



2025 June 6th

### Introduction to FAIR software environnement Conda, Docker and Apptainer

Julien Seiler



In bioinformatics, **reproducibility** refers to the ability to replicate the results of an experiment or analysis using the same data and computational steps.

### **Challenges in Achieving Reproducibility**

- **Complexity of Workflows**: Bioinformatics workflows often involve multiple steps, **tools**, and data sources, making them complex and difficult to replicate.
- **Dependency Management**: Different software versions, libraries, and dependencies can lead to inconsistencies when trying to reproduce results.
- Environment Variability: Differences in computing environments, such as operating systems, hardware, and installed software, can affect the reproducibility of results.





Imagine I want to share a Python script with all of you...

Which version of Python is available on your laptop ?

- 1. Open a Terminal console
- 2. Type:python --version or python3 --version







### Main changes introduced in latest Python 3 versions

Version	Main Changes
3.5	Asyncio introduced for asynchronous programming, Type hints, and the @ operator for matrix multiplication.
3.6	Formatted string literals (f-strings), Underscores in numeric literals, Asynchronous generators and comprehensions.
3.7	Data classes, Context variables, Guaranteed dictionary ordering, and Performance improvements.
3.8	Assignment expressions (walrus operator :=), Positional-only parameters, and Improved type hints.
3.9	Dictionary union operators, Type hinting generics in standard collections, and Flexible function and variable annotations.
3.10	Structural pattern matching, Parenthesized context managers, and Precise error messages.
3.11	Faster runtime due to optimizations, Improved error messages with more detail, and Introduction of the ExceptionGroup exception type.
3.12	More flexible f-strings, Improved error messages, and Typing and performance enhancements.



### Installing software is hard





### Installing software can be a tricky operation:



To facilitate software deployment, a number of package systems are available:

- offers pre-compiled software versions
- automatically manages dependencies
- offers updates

However, these package systems have a number of **drawbacks**:

- require administrative rights
- often only one version of a tool can be installed (sometimes without choice)
- poor support of scientific software
- hard to contribute



# Introduction to CONDA







### What is Conda?

Conda is an open-source package management system and environment management system.

It helps in installing, running, and updating tools and their dependencies.





### **Key Features**

- Does not require administrator rights
- Allows you to **choose the version** of the tool you want to install
- Allows multiple versions of the same tool to coexist (through isolated environments)
- Proposes many tools in the scientific field (bioconda)
- A very open contribution model
- Compatible with any Linux distribution
- Also compatible with **MacOS** and **Windows** (in theory)







### How does it work ?

- Each software and library is available in the form of a package (a tar bz2 archive).
- A package contains the compiled version of the software and the list of packages on which it depends.







### How does it work ?

• The packages are hosted on a central server: anaconda.org





Free

Organised by channels





### By the way, what is a channel ?

A channel is a set of packages available on the conda repositories and each channel is managed by a dedicated organisation.

The creators of conda offer a set of packages in the ANACONDA channel.

For several years, two main channels propose most of the packages:

- CONDA-FORCE : for all user tools, languages and libraries
- BIOCONDA: for bioinformatics tools and libraries.

There are also numerous channels proposed by conda users. On these channels you can sometimes find tools in more recent versions, **but reliability is not guaranteed**.





### How does it work ?

• The anaconda.org website has a search engine that lets you look for packages across all the channels.

ANACO	NDA. ORG	About Anaco	onda Help	Download Anacon	ida Sign
samtools					۹
▼ Filters — Type: All \	~	Access: All ~	Platfor	m: All ~	
					Platform
25	4977181	<b>O bioconda / samtools</b> 1.19.2 Tools for dealing with SAM, BAM and CRAM files		conda	linux-64 linux- aarch64 osx-64
2	878734	O bioconda / bioconductor-rsamtools 218.0 Binary alignment (BAM), FASTA, variant call (BCF), and tabix file impo	ort	conda	linux-64 osx-64
3	215435	Soil / samtools 1.6 Tools for dealing with SAM, BAM and CRAM files		copy conda	linux-64 osx-64
0	193215	O bioconda / perl-bio-samtools 143 Read SAM/BAM files		conda	linux-64 osx-64
0	8097	O bioconda / msamtools 11.3 microbiome-related extension to samtools		conda	linux-64 osx-64
0	5555	O BioBuilds / samtools 16.0			linux-64 linux-





### How does it work ?

• For a given software version, there are several archives corresponding to the different OS and CPU architectures for which the software is available.

bic	bioconda / packages / samtools							
Tools	for dealing	g with SAM	, BAM and CRA	M files				
C	onda	Files	Labels	Badges	]			
<b>▼</b> Fi	lters ype: All ~			Versio	on: 1.19.1 ~	Label: All ~		
	\$ Туре	\$ Size	<b>≑</b> Name			→ Uploaded	Downloads	Labels
	conda	465.5 kB	1 osx-64/san	ntools-1.19.1-ha	d510865_0.tar.bz2	🛗 1 month and 22 days ago	164	main
	conda	462.2 kB	Ilinux-64/sa	mtools-1.19.1-h	n50ea8bc_0.tar.bz2	🛗 1 month and 22 days ago	1018	main





### **Conda environments**

### With conda, each piece of software must be installed in an environment.

An environment is a folder containing all the files needed for the software to work.

This folder looks like a miniature operating system.

You can install multiple software in an environment but **only one version of a software in a given environment**.

You can create **as many environments as you want** each containing their own set of software.





### Creating a conda environment :

\$ conda create -n my\_env
Empty environment created at prefix: ~/.conda/envs/my\_env

### Activate a Conda environment :

To choose which conda environment you want to use, you have to activate it

```
$ conda activate my_env
(my_env) $
```





### Installing software





<pre>pipeline)\$ conda install -c b</pre>	ioconda f	astqc=0.12.1			
Updating specs:					
- fastqc=0.12.1					
Package	Version	Build	Channel	Size	
Install:					
<pre>+ _libgcc_mutex + _openmp_mutex + fastqc + font-ttf-dejavu-sans-mono + fontconfig + freetype + libexpat + libfreetype6 + libgcc + libgcc-ng + libgomp + libpng + libuuid + libxcrypt + libzlib + openjdk + perl + zlib Summary: Install: 19 packages Total download: 209MB</pre>	0.1 4.5 0.12.1 2.37 2.15.0 2.13.3 2.7.0 2.13.3 2.7.0 15.1.0 15.1.0 15.1.0 15.1.0 15.1.0 15.1.0 15.1.0 1.6.47 2.38.1 1.3.1 1.3.1	conda_forge 2_gnu hdf78af_0 hab24e00_0 h7230c49_1 ha770c72_1 h4886fc4_1 h767d61c_2 h69a702a_2 h767d61c_2 h943b412_0 h0b41bf4_0 hd590300_1 hb9d3cd8_2 h516909a_1016 7_hd590300_per15 hb9d3cd8_2	conda-forge conda-forge bioconda conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge conda-forge	Cached Cached	

Confirm changes: [Y/n]

Conda proposes a large list of packages to install.

These are all dependencies required to run fastqc.

For each package it indicates the source channel and the version it will deploy.

We can see that it has found fastqc in the bioconda channel but that most of the dependencies come from the conda-forge channel.

The **bioconda** channel offers tools specifically for bioinformatics.

The conda-forge channel offers genetic tools and libraries.





### **Benefits for Reproducibility**

- <u>Software catalog</u>: Anaconda offers a vast catalog of software and maintains all known versions for each software package.
- <u>Isolation</u>: Environments can be isolated to avoid conflicts between package versions.
- <u>Sharing Environments</u>: Environments can be shared using environment files (environment.yml), ensuring that others can replicate the exact setup.





### Conda is available in different flavors





A faster solver as been developed for conda : libmamba This library has been developed by a french company called **QuantStack** 

A new package manager based on libmamba is now available : mamba *It uses the same channels as conda.* 













#### mamba, the super powered solver







conda or mamba run natively on many operating systems without the need to install special system packages. They can therefore be used as a non-administrator user.

Similarly, the same installation of conda or mamba can run on several OS. This is the case on the Core cluster where some of the compute nodes are running Centos and others Ubuntu.

conda or mamba also support OS updates very well.



For the most light and simple experience, we recommend using **micromamba** 

It can be installed with one single command :

\$ "\${SHELL}" <(curl -L micro.mamba.pm/install.sh)</pre>

After installation, the micromamba command will be available on your computer.

### Same tool, different names :

If you rather install miniconda or anaconda, the command will be conda

If you choose to install mamba, the command will be mamba

micromamba, mamba and conda proposes all the same command line interface.



# A few words about encapsulation







### Let's say we want to install RStudio...

Denter Scherken all all a		3.mm
Studio	Poduts - Solutions - Co	samoo umar ses canaar Q domes Resoutes - Alaud - Prong
	Download the RStud	io IDE
		N17
Choose Your Ver	rsion	Studio Team
Choose Your Ver The Ritudio IDE is a set of integral productive with II and Pythem. It is	rsion ted tush designed to help yas be more visiden a connecter, spetia highlighting	Studio Team
Choose Your Ver The Ritudio IDC is a set of enegrat productive with R and Pythem. It is addlar that supports direct code to for platting, eleving history, difes	rsion within the large star is being year. He more relation a control, myrtas highlighting macrature, and a variety of robust tubits gang and managing pour sonitapore.	Studio Team
Choose Your Ver The Ritudio IDC is a set of integra productive with Yan Pythons 11 editor that supports divect code es for platting, viewing bistory, debu	rsion bet taalt designed to help yes be more ecidera a cermin, protein highlighting escature, and a variety of rebeat tools geing and managing your workspace.	Studio Team

#### MacOS



ncy tree

Done Done





studio-1.2.5001-and64.deb

rently installed.

a set of integrated tools designed to help you be more productive with R. It. re package? [v/N]:v 181189 files and directories //rstudio-1.2.5001-and64.deb

> nome-menus (3.13.3-llubuntul) ... Mesktop-file-utils (0.23-lubuntu3.18.04.1) Nime-support (3.60ubuntu1) ... olor-icon-theme (0.17-red-mime-info (1.9-2)

### **Use Rstudio**

A FILO	All safety		
Contraction Contraction			E Paget Dave
Cenaile	C Internet Holes		
R version 3.0.1 (2013-05-16) "Good 5 Copyright (C) 2013 The R Poundation for Platform: x86_64-apple-darwin10.8.0 (64	port" Statistical Computing -bit)	n 🖌	
R is free software and comes with ABSC Tow are welcome to redistribute it unde Type 'licensel' or 'licence'' for dis Matural language support but running R is a collaborative project with many Type 'contributors()' for more informat 'citation()' on how to cita R or R pack Type 'demo()' for some demus, 'help()'	UTELY NO BABBLARTY. r certain conditions. tribution decisia. in an English locale contributors. ion and ages in publications. for on-line help, or		
'help.stort()' for an HTML brokser inte Type 'q()' to quit R.	rface to help.	Inde	-
Discharges January From of Effettal	Constant C. See	in house i @ Herer	
functional connection of teneroly	( ) ( ) ( ) ( ) ( ) ( ) ( ) ( ) ( ) ( )		and the second se
1	at the	1118	Ann DA. DOLLA, AD. DA. War.
-1	0.00.0000	27.2.48	And DA THESE LT IN MAL
	O T Bandle	212 holes	100 10. 0013. 1:04 PM
	0 0 m		
	C. C. Destroy		
	in the second		
	C. C. Destant		
	C C Branker		
	C C fam		
	C Co fee cothe		
	C C Inclust		





We started with a computer using a specific OS...

Host OS

Computer







We started with a computer using a specific OS... And inside this environment, we installed a new application.





Usually dependencies of different applications don't interfere. But what if we want to test the latest version of our favourite tool? There might be conflicts...





Usually dependencies of different applications don't interfere. But what if we want to test the latest version of our favourite tool? There might be conflicts...





Idea : create separated environments for each application.





Idea : create separated environments for each application. More versatile: create a new environment per analysis.





But what if we want to install a software from a different OS?





Idea: use virtual machines Pros:

- Each application gets a completely different and independent environment
- Virtual machines can be transferred to another computer (using the same manager)









Idea: use virtual machines Pros: transferable independent environments Cons:

- Redundancy between VMs
- Heavy to set up
- No automation







•



Idea: "trick" applications into believing that they are in a different OS than the host's Avoid redundancy.







OS virtualisation vs hardware virtualisation Pros:

- Speed
  - Installation is faster
  - No boot time
- Lightweight
  - Minimal base OS
  - Minimal libraries and application set
- Easy sharing of applications





Cons:

 Singularity to use images on a cluster



### Encapsulation

RStudio v1	RStudio v1.2			
R Packages	R Packages			
Environment 1	Environment 2			
Conda				
Host OS				
Computer				



Application1	Application2			
Libraries	Libraries			
Guest OS 1	Guest OS 2			
VM manager				
Host OS				
Computer				



![](_page_42_Picture_5.jpeg)

.

# Introduction to docker

![](_page_43_Picture_1.jpeg)

![](_page_44_Picture_1.jpeg)

### What is Docker?

Docker is a platform designed to help developers build, share, and run container applications.

Containers allow a developer to package up an application with all the parts it needs, such as libraries and other dependencies, and ship it all out as one package.

It is often defined as lightweight virtualization.

![](_page_44_Picture_6.jpeg)

![](_page_45_Picture_1.jpeg)

### **Key Features**

- Docker uses containerization to ensure that applications run consistently across different computing environments.
- Pre-built images can be used to quickly set up environments.
- Registries are available for sharing and accessing Docker images. The most popular is the Docker Hub.

![](_page_45_Picture_6.jpeg)

![](_page_46_Picture_1.jpeg)

### How does it work ?

![](_page_46_Figure_3.jpeg)

(https://docs.docker.com/get-started/overview/)

![](_page_46_Picture_5.jpeg)

![](_page_47_Picture_1.jpeg)

### **Benefits for Reproducibility**

- <u>Consistency</u>: Ensures that the application runs the same way regardless of where it is deployed.
- <u>Isolation</u>: Containers provide process and filesystem isolation, which helps in avoiding conflicts.
- <u>Version Control</u>: Docker images can be versioned and shared, making it easy to replicate environments.

![](_page_47_Picture_6.jpeg)

![](_page_48_Picture_0.jpeg)

![](_page_48_Picture_1.jpeg)

![](_page_48_Picture_2.jpeg)

Docker requires a **Docker Host** to run containers.

The Docker Host is a system daemon that run as root and can access to a reserved part of the hardware resources.

This is not compatible with an HPC cluster where hardware resources are already managed by a job scheduler (SLURM)

![](_page_49_Figure_5.jpeg)

(https://docs.docker.com/get-started/overview/)

![](_page_49_Picture_7.jpeg)

![](_page_50_Picture_1.jpeg)

Apptainer is an open source container platform designed to run complex applications on high-performance computing (HPC) clusters in a simple, portable, and reproducible way.

### An Apptainer container image is **a file** An Apptainer running container is **a user process**

![](_page_50_Picture_4.jpeg)

![](_page_50_Picture_5.jpeg)

### **Key Features**

- Optimized for use in HPC environments where Docker might not be suitable.
- Designed with security in mind, allowing users to run containers without requiring root privileges.
- Can run Docker images, making it easy to transition from Docker.

![](_page_51_Picture_5.jpeg)

![](_page_52_Picture_1.jpeg)

### Running a Docker image on an HPC cluster with Apptainer

- # load apptainer
- \$ module load apptainer

![](_page_52_Picture_5.jpeg)

![](_page_53_Picture_1.jpeg)

### Running a Docker image on an HPC cluster with Apptainer

```
$ apptainer build cowsay.sif docker://rancher/cowsay
INFO:
        Starting build...
Copying blob 34d5e986f175 done
Copying blob dd05e66d8cea done
Copying blob cbdbe7a5bc2a done
Copying blob 13eefd6dff68 done
Copying config 223a921ebc done
Writing manifest to image destination
2025/06/05 22:39:16 info unpack layer: sha256:cbdbe7a5bc2a134ca8ec91be58565ec07d037386d1f1d8385
2025/06/05 22:39:16 info unpack layer: sha256:dd05e66d8ceaaf3d7cf6712075c8c28099aa003d0c625be8b
2025/06/05 22:39:17 info unpack layer: sha256:34d5e986f1757ac7e1b094d89b7ba9e40ee5b614fc91ad2a9
2025/06/05 22:39:17 info unpack layer: sha256:13eefd6dff6843c9689bb227b3eeb54c1f427c1397e029866
INFO: Creating SIF file...
        Build complete: cowsay.sif
INFO:
```

![](_page_53_Picture_4.jpeg)

![](_page_54_Picture_1.jpeg)

### Running a Docker image on an HPC cluster with Apptainer

![](_page_54_Figure_3.jpeg)

![](_page_54_Picture_4.jpeg)

### **Benefits for Reproducibility**

- <u>Portability</u>: Containers can be easily moved and run on different HPC systems.
- <u>Reproducibility</u>: Ensures that the computational environment is consistent and can be replicated.
- <u>Security</u>: Enhanced security features make it suitable for multi-user HPC environments.

![](_page_55_Picture_6.jpeg)

## A brief summary

![](_page_56_Figure_1.jpeg)

![](_page_56_Picture_2.jpeg)

ΤοοΙ	Primary Use Case	Key Strengths	Potential Weaknesses
Conda	Package and environment management	Easy to use, cross-platform, no root access required	Limited to package built by the community, less isolation compared to containers
Docker	Containerization	Highly portable, strong isolation, extensive image repository	Requires root access, not ideal for HPC environments
Apptainer	Containerization for HPC	Secure, no root access required, compatible with Docker images	Less user-friendly for non-HPC use cases

![](_page_57_Picture_3.jpeg)

![](_page_58_Picture_1.jpeg)

![](_page_58_Picture_2.jpeg)

![](_page_58_Picture_3.jpeg)

![](_page_59_Picture_0.jpeg)

# Thank you for your attention !

![](_page_59_Picture_2.jpeg)

![](_page_59_Picture_3.jpeg)

![](_page_59_Picture_4.jpeg)

![](_page_59_Picture_5.jpeg)