# Projet Burkinabioinfo

# Projet burkinabioinfo



2018 - CERAAS
2019 - UJKZ

2018 - Formation niveau 1 - CERAAS
2019 - Formation niveau 1 & 2 - UJKZ

2018 - Formation niveau 1 CERAAS
2019 - Formation niveau 1 - UJKZ

2019 - 2020 - 2021

2020 - 2021
Formation niveau 1 - UJKZ
Formation niveau 2 - UJKZ

**Installation de clusters locaux** → **Formation Administration de cluster de calcul** → **Formation bioinformatique** → **groupes de travail thématiques projets tutorés** → **Formations bioinformatiques**

Administrateurs d'Afrique de l'Ouest (Burkina Faso, Côte d'Ivoire, Sénégal)

Chercheurs d'Afrique de l'Ouest (Burkina Faso, Bénin, Côte d'Ivoire, Mali, Sénégal)

## A Moyen / long terme

- Clusters installés localement
- Réseau d'experts en administration de cluster
- Réseau de chercheurs experts en bioinformatique
- Intégration de modules de formation progressivement dans les cursus universitaires
- Mise en ligne des supports et cours, développement de cours accessibles en ligne

**13 apprenants**

**Objectifs :** Initiation à la bioinformatique avec 2 cas pratique détection de variants structurants
- SNP : données short reads illumina
  - Détections SNP > post Analyses
- Génomique comparatique : données long reads
  - Détections de SNPs et de variants structuraux
  - Assemblage et Comparaison de génomes assemblés

→ linux, Jupyter book, bash, python, protocoles bioinfos et bioanalyses

**9 apprenants**

**Objectifs :** 2 projets à mener
- Métagénomique à partir d'un échantillon prélevé dans un champs (virus, bactérie)
- Détection de SNPs à partir d'échantillons d'Ignames séquencés

→ Terminal Linux, Cluster UZB

# Programme

| | Débutant | Veteran SNP / diversite plantes | Veteran Metagenomique | | |
|---|---|---|---|---|---|
| **LUNDI** | | | | | **COMMUN** |
| 8h30-10h00 | Accueil + Presentation formation / etudiants babies / etudiants veterans | | | | **COURS** |
| 10h30-12h00 | Cours NGS | Description jeu de données | Description jeu de données | | Autonomie |
| | Pause déjeuner | | | | |
| 14h00-15h30 | Cours | Autonomie / Bibliographie / Veille technologique | | | |
| 16h-17h30 | Contrôle qualité + mapping | | | | |
| | | | | | |
| **MARDI** | | | | | |
| 8h-10h00 | Autonomie | Accompagnement mini-projets | Accompagnement mini-projets | | |
| 10h30-12h | Contrôle qualité + mapping | | | | |
| | Pause déjeuner | | | | |
| 14h00-15h30 | Cours | Autonomie | Autonomie | | |
| 16h-17h30 | mapping + SNP calling | | | | |
| | | | | | |
| **MERCREDI** | | | | | |
| 9h-10h30 | Autonomie | Accompagnement mini-projets | Accompagnement mini-projets | | |
| 10h30-12h | mapping + SNP calling | | | | |
| | Pause déjeuner | | | | |
| 14h00-15h30 | Cours | Autonomie | Autonomie | | |
| 16h-17h30 | Analyse diversité | | | | |
| | | | | | |
| **JEUDI** | | | | | |
| 9h-10h30 | Autonomie | Accompagnement mini-projets | Accompagnement mini-projets | | |
| 10h30-12h | Analyse diversité | | | | |
| | Pause déjeuner | | | | |
| 14h00-15h30 | Cours | Autonomie | Autonomie | | |
| 16h-17h30 | Génomique comparative | | | | |
| | | | | | |
| **VENDREDI** | | | | | |
| 9h-10h30 | Restitution des mini-projets Veterans - Questions et discussions diverses | | | | |
| 10h30-12h | Questions et discussions autour des projets/données des participants | | | | |
| | Pause déjeuner | | | | |
| 14h00-15h30 | Questions et discussions autour des projets/données des participants | | | | |
| 16h-17h30 | | | | | |

KIENDREBEOGO Touwendpoulimdé Isabelle

Génétique, Bm, Cancer du sein, Mutation, BRCA

AHONON Awovi Selom

S. rotundifolius, S. stenocarpa, plantes mineures, diversité

GBEKLEY Efui Holaly

TUINA Sévérin

Flux de gènes, Dynamique de la diversité, Sorghum bicolor,

DOSSIM Sika

BA Aminata Hamidou

Diversité S. rotundifolius: Profil morphologique/ génétique des morphotypes cultivés ( BF, Ghana)

TONDE Ignace

Solenostemon rotundifolius, interactions génotype x environnement , profil génitique,

SIRIMA Constant

RNA Seq, Plasmodium falciparum ACT-sensible/ ACT-résistant

BADOUM Emilie Salimata

ZOURE Abdou Azaque

Microbiome instestianal du moustique (Illumina) , Gène BRCA (Sanger)

ADAMOU IBRAHIM Maman Laouali

Analyse de la distribution génétiques et des régions liées au sexe des palmiers du Sahel

13 Apprenants

SANOU Estèle Pélagie

tilapia du Nil, déterminisme du sexe, contrôle du sexe, population monosexe mâle, marqueurs chromosomiques

PALANGA Essowè

Metagenomique, virus, interaction plante-parasite, phytopathologie

**OUEDRAOGO Jacques**

**DANOU-KODJO Kodjovi Atassé**

Métagénomique-Variabilité génétique-Phytopathologie

**SAGNON Adama**

phosphate, solubilisation, bactéries, champignons

**SORY Siedou**

Diversité génétique/biochimique des cultivars d'ignames cultivées au Burkina Faso.

**NAME Pakyendou Estel**

Epidémiologie; Virus; ADN; CRESS; Séquençage

**ZONGO Saïdou**

Oxford Nanopore Technologie, Séquençage, Geminivirus, longs reads, NGS

**LALLOGO P. Doriane Tatiana**

SARS-CoV 2, facteurs génétiques, clairance, l'hôte humain, formes sévères.

**SAWADOGO Seydou**

Surveillance participative, maladies virales, racines et tubercules, séquençage

*8 Apprenants*

**Dereeper Alexis**

Interaction plantes-pathogène, pangénomique des
xanthomonas, diversité du riz, Kite surfeur

PHIM | IRD Institut de Recherche pour le Développement FRANCE

**Tibiri B. Ezéchiel**

Interaction plantes-pathogène,

UNIVERSITE DE OUAGADOUGOU BURKINA FASO POPULI SAPIENTIA POPULO | INERA

**Orjuela-Bouniol Julie**

Assemblages de génomes, annotation, diversité, métagénomique
Développement des méthodes pour l'analyse de la diversité des plantes sans
référence

diade | IRD Institut de Recherche pour le Développement FRANCE

**Tranchant-Dubreuil
Christine**

diversité des riz africains, mécanismes
d'adaptation et de sélection, pangénomique,
variants structuraux (SNP mais pas que)

diade | IRD Institut de Recherche pour le Développement FRANCE

**Brunel Dominque
Centre Nationale de
Génotypage - INRAE**

2 formations    …

2 ambiances

### Mode "training"

- Session cours suivi par
- Session pratique en autonomie (individuel ou en groupe)
- Correction en groupe

### Mode "projet"

- brainstroming en groupe, avec les formateurs
- projet en autonomie…
- debriefing collectif
- 2 projets en parallèle !

Des données différentes pour les 2 groupes avec des analyses différentes !!!

Apprendre à réaliser une analyse bioinformatique
Avoir un oeil critique lors de la (bio)analyse des données
Maîtriser linux, les outils bioinformatiques

**auton**

**omie**

- Nos Adminsys burkinabé : Ousmane Barra, Seydou Konsimbo, Ndomassi Tando… Ezéchiel

- Le comité d'organisation : Ezéchiel, Fidéle Tiendrebeogo, Romaric Nanema, Isidore Boungoungou…

- Toutes nos tutelles, l'université JZK

- Le LMI Patho Bios : James Neya, Charlotte Tollenaere et Christophe Brugidou

# Introduction
# Bioinformatics & Sequencing

**A interdisciplanary science**

**1866**

**Lois de l'hérédité**

**1944**

**Nature chimique de l'ADN, matériel héréditaire**

*O. Avery, C. McLeod & McCarthy*

**1954**

**Structure en double hélice de l'ADN**

*J. Watson & F. Cricks & franklin*

**1961**

**Code génétique et règle de correspondance gènes-protéines**

*M. Nirenberg & H. Matthaei*

**1965**

**Mécanismes de la régulation génétique**

*André Lwoff, F. Jacob & J. Monod*

**1970**

**Algo *Alignement global de séquence***

*Needman, & Wunsh*

**1972**
**8008**

**1er microprocesseur intel**

**1977**

**Micro-ordinateurs**

**Séquençage ADN**

*P. Berg, W. Gilbert & F. Sanger*

**1980**

**Banque EMBL, GenBank, PIR**

**1984**

**Amplification ADN - PCR**

*Karry Mullis*

**1985**

**Algo Alignement local de séquence FASTA**

*Person & Lipman*

**1987**

**1er séquenceur automaisé**

*L. Hood Société Applied Biosystems*

**1990**

**Algo Alignement local de séquence BLAST**

*Altschul & al.*

1 - How many base pairs (bp) are there in a human genome?

2 - How much did it cost to sequence the first human genome?

3 - How long did it take to sequence the first human genome?

1 - How many base pairs (bp) are there in a human genome?

**3 billion**

2 - How much did it cost to sequence the first human genome?

**2.7 billions**

3 - How long did it take to sequence the first human genome?

**> 10 years**

*DNA sequencing : determining the order of the four bases or nucleotides that make up a given molecule of DNA

# A little history of sequencing...

Next-generation sequencing

| Virus phiX174 | H. influenza | S. cerevisiae | C. elegans | A. thaliana | H. sapiens | O. sativa | |
|---|---|---|---|---|---|---|---|
| **1977** | **1995** | **1996** | **1998** | **2000** | **2001** | **2002** | **2008** |
| 5kb<br>11 | 1,8 Mb<br>1,740 | 12 Mb<br>300 | 100 Mb<br>22,000 | 125 Mb<br>~30,000 | 3,3 Gb<br>~21,000 | 289 Mb<br>~35,000 | |

**1** FIRST GENERATION
From 1977

sanger

SEQUENCING TECHNOLOGY

**2** SECOND GENERATION
From 2007

solexa, 454, illumina

**3** THIRD GENERATION
From 2011

PacBio, oxford nanopore

**SEQUENCING TECHNOLOGY**

**1 FIRST GENERATION**
From 1977

sanger

**2 SECOND GENERATION**
From 2007

solexa, 454, illumina

**3 THIRD GENERATION**
From 2011

PacBio, oxford nanopore

➡ Sequencing output, price, reads size, sequencing quality

# From Sanger to 3rd sequencing technology



Cost per Genome

Moore's Law

NIH National Human Genome Research Institute

genome.gov/sequencingcosts

1 FIRST GENERATION
From 1977

sanger

2 SECOND GENERATION
From 2007

Illumina

3 THIRD GENERATION
From 2011

PacBio, ONT

Une augmentation du débit de séquençage

Une augmentation du débit de séquençage

**Short Reads ?**
**Long Reads ?**

Vijender Chaitankar et al 2016

Shotgun

Capture

Amplicon

https://dridk.me/ngs.html

**2** SECOND GENERATION
From 2007

**YES**

✓ **Output volume** 20 billions of 150b reads, 6T

  *NovaSeq6000*

✓ **Accuracy** ~99 %

✓ **Run is cheap**

✓ **MySeq is cheap** ~60 000 USD per machine

**NO**

● **Size** 150 + 150, *NovaSeq*

  but 400 pb, *MySeq*

**3** **THIRD GENERATION**
From 2011

**A**



5' C G G A C T C 3'

electric signal

time

**ONT** is based on the translocation of a DNA or RNA strand through a nanopore located in an artificial membrane. Multiple nucleotides located in the nanopore determine the flow of ions through this nanopore in a specific way by physically blocking the space. This change in ion flux is recorded as an electric signal and further converted into sequence information.

**B**



excitation emission

signal intensity

'C' pulse    'A' pulse    'T' pulse    'G' pulse

time

Single-Molecule Real Time (**SMRT**) sequencing detects fluorescent light emitted from nucleotides upon incorporation into a DNA strand. The DNA polymerase is located at the bottom of a well and synthesises a new DNA strand. The integration into the new DNA strand keeps the nucleotide for a sufficiently long time in the well to allow detection.

Boas Pucker et al. 2022

# Two technologies

**Oxford Nanopore**

MinION

GridION

PromethION

**Pacific BioScience**

RSII

Sequel

from Elixir GAAS 2018

**3** THIRD GENERATION From 2011

Plant genome project workflow from DNA extraction over ONT sequencing to data submission

| | | task | consumed time | hands-on time | equipment | estimated costs of consumables | estimated costs of lab equipment |
|---|---|---|---|---|---|---|---|
| A | | plant incubation in darkness | 2-3d | 1h | | | |
| B | | non-destructive sampling | - | 1h | | | |
| C | | DNA extraction | 1d | 8h | waterbath, centrifuge | $50 | $1000 $8000 |
| D | | quality control | 1h | 1h | NanoDrop, Qubit | $20 | |
| E | | short fragment depletion | 2h | 1h | centrifuge | $50 | |
| F | | quality control | 1h | 1h | NanoDrop, Qubit | $20 | $5000 $5000 |
| G | | library preparation & sequencing | 1-5d | 4-16h | centrifuge, magnetic rack, sequencer | $3000 | $250 $1000 |
| H | | basecalling | 1d | 1h | computer with GPU | | $3000 |
| I | | assembly | 1-15d | 1h | | | |
| J | | polishing | 1-5d | 1h | compute cluster / cloud | | |
| K | | annotation | 1-5d | 1h | | | |
| L | | data submission | 2h | 2h | fast internet connection | | |

Boas Pucker et al. 2022

Triticum aestivum
16 Gb

Homo sapiens
3.2 Gb

Mus musculus
2.7 Gb

Danio rerio
1.4 Gb

Drosophila
melanogaster
144 Mb

Arabidopsis
thaliana
119 Mb

Saccharomyces
cerevisiae
12 Mb

Escherichia
coli K-12
4.6 Mb

Mycobacterium
tuberculosis
4.4 Mb

Ebola
19 kb

Influenza A
13.5 kb

| Microbial genomes | Human genomes | Animal genomes | Plant genomes |
| --- | --- | --- | --- |

- Simplify de novo assembly and correct existing genomes
- They bridge repetitions and build less fragmented genomes. SV, repeats, phasing
- They come from technologies which do not amplify the DNA fragments and therefore have less coverage bias.
- They are affordable.
- Detecting base modifications : they provide methylation information
- Analysing long-read transcriptomes

10 million 'pieces' (short reads)     2,000 'pieces' (long reads)

**3** THIRD GENERATION
From 2011

From Circulation
Research

- No Amplification

- NO SYNTHESIS

- Very Long Length

**3** THIRD GENERATION
From 2011

**YES**

✓ No Amplification, NO SYNTHESIS, Very Long Length

✓ Single strand direct sequencing

✓ Bases Modification detection in real-time

✓ Native RNA!

✓ **Read length** ~ 10-50kb more than 2Mb rep

✓ **Run cheap** 1,000 USD for 30Gb by now minimur

✓ **Machine cheap** 1,000 USD for Minion

✓ **Fast** 15mn library, 48-72h run

**NO**

● Error Rate 3-8%, can be corrected, 1-2% in tests

● Quality of DNA/RNA limits the sequencing

**3** **THIRD GENERATION**
From 2011



**Research areas**

- ✴ Microbiology
- 🐾 Microbiome
- 🌱 Environmental
- 🌳 Plant
- 🐄 Animal

- 👤 Human genomics
- 👥 Clinical research
- Cancer
- Transcriptome
- Populations genomics

From Nanopore website

**3** THIRD GENERATION
From 2011



**Research areas**

- ☀ Microbiology
- 🧫 Microbiome
- 🌱 Environmental
- 🌳 Plant
- 🐮 Animal

- 👤 Human genomics

**Investigations**

- 📏 Structural variation
- 🧬 SNVs and phasing
- ⚗ Gene expression
- 🔍 Identification
- ⠿ Splice variation

- ✎ Assembly
- ⇒ Fusion transcripts
- ⚘ Chromatin conformation
- ⌇ Epigenetics
- ⊙ Single cell

From Nanopore website

**3** **THIRD GENERATION** From 2011

## Research areas

- ✱ Microbiology
- 🦠 Microbiome
- 🌱 Environmental
- 🌳 Plant
- 🐄 Animal
- 👤 Human genomics

## Investigations

- Structural variation
- SNVs and phasing
- Gene expression
- Identification
- Splice variation
- Assembly

## Techniques

- Whole genome
- Targeted
- Whole transcriptome
- Metagenomics

From Nanopore website

Boas Pucker et al. 2022

# Evolution of sequencing technologies

| Sequencers | 1st generation : Sanger 3730xl (~ 2000) | 2nd generation : Illumina HiSeq 2000 (~ 2006) | 3rd generation : Pacific Biosciences RS II (~ 2012) |
|---|---|---|---|
| Method | Terminaison with dideoxynucleotides | Sequencing by synthesis with a polymerase | Real-time single molecule |
| Read length | 400-900 nt | 100 nt | 10 000 nt |
| Error rate | 0,001% | 0,1% | 15% |
| Amount of data produced at once | $10^5$ nt | $10^{12}$ nt | $10^{10}$ nt |
| Cost for $10^9$ nt | 2 000 000 € | 80 € | 1 000 € |

Adapted from http://data-science-sequencing.github.io/lectures/lecture1/

From Camille Rustenholz (Univ. Strasbourg, inrae) - Methods for plant genome assembly

The Great Wave off Kanagawa, Hokusa          @amitechsolutions.com

# Genome Totals by year and status



**Genome Sequencing Projects**

from https://gold.jgi.doe.gov/statistics

Phylogenetic distribution of Bacterial Genome Projects

from https://gold.jgi.doe.gov/statistics

# Biological Databases

✓ **Sequence**

    ○ Nucleic :     Genbank   EMBL-EBI      DDBJ *DNA Data Bank of Japan*

    ○ Proteic :     swissprot         TrEMBL

            PIR, Pfam, Prosite

✓ **Structure**      PDB        SCOP        CATH

✓ **Specialized**     by organism, by sequence type

https://www.ncbi.nlm.nih.gov/

# https://www.ncbi.nlm.nih.gov/

https://www.ncbi.nlm.nih.gov/

https://www.ncbi.nlm.nih.gov/

https://www.ncbi.nlm.nih.gov/

https://www.ncbi.nlm.nih.gov/datasets/genomes

# SRA (Sequence Reads Archive) / ENA (European Nucleotide Archive)

## Available Tools

The Rice Genome Hub provides a serie of tools to browse, visualize and search among all data sets available.

### DIANE

Tool for RNA-seq data analyses, from raw count to gene regulatory network. Allow the user to...

### Gene Search

Search for a gene by name, location, functional annotation keywords...

### Primer Designer

Primer Designer allows users to design new target-specific primers in one step as well as to...

### Primer Blaster

Check PCR primer specificity on any Rice Genome

# Sequencing project

**Design expérimental**

- Question scientifique => quelle stratégie ? Quel échantillonnage ?

    Quelle stratégie bioinfo ?

**Design expérimental**

- Question scientifique => quelle stratégie ? Quel échantillonnage ?

  Quelle stratégie bioinfo ?

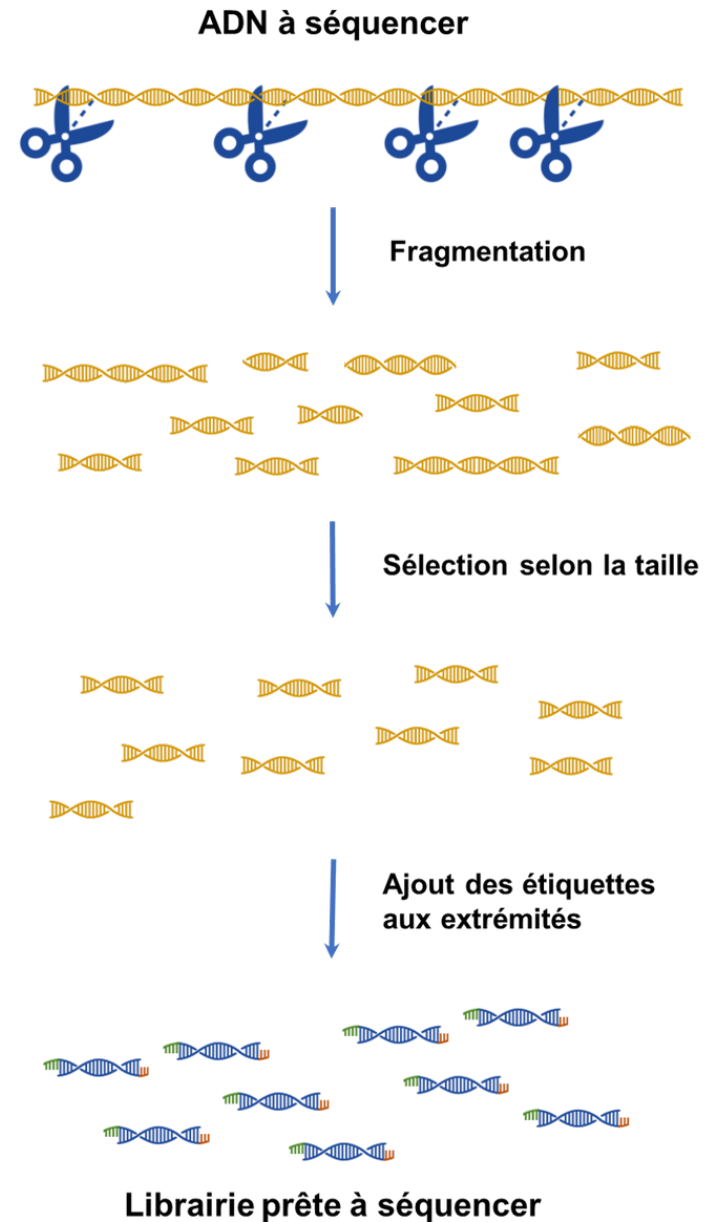- Quel méthodo de séquençage ? Quelle couverture de séquençage ?

**Design expérimental**

- Question scientifique => quelle stratégie ? Quel échantillonnage ?

    Quelle stratégie bioinfo ?

- Quel méthodo de séquençage ? Quelle couverture de séquençage ?

- Quel volume de données brut? Sur quel cluster les analyses bioinformatiques vont-elles être tournées ?
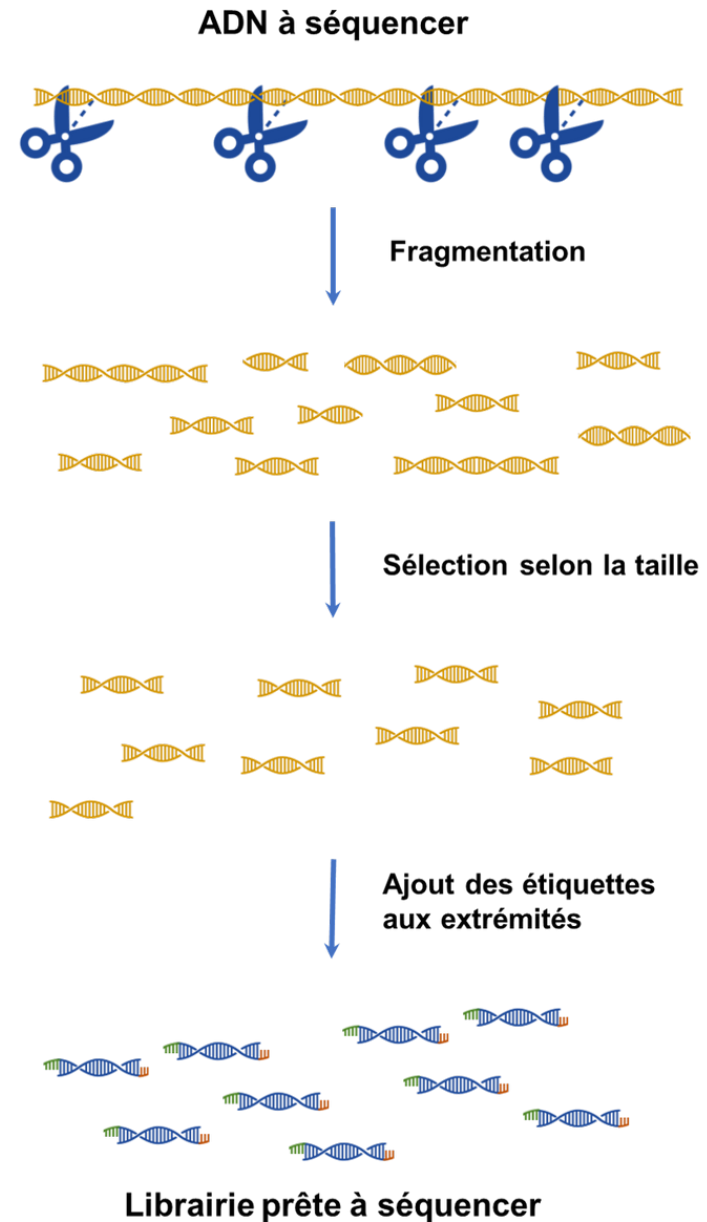
**Design expérimental**

- Question scientifique => quelle stratégie ? Quel échantillonnage ?

  Quelle stratégie bioinfo ?

- Quel méthodo de séquençage ? Quelle couverture de séquençage ?

- Quel volume de données brut? Sur quel cluster les analyses bioinformatiques vont-elles être tournées ?

- Qui va analyser mes données ?

- Où est ce que je vais stocker mes données?

Design expérimental

Préparation banque

ADN à séquencer

Fragmentation

Sélection selon la taille

Ajout des étiquettes aux extrémités

Librairie prête à séquencer

Design expérimental → Préparation banque

- Adaptateurs
- Contamination

**ADN à séquencer**

Fragmentation

Sélection selon la taille

Ajout des étiquettes aux extrémités

**Librairie prête à séquencer**

| Design expérimental | Préparation banque | Séquençage | .fastq |

- Qualité de séquençage
- Profondeur de séquençage

Genomic DNA is fragmented (not Nanopore) and sequenced -> millions of small sequences (reads) from random parts of the genome
Depending on sequence technology, reads can be from 100 bp up to 100kb in length

From Camille Rustenholz (Univ. Strasbourg, inrae) - Methods for plant genome assembly

CCGAAGT GTCAAA ATCGAGG
ATCGAGGTTC AT TTCCCG
AAGTCAAA ATCGA
TCCCGAAGTCAAA CGAGGT
GGTTCCCGAA CAAA

From Camille Rustenholz (Univ. Strasbourg, inrae) - Methods for plant genome assembly

From Camille Rustenholz (Univ. Strasbourg, inrae) - Methods for plant genome assembly

**Puzzle 400 pièces  "petite taille"**

**Puzzle 400 pièces "petite taille"**



+ 100 pièces "ciel" + …

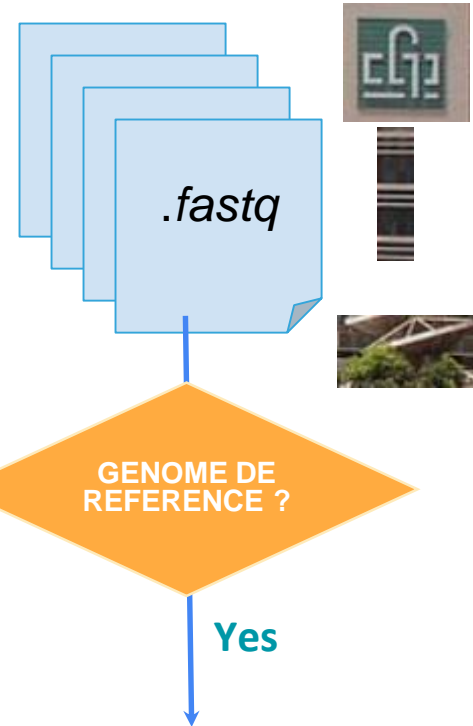**Puzzle 100 pièces - "grande taille"**

**Puzzle 100 pièces "grande taille"**



+ ~ 20 pièces "ciel"

Design expérimental → Préparation banque → Séquençage

.fastq

GENOME DE REFERENCE ?

Yes

Design expérimental → Préparation banque → Séquençage → .fastq

GENOME DE REFERENCE ?

Yes

Adapted from Ross Whetten…

Adapted from Ross Whetten…

Design expérimental → Préparation banque → Séquençage → .fastq

GENOME DE REFERENCE ?

**No** → ASSEMBLAGE *de novo*

Annotation

Yes → MAPPING CONTRE LE GENOME → ANALYSES POST MAPPING

SNP, GWAS? expression différentielle

Adapted from Ross Whetten…

# What metagenomics is ?

Metagenomics ( Environmental Genomics or Community Genomics) is the study of genomes recovered from environmental samples without the need for culturing them

Metagenomics processes data using bioinformatics tools

=> Organisms can be studied directly in their environments bypassing the need to isolate each species

=> There are significant advantages for viral metagenomics, because of difficulties cultivating the appropriate host
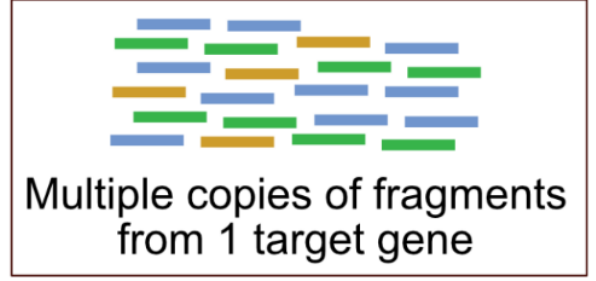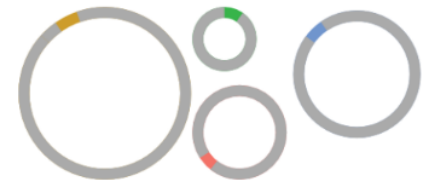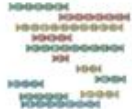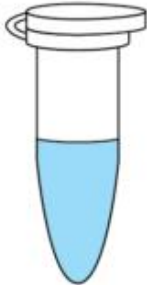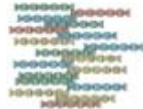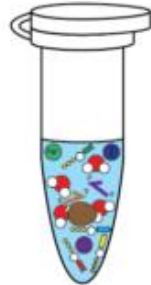
1. Collect an environmental sample
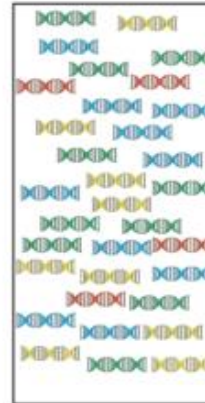2. DNA extraction from environmental sample
3. Amplify DNA markers
4. High-throughput sequencing
5. Bioinformatic processing
6. Species identification
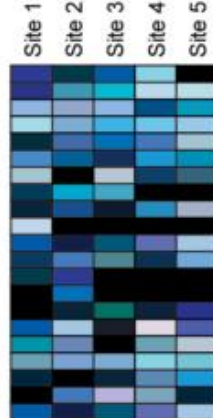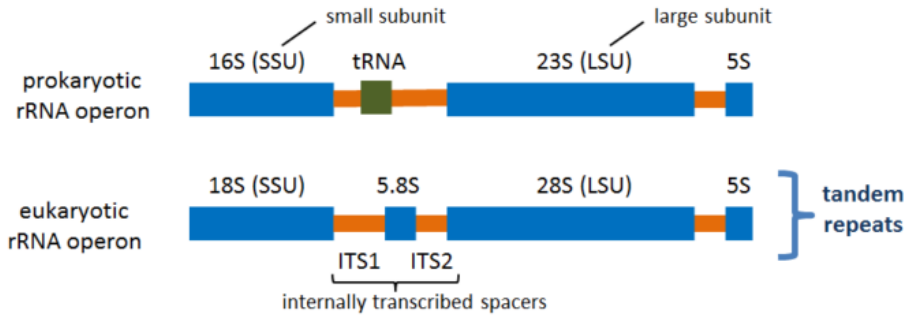7. Ecological analysis

# Markers genes vs Shotgun metagenomics

| Marker Gene Profiling | Shotgun Metagenomics Profiling |
| --- | --- |
| Less expensive (~$100 per sample) | Still very expensive (~$1000 per sample) |
| Computational needs can be met by desktop / small server computers | Usually requires huge computational resources (cluster of computers) |
| Provides mainly taxonomic profiling | Provides both taxonomic and functional profiling |
| For 16S, majority of genes can be assigned at least to phylum level | Many more unassigned gene fragments ("wasted" data) |
| Relatively free of host DNA contamination | Prone to host DNA contamination |

# Pourquoi faire du RNAseq ?

- **L'analyse d'expression différentielle** (différence d'expression dans des conditions précises) au niveau transcriptomique.

- Etude de **l'épissage alternatif** (isoformes) et recherche de nouveaux transcrits.

- **Recherche d'allèles spécifiques** et quantification de leur expression.

- **Construction d'un transcriptome** de novo pour les organismes non modèles.

# RNA sequencing

From Abims RNAseq Training oct 2018

$\Rightarrow$ Comparaison entre conditions expérimentales différentes
Ex:

- Comparaison plante infectée/saine
- Comparaison d'expression à différentes altitudes
- Comparaison ombre/soleil

$\Rightarrow$ Comparaison dans le temps (time series): cinétique
Ex:

- Cinétique d'infection de pathogènes
- Étude du rythme circadien sur l'expression de gènes

=> logiciels dédiés pour ce type de problèmatique

Origin of domestication and evolutionary history of African crop?

Where, when, how, (why) ?

| African rice | Pearl millet | Yam | Fonio | Sorghum |
|---|---|---|---|---|

From Y. Vigouroux

246 fully resequenced genomes
3 051 681 SNPs



Cubry P, Tranchant-Dubreuil C, Thuillet AC, Monat C, *et al*. Current Biol 2018

From Y. Vigouroux

**WHEN ?**

Pairwise Sequentially Markovian Coalescent (PSMC)



**The Rise and Fall of African Rice Cultivation Revealed by Analysis of 246 New Genomes**

From Cubry et al, 2018

## WHERE ?



Simulated data

Observed data

SPLATCH
Simulated diversity statistics
(SFS, rare variants)

Observed diversity statistics
(SFS, rare variants)

ABC estimation of probability of origin

From Cubry et al, 2018

# Prostrate growth 1



Prog1 deletion

# Example in GWAs & Population Genomics

# Studying Genetic Diversity using Single Nucleotidic Polymorphism (SNP)

**Understanding how individuals of a same species vary**

✓ **Variations** between individuals

✓ **Natural selection** in a population

- Each individual = unique combination of traits

- Inherited varaitions that confer an advantage (increasing an organism's chance of survival) will be pass to offspring

**Single Nucleotide Polymorphism**

## Single Nucleotide Polymorphism



## Structural Variations

**Single Nucleotide Polymorphism**

**Structural Variations**



**Presence Absence Variation (PAV)**

Deletion, duplication, copy number variation, mobile element insertion

From Yang et al., 2013

From Li et al. 2012

From Qi et al. 2014

From Yang et al., 2014

From Lin et al. 2012

From Xiao et al. 2008

From Hattori et al. 2009

From Xu et al. 2006

From Gamuyao et al. 2012

From Wang et al. 2015

From Bai et al. 2017

FromYang et al., 2013

From Li et al. 2012

From Qi et al. 2014

FromYang et al., 2014

# Is One Reference genome enough to capture all genetic diversity ?

From Lin et al. 2012

From Xiao et al. 2008

From Hattori et al. 2009

From Xu et al. 2006

From Gamuyao et al. 2012

From Wang et al. 2015

From Bai et al. 2017

*Streptococcus agalactiae*

## Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the microbial ''pan-genome''

Hervé Tettelin[a,b], Vega Masignani[b,c], Michael J. Cieslewicz[b,d,e], Claudio Donati[c], Duccio Medini[c], Naomi L. Ward[a,f], Samuel V. Angiuoli[a], Jonathan Crabtree[a], Amanda L. Jones[g], A. Scott Durkin[a], Robert T. DeBoy[a], Tanja M. Davidsen[a],

*Tetellin et al., 2005*

▸ 8 strains sequenced

▸ SNP variations

Large number of genes not shared between isolates

20% genome variability and 80 % shared by all isolates

***Pangenome concept***

## Pangenome

Collection of genes or sequences found in all individuals of a population (intra or inter species)

▸ **Core genome** : present in all individuals

▸ **Dispensable genome** : absent from one or several individuals (also called variable, accessory,...)

▸ 12,150 genes absent from the reference (18 cultivars)

▸ 27,175 genes absent from the reference (26 cultivars)

Over 20 eukaryotic pangenomes
constructed (12 Mb to 17 Gb)



**2005**
The pangenome
introduced by
Tettelin *et al.* [2]

**2007**
Plant pangenome
concept proposed
by Morgante *et al.*
[31]

**2009**
Human pangenome
built using three
genomes [50]

Bacterial super-
kingdom
pangenome built
using 573 genomes
[5]

**2013**
Pangenome of
phytoplankton
*Emiliania huxleyi*
[97]

**2015**
Review of analytical
tools and models
developed over 10
years of pangenome
research [16]

*E. coli* pangenome
built using 2085
genomes [12]

Rice accessory
genome characterized
[99]

**2017**
Pangenomes of
bread wheat [40]
and stiff brome [33]

**2019**
Human pangenome
built using 910
genomes from
humans of African
descent [51]

Human pangenome
built using 275 Han
Chinese genomes
[3]

Pig pangenome
[53]

**2006**
*Streptococcus
pneumoniae*
pangenome [17]

**2008**
*Escherichia coli*
pangenome [18]

**2014**
Soybean wild relatives
pangenome [39]

Maize pan-transcriptome
[101]

**2016**
*Brassica oleracea*
pangenome [34]

Poplar pangenome
[100]

**2018**
Rice pangenome
built using 3010
accessions [56]

Oilseed rape
pangenome [35]

*Saccharomyces
cerevisiae*
pangenome built
using 1011 isolates
[27]

using 1011 isolates
[27]

**Trends in Genetics**

**Trends in Genetics**

Golicz et al., 2020

2 formations
2 ambiances …

## Mode "training"

- Session cours suivi par
- Session pratique en autonomie (individuel ou en groupe)
- Correction

## Mode "projet"

- brainstroming en groupe, avec les formateurs
- projet en autonomie…
- debriefing collectif
- 2 projets en parallèle !

Des données différentes pour les 2 groupes avec des analyses différentes !!!

# Détection de variants
## à partir de données de séquençage short & long reads

# #data

**Genre Oryza**

● 21 espèces sauvages



From Wikimedia



From Kellogg, 2009

**Genre Oryza**

- 21 espèces sauvages

- 2 espèces domestiquées

  - *Oryza sativa*

  - *Oryza glaberrima*



From Wikimedia

From
Wikimedia

*Oryza sativa*

- Culture céréalière importante

- aliment de base de plus de la moitié de la population humaine

- Espèces diploïdes, 2n = 24 (genome AA)

- Céréale avec un génome petit

- Plante modèle

- Domestiquée ~10000 ans - O. rufipogon

*20 indidividus d'O. sativa ⇔ 20 clones*
*avec une diversité intéressante*



**Séquençage short and long reads**

*20 indidividus d'O. sativa ⇔ 20 clones*
*avec une diversité intéressante*

**Séquençage short and long reads**

*20 indidividus d'O. sativa ⇔ 20 clones*
*avec une diversité intéressante*

1. Extract 1 Mb from the Chromosome 1

2. Create 20 exact clones

3. Introduce mutations with bioinformatics program

   a. SNP : from 1 to 10%
   b. indel : between 10bp and 10kb

   c. duplications

4. Getting 20 clones with different mutations that were sequenced in silico (short & long reads)

# Projet SNP

MISSION IMPOSSIBLE
NOM DE CODE : "PROJET SNP"

Votre mission si vous l'acceptez…

# #LIEU : Burkina Faso

# #MISSION :

Le Docteur kezako, chercheuse non spécialiste en bioinformatique a réalisé une longue prospection **en Afrique.**

Elle a notamment ramené des échantillons d'ignames (elle pense que c'est de l'igname) qui présentent une diversité phénotypique particulièrement intéressante dans le contexte climatique actuel.



- Avons nous collecté une nouvelle espèce d'igname ?

- Ou avons nous collecté des ignames domestiqués ? sauvages ?

# #MISSION :

Malgrè son emploi du temps très chargée, **elle a séquencé 10 individus**

Elle met à votre disposition ces données de séquençage ainsi 5 collègues qui pourront vous assister mais leur temps est précieux car ils ont une autre mission à mener en parallèle…

Dominique

Je compte sur vous !!!!

# #DATA :

Décrire où seront les données à partir de mardi…

# A vous de jouer !

# Metagenomic

# #MISSION :

La productrice Mme. BOBODOU voudrais savoir pourquoi son champ d'ananas est peu productif

Elle a vu que les feuilles de la plante etaient plus jaunes que d'habitude… Elle s'inquiète!
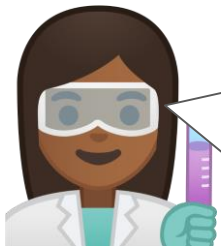
Elle a collecté quelques feuilles et a fait séquencer son échantillon mystère par la technologie Oxford Nanopore à un collègue de l'UJKZ.

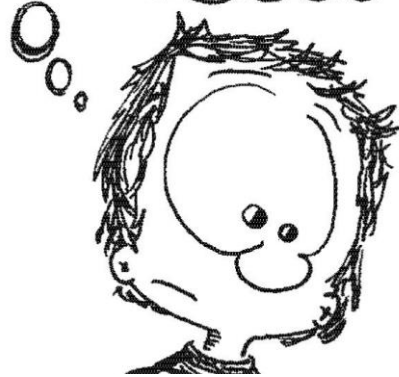**=> Aidez-lui à caractériser cet échantillon !**

# #MISSION :

**- Aidez-lui à caractériser cet échantillon !**

- Au laboratoire quelques marqueurs PCR sont négatifs pour les bactéries et les champignons pathogènes? aussi pour certains virus. Il s'agit d'une nouvelle espèce ?

# A vous de jouer !

# Bioinformatics resources

# On va travailler sous Linux !
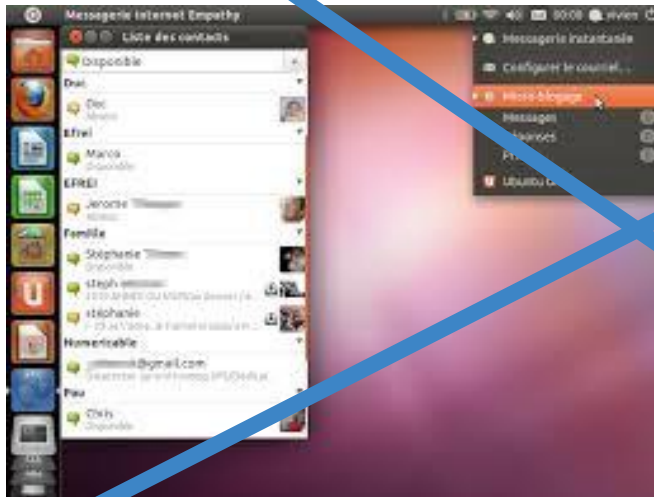
- 2 façons d'utiliser linux :
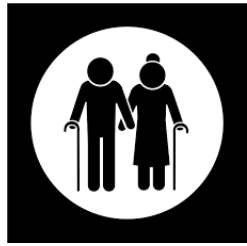
  en *mode graphique*

# On va travailler sous Linux !

- 2 façons d'utiliser linux :

    en *mode graphique*

# En mode terminal

- 2 façons d'utiliser linux :

  en *mode console*

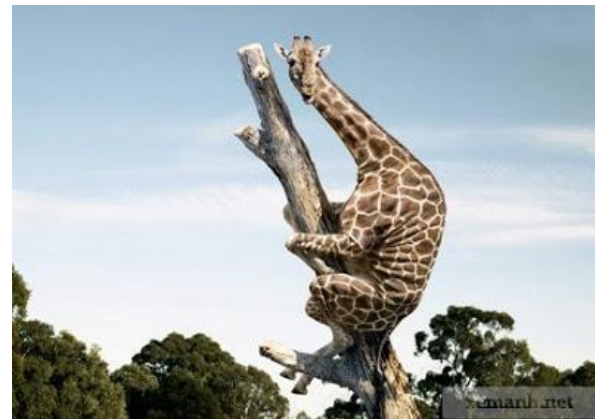# En mode terminal

- 2 façons d'utiliser linux :

  en ***mode console***



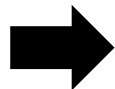Sur le cluster de l'université !

# En mode jupyter book

- Une troisième façon d'utiliser linux :

  en ***mode jupyter bool***

Sur le cloud IFB!

# *Let's discover Jupyter !*
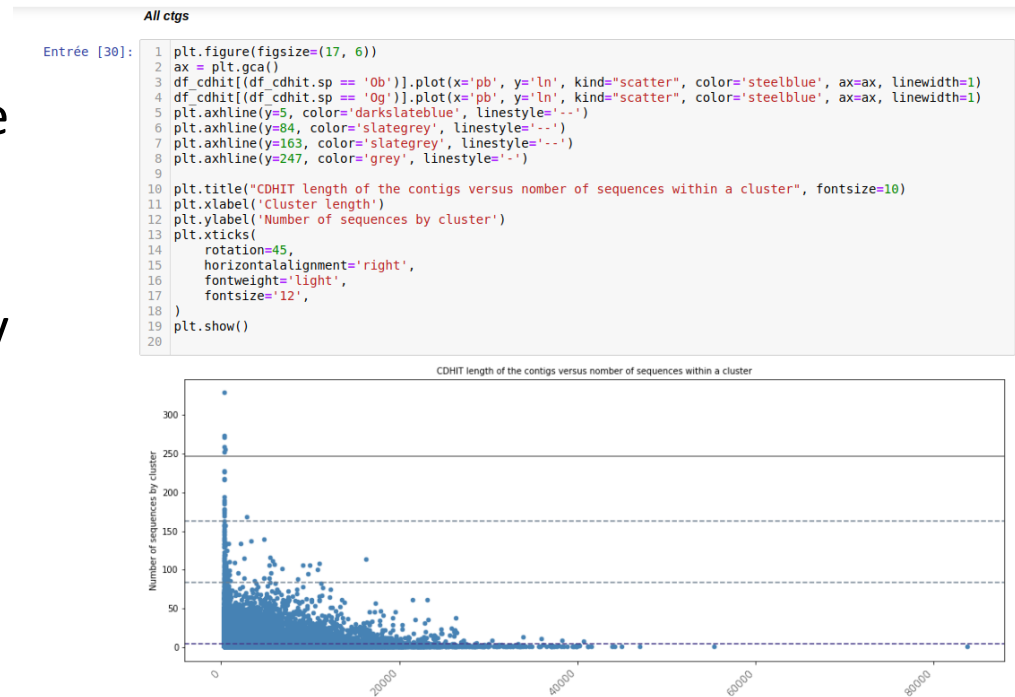
## *Working environment*

- One of the most popular tool among data scientists to perform data analysis

- Provides a complete environment in which numerous programming languages can be use through a simple web browser

ex : Bash (Linux), Python, Java, R, Julia, Matlab, Octave, Scheme, Processing, Scala

An unique interface/file where text,code and output codes can be mixed :

- code can be executed inside each cell of the notebook

- code output is directly displayed in the notebook
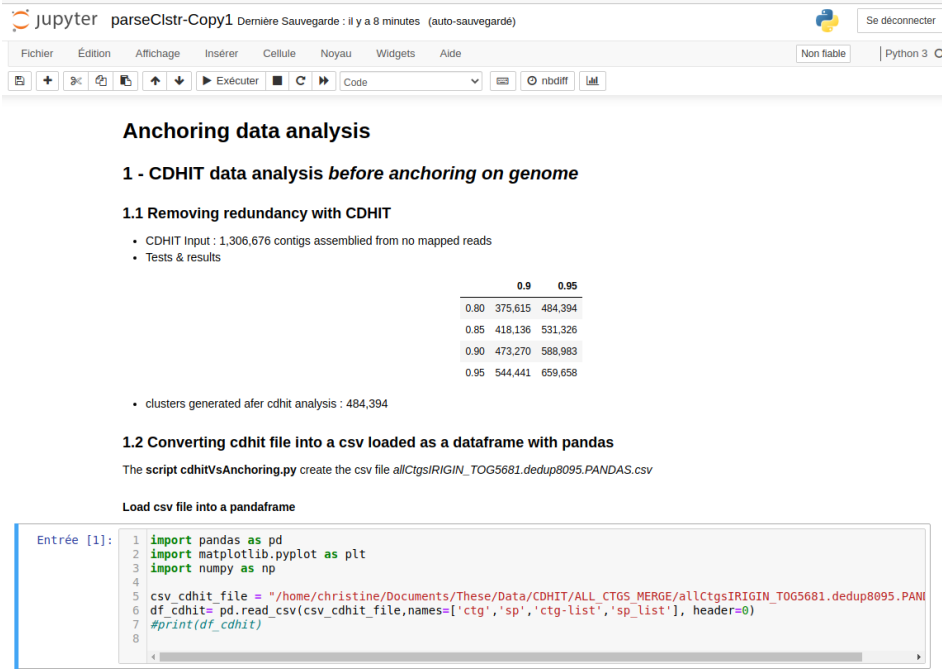
An unique interface/file where text,code and output codes can be mixed :

- code can be executed inside each cell of the notebook

- code output is directly displayed in the notebook

- explanations, formulas, charts can be added

- One file to analyze data and generate reports

- Can be exported to many formats, including PDF and HTML, which makes it easy to share your project with anyone.

- Analysis are more transparent, repeatable and shareable

- facilement importer des fichiers tabulés dans des dataframes, similaires aux dataframes sous R.

  (et exporter)

- manipuler ces tableaux de données / DataFrames

- facilement tracer des graphes à partir de ces DataFrames grâce à matplotlib

- Launch our analyses through a jupyter book within a virtual machine launched via the IFB cloud "BIOSPHERE"

- Launch our analyses through a jupyter book within a virtual machine launched via the IFB cloud "BIOSPHERE"

- Through this virtual machine, we will create jupyter books and execute all our analysis
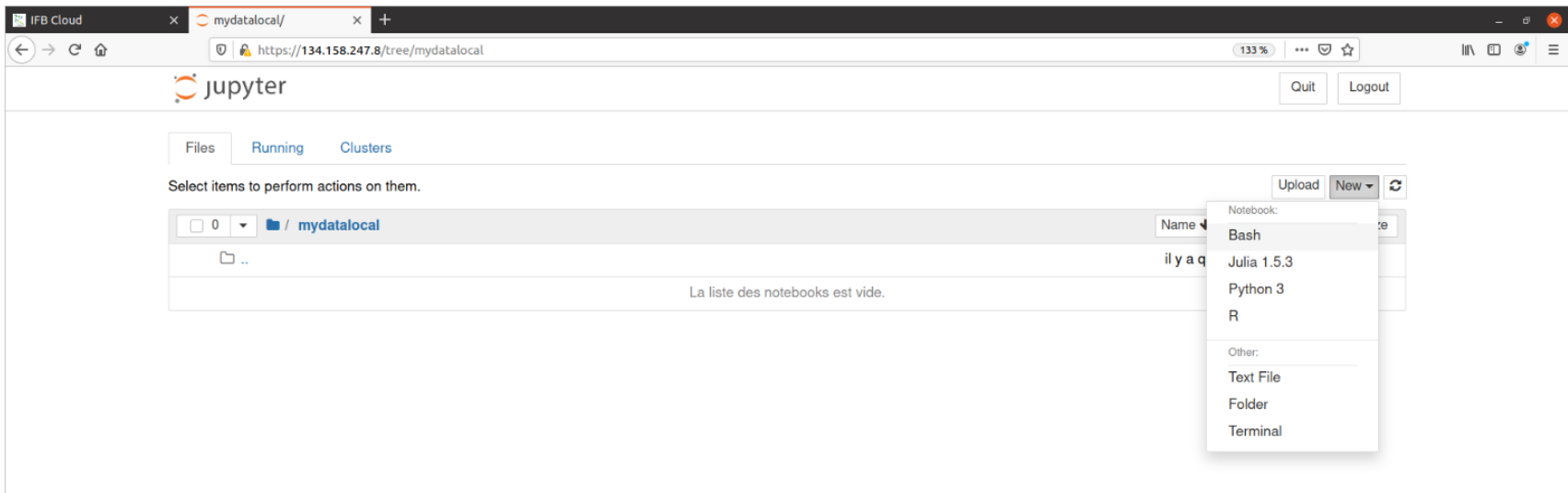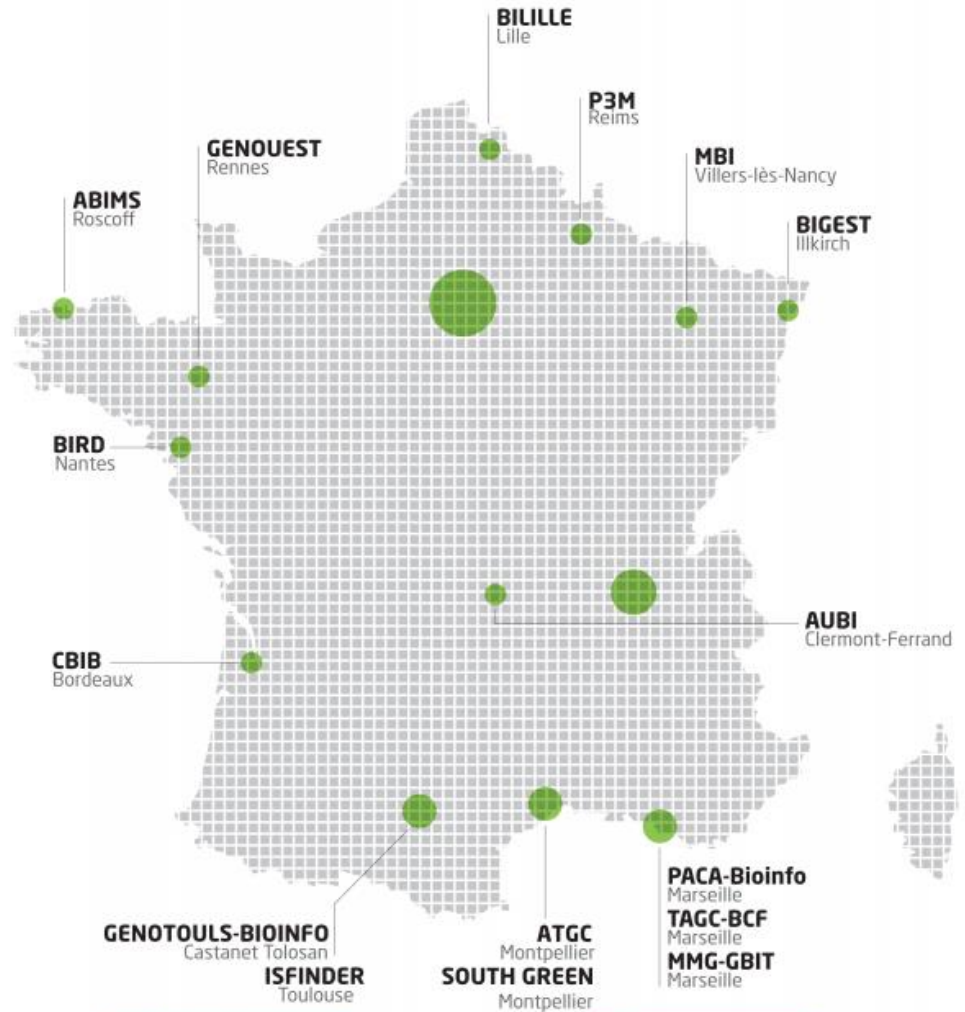
22 plateformes-membres
7 plateformes contributrices
8 équipes associées
>400 experts (~200 FTE)

- A federation of clouds, which relies on interconnected IFB's infrastructures, providing distributed services to analyze life science data

- .Access to a large set of virtual machines (computing ressources, bioinformatics tool)

- Used for scientific production in the life sciences, developments, and also to support events like cloud and scientific training sessions, hackathons or workshops.

- Open the biosphere website : https://biosphere.france-bioinformatique.fr/cloud/ and sign in

RAINBIO catalog to access our Virtual Machine (VM)

Searching for the vm we will use

Loading…

ready !

get the url… link "https"

Open your vm (https link) to access to your own jupyter lab

Go into the directory "work" and create a new jupyter book
-> kernel : bash

myFirstJupyterBook

- All jupyterbook used for practice are here :
  https://github.com/SouthGreenPlatform/training_SV_teaching/tree/2022

- Download all the jupyter books with the command *git clone*

git clone --branch 2022_burkina
https://github.com/SouthGreenPlatform/training_SV_teaching.git

Nécessité de la pratique et de l'expérience

⇔ **Investissement non négligeable pour de bons résultats rapidement**

# Détection de variants
# à partir de données de
# séquençage short & long reads

**Alexis Dereeper** - UMR PHIM

**Julie Orjuela** - UMR DIADE

**Christine Tranchant-Dubreuil** - UMR DIADE

**Détecter des variants (SNP, variants structuraux) à partir
de données de séquençage short et long reads.**

**Applications :**

- Mapper des reads contre un génome *bwa*

- Détecter des SNPs à partir du mapping de reads - *bcftools*

- Analyser les données SNPs brutes (ex: stats, filtres) - *vcftools, bcftools*

- Exemples d'études possibles à partir de SNPs - *SNIPlay*

**Avec jupyter book** : lancer les commandes + analyser les résultats

=> Avoir un plan de bataille opérationnel

# RAW SEQUENCING DATA

Design expérimental → Préparation banque → Séquençage → .*fastq*

- Statistics
- Sequencing quality ?

  Adaptators ? Contaminants ?

# fastq format

```
@H4:C7C99ACXX:6:1101:1360:74584/2
CTGTTTCTTAGTATTTTTGTAGTCATTCCGTGTTGGTTTAGTTGCAAGGT
+
@@@DADFFHHFFHIIEFEIGJGGHI4FFIEIGHI<FHGAHGGGB@3?BDB9D
@H4:C7C99ACXX:6:1101:1452:19906/2
CTGAGATCAATTGGATCCTGATGATACTGTGCTTAGCTATTACCTTTGGT
+
@@@DDDD>FFFAFBEABB4C+3?:CBB@<<A?E4A???9C@CFF*9*B3D?B
@H4:C7C99ACXX:6:1101:1476:35220/2
CATGTGCTATTACCAAAAGTGCAGTAACGACCTATAAATTTTAAAGTAGC
+
@CFFFFFGGHHHHIJJJJIEE<HHHIJJIGBHGGEEIIJJEIEIJIHHJFIIJJJGHJJ
@H4:C7C99ACXX:6:1101:1491:94128/2
AGAAGTCTTCGGAAAAGTTCGGGTATGGCTCTAGTAGCTTTTGTCTTAT
+
@C@FFFFFGGHHDHGIIEEHIII<CGHIJIJIIJ:?FC9DGAFGHII?DGBFIIJHBI
@H4:C7C99ACXX:6:1101:1538:34462/2
ACAAAAAGCTAAAAGAACACAGTTGCTTGAAGCAGCAAACACAAGAAC
+
B@@DFFFFGHHHHJIIIIIJJJIIGJCHHEIII>GHIG@GHIDHGJIIFIFHIJJIJJJG
@H4:C7C99ACXX:6:1101:1568:67898/2
ACAAATGGGTGTGTAAGAGTTAAAAAACAATTAATGAGCAACTGAGTTC
+
@@CFFFFFHFFHFGIJJIHHIIJJIIIHIJJJECGHIJJCHGICDGGGHJ<FGGIJJ
@H4:C7C99ACXX:6:1101:1575:18963/2
AACATGTTTGTCGGGGGGTTGGGAAATTGTCACTTTCTGCTACAATGCCG
+
@<@DDDDDHFFFFDIIBDFGHHGG;FGGCHHAGGGIIH@E>AEDDEECAB>
```

## 1 sequence/read = 4 lines

- read id, starting by **@**
- read sequence
- Comment line starting by **+**
(usually contains read id).
- read Quality for each base

# PHRED SCORE

- Séquenceur assigne à chaque base séquencée un score lié à la probabilité que la base appelée soit fausse

$$Q = -10 \ \log_{10} P$$

or

$$P = 10^{\frac{-Q}{10}}$$

*Ewing 1998*

- Ce score (PHRED score) varie entre 0 et 50

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---------------------|-----------------------------------|--------------------|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10000 | 99,99% |
| 50 | 1 in 100000 | 99.999 % |

# How to code quality score for each base with one letter ?

```
@H4:C7C99ACXX:6:1101:1360:74584/2
CTGTTTCTTAGTATTTTTGTAGTCATTCCGTGTTGGTTTAGTTGCAAGGT
+
@@@DADFFHHFFHIIEFEIGJGGHI4FFIEIGHI<FHGAHGGGB@3?BDB9D
@H4:C7C99ACXX:6:1101:1452:19906/2
CTGAGATCAATTGGATCCTGATGATACTGTGCTTAGCTATTACCTTTGGT
+
@@@DDDD>FFFAFBEABB4C+3?:CBB@<<A?E4A???9C@CFF*9*B3D?B
@H4:C7C99ACXX:6:1101:1476:35220/2
CATGTGCTATTACCAAAAGTGCAGTAACGACCTATAAATTTTAAAGTAGC
+
@CFFFFFGGHHHHIJJJJIEE<HHHIJJIGBHGGEEIIJJEIEIJIHHJFIIJJJGHJJ
@H4:C7C99ACXX:6:1101:1491:94128/2
AGAAGTCTTCGGAAAAGTTCGGGTATGGCTCTAGTAGCTTTTGTCTTAT
+
@C@FFFFFGGHHDHGIIEEHIII<CGHIJIJIIJ:?FC9DGAFGHII?DGBFIIJHBI
@H4:C7C99ACXX:6:1101:1538:34462/2
ACAAAAAGCTAAAAGAACACAGTTGCTTGAAGCAGCAAACACAAGAAC
+
B@@DFFFFGHHHHJIIIIIJJJIIGJCHHEIII>GHIG@GHIDHGJIIFIFHIJJIJJJG
@H4:C7C99ACXX:6:1101:1568:67898/2
ACAAATGGTGTGTAAGAGTTAAAAAACAATTAATGAGCAACTGAGTTC
+
@@CFFFFFHFFHFGIJJIHHIIJJIIIHIJJJECGHIJJCHGICDGGGHJ<FGGIJJ
@H4:C7C99ACXX:6:1101:1575:18963/2
AACATGTTTGTCGGGGGGTTGGGAAATTGTCACTTTCTGCTACAATGCCG
+
@<@DDDDDHFFFFDIIBDFGHHGG;FGGCHHAGGGIIH@E>AEDDEECAB>
```

## 1 sequence/read = 4 lines

- read id, starting by **@**
- read sequence
- Comment line starting by **+** (usually contains read id).
- read Quality for each base

**Code ASCII**

## ASCII Table

| | Code Char | Code | Char | Code | Char | Code | Char |
|---|---|---|---|---|---|---|---|
| 0 | NUL (null) | 32 | SPACE | 64 | @ | 96 | ` |
| 1 | SOH (start of heading) | 33 | ! | 65 | A | 97 | a |
| 2 | STX (start of text) | 34 | " | 66 | B | 98 | b |
| 3 | ETX (end of text) | 35 | # | 67 | C | 99 | c |
| 4 | EOT (end of transmission) | 36 | $ | 68 | D | 100 | d |
| 5 | ENQ (enquiry) | 37 | % | 69 | E | 101 | e |
| 6 | ACK (acknowledge) | 38 | & | 70 | F | 102 | f |
| 7 | BEL (bell) | 39 | ' | 71 | G | 103 | g |
| 8 | BS (backspace) | 40 | ( | 72 | H | 104 | h |
| 9 | TAB (horizontal tab) | 41 | ) | 73 | I | 105 | i |
| 10 | LF (NL line feed, new line) | 42 | * | 74 | J | 106 | j |
| 11 | VT (vertical tab) | 43 | + | 75 | K | 107 | k |
| 12 | FF (NP form feed, new page) | 44 | , | 76 | L | 108 | l |
| 13 | CR (carriage return) | 45 | - | 77 | M | 109 | m |
| 14 | SO (shift out) | 46 | . | 78 | N | 110 | n |
| 15 | SI (shift in) | 47 | / | 79 | O | 111 | o |
| 16 | DLE (data link escape) | 48 | 0 | 80 | P | 112 | p |
| 17 | DC1 (device control 1) | 49 | 1 | 81 | Q | 113 | q |
| 18 | DC2 (device control 2) | 50 | 2 | 82 | R | 114 | r |
| 19 | DC3 (device control 3) | 51 | 3 | 83 | S | 115 | s |
| 20 | DC4 (device control 4) | 52 | 4 | 84 | T | 116 | t |
| 21 | NAK (negative acknowledge) | 53 | 5 | 85 | U | 117 | u |
| 22 | SYN (synchronous idle) | 54 | 6 | 86 | V | 118 | v |
| 23 | ETB (end of trans. block) | 55 | 7 | 87 | W | 119 | w |
| 24 | CAN (cancel) | 56 | 8 | 88 | X | 120 | x |
| 25 | EM (end of medium) | 57 | 9 | 89 | Y | 121 | y |
| 26 | SUB (substitute) | 58 | : | 90 | Z | 122 | z |
| 27 | ESC (escape) | 59 | ; | 91 | [ | 123 | { |
| 28 | FS (file separator) | 60 | < | 92 | \ | 124 | | |
| 29 | GS (group separator) | 61 | = | 93 | ] | 125 | } |
| 30 | RS (record separator) | 62 | > | 94 | ^ | 126 | ~ |
| 31 | US (unit separator) | 63 | ? | 95 | _ | 127 | DEL |

# How to code quality score for each base with one letter ?

*Code ASCII*

| Code Char |
|---|
| 64 @ |
| 65 A |
| 66 B |
| 67 C |
| 68 D |
| 69 E |
| 70 F |
| 71 G |
| 72 H |
| 73 I |
| 74 J |
| 75 K |
| 76 L |
| 77 M |
| 78 N |
| 79 O |
| 80 P |
| 81 Q |
| 82 R |
| 83 S |
| 84 T |
| 85 U |
| 86 V |
| 87 W |
| 88 X |
| 89 Y |
| 90 Z |
| 91 [ |
| 92 \ |
| 93 ] |
| 94 ^ |
| 95 _ |

| Phred Quality Score |
|---|
| 0 .. 50 |

| Code Char |
|---|
| 96 ` |
| 97 a |
| 98 b |
| 99 c |
| 100 d |
| 101 e |
| 102 f |
| 103 g |
| 104 h |
| 105 i |
| 106 j |
| 107 k |
| 108 l |
| 109 m |
| 110 n |
| 111 o |
| 112 p |
| 113 q |
| 114 r |
| 115 s |
| 116 t |
| 117 u |
| 118 v |
| 119 w |
| 120 x |
| 121 y |
| 122 z |
| 123 { |
| 124 | |
| 125 } |
| 126 ~ |
| 127 DEL |

Design expérimental → Préparation banque → Séquençage → *.fastq*

- Statistics

- Sequencing quality ? Adaptators ?

  Contaminants ?

Design expérimental → Préparation banque → Séquençage → .fastq

- Statistics
- Sequencing quality ? Adaptators ?
  Contaminants ?

➡ Basic statistics and quality control checks using **fastqc**

**fastqc** to get some basic statistics and to do some quality control checks

```
# fastqc command
fastqc /path2fastq/AX8798.fastq  -o path2fastqcDIR

fastqc /path2fastq/*  -o path2fastqcDIR
```

https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

**[command line] manuel :**
https://manpages.ubuntu.com/manpages/trusty/man1/fastqc.1.html#:~:text=DESCRIPTION,of%20proble
m%20in%20your%20data

This plot shows the base quality score distribution for all reads in a lane, with each read position considered independently.

- x-axis = position in read (bp)
- y-axis = Phred-like base quality score [pink=0-20, tan=20-30, green=30-40]
- red bar = median score, blue line = mean score
- yellow box = 25th to 75th percentile, black whiskers = 10th to 90th percentile



**GOOD/NORMAL LANE**

**SALVAGEABLE LANE**



**FAILED LANE**

This plot shows the nucleotide distribution per read position for all reads in a lane.

- x-axis = position in read (bp)
- y-axis = % of all reads in the lane
- colors refer to individual nucleotides: **A, C, G, T**

**GOOD LANE**



**BAD LANE**



**Can this be fixed?** No.

- A contamination ?

- A contamination ?



**Can this be fixed ?**  Maybe…

GC distribution over all sequences

GC distribution over all sequences

*Sabellaria alveolata* : mantle transcriptome

This plot shows the degree of duplication for a subset of reads in a lane.

- x-axis = sequence duplication level
- y-axis = % duplicates relative to unique reads

**GOOD LANE**



**BAD LANE**



**Can this be fixed?** Maybe.

**Can this be fixed?  Hem...**

## Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| AGAGTTTTATCGCTTCCATGACGCAGAAGTTAACACTTTC | 2065 | 0.5224039181558763 | No Hit |
| GATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCCATG | 2047 | 0.5178502762542754 | No Hit |
| ATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCCATGA | 2014 | 0.5095019327680071 | No Hit |
| CGATAAAAATGATTGGCGTATCCAACCTGCAGAGTTTTAT | 1913 | 0.4839509420979134 | No Hit |
| GTATCCAACCTGCAGAGTTTTATCGCTTCCATGACGCAGA | 1879 | 0.47534961850600066 | No Hit |
| AAAAATGATTGGCGTATCCAACCTGCAGAGTTTTATCGCT | 1846 | 0.4670012750197325 | No Hit |

Adapter dimers
rRNA
Satellite sequences

| | | | |
|---|---|---|---|
| TCATGGAAGCGATAAAACTCTGCAGGTTGGATACGCCAAT | 665 | 0.16823177025358726 | No Hit |
| TCTGCGTCATGGAAGCGATAAAACTCTGCAGGTTGGATAC | 627 | 0.15861852623909656 | No Hit |
| GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCT | 624 | 0.1578595859221631 | Illumina Paired End PCR Primer 2 (100% over 40bp) |
| CCTGCAGAGTTTTATCGCTTCCATGACGCAGAAGTTAACA | 613 | 0.15507680476007366 | No Hit |
| CGGTTCAGCAGGAATGCCGAGATCGGAAGAGCGGTTCAGC | 599 | 0.15153508328105078 | Illumina Paired End PCR Primer 2 (96% over 25bp) |
| TCTGCAGGTTGGATACGCCAATCATTTTTATCGAAGCGCG | 585 | 0.1479933618020279 | No Hit |
| CGCTTAAAGCTACCAGTTATATGGCTGGGGGGTTTTTTTT | 552 | 0.13964501831575965 | No Hit |
| CTCTGCAGGTTGGATACGCCAATCATTTTTATCGAAGCGC | 532 | 0.1345854162028698 | No Hit |
| CTGCGTCATGGAAGCGATAAAACTCTGCAGGTTGGATACG | 515 | 0.13028475440691342 | No Hit |
| CTGCAGGTTGGATACGCCAATCATTTTTATCGAAGCGCGC | 505 | 0.12775495335046852 | No Hit |
| GCTTAAAGCTACCAGTTATATGGCTGGGGGGTTTTTTTTG | 411 | 0.10397482341988626 | No Hit |

## Kmer Content



Relative enrichment over read length

| Sequence | Count | Obs/Exp Overall | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|-----------------|-------------|----------------------|
| TTTTT | 192940 | 8.590186 | 21.06293 | 29 |
| CTGCA | 90975 | 7.7906475 | 12.251836 | 10 |
| GCAGA | 84910 | 7.163295 | 13.539302 | 23 |
| TGCAG | 92470 | 7.002405 | 10.671717 | 11 |
| CCTGC | 57235 | 5.4987235 | 8.729035 | 16 |
| GTTTT | 108205 | 5.324498 | 10.243909 | 28 |
| CAACC | 49005 | 5.2869425 | 9.85526 | 13 |
| ATCGC | 58320 | 4.9942355 | 8.029807 | 29 |
| CCAAC | 46220 | 4.9864807 | 9.408141 | 12 |
| AAAAA | 62285 | 4.7588468 | 8.0126295 | 5 |
| CAGAG | 56370 | 4.7555633 | 7.148592 | 20 |
| ACCTG | 55315 | 4.736902 | 7.919266 | 15 |
| CGCCA | 44035 | 4.7130895 | 8.830201 | 35 |
| GGGGG | 63675 | 4.67525 | 16.94222 | 27 |
| GCAGG | 55380 | 4.6350074 | 17.521912 | 19 |
| AAAAC | 51945 | 4.452569 | 8.159592 | 24 |
| TATCG | 64615 | 4.4271946 | 8.394971 | 34 |
| GCTGG | 58505 | 4.3952427 | 10.37436 | 18 |
| AACCT | 50775 | 4.382863 | 7.691214 | 14 |
| TTATC | 70080 | 4.3444843 | 7.810299 | 33 |
| TTTTA | 87340 | 4.332125 | 7.8541703 | 28 |
| TTTAT | 86645 | 4.297653 | 7.9511886 | 35 |
| CGCTT | 54695 | 4.2042785 | 6.9374876 | 31 |

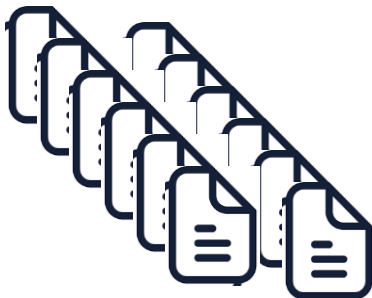**fastqc** to get some basic statistics and to do some quality control checks

```
# fastqc command
fastqc /path2fastq/AX8798.fastq  -o path2fastqcDIR

fastqc /path2fastq/*  -o path2fastqcDIR
```

**fastqc** generate one report by fastq file

**With numerous fastq and fastqc report => use MultiQC**

**MultiQC** : a modular tool to summarise results from a bioinformatics analyse performed on many samples into a single report

```
# MultiQC command
multiqc path2fastqcDIR
```

https://multiqc.info/

A modular tool to aggregate results from bioinformatics analyses across many samples into a single report.

Report generated on 2020-10-29, 16:10 based on data in: `/work_home/orue/FROGS_16S/FASTQC`

## QUALITE DE SEQUENÇAGE & « NETTOYAGE »

### cutadapt, trimmomatics

- Détection et retrait des adaptateurs et primers

- Retrait des queue polyA/T

- Détection des séquences contaminantes, ARN ribosomal

- Masquage des bases avec phred score bas par N

- Séquences courtes après retriat des adaptateurs