Google Cloud

# Data Analytics on GCP

Spyrales.fr

Pascal Rabier
2020-04-30

# 15+ Years of Tackling Big Data Problems

**Open Source**

**Google Papers**

GFS | Map Reduce | BigTable | Dremel | PubSub | Flume Java | Millwheel | Tensorflow Dataflow

**Google Cloud Products**
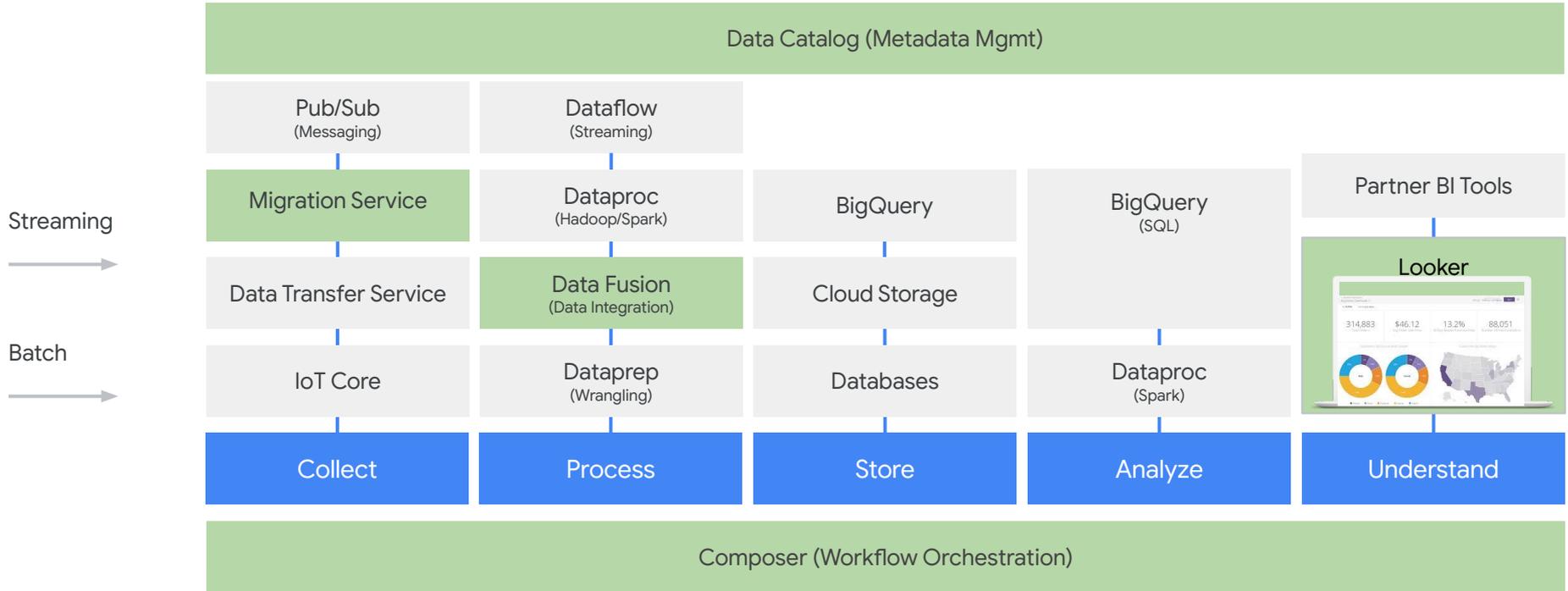
BigQuery | Pub/Sub | Dataflow | Bigtable

2002 2004 2005 2006 2008 2010 2012 2014 2015 2016

# Fully managed storage & database services

| Object | Key-value | Non-relational | | Relational | | Warehouse |
|--------|-----------|----------------|---|-----------|---|-----------|
| **Cloud Storage** | **App Engine Memcache** | **Cloud Firestore** | **Cloud Bigtable** | **Cloud SQL** | **Cloud Spanner** | **BigQuery** |
| Binary or object data | Web/mobile applications, gaming | Hierarchical, mobile, web | Heavy read + write, events | Web frameworks | RDBMS+scale, HA, HTAP | Enterprise Data Warehouse |
| Images, media serving, backups | Game state, user sessions | User profiles, Game State | AdTech, financial, IoT | CMS, eCommerce | Transactions, Ad/Fin/MarTech | Analytics, Dashboards |

Google Cloud

# Google's Smart Analytics Platform

Collect, process, store, analyze and visualize data and insights

| Data Catalog (Metadata Mgmt) | | | | |
|---|---|---|---|---|
| Pub/Sub (Messaging) | Dataflow (Streaming) | | | |
| Migration Service | Dataproc (Hadoop/Spark) | BigQuery | BigQuery (SQL) | Partner BI Tools |
| Data Transfer Service | Data Fusion (Data Integration) | Cloud Storage | | Looker |
| IoT Core | Dataprep (Wrangling) | Databases | Dataproc (Spark) | |
| Collect | Process | Store | Analyze | Understand |

Streaming →

Batch →

Composer (Workflow Orchestration)

**Smart Analytics as a Service**: Fully Managed. Serverless. Enterprise class. Globally Distributed. Secure

# Providing choice to customers

## Cloud Native Services

### Differentiation

BigQuery

Dataflow

Pub/Sub

Data Catalog

looker

## Managed Open Source Services

### Familiarity

Composer

beam

Dataproc
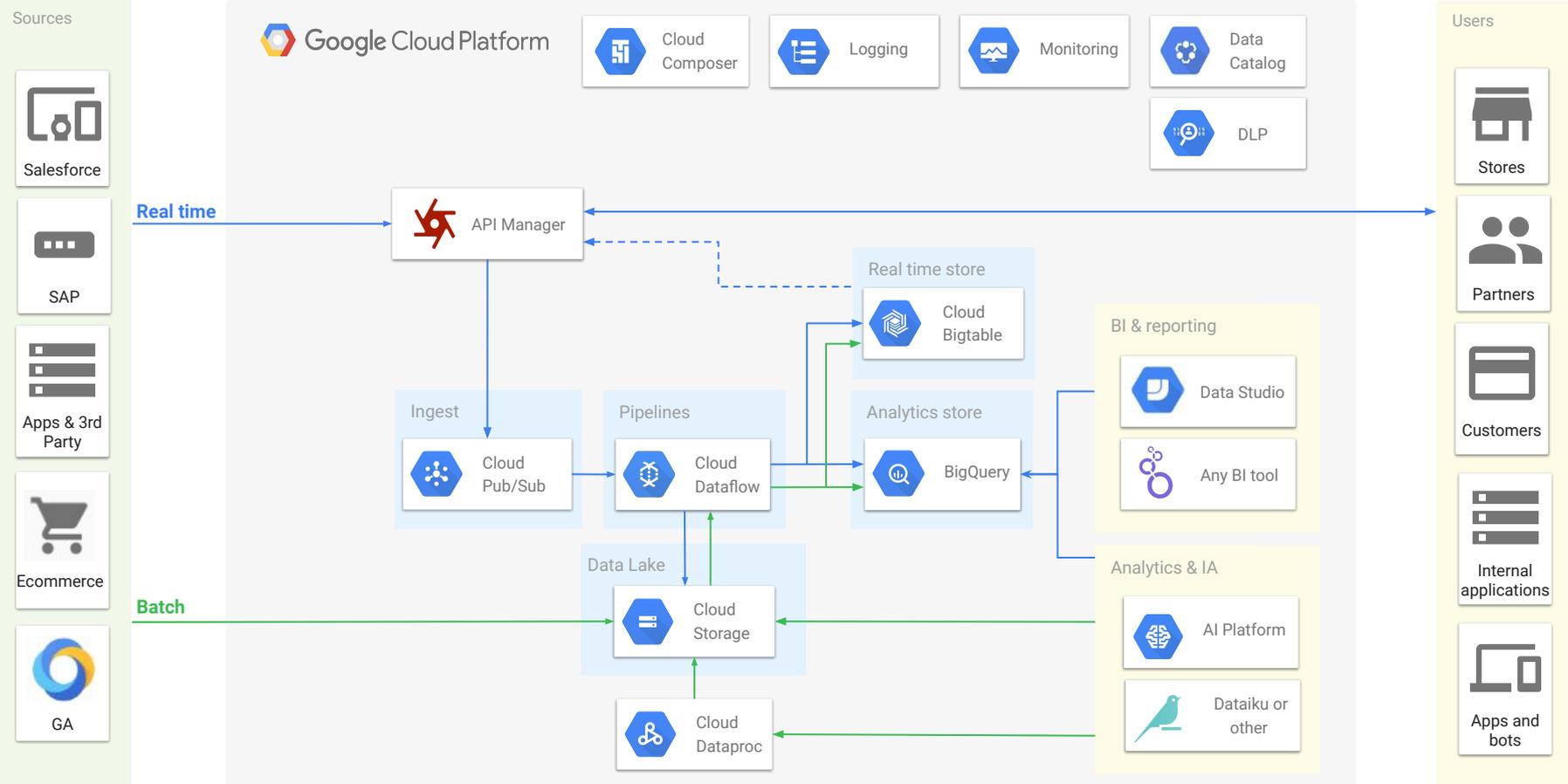
Data Fusion

## Partner Services

### Completeness

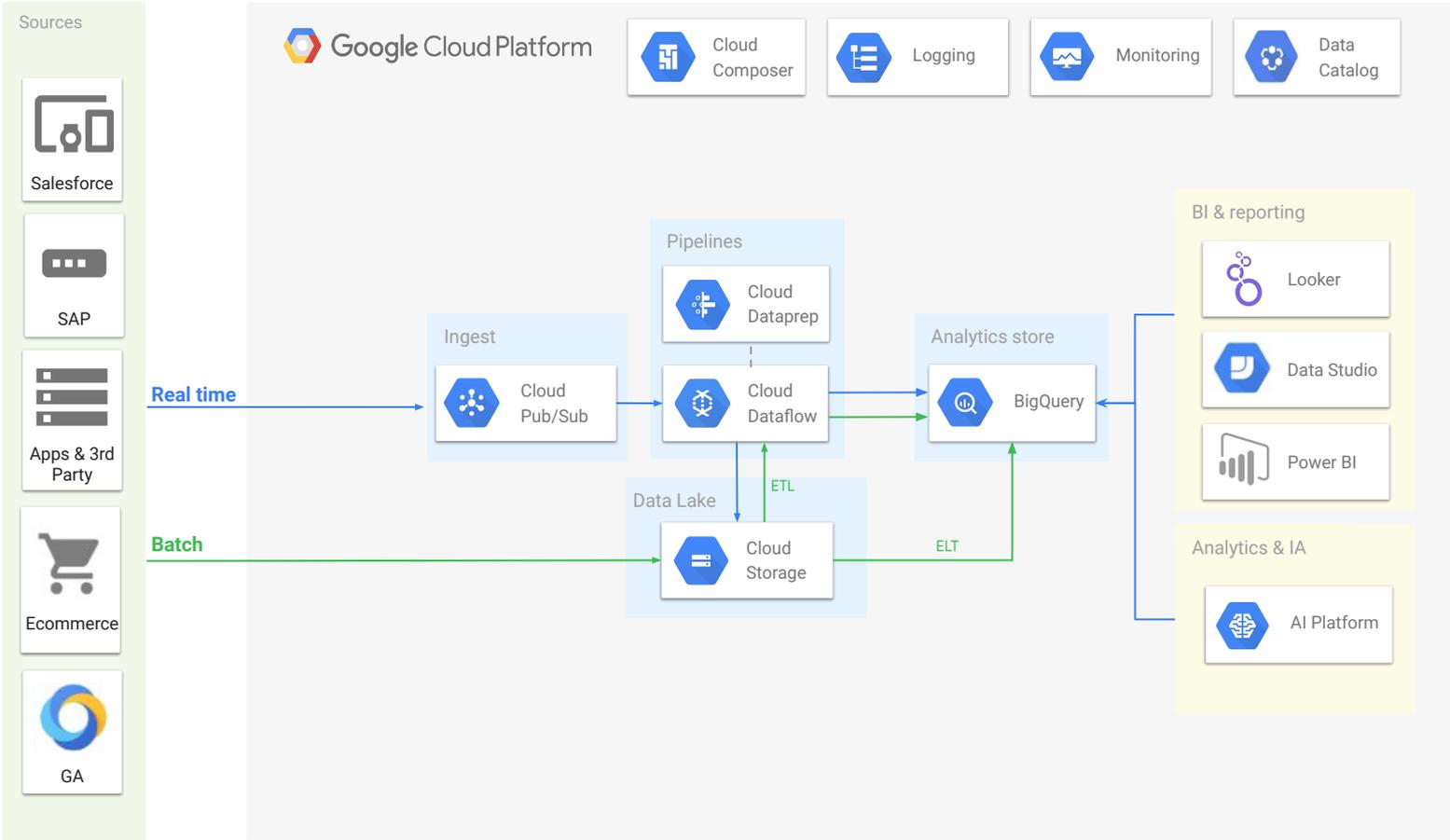confluent

Dataprep

Informatica

collibra

Qubole

Google

# An example of a Big Data architecture with GCP

# First step

# BigQuery

**BigQuery**

Google Cloud Platform's
**enterprise data warehouse**
for analytics

Gigabyte- to **petabyte-scale**
storage and SQL queries

**Encrypted,** durable,
And highly available

Fully managed and **serverless**
for maximum agility and scale

**Unique**

**Real-time** insights from streaming data

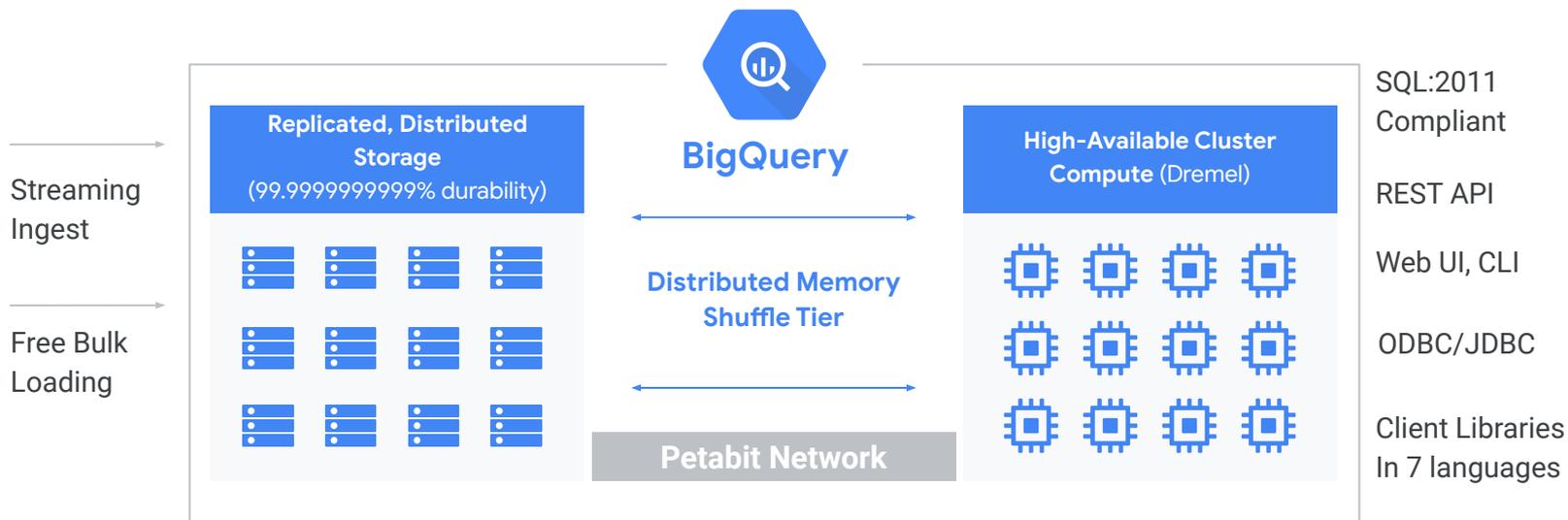**Unique**

Built-in **ML** for out-of-the-box
predictive insights

**Unique**

High-speed, in-memory **BI Engine**
for faster reporting and analysis

**Unique**

Google Cloud

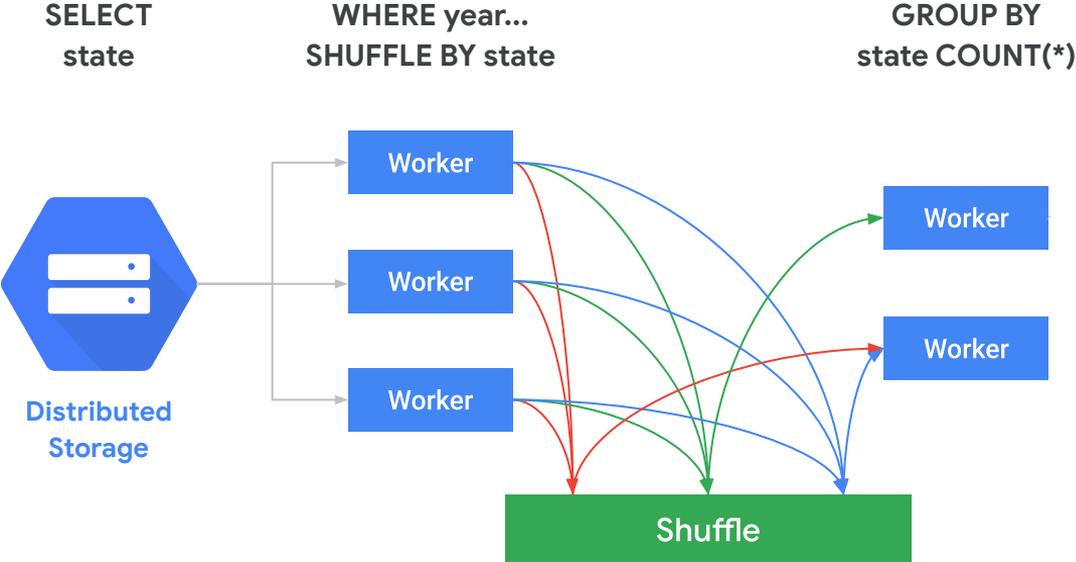# BigQuery | Architecture

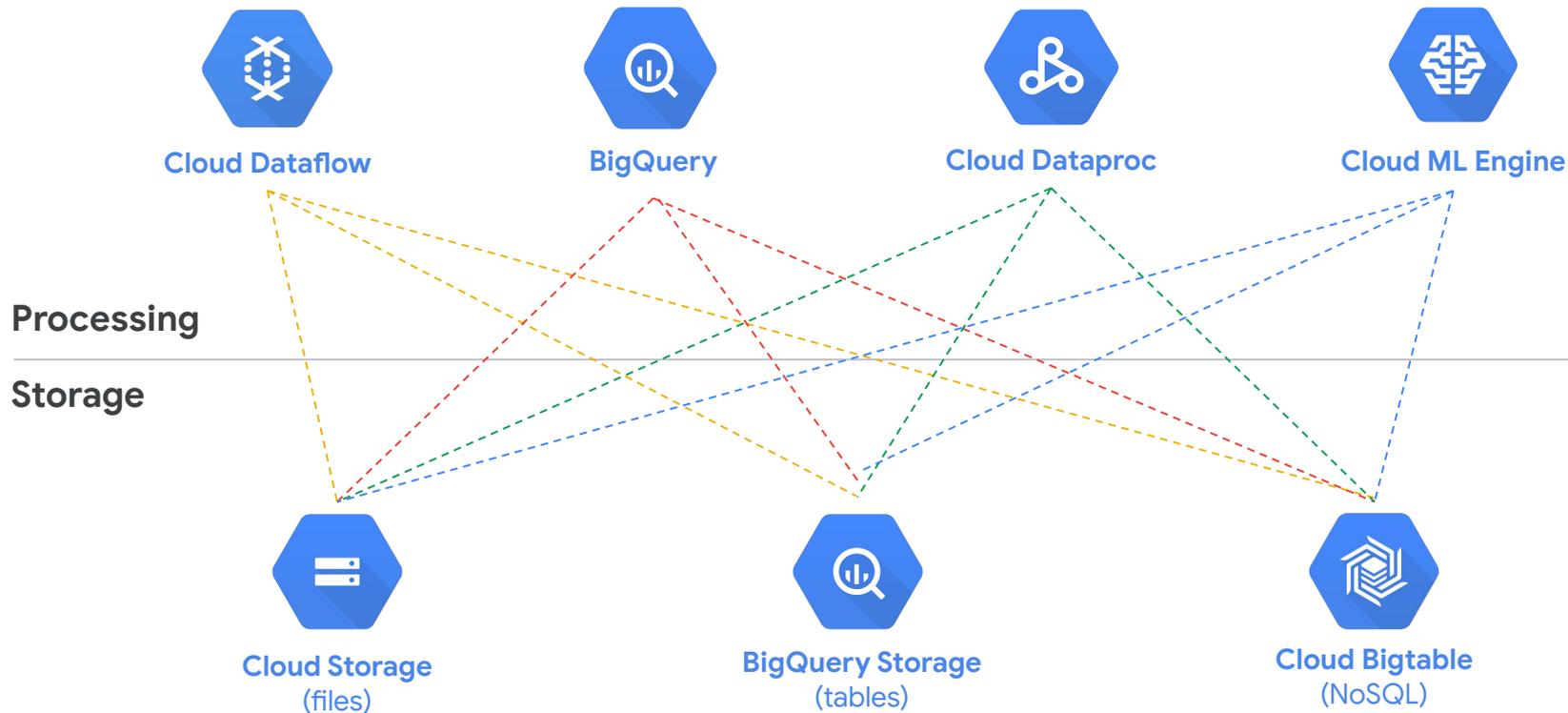Decoupled storage and compute for maximum flexibility

# BigQuery remote memory shuffle

Faster performance for complex queries

Join and aggregate more data

Better scalability

**SELECT state**

**WHERE year... SHUFFLE BY state**

**GROUP BY state COUNT(*)**

Distributed Storage

Worker

Worker

Worker

Worker

Worker

Shuffle

Google Cloud

# Separation of storage and compute

**Cloud Dataflow**

**BigQuery**

**Cloud Dataproc**

**Cloud ML Engine**

**Processing**

**Storage**

**Cloud Storage**
(files)

**BigQuery Storage**
(tables)

**Cloud Bigtable**
(NoSQL)

# BigQuery platform interoperability

## BigQuery Storage API

Use BigQuery Storage like GCS for Dataflow and Dataproc, break down the Data Warehouse storage wall

Run high-performance **dataframes** on BigQuery

## Cloud SQL and Cloud Bigtable Federation

Query your Cloud SQL and Cloud Bigtable instances directly from BigQuery, without moving data around.

## Parquet & ORC Federation

Query Parquet and ORC files directly in GCS



**Cloud Dataflow**

**Cloud Dataproc**

BQ Storage API

**BigQuery**

Query Federation

**Parquet & ORC in GCS**

**Cloud SQL**

**Cloud Bigtable**

# You don't have to take my word for it

# Enterprise-grade Workload management With **Reservations**

**BigQuery Reservations allows customers to:**

- Control flat-rate spend
- Buy slots in Web UI in seconds
- Efficiently manage workloads in BigQuery
- Automatically share any unused capacity



Google Cloud

# Introducing Flex slots

- A new <u>commitment</u> type
  - Alongside monthly & annual
- Pricing
  - $30 per slot per month*
- More flexible
  - 60 second minimum
- Combine with monthly/annual
- Available in all BQ Reservations regions!
- Available in **BigQuery Reservations** today!

Google Cloud

*May vary per region

---

BigQuery          ← Buy Slots BETA

**1 Configure**

Configure your BigQuery slot commitment.

BigQuery offers flat-rate pricing as a predictable, fixed budget option. Flat rate customers purchase dedicated BigQuery slot commitments for query execution, and associated projects, folders, or organization are not subject to per-query charges.

BigQuery commitments are offered at commitment durations of one month (30 days) or one year. You cannot cancel until your commitment end date.

˅ SHOW DETAILS

Commitment duration *

Monthly (default duration)
Annual
Flex

Number of slots *

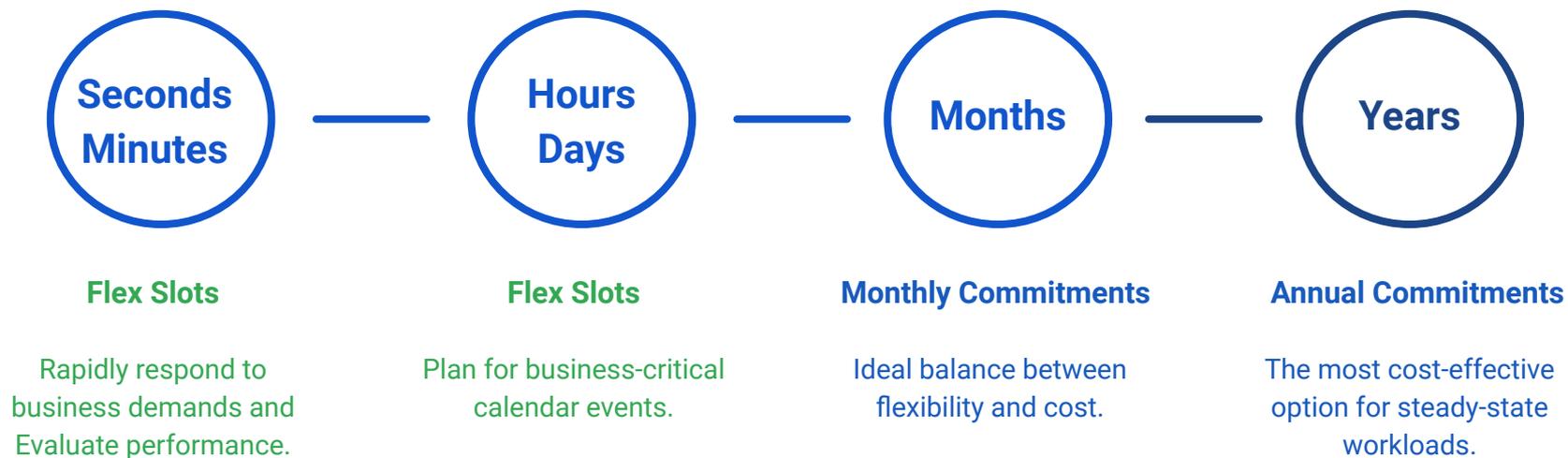Slots can be reserved in increments of 500.

NEXT

**2 Confirm and submit**

**3 Confirmation**

# **BigQuery** Commitment Types and Use Cases

**Seconds Minutes** — **Hours Days** — **Months** — **Years**

**Flex Slots**

Rapidly respond to business demands and Evaluate performance.

**Flex Slots**

Plan for business-critical calendar events.

**Monthly Commitments**

Ideal balance between flexibility and cost.

**Annual Commitments**

The most cost-effective option for steady-state workloads.

Google Cloud

# **BigQuery** workload management

**Customers can programmatically perform workload management using Reservations:**

Create and delete reservations

Move projects between reservations

Move slots between reservations

Idle slots are seamlessly and automatically shared in real-time

**Example**

At 3am an important workload in project_d needs to run

**At 6am we delete the**
**At 3am we create a reservation**reservation
Move 1000 slots to the reservation    Move 1000 slots back
Move project_d into reservation       Move project_d back

**Project_d was guaranteed 1000 slots 3am-6am**
**Idle slots seamlessly shared**

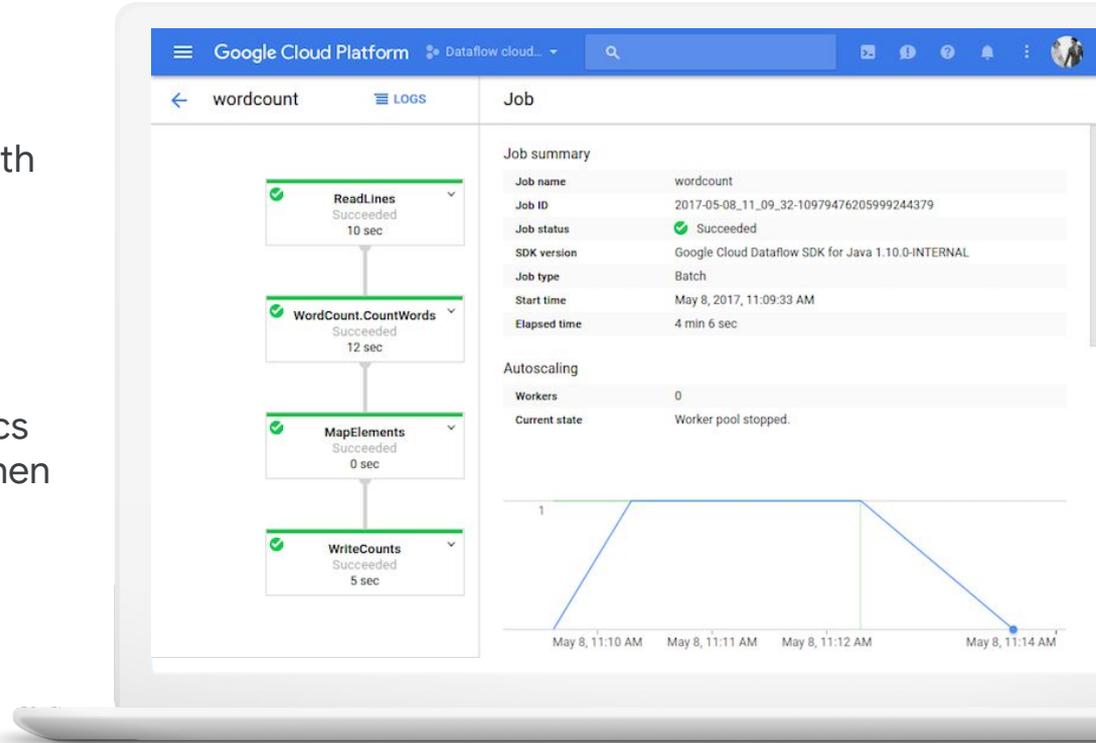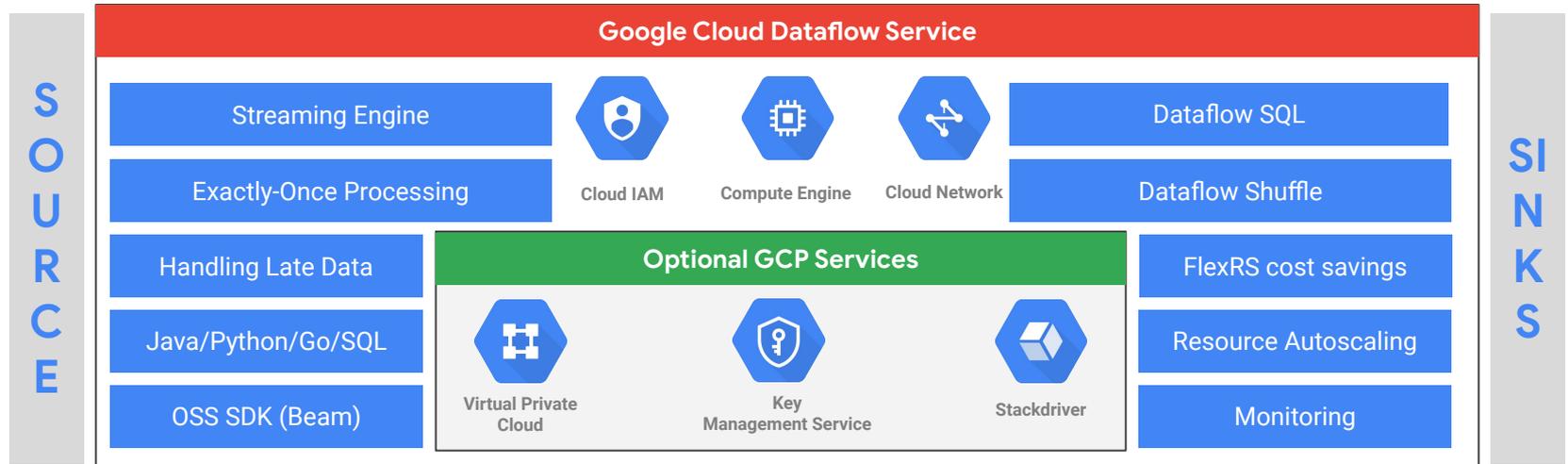| On-Demand | Default<br>*1000 slots* | BI<br>*1500 slots* | best-effort<br>*30 slots* |
|---|---|---|---|
| *Project _f* | Project_a<br>Project_b<br>Project_c | Project_d | Project_e |

Google Cloud

# ETL / Scale streaming analytics pipelines with
# Cloud Dataflow

Streaming analytics service that minimizes processing time and cost with autoscaling while blending **batch and stream** processing.

- Fastest stream and batch processing on one service
- Lower TCO for streaming analytics
- Automatically burst resources when data spikes
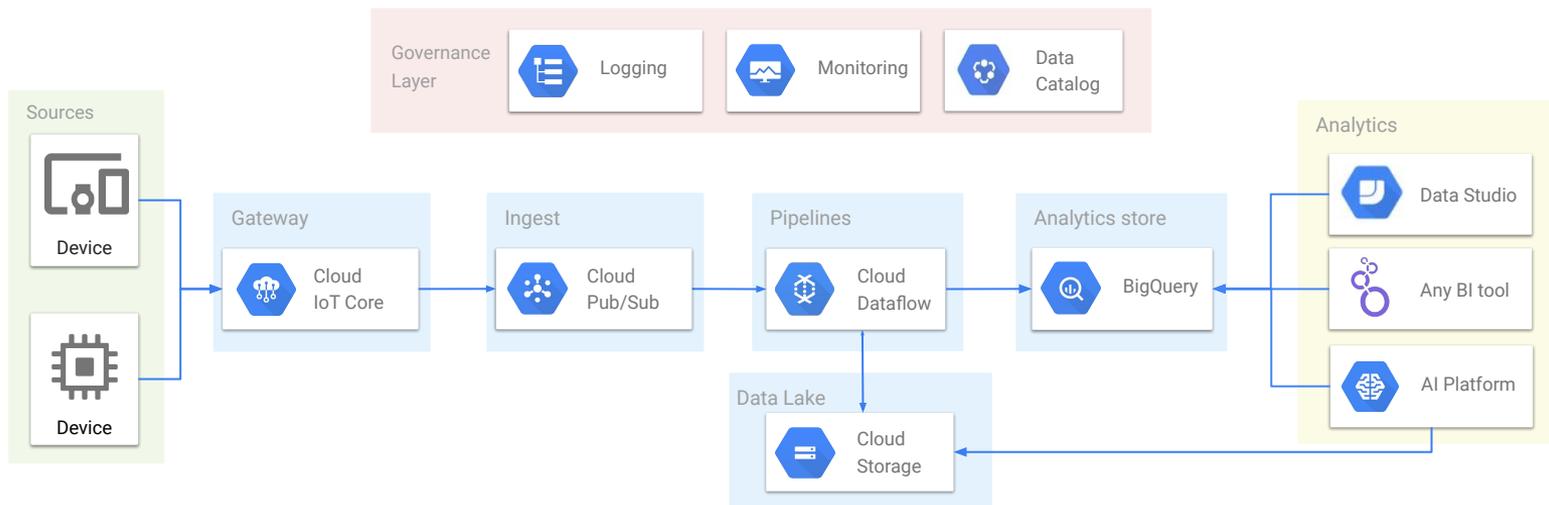- Build and monitor Apache Beam pipelines



Google Cloud

# Dataflow: Stream Analytics as a managed service



**Google Cloud Dataflow Service**

**SOURCE**

Streaming Engine

Exactly-Once Processing

Cloud IAM

Compute Engine

Cloud Network

Dataflow SQL

Dataflow Shuffle

Handling Late Data

**Optional GCP Services**

FlexRS cost savings

Java/Python/Go/SQL

Virtual Private Cloud

Key Management Service

Stackdriver

Resource Autoscaling

OSS SDK (Beam)

Monitoring
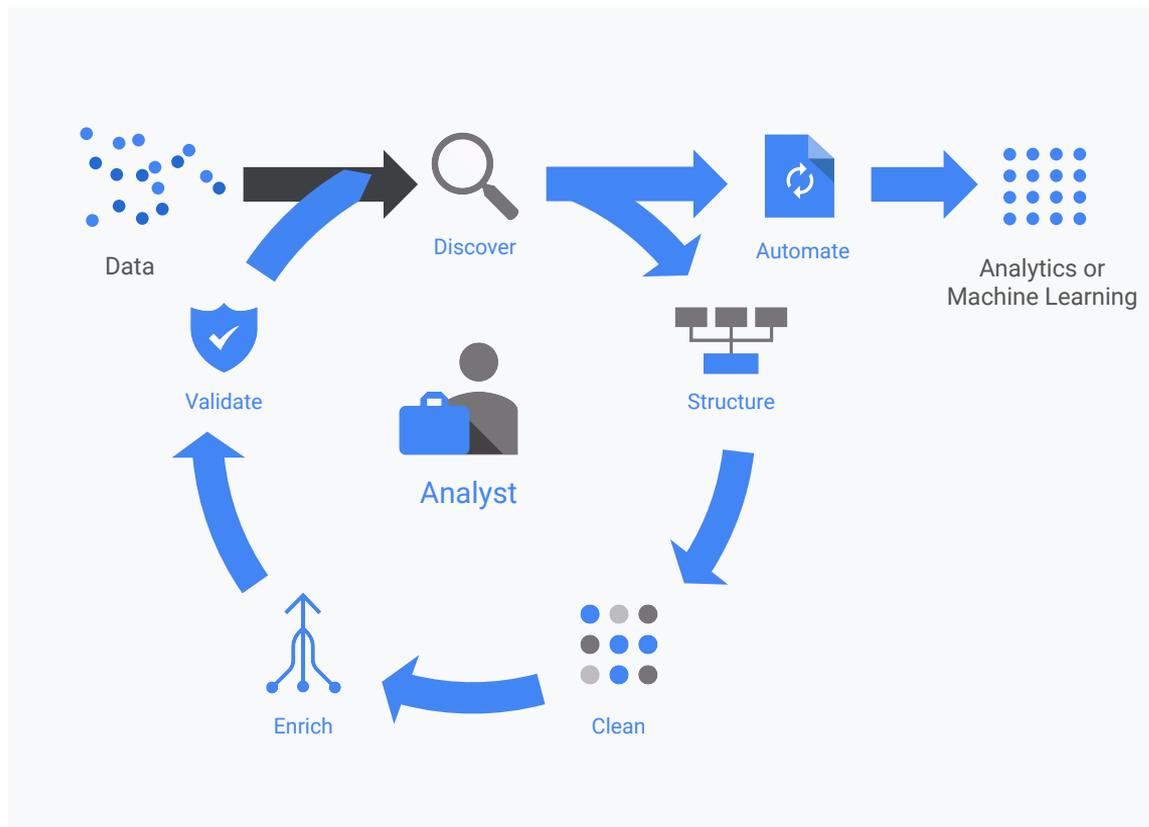
**SINKS**

Google Cloud

# Demo: A simple Streaming reference architecture

Scales seamlessly to petabytes to let you focus on bringing actual value

# Simplify
# the data lifecycle
# with **Cloud**
# **Dataprep**



Data

Discover

Automate

Analytics or
Machine Learning

Validate

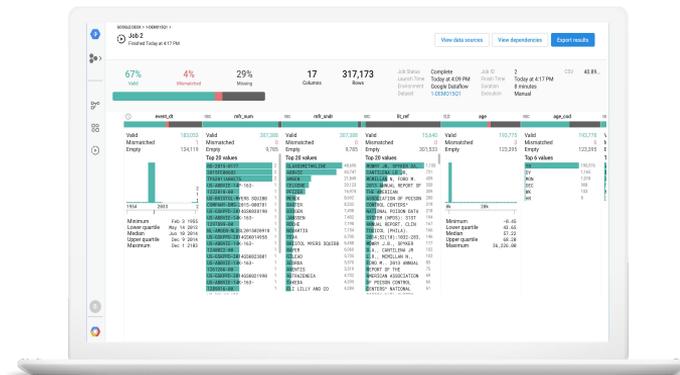Analyst

Structure

Enrich

Clean

# Serverless and cloud-native

## Legacy data preparation

❌ Business users not empowered to transform data samples

❌ Must hire an IT/Data ops team and manage a Hadoop cluster

❌ Negotiate org-wide software licenses, arrange billing and manage seats

❌ Integrate application permissions with infrastructure permissions

## Modern data preparation on Cloud Dataprep:

✅ Business users push the "Run Job" button to apply transformations to datasets of any size

✅ No need to create or manage infrastructure

✅ No need to provision software licenses
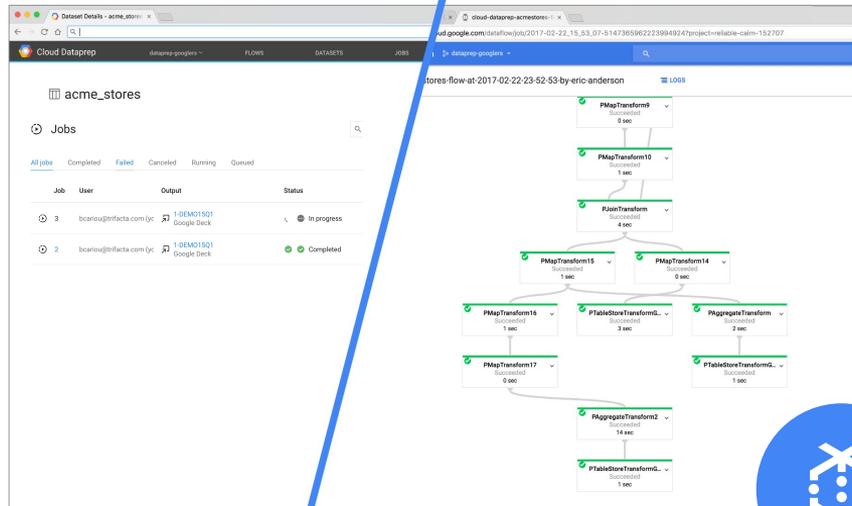
✅ Integrated, and highly scalable

| **Serverless Simplicity** | Fast Exploration | Easy Preparation |

Google Cloud

# Powerful & easy processing with Cloud Dataflow under the hood

✓ Process diverse datasets - structured or unstructured

✓ Prepare datasets of any size, PB or MB, with equal ease

✓ Leverages Cloud Dataflow without needing to write any scripts

✓ Auto-scalable and can easily handle processing massive data sets



| Serverless Simplicity | Fast Exploration | **Easy Preparation** |

Google Cloud

# Supports common data types of any size

## Sources

**BigQuery** tables

**Cloud Storage** or local upload using common file formats:

- CSV
- JSON
- TXT
- LOG
- GZIP
- BZIP

## Targets

**BigQuery** tables

**Cloud Storage**:

- CSV (compressed or not)
- JSON (compressed or not)
- Avro

| Serverless Simplicity | Fast Exploration | **Easy Preparation** |
|---|---|---|

# Cloud Dataproc
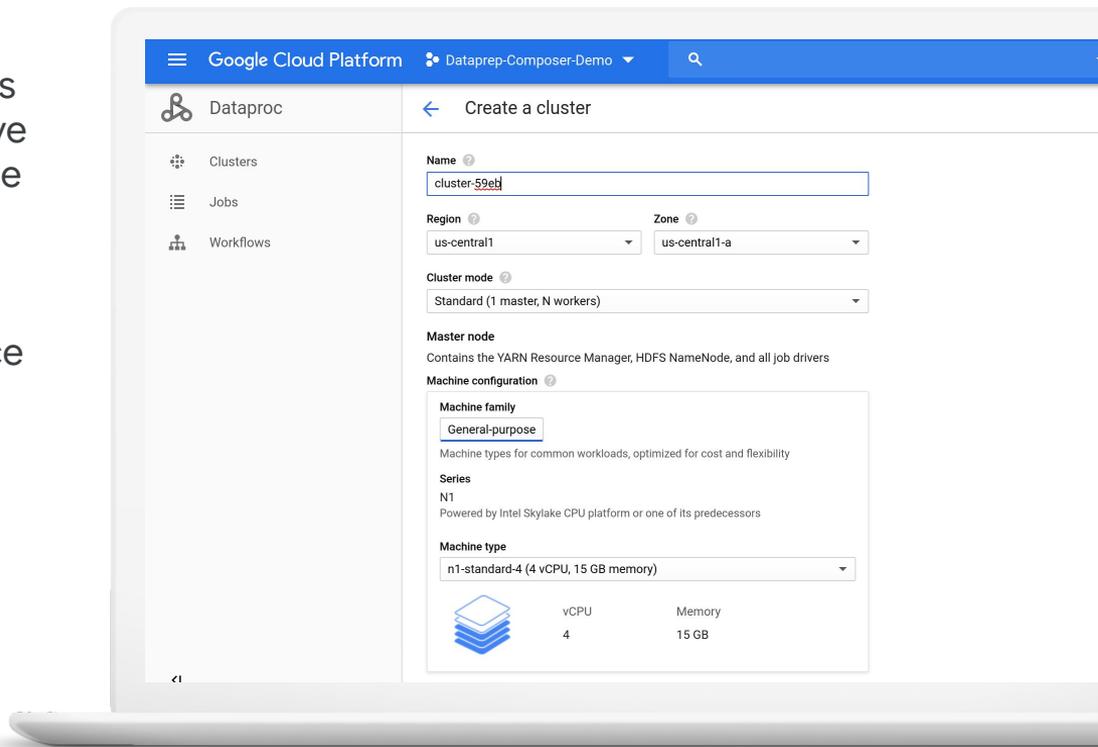
Combining the best of open source and cloud.

# Open source data and analytics processing at scale on Cloud Dataproc

Build data and analytics processing jobs using the open source software you love with the scale, security, and governance of the cloud.

- Autoscale SQL, batch, streaming, and machine learning open source processing (Apache MapReduce, Apache Spark, Presto, etc.)
- Lower TCO of running OSS
- Build Spark jobs on Kubernetes



Google Cloud

# The benefits of Hadoop/Spark on Cloud

|  | On premises | On compute engine | Cloud Dataproc |
|--|-------------|-------------------|----------------|
| Custom code | Custom code | Custom code | Custom code |
| Monitoring/Health | Monitoring/Health | Monitoring/Health | Monitoring/Health |
| Dev integration | Dev integration | Dev integration | Dev integration |
| Scaling | Scaling | Scaling | Scaling |
| Job submission | Job submission | Job submission | Job submission |
| GCP connectivity | GCP connectivity | GCP connectivity | GCP connectivity |
| Deployment | Deployment | Deployment | Deployment |
| Creation | Creation | Creation | Creation |

■ Self-managed   ■ Google managed
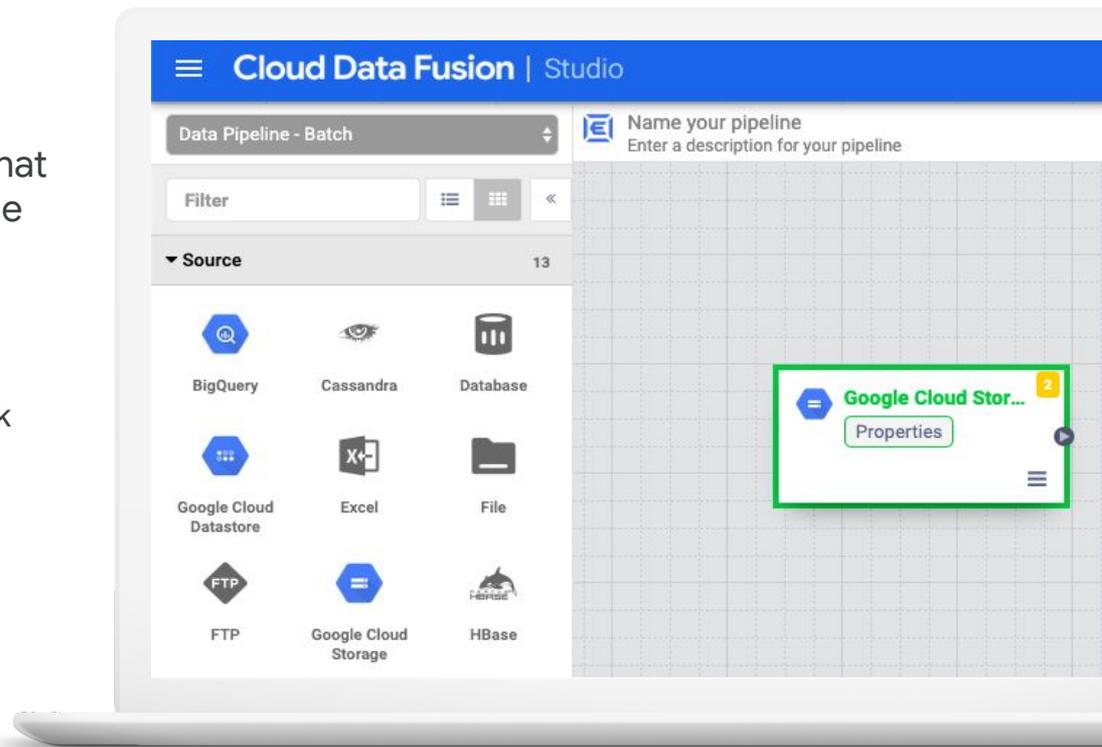
Google Cloud

# Build code free data pipelines with Data Fusion

Cloud Data Fusion is a fully managed, cloud-native data integration service that helps users efficiently build and manage ETL/ELT data pipelines.

- Use pre build open source library of connectors
- Execute data pipelines in Apache Spark
- Metadata and lineage integrations
- Build Apache Kakfa pipelines



Google Cloud

That's a wrap.

pascalr@google.com

Google Cloud