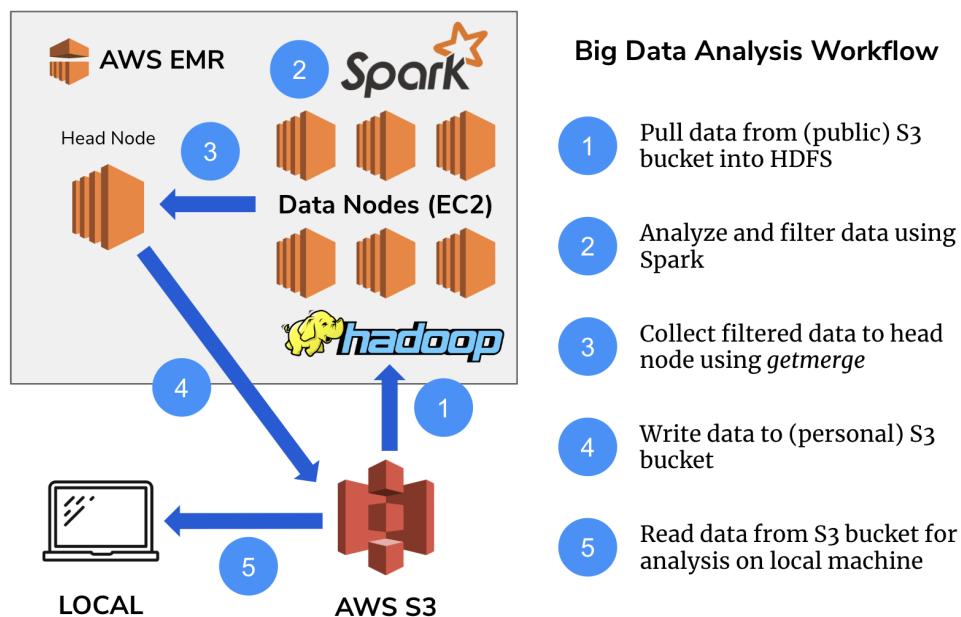# Big Data Wrangling With Google Books Ngrams

In this assignment, I will apply the skills I've learned in the Big Data Fundamentals unit to load, filter, and visualize a large real-world dataset in a cloud-based distributed computing environment using Hadoop, Spark, Hive, and the S3 filesystem. I will prepare a professional report to summarize my findings and include an appendix with screenshots of the steps completed to setup EMR Cluster

The Google Ngrams dataset was created by Google's research team by analyzing all of the content in Google Books. These digitized texts represent approximately 4% of all books ever printed, spanning a time period from the 1800s into the 2000s.

The dataset is hosted in a public S3 bucket as part of the Amazon S3 Open Data Registry. For this assignment, the data has been converted to CSV and hosted on a public S3 bucket, which can be accessed here: s3://brainstation-dsft/eng_1M_1gram.csv.

As part of this workflow, I will filter and reduce the data down to a manageable size, and then perform some analysis locally on my machine after extracting data from the cloud and processing it using Big Data tools. The workflow and steps in the process are illustrated below:
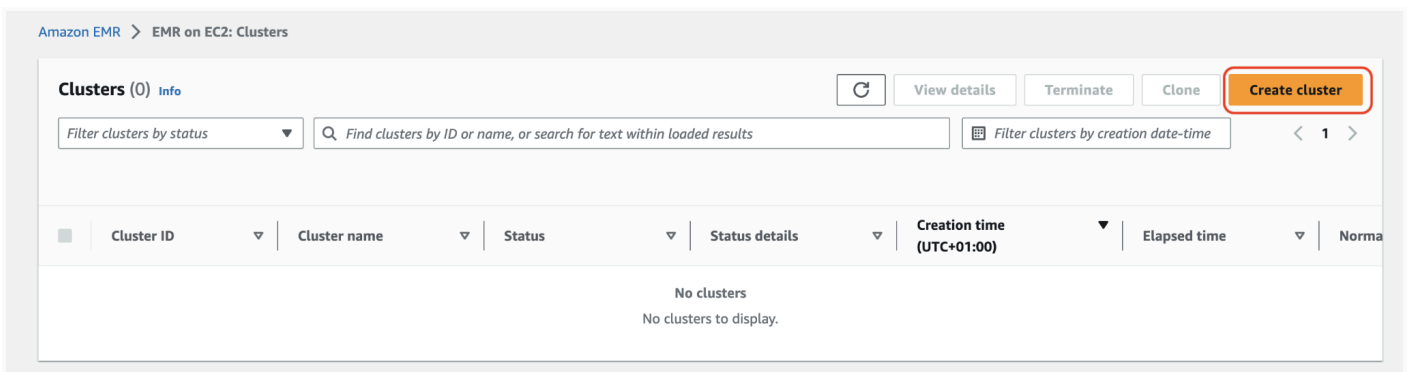
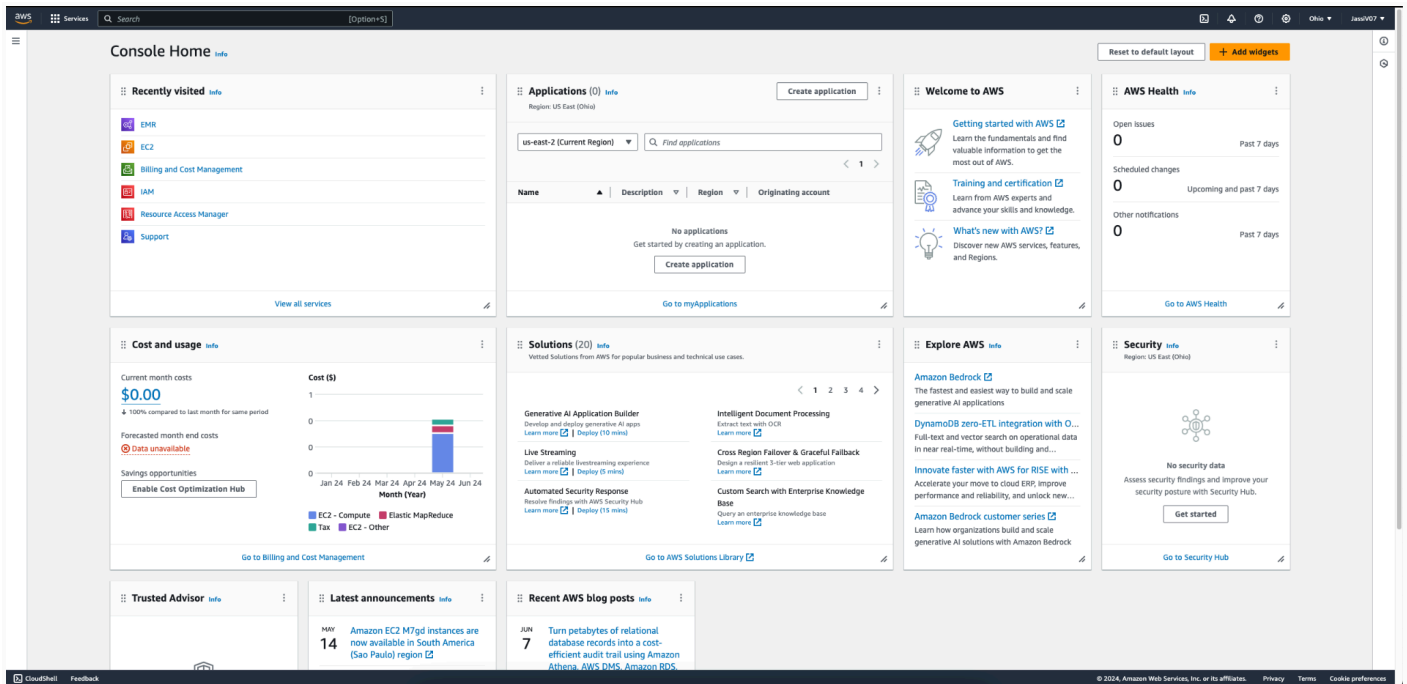O1. <u>Spin up a new EMR cluster on AWS for using Spark and EMR notebooks</u>
      a. Go to https://aws.amazon.com and sign in to your account.
      b. Navigate to the EMR panel in AWS and click 'Create Cluster':

c. Name the Cluster and Select Applications.
   i. Give your cluster a name.
   ii. In the 'Release' dropdown, select emr-6.10.0.
   iii. Select the Custom application bundle, and tick the boxes for Hadoop, Hue, JupyterHub, Livy, Hive, and Spark.

Name

my_Spark_cluster

Amazon EMR release   Info
A release contains a set of applications which can be installed on your cluster.

emr-6.10.0 ▼

Application bundle

| Spark | Core Hadoop | HBase | Presto | Trino | Custom |
|-------|-------------|-------|--------|-------|--------|
| Spark | hadoop | HBASE | presto | trino | aws |

▼ Customise your application bundle

Applications included in bundle

| ☐ Flink 1.16.0 | ☐ Ganglia 3.7.2 | ☐ HBase 2.4.15 |
| ☐ HCatalog 3.1.3 | ☑ Hadoop 3.3.3 | ☑ Hive 3.1.3 |
| ☑ Hue 4.10.0 | ☐ JupyterEnterpriseGateway 2.6.0 | ☑ JupyterHub 1.5.0 |
| ☑ Livy 0.7.1 | ☐ MXNet 1.9.1 | ☐ Oozie 5.2.1 |
| ☐ Phoenix 5.1.2 | ☐ Pig 0.17.0 | ☐ Presto 0.278 |
| ☑ Spark 3.3.1 | ☐ Sqoop 1.4.7 | ☐ TensorFlow 2.11.0 |
| ☐ Tez 0.10.2 | ☐ Trino 403 | ☐ Zeppelin 0.10.1 |
| ☐ ZooKeeper 3.5.10 | | |

> d. Instance Group
>> i. Remove the Task instance group

○ **Instance groups**
Choose one instance type per node group

○ **Instance fleets**
Choose any combination of instance types within each node group

## Instance groups

**Primary**

Choose EC2 instance type

m5.xlarge
4 vCore    16 GiB memory    EBS only storage
On-demand price: -
Lowest spot price: $0.099 (us-east-1d)

Actions ▼

☐ Use multiple primary nodes
To improve cluster availability, use three primary nodes with the same configuration and bootstrap actions. You cannot use multiple primary nodes with instance fleets.

▶ **Node configuration - *optional***

**Core**

Choose EC2 instance type

m5.xlarge
4 vCore    16 GiB memory    EBS only storage
On-demand price: -
Lowest spot price: $0.099 (us-east-1d)

Actions ▼

▶ **Node configuration - *optional***

**Task 1 of 1**

**Remove instance group**

Name

Task - 1

Choose EC2 instance type

ii. Allocate 2 Nodes to the core instance group

**Cluster scaling and provisioning option** Info

The Amazon EMR console only supports EMR-managed scaling. To create a cluster with auto-scaling, use CLI or SDK.

○ **Set cluster size manually**
Use this option if you know your workload patterns in advance.

○ **Use EMR-managed scaling**
Monitor key workload metrics so that EMR can optimise the cluster size and optimise its resource utilisation.

| Name | Instance type | Size | | Use spot purchasing option |
|------|---------------|------|--|---------------------------|
| Core | m5.xlarge | 2 | instance(s) | ☐ |

e. Cluster Termination
  i. Set cluster termination to 4h idle time.
  ii. Turn termination protection off.

**Cluster termination** Info

○ Manually terminate cluster
● Terminate cluster after idle time (recommended)

Idle time
Enter the time until your cluster terminates.

| 0 days ▼ | 04:00:00 |

Choose a time that is greater than 1 minute (00:01:00) and less than 7 days. The time is in hh:mm:ss (24-hour) format.

☐ Terminate cluster after last step has been completed
You cannot edit this selection after you have created your cluster.

☐ Use termination protection
Protect your EC2 instances from accidental termination.

f.  Security and Access Management
    i.   Select your key pair (keys are associated to geographies so if you switched recently, you might need to create a new key pair).
    ii.  In Identity and Access Management, choose the EMR_DefaultRole and EMR_EC2_DefaultRole

Unit 4 Deliverable 1
Jashkirat Virdi

**Note** Once you create the cluster the process will start but it will take some time for it to get fully loaded and all of the things to get configured, the state will change from Starting to Waiting.

    02.Connect to the head node of the cluster using SSH
        a.  Connect to the Primary Node of the cluster, and access JupyterHub in a browser window, edit the following bash command:

        b.  Replace 'xxxxxxxxx' with your 'Primary node public DNS' found on the overview page of your cluster.

```
ssh -i mykey.pem -L 9995:localhost:9443
hadoop@xxxxxxxxxxxxxx.compute.amazonaws.com
```

## unit_4_del

Updated 7 minutes ago | Terminate | Clone in AWS CLI | Clone

▼ **Summary**

| Cluster info | Applications | Cluster management | Status and time |
|---|---|---|---|
| **Cluster ID** <br> j-W83ESHBFJRN6 | **Amazon EMR version** <br> emr-6.10.0 | **Log destination in Amazon S3** <br> aws-logs-992382776059-us-east-2/elasticmapreduce | **Status** <br> ⊖ Terminated |
| **Cluster configuration** <br> Instance groups | **Installed applications** <br> Hadoop 3.3.3, Hive 3.1.3, Hue 4.10.0, JupyterHub 1.5.0, Livy 0.7.1, Spark 3.3.1 | **Persistent application UIs** <br> Spark History Server ↗ | **Creation time** <br> June 09, 2024, 17:39 (UTC-04:00) |
| **Capacity** <br> 1 Primary 1 Core 1 Task | | YARN timeline server ↗ <br><br> Tez UI ↗ | **Elapsed time** <br> 4 hours, 41 minutes |
| | | **Primary node public DNS** <br> ▢ ec2-18-221-130-68.us-east-2.compute.amazonaws.com <br><br> Connect to the Primary node using SSH | **End time** <br> June 09, 2024, 22:20 (UTC-04:00) |

c. Paste the command in the terminal

```
sh -i brainstation.pem -L 9995:localhost:9443 hadoop@ec2-18-221-130-68.us-east-2.compute.amazonaws.com
t login: Sun Jun  9 22:51:49 2024
      #_
~\_  ####_          Amazon Linux 2
~  \_#####\
~    \###|          AL2 End of Life is 2025-06-30.
~     \#/ ___
~~     V~' '->
 ~~~        /       A newer version of Amazon Linux is available!
  ~._.   _/
   _/ _/           Amazon Linux 2023, GA and supported until 2028-03-15.
  _/m/'               https://aws.amazon.com/linux/amazon-linux-2023/


EEEEEEEEEEEEEEEEEE MMMMMMMM            MMMMMMMM RRRRRRRRRRRRRRRR
::::::::::::::::E M:::::::M          M:::::::M R::::::::::::::R
::::EEEEEEEEE:::E M:::::::M         M:::::::M R:::::RRRRRR:::::R
::::E       EEEEE M:::::::::M      M:::::::::M RR::::R     R::::R
::::E            M:::::M:::M    M:::M:::::::M  R:::R      R::::R
::::EEEEEEEEEE   M:::::M M:::M M:::M M::::::M  R:::RRRRRR:::::R
::::::::::::::E   M:::::M  M:::M:::M  M:::::M  R:::::::::::RR
::::EEEEEEEEEE   M:::::M   M:::::M   M:::::M   R:::RRRRRR::::R
::::E           M:::::M    M:::M    M:::::M   R:::R     R::::R
::::E     EEEEE M:::::M     MMM     M:::::M   R:::R     R::::R
::::EEEEEEEE::::E M:::::M           M:::::M   R:::R     R::::R
::::::::::::::::E M:::::M           M:::::M RR::::R     R::::R
EEEEEEEEEEEEEEEEEE MMMMMMM          MMMMMMM RRRRRRR     RRRRRR

doop@ip-172-31-42-62 ~]$ █
```

03. Copy the data folder from the S3 bucket directly into a directory on the Hadoop File System (HDFS) named /user/hadoop/eng_1M_1gram.

a. Paste the following command in the terminal

```
hadoop distcp s3://brainstation-dsft/eng_1M_1gram.csv /user/hadoop/eng_1M_1gram
```

b. Check /user/hadoop/eng_1M_1gram

```
[hadoop@ip-172-31-42-62 ~]$ hadoop fs -ls
Found 3 items
drwxr-xr-x   - hadoop hdfsadmingroup          0 2024-06-09 23:32 .sparkStaging
-rw-r--r--   1 hadoop hdfsadmingroup 5292105197 2024-06-09 22:24 eng_1M_1gram
```

c. Access JupyterHub in a browser window at `*https://localhost:9995*`

d. Login using username `jovyan` and password `jupyter`



e. Open a New PySpark Notebook

**04.  Check Step_4_Big_Data.ipynb Notebook for the Below.**

    **a.** Describe the dataset (examples include size, shape, schema) in pyspark
    b. Create a new DataFrame from a query using Spark SQL, filtering to include only
        the rows where the token is "data" and describe the new dataset
    c. Write the filtered data back to a directory in the HDFS from Spark using
        df.write.csv().

05.Collect the contents of the directory into a single file on the local drive of the head
    node using getmerge and move this file into a S3 bucket in your account.
    a. Check hadoop directory using hadoop fs –ls

```
[hadoop@ip-172-31-42-62 ~]$ hadoop fs -ls /user/hadoop/
Found 3 items
drwxr-xr-x   - hadoop hdfsadmingroup          0 2024-06-09 23:32 /user/hadoop/.sparkStaging
-rw-r--r--   1 hadoop hdfsadmingroup 5292105197 2024-06-09 22:24 /user/hadoop/eng_1M_1gram
drwxr-xr-x   - livy   hdfsadmingroup          0 2024-06-09 23:25 /user/hadoop/filtered_data
```

    b. Use hadoop fs –getmerge to merge the contents into a single file on the local
        drive of the head node.
    c. Move the file into an S3 bucket in your account using aws s3 cp.

```
[hadoop@ip-172-31-42-62 ~]$ aws s3 cp filtered_data.csv  s3://aws-emr-studio-992382776059-us-east-2/171796
5722053/e-9BBQW3CFDNRYZC7JG8NQ3YHG9/
upload: ./filtered_data.csv to s3://aws-emr-studio-992382776059-us-east-2/1717965722053/e-9BBQW3CFDNRYZC7J
G8NQ3YHG9/filtered_data.csv
```
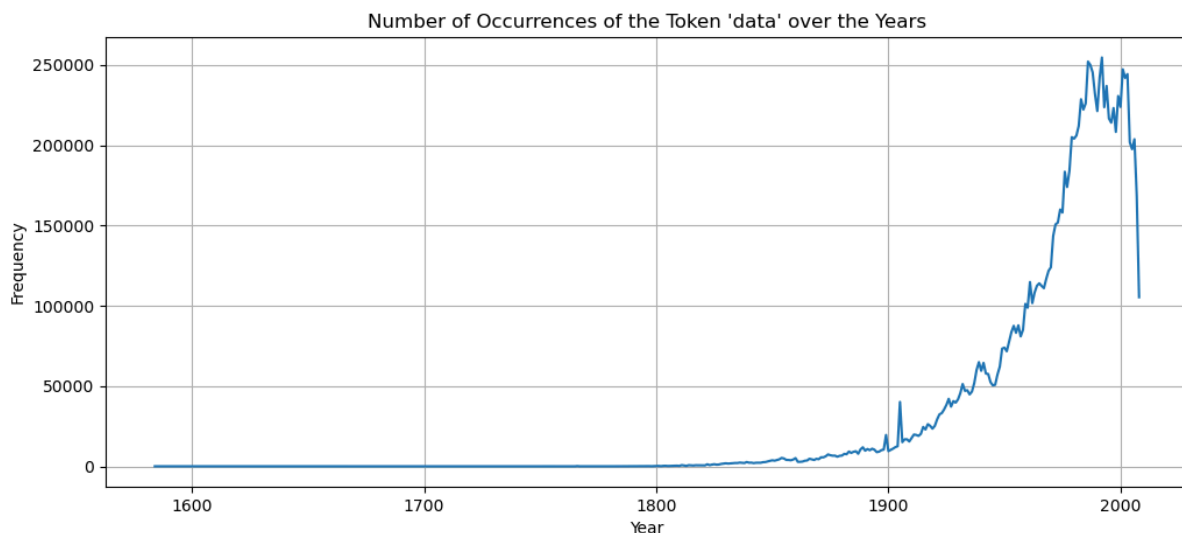
Unit 4 Deliverable 1
Jashkirat Virdi

## Check Step_6_Big_Data.ipynb Notebook for the Below.

06. On your local machine in python, read the CSV data from the S3 folder into a pandas DataFrame Note you must have first authenticated on your machine using aws configure on the command line to complete this step.

    a. Open Jupyter Notebook and import necessary libraries
    b. Read Data into dataframe

07. Plot the number of occurrences of the token (the frequency column) of data over the years using matplotlib.

    a. Data Wrangling
    b. Plot number of occurrences of the token `data` over the years



*The heavy usage of the word "data" in books began in the mid-1800s and increased exponentially, reaching a peak between the 1990s and 2000s, likely due to the absence of more recent data.*

O8.Compare Hadoop and Spark as distributed file systems.

    a. What are the advantages/ differences between Hadoop and Spark? List two advantages for each.

        i. **Advantages of Hadoop:**

            1. **Storage Capability (HDFS):** Hadoop Distributed File System (HDFS) allows for the storage and processing of large data sets across distributed clusters. It ensures high fault tolerance and reliability, making it suitable for storing vast amounts of unstructured data.

            2. **Cost-Effective:** Hadoop is often more cost-effective for large-scale batch processing due to its open-source nature and compatibility with commodity hardware, reducing the need for expensive proprietary solutions

        ii. **Advantages of Spark:**

            1. **Speed:** Spark processes data much faster than Hadoop MapReduce due to its in-memory computation capabilities. It can cache data in memory, which significantly speeds up iterative algorithms and interactive data analysis.

            2. **Ease of Use:** Spark provides high-level APIs in Java, Scala, Python, and R, making it more accessible for developers. It also includes built-in libraries for streaming, machine learning, and graph processing, simplifying complex data processing tasks.

    b. Explain how the HDFS stores the data.

HDFS (Hadoop Distributed File System) stores data by splitting large files into smaller blocks, typically 128MB or 256MB each. These blocks are distributed across multiple nodes in a cluster for parallel processing. Each block is replicated across multiple nodes (usually three) to ensure fault tolerance and reliability. The NameNode manages the metadata and directory structure, while DataNodes store the actual data blocks. This design allows HDFS to handle large-scale data storage with high availability and resilience.