# **Overview**

1    Intro to HPCs
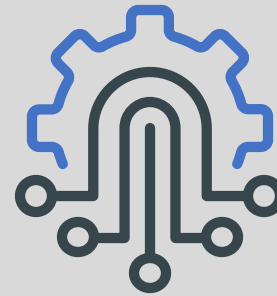
2    UCloud platform & setup

3    Pipelines, workflows & nf-core

HeaDS

# HPC

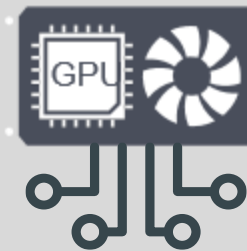## What is high performance computing (HPC)?

Using a supercomputer or a cluster of computers to perform jobs too computationally intensive for a personal computer (laptop or desktop)
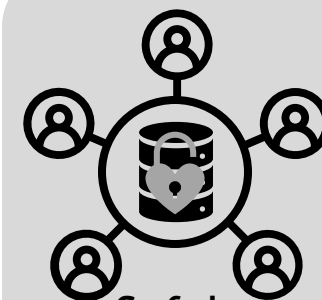


Compute-heavy tasks

RAM-heavy tasks

GPU-based tasks

Safely shared tasks

HeaDS

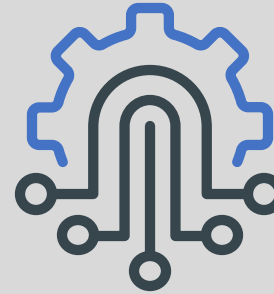# When would an HPC platform help you in biotech / health data science?

# Tasks that benefit from HPC

**Compute-heavy**: many easily broken down tasks to be run – sequential or parallel

**RAM-heavy:** large datasets need to be processed in close proximity – i.e. genome alignment

**GPU-based**: large ML models (neural networks) and multi-dim models (biomechanics, weather)
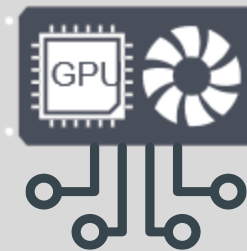
**Safely shared**: you're working with sensitive data and/or many users need to access the same private datasets in parallel
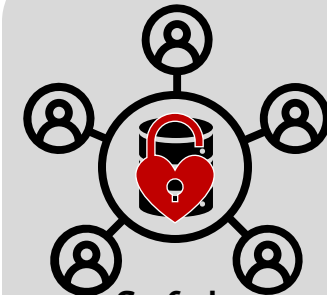
Compute-heavy tasks

RAM-heavy tasks

GPU-based tasks

Safely shared tasks
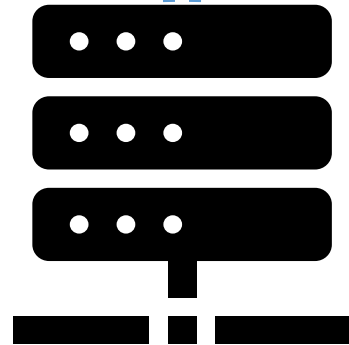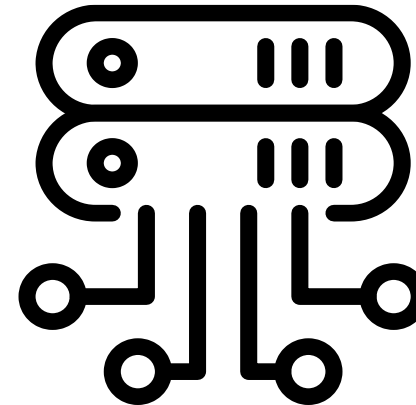
HeaDS

**HPC**

**An HPC poll – computing experience?**

In-house mgmt → ← Research IT

computerome

Microsoft Azure
amazon web services™

Juiced up lab PC

Local server

Uni HPC platform

Commercial cloud

HeaDS

# HPC

**Important considerations**

In-house mgmt    Research IT

computerome

Microsoft Azure

amazon web services™

| Juiced up lab PC | Local server | Uni HPC platform | Commercial cloud |

**AUTONOMY**

**HPC**

**Important considerations**

In-house mgmt ← → Research IT

computerome

Microsoft Azure

amazon web services™

Juiced up lab PC

Local server

Uni HPC platform

Commercial cloud

AUTONOMY                    SCALABILITY

**HPC**

**Important considerations**

In-house mgmt · Research IT

computerome

Microsoft Azure

amazon web services™

Juiced up lab PC

Local server

Uni HPC platform

Commercial cloud

AUTONOMY · SECURITY* · SCALABILITY

*GDPR

# Components of an HPC

**HPC CLUSTER ARCHITECTURE**

User

Login or Head Node

COMPUTING CLUSTER

Interconnected Computing Nodes

PARALLEL FILE SYSTEM

Storage System

https://www.weka.io/learn/hpc/what-are-hpc-and-hpc-clusters/

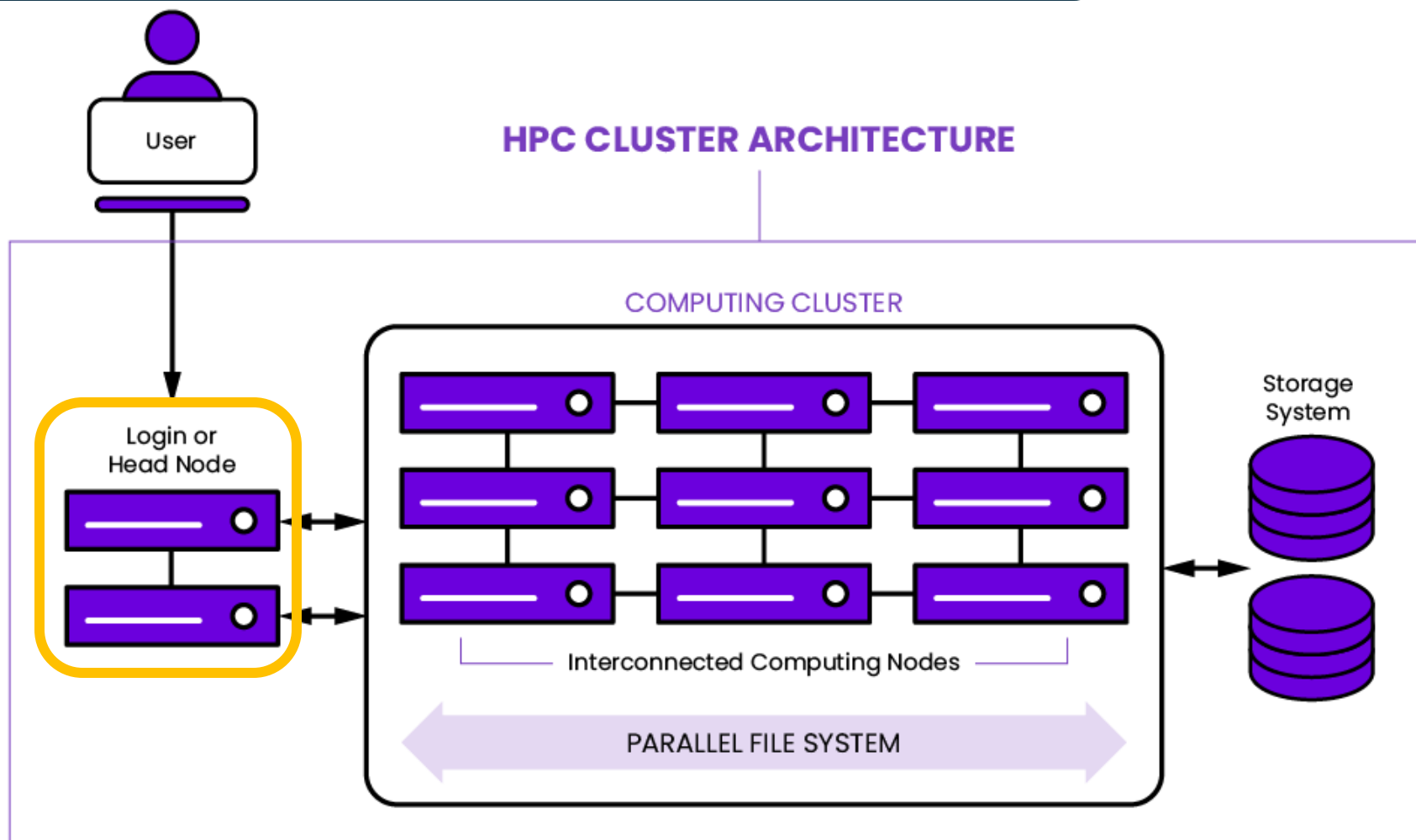**Nodes**
- ➢ Any physical device that can send, receive, or pass information
- ➢ In HPC, term usually references a compute node or a login node
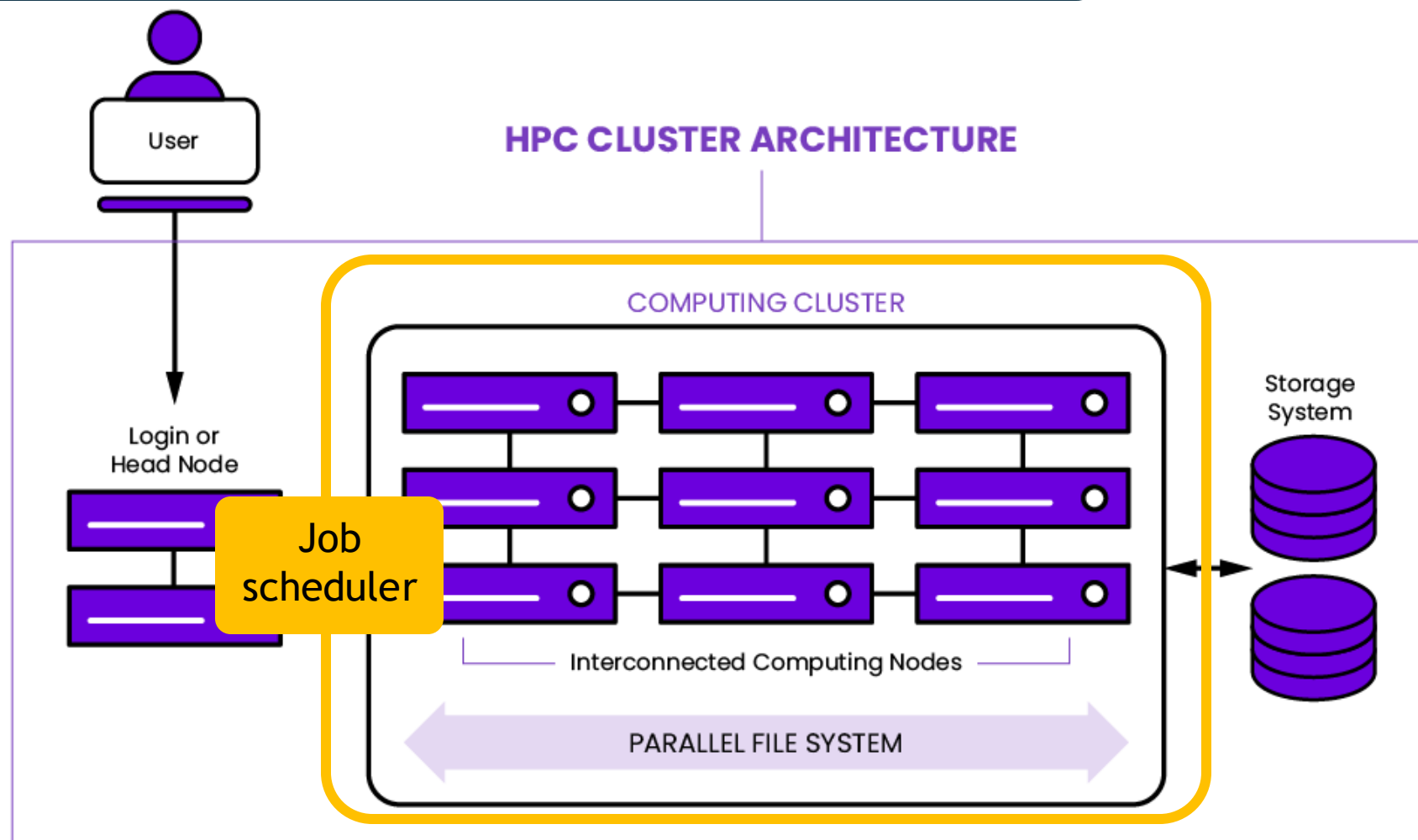- ➢ Nodes are interconnected to facilitate communication and data transfer between them

HeaDS

# Components of an HPC



HPC CLUSTER ARCHITECTURE

COMPUTING CLUSTER

Login or Head Node

Interconnected Computing Nodes

PARALLEL FILE SYSTEM

Storage System

User

**Login Node**

➢ A computer that acts as the front end to the HPC system, where users access (request) cluster resources and submit tasks for the computing nodes to perform
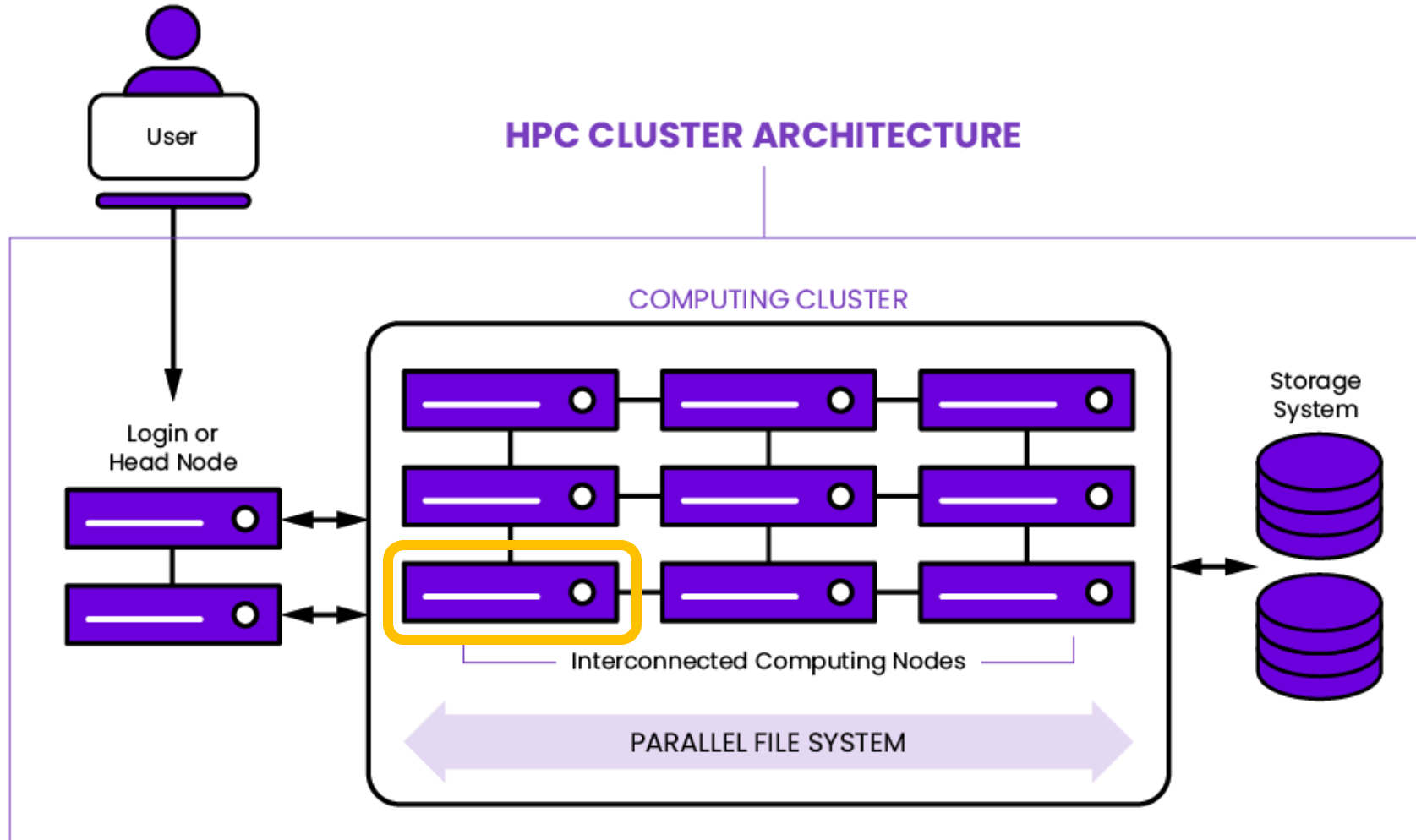
https://www.weka.io/learn/hpc/what-are-hpc-and-hpc-clusters/

HeaDS

# Components of an HPC



HPC CLUSTER ARCHITECTURE

User

Login or Head Node

Job scheduler

COMPUTING CLUSTER

Interconnected Computing Nodes

PARALLEL FILE SYSTEM

Storage System

**Computing Cluster**

➤ A (large) group of closely interconnected computers that work together as a single system to complete jobs

➤ Job scheduler manages and schedules jobs (tasks) across compute nodes allocating resources to complete the job

HeaDS

# Components of an HPC



HPC CLUSTER ARCHITECTURE

User

Login or Head Node

COMPUTING CLUSTER

Interconnected Computing Nodes

PARALLEL FILE SYSTEM

Storage System

**Computing Node**

➢ an individual computer within the compute cluster made up of a set of processors and their local memory

➢ 'Size' of the node traditionally varies by number of processors and amount of memory

HeaDS

# Components of an HPC

**Storage system**
➢ Provides persistent storage for data and programs used by the HPC system



HPC CLUSTER ARCHITECTURE

User

Login or Head Node

COMPUTING CLUSTER

Interconnected Computing Nodes

PARALLEL FILE SYSTEM

Storage System

HeaDS

# HPC interfaces you might encounter



Terminal /
command line

Virtual Machine
(usually linux OS)

Custom Graphical
User Interface

HeaDS

# HPC interfaces you might encounter



Terminal /
command line



Virtual Machine
(usually linux OS)



Custom Graphical
User Interface

# A standard HPC workflow

Access login node → Transfer data to platform storage → Configure tools → Set up scripts / test scaling → Write job & submit to scheduler → After queuing, monitor job run → Check output / revise

HeaDS

# A standard HPC workflow

| Access login node | → | Transfer data to platform storage | → | Configure tools | → | Set up scripts / test scaling | → | Write job & submit to scheduler | → | After queuing, monitor job run | → | Check output / revise |

```
[bartell@fe-open-01 HPC-Pipes]$ cat SBATCH_ex.sh
#!/bin/bash
#SBATCH --account HDSSandbox
#SBATCH --cpus-per-task=2
#SBATCH --mem 12g
#SBATCH --time 03:00:00
#SBATCH --output=std.out
#SBATCH --error=std.err

#activate environment
eval "$(conda shell.bash hook)"
conda activate HPCpipes_workshop

echo hello world
[bartell@fe-open-01 HPC-Pipes]$ sbatch SBATCH_ex.sh
```

HeaDS

# A standard HPC workfl[ow]

Access login node → Transfer data to platform storage → Configure tools → Set up / test s[...]

```
[bartell@fe-open-01 HPC-Pipes]$ cat SBATCH_ex.sh
#!/bin/bash
#SBATCH --account HDSSandbox
#SBATCH --cpus-per-task=2
#SBATCH --mem 12g
#SBATCH --time 03:00:00
#SBATCH --output=std.out
#SBATCH --error=std.err

#activate environment
eval "$(conda shell.bash hook)"
conda activate HPCpipes_workshop

echo hello world
[bartell@fe-open-01 HPC-Pipes]$ sbatch SBATCH_ex.sh
```

HeaDS

---

## Coder ⭐

Base ⌄   1.93.1 ⌄

⤢ Documentation      Sandbox_workshop ⌄

Run Visual Studio Code on UCloud and access it through your browser. For more information, check here.

↕ Import parameters      ▶ Submit

E-mail notification settings

Do not notify me ⌄

Estimated cost          6 Core-hours
Current balance     23,16K Core-hours

**Job name**
test_JAB

**Hours** *
3      +1   +8   +24

**Machine type** *

u1-standard-2

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 2 (Intel Xeon Gold 6130) | 12 | None | 2 Core-hours/hour |

**Select folders to use**      Add folder

Your files will be available at /work/.

Remove ✕
/Member Files: kcs305kcs305#7929/work_JAB

Remove ✕
/shared/HPCLab_workshop

**Additional Parameters**

Initialization                                    Remove ✕
/shared/HPCLab_workshop/setup.sh

*Run a Bash script (*.sh) for initialization.*

**Optional Parameters**                     Search

Modules path                                      Use

**Configure SSH access**

This application has optional support for SSH. In order to use SSH access, you must configure at least one SSH key. You can configure your SSH keys here.

☑ Enable SSH server

# **UCloud**

- UCloud is a danish High Performance Computing environment
  - Lots of storage, lots of cpus and RAM (computing power)

- Danish institutions have access to it
  - You personally have 1000dkk in computing resources

- UCloud works in apps, giving you access to different programs
  - All apps have documentation on how to use them!

- This means everyone is using the same versions of software
  - Makes teaching much much easier as results are reproducible

HeaDS

# UCloud access

https://hds-sandbox.github.io/bulk_RNAseq_course/develop/

# UCloud log-in

To access *UCloud* please choose your login provider

**SDU**

University of Copenhagen

☐ Always use the login provider that I choose now. At my.wayf.dk I can res... use a different login provider.

Search here 🔍

**1. Search for your uni & then click on link**

**2. Sign-in via your uni portal**

https://cloud.sdu.dk/app/login

...cholar | KU library login | HeaDS SharePoint | hds-sandbox Share... | Sandbox Planner

**DeiC**

Integration Portal

WAYF    Login

Other login options →

**UCloud**

# UCloud log-in

Back to the Info Nov '24 page...

## Access Sandbox resources

Our first choice is to provide all the **training materials, tutorials, and tools as interactive apps on UCloud**, the supercomputer located at the University of Southern Denmark. Anyone using these resources needs the following:

1. a Danish university ID so you can sign on to UCloud via WAYF[1].

> for UCloud Access click here

2. basic ability to navigate in Linux/RStudio/Jupyter. **You don't need to be an expert**, but it is beyond our ambitions (and course material) to teach you how to code from zero and how to run analyses simultaneously. We recommend a basic R or Python course before diving in.

3. **For workshop participants:** Use our invite link to the correct UCloud workspace that will be shared on the day of the workshop. This way, we can provide you with compute resources for the active sessions of the workshop[2] Click the link below after your first uCloud access and accept the invite that shows.

> Invite link to uCloud workspace

d67-3c40-4ccb-8a52-1f289c7e3df0

Search files and applications...

You have been invited to join Sandbox RNAseq workshop

Join project    Ignore

# 10 min break!

# Workspaces



Virtual workspaces allow you to share resources and work together with project collaborators

# UCloud usage

During this course, utilize **Sandbox RNAseq workshop workspace** (resources have been requested for this purpose)

Following the workshop, switch to "My workspace

# Drives

Project folders, files, etc. that only belong to the **active workspace** will be accessible from the menu at the left

↓

**Drives**

⊞ | ⇵ Create drive ⌥ Q

👤 View member files ☐

🔍 🔄 **Sandbox RNAseq works...** ⌄

Your username: should be FirstLast#0000...

| Drive name | Provider | Created by | Created at |
|---|---|---|---|
| Member Files: AlbaRefoyoMartínez#0753 | SDU/K8 | AlbaRefoyoMartínez#0753 | 15:22 05/02/2024 |
| sandbox_bulkRNAseq | SDU/K8 | JoseAlejandroHerreraRomer... | 10:08 08/08/2022 |
| sequencing_data | SDU/K8 | JoseAlejandroHerreraRomer... | 10:37 17/05/2023 |

You have a personal drive

You have shared drives

# Drives

Access file structure (shared and generated in previous jobs)



- Personal workspace folder "Member Files:username": results will go here
  - Jobs folder
    - Subfolders with Apps names
    - App name: All runs (the job's name) results
- "sandbox_bulkRNAseq" : contains some course material for teachers
- "sequencing_data" : contains fastq files for preprocessing (nf-core RNAseq)

HeaDS

# **Applications**

There is a wide variety of applications.

Here are some of my favorites!

# Apps

Search for Sandbox apps

# Submitting a job with a Sandbox app

1. **App & version (dropdown menu to change it)**

2. Read documentation before using it

3. Import parameters (if wanted from a previous job)

4. Job name, hours, and machine type (resources set-up)

5. Folders to access while running this particular job

6. Module to use (which includes Notebooks & Data)

7. ▶ Submit



**Transcriptomics Sandbox** ⭐
Default ⌄    2024.06 ⌄    **version 2024.06**

⧉ Documentation

Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

⇅ Import parameters    ▶ Submit

E-mail notification settings
Do not notify me ⌄

Estimated cost         1 Core-hours
Current balance        19,53K Core-hours

Job name
test_JAB

Hours *
1    +1  +8  +24

Machine type *
u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour ⌄ |

Select folders to use                    Add folder

Your files will be available at /work/.

Remove ✕
/Member Files: kcs305kcs305#7929/work_JAB

Mandatory Parameters

Select a module *
Introduction to bulk RNAseq analysis in R    ⌄

# Submitting a job
# with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** ▶ Submit

# Submitting a job with a Sandbox app

1. App & version (dropdown menu to change it)

2. Read documentation before using it

3. **Import parameters (if wanted from a previous job**

4. Job name, hours, and machine type (resources set-up)

5. Folders to access while running this particular job

6. Module to use (which includes Notebooks & Data)

7. ▶ Submit

---

## Transcriptomics Sandbox ⭐

Default ⌄    2024.06 ⌄

⧉ Documentation

Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

⇅ Import parameters    ▶ Submit

Estimated cost        1 Core-hours
Current balance       19,53K Core-hours

**E-mail notification settings**

Do not notify me ⌄

### Job name
test_JAB

**Hours** *
1    +1  +8  +24

### Machine type *
u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|---|---|---|---|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour |

### Select folders to use
Add folder

Your files will be available at /work/.

Remove ✕

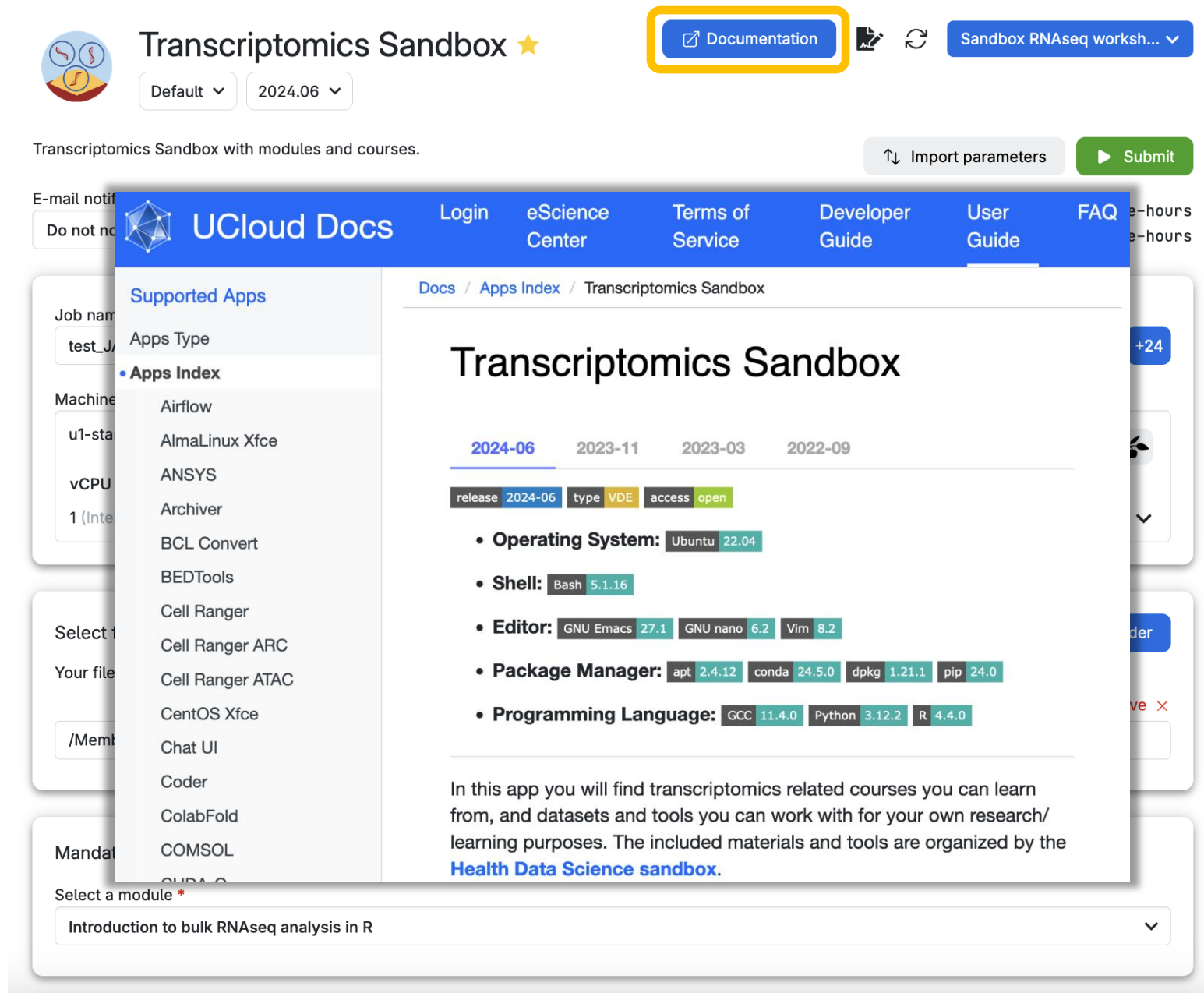/Member Files: kcs305kcs305#7929/work_JAB

### Mandatory Parameters

Select a module *

Introduction to bulk RNAseq analysis in R

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** **Import parameters (if wanted from a previous job**

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** ▶ Submit



## Transcriptomics Sandbox ⭐

Default ⌄    2024.06 ⌄

🔗 Documentation

Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

⇅ Import parameters    ▶ Submit

Estimated cost    1 Core-hours
Current balance   19,53K Core-hours

E-mail notification settings
Do not notify me ⌄

### Jobs

🔍  🔄  Sandbox RNAseq worksh... ⌄

⬆ Upload JobParameters.json    📄 Select file from UCloud

📅 Created ⌄    ○ Status ⌄    👤 Created by

test_JAB    12:58    Import

test_JAB    12:01    Import

Add folder

Select folders to use

Your files will be available at /work/.

Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB

### Mandatory Parameters

Select a module *
Introduction to bulk RNAseq analysis in R    ⌄

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** ▶ Submit

---

### Transcriptomics Sandbox ⭐
Default ⌄   2024.06 ⌄

⤢ Documentation

Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

⇅ Import parameters   ▶ Submit

**E-mail notification settings**

Do not notify me   ⌄

Estimated cost   1 Core-hours
Current balance   19,53K Core-hours

---

**Job name**

test_JAB

**Use your initials / a unique name!**

**Hours** *
1 ⌄   +1   +8   +24

**Machine type** *

u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour ⌄ |

---

**Select folders to use**                                   Add folder

Your files will be available at /work/.

                                                            Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB

---

**Mandatory Parameters**

Select a module *

Introduction to bulk RNAseq analysis in R   ⌄

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** ▶ Submit



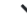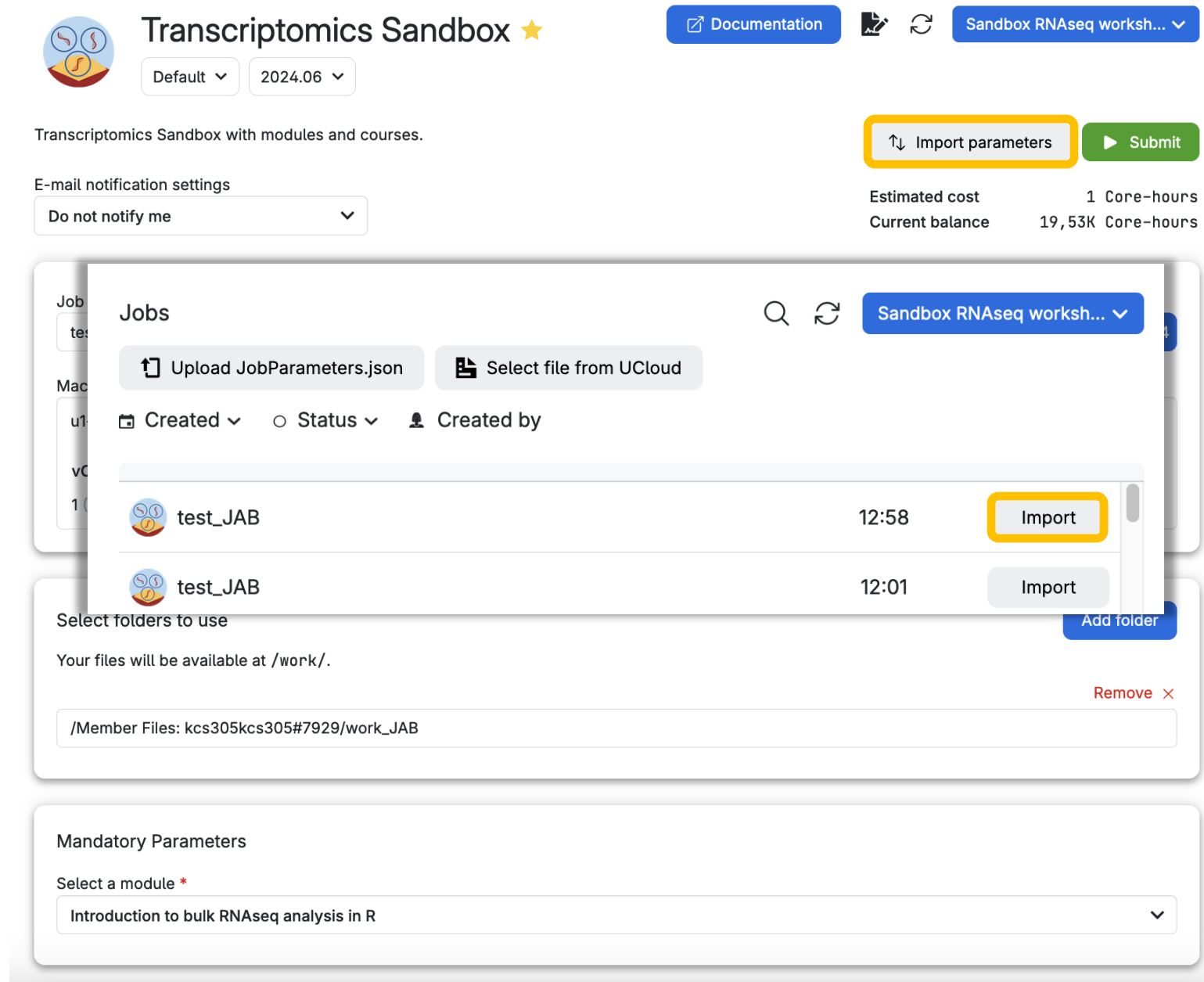## Transcriptomics Sandbox ★

Default ⌄  2024.06 ⌄

⧉ Documentation

Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

**E-mail notification settings**

Do not notify me ⌄

⇅ Import parameters    ▶ Submit

Estimated cost          1 Core-hours
Current balance     19,53K Core-hours

**Job name**

test_JAB

**Machine type** *

u1-standard-1

| vCPU | Mer |
|------|-----|
| 1 (Intel Xeon Gold 6130) | 6 |

**Hours** *

1 ⌄  +1  +8  +24

**Price**

1 Core-hours/hour ⌄

**Today, ask for an hour (for testing)**

**Billing is in hour increments, and you can ask for extra hours while your job is running**

**Select folders to use**

Your files will be available at /work/.

Add folder

Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB

**Mandatory Parameters**

Select a module *

Introduction to bulk RNAseq analysis in R ⌄

# Submitting a job with a Sandbox app

1. App & version (dropdown menu to change it)
2. Read documentation before using it
3. Import parameters (if wanted from a previous job)
4. **Job name, hours, and machine type (resources set-up)**
5. Folders to access while running this particular job
6. Module to use (which includes Notebooks & Data)
7. ▶ Submit



Transcriptomics Sandbox ⭐
Default ▾   2024.06 ▾

⎘ Documentation

Sandbox RNAseq worksh... ▾

Transcriptomics Sandbox with modules and courses.

⇅ Import parameters   ▶ Submit

Estimated cost        1 Core-hours
Current balance    19,53K Core-hours

E-mail notification settings
Do not notify me ▾

Job name
test_JAB

Hours *
1 ⇕   +1  +8  +24

Machine type *
u1-standard-1

Search machine types...

| Name | vCP... | | | |
|---|---|---|---|---|
| 🐇 u1-standard-1 | 1 (In... | | | |
| 🐇 u1-standard-2 | 2 (In... | | | |
| 🐇 u1-standard-4 | 4 (Intel Xeon Gold 6130) | 24 | None | 4 Core-hours/hour |
| 🐇 u1-standard-8 | 8 (Intel Xeon Gold 6130) | 48 | None | 8 Core-hours/hour |
| 🐇 u1-standard-16 | 16 (Intel Xeon Gold 6130) | 96 | None | 16 Core-hours/hour |
| 🐇 u1-standard-32 | 32 (Intel Xeon Gold 6130) | 192 | None | 32 Core-hours/hour |
| 🐇 u1-standard-64 | 64 (Intel Xeon Gold 6130) | 384 | None | 64 Core-hours/hour |

Introduction to bulk RNAseq analysis in R ▾

**1 core is sufficient for the data analysis we'll do together**

**(You would need more cores to analyze raw data via nf-core pipeline)**

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** ▶ Submit

---

## Transcriptomics Sandbox ⭐

Default ⌄   2024.06 ⌄

⬈ Documentation     Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

↑↓ Import parameters     ▶ Submit

E-mail notification settings

Do not notify me ⌄

Estimated cost          1 Core-hours
Current balance     19,53K Core-hours

Job name                                                                Hours *
test_JAB                                                        1 ⇕   +1  +8  +24

Machine type *
u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour |

Select folders to use                                        Add folder

Your files will be available at /work/.

Remove

/Member Files: kcs305kcs305#7929/work_JAB

**Add the custom working directory we made earlier**

Mandatory Parameters

Select a module *
Introduction to bulk RNAseq analysis in R ⌄

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** ▶ Submit



**Transcriptomics Sandbox** ⭐

Default ⌄    2024.06 ⌄

Transcriptomics Sandbox with modules and courses.

◻ Documentation    Sandbox RNAseq worksh... ⌄

⇅ Import parameters    ▶ Submit

E-mail notification settings

Do not notify me    ⌄

Estimated cost    1 Core-hours
Current balance    19,53K Core-hours

Job name

test_JAB

Hours *

1    +1  +8  +24

Machine type *

u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour |

Select folders to use    Add folder

Your files will be available at /work/.

Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB

Mandatory Parameters

Select a module *

Introduction to bulk RNAseq analysis in R

**The app contains a few modules – this fits the current workshop and will load the necessary tools and notebooks**

# Submitting a job with a Sandbox app

For those with some HPC experience...

Do these steps look familiar?

Perhaps similar to a job's bash script that you submit using a workload manager like SLURM or PBS?



## Transcriptomics Sandbox ⭐

Default ⌄   2024.06 ⌄

📄 Documentation    Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

⇅ Import parameters    ▶ Submit

Estimated cost        1 Core-hours
Current balance       19,53K Core-hours

E-mail notification settings
Do not notify me    ⌄

Job name
test_JAB

Hours *
1    ⌃⌄    +1   +8   +24

Machine type *
u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour |

Select folders to use                                    Add folder

Your files will be available at /work/.

                                                         Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB

Mandatory Parameters

Select a module *
Introduction to bulk RNAseq analysis in R    ⌄

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** Review & Submit



**Transcriptomics Sandbox** ⭐

Default ⌄   2024.06 ⌄

🗗 Documentation    📝    🔄    Sandbox RNAseq worksh... ⌄

Transcriptomics Sandbox with modules and courses.

↕ Import parameters    ▶ Submit

Estimated cost        1 Core-hours
Current balance       19,53K Core-hours

**E-mail notification settings**

Do not notify me                                          ⌄

---

**Job name**

test_JAB

**Hours** *

1    ⇅    +1    +8    +24

**Machine type** *

u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|------|-------------|-----|-------|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour |

---

**Select folders to use**                                Add folder

Your files will be available at /work/.

Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB

---

**Mandatory Parameters**

**Select a module** *

Introduction to bulk RNAseq analysis in R                        ⌄

# Jobs

**Test is now running** (ID: 5049428)

🗑 Stop application

**Hold to stop job**

⧉ Open terminal   ⧉ Open interface

🕐 **Time allocation**

**Job start:** 13:41 17/05/2024
**Job expiry:** 14:41 17/05/2024
**Time remaining:** 00:57:01
Extend allocation (hours):

+1   +8   +24

💬 **Messages**

[13:40] AlbaRefoyoMartínez#0753 has requested 1x u1-standard-1 from
DeiC Interactive HPC (SDU)
[13:40] Assigned to nodeaa-05
[13:40] Job is starting soon
[13:41] Job has started

**Before time remaining is over,
you can add extra hours**

🖴 **Node 1**

extracting: /w_____.gz
 inflating: /work/Intro_to_bulkRNAseq/Data/salmon/control_2/aux_info/observed_bias_3p.gz

**If your 'Open interface button' doesn't go dark blue after you get a message that your 'Job has started', then hit refresh (browser)!**

**Do the same on the new tab that pops up if it just spins, too**

HeaDS

# Jobs



Our apps are made to be used by any UCloud user with their own compute resources, so module materials (data & notebooks) are downloaded fresh with each app run

Do you have another familiar folder here? >>>

HeaDS

# Jobs

Our apps are made to be used by any UCloud user with their own compute resources, so module materials (data & notebooks) are downloaded fresh with each app run

We want to edit and save our notebooks over the course, so…
**we're going to copy this directory somewhere writeable**

```
R    File    Edit    Code    View    Plots    Session    Build    De

Console    Terminal ×    Background Jobs ×

R    R 4.4.0 · /work/

R version 4.4.0 (2024-04-24) -- "Puppy Cup"
Copyright (C) 2024 The R Foundation for Statistical Comp
Platform: x86_64-pc-linux-gnu

R is free software and comes with ABSOLUTELY NO WARRANTY
You are welcome to redistribute it under certain conditi
Type 'license()' or 'licence()' for distribution details

  Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

>
```

Files    Plots    Packages    Help    Viewer    Presentation

/ > work

| Name | Size | Modified |
|---|---|---|
| Intro_to_bulkRNAseq | | |
| JobParameters.json | 1.1 KB | Nov 13, 2024, 3: |

# Jobs

**Test is now running** (ID: 5049428)

[⊟ Open terminal]  [↗ Open interface]

🗑 **Stop application**

**Stop the app by holding down this button**

## 🕐 Time allocation

**Job start:** 13:41 17/05/2024
**Job expiry:** 14:41 17/05/2024
**Time remaining:** 00:57:01
Extend allocation (hours):

[+1] [+8] [+24]

## 💬 Mes...

```
[13:4...
DeiC Interactive HPC (SDU)
[13:40] Assigned to nodeaa-05
[13:40] Job is starting soon
[13:41] Job has started
```

## 🖴 Node 1

```
extracting: /work/Intro_to_bulkRNAseq/Data/salmon/control_2/aux_info/exp_gc.gz
  inflating: /work/Intro_to_bulkRNAseq/Data/salmon/control_2/aux_info/observed_bias_3p.gz
```

HeaDS

50

# Jobs

1. Go into Transcriptomics Sandbox and pick the job we just ran

2. Click ONCE to select Intro_to_bulkRNAseq

3. Click 'Copy to...'

4. Click on your root member files directory

5. Copy to the custom working directory you made

**Now you will always have your own, writeable directory when using the app**



HeaDS

# Jobs

1. Go into Transcriptomics Sandbox and pick the job we just ran

2. Click ONCE to select Intro_to_bulkRNAseq

3. Click 'Copy to...'

4. Click on your root member files directory

5. Copy to the custom working directory you made

**Now you will always have your own, writeable directory when using the app**



/Member Files: kcs305kcs305#.../.../Transcriptomics Sandbox/test_JAB (5117705) **1**

Launch with... ⌥ O    Rename ⌥ F    Copy to... ⌥ C **3**    Move to... ⌥ M

Show hidden files ☐

**Name**

stdout.txt

JobParameters.json

Intro_to_bulkRNAseq **2**

If you add **Member.../Intro_to_bulkRNAseq** as a folder in your job setup, you will get your edited files

If you don't, you'll get clean files from the app

/Member Files: kcs305kcs305#... **4**    Sandbox RNAseq worksh... ∨

✓ Use this folder    ⤧ Create folder

Show hidden files ☐

**Name**    **Modified at**

work_JAB    P    11:49    Copy to **5**

work    P    11/01/2024    Copy to

Trash    P    11:38    Copy to

HeaDS

53

# Submitting a job with a Sandbox app

**1** App & version (dropdown menu to change it)

**2** Read documentation before using it

**3** Import parameters (if wanted from a previous job)

**4** Job name, hours, and machine type (resources set-up)

**5** Folders to access while running this particular job

**6** Module to use (which includes Notebooks & Data)

**7** Review & Submit

---

## Transcriptomics Sandbox ⭐

Default ⌄   2024.06 ⌄

Transcriptomics Sandbox with modules and courses.

⧉ Documentation       ↻       Sandbox RNAseq worksh... ⌄

⇅ Import parameters      ▶ Submit

Estimated cost         1 Core-hours
Current balance       19,51K Core-hours

**E-mail notification settings**

Do not notify me ⌄

| Job name | Hours * |
|---|---|
| test_JAB | 4 ⇅  +1  +8  +24 |

**Machine type** *

u1-standard-1

| vCPU | Memory (GB) | GPU | Price |
|---|---|---|---|
| 1 (Intel Xeon Gold 6130) | 6 | None | 1 Core-hours/hour |

**Select folders to use**        Add folder

Your files will be available at /work/.

Remove ✕

/Member Files: kcs305kcs305#7929/work_JAB/Intro_to_bulkRNAseq

**Mandatory Parameters**

Select a module *

Introduction to bulk RNAseq analysis in R ⌄

# Sandbox RStudio-based app

In the transcriptomics app, you will get RStudio in your web browser.

Click on the *Intro_to_bulkRNAseq* folder and navigate through the R Markdown notebooks to run the analyses.

HeaDS

# UCloud jobs

**Access past and current running jobs**

**A** Click on a job to stop it or rerun (ensuring the use of the same parameters)

or

**B** Access old jobs from 'Import parameters' button on job setup page



HeaDS

# 15 min break!

HeaDS

# Why use pipelines/ workflows

A pipeline (workflow) is a series of programmatic steps to transform raw data into processed results, figures, and insights.

# Why use pipelines/ workflows

Each step involving **different tools, parameters, reference databases, and specific requirements.**



**Let me do this by hand via single tool calls in the terminal...**

HeaDS

# Why use pipelines/ workflows

**Now apply the same analysis to new data...** 😣

Dataset 2 INPUT

Dataset 3 INPUT

Dataset 4 INPUT

Dataset 5 INPUT

Dataset N INPUT

Dataset 1 FASTQ

ALIGNED READS

Quality control → Trimming → Aligning → Indexing

Intermediate files

Intermediate files

Intermediate files

HeaDS

# Wait, should I use a different mapping tool? Which one?



BWA

BOWTIE2

VG

MOSAIK

Dataset 2
INPUT

Dataset 3
INPUT

Dataset 4
INPUT

Dataset 5
INPUT

Dataset N
INPUT

Dataset 1
FASTQ

ALIGNED
READS

Quality
control

Trimming

?

Indexing

Intermediate
Intermediate
Intermediate
Intermediate
Intermediate
files

Intermediate
Intermediate
Intermediate
Intermediate
Intermediate
files

Intermediate
Intermediate
Intermediate
Intermediate
Intermediate
files

HeaDS

# Why use pipelines/ workflows

**SOFTWARE**

- **Multiple softwares** required, and sometimes even more **parameters** to tweak

- Small changes in the parameters software can cause a large difference in the results

- Differences in **program resource** needs at each step (compute power, data inputs, software dependencies, etc.)

**DATA**

- Many files are being generated (also intermediate files) and the size of the data files can be large

- Differences in **data type, shape and scale**

**Bioinformatics workflows are complex... and reproducibility can be very challenging**

HeaDS

# Workflow managers (languages)

Based on Groovy (Java dialect)

Pipeline order controlled by channel flow (allows diff outputs, FIFO)

Maintained by venture-funded co.

Native cloud support ($) +++ many plug-ins

Standardized portable modules

X Excessive temp dirs / symlinks

HeaDS

# WfM languages



Based on Groovy (Java dialect)

Pipeline order controlled by channel flow (allows diff outputs, FIFO)

Maintained by venture-funded co.

Native cloud support ($) +++ many plug-ins

Standardized portable modules

X Excessive temp dirs / symlinks

```
params.samples = 'data/samples/*.fastq'
params.outdir = 'output'
params.genome = 'data/genome.fa'

process FASTP {

    publishDir "${params.outdir}/QC", pattern: "*.html"

    input:
    tuple val(id), path(reads)

    output:
    path '*.json', emit: json
    path '*.html', emit: html

    """
    fastp -i ${reads} \ --json ${id}.fastp.json \   --html ${id}.fastp.html
    """
}

process MINIMAP2 {

    cpus 2

    input:
    tuple val(id), path(reads)
    path genome

    output:
    tuple val(id), path("*.sam")

    """
    minimap2 -t ${task.cpus} \
```

HeaDS

# WfM languages

nextflow

Based on Groovy (Java dialect)

Pipeline order controlled by channel flow (allows diff outputs, FIFO)

Maintained by venture-funded co.

Native cloud support ($) +++ many plug-ins

Standardized portable modules

X Excessive temp dirs / symlinks

HeaDS

```
params.samples = 'data/samp
params.outdir = 'output'
params.genome = 'data/geno

process FASTP {

    publishDir "${params.outdi

    input:
    tuple val(id), path(reads)

    output:
    path '*.json', emit: json
    path '*.html', emit: html

    """
    fastp -i ${reads} \ --json $
    """
}


process MINIMAP2 {

    cpus 2

    input:
    tuple val(id), path(reads)
    path genome

    output:
    tuple val(id), path("*.sam"

    """
    minimap2 -t ${task.cpus} \
```

```
workflow {

    // Get our samples into a channel
    ch_samples = channel.fromPath(
params.samples )
        | map { [ it.simpleName, it ] }

    // Invoke fastp and put output into a
channel
    ch_fastp = ch_samples | FASTP

    ch_genome = channel.value(
file(params.genome, checkIfExists: true) )

    ch_flagstat = MINIMAP2( ch_samples,
ch_genome )
        | SAMTOOLS_VIEW
        | SAMTOOLS_FLAGSTAT


    ch_files = ch_flagstat
        | map { it[1] }
        | mix( ch_fastp.json )
        | collect

    MULTIQC ( ch_files,
file("${workflow.projectDir}/assets/multiq
c_config.yml", checkIfExists: true) )


}
```

# Workflow managers

- WfMS are software tools designed to **automate** and **optimize complex data analysis workflows.**

- WfMS allow users to **define, run, and track tasks and dependencies** in a modular, scalable way, simplifying the **management and reproducibility** of analyses.

- Why do we want to use them?

- Using WfMS helps **reduce errors and inconsistencies** in analyses, while enhancing the efficiency and reproducibility of research.

HeaDS

# Workflow managers

Execution without manual intervention

**Choose the pipeline that does the work for you! (or define them yourself)**

FASTQ

ALIGNED READS

Quality control — Trimming — Aligning — Indexing

Specify input files

WORKFLOW

Generate output files

HeaDS

**nf-core**

https://nf-co.re/pipelines

A community effort to collect a curated set of analysis pipelines built using Nextflow.

**For facilities**
Highly optimised pipelines with excellent reporting.
Validated releases ensure reproducibility.

**For users**
Portable, documented and easy to use workflows.
Pipelines that you can trust.

**For developers**
Companion templates and tools help to validate your code and simplify common tasks.

We will have a look at the results folder in a sec!

Big community, routinely maintained/updated.
Slack channel for questions & discussions (get helped and advised!)

Use their template!

HeaDS

[https://nf-co.re/pipelines](https://nf-co.re/pipelines)

- Bioinformatics community for curated pipelines written in **nextflow**
  - **Bulk RNAseq**
  - **Single Cell RNAseq**
  - **ATACseq**
  - **ChIPseq**
  - **HICseq**

- Completely reproducible, following gold standards (best practices) and open source

- Easy to implement (packaged software). They work on any computational infrastructures

- Very well documented

- More and more pipelines are being introduced and updated

**HeaDS**

# Pipelines on UCloud

Search for nf-core apps

# Your problem:

To perform quality control and quantify the expression of genes in a genome (bulk RNA sequencing)

# The solution:

## nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

It performs quality control (QC), trimming and (pseudo-)alignment, and produces a gene expression matrix and extensive QC report

HeaDS

## All you need is a samplesheet and FASTQ files!

```
> cat samplesheet.csv

sample,fastq_1,fastq_2,strandedness
CONTROL_REP1,AEG588A1_S1_L002_R1_001.fastq.gz,AEG588A1_S1_L002_R2_001.fastq.gz,auto
CONTROL_REP1,AEG588A1_S1_L003_R1_001.fastq.gz,AEG588A1_S1_L003_R2_001.fastq.gz,auto
CONTROL_REP1,AEG588A1_S1_L004_R1_001.fastq.gz,AEG588A1_S1_L004_R2_001.fastq.gz,auto
```

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

## All you need is a samplesheet and FASTQ files!

STAGE
1. Pre-processing
2. Genome alignment & quantification
3. Pseudo-alignment & quantification
4. Post-processing
5. Final QC

METHOD
— Aligner: STAR, Quantification: Salmon (default)
— Aligner: STAR, Quantification: RSEM
— Aligner: HISAT2, Quantification: None
— Pseudo-aligner: Salmon, Quantification: Salmon

License: ©①

1. Merge re-sequenced FastQ files (cat)
2. Sub-sample FastQ files and auto-infer strandedness (fq, Salmon)
3. Read QC (FastQC)
4. UMI extraction (UMI-tools)
5. Adapter and quality trimming (Trim Galore!)
6. Removal of genome contaminants (BBSplit)
7. Removal of ribosomal RNA (SortMeRNA)
8. Choice of multiple alignment and quantification routes:
    1. STAR -> Salmon
    2. STAR -> RSEM
    3. HiSAT2 -> **NO QUANTIFICATION**
9. Sort and index alignments (SAMtools)
10. UMI-based deduplication (UMI-tools)
11. Duplicate read marking (picard MarkDuplicates)
12. Transcript assembly and quantification (StringTie)
13. Create bigWig coverage files (BEDTools, bedGraphToBigWig)
14. Extensive quality control:
    1. RSeQC
    2. Qualimap
    3. dupRadar
    4. Preseq
    5. DESeq2
15. Pseudoalignment and quantification (Salmon or 'Kallisto'; *optional*)
16. Present QC for raw read, alignment, gene biotype, sample similarity, and strand-specificity checks (MultiQC, R)

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

## The nf-core/rnaseq pipeline will perform:

1. **Pre-processing**: trim and clean your reads

2. **Alignment and quantification of expression levels**
   - Alignment **post-processing**: mark duplicates

3. **Pseudoalignment and quantification**: align your reads and create count matrix

4. Postprocessing
   - Create BigWigs (coverage files)

5. **QC (for all the steps)**

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

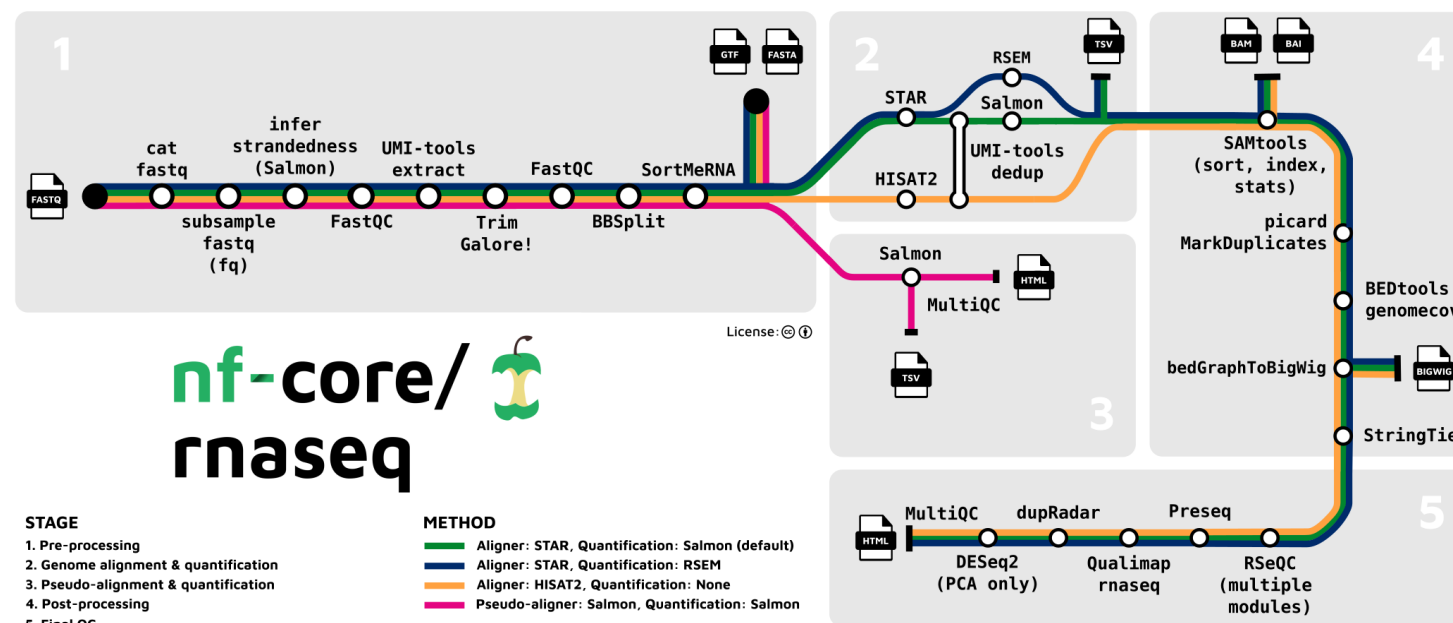## This takes some time, so we have run the pipeline for you!!

You have looked at some of the output files already

1. **Pre-processing**: trim and clean your reads

2. **Alignment and quantification of expression levels**
   - Alignment **post-processing**: mark duplicates

3. **Pseudoalignment and quantification**: align your reads and create count matrix

4. Postprocessing
   - Create BigWigs (coverage files)

5. **QC (for all the steps)**



License: ⓒ ⓘ

nf-core/ rnaseq

**STAGE**
1. Pre-processing
2. Genome alignment & quantification
3. Pseudo-alignment & quantification
4. Post-processing
5. Final QC

**METHOD**
— Aligner: STAR, Quantification: Salmon (default)
— Aligner: STAR, Quantification: RSEM
— Aligner: HISAT2, Quantification: None
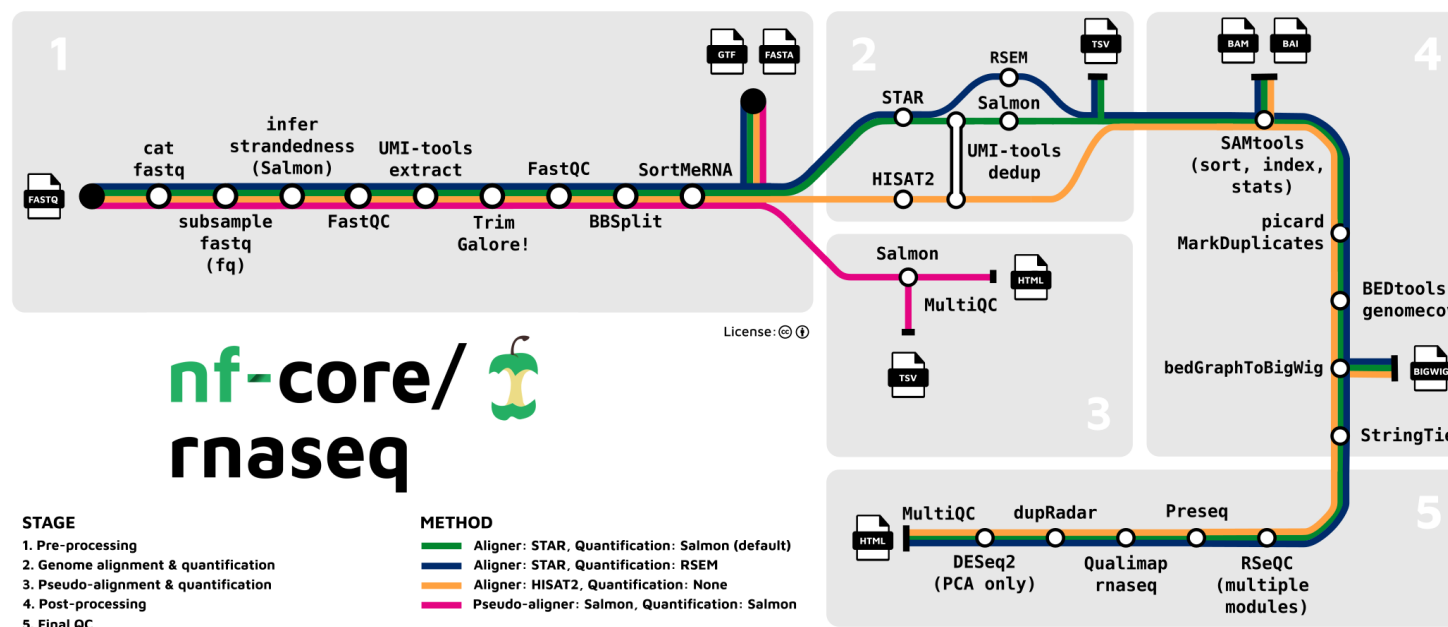— Pseudo-aligner: Salmon, Quantification: Salmon

HeaDS

# Pipelines and workflows

The pipeline generates:

- Results (e.g. count matrix, QC plots...)
- Other output files and reports:
  - Checkpoints, logs and progress

⬇

Why workflow reporting files?
- troubleshoot errors with the running
- For every launch commands -> run times and resource usage.

**nextflow**

```
RNASEQ:PREPARE_GENOME:GTF_GENE_FILTER                        -
RNASEQ:PREPARE_GENOME:MAKE_TRANSCRIPTS_FASTA                 -
RNASEQ:PREPARE_GENOME:CUSTOM_GETCHROMSIZES                  -
RNASEQ:PREPARE_GENOME:SALMON_INDEX                          -
RNASEQ:INPUT_CHECK:SAMPLESHEET_CHECK (samplesheet.csv) [100%] 1 of 1
RNASEQ:CAT_FASTQ                                            -
RNASEQ:FASTQC_UMITOOLS_TRIMGALORE:FASTQC (Control_1)    [  0%] 0 of 8
RNASEQ:FASTQC_UMITOOLS_TRIMGALORE:TRIMGALORE            [  0%] 0 of 8
RNASEQ:QUANTIFY_SALMON:SALMON_QUANT                         -
RNASEQ:QUANTIFY_SALMON:SALMON_TX2GENE                       -
RNASEQ:QUANTIFY_SALMON:SALMON_TXIMPORT                      -
RNASEQ:QUANTIFY_SALMON:SALMON_SE_GENE                       -
RNASEQ:QUANTIFY_SALMON:SALMON_SE_GENE_LENGTH_SCALED         -
RNASEQ:QUANTIFY_SALMON:SALMON_SE_GENE_SCALED                -
RNASEQ:QUANTIFY_SALMON:SALMON_SE_TRANSCRIPT                 -
RNASEQ:CUSTOM_DUMPSOFTWAREVERSIONS                          -
RNASEQ:MULTIQC                                             -
```
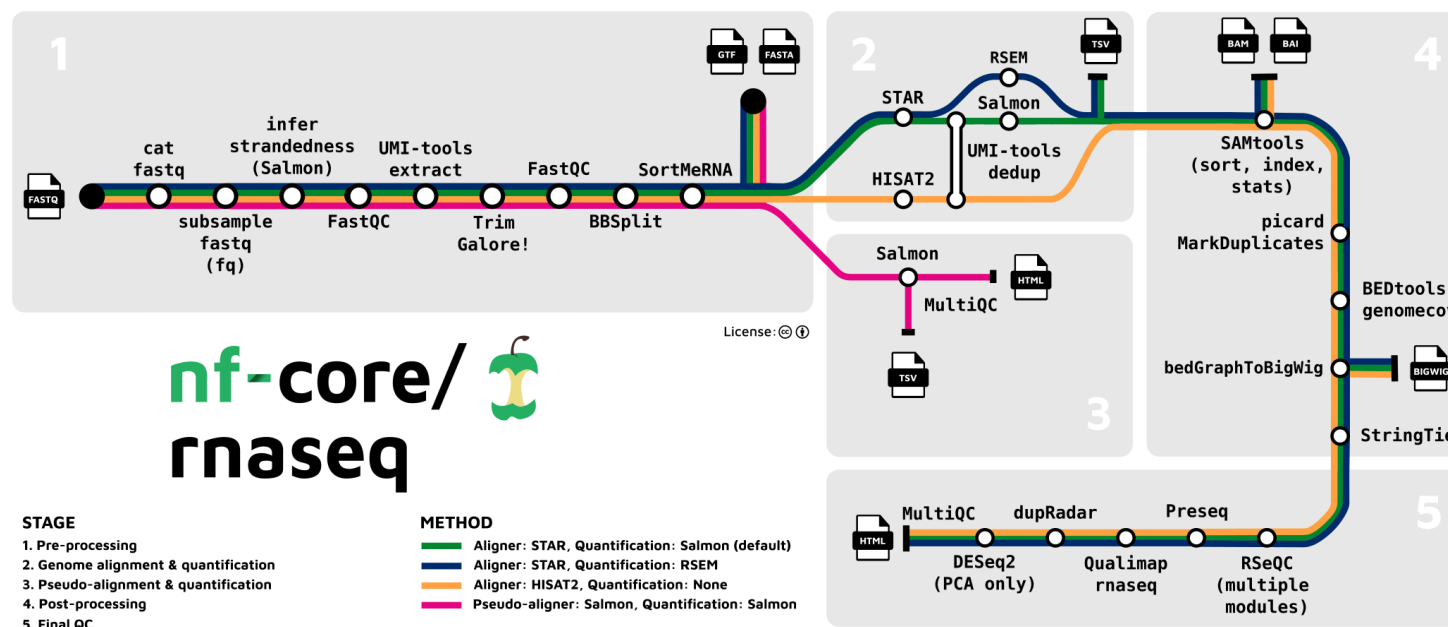
HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

All pipelines are run like this (test example):

```
nextflow run nf-core/rnaseq -r 3.17.0 -profile test --outdir <OUTDIR>
```

Arguments from **Nextflow** start with `-`
-r: pipeline version
-profile: docker, conda, etc
-resume: restart failed job

Arguments from **nf-core** start with `--`
--input: csv file with sample and read metadata
--outdir: results folder
--genome: reference genome to use
--aligner: select aligner
--skip_<X>: skip X process from pipeline

https://nf-co.re/rnaseq/3.17.0/parameters

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

After starting the workflow, you will see this:

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

Let's check out the results folder:

| | |
|---|---|
| 📁 fastqc/ | Raw reads quality control |
| 📁 multiqc/ | Full quality control report |
| 📁 pipeline_info/ | Pipeline information |
| 📁 salmon/ | Results from salmon pseudoaligner |
| 📁 star_salmon/ | Results from STAR aligner and quantification with salmon |
| 📁 trimgalore/ | Trimming and cleaning of reads + fastqc |

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.
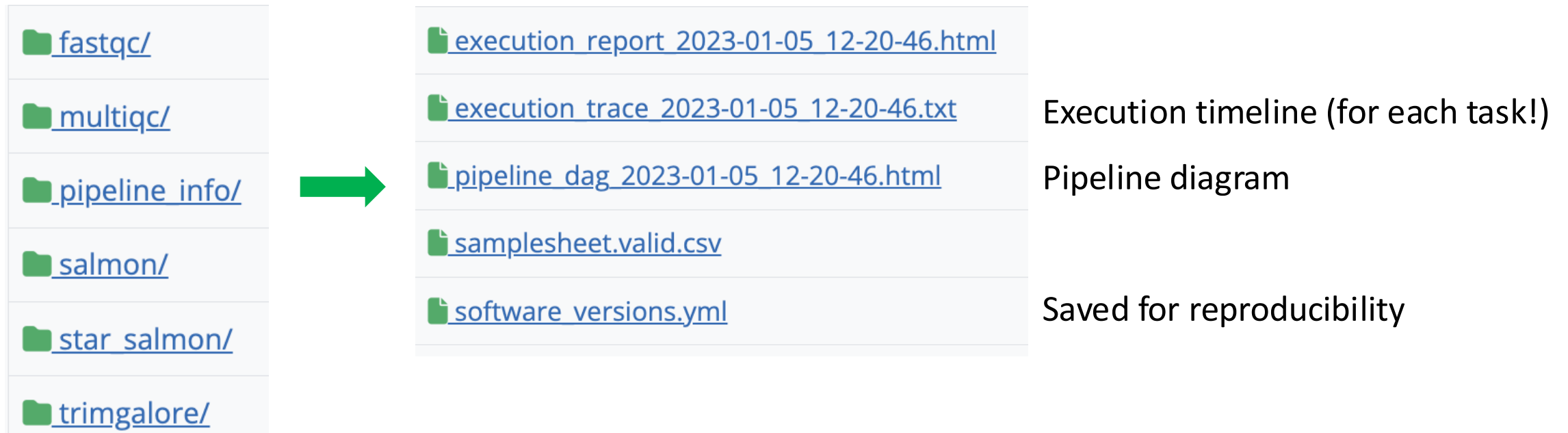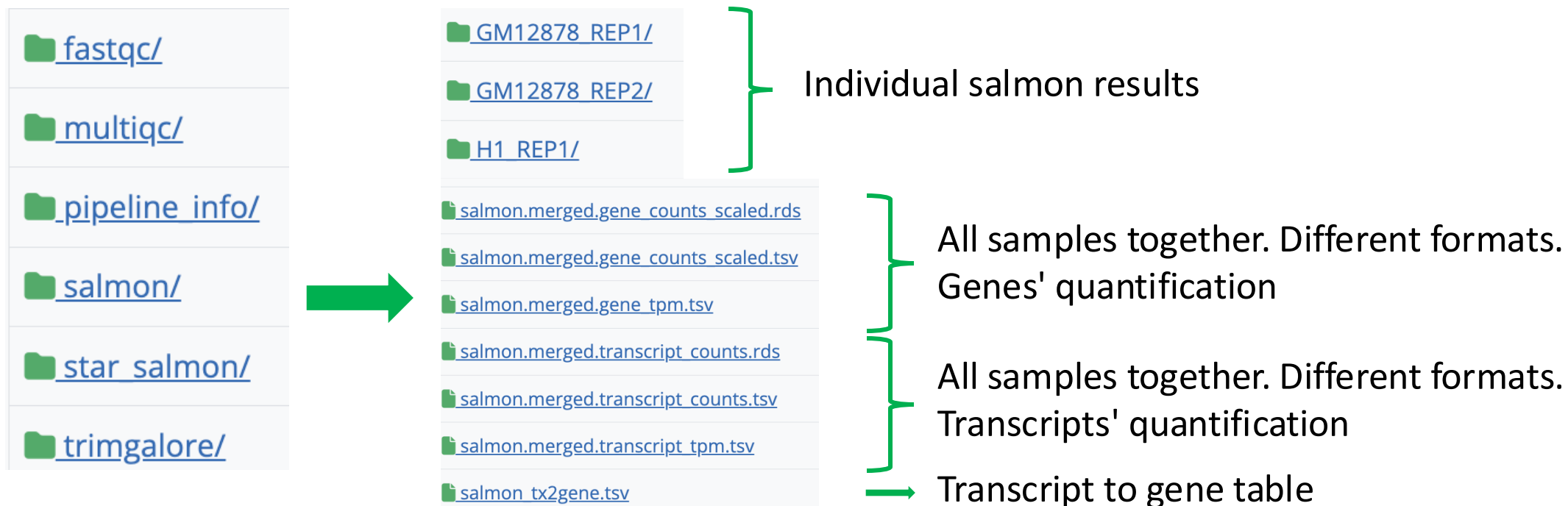
Let's check out the results folder: **fastqc**

📁 fastqc/

📁 multiqc/

📁 pipeline_info/

📁 salmon/

📁 star_salmon/

📁 trimgalore/

➡️

📄 GM12878_REP1_1_fastqc.html

📄 GM12878_REP1_1_fastqc.zip

📄 GM12878_REP1_2_fastqc.html

📄 GM12878_REP1_2_fastqc.zip

📄 GM12878_REP2_1_fastqc.html

📄 GM12878_REP2_1_fastqc.zip

📄 GM12878_REP2_2_fastqc.html

- Individual fastqc reports for raw reads
- Also in zip form

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

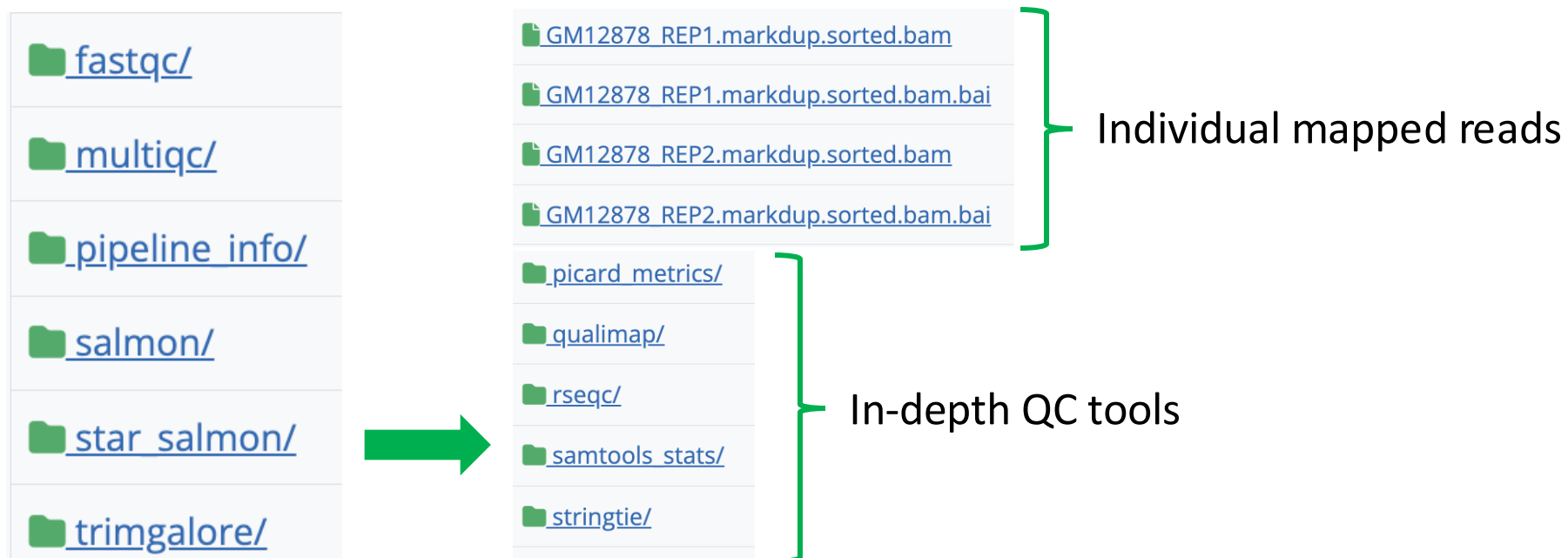Let's check out the results folder: **multiqc**

📁 fastqc/

📁 multiqc/ ➡️

📁 pipeline_info/

📁 salmon/

📁 star_salmon/

📁 trimgalore/

📁 multiqc_data/     Multiqc report data

📁 multiqc_plots/     Multiqc individual plots

📄 multiqc_report.html     Full MultiQC report

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

Let's check out the results folder: **pipeline_info**

📁 fastqc/

📁 multiqc/

📁 pipeline_info/  ➡️

📁 salmon/

📁 star_salmon/

📁 trimgalore/

📄 execution_report_2023-01-05_12-20-46.html

📄 execution_trace_2023-01-05_12-20-46.txt    Execution timeline (for each task!)

📄 pipeline_dag_2023-01-05_12-20-46.html    Pipeline diagram

📄 samplesheet.valid.csv

📄 software_versions.yml    Saved for reproducibility

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

## Let's check out the results folder: **salmon**

📁 fastqc/

📁 multiqc/

📁 pipeline_info/

📁 salmon/

📁 star_salmon/

📁 trimgalore/

📁 GM12878_REP1/

📁 GM12878_REP2/

📁 H1_REP1/

Individual salmon results

📄 salmon.merged.gene_counts_scaled.rds

📄 salmon.merged.gene_counts_scaled.tsv

📄 salmon.merged.gene_tpm.tsv

All samples together. Different formats. Genes' quantification

📄 salmon.merged.transcript_counts.rds

📄 salmon.merged.transcript_counts.tsv

📄 salmon.merged.transcript_tpm.tsv

All samples together. Different formats. Transcripts' quantification

📄 salmon_tx2gene.tsv

Transcript to gene table

HeaDS

# nf-core/rnaseq

RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/isoform counts and extensive quality control.

Let's check out the results folder: **star_salmon**

📁 fastqc/

📁 multiqc/

📁 pipeline_info/

📁 salmon/

📁 star_salmon/ ➡

📁 trimgalore/

📄 GM12878_REP1.markdup.sorted.bam

📄 GM12878_REP1.markdup.sorted.bam.bai

📄 GM12878_REP2.markdup.sorted.bam

📄 GM12878_REP2.markdup.sorted.bam.bai

Individual mapped reads

📁 picard_metrics/

📁 qualimap/

📁 rseqc/

📁 samtools_stats/

📁 stringtie/

In-depth QC tools

star_salmon also contains same results as salmon

HeaDS

# Run nf-core on UCloud

# UCloud Files

Salmon multiqc report

In **Files:**

- *Member Files:username:* your personal space
  - Work results will be here

  `Path: Member Files:username/nf-core rnaseq/<runName>/results_salmon`

- *sandbox_bulkRNAseq*: contains some course material for teachers

- *sequencing_data*: contains fastq files for preprocessing **(Don't try to modify... write-protected!)**

  `Path: sequencing_data/preprocessing_results_salmon/results_salmon`

# sequencing_data (778339) > preprocessing_results_salmon > results_salmon/multiqc > multiqc_report.html

# More pipeline help

**Git & Github**
- Code management
- Version control

**Bash & Unix**
- Operating from the terminal

**HPC-Launch**
- Omics data management
- Using DK HPCs

**HPC-Pipes**
- Software envs
- Pipeline management

> git annex

**VERSION CONTROL**

git

**DATA**

Input  COOKIECUTTER  Output

**COMPUTATIONS**

R  C++  python  snakemake  nextflow

**ENVIRONMENT**

Env 1  Env 2  Env 3

Mamba  docker

**HPC**

slurm workload manager

HeaDS