# Sequencing technologies

Principal technologies:

## 454 Life Sciences/Roche
Reads size: 0.5-1kb
Reads nb: ~$10^6$
Total seq: 0.7 Gb
https://en.wikipedia.org/wiki/DNA_sequencer
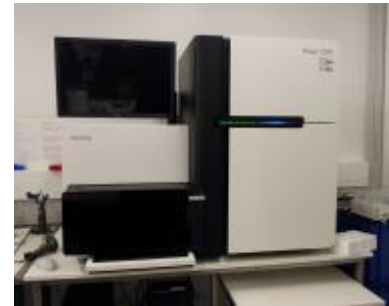
## Illumina
Reads size: 2*150
Reads nb: ~$6*10^9$-$20*10^9$
Total seq: 600Gb
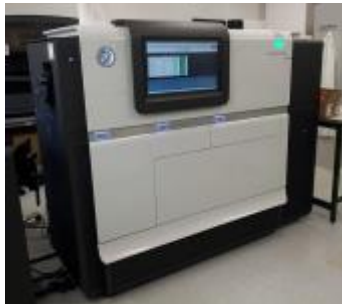https://emea.illumina.com/systems/sequencing-platforms.html

## PacBio
Reads size: 30 kb
Total seq: 20Gb
https://www.pacb.com/products-and-services/sequel-system/

## Oxford nanopore
Reads size: 30 kb
Total seq: 15Tb
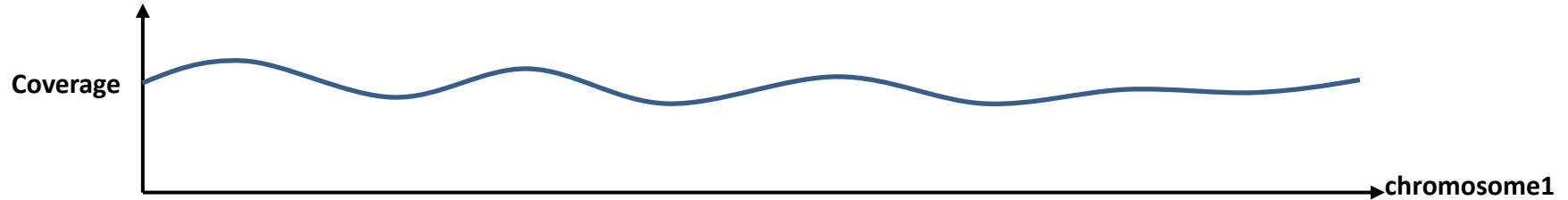https://nanoporetech.com/products/promethion

# From the output of sequencing to the variant calling file
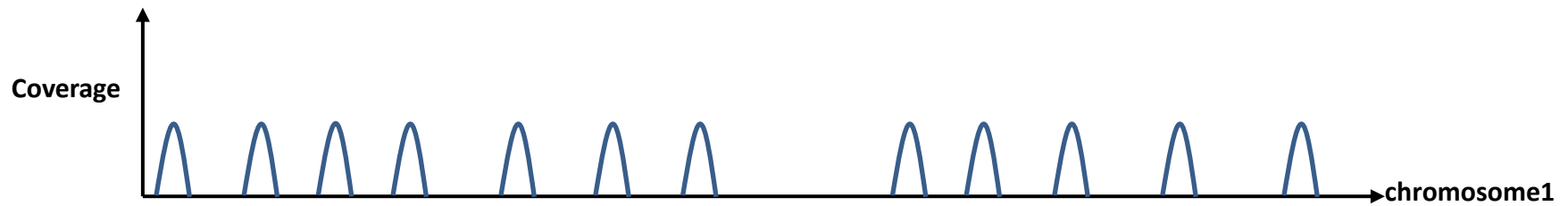
Principles types of sequencing:

- **W**hole **G**enome **S**equencing (WGS): All the genome is "uniformly" sampled (some biases exist depending to sequencing technologies).

Coverage

chromosome1

- messenger **RNA seq**uencing (RNAseq): mRNA are sequenced after a step of cDNA complementation

Coverage

chromosome1

- **G**enotyping **B**y **S**equencing (GBS): The genome is sampled and only part of it is sequenced.

Coverage

chromosome1

# From the output of sequencing to the variant calling file

Standard workflow:

Sequences obtained from the sequencer = reads. (generally short: 100-250 bases

**Alignment against a reference sequence**

**Raw alignment**

reference

**Post alignment processing steps**

**Cured alignment**

reference

**Polymorphism identification**

```
ACAGGTGTCCACTGACTTTGCAAC      AAGCTTCCGTACTGTACCT
ATTGGACAGGTGTCCACTGAC    TGCAACTCCAAGGTTCCGTACT
ATGCATTGGACAGGTGTCCAC    ACTTTGCAACTCCAAGCTTCCGTA
ATGCATTGGACTGGTGTCCACTGACTTTGCAACTCCAAGGTTCCGTACT  reference
```

**Variant calling**

| CHROM | POS | Genotype |
|-------|-----|----------|
| chr01 | 12  | A/A      |
| chr01 | 40  | C/G      |

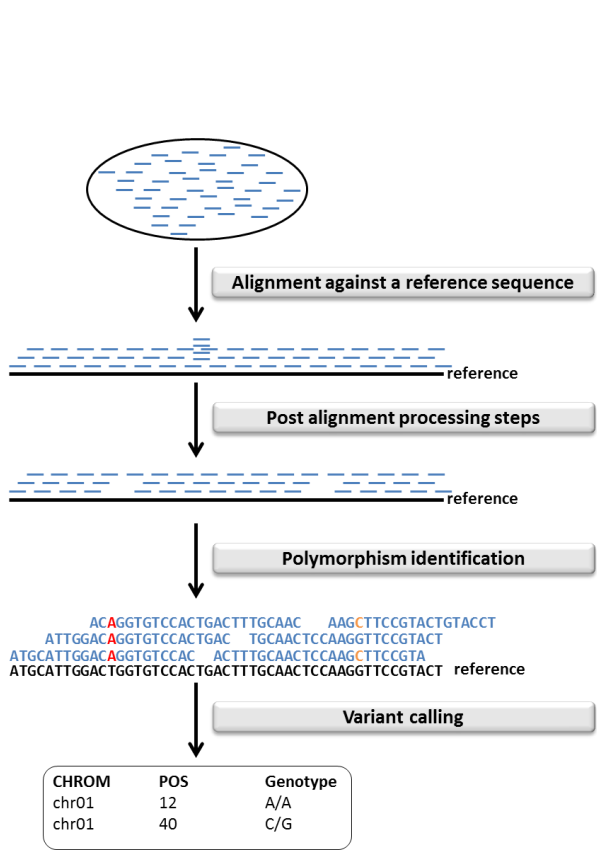# From the output of sequencing to the variant calling file

**Depending on the sequencing technologies: steps from the sequencing data to the variant calling format are distinct**
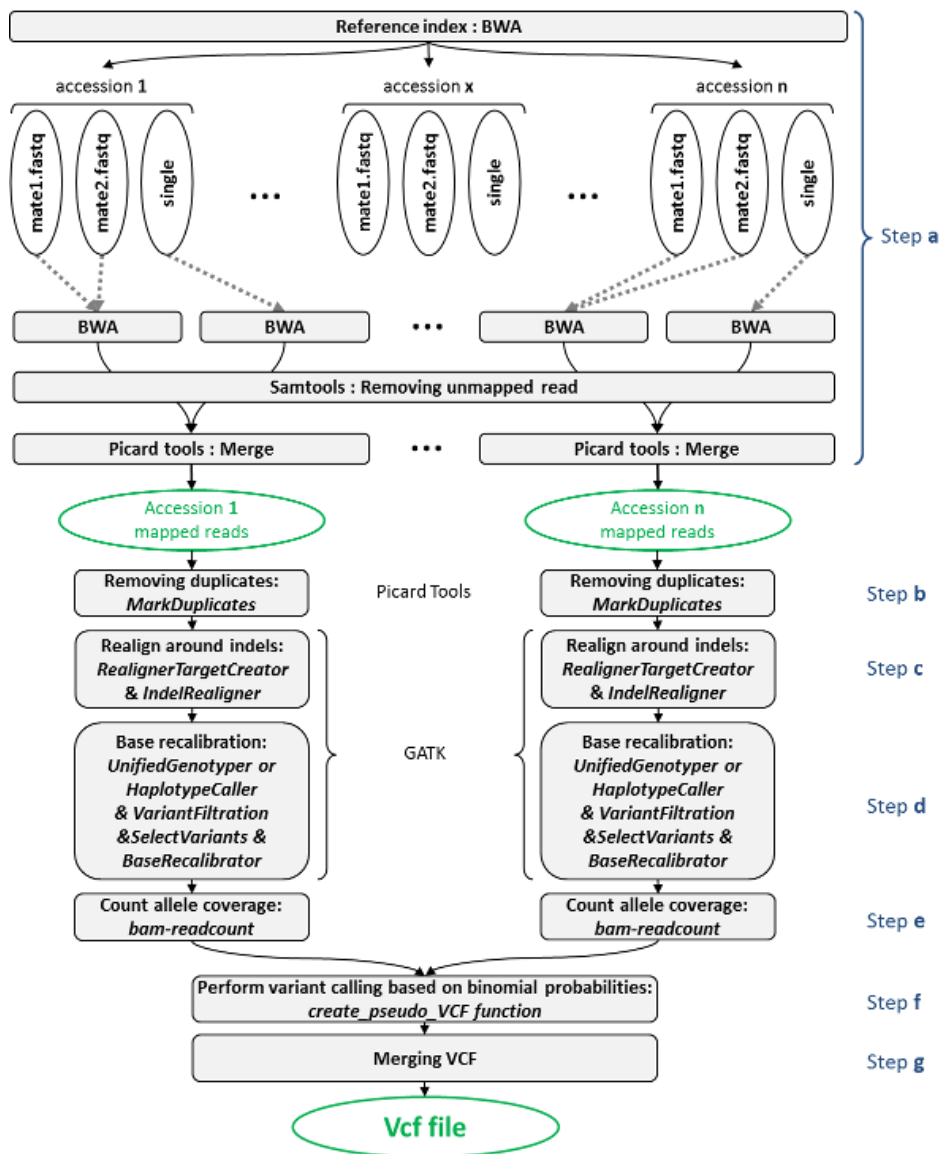


- RNAseq: Aligner should take into account mRNA splicing.

- PCR duplicates are usually removed because they biased allelic ratio. It is not possible for GBS du to the approach... (see latter)

- RNAseq: Read overlapping splicing sites should be split.

# From the output of sequencing to the variant calling file

**Depending on the sequencing technologies: steps from the sequencing data to the variant calling format are distinct**



- RNAseq: Aligner should take into account mRNA splicing.

- PCR duplicates are usually removed because they biased allelic ratio. It is not possible for GBS du to the approach... (see latter)

- RNAseq: Read overlapping splicing sites should be split.

**Several workflow exists:**
- TOGGLe: https://github.com/SouthGreenPlatform/TOGGLE

- GATK best practice: https://software.broadinstitute.org/gatk/best-practices/

- VcfHunter: https://github.com/SouthGreenPlatform/VcfHunter

VcfHunter detailed workflow (Developped under GenomeHarvest) for WGS and GBS:



Step b ← For WGS only

Possible but not recommended:
- High computation time
- Result not so good

# The Genotyping By Sequencing in detail

- Principle: sequencing a constant part of the genome in several accessions

- Why?
  - ✓ The amount of reads obtained per sequencing run is constant
  - ✓ Necessity to have enough coverage to have a confident genotype calling
  - ✓ Several accessions can be sequenced in one run

> **Sequencing a sample of the genome which is a constant part ➔ allow to sequence more accessions in a run and to keep the same coverage**

**WGS** 54 reads
~ 3x

reference

**GBS** 54 reads
~ 9x

reference

**GBS** 54 reads
~ 3x/accessions

reference

- Cutting the genome with restriction enzymes
- Selection of "short" fragments (<500)
- Sequencing of extremities of selected fragments
- Relative constant sampling of regions in distinct samples (exception if mutation in restriction sites)
- Single or combination of restriction enzyme(s)

*pstI*        *mseI*

5'...TCCTCTTACAGGATCCTGCAGCAACAAGGGTTAAGAATTATAAGCA...3'
3'...AGGAGAATGTCCTAGGACGTCGTTGTTCCCAATTCTTAATATTCGT...5'

**Enzymatic restriction**

**DNA insert**
5'...        GCAACAAGGGT        ...3'
3'...    ACGTCGTTGTTCCCAAT    ...5'

+

**Barcode adapter**                                                **Common adapter**

5'-ACACTCTTTCCTACACGACGCTCTTCCGATCTXXXXTGCA        TACAGATCGGAAGAGCGGTTCAGCAGGAATGCCGA-3'
3'-TGTGAGAAAGGGATGTGCTGCGAGAAGGTAGAYYYY        GTCTAGCCTTCTCGCCAACTCGTCCTTACGGCT-5'

Barcode (unique to each individual)

**Ligation**

5'-ACACTCTTTCCTACACGACGCTCTTCCGATCTXXXXTGCAGCAACAAGGGTTACAGATCGGAAGAGCGGTTCAGCAGGAATGCCGA-3'
3'-TGTGAGAAAGGGATGTGCTGCGAGAAGGTAGAYYYYACGTCGTTGTTCCCAATGTCTAGCCTTCTCGCCAACTCGTCCTTACGGCT-5'

**5' -> 3' fragment selection for sequencing**

Illumina Primer

5'-ACACTCTTTCCTACACGACGCTCTTCCGATCTXXXXTGCAGCAACAAGGGTTACAGATCGGAAGAGCGGTTCAGCAGGAATGCCGA-3'

**Sequencing**

Barcode  Restriction site  Sequence  Restriction site  Common adapter

5'-XXXXTGCAGCAACAAGGGTTACAGATCGGAAGAGCGGTTCAGCAGGAATGCCGA-3'

- We have generated a small GBS dataset comprising 12 samples for which *pstI* and *mseI* enzymes have been used and a sample specific barcode have been used.

| | |
|---|---|
| Sample1 | + AACT |
| Sample2 | + CCAG |
| Sample3 | + TTGA |
| Sample4 | + GGTA |
| Sample5 | + ATTG |
| Sample6 | + CGGT |
| Sample7 | + TGCG |
| Sample8 | + GTAT |
| Sample9 | + AACCA |
| Sample10 | + CCACG |
| Sample11 | + TATAA |
| Sample12 | + GAGCG |

Sequenced

Unique fastq file containing reads from all accessions

South Green bioinformatics platform

GenomeHarvest diversity, organization and dynamics

# From the output of GBS sequencing to the variant calling file: in command line

- Obtaining the datasets:
    1. Log onto the cluster

    2. Go to your "work" directory

        **c**hange **d**irectory (with nothing else: it goes to your home - /home/Your_ID)

        `cd`
        `cd work` ← **c**hange **d**irectory to work

        Double clic

    3. Create a folder in which we will be working and go into:

        **m**ake **di**rectory vcfhunterGBS

        `mkdir vcfhunterGBS`
        `cd vcfhunterGBS`

    4. Copy the folder containing the sequencing information:

        `cp -R /home/gmartin/WorkShop/VCFHUNTER/data/WorkShopDataset .`

        Copy

        The folder and all it contains

        Location of the folder

        Copy it here
        *i.e*
        /home/Your_ID/work/vcfhunterGBS

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the datasets:

  `ll WorkShopDataset`



① A compressed file (.gz) containing all reads from all accessions obtained from the sequencer

- To have a look at this file

  `zmore WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz`

  ↑                                                          ↑

  Reading line by line a zipped file (.gz)        Path to the fastq file

- Because zmore will list the file until its end using the "enter" key, and we do not want that because the file is big, we can "kill" the command with a combination of key:

  `"Ctrl" + "C"`

# From the output of GBS sequencing to the variant calling file: in command line

- To have a look at this file
  - `zmore WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz`

**Fatsq format (for each read)**

Read name

Read nucleotides sequence

Separator

Read nucleotides quality



Read1

Read2

Base quality encoding: for base "C" = A

But what does "A" mean?

- Each letter has informatically a numeric value. For example "A" is equal to 65
- We should remove 33 to this value and thus "A" = 65-33 = 32!

```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJ
|                              |   |        |
33                             59  64       73
 0.2......................26...31........41
```

# From the output of GBS sequencing to the variant calling file: in command line

- To have a look at this file
  - `zmore WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz`



Sample8 tag
➔ **Read from sample8**

Sample2 tag
➔ **Read from sample2**

Sample11 tag
➔ **Read from sample11**

pstI restriction site

mseI restriction site

Adapter sequence

## From the output of GBS sequencing to the variant calling file: in command line

- Listing the datasets:

  `ll WorkShopDataset`



① A compressed file (.gz) containing all reads from all accessions obtained from the sequencer

② **A file that will be used to separate reads in distinct file according to the accession they belong**

- To have a look at this file

  `more WorkShopDataset/DemultiplexingFile.tab`

Reading line by line a file                    Path to the file

# From the output of GBS sequencing to the variant calling file: in command line

- To have a look at this file

  `more WorkShopDataset/DemultiplexingFile.tab`

Sample name

Sample tag

Restriction enzyme1

Restriction enzyme2



```
Multi PuTTY Manager
File   View   Tools   Help
       Import Database      Close All Sessions
Protocol  SSH          Host                 Login as              Password
Multi Sessions Command                                               Sessions
  cc2-gmartin   cc2-gmartin   cc2-gmartin   cc2-gmartin
[gmartin@cc2-login vcfhunterGBS]$ more WorkShopDataset/DemultiplexingFile.tab
sample1 AACT     PstI    MseI
sample2 CCAG     PstI    MseI
sample3 TTGA     PstI    MseI
sample4 GGTA     PstI    MseI
sample5 ATTG     PstI    MseI
sample6 CGGT     PstI    MseI
sample7 TGCG     PstI    MseI
sample8 GTAT     PstI    MseI
sample9 AACCA    PstI    MseI
sample10         CCACG   PstI    MseI
sample11         TATAA   PstI    MseI
sample12         GAGCG   PstI    MseI
[gmartin@cc2-login vcfhunterGBS]$
```

| Sample1 | + AACT |
| Sample2 | + CCAG |
| Sample3 | + TTGA |
| Sample4 | + GGTA |
| Sample5 | + ATTG |
| Sample6 | + CGGT |
| Sample7 | + TGCG |
| Sample8 | + GTAT |
| Sample9 | + AACCA |
| Sample10 | + CCACG |
| Sample11 | + TATAA |
| Sample12 | + GAGCG |

# From the output of GBS sequencing to the variant calling file: in command line

- Now it is time to demultiplex! *i.e.* parse reads in files corresponding to sample.

- For that we will use GBSX (https://github.com/GenomicsCoreLeuven/GBSX, https://doi.org/10.1186/s12859-015-0514-3)

- A small parenthesis: On the AGAP cluster, several modules are already available. To access the list of available modules, use the following command line:

```
module avail
```

A list of modules appears and we can find "GBSX" program in this list!



➔ Two versions are available! We will take the 1.2 version

## From the output of GBS sequencing to the variant calling file: in command line

- To load this module run the command line:

  `module load bioinfo/GBSX/1.2`

- The module is now loaded. This can be verified by listing the loaded modules with de following command line:

  `module list`



The GBSX module is loaded. But what you don't know, is that GBSX need another program to be used! This program is JAVA. To load java we will run the command line:

`module load system/java/jre8`

You can try again `module list` to verify that java has been loaded

# From the output of GBS sequencing to the variant calling file: in command line

- At this point all is ready to demultiplex the fastq file! All we have to do is to run the following command line (in one single line):

```
qsub -q normal.q -l mem_free=12G -b yes -V -N DEMULT java -XX:ParallelGCThreads=1 -Xmx8G
-jar /usr/local/bioinfo/GBSX/1.2/GBSX_v1.1.2.jar --Demultiplexer
-f1 WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz
-i WorkShopDataset/DemultiplexingFile.tab -o Demultiplexed -gzip true -mb 0
```

- Now a little piece of explanation:
  - ✓ We are working on a cluster.
    - ❖ This means that we have several computers which are connected so that they can work together.
    - ❖ It also allows that several people can run huge calculation at the same time!
    - ❖ It also means that there is a strict procedure to perform calculation on the cluster and this procedure is associated to the way a cluster work:



*A single computer to rule them all*

## How the cluster works?

*User command Line* ①

④

*Result of the command line*

② 

Master Computer1 → ③ → Computer1 / Computer2 ④ / ... / ComputerX

1. The user tip a command line

2. Which is sent to the master computer

3. Based on this command line, the master computer identify which computer it rules match the command requirements and which of them are available

4. The command line is executed on the chosen computer (in this example **Computer2**)

5. Which returns the result of the command line

- Back to the command line:

```
qsub -q normal.q -l mem_free=12G -b yes -V -N DEMULT "java -XX:ParallelGCThreads=1 -Xmx8G
-jar /usr/local/bioinfo/GBSX/1.2/GBSX_v1.1.2.jar --Demultiplexer
-f1 WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz
-i WorkShopDataset/DemultiplexingFile.tab -o Demultiplexed -gzip true -mb 0"
```

- The first part of the command line (in bold) is **used by the master computer**:
  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer
  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue. Several queues exist depending on computation requirement:
    - ✓ **normal.q**: access to computers of 48 processors with 192Go shared memory (RAM) and a command line cannot exceed 48hours of running time.
    - ✓ **long.q**: access to computers of 48 processors with 192Go shared memory but there is not running time limit
    - ✓ **bigmem.q**: access to a unique computer of 96 processors with 2,6To shared memory and no time limit
  - ❖ **-l mem_free=12G**: precise that the program will use 12G of RAM (so the master computer will check that it is available on the computers). This is a facultative option but necessary when using **java** program to prevent errors…
  - ❖ **-b yes**: it is not important, but put it.
  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose
  - ❖ **-N DEMULT**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

# From the output of GBS sequencing to the variant calling file: in command line

- Back to the command line:

```
qsub -q normal.q -l mem_free=12G -b yes -V -N DEMULT "java -XX:ParallelGCThreads=1 -Xmx8G
-jar /usr/local/bioinfo/GBSX/1.2/GBSX_v1.1.2.jar --Demultiplexer
-f1 WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz
-i WorkShopDataset/DemultiplexingFile.tab -o Demultiplexed -gzip true -mb 0"
```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**.
  - ❖ **/usr/local/bioinfo/GBSX/1.2/GBSX_v1.1.2.jar**: is the program that is used to demultiplex the fastq file. Element in black are options/argument passed to this program to make it work (as a function and its arguments in Excel!).
  - ❖ **--Demultiplexer**: Tell the program that we want to demultiplex the fastq
  - ❖ **-f1 WorkShopDataset/ReadFromTheSequencer_R1.fastq.gz**: locate the fastq file to demultiplex
  - ❖ **-i WorkShopDataset/DemultiplexingFile.tab**: loacte the file containing the multiplexing informations (which tags correspond to which samples and restriction enzymes used)
  - ❖ **-o Demultiplexed**: The name of the output folder (this folder will be created by the program).
  - ❖ **-gzip true**: Tells the program that output should be compressed to gain space (equivalent to .zip files on Windows)
  - ❖ **-mb 0**: Tells the program that 0 mismatch are allowed in the tag to attribute a read to an accession

  - ❖ **java -XX:ParallelGCThreads=1 -Xmx8G -jar**: Tells to the computer that the program */usr/local/bioinfo/GBSX/1.2/GBSX_v1.1.2.jar* is written in java language (**java**), that java should only use one processor (**-XX:ParallelGCThreads=1**) and that 8G memory are available for java (**-Xmx8G -jar**). **-jar** indicate to java that the program is directly after.

# From the output of GBS sequencing to the variant calling file: in command line

- One can check the status of job(s) with the following command line:

  `qstat`

- Because the job we have sent is a very short one it is likely that it will be finished before you run this command line… Here is an example of the what we can observe:



**Computer used**

**Processor number used**

Name of the job (-N option in the qsub)

Owner of the job

Job ID (unique)

Priority of the job

Job status
r = running
qw = waiting to run (no computer available)
other (Eqw, dt, …) = there is a problem

# From the output of GBS sequencing to the variant calling file: in command line

- Output of the demultiplexing command line. Listing the current directory:

  `ll`

- One file and one folder are generated:

  **①** A file named DEMULT.o7157685
  - Correspond to the Name of the job passed to the qsub (`-N DEMULT`) concatenated with the unique job ID attributed by the master computer to the command line (here: **7157685**).
  - Because some programs "speak": this file contained what they say. We can have a look at what the program say with the more command:
  `more DEMULT.o7157685`

```
[gmartin@cc2-login vcfhunterGBS]$ more DEMULT.o7157685
Start the demultiplexing.
100000 reads demultiplexed
200000 reads demultiplexed
300000 reads demultiplexed
400000 reads demultiplexed
500000 reads demultiplexed
538230 reads demultiplexed
Demultiplexing ended.
[gmartin@cc2-login vcfhunterGBS]$ 
```

  **②** A folder named Demultiplexed
  This folder was created by GBSX as we tell him to do it with the (`-o Demultiplexed`) argument.

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the demultiplexed folder:

  `ll Demultiplexed`

A file summarizing demultiplexing options

A file with demultiplexing statistics

```
[gmartin@cc2-login vcfhunterGBS]$ ll Demultiplexed
total 4576
-rw-r--r-- 1 gmartin users   1081 Jan 11 09:19 gbsDemultiplex.log
-rw-r--r-- 1 gmartin users   1033 Jan 11 09:19 gbsDemultiplex.stats
-rw-r--r-- 1 gmartin users 392909 Jan 11 09:19 sample10.R1.fastq.gz
-rw-r--r-- 1 gmartin users 383291 Jan 11 09:19 sample11.R1.fastq.gz
-rw-r--r-- 1 gmartin users 393619 Jan 11 09:19 sample12.R1.fastq.gz
-rw-r--r-- 1 gmartin users 335532 Jan 11 09:19 sample1.R1.fastq.gz
-rw-r--r-- 1 gmartin users 373870 Jan 11 09:19 sample2.R1.fastq.gz
-rw-r--r-- 1 gmartin users 352415 Jan 11 09:19 sample3.R1.fastq.gz
-rw-r--r-- 1 gmartin users 378318 Jan 11 09:19 sample4.R1.fastq.gz
-rw-r--r-- 1 gmartin users 363557 Jan 11 09:19 sample5.R1.fastq.gz
-rw-r--r-- 1 gmartin users 381574 Jan 11 09:19 sample6.R1.fastq.gz
-rw-r--r-- 1 gmartin users 369568 Jan 11 09:19 sample7.R1.fastq.gz
-rw-r--r-- 1 gmartin users 392967 Jan 11 09:19 sample8.R1.fastq.gz
-rw-r--r-- 1 gmartin users 378680 Jan 11 09:19 sample9.R1.fastq.gz
-rw-r--r-- 1 gmartin users 152582 Jan 11 09:19 undetermined.fastq.gz
[gmartin@cc2-login vcfhunterGBS]$
```
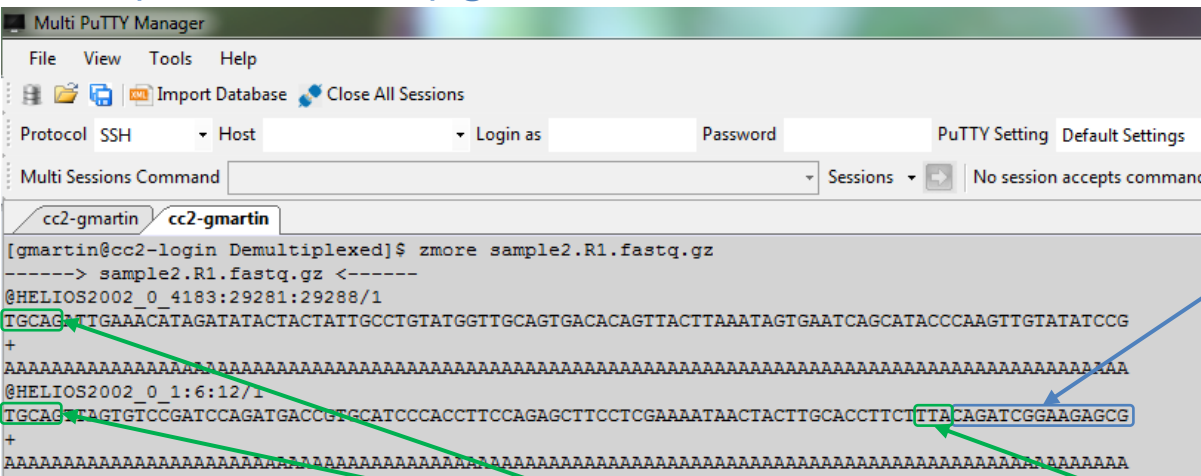
Reads parsed according to the accession they belong to

A file containing reads that could not be attributed to an accession (*i.e.* sequencing error in the tag)

- To have a look at these files:
  `more Demultiplexed/gbsDemultiplex.log` (for example)

- But because it is boring to always put Demultiplexed/ for all file which are in the directory, we will directly go into this directory:
  `cd Demultiplexed`

# From the output of GBS sequencing to the variant calling file: in command line

- The gbsDemultiplex.log file:

  ```
  more gbsDemultiplex.log
  ```

# From the output of GBS sequencing to the variant calling file: in command line

- The gbsDemultiplex.log file:

    `more gbsDemultiplex.stats`



- Not very easy to read… We will load this file on our computer.
    - ✓ For that we need FileZilla: https://filezilla-project.org/
    - ✓ Install it
    - ✓ Connect to your cluster account:

- The gbsDemultiplex.log file:



- Go to the Demultiplexed folder: work ➔ vcfhunterGBS ➔ Demultiplexed

- The gbsDemultiplex.log file:



- Your file has been copied to your desktop.
- Open it with Excel!

# From the output of GBS sequencing to the variant calling file: in command line

- The gbsDemultiplex.log file:



| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | sampleID | barcode | enzyme | total.count | total.perc | mismatch.0. | mismatch.0. | basecall.cou | basecall.abo | basecall.qual.avg | |
| 2 | sample1 | AACT | PstI | 44566 | 0.08280103 | 44566 | 1 | 4230367 | 1 | 32 | |
| 3 | sample10 | CCACG | PstI | 44228 | 0.08217305 | 44228 | 1 | 4198304 | 1 | 32 | |
| 4 | sample11 | TATAA | PstI | 44294 | 0.08229567 | 44294 | 1 | 4204723 | 1 | 32 | |
| 5 | sample12 | GAGCG | PstI | 44376 | 0.08244802 | 44376 | 1 | 4212401 | 1 | 32 | |
| 6 | sample2 | CCAG | PstI | 44579 | 0.08282519 | 44579 | 1 | 4231636 | 1 | 32 | |
| 7 | sample3 | TTGA | PstI | 44218 | 0.08215447 | 44218 | 1 | 4197408 | 1 | 32 | |
| 8 | sample4 | GGTA | PstI | 44408 | 0.08250748 | 44408 | 1 | 4215298 | 1 | 32 | |
| 9 | sample5 | ATTG | PstI | 44553 | 0.08277688 | 44553 | 1 | 4229275 | 1 | 32 | |
| 10 | sample6 | CGGT | PstI | 44219 | 0.08215633 | 44219 | 1 | 4197549 | 1 | 32 | |
| 11 | sample7 | TGCG | PstI | 44461 | 0.08260595 | 44461 | 1 | 4220581 | 1 | 32 | |
| 12 | sample8 | GTAT | PstI | 44491 | 0.08266169 | 44491 | 1 | 4223375 | 1 | 32 | |
| 13 | sample9 | AACCA | PstI | 44170 | 0.08206529 | 44170 | 1 | 4192803 | 1 | 32 | |
| 14 | undetermined | | | 5667 | 0.01052896 | | | | | | |
| 15 | | | | | | | | | | | |

# From the output of GBS sequencing to the variant calling file: in command line

- The sample***X***.R1.fastq.gz files: For example sample2.R1.fastq.gz

  ```
  zmore sample2.R1.fastq.gz
  ```



pstI restriction site

mseI restriction site

Adapter sequence

- Sample tags were removed from reads
- Illumina adapters are still present at the end of some read (i.e. when sequenced fragments are shorter than illumina reads) ➔ These adapters should be removed as they do not belong to the sample!

# From the output of GBS sequencing to the variant calling file: in command line

- Removing adapters and quality trimming of read.

The quality trimming is not necessary here as this is simulated reads with top quality but in reality as sequencing quality decrease along a read this is necessary.

- For that we will use cutadapt (https://cutadapt.readthedocs.io/en/stable/guide.html, https://doi.org/10.14806/ej.17.1.200)

- To load cutadapt:

```
module purge                          ← To remove already loaded modules (prevent conflicts)
module load bioinfo/cutadapt/1.8.1    ← The cutadapt module
module load system/python/3.4.3       ← cutadapt also required python module
```

- To use cutadapt on sample2, run the command line:

```
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG
-O 10 -q 20,20 -f fastq -m 30 -o sample2.R1.fastq.gz.cut.gz
sample2.R1.fastq.gz
```

# From the output of GBS sequencing to the variant calling file: in command line

- Command line explanation

  ```
  qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o
  sample2.R1.fastq.gz.cut.gz sample2.R1.fastq.gz
  ```

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  - ❖ **-b yes**: it is not important, but put it.

  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose to run the program

  - ❖ **-N CUTADAPT**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

# From the output of GBS sequencing to the variant calling file: in command line

- Command line explanation

```
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o
sample2.R1.fastq.gz.cut.gz sample2.R1.fastq.gz
```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**.

  ❖ **cutadapt**: tell that we will be using cutadapt program

  ❖ **-a CAGATCGGAAGAGCG**: tells cutadapt that it should look for adapter sequence at 3' end and that it should remove this sequence and all that follows.

  ❖ **-O 10**: If the overlap between the read and the adapter is shorter than 10, the read is not modified. This reduces the no. of bases trimmed purely due to short random adapter matches

  ❖ **-q 20,20**: Trim the 5' and the 3' until a base quality of 20 is reached

  ❖ **-f fastq** : The input format file is fastq

  ❖ **-m 30**  : only read equal or greater than 30 bases will be conserved

  ❖ **-o sample2.R1.fastq.gz.cut.gz**: Name of the output file

  ❖ **sample2.R1.fastq.gz**: Name of the file processed by cutadapt

# From the output of GBS sequencing to the variant calling file: in command line

- Outputs: To visualize new file generated, list the files in the repository:

```
ll
```



**Two files have been generated**:

① The CUTADAPT.oxxxxxxx file containing what cutadapt told us while it was executing

② The sample2.R1.fastq.gz.cut.gz containing filtered read

# From the output of GBS sequencing to the variant calling file: in command line

- The sample2.R1.fastq.gz file before cutadapt:
  `zmore sample2.R1.fastq.gz`



Adapter sequence

pstI restriction site

mseI restriction site

- And After cutadapt
  `zmore sample2.R1.fastq.gz.cut.gz`

- The CUTADAPT.oxxxxxxx file:

  `zmore CUTADAPT.oxxxxxxx`



There is a warning saying that maybe the adapter sequence is incomplete because very often (99.8% of cases), when an adapter is found, the "A" base was found just before…

This is normal because just before the adapter we have our *mseI* restriction site

# From the output of GBS sequencing to the variant calling file: in command line

- This command line should be adapted and executed for each sample:

```
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample10.R1.fastq.gz.cut.gz sample10.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample11.R1.fastq.gz.cut.gz sample11.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample12.R1.fastq.gz.cut.gz sample12.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample1.R1.fastq.gz.cut.gz sample1.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample2.R1.fastq.gz.cut.gz sample2.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample3.R1.fastq.gz.cut.gz sample3.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample4.R1.fastq.gz.cut.gz sample4.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample5.R1.fastq.gz.cut.gz sample5.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample6.R1.fastq.gz.cut.gz sample6.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample7.R1.fastq.gz.cut.gz sample7.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample8.R1.fastq.gz.cut.gz sample8.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample9.R1.fastq.gz.cut.gz sample9.R1.fastq.gz
```

- This is relatively easy when we have few files but when this should be done on hundreds of files it is a bit annoying… This can be solved with "for" loop in bash programing!

- Here is the command line for our example (advanced programing!):

```
for i in *.fastq.gz

do qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a
CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o $i.cut.gz $i

done
```

# From the output of GBS sequencing to the variant calling file: in command line

- This command line should be adapted and executed for each sample:

```
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample10.R1.fastq.gz.cut.gz sample10.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample11.R1.fastq.gz.cut.gz sample11.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample12.R1.fastq.gz.cut.gz sample12.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample1.R1.fastq.gz.cut.gz sample1.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample2.R1.fastq.gz.cut.gz sample2.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample3.R1.fastq.gz.cut.gz sample3.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample4.R1.fastq.gz.cut.gz sample4.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample5.R1.fastq.gz.cut.gz sample5.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample6.R1.fastq.gz.cut.gz sample6.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample7.R1.fastq.gz.cut.gz sample7.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample8.R1.fastq.gz.cut.gz sample8.R1.fastq.gz
qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o sample9.R1.fastq.gz.cut.gz sample9.R1.fastq.gz
```

- This is relatively easy when we have few files but when this should be done on hundreds of files it is a bit annoying… This can be solved with "for" loop in bash programing!

- Here is the command line for our example (advanced programing!):

**Initiation of a loop: For all files in the folder finishing by ".fastq.gz"…**

```
for i in *.fastq.gz
```

**Their name is sequentially stored in a variable "i", and, for each values "i" (each read sample files), the cutadapt command line is executed on the file recorded in the variable i ($i) and the output is stored in a file called i+".cut.gz" ($i.cut.gz) .**
**For example when i = sample10.R1.fastq.gz : $i.cut.gz = sample10.R1.fastq.gz.cut.gz**

```
do qsub -q normal.q -b yes -V -N CUTADAPT cutadapt -a
CAGATCGGAAGAGCG -O 10 -q 20,20 -f fastq -m 30 -o $i.cut.gz $i
```

**Tell that this is the end of the loop**

```
done
```

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the files in the folder:

  `ll`



- A CUTADAPT.oxxxxxxx file has been generated per sample

- A filtered fastq file per sample has been generated per accessions

# From the output of GBS sequencing to the variant calling file: in command line

- We will use vcfHunter program which is installed on the AGAP cluster under module "*vcfhunter*"

- To load this module run the command line:

  ```
  module purge
  module load bioinfo/vcfhunter/1.0.0
  ```

- The module is now loaded. This can be verified with de following command line:

  ```
  module list
  ```



- We can see that the vcfhunter module is loaded as well as several other modules which will be used by vcfhunter

# From the output of GBS sequencing to the variant calling file: in command line

- We are going to work in a new folder for vcfHunter. This is not necessary but for file ordering, this will be better. But first where are we? To answer this question we use a simple command:

  pwd



This locate the path where you are when you execute the pwd command. Instead of "gmartin", you should have your login ID

- From there we want to go back to vcfhunterGBS folder. There are two possibility:

  cd /home/**Your_ID**/work/vcfhunterGBS

  **c**hange **d**irectory to
  /home/**Your_ID**/work/vcfhunterGBS

  Or

  cd ..

  **c**hange **d**irectory to one folder before. And
  one folder before there is vcfhunterGBS

# From the output of GBS sequencing to the variant calling file: in command line

- Where are we now?

pwd



```
[gmartin@cc2-n16 vcfhunterGBS]$ pwd
/home/gmartin/work/vcfhunterGBS
[gmartin@cc2-n16 vcfhunterGBS]$
```

- Now we create the new folder

  mkdir Mapping

- And we go into this folder

  cd Mapping

# From the output of GBS sequencing to the variant calling file: in command line

- At this stage, we have 12 fastq files:
  - ✓ One for each samples, which comprised cleaned/filtered reads.
  - ✓ These files are located in a folder named `Demultiplexed`, located `/home/Your_ID/work/vcfhunterGBS`

- To run vcfHunter program, we also need an additional file which contained the reference sequence (in fasta format), on which we will align the reads. This file is already present in the `WorkShopDataset` folder located here: `/home/Your_ID/work/vcfhunterGBS/WorkShopDataset`. This file is named `Ref.fasta` (It is the folder you copied at the beginning of this exercise).

- Because at this time we are in the `Mapping` folder loacted `/home/Your_ID/work/vcfhunterGBS/Mapping`, to have a look at this file we should go back from one folder (`..`) to enter the `WorkShopDataset` folder and then access to `Ref.fasta` file. Thus, to have a look at this file:

  `more ../WorkShopDataset/Ref.fasta`



Standard fasta format with each sequences beginning with a ">"+sequence name ①, followed by DNA sequence ②.

# From the output of GBS sequencing to the variant calling file: in command line

- The sample fastq read file and reference fasta files should be passed recorded in a configuration unique file which will be given to *vcfHunter* program.

- For this example, the configuration file (`GBSCalling.conf`) has already been created can be found here: `/home/Your_ID/work/vcfhunterGBS/WorkShopDataset`. To have a look at this file and because we are in the Mapping folder we just created:

  `more ../WorkShopDataset/GBSCalling.conf`

**A [Reference] section locating how to access reference fasta.**

**A [Libraries] section locating how to access sample fastq reads files and additional information to sample:**
① **Unique ID for each fastq**
② **Sample Name (Name that will appear in the vcf)**
③ **How to access to the fastq read file**
④ **Accession ploidy**

- Possible to generate this file with a loop for Those who want to try!

```
Multi PuTTY Manager
File   View   Tools   Help
[XML] Import Database   Close All Sessions
Protocol SSH     ▾ Host                    ▾ Login as              Password
Multi Sessions Command

  cc2-gmartin   cc2-gmartin   cc2-gmartin
[gmartin@cc2-n16 Mapping]$ more ../WorkShopDataset/GBSCalling.conf
[Reference]
genome = ../WorkShopDataset/Ref.fasta
[Libraries]
Lib01 = S1 ../Demultiplexed/sample1.R1.fastq.gz.cut.gz 2
Lib02 = S2 ../Demultiplexed/sample2.R1.fastq.gz.cut.gz 2
Lib03 = S3 ../Demultiplexed/sample3.R1.fastq.gz.cut.gz 2
Lib04 = S4 ../Demultiplexed/sample4.R1.fastq.gz.cut.gz 2
Lib05 = S5 ../Demultiplexed/sample5.R1.fastq.gz.cut.gz 2
Lib06 = S6 ../Demultiplexed/sample6.R1.fastq.gz.cut.gz 2
Lib07 = S7 ../Demultiplexed/sample7.R1.fastq.gz.cut.gz 2
Lib08 = S8 ../Demultiplexed/sample8.R1.fastq.gz.cut.gz 2
Lib09 = S9 ../Demultiplexed/sample9.R1.fastq.gz.cut.gz 2
Lib10 = S10 ../Demultiplexed/sample10.R1.fastq.gz.cut.gz 2
Lib11 = S11 ../Demultiplexed/sample11.R1.fastq.gz.cut.gz 2
Lib12 = S12 ../Demultiplexed/sample12.R1.fastq.gz.cut.gz 2
```
① ② ③ ④

- One last thing before using vcfHunter module: This program has several programs, we will use process_reseq_1.0.py  program which have several options, to have access to a description of these options, you can try the following command line:

```
process_reseq_1.0.py -h
```



**Several options**

**Distinct steps: which will be performed sequentially for better explanation**

# From the output of GBS sequencing to the variant calling file: in command line

- Running read mapping process
  **`qsub -q normal.q -l mem_free=12G -b yes -V -N GBSa`** `"process_reseq_1.0.py -c ../WorkShopDataset/GBSCalling.conf -p GBSset -s a -t 1"`

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  - ❖ **-l mem_free=12G**: precise that the program will use 12G of RAM (so the master computer will check that it is available on the computers). This is necessary because this step will use **java** program and this will prevent errors…

  - ❖ **-b yes**: it is not important, but put it.

  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose

  - ❖ **-N GBSa**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

# From the output of GBS sequencing to the variant calling file: in command line
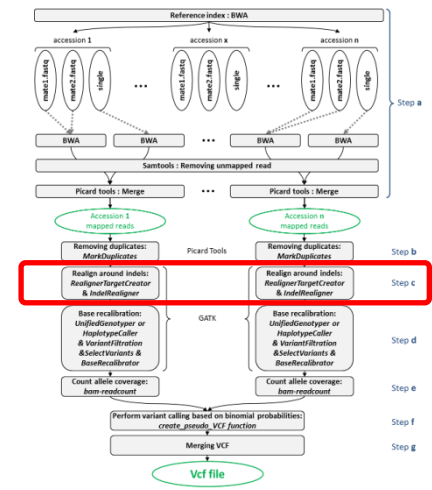
- Running read mapping process
  ```
  qsub -q normal.q -l mem_free=12G -b yes -V -N GBSa "process_reseq_1.0.py -c ../WorkShopDataset/GBSCalling.conf -p GBSset -s a -t 1"
  ```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**:
  - ❖ **process_reseq_1.0.py**: We will use process_reseq_1.0.py program
  - ❖ **-c ../WorkShopDataset/GBSCalling.conf**: Locates the configuration file
  - ❖ **-p GBSset**: A prefix for final output file
  - ❖ **-s a**: Tell the program that we will perform step "a" of the workflow



  - ❖ **-t 1**: Tell the program that only one processor is available. This means that each accessions will be treated sequencially

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the files generated:

  `ll`



The GBSa.oxxxxxxx file containing what process_reseq_1.0.py told us while it was executing

A folder for each accession. Which contained several items. To have a look at these items, for example for S1 accession:

`ll S1`

Read mapping statistics

A .bai file: is an index of a bam file for computation performance

A .bam file: contained sample 1 reads aligned onto the reference

A folder containing read and alignment statistics

- Listing one of the stat folder:
  ```
  ll S1/STATS/
  ```



- Several files are generated but one summarize all of them: the one named: index.html
- This is an html file readable by firefox. To have a look at this file:
  ```
  firefox S1/STATS/index.html
  ```
- This command open a firefox window:

- Listing one of the stat folder:
  ```
  ll S1/STATS/
  ```



- Several files are generated but one summarize all of them: the one named: index.html
- This is an html file readable by firefox. To have a look at this file:
  ```
  firefox S1/STATS/index.html
  ```
- This command open a firefox window:





| Reads | | |
|---|---|---|
| total: | 44,566 | |
| filtered: | 0 | (0.0%) |
| non-primary: | 0 | |
| duplicated: | 0 | (0.0%) |
| mapped: | 44,565 | (100.0%) |
| zero MQ: | 0 | (0.0%) |
| avg read length: | 93 | |

| Bases | | |
|---|---|---|
| total: | 4,170,151 | (99.2%) |
| mapped: | 4,135,266 | |
| error rate: | 1.73% | |

51

# From the output of GBS sequencing to the variant calling file: in command line

- The alignment file (bam format): These file are compressed binary files (easier to use by programs) but not directly readable for human... These file can still be observed with the **samtools** program with the command line:

```
samtools view -h S1/S1_merged.bam | more
```

Convert the bam is sam format (readable by human)

Read this converted file line by line

**Mapping quality**

**Read position**

**Read sequence**

**Read quality**

**Header containing information on:**
- **reference sequences**
- **Aligner used**



**Read name**

**Tag regarding read mapping information**
https://broadinstitute.github.io/picard/explain-flags.html

**Read mapping chromosome**

**CIGAR (80M = 80 Match)**

+ Other informations (see link for more information)
https://samtools.github.io/hts-specs/SAMv1.pdf

- **To quit:** "Ctrl" + "C"

# From the output of GBS sequencing to the variant calling file: in command line

- The GBSa.oxxxxxxx file:
  `more GBSa.oxxxxxxxx`
- List steps performed during step "a"

**Reference indexation**

**Read alignment with "bwa"**

**Calculating alignment statistics**

**Ploting alignment stats**

**Removing unmapped read and secondary alignment**

**Starting a new accession**

**Concatenating reads from the same accessions but from several libraries (not necessary here but performed anyway)**

- **To quit:** "Ctrl" + "C" or "enter" until the end of file

# From the output of GBS sequencing to the variant calling file: in command line

- Running read indel realignment:
  ```
  qsub -q normal.q -l mem_free=12G -b yes -V -N GBSc "process_reseq_1.0.py -c
  ../WorkShopDataset/GBSCalling.conf -p GBSset -s c -t 1"
  ```

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  ❖ **-l mem_free=12G**: precise that the program will use 12G of RAM (so the master computer will check that it is available on the computers). This is necessary because this step will use **java** program and this will prevent errors...

  ❖ **-b yes**: it is not important, but put it.

  ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose

  ❖ **-N GBSc**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

# From the output of GBS sequencing to the variant calling file: in command line
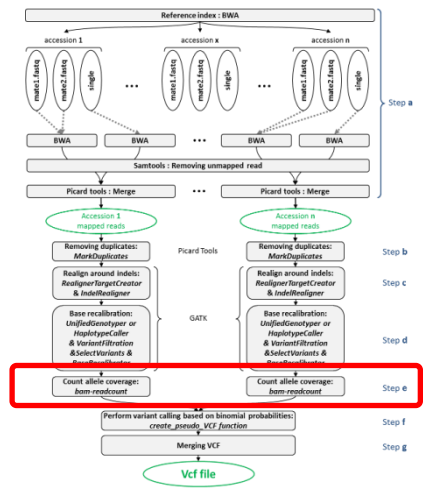
- Running read indel realignment:
  ```
  qsub -q normal.q -l mem_free=12G -b yes -V -N GBSc "process_reseq_1.0.py -c
  ../WorkShopDataset/GBSCalling.conf -p GBSset -s c -t 1"
  ```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**:
  - ❖ **process_reseq_1.0.py**: We will use process_reseq_1.0.py program
  - ❖ **-c ../WorkShopDataset/GBSCalling.conf**: Locates the configuration file
  - ❖ **-p GBSset**: A prefix for final output file
  - ❖ **-s c**: Tell the program that we will perform step "c" of the workflow

  

  - ❖ **-t 1**: Tell the program that only one processor is available. This means that each accessions will be treated sequencially

- Why performing indel realalignment?
  - ✓ Because the alignment around indel can be problematic...

      Reference   GCAACAAGGGTTACAGATCGGAAAAGAGCGGTTCAGCAGGAATGCCG
                  CAAGGGTTACAGATCGGAAA-TAGCGGTTCAGCA
                          GGGTTACAGATCGGAAAT-AGCGGTTCAGCAGGAATGCCG
                          AGGGTTACAGATCGGAA-ATAGCGGTTCAGCAGGAATGCCG
                                          ***

      indel        Indel+SNP        SNP+indel

  - ✓ ➜ several polimorphism with the same sequence!

➜ Realignment around indel:

      Reference   GCAACAAGGGTTACAGATCGGAAAAGAGCGGTTCAGCAGGAATGCCG
                  CAAGGGTTACAGATCGGAAA-TAGCGGTTCAGCA
                          GGGTTACAGATCGGAAA-TAGCGGTTCAGCAGGAATGCCG
                          AGGGTTACAGATCGGAAA-TAGCGGTTCAGCAGGAATGCCG
                                          **

      indel                      SNP

56

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the files generated:

`ll`



The GBSc.oxxxxxxx file containing what process_reseq_1.0.py told us while it was executing

A folder for each accession which contained realigned reads. To have a look at these files, for example for S1 accession:

`ll S1`

A .bai file: is an index of the realigned bam file for computation performance

A realigned.bam file: contained sample 1 reads realigned around indels

# From the output of GBS sequencing to the variant calling file: in command line

- The GBSa.oxxxxxxx file:
  `more GBSc.oxxxxxxx`
- List steps performed during step "c"



- Indel realignment was performed using GATK (https://software.broadinstitute.org/gatk/) in two steps (https://software.broadinstitute.org/gatk/documentation/tooldocs/3.8-0/org_broadinstitute_gatk_tools_walkers_indels_RealignerTargetCreator.php):

  ① "Determining (small) suspicious intervals which are likely in need of realignment"

  ② "Running the realigner over those intervals"

**From the output of GBS sequencing to the variant calling file: in command line**

- Running allele count:
  **qsub -q normal.q -b yes -V -N GBSe** "process_reseq_1.0.py -c
  ../WorkShopDataset/GBSCalling.conf -p GBSset -s e -t 1"

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  - ❖ **-b yes**: it is not important, but put it.

  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose

  - ❖ **-N GBSe**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

# From the output of GBS sequencing to the variant calling file: in command line

- Running allele count:
  ```
  qsub -q normal.q -b yes -V -N GBSe "process_reseq_1.0.py -c
  ../WorkShopDataset/GBSCalling.conf -p GBSset -s e -t 1"
  ```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**:
  - ❖ **process_reseq_1.0.py**: We will use process_reseq_1.0.py program
  - ❖ **-c ../WorkShopDataset/GBSCalling.conf**: Locates the configuration file
  - ❖ **-p GBSset**: A prefix for final output file
  - ❖ **-s e**: Tell the program that we will perform step "e" of the workflow



  - ❖ **-t 1**: Tell the program that only one processor is available. This means that each accessions will be treated sequencially

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the files generated:

  `ll`



The GBSe.oxxxxxxx file containing what process_reseq_1.0.py told us while it was executing

A folder for each accession which contained realigned reads. To have a look at these files, for example for S1 accession:

`ll S1`

Three files (one for each chromosomes) which count for each covered position by reads, the number of read supporting each possible alleles

- Example of **_S1_allele_count_chr01.gz_** file:
  zmore S1/S1_allele_count_chr01.gz



Chromosome

Position

Reference base

Total read coverage

Reads with A alleles

Reads with C alleles ...

Reads with deletion

- **To quit:** "Ctrl" + "C" or "enter"
  **until the end of file**

# From the output of GBS sequencing to the variant calling file: in command line

- Creating the variant calling file (VCF):
  ```
  qsub -q normal.q -pe parallel_smp 3 -b yes -V -N GBSf "process_reseq_1.0.py
  -c ../WorkShopDataset/GBSCalling.conf -p GBSset -s f -t 3"
  ```

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  - ❖ **-pe parallel_smp 3**: tells the master computer that we need 3 processor (this can be used to gain speed in computation time if the program allowed it)

  - ❖ **-b yes**: it is not important, but put it.

  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose

  - ❖ **-N GBSf**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

# From the output of GBS sequencing to the variant calling file: in command line

- Creating the variant calling file (VCF):

  ```
  qsub -q normal.q -pe parallel_smp 3 -b yes -V -N GBSf "process_reseq_1.0.py
  -c ../WorkShopDataset/GBSCalling.conf -p GBSset -s f -t 3"
  ```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**:
  - ❖ **process_reseq_1.0.py**: We will use process_reseq_1.0.py program
  - ❖ **-c ../WorkShopDataset/GBSCalling.conf**: Locates the configuration file
  - ❖ **-p GBSset**: A prefix for final output file
  - ❖ **-s f**: Tell the program that we will perform step "f" of the workflow



  - ❖ **-t 3**: Tell the program that only three processors are available (allowed by -pe parallel_smp 3). With this option, all three chromosomes will be treated independently at the same time by one processor each. This allowed to gain computation time

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the files generated:

  `ll`



**The GBSf.oxxxxxxx file containing what process_reseq_1.0.py told us while it was executing**

**An always empty file associated to -pe parallel_smp 3 options**

**Three vcf files containing genotyping informations, one for each chromosomes**

# From the output of GBS sequencing to the variant calling file: in command line

- What can be found in a vcf format:

  `more GBSset_chr01_all_allele_count.vcf`



① **Real header of variant calling file**

② **Variant line 1**

③ **Variant line 3**

**Header of the vcf file containing information about:**
- ✓ **Reference file location**
- ✓ **Genotype format description**
- ✓ **Reference sequence name and size**

- **To quit:** "Ctrl" + "C" or "enter"
  **until the end of file**

# From the output of GBS sequencing to the variant calling file: in command line

- Looking at the vcf file with excel because it is easier (Not to do on real dataset!):

- Using filezilla to get the data on our computer:



- Open it with Excel!

# From the output of GBS sequencing to the variant calling file: in command line

- The vcf file format:

# From the output of GBS sequencing to the variant calling file: in command line

- The vcf file format:

Describe the way the genotype is formatted for each accessions:

- ✓ GT = genotype
- ✓ AD = allele depth
- ✓ DP = depth

| 9 | #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | S1 | S10 | S11 | S12 | S2 |
|---|--------|-----|----|-----|-----|------|--------|------|--------|-----|-----|-----|-----|-----|
| 10 | chr01 | 30 | . | A | T | . | . | . | GT:AD:DP | ./.:0,0:0 | ./.:0,1:1 | ./.:0,0:0 | ./.:0,0:0 | ./.:0,0:0 |
| 11 | chr01 | 36 | . | T | A,C | . | . | . | GT:AD:DP | 0/0:17,0,0:17 | 0/0:19,0,0:19 | 0/0:24,0,0:24 | 0/0:18,0,0:18 | 0/0:16,1,0:17 |

- ✓ GT = 0/0
- ✓ AD = 17,0,0
- ✓ DP = 17

Based on these allelic depths, calculation of the likelihood of each haplotypes:

- ✓ 0/0 = T/T
- ✓ 0/1 = T/A
- ✓ 0/2 = T/C
- ✓ 1/2 = A/C
- ✓ 1/1 = A/A
- ✓ 2/2 = C/C

# From the output of GBS sequencing to the variant calling file: in command line

- Because it is sometime easier to have only one file for all chromosomes, this unique file can be produced with this last command line:
  ```
  qsub -q normal.q -b yes -V -N GBSg "process_reseq_1.0.py -c
  ../WorkShopDataset/GBSCalling.conf -p GBSset -s g -t 1"
  ```

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  - ❖ **-b yes**: it is not important, but put it.

  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose

  - ❖ **-N GBSf**: A name passed to the command line to look at its status (waiting, running or error) on the cluster
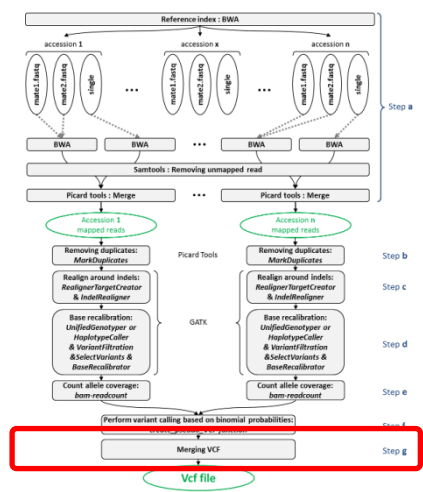
# From the output of GBS sequencing to the variant calling file: in command line

- Because it is sometime easier to have only one file for all chromosomes, this unique file can be produced with this last command line:

  ```
  qsub -q normal.q -b yes -V -N GBSg "process_reseq_1.0.py -c
  ../WorkShopDataset/GBSCalling.conf -p GBSset -s g -t 1"
  ```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**:
  - ❖ **process_reseq_1.0.py**: We will use process_reseq_1.0.py program
  - ❖ **-c ../WorkShopDataset/GBSCalling.conf**: Locates the configuration file
  - ❖ **-p GBSset**: A prefix for final output file
  - ❖ **-s g**: Tell the program that we will perform step "g" of the workflow



  - ❖ **-t 1**: Tell the program that only one processor is available.

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the file generated:

  `ll`



**The GBSg.oxxxxxxx file containing what process_reseq_1.0.py told us while it was executing**

**A vcf file containing all chromosomes**

# From the output of GBS sequencing to the variant calling file: in command line

- To discriminate between sequencing errors and true variant site we developed an additional program which allowed to select true polymorphous SNP according to selected parameters. This program is called *VcfPreFilter.1.0.py* and can be executed with the following command line:

```
qsub -q normal.q -b yes -V -N PREFLTR "VcfPreFilter.1.0.py -v
GBSset_all_allele_count.vcf -m 10 -M 10000 -f 0.05 -c 3 -o
GBSset_prefiltered.vcf -d y"
```

- The first part of the command line (in bold) is **used by the master computer** (as previously described):

  - ❖ **qsub**: Means that we will send a command that the master computer needs to analyze to choose the best computer

  - ❖ **-q normal.q**: tells the master computer that we will use computer from normal queue.

  - ❖ **-b yes**: it is not important, but put it.

  - ❖ **-V**: Tell the master computer to load the module previously loaded on the computer it will choose

  - ❖ **-N PREFLTR**: A name passed to the command line to look at its status (waiting, running or error) on the cluster

```
qsub -q normal.q -b yes -V -N PREFLTR "VcfPreFilter.1.0.py -v
GBSset_all_allele_count.vcf -m 10 -M 10000 -f 0.05 -c 3 -o
GBSset_prefiltered.vcf -d y"
```

- The part of the command line between quotation marks (in bold) is the command line that is executed on the **computer chosen by the master computer**:
  - ❖ **VcfPreFilter.1.0.py**: We will use VcfPreFilter.1.0.py program
  - ❖ **-v GBSset_all_allele_count.vcf**: Locates the vcf file to filter
  - ❖ **-m 10** : Only datapoint with coverage supported by more than 10 reads will be considered
  - ❖ **-M 10000**: Only datapoint with coverage supported by less than 10000 reads will be considered (to manage very high repeats)
  - ❖ **-f 0.05**: An allele is kept if it is present in at least this proportion in at least one accession.
  - ❖ **-c 3**: An allele is kept if it is supported by at least **3** reads in at least one accession.
  - ❖ **-o GBSset_prefiltered.vcf**: Name of the output file
  - ❖ **-d y**: Perform only diallelic calling (i.e. for triploid accessions, A/C/G genotype is not possible because only two alleles are allowed in a genotype: A/A/C or A/G/G, or … genotype are tested).

- According to -m, -M, -f and -c parameters the number of possible alleles is counted (including the reference sequence allele, and if this number is strictly greater than 1, the line is identified as a polymorphous line that should be reported)

# From the output of GBS sequencing to the variant calling file: in command line

- Prefiltering example: `-m 10 -M 10000 -f 0.05 -c 3`

**No allele pass the -m 10 cutoff**

**Number of alleles reported = 0**
➔ **Not a reported variant line**

# From the output of GBS sequencing to the variant calling file: in command line

- Prefiltering example: `-m 10 -M 10000 -f 0.05 -c 3`



Allele passing cutoffs: A  A  A  A  A  A  A  A  A  A  A̶  A

**Number of alleles reported = 1 < 2**
**➔ Not a reported variant line because homozygous.**

**Reported first because sequencing error in S2 with one read having "C"**

# From the output of GBS sequencing to the variant calling file: in command line

- Prefiltering example:  `-m 10 -M 10000 -f 0.05 -c 3`

GBSset_chr01_all_allele_count.vcf - Microsoft Excel

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ##fileformat=VCFv4.2 | | | | | | | | | | | | | | | | | | | | |
| 2 | ##reference=file:///../WorkShopDataset/Ref.fasta | | | | | | | | | | | | | | | | | | | | |
| 3 | ##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype"> | | | | | | | | | | | | | | | | | | | | |
| 4 | ##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth"> | | | | | | | | | | | | | | | | | | | | |
| 5 | ##FORMAT=<ID=AD,Number=.,Type=Integer,Description="Allelic depths for the ref and alt alleles in the order listed"> | | | | | | | | | | | | | | | | | | | | |
| 6 | ##contig=<ID=chr03,length=100001> | | | | | | | | | | | | | | | | | | | | |
| 7 | ##contig=<ID=chr02,length=100001> | | | | | | | | | | | | | | | | | | | | |
| 8 | ##contig=<ID=chr01,length=100001> | | | | | | | | | | | | | | | | | | | | |
| 9 | #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | S1 | S10 | S11 | S12 | S2 | S3 | S4 | S5 | S6 | S7 | S8 | S9 |
| 10 | chr01 | 30 | . | A | T | . | . | . | GT:AD:DP | ./.:0,0:0 | ./.:0,1:1 | ./.:0,0:0 | ./.:0,0:0 | ./.:0,0:0 | ./.:0,0:0 | ./.:1,0:1 | ./.:0,0:0 | ./.:0,0:0 | ./.:0,0:0 | ./.:0,0:0 |
| 11 | chr01 | 36 | . | T | A,C | . | . | . | GT:AD:DP | 0/0:17,0,0:17 | 0/0:19,0,0:19 | 0/0:24,0,0:24 | 0/0:18,0,0:18 | 0/0:16,1,0:17 | 0/0:16,0,0:16 | 0/0:19,0,0:19 | 0/0:16,0,0:16 | 0/0:14,0,0:14 | 0/0:14,0,1:15 | 0/0:18,0,0:18 | 0/0:17,0,0:17 |
| 12 | chr01 | 39 | . | A | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:17,0:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:15,0:15 | 0/0:17,1:18 | 0/0:17,0:17 |
| 13 | chr01 | 42 | . | A | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:17,0:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:14,1:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 14 | chr01 | 49 | . | G | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:17,0:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:14,1:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 15 | chr01 | 50 | . | T | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:18,1:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:17,0:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:15,0:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 16 | chr01 | 51 | . | A | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:16,1:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:15,0:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 17 | chr01 | 52 | . | C | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | | | | |
| 18 | chr01 | 53 | . | A | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 19 | chr01 | 55 | . | T | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:17,1:18 | 0/0:17,0:17 |
| 20 | chr01 | 59 | . | A | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:16,1:17 |
| 21 | chr01 | 63 | . | A | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 22 | chr01 | 67 | . | A | C | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:17,1:18 | 0/0:17,0:17 |
| 23 | chr01 | 68 | . | A | T | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:16,1:17 |
| 24 | chr01 | 71 | . | T | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 25 | chr01 | 73 | . | C | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:18,1:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 26 | chr01 | 74 | . | T | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 27 | chr01 | 75 | . | C | T | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:16,1:17 |
| 28 | chr01 | 76 | . | T | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/ | | | | | | :15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 29 | chr01 | 77 | . | T | A | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:16,1:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:15,0:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 30 | chr01 | 79 | . | C | T | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:16,1:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:14,0:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 31 | chr01 | 87 | . | C | G | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:16,1:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:14,0:14 | 0/0:15,0:15 | 0/0:17,0:18 | 0/0:17,0:17 |
| 32 | chr01 | 89 | . | C | T | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:17,0:18 | 0/0:17,0:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:15,1:16 | 0/0:14,0:14 | 0/0:15,0:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 33 | chr01 | 90 | . | G | A | . | . | . | GT:AD:DP | 0/0:17,0:17 | 0/0:19,0:19 | 0/0:24,0:24 | 0/0:18,0:18 | 0/0:17,0:17 | 0/0:16,0:16 | 0/0:19,0:19 | 0/0:16,0:16 | 0/0:13,1:14 | 0/0:15,0:15 | 0/0:18,0:18 | 0/0:17,0:17 |
| 38 | chr01 | 107 | . | A | C,G,T | . | . | . | GT:AD:D | 0/0:17,0,0:17 | 3/3:0,0,0,19:19 | 0/3:11,0,1,12:24 | 3/3:0,0,0,18:18 | 0/0:17,0,0:17 | 3/3:0,0,0,16:16 | 3/3:0,1,0,18:19 | 3/3:0,0,0,16:16 | 3/3:0,0,0,14:14 | 3/3:0,0,0,15:15 | 3/3:0,0,0,18:18 | 3/3:0,0,0,17:17 |

**Number of alleles reported = 2**
**➔ Reported variant line because polymorphism was detected (according to passed parameters).**

Allele passing cutoffs:   A   T   A~~G~~T   T   A   T   ~~C~~T   T   T   T   T   T

- Listing the files generated:
    `ll`



**A vcf file prefiltered**

**PREFLTR.oxxxxxxx file containing what VcfPreFilter.1.0.py told us while it was executing**

# From the output of GBS sequencing to the variant calling file: in command line
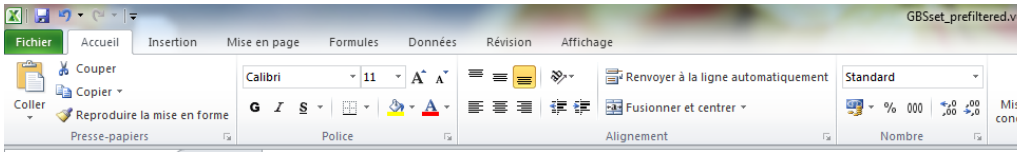
- Download this file with filezilla:

# From the output of GBS sequencing to the variant calling file: in command line

- Open the vcf with excel: (less missing data)



**An additional tag (GC) appeared: the ratio between the best genotype probability found and the second best genotype probability found.**

# From the output of GBS sequencing to the variant calling file: in command line

- This prefiltering step was designed to discriminate between variant lines resulting from sequencing errors and true variant line.

- However, one can want to apply additional filters such as reporting only diallelic polymorphous SNP, minimal coverage confidence to call a variant, missing data proportion, etc…

- For that we first need to generate a file containing a list of accessions we want to apply filter on. If we want to apply this filter on all accessions of the vcf, this file can be generated by a "simple" command line that will work on any vcf files you have!

```
head -n 10000 GBSset_prefiltered.vcf | grep "#CHROM" | sed 's/\t/\n/g' |
tail -n +10 > all_names.tab
```

- We take the first 10000 lines of the vcf: `head -n 10000 GBSset_prefiltered.vcf`
- Of these 10000 lines, we get the line with the accessions names which also contained "#CHROM": `grep "#CHROM"`
- Of this line we convert tabulation into carriage return: `sed 's/\t/\n/g'`
- And we take all lines from the result, but only from line number 10 to the end: `tail -n +10`
- The selected lines are written to a file named: `all_names.tab`

# From the output of GBS sequencing to the variant calling file: in command line

- Once the name file as been created: this can be verified with `ll` command:

```
[gmartin@cc2-login Mapping]$ ll
total 28544
-rw-r--r-- 1 gmartin users        39 Jan 16 10:26 all_names.tab
-rw-r--r-- 1 gmartin users     13345 Jan 15 08:07 GBSa.o7186205
-rw-r--r-- 1 gmartin users     13722 Jan 15 09:35 GBSc.o7186541
-rw-r--r-- 1 gmartin users      2571 Jan 15 10:08 GBSe.o7186786
```

```
[gmartin@cc2-login Mapping]$ more all_names.tab
S1
S10
S11
S12
S2
S3
S4
S5
S6
S7
S8
S9
```

- This file contained accession names: `more all_names.tab`

- A third script, called vcfFilter.1.0.py as been designed to filter the vcf (GBSset_prefiltered.vcf). For example, we may want to:

(1) convert to missing data:

- ✓ all datapoints which are not supported by at least **15** reads (no sufficient coverage to call good genotype)
- ✓ all datapoints which are not supported by more than **300** reads (probably repeat sequences)
- ✓ all datapoints for which each alleles is not supported by **3** read and a minimal read proportion of 0.2

(2) remove all line which contained missing data,

(3) remove mono, tri and tetra allelic sites,

(4) write the output in a file which prefix is **GBSset_Filtered.**

To apply these filters do not try the command following command line:

```
qsub -q normal.q -b yes -V -N FLTR "vcfFilter.1.0.py --vcf
GBSset_prefiltered.vcf --names all_names.tab --MinCov 15 --MaxCov 300 --
MinAl 3 --MinFreq 0.2 --nMiss 0 --RmAlAlt 1:3:4 --prefix GBSset_Filtered"
```

# From the output of GBS sequencing to the variant calling file: in command line

- Listing the files generated:

  `ll`

**FLTR.oxxxxxxx file containing what vcfFilter.1.0.py told us while it was executing**

`more FLTR.oxxxxxxx`

- A tutorial for variant calling of WGS data is also available here:

  https://github.com/SouthGreenPlatform/VcfHunter/blob/master/tutorial_VariantCalling.md

- Vcfhunter module contained additional tools for genetic mapping analysis and chromosome painting described and available here:

  https://github.com/SouthGreenPlatform/VcfHunter