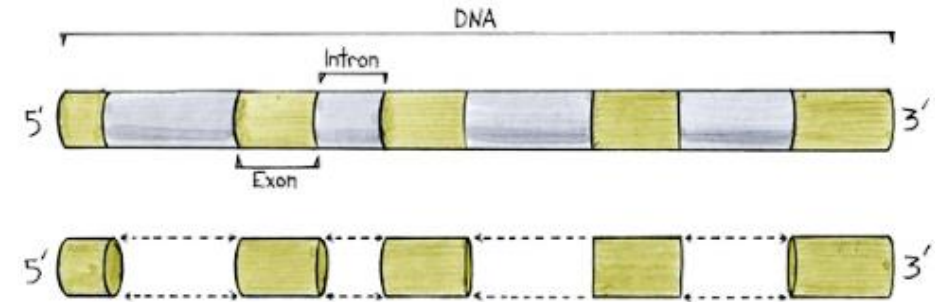
channels

containers

# My « future » pipeline

- **Robust** (managed to process 387 samples during a single week-end)

- Every pipeline version will be « **Vmized** » with all its dependencies (reproductibility for each version)

- **Fully fonctionnal** and **documented exome version** coming **Q1 2019**
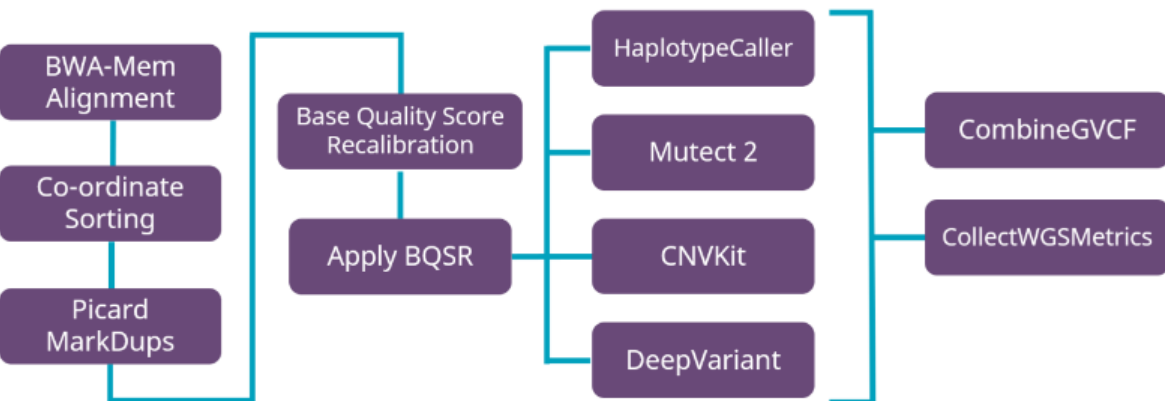
- **Coming next, genome version !**

# And soon … The Human Genome !

- We need **novel approaches** to **speed up** the analysis :

  - **Splitting** operation per **logical unit** (chromosomes) :

    => **Reformatting** my **code**

  - Speeding up by **massively parallelizing processes** :

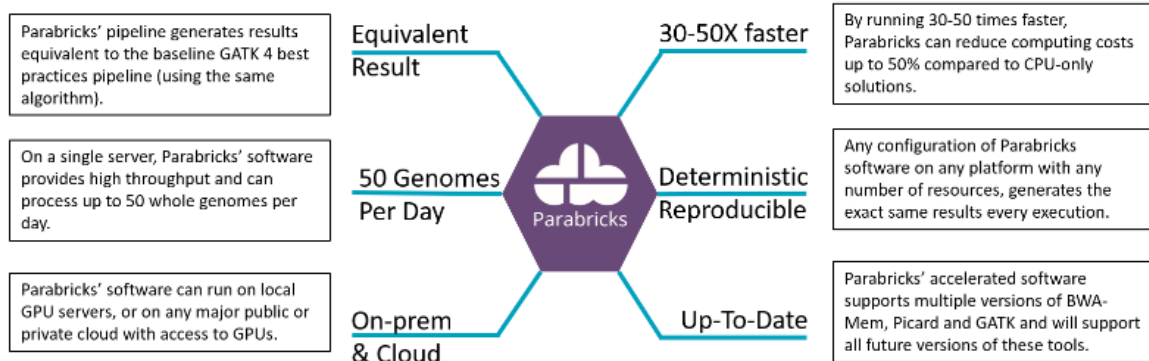    => **GPU programming**

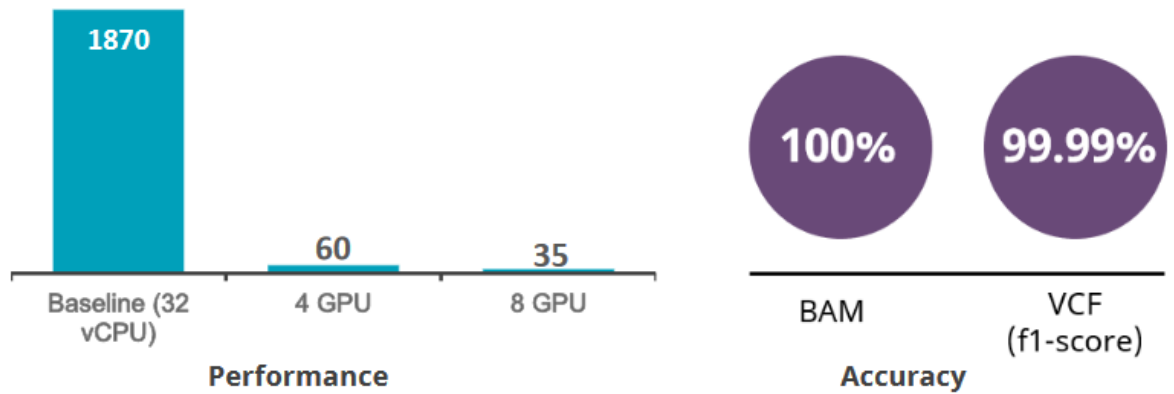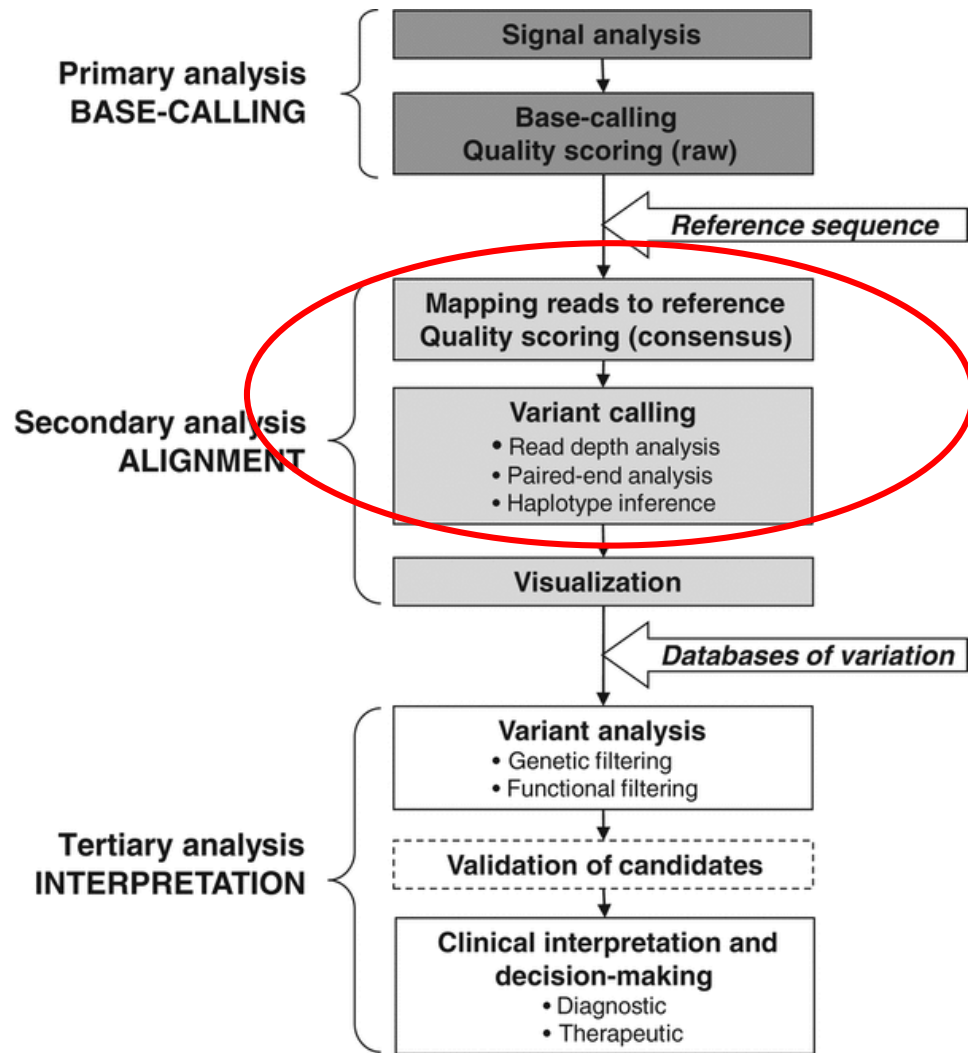    => **MapReduce Spark**

**Tools already exist !**



DNA

Intron

5' 3'

Exon

5' 3'

Whole genome

Exome = 1%

Copyright © 2012 University of Washington

# GPU
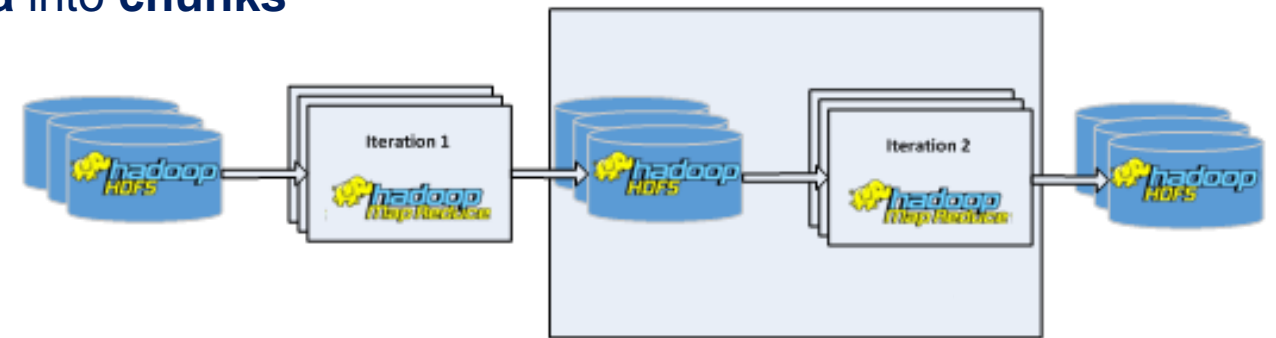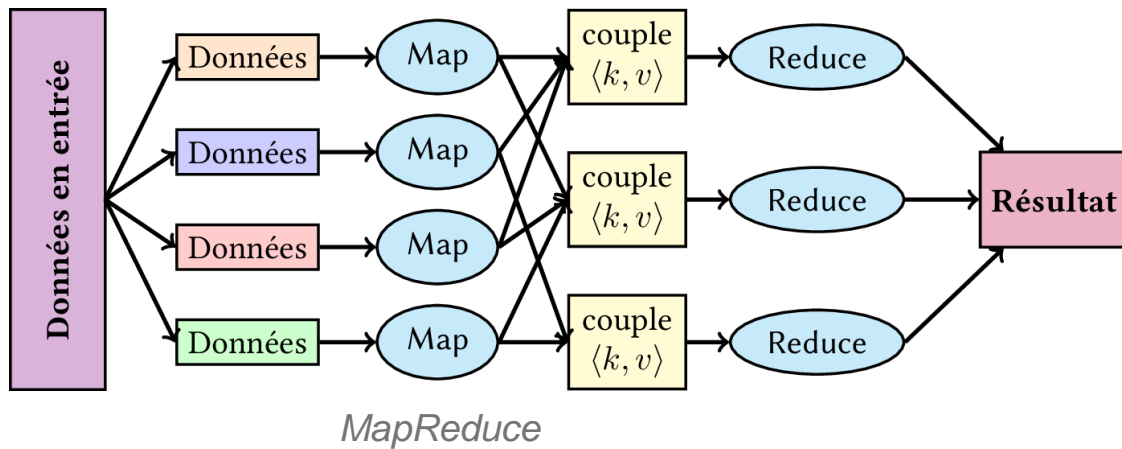
# GATK4 and Spark



**Native Spark tools** are **included** into **GATK4** for **mapping** and **variant calling** (two of the most paralellizable steps)
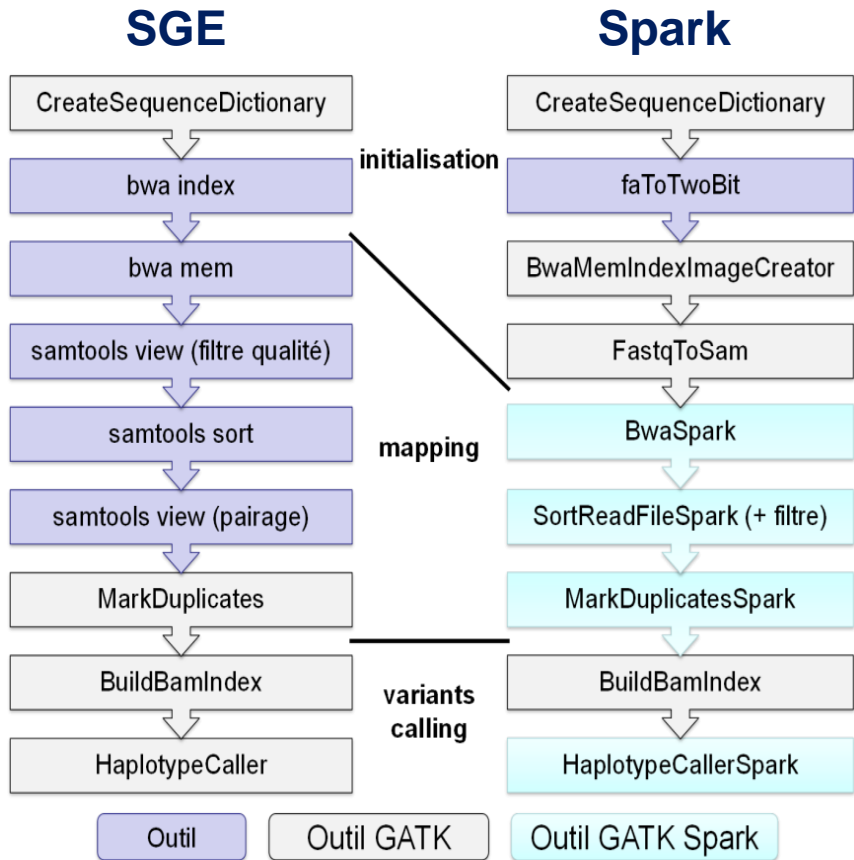
# Why Spark?

- Permit to **massively parralelize** processes

- Spark Bioinformatics tools « cut » **biological data** into **chunks**

and do « **bioinformatic** » **operations** on it



*MapReduce*

*Hadoop MapReduce and Spark*

# GATK4 Spark benchmarking



**SGE**

- CreateSequenceDictionary
- bwa index
- bwa mem
- samtools view (filtre qualité)
- samtools sort
- samtools view (pairage)
- MarkDuplicates
- BuildBamIndex
- HaplotypeCaller

*initialisation*

*mapping*

*variants calling*

**Spark**

- CreateSequenceDictionary
- faToTwoBit
- BwaMemIndexImageCreator
- FastqToSam
- BwaSpark
- SortReadFileSpark (+ filtre)
- MarkDuplicatesSpark
- BuildBamIndex
- HaplotypeCallerSpark

Outil | Outil GATK | Outil GATK Spark

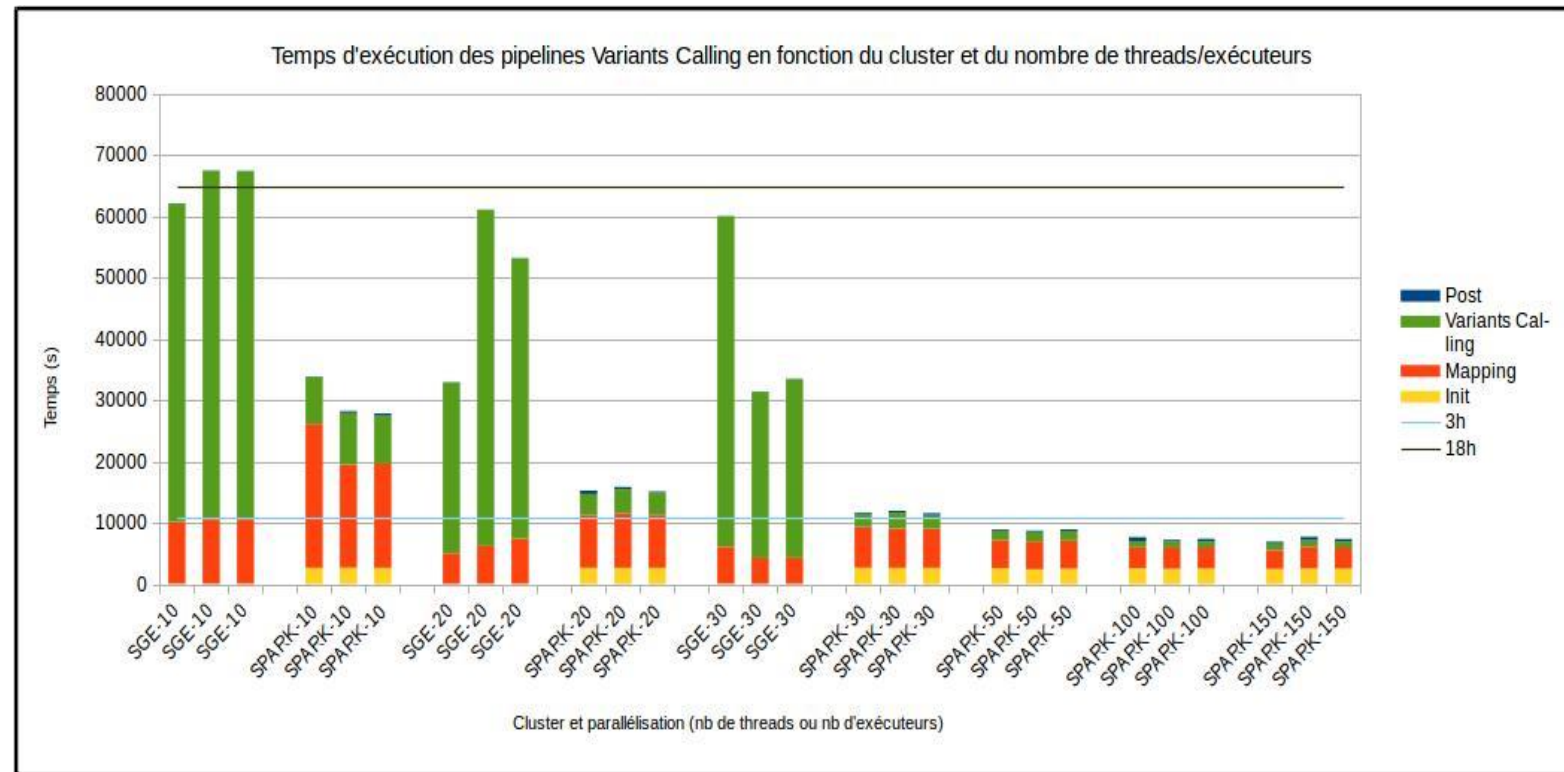- **Toulouse LIPM Cluster** :

  - 8 physical nodes

  - 252 cores

  - 642 Go RAM memory

  - 34To HDD

- One supplementary frontal node to **distribute Spark jobs**

- The **same cluster** have been used for **Spark** and **SGE jobs**, but **not at the same time**

# Spark benchmarking

- The **GATK4 non-Spark tools** are **poorly parallelized**

- The **GATK4 Spark tools** are still in **Beta**

- There is a **significative speed improvement** with **Spark tools**



*Genome : Plasmopara halstedii (75 Mbases), fastq.gz : 2x12Go, 3 repeats*

*Source : Axel Verdier INRA UMR LIPM Toulouse*

# Conclusion

- Currently **CIMENT luke/Dahu** HPC nodes are **sufficient** for our **exome analyses**

- But for **high-scale genome analysis** we will need to **step up** :

    - Produce a **better code**

    - Use **bigger HPC infrastructures**

    - Use **Spark** or **GPU** programmed **tools** and **infrastructures**

- It's not that hard ! **Tools already exist** and are indeed **speeding up genome analyses**

- We need to **proceed** with **caution**. We are working with **health datasets**. **No mistake allowed**.

# Thank you!