# Data Analytics Platforms in the Cloud

# Large Scale Data Platforms in the Cloud

Two marketing customers

zühlke
empowering ideas

| 60TB | 55B | 300 | 800 | Real-time processing of click stream data |
|------|-----|-----|-----|-------------------------------------------|
| Raw data on linear blob storage | Number of rows in MPP Database | Number of nodes in transparent elastic scaling | Websites with simultanious click stream collection | Automated Machine Learning Lifecycle |
| | | | | Collaboration with Google |

Joel Akeret
Expert Data Scientist
Zühlke Engineering AG



Hendrik Schöneberg
Lead Software Architect
Zühlke Engineering AG

# Why data analytics platforms?

# Drivers for Data Analytics Platforms

zühlke
empowering ideas

## Business Use Cases

| Personalized Content | Website Optimization | Churn Prediction | Targeted Marketing | Automated Lead Scoring |
|---|---|---|---|---|

## Capabililties / Enablers

| Operationalized Machine Learning | A/B/n Testing | Data Exploration |
|---|---|---|

**Data Platform / 360° View**

Scalability

Workload

Volume

Costs

BIG DATA & AI LANDSCAPE 2018

http://mattturck.com/bigdata2018

How to choose the right components?

# Data platform architecture

# Platform Architecture Overview



**zühlke**

**Source Systems**
- Databases
- External Systems
- Files
- Enterprise Service Bus

**Ingestion Layer**
- Batch
- Streaming

**Staging Area / Raw Storage**
- Blob Storage

**Data Pipelines**
- ETL

**Machine Learning Operationalization**
- Model Training
- Model Serving

**Refined Storage**
- DWH
- Catalogue & Metadata
- Relational Data
- Non-Relational Data

**Data Driven Applications**
- Data Analytics
- Machine Learning
- Exploration
- Business Intelligence
- Reporting

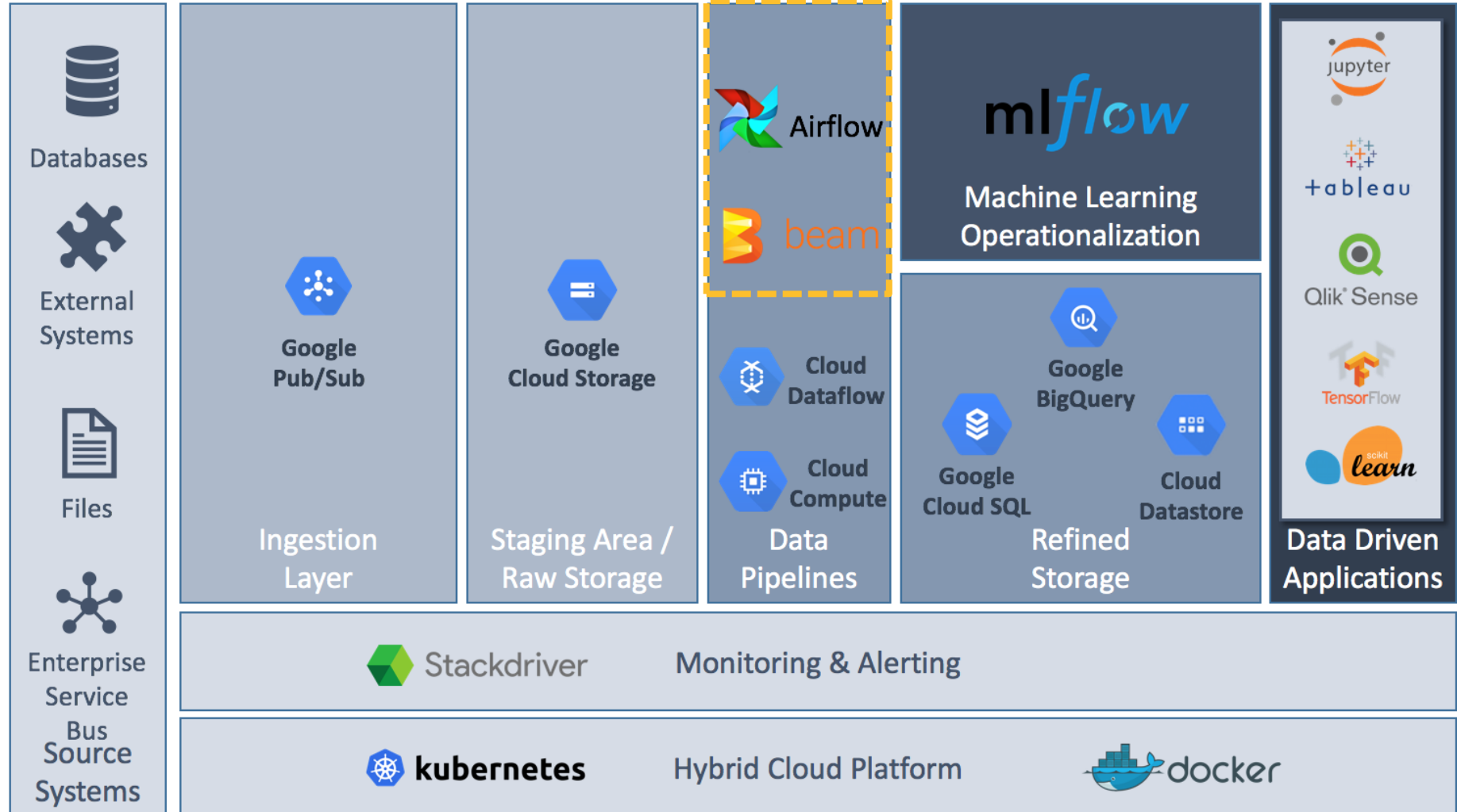**Data Governance & Security**

**Monitoring, Alerting**

# Technological choices

Amazon Web Services

# Technological choices

## Google Cloud Platform

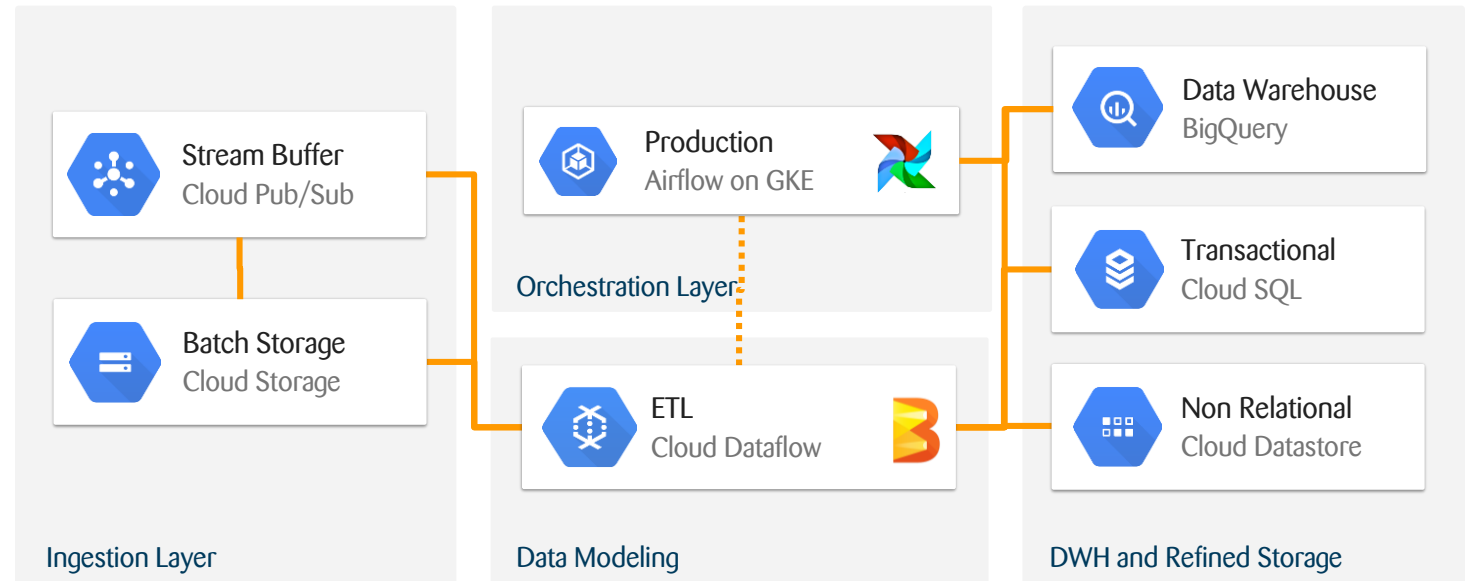# Deepdive
# Data Processing Pipelines

# Data Platform Architecture

Core Data Ingestion and Processing

■ Workflows orchestrated by **Apache Airflow**

■ Execution of parallelizable data transformations done with **Apache Beam** on **Google Cloud Dataflow**

beam

**Ingestion Layer**
- Stream Buffer — Cloud Pub/Sub
- Batch Storage — Cloud Storage

**Orchestration Layer**
- Production — Airflow on GKE

**Data Modeling**
- ETL — Cloud Dataflow

**DWH and Refined Storage**
- Data Warehouse — BigQuery
- Transactional — Cloud SQL
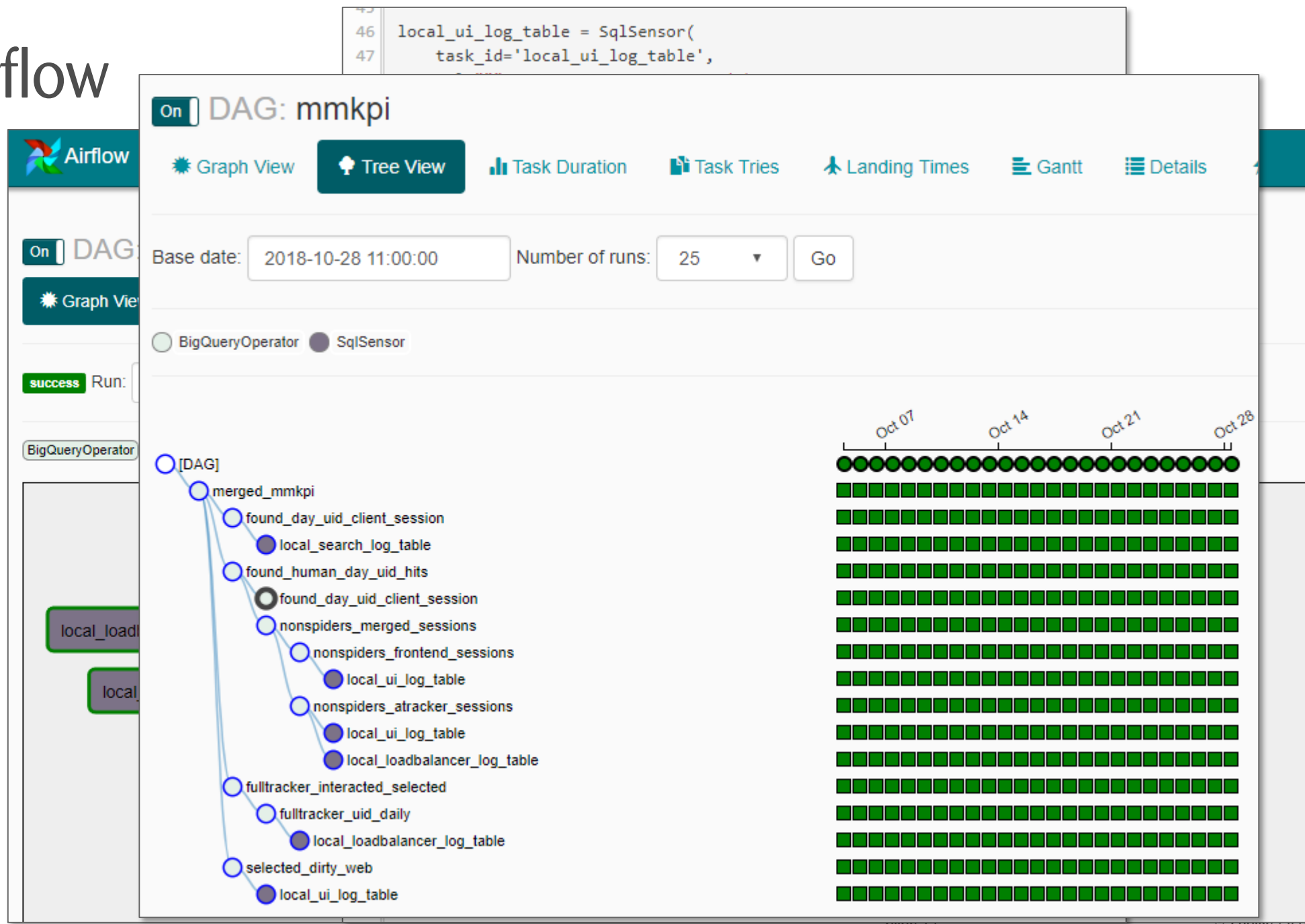- Non Relational — Cloud Datastore

# Apache Airflow

- Originally created at Airbnb, open sourced 2015

- Author, schedule and monitor workflows

- "Cron on steroids"

- Workflows are part of the codebase (Python)

- Workflows defined as DAGs of tasks

- Clear and transparent

- Easy to rerun or reproduce historical jobs by date →backfilling
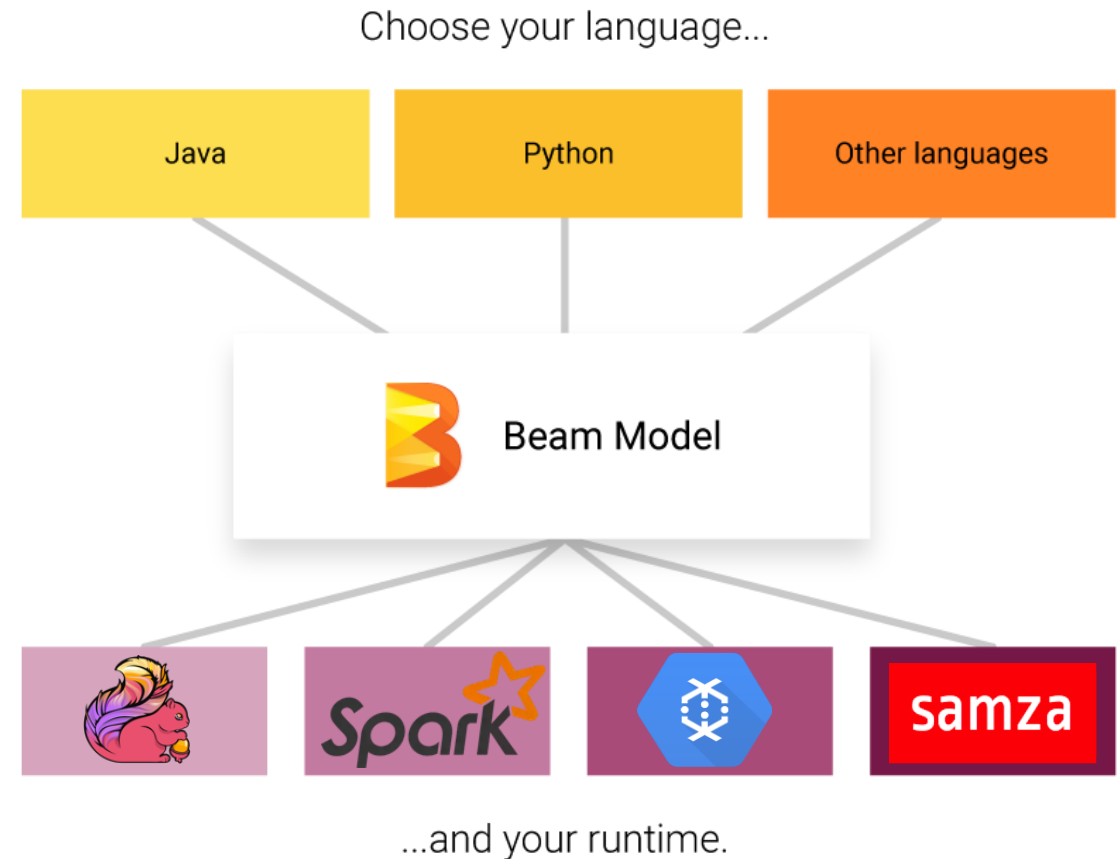
- Thriving community

# Apache Airflow

## Example

# Apache Beam

- Unified **model** for batch and streaming

- Executes on a broad variety of runners (no vendor lock-in)

- Decouples data processing from the executor

- Comprehensive set of windowing, timing, lateness and triggering primitives

Choose your language...

| Java | Python | Other languages |

Beam Model

Spark    samza

...and your runtime.

**Source**
- http://beam.apache.org

# Apache Beam

Python SDK



```python
61    with beam.Pipeline(options=options) as p:
62      messages = (p | 'ReadPubSub' >> beam.io.ReadFromPubSub(args.topic)
63                    | "ParseRawMessage" >> beam.Map(parse_message)
64                    | "ConvertRadiogram" >> beam.Map(convert_to_dict)
65                    | "ProcessMessage" >> beam.Map(process_message)
66                 )
67
68      (messages | "ConvertForPubSub" >> beam.Map(to_json)
69                | "WriteToPubSub" >> beam.io.WriteToPubSub(args.response_topic)
70      )
71
72      messages | "WriteToBigQuery" >> beam.io.WriteToBigQuery(
73                    args.table_name, args.dataset,
74                    schema=SCHEMA,
75                    create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED,
76                    write_disposition=beam.io.BigQueryDisposition.WRITE_APPEND)
77
78
```

# Apache Beam

Resulting pipeline

# Apache Beam

## Unified model for stream and batch processing

Example: Calculate Team scores by the hour



**Apache Beam**

Where?   When?

```
gameEvents
  [... input ...]
  .apply("LeaderboardTeamFixedWindows", Window
    .<GameActionInfo>into(FixedWindows.of(
      Duration.standardMinutes(Durations.minutes(60))))   ] Window
    .triggering(AfterWatermark.pastEndOfWindow()           - Watermark Trigger
      .withEarlyFirings(AfterProcessingTime.pastFirstElementInPane()  ] Early Trigger
        .plusDelayOf(Durations.minutes(5)))
      .withLateFirings(AfterProcessingTime.pastFirstElementInPane()   ] Late Trigger
        .plusDelayOf(Durations.minutes(10))))
    .withAllowedLateness(Duration.standardMinutes(120)     - Garbage Collection
    .accumulatingFiredPanes())                             - Accumulation
  .apply("ExtractTeamScore", new ExtractAndSumScore("team"))  - Sum
  [... output ...]
```

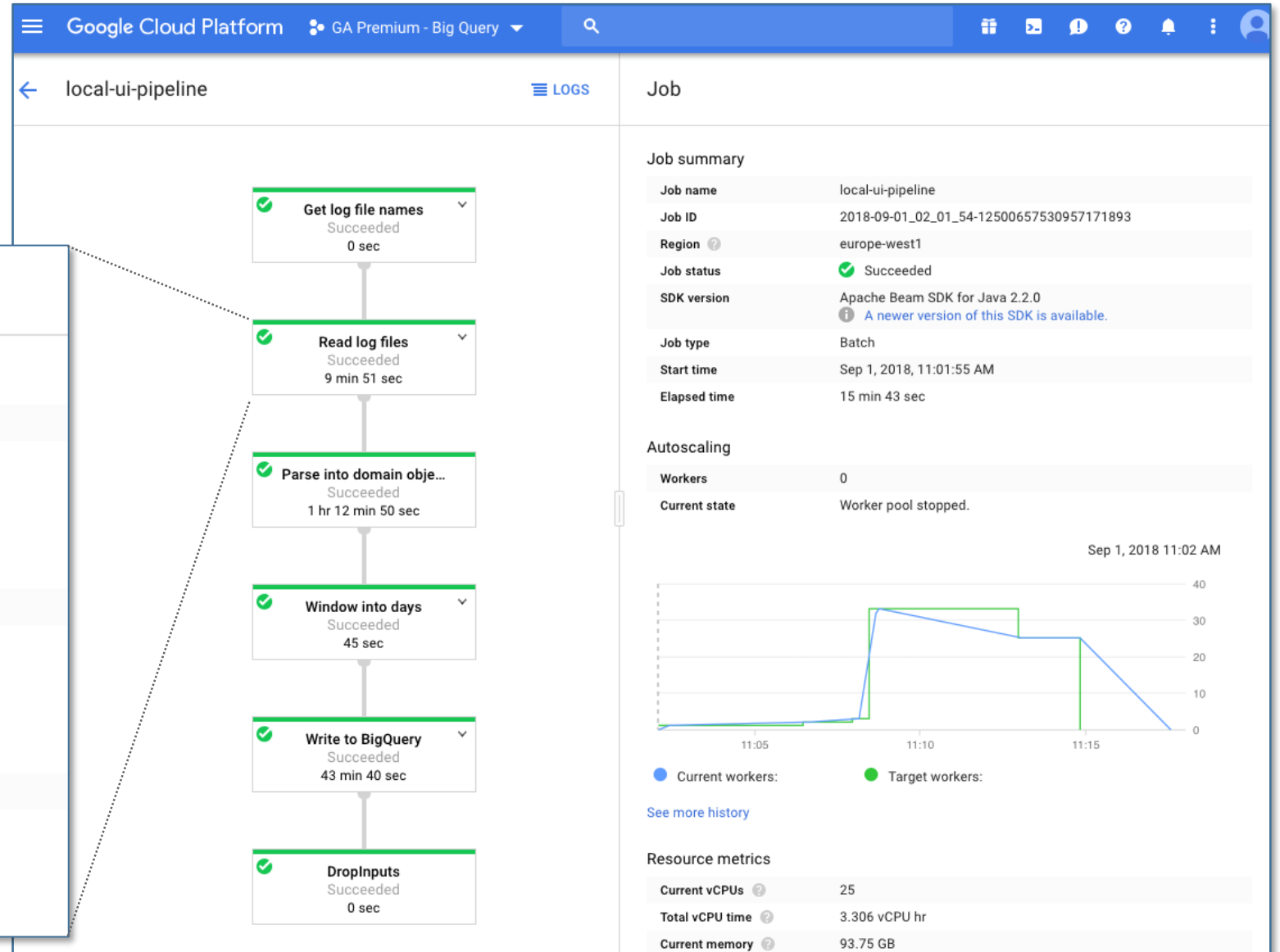How?   What?

**Spark**

```
gameEvents
  [... input ...]
  .window(Durations.minutes(10), Durations.minutes(10))
  .mapToPair(new ExtractUserScore())
  .reduceByKey(new SumScore())
  .transformToPair((rdd, timestamp) -> {
    userWindowTimestamp.set(Math.max(
      userWindowTimestamp.get(), timestamp.milliseconds()));
    return rdd;
  })
  .updateStateByKey(new SumAggregator())
  .filter(x -> x._2().timestamp() >= userWindowTimestamp.get())
  [... output ...]

private static class SumAggregator implements Function2<
    List, Optional,
    Optional> {
  final private static Integer INITIAL_STATE = 0;

  public Optional call(
      List scores, Optional state) {
    if (scores.size() == 0) return state;

    Integer sumWithTimestamp = state.or(INITIAL_STATE);
    Integer sum = sumWithTimestamp.val() +
      scores.stream().mapToInt(Integer::intValue).sum();
    return Optional.of(sum);
  }
}
```

Window, Trigger
Accumulation,
& Sum (but no
Lateness), all
mixed together

Source:
https://cloud.google.com/dataflow/blog/dataflow-beam-and-spark-comparison
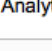
# Apache Beam

## Autoscaling

# Cloud Vendor comparison

# Data Platform Monthly Breakdown



| | | | |
|---|---|---|---|
| Cloud Storage | Multi-Regional storage | 15360 GB | $399.36 |
| Cloud Storage | Regional storage | 5120 GB | $117.76 |
| Cloud Storage | Coldline storage | 10240 GB | $102.40 |
| Analytics | BigQuery | 10 GB | $949.80 |
| Cloud Dataflow | 50 x n1-standard-2 workers in | 3000 | $232.03 |
| 5 x Application Abstraction Layer | n1-standard-4

Sustained Usage Discount Monthly Breakdown:

• 1st ¼ - 912.5 hrs @ 0.0% off: $223.38
• 2nd ¼ - 912.5 hrs @ 20.0% off: $178.70 ($44.68 saved)
• 3rd ¼ - 912.5 hrs @ 40.0% off: $134.03 ($89.35 saved)
• 4th ¼ - 912.5 hrs @ 60.0% off: $89.35 ($134.03 saved) | 3650 total hours per month | $625.46 |
| Persistent disk | Storage | 500 GB | $24.00 |
| **Total Estimated Monthly Cost** | | | **$2,450.82** |

Storage
Database
Elastic Processing
Compute

*Google Cloud Platform*

| Service Type | Components | Region | Component Price | Service Price |
|---|---|---|---|---|
| Amazon EC2 Service (EU (Frankfurt)) | | | | $612.32 |
| | Compute: | EU (Frankfurt) | $439.2 | |
| | EBS Volumes: | EU (Frankfurt) | $147.5 | |
| | Elastic IPs: | EU (Frankfurt) | $3.66 | |
| | Classic LBs: | EU (Frankfurt) | $21.96 | |
| Amazon S3 Service (EU (Frankfurt)) | | | | $445.45 |
| | S3 Standard Storage: | EU (Frankfurt) | $376.32 | |
| | S3 Standard Other Requests: | EU (Frankfurt) | $0.01 | |
| | S3 Standard - IA Storage: | EU (Frankfurt) | $69.12 | |
| Amazon Redshift Service (EU (Frankfurt)) | | | | $1055.81 |
| | Compute: | EU (Frankfurt) | $1055.81 | |
| Amazon Glacier Service (EU (Frankfurt)) | | | | $46.08 |
| | Storage: | EU (Frankfurt) | $46.08 | |
| Amazon Elastic MapReduce Service (EU (Frankfurt)) | | | | $189.45 |
| | Compute: | EU (Frankfurt) | $189.45 | |
| AWS Support (Business) | | | | $232.53 |
| | Support for all AWS services: | | $232.53 | |
| | Free Tier Discount: | | | $-23.87 |
| **Total Monthly Payment:** | | | | **$2557.77** |

Compute
Storage
Database
Elastic Processing
Support Fee

*amazon web services*

*zühlke empowering ideas*

# Conclusion

- We created **reproducible** and **traceable** data pipelines using Apache Airflow and Apache Beam

- We established an essential building block to **incrementally** create **trustworthy** data-driven applications

- The tech stack relies on OSS and avoids a vendor lock-in (runs on own laptop, on-prem, AWS, GCP, …)

- Great collaboration with Google

# Thank you