# SETTING UP APACHE SPARK CLUSTER ON GOOGLE CLOUD

## DR. ALVIN ANG

# CONTENTS

https://cloud.google.com/

## Step 2 of 2 Payment Information Verification

Your payment information helps us reduce fraud and abuse. You won't be charged unless you turn on automatic billing.

### Payments profile ⓘ

Choose the payments profile that will be associated with this account or transaction. A payments profile is shared and used across all Google products.

> **Dr. Alvin Ang**
> Individual profile for Ads
> Payments profile ID: 5375-7451-3343                          ⌄

### Payment method

**VISA** ▮▮▮▮▮▮▮▮         pls put your credit card here         ⌄

You'll be charged automatically on the 1st of each month. If your balance reaches your payment threshold before then, you'll be charged immediately. Learn more

### Tax information ⓘ  ✎

Tax status : Individual

The personal information you provide here will be added to your payments profile. It will be stored securely and treated in accordance with the Google Privacy Policy.

[ **START MY FREE TRIAL** ]
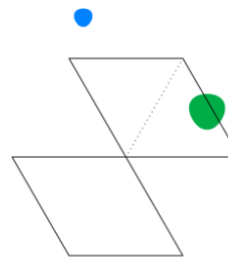
### Access to all Cloud Platform Products

Get everything you need to build and run your apps, websites and services, including Firebase and the Google Maps API.

### $300 credit for free

Put Google Cloud to work with $300 in credit to spend over the next 90 days.

### No autocharge after free trial ends

We ask you for your credit card to make sure you are not a robot. You won't be charged unless you manually upgrade to a paid account.

---

### Google Cloud Platform

## Welcome Dr Alvin Ang!

✓ What brought you to Google Cloud?

③ **What are you interested in doing with Google Cloud?**

| Websites | Mobile apps | ✓ Storage / backup |
| ✓ Data analytics | ✓ Artificial intelligence / machine learning |
| Game development | Containerization | ✓ Data management |
| Virtual machines (VMs) | Google Maps |
| Other APIs (e.g. Text-to-Speech, Speech-to-Text, Vision) |
| Google Photos or Google Workspace | Other | I'm not sure yet |

**NEXT**

④ What best describes your role?

CLOSE     **DONE**

**Google** Cloud Platform

## Get started with an interactive tutorial

### Try Compute Engine
🕐 4 minutes

Learn how to create a highly configurable Linux VM instance for running workloads on Compute Engine.

✅ Custom machine types to optimize vCPU and memory while balancing cost

✅ Preemptible machines to reduce computing costs

✅ Rightsizing recommendations to optimize resource utilization

Monthly estimate: $25.46 credits ⓘ

### Try Cloud Storage
🕐 5 minutes

Learn how to work with object storage for all-sized needs. Store any amount of data. Retrieve as often as you'd like.

✅ Store files and objects remotely and retrieve from anywhere

✅ Worldwide access and storage locations

✅ High availability and durability

Monthly estimate: $5.10 credits ⓘ

### Try Cloud SQL
🕐 10 minutes

Learn how to get started with a fully managed relational database service for MySQL, PostgreSQL, and SQL Server.

✅ Fully managed database set up in minutes

✅ Easily migrate from existing databases

✅ Integrate with any application with full database compatibility

START TUTORIAL

SKIP FOR NOW

**Google** Cloud Platform

# Welcome Dr Alvin Ang!

Your free trial includes $300 in credit to spend over the next 90 days. To help us serve you better, please answer 4 questions.

✓ What best describes your organization or needs?

✓ What brought you to Google Cloud?

✓ What are you interested in doing with Google Cloud?

4 What best describes your role?

Please select *
Data Scientist / Data Engineer ▼

CLOSE     DONE

# A. SETUP CLUSTER

## 1. SETUP NAME / LOCATION / CLUSTER TYPE



## 2. IGNORE AUTOSCALING AND FLEXIBILITY MODE

3. CHOOSE VERSIONING

← Create a Dataproc cluster on Compute Engine

- **Set up cluster**
  Begin by providing basic information.

- **Configure nodes** (optional)
  Change node compute and storage capabilities.

- **Customize cluster** (optional)
  Add cluster properties, features, and actions.

- **Manage security** (optional)
  Change access, encryption, and security settings.

**CREATE**    CANCEL

Compatible File System (HCFS) shuffle. Learn more

ℹ An autoscaling policy must be selected to configure EFM.

**Versioning**

Use a custom image to load pre-installed packages. Learn more

**Image Type and Version**
1.5-debian10

**Release Date**
First released on 3/25/20?

**CHANGE**

**we select this version but it doesn't really matter**

**Components**

**Component Gateway**

☑ Enable component gateway
Provides access to the web interfaces of default and selected optional components on the

---

?1 days remaining - with a full account, you'll get unlimited ac

My First Project ▾    Search (/) for resources,

← Create a Dataproc cluster on Con

- **Set up cluster**
  Begin by providing basic information.

- **Configure nodes** (optional)
  Change node compute and storage capabilities.

- **Customize cluster** (optional)
  Add cluster properties, features, and actions.

- **Manage security** (optional)
  Change access, encryption, and security settings.

**CREATE**    CANCEL

**Choose Image Version**

| STANDARD DATAPROC IMAGE | CUSTOM IMAGE |

Cloud Dataproc uses versioned images to bundle the operating system, big data compo
Platform connectors into one package that is deployed on your cluster. Learn more

○ 2.1 (Debian 11, Hadoop 3.3, Spark 3.3)
First released on 1/22/2021.

○ 2.1 (RockyLinux 8, Hadoop 3.3, Spark 3.3)
First released on 2/18/2022.

○ 2.1 (Ubuntu 20.04 LTS, Hadoop 3.3, Spark 3.3)
First released on 1/22/2021.

○ 2.0 (Debian 10, Hadoop 3.2, Spark 3.1)
First released on 1/22/2021.

○ 2.0 (RockyLinux 8, Hadoop 3.2, Spark 3.1)
First released on 2/18/2022.

○ 2.0 (Ubuntu 18.04 LTS, Hadoop 3.2, Spark 3.1)
First released on 1/22/2021.

◉ 1.5 (Debian 10, Hadoop 2.10, Spark 2.4)
First released on 3/25/2020.

○ 1.5 (RockyLinux 8, Hadoop 2.10, Spark 2.4)
First released on 2/18/2022.

---

**14 | PAGE**

4. SELECT COMPONENTS

## B. CONFIGURE NODES

### 1. CONFIGURE MASTER NODE



### 2. CONFIGURE SLAVE NODES

3. OBSERVE THE TOTAL YARN USAGE



**Create a Dataproc cluster on Compute Engine**

- Set up cluster
  Begin by providing basic information.
- **Configure nodes** (optional)
  Change node compute and storage capabilities.
- Customize cluster (optional)
  Add cluster properties, features, and actions.
- Manage security (optional)
  Change access, encryption, and security settings.

**CREATE**   CANCEL

EQUIVALENT COMMAND LINE ▾

**Secondary worker nodes**

Each contains a YARN NodeManager. HDFS does not run on secondary worker nodes. Secondary worker VMs are preemptible by default. Spot and preemptible VMs costs less, but can be terminated at any time due to system demands. Learn more

**Sole-tenancy**

Enable to create this cluster on sole-tenant nodes. This grants exclusive access to a physical Compute Engine server that is dedicated to hosting only your project's VMs. If you are creating a cluster with an autoscaling policy, it is recommended that the node group you select also uses an autoscaling policy. Learn more

Enable

**Shielded VM**

Turn on all settings for the most secure configuration. Learn more
☐ Turn on Secure Boot ❓
☐ Turn on vTPM ❓
☐ Turn on Integrity Monitoring ❓

**we ignore these 3**

**Total YARN usage**

The number of worker nodes times the amount of memory on each node times the fraction given to YARN (0.8)  ✕

YARN cores ❓        YARN memory ❓
3                    9 GB

**3 x 3.75GB x 0.8**

---

**Create a Dataproc cluster on Compute Engine**

- Set up cluster
  Begin by providing basic information.
- **Configure nodes** (optional)
  Change node compute and storage capabilities.
- Customize cluster (optional)
  Add cluster properties, features, and actions.
- Manage security (optional)
  Change access, encryption, and security settings.

**CREATE**   CANCEL

an autoscaling policy. Learn more

Enable

**Shielded VM**

Turn on all settings for the most secure configuration. Learn more
☐ Turn on Secure Boot ❓
☐ Turn on vTPM ❓
☐ Turn on Integrity Monitoring ❓

The number of worker nodes times the number of vCPUs per node  ✕

YARN cores ❓        YARN memory ❓
3                    9 GB

**3 x 1vCPU (virtual CPU... means virtual computer)**

## C. CUSTOMIZE CLUSTER

### 1. SCHEDULED DELETION OF CLUSTER



schedule a deletino of your cluster after 1 hour if idle to save you money in case overcharge

### 2. EDIT BUCKET NOW…

*a)  Create a New Bucket*

Select bucket          let's click here to try creating
                        a new bucket...

‹  Buckets ▾                    🗑  🔍

📦  dataproc-staging-us-central1-699946211312-8lvofagg    ›

📦  dataproc-temp-us-central1-699946211312-bdfhm5ku      ›

u realize that 2 buckets have already
been DEFAULT created for you
and will be used if you dun change
anything....

SELECT     CANCEL

*b)  Name Your New Bucket*

Create a bucket

✓  **Name your bucket**

Pick a **globally unique**, permanent name. Naming guidelines ⧉

alvin-yarn-cluster-bucket        rename

Tip: Don't include any sensitive information

⌄ LABELS (OPTIONAL)

CONTINUE

•  **Choose where to store your data**

This choice defines the geographic placement of your data and affects cost,
performance, and availability. Cannot be changed later. Learn more ⧉

**Location type**

◉  Multi-region
    Highest availability across largest area

us (multiple regions in United States)        ▾

*c) Where to Store Your Data?*



# Choose where to store your data

This choice defines the geographic placement of your data and affects cost, performance, and availability. Cannot be changed later. Learn more

**Location type**    we don't want multi-region

○ Multi-region
   Highest availability across largest area

○ Dual-region
   High availability and low latency across 2 regions

● Region
   Lowest latency within a single region

| us-central1 (Iowa) |

we want to be in the same data center as the cluster we created earlier... same location

CONTINUE

*d) Leave the Rest of the Default Settings*

- **Choose how to protect object data** IGNORE

Your data is always protected with Cloud Storage but you can also choose from these additional data protection options to prevent data loss. Note that object versioning and retention policies cannot be used together.

**Protection tools**

( ) None

( ) Object versioning (best for data recovery)
For restoring deleted or overwritten objects. To minimize the cost of storing versions, we recommend limiting the number of noncurrent versions per object and scheduling them to expire after a number of days. Learn more ⬚

( ) Retention policy (best for compliance)
For preventing the deletion or modification of the bucket's objects for a specified minimum duration of time after being uploaded. Learn more ⬚

⌄ DATA ENCRYPTION

CREATE THE BUCKET NOW

**CREATE**   CANCEL

## Public access will be prevented

This bucket is set to prevent exposure of its data on the public internet.

Keep this setting enabled unless you have a use case that requires public access (such as static website hosting). You can change it now or later. Learn more ⬚

☑ Enforce public access prevention on this bucket

☐ Don't show this message again

CANCEL   **CONFIRM**

# Select bucket



<    Buckets ▾

🗑    🔍

| | |
|---|---|
| 🛒   alvin-yarn-cluster-bucket | > |
| 🛒   dataproc-staging-us-central1-699946211312-8lvofagg | > |
| 🛒   dataproc-temp-us-central1-699946211312-bdfhm5ku | > |

**select your newly created bucket**

**SELECT**    CANCEL

← Create a Dataproc cluster on Compute Engine

- **Set up cluster**
  Begin by providing basic information.

- **Configure nodes** (optional)
  Change node compute and storage capabilities.

- **Customize cluster** (optional)
  Add cluster properties, features, and actions.

- **Manage security** (optional)
  Change access, encryption, and security settings.

ignore this completely

CREATE    CANCEL

### Project access

☐ Enables the cloud-platform scope for this cluster   Learn more

**Encryption**

◉ Google-managed encryption key
   No configuration required

◯ Customer-managed encryption key (CMEK)
   Manage via Google Cloud Key Management Service

### Personal Cluster Authentication

Enable Dataproc Personal Cluster Authentication to allow interactive workloads on the cluster to securely run as your end user identity. Learn more

◯ Enable

### Secure Multi Tenancy

## A. NO ANACONDA ISSUE....

When u try to create a new cluster....

**B. CREATE YOUR CLUSTER!**

← Create a Dataproc cluster on Compute Engine

- **Set up cluster**
  Begin by providing basic information.

- **Configure nodes** (optional)
  Change node compute and storage capabilities.

- **Customize cluster** (optional)
  Add cluster properties, features, and actions.

- **Manage security** (optional)
  Change access, encryption, and security settings.

+ ADD METADATA

## Scheduled deletion

Use Scheduled Deletion to help avoid incurring Google Cloud charges for an inactive cluster.
Learn more

☐ Delete on a fixed time schedule

☑ Delete after a cluster idle time period without submitted jobs

Timeout *
15                                    Minutes ▾

The cluster will be deleted when idle for more than 15 minutes

## Cloud Storage staging bucket

Storage staging bucket
🗑 alvin-yarn-cluster-bucket                         BROWSE

Cloud Storage staging bucket to be used for storing cluster job dependencies, job driver output, and cluster config files.

NOW CREATE YOUR CLUSTER~!!!

**CREATE**    CANCEL

**C. IT TAKES SUPER LONG….**

Free trial status: $394.09 credit and 90 days remaining - with a full account, you'll get unlimited access to all of Google Cloud Platform.    DISMISS    ACTIVATE

≡  Google Cloud    My First Project ▾    Search (/) for resources, docs, products, and more    🔍 Search

⚙ Dataproc    Clusters    ⊞ CREATE CLUSTER    · REFRESH    ▶ START    ■ STOP    🗑 DELETE    REGIONS ▾    + 5 RECOMMENDED ALERTS    H

Jobs on Clusters
  ⊕ Clusters
  ☰ Jobs
  ⊓ Workflows
  ⊪ Autoscaling policies

Serverless
  ☰ Batches

Metastore Services

  ▤ Release Notes

⫷

≡ Filter   Search clusters, press Enter                    ❓ ▥

☐  Name ↑    Status    Region    Zone    Total worker nodes
☐  yarn-cluster    ⟳ Provisioning    us-central1    us-central1-c    3

this thing takes super long and takes forever to create your cluster….

Request to create cluster yarn-cluster submitted    ✕

No clusters selected

PERMISSIONS    LABELS

ⓘ  Please select at least one resource.

## D. MEANWHILE YOU MAY CHECK YOUR DASHBOARD…



yes is now up!

click this



u created 1 master and 3 workers

let's take a look at the spark shell (which looks like the cmd prompt in laptop)

## E. SPARK SHELL

## 1. SPARK HISTORY SERVER





we need to stop the Spark Shell to see this....

However, we are not likely to use the Spark Shell... we will most probably use the Jupyter Notebook...

## F.   JUPYTER NOTEBOOK

the notebook is not like a shell which immediately establishes the driver and executors

it is just a web interface which is not yet connected to the cluster

we must run some spark commands to initiate the connection

1. YARN RESOURCE MANAGER



this was our previous spark shell which we closed already

this is our current Jupyer Notebook running actively

click here

2.  CHECK PREVIOUS JOBS



3.  CHECK EXECUTORS

https://www.alvinang.sg/s/wordcount.py

https://www.alvinang.sg/s/pi.py

### A.   WE WILL TRY SUBMITTING AN APPLICATION TO THE CLUSTER NOW

upload your wordcount.py
to the cloud

I was trying out wordcount.py… but it seems to get stuck forever… so I decided to try using pi.py…

## C. RUN THE CODE (SPARK-SUBMIT)

**D.   CHECK OUT THE SPARK HISTORY…**



**E.   FAILED…..SIGH….**



it just shows me back the old 'app ID' which i did earlier on.. its supposed to show me
all the new 'app IDs' like 0005 etc...

perhaps i need to delete and restart my entire cluster and run everything again
and hopefully it will work smoothly...

## A. TO PREVENT OVERCHARGING



This process takes super super long…..

**B. DELETE YOUR STORAGE BUCKET TOO**



let's go find our storage buckets



delete them all away!

Dr. Alvin Ang earned his Ph.D., Masters and Bachelor degrees from NTU, Singapore. He is a scientist, entrepreneur, as well as a personal/business advisor. More about him at www.AlvinAng.sg.