# EECS E6893 Big Data Analytics
# Intro to Big Data Analytics on GCP

Frank Ou Yang

# Agenda
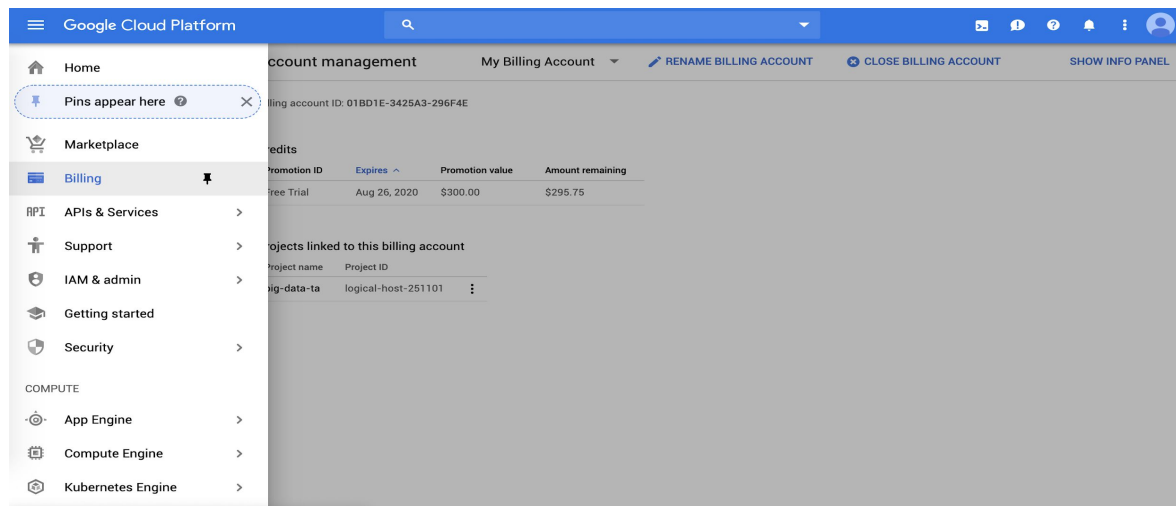
- GCP
- Cloud Storage
- BigQuery
- Dataproc

# Google Cloud Platform (GCP)

# GCP

- Cloud computing platform
  - Flexibility: on-demand and scale as you want
  - Efficiency: no need to maintain infra
- Services (relevant to this assignment)
  - Compute
    - Compute Engines: VMs
  - Big data products
    - BigQuery: Data warehouse for analytics
    - Dataproc: Hadoop and Spark
  - Storage & DBs
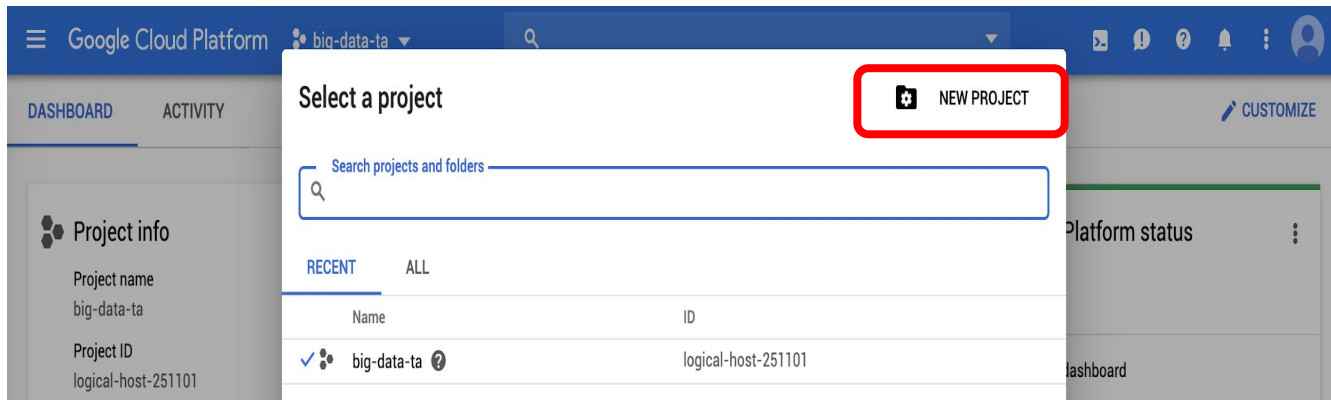    - Cloud Storage: Object storage system
  - More at https://cloud.google.com/products/

# GCP Setup

- Create a google account, you could use your Columbia account
- Apply for $300 credit for the first year: https://cloud.google.com/free/
- Go to Console dashboard -> Billing to check credit is there

# GCP: Create project

- Project: basic unit for creating, enabling, and using all GCP services
  - managing APIs, billing, permissions
  - adding and removing collaborators
- Visit console dashboard or cloud resource manager
- Click on "create project / new project" and complete the flow
- Ensure billing is pointing to the $300 credit

# GCP: Interaction

- Graphical UI / console: Useful to create VMs, set up clusters, provision resources, manage teams, etc
- Command line tools / Cloud SDK: Useful for interacting from local host and using the resources once provisioned. E.x. ssh into instances, submit jobs, copy files, etc
- Cloud shell: Same as command line, but web-based and pre-installed with SDK and tools

# GCP: console

Search for services here



Manage / Enable APIs

# GCP: Cloud SDK setup

- Install the SDK that is suitable for your local environment:
  https://cloud.google.com/sdk/docs/quickstarts

- Some testing after installation:
  - `gcloud auth list`
  - `gcloud components list`

```
[dyn-129-236-216-148:~ frank$ gcloud components list

Your current Cloud SDK version is: 259.0.0
The latest available version is: 261.0.0

┌─────────────────────────────────────────────────────────────────────────────────────────────┐
│                                          Components                                            │
├─────────────────┬───────────────────────────────────────────────┬──────────────────────┬──────────┤
│      Status     │                      Name                       │          ID          │   Size   │
├─────────────────┼───────────────────────────────────────────────┼──────────────────────┼──────────┤
│ Update Available │ BigQuery Command Line Tool                     │ bq                   │ < 1 MiB  │
│ Update Available │ Cloud SDK Core Libraries                       │ core                 │ 11.5 MiB │
│ Not Installed    │ App Engine Go Extensions                       │ app-engine-go        │ 56.4 MiB │
│ Not Installed    │ Cloud Bigtable Command Line Tool               │ cbt                  │ 7.3 MiB  │
│ Not Installed    │ Cloud Bigtable Emulator                        │ bigtable             │ 6.6 MiB  │
│ Not Installed    │ Cloud Datalab Command Line Tool                │ datalab              │ < 1 MiB  │
│ Not Installed    │ Cloud Datastore Emulator                       │ cloud-datastore-emulator │ 18.4 MiB │
│ Not Installed    │ Cloud Datastore Emulator (Legacy)              │ gcd-emulator         │ 38.1 MiB │
│ Not Installed    │ Cloud Firestore Emulator                       │ cloud-firestore-emulator │ 36.8 MiB │
│ Not Installed    │ Cloud Pub/Sub Emulator                         │ pubsub-emulator      │ 34.8 MiB │
│ Not Installed    │ Cloud SQL Proxy                                │ cloud_sql_proxy      │ 3.7 MiB  │
│ Not Installed    │ Emulator Reverse Proxy                         │ emulator-reverse-proxy │ 14.5 MiB │
│ Not Installed    │ Google Cloud Build Local Builder               │ cloud-build-local    │ 5.9 MiB  │
│ Not Installed    │ Google Container Registry's Docker credential helper │ docker-credential-gcr │ 1.8 MiB  │
│ Not Installed    │ gcloud Alpha Commands                          │ alpha                │ < 1 MiB  │
│ Not Installed    │ gcloud app Java Extensions                     │ app-engine-java      │ 85.9 MiB │
│ Not Installed    │ gcloud app PHP Extensions                      │ app-engine-php       │ 21.9 MiB │
│ Not Installed    │ gcloud app Python Extensions                   │ app-engine-python    │ 6.0 MiB  │
│ Not Installed    │ gcloud app Python Extensions (Extra Libraries) │ app-engine-python-extras │ 28.5 MiB │
│ Not Installed    │ kubectl                                        │ kubectl              │ < 1 MiB  │
│ Installed        │ Cloud Storage Command Line Tool                │ gsutil               │ 3.6 MiB  │
│ Installed        │ gcloud Beta Commands                           │ beta                 │ < 1 MiB  │
└─────────────────┴───────────────────────────────────────────────┴──────────────────────┴──────────┘

To install or remove components at your current SDK version [259.0.0], run:
  $ gcloud components install COMPONENT_ID
  $ gcloud components remove COMPONENT_ID

To update your SDK installation to the latest version [261.0.0], run:
  $ gcloud components update
```

# Cloud Storage

# Cloud Storage

- Online file storage system
- Graphical UI through console



- Command line tool: `gsutil`

# Cloud Storage - graphical UI

# Cloud Storage - graphical UI (cont')



Like *filepath* on GCP,
use this in your program

# Cloud Storage - `gsutil`

- Interact with Cloud Storage through command line
- Works similar to unix command line
- Useful commands:
    - Concatenate object content to stdout:
      ```
      gsutil cat [-h] url…
      ```
    - Copy file:
      ```
      gsutil cp [OPTION]... src_url dst_url
      ```
    - List files:
      ```
      gsutil ls [OPTION]... url…
      ```
- Explore more at https://cloud.google.com/storage/docs/gsutil

# BigQuery

# BigQuery

- Data warehouse for analytics
- SQL-like languages to interact with DB
- RESTful APIs for programmatic access
- Graphical UI

# Dataproc

# Dataproc

- On-demand, fully managed cloud service for running Apache Hadoop and Spark on GCP
- Cluster creation (using Cloud SDK):
  - Automatically creates VMs with Spark pre-installed
  - `gcloud dataproc clusters create <cluster-name>`
  - `gcloud beta dataproc clusters create <cluster-name>`
    `--optional-components=ANACONDA,JUPYTER --image-version=1.3`
    `--enable-component-gateway --bucket <bucket-name> --project`
    `<project-id> --single-node --metadata`
    `'PIP_PACKAGES=graphframes==0.6' --initialization-actions`
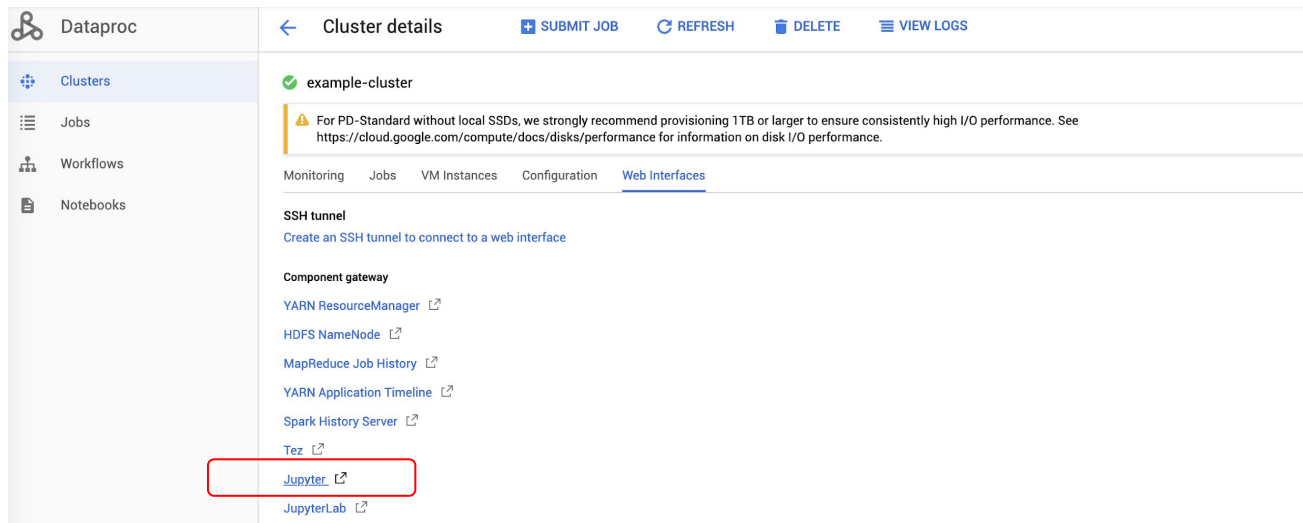    `gs://dataproc-initialization-actions/python/pip-install.sh`

Install Jupyter Notebook

Cloud Storage bucket: where your jupyter notebooks are saved

Works like pip install <your package>

# Dataproc - Spark execution / submit jobs

- Jupyter notebook:



- Cloud SDK:
  - ```
    gcloud dataproc jobs submit pyspark <your_program.py>
    --cluster=<cluster-name>
    ```
  - [View your jobs in console](#)

- program could be Cloud Storage URI or local path
- Data should be on Cloud storage

19

# Dataproc - Spark execution / submit jobs (cont')

- ## Spark shell
    - ssh into master node



    - pyspark

# HW0

1. Read documentations and tutorials
   a. Setup GCP and Cloud SDK
   b. Run Spark examples on Dataproc - Pi calculation and word count
   c. Familiar yourself with BigQuery
2. Two light programming questions
   a. BigQuery
   b. Spark program - Find top k most frequent words

**Remember to delete your dataproc clusters when you finish executions to save money.**