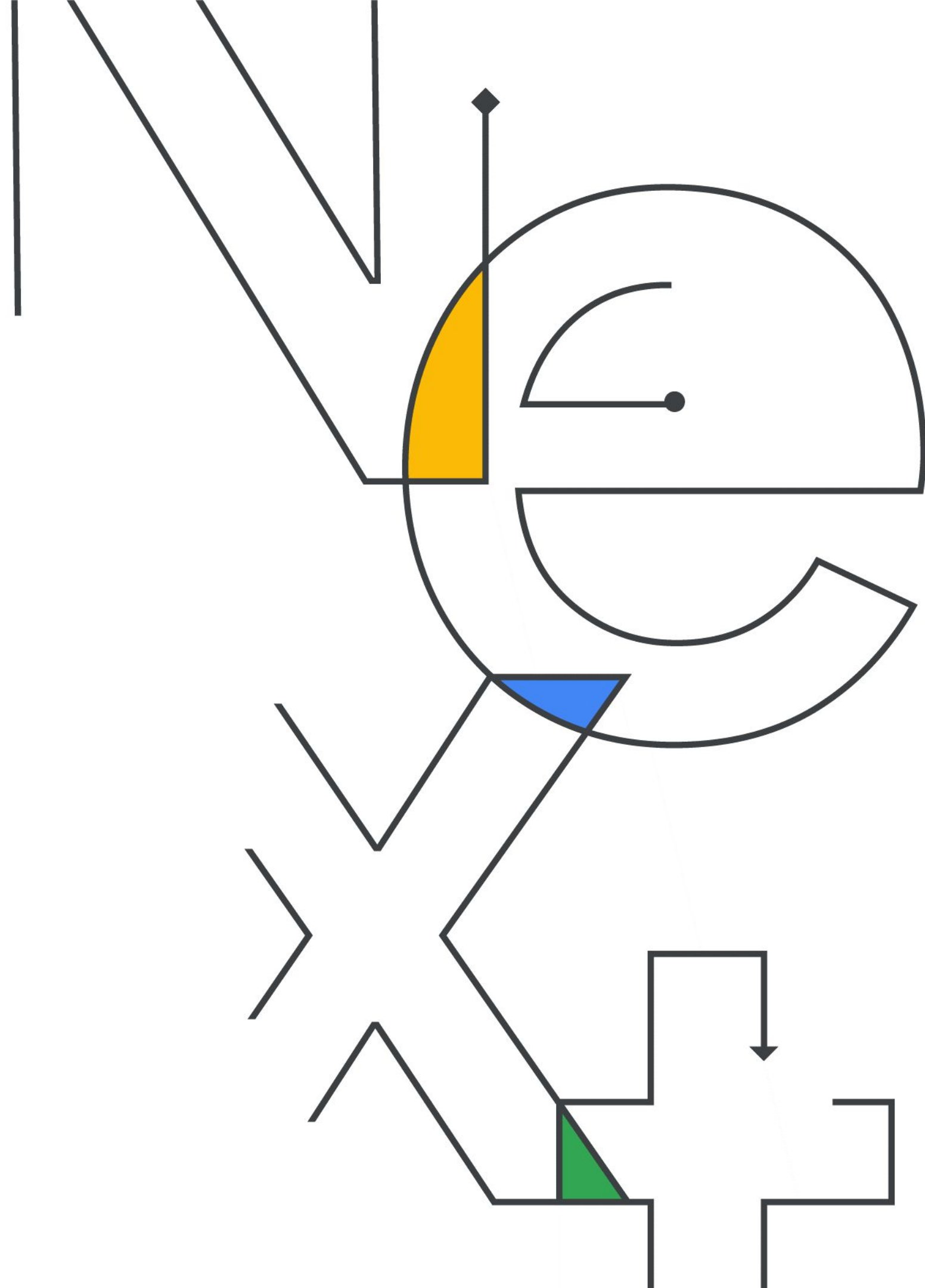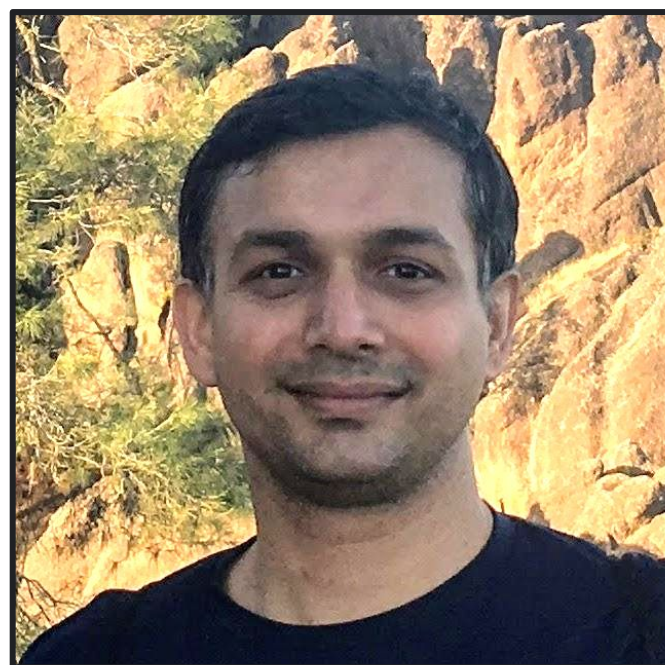Google Cloud

Next '22

# How to scale data analytics securely with Spark on Google Cloud

**Abhishek Kashyap**

Group Product Manager
**Google Cloud**

**Mithun Bondugula**

Sr Engineering Manager
**LiveRamp**

# Agenda

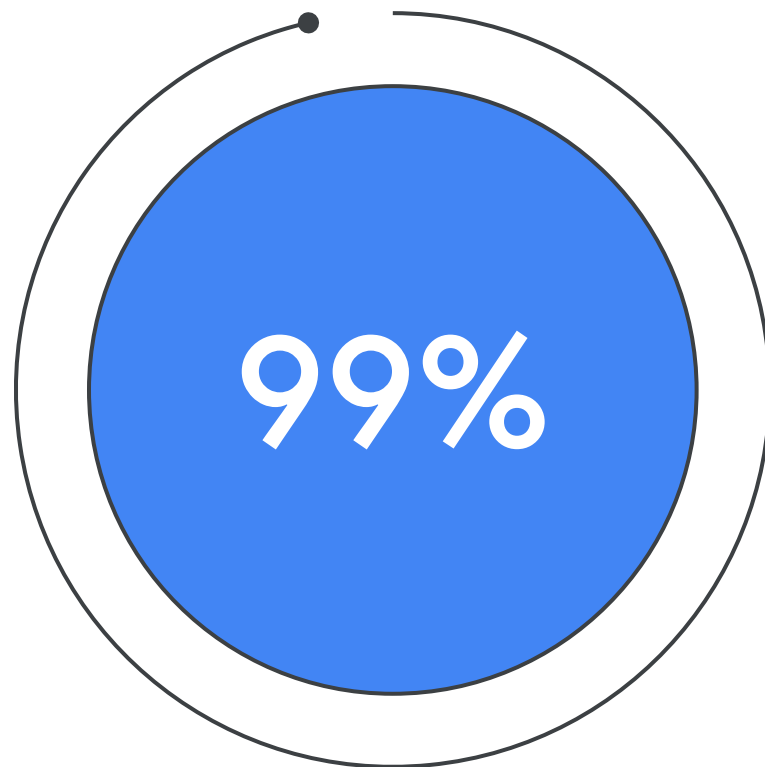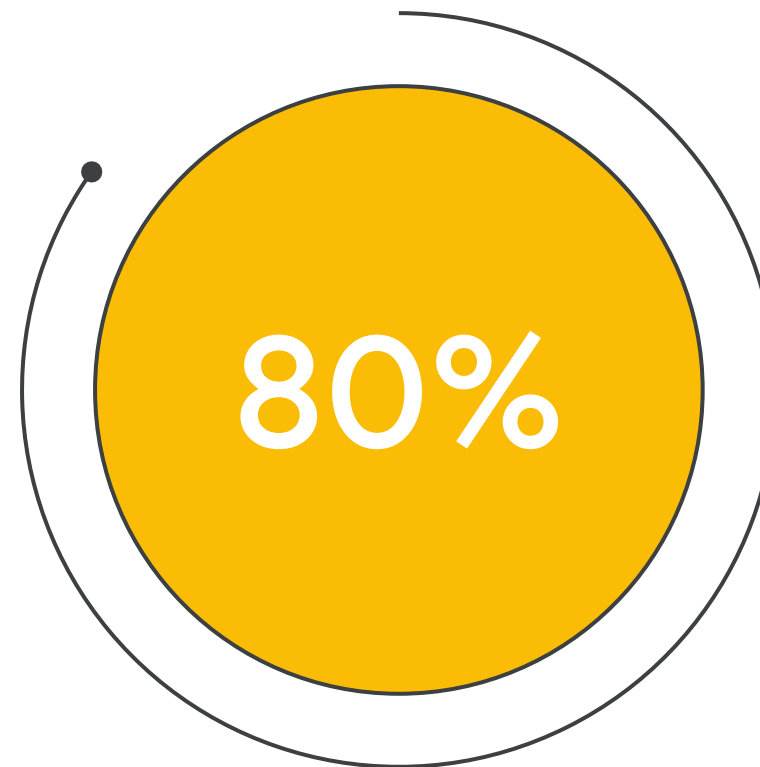Google Cloud

Organizations are doubling down on their open source investments as part of the overall data architecture.
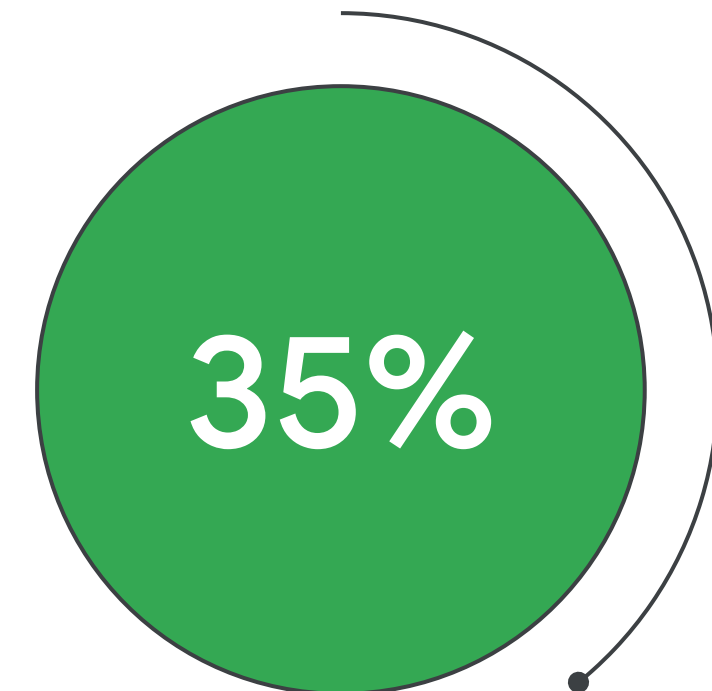
# The rise of open source

**99%**

of Fortune 500 companies currently use open source software

**80%**

of IT departments will increase use of open source in 2021

**35%**

of all enterprise software is based on open source code

Google Cloud

# At Google, we're committed to helping customers create an open and integrated data platform

**METRO**

**80%**
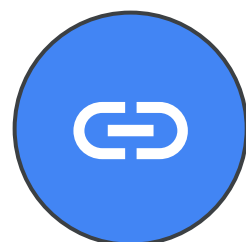Increase in ecommerce platform stability

**30-50%**
Reduction in infrastructure cost
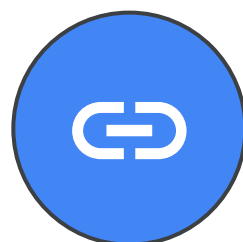
🔗 Learn more

**vodafone**

**17PB**
Of data migrated over to Google

**600**
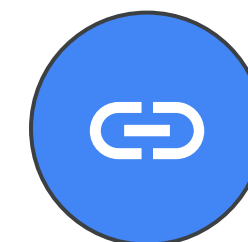Servers of Hadoop eliminated

🔗 Learn More

**/LiveRamp**

**30PB**
Of storage in Hadoop deployments migrated over to Google

**100,000+**

**YARN applications**
To deliver billions in records per day

🔗 Learn More

Google Cloud

# Spark on GCP

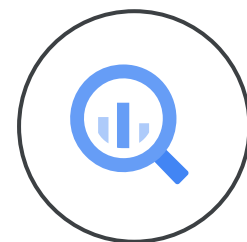| Scale without managing infrastructure | Work with tools you already know | Choose the right deployment model |
|---|---|---|
| • Auto-scale, without manual infrastructure provisioning or tuning<br><br>• Comes with latest OSS frameworks | • Connect, analyze, and execute Spark jobs from BigQuery, Vertex AI or Dataplex in 2 clicks<br><br>• No custom integrations, using Google-native and Open Source tools | **Choose between:**<br><br>• Serverless<br><br>• Google Kubernetes Engine (GKE)<br><br>• Compute clusters for your Spark applications |

BigQuery

Vertex AI

Dataplex

Composer

# Serverless Spark for ETL

> "
>
> OpenX used serverless Spark to **abstract** away all the **cluster resources** and just focus on the job itself. This significantly helped to **boost** the team's productivity, **while** reducing infrastructure costs.
>
> Marek Wolczanski, Data Platform Engineer, OpenX

### Customers

- PayPal
- ASML
- IAC
- LEVI STRAUSS & CO.
- GENERALI
- orange
- HKT

### Partner Ecosystem

- accenture
- Cognizant
- HCL
- TEKsystems Own change

Google Cloud

# What's new

## Serverless Spark

- Native Spark support in BigQuery
- Custom executor shapes (CPU:RAM)
- Customizable autoscaling speed
- Docker container streaming

## Security & Governance

- Fine grained governance through Big Lake
- Automated Dataproc policy management
- Dataproc Metastore Hive and BigQuery federation

# BigQuery stored procedures for Apache Spark
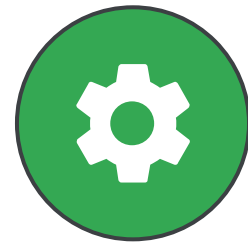
**BigQuery moves beyond SQL with new developer extensions**
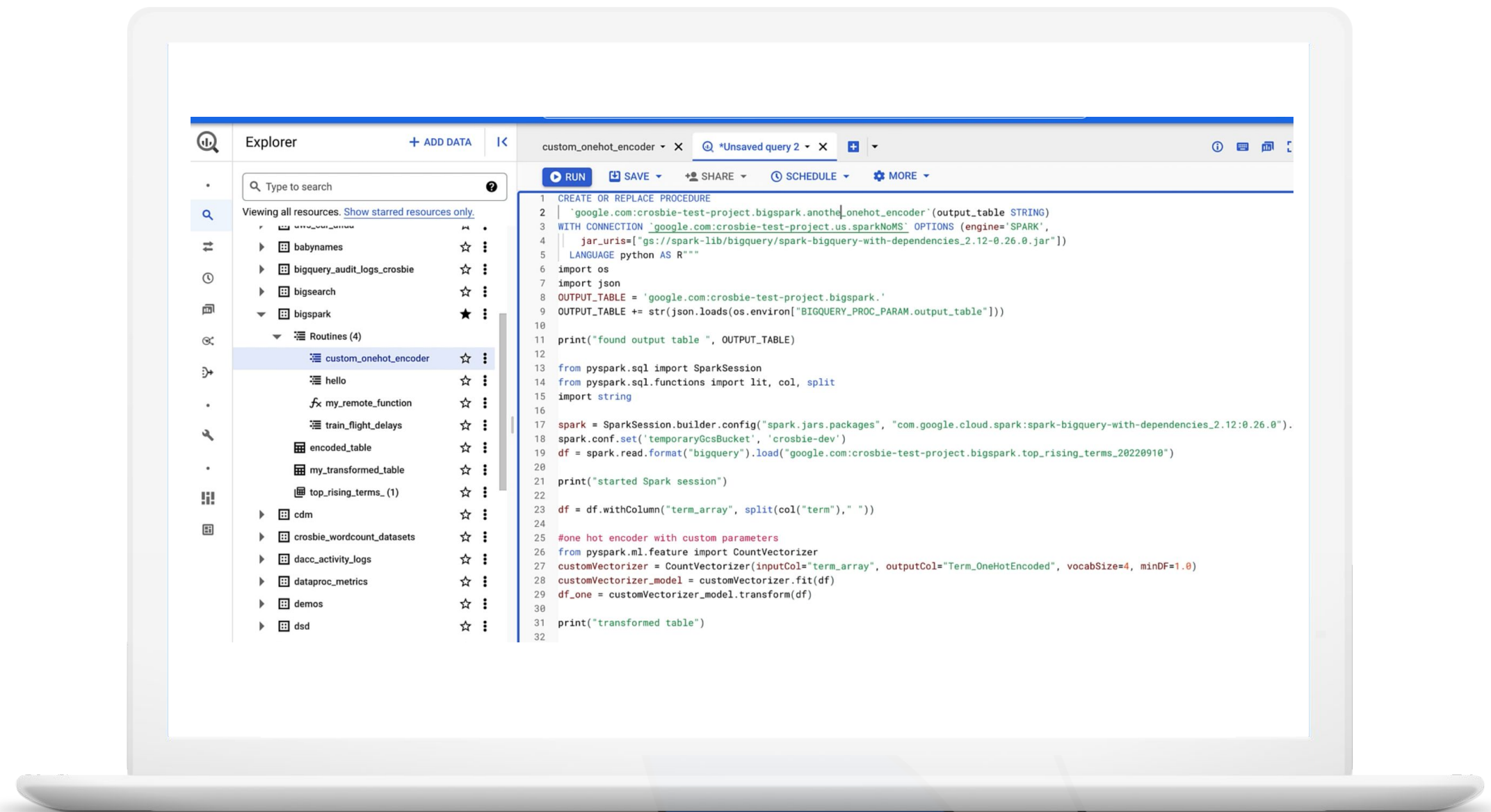
**Spark as a first class citizen in BigQuery**
Execute Spark optimally with BigQuery SQL, as a stored procedure

**Integrated BigQuery billing**
Use BigQuery reservations to execute Spark

**Integrated security and governance**
Manage access through BigQuery, no data analyst access to underlying Spark infrastructure



Google Cloud

# DSN201_Demo1_AbhishekKashyap

# Serverless Spark
# Interactive + Vertex AI for Data Science

## Accelerate data science development and MLOps pipelines

### Spark for Data Science in 1 click
Data scientists can use Spark for development from notebooks and Vertex AI workbench seamlessly

No cluster creation needed

### Built-in security and authentication
GCP security and user access are automatically applied from Vertex AI to Spark

### Integrate Spark with MLOps
Execute Spark code through Kubeflow pipelines

Google Cloud

# Open source templates

## tinyurl.com/dtproc-templates



Easily get started with serverless Spark for your use cases

### Templates

- 16+ Java templates
- 16+ Pyspark templates
- Notebooks

### Easy to use

- Open source
- Launch from cloud shell using inbuilt scripts

Google Cloud

# Fine grained governance for Spark through BigLake

## OSS Analytics on all of your data

**Any OSS engine on any data, anywhere, with a unified governance and access layer**

- Dataproc runs your OSS workloads

- Dataplex scales your data governance

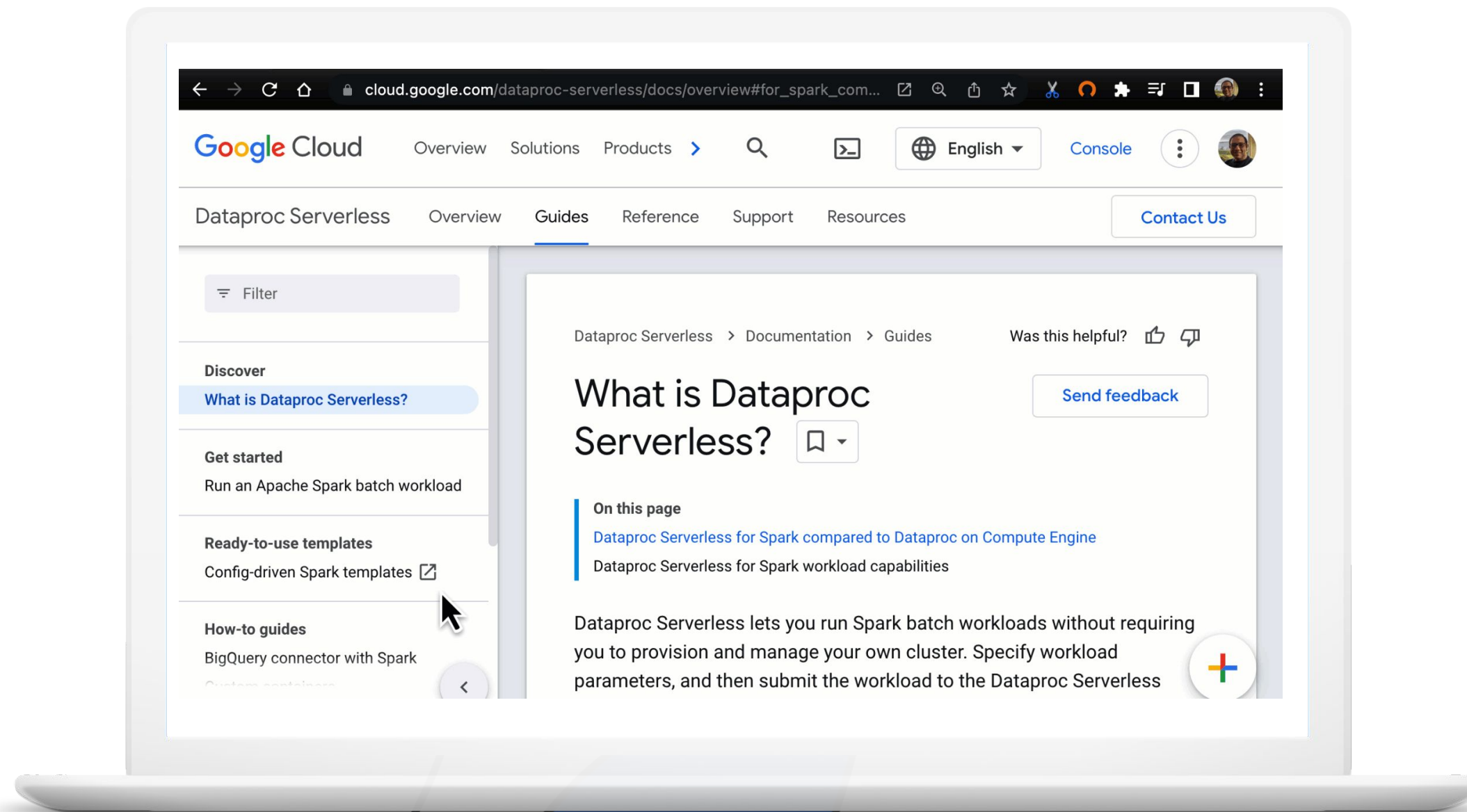- BigLake provides a standardized access layer with fine grained access control to any data

Multi-engine OSS Data Analytics with Dataproc

Governance at scale with Dataplex

Big Lake read/write with fine grained access controls

| BigQuery data | Open file and table formats, anywhere |
|---|---|

BigQuery managed storage

Google Cloud Storage

Amazon S3

Azure Data Lake Storage

**Google Cloud**

**Cross Cloud**

Google Cloud

# Automated Dataproc policy management

**GCP Folder**

**GCP Project**

Org

Marketing

Marketing
Analysis

Marketing
Data science

Finance

Finance
Reporting

Finance
Forecast

**Standardize config per Org, Folder, or Project**

- Resource policies for cost management, e.g., GPUs restricted to data science, VM configs

- Security policies, e.g., more stringent for projects dealing with PII data

- Network policies, e.g., internal IP only
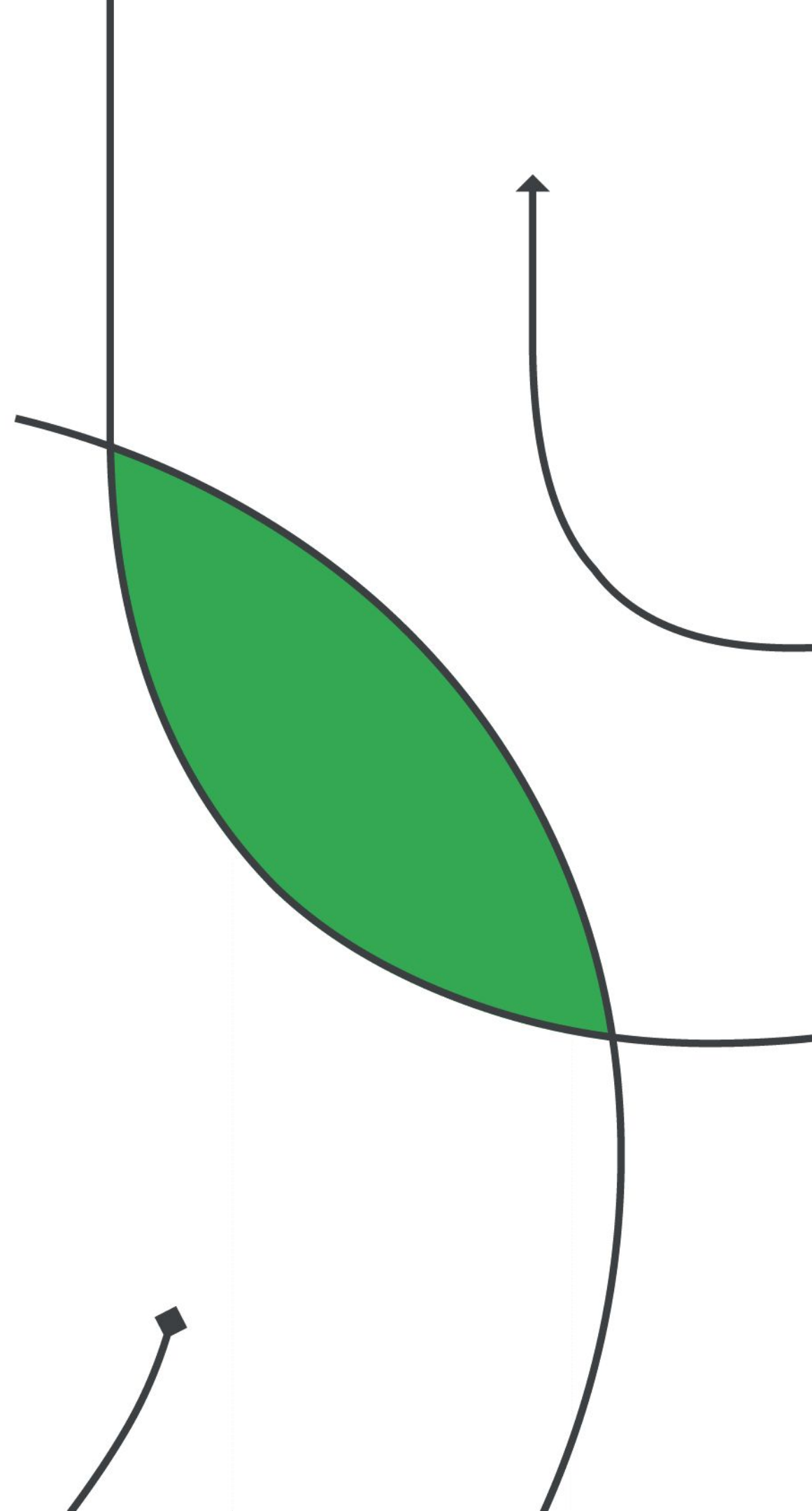
- Images and components

- Metastore configuration

Google Cloud

# Dataproc Metastore BigQuery Federation

- Read and write access to BQ tables from Hive metastore

- Fully integrated with BQ permissions

- Support for both DDL and DML statements from Spark

```scala
scala > spark.sql("create database bqdataset");
scala > spark.sql("show databases").show();
+-----------------+
|        namespace|
+-----------------+
|          default |
|         bqdataset|
+-----------------+
scala > spark.sql("create table bqtable(id int, name string);
scala > spark.sql("desc bqtable").show();
+--------+---------+-------+
|col_name|data_type|comment|
+--------+---------+-------+
|      id|   bigint|   null|
|    name|   string|   null|
+--------+---------+-------+
```

Google Cloud

/LiveRamp

We connect consumer data with durable privacy-conscious post-cookie identifiers for more accurate customers views, improved measurement, and secure data collaborations.

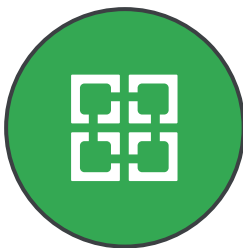# At LiveRamp, we make it safe and easy for companies to use data effectively.
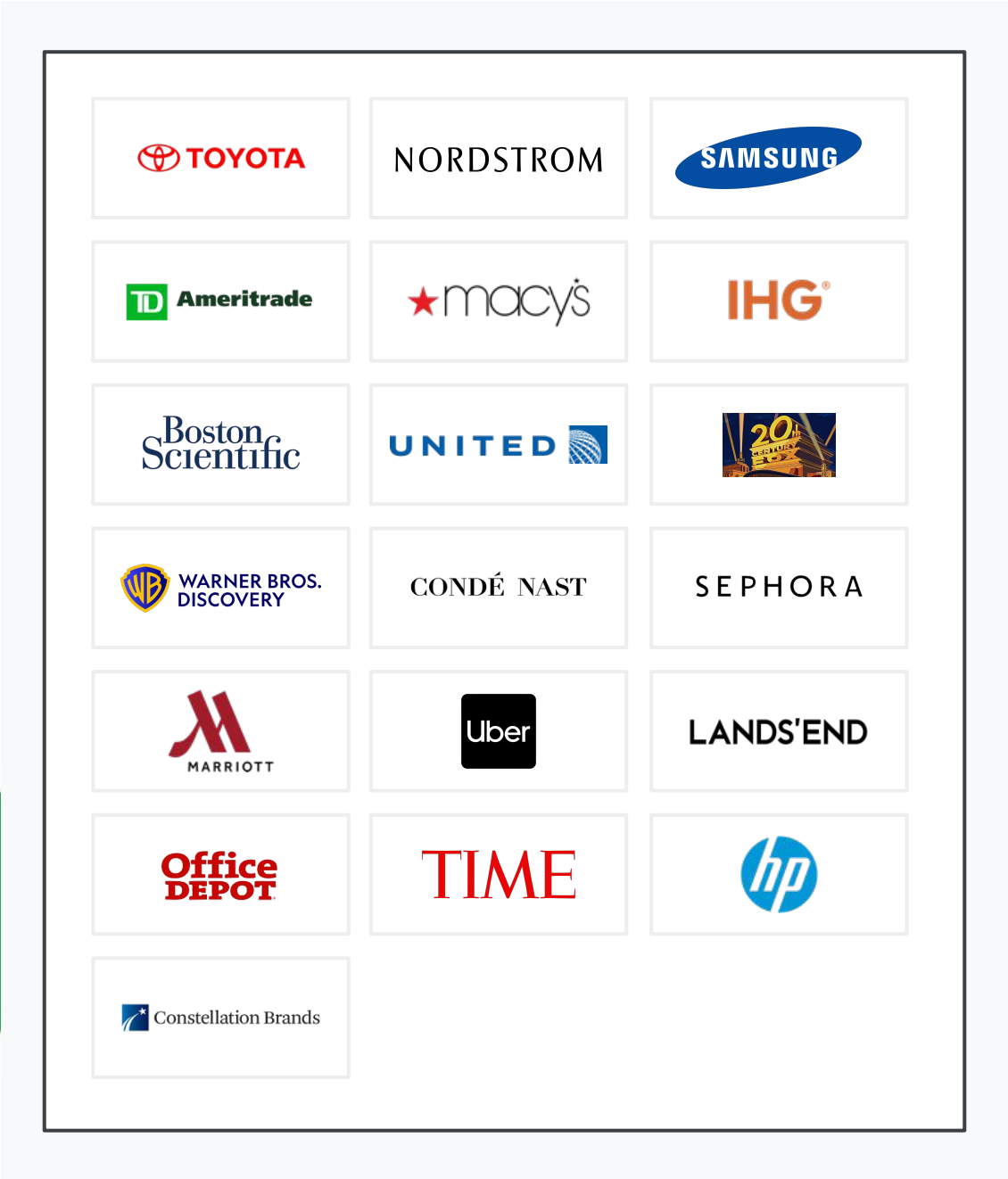
**Identity resolution**

**Data enrichment and audience activation**

**Data collaboration and media networks**

**LiveRamp identity platform with embedded cloud services**

# LiveRamp engineering

## Team

Spans 11 countries with a total of > 450 engineers. LiveRamp is proud to be named in Fortune's Best Workplaces in Technology™ 2021.
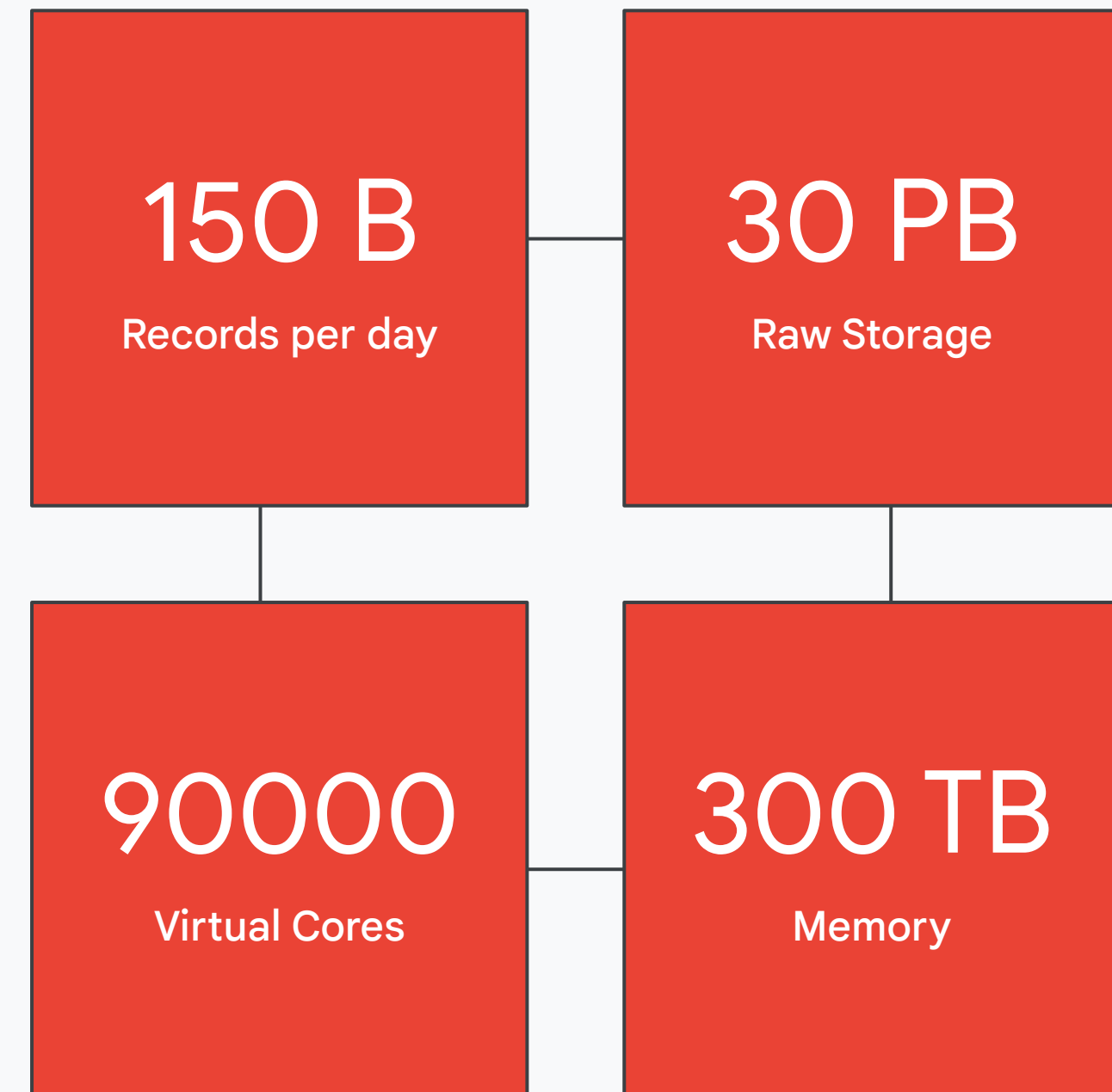
## Data processing

Not map-reduce but map-join

Biggest Identity graph

## Technology stack

Open source technologies

Multi-cloud support (GCP, AWS, Snowflake)

## Problem and infrastructure scale

| | |
|---|---|
| **150 B**<br>Records per day | **30 PB**<br>Raw Storage |
| **90000**<br>Virtual Cores | **300 TB**<br>Memory |

Google Cloud

# Migrating to GCP

## Architectural Decisions

Decentralized team ownership

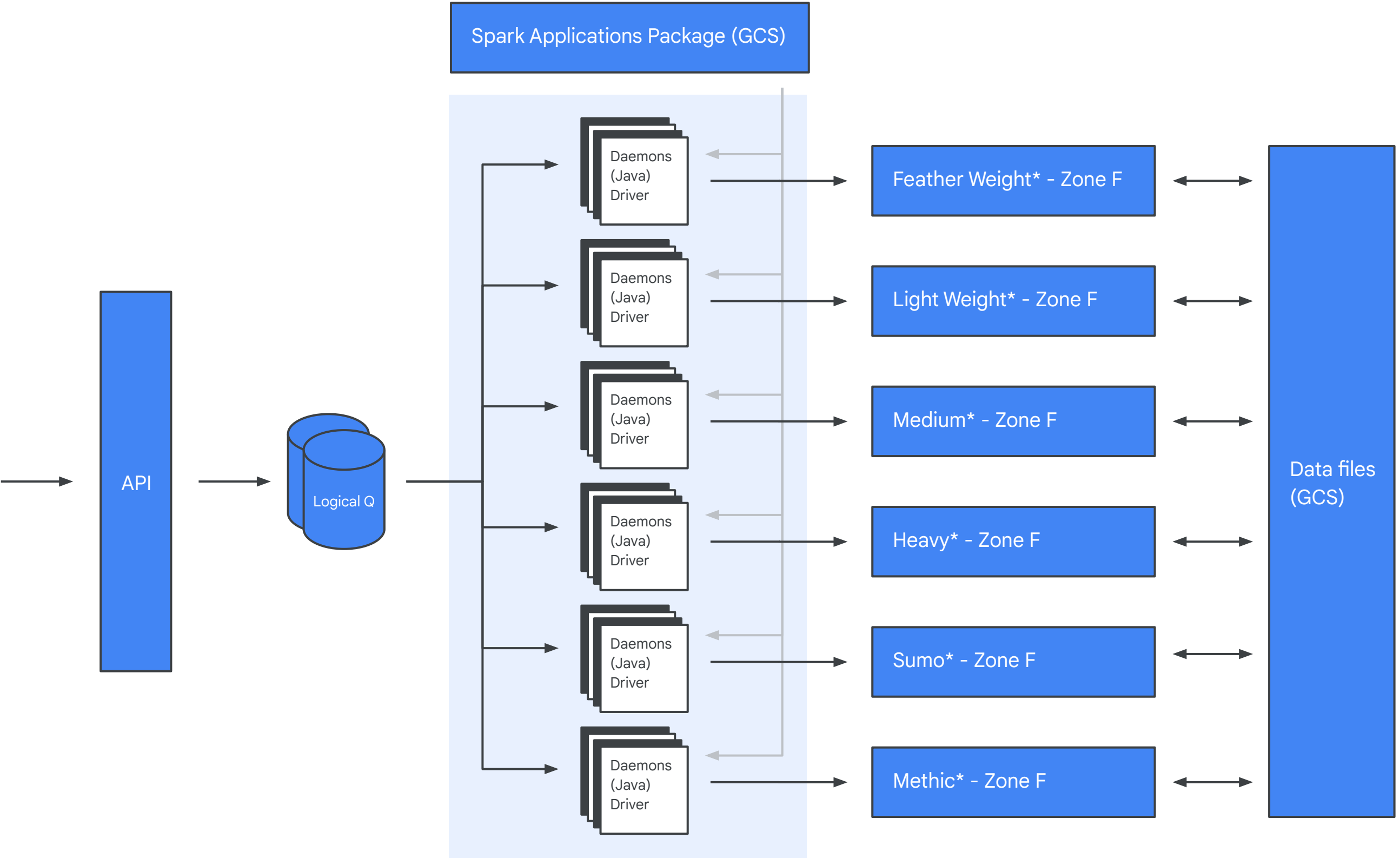HDFS -> GCS

Autoscaling clusters

## Infrastructure as Code

Self-service tooling for data engineers to deploy infrastructure as easily and safely as possible

## Map Side Join (MSJ) Library

Schema that defines virtual partitions, a strategy for assigning records to those partitions, and a library that handles writing/reading that data in a distributed fashion.

## LiveRamp's largest workload Infrastructure

Spark Applications Package (GCS)

API

Logical Q

Daemons (Java) Driver

Daemons (Java) Driver

Daemons (Java) Driver

Daemons (Java) Driver

Daemons (Java) Driver

Daemons (Java) Driver

Feather Weight* - Zone F

Light Weight* - Zone F

Medium* - Zone F

Heavy* - Zone F

Sumo* - Zone F

Methic* - Zone F

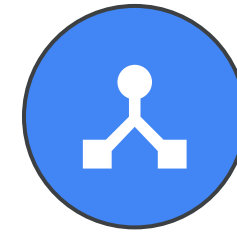Data files (GCS)

Google Cloud

# Key benefits on GCP

## Support

- Weekly sync-ups
- Roadmap collaboration
- Clear escalation path

## Cost attribution

- Cost attributed at Project/Cluster/Asset (VM etc.) using tagging
- Monitoring and alerting to protect against the cost overruns

## Flexibility

- Self-service Terraform module
- Agility in cluster management - create, delete etc.
- Environments for A/B testing

## Cost savings

- ~ 30% cost savings in some clusters
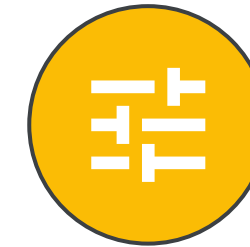- Some applications is now 10x faster

Google Cloud

# Where we're headed

**Performance**

GPUs and VM shape

**Scale**

Driver Pools

**Management**

Enhanced Monitoring

Alerting

**Agility**

VertexAI

Serverless DataProc

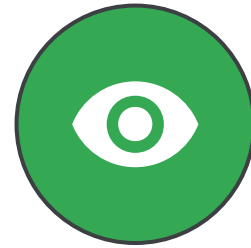Google Cloud

# Lessons learned on migrating to the Cloud

### Benchmark

Ensure current behavior of key workloads is clearly understood before migration

### Focus

For us, reliability first then cost optimization post-migration

### Preview < > GA

Understand feature stability to judge risk and time for adoption

### Quotas

We had to increase IP space and change quotas as we tried out new VM and disk settings

### Discounts

Understand impact on Committed Use Discounts (CUDs) when you start to migrate to Spot VMs

Google Cloud

Google Cloud

Next '22

How to scale data analytics securely with Spark on Google Cloud