

## Distributed Computing with Dataproc Assignment

As part of my Distributed Computing with Dataproc assignment, I completed two hands-on labs to gain practical experience with Google Cloud services. I worked on the following labs:

- Introduction to Cloud Dataproc: Hadoop and Spark on Google Cloud
- Dataproc: Qwik Start Console.

These labs allowed me to dive into Google Cloud's Dataproc platform and learn how to set up and manage clusters while performing distributed computing tasks. I investigated the platform's capabilities and learned how its services and tools perform in real-world scenarios.

### *Introduction to Cloud Dataproc: Hadoop and Spark on Google Cloud*

The screenshot displays the Google Cloud Qwiklabs interface. At the top, the lab title 'Introduction to Cloud Dataproc: Hadoop and Spark on Google Cloud' is shown with a '20 pts' rating. The main content area features a 'Congratulations!' message and a 'Next steps / Learn more' section with a link to 'Dataproc Documentation'. On the right, a sidebar lists lab instructions and tasks, including 'Setup and requirements', 'Task 1. Create a Cloud Dataproc cluster', 'Task 2. Submit a Spark job to your cluster', 'Task 3. Shut down your cluster', and 'Task 4. Test your understanding'. Below the lab completion message, the user's profile is shown, including the name 'Vaishnavi Pawar', 'Member since 2024', and '470 points'. A notification states 'Your profile is not public and accessible. Make profile public'. The bottom section shows a table of activities, including the completed lab and 'Getting Started with Cloud Shell and gcloud'.

Activity	Type	Date started	Date finished	Score	Passed
<a href="#">Introduction to Cloud Dataproc: Hadoop and Spark on Google Cloud</a>	Lab	18 minutes ago	0 minutes ago	Assessment: 100%	✓
<a href="#">Getting Started with Cloud Shell and gcloud</a>	Lab	12 days ago	12 days ago	Assessment: 100%	✓

Dataproc: Qwik Start - Console

← Dataproc: Qwik Start - Console

Share

Heart

Help

Global

★ 150 pts

🕒 13th

End Lab

00:02:26

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more.](#)

Open Google Cloud console

Username

student-01-3cd9fe746i

📄

Password

xkfh7yZ7pDg7

📄

Project ID

qwiklabs-gcp-08-9712

📄

Student Resources

Congratulations!

Now you know how to use the Google Cloud console to create and update a Dataproc cluster and then submit a job in that cluster.

Next steps / Learn more

This lab is also part of a series of labs called Qwik Starts. These labs are designed to give you a little taste of the many features available with Google Cloud. Search for "Qwik Starts" in the [lab catalog](#) to find the next lab you'd like to take!

Lab instructions and tasks

GSP103

100/100

Overview

Setup and requirements

Task 1. Create a cluster

Task 2. Submit a job

Task 3. View the job output

Task 4. Update a cluster to modify the number of workers

Task 5. Test your understanding

Congratulations!

Vaishnavi Pawar

Member since 2024

600 points

Your profile is not public and accessible. [Make profile public](#)

Paths

Activities

Leaderboard

Badges

Course

Lab

Quiz

Game

In progress

Finished

Activity	Type	Date started	Date finished	Score	Passed
<a href="#">Dataproc: Qwik Start - Console</a>	Lab	27 minutes ago	0 minutes ago	Assessment: 100%	✓
<a href="#">Introduction to Cloud Dataproc: Hadoop and Spark on Google Cloud</a>	Lab	52 minutes ago	33 minutes ago	Assessment: 100%	✓

## **Report:**

### **My Understanding of Cloud Dataproc Capabilities**

According to what I've learned, Cloud Dataproc is essentially a tool for simplifying Hadoop and Spark cluster management. Dataproc takes care of everything for you, so you don't have to worry about manually configuring each node or balancing workloads. It's like having an automated assistant who not only configures your cluster but also ensures that it's compatible with other GCP services such as Google Cloud Storage and BigQuery.

The coolest part for me is how adaptable Dataproc is; you can scale up or down depending on your job requirements. This means you're not committing resources (or costs) that you don't need. It's also very convenient that everything is integrated into the GCP ecosystem, making it easier to manage data flows without switching between platforms.

### **What I Learned in Each Lab**

#### Introduction to Cloud Data Dataproc: Hadoop and Spark on Google Cloud.

This lab gave me firsthand experience with creating and configuring a Dataproc cluster. What caught my interest the most was how easy it was to run jobs using both Hadoop and Spark without getting bogged down in configuration details. I saw firsthand how Spark jobs can be distributed across multiple nodes, which was eye-opening when considering how large datasets are handled.

Main Takeaways:

- I learned how to configure a Dataproc cluster and select the appropriate machine types and worker nodes for each task.
- I ran a sample Spark job, and it was fascinating to see how the cluster handled it.
- The integration with Google Cloud Storage made it simple to pull in data, process it, and save the results.

#### Dataproc: Qwik Start - Console.

This lab focused more on using the GCP console, which I appreciated because I didn't have to use the command line to complete tasks. Navigating the console was simple, and I was able to quickly launch a Dataproc cluster and manage jobs directly from the website. I spent a lot of time looking through the logs and tracking how jobs were progressing, which helped me understand how Dataproc manages jobs behind the scenes.

#### Main Takeaways:

- The GCP console is extremely user-friendly for managing Dataproc clusters.
- I learned how to launch jobs directly from the interface and monitor their status in real time.
- The lab demonstrated how easy and manageable logs and error tracking are with GCP's built-in tools.
- I also studied into what happens when you change the number of worker nodes in the cluster. Adding more nodes increased workload distribution efficiency, which improved job performance, particularly when dealing with large datasets.

### **How Changing the Number of Worker Nodes Affects Performance**

My experience in these labs taught me that the number of worker nodes has a significant impact on performance. More worker nodes result in more parallel processing, which speeds up the job. It's pretty simple: if you have a large dataset or a complex Spark job, adding more worker nodes will speed up and smooth out the workload by distributing it across more resources.

However, there is a balance. While adding nodes boosts performance, it also raises costs, so it's critical to find the sweet spot. I discovered that Dataproc's autoscaling feature is a lifesaver here; it automatically adjusts the number of nodes based on the job's requirements, allowing you to avoid paying for resources that aren't required. In short, more nodes equal faster processing, but you must consider cost efficiency.