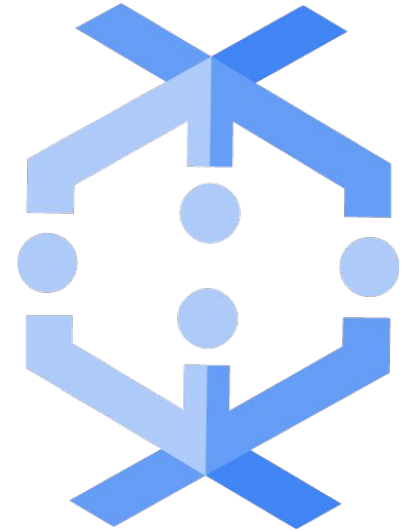


# Introduction

- Dataflow
- Dataproc
- Dataflow VS. Dataproc



# Dataflow

# Dataflow

Dataflow is a Google Cloud service that provides ***unified stream*** and ***batch*** data processing at scale

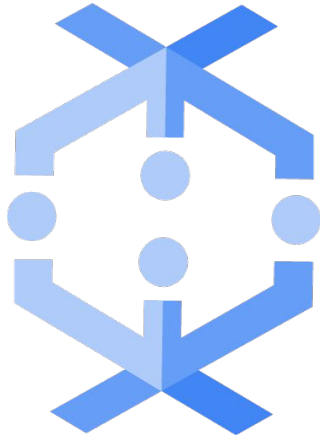


Cloud  
DataFlow



beam

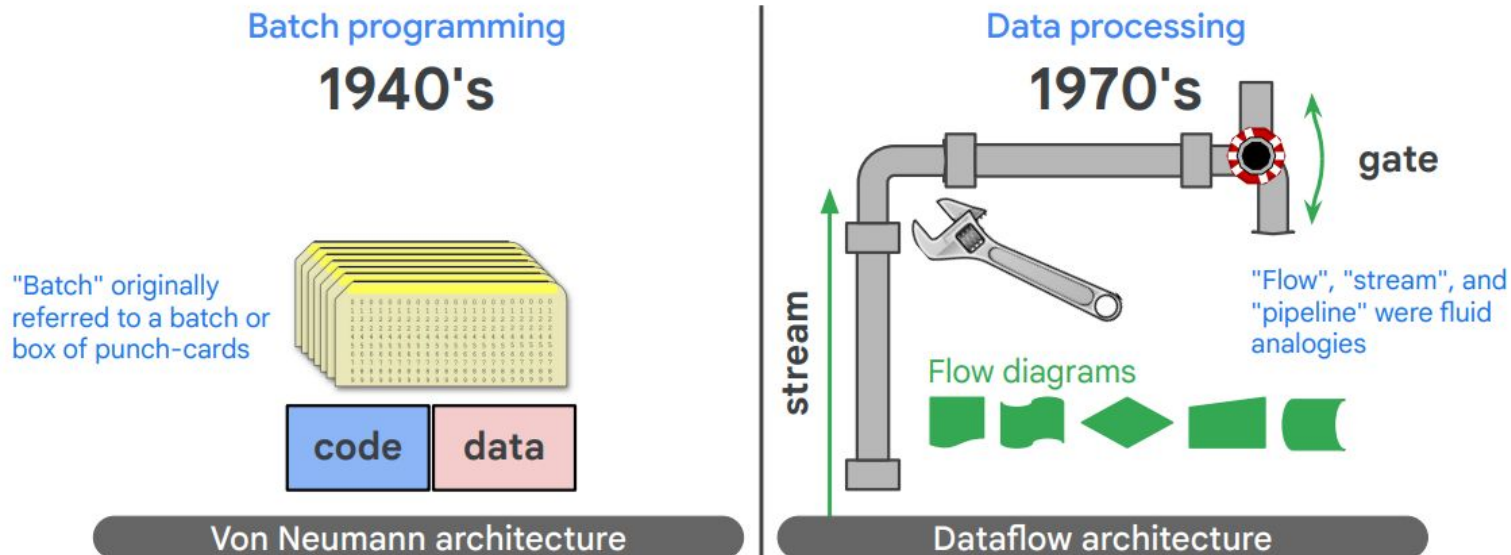
# Dataflow



Qualities that Dataflow contributes to data engineering solutions:

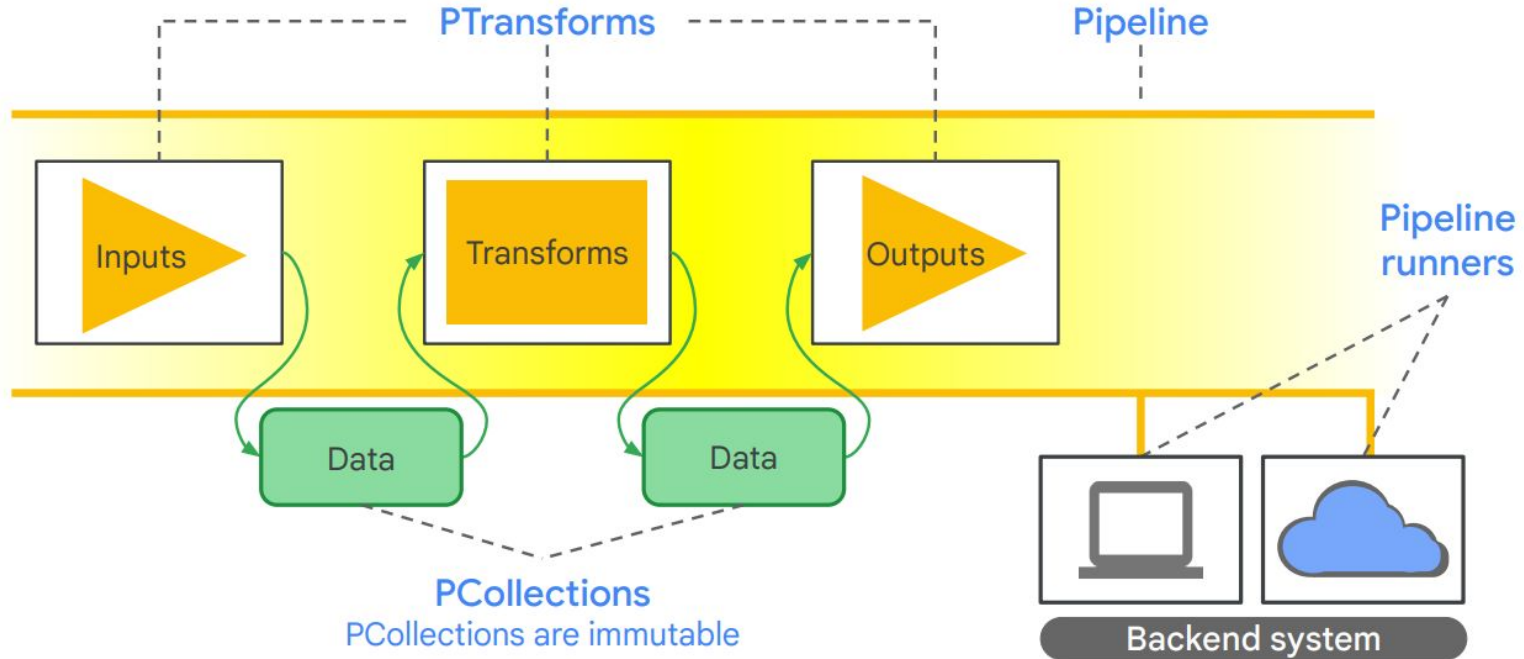
- Scalability
- Low latency

# Batch programming and data processing used to be two very separate and different things

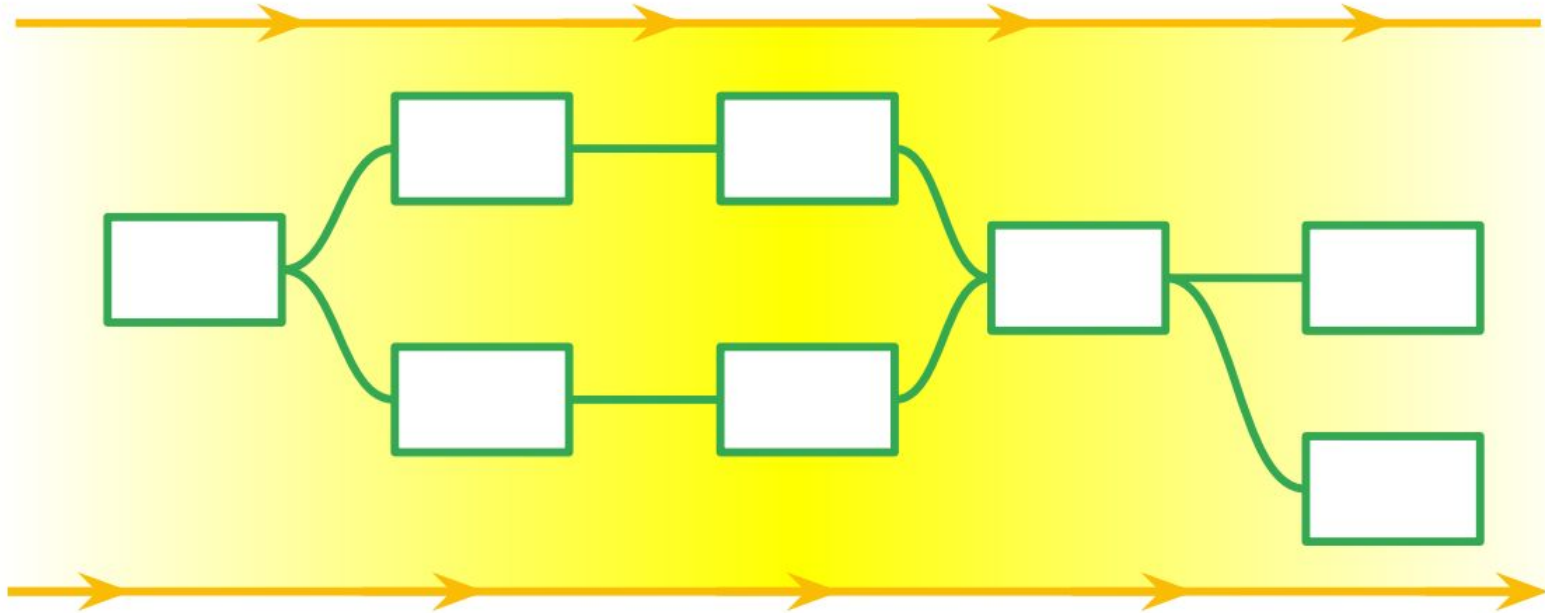


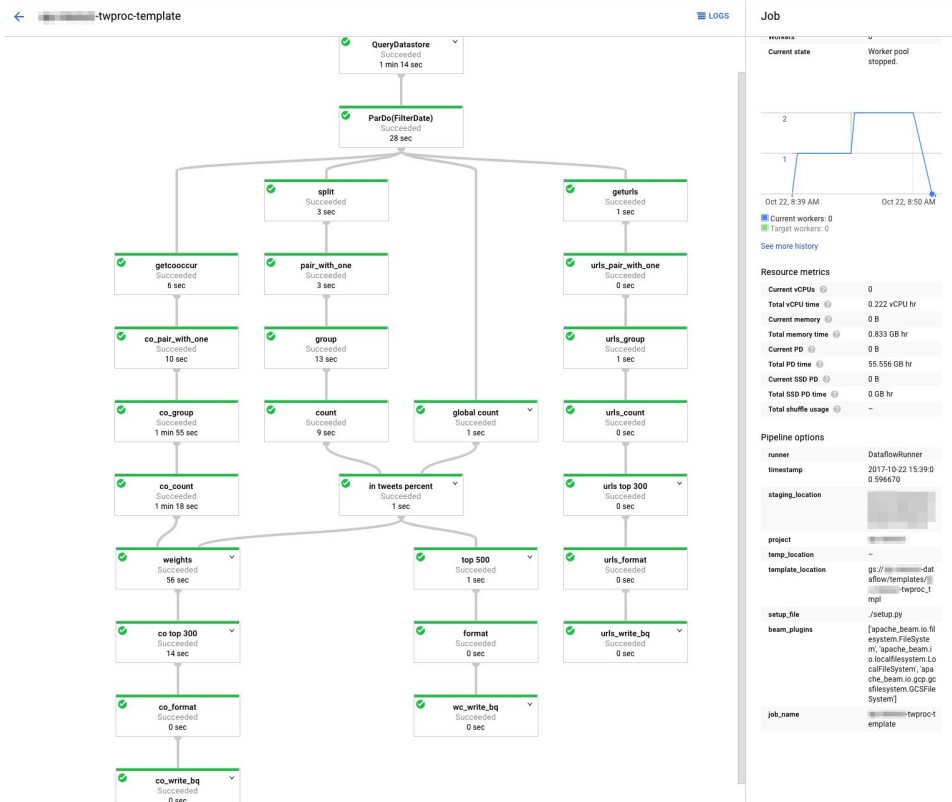
*Different tools, different platforms, different concepts, different methods.*

# Apache BEAM = Batch + strEAM



**A Dataflow pipeline is a directed graph of steps**



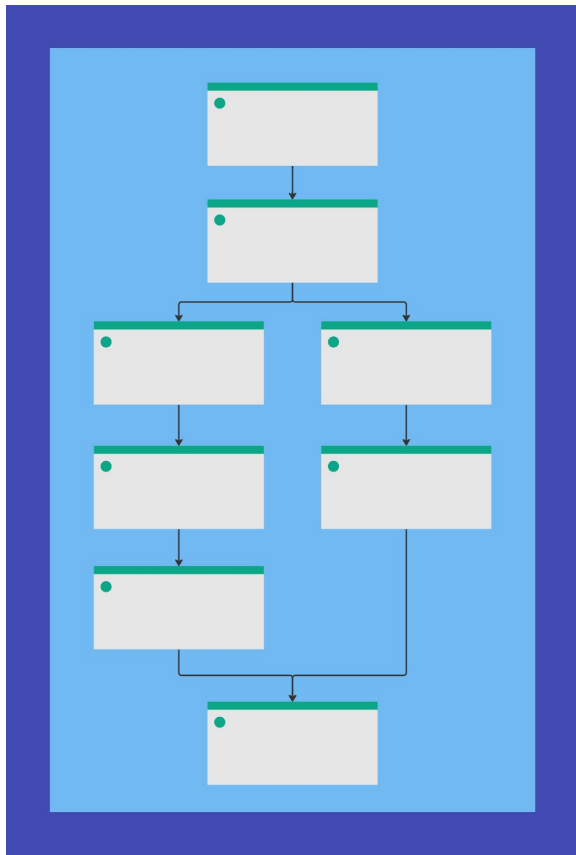




# Dataflow templates

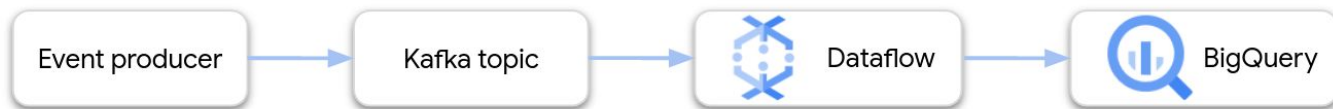
Dataflow templates allow you to package a Dataflow pipeline for deployment.

You can create your own *custom Dataflow templates*, and Google provides *pre-built templates* for common scenarios.

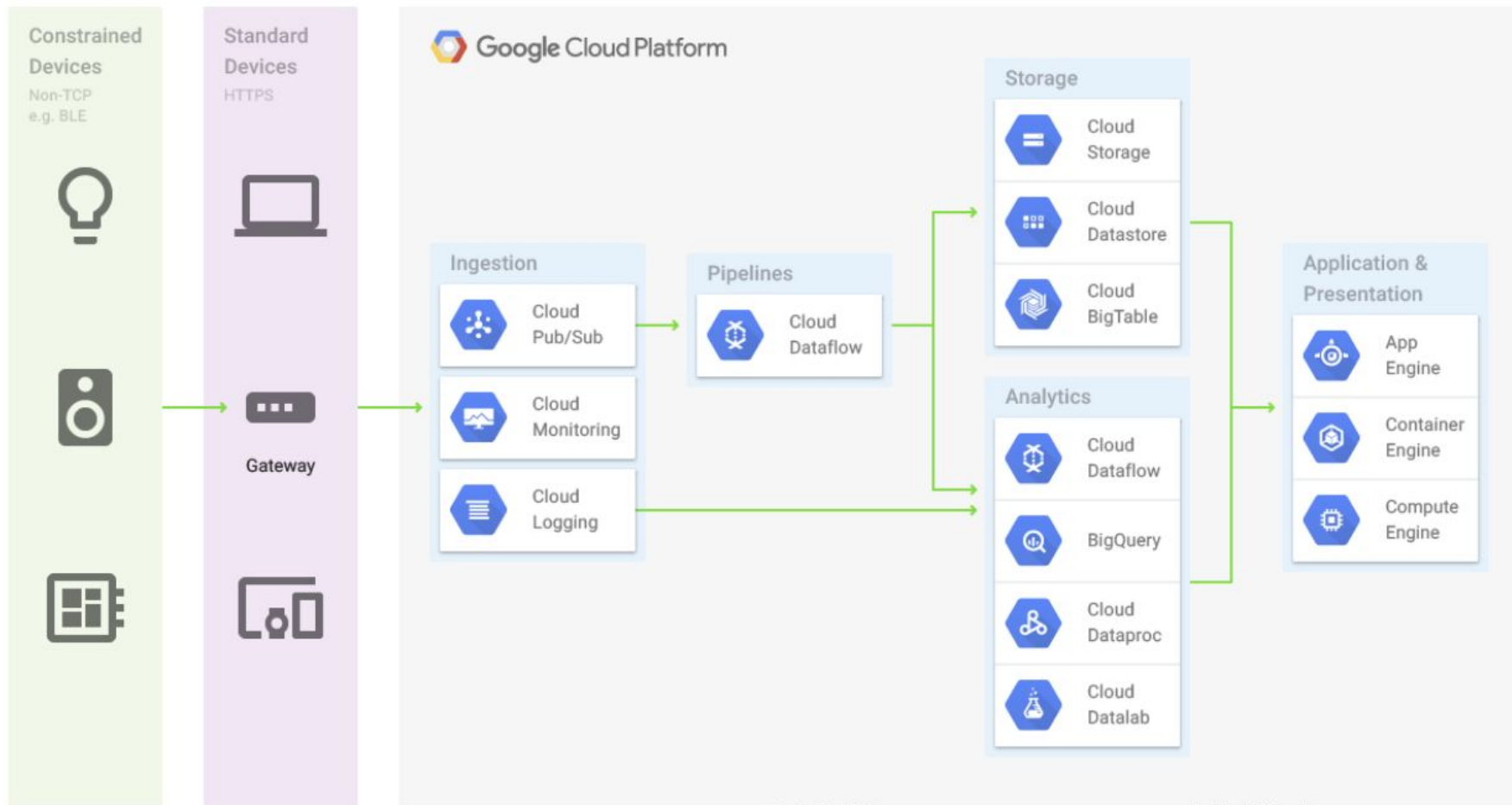


# Dataflow templates Pre-built templates

- Apache Kafka to BigQuery
- Change Data Capture from MySQL to BigQuery (Stream)
- MongoDB to BigQuery (CDC)
- Pub/Sub to Pub/Sub

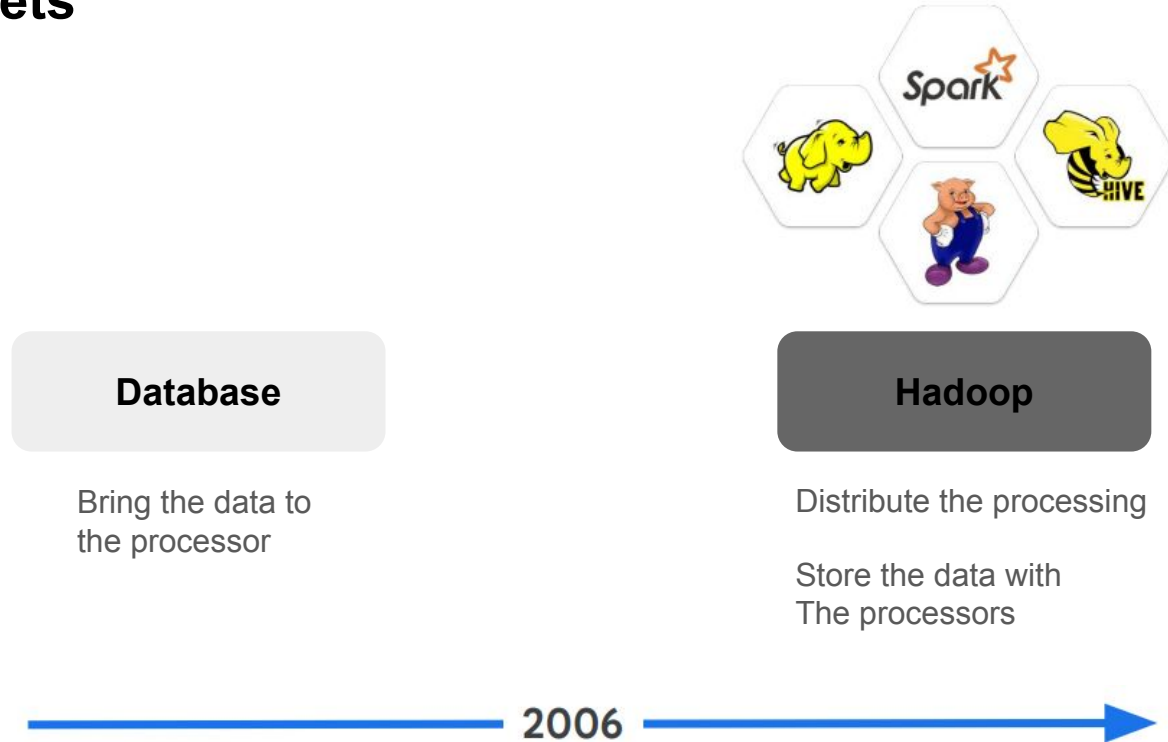


# Reference Architecture (Real Time Streaming IoT)

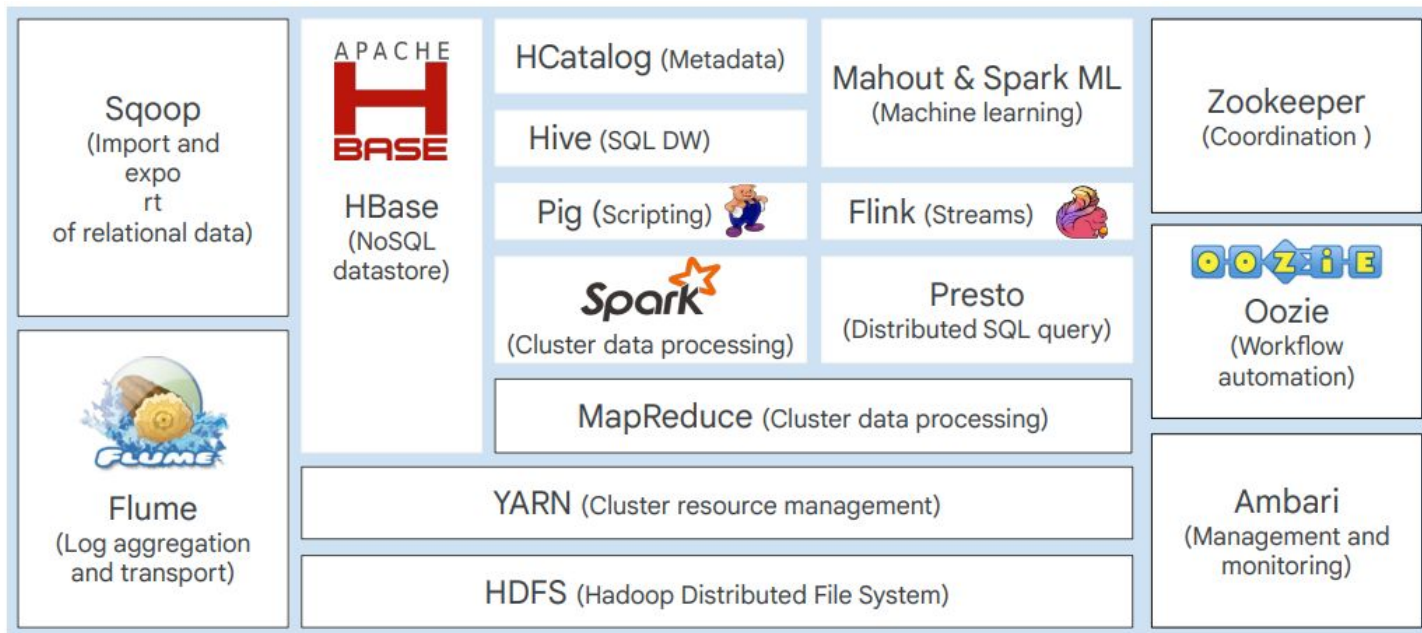


# Dataprocc

# The Hadoop ecosystem developed because of a need to analyze large datasets



# The Hadoop ecosystem is very popular for Big Data workloads



**Apache Spark is a popular, flexible, powerful way to process large datasets**

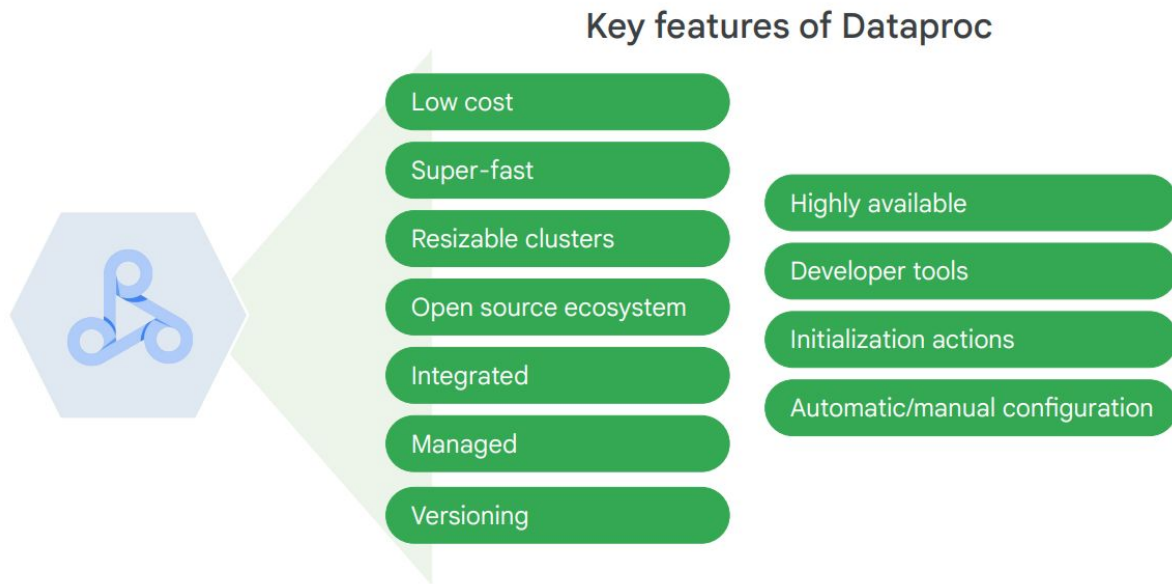


etc.



[spark.apache.org](https://spark.apache.org)

# Dataproc is a managed service for running Hadoop and Spark data processing workloads







# There are other OSS options available in Dataproc

- Spark (default)
- Pig (default)
- Kafka
- Presto
- Jupyter
- IPython
- Much more...
- Hive (default)
- Zeppelin
- Hue
- Anaconda
- Apache Flink
- Oozie
- HDFS (default)
- Zookeeper
- Tez
- Cloud SQL Proxy
- Datalab
- Sqoop

# Comparison between Dataproc & Dataflow

# Dataflow VS. Dataproc

	 Dataflow	 Dataproc
Recommended for:	New data processing pipelines, unified batch and streaming	Existing Hadoop/Spark applications, machine learning/data science ecosystem, large-batch jobs, preemptible VMs
Serverless:	Yes	No
Auto-scaling:	Yes, transform-by-transform (adaptive)	Yes, based on cluster utilization (reactive)
Expertise:	Apache Beam	Hadoop, Hive, Pig, Apache Big Data ecosystem, Spark, Flink, Presto, Druid

# Choosing between Dataflow and Dataproc

