



*David Roquis, Rémi Emans, Guillaume Mitta &  
Christoph Grunau*

# Galaxy for the rest of us

*Ph.D. student, 2<sup>nd</sup> year*



Galaxy Day  
UPMC  
December 4<sup>th</sup>, 2013





# Our laboratory



UPVD

Université de Perpignan **Via Domitia**



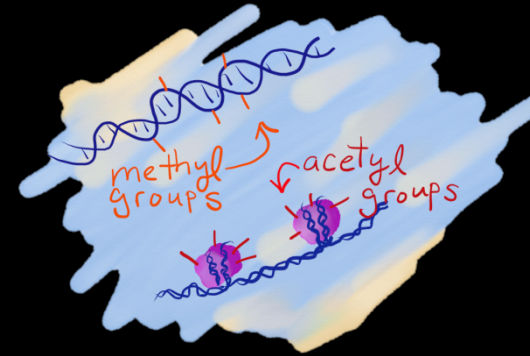


# Our laboratory



UPVD

Université de Perpignan Via Domitia



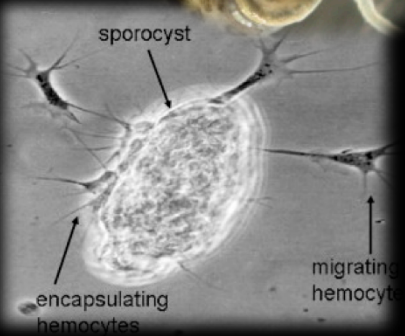
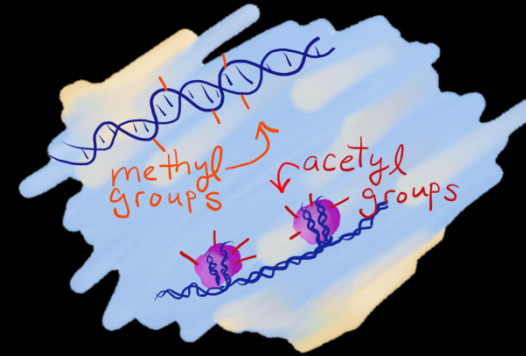


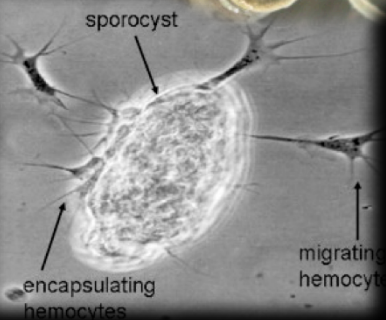
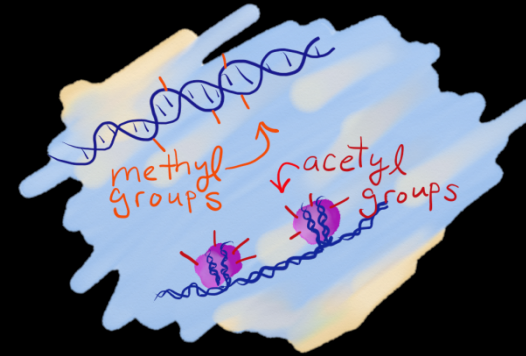
# Our laboratory



UPVD

Université de Perpignan Via Domitia





*A human parasite*



*Schistosoma mansoni*

*A freshwater snail*



*Biomphalaria glabrata*

*A tropical coral*



*Pocillopora damicornis*

*An insect*



*Dinocampus coccinellae*

+ Previously → Candidate gene approach.



+ Previously → Candidate gene approach.

+ Next Generation Sequencing more **available** and **affordable**





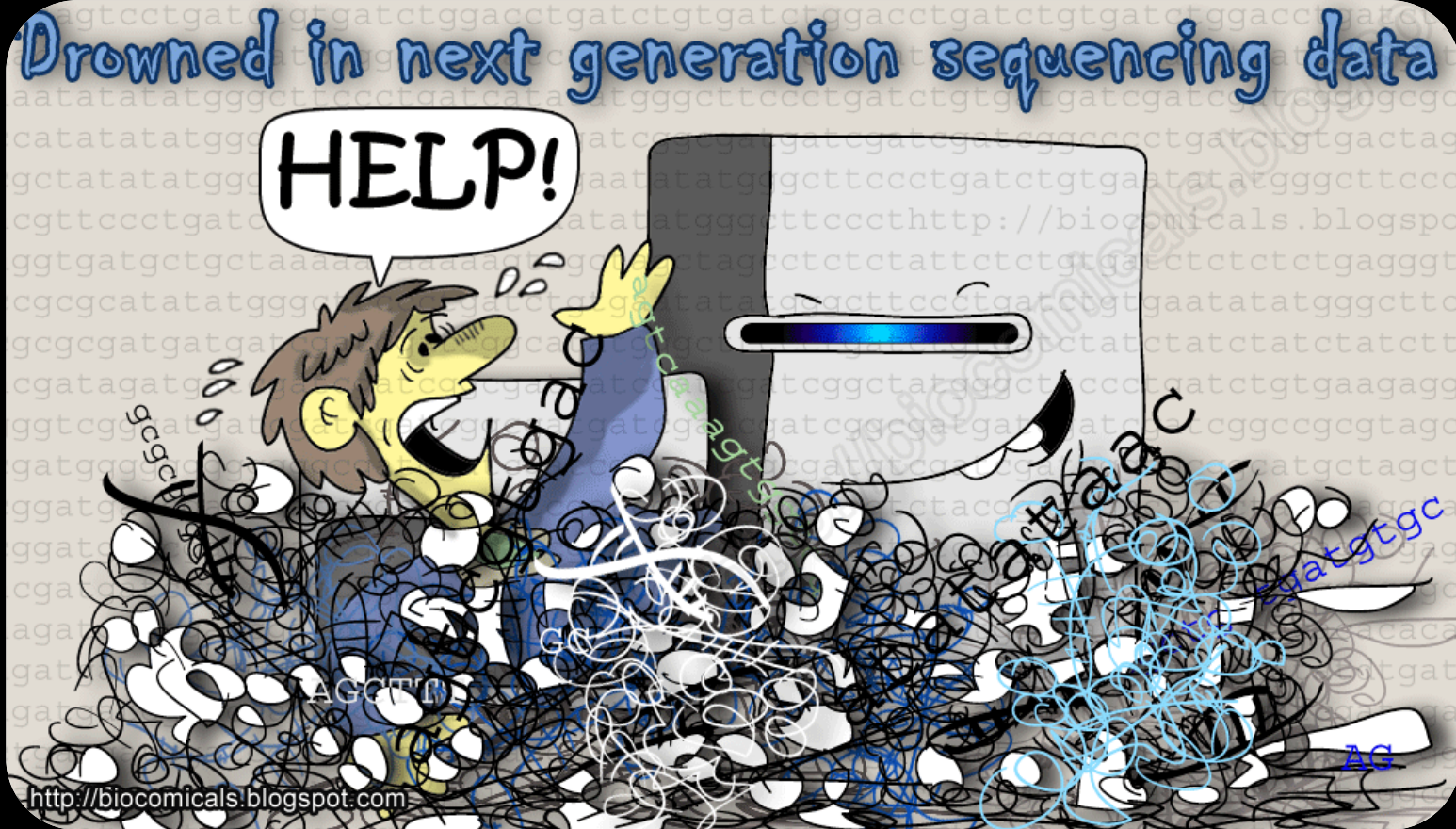
+ Previously → Candidate gene approach.

+ Next Generation Sequencing more **available** and **affordable**



+ Will to switch to genome-wide approaches to study **Genomes, Epigenomes, Transcriptomes.**

- + No bioinformatician.
- + No dedicated platform or hardware on the campus.



- + But, a few people with some bioinformatics knowledge...

 Galaxy

## But why a local instance?

- + Very different and specific needs.
- + Working on non-model organisms means that most of the tools need to be modified/adapted.
- + Flexibility.





Network

*2 x 100 mb/sec*



**Dell PowerEdge R820**  
(~5 000 €)

- + 16 cores, Xeon 2,4 GHz
- + 96 Gb RAM
- + 1.67 Tb hard drive



Network

$2 \times 100 \text{ mb/sec}$



**Dell PowerEdge R820**  
(~5 000 €)

- + 16 cores, Xeon 2,4 GHz
- + 96 Gb RAM
- + 1.67 Tb hard drive

Network

$2 \times 100 \text{ mb/sec}$

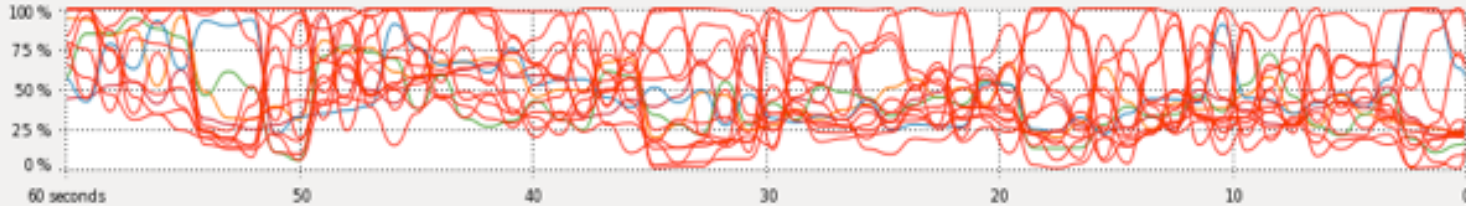
**Buffalo TerraStation 5400 WSS**

(~2 000 € X 2)

- + 20 Tb hard drive
- + Main and backup

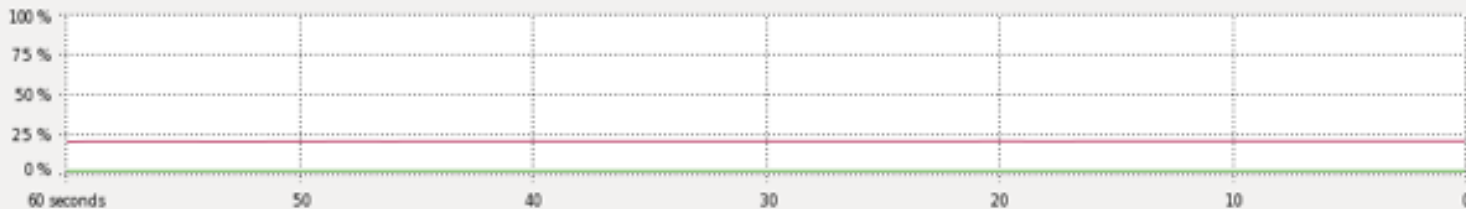


## CPU History



|             |             |             |             |
|-------------|-------------|-------------|-------------|
| CPU1 21.7%  | CPU2 34.2%  | CPU3 26.1%  | CPU4 21.8%  |
| CPU5 8.5%   | CPU6 78.8%  | CPU7 10.8%  | CPU8 19.1%  |
| CPU9 98.1%  | CPU10 5.7%  | CPU11 27.8% | CPU12 22.3% |
| CPU13 81.1% | CPU14 17.1% | CPU15 91.3% | CPU16 21.9% |

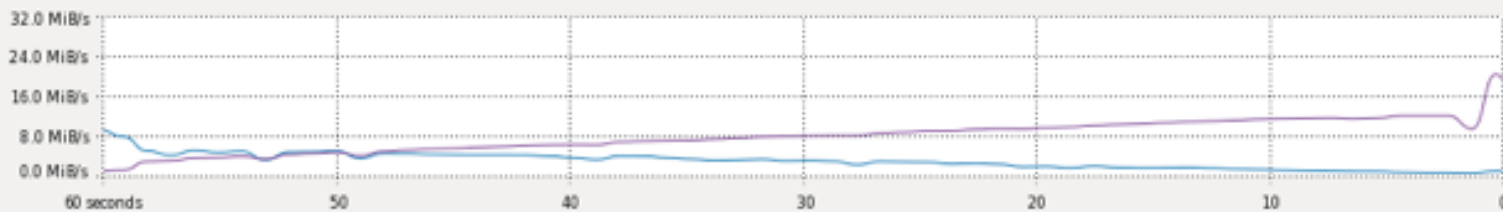
## Memory and Swap History



**Memory**  
18.5 GiB (19.6 %) of 94.6 GiB

**Swap**  
54.0 MiB (1.1 %) of 4.9 GiB

## Network History



**Receiving**  
Total Received

376.0 KiB/s  
769.0 GiB

**Sending**  
Total Sent

11.2 MiB/s  
602.1 GiB



# Data management



UPVD

Université de Perpignan Via Domitia







Local storage by  
project manager





Local storage by project manager



Data libraries by Galaxy Admin



## Data Libraries

[Advanced Search](#)

| <u>Data library name</u> ↓                          | <u>Data library description</u>  |
|---|--|
| <a href="#">B.glabrata annotation - genome</a>      | FASTA and gff files from VectorBase                                      |
| <a href="#">B.glabrata annotation - results</a>     | B.glabrata analysis results  |
| <a href="#">B.glabrata annotation - RNA-Seq I</a>   | RNA-Seq data provided by WashU (fastq f/r)                               |
| <a href="#">B.glabrata annotation - RNA-Seq II</a>  | RNA-Seq data provided by WashU (aligned BAM and assembled fasta)         |
| <a href="#">B.glabrata annotation - RNA-Seq III</a> | sorted BAM files of RNA-Seq suitable for upload to IGV (from G.Oliveira) |
| <a href="#">B.glabrata - assemblies</a>             | assemblies for read dependency   |
| <a href="#">B.glabrata RNA-Seq (clean)</a>          | RNA-Seq without adaptors, quality clipped                                |
| <a href="#">B.glabrata RNA-Seq (paired)</a>         | fastq custom upload N.Dheilly 16/01/13                                   |



Local storage by project manager



Data libraries by Galaxy Admin

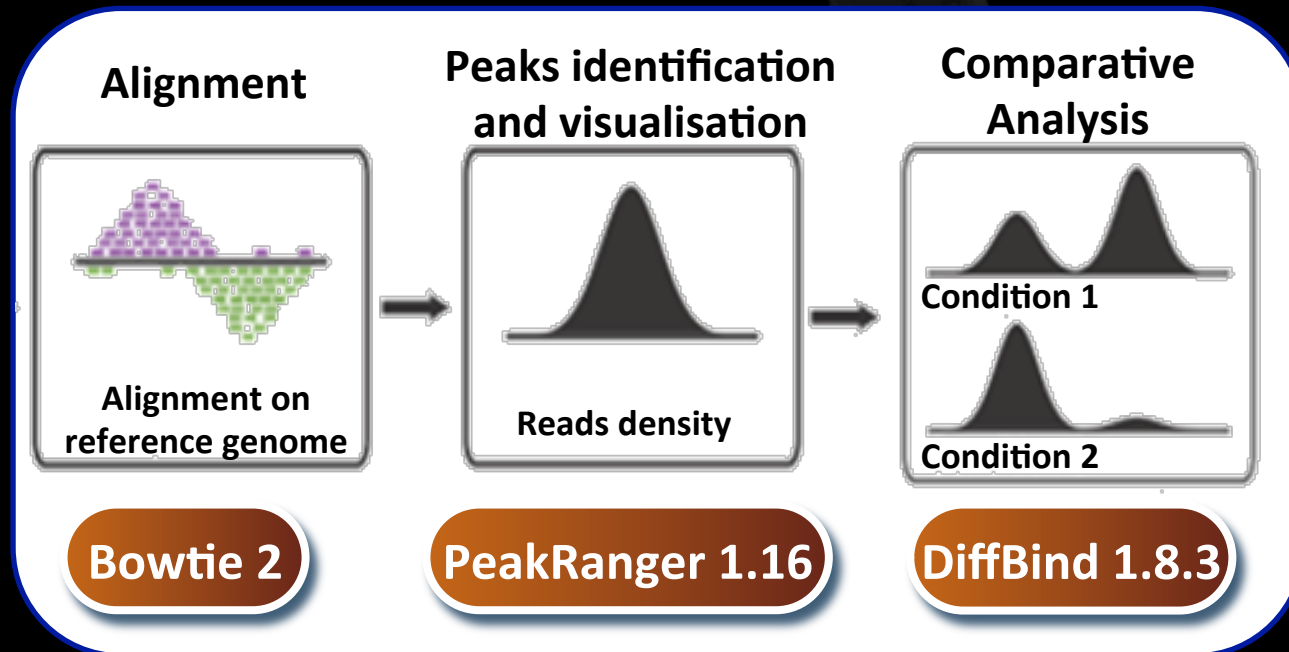


SFTP?

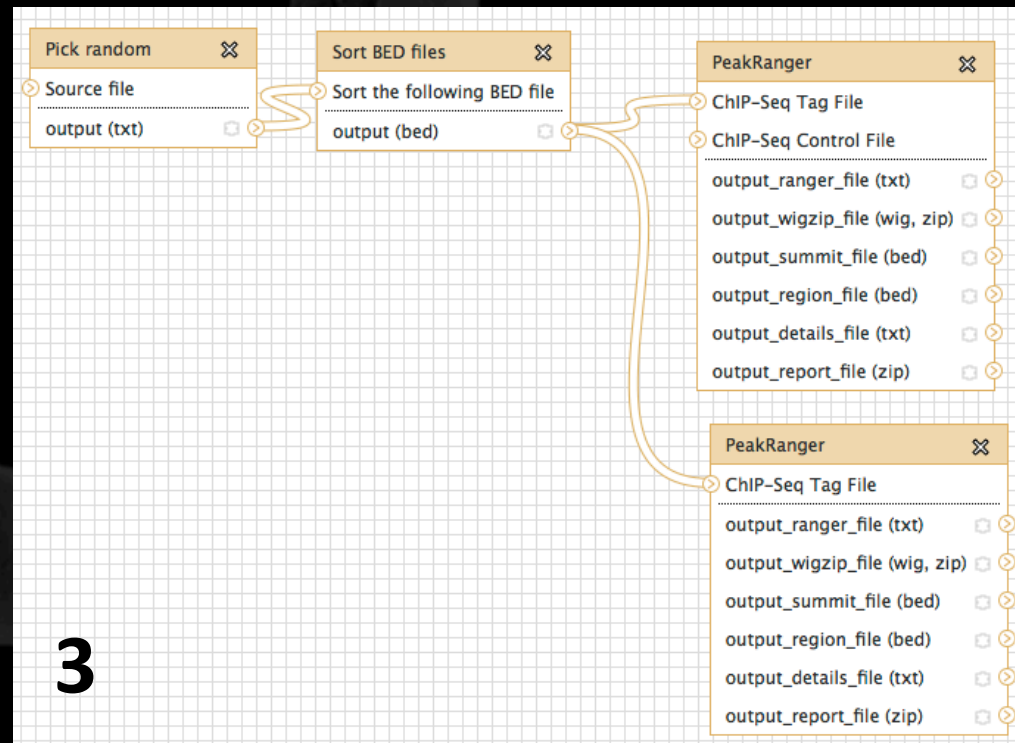
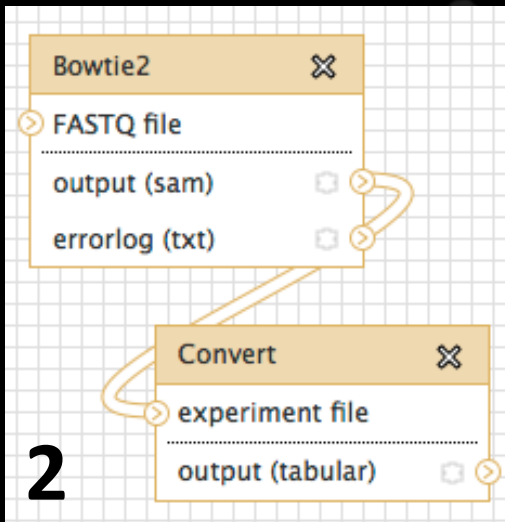
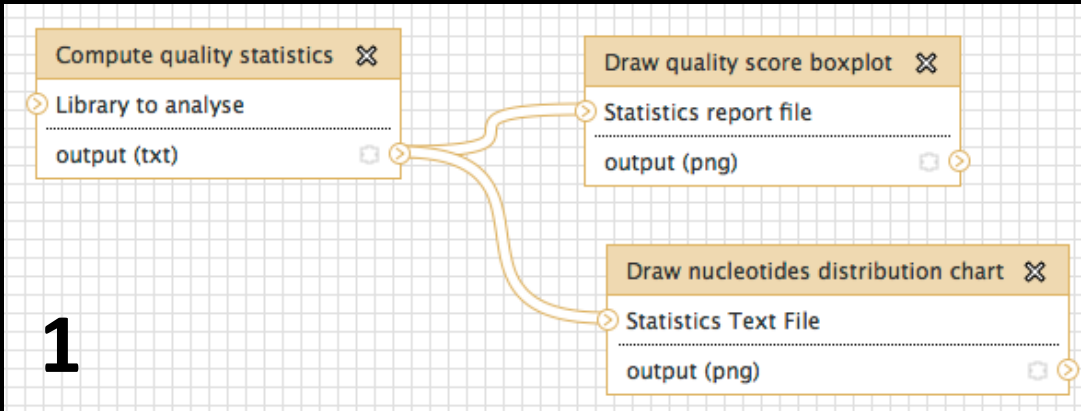
## Data Libraries

[Advanced Search](#)

| <u>Data library name</u> ↓                          | <u>Data library description</u>  |
|---|--|
| <a href="#">B.glabrata annotation - genome</a>      | FASTA and gff files from VectorBase                                      |
| <a href="#">B.glabrata annotation - results</a>     | B.glabrata analysis results  |
| <a href="#">B.glabrata annotation - RNA-Seq I</a>   | RNA-Seq data provided by WashU (fastq f/r)                               |
| <a href="#">B.glabrata annotation - RNA-Seq II</a>  | RNA-Seq data provided by WashU (aligned BAM and assembled fasta)         |
| <a href="#">B.glabrata annotation - RNA-Seq III</a> | sorted BAM files of RNA-Seq suitable for upload to IGV (from G.Oliveira) |
| <a href="#">B.glabrata - assemblies</a>             | assemblies for read dependency   |
| <a href="#">B.glabrata RNA-Seq (clean)</a>          | RNA-Seq without adaptors, quality clipped                                |
| <a href="#">B.glabrata RNA-Seq (paired)</a>         | fastq custom upload N.Dheilly 16/01/13                                   |



Adapted from Bardet et al. (2012)





# Example of application



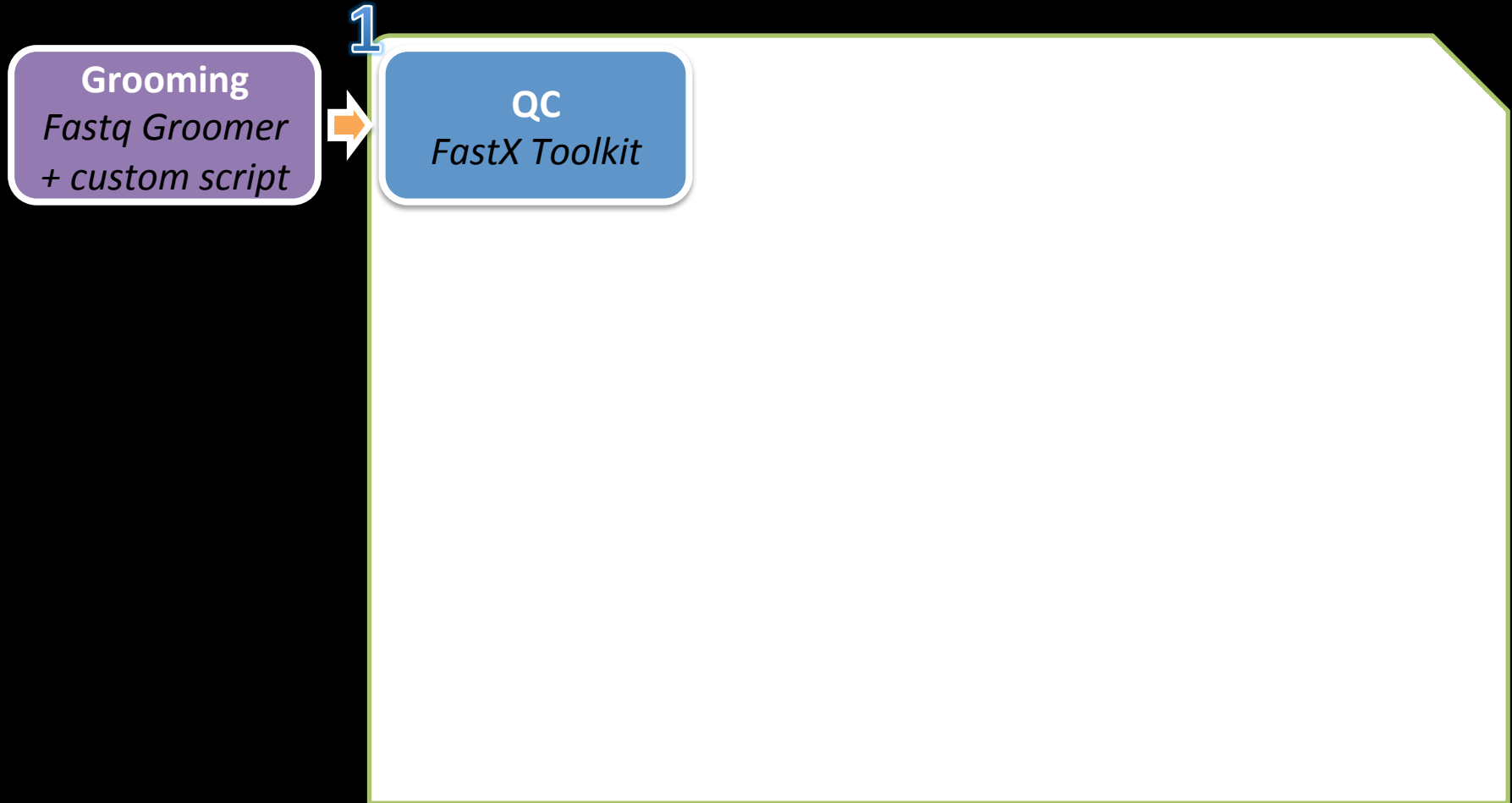
UPVD

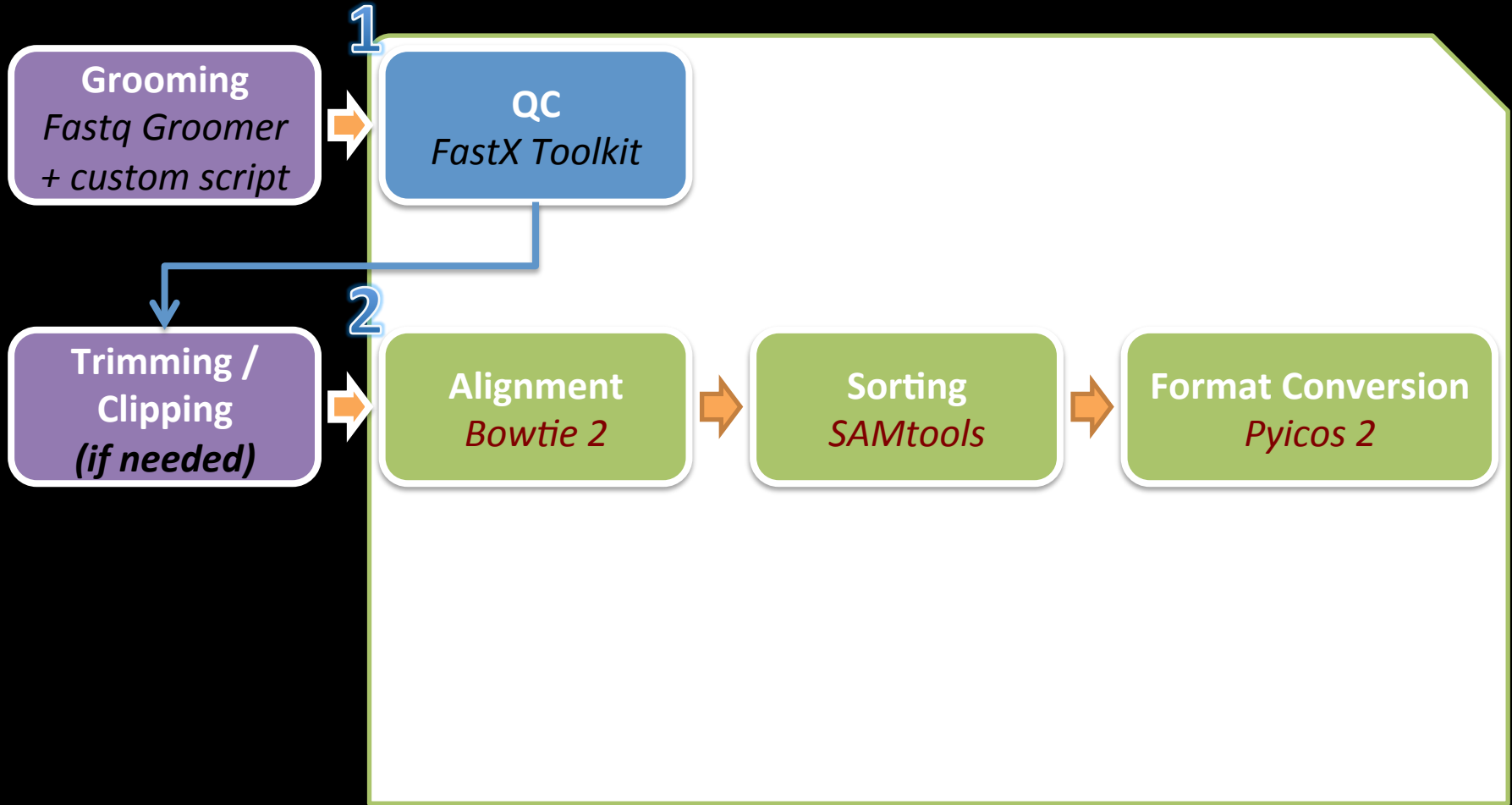
Université de Perpignan Via Domitia

## Grooming

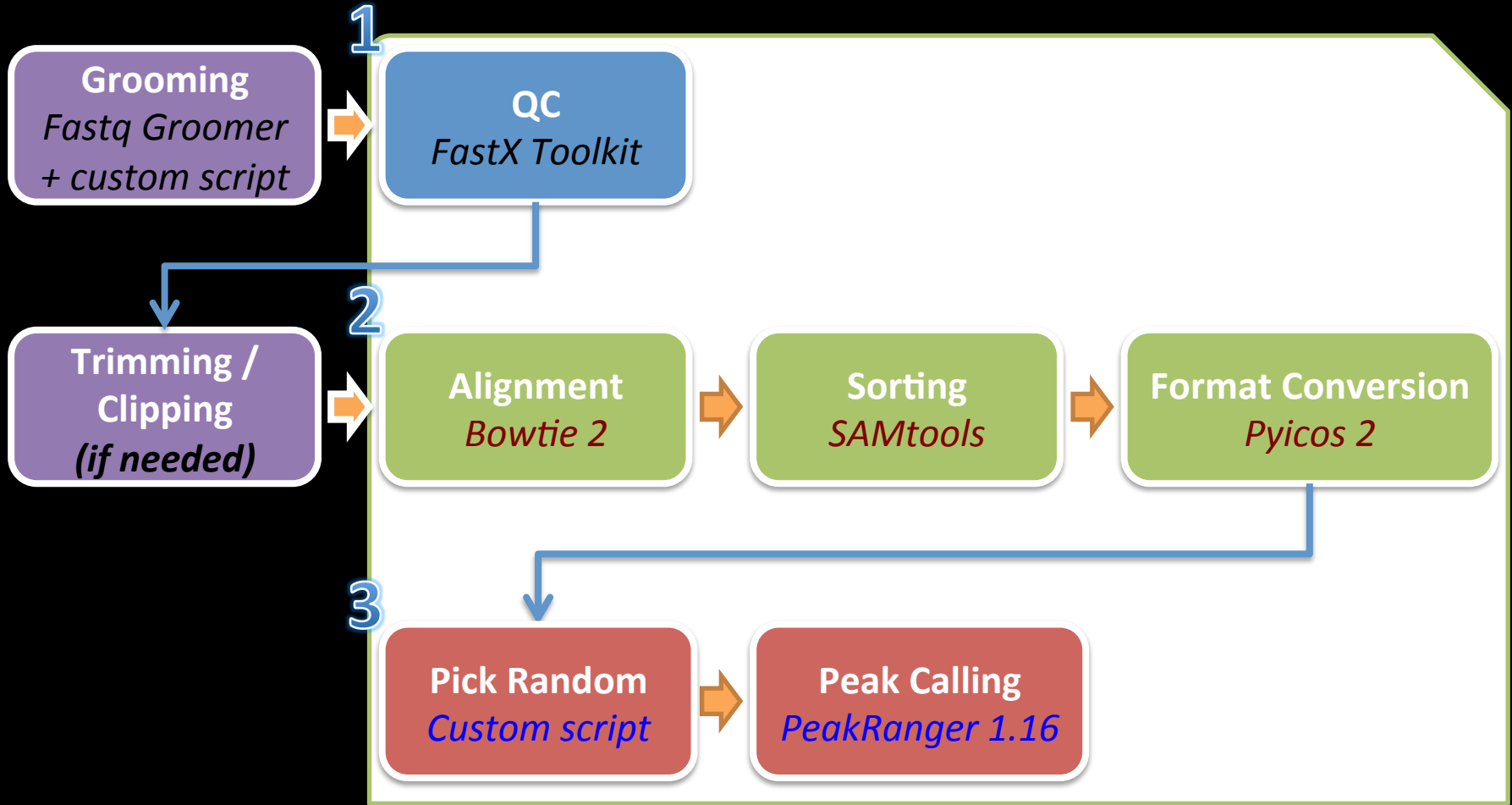
*Fastq Groomer  
+ custom script*

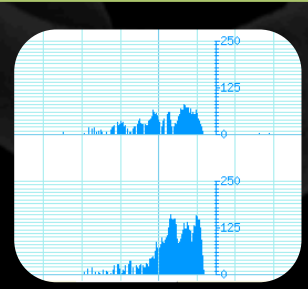
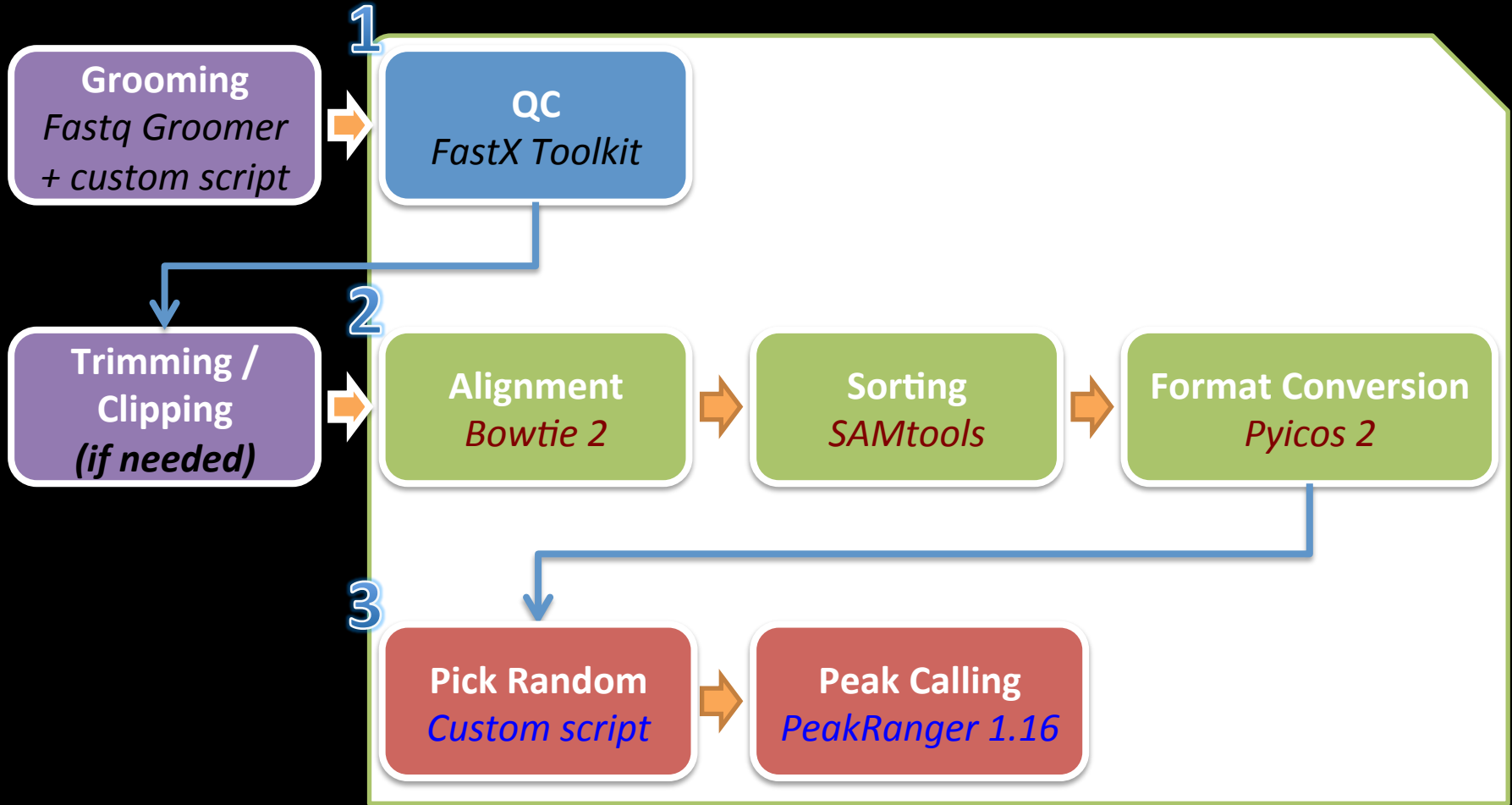


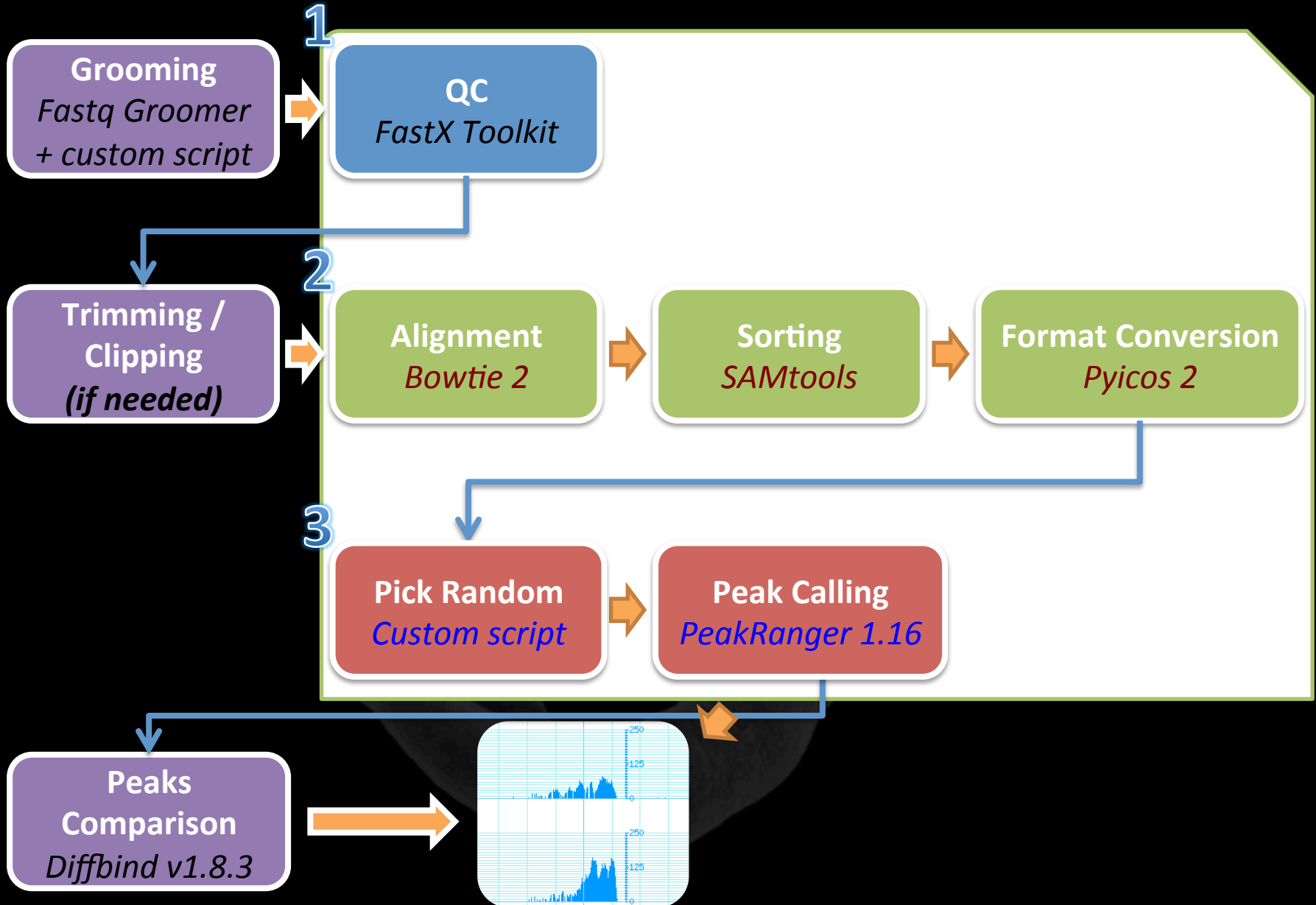






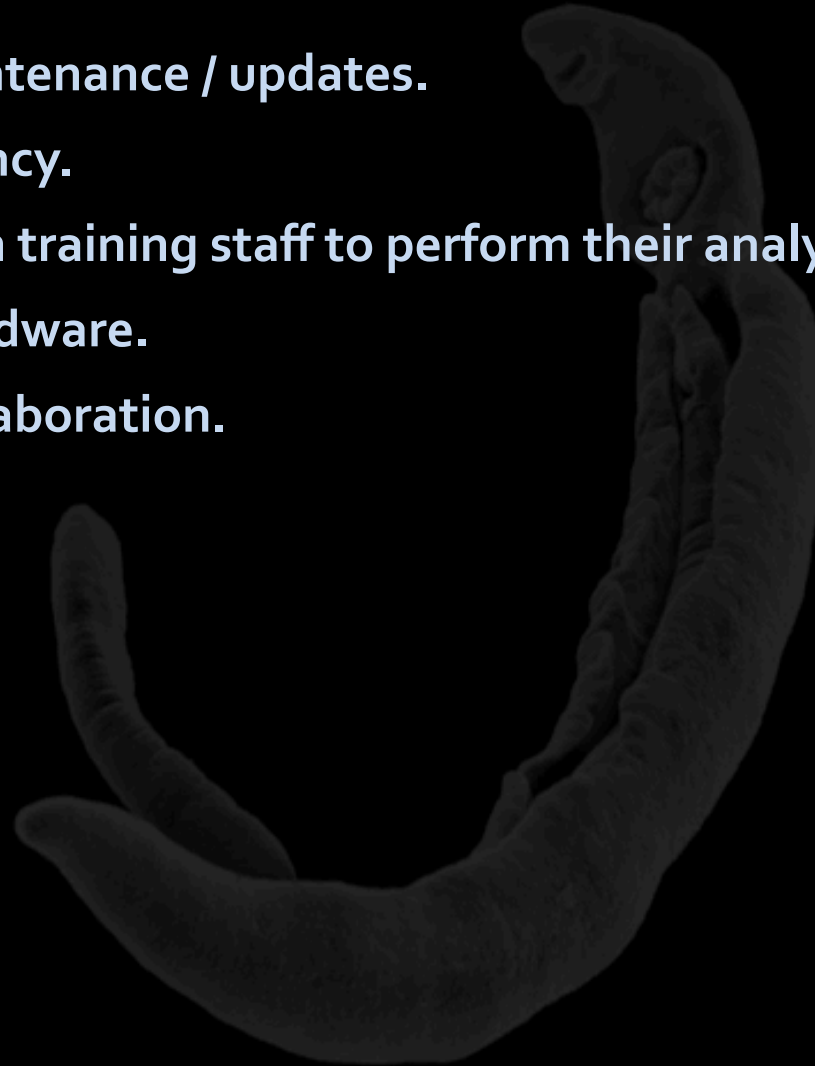






## Pros

- + Very flexible maintenance / updates.
- + No data redundancy.
- + Less time spent in training staff to perform their analyses.
- + No expensive hardware.
- + International collaboration.



## Pros

- + Very flexible maintenance / updates.
- + No data redundancy.
- + Less time spent in training staff to perform their analyses.
- + No expensive hardware.
- + International collaboration.

## Cons

- + **Dependant of the general network speed.**
- + **No data integrity verification.**



# Acknowledgments



UPVD

Université de Perpignan Via Domitia



**System Administrator**

*Rémi Emans*



**Supervisor**

*Céline Cosseau*



**Supervisor**

*Christoph Grunau*



**Lab Director**

*Guillaume Mitta*

AGENCE NATIONALE DE LA RECHERCHE



**EPIGEVOL** *ÉPIGÉnétique  
et EVOution*

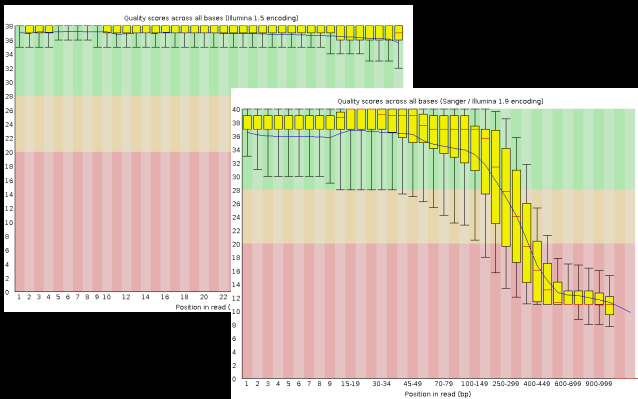


**Ecole Doctorale UPVD**  
Université de Perpignan Via Domitia

```
@HWI-ST609:177:D16HYACXX:3:1308:14623:59114 1:N:0:AGTTCC
TGAGACTAAGGCATATTATAATAATGCGCGGAAGTGTATAGTAAAAGCAC
+
@CCFFFFFGHHGHIIJJIJJJJJJJIJJIGGJIIIIIGFIIIIIGIJJIIIBHI
```



+ FastQC



+ Store clean data

**Data Libraries**

search dataset name, info, message, dbkey

Advanced Search

| Data library name ↓                 | Data library description   |
|-------------------------------------|--|
| B.glabrata annotation - genome      | FASTA and gff files from VectorBase                                      |
| B.glabrata annotation - results     | B.glabrata analysis results  |
| B.glabrata annotation - RNA-Seq I   | RNA-Seq data provided by WashU (fastq f/r)                               |
| B.glabrata annotation - RNA-Seq II  | RNA-Seq data provided by WashU (aligned BAM and assembled fasta)         |
| B.glabrata annotation - RNA-Seq III | sorted BAM files of RNA-Seq suitable for upload to IGV (from G.Oliveira) |
| B.glabrata - assemblies             | assemblies for read dependency   |
| B.glabrata RNA-Seq (clean)          | RNA-Seq without adaptors, quality clipped                                |
| B.glabrata RNA-Seq (paired)         | fastq custom upload N.Dhelly 16/01/13                                    |

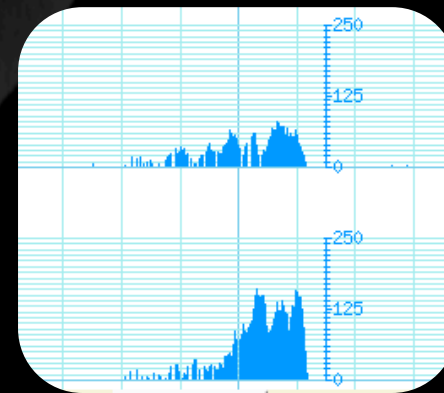


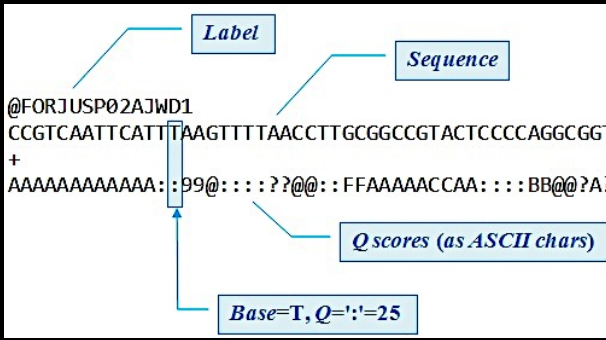
+ Align Reads

- + Format Conversion
- + Sorting



+ Find Enriched Regions





- + Format Conversion
- + Trimming depending on QC



+ Align Reads

- + Format Conversion
- + Sorting



+ Find Enriched Regions



+ FastQC

