

# Open source serverless spark data pipelines





## Ajay Yadav

Data Platform Practice Lead  
Google Cloud

[linkedin.com/in/ajayydv](https://www.linkedin.com/in/ajayydv)



## Shashank Agarwal

Staff Strategic Cloud Eng.  
Google PSO

[linkedin.com/in/shashank181](https://www.linkedin.com/in/shashank181)

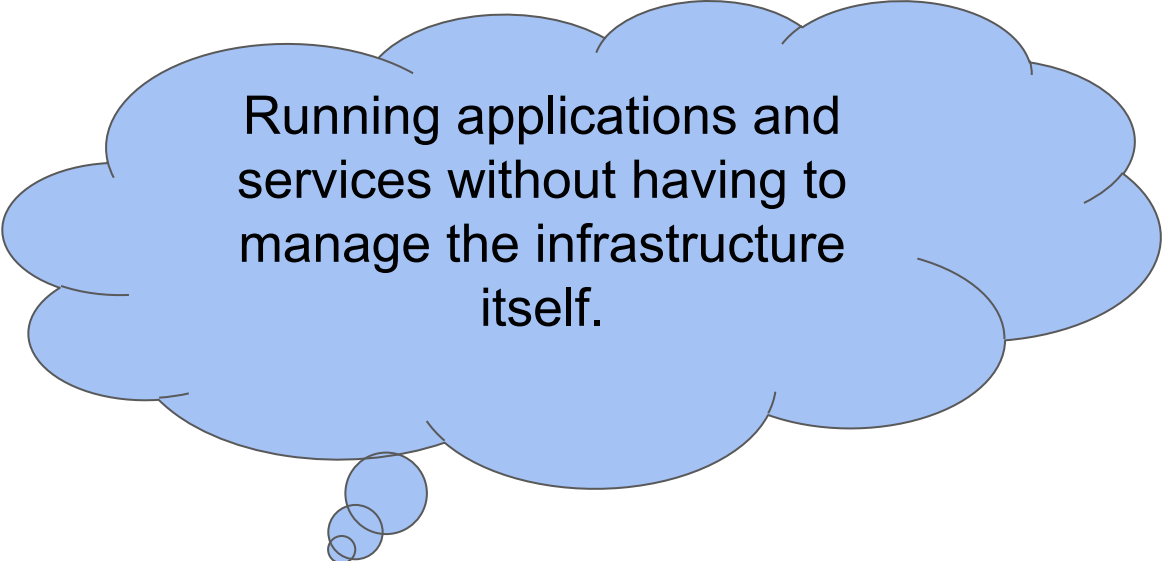


# Agenda



- 01 Serverless computing
- 02 What is Serverless Spark ?
- 03 Serverless Spark benefits
- 04 Open source spark templates
- 05 Demo
- 06 Questions

What is serverless computing?



Running applications and services without having to manage the infrastructure itself.

**“ Serverless architectures enable developers to focus on what they should be doing — writing code and optimizing application design — making way for business agility”**

Gartner

# Serverless on the rise



Developer Speed  
& Efficiency



Security & Flexibility



Data Governance  
& Compliance



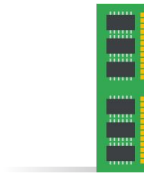
Cost Savings

What is Serverless Spark ?

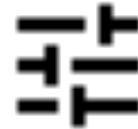
Offering by GCP Dataproc which lets you run **Spark workloads in Serverless mode.**



Autoscaling



Auto provisioning  
and cleanup



Auto tuning

# Benefits of Serverless Spark?



- Accelerate spark development and deployment
- Use ready to use Dataproc templates to bootstrap your project.
- Integration with Vertex AI, Jupyter notebooks and Dataplex.
- Submit jobs in Pyspark, Spark SQL, Spark R or Spark Java/Scala



- Reduced Operational overhead
- Simplified Governance & Security
- Accelerate Data lake modernization & migration efforts.
- Utilize existing skill set with a cloud native, serverless environment.

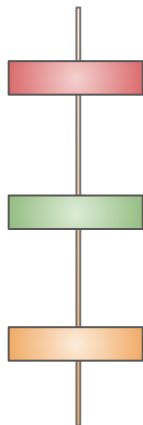


# Serverless Spark Templates

Ready to use,  
open sourced  
([github](#)),  
customizable  
templates and  
notebooks.

The screenshot shows a web browser window displaying the Google Cloud documentation page for Dataproc Serverless Spark templates. The browser's address bar shows the URL: `cloud.google.com/dataproc-serverless/docs/overview#for_spark_com...`. The page header includes the Google Cloud logo and navigation links for Overview, Solutions, Products, Search, and a language dropdown set to English. Below the header, the breadcrumb trail reads: Dataproc Serverless > Overview > Guides > Reference > Support > Resources. A 'Contact Us' button is visible in the top right. The main content area features a left-hand navigation menu with sections: Discover (with a link to 'What is Dataproc Serverless?'), Get started (with a link to 'Run an Apache Spark batch workload'), Ready-to-use templates (with a link to 'Config-driven Spark templates'), and How-to guides (with a link to 'BigQuery connector with Spark'). The main article title is 'What is Dataproc Serverless?' with a 'Send feedback' button and a 'Was this helpful?' feedback widget. Below the title, the 'On this page' section lists two links: 'Dataproc Serverless for Spark compared to Dataproc on Compute Engine' and 'Dataproc Serverless for Spark workload capabilities'. The main text of the article begins with: 'Dataproc Serverless lets you run Spark batch workloads without requiring you to provision and manage your own cluster. Specify workload parameters, and then submit the workload to the Dataproc Serverless'.

# How to use Serverless spark templates?



Step 1: Clone github repo in cloud shell (or preferred location).

<https://github.com/GoogleCloudPlatform/dataproc-templates>

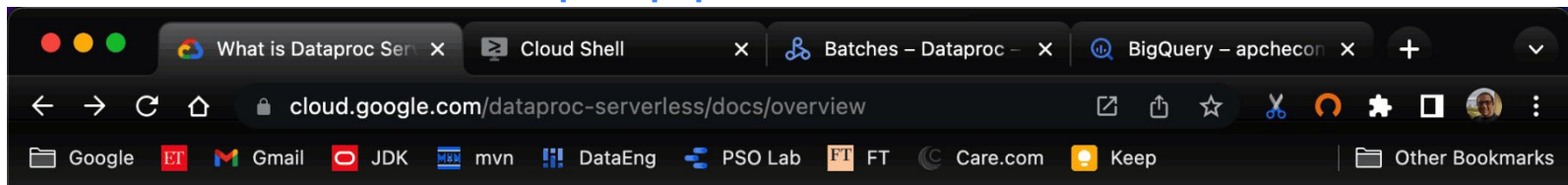
Step 2: Authenticate with your GCP identity or use service account.

```
gcloud auth application-default login
```

Step 3: Launch respective template from the shell.

It will automatically build, deploy and run the spark code.

# Example pipeline execution



Google Cloud Overview Solutions Products Pric > English Console

Dataproc Serverless Overview Guides Reference Support Resources [Contact Us](#)

Filter

Discover

- [What is Dataproc Serverless?](#)

Get started

- Run an Apache Spark batch workload

Ready-to-use templates

- Config-driven Spark templates


## What is Dataproc Serverless?

[Send feedback](#)

On this page

- Dataproc Serverless for Spark compared to Dataproc on Compute Engine
- Dataproc Serverless for Spark workload capabilities

Dataproc Serverless lets you run Spark batch workloads without requiring you to provision and manage your own cluster. Specify workload parameters, and then submit the workload to the Dataproc Serverless service. The service will run the workload on a managed compute infrastructure, autoscaling resources as needed. Dataproc Serverless charges apply only to the time when the workload is executing.



Cloud

Demo

Java	Python	Notebooks
HiveToBigQuery	BigQueryToGCS ( <a href="#">blogpost</a> )	<a href="#">Hive to Bigquery Notebook</a>
HiveToGCS	GCSToBigTable	<a href="#">MySQL to Cloud Spanner Notebook</a>
GCSToBigQuery	GCSToJDBC	SQL Server to Postgresql Notebook (on roadmap)
GCSToGCS	GCSToMongo ( <a href="#">blogpost</a> )	SQL Server to BigQuery Notebook (on roadmap)
GCSToSpanner ( <a href="#">blogpost</a> )	GCSToGCS	Oracle to Spanner Notebook (on roadmap)
HBaseToGCS	HiveToBigQuery ( <a href="#">blogpost</a> )	Oracle to BigQuery Notebook (on roadmap)
SpannerToGCS ( <a href="#">blogpost</a> )	HiveToGCS ( <a href="#">blogpost</a> )	
JDBCToBigQuery	HbaseToGCS	
JDBCToGCS ( <a href="#">blogpost</a> )	MongoToGCS ( <a href="#">blogpost</a> )	
PubSubToGCS ( <a href="#">blogpost</a> )	SnowflakeToGCS	
GCSToJDBC ( <a href="#">blogpost</a> )	JDBCToGCS	
SnowflakeToGCS	JDBCToBigQuery	
KafkaToBQ ( <a href="#">blogpost</a> )	RedshiftToGCS	

# Community engagement and learning resources

Running Spark Jobs in a Serverless Environment ([Link](#))

How to use Kubernetes with Spark ([Link](#))

How to use Notebooks with Spark ([Link](#))

How to use Iceberg and Delta with Serverless Spark ([Link](#))



# Contact Us

Feedback, ideas, thoughts [feedback-form](#)

Questions, issues, and comments - [dataproc-templates-support-external@googlegroups.com](mailto:dataproc-templates-support-external@googlegroups.com)



**Thank you.**