Work on making a CI-CD (Continuous Integration and Continuous Delivery) the task to submit a spark word count program from S3 to the EMR cluster using AWS code pipeline.

**DevOps**

A set of practices intended to reduce the time between committing a change to a system and the change being placed into normal production, while ensuring high quality. is the combination of practices and tools designed to increase an organization's ability to deliver applications and services faster than traditional software development processes.

**Continuous Integration and Delivery**

Step 1 – Split the entire chunk of codes into segments

Step 2 – keep small segments of manageable code

Step 3 – integrate the segmented code, multiple times a day

Step 4 – Adopt a continuous integration methodology to coordinate with your team

we have a source code repository where the developers work continuously submit their pieces of the code repository. Such that a central place where the changes are constantly committed then we have a belt server where everything is compiled reviewed tested integrated and then packets as well finally started test final test goes to the minorities and then it goes to the production environment, where this process the building the Staging takes palace and the committee process it

automating this process becomes very important and if we will do it manually  we will suffer a lot

**AWS Codepipeline**

AWS CodePipeline is a continuous delivery service you can use to model, visualize, and automate the steps required to release your software.
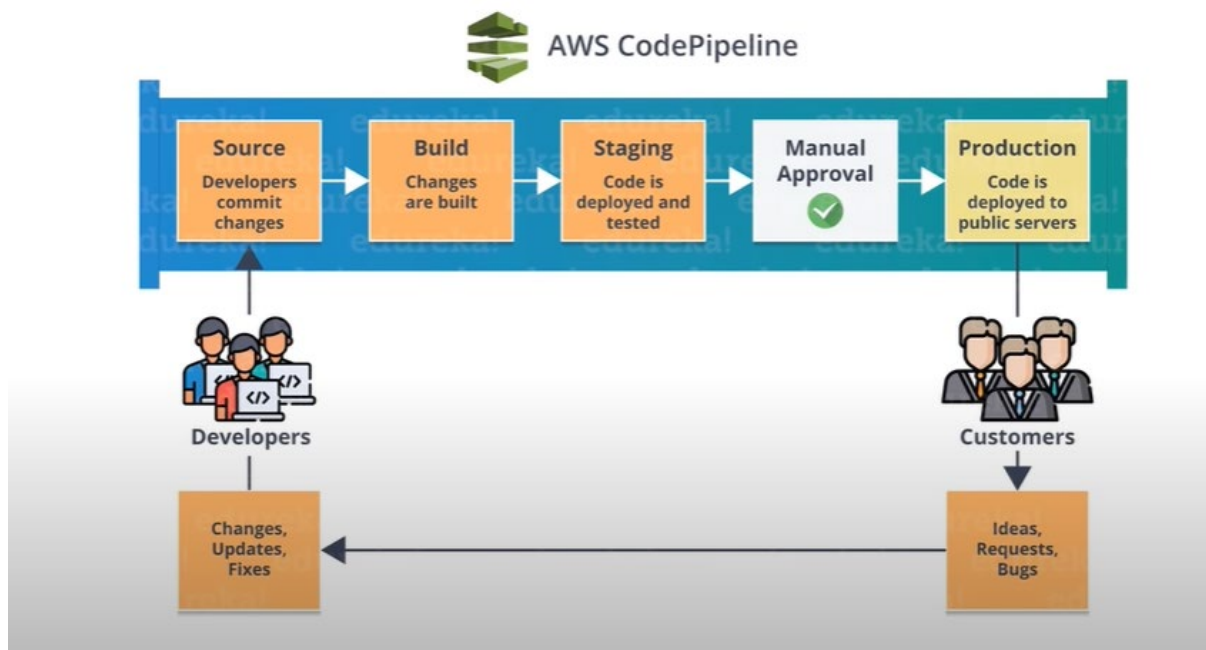
Its features -

Monitor your processes in real-time

Ensure Consistent Release Process

Speed up delivery while improving quality

View pipeline history details

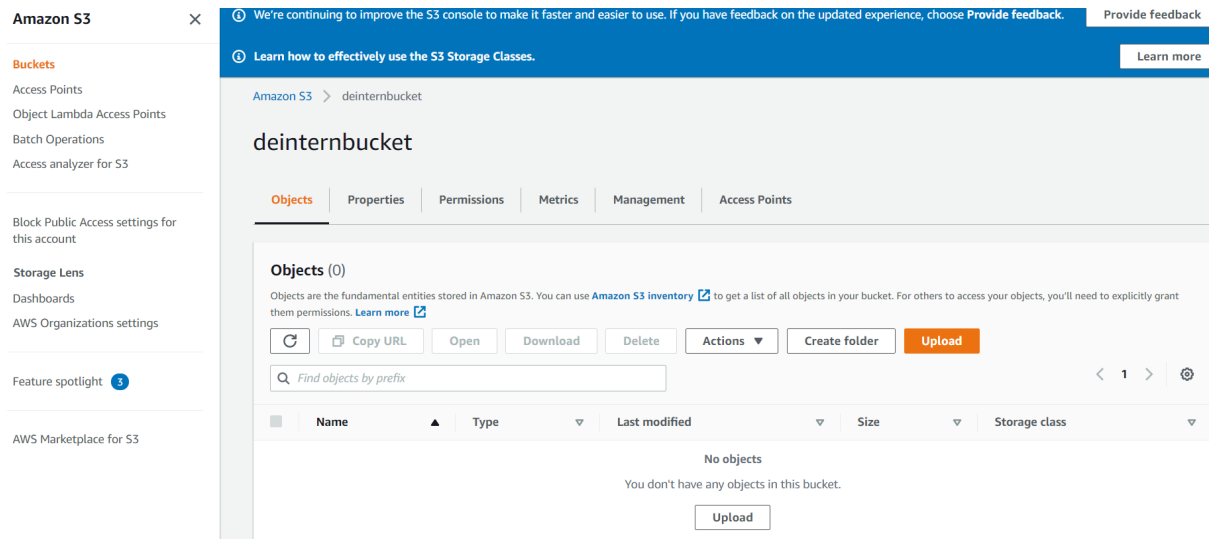The architecture of the code pipeline looks like this

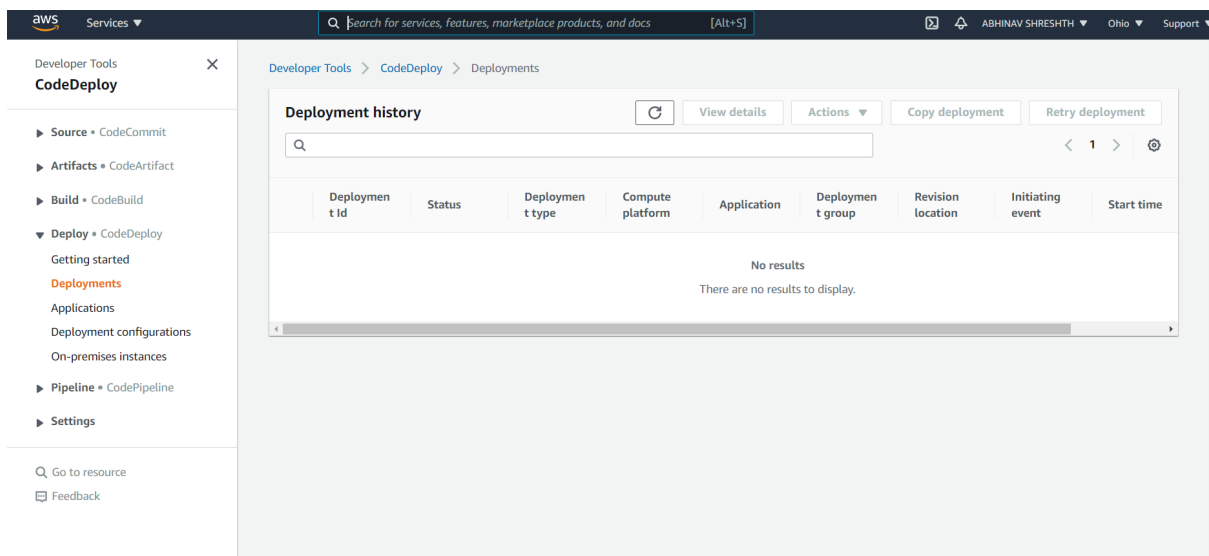Further, I have signed up for the Amazon AWS console.

Then I opened AWS EMR and then as I don't know how to proceed further then I YouTube it to know how to proceed



On AWS s3  I created bucket there

further



## Amazon EMR

Amazon EMR cluster provides a managed Hadoop framework that makes it easy, fast, and cost-effective to process vast amounts of data across dynamically scalable Amazon EC2 instances.

We can also run other popular distributed frameworks such as Apache **Spark** and HBase in Amazon EMR, and interact with data in other AWS data stores such as Amazon S3 and Amazon DynamoDB.

## S3(Amazon Simple Storage Service)

*We are going to run our Spark application on top of the Hadoop cluster, and we will put the input data source into the S3.*

S3 is a distributed storage system and AWS's equivalent to HDFS. We because we want to make sure that

- Our data is coming from some distributed file system that can be accessed by every node on our Spark cluster.

- Our Spark application doesn't assume that our input data sits somewhere on our local disk because that will not scale. By saving our input data source into S3, each spark node deployed on the EMR cluster can read the input data source from S3.



To create an Amazon EC2 key pair:

1. Go to the Amazon EC2 console
2. In the Navigation pane, click Key Pairs

3. On the Key Pairs page, click Create Key Pair
4. In the Create Key Pair dialogue box, enter a name for your key pair, such as mykeypair
5. Click Create
6. Save the resulting PEM file in a safe location





After the activation of the account under **Security and access -> Security groups for Master**

Click on it then this window will open

There create a new rule of an inbound connection



Next step – download putty application

Open the Putty key generator then load the key you have downloaded
before from aws

**Waiting** Cluster ready after last step completed.

er interfaces | Monitoring | Hardware | Configurations | Events | Steps | Bootstrap actions

**SSH**

## Connect to the Master Node Using SSH

j-1DAIAF16L3Y

: 2021-05-22 18

: 1 hour

: Cluster waits

: Off Change

: -- View All / E

compute.amazo
sing SSH

: YARN timeline

Not Enabled E

You can connect to the Amazon EMR master node using SSH to run interactive q
Learn more 🔗.

| Windows | Mac / Linux |

1. Download PuTTY.exe to your computer from:
   http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html 🔗
2. Start PuTTY.
3. In the Category list, click Session.
4. In the Host Name field, type **hadoop@ec2-18-220-249-246.us-east-2.co**
5. In the Category list, expand Connection > SSH, and then click Auth.
6. For Private key file for authentication, click Browse and select the private
7. Click Open.
8. Click Yes to dismiss the security alert.

**PuTTY Configuration**                     ?

Category:
- Features
- Window
  - Appearance
  - Behaviour
  - Translation
  - Selection
  - Colours
- Connection
  - Data
  - Proxy
  - SSH
    - Kex
    - Host keys
    - Cipher
    - Auth
    - TTY
    - X11
    - Tunnels
    - Bugs
    - More bugs
  - Serial
  - Telnet
  - Rlogin
  - SUPDUP

Basic options for your PuTTY session

Specify the destination you want to connect to
Host Name (or IP address)                     Port
'49-246.us-east-2.compute.amazonaws.com   22

Connection type:
◉ SSH  ○ Serial  ○ Other:  Telnet

Load, save or delete a stored session
Saved Sessions

Default Settings                              Loa
                                             Sav
                                             Dele

Close window on exit
○ Always  ○ Never  ◉ Only on clean exit

About    Help              Open    Canc

: mykeypair

```
Now I am inside EMR cluster
```



```
[hadoop@ip-172-31-21-213 ~]$ ls
[hadoop@ip-172-31-21-213 ~]$ aws s3 cp s3:
Note: AWS CLI version 2, the latest major version of the AWS CLI, is now stable
and recommended for general use. For more information, see the AWS CLI version 2
 installation instructions at: https://docs.aws.amazon.com/cli/latest/userguide/
install-cliv2.html


usage: aws [options] <command> <subcommand> [<subcommand> ...] [parameters]
To see help text, you can run:

  aws help
  aws <command> help
  aws <command> <subcommand> help
aws: error: too few arguments
[hadoop@ip-172-31-21-213 ~]$
```