

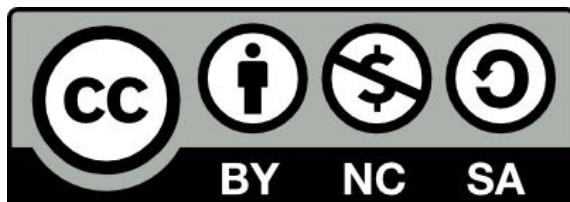
TOGGLE

Toolbox for Generic NGS analyses

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

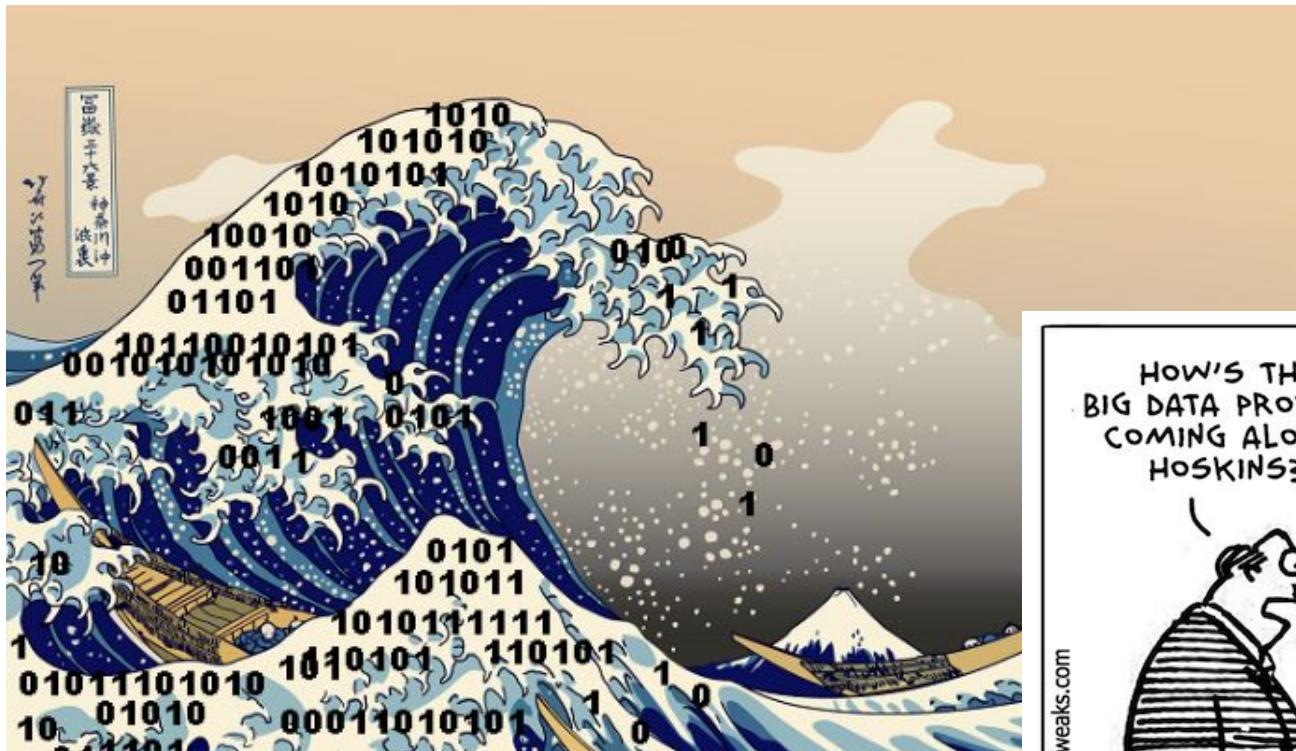
www.southgreen.fr

<https://southgreenplatform.github.io/trainings>



TOGGLE Team
toggle@ird.fr

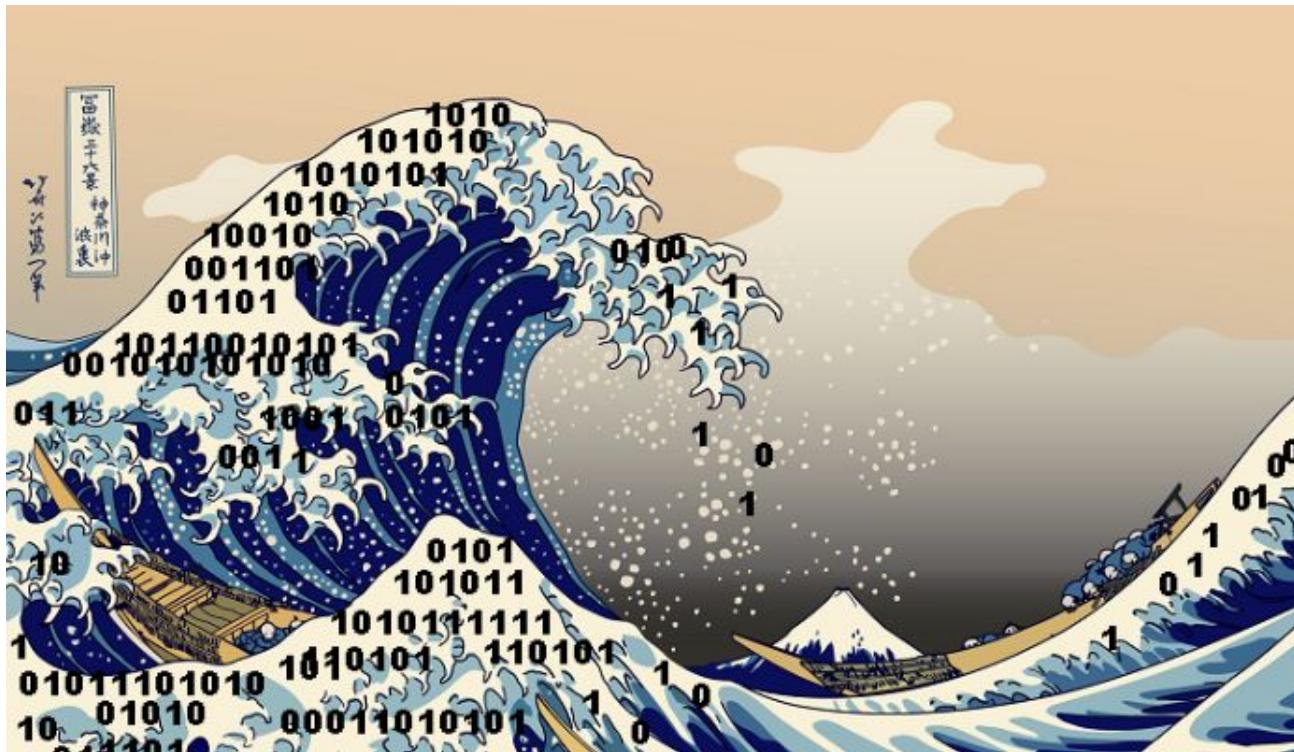
Why using workflow manager?



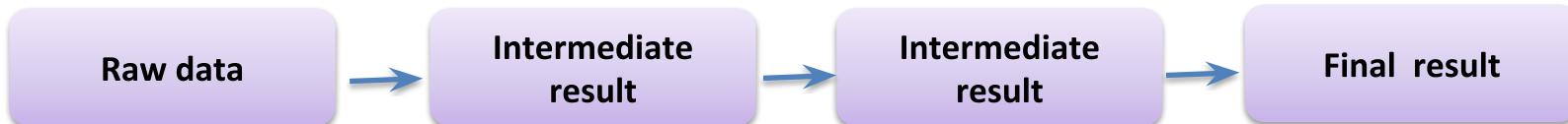
The Great Wave off Kanagawa, Hokusai @amitechsolutions.com



Why using workflow manager?



To create his own pipeline through an easy and user-friendly approach



- 3 solutions used and implemented by

GUI tools



CLI tools

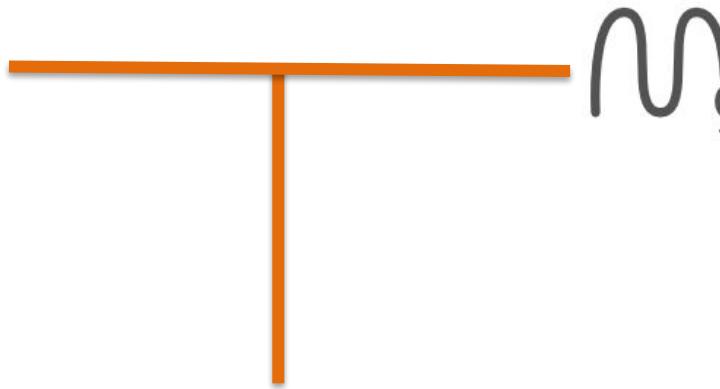


TOGGLE

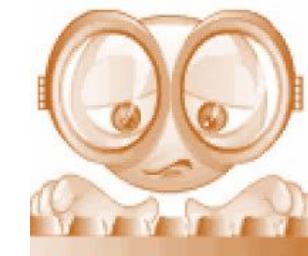


Why using workflow manager?

- 3 solutions used and implemented by



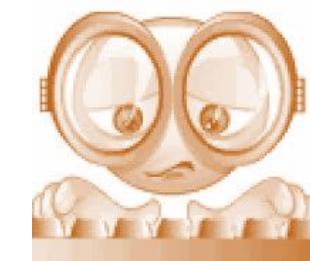
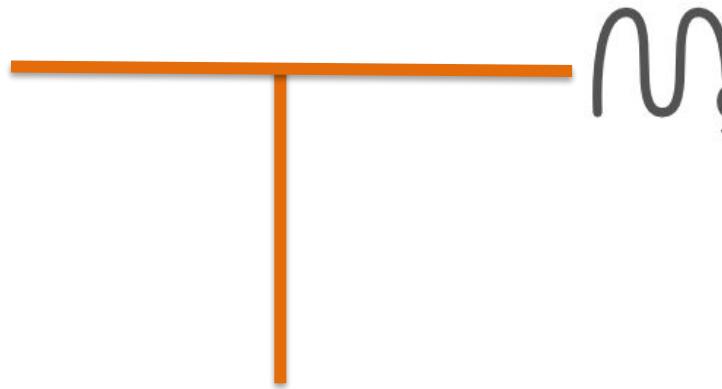
Snakemake



- Targets both biologists & bioinformaticians

Why using workflow manager?

- 3 solutions used and implemented by

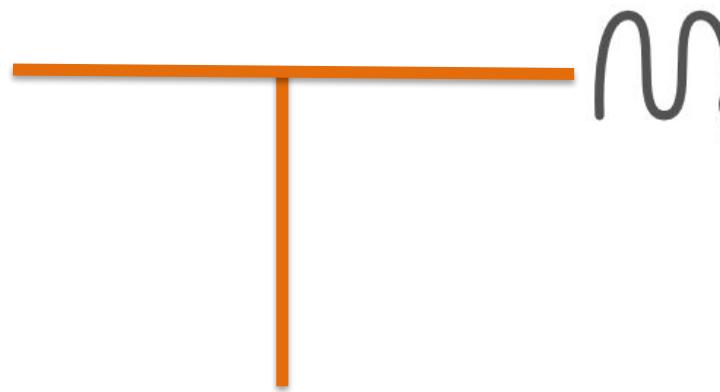


Ease of use
Well-documented
manual & workflow
examples

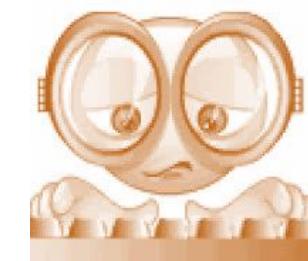
Ease of development
&
evolution

Why using workflow manager?

- 3 solutions used and implemented by



Snakemake



TOGGLE

GWAS
transcriptome assembly
structural variant detection
phylogeny

SNP detection
population genetics
genome assembly
transcriptomics
differential expression

Why using TOGGLE?

Pipeline & data
sanity controls

A robust bioinformatics
framework



File format & content
Pipeline content

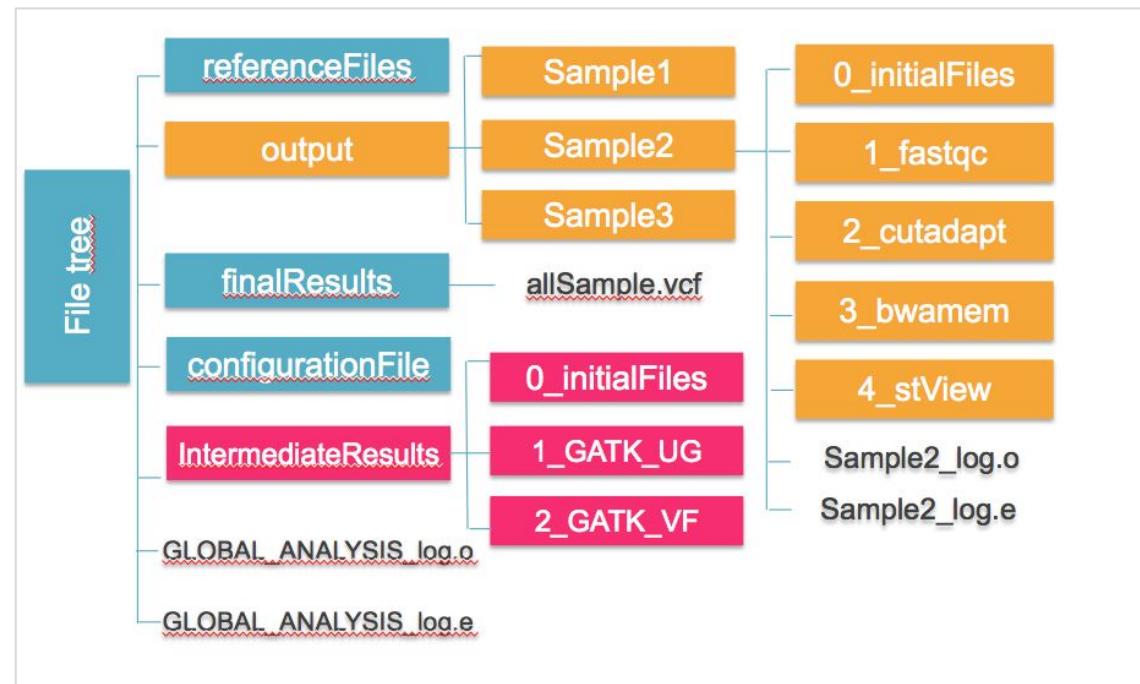


Missing but requested steps for ensuring
the pipeline running

Why using workflow manager?

Pipeline & data
sanity controls

Reproducibility
& Traceability



Why using workflow manager?

**Pipeline & data
sanity controls**

**A robust bioinformatics
framework**

**Reproducibility
& Traceability**

**Error tracking &
reentrancy**

Why using workflow manager?

Pipeline & data
sanity controls



A robust bioinformatics
framework

Reproducibility
& Traceability

Large numbers
of sample
analyzed

Error tracking &
reentrancy

Why using workflow manager?

Pipeline & data
sanity controls

HPC & Parallel
execution

A robust bioinformatics
framework

Reproducibility
& Traceability

Large numbers
of sample
analyzed

Error tracking &
reentrancy

TOGGLE



| Interface | Command line | GUI (Web interface) | |
|--|--|--|------------------------------|
| Predefined Pipelines | SNP calling, RNASeq and WGS large scale | Metagenomics, RNASeq, SNP calling, post-analyses | |
| Number of Samples | 1 to 10000 | 1 to 50 | |
| Quota (related to infra) | Disk space “/data/projects” 500Go to 1T | IRD Cirad | 100Go data 100Go => 300Go |
| Parallelization (related to infra conf) | IRD Cirad | 300 cores 600 cores | IRD Cirad |
| Number of tools available | 120 | 16 cores / one node 200 cores | |
| Post-analyses Graphical figures | No | 500 installed (total : 5500) Yes | |

What's TOGGLE ?

TOGGLE

What is TOGGLE ?



- A toolbox to perform large-scale NGS analyses

19 modules, 120 functions
120 open-source tools

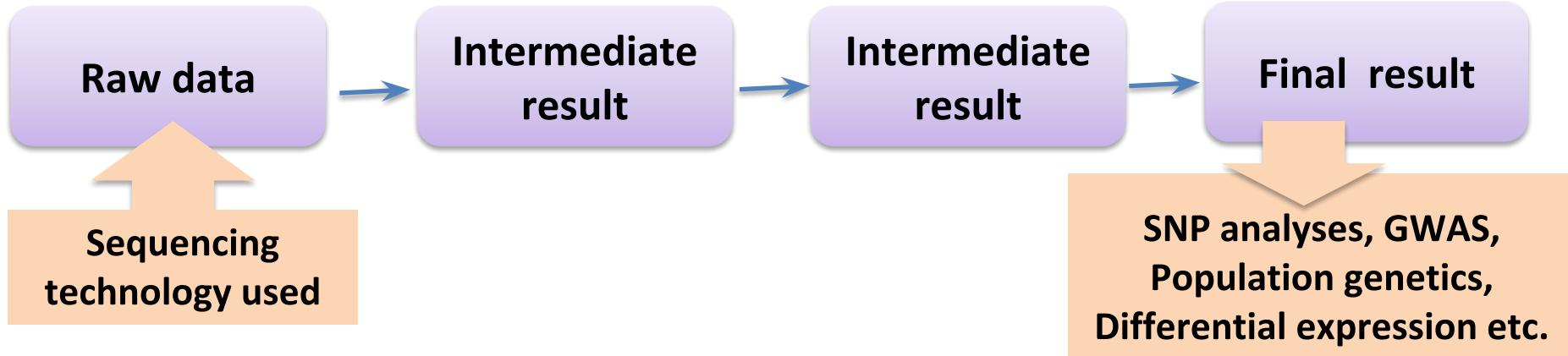


What is TOGGLE ?



- A toolbox to perform large-scale NGS analyses

19 modules, 120 functions
120 open-source tools



What is TOGGLE ?

GBS

RADSeq

RNASeq

WGS

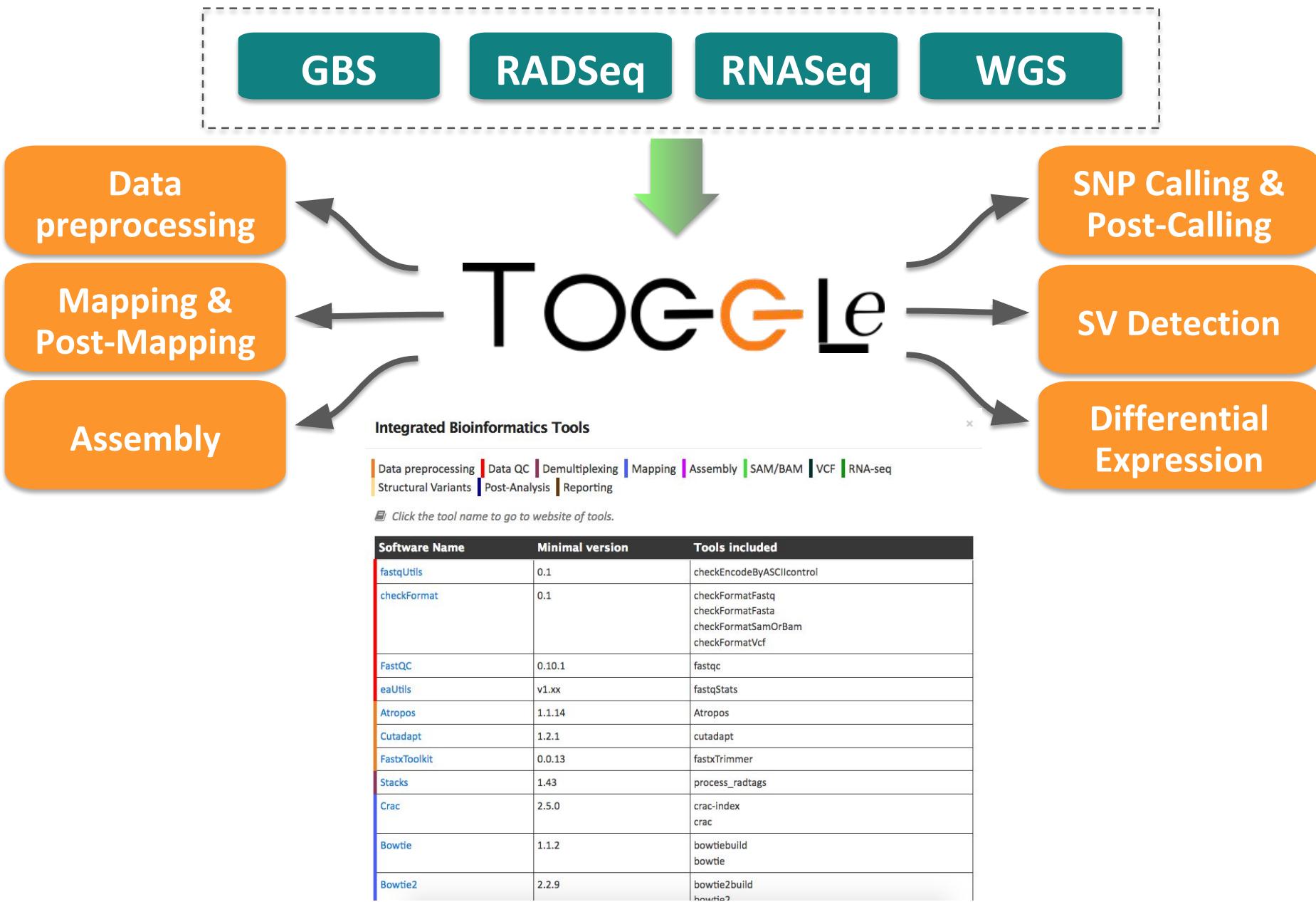


TOGGLE

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



Use TOGGLE for ?



Tools includes

Data preprocessing

Atropos, Stacks
FastQC, Cutadapt
FASTX-Toolkit

Mapping

Bowtie, Bowtie2
BWA, crac
Tophat2

Post Mapping

NGSUtils,
picardTools
SAMtools, GATK

SV Detection

DuplicationDetector,
BreakDancer, Pindel

Differential Expression

Htseq-count,
cufflinks

SNP Calling

SAMtools, GATK
VarScan, SNPEff

Assembly

Trinity
TGI-CL
Abyss

Post Calling

BEDtools, plink, sNMF
FastME, plink, readseq

TOGGLE

<https://toggle.southgreen.fr>

A command-line based pipeline framework



A single command line

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

What does TOGGLE need to run ?

```
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 <i>software.config.txt</i> 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. |
| -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...

A simple configuration file ...

\$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
```

TOGGLE

1=fastqc
2=cutadapt
3=bwa mem
4=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

To create your own workflow

\$order

1=fastqc

2=cutadapt

3=bwa mem

4=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration

Create your own workflow

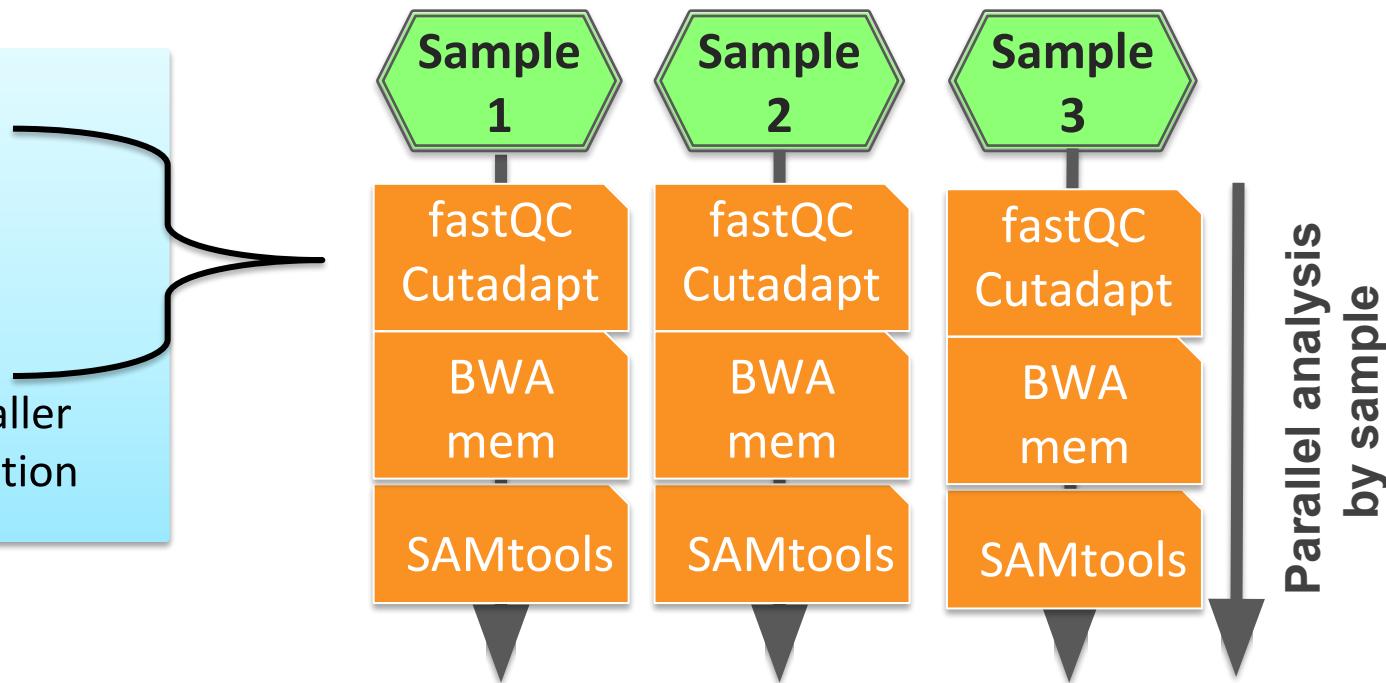
Step number < 1000

Parallel analysis by sample

To create your own workflow

\$order

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



To create your own workflow

\$order

1=fastqc

2=cutadapt

3=bwa mem

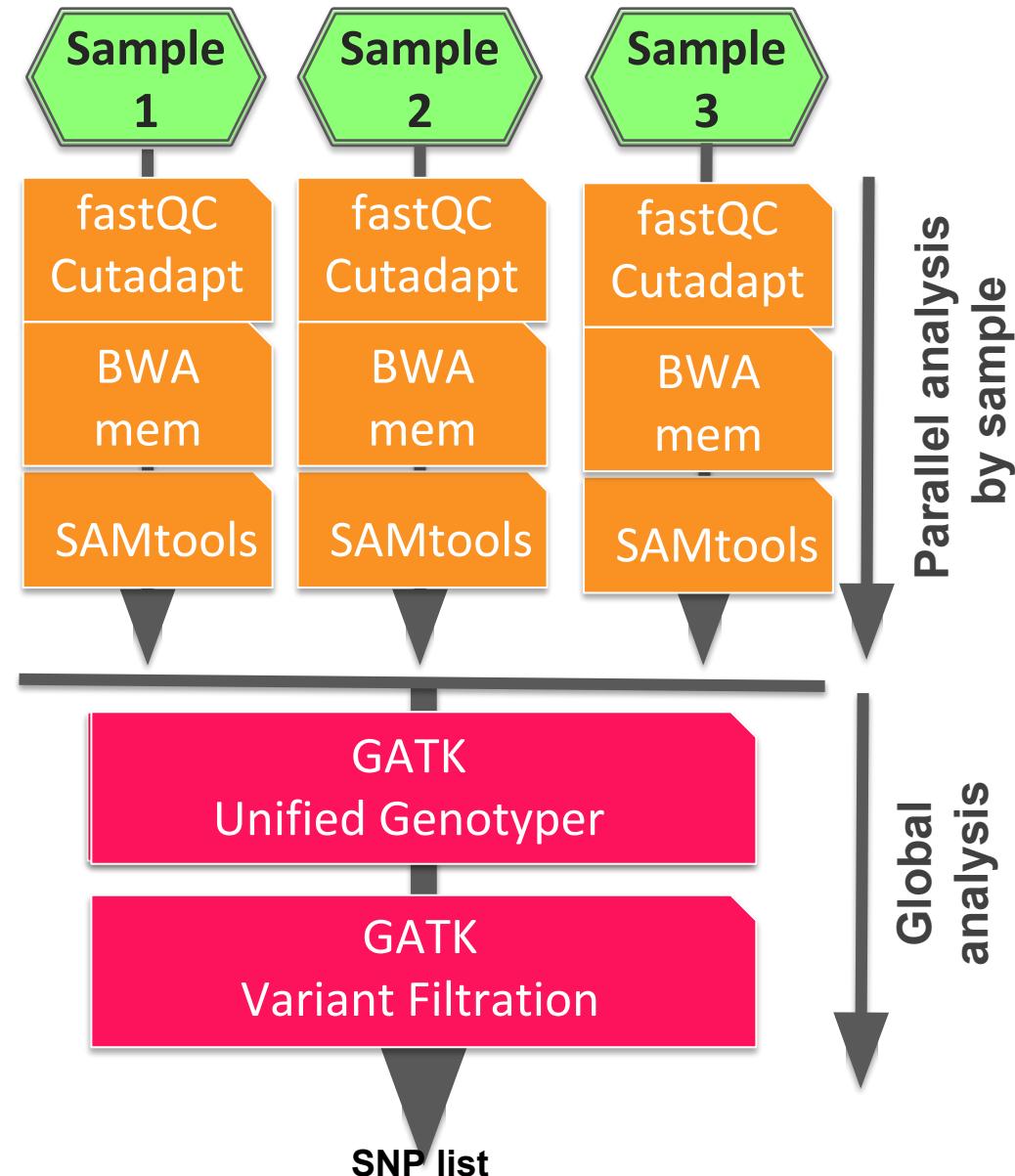
4=samToolsView

1000=gatkHaplotypeCaller

1001=gatkVariantFiltration

Step number >= 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

\$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
```

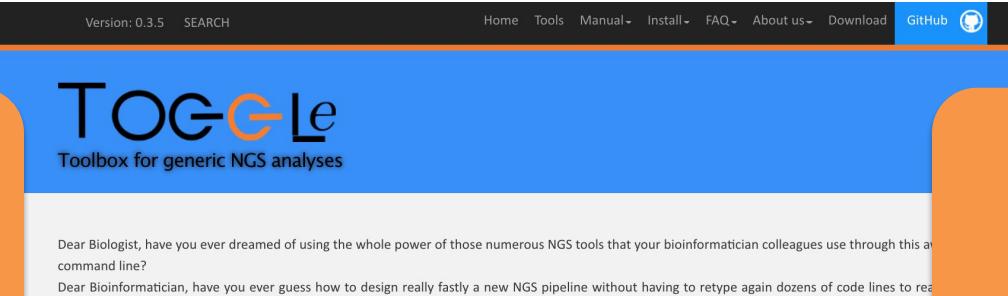
Job schedulers

LSF, MPRUN, SLURM, SGE



<http://toggle.southgreen.fr/>

User Manuals



Version: 0.3.5 SEARCH Home Tools Manual Install FAQ About us Download GitHub

TOGGLE

Toolbox for generic NGS analyses

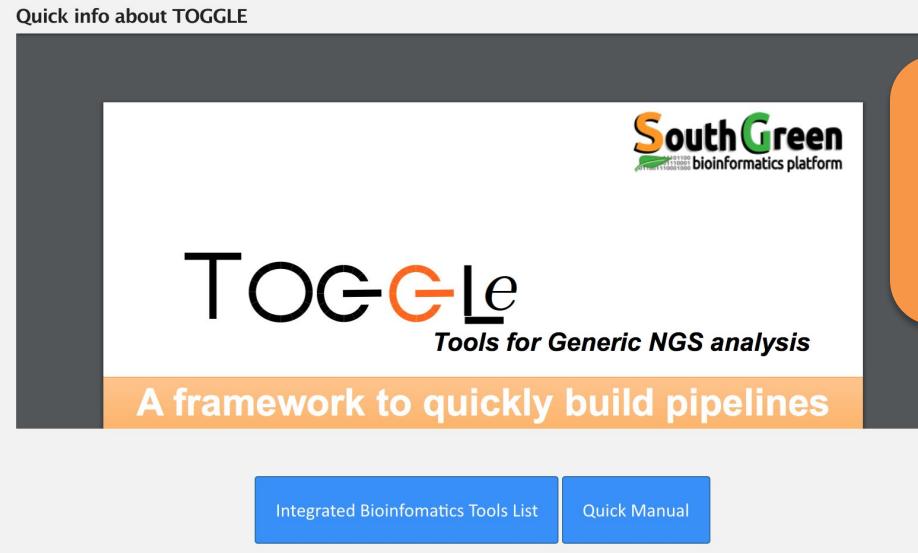
Dear Biologist, have you ever dreamed of using the whole power of those numerous NGS tools that your bioinformatician colleagues use through this command line?

Dear Bioinformatician, have you ever guess how to design really fastly a new NGS pipeline without having to retype again dozens of code lines to read scripts or starting from scratch?

So, be Happy! TOGGLE is for you!!

Screencast

Developer manual



Quick info about TOGGLE

TOGGLE

Tools for Generic NGS analysis

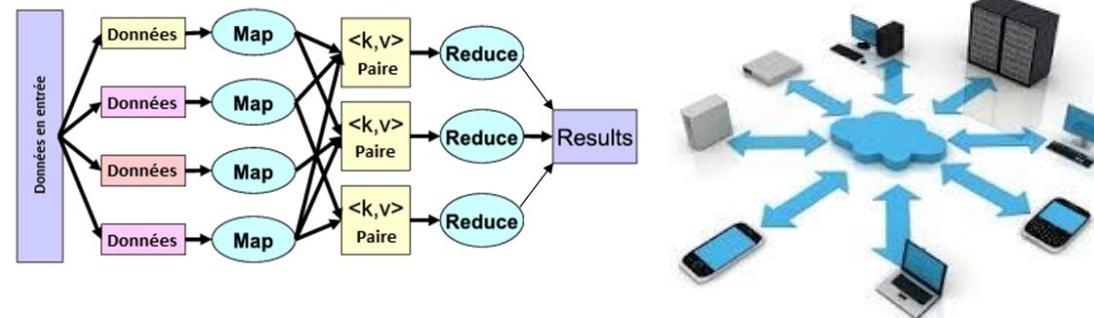
A framework to quickly build pipelines

Integrated Bioinformatics Tools List Quick Manual



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

- **New data analysis** : metagenomics, pacbio assembly, pangenome
- **New features** : complete reentrancy, non-sequential pipelines, embarrassingly-parallel approaches



TOG-ELe's team



UMR DIADE

François Sabot

Christine

Tranchant-Dubreuil

Ndomassi Tando

Alexis Dereeper

Cécile Monat

Mawussé Agbessi

Souhila Amazougarene

Abdoulaye Diallo

Laura Helou

Ayité Kougbeadjo



 **cirad**

UMR BGPI

Sébastien Ravel




ADN
Agence
d'information
génétique

Julie Orjuela-Bouniol



UMR AGAP

Cédric Farcy

Enrique Ortega-Abboud

Gautier Sarah

Maryline Summo

 **cirad**

 **INRA**
SCIENCE & IMPACT


South Green
bioinformatics platform
1101100
01110001
01100110001000