# Data Pipeline &
# Data Lake
## (in Google Cloud)

Megazone Google Cloud Team

Jung hoo Park

Google Cloud

Cloud Innovator
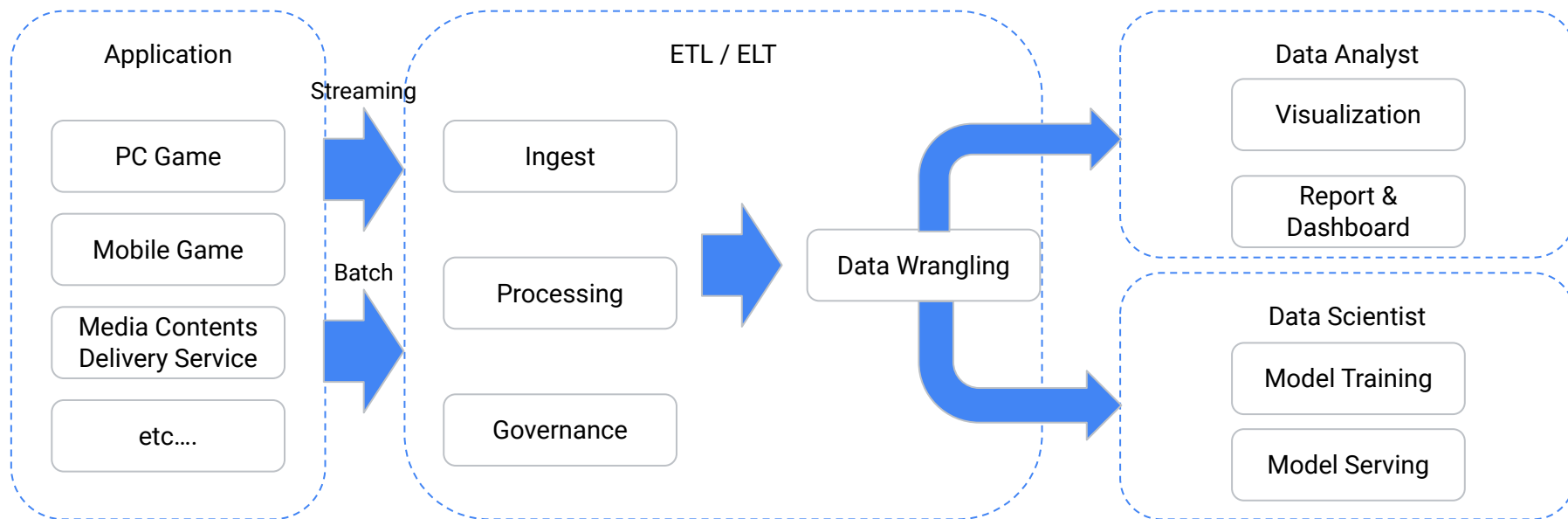MEGAZONE

# Agenda

- **What is a Data Pipeline?**

- **Data Engineering in Google Cloud**

Google Cloud

# What is a Data Pipeline?

01

# Conceptual Data Pipeline

DATA & AI LANDSCAPE 2019

July 16, 2019 - FINAL 2019 VERSION

© Matt Turck (@mattturck), Lisa Xu (@lisaxu92), & FirstMark (@firstmarkcap)    mattturck.com/data2019

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

# 일반적인 데이터 처리 환경 프로비저닝

스케일

재설정

서버 구매 → 서버설정 → OS 인스톨 → OS 설정 → OS 최적화 → OS 디버그 → 프로비저
닝

Google Cloud

Cloud Innovator
MEGAZONE

# Google Cloud
# Service for
# Data Pipeline

Google Cloud

02

| Ingest | Process | Store | Analyze | Visualize |
|--------|---------|-------|---------|-----------|
| Stackdriver Logging | Cloud Dataflow | Cloud Storage | BigQuery | Data Studio |
| BigQuery Streaming API | Cloud Dataproc | BigQuery | | |
| Cloud Pub/Sub | Cloud Composer | Cloud SQL | Cloud Dataprep | Cloud Datalab |
| Cloud IoT Core | | Cloud Datastore | | |
| Transfer Appliance | Cloud Dataprep | Cloud Bigtable | | |
| Transfer Service | Cloud Datalab | Cloud Spanner | Cloud Datalab | Google Cloud Partner 3rd Party |

Google Cloud

# Ingest　　Process　　Store　　Analyze　　Visualize

**Cloud Dataflow**

**Cloud Pub/Sub**

**Cloud Composer**

**Transfer Service**

**Cloud Storage**

**BigQuery**

**BigQuery**

Google Cloud

Ingest

**Cloud
Pub/Sub**

- Exactly-once processing
- Global by default
- No provisioning,
auto-everything
- Seek and replay

Google Cloud

**Storage Transfer Service**

Google Cloud

- Transfer data from cloud to cloud
- Transfer data from bucket to bucket
- Centralized job management
- High-performance copies
- Data security

Cloud Innovator
MEGAZONE

# Processing

beam

Cloud
Dataflow

Auto-Scaling

Streaming Engine

Dataflow Shuffle

Dataflow SQL

Dataflow Template

Inline Monitoring

Google Cloud

Cloud Innovator
MEGAZONE

Apache Airflow

Cloud
Composer

Open Source

Multi-Cloud

Hybrid

Integration

Python Language

Fully-Managed

Google Cloud

Cloud Innovator
MEGAZONE

Cloud
Storage

- Object Life-cycle Management
- Object Version Management
- Retention Policies
- Pub/Sub Notifications for Cloud Storage
- Customer-managed encryption keys
- Cloud Audit Logs with Cloud Storage

Google Cloud

Cloud Innovator
MEGAZONE

# Google Cloud Storage Class

High-performance object storage

Backup and archival storage

| High frequency access | Less frequent access | Low frequency access | Lowest frequency access |
|---|---|---|---|
| **Standard** | **Nearline** | **Coldline** | **Archive** |
| Most projects start with our Standard class of storage, which is optimized for performance and high frequency access. | Our Nearline class of storage is fast, highly durable storage for data accessed less than once a month. | Our Coldline class of storage is fast, highly durable storage for data accessed less than once a quarter. | Our Archive class of storage is designed for cost-effective, long-term preservation of data accessed less than once a year. |

A single API for all storage classes

Google Cloud

**BigQuery**

- Serverless
- Petabyte Scale
- Data Transfer Service
- Foundation for AI & BI
- Big data ecosystem integration
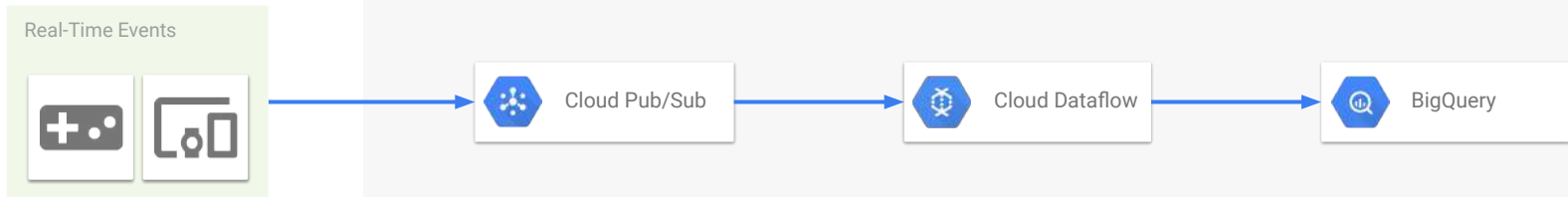- Geo-expansion

Google Cloud

Cloud Innovator
MEGAZONE

# Demo

Google Cloud

Architecture: Demo Scenario

Google Cloud Platform

Real-Time Events

Cloud Pub/Sub → Cloud Dataflow → BigQuery

# Thank you

Google Cloud

Cloud Innovator
MEGAZONE