South Green
bioinformatics platform

Charte

**diade**
Larmande Pierre
Sabot François
Tando Ndomassi
**Tranchant-Dubreuil**
**Christine**
**IPME**
Comte Aurore
Dereeper Alexis
**diade** **IPME**
**Orjuela-Bouniol Julie**

**agap**
Bocs Stephanie
De Lamotte Fredéric
**Droc Gaetan**
Dufayard Jean-François
Hamelin Chantal
Martin Guillaume
Pitollat Bertrand
**Ruiz Manuel**
**Sarah Gautier**
Summo Marilyne

**Bioversity International**
**Rouard Mathieu**
Guignon Valentin
Catherine Breton

**UMR BGPI**
**Mahé Frédéric**
**Ravel Sébastien**
**intertryp**

Sempere Guilhem

Bioversity International    cirad    IRD Institut de Recherche pour le Développement FRANCE    INRA SCIENCE & IMPACT
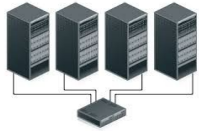
South Green bioinformatics platform

**Workflow manager**

TOGGLe — Toolbox for generic NGS analyses
SNAKEMAKE
Galaxy

**HPC and trainings….**

37 courses organized last 7 years

Trainees number
100
200
300
400

IRD — Institut de Recherche pour le Développement
cirad

**Genome Hubs & Information System**

Gigwa

*SNPs and Indels*

GreenPhyl

| Family Id | Family Name | Number of sequences | Status |
|---|---|---|---|
| GP000010 | Cytochrome P450 superfamily | 6942 | |
| GP000017 | AP2/EREBP transcription factor family: ERF/DREB group (partial) | 5142 | |
| GP000020 | NAC transcription factor family | 4574 | |
| GP000028 | MADS transcription factor family | | |
| GP000031 | Haem peroxidase superfamily | | |
| GP000066 | General substrate transporter superfamily | | |
| GP000022 | Subtilisin-like Serine Proteases family | | |
| GP000019 | NPF, NRT1/PTR FAMILY | | |

*Gene families*

SNiPlay

https://github.com/SouthGreenPlatform

@green_bioinfo

*The South Green portal: a comprehensive resource for tropical and Mediterranean crop genomics*, Current Plant Biology, 2016

SouthGreen bioinformatics platform

ABiMS — Erwan Corre

ASTRE — Etienne Loire
Julie Reveillaud

diade — Valentin Klein

TransVIHMI — Emmanuelle Beyne

IPME — Marie Simonin
Sébastien Cunnac

UMR QualiSud — Florentin Constancias

MiVEGEC / i-Trop bioinformatics — Valérie Noël

**And more collaborators !**

# Modules de formation 2019

- Toutes nos formations :

    **https://southgreenplatform.github.io/trainings/**

- Topo & TP : **Workflow managers**

- Environnement de travail : **Logiciels à installer**

# Formateurs

- Christine Tranchant-Dubreuil
- **Sebastien Ravel**
- Alexis Dereeper
- **Jean-François Dufayard**
- Ndomassi Tando
- Bertrand Pitollat
- **François Sabot**
- **Julie Orjuela-Bouniol**
- Gautier Sarah
- **Aurore Comte**
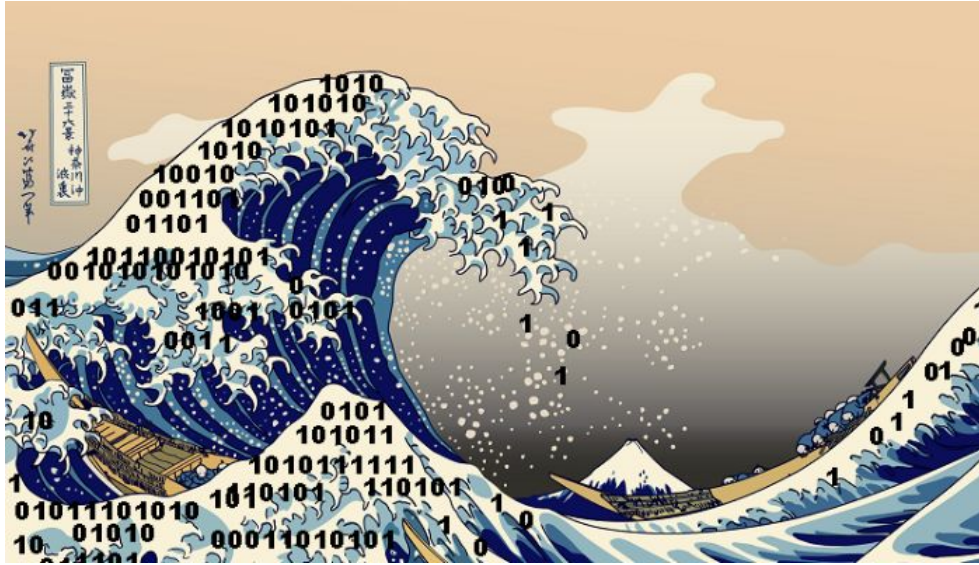- **Marilyne Summo**
- **Guilhem Sempere**

**objectifs:**

**Utiliser les gestionnaires de Workflow de South Green afin de construire de manière automatique vos propres pipelines.**

**Applications**

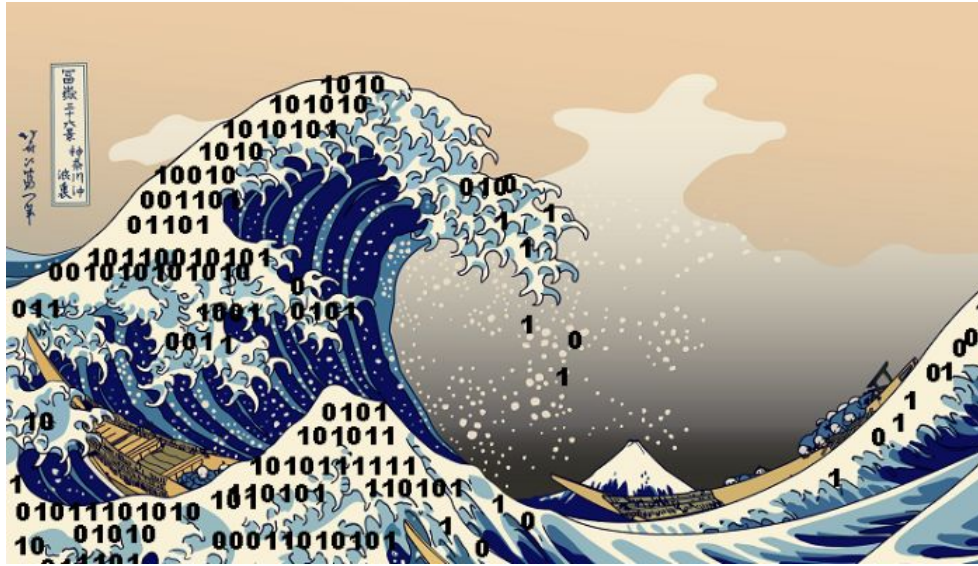Tout savoir sur les 2 principaux gestionnaires de workflow

- Utiliser les outils
- Construire son propre workflow
- Pratiquer sur un même cas d'utilisation : Appel de SNPs à partir de reads Illumina de 3 échantillons

The Great Wave off Kanagawa, Hokusa        @amitechsolutions.com

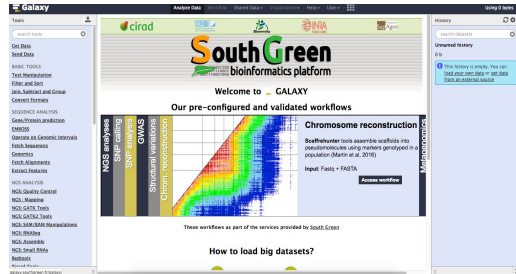Créer son propre pipeline via une méthode facile et conviviale

| Données brutes | → | Résultats Intermédiaires | → | Résultats Intermédiaires | → | Résultat Final |

- 3 solutions proposées par

Galaxy

*Snakemake*

TOGGLe

**Facilité d'utilisation
Bonne documentation**

**Facilité de
développement**

- 3 solutions proposées par

**Contrôle du pipeline et des données**

# Apporte un cadre robuste

✓ Vérifie le format des fichiers

Valide l'enchaînement des outils

➕ Automatisation de certaines étapes clefs
( ex : indexation de la référence )

| -add / –add: | use if you want to add new samples to an already run analysis. |
|---|---|
| -rerun / –rerun: | use if you want to re-run samples that have encountered error previously. |

**Contrôle du pipeline et des données**

**Apporte un cadre robuste**

**Reproductibilité & traçabilité**

**Suivi des erreurs & reprise en cours**

**Contrôle du pipeline et des données**

**Apporte un cadre robuste**

**Reproductibilité & traçabilité**

**Analyse de gros jeu de données**

**Suivi des erreurs & reprise en cours**

**Contrôle du pipeline et des données**

**Connection HPC Parallélisation**

**Apporte un cadre robuste**

**Reproductibilité & traçabilité**

**Analyse de gros jeu de données**

**Suivi des erreurs & reprise en cours**

| | TOGGLe | Galaxy |
|---|---|---|
| **Interface** | Command line | GUI (Web interface) |
| **Predefined Pipelines** | SNP calling, RNASeq and WGS large scale | Metagenomics, RNASeq, SNP calling, post-analyses |
| **Number of Samples** | +++ | ++ |
| **Quota (related to infra)** | Disk space "/data/projects" | IRD     100Go data<br>Cirad    100Go => 300Go |
| **Parallelization (related to infra conf)** | IRD     300 cores<br>Cirad    600 cores | IRD     *16 cores / one node*<br>Cirad    *200 cores* |
| **Number of tools available** | ++ (120) | ++++ (5500 avail) |
| **Post-analyses Graphical figures** | Not yet | Yes |

- An alternate solution between GUI tools (Galaxy) & CLI tools (SnakeMake)

- Numerous tools integrated to perform post analysis

- Targets both biologists & bionformaticians

19 modules, 120 functions 120 open-source tools

GBS   RADSeq   RNASeq   WGS

TOGGLe

**Various data format** : fasta, fastq, sam, bam, bed, vcf (compressed or not)

A command-line based pipeline framework

A single command line

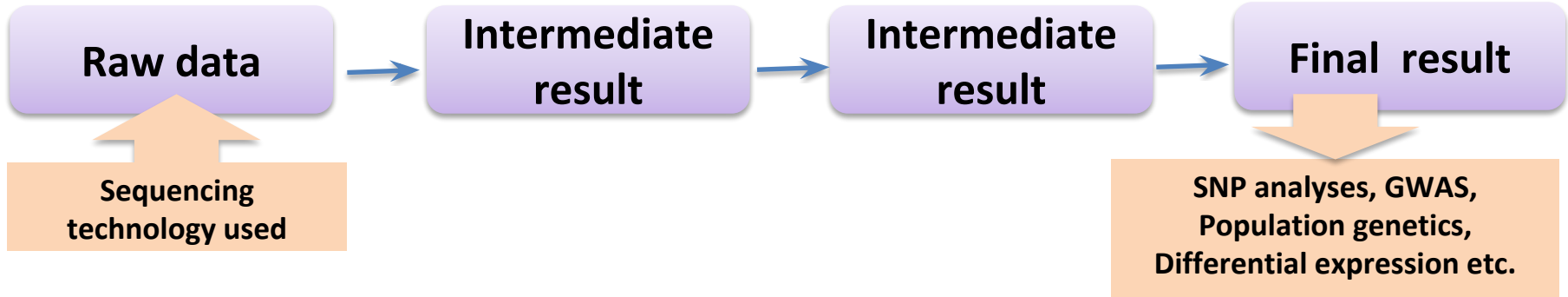**toggleGenerator.pl -d DIR -o DIR -c FILE**

toggleGenerator.pl -d|--directory DIR -c|--config FILE -o|--outputdir DIR [-r|--reference FILE] [-k|--keyfile FILE] [-g|--gff FILE] [-nocheck|--nocheckFastq] [--help|-h]

**Required named arguments:**

| | |
|---|---|
| -d / –directory DIR | a folder with raw data to be treated (FASTA, FASTQ, SAM, BAM, BED, GFF, VCF) |
| -c / –config FILE | it is the *software.config.txt* file but it can be any text file (Unix format). |
| -o / –outputdir DIR | the current version of TOGGLe will not modify the initial data folder but will create an output directory with all analyses in. This module must be empty (TOGGLe will stop if not). |

**Optional named arguments:**

| | |
|---|---|
| -r / –reference FILE | a reference FASTA file to be used. (1) |
| -g / -gff FILE | a GFF file to be used for some tools . Be careful the gff name must be different than the FASTA. |
| -k / –keyfile FILE | a keyfile use for demultiplexing step. |
| -add / –add | use if you want to add new samples to an already run analysis. |
| -rerun / –rerun: | use if you want to re-run samples that have encountered error previously. |
| -nocheck | by default checks if given formats for input files are correct. This option allows to skip this step. |
| -report / –report | generate pdf report (more info) |
| -h / –help | show help message and exit |

- 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=samToolsView
1000=gatkHaplotypeCaller
1001=gatkVariantFiltration

**$cutadapt**
-q 30
-m 35
**$bwa mem**

-n 5
**…**

**$sge**
-q bioinfo.q
-b Y

## Create your own workflow

- The workflow order

- The list of softwares to run

```
$order
1=fastqc
2=cutadapt
3=bwa mem
4=samToolsView
1000=gatkHaplotypeCaller
1001=gatkVariantFiltration
```

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

**$order**
1=fastqc
2=cutadapt
3=bwa mem
4=samToolsView
1000=gatkHaplotypeCaller
1001=gatkVariantFiltration

## Create your own workflow

Step number < 1OOO

**Parallel analysis by sample**

```
$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

**$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

**Job schedulers**

LSF, MPRUN, SLURM, SGE

# More info on website

http://toggle.southgreen.fr/

**User Manuals**

**Screencast**

**Developer manual**

**Pre-defined Workflow files**

https://github.com/SouthGreenPlatform/TOGGLE

# TOGGLE 's team



Christine Tranchant-Dubreuil
christine.tranchant@ird.fr
IRD Institut de Recherche pour le Développement FRANCE • diade

Julie Orjuela-Bouniol
julie.orjuela@ird.fr
IRD Institut de Recherche pour le Développement FRANCE • BOREA

François Sabot
francois.sabot@ird.fr
IRD Institut de Recherche pour le Développement FRANCE • diade

Sébastien Ravel
sebastien.ravel@cirad.fr
cirad • BGPI

Alexis Dereeper
alexis.dereeper@ird.fr
IRD Institut de Recherche pour le Développement FRANCE • IPME

Ndomassi Tando
ndomassi.tando@ird.fr
IRD Institut de Recherche pour le Développement FRANCE • diade

Valérie Noel
valerie.noel@ird.fr
cnrs • MiVEGEC

Valentin Klein
valentin.klein@ird.fr
IRD Institut de Recherche pour le Développement FRANCE • diade

@ toggle@ird.fr

## Former collaborators

- Cécile Monat (UMR DIADE - IRD)
- Gautier Sarah (UMR AGAP - CIRAD)
- Abdoulaye Diallo (UMR DIADE - IRD)
- Laura Helou (UMR DIADE - IRD)
- Souhila Amazougarene (UMR DIADE - IRD)

- Mawussé Agbessi (UMR DIADE - IRD)
- Enrique Ortega-Abboud (UMR AGAP - CIRAD)
- Cédric Farcy (UMR AGAP - CIRAD)
- Maryline Summo (UMR AGAP - CIRAD)
- Ayité Kougbeadjo (UMR DIADE - IRD)

## Comment citer TOGGLe?

Tranchant-Dubreuil, C., Ravel, S., Monat, C., Sarah, G., Diallo, A., Helou, L., … Sabot, F. (2018). [TOGGLe, a flexible framework for easily building complex workflows and performing robust large-scale NGS analyses.](#) *BioRxiv*. https://doi.org/10.1101/245480

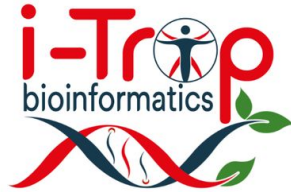→ N'oubliez pas de citer aussi les outils utilisés !

## Comment citer les clusters?

"The authors acknowledge the IRD i-Trop HPC at IRD Montpellier for providing HPC resources that have contributed to the research results reported within this paper. URL: http://bioinfo.ird.fr/ "

"The authors acknowledge the CIRAD UMR-AGAP HPC (South Green Platform) at CIRAD montpellier for providing HPC resources that have contributed to the research results reported within this paper. URL: http://www.southgreen.fr"

South Green : **@green_bioinfo**

I-Trop : **@ItropBioinfo**

Supports VCF with functional annotations

**Genotypes**
(VCF, PLINK, HapMap)

Sample metadata
(for easy selection)

**All using your web browser!**

EASY TO USE
CERTIFIED

Easy to install, too!

Gigwa

**NoSQL DB**
(for big data)

Filtering

**Refined datasets**

**Advanced filtering:**
- Functional annotations
- Data quality
- Missing data threshold
- Minor Allele Frequency
- Various genotype patterns
- Phenotype-based discrimination

**Open-source**

**Multi-platform**

**Species-agnostic**

Automated data transmission

Automated data transmission

**Visualization tools**

JBrowse

Flapjack

GBrowse

igv

Numerous **export** formats,
real-time **distribution graph**

**Further analysis**

SNiPlay

Galaxy

Sempéré G, Philippe F, Dereeper A, Ruiz M, Sarah G, Larmande P. *Gigwa-Genotype investigator for genome-wide analyses.*

# Merci pour votre attention !