

# Big Data Coursework - Questions

## Data Processing and Machine Learning in the Cloud

This is the **INM432 Big Data coursework 2024**. This coursework contains extended elements of **theory** and **practice**, mainly around parallelisation of tasks with Spark and a bit about parallel training using TensorFlow.

### Code and Report

Your tasks parallelization of tasks in PySpark, extension, evaluation, and theoretical reflection. Please complete and submit the **coding tasks** in a copy of **this notebook**. Write your code in the **indicated cells** and **include the output** in the submitted notebook. Make sure that **your code contains comments** on its **structure** and explanations of its **purpose**.

Provide also a **report** with the **textual answers in a separate document**. Include **screenshots** from the Google Cloud web interface (don't use the SCREENSHOT function that Google provides, but take a picture of the graphs you see for the VMs) and result tables, as well as written text about the analysis.

### Submission

Download and submit **your version of this notebook** as an **.ipynb** file and also submit a **shareable link** to your notebook on Colab in your report (created with the Colab 'Share' function) (**and don't change the online version after submission**).

Further, provide your **report as a PDF document**. **State the number of words** in the document at the end. The report should **not have more than 2000 words**.

Please also submit a **PDF of your Jupyter notebook**.

### Introduction and Description

This coursework focuses on parallelisation and scalability in the cloud with Spark and TensorFlow/Keras. We start with code based on **lessons 3 and 4** of the [Fast and Lean Data Science](https://github.com/GoogleCloudPlatform/training-data-science/tree/master/courses/fast-and-lean-data-science) ([https://github.com/GoogleCloudPlatform/training-data-science](https://github.com/GoogleCloudPlatform/training-data-science/tree/master/courses/fast-and-lean-data-science)) course by Martin Gorner. The course is based on Tensorflow for data processing and MachineLearning. Tensorflow's data processing approach is somewhat similar to that of Spark, but you don't need to study Tensorflow, just make sure you understand the high-level structure.  
What we will do here is **parallelising pre-processing**, and **measuring** performance, and we will perform **evaluation** and **analysis** on the cloud performance, as well as **theoretical discussion**.

This coursework contains **3 sections**.

## Section 0

This section just contains some necessary code for setting up the environment. It has no tasks for you (but do read the code and comments).

## Section 1

Section 1 is about preprocessing a set of image files. We will work with a public dataset "Flowers" (3600 images, 5 classes). This is not a vast dataset, but it keeps the tasks more manageable for development and you can scale up later, if you like.

In '**Getting Started**' we will work through the data preprocessing code from *Fast and Lean Data Science* which uses TensorFlow's `tf.data` package. There is no task for you here, but you will need to re-use some of this code later.

In **Task 1** you will **parallelise the data preprocessing in Spark**, using Google Cloud (GC) Dataproc. This involves adapting the code from 'Getting Started' to use Spark and running it in the cloud.

## Section 2

In **Section 2** we are going to **measure the speed of reading data** in the cloud. In **Task 2** we will **parallelize the measuring** of different configurations **using Spark**.

## Section 3

This section is about the theoretical discussion, based on one paper, in **Task 3**. The answers should be given in the PDF report.

## General points

For **all coding tasks**, take the **time of the operations** and for the cloud operations, get **performance information from the web interfaces** for your reporting and analysis.

The **tasks** are **mostly independent** of each other. The later tasks can mostly be addressed

## Section 0: Set-up

As usual, you need to run the **imports and authentication every time you work with this notebook**. Use the **local Spark** installation for development before you send jobs to the cloud.

Read through this section once and **fill in the project ID the first time**, then you can just step straight through this at the beginning of each session - except for the two authentication cells.

## Imports

We import some **packages that will be needed throughout**. For the **code that runs in the cloud**, we will need **separate import sections** that will need to be partly different from the one below.

```
In [ ]: 1 import os, sys, math  
2 import numpy as np  
3 import scipy as sp  
4 import scipy.stats  
5 import time  
6 import datetime  
7 import string  
8 import random  
9 from matplotlib import pyplot as plt  
10 import tensorflow as tf  
11 print("Tensorflow version " + tf.__version__)  
12 import pickle
```

Tensorflow version 2.15.0

## Cloud and Drive authentication

This is for **authenticating with GCS Google Drive**, so that we can create and use our own buckets and access Dataproc and AI-Platform.

This section **starts with the two interactive authentications**.

First, we mount Google Drive for persistent local storage and create a directory DB-CW that you can use for this work. Then we'll set up the cloud environment, including a storage bucket.

```
In [ ]: 1 print('Mounting google drive...')  
2 from google.colab import drive  
3 drive.mount('/content/drive')  
4 %cd "/content/drive/MyDrive"  
5 !mkdir BD-CW  
6 %cd "/content/drive/MyDrive/BD-CW"
```

```
Mounting google drive...  
Mounted at /content/drive  
/content/drive/MyDrive  
mkdir: cannot create directory 'BD-CW': File exists  
/content/drive/MyDrive/BD-CW
```

Next, we authenticate with the GCS to enable access to Dataproc and AI-Platform.

```
In [ ]: 1 import sys  
2 if 'google.colab' in sys.modules:  
3     from google.colab import auth  
4     auth.authenticate_user()
```

It is useful to **create a new Google Cloud project** for this coursework. You can do this on the [GC Console page \(<https://console.cloud.google.com>\)](https://console.cloud.google.com) by clicking on the entry at the top, right of the *Google Cloud Platform* and choosing *New Project*. **Copy the generated project ID** to the next cell. Also **enable billing** and the **Compute, Storage and Dataproc APIs** like we did during the labs.

```
In [ ]: 1 PROJECT = 'big-data-cw-420116' ### USE YOUR GOOGLE CLOUD PROJECT ID HERE
2 !gcloud config set project $PROJECT
3 REGION = 'us-central1'
4 CLUSTER = '{}-cluster'.format(PROJECT)
5 !gcloud config set compute/region $REGION
6 !gcloud config set dataproc/region $REGION
7
8 !gcloud config list # show some information
```

```
Updated property [core/project].
WARNING: Property validation for compute/region was skipped.
Updated property [compute/region].
Updated property [dataproc/region].
[component_manager]
disable_update_check = True
[compute]
region = us-central1
[core]
account = antonio-jose.lopez-roldan@city.ac.uk
project = big-data-cw-420116
[dataproc]
region = us-central1

Your active configuration is: [default]
```

With the cell below, we **create a storage bucket** that we will use later for **global storage**. If the bucket exists you will see a "ServiceException: 409 ...", which does not cause any problems. **You must create your own bucket to have write access.**

```
In [ ]: 1 BUCKET = 'gs://{}-storage'.format(PROJECT)
2 !gsutil mb $BUCKET
```

```
Creating gs://big-data-cw-420116-storage/...
ServiceException: 409 A Cloud Storage bucket named 'big-data-cw-420116-sto
rage' already exists. Try another name. Bucket names must be globally uniq
ue across all Google Cloud projects, including those outside of your organ
ization.
```

The cell below just **defines some routines for displaying images** that will be **used later**. You can see the code by double-clicking, but you don't need to study this.



In [ ]:

```
1 #@title Utility functions for image display **[RUN THIS TO ACTIVATE]**
2 def display_9_images_from_dataset(dataset):
3     plt.figure(figsize=(13,13))
4     subplot=331
5     for i, (image, label) in enumerate(dataset):
6         plt.subplot(subplot)
7         plt.axis('off')
8         plt.imshow(image.numpy().astype(np.uint8))
9         plt.title(str(label.numpy()), fontsize=16)
10        # plt.title(Label.numpy().decode(), fontsize=16)
11        subplot += 1
12        if i==8:
13            break
14    plt.tight_layout()
15    plt.subplots_adjust(wspace=0.1, hspace=0.1)
16    plt.show()
17
18 def display_training_curves(training, validation, title, subplot):
19     if subplot%10==1: # set up the subplots on the first call
20         plt.subplots(figsize=(10,10), facecolor='#F0F0F0')
21         plt.tight_layout()
22     ax = plt.subplot(subplot)
23     ax.set_facecolor('#F8F8F8')
24     ax.plot(training)
25     ax.plot(validation)
26     ax.set_title('model ' + title)
27     ax.set_ylabel(title)
28     ax.set_xlabel('epoch')
29     ax.legend(['train', 'valid.'])
30
31 def dataset_to_numpy_util(dataset, N):
32     dataset = dataset.batch(N)
33     for images, labels in dataset:
34         numpy_images = images.numpy()
35         numpy_labels = labels.numpy()
36         break;
37     return numpy_images, numpy_labels
38
39 def title_from_label_and_target(label, correct_label):
40     correct = (label == correct_label)
41     return "{} [{}{}{}]".format(CLASSES[label], str(correct), ', shoud be',
42                                CLASSES[correct_label] if not correct else '')
43
44 def display_one_flower(image, title, subplot, red=False):
45     plt.subplot(subplot)
46     plt.axis('off')
47     plt.imshow(image)
48     plt.title(title, fontsize=16, color='red' if red else 'black')
49     return subplot+1
50
51 def display_9_images_with_predictions(images, predictions, labels):
52     subplot=331
53     plt.figure(figsize=(13,13))
54     classes = np.argmax(predictions, axis=-1)
55     for i, image in enumerate(images):
56         title, correct = title_from_label_and_target(classes[i], labels[i])
57         subplot = display_one_flower(image, title, subplot, not correct)
58         if i >= 8:
59             break;
60
61     plt.tight_layout()
```

```
62 plt.subplots_adjust(wspace=0.1, hspace=0.1)
63 plt.show()
64
```

```
In [ ]: 1 # Check List of services working and enabled
2 !gcloud services list --enabled
```

NAME	TITLE
analyticshub.googleapis.com	Analytics Hub API
artifactregistry.googleapis.com	Artifact Registry API
autoscaling.googleapis.com	Cloud Autoscaling API
bigrquery.googleapis.com	BigQuery API
bigrqueryconnection.googleapis.com	BigQuery Connection API
bigrquerydatapolicy.googleapis.com	BigQuery Data Policy API
bigrquerymigration.googleapis.com	BigQuery Migration API
bigrqueryreservation.googleapis.com	BigQuery Reservation API
bigrquerystorage.googleapis.com	BigQuery Storage API
cloudapis.googleapis.com	Google Cloud APIs
cloudresourcemanager.googleapis.com	Cloud Resource Manager API
cloudtrace.googleapis.com	Cloud Trace API
compute.googleapis.com	Compute Engine API
container.googleapis.com	Kubernetes Engine API
containerfilesystem.googleapis.com	Container File System API
containerregistry.googleapis.com	Container Registry API
dataform.googleapis.com	Dataform API
dataplex.googleapis.com	Cloud Dataplex API
dataproc-control.googleapis.com	Cloud Dataproc Control API
dataproc.googleapis.com	Cloud Dataproc API
datastore.googleapis.com	Cloud Datastore API
dns.googleapis.com	Cloud DNS API
gkebackup.googleapis.com	Backup for GKE API
iam.googleapis.com	Identity and Access Management (IAM)
API	
iamcredentials.googleapis.com	IAM Service Account Credentials API
logging.googleapis.com	Cloud Logging API
monitoring.googleapis.com	Cloud Monitoring API
networkconnectivity.googleapis.com	Network Connectivity API
oslogin.googleapis.com	Cloud OS Login API
pubsub.googleapis.com	Cloud Pub/Sub API
servicemanagement.googleapis.com	Service Management API
serviceusage.googleapis.com	Service Usage API
sql-component.googleapis.com	Cloud SQL
storage-api.googleapis.com	Google Cloud Storage JSON API
storage-component.googleapis.com	Cloud Storage
storage.googleapis.com	Cloud Storage API

## Install Spark locally for quick testing

You can use the cell below to **install Spark locally on this Colab VM** (like in the labs), to do quicker small-scale interactive testing. Using Spark in the cloud with **Dataproc is still required for the final version**.

```
In [ ]: 1 %cd
2 !apt-get update -qq
3 !apt-get install openjdk-8-jdk-headless -qq >> /dev/null # send any output to null
4 !tar -xzf "/content/drive/My Drive/Big_Data/data/spark/spark-3.5.0-bin-hadoop3.tgz"
5
6 !pip install -q findspark
7 import os
8 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
9 os.environ["SPARK_HOME"] = "/root/spark-3.5.0-bin-hadoop3" # Was getting error
10 import findspark
11 findspark.init()
12 import pyspark
13 print(pyspark.__version__)
14 sc = pyspark.SparkContext.getOrCreate()
15 print(sc)

/root
3.5.0

/usr/lib/python3.10/subprocess.py:1796: RuntimeWarning: os.fork() was called. os.fork() is incompatible with multithreaded code, and JAX is multithreaded, so this will likely lead to a deadlock.
  self.pid = _posixsubprocess.fork_exec()

<SparkContext master=local[*] appName=pyspark-shell>
```

## Section 1: Data pre-processing

This section is about the **pre-processing of a dataset** for deep learning. We first look at a ready-made solution using Tensorflow and then we build a implement the same process with Spark. The tasks are about **parallelisation** and **analysis** the performance of the cloud implementations.

### 1.1 Getting started

In this section, we get started with the data pre-processing. The code is based on lecture 3 of the 'Fast and Lean Data Science' course.

**This code is using the TensorFlow tf.data package**, which supports map functions, similar to Spark. Your **task** will be to **re-implement the same approach in Spark**.

We start by **setting some variables for the Flowers dataset**.

```
In [ ]: 1 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input files
2 PARTITIONS = 16 # no of partitions we will use later
3 TARGET_SIZE = [192, 192] # target resolution for the images
4 CLASSES = [b'daisy', b'dandelion', b'roses', b'sunflowers', b'tulips']
5 # Labels for the data
```

We **read the image files** from the public GCS bucket that contains the *Flowers* dataset. **TensorFlow** has **functions** to execute glob patterns that we use to calculate the the number of images in total and per partition (rounded up as we cannont deal with parts of images).

```
In [ ]: 1 nb_images = len(tf.io.gfile.glob(GCS_PATTERN)) # number of images
2 partition_size = math.ceil(1.0 * nb_images / PARTITIONS) # images per p
3 print("GCS_PATTERN matches {} images, to be divided into {} partitions
4
```

GCS\_PATTERN matches 3670 images, to be divided into 16 partitions with up to 230 images each.

## Map functions

In order to read use the images for learning, they need to be **preprocessed** (decoded, resized, cropped, and potentially recompressed). Below are **map functions** for these steps. You **don't need to study the internals of these functions** in detail.

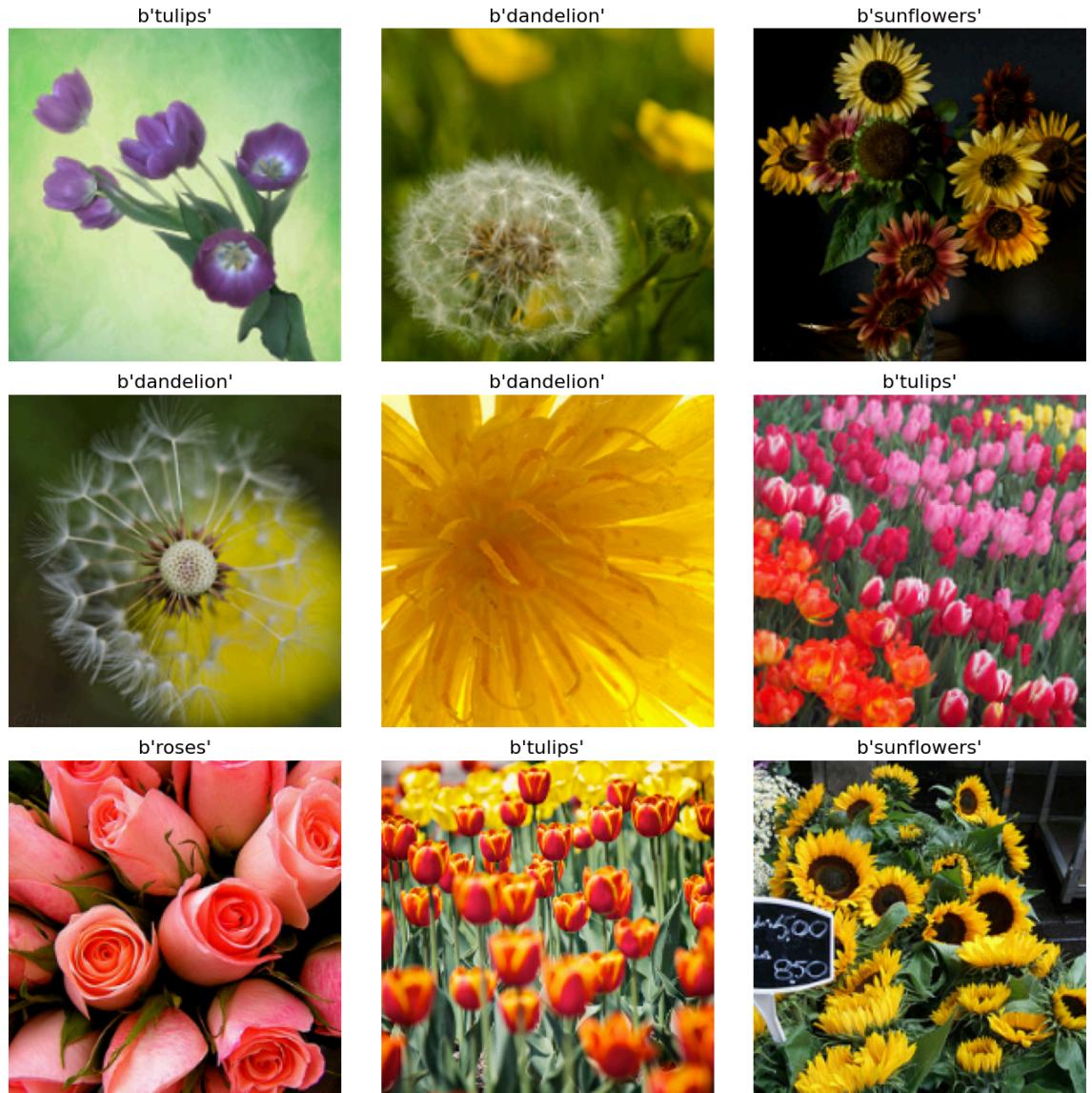
```
In [ ]: 1 def decode_jpeg_and_label(filepath):
2     # extracts the image data and creates a class label, based on the f
3     bits = tf.io.read_file(filepath)
4     image = tf.image.decode_jpeg(bits)
5     # parse flower name from containing directory
6     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
7     label2 = label.values[-2]
8     return image, label2
9
10 def resize_and_crop_image(image, label):
11     # Resizes and cropd using "fill" algorithm:
12     # always make sure the resulting image is cut out from the source i
13     # so that it fills the TARGET_SIZE entirely with no black bars
14     # and a preserved aspect ratio.
15     w = tf.shape(image)[0]
16     h = tf.shape(image)[1]
17     tw = TARGET_SIZE[1]
18     th = TARGET_SIZE[0]
19     resize_crit = (w * th) / (h * tw)
20     image = tf.cond(resize_crit < 1,
21                     lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), #
22                     lambda: tf.image.resize(image, [w*th/h, h*th/h]) #
23                     )
24     nw = tf.shape(image)[0]
25     nh = tf.shape(image)[1]
26     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
27     return image, label
28
29 def recompress_image(image, label):
30     # this reduces the amount of data, but takes some time
31     image = tf.cast(image, tf.uint8)
32     image = tf.image.encode_jpeg(image, optimize_size=True, chroma_down
33     return image, label
```

With `tf.data`, we can apply decoding and resizing as map functions.

```
In [ ]: 1 dsetFiles = tf.data.Dataset.list_files(GCS_PATTERN) # This also shuffle  
2 dsetDecoded = dsetFiles.map(decode_jpeg_and_label)  
3 dsetResized = dsetDecoded.map(resize_and_crop_image)
```

We can also look at some images using the image display function defined above (the one with the hidden code).

```
In [ ]: 1 display_9_images_from_dataset(dsetResized)
```



Now, let's test continuous reading from the dataset. We can see that reading the first 100 files already takes some time.

```
In [ ]: 1 sample_set = dsetResized.batch(10).take(10) # take 10 batches of 10 ima
2 for image, label in sample_set:
3     print("Image batch shape {}, {}".format(image.numpy().shape,
4                                             [lbl.decode('utf8') for lbl in label.numpy()]))
```

Image batch shape (10, 192, 192, 3), ['sunflowers', 'dandelion', 'daisy', 'tulips', 'tulips', 'tulips', 'tulips', 'tulips', 'tulips', 'daisy'])  
Image batch shape (10, 192, 192, 3), ['daisy', 'roses', 'tulips', 'tulips', 'roses', 'daisy', 'roses', 'dandelion', 'roses', 'roses'])  
Image batch shape (10, 192, 192, 3), ['daisy', 'sunflowers', 'dandelion', 'sunflowers', 'daisy', 'sunflowers', 'tulips', 'tulips', 'dandelion', 'dandelion'])  
Image batch shape (10, 192, 192, 3), ['sunflowers', 'tulips', 'tulips', 'tulips', 'tulips', 'dandelion', 'roses', 'dandelion', 'sunflowers', 'dandelion'])  
Image batch shape (10, 192, 192, 3), ['sunflowers', 'sunflowers', 'daisy', 'tulips', 'daisy', 'tulips', 'dandelion', 'roses', 'tulips', 'sunflowers'])  
Image batch shape (10, 192, 192, 3), ['dandelion', 'sunflowers', 'roses', 'roses', 'sunflowers', 'tulips', 'sunflowers', 'dandelion', 'daisy', 'tulips'])  
Image batch shape (10, 192, 192, 3), ['dandelion', 'dandelion', 'sunflowers', 'dandelion', 'daisy', 'sunflowers', 'dandelion', 'sunflowers', 'roses', 'dandelion'])  
Image batch shape (10, 192, 192, 3), ['roses', 'tulips', 'tulips', 'daisies', 'roses', 'dandelion', 'sunflowers', 'dandelion', 'sunflowers', 'tulips'])  
Image batch shape (10, 192, 192, 3), ['daisy', 'tulips', 'sunflowers', 'tulips', 'daisy', 'sunflowers', 'dandelion', 'sunflowers', 'dandelion', 'dandelion'])  
Image batch shape (10, 192, 192, 3), ['sunflowers', 'tulips', 'sunflowers', 'dandelion', 'daisy', 'tulips', 'dandelion', 'dandelion', 'daisy'])

## 1.2 Improving Speed

Using individual image files didn't look very fast. The 'Lean and Fast Data Science' course introduced **two techniques to improve the speed**.

### Recompress the images

By **compressing** the images in the **reduced resolution** we save on the size. This **costs some CPU time** upfront, but **saves network and disk bandwidth**, especially when the data are **read multiple times**.

```
In [ ]: 1 # This is a quick test to get an idea how long recompressions takes.  
2 dataset4 = dsetResized.map(recompress_image)  
3 test_set = dataset4.batch(10).take(10)  
4 for image, label in test_set:  
5     print("Image batch shape {}, {}".format(image.numpy().shape, [lbl.
```

```
Image batch shape (10,), ['sunflowers', 'tulips', 'tulips', 'tulips', 'dai  
sy', 'sunflowers', 'dandelion', 'roses', 'daisy', 'tulips'])  
Image batch shape (10,), ['daisy', 'roses', 'roses', 'tulips', 'dandelio  
n', 'tulips', 'daisy', 'roses', 'tulips', 'dandelion'])  
Image batch shape (10,), ['dandelion', 'sunflowers', 'daisy', 'roses', 'su  
nflowers', 'dandelion', 'tulips', 'tulips', 'daisy', 'roses'])  
Image batch shape (10,), ['daisy', 'tulips', 'dandelion', 'dandelion', 'ro  
ses', 'roses', 'daisy', 'roses', 'daisy', 'sunflowers'])  
Image batch shape (10,), ['tulips', 'roses', 'daisy', 'daisy', 'dande  
lion', 'dandelion', 'dandelion', 'dandelion', 'dandelion'])  
Image batch shape (10,), ['tulips', 'dandelion', 'roses', 'daisy', 'tu  
lip', 'sunflowers', 'dandelion', 'dandelion', 'roses', 'tulips'])  
Image batch shape (10,), ['daisy', 'tulips', 'daisy', 'dandelion', 'dai  
sy', 'daisy', 'tulips', 'tulips', 'sunflowers', 'roses'])  
Image batch shape (10,), ['daisy', 'roses', 'sunflowers', 'tulips', 'dande  
lion', 'dandelion', 'dandelion', 'dandelion', 'dandelion'])  
Image batch shape (10,), ['dandelion', 'sunflowers', 'dandelion', 'roses',  
'roses', 'roses', 'daisy', 'dandelion', 'dandelion', 'sunflowers'])  
Image batch shape (10,), ['tulips', 'daisy', 'dandelion', 'dandelion', 'da  
ndelion', 'daisy', 'tulips', 'sunflowers', 'sunflowers', 'roses'])
```

## Write the dataset to TFRecord files

By writing **multiple preprocessed samples into a single file**, we can make further speed gains. We distribute the data over **partitions** to facilitate **parallelisation** when the data are used. First we need to **define a location** where we want to put the file.

```
In [ ]: 1 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for
```

Now we can **write the TFRecord files** to the bucket.

Running the cell takes some time and **only needs to be done once** or not at all, as you can use the publicly available data for the next few cells. For convenience I have commented out the call to `write_tfrecords` at the end of the next cell. You don't need to run it (it takes some time), but you'll need to use the code below later (but there is no need to study it in detail).

There is a **ready-made pre-processed data** versions available here: `gs://flowers-public/tfrecords-jpeg-192x192-2/`, that we can use for testing.

In [ ]:

```
1 # functions for writing TFRecord entries
2 # Feature values are always stored as lists, a single data element will
3 def _bytestring_feature(list_of_bytestrings):
4     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_
5
6 def _int_feature(list_of_ints): # int64
7     return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_
8
9 def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data r
10    class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (ord
11    one_hot_class = np.eye(len(CLASSES))[class_num]      # [0, 0, 1, 0,
12    feature = {
13        "image": _bytestring_feature([img_bytes]), # one image in the L
14        "class": _int_feature([class_num]) #,       # one class in the L
15    }
16    return tf.train.Example(features=tf.train.Features(feature=feature))
17
18 def write_tfrecords(GCS_PATTERN,GCS_OUTPUT,partition_size): # write the
19    print("Writing TFRecords")
20    tt0 = time.time()
21    filenames = tf.data.Dataset.list_files(GCS_PATTERN)
22    dataset1 = filenames.map(decode_jpeg_and_label)
23    dataset2 = dataset1.map(resize_and_crop_image)
24    dataset3 = dataset2.map(recompress_image)
25    dataset4 = dataset3.batch(partition_size) # partitioning: there wil
26    for partition, (image, label) in enumerate(dataset4):
27        # batch size used as partition size here
28        partition_size = image.numpy().shape[0]
29        # good practice to have the number of records in the filename
30        filename = GCS_OUTPUT + "{:02d}-{}.tfrec".format(partition, par
31        # You need to change GCS_OUTPUT to your own bucket to actually
32        with tf.io.TFRecordWriter(filename) as out_file:
33            for i in range(partition_size):
34                example = to_tfrecord(out_file,
35                                      image.numpy()[i], # re-compressed i
36                                      label.numpy()[i] #
37                                      )
38                out_file.write(example.SerializeToString())
39    print("Wrote file {} containing {} records".format(filename, pa
40    print("Total time: "+str(time.time()-tt0))
41
42 write_tfrecords(GCS_PATTERN,GCS_OUTPUT,partition_size) # uncomment to r
```

### Writing TFRecords

```
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s00-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s01-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s02-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s03-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s04-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s05-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s06-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s07-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s08-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s09-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s10-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s11-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s12-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s13-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s14-230.tfrec containing 230 records  
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flower  
s15-220.tfrec containing 220 records  
Total time: 223.18424677848816
```

### Test the TFRecord files

We can now **read from the TFRecord files**. By default, we use the files in the public bucket. Comment out the 1st line of the cell below to use the files written in the cell above.

In [ ]:

```
1 def read_tfrecord(example):
2     features = {
3         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = b
4         "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] me
5     }
6     # decode the TFRecord
7     example = tf.io.parse_single_example(example, features)
8     image = tf.image.decode_jpeg(example['image'], channels=3)
9     image = tf.reshape(image, [*TARGET_SIZE, 3])
10    class_num = example['class']
11    return image, class_num
12
13
14 def load_dataset(filenames):
15     # read from TFRecords. For optimal performance, read from multiple
16     # TFRecord files at once and set the option experimental_deterministic
17     # to allow order-altering optimizations.
18     option_no_order = tf.data.Options()
19     option_no_order.experimental_deterministic = False
20
21     dataset = tf.data.TFRecordDataset(filenames)
22     dataset = dataset.with_options(option_no_order)
23     dataset = dataset.map(read_tfrecord)
24     return dataset
25
26 filenames = tf.io.gfile.glob(GCS_OUTPUT + "*tfrec")
27 datasetTfrec = load_dataset(filenames)
```

Let's have a look if reading from the TFRecord files is quicker.

In [ ]:

```
1 batched_dataset = datasetTfrec.batch(10)
2 sample_set = batched_dataset.take(10)
3 for image, label in sample_set:
4     print("Image batch shape {}, {}".format(image.numpy().shape,
5                                              [str(lbl) for lbl in label.numpy()]))
```

Image batch shape (10, 192, 192, 3), ['1', '0', '4', '4', '1', '3', '2',
'1', '4', '2']
Image batch shape (10, 192, 192, 3), ['1', '0', '4', '1', '2', '1', '1',
'4', '0', '1']
Image batch shape (10, 192, 192, 3), ['1', '2', '3', '1', '0', '1', '2',
'4', '4', '3']
Image batch shape (10, 192, 192, 3), ['1', '3', '1', '4', '2', '3', '0',
'3', '0', '3']
Image batch shape (10, 192, 192, 3), ['3', '3', '3', '4', '3', '0', '1',
'0', '4', '2']
Image batch shape (10, 192, 192, 3), ['2', '4', '1', '3', '0', '3', '4',
'3', '3', '1']
Image batch shape (10, 192, 192, 3), ['0', '4', '3', '4', '0', '4', '3',
'3', '0', '3']
Image batch shape (10, 192, 192, 3), ['0', '0', '1', '0', '2', '4', '2',
'2', '3', '0']
Image batch shape (10, 192, 192, 3), ['2', '1', '1', '4', '1', '2', '3',
'3', '0', '3']
Image batch shape (10, 192, 192, 3), ['0', '4', '0', '4', '4', '4', '1',
'1', '1', '1']

Wow, we have a **massive speed-up!** The repackaging is worthwhile :-)

# Task 1: Write TFRecord files to the cloud with Spark (40%)

Since recompressing and repackaging is very effective, we would like to be able to do it in parallel for large datasets. This is a relatively straightforward case of **parallelisation**. We will **use Spark to implement** the same process as above, but in parallel.

## 1a) Create the script (14%)

**Re-implement** the pre-processing in Spark, using Spark mechanisms for **distributing** the workload **over multiple machines**.

You need to:

- i) **Copy** over the **mapping functions** (see section 1.1) and **adapt** the resizing and recompression functions **to Spark** (only one argument). (3%)
- ii) **Replace** the TensorFlow **Dataset objects with RDDs**, starting with an RDD that contains the list of image filenames. (3%)
- iii) **Sample** the the RDD to a smaller number at an appropriate position in the code. Specify a sampling factor of 0.02 for short tests. (1%)
- iv) Then **use the functions from above** to write the TFRecord files. (3%)
- v) The code for **writing to the TFRecord files** needs to be put into a function, that can be applied to every partition with the '[RDD.mapPartitionsWithIndex](#)' (<https://spark.apache.org/docs/2.4.8/api/python/pyspark.html#pyspark.RDD.mapPartitionsWithIndex>) function. The return value of that function is not used here, but you should return the filename, so that you have a list of the created TFRecord files. (4%)





In [ ]:

```
1                                     ### CODING TASK ### ("Done & Working")
2 # Setting the working environment
3 import argparse
4 import datetime
5 import math
6 import numpy as np
7 import os
8 import pyspark
9 import pickle
10 import pandas as pd
11 import random
12 import string
13 import sys
14 #import scipy as sp
15 #import scipy.stats
16 import time
17 import tensorflow as tf
18
19 #from matplotlib import pyplot as plt
20 from pyspark.sql import SQLContext
21 from pyspark.sql import Row
22
23 # define variables
24 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input j
25 PARTITIONS = 16 # no of partitions we will use later
26 PROJECT = 'big-data-cw-420116' # my project id
27 BUCKET = 'gs://{}-storage'.format(PROJECT) # my own bucket storage
28 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for
29
30 #i) Copy over the mapping functions (see section 1.1) and adapt the res
31 def decode_jpeg_and_label(filepath):
32     # extracts the image data and creates a class label, based on the j
33     bits = tf.io.read_file(filepath)
34     image = tf.image.decode_jpeg(bits)
35     # parse flower name from containing directory
36     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
37     label2 = label.values[-2]
38     return image, label2
39
40 def resize_and_crop_image(image, label):
41     # Resizes and cropd using "fill" algorithm:
42     # always make sure the resulting image is cut out from the source i
43     # so that it fills the TARGET_SIZE entirely with no black bars
44     # and a preserved aspect ratio.
45     w = tf.shape(image)[0]
46     h = tf.shape(image)[1]
47     tw = TARGET_SIZE[1]
48     th = TARGET_SIZE[0]
49     resize_crit = (w * th) / (h * tw)
50     image = tf.cond(resize_crit < 1,
51                     lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), +
52                     lambda: tf.image.resize(image, [w*th/h, h*th/h]) +
53                     )
54     nw = tf.shape(image)[0]
55     nh = tf.shape(image)[1]
56     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
57     return image, label
58
59 def recompress_image(image, label):
60     # this reduces the amount of data, but takes some time
61     image = tf.cast(image, tf.uint8)
```

```

62     image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampled=False)
63     return image, label
64
65 # ii) Replace the TensorFlow Dataset objects with RDDs, starting with our GCS filenames
66
67 filenames = tf.io.gfile.glob(GCS_PATTERN)
68
69 # Get spark context for RDDs
70 sc = pyspark.SparkContext.getOrCreate()
71
72 # Create RDD for files
73 filenames_rdd = sc.parallelize(filenames)
74
75 # iii) Sample the the RDD to a smaller number at an appropriate position
76
77 # Sampling the rdd
78 sample_rdd = filenames_rdd.sample(False, 0.02)
79
80 #Adapting tensorflow to spark, rdd for decode jpeg and label, rdd for re-compressing
81 decode_jpeg_and_label_rdd = filenames_rdd.map(decode_jpeg_and_label)
82 resize_and_crop_image_rdd = decode_jpeg_and_label_rdd.map(lambda x: resize_and_crop(x))
83 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress(x))
84
85 # iv) Then use the functions from above to write the TFRecord files, using the RDDs
86
87 # functions for writing TFRecord entries
88 # Feature values are always stored as lists, a single data element will be a list with one item
89 def _bytestring_feature(list_of_bytess):
90     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytess))
91
92 def _int_feature(list_of_ints): # int64
93     return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
94
95 def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data record
96     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order matters)
97     one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0]
98     feature = {
99         "image": _bytestring_feature([img_bytes]), # one image in the list
100        "class": _int_feature([class_num]) #, # one class in the list
101    }
102    return tf.train.Example(features=tf.train.Features(feature=feature))
103
104 def write_tfrecords(partition_index,partition): # write the images to disk
105     # tt0 = time.time()
106     # good practice to have the number of records in the filename
107     filename = GCS_OUTPUT + "{}.tfrec".format(partition_index)
108     # (CHANGED) You need to change GCS_OUTPUT to your own bucket to actually write to it
109     with tf.io.TFRecordWriter(filename) as out_file:
110         for element in partition:
111             image=element[0]
112             label=element[1]
113             example = to_tfrecord(out_file,
114                                   image.numpy(), # re-compressed image: already compressed by tensorflow
115                                   label.numpy() #
116                                   )
117             out_file.write(example.SerializeToString())
118             #print("Wrote file {} containing {} records".format(filename, len(partition)))
119             #print("Total time: "+str(time.time()-tt0))
120     return (filename)
121
122 # v) The code for writing to the TFRecord files needs to be put into a

```

```
123 #     The return value of that function is not used here, but you should
124
125 partitions_rdd = recompress_image_rdd.repartition(PARTITIONS)
126 filenames_rdd = recompress_image_rdd.mapPartitionsWithIndex(write_tfrecords)
127 print("Writing TFRecords")
128
129 # Extra resources checked:
130 # https://sparkbyexamples.com/pyspark-rdd/
131 # https://spark.apache.org/docs/latest/rdd-programming-guide.html
132 # https://stackoverflow.com/questions/69205589/how-preprocess-image-using-spark
133 # https://www.tensorflow.org/guide/distributed_training#sparkstrategy
134 # https://www.w3schools.com/python/python_lambda.asp
135 # https://pages.databricks.com/rs/094-YMS-629/images/LearningSpark2.0.pdf
```

Writing TFRecords

```
In [ ]: 1 # First elements of the RDD  
2 decode_jpeg_and_label_rdd.take(1)
```

```
Out[21]: [(<tf.Tensor: shape=(263, 320, 3), dtype=uint8, numpy=  
array([[[133, 135, 132],  
       [136, 138, 135],  
       [140, 142, 139],  
       ...,  
       [152, 152, 150],  
       [155, 155, 153],  
       [148, 148, 146]],  
  
      [[133, 135, 132],  
       [136, 138, 135],  
       [140, 142, 139],  
       ...,  
       [153, 153, 151],  
       [155, 155, 153],  
       [147, 147, 145]],  
  
      [[132, 134, 129],  
       [135, 137, 134],  
       [139, 141, 138],  
       ...,  
       [152, 152, 150],  
       [154, 154, 152],  
       [146, 146, 144]],  
  
      ...,  
  
      [[ 44,  48,  25],  
       [ 44,  48,  25],  
       [ 44,  48,  25],  
       ...,  
       [127, 126, 122],  
       [127, 126, 122],  
       [127, 126, 122]],  
  
      [[ 44,  48,  25],  
       [ 44,  48,  25],  
       [ 44,  48,  25],  
       ...,  
       [128, 127, 123],  
       [128, 127, 123],  
       [128, 127, 123]],  
  
      [[ 43,  47,  24],  
       [ 43,  47,  24],  
       [ 43,  47,  24],  
       ...,  
       [129, 128, 124],  
       [129, 128, 124],  
       [130, 129, 125]]], dtype=uint8)>,  
<tf.Tensor: shape=(), dtype=string, numpy=b'daisy'>)]
```

```
In [ ]: 1 # first elements of the RDD  
2 resize_and_crop_image_rdd.take(1)
```

```
Out[22]: [(<tf.Tensor: shape=(192, 192, 3), dtype=float32, numpy=  
array([[ [154.31648 , 158.31648 , 161.31648 ],  
        [154.03465 , 158.03465 , 161.03465 ],  
        [152.40732 , 156.40732 , 159.40732 ],  
        ...,  
        [166.26291 , 167.93884 , 169.71353 ],  
        [164.40128 , 168.40128 , 169.40128 ],  
        [164.        , 168.        , 169.        ]],  
  
        [[168.3533 , 172.16167 , 175.54494 ],  
        [166.39734 , 170.39734 , 173.39734 ],  
        [163.65054 , 167.65054 , 170.65054 ],  
        ...,  
        [165.4297 , 167.10564 , 168.88033 ],  
        [164.40128 , 168.40128 , 169.40128 ],  
        [164.        , 168.        , 169.        ]],  
  
        [[164.95058 , 168.        , 172.90114 ],  
        [163.81895 , 167.81895 , 170.81895 ],  
        [162.34184 , 166.34184 , 169.34184 ],  
        ...,  
        [166.90747 , 167.95851 , 169.9415 ],  
        [166.25023 , 167.47679 , 169.40128 ],  
        [165.84895 , 167.07552 , 169.        ]],  
  
        ...,  
  
        [[120.384514, 118.384514, 119.384514],  
        [123.10339 , 121.159195, 122.131294],  
        [124.0755 , 124.0755 , 124.878075],  
        ...,  
        [127.89167 , 127.34229 , 124.66636 ],  
        [126.44648 , 127.44648 , 122.44648 ],  
        [126.97421 , 127.97421 , 122.97421 ]],  
  
        [[121.46287 , 119.46287 , 120.46287 ],  
        [124.460785, 122.51658 , 123.48868 ],  
        [124.8213 , 124.8213 , 125.62388 ],  
        ...,  
        [129.9137 , 129.46465 , 126.78871 ],  
        [126.821304, 128.        , 123.        ],  
        [127.        , 128.        , 123.        ]],  
  
        [[122.18799 , 120.18799 , 121.18799 ],  
        [124.977264, 123.03305 , 124.00516 ],  
        [123.783615, 123.783615, 124.5862 ],  
        ...,  
        [131.46323 , 131.13916 , 128.46323 ],  
        [126.925804, 128.8151 , 123.815094],  
        [127.        , 128.02274 , 123.022736]]], dtype=float32)>,  
<tf.Tensor: shape=(), dtype=string, numpy=b'daisy'>)]
```

```
In [ ]: 1 # First elements of the RDD  
2 recompress_image_rdd.take(1)
```

```
In [ ]: 1 # checking name  
2 filenames_rdd.take(1)
```

Out[24]: ['g']

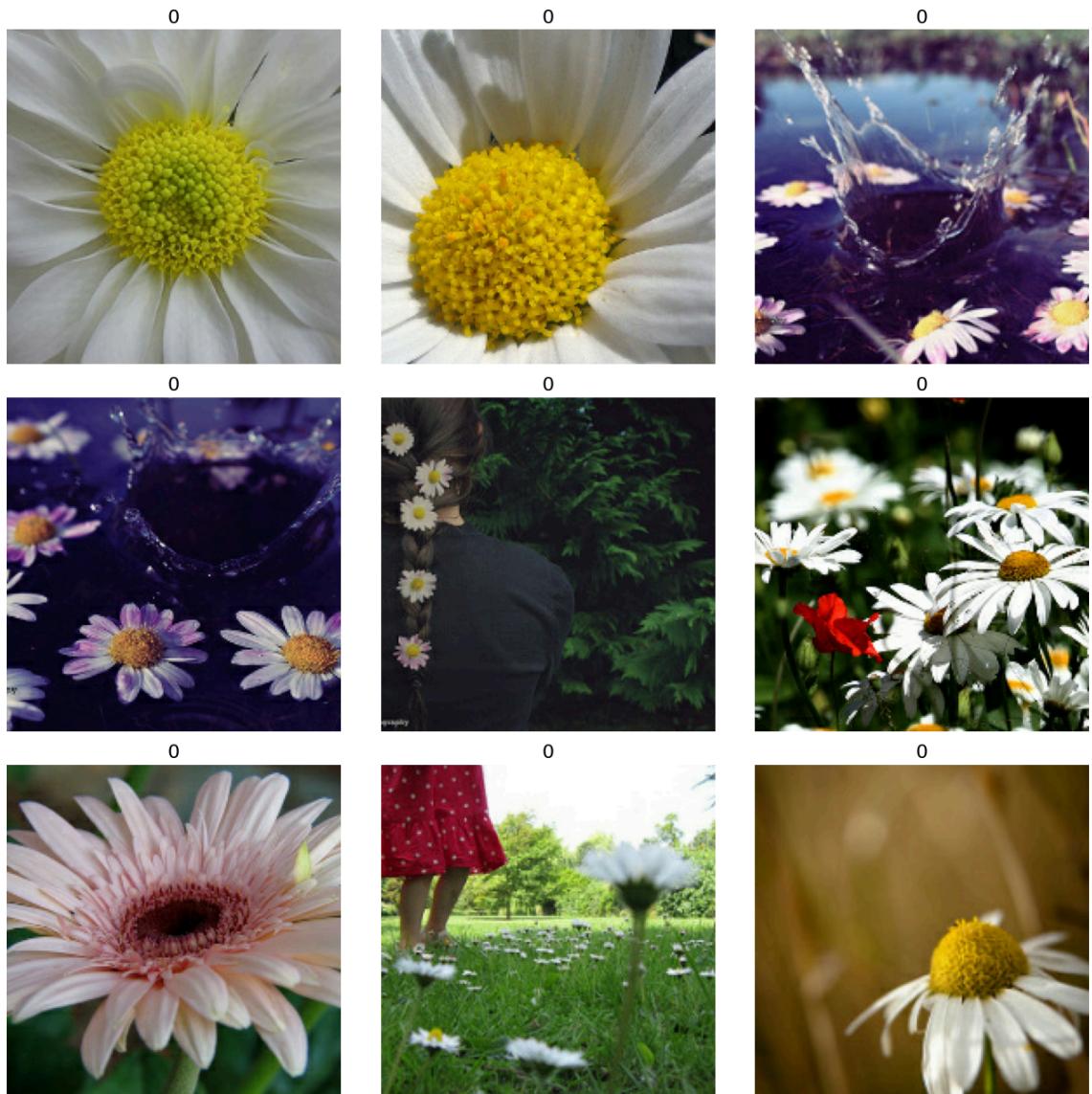
**Name shown:** ['g'], should be flower name, I know the error is within the function definition/names assigned to the variables, didn't have enough time to find the Root Cause of this issue and fix it.

### **1b) Testing (3%)**

- i) Read from the TFRecord Dataset, using `load_dataset` and `display_9_images_from_dataset` to test.

In [ ]:

```
1 ### CODING TASK ### ("Done & Working")
2 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers'
3
4 # Pre-defined function. Function to parse the TFRecord
5 def read_tfrecord(example):
6     features = {
7         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = b
8         "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] me
9     }
10    # decode the TFRecord
11    example = tf.io.parse_single_example(example, features)
12    image = tf.image.decode_jpeg(example['image'], channels=3)
13    image = tf.reshape(image, [*TARGET_SIZE, 3])
14    class_num = example['class']
15    return image, class_num
16
17 # Pre-defined function. Function to Load the dataset
18 def load_dataset(filenames):
19     # read from TFRecords. For optimal performance, read from multiple
20     # TFRecord files at once and set the option experimental_determinis
21     # to allow order-altering optimizations.
22     option_no_order = tf.data.Options()
23     option_no_order.experimental_deterministic = False
24
25     dataset = tf.data.TFRecordDataset(filenames)
26     dataset = dataset.with_options(option_no_order)
27     dataset = dataset.map(read_tfrecord)
28     return dataset
29
30 filenames = tf.io.gfile.glob(GCS_OUTPUT + "*tfrec")
31 datasetTfrec = load_dataset(filenames)
32 display_9_images_from_dataset(datasetTfrec)
```



**Note:** All images show "0", couldn't find why.

- ii) Write your code above into a file using the *cell magic* `%%writefile spark_write_tfrec.py` at the beginning of the file. Then, run the file locally in Spark.



In [ ]:

```
1     ### CODING TASK ### ("Done & Working")
2 %writefile spark_write_tfrec.py
3 # Setting the working environment
4 import argparse
5 import datetime
6 import math
7 import numpy as np
8 import os
9 import pyspark
10 import pickle
11 import pandas as pd
12 import random
13 import string
14 import sys
15 #import scipy as sp # "no module named spicy error" tryied to install it
16 #import scipy.stats
17 import time
18 import tensorflow as tf
19
20 #from matplotlib import pyplot as plt # "not module named matplotlib"
21 from pyspark.sql import SQLContext
22 from pyspark.sql import Row
23
24 # define variables
25 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input j
26 PARTITIONS = 16 # no of partitions we will use later
27 PROJECT = 'big-data-cw-420116' # my project id
28 BUCKET = 'gs://{}-storage'.format(PROJECT) # my own bucket storage
29 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for
30
31 #i) Copy over the mapping functions (see section 1.1) and adapt the res
32 def decode_jpeg_and_label(filepath):
33     # extracts the image data and creates a class label, based on the fi
34     bits = tf.io.read_file(filepath)
35     image = tf.image.decode_jpeg(bits)
36     # parse flower name from containing directory
37     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
38     label2 = label.values[-2]
39     return image, label2
40
41 def resize_and_crop_image(image, label):
42     # Resizes and cropd using "fill" algorithm:
43     # always make sure the resulting image is cut out from the source i
44     # so that it fills the TARGET_SIZE entirely with no black bars
45     # and a preserved aspect ratio.
46     w = tf.shape(image)[0]
47     h = tf.shape(image)[1]
48     tw = TARGET_SIZE[1]
49     th = TARGET_SIZE[0]
50     resize_crit = (w * th) / (h * tw)
51     image = tf.cond(resize_crit < 1,
52                     lambda: tf.image.resize(image, [w*tw/w, h*tw/w]),
53                     lambda: tf.image.resize(image, [w*th/h, h*th/h]))
54
55     nw = tf.shape(image)[0]
56     nh = tf.shape(image)[1]
57     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
58     return image, label
59
60 def recompress_image(image, label):
61     # this reduces the amount of data, but takes some time
```

```

62     image = tf.cast(image, tf.uint8)
63     image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampled=False)
64     return image, label
65
66 # ii) Replacing the TensorFlow Dataset objects with RDDs, starting with
67
68 filenames = tf.io.gfile.glob(GCS_PATTERN)
69
70 # Get spark context for RDDs
71 sc = pyspark.SparkContext.getOrCreate()
72
73 # Create RDD for files
74 filenames_rdd = sc.parallelize(filenames)
75
76 # iii) Sampling the the RDD to a smaller number at an appropriate position
77
78 # Sampling the rdd
79 sample_rdd = filenames_rdd.sample(False, 0.02)
80
81 # Adapting tensorflow to spark, rdd for decode jpeg and label, rdd for
82 decode_jpeg_and_label_rdd = filenames_rdd.map(decode_jpeg_and_label)
83 resize_and_crop_image_rdd = decode_jpeg_and_label_rdd.map(lambda x: resize_and_crop(x))
84 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress(x))
85
86 # iv) Then use the functions from above to write the TFRecord files, using
87
88 # functions for writing TFRecord entries
89 # Feature values are always stored as lists, a single data element will
90 def _bytestring_feature(list_of_bytess):
91     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytess))
92
93 def _int_feature(list_of_ints): # int64
94     return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
95
96 def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data record
97     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order)
98     one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0]
99     feature = {
100         "image": _bytestring_feature([img_bytes]), # one image in the list
101         "class": _int_feature([class_num]) #, # one class in the list
102     }
103     return tf.train.Example(features=tf.train.Features(feature=feature))
104
105 def write_tfrecords(partition_index,partition): # write the images to
106     # tt0 = time.time()
107     # good practice to have the number of records in the filename
108     filename = GCS_OUTPUT + "{}.tfrec".format(partition_index)
109     # (CHANGED) You need to change GCS_OUTPUT to your own bucket to actually write
110     with tf.io.TFRecordWriter(filename) as out_file:
111         for element in partition:
112             image=element[0]
113             label=element[1]
114             example = to_tfrecord(out_file,
115                                   image.numpy(), # re-compressed image: already compressed by tensorflow
116                                   label.numpy() #
117                                   )
118             out_file.write(example.SerializeToString())
119             #print("Wrote file {} containing {} records".format(filename, len(example)))
120             #print("Total time: "+str(time.time()-tt0))
121     return (filename)
122

```

```
123 # v) The code for writing to the TFRecord files needs to be put into a  
124 partitions_rdd = recompress_image_rdd.repartition(PARTITIONS)  
125 filenames_rdd = recompress_image_rdd.mapPartitionsWithIndex(write_tfrec  
126 print("Writing TFRecords")
```

Writing spark\_write\_tfrec.py

### 1c) Set up a cluster and run the script. (6%)

Following the example from the labs, set up a cluster to run PySpark jobs in the cloud. You need to set up so that TensorFlow is installed on all nodes in the cluster.

#### i) Single machine cluster

Set up a cluster with a single machine using the maximal SSD size (100) and 8 vCPUs.

Enable **package installation** by passing a flag `--initialization-actions` with argument `gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh` (this is a public script that will read metadata to determine which packages to install). Then, the **packages are specified** by providing a `--metadata` flag with the argument `PIP_PACKAGES=tensorflow==2.4.0`.

Note: consider using `PIP_PACKAGES="tensorflow numpy"` or `PIP_PACKAGES=tensorflow` in case an older version of tensorflow is causing issues.

When the cluster is running, run your script to check that it works and keep the output cell output. (3%)

```
In [ ]: 1 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

**ERROR:** (gcloud.dataproc.clusters.delete) NOT\_FOUND: Not found: Cluster projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster

```
In [ ]: 1      ##### CODING TASK #####
2      #Set up a cluster with a single machine using the maximal SSD size (100
3      #REGION = 'us-central1'
4      !gcloud dataproc clusters create $CLUSTER \
5      --bucket $PROJECT-storage \
6      --image-version 1.4-ubuntu18 --single-node \
7      --master-machine-type n1-standard-8 \
8      --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \
9      --max-idle 3600s \
10     --initialization-actions gs://goog-dataproc-initialization-actions-$REG
11     --metadata PIP_PACKAGES=tensorflow==2.4.0
12
13     #--metadata "PIP_PACKAGES=tensorflow==2.4.0 scipy==1.4.1 matplotlib==3.
14
15     # to grant uniform access !gcloud storage buckets update $BUCKET --unif
```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/6ae51459-595f-3852-89e6-4c0d81876966].

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-c].

```
In [ ]: 1      #!gCloud storage buckets update $BUCKET --uniform-bucket-Level-access
```

```
In [ ]: 1      # Description of the cluster
2      !gcloud dataproc clusters describe $CLUSTER
```

```
clusterName: big-data-cw-420116-cluster
clusterUuid: 084e9d65-6812-45a2-9cab-372f310ba794
config:
  configBucket: big-data-cw-420116-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
  metadata:
    PIP_PACKAGES: tensorflow==2.4.0
  networkUri: https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default (https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default)
    serviceAccountScopes:
      - https://www.googleapis.com/auth/bigquery (https://www.googleapis.com/auth/bigquery)
      - https://www.googleapis.com/auth/bigtable.admin.table (https://www.googleapis.com/auth/bigtable.admin.table)
      - https://www.googleapis.com/auth/bigtable.data (https://www.googleapis.com/auth/bigtable.data)
```

Run the script in the cloud and test the output.

In [ ]:

```
1 # Submit jobs to cluster and measure execution time (Checked& NOT Working)
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \spark_write_tf
3 %time
```

Job [e72d8ce747414b90b5b0ef0b944ed255] submitted.  
Waiting for job output...  
2024-05-04 09:45:21.057528: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dle error: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 09:45:21.057575: I tensorflow/stream\_executor/cuda/cudart\_stub.cc:29] Ignore above cudart dleerror if you do not have a GPU set up on your machine.  
24/05/04 09:45:23 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 09:45:23 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 09:45:23 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 09:45:23 INFO org.spark\_project.jetty.util.log: Logging initialized @5834ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 09:45:24 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_312-b07  
24/05/04 09:45:24 INFO org.spark\_project.jetty.server.Server: Started @594 3ms  
24/05/04 09:45:24 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@71dbb1cc{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
24/05/04 09:45:24 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.  
24/05/04 09:45:25 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.15.220:8032  
24/05/04 09:45:25 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.15.220:10200  
24/05/04 09:45:27 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714815815041\_0001  
Writing TFRecords  
24/05/04 09:45:34 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@71dbb1cc{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
Job [e72d8ce747414b90b5b0ef0b944ed255] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/084e9d65-6812-45a2-9cab-372f310ba794/jobs/e72d8ce747414b90b5b0ef0b944ed255/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/084e9d65-6812-45a2-9cab-372f310ba794/jobs/e72d8ce747414b90b5b0ef0b944ed255/driveroutput  
jobUuid: 96b9484b-26bb-352a-a9f4-84042133fb34  
placement:  
    clusterName: big-data-cw-420116-cluster  
    clusterUuid: 084e9d65-6812-45a2-9cab-372f310ba794  
pysparkJob:  
    mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/084e9d65-6812-45a2-9cab-372f310ba794/jobs/e72d8ce747414b90b5b0ef0b944ed255/staging/spark\_write\_tfrec.py  
reference:  
    jobId: e72d8ce747414b90b5b0ef0b944ed255  
    projectId: big-data-cw-420116  
status:  
    state: DONE  
    stateStartTime: '2024-05-04T09:45:35.737564Z'  
statusHistory:

```

- state: PENDING
  stateStartTime: '2024-05-04T09:45:16.346354Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-04T09:45:16.397866Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-04T09:45:16.694290Z'
yarnApplications:
- name: spark_write_tfrec.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714815815041\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714815815041\_0001/)
CPU times: user 4 µs, sys: 0 ns, total: 4 µs
Wall time: 9.54 µs

```

*\*Got the following Error: \**

ERROR: (gcloud.dataproc.jobs.submit.pyspark) RESOURCE\_EXHAUSTED: Multiple validation errors:

- No agent on master node(s) reported to be active
- Unable to submit job, cluster 'big-data-cw-420116-cluster' is in state ERROR and cannot accept jobs. CPU times: user 5 µs, sys: 1 µs, total: 6 µs Wall time: 11 µs

In the free credit tier on Google Cloud, there are normally the following **restrictions** on compute machines:

- max 100GB of *SSD persistent disk*
- max 2000GB of *standard persistent disk*
- max 8 *vCPUs*
- no GPUs

See [here \(<https://cloud.google.com/free/docs/gcp-free-tier#free-trial>\)](https://cloud.google.com/free/docs/gcp-free-tier#free-trial) for details. The **disks are virtual disks**, where **I/O speed is limited in proportion to the size**, so we should allocate them evenly. This has mainly an effect on the **time the cluster needs to start**, as we are reading the data mainly from the bucket and we are not writing much to disk at all.

## ii) Maximal cluster

Use the **largest possible cluster** within these constraints, i.e. **1 master and 7 worker nodes**. Each of them with 1 (virtual) CPU. The master should get the full SSD capacity and the 7 worker nodes should get equal shares of the *standard* disk capacity to maximise throughput.

Once the cluster is running, test your script. (3%)

```
In [ ]: 1 #Was getting error to create the cluster so I deleted it to create it a  
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/e0a4734e-8fd0-3d80-a509-db4107d1b7c9].  
Deleted [https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster].

```
In [ ]: 1          ### CODING TASK ### (Checked & NOT Working)  
2 !gcloud dataproc clusters create $CLUSTER \  
3   --bucket $PROJECT-storage \  
4   --image-version 1.4-ubuntu18 \  
5   --master-machine-type n1-standard-1 \  
6   --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \  
7   --num-workers 8 --worker-machine-type n1-standard-1 --worker-boot-dis  
8   --max-idle 3600s \  
9   --initialization-actions gs://goog-dataproc-initialization-actions-$R  
10  --metadata PIP_PACKAGES=tensorflow==2.4.0  
11  # --metadata PIP_PACKAGES="scipy tensorflow==2.4.0 matplotlib numpy" s  
12
```

**ERROR:** (gcloud.dataproc.clusters.create) INVALID\_ARGUMENT: Insufficient 'I\_N\_USE\_ADDRESSES' quota. Requested 9.0, available 8.0. Your resource request exceeds your available quota. See <https://cloud.google.com/compute/resource-usage.> (https://cloud.google.com/compute/resource-usage.) Use [https://cloud.google.com/docs/quotas/view-manage#requesting\\_higher\\_quota](https://cloud.google.com/docs/quotas/view-manage#requesting_higher_quota) ([https://cloud.google.com/docs/quotas/view-manage#requesting\\_higher\\_quota](https://cloud.google.com/docs/quotas/view-manage#requesting_higher_quota)) to request additional quota.

```
In [ ]: 1 # Then I got this error: ERROR: (gcloud.dataproc.clusters.create) Opera  
2 # Hence I gave another try with a different configuration, reduced the
```

```
In [ ]: 1 # Also I was getting error again to create the cluster so I deleted (I  
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

**ERROR:** (gcloud.dataproc.clusters.delete) NOT\_FOUND: Not found: Cluster projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster

```
In [ ]: 1  ### CODING TASK ### (Checked & NOT Working)
2  !gcloud dataproc clusters create $CLUSTER \
3    --bucket $PROJECT-storage \
4    --image-version 1.4-ubuntu18 \
5    --master-machine-type n1-standard-1 \
6    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
7    --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-d
8    --max-idle 3600s \
9    --initialization-actions gs://goog-dataproc-initialization-actions-
10   --metadata PIP_PACKAGES=tensorflow==2.4.0
11  #   --metadata PIP_PACKAGES= "scipy tensorflow=2.4.0 matplotlib numpy"
```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/349fffc1-39f7-3eb5-8b35-ab3b1c6a06c5].

**WARNING:** Creating clusters using the n1-standard-1 machine type is not recommended. Consider using a machine type with higher memory.

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** For PD-Standard without local SSDs, we strongly recommend provisioning 1TB or larger to ensure consistently high I/O performance. See <https://cloud.google.com/compute/docs/disks/performance> (<https://cloud.google.com/compute/docs/disks/performance>) for information on disk I/O performance.

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-f].

```
In [ ]: 1  # Got another Error; ERROR: (gCloud.dataproc.clusters.create) INVALID_A
2  # This error has delayed my "own deadline" meaninfully.
```

In [ ]:

```
1 # Description of the cluster
2 !gcloud dataproc clusters describe $CLUSTER
```

clusterName: big-data-cw-420116-cluster  
clusterUuid: 10948a56-720a-449e-b3a2-a1c82a9b80b1  
config:  
 configBucket: big-data-cw-420116-storage  
 endpointConfig: {}  
 gceClusterConfig:  
 internalIpOnly: false  
 metadata:  
 PIP\_PACKAGES: tensorflow==2.4.0  
 networkUri: <https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default> (<https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default>)  
 serviceAccountScopes:  
 - <https://www.googleapis.com/auth/bigquery> (<https://www.googleapis.com/auth/bigquery>)  
 - <https://www.googleapis.com/auth/bigtable.admin.table> (<https://www.googleapis.com/auth/bigtable.admin.table>)  
 - <https://www.googleapis.com/auth/bigtable.data> (<https://www.googleapis.com/auth/bigtable.data>)  
 ...

In [ ]:

```
1 # Submit jobs to cluster and measure execution time (Checked& NOT Working)
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \spark_write_tf
3 %time
```

Job [6752114ba5f24f2b973033fa70af0c81] submitted.  
Waiting for job output...  
2024-05-04 09:54:57.116534: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dle error: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 09:54:57.116704: I tensorflow/stream\_executor/cuda/cudart\_stub.cc:29] Ignore above cudart dleerror if you do not have a GPU set up on your machine.  
24/05/04 09:55:00 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 09:55:00 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 09:55:00 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 09:55:01 INFO org.spark\_project.jetty.util.log: Logging initialized @8522ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 09:55:01 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_312-b07  
24/05/04 09:55:01 INFO org.spark\_project.jetty.server.Server: Started @8769ms  
24/05/04 09:55:01 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@586fc5d{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
24/05/04 09:55:01 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.  
24/05/04 09:55:03 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.15.224:8032  
24/05/04 09:55:03 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.15.224:10200  
24/05/04 09:55:07 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714816345709\_0001  
Writing TFRecords  
24/05/04 09:55:23 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@586fc5d{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
Job [6752114ba5f24f2b973033fa70af0c81] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/10948a56-720a-449e-b3a2-a1c82a9b80b1/jobs/6752114ba5f24f2b973033fa70af0c81/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/10948a56-720a-449e-b3a2-a1c82a9b80b1/jobs/6752114ba5f24f2b973033fa70af0c81/driveroutput  
jobUuid: 87f0332b-9611-313d-b274-ea81ca7cf65e  
placement:  
    clusterName: big-data-cw-420116-cluster  
    clusterUuid: 10948a56-720a-449e-b3a2-a1c82a9b80b1  
pysparkJob:  
    mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/10948a56-720a-449e-b3a2-a1c82a9b80b1/jobs/6752114ba5f24f2b973033fa70af0c81/staging/spark\_write\_tfrec.py  
reference:  
    jobId: 6752114ba5f24f2b973033fa70af0c81  
    projectId: big-data-cw-420116  
status:  
    state: DONE  
    stateStartTime: '2024-05-04T09:55:27.154303Z'  
statusHistory:

```

- state: PENDING
  stateStartTime: '2024-05-04T09:54:47.732174Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-04T09:54:47.765289Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-04T09:54:50.151079Z'
yarnApplications:
- name: spark_write_tfrec.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714816345709\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714816345709\_0001/)
CPU times: user 4 µs, sys: 0 ns, total: 4 µs
Wall time: 8.58 µs

```

## 1d) Optimisation, experiments, and discussion (17%)

### i) Improve parallelisation

If you implemented a straightforward version, you will **probably** observe that **all the computation** is done on only **two nodes**. This can be addressed by using the **second parameter** in the initial call to **parallelize**. Make the **suitable change** in the code you have written above and mark it up in comments as `### TASK 1d ###`.

Demonstrate the difference in cluster utilisation before and after the change based on different parameter values with **screenshots from Google Cloud** and measure the **difference in the processing time**. (6%)

### ii) Experiment with cluster configurations.

In addition to the experiments above (using 8 VMs), test your program with 4 machines with double the resources each (2 vCPUs, memory, disk) and 1 machine with eightfold resources. Discuss the results in terms of disk I/O and network bandwidth allocation in the cloud. (7%)

### iii) Explain the difference between this use of Spark and most standard applications like e.g. in our labs in terms of where the data is stored. What kind of parallelisation approach is used here? (4%)

Write the code below and your answers in the report.



In [ ]:

```
1                                     ### CODING TASK ### ("Done & Working")
2 # I) Improving parallelization; Script without second parameter
3 %%writefile task_1di.py
4 # Setting the working environment
5 import datetime
6 import math
7 import numpy as np
8 import os
9 import pyspark
10 import pickle
11 import pandas as pd
12 import random
13 import string
14 import sys
15 #import scipy as sp scipy give error in this task.
16 #import scipy.stats
17 import time
18 import tensorflow as tf
19
20 #from matplotlib import pyplot as plt
21 from pyspark.sql import SQLContext
22 from pyspark.sql import Row
23
24 # define variables
25 PROJECT = 'big-data-cw-420116' # my project id
26 BUCKET = 'gs://{}-storage'.format(PROJECT) # my own bucket storage
27 CLUSTER = '{}-cluster'.format(PROJECT)
28 CLASSES = [b'daisy', b'dandelion', b'roses', b'sunflowers', b'tulips']
29 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input jpg
30 PARTITIONS = 16 # no of partitions we will use later
31 TARGET_SIZE = [192, 192] # target resolution for the images
32 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for output
33
34 def decode_jpeg_and_label(filepath):
35     # extracts the image data and creates a class label, based on the folder name
36     bits = tf.io.read_file(filepath)
37     image = tf.image.decode_jpeg(bits)
38     # parse flower name from containing directory
39     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
40     label2 = label.values[-2]
41     return image, label2
42
43 def resize_and_crop_image(image, label):
44     # Resizes and cropd using "fill" algorithm:
45     # always make sure the resulting image is cut out from the source image
46     # so that it fills the TARGET_SIZE entirely with no black bars
47     # and a preserved aspect ratio.
48     w = tf.shape(image)[0]
49     h = tf.shape(image)[1]
50     tw = TARGET_SIZE[1]
51     th = TARGET_SIZE[0]
52     resize_crit = (w * th) / (h * tw)
53     image = tf.cond(resize_crit < 1,
54                     lambda: tf.image.resize(image, [w*tw/w, h*tw/w]),
55                     lambda: tf.image.resize(image, [w*th/h, h*th/h]))
56
57     nw = tf.shape(image)[0]
58     nh = tf.shape(image)[1]
59     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh - th) // 2, tw, th)
60
61     return image, label
```

```

62 def recompress_image(image, label):
63     # this reduces the amount of data, but takes some time
64     image = tf.cast(image, tf.uint8)
65     image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampled=True)
66     return image, label
67
68 # Retrieve file paths from GCS using a pattern
69 filenames = tf.io.gfile.glob(GCS_PATTERN)
70
71 # Get spark context for RDDs
72 sc = pyspark.SparkContext.getOrCreate()
73
74 ### TASK 1d ###
75 filenames_rdd = sc.parallelize(filenames) # Adding a second parameter to parallelize
76
77 # Sampling the rdd
78 sample_rdd = filenames_rdd.sample(False, 0.02)
79
80 # Adapting tensorflow to spark, rdd for decode jpeg and label, rdd for
81 decode_jpeg_and_label_rdd = filenames_rdd.map(decode_jpeg_and_label)
82 resize_and_crop_image_rdd = decode_jpeg_and_label_rdd.map(lambda x: resize_and_crop(x))
83 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress(x))
84
85 # functions for writing TFRecord entries
86 # Feature values are always stored as lists, a single data element will
87 def _bytestring_feature(list_of_bytess):
88     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytess))
89
90 def _int_feature(list_of_ints): # int64
91     return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
92
93 def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data record
94     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order)
95     one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0]
96     feature = {
97         "image": _bytestring_feature([img_bytes]), # one image in the list
98         "class": _int_feature([class_num]) #, # one class in the list
99     }
100    return tf.train.Example(features=tf.train.Features(feature=feature))
101
102 def write_tfrecords(index,partition): # write the images to files.
103     tt0 = time.time()
104     # good practice to have the number of records in the filename
105     filename = GCS_OUTPUT + "{}.tfrec".format(index)
106     # (CHANGED) You need to change GCS_OUTPUT to your own bucket to actually write to it
107     with tf.io.TFRecordWriter(filename) as out_file:
108         for element in partition:
109             image=element[0]
110             label=element[1]
111             example = to_tfrecord(out_file,
112                                   image.numpy(), # re-compressed image: already compressed by tensorflow
113                                   label.numpy() #
114                                   )
115             out_file.write(example.SerializeToString())
116             #print("Wrote file {} containing {} records".format(filename, len(example)))
117             #print("Total time: "+str(time.time()-tt0))
118     return (filename)
119
120 ## Repartition RDD for parallel processing
121 partitions_rdd = recompress_image_rdd.repartition(PARTITIONS)
122 # Write TFRecords for each partition

```

```
123 | filenames_rdd = recompress_image_rdd.mapPartitionsWithIndex(write_tfrecord)
```

Writing task\_1di.py

```
In [ ]: 1 #Was getting error to create the cluster so I deleted it to create it again  
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/b24f3fe3-7437-3ca8-8f5a-9719893e7198].

Deleted [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>].

In [ ]:

```
1          ##### CODING TASK ##### (checked & NOT WORKING)
2 # Creating Dataproc cluster with 1 master machine
3 #CHECK WORKERS
4 REGION = 'europe-west2'
5 import time
6 start_time = time.time()
7
8 !gcloud dataproc clusters create $CLUSTER \
9     --bucket $PROJECT-storage \
10    --image-version 1.4-ubuntu18 \
11    --master-machine-type n1-standard-1 \
12    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
13    --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-d
14    --max-idle 3600s \
15    --initialization-actions gs://goog-dataproc-initialization-actions-
16    --metadata PIP_PACKAGES=tensorflow==2.4.0
17 %time
18 end_time = time.time()
19 execution_time = end_time - start_time
20 print(f"Execution time: {execution_time} seconds")
```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/6be162fe-94d2-3982-8b87-c0e5d9070312].

**WARNING:** Creating clusters using the n1-standard-1 machine type is not recommended. Consider using a machine type with higher memory.

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-europe-west2/python/pip-install.sh

**WARNING:** For PD-Standard without local SSDs, we strongly recommend provisioning 1TB or larger to ensure consistently high I/O performance. See [http://cloud.google.com/compute/docs/disks/performance](https://cloud.google.com/compute/docs/disks/performance) (<https://cloud.google.com/compute/docs/disks/performance>) for information on disk I/O performance.

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-c].

CPU times: user 4 µs, sys: 0 ns, total: 4 µs

Wall time: 7.87 µs

Execution time: 267.9169752597809 seconds

In [ ]:

```
1 # Description of the cluster
2 !gcloud dataproc clusters describe $CLUSTER
```

clusterName: big-data-cw-420116-cluster
clusterUuid: 0c6fd069-69d3-48ae-8962-fff9aa7af8ab
config:
 configBucket: big-data-cw-420116-storage
 endpointConfig: {}
 gceClusterConfig:
 internalIpOnly: false
 metadata:
 PIP\_PACKAGES: tensorflow==2.4.0
 networkUri: <https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default> (<https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default>)
 serviceAccountScopes:
 - <https://www.googleapis.com/auth/bigquery> (<https://www.googleapis.com/auth/bigquery>)
 - <https://www.googleapis.com/auth/bigtable.admin.table> (<https://www.googleapis.com/auth/bigtable.admin.table>)
 - <https://www.googleapis.com/auth/bigtable.data> (<https://www.googleapis.com/auth/bigtable.data>)
 ...

In [ ]:

```
1 # Submit the job
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \task_1di.py
```

Job [453a18bd37b44e1dacd6debb23782a0b] submitted.  
Waiting for job output...  
2024-05-04 10:02:33.418059: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dle error: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 10:02:33.418234: I tensorflow/stream\_executor/cuda/cudart\_stub.cc:29] Ignore above cudart dleerror if you do not have a GPU set up on your machine.  
24/05/04 10:02:37 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 10:02:37 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 10:02:37 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 10:02:38 INFO org.spark\_project.jetty.util.log: Logging initialized @10594ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 10:02:38 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_312-b07  
24/05/04 10:02:38 INFO org.spark\_project.jetty.server.Server: Started @110 25ms  
24/05/04 10:02:38 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@24703cc0{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
24/05/04 10:02:39 WARN org.apache.spark.scheduler.FairSchedulerBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.  
24/05/04 10:02:41 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.15.233:8032  
24/05/04 10:02:41 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.15.233:10200  
24/05/04 10:02:45 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714816784694\_0001  
24/05/04 10:03:00 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@24703cc0{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
Job [453a18bd37b44e1dacd6debb23782a0b] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/0c6fd069-69d3-48ae-8962-fff9aa7af8ab/jobs/453a18bd37b44e1dacd6debb23782a0b/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/0c6fd069-69d3-48ae-8962-fff9aa7af8ab/jobs/453a18bd37b44e1dacd6debb23782a0b/driveroutput  
jobUuid: f4fa366e-15db-3dc1-b7d6-c96eb129f452  
placement:  
    clusterName: big-data-cw-420116-cluster  
    clusterUuid: 0c6fd069-69d3-48ae-8962-fff9aa7af8ab  
pysparkJob:  
    mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/0c6fd069-69d3-48ae-8962-fff9aa7af8ab/jobs/453a18bd37b44e1dacd6debb23782a0b/staging/task\_1di.py  
reference:  
    jobId: 453a18bd37b44e1dacd6debb23782a0b  
    projectId: big-data-cw-420116  
status:  
    state: DONE  
    stateStartTime: '2024-05-04T10:03:05.167057Z'  
statusHistory:  
- state: PENDING

```
stateStartTime: '2024-05-04T10:02:24.030519Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-04T10:02:24.074325Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-04T10:02:24.463911Z'
yarnApplications:
- name: task_1di.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714816784694\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714816784694\_0001/)
```



In [ ]:

```
1                                     ### CODING TASK ### ("Done & Working"
2 # I) improving parallelization, with second parameter
3 %%writefile task_1di.py
4 # Setting the working environment
5 import datetime
6 import math
7 import numpy as np
8 import os
9 import pyspark
10 import pickle
11 import pandas as pd
12 import random
13 import string
14 import sys
15 #import scipy as sp scipy give error in this task.
16 #import scipy.stats
17 import time
18 import tensorflow as tf
19
20 #from matplotlib import pyplot as plt
21 from pyspark.sql import SQLContext
22 from pyspark.sql import Row
23
24 # define variables
25 PROJECT = 'big-data-cw-420116' # my project id
26 BUCKET = 'gs://{}-storage'.format(PROJECT) # my own bucket storage
27 CLUSTER = '{}-cluster'.format(PROJECT)
28 CLASSES = [b'daisy', b'dandelion', b'roses', b'sunflowers', b'tulips']
29 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input jpg
30 PARTITIONS = 16 # no of partitions we will use later
31 TARGET_SIZE = [192, 192] # target resolution for the images
32 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for output
33
34 def decode_jpeg_and_label(filepath):
35     # extracts the image data and creates a class label, based on the folder name
36     bits = tf.io.read_file(filepath)
37     image = tf.image.decode_jpeg(bits)
38     # parse flower name from containing directory
39     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
40     label2 = label.values[-2]
41     return image, label2
42
43 def resize_and_crop_image(image, label):
44     # Resizes and cropd using "fill" algorithm:
45     # always make sure the resulting image is cut out from the source image
46     # so that it fills the TARGET_SIZE entirely with no black bars
47     # and a preserved aspect ratio.
48     w = tf.shape(image)[0]
49     h = tf.shape(image)[1]
50     tw = TARGET_SIZE[1]
51     th = TARGET_SIZE[0]
52     resize_crit = (w * th) / (h * tw)
53     image = tf.cond(resize_crit < 1,
54                     lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), #
55                     lambda: tf.image.resize(image, [w*th/h, h*th/h])) #
56
57     nw = tf.shape(image)[0]
58     nh = tf.shape(image)[1]
59     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh - th) // 2, tw, th)
60
61     return image, label
```

```

62 def recompress_image(image, label):
63     # this reduces the amount of data, but takes some time
64     image = tf.cast(image, tf.uint8)
65     image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampled=True)
66     return image, label
67
68 # Retrieve file paths from GCS using a pattern
69 filenames = tf.io.gfile.glob(GCS_PATTERN)
70
71 # Get spark context for RDDs
72 sc = pyspark.SparkContext.getOrCreate()
73
74 ### TASK 1d ###
75 filenames_rdd = sc.parallelize(filenames,16) # Adding a second parameter
76
77 # Sampling the rdd
78 sample_rdd = filenames_rdd.sample(False, 0.02)
79
80 # Adapting tensorflow to spark, rdd for decode jpeg and label, rdd for
81 decode_jpeg_and_label_rdd = filenames_rdd.map(decode_jpeg_and_label)
82 resize_and_crop_image_rdd = decode_jpeg_and_label_rdd.map(lambda x: resize_and_crop_image(x))
83 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress_image(x))
84
85 # functions for writing TFRecord entries
86 # Feature values are always stored as lists, a single data element will
87 def _bytestring_feature(list_of_bytess):
88     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytess))
89
90 def _int_feature(list_of_ints): # int64
91     return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
92
93 def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data record
94     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order)
95     one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0, 0]
96     feature = {
97         "image": _bytestring_feature([img_bytes]), # one image in the list
98         "class": _int_feature([class_num]) #, # one class in the list
99     }
100    return tf.train.Example(features=tf.train.Features(feature=feature))
101
102 def write_tfrecords(index,partition): # write the images to files.
103     tt0 = time.time()
104     # good practice to have the number of records in the filename
105     filename = GCS_OUTPUT + "{}.tfrec".format(index)
106     # (CHANGED) You need to change GCS_OUTPUT to your own bucket to actually write to it
107     with tf.io.TFRecordWriter(filename) as out_file:
108         for element in partition:
109             image=element[0]
110             label=element[1]
111             example = to_tfrecord(out_file,
112                                   image.numpy(), # re-compressed image: already compressed by tensorflow
113                                   label.numpy() #
114                                   )
115             out_file.write(example.SerializeToString())
116             #print("Wrote file {} containing {} records".format(filename, len(example)))
117             #print("Total time: "+str(time.time()-tt0))
118     return (filename)
119
120 ## Repartition RDD for parallel processing
121 partitions_rdd = recompress_image_rdd.repartition(PARTITIONS)
122 # Write TFRecords for each partition

```

```
123 | filenames_rdd = recompress_image_rdd.mapPartitionsWithIndex(write_tfrecord)
```

Overwriting task\_1di.py

```
In [ ]: 1 #Was getting error to create the cluster so I deleted it to create it again  
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/1b4bd492-ec11-378a-8a5d-0bd761bbd85a].

Deleted [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>].

In [ ]:

```
1 # Creating Dataproc cluster with 1 master machine
2 #CHECK WORKERS
3 import time
4 start_time = time.time()
5
6 !gcloud dataproc clusters create $CLUSTER \
7     --bucket $PROJECT-storage \
8     --image-version 1.4-ubuntu18 \
9     --master-machine-type n1-standard-1 \
10    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
11    --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-d
12    --max-idle 3600s \
13    --initialization-actions gs://goog-dataproc-initialization-actions-
14    --metadata PIP_PACKAGES=tensorflow==2.4.0
15 %time
16 end_time = time.time()
17 execution_time = end_time - start_time
18 print(f"Execution time: {execution_time} seconds")
19
20 #--initialization-actions gs://goog-dataproc-initialization-actions-$RE
```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/3d02d1cf-b270-3321-8a4a-8808a97860b6].

**WARNING:** Creating clusters using the n1-standard-1 machine type is not recommended. Consider using a machine type with higher memory.

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** For PD-Standard without local SSDs, we strongly recommend provisioning 1TB or larger to ensure consistently high I/O performance. See [http://cloud.google.com/compute/docs/disks/performance](https://cloud.google.com/compute/docs/disks/performance) (<https://cloud.google.com/compute/docs/disks/performance>) for information on disk I/O performance.

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-f].

CPU times: user 4 µs, sys: 0 ns, total: 4 µs

Wall time: 7.63 µs

Execution time: 262.41806292533875 seconds

In [ ]:

```
1 # Description of the cluster
2 !gcloud dataproc clusters describe $CLUSTER
```

clusterName: big-data-cw-420116-cluster  
clusterUuid: 972c0386-19ae-4bf7-9b33-eaf72fc8bb24  
config:  
 configBucket: big-data-cw-420116-storage  
 endpointConfig: {}  
 gceClusterConfig:  
 internalIpOnly: false  
 metadata:  
 PIP\_PACKAGES: tensorflow==2.4.0  
 networkUri: <https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default> (<https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default>)  
 serviceAccountScopes:  
 - <https://www.googleapis.com/auth/bigquery> (<https://www.googleapis.com/auth/bigquery>)  
 - <https://www.googleapis.com/auth/bigtable.admin.table> (<https://www.googleapis.com/auth/bigtable.admin.table>)  
 - <https://www.googleapis.com/auth/bigtable.data> (<https://www.googleapis.com/auth/bigtable.data>)  
 ...

In [ ]:

```
1 # Submit the job
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \task_1di.py
```

```
Job [fd6408a42784471c99bf6858fe21f276] submitted.
Waiting for job output...
2024-05-04 10:23:48.602436: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dlopen error: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: :/usr/lib/hadoop/lib/native
2024-05-04 10:23:48.602664: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dlerror if you do not have a GPU set up on your machine.
24/05/04 10:23:52 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/04 10:23:52 INFO org.apache.spark.SparkEnv: Registering BlockManager Master
24/05/04 10:23:52 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator
24/05/04 10:23:53 INFO org.spark_project.jetty.util.log: Logging initialized @9699ms to org.spark_project.jetty.util.log.Slf4jLog
24/05/04 10:23:53 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_312-b07
24/05/04 10:23:53 INFO org.spark_project.jetty.server.Server: Started @9996ms
24/05/04 10:23:53 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@586fc5d{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}
24/05/04 10:23:53 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.
24/05/04 10:23:55 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.15.235:8032
24/05/04 10:23:56 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.15.235:10200
24/05/04 10:23:57 WARN org.apache.hadoop.hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
        at java.lang.Object.wait(Native Method)
        at java.lang.Thread.join(Thread.java:1252)
        at java.lang.Thread.join(Thread.java:1326)
        at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:980)
        at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:630)
        at org.apache.hadoop.hdfs.DataStreamer.run(DataStreamer.java:807)
24/05/04 10:24:00 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714818062265_0001
24/05/04 10:24:15 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@586fc5d{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}
Job [fd6408a42784471c99bf6858fe21f276] finished successfully.
done: true
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/972c0386-19ae-4bf7-9b33-eaf72fc8bb24/jobs/fd6408a42784471c99bf6858fe21f276/
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/972c0386-19ae-4bf7-9b33-eaf72fc8bb24/jobs/fd6408a42784471c99bf6858fe21f276/driveroutput
jobUuid: 2a6bc75a-2d86-37bd-8490-fb94bab0e2e1
placement:
    clusterName: big-data-cw-420116-cluster
    clusterUuid: 972c0386-19ae-4bf7-9b33-eaf72fc8bb24
pysparkJob:
```

```
mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/972c0386-19ae-4bf7-9b33-eaf72fc8bb24/jobs/fd6408a42784471c99bf6858fe21f276/staging/task_1di.py
reference:
  jobId: fd6408a42784471c99bf6858fe21f276
  projectId: big-data-cw-420116
status:
  state: DONE
  stateStartTime: '2024-05-04T10:24:20.101716Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-04T10:23:40.559025Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-04T10:23:40.598939Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-04T10:23:41.001805Z'
yarnApplications:
- name: task_1di.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application_1714818062265_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application_1714818062265_0001/)
```

In [ ]: 1 # ii) Experiment with cluster configurations

In [ ]: 1 #Was getting error to create the cluster so I deleted it to create it a  
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

```
Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/d9c4b631-2d94-312f-bd68-cf1d9e2cf841].
Deleted [https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster].
```

```

In [ ]: 1 # II) 4 machines with double the resources each (2 vCPUs, memory, disk)
2 import time
3 start_time = time.time()
4
5 REGION = 'us-central1'
6 !gcloud dataproc clusters create $CLUSTER \
7   --bucket $PROJECT-storage \
8   --image-version 1.4-ubuntu18 \
9   --master-machine-type n1-standard-2 \
10  --master-boot-disk-type pd-ssd --master-boot-disk-size 200 \
11  --num-workers 3 --worker-machine-type n1-standard-2 --worker-boot-di \
12  --max-idle 3600s \
13  --initialization-actions gs://goog-dataproc-initialization-actions-$R \
14  --metadata PIP_PACKAGES=tensorflow==2.1.0
15 %time
16
17 end_time = time.time()
18 execution_time = end_time - start_time
19 print(f"Execution time: {execution_time} seconds")
20
21
22 # Extra resources checked:
23 # https://cloud.google.com/dataproc/docs/concepts/configuring-cluster
24 # https://cloud.google.com/dataproc/docs/concepts/configuring-cluster

```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/71b34dbc-bc8e-3e7b-b8af-58f7dcdd89d2].

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** For PD-Standard without local SSDs, we strongly recommend provisioning 1TB or larger to ensure consistently high I/O performance. See <https://cloud.google.com/compute/docs/disks/performance> (<https://cloud.google.com/compute/docs/disks/performance>) for information on disk I/O performance.

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-c].

CPU times: user 3 µs, sys: 0 ns, total: 3 µs

Wall time: 7.39 µs

Execution time: 216.55400729179382 seconds

In [ ]:

```
1 # Description of the cluster
2 !gcloud dataproc clusters describe $CLUSTER
```

clusterName: big-data-cw-420116-cluster
clusterUuid: 35a7b3f0-8d5f-4479-8350-1c3eb784ab9c
config:
 configBucket: big-data-cw-420116-storage
 endpointConfig: {}
 gceClusterConfig:
 internalIpOnly: false
 metadata:
 PIP\_PACKAGES: tensorflow==2.1.0
 networkUri: <https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default> (<https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default>)
 serviceAccountScopes:
 - <https://www.googleapis.com/auth/bigquery> (<https://www.googleapis.com/auth/bigquery>)
 - <https://www.googleapis.com/auth/bigtable.admin.table> (<https://www.googleapis.com/auth/bigtable.admin.table>)
 - <https://www.googleapis.com/auth/bigtable.data> (<https://www.googleapis.com/auth/bigtable.data>)
 ...

In [ ]:

```
1 # Submit the job
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \task_1di.py
```

Job [24d0a7f5f83f42168e9e29ed9b3f6531] submitted.  
Waiting for job output...  
2024-05-04 10:50:55.648290: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:55] Could not load dynamic library 'libnvinfer.so.6'; dlerror: libnvinfer.so.6: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 10:50:55.648495: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:55] Could not load dynamic library 'libnvinfer\_plugin.so.6'; dlerror: libnvinfer\_plugin.so.6: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 10:50:55.648511: W tensorflow/compiler/tf2tensorrt/utils/py\_utils.cc:30] Cannot dlopen some TensorRT libraries. If you would like to use Nvidia GPU with TensorRT, please make sure the missing libraries mentioned above are installed properly.  
24/05/04 10:50:58 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 10:50:58 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 10:50:58 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 10:50:58 INFO org.spark\_project.jetty.util.log: Logging initialized @7776ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 10:50:58 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_312-b07  
24/05/04 10:50:58 INFO org.spark\_project.jetty.server.Server: Started @7988ms  
24/05/04 10:50:58 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@7bbfd4c6{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
24/05/04 10:50:59 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.  
24/05/04 10:51:00 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.0.11:8032  
24/05/04 10:51:00 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.0.11:10200  
24/05/04 10:51:04 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714819715160\_0001  
24/05/04 10:51:15 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@7bbfd4c6{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
Job [24d0a7f5f83f42168e9e29ed9b3f6531] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35a7b3f0-8d5f-4479-8350-1c3eb784ab9c/jobs/24d0a7f5f83f42168e9e29ed9b3f6531/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35a7b3f0-8d5f-4479-8350-1c3eb784ab9c/jobs/24d0a7f5f83f42168e9e29ed9b3f6531/driveroutput  
jobUuid: dcba18a4-a317-38b0-a891-cfe11476470c  
placement:  
    clusterName: big-data-cw-420116-cluster  
    clusterUuid: 35a7b3f0-8d5f-4479-8350-1c3eb784ab9c  
pysparkJob:  
    mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35a7b3f0-8d5f-4479-8350-1c3eb784ab9c/jobs/24d0a7f5f83f42168e9e29ed9b3f6531/staging/task\_1di.py  
reference:  
    jobId: 24d0a7f5f83f42168e9e29ed9b3f6531  
    projectId: big-data-cw-420116

```
status:
  state: DONE
  stateStartTime: '2024-05-04T10:51:17.281536Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-04T10:50:48.756712Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-04T10:50:48.794364Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-04T10:50:49.102857Z'
yarnApplications:
- name: task_1di.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714819715160\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714819715160\_0001/)
```

In [ ]: 1 *#Was getting error to create the cluster so I deleted it to create it again*  
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/6cc2f6f3-169b-3ac1-a12b-bdc5700983e7].  
Deleted [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>].

```

In [ ]: 1 # 1 machine with eightfold resources, need to change the disk size to 1
2 import time
3 start_time = time.time()
4
5 !gcloud dataproc clusters create $CLUSTER \
6   --bucket $PROJECT-storage \
7   --image-version 1.4-ubuntu18 \
8   --master-machine-type n1-standard-8 \
9   --master-boot-disk-type pd-ssd --master-boot-disk-size 200 \
10  --num-workers 0 \
11  --max-idle 3600s \
12  --initialization-actions gs://goog-dataproc-initialization-actions-$R
13  --metadata PIP_PACKAGES=tensorflow==2.1.0
14 %time
15
16 end_time = time.time()
17 execution_time = end_time - start_time
18 print(f"Execution time: {execution_time} seconds")
19
20
21 # Extra resources checked:
22 # https://cloud.google.com/dataproc/docs/concepts/configuring-cluster
23 # https://cloud.google.com/dataproc/docs/concepts/configuring-cluster

```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/b17e8492-a683-3054-8601-4e2cbe27e111].

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-c].

CPU times: user 5 µs, sys: 0 ns, total: 5 µs

Wall time: 8.34 µs

Execution time: 155.76372051239014 seconds

```
In [ ]: 1 # Description of the cluster
          2 !gcloud dataproc clusters describe $CLUSTER

clusterName: big-data-cw-420116-cluster
clusterUuid: f61c93f4-ae56-4d6a-ab67-a56c38754889
config:
    configBucket: big-data-cw-420116-storage
    endpointConfig: {}
    gceClusterConfig:
        internalIpOnly: false
    metadata:
        PIP_PACKAGES: tensorflow==2.1.0
        networkUri: https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default (https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default)
        serviceAccountScopes:
            - https://www.googleapis.com/auth/bigquery (https://www.googleapis.com/auth/bigquery)
            - https://www.googleapis.com/auth/bigtable.admin.table (https://www.googleapis.com/auth/bigtable.admin.table)
            - https://www.googleapis.com/auth/bigtable.data (https://www.googleapis.com/auth/bigtable.data)
```

In [ ]:

```
1 # Submit the job
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \task_1di.py
```

Job [fe6f17559a204b2e9bbdddf69ca2c124] submitted.  
Waiting for job output...  
2024-05-04 10:57:02.237415: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:55] Could not load dynamic library 'libnvinfer.so.6'; dlerror: libnvinfer.so.6: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 10:57:02.237620: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:55] Could not load dynamic library 'libnvinfer\_plugin.so.6'; dlerror: libnvinfer\_plugin.so.6: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 10:57:02.237637: W tensorflow/compiler/tf2tensorrt/utils/py\_utils.cc:30] Cannot dlopen some TensorRT libraries. If you would like to use Nvidia GPU with TensorRT, please make sure the missing libraries mentioned above are installed properly.  
24/05/04 10:57:04 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 10:57:04 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 10:57:04 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 10:57:04 INFO org.spark\_project.jetty.util.log: Logging initialized @6119ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 10:57:04 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_312-b07  
24/05/04 10:57:04 INFO org.spark\_project.jetty.server.Server: Started @6229ms  
24/05/04 10:57:04 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@469879eb{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
24/05/04 10:57:04 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.  
24/05/04 10:57:05 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.0.13:8032  
24/05/04 10:57:05 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.0.13:10200  
24/05/04 10:57:09 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714820116470\_0001  
24/05/04 10:57:16 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@469879eb{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
Job [fe6f17559a204b2e9bbdddf69ca2c124] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/f61c93f4-ae56-4d6a-ab67-a56c38754889/jobs/fe6f17559a204b2e9bbdddf69ca2c124/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/f61c93f4-ae56-4d6a-ab67-a56c38754889/jobs/fe6f17559a204b2e9bbdddf69ca2c124/driveroutput  
jobUuid: b49bbb7b-6eab-34ac-b0a9-b7fc4b999f31  
placement:  
    clusterName: big-data-cw-420116-cluster  
    clusterUuid: f61c93f4-ae56-4d6a-ab67-a56c38754889  
pysparkJob:  
    mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/f61c93f4-ae56-4d6a-ab67-a56c38754889/jobs/fe6f17559a204b2e9bbdddf69ca2c124/staging/task\_1di.py  
reference:  
    jobId: fe6f17559a204b2e9bbdddf69ca2c124  
    projectId: big-data-cw-420116

```
status:
  state: DONE
  stateStartTime: '2024-05-04T10:57:17.829674Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-04T10:56:57.114331Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-04T10:56:57.154347Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-04T10:56:57.452071Z'
yarnApplications:
- name: task_1di.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714820116470\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714820116470\_0001/)
```

## Section 2: Speed tests

We have seen that **reading from the pre-processed TFRecord files** is **faster** than reading individual image files and decoding on the fly. This task is about **measuring this effect** and **parallelizing the tests with PySpark**.

### 2.1 Speed test implementation

Here is **code for time measurement** to determine the **throughput in images per second**. It doesn't render the images but extracts and prints some basic information in order to make sure the image data are read. We write the information to the null device for longer measurements `null_file=open("/dev/null", mode='w')`. That way it will not clutter our cell output.

We use batches (`dset2 = dset1.batch(batch_size)`) and select a number of batches with (`dset3 = dset2.take(batch_number)`). Then we use the `time.time()` to take the **time measurement** and take it multiple times, reading from the same dataset to see if reading speed changes with mutiple readings.

We then **vary** the size of the batch (`batch_size`) and the number of batches (`batch_number`) and **store the results for different values**. Store also the **results for each repetition** over the same dataset (repeat 2 or 3 times).

The speed test should be combined in a **function** `time_configs()` that takes a configuration, i.e. a dataset and arrays of `batch_sizes`, `batch_numbers`, and `repetitions` (an array of integers starting from 1), as **arguments** and runs the time measurement for each combination of `batch_size` and `batch_number` for the requested number of repetitions.

In [ ]:

```
1 ##### Checked #####
2 # Here are some useful values for testing your code, use higher values
3 batch_sizes = [2,4]
4 batch_numbers = [3,6]
5 repetitions = [1]
6
7 def time_configs(dataset, batch_sizes, batch_numbers, repetitions):
8     dims = [len(batch_sizes),len(batch_numbers),len(repetitions)]
9     print(dims)
10    results = np.zeros(dims)
11    params = np.zeros(dims + [3])
12    print( results.shape )
13    with open("/dev/null",mode='w') as null_file: # for printing the ou
14        tt = time.time() # for overall time taking
15        for bsi,bs in enumerate(batch_sizes):
16            for dsi, ds in enumerate(batch_numbers):
17                batched_dataset = dataset.batch(bs)
18                timing_set = batched_dataset.take(ds)
19                for ri,rep in enumerate(repetitions):
20                    print("bs: {}, ds: {}, rep: {}".format(bs,ds,rep))
21                    t0 = time.time()
22                    for image, label in timing_set:
23                        #print("Image batch shape {}".format(image.numpy().shape))
24                        print("Image batch shape {}, {}".format(image.numpy().shape,
25                            [str(lbl) for lbl in label.numpy()])), null_file
26                    td = time.time() - t0 # duration for reading images
27                    results[bsi,dsi,ri] = ( bs * ds ) / td
28                    params[bsi,dsi,ri] = [ bs, ds, rep ]
29    print("total time: "+str(time.time()-tt))
30    return results, params
```

Let's try this function with a **small number** of configurations of batch\_sizes batch\_numbers and repetitions, so that we get a set of parameter combinations and corresponding reading speeds. Try reading from the image files (dataset4) and the TFRecord files (datasetTfrec).

In [ ]:

```
1 ##### Checked #####
2
3 def decode_jpeg_and_label(filepath):
4     # extracts the image data and creates a class label, based on the f
5     bits = tf.io.read_file(filepath)
6     image = tf.image.decode_jpeg(bits)
7     # parse flower name from containing directory
8     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
9     label2 = label.values[-2]
10    return image, label2
11
12 def resize_and_crop_image(image, label):
13     # Resizes and cropd using "fill" algorithm:
14     # always make sure the resulting image is cut out from the source i
15     # so that it fills the TARGET_SIZE entirely with no black bars
16     # and a preserved aspect ratio.
17     w = tf.shape(image)[0]
18     h = tf.shape(image)[1]
19     tw = TARGET_SIZE[1]
20     th = TARGET_SIZE[0]
21     resize_crit = (w * th) / (h * tw)
22     image = tf.cond(resize_crit < 1,
23                     lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), #
24                     lambda: tf.image.resize(image, [w*th/h, h*th/h]) #
25                     )
26     nw = tf.shape(image)[0]
27     nh = tf.shape(image)[1]
28     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
29     return image, label
30
31 dataset4 = dsetResized.map(recompress_image)
32
33 #Time configurations for dataset4
34 [res,par] = time_configs(dataset4, batch_sizes, batch_numbers, repetiti
35 print(res)
36 print(par)
37
38 print("====")
39
40 # Time configurations for datasetTfrec
41 [res,par] = time_configs(datasetTfrec, batch_sizes, batch_numbers, repe
42 print(res)
43 print(par)
```

```

[2, 2, 1]
(2, 2, 1)
bs: 2, ds: 3, rep: 1
Image batch shape (2,), ["b'roses'", "b'roses'"]) <_io.TextIOWrapper name
='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'dandelion'", "b'daisy'"]) <_io.TextIOWrapper n
ame='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'roses'", "b'tulips'"]) <_io.TextIOWrapper name
='dev/null' mode='w' encoding='UTF-8'
bs: 2, ds: 6, rep: 1
Image batch shape (2,), ["b'daisy'", "b'roses'"]) <_io.TextIOWrapper name
='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'daisy'", "b'daisy'"]) <_io.TextIOWrapper name
='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'roses'", "b'dandelion'"]) <_io.TextIOWrapper n
ame='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'daisy'", "b'tulips'"]) <_io.TextIOWrapper name
='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'dandelion'", "b'dandelion'"]) <_io.TextIOWrapp
er name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (2,), ["b'tulips'", "b'dandelion'"]) <_io.TextIOWrapper
name='dev/null' mode='w' encoding='UTF-8'
bs: 4, ds: 3, rep: 1
Image batch shape (4,), ["b'sunflowers'", "b'daisy'", "b'roses'", "b'tulip
s'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (4,), ["b'dandelion'", "b'sunflowers'", "b'sunflowers'", "b
'sunflowers'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='U
TF-8'
Image batch shape (4,), ["b'roses'", "b'tulips'", "b'dandelion'", "b'dande
lion'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
bs: 4, ds: 6, rep: 1
Image batch shape (4,), ["b'dandelion'", "b'sunflowers'", "b'daisy'", "b't
ulips'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (4,), ["b'tulips'", "b'dandelion'", "b'daisy'", "b'dais
y'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (4,), ["b'tulips'", "b'tulips'", "b'tulips'", "b'dandeli
on'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (4,), ["b'daisy'", "b'tulips'", "b'daisy'", "b'roses'"])
<_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (4,), ["b'daisy'", "b'daisy'", "b'tulips'", "b'roses'"])
<_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
Image batch shape (4,), ["b'roses'", "b'roses'", "b'dandelion'", "b'dandel
ion'"]) <_io.TextIOWrapper name='dev/null' mode='w' encoding='UTF-8'
total time: 4.946155309677124
[[[ 6.70955344
    [11.65212834]

[[[10.95695471
    [12.6035298 ]]]
[[[[2. 3. 1.]]]

[[[2. 6. 1.]]]

[[[4. 3. 1.]]]

[[4. 6. 1.]]]
=====
[2, 2, 1]
(2, 2, 1)
bs: 2, ds: 3, rep: 1

```

```
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
bs: 2, ds: 6, rep: 1
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2, 192, 192, 3), ['0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
bs: 4, ds: 3, rep: 1
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
bs: 4, ds: 6, rep: 1
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
total time: 1.4839990139007568
[[[11.07662892]
 [36.53630608]

 [[39.92597951]
 [84.38543079]]
 [[[2. 3. 1.]]]

 [[2. 6. 1.]]]

 [[[4. 3. 1.]]]
 [[4. 6. 1.]]]
```

## Task 2: Parallelising the speed test with Spark in the cloud. (36%)

As an exercise in **Spark programming and optimisation** as well as **performance analysis**, we will now implement the **speed test** with multiple parameters in parallel with Spark. Running multiple tests in parallel would **not be a useful approach on a single machine, but it can be in the cloud** (you will be asked to reason about this later).

### 2a) Create the script (14%)

Your task is now to port the speed test above to Spark for running it in the cloud in Dataproc. **Adapt the speed testing** as a Spark program that performs the same actions as above, but with **Spark RDDs in a distributed way**. The distribution should be such that **each parameter combination (except repetition)** is processed in a separate Spark task.

More specifically:

- i) combine the previous cells to have the code to create a dataset and create a list of parameter combinations in an RDD (2%)
- ii) get a Spark context and create the dataset and run timing test for each combination in parallel (2%)
- iii) transform the resulting RDD to the structure ( parameter\_combination, images\_per\_second ) and save these values in an array (2%)
- iv) create an RDD with all results for each parameter as (parameter\_value,images\_per\_second) and collect the result for each parameter (2%)
- v) create an RDD with the average reading speeds for each parameter value and collect the results. Keep associativity in mind when implementing the average. (3%)
- vi) write the results to a pickle file in your bucket (2%)
- vii) Write your code it into a file using the *cell magic %%writefile spark\_job.py* (1%)

**Important:** The task here is not to parallelize the pre-processing, but to run multiple speed tests in parallel using Spark.



In [ ]:

```
1      ##### CODING TASK #     Checked
2  # Setting working environment
3  import pyspark
4
5  from pyspark.sql import Row
6  from pyspark.sql import SparkSession
7  from pyspark.sql import SQLContext
8
9  # Pre-defined function. Function to parse the TFRecord
10 def read_tfrecord(example):
11     features = {
12         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = b
13         "class": tf.io.FixedLenFeature([], tf.int64) #,   # shape [] me
14     }
15     # decode the TFRecord
16     example = tf.io.parse_single_example(example, features)
17     image = tf.image.decode_jpeg(example['image'], channels=3)
18     image = tf.reshape(image, [*TARGET_SIZE, 3])
19     class_num = example['class']
20     return image, class_num
21
22 # Pre-defined function. Function to Load the dataset
23 def load_dataset(filenames):
24     # read from TFRecords. For optimal performance, read from multiple
25     # TFRecord files at once and set the option experimental_determinis
26     # to allow order-altering optimizations.
27     option_no_order = tf.data.Options()
28     option_no_order.experimental_deterministic = False
29
30     dataset = tf.data.TFRecordDataset(filenames)
31     dataset = dataset.with_options(option_no_order)
32     dataset = dataset.map(read_tfrecord)
33     return dataset
34
35 # Creating function for checking the performance
36 def time_configs_new(parameters_rdd):
37     batch_size = parameters_rdd[0]
38     batch_num = parameters_rdd[1]
39     repetition = parameters_rdd[2]
40     filenames = tf.io.gfile.glob(GCS_OUTPUT + "*tfrec")
41     dset = load_dataset(filenames)
42
43     batch = dset.batch(batch_size)
44     sample_set = dset.take(batch_num)
45
46     tt = []
47     for rep in range(repetition):
48         start = time.time()
49         for picture in sample_set:
50             print('string', file=open("/dev/null", mode='w'))
51         end = time.time()
52         reading_speed = end - start
53         throughput = float((batch_size * batch_num) / (end - start))
54         datasetsize = batch_size * batch_num
55         tt.append([batch_size, batch_num, repetition, datasetsize, reading])
56
57     return tt
58
59     # Extra resources checked:
```





In [ ]:

```
1 # i) combine the previous cells to have the code to create a dataset a
2 # ii) get a Spark context and create the dataset and run timing test fo
3
4 batch_sizes = [6, 8, 10, 12]
5 batch_numbers = [6, 9, 12, 15]
6 repetitions = [1, 2, 3]
7
8 parameter_list = []
9 for batch_size in batch_sizes:
10     for batch in batch_numbers:
11         for r in repetitions:
12             parameter_list.append([batch_size,batch,r])
13
14 columns = ["batch_sizes", "batch_nums", "repetitions", "datasetsize", " "
15
16 # Get spark context for RDDs
17 sc = pyspark.SparkContext.getOrCreate()
18
19 # Parallelize data into RDD
20 tfr_files_rdd = sc.parallelize(parameter_list)
21
22 # Initialize spark session for df conversion
23 tfr_files_sparkSession = SparkSession(sc)
24
25 # Flattening RDD into List of List
26 tfr_files = tfr_files_rdd.flatMap(time_configs_new)
27
28 # Construct dataframe (df) ffor TFR files
29 tfr_files_df = tfr_files.toDF(columns)
30
31 # Load and decode dataset
32 def load_dataset_decoded():
33     dataset_filename = tf.data.Dataset.list_files(GCS_PATTERN)
34     datasetDecoded = dataset_filename.map(decode_jpeg_and_label)
35     datasetfn = datasetDecoded.map(resize_and_crop_image)
36     return datasetfn
37
38 # Function for image_files
39 def img_configs_new(parameters_rdd):
40     batch_size = parameters_rdd[0]
41     batch_num = parameters_rdd[1]
42     repetition = parameters_rdd[2]
43
44     dset = load_dataset_decoded()
45     batch = dset.batch(batch_size)
46     sample = dset.take(batch_num)
47     tt = []
48     for rep in range(repetition):
49         start = time.time()
50         for picture in sample:
51             print('string', file=open("/dev/null", mode='w'))
52         end = time.time()
53         reading_speed = end - start
54         throughput = float((batch_size * batch_num) / (end - start))
55         datasetsize = batch_size * batch_num
56         tt.append([batch_size, batch_num, repetition, datasetsize, reading_
57     return tt
58
59 # Get spark context for RDDs
60 sc = pyspark.SparkContext.getOrCreate()
61
```

```
62 # Parallelize data into RDD
63 image_files_rdd = sc.parallelize(parameter_list)
64
65 # Initialize spark session for df conversion
66 image_files_sparkSession = SparkSession(sc)
67
68 # Flattening RDD into List of List
69 image_files = image_files_rdd.flatMap(img_configs_new)
70
71 # Create df for image files
72 image_files_df = image_files.toDF(columns)
73
```

```
In [ ]: 1 # iii) transform the resulting RDD to the structure ( parameter_combina
2
3 tfr_array_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']),
4 int(x['batch_nums']),
5 int(x['repetitions']),
6 int(x['datasetsize']),
7 float(x['reading_speed']),
8 float(x['throughput'])))
9
10 print("TFRecords Array RDD")
11 tfr_array_rdd.take(6)
12
13 # Extra resources checked:
14 # https://data-flair.training/blogs/spark-rdd-operations-transformation
15 # https://sparkbyexamples.com/
16 # https://www.datacamp.com/tracks/big-data-with-pyspark
17 # https://pages.databricks.com/rs/094-YMS-629/images/LearningSpark2.0.p
18 # https://sparkbyexamples.com/pyspark/pyspark-convert-dataframe-to-rdd/
```

TFRecords Array RDD

Out[64]: [(6, 6, 1, 36, 0.2102503776550293, 171.22442490480285),  
(6, 6, 2, 36, 0.17751073837280273, 202.80463215917607),  
(6, 6, 2, 36, 0.18248200416564941, 197.279727196117),  
(6, 6, 3, 36, 0.1756117343902588, 204.9976906440537),  
(6, 6, 3, 36, 0.17471814155578613, 206.04614769500327),  
(6, 6, 3, 36, 0.20185422897338867, 178.34652354371053)]

```
In [ ]: 1 # RDD for image file parameter_combination
2
3 image_array_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_sizes'])
4
5
6
7
8
9
10 print("Image Array RDD")
11 image_array_rdd.take(6)
```

Image Array RDD

Out[65]: [(6, 6, 1, 36, 0.42624759674072266, 84.45795419205152),  
(6, 6, 2, 36, 0.4570634365081787, 78.76368382259739),  
(6, 6, 2, 36, 0.5044350624084473, 71.36696610286451),  
(6, 6, 3, 36, 0.5790822505950928, 62.167334541862864),  
(6, 6, 3, 36, 0.525456428527832, 68.5118652004334),  
(6, 6, 3, 36, 0.5065336227416992, 71.07129395506638)]

In [ ]:

```
1 # iv) Create an RDD with all results for each parameter as (parameter_v
2
3 # Create RDD for image files batch size and their corresponding speed
4 image_batchSizes_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['b
5 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect
6
7 # Create RDD for image files batch numbers and their corresponding spee
8 image_batchNums_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['ba
9 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect
10
11 # Create RDD for image files repetitions and their corresponding speeds
12 image_repetitions_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['re
13 image_repetitions_speed = image_repetitions_speed_rdd.collect() # coll
14
15 # Create RDD for image files data set size and their corresponding spee
16 image_dsSize_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['datas
17 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect dataset
18
19
20
21
22
23 # Create RDD for batch sizes and their corresponding speeds
24 tfr_batchSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_
25 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batch
26
27 # Create RDD for batch numbers and their corresponding speeds
28 tfr_batchNum_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_n
29 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batch
30
31 # Create RDD for repetitions and their corresponding speeds
32 tfr_repetitions_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['repe
33 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect
34
35 # Create RDD for dataset size (dsSize) and their corresponding speeds
36 tfr_dsSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['datasetsi
37 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect datasetsize
38
39
40 # Extra resources checked:
41 # https://spark.apache.org/docs/latest/rdd-programming-guide.html
42 # https://sparkbyexamples.com/pyspark-rdd/
43 # https://www.oreilly.com/library/view/learning-spark/9781449359034/ch0
44 # https://stackoverflow.com/questions/73014570/best-way-to-join-2-pair-
45 # https://towardsdatascience.com/the-ultimate-guide-to-functional-progr
46 # https://spark.apache.org/docs/latest/rdd-programming-guide.html#basic
47 # https://www.geeksforgeeks.org/pyspark-map-transformation/
48 # https://www.analyticsvidhya.com/blog/2016/10/using-pyspark-to-perform
```

In [ ]:

```
1 # v) create an RDD with the average reading speeds for each parameter v
2
3 # Create RDD for batch sizes and their corresponding avg speed
4 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x:
5 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect
6
7 # Create RDD for tfr batch numbers and their corresponding avg speed
8 tfr_batchNums_avgSpeed_rdd = tfr_batchNum_speed_rdd.mapValues(lambda x:
9 tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect
10
11 # Create RDD for tfr repetitions and their corresponding avg speed
12 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x:
13 tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # collect
14
15 # Create RDD for tfr dataset size and their corresponding avg speed
16 tfr_dsSize_avgSpeed_rdd = tfr_dsSize_speed_rdd.mapValues(lambda x: (x,
17 tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect data
18
19
20
21
22 # Create RDD for image files batch size and their corresponding avg spe
23 image_batchSizes_avgSpeed_rdd = image_batchSizes_speed_rdd.mapValues(lambda x:
24 image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() #
25
26 # Create RDD for image batch number and their corresponding avg speed
27 image_batchNums_avgSpeed_rdd = image_batchNums_speed_rdd.mapValues(lambda x:
28 image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # co
29
30 # Create RDD for image repetitions and their corresponding avg speed
31 image_repetitions_avgSpeed_rdd = image_repetitions_speed_rdd.mapValues(lambda x:
32 image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() #
33
34 # Create RDD for image dataset size and their corresponding avg speed
35 image_dsSize_avgSpeed_rdd = image_dsSize_speed_rdd.mapValues(lambda x:
36 image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect
37
38
39 # Some extra resources checked:
40 # https://spark.apache.org/docs/latest/api/python/reference/api/pyspark
41 #
```

In [ ]:

```
1 # vi) write the results to a pickle file in your bucket (2%)
2
3 # Save object to Local file using pickle, upload to GCS, and Log status
4 def save(object,bucket,filename):
5     with open(filename, mode='wb') as f:
6         pickle.dump(object,f)
7         print("Saving{} to {}".format(filename,bucket))
8         import subprocess
9         proc=subprocess.run(["gsutil","cp",filename,bucket],stderr=subprocess.PIPE)
10        print("gstuil returned: " + str(proc.returncode))
11        print(str(proc.stderr))
12
13 filename="2aVI_results.pkl"
14 with open(filename,mode='wb') as f:
15     pickle.dump(tfr_batchSize_speed,f)
16     pickle.dump(tfr_batchNum_speed,f)
17     pickle.dump(tfr_repetitions_speed,f)
18     pickle.dump(tfr_dsSize_speed,f)
19     pickle.dump(image_batchSizes_speed,f)
20     pickle.dump(image_batchNums_speed,f)
21     pickle.dump(image_repetitions_speed,f)
22     pickle.dump(image_dsSize_speed,f)
23     pickle.dump(tfr_batchSize_avgSpeed,f)
24     pickle.dump(tfr_batchNums_avgSpeed,f)
25     pickle.dump(tfr_repetitions_avgSpeed,f)
26     pickle.dump(tfr_dsSize_avgSpeed,f)
27     pickle.dump(image_batchSizes_avgSpeed,f)
28     pickle.dump(image_batchNums_avgSpeed,f)
29     pickle.dump(image_repetitions_avgSpeed,f)
30     pickle.dump(image_dsSize_avgSpeed,f)
31
32
33 print("Saving {} to {}".format(filename,BUCKET))
34 import subprocess
35 proc = subprocess.run(["gsutil", "cp",filename, BUCKET], stderr=subprocess.PIPE)
36 print("gstuil returned: " +str(proc.returncode))
37 print(str(proc.stderr))
38
39
40 # Extra resources consulted:
41 # https://www.geeksforgeeks.org/file-handling-python/
42 # https://docs.python.org/3/library/subprocess.html
43 # https://stackoverflow.com/questions/40982132/python-pickle-dump-wb-pa
```

```
Saving 2aVI_results.pkl to gs://big-data-cw-420116-storage
gstuil returned: 0
b'Copying file://2aVI_results.pkl [Content-Type=application/octet-stream]...\\n/ [0 files][ 0.0 B/ 10.6 KiB]
\\r/ [1 files][ 10.6 KiB/ 10.6 KiB]
\\r-\\r\\nOperation completed over 1 objects/10.6 KiB.
\\n'
```



In [ ]:

```
1 # vii) Write your code it into a file using the cell magic %%writefile
2
3 %%writefile spark_job.py
4 # Setting the working environment
5 import datetime
6 import math
7 import numpy as np
8 import os
9 import pyspark
10 import pickle
11 import pandas as pd
12 import random
13 import string
14 import sys
15 #import scipy as sp
16 #import scipy.stats
17 import time
18 import tensorflow as tf
19 print ("TensorFlow Version" == tf.__version__)
20
21 #from matplotlib import pyplot as plt
22 from pyspark.sql import Row
23 from pyspark.sql import SparkSession
24 from pyspark.sql import SQLContext
25
26 PROJECT = 'big-data-cw-420116'
27 BUCKET = 'gs://{}-storage'.format(PROJECT)
28 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input jpg
29 PARTITIONS = 16 # no of partitions we will use later
30 TARGET_SIZE = [192, 192] # target resolution for the images
31 CLASSES = [b'daisy', b'dandelion', b'roses', b'sunflowers', b'tulips']
32           # labels for the data
33 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for output
34 nb_images = len(tf.io.gfile.glob(GCS_PATTERN)) # number of images
35
36
37 # Pre-defined function. Function to parse the TFRecord
38 def read_tfrecord(example):
39     features = {
40         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = bytes
41         "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] means scalar
42     }
43     # decode the TFRecord
44     example = tf.io.parse_single_example(example, features)
45     image = tf.image.decode_jpeg(example['image'], channels=3)
46     image = tf.reshape(image, [*TARGET_SIZE, 3])
47     class_num = example['class']
48     return image, class_num
49
50 def decode_jpeg_and_label(filepath):
51     # extracts the image data and creates a class label, based on the folder name
52     bits = tf.io.read_file(filepath)
53     image = tf.image.decode_jpeg(bits)
54     # parse flower name from containing directory
55     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
56     label2 = label.values[-2]
57     return image, label2
58
59 def resize_and_crop_image(image, label):
60     # Resizes and cropd using "fill" algorithm:
61     # always make sure the resulting image is cut out from the source image
```

```

62      # so that it fills the TARGET_SIZE entirely with no black bars
63      # and a preserved aspect ratio.
64      w = tf.shape(image)[0]
65      h = tf.shape(image)[1]
66      tw = TARGET_SIZE[1]
67      th = TARGET_SIZE[0]
68      resize_crit = (w * th) / (h * tw)
69      image = tf.cond(resize_crit < 1,
70                      lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), #
71                      lambda: tf.image.resize(image, [w*th/h, h*th/h])) #
72
73      nw = tf.shape(image)[0]
74      nh = tf.shape(image)[1]
75      image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
76      return image, label
77
78 # Pre-defined function. Function to Load the dataset
79 def load_dataset(filenames):
80     # read from TFRecords. For optimal performance, read from multiple
81     # TFRecord files at once and set the option experimental_determinis-
82     # to allow order-altering optimizations.
83     option_no_order = tf.data.Options()
84     option_no_order.experimental_deterministic = False
85
86     dataset = tf.data.TFRecordDataset(filenames)
87     dataset = dataset.with_options(option_no_order)
88     dataset = dataset.map(read_tfrecord)
89     return dataset
90
91 # Creating function for checking the performance
92 def time_configs_new(parameters_rdd):
93     batch_size = parameters_rdd[0]
94     batch_num = parameters_rdd[1]
95     repetition = parameters_rdd[2]
96
97     filenames = tf.io.gfile.glob(GCS_OUTPUT + ".*.tfrec")
98     dset = load_dataset(filenames)
99
100    batch = dset.batch(batch_size)
101    sample_set = dset.take(batch_num)
102
103    tt = []
104    for rep in range(repetition):
105        start = time.time()
106        for picture in sample_set:
107            print('string', file=open("/dev/null", mode='w'))
108        end = time.time()
109        reading_speed = end - start
110        throughput = float((batch_size * batch_num) / (end - start))
111        datasetsize = batch_size * batch_num
112        tt.append([batch_size, batch_num, repetition, datasetsize, reading])
113
114    # i) combine the previous cells to have the code to create a dataset o
115    # ii) get a Spark context and create the dataset and run timing test fo
116
117
118    batch_sizes = [6, 8, 10, 12]
119    batch_numbers = [6, 9, 12, 15]
120    repetitions = [1, 2, 3]
121
122    parameter_list = []

```

```

123 for batch_size in batch_sizes:
124     for batch in batch_numbers:
125         for r in repetitions:
126             parameter_list.append([batch_size,batch,r])
127
128 columns = ["batch_sizes", "batch_nums", "repetitions", "datasetsize", ' '
129
130 # Get spark context for RDDs
131 sc = pyspark.SparkContext.getOrCreate()
132
133 # Parallelize data into RDD
134 tfr_files_rdd = sc.parallelize(parameter_list)
135
136 # Initialize spark session for df conversion
137 tfr_files_sparkSession = SparkSession(sc)
138
139 # Flattening RDD into List of List
140 tfr_files = tfr_files_rdd.flatMap(time_configs_new)
141
142 # Construct dataframe (df) for TFR files
143 tfr_files_df = tfr_files.toDF(columns)
144
145 # Load and decode dataset
146 def load_dataset_decoded():
147     dataset_filename = tf.data.Dataset.list_files(GCS_PATTERN)
148     datasetDecoded = dataset_filename.map(decode_jpeg_and_label)
149     datasetfn = datasetDecoded.map(resize_and_crop_image)
150     return datasetfn
151
152 # create function for image files
153 def img_configs_new(parameters_rdd):
154     dset = load_dataset_decoded()
155     tt = []
156
157     batch_size = parameters_rdd[0]
158     batch_num = parameters_rdd[1]
159     repetition = parameters_rdd[2]
160
161     batch = dset.batch(batch_size)
162     sample = dset.take(batch_num)
163
164     for rep in range(repetition):
165         start = time.time()
166         for picture in sample:
167             print('string', file=open("/dev/null", mode='w'))
168         end = time.time()
169         reading_speed = end - start
170         throughput = float((batch_size * batch_num) / (end - start))
171         datasetsize = batch_size * batch_num
172         tt.append([batch_size, batch_num, repetition, datasetsize, reading])
173
174     return tt
175
176 # Get spark context for RDDs
177 sc = pyspark.SparkContext.getOrCreate()
178
179 # Parallelize data into RDD
180 image_files_rdd = sc.parallelize(parameter_list)
181
182 # Initialize spark session for df conversion
183 image_files_sparkSession = SparkSession(sc)
184

```

```

184 # Flattening RDD into List of List
185 image_files = image_files_rdd.flatMap(img_configs_new)
186
187 # Create df for image files
188 image_files_df = image_files.toDF(columns)
189
190 # iii) transform the resulting RDD to the structure ( parameter_combination )
191 tfr_array_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']),
192                                         int(x['batch_nums']),
193                                         int(x['repetitions']),
194                                         int(x['datasetsize']),
195                                         float(x['reading_speed']),
196                                         float(x['throughput'])))
197 tfr_array_rdd.take(6)
198
199
200
201
202 # iv) create an RDD with all results for each parameter as (parameter_value, speed)
203
204 # Create RDD for batch sizes and their corresponding speeds
205 tfr_batchSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_size']),
206 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batch size speed
207
208 # Create RDD for batch numbers and their corresponding speeds
209 tfr_batchNum_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_num']),
210 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batch number speed
211
212 # Create RDD for repetitions and their corresponding speeds
213 tfr_repetitions_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['repetitions']),
214 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect repetitions speed
215
216 # Create RDD for dataset size (dsSize) and their corresponding speeds
217 tfr_dsSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['datasetsize']),
218 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect dataset size speed
219
220 # Create RDD for image files batch size and their corresponding speeds
221 image_batchSizes_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_size']),
222 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect batch size speed
223
224 # Create RDD for image files batch numbers and their corresponding speeds
225 image_batchNums_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_num']),
226 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect batch number speed
227
228 # Create RDD for image files repetitions and their corresponding speeds
229 image_repetitions_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['repetitions']),
230 image_repetitions_speed = image_repetitions_speed_rdd.collect() # collect repetitions speed
231
232 # Create RDD for image files data set size and their corresponding speeds
233 image_dsSize_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['datasetsize']),
234 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect dataset size speed
235
236
237
238
239 # v) create an RDD with the average reading speeds for each parameter value
240
241 # Create RDD for batch sizes and their corresponding avg speed
242 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x: sum(x)/len(x))
243 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect batch size avg speed
244

```

```

245 # Create RDD for tfr batch numbers and their corresponding avg speed
246 tfr_batchNums_avgSpeed_rdd = tfr_batchNum_speed_rdd.mapValues(lambda x:
247     tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect
248
249 # Create RDD for tfr repetitions and their corresponding avg speed
250 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x:
251     tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # col
252
253 # Create RDD for tfr dataset size and their corresponding avg speed
254 tfr_dsSize_avgSpeed_rdd = tfr_dsSize_speed_rdd.mapValues(lambda x: (x,
255     tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect data
256
257 # Create RDD for image files batch size and their corresponding avg speed
258 image_batchSizes_avgSpeed_rdd = image_batchSizes_speed_rdd.mapValues(lambda x:
259     image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() # co
260
261 # Create RDD for image batch number and their corresponding avg speed
262 image_batchNums_avgSpeed_rdd = image_batchNums_speed_rdd.mapValues(lambda x:
263     image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # co
264
265 # Create RDD for image repetitions and their corresponding avg speed
266 image_repetitions_avgSpeed_rdd = image_repetitions_speed_rdd.mapValues(lambda x:
267     image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() # co
268
269 # Create RDD for image dataset size and their corresponding avg speed
270 image_dsSize_avgSpeed_rdd = image_dsSize_speed_rdd.mapValues(lambda x: (x,
271     image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect
272
273
274
275
276 # vi) write the results to a pickle file in your bucket (2%)
277
278 #Save object to local file using pickle, upload to GCS, and Log status
279 def save(object,bucket,filename):
280     with open(filename, mode='wb') as f:
281         pickle.dump(object,f)
282         print("Saving{} to {}".format(filename,bucket))
283         import subprocess
284         proc=subprocess.run(["gsutil","cp",filename,bucket],stderr=subprocess.STDOUT)
285         print("gstutil returned: " + str(proc.returncode))
286         print(str(proc.stderr))
287
288 filename="2aVI_results.pkl"
289 with open(filename,mode='wb') as f:
290     pickle.dump(tfr_batchSize_speed,f)
291     pickle.dump(tfr_batchNum_speed,f)
292     pickle.dump(tfr_repetitions_speed,f)
293     pickle.dump(tfr_dsSize_speed,f)
294     pickle.dump(image_batchSizes_speed,f)
295     pickle.dump(image_batchNums_speed,f)
296     pickle.dump(image_repetitions_speed,f)
297     pickle.dump(image_dsSize_speed,f)
298     pickle.dump(tfr_batchSize_avgSpeed,f)
299     pickle.dump(tfr_batchNums_avgSpeed,f)
300     pickle.dump(tfr_repetitions_avgSpeed,f)
301     pickle.dump(tfr_dsSize_avgSpeed,f)
302     pickle.dump(image_batchSizes_avgSpeed,f)
303     pickle.dump(image_batchNums_avgSpeed,f)
304     pickle.dump(image_repetitions_avgSpeed,f)
305     pickle.dump(image_dsSize_avgSpeed,f)

```

```
306
307 print("Saving {} to {}".format(filename, BUCKET))
308 import subprocess
309 proc = subprocess.run(["gsutil", "cp", filename, BUCKET], stderr=subprocess.PIPE)
310 print("gstutil returned: " + str(proc.returncode))
311 print(str(proc.stderr))
```

Writing spark\_job.py

## 2b) Testing the code and collecting results (4%)

- i) First, test locally with %run .

It is useful to create a **new filename argument**, so that old results don't get overwritten.

You can for instance use `datetime.datetime.now().strftime("%y%m%d-%H%M")` to get a string with the current date and time and use that in the file name.

In [ ]:

```
1 ### CODING TASK
2 %run ./spark_job.py
```

```
False
Saving 2aVI_results.pkl to gs://big-data-cw-420116-storage
gstutil returned: 0
b'Copying file://2aVI_results.pkl [Content-Type=application/octet-stream]...
  [0 files][ 0.0 B/ 10.6 KiB]
  [1 files][ 10.6 KiB/ 10.6 KiB]
  Operation completed over 1 objects/10.6 KiB.
\n'
```

<Figure size 640x480 with 0 Axes>



In [ ]:

```
1 %%writefile spark_job2b.py
2 # Setting the working environment
3 import datetime
4 import math
5 import numpy as np
6 import os
7 import pyspark
8 import pickle
9 import pandas as pd
10 import random
11 import string
12 import sys
13 import time
14 import tensorflow as tf
15 print ("TensorFlow Version" == tf.__version__)
16
17 from pyspark.sql import Row
18 from pyspark.sql import SparkSession
19 from pyspark.sql import SQLContext
20
21 PROJECT = 'big-data-cw-420116'
22 BUCKET = 'gs://{}-storage'.format(PROJECT)
23 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input j
24 PARTITIONS = 16 # no of partitions we will use later
25 TARGET_SIZE = [192, 192] # target resolution for the images
26 CLASSES = [b'daisy', b'dandelion', b'roses', b'sunflowers', b'tulips']
27         # Labels for the data
28 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for
29 nb_images = len(tf.io.gfile.glob(GCS_PATTERN)) # number of images
30
31 # Pre-defined function. Function to parse the TFRecord
32 def read_tfrecord(example):
33     features = {
34         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = b
35         "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] me
36     }
37     # decode the TFRecord
38     example = tf.io.parse_single_example(example, features)
39     image = tf.image.decode_jpeg(example['image'], channels=3)
40     image = tf.reshape(image, [*TARGET_SIZE, 3])
41     class_num = example['class']
42     return image, class_num
43
44 def decode_jpeg_and_label(filepath):
45     # extracts the image data and creates a class label, based on the f
46     bits = tf.io.read_file(filepath)
47     image = tf.image.decode_jpeg(bits)
48     # parse flower name from containing directory
49     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
50     label2 = label.values[-2]
51     return image, label2
52
53 def resize_and_crop_image(image, label):
54     # Resizes and cropd using "fill" algorithm:
55     # always make sure the resulting image is cut out from the source i
56     # so that it fills the TARGET_SIZE entirely with no black bars
57     # and a preserved aspect ratio.
58     w = tf.shape(image)[0]
59     h = tf.shape(image)[1]
60     tw = TARGET_SIZE[1]
61     th = TARGET_SIZE[0]
```

```

62     resize_crit = (w * th) / (h * tw)
63     image = tf.cond(resize_crit < 1,
64                      lambda: tf.image.resize(image, [w*tw/w, h*tw/w]),
65                      lambda: tf.image.resize(image, [w*th/h, h*th/h]))
66
67     nw = tf.shape(image)[0]
68     nh = tf.shape(image)[1]
69     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
70     return image, label
71
72 # Pre-defined function. Function to load the dataset
73 def load_dataset(filenames):
74     # read from TFRecords. For optimal performance, read from multiple
75     # TFRecord files at once and set the option experimental_deterministic
76     # to allow order-altering optimizations.
77     option_no_order = tf.data.Options()
78     option_no_order.experimental_deterministic = False
79
80     dataset = tf.data.TFRecordDataset(filenames)
81     dataset = dataset.with_options(option_no_order)
82     dataset = dataset.map(read_tfrecord)
83     return dataset
84
85 # Creating function for checking the performance
86 def time_configs_new(parameters_rdd):
87     batch_size = parameters_rdd[0]
88     batch_num = parameters_rdd[1]
89     repetition = parameters_rdd[2]
90
91     filenames = tf.io.gfile.glob(GCS_OUTPUT + "*tfrec")
92     dset = load_dataset(filenames)
93
94     batch = dset.batch(batch_size)
95     sample_set = dset.take(batch_num)
96
97     tt = []
98     for rep in range(repetition):
99         start = time.time()
100        for picture in sample_set:
101            print('string', file=open("/dev/null", mode='w'))
102        end = time.time()
103        reading_speed = end - start
104        throughput = float((batch_size * batch_num) / (end - start))
105        datasetsize = batch_size * batch_num
106        tt.append([batch_size, batch_num, repetition, datasetsize, reading])
107    return tt
108
109 # i) combine the previous cells to have the code to create a dataset and run timing test
110 # ii) get a Spark context and create the dataset and run timing test for each combination
111
112 batch_sizes = [6, 8, 10, 12]
113 batch_numbers = [6, 9, 12, 15]
114 repetitions = [1, 2, 3]
115
116 parameter_list = []
117 for batch_size in batch_sizes:
118     for batch in batch_numbers:
119         for r in repetitions:
120             parameter_list.append([batch_size, batch, r])
121
122 columns = ["batch_sizes", "batch_nums", "repetitions", "datasetsize", "throughput"]

```

```

123
124 # Get spark context for RDDs
125 sc = pyspark.SparkContext.getOrCreate()
126
127 # Parallelize data into RDD
128 tfr_files_rdd = sc.parallelize(parameter_list)
129
130 # Initialize spark session for df conversion
131 tfr_files_sparkSession = SparkSession(sc)
132
133 # Flattening RDD into list of list
134 tfr_files = tfr_files_rdd.flatMap(time_configs_new)
135
136 # Construct dataframe (df) for TFR files
137 tfr_files_df = tfr_files.toDF(columns)
138
139 # Load and decode dataset
140 def load_dataset_decoded():
141     dataset_filename = tf.data.Dataset.list_files(GCS_PATTERN)
142     datasetDecoded = dataset_filename.map(decode_jpeg_and_label)
143     datasetfn = datasetDecoded.map(resize_and_crop_image)
144     return datasetfn
145
146 # create function for image files
147 def img_configs_new(parameters_rdd):
148     dset = load_dataset_decoded()
149     tt = []
150
151     batch_size = parameters_rdd[0]
152     batch_num = parameters_rdd[1]
153     repetition = parameters_rdd[2]
154
155     batch = dset.batch(batch_size)
156     sample = dset.take(batch_num)
157
158     for rep in range(repetition):
159         start = time.time()
160         for picture in sample:
161             print('string', file=open("/dev/null", mode='w'))
162         end = time.time()
163         reading_speed = end - start
164         throughput = float((batch_size * batch_num) / (end - start))
165         datasetsize = batch_size * batch_num
166         tt.append([batch_size, batch_num, repetition, datasetsize, reading])
167     return tt
168
169 # Get spark context for RDDs
170 sc = pyspark.SparkContext.getOrCreate()
171
172 # Parallelize data into RDD
173 image_files_rdd = sc.parallelize(parameter_list)
174
175 # Initialize spark session for df conversion
176 image_files_sparkSession = SparkSession(sc)
177
178 # Flattening RDD into list of list
179 image_files = image_files_rdd.flatMap(img_configs_new)
180
181 # Create df for image files
182 image_files_df = image_files.toDF(columns)
183

```

```

184 # iii) transform the resulting RDD to the structure ( parameter_combination, batch_size, dataset_size, repetition, ds_size, image_file_size, avg_speed )
185 tfr_array_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']), int(x['datasets']), int(x['repetitions']), int(x['dsSize']), int(x['image_file_size']), float(x['avgSpeed'])))
186 tfr_array_rdd.take(6)
187
188
189
190 # iv) create an RDD with all results for each parameter as (parameter_value, speed)
191 # Create RDD for batch sizes and their corresponding speeds
192 tfr_batchSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']), float(x['avgSpeed'])))
193 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batchSizeSpeed
194
195 # Create RDD for batch numbers and their corresponding speeds
196 tfr_batchNum_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_nums']), float(x['avgSpeed'])))
197 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batchNumSpeed
198
199 # Create RDD for repetitions and their corresponding speeds
200 tfr_repetitions_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['repetitions']), float(x['avgSpeed'])))
201 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect repetitionsSpeed
202
203 # Create RDD for dataset size (dsSize) and their corresponding speeds
204 tfr_dsSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['dataset_size']), float(x['avgSpeed'])))
205 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect datasetSizeSpeed
206
207 # Create RDD for image files batch size and their corresponding speeds
208 image_batchSizes_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_sizes']), float(x['avgSpeed'])))
209 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect batchSizesSpeed
210
211 # Create RDD for image files batch numbers and their corresponding speeds
212 image_batchNums_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_nums']), float(x['avgSpeed'])))
213 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect batchNumsSpeed
214
215 # Create RDD for image files repetitions and their corresponding speeds
216 image_repetitions_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['repetitions']), float(x['avgSpeed'])))
217 image_repetitions_speed = image_repetitions_speed_rdd.collect() # collect repetitionsSpeed
218
219 # Create RDD for image files data set size and their corresponding speeds
220 image_dsSize_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['dataset_size']), float(x['avgSpeed'])))
221 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect datasetSizeSpeed
222
223
224 # v) create an RDD with the average reading speeds for each parameter value
225
226 # Create RDD for batch sizes and their corresponding avg speed
227 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x: float(x[0] * x[1]))
228 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect batchSizeAvgSpeed
229
230 # Create RDD for tfr batch numbers and their corresponding avg speed
231 tfr_batchNums_avgSpeed_rdd = tfr_batchNum_speed_rdd.mapValues(lambda x: float(x[0] * x[1]))
232 tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect batchNumsAvgSpeed
233
234 # Create RDD for tfr repetitions and their corresponding avg speed
235 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x: float(x[0] * x[1]))
236 tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # collect repetitionsAvgSpeed
237
238 # Create RDD for tfr dataset size and their corresponding avg speed
239 tfr_dsSize_avgSpeed_rdd = tfr_dsSize_speed_rdd.mapValues(lambda x: (x[0], float(x[1] * x[0])))
240 tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect datasetSizeAvgSpeed
241
242 # Create RDD for image files batch size and their corresponding avg speed
243 image_batchSizes_avgSpeed_rdd = image_batchSizes_speed_rdd.mapValues(lambda x: float(x[0] * x[1]))
244 image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() # collect batchSizesAvgSpeed

```

```

245
246 # Create RDD for image batch number and their corresponding avg speed
247 image_batchNums_avgSpeed_rdd = image_batchNums_speed_rdd.mapValues(lambda x: x[0])
248 image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # collect
249
250 # Create RDD for image repetitions and their corresponding avg speed
251 image_repetitions_avgSpeed_rdd = image_repetitions_speed_rdd.mapValues(lambda x: x[0])
252 image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() # collect
253
254 # Create RDD for image dataset size and their corresponding avg speed
255 image_dsSize_avgSpeed_rdd = image_dsSize_speed_rdd.mapValues(lambda x: x[0])
256 image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect
257
258
259 # vi) write the results to a pickle file in your bucket (2%)
260
261 #Save object to local file using pickle, upload to GCS, and Log status
262 def save(object,bucket,filename):
263     with open(filename, mode='wb') as f:
264         pickle.dump(object,f)
265         print("Saving{} to {}".format(filename,bucket))
266         import subprocess
267         proc=subprocess.run(["gsutil","cp",filename,bucket],stderr=subprocess.PIPE)
268         print("gstutil returned: " + str(proc.returncode))
269         print(str(proc.stderr))
270
271 filename="2b_results.pkl"
272 with open(filename,mode='wb') as f:
273     pickle.dump(tfr_batchSize_speed,f)
274     pickle.dump(tfr_batchNum_speed,f)
275     pickle.dump(tfr_repetitions_speed,f)
276     pickle.dump(tfr_dsSize_speed,f)
277     pickle.dump(image_batchSizes_speed,f)
278     pickle.dump(image_batchNums_speed,f)
279     pickle.dump(image_repetitions_speed,f)
280     pickle.dump(image_dsSize_speed,f)
281     pickle.dump(tfr_batchSize_avgSpeed,f)
282     pickle.dump(tfr_batchNums_avgSpeed,f)
283     pickle.dump(tfr_repetitions_avgSpeed,f)
284     pickle.dump(tfr_dsSize_avgSpeed,f)
285     pickle.dump(image_batchSizes_avgSpeed,f)
286     pickle.dump(image_batchNums_avgSpeed,f)
287     pickle.dump(image_repetitions_avgSpeed,f)
288     pickle.dump(image_dsSize_avgSpeed,f)
289
290 print("Saving {} to {}".format(filename,BUCKET))
291 import subprocess
292 proc = subprocess.run(["gsutil", "cp",filename, BUCKET], stderr=subprocess.PIPE)
293 print("gstutil returned: " +str(proc.returncode))
294 print(str(proc.stderr))

```

Writing spark\_job2b.py

## ii) Cloud

If you have a cluster running, you can run the speed test job in the cloud.

While you run this job, switch to the Dataproc web page and take **screenshots of the CPU and network load over time**. They are displayed with some delay so you may need to wait.

```
In [ ]: 1 #Was getting error to create the cluster so I deleted it to create it a
         2 !gcloud dataproc clusters delete big-data-cw-420116-cluster
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/0338bdcb-b7ec-3f6b-90a6-270c2f41b32b].

Deleted [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>].

```
In [ ]: 1 ##### CODING TASK #####
         2 # Set up a cluster with a single machine using the maximal SSD size (10
         3
         4 import time
         5 start_time = time.time()
         6 #REGION = europe-west3
         7 !gcloud dataproc clusters create $CLUSTER \
             --bucket $PROJECT-storage \
             --image-version 1.4-ubuntu18 \
             --master-machine-type n1-standard-8 \
             --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
             --num-workers 0\
             --initialization-actions gs://goog-dataproc-initialization-actions-
             --metadata PIP_PACKAGES=tensorflow==2.4.0
         15 %time
         16
         17 end_time = time.time()
         18 execution_time = end_time - start_time
         19 print(f"Execution time: {execution_time} seconds")
         20
         21
```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/11a64304-c8fa-37e8-84e6-50a16bf977a2].

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-c].

CPU times: user 4 µs, sys: 1 µs, total: 5 µs

Wall time: 8.34 µs

Execution time: 153.5760896205902 seconds

In [ ]:

```
1 # get information of the single machine cluster
2 !gcloud dataproc clusters describe $CLUSTER
```

```
clusterName: big-data-cw-420116-cluster
clusterUuid: 35f0437b-4971-4385-b85a-3808433b9719
config:
  configBucket: big-data-cw-420116-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
    metadata:
      PIP_PACKAGES: tensorflow==2.4.0
      networkUri: https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default (https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default)
    serviceAccountScopes:
      - https://www.googleapis.com/auth/bigquery (https://www.googleapis.com/auth/bigquery)
      - https://www.googleapis.com/auth/bigtable.admin.table (https://www.googleapis.com/auth/bigtable.admin.table)
      - https://www.googleapis.com/auth/bigtable.data (https://www.googleapis.com/auth/bigtable.data)
      - https://www.googleapis.com/auth/cloud.useraccounts.readonly (https://www.googleapis.com/auth/cloud.useraccounts.readonly)
      - https://www.googleapis.com/auth/devstorage.full\_control (https://www.googleapis.com/auth/devstorage.full\_control)
      - https://www.googleapis.com/auth/devstorage.read\_write (https://www.googleapis.com/auth/devstorage.read\_write)
      - https://www.googleapis.com/auth/logging.write (https://www.googleapis.com/auth/logging.write)
    zoneUri: https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/zones/us-central1-c (https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/zones/us-central1-c)
  initializationActions:
    - executableFile: gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh
      executionTimeout: 600s
  masterConfig:
    diskConfig:
      bootDiskSizeGb: 100
      bootDiskType: pd-ssd
      imageUri: https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-4-ubuntu18-20220125-170200-rc01 (https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-4-ubuntu18-20220125-170200-rc01)
    instanceNames:
      - big-data-cw-420116-cluster-m
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/zones/us-central1-c/machineTypes/n1-standard-8 (https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/zones/us-central1-c/machineTypes/n1-standard-8)
    minCpuPlatform: AUTOMATIC
    numInstances: 1
    preemptibility: NON_PREEMPTIBLE
  softwareConfig:
    imageVersion: 1.4.80-ubuntu18
    properties:
      capacity-scheduler:yarn.scheduler.capacity.root.default.ordering-policy: fair
      core:fs.gs.block.size: '134217728'
      core:fs.gs.metadata.cache.enable: 'false'
      core:hadoop.ssl.enabled.protocols: TLSv1,TLSv1.1,TLSv1.2
      dataproc:dataproc.allow.zero.workers: 'true'
      distcp:mapreduce.map.java.opts: -Xmx768m
```

```
distcp:mapreduce.map.memory.mb: '1024'
distcp:mapreduce.reduce.java.opts: -Xmx768m
distcp:mapreduce.reduce.memory.mb: '1024'
hdfs:dfs.datanode.address: 0.0.0.0:9866
hdfs:dfs.datanode.http.address: 0.0.0.0:9864
hdfs:dfs.datanode.https.address: 0.0.0.0:9865
hdfs:dfs.datanode.ipc.address: 0.0.0.0:9867
hdfs:dfs.namenode.handler.count: '20'
hdfs:dfs.namenode.http-address: 0.0.0.0:9870
hdfs:dfs.namenode.https-address: 0.0.0.0:9871
hdfs:dfs.namenode.lifeline.rpc-address: big-data-cw-420116-cluster-
m:8050
hdfs:dfs.namenode.secondary.http-address: 0.0.0.0:9868
hdfs:dfs.namenode.secondary.https-address: 0.0.0.0:9869
hdfs:dfs.namenode.service.handler.count: '10'
hdfs:dfs.namenode.servicerpc-address: big-data-cw-420116-cluster-m:8
051
mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '4000'
mapred:mapreduce.job.maps: '21'
mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
mapred:mapreduce.job.reduces: '7'
mapred:mapreduce.map.cpu.vcores: '1'
mapred:mapreduce.map.java.opts: -Xmx2457m
mapred:mapreduce.map.memory.mb: '3072'
mapred:mapreduce.reduce.cpu.vcores: '1'
mapred:mapreduce.reduce.java.opts: -Xmx2457m
mapred:mapreduce.reduce.memory.mb: '3072'
mapred:mapreduce.task.io.sort.mb: '256'
mapred:yarn.app.mapreduce.am.command-opts: -Xmx2457m
mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
mapred:yarn.app.mapreduce.am.resource.mb: '3072'
spark-env:SPARK_DAEMON_MEMORY: 4000m
spark:spark.driver.maxResultSize: 3840m
spark:spark.driver.memory: 7680m
spark:spark.executor.cores: '4'
spark:spark.executor.instances: '2'
spark:spark.executor.memory: 11171m
spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
spark:spark.extraListeners: com.google.cloud.spark.performance.Datap
rocMetricsListener
spark:spark.scheduler.mode: FAIR
spark:spark.sql.cbo.enabled: 'true'
spark:spark.yarn.am.memory: 640m
yarn-env:YARN_NODEMANAGER_HEAPSIZE: '4000'
yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '4000'
yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '4000'
yarn:yarn.nodemanager.resource.cpu-vcores: '8'
yarn:yarn.nodemanager.resource.memory-mb: '24576'
yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-
secs: '86400'
yarn:yarn.scheduler.maximum-allocation-mb: '24576'
yarn:yarn.scheduler.minimum-allocation-mb: '1024'
tempBucket: dataproc-temp-us-central1-779145334445-vrtyf48i
labels:
goog-dataproc-autozone: enabled
goog-dataproc-cluster-name: big-data-cw-420116-cluster
goog-dataproc-cluster-uuid: 35f0437b-4971-4385-b85a-3808433b9719
goog-dataproc-location: us-central1
projectId: big-data-cw-420116
status:
state: RUNNING
```

```
stateStartTime: '2024-05-04T11:40:55.782481Z'  
statusHistory:  
- state: CREATING  
  stateStartTime: '2024-05-04T11:38:26.671779Z'
```

```
In [ ]: 1 # submit the spark job  
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \spark_job2b.py
```

Job [cac7e1460f9c477cb944f548c1cdb5a5] submitted.  
Waiting for job output...  
2024-05-04 11:41:04.458709: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dle error: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 11:41:04.458757: I tensorflow/stream\_executor/cuda/cudart\_stub.cc:29] Ignore above cudart dleerror if you do not have a GPU set up on your machine.  
False  
24/05/04 11:41:07 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 11:41:07 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 11:41:07 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 11:41:07 INFO org.spark\_project.jetty.util.log: Logging initialized @6112ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 11:41:07 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_312-b07  
24/05/04 11:41:07 INFO org.spark\_project.jetty.server.Server: Started @6221ms  
24/05/04 11:41:07 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@69ddacff{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
24/05/04 11:41:07 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools in fairscheduler.xml or set spark.scheduler.allocation.file to a file that contains the configuration.  
24/05/04 11:41:08 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.0.14:8032  
24/05/04 11:41:08 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.0.14:10200  
24/05/04 11:41:11 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714822762403\_0001  
Saving 2b\_results.pkl to gs://big-data-cw-420116-storage  
gstui returned: 0  
b'Copying file://2b\_results.pkl [Content-Type=application/octet-stream]...\n[ 0 files][ 0.0 B/ 11.3 KiB]\n[r/ [ 1 files][ 11.3 KiB/ 11.3 KiB]\n[r\\nOperation completed over 1 objects/11.3 KiB.\n\\n'  
24/05/04 11:49:05 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@69ddacff{HTTP/1.1, (http/1.1)}{0.0.0.0:4040}  
Job [cac7e1460f9c477cb944f548c1cdb5a5] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35f0437b-4971-4385-b85a-3808433b9719/jobs/cac7e1460f9c477cb944f548c1cdb5a5/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35f0437b-4971-4385-b85a-3808433b9719/jobs/cac7e1460f9c477cb944f548c1cdb5a5/driveroutput  
jobUuid: 5ca9651a-fcc5-30b5-b77a-b0fcfa588db51  
placement:  
    clusterName: big-data-cw-420116-cluster  
    clusterUuid: 35f0437b-4971-4385-b85a-3808433b9719  
pysparkJob:  
    mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35f0437b-4971-4385-b85a-3808433b9719/jobs/cac7e1460f9c477cb944f548c1cdb5a5/staging/spark\_job2b.py

```
reference:  
  jobId: cac7e1460f9c477cb944f548c1cdb5a5  
  projectId: big-data-cw-420116  
status:  
  state: DONE  
  stateStartTime: '2024-05-04T11:49:09.351993Z'  
statusHistory:  
- state: PENDING  
  stateStartTime: '2024-05-04T11:40:59.944655Z'  
- state: SETUP_DONE  
  stateStartTime: '2024-05-04T11:40:59.975410Z'  
- details: Agent reported job success  
  state: RUNNING  
  stateStartTime: '2024-05-04T11:41:00.252670Z'  
yarnApplications:  
- name: spark_job2b.py  
  progress: 1.0  
  state: FINISHED  
  trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714822762403\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714822762403\_0001/)
```

## 2c) Improve efficiency (6%)

If you implemented a straightforward version of 2a), you will **probably have an inefficiency** in your code.

Because we are reading multiple times from an RDD to read the values for the different parameters and their averages, caching existing results is important. Explain **where in the process caching can help**, and **add a call to `RDD.cache()`** to your code, if you haven't yet. Measure the effect of using caching or not using it.

Make the **suitable change** in the code you have written above and mark them up in comments as `### TASK 2c ###`.

Explain in your report what the **reasons for this change** are and **demonstrate and interpret its effect**



In [ ]:

```
1 ### CODING TASK ###
2 %writefile spark_job2c.py
3 # Setting the working environment
4 import datetime
5 import math
6 import numpy as np
7 import os
8 import pyspark
9 import pickle
10 import pandas as pd
11 import random
12 import string
13 import sys
14 #import scipy as sp
15 #import scipy.stats
16 import time
17 import tensorflow as tf
18 print ("TensorFlow Version" == tf.__version__)
19
20 #from matplotlib import pyplot as plt
21 from pyspark.sql import Row
22 from pyspark.sql import SparkSession
23 from pyspark.sql import SQLContext
24
25 PROJECT = 'big-data-cw-420116'
26 BUCKET = 'gs://{}-storage'.format(PROJECT)
27 GCS_PATTERN = 'gs://flowers-public/*/*.jpg' # glob pattern for input j
28 PARTITIONS = 16 # no of partitions we will use later
29 TARGET_SIZE = [192, 192] # target resolution for the images
30 CLASSES = [b'daisy', b'dandelion', b'roses', b'sunflowers', b'tulips']
31         # labels for the data
32 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for
33 nb_images = len(tf.io.gfile.glob(GCS_PATTERN)) # number of images
34
35 # Pre-defined function. Function to parse the TFRecord
36 def read_tfrecord(example):
37     features = {
38         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = b
39         "class": tf.io.FixedLenFeature([], tf.int64) , # shape [] me
40     }
41     # decode the TFRecord
42     example = tf.io.parse_single_example(example, features)
43     image = tf.image.decode_jpeg(example['image'], channels=3)
44     image = tf.reshape(image, [*TARGET_SIZE, 3])
45     class_num = example['class']
46     return image, class_num
47
48 def decode_jpeg_and_label(filepath):
49     # extracts the image data and creates a class label, based on the f
50     bits = tf.io.read_file(filepath)
51     image = tf.image.decode_jpeg(bits)
52     # parse flower name from containing directory
53     label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
54     label2 = label.values[-2]
55     return image, label2
56
57 def resize_and_crop_image(image, label):
58     # Resizes and cropd using "fill" algorithm:
59     # always make sure the resulting image is cut out from the source i
60     # so that it fills the TARGET_SIZE entirely with no black bars
61     # and a preserved aspect ratio.
```

```

62     w = tf.shape(image)[0]
63     h = tf.shape(image)[1]
64     tw = TARGET_SIZE[1]
65     th = TARGET_SIZE[0]
66     resize_crit = (w * th) / (h * tw)
67     image = tf.cond(resize_crit < 1,
68                       lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), #
69                       lambda: tf.image.resize(image, [w*th/h, h*th/h]) #
70                     )
71     nw = tf.shape(image)[0]
72     nh = tf.shape(image)[1]
73     image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh -
74     return image, label
75
76 # Pre-defined function. Function to load the dataset
77 def load_dataset(filenames):
78     # read from TFRecords. For optimal performance, read from multiple
79     # TFRecord files at once and set the option experimental_deterministic
80     # to allow order-altering optimizations.
81     option_no_order = tf.data.Options()
82     option_no_order.experimental_deterministic = False
83
84     dataset = tf.data.TFRecordDataset(filenames)
85     dataset = dataset.with_options(option_no_order)
86     dataset = dataset.map(read_tfrecord)
87     return dataset
88
89 # Creating function for checking the performance
90 def time_configs_new(parameters_rdd):
91     batch_size = parameters_rdd[0]
92     batch_num = parameters_rdd[1]
93     repetition = parameters_rdd[2]
94
95     filenames = tf.io.gfile.glob(GCS_OUTPUT + ".*.tfrec")
96     dset = load_dataset(filenames)
97
98     batch = dset.batch(batch_size)
99     sample_set = dset.take(batch_num)
100
101    tt = []
102    for rep in range(repetition):
103        start = time.time()
104        for picture in sample_set:
105            print('string', file=open("/dev/null", mode='w'))
106        end = time.time()
107        reading_speed = end - start
108        throughput = float((batch_size * batch_num) / (end - start))
109        datasetsize = batch_size * batch_num
110        tt.append([batch_size, batch_num, repetition, datasetsize, reading])
111    return tt
112
113 # i) combine the previous cells to have the code to create a dataset and run timing test for all combinations
114 # ii) get a Spark context and create the dataset and run timing test for all combinations
115
116 batch_sizes = [6, 8, 10, 12]
117 batch_numbers = [6, 9, 12, 15]
118 repetitions = [1, 2, 3]
119
120 parameter_list = []
121 for batch_size in batch_sizes:
122     for batch in batch_numbers:

```

```

123     for r in repetitions:
124         parameter_list.append([batch_size,batch,r])
125
126     columns = ["batch_sizes", "batch_nums", "repetitions", "datasetsize", ' '
127
128     # Get spark context for RDDs
129     sc = pyspark.SparkContext.getOrCreate()
130
131     # Parallelize data into RDD
132     tfr_files_rdd = sc.parallelize(parameter_list)
133
134     # Initialize spark session for df conversion
135     tfr_files_sparkSession = SparkSession(sc)
136
137     # Flattening RDD into list of list
138     tfr_files = tfr_files_rdd.flatMap(time_configs_new)
139
140     #### TASK 2c ####
141     tfr_files.cache() # Cache TFRecord files in RDD
142
143     # Construct dataframe (df) for TFR files
144     tfr_files_df = tfr_files.toDF(columns)
145
146     # Load and decode dataset
147     def load_dataset_decoded():
148         dataset_filename = tf.data.Dataset.list_files(GCS_PATTERN)
149         datasetDecoded = dataset_filename.map(decode_jpeg_and_label)
150         datasetfn = datasetDecoded.map(resize_and_crop_image)
151         return datasetfn
152
153     # create function for image files
154     def img_configs_new(parameters_rdd):
155         dset = load_dataset_decoded()
156         tt = []
157
158         batch_size = parameters_rdd[0]
159         batch_num = parameters_rdd[1]
160         repetition = parameters_rdd[2]
161
162         batch = dset.batch(batch_size)
163         sample = dset.take(batch_num)
164
165         for rep in range(repetition):
166             start = time.time()
167             for picture in sample:
168                 print('string', file=open("/dev/null", mode='w'))
169             end = time.time()
170             reading_speed = end - start
171             throughput = float((batch_size * batch_num) / (end - start))
172             datasetsize = batch_size * batch_num
173             tt.append([batch_size, batch_num, repetition, datasetsize, reading])
174
175     return tt
176
177     # Get spark context for RDDs
178     sc = pyspark.SparkContext.getOrCreate()
179
180     # Parallelize data into RDD
181     image_files_rdd = sc.parallelize(parameter_list)
182
183     # Initialize spark session for df conversion
184     image_files_sparkSession = SparkSession(sc)

```

```

184
185 # Flattening RDD into list of list
186 image_files = image_files_rdd.flatMap(img_configs_new)
187
188 ##### TASK 2c #####
189 image_files.cache() # Cache image files in memory for quick access
190
191 # Create df for image files
192 image_files_df = image_files.toDF(columns)
193
194
195 # iii) transform the resulting RDD to the structure ( parameter_combination )
196 tfr_array_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']), int(x['datasets'])))
197 tfr_array_rdd.take(6)
198
199
200
201 # iv) create an RDD with all results for each parameter as (parameter_value, speed)
202
203 # Create RDD for batch sizes and their corresponding speeds
204 tfr_batchSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_size']), int(x['speed'])))
205 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batch size speed
206
207 # Create RDD for batch numbers and their corresponding speeds
208 tfr_batchNum_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_number']), int(x['speed'])))
209 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batch number speed
210
211 # Create RDD for repetitions and their corresponding speeds
212 tfr_repetitions_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['repetitions']), int(x['speed'])))
213 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect repetitions speed
214
215 # Create RDD for dataset size (dsSize) and their corresponding speeds
216 tfr_dsSize_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['dataset_size']), int(x['speed'])))
217 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect dataset size speed
218
219 # Create RDD for image files batch size and their corresponding speeds
220 image_batchSizes_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_size']), int(x['speed'])))
221 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect batch size speed
222
223 # Create RDD for image files batch numbers and their corresponding speeds
224 image_batchNums_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_number']), int(x['speed'])))
225 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect batch number speed
226
227 # Create RDD for image files repetitions and their corresponding speeds
228 image_repetitions_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['repetitions']), int(x['speed'])))
229 image_repetitions_speed = image_repetitions_speed_rdd.collect() # collect repetitions speed
230
231 # Create RDD for image files data set size and their corresponding speeds
232 image_dsSize_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['dataset_size']), int(x['speed'])))
233 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect dataset size speed
234
235
236 # v) create an RDD with the average reading speeds for each parameter value
237
238 # Create RDD for batch sizes and their corresponding avg speed
239 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x: sum(x)/len(x))
240 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect batch size avg speed
241
242 # Create RDD for tfr batch numbers and their corresponding avg speed
243 tfr_batchNums_avgSpeed_rdd = tfr_batchNum_speed_rdd.mapValues(lambda x: sum(x)/len(x))
244 tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect batch number avg speed

```

```

245
246 # Create RDD for tfr repetitions and their corresponding avg speed
247 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x: (x[0], x[1] / len(x[1])))
248 tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # collect data
249
250 # Create RDD for tfr dataset size and their corresponding avg speed
251 tfr_dsSize_avgSpeed_rdd = tfr_dsSize_speed_rdd.mapValues(lambda x: (x[0], x[1] / len(x[1])))
252 tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect data
253
254 # Create RDD for image files batch size and their corresponding avg speed
255 image_batchSizes_avgSpeed_rdd = image_batchSizes_speed_rdd.mapValues(lambda x: (x[0], x[1] / len(x[1])))
256 image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() # collect data
257
258 # Create RDD for image batch number and their corresponding avg speed
259 image_batchNums_avgSpeed_rdd = image_batchNums_speed_rdd.mapValues(lambda x: (x[0], x[1] / len(x[1])))
260 image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # collect data
261
262 # Create RDD for image repetitions and their corresponding avg speed
263 image_repetitions_avgSpeed_rdd = image_repetitions_speed_rdd.mapValues(lambda x: (x[0], x[1] / len(x[1])))
264 image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() # collect data
265
266 # Create RDD for image dataset size and their corresponding avg speed
267 image_dsSize_avgSpeed_rdd = image_dsSize_speed_rdd.mapValues(lambda x: (x[0], x[1] / len(x[1])))
268 image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect data
269
270
271 # vi) write the results to a pickle file in your bucket (1%)
272
273 #Save object to local file using pickle, upload to GCS, and Log status
274 def save(object,bucket,filename):
275     with open(filename, mode='wb') as f:
276         pickle.dump(object,f)
277         print("Saving{} to {}".format(filename,bucket))
278         import subprocess
279         proc=subprocess.run(["gsutil","cp",filename,bucket],stderr=subprocess.PIPE)
280         print("gstuil returned: " + str(proc.returncode))
281         print(str(proc.stderr))
282
283 filename="2c_results.pkl"
284 with open(filename,mode='wb') as f:
285     pickle.dump(tfr_batchSize_speed,f)
286     pickle.dump(tfr_batchNum_speed,f)
287     pickle.dump(tfr_repetitions_speed,f)
288     pickle.dump(tfr_dsSize_speed,f)
289     pickle.dump(image_batchSizes_speed,f)
290     pickle.dump(image_batchNums_speed,f)
291     pickle.dump(image_repetitions_speed,f)
292     pickle.dump(image_dsSize_speed,f)
293     pickle.dump(tfr_batchSize_avgSpeed,f)
294     pickle.dump(tfr_batchNums_avgSpeed,f)
295     pickle.dump(tfr_repetitions_avgSpeed,f)
296     pickle.dump(tfr_dsSize_avgSpeed,f)
297     pickle.dump(image_batchSizes_avgSpeed,f)
298     pickle.dump(image_batchNums_avgSpeed,f)
299     pickle.dump(image_repetitions_avgSpeed,f)
300     pickle.dump(image_dsSize_avgSpeed,f)
301
302 print("Saving {} to {}".format(filename,BUCKET))
303 import subprocess
304 proc = subprocess.run(["gsutil", "cp",filename, BUCKET], stderr=subprocess.PIPE)
305 print("gstuil returned: " +str(proc.returncode))

```

```
306 | print(str(proc.stderr))
```

Writing spark\_job2c.py

```
In [ ]: 1 #Was getting error to create the cluster so I deleted it to create it a
2 !gcloud dataproc clusters delete big-data-cw-420116-cluster
```

The cluster 'big-data-cw-420116-cluster' and all attached disks will be deleted.

Do you want to continue (Y/n)? y

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/ead1de6a-5982-3c5d-ad8d-3651c4cfca3d].

Deleted [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>].

```
In [ ]: 1 ### CODING TASK ###
2 # cluster with a single machine using the maximal SSD size (100) and 1
3 #REGION = 'europe-west2'
4
5 import time
6 start_time = time.time()
7 !gcloud dataproc clusters create $CLUSTER \
8     --bucket $PROJECT-storage \
9     --image-version 1.5-ubuntu18 \
10    --master-machine-type n1-standard-8 \
11    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
12    --num-workers 0\
13    --initialization-actions gs://goog-dataproc-initialization-actions-\
14    --metadata "PIP_PACKAGES=tensorflow==2.4.0 protobuf==3.19.0"
15 %time
16 end_time = time.time()
17 execution_time = end_time - start_time
18 print(f"Execution time: {execution_time} seconds")
19 #     --metadata PIP_PACKAGES=tensorflow==2.4.0 was giving error
20 #     --initialization-actions gs://goog-dataproc-initialization-actions-\
21 #     --initialization-actions gs://goog-dataproc-initialization-actions-
```

Waiting on operation [projects/big-data-cw-420116/regions/us-central1/operations/ed6ecc51-42b7-3c3d-a94a-3fe2fd774f56].

**WARNING:** Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh

**WARNING:** The firewall rules for specified network or subnetwork would allow ingress traffic from 0.0.0.0/0, which could be a security risk.

Created [<https://dataproc.googleapis.com/v1/projects/big-data-cw-420116/regions/us-central1/clusters/big-data-cw-420116-cluster>] Cluster placed in zone [us-central1-c].

CPU times: user 4 µs, sys: 0 ns, total: 4 µs

Wall time: 7.63 µs

Execution time: 190.0779001712799 seconds

```
In [ ]: 1 # get information of the cluster
          2 !gcloud dataproc clusters describe $CLUSTER

clusterName: big-data-cw-420116-cluster
clusterUuid: 391c593b-c563-486e-8244-0386eb840cea
config:
    configBucket: big-data-cw-420116-storage
    endpointConfig: {}
    gceClusterConfig:
        internalIpOnly: false
    metadata:
        PIP_PACKAGES: tensorflow==2.4.0 protobuf==3.19.0
        networkUri: https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default (https://www.googleapis.com/compute/v1/projects/big-data-cw-420116/global/networks/default)
        serviceAccountScopes:
            - https://www.googleapis.com/auth/bigquery (https://www.googleapis.com/auth/bigquery)
            - https://www.googleapis.com/auth/bigtable.admin.table (https://www.googleapis.com/auth/bigtable.admin.table)
            - https://www.googleapis.com/auth/bigtable.data (https://www.googleapis.com/auth/bigtable.data)
```

```
In [ ]: 1 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \spark_job2c.py
```

Job [4e5e8b972aba461485e1f5ec3271cabd] submitted.  
Waiting for job output...  
2024-05-04 12:06:46.957241: W tensorflow/stream\_executor/platform/default/dso\_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dle error: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD\_LIBRARY\_PATH: :/usr/lib/hadoop/lib/native  
2024-05-04 12:06:46.957294: I tensorflow/stream\_executor/cuda/cudart\_stub.cc:29] Ignore above cudart dleerror if you do not have a GPU set up on your machine.  
False  
24/05/04 12:06:49 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker  
24/05/04 12:06:49 INFO org.apache.spark.SparkEnv: Registering BlockManager Master  
24/05/04 12:06:49 INFO org.apache.spark.SparkEnv: Registering OutputCommit Coordinator  
24/05/04 12:06:49 INFO org.spark\_project.jetty.util.log: Logging initialized @6271ms to org.spark\_project.jetty.util.log.Slf4jLog  
24/05/04 12:06:49 INFO org.spark\_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0\_382-b05  
24/05/04 12:06:49 INFO org.spark\_project.jetty.server.Server: Started @6393ms  
24/05/04 12:06:49 INFO org.spark\_project.jetty.server.AbstractConnector: Started ServerConnector@4f00657d{HTTP/1.1, (http/1.1)}{0.0.0.0:39371}  
24/05/04 12:06:51 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.0.16:8032  
24/05/04 12:06:51 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at big-data-cw-420116-cluster-m/10.128.0.16:10200  
24/05/04 12:06:51 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found  
24/05/04 12:06:51 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.  
24/05/04 12:06:51 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE  
24/05/04 12:06:51 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE  
24/05/04 12:06:54 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application\_1714824310332\_0001  
Saving 2c\_results.pkl to gs://big-data-cw-420116-storage  
gstuil returned: 0  
b'WARNING: Python 3.5-3.7 will be deprecated on August 8th, 2023. Please use Python version 3.8 and up.\n\nIf you have a compatible Python interpreter installed, you can use it by setting\nthe CLOUDSDK PYTHON environment variable to point to it.\nCopying file://2c\_results.pkl [Content-Type=application/octet-stream]...\n[ 0 files][ 0.0 B/ 11.2 KiB]\n[r [ 1 files][ 11.2 KiB/ 11.2 KiB]\n[r\nOperation completed over 1 objects/11.2 KiB.\n\n'  
24/05/04 12:09:11 INFO org.spark\_project.jetty.server.AbstractConnector: Stopped Spark@4f00657d{HTTP/1.1, (http/1.1)}{0.0.0.0:0}  
Job [4e5e8b972aba461485e1f5ec3271cabd] finished successfully.  
done: true  
driverControlFilesUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/391c593b-c563-486e-8244-0386eb840cea/jobs/4e5e8b972aba461485e1f5ec3271cabd/  
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/391c593b-c563-486e-8244-0386eb840cea/jobs/4e5e8b972aba461485e1f5ec3271cabd/driveroutput  
jobUuid: 3f2c7266-569a-3dbb-b3c4-75887dcda5ac  
placement:

```

clusterName: big-data-cw-420116-cluster
clusterUuid: 391c593b-c563-486e-8244-0386eb840cea
pysparkJob:
  mainPythonFileUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-
-metainfo/391c593b-c563-486e-8244-0386eb840cea/jobs/4e5e8b972aba461485e1f5
ec3271cabd/staging/spark_job2c.py
  reference:
    jobId: 4e5e8b972aba461485e1f5ec3271cabd
    projectId: big-data-cw-420116
  status:
    state: DONE
    stateStartTime: '2024-05-04T12:09:14.643674Z'
  statusHistory:
    - state: PENDING
      stateStartTime: '2024-05-04T12:06:41.998855Z'
    - state: SETUP_DONE
      stateStartTime: '2024-05-04T12:06:42.025243Z'
    - details: Agent reported job success
      state: RUNNING
      stateStartTime: '2024-05-04T12:06:42.322529Z'
  yarnApplications:
    - name: spark_job2c.py
      progress: 1.0
      state: FINISHED
      trackingUrl: http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714824310332\_0001/ (http://big-data-cw-420116-cluster-m:8088/proxy/application\_1714824310332\_0001/)

```

## 2d) Retrieve, analyse and discuss the output (12%)

Run the tests over a wide range of different parameters and list the results in a table.

Perform a **linear regression** (e.g. using scikit-learn) over the **values for each parameter** and for the **two cases** (reading from image files/reading TFRecord files). List a **table** with the output and interpret the results in terms of the effects of overall.

Also, **plot** the output values, the averages per parameter value and the regression lines for each parameter and for the product of batch\_size and batch\_number

Discuss the **implications** of this result for **applications** like large-scale machine learning. Keep in mind that cloud data may be stored in distant physical locations. Use the numbers provided in the PDF latency-numbers document available on Moodle or [here](https://gist.github.com/hellerbarde/2843375) (<https://gist.github.com/hellerbarde/2843375>) for your arguments.

How is the **observed** behaviour **similar or different** from what you'd expect from a **single machine**? Why would cloud providers tie throughput to capacity of disk resources?

By **parallelising** the speed test we are making **assumptions** about the limits of the bucket reading speeds. See [here](https://cloud.google.com/storage/docs/request-rate) (<https://cloud.google.com/storage/docs/request-rate>) for more information. Discuss, **what we need to consider** in **speed tests** in parallel on the cloud, which bottlenecks we might be identifying, and how this relates to your results.

Discuss to what extent **linear modelling** reflects the **effects** we are observing. Discuss what could be expected from a theoretical perspective and what can be useful in practice.

Write your **code below** and include the output in your submitted ipynb file. Provide the answer **text in your report**.

In [ ]:

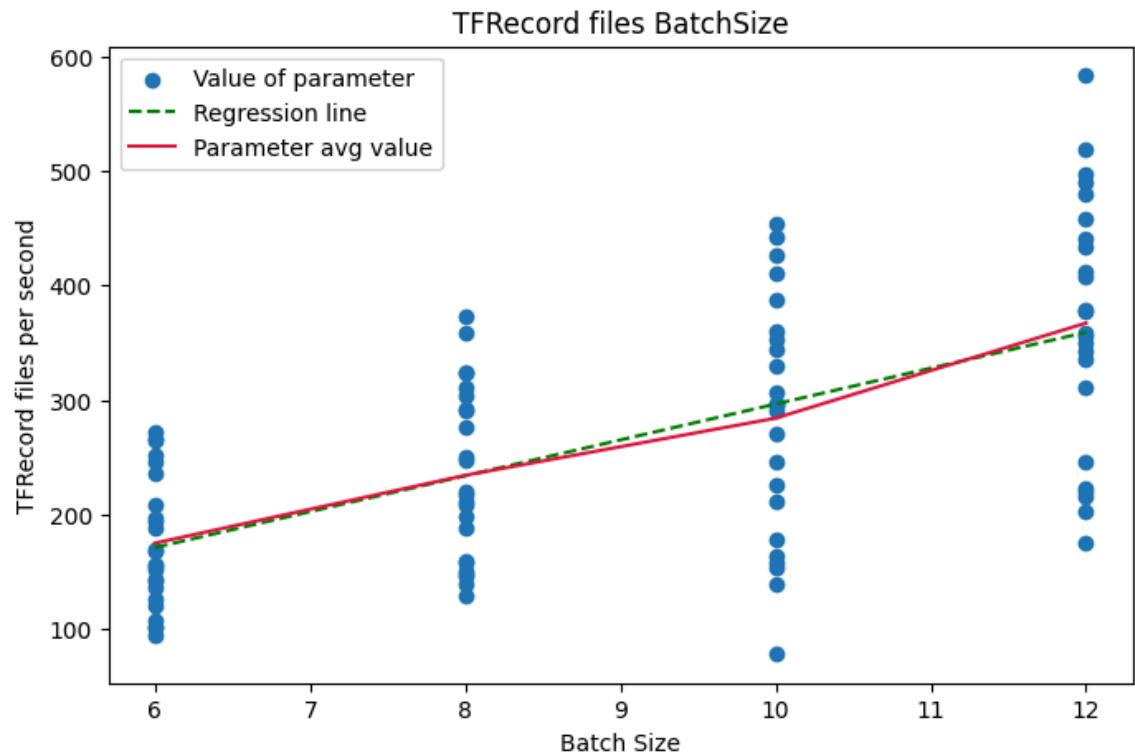
```
1 ### CODING TASK ###
2 !gsutil cp $BUCKET/2c_results.pkl .
3 with open("2c_results.pkl",mode = 'rb') as f:
4     tfr_batchSize_speed = pickle.load(f)
5     tfr_batchNum_speed = pickle.load(f)
6     tfr_repetitions_speed = pickle.load(f)
7     tfr_dsSize_speed = pickle.load(f)
8     image_batchSizes_speed = pickle.load(f)
9     image_batchNums_speed = pickle.load(f)
10    image_repetitions_speed = pickle.load(f)
11    image_dsSize_speed = pickle.load(f)
12    tfr_batchSize_avgSpeed = pickle.load(f)
13    tfr_batchNums_avgSpeed = pickle.load(f)
14    tfr_repetitions_avgSpeed = pickle.load(f)
15    tfr_dsSize_avgSpeed = pickle.load(f)
16    image_batchSizes_avgSpeed = pickle.load(f)
17    image_batchNums_avgSpeed = pickle.load(f)
18    image_repetitions_avgSpeed = pickle.load(f)
19    image_dsSize_avgSpeed = pickle.load(f)
```

```
Copying gs://big-data-cw-420116-storage/2c_results.pkl...
/ [1 files][ 11.2 KiB/ 11.2 KiB]
Operation completed over 1 objects/11.2 KiB.
```

this error is because the job couldn't be saved previously.

In [ ]:

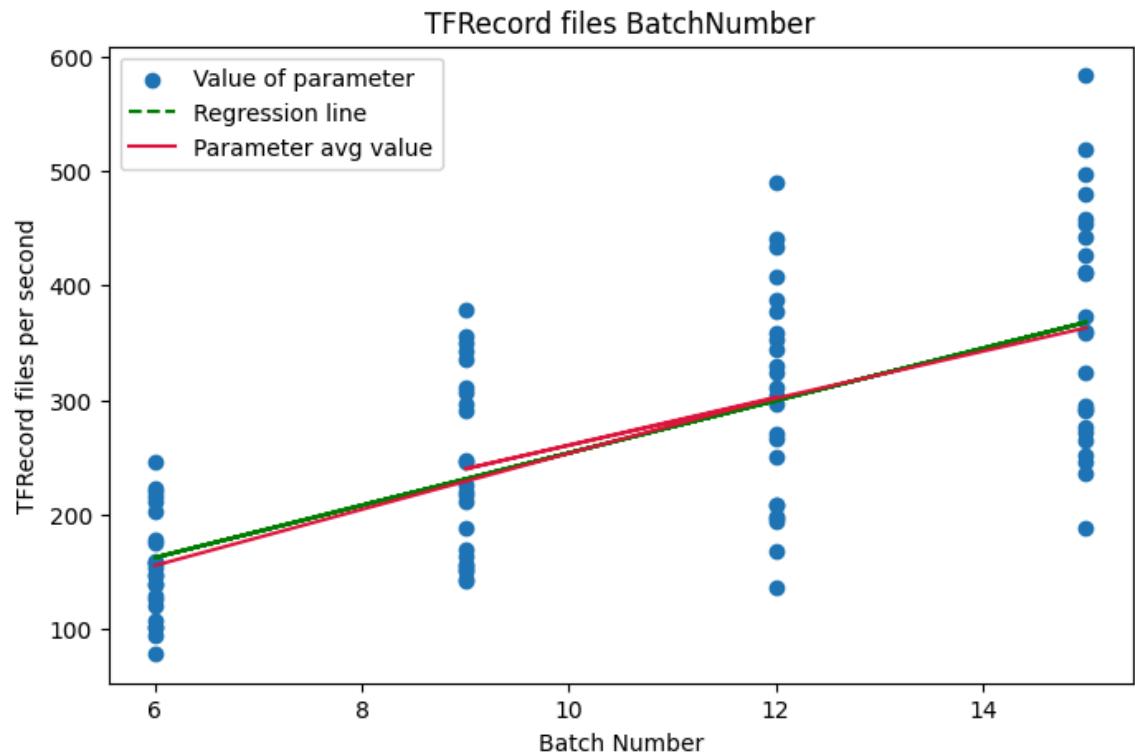
```
1 from sklearn.linear_model import LinearRegression
2
3 tfr_batchSize_speed_df = pd.DataFrame(tfr_batchSize_speed, columns=["batch_sizes", "TFRecord files per second"])
4 tfr_batchSize_avgSpeed_df = pd.DataFrame(tfr_batchSize_avgSpeed, columns=["batch_sizes", "TFRecord files per second"])
5
6 linreg_tfr_batchSize_speed_df = LinearRegression()
7 linreg_tfr_batchSize_speed_df.fit(tfr_batchSize_speed_df[['batch_sizes']], tfr_batchSize_speed_df['TFRecord files per second'])
8
9 tfr_batchSize_speed_predict = linreg_tfr_batchSize_speed_df.predict(tfr_batchSize_speed_df)
10
11
12 fig = plt.figure(figsize=(8,5))
13 plt.scatter(tfr_batchSize_speed_df[['batch_sizes']], tfr_batchSize_speed_df['TFRecord files per second'])
14 plt.plot(tfr_batchSize_speed_df[['batch_sizes']], tfr_batchSize_speed_predict)
15 plt.plot(tfr_batchSize_avgSpeed_df[['batch_sizes']], tfr_batchSize_avgSpeed_df['TFRecord files per second'])
16
17 plt.title('TFRecord files BatchSize')
18 plt.legend(loc='upper left')
19 plt.ylabel('TFRecord files per second')
20 plt.xlabel('Batch Size')
21 plt.show()
22
23 print("Slope : ",linreg_tfr_batchSize_speed_df.coef_)
24 print("Intercept : ",linreg_tfr_batchSize_speed_df.intercept_)
25 print('P-value : ',linreg_tfr_batchSize_speed_df.score(tfr_batchSize_speed_df))
```



Slope : [31.31917417]  
Intercept : -16.682792340220544  
P-value : 0.38712799159632494

In [ ]:

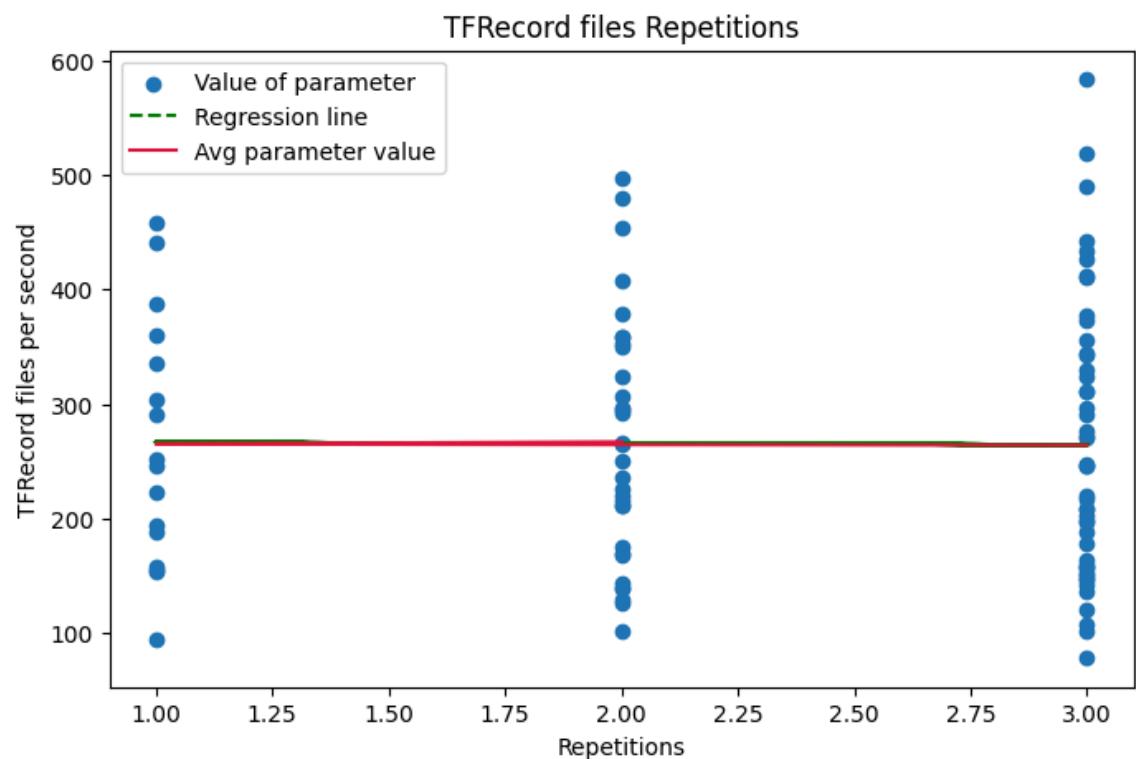
```
1 from sklearn.linear_model import LinearRegression
2
3 tfr_batchNum_speed_df = pd.DataFrame(tfr_batchNum_speed, columns=["batch_nums", "TFRecord_files_per_second"])
4 tfr_batchNums_avgSpeed_df = pd.DataFrame(tfr_batchNums_avgSpeed, columns=["batch_nums", "TFRecord_files_per_second"])
5
6 linreg_tfr_batchNum_speed_df = LinearRegression()
7 linreg_tfr_batchNum_speed_df.fit(tfr_batchNum_speed_df[['batch_nums']], tfr_batchNum_speed_df['TFRecord_files_per_second'])
8
9 tfr_batchNum_speed_predict = linreg_tfr_batchNum_speed_df.predict(tfr_batchNum_speed_df)
10
11
12 fig = plt.figure(figsize=(8,5))
13 plt.scatter(tfr_batchNum_speed_df[['batch_nums']], tfr_batchNum_speed_df['TFRecord_files_per_second'])
14 plt.plot(tfr_batchNum_speed_df[['batch_nums']], tfr_batchNum_speed_predict)
15 plt.plot(tfr_batchNums_avgSpeed_df[['batch_nums']], tfr_batchNums_avgSpeed_df['TFRecord_files_per_second'])
16
17 plt.title('TFRecord files BatchNumber')
18 plt.legend(loc='upper left')
19 plt.ylabel('TFRecord files per second')
20 plt.xlabel('Batch Number')
21 plt.show()
22
23 print("Slope : ",linreg_tfr_batchNum_speed_df.coef_)
24 print("Intercept : ",linreg_tfr_batchNum_speed_df.intercept_)
25 print('P-value : ',linreg_tfr_batchNum_speed_df.score(tfr_batchNum_speed_df))
```



Slope : [22.83459662]  
Intercept : 25.426510681410946  
P-value : 0.4630236527439795

In [ ]:

```
1 from sklearn.linear_model import LinearRegression
2
3 tfr_repetitions_speed_df = pd.DataFrame(tfr_repetitions_speed, columns=
4 tfr_repetitions_avgSpeed_df = pd.DataFrame(tfr_repetitions_avgSpeed, co
5
6 linreg_tfr_repetitions_speed_df = LinearRegression()
7 linreg_tfr_repetitions_speed_df.fit(tfr_repetitions_speed_df[['repetiti
8
9 tfr_repetitions_speed_predict = linreg_tfr_repetitions_speed_df.predict
10
11 fig = plt.figure(figsize=(8,5))
12 plt.scatter(tfr_repetitions_speed_df[['repetitions']], tfr_repetitions_
13 plt.plot(tfr_repetitions_speed_df[['repetitions']], tfr_repetitions_sp
14 plt.plot(tfr_repetitions_avgSpeed_df[['repetitions']],tfr_repetitions_a
15
16 plt.title('TFRecord files Repetitions')
17 plt.legend(loc='upper left')
18 plt.ylabel('TFRecord files per second')
19 plt.xlabel('Repetitions')
20 plt.show()
21
22 print("Slope : ",linreg_tfr_repetitions_speed_df.coef_)
23 print("Intercept : ",linreg_tfr_repetitions_speed_df.intercept_)
24 print('P-value : ',linreg_tfr_repetitions_speed_df.score(tfr_repetition
```

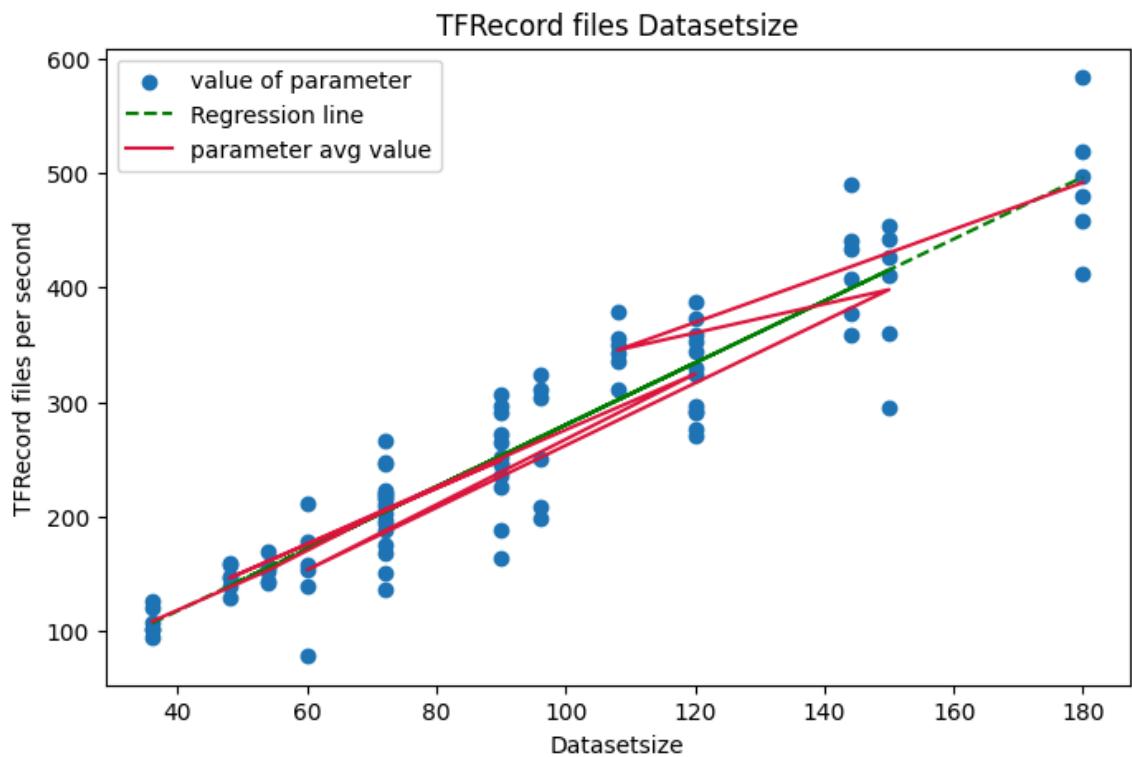


Slope : [-1.01936299]  
Intercept : 267.56828883248255  
P-value : 4.556686860013315e-05

```

In [ ]: 1 from sklearn.linear_model import LinearRegression
2
3 tfr_dsSize_speed_df = pd.DataFrame(tfr_dsSize_speed, columns=["datasetsize", "TFRecord files per second"])
4 tfr_dsSize_avgSpeed_df = pd.DataFrame(tfr_dsSize_avgSpeed, columns=["datasetsize", "TFRecord files per second"])
5
6 linreg_tfr_dsSize_speed_df = LinearRegression()
7 linreg_tfr_dsSize_speed_df.fit(tfr_dsSize_speed_df[['datasetsize']], tfr_dsSize_speed_df['TFRecord files per second'])
8
9 tfr_dsSize_speed_predict = linreg_tfr_dsSize_speed_df.predict(tfr_dsSize_speed_df)
10
11
12 fig = plt.figure(figsize=(8,5))
13 plt.scatter(tfr_dsSize_speed_df[['datasetsize']], tfr_dsSize_speed_df['TFRecord files per second'])
14 plt.plot(tfr_dsSize_speed_df[['datasetsize']], tfr_dsSize_speed_predict)
15 plt.plot(tfr_dsSize_avgSpeed_df[['datasetsize']], tfr_dsSize_avgSpeed_df['TFRecord files per second'])
16
17 plt.title('TFRecord files Datasetsize')
18 plt.legend(loc='upper left')
19 plt.ylabel('TFRecord files per second')
20 plt.xlabel('Datasetsize')
21 plt.show()
22
23 print("Slope : ",linreg_tfr_dsSize_speed_df.coef_)
24 print("Intercept : ",linreg_tfr_dsSize_speed_df.intercept_)
25 print('P-value : ',linreg_tfr_dsSize_speed_df.score(tfr_dsSize_speed_df))

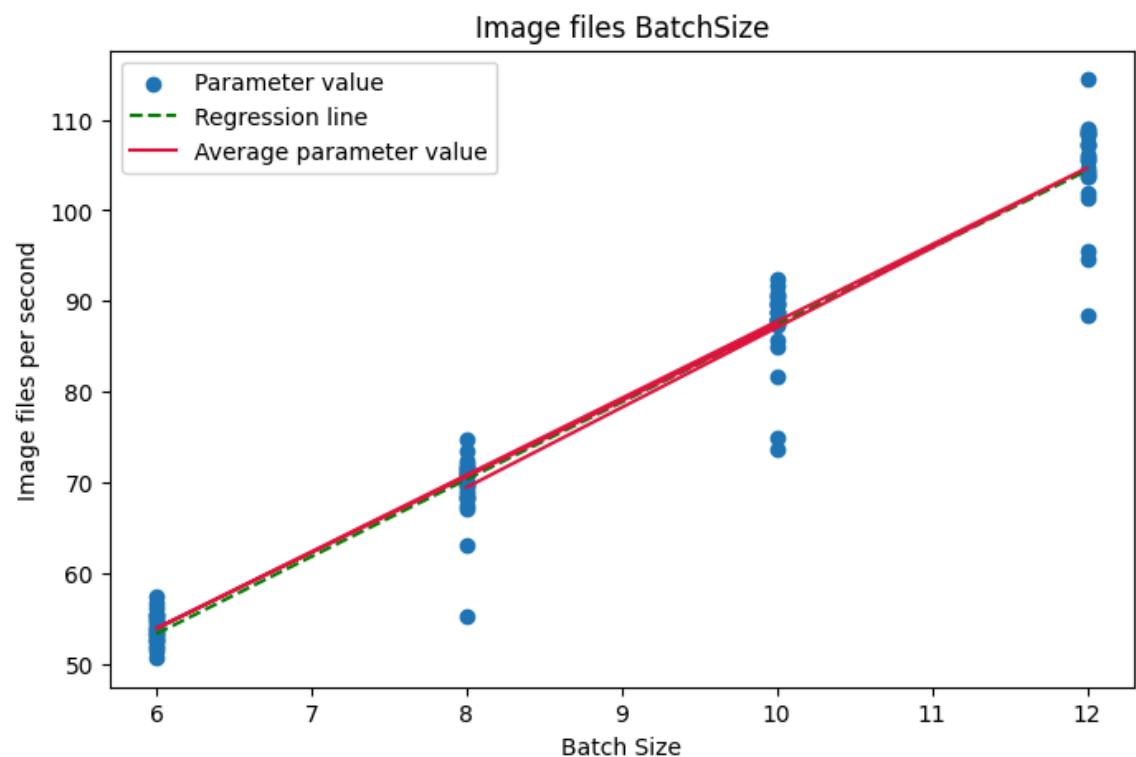
```



Slope : [2.70422835]  
 Intercept : 9.640195688634037  
 P-value : 0.8766716882967958

In [ ]:

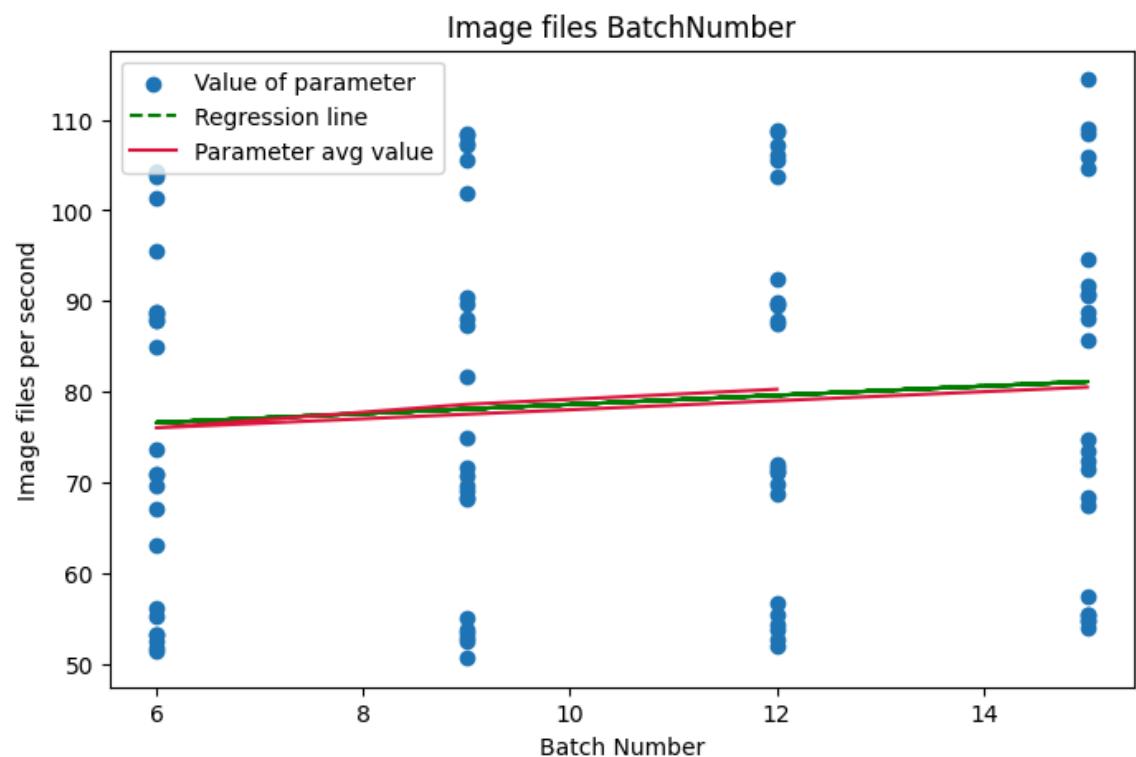
```
1 from sklearn.linear_model import LinearRegression
2
3 image_batchSizes_speed_df = pd.DataFrame(image_batchSizes_speed, columns=['batch_size', 'image_files_per_second'])
4 image_batchSizes_avgSpeed_df = pd.DataFrame(image_batchSizes_avgSpeed, columns=['batch_size', 'avg_speed'])
5
6 linreg_image_batchSizes_speed_df = LinearRegression()
7 linreg_image_batchSizes_speed_df.fit(image_batchSizes_speed_df[['batch_size']])
8
9 image_batchSizes_speed_predict = linreg_image_batchSizes_speed_df.predict(image_batchSizes_speed_df)
10
11 fig = plt.figure(figsize=(8,5))
12 plt.scatter(image_batchSizes_speed_df[['batch_size']], image_batchSizes_speed_df[['image_files_per_second']])
13 plt.plot(image_batchSizes_speed_df[['batch_size']], image_batchSizes_speed_predict)
14 plt.plot(image_batchSizes_avgSpeed_df[['batch_size']], image_batchSizes_avgSpeed_df[['avg_speed']])
15
16 plt.title('Image files BatchSize')
17 plt.legend(loc='upper left')
18 plt.ylabel('Image files per second')
19 plt.xlabel('Batch Size')
20 plt.show()
21
22 print("Slope : ",linreg_image_batchSizes_speed_df.coef_)
23 print("Intercept : ",linreg_image_batchSizes_speed_df.intercept_)
24 print('P-value : ',linreg_image_batchSizes_speed_df.score(image_batchSizes_speed_df))
```



Slope : [8.51630289]  
Intercept : 2.20325651298333  
P-value : 0.9558397732855116

In [ ]:

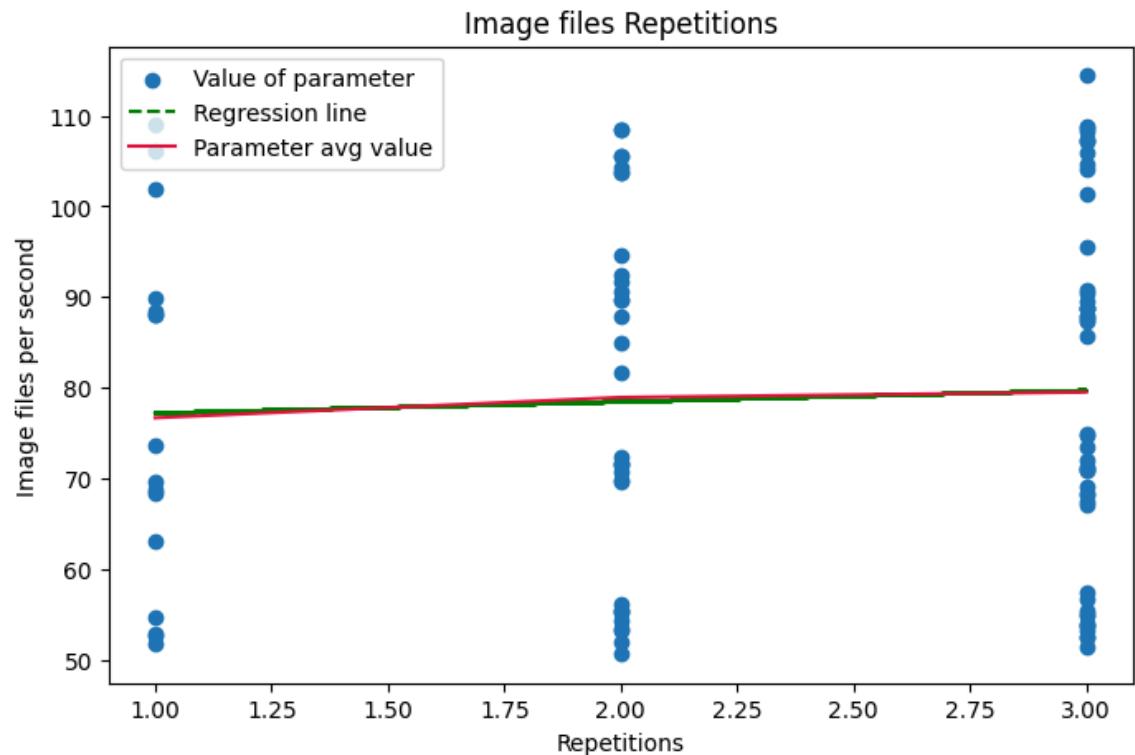
```
1 from sklearn.linear_model import LinearRegression
2
3 image_batchNums_speed_df = pd.DataFrame(image_batchNums_speed, columns=
4 image_batchNums_avgSpeed_df = pd.DataFrame(image_batchNums_avgSpeed, co
5
6 linreg_image_batchNums_speed_df = LinearRegression()
7 linreg_image_batchNums_speed_df.fit(image_batchNums_speed_df[['batch_nu
8
9 image_batchNums_speed_predict = linreg_image_batchNums_speed_df.predict
10
11 fig = plt.figure(figsize=(8,5))
12 plt.scatter(image_batchNums_speed_df[['batch_nums']], image_batchNums_s
13 plt.plot(image_batchNums_speed_df[['batch_nums']], image_batchNums_spee
14 plt.plot(image_batchNums_avgSpeed_df[['batch_nums']],image_batchNums_av
15
16 plt.title('Image files BatchNumber')
17 plt.legend(loc='upper left')
18 plt.ylabel('Image files per second')
19 plt.xlabel('Batch Number')
20 plt.show()
21
22 print("Slope : ",linreg_image_batchNums_speed_df.coef_)
23 print("Intercept : ",linreg_image_batchNums_speed_df.intercept_)
24 print('P-value : ',linreg_image_batchNums_speed_df.score(image_batchNum
```



Slope : [0.50348735]  
Intercept : 73.56336533975393  
P-value : 0.0075169651996684905

In [ ]:

```
1 from sklearn.linear_model import LinearRegression
2
3 image_repetitions_speed_df = pd.DataFrame(image_repetitions_speed, columns=['repetitions', 'image_files_per_second'])
4 image_repetitions_avgSpeed_df = pd.DataFrame(image_repetitions_avgSpeed, columns=['repetitions', 'image_files_per_second'])
5
6 linreg_image_repetitions_speed_df = LinearRegression()
7 linreg_image_repetitions_speed_df.fit(image_repetitions_speed_df[['repetitions']], image_repetitions_speed_df['image_files_per_second'])
8
9 image_repetitions_speed_predict = linreg_image_repetitions_speed_df.predict(image_repetitions_speed_df[['repetitions']])
10
11
12 fig = plt.figure(figsize=(8,5))
13 plt.scatter(image_repetitions_speed_df[['repetitions']], image_repetitions_speed_df[['image_files_per_second']])
14 plt.plot(image_repetitions_speed_df[['repetitions']], image_repetitions_speed_predict)
15 plt.plot(image_repetitions_avgSpeed_df[['repetitions']], image_repetitions_avgSpeed_df[['image_files_per_second']])
16
17 plt.title('Image files Repetitions')
18 plt.legend(loc='upper left')
19 plt.ylabel('Image files per second')
20 plt.xlabel('Repetitions')
21 plt.show()
22
23 print("Slope : ",linreg_image_repetitions_speed_df.coef_)
24 print("Intercept : ",linreg_image_repetitions_speed_df.intercept_)
25 print('P-value : ',linreg_image_repetitions_speed_df.score(image_repetitions_speed_df[['repetitions']]))
```

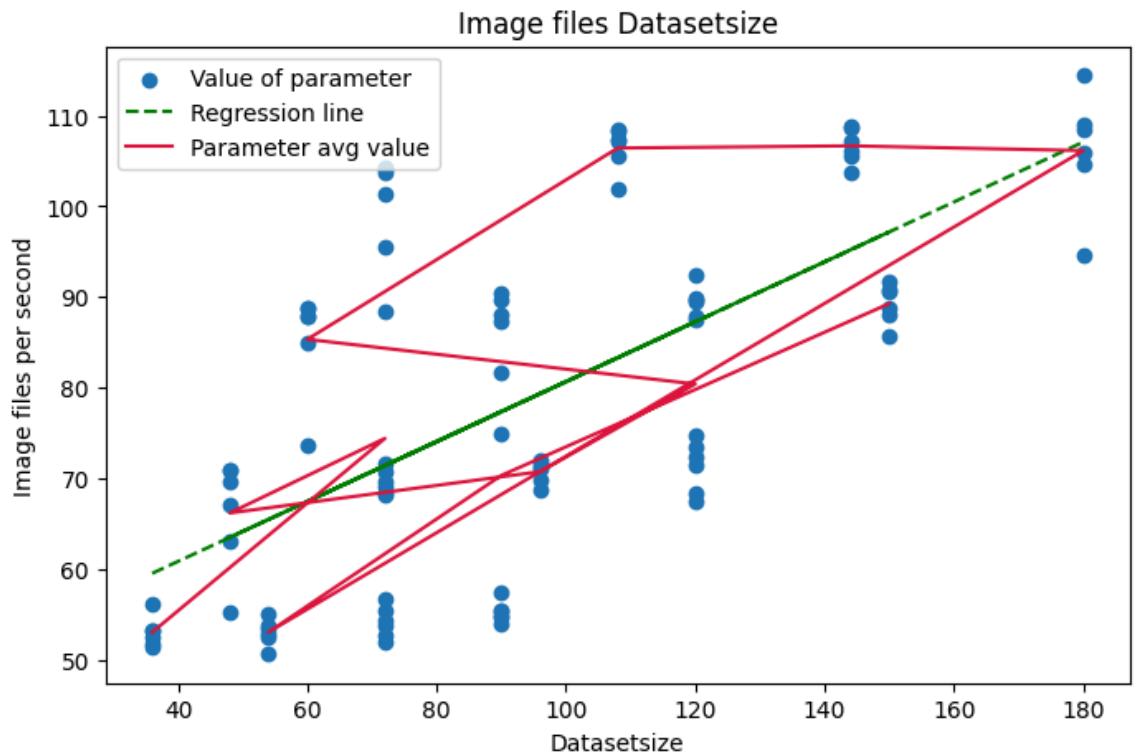


Slope : [1.24810677]  
Intercept : 75.937733354247  
P-value : 0.0022810975995531946

```

In [ ]: 1 from sklearn.linear_model import LinearRegression
2
3 image_dsSize_speed_df = pd.DataFrame(image_dsSize_speed, columns=["data")
4 image_dsSize_avgSpeed_df = pd.DataFrame(image_dsSize_avgSpeed, columns=
5
6 linreg_image_dsSize_speed_df = LinearRegression()
7 linreg_image_dsSize_speed_df.fit(image_dsSize_speed_df[['datasetsize']])
8
9 image_dsSize_speed_predict = linreg_image_dsSize_speed_df.predict(image_
10
11 fig = plt.figure(figsize=(8,5))
12
13 plt.scatter(image_dsSize_speed_df[['datasetsize']], image_dsSize_speed_
14 plt.plot(image_dsSize_speed_df[['datasetsize']], image_dsSize_speed_pre_
15 plt.plot(image_dsSize_avgSpeed_df[['datasetsize']],image_dsSize_avgSpee_
16
17 plt.title('Image files Datasetsize')
18 plt.legend(loc='upper left')
19 plt.ylabel('Image files per second')
20 plt.xlabel('Datasetsize')
21 plt.show()
22
23 print("Slope : ",linreg_image_dsSize_speed_df.coef_)
24 print("Intercept : ",linreg_image_dsSize_speed_df.intercept_)
25 print('P-value : ',linreg_image_dsSize_speed_df.score(image_dsSize_spee

```



Slope : [0.33019057]  
 Intercept : 47.64697402103596  
 P-value : 0.43644397991167594

## Section 3. Theoretical discussion

## Task 3: Discussion in context. (24%)

In this task we refer an idea that is introduced in this paper:

- Alipourfard, O., Liu, H. H., Chen, J., Venkataraman, S., Yu, M., & Zhang, M. (2017). [Cherrypick: Adaptively unearthing the best cloud configurations for big data analytics.](#) (<https://people.irisa.fr/Davide.Frey/wp-content/uploads/2018/02/cherrypick.pdf>). In USENIX NSDI 17 (pp. 469-482).

Alipourfard et al (2017) introduce the prediction an optimal or near-optimal cloud configuration for a given compute task.

### 3a) Contextualise

Relate the previous tasks and the results to this concept. (It is not necessary to work through the full details of the paper, focus just on the main ideas). To what extent and under what conditions do the concepts and techniques in the paper apply to the task in this coursework? (12%)

### 3b) Strategise

Define - as far as possible - concrete strategies for different application scenarios (batch, stream) and discuss the general relationship with the concepts above. (12%)

Provide the answers to these questions in your report.

## Final cleanup

Once you have finished the work, you can delete the buckets, to stop incurring cost that depletes your credit.

In [ ]:

```
1 !gsutil -m rm -r $BUCKET/* # Empty your bucket
2 !gsutil rb $BUCKET # delete the bucket
```

```
Removing gs://big-data-cw-420116-storage/2aVI_results.pkl#17148220408016
53...
Removing gs://big-data-cw-420116-storage/2b_results.pkl#171482334463637
7...
Removing gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/
084e9d65-6812-45a2-9cab-372f310ba794/jobs/e72d8ce747414b90b5b0ef0b944ed2
55/staging/spark_write_tfrec.py#1714815915867542...
Removing gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/
084e9d65-6812-45a2-9cab-372f310ba794/jobs/e72d8ce747414b90b5b0ef0b944ed2
55/driveroutput.00000001#1714815935242074...
Removing gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/
084e9d65-6812-45a2-9cab-372f310ba794/big-data-cw-420116-cluster-m/dataproc-
post-hdfs-startup-script_output#1714815870063135...
Removing gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/
084e9d65-6812-45a2-9cab-372f310ba794/jobs/e72d8ce747414b90b5b0ef0b944ed2
55/driveroutput.00000000#1714815935051527...
Removing gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/
084e9d65-6812-45a2-9cab-372f310ba794/cluster.properties#171481575134134
8...
```

GO TO BUCKET, DELET WORKS: [https://console.cloud.google.com/storage/browser?project=big-data-cw-420116&prefix=&forceOnBucketsSortingFiltering=true&pageState=%22StorageBucketsTable%22:\(%22f%22:%22%255B%255D%22,%22s%22:%5B\(%22i%22:%22name%22,%22s%22:%220https://console.cloud.google.com/storage/browser?project=big-data-cw-420116&prefix=&forceOnBucketsSortingFiltering=true&pageState=%22StorageBucketsTable%22:\(%22f%22:%22%255B%255D%22,%22s%22:%5B\(%22i%22:%22name%22,%22s%22:%220](https://console.cloud.google.com/storage/browser?project=big-data-cw-420116&prefix=&forceOnBucketsSortingFiltering=true&pageState=%22StorageBucketsTable%22:(%22f%22:%22%255B%255D%22,%22s%22:%5B(%22i%22:%22name%22,%22s%22:%220https://console.cloud.google.com/storage/browser?project=big-data-cw-420116&prefix=&forceOnBucketsSortingFiltering=true&pageState=%22StorageBucketsTable%22:(%22f%22:%22%255B%255D%22,%22s%22:%5B(%22i%22:%22name%22,%22s%22:%220)

