# Make the cloud work for you

Easy, fast, and low-cost streaming with Apache Flink on Google Cloud Platform
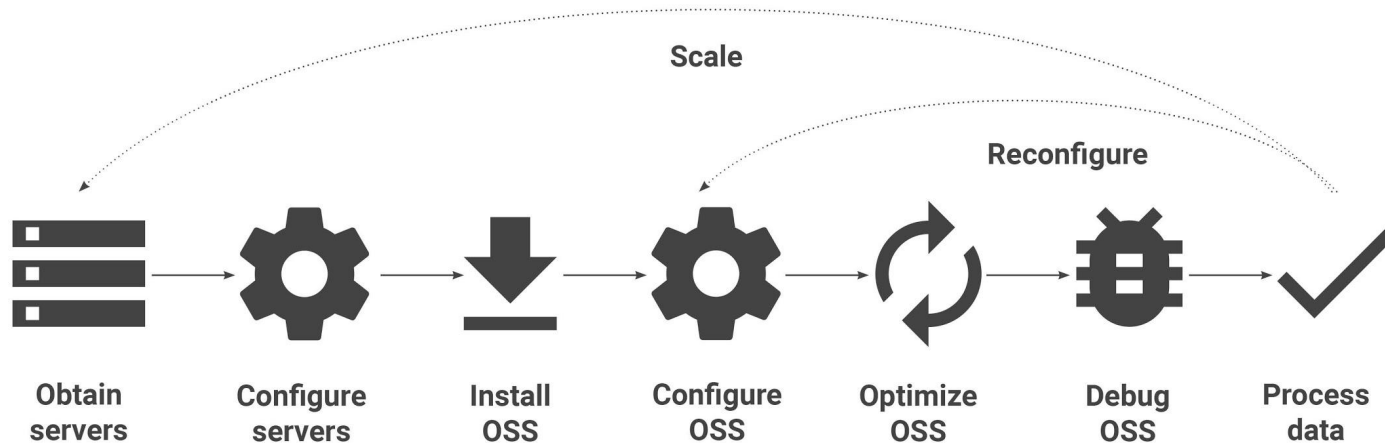
James Malone [ jamesmalone@google.com ]

Open source
**powerful** but **complex**

# The Apache ecosystem

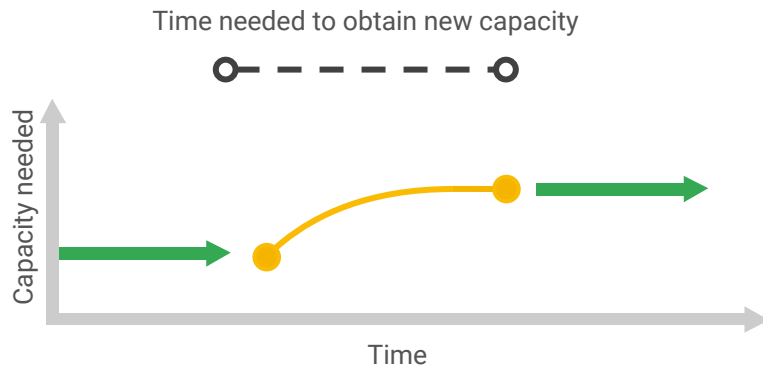# Typical OSS deployments



Scale

Reconfigure

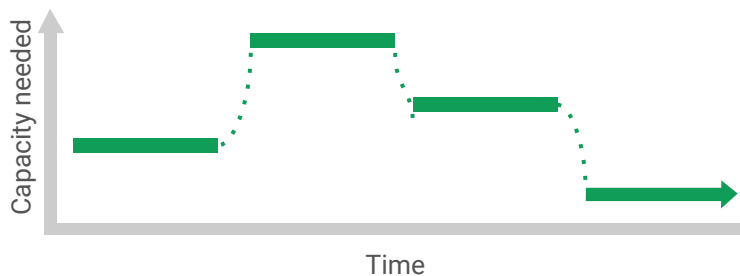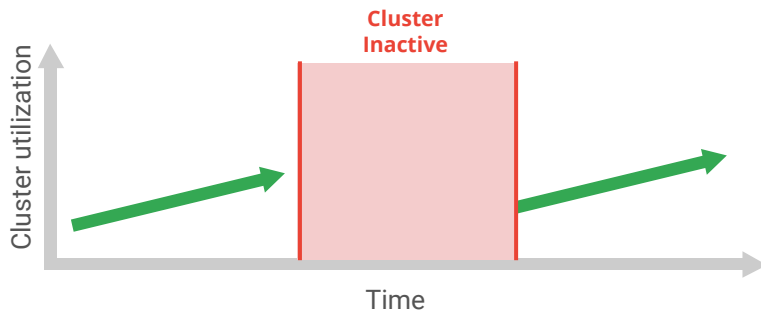| Obtain servers | Configure servers | Install OSS | Configure OSS | Optimize OSS | Debug OSS | Process data |

# Scaling makes your life difficult



Scaling can take hours, days, or weeks to perform which may delay needed data processing

Google Cloud

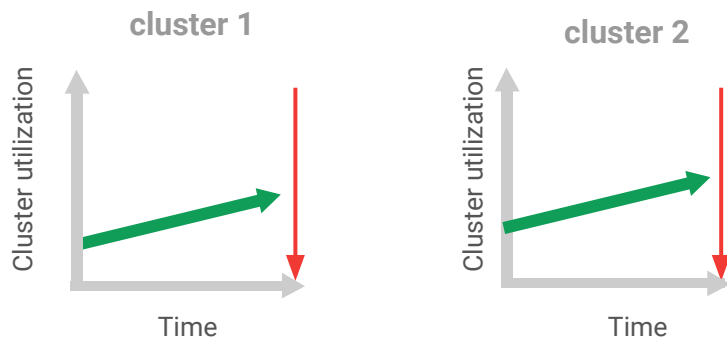# Scaling should be painless and fast

Capacity needed

Time

Things take seconds to minutes,
not hours or weeks.

# You have to babysit utilization



Requires effort to pack clusters
so the it does not have periods of
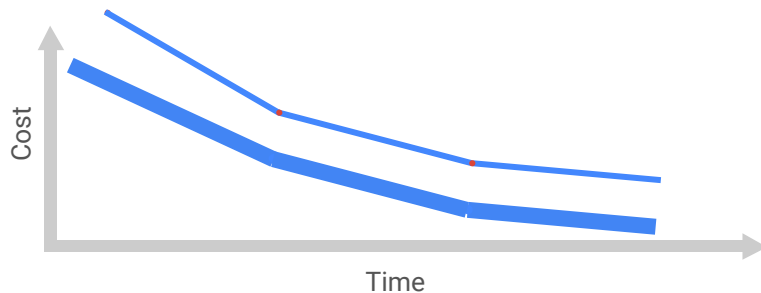inactivity and wasted resources

# Only use clusters when you need them



cluster 1      cluster 2

Cluster utilization

Time

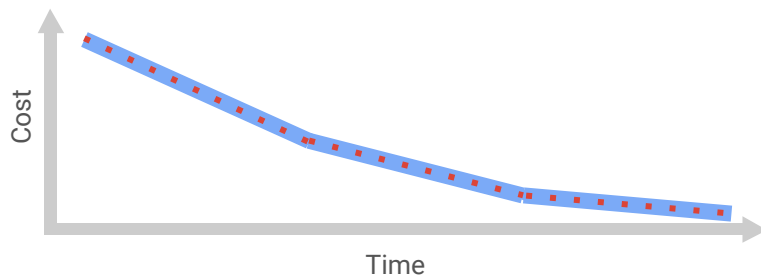Be an expert with your data, not your infrastructure

# You are not paying for what you use



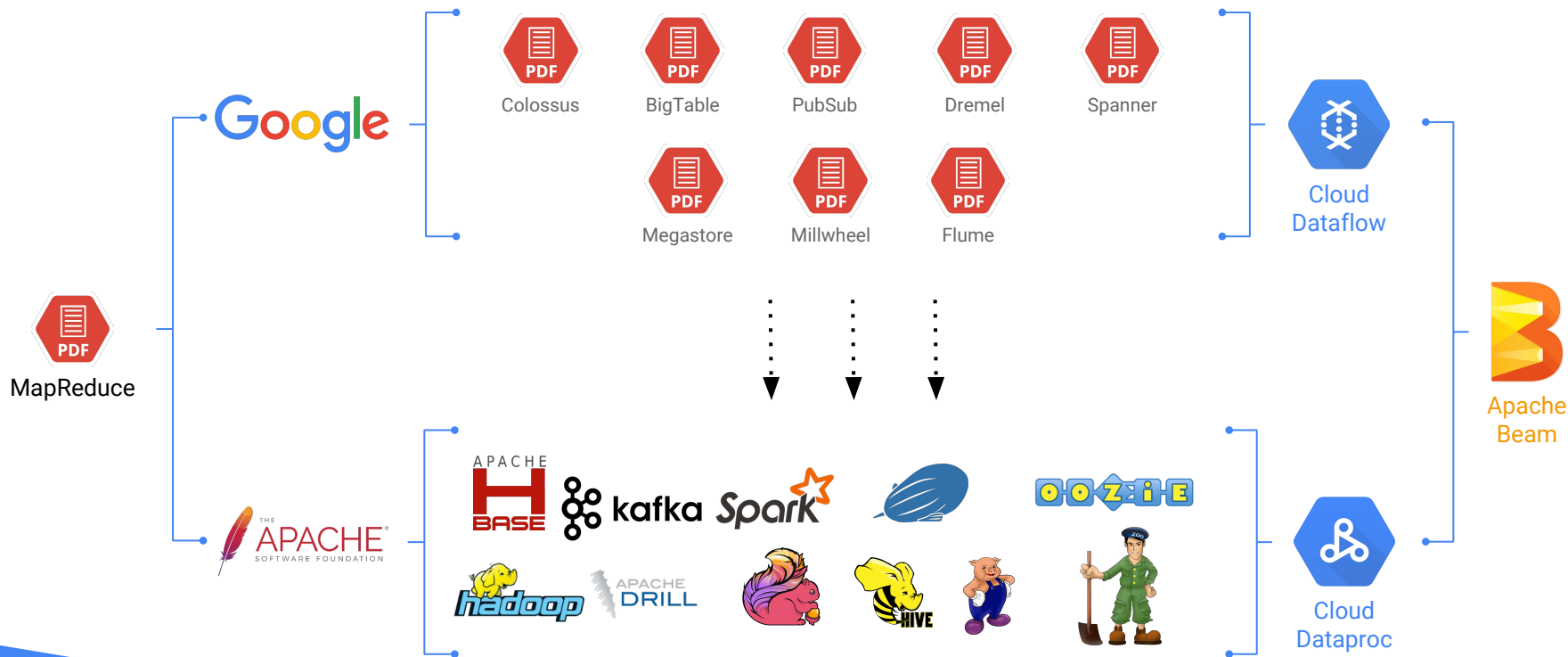You are paying for more (spare) capacity than you actually need to process your data

Google Cloud

# Pay for exactly what you use



You only pay for resources when
you need them

# Open source
## on **Google Cloud**

# Google is passionate about open source

# What is Cloud Dataflow?

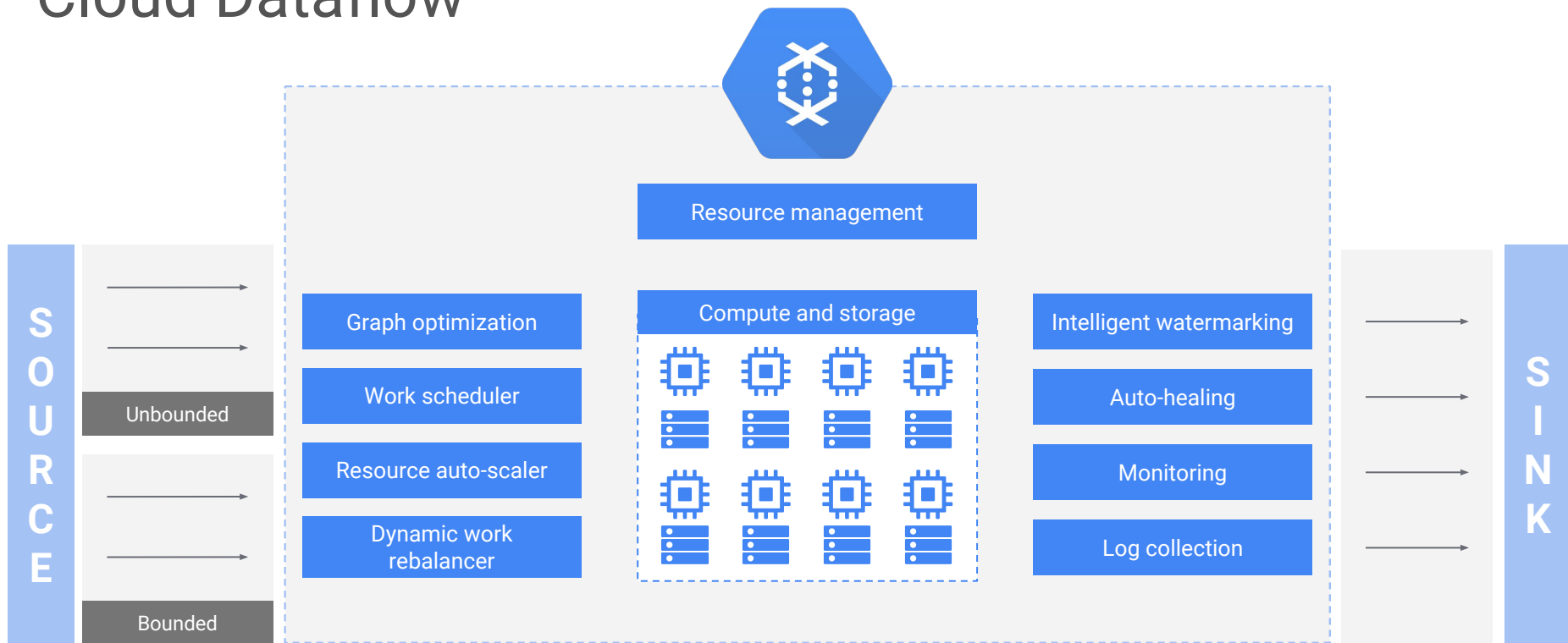- Unified batch and streaming processing

- Fully managed, no-ops data processing

- Open source programming model

- Intelligently scales to millions of QPS

beam

# Cloud Dataflow

# Cloud Dataproc offers a spectrum

**Cloud Dataflow**

**Cloud Dataproc**

**Cloud Dataflow**

Cloud Dataflow is a real-time data processing service for batch and stream data processing.

- Fully managed
- Unified programming model
- Integrated and open source
- Resource management
- Autoscaling
- Monitoring

# What is Cloud Dataproc?

Google Cloud Dataproc is a fast, easy to use, low cost and fully-managed service, powered by Google Cloud Platform, that helps you take advantage of the Spark, Flink, and Hadoop ecosystem.

# Google Cloud Dataproc

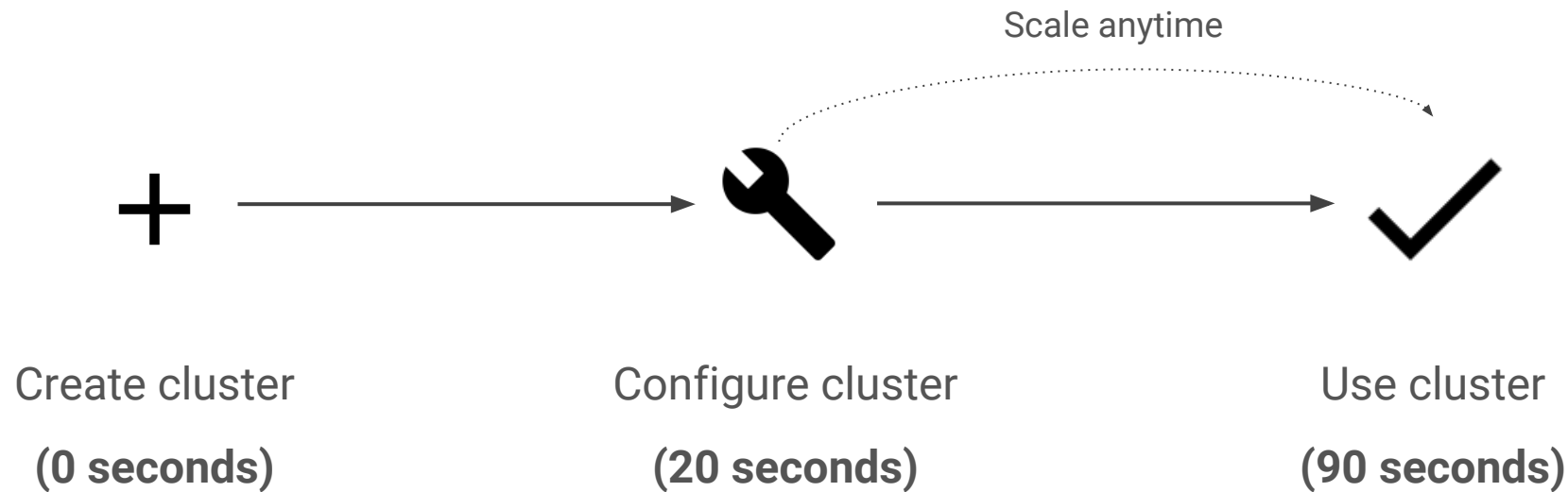## Fast

Things take seconds to minutes, not hours or weeks

## Easy

Be an expert with your data, not your data infrastructure

## Cost-effective

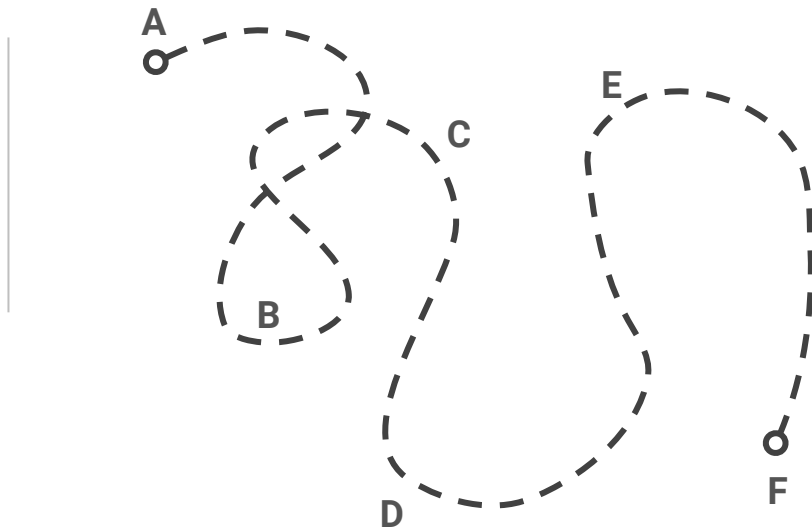Pay for exactly what you use to process your data, not more

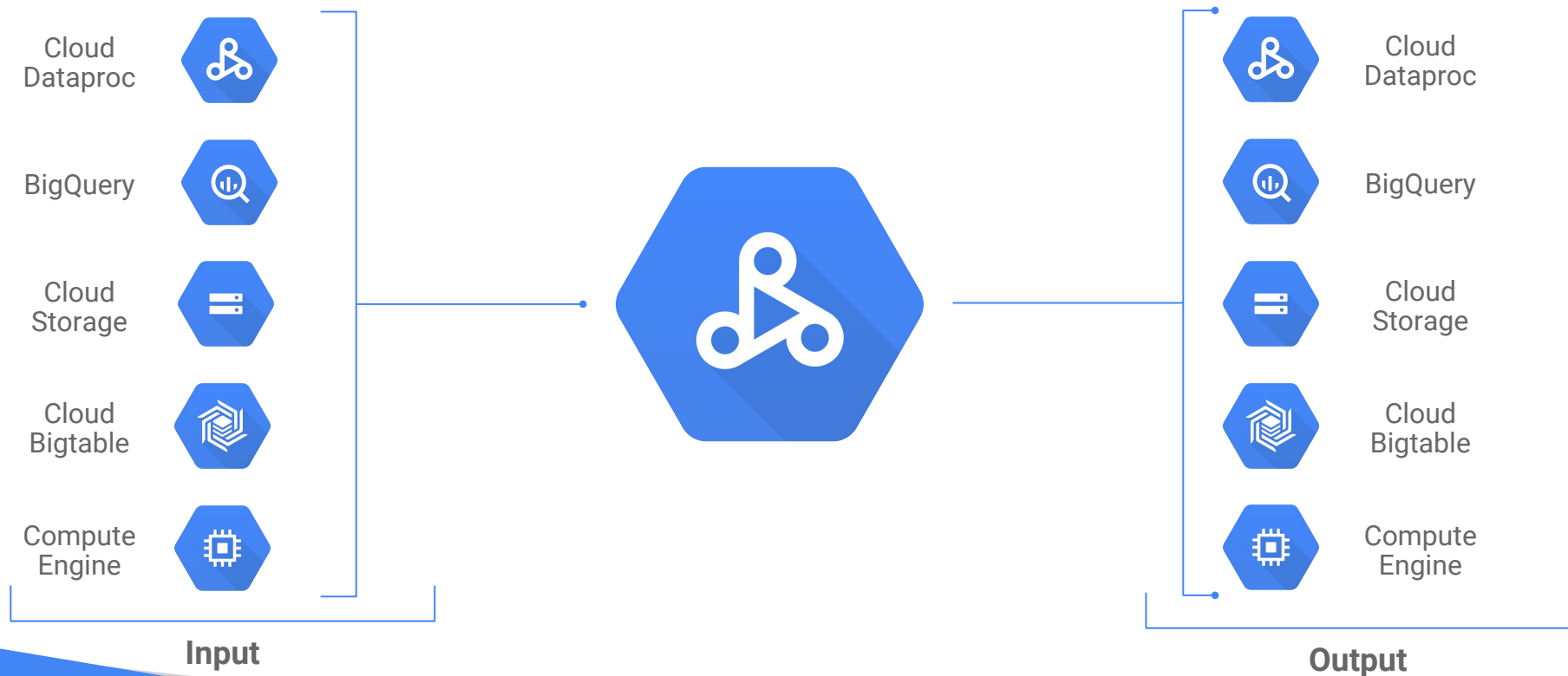Google Cloud

# Cloud Dataproc clusters



Scale anytime

Create cluster

**(0 seconds)**

Configure cluster

**(20 seconds)**

Use cluster

**(90 seconds)**

# Cloud Dataproc offers a spectrum

# Connecting OSS to Cloud Platform



Cloud Dataproc

BigQuery

Cloud Storage

Cloud Bigtable

Compute Engine

**Input**

Cloud Dataproc

BigQuery

Cloud Storage

Cloud Bigtable

Compute Engine

**Output**

Google Cloud

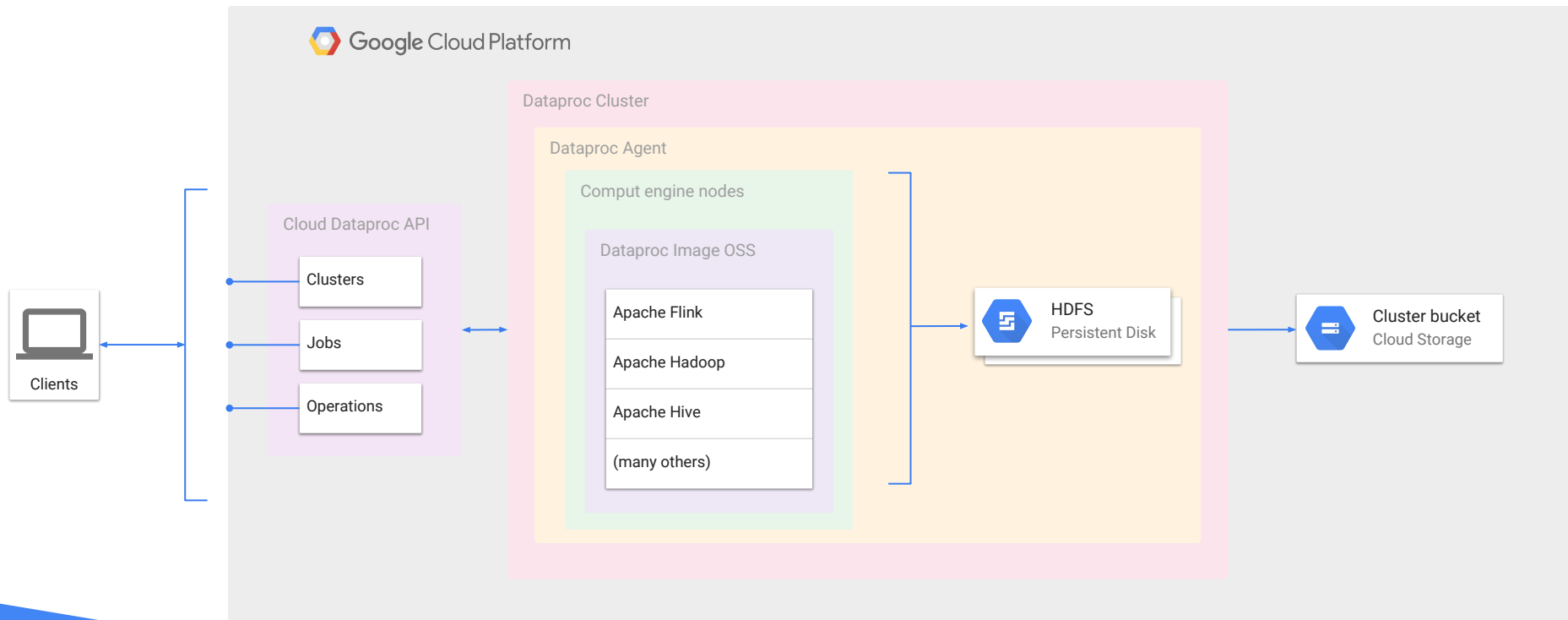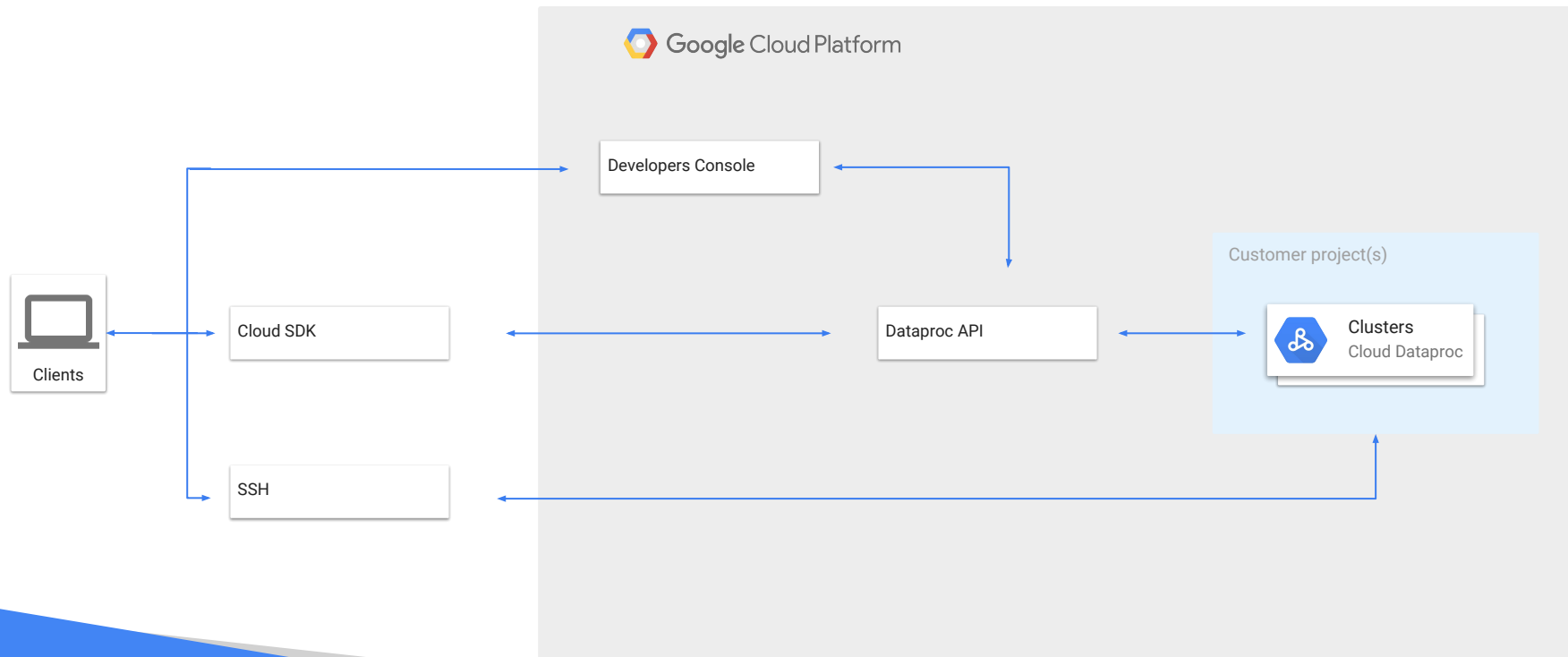# Cloud Dataproc - under the hood

# Cloud Dataproc - under the hood

# Cloud Dataproc - under the hood

# Cloud Storage performance (and improvement)



Legend:
- Old method (light blue)
- BoringSSL-based (dark blue)

Reads
Writes

0    100    200    300

Throughput (mb/sec)

Google Cloud

Cloud Dataproc **demo**

# Example new features and their impact

# Restartable jobs (beta)

- Any job submitted through the Cloud Dataproc Jobs API can now be set to automatically restart on failure

- Very useful for both batch **and** streaming jobs. Jobs which checkpoint can also be automatically restarted

- Specified with the switch `--max_retries_per_hour` when using the Cloud SDK (`gcloud`) the `max_failures_per_hour` in the Jobs API

Google Cloud

# Clusters with GPUs (beta)

- Cloud Dataproc clusters support Compute Engine nodes with Nvidia Tesla K80 GPUs attached to them

- We expect GPU support will continue to grow in the open source data processing ecosystem throughout 2017

- Easily add GPUs to a Cloud Dataproc cluster with the switch `--master/worker_accelerator` with the Cloud SDK (`gcloud`)

Google Cloud

# Single-node clusters (beta)

- Create a sandbox Cloud Dataproc "cluster" with only one node instead of the typical three node design (1 master, 2 workers)

- Great for lightweight data science, small-scale testing, proof of concept building, and education

- Use the `--single-node` argument in the Cloud SDK or select "Single node" when creating a cluster in the Google Cloud Console
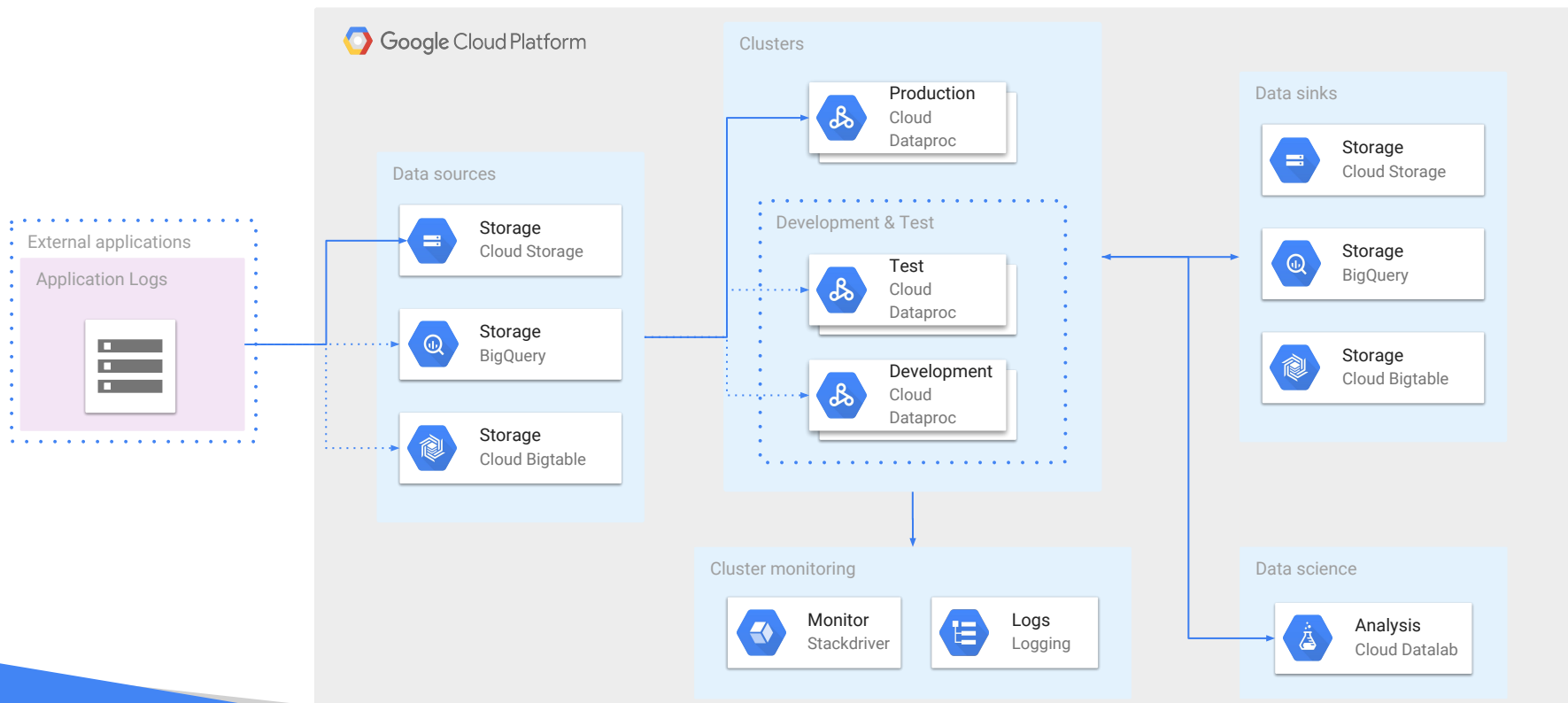
# Regional endpoints & private IP clusters (beta)

- Cloud Dataproc supports a "global" endpoint and "regional endpoints" in each Compute Engine region. This allows you to isolate Cloud Dataproc interactions to one specific region

- Traditionally clusters have needed a public IP attached to them. Cloud Dataproc now supports (easy to setup) "private IP only" clusters which do not require a public IP address on Compute Engine nodes
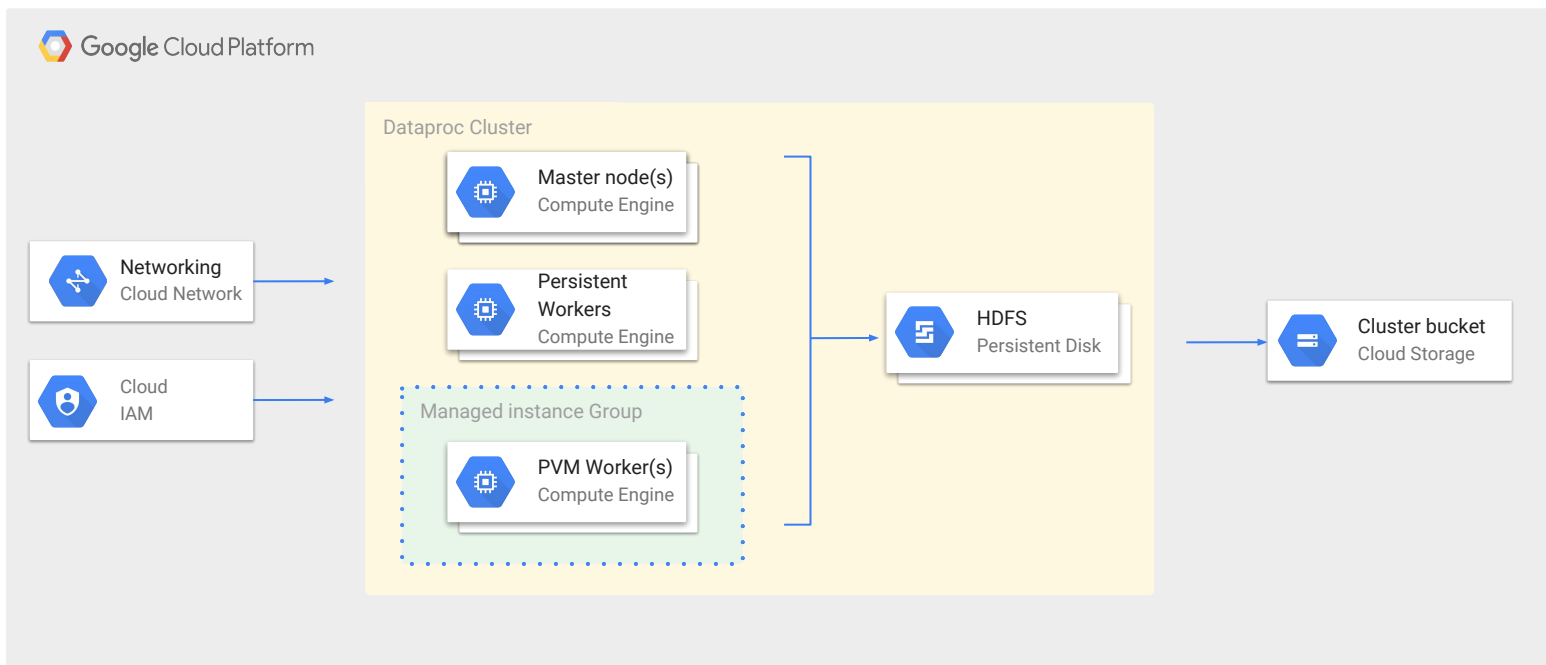
# Cloud platform
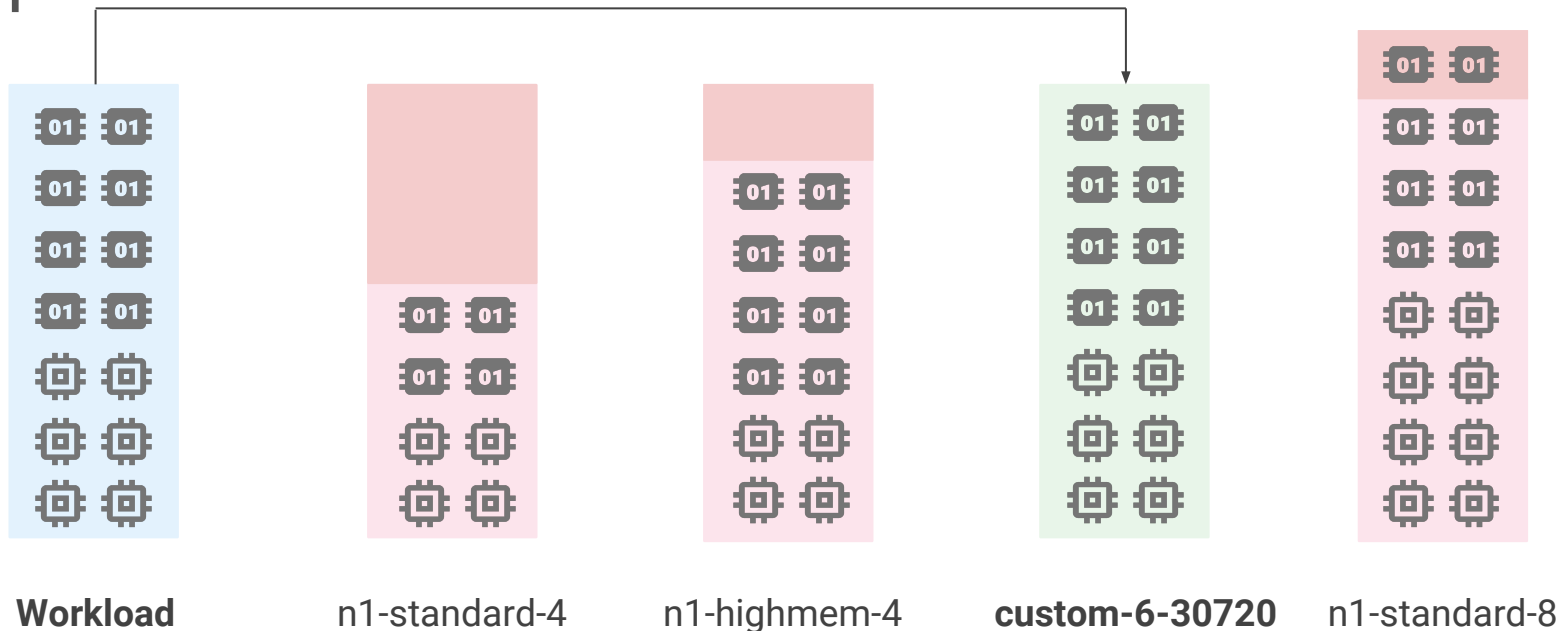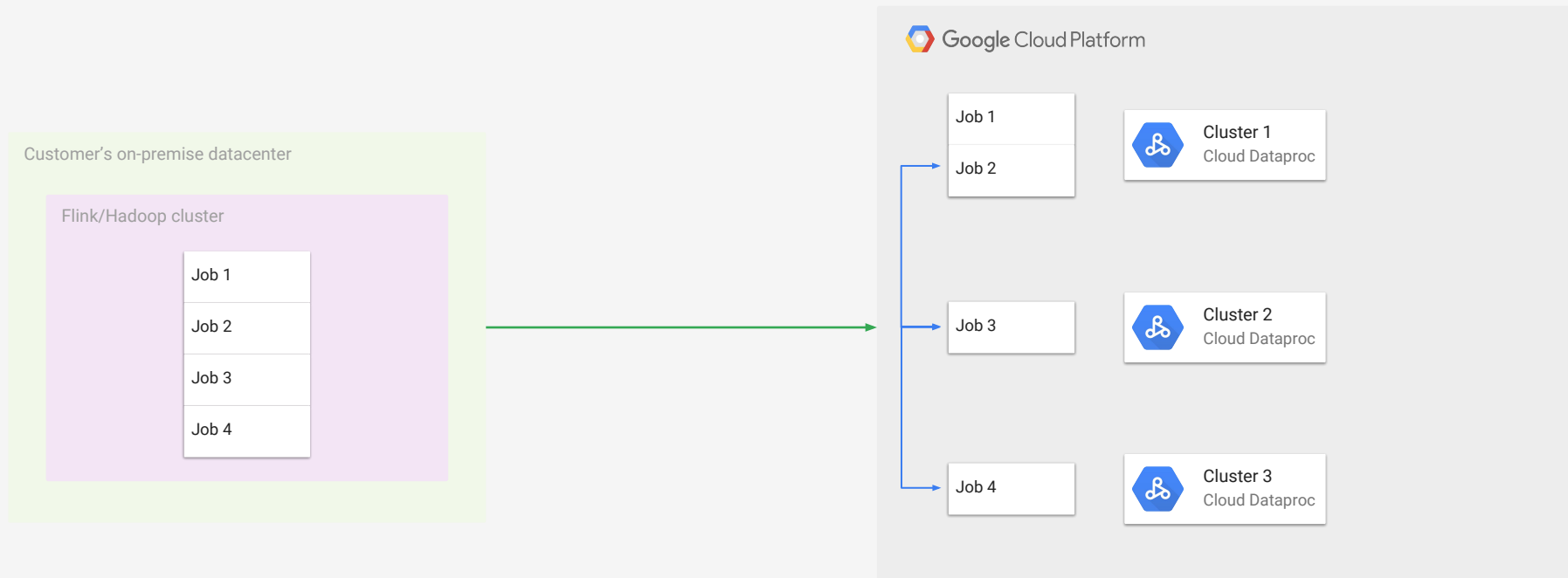**architecture concepts**

# Disaggregation of storage and compute

# Cost savings through preemptible VMs

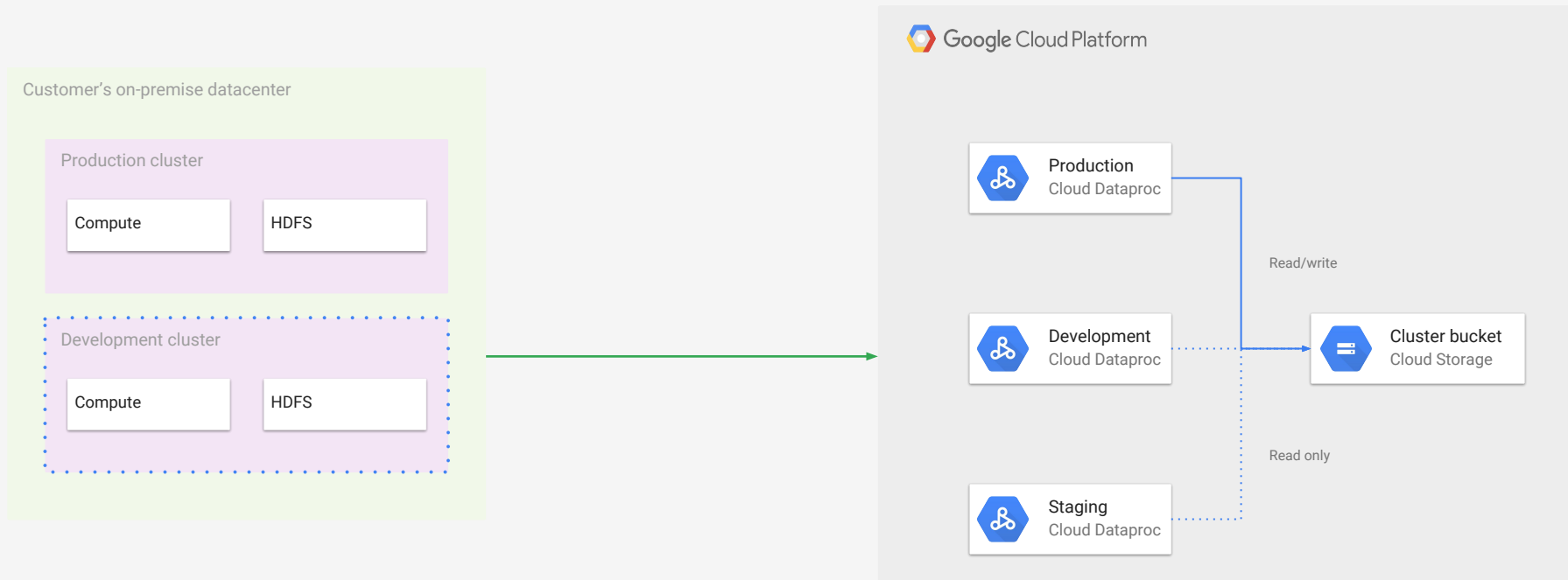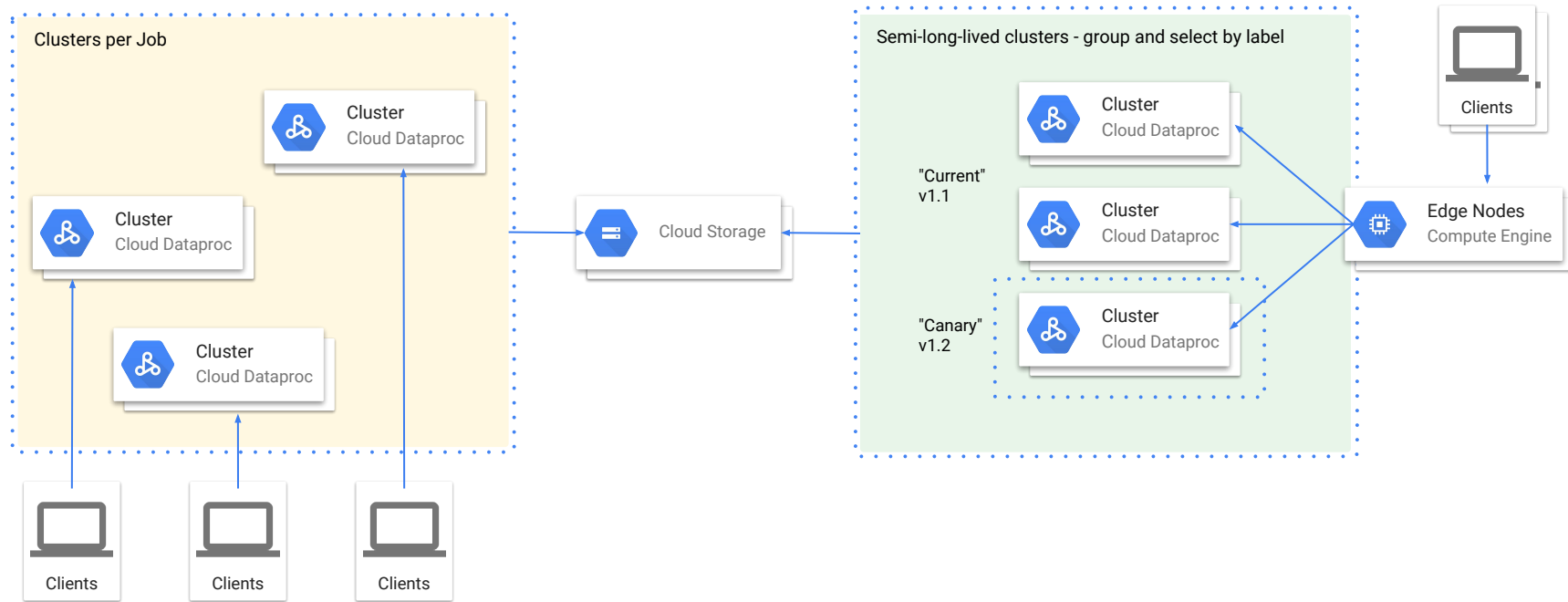# Right-sizing your hardware with custom machine types



**Workload**   n1-standard-4   n1-highmem-4   **custom-6-30720**   n1-standard-8

Google Cloud

# Split clusters and jobs



Customer's on-premise datacenter

Flink/Hadoop cluster

| Job 1 |
| Job 2 |
| Job 3 |
| Job 4 |

Google Cloud Platform

| Job 1 |
| Job 2 |

Cluster 1
Cloud Dataproc

| Job 3 |

Cluster 2
Cloud Dataproc

| Job 4 |

Cluster 3
Cloud Dataproc

Next

Google Cloud

# Separate development and production

# Ephemeral and semi-long-lived clusters

# Questions and **next steps**

# Getting started

**Codelabs** - codelabs.developers.google.com/codelabs/cloud-dataproc-starter/

**Cloud Dataproc quickstarts** - cloud.google.com/dataproc/docs/quickstarts

**Cloud Dataproc tutorials** - cloud.google.com/dataproc/docs/tutorials

**Cloud Dataproc initialization actions** - github.com/GoogleCloudPlatform/dataproc-initialization-actions

Google Cloud

# Getting help

**Cloud Dataproc documentation** - cloud.google.com/dataproc/docs

**Cloud Dataproc release notes** - cloud.google.com/dataproc/docs

**Stack Overflow** - google-cloud-dataproc

**Cloud Dataproc email discussion** - cloud-dataproc-discuss@googlegroups.com

**Google Cloud Support** - cloud.google.com/support

Thank you

**https://goo.gl/Qyf5U7**