

# Mapping, SNP & SV

---

Christine Tranchant-Dubreuil & Francois Sabot

October, 2018

IRD - UMR DIADE

- Quality control of NGS data

- Quality control of NGS data
- Learn to manipulate NGS data

- Quality control of NGS data
- Learn to manipulate NGS data
- Having a critical look on *Mapping*

- Quality control of NGS data
- Learn to manipulate NGS data
- Having a critical look on *Mapping*
- Learn to launch a *Calling* and having a critical look

- Quality control of NGS data
- Learn to manipulate NGS data
- Having a critical look on *Mapping*
- Learn to launch a *Calling* and having a critical look
- Learn the basic of Structural Variations

What is a Cluster ?



What is a Cluster ?

- A logical unit composed of multiple servers





What is a Cluster ?



- A logical unit composed of multiple servers
- Will work as a unique powerful machine

## What is a Cluster ?



- A logical unit composed of multiple servers
- Will work as a unique powerful machine
- HPC: High-Performance Computing

## What is a Cluster ?



- A logical unit composed of multiple servers
- Will work as a unique powerful machine
- HPC: High-Performance Computing
- Higher storage capacity

## What is a Cluster ?



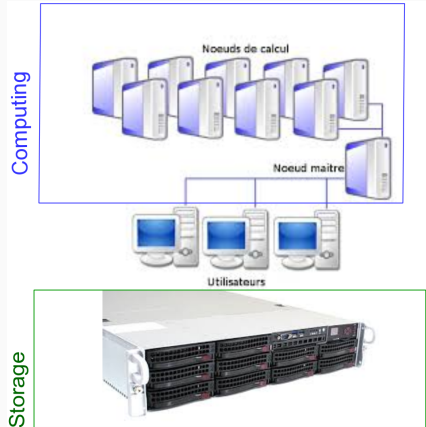
- A logical unit composed of multiple servers
- Will work as a unique powerful machine
- HPC: High-Performance Computing
- Higher storage capacity
- Higher reliability

## What is a Cluster ?

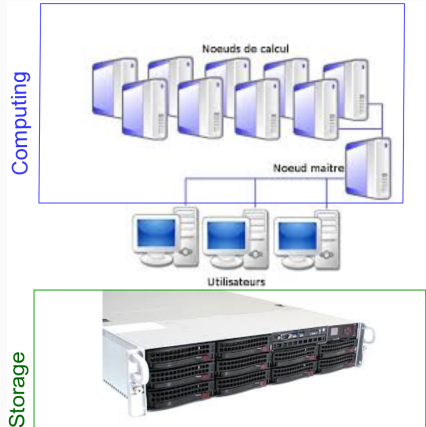


- A logical unit composed of multiple servers
- Will work as a unique powerful machine
- HPC: High-Performance Computing
- Higher storage capacity
- Higher reliability
- Higher resources availability

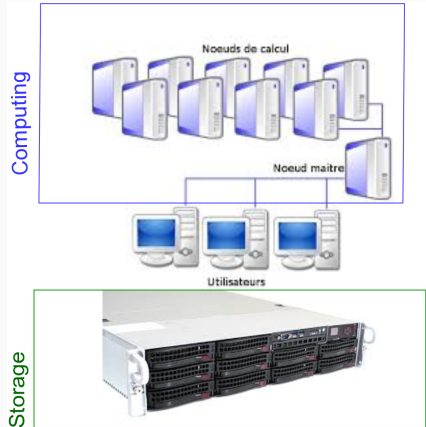
- The Master node:



- The Master node:
  - Users connect on it

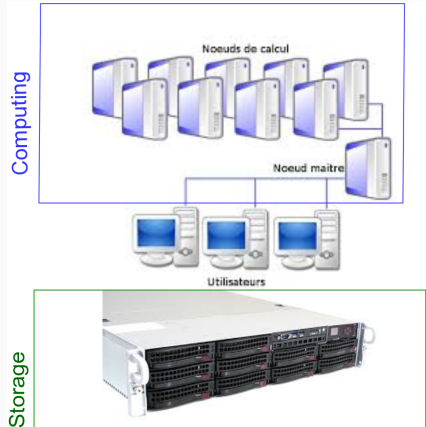


- The Master node:
  - Users connect on it
  - Schedules jobs

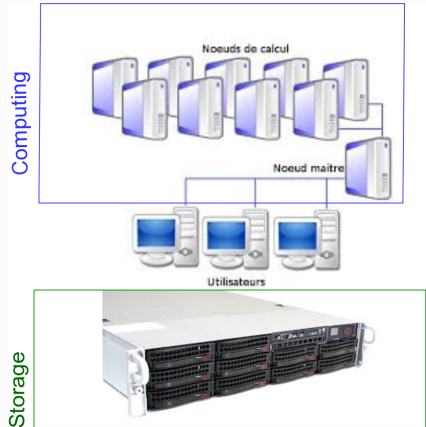




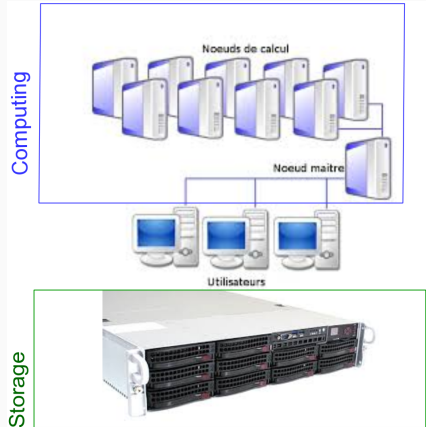
- The Master node:
  - Users connect on it
  - Schedules jobs
  - Manages resources and priorities



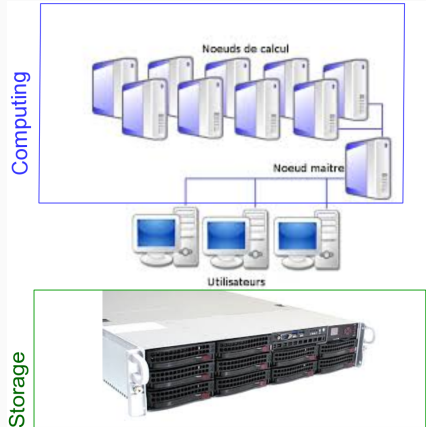
- The Master node:
  - Users connect on it
  - Schedules jobs
  - Manages resources and priorities
- The Computing Nodes:



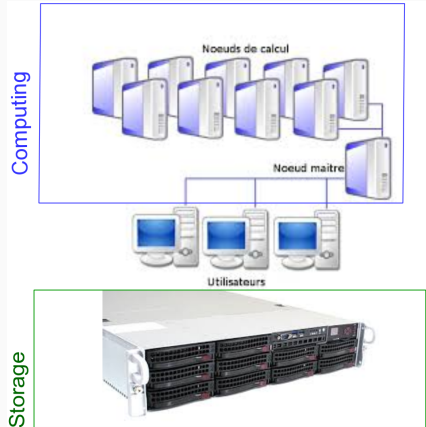
- The Master node:
  - Users connect on it
  - Schedules jobs
  - Manages resources and priorities
- The Computing Nodes:
  - Receive job instructions



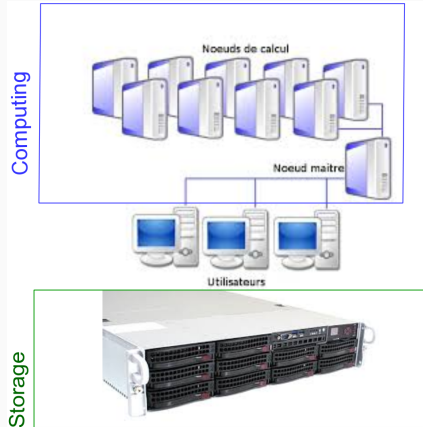
- The Master node:
  - Users connect on it
  - Schedules jobs
  - Manages resources and priorities
- The Computing Nodes:
  - Receive job instructions
  - Perform jobs



- The Master node:
  - Users connect on it
  - Schedules jobs
  - Manages resources and priorities
- The Computing Nodes:
  - Receive job instructions
  - Perform jobs
  - Send to Master the results
- The NAS:



- The Master node:
  - Users connect on it
  - Schedules jobs
  - Manages resources and priorities
- The Computing Nodes:
  - Receive job instructions
  - Perform jobs
  - Send to Master the results
- The NAS:
  - Store data for computing



```
@HWI-EAS236_3_FC_20BTNAAXX:2:1:215:593¶  
GAGAAGTTCAACAGCTGGTATTATTTTTGTTAACAT¶  
+HWI-EAS236_3_FC_20BTNAAXX:2:1:215:593¶  
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhUhhE¶  
@HWI-EAS236_3_FC_20BTNAAXX:2:1:234:551¶  
TGGGACTTTATCTGGAGGAGTGTGGAAAGCCATT¶  
+HWI-EAS236_3_FC_20BTNAAXX:2:1:234:551¶  
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh¶  
@HWI-EAS236_3_FC_20BTNAAXX:2:1:338:194¶  
TGGTTTATGCAGAAATTTCTAGAATAAGGGTAACTT¶  
+HWI-EAS236_3_FC_20BTNAAXX:2:1:338:194¶  
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh¶  
@HWI-EAS236_3_FC_20BTNAAXX:2:1:363:717¶  
TCTCAGAACTTGTTTGTGATGTGTGATTCAACTA¶  
+HWI-EAS236_3_FC_20BTNAAXX:2:1:363:717¶  
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh¶  
@HWI-EAS236_3_FC_20BTNAAXX:2:1:208:209¶  
TTGATTTAACTCTGACAAAATAAACAAAGTCTTAGG¶  
+HWI-EAS236_3_FC_20BTNAAXX:2:1:208:209¶  
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhGh¶
```

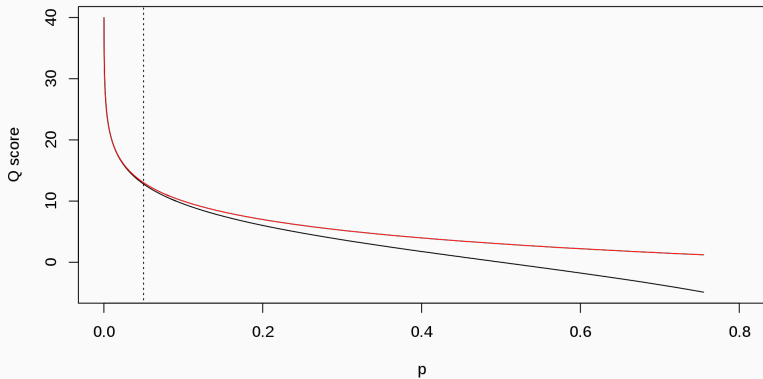
Sequencing info

Nucleotide sequence

Quality score in ASCII







Data are on the CERAAS cluster, located on the storage at  
**/data/FORMATION/2018/TPsnpSV**

1. Connect to the cluster:

```
ssh -X formationX@41.82.52.216
```

Data are on the CERAAS cluster, located on the storage at  
**/data/FORMATION/2018/TPsnpSV**

1. Connect to the cluster:

```
ssh -X formationX@41.82.52.216
```

2. Launch a QRSB command:

```
qrsb
```

Data are on the CERAAS cluster, located on the storage at  
**`/data/FORMATION/2018/TPsnpSV`**

1. Connect to the cluster:

```
ssh -X formationX@41.82.52.216
```

2. Launch a QRSB command:

```
qrsb
```

3. Create your folder in scratch and go in it:

```
mkdir /scratch/formationX  
cd /scratch/formationX
```

Data are on the CERAAS cluster, located on the storage at  
**`/data/FORMATION/2018/TPsnpSV`**

1. Connect to the cluster:

```
ssh -X formationX@41.82.52.216
```

2. Launch a QRSH command:

```
qrsh
```

3. Create your folder in scratch and go in it:

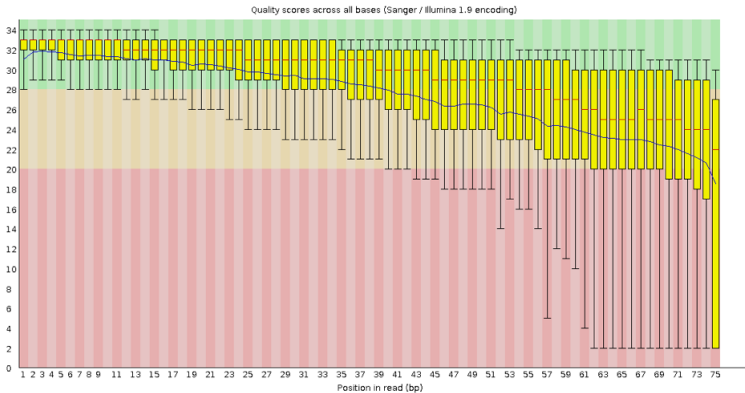
```
mkdir /scratch/formationX
```

```
cd /scratch/formationX
```

4. Transfer the data from nas using SCP:

```
scp -r master:/data/FORMATION/2018/TPsnpSV .
```

## ❌ Per base sequence quality



- Use FastQC on each data  
*fastqc FILE.fastq*

- Use FastQC on each data  
*fastqc FILE.fastq*
- What is the global quality of data ?



- Use FastQC on each data  
*fastqc FILE.fastq*
- What is the global quality of data ?
- For that you have to download it from your scratch to master then your own computer...

## Why Cleaning Data ?

- Removing Adapters and Tags

## Why Cleaning Data ?

- Removing Adapters and Tags
- Removing low quality bases (only conserving 20+ QPHRED value bases)

## Why Cleaning Data ?

- Removing Adapters and Tags
- Removing low quality bases (only conserving 20+ QPHRED value bases)
- Removing contaminants (rRNA genes, organelle data,...)

## Tools for Cleaning Data

- FASTX-tools toolbox

## Tools for Cleaning Data

- FASTX-tools toolbox
- CutAdapt

## Tools for Cleaning Data

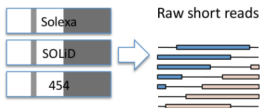
- FASTX-tools toolbox
- CutAdapt
- Trimmomatic

## Tools for Cleaning Data

- FASTX-tools toolbox
- CutAdapt
- Trimmomatic
- Home-made Scripts based on QPHRED scale encoding

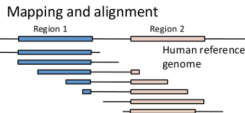


## From unmapped reads to true genetic variation in next-generation sequencing data

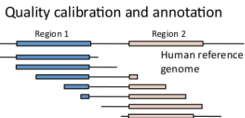


A single run of a sequencer generates  
~50M ~75bp short reads for analysis

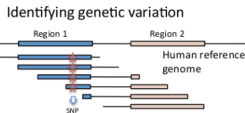
From 1000 genomes projects



The origin of each read from the  
human genome sequence is found

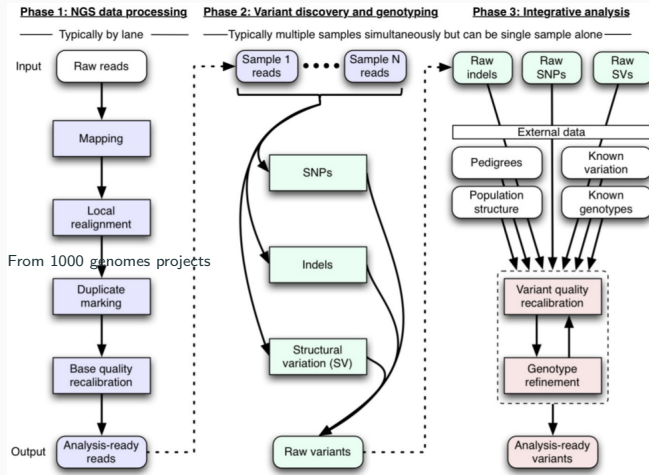


The quality of each read is calibrated  
and additional information annotated  
for downstream analyses



SNPs and indels from the reference  
are found where the reads collectively  
provide evidence of a variant

# Principle of SNP calling



“Classic” launch

1. *Mapping*: bwa aln/sampe, bwa mem, bowtie2, ...

## “Classic” launch

1. *Mapping*: bwa aln/sampe, bwa mem, bowtie2, ...
2. *Cleaning mapping*: samtools, picard-tools,...

## “Classic” launch

1. *Mapping*: bwa aln/sampe, bwa mem, bowtie2, ...
2. *Cleaning mapping*: samtools, picard-tools,...
3. *Realigning and Duplicates*: GATK, picard-tools,...

## “Classic” launch

1. *Mapping*: bwa aln/sampe, bwa mem, bowtie2, ...
2. *Cleaning mapping*: samtools, picard-tools,...
3. *Realigning and Duplicates*: GATK, picard-tools,...
4. *SNP calling and Cleaning*: GATK,...

“Classic” launch

1. *Mapping*: bwa aln/sampe, bwa mem, bowtie2, ...
2. *Cleaning mapping*: samtools, picard-tools,...
3. *Realigning and Duplicates*: GATK, picard-tools,...
4. *SNP calling and Cleaning*: GATK,...

**Between 8 and 15 different commands...**

## Problems with manual launches

- Long
- Fastidious
- Error prone
- Tracability and reproducibility not ensured



## Problems with manual launches

- Long
- Fastidious
- Error prone
- Tracability and reproducibility not ensured

**Solution**  $\implies$  Workflow Manager

# TOGGLE

*Tools for Generic NGS analysis*

A framework to quickly build pipelines  
and to perform large-scale NGS analysis

**Christine Tranchant-Dubreuil**  
[christine.tranchant@ird.fr](mailto:christine.tranchant@ird.fr)

# TOGGLE



- A toolbox to perform large-scale NGS analyses

19 modules, 120 functions  
40 open-source tools



# TOGGLE



- A toolbox to perform large-scale NGS analyses

19 modules, 120 functions  
40 open-source tools



- **Version 2 published in BMC bioinformatics**

RESEARCH

TOGGLE: Toolbox for generic NGS analyses

Cecile Monat<sup>1\*</sup>, Christine Tranchant-Dubreuil<sup>1\*</sup>, Ayité Kougbéadjò<sup>2</sup>, Cédric Farcy<sup>2</sup>, Mawussé Agbessi<sup>1</sup>, Maryline Summo<sup>2</sup> and François Sabot<sup>1\*</sup>

## Data preprocessing

Fastqc, Cutadapt  
FastxTrimmer  
Stack process\_radstats

## Structural Variations

MindTheGap,  
BreakDancer, Pindel

## RNA-seq Assembly

Trinity  
TGI-CL

# TOGGLE

## ReadCount

Htseq-count

## Mapping

Bwa aln, sampe/ samse  
Bwa mem  
Tophat2



## SNP calling/ filtering

SAMtools, GATK, VarScan, SNPEff

## SAM/BAM management

picardTools, SAMtools, GATK

## Data preprocessing

Fastqc, Cutadapt  
FastxTrimmer  
Stack process\_radstats

## Structural Variations

MindTheGap,  
BreakDancer, Pindel

## RNA-seq Assembly

TransAbyss, Trinity  
TGI-CL

## ReadCount

Htseq-count

# TOGGLE



## Mapping

Bwa aln, sampe/ samse  
Bwa mem  
Tophat2

## SNP calling/ filtering

SAMtools, GATK, VarScan, SNPEff

## SAM/BAM management

picardTools, SAMtools, GATK



<https://github.com/SouthGreenPlatform/TOGGLE>

# TOGGLE



Version 2



Version 3

*From hard-coded pipelines  
To a bioinformatic pipeline framework*

# TOGGLE



Version 2



Version 3

*From hard-coded pipelines  
To a bioinformatic pipeline framework*

Biologists create their own pipeline through an easy and user-friendly approach





# How to perform an analysis with TOGGLE ?

A command-line based pipeline framework



A single command line

```
toggleGenerator.pl -d DIR-c FILE -o DIR
```

# What does TOGGLE need to run ?

- An input directory (with fastq, sam/bam, vcf files)
- The name of output directory used to store the data generated by the analyses
- A unique and simple configuration file to design the pipeline and define software parameters.
- Optional arguments : reference file, annotation...

**\$order**

1=fastqc

2=cutadapt

3=bwa mem

4=picardToolsSortSam

5=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration

**\$cutadapt**

-q 30

-m 35

**\$bwa mem**

-n 5

...

**\$sge**

-q bioinfo.q

-b Y

## \$order

```
1=fastqc  
2=cutadapt  
3=bwa mem  
4=picardToolsSortSam  
5=samToolsView  
1000=gatkHaplotypeCaller  
1001=gatkVariantFiltration
```

## Create your own workflow

- The workflow order
- The list of softwares to run

One line = the step followed by the software's name

**\$order**

1=fastqc

2=cutadapt

3=bwa mem

4=picardToolsSortSam

5=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration

**Create your own workflow**

Step number &lt; 1000

Step number &gt;= 1000

```
$order
```

```
1=fastqc
```

```
2=cutadapt
```

```
3=bwa mem
```

```
4=picardToolsSortSam
```

```
5=samToolsView
```

```
1000=gatkHaplotypeCaller
```

```
1001=gatkVariantFiltration
```

Create your own workflow

Step number < 1000

Parallel analysis by sample

**\$order**

1=fastqc

2=cutadapt

3=bwa mem

4=picardToolsSortSam

5=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration



Input directory



Sample 1

Sample 2

Sample 3

fastQC

cutadapt

bwa mem

pT sortSam

sT view

fastQC

cutadapt

bwa mem

pT sortSam

sT view

fastQC

cutadapt

bwa mem

pT sortSam

sT view

**\$order**

1=fastqc

2=cutadapt

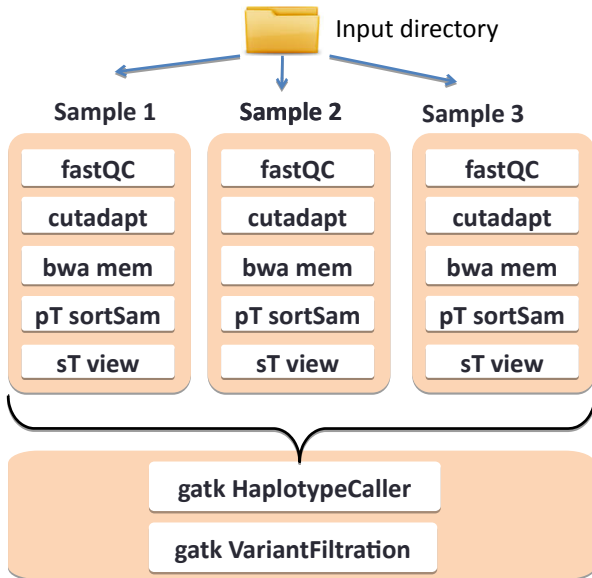
3=bwa mem

4=picardToolsSortSam

5=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration

Step number  $\geq$  1000**Global analysis  
(all samples)**



**\$order**

1=fastqc

2=cutadapt

3=bwa mem

4=picardToolsSortSam

5=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration

**\$cutadapt**

-q 30

-m 35

**\$bwa mem**

-n 5

...

**\$sge**

-q bioinfo.q

-b Y

**Software parameters**

One tag per software (\$softwareName)  
followed by the list of options

- Create a **TOGGLE** configuration file with as order:
  1. bwa aln
  2. bwa sampe
  3. picardtools sortsam
  4. samtools view
  5. samtools index
  6. gatkindelrealigner
  7. picardtools markduplicates

- Create a **TOGGLE** configuration file with as order:
  1. bwa aln
  2. bwa sampe
  3. picardtools sortsam
  4. samtools view
  5. samtools index
  6. gatkindelrealigner
  7. picardtools markduplicates
- Add the options...

- Create a **TOGGLE** configuration file with as order:
  1. bwa aln
  2. bwa sampe
  3. picardtools sortsam
  4. samtools view
  5. samtools index
  6. gatkindelrealigner
  7. picardtools markduplicates
- Add the options...
- Ok, we help... Look to configuration files in */data/FORMATION/2018/TPsnpSV*

- Download Tablet (use Google and Tablet+NGS)

- Download Tablet (**use Google and Tablet+NGS**)
- Transfer the BAMs and the reference from the node to the master then to your local computer (use scp at each point)

- Download Tablet (**use Google and Tablet+NGS**)
- Transfer the BAMs and the reference from the node to the master then to your local computer (use scp at each point)
- Open Tablet, look at the mapping and try to find SNPs

- Pick up all last BAM file (*MarkDuplicates* steps)



- Pick up all last BAM file (*MarkDuplicates* steps)
- Launch a **TOGGLE** with *1000=gatkUnifiedGenotyper*

- Pick up all last BAM file (*MarkDuplicates* steps)
- Launch a **TOGGLE** with *1000=gatkUnifiedGenotyper*
- Try to launch a **TOGGLE** from *FASTQ* to *VCF*

SAM format : <http://samtools.sourceforge.net/samtools.shtml>

Type	Tag	Description
HD - header	VN*	File format version.
	SO	Sort order. Valid values are: <i>unsorted</i> , <i>queryname</i> or <i>coordinate</i> .
	GO	Group order (full sorting is not imposed in a group). Valid values are: <i>none</i> , <i>query</i> or <i>reference</i> .
SQ - Sequence dictionary	SN*	Sequence name. Unique among all sequence records in the file. The value of this field is used in alignment records.
	LN*	Sequence length.
	AS	Genome assembly identifier. Refers to the reference genome assembly in an unambiguous form. Example: HG18.
	M5	MD5 checksum of the sequence in the uppercase (gaps and space are removed)
	UR	URI of the sequence
	SP	Species.
RG - read group	ID*	Unique read group identifier. The value of the ID field is used in the RG tags of alignment records.
	SM*	Sample (use pool name where a pool is being sequenced)
	LB	Library
	DS	Description
	PU	Platform unit (e.g. lane for Illumina or slide for SOLiD); should be a full, unambiguous identifier
	PI	Predicted median insert size (maybe different from the actual median insert size)
	CN	Name of sequencing center producing the read.
	DT	Date the run was produced (ISO 8601 date or date/time).
PG - Program	PL	Platform/technology used to produce the read.
	ID*	Program name
	VN	Program version
CO - comment	CL	Command line
		One-line text comments

SAM format : <http://samtools.sourceforge.net/samtools.shtml>

Type	Tag	Description
HD - header	VN*	File format version.
	SO	Sort order. Valid values are: <i>unsorted</i> , <i>queryname</i> or <i>coordinate</i> .
	GO	Group order (full sorting is not imposed in a group). Valid values are: <i>none</i> , <i>query</i> or <i>reference</i> .
SQ - Sequence dictionary	SN*	Sequence name. Unique among all sequence records in the file. The value of this field is used in alignment records.
	LN*	Sequence length.
	AS	Genome assembly identifier. Refers to the reference genome assembly in an unambiguous form. Example: HG18.
	M5	MD5 checksum of the sequence in the uppercase (gaps and space are removed)
	UR	URI of the sequence
	SP	Species.
RG - read group	ID*	Unique read group identifier. The value of the ID field is used in the RG tags of alignment records.
	SM*	Sample (use pool name where a pool is being sequenced)
	LB	Library
	DS	Description
	PU	Platform unit (e.g. lane for Illumina or slide for SOLiD); should be a full, unambiguous identifier
	PI	Predicted mapping quality
	CN	Name of sequencing center
PG - Program	DT	Date the run was made
	PL	Platform/technology
	ID*	Program name
	VN	Program version
CO - comment	CL	Command line
		One-line text comment

```

@HD VN:1.3 SO:coordinate
@SQ SN:ref LN:45
r001 163 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5H6M * 0 0 AGCTAA * NM:i:1
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTCAGC *
r003 16 ref 29 30 6H5M * 0 0 TAGGC * NM:i:0
r001 83 ref 37 30 9M = 7 -39 CAGCGCCAT *
```

SAM format : <http://samtools.sourceforge.net/samtools.shtml>

Col	Name	Description
1	<b>QNAME</b>	Query NAME of the read or the read pair
2	<b>FLAG</b>	bitwise FLAG (pairing, strand, mate strand, etc.)
3	<b>RNAME</b>	Reference sequence NAME
4	<b>POS</b>	1-based leftmost POSition of clipped alignment
5	<b>MAPQ</b>	MAPping Quality (Phred-scaled)
6	<b>CIGAR</b>	extended CIGAR string (operations: MIDNSHP)
7	<b>NRNM</b>	Mate Reference NaMe ('=' if same as RNAME)
8	<b>MPOS</b>	1-based leftmost Mate POSition
9	<b>ISIZE</b>	inferred Insert SIZE
10	<b>SEQ</b>	query SEQUENCE on the reference
11	<b>QUAL</b>	query QUALITY (ASCII-33)

```

@HD VN:1.3 SO:coordinate
@SQ SN:ref LN:45
r001 163 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5H6M * 0 0 AGCTAA * NM:i:1
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 16 ref 29 30 6H5M * 0 0 TAGGC * NM:i:0
r001 83 ref 37 30 9M = 7 -39 CAGCGCCAT *
```

**SAM format :** <http://samtools.sourceforge.net/samtools.shtml>

Suppose we have the following alignment with bases in lower cases clipped from the alignment. Read r001/1 and r001/2 constitute a read pair; r003 is a chimeric read; r004 represents a split alignment.

```

Coor      12345678901234  5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGGCCAT

+r001/1           TTAGATAAAGGATA*CTG
+r002             aaaAGATAA*GGATA
+r003             gcctaAGCTAA
+r004             ATAGCT.....TCAGC
-r003             ttagctTAGGC
-r001/2                               CAGGCCAT
  
```

**SAM format :** <http://samtools.sourceforge.net/samtools.shtml>

Suppose we have the following alignment with bases in lower cases clipped from the alignment. Read r001/1 and r001/2 constitute a read pair; r003 is a chimeric read; r004 represents a split alignment.

```

Coord      12345678901234  5678901234567890123456789012345
ref        AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGGCCAT

+r001/1      TTAGATAAAGGATA*CTG
+r002        aaaAGATAA*GGATA
+r003        gcctaAGCTAA
+r004                ATAGCT.....TCAGC
-r003                ttagctTAGGC
-r001/2                CAGGCCAT
  
```

The corresponding SAM format is:

```

@HD VN:1.3 SO:coordinate
@SQ SN:ref LN:45
r001 163 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAAGGATA *
r003 0 ref 9 30 5H6M * 0 0 AGCTAA * NM:i:1
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 16 ref 29 30 6H5M * 0 0 TAGGC * NM:i:0
r001 83 ref 37 30 9M = 7 -39 CAGGCCAT *
  
```

## SAM format: FLAG field

numeric	binary	description
1	00000001	template has multiple fragments in sequencing
2	00000010	each fragment properly mapped according to aligner
4	00000100	fragment is unmapped
8	00001000	mate is unmapped
16	00010000	sequence is reverse complemented
32	00100000	sequence of mate is reversed
64	01000000	is first fragment in template
128	10000000	is second fragment in template

From J. Aerts, in SlideShare



## The Variant Call Format (VCF) used in bioinformatics for storing gene sequence variations

**VCF header**

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

**Mandatory header lines** (indicated by a red arrow pointing to ##fileformat=VCFv4.0)

**Optional header lines** (meta-data about the annotations in the VCF body) (indicated by a red arrow pointing to the remaining header lines)

**Body**

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

**Reference alleles (GT=0)** (indicated by a blue arrow pointing to the REF column)

**Alternate alleles (GT>0 is an index to the ALT column)** (indicated by a blue arrow pointing to the ALT column)

**Phased data** (G and C above are on the same chromosome) (indicated by a blue arrow pointing to the 0|1 in the GQ field)

**Deletion** (indicated by a blue arrow pointing to the <DEL> in the ALT column)

**SNP** (indicated by a blue arrow pointing to the C/T variant)

**Large SV** (indicated by a blue arrow pointing to the SVTYPE=DEL)

**Insertion** (indicated by a blue arrow pointing to the A/G variant)

**Other event** (indicated by a blue arrow pointing to the H2;AA=T in the INFO field)

```
##fileformat=VCFv4.1
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,m5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA000001 NA000002 NA000003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:..
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ 0|0:49:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTC G,GTCT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP 0/1:35:4 0/2:17:2 1/1:40:3'
```

- **Variation 1** : a good SNP
- **Variation 2** : a possible SNP that has been filtered out because its quality is below 10
- **Variation 3** : a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error)
- **Variation 4** : a site that is called monomorphic reference (i.e. with no alternate alleles)
- **Variation 5** : a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T).

- Open Tablet, load BAM, reference and VCF (pick it on your local computer)

- Open Tablet, load BAM, reference and VCF (pick it on your local computer)
- Look for SNP and their mapping

- Open Tablet, load BAM, reference and VCF (pick it on your local computer)
- Look for SNP and their mapping
- Relaunch the same analyses removing one step or changing parameters...

- Using GATK Variant Filtration, a flag per filter

- Using GATK Variant Filtration, a flag per filter
- Depth filter:  
*DP<10 or DP>20000*

- Using GATK Variant Filtration, a flag per filter
- Depth filter:  
 *$DP < 10$  or  $DP > 20000$*
- MQ0 filter:  
 *$MQ0 < 4$  or  $MQ0 < 0.1 DP$*



- Using GATK Variant Filtration, a flag per filter
- Depth filter:  
 *$DP < 10$  or  $DP > 20000$*
- MQ0 filter:  
 *$MQ0 < 4$  or  $MQ0 < 0.1 DP$*
- QUAL filter:  
 *$QUAL < 200$*

- Using GATK Variant Filtration, a flag per filter
- Depth filter:  
 *$DP < 10$  or  $DP > 20000$*
- MQ0 filter:  
 *$MQ0 < 4$  or  $MQ0 < 0.1 DP$*
- QUAL filter:  
 *$QUAL < 200$*
- SNPcluster filter:  
*more than 3 SNP per 10b*

The command is then...

```
java -jar GenomeAnalysisTK.jar -T VariantFiltration  
-R REFERENCE.FASTA -V INPUT.VCF --filterExpression  
'QUAL<200' --filterName 'LOW-QUAL' --filterExpression  
'MQ0>=4 && ((MQ0/(1.0DP)) > 0.1)' --filterName  
'LOW-MQ0' --filterExpression 'DP<10' --filterName  
'LOWDP' --clusterSize 3 --clusterWindowSize 10 --  
filterExpression 'DP>20000' --filterName 'HIGH-DP' -o  
OUTPUT.VCF
```

- Recoding with SNP PASS:

```
vcftools --vcf FILEIN --remove-filtered-all  
--recode --recode-INFO-all --out FILEOUT
```

- Recoding with SNP PASS:

```
vcftools --vcf FILEIN --remove-filtered-all  
--recode --recode-INFO-all --out FILEOUT
```

- Missing data filtering:

```
vcftools --vcf FILEIN --max-missing-count 12  
--recode --recode-INFO-all --out FILEOUT
```

- Recoding with SNP PASS:

```
vcftools --vcf FILEIN --remove-filtered-all  
--recode --recode-INFO-all --out FILEOUT
```

- Missing data filtering:

```
vcftools --vcf FILEIN --max-missing-count 12  
--recode --recode-INFO-all --out FILEOUT
```

- MAF:

```
vcftools --vcf FILEIN --output FILEOUT --maf 0.1
```

- Recoding with SNP PASS:

```
vcftools --vcf FILEIN --remove-filtered-all  
--recode --recode-INFO-all --out FILEOUT
```

- Missing data filtering:

```
vcftools --vcf FILEIN --max-missing-count 12  
--recode --recode-INFO-all --out FILEOUT
```

- MAF:

```
vcftools --vcf FILEIN --output FILEOUT --maf 0.1
```

- Biallelic SNPs:

```
vcftools --vcf FILEIN --min-alleles 2  
--max-alleles 2 --recode --recode-INFO-all --out  
FILEOUT
```

- *sNMF*: tool to estimate ancestry coefficients



- *sNMF*: tool to estimate ancestry coefficients
- Developed by E. Frichot and O. Francois, TIMC-IMAG

- *sNMF*: tool to estimate ancestry coefficients
- Developed by E. Frichot and O. Francois, TIMC-IMAG
- R and Command-line version

- *sNMF*: tool to estimate ancestry coefficients
- Developed by E. Frichot and O. Francois, TIMC-IMAG
- R and Command-line version
- Much faster than ADMIXTURE or STRUCTURE, as efficient

In your QRSB session, copy the files from  
/data/FORMATION/2018/TpPop in your /scratch folder, and  
follow these commands:

1. `scp master:/data/FORMATION/2018/TpPop/sample.vcf`

.

In your QRSB session, copy the files from /data/FORMATION/2018/TpPop in your /scratch folder, and follow these commands:

1. `scp master:/data/FORMATION/2018/TpPop/sample.vcf`  
`.`
2. `grep "#" sample.vcf > subsamples.vcf`

In your QRSB session, copy the files from /data/FORMATION/2018/TpPop in your /scratch folder, and follow these commands:

1. `scp master:/data/FORMATION/2018/TpPop/sample.vcf`  
`.`
2. `grep "#" sample.vcf > subsamples.vcf`
3. `shuf -n 5000 sample.vcf | grep -v "#" > tmp.vcf`

In your QRSB session, copy the files from /data/FORMATION/2018/TpPop in your /scratch folder, and follow these commands:

1. `scp master:/data/FORMATION/2018/TpPop/sample.vcf`  
`.`
2. `grep "#" sample.vcf > subsamples.vcf`
3. `shuf -n 5000 sample.vcf | grep -v "#" > tmp.vcf`
4. `cat tmp.vcf » subsamples.vcf`

The VCF is now submapped and ready for pop analysis

1. `vcf2geno subsamples.vcf test.geno`



The VCF is now submapped and ready for pop analysis

1. `vcf2geno subsamples.vcf test.geno`
2. `for K in {1..10}; do echo K=$K; sNMF -x test.geno  
-K $K -c > test.$K.log; done`

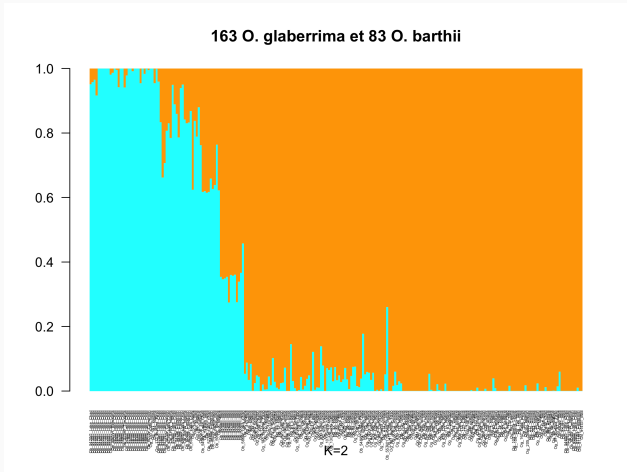
The VCF is now submapped and ready for pop analysis

1. `vcf2geno subsamples.vcf test.geno`
2. `for K in {1..10}; do echo K=$K; sNMF -x test.geno  
-K $K -c > test.$K.log; done`
3. `grep "Cross-Entropy (masked data):" test.*.log`

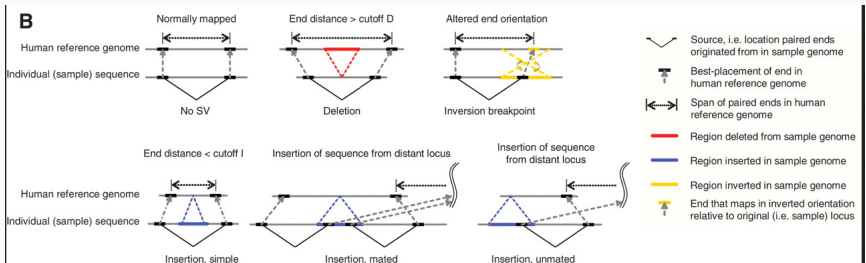
The VCF is now submapped and ready for pop analysis

1. `vcf2geno subsamples.vcf test.geno`
2. `for K in {1..10}; do echo K=$K; sNMF -x test.geno  
-K $K -c > test.$K.log; done`
3. `grep "Cross-Entropy (masked data):" test.*.log`
4. Look at all files for the best Cross-Entropy value

Use this with e.g. CLUMPP to extract the figure



From Cubry et al, 2018



From Korbelt et al, 2007

We can use **BreakDancer** (*Chen et al, 2009, Nature Methods*)

**BE CAREFUL: if you do it, do it through QRSH and scp copy on /scratch!!**

- From the raw BAMs, launch:

```
/usr/local/breakdancer-1.1.2/perl/bam2cfg.pl *.bam >  
bamConf.cfg
```

We can use **BreakDancer** (*Chen et al, 2009, Nature Methods*)

**BE CAREFUL: if you do it, do it through QRSH and scp copy on /scratch!!**

- From the raw BAMs, launch:

```
/usr/local/breakdancer-1.1.2/perl/bam2cfg.pl *.bam >  
bamConf.cfg
```

- Launch the **BreakDancer** command:

```
/usr/local/breakdancer-1.1.2/cpp/breakdancer-max  
bamConf.cfg > breakDancer.out
```

We can use **BreakDancer** (*Chen et al, 2009, Nature Methods*)

**BE CAREFUL: if you do it, do it through QRSH and scp copy on /scratch!!**

- From the raw BAMs, launch:

```
/usr/local/breakdancer-1.1.2/perl/bam2cfg.pl *.bam >  
bamConf.cfg
```

- Launch the **BreakDancer** command:

```
/usr/local/breakdancer-1.1.2/cpp/breakdancer-max  
bamConf.cfg > breakDancer.out
```

- Check the results...



## 1. Fragment DNA and sequence

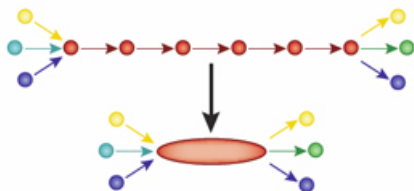


## 2. Find overlaps between reads

```
...AGCCTAGACCTACAGGATGCGCGACACGT  
                GGATGCGCGACACGTCGCATATCCGGT...
```

From Baker, 2012

### 3. Assemble overlaps into contigs



### 4. Assemble contigs into scaffolds



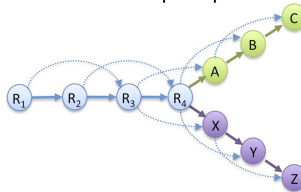
*Michael Schatz, Cold Spring Harbor*

From Baker, 2012

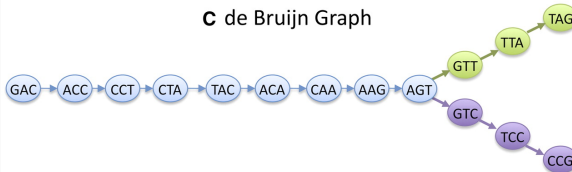
## A Read Layout

R<sub>1</sub>: GACCTACA  
 R<sub>2</sub>: ACCTACAA  
 R<sub>3</sub>: CCTACAAG  
 R<sub>4</sub>: CTACAAGT  
 A: TACAAGTT  
 B: ACAAGTTA  
 C: CAAGTTAG  
 X: TACAAGTC  
 Y: ACAAGTCC  
 Z: CAAGTCCG

## B Overlap Graph



## C de Bruijn Graph



From Schatz, 2010

1. Download using *scp* the folder *TPassembly* in  
*/data/FORMATION/2018/TPassembly*

1. Download using *scp* the folder *TPassembly* in */data/FORMATION/2018/TPassembly*
2. Go in this folder and check the data in the *Ebola* folder.  
What are they ?

1. Download using *scp* the folder *TPassembly* in */data/FORMATION/2018/TPassembly*
2. Go in this folder and check the data in the *Ebola* folder.  
What are they ?
3. Check them using *FASTQC*

1. Download using *scp* the folder *TPassembly* in */data/FORMATION/2018/TPassembly*
2. Go in this folder and check the data in the *Ebola* folder.  
What are they ?
3. Check them using *FASTQC*
4. Use the *Abyss* tool to assemble the pair-end data

1. Download using **scp** the folder **TPassembly** in ***/data/FORMATION/2018/TPassembly***
2. Go in this folder and check the data in the **Ebola** folder.  
What are they ?
3. Check them using **FASTQC**
4. Use the **Abyss** tool to assemble the pair-end data
5. Use the **BLAST** tool to compare it to the **ebola.fasta** file as database



1. Download using *scp* the folder *TPassembly* in */data/FORMATION/2018/TPassembly*
2. Go in this folder and check the data in the *Ebola* folder.  
What are they ?
3. Check them using *FASTQC*
4. Use the *Abyss* tool to assemble the pair-end data
5. Use the *BLAST* tool to compare it to the *ebola.fasta* file as database
6. Use the *MUSCLE* tool to align your assembly to the *ebola.fasta*

Thanks for your attention

