

FAIR_bioinfo : Open Science and FAIR principles in a bioinformatics project

How to make a bioinformatics project more reproducible

C. Hernandez¹ T. Denecker² J. Sellier² G. Le Corguillé²
C. Toffano-Nioche¹

¹Institute for Integrative Biology of the Cell (I2BC)
UMR 9198, Université Paris-Sud, CNRS, CEA
91190 - Gif-sur-Yvette, France

²IFB Core Cluster taskforce

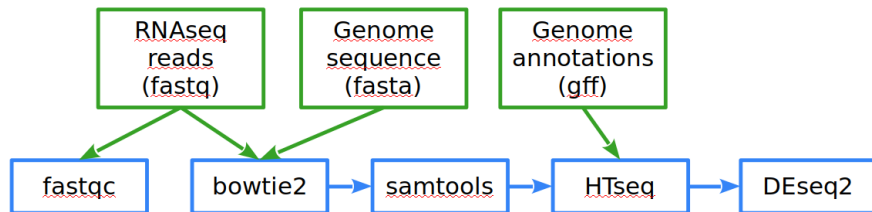
June 2021



From a bash script to a complete snakefile

A classical RNAseq analysis

Analysis workflow



green=input, blue=tool

fastqc control quality of the input reads

bowtie2 reads mapping on the genome sequence

samtools mapped reads selection & formatting

HTseq count table of mapped reads on genes (annotations)

DEseq2 statistical analysis: genes list having differential expression

Data, a bash script and its command line

Data

-g genome sequence acces (including extention .fna, .fasta)

-a genome annotation acces (inluding extention .gff)

-d RNAseq sample prefix

next args: RNAseq sample prefix, no .fastq.gz extention

Bash command line

```
1 FAIR_initial_script.sh -g ../0.tauri_genome.fna -a ../0.  
   tauri_annotation.gff -d ../ SRR3099585_chr18 S*86_chr18  
   S*87_chr18 S*97_chr18 S*98_chr18 S*99_chr18
```

Script in 3 main blocks

```
1 1) while getopts do ... done
```

```
2 2) for sample in $* ; do ... done
```

```
3 3) creation of the result file, counts.txt, with paste, awk,  
   and sed bash commands
```

Complete bash script, 1/3

getopts block

```
1 while getopts g:a:d: flag do
2     case $flag in
3         g) genome=$OPTARG
4           echo genome is $genome ;;
5         a) annots=$OPTARG
6           echo annotation is $annots ;;
7         d) rnadir=$OPTARG
8           echo RNAseq path is $rnadir ;;
9         :) echo "L'option $OPTARG requiert un argument"
10          exit 1 ;;
11        \?) echo "$OPTARG : option invalide"
12          exit 1 ;;
13     esac
14 done
15 shift $(( OPTIND - 1 )) # shift past the last flag or
16 echo samples are $*
```

Complete bash script, 2/3

for block

```
1 nbs=0;
2 for sample in $* ; do
3     nbs=$(expr ${nbs} + 1)
4     echo traitement of sample ${sample}
5     # ----- quality control of reads
6     if [ ! -d FastQC ]; then
7         mkdir FastQC
8     fi
9     fastqc --outdir FastQC ${rnadir}${sample}.fastq.gz >
10    FastQC/${sample}.log 2>&1
11    #----- reads mapping
12    if [ ! -d Bwt2_index ]; then
13        mkdir Bwt2_index
14        bowtie2-build ${genome} Bwt2_index/tauri > Bwt2_index
15        /Bwt2_index.log 2>&1
16    fi
```

Complete bash script, 3/3

for block, continuation

```
1 bowtie2 -x Bwt2_index/tauri -U ${rnadir}${sample}.fastq.  
gz -S ${sample}.sam > ${sample}_bowtie2.log 2>&1  
2 #----- selection and format modification  
3 samtools view -b ${sample}.sam -o ${sample}.bam  
4 samtools sort ${sample}.bam -o ${sample}_sort.bam  
5 samtools index ${sample}_sort.bam  
6 #----- counting of mapped reads by gene  
7 featureCounts -t gene -g ID -a ${annots} -s 2 -o ${sample  
}_ftc.txt ${sample}_sort.bam > ${sample}_ftc.log 2>&1  
8 done
```

Count table block

```
1 paste *_ftc.txt > ftc_tmp.txt  
2 awk -v nb=${nbs} -v col=7 'BEGIN{FS="\t"}{ctmp=$1; for(i=col  
;i<=nb*col;i=i+col){count=sprintf("%s\t%s",ctmp,$i);ctmp  
=count};print count}' ftc_tmp.txt | sed 1d > counts.txt
```

Exercise 2

Continue the snakefile of the previous exercise in order to replace the bash script.

We will:

Objectives

- add a configuration file
- use a builtin snakemake function to get filenames of the input RNAseq data
- add rules to replace the mapping, formatting, counting, and counts aggregating steps of the bash script

ex2_o1.smk

```
1 cp ex1_o7.smk ex2_o1.smk
```


The getopts block

Shell script

```
1 while getopts g:a:d: flag do
2     case $flag in
3         g) genome=$OPTARG
4         ...
```

We will use a configuration file:

Objective 1

Add a configuration file, named `RNAseq.yml`, containing both the genome sequence and the annotation files `manes`, and the access to the `Data` directory.

In the snakefile, change the configured variables (ex. replace `Data/` and `genome` by their `config[]` values). The Python strings concatenation is `+`. Then, run `snakemake` with the `--configfile` option.

Adding a configuration file

ex2_o1.yml

```
1 genome:
2   O.tauri.fna
3 annots:
4   O.tauri.gff
5 dataDir:
6   Data/
```

ex2_o1.smk: "Data/..." in inputs replaced by a config call:

```
1 rule genome_bwt2_index:  config["dataDir"]+config["genome"]
2 rule fastqc:            config["dataDir"]+"{sample}.fastq.gz"
```

snakemake run:

```
1 rm -Rf FastQC/ Result/ Tmp/ Logs/ ; snakemake -s ex2_o1.smk
   --configfile ex2_o1.yml
```

The for block

Shell script

```
1 nbs=0;
2 for sample in $* ; do
3     ...
4 done
```

To manage all *.fastq.gz files in a directory, use the `glob_wildcards()` function. In `ex2_o2.smk`, replace the `SAMPLES` definition by:

ex2_o2.smk

```
1 SAMPLES , = glob_wildcards (config ["dataDir"]+"{sample}.fastq.
   gz")
```

and run `snakemake`.



Quality control, fastqc

```
1 if [ ! -d FastQC ]; then
2     mkdir FastQC
3 fi
4 fastqc --outdir FastQC ${sample}.fastq.gz > FastQC/${sample}
   }.log 2>&1
```

No more need to test the existence of a directory, it is created as needed.

rule fastqc:

This rule was already present in the snakefile

Reads mapping, bowtie2

```
1 if [ ! -d Bwt2_index ]; then
2     mkdir Bwt2_index
3     bowtie2-build ${genome} Bwt2_index/tauri > Bwt2_index/
      Bwt2_index.log 2>&1
4 fi
5 bowtie2 -x Bwt2_index/tauri -U ${sample}.fastq.gz -S ${
      sample}.sam > ${sample}_bowtie2.log 2>&1
```

2 rules: genome_bwt2_index (cf. previous ex.) and bwt2_mapping

ex2_o3.smk, rule bwt2_mapping (no run):

```
1 output: "results/{sample}.sam"
2 input: config["dataDir"]+"{sample}.fastq.gz",
3         expand("Tmp/0tauri.{ext}.bt2", ext=BIDX)
4 log: "Logs/{sample}_bwt2_mapping.log"
5 shell: "bowtie2 -x Tmp/0tauri -U {input[0]} -S {output} 2>
      {log} "
```

Reads mapping, bowtie2

Why some troubles?

The snakemake launch probably didn't do what was expected. What have we forgotten?

We added a new rule to a snakemake but we didn't manage the rule tree, there is no input-output link to include the new rule to the workflow.

We will do that by completing the input directive of the target rule (caution to respect the Python "list" structure, coma-separated).

ex2_o2.smk, target rule:

```
1 rule all:
2     input:
3         expand("FastQC/{sample}_fastqc.html", sample=SAMPLES),
4         expand("Tmp/0tauri.{ext}.bt2", ext=BIDX),
5         expand("Tmp/{sample}.sam", sample=SAMPLES)
```

samtools

Shell script

```
1 samtools sort -O bam -o ${sample}_sort.bam ${sample}.sam
2 samtools index ${sample}_sort.bam
```

ex2_o4.smk, rule sam2bam_sort (no run):

```
1 output:
2   bam="Result/{sample}_sort.bam",
3   bai="Result/{sample}_sort.bam.bai"
4 input: "Tmp/{sample}.sam"
5 log:
6   sort="Logs/{sample}_sam2bam_sort.log",
7   index="Logs/{sample}_bam2bai.log"
8 shell:
9   "samtools sort -O bam -o {output.bam} {input} 2> {log.
10  sort} ;"
   "samtools index {output.bam} 2> {log.index}"
```

FeatureCount

Shell script

```
1 featureCounts -t gene -g ID -a ${annots} -s 2 -o ${sample}_ftc.txt ${sample}_sort.bam > ${sample}_ftc.log 2>&1
```

ex2_o5.smk, rule counting (params; no run):

```
1 output: "Tmp/{sample}_ftc.txt"  
2 input:  
3     bam="Result/{sample}_sort.bam",  
4     annot=config["dataDir"]+config["annots"]  
5 params: t="gene", g="ID", s="2"  
6 log: "Logs/{sample}_counts.log"  
7 shell: "featureCounts -t {params.t} -g {params.g} -a {  
     input.annot} -s {params.s} -o {output} {input.bam} &> {  
     log}"
```


Counts matrix creation

Shell script

```
1 paste *_ftc.txt > counts_tmp.txt
2 awk -v nb=${nb_sample} 'BEGIN{FS="\t"}{count_tmp=$1; for(i
   =7;i<=nb*7;i=i+7){count=sprintf("%s\t%s",count_tmp,$i);
   count_tmp=count};print count}' counts_tmp.txt | sed 1d >
   counts.txt
```

Hint

Create 2 rules to manage some files aggregation to one result file:

- rule `extract_counts`: extract geneID and counts in individual files
- rule `matrix_counts`: paste these files



Counts matrix creation

ex2_o6.smk (2 rules, shell 3", copy, run):

```
1 rule matrix_counts:
2   output: "Result/counts_matrix.txt"
3   input: countfile=expand("Tmp/{sample}_ftc7.txt", sample=
4         SAMPLES), geneID=expand("Tmp/{sample}_ftc1.txt", sample=
5         SAMPLES)
6   log: "Logs/matrix_counts.log"
7   shell: ""cp {input.geneID[0]} Tmp/ftc_geneID.txt > {log}
8         ; paste Tmp/ftc_geneID.txt {input.countfile} > {output}
9         > {log}""
10
11 rule extract_counts:
12   output: col7="Tmp/{sample}_ftc7.txt",
13          col1="Tmp/{sample}_ftc1.txt"
14   input: "Tmp/{sample}_ftc.txt"
15   log: "Logs/{sample}_extract_counts.log"
16   shell: ""cut -f 7 {input} | sed 1d > {output.col7} > {log}
17         ; cut -f 1 {input} | sed 1d > {output.col1} ""
```

DESeq2

The DESeq2 step is the statistical analysis. From the count matrix, the statistical analysis is managed by a non parallelizable R script, DESeq2.

So, up to date, the workflow is complete and the only thing left is this DESeq2 statistical analysis. We will see this analysis through the notebooks session.



Bonus

Add a help rule

```
https://lachlandeer.github.io/snakemake-econ-r-tutorial/  
self-documenting-help.html#a-help-rule
```