# BIG DATA

Hands on AWS

# Identity and Access Management

## Identity and Access Management (IAM)
- Web service that controls fine-grained access to AWS resources
- IAM controls who is authenticated and authorized to use resources

## IAM user
- Unique identity recognized by AWS services and applications
- Similar to user in an operating system like Windows or UNIX

# Identity and Access Management

## IAM role

- Set of policies for making AWS service requests
- Trusted entities (e.g., such as IAM users) assume roles
  - Delegate access with defined permissions to trusted entities
  - There is no limit to the number of IAM roles a user can assume

## User vs role

- User has permanent long-term credentials and is used to directly interact with AWS services
- Role does not have credentials and cannot make direct requests to AWS services
- Roles are assumed by authorized entities, such as IAM users

# Identity and Access Management

Alice (i.e., an IAM user) is a firewoman

- She is the same person with or without her turnout gear
- As a firewoman (i.e., a role)
  - If she speeds to a house fire and passes a police officer, he isn't going to give her a ticket
  - In her role as a *firewoman*, she is allowed to speed to the house fire
- As a private citizen (i.e., another role)
  - When she is off duty, if she speeds past that same police officer, he's going to give her a ticket
  - In her role as a *private citizen*, she is not allowed to speed

# AWS

Amazon Web Services (AWS) is a public-cloud platform
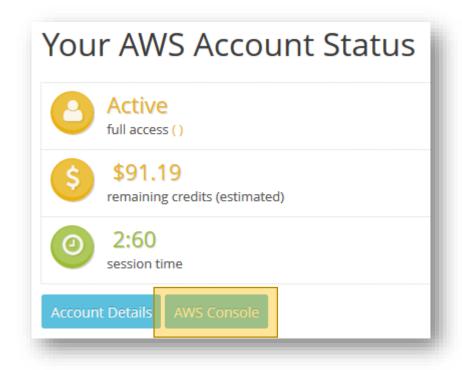
Services can be accessed in multiple ways
- Web GUI: intuitive point and click access without any programming
  - Intuitive interfaces is part of the attraction of cloud services
  - Tedious if the same actions must be performed repeatedly
- (REST) Application programming interface (API)
  - Permits requests to be transmitted via Hypertext Transfer Protocol (HTTPS)
- Software development kits (SDKs) that you install on your computer
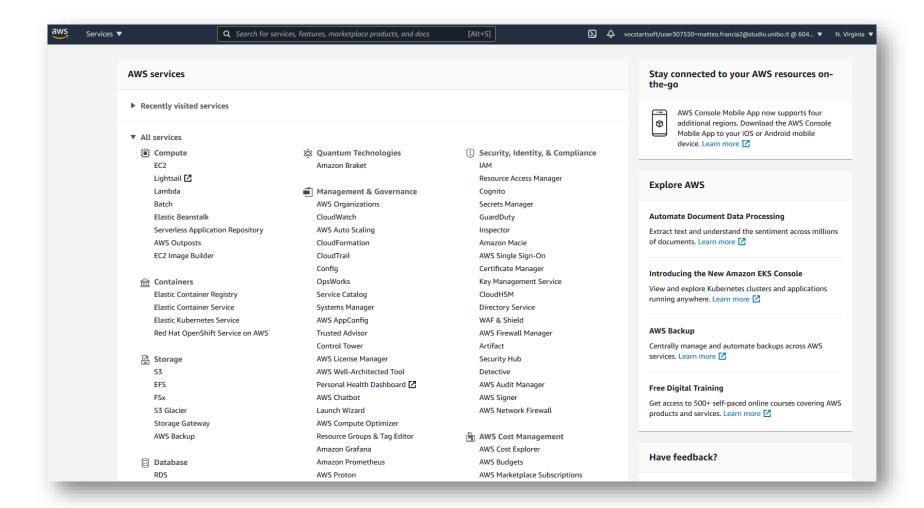  - Access from programming languages such as Python, Java, etc.

# AWS Web console

We use the AWS Educate program
- Login with the provided account
- You got 150$ to work on AWS services
- Provisioned services charge even if not used

https://www.awseducate.com/signin/SiteLogin



Your AWS Account Status

Active
full access ( )

$91.19
remaining credits (estimated)

2:60
session time

Account Details    AWS Console

# AWS Web console

# AWS CLI

CLI interface
- Necessary to install the CLI (version 2)
- See https://docs.aws.amazon.com/cli/latest/userguide/install-cliv2.html

```
Synopsis

********

aws [options] <command> <subcommand> [parameters]


Description

***********

A unified tool to manage your AWS services.
```
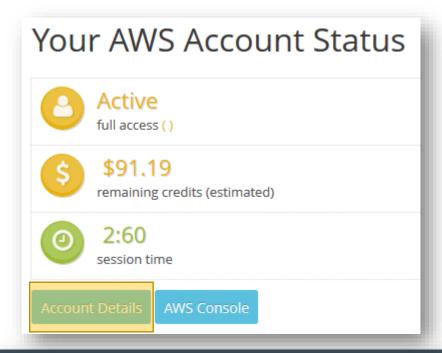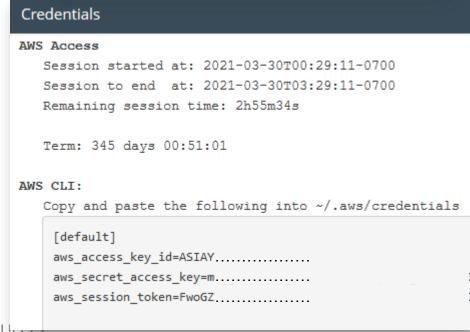
https://docs.aws.amazon.com/cli/latest/userguide/install-cliv2-linux.html

# AWS CLI

CLI needs credentials to work

- Go back to AWS Educate
- Click on "Account Details"
- Copy the content into the file `~/.aws/credentials`
- Henceforth, we assume that you have set up the credentials file
- Credentials expire after some time; you need a manually refresh

## Your AWS Account Status

**Active**
full access ( )

**$91.19**
remaining credits (estimated)

**2:60**
session time

Account Details    AWS Console

## Credentials

**AWS Access**
    Session started at: 2021-03-30T00:29:11-0700
    Session to end  at: 2021-03-30T03:29:11-0700
    Remaining session time: 2h55m34s

    Term: 345 days 00:51:01

**AWS CLI:**
    Copy and paste the following into ~/.aws/credentials

    [default]
    aws_access_key_id=ASIAY.................
    aws_secret_access_key=m.................
    aws_session_token=FwoGZ.................

# AWS CLI

Run `aws configure`
- Confirm AWS Access Key ID (press enter)
- Confirm AWS Secret Access Key (press enter)
- Set Default region name to `us-east-1`
- Set Default output format to `json`

It is also possible to configure an AWS profile
- A (named) profile is a collection of settings and credentials
- If profile is specified, its settings and credentials are used to run a command
- When no profile is explicitly referenced, use `default`
  - We stick to `default`

# Object storage: S3

Create S3 bucket, the following rules apply for naming buckets

- Must be between 3 and 63 characters long
- Can consist only of lowercase letters, numbers, dots (.), and hyphens (-)
- Must be unique within a partition (i.e., a group of regions)

```
$ git clone https://github.com/w4bo/bigdata-aws/

$ cd bigdata-aws/lab01-lambda

$ aws s3api create-bucket --bucket aws-bucket-bigdata2021

$ aws s3 cp datasets/inferno.txt s3://aws-bucket-bigdata2021/inferno.txt

$ aws s3api list-objects --bucket aws-bucket-bigdata2021
```

https://s3.console.aws.amazon.com/s3/home?region=us-east-1#

# BIG DATA

Data pipelines on AWS Lambda

# Requirements

To start this lecture, you need to

- Activate your AWS Educate account
- Either
  - Install the necessary software
    - git
    - IntelliJ IDEA (with AWS Toolkit and Scala plugins)
    - python
    - java 1.8
    - Docker
    - AWS CLI, AWS SAM CLI
  - Be able to download and run the VM

# AWS SAM CLI

Serverless Application Model is a framework to build serverless applications

- A serverless application is a combination of Lambda functions, event sources, etc.
- Install AWS SAM CLI (on Linux)

```
sudo group add docker

sudo usermod –aG docker $USER

newgrp docker

sudo chmod 666 /var/run/docker.sock

wget https://github.com/aws/aws-sam-cli/releases/latest/download/aws-sam-cli-linux-x86_64.zip

unzip aws-sam-cli-linux-x86_64.zip -d sam-installation

sudo ./sam-installation/install

sam --version
```

https://docs.aws.amazon.com/serverless-application-model/latest/developerguide/serverless-sam-cli-install.html

# AWS services

AWS Educate (and AWS console)

- https://aws.amazon.com/it/education/awseducate/
- https://console.aws.amazon.com/console/home?region=us-east-1

IAM (authentication)

- https://docs.aws.amazon.com/IAM/latest/UserGuide/iam-ug.pdf

SDK (software API)

- https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/home.html

Lambda (serverless computing and processing)

- https://docs.aws.amazon.com/lambda/latest/dg/getting-started.html
- https://console.aws.amazon.com/lambda/home?region=us-east-1#/functions

DynamoDB (key-value database)

- https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/Introduction

S3 (object storage)

- https://s3.console.aws.amazon.com/s3/home?region=us-east-1

# Case study

Given a dataset of sales per customer
find the products frequently bought together

```
Dataset sample
%%%%%%%%%%%%%%


[ { customerName: Alice, products: [Pizza, Beer, Diaper] },
  { customerName: Bob, products: [Pizza, Beer, Diaper] },
  { customerName: Charlie, products: [Pizza, Cola] } ]
```

# Case study

The pipeline involves a single transformation

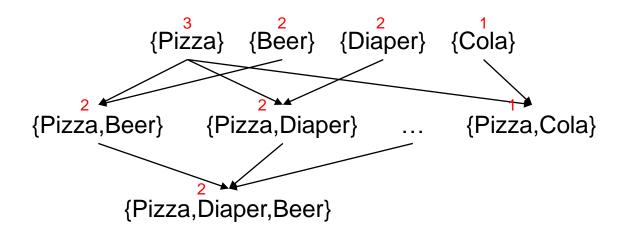- A classic mining problem, which one?

# Frequent itemset mining

Find sets of items (i.e., itemsets) frequently appearing together

- **Item**: a product
- **Itemset**: a set of products
- **Frequently**: support above threshold
- **Support**: number of clients buying a set of products

Complexity: $O(2^{|items|})$

```
Dataset sample

%%%%%%%%%%%%%%

[[Pizza, Beer, Diaper],
 [Pizza, Beer, Diaper],
 [Pizza, Cola]]
```

# Case study

$$FIM: List[List[String]] \rightarrow List[Set[String]]$$

- FIM requires a list of lists as input, but we have nested JSON objects
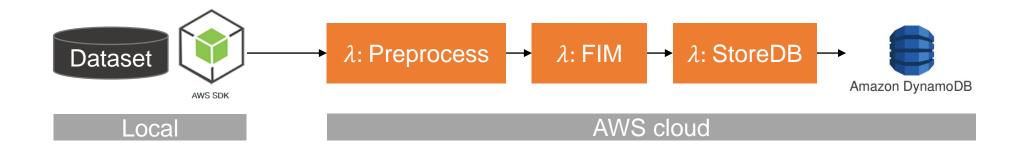- We need a pre-processing step

```
Raw dataset sample

%%%%%%%%%%%%%%

[ { customerName: Alice, products: [Pizza, Beer, Diaper] },

  { customerName: Bob, products: [Pizza, Beer, Diaper] },

  { customerName: Charlie, products: [Pizza, Cola] } ]
```

```
Processed dataset sample

%%%%%%%%%%%%%%

[[Pizza, Beer, Diaper],

 [Pizza, Beer, Diaper],

 [Pizza, Cola]]
```

Finally, we need to store the itemsets in the database

# Reference pipeline

# NOSQL storage: DynamoDB

Basic DynamoDB components: tables and items

**Tables**, collection of (data) items

**Items**, a group of attributes that is uniquely identifiable

- Each table contains zero or more items
  - No limit to the number of items you can store in a table
- Each item in the table has a unique identifier, or primary key
- E.g., in the table `people`, each item represents a `person`
  - The primary key consists of one attribute (`fiscalCode`)

# NOSQL storage: DynamoDB

## Attributes

- A data element that is not broken down any further
  - E.g., an item in the `people` table contains attributes `fiscalCode` and `lastName`
- Most of the attributes are scalar (have only one value)
- Some of the items have a nested attribute (`address`) up to 32 levels deep

## Schemaless

- Other than the primary key, a table is schemaless
  - Neither the attributes nor their data types need to be defined beforehand
  - Each item can have its own distinct attributes

# NOSQL storage: DynamoDB

Primary Key
- To create a table, you must specify the primary key of the table
- No two items can have the same key

Two types of primary keys
- Partition key: a simple primary key composed of one attribute (partition key)
  - Keys are inputs to an internal hash function
  - The hash function determines the physical partition in which the item will be stored
  - E.g., access any item in the `people` table directly by providing the `fiscalCode`
- Composite primary key: partition key and sort key (two attributes)
  - First attribute is the partition key
  - Second attribute is the sort key
  - Items in same partition key value are stored together and sorted by sort key

# NOSQL storage: DynamoDB

| Primary Key | | Data-Item Attributes... | | |
|---|---|---|---|---|
| **Partition Key** | **Sort Key** | Attribute 1 | Attribute 2 | ... |
| **HR-974**<br>*(employee ID)* | Employee_Name | **Data:** Murphy, John<br>*(employee name)* | **Start:** 2008-11-08<br>*(start date)* | ...etc. |
| | YYYY-Q1 | **Data:** $5,477<br>*(order totals in USD)* | **Name:** Murphy, John<br>*(employee name)* | |
| | HR_confidential | **Data:** 2008-11-08<br>*(hire date)* | **Name:** Murphy, John<br>*(employee name)* | ...etc. |
| | Warehouse_01 | **Data:** Murphy, John<br>*(employee name)* | | |
| | v0_Job_title | **Data:** Operator-1<br>*(job title)* | **Start:** 2008-11-08<br>*(start date)* | ...etc. |
| | v1_Job_title | **Data:** Operator-2<br>*(job title)* | **Start:** 2016-11-04<br>*(start date)* | ...etc. |
| | v2_Job_title | **Data:** Supervisor-1<br>*(job title)* | **Start:** 2017-11-01<br>*(start date)* | ...etc. |

https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/bp-gsi-overloading.html

# NOSQL storage: DynamoDB

Create a table `frequent-sales` with a composite key

- `dataset`: String
- `timestamp`: String

```
$ aws dynamodb create-table \

    --table-name frequent-sales \

    --attribute-definitions AttributeName=dataset,AttributeType=S AttributeName=timestamp,AttributeType=S \

    --key-schema AttributeName=dataset,KeyType=HASH AttributeName=timestamp,KeyType=RANGE \

    --provisioned-throughput ReadCapacityUnits=1,WriteCapacityUnits=1


$ aws dynamodb list-tables


$ aws dynamodb delete-table --table-name frequent-sales
```

# NOSQL storage: DynamoDB

Reading data from DynamoDB might not reflect the results of a recent write

Eventually Consistent Reads (default)
- Response might include stale data
- After short time, the response should return the latest data

Strongly Consistent Reads
- Response includes the most up-to-date data
- A strongly consistent read might not be available if there is a network delay or outage
  - In this case, DynamoDB may return a server error (HTTP 500)
- Strongly consistent reads may have higher latency than eventually consistent reads
- Strongly consistent reads are not supported on global secondary indexes

# NOSQL storage: DynamoDB

Provisioned mode: specify the #reads and #writes per second

- You have predictable application traffic or traffic ramps gradually
- You can forecast capacity requirements to control costs

One read capacity unit

- One strongly consistent read per second, two eventually consistent reads per second
- RCUs also depend on the item size (a read is up to 4 KB in size), if item size is 8 KB
  - 2 RCUs to sustain one strongly consistent read per second
  - 1 RCU if you choose eventually consistent reads

One write capacity unit represents one write per second for an item up to 1 KB in size

# NOSQL storage: DynamoDB

Put a new item and get it back

```
$ aws dynamodb put-item
    --table-name frequent-sales
    --item '{"dataset": {"S": "sales"}, "timestamp": {"S": "1611226870"}, "bar": {"S": "foobar"}}'


$ aws dynamodb query
    --table-name frequent-sales
    --key-condition-expression "dataset = :n"
    --expression-attribute-values '{":n":{"S":"sales"}}'
```
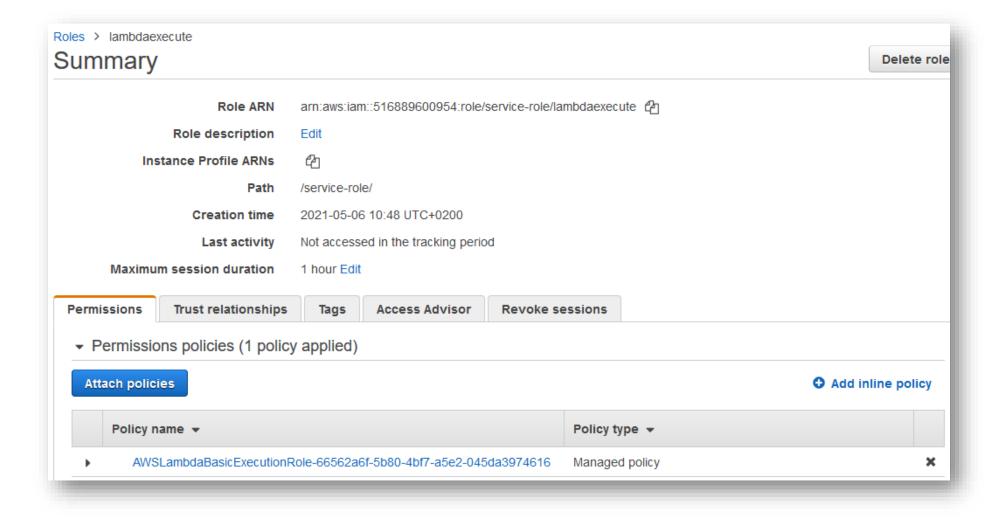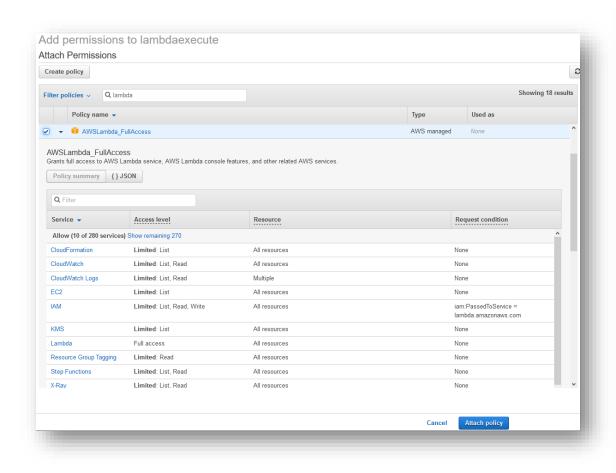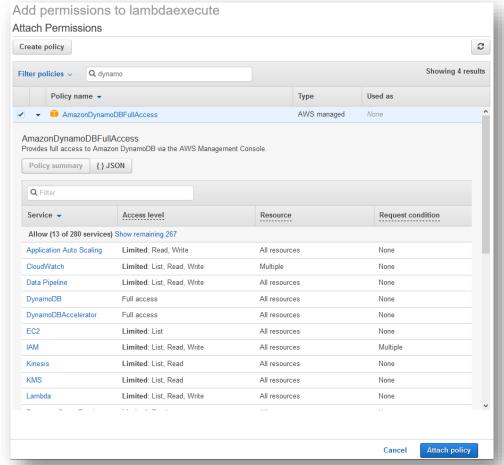
# Lambda: create a function



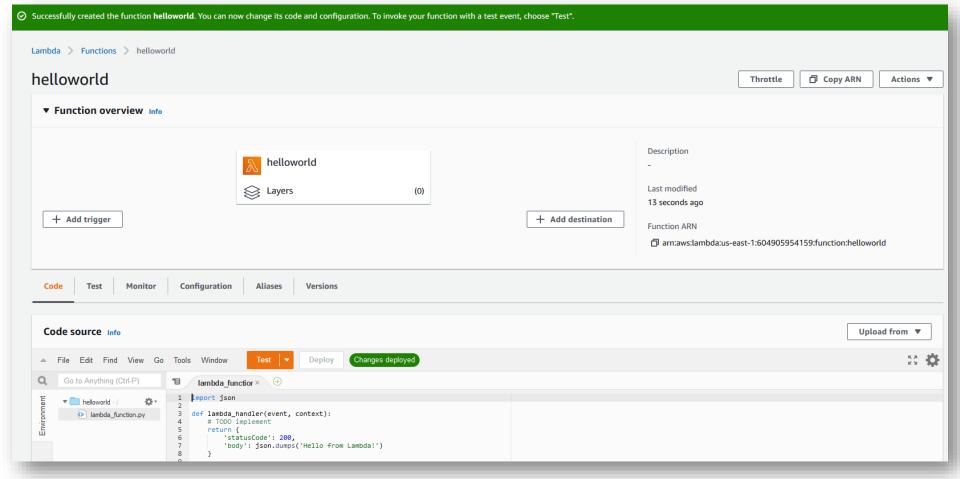https://console.aws.amazon.com/lambda/home?region=us-east-1#/functions

# Lambda: attaching a role

# Lambda: attaching a role

# Lambda: attaching a role

# Lambda: attaching a role

# Lambda: create a function

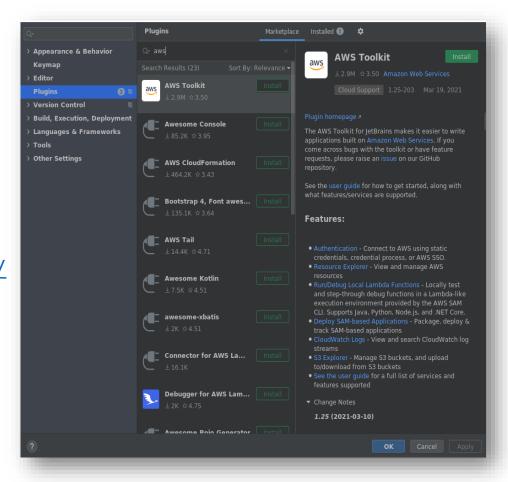https://console.aws.amazon.com/lambda/home?region=us-east-1#/functions

# Lambda: create a function

Manually creating the functions is cumbersome

- We must copy and paste code
- No automatic testing
- No debugging
- No IDE support (and not all languages are supported)

Switch to IntelliJ IDEA + AWS Toolkit

# AWS Toolkit

- Get the latest IntelliJ IDEA
- Install the `AWS Toolkit`
- Copy the credentials
  `cp ~/.aws/credentials ~/.aws/config`
- Clone the repo
  `git clone` https://github.com/w4bo/bigdata-aws/
- Import `lab01-lambda` as a Gradle project
- Verify that the project builds
  `./gradlew`

# AWS Toolkit

## Click on `AWS Explorer`

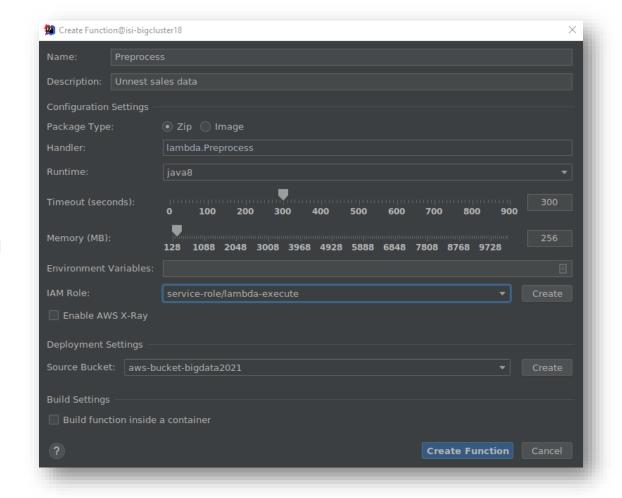- You can see the `helloworld` function
- Plus `CloudWatch Logs` and `S3`

# AWS Toolkit

Test the existing code locally
- With Gradle
- Or with local Lambda execution

Deploy a new Lambda function from the existing code
- Right click on AWS Explorer > Lambda
- Select `Create new AWS Lambda…`
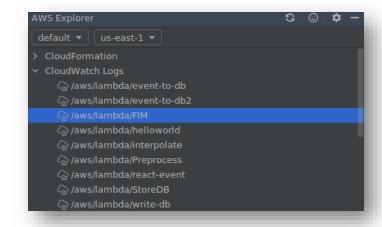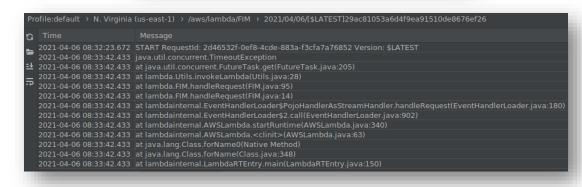- Populate the settings
- `Create the function`

https://aws.amazon.com/lambda/pricing/

# AWS Toolkit



Check the log for errors and pricing

- AWS Toolkit > CloudWatch Logs
- Double click on the function name
- Double click on the log entry

# Data pipeline

Deploy and execute the HelloWorld.java lambda function

Given the created storage: S3 and DynamoDB

- Deploy the function `FIM`
- Deploy the function `Preprocess`
- Run ReadDataset.java
- Check that the table `frequent-sales` has the FIs for the dataset `sales`

Some hints

- Function names are case sensitive
- Some function need more than 128MB of memory
  - Behold! The higher the RAM, the higher the costs

# BIG DATA

Amazon EMR

# EC2

AWS uses public-key cryptography to secure the login

You can create one using the Amazon EC2 console
- Open the Amazon EC2 console at https://console.aws.amazon.com/ec2/
- In the navigation pane, choose `Key Pairs`
- Choose `Create key pair`
- For `Name`, enter a descriptive name for the key pair
- For `File format`, choose the format in which to save the private key
  - OpenSSH, choose `pem` (`chmod 400 *my-key-pair*.pem`)
  - PuTTY, choose `ppk`
- Choose `Create key pair`
- The private key file is automatically downloaded by your browser

# Creating the cluster

Choose the frameworks and applications to install

Data process
- Submit jobs or queries directly to installed applications
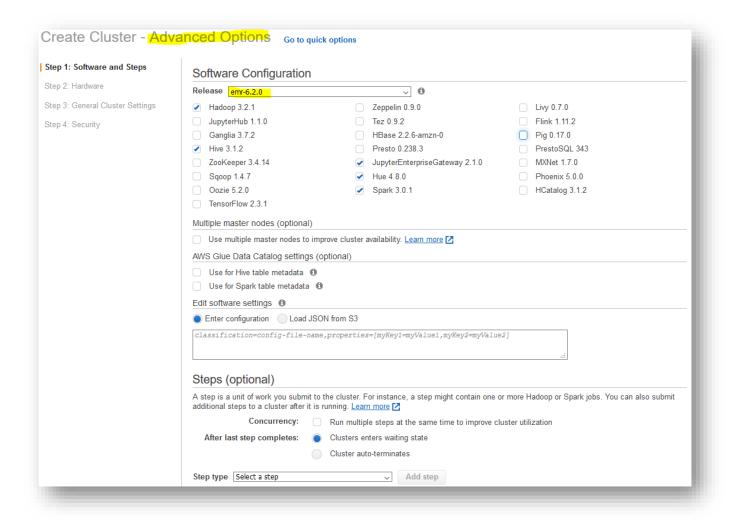- Run steps in the cluster

Submitting jobs
- Connect to the master node over a secure connection
- Access the interfaces and tools that are available on your cluster

# Creating the cluster

Using CLI (command line interface)

```
aws emr create-cluster \
    --name "My First EMR Cluster" \
    --release-label emr-5.32.0 \
    --applications Name=Spark \
    --ec2-attributes KeyName=myEMRKeyPairName \
    --instance-type m5.xlarge \
    --instance-count 3 \
    --use-default-roles
```

This is more pragmatic, but there are many options to explore
- Let's stick to AWS Console
- https://console.aws.amazon.com/elasticmapreduce/

# Creating the cluster

# Creating the cluster

# Creating the cluster

# Creating the cluster



Create Cluster - Advanced Options    Go to quick options

Step 1: Software and Steps

Step 2: Hardware

Step 3: General Cluster Settings

| Step 4: Security

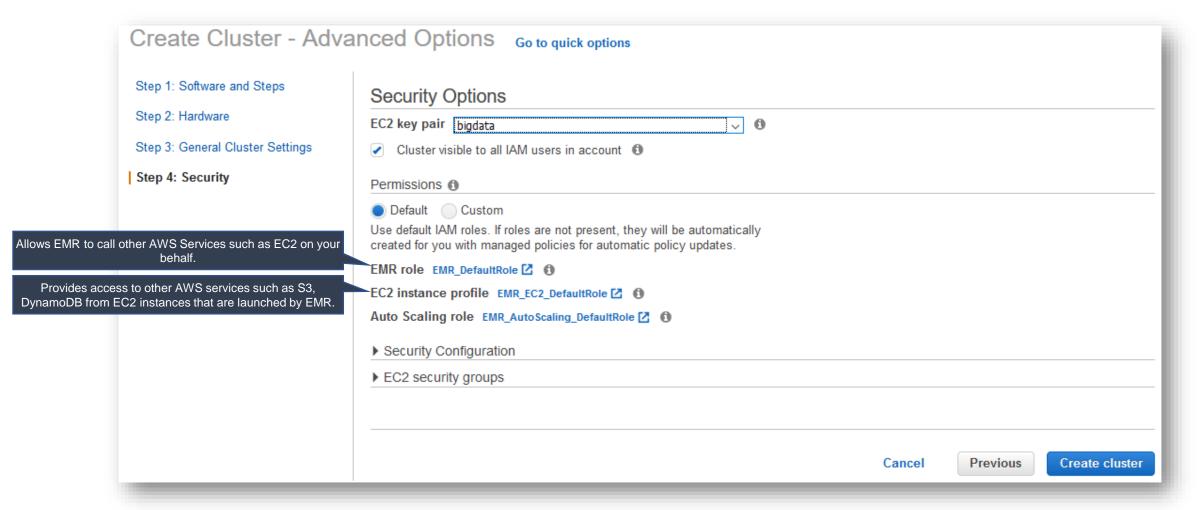**Security Options**

EC2 key pair    bigdata    ⌄    ⓘ

☑ Cluster visible to all IAM users in account    ⓘ

Permissions    ⓘ

🔘 Default    ⚪ Custom
Use default IAM roles. If roles are not present, they will be automatically
created for you with managed policies for automatic policy updates.

> Allows EMR to call other AWS Services such as EC2 on your behalf.

EMR role    EMR_DefaultRole ↗    ⓘ

> Provides access to other AWS services such as S3, DynamoDB from EC2 instances that are launched by EMR.

EC2 instance profile    EMR_EC2_DefaultRole ↗    ⓘ

Auto Scaling role    EMR_AutoScaling_DefaultRole ↗    ⓘ

▶ Security Configuration

▶ EC2 security groups

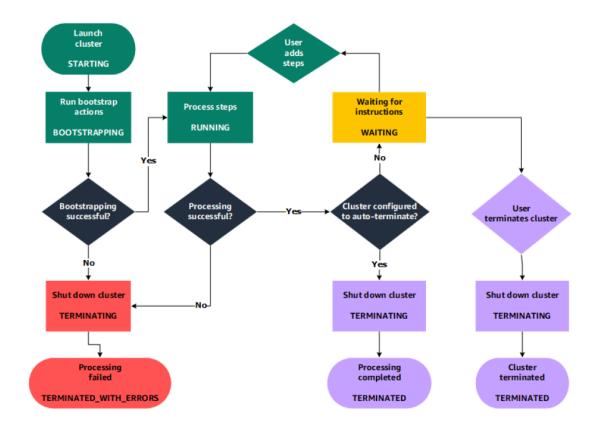Cancel    Previous    **Create cluster**

# Creating the cluster

Using CLI (command line interface)

```
aws emr create-cluster --auto-scaling-role EMR_AutoScaling_DefaultRole --termination-protected --
applications Name=Hadoop Name=Hive Name=Hue Name=JupyterEnterpriseGateway Name=Spark --ebs-root-volume-
size 10 --ec2-attributes
'{"KeyName":"bigdata","InstanceProfile":"EMR_EC2_DefaultRole","SubnetId":"subnet-
5fa2f912","EmrManagedSlaveSecurityGroup":"sg-07818b5690a50b3f1","EmrManagedMasterSecurityGroup":"sg-
0e2f5550a2cb98f79"}' --service-role EMR_DefaultRole --enable-debugging --release-label emr-6.2.0 --log-
uri 's3n://aws-logs-604905954159-us-east-1/elasticmapreduce/' --name 'BigData' --instance-groups
'[{"InstanceCount":1,"BidPrice":"OnDemandPrice","EbsConfiguration":{"EbsBlockDeviceConfigs":[{"VolumeSpe
cification":{"SizeInGB":32,"VolumeType":"gp2"},"VolumesPerInstance":2}]},"InstanceGroupType":"MASTER","I
nstanceType":"m4.xlarge","Name":"Master -
1"},{"InstanceCount":1,"BidPrice":"OnDemandPrice","EbsConfiguration":{"EbsBlockDeviceConfigs":[{"VolumeS
pecification":{"SizeInGB":32,"VolumeType":"gp2"},"VolumesPerInstance":2}]},"InstanceGroupType":"CORE","I
nstanceType":"m4.xlarge","Name":"Core - 2"}]' --scale-down-behavior TERMINATE_AT_TASK_COMPLETION --
region us-east-1
```

# Cluster lifecycle

Creating a cluster (it takes ~10 minutes)

- A cluster cannot be stopped
- It can only be terminated

# Cluster lifecycle

`STARTING`: EMR provisions EC2 instances for each required instance

`BOOTSTRAPPING`: EMR runs actions that you specify on each instance

- E.g., install custom applications and perform customizations

Amazon EMR installs the native applications

- E.g., Hive, Hadoop, Spark, and so on

`RUNNING`: a step for the cluster is currently being run

- Cluster sequentially runs any steps that you specified when you created the cluster

`WAITING`: after steps run successfully

`TERMINATING`: after manual shut down
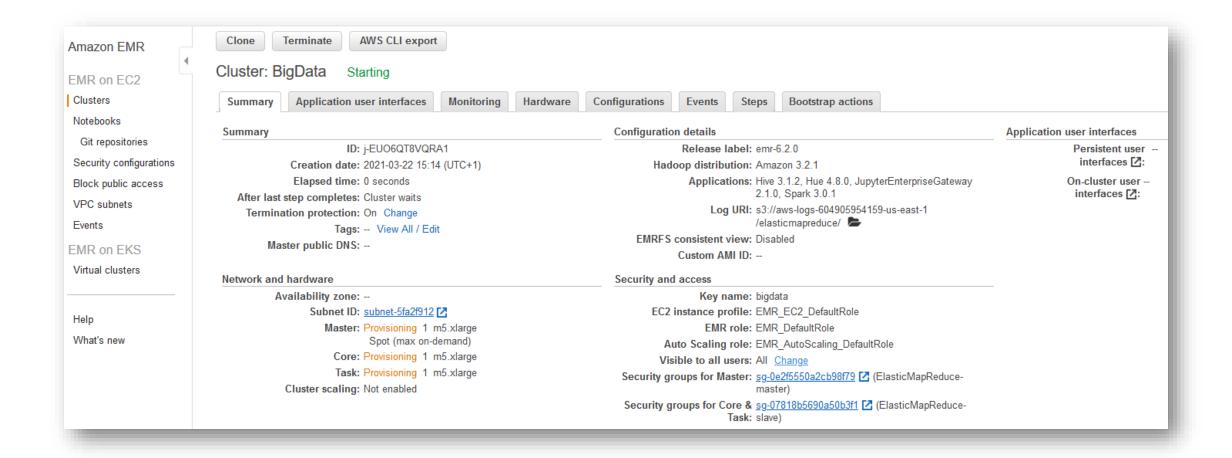
- Any data stored on the cluster is deleted

# Cluster: EMR

A **step** is a user-defined unit of processing

- E.g., one algorithm that manipulates the data

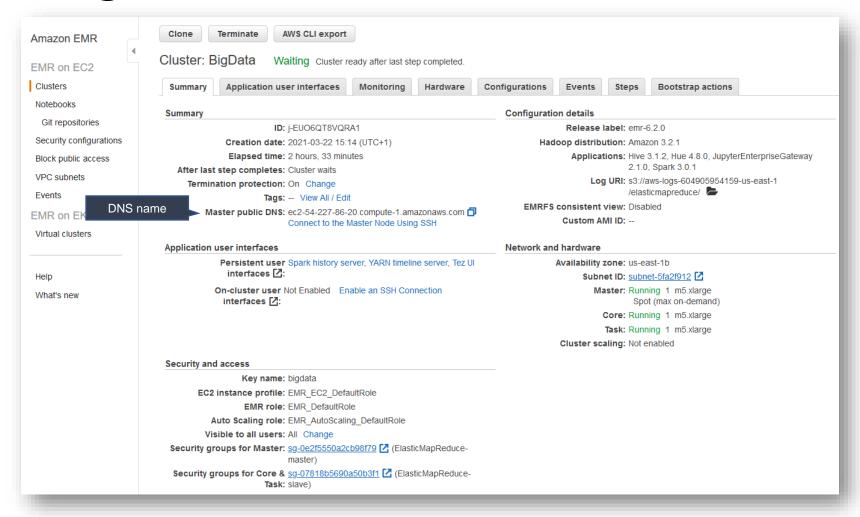Step states

- `PENDING:` The step is waiting to be run
- `RUNNING:` The step is currently running
- `COMPLETED:` The step completed successfully
- `CANCELLED:` The step was cancelled before running because an earlier step failed
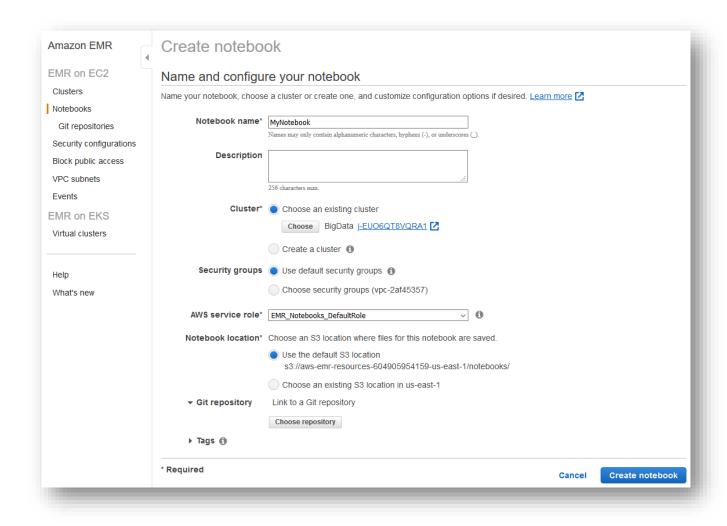- `FAILED:` The step failed while running
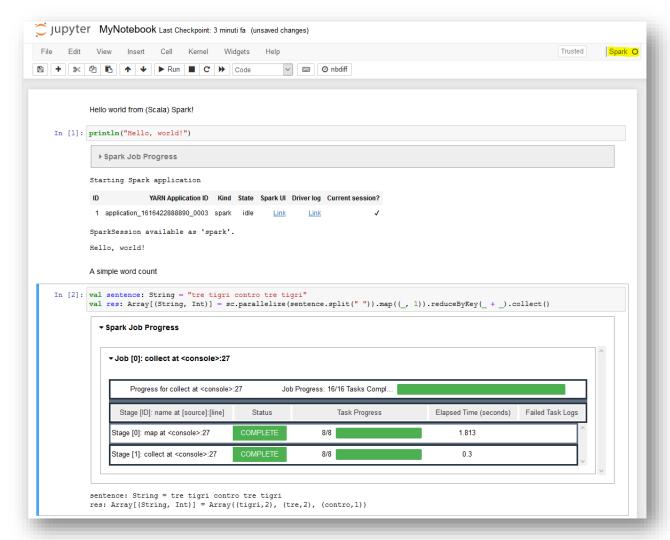
# Running the cluster

# Running the cluster

# Creating a notebook

# Hello, world!

# Add some storage

Save the result to HDFS

```scala
import org.apache.hadoop.fs.{FileSystem, Path}
val fs = FileSystem.get(sc.hadoopConfiguration) // get the file system
val outPutPath = new Path(path)
if (fs.exists(outPutPath)) { // delete the HDFS folder if exists
    fs.delete(outPutPath, true)
}

val hdfspath: String = "wordcount" // HDFS path
def writeandread(path: String) = {
    sc.parallelize(res).saveAsTextFile(path) // save the RDD
    val rdd = sc.textFile(path) // read it back
    rdd.collect() // print it
}

writeandread(hdfspath)
```

▸ Spark Job Progress

```
import org.apache.hadoop.fs.{FileSystem, Path}
fs: org.apache.hadoop.fs.FileSystem = DFS[DFSClient[clientName=DFSClient_NONMAPREDUCE_1600703682_22, ugi=livy (auth:SIMPL
E)]]
outPutPath: org.apache.hadoop.fs.Path = wordcount
res28: AnyVal = true
hdfspath: String = wordcount
writeandread: (path: String)Array[String]
res32: Array[String] = Array((tigri,2), (tre,2), (contro,1))
```

... and to S3 as well

```scala
val s3bucket: String = "s3://aws-emr-resources-604905954159-us-east-1/wordcount"
writeandread(s3bucket)
```

# Running a Spark Job

Connect using SSH

Install git

Clone & build the project

```
ssh -i ~/bigdata.pem hadoop@ec2-54-242-176-32.compute-1.amazonaws.com

sudo yum install git -y

git clone https://github.com/w4bo/spark-word-count.git

cd spark-word-count

./gradlew

spark-submit --class it.unibo.big.WordCount build/libs/WordCount-all.jar
                     s3://aws-bucket-bigdata2021/inferno.txt
```
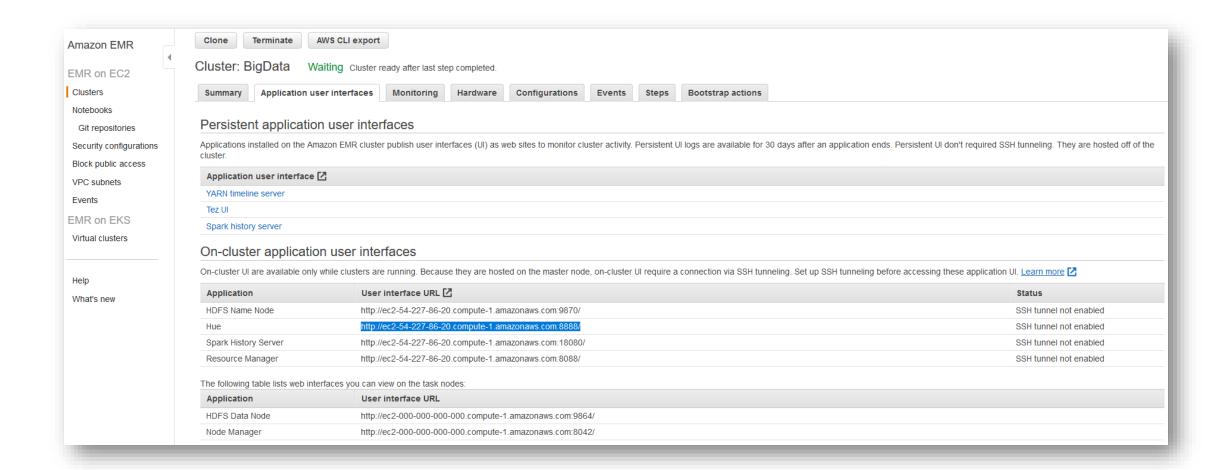
# Other services: HUE

## Connecting to Hue

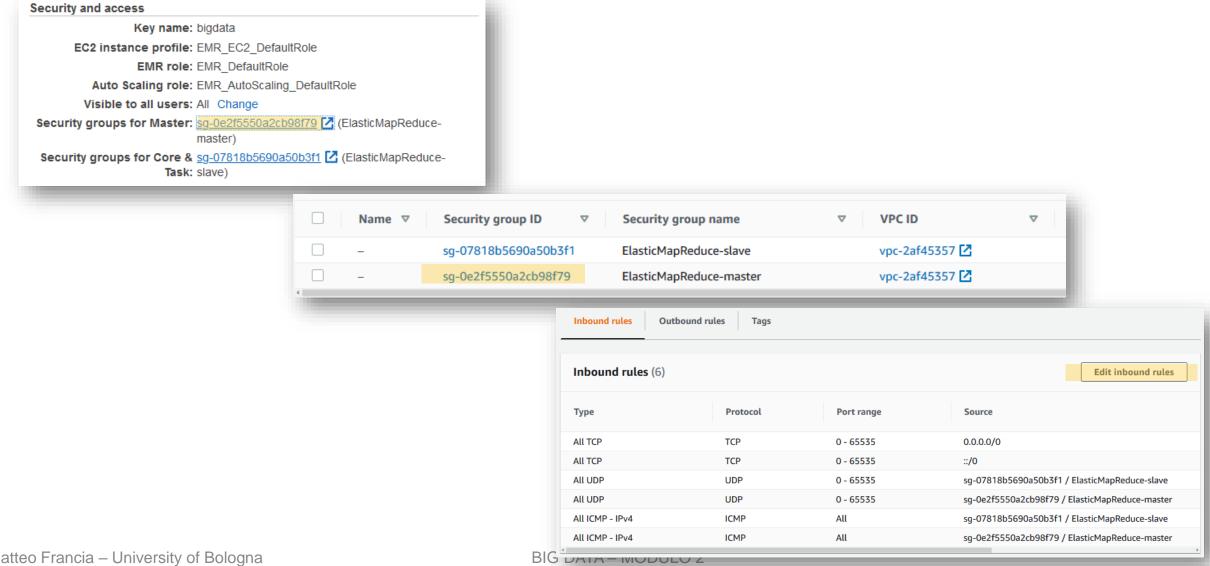- I.e., connecting to any HTTP interface hosted on the master node of a cluster

## To view the Hue web user interface

- Set Up an SSH Tunnel to the Master Node Using Dynamic Port Forwarding
- Type the following address in your browser to open the Hue web interface
  - [http://master-public-DNS:8888](http://master-public-DNS:8888)
  - Where master-public-DNS is the public DNS name of the master node
- If you are the administrator logging in for the first time
  - Enter a username and password to create your Hue superuser account
  - Otherwise, type your username and password and select Create account

# Other services: HUE

# Set Up an SSH Tunnel

# Connect to HUE

**Application user interfaces**

> **Persistent user** Spark history server, YARN timeline server, Tez UI
> interfaces ↗:
>
> **On-cluster user** HDFS Name Node, Hue, Spark History Server,
> interfaces ↗: Resource Manager

## On-cluster application user interfaces

On-cluster UI are available only while clusters are running. Because they are hosted on the master node, on-cluster UI require a connection via SSH tunneling. Set up SSH tunneling before accessing these application UI. Learn more ↗

| Application | User interface URL ↗ | Status |
|---|---|---|
| HDFS Name Node | http://ec2-54-242-176-32.compute-1.amazonaws.com:9870/ | Available |
| Hue | http://ec2-54-242-176-32.compute-1.amazonaws.com:8888/ | Available |
| Spark History Server | http://ec2-54-242-176-32.compute-1.amazonaws.com:18080/ | Available |
| Resource Manager | http://ec2-54-242-176-32.compute-1.amazonaws.com:8088/ | Available |

# Connect using SSH