

REPRODUCIBILITY

TESTIMONIAL



Testimonial as a learner

What I thought I made



What I actually did



What I actually did



Difficulty

- to track bug
- to update tools
- to collaborate on the code
- to share (image >20Go)
- etc

Reproducibility is nothing without comprehensiveness

Version control system was quiet new. (Few around was using it, not learnt at school)
Conda / Container didn't exist.

=> naming/
versioning

```
script.py  
scriptFinal.py  
scriptFinalV2.py  
scriptUSE_THIS_ONE.py  
scriptV2.py
```

=> comment

```
#####  
# Script v1. 07/04/2007  
# Author: Jacques Dainat  
# Modified by JD 15/02/2008 Add feature1  
# Modified by JD 27/11/2008 Add feature2  
# Modified by Colleague1 11/08/2009 modify feature1 to fix error
```

=> readme

```
----- Script -----  
Purpose  
Installation  
Usage
```



- naming 

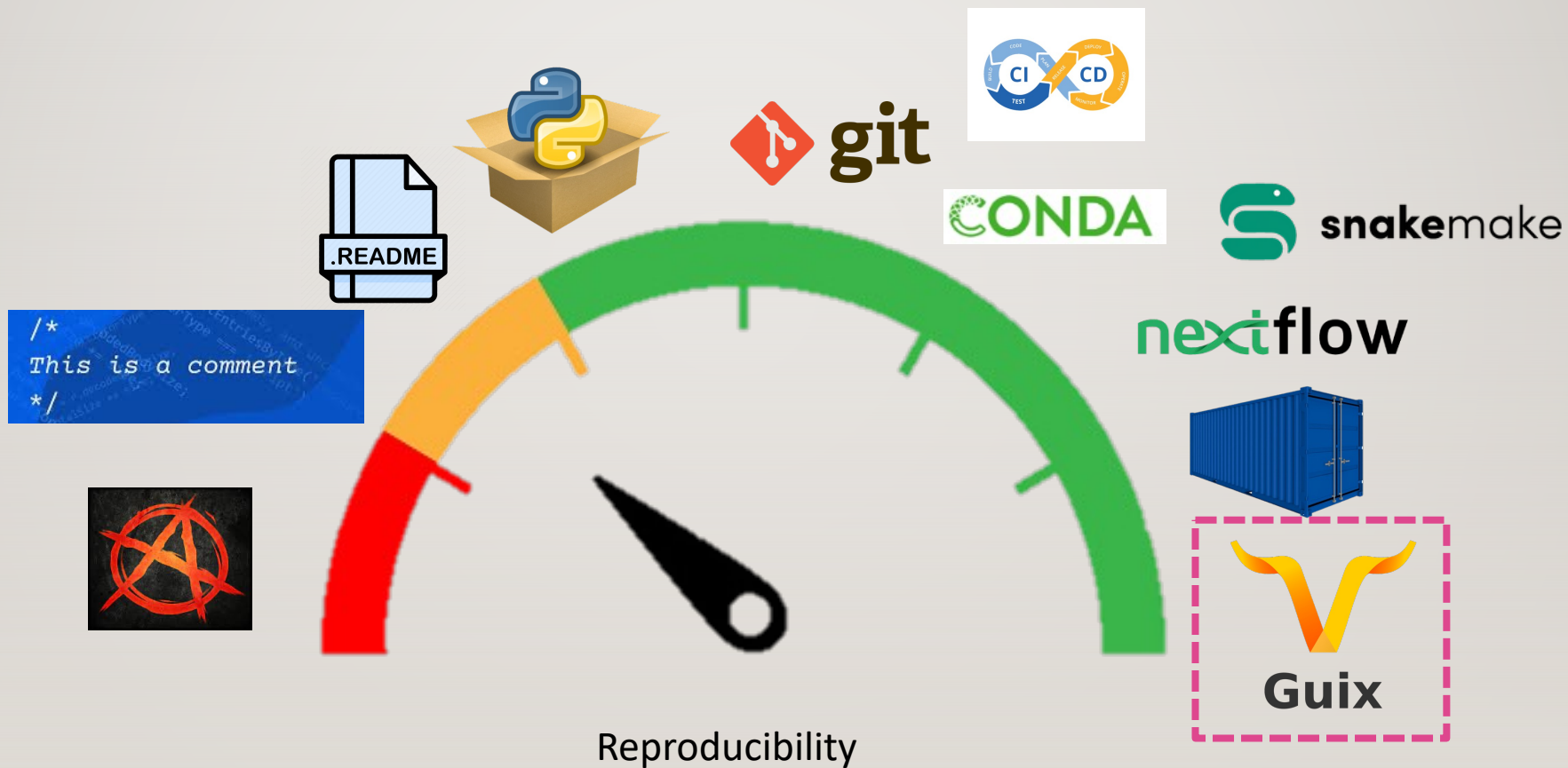
- comment  => 

- Readme
Doc

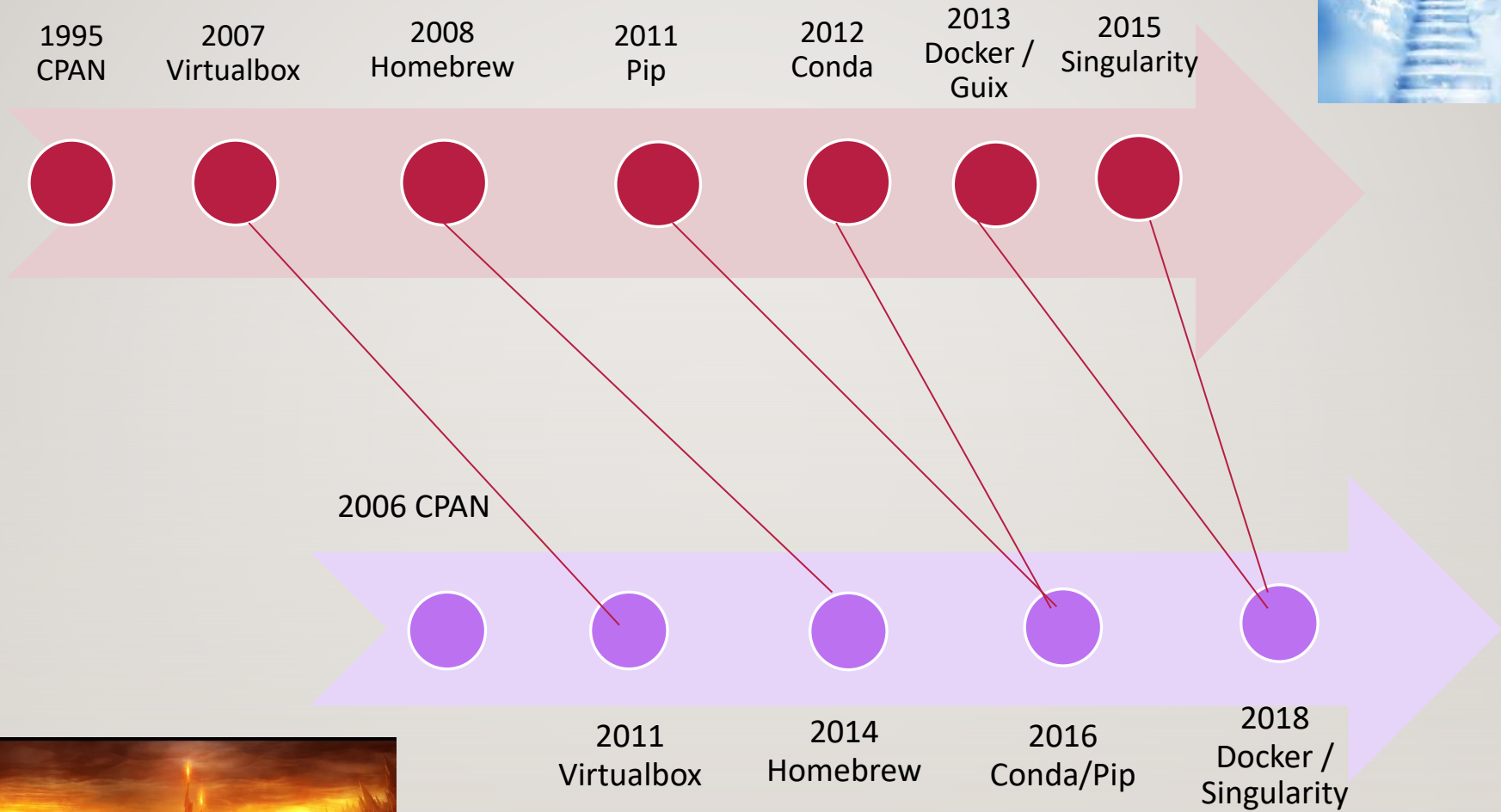


Programming / Development

- Many tool available depending your needs (level of reproducibility)
- Keep an eye on new tools (future game changer), and decide when to make the step



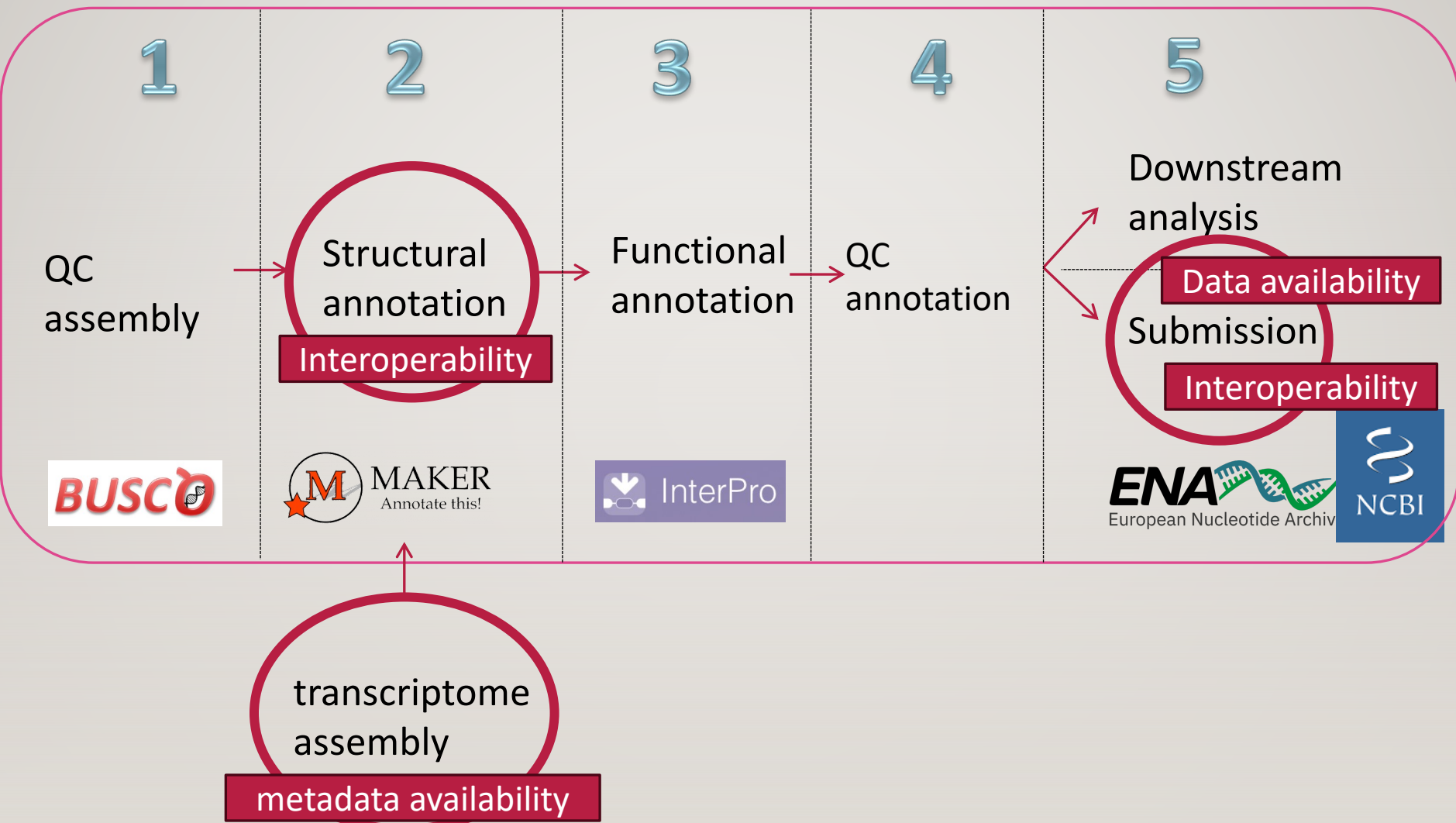
Dependencies HELL





- on demand
- on hard drive

The main steps in genome annotation



install with **bioconda** container **none** JOSS **10.21105/joss.01344** DOI **10.5281/zenodo.3268394** pypi package **0.2.5**
 downloads **1k total** license **GPLv3**

GUESSmyLT

Software to guess the RNA-Seq library type of paired and single end read files using mapping and gene annotation.

Languages

Python 95.7% TeX 4.3%

! Data submitted to ENA was not exactly what ended up in the DB...

Discrepancies between publication and data in public archive for the functional annotation.

CI **passing** DOI **10.1186/s13104-018-3686-x** install with **bioconda** downloads **6.1k** container **Docker** container **Singularity**
 license **GPLv3**

EMBLmyGFF3

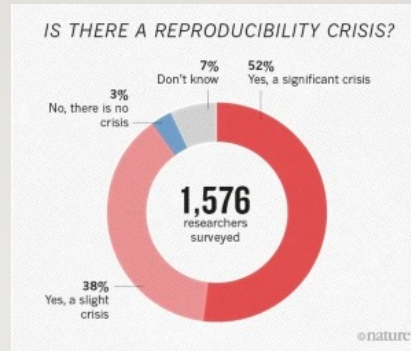
GFF3 to EMBL conversion tool

EMBLmyGFF3 converts an assembly in FASTA format along with associated annotation in GFF3 format into the EMBL flat file format which is the required format for submitting annotated assemblies to ENA.



DATA

Guiding Principles for scientific data management and stewardship



(<https://doi.org/10.1038/533452a>)

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



D M P

Data Management Planning

(ANR 2019)

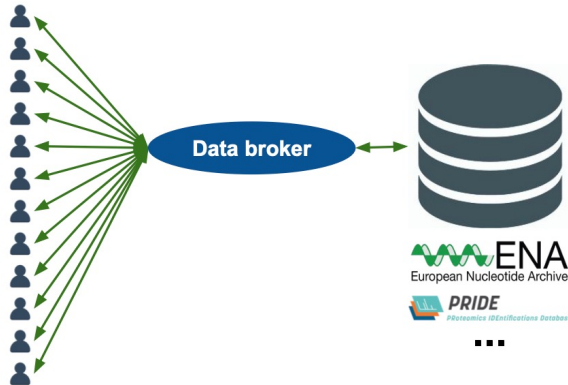




4 FTE staff dedicated to data management

13 FTE Systems development

Un data broker



Intermédiaire entre le producteur de données et la ressource internationale de stockage/archivage
 Rationalisation des échanges



35 FTE Support

12 FTE Training

...

Testimonial as a teacher

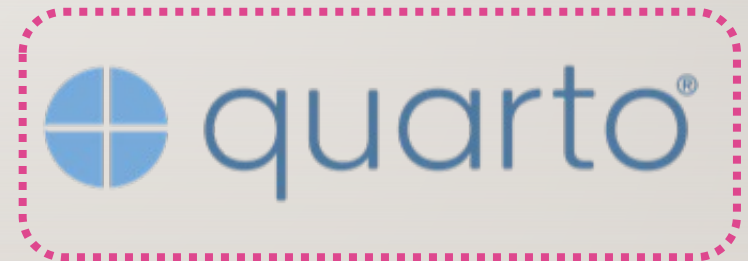
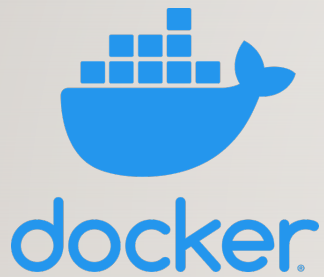


Tools for Reproducible Research

10 courses since may 2019



Tools covered



Tools are one thing, but you need to know how to you use them properly...

(check license!)

- Git hosted static website
Rendering: ReadTheDocs => Mkdocs
- Lectures: .Rmd compiled => PPTX
- Nextflow Snakemake 50/50 => 20/80 => 0/100
- We choose to ask help from IFB via Thomas Denecker
- We added Quarto
- We choose to install tools during the course (conda, then conda env)

It is time expensive

- numerous meeting
- split course according to knowledge of each teacher
- check material
- update slides
- adapt material (command git, problem with conda)

Surprised by the appeal of data management (Round table)



We did not have expert for that and questions were beyond our expertise.

Conda was a mess!!!
(Remind: Conda share pre-compiled binaries.)



=> We should have prepared a environment file and test them for each OS (windows, linux, macOS).



x86

architecture

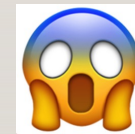
VS



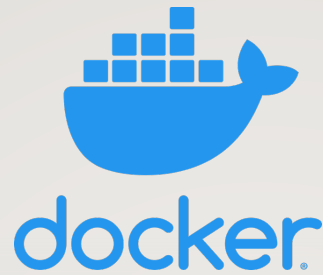
ARM (Apple Silicon)

=> We checked conda was working on ARM but didn't check that all dependencies were available for this architecture.

=> All the course was based on conda environment



What we experienced



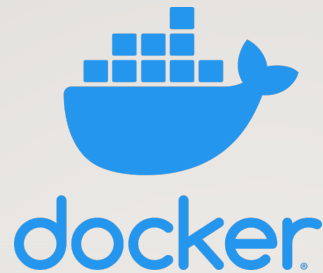
- Version beta but worked fine



ARM (Apple Silicon)



- Must install UTM (VM) to install a Rocky Linux (Linux OS ARM architecture) to install singularity



ARM (Apple Silicon)

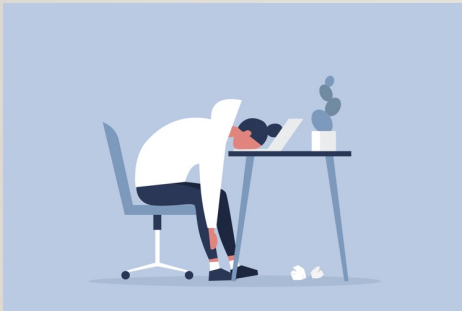
Images from biocontainer are not all ARM compatible (cross-plattform compilation)

Emulation of x86 platform can make tools very slow
Rosetta2 for emulation works great (but need to be present on the macOS computer)

- teach conda but avoid to use it outside this part of the course
- Teach container early in the course
=> Use container for all part of the course
- Introduce computer processors architectures?
- Teach Cross-Compilation container (or check rosetta2 presence for macOS users)
- Test (CI?) on all

- Greater cohesion within the team that set up the course
- Networking with students, useful feedback
- Transfer of skills between teachers
- Awareness of reproducibility limits

Rewards



Is setting a second session of the same course as demanding?

The End

Thank you for your attention

Jacques.dainat@ird.fr