# Google Dataflow 小試

Simon Su @ LinkerNetworks
*{Google Developer Expert}*

# I'm Simon Su…

var simon = {/** I am at GCPUG.TW **/};

simon.aboutme = 'http://about.me/peihsinsu';

simon.nodejs = 'http://opennodes.arecord.us';

simon.googleshare = 'http://gappsnews.blogspot.tw'

simon.nodejsblog = 'http://nodejs-in-example.blogspot.tw';

simon.blog = 'http://peihsinsu.blogspot.com';

simon.slideshare = 'http://slideshare.net/peihsinsu/';

simon.email = 'simonsu.mail@gmail.com';

simon.say('Good luck to everybody!');

https://www.facebook.com/groups/GCPUG.TW/

https://plus.google.com/u/0/communities/116100913832589966421

- Data scientist
- Data engineer
- Frontend engineer

# Google Cloud in Big Data Solution

# Google Focused Cloud

## Now

### Assembly required

1st Wave
**Colocation**

Your kit, someone else's building. Yours to manage.

2nd Wave
**Virtualized Data Centers**

Standard virtual kit for Rent. Still yours to manage.

Storage    Processing    Memory    Network

## Next

### True On Demand Cloud

3rd Wave
**An actual, global elastic cloud**

Invest your energy in great apps.

Clusters

Containers

Distributed Storage, Processing & Machine Learning

# Google Cloud Family

Application
Runtime Services

**Enabling No-Touch Operations**

Data Services

**Breakthrough Insights,
Breakthrough Applications**

Foundation
Infrastructure & Operations

**The Gear that Powers Google**

# GCP tools for data processing and analysis

**Capture**

Pub/Sub
Logs
App Engine
BigQuery streaming

**Store**

Cloud Storage

BigQuery Storage

Cloud SQL (mySQL)

Cloud Datastore (NoSQL)

**Process**
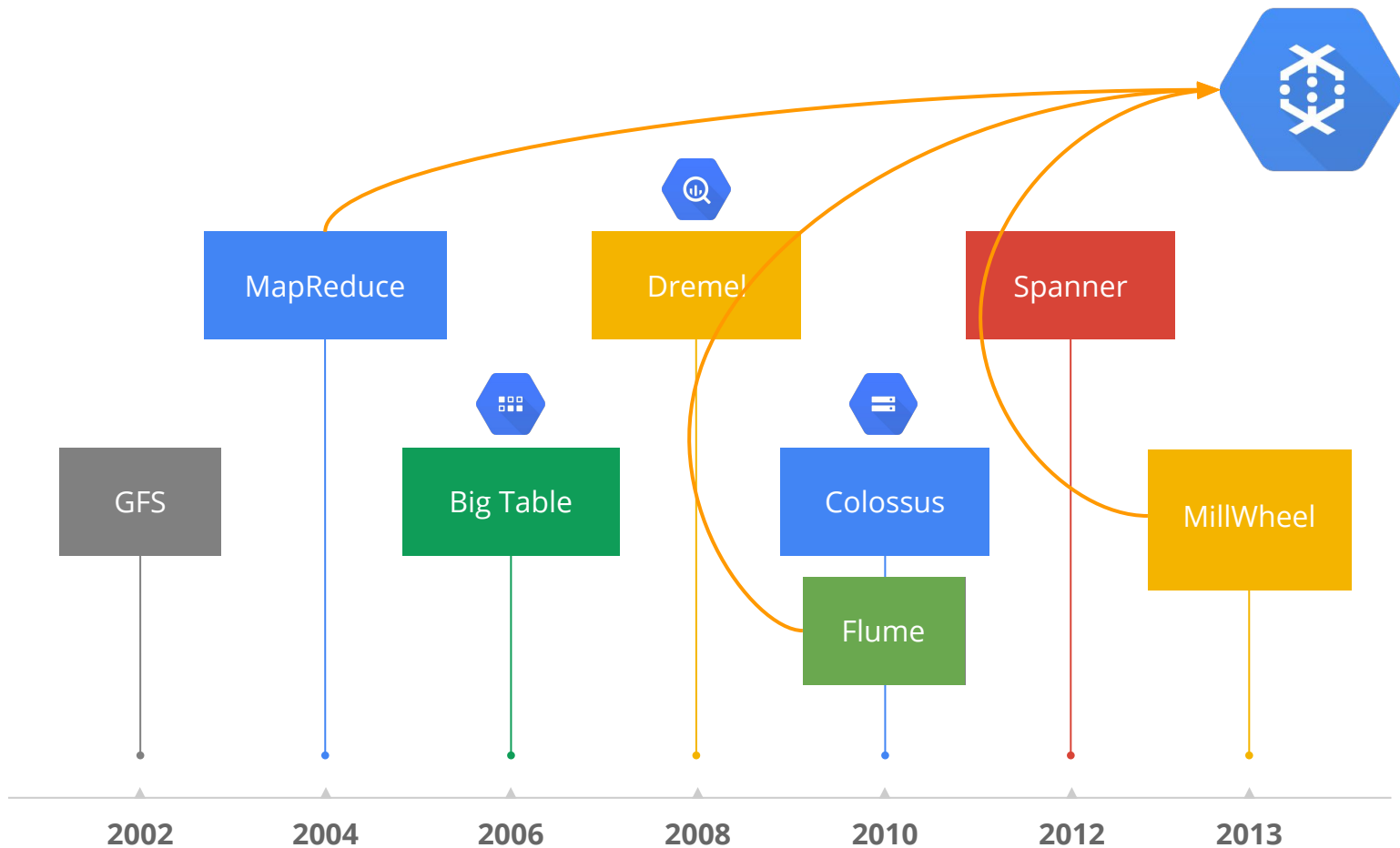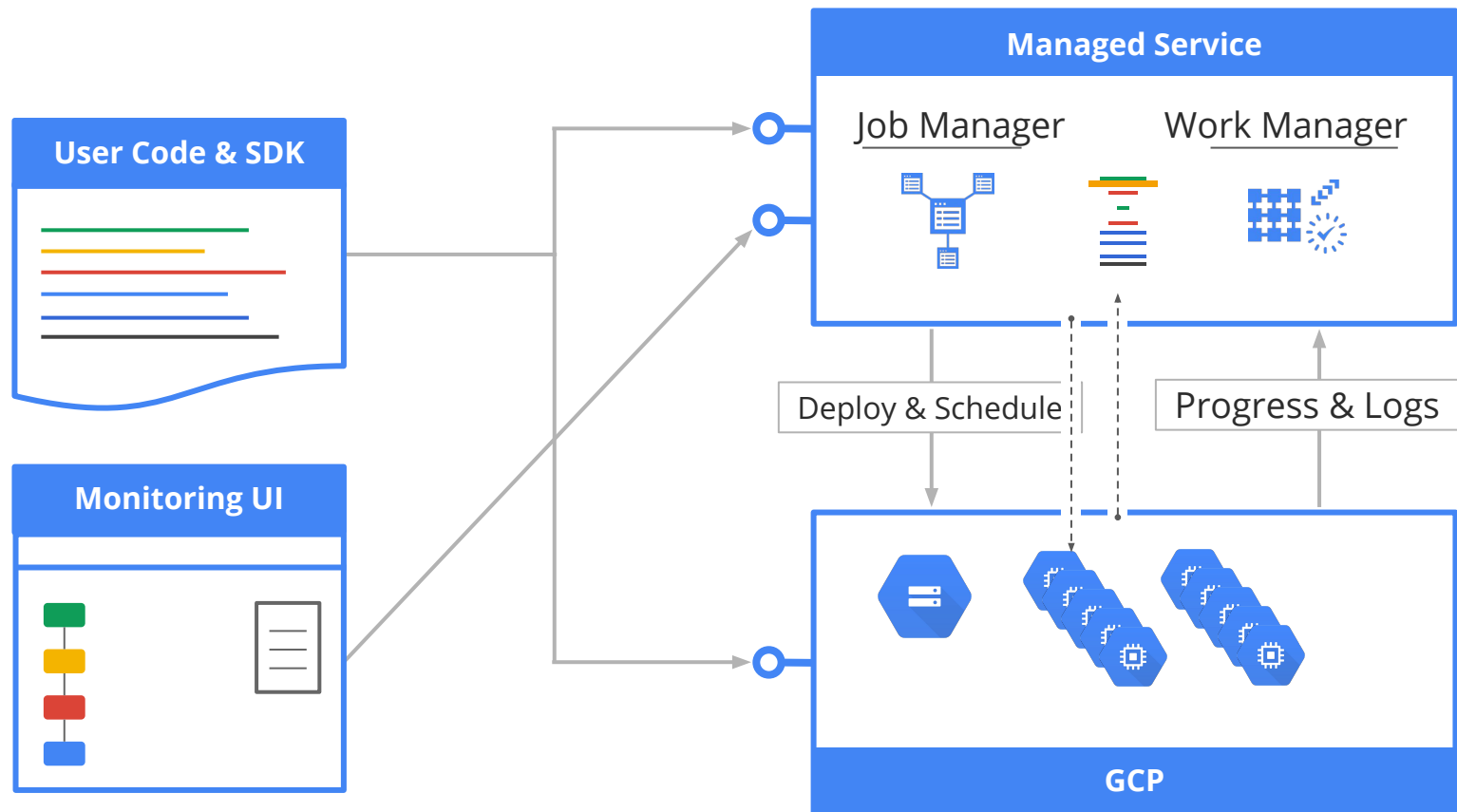
Dataflow

Dataproc

**Analyze**

BigQuery

Dataproc

Larger Hadoop Ecosystem

# Common big data processing flow



Devices → [server] → kafka → Spark → cassandra

1M Devices → → 16.6K Events/sec → → 16.6K Events/sec → → 43B Events/month →

# Look into Dataflow

# Dataflow use case

### ETL

- Movement
- Filtering
- Enrichment
- Shaping

### Analysis

- Reduction
- Batch computation
- Continuous computation

### Orchestration

- Composition
- External orchestration
- Simulation

# Getting Start - Installation

# Get your GCP project

# Install Eclipse Plugin for Dataflow

# Verify your installation

# Getting start with GCS

# Google Cloud Storage Features



Regional buckets

Object versioning

Offline import (third party)

ACLs

Object lifecycle management

Online cloud import (Cloud Storage Transfer Service)

Object change notification

# Create your bucket for Dataflow use

# Run Dataflow in Local

# Create dataflow project

# Dataflow Sample

```
    Project Explorer   [H] Package Explorer  ⊠       ⊟ ⊟

  ▼ PipelineSample
    ▼ src/main/java
      ▼ com.gcpdemo.dataflow
        ▶ J StarterPipeline.java
    ▶ JRE System Library [JavaSE-1.7]
    ▶ Maven Dependencies
    ▶ bin
    ▶ src
      target
    M pom.xml
```

```java
46  @SuppressWarnings("serial")
47  public class StarterPipeline {
48    private static final Logger LOG = LoggerFactory.getLogger(StarterPipeline.class);
49
50⊖   public static void main(String[] args) {
51      Pipeline p = Pipeline.create(
52          PipelineOptionsFactory.fromArgs(args).withValidation().create());
53
54      p.apply(Create.of("Hello", "World"))
55⊖     .apply(ParDo.of(new DoFn<String, String>() {
56⊖       @Override
△57       public void processElement(ProcessContext c) {
58           c.output(c.element().toUpperCase());
59         }
60      }))
61⊖     .apply(ParDo.of(new DoFn<String, Void>() {
62⊖       @Override
△63       public void processElement(ProcessContext c)  {
64           LOG.info(c.element());
65         }
66      }));
67
68      p.run();
69    }
70  }
71
```

# Execute Dataflow

# Lab 1: Ready your dataflow environment and create your first dataflow project

- **Create GCP project**
- **Install Eclipse and Dataflow plugin**
- **Create first Dataflow project**
- **Run your project**

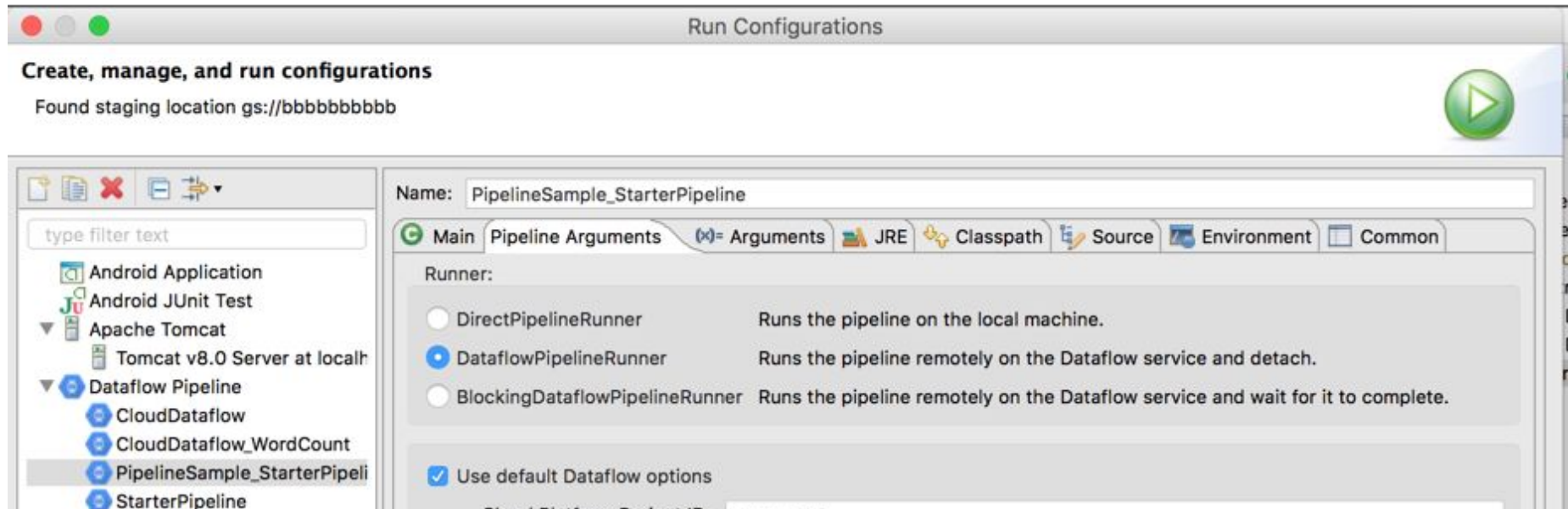# After lab - What thing in GCS bucket

# After lab - The Run Configuration

# Dataflow in Batch Mode

# What we do in Big Data process...

| Map | | ParDo |
| --- | --- | --- |
| ↓ | | ↓ |
| Shuffle | **=** | GroupByKey |
| ↓ | | ↓ |
| Reduce | | ParDo |

# Pipeline

- A Direct Acyclic Graph of data processing transformations
- Can be submitted to the Dataflow Service for optimization and execution or executed on an alternate runner e.g. Spark
- May include multiple inputs and multiple outputs
- May encompass many logical MapReduce operations
- PCollections flow through the pipeline

# Inputs & Outputs

› Read from standard Google Cloud Platform
  data sources

  • GCS, Pub/Sub, BigQuery, Datastore

› Write your own custom source by teaching
  Dataflow how to read it in parallel

  • Currently for bounded sources only

› Write to GCS, BigQuery, Pub/Sub

  • More coming…

› Can use a combination of text, JSON, XML,
  Avro formatted data

Your

Source/Sink

Here

# PCollection

› A collection of data of type T in a pipeline
   - PCollection<K,V>

› Maybe be either *bounded* or *unbounded* in size

› Created by using a PTransform to:
   • Build from a java.util.Collection
   • Read from a backing data store
   • Transform an existing PCollection

› Often contain the key-value pairs using KV

```
{Seahawks, NFC, Champions, Seattle,
...}
```

```
{...,
 "NFC Champions #GreenBay",
 "Green Bay #superbowl!",
 ...
 "#GoHawks",
  ...}
```

# Transforms

- A step, or a processing operation that transforms data

    - convert format , group , filter data

- Type of Transforms

    - ParDo

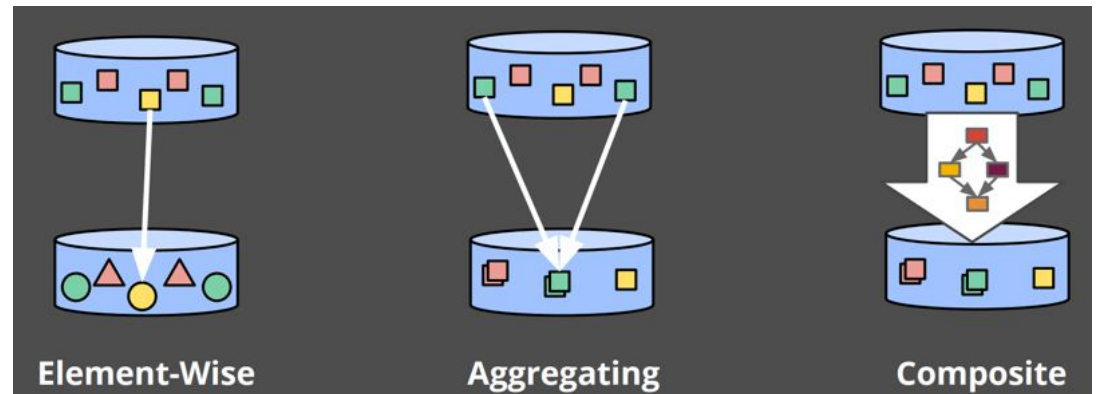    - GroupByKey

    - Combine

    - Flatten



Element-Wise    Aggregating    Composite

    - Multiple PCollection objects that contain the same data type, you can merge them into a single logical PCollection using the Flatten transform

# Pardo (Parallel do)

› Processes each element of a PCollection independently using a user-provided DoFn

› Corresponds to both the Map and Reduce phases in Hadoop i.e. ParDo->GBK->ParDo

› Useful for

**Filtering a data set.**

**Formatting or converting the type of each element in a data set.**

**Extracting parts of each element in a data set.**

**Performing computations on each element in a data set.**

{Seahawks, NFC, Champions, Seattle, ...}

↓

KeyBySessionId

↓

{
  KV<S, Seahawks>,
  KV<C,Champions>,
  <KV<S, Seattle>,
  KV<N, NFC>, …
}

# Group by key

- Takes a PCollection of key-value pairs and gathers up all values with the same key
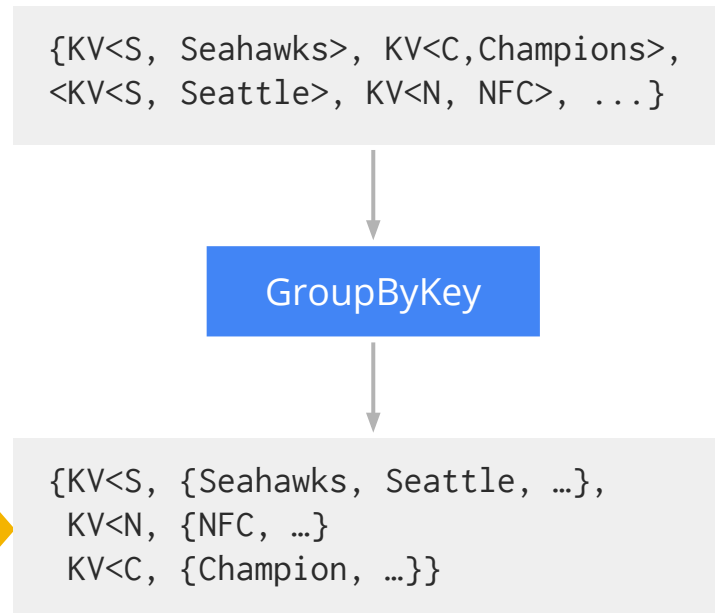
- Corresponds to the shuffle phase in Hadoop

```
{KV<S, Seahawks>, KV<C,Champions>,
<KV<S, Seattle>, KV<N, NFC>, ...}
```

GroupByKey

```
{KV<S, {Seahawks, Seattle, …},
 KV<N, {NFC, …}
 KV<C, {Champion, …}}
```

*Wait a minute…*
*How do you do a GroupByKey on an unbounded PCollection?*

# Group by key sample

```java
@Override
public PCollection<KV<T, Long>> apply(PCollection<T> in) {
  return
      in
      .apply(ParDo.named("Init")
              .of(new DoFn<T, KV<T, Long>>() {
                  @Override
                  public void processElement(ProcessContext c) {
                    c.output(KV.of(c.element(), 1L));
                  }
                }))

      .apply(Combine.<T, Long>perKey(
              new SerializableFunction<Iterable<Long>, Long>() {
                  @Override
                  public Long apply(Iterable<Long> values) {
                    long sum = 0;
                    for (Long value : values) {
                      sum += value;
                    }
                    return sum;
                  }
                }));
}
```

# Lab 2: Deploy your first project to Google Cloud Platform

- **Checking the Lab 1 project working well**
- **Deploy to cloud and watch the dataflow task dashboard**
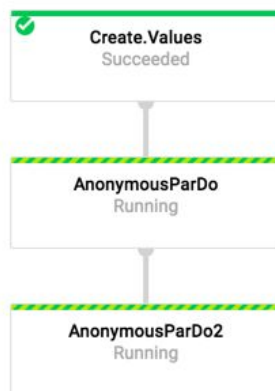
# Dataflow Task Dashboard

# Dataflow in Streaming Mode

# Pub/Sub working model

# Pub/Sub Operations - Topics

```
$ gcloud --format=json alpha pubsub topics create my-topic


[
  {
    "reason": "",
    "success": true,
    "topicId": "projects/sunny-573/topics/my-topic"
  }
]
```

```
$ gcloud alpha pubsub topics publish my-topic '{"aaa":123,"bbb":223}'
```

# Pub/Sub Operations - Subscriber

```
$ gcloud alpha pubsub subscriptions create sub002 --topic my-topic
---
ackDeadlineSeconds: 10
pushEndpoint: null
reason: ''
subscriptionId: projects/sunny-573/subscriptions/sub002
success: true
topic: projects/sunny-573/topics/my-topic
type: push
```

```
$ gcloud alpha pubsub subscriptions pull sub001
```

| DATA | MESSAGE_ID | ATTRIBUTES | ACK_ID |
|------|-----------|-----------|--------|
| {"aaa":123,"bbb":223} | 43961024144056 | | MTJFQV5AEkw6...4cqZhg9XxJLLD5- |

```
$ gcloud alpha pubsub subscriptions ack sub001 MkVBXkASTDo...JLLD5-MQ
ackIds:
- MkVBXkASTDo...JLLD5-MQ
subscriptionId: projects/sunny-573/subscriptions/sub001
```

# Simple Guide for Pub/Sub

# Dataflow with Cloud Pub/Sub Use Case

Globally redundant
Low latency (sub sec.)
N to N coupling
Batched read/write
Push & Pull
Guaranteed Delivery
Auto expiration

fluentd

YOUR LOGO HERE

| Publisher A | Publisher B | | Publisher C |
|---|---|---|---|

Message 1    Message 2    Message 3

Topic A    Topic B    Topic C

Subscription XA    Subscription XB    Subscription YC    Subscription ZC

Cloud Pub/Sub

Message 1    Message 2    Message 3    Message 3

Subscriber X    Subscriber Y    Subscriber Z

# Lab 3: Create a Streaming Dataflow model

- Create PubSub topic
- Deploy Dataflow streaming sample
- Watch Dataflow task dashboard

# After Dataflow



**Datalab**
An easy tool for analysis and report

**BigQuery**
An interactive analysis service

# Google Data Studio