

Section 18: overview of Bigdata Services

Cloud Pub/Sub

- ↳ Messaging Service
- ↳ applications can write/read messages in an asynchronous way.

Topics are structures for storing messages

- ↳ applications subscribe to topics.

Setup Pub Sub

create subscription to topic

- ↳ Delivery type
 - push vs pull

- ↳ expiration delay

- ↳ acknowledgement deadline

Can save a Message state by making a snapshot

Cloud dataproc

Manages Hadoop clusters

Setup

- ↳ 3 Modes

- Single Node 1 master 0 worker
- Standalone Node 1 Master N workers
- High availability 3 Masters N workers

- ↳ staging Bucket used for data store

Cloud Dataflow

- ↳ Based on apache beam python or Java

- ↳ Template based Jobs

- ↳ Manages clusters automatically

Cloud Transfer

copy large amount of data to buckets

From webaccessible Source

BigQuery

- Large volumes of data
- Data Warehousing
- Data Transfer Service

BigQuery ML

- Create Models in SQL

- ↳ Basic type Models

- ↳ Tensorflow

- ↳ autoML

- BigQuery can do data engineering without exporting data
- Different limited options
- evaluation can be done in SQL too

BigQuery IAM

Scopes

- Datasets
- operations
- Jobs
- Tables
- ML Models

Policy can be set through bq update --cmd.

CloudCo-poser

- ↳ apache airflow

- ↳ Manages workflows for data engineering

Workflows

- ↳ Tasks with dependencies

- ↳ Python Scripts based

- ↳ can use python dependencies

Architecture

- deployed in environments

Tenant Project

- ↳ Managed by google

Resources

cloud sql - Airflow Metadata

App engine - Airflow webserver

Customer project

Resources

- Cloud Storage
- Kubernetes engine - airflow deploy
- Redis - Message broker
 - ↳ persistent cache

. Cloud data fusion

- * ETL deployment/development

- ↳ Code-Free

- * Versions

- Basic
- enterprise

- * Minimal Set-up