



STYLIGHT

# Building data pipelines: from simple to more advanced - hands-on experience

SERGII KHOMENKO, DATA SCIENTIST,  
SERGII.KHOMENKO@STYLIGHT.COM, @lcod3r

STYLIGHT.COM

# AGENDA

Who? What? Why?

The Good, The Bad And The Legacy

Open Source stack

Amazon AWS

Google Cloud

Tips, tricks and best practices

In computing, a pipeline is a set of data processing elements connected in series, where the output of one element is the input of the next one.



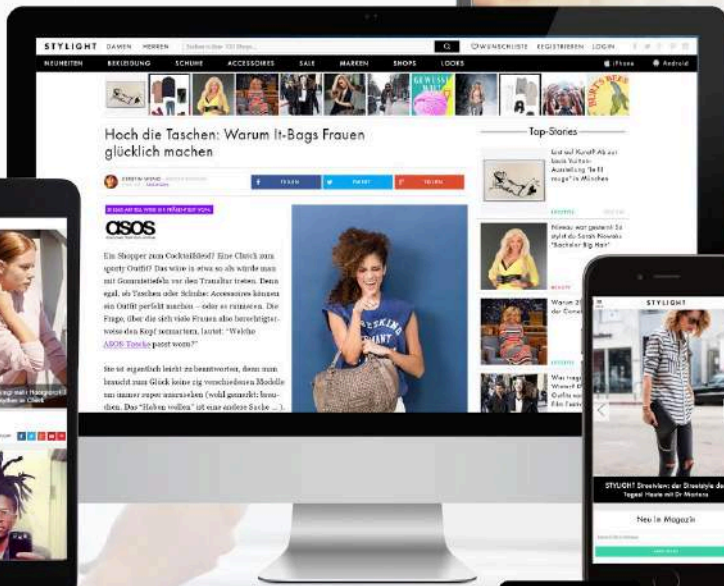
## **Sergii Khomenko**

Data Scientist at one of the biggest fashion communities, STYLIGHT. Data analysis and visualization hobbyist.

Speaker at Berlin Buzzwords 2014, ApacheCon Europe 2014, Puppet Camp London 2015

Founder and speaker at Munich Golang UG, Munich Tableau UG.  
Speaker at Munich UserR Group, Munich Search UG, Munich Quantified Self UG, Munich DataGeeks.

# STYLIGHT.de



STYLIGHT DAMEN HERREN Suchen in über 100 Shops...

NEUHEITEN BEKLEIDUNG SCHUHE ACCESSOIRES SALE MARKEN SHOPS LOOKS iPhone Android

DAMEN HERREN

MODE Mode 70836 Produkte gefunden

DAMEN > HERREN >

Ähnliche Suchen  
 adidas Originals  
 Aigle  
 Ash  
 Barbour  
 Bench  
 Ben Sherman  
 Birkenstock  
 Broer Bonani  
 Buffalo  
 Bugatti

Nike - GYM VINTAGE Joggingh...  
 24,73 € 39,99 €  
 Versand: kostenlos

Bodyfit - Strickjacke 3/4 Arm L...  
 22,99 € 24,99 €  
 Versand: 3,99 €

Nike - LEGEND 2.0 TRAINING ...  
 29,95 €  
 Versand: kostenlos

Cheap Monday - SLIM Jeans SL...  
 34,92 € 56,99 €  
 Versand: kostenlos

EvenOld - T-Shirt print black  
 9,92 € 14,99 €  
 Versand: kostenlos

Venice Beach - Hot Pants Damen  
 9,92 € 24,95 €  
 Versand: 3,99 €

Urban Classics - College Socks ...  
 4,90 €  
 Versand: 3,30 €

MintBerry - jumpsuit light den...  
 24,92 € 49,99 €  
 Versand: kostenlos

Moderno - Modorno Stoffhose ...  
 14,99 € 24,99 €  
 Versand: 3,30 €

sOliver - Damen Schal 19.409...  
 12,14 € 29,99 €  
 Versand: zzgl. Versandkosten

Zolania Essentials - T-Shirt basic ...  
 6,95 € 14,95 €  
 Versand: kostenlos

Closet - Cocktailkleid / festliche...  
 34,95 € 49,99 €  
 Versand: kostenlos

Auss - Ridley - Röhrhosen in ve...  
 30,00 €  
 Versand: kostenlos

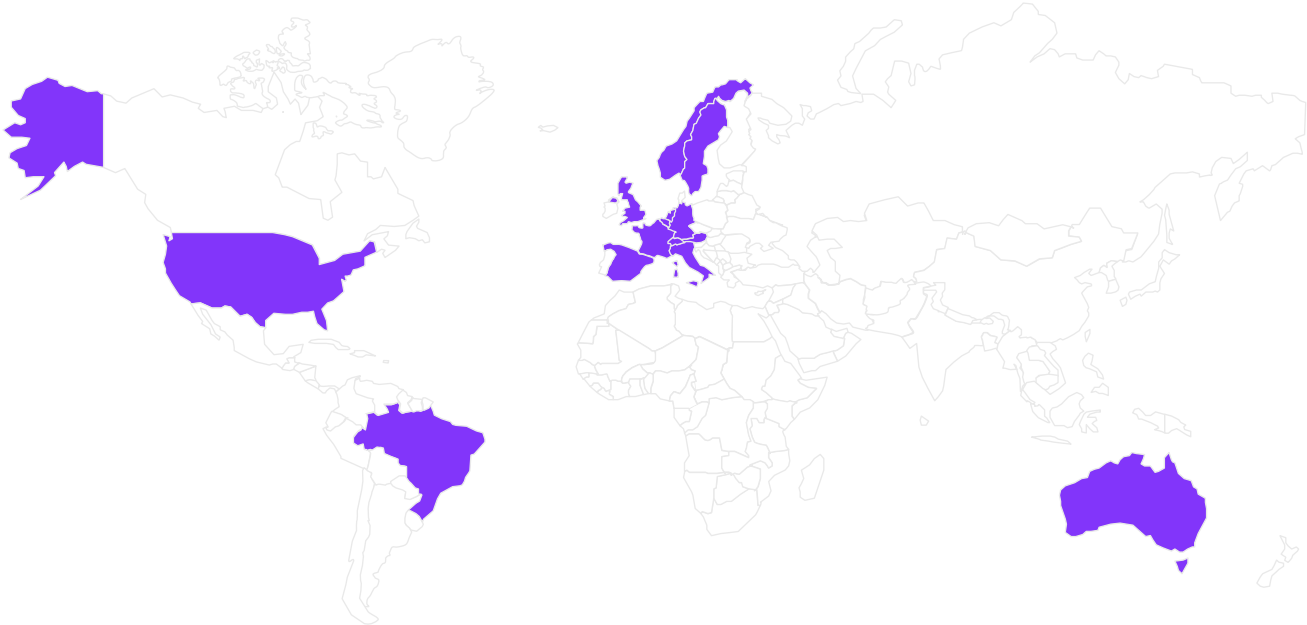
Adidas Originals - Sneaker

Inspiration gefällig?  
 Entdecke Top Blogger Looks auf STYLIGHT



# STYLIGHT – international community

Live in 14 countries





STYLIGHT

# The Good, The Bad And The Legacy

STYLIGHT.COM



# Data sources

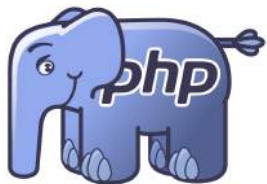
Where your data coming form

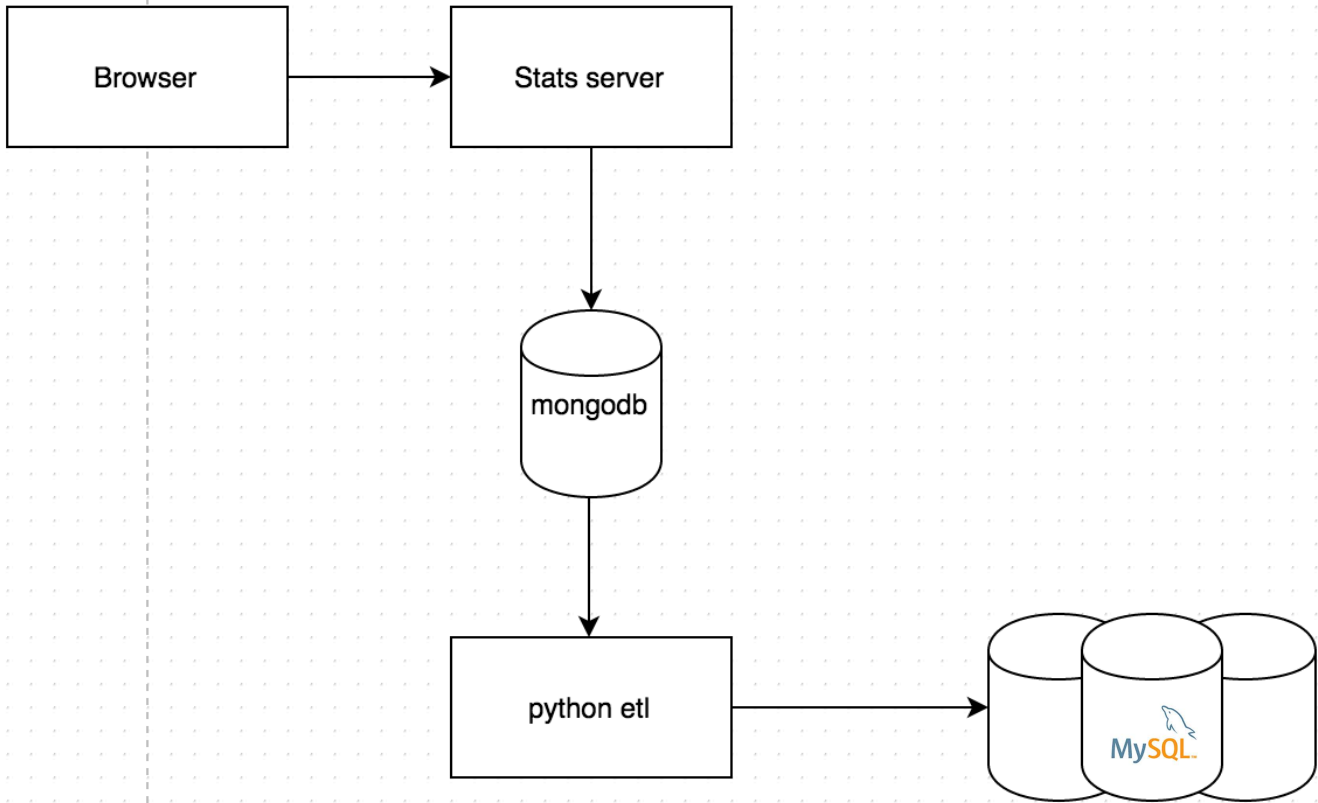
## Sources of data:

- Web tracking
  - Metrics tracking
  - Behavior tracking
- Business intelligence ETL
- Internal Services
  - ML tagging service

# Access patterns

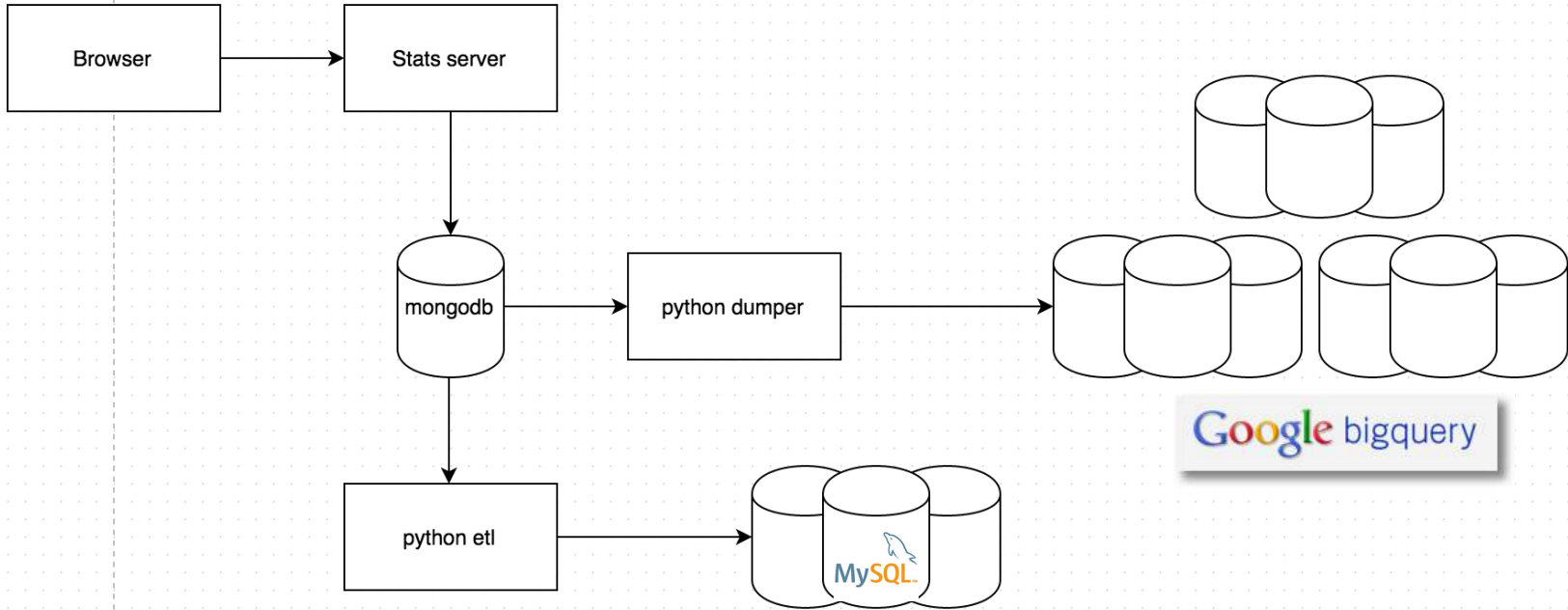
- Real-time
- Nearly real-time
- Daily batches

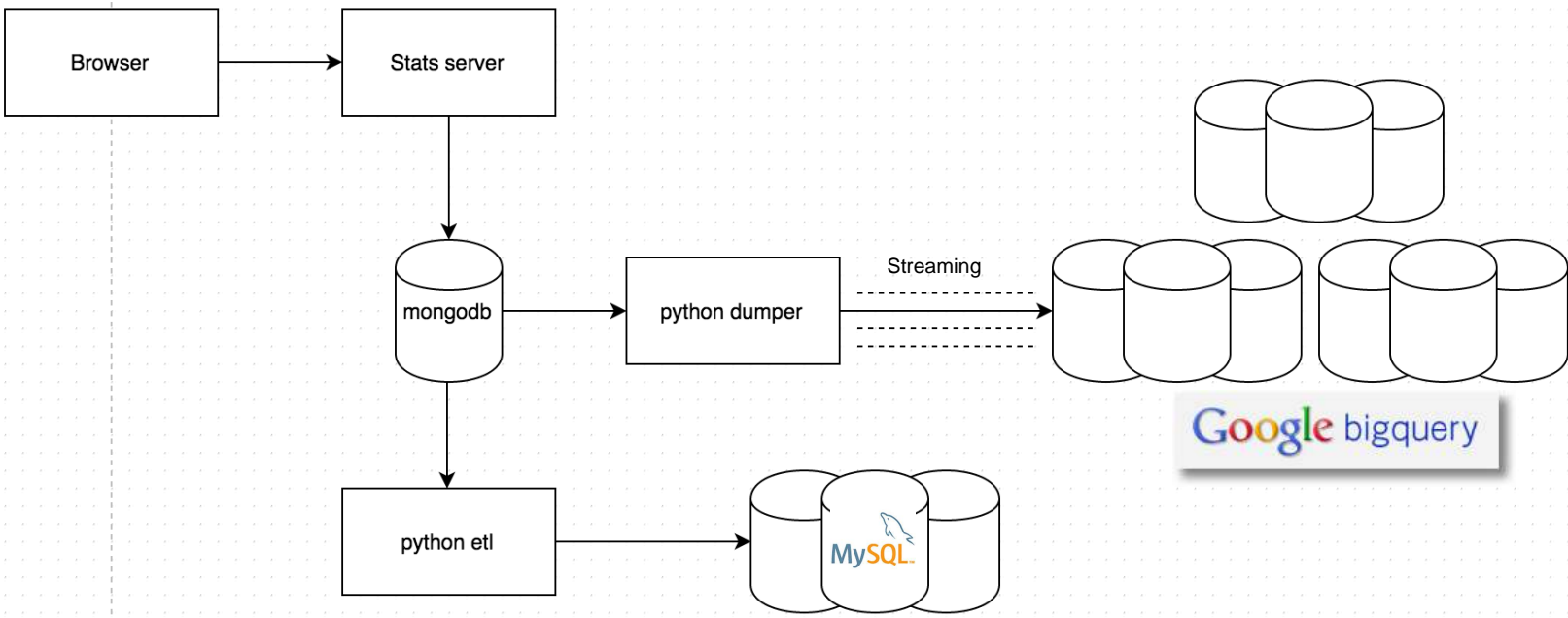




# Properties

- Data consistency
  - Doesn't scale
  - Hard to add new sources
  - Complex system
  - Many interfaces
- 
- As lean and legacy as possible
  - No need for special services





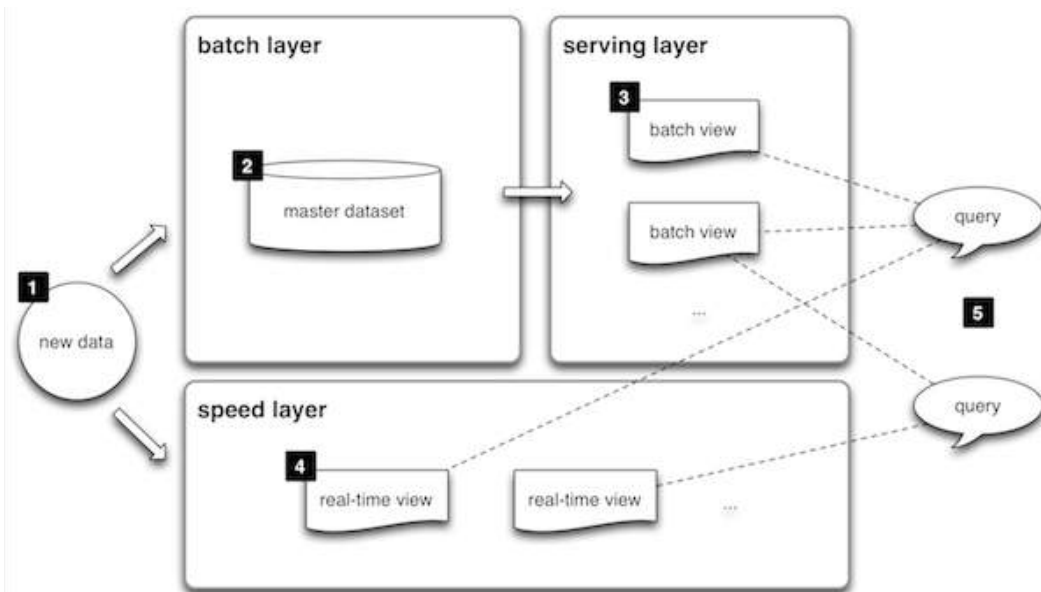


STYLIGHT

# Open Source Stack

STYLIGHT.COM

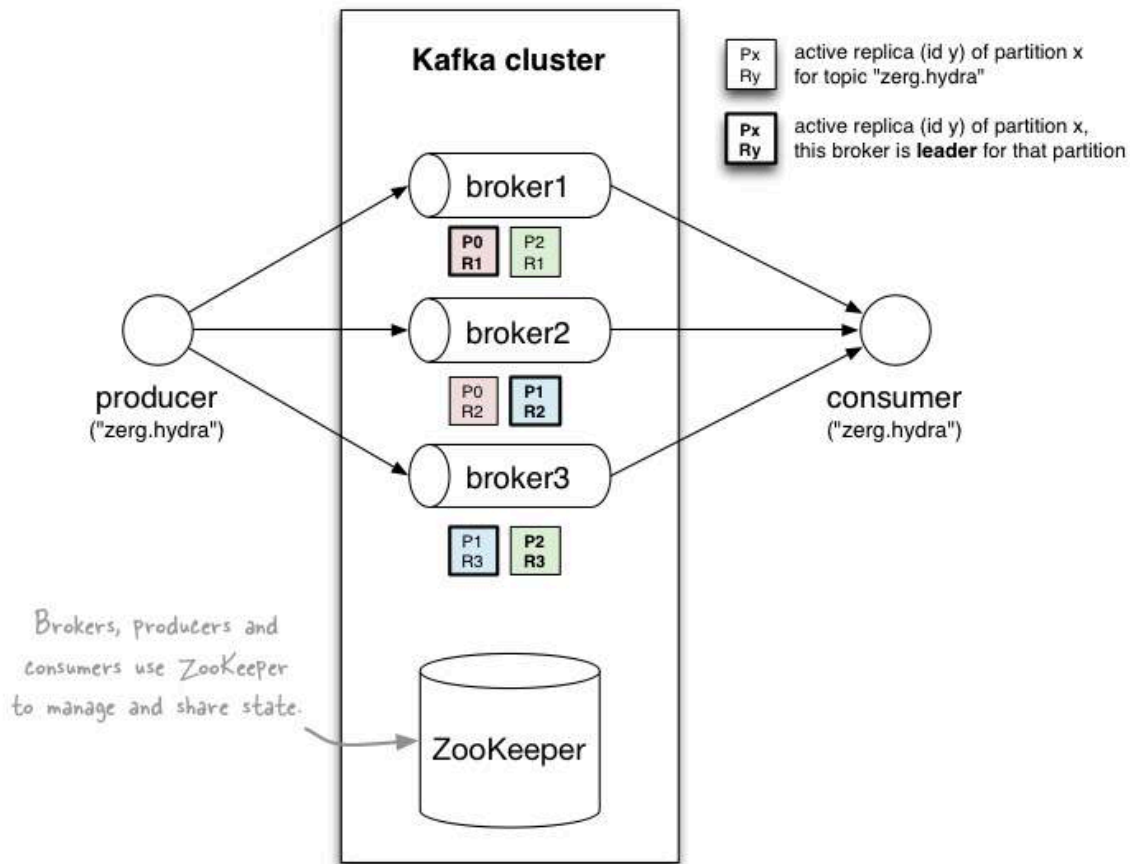




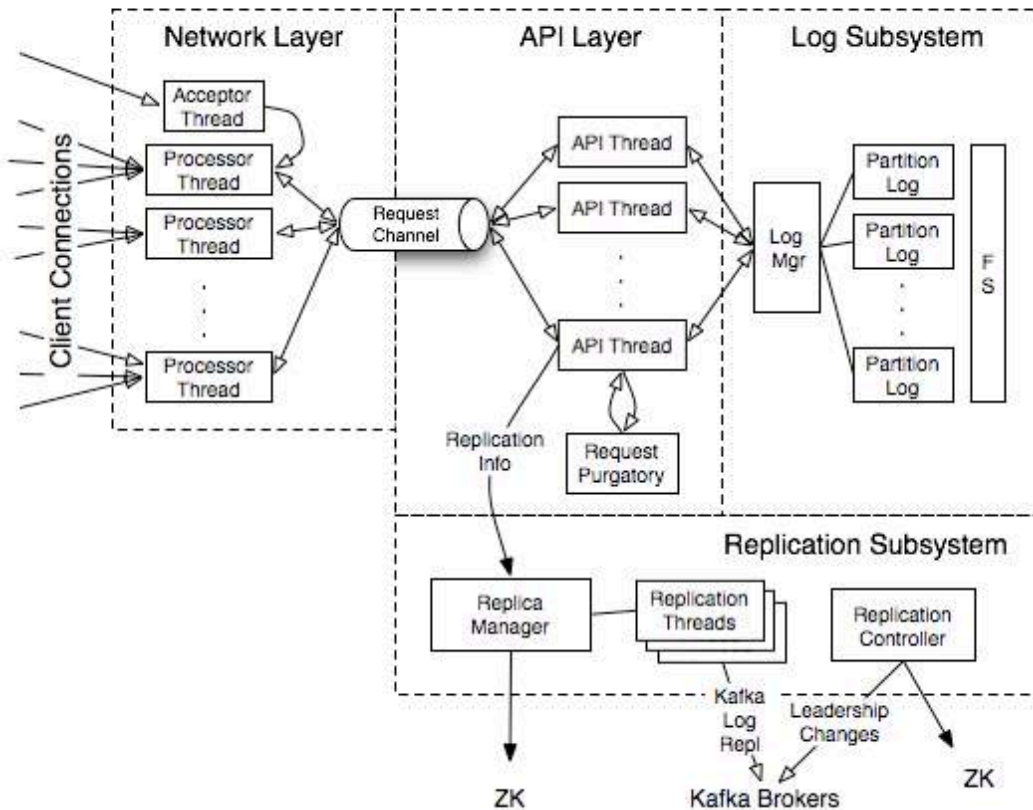
<http://lambda-architecture.net/>

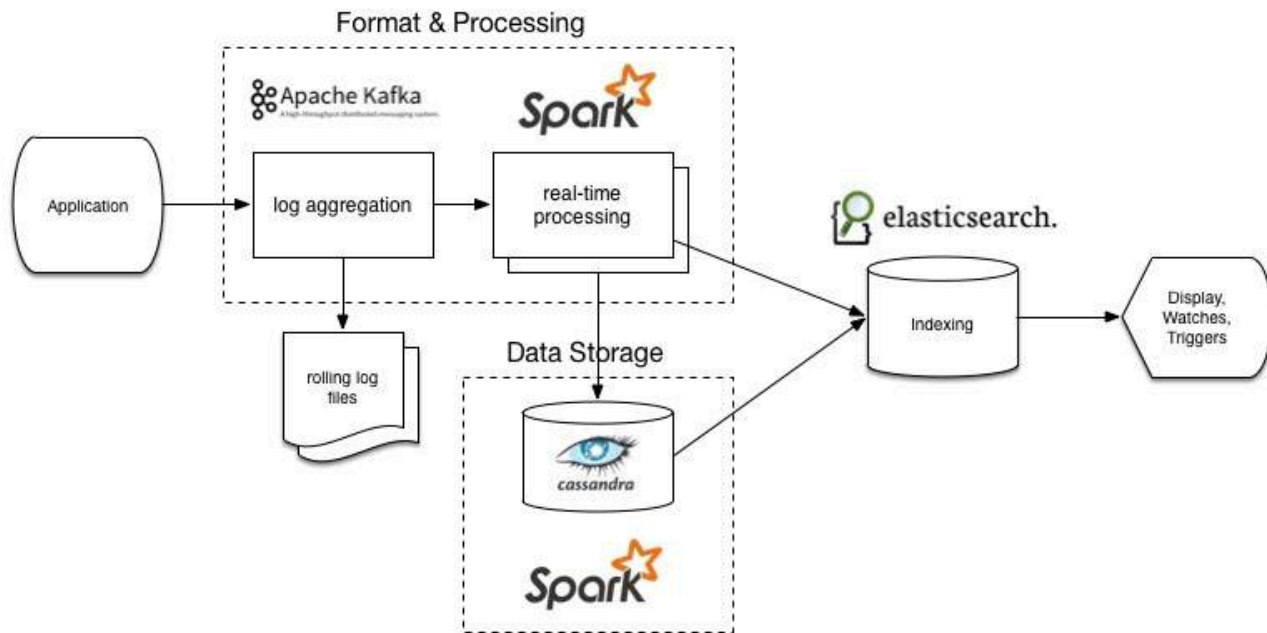


Apache Kafka is publish-  
subscribe messaging rethought  
as a distributed commit log.



# Kafka Broker Internals





<http://www.ipponusa.com/wp-content/uploads/2014/10/spark-architecture.jpg>

# Results

- Scalable
- Flexible
  
- High costs of maintenance
- Not so easy to setup



STYLIGHT

Amazon AWS

STYLIGHT.COM

# Amazon Web Services

## Compute & Networking

-  **Direct Connect**  
Dedicated Network Connection to AWS
-  **EC2**  
Virtual Servers in the Cloud
-  **Route 53**  
Scalable Domain Name System
-  **VPC**  
Isolated Cloud Resources

## Storage & Content Delivery

-  **CloudFront**  
Global Content Delivery Network
-  **Glacier**  
Archive Storage in the Cloud
-  **S3**  
Scalable Storage in the Cloud
-  **Storage Gateway**  
Integrates On-Premises IT Environments with Cloud Storage

## Database

-  **DynamoDB**  
Predictable and Scalable NoSQL Data Store
-  **ElastiCache**  
In-Memory Cache
-  **RDS**  
Managed Relational Database Service
-  **Redshift**  
Managed Petabyte-Scale Data Warehouse Service




## Deployment & Management

-  **CloudFormation**  
Templated AWS Resource Creation
-  **CloudTrail**  
User Activity and Change Tracking
-  **CloudWatch**  
Resource and Application Monitoring
-  **Elastic Beanstalk**  
AWS Application Container
-  **IAM**  
Secure AWS Access Control
-  **OpsWorks**  
DevOps Application Management Service
-  **Trusted Advisor**  
AWS Cloud Optimization Expert

## Analytics

-  **Data Pipeline**  
Orchestration for Data-Driven Workflows
-  **Elastic MapReduce**  
Managed Hadoop Framework
-  **Kinesis**  
Real-time Processing of Streaming Big Data

## Mobile Services

-  **Cognito**  
User Identity and App Data Synchronization
-  **Mobile Analytics**  
Understand App Usage Data at Scale
-  **SNS**  
Push Notification Service

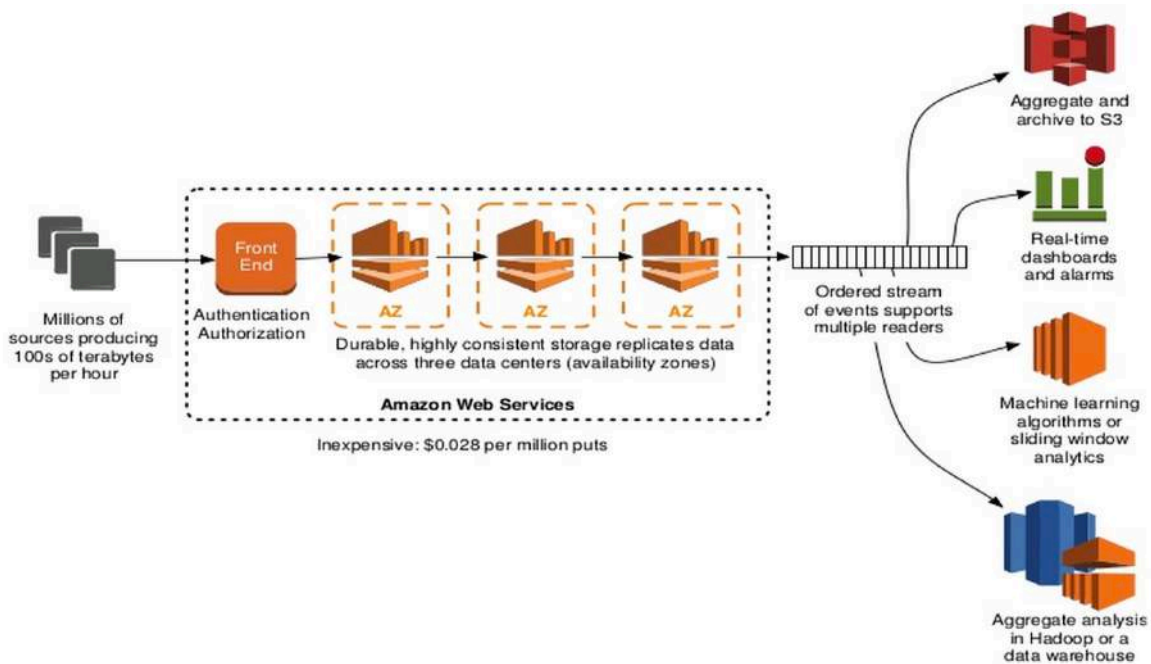
## App Services

-  **AppStream**  
Low Latency Application Streaming
-  **CloudSearch**  
Managed Search Service
-  **Elastic Transcoder**  
Easy-to-use Scalable Media Transcoding
-  **SES**  
Email Sending Service
-  **SQS**  
Message Queue Service
-  **SWF**  
Workflow Service for Coordinating Application Components

## Applications

-  **WorkSpaces**  
Desktops in the Cloud
-  **Zocalo**  
Secure Enterprise Storage and Sharing Service





## Real-time Ingest

Highly Scalable

Durable

Elastic

Replay-able Reads



## Continuous Processing FX

Elastic

Load-balancing incoming streams

Fault-tolerance, Checkpoint / Replay

Enable multiple processing apps in parallel

---

**Managed Service**

---

**Low end-to-end latency**

---

**Enable data movement into Stores/ Processing Engines**

---

# AWS Lambda: A compute service that runs your code in response to events

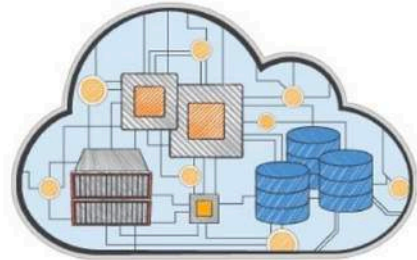
Lambda functions: Stateless, event-driven code execution

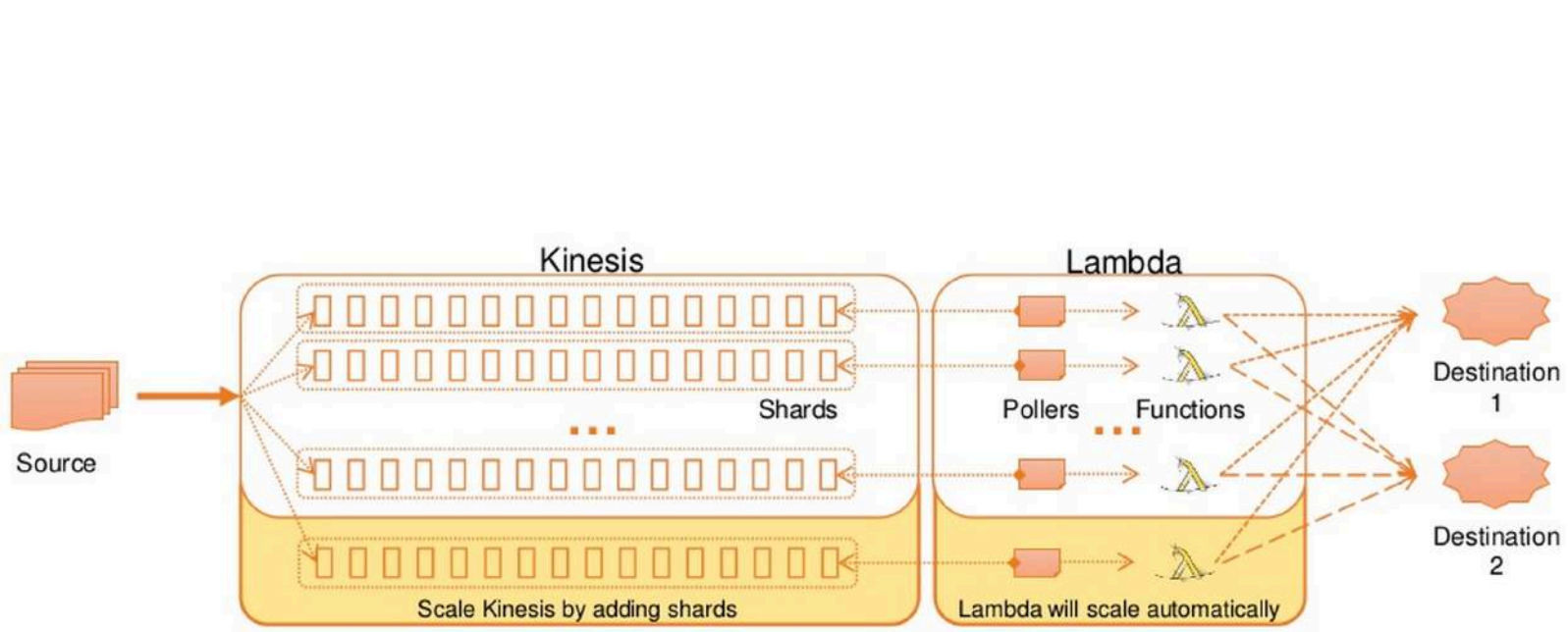
Triggered by events:

- Put to an Amazon S3 bucket
- Record in an Amazon Kinesis stream
- Direct sync and async invocations

Makes it easy to

- Build back-end services that perform at scale
- Perform data-driven auditing, analysis, and notification





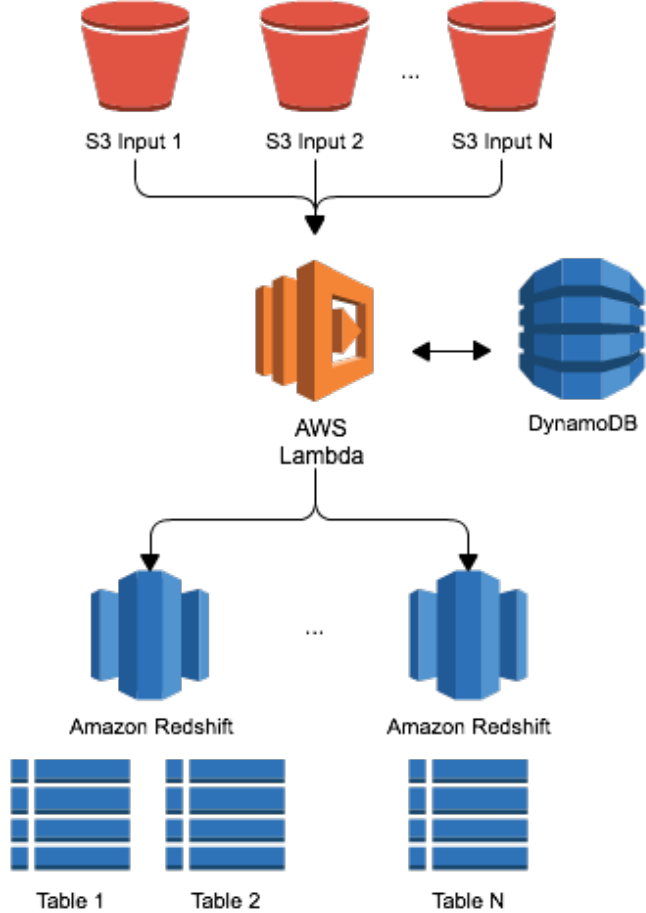
# Performance tuning Kinesis as an event source

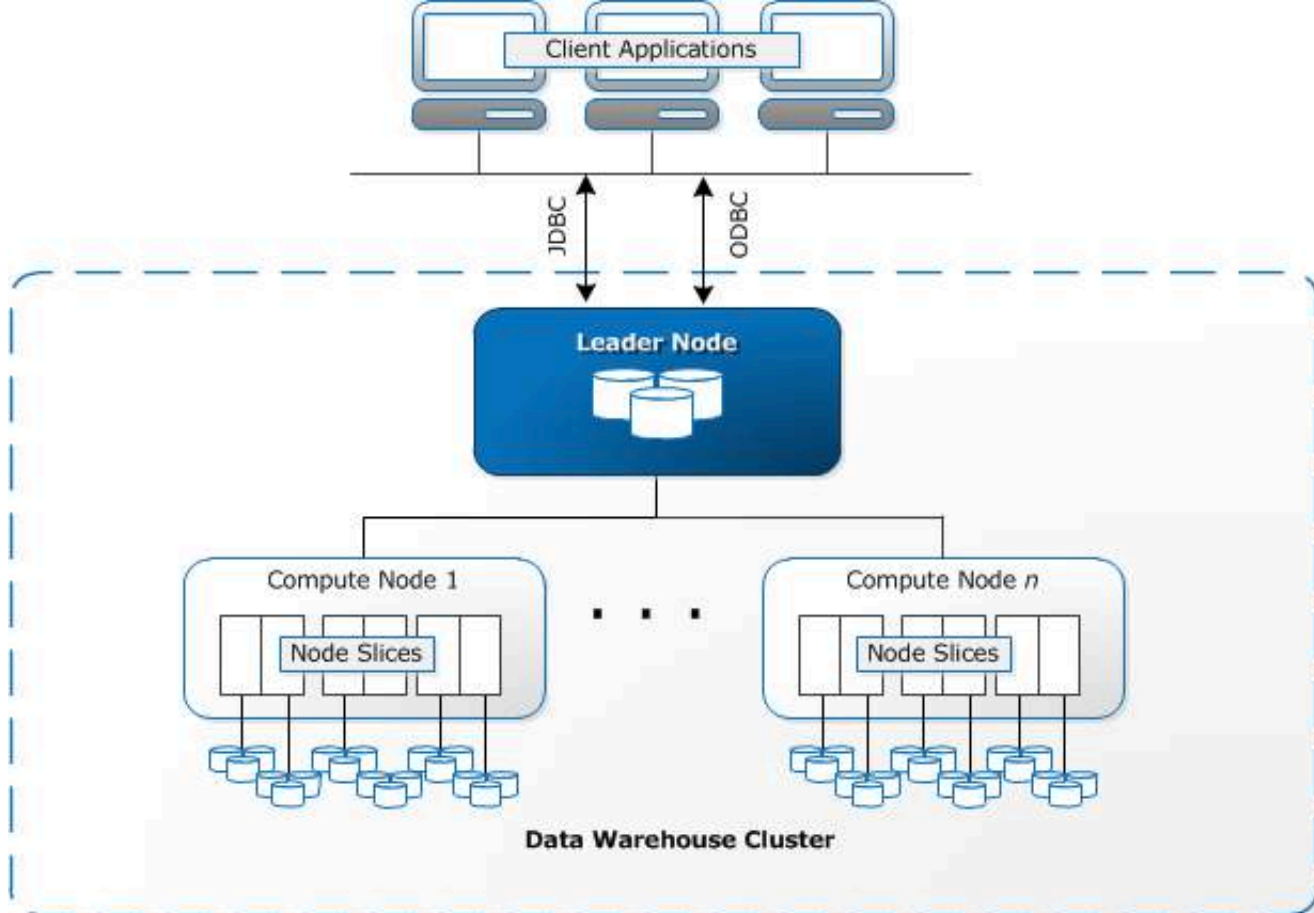
The image shows a configuration interface for an AWS Lambda function with a Kinesis event source. The configuration includes the following fields:

- Event source type:** Kinesis
- Kinesis stream:** demo-Kinesis
- Batch size:** 100
- Starting position:** Trim horizon

The 'Batch size' and 'Starting position' fields are highlighted with red ovals.

- **Batch size:** Number of records that AWS Lambda will retrieve from Kinesis at the time of invoking your function
- Increasing batch size will cause fewer Lambda function invocations with more data processed per function
- **Starting Position:** The position in the stream where Lambda starts reading
- Set to "Trim Horizon" for ordered processing (FIFO)
- Set to "Latest" for reading most recent data (LIFO)







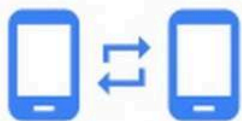
STYLIGHT

Google Cloud

STYLIGHT.COM



## Data on Google Cloud



Capture

Pub/Sub



Store

Storage

SQL

Datastore



Process

Dataflow



Analyze

BigQuery

Dataflow

Open Source Tools



---

### Pub/Sub

Scalable, flexible, and globally available messaging



---

### Dataflow

Stream & batch processing, unified and simplified



---

### BigQuery

Ingest data at 100,000 rows per second

---

## Fully Managed, No-Ops Services



**Web servers:**

- Google Cloud Platform apps
- 3rd-party network apps



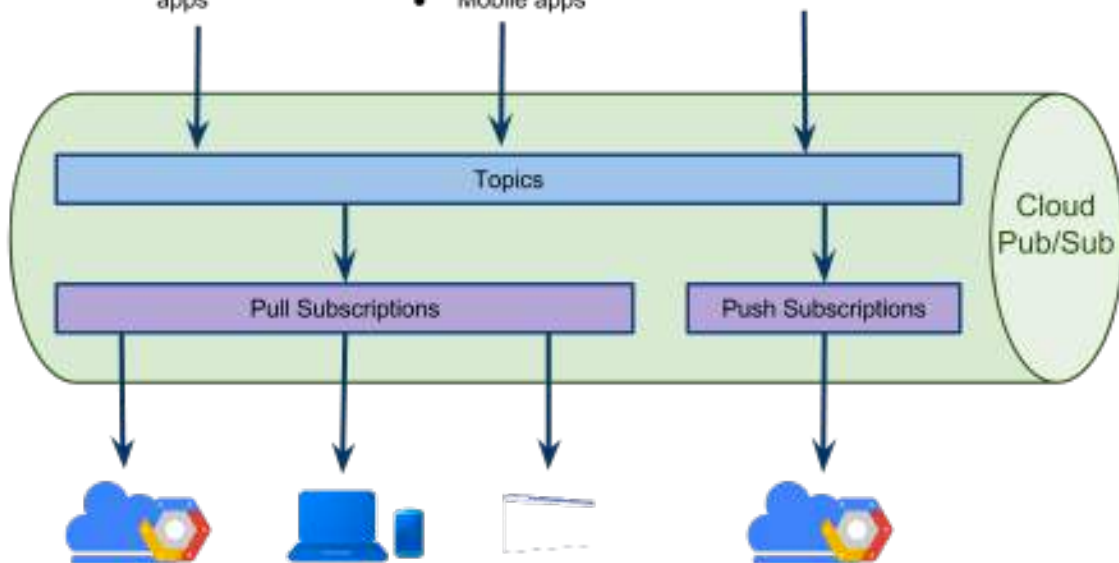
**Native installed apps:**

- Desktop / command-line apps
- Mobile apps

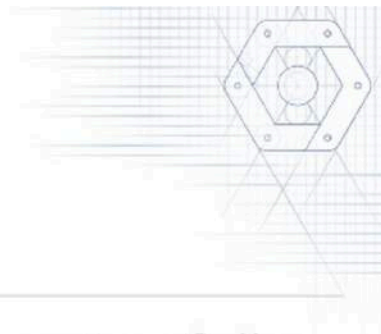


**Browsers:**

- JavaScript / HTML 5 clients



## Cloud Dataflow



---

Cloud Dataflow is a collection of SDKs for **building** batch or streaming parallelized data processing pipelines.

---

---

Cloud Dataflow is a fully managed service for **executing** optimized parallelized data processing pipelines.

---

```
Pipeline{
```

```
  Who => Inputs
```

```
  What => Transforms
```

```
  Where => Windows
```

```
  When => Watermarks + Triggers
```

```
  To => Outputs
```

```
}
```

<- At once guarantee (modulo completeness thresholds)

<- GCS, Pub/Sub, BigQuery, w/Avro, XML, JSON, ...

<- Aggregations, Filters, Joins, ...

<- Time space Fixed, Sliding, Sessions, ...

<- Correctness

<- GCS, Pub/Sub, BigQuery, ...

# Benefits of Cloud Dataflow

## › No Ops - truly elastic data processing for the cloud

- On demand resource allocation w/intelligent auto-scaling
- Automated worker lifetime management
- Automated work optimization

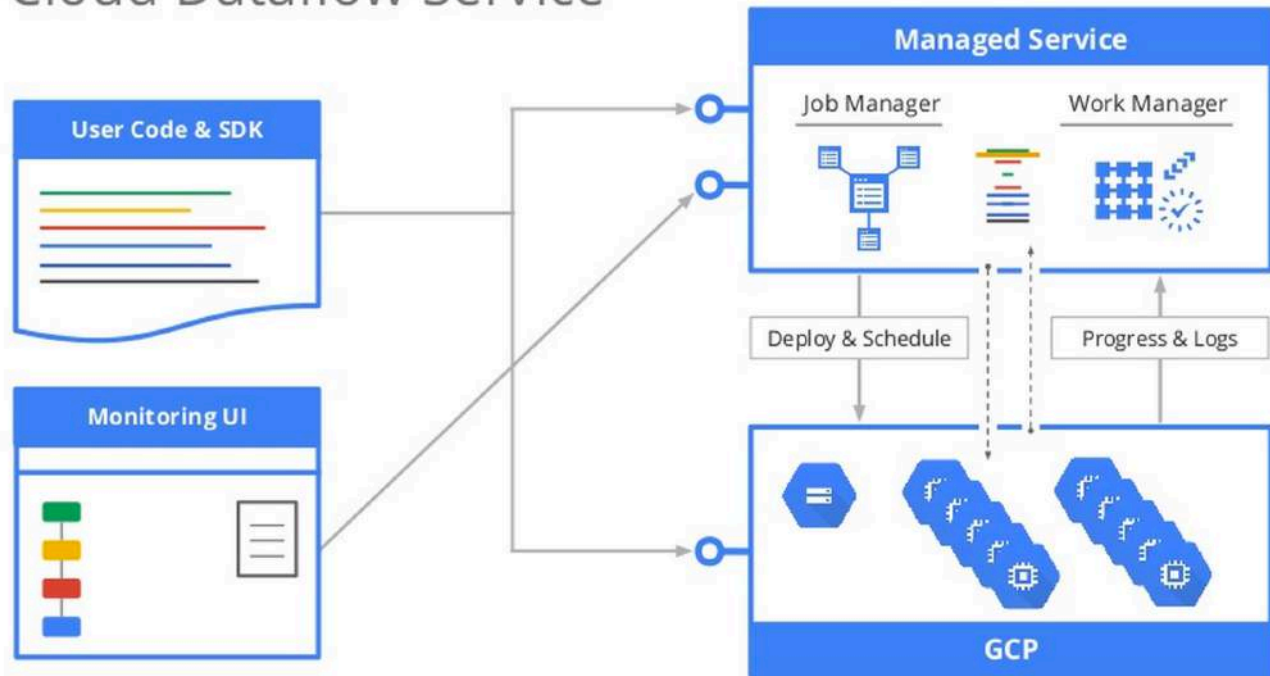
## › Unified model - for batch & stream based processing

- Functional programming model
- Fine grained correctness primitives

## › Open sourced SDK @ github

- Java 7 today @ /GoogleCloudPlatform/DataflowJavaSDK
- Python 2 in progress
- Scala @/darkjh/scalaflow & /jhlch/scala-dataflow-dsl
- Spark runner@ /cloudera/spark-dataflow
- Flink runner @ /dataArtisans/flink-dataflow

# Cloud Dataflow Service





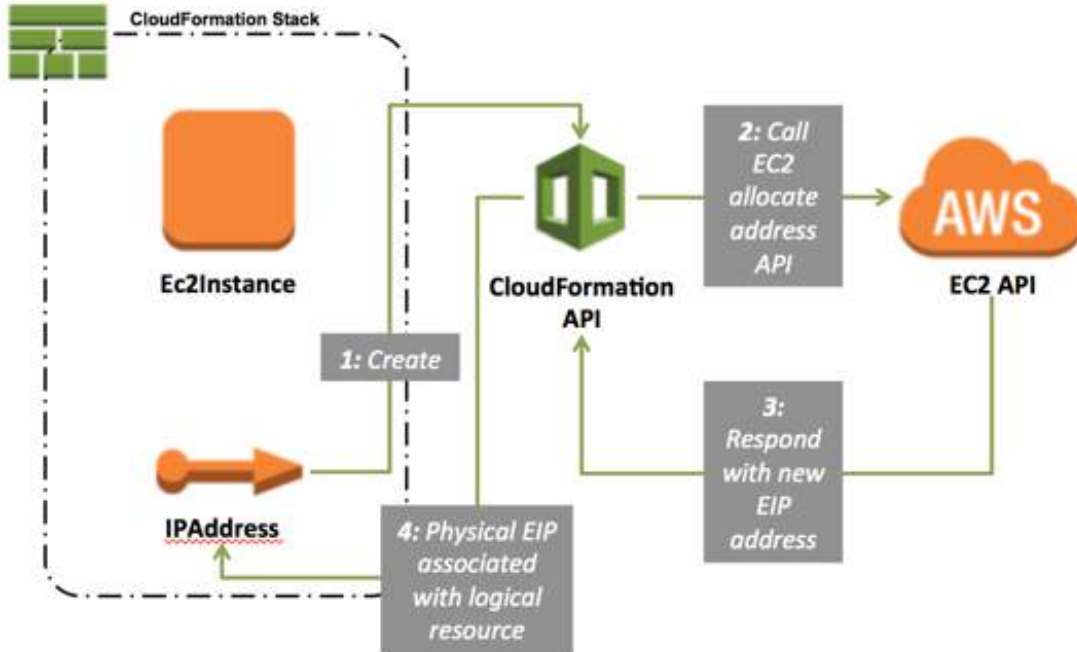
STYLIGHT

# Tips, tricks and best practices

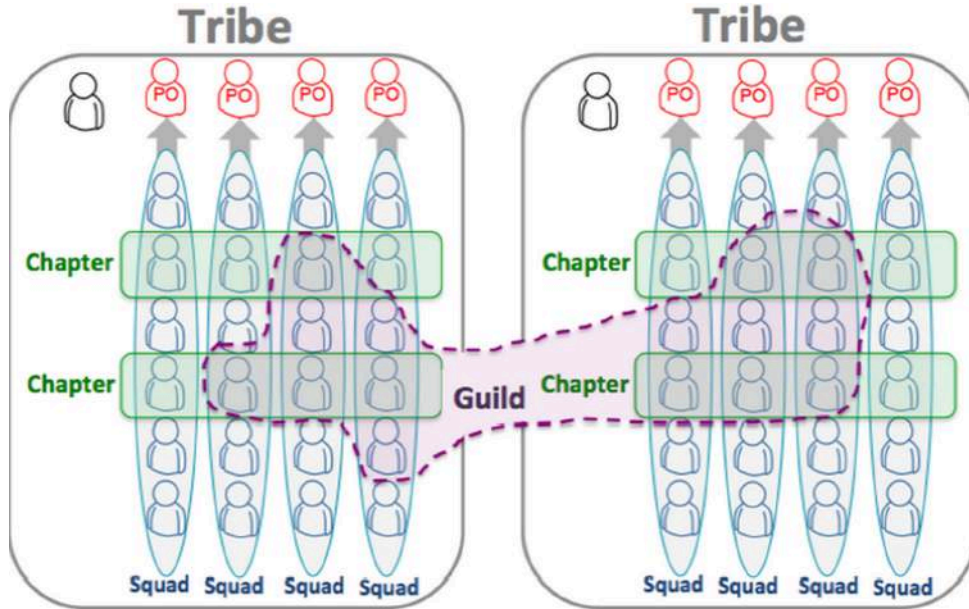
STYLIGHT.COM



## CLOUDFORMATION STACK ON AMAZON AWS



MAKE YOUR TEAM THE OWNERS OF PRODUCT AND DATA



Show 7d Feb 19, 7:00PM - Feb 26, 7:00PM



855 matching events from Feb 19, 7:00PM - Feb 26, 7:00PM



Leave a status update...

Post



**Datasource has been refreshed** #app:tableau\_refresh #source:busdev\_rank\_shop\_values #tableau\_vm

Refreshed BUSDEV rank shop values datasource

∞ · 2 hours ago · [Add comment](#) · [Lower priority](#)



**Refresh finished** #app:tableau\_refresh #source:cg\_pdots\_can\_go\_live #tableau\_vm

Refreshed 613 rows in 63.305000 second

∞ · 5 hours ago · [Add comment](#) · [Lower priority](#)



**Refresh started** #app:tableau\_refresh #department:tam-i18n #source:cg\_pdots\_can\_go\_live #tableau\_vm

Refreshing CG pdots can go live datasource

∞ · 5 hours ago · [Add comment](#) · [Lower priority](#)



**Refresh finished** #app:tableau\_refresh #source:tam\_i18n\_unchecked\_brands #tableau\_vm

Refreshed 18439 rows in 41.626000 second

∞ · 5 hours ago · [Add comment](#) · [Lower priority](#)



**Refresh started** #app:tableau\_refresh #department:tam-i18n #source:tam\_i18n\_unchecked\_brands #tableau\_vm

Refreshing TAM i18n Unchecked brands datasource

∞ · 5 hours ago · [Add comment](#) · [Lower priority](#)



**Refresh finished** #app:tableau\_refresh #source:tam\_i18n\_products\_data #tableau\_vm

Refreshed 826 rows in 96.572000 second

∞ · 5 hours ago · [Add comment](#) · [Lower priority](#)



**Refresh started** #app:tableau\_refresh #department:tam-i18n #source:tam\_i18n\_products\_data #tableau\_vm

Refreshing TAM i18n Products data datasource

∞ · 5 hours ago · [Add comment](#) · [Lower priority](#)

All Sources

syslog.host:"WIN-R2T9EJJ3KBT"

Last day

Search

★



Applied Filters: json.SourceName : Python ETL ✕ AND syslog.severity : Error ✕

Field Explorer

17 events

Feb 25 - Feb 26



View:

Collapsed Events



Sorting:

Descending



2015-02-26 09:29:00.000 UTC

View Surrounding Events

host : WIN-R2T9EJJ3KBT log type : syslog log type : json tag : tableau\_vm

LogglyNotifications:

syslog:

```
severity: Error
appName: Python_ETL
timestamp: 2015-02-26T09:29:00.000000+00:00
facility: user-level messages
pid: 0
priority: 11
host: WIN-R2T9EJJ3KBT
```

json:

```
EventID: 1
ProcessID: 0
EventTime: 2015-02-26 09:29:00
Task: 0
Severity: ERROR
SourceModuleType: im_msvistalog
SourceName: Python ETL
EventType: ERROR
SourceModuleName: eventlog
Hostname: WIN-R2T9EJJ3KBT
Opcode: Info
ThreadID: 0
Keywords: 36028797018963968
SeverityValue: 4
Message:
```

Thanks for your  
attention!

# References

- [Google's Cloud Pub/Sub Real-Time Messaging Service Is Now In Public Beta](#)
- [Streaming Data Processing with Amazon Kinesis and AWS Lambda](#)
- [Google Cloud Dataflow Two Worlds Become a Much Better One](#)



STYLIGHT

Sergii Khomenko

Data Scientist

STYLIGHT GmbH

[sergii.khomenko@stylight.com](mailto:sergii.khomenko@stylight.com)

@lc0d3r

Nymphenburger Straße 86

80636 Munich, Germany

STYLIGHT.COM

<http://dahho.am/>

DAH. AM

12.06.15 Munich Developer  
Conference

BY STYLIGHT