

# Development of a workflow for SNPs detection in grapevine species: *MAPHiTS*

**-> *Integration of MAPHiTS in Galaxy***



## A. Galaxy Presentation



# A.1. Galaxy Homepage



The screenshot shows the Galaxy web interface. At the top, there is a navigation bar with the 'Galaxy' logo and menu items: 'Analyze Data', 'Workflow', 'Shared Data', 'Help', and 'User'. On the left side, there is a 'Tools' panel with a list of various bioinformatics tools such as 'Get Data', 'Send Data', 'ENCODE Tools', 'Lift-Over', 'Text Manipulation', 'Filter and Sort', 'Unix Tools', 'Join, Subtract and Group', 'Convert Formats', 'Extract Features', 'Fetch Sequences', 'Fetch Alignments', 'Get Genomic Scores', 'Operate on Genomic Intervals', 'Statistics', 'Wavelet Analysis', 'Graph/Display Data', 'Regional Variation', 'Multiple regression', 'Multivariate Analysis', 'Evolution', 'Metagenomic analyses', 'FASTA manipulation', 'NGS: QC and manipulation', 'NGS: Indel Analysis', 'NGS: SAM Tools', 'FastX Toolkit', 'MAPHITS', 'S-MART', and 'Workflows'. The main content area features a central banner for 'Unité Recherche Génomique Info' with the URGI logo and several DNA double helix illustrations. Below the banner, there is a text block: 'The Galaxy team is a part of BX at Penn State. This project is supported in part by NSF, NHGRI, and the Huck Institutes of the Life Sciences.' To the right of the banner is the INRA logo. On the far right, there is a 'History' panel showing 'Unnamed history' and a message: 'Your history is empty. Click "Get Data" on the left pane to start'. At the bottom of the main content area, the URL <http://urgi.versailles.inra.fr/galaxy/> is displayed in blue text.

## A.2. Installation of URGI Galaxy

**Galaxy is installed on URGI cluster with:**

- CPU: **704** (Intel Xeon)
- RAM max: **96 Gb** per job
- Storage: **60 Tb**



Using Sun Grid Engine (for job management) and a PostgreSQL Database (for Galaxy).

# A.3. Homepage presentation

Return to homepage

Tools

Management of Galaxy

History



The screenshot shows the Galaxy web interface. At the top left, the 'Galaxy' logo is highlighted with a red box, with a red arrow pointing to the text 'Return to homepage'. Below it is a blue-bordered 'Tools' sidebar. At the top center, a dark navigation bar contains 'Analyze Data', 'Workflow', 'Data Libraries', 'Help', and 'User', with an orange box around it and an orange arrow pointing to 'Management of Galaxy'. On the right, a 'History' sidebar is highlighted with a green border and a green arrow pointing to 'History'. The main content area features a banner for 'Unité Recherche Génomique Info' with the URL <http://urgj.versailles.inra.fr/> and the URGI logo. The background of the banner is decorated with DNA double helix illustrations.

# A.4. How to upload your data ?

a. Tools → Get Data → Upload File

**Tools**

- Get Data
  - Upload File from your computer

**File Format:**  
 Auto-detect (dropdown menu)  
 Which format? See help below

**File:**  
 Parcourir... (button)

**URL/Text:**  
 (text area)

Here you may specify a list of URLs (one per line) or paste the contents of a file.

**Convert spaces to tabs:**  
 Yes  
 Use this option if you are entering intervals by hand.

**Genome:**  
 Click to Search or Select Build

Execute (button)      Help

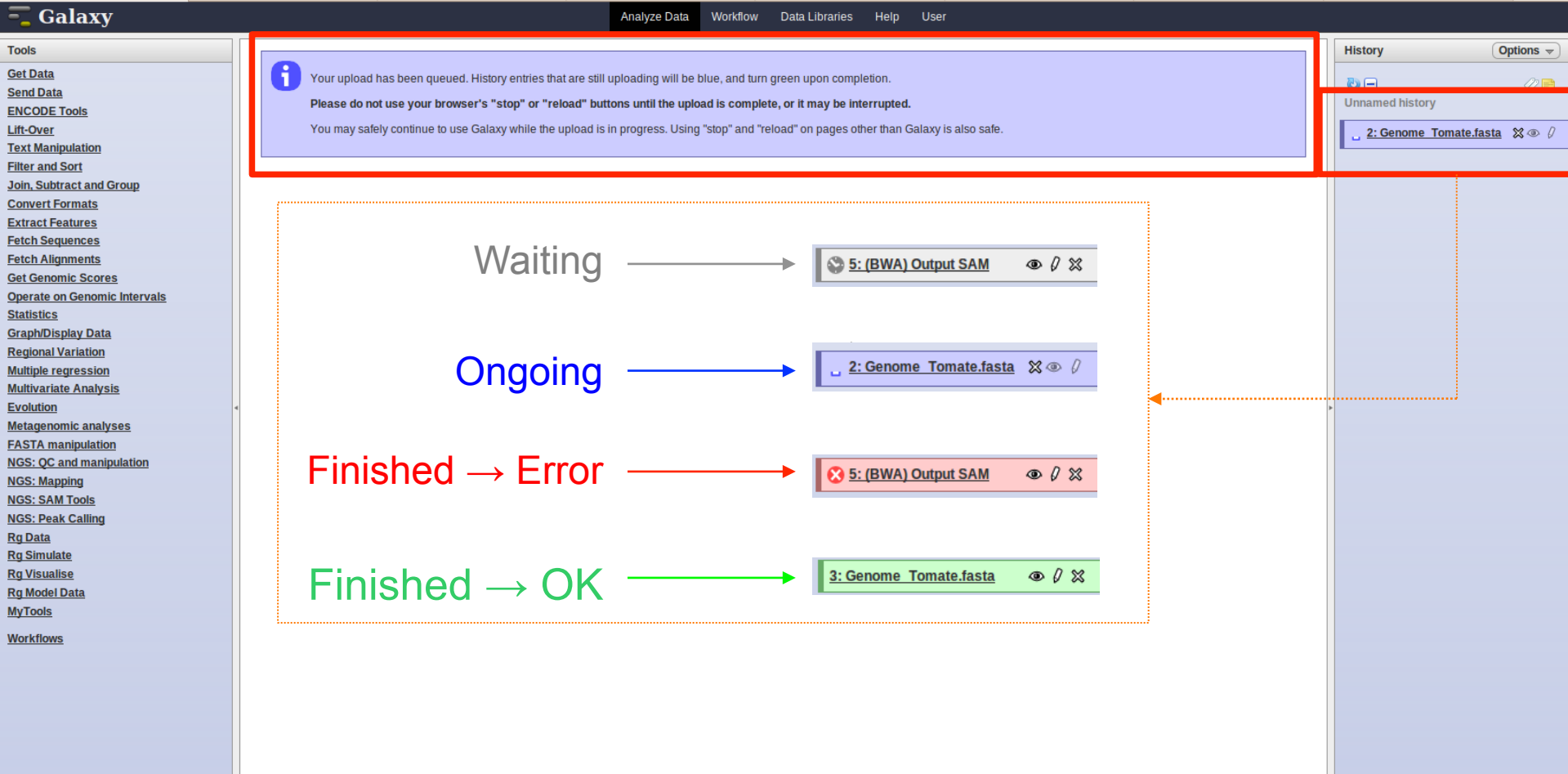
**Auto-detect**  
 The system will attempt to detect Axt, Fasta, Fastqsolexa, Gff, Gff3, Html, Lav, Maf, Tabular, Wiggle, Bed and Interval (Bed with headers) formats. If your file is not detected properly as one of the known formats, it most likely means that it has some format problems (e.g., different number of columns on different rows). You can still coerce the system to set your data to the format you think it should be. You can also upload compressed files, which will automatically be decompressed.

**Ab1**  
 A binary sequence file in 'ab1' format with a '.ab1' file extension. You must manually select this 'File Format' when uploading the file.

**Axt**  
 blastz pairwise alignment format. Each alignment block in an axt file contains three lines: a summary line and 2 sequence lines. Blocks are separated from one another by blank lines. The summary line contains chromosomal position and size information about the alignment. It consists of 9 required fields.

**Bam**  
 A binary file compressed in the BGZF format with a '.bam' file extension.

# A.4. How to upload your data ?



The screenshot shows the Galaxy web interface with a navigation menu on the left and a main workspace. A red box highlights an information message at the top: "Your upload has been queued. History entries that are still uploading will be blue, and turn green upon completion. Please do not use your browser's 'stop' or 'reload' buttons until the upload is complete, or it may be interrupted. You may safely continue to use Galaxy while the upload is in progress. Using 'stop' and 'reload' on pages other than Galaxy is also safe." A legend in the center explains the status of history items: "Waiting" (grey), "Ongoing" (blue), "Finished → Error" (red), and "Finished → OK" (green). The right sidebar shows a history list with "2: Genome\_Tomate.fasta" highlighted in blue, corresponding to the "Ongoing" status in the legend.

**Galaxy** Analyze Data Workflow Data Libraries Help User

**Tools**

- Get Data
- Send Data
- ENCODE Tools
- Lift-Over
- Text Manipulation
- Filter and Sort
- Join, Subtract and Group
- Convert Formats
- Extract Features
- Fetch Sequences
- Fetch Alignments
- Get Genomic Scores
- Operate on Genomic Intervals
- Statistics
- Graph/Display Data
- Regional Variation
- Multiple regression
- Multivariate Analysis
- Evolution
- Metagenomic analyses
- FASTA manipulation
- NGS: QC and manipulation
- NGS: Mapping
- NGS: SAM Tools
- NGS: Peak Calling
- Rg Data
- Rg Simulate
- Rg Visualise
- Rg Model Data
- MyTools
- Workflows

**History** Options

Unnamed history

- 2: Genome\_Tomate.fasta

**Waiting** → 5: (BWA) Output SAM

**Ongoing** → 2: Genome\_Tomate.fasta

**Finished → Error** → 5: (BWA) Output SAM

**Finished → OK** → 3: Genome\_Tomate.fasta

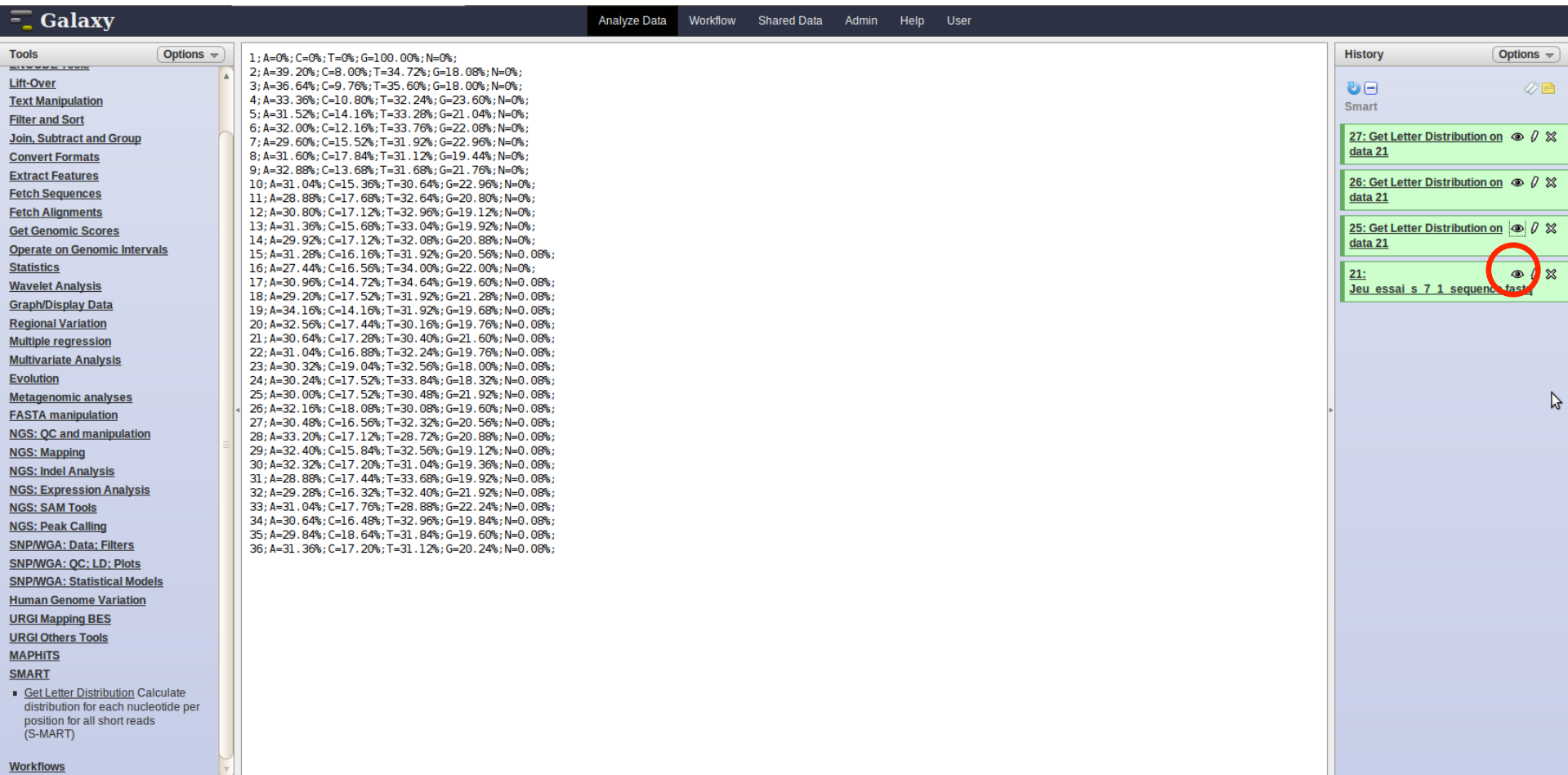
# A.5. How to use a tool ?

a. Choose a tool from the list

The screenshot shows the Galaxy web interface for the 'Map with Bowtie for Illumina' tool. The left-hand navigation menu is visible, with 'Map with Bowtie for Illumina' highlighted in a red box. The main configuration area contains several sections: 'Will you select a reference genome from your history or use a built-in index?', 'Select a reference genome:', 'Is this library mate-paired?', 'FASTQ file:', 'Bowtie settings to use:', and 'Suppress the header in the output SAM file:'. A green bracket and arrow point to these configuration options, labeled 'b. Set options and parameters'. An orange arrow points to the 'Execute' button, labeled 'c. Execute Help'. The right-hand history panel shows 'Unnamed history' and a message: 'Your history is empty. Click 'Get Data' on the left pane to start'.



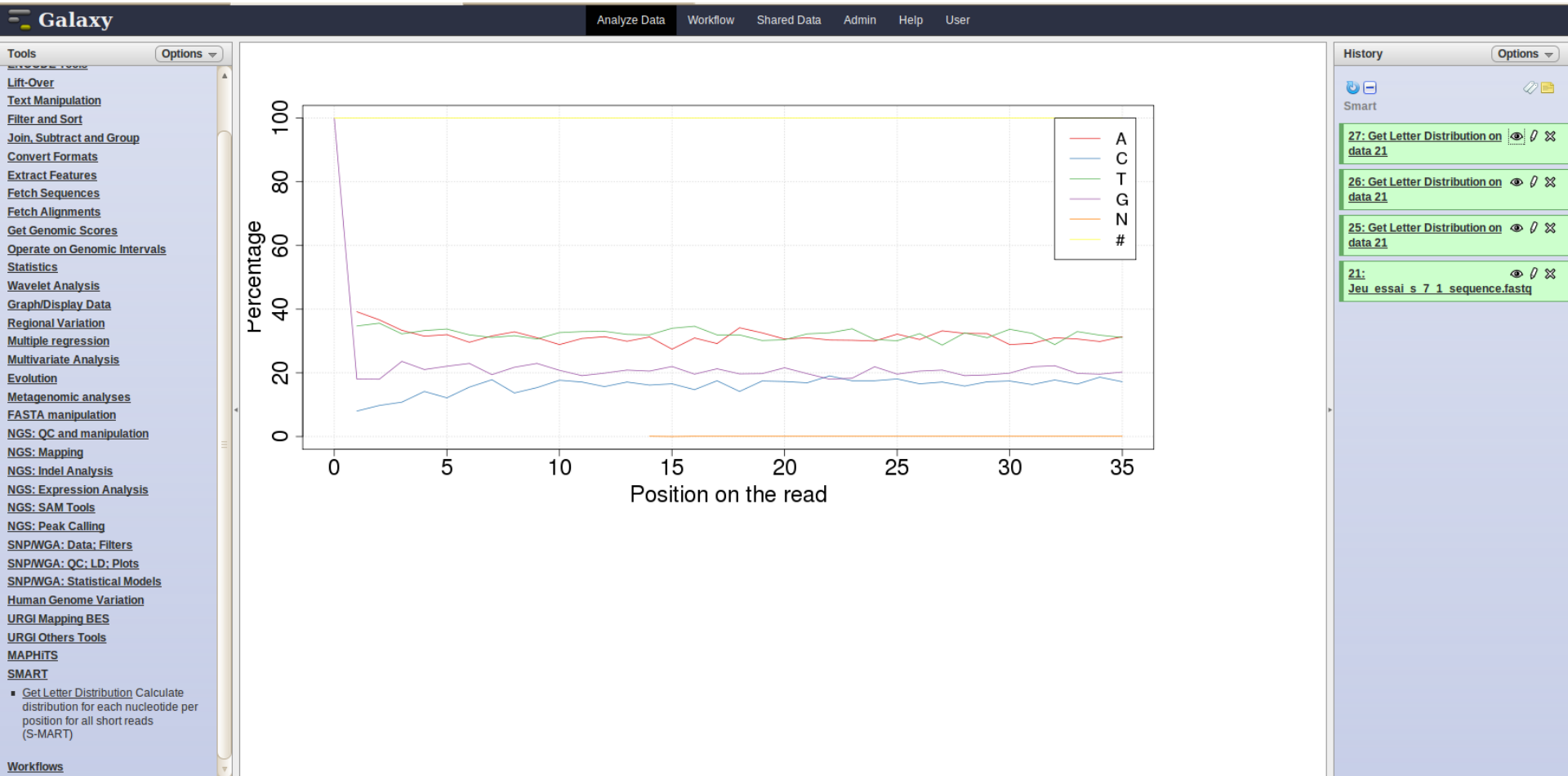
# A.6. Example



The screenshot displays the Galaxy web interface. On the left is a 'Tools' sidebar with a search bar and a list of tool categories including 'Text Manipulation', 'Filter and Sort', 'Join, Subtract and Group', 'Convert Formats', 'Extract Features', 'Fetch Sequences', 'Fetch Alignments', 'Get Genomic Scores', 'Operate on Genomic Intervals', 'Statistics', 'Wavelet Analysis', 'Graph/Display Data', 'Regional Variation', 'Multiple regression', 'Multivariate Analysis', 'Evolution', 'Metagenomic analyses', 'FASTA manipulation', 'NGS: QC and manipulation', 'NGS: Mapping', 'NGS: Indel Analysis', 'NGS: Expression Analysis', 'NGS: SAM Tools', 'NGS: Peak Calling', 'SNP/WGA: Data: Filters', 'SNP/WGA: QC: LD: Plots', 'SNP/WGA: Statistical Models', 'Human Genome Variation', 'URGI Mapping BES', 'URGI Others Tools', 'MAPHITS', and 'SMART'. The 'SMART' tool is expanded, showing a sub-item 'Get Letter Distribution Calculate distribution for each nucleotide per position for all short reads (S-MART)'. The main panel shows a workflow with 36 steps, each displaying nucleotide distribution percentages for A, C, G, T, and N. The right sidebar shows a 'History' panel with a list of workflow instances, including '27: Get Letter Distribution on data 21', '26: Get Letter Distribution on data 21', '25: Get Letter Distribution on data 21', and '21: Jeu essai s 7 1 sequenc fastq'. The '21: Jeu essai s 7 1 sequenc fastq' entry has a red circle around its eye icon, indicating it is the current view.

## SMART - Get Letter Distribution

# A.6. Example



## SMART - Get Letter Distribution

## B. MAPHiTS Presentation



# B.1. Background and objectives of the pipeline

## ▪ Objectives:

Detect a set of SNPs between various species of Grape after mapping short reads against a reference genome.

## ▪ Data:

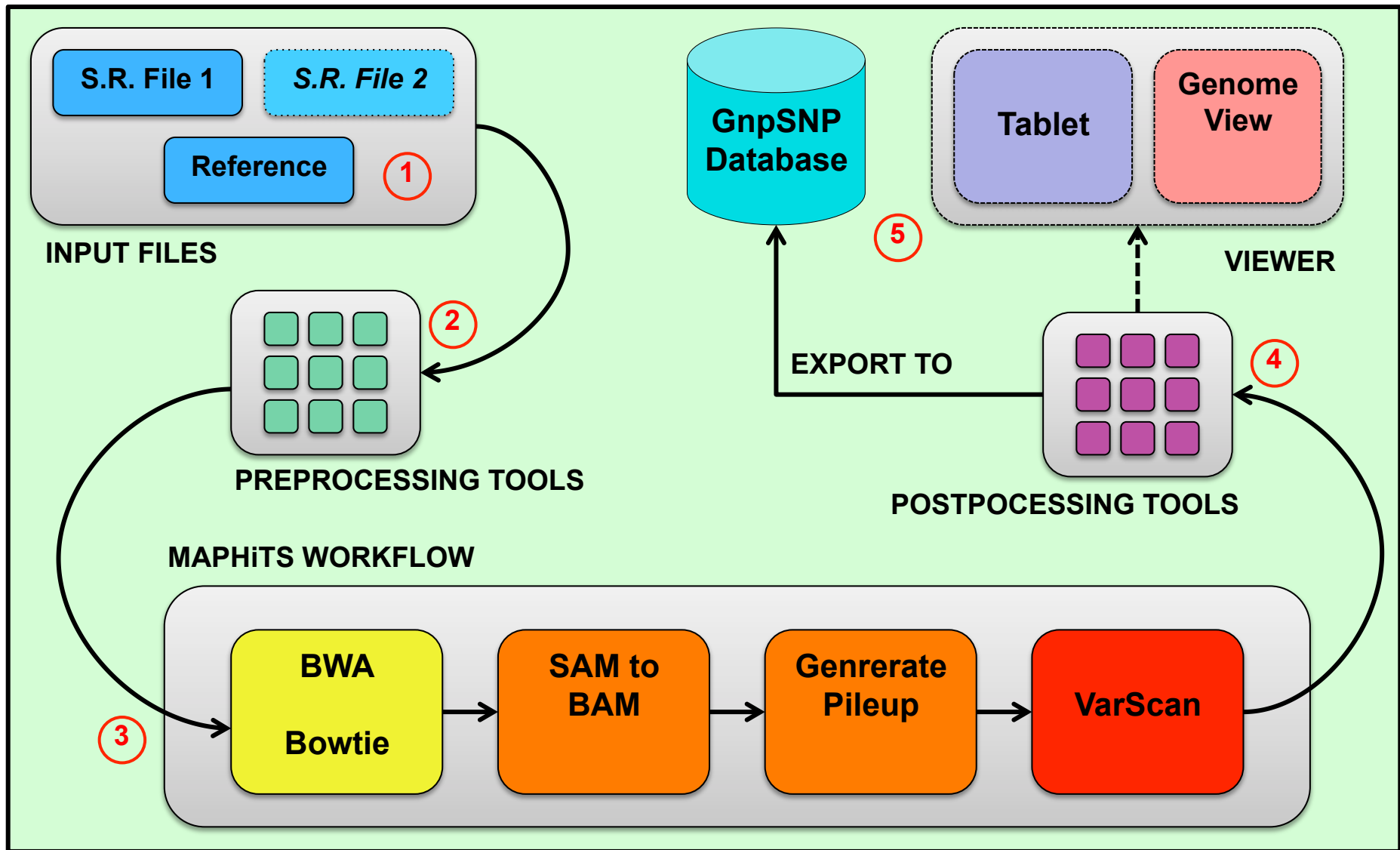
- Muscares : 6 genotypes

- GrapeReSeq : 16 genotypes

Short reads are in paired-ends with 76, 101 or 114 bp (*Illumina GAII*).

➔ Other projects are also in progress with others species.

# B.2. MAPHiTS: Resume



## B.3. MAPHiTS Development Tools

- **Optimization tools:**
  - BWA in parallel
  - SAM-to-BAM in parallel

**Time Saving: 10x average !**

Exemple:

- Before: 11 -12 hours
- Now: 1 - 2 hours

**Fewer Ressources Required !**

Exemple:

- Before: 1 big job in 1 cluster node
- Now: N little jobs in N cluster nodes

## B.3. MAPHiTS Development Tools

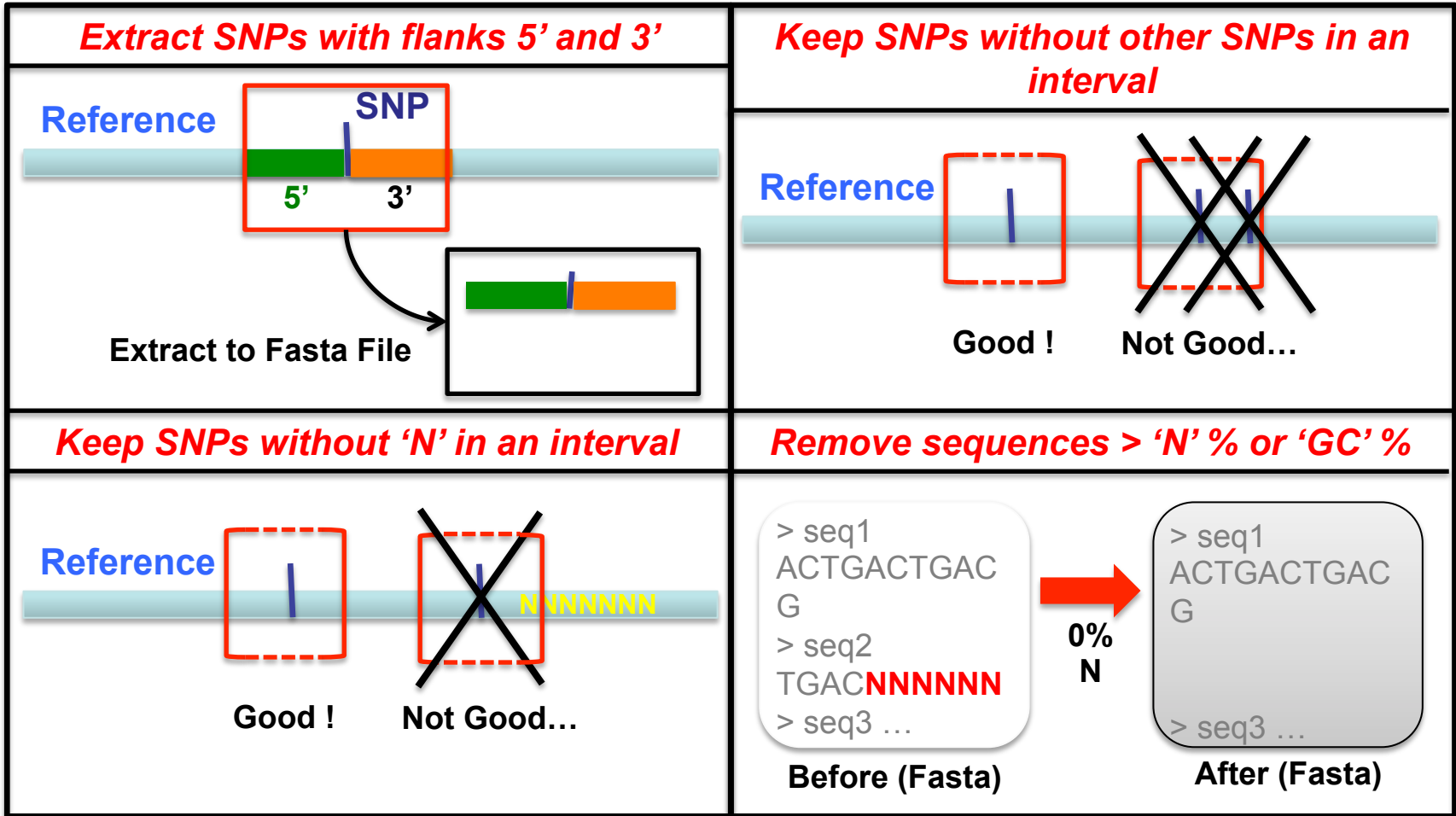
### ■ Preprocessing tools:

- Remove duplicated short-reads
- Remove short reads not in paired-ends
- Remove short reads > 'N'%
- Remove informations in each FASTA file header

### ■ Postprocessing tools:

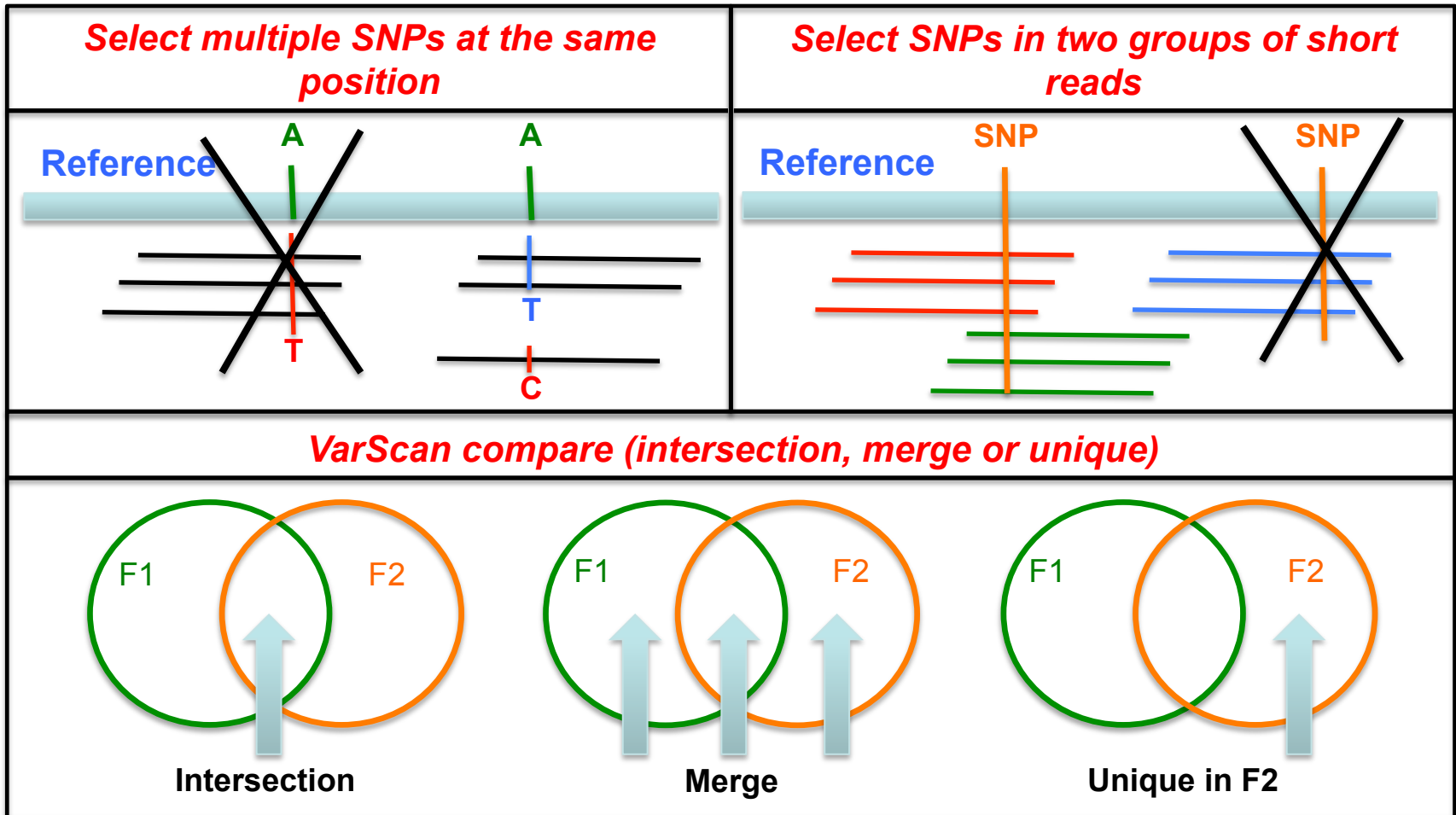
- Count multiple hits from the results of BWA
- Extract short reads from SAM file
- VarScan compare (multiple files)
- VarScan filter
- VarScan to Gff3

# B.3. MAPHiTS Development Tools





# B.3. MAPHiTS Development Tools



# B.3. MAPHiTS Development Tools

**Get SNPs distribution by chromosome**

VarScan File

Graph chr1

Graph chr2

Graph chr3

**Give consensus nucleotides (with different VarScan analysis)**

VarScan Files

HTML File

INRA - URGI Team

MAPHiTS Genotype Card									
Ref	Position	RefAllele	Trayshed	Vitis-Rotundifolia	Regale_run2	Regale_run1	Fry	Carlos	Consensus
chr18	801	G	T	T	G	G	T	T	NA
chr18	807	C	T	T	C	C	T	T	NA
chr18	977	A	A	A	A	A	A	T	NA
chr18	1060	T	T	T	T	T	T	A	NA
chr18	1751	A	A	A	A	A	A	G	NA
chr18	1752	A	A	A	A	A	A	T	NA
chr18	1786	A	A	A	A	A	A	G	NA
chr18	1890	T	T	T	T	T	T	G	NA
chr18	2586	T	T	T	T	T	T	G	NA
chr18	2591	T	G	G	T	T	T	G	NA
chr18	4499	A	G	G	G	G	G	G	G 66
chr18	4534	T	A	A	A	T	A	A	NA
chr18	4585	G	C	G	C	G	C	C	NA
chr18	4612	C	T	C	T	T	T	C	NA
chr18	4618	T	C	T	C	C	C	C	NA
chr18	4665	A	G	G	G	G	G	G	G 66
chr18	4754	T	C	C	C	C	C	C	C 66
chr18	4785	C	A	C	A	A	A	A	NA
chr18	4813	A	G	G	G	G	G	G	G 66
chr18	4824	A	T	T	T	T	T	T	T 66
chr18	4910	A	G	G	G	A	A	G	NA

# B.3. MAPHiTS Development Tools

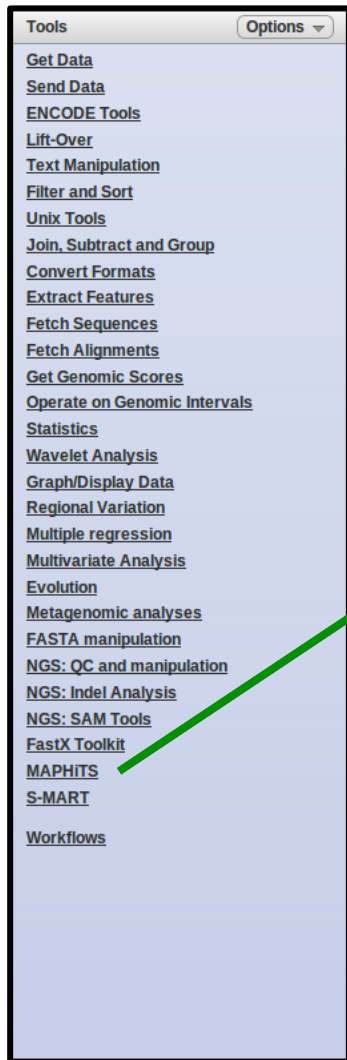
MAPHiTS Genotype Card

Ref	Position	RefAllele	Trayshed	Vitis-Rotundifolia	Regale_run2	Regale_run1	Fry	Carlos	Consensus
chr18	801	G	T	T	G	G	T	T	NA
chr18	807	C	T	T	C	C	T	T	NA
chr18	977	A	A	A	A	A	A	T	NA
chr18	1060	T	T	T	T	T	T	A	NA
chr18	1751	A	A	A	A	A	A	G	NA
chr18	1752	A	A	A	A	A	A	T	NA
chr18	1786	A	A	A	A	A	A	G	NA
chr18	1890	T	T	T	T	T	T	G	NA
chr18	2586	T	T	T	T	T	T	G	NA
chr18	2591	T	G	G	T	T	T	G	NA
chr18	4499	A	G	G	G	G	G	G	G 6/6
chr18	4534	T	A	A	A	T	A	A	NA
chr18	4585	G	C	G	C	G	C	C	NA
chr18	4612	C	T	C	T	T	T	C	NA
chr18	4618	T	C	T	C	C	C	C	NA
chr18	4665	A	G	G	G	G	G	G	G 6/6
chr18	4754	T	C	C	C	C	C	C	C 6/6
chr18	4785	C	A	C	A	A	A	A	NA
chr18	4813	A	G	G	G	G	G	G	G 6/6
chr18	4824	A	T	T	T	T	T	T	T 6/6
chr18	4910	A	G	G	G	A	A	G	NA

## C. MAPHiTS in Galaxy



# C.2. New URGI Integrated tools



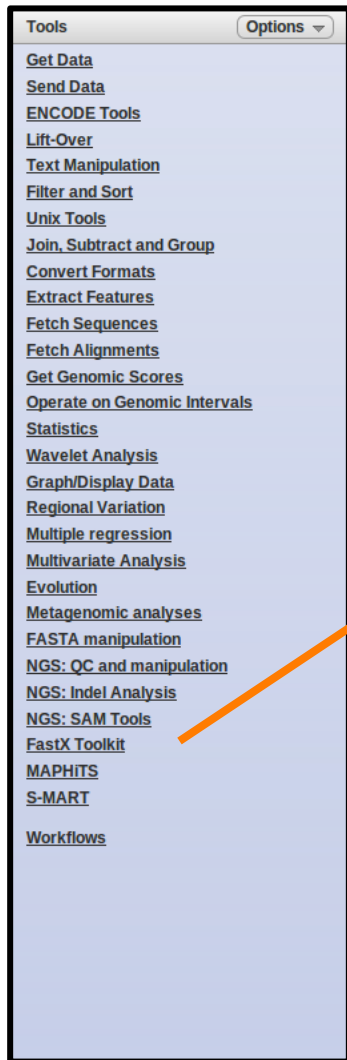
## MAPHiTS

### MAPHiTS

#### PREPROCESS TOOLS

- Header fasta filter Remove all informations in each header of fasta file.
- Remove duplicate short reads
- Remove duplicate short reads for big files (> 2Go)
- Remove short reads not in paired-ends
- Remove short reads not in paired-ends for big files (>2Go)
- Remove short reads > N %
- Remove short reads > N % for big files (>2Go)

## C.2. New Integrated tools



**FASTX-Toolkit**

FastX Toolkit

TOOLS

- Barcode Splitter
- Clip adapter sequences
- Collapse sequences
- Compute quality statistics
- FASTA Width formatter
- FASTQ to FASTA converter
- Filter by quality
- Mask nucleotides (based on quality)
- Quality format converter (ASCII-Numeric)
- Remove sequencing artifacts

# C.2. New Others URGI

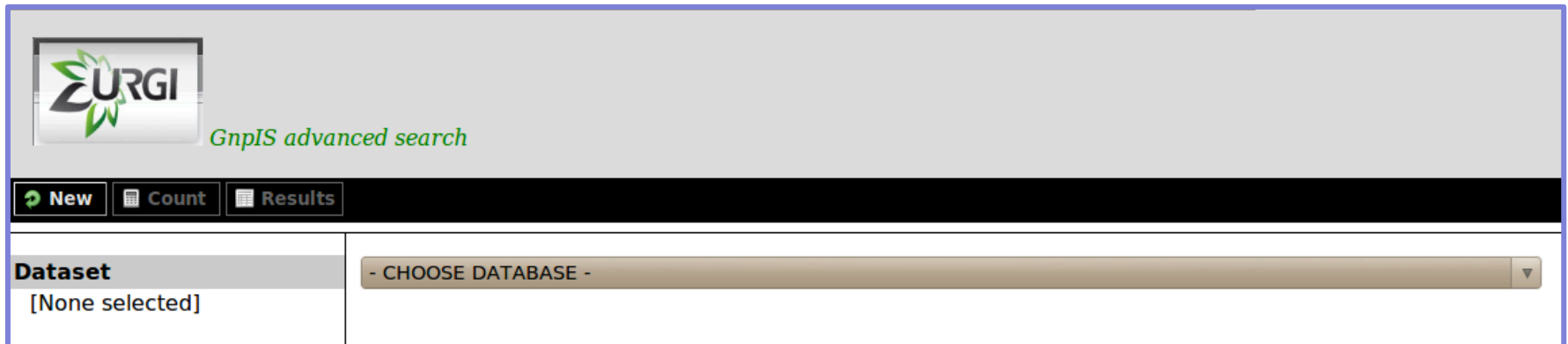
## Integrated tools

Access to URGI  
Information System  
via **BioMart** software

**Get Data:**

- [Upload File](#) from your computer
- [UCSC Main](#) table browser
- [UCSC Test](#) table browser
- [UCSC Archaea](#) table browser
- [BX main](#) browser
- [Get Microbial Data](#)
- [BioMart](#) Central server
- [BioMart INRA URGI GnpIs](#)
- [CBI Rice Mart](#) rice mart
- [GrameneMart](#) Central server

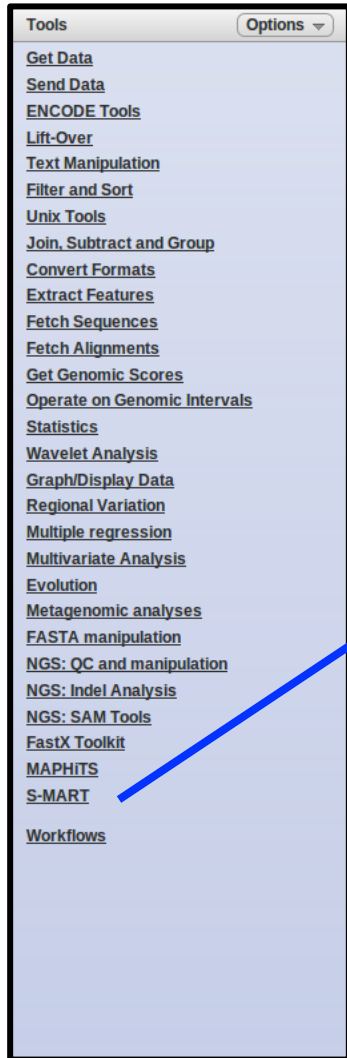
**BioMart  
URGI  
GnpIs**



The screenshot shows the BioMart interface. At the top left is the URGI logo and the text "GnpIS advanced search". Below this is a navigation bar with buttons for "New", "Count", and "Results". The main area features a "Dataset" section with "[None selected]" and a dropdown menu labeled "- CHOOSE DATABASE -". A blue arrow points from the "BioMart INRA URGI GnpIs" option in the "Get Data" menu to the BioMart interface.

# C.2. New Others URGI

## Integrated tools



S-MART

### S-MART

#### FILES CONVERTER

- Bed -> Csv Convert Bed File to Csv File.
- Bed -> Gff2 Convert Bed File to Gff2 File.
- Bed -> Gff3 Convert Bed File to Gff3 File.
- Bed -> Sam Convert Bed File to Sam File.
- Blast (-m 8) -> Csv Convert Blast (-m 8) File to Csv File.
- Blast (-m 8) -> Gff2 Convert Blast (-m 8) File to Gff2 File.



## C.2. New Others URGI

# Integrated tools

What can I do with all this RNA-Seq data ?

S-MART:

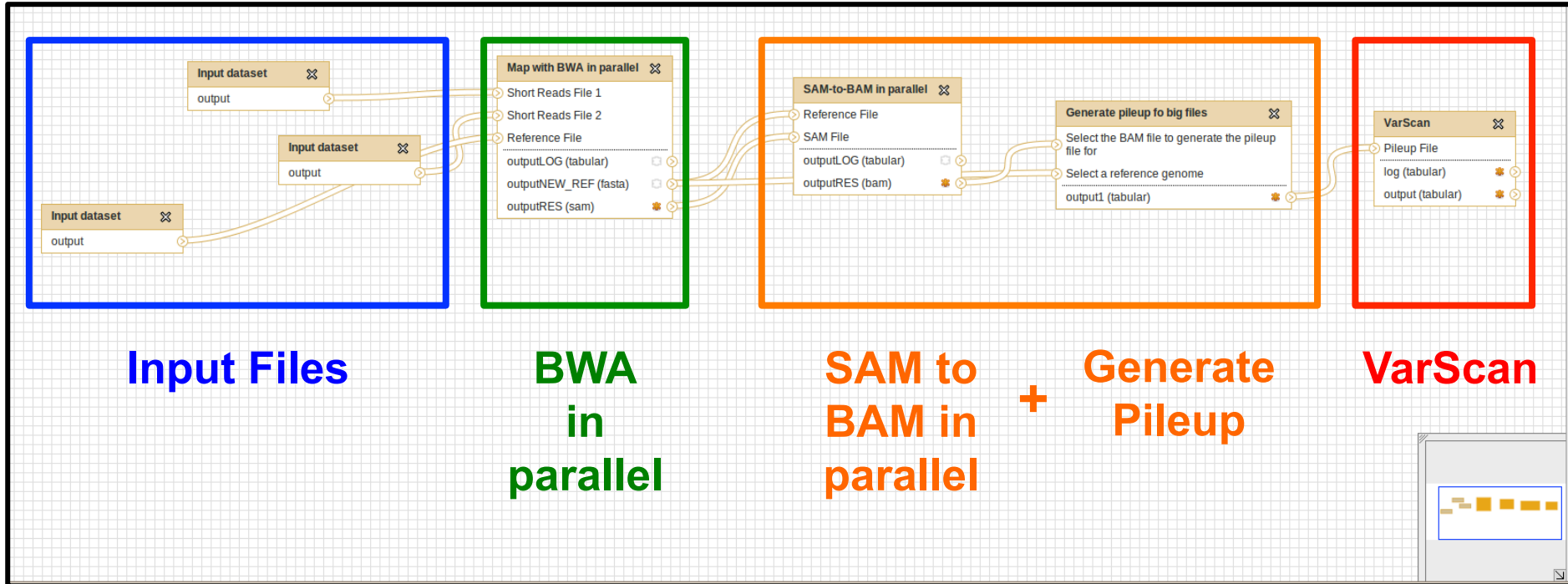
- is a set of independant tools
- works on a standard PC and with Galaxy
- can be installed and used easily

Use S-MART for data manipulation, data visualization, differential expression, ...

**Link:** <http://urgi.versailles.inra.fr/Tools/S-MART>

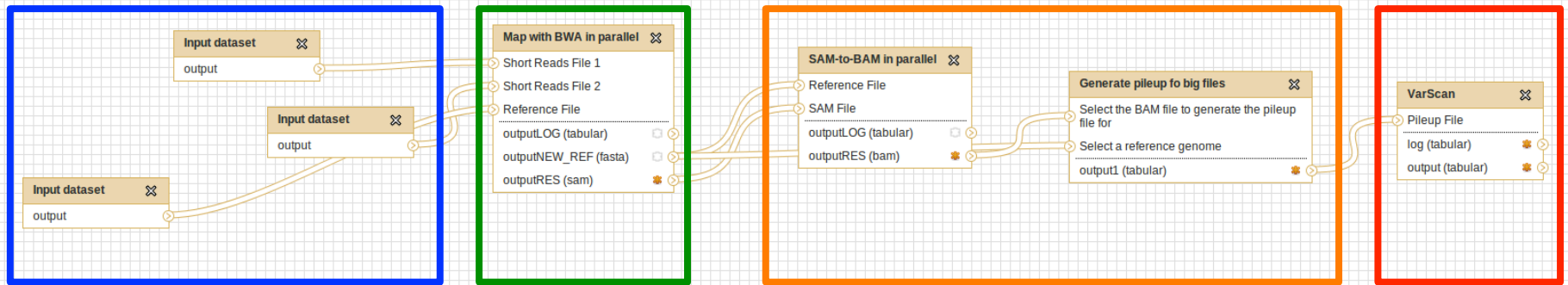
**Contact:** [matthias.zytnicki@versailles.inra.fr](mailto:matthias.zytnicki@versailles.inra.fr)

# C.3. MAPHiTS: Build



MAPHiTS is build using the graphical interface of Galaxy.

# C.3. MAPHiTS: Build



**Input Files**

**BWA  
in  
parallel**

**SAM to  
BAM in  
parallel**

+

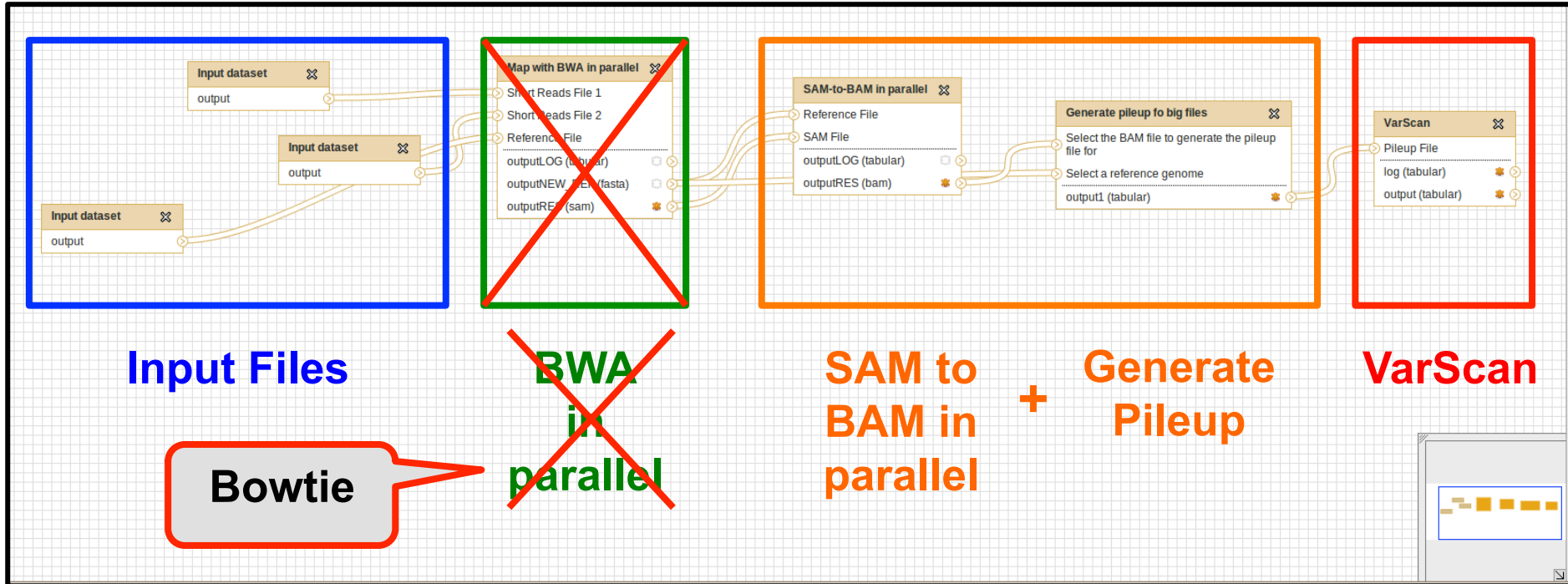
**Generate  
Pileup**

**VarScan**

If you want  
you can add  
some  
preprocessing  
tools ...

... or some  
postprocessing tools

# C.3. MAPHiTS: Build



You can remove one tool and replace it by an other tool very quickly.

# C.4. MAPHiTS: Launch

Running workflow "MAPHiTS Parallel (paired)"

**Step 1: Input dataset**

Reference File (.fasta) **STEP 1**

**Step 2: Input dataset**

Short Reads File 1 (.fastq) **STEP 2**

**Step 3: Input dataset**

Short Reads File 2 (.fastq) **STEP 3**

**Step 4: Map with BWA in parallel** **STEP 4**

Type of Short Reads  
Paired-ends

Short Reads File 1  
Output dataset 'output' from step 2

Short Reads File 2  
Output dataset 'output' from step 3

Reference File  
Output dataset 'output' from step 1

Use default parameters for Bwa  
No

Maximum edit distance if the value is INT, or the fraction of missing alignments given 2% uniform base error rate if FLOAT. In the latter case, the maximum edit distance is automatically chosen for different read lengths. (-n)  
 ← Parameter

Maximum number of gap opens (-o)  
1

Maximum number of gap extensions (-e)  
-1

Disallow long deletion within [value] bp towards the 3'-end (-d)  
16

# C.4. MAPHiTS: Launch



The screenshot shows the Galaxy web interface. At the top, there is a navigation bar with 'Galaxy' and tabs for 'Analyze Data', 'Workflow', 'Shared Data', 'Help', and 'User'. On the left, a 'Tools' sidebar lists various analysis tools, with 'MAPHiTS' highlighted. The main workspace displays a green message box indicating a successful workflow run. The message lists the datasets added to the queue: 1: Genome.fasta, 2: SR\_1.fastq, 3: SR\_2.fastq, 4: [HeaderFastaFilter] Output Fasta File, 5: [MAPHiTS] SAM file, 6: [MAPHiTS] BAM file, 7: [MAPHiTS] PILEUP file, 8: [MAPHiTS] RESUME file, and 9: [MAPHiTS] VARSCAN file. On the right, a 'History' panel shows a list of recent jobs, with the first three jobs (1: Genome.fasta, 2: SR\_1.fastq, 3: SR\_2.fastq) highlighted in green, corresponding to the datasets listed in the message box. A red callout box points to the message box with the text: 'When you run the workflow, this message appears !'.

Galaxy

Analyze Data Workflow Shared Data Help User

Tools Options

- Convert Formats
- Extract Features
- Fetch Sequences
- Fetch Alignments
- Get Genomic Scores
- Operate on Genomic Intervals
- Statistics
- Wavelet Analysis
- Graph/Display Data
- Regional Variation
- Multiple regression
- Multivariate Analysis
- Evolution
- Metagenomic analyses
- FASTA manipulation
- NGS: QC and manipulation
- NGS: Indel Analysis
- NGS: SAM Tools
- FastX Toolkit
- MAPHiTS
- S-MART
- Workflows
  - Trim And Compare ALL Short Reads (paired)
  - MAPHiTS Not Parallel (single)
  - MAPHiTS Not Parallel (paired)
  - MAPHiTS Parallel (single)
  - MAPHiTS Parallel (paired)
  - Trim And Compare EPGV Short Reads (paired)
  - All workflows

Options

Successfully ran workflow "MAPHiTS Not Parallel (paired)", the following datasets have been added to the queue.

- 1: Genome.fasta
- 2: SR\_1.fastq
- 3: SR\_2.fastq
- 4: [HeaderFastaFilter] Output Fasta File
- 5: [MAPHiTS] SAM file
- 6: [MAPHiTS] BAM file
- 7: [MAPHiTS] PILEUP file
- 8: [MAPHiTS] RESUME file
- 9: [MAPHiTS] VARSCAN file

History Options

Workshop 6

- 9: [MAPHiTS] VARSCAN file
- 8: [MAPHiTS] RESUME file
- 7: [MAPHiTS] PILEUP file
- 6: [MAPHiTS] BAM file
- 5: [MAPHiTS] SAM file
- 4: [HeaderFastaFilter] Output Fasta File
- 3: SR\_2.fastq
- 2: SR\_1.fastq
- 1: Genome.fasta

When you run the workflow, this message appears !

# C.4. MAPHiTS: Launch

**Galaxy** Analyze Data Workflow Shared Data Help User

**Tools** Options ▾

- [Convert Formats](#)
- [Extract Features](#)
- [Fetch Sequences](#)
- [Fetch Alignments](#)
- [Get Genomic Scores](#)
- [Operate on Genomic Intervals](#)
- [Statistics](#)
- [Wavelet Analysis](#)
- [Graph/Display Data](#)
- [Regional Variation](#)
- [Multiple regression](#)
- [Multivariate Analysis](#)
- [Evolution](#)
- [Metagenomic analyses](#)
- [FASTA manipulation](#)
- [NGS: QC and manipulation](#)
- [NGS: Indel Analysis](#)
- [NGS: SAM Tools](#)
- [FastX Toolkit](#)
- [MAPHiTS](#)
- [S-MART](#)
- Workflows**
- [Trim And Compare ALL Short Reads \(paired\)](#)
- [MAPHiTS Not Parallel \(single\)](#)
- [MAPHiTS Not Parallel \(paired\)](#)
- [MAPHiTS Parallel \(single\)](#)
- [MAPHiTS Parallel \(paired\)](#)
- [Trim And Compare EPGV Short Reads \(paired\)](#)
- [All workflows](#)

✓ Successfully ran workflow "MAPHiTS Not Parallel (paired)", the following datasets have been added to the queue.

- 1: Genome.fasta
- 2: SR\_1.fastq
- 3: SR\_2.fastq
- 4: [HeaderFastaFilter] Output Fasta File
- 5: [MAPHiTS] SAM file
- 6: [MAPHiTS] BAM file
- 7: [MAPHiTS] PILEUP file
- 8: [MAPHiTS] RESUME file
- 9: [MAPHiTS] VARSCAN file

**History** Options ▾

Workshop 6

- 9: [MAPHiTS] VARSCAN file
- 8: [MAPHiTS] RESUME file
- 7: [MAPHiTS] PILEUP file
- 6: [MAPHiTS] BAM file
- 5: [MAPHiTS] SAM file
- 4: [HeaderFastaFilter] Output Fasta File
- 3: SR 2.fastq
- 2: SR 1.fastq
- 1: Genome.fasta

**VarScan results**


**Generate Pileup  
SAM-to-BAM**

**BWA**

**PreProcessing tool**

**Input files**

## C.5. Shared Workflows

Published Workflows	
<input type="text" value="search"/>    <a href="#">Advanced Search</a>	
Name	Annotation
<a href="#">MAPHiTS Parallel (paired)</a>	Workflow of SNPs detection, in parallel, for paired-end short reads.
<a href="#">Trim And Compare EPGV Short Reads (paired)</a>	
<a href="#">Trim And Compare ALL Short Reads (paired)</a>	This workflow can filter your short reads (remove short reads with 'N' and short reads not in paired-ends) and generates graphs before and after this...

Some workflows are available for logged users in ‘*Shared Data*’ and ‘*Published Workflows*’ section.



## C.6. Shared your History

If a user wants to share its results with other users or a specific user, it's possible !



The screenshot shows a web interface titled "Published Histories". It features a search bar with the placeholder text "search" and a magnifying glass icon, followed by a link to "Advanced Search". Below the search bar is a table with two columns: "Name" and "Annotation". The table contains two entries, both with underlined text:


Name	Annotation
<u>VarScan compare Muscares</u>	
<u>VarScan compare Muscares v2</u>	

All this histories are in *'Shared Data'* and *'Published Histories'*.

## C.7. Shared Data


- In 'Shared Data' and 'Data Libraries' section, logged users can see [1 directory per Project.](#)
- Users [can only see their](#) projects.

### Data Libraries


  | [Advanced Search](#)

<u>Name</u> ↓
<u>grapereseq</u>
<u>magictomsnps</u>
<u>muscares</u>
<u>poplar</u>

## C.7. Shared Data



**Data Library "grapereseq"**

Name	
<input type="checkbox"/>	 <u>short reads</u> ▼
<input type="checkbox"/>	<u>Vvinifera_v5.1_chr_05Jan2010.fasta</u> ▼

For selected items:  ▼

**All short reads** (callout for 'short reads')

**Reference Genome** (callout for 'Vvinifera\_v5.1\_chr\_05Jan2010.fasta')

They can import their data into the history quickly.

→ Usefull for **NGS** !

## D. Perspectives



## D. Perspectives

- **Add new tools** (all tools used in all our pipelines)
- **Link Galaxy to a visualization software** (Gbrowse 2, Tablet, GenomeView, ...)
- **Application Note in progress (2011)**

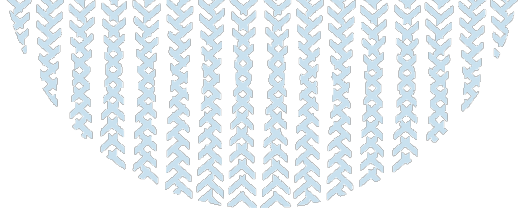
# Acknowledgements



- **EPGV Team**
- **URGI Team**



- **Galaxy developers**
- **Galaxy community**



# Thank you for your attention !!!

