

✓ 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**](#) 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 without needing the solution to the earlier ones.

▼ 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.

```
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.

```
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.

```
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](#) 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.

We also specify the **default project and region**. The REGION should be us-central1 as that seems to be the only one that reliably works with the free credit. This way we don't have to specify this information every time we access the cloud.

```
1 PROJECT = 'big-data-cw-420116' ### USE YOUR GOOGLE CLOUD PROJECT ID HERE. ### (DONE!)
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**.

```
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-storage' already exists. Try another name. Bucket names must be g
```

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.

➤ Utility functions for image display [RUN THIS TO ACTIVATE]

[Show code](#)

```
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
bigquery.googleapis.com       BigQuery API
bigqueryconnection.googleapis.com BigQuery Connection API
bigquerystorage.googleapis.com BigQuery Storage API
bigquerystorage.googleapis.com BigQuery Storage API
bigquerystorage.googleapis.com Google Cloud APIs
cloudapis.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.

[] ↓ 1 cell hidden

✓ 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**.

[] ↓ 12 cells hidden

✓ 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**.

```
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.decode('utf8') for lbl in label.numpy()]))

Image batch shape (10,), ['sunflowers', 'tulips', 'tulips', 'tulips', 'daisy', 'sunflowers', 'dandelion', 'roses', 'daisy', 'tulips']
Image batch shape (10,), ['daisy', 'roses', 'roses', 'tulips', 'dandelion', 'tulips', 'daisy', 'roses', 'tulips', 'dandelion']
Image batch shape (10,), ['dandelion', 'sunflowers', 'daisy', 'roses', 'sunflowers', 'dandelion', 'tulips', 'tulips', 'daisy', 'roses']
Image batch shape (10,), ['daisy', 'tulips', 'dandelion', 'roses', 'roses', 'daisy', 'roses', 'daisy', 'sunflowers']
Image batch shape (10,), ['tulips', 'roses', 'daisy', 'dandelion', 'dandelion', 'dandelion', 'dandelion', 'dandelion', 'dandelion']
Image batch shape (10,), ['tulips', 'dandelion', 'roses', 'daisy', 'tulips', 'sunflowers', 'dandelion', 'dandelion', 'roses', 'tulips']
Image batch shape (10,), ['daisy', 'tulips', 'daisy', 'dandelion', 'daisy', 'tulips', 'tulips', 'sunflowers', 'roses']
Image batch shape (10,), ['daisy', 'roses', 'sunflowers', 'tulips', 'dandelion', 'dandelion', 'dandelion', 'dandelion', 'dandelion']
Image batch shape (10,), ['dandelion', 'sunflowers', 'dandelion', 'roses', 'roses', 'daisy', 'dandelion', 'dandelion', 'sunflow
Image batch shape (10,), ['tulips', 'daisy', 'dandelion', 'dandelion', 'dandelion', 'daisy', 'tulips', 'sunflowers', 'rose
```

✓ 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.

```
1 GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # prefix for output file names
```

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.

```
1 # functions for writing TFRecord entries
2 # Feature values are always stored as lists, a single data element will be a list of size 1
3 def _bytestring_feature(list_of_bytestrings):
4     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))
5
6 def _int_feature(list_of_ints): # int64
7     return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
8
9 def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data records
10    class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
11    one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
12    feature = {
13        "image": _bytestring_feature([img_bytes]), # one image in the list
14        "class": _int_feature([class_num]) #, # one class in the list
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 images to files.
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 will be one "batch" of images per file
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, partition_size)
31         # You need to change GCS_OUTPUT to your own bucket to actually create new files
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 image: already a byte string
36                     label.numpy()[i] #
37                 )
38                 out_file.write(example.SerializeToString())
39     print("Wrote file {} containing {} records".format(filename, partition_size))
40     print("Total time: "+str(time.time()-tt0))
41
42 write_tfrecords(GCS_PATTERN,GCS_OUTPUT,partition_size) # uncomment to run this cell
```

```
Writing TFRecords
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers00-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers01-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers02-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers03-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers04-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers05-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers06-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers07-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers08-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers09-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers10-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers11-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers12-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers13-230.tfrec containing 230 records
```

```
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers14-230.tfrec containing 230 records
Wrote file gs://big-data-cw-420116-storage/tfrecords-jpeg-192x192-2/flowers15-220.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.

```
1
2 def read_tfrecord(example):
3     features = {
4         "image": tf.io.FixedLenFeature([], tf.string), # tf.string = bytestring (not text string)
5         "class": tf.io.FixedLenFeature([], tf.int64) #,   # shape [] means scalar
6     }
7     # decode the TFRecord
8     example = tf.io.parse_single_example(example, features)
9     image = tf.image.decode_jpeg(example['image'], channels=3)
10    image = tf.reshape(image, [*TARGET_SIZE, 3])
11    class_num = example['class']
12    return image, class_num
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 = False
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**.

```
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()]))
6
7 Image batch shape (10, 192, 192, 3), ['1', '0', '4', '4', '1', '3', '2', '1', '4', '2']
8 Image batch shape (10, 192, 192, 3), ['1', '0', '4', '1', '2', '1', '1', '4', '0', '1']
9 Image batch shape (10, 192, 192, 3), ['1', '2', '3', '1', '0', '1', '2', '4', '4', '3']
10 Image batch shape (10, 192, 192, 3), ['1', '3', '1', '4', '2', '3', '0', '3', '0', '3']
11 Image batch shape (10, 192, 192, 3), ['3', '3', '3', '4', '3', '0', '1', '0', '4', '2']
12 Image batch shape (10, 192, 192, 3), ['2', '4', '1', '3', '0', '3', '4', '3', '3', '1']
13 Image batch shape (10, 192, 192, 3), ['0', '4', '3', '4', '0', '4', '3', '0', '3']
14 Image batch shape (10, 192, 192, 3), ['0', '0', '1', '0', '2', '4', '2', '2', '3', '0']
15 Image batch shape (10, 192, 192, 3), ['2', '1', '1', '4', '1', '2', '3', '3', '0', '3']
16 Image batch shape (10, 192, 192, 3), ['0', '4', '0', '4', '4', '1', '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'**](#) 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%)

```

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 files
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 output file names in my bucket
29
30 #i) Copy over the mapping functions (see section 1.1) and adapt the resizing and recompression functions to Spark (only one argument). (3%
31 def decode_jpeg_and_label(filepath):
32     # extracts the image data and creates a class label, based on the filepath
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 image
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]), # if true
52                     lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
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 - th) // 2, tw, th)
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_downsampling=False)
63     return image, label
64
65 # ii) Replace the TensorFlow Dataset objects with RDDs, starting with an RDD that contains the list of image filenames. (3%
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 in the code. Specify a sampling factor of 0.02 for short tests.
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 resize and crop image, rdd for recompress image
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[0],x[1]))
83 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress_image(x[0],x[1]))
84
85 # iv) Then use the functions from above to write the TFRecord files, using an RDD as the vehicle for parallelisation but not for storing
86
87 # functions for writing TFRecord entries
88 # Feature values are always stored as lists, a single data element will be a list of size 1
89 def _bytestring_feature(list_of_bytestrings):
90     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))
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 records
96     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
97     one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
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 files.
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 create new files
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 a byte string
115                 label.numpy() #
116             )
117             out_file.write(example.SerializeToString())
118             #print("Wrote file {} containing {} records".format(filename, partition_size))
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 function, that can be applied to every partition with the 'RDD.map'
123 #   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 TFRec
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-pyspark
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

```

1 # First elements of the RDD
2 decode_jpeg_and_label_rdd.take(1)

[(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],
       ...,
       [152, 152, 150],
       [155, 155, 153],
       [148, 148, 146]]),

```

```
[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'>]
```

```
1 # first elements of the RDD
2 resize_and_crop_image_rdd.take(1)
```

```
((<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'>)]
```

```
1 # First elements of the RDD  
2 recompress_image_rdd.take(1)
```

```
1 # checking name  
2 filenames rdd.take(1)
```

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

> 1b) Testing (3%)

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

```
[ ] ↓ 4 cells hidden
```

✓ 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%)

```
1 !gcloud dataproc clusters delete big-data-cw-420116-cluster --region us-central1
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

1           ##### CODING TASK #####
2 #Set up a cluster with a single machine using the maximal SSD size (100) and 8 vCPUs.
3 #REGION = 'us-central1'
4 !gcloud dataproc clusters create $CLUSTER \
5 --bucket $PROJECT-storage \
6 --image-version 1.4-ubuntu18 --single-node \
7 --master-machine-type n1-standard-8 \
8 --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \
9 --max-idle 3600s \
10 --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
11 --metadata PIP_PACKAGES=tensorflow==2.4.0
12
13 #--metadata "PIP_PACKAGES=tensorflow==2.4.0 scipy==1.4.1 matplotlib==3.7.0 numpy==1.25.2" # Tried this but didn't work
14
15 # to grant uniform access !gcloud storage buckets update $BUCKET --uniform-bucket-level-access

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
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

1 #!gcloud storage buckets update $BUCKET --uniform-bucket-level-access

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
    serviceAccountScopes:
```

```
- https://www.googleapis.com/auth/bigquery
- https://www.googleapis.com/auth/bigtable.admin.table
- https://www.googleapis.com/auth/bigtable.data
- https://www.googleapis.com/auth/cloud.useraccounts.readonly
- https://www.googleapis.com/auth/devstorage.full_control
- https://www.googleapis.com/auth/devstorage.read_write
- https://www.googleapis.com/auth/logging.write
zoneUri: 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
lifecycleConfig:
  idleDeleteTtl: 3600s
  idleStartTime: '2024-05-04T09:42:30.961969Z'
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
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
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
```

Run the script in the cloud and test the output.

```
1 # Submit jobs to cluster and measure execution time (Checked& NOT Working)
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \spark_write_tfrec.py
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'
2024-05-04 09:45:21.057575: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dlerror if you do not have a GPU se
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 BlockManagerMaster
24/05/04 09:45:23 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
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 @5943ms
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 sc
24/05/04 09:45:25 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at big-data-cw-420116-cluster-m/10.128.15.22
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/
24/05/04 09:45:27 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714915815041_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
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
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/
CPU times: user 4 µs, sys: 0 ns, total: 4 µs
Wall time: 9.54 µs

```

Double-click (or enter) to edit

*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](#) 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%)

[] ↓ 8 cells hidden

➤ 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.


```

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 files
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 file names in my bucket
33
34 def decode_jpeg_and_label(filepath):
35     # extracts the image data and creates a class label, based on the filepath
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]), # if true
55                     lambda: tf.image.resize(image, [w*th/h, h*th/h])) # if false
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     return image, label
61
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_downsampling=False)
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 in the initial call 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 resize and crop image, rdd for recompress image
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[0],x[1]))
83 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress_image(x[0],x[1]))
84
85 # functions for writing TFRecord entries
86 # Feature values are always stored as lists, a single data element will be a list of size 1
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 records
94     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
95     one_hot_class = np.eye(len(CLASSES))[class_num]      # [0, 0, 1, 0, 0] for class #2, roses
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 create new files
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 a byte string
113                                 label.numpy() #
114                                 )
115             out_file.write(example.SerializeToString())
116             #print("Wrote file {} containing {} records".format(filename, partition_size))
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_tfrecords)

```

Writing task_1di.py

```

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-central1

```

ERROR: (gcloud.dataproc.clusters.delete) Error parsing [cluster].
The [cluster] resource is not properly specified.
Failed to find attribute [project]. The attribute can be set in the following ways:
- provide the argument `cluster` on the command line with a fully specified name
- provide the argument `--project` on the command line
- set the property `core/project`

```

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 $PROJECT-storage \
10  --image-version 1.4-ubuntu18 \
11  --master-machine-type n1-standard-1 \
12  --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
13  --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-disk-size 100\
14  --max-idle 3600s \
15  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
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")

```

ERROR: (gcloud.dataproc.clusters.create) argument --bucket: expected one argument
Usage: gcloud dataproc clusters create (CLUSTER : --region=REGION) [optional flags]
optional flags may be --action-on-failed-primary-workers | --no-address |
--async | --autoscaling-policy | --bucket |
--confidential-compute | --dataproc-metastore |
--driver-pool-accelerator |
--driver-pool-boot-disk-size |
--driver-pool-boot-disk-type | --driver-pool-id |
--driver-pool-local-ssd-interface |
--driver-pool-machine-type |
--driver-pool-min-cpu-platform | --driver-pool-size |
--enable-component-gateway | --enable-kerberos |
--expiration-time | --gce-pd-kms-key |
--gce-pd-kms-key-keyring | --gce-pd-kms-key-location |
--gce-pd-kms-key-project | --help |
--identity-config-file | --image | --image-version |
--initialization-action-timeout |
--initialization-actions | --kerberos-config-file |
--kerberos-kms-key | --kerberos-kms-key-keyring |
--kerberos-kms-key-location |
--kerberos-kms-key-project |
--kerberos-root-principal-password-uri | --kms-key |
--kms-keyring | --kms-location | --kms-project |
--labels | --master-accelerator |
--master-boot-disk-size | --master-boot-disk-type |
--master-local-ssd-interface | --master-machine-type |
--master-min-cpu-platform | --max-age | --max-idle |
--metadata | --metric-overrides |
--metric-overrides-file | --metric-sources |
--min-num-workers | --min-secondary-worker-fraction |
--network | --node-group |
--num-driver-pool-local-ssds |
--num-master-local-ssds | --num-masters |
--num-secondary-worker-local-ssds |
--num-secondary-workers | --num-worker-local-ssds |
--num-workers | --optional-components |
--private-ipv6-google-access-type | --properties |
--public-ip-address | --region | --reservation |
--reservation-affinity | --scopes |
--secondary-worker-accelerator |
--secondary-worker-boot-disk-size |
--secondary-worker-boot-disk-type |
--secondary-worker-local-ssd-interface |
--secondary-worker-machine-types |
--secondary-worker-type |
--secure-multi-tenancy-user-mapping |
--service-account | --shielded-integrity-monitoring |
--shielded-secure-boot | --shielded-vtpm |
--single-node | --subnet | --tags | --temp-bucket |
--worker-accelerator | --worker-boot-disk-size |
--worker-boot-disk-type |
--worker-local-ssd-interface | --worker-machine-type |
--worker-min-cpu-platform | --zone

For detailed information on this command and its flags, run:

```

gcloud dataproc clusters create --help
CPU times: user 2 µs, sys: 0 ns, total: 2 µs
Wall time: 6.44 µs

```

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

ERROR: (gcloud.dataproc.clusters.describe) argument (CLUSTER : --region=REGION): Must be specified.
Usage: gcloud dataproc clusters describe (CLUSTER : --region=REGION) [optional flags]
optional flags may be --help | --region

For detailed information on this command and its flags, run:
gcloud dataproc clusters describe --help

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

ERROR: (gcloud.dataproc.jobs.submit.pyspark) argument PY_FILE: Must be specified.
Usage: gcloud dataproc jobs submit pyspark PY_FILE (--cluster=CLUSTER | --cluster-labels=[KEY=VALUE,...]) [optional flags] [--JOB_ARGS optional flags may be --archives | --async | --bucket | --cluster |
--cluster-labels | --driver-log-levels |
--driver-required-memory-mb |
--driver-required-vcores | --files | --help | --jars |
--labels | --max-failures-per-hour |
--max-failures-total | --properties |
--properties-file | --py-files | --region

For detailed information on this command and its flags, run:
gcloud dataproc jobs submit pyspark --help

```

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 files
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 file names in my bucket
33
34 def decode_jpeg_and_label(filepath):
35     # extracts the image data and creates a class label, based on the filepath
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]), # if true
55                     lambda: tf.image.resize(image, [w*th/h, h*th/h])) # if false
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     return image, label
61
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_downsampling=False)
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 in the initial call 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 resize and crop image, rdd for recompress image
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[0],x[1]))
83 recompress_image_rdd = resize_and_crop_image_rdd.map(lambda x: recompress_image(x[0],x[1]))
84
85 # functions for writing TFRecord entries
86 # Feature values are always stored as lists, a single data element will be a list of size 1
87 def _bytestring_feature(list_of_bytestrings):
88     return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))
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 records
94     class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
95     one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
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 create new files
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 a byte string
113                                 label.numpy() #
114                                 )
115             out_file.write(example.SerializeToString())
116             #print("Wrote file {} containing {} records".format(filename, partition_size))
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_tfrecords)

```

Overwriting task_1di.py

```

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-central1

```

```

ERROR: (gcloud.dataproc.clusters.delete) Error parsing [cluster].
The [cluster] resource is not properly specified.
Failed to find attribute [project]. The attribute can be set in the following ways:
- provide the argument `cluster` on the command line with a fully specified name
- provide the argument `--project` on the command line
- set the property `core/project`

```

```

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-disk-size 100\
12    --max-idle 3600s \
13    --initialization-actions gs://goog-dataproc-initialization-actions-us-central1/python/pip-install.sh \
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-$REGION/python/pip-install.sh \
ERROR: (gcloud.dataproc.clusters.create) argument --bucket: expected one argument
Usage: gcloud dataproc clusters create (CLUSTER : --region=REGION) [optional flags]
  optional flags may be --action-on-failed-primary-workers | --no-address |
    --async | --autoscaling-policy | --bucket |
    --confidential-compute | --dataproc-metastore |
    --driver-pool-accelerator |
    --driver-pool-boot-disk-size |
    --driver-pool-boot-disk-type | --driver-pool-id |
    --driver-pool-local-ssd-interface |
    --driver-pool-machine-type |
    --driver-pool-min-cpu-platform | --driver-pool-size |
    --enable-component-gateway | --enable-kerberos |
    --expiration-time | --gce-pd-kms-key |
    --gce-pd-kms-key-keyring | --gce-pd-kms-key-location |
    --gce-pd-kms-key-project | --help |
    --identity-config-file | --image | --image-version |
    --initialization-action-timeout |
    --initialization-actions | --kerberos-config-file |
    --kerberos-kms-key | --kerberos-kms-key-keyring |
    --kerberos-kms-key-location |
    --kerberos-kms-key-project |
    --kerberos-root-principal-password-uri | --kms-key |
    --kms-keyring | --kms-location | --kms-project |
    --labels | --master-accelerator |
    --master-boot-disk-size | --master-boot-disk-type |
    --master-local-ssd-interface | --master-machine-type |
    --master-min-cpu-platform | --max-age | --max-idle |
    --metadata | --metric-overrides |
    --metric-overrides-file | --metric-sources |
    --min-num-workers | --min-secondary-worker-fraction |
    --network | --node-group |
    --num-driver-pool-local-ssds |
    --num-master-local-ssds | --num-masters |
    --num-secondary-worker-local-ssds |
    --num-secondary-workers | --num-worker-local-ssds |
    --num-workers | --optional-components |
    --private-ipv6-google-access-type | --properties |
    --public-ip-address | --region | --reservation |
    --reservation-affinity | --scopes |
    --secondary-worker-accelerator |
    --secondary-worker-boot-disk-size |
    --secondary-worker-boot-disk-type |
    --secondary-worker-local-ssd-interface |
    --secondary-worker-machine-types |
    --secondary-worker-type |
    --secure-multi-tenancy-user-mapping |
    --service-account | --shielded-integrity-monitoring |
    --shielded-secure-boot | --shielded-vtpm |
    --single-node | --subnet | --tags | --temp-bucket |
    --worker-accelerator | --worker-boot-disk-size |
    --worker-boot-disk-type |
    --worker-local-ssd-interface | --worker-machine-type |
    --worker-min-cpu-platform | --zone

```

For detailed information on this command and its flags, run:

```

gcloud dataproc clusters create --help
CPU times: user 4 µs, sys: 1 µs, total: 5 µs
Wall time: 7.87 µs
... --- - - - - -

```

1 # Description of the cluster

```
2 !gcloud dataproc clusters describe $CLUSTER
```

```
ERROR: (gcloud.dataproc.clusters.describe) argument (CLUSTER : --region=REGION): Must be specified.
Usage: gcloud dataproc clusters describe (CLUSTER : --region=REGION) [optional flags]
  optional flags may be --help | --region
```

For detailed information on this command and its flags, run:

```
gcloud dataproc clusters describe --help
```

1 # Submit the job

```
2 !gcloud dataproc jobs submit pyspark --cluster $CLUSTER \task_1di.py
```

```
ERROR: (gcloud.dataproc.jobs.submit.pyspark) argument PY_FILE: Must be specified.
Usage: gcloud dataproc jobs submit pyspark PY_FILE (--cluster=CLUSTER | --cluster-labels=[KEY=VALUE,...]) [optional flags] [-- JOB_ARGS]
  optional flags may be --archives | --async | --bucket | --cluster |
    --cluster-labels | --driver-log-levels |
    --driver-required-memory-mb |
    --driver-required-vcores | --files | --help | --jars |
```

```
--labels | --max-failures-per-hour |
--max-failures-total | --properties |
--properties-file | --py-files | --region
```

For detailed information on this command and its flags, run:
`gcloud dataproc jobs submit pyspark --help`

```
1 # ii) Experiment with cluster configurations
```

```
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-central1
```

```
ERROR: (gcloud.dataproc.clusters.delete) Error parsing [cluster].
The [cluster] resource is not properly specified.
Failed to find attribute [project]. The attribute can be set in the following ways:
- provide the argument `cluster` on the command line with a fully specified name
- provide the argument `--project` on the command line
- set the property `core/project`
```

```
1 # II) 4 machines with double the resources each (2 vCPUs, memory, disk) I need to change to 100 and take screenshot
2 import time
3 start_time = time.time()
4
5 REGION = 'us-central1'
6 !gcloud dataproc clusters create $CLUSTER \
7   --bucket $PROJECT-storage \
8   --image-version 1.4-ubuntu18 \
9   --master-machine-type n1-standard-2 \
10  --master-boot-disk-type pd-ssd --master-boot-disk-size 200 \
11  --num-workers 3 --worker-machine-type n1-standard-2 --worker-boot-disk-size 200 \
12  --max-idle 3600s \
13  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
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-clusters/init-actions
24 # https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/cluster-properties
```

```
ERROR: (gcloud.dataproc.clusters.create) argument --bucket: expected one argument
Usage: gcloud dataproc clusters create (CLUSTER : --region=REGION) [optional flags]
optional flags may be --action-on-failed-primary-workers | --no-address |
--async | --autoscaling-policy | --bucket |
--confidential-compute | --dataproc-metastore |
--driver-pool-accelerator |
--driver-pool-boot-disk-size |
--driver-pool-boot-disk-type | --driver-pool-id |
--driver-pool-local-ssd-interface |
--driver-pool-machine-type |
--driver-pool-min-cpu-platform | --driver-pool-size |
--enable-component-gateway | --enable-kerberos |
--expiration-time | --gce-pd-kms-key |
--gce-pd-kms-key-keyring | --gce-pd-kms-key-location |
--gce-pd-kms-key-project | --help |
--identity-config-file | --image | --image-version |
--initialization-action-timeout |
--initialization-actions | --kerberos-config-file |
--kerberos-kms-key | --kerberos-kms-key-keyring |
--kerberos-kms-key-location |
--kerberos-kms-key-project |
--kerberos-root-principal-password-uri | --kms-key |
--kms-keyring | --kms-location | --kms-project |
--labels | --master-accelerator |
--master-boot-disk-size | --master-boot-disk-type |
--master-local-ssd-interface | --master-machine-type |
--master-min-cpu-platform | --max-age | --max-idle |
--metadata | --metric-overrides |
--metric-overrides-file | --metric-sources |
--min-num-workers | --min-secondary-worker-fraction |
--network | --node-group |
--num-driver-pool-local-ssds |
--num-master-local-ssds | --num-masters |
--num-secondary-worker-local-ssds |
```

```
--num-secondary-workers | --num-worker-local-ssds |
--num-workers | --optional-components |
--private-ipv6-google-access-type | --properties |
--public-ip-address | --region | --reservation |
--reservation-affinity | --scopes |
--secondary-worker-accelerator |
--secondary-worker-boot-disk-size |
--secondary-worker-boot-disk-type |
--secondary-worker-local-ssd-interface |
--secondary-worker-machine-types |
--secondary-worker-type |
--secure-multi-tenancy-user-mapping |
--service-account | --shielded-integrity-monitoring |
--shielded-secure-boot | --shielded-vtpm |
--single-node | --subnet | --tags | --temp-bucket |
--worker-accelerator | --worker-boot-disk-size |
--worker-boot-disk-type |
--worker-local-ssd-interface | --worker-machine-type |
--worker-min-cpu-platform | --zone
```

For detailed information on this command and its flags, run:

```
gcloud dataproc clusters create --help
CPU times: user 3 µs, sys: 1e+03 ns, total: 4 µs
Wall time: 6.68 µs
```

1 # Description of the cluster

2 !gcloud dataproc clusters describe \$CLUSTER

```
ERROR: (gcloud.dataproc.clusters.describe) argument (CLUSTER : --region=REGION): Must be specified.
Usage: gcloud dataproc clusters describe (CLUSTER : --region=REGION) [optional flags]
    optional flags may be --help | --region
```

For detailed information on this command and its flags, run:

```
gcloud dataproc clusters describe --help
```

1 # Submit the job

2 !gcloud dataproc jobs submit pyspark --cluster \$CLUSTER \task_1di.py

```
ERROR: (gcloud.dataproc.jobs.submit.pyspark) argument PY_FILE: Must be specified.
Usage: gcloud dataproc jobs submit pyspark PY_FILE (--cluster=CLUSTER | --cluster-labels=[KEY=VALUE,...]) [optional flags] [-- JOB_ARGS
    optional flags may be --archives | --async | --bucket | --cluster |
        --cluster-labels | --driver-log-levels |
        --driver-required-memory-mb |
        --driver-required-vcores | --files | --help | --jars |
        --labels | --max-failures-per-hour |
        --max-failures-total | --properties |
        --properties-file | --py-files | --region
```

For detailed information on this command and its flags, run:

```
gcloud dataproc jobs submit pyspark --help
```

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-central1

```
ERROR: (gcloud.dataproc.clusters.delete) Error parsing [cluster].
The [cluster] resource is not properly specified.
Failed to find attribute [project]. The attribute can be set in the following ways:
- provide the argument `cluster` on the command line with a fully specified name
- provide the argument `--project` on the command line
- set the property `core/project`
```

```

1 # 1 machine with eightfold resources, need to change the disk size to 100 after running.
2 import time
3 start_time = time.time()
4
5 !gcloud dataproc clusters create $CLUSTER \
6   --bucket $PROJECT-storage \
7   --image-version 1.4-ubuntu18 \
8   --master-machine-type n1-standard-8 \
9   --master-boot-disk-type pd-ssd --master-boot-disk-size 200 \
10  --num-workers 0 \
11  --max-idle 3600s \
12  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
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-clusters/init-actions
23 # https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/cluster-properties

```

ERROR: (gcloud.dataproc.clusters.create) argument --bucket: expected one argument
Usage: gcloud dataproc clusters create (CLUSTER : --region=REGION) [optional flags]
optional flags may be --action-on-failed-primary-workers | --no-address |
--async | --autoscaling-policy | --bucket |
--confidential-compute | --dataproc-metastore |
--driver-pool-accelerator |
--driver-pool-boot-disk-size |
--driver-pool-boot-disk-type | --driver-pool-id |
--driver-pool-local-ssd-interface |
--driver-pool-machine-type |
--driver-pool-min-cpu-platform | --driver-pool-size |
--enable-component-gateway | --enable-kerberos |
--expiration-time | --gce-pd-kms-key |
--gce-pd-kms-key-keyring | --gce-pd-kms-key-location |
--gce-pd-kms-key-project | --help |
--identity-config-file | --image | --image-version |
--initialization-action-timeout |
--initialization-actions | --kerberos-config-file |
--kerberos-kms-key | --kerberos-kms-key-keyring |
--kerberos-kms-key-location |
--kerberos-kms-key-project |
--kerberos-root-principal-password-uri | --kms-key |
--kms-keyring | --kms-location | --kms-project |
--labels | --master-accelerator |
--master-boot-disk-size | --master-boot-disk-type |
--master-local-ssd-interface | --master-machine-type |
--master-min-cpu-platform | --max-age | --max-idle |
--metadata | --metric-overrides |
--metric-overrides-file | --metric-sources |
--min-num-workers | --min-secondary-worker-fraction |
--network | --node-group |
--num-driver-pool-local-ssds |
--num-master-local-ssds | --num-masters |
--num-secondary-worker-local-ssds |
--num-secondary-workers | --num-worker-local-ssds |
--num-workers | --optional-components |
--private-ipv6-google-access-type | --properties |
--public-ip-address | --region | --reservation |
--reservation-affinity | --scopes |
--secondary-worker-accelerator |
--secondary-worker-boot-disk-size |
--secondary-worker-boot-disk-type |
--secondary-worker-local-ssd-interface |
--secondary-worker-machine-types |
--secondary-worker-type |
--secure-multi-tenancy-user-mapping |
--service-account | --shielded-integrity-monitoring |
--shielded-secure-boot | --shielded-vtpm |
--single-node | --subnet | --tags | --temp-bucket |
--worker-accelerator | --worker-boot-disk-size |
--worker-boot-disk-type |
--worker-local-ssd-interface | --worker-machine-type |
--worker-min-cpu-platform | --zone

For detailed information on this command and its flags, run:

```
gcloud dataproc clusters create --help
```

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

Wall time: 6.68 μs

```

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

ERROR: (gcloud.dataproc.clusters.describe) argument (CLUSTER : --region=REGION): Must be specified.
Usage: gcloud dataproc clusters describe (CLUSTER : --region=REGION) [optional flags]
optional flags may be --help | --region

For detailed information on this command and its flags, run:
gcloud dataproc clusters describe --help

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

ERROR: (gcloud.dataproc.jobs.submit.pyspark) argument PY_FILE: Must be specified.
Usage: gcloud dataproc jobs submit pyspark PY_FILE (--cluster=CLUSTER | --cluster-labels=[KEY=VALUE,...]) [optional flags] [-- JOB_ARGS
optional flags may be --archives | --async | --bucket | --cluster |
--cluster-labels | --driver-log-levels |
--driver-required-memory-mb |
--driver-required-vcores | --files | --help | --jars |
--labels | --max-failures-per-hour |
--max-failures-total | --properties |
--properties-file | --py-files | --region

For detailed information on this command and its flags, run:
gcloud dataproc jobs submit pyspark --help

```

▼ 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 multiple 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.

```

1 ### Checked ####
2 # Here are some useful values for testing your code, use higher values later for actually testing throughput
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 output without showing it
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).

```

1 ### Checked ###
2
3 def decode_jpeg_and_label(filepath):
4     # extracts the image data and creates a class label, based on the filepath
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 image
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]), # if true
24                     lambda: tf.image.resize(image, [w*th/h, h*th/h])) # if false
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 - th) // 2, tw, th)
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, repetitions)
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, repetitions)
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 name='/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 name='/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.TextIOWrapper 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'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding
Image batch shape (4,), ["b'dandelion'", "b'sunflowers'", "b'sunflowers'", "b'sunflowers'"]) <_io.TextIOWrapper name='/dev/null' mode
Image batch shape (4,), ["b'roses'", "b'tulips'", "b'dandelion'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encod
bs: 4, ds: 6, rep: 1
Image batch shape (4,), ["b'dandelion'", "b'sunflowers'", "b'daisy'", "b'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' enco
Image batch shape (4,), ["b'tulips'", "b'dandelion'", "b'daisy'", "b'daisy'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding=
Image batch shape (4,), ["b'tulips'", "b'tulips'", "b'dandelion'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' mode= 'w' encoding=
Image batch shape (4,), ["b'daisy'", "b'tulips'", "b'daisy'", "b'roses'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF
Image batch shape (4,), ["b'daisy'", "b'daisy'", "b'tulips'", "b'roses'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF
Image batch shape (4,), ["b'roses'", "b'roses'", "b'dandelion'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding
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'>
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'>
```

▼ 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.

```

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 = bytestring (not text string)
13         "class": tf.io.FixedLenFeature([], tf.int64) #,   # shape [] means scalar
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_deterministic = False
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_speed, throughput])
56     return tt
57
58 # Extra resources checked:
59 # https://www.tensorflow.org/io/api_docs/python/tfio/IODataset

```

```

1 # i) combine the previous cells to have the code to create a dataset and create a list of parameter combinations in an RDD (2%)
2 # ii) get a Spark context and create the dataset and run timing test for each combination in parallel (2%)
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", "reading_speed", "throughput"]
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_speed, throughput])
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

```

```

1 # iii) transform the resulting RDD to the structure ( parameter_combination, images_per_second ) and save these values in an array (2%)
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-transformations-actions/
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.pdf (This book in particular has been really helpful)
18 # https://sparkbyexamples.com/pyspark/pyspark-convert-dataframe-to-rdd/
19
20 TFRecords Array RDD
21 [(6, 6, 1, 36, 0.2102503776550293, 171.22442490480285),
22  (6, 6, 2, 36, 0.17751073837280273, 202.80463215917607),
23  (6, 6, 2, 36, 0.18248200416564941, 197.279727196117),
24  (6, 6, 3, 36, 0.1756117343902588, 204.9976906440537),
25  (6, 6, 3, 36, 0.17471814155578613, 206.04614769500327),
26  (6, 6, 3, 36, 0.20185422897338867, 178.34652354371053)]
27
28
29 # RDD for image file parameter_combination
30
31 image_array_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_sizes']),
32                                              int(x['batch_nums']),
33                                              int(x['repetitions']),
34                                              int(x['datasetsize']),
35                                              float(x['reading_speed']),
36                                              float(x['throughput'])))
37
38 print("Image Array RDD")
39 image_array_rdd.take(6)
40
41 Image Array RDD
42 [(6, 6, 1, 36, 0.42624759674072266, 84.45795419205152),
43  (6, 6, 2, 36, 0.4570634365081787, 78.76368382259739),
44  (6, 6, 2, 36, 0.5044350624084473, 71.36696610286451),
45  (6, 6, 3, 36, 0.5790822505950928, 62.167334541862864),
46  (6, 6, 3, 36, 0.525456428527832, 68.5118652004334),
47  (6, 6, 3, 36, 0.5065336227416992, 71.07129395506638)]

```

```
1 # iv) Create an RDD with all results for each parameter as (parameter_value,images_per_second) and collect the result for each parameter
2
3 # Create RDD for image files batch size and their corresponding speeds
4 image_batchSizes_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_sizes']), float(x['throughput']))))
5 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect batch size results
6
7 # Create RDD for image files batch numbers and their corresponding speeds
8 image_batchNums_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['batch_nums']), float(x['throughput']))))
9 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect batch number results
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['repetitions']), float(x['throughput']))))
13 image_repetitions_speed = image_repetitions_speed_rdd.collect() # collect repetitions results
14
15 # Create RDD for image files data set size and their corresponding speeds
16 image_dsSize_speed_rdd = image_files_df.rdd.map(lambda x: (int(x['datasetsize']), float(x['throughput']))))
17 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect datasetsize results
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_sizes']), float(x['throughput']))))
25 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batch size results
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_nums']), float(x['throughput']))))
29 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batch number results
30
31 # Create RDD for repetitions and their corresponding speeds
32 tfr_repetitions_speed_rdd = tfr_files_df.rdd.map(lambda x: (int(x['repetitions']), float(x['throughput']))))
33 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect repetitions results
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['datasetsize']), float(x['throughput']))))
37 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect datasetsize results
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/ch04.html (best resource)
44 # https://stackoverflow.com/questions/73014570/best-way-to-join-2-pair-rdds-with-exactly-the-same-keys-and-where-all-keys-are-u
45 # https://towardsdatascience.com/the-ultimate-guide-to-functional-programming-for-big-data-1e57b0d225a3
46 # https://spark.apache.org/docs/latest/rdd-programming-guide.html#basics
47 # https://www.geeksforgeeks.org/pyspark-map-transformation/
48 # https://www.analyticsvidhya.com/blog/2016/10/using-pyspark-to-perform-transformations-and-actions-on-rdd/
```

```
1 # v) create an RDD with the average reading speeds for each parameter value and collect the results. Keep associativity in mind when impl
2
3 # Create RDD for batch sizes and their corresponding avg speed
4 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
5 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect batch size results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
9 tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect batch number results
10
11 # Create RDD for tfr repetitions and their corresponding avg speed
12 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
13 tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # collect repetitions results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
17 tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect datasetsize results
18
19
20
21
22 # Create RDD for image files batch size and their corresponding avg speed
23 image_batchSizes_avgSpeed_rdd = image_batchSizes_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
24 image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() # collect batch size results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
28 image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # collect batch number results
29
30 # Create RDD for image repetitions and their corresponding avg speed
31 image_repetitions_avgSpeed_rdd = image_repetitions_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
32 image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() # collect repetitions results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
36 image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect datasetsize results
37
38
39 # Some extra resources checked:
40 # https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.RDD.mapValues.html
41 #
```

```
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("gstutil 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("gstutil 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-parameter
```

```
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]...\\n/ [0 files][ 0.0 B/ 10.6 KiB]
```

```

1 # vii) Write your code it into a file using the cell magic %%writefile spark_job.py (1%)
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 files
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 file names in my bucket
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 = bytestring (not text string)
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 filepath
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]), # if true
71                     lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
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 - th) // 2, tw, th)
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_deterministic = False
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_speed, throughput])
113    return tt
114
115 # i) combine the previous cells to have the code to create a dataset and create a list of parameter combinations in an RDD (2%)
116 # ii) get a Spark context and create the dataset and run timing test for each combination in parallel (2%)
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", "reading_speed", "throughput"]
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_speed, throughput])
173 return tt
174
175 # Get spark context for RDDs
176 sc = pyspark.SparkContext.getOrCreate()
177
178 # Parallelize data into RDD
179 image_files_rdd = sc.parallelize(parameter_list)
180
181 # Initialize spark session for df conversion
182 image_files_sparkSession = SparkSession(sc)
183
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, images_per_second ) and save these values in an array (2%)
191 tfr_array_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']), int(x['batch_nums']), int(x['repetitions']), int(x['datasetsize']), float(x['reading_speed']), float(x['throughput'])))
192
193
194
195
196
197 tfr_array_rdd.take(6)
198
199
200
201
202 # iv) create an RDD with all results for each parameter as (parameter_value,images_per_second) and collect the result for each parameter
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_sizes']), float(x['throughput'])))
206 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batch size results
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_nums']), float(x['throughput'])))
210 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batch number results
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']), float(x['throughput'])))
214 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect repetitions results
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']), float(x['throughput'])))
218 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect datasetsize results
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_sizes']), float(x['throughput'])))
222 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect batch size results
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_nums']), float(x['throughput'])))
226 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect batch number results
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']), float(x['throughput'])))
230 image_repetitions_speed = image_repetitions_speed_rdd.collect() # collect repetitions results
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']), float(x['throughput']))))
234 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect datasetsize results
235
236
237
238
239 # v) create an RDD with the average reading speeds for each parameter value and collect the results. Keep associativity in mind when impl
240
241 # Create RDD for batch sizes and their corresponding avg speed
242 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
243 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect batch size results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues
247 tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect batch number results
248
249 # Create RDD for tfr repetitions and their corresponding avg speed
250 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapV
251 tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # collect repetitions results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lamk
255 tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect datasetsize results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).ma
259 image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() # collect batch size results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapv
263 image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # collect batch number results
264
265 # Create RDD for image repetitions and their corresponding avg speed
266 image_repetitions_avgSpeed_rdd = image_repetitions_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).ma
267 image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() # collect repetitions results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(
271 image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect datasetsize results
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.PIPE)
285     print("gstui 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("gstuil 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.

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

False
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]
<Figure size 640x480 with 0 Axes>
```

```

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 files
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 output file names in my bucket
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 = bytestring (not text string)
35         "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] means scalar
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 filepath
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 image
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]), # if true
65                     lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
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 - th) // 2, tw, th)
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 = False
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_speed, throughput])
107    return tt
108
109 # i) combine the previous cells to have the code to create a dataset and create a list of parameter combinations in an RDD (2%)
110 # ii) get a Spark context and create the dataset and run timing test for each combination in parallel (2%)
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", "reading_speed", "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_speed, throughput])
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, images_per_second ) and save these values in an array (2%)
185 tfr_array_rdd = tfr_files_df.rdd.map(lambda x: (int(x['batch_sizes']),int(x['batch_nums']),int(x['repetitions']),
186                                         int(x['datasetsize']),float(x['reading_speed']), float(x['throughput'])))
187 tfr_array_rdd.take(6)
188
189
190 # iv) create an RDD with all results for each parameter as (parameter_value,images_per_second) and collect the result for each parameter
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['throughput'])))
193 tfr_batchSize_speed = tfr_batchSize_speed_rdd.collect() # collect batch size results
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['throughput'])))
197 tfr_batchNum_speed = tfr_batchNum_speed_rdd.collect() # collect batch number results
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['throughput'])))
201 tfr_repetitions_speed = tfr_repetitions_speed_rdd.collect() # collect repetitions results
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['datasetsize']), float(x['throughput'])))
205 tfr_dsSize_speed = tfr_dsSize_speed_rdd.collect() # collect datasetsize results
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['throughput'])))
209 image_batchSizes_speed = image_batchSizes_speed_rdd.collect() # collect batch size results
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['throughput'])))
213 image_batchNums_speed = image_batchNums_speed_rdd.collect() # collect batch number results
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['throughput'])))
217 image_repetitions_speed = image_repetitions_speed_rdd.collect() # collect repetitions results
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['datasetsize']), float(x['throughput'])))
221 image_dsSize_speed = image_dsSize_speed_rdd.collect() # collect datasetsize results
222
223
224 # v) create an RDD with the average reading speeds for each parameter value and collect the results. Keep associativity in mind when impl
225
226 # Create RDD for batch sizes and their corresponding avg speed
227 tfr_batchSize_avgSpeed_rdd = tfr_batchSize_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValue
228 tfr_batchSize_avgSpeed = tfr_batchSize_avgSpeed_rdd.collect() # collect batch size results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues

```

```

232 tfr_batchNums_avgSpeed = tfr_batchNums_avgSpeed_rdd.collect() # collect batch number results
233
234 # Create RDD for tfr repetitions and their corresponding avg speed
235 tfr_repetitions_avgSpeed_rdd = tfr_repetitions_speed_rdd.mapValues(lambda x: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lambda x: x[0])
236 tfr_repetitions_avgSpeed = tfr_repetitions_avgSpeed_rdd.collect() # collect repetitions results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lambda x: x[0])
240 tfr_dsSize_avgSpeed = tfr_dsSize_avgSpeed_rdd.collect() # collect datasetsize results
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: (x, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lambda x: x[0])
244 image_batchSizes_avgSpeed = image_batchSizes_avgSpeed_rdd.collect() # collect batch size results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lambda x: x[0])
248 image_batchNums_avgSpeed = image_batchNums_avgSpeed_rdd.collect() # collect batch number results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lambda x: x[0])
252 image_repetitions_avgSpeed = image_repetitions_avgSpeed_rdd.collect() # collect repetitions results
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, 1)).reduceByKey(lambda x,y: (x[0]+y[0], x[1]+y[1])).mapValues(lambda x: x[0])
256 image_dsSize_avgSpeed = image_dsSize_avgSpeed_rdd.collect() # collect datasetsize results
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 a little. These images will be useful in the next task. Again, don't use the SCREENSHOT function that Google provides, but just take a picture of the graphs you see for the VMs.

```

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

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].
```

```

1 ### CODING TASK ###
2 # Set up a cluster with a single machine using the maximal SSD size (100) and 1 machine with eightfold resources.
3
4 import time
5 start_time = time.time()
6 #REGION = europe-west3
7 !gcloud dataproc clusters create $CLUSTER \
8     --bucket $PROJECT-storage \
9     --image-version 1.4-ubuntu18 \
10    --master-machine-type n1-standard-8 \
11    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
12    --num-workers 0\
13    --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
14    --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 bucket. 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

CPU times: user 4 µs, sys: 1 µs, total: 5 µs
Wall time: 8.34 µs
Execution time: 153.5760896205902 seconds

```

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
    serviceAccountScopes:
      - https://www.googleapis.com/auth/bigquery
      - https://www.googleapis.com/auth/bigtable.admin.table
      - https://www.googleapis.com/auth/bigtable.data
      - https://www.googleapis.com/auth/cloud.useraccounts.readonly
      - https://www.googleapis.com/auth/devstorage.full\_control
      - https://www.googleapis.com/auth/devstorage.read\_write
      - https://www.googleapis.com/auth/logging.write
    zoneUri: 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
    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
    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.servicercp-address: big-data-cw-420116-cluster-m:8051
mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '4000'

```

```

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'
2024-05-04 11:41:04.458757: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dlerror if you do not have a GPU se
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 BlockManagerMaster
24/05/04 11:41:07 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
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: jetty-9.4.2-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:
24/05/04 11:41:07 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be sc
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:
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/
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]
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/cac7e146
driverOutputResourceUri: gs://big-data-cw-420116-storage/google-cloud-dataproc-metainfo/35f0437b-4971-4385-b85a-3808433b9719/jobs/cac7e1
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/cac7e1460f
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/

```

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**

```

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 files
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 output file names in my bucket
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 = bytestring (not text string)
39         "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] means scalar
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 filepath
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 image
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]), # if true
69                     lambda: tf.image.resize(image, [w*th/h, h*th/h])) # if false
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 - th) // 2, tw, th)
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 = False
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]
```