# Data Processing Architecture

Google Cloud Platform Fundamentals: Big Data and Machine Learning

Version #1.1

Google Cloud

# Agenda

1. Introduction
2. Fundamentals of GCP
3. Data Analysis on the Cloud
4. Scaling Data Analysis
5. Machine Learning
6. Data Processing Architecture
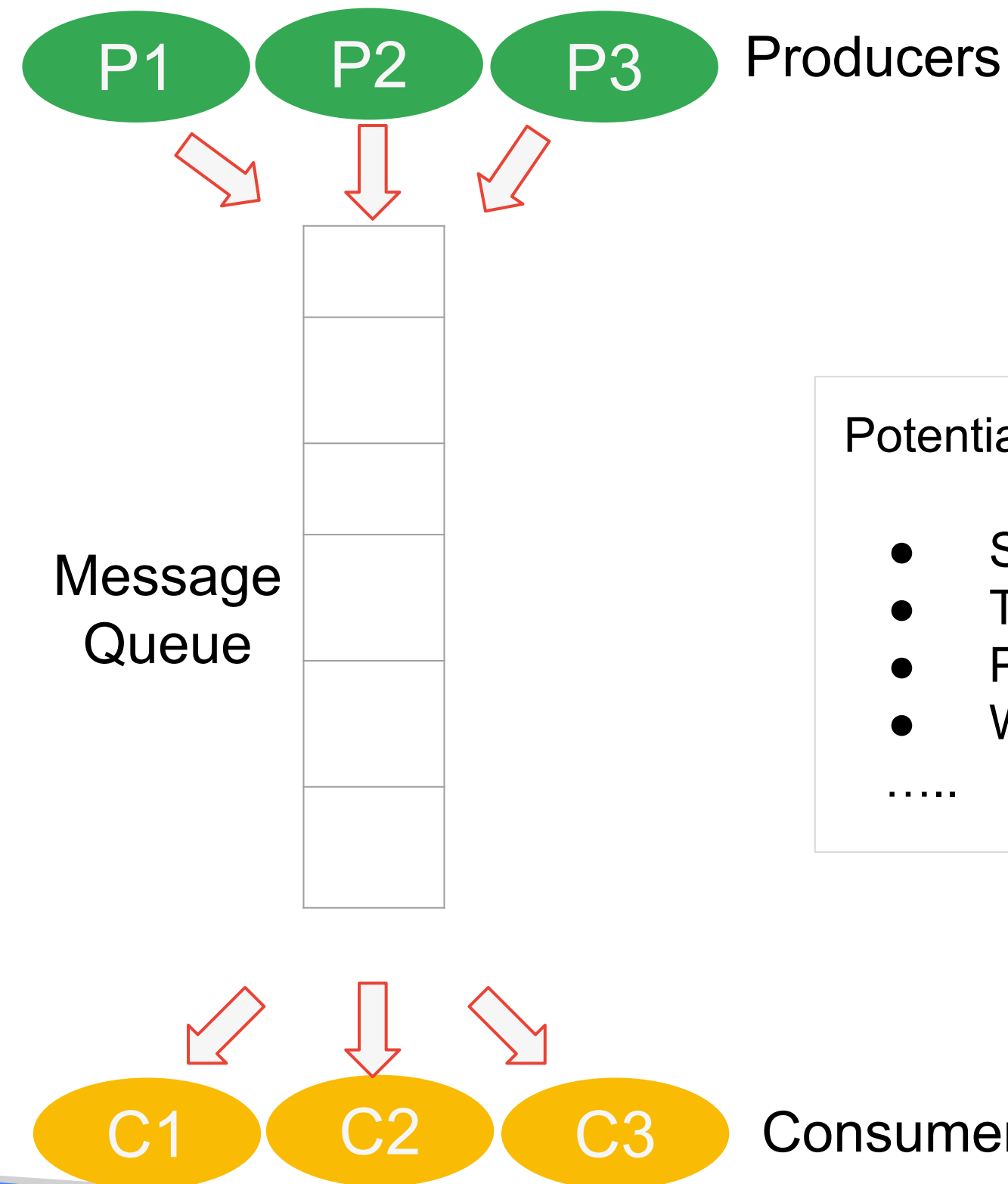
Google Cloud

# Agenda

Message-oriented architectures

Serverless data pipelines

GCP Reference Architecture

Google Cloud

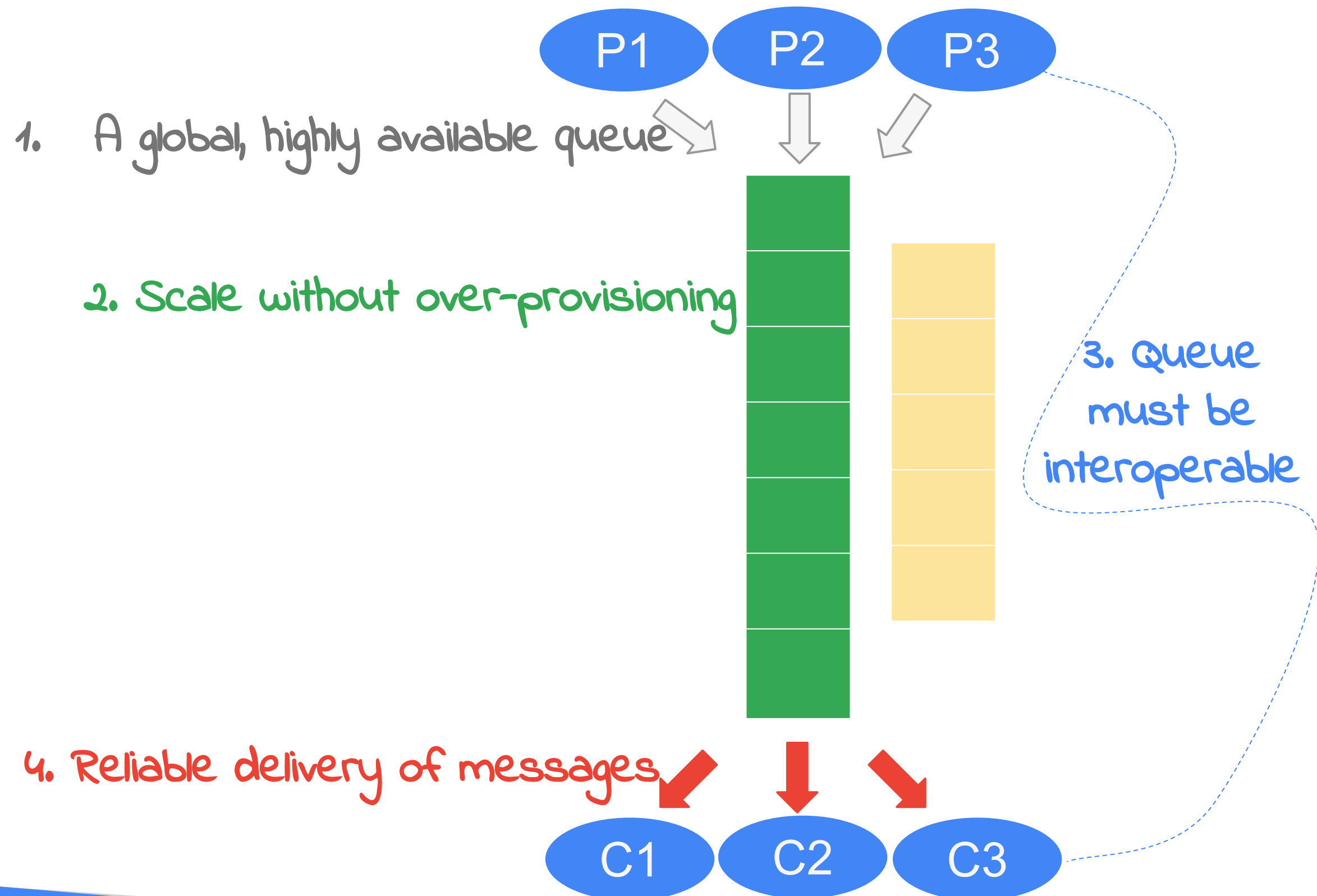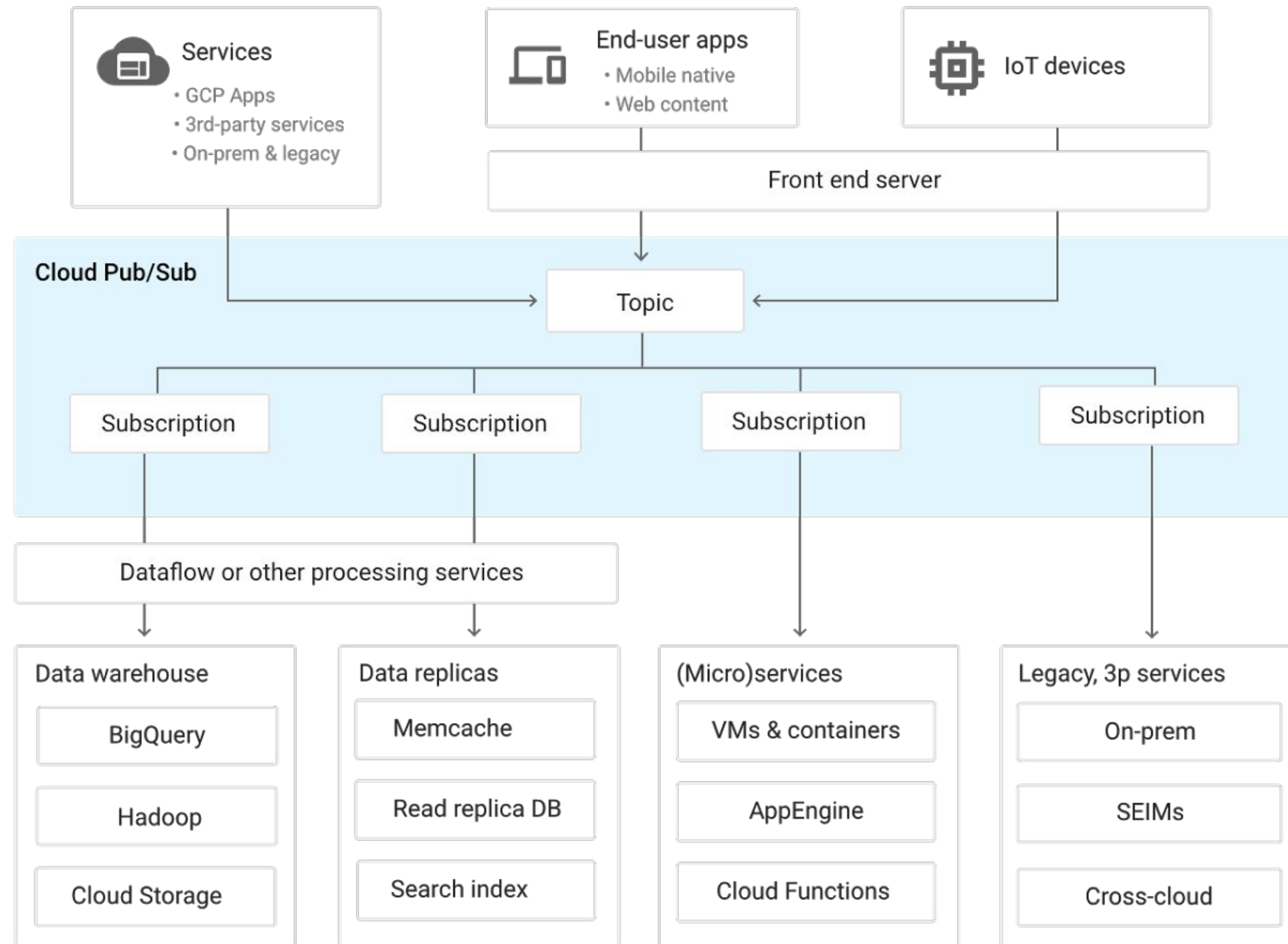# Asynchronous processing is useful for long-lived tasks or to have loose coupling between two systems

P1 P2 P3 **Producers**

**Message Queue**

Potential use cases:

- Send an SMS
- Train ML model
- Process data from multiple sources
- Weekly reports

…..

C1 C2 C3 **Consumers**

Google Cloud

# For robust asynchronous processing, you need:

P1  P2  P3

1. A global, highly available queue

2. Scale without over-provisioning

3. Queue must be interoperable

4. Reliable delivery of messages

C1  C2  C3

Google Cloud

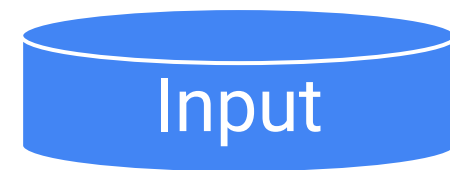# Pub/Sub provides a no-ops, serverless global message queue

# Agenda

Message-oriented architectures

Serverless data pipelines

GCP Reference Architecture

Google Cloud

# Dataflow offers NoOps data pipelines in Java and Python



```python
p = beam.Pipeline(options=options)

lines = p | beam.io.ReadFromText('gs://…')

traffic = lines | beam.Map(parse_data).with_output_types(unicode)

          | beam.Map(get_speedbysensor) # (sensor, speed)

          | beam.GroupByKey() # (sensor, [speed])

          | beam.Map(avg_speed) # (sensor, avgspeed)

          | beam.Map(lambda tup: '%s: %d' % tup))

output = traffic | beam.io.WriteToText('gs://...]')

p.run();
```

Open-source API (Apache Beam) can be executed on Flink, Spark, etc. also

Map

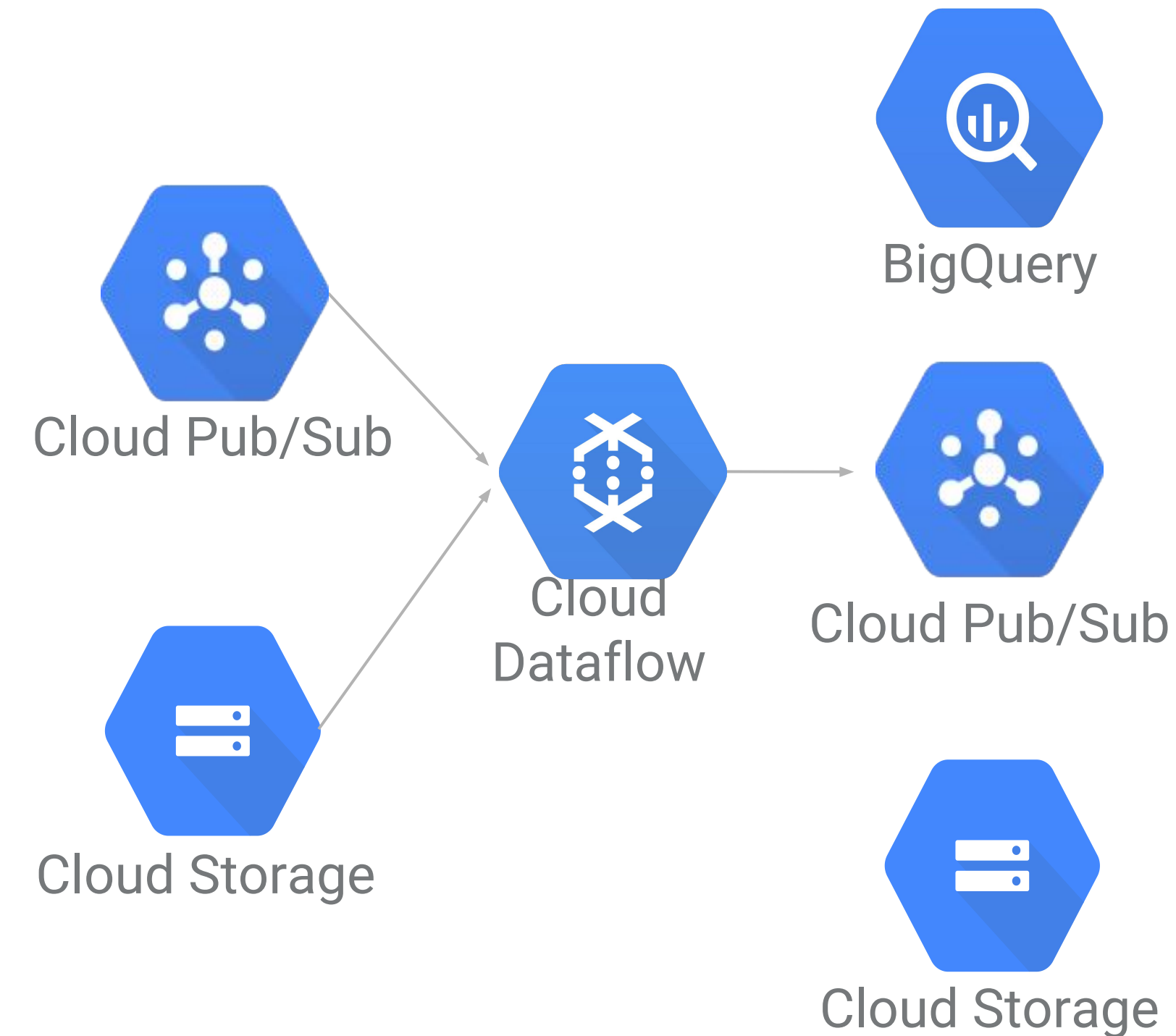Group-By

Reduce

Each of these steps is run in parallel and autoscaled by execution framework

**Input**

**Read**

**Transform 1**

**Transform 2**

**Group**

**Transform 3**

**Transform 4**

**Write**

Output

Google Cloud

# Same code does real-time and batch

BigQuery

Cloud Pub/Sub

Cloud Storage
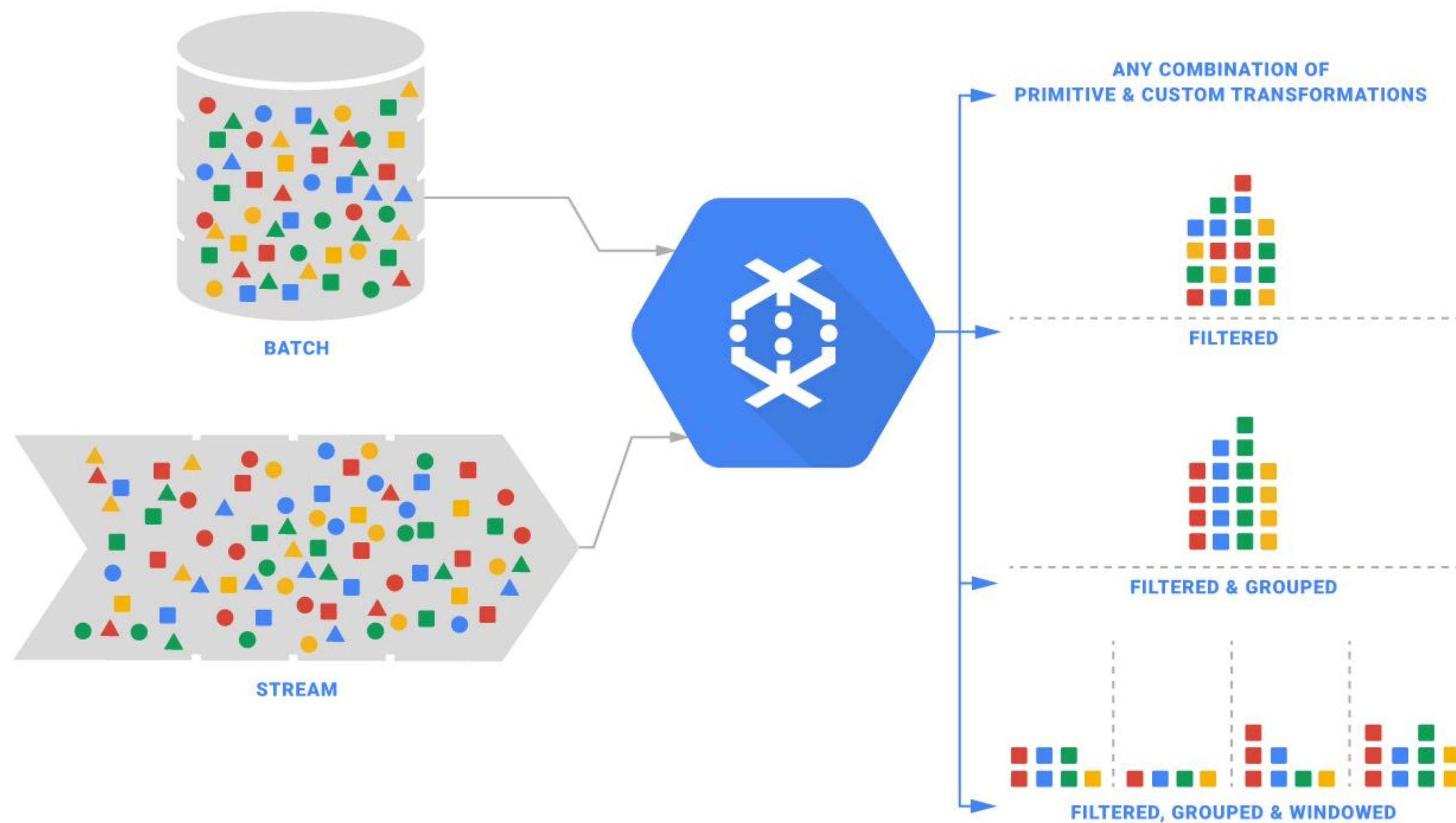
Cloud Dataflow

Cloud Pub/Sub

Cloud Storage

```python
options = PipelineOptions(pipeline_args)
options.view_as(StandardOptions).streaming = True
p = beam.Pipeline(options=options)
    lines = p | beam.io.ReadStringsFromPubSub(input_topic)
    traffic = (lines
        | beam.Map(parse_data).with_output_types(unicode)
        | beam.Map(get_speedbysensor) # (sensor, speed)
        | beam.WindowInto(window.FixedWindows(15, 0))
        | beam.GroupByKey() # (sensor, [speed])
        | beam.Map(avg_speed) # (sensor, avgspeed)
        | beam.Map(lambda tup: '%s: %d' % tup))
    traffic | beam.io.WriteStringsToPubSub(output_topic)

p.run()
```

Google Cloud

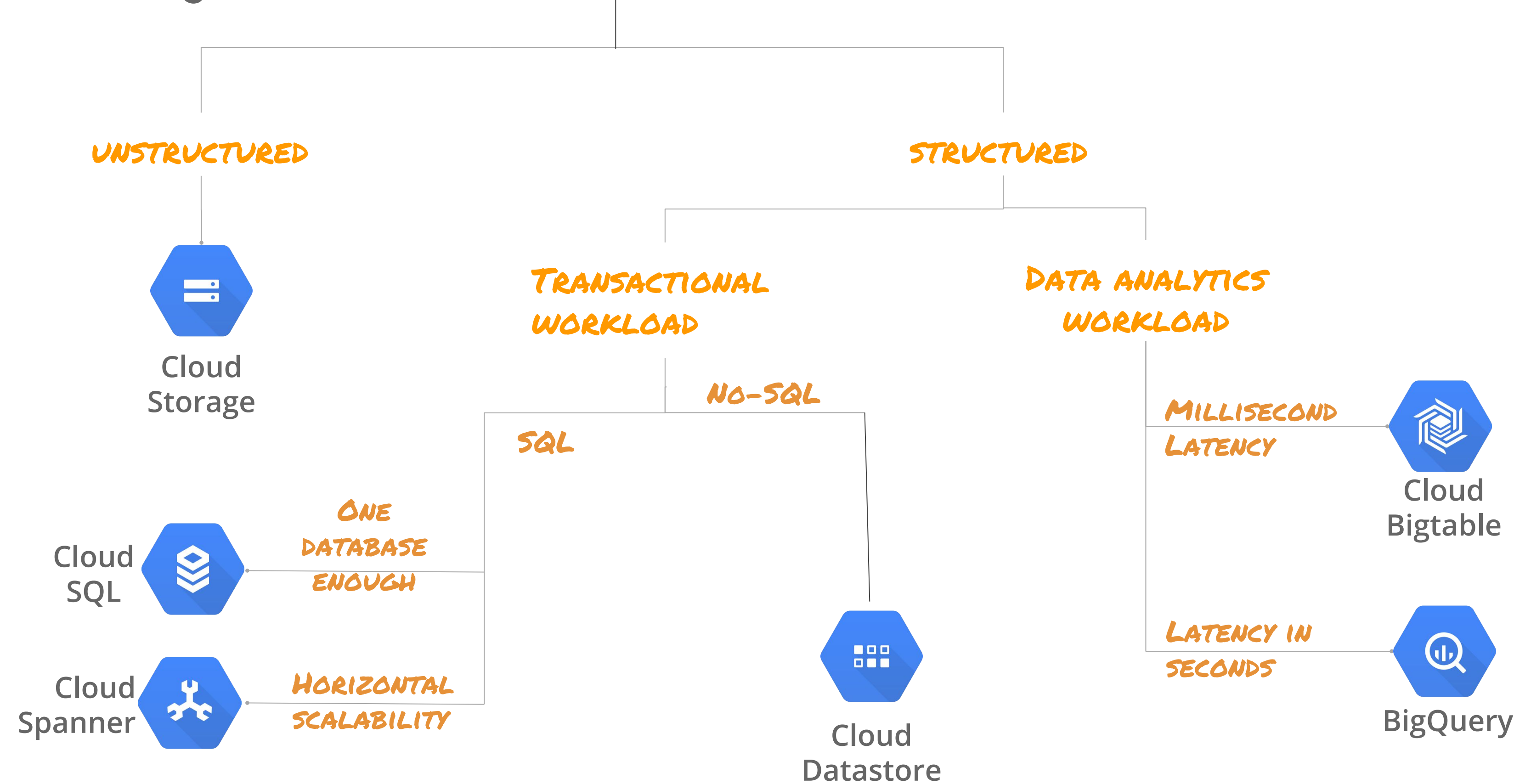# Dataflow does ingest, transform, and load; consider using it instead of Spark

# Agenda

Message-oriented architectures

Serverless data pipelines

GCP Reference Architecture

Google Cloud

# Choosing where to store **data** on GCP



UNSTRUCTURED

STRUCTURED

Cloud Storage

TRANSACTIONAL WORKLOAD

DATA ANALYTICS WORKLOAD

NO-SQL

SQL

MILLISECOND LATENCY

Cloud Bigtable

ONE DATABASE ENOUGH

Cloud SQL

HORIZONTAL SCALABILITY

Cloud Spanner

LATENCY IN SECONDS

BigQuery

Cloud Datastore

Google Cloud

# Run Spark/Hadoop jobs on Cloud Dataproc

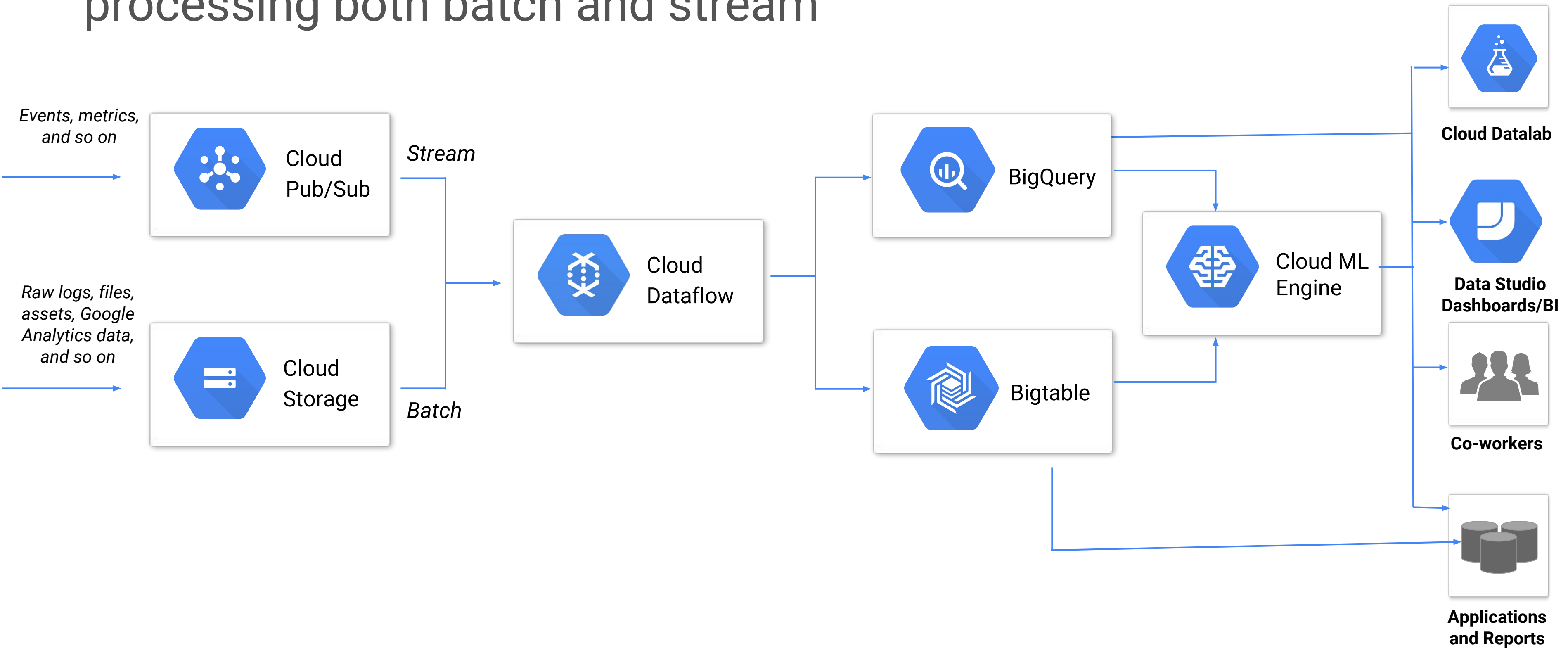# On GCP, you can have the same data processing pipeline for processing both batch and stream

*Events, metrics, and so on*

**Cloud Pub/Sub**

*Stream*

*Raw logs, files, assets, Google Analytics data, and so on*

**Cloud Storage**

*Batch*

**Cloud Dataflow**

**BigQuery**

**Bigtable**

**Cloud ML Engine**

**Cloud Datalab**

**Data Studio Dashboards/BI**

**Co-workers**

**Applications and Reports**

Google Cloud

# Module Review

Google Cloud

# Module review

Match the use case on the left with the product
on the right

A.  Decoupling producers and consumers of
    data in large organizations and complex
    systems

B.  Scalable, fault-tolerant multi-step processing
    of data

1.  Cloud Dataflow

2.  Cloud Pub/Sub

Google Cloud

# Module review answers

Match the use case on the left with the product
on the right

A.   Decoupling producers and consumers of
     data in large organizations and complex
     systems

B.   Scalable, fault-tolerant multi-step processing
     of data

1.   Cloud Dataflow

2.   Cloud Pub/Sub

Google Cloud

# Resources (1 of 2)

| | |
|---|---|
| Cloud Pub/Sub | https://cloud.google.com/pubsub/ |
| Cloud Dataflow | https://cloud.google.com/dataflow/ |
| Processing media using Cloud Pub/Sub and Compute Engine | https://cloud.google.com/solutions/media-processing-pub-sub-compute-engine |

Google Cloud

# Resources (2 of 2)

| | |
|---|---|
| Reverse Geocoding of Geolocation Telemetry in the Cloud Using the Maps API | https://cloud.google.com/solutions/reverse-geocoding-geolocation-telemetry-cloud-maps-api |
| Using Cloud Pub/Sub for Long-running Tasks | https://cloud.google.com/solutions/using-cloud-pub-sub-long-running-tasks |

Google Cloud

cloud.google.com

Images by Connie Zhou