# Making Bioinformatics Accessible

Andrew Wheeler

11/27/2024

# Roots for Resilience Program


Arizona Institute for Resilience


RESEARCH, INNOVATION & IMPACT
Data Science Institute


CYVERSE®

- Fellowship and outreach program
- Build data science skills
- Connect people across disciplines

# Open Science Principles

- Accessibility
  - Data management
  - Code availability
  - Documentation
- Reproducibility
  - Software Environments
  - Containers
  - Pipelines

# Data Management

- Can users find and access your data?
- Can they understand, use, and reuse it consistently?
- Are you ethically protecting data where needed?
- Data management plans

- Version controlled code development platform
  - Make code available to public
  - Track changes
  - Work with team
- README
  - Where to find everything
  - How to run code

# GitHub

- Basic Commands
  - clone: copy a repository locally
  - pull: update local repository
  - branch: a version history. Multiple parallel branches can be active
  - fork: copy of someone else's repository stored on your account
  - commit: finalize a change
  - push: add change back to remote repository
  - merge: apply changes from a branch or fork to the main branch
  - pull request: submit changes to be added t repository
  - issue: flag suggestions or tasks

# Reproducibility

- Can users run your code?
- Will they get the same outcome as you if they do?
- Can users apply your code to their own projects easily?



**Reproducibility Spectrum**

Publication only | Publication + (Code · Code and data · Linked and executable code and data) | Full replication
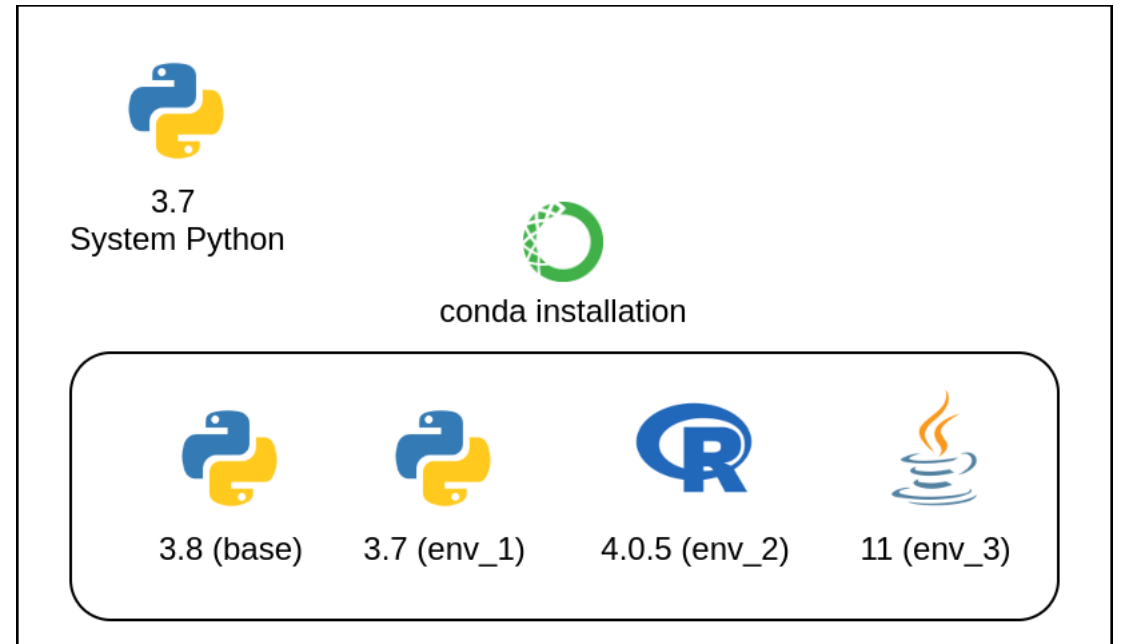
Not reproducible ← → Gold standard

# Computing Environments

- All the hardware, software, and resources you are using
  - Hardware: CPUs, GPUs, RAM
  - Operating System: Windows, Mac, Linux
  - Software Versions: R, Python, etc.
  - Packages and Package Versions: specific software packages
- Software Dependency Hell:
  - Incorrect versions
  - Missing dependencies
  - Obsolete code

# Environment Managers

Environments can be exported and shared so all users have the same versions

- Conda – most popular
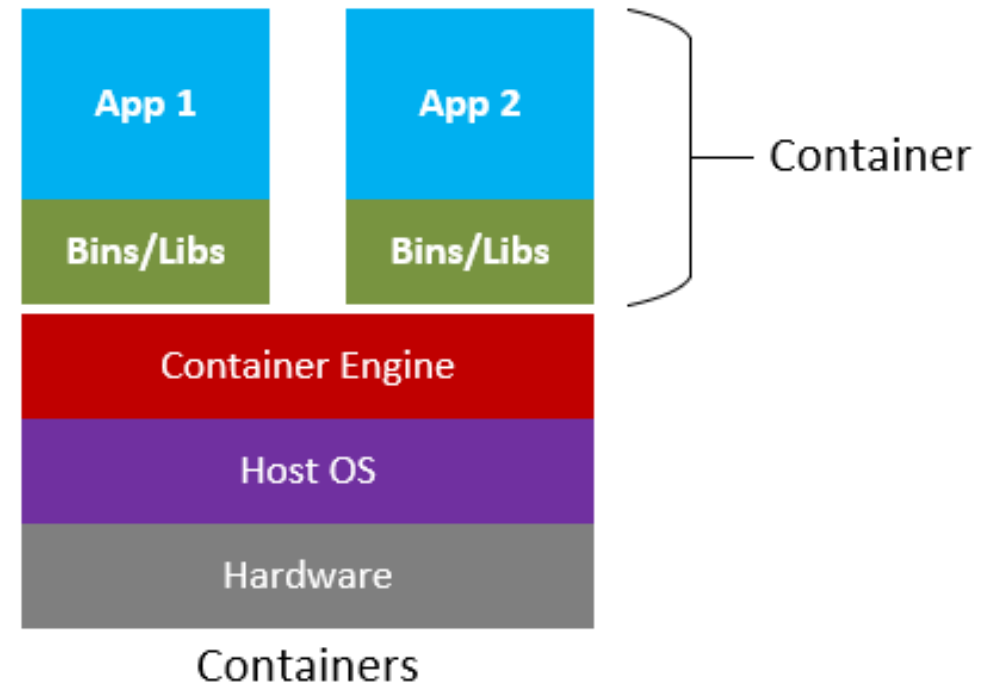- Mamba – implemented in C, faster than Conda
- Pip
- Renv



Export: conda env export > my_conda_env.yml

Reproduce: conda env create --file environment.yml
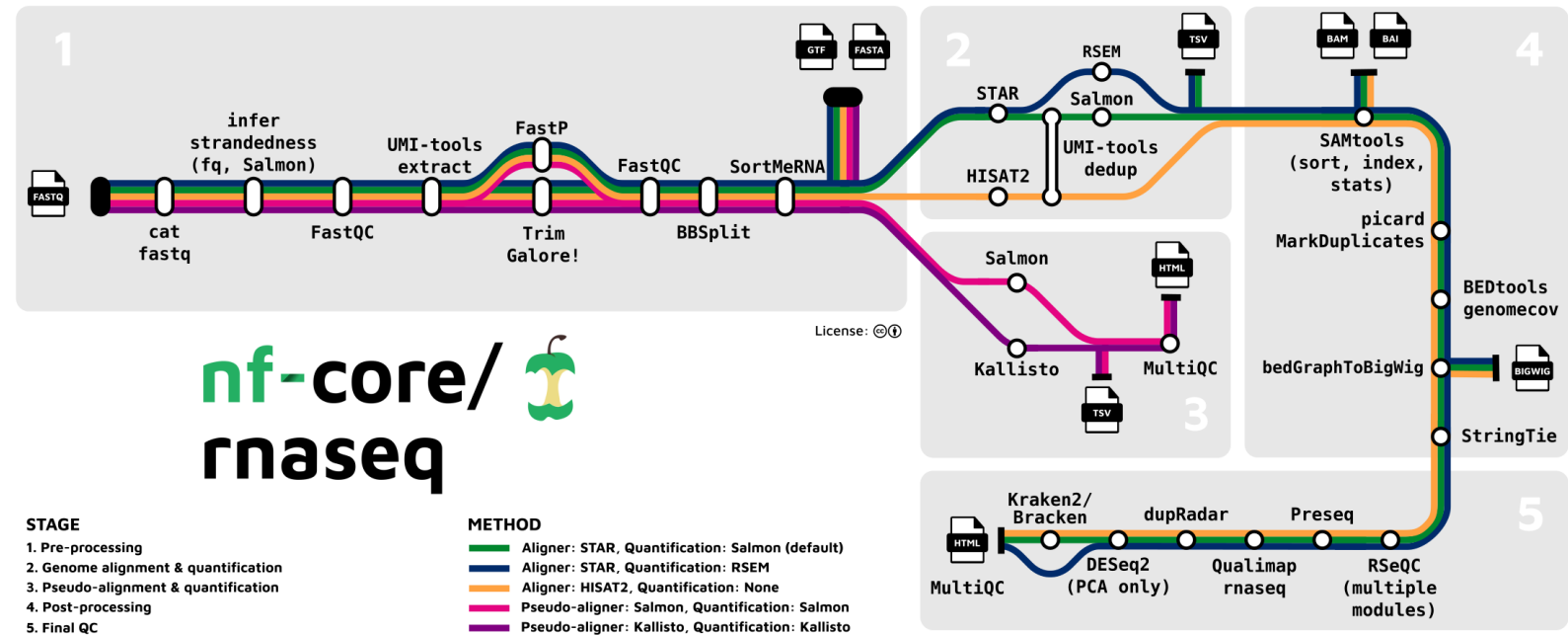
# Containers

- Contains everything you need to run the code in a single unit
  - Easy to share
  - Can run on any machine
  - Isolated from the rest of the computer
- Docker – most popular container management software



Containers

# Pipelines

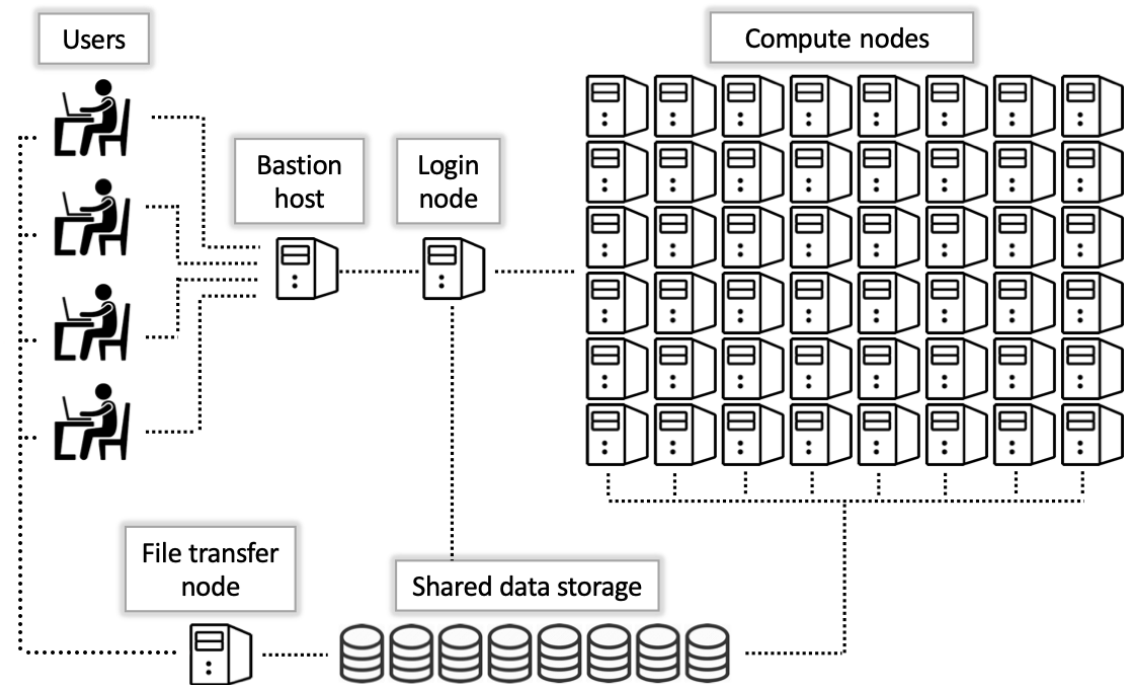Workflow managers can let you standardize complex, multistep analyses

- NextFlow
- Snakemake

# Remote Computing

Working beyond your local machine

- HPC – high performance computing
- CyVerse – cloud computing service hosted by UA

- Computing clusters
  - Large number of machines linked together
  - Power from parallelization

# Summary

How can we make apply open science principles to bioinformatics?

- Well documented, easily accessible code and data on github or other cloud services

- Software environments to avoid dependency hell

- Containers and Pipeline managers to handle larger projects

- Making use of available remote computing resources

# FOSS Course

Foundational Open Science Skills
https://foss.cyverse.org/

Consider attending in future terms if you want to learn more!

# Thanks!