**NIH STRIDES** Accelerating biomedical research

## Accelerating Bioinformatics Workflows with Nextflow

Zelaikha Yosufzai, MS [C]



## **Introduction to Cloud Lab**



## **Opportunities and Challenges of the Cloud**

Cloud is the new frontier of frontier of biomedical research: the potential benefits are considerable, but so are the challenges.

# THE OPPORTUNITIES



Simplifies testing and iteration



Always available on demand



**Pay as you go** for only what you use



**Democratizes** access to scientific data





**Provides a rich set** of tools & services

# THE CHALLENGES

Acquisition is complex and time-consuming

Security protections are unclear

**Costs** are unpredictable

Training deficit is substantial



## **Gap Between Cloud Interest and Cloud Adoption**

It's very common for researchers to face significant barriers in taking the first step toward cloud adoption. Some of the barriers we hear about most frequently are included below.

- "I don't know which cloud platform to use or how different services compare to one another within or across cloud platforms."
- "Coordinating internal funding at NIH (MOUs and DCCs) and/or establishing contractual agreements with resellers for extramural institutions seems complicated and timeconsuming."
- "I'm not sure how much I need to budget for commercial cloud services and am concerned about overspending."

- "I'm not confident that I understand all of the roles and responsibilities necessary to support work in the cloud."
- "I'm uncertain about what the transition to cloud will entail and am hesitant to jump in with both feet."



## **NIH Cloud Lab** | Experiment in the Cloud

NIH Cloud Lab is a no-cost, 90-day program for NIH intra- and extramural researchers to try commercial cloud services in an NIH-approved environment. Cloud Lab provides training and guardrails to protect against financial and security risks.

#### **How It Works**



	NIHCloudLabGCP (Patie)	* 17 fek 50 + 12 Star 20 +			
	2' main - 2' 6 Branches © 0 Tags	Q, Go to file () Add file +	O Cede +	About (1)	
	😨 zbyosufizal. Merge pull request #68 from hadip	Documentation and tutorials on using GCP for biomedical research			
	📄 .gthub	update gh-action	2 weeks ago	P cloud.nih.gov/resources/cloudlab/	
	🖿 docs	Updated image path in vertexal red	2 months ago	bioinformatics gep people-cloud	
	images	added images for vertexai.md	2 months ago	the-sciences nin genal	
	in notebooks	Github Action: Lint Notebooks	2 weeks ago	🖾 Readme	
	in testing	Update test_notebooks.py	3 weeks ago		
	Inarkdown-link-check.json	reformatted readme status, rm tutorials dir, updated $mark_{-}$	9 months ago	Custom properties  21 stars	
	README.md	fixed pubmed tutorial link in README.md	9 months ago	4 watching	
	II README		/ =	V 18 forks Report repository	
	GCP Tutorial Resource	Releases No minute published Create a new minute			
	We have pulled together a variety of tutorials here from disparate sources. Some use Compute Engine, others use Vertex AI notebooks, and others use only managed services. Tutorials are organized by research method, but we try to designate what GCP pervices are used to halp work unavitate			Packages No packages published Publish your first package	

NIH Cloud Lab GCP Tutorial Repository





Example of NIH Cloud Lab Use Case

## NIH STRIDES

Accelerating biomedical research

#### **NIH Use Cases**

#### **Evaluate Utility & Cost**

Provides an easy route to evaluate the cloud's utility/cost for a project without major time or financial commitments

#### **Develop New Tools**

Allows experienced teams to prototype new architectures and evaluate software and hardware combinations

#### **Share Ideas**

Connects NIH'ers from across ICs to share ideas on how to conduct biomedical research in the cloud

#### Learn New Skills

Simplifies access to tools and cloud environments that participants can use for training purposes

\*Step 2 will be provided to you by the Cloud Lab team.

## **Training Materials**

We provide jumpstart pages, in-depth tutorials, and how-to documents that help users get ready for the cloud.

#### Azure

- Jumpstart page: <u>https://cloud.nih.gov/resour</u> <u>ces/cloudlab/azure-</u> jumpstart/
- GitHub repo: <u>https://github.com/STRIDES/</u> <u>NIHCloudLabAzure</u>





Accelerating biomedical research

#### **Google Cloud**

- Jumpstart page: <u>https://cloud.nih.gov/resour</u> <u>ces/cloudlab/google-cloud-jumpstart/</u>
- GitHub repo: <u>https://github.com/STRIDES/</u> <u>NIHCloudLabGCP</u>



#### AWS

- Jumpstart page: <u>https://cloud.nih.gov/resour</u> <u>ces/cloudlab/aws-jumpstart/</u>
- GitHub repo: <u>https://github.com/STRIDES/</u> <u>NIHCloudLabAWS</u>



## **Publicly Available Interactive Tutorials Across All Three CSPs!**

Cloud Lab offers a <u>suite of interactive tutorials</u> designed to help participants perform viable research in the cloud. Open to all who wish to gain skills, not just Cloud Lab participants!

#### Scientific

- Accessing Open Datasets
- AI/ML
- BLAST
- FHIR
- Genomic Variant Calling
- GWAS
- Long Read Genomics
- Medical Image Segmentation
- Proteomics
- RNAseq
- scRNAseq
- Radiogenomics
- SRA Data Interaction
- VCF Query with SQL
- ...and more



#### GenAl

- Model Training
- Guidance & Best Practices
- Introduction to tools across all platforms
- Leverage Large Language Models (LLMs)
- Chatbot Implementation
- Retrieval- Augmented Generation (RAG)

#### **General Cloud Ops**

- Access Marketplace Offerings
- Access Public Datasets
- Command Line Tools / SDKs
- Disk Images
- Ingest and Store Data
- Introduction to Cloud
- Jupyter Notebooks
- Serverless Functionality
- Virtual Machines

#### **Computing & Code**

- Conda Environments
- Container Registries
- Git Repos
- HPC Clusters
- Kubeflow
- Serverless Functionality

Code O Inser i I huirequests O Actors Projects ED Will O Security L Insights  P main - NHIRCould.Lab.AWS / tutorials / Go to fie Add file -  Melocomedia HAH Update RADAE.md  Arran I tutorial Resources  AWS Tutorial Resources  Overview of Page Contents  BacAMAE.md  BacAM	STRIDES / NIHCloudLabAWS (Public)				⊙Watch 3 • ¥ Fork 6 • ☆ Star 3 •			
P     Image: NBHCloudLabAWS / tutorials /     Image: Med Here       Image: NBHCloudLabAWS / tutorials /     Image: Med Here     Image: Med Here       Image: NBHCloudLabAWS / tutorials /     Image: Med Here     Image: Med Here       Image: NBHCloudLabAWS / tutorials // tutorial     Image: NBHCloudLabAWS // tutorials // tutorials // tutorials // tutorials // tutorials // tutorials // tutorial       Image: NBHCloudLabAWS // tutorials	Code 💿 Issues 👔 🕮 Pull	requests 🕑 Actions 🖽 Projects	🛛 Wiki 🔇	Security	∠ Insights			
	۶ main + NIHCloudLabAWS	/ tutorials /				Go to file Add file *		
Image: State	kyleoconnell-NIH Update READ!	ME.md				d1f5873 last month 🕥 Histe		
rotebooks      Update SRA-Download symb      second      RADME.md      Update RRADME.md      text      RADME.md      Coverview of Page Contents      Coverview of Page Contents      Sourcedual Workflows on AWS      Download SRA Data      CwsS      Submedical Workflows on AWS      Submedical								
	notebooks	Update SRA-Download.jpynb				5 months ago		
	README.md	.md Update README.md			last month			
AWS Tutorial Resources Overview of Page Contents   Bornedical Workflows on AWS  Cownbad SRA Data  GWAS  Wedical Imaging  RWAdeq  SWAS  BLAST  Long Read Sequencing Analysis  AJMAL Repline	= README and							
Biomedical Workflows on AWS Downhoad SRA Data GWAS GWAS GWAS Wedical Imaging RWAseq USNAseq BLAST BLAST Corp Read Sequencing Analysis AJML Pensine	AWS Tutorial R	Resources						
Domnbad SRA Data     WMS	AWS Tutorial R	Resources						
CWAS     Medical Insging     FNAseq     stRVAreq     LRAT     Long Read Sequencing Analysis     AJML Previne	AWS Tutorial R Overview of Page C	Contents						
Medical Imaging     RVAceq     scRVAceq     scRVAceq     ELAST     Long Reud Sequencing Analysis     AJML Popeline	AWS Tutorial R Overview of Page C Biomedical Workflows on Download SRA Data	Contents						
RNAcq     scRNAcq     scRNAcq     scRNAcq     BLAST     Long face Sequencing Analysis     AJML Popline	AWS Tutorial R Overview of Page C • Biomedical Workflows on • Download SRA Data	iontents Aws						
scRNAxeq     BLAST     Long Read Sequencing Analysis     Long Read Sequencing Analysis     AJML Papeline	AWS Tutorial R Overview of Page C Biomedical Workflows on Download SRA Data GMAS Medical Imaging	Contents AWS						
BLAST     Long Read Sequencing Analysis     AIML Pipeline	AWS Tutorial R Overview of Page C Biomedical Workflows on Download SRA Data GWAS Medical Imaging RNAueg	Contents ANS						
Long Read Sequencing Analysis     Al/ML Pipeline	AWS Tutorial R Overview of Page C Bornedical Workflow on Download SRA Data GWAS Medical Imaging BRIAnce scRNAseq	Contents ANNS						
<ul> <li>Al/ML Pipeline</li> </ul>	AWS Tutorial R Overview of Page C - Biamedical Workflows on - Download SKA Data - GWAS - Medical Imaging - RIAAreq - scRVAdrq - BLAST	iontents AWS						
- Once Date	AWS Tutorial R Overview of Page C Biomedical Workflows on Download SRA Data GWAS Medical Imaging Rikkang sRARAng BLAST Long Read Sequencing A	contents Aws						

#### NIH Cloud Lab AWS Tutorial Repository



## Introduction to Nextflow



## What is Nextflow?

Nextflow is a powerful workflow management system designed specifically for bioinformatics and data science that can be run locally or in the cloud.

Imagine you have a complex analysis task that involves many different programs and steps, potentially running on different computers or cloud services. Nextflow lets you define this entire process as a program, specifying the order of operations, how data flows between steps, and how to handle errors or resource limitations.

#### **Nextflow consists of two main pieces:**

**Main file:** This is the core of your workflow. It's written in the Nextflow DSL (Domain-Specific Language) and defines the processes, data channels, and execution logic of your analysis. Files are created using the .nf extension

• You can also break out processes into modules and sub workflows

**Config file:** Allows you to customize various aspects of your workflow's execution without modifying the main pipeline script improving reproducibility. Here you can set parameters like input/output paths, resource allocation (CPU, memory), execution environments (e.g., Docker images, Singularity containers), and other settings affecting the pipeline's behavior across different runs and environments. Config files are created with the .config extension.



## **Scripting in Nextflow**

A Nextflow script, typically with a .nf extension, is composed of several key parts working together to define your data processing workflow:

- **1. params** block: Defines parameters that control the workflow's behavior.
  - These can be things like **input file paths**, **output directories**, **specific program options**, **or thresholds**. You can set default values, allowing users to run the workflow with minimal configuration, or override them at runtime. This allows for flexibility and reproducibility across different runs.
- 2. processes block: This is the heart of the script. It defines individual processing steps, each encapsulated as a "process." Each process consists of:
  - 1. name: A unique identifier for the process.

Accelerating biomedical research

- 2. script: The commands to execute (often shell commands). This is where you specify the software and its options.
- 3. input: Specifies the data the process requires. This is often defined using channels, which manage the flow of data between processes.
- 4. output: Specifies the data the process produces. Again, often uses channels.



## **Scripting in Nextflow continued**

- 3. **channels** (implicit or explicit): Channels are the pipelines' arteries. They manage the flow of data between processes. Data is passed between processes via channels. While not always explicitly defined, they are the fundamental mechanism for data transfer. If not explicitly named, Nextflow will infer channels based on input and output statements within processes.
- 4. **workflow block**: This section assembles the individual processes into a cohesive workflow. It defines the execution order and how the output of one process becomes the input of another. This is often achieved by connecting processes through channels.
- 5. **publishDir** (optional): This parameter specifies a directory where the final results of the workflow will be saved.

6. **workDir** (optional): This parameter specifies a directory where intermediate files will be stored during the workflow's execution.



## **Using Nf-Core**

Let's not reinvent the wheel! One of the great advantages to using Nextflow is that it is integrated with a community driven effort named <u>nf-core</u>.

Nf-core's goal is to curate and collect sets of analysis workflows or pipelines built using Nextflow. These workflows/pipelines are bioinformatic focused and allow the user to run well known analyses like RNAseq, Methylseq, etc.

Running any of the nf-core workflows automatically pulls the workflow scripts from GitHub. As depicted from the image to the right Nextflow print outs the Nextflow version and that it is pulling the Methylseq scripts from the nf-core GitHub.





## **Installing Nextflow and Using nf-core**

The code below shows how to install Nextflow and run a nf-core pipeline.

#### 1. Install Java

In order to successfully install Nextflow your system must have Java preinstalled.

```
sudo apt update
sudo apt-get install default-jdk -y
java -version
```

#### 2. Specify Plug-ins

If you are using a cloud platform to run your Nextflow analysis you can specify the associated Nextflow plug-in, in the example below we are installing the google plug-in.

export NXF\_MODE=google

#### 3. Install Nextflow

Install Nextflow, make it executable, and update it.

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

chmod +x nextflow

./nextflow self-update

#### 4. Run a nf-core pipeline

```
./nextflow run nf-core/methylseq \
    --input 'SRR067701.fastq.gz' \
    --genome GRCh37 \
    --single_end \
    --max_cpus 32 \
    --max_memory '110.GB'
```



## **Utilizing Schedulers**

Cloud platforms like Azure, Google Cloud, and AWS all have batch systems that provides users a serverless, scalable, and cost-effective way to run Nextflow pipelines.

- 1. Configure specified Batch system, eliminating manual cluster management
- 2. Auto provision and deallocates resources, handling parallelism and scaling
- 3. Store results in cloud storage for long term access











## **Poll and Questions, Call to Action Banner**

The Cloud Lab team can help users navigate the complexities of cloud computing by offering informative sessions tailored to their needs providing clear explanations of how various cloud services function, focusing on practical application and addressing specific user questions.

### **Support We Provide**

- 1. Technical support related to tutorial challenges
- 2. General Troubleshooting
- 3. Ensuring work done within Cloud Lab can be migrated to a production ready environment



Accelerating biomedical research

#### **Collaboration Opportunities**

Cloud Services is creating tutorials monthly and would love to collaborate with your IC to develop scientifically focused content, if interested please email us at <u>CloudLab@nih.gov</u>.

I.	

#### **User Poll: Tutorial Content**

Vote on what tutorial topic you find most beneficial to you.

Join at: Slido.com











## **NIH STRIDES**

Accelerating biomedical research

For more information, please visit <u>cloud.nih.gov</u> or send an email to <u>STRIDES@nih.gov</u>.