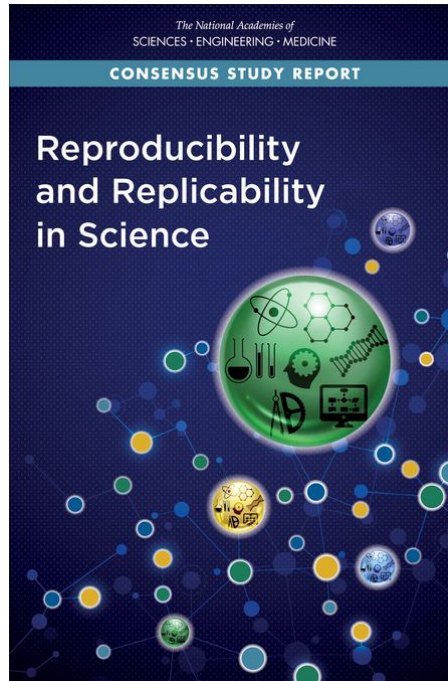# Reproducibility at the service of Computational Biology

Thomas DENECKER

June 11, 2019
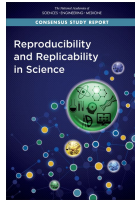
# Reproducibility in Science in 2019



*National Academies of Sciences, Engineering, and Medicine. 2019. Reproducibility and Replicability in Science. Washington, DC: The National Academies Press. https://doi.org/10.17226/25303.*
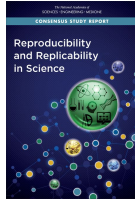
# What is reproducibility?



**Reproducibility**

Obtaining consistent computational results using the same input data, computational steps, methods, codes, and conditions of analysis

# Reproducibility *versus* Replicability
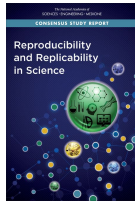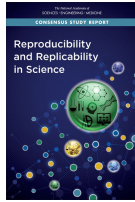
**Reproducibility**

Obtaining consistent computational results using the same input data, computational steps, methods, codes, and conditions of analysis

≠

**Replicability**

Obtaining consistent results across studies aimed at answering the same scientific question, each of which has obtained its own data

# Reproducibility *versus* Repeatability
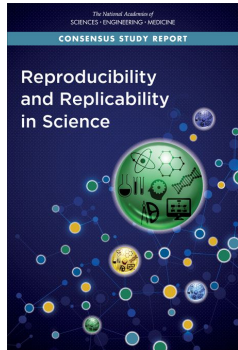
**Reproducibility**
Obtaining consistent computational results using the same input data, computational steps, methods, code, and conditions of analysis

≠

**Repeatability** (adapted from de Hans E. Plesser, *Front. Neuroinf.* , 2017)
Obtaining the closest possible results by performing an identical experiment (methods, equipment, experimenters, laboratory and conditions)

# Recommendations for reproducibility in 2019

**Description of the experimental part**
Methods, instruments, procedures, measurements, experimental conditions

**Description of the computational part**
Steps in data analysis and technical choices

**Description of the statistical part**
Analytical decisions: when, how, why

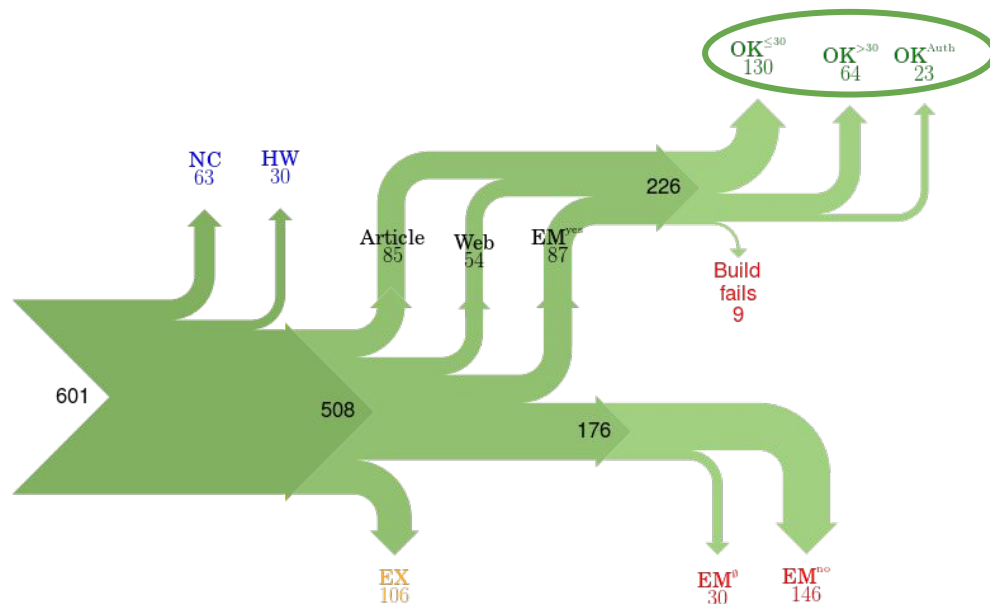**Discussion of the choices and results obtained**

# In reality ?

# A problem of reproducibility in Biology

# 70 %

analyses in Experimental Biology **are not reproducible**

(Monya Baker, *Nature*, 2016)

# In computer science



Collberg *et al*, University of Arizona TR 14-04, 2015

# Bioinformatics is also affected...

Problems still too frequent

- **Impossibility to install tools**
    - OS not compatible
    - Dependency no longer available
- **Update of the tool rendering the codes unusable**
    - Python 2 and Python 3 !
    - Changing the arguments of the functions used ( R )
- **Impossibility to reproduce the results of the computational analysis**
    - IDE : stable version of the language different according to the OS (Rstudio)
    - Package versions

# Fortunately, there are solutions!

# When bioinformatics meets development solutions



SNAKEMAKE

...

# An example of reproducibility in Bioinformatics

I2BC training with Claire Toffano-Nioche

**Use the FAIR data principles to make an analysis protocol reproducible and always obtain the same results from the same data**

All courses and codes are open source

https://github.com/thomasdenecker/FAIR_Bioinfo

# Proposal of a solution

**Equivalents for each choice**



...

# The FAIR data principles

**F**indable          **A**ccessible          **I**nteroperable          **R**eusable

# How to use the FAIR data principles?

**F**

**A**

**I**

**R**

# How to use the FAIR data principles?

🔍 **Findable**

- Tools used = references in their field
- Easy to find analysis protocol (Github pages)

👆 **A**

⚙️ **I**

♻️ **R**

# How to use the FAIR data principles?

🔍 **Findable**

- Tools used = references in their field
- Easy to find analysis protocol (Github pages)

👆 **Accessible**

- Available resources (Github, dockerhub)
- Open source tools (conda)

⚙️ **I**


♻️ **R**

# How to use the FAIR data principles?

⌕ **Findable**

- Tools used = references in their field
- Easy to find analysis protocol (Github pages)

✋ **Accessible**

- Available resources (Github, dockerhub)
- Open source tools (conda)

⚙ **Interoperable**

- Cooperation of tools (snakemake, docker) both locally and on servers (cloud or cluster)

♺ **R**

# How to use the FAIR data principles?

## 🔍 Findable

- Tools used = references in their field
- Easy to find analysis protocol (Github pages)

## 👆 Accessible

- Available resources (Github, dockerhub)
- Open source tools (conda)

## ⚙️ Interoperable

- Cooperation of tools (snakemake, docker) both locally and on servers (cloud or cluster)

## ♻️ Reusable

- Protocol that can be simply replayed (snakemake) in the same way (Rmarkdown) in a virtual environment (docker)

# Our creed : So FAIR !

**FAIR raw data**

**+**

**FAIR_bioinfo scripts/protocols**

**=**

**FAIR processed data**

# The data



Stress

No stress

**Condition 1 (3x)**     *Repeatability*     **Condition 2 (3x)**

**Measurements 1** (3x)     **Measurements 2** (3x)

**Differences ?**

# Level of gene expression

# Differences between conditions

| | Condition 1 | | | Condition 2 | | | |
|---|---|---|---|---|---|---|---|
| | **1** | **2** | **3** | **4** | **5** | **6** | |
| **Gene 1** | 1 | 1 | 1 | 3 | 3 | 3 | ≠ |
| **Gene 2** | 3 | 3 | 3 | 1 | 1 | 1 | ≠ |
| **...** | … | … | … | … | … | … | |
| **Gene n** | 2 | 2 | 2 | 2 | 2 | 2 | = |

**Conclusion**
Genes 1 and 2 differentially expressed between condition 1 and 2

# How to process this data?

# Analysis pipeline

# What is done mostly

1. Install the tools locally (sometimes writing an installation script)
2. Write a script to run all the analyses (not always...)
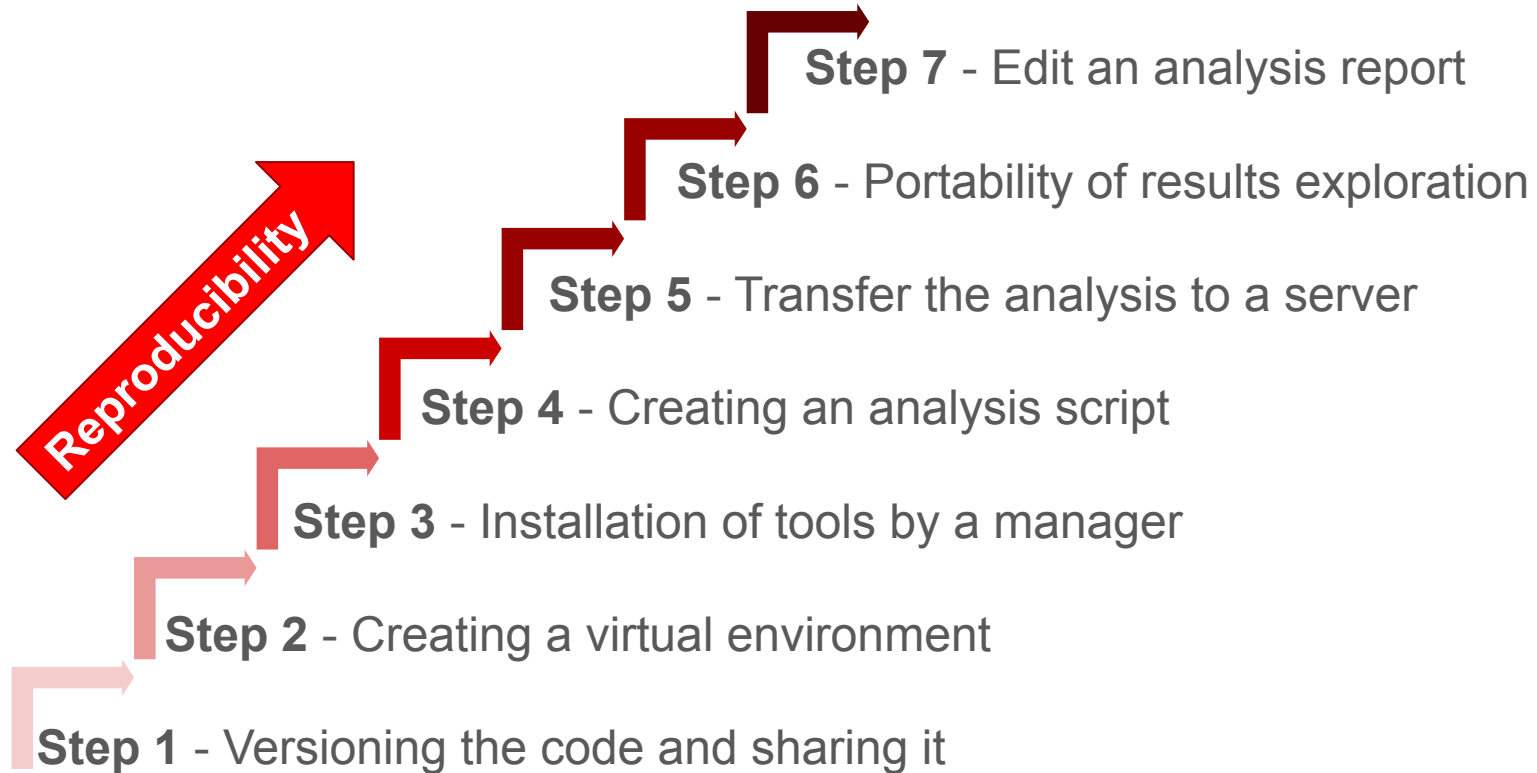3. Sharing the script (by publishing, by email, USB key,...)

But low reproducibility ...

# What are the solutions to be more reproducible in Bioinformatics?

# A 7-step solution

**Step 7** - Edit an analysis report

**Step 6** - Portability of results exploration

**Step 5** - Transfer the analysis to a server

**Step 4** - Creating an analysis script

**Step 3** - Installation of tools by a manager

**Step 2** - Creating a virtual environment

**Step 1** - Versioning the code and sharing it

# A 7-step solution



**Step 7** - Edit an analysis report

**Step 6** - Portability of results exploration

**Step 5** - Transfer the analysis to a server

**Step 4** - Creating an analysis script

**Step 3** - Installation of tools by a manager

**Step 2** - Creating a virtual environment

**Step 1** - Versioning the code and sharing it

Reproducibility

# Step 1 - Versioning the code and sharing it

**Why?**

- Have the right version of the code
- Vision over time
- Open to the community

# Step 1 - Versioning the code and sharing it

**Before**

# Step 1 - Versioning the code and sharing it



**Before**



**Advantages**

- Saving the code
- Simple to share
- Automatic version management

**Disadvantages**

- Not easy for novices

# Step 1 - Versioning the code and sharing it

**Before**

**After**

# Step 2 - Creating a virtual environment

**Why ?**
- Same the environment for all
- Sharing the environment

# Step 2 - Creating a virtual environment
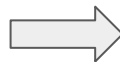
**Before**

# Step 2 - Creating a virtual environment

**Before**



**Advantages**

- Fast and lightweight
- Portable
- Easy to share and deploy

**Disadvantages**

- root
- Up-to-date system

# Step 2 - Creating a virtual environment

**Before**

**After : Ubuntu 16.04**



```
$ cat > dockerfile
FROM ubuntu:16.04
RUN apt-get update

# Set environment variables
ENV HOME /root

# Define working directory
WORKDIR /root

# Define default command
CMD ["bash"]

$ docker --tag=toto build .
$ docker run toto
```

# Step 3 - Installation of the tools by a manager

**Why ?**
- Good version
- Simply install

# Step 3 - Installation of the tools by a manager

**Before : FastQC**

1) Download the source
2) Unzip the file
3) Install and update Java (many problems)
4) Changing rights

# Step 3 - Installation of the tools by a manager

**Before : FastQC**

1) Download the source
2) Unzip the file
3) Install and update Java (many problems)
4) Changing rights

CONDA

**Advantages**

- Easy to install manager
- Easy package installation
- Version management

**Disadvantages**

- May be heavy (miniconda solution)
- Missing packages (R)

# Step 3 - Installation of the tools by a manager

**Before : FastQC**

1) Download the source
2) Unzip the file
3) Install and update Java (many problems)
4) Changing rights

CONDA

**Après**

```
$ conda install -c bioconda -y
fastqc=0.11.2
```

All the tools used in the protocol are available on Conda (https://anaconda.org/) : bowtie2, samtools, htseqcount, aspera, snakemake, ...

Installation as easy as that

# Step 4 - Creating an analysis script

**Why ?**

- Have a reproducible analysis script
- Do not repeat what has already been done
- Parallelize

# Step 4 - Creating an analysis script

**Before (Shell script)**

```
for sample in `ls *.fastq.gz`
do
  fastqc ${sample}
done
```

# Step 4 - Creating an analysis script

**Before (Shell script)**

```
for sample in `ls *.fastq.gz`
do
  fastqc ${sample}
done
```
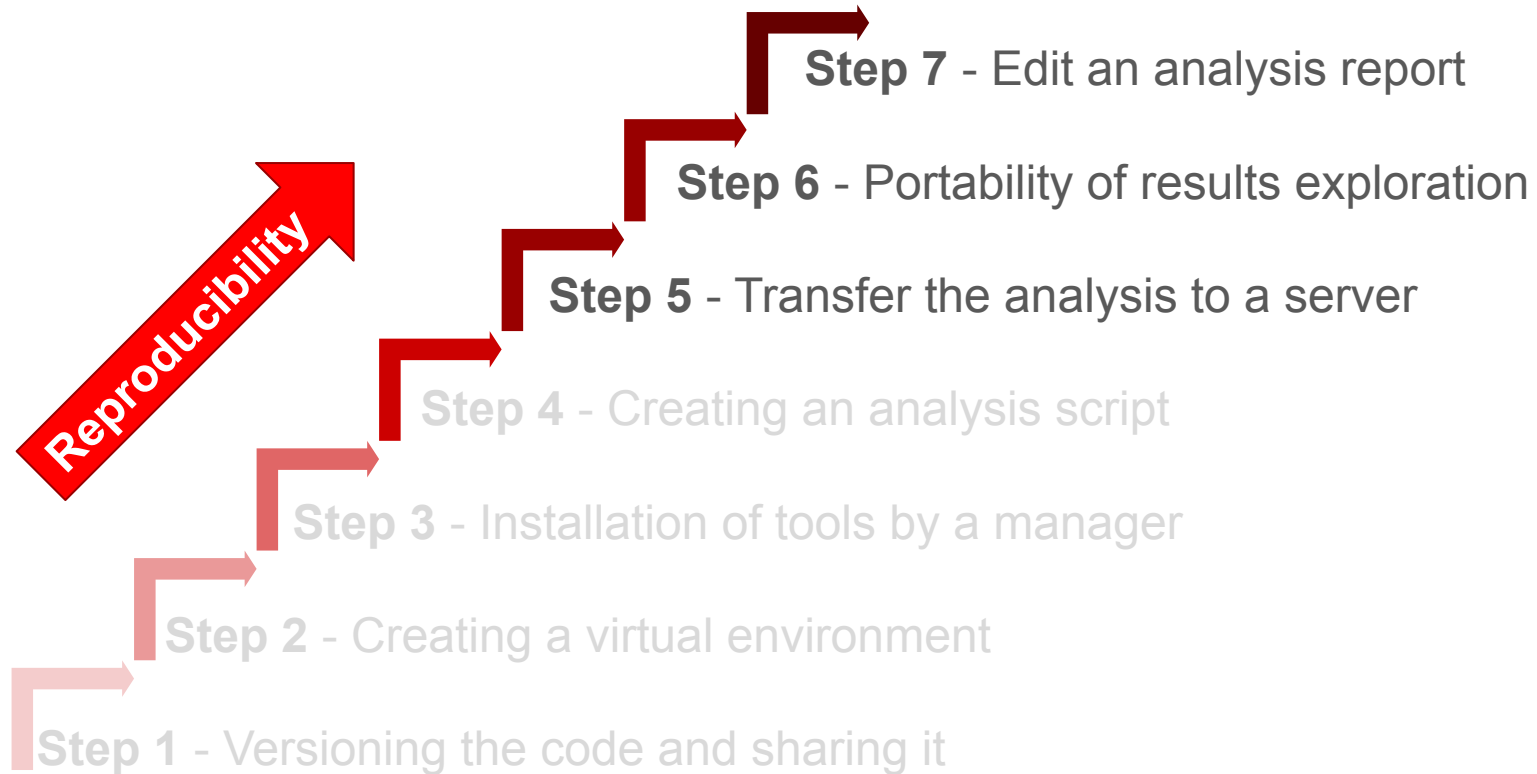
SNAKEMAKE

**Advantages**
- Workflow (job management)
- Powerful and fast

**Disadvantages**
- A logic to be taken
- Syntax less simple than shell script

# Step 4 - Creating an analysis script

**Before (Shell script)**

```
for sample in `ls *.fastq.gz`
do
  fastqc ${sample}
done
```

Shorter to write but not to execute

Snakemake = Parallel

**After (Snakefile)**

```
$ cat > Snakefile
SAMPLES, =
glob_wildcards("./samples/{smp}.fastq.gz")

rule final:
input:expand("fastqc/{smp}/{smp}_fastqc.zip"
,smp=SAMPLES)

rule fastqc:
    input:  "samples/{smp}.fastq.gz"
    output:"fastqc/{smp}/{smp}_fastqc.zip"
    message: """Quality check"""
    shell: """fastqc {input} --outdir
fastqc/{wildcards.smp}"""
$ snakemake
```

# A 7-step solution

**Step 7** - Edit an analysis report

**Step 6** - Portability of results exploration

**Step 5** - Transfer the analysis to a server

Step 4 - Creating an analysis script

Step 3 - Installation of tools by a manager

Step 2 - Creating a virtual environment

Step 1 - Versioning the code and sharing it

Reproducibility

# Step 5 - Transfer the analysis to a server

**Why ?**
- Controlled environment
- Offset of the analysis

# Step 5 - Transfer the analysis to a server

**Before**

**Adaptation locally and on servers that is difficult or even unmanaged ...**

# Step 5 - Transfer the analysis to a server

**Before**

**Adaptation locally and on servers that is difficult or even unmanaged ...**

**Advantages**
- Easy to set up
- Increase in power (cloud or cluster)
- For everyone

**Disadvantages**
- Not easy for novices

# Step 5 - Transfer the analysis to a server

**Before**

**After**

**Adaptation locally and on servers that is difficult or even unmanaged ...**

```
$ git clone
https://github.com/thomasdenecker/FAIR_Bioinfo

$ cd FAIR_Bioinfo

$ sudo docker run --rm -d -p 80:8888 --name
fair_bioinfo -v ${PWD}:/home/rstudio
tdenecker/fair_bioinfo bash ./FAIR_script.sh
```

**The protocol is running !**

# Step 6 - Portability of the results exploration

**Why ?**
- Make it easy to explore
- Easy to share

# Step 6 - Portability of the results exploration

**Before : R terminal**

```
dds <- DESeqDataSetFromMatrix(countData =
cts,colData = coldata, design= ~ batch +
condition)

dds <- DESeq(dds)
resultsNames(dds) # lists the
coefficients
res <- results(dds, name =
"condition_trt_vs_untrt")

# or to shrink log fold changes
# association with condition:
res <- lfcShrink(dds,
coef="condition_trt_vs_untrt",
type="apeglm")
```

# Step 6 - Portability of the results exploration

**Before : R terminal**

```
dds <- DESeqDataSetFromMatrix(countData =
cts,colData = coldata, design= ~ batch +
condition)

dds <- DESeq(dds)
resultsNames(dds) # lists the
coefficients
res <- results(dds, name =
"condition_trt_vs_untrt")

# or to shrink log fold changes
# association with condition:
res <- lfcShrink(dds,
coef="condition_trt_vs_untrt",
type="apeglm")
```

**Advantages**
- Portable (HTML)
- Accessible everywhere
- Interactive (configurable, dynamic graphs,...)

**Disadvantages**
- Mixing R and HTML

# Step 6 - Portability of the results exploration

**Before : R terminal**

```
dds <- DESeqDataSetFromMatrix(countData =
cts,colData = coldata, design= ~ batch +
condition)

dds <- DESeq(dds)
resultsNames(dds) # lists the
coefficients
res <- results(dds, name =
"condition_trt_vs_untrt")

# or to shrink log fold changes
# association with condition:
res <- lfcShrink(dds,
coef="condition_trt_vs_untrt",
type="apeglm")
```

**After**

# Step 7 - Edit an analysis report

**Why ?**

- Have a trace of the analysis
  (date, time, parameters,...)
- Store tool versions

# Step 7 - Edit an analysis report

**Before**

# Step 7 - Edit an analysis report

**Before**



**Advantages**
- Simple syntax (Markdown)
- Sharing (PDF, HTML,...)

**Disadvantages**
- Rare visualization problems in LaTeX

# Step 7 - Edit an analysis report
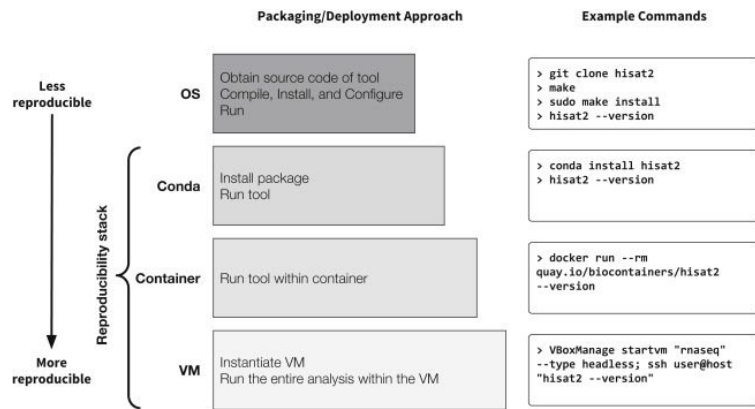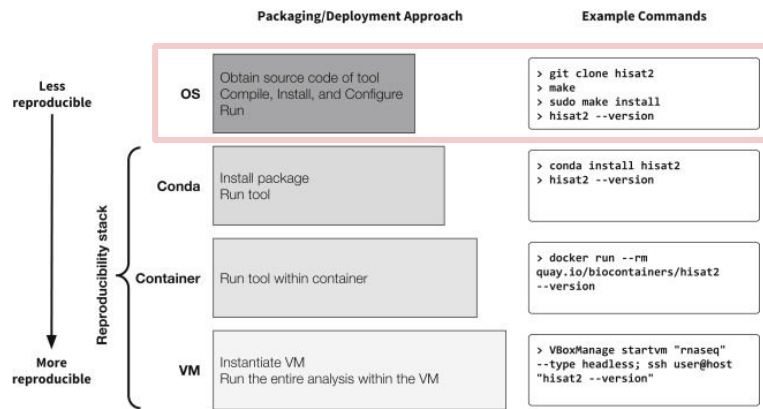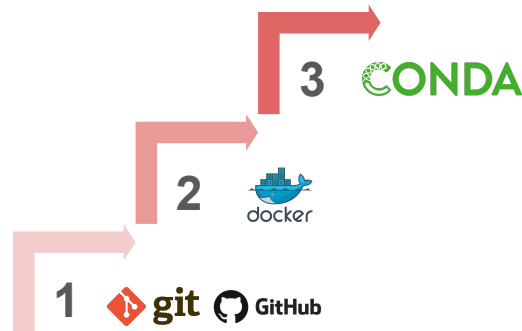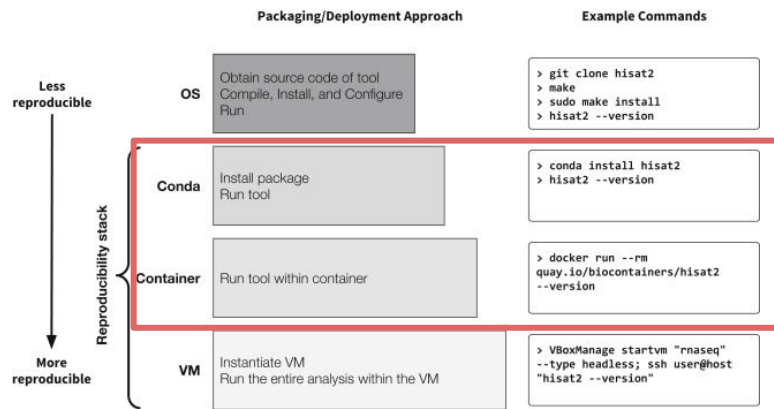
**Before**

**After**

# Conclusion

# What is our level of reproducibility?



**Practical Computational Reproducibility in the Life Sciences,** Björn Grüning *et al*, 2018

# What is our level of reproducibility?



**Practical Computational Reproducibility in the Life Sciences,** Björn Grüning *et al*, 2018
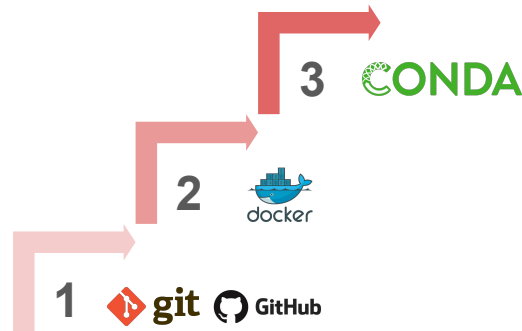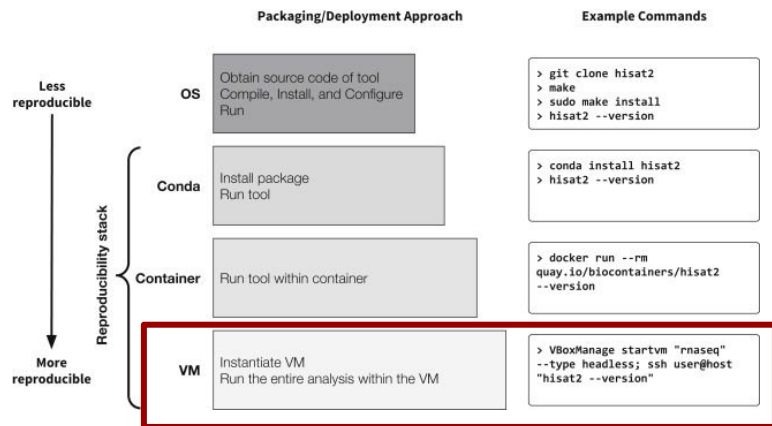
# What is our level of reproducibility?



**Practical Computational Reproducibility in the Life Sciences,** Björn Grüning *et al*, 2018
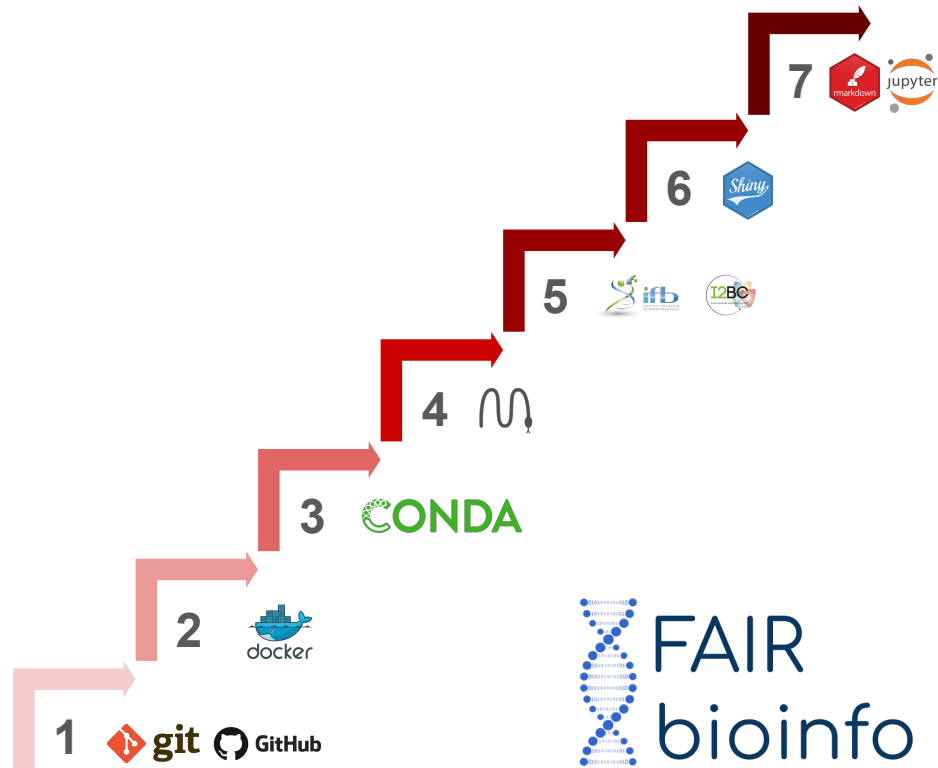
# What is our level of reproducibility?



**Practical Computational Reproducibility in the Life Sciences,** Björn Grüning *et al*, 2018
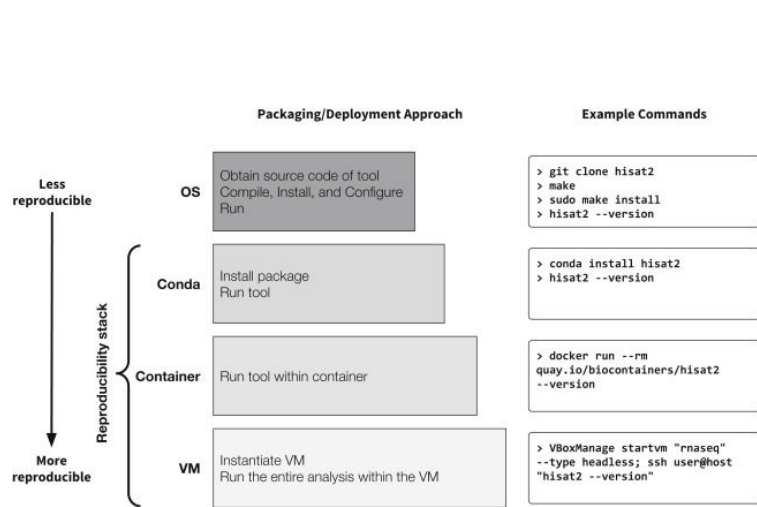
# What is our level of reproducibility?



**Practical Computational Reproducibility in the Life Sciences,** Björn Grüning *et al*, 2018

# What is our level of reproducibility?



**Practical Computational Reproducibility in the Life Sciences,** Björn Grüning *et al*, 2018
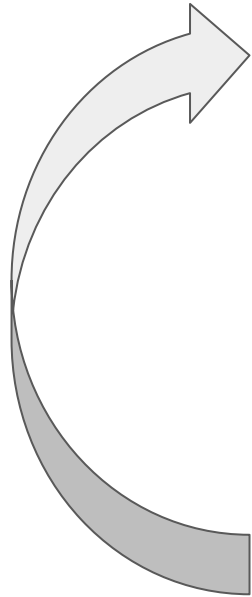
# Take home messages

A real reflection on the reproducibility of analyses in Bioinformatics

Proposal of a solution that helps to make any analytical protocol reproducible

**Reproducibility is an real plus value for Bioinformatics!**

# A virtuous circle



**FAIR raw data**

**+**

**FAIR_bioinfo scripts/protocols**

**=**

**FAIR processed data**