

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 Gerner. 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.

```
In [ ]: import os, sys, math
import numpy as np
import scipy as sp
import scipy.stats
import time
import datetime
import string
import random
from matplotlib import pyplot as plt
import tensorflow as tf
print("Tensorflow version " + tf.__version__)
import pickle
```

Tensorflow version 2.15.0

Cloud and Drive authentication

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

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

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

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

```
Mounting google drive...
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", for
ce_remount=True).
/content/drive/MyDrive
mkdir: cannot create directory 'BD-CW': File exists
/content/drive/MyDrive/BD-CW
```

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

```
In [ ]: import sys
if 'google.colab' in sys.modules:
    from google.colab import auth
    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.

```
In [ ]: PROJECT = 'daring-healer-421511'  ### USE YOUR GOOGLE CLOUD PROJECT ID HERE. ###
!gcloud config set project $PROJECT
REGION = 'us-west1'
CLUSTER = '{}-cluster'.format(PROJECT)
```

```
!gcloud config set compute/region $REGION
!gcloud config set dataproc/region $REGION

!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-west1
[core]
account = Nicholas.Tsioras@city.ac.uk
project = daring-healer-421511
[dataproc]
region = us-west1
```

Your active configuration is: [default]

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

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

Creating gs://daring-healer-421511-storage/...
ServiceException: 409 A Cloud Storage bucket named 'daring-healer-421511-storage' already exists. Try another name. Bucket names must be globally unique across all Google Cloud projects, including those outside of your organization.

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

```
In [ ]: #@title Utility functions for image display **[RUN THIS TO ACTIVATE]** { display-mode: "form" }
def display_9_images_from_dataset(dataset):
    plt.figure(figsize=(13,13))
    subplot=331
    for i, (image, label) in enumerate(dataset):
        plt.subplot(subplot)
        plt.axis('off')
        plt.imshow(image.numpy().astype(np.uint8))
        plt.title(str(label.numpy()), fontsize=16)
        # plt.title(label.numpy().decode(), fontsize=16)
        subplot += 1
        if i==8:
            break
    plt.tight_layout()
    plt.subplots_adjust(wspace=0.1, hspace=0.1)
    plt.show()

def display_training_curves(training, validation, title, subplot):
    if subplot%10==1: # set up the subplots on the first call
        plt.subplots(figsize=(10,10), facecolor='#F0F0F0')
        plt.tight_layout()
        ax = plt.subplot(subplot)
        ax.set_facecolor('#F8F8F8')
        ax.plot(training)
        ax.plot(validation)
        ax.set_title('model ' + title)
        ax.set_ylabel(title)
        ax.set_xlabel('epoch')
        ax.legend(['train', 'valid.'])

def dataset_to_numpy_util(dataset, N):
    dataset = dataset.batch(N)
    for images, labels in dataset:
        numpy_images = images.numpy()
        numpy_labels = labels.numpy()
        break;
    return numpy_images, numpy_labels

def title_from_label_and_target(label, correct_label):
    correct = (label == correct_label)
    return "{} [{}{}{}]" .format(CLASSES[label], str(correct), ', should be ' if not correct else '',
                                CLASSES[correct_label] if not correct else ''), correct

def display_one_flower(image, title, subplot, red=False):
    plt.subplot(subplot)
    plt.axis('off')
    plt.imshow(image)
    plt.title(title, fontsize=16, color='red' if red else 'black')
    return subplot+1

def display_9_images_with_predictions(images, predictions, labels):
    subplot=331
    plt.figure(figsize=(13,13))
```

```

classes = np.argmax(predictions, axis=-1)
for i, image in enumerate(images):
    title, correct = title_from_label_and_target(classes[i], labels[i])
    subplot = display_one_flower(image, title, subplot, not correct)
    if i >= 8:
        break;

plt.tight_layout()
plt.subplots_adjust(wspace=0.1, hspace=0.1)
plt.show()

```

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 **Dataprocc** is still required for the final version.

```

In [ ]: %cd
!apt-get update -qq
!apt-get install openjdk-8-jdk-headless -qq >> /dev/null # send any output to null device
!tar -xzf "/content/drive/My Drive/Big_Data/data/spark/spark-3.5.0-bin-hadoop3.tgz" # unpack

!pip install -q findspark
import os
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
os.environ["SPARK_HOME"] = "/root/spark-3.5.0-bin-hadoop3"
import findspark
findspark.init()
import pyspark
print(pyspark.__version__)
sc = pyspark.SparkContext.getOrCreate()
print(sc)

/root
3.5.0
<SparkContext master=local[*] appName=pyspark-shell>

```

Section 1: Data pre-processing

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

1.1 Getting started

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

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

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

```

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

```

We **read the image files** from the public GCS bucket that contains the *Flowers* dataset. **TensorFlow** has **functions** to execute glob patterns that we use to calculate the the number of images in total and per partition (rounded up as we cannot deal with parts of images).

```

In [ ]: nb_images = len(tf.io.gfile.glob(GCS_PATTERN)) # number of images
partition_size = math.ceil(1.0 * nb_images / PARTITIONS) # images per partition (float)
print("GCS_PATTERN matches {} images, to be divided into {} partitions with up to {} images each.".format(nb_im

GCS_PATTERN matches 3670 images, to be divided into 16 partitions with up to 230 images each.

```

Map functions

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

```

In [ ]: def decode_jpeg_and_label(filepath):
          # extracts the image data and creates a class label, based on the filepath
          bits = tf.io.read_file(filepath)
          image = tf.image.decode_jpeg(bits)
          # parse flower name from containing directory

```

```

label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
label2 = label.values[-2]
return image, label2

def resize_and_crop_image(image, label):
    # Resizes and cropd using "fill" algorithm:
    # always make sure the resulting image is cut out from the source image
    # so that it fills the TARGET_SIZE entirely with no black bars
    # and a preserved aspect ratio.
    w = tf.shape(image)[0]
    h = tf.shape(image)[1]
    tw = TARGET_SIZE[1]
    th = TARGET_SIZE[0]
    resize_crit = (w * th) / (h * tw)
    image = tf.cond(resize_crit < 1,
                    lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), # if true
                    lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
                    )
    nw = tf.shape(image)[0]
    nh = tf.shape(image)[1]
    image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh - th) // 2, tw, th)
    return image, label

def recompress_image(image, label):
    # this reduces the amount of data, but takes some time
    image = tf.cast(image, tf.uint8)
    image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampling=False)
    return image, label

```

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

```

In [ ]: dsetFiles = tf.data.Dataset.list_files(GCS_PATTERN) # This also shuffles the images
dsetDecoded = dsetFiles.map(decode_jpeg_and_label)
dsetResized = dsetDecoded.map(resize_and_crop_image)

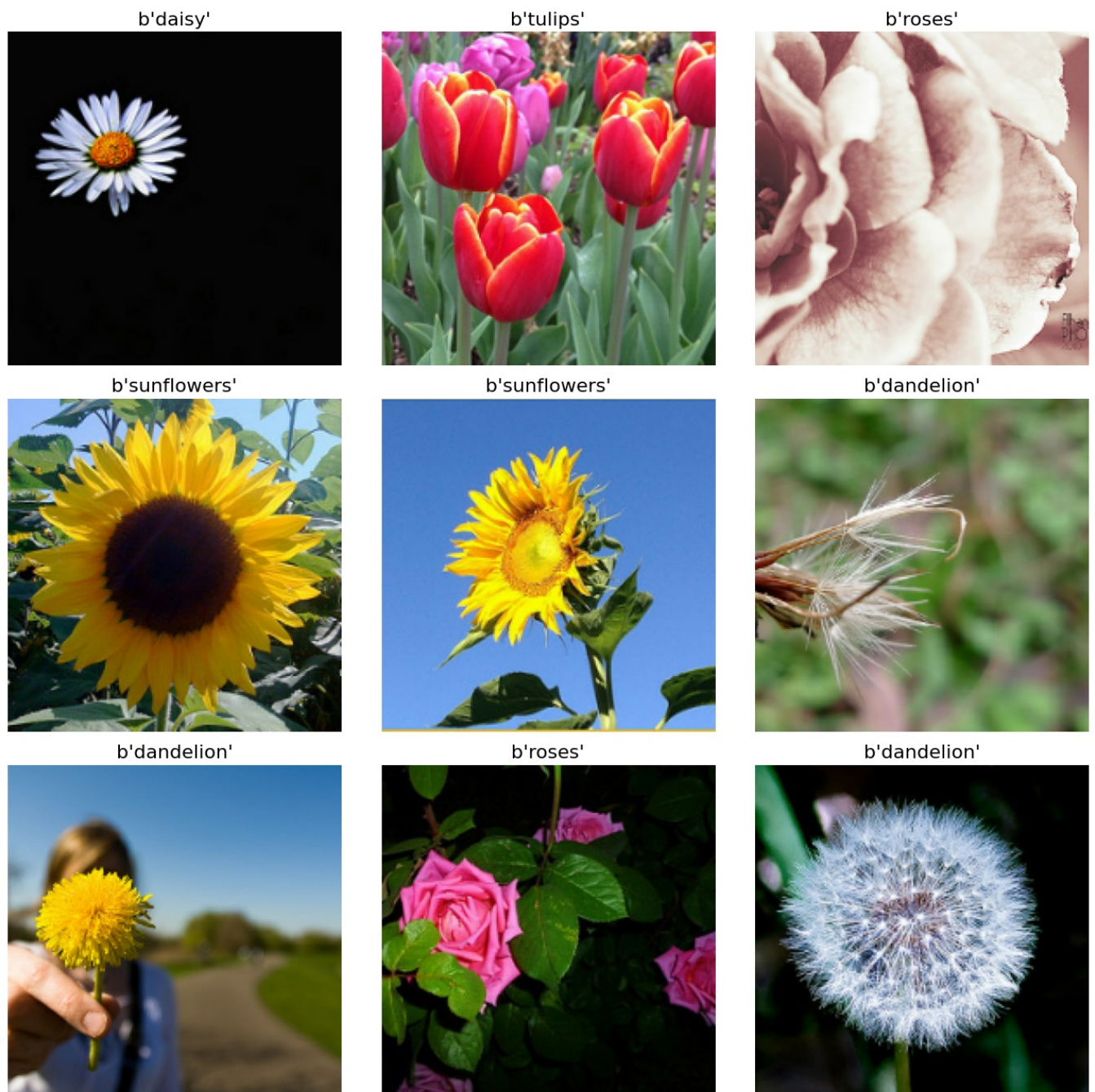
```

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

```

In [ ]: display_9_images_from_dataset(dsetResized)

```

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

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

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

1.2 Improving Speed

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

speed.

Recompress the images

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

```
In [ ]: # This is a quick test to get an idea how long recompressions takes.
dataset4 = dsetResized.map(recompress_image)
test_set = dataset4.batch(10).take(10)
for image, label in test_set:
    print("Image batch shape {}, {}".format(image.numpy().shape, [lbl.decode('utf8') for lbl in label.numpy()])

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

Write the dataset to TFRecord files

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

```
In [ ]: 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.

```
In [ ]: # functions for writing TFRecord entries
# Feature values are always stored as lists, a single data element will be a list of size 1
def _bytestring_feature(list_of_bytestrings):
    return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))

def _int_feature(list_of_ints): # int64
    return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))

def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data records
    class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
    one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
    feature = {
        "image": _bytestring_feature([img_bytes]), # one image in the list
        "class": _int_feature([class_num]) #, # one class in the list
    }
    return tf.train.Example(features=tf.train.Features(feature=feature))

def write_tfrecords(GCS_PATTERN, GCS_OUTPUT, partition_size): # write the images to files.
    print("Writing TFRecords")
    tt0 = time.time()
    filenames = tf.data.Dataset.list_files(GCS_PATTERN)
    dataset1 = filenames.map(decode_jpeg_and_label)
    dataset2 = dataset1.map(resize_and_crop_image)
    dataset3 = dataset2.map(recompress_image)
    dataset4 = dataset3.batch(partition_size) # partitioning: there will be one "batch" of images per file
    for partition, (image, label) in enumerate(dataset4):
        # batch size used as partition size here
        partition_size = image.numpy().shape[0]
        # good practice to have the number of records in the filename
        filename = GCS_OUTPUT + "{:02d}-{}.tfrec".format(partition, partition_size)
        # You need to change GCS_OUTPUT to your own bucket to actually create new files
        with tf.io.TFRecordWriter(filename) as out_file:
```

```

    for i in range(partition_size):
        example = to_tfrecord(out_file,
                               image.numpy()[i], # re-compressed image: already a byte string
                               label.numpy()[i] #
                            )
        out_file.write(example.SerializeToString())
    print("Wrote file {} containing {} records".format(filename, partition_size))
    print("Total time: "+str(time.time()-tt0))

```

```
write_tfrecords(GCS_PATTERN,GCS_OUTPUT,partition_size) # uncomment to run this cell
```

Writing TFRecords

```

Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers00-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers01-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers02-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers03-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers04-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers05-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers06-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers07-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers08-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers09-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers10-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers11-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers12-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers13-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers14-230.tfrec containing 230 records
Wrote file gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers15-220.tfrec containing 220 records
Total time: 137.7730758190155

```

Test the TFRecord files

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

```

In [ ]: #GCS_OUTPUT = 'gs://flowers-public/tfrecords-jpeg-192x192-2/'
        # remove the line above to use your own files that you generated above

def read_tfrecord(example):
    features = {
        "image": tf.io.FixedLenFeature([], tf.string), # tf.string = bytestring (not text string)
        "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] means scalar
    }
    # decode the TFRecord
    example = tf.io.parse_single_example(example, features)
    image = tf.image.decode_jpeg(example['image'], channels=3)
    image = tf.reshape(image, [*TARGET_SIZE, 3])
    class_num = example['class']
    return image, class_num

def load_dataset(filenamees):
    # read from TFRecords. For optimal performance, read from multiple
    # TFRecord files at once and set the option experimental_deterministic = False
    # to allow order-altering optimizations.
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False

    dataset = tf.data.TFRecordDataset(filenamees)
    dataset = dataset.with_options(option_no_order)
    dataset = dataset.map(read_tfrecord)
    return dataset

filenamees = tf.io.gfile.glob(GCS_OUTPUT + "/*.tfrec")
datasetTfrec = load_dataset(filenamees)

```

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

```

In [ ]: batched_dataset = datasetTfrec.batch(10)
        sample_set = batched_dataset.take(10)

```



```
for image, label in sample_set:
    print("Image batch shape {}, {}".format(image.numpy().shape, \
                                              [str(lbl) for lbl in label.numpy()])))
```

```
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
Image batch shape (10, 192, 192, 3), ['0', '0', '0', '0', '0', '0', '0', '0', '0', '0'])
```

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:

- Copy** over the **mapping functions** (see section 1.1) and **adapt** the resizing and recompression functions **to Spark** (only one argument). (3%)
- Replace** the TensorFlow **Dataset objects with RDDs**, starting with an RDD that contains the list of image filenames. (3%)
- 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%)
- Then **use the functions from above** to write the TFRecord files. (3%)
- 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%)

```
In [ ]: #We first import the necessary libraries
import os, sys, math
import numpy as np
import scipy as sp
import scipy.stats
import time
import string
import datetime
import random
from matplotlib import pyplot as plt
import tensorflow as tf
print("Tensorflow version " + tf.__version__)
import pickle
import pyspark
from pyspark.sql import SQLContext
from pyspark.sql import Row

#We define the variables that will be needed for this project.
GCS_PATTERN = 'gs://flowers-public/*/*.jpg' #glob pattern for input files.
PROJECT = 'daring-healer-421511'
BUCKET = 'gs://{}-storage'.format(PROJECT) # This is the bucket storage.
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # We prefix for output file names.
PARTITIONS = 16 # This is the number of partitions that will be used later.

#Task 1a) - i) Copy over the mapping functions and adapt the resizing and recompression functions to Spark (only one argument)
def decode_jpeg_and_label(filepath):
    #This function reads and decodes the JPEG image from the file.
    bits = tf.io.read_file(filepath)
    image = tf.image.decode_jpeg(bits) #We decode the read data as JPEG images.
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/') #We extract the label (flower name) of the file
    label2 = label.values[-2]
    return image, label2

def resize_and_crop_image(data): #We specify one argument
    image, label = data #The input is basically a tuple containing the image data and the label.
    #We determine the dimensions of the input images and the target dimensions.
    w = tf.shape(image)[0]
    h = tf.shape(image)[1]
    tw = TARGET_SIZE[1] #Target Width
    th = TARGET_SIZE[0] #Target Height
    #Calculate resize criterion to maintain the aspect ratio.
```

```

resize_crit = (w * th) / (h * tw)
#We resize the images to fill the target dimensions but also preserving the aspect ratio.
image = tf.cond(resize_crit < 1,
                 lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), # if true
                 lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
                )
nw = tf.shape(image)[0]
nh = tf.shape(image)[1]
image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh - th) // 2, tw, th)
return (image, label)

def recompress_image(data):
    image, label = data #The input is basically a tuple containing the image data and the label.
    #This function is used to reduce the amount of data, but takes some time
    image = tf.cast(image, tf.uint8) #The image is transformed to a uint8 so the images can be compressed to JP
    image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampling=False) #We compress the images
    return (image, label)

```

Tensorflow version 2.15.0

```

In [ ]: #Task 1a) - ii) Replace the TensorFlow Dataset objects with RDDs, starting with an RDD that contains the list
from pyspark.sql import SparkSession
sc = pyspark.SparkContext.getOrCreate() #Create or get the existing Spark context
spark = SparkSession.builder.getOrCreate() #Initialize the Spark session

image_paths = tf.io.gfile.glob(GCS_PATTERN) #We retrieve the list of image filenames
image_rdd = sc.parallelize(image_paths) #We create an RDD from the list of image filenames

```

```

In [ ]: #Task 1a) - iii) Sample the RDD to a smaller number at an appropriate position in the code. Specify a sampling
rdd1_sample = image_rdd.sample(False, 0.02) #We first sample the RDD with a sampling factor of 0.02
rdd2_decode_jpeg_and_label = image_rdd.map(decode_jpeg_and_label) #We decode the jpeg and label of the images o
rdd3_resize_and_crop_image = rdd2_decode_jpeg_and_label.map(resize_and_crop_image) #We resize and crop the imag
rdd4_recompress_image = rdd3_resize_and_crop_image.map(recompress_image) #We recompress the cropped images of t

```

```

In [ ]: #Task 1a) - iv) Then use the functions from above to write the TFRecord files.
#This function is used to create a TensorFlow feature for storing byte strings. This is for the image data.
def _bytestring_feature(list_of_bytestrings):
    return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))
#This function is used for storing integer values. This is for the labels.
def _int_feature(list_of_ints):
    return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
#This function is used to create a TFRecord from image bytes and a label.
def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data records
    class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
    one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
    feature = {
        "image": _bytestring_feature([img_bytes]), # one image in the list
        "class": _int_feature([class_num]) #, # one class in the list
    }
    return tf.train.Example(features=tf.train.Features(feature=feature))

print("Writing TFRecords")
#This function is used to write TFRecords, which will be applied to partitions of an RDD.
def write_tfrecords(partition_index, partition):
    filename = GCS_OUTPUT + "{}.tfrec".format(partition_index) #We define the filename for this TFRecord file bas
    #We Open a TFRecordWriter object for the filename.
    with tf.io.TFRecordWriter(filename) as out_file:
        for element in partition:
            image=element[0]
            label=element[1]
            example = to_tfrecord(out_file,
                                image.numpy(), # The re-compressed image is already a byte string
                                label.numpy()
                               )
            out_file.write(example.SerializeToString()) #We convert the image and label to TFRecord format and write
    return [filename] #We return the filename of the written TFRecord.

```

Writing TFRecords

```

In [ ]: #Task 1a) - v) The code for writing to the TFRecord files needs to be put into a function, that can be applied
#should return the filename, so that you have a list of the created TFRecord files.

rdd5_partitions = rdd4_recompress_image.repartition(PARTITIONS) #We repartition the RDD to have the number of p
rdd1_filenames = rdd4_recompress_image.mapPartitionsWithIndex(write_tfrecords) #We specify an index for each pa

```

```

In [ ]: rdd1_filenames.take(1) #We check to see the filename of the first TFRecord.

```

```

Out[ ]: ['gs://daring-healer-421511-storage/tfrecords-jpeg-192x192-2/flowers0.tfrec']

```

```

In [ ]: rdd2_decode_jpeg_and_label.take(1) #This is the first element of the RDD2 showing the decoded JPEG na labels of t

```

```

Out[ ]: [(<tf.Tensor: shape=(263, 320, 3), dtype=uint8, numpy=
array([[133, 135, 132],
       [136, 138, 135],
       [140, 142, 139],
       ...,
       [152, 152, 150],
       [155, 155, 153],
       [148, 148, 146]],

       [[133, 135, 132],
       [136, 138, 135],
       [140, 142, 139],
       ...,
       [153, 153, 151],
       [155, 155, 153],
       [147, 147, 145]],

       [[132, 134, 129],
       [135, 137, 134],
       [139, 141, 138],
       ...,
       [152, 152, 150],
       [154, 154, 152],
       [146, 146, 144]],

       ...,

       [[ 44,  48,  25],
       [ 44,  48,  25],
       [ 44,  48,  25],
       ...,
       [127, 126, 122],
       [127, 126, 122],
       [127, 126, 122]],

       [[ 44,  48,  25],
       [ 44,  48,  25],
       [ 44,  48,  25],
       ...,
       [128, 127, 123],
       [128, 127, 123],
       [128, 127, 123]],

       [[ 43,  47,  24],
       [ 43,  47,  24],
       [ 43,  47,  24],
       ...,
       [129, 128, 124],
       [129, 128, 124],
       [130, 129, 125]]], dtype=uint8)>,
  <tf.Tensor: shape=(), dtype=string, numpy=b'daisy'>)]

```

```

In [ ]: rdd3_resize_and_crop_image.take(1) #Take the first element of the RDD3 showing the resized and cropped images.

```



```

x8Eh\xd6\ddr\|x98\|x9e0V\|x1b\|xbe; \|x8e\|xad\}\|xal\|x90\|x99\|h\|xa2\|xb9\|xc0V\|\|x:b0\|xf8\|x0c\|x91\|xf8\|x1dK\|x95Au\|xac\|xc7\|wY\|xaa-s\|xb3p\|xcac\|xb15\|xdb\|x1aZaH\|xa!\|xa2\|xdckAX)\|%\|x9f\|xd1W0\|x86F\|x19U?|\|xcd\|xf3\|x9f\|xd3\|xbf\|xb6\|x96-M.<#|x6a\|xd2*\|xd42\|x88#<(c8\|x9a\|x12Fc0\|x97\|xd4\|x1f\|x\|xc\|xb5\|x9c\|xf1\|xb4\|xa9D\|xec\|xd4\|xe5\|_xb20\|xf1\|xdc\|xe92T\|x10'\|_x84\|x\|f8u\|xf7\|x\|c\|xfz\|x985\|x067\|xd8\|xe1q\|x00\|x85\|xb8\|xcd\|xb\|xf6\|xadU\|xba\|x0b\|xb55j\|4UQ+\|xac'd\|xaeGp->\|x7d\|s\|xc3\|x\|db\|xf4B\|r\|xa4\|x82\|xb\|xa86j\|W+0\|x1b\|xc8A\|xed\|x98\|x81\|x18\|x\|c\|x\|ce\|x80Z{\|x14@\|x02\|xf1\|x\|fe\|xd1\|xecu\|xc3\|x15\|xb6\|xe0\|xb8\|x03\|x05\|xa0Q\|x\|feCB>\|x\|fb\|x\|ba\|x826\|x82\|xe5o\|x9f\|x15v\|xc5G*\|x\|c1\|x82\|x\|af\|xc2r-<\|x8e\|xd9\|x\|fc\|F\|xa9\|xb9\|xa1\|x\|de\|xe0\|xb8\|x8b0\|x16\|x\|fd\|x\|x\|ab\|x\|c6\|x\|fe\|LJJ\|x8b\|x89'\|x93\|x\|ef\|x021\|x\|fd\|x\|cb\|x\|ddm\|x07\|xb2\|xa0\|x05M\|_x93\|x\|f6\|x\|ce\|xe8\|x07\|x8f\|x1a\|x\|cd\|\|xad\|xc5\|x11MU|\|_x3\|xa8\|xabW\|x\|f1\|x9e\|xec51\|;|xdd->\|x17.\|x1b\|x86\|xb1R1J\|x5\|x\|ce\|x1ed\|x\|ff\|x00\|x1a\|\|x\|c\|f\|xd\|f0\|x\|fb\|xaeVw*\|xf1\|x1d\|x03A!\|n\|xae\|x\|bd\|x\|bf\|x\|eb\|x\|fd\|x\|f5\|xae\|x\|f7Z\|xa6B\|x\|cew\|x05\|xc9j\|xd5\|x84J->\|x169\|'K\|x\|bd\|xe3n\|x15\|x85\|xa51\|xa27+\|x85\|xc6\|xa5\|\|xe9\|xe0\|x8d\|x9eY\|xe5p\|xaa\|x88\|xa3\|_xccc0\|_x00\|x19\|xc9\|xd7F\|xe79\|xd\|fE%\|x13\|xd1\|xfaz\|xc\|d\|xca*~\|x97u\|xa6\|xf4\|xd2\|xf2\|xa8\|xb6z\|FB\|xa\|xe\|x\|dc\|x99\|xf4Cg\|xc3I\|xc8\|x\|cc\|x\|cb\|x\|fb\|x\|bc\|xd1\|x0eLy\|xd3\|xd1\|xd5*e*\|TS4\|x90<\|x2F\|x\|cc\|x07\|'W\|xb9\|xc9\|xec\|x064\|xb\|xae6M\|x00\|xb\|xc\|x2\|xb\|E\|xa1\|xae\|x14\|x92A\|xc66h\|xc\|x7\|x15\|xc7\|x10\|\|_x9dQ\|x8f%\|x80\|x95\|xc88\|xe1XY\|xa4\|x90\|x8c\|x8e\|xc7\|x1f\|xe5\|xa2\|x85\|xc8\|x98\|xb8\|x14R\|x1f\|x05[-\|x83\|x\|fd\|x\|b4f.A\|x\|de\|x\|ef\|xd1\|xc0\|x1a\|x9a\|x9c\|x8e\|x\|fd\|xd9\|x\|bd\|x\|b4n\|xd8\|\|\|x92\|x8b1\|xf4\|x\|bd\|x\|c2\|x\|fe\|x\|f6\|x08\|x\|f6\|xd5\|xa8\|xa8\|x\|dc\|x17\|x89\|Lg\|x\|ed\|x14\|\|x\|c8\|x\|c\|x\|c4\|x9e\|x\|eb\|xf8\|x\|c\|x97r\|x\|c\|x\|d9w\|_x\|f0\|xe2\|x\|f\|x\|b\|x\|c6\|\|x920F\|xb8-Aj\|xf1\|xb9\|x16:\|x98\|x99\|x98\|xf2G\|x1fw\|x1e\|x0e\|xaa\|x\|f\|c\|x\|e5\|x19\|x0b4\|xddV\|x8F\|H\|x4\|x94\|xeeH\|x1e\|x01\|x1d\|x\|d\|x\|f4\|x\|f\|x\|a4\|x\|b0\|x82\|xad\|x\|_;\|x92H\|x82E;\|_xd2%\|xe2\|x8c0\|xde0\|xe0o\|xa7\|x\|cb\|xe5\|xa4\|xe4\|xa7'+\|x01i\|x95\|xea\|xe1M\|x1c\|t\|\|x11\|x\|f8\|x\|9c\|x90#\|x\|ff\|x00\|xae\|x83\|x14m\|x\|be\|x\|ca\|x0e\|x13\|\|xb7\|x\|ba\|x\|aa\|x\|a5Zy\|x98\|x80Ga\|x9c\|xe3\|x\|f0\|x\|fe\|x\|fa\|xd2k\|xc2\|xee\|x\|f4\|x\|b8\|xa9\|xb8\|x05\|x05\|x85r\|xa60'0\|x8dT\|xb8\|x05\|x16\|x877Z\|xc4\|x\|f8\|xe2\|x99e\|x03!\|x95\|T\|xb0WX\|\|ad\|x\|db\|x\|cc6zI\|xb9D\|x\|ef\|x04\|x87\|x\|eb\|x\|d8\|x\|ff\|x00\|xa6\|xaa\|\|x\|db\|xa2\|xa4eU\|yb7\|x90\|xb\|x\|cb\|x\|b1\|x89\|x\|fd\|x\|b\|x94#\|x\|cacp\|x1e\|xe3\|x\|f1\|x\|d5ZK\|x0e\|x17*\|x07\|xb9Gw\|x89+\|xa3B\|x0e\|n\|x\|ca\|x87\|xf7[\|x\|fe\|x\|fa0R\|xd3\|x1b\|xb77\|x82\|xb9\|_xb8\|xb4t\|xb4\|xed3\|x\|de\|xec>\|x9a\|x03e\|xd9\|x19+\|x94\|x\|bd\|\|xc1g\|x\|dc\|x\|f6\|x\|fab\|xb9\|x1f\|x\|b\|x90\|x\|f\|c\|x18i\|x\|f\|x\|b5\|x00\|?-\|xaa\|x\|ed\|x\|cd\|x\|x88\|x87\|x93\|x1f\|x80\|x11\|xc7\|xe7\|xad\|i\|xb7*\|x0c\|x05\|x19ZfW\|_xa3\|x00\|x8c\|x\|aa\|x01\|x\|e\|x\|f5\|x\|d0\|x8c\|x96\|[\|x94\|xa7q\|x\|af\|x\|b6\|x81\|x0c\|\|x1a\|x\|b0+\|x\|db\&\|x91%\|x\|bf+\|x\|b7\|xc3Q0N\|x\|bd\|xc5\|_>n<|\|x87\|x8f\|LH\|xa7\|x\|bb\|xae\|x\|d\|x885\|_xaeZm\|x0c\|x1\|xa4h\|x\|be\|x\|x0\|x1\|x01B\|x8d2\|xd1\|reMW\|r\|x00v\|x97\|x1\|x01\|xee=\|xd8\|x\|fd4\|tH'\|x97s\|xb7\|x\|b\|x84\|4\|x93\|x\|cb\|x18r\|x9e\|x\|ca;\|xe3Ch\|x\|da\|x03B\|xe4\|x\|bf\|x\|ed\|xb19\|xe6T\|x020\|x\|e8\|x\|ed\|x1eW\|n\|x\|f7\|x\|f8\|x92\|x16\|x82\|x19\|?hs\|x\|f7F\|x8c\|x1c8\|\|\|x91\|xa2U\|v\|b9\|xee\|x\|ccI\|xee\|x07\|xb9\|x\|f9h\|x8d"\|x89\|\|x\|b4\|xa6\|x\|dbi\|\|x18u\|x88\|x06#\|xe1c\|xe0\|xe9\|xb2\|x\|x84\|x0b!\|\|xbam:\|x98\|xb9\|x95\|xa4\|n\|xe3\|xc3\|\|xec\|xc3Bp\|H5\|x\|d2*\|x\|eb\|\|x91;\|x19\|\|x1e\|=\|x\|c1\|x8e\|x86\|xe0A\|xb5rm\|t0U\|I\|xe9\|xb9\|x93\|x04W\|x\|f6\|xd0\|x9c\|xe4T\|x81a\|NnA\|x92\|x\|ce\|x99e\|a\|x\|d8g\|xc6\|x83!\|x\|b2\|x\|aacr\|xa2\|xea\|xe2hf\|x\|f5\|xe2P\|x00\|x\|f6\|x1a\|JL\|x1b\|n\|x0\|d\|e\|x\|ae\|x\|a6\|xe2\|t\|fN\|xae\|x\|ba\|xa\|x02\|x\|ca\|x\|ea5\|xa2-\|x\|d3K\|x13\|xe4\|x16\|xe5\|x8c*\|x8c\|xea\|xb2j\|xe3\|x8c\|xd9+\|xb6\|x94\|x\|f6\|x8e\|x07\|xb8\|xc4\|xd3\|x\|bcA@\|x18=\|x\|bb\|x\|e99\|xb5\|x\|b2m\|x\|dc\|x\|d1aue\|x01p\|xa4J\|YK%\|xec\|x87=\|x8a\|x\|f6\|x\|d6\|[\|x\|fa\|xa4\|x834\|x8a\|x19\|x84\|\|xae\|x9a5\|xc4\|xc7\|x\|d6Q\|xe1\|x\|b3\|x\|dfB\|x8f\|xae\|xc
```



```
x1AIC\r\b2A-\xae\x82\xa3\xec\xea\xb5\xee\xe0\xbc\x91\x86*LG\x93\xb9\x1f\x03c\xe1\x00d\x1d\|xa9\x97I\x0b\x1dd\
xee.\$8\xbb\xf6\x0b\xc1\xf0p\x05w\xe7xc8B\lx96\xeb\x1c\x0f\xee\|xa9\|x:\j|\xb9n\|xaj|\xfef|x15\x7f\xdbb\xed\xfd\x40\
\xa40\|\xe3c\x11Y\x8bf\xhf\xfb\x96\xae1\x92q\|xc7\''\|'\|xe2#\|x199_\|x8d-\|xd2\xba-\|xb4\xae2Nn\|xf1\|xf4\xfd\xfe\xcx9\x8b\
e4\|\l\x03\x82\x8c\|x8b\|xdc\|x92\|f0\!|\x10\x96T\xbd\xad2\|x9e\x9df\x02\|xb2\|x9f\xdb1nm\|xae,\|xc4\xfr8\|\xf2\x00\xfc\xcf\x4m75\|x
8bl\x8c$>|\x88'\|\x90\x03\x8e\x08\|xee\|x062>|\xea(H\|xda\|xa4\|x93n\|xef\|xab<\|x10H)\|xa9\|xe7\|xb7S,&|\xa6j\|x89\x91\|x832\|b
3\|x00\|x1c\|xe7\|xb9\|x05\|x98_\|x0f88\|xf0q\|xa9\|x07P\|xeb\|x1a\|x19\|xdb5R\|\xb7gos\|xe0\|xf1\|x7F\|xa5\|xf0\|xc5\|xe8\|xd9\|xee\
8Cz\|xadk\|xc8f\|x1fKY-I\|xc9""R\|x00\|xecG\|x8f\|xa8?|\xdb\x1a\|xf6\|xa9\|xfcc\|xeBR\|xea\|x19\|x06\|xa0\|x02l\|x0b7u\|xe1\|fx\|db\
3\|xe2km\|xaa-|\xfb\|yb9\|x8e\|xa3\|x84\|xc5\|x57\|x08\|x0b\|xf1\|xaf\|xb1\|xf4n\|xb2\|xc9\|xc6\|xf6\|fb\|x8a\|xc8\|x92\|x1cd\|xa3(\
zz\|x9a\|xbb\|xa2\|x83?(\|xa2?\|xd4\|xd8f\|4\|xefx\|xb9\|r%\|x0bE\|xaf\|x17\|x1b\|xf5\|x1d\|xa2\|x98\|xd3\|xd1\|xa0g\|xf7o\|xae\|xb6c\|xd
44\|x8c!\|xed)\|x0c\|xf55\|x17i\|x83\|xd6\|xd5\|xb1\|xc9\|xec\|xa0v\|xlan7=\|xc6\|xadT\|x9aL\|xed\|xb6k\|x83*\|xfd\|x85J\|xae-|\xf0\|xf
7\|xd3\|xad\|x07\|x81\|xca\|xa5\|xdf)\|xed\|x0e\|xcg\|xab\|x9c\|xabTM)c\|xf3\|\xdai\|x99P\|x9c\|xc1\|xb5_\|x83\|xf6\|x9c\|xf5P(\|xf91\|x
d1\|xac\|x8c.Y}\|U=U\|xb2_V\|\'|x18\|xf6#M]'.^|\xc9e\|\xb9\|xb9I\|x8a\|xf0\|x1a\|xa\|xac\|xae#\|xac\|xb8\|fefg\|xc7\|xc4<\|x8\|xc7o\|xf\|xd
7\|xdbKH\|xc8\|xe49\|x0b\|xb8J\|xef\|bdb->\|xe3g0\|xb4\|xd4\|xd3\|xfd\|xa2\|xd9d\|xf\|x00\|xdd\|xd6\|xc1\|x96\|x8d\|xfe\|x84\|xf95\|xfc\|x
ad\|x83\|xee24\|xa4\|x9a\}|xa2\|xc2\|x90\|xe4\|x92\|xa2\|xdb4\|xb1\|xfa\|xf1@|x00\|x1e\|xf8\|x1d\|xb4\|xa3\|xe2\|x7f\!|\xaf\|xa2\|x95\|x
d6\|xda\|xaa\|xe3&I\|x1b+|\xe0\|xb0Q\|xdbJ\|xb9\|xae\|n(6\|x10\|xf4\|xf7K\|x8d\|x9aQ\|xe8\|xd6\|xb8+|\xde6Q\|xed\|xa5\|x9f)iD\|xda\|yd
7f\|xd4\|x1dQ\|xdd\|x16\|xb8\|x83X\|xaeTfA0ymw\|x80~\|xc7R\|xbf\|xbe\|xa1\|xd4\|x17\|xa7s\|xed\|\xf2P\|\{xb4o\|xc9\|x80\|x1e\|xa1\|xe8\
|xfc\|xfc*\|xba\|x10x\|\|x9e\|xe1\|xb6\|fa\|xb8\|x96\|x95W\|\}xb9q\|x82E\|x4\|x8a\|x0b\|xc9J\|xf1\|xc88\|xb7\|xe5\|x0c\|xd1\|x9e2\|xa1\|x
1e\|x1b\|x03\|xe5\|xd8\|xf6\|xd7\|xce\|xfe\|\xf8\|x8f\|xa7h#p\|x07\|xdcA\|xab\|xfd9\|xed\|x94\|xfe\|x8fk)\|xc9\|xe1V\|x7f\|x86\|xb6\|xdd
5\|Mt\|xdd5iGKNj\|c\|x8e\|x9e!\|xea\|xc2\|xaa\|xc0\|x00\|xc4\|x86\|xe5\|x1eG\|x13\|x903\|xdb\|xb8\|xec0\|xe7g\|xd5n_xb0o_.\|xa3g\|x04\
|xc9\|xc8\|xcf#\|x91\|xd9z\|x00\|xc6\|xb8\|x84\|xef\|xa6\|x97\|xbd\|x87\|xbcx\|xeem\|xbfZ\|xb6\|xea\|abb-\|xbaj*\|j\|xfb\|x94i_e\|xe2\|t\|x
95\|x80\|xe3\|xe0\|x10\|xa3\|x1f\|x17\|xc4\|\{xb1\|xc6\|x15\|xd7\|tc3\|xaa\|ydz\|\}xd1\|x7\|xcb\|\}x04\|x86\|xd85\|s\|x9b\|xe4-\|x06\|x00\|x
b0\|xa1\|xc0\|x84\|xfe\|xabs[\|xaav\|r\|xae\|xbe\|xe5\}|\x8a\|xe2\|xd4\|xd7@\|xc6\|x11\|x03\|x85\|x902\|xb1c\|x13\|xa1\|xe4\|xc0d\|x96\|x0
7\|xca\|x8c\|xe3_i6\|xc6\|xe7N\|xe8\|x9a\|x08\|xc7z\|xc0\|xc1\|xcf\|xf1\|xf5\|x05\|xee_j\|faf\|x14\|xb77\|x16\|x99\|xf7\|x17L\|xb7\|nz\
|x12\|xd4\|xc6\|x8bew\|xe4\|xf2\|fec\|x00bpyXrnD\|xbe\|x0ep1\|x96\|xc6\|xac\|xed8\|x91\|xdbg\|xf6\|xb8w\|xae\|x0f\|xd8_\|xdb\|x87A82Riw\
|xdc\|xbby\|x9e\|x9bq\|xd0\|zd2n\|\}ec[-t\|xad0\|xb6U2gr\|xb7\|x95w\|xe3_xe2\|xab\|xaf\|x7q\|xf7\|fb0\|xf9\|x00\|ys\|x1f\|x85\|x87\|xd3\}
|xdc\|x83\|xb8\|x06\|x8bG\|x8f\|xe7\|xf6*\|xa2P8\|tu\|xa7\|x7f\|xd9\|xae\|x16jK5\|x81\|x8d%|\xcaj\|xb6ZZ\|x9a\|x96\|xa8h\|x98'\|.e\|xc98
q\|xcb\|xc3'\|x03\|x9c\)|\|x18\|xc6\|xaaM\|x0c\|xac{\|x9d\(|xb6\|x81\|xda\|xaf5@W\|xdf\|\}x18J\|x15\}%\|xe2\|xacWU\|xdb\|xf7\|rj\|xcdX\)|\xf
e\|xcFCs\|xa3\|xf5\|x19\|x10\|xc8\|x84\|xc3\|xc4'_\|ya\|x95n\|xd7\|x00\|xf8_xeb4\|xc2\|x1a\|xd6\|xbe1\|xb0h\|xfd\|x0ew_a\|xe4\|n\|xe4r
\|x8e%\|xdc,\|xiaG\|xdc\|xe4\|af\|bfl\|xd2\|x94\|x94\|f\|x91\|xd7\|xd2\|xc4XQ\|cdc\|xe9@s\|x84#\|xd0\|x95\|x88fe?\|x17\|x6\|xca\|f
9\|xf7\|xd0\|xe3\|r\|x87QN\,'i<|\xb7<-|\xf7\|xfc\|xd5\|xc1\|x04\|x8caF\|xeeK\|cce\|xe0\|xb4\|xd2\|xd1\|xd9\|xb7\|r\|xb2\|xe8k\|xea\|xe02\|x
b5\|r%\|x0f#\|x06H\|xe2=Q\|xf0\|x01\|xc5\|\|\|xb7\|xde\|x03\|xe68\|xe7Zh#\|x9c\|xf8\|xc8\|xda\|7\|xc8\|xac\|xe3\|xfb+\|xb8\|xb5\|xd25W\
|x1a\|x94\|x86\|xc5\|x04"\|x08\|xe3\|x86\|xa6\|xacs\|xaa2I\|xcd\|xd5\|x8f--\|n\|xc7\|xbf\|x15\|fb\|bbe;\|xe3\|bbeN\|bbe\|xf7\|xd3\|xba\|xd
70\|x8\|xa7@|\|xa7@|\|x9f\|x96@|\|x07\|xe8\|x02\|xc\|fe\|x808\|x03\|x8XB\|xfdc\|xcd\|\}x0e\|x17\|x99j\|ccc\|x91\|xb4un\|x18\|x05\|xfb\
|xd1\|x9c\|x91\|xf5\|x04k\|xe9j\|\}x07\|xe2n\|x99\|xf1<-|\}xdf\|xedD,\|xe9\|xb4\|xf2i\|xdd\|x94\|x8et\|xb5\|xd1\|fb\|xb\|x84\|xa6c\|xe5YW
\|xd8\|xfd5\|xe9!\|x84B\|xfa$Wd\|x0c\|x91\|x84\|xca\|xd3\}\|xb5\|xca\|xf1\|x12q\|xc7\|x91$\}\|xbf\|xb6\|xb6a\|xd9\|xe4_xba\|xc10\|xe9n\
|xd6\|xa8T-%\|x02\|x96\|xce\|x15\|x83'g\|xf0\|xc6\|x9e\|x8c6\|xee\|xd5\|x95\|xdc0x\|xac^p\|~/\|xe1\|x8d\|x80\|x00i\|x88\|xd8G\|x0b\|x9
7H,\|x17\|xe9\|xfe\)|xa9\|xdd\|bbe\|xbd\|xff\|x00\|r\|x15\|xac\|bbeUw\|x05\|x01x\|x06\|x9f\|xb8\|xc6\|x08\|fb\|xa7\|bbe\|x99-VR\|xb7Z\|
xe8\|xa9$%\|xe9X\|x8c\|xf6#\|xb6\|x96\|x7\|xd1p\|xca\|xf7\|xb7z\|x93j\|xb7\|aaa\|xa4\|xad\|fb4\|xb9\|x0e2E*\|x86Y\
```

[illegible]

```

a\|x8bW\xd3\x7f\x02\xe1\x87\x3b\x03\xed\xad\x08\x9b\x8eU\r\x04\xc6\x8f\p\cc\x9f\x0c\x87\x19\xff7:\|xdl\x8c<r\|a8
h\|a7v\|xcd\x05\x4c4q\|x5\xff\|x00\x0e\|x9\|xd3m\|x95\|xa\|x85\|xb8V4\|xa2\|xd4\x05\|xa68>0\|xd6\|xf8\|x10i\|x08\|x87Z\|xba
\|xac\|x9fNL\|x9f\|x19\|xe3U*\|x84Q\|xd1m\|xba\|xa\|xa0\|x0cq6I\|xee\|xa7\|xbe\|xa8\|xe8\|xfc.\|n*\|x1f\|xa9\|x955.\|xd02\|x83\|xe
4/\|x8bN\|x8e\|xc6T\|xeff\|xc5\|xd2H\|x8c\|x89\|x1c\|x80=\|x98hN\|x00\|x91\|x85;\|xa\|x06\|xe.\|nd\|x1\|xaa\|x17\|\x91r\|xc0\|xef\|x95
\|xe\|x92\|x97Mb\|x81R\|x1eVI\|xd6j\|xad\|xb0\|xfaq@*\|xb7%\|xc28\|x8c\|x991SE\|x86\|x9ab7\|x85s\|x9c}\|N\|x00\|xf9\|xeb\|xceu-\|x9f\|x
a4\|x8a-\|xd2\|x1f\|x2\|x9b\|x81\|xd28\|xe1b\|xf4;\|xe7jT\|xdd!\|xbeP0r\|xa0\|xa9\|xa2--\|x11\|x87\|x8a$G\|x88\|x1c\|x9b\|xbf6l\|xb1\|
xed\|x820<c_\|x9d>#\|xd3\|xf5f\|xaa\|tm\|xd8q\|\|xd8\|xe7\|xb8\|xac\|x8c\|x95\|xe9\|xb4\|xe\|x887ot\|x07\|fbbD\|xb8\|xdb7\|\|xb7pn\|xc
b\|xfd\|xc6\|xf6\|f\|xb\|db\,|x94\|x14q\|xc4\|x8b4\|x84\|x01\|nR>-\|xde\|t\|x01\|x9b\|xeb\|xac\|b1\|xa1\|xf5t\|xaf\|xd3\|xc4\|xc6\|xc7\|
xb7\|xe6q5\|db8\|x1f\|xd3d\|x13U\|xf5\|x03\|x14\|x8a\|\|xc3r\|xa9\|xbd\|xde\|xfasw\|xe9\|xe4\|xdbR\|xe\|fb2\|xe2\|xa0\|xdc\|xf2W\|xcbUh
4\|x10\|x94\|xaa\|x9agb\|xbc=H\|x80\|x0e\|xbe\|x89\|x19\|x04\|x95\|x04\|x02G\|xc2\|x0e\|xb3\|xa0f\|xb5\|xba\|xafU\|x92n\|x85\|xa0\|x07Y\|
x14\|x00\|xc94n\|xb3\|x91B\|xc8\|xc5\|xa1=\|x8d-\|xc\|f+\|xc0\|xd94\|bdb\|x8e\|xa2\|xc7e\|xdb\|\|x8a\|xb6\|x8e\|(\|x89\|xe4\|xa1\|xba\|xc\|
\|x90H\|x92\|x0fI\|x03\|xc8\|xb7r\|x83\|x91.9*\|xe3\|xeaF4\|xea\|xdd\|x13%\|x2V\|x02\|x1e\|x0e\|E\|x02\|xc7\|x14\|xaa\|xe\|xc2Y\|xcd\|
xb0\|xa6S\|xabJ\|w\|xdc[\|xff\|x00y\|xf4J\|xdd1Z-\|xc6\|x14aK\|\|x96\|x9d\|xa5I\|x90\|x94\|x0c\|xe\|f\|h\|xd2\|(\|x4\|xcarc\|xc7\|x88c\|x825
\|xae:7M\|x83M\|xa5\|xea\|x0fu\|xc4\|xeb\|xf7:\|x80\|xb3c\|x02\|x89%\|xb9\|xaa\|xb\|xc8b\|(\|xf6\|x92\|x0f\)|\x85\|xabm\|xf5\|\xad\|t\|x12.
\|xda\|xda\|x16\|r\|xc3u\|xa5\|xa5Z\|x1a\|xbfZ\|xb4Q\|xbd\|x82\|xb1ZB\|x81$dy=0\|xde\|xe2\|x00\|xf8\|xfb\|x10\|x00V"\|t\|xbd\|/xa4J\|xd6M
+\|xda\|xc2I\|x9fwon\|x82\|x06\|xda\|xc17\|xc7\|x9c\|xa9\|xb7\|x91\|x85\|x1b\|x0d\|xb6+\|xa4\|xb7\|xda\|x8f\|x1L\|xf4\|xd2\|xc9N\|x94
\|x15\|x95\|xb7K\|\|x99\|x05'FGP\|x8c\|x15HA\|xc4\|x84RNT\|xf9\|xc8:\|xd0\|x9fM,\|N\|x80\|x18\|xee\|x8d\|x80\|xd3\|x8d\|xb6\|x08\|xe0\|x9e
s\|x93C?z\|xe6\|xdb\|x81\|xab\|oE:\|x7f\|xd2\|x8d\|xedo\|xb\|cee\|xc\|b\|xe5\|x92\|xebf\|x96\|xc3Q\|x0c\|xb0\|xd7\|xc5\|x14+\|xf6\|x97t\|
|x89KF\|x94\|x8eW\|x99\|x05T\|xb0,\|xa7#\|xb1\|xd6#\|x\|fa\|x\|be\|x\|b\|f\|xa9\|xe8\|x\|dd3\|x1c\|xd9\|x03\|x81\|xb0I\|xa2\|xcdj\|n\|xed\|x82x
\|xe0"\|xb1\|x94\|xec\|xad\|x02\|xe3t\|xea\|x8d?S!\|xc\|f\|xa\|x7\|d5\|r\|xed\|x\|ef\|x\|b\|x\|x19\|\|xa2\|xb9\|xd7\|x8d%\|x90\|x82\|xa3\|x9c\|x
92'\|*t\|r\|x\|da\|xa6\|x11\|xf\|c\|x83\|x92\|xc0\|x05#\|xc\|f\|fa=0\|x\|f4\|xe3>+\|xa6P$'\|;H$\|x0f\|xb05\|x93VF--\|x89\|xc8X\|xe\|fP\|x808\|x1d\|xd
1u;\|xd7\|xa8\|xddA\|xdc\|x14\|xbb\|x93e\|xee\|xda*\|xfbE\|xacIGp\|xb6<-\|x12U\|xc0\|x80\|xcc\|xa6*\|x85>+\|xa69\|xaa\|x8e\|xed\|x9f\|x8
7\|xc1\|xf1Y\|xba--\|x8f\|xa6i\|xb6j\|x190w\|xb8;\|x9a\|x0e\|xe2\|xc8\|xe3\|x8e\|xc6\|x95\|xa1\|x94\|xbAM\|xd8-VW[M>+\|xe0\|xbb\|S\|x16\|x
9a\|x9fqG\|x05\|x1b-\|J8\|xa62\|xcb0"\|x07Y\|x128\|x00#\|x89\|x00r\|x84\|xe1\|xb9\|xa0<\|\|x9e\|xdd\|xd3z\|x1fY\|xd4A\|x7f\|>|.npG\|xcaA
"\|x89<\|x87\|x03\|xc1\|xa3m\|\|x8a\|xb2W\|xea--\|x7\|xd7\|xfb\|x8f\|xd3\|xf9YM\|xe3\|xb3>\|x\|ec\|xaa\|xa9\|xde\|x9d\|x0c\|xdUu\|xd1\|x
c14\|x82\|xba\|x\|ce\|xb5\|xb2\|x1a\|x98\|x1dX\|xf2*\|x0e\|x1aU\|x07\|x97$\|x7f\|xda)\|x\|ef\|x86\|xc9+\|xf60\|x86\|xba\|x0c\|xed\|xd3\|xb7\|
xf1Qm\|x9b\|xb88\|xba<\|x81\|xc0\|xbf\|x03\|x1e0\|c\|j\|xa5\|xd8\|xea\|x07\|xda\|xbf\|xb6\|xee\|xe7\|x\|ac\|x\|af\|xa7Y\|xeba\|x02\|xa9<\|xbf\|
x1e,\|xbf\|C\|x07\|\|xf7K\|x06\|xc3UT\|x93\|x7f\|xb8+>\|xdd\|xb\|bd\|xd4\|x96j\|qn\|xdc\|xb4\|x0bQ\|xe3\|x0cd\|x8c7\|x7f\|x99\|x07\|xd\|x\|ea;\|
x\|8f\|xae\|xbd>+\|x93V\|xf6\|x8d\|xaf\|xc8J>\|x3\|xd9>+\|x87lm\|x9b\|xd5\|x18\|x9a\|xc9P\|xa8\|xac->\|x1eG\|x92\|x0f\|xc\|f\|x\|ca\|x\|af\|xa1\|xd6
\|xcqcE\|x\|xb6\|x1b\|x08:\|x\|fe\|x9f\|x\|dc\|xe0n\|x13\|xc62\|xa7\|dx8\|xa\|x2\|x08\|x1c\|xd3\|x95M\|xc5\|x03.\|xda\|xa8\|xa7\|x1c\|xd6\|'A\|
x9e\|x\|fd\|xf4v\|xb0\|x05k\|xb5\|xa1\|xd9n\|xfb-$\|x0b\|x14j\|xa3\|x1d\|xdd\|x9b\|ChCv\|xe5Km\|xbY\|xe4\|xe2\|x80\|x07R?\|xf0\|xfb\|xe7
\|xfah\|x84\|xb1\|x\|ca\|xb45j\|j\|x9a\|x0c\|xf2Jw\|x18\|xf1\|xc8\|x81\|x8f\|xeas\|x\|db\|xb5\|x05\|xb5\|xc2\|xa7t\|x\|de\|x8a\|xb9\|x0c\|\|B\|x0c
\|xe7\|x00\|xf9\|xd4\|x1ap\|J\|X\|xae\|x9d\|xf0\|xb46e?\|xad->\|xbd,\|n\|xc3+\|x02\|xf\|cR?\|xf\|c*20\|xe3\|xe3\|xeb\|xa5\|xa5\|x9a\|(\|x7\|xb8\|x
ab\|x06\|x92\|xb1\|x\|a6\|x\|fe\|x92\|\|xb6\|x\|f0\|xd2\|\|xf65\|x19\|xb4\|xd2\|x91\|xc54*\|x91Z\|xa1\|x\|be\|xa\|x0eV?\|xa7\|x\|de\|x\|c\|f\|mbj\|x\|fa
\|xb\|x\|dc\|Kc\|x14<\|xa3\|xb2.\|xe5bU\|xb\|b;\|u\|x\|ef\|xaa\|xb9\|x\|be\|xc1\|rE\|xc6w\|xc\|f\|xdakjY\|x\|dc\|x9f\|xc5\|x\|ce
```

s\xca\x2e\x9u\xea\x5e\x5d\x16:\xcb\x4cW;S\xb6S\x9eRjs\xeo\x95S\x90\xcb\x4f\x5e\x8cM\x9e\x2||\x9b\n\x8fir\xfc
c\xd1\x054\x20\x86\x91\x0eq\x7f\x89\xed\x4aH}\x80\xef\x08\xfaX\x2a\x0f\x87U\xcc0\x00\x9f\x8c\x81\xff\x00\xcf\xfa
x\ea\x94\x2h.\x0e\x7f2\xaf6'\xL.;\x8a5\x9b9\xdd\x8bP\xdd\x8bE\x5f\x18y\x07\x7f2\x9a9\x07\x98\x9f\x8e\x9a\x98tFA0\xel1\t\x
cf7\x85\x2a[\xad6\xbb-\x1a\xda\x2b\xfd\x02\x08\x87c\x81\xdd\x8f\xcc\x93\xdd\x8f\x2t\x8cX\x03Z\x15\t\$xae\xedi
\x8c\x1f^\xe1M\x83\x8c\xcf0Q\xd8~:\xa1\x8d\xbc\x95w\x85\x8c6Z(U\xcb\xad9\x00wPF2?=\x03kn\x24\x82J\x9c\xdc2\x8c9,\x
82%E\x2eI\x08\x91yb}\x87\xcc\x9eIF\x7e7\x00\x9ae\x0e\x15f\xcb\x9e9\x85\xabd\x2c\x9\x8b\xee\x2b4\x8e\x8f\x7\x89\x2a2\x
cb\x9b\xf8\xbe\xcc\xbe\x8e\x8a7\xcc0\xccfnM\x4e4\x9e\xde\x07r\x08\x19\x8a\x98\x93K9\xea\x855^\xe9\xaf7\x8z\xde
ff\x0b\x0c@g\x93\x13\x2d\x8\x0f\xae\xbc\xde\xbe\x17\xea\x5e\x2a2\x3J}\xb7*\x8a\x8c5\x2d3\x8a8\x8fag\x2b3\xda\xcfiU7\x1a\x
xa1\xce\x2e1XG\x87#\xbaa\x8f\x2a6;~\x1f\x89\x2d3M\x2d13A\x2a6\x4c\x140'\x2e\xab\xbc\x9\x2d6V\x13\x2d4~\x9b\x2d6\x2d3}\x
c9\x2c34\x820Vi\x9b<\x8b\x7e7\xce0\x2d7\xbf\x4f\x9k\x7e7\x1dK\x8aF\x9f9K\x9c\xdb\x7e7\x9b\xef\xcfu\x2a5\x06\x2a2\x9bI
\x16\x9f\x2U\x2dU\x2d6\xab\\\x4Ih\x2eU\x2d2\x2aax\x9d\x2a0(\x9d\x0be+\x2e3\xefw\xf1\x83\x2d0\x0eq\xaf\x01\xab\x
f8nM\x0bd\x2D-\x2a4\x02I\x2c9\x2a3Go\x2db\x21\x9f9\x2a7\x2e3\x2d42\x84=\x2dc\xab\r\x2ada\x2b\x2d7Z\x13\x2a8;\x8c7\x2d9%\x9
z\x2b4\x2b4\x2d4\x91\x04\x13\x8c8X\x89*\xcfc\xcd\x2b9s\x0b\x8c\x0c\x2e4\xfc\x88\x14?'\x05ju\x1d\x1d\x2a\x8a!\x7e7\x0c\x0
7\x2c0\x2e\x2f4\x06\x80\x2f6eL\x9a\x2d66'\x2d0x\x2e5G\l\x2b7;\N\x2f8\x96\x2e\x2a8\x8cMqy}\x17\x99\x2c0x\x2d\x2a\x98J9
p*\xc68\x2c9\x185\x95S\x9c\x80t\x2cfn\x8e8z\x2b6\x2c1\x2a8\xea#\xf9v\x8e0H\x1c\x83\x2c9\x2a4\x8f\xcc\x2a8}\x2d.\x2d\x2c2\x
9eu\x06\x2dfu\x2e9\x2e5\x86\x2d3\x2b8w\x14\x93\x2d22\x2aax\x2a0\x2c0\l\x2a9\x2e5d\x0f#\x0eK\x1b\x2b2\x2f7\xce\x04\x85\t\x2e
8\x2f8\x2f8@\x2d3z\x8f\x83^5-\x94\x8d\x2d1\x2b8\x0c\x90\ta\x1d\x2b3\x7e7\x2b7\x9c\x83\x2e5@\x2d6na\x03\x91\x2fb\x2ad\x0fat\
x8e\x2cd\x2d4\x8e\x91WwY\x2a1\x8d\x2d95n\x2e53\x2c5\x18\x04!\x89;\x11\x2eCA\n\x08\x2f9\x2a5\x2ee~\x19\x2f8b8Z\n\x2d7\x12\\\x
2d7r~\x2cb7U\x2a9q\x2b\x1e\x146\x2f0\x2d8R\x2da\x84KQ0\x2c6H\x2fb\x17\x2c7\x9f1\x2ebvN\x98\x2d8c\x2aCd\x2a5\x2cb\x2b70wU\x2c7
c\x2ee\x04\x2af\x2a6\x2a7\x12"\x8e\x154\x2fc\x2b0\x2b3\x2c7\x2fd\x2f0}\x2c1\x2f6?N\x2dak\x2a6\x8\x2ed5\x2e0\x81\x2f7Q3C\x2db\x2f5
[]v\x2da\x2db\x2ddC\x2b0&\x2e6\x2db\x20e'\x2a7\x91Ha\x2e1\x2e2otq\x2ec\x2c3\x2e5\x2fe\x84\x1dz\x93\x204S\x2b3{Ra\x2c4\x1a+=\x2e
m\n\x2cb<\x2e7\x94/\x2c0\x2f1\x2eb\x2a0\x2d9/\x2d0\x2ad\x2b8\x2a5\x2f1#\x2d3\x2b7\x2c2\x18c\x2c8\x2c6\x2aex\x2d6Q\x2cfe. \x2a6\x96\x2e
b\x2bdL\$DX\x95'\x2b6[\x2c6\x9ci\x203\x2e5T\x2da\x20b\x2b8\x2c88K\x2cax\x2c\x2b9\x07\x2b6r4f\x2da\x2e4m-\x2d8D0\x9000\x1d\x2c7\x2e7
\x203G\x16\x17,\x2da2X\l\x2d1\x2b9\x200\x2de<y\l\x2fe\x2daP\x8f\x2c8nM+\x2c9\x2b\x2f4\x2f2\x9e\x9dc\x2dc(\x8e\x2b9\x93\x1c5-\$\x83\x2f
a;\x0f\x2aex{\x0f\x2cc\x2f9\x203M\x2c3\x200o\x2b8\x85B\x2ef\x20bG\x83\x2edU\x2ec\x2a1\x2a1|\x2020\x2be\x2dai\x2c0*\x2d8\x2e5:\x2a7\x
a0\x2a7\x2a3\x2e3\x2fb\x2f1\x2da7\x81\x2db\x2b6\x2aax\x2a0\x2b1}\$\x2a0\x2a8\x2e7\x2aex\x2cb\x2a1T~\x2eb\rU\x2d2\x2a1\x2c1\x2c2M~\x2b8F\x
d2\x2b292.0\x2a1W\x2ef\x2f1\x2faif\x19\x1e\tW\x20b\x2b7N\x2b6\x2ccw+\x99\x2dd3\x2c2\x2044\x2e4\x2ad8#\x2ef8\x2f2\x2c0\x2fbq\x2ff\x20
03\x2f4\x2ef\x2d6\x2db\x2f5\n\x2e7\x13T\x11=A\x2bdMR\x8dm\x2a5r#\x2fbF\x2f1\x9d!\x2ad\x91\x2ce;Z\x2a5\x2a2\x2c2\x2b7d\x2ec\
x99\x2a9*5rVBWVC\x2c6\x14\x90vq\x2f7\x2f7\x2faH\x2f0}\x2b0\x2f7\$\x11h4\x2a21\x2bc\x2a9q\x216\x2ab'\x2da\x14\x2c6\x114\x2b4\x2af
\x2dc\x2b1\x90g\x2bf\x2e1\x2ab\x2ba&\x2d5\x2bb\x95P@Y\x2afR\x2f6J=s04 \x2e6@\x2cc\x206?-ak4\x2be\x2b3\x2a8\x2a2\x2b4\x2ed\x16\x2b2\
x2a\x2bd\x83&\x2f9\x2eb\x2b1\x2b5[\x2f2\x84\x2c9r\x2b5d\x203\x2ee*S\x2aax\x2cc\x2ff\x200\x2f2\x2b0\x94\x2feGxst\x2e6\x2ebz\x2a3b\x2a
f\x2bf\x2d8&\x2b9b'\x2a7\x12\x2fd\x205\x2f7\x2e9%\x2046\x2da{E\l'r'\x2a7KIN\x91C\x1d\x2a0\x200\x201\x2fb\x9e\x206\x2bd1\x2fd:7
44\x20c\x20c%~'\x2eb\x2f1\x2faex\x2e7I\x2f7V\x2dc>\x2d9M\x201\x2e2\x2b9'@\x2c7\x83\x2acmGI\x2d8\x2a6\x2da\x11\x99+\x80AZ\x2b6\x2dd\x
16\x2ee\x2b1\x2d5l+\x2b4\x2cb\x208\x2b8A\x8ai\x2a5\x2fb\x910\x90cf\x2fa\x207\x203?x2ca[D\x2b8D\x1f\x11\x8c\x8eUw\x2d3\x2b7\x204
\x2d3\x2f4,\x2bf\x2aex\x2d8\x2dc\x295\x2fd3\x2b8D\x2f0\x2a5\x2d5\x203S\x2a4\x9d\x2bd*\x2a8\x2b2\x2b2FF;3\x203\x9f\x2c6#\x2f3\x2d1\x2ba0
\x2f4^\x2f8\x1f\x2c9\x2ca\x2b6\x2a3\x2fe\x2e0\x2af\x2ba\x2d1\x2d2\x2b13\x2ad}\x2ba\x94\x14?yx\x2f69\x1f\x2b6\x2b5\x2e4\x2d8\x2b7}\x85
0\x203\x97\x2e7}\x2c5\x2b6\x2eal\x2b7Y*!\x2ca\x85'\x2w7\x96\x2b2\x2d2\x2d16\x2cb\x2b8b\x8e\x2d9\x01\x14\x9et\x2f7\x2a8\x2b7}\x85
v[\x95\x1cm%<\x2c3\x15\x2b4D\x2fc3'\x2b7\x2e0\x2c3\x2d8\x2fe^\t\x2d3\x1aI}\x203\x2be\x2a8\x92S\x96\x2e15k\x2db\x2fb\x2d6\x2cd\x
1d\x2fe\x2c9(\x9a\x9a\x2a9y+\x2af\x959\x2c1R=\x889\x04)5\x2e8X\x2c6H\x2cd\x2cd(\x1b\x2b8)K<\x2dd{\x2a2\x2d8E\lD\x85\x2ce9\x2e
3\x2b6\x20e\x2cfcTde\l8\x15:\x2f4F\hf8u\x2ee\x206<\x8f\x2a6\x86\x1aU\x91t'\x89\n\x13 \x2be\x20e\x2eamr\l\x92\x2a2\x20f\x8d
\x2d4\x2f1a\x95\x2d3\x15\x206[\x17i\x2c7q\x2f6}\x95t%\x2e9\x19Tc\x2d9\x2df\x2d8\x1f\x98\x2cfs\x2f8\x0f\x2aex\x2a6(\{\x95\x04\x
x0Z}\x2b2\x8e\x2a6\x2ad\x879\x2c60uP\x2bar\x2a8\x2a1\x2e1SZ\x2ad\x2d5\x14\x2d8S0,\x2008\x03\x2c6\x2b8\x144\x2d2*9~\x2cxe\x80\
x81\x92\x2d8\x07\x2f2\x2d5r\x2b9\x01}\x2b8U+}\x9e9c\x1c\x8e0/\x2eb\x2a5\x2a4\x2de\x2e3JG*'\x2a8\x2cb\x2bd\x2d5m\x2d1l\x2a6Y\x
9f\x87.>.\x2bf\x80\x197\x96\x96\r\x2deU\x89\x2a5_p\x2b8RX\x2adi1\x2a1!R(\x2f8\x2c6\x2a0{\x2cf\x2e3\x2fez\x2bc\x2cf\rM\x05\rh
*\x2cdC6\x2e2\x2be\x2ac\x15NE4x\x92r;\x2fc9\x2f1\x2f8\x93\x2dbY\x2b1D\x2e9\x206r<\x2c0\x95w\x2bbhZ\x8d\x2b2\x2daex\x2e0\x2d1\x2c
6\x204\x8d>\x200Gd\l\x2f1\x2f6\x2d6\x2ab\x2c6\x20}\x204e;\x9a\x2d1\x205=\x2b8D0\x2e2.ZC\x2ef\x8f\l\x9a\x91\x8dcB\x91\x92\x2b2\x2ed\
x2e9\x14\x15W\x99**\x2a4\x2c4\x20b\x99e }\x2d8\x91r\x2c7\x2ff\x200Ho\x2eb\x2ac\x2d7\x2b0\x97\x2d9E\x1cR\x94\x2fd\x19\x2b6\x2b1
x2bcuY\x2ef\x2b5\x11\x2f7\x2b8\x2a6\x9a\x2a5\x2f2\x2be5\x91\x80\x1f\x2d43\x2ff\x200}!\x2d1\x2a1\x2f5:\x83\x2e6#\x2b8\x2fd\x2d5\x2e6
?x2f4\x2e0\x2d7\x2ab<f'\x292/\x85\x2cf\x2c3\x2afJ\x2e8\x9aE\x2a5\x2c15J\x1b\x2a8\x2bbU/\x16\x2c7\x94\x2d8\x2c7fV68Q\x2b7\x2c9\x2a
5_\x208'\x2e1\x2d4\x2b0\x8d\x2c5m\x2a9\x2b1}\x2d8W\x20c'\x2b0\nq\x2e7\x2e9\x2a4\x1d\x1b\x2b9\x2f2"\x203iW kv\x2e4\x2b5\x2f50k4
\x94\x2d2\x2d7?\x2a93\x2a8\x2ca\x2c3p\x87\x1c\x8f\x2d0H\x2b8\x93\x1e\x2e4\x2c9\x2f2\x2ec\rLF9\x9b;Q#;\x98ZW\x2e9\x2cd\x93\x2b9
v\x2ff\x200W\x2f6%=\x2ee\x87\x83\x19\x2a3\x2e3Q\x20f,\x2b4\x13\x20f\x2bf\x19\x2fc\x20f\x2a8 \x2fb\x2ebU\x85\x93\x2c6\x1c;\x2a5\
x2d\x2bd\x2d5\x2e9d\x91\x2c9Y\x08\x2a7\x20cyE/p}\x2e8e\x81," \x2e4\x15\x9a\x2e2\x167z\x2da\x2d5\x2b6j\x2b5u\x84\x80\x2a3\x0c
x84\x2d9\x202\x83i\x2c2\x8b6\x2a8\x2baQ\x2d4\x2da\x2d81t\x2f4\x11e\x2a6\x2b6U8\x15\x2b4\x2e3\x2ca\x2fbz\x2b8\x2c3\x0c\x07\x2f5\
x1f\x91\x06\x2d3J}\x2e8_G\x85\x20ehr\x2d8o1Z\x2eb\x2edks\x2b7\x2b2TQ\x2d5F\x1e9\x17\x2ba\x2b0:\x2d7\x14\x2e6\x2fd6&\x9d\x2a5An=
\x2afE*5M\x1ca\x2e3+\x2f1*\x9f\x8a3\x2ff\x200M\x05\x2f1\x2a8\x2b4@\x2eb5S\x2d1\x2cbn\x94:\x206x\x2db\x2dc\x20f\x2bas\x2ef\x2a0\x9
6\x200\x2aex\x99[\x2aax)\x2ab \x2e2\x2c9\x8e\x2dd\x2f2upmr\x2a4\x2b1RE\x1d<p\x2d2\x2c2\x2aax\x91\x8c"\x91\x2e3\x2fe\x2fau\x08\x2b9S
[\x98FR8Uy\x11\x2f18\x2ec\x06\x2adb\x94\x13\x84\x2da\x9a9aY\x9a\x2af\x20b\x2e4\x2f1\x2d5K\x8d*\x2afS\x2d5H\x2e8t\x19'\x2c9
-\x2f7\x84I\x05r\x9b\x2bd\x2d4F\x2b5M\x2ff\x200\x2ddq\x20b\x2e4\x2e7\x2cf\x2d0i1\x1c1\x2b8v\x201x\x2da\x2dc\x94\x2b5\x12\x2de'\x
x94\x82\x2b1\x94\x89\x0f\x2ee\x8f\x2deo\x2f4\x2fe\x2ba\x2ab=\x82\x2ca\x92-\x2f7f\l\x2ee\x2d5W\x2cb\x98\x2a7\x2a3No6#\x85\x17\x
c0\x2f9\x2feZ\x1x87K6\x2d0\x2ac\x2da\n\x2ff\x200dmZ[\l\nIi\x18\x00}\x2bf\x8d\x2bd\x2c9\x2ff\x200'\x2ad\x16\x2c6\x18)\t\x2d6\x2e2\
x2af\x2acpG\x2e9\x87\x8e<\x05\x1d\x2f2;\x9f\x9e\x2admM\x2ac(\x202\x97\x2dd\x2d5 \x2a7\x82<\x2c02)\x2f1\x964\x19\x9b\x2ecR\x2d
5\x2a8\x15}\x9b\x2f5\x1dH;\x2d4\x2ba5\x20fb\x20b\x2b7&\x19\x2f4\x2aao7?\x2f1k7TV\x2c7\x84\x2d1T?\x2d2\x2cd\x2a9\r-\x2de\x2f6\x2a
9\x87z\x98\x2e0F\x2feU\x04\x2ff\x200\x2a8:/I\x8e\x2a2s\x2c0\x2e5V^\x2cbR\x2bc\x18\x1c\x88\x2cb\x1e\>\x2f7\x2fakY\x2cc4\x84AK
o4\x2d4\x2f5\x10(hG6\x2c2\x82=\x86;\x2f7}\x20c\x2b0\x13\x95\x2c0\x95\x8f\x2f5\x8b'\x2acN%0\x2f1S\x19!\x80\x2fd\x2ef\x9e\x83
,vp\x88\x2c7W+\x2a3\x2a2\x2fdcGY\x2b3\x2a2,q\x2d6\x2fb6\x8eIG\x2c3\x1c\x2eb\x9fM\x2cf\x2c8\x202H\x2fd\x2d6m*\x2fb\x2f7\x2b4\x2b5
\x2c8\x80\x93\x208\x19\x2d1\x9f\x2a8U=;\x2ean\x2ed\x2fb6\x84\x2be\x95\x05\x2d2aM[\x20c\x2bf\x2f83\x83\x2c5_\x2e9\x2df*~\x2b4\x1f!
x2a5\x2ba|\x8e\x86M\x8eV\x95\x2a0\x2b7\x20b\x2f4v\x2fc\x2db\x202\x2e3F\x2d5T\x2f0\x2f3n8\x90/\x96\x1a\x2dc-\x2b0\x97X\x2ee\x2f6\
x2d9\x2f4\x2b3\x2fa\x2f7fh\x2a1r\x2dc\x2a8\x2bbc\x2fb\x2ff\x200\x2f3\x2e5\x2a1l\x2010*AY\x2b6\x2e8\x9e9\x95E\x2beg\x9e\x2df\x2c9\x979\t
x2ee\x206\x95[\x207\x2b2\x2bbj}\x2d9u\x2e9\x2c7R.;\n\x2aamM\x2b7\x2b8\x92I-\x93\x2bf\x2c6\x2a5r'o\x2e3_\x98\x2f9\x8f\x2c7e7\x92
A)\x84\x2edr\x967\x2c6\x15U\x2eeh\x2a3\x8c\lM5I;\x12!hg\x8d\x2b2\x18\x7f\x2aex\x9a\x94\x2d2\x2dbb\x2e0\x14\x2a7\x2eb*\x2e8\x
2eb\x2e3\x2fb8\x2c0\x8aU\x2f0\x2e3\x2ee\x9f\x2c4h~%\x2e3<\x2a3!\x2d23\x04\x2d9\x91\x2b8\x86\x2f0\x2e9\x2e0\x2ea\x2d5K\x95\r\x2aax\x2e
3\x202 \x2a8a\x2df\x1e\x2det\x2c06\x86k\x2ba\x2a2\x82\x2eb\x14H\x8948\x2e0\x07'<\x292}\x2f5\x2d6\x17R\x2e9\r\x2c9\x12X\x2e4\x
9e\x2dd\x2bb\r\r\x2c4\x2aex,\x2eb\x2ab\x2e1\x2a5\x2a3\x2caLr\x2e7\x2b0'\x2b9\x2d50\n6\x14\x9e' \x2b1\x2d4\x2a6i\x14\x2c8\x2c3\x2b0
S\x9c\r\x04\x2b7\x2b9R\x86\x2b9\x2de\x2a3\x2a5\x02\x96)\x2fb\x2a8\x2c326\h\x2a\x12>\x2f0\x2aexUJ]-\x2db\x2cfTM\x2deuax\x2eb\x20eU\x86
\x2fd\x2d4\x2f7?\x2f3\x1f\x2ec;y\x2d3Zxv\x8d\x2c5r\x2f4\x2e5K0u%<W\x2cew\x2f7}\x2b4\x2d3\x9bd\x2aa*+tT91\x2c6\x2a2\$<\x2f1\x2f84
\x01\x2e7U\x2a0\x202\x2e57\x2be.q\x8a|\x89\x08\x2e2H\x2e3\x2a0J\x2db\x16\x2b8wY\x207P*L\x2cd\x204,\x2fd\x2a3\r&\x201\x2f2\x2c7\n?
\x2a6\x20e\x2b2uc\x2da\x202*\x202\x2d0\x2ff\x200F\x2f5h6\x1a\x2cc\x2a7\x2bdM[\x2ce@\x2feW\x2e0?\x2f6kW\x2a703N(*)\x2c9\x92\x2ac\x2af\
x15\x08\x2b3\x16\x85\x2c1*\x87\x2d4R}\x2b3\x2e7N<\x2d9\x2ec\x2ab\x2d9z\x824\x2ad\x2a9\x2b7\x202[\x2e\x2f1cjo\l\x2b83\x85\x85\x2c9F\
x2a\x2b3%U\x2bd\x2e1\x2aax\x84\x2b2H\x2b8\x07\x2f80\x2cbU\x2a0\x1b\x95\x2c8\x04\x2d9\x2cb\x204\x2ad\x2fbzKMa\x88\x1c\x2f\x18\x2b4n\x0fr\x2f\x2d
\x2f4\x2bb\x99f\x2d5\x9b\x2c2\x85\x2aexE\n\x2d3\x2ee\x1f\x2b6F\n}\x2b6\x15\x9f\x9eq\x2c5\x2f0C\x1f\x2cd\x87/\x2ff\x200\x2ddfj\x
d9\x2b2P\x2e4f\x9e\x2cb\x2f5WF\x2b7\x2ac\x1d0\x2e9\x95\r\x2cex\x2a2' \x2d3\x2a4\x02\x1a\x2c0?\x2f3\x14' \x2ff\x200}\l\x2e9\x2de%\x88
\x14\x17\x82\n_\x2bb\x2aI\x20b5<\x2a0\x15\x2f59F\x2e0{\x2b4B+\x95U3-\x2a2\x96\x2a2c\x15D*\x2c5\x2b3\x96\x2e3\x2f7U-\x204\x2a
e\x2aahH\x27\x20fI,\x97\x2e8\x8a\x2faA}\x2fe\x2e3\x01\x2df8\x2f1\x2aax~\x19\x85\x2c4\x2a9\x2deB\x89\x9b\x2a7W\x2f3\x2ac\x2cd\x08\x
x15\x2d1\x96\x02\x2e7\x9c\x12/%\x2cf\x2e1\x2ec\x2f7r@x8d\x2cc\n\x2fb\x85\$W\x2ebU\x2fa\x2da}S\n\x2aex\x2b9\x2c24M\x2db\x2f0\x2ef\
x2a1\x2be\x2d9\x93\x2c2\x2b5\x2da\x202\x92\x2ff\x200:?\x2d9\x2ea\x2a21\x95\x2f3\x1b\x2aex\x20f\x2f7\x2d46Pp\x2a7i\n\x2edW\x2d8\x2a9\x2d
5cG\x206Bf\x93\x2ed\x2a2\x2d9\lE\x2a76\x2eb\x83LH\x91\x2b9r?x\x9f:\x9bR\x9a\x2a5J\x2fa^\x2c9q\x2a9|\x2e7\x03\x2dbR2\x2b9r\x2aax
\x2ac\x2f4\x81\x2e5\x97\x91\x88R\x2df!\x7e7"\x94^R\x14\x2f1e\x15MmD\x2b4\x2f7\x2e0Jf\x2e0\x2d2\x03\x2d8\x2b6\x2f8\x1f\x2e87o\x
2eb\x2a5\x2cb\x2c1\x18P\x01\x05q\x2db\x2abS\x2bb7\x146\x2c4rjg\x2f8\x201\x94\x2f2\x01\x9f\x2e0\x2f9\x0d\x2cchq\x2b3\x2d4\x92\x978\x2af\x2d

[illegible]

[illegible]

\x7fKd\xa6^Gg\xfb3\xfb9\xfb9\xfb9\xdb6\x0b\xe7\xcb6\xe7\xcb6\x9aV\x9e9N\x84\xff\x00\xcc\xa9\x8eag\x96\x12z\xfe\xfe\xfa\x
x6\x92\xadv\x04\x8e\x9d95\xa3. \xa6\xef\xfb4\x0fM, \xc9\x84\xbc4\xa8\x0f\x03\xbc0\x00\xdc4\x8b\x4d4b\xe4wb\x04n{\xa2
x8f _xe4\xfb4c\x85+\xfea\x8e\xfb4\x8c\x9c\x0e\xdb6\x10}J\x1c\xec\xbb-xpH\x03\xff\x00\xdc5\xcb7\x0b: 2\xa2\xe6\x8f\x2f\x
a1<\x96{8\xcb0\xfb\xea, \xd8\x96\xec\xcaR\x94T\xe5M\xa2\xdb4&D\xf5\x9b`q\x1e\xe7\xadyq\xbbu0\x95&36\xcbia}\x94\xf5
\x05\x12\x15\x8c}I _x83\xaaC\x11\xdc\xce5yY}\x1fCS\xaw5wk\x7f\xdb3\x97\x8c\x87^` . \x1c\x8eG\x94\xbb6QK\x17\x87; 1\xd
8\x93\xed\xeb\x92\x92\xfb\x00)\x86\x1c#\x1d\rf|\xef\x98\xden\x02x\xcb7\x9b=0\xed\xdc5r\xbb1\x1e\x9c\xe7}\\$ \x10A\x03
\xdbb\xfb5\xdb4\x9e\x99\x99\x5d\xcbF\x00\x0c\x18\|\xba\x5\x1c\x19\x1n\xa3r\xdb7\x8a\x8e\x9f5\x8a\x05\x02\x9a\Ou>\xb8\
xc0J0\x12\x94\xe7>F\x15\x95\x0cg\x19\x07_ \xaa\x05a\xfb1"%XP\xfb8\x87\xceT\xbb6f\xdbf1\x99v\xad\\\xa5\xdc\x94\xel\x1
d\xdb9l\xca}\xb8\xebp\xbb4\xeb1)PB\x8a\x12\xfb\xdb3\x92@\xceF@xe4\x0f\x9a-\x06\xbeb0\x8a\x04\xef\xdbdw\xdb\xacX\xa6
\xa2\x0c\x00\x06e\xado\xdb4\x1a\x90\x97[_ \xae[r3\x05%G\x82\x82[*\x19\xe4\x95\x820\xbb2R\x90E\x9c\x1e\x9a\x90\x1b
\xdcel\x17\xcb5\x1d\xcf2\xbb1T\xbb2n+RC\xfb7 e\xcb5|\xfcc\xdb6\xfb7x86"? \xb4\xal1\xa7\x01I _02y\x90\x92FF0\xbb3\xde\x1
9\x8b0f\xeb3\xf33. "G!\xf1*~ \xb6\xdau\xdb1x1\x96\xbd\x1a\x84\xdb4i\x10\xe3\xcb6\x9a\xcb4\xbb7\x96\xe8J\xe4\x8c\xae1\x
d4\xbb6T\x8c+\x9a\x02U\x82F\n\xfaQ\xfb1\xa6: \x96j\x11h\xdb4\xa0\x18\th\|\&\x99Uv\xe6\xa7R\x90\x8fA\|t\|ec\x3q! \x05
\x07\xeaK\xao\xeo\x92\xebf}\xb8f}\xe99. \xc1?\x11\xcb2\xcb9\x8c\xbb; \xddT\xadR\xeb6\xcb6r\xe7\xaa\x84\xbdNR=E\xa0-.n%
9B\xceI?\xb6!Y\xeb\xbb1\xa2\x0e\xdb\x9e\xe2\x03\xcc\xcb4\xcb5\xdbf\x99\xdbf\x10\xef\x10\xdc\x97\x1a\x9a\x9\x17\x1b\
xa9\xbb\xbb4c\x8a\x8106\xbb0\xae\x81\xcb9\x0f\xdb3P\x8c\x9e\x9b\xf22\xa6\x1d8\x89\xbb5>\x1e\x8d\xeiQui\x11t\t\x
8c\xbbdG\x9f\x99M\xca\x99\xf5\xf8\x1d4\xe2=\xcfx\xce\xba\x00\xbb6\x80\xeas\xdb9M\x9b\xee*-96. \xa1\xbb5[\xa2k\xfb2\xe
2\xcb6K\x15\xdb6\x8c\x84. +<\x1a\x0e\xe4\xfb3H\x1e\xdeA\xfe\xba\xeb6y\xdb8\xdb\x1e`GFt<<|\xbcf\xfb1\xfb1n\xcb4\xcc\xbb7u\xcb
3\xfb8q\rF\x8f\xcb9\xfb7\x16\x11\x19D\x91\xcb1D\x1c\xaf\xae\xfb2\x94\xfb2\xcb7c\xbb2\x0e\xa5, H\xdb4\xa8\x8e5\xa6\xcb0)\xe
3i7r\x85\xba\xaf\xdb0QW\xfbcnp\x94\xdb55\xe9E\xbb0\xeb\x88B\xbd. K\t\$ \x04\xac\xa5D\x81\x93\xcb7\x00\x8c\xe4{\xb1b\x00
\xdb7\x04\xal1*@\x8c\xbb\xcbf\xe3\x1e\xeb\xbbf\xed-\xc4\xbb7791\xdb6\xe5\xe7R\x87p6\x9ad\x10\x86\xda\xaf\xbb3\xe9\xbb6\
xe4\xdb4\xfbYl\xba\xdb7\xaa\\\t8R\xdb6: \x00\x00\x18s\x06fc\xfb3\x17\xe8\x85\x0b]\x0f\xfb8\xfb8\x90m\xbb\xcbf\xe2\x83|\x
/ad\xdb6m: \xcce\x9\xdb5\x8d\x14SE9\xcb6]\x90\xa2c\x0b\xeb\xcb6\x96^>\x0s\x15\xfb99\xe7\x86HN\x07Zn. yuz\x80B\xa9\xbb5\x14a=
/x0E0\x9a\x84! "%\x1d\xdb5\x9c\xfb2t\x1c\xfb2c\xfb8d1\xfb6\x1e\x988\x9d. \xcfc\xfb0\xcb8\xdb9\x149Q0\xbb6\x2t\x15\xcb1+G
G\xfb\x1dM\x93\xcb1c\xfbfb\xfb8f, \xb8f\xdc"\x8a\xad\xfb0\x9b\xbb\xdb6D\xa5\xcb2\xaf\xdb2\t\x01]\x04\x05d\xfe\xbd\xeaV\xfb
1\xbb3\`bT\x99\xfb1: \xcfbh\xdb9[\xae\x14c!\xdba[\x880\xea\x8c\x9f\xfb6\xdb0\x9c/]F\xfb\xeb2\xfb4e\xfb5\x9b\xfb0\xeb\n\xfb5R
\rA\xbb4\xa5}\xeb3?\xfcc\xeb\x9d\xe4\xab\x05&?\x13\x860\xbe\xbb5\xfb0\xedG\xbb4hd"? \xa8J{\x18\xeb: \xe4b\xfb3\xbb3\x06\x
d8\x9a\xfb9Q\x00\xfbf\xfbal\xcbZ\xab\x84\xadl\xa0\xaf(R\x0fb\xdb2\x7fc\za200\xber\x1aA\x93\xcb8G4\xdb2\xee\xebjM\xcb1\
x03\x13\x16\xda\x8a\xbc\xab=c]\x8fY\x14\n\x80\x8a\x08\xfb6\xcb1#@ \xa4\xcb5[\xadb\xcb9<\x00\xdb2z\xdb4\x1eve\xfb4}\xa5
\xdb4\xde\xcb4\x0b\xbe\xee\x18\x94\xeaZ\xfb9<H#\xfba\xea\x1f\xa7"\xb3\x92F\xe5\x8c\x14-K\xbbag\xcb4t; \x1a\xe7hV\xa8H\
xacQ(\xf9o\x0eH\xfb4\x89R\x00F\x12|\xfbd8\xe8\xfb7\xdb7\xdbf. |]\x82\xa75\xa81{ lve\xa3\xfb7\x13\xdb6\x96\xe7\xdbf\x92\x
eW\xad\x03p|\xa4\x10\x12\x85pi\x96=. 9\xeo\x95+\x0eq@Rr~\xb5\x03\xdb0\x18\xcb6\xbb7\xfb8\x84g\xbb8>\x8dE\x1e\x19\xdb0\x0b
t6\xbb2\xec\xcb6T\x9d\xbb0\xbb1\x1a\x80\xcb3\x0c\x06\xa6\x04QD4\x14V\xa5\x9c. 6]\x04\x10=\xa5` 4\x82q\xdb1H=\x14\x8dce\xcb6\xe3B
\xa2\xbbd7V\xbb67\n\xed\x9b\x12\xdb1\xdc\xdb9\xab\xfb2\xdd\xa8?\x06\xe9\xbb5\x1a5\xad\xfbf\xcb3hzCJ. -*\xc8\xfbf\xcb9\n\
xeb\x91\x03\x05g\x1a5\n}\xc3\xbb1<I]|\x18M\xbb8\xbb\xbb7&\xbfb\xbb3\xfb1w. \xdb9\x83&%A\x0e\xa2\x97Z\x8e\xbb9. 4\xab6\x14\
x12\x82\xcb2\xca[<V\xa5\x80q\x8f\x05\$w\xdb6\x1e23 z\x8a\xa18\xed\xcbf\xda\x04]; \xb8c\x8b\x963NS\x9d\x8a\x85\x03}\FKP\
x1d\xe3\xeb\xcb8R\x89\xdb3\x8ce*\x18\x1fl\xfb7\xef\xdb0\xdb8\x1c1j\x1a)5<Zu\xfbaf\xccj\xdb3\xfb4XP\x1fJ. \x8a\x08A\n\x0
5A\x1d\x85\xbb4r{\x05_ \x7f\xbb0\xdb2\x03\x91\xa8\xdb2\x00\x17\x0e*6. \xdbb^v\xddn\x9dY\xbb3WPV\xao\xfb3i\x872#\xbcf\xbd\
x07\x1a0\xdb4\xbb5c\xae\x7fP\xcbf\xdbf\x1aaT \xdb8\x8bVp1\x18\xdb3\xfb8[\x81P\xda\xcb5\xcb6\xbb0\xaa\xfb57\x1a\xfbcf6(\x97N\
xa86r\$4\xbc\xeb1\n\x1f\xfb3\x0e\xbb3\xa7\xa0\xe2\xbb5cd\xee\xfbf\xdb8\x88K\xfb1\xe5)\xad\xcbf\xdb8\xe4*\x94\xcb7\xce\xdb4\
x06&Ij\x07\x00#\x8a\xcb7\xdbf\xcb0\xff\x00m\` \xcxbN\xdb8\x8c\xfb1\x9c&M\xdb8\x9b5\x9a\x9d[\x8bZ\x83\x1d\xcb0\x14\x54\xca
\xc96\xbb4\xfb2\xed%\$0\$1c~\xc1C\|\x12h\xfb1\x9d\xbb0?\xdb3,>f\x83\xdb9\xad\x84\xbb1\x1f\xa52\xfb5\xcbf5k\xdb6_g\xdb5\xeb\x
ce\xa8\xcb6\xa9\xfb3\x17n; \x967\xbb6\xcb6m\xdb5\x16\xa1\xfb3p[a\xdc\xa7\xa2FI\xfe\x9asb\xcb6\x16\xdb6)I-FE\xa5J\xaaR[E*
\x8b\x11L\xa0\xfb5\xcb8\x8e' \xf6\x1a\x00Yz\x8c\xfb6\xcb3+Mu\xfbDj}\xa8\xa9a\xcb0\x97\x15\x94\xbaVH\xefF=" \xdb9q\x99\
xa2, ?\x89z\x9d\xa7\x01\x8a\x1dZ\x98\xe3\x84\x802\x07@j\x9fP\xa8\xdc\x99\xfb1\x83\xdb4\xbb9\xdc\xa6\xa5\xbb3}R\x1cb5\
95\x1e\x90\x85\xbe\xeb0nJB<\xfdf\xfb\xdb5"\xdbc\|\x9c\xdb2\x1fb\x81\xbb4\xba\r\x8f \x1aEv\x94\x86a]!)\$xe3B\xdb5\xdb5\xdb
1\x9b\xcc\x91bQK\xdb9\xaaKUw\x05\x11e(\$\x14\x94\x8f\xcb\xa9\xfb3\xfb8X\xfb3-T~/ \xa1\xdc\x06\xdd\xcb9\xfb1mv\|\xa0V\x
1dG2\x93\xcb1C\xdbf\xfbf\x00\x9d|\xfbb\xfb8)\x85\xe9\xdb6\|Tf|w3\xcb6\xea\xdb7\x97o\xdb3\x9c\x9e\xdc\x1em\x04\xe5*\x039\
xdb77?\xdb3\x8f; Y1\xcb0\xcb0\xcb1\xab?r\xa3\xdc\xdb6\xfb0n \xcbb\xcb7#\x87\x1e\xcb9\xfb\x9a\xfb1-e\xfb6&\x02\xa4\xfb9\x96\
xf5\xdb2\xe5r\xfb1c1c\x1c\` \x0b0?B\x99FM\xfb2\xbc\x86\xbb64%9c` \xdbb\x14\x17\x99w\xdc\x13\xdb7rQ|\x02-\x85-d\xa7\
xce\xbb>\x1f\xdb3\xce#\xdb4\x81\xfb3\xdb9\xa1\x17\xdb4\x8d\xbb6\xa8M\x8e\x995\xa7V\xfb2T\xa3\xcb1\xbb59\xcb5\x19\xef\xbb2\
xa1\x90A\xfb3\x9f\xdb3\xfb5\xdb7i1\x92\xa0\x98\xbb2\xcaF\xa3R\xcb0\xbb2\xef[B\x98\xcb5\xcb7E\xfb1c\x96\x15\xcb5ap\x9cJ\x9d
\x86\x90\xa0PT1\xdb7\xf78#\x19\xefL\xeo\xcb\xbb1\xfb1\x04\xba\xbb0\xe2a\xad\x03s\xfb6\xdb0\xdb5*3\xef{\x1bb: \xe6a\x86\
x1e\x99\x11\xcb5F*`xdbb\x19T\x82\x11\xdb1*\X\x04d\x8c\x16\xcb9\xdb0\xde\x9a\x99\x10\xdb5\x9e\xeb2|\x16@\xb4#\xdb7\xe1z\x9f\
x07v, \xe9\xfb4\xaa\x15L\xdb6j\x11\x80\x02\|. HT\x97\n\xca\x92\x16\x08\xcb0Y\xca5\xcb0\xfb1)Q\x1d\xa8\xbb3\x0b\x07\xdb
0\x93\xeb4\x1c1r\x99EsU)\x16\xec: \xedr\xcb8\xba\xa14\xfb2k\x11!\x06b0\xeb1\xfb4p\xbb4\x14\xcb7ymd\x84(\x04\x94\xfb4\x12
\t\xec\r\x11\xa4\$ \x8e\xef\x07\xbb8T]\xdb9\xdb4\xdbb\x86\x0c\xfb1L, %n, %)\xf5\x00\xeb6\x85{\x04c\xbb3\xfb6\xcb6=\xf4\xbb1\
x5\x18\xdbd^\xa5\xef\xeb0w\x1d \xc0\xbb8-\xe7\x1dzW\x00\xff\x00\xa3!\|\xb2\xa0H\xeb0A8\tg\xfb\x8d\x0b\x03z\x87\xa
2\xbb54N\xdb5V\x19\xbbd|\xdbf\xcb5\xec\x86\xeb3R\x1c|\xaea\xdaB\x96\x8c4\xcb\x98\xcb1P\x03\x1d\x9cc\xfbf\xdbfM\xcb6\xdc\x8
5\x8e\xfbf\xce\xa5N\xcb4\x17\xbb4o\xcb\xdb6\xdb6\xbe\xda\xa2]\xdb5\xa8Ri,) #;\xdb6'\x01g\xa4\xfb\x8ar: \x1a\$ r\x998\x9
3\xa9\x8e\xa0\x8b\x9a\x1e\x8dv\xda\xfb3LY\xdb4\xea\xa4\x98\xaf; . \x9c\xeb\x11\xa6\x14)*P\x1c\x8aN: 9! \xad5\xa8\x82
L\x1eQa\xfb0\xbb7\xbb4\x1b\` \xb9w\xbb6\xe8^x5\xea3m\xfb\x18\xe52\xfbfLCH\x05! \xc4\xcbZR\xdbf_ \xfbe\xfb7\xbe\xa3\xcb3
x8b\x16L\xaeH\xeb[\xae4e8j\x8a\x10\x7f5"\xeffW\xcb3\xbb5\x90\xbb4\xfb\xcc\xfb4\xfbfH\x07\` \xdb2=\x7f\xdbDj\xfb\xdbbK\xfb2-B\x
9e\xe5\xea0\x8f\xe4\x064\xdb09\xfaeBT6\xcb5\x19\xcb7\$ \xcb0\xcb6\x163\x9f\xbe\x96\x83\x92T{Z\x9b\x87\x9b0\xbb2\xfb5J\x
b5M\x0f\xdbd\xdb3\x12\xcb2\xdb0\x9ea\x07\xcb7\xe8; \xdb1\xa62\xa7pK\xfb2\x97i\xaa\xce\xab\xdbf0Y\x94\xcb8`0\t.p\x1cs\xeb3
N\xfbf\xcb7\x8cXf\xbb9Gw\xfb7n\xdb2\xba\x13*\xabAK\x90\x82\x81K\xcb8\x07\x18\x1a\x00\x8c\x1bs[\` -F]#r*fW\x05\x1d\x943
C, : \x94\xe08\xa4\x0c+\x03\xcb6t\xfbce\x81\xdb0\x92\xe4\x1e\xe9EQ\xae2\xdb5E\xdb97\x13\x89m\xdb9\x10P\x93\xeb0\xfbf\xdb
4\xfb6(E\x9e\xe7\x80\xe3\xa3<\xbdb\xbe\xdb4X\x14\xdb13\xdb7\x1e\xab\` \x8a\x94\x9e\xcb2\x80\xdb03\x80\xbb3\xcb1Kj- \xf7f\x
ea\xbb5\xbb7R\x9e\xaa\x81\x90\x94>\xc8Q\x04t\t\xdb4y\x97\x16m\x9f\x89f2\xe8(DL\xbb\xca\x80\xfb5\x16\xa5I. Fo\xa2\xe
6\x01\x1f\xdbb\|\x8c\xdb8\xdb7h\xbb1: 8]\x8ce\xbb4\xad\xbb4\$ m\xfbf\xdbdXun#(\x01\$\x92\x94{\x1f\xdbb\xdbDR\x15\x16\x0f\xdb
aR\x118\xfbfL\xeb3wK\xbb7r\rcd2\x98\x94T\xfb0\x8a\x07\x17\x1c>V>\x04\xaf\x1b\x9eJg[\xb8\x03\xfb3\|tm\n\xcd\x08\xa1H~
\x19o\x05I\x19#W\xa9R5%\xeb1\xee\x8bwm\xbam\xad\x1ae\xdb3FC\xfb3\xeb1z%M%\xdb5\x90\xbb0\xdb7\x1f\xcb\xdb0\`#\xa4\xfb41\x
e4\x82z\xcb3\x98\x15\xdcAc\x16\xfb2. \xdd\xfb9\xdbce\xbbbt\x10\$A\x8f)H%\xb6\$ \xb8\xdb6R\x9f\x03%W\x93\xe4\xfbfFmW\xbb8
\xfbf\xfb4\xcb0\xbb1\xbb9mK\xa2*\x95>56\xfbf\xbb7\x94\xa4\x95\x96\x956k\xa1\xa5\xcb7X)Qq\` \xa0\xe18#\xc8\xcb7, \xe1Y\xfbfA\
\x82\x850\x9e\xcb8MX3S]\` \n\x8d\x03h\xdbf\x8f\x12\xe7\xabIL\xca\x8b\xad\xbe\xda\x1b@\x8c\xcb0R\x909\x12\xa0J\x943\
x95\x14\x91\x9e)\x1a\x1a\x96<\x94\x0c\x9b*\xb3\xa7#Z\x85\xfb6\xdbd\x936U\xbb5Pa\xba\xe6\xba\x8b2X). \x16\xa3\x
eb\xbb5\$!YK\xccq\x00\xbb4\xbb0\xdb2\xcb0\x089\xcb1t)\` \x8\x17\x81hA\xeeN\xcb6\x9b\x90\x8a{J\xdb9\x15J; \|\|\xf1>R\xad\x06{
\xdb0*\xcce\xa1\` \xa4\x8f\xdb2\xdb3\xe9\xcb1>T\x0eq\x81\x92: \x1djTj\xeb4\x9cM\xdc-B\x15\xbb5\`t#U\x95F\xdb4\x9aa\xca[\|N\x
14\x85\xcbj\x01u\xeb5\x95\x1f\xcc\xaf\$ \xa4\x8f\x7f: \xc6\x01D\xdb4`[\x9f\xdd\xaa, Rb\xcaG\xfb9\xbb8\x89"\xa2\x87\x9b\
W\xce\x9f\xcb8\x9f\xeb4\t>\x14\t\cb8\xfbZR\x90\xad\xfb3c1aFF\xbb9\xad\x0e\x7\xdbd\x0c\xca\x15\x14T#\x16c6p\xfbfA\x90
\x98\xfbf\x1a\xdbfFp\x9eG\xcb2\x8et\xdb7%\x9a\xcb4\x105_ \x10\xa2\xcb\xbbn\xfbjK\xfb61\xa7\x890j\xaa*\x88\xcb\x87\xbb8k
J\xcb1Rp\xfb70\xcb1\xdb3\x10\xbb0\x05\x07\xcc\x03@\xdbc\xdbf\xfb8\x1b\xde\x9b~\x1c\xbb\x8d\x8a\xcc\xdb4\xbb5\xfb8\xe5\xdb12
Aqq\t[\^xa7\xa4\x9c\x0f\xbb7\x16\xcb64\x1e>\x86\xcbf2I\x8c\xfb2q\xdbf\x15\xfb\x01\xff\x00\x11\x87\xba; \xab\xfbfI\xa8I
I\xbb2m\xa878K)" \x1a\x95\x8b\xfb3\xca\xeb9D\xfb9Udp\x06\xbed\xdb8\x92\xacAI\xfb1\` \xf4f\xcb7T\xfbf4\x9a\x88\xe4\xfb4\
x9f\xca0\xbb1\xdb2\x1d8JW)}\x05\xa7\xbbdt\x0b\xa9\x85\x11J\x1e\xa6\x08i\xee\xbb1\xfb\xfbf\x9a\x10n\x19E\xfb9\x95\x94
\x8d\x9e\xdbd9U\xbb7o\x8bB!\xc1Y\$\xf2\xcb62|\x83\xef\xfbfQ\` PK\x00: \x8d\x0b\n\xeb8\xa7\xde<\xf6\xfbf\x00r\xdag\xeb6\x8
2BK\x9c|\xfbf\x00\xed\x8a6\x86\xed00L\xbc\x0f\xedR\xa5\xed\xfb5i\x14; RB\x16\xcb1QP\x1c\xbb2: \xfdb~\xfba\xcb6FSW4r\n\x0
1\x82\x15G!\` \x8e\x8aK\xa8a\x02Z\x13\x85\x0eC#\x8a\xfbf\xfb7\x8a\xabq\xdbVMN\xdb8\x88\x9e\x120\x11\x88\xebI\n\xceF\
xb3(\xa5\x19a\x8a\x16\xbb4k\x1e\xbbY\x9aT\x97\xfb6\x86\x942\x018\x07b7>6v\xdb4\x7f\xa9\x94\xdc\xba\x99c~n\xdb\xaf\x

[illegible]

[illegible]

10) Testing (3%)

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

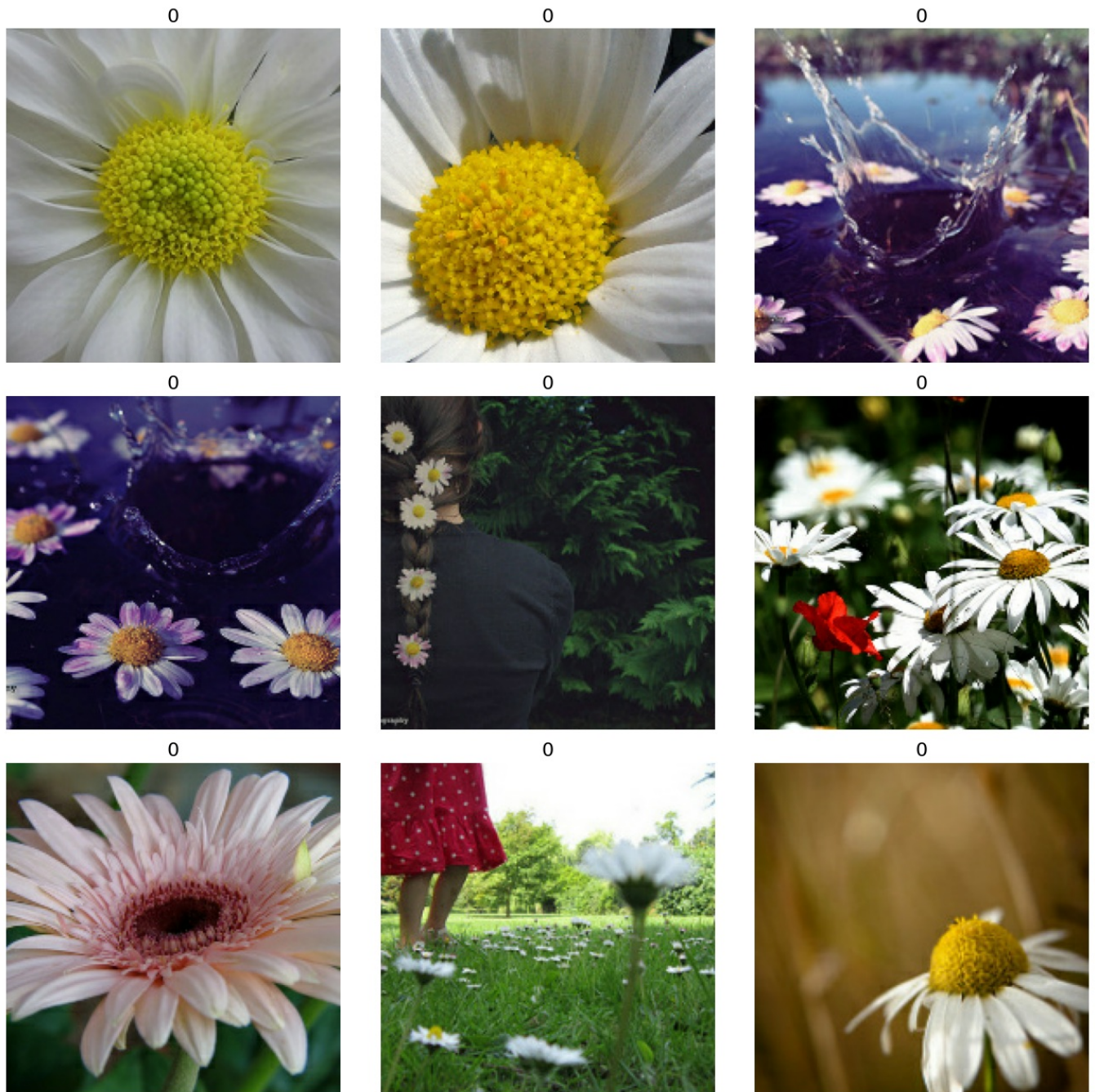
```
In [ ]: #GCS_OUTPUT = 'gs://flowers-public/tfrecords-jpeg-192x192-2/' # We should remove this line to use our own files
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' #These are the files that we generated above.

def read_tfrecord(example):
    features = {
        "image": tf.io.FixedLenFeature([], tf.string), # tf.string = bytestring (not text string)
        "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] means scalar
    }
    # decode the TFRecord
    example = tf.io.parse_single_example(example, features)
    image = tf.image.decode_jpeg(example['image'], channels=3)
    image = tf.reshape(image, [*TARGET_SIZE, 3])
    class_num = example['class']
    return image, class_num

def load_dataset(filenamees):
    # read from TFRecords. For optimal performance, read from multiple
    # TFRecord files at once and set the option experimental_deterministic = False
    # to allow order-altering optimizations.
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False

    dataset = tf.data.TFRecordDataset(filenamees)
    dataset = dataset.with_options(option_no_order)
    dataset = dataset.map(read_tfrecord)
    return dataset

filenamees = tf.io.gfile.glob(GCS_OUTPUT + "/*.tfrec")
DatasetDec = load_dataset(filenamees)
display_9_images_from_dataset(DatasetDec)
```



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

```
In [ ]: %%writefile spark_write_tfrec.py
#We first import the necessary libraries
import os
os.environ['PROTOCOL_BUFFERS_PYTHON_IMPLEMENTATION'] = 'python'
import os, sys, math
import numpy as np
#import scipy as sp
#import scipy.stats
import time
import string
import datetime
import random
from matplotlib import pyplot as plt
#import matplotlib.pyplot as plt
import tensorflow as tf
print("Tensorflow version " + tf.__version__)
import pickle
import pyspark
from pyspark.sql import SQLContext
from pyspark.sql import Row

#We define the variables that will be needed for this project.
GCS_PATTERN = 'gs://flowers-public/*/*.jpg' #glob pattern for input files.
PROJECT = 'daring-healer-421511'
BUCKET = 'gs://{}-storage'.format(PROJECT) # This is the bucket storage.
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # We prefix for output file names.
PARTITIONS = 16 # This is the number of partitions that will be used later.
TARGET_SIZE = [192, 192]
```

```

#Task 1a) - i) Copy over the mapping functions and adapt the resizing and recompression functions to Spark (only)
def decode_jpeg_and_label(filepath):
    #This function reads and decodes the JPEG image from the file.
    bits = tf.io.read_file(filepath)
    image = tf.image.decode_jpeg(bits) #We decode the read data as JPEG images.
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/') #We extract the label (flower name) of the image
    label2 = label.values[-2]
    return image, label2

def resize_and_crop_image(data): #We specify one argument
    image, label = data #The input is basically a tuple containing the image data and the label.
    #We determine the dimensions of the input images and the target dimensions.
    w = tf.shape(image)[0]
    h = tf.shape(image)[1]
    tw = TARGET_SIZE[1] #Target Width
    th = TARGET_SIZE[0] #Target Height
    #Calculate resize criterion to maintain the aspect ratio.
    resize_crit = (w * th) / (h * tw)
    #We resize the images to fill the target dimensions but also preserving the aspect ratio.
    image = tf.cond(resize_crit < 1,
                    lambda: tf.image.resize(image, [w*tw/w, h*tw/w]), # if true
                    lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
                    )
    nw = tf.shape(image)[0]
    nh = tf.shape(image)[1]
    image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh - th) // 2, tw, th)
    return (image, label)

def recompress_image(data):
    image, label = data #The input is basically a tuple containing the image data and the label.
    #This function is used to reduce the amount of data, but takes some time
    image = tf.cast(image, tf.uint8) #The image is transformed to a uint8 so the images can be compressed to JPEG
    image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampling=False) #We compress the images
    return (image, label)

#Task 1a) - ii) Replace the TensorFlow Dataset objects with RDDs, starting with an RDD that contains the list
from pyspark.sql import SparkSession
sc = pyspark.SparkContext.getOrCreate() #Create or get the existing Spark context
spark = SparkSession.builder.getOrCreate() #Initialize the Spark session

image_paths = tf.io.gfile.glob(GCS_PATTERN) #We retrieve the list of image filenames
image_rdd = sc.parallelize(image_paths) #We create an RDD from the list of image filenames

#Task 1a) - iii) Sample the RDD to a smaller number at an appropriate position in the code. Specify a sampling
rdd1_sample = image_rdd.sample(False, 0.02) #We first sample the RDD with a sampling factor of 0.02
rdd2_decode_jpeg_and_label = image_rdd.map(decode_jpeg_and_label) #We decode the jpeg and label of the images
rdd3_resize_and_crop_image = rdd2_decode_jpeg_and_label.map(resize_and_crop_image) #We resize and crop the images
rdd4_recompress_image = rdd3_resize_and_crop_image.map(recompress_image) #We recompress the cropped images

#Task 1a) - iv) Then use the functions from above to write the TFRecord files.
#This function is used to create a TensorFlow feature for storing byte strings. This is for the image data.
def _bytestring_feature(list_of_bytestrings):
    return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))
#This function is used for storing integer values. This is for the labels.
def _int_feature(list_of_ints):
    return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
#This function is used to create a TFRecord from image bytes and a label.
def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data records
    class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
    one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
    feature = {
        "image": _bytestring_feature([img_bytes]), # one image in the list
        "class": _int_feature([class_num]) #, # one class in the list
    }
    return tf.train.Example(features=tf.train.Features(feature=feature))

print("Writing TFRecords")
#This function is used to write TFRecords, which will be applied to partitions of an RDD.
def write_tfrecords(partition_index, partition):
    filename = GCS_OUTPUT + "{}.tfrec".format(partition_index) #We define the filename for this TFRecord file based on the partition index
    #We Open a TFRecordWriter object for the filename.
    with tf.io.TFRecordWriter(filename) as out_file:
        for element in partition:
            image=element[0]
            label=element[1]
            example = to_tfrecord(out_file,
                                image.numpy(), # The re-compressed image is already a byte string
                                label.numpy()
                                )
            out_file.write(example.SerializeToString()) #We convert the image and label to TFRecord format and write it to the file
    return [filename] #We return the filename of the written TFRecord.

#Task 1a) - v) The code for writing to the TFRecord files needs to be put into a function, that can be applied
#should return the filename, so that you have a list of the created TFRecord files.
rdd5_partitions = rdd4_recompress_image.repartition(PARTITIONS) #We repartition the RDD to have the number of partitions
rdd1_filenames = rdd5_partitions.mapPartitionsWithIndex(write_tfrecords) #We specify an index for each partition

#Task 1b) i) Read from the TFRecord Dataset, using load_dataset and display_9_images_from_dataset to test.

```

```
#GCS_OUTPUT = 'gs://flowers-public/tfrecords-jpeg-192x192-2/' # We should remove this line to use our own files
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' #These are the files that we generated above.
```

```
def read_tfrecord(example):
    features = {
        "image": tf.io.FixedLenFeature([], tf.string), # tf.string = bytestring (not text string)
        "class": tf.io.FixedLenFeature([], tf.int64) #, # shape [] means scalar
    }
    # decode the TFRecord
    example = tf.io.parse_single_example(example, features)
    image = tf.image.decode_jpeg(example['image'], channels=3)
    image = tf.reshape(image, [*TARGET_SIZE, 3])
    class_num = example['class']
    return image, class_num

def load_dataset(filenamees):
    # read from TFRecords. For optimal performance, read from multiple
    # TFRecord files at once and set the option experimental_deterministic = False
    # to allow order-altering optimizations.
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False

    dataset = tf.data.TFRecordDataset(filenamees)
    dataset = dataset.with_options(option_no_order)
    dataset = dataset.map(read_tfrecord)
    return dataset

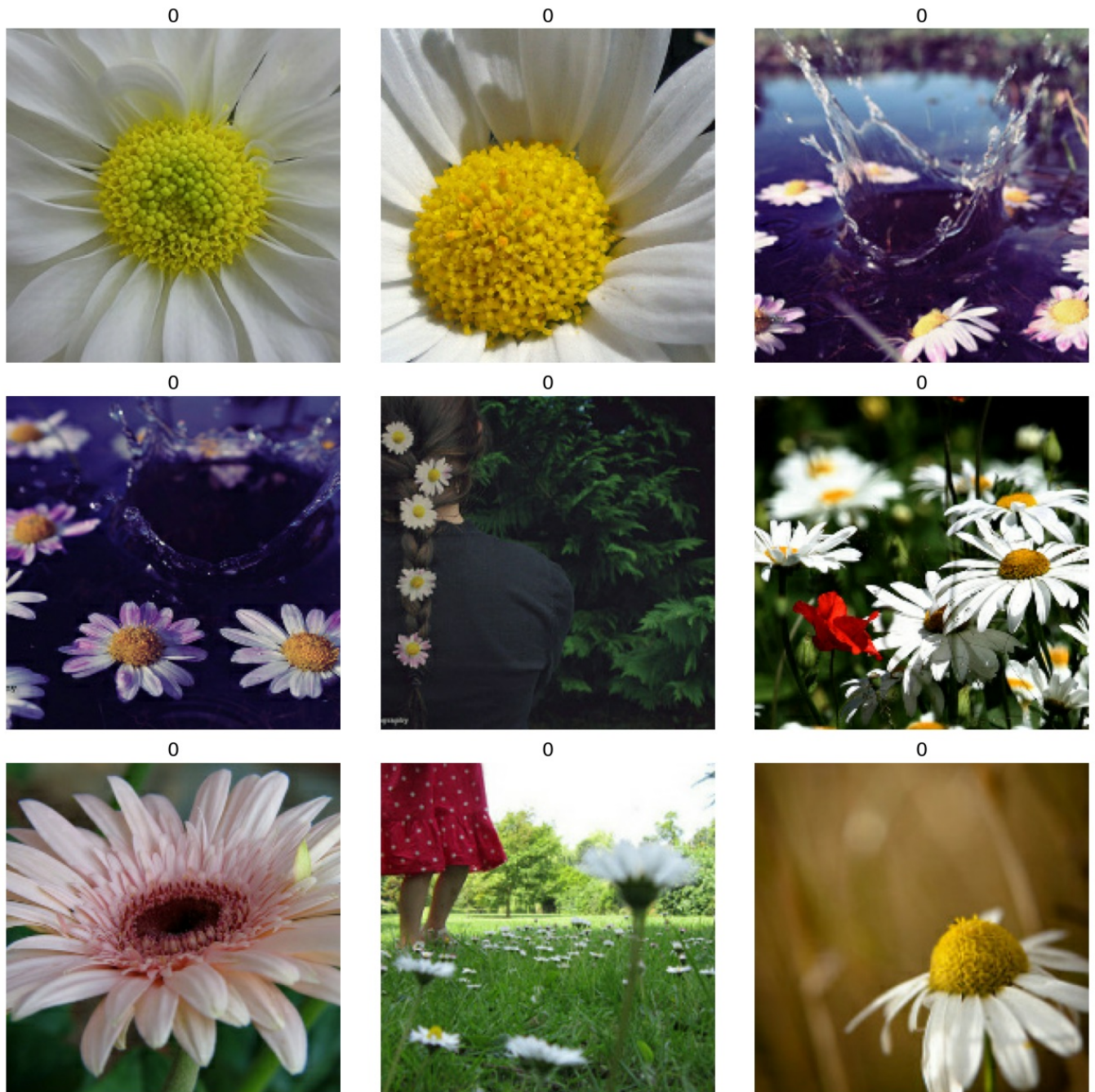
def display_9_images_from_dataset(dataset):
    plt.figure(figsize=(13,13))
    subplot=331
    for i, (image, label) in enumerate(dataset):
        plt.subplot(subplot)
        plt.axis('off')
        plt.imshow(image.numpy().astype(np.uint8))
        plt.title(str(label.numpy()), fontsize=16)
        # plt.title(label.numpy().decode(), fontsize=16)
        subplot += 1
        if i==8:
            break
    plt.tight_layout()
    plt.subplots_adjust(wspace=0.1, hspace=0.1)
    plt.show()

filenamees = tf.io.gfile.glob(GCS_OUTPUT + "/*.tfrec")
DatasetDec = load_dataset(filenamees)
display_9_images_from_dataset(DatasetDec)
```

Writing spark_write_tfrec.py

```
In [ ]: %run spark_write_tfrec.py #We run the script locally on spark
```

Tensorflow version 2.15.0
Writing TFRecords



<Figure size 640x480 with 0 Axes>

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

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

i) Single machine cluster

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

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

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

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

```
In [ ]: #In this code we create a cluster in google cloud with a single machine using the maximal SSD size (100) and 8
CLUSTER = '{}-cluster'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0'
!gcloud dataproc clusters create $CLUSTER \
  --region $REGION \
  --bucket $PROJECT-storage \
```



```
--image-version 1.5-ubuntu18 --single-node \
--master-machine-type n1-standard-8 \
--master-boot-disk-type pd-ssd --master-boot-disk-size 100\
--initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
--metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib'
```

#Reference link for code: <https://cloud.google.com/dataproc/docs/guides/create-cluster>

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/8c9777f5-feba-342a-b726-49320764b1b3].

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting 'https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453'.

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.

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See <https://cloud.google.com/storage/docs/uniform-bucket-level-access>.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-cluster] Cluster placed in zone [us-west1-c].

In []: !gcloud dataproc clusters describe \$CLUSTER --region \$REGION #This is the main information of the single machine

```
clusterName: daring-healer-421511-cluster
clusterUuid: bbfc6d78-16be-44aa-b3ee-b7f62780dc7a
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
    metadata:
      PIP_PACKAGES: tensorflow==2.4.0 matplotlib
    networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
      - https://www.googleapis.com/auth/monitoring.write
    zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c
  initializationActions:
    - executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubu18-20
230909-165100-rc01
    instanceNames:
      - daring-healer-421511-cluster-m
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machineTypes/n1-standard-8
    minCpuPlatform: AUTOMATIC
    numInstances: 1
    preemptibility: NON_PREEMPTIBLE
  softwareConfig:
    imageVersion: 1.5.90-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: daring-healer-421511-cluster-m:8050
      hdfs:dfs.namenode.secondary.http-address: 0.0.0.0:9868
```

```
hdfs:dfs.namenode.secondary.https-address: 0.0.0.0:9869
hdfs:dfs.namenode.service.handler.count: '10'
hdfs:dfs.namenode.servicerpc-address: daring-healer-421511-cluster-m:8051
hive:hive.fetch.task.conversion: none
mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '4000'
mapred:mapreduce.job.maps: '21'
mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
mapred:mapreduce.job.reduces: '7'
mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLevelDbS
tateStoreService
mapred:mapreduce.map.cpu.vcores: '1'
mapred:mapreduce.map.java.opts: -Xmx2457m
mapred:mapreduce.map.memory.mb: '3072'
mapred:mapreduce.reduce.cpu.vcores: '1'
mapred:mapreduce.reduce.java.opts: -Xmx2457m
mapred:mapreduce.reduce.memory.mb: '3072'
mapred:mapreduce.task.io.sort.mb: '256'
mapred:yarn.app.mapreduce.am.command-opts: -Xmx2457m
mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
mapred:yarn.app.mapreduce.am.resource.mb: '3072'
spark-env:SPARK_DAEMON_MEMORY: 4000m
spark:spark.driver.maxResultSize: 3840m
spark:spark.driver.memory: 7680m
spark:spark.executor.cores: '4'
spark:spark.executor.instances: '2'
spark:spark.executor.memory: 11171m
spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
spark:spark.extraListeners: com.google.cloud.spark.performance.DataprocMetricsListener
spark:spark.scheduler.mode: FAIR
spark:spark.sql.cbo.enabled: 'true'
spark:spark.ui.port: '0'
spark:spark.yarn.am.memory: 640m
yarn-env:YARN_NODEMANAGER_HEAPSIZE: '4000'
yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '4000'
yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '4000'
yarn:yarn.nodemanager.address: 0.0.0.0:8026
yarn:yarn.nodemanager.resource.cpu-vcores: '8'
yarn:yarn.nodemanager.resource.memory-mb: '24576'
yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
yarn:yarn.scheduler.maximum-allocation-mb: '24576'
yarn:yarn.scheduler.minimum-allocation-mb: '1024'
tempBucket: dataproc-temp-us-west1-869525608453-ejfnpe5x
labels:
goog-dataproc-autozone: enabled
goog-dataproc-cluster-name: daring-healer-421511-cluster
goog-dataproc-cluster-uuid: bbfc6d78-16be-44aa-b3ee-b7f62780dc7a
goog-dataproc-location: us-west1
metrics:
hdfsMetrics:
dfs-blocks-corrupt: '0'
dfs-blocks-default-replication-factor: '1'
dfs-blocks-missing: '0'
dfs-blocks-missing-repl-one: '0'
dfs-blocks-pending-deletion: '0'
dfs-blocks-under-replication: '0'
dfs-capacity-present: '92799987712'
dfs-capacity-remaining: '92799963136'
dfs-capacity-total: '103865303040'
dfs-capacity-used: '24576'
dfs-nodes-decommissioned: '0'
dfs-nodes-decommissioning: '0'
dfs-nodes-running: '1'
yarnMetrics:
yarn-apps-completed: '0'
yarn-apps-failed: '0'
yarn-apps-killed: '0'
yarn-apps-pending: '0'
yarn-apps-running: '0'
yarn-apps-submitted: '0'
yarn-containers-allocated: '0'
yarn-containers-pending: '0'
yarn-containers-reserved: '0'
yarn-memory-mb-allocated: '0'
yarn-memory-mb-available: '24576'
yarn-memory-mb-pending: '0'
yarn-memory-mb-reserved: '0'
yarn-memory-mb-total: '24576'
yarn-nodes-active: '1'
yarn-nodes-decommissioned: '0'
yarn-nodes-decommissioning: '0'
yarn-nodes-lost: '0'
yarn-nodes-new: '0'
yarn-nodes-rebooted: '0'
yarn-nodes-shutdown: '0'
yarn-nodes-unhealthy: '0'
yarn-vcores-allocated: '0'
yarn-vcores-available: '8'
yarn-vcores-pending: '0'
yarn-vcores-reserved: '0'
```

```
    yarn-vcores-total: '8'
projectId: daring-healer-421511
status:
  state: RUNNING
  stateStartTime: '2024-05-03T14:43:08.195104Z'
statusHistory:
- state: CREATING
  stateStartTime: '2024-05-03T14:40:09.247203Z'
```

Run the script in the cloud and test the output.

```
In [ ]: !gcloud dataproc jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_write_tfrec.py
        %time
        #We submit the python script ('spark_write_tfrec.py') which we created above as a pyspark job to the cluster th
```

```

Job [1c9eb9587aaf4ea483df1b2951b88af1] submitted.
Waiting for job output...
2024-05-03 10:23:50.231590: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dlderror: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: :/usr/lib/hadoop/lib/native
2024-05-03 10:23:50.231632: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dlerror if you do not have a GPU set up on your machine.
Tensorflow version 2.4.0
24/05/03 10:23:54 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 10:23:54 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 10:23:54 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 10:23:54 INFO org.spark_project.jetty.util.log: Logging initialized @11232ms to org.spark_project.jetty.util.log.Slf4jLog
24/05/03 10:23:54 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_382-b05
24/05/03 10:23:54 INFO org.spark_project.jetty.server.Server: Started @11358ms
24/05/03 10:23:54 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@380d866e{HTTP/1.1, (http/1.1)}{0.0.0.0:33907}
24/05/03 10:23:55 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-421511-cluster-m/10.138.0.24:8032
24/05/03 10:23:56 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at daring-healer-421511-cluster-m/10.138.0.24:10200
24/05/03 10:23:56 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 10:23:56 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 10:23:56 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory -mb, units = Mi, type = COUNTABLE
24/05/03 10:23:56 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
24/05/03 10:23:59 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714731723280_0001
Writing TFRecords
2024-05-03 10:24:09.241539: I tensorflow/compiler/jit/xla_cpu_device.cc:41] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 10:24:09.241879: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcuda.so.1'; dlderror: libcuda.so.1: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: :/usr/lib/hadoop/lib/native
2024-05-03 10:24:09.241909: W tensorflow/stream_executor/cuda/cuda_driver.cc:326] failed call to cuInit: UNKNOWN ERROR (303)
2024-05-03 10:24:09.241944: I tensorflow/stream_executor/cuda/cuda_diagnostics.cc:156] kernel driver does not appear to be running on this host (daring-healer-421511-cluster-m): /proc/driver/nvidia/version does not exist
2024-05-03 10:24:09.243069: I tensorflow/compiler/jit/xla_gpu_device.cc:99] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 10:24:09.385526: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:116] None of the MLIR optimization passes are enabled (registered 2)
2024-05-03 10:24:09.386110: I tensorflow/core/platform/profile_utils/cpu_utils.cc:112] CPU Frequency: 2200170000 Hz
24/05/03 10:24:11 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@380d866e{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
Job [1c9eb9587aaf4ea483df1b2951b88af1] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/c1a54c08-bd2b-46c7-82c7-f2b3332be0a3/jobs/1c9eb9587aaf4ea483df1b2951b88af1/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/c1a54c08-bd2b-46c7-82c7-f2b3332be0a3/jobs/1c9eb9587aaf4ea483df1b2951b88af1/driveroutput
jobUuid: ad1d1035-1936-3399-a2bc-9e50ece128a0
placement:
  clusterName: daring-healer-421511-cluster
  clusterUuid: c1a54c08-bd2b-46c7-82c7-f2b3332be0a3
pysparkJob:
  mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/c1a54c08-bd2b-46c7-82c7-f2b3332be0a3/jobs/1c9eb9587aaf4ea483df1b2951b88af1/staging/spark_write_tfrec.py
reference:
  jobId: 1c9eb9587aaf4ea483df1b2951b88af1
  projectId: daring-healer-421511
status:
  state: DONE
  stateStartTime: '2024-05-03T10:24:15.184231Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-03T10:23:41.595317Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-03T10:23:41.631552Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-03T10:23:41.904189Z'
yarnApplications:
- name: spark_write_tfrec.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://daring-healer-421511-cluster-m:8088/proxy/application_1714731723280_0001/
CPU times: user 4 µs, sys: 1 µs, total: 5 µs
Wall time: 9.78 µs

```

```

In [ ]: CLUSTER = '{}-cluster'.format(PROJECT)
        REGION = 'us-west1'
        !gcloud dataproc clusters delete $CLUSTER --region $REGION
        #Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c

```

The cluster 'daring-healer-421511-cluster' and all attached disks will be deleted.

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

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/db4513c4-e753-3ed7-8ea1-54f2c190d8db].

Deleted [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-cluster].

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%)

```
In [ ]: #In this code cell we we create a cluster in google cloud with the maximal SSD size (100) and 7 worker nodes wi
CLUSTER = '{}-maximalcluster7'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0 matplotlib'

!gcloud dataproc clusters create $CLUSTER \
  --bucket $PROJECT-storage \
  --image-version 1.5-ubuntu18 \
  --master-machine-type n1-standard-1 \
  --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \
  --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-disk-size 100 \
  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
  --metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib'
#Reference link for code: https://cloud.google.com/dataproc/docs/guides/create-cluster
```

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/417b8a18-0c87-3066-8688-14120e5e2674].

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

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting 'https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453'.

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

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

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See <https://cloud.google.com/storage/docs/uniform-bucket-level-access>.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-maximalcluster7] Cluster placed in zone [us-west1-c].

```
In [ ]: !gcloud dataproc clusters describe $CLUSTER --region $REGION #This is the main information of the maximal machine
clusterName: daring-healer-421511-maximalcluster7
clusterUuid: 3774a390-9e84-4997-9102-8e033a20bf5d
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
  metadata:
    PIP_PACKAGES: tensorflow==2.4.0 matplotlib
  networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
- https://www.googleapis.com/auth/monitoring.write

zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c

initializationActions:

- executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubuntu18-20230909-165100-rc01

instanceNames:

- daring-healer-421511-maximalcluster7-m

machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machineTypes/n1-standard-1

minCpuPlatform: AUTOMATIC

numInstances: 1

preemptibility: NON_PREEMPTIBLE

softwareConfig:

imageVersion: 1.5.90-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
- distcp:mapreduce.map.java.opts: -Xmx576m
- distcp:mapreduce.map.memory.mb: '768'
- distcp:mapreduce.reduce.java.opts: -Xmx576m
- distcp:mapreduce.reduce.memory.mb: '768'
- 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: '60'
- 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: daring-healer-421511-maximalcluster7-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: '30'
- hdfs:dfs.namenode.servicerpc.address: daring-healer-421511-maximalcluster7-m:8051
- hive:hive.fetch.task.conversion: none
- mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '1000'
- mapred:mapreduce.job.maps: '60'
- mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
- mapred:mapreduce.job.reduces: '7'
- mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLevelDBStateStoreService
- mapred:mapreduce.map.cpu.vcores: '1'
- mapred:mapreduce.map.java.opts: -Xmx819m
- mapred:mapreduce.map.memory.mb: '1024'
- mapred:mapreduce.reduce.cpu.vcores: '1'
- mapred:mapreduce.reduce.java.opts: -Xmx1638m
- mapred:mapreduce.reduce.memory.mb: '2048'
- mapred:mapreduce.task.io.sort.mb: '256'
- mapred:yarn.app.mapreduce.am.command-opts: -Xmx819m
- mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
- mapred:yarn.app.mapreduce.am.resource.mb: '1024'
- spark-env:SPARK_DAEMON_MEMORY: 1000m
- spark:spark.driver.maxResultSize: 480m
- spark:spark.driver.memory: 960m
- spark:spark.executor.cores: '1'
- spark:spark.executor.instances: '2'
- spark:spark.executor.memory: 2688m
- spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
- spark:spark.extraListeners: com.google.cloud.spark.performance.DataprocMetricsListener
- spark:spark.scheduler.mode: FAIR
- spark:spark.sql.cbo.enabled: 'true'
- spark:spark.ui.port: '0'
- spark:spark.yarn.am.memory: 640m
- yarn-env:YARN_NODEMANAGER_HEAPSIZE: '1000'
- yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '1000'
- yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '1000'
- yarn:yarn.nodemanager.address: 0.0.0.0:8026
- yarn:yarn.nodemanager.resource.cpu-vcores: '1'
- yarn:yarn.nodemanager.resource.memory-mb: '3072'
- yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
- yarn:yarn.scheduler.maximum-allocation-mb: '3072'
- yarn:yarn.scheduler.minimum-allocation-mb: '256'

tempBucket: dataproc-temp-us-west1-869525608453-ejfnpe5x

workerConfig:

diskConfig:


```

    bootDiskSizeGb: 100
    bootDiskType: pd-standard
    imageUri: https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-5-ubu18-20
230909-165100-rc01
    instanceNames:
    - daring-healer-421511-maximalcluster7-w-0
    - daring-healer-421511-maximalcluster7-w-1
    - daring-healer-421511-maximalcluster7-w-2
    - daring-healer-421511-maximalcluster7-w-3
    - daring-healer-421511-maximalcluster7-w-4
    - daring-healer-421511-maximalcluster7-w-5
    - daring-healer-421511-maximalcluster7-w-6
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machin
eTypes/n1-standard-1
    minCpuPlatform: AUTOMATIC
    numInstances: 7
    preemptibility: NON_PREEMPTIBLE
labels:
    goog-dataproc-autozone: enabled
    goog-dataproc-cluster-name: daring-healer-421511-maximalcluster7
    goog-dataproc-cluster-uuid: 3774a390-9e84-4997-9102-8e033a20bf5d
    goog-dataproc-location: us-west1
projectId: daring-healer-421511
status:
    state: RUNNING
    stateStartTime: '2024-05-03T10:37:45.555871Z'
statusHistory:
- state: CREATING
  stateStartTime: '2024-05-03T10:33:06.890157Z'

```

```

In [ ]: #We run the script in the cloud to test the output of the script
!gcloud dataproc jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_write_tfrec.py
%time
#We submit the python script ('spark_write_tfrec.py') which we created above as a pyspark job to the cluster th

```

```

Job [eccf04a36ba14d34b51dd455c1b14992] submitted.
Waiting for job output...
2024-05-03 10:38:03.407410: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dLError: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 10:38:03.407568: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dLError if you do not have a GPU set up on your machine.
Tensorflow version 2.4.0
24/05/03 10:38:08 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 10:38:08 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 10:38:08 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 10:38:08 INFO org.spark_project.jetty.util.log: Logging initialized @13639ms to org.spark_project.jetty.util.log.Slf4jLog
24/05/03 10:38:08 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_382-b05
24/05/03 10:38:08 INFO org.spark_project.jetty.server.Server: Started @13871ms
24/05/03 10:38:08 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@45915c69{HTTP/1.1, (http/1.1){0.0.0.0:45707}}
24/05/03 10:38:10 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-421511-maximalcluster7-m/10.138.0.27:8032
24/05/03 10:38:11 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at daring-healer-421511-maximalcluster7-m/10.138.0.27:10200
24/05/03 10:38:11 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 10:38:11 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 10:38:11 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory -mb, units = Mi, type = COUNTABLE
24/05/03 10:38:11 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
24/05/03 10:38:16 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714732543370_0001
Writing TFRecords
2024-05-03 10:38:34.435846: I tensorflow/compiler/jit/xla_cpu_device.cc:41] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 10:38:34.438248: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcuda.so.1'; dLError: libcuda.so.1: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 10:38:34.438489: W tensorflow/stream_executor/cuda/cuda_driver.cc:326] failed call to cuInit: UNKNOWN ERROR (303)
2024-05-03 10:38:34.438566: I tensorflow/stream_executor/cuda/cuda_diagnostics.cc:156] kernel driver does not appear to be running on this host (daring-healer-421511-maximalcluster7-m): /proc/driver/nvidia/version does not exist
2024-05-03 10:38:34.441707: I tensorflow/compiler/jit/xla_gpu_device.cc:99] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 10:38:34.607747: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:116] None of the MLIR optimization passes are enabled (registered 2)
2024-05-03 10:38:34.611708: I tensorflow/core/platform/profile_utils/cpu_utils.cc:112] CPU Frequency: 2199995000 Hz
24/05/03 10:38:36 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@45915c69{HTTP/1.1, (http/1.1){0.0.0.0:0}}
Job [eccf04a36ba14d34b51dd455c1b14992] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/3774a390-9e84-4997-9102-8e033a20bf5d/jobs/eccf04a36ba14d34b51dd455c1b14992/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/3774a390-9e84-4997-9102-8e033a20bf5d/jobs/eccf04a36ba14d34b51dd455c1b14992/driveroutput
jobUuid: b5709024-1cf6-3814-91a2-a60770ba687c
placement:
  clusterName: daring-healer-421511-maximalcluster7
  clusterUuid: 3774a390-9e84-4997-9102-8e033a20bf5d
pysparkJob:
  mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/3774a390-9e84-4997-9102-8e033a20bf5d/jobs/eccf04a36ba14d34b51dd455c1b14992/staging/spark_write_tfrec.py
reference:
  jobId: eccf04a36ba14d34b51dd455c1b14992
  projectId: daring-healer-421511
status:
  state: DONE
  stateStartTime: '2024-05-03T10:38:37.377671Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-03T10:37:51.316818Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-03T10:37:51.357006Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-03T10:37:51.729851Z'
yarnApplications:
- name: spark_write_tfrec.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://daring-healer-421511-maximalcluster7-m:8088/proxy/application_1714732543370_0001/
CPU times: user 3 µs, sys: 0 ns, total: 3 µs
Wall time: 7.87 µs

```

```

In [ ]: CLUSTER = '{}-maximalcluster7'.format(PROJECT)
        REGION = 'us-west1'
        !gcloud dataproc clusters delete $CLUSTER --region $REGION
#Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c

```

The cluster 'daring-healer-421511-maximalcluster3' and all attached disks will be deleted.

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

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/a60de350-2357-347e-b122-078365bc2121].
Deleted [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-maximalcluster3].

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

i) Improve parallelisation

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

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

ii) Experiment with cluster configurations.

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

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

Write the code below and your answers in the report.

```
In [ ]: #SUBTASK i)
```

```
In [ ]: #For this task, We first create a cluster with 8 machines with 1 master with 7 workers type n1 machine using t
CLUSTER = '{}-cluster-task1d-7workers'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0'
!gcloud dataproc clusters create $CLUSTER \
  --region $REGION \
  --bucket $PROJECT-storage \
  --image-version 1.5-ubuntu18 \
  --master-machine-type n1-standard-8 \
  --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \
  --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-disk-size 100 \
  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
  --metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib'

#This is the cluster that is going to be used to test parallelization.

#Reference link for code: https://cloud.google.com/dataproc/docs/guides/create-cluster
```

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/01117f6d-e733-35b0-8974-a3d4d76aa9bb].

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

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting 'https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453'.

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

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

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See https://cloud.google.com/storage/docs/uniform-bucket-level-access.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-cluster-task1d-7workers] Cluster placed in zone [us-west1-c].

```
In [ ]: %%writefile spark_write_tfrec_parallel.py
#We improve parallelization
#We first import the necessary libraries
import os
```

```

os.environ['PROTOCOL_BUFFERS_PYTHON_IMPLEMENTATION'] = 'python'
import os, sys, math
import numpy as np
#import scipy as sp
#import scipy.stats
import time
import string
import datetime
import random
from matplotlib import pyplot as plt
import tensorflow as tf
print("Tensorflow version " + tf.__version__)
import pickle
import pyspark
from pyspark.sql import SQLContext
from pyspark.sql import Row

#We define the variables that will be needed for this project.
GCS_PATTERN = 'gs://flowers-public/*/*.jpg' #glob pattern for input files.
PROJECT = 'daring-healer-421511'
BUCKET = 'gs://{}-storage'.format(PROJECT) # This is the bucket storage.
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # We prefix for output file names.
PARTITIONS = 16 # This is the number of partitions that will be used later.
TARGET_SIZE = [192, 192]

#Task 1a) - i) Copy over the mapping functions and adapt the resizing and recompression functions to Spark (only)
def decode_jpeg_and_label(filepath):
    #This function reads and decodes the JPEG image from the file.
    bits = tf.io.read_file(filepath)
    image = tf.image.decode_jpeg(bits) #We decode the read data as JPEG images.
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/') #We extract the label (flower name) of the file
    label2 = label.values[-2]
    return image, label2

def resize_and_crop_image(data): #We specify one argument
    image, label = data #The input is basically a tuple containing the image data and the label.
    #We determine the dimensions of the input images and the target dimensions.
    w = tf.shape(image)[0]
    h = tf.shape(image)[1]
    tw = TARGET_SIZE[1] #Target Width
    th = TARGET_SIZE[0] #Target Height
    #Calculate resize criterion to maintain the aspect ratio.
    resize_crit = (w * th) / (h * tw)
    #We resize the images to fill the target dimensions but also preserving the aspect ratio.
    image = tf.cond(resize_crit < 1,
                    lambda: tf.image.resize(image, [w*tw/w, h*tw/h]), # if true
                    lambda: tf.image.resize(image, [w*th/h, h*th/h]) # if false
                    )
    nw = tf.shape(image)[0]
    nh = tf.shape(image)[1]
    image = tf.image.crop_to_bounding_box(image, (nw - tw) // 2, (nh - th) // 2, tw, th)
    return (image, label)

def recompress_image(data):
    image, label = data #The input is basically a tuple containing the image data and the label.
    #This function is used to reduce the amount of data, but takes some time
    image = tf.cast(image, tf.uint8) #The image is transformed to a uint8 so the images can be compressed to JP
    image = tf.image.encode_jpeg(image, optimize_size=True, chroma_downsampling=False) #We compress the images
    return (image, label)

#Task 1a) - ii) Replace the TensorFlow Dataset objects with RDDs, starting with an RDD that contains the list of image filenames
from pyspark.sql import SparkSession
sc = pyspark.SparkContext.getOrCreate() #Create or get the existing Spark context
spark = SparkSession.builder.getOrCreate() #Initialize the Spark session

image_paths = tf.io.gfile.glob(GCS_PATTERN) #We retrieve the list of image filenames
### TASK 1d ###
image_rdd = sc.parallelize(image_paths, 16) #We create an RDD from the list of image filenames with 16 partitions

#Task 1a) - iii) Sample the RDD to a smaller number at an appropriate position in the code. Specify a sampling rate
rdd1_sample = image_rdd.sample(False, 0.02) #We first sample the RDD with a sampling factor of 0.02
rdd2_decode_jpeg_and_label = image_rdd.map(decode_jpeg_and_label) #We decode the jpeg and label of the images of the RDD
rdd3_resize_and_crop_image = rdd2_decode_jpeg_and_label.map(resize_and_crop_image) #We resize and crop the images of the RDD
rdd4_recompress_image = rdd3_resize_and_crop_image.map(recompress_image) #We recompress the cropped images of the RDD

#Task 1a) - iv) Then use the functions from above to write the TFRecord files.
#This function is used to create a TensorFlow feature for storing byte strings. This is for the image data.
def _bytestring_feature(list_of_bytestrings):
    return tf.train.Feature(bytes_list=tf.train.BytesList(value=list_of_bytestrings))
#This function is used for storing integer values. This is for the labels.
def _int_feature(list_of_ints):
    return tf.train.Feature(int64_list=tf.train.Int64List(value=list_of_ints))
#This function is used to create a TFRecord from image bytes and a label.
def to_tfrecord(tfrec_filewriter, img_bytes, label): # Create tf data records
    class_num = np.argmax(np.array(CLASSES)==label) # 'roses' => 2 (order defined in CLASSES)
    one_hot_class = np.eye(len(CLASSES))[class_num] # [0, 0, 1, 0, 0] for class #2, roses
    feature = {
        "image": _bytestring_feature([img_bytes]), # one image in the list

```

```

        "class": _int_feature([class_num]) #,      # one class in the list
    }
    return tf.train.Example(features=tf.train.Features(feature=feature))

print("Writing TFRecords")
#This function is used to write TFRecords, which will be applied to partitions of an RDD.
def write_tfrecords(partition_index,partition):
    filename = GCS_OUTPUT + "{}.tfrec".format(partition_index) #We define the filename for this TFRecord file bas
    #We Open a TFRecordWriter object for the filename.
    with tf.io.TFRecordWriter(filename) as out_file:
        for element in partition:
            image=element[0]
            label=element[1]
            example = to_tfrecord(out_file,
                                image.numpy(), # The re-compressed image is already a byte string
                                label.numpy()
                                )
            out_file.write(example.SerializeToString()) #We convert the image and label to TFRecord format and write
    return [filename] #We return the filename of the written TFRecord.

#Task 1a) - v) The code for writing to the TFRecord files needs to be put into a function, that can be applied
#should return the filename, so that you have a list of the created TFRecord files.
rdd5_partitions = rdd4_recompress_image.repartition(PARTITIONS) #We repartition the RDD to have the number of p
rdd1_filenames = rdd4_recompress_image.mapPartitionsWithIndex(write_tfrecords) #We specify an index for each pa

#Task 1b) i) Read from the TFRecord Dataset, using load_dataset and display_9_images_from_dataset to test.

#GCS_OUTPUT = 'gs://flowers-public/tfrecords-jpeg-192x192-2/' # We should remove this line to use our own files
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' #These are the files that we generated above.

def read_tfrecord(example):
    features = {
        "image": tf.io.FixedLenFeature([], tf.string), # tf.string = bytestring (not text string)
        "class": tf.io.FixedLenFeature([], tf.int64) #,      # shape [] means scalar
    }
    # decode the TFRecord
    example = tf.io.parse_single_example(example, features)
    image = tf.image.decode_jpeg(example['image'], channels=3)
    image = tf.reshape(image, [*TARGET_SIZE, 3])
    class_num = example['class']
    return image, class_num

def load_dataset(filenames):
    # read from TFRecords. For optimal performance, read from multiple
    # TFRecord files at once and set the option experimental_deterministic = False
    # to allow order-altering optimizations.
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False

    dataset = tf.data.TFRecordDataset(filenames)
    dataset = dataset.with_options(option_no_order)
    dataset = dataset.map(read_tfrecord)
    return dataset

def display_9_images_from_dataset(dataset):
    plt.figure(figsize=(13,13))
    subplot=331
    for i, (image, label) in enumerate(dataset):
        plt.subplot(subplot)
        plt.axis('off')
        plt.imshow(image.numpy().astype(np.uint8))
        plt.title(str(label.numpy()), fontsize=16)
        # plt.title(label.numpy().decode(), fontsize=16)
        subplot += 1
        if i==8:
            break
    plt.tight_layout()
    plt.subplots_adjust(wspace=0.1, hspace=0.1)
    plt.show()

filenames = tf.io.gfile.glob(GCS_OUTPUT + "*.tfrec")
DatasetDec = load_dataset(filenames)
display_9_images_from_dataset(DatasetDec)

```

Writing spark_write_tfrec_paral.py

In []: !gcloud dataproc clusters describe \$CLUSTER --region \$REGION #We get the main information of the cluster

```

clusterName: daring-healer-421511-cluster-task1d-7workers
clusterUuid: 4085db72-fcba-4bd1-af18-508d2b7a67be
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
  metadata:
    PIP_PACKAGES: tensorflow==2.4.0 matplotlib
  networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
- https://www.googleapis.com/auth/monitoring.write
zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c
initializationActions:
- executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubuntu18-20
230909-165100-rc01
  instanceNames:
    - daring-healer-421511-cluster-task1d-7workers-m
  machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machineTypes/n1-standard-8
  minCpuPlatform: AUTOMATIC
  numInstances: 1
  preemptibility: NON_PREEMPTIBLE
softwareConfig:
  imageVersion: 1.5.90-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
    distcp:mapreduce.map.java.opts: -Xmx576m
    distcp:mapreduce.map.memory.mb: '768'
    distcp:mapreduce.reduce.java.opts: -Xmx576m
    distcp:mapreduce.reduce.memory.mb: '768'
    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: '60'
    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: daring-healer-421511-cluster-task1d-7workers-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: '30'
    hdfs:dfs.namenode.servicerpc-address: daring-healer-421511-cluster-task1d-7workers-m:8051
    hive:hive.fetch.task.conversion: none
    mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '4000'
    mapred:mapreduce.job.maps: '60'
    mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
    mapred:mapreduce.job.reduces: '7'
    mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLevelDbStateStoreService
    mapred:mapreduce.map.cpu.vcores: '1'
    mapred:mapreduce.map.java.opts: -Xmx819m
    mapred:mapreduce.map.memory.mb: '1024'
    mapred:mapreduce.reduce.cpu.vcores: '1'
    mapred:mapreduce.reduce.java.opts: -Xmx1638m
    mapred:mapreduce.reduce.memory.mb: '2048'
    mapred:mapreduce.task.io.sort.mb: '256'
    mapred:yarn.app.mapreduce.am.command.opts: -Xmx819m
    mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
    mapred:yarn.app.mapreduce.am.resource.mb: '1024'
    spark-env:SPARK_DAEMON_MEMORY: 4000m
    spark:spark.driver.maxResultSize: 3840m
    spark:spark.driver.memory: 7680m
    spark:spark.executor.cores: '1'
    spark:spark.executor.instances: '2'
    spark:spark.executor.memory: 2688m
    spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
    spark:spark.extraListeners: com.google.cloud.spark.performance.DataprocMetricsListener
    spark:spark.scheduler.mode: FAIR
    spark:spark.sql.cbo.enabled: 'true'
    spark:spark.ui.port: '0'
    spark:spark.yarn.am.memory: 640m
    yarn-env:YARN_NODEMANAGER_HEAPSIZE: '1000'
    yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '4000'
    yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '4000'
    yarn:yarn.nodemanager.address: 0.0.0.0:8026
    yarn:yarn.nodemanager.resource.cpu-vcores: '1'
    yarn:yarn.nodemanager.resource.memory.mb: '3072'
    yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
    yarn:yarn.scheduler.maximum-allocation-mb: '3072'
    yarn:yarn.scheduler.minimum-allocation-mb: '256'
tempBucket: dataproc-temp-us-west1-869525608453-ejfnpe5x
workerConfig:

```



```

    diskConfig:
      bootDiskSizeGb: 100
      bootDiskType: pd-standard
    imageUri: https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-5-ubu18-20
230909-165100-rc01
    instanceNames:
      - daring-healer-421511-cluster-task1d-7workers-w-0
      - daring-healer-421511-cluster-task1d-7workers-w-1
      - daring-healer-421511-cluster-task1d-7workers-w-2
      - daring-healer-421511-cluster-task1d-7workers-w-3
      - daring-healer-421511-cluster-task1d-7workers-w-4
      - daring-healer-421511-cluster-task1d-7workers-w-5
      - daring-healer-421511-cluster-task1d-7workers-w-6
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machin
eTypes/n1-standard-1
    minCpuPlatform: AUTOMATIC
    numInstances: 7
    preemptibility: NON_PREEMPTIBLE
  labels:
    goog-dataproc-autozone: enabled
    goog-dataproc-cluster-name: daring-healer-421511-cluster-task1d-7workers
    goog-dataproc-cluster-uuid: 4085db72-fcba-4bd1-af18-508d2b7a67be
    goog-dataproc-location: us-west1
  projectId: daring-healer-421511
  status:
    state: RUNNING
    stateStartTime: '2024-05-03T10:50:07.233406Z'
  statusHistory:
    - state: CREATING
      stateStartTime: '2024-05-03T10:46:26.816792Z'

```

```

In [ ]: !gcloud dataproc jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_write_tfrec_paral.py
%time

#We submit the python script as a job to our cluster

```

```

Job [e173d9441aff4f42a9bf94c81018d4f4] submitted.
Waiting for job output...
2024-05-03 10:50:26.828368: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dLError: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 10:50:26.828415: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dLError if you do not have a GPU set up on your machine.
Tensorflow version 2.4.0
24/05/03 10:50:30 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 10:50:30 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 10:50:30 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 10:50:31 INFO org.spark_project.jetty.util.log: Logging initialized @11432ms to org.spark_project.jetty.util.log.Slf4jLog
24/05/03 10:50:31 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_382-b05
24/05/03 10:50:31 INFO org.spark_project.jetty.server.Server: Started @11576ms
24/05/03 10:50:31 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@600eaa78{HTTP/1.1, (http/1.1)}{0.0.0.0:43279}
24/05/03 10:50:32 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-421511-cluster-taskId-7workers-m/10.138.0.32:8032
24/05/03 10:50:33 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at daring-healer-421511-cluster-taskId-7workers-m/10.138.0.32:10200
24/05/03 10:50:33 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 10:50:33 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 10:50:33 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory -mb, units = Mi, type = COUNTABLE
24/05/03 10:50:33 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
24/05/03 10:50:37 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714733279670_0001
Writing TFRecords
2024-05-03 10:50:53.187454: I tensorflow/compiler/jit/xla_cpu_device.cc:41] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 10:50:53.187792: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcuda.so.1'; dLError: libcuda.so.1: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 10:50:53.187818: W tensorflow/stream_executor/cuda/cuda_driver.cc:326] failed call to cuInit: UNKNOWN ERROR (303)
2024-05-03 10:50:53.187847: I tensorflow/stream_executor/cuda/cuda_diagnostics.cc:156] kernel driver does not appear to be running on this host (daring-healer-421511-cluster-taskId-7workers-m): /proc/driver/nvidia/version does not exist
2024-05-03 10:50:53.188967: I tensorflow/compiler/jit/xla_gpu_device.cc:99] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 10:50:53.307246: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:116] None of the MLIR optimization passes are enabled (registered 2)
2024-05-03 10:50:53.307686: I tensorflow/core/platform/profile_utils/cpu_utils.cc:112] CPU Frequency: 2199995000 Hz
24/05/03 10:50:54 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@600eaa78{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
Job [e173d9441aff4f42a9bf94c81018d4f4] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/4085db72-fcba-4bd1-af18-508d2b7a67be/jobs/e173d9441aff4f42a9bf94c81018d4f4/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/4085db72-fcba-4bd1-af18-508d2b7a67be/jobs/e173d9441aff4f42a9bf94c81018d4f4/driveroutput
jobUuid: eaffcb97-d15f-3f07-bcc5-21bbfaeca036
placement:
  clusterName: daring-healer-421511-cluster-taskId-7workers
  clusterUuid: 4085db72-fcba-4bd1-af18-508d2b7a67be
pysparkJob:
  mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/4085db72-fcba-4bd1-af18-508d2b7a67be/jobs/e173d9441aff4f42a9bf94c81018d4f4/staging/spark_write_tfrec_parallel.py
reference:
  jobId: e173d9441aff4f42a9bf94c81018d4f4
  projectId: daring-healer-421511
status:
  state: DONE
  stateStartTime: '2024-05-03T10:50:59.474905Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-03T10:50:17.217414Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-03T10:50:17.253100Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-03T10:50:17.548975Z'
yarnApplications:
- name: spark_write_tfrec_parallel.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://daring-healer-421511-cluster-taskId-7workers-m:8088/proxy/application_1714733279670_0001/
CPU times: user 2 µs, sys: 0 ns, total: 2 µs
Wall time: 5.72 µs

```

```

In [ ]: CLUSTER = '{}-cluster-taskId-7workers'.format(PROJECT)
        REGION = 'us-west1'
        !gcloud dataproc clusters delete $CLUSTER --region $REGION
#Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c

```

The cluster 'daring-healer-421511-cluster-taskId' and all attached disks will be deleted.

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

Waiting on operation [projects/daring-healer-421511/regions/us-central1/operations/e28e3fc9-4b65-39c4-bc2e-bcb609da03bd].

Deleted [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-central1/clusters/daring-healer-421511-cluster-taskId].

```
In [ ]: #SUBTASK ii)
```

```
In [ ]: #Experiment with different cluster configurations
#In this subtask, we aim to test and compare the performance of the cluster created above with 8 machines with
#type n - 2 and 1 cluster with 1 master and 0 workers type n-8). The disk sizes for all clusters are set the sa

#This cluster has 1 master with 3 workers type n - 2 using the max SSD size of 100.
CLUSTER = '{}-cluster-taskId-3workers'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0'
!gcloud dataproc clusters create $CLUSTER \
  --region $REGION \
  --bucket $PROJECT-storage \
  --image-version 1.5-ubuntu18 \
  --master-machine-type n1-standard-8 \
  --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \
  --num-workers 3 --worker-machine-type n1-standard-2 --worker-boot-disk-size 100 \
  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
  --metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib'
#Reference link for code: https://cloud.google.com/dataproc/docs/guides/create-cluster
```

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/86da704f-37d6-35f1-aa76-b968fe224bde].

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of an initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting '<https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453>'.

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

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

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See <https://cloud.google.com/storage/docs/uniform-bucket-level-access>.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-cluster-taskId-3workers] Cluster placed in zone [us-west1-c].

```
In [ ]: #We get information on the cluster
!gcloud dataproc clusters describe $CLUSTER --region $REGION
```

```
clusterName: daring-healer-421511-cluster-taskId-3workers
clusterUuid: c50abefa-bdd7-4df3-8d73-ef0060e0b8f5
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
    metadata:
      PIP_PACKAGES: tensorflow==2.4.0 matplotlib
    networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
      - https://www.googleapis.com/auth/monitoring.write
    zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c
  initializationActions:
    - executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubu18-20230909-165100-rc01
    instanceNames:
```

```

- daring-healer-421511-cluster-task1d-3workers-m
machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machin
eTypes/n1-standard-8
minCpuPlatform: AUTOMATIC
numInstances: 1
preemptibility: NON_PREEMPTIBLE
softwareConfig:
  imageVersion: 1.5.90-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
    distcp:mapreduce.map.java.opts: -Xmx576m
    distcp:mapreduce.map.memory.mb: '768'
    distcp:mapreduce.reduce.java.opts: -Xmx576m
    distcp:mapreduce.reduce.memory.mb: '768'
    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: '40'
    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: daring-healer-421511-cluster-task1d-3workers-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: '20'
    hdfs:dfs.namenode.servicerpc-address: daring-healer-421511-cluster-task1d-3workers-m:8051
    hive:hive.fetch.task.conversion: none
    mapred:env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '4000'
    mapred:mapreduce.job.maps: '24'
    mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
    mapred:mapreduce.job.reduces: '8'
    mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLeveldbS
tateStoreService
    mapred:mapreduce.map.cpu.vcores: '1'
    mapred:mapreduce.map.java.opts: -Xmx1638m
    mapred:mapreduce.map.memory.mb: '2048'
    mapred:mapreduce.reduce.cpu.vcores: '1'
    mapred:mapreduce.reduce.java.opts: -Xmx1638m
    mapred:mapreduce.reduce.memory.mb: '2048'
    mapred:mapreduce.task.io.sort.mb: '256'
    mapred:yarn.app.mapreduce.am.command.opts: -Xmx1638m
    mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
    mapred:yarn.app.mapreduce.am.resource.mb: '2048'
    spark:env:SPARK_DAEMON_MEMORY: 4000m
    spark:spark.driver.maxResultSize: 3840m
    spark:spark.driver.memory: 7680m
    spark:spark.executor.cores: '1'
    spark:spark.executor.instances: '2'
    spark:spark.executor.memory: 2688m
    spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
    spark:spark.extraListeners: com.google.cloud.spark.performance.DataproMetricsListener
    spark:spark.scheduler.mode: FAIR
    spark:spark.sql.cbo.enabled: 'true'
    spark:spark.ui.port: '0'
    spark:spark.yarn.am.memory: 640m
    yarn:env:YARN_NODEMANAGER_HEAPSIZE: '1920'
    yarn:env:YARN_RESOURCEMANAGER_HEAPSIZE: '4000'
    yarn:env:YARN_TIMELINESERVER_HEAPSIZE: '4000'
    yarn:yarn.nodemanager.address: 0.0.0.0:8026
    yarn:yarn.nodemanager.resource.cpu-vcores: '2'
    yarn:yarn.nodemanager.resource.memory-mb: '6144'
    yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
    yarn:yarn.scheduler.maximum-allocation-mb: '6144'
    yarn:yarn.scheduler.minimum-allocation-mb: '512'
tempBucket: dataproc-temp-us-west1-869525608453-ejfnpe5x
workerConfig:
  diskConfig:
    bootDiskSizeGb: 100
    bootDiskType: pd-standard
  imageUri: https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-5-ubu18-20
30909-165100-rc01
instanceNames:
- daring-healer-421511-cluster-task1d-3workers-w-0
- daring-healer-421511-cluster-task1d-3workers-w-1
- daring-healer-421511-cluster-task1d-3workers-w-2
machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machin
eTypes/n1-standard-2
minCpuPlatform: AUTOMATIC
numInstances: 3
preemptibility: NON_PREEMPTIBLE
labels:
  goog-dataproc-autozone: enabled
  goog-dataproc-cluster-name: daring-healer-421511-cluster-task1d-3workers
  goog-dataproc-cluster-uuid: c50abefa-bdd7-4df3-8d73-ef0060e0b8f5
  goog-dataproc-location: us-west1
projectId: daring-healer-421511

```

```
status:
  state: RUNNING
  stateStartTime: '2024-05-03T11:02:21.732121Z'
statusHistory:
- state: CREATING
  stateStartTime: '2024-05-03T10:58:51.411629Z'
```

```
In [ ]: #We submit the job to the cluster
!gcloud dataproc jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_write_tfrec_parallel.py
%time
```

```

Job [f380694063b348f1950d63a4568b9a85] submitted.
Waiting for job output...
2024-05-03 11:02:37.277570: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dLError: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 11:02:37.277620: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dLError if you do not have a GPU set up on your machine.
Tensorflow version 2.4.0
24/05/03 11:02:40 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 11:02:40 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 11:02:40 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 11:02:41 INFO org.spark_project.jetty.util.log: Logging initialized @10563ms to org.spark_project.jetty.util.log.Slf4jLog
24/05/03 11:02:41 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_382-b05
24/05/03 11:02:41 INFO org.spark_project.jetty.server.Server: Started @10700ms
24/05/03 11:02:41 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@77391f2f{HTTP/1.1, (http/1.1)}{0.0.0.0:41511}
24/05/03 11:02:42 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-421511-cluster-taskId-3workers-m/10.138.0.42:8032
24/05/03 11:02:42 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at daring-healer-421511-cluster-taskId-3workers-m/10.138.0.42:10200
24/05/03 11:02:42 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 11:02:42 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 11:02:42 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory -mb, units = Mi, type = COUNTABLE
24/05/03 11:02:42 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
24/05/03 11:02:45 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714734019364_0001
Writing TFRecords
2024-05-03 11:02:58.967230: I tensorflow/compiler/jit/xla_cpu_device.cc:41] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 11:02:58.967545: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcuda.so.1'; dLError: libcuda.so.1: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 11:02:58.967577: W tensorflow/stream_executor/cuda/cuda_driver.cc:326] failed call to cuInit: UNKNOWN ERROR (303)
2024-05-03 11:02:58.967601: I tensorflow/stream_executor/cuda/cuda_diagnostics.cc:156] kernel driver does not appear to be running on this host (daring-healer-421511-cluster-taskId-3workers-m): /proc/driver/nvidia/version does not exist
2024-05-03 11:02:58.968611: I tensorflow/compiler/jit/xla_gpu_device.cc:99] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 11:02:59.085344: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:116] None of the MLIR optimization passes are enabled (registered 2)
2024-05-03 11:02:59.085877: I tensorflow/core/platform/profile_utils/cpu_utils.cc:112] CPU Frequency: 2199995000 Hz
24/05/03 11:03:00 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@77391f2f{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
Job [f380694063b348f1950d63a4568b9a85] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/c50abefa-bdd7-4df3-8d73-ef0060e0b8f5/jobs/f380694063b348f1950d63a4568b9a85/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/c50abefa-bdd7-4df3-8d73-ef0060e0b8f5/jobs/f380694063b348f1950d63a4568b9a85/driveroutput
jobUuid: bbf2efb6-571a-3009-9e0b-154bcc19526e
placement:
  clusterName: daring-healer-421511-cluster-taskId-3workers
  clusterUuid: c50abefa-bdd7-4df3-8d73-ef0060e0b8f5
pysparkJob:
  mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/c50abefa-bdd7-4df3-8d73-ef0060e0b8f5/jobs/f380694063b348f1950d63a4568b9a85/staging/spark_write_tfrec_parallel.py
reference:
  jobId: f380694063b348f1950d63a4568b9a85
  projectId: daring-healer-421511
status:
  state: DONE
  stateStartTime: '2024-05-03T11:03:02.187354Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-03T11:02:29.004432Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-03T11:02:29.053704Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-03T11:02:29.334748Z'
yarnApplications:
- name: spark_write_tfrec_parallel.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://daring-healer-421511-cluster-taskId-3workers-m:8088/proxy/application_1714734019364_0001/
CPU times: user 4 µs, sys: 0 ns, total: 4 µs
Wall time: 8.82 µs

```

```

In [ ]: CLUSTER = '{}-cluster-taskId-3workers'.format(PROJECT)
        REGION = 'us-west1'
        !gcloud dataproc clusters delete $CLUSTER --region $REGION
#Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c

```


The cluster 'daring-healer-421511-cluster-task1d-3vms' and all attached disks will be deleted.

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

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/25c44884-0470-3592-884e-a3bb2c8ba58f].

Deleted [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-cluster-task1d-3vms].

```
In [ ]: #This is the cluster with 1 master and 0 workers type n1-standard-8 using the max SSD size of 100.
CLUSTER = '{}-cluster-task1d-0workers'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0'
!gcloud dataproc clusters create $CLUSTER \
  --region $REGION \
  --bucket $PROJECT-storage \
  --image-version 1.5-ubuntu18 \
  --single-node \
  --master-machine-type n1-standard-8 \
  --master-boot-disk-type pd-ssd --master-boot-disk-size 100 \
  --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
  --metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib'
#Reference link for code: https://cloud.google.com/dataproc/docs/guides/create-cluster
```

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/4d0b5784-02a0-348f-9cb8-6189e677ef18].

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting '<https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453>'.

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.

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See <https://cloud.google.com/storage/docs/uniform-bucket-level-access>.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-cluster-task1d-0workers] Cluster placed in zone [us-west1-c].

```
In [ ]: #We get information on the cluster
!gcloud dataproc clusters describe $CLUSTER --region $REGION
```

```
clusterName: daring-healer-421511-cluster-task1d-0workers
clusterUuid: 830a382c-b5c5-4426-834e-622eb4cdd050
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
  metadata:
    PIP_PACKAGES: tensorflow==2.4.0 matplotlib
  networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
  - https://www.googleapis.com/auth/monitoring.write
  zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c
  initializationActions:
  - executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubu18-20230909-165100-rc01
    instanceNames:
    - daring-healer-421511-cluster-task1d-0workers-m
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machineTypes/n1-standard-8
    minCpuPlatform: AUTOMATIC
    numInstances: 1
    preemptibility: NON_PREEMPTIBLE
  softwareConfig:
    imageVersion: 1.5.90-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
dataproц:dataproц.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: daring-healer-421511-cluster-task1d-0workers-m:8050
hdfs:dfs.namenode.secondary.http-address: 0.0.0.0:9868
hdfs:dfs.namenode.secondary.https-address: 0.0.0.0:9869
hdfs:dfs.namenode.service.handler.count: '10'
hdfs:dfs.namenode.servicerpc-address: daring-healer-421511-cluster-task1d-0workers-m:8051
hive:hive.fetch.task.conversion: none
mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '4000'
mapred:mapreduce.job.maps: '21'
mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
mapred:mapreduce.job.reduces: '7'
mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLevelDBS
tateStoreService
mapred:mapreduce.map.cpu.vcores: '1'
mapred:mapreduce.map.java.opts: -Xmx2457m
mapred:mapreduce.map.memory.mb: '3072'
mapred:mapreduce.reduce.cpu.vcores: '1'
mapred:mapreduce.reduce.java.opts: -Xmx2457m
mapred:mapreduce.reduce.memory.mb: '3072'
mapred:mapreduce.task.io.sort.mb: '256'
mapred:yarn.app.mapreduce.am.command-opts: -Xmx2457m
mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
mapred:yarn.app.mapreduce.am.resource.mb: '3072'
spark-env:SPARK_DAEMON_MEMORY: 4000m
spark:spark.driver.maxResultSize: 3840m
spark:spark.driver.memory: 7680m
spark:spark.executor.cores: '4'
spark:spark.executor.instances: '2'
spark:spark.executor.memory: 11171m
spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
spark:spark.extraListeners: com.google.cloud.spark.performance.DataproцMetricsListener
spark:spark.scheduler.mode: FAIR
spark:spark.sql.cbo.enabled: 'true'
spark:spark.ui.port: '0'
spark:spark.yarn.am.memory: 640m
yarn-env:YARN_NODEMANAGER_HEAPSIZE: '4000'
yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '4000'
yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '4000'
yarn:yarn.nodemanager.address: 0.0.0.0:8026
yarn:yarn.nodemanager.resource.cpu-vcores: '8'
yarn:yarn.nodemanager.resource.memory-mb: '24576'
yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
yarn:yarn.scheduler.maximum-allocation-mb: '24576'
yarn:yarn.scheduler.minimum-allocation-mb: '1024'
tempBucket: dataproц-temp-us-west1-869525608453-ejfnpe5x
labels:
  goog-dataproц-autozone: enabled
  goog-dataproц-cluster-name: daring-healer-421511-cluster-task1d-0workers
  goog-dataproц-cluster-uuid: 830a382c-b5c5-4426-834e-622eb4cdd050
  goog-dataproц-location: us-west1
projectId: daring-healer-421511
status:
  state: RUNNING
  stateStartTime: '2024-05-03T11:13:19.818370Z'
statusHistory:
  - state: CREATING
    stateStartTime: '2024-05-03T11:10:06.777603Z'

```

```

In [ ]: #We submit the job to the cluster
!gcloud dataproц jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_write_tfrec_parallel.py
%time

```

```

Job [cfe4265100f34c1190312be7992be102] submitted.
Waiting for job output...
2024-05-03 11:13:39.247642: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcudart.so.11.0'; dLError: libcudart.so.11.0: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 11:13:39.247702: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dLError if you do not have a GPU set up on your machine.
Tensorflow version 2.4.0
24/05/03 11:13:43 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 11:13:43 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 11:13:43 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 11:13:43 INFO org.spark_project.jetty.util.log: Logging initialized @11482ms to org.spark_project.jetty.util.log.Slf4jLog
24/05/03 11:13:43 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_382-b05
24/05/03 11:13:43 INFO org.spark_project.jetty.server.Server: Started @11632ms
24/05/03 11:13:43 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@4a63dd5d{HTTP/1.1, (http/1.1)}{0.0.0.0:46221}
24/05/03 11:13:45 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-421511-cluster-taskId-0workers-m/10.138.0.45:8032
24/05/03 11:13:45 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at daring-healer-421511-cluster-taskId-0workers-m/10.138.0.45:10200
24/05/03 11:13:45 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 11:13:45 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 11:13:45 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory -mb, units = Mi, type = COUNTABLE
24/05/03 11:13:45 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
24/05/03 11:13:49 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1714734718789_0001
Writing TFRecords
2024-05-03 11:14:00.851941: I tensorflow/compiler/jit/xla_cpu_device.cc:41] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 11:14:00.852373: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic library 'libcuda.so.1'; dLError: libcuda.so.1: cannot open shared object file: No such file or directory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 11:14:00.852414: W tensorflow/stream_executor/cuda/cuda_driver.cc:326] failed call to cuInit: UNKNOWN ERROR (303)
2024-05-03 11:14:00.852460: I tensorflow/stream_executor/cuda/cuda_diagnostics.cc:156] kernel driver does not appear to be running on this host (daring-healer-421511-cluster-taskId-0workers-m): /proc/driver/nvidia/version does not exist
2024-05-03 11:14:00.853899: I tensorflow/compiler/jit/xla_gpu_device.cc:99] Not creating XLA devices, tf_xla_enable_xla_devices not set
2024-05-03 11:14:00.995579: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:116] None of the MLIR optimization passes are enabled (registered 2)
2024-05-03 11:14:00.996035: I tensorflow/core/platform/profile_utils/cpu_utils.cc:112] CPU Frequency: 2200165000 Hz
24/05/03 11:14:02 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@4a63dd5d{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
Job [cfe4265100f34c1190312be7992be102] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/830a382c-b5c5-4426-834e-622eb4cdd050/jobs/cfe4265100f34c1190312be7992be102/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/830a382c-b5c5-4426-834e-622eb4cdd050/jobs/cfe4265100f34c1190312be7992be102/driveroutput
jobUuid: e0851ff2-fea4-3015-9428-35a3a45039d0
placement:
  clusterName: daring-healer-421511-cluster-taskId-0workers
  clusterUuid: 830a382c-b5c5-4426-834e-622eb4cdd050
pysparkJob:
  mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/830a382c-b5c5-4426-834e-622eb4cdd050/jobs/cfe4265100f34c1190312be7992be102/staging/spark_write_tfrec_parallel.py
reference:
  jobId: cfe4265100f34c1190312be7992be102
  projectId: daring-healer-421511
status:
  state: DONE
  stateStartTime: '2024-05-03T11:14:04.262895Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-03T11:13:29.718852Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-03T11:13:29.752786Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-03T11:13:30.186087Z'
yarnApplications:
- name: spark_write_tfrec_parallel.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://daring-healer-421511-cluster-taskId-0workers-m:8088/proxy/application_1714734718789_0001/
CPU times: user 6 µs, sys: 1e+03 ns, total: 7 µs
Wall time: 22.4 µs

```

```

In [ ]: CLUSTER = '{}-cluster-taskId-0workers'.format(PROJECT)
        REGION = 'us-west1'
        !gcloud dataproc clusters delete $CLUSTER --region $REGION
#Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c

```

The cluster 'daring-healer-421511-cluster-task1d-6vms' and all attached disks will be deleted.

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

Waiting on operation [projects/daring-healer-421511/regions/us-central1/operations/bf22992a-d60e-37b9-851f-2a1e6bf233b8].

Deleted [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-central1/clusters/daring-healer-421511-cluster-task1d-6vms].

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.

```
In [ ]: # Here are some useful values for testing your code, use higher values later for actually testing throughput
batch_sizes = [2,4]
batch_numbers = [3,6]
repetitions = [1]

def time_configs(dataset, batch_sizes, batch_numbers, repetitions):
    dims = [len(batch_sizes), len(batch_numbers), len(repetitions)]
    print(dims)
    results = np.zeros(dims)
    params = np.zeros(dims + [3])
    print( results.shape )
    with open("/dev/null", mode='w') as null_file: # for printing the output without showing it
        tt = time.time() # for overall time taking
        for bsi, bs in enumerate(batch_sizes):
            for dsi, ds in enumerate(batch_numbers):
                batched_dataset = dataset.batch(bs)
                timing_set = batched_dataset.take(ds)
                for ri, rep in enumerate(repetitions):
                    print("bs: {}, ds: {}, rep: {}".format(bs, ds, rep))
                    t0 = time.time()
                    for image, label in timing_set:
                        #print("Image batch shape {}".format(image.numpy().shape),
                        print("Image batch shape {}, {}".format(image.numpy().shape,
                            [str(lbl) for lbl in label.numpy()] ), null_file)
                    td = time.time() - t0 # duration for reading images
                    results[bsi, dsi, ri] = ( bs * ds ) / td
                    params[bsi, dsi, ri] = [ bs, ds, rep ]
    print("total time: "+str(time.time()-tt))
    return results, params
```

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

```
In [ ]: [res, par] = time_configs(dataset4, batch_sizes, batch_numbers, repetitions)
print(res)
print(par)

print("=====")

[res, par] = time_configs(datasetTfrec, batch_sizes, batch_numbers, repetitions)
print(res)
print(par)

[2, 2, 1]
(2, 2, 1)
bs: 2, ds: 3, rep: 1
```

```
Image batch shape (2,), ["b'dandelion'", "b'sunflowers'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2,), ["b'tulips'", "b'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2,), ["b'sunflowers'", "b'sunflowers'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
bs: 2, ds: 6, rep: 1
Image batch shape (2,), ["b'dandelion'", "b'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2,), ["b'dandelion'", "b'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2,), ["b'roses'", "b'roses'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2,), ["b'sunflowers'", "b'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (2,), ["b'daisy'", "b'dandelion'"]) <_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'>
bs: 4, ds: 3, rep: 1
Image batch shape (4,), ["b'daisy'", "b'daisy'", "b'roses'", "b'tulips'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'tulips'", "b'dandelion'", "b'daisy'", "b'daisy'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'tulips'", "b'roses'", "b'dandelion'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
bs: 4, ds: 6, rep: 1
Image batch shape (4,), ["b'dandelion'", "b'tulips'", "b'daisy'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'tulips'", "b'roses'", "b'daisy'", "b'sunflowers'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'tulips'", "b'tulips'", "b'sunflowers'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'tulips'", "b'daisy'", "b'tulips'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'daisy'", "b'dandelion'", "b'daisy'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4,), ["b'tulips'", "b'sunflowers'", "b'dandelion'", "b'dandelion'"]) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
total time: 3.4361982345581055
[[[ 9.77924234]
   [14.86563858]]

 [[17.00530245]
  [19.15936353]]]
[[[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'>
bs: 4, ds: 6, rep: 1
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
Image batch shape (4, 192, 192, 3), ['0', '0', '0', '0']) <_io.TextIOWrapper name='/dev/null' mode='w' encoding='UTF-8'>
total time: 0.8428754806518555
[[[ 24.65395851]
```



```
[ 68.53605889]]

[[ 58.35654881]
 [118.46587362]]]
[[[2. 3. 1.]]

 [2. 6. 1.]]]

[[[4. 3. 1.]]

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

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

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

2a) Create the script (14%)

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

More specifically:

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

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

```
In [ ]: #Task 2a) --> i) Combine the previous cells to have the code to create a dataset and create a list of parameter

import itertools
import time
import tensorflow as tf
import pyspark
from pyspark import SparkContext

def read_tfrecord(record): #This function is created to read and parse a TFRecord file
    feature_description = {
        "image": tf.io.FixedLenFeature([], tf.string),
        "class": tf.io.FixedLenFeature([], tf.int64)
    }
    example = tf.io.parse_single_example(record, feature_description)
    image_decoded = tf.image.decode_jpeg(example['image'], channels=3)
    image_resized = tf.reshape(image_decoded, [*TARGET_SIZE, 3])
    class_label = example['class']
    return image_resized, class_label

def load_dataset(filenames): #This function is created to create a dataset from the TFRecord files
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False
    dataset = tf.data.TFRecordDataset(filenames)
    dataset = dataset.with_options(option_no_order)
    dataset_new = dataset.map(read_tfrecord)
    return dataset_new

def decode_jpeg_and_label(filepath): #This function is created to decode the JPEG images from the filepath and
    image_data = tf.io.read_file(filepath)
    decoded_image = tf.image.decode_jpeg(image_data)
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
    class_label = label.values[-2]
    return decoded_image, class_label

def resize_and_crop_image(image, label): #This function is created to resize and crop the images to the specifi
    image_width = tf.shape(image)[0]
    image_height = tf.shape(image)[1]
    target_width = TARGET_SIZE[1]
    target_height = TARGET_SIZE[0]
```

```

resize_crit = (image_width * target_height) / (image_height * target_width)
image = tf.cond(resize_crit < 1,
                lambda: tf.image.resize(image, [image_width*target_width/image_height, image_height*target_width]),
                lambda: tf.image.resize(image, [image_width*target_height/image_height, image_height*target_width]))
new_width = tf.shape(image)[0]
new_height = tf.shape(image)[1]
image = tf.image.crop_to_bounding_box(image, (new_width - target_width) // 2, (new_height - target_height) // 2, target_width, target_height)
return image, label

def recompress_image(image, label): #This function then recompresses the processed images to JPEG
uint8_image = tf.cast(image, tf.uint8)
jpeg_image = tf.image.encode_jpeg(uint8_image, optimize_size=True, chroma_downsampling=False)
return jpeg_image, label

def load_image_jpeg_dataset(): #We load the dataset of the images from the JPEG files which are stored in a Google Cloud Storage bucket
GCS_PATTERN = 'gs://flowers-public/*.jpg'
files_dataset = tf.data.Dataset.list_files(GCS_PATTERN)
decoded_images_dataset = files_dataset.map(decode_jpeg_and_label)
resized_dataset = decoded_images_dataset.map(resize_and_crop_image)
image_jpeg_dataset = resized_dataset.map(recompress_image)
return image_jpeg_dataset

def parameter_combinations(batch_sizes, batch_numbers, repetitions): #This function is used to generate all of the combinations of batch sizes, batch numbers and repetitions
return list(itertools.product(batch_sizes, batch_numbers, repetitions))

```

In []: #Task 2a) --> ii) Get a Spark context and create the dataset and run timing test for each combination in parallel

```

sc = SparkContext.getOrCreate() #We set up the spark context and then define the parameter combinations

#These are the batch sizes, batch numbers and repetitions that we will be testing
batch_sizes = [2, 4]
batch_numbers = [3, 6, 9, 12]
repetitions = [1, 2, 3]

# We create an RDD from the parameter combinations
parameters_rdd = sc.parallelize(parameter_combinations(batch_sizes, batch_numbers, repetitions))

#This function is created to time the processing of the TFRecord batches
def time_configs(params):
    batch_size, batch_number, repetition = params
    TRF_dataset = load_dataset(filename)
    results = []
    for repet in range(repetition):
        start_time = time.time()
        for _ in TRF_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        results.append(('TFRecord', params, images_per_second))
    return results

#This function is created to time the processing of the image data batches
def image_configs(params):
    batch_size, batch_number, repetition = params
    img_dataset = load_image_jpeg_dataset()
    image_processing_results = []
    for rep in range(repetition):
        start_time = time.time()
        for _ in img_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        image_processing_results.append(('Image', params, images_per_second))
    return image_processing_results

#These are the results from the time processing of the TFRecord and image data batches
resulting_rdd = parameters_rdd.flatMap(time_configs)
resulting_rdd_2 = parameters_rdd.flatMap(image_configs)

```

In []: #Task 2a) iii) --> Transform the resulting RDD to the structure (parameter_combination, images_per_second) and

```

# We merge the 2 resulting RDDs (TFRecord and image results) into a combined RDD using the .union function
results_rdd_combined = resulting_rdd.union(resulting_rdd_2)
final_results = results_rdd_combined.collect()

```

In []: #We print the results

```

print('Combined Results:')
for result in final_results:
    print(f'Dataset: {result[0]}, Parameters: (Batch Size, Batch Number, Repetitions): {result[1]}, Average Image

```

```

Combined Results:
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 8
.672589037158751
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 1
7.993788002722752
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 1
7.54538818553928
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1

```

7.0174340997521
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
6.36668498496057
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
7.577459122308287
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 3
5.06518022480552
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 6
2.61775185620161
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 3
1.016898202212708
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 3
3.80220214169865
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 5
8.11136511631129
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 6
5.3718044131342
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 5
2.836780669628894
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 2
6.943425476985976
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 5
2.70142005415514
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 6
2.243933486626176
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 8
5.23035205265268
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 5
3.89746859028771
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second: 124.76518488588562
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 70.49156628868383
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 87.7592618908337
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 71.84318629014818
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 87.86152095131028
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 71.25399383750289
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 1
7.126621922719373
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 4
3.010672416047335
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 4
0.26013289488314
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 4
4.41032989450576
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
2.35012318096772
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
5.68059083365648
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 6
1.20923252605381
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 6
0.09684433241912
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 9
5.43860500747101
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 8
5.0917593197616
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 8
1.96495799290952
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 7
3.74270893388648
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 1
37.92969186753922
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 1
78.86751041256994
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 2
11.33602994073996
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 1
50.4910997978767
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 1
55.1860537537269
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 1
84.60055137507564
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second: 214.45601835577742
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 234.73556619414836
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 198.80063118022406
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 224.66301316666053
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 239.5821040455824
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 193.4701559088533

Dataset: Image, 39854168743527	Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 10.2
Dataset: Image, 10985436568447	Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 11.6
Dataset: Image, 21079467266252	Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 11.1
Dataset: Image, 6299562818131	Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 10.8
Dataset: Image, 04695732433617	Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 12.8
Dataset: Image, 02564193367723	Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 11.0
Dataset: Image, 94775417772631	Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 15.4
Dataset: Image, 03802045956633	Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 11.8
Dataset: Image, 15869552374013	Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 13.2
Dataset: Image, 48791291803045	Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 11.2
Dataset: Image, 96541718784673	Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 13.9
Dataset: Image, 3376085680006	Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 12.5
Dataset: Image, 6227155044353	Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 21.3
Dataset: Image, 10168169671876	Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 21.3
Dataset: Image, 68539749879625	Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 20.4
Dataset: Image, 0781641243588	Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 22.7
Dataset: Image, 51567605804203	Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 21.8
Dataset: Image, 05933665155454	Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 20.8
Dataset: Image, 41216905657873	Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second: 21.
Dataset: Image, 367231355802254	Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 19.
Dataset: Image, 21013999903443	Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 14.
Dataset: Image, 971485769016976	Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 16.
Dataset: Image, 76585315748315	Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 17.
Dataset: Image, 49439215583462	Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 16.
Dataset: Image, 12642334744724	Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 17.6
Dataset: Image, 2172885974805	Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 14.9
Dataset: Image, 66155180715782	Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 15.4
Dataset: Image, 6793520968833	Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 14.5
Dataset: Image, 00325121704469	Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 13.6
Dataset: Image, 24722610910998	Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 11.8
Dataset: Image, 94989917476793	Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 16.1
Dataset: Image, 6617579423364	Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 17.1
Dataset: Image, 8909301748649	Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 23.9
Dataset: Image, 83141211436678	Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 22.4
Dataset: Image, 2113530106908	Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 23.1
Dataset: Image, 79811719823385	Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 25.4
Dataset: Image, 85937601669412	Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 24.6
Dataset: Image, 76991001014675	Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 23.4
Dataset: Image, 90310381164603	Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 18.5
Dataset: Image, 57758419046294	Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 13.2
Dataset: Image, 97971014366842	Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 16.3
Dataset: Image, 09736781367477	Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 20.1
Dataset: Image, 02988234016515	Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second: 28.
Dataset: Image, 78318570092852	Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 32.
Dataset: Image,	Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 33.

```
119408568453366
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 31.423293743860157
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 33.204511786621985
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 30.22334579437569
```

```
In [ ]: #Task 2a --> iv) Create an RDD with all results for each parameter as (parameter_value,images_per_second) and c
```

```
#This function extracts the results for each parameter
def param_results(index_param):
    return results_rdd_combined.map(lambda x: (x[1][index_param], x[2]))

#We create separate RDDs for each parameter
rdd_batch_size = param_results(0)
rdd_batch_number = param_results(1)
rdd_repetition = param_results(2)

#This function is created to calculate the average images per second for each parameter separately
def average_images_param(data):
    total_images = sum(images for images in data)
    count = len(data)
    return total_images / count

avg_batch_size_rdd = rdd_batch_size.groupByKey().mapValues(average_images_param).collect()
avg_batch_number_rdd = rdd_batch_number.groupByKey().mapValues(average_images_param).collect()
avg_repetition_rdd = rdd_repetition.groupByKey().mapValues(average_images_param).collect()

#We print the average images per second for each unique parameter
print("Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec):", avg_batch_size_rdd)
print("Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec):", avg_batch_number_rdd)
print("Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec):", avg_repetition_rdd)
```

```
Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec): [(4, 63.0808506251448), (2, 29.66050117001694)]
Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec): [(12, 85.01437109141926), (9, 72.21338865151826), (6, 46.229904480424985), (3, 23.234236315690193)]
Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec): [(1, 43.104113930752455), (2, 51.080036857344126), (3, 48.66364517169166)]
```

```
In [ ]: #Task 2a --> v) Create an RDD with the average reading speeds for each parameter value and collect the results
```

```
#This function calculates the average reading speed using a tuple (total, count) for the associativity
def seq_op(accumu, n):
    return (accumu[0] + n, accumu[1] + 1)

def comb_op(accumu1, accumu2):
    return (accumu1[0] + accumu2[0], accumu1[1] + accumu2[1])

#Now we can create a function that calculates the average images per second while also keeping associativity in
def average_assoc(rdd_params):
    return rdd_params.aggregateByKey((0,0), seq_op, comb_op).mapValues(lambda x: x[0] / x[1]).collect()

avg_batch_size = average_assoc(rdd_batch_size)
avg_batch_number = average_assoc(rdd_batch_number)
avg_repetition = average_assoc(rdd_repetition)

#Now we print the average reading speed for the associativity-respecting method
print("Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec):",
print("Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec):",
print("Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec):", avg_
```

```
Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec): [(4, 71.50003330036951), (2, 33.94330722930943)]
Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec): [(12, 95.15781559985273), (9, 68.57971692643322), (6, 48.87393143768083), (3, 25.32879120167082)]
Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec): [(1, 46.9615611641283), (2, 54.5184233486253), (3, 49.760002537199995)]
```

```
In [ ]: #Task 2a --> vi) Write the results to a pickle file in your bucket
```

```
import pickle
from google.cloud import storage

#We store the results from previous tasks to a dictionary
results_task2a = {'Final Results Task iii': final_results,
                  'AVG Batch Size Task iv': avg_batch_size_rdd,
                  'AVG Batch Number Task iv': avg_batch_number_rdd,
                  'AVG Repetition Task iv': avg_repetition_rdd,
                  'AVG Batch Size Task v': avg_batch_size,
                  'AVG Batch Number Task v': avg_batch_number,
                  'AVG Repetition Task v': avg_repetition}

#We write the results to a local pickle file
with open('/tmp/all_results.pkl', 'wb') as f:
    pickle.dump(results_task2a, f)

#This function is used to upload the pickle file to the bucket that we specify
def upload_to_bucket(blob_name, path_to_file, bucket_name):
```

```

try:
    client = storage.Client()
    bucket = client.get_bucket(bucket_name)
    blob = bucket.blob(blob_name)
    blob.upload_from_filename(path_to_file)
    print("File uploaded successfully")
except Exception as e:
    print(f"An error occurred: {e}")

```

```
upload_to_bucket('results/results_task2a.pkl', '/tmp/all_results.pkl', 'daring-healer-421511-storage')
```

File uploaded successfully

```

In [ ]: %%writefile spark_job.py
#Task 2a --> vii) Write your code it into a file using the cell magic %%writefile spark_job.py
import itertools
import time
import tensorflow as tf
import pyspark
from pyspark import SparkContext
#We declare the variables
GCS_PATTERN = 'gs://flowers-public/*/*.jpg' #glob pattern for input files.
PROJECT = 'daring-healer-421511'
BUCKET = 'gs://{}-storage'.format(PROJECT) # This is the bucket storage.
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # We prefix for output file names.
filenames = tf.io.gfile.glob(GCS_OUTPUT + "/*.tfrec")
TARGET_SIZE = [192,192]
PARTITIONS = 16 # This is the number of partitions that will be used later.

#Task 2a) --> i) Combine the previous cells to have the code to create a dataset and create a list of parameter
import itertools
import time
import tensorflow as tf
import pyspark
from pyspark import SparkContext

def read_tfrecord(record): #This function is created to read and parse a TFRecord file
    feature_description = {
        "image": tf.io.FixedLenFeature([], tf.string),
        "class": tf.io.FixedLenFeature([], tf.int64)
    }
    example = tf.io.parse_single_example(record, feature_description)
    image_decoded = tf.image.decode_jpeg(example['image'], channels=3)
    image_resized = tf.image.resize(image_decoded, [*TARGET_SIZE, 3])
    class_label = example['class']
    return image_resized, class_label

def load_dataset(filenames): #This function is created to create a dataset from the TFRecord files
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False
    dataset = tf.data.TFRecordDataset(filenames)
    dataset = dataset.with_options(option_no_order)
    dataset_new = dataset.map(read_tfrecord)
    return dataset_new

def decode_jpeg_and_label(filepath): #This function is created to decode the JPEG images from the filepath and
    image_data = tf.io.read_file(filepath)
    decoded_image = tf.image.decode_jpeg(image_data)
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
    class_label = label.values[-2]
    return decoded_image, class_label

def resize_and_crop_image(image, label): #This function is created to resize and crop the images to the speci
    image_width = tf.shape(image)[0]
    image_height = tf.shape(image)[1]
    target_width = TARGET_SIZE[1]
    target_height = TARGET_SIZE[0]
    resize_crit = (image_width * target_height) / (image_height * target_width)
    image = tf.cond(resize_crit < 1,
        lambda: tf.image.resize(image, [image_width*target_width/image_height, image_height*target_w
        lambda: tf.image.resize(image, [image_width*target_height/image_height, image_height*target
    new_width = tf.shape(image)[0]
    new_height = tf.shape(image)[1]
    image = tf.image.crop_to_bounding_box(image, (new_width - target_width) // 2, (new_height - target_height)
    return image, label

def recompress_image(image, label): #This function then recompresses the processed images to JPEG
    uint8_image = tf.cast(image, tf.uint8)
    jpeg_image = tf.image.encode_jpeg(uint8_image, optimize_size=True, chroma_downsampling=False)
    return jpeg_image, label

def load_image_jpeg_dataset(): #We load the dataset of the images from the JPEG files which are stored in a Goo
    GCS_PATTERN = 'gs://flowers-public/*/*.jpg'
    files_dataset = tf.data.Dataset.list_files(GCS_PATTERN)
    decoded_images_dataset = files_dataset.map(decode_jpeg_and_label)
    resized_dataset = decoded_images_dataset.map(resize_and_crop_image)
    image_jpeg_dataset = resized_dataset.map(recompress_image)
    return image_jpeg_dataset

```



```

def parameter_combinations(batch_sizes, batch_numbers, repetitions): #This function is used to generate all of
    return list(itertools.product(batch_sizes, batch_numbers, repetitions))

#Task 2a) --> ii) Get a Spark context and create the dataset and run timing test for each combination in parallel

sc = SparkContext.getOrCreate() #We set up the spark context and then define the parameter combinations

#These are the batch sizes, batch numbers and repetitions that we will be testing
batch_sizes = [2, 4]
batch_numbers = [3, 6, 9, 12]
repetitions = [1, 2, 3]

# We create an RDD from the parameter combinations
parameters_rdd = sc.parallelize(parameter_combinations(batch_sizes, batch_numbers, repetitions))

#This function is created to time the processing of the TFRecord batches
def time_configs(params):
    batch_size, batch_number, repetition = params
    TRF_dataset = load_dataset(filenamees)
    results = []
    for repet in range(repetition):
        start_time = time.time()
        for _ in TRF_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        results.append(('TFRecord', params, images_per_second))
    return results

#This function is created to time the processing of the image data batches
def image_configs(params):
    batch_size, batch_number, repetition = params
    img_dataset = load_image_jpeg_dataset()
    image_processing_results = []
    for rep in range(repetition):
        start_time = time.time()
        for _ in img_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        image_processing_results.append(('Image', params, images_per_second))
    return image_processing_results

#These are the results from the time processing of the TFRecord and image data batches
resulting_rdd = parameters_rdd.flatMap(time_configs)
resulting_rdd_2 = parameters_rdd.flatMap(image_configs)

#Task 2a) iii) --> Transform the resulting RDD to the structure ( parameter_combination, images_per_second ) and collect

# We merge the 2 resulting RDDs (TFRecord and image results) into a combined RDD using the .union function
results_rdd_combined = resulting_rdd.union(resulting_rdd_2)
final_results = results_rdd_combined.collect()

#We print the results
print('Combined Results:')
for result in final_results:
    print(f'Dataset: {result[0]}, Parameters: (Batch Size, Batch Number, Repetitions): {result[1]}, Average Image

#Task 2a --> iv) Create an RDD with all results for each parameter as (parameter_value, images_per_second) and collect

#This function extracts the results for each parameter
def param_results(index_param):
    return results_rdd_combined.map(lambda x: (x[1][index_param], x[2]))

#We create separate RDDs for each parameter
rdd_batch_size = param_results(0)
rdd_batch_number = param_results(1)
rdd_repetition = param_results(2)

#This function is created to calculate the average images per second for each parameter separately
def average_images_param(data):
    total_images = sum(images for images in data)
    count = len(data)
    return total_images / count

avg_batch_size_rdd = rdd_batch_size.groupByKey().mapValues(average_images_param).collect()
avg_batch_number_rdd = rdd_batch_number.groupByKey().mapValues(average_images_param).collect()
avg_repetition_rdd = rdd_repetition.groupByKey().mapValues(average_images_param).collect()

#We print the average images per second for each unique parameter
print("Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec):", avg_batch_size_rdd)
print("Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec):", avg_batch_number_rdd)
print("Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec):", avg_repetition_rdd)

#Task 2a --> v) Create an RDD with the average reading speeds for each parameter value and collect the results

#This function calculates the average reading speed using a tuple (total, count) for the associativity
def seq_op(accumu, n):
    return (accumu[0] + n, accumu[1] + 1)

```

```

def comb_op( accumu1, accumu2):
    return (accumu1[0] + accumu2[0], accumu1[1] + accumu2[1])

#Now we can create a function that calculates the average images per second while also keeping associativity in
def average_assoc(rdd_params):
    return rdd_params.aggregateByKey((0,0), seq_op, comb_op).mapValues(lambda x: x[0] / x[1]).collect()

avg_batch_size = average_assoc(rdd_batch_size)
avg_batch_number = average_assoc(rdd_batch_number)
avg_repetition = average_assoc(rdd_repetition)

#Now we print the average reading speed for the associativity-respecting method
print("Average Images per Second By Each Batch Size (Associativity-method), (Batch Size: AVG_Images_per_Sec):",
print("Average Images per Second By Each Batch Number (Associativity-method), (Batch Number: AVG_Images_per_Sec)",
print("Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec):", avg_

#Task 2a --> vi) Write the results to a pickle file in your bucket
import pickle
from google.cloud import storage

#We store the results from previous tasks to a dictionary
results_task2a = {'Final Results Task iii': final_results,
                  'AVG Batch Size Task iv': avg_batch_size_rdd,
                  'AVG Batch Number Task iv': avg_batch_number_rdd,
                  'AVG Repetition Task iv': avg_repetition_rdd,
                  'AVG Batch Size Task v': avg_batch_size,
                  'AVG Batch Number Task v': avg_batch_number,
                  'AVG Repetition Task v': avg_repetition}

#We write the results to a local pickle file
with open('/tmp/all_results.pkl', 'wb') as f:
    pickle.dump(results_task2a, f)

#This function is used to upload the pickle file to the bucket that we specify
def upload_to_bucket(blob_name, path_to_file, bucket_name):
    try:
        client = storage.Client()
        bucket = client.get_bucket(bucket_name)
        blob = bucket.blob(blob_name)
        blob.upload_from_filename(path_to_file)
        print("File uploaded successfully")
    except Exception as e:
        print(f"An error occurred: {e}")

upload_to_bucket('results/results_task2a.pkl', '/tmp/all_results.pkl', 'daring-healer-421511-storage')

```

Writing spark_job.py

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

i) First, test locally with `%run`.

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

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

In []: `#Task 2b) i) --> We test locally on spark with the %run command`
`%run ./spark_job.py`

```

Combined Results:
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 1
6.01438670375341
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 1
6.577173719008155
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 2
2.828628966418112
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
6.226321565310016
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
6.18267529800162
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 2
4.96795297269142
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 3
0.460940705692675
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 3
3.01654840442734
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 2
8.606320465552194
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 3
0.931654528146392
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 2
9.76851229737019
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 2
9.99403952790582
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 4

```

9.307337373404884
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 4
2.828155207624235
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 4
7.21232939675631
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 6
3.09369351149472
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 5
1.9894172169925
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 2
5.908899885001812
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second:
87.57145156281128
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second:
77.0613944700729
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second:
66.78145677790353
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
56.90651265967334
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
67.97688348417053
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
69.47974448083434
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 3
1.52193937814943
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 3
1.43246799557599
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 4
4.46396134156091
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 4
4.04590144997545
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 4
4.619029086966584
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
0.213648462233877
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 7
2.62733067585361
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7
0.3238844674953
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7
5.87460673276522
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 6
6.30935006317165
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 7
9.35835517732062
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 9
5.28710239573049
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 1
07.36883781785085
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 1
19.71241419662353
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 1
27.62555605611993
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 1
27.85151792775385
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 9
5.42045112002715
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 1
25.39348214411764
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second:
146.6222648509173
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second:
170.81595802538058
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second:
167.11498033981368
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second:
140.75947781146175
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second:
130.56851217602954
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second:
156.44319561985827
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 7.54
0068767729742
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 8.61
9339088047957
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 7.73
3997311552028
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 6.97
5080114480141
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 8.03
8260419611245
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 6.84
2594483306547
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 13.0
03755799890817
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 21.5
66938233505773
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 18.4
6983475305065

Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 19.2
4964708861011
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 18.2
3090589748503
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 16.5
52477495550413
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 20.1
96836668196863
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 20.7
55415764142402
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 20.4
63712962827437
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 20.7
8452852503642
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 18.9
06466453421295
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 13.8
6166440282501
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second: 13.
992706960803215
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 17.
81707072947241
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 16.
35004650986014
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 15.
499468252920991
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 15.
482644963168015
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 18.
913996368766885
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 10.6
60236987121296
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 13.5
89437467549631
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 10.8
26890508292205
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 12.2
23032315298738
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 10.8
23679793004171
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 13.4
58764663302576
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 25.7
7477232731199
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 22.9
00370497935942
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 22.9
63110580660704
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 22.5
14957888163416
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 21.2
68000355793163
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 23.0
65777244955758
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 21.0
5119878741035
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 18.4
02186301685063
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 16.5
50284022248274
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 18.6
81516371309428
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 18.3
06123057637123
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 17.2
21007093125543
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second: 28.
068570281359662
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 33.
269154544638816
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 35.
24326580170728
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 35.
524408467914476
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 34.
96797752340617
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 34.
515883438473146
Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec): [(4, 63.135423818380644), (2, 31.
.324376398417105)]
Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec): [(12, 83.40902202446058), (9
, 53.827162916524486), (6, 38.74526696990669), (3, 22.04843169570539)]
Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec): [(1, 60.39748965139749), (2, 60.
39259585420249), (3, 56.14607246572556)]
Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec): [(4, 63.
02181885053674), (2, 29.031656682472033)]
Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec): [(12
, 74.2086218855064), (9, 50.873770902546404), (6, 42.07954322345875), (3, 21.54856210188757)]
Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec): [(1, 50.39164

548695515), (2, 56.650266009866115), (3, 55.550553089101506)]
File uploaded successfully

In []: #New filename argument

```
%%writefile spark_job_2b.py
#Task 2a --> vii) Write your code it into a file using the cell magic %%writefile spark_job.py
import os
os.environ['PROTOCOL_BUFFERS_PYTHON_IMPLEMENTATION'] = 'python'
import itertools
import time
import tensorflow as tf
import google.protobuf
import pyspark
from pyspark import SparkContext
#We declare the variables
GCS_PATTERN = 'gs://flowers-public/*/*.jpg' #glob pattern for input files.
PROJECT = 'daring-healer-421511'
BUCKET = 'gs://{}-storage'.format(PROJECT) # This is the bucket storage.
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # We prefix for output file names.
filenames = tf.io.gfile.glob(GCS_OUTPUT + "/*.tfrec")
TARGET_SIZE = [192,192]
PARTITIONS = 16 # This is the number of partitions that will be used later.

#Task 2a) --> i) Combine the previous cells to have the code to create a dataset and create a list of parameter
import itertools
import time
import tensorflow as tf
import pyspark
from pyspark import SparkContext

def read_tfrecord(record): #This function is created to read and parse a TFRecord file
    feature_description = {
        "image": tf.io.FixedLenFeature([], tf.string),
        "class": tf.io.FixedLenFeature([], tf.int64)
    }
    example = tf.io.parse_single_example(record, feature_description)
    image_decoded = tf.image.decode_jpeg(example['image'], channels=3)
    image_resized = tf.reshape(image_decoded, [*TARGET_SIZE, 3])
    class_label = example['class']
    return image_resized, class_label

def load_dataset(filenames): #This function is created to create a dataset from the TFRecord files
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False
    dataset = tf.data.TFRecordDataset(filenames)
    dataset = dataset.with_options(option_no_order)
    dataset_new = dataset.map(read_tfrecord)
    return dataset_new

def decode_jpeg_and_label(filepath): #This function is created to decode the JPEG images from the filepath and
    image_data = tf.io.read_file(filepath)
    decoded_image = tf.image.decode_jpeg(image_data)
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
    class_label = label.values[-2]
    return decoded_image, class_label

def resize_and_crop_image(image, label): #This function is created to resize and crop the images to the specifi
    image_width = tf.shape(image)[0]
    image_height = tf.shape(image)[1]
    target_width = TARGET_SIZE[1]
    target_height = TARGET_SIZE[0]
    resize_crit = (image_width * target_height) / (image_height * target_width)
    image = tf.cond(resize_crit < 1,
        lambda: tf.image.resize(image, [image_width*target_width/image_height, image_height*target_w
        lambda: tf.image.resize(image, [image_width*target_height/image_height, image_height*target
    new_width = tf.shape(image)[0]
    new_height = tf.shape(image)[1]
    image = tf.image.crop_to_bounding_box(image, (new_width - target_width) // 2, (new_height - target_height)
    return image, label

def recompress_image(image, label): #This function then recompresses the processed images to JPEG
    uint8_image = tf.cast(image, tf.uint8)
    jpeg_image = tf.image.encode_jpeg(uint8_image, optimize_size=True, chroma_downsampling=False)
    return jpeg_image, label

def load_image_jpeg_dataset(): #We load the dataset of the images from the JPEG files which are stored in a Goo
    GCS_PATTERN = 'gs://flowers-public/*/*.jpg'
    files_dataset = tf.data.Dataset.list_files(GCS_PATTERN)
    decoded_images_dataset = files_dataset.map(decode_jpeg_and_label)
    resized_dataset = decoded_images_dataset.map(resize_and_crop_image)
    image_jpeg_dataset = resized_dataset.map(recompress_image)
    return image_jpeg_dataset

def parameter_combinations(batch_sizes, batch_numbers, repetitions): #This function is used to generate all of
    return list(itertools.product(batch_sizes, batch_numbers, repetitions))

#Task 2a) --> ii) Get a Spark context and create the dataset and run timing test for each combination in parall
```

```

sc = SparkContext.getOrCreate() #We set up the spark context and then define the parameter combinations

#These are the batch sizes, batch numbers and repetitions that we will be testing
batch_sizes = [2, 4]
batch_numbers = [3, 6, 9, 12]
repetitions = [1, 2, 3]

# We create an RDD from the parameter combinations
parameters_rdd = sc.parallelize(parameter_combinations(batch_sizes, batch_numbers, repetitions))

#This function is created to time the processing of the TFRecord batches
def time_configs(params):
    batch_size, batch_number, repetition = params
    TRF_dataset = load_dataset(filenamees)
    results = []
    for repet in range(repetition):
        start_time = time.time()
        for _ in TRF_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        results.append(('TFRecord', params, images_per_second))
    return results

#This function is created to time the processing of the image data batches
def image_configs(params):
    batch_size, batch_number, repetition = params
    img_dataset = load_image_jpeg_dataset()
    image_processing_results = []
    for rep in range(repetition):
        start_time = time.time()
        for _ in img_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        image_processing_results.append(('Image', params, images_per_second))
    return image_processing_results

#These are the results from the time processing of the TFRecord and image data batches
resulting_rdd = parameters_rdd.flatMap(time_configs)
resulting_rdd_2 = parameters_rdd.flatMap(image_configs)

#Task 2a) iii) --> Transform the resulting RDD to the structure ( parameter_combination, images_per_second ) and c

# We merge the 2 resulting RDDs (TFRecord and image results) into a combined RDD using the .union function
results_rdd_combined = resulting_rdd.union(resulting_rdd_2)
final_results = results_rdd_combined.collect()

#We print the results
print('Combined Results:')
for result in final_results:
    print(f'Dataset: {result[0]}, Parameters: (Batch Size, Batch Number, Repetitions): {result[1]}, Average Image

#Task 2a --> iv) Create an RDD with all results for each parameter as (parameter_value, images_per_second) and c

#This function extracts the results for each parameter
def param_results(index_param):
    return results_rdd_combined.map(lambda x: (x[1][index_param], x[2]))

#We create separate RDDs for each parameter
rdd_batch_size = param_results(0)
rdd_batch_number = param_results(1)
rdd_repetition = param_results(2)

#This function is created to calculate the average images per second for each parameter separately
def average_images_param(data):
    total_images = sum(images for images in data)
    count = len(data)
    return total_images / count

avg_batch_size_rdd = rdd_batch_size.groupByKey().mapValues(average_images_param).collect()
avg_batch_number_rdd = rdd_batch_number.groupByKey().mapValues(average_images_param).collect()
avg_repetition_rdd = rdd_repetition.groupByKey().mapValues(average_images_param).collect()

#We print the average images per second for each unique parameter
print("Average Images per Second By Each Batch Size (Batch Size: AVG_Images_per_Sec):", avg_batch_size_rdd)
print("Average Images per Second By Each Batch Number (Batch Number: AVG_Images_per_Sec):", avg_batch_number_rdd)
print("Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec):", avg_repetition_rdd)

#Task 2a --> v) Create an RDD with the average reading speeds for each parameter value and collect the results

#This function calculates the average reading speed using a tuple (total, count) for the associativity
def seq_op(accumu, n):
    return (accumu[0] + n, accumu[1] + 1)

def comb_op(accumu1, accumu2):
    return (accumu1[0] + accumu2[0], accumu1[1] + accumu2[1])

#Now we can create a function that calculates the average images per second while also keeping associativity in

```



```

def average_assoc(rdd_params):
    return rdd_params.aggregateByKey((0,0), seq_op, comb_op).mapValues(lambda x: x[0] / x[1]).collect()

avg_batch_size = average_assoc(rdd_batch_size)
avg_batch_number = average_assoc(rdd_batch_number)
avg_repetition = average_assoc(rdd_repetition)

#Now we print the average reading speed for the associativity-respecting method
print("Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec):",
print("Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec
print("Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec):", avg_

#Task 2a --> vi) Write the results to a pickle file in your bucket
import pickle
from google.cloud import storage

#We store the results from previous tasks to a dictionary
results_task2b = {'Final Results Task iii': final_results,
                  'AVG Batch Size Task iv': avg_batch_size_rdd,
                  'AVG Batch Number Task iv': avg_batch_number_rdd,
                  'AVG Repetition Task iv': avg_repetition_rdd,
                  'AVG Batch Size Task v': avg_batch_size,
                  'AVG Batch Number Task v': avg_batch_number,
                  'AVG Repetition Task v': avg_repetition}

#We write the results to a local pickle file
with open('/tmp/all_results.pkl', 'wb') as f:
    pickle.dump(results_task2b, f)

#This function is used to upload the pickle file to the bucket that we specify
def upload_to_bucket(blob_name, path_to_file, bucket_name):
    try:
        client = storage.Client()
        bucket = client.get_bucket(bucket_name)
        blob = bucket.blob(blob_name)
        blob.upload_from_filename(path_to_file)
        print("File uploaded successfully")
    except Exception as e:
        print(f"An error occurred: {e}")

upload_to_bucket('results/results_task2b.pkl', '/tmp/all_results.pkl', 'daring-healer-421511-storage')

```

Writing spark_job_2b.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.

```

In [ ]: #We will set up the maximal cluster which is a cluster with 1 machine and 7 workers using maximal SSD size (100
CLUSTER = '{}-maximalcluster7task2b'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0'
!gcloud dataproc clusters create $CLUSTER \
    --region $REGION \
    --bucket $PROJECT-storage \
    --image-version 1.5-ubuntu18 \
    --master-machine-type n1-standard-1 \
    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
    --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-disk-size 100 \
    --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
    --metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib google-cloud-storage==1.42.0 google-api-core==1.31.0
#Reference link for code: https://cloud.google.com/dataproc/docs/guides/create-cluster

```

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/9dbcbb79-60f8-3463-a9a4-10ba23354652].

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

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting 'https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453'.

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

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

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See <https://cloud.google.com/storage/docs/uniform-bucket-level-access>.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-maximalcluster7task2b] Cluster placed in zone [us-west1-c].

```
In [ ]: #We get information on the single machine cluster we created
!gcloud dataproc clusters describe $CLUSTER --region $REGION
```

```
clusterName: daring-healer-421511-maximalcluster7task2b
clusterUuid: 8abd7648-fc9d-4875-9486-ada5038b61b3
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
  metadata:
    PIP_PACKAGES: tensorflow==2.4.0 matplotlib google-cloud-storage==1.42.0 google-api-core==1.31.0
    protobuf==3.17.3
  networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
    - https://www.googleapis.com/auth/monitoring.write
  zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c
  initializationActions:
    - executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubuntu18-2030909-165100-rc01
    instanceNames:
      - daring-healer-421511-maximalcluster7task2b-m
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machineTypes/n1-standard-1
    minCpuPlatform: AUTOMATIC
    numInstances: 1
    preemptibility: NON_PREEMPTIBLE
  softwareConfig:
    imageVersion: 1.5.90-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
      distcp:mapreduce.map.java.opts: -Xmx576m
      distcp:mapreduce.map.memory.mb: '768'
      distcp:mapreduce.reduce.java.opts: -Xmx576m
      distcp:mapreduce.reduce.memory.mb: '768'
      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: '60'
      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: daring-healer-421511-maximalcluster7task2b-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: '30'
```

```

hdfs:dfs.namenode.servicerpc-address: daring-healer-421511-maximalcluster7task2b-m:8051
hive:hive.fetch.task.conversion: none
mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '1000'
mapred:mapreduce.job.maps: '60'
mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
mapred:mapreduce.job.reduces: '7'
mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLevelDbS
tateStoreService
mapred:mapreduce.map.cpu.vcores: '1'
mapred:mapreduce.map.java.opts: -Xmx819m
mapred:mapreduce.map.memory.mb: '1024'
mapred:mapreduce.reduce.cpu.vcores: '1'
mapred:mapreduce.reduce.java.opts: -Xmx1638m
mapred:mapreduce.reduce.memory.mb: '2048'
mapred:mapreduce.task.io.sort.mb: '256'
mapred:yarn.app.mapreduce.am.command-opts: -Xmx819m
mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
mapred:yarn.app.mapreduce.am.resource.mb: '1024'
spark-env:SPARK_DAEMON_MEMORY: 1000m
spark:spark.driver.maxResultSize: 480m
spark:spark.driver.memory: 960m
spark:spark.executor.cores: '1'
spark:spark.executor.instances: '2'
spark:spark.executor.memory: 2688m
spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
spark:spark.extraListeners: com.google.cloud.spark.performance.DataprocMetricsListener
spark:spark.scheduler.mode: FAIR
spark:spark.sql.cbo.enabled: 'true'
spark:spark.ui.port: '0'
spark:spark.yarn.am.memory: 640m
yarn-env:YARN_NODEMANAGER_HEAPSIZE: '1000'
yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '1000'
yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '1000'
yarn:yarn.nodemanager.address: 0.0.0.0:8026
yarn:yarn.nodemanager.resource.cpu-vcores: '1'
yarn:yarn.nodemanager.resource.memory-mb: '3072'
yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
yarn:yarn.scheduler.maximum-allocation-mb: '3072'
yarn:yarn.scheduler.minimum-allocation-mb: '256'
tempBucket: dataproc-temp-us-west1-869525608453-ejfnpe5x
workerConfig:
  diskConfig:
    bootDiskSizeGb: 100
    bootDiskType: pd-standard
    imageUri: https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-5-ubu18-20
230909-165100-rc01
  instanceNames:
    - daring-healer-421511-maximalcluster7task2b-w-0
    - daring-healer-421511-maximalcluster7task2b-w-1
    - daring-healer-421511-maximalcluster7task2b-w-2
    - daring-healer-421511-maximalcluster7task2b-w-3
    - daring-healer-421511-maximalcluster7task2b-w-4
    - daring-healer-421511-maximalcluster7task2b-w-5
    - daring-healer-421511-maximalcluster7task2b-w-6
  machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machin
eTypes/n1-standard-1
  minCpuPlatform: AUTOMATIC
  numInstances: 7
  preemptibility: NON_PREEMPTIBLE
labels:
  goog-dataproc-autozone: enabled
  goog-dataproc-cluster-name: daring-healer-421511-maximalcluster7task2b
  goog-dataproc-cluster-uid: 8abd7648-fc9d-4875-9486-ada5038b61b3
  goog-dataproc-location: us-west1
projectId: daring-healer-421511
status:
  state: RUNNING
  stateStartTime: '2024-05-03T11:39:06.978113Z'
statusHistory:
  - state: CREATING
    stateStartTime: '2024-05-03T11:34:20.326499Z'

```

```

In [ ]: #We submit the job to the cluster
!gcloud dataproc jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_job_2b.py
%time

```

```

Job [8d81a4061c7c48d394aff80aaa97b886] submitted.
Waiting for job output...
2024-05-03 11:39:24.770395: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dyna
mic library 'libcudart.so.11.0'; dLError: libcudart.so.11.0: cannot open shared object file: No such file or di
rectory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 11:39:24.770563: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dlerror if
you do not have a GPU set up on your machine.
24/05/03 11:39:30 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 11:39:30 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 11:39:30 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 11:39:30 INFO org.spark_project.jetty.util.log: Logging initialized @10831ms to org.spark_project.jett
y.util.log.Slf4jLog
24/05/03 11:39:31 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknow

```

```
n; jvm 1.8.0_382-b05
24/05/03 11:39:31 INFO org.spark_project.jetty.server.Server: Started @11135ms
24/05/03 11:39:31 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@7d7e2198{HTTP/
1.1, (http/1.1)}{0.0.0.0:40305}
24/05/03 11:39:33 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-42
1511-maximalcluster7task2b-m/10.138.0.47:8032
24/05/03 11:39:33 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at dari
ng-healer-421511-maximalcluster7task2b-m/10.138.0.47:10200
24/05/03 11:39:34 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 11:39:34 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 11:39:34 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory
-mb, units = Mi, type = COUNTABLE
24/05/03 11:39:34 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores
, units = , type = COUNTABLE
24/05/03 11:39:39 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application
_1714736218059_0001
Combined Results:
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 1
1.369661144139451
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 7
.8505496452628565
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 1
7.323960183387715
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
7.188889913426568
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
7.184699596840275
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
6.63440409603418
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 3
8.29543201356769
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 2
3.905252238936406
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 3
9.82397319941481
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 3
4.056861525802994
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 3
4.43764274424747
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 9
.214617020065164
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 6
1.179955316878214
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 4
9.82233091846418
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 5
4.25165562913907
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 5
5.24701143907551
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 5
7.78529632603326
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 4
8.66346443903005
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second:
71.11857677071757
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second:
61.74181285348116
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second:
78.01106657615013
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
67.78056944318193
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
75.69579498285508
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
79.48367779607017
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 2
0.401973409798707
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 3
7.67209212856444
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 3
6.63630216207521
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
5.79483043705529
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
7.14337330192522
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
6.850271114814156
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 7
2.94346595479189
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7
1.89911882536714
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7
3.626529843514
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 5
8.701036302364
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 6
3.07010665052686
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 6
6.692480612708
```

Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 109.9235488877694
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 103.66651493774651
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 117.5995065328324
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 97.45054638672254
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 116.02910814428664
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 93.80164326102704
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second: 151.0340618220119
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 156.2462637395898
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 146.74174458938512
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 109.30903985262331
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 150.42929225380263
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 145.95218649254278
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 6.9876417335559085
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 6.543571635657332
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 7.5447910267745115
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 7.30519625849185
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 7.38000920234874
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 6.076178514283562
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 7.9265932986021435
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 8.55937174601895
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 8.080641292603397
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 7.711285916325252
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 7.539569535084589
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 7.801167114344657
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 8.123648873594158
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 7.4021766657734736
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 7.598119105662378
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 7.7753561567074065
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 6.667990181916812
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 7.364068832061518
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second: 7.90195553258504
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 7.961398683209109
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 8.270428402544212
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 8.112498439567775
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 8.03660792161969
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 7.9943956880280105
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 8.223015651809085
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 7.263332569936329
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 7.211954745426906
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 8.30895283878551
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 7.521007175281098
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 7.788392904670593
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 7.7819173982622605
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7.715355177154957
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 8.52

```

4045132435072
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 8.08
0181415244072
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 7.56
6064553118483
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 7.98
7819806843288
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 7.98
1031782166987
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 8.11
7958423155816
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 7.97
0156202345711
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 7.94
1531540401554
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 8.10
40069806282
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 7.85
8210241537205
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second: 7.9
36212787606353
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 7.9
046289761598745
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 8.0
27201001636847
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 7.9
51445479009535
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 8.1
80730524417962
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 8.5
06152393386104
Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec): [(4, 47.23018529843363), (2, 27.
13569634339385)]
Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec): [(12, 56.62886134251874), (9
, 44.842599084098005), (6, 31.829199036760468), (3, 16.58136666064826)]
Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec): [(1, 39.0461379275475), (2, 36.4
6443658775471), (3, 38.2961960396541)]
Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec): [(4, 48.
56638437485924), (2, 27.83507904407864)]
Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec): [(12
, 57.78135044340478), (9, 47.186230283853604), (6, 33.118208895422974), (3, 18.099416014608707)]
Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec): [(1, 39.06135
827504427), (2, 38.7484547386544), (3, 37.948242768235396)]
File uploaded successfully
24/05/03 11:52:03 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@7d7e2198{HTTP/1.1, (http
/1.1)}{0.0.0.0:0}
Job [8d81a4061c7c48d394aff80aaa97b886] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/8abd7648-fc9d-4875-9486
-ada5038b61b3/jobs/8d81a4061c7c48d394aff80aaa97b886/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/8abd7648-fc9d-4875-94
86-ada5038b61b3/jobs/8d81a4061c7c48d394aff80aaa97b886/driveroutput
jobUuid: 68f7a417-e99e-3728-8216-6f871b03c190
placement:
  clusterName: daring-healer-421511-maximalcluster7task2b
  clusterUuid: 8abd7648-fc9d-4875-9486-ada5038b61b3
pysparkJob:
  mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/8abd7648-fc9d-4875-9486-a
da5038b61b3/jobs/8d81a4061c7c48d394aff80aaa97b886/staging/spark_job_2b.py
reference:
  jobId: 8d81a4061c7c48d394aff80aaa97b886
  projectId: daring-healer-421511
status:
  state: DONE
  stateStartTime: '2024-05-03T11:52:07.702944Z'
statusHistory:
- state: PENDING
  stateStartTime: '2024-05-03T11:39:16.693589Z'
- state: SETUP_DONE
  stateStartTime: '2024-05-03T11:39:16.724107Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2024-05-03T11:39:17.171750Z'
yarnApplications:
- name: spark_job_2b.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://daring-healer-421511-maximalcluster7task2b-m:8088/proxy/application_1714736218059_0001/
CPU times: user 4 µs, sys: 0 ns, total: 4 µs
Wall time: 8.58 µs

```

```

In [ ]: CLUSTER = '{}-maximalcluster7task2b'.format(PROJECT)
        REGION = 'us-west1'
        !gcloud dataproc clusters delete $CLUSTER --region $REGION
#Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c

```

2c) Improve efficiency (6%)

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

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

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

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

```
In [ ]: #New filename argument

%%writefile spark_job_2c.py
#Task 2a --> vii) Write your code it into a file using the cell magic %%writefile spark_job.py
import os
os.environ['PROTOCOL_BUFFERS_PYTHON_IMPLEMENTATION'] = 'python'
import itertools
import time
import tensorflow as tf
import google.protobuf
import pyspark
from pyspark import SparkContext
#We declare the variables
GCS_PATTERN = 'gs://flowers-public/*/*.jpg' #glob pattern for input files.
PROJECT = 'daring-healer-421511'
BUCKET = 'gs://{}-storage'.format(PROJECT) # This is the bucket storage.
GCS_OUTPUT = BUCKET + '/tfrecords-jpeg-192x192-2/flowers' # We prefix for output file names.
filenames = tf.io.gfile.glob(GCS_OUTPUT + "/*.tfrec")
TARGET_SIZE = [192,192]
PARTITIONS = 16 # This is the number of partitions that will be used later.

#Task 2a) --> i) Combine the previous cells to have the code to create a dataset and create a list of parameter
import itertools
import time
import tensorflow as tf
import pyspark
from pyspark import SparkContext

def read_tfrecord(record): #This function is created to read and parse a TFRecord file
    feature_description = {
        "image": tf.io.FixedLenFeature([], tf.string),
        "class": tf.io.FixedLenFeature([], tf.int64)
    }
    example = tf.io.parse_single_example(record, feature_description)
    image_decoded = tf.image.decode_jpeg(example['image'], channels=3)
    image_resized = tf.reshape(image_decoded, [*TARGET_SIZE, 3])
    class_label = example['class']
    return image_resized, class_label

def load_dataset(filenames): #This function is created to create a dataset from the TFRecord files
    option_no_order = tf.data.Options()
    option_no_order.experimental_deterministic = False
    dataset = tf.data.TFRecordDataset(filenames)
    dataset = dataset.with_options(option_no_order)
    dataset_new = dataset.map(read_tfrecord)
    return dataset_new

def decode_jpeg_and_label(filepath): #This function is created to decode the JPEG images from the filepath and
    image_data = tf.io.read_file(filepath)
    decoded_image = tf.image.decode_jpeg(image_data)
    label = tf.strings.split(tf.expand_dims(filepath, axis=-1), sep='/')
    class_label = label.values[-2]
    return decoded_image, class_label

def resize_and_crop_image(image, label): #This function is created to resize and crop the images to the specifi
    image_width = tf.shape(image)[0]
    image_height = tf.shape(image)[1]
    target_width = TARGET_SIZE[1]
    target_height = TARGET_SIZE[0]
    resize_crit = (image_width * target_height) / (image_height * target_width)
    image = tf.cond(resize_crit < 1,
        lambda: tf.image.resize(image, [image_width*target_width/image_height, image_height*target_w
        lambda: tf.image.resize(image, [image_width*target_height/image_height, image_height*target
    new_width = tf.shape(image)[0]
    new_height = tf.shape(image)[1]
    image = tf.image.crop_to_bounding_box(image, (new_width - target_width) // 2, (new_height - target_height)
    return image, label

def recompress_image(image, label): #This function then recompresses the processed images to JPEG
    uint8_image = tf.cast(image, tf.uint8)
    jpeg_image = tf.image.encode_jpeg(uint8_image, optimize_size=True, chroma_downsampling=False)
    return jpeg_image, label

def load_image_jpeg_dataset(): #We load the dataset of the images from the JPEG files which are stored in a Goo
    GCS_PATTERN = 'gs://flowers-public/*/*.jpg'
    files_dataset = tf.data.Dataset.list_files(GCS_PATTERN)
```

```

        decoded_images_dataset = files_dataset.map(decode_jpeg_and_label)
        resized_dataset = decoded_images_dataset.map(resize_and_crop_image)
        image_jpeg_dataset = resized_dataset.map(recompress_image)
        return image_jpeg_dataset

def parameter_combinations(batch_sizes, batch_numbers, repetitions): #This function is used to generate all of
    return list(itertools.product(batch_sizes, batch_numbers, repetitions))

#Task 2a) --> ii) Get a Spark context and create the dataset and run timing test for each combination in parallel
sc = SparkContext.getOrCreate() #We set up the spark context and then define the parameter combinations

#These are the batch sizes, batch numbers and repetitions that we will be testing
batch_sizes = [2, 4]
batch_numbers = [3, 6, 9, 12]
repetitions = [1, 2, 3]

##TASK 2C --> We add the cache function to the RDD
parameters_rdd = sc.parallelize(parameter_combinations(batch_sizes, batch_numbers, repetitions)).cache()

#This function is created to time the processing of the TFRecord batches
def time_configs(params):
    batch_size, batch_number, repetition = params
    TRF_dataset = load_dataset(filenames)
    results = []
    for repet in range(repetition):
        start_time = time.time()
        for _ in TRF_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        results.append(('TFRecord', params, images_per_second))
    return results

#This function is created to time the processing of the image data batches
def image_configs(params):
    batch_size, batch_number, repetition = params
    img_dataset = load_image_jpeg_dataset()
    image_processing_results = []
    for rep in range(repetition):
        start_time = time.time()
        for _ in img_dataset.batch(batch_size).take(batch_number):
            pass
        total_time = time.time() - start_time
        images_per_second = (batch_size * batch_number) / total_time
        image_processing_results.append(('Image', params, images_per_second))
    return image_processing_results

#These are the results from the time processing of the TFRecord and image data batches
resulting_rdd = parameters_rdd.flatMap(time_configs)
resulting_rdd_2 = parameters_rdd.flatMap(image_configs)

#Task 2a) iii) --> Transform the resulting RDD to the structure ( parameter_combination, images_per_second ) and collect

# We merge the 2 resulting RDDs (TFRecord and image results) into a combined RDD using the .union function
results_rdd_combined = resulting_rdd.union(resulting_rdd_2)
final_results = results_rdd_combined.collect()

#We print the results
print('Combined Results:')
for result in final_results:
    print(f'Dataset: {result[0]}, Parameters: (Batch Size, Batch Number, Repetitions): {result[1]}, Average Image

#Task 2a --> iv) Create an RDD with all results for each parameter as (parameter_value, images_per_second) and collect

#This function extracts the results for each parameter
def param_results(index_param):
    return results_rdd_combined.map(lambda x: (x[1][index_param], x[2]))

#We create separate RDDs for each parameter
rdd_batch_size = param_results(0)
rdd_batch_number = param_results(1)
rdd_repetition = param_results(2)

#This function is created to calculate the average images per second for each parameter separately
def average_images_param(data):
    total_images = sum(images for images in data)
    count = len(data)
    return total_images / count

avg_batch_size_rdd = rdd_batch_size.groupByKey().mapValues(average_images_param).collect()
avg_batch_number_rdd = rdd_batch_number.groupByKey().mapValues(average_images_param).collect()
avg_repetition_rdd = rdd_repetition.groupByKey().mapValues(average_images_param).collect()

#We print the average images per second for each unique parameter
print("Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec):", avg_batch_size_rdd)
print("Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec):", avg_batch_number_rdd)
print("Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec):", avg_repetition_rdd)

```

```

#Task 2a --> v) Create an RDD with the average reading speeds for each parameter value and collect the results

#This function calculates the average reading speed using a tuple (total, count) for the associativity
def seq_op(accumu, n):
    return (accumu[0] + n, accumu[1] + 1)

def comb_op(accumu1, accumu2):
    return (accumu1[0] + accumu2[0], accumu1[1] + accumu2[1])

#Now we can create a function that calculates the average images per second while also keeping associativity in
def average_assoc(rdd_params):
    return rdd_params.aggregateByKey((0,0), seq_op, comb_op).mapValues(lambda x: x[0] / x[1]).collect()

avg_batch_size = average_assoc(rdd_batch_size)
avg_batch_number = average_assoc(rdd_batch_number)
avg_repetition = average_assoc(rdd_repetition)

#Now we print the average reading speed for the associativity-respecting method
print("Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec):",
print("Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec)
print("Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec):", avg_

#Task 2a --> vi) Write the results to a pickle file in your bucket
import pickle
from google.cloud import storage

#We store the results from previous tasks to a dictionary
results_task2c = {'Final Results Task iii': final_results,
                  'AVG Batch Size Task iv': avg_batch_size_rdd,
                  'AVG Batch Number Task iv': avg_batch_number_rdd,
                  'AVG Repetition Task iv': avg_repetition_rdd,
                  'AVG Batch Size Task v': avg_batch_size,
                  'AVG Batch Number Task v': avg_batch_number,
                  'AVG Repetition Task v': avg_repetition}

#We write the results to a local pickle file
with open('/tmp/all_results.pkl', 'wb') as f:
    pickle.dump(results_task2c, f)

#This function is used to upload the pickle file to the bucket that we specify
def upload_to_bucket(blob_name, path_to_file, bucket_name):
    try:
        client = storage.Client()
        bucket = client.get_bucket(bucket_name)
        blob = bucket.blob(blob_name)
        blob.upload_from_filename(path_to_file)
        print("File uploaded successfully")
    except Exception as e:
        print(f"An error occurred: {e}")

upload_to_bucket('results/results_task2c.pkl', '/tmp/all_results.pkl', 'daring-healer-421511-storage')

```

Writing spark_job_2c.py

```

In [ ]: #We will set up the maximal cluster which is a cluster with 1 machine and 7 workers using maximal SSD size (100
CLUSTER = '{}-maximalcluster7task2c'.format(PROJECT)
REGION = 'us-west1'
PIP_PACKAGES = 'tensorflow==2.4.0'
!gcloud dataproc clusters create $CLUSTER \
    --region $REGION \
    --bucket $PROJECT-storage \
    --image-version 1.5-ubuntu18 \
    --master-machine-type n1-standard-1 \
    --master-boot-disk-type pd-ssd --master-boot-disk-size 100\
    --num-workers 7 --worker-machine-type n1-standard-1 --worker-boot-disk-size 100 \
    --initialization-actions gs://goog-dataproc-initialization-actions-$REGION/python/pip-install.sh \
    --metadata PIP_PACKAGES='tensorflow==2.4.0 matplotlib google-cloud-storage==1.42.0 google-api-core==1.31.0
#Reference link for code: https://cloud.google.com/dataproc/docs/guides/create-cluster

```

Waiting on operation [projects/daring-healer-421511/regions/us-west1/operations/de2f5a46-f3e0-3ee6-b3dc-0387d831ed04].

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

WARNING: Don't create production clusters that reference initialization actions located in the gs://goog-dataproc-initialization-actions-REGION public buckets. These scripts are provided as reference implementations, and they are synchronized with ongoing GitHub repository changes—a new version of a initialization action in public buckets may break your cluster creation. Instead, copy the following initialization actions from public buckets into your bucket : gs://goog-dataproc-initialization-actions-us-west1/python/pip-install.sh

WARNING: Failed to validate permissions required for default service account: '869525608453-compute@developer.gserviceaccount.com'. Cluster creation could still be successful if required permissions have been granted to the respective service accounts as mentioned in the document https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/service-accounts#dataproc_service_accounts_2. This could be due to Cloud Resource Manager API hasn't been enabled in your project '869525608453' before or it is disabled. Enable it by visiting 'https://console.developers.google.com/apis/api/cloudresourcemanager.googleapis.com/overview?project=869525608453'.

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

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

WARNING: The specified custom staging bucket 'daring-healer-421511-storage' is not using uniform bucket level access IAM configuration. It is recommended to update bucket to enable the same. See <https://cloud.google.com/storage/docs/uniform-bucket-level-access>.

Created [https://dataproc.googleapis.com/v1/projects/daring-healer-421511/regions/us-west1/clusters/daring-healer-421511-maximalcluster7task2c] Cluster placed in zone [us-west1-c].

```
In [ ]: #We get information on the single machine cluster we created
!gcloud dataproc clusters describe $CLUSTER --region $REGION
```

```
clusterName: daring-healer-421511-maximalcluster7task2c
clusterUuid: e629e54c-5347-4af4-ab88-888fe4773d32
config:
  configBucket: daring-healer-421511-storage
  endpointConfig: {}
  gceClusterConfig:
    internalIpOnly: false
  metadata:
    PIP_PACKAGES: tensorflow==2.4.0 matplotlib google-cloud-storage==1.42.0 google-api-core==1.31.0
    protobuf==3.17.3
  networkUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/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
    - https://www.googleapis.com/auth/monitoring.write
  zoneUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c
  initializationActions:
    - executableFile: gs://goog-dataproc-initialization-actions-us-west1/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-1-5-ubuntu18-2030909-165100-rc01
    instanceNames:
      - daring-healer-421511-maximalcluster7task2c-m
    machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machineTypes/n1-standard-1
    minCpuPlatform: AUTOMATIC
    numInstances: 1
    preemptibility: NON_PREEMPTIBLE
  softwareConfig:
    imageVersion: 1.5.90-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
      distcp:mapreduce.map.java.opts: -Xmx576m
      distcp:mapreduce.map.memory.mb: '768'
      distcp:mapreduce.reduce.java.opts: -Xmx576m
      distcp:mapreduce.reduce.memory.mb: '768'
      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: '60'
      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: daring-healer-421511-maximalcluster7task2c-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: '30'
```

```

hdfs:dfs.namenode.servicerpc-address: daring-healer-421511-maximalcluster7task2c-m:8051
hive:hive.fetch.task.conversion: none
mapred-env:HADOOP_JOB_HISTORYSERVER_HEAPSIZE: '1000'
mapred:mapreduce.job.maps: '60'
mapred:mapreduce.job.reduce.slowstart.completedmaps: '0.95'
mapred:mapreduce.job.reduces: '7'
mapred:mapreduce.jobhistory.recovery.store.class: org.apache.hadoop.mapreduce.v2.hs.HistoryServerLevelDbS
tateStoreService
mapred:mapreduce.map.cpu.vcores: '1'
mapred:mapreduce.map.java.opts: -Xmx819m
mapred:mapreduce.map.memory.mb: '1024'
mapred:mapreduce.reduce.cpu.vcores: '1'
mapred:mapreduce.reduce.java.opts: -Xmx1638m
mapred:mapreduce.reduce.memory.mb: '2048'
mapred:mapreduce.task.io.sort.mb: '256'
mapred:yarn.app.mapreduce.am.command-opts: -Xmx819m
mapred:yarn.app.mapreduce.am.resource.cpu-vcores: '1'
mapred:yarn.app.mapreduce.am.resource.mb: '1024'
spark-env:SPARK_DAEMON_MEMORY: 1000m
spark:spark.driver.maxResultSize: 480m
spark:spark.driver.memory: 960m
spark:spark.executor.cores: '1'
spark:spark.executor.instances: '2'
spark:spark.executor.memory: 2688m
spark:spark.executorEnv.OPENBLAS_NUM_THREADS: '1'
spark:spark.extraListeners: com.google.cloud.spark.performance.DataprocMetricsListener
spark:spark.scheduler.mode: FAIR
spark:spark.sql.cbo.enabled: 'true'
spark:spark.ui.port: '0'
spark:spark.yarn.am.memory: 640m
yarn-env:YARN_NODEMANAGER_HEAPSIZE: '1000'
yarn-env:YARN_RESOURCEMANAGER_HEAPSIZE: '1000'
yarn-env:YARN_TIMELINESERVER_HEAPSIZE: '1000'
yarn:yarn.nodemanager.address: 0.0.0.0:8026
yarn:yarn.nodemanager.resource.cpu-vcores: '1'
yarn:yarn.nodemanager.resource.memory-mb: '3072'
yarn:yarn.resourcemanager.nodemanager-graceful-decommission-timeout-secs: '86400'
yarn:yarn.scheduler.maximum-allocation-mb: '3072'
yarn:yarn.scheduler.minimum-allocation-mb: '256'
tempBucket: dataproc-temp-us-west1-869525608453-ejfnpe5x
workerConfig:
  diskConfig:
    bootDiskSizeGb: 100
    bootDiskType: pd-standard
    imageUri: https://www.googleapis.com/compute/v1/projects/cloud-dataproc/global/images/dataproc-1-5-ubu18-20
230909-165100-rc01
  instanceNames:
    - daring-healer-421511-maximalcluster7task2c-w-0
    - daring-healer-421511-maximalcluster7task2c-w-1
    - daring-healer-421511-maximalcluster7task2c-w-2
    - daring-healer-421511-maximalcluster7task2c-w-3
    - daring-healer-421511-maximalcluster7task2c-w-4
    - daring-healer-421511-maximalcluster7task2c-w-5
    - daring-healer-421511-maximalcluster7task2c-w-6
  machineTypeUri: https://www.googleapis.com/compute/v1/projects/daring-healer-421511/zones/us-west1-c/machin
eTypes/n1-standard-1
  minCpuPlatform: AUTOMATIC
  numInstances: 7
  preemptibility: NON_PREEMPTIBLE
labels:
  goog-dataproc-autozone: enabled
  goog-dataproc-cluster-name: daring-healer-421511-maximalcluster7task2c
  goog-dataproc-cluster-uuid: e629e54c-5347-4af4-ab88-888fe4773d32
  goog-dataproc-location: us-west1
projectId: daring-healer-421511
status:
  state: RUNNING
  stateStartTime: '2024-05-03T12:26:43.751203Z'
statusHistory:
  - state: CREATING
    stateStartTime: '2024-05-03T12:21:01.701269Z'

```

```

In [ ]: #We submit the job to the cluster
!gcloud dataproc jobs submit pyspark --cluster $CLUSTER --region $REGION ./spark_job_2c.py
%time

```

```

Job [69593d22d7bd43898ae6b01d05cf64f2] submitted.
Waiting for job output...
2024-05-03 12:27:00.352146: W tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dyna
mic library 'libcudart.so.11.0'; dLError: libcudart.so.11.0: cannot open shared object file: No such file or di
rectory; LD_LIBRARY_PATH: /usr/lib/hadoop/lib/native
2024-05-03 12:27:00.352390: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dLError if
you do not have a GPU set up on your machine.
24/05/03 12:27:07 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
24/05/03 12:27:07 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
24/05/03 12:27:07 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
24/05/03 12:27:08 INFO org.spark_project.jetty.util.log: Logging initialized @13032ms to org.spark_project.jett
y.util.log.Slf4jLog
24/05/03 12:27:08 INFO org.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknow

```

```
n; jvm 1.8.0_382-b05
24/05/03 12:27:08 INFO org.spark_project.jetty.server.Server: Started @13391ms
24/05/03 12:27:08 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@52e43b3d{HTTP/
1.1, (http/1.1)}{0.0.0.0:40269}
24/05/03 12:27:11 INFO org.apache.hadoop.yarn.client.RMPProxy: Connecting to ResourceManager at daring-healer-42
1511-maximalcluster7task2c-m/10.138.0.61:8032
24/05/03 12:27:11 INFO org.apache.hadoop.yarn.client.AHSPProxy: Connecting to Application History server at dari
ng-healer-421511-maximalcluster7task2c-m/10.138.0.61:10200
24/05/03 12:27:12 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
24/05/03 12:27:12 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
24/05/03 12:27:12 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory
-mb, units = Mi, type = COUNTABLE
24/05/03 12:27:12 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores
, units = , type = COUNTABLE
24/05/03 12:27:14 WARN org.apache.hadoop.hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
    at java.lang.Object.wait(Native Method)
    at java.lang.Thread.join(Thread.java:1257)
    at java.lang.Thread.join(Thread.java:1331)
    at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:980)
    at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:630)
    at org.apache.hadoop.hdfs.DataStreamer.run(DataStreamer.java:807)
24/05/03 12:27:17 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application
_1714739042107_0001
Combined Results:
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 1
7.953838778165878
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 1
8.369925259882315
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 2
0.27824074596342
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
9.35515488616096
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
7.48685592942393
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 1
9.704195411613547
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 3
7.18915084720215
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 3
7.84741366928925
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 3
4.96822046556997
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 4
0.629882207067574
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 3
5.216830756937654
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 3
8.56707140459847
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 4
7.97504959705964
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 5
5.96924331497283
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 6
0.37435915070024
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 2
3.752435335368876
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 5
8.262224286977684
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 5
1.887690659284395
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second:
62.34862928555369
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second:
81.52854620555601
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second:
76.81682008923768
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
77.89962374769583
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
71.26241777410284
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second:
81.78887688195195
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 3
2.09494640708425
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 3
7.04462874151012
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 3
1.921356673488763
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 4
0.088799398488575
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
4.3754318985094
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 3
7.726562213059275
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 7
0.82470579792331
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7
8.82591601000131
```


Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 7
3.55718068781677
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 7
3.83195236066014
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 6
4.7188819803098
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 7
7.03792229419471
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 8
7.67362950206098
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 1
22.97907982505437
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 1
09.3887240356083
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 9
8.46721453659516
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 1
15.09173270688812
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 9
9.85797471987577
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second:
133.80347244929698
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second:
140.67086317740004
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second:
139.4111656848699
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second:
136.45633973886262
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second:
125.38057100722166
Dataset: TFRecord, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second:
144.78969216297153
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 1), Average Images per Second: 4.06
4098414980844
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 7.81
5475719091838
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 2), Average Images per Second: 8.50
1094482188316
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 8.12
7611432476504
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 7.92
70627012806365
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 3, 3), Average Images per Second: 7.74
3485331245086
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 1), Average Images per Second: 8.31
7956662435384
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 8.39
149607323187
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 2), Average Images per Second: 8.91
3268221794803
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 8.39
626679469456
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 8.70
0475079132318
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 6, 3), Average Images per Second: 8.68
4269753263008
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 1), Average Images per Second: 8.04
3541842414175
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 8.48
4966103392155
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 2), Average Images per Second: 8.24
7591595302582
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 8.49
7821051957418
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 6.82
9483878169507
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 9, 3), Average Images per Second: 8.51
9150732970521
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 1), Average Images per Second: 8.0
66017058113966
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 8.2
92142739464232
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 2), Average Images per Second: 9.3
50642633595742
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 9.0
19686124125132
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 8.6
63754400158362
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (2, 12, 3), Average Images per Second: 8.3
21030140551866
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 1), Average Images per Second: 7.99
7422095615084
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 8.43
0377711434254
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 2), Average Images per Second: 9.23
6254479198188
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 8.32
24005296531
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 9.12

8501993495815
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 3, 3), Average Images per Second: 9.08
1366414911935
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 1), Average Images per Second: 9.23
8312577331396
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 9.41
9957547009579
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 2), Average Images per Second: 9.37
4353276013142
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 8.95
8867542238396
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 9.89
7539923408175
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 6, 3), Average Images per Second: 9.22
5686013058693
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 1), Average Images per Second: 8.92
81531597319
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 8.83
297536160838
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 2), Average Images per Second: 9.62
8148701342425
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 9.31
1832188998288
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 9.76
4547156549492
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 9, 3), Average Images per Second: 9.16
865751604007
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 1), Average Images per Second: 9.6
24447629435906
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 9.6
36691213608945
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 2), Average Images per Second: 9.8
95435430121797
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 9.9
55975575525757
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 9.8
23049979595346
Dataset: Image, Parameters: (Batch Size, Batch Number, Repetitions): (4, 12, 3), Average Images per Second: 9.4
7422823238864
Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec): [(4, 50.779586230278774), (2, 28.62060671733013)]
Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec): [(12, 58.10352390521206), (9, 46.57077368978701), (6, 33.20917095792271), (3, 18.672581908502966)]
Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec): [(1, 39.902618464093095), (2, 38.518019118374426), (3, 39.9435729524705)]
Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec): [(4, 48.79294820323631), (2, 28.13458388555063)]
Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec): [(12, 56.804326011146564), (9, 44.49631314070965), (6, 31.256141239586494), (3, 18.08197688324735)]
Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec): [(1, 37.88179564401834), (2, 37.09477353654646), (3, 39.47990942038509)]
File uploaded successfully
24/05/03 12:38:21 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@52e43b3d{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
Job [69593d22d7bd43898ae6b01d05cf64f2] finished successfully.
done: true
driverControlFilesUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/e629e54c-5347-4af4-ab88-888fe4773d32/jobs/69593d22d7bd43898ae6b01d05cf64f2/
driverOutputResourceUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/e629e54c-5347-4af4-ab88-888fe4773d32/jobs/69593d22d7bd43898ae6b01d05cf64f2/driveroutput
jobUuid: 215ffd5d-4bc0-381a-a2d5-c18c0e27bd29
placement:
 clusterName: daring-healer-421511-maximalcluster7task2c
 clusterUuid: e629e54c-5347-4af4-ab88-888fe4773d32
pysparkJob:
 mainPythonFileUri: gs://daring-healer-421511-storage/google-cloud-dataproc-metainfo/e629e54c-5347-4af4-ab88-888fe4773d32/jobs/69593d22d7bd43898ae6b01d05cf64f2/staging/spark_job_2c.py
reference:
 jobId: 69593d22d7bd43898ae6b01d05cf64f2
 projectId: daring-healer-421511
status:
 state: DONE
 stateStartTime: '2024-05-03T12:38:24.247580Z'
statusHistory:
 - state: PENDING
 stateStartTime: '2024-05-03T12:26:51.422476Z'
 - state: SETUP_DONE
 stateStartTime: '2024-05-03T12:26:51.455865Z'
 - details: Agent reported job success
 state: RUNNING
 stateStartTime: '2024-05-03T12:26:51.874547Z'
yarnApplications:
 - name: spark_job_2c.py
 progress: 1.0
 state: FINISHED
 trackingUrl: http://daring-healer-421511-maximalcluster7task2c-m:8088/proxy/application_1714739042107_0001/
CPU times: user 82 µs, sys: 0 ns, total: 82 µs
Wall time: 181 µs

```
In [ ]: CLUSTER = '{}-maximalcluster7task2c'.format(PROJECT)
REGION = 'us-west1'
!gcloud dataproc clusters delete $CLUSTER --region $REGION
#Due to the fact that we have quota related problems as we seem to have few quota available, after creating a c
```

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

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

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

Also, **plot** the output values, the averages per parameter value and the regression lines for each parameter and for the product of batch_size and batch_number

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

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

By **parallelising** the speed test we are making **assumptions** about the limits of the bucket reading speeds. See [here](#) for more information. Discuss, **what we need to consider in speed tests** in parallel on the cloud, which bottlenecks we might be identifying, and how this relates to your results.

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

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

```
In [ ]: !gsutil cp gs://daring-healer-421511-storage/results/results_task2c.pkl ./results_task2c.pkl #We first read the
Copying gs://daring-healer-421511-storage/results/results_task2c.pkl...
/ [1 files][ 2.5 KiB/ 2.5 KiB]
Operation completed over 1 objects/2.5 KiB.
```

```
In [ ]: import pandas as pd

#We load the pickle file
final_res = pd.read_pickle('results_task2c.pkl')
print(final_res)
```

```
{'Final Results Task iii': [('TFRecord', (2, 3, 1), 17.953838778165878), ('TFRecord', (2, 3, 2), 18.36992525988
2315), ('TFRecord', (2, 3, 2), 20.27824074596342), ('TFRecord', (2, 3, 3), 19.35515488616096), ('TFRecord', (2,
3, 3), 17.48685592942393), ('TFRecord', (2, 3, 3), 19.704195411613547), ('TFRecord', (2, 6, 1), 37.189150847202
15), ('TFRecord', (2, 6, 2), 37.84741366928925), ('TFRecord', (2, 6, 2), 34.96822046556997), ('TFRecord', (2, 6
, 3), 40.629882207067574), ('TFRecord', (2, 6, 3), 35.216830756937654), ('TFRecord', (2, 6, 3), 38.567071404598
47), ('TFRecord', (2, 9, 1), 47.97504959705964), ('TFRecord', (2, 9, 2), 55.96924331497283), ('TFRecord', (2, 9
, 2), 60.37435915070024), ('TFRecord', (2, 9, 3), 23.752435335368876), ('TFRecord', (2, 9, 3), 58.2622242869776
84), ('TFRecord', (2, 9, 3), 51.887690659284395), ('TFRecord', (2, 12, 1), 62.34862928555369), ('TFRecord', (2,
12, 2), 81.52854620555601), ('TFRecord', (2, 12, 2), 76.81682008923768), ('TFRecord', (2, 12, 3), 77.8996237476
9583), ('TFRecord', (2, 12, 3), 71.26241777410284), ('TFRecord', (2, 12, 3), 81.78887688195195), ('TFRecord', (
4, 3, 1), 32.09494640708425), ('TFRecord', (4, 3, 2), 37.04462874151012), ('TFRecord', (4, 3, 2), 31.9213566734
88763), ('TFRecord', (4, 3, 3), 40.088799398488575), ('TFRecord', (4, 3, 3), 34.3754318985094), ('TFRecord', (4
, 3, 3), 37.726562213059275), ('TFRecord', (4, 6, 1), 70.82470579792331), ('TFRecord', (4, 6, 2), 78.8259160100
0131), ('TFRecord', (4, 6, 2), 73.55718068781677), ('TFRecord', (4, 6, 3), 73.83195236066014), ('TFRecord', (4,
6, 3), 64.7188819803098), ('TFRecord', (4, 6, 3), 77.03792229419471), ('TFRecord', (4, 9, 1), 87.67362950206098
), ('TFRecord', (4, 9, 2), 122.97907982505437), ('TFRecord', (4, 9, 2), 109.3887240356083), ('TFRecord', (4, 9,
3), 98.46721453659516), ('TFRecord', (4, 9, 3), 115.09173270688812), ('TFRecord', (4, 9, 3), 99.85797471987577)
, ('TFRecord', (4, 12, 1), 133.80347244929698), ('TFRecord', (4, 12, 2), 140.67086317740004), ('TFRecord', (4,
12, 2), 139.4111656848699), ('TFRecord', (4, 12, 3), 136.45633973886262), ('TFRecord', (4, 12, 3), 125.38057100
722166), ('TFRecord', (4, 12, 3), 144.78969216297153), ('Image', (2, 3, 1), 4.064098414980844), ('Image', (2, 3
, 2), 7.815475719091838), ('Image', (2, 3, 2), 8.501094482188316), ('Image', (2, 3, 3), 8.127611432476504), ('I
mage', (2, 3, 3), 7.9270627012806365), ('Image', (2, 3, 3), 7.743485331245086), ('Image', (2, 6, 1), 8.31795666
2435384), ('Image', (2, 6, 2), 8.39149607323187), ('Image', (2, 6, 2), 8.913268221794803), ('Image', (2, 6, 3),
8.39626679469456), ('Image', (2, 6, 3), 8.700475079132318), ('Image', (2, 6, 3), 8.684269753263008), ('Image',
(2, 9, 1), 8.043541842414175), ('Image', (2, 9, 2), 8.484966103392155), ('Image', (2, 9, 2), 8.247591595302582)
, ('Image', (2, 9, 3), 8.497821051957418), ('Image', (2, 9, 3), 6.829483878169507), ('Image', (2, 9, 3), 8.5191
50732970521), ('Image', (2, 12, 1), 8.066017058113966), ('Image', (2, 12, 2), 8.292142739464232), ('Image', (2,
12, 2), 9.350642633595742), ('Image', (2, 12, 3), 9.019686124125132), ('Image', (2, 12, 3), 8.663754400158362),
('Image', (2, 12, 3), 8.321030140551866), ('Image', (4, 3, 1), 7.997422095615084), ('Image', (4, 3, 2), 8.43037
7711434254), ('Image', (4, 3, 2), 9.236254479198188), ('Image', (4, 3, 3), 8.3224005296531), ('Image', (4, 3, 3
), 9.128501993495815), ('Image', (4, 3, 3), 9.081366414911935), ('Image', (4, 6, 1), 9.238312577331396), ('Imag
e', (4, 6, 2), 9.419957547009579), ('Image', (4, 6, 2), 9.374353276013142), ('Image', (4, 6, 3), 8.958867542238
396), ('Image', (4, 6, 3), 9.897539923408175), ('Image', (4, 6, 3), 9.225686013058693), ('Image', (4, 9, 1), 8.
9281531597319), ('Image', (4, 9, 2), 8.83297536160838), ('Image', (4, 9, 2), 9.628148701342425), ('Image', (4,
9, 3), 9.311832188998288), ('Image', (4, 9, 3), 9.764547156549492), ('Image', (4, 9, 3), 9.16865751604007), ('I
mage', (4, 12, 1), 9.624447629435906), ('Image', (4, 12, 2), 9.636691213608945), ('Image', (4, 12, 2), 9.895435
430121797), ('Image', (4, 12, 3), 9.955975575525757), ('Image', (4, 12, 3), 9.823049979595346), ('Image', (4, 1
2, 3), 9.47422823238864)], 'AVG Batch Size Task iv': [(4, 50.779586230278774), (2, 28.62060671733013)], 'AVG Ba
tch Number Task iv': [(12, 58.10352390521206), (9, 46.57077368978701), (6, 33.20917095792271), (3, 18.672581908
502966)], 'AVG Repetition Task iv': [(1, 39.902618464093095), (2, 38.518019118374426), (3, 39.9435729524705)],
'AVG Batch Size Task v': [(4, 48.79294820323631), (2, 28.13458388555063)], 'AVG Batch Number Task v': [(12, 56.
804326011146564), (9, 44.49631314070965), (6, 31.256141239586494), (3, 18.08197688324735)], 'AVG Repetition Tas
k v': [(1, 37.88179564401834), (2, 37.09477353654646), (3, 39.47990942038509)]}
```

```
In [ ]: #We store the results into variables to then put them in dataframes
final_res_task_iii = final_res['Final Results Task iii']
avg_batch_size_task_iv = final_res['AVG Batch Size Task iv']
avg_batch_number_task_iv = final_res['AVG Batch Number Task iv']
avg_repetition_task_iv = final_res['AVG Repetition Task iv']
avg_batch_size_task_v = final_res['AVG Batch Size Task v']
avg_batch_number_task_v = final_res['AVG Batch Number Task v']
avg_repetition_task_v = final_res['AVG Repetition Task v']

#We print out all of the results we got
print('Final Results')
print(final_res_task_iii)
print("Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec)")
print(avg_batch_size_task_iv)
print("Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec)")
print(avg_batch_number_task_iv)
print("Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec)")
print(avg_repetition_task_iv)
print("Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec)")
print(avg_batch_size_task_v)
print("Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec)")
print(avg_batch_number_task_v)
print("Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec)")
print(avg_repetition_task_v)
```

Final Results

```
[('TFRecord', (2, 3, 1), 17.953838778165878), ('TFRecord', (2, 3, 2), 18.369925259882315), ('TFRecord', (2, 3, 2), 20.27824074596342), ('TFRecord', (2, 3, 3), 19.35515488616096), ('TFRecord', (2, 3, 3), 17.48685592942393), ('TFRecord', (2, 3, 3), 19.704195411613547), ('TFRecord', (2, 6, 1), 37.18915084720215), ('TFRecord', (2, 6, 2), 37.84741366928925), ('TFRecord', (2, 6, 2), 34.96822046556997), ('TFRecord', (2, 6, 3), 40.629882207067574), ('TFRecord', (2, 6, 3), 35.216830756937654), ('TFRecord', (2, 6, 3), 38.56707140459847), ('TFRecord', (2, 9, 1), 47.97504959705964), ('TFRecord', (2, 9, 2), 55.96924331497283), ('TFRecord', (2, 9, 2), 60.37435915070024), ('TFRecord', (2, 9, 3), 23.752435335368876), ('TFRecord', (2, 9, 3), 58.262224286977684), ('TFRecord', (2, 9, 3), 51.887690659284395), ('TFRecord', (2, 12, 1), 62.34862928555369), ('TFRecord', (2, 12, 2), 81.52854620555601), ('TFRecord', (2, 12, 2), 76.81682008923768), ('TFRecord', (2, 12, 3), 77.89962374769583), ('TFRecord', (2, 12, 3), 71.26241777410284), ('TFRecord', (2, 12, 3), 81.78887688195195), ('TFRecord', (4, 3, 1), 32.0949464070842), ('TFRecord', (4, 3, 2), 37.04462874151012), ('TFRecord', (4, 3, 2), 31.921356673488763), ('TFRecord', (4, 3, 3), 40.088799398488575), ('TFRecord', (4, 3, 3), 34.3754318985094), ('TFRecord', (4, 3, 3), 37.72656221305927), ('TFRecord', (4, 6, 1), 70.82470579792331), ('TFRecord', (4, 6, 2), 78.82591601000131), ('TFRecord', (4, 6, 2), 73.55718068781677), ('TFRecord', (4, 6, 3), 73.83195236066014), ('TFRecord', (4, 6, 3), 64.7188819803098), ('TFRecord', (4, 12, 3), 136.45633973886262), ('TFRecord', (4, 9, 1), 87.67362950206098), ('TFRecord', (4, 9, 2), 122.97907982505437), ('TFRecord', (4, 9, 2), 109.3887240356083), ('TFRecord', (4, 9, 3), 98.46721453659516), ('TFRecord', (4, 9, 3), 115.09173270688812), ('TFRecord', (4, 9, 3), 99.85797471987577), ('TFRecord', (4, 12, 1), 133.80347244929698), ('TFRecord', (4, 12, 2), 140.67086317740004), ('TFRecord', (4, 12, 2), 139.4111656848699), ('TFRecord', (4, 12, 3), 125.38057100722166), ('TFRecord', (4, 12, 3), 144.78969216297153), ('Image', (2, 3, 1), 4.064098414980844), ('Image', (2, 3, 2), 7.815475719091838), ('Image', (2, 3, 2), 8.501094482188316), ('Image', (2, 3, 3), 8.127611432476504), ('Image', (2, 3, 3), 7.9270627012806365), ('Image', (2, 3, 3), 7.743485331245086), ('Image', (2, 6, 1), 8.317956662435384), ('Image', (2, 6, 2), 8.39149607323187), ('Image', (2, 6, 2), 8.913268221794803), ('Image', (2, 6, 3), 8.39626679469456), ('Image', (2, 6, 3), 8.700475079132318), ('Image', (2, 6, 3), 8.684269753263008), ('Image', (2, 9, 1), 8.04354184241417), ('Image', (2, 9, 2), 8.484966103392155), ('Image', (2, 9, 2), 8.247591595302582), ('Image', (2, 9, 3), 8.497821051957418), ('Image', (2, 9, 3), 6.829483878169507), ('Image', (2, 9, 3), 8.519150732970521), ('Image', (2, 12, 1), 8.066017058113966), ('Image', (2, 12, 2), 8.292142739464232), ('Image', (2, 12, 2), 9.350642633595742), ('Image', (2, 12, 3), 9.019686124125132), ('Image', (2, 12, 3), 8.663754400158362), ('Image', (2, 12, 3), 8.321030140551866), ('Image', (4, 3, 1), 7.997422095615084), ('Image', (4, 3, 2), 8.430377711434254), ('Image', (4, 3, 2), 9.236254479198188), ('Image', (4, 3, 3), 8.3224005296531), ('Image', (4, 3, 3), 9.128501993495815), ('Image', (4, 3, 3), 9.081366414911935), ('Image', (4, 6, 1), 9.238312577331396), ('Image', (4, 6, 2), 9.419957547009579), ('Image', (4, 6, 2), 9.374353276013142), ('Image', (4, 6, 3), 8.958867542238396), ('Image', (4, 6, 3), 9.897539923408175), ('Image', (4, 6, 3), 9.225686013058693), ('Image', (4, 9, 1), 8.9281531597319), ('Image', (4, 9, 2), 8.83297536160838), ('Image', (4, 9, 2), 9.628148701342425), ('Image', (4, 9, 3), 9.311832188998288), ('Image', (4, 9, 3), 9.764547156549492), ('Image', (4, 9, 3), 9.16865751604007), ('Image', (4, 12, 1), 9.624447629435906), ('Image', (4, 12, 2), 9.636691213608945), ('Image', (4, 12, 2), 9.895435430121797), ('Image', (4, 12, 3), 9.955975575525757), ('Image', (4, 12, 3), 9.823049979595346), ('Image', (4, 12, 3), 9.47422823238864)]
Average Images per Second By Each Batch Size (Batch_Size: AVG_Images_per_Sec)
[(4, 50.779586230278774), (2, 28.62060671733013)]
Average Images per Second By Each Batch Number (Batch_Number: AVG_Images_per_Sec)
[(12, 58.10352390521206), (9, 46.57077368978701), (6, 33.20917095792271), (3, 18.672581908502966)]
Average Images per Second By Each Repetition (Repetition: AVG_Images_per_Sec)
[(1, 39.902618464093095), (2, 38.518019118374426), (3, 39.9435729524705)]
Average Images per Second By Each Batch Size (Associativity-method), (Batch_Size: AVG_Images_per_Sec)
[(4, 48.79294820323631), (2, 28.13458388555063)]
Average Images per Second By Each Batch Number (Associativity-method), (Batch_Number: AVG_Images_per_Sec)
[(12, 56.804326011146564), (9, 44.49631314070965), (6, 31.256141239586494), (3, 18.08197688324735)]
Average Images per Second By Repetition (Associativity-method), (Repetition: AVG_Images_per_Sec)
[(1, 37.88179564401834), (2, 37.09477353654646), (3, 39.47990942038509)]
```

In []: **import** pandas **as** pd

```
#We create a dataset with each combination of parameters and the images per second for each dataset separately
data = {
    'Dataset': [item[0] for item in final_res_task_iii],
    'Batch Size': [item[1][0] for item in final_res_task_iii],
    'Batch Number': [item[1][1] for item in final_res_task_iii],
    'Repetitions': [item[1][2] for item in final_res_task_iii],
    'Images per Second': [item[2] for item in final_res_task_iii],
}

dataframe = pd.DataFrame(data)
dataframe['Dataset Volume (per Batch Setup)'] = dataframe['Batch Size'] * dataframe['Batch Number']
print(dataframe)

#Reference link for code: https://colab.research.google.com/github/google/eng-edu/blob/main/ml/cc/exercises/pan
```

	Dataset	Batch Size	Batch Number	Repetitions	Images per Second \
0	TFRecord	2	3	1	17.953839
1	TFRecord	2	3	2	18.369925
2	TFRecord	2	3	2	20.278241
3	TFRecord	2	3	3	19.355155
4	TFRecord	2	3	3	17.486856
...
91	Image	4	12	2	9.636691
92	Image	4	12	2	9.895435
93	Image	4	12	3	9.955976
94	Image	4	12	3	9.823050
95	Image	4	12	3	9.474228

	Dataset Volume (per Batch Setup)
0	6
1	6
2	6
3	6
4	6
...	...
91	48
92	48
93	48
94	48
95	48

[96 rows x 6 columns]

In []: *#Now we are going to perform a linear regression (e.g. using scikit-learn) over the values for each parameter a*

```
#We import the necessary libraries
from sklearn.linear_model import LinearRegression
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
from scipy.stats import t
import numpy as np
#This function is used to create the plots of the images per second by each feature and also shows the regressi
def regression_by_parameter(dataset, parameter, title):
    filtered_data = dataframe[dataframe['Dataset']==dataset] #We filter out the data for each dataset (TFRecord,

    #We fit the linear regression model on each parameter and make predictions
    X = filtered_data[[parameter]]
    y = filtered_data['Images per Second']
    linear_model = LinearRegression()
    linear_model.fit(X, y)
    model_predictions = linear_model.predict(X)

    #We calculate the metrics for the linear regression model
    slope = float(linear_model.coef_[0])
    intercept = float(linear_model.intercept_)
    residuals = y - model_predictions
    mse = mean_squared_error(y, model_predictions)
    n = len(y)
    d_o_f = n - 2 #Degrees of freedom
    t_statistic = slope / np.sqrt(mse / (n * np.var(X, ddof=1)))
    p_value = (1 - t.cdf(abs(t_statistic), d_o_f)) * 2 #We perform a two-tailed test

    #We plot the actual values vs the predicted values from the linear regression model for each parameter
    plt.figure(figsize=(8,4))
    plt.scatter(X[parameter], y, color='blue', label='Actual Values')
    plt.plot(X[parameter], model_predictions, color='red', label='Regression Line')
    plt.title(f'{title} ({dataset})')
    plt.xlabel(parameter)
    plt.ylabel('Images per Second')
    plt.legend()
    plt.grid(True)
    print(f'Slope: {slope:.2f}')
    print(f'Intercept: {intercept:.2f}')
    print('P-Value:', p_value)
    plt.show()

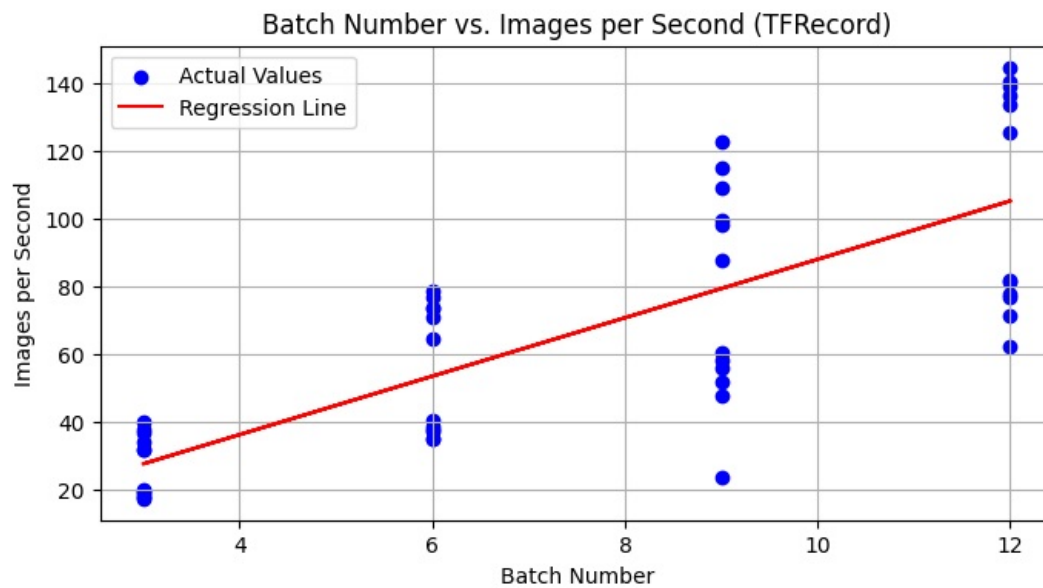
#We specify the parameters
parameters = ['Batch Number', 'Batch Size', 'Repetitions', 'Dataset Volume (per Batch Setup)']
titles = ['Batch Number vs. Images per Second', 'Batch Size vs. Images per Second',
          'Repetitions vs. Images per Second', 'Dataset Volume (per Batch Setup) vs. Images per Second']

#We generate the plots for the TFRecord Dataset
for parameter, title in zip(parameters, titles):
    regression_by_parameter('TFRecord', parameter, title)

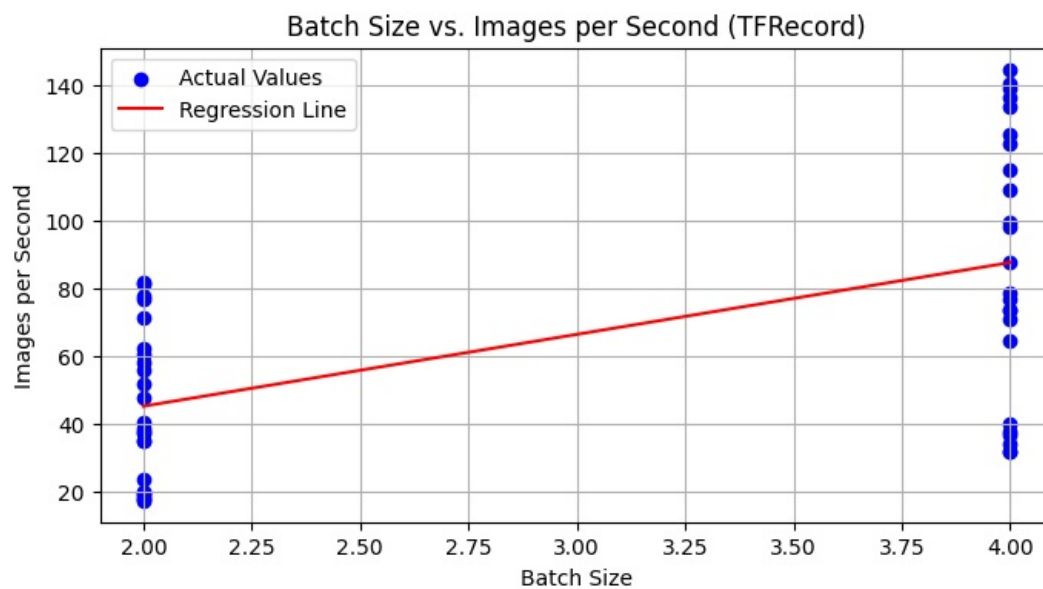
#We generate the plots for the Image Dataset
for parameter, title in zip(parameters, titles):
    regression_by_parameter('Image', parameter, title)
```

#Reference link for code: https://colab.research.google.com/github/csmastersUH/data_analysis_with_python_2020/b

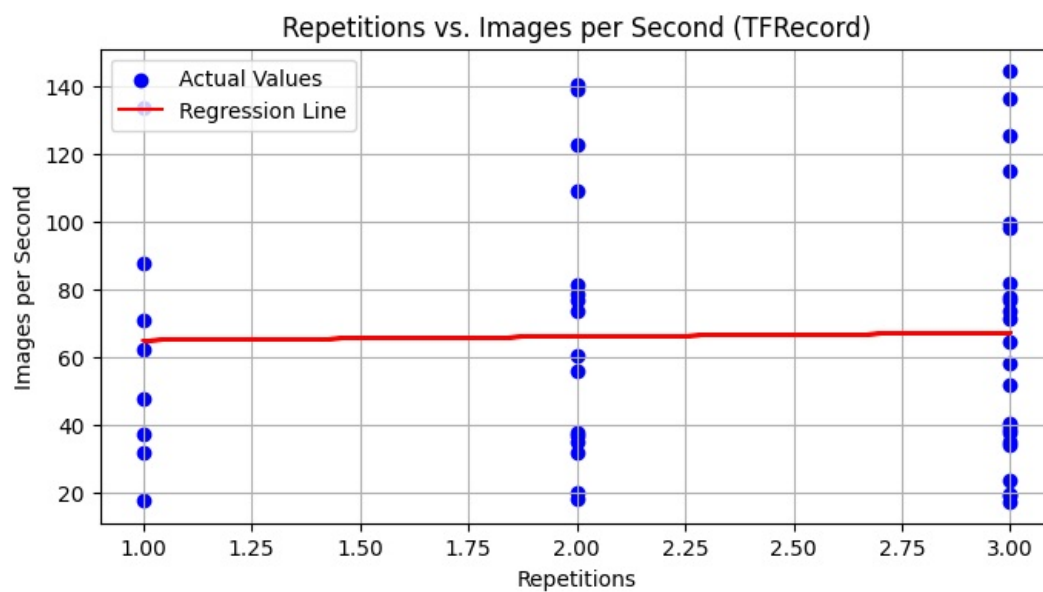
Slope: 8.63
Intercept: 1.83
P-Value: [7.48789919e-11]



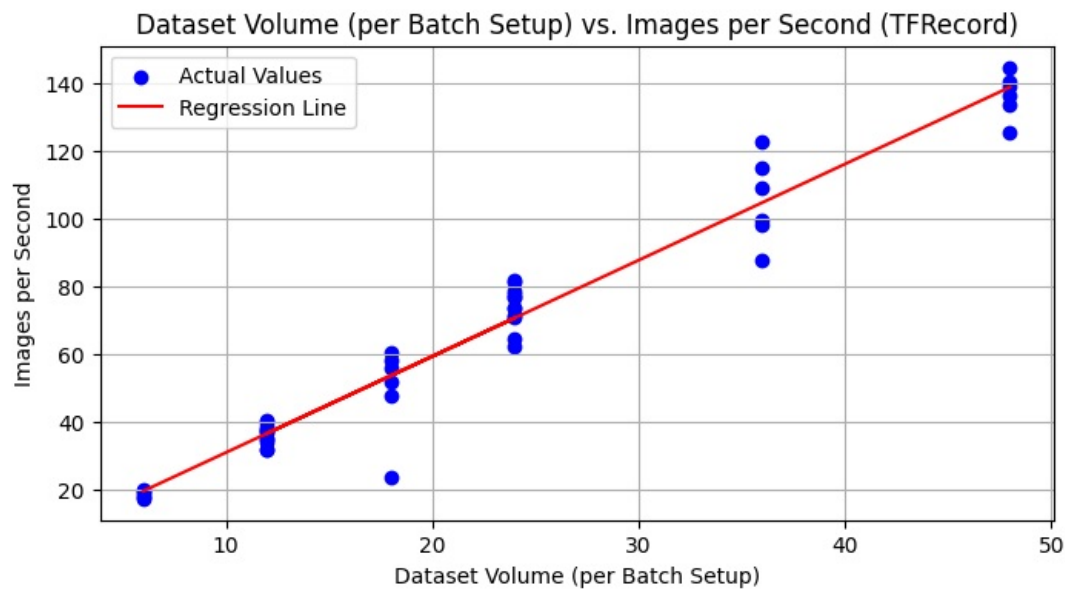
Slope: 21.22
Intercept: 2.87
P-Value: [1.85617356e-05]



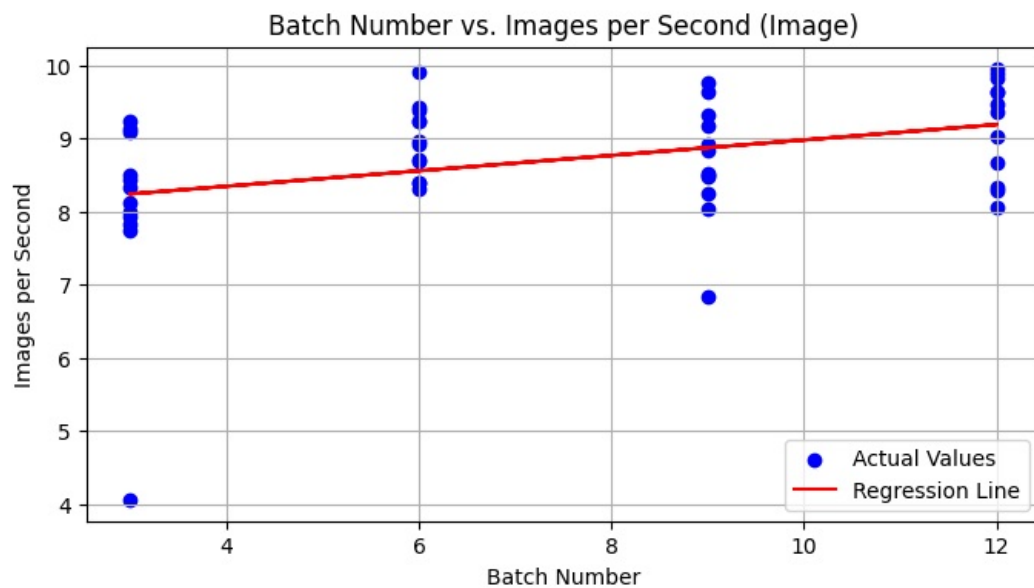
Slope: 1.10
Intercept: 63.97
P-Value: [0.8796578]



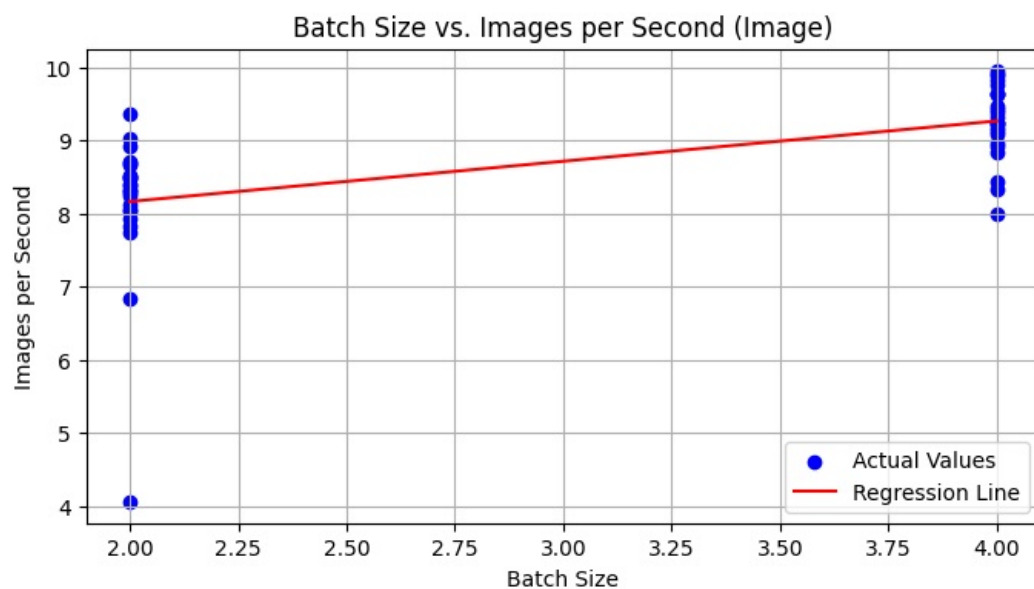
Slope: 2.84
Intercept: 2.62
P-Value: [0.]



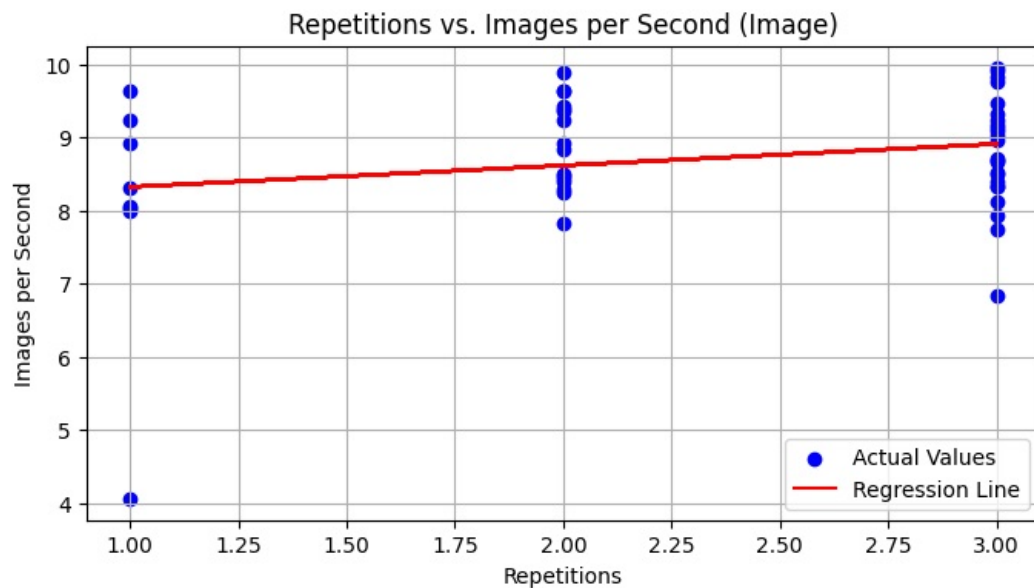
Slope: 0.11
Intercept: 7.92
P-Value: [0.00747793]



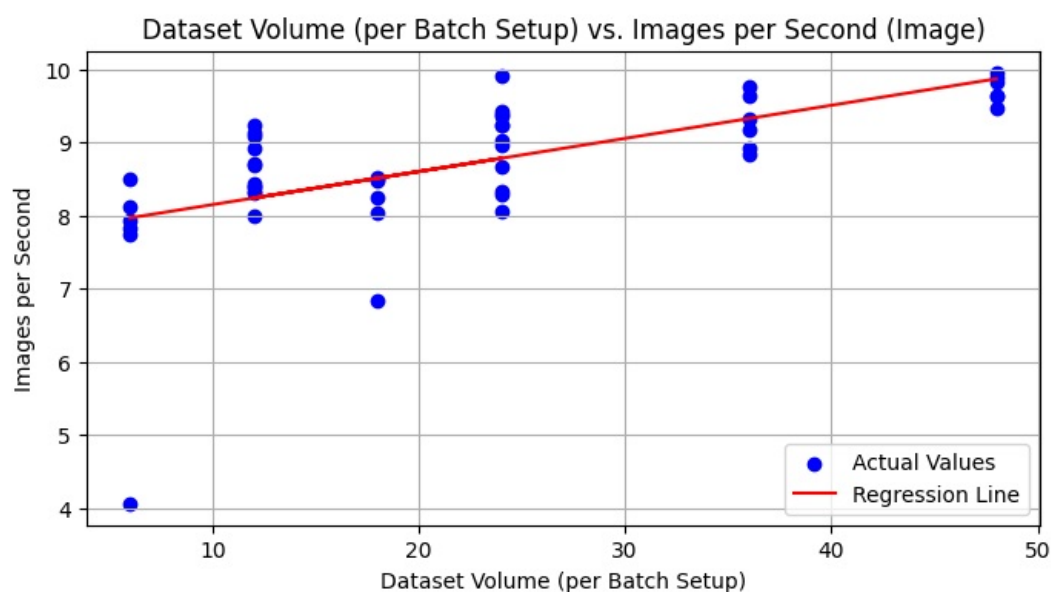
Slope: 0.55
Intercept: 7.06
P-Value: [1.02428895e-05]



Slope: 0.29
Intercept: 8.03
P-Value: [0.10514642]



Slope: 0.05
Intercept: 7.70
P-Value: [1.6921853e-06]



Section 3. Theoretical discussion

Task 3: Discussion in context. (24%)

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

- Alipourfard, O., Liu, H. H., Chen, J., Venkataraman, S., Yu, M., & Zhang, M. (2017). [Cherrypick: Adaptively unearthing the best cloud configurations for big data analytics.](#) In USENIX NSDI 17 (pp. 469-482).

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

3a) Contextualise

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

3b) Strategise

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

Provide the answers to these questions in your report.

Final cleanup

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

```
In [ ]: !gsutil -m rm -r $BUCKET/* # Empty your bucket
!gsutil rb $BUCKET # delete the bucket
```