



Introduction to the Workflow Management System -- Nextflow & nf-core

Yanni Chen

*Graduate Research Assistant
High Performance Computing Center*

Fall 2022



Before Workshop

- **Check your conda environment**
 - [Get yourself start with conda](#)
 - [Manage your conda environment](#)
- [Pre-workshop survey](#)



Target Audience

This session is intended for an audience who ...

- currently runs, plans to run workflow / pipeline / series analysis
- has to analyze data using different platform (personal computer, lab cluster, HPC, or cloud computing)
- tries to improve reproducibility of data analysis
- writes proposals for workflow development
- is simply curious about workflow management



Purposes of Training

Purposes of this training are ...

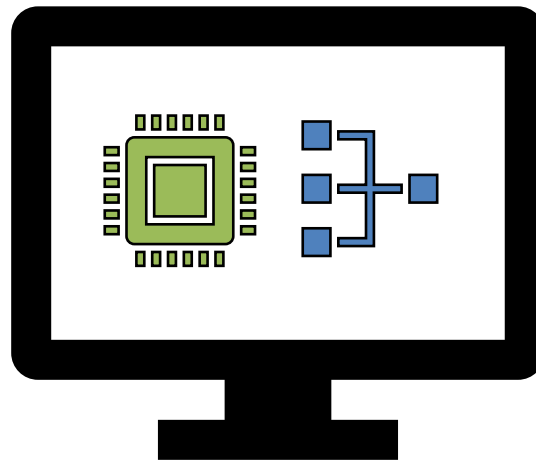
- To introduce workflow management system to cluster users
- To introduce some existing workflows and their development community
- To demonstrate how to run the workflow on RedRaider Cluster



- ❖ Part I: Workflow Management System
- ❖ Part II: Nextflow, data-driven computational pipelines
- ❖ Part III: nf-core pipelines, a community curation of bioinformatic pipelines
- ❖ Part IV: Getting Help



Computation Workflow Management System





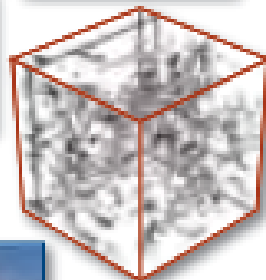
Part I: Workflow Management System

Science Paradigms

- Thousand years ago:
science was **empirical**
describing natural phenomena
- Last few hundred years:
theoretical branch
using models, generalizations
- Last few decades:
a **computational** branch
simulating complex phenomena
- Today: **data exploration** (eScience)
unify theory, experiment, and simulation
 - Data captured by instruments
or generated by simulator
 - Processed by software
 - Information/knowledge stored in computer
 - Scientist analyzes database/files
using data management and statistics



$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$



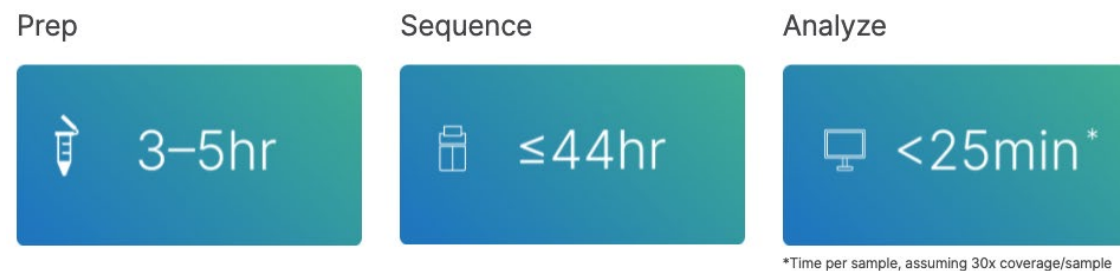
Fourth paradigm is using the computational power to generate, curate, analyze, archive a massive scientific data

Liew et al. 2016



Part I: Workflow Management System

Biological Experiment Sequencing Data Prep and Data Size (Illumina NovaSeq 6000)



System specifications range

Output range	80 - 6000 Gb
Paired end reads per run	1.6 - 40B
Max read length	2 × 250 bp
Run time	13 - 44 hours

In a single sequencing run on NovaSeq 6000:
about 4 days library prep and sequencing work could generate 6000 Gb data.



Part I: Workflow Management System

Challenges of scientific data explosion

- Large amount of data are generated from sequencing steps for different experiments
- Data need to go through a series of processing steps based on analysis
- Data processing may be executed on different platforms



TEXAS TECH UNIVERSITY

Information Technology Division™

Part I: Workflow Management System

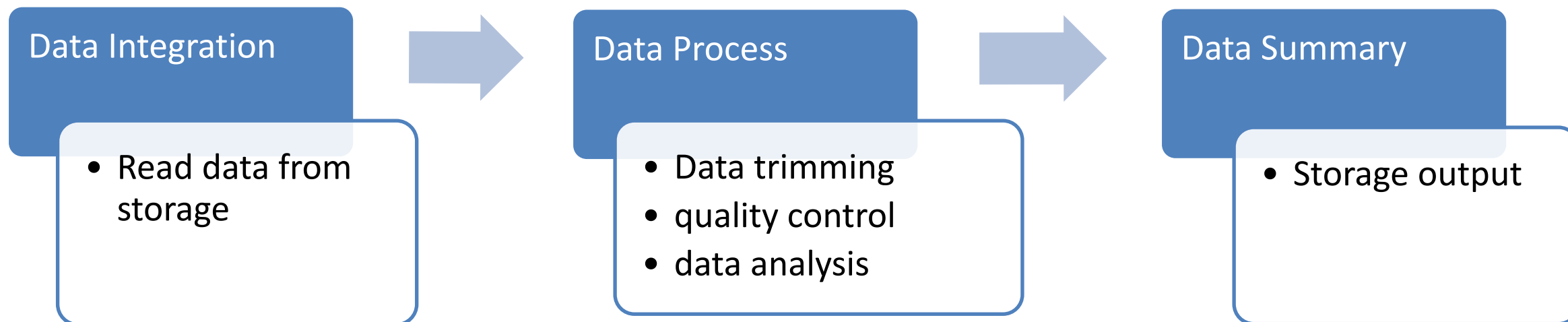
The bioinformatic researchers are trying to automate the computational data processes to assist the data preparation process



Part I: Workflow Management System

TEXAS TECH UNIVERSITY
Information Technology Division™

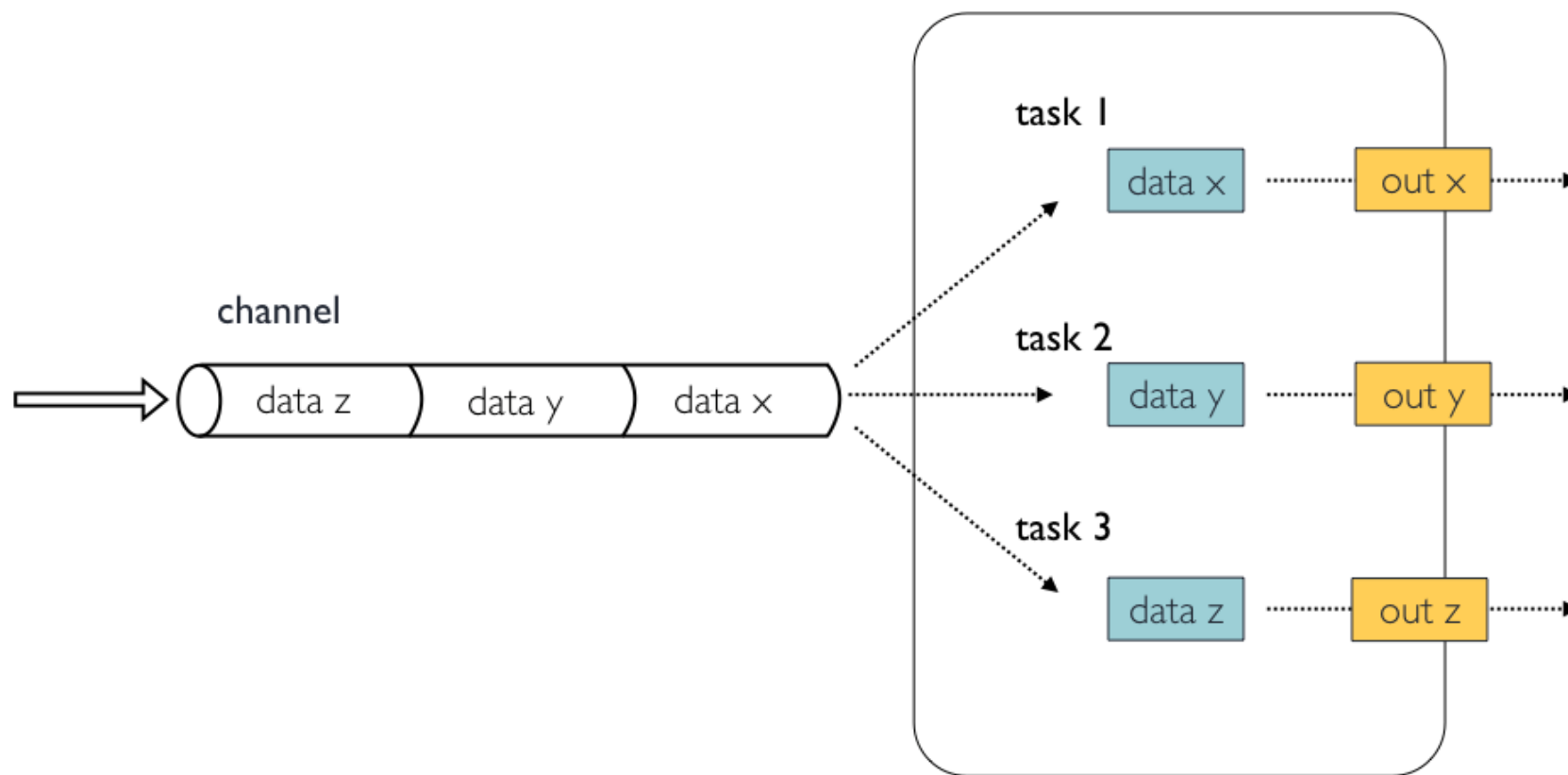
Bioinformatic Data Analysis Pipeline





Part I: Workflow Management System

Nextflow Processes and Channel





TEXAS TECH UNIVERSITY

Information Technology Division

Part I: Workflow Management System

Advance Solution

Model biological data as workflows and use a workflow management system to organize their execution.

Pipeline + HPC infrastructure -> workflow management



Part I: Workflow Management System

Dictionary: the sequence of industrial, administrative, or other processes through which a piece of work passes from initiation to completion.

Wikipedia:

Workflows may be viewed as one fundamental building block to be combined with other parts of an organization's structure such as information technology, teams, projects and hierarchies.



Part I: Workflow Management System

A [workflow management system](#) (WfMS) is a software system for setting up, performing, and monitoring of a defined sequence of processes and tasks, with the broad goals of increasing productivity, reducing costs, becoming more agile, and improving information exchange within an organization.



Data intensive workflow often involves:

- Moving data from data sources to computational resources
- Cleaning, calibrating and normalizing data
- Constructing a model using part of that preprocessed data
- Validating the model with the remaining data
- Visualizing the results
- Moving the results to a storage system



Part I: Workflow Management System

TEXAS TECH UNIVERSITY
Information Technology Division

Common roles for workflows

- Support for collaborative research by enabling scientific communities to share automated and formalized processes such as data analysis
- Construction free from distracting details about workflow management and execution
- The ability to automat workflow steps, that is, their mapping and execution, and to repeat in silicon experiments



Common roles for workflows (cont.)

- Integrating resources from distributed and heterogeneous enactment platforms
- Handling large volumes of data and complex computations
- Improving the execution through various optimization strategies



Part I: Workflow Management System

TEXAS TECH UNIVERSITY
Information Technology Division™

The benefit of using workflow

Use Domain Language

- handling I/O file between flows
- provide standard log file along with analysis

Work with scheduler

- Each process get its job submission
- Build-in checkpoint

Work with container

- Reproducibility / version control
- Avoid install conflicts

Part I: Workflow Management System



TEXAS TECH UNIVERSITY
Information Technology Division™

OPEN ACCESS

EDUCATION

Using prototyping to choose a bioinformatics workflow management system

Michael Jackson , Kostas Kavoussanakis, Edward W. J. Wallace

Published: February 25, 2021 • <https://doi.org/10.1371/journal.pcbi.1008622>

Article	Authors	Metrics	Comments	Media Coverage
---------	---------	---------	----------	----------------



Abstract

Author summary

Introduction

Conclusions

Supporting information

Acknowledgments

References

Reader Comments

Figures

Abstract

Workflow management systems represent, manage, and execute multistep computational analyses and offer many benefits to bioinformaticians. They provide a common language for describing analysis workflows, contributing to reproducibility and to building libraries of reusable components. They can support both incremental build and re-entrancy—the ability to selectively re-execute parts of a workflow in the presence of additional inputs or changes in configuration and to resume execution from where a workflow previously stopped. Many workflow management systems enhance portability by supporting the use of containers, high-performance computing (HPC) systems, and clouds. Most importantly, workflow management systems allow bioinformaticians to delegate how their workflows are run to the workflow management system and its developers. This frees the bioinformaticians to focus on what these workflows should do, on their data analyses, and on their science.

Part II



TEXAS TECH UNIVERSITY
Information Technology Division™

nextflow

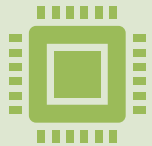
Part II: Nextflow



TEXAS TECH UNIVERSITY
Information Technology Division™



A reactive workflow framework



A programming DSL /DSL2 that eases the writing of data-intensive computational pipelines.

Part II: Nextflow



Nextflow Highlights

- Domain specific language
- Work with cluster scheduler
- Error handling
- Reproducibility



Scripting Language

- Groovy
- Bash
- Python
- R
- Perl

Mixing scripting language of Nextflow:
<https://www.nextflow.io/example2.html>

```
1  #!/usr/bin/env nextflow
2
3  params.range = 100
4
5  /*
6   * A trivial Perl script producing a list of numbers pair
7   */
8  process perlTask {
9      output:
10         stdout into randNums
11
12     shell:
13     '''
14     #!/usr/bin/env perl
15     use strict;
16     use warnings;
17
18     my $count;
19     my $range = !{params.range};
20     for ($count = 0; $count < 10; $count++) {
21         print rand($range) . ', ' . rand($range) . "\n";
22     }
23     '''
24 }
25
26 /*
27 * A Python script task which parses the output of the previous script
28 */
29 process pyTask {
30     echo true
31
32     input:
33         stdin from randNums
34
35     '''
36     #!/usr/bin/env python
37     import sys
38
39     x = 0
40     y = 0
41     lines = 0
42     for line in sys.stdin:
43         items = line.strip().split(",")
44         x = x+ float(items[0])
45         y = y+ float(items[1])
46         lines = lines+1
47
48     print "avg: %s - %s" % ( x/lines, y/lines )
49     '''
50 }
```




Part II: Nextflow

HPC Systems

- Slurm
- SGE
- LSE

Cloud Systems

- AW Batch
- Google Life Science



Part II: Nextflow

Benefit of Work with Different Scheduler:

- adaptability;
- parallelism of jobs:
 - create job channels
 - split job using scheduler

[Scatter Execution of Nextflow Examples](#)

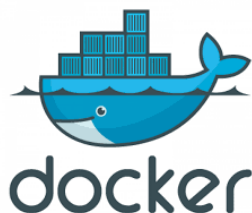
```
1  #!/usr/bin/env nextflow
2  // Author: Bogdan Kirilenko, 2020
3  // Nextflow procedure to execute chain feature extraction jobs
4  // Joblist contains a file where each line is a separate command
5  // We just call these lines in parallel
6
7  // params section: basically command line arguments
8  params.joblist = 'NONE' // file containing jobs
9
10 // if still default -> nothing assigned: show usage message and quit
11 if (params.joblist == "NONE"){
12     printf("Usage: nextflow execute_joblist.nf --joblist [joblist file] -c
13         System.exit(2);
14 }
15
16 // create channel lines -> we need to execute lines in parallel
17 joblist = file(params.joblist)
18 lines = Channel.from(joblist.readLines())
19
20 process execute_jobs {
21
22     // allow each process to fail 3 times
23     errorStrategy 'retry'
24     maxRetries 3
25
26     input:
27     val line from lines
28
29     "${line}" // one line represents an independent command
30 }
31
```



Part II: Nextflow

Version Control

- Docker
- Apptainer / Singularity
- Shifter
- Podman



Part II: Nextflow



TEXAS TECH UNIVERSITY
Information Technology Division™

Exercise #0 Install NextFlow



Part II: Nextflow

Exercise #0 Install Nextflow ([Solution](#))

Check prerequisites

Java 8 or later is required

1

Make sure 8 or later is installed
on your computer by using the command:
`java -version`

Note: version numbers "1.8.y_z" and "8" identify the same
[Java release](#).

Set up

Dead easy to install

2

Enter this command in your terminal:
`curl -s https://get.nextflow.io | bash`
(it creates a file `nextflow` in the current dir)

Note: it can also be downloaded from [GitHub](#) or installed by
using [Bioconda](#) package manager.

Launch

Try a simple demo

3

Run the classic *Hello world*
by entering the following command:
`./nextflow run hello`



Part II: Nextflow

Exercise #0 Install Nextflow ([Solution](#))

- `curl -s https://get.nextflow.io | bash`

Or

- `conda install -c bioconda nextflow`

Test Installation

- `./nextflow run hello`

Or

- `nextflow run hello`



TEXAS TECH UNIVERSITY

Information Technology Division™

Part II: Nextflow

Exercise #0 Install Nextflow ([Solution](#))

- Recommendation:

Install Nextflow in your conda environment



Part II: Nextflow

Potential Problem:

1. nextflow installation needs java 8+ version:

- `Java -version`

2. conda install needs Bioconda package manager:

- `conda config --add channels defaults`
- `conda config --add channels bioconda`
- `conda config --add channels conda-
forge`

3. [Install miniconda](#)



Part II: Nextflow

Your First Script:

- [tutorial.nf](#) (DSL2)

```
params.str = 'Hello world!'

process splitLetters {
    output:
        path 'chunk_*'

    """
    printf '${params.str}' | split -b 6 - chunk_
    """
}

process convertToUpper {
    input:
        path x
    output:
        stdout

    """
    cat $x | tr '[a-z]' '[A-Z]'
    """
}

workflow {
    splitLetters | flatten | convertToUpper | view { it.trim() }
}
```



Part II: Nextflow

Your First Script:

- [tutorial.nf](#) (DSL1)

```
#!/usr/bin/env nextflow

params.str = 'Hello world!'

process splitLetters {

    output:
    file 'chunk_*' into letters

    """
    printf '${params.str}' | split -b 6 - chunk_
    """
}

process convertToUpper {

    input:
    file x from letters.flatten()

    output:
    stdout result

    """
    cat $x | tr '[a-z]' '[A-Z]'
    """
}

result.view { it.trim() }
```



Exercise #1 Run first nextflow script

- Copy existing script to personal directory
- Run the script
- Modify script and run



Part II: Nextflow

Exercise #1 Run first nextflow script ([solution](#))

```
nextflow run tutorial.nf
```

```
#!/usr/bin/env nextflow

params.str = 'Hello world!'

process splitLetters {

    output:
    file 'chunk_*' into letters

    """
    printf '${params.str}' | split -b 6 - chunk_
    """
}

process convertToUpper {

    input:
    file x from letters.flatten()

    output:
    stdout result

    """
    cat $x | tr '[a-z]' '[A-Z]'
    """
}

result.view { it.trim() }
```



Part II: Nextflow

Exercise #1 Run first nextflow script ([solution](#))

```
nextflow run tutorial.nf
```

```
nextflow run tutorial.nf  
-resume
```

```
nextflow run tutorial.nf  
--str 'Bonjour le monde'
```

```
#!/usr/bin/env nextflow  
  
params.str = 'Hello world!'  
  
process splitLetters {  
  
    output:  
    file 'chunk_*' into letters  
  
    """  
    printf '${params.str}' | split -b 6 - chunk_  
    """  
}  
  
process convertToUpper {  
  
    input:  
    file x from letters.flatten()  
  
    output:  
    stdout result  
  
    """  
    cat $x | tr '[a-z]' '[A-Z]'  
    """  
}  
  
result.view { it.trim() }
```



Part II: Nextflow

DSL2

```
nextflow.enable.dsl=2

process foo {
    output:
        path 'foo.txt'
    script:
        """
        your_command > foo.txt
        """
}

process bar {
    input:
        path x
    output:
        path 'bar.txt'
    script:
        """
        another_command $x > bar.txt
        """
}
```

Part III



TEXAS TECH UNIVERSITY
Information Technology Division™



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



Part III: nf-core pipelines

Use nf-core pipelines

check available nf-core pipeline, through [website](#)

The screenshot shows the nf-core Pipelines website. At the top, there is a navigation bar with links for Home, Pipelines, Tools, Usage, Developers, Events, and About. Below the navigation bar is a green header with the word "Pipelines" in large white text. Underneath the header, it says "Browse the 53 pipelines that are currently available as part of nf-core." Below this is a section titled "Available Pipelines" with a link to "Let us know!" if you have a suggestion. There is a search bar and filter options: Released (33), Under development (15), and Archived (5). The sort options are Last Release, Alphabetical, and Stars. The display options are a grid and a list. The main content area shows four pipeline cards:

- nf-core/fetchngs** (45 stars): Pipeline to fetch metadata and raw FastQ files from public and private databases. Version 1.5, Published 7 days ago. Tags: ddbj, download, ena, fastq, geo, sra, synapse.
- nf-core/eager** (61 stars): A fully reproducible and state-of-the-art ancient DNA analysis pipeline. Version 2.4.1, Published 1 week ago. Tags: adna, ancient-dna-analysis, ancientdna, genome, metagenomics, pathogen-genomics, population-genetics.
- nf-core/nanoseq** (54 stars): Nanopore demultiplexing, QC and alignment pipeline. Version 2.0.1, Published 1 week ago. Tags: alignment, demultiplexing, nanopore, qc.
- nf-core/mag** (69 stars): Assembly and binning of metagenomes. Version 2.1.1, Published 2 weeks ago. Tags: annotation, assembly, binning, long-read-sequencing, metagenomes, metagenomics, nanopore, nanopore-sequencing.



Part III: nf-core pipelines

Exercise #2 Check available pipelines

- using `nf-core list`

Pipeline Name	Stars	Latest Release	Released
fetchngs	45	1.5	6 days ago
eager	61	2.4.1	1 week ago
nanoseq	54	2.0.1	1 week ago
mag	69	2.1.1	2 weeks ago
cutandrun	16	1.0.0	1 months ago
ampliseq	75	2.1.1	1 months ago
rnaseq	399	3.4	2 months ago
mhcquant	16	2.0.0	3 months ago
bacass	33	2.0.0	3 months ago
viralrecon	51	2.2	4 months ago
bcellmagic	16	2.0.0	5 months ago
bactmap	25	1.0.0	6 months ago
smrnaseq	31	1.1.0	6 months ago
sarek	129	2.7.1	6 months ago
hic	32	1.3.0	7 months ago



Part III: nf-core pipelines

Exercise #2 Check available pipelines (Solution)

- [Install nf-core](#)
 - `conda install nf-core` or
 - `conda update nf-core`
- [list all pipelines](#)
 - `nf-core list`



Part III: nf-core pipelines

Run nf-core pipelines on RedRaider Cluster

- Ask scheduler for an interactive session
- [Download nf-core pipelines](#)
- Adapt the config file
 - Enable singularity container
 - Define the executor

- Run the job

```
nextflow run <pipeline>/workflow/ -c <config file>
```

or

```
nextflow run <pipeline>/workflow/ -c <config file> -resume
```



Part III: nf-core pipelines

Config File

- Enable singularity container

```
singularity {  
  enabled = true  
  cacheDir = "<singularity image location>" }
```

- Define the executor

```
process {  
  executor = 'slurm'  
  clusterOptions = '-p nocona -N 1 -n 2' }
```

- Define the input files




Part III: nf-core pipelines

Exercise #3

- Download nf-core/RNASeq
- Run nf-core/RNASeq

```
(nxtfl_env) cpu-25-39:~$ nf-core download

NF-CORE 
nf-core/tools version 2.6 - https://nf-co.re

Specify the name of a nf-core pipeline or a GitHub repository name (user/repo).
[?] Pipeline name: rnaseq
[?] Select release / branch: 3.9 [release]

In addition to the pipeline code, this tool can download software containers.
[?] Download software container images: singularity

If transferring the downloaded files to another system, it can be convenient to have everything compressed in a single file.
This is not recommended when downloading Singularity images, as it can take a long time and saves very little space.
[?] Choose compression type: tar.gz
INFO Saving 'nf-core/rnaseq' download.py:159
    Pipeline revision: '3.9'
    Pull containers: 'singularity'
    Using $NXF_SINGULARITY_CACHEDIR: '/home/yannchen/nextflow'
    Output file: 'nf-core-rnaseq-3.9.tar.gz'
INFO Downloading workflow files from GitHub download.py:162
INFO Downloading centralised configs from GitHub download.py:166
INFO Found 30 containers download.py:493
Downloading singularity images ██████████ 100% • 30/30 completed
INFO Compressing download.. download.py:186
```



Part III: nf-core pipelines

Exercise #3

- Download nf-core/RNASeq
- Run nf-core/RNASeq

- ```
nf-core download
<rnaseq> -r <3.9> --
outdir <nf-core-rnaseq-
3.9> -x <none> -c
<singularity> <--force>
```

- ```
nextflow run <nf-core-  
rnaseq-3.9/workflow/> -c  
<test/test.config> --  
outdir <rnaseq-3.9-  
outdir>
```



Part III: nf-core pipelines

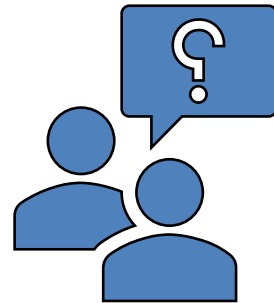
Develop nf-core pipelines

- Discuss with nf-core
- Download nf-core template
- Check available modules
- Write your own pipelines
- Push to nf-core github

```
nf-core create
```



Getting Help



Part IV: Getting Help



TEXAS TECH UNIVERSITY
Information Technology Division™

nf-core 

-  **slack** : communicate with developers, trouble shooting
-  **GitHub** : see source code
-  **twitter** : milestone announcement
-  **YouTube** : Tutorial videos

Part IV: Getting Help



TEXAS TECH UNIVERSITY
Information Technology Division™



- Visit Our Website:
 - [Nextflow](#)
 - [nf-core](#)
 - TTU Training Material
- Submit a support ticket:
 - Send an email to hpccsupport@ttu.edu



Part IV: Getting Help

Please help us improve this workshop and develop further training sessions:

- Ask questions at the end of this session
- Fill out survey forms
- Email any questions you have
- [Post-workshop survey](#)



- A Long run a partition RedRaider cluster
 - Hardware: cores, clock speed, RAM
 - Scheduler: job length
- Other workflow management systems





TEXAS TECH UNIVERSITY
Information Technology Division™



TEXAS TECH UNIVERSITY
Information Technology Division™