

Setting up the GCP environment for Dataflow and running wordcount examples:

https://youtu.be/re6c_ee7uTc

Screenshots:

```
pranjal_saloni612@cloudshell:~ (projectmill1)$ python3 --version
Python 3.9.2
pranjal_saloni612@cloudshell:~ (projectmill1)$ python3 -m venv env
pranjal_saloni612@cloudshell:~ (projectmill1)$ ls
env index.html README-cloudshell.txt SOFE4630U-tut3
pranjal_saloni612@cloudshell:~ (projectmill1)$ source ~/env/bin/activate
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ deactivate
pranjal_saloni612@cloudshell:~ (projectmill1)$ source ~/env/bin/activate
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ python --version
Python 3.9.2
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ pip install pip --upgrade
Requirement already satisfied: pip in ./env/lib/python3.9/site-packages (20.3.4)
Collecting pip
  Downloading pip-22.0.4-py3-none-any.whl (2.1 MB)
    |████████████████████████████████| 2.1 MB 5.1 MB/s
Installing collected packages: pip
  Attempting uninstall: pip
    Found existing installation: pip 20.3.4
    Uninstalling pip-20.3.4:
      Successfully uninstalled pip-20.3.4
Successfully installed pip-22.0.4
```

```
(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ pip install 'apache-beam[gcp]'

Collecting apache-beam[gcp]
  Downloading apache_beam-2.37.0-cp39-cp39-manylinux2010_x86_64.whl (11.1 MB)
    11.1/11.1 MB 42.2 MB/s eta 0:00:00

Collecting pytz>=2018.3
  Downloading pytz-2022.1-py2.py3-none-any.whl (503 kB)
    503.5/503.5 kB 39.6 MB/s eta 0:00:00

Collecting requests<3.0.0,>=2.24.0
  Downloading requests-2.27.1-py2.py3-none-any.whl (63 kB)
    63.1/63.1 kB 8.9 MB/s eta 0:00:00

Collecting protobuf<4,>=3.12.2
  Downloading protobuf-3.19.4-cp39-cp39-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (1.1 MB)
    1.1/1.1 MB 59.9 MB/s eta 0:00:00

Collecting crcmod<2.0,>=1.7
  Downloading crcmod-1.7.tar.gz (89 kB)
    89.7/89.7 kB 12.2 MB/s eta 0:00:00

  Preparing metadata (setup.py) ... done

Collecting pymongo<4.0.0,>=3.8.0
  Downloading pymongo-3.12.3-cp39-cp39-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (516 kB)
    516.3/516.3 kB 39.3 MB/s eta 0:00:00

Collecting typing-extensions>=3.7.0
  Downloading typing_extensions-4.1.1-py3-none-any.whl (26 kB)

Collecting pyarrow<7.0.0,>=0.15.1
  Downloading pyarrow-6.0.1-cp39-cp39-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (25.6 MB)
    25.6/25.6 MB 31.1 MB/s eta 0:00:00

Collecting fastavro<2,>=0.23.6
  Downloading fastavro-1.4.10-cp39-cp39-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (2.5 MB)
    2.5/2.5 MB 77.7 MB/s eta 0:00:00

Collecting numpy<1.22.0,>=1.14.3
  Downloading numpy-1.21.5-cp39-cp39-manylinux_2_12_x86_64.manylinux2010_x86_64.whl (15.7 MB)
    15.7/15.7 MB 45.1 MB/s eta 0:00:00

Collecting oauth2client<5,>=2.0.1
  Downloading oauth2client-4.1.3-py2.py3-none-any.whl (99 kB)
```

```
(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ python -m apache_beam.examples.wordcount --output outputs
INFO:root:Missing pipeline option (runner). Executing pipeline using the default runner: DirectRunner.
INFO:apache_beam.internal.gcp.auth:Setting socket default timeout to 60 seconds.
INFO:apache_beam.internal.gcp.auth:socket default timeout is 60.0 seconds.
INFO:oauth2client.transport:Attempting refresh to obtain initial access token
WARNING:root:Make sure that locally built Python SDK docker image has Python 3.9 interpreter.
INFO:root:Default Python SDK image for environment is apache/beam_python3.9_sdk:2.37.0
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function annotate_downstream_side_inputs at 0x7f9ff335bdc0> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function fix_side_input_pcollections at 0x7f9ff335bee0> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function pack_combiners at 0x7f9ff335d430> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function lift_combiners at 0x7f9ff335d4c0> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function expand_sdf at 0x7f9ff335d670> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function expand_gbk at 0x7f9ff335d700> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:===== <function sink_flattens at 0x7f9ff335d820> =====
```

```
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ ls
env  index.html  outputs-00000-of-00001  README-cloudshell.txt  SOFE4630U-tut3
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ more outputs-00000-of-00001
KING: 243
LEAR: 236
DRAMATIS: 1
PERSONAE: 1
king: 65
of: 447
Britain: 2
OF: 15
FRANCE: 10
DUKE: 3
BURGUNDY: 8
CORNWALL: 63
ALBANY: 67
EARL: 2
KENT: 156
GLOUCESTER: 141
EDGAR: 126
son: 29
to: 438
Gloucester: 26
EDMUND: 99
bastard: 7
CURAN: 6
a: 366
courtier: 1
```

```
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ PROJECT=projectmill1
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ echo $PROJECT
projectmill1
```

← Bucket details

projectmill1-gs

Location	Storage class	Public access	Protection
northamerica-northeast2 (Toronto)	Standard	Not public	None

```
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ BUCKET=gs://projectmill1-gs
(env) pranjal_saloni612@cloudshell:~ (projectmill1)$ echo $BUCKET
gs://projectmill1-gs
```

```
(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ python -m apache_beam.examples.wordcount --project $PROJECT --region northamerica-northeast2 --runner Dataflowrunner --temp location $BUCKET/temp/ --output $BUCKET/result/outputs
INFO:apache_beam.internal.gcp.auth:Setting socket default timeout to 60 seconds.
INFO:apache_beam.internal.gcp.auth:socket default timeout is 60.0 seconds.
INFO:oauth2client.transport:Attempting refresh to obtain initial access_token
INFO:apache_beam.runners.portability.stager:Downloading source distribution of the SDK from PyPi
INFO:apache_beam.runners.portability.stager:Executing command: ['/home/pranjal_salon1612/env/bin/python', '-m', 'pip', 'download', '--dest', '/tmp/tmp_sxaw0w7', 'apache-beam==37.0', '--no-deps', '--no-binary', ':all:']
INFO:apache_beam.runners.portability.stager:Staging SDK sources from PyPI: dataflow_python_sdk.tar
INFO:apache_beam.runners.portability.stager:Downloading binary distribution of the SDK from PyPi
INFO:apache_beam.runners.portability.stager:Executing command: ['/home/pranjal_salon1612/env/bin/python', '-m', 'pip', 'download', '--dest', '/tmp/tmp_sxaw0w7', 'apache-beam==37.0', '--no-deps', '--only-binary', ':all:', '--python-version', '39', '--implementation', 'cp', '--abi', 'cp39', '--platform', 'manylinux1_x86_64']
INFO:apache_beam.runners.portability.stager:Staging binary distribution of the SDK from PyPI: apache-beam-2.37.0-cp39-cp39-manylinux1_x86_64.whl
WARNING:root:Make sure that locally built Python SDK docker image has Python 3.9 interpreter.
INFO:root:Default Python SDK image for environment is apache/beam_python3.9_sdk:2.37.0
```

```
prefer: 1
Beat: 1
pight: 1
beaten: 2
summit: 1
grossly: 1
striving: 1
Fairest: 1
meats: 1
glove: 2
notice: 2
encounter: 1
bold: 4
Messenger: 10
knaves: 3
passion: 4
swaggered: 1
meeting: 2
garb: 1
headlong: 1
cage: 1
needless: 1
patron: 2
spaniel: 1
FRANCE: 10
condemn'd: 1
corky: 1
dissuaded: 1
smile: 2
buz: 1
Wherefore: 5
egg: 4
despised: 2
football: 1
gracious: 1
```

```
pranjal_salon1612@cloudshell:~ (projectmill1)$ find ~/env -name 'wordcount.py'
~/home/pranjal_salon1612/env/lib/python3.9/site-packages/apache_beam/examples/dataframe/wordcount.py
~/home/pranjal_salon1612/env/lib/python3.9/site-packages/apache_beam/examples/wordcount.py
pranjal_salon1612@cloudshell:~ (projectmill1)$ ls
env index.html outputs-00000-of-00001 README-cloudshell.txt SOPF4630U-tut3
pranjal_salon1612@cloudshell:~ (projectmill1)$ cp ~/home/pranjal_salon1612/env/lib/python3.9/site-packages/apache_beam/examples/wordcount.py
cp: missing destination file operand after '/home/pranjal_salon1612/env/lib/python3.9/site-packages/apache_beam/examples/wordcount.py'
Try 'cp --help' for more information.
pranjal_salon1612@cloudshell:~ (projectmill1)$ cp /home/pranjal_salon1612/env/lib/python3.9/site-packages/apache_beam/examples/wordcount.py ~/wordcount.py
pranjal_salon1612@cloudshell:~ (projectmill1)$ ls
env outputs-00000-of-00001 SOPF4630U-tut3
index.html README-cloudshell.txt wordcount.py
pranjal_salon1612@cloudshell:~ (projectmill1)$ source ~/env/bin/activate
(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ cp ~/wordcount2.py --output outputs
cp: unrecognized option '--output'
Try 'cp --help' for more information.
(env) pranjal_salon1612@cloudshell:~ (projectmill1)$
```

```
l-WriteBundles_50))+ref_AppliedPTransform_Write2-Write-WriteImpl-Pair_51))+Write2/Write/WriteImpl/GroupByKey/Write)
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running ((Write2/Write/WriteImpl/GroupByKey/Read)+(ref_AppliedPTransform_Write2-Write-WriteImpl-Extract_53))+ref_PCollection_PCollection_35/Write)
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running ((ref_PCollection_PCollection_29/Read)+(ref_AppliedPTransform_Write2-Write-WriteImpl-Prefinalize_54))+ref_PCollection_PCollection_36/write)
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running (ref_PCollection_PCollection_29/Read)+(ref_AppliedPTransform_Write2-Write-WriteImpl-FinalizeWrite_55)
INFO:apache_beam.io.filebasedsink:Starting finalize_write threads with num_shards: 1 (skipped: 0), batches: 1, num_threads: 1
INFO:apache_beam.io.filebasedsink:Renamed 1 shards in 0.00 seconds.
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running (((ref_AppliedPTransform_Write-Write-WriteImpl-DoOnce-Impulse_28)+(ref_AppliedPTransform_Write-Write-WriteImpl-DoOnce-FlatMap-lambda-at-core-py-3228-29))+ref_AppliedPTransform_Write-Write-WriteImpl-DoOnce-Map-decode_31))+ref_AppliedPTransform_Write-Write-WriteImpl-InitializeWrite_32)+(ref_PCollection_PCollection_18/Write)+(ref_PCollection_PCollection_19/Write)
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running (((((GroupAndSum/Group/Read)+(GroupAndSum/Merge)+(GroupAndSum/ExtractOutputs))+(ref_AppliedPTransform_Format_22))+ref_AppliedPTransform_Write-Write-WriteImpl-WindowIntoFn_33)+(ref_AppliedPTransform_Write-Write-WriteImpl-WritesBundles_34))+ref_AppliedPTransform_Write-Write-WriteImpl-Pair_35)+(ref_PCollection_PCollection_25/Write)
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running ((ref_PCollection_PCollection_18/Read)+(ref_AppliedPTransform_Write-Write-WriteImpl-Prefinalize_38))+ref_PCollection_PCollection_25/Write)
INFO:apache_beam.runners.portability.fn_api_runner.fn_runner:Running (ref_PCollection_PCollection_18/Read)+(ref_AppliedPTransform_Write-Write-WriteImpl-FinalizeWrite_39)
INFO:apache_beam.io.filebasedsink:Starting finalize_write threads with num_shards: 1 (skipped: 0), batches: 1, num_threads: 1
INFO:apache_beam.io.filebasedsink:Renamed 1 shards in 0.00 seconds.
```

```
(env) pranjali_salon1612@cloudshell:~ (projectmill1)$ gcloud config list project
[core]
project = projectmill1

Your active configuration is: [cloudshell-15651]
(env) pranjali_salon1612@cloudshell:~ (projectmill1)$ gcloud config list project --format 'value(core.project)'

(env) pranjali_salon1612@cloudshell:~ (projectmill1)$
```

```
(env) goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ gcloud config list project --format "value(core.project)"
nomadic-bedrock-342622
(env) goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ PROJECT=$(_gcloud config list project --format 'value(core.project)')
(env) goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ echo $PROJECT
nomadic-bedrock-342622
(env) goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ BUCKET=$gs://$PROJECT-gs
echo $BUCKET
gs://nomadic-bedrock-342622-gs
(env) goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ python wordcount2.py \
--region northamerica-northeast2 \
--runner DataflowRunner \
--project $PROJECT \
--temp_location $BUCKET/tmp/ \
--output $BUCKET/result/output1 \
--output2 $BUCKET/result/output2 \
--experiment use unsupported python version
```

```
gs://dataflow-samples/shakespeare/macbeth.txt
gs://dataflow-samples/shakespeare/measureforemeasure.txt
gs://dataflow-samples/shakespeare/merchantofvenice.txt
gs://dataflow-samples/shakespeare/merrywivesofwindsof.txt
gs://dataflow-samples/shakespeare/midsummersnightsdream.txt
gs://dataflow-samples/shakespeare/muchadoaboutnothing.txt
gs://dataflow-samples/shakespeare/othello.txt
gs://dataflow-samples/shakespeare/particlesprinceoftyrr.txt
gs://dataflow-samples/shakespeare/rapeoflucrece.txt
gs://dataflow-samples/shakespeare/romeoandjuliet.txt
gs://dataflow-samples/shakespeare/sonnets.txt
gs://dataflow-samples/shakespeare/tamingofthebsrew.txt
gs://dataflow-samples/shakespeare/tempo.txt
gs://dataflow-samples/shakespeare/timesofethens.txt
gs://dataflow-samples/shakespeare/titusandronicus.txt
gs://dataflow-samples/shakespeare/troilusandcressida.txt
gs://dataflow-samples/shakespeare/twelfthnight.txt
gs://dataflow-samples/shakespeare/twogentlemenofverona.txt
gs://dataflow-samples/shakespeare/various.txt
gs://dataflow-samples/shakespeare/venusandadonis.txt
gs://dataflow-samples/shakespeare/winterstale.txt
goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs/result/output*-00000-of-00003
gs://nomadic-bedrock-342622-gs/result/output*-00000-of-00003
goergedao2@cloudshell:~ (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs/result/output*-00000-of-00003
```

```

aelse: 11
boreast: 1
countries: 1
digest: 1
becomes: 1
battles: 1
applied: 1
die like: 3
goergoed12@cloudshell:~ (nomadic-bedrock-342622)$ gutil cat gs://nomadic-bedrock-342622-pa/result/output2-00000-of-00003
b: 1231
q: 1353
t: 1214
u: 254
q1: 44
q2: 1006
q3: 821
t1: 492
u1: 914
p1: 718
t2: 3900
y1: 813
l: 1771
r: 156

```

Various Dataflow examples for Batch and stream processing for the mnist dataset for various source and destination types; text file, MySQL database, and Kafka topics.

<https://youtu.be/9ZDj9KDGtEs>

Screenshots:

```

(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ git clone https://github.com/goerge
aoud/SOFE4630U-tut4.git
Cloning into 'SOFE4630U-tut4'...
remote: Enumerating objects: 122, done.
remote: Counting objects: 100% (122/122), done.
remote: Compressing objects: 100% (85/85), done.
remote: Total 122 (delta 58), reused 96 (delta 37), pack-reused 0
Receiving objects: 100% (122/122), 14.85 MiB | 22.16 MiB/s, done.
Resolving deltas: 100% (58/58), done.

```

```

(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ ls
env      outputs-00000-of-00001  SOFE4630U-tut3  wordcount.py
index.html README-cloudshell.txt  SOFE4630U-tut4
(env) pranjal_salon1612@cloudshell:~ (projectmill1)$ cd SOFE4630U-tut4
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4 (projectmill1)$ ls
mnist  wordcount
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4 (projectmill1)$ cd mnist
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ ls
data      mysql-app.yaml  predictV1.py  predictV3.py  setup.py
mnist.txt  mysql-pvc.yaml  predictV2.py  predictV4.py
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ source ~/env/b
n/activate
sorflow-cpu==2.8.0
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ pip install te
Collecting tensorflow-cpu==2.8.0
  Downloading tensorflow_cpu-2.8.0-cp39-cp39-manylinux2010_x86_64.whl (190.6 MB)
    190.6/190.6 MB 5.0 MB/s eta 0:00:00
Requirement already satisfied: protobuf>=3.9.2 in /home/pranjal_salon1612/env/lib/python
3.9/site-packages (from tensorflow-cpu==2.8.0) (3.19.4)
Collecting google-pasta>=0.1.1
  Downloading google_pasta-0.2.0-py3-none-any.whl (57 kB)
    57.5/57.5 KB 9.7 MB/s eta 0:00:00
Collecting astunparse>=1.6.0
  Downloading astunparse-1.6.3-py2.py3-none-any.whl (12 kB)

```

```
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectml1)$ df
Filesystem           1K-blocks      Used Available Use% Mounted on
overlay                  62742040  48102400  14623256  77% /
tmpfs                      65536         0    65536   0% /dev
tmpfs                      8196724         0   8196724   0% /sys/fs/cgroup
/dev/sda1                  62742040  48102400  14623256  77% /root
/dev/disk/by-id/google-home-part1  5028480  1552972   3197032  33% /home
/dev/root                  2003760  1063960   939800  54% /usr/lib/modules
shm                         65536         0    65536   0% /dev/shm
tmpfs                      3278692       904   3277788   1% /google/host/var/ru
```

```
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectml1)$ ls -la
total 56
drwxr-xr-x 3 pranjal_salon1612 pranjal_salon1612 4096 Mar 29 04:27 .
drwxr-xr-x 5 pranjal_salon1612 pranjal_salon1612 4096 Mar 29 04:27 ..
drwxr-xr-x 3 pranjal_salon1612 pranjal_salon1612 4096 Mar 29 04:27 data
-rw-r--r-- 1 pranjal_salon1612 pranjal_salon1612 3940 Mar 29 04:27 mnist.txt
-rw-r--r-- 1 pranjal_salon1612 pranjal_salon1612 1053 Mar 29 04:27 mysql-app.yaml
-rw-r--r-- 1 pranjal_salon1612 pranjal_salon1612 166 Mar 29 04:27 mysql-pvc.yaml
-rw-r--r-- 1 pranjal_salon1612 pranjal_salon1612 3575 Mar 29 04:27 predictV1.py
-rw-r--r-- 1 pranjal_salon1612 pranjal_salon1612 4706 Mar 29 04:27 predictV2.py
```

```
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectml1)$ python ./predictV1.py \
>   --staging_location ./staging \
  --temp_location ./temp \
  --model ./data \
  --source text \
    --setup_file ./setup.py \
  --input ./data/images.txt \
  --staging_location ./staging \
  --temp_location ./temp \
  --model ./data \
  --source text \
    --setup_file ./setup.py \
  --input ./data/images.txt \
  --output ./predict
INFO:root:Missing pipeline option (runner). Executing pipeline using the default runner
irectRunner.
WARNING:root:Make sure that locally built Python SDK docker image has Python 3.9 inter-
preter.
INFO:root:Default Python SDK image for environment is apache/beam_python3.9_sdk:2.37.0
INFO:apache_beam.runners.portability.fn_api_runner.translations:=====
  connection annotate_downstream_side_inputs at 0x7f2b4e3311f0> =====
INFO:apache_beam.runners.portability.fn_api_runner.translations:=====
  connection annotate_downstream_side_inputs at 0x7f2b4e3311f0> =====
INFO:apache_beam.io.filebasedsink:Renamed 1 shards in 0.00 seconds.
```

```
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectml1)$ ls
data      mysql-app.yaml  predict-00000-of-00001  predictV2.py  predictV4.py
mnist.txt  mysql-pvc.yaml  predictV1.py          predictV3.py  setup.py
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectml1)$ 
```

```
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ PROJECT=$(gcloud config list project --format "value(core.project)")  
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ echo PROJECT  
PROJECT  
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ echo $PROJECT  
projectmill1  
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ PROJECT=$(gcloud config list project --format "value(core.project)")  
echo $PROJECT  
BUCKET=gs://$PROJECT-gs  
echo $BUCKET  
projectmill1  
gs://projectmill1-gs
```

```
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ gsutil cp data/export* $BUCKET/model/  
Copying file://data/export.data-00000-of-00001 [Content-Type=application/octet-stream]...  
.  
Copying file://data/export.index [Content-Type=application/octet-stream]...  
Copying file://data/export.meta [Content-Type=application/octet-stream]...  
- [3 files] [ 12.5 MiB/ 12.5 MiB]  
Operation completed over 3 objects/12.5 MiB.  
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ gsutil ls $PROJECT  
CommandException: "ls" command does not support "file://" URLs. Did you mean to use a g:// URL?  
(env) pranjal_salon1612@cloudshell:~/SOFE4630U-tut4/mnist (projectmill1)$ []
```

```
(env) goergedaoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadic-bedrock-342622)$ gsutil ls $BUCKET  
ServiceException: 401 Anonymous caller does not have storage.objects.list access to the Google Cloud Storage bucket.  
(env) goergedaoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs  
ServiceException: 401 Anonymous caller does not have storage.objects.list access to the Google Cloud Storage bucket.  
(env) goergedaoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs  
ServiceException: 401 Anonymous caller does not have storage.objects.list access to the Google Cloud Storage bucket.  
(env) goergedaoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs  
gs://nomadic-bedrock-342622-gs/result  
gs://nomadic-bedrock-342622-gs/tmp/  
(env) goergedaoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadic-bedrock-342622)$ gsutil cp ./data/export* $BUCKET/model/  
Copying file://./data/export.data-00000-of-00001 [Content-Type=application/octet-stream]...  
Copying file://./data/export.index [Content-Type=application/octet-stream]...  
Copying file://./data/export.meta [Content-Type=application/octet-stream]...  
- [3 files] [ 12.5 MiB/ 12.5 MiB]  
Operation completed over 3 objects/12.5 MiB.  
(env) goergedaoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadic-bedrock-342622)$ gsutil cp ./data/images.txt $BUCKET/input/  
Copying file://./data/images.txt [Content-Type=text/plain]...  
- [1 files] [ 45.2 MiB/ 45.2 MiB]  
Operation completed over 1 objects/45.2 MiB.
```

```

nomadic-bedrock-342622
gs://nomadic-bedrock-342622-gs
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ gsutil ls $BUCKET
ServiceException: 403 Anonymous caller does not have storage.objects.list access to the Google Cloud Storage bucket.
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs
ServiceException: 403 Anonymous caller does not have storage.objects.list access to the Google Cloud Storage bucket.
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs
ServiceException: 403 Anonymous caller does not have storage.objects.list access to the Google Cloud Storage bucket.
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ gsutil ls gs://nomadic-bedrock-342622-gs
gs://nomadic-bedrock-342622-gs/export/
gs://nomadic-bedrock-342622-gs/tmp/
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ gsutil cp ./data/export $BUCKET/model/
Copying file://./data/export.data-00000-of-00001 [Content-Type=application/octet-stream]...
copying file://./data/export.index [Content-Type=application/octet-stream]...
copying file://./data/export.meta [Content-Type=application/octet-stream]...
- [3 files] 12.5 MiB/ 12.5 MiB
Operation completed over 3 objects/12.5 MiB.
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ gsutil cp ./data/images.txt $BUCKET/input/
Copying file://./data/images.txt [Content-Type=text/plain]...
- [1 files] 45.2 MiB/ 45.2 MiB
Operation completed over 1 objects/45.2 MiB.
(env) george@aoud2:~/.SOFR4630U-tut4/minst (nomadic-bedrock-342622)$ python ./predictVI.py \
--model $BUCKET/model \
--project $PROJECT \
--staging_location $BUCKET/staging \
--temp_location $BUCKET/temp \
--model $BUCKET/model \
--source test \
--setup_file ./setup.py \
--input $BUCKET/input/images.txt \
--output $BUCKET/output/predict \

```

```
(env) pranjal_salonie12@cloudshell:~/SOFE4630U-tut4/mnist (projectmill)$ gcloud config set compute/zone northamerica-northeast2-a
Updated property [compute/zone].
(env) pranjal_salonie12@cloudshell:~/SOFE4630U-tut4/mnist (projectmill)$ gcloud container clusters create gk-cluster --num-nodes=3
Default change: VPC-native is the default mode during cluster creation for versions greater than 1.21.0-gke.1500. To create advanced routes based clusters, please pass the `--no-enable-ip-alias` flag.
Note: Your Pod address range ('--cluster-ip-range') can accommodate at most 1008 node(s).
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being configured..working..
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being configured..working..
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being configured..working..
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being configured..working..
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being configured..working..
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being health-checked..working..
Creating cluster gk-cluster in northamerica-northeast2-a.... Cluster is being health-checked (master is healthy) ..done.
Created [https://container.googleapis.com/v1/projects/projectmill/locations/northamerica-northeast2-a/clusters/gk-cluster].
To inspect the contents of your cluster, go to: https://console.cloud.google.com/kubernetes/workload_gcloud/northamerica-northeast2-a/gk-cluster?project=projectmill
kubernetes config entry generated for gk-cluster.
NAME: gk-cluster
LOCATION: northamerica-northeast2-a
MASTER_VERSION: 1.21.9-gke.1002
MASTER_IP: 34.130.252.229
MACHINE_TYPE: e2-medium
NODE_VERSION: 1.21.9-gke.1002
NUM_NODES: 3
STATUS: RUNNING
(env) pranjal_salonie12@cloudshell:~/SOFE4630U-tut4/mnist (projectmill)$

NAME_VERSION: 1.21.6-gke.1500
NUM_NODES: 3
STATUS: RUNNING
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ gcloud container clusters get-credentials gk-cluster
Fetching cluster endpoint and auth data.
kubernetes config entry generated for gk-cluster.
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl apply -f mysql-pvc.yaml
persistentvolumeclaim/mysql-pvc created
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl apply -f mysql-app.yaml
service/mysql created
deployment.apps/mysql created
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
kubernetes    ClusterIP   10.112.0.1    <none>          443/TCP       2m31s
mysql         LoadBalancer 10.112.10.46   <pending>       3306:31377/TCP  10s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
kubernetes    ClusterIP   10.112.0.1    <none>          443/TCP       2m43s
mysql         LoadBalancer 10.112.10.46   <pending>       3306:31377/TCP  19s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
kubernetes    ClusterIP   10.112.0.1    <none>          443/TCP       2m45s
mysql         LoadBalancer 10.112.10.46   <pending>       3306:31377/TCP  20s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
kubernetes    ClusterIP   10.112.0.1    <none>          443/TCP       2m56s
mysql         LoadBalancer 10.112.10.46   <pending>       3306:31377/TCP  20s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
kubernetes    ClusterIP   10.112.0.1    <none>          443/TCP       3m1s
mysql         LoadBalancer 10.112.10.46   <pending>       3306:31377/TCP  38s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
kubernetes    ClusterIP   10.112.0.1    <none>          443/TCP       3m13s
mysql         LoadBalancer 10.112.10.46   34.130.188.137  3306:31377/TCP  45s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get service mysql
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
mysql         LoadBalancer 10.112.10.46   34.130.188.137  3306:31377/TCP  67s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get service mysql
NAME          TYPE        CLUSTER-IP      EXTERNAL-IP      PORT(S)        AGE
mysql         LoadBalancer 10.112.10.46   34.130.188.137  3306:31377/TCP  103s
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ kubectl get services mysql -o jsonpath='{.status.loadBalancer.ingress[0].ip}'
(env) georgedasoud2@cloudshell:~/SOFE4630U-tut4/mnist (nomadio-bedrock-342422)$ K8SQQIP=$ (kubectl get services mysql -o jsonpath='{.status.loadBalancer.ingress[0].ip}')
echo $K8SQQIP
34.130.188.137
```

```
■ Command Prompt - py consumerMnist.py
Volume Serial Number is 3826-8548

Directory of F:\SOFE4630U\SOFE4630U-tut4\mnist\data

02/26/2022 11:43 AM <DIR> .
02/26/2022 11:43 AM <DIR> ..
02/26/2022 09:51 AM 1,057 consumerMnist.py
02/23/2022 12:21 PM 184 cred.json
08/30/2018 04:28 AM 13,098,536 export.data-00000-of-00001
08/30/2018 04:28 AM 405 export.index
08/30/2018 04:28 AM 20,761 export.meta
02/14/2022 03:21 PM 221,968 images.sql
01/19/2017 05:48 PM 47,415,498 images.txt
02/26/2022 02:20 PM <DIR> opt
02/26/2022 09:50 AM 830 producerMnist.py
    8 File(s) 60,259,239 Bytes
    3 0Dir(s) 130,786,680,832 bytes free

F:\SOFE4630U\SOFE4630U-tut4\mnist\data>py consumerMnist.py
C:\Python39\python.exe: can't open file 'F:\SOFE4630U\SOFE4630U-tut4\mnist\data\consumerMnist.py': [Errno 2] No such file or directory

F:\SOFE4630U\SOFE4630U-tut4\mnist\data>py consumerMnist.py

■ Select Command Prompt - py producerMnist.py
Image with key 1350 is sent
Image with key 1351 is sent
Image with key 1352 is sent
Image with key 1353 is sent
Image with key 1354 is sent
Image with key 1355 is sent
Image with key 1356 is sent
Image with key 1357 is sent
Image with key 1358 is sent
Image with key 1359 is sent
Image with key 1360 is sent
Image with key 1361 is sent
Image with key 1362 is sent
Image with key 1363 is sent
Image with key 1364 is sent
Image with key 1365 is sent
Image with key 1366 is sent
Image with key 1367 is sent
Image with key 1368 is sent
Image with key 1369 is sent
Image with key 1370 is sent
Image with key 1371 is sent
Image with key 1372 is sent
Image with key 1373 is sent
```

Q&A:

Google Cloud has another processing service called DataProc. Name another processing service that is usually used in the cloud environment (not necessarily GCP). Compare between it and both Dataflow and DataProc. Your comparison may include but is not limited to the major differences, advantages, disadvantages, and limitations.

Another processing service used in cloud environments is Google Cloud Dataprep.

	Dataflow	Dataproc	Dataprep
Main use	Executes Apache Beam pipelines within the GCP environment.	Highly scalable service that runs Apache Spark, Apache Flink, and various other open source tools and frameworks.	Visually explores, cleans and transforms structured and unstructured data for analytics and reporting.
Major differences	Follows batch and stream processing of data. Can only observe and stream processing data.	Dataproc is a Google Cloud product with Data Science/ML service for Spark and Hadoop.	Dataprep is UI-driven, scales on-demand and fully automated.
When to use?	Dataflow provides a clear separation between processing logic and the underlying execution engine.	Dataproc should be used if the processing has any dependencies to tools in the Hadoop ecosystem.	Dataprep is used only as a medium of processing data for further use, such as in BigQuery.
Advantages	Helps understand the functioning and the limits of a system.	Utilizes Spark and Hadoop; help create and manage clusters quickly.	Allows quick exploration of new datasets.
Disadvantages/Limitations	Limited job computation.	No choice of selecting a specific version of Hadoop/hive/spark stack.	Limited access to APIs.

Suggest a practical application using both stream and batch processing that can be applied to a given dataset. It's expected to use the dataset uploaded in the third milestone but you can use any other dataset. If you decide to use another dataset, It should maintain both variety and huge volume. Your report should include but not limited to:

- The application.
- Its impact.
- The used dataset (size, schema/structure).

- A graph showing the proposed pipeline(s).
- List of other tools (AI, clustering,...) needed to implement that application.

An example of batch processing is monthly bill generations. Batch processing is a technique for automating and processing multiple transactions as a single group. While batch processing can be carried out at any time, it's particularly suited to end-of-cycle processing, such as for processing a bank's reports at the end of a day or generating monthly or biweekly payrolls. Repeated jobs are done fast in batch systems without user interaction. ATM transactions use stream processing as data needs to be collected and processed in real time. BigQuery is a tool that can be implemented for these applications. BigQuery is a fully-managed, serverless data warehouse that enables scalable analysis over petabytes of data.