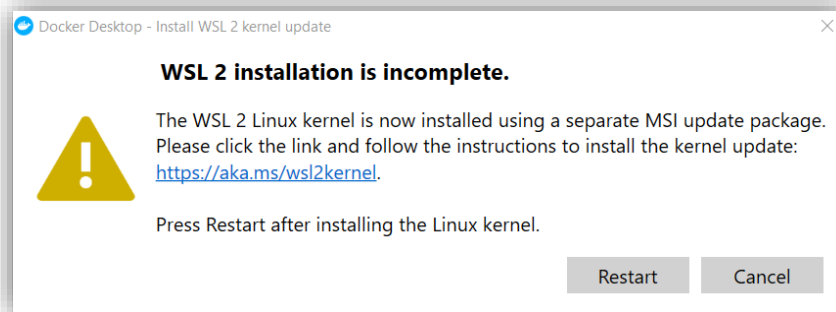


## 1) HDFS

- a) Docker - download & install from this URL:

[Docker Desktop for Mac and Windows | Docker](#)

*Windows: In case the following exception occurs, follow the link in the message and then download and install WSL2:*



- b) Clone Hadoop+Hive repository: [GitHub - FreeUniDataEngineering/hadoop\\_hive](#)

- c) Commands **docker-compose up** or **docker-compose up -d** to build, create and launch containers

- d) Hadoop components - short overview

<https://aws.amazon.com/emr/details/hadoop/what-is-hadoop/>

<https://www.geeksforgeeks.org/hadoop-history-or-evolution/>

- e) Docker commands: **docker ps**, **docker container list**, **docker container stats**

- f) **docker exec -ti CONTAINER bash**

**ls**, **pwd**, **cd**, **rm**, **cat**, **echo 'hello, world!'** > **hello.txt** commands

- g) Copy files from container to local fs:

**docker cp [OPTIONS] CONTAINER:SRC\_PATH DEST\_PATH**

**docker cp 059af7d2f810:./etc/hosts C:/Users/Gigi/Desktop/hd/hosts**

Copy files from local fs to container:

**docker cp [OPTIONS] SRC\_PATH CONTAINER:DEST\_PATH**

**docker cp C:/Users/Gigi/Desktop/hd/hosts 059af7d2f810:./etc/hosts**

d) HDFS from command-line:

- **help** :)))
- **version**
- **env**
- **dfs**:
  - List all available commands - **help**
  - **hdfs dfs -ls /**
  - **hdfs dfs -mkdir /new\_dir**
  - no **cd** here :)))
  - **hdfs dfs -put /hello.txt /new\_dir/hello\_again.txt**
  - **hdfs dfs -count -h -v /new\_dir**
  - **hdfs dfs -mv /new\_dir/hello\_again.txt /new\_dir/new\_subdir/hello\_again\_moved.txt**
  - **hdfs dfs -rm -r /new\_dir**
  - **head/tail: hdfs dfs -head /new\_dir/hello.txt & hdfs dfs -tail /new\_dir/hello.txt**

e) HDFS Python Client:

- run **pip install hdfs** inside the container named 'edge' (running on port 7777)
- **from hdfs import InsecureClient**
- **InsecureClient('http://namenode:50070', root='/')**
- **list(hdfs\_path, status=False)**
- **walk(hdfs\_path, depth=0, status=False)**
- **content(hdfs\_path)**
- **read(hdfs\_path, offset=0, encoding=None)**
- **makedirs(hdfs\_path)**
- **rename(hdfs\_src\_path, hdfs\_dst\_path)**
- **delete(hdfs\_path, recursive=False)**
- **write(hdfs\_path, data=None, overwrite=False, append=False, blocksize=None, replication=None, encoding=None)**
- **download(hdfs\_path, local\_path, overwrite=False)**
- **upload(hdfs\_path, local\_path)**

More about HDFS Client: [API reference — HdfsCLI 2.5.8 documentation](#)

## 2) Google Cloud Storage:

- UI:
  - Create Buckets
  - Upload, Make Public, and Delete Objects in Your Bucket
  - Enable Version Control - display versions & restore
- Command Line: Gsutil:
  - Needs Google Cloud SDK installed:  
[Installing Cloud SDK | Cloud SDK Documentation | Google Cloud](#)
  - **gsutil ls**
  - **gsutil ls gs://buckety-bucket/**
  - **gsutil cat gs://buckety-bucket/DataLake/testv2.txt**
  - **gsutil cp D:\Workspace\test.txt gs://buckety-bucket/DataLake/test.txt**
  - **gsutil cp gs://buckety-bucket/DataLake/test.txt D:\Workspace\hello\_from\_GCP.txt**
  - **gsutil compose gs://buckety-bucket/DataLake/test1.txt gs://buckety-bucket/DataLake/test2.txt gs://buckety-bucket/DataLake/test3.txt**

## 3) Python Client:

- Python 3.6+ with Pip
- **pip install --upgrade google-cloud-storage**
- Set up env var GOOGLE\_APPLICATION\_CREDENTIALS. Value: location of Private Key
- **from google.cloud import storage**
- **storage.Client()**
- Bucket constructor: **storage\_client.bucket(bucket\_name)**
- **storage\_client.create\_bucket(bucket, location="us")**
- **list\_buckets()**
- **get\_bucket(bucket\_name)**
- **list\_blobs(bucket\_name)**
- **list\_blobs(bucket\_name, prefix=prefix, delimiter=delimiter)**
- **bucket.blob(source\_blob\_name)**

- `blob.download_to_filename(destination_file_name)`
- `blob.upload_from_filename(source_file_name)`
- `bucket.labels`
- `bucket.patch()` = update
- `blob.delete()`
- `bucket.delete()`
- `blob_to_be_created.compose(list of source blobs [blob1, blob2...])`

Doc: [Storage Client — google-cloud-storage documentation \(googleapis.dev\)](https://googleapis.dev/python/google-cloud-storage/latest/storage.html)