



Institut Pasteur

Sequana: a set of flexible genomic pipelines for processing and reporting NGS analysis

Thomas Cokelaer - Dimitri Desvillechabrol

Institut Pasteur

Oct 19th 2017, ENS, Paris (séminaire bioinfo)

Motivation

Jan 2015: provide NGS pipelines to Biomics sequencing platform
<https://research.pasteur.fr/en/team/biomics/> (Institut Pasteur)

- Genomics: QC + variant calling + de-novo
- Transcriptomics: RNA-seq + ChIP-seq
- Metagenomics
- Illumina but also Pacbio long reads technologies

How ?



a glue language, a scientific language



a pipeline framework mixing Python and Makefile

Köster, Johannes and Rahmann, Sven. Snakemake - A scalable bioinformatics workflow engine. Bioinformatics 2012.



Dedicated standalone such as genome coverage characterisation or a graphical user interface for Snakemake pipelines (Sequanix).

Snakemake as a workflow manager



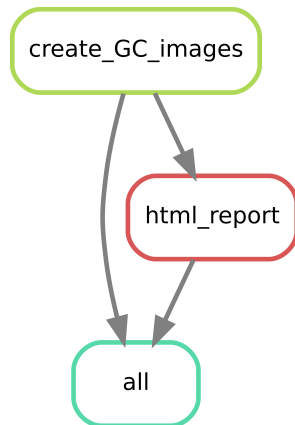
A minimalist snakemake workflow example (GC content)

```
# list of FastQ files without extension
files = ["A", "B", ....]

rule all:
    input: expand("{data}.png", data=files)
           "index.html"

rule create_GC_images:
    input: "{data}.fastq.gz"
    output: "{data}.png"
    run: # CODE TO COMPUTE IMAGES

rule html_report:
    input: expand("{data}.png", data=files)
    output: "index.html"
    run: # CODE TO CREATE HTML
```



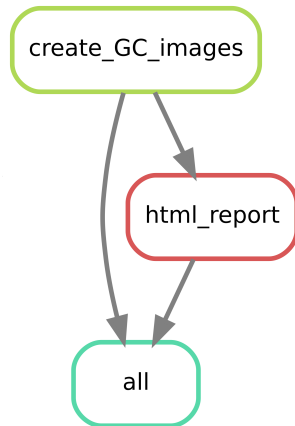
A minimalist snakemake workflow example (GC content)

```
# Use an external configuration file
import glob
files = glob.glob(config['gc']['directory'])
files = [f.replace(".fastq.gz", "")
         for f in filenames]

rule all:
    input: expand("{data}.png", data=files)
           "index.html"

rule create_GC_images:
    input: "{data}.fastq.gz"
    output: "{data}.png"
    run: # CODE TO COMPUTE IMAGES

rule html_report:
    input: expand("{data}.png", data=files)
    output: "index.html"
    run: # CODE TO CREATE HTML
```



configuration file example in YAML format

```
#####  
# Input parameters for the fractal analysis  
#  
# :Parameters:  
#  
# - size: output image size formatted as NxM where N and M  
#       are integers  
# - depth: a integer (e.g. 200)  
# - zoom: a positive value e.g. 0.5  
# - N: number of random sets  
gc:  
  - window: 100  
  - directory: /home/user/fastq_files
```


Execution

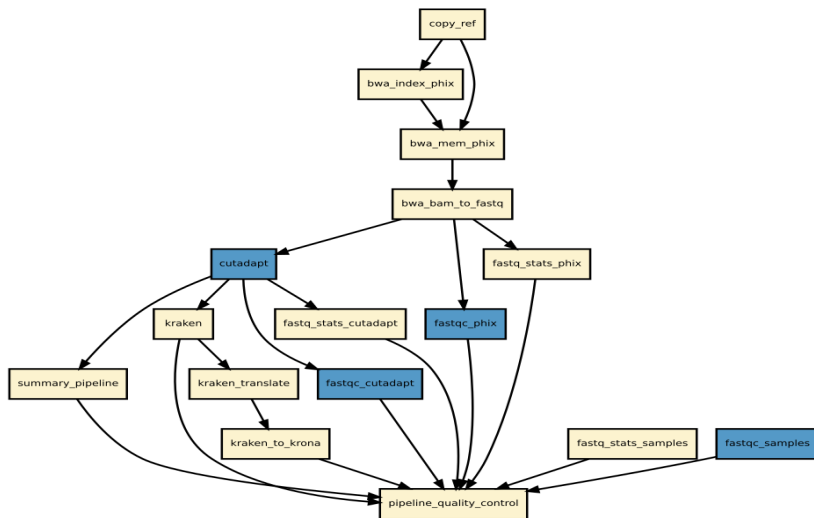
From a shell:

```
snakemake -s gc_minimalist.rules
```

More options if a configuration file is required, or execution is on a cluster, or ... something goes wrong.

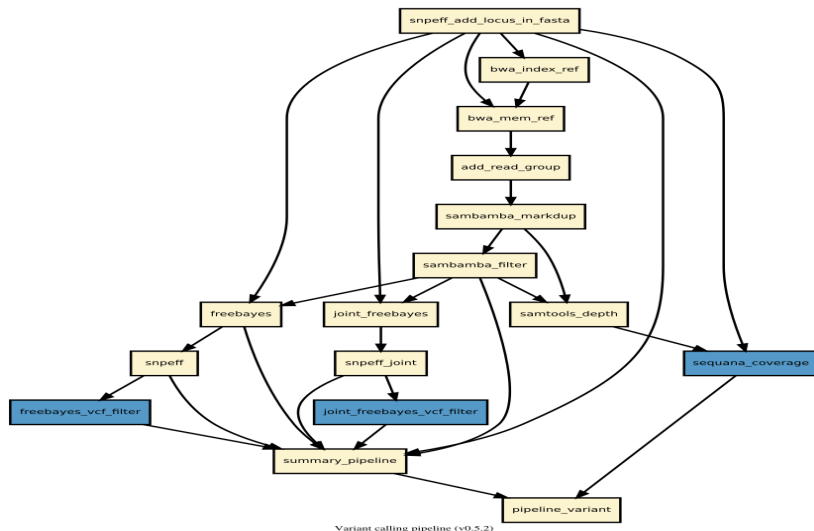
Sequana pipelines (an overview)

Pipeline example: quality control pipeline

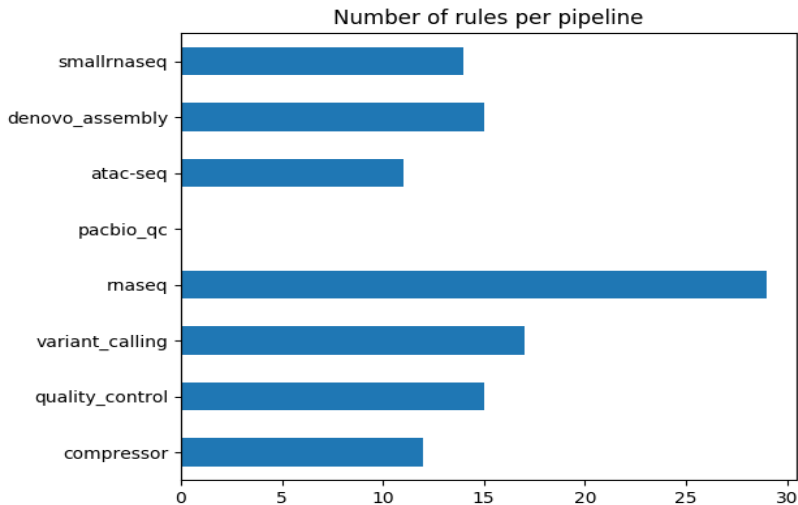


Quality control pipeline (sequana v0.5.2)

Pipeline example: variant calling



Pipeline complexity



Modularization: Factorise and reuse rules

Local standard rules

```
include: "path_to_rule_file"
```

Sequana rules

```
from sequana import snakertools as sm  
include: sm.modules['rulegraph']
```

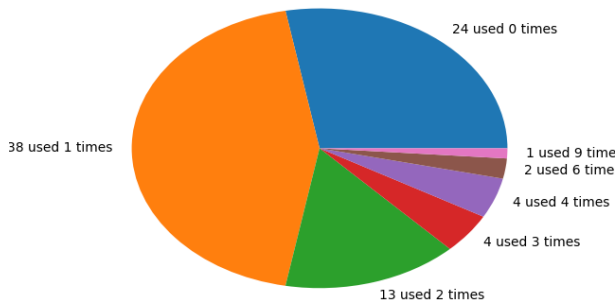
Dynamic rules:

```
sm.init("quality_control.rules", globals())  
with open(sequana.modules["fastqc_dynamic"], "r")  
    exec(dynrule.read())
```

```
manager = sm.PipelineManager("quality_control",  
                             config)
```

```
include: fastqc_dynamic("example1", manager)  
include: fastqc_dynamic("example2", manager)
```

Factorization



Once upon time there was a pipeline ...

Once upon time there was a pipeline ...

The Snakemake file assembles the rules together but all parameters can be externalised in a configuration file.

Once upon time there was a pipeline ...

The Snakemake file assembles the rules together but all parameters can be externalised in a configuration file.

One need to edit the configuration file ... without typos

Once upon time there was a pipeline ...

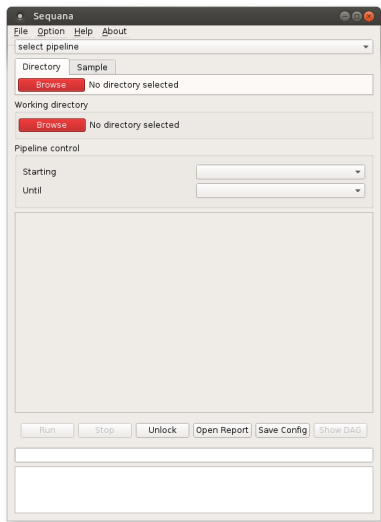
The Snakemake file assembles the rules together but all parameters can be externalised in a configuration file.

One need to edit the configuration file ... without typos

One need to launch the Snakemake command ... without typos

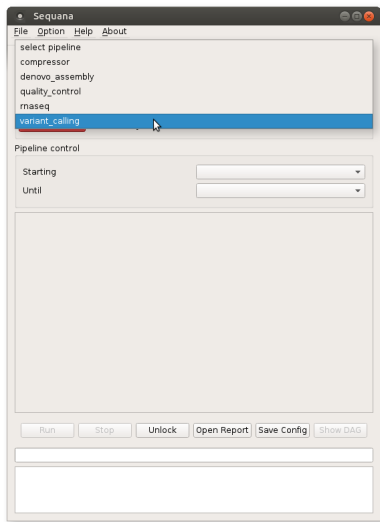
Sequanix

GUI to simplify the usage of snakemake



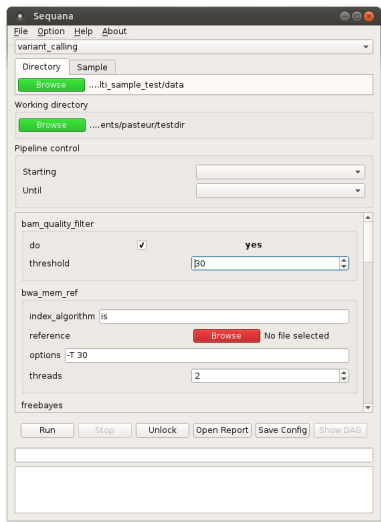
- Interface developed with PyQT5 and python
- Wrap our snakemake pipelines to ease the usage
- Usable on our cluster, which allows X11

GUI to simplify the usage of snakemake



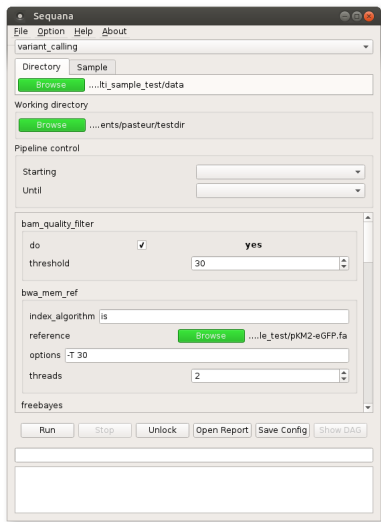
① Choose a pipeline

GUI to simplify the usage of snakemake



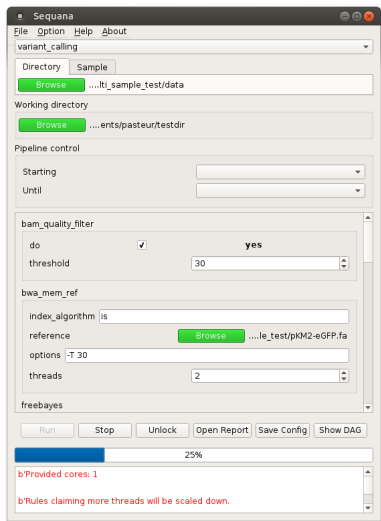
- 1 Choose a pipeline
- 2 Set input and output

GUI to simplify the usage of snakemake



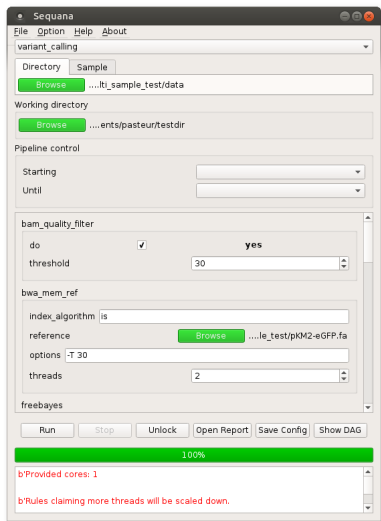
- 1 Choose a pipeline
- 2 Set input and output
- 3 Fill the config formular

GUI to simplify the usage of snakemake



- ① Choose a pipeline
- ② Set input and output
- ③ Fill the config formular
- ④ Run the pipeline

GUI to simplify the usage of snakemake



- ① Choose a pipeline
- ② Set input and output
- ③ Fill the config formular
- ④ Run the pipeline
- ⑤ Finished !

Reference

Sequanix: A Dynamic Graphical Interface for Snakemake Workflows

Dimitri Desvillechabrol, Rachel Legendre, Claire Rioualen, Christiane Bouchier, Jacques van Helden, Sean Kennedy, Thomas Cokelaer

<https://www.biorxiv.org/content/early/2017/07/12/162701>

Sequana coverage

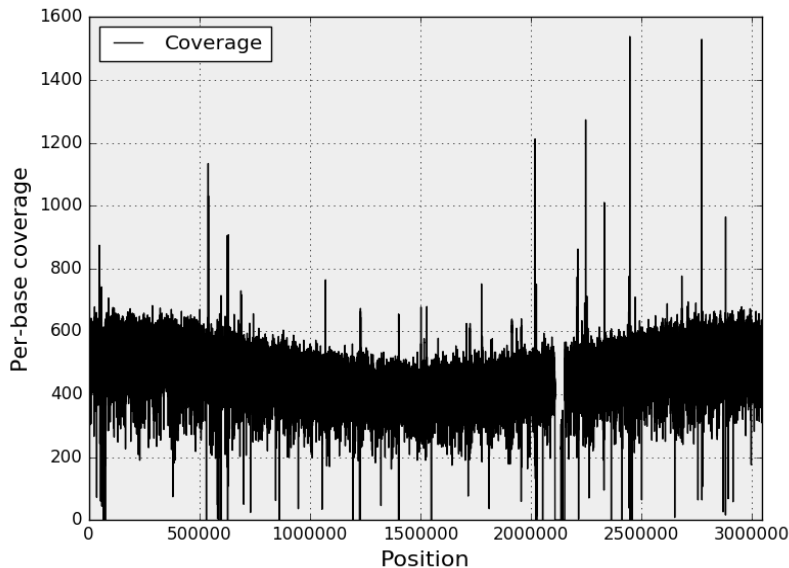
Genome coverage

Definition: The number of reads mapped to a specific position, b , within the reference genome.

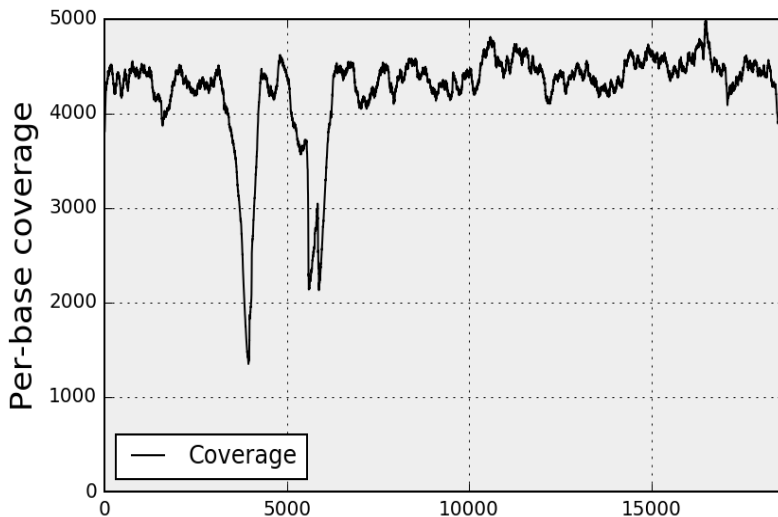
Notation: $C(b)$ also denoted C_b

Theoretical distribution: Poisson distribution but in practice over dispersed. The poisson parameter is distributed according to a Gamma hence leading to a negative binomial (See e.g., Linder et al 2013).

Bacteria case (low/high μ components and del. region)



Virus case



Question: how to automatically detect and characterise under and over covered genomic regions

1. The algorithm

1. Detrending (running median)
2. Mixture model estimation (Gaussian approximation)
3. Set a statistics (z-score)
4. Clustering (double threshold)

1. Detrending

We divide C_b by its moving average (MA), or even better its running median (RM) defined as

$$RM_W(b) = \text{median}(C(b - V), \dots, C(b + V))$$

W is the running window and $V = (W - 1)/2$.

The normalised genome coverage

$$\tilde{C}_b = \frac{C_b}{RM_W(b)}$$

1. Detrending

We divide C_b by its moving average (MA), or even better its running median (RM) defined as

$$RM_W(b) = \text{median}(C(b - V), \dots, C(b + V))$$

W is the running window and $V = (W - 1)/2$.

The normalised genome coverage

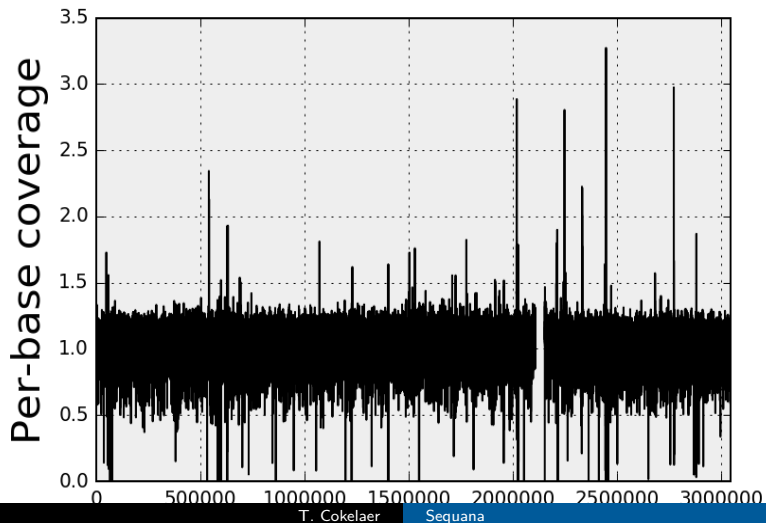
$$\tilde{C}_b = \frac{C_b}{RM_W(b)}$$

Computational note

Running median is slow due to the sorting task.

Solution: a rolling window + a bisection method to insert new element in a sorted list + efficient insert/deletion in a list (B-tree) helps: **Only a few seconds to scan a 3Mb-length genome.**

Normalised coverage



2. Building a statistics

Definition: the normalised data can be decomposed into a **central** distribution \tilde{C}^0 and a set of **outliers** \tilde{C}^1 (above and below the central distribution)

$$\tilde{C}_b = \left\{ \tilde{C}_b^0, \{ \tilde{C}_b^+, \tilde{C}_b^- \} \right\}$$

Hypothesis 1: Central distribution is predominant.

$$\left| \tilde{C}_b^0 \right| > \left| \tilde{C}_b^1 \right|$$

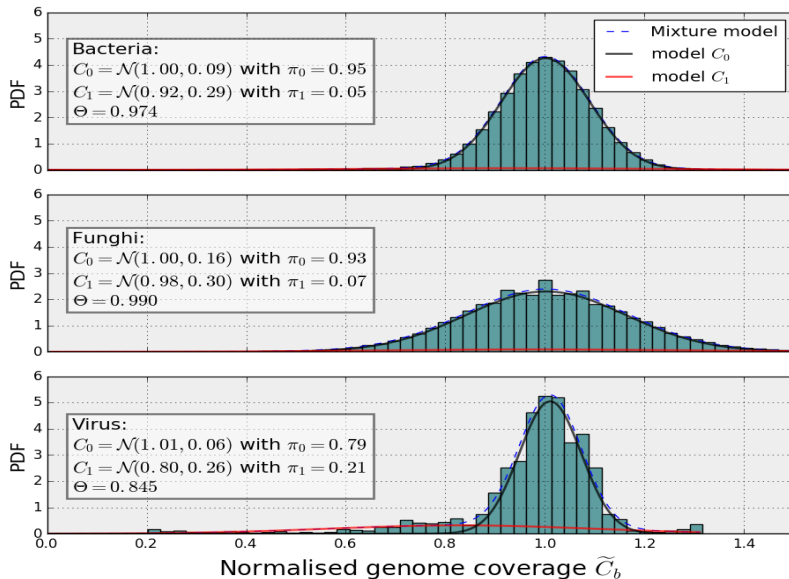
Hypothesis 2: The normalised genome coverage follows a Gaussian distribution in particular the central distribution

$$PDF(\tilde{C}_b^0) \sim \mathcal{N}(\mu_0, \sigma_0)$$

We consider regime with $\delta \gg 1$. However, the algorithm works for low values: $\delta \sim 5 - 10$.

The central distribution C_0 will be fitted to a Gaussian distribution while the outliers will be fitted to another distribution.

- With an EM algorithm using $k = 2$ distributions we can estimate the parameters.
- The parameters of the outliers components are not used.
- The average of the central distribution is 1 (by definition)
- Note that the average of the outliers can be around 1 as well if the weights of C_+ and C_- are equivalent



C. From a constant to adaptative z-score

Assuming that the central distribution is the null hypothesis, we can now assign a **z-score in the normalised space**:

$$z(b) = \frac{\tilde{C}(b) - \tilde{\mu}_0}{\tilde{\sigma}_0}$$

We can replace \tilde{C}_b by its expression ($C_b/RM_W(b)$) and express C_b as a function of the running median, the $z(b)$ and the parameters of the central distribution:

$$C(b) = (\tilde{\mu}_0 + z(b)\tilde{\sigma}_0) RM_W(b).$$

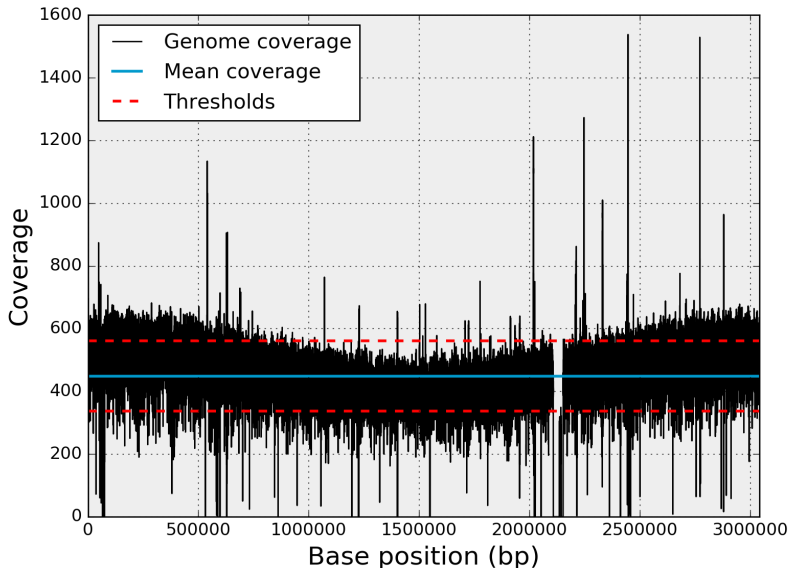
C. From a constant to adaptative z-score

We can now set a constant threshold in the normalised space (e.g. $\lambda \pm 3$) and get an adaptative threshold in the original space.

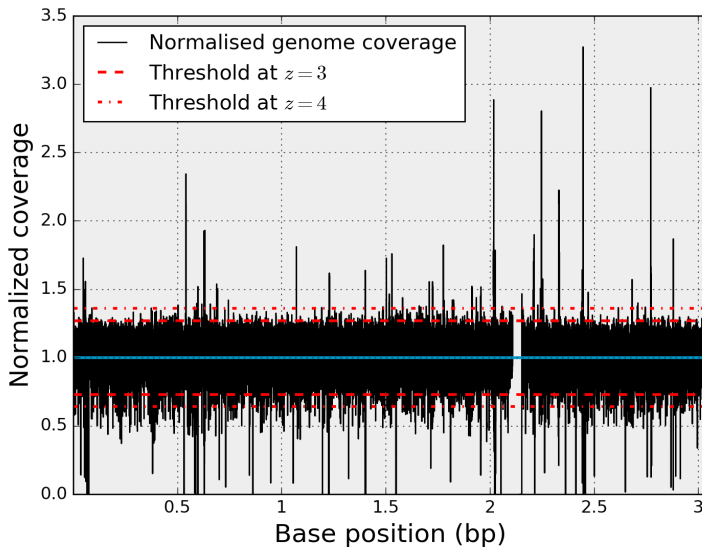
We can derive an **adaptative threshold** in the original space that is function of the genome position:

$$\tilde{\delta}^{\pm}(b) = (\tilde{\mu}_0 \pm \lambda^{\pm} \times \tilde{\sigma}_0) RM_W(b). \quad (1)$$

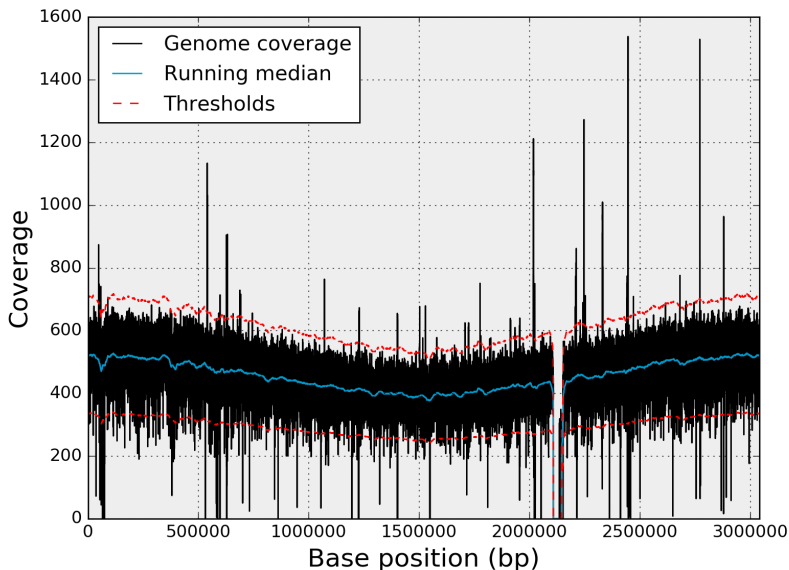
constant thresholds in the original space



constant thresholds in the normalised space



adaptative thresholds in the original space



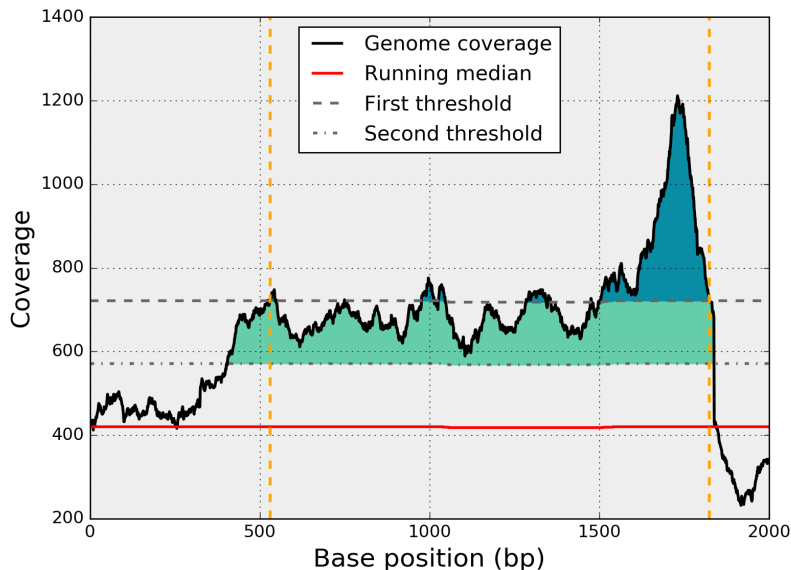
D. Regions of interest (clustering)

On a 3Mb genome with low thresholds the number of outliers are high. $\lambda = 3$ means about 4000 events by pure chance.

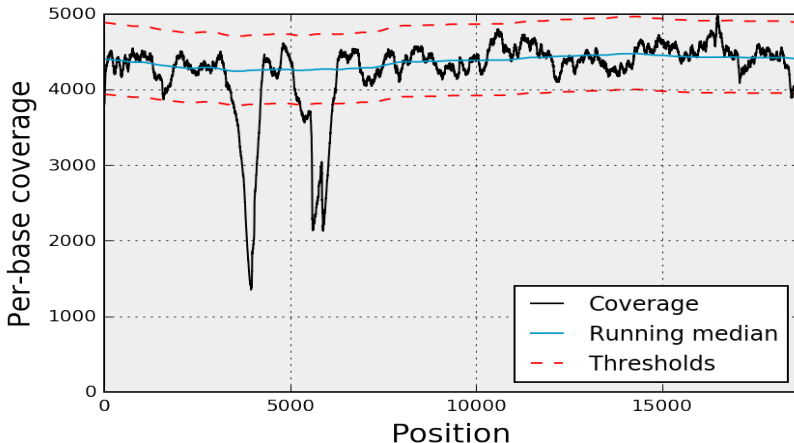
We need a strategy to reduce the number of interesting events. This is achieved by clustering the data.

Besides, to cluster close-by events further, we proceed with a double-thresholds approach.

Double thresholds



```
sequana_coverage --download-genbank JB409847  
sequana_coverage --input JB409847.bed -w 3001 --circular  
                  --genbank JB409847.gbk --show-html
```



Summary

- A robust and fast algo. to detect under/over covered regions
- The algorithm is made of 3 steps:
 1. Normalisation (running median)
 2. Set a statistics using EM to estimate mixture model
 3. Clustering of events in original space
- Implementation in Sequana as a standalone
 - HTML reports with ROIs provided as sortable tables
 - Identify repeated regions
 - A genbank can be provided to annotate the report

Dimitri Desvillechabrol, Christiane Bouchier, Sean Kennedy, Thomas Cokelaer
Detection and characterization of low and high genome coverage regions using an efficient running median and a double threshold approach. bioRxiv 092478;
 doi: <http://dx.doi.org/10.1101/092478>

Sequana: Continuous integration

Versioning, Test and Documentation



<https://github.com/sequana/sequana>



Travis CI

Continuous Integration on Travis with 185 tests with 85% coverage



Uses Sphinx (RST syntax) to document the source code and provides user guide.

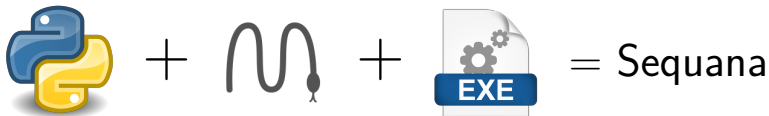


Read *the* Docs

Updated after each commits on sequana.readthedocs.io

Summary

Sequana in a nutshell



<http://sequana.readthedocs.io>

Summary

Sequana is a versatile tool that provides

- ① A Python library dedicated to NGS analysis (e.g., tools to visualise standard NGS formats).
- ② A set of snakemake workflows and rules dedicated to NGS
- ③ A GUI to execute them easily with Sequanix
- ④ HTML reports
- ⑤ Standalone applications:
 - **sequana_coverage** ease the extraction of genomic regions of interest and genome coverage information
 - **sequana_taxonomy** get a quick overview of read contents
 - ...

You like it ? Please, add a star on our github

<https://github.com/sequana/sequana/stargazers>

Stargazers

All 20

You know 7



Zhang (Frank) Che...

Joined on Mar 31, 2012

Follow



Peter Clarke

Joined on Jul 2, 2011

Unfollow



Firas

University of Cambridge

Unfollow



Hyeshek Chang

Institute for Basic Science

Follow



Peter Diakumis

@UMCCR

Follow



Raony Guimarães ...

Genomics Medicine Ireland / Torc...

Follow



Sourav Singh

Joined on May 1, 2013

Follow



EttoreZ

Joined on Sep 22, 2015

Follow



matrs

Joined on May 18, 2015

Follow



venu

Joined on Dec 19, 2014

Follow



fredericlemoine

Institut Pasteur @evolbioinfo @C...

Unfollow



mcardon

Joined on Feb 7, 2017

Follow



rlegendre

Pasteur Institut

Unfollow



Hugo Pereira

LABGeM

Follow



Kevin Murray

Australian National University (A...

Follow



Justin Fear

National Institute of Diabetes and...

Follow



Ryan Dale

National Institute of Diabetes and...

Unfollow



Brad Chapman

Harvard Chan Bioinformatics Core

Unfollow



Marco Galardini



AllanZhang

You like it ?

Join us ! add rules and pipelines !

Acknowledgements

- Dimitri Desvillechabrol (variant calling, denovo, sequana, sequanix)
- Rachel (Transcriptomics)
- Mélissa Cardon (pacbio)
- Biomix users (Institut Pasteur)

Thank you

Slides available on
http://github.com/sequana/sequana_presentations/2017/ENS