

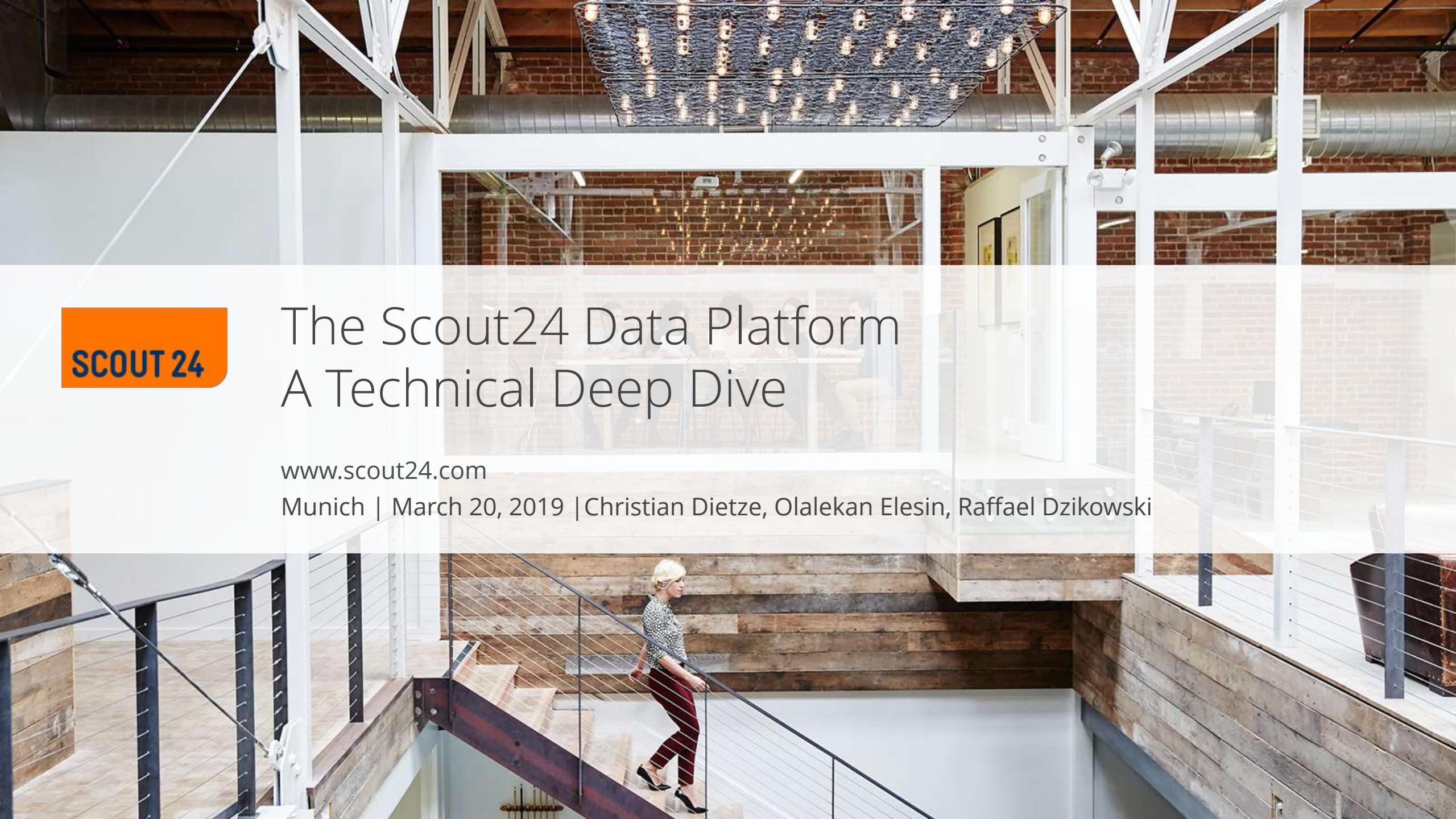


# The Scout24 Data Platform

## A Technical Deep Dive

[www.scout24.com](http://www.scout24.com)

Munich | March 20, 2019 | Christian Dietze, Olalekan Elesin, Raffael Dzikowski



# Scout24 AG

- MDAX
- € 531.7 million revenue (2018)

2

Major Household Brand Names



5

Core Geographies  
and an overall presence  
in 18 countries

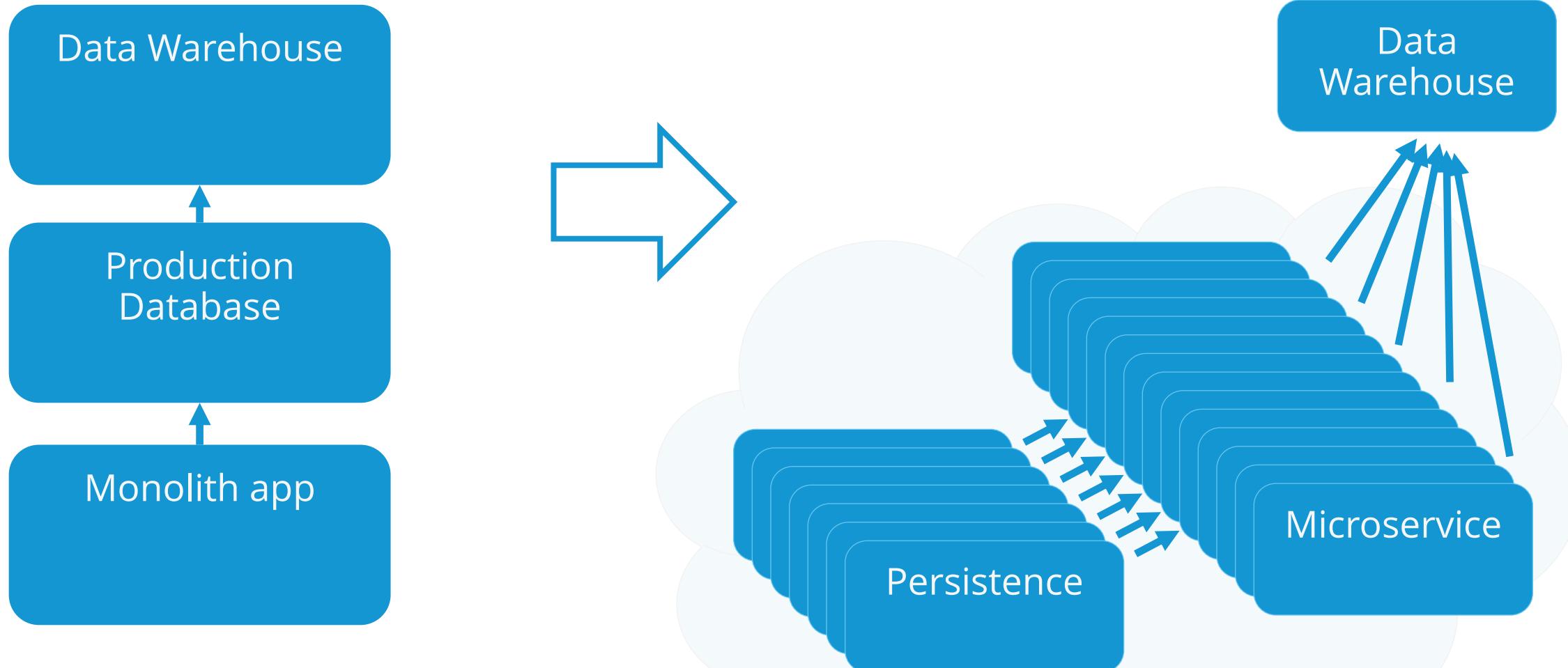
80m

Household Reach



SCOUT24

# Our technical evolution





Our data warehouse was a bottle neck

# Scout24 wants to become a truly data-driven company

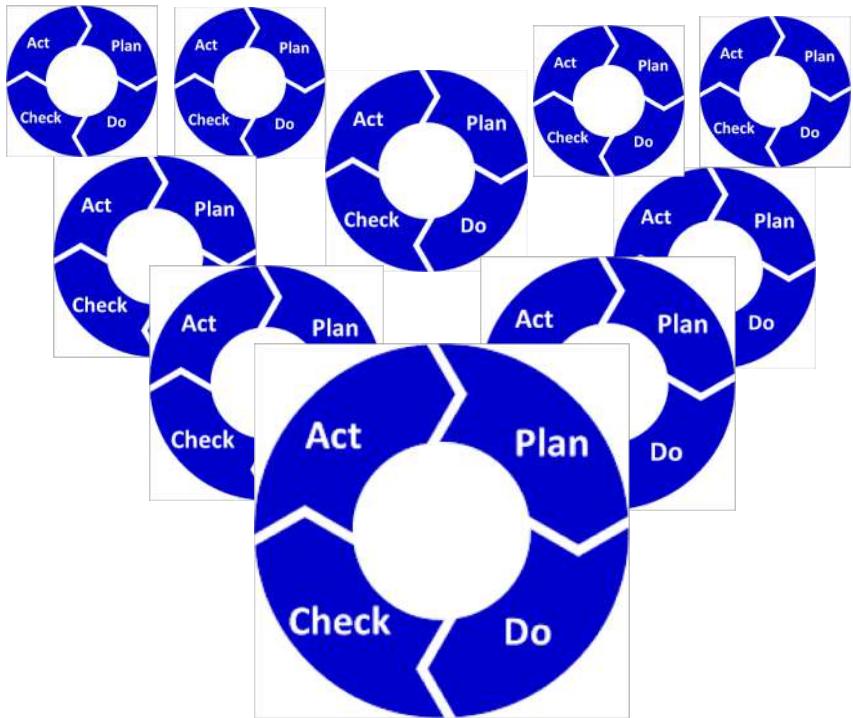
Fast & easy data-driven  
product development...

...supported by  
Data & Analytics



# Scout24 wants to become a truly data-driven company

Everywhere in the company...

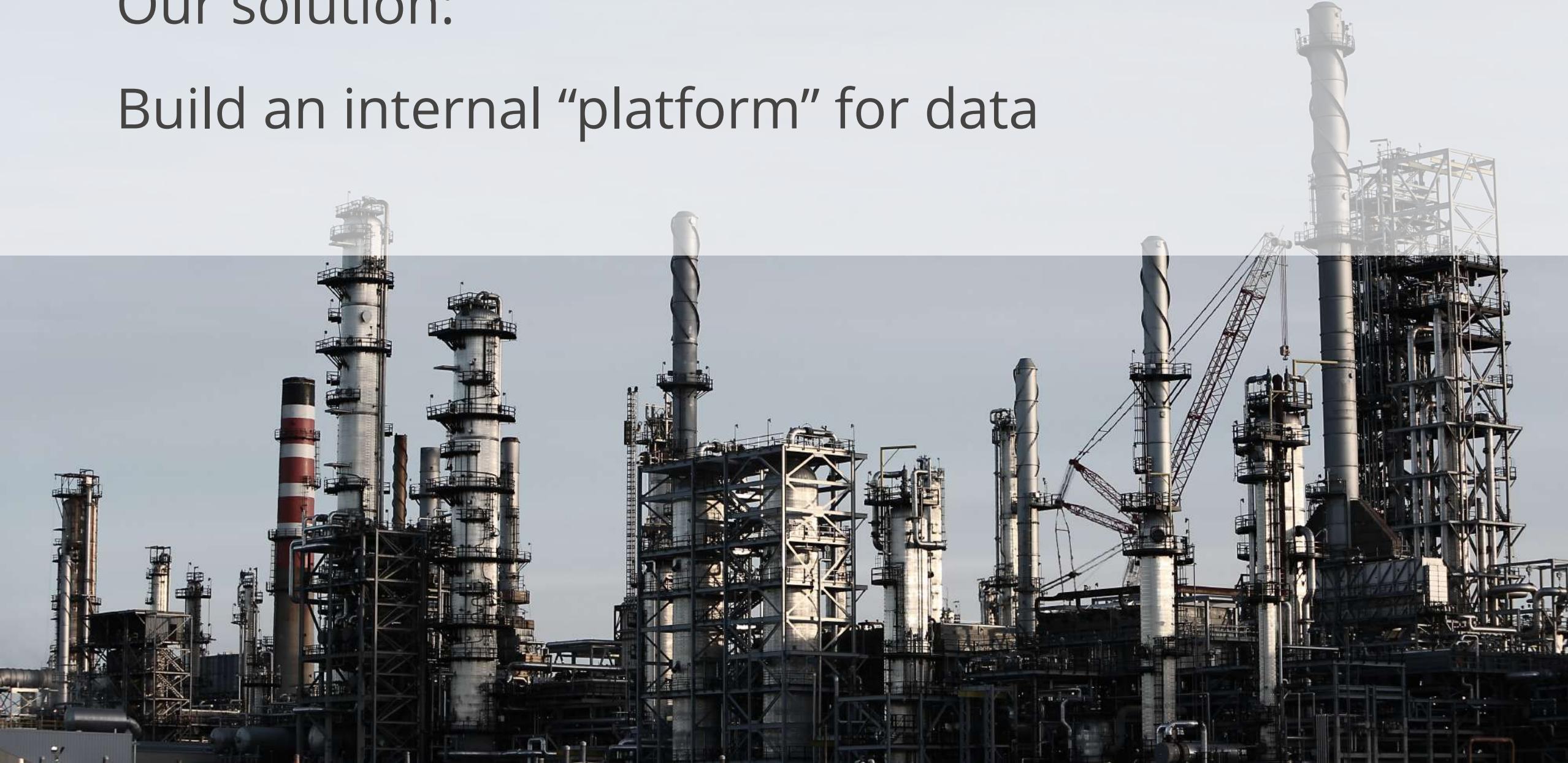


...without bloating up  
Data & Analytics



# Our solution:

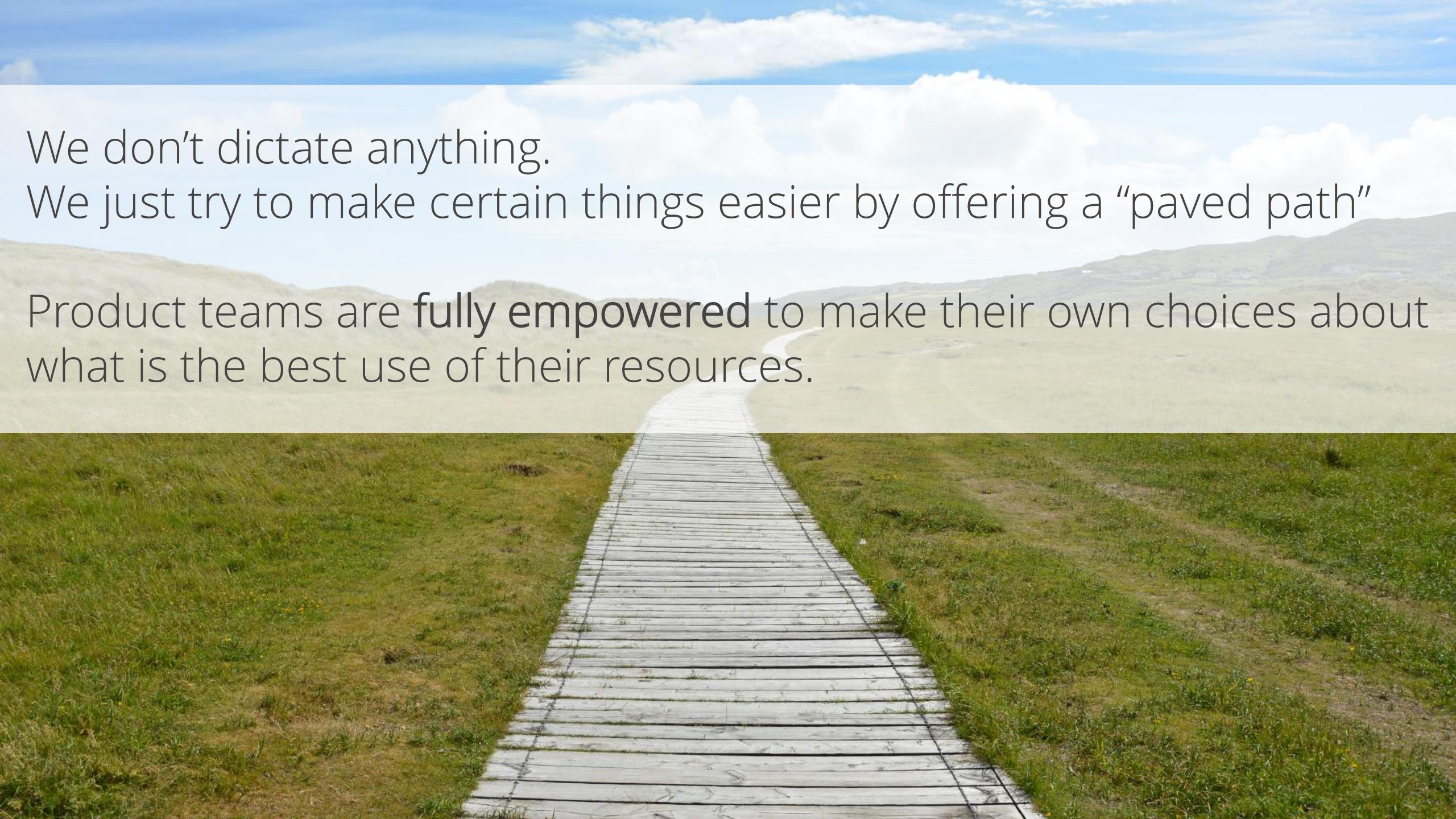
Build an internal “platform” for data



# What is the Scout24 Data “Platform”?

# We think of our Data Platform as a Product

- Just like AWS, Salesforce, etc. – the platform is a **generic layer** upon which Scout24's products can be built
- BUT, we have a very, very small number of customers.
- That means, product teams get **personalized support** and there is lots of **opportunity for collaboration**.

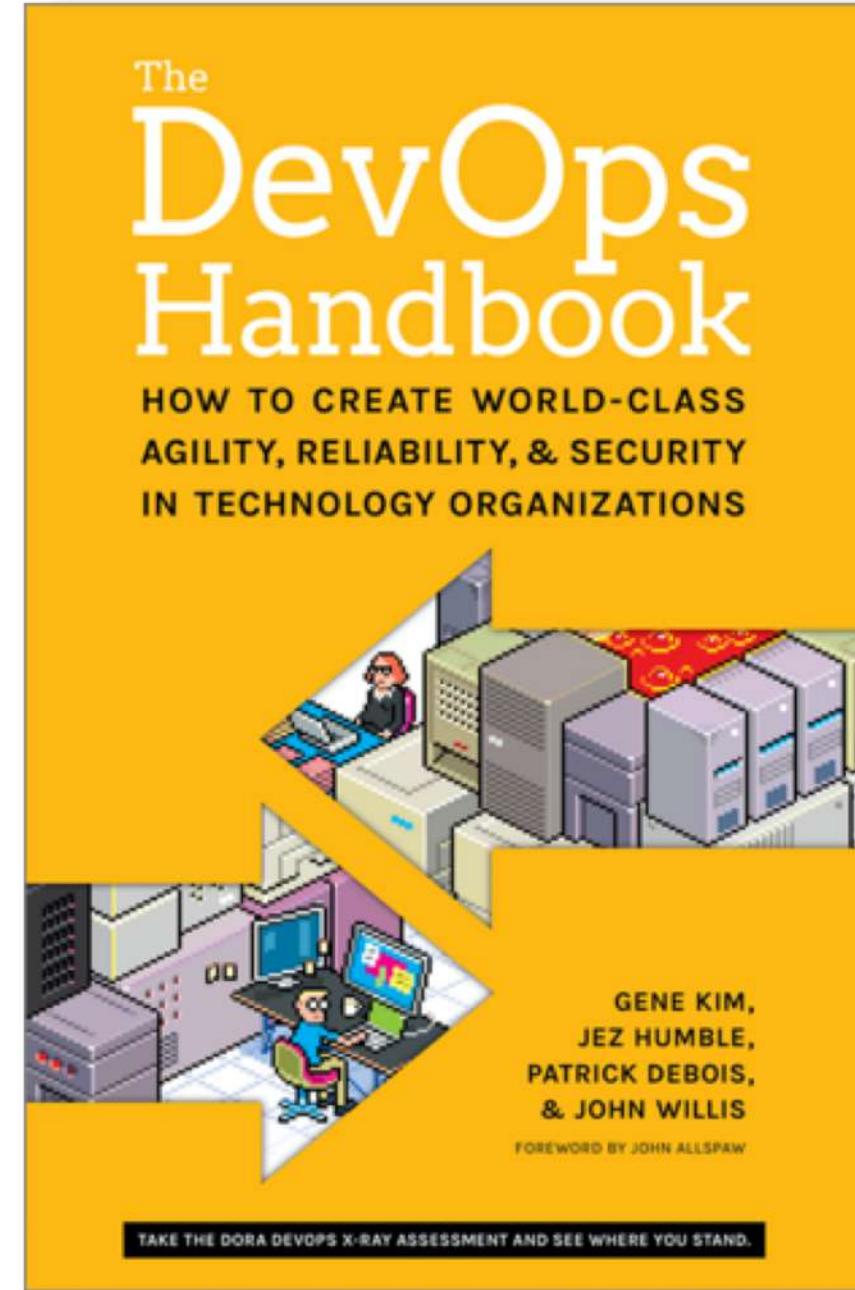
A photograph of a long, straight wooden boardwalk path made of light-colored planks. The path leads from the foreground into a green, hilly landscape under a bright blue sky with scattered white clouds.

We don't dictate anything.

We just try to make certain things easier by offering a "paved path"

Product teams are **fully empowered** to make their own choices about what is the best use of their resources.

“In almost all cases, we will not mandate that internal team use these platforms and services—these platform teams will have to win over and satisfy their internal customers, sometimes even competing with external vendors.”



# | Guiding principle of the platform

- Autonomy for producers and consumers

Self-service Analytics

Self-service Data Ingestion

Self-service ETL

# Data Landscape Manifesto



# Data Landscape Manifesto

## Data Landscape Manifesto

A federal landscape of data producers and consumers with just enough rules to ensure seamless cooperation without severely impeding autonomy

# Data Landscape Manifesto



## Data Landscape Manifesto

A federal landscape of data producers and consumers with just enough rules to ensure seamless cooperation without severely impeding autonomy

# Data Landscape Manifesto

#1  
Preamble

#2  
Responsibility of  
Data and  
Analytics

## Data Landscape Manifesto

A federal landscape of data producers and consumers with just enough rules to ensure seamless cooperation without severely impeding autonomy

# Data Landscape Manifesto

#1  
Preamble

#2  
Responsibility of  
Data and  
Analytics

#3  
Data Autonomy,  
not Anarchy

## Data Landscape Manifesto

A federal landscape of data pro-  
ducers and consumers with just  
enough rules to ensure seamless  
cooperation without severely  
impeding autonomy

# Data Landscape Manifesto

#1  
Preamble

#2  
Responsibility of  
Data and  
Analytics

#3  
Data Autonomy,  
not Anarchy

## Data Landscape Manifesto

A federal landscape of data pro-  
ducers and consumers with just  
enough rules to ensure seamless  
cooperation without severely  
impeding autonomy

#4  
Producer's  
Responsibility

# Data Landscape Manifesto

#1  
Preamble

#2  
Responsibility of  
Data and  
Analytics

#3  
Data Autonomy,  
not Anarchy

## Data Landscape Manifesto

A federal landscape of data pro-  
ducers and consumers with just  
enough rules to ensure seamless  
cooperation without severely  
impeding autonomy

#4  
Producer's  
Responsibility

#5  
Consumer's  
Responsibility

# Data Landscape Manifesto

#1  
Preamble

#2  
Responsibility of  
Data and  
Analytics

#3  
Data Autonomy,  
not Anarchy

## Data Landscape Manifesto

A federal landscape of data pro-  
ducers and consumers with just  
enough rules to ensure seamless  
cooperation without severely  
impeding autonomy

#4  
Producer's  
Responsibility

#5  
Consumer's  
Responsibility

#6  
Exception:  
Core KPIs

# Data Landscape Manifesto

#1  
Preamble

#2  
Responsibility of  
Data and  
Analytics

#3  
Data Autonomy,  
not Anarchy

## Data Landscape Manifesto

A federal landscape of data pro-  
ducers and consumers with just  
enough rules to ensure seamless  
cooperation without severely  
impeding autonomy

#7  
Transparency  
over  
Continuity

#4  
Producer's  
Responsibility

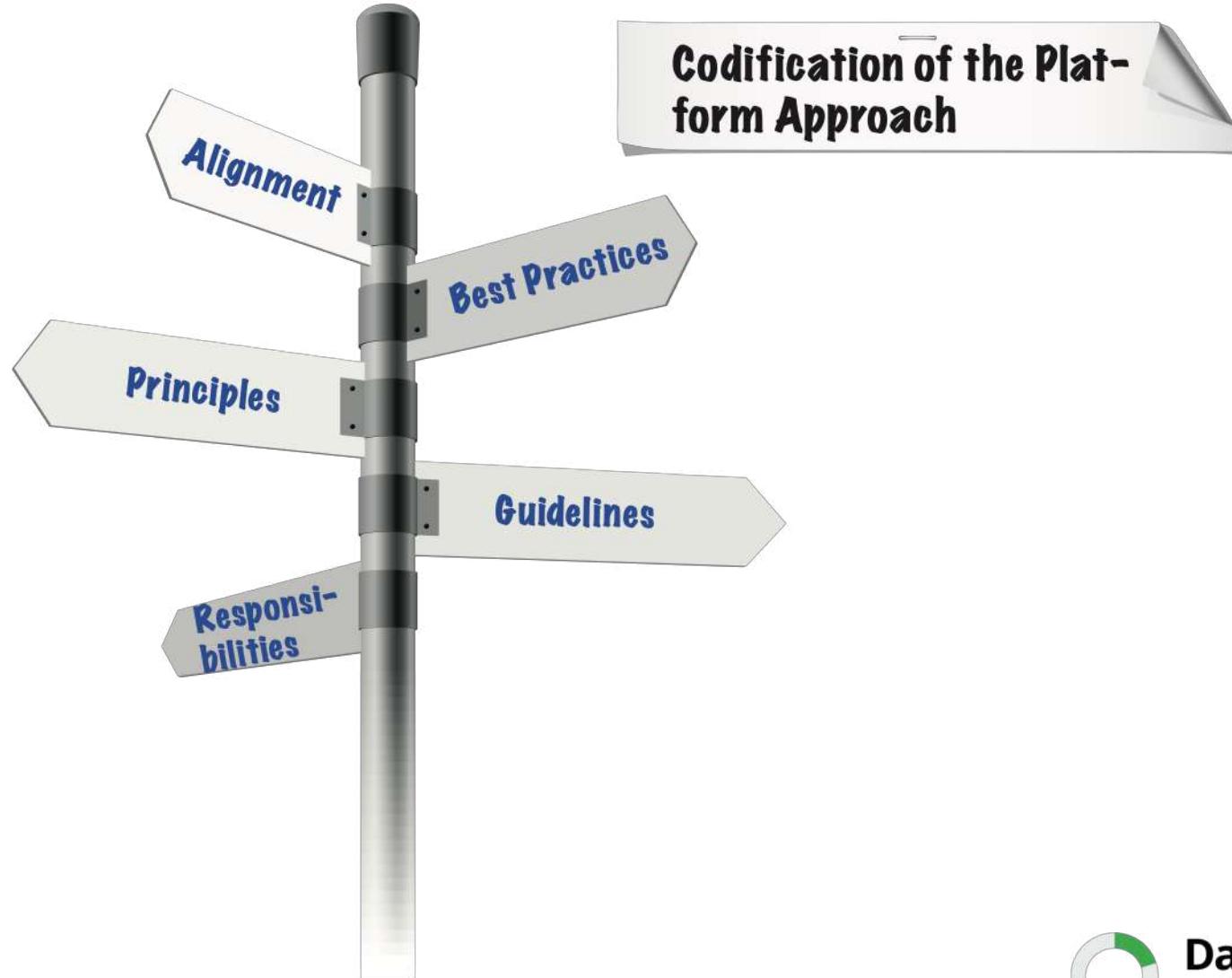
#5  
Consumer's  
Responsibility

#6  
Exception:  
Core KPIs

# Data Landscape Manifesto – Purpose



# Data Landscape Manifesto – Purpose



# Necessary Cultural Changes during Migration from Data Warehouse to Cloud Data Platform

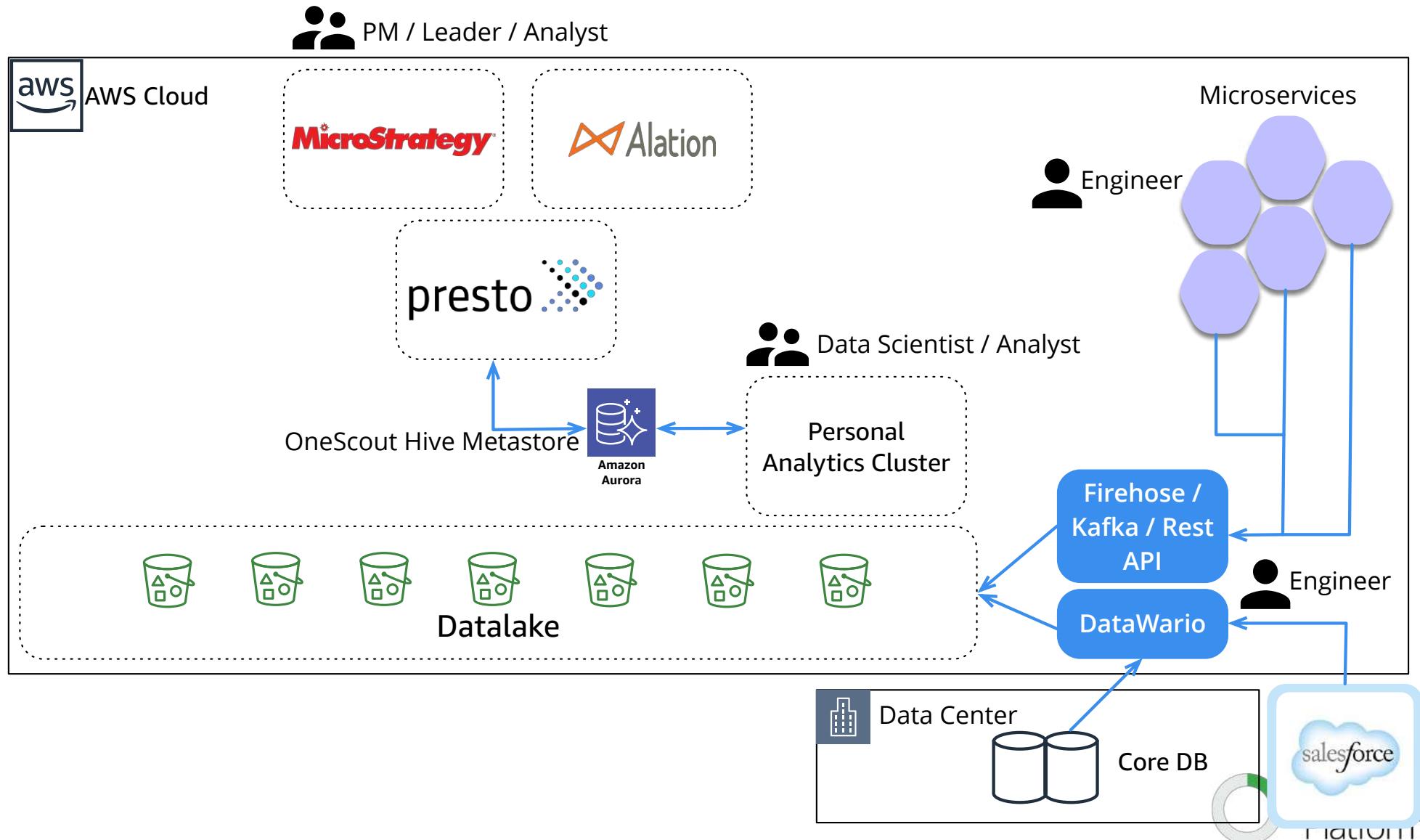
*Dr. Sean Gustafson, Scout24 AG*

*Data Festival 2018*

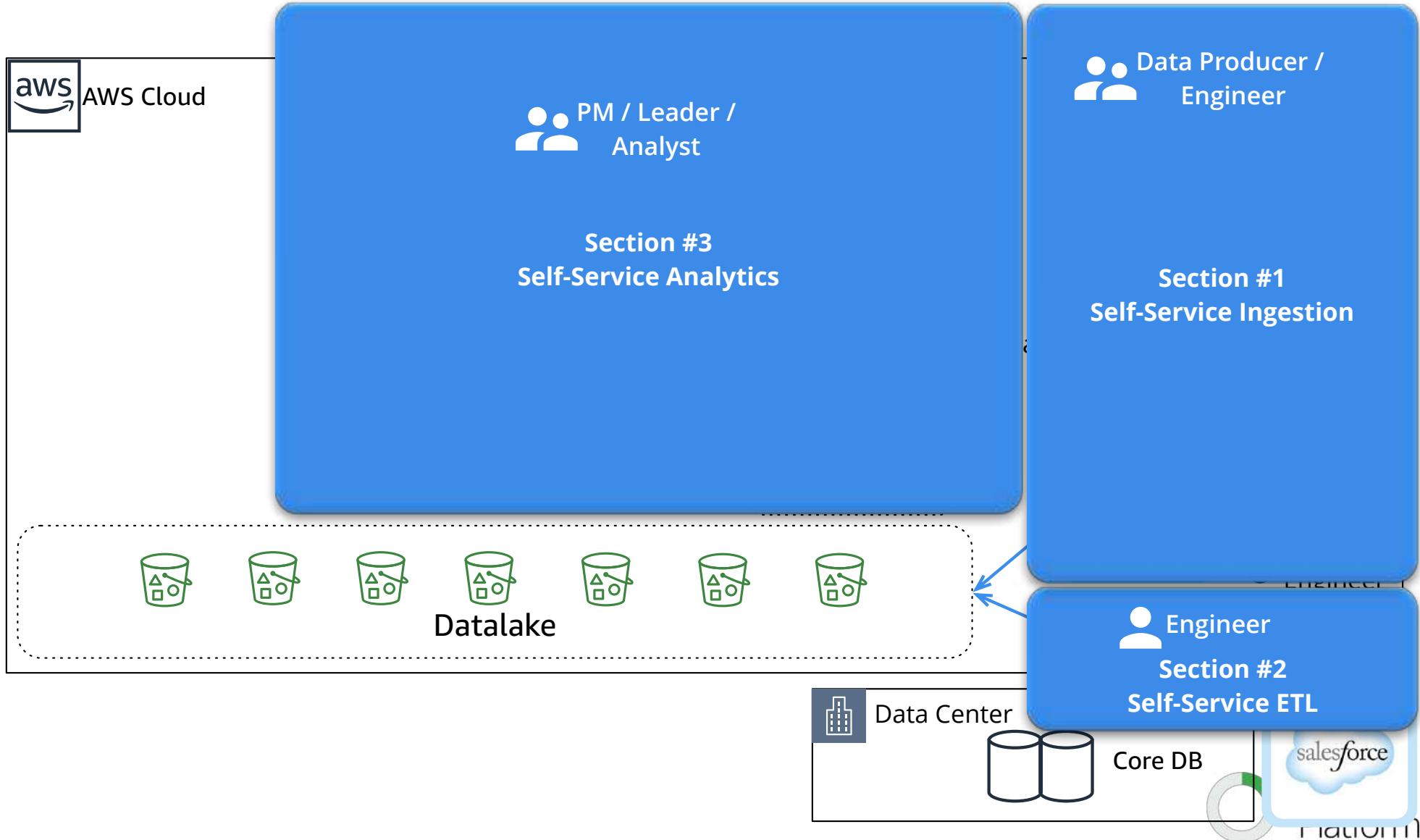


[Watch the Video](#) | [Read the Blog Post on Medium](#)

# 10,000 Foot Architecture Overview



# 10,000 Foot Architecture Overview



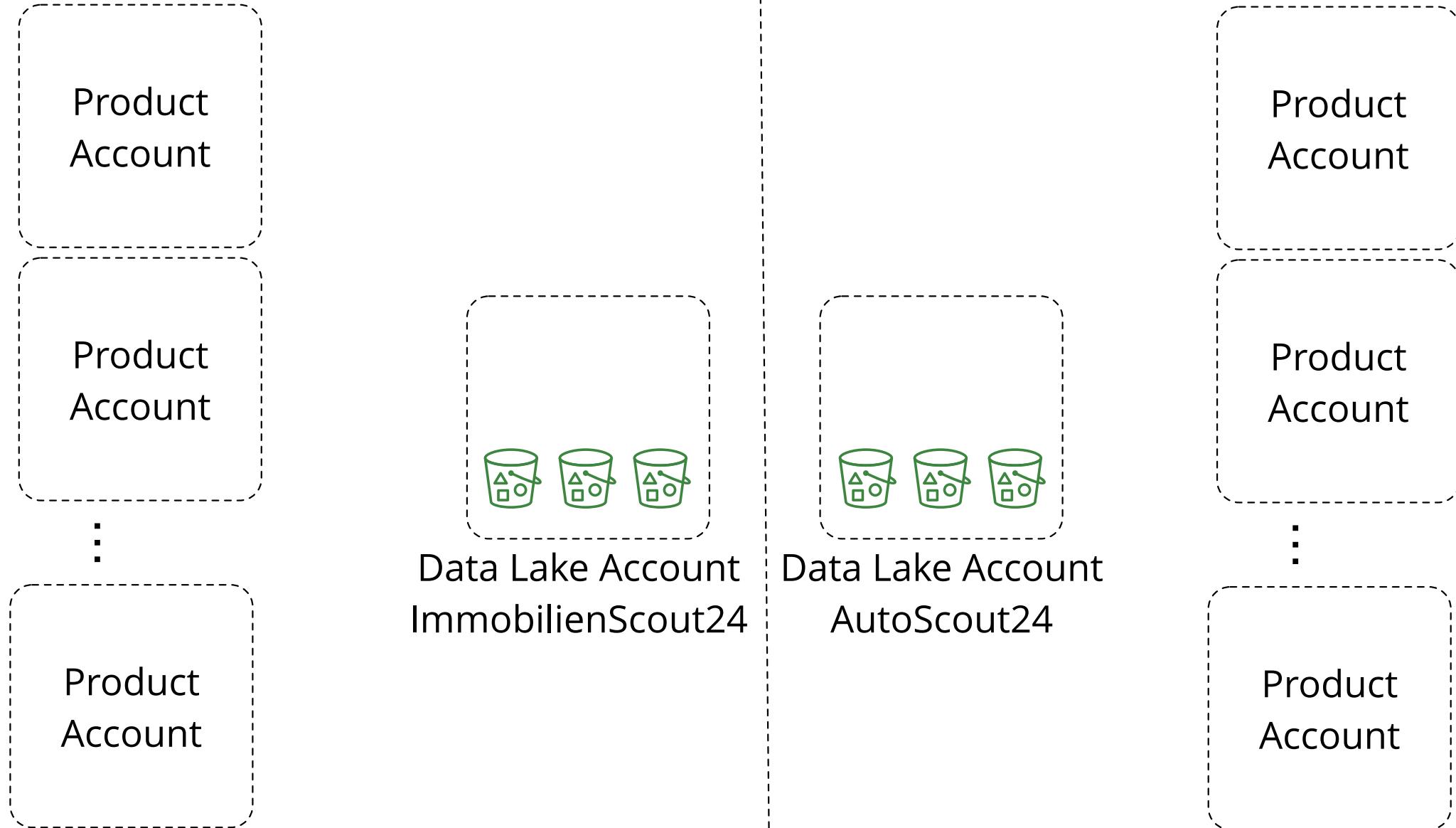
# Self-Service Ingestion



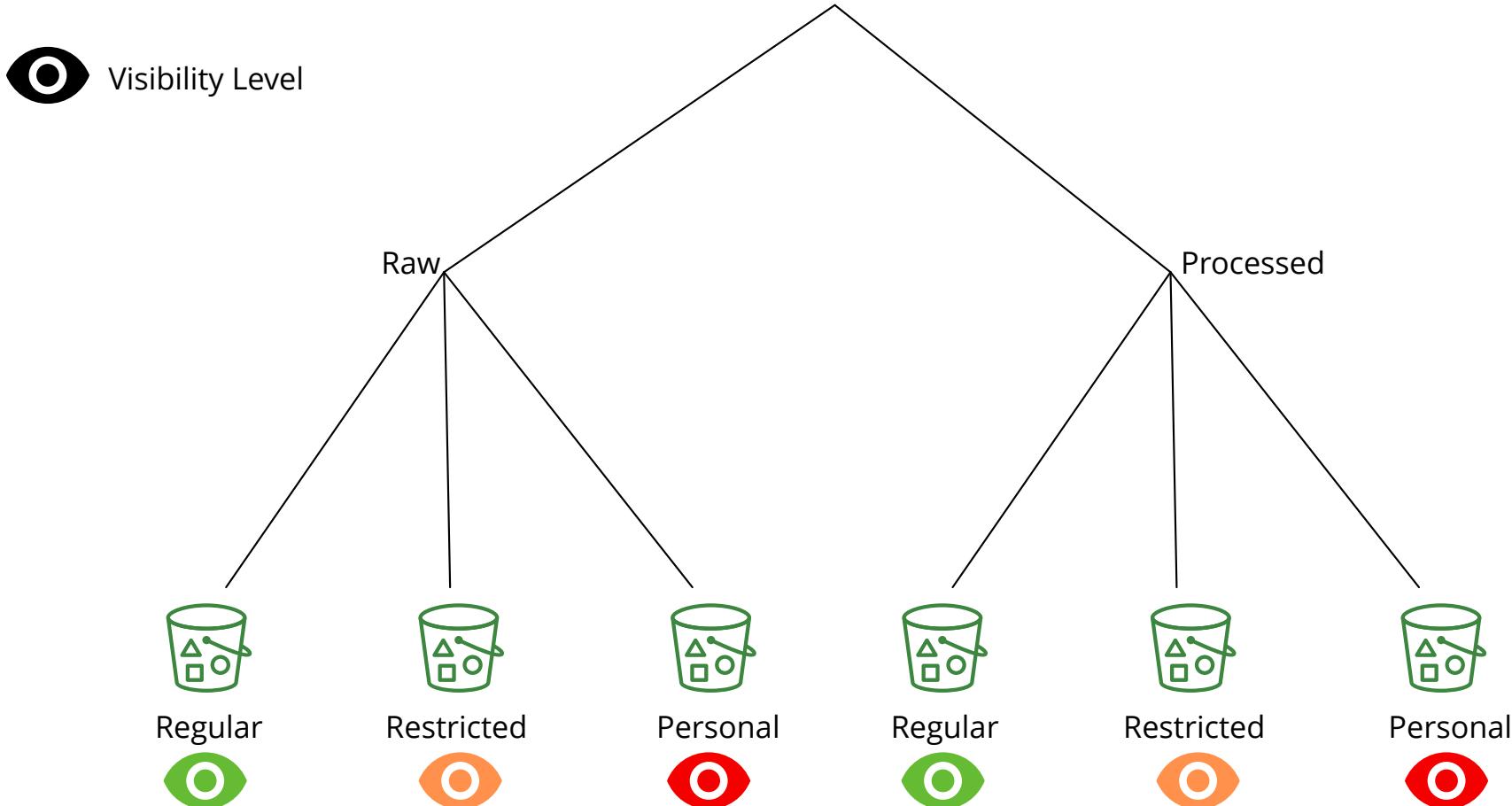
Data  
Platform

SCOUT24

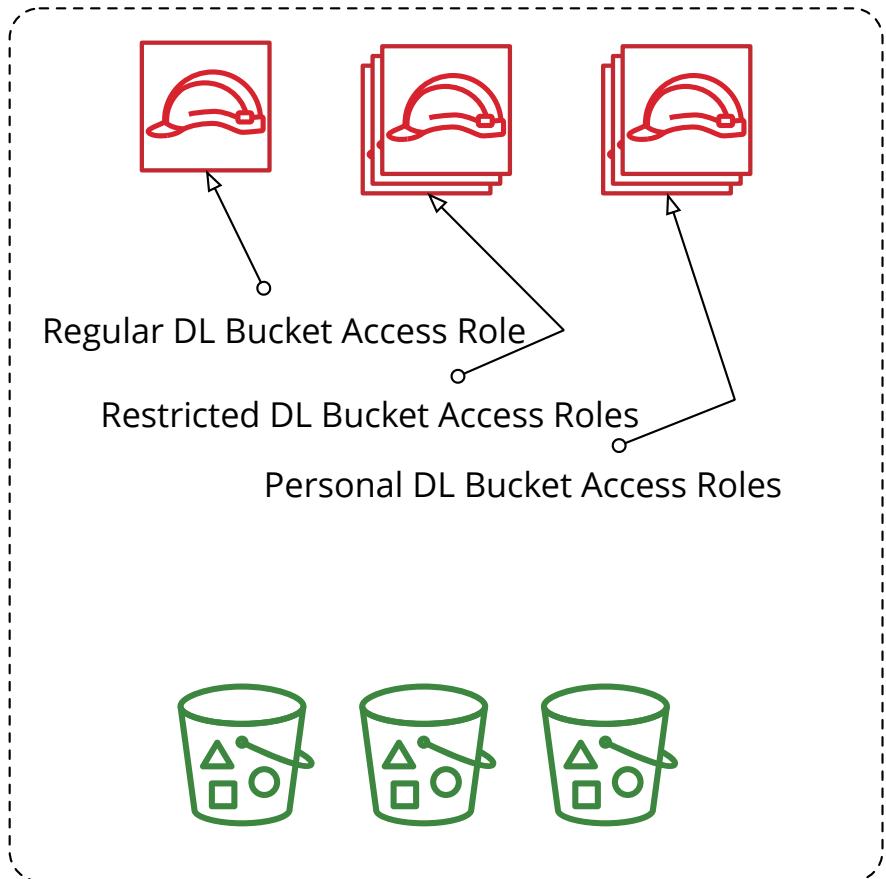
# Multi-Account Setting



# Data Lake Bucket Types



# Data Lake Access Roles



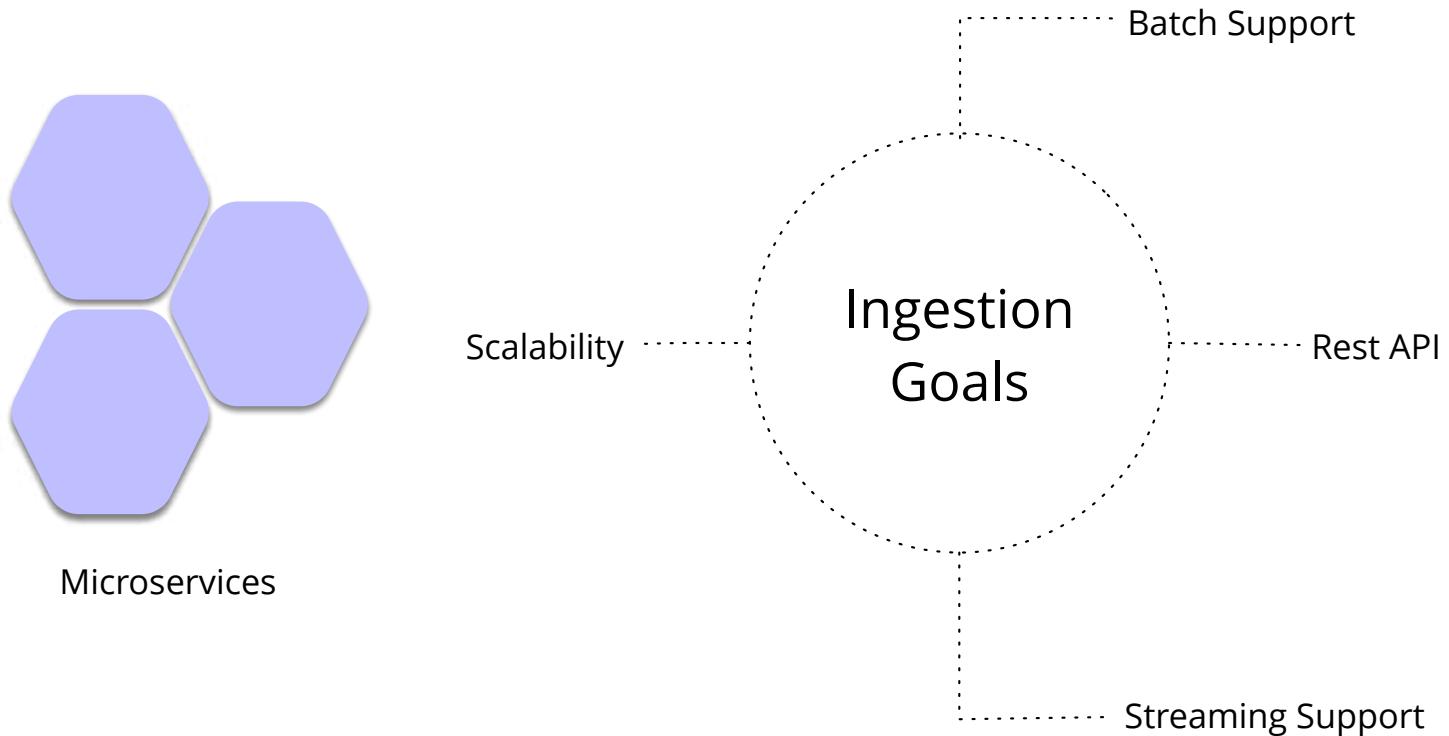
Data Lake Account  
ImmobilienScout24



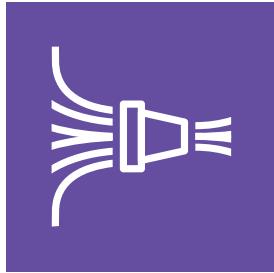
Data Lake Account  
AutoScout24



# Data Lake Ingestion Goals



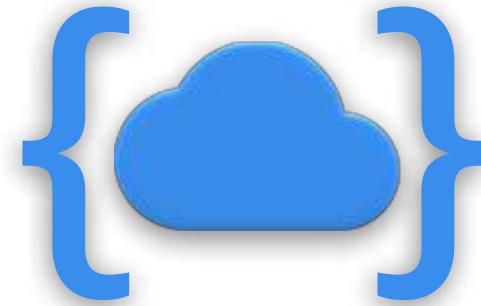
# Ingestion Options



Amazon  
Kinesis Data  
Firehose

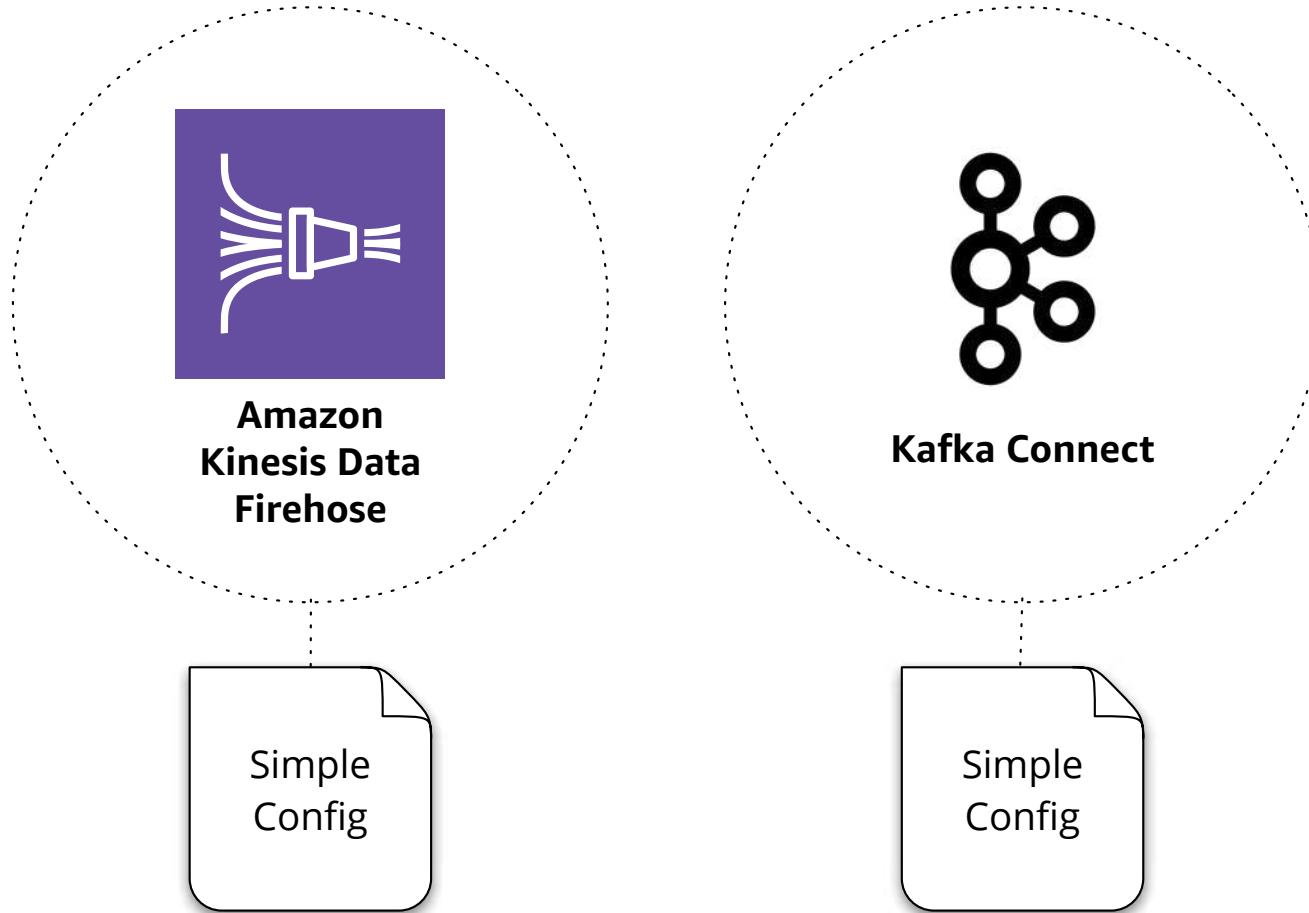


Kafka Connect



Hadoop Rest API

# Ingestion Options



**Hadoop Rest API**



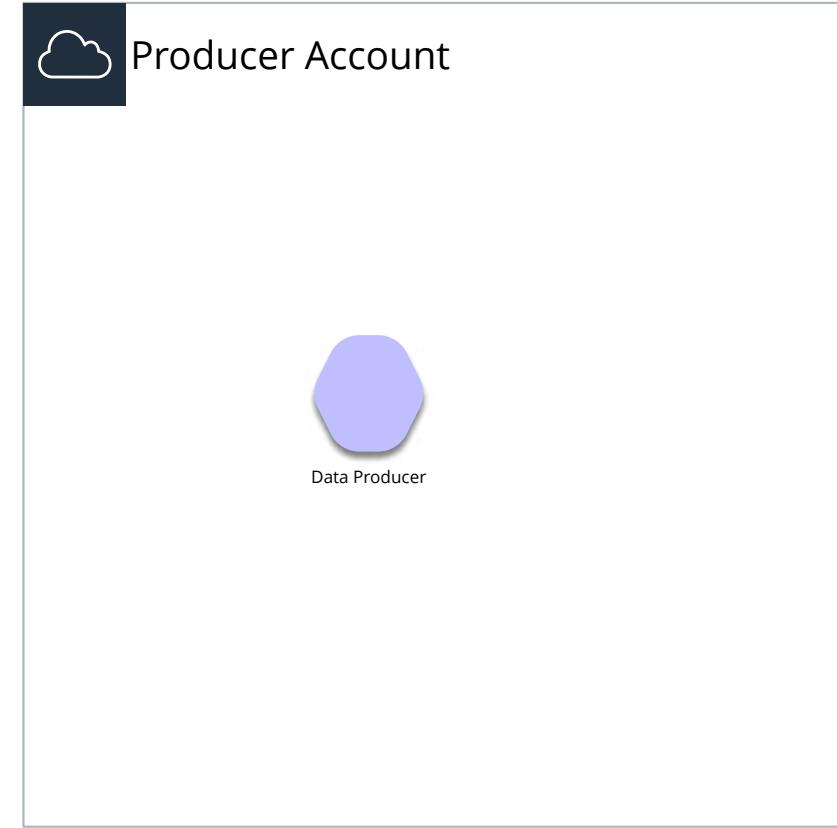
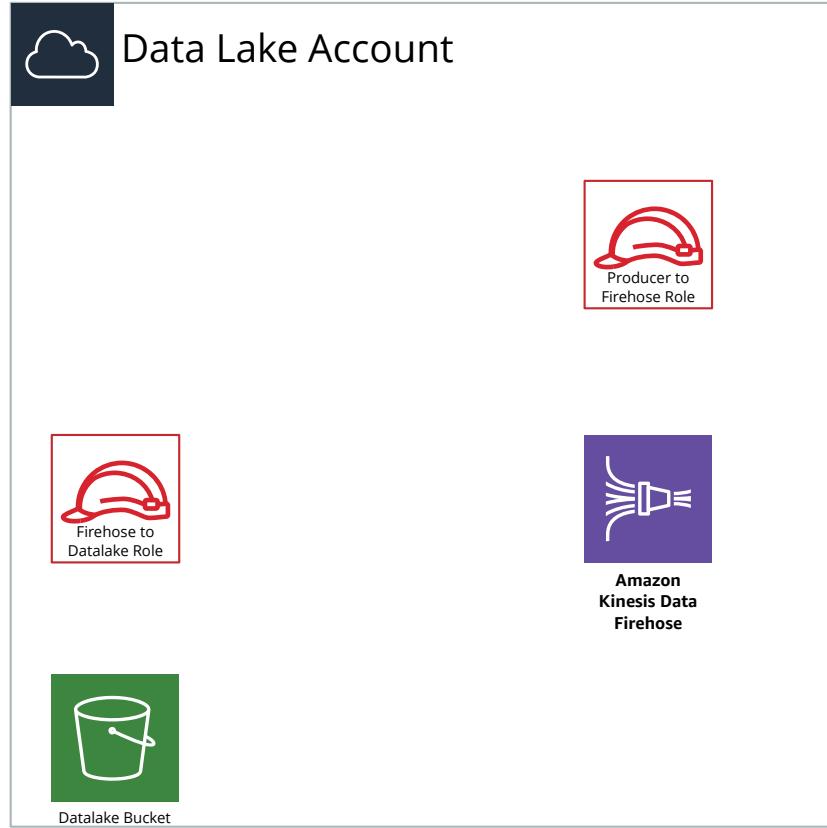
Configuration Free  
**Data**  
Platform

**SCOUT 24**

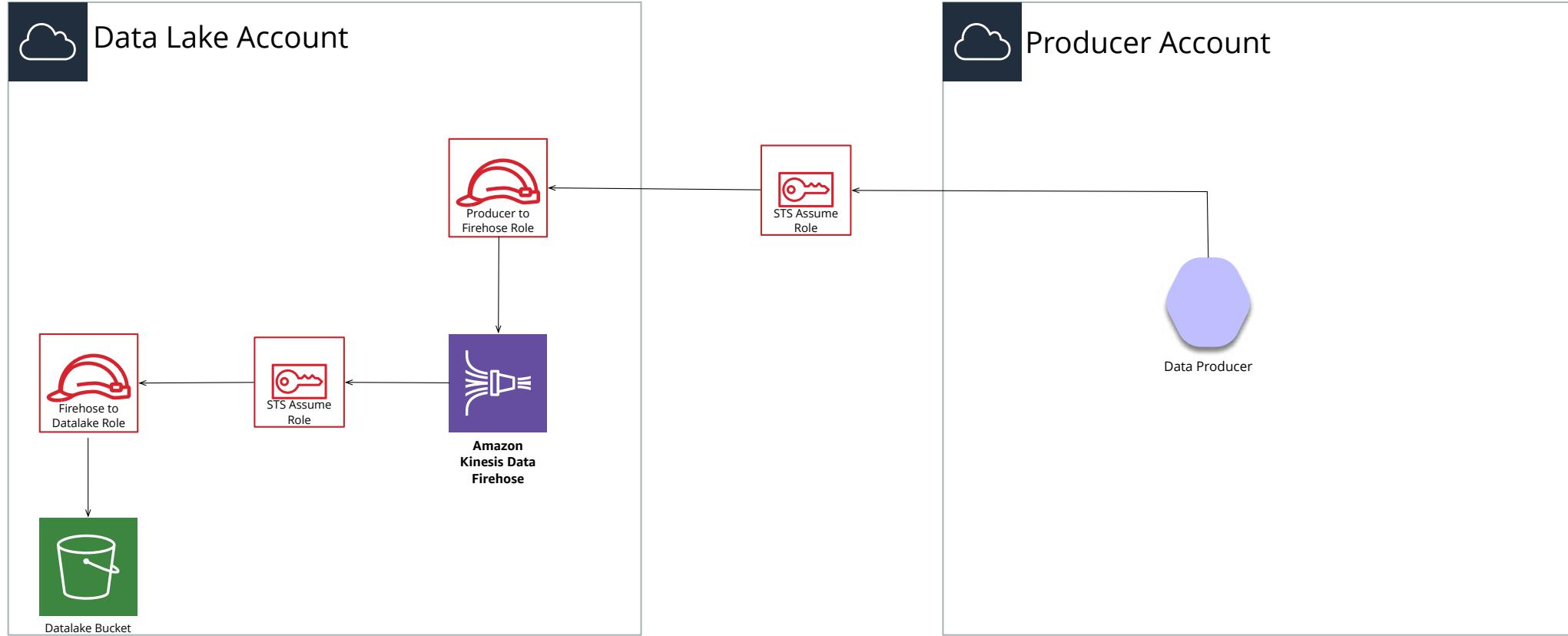
# Kinesis Data Firehose Ingestion Mechanism

---

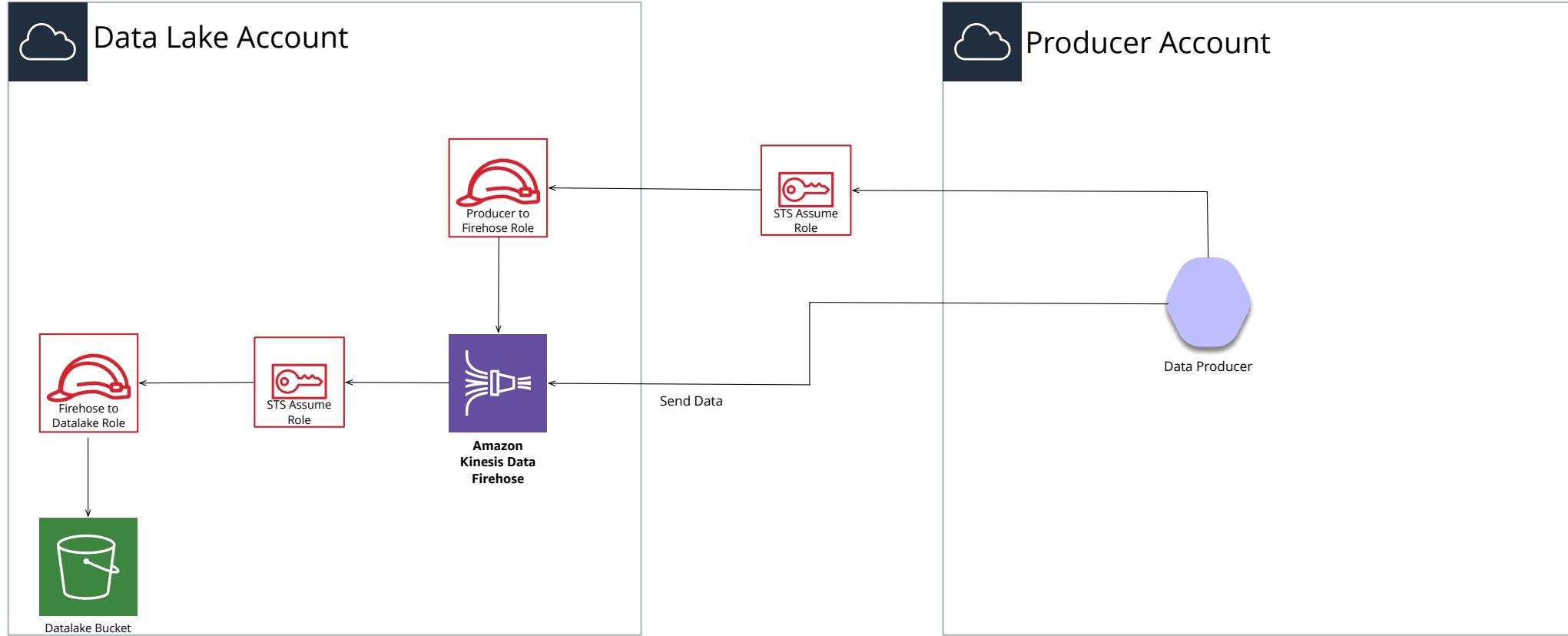
# Firehose Ingestion Architecture



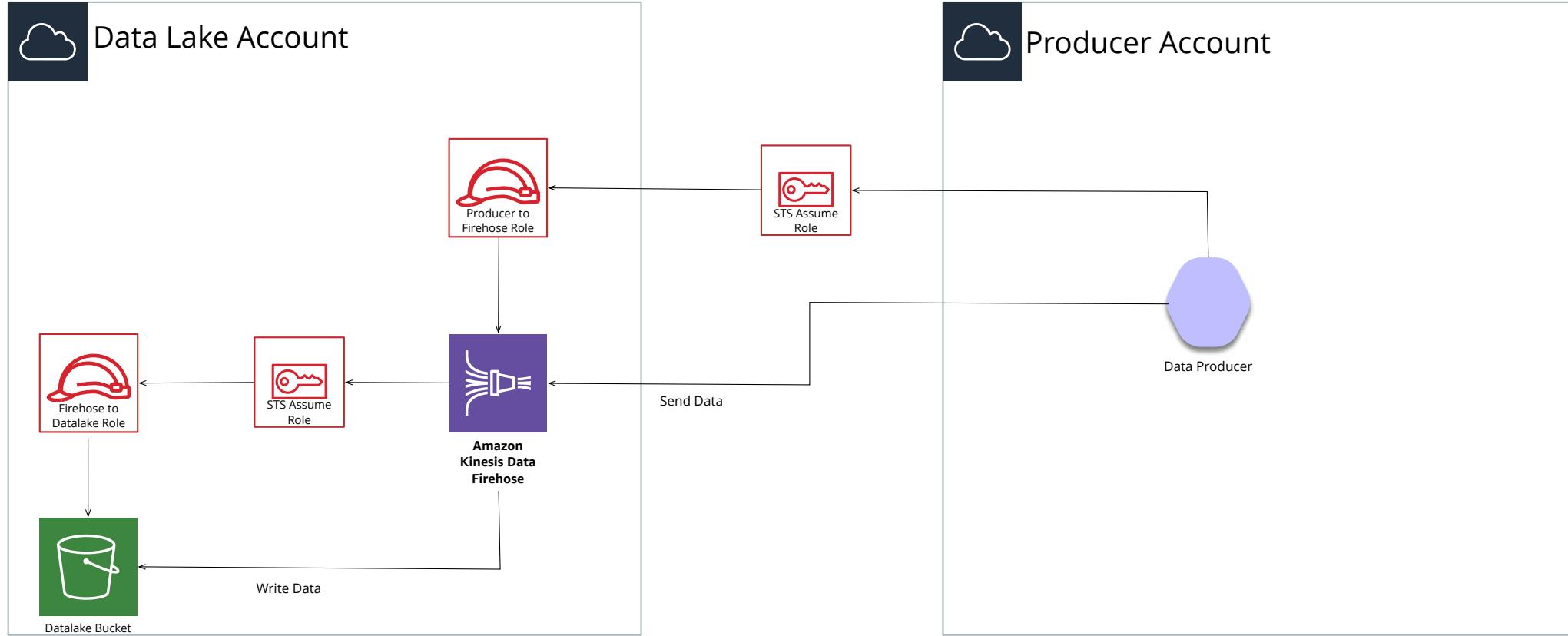
# Firehose Ingestion Architecture



# Firehose Ingestion Architecture



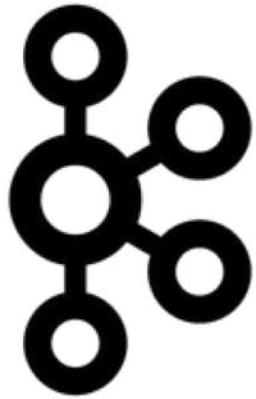
# Firehose Ingestion Architecture



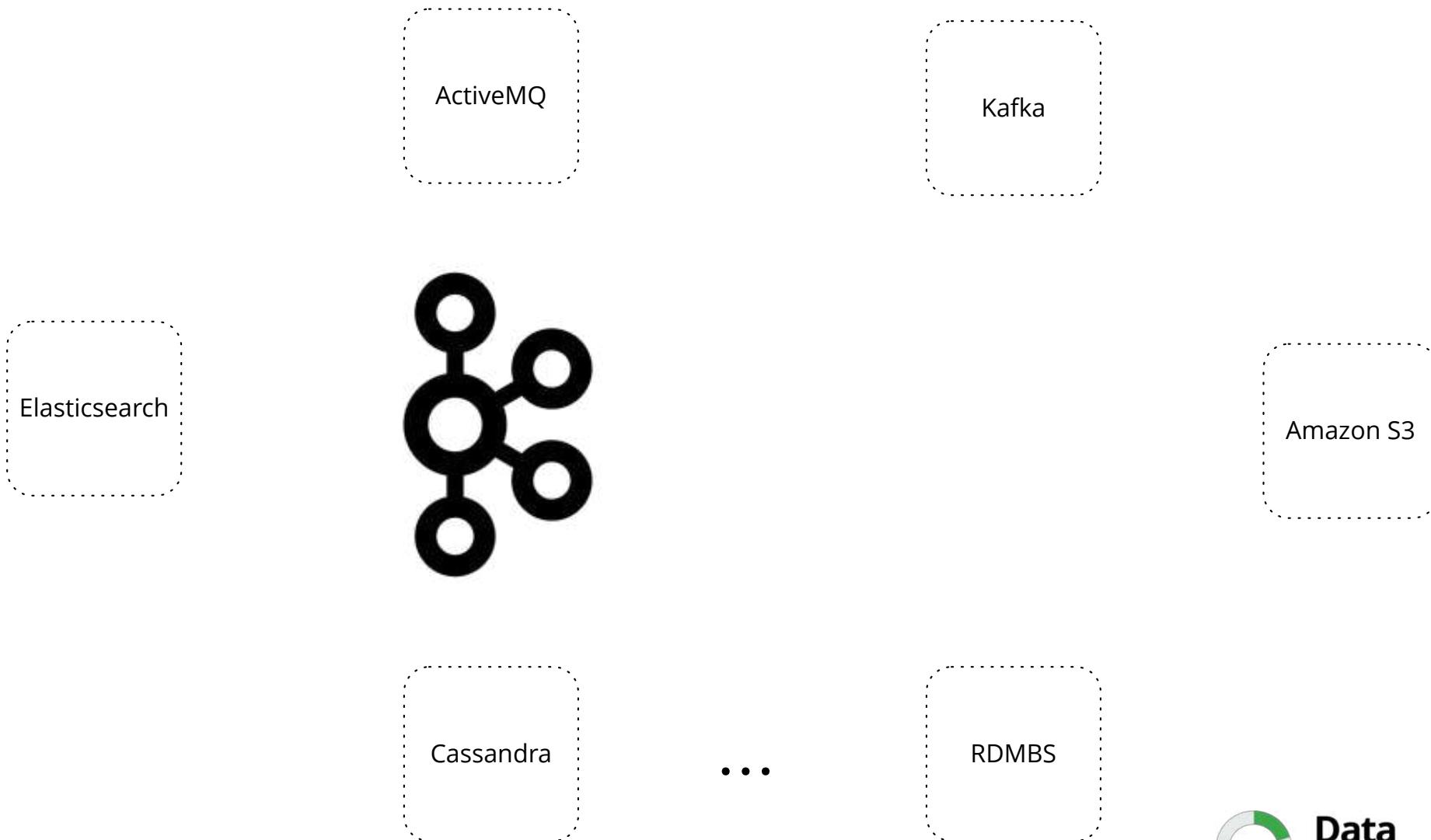
# Kafka Connect Ingestion Mechanism

---

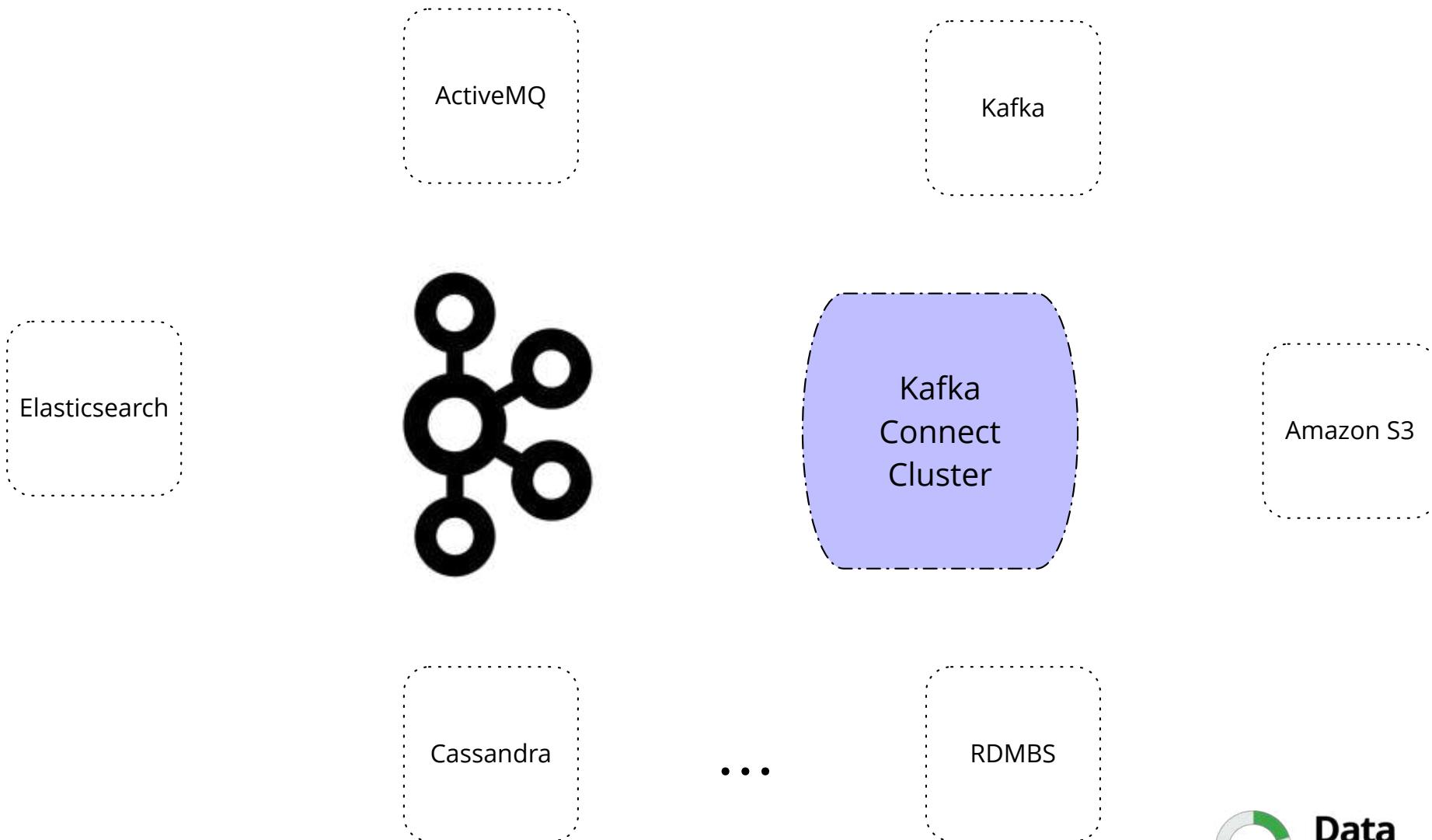
# Kafka Connect



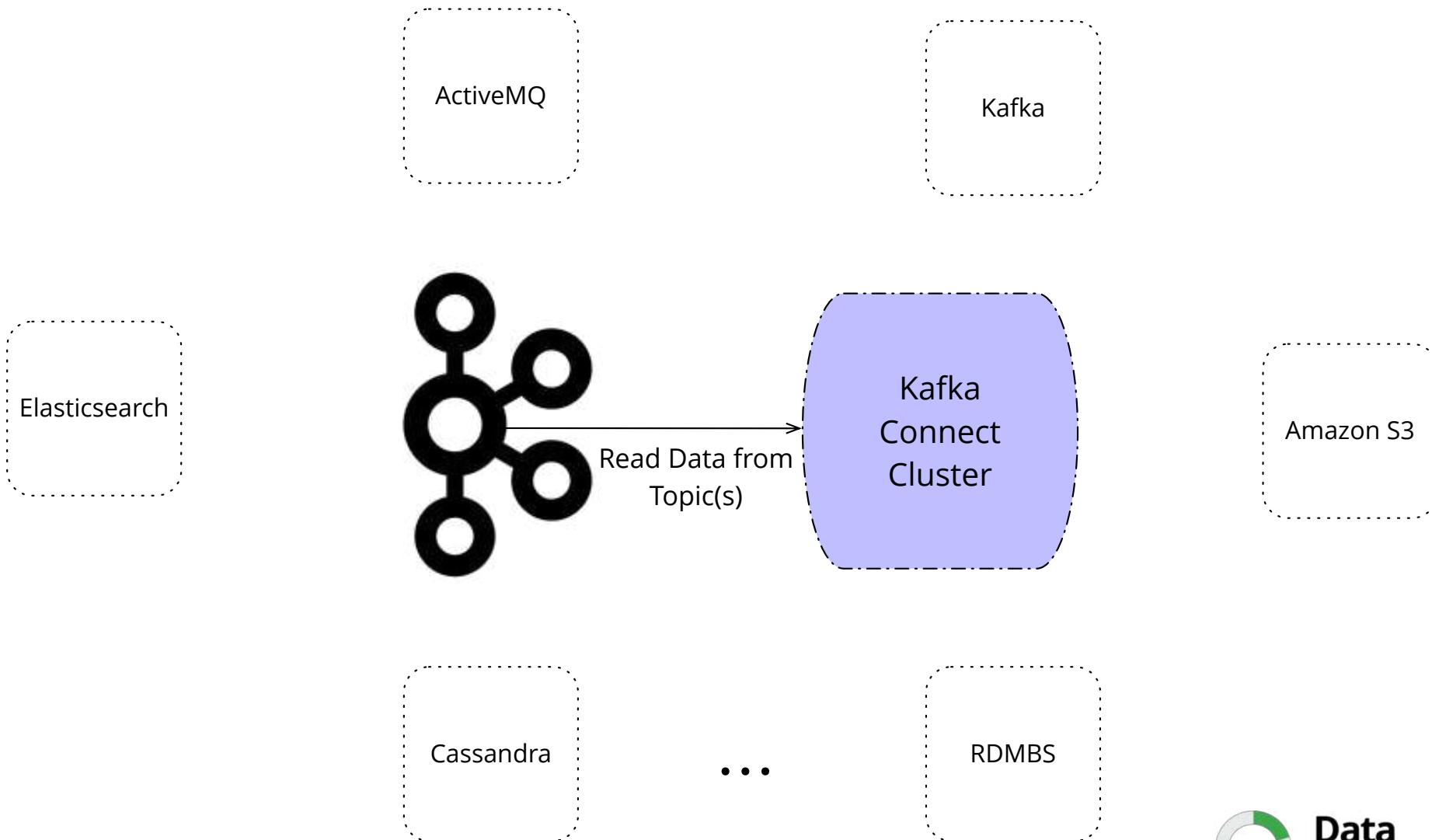
# Kafka Connect



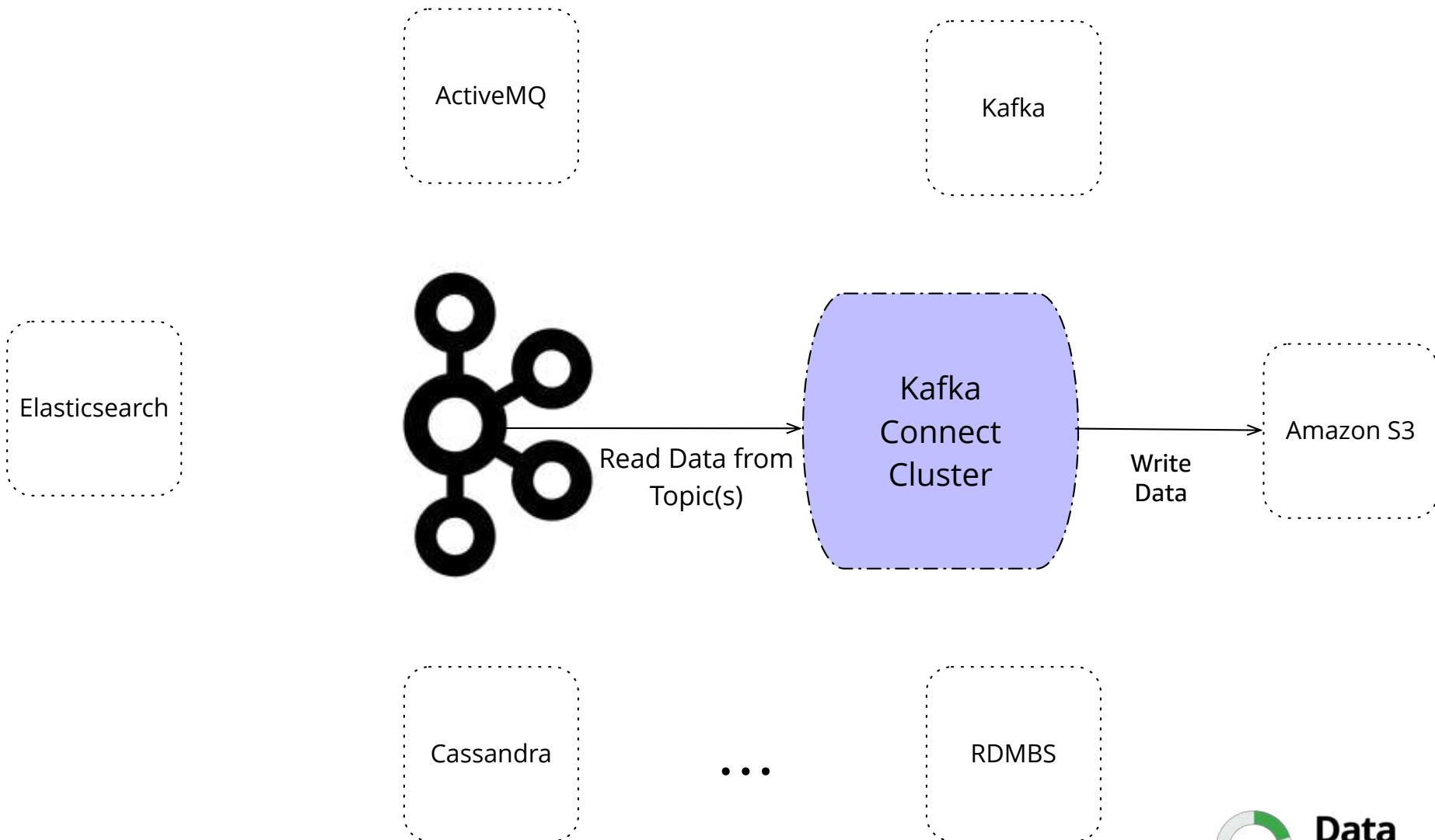
# Kafka Connect



# Kafka Connect



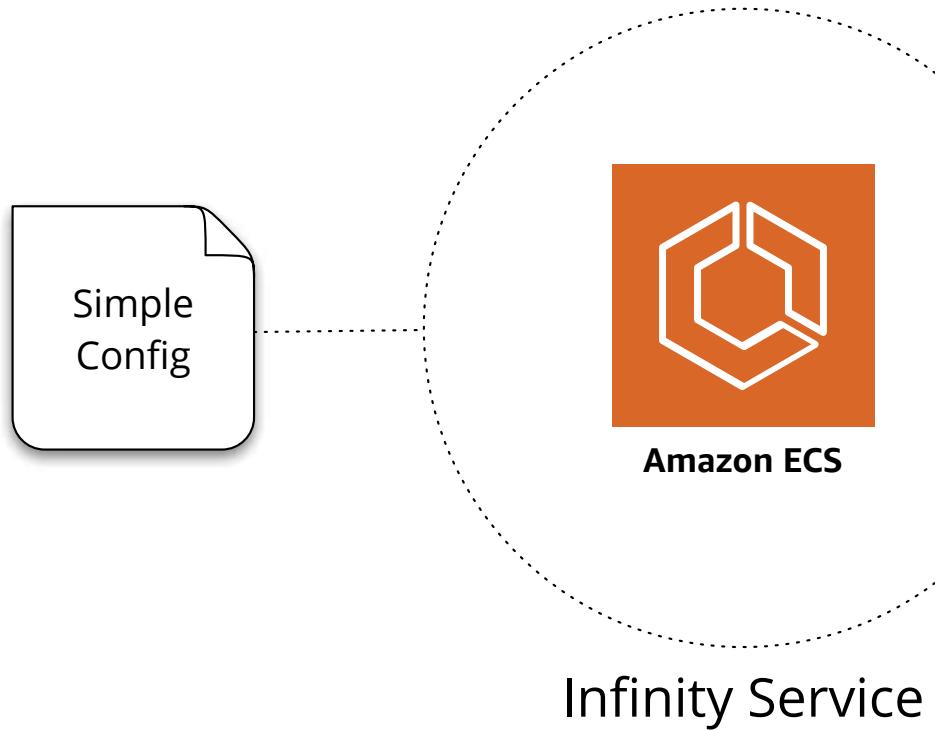
# Kafka Connect



# Scout24 Infinity Cluster



# Scout24 Infinity Cluster



# Scout24 Infinity Cluster



# To Infinity and Beyond Handling Heterogenous Container Clusters in AWS

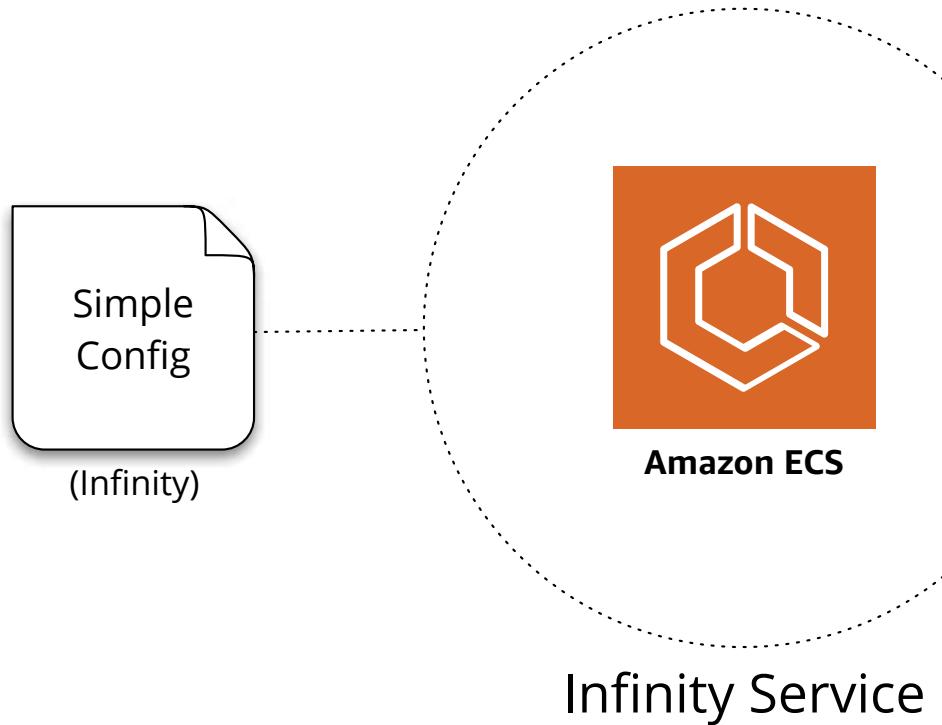
*Christine Trahe, Scout24 AG*

*AWS Summit Berlin 2019*

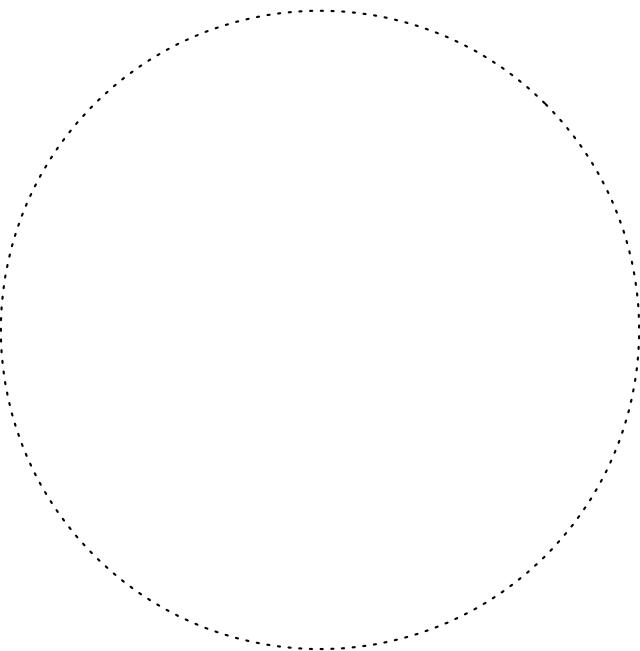


[Get the Slide Deck](#)

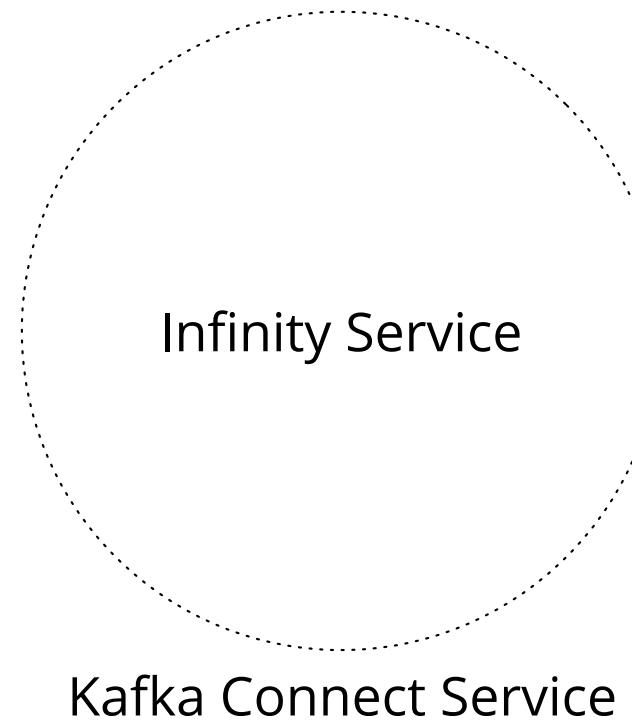
# Kafka Connect on Infinity



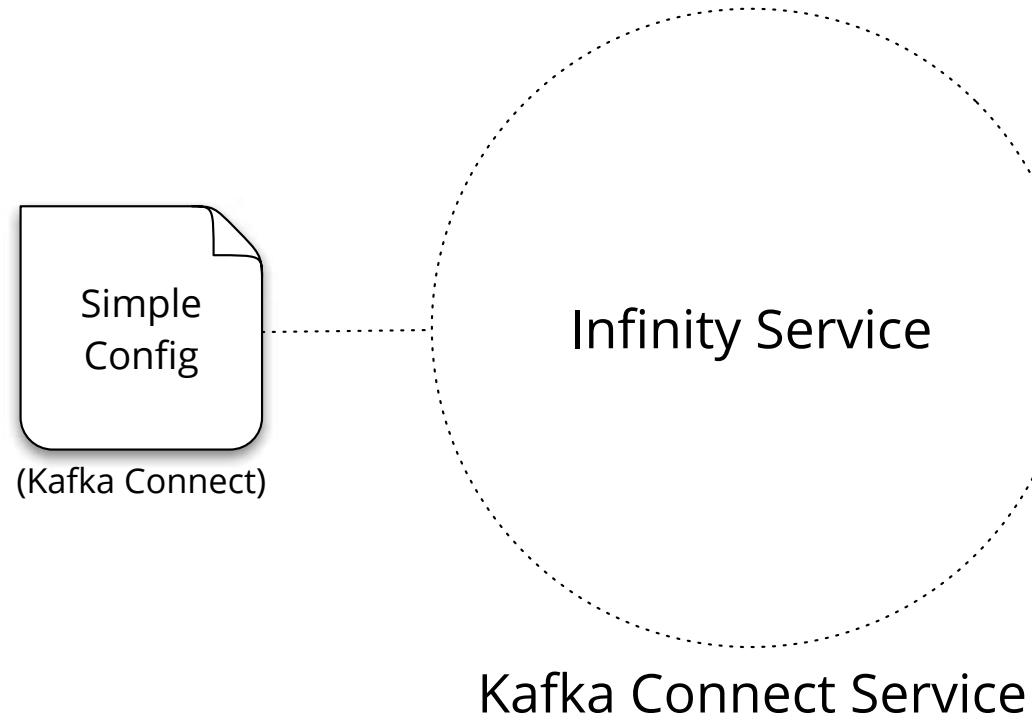
# Kafka Connect on Infinity



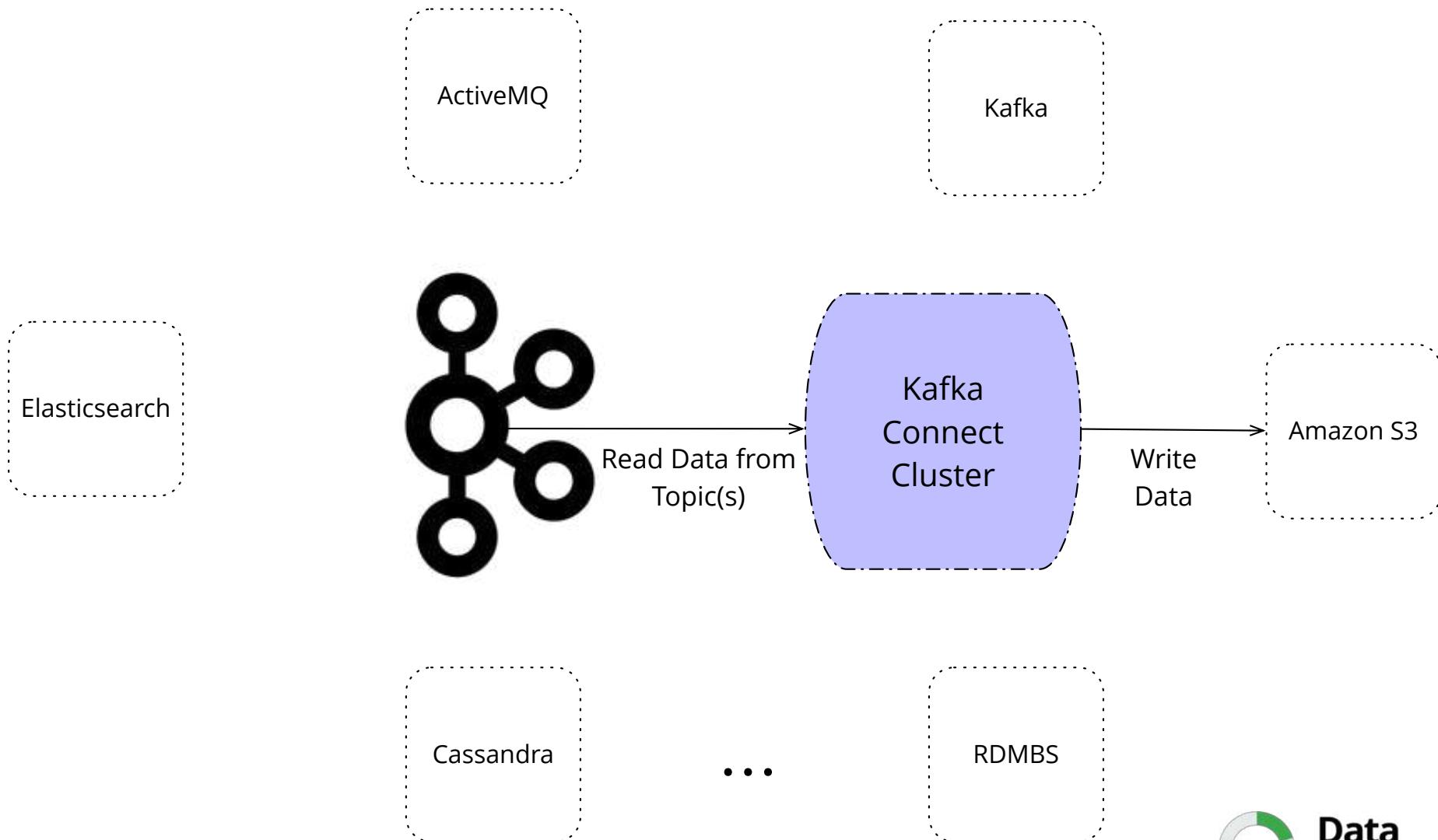
# Kafka Connect on Infinity



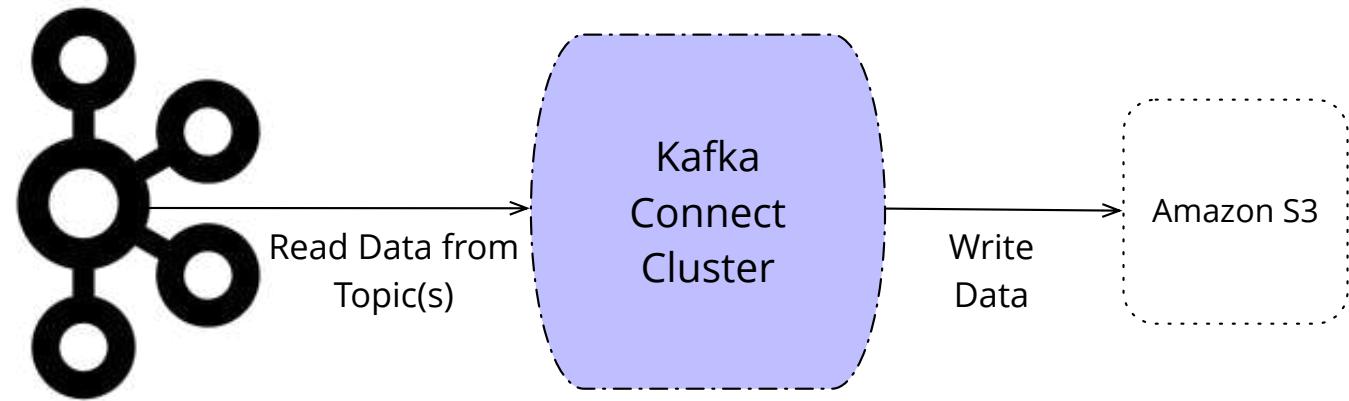
# Kafka Connect on Infinity



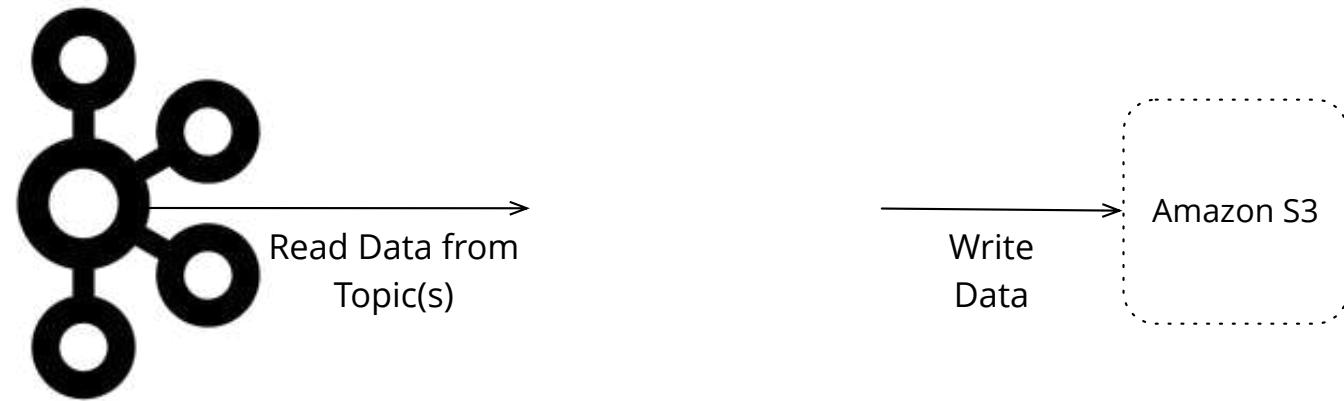
# Kafka Connect Deployment



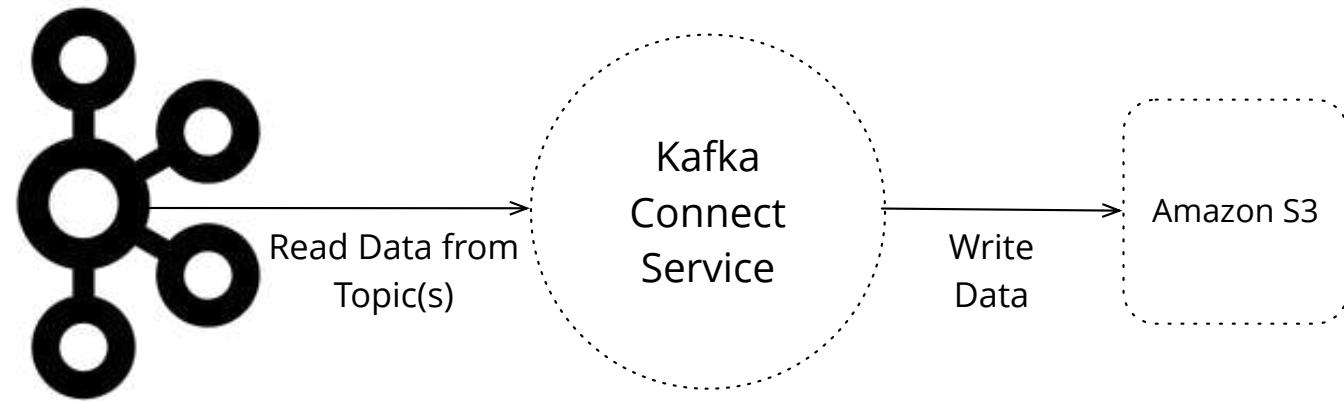
# Kafka Connect Deployment



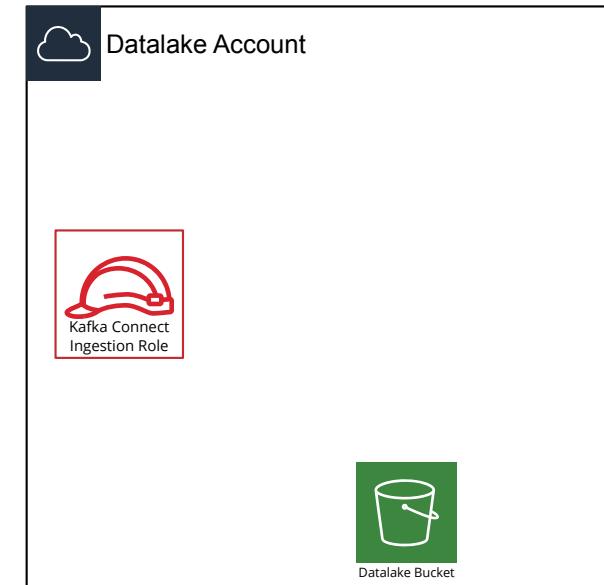
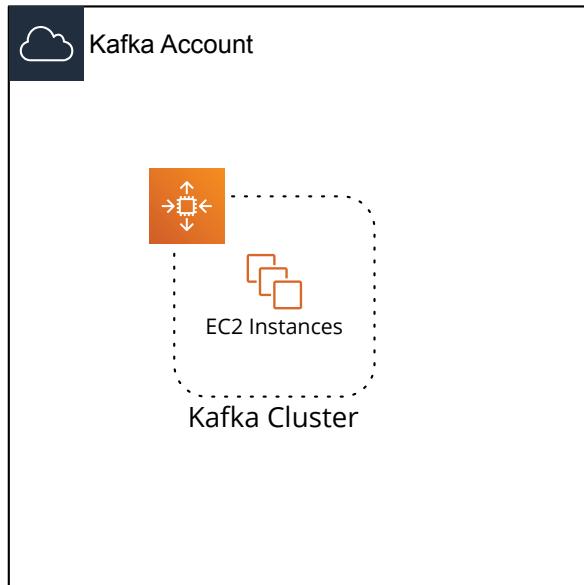
# Kafka Connect Deployment



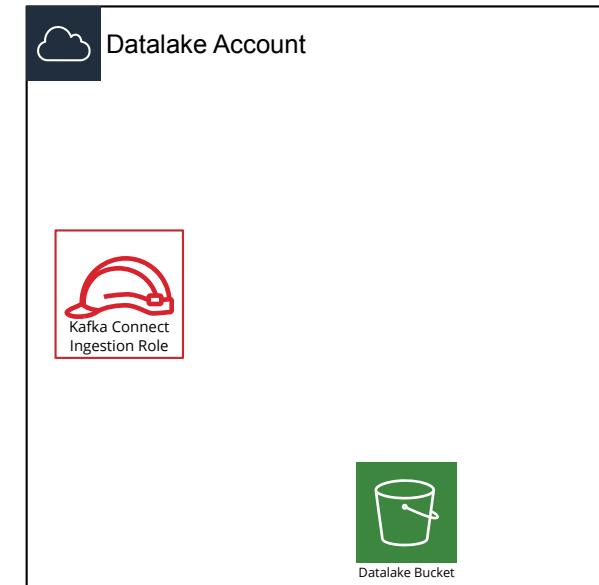
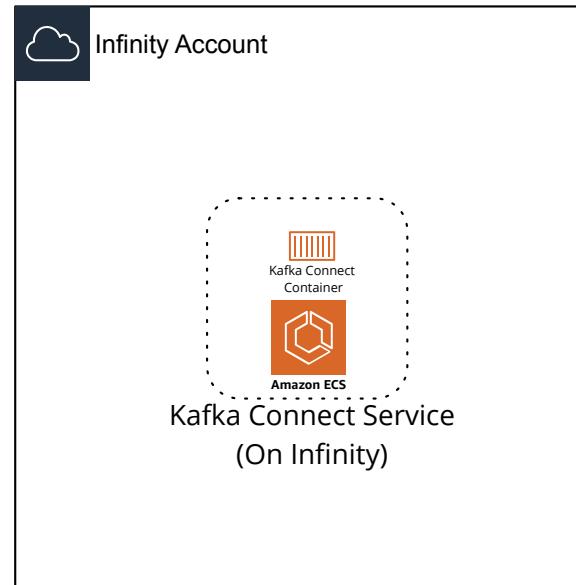
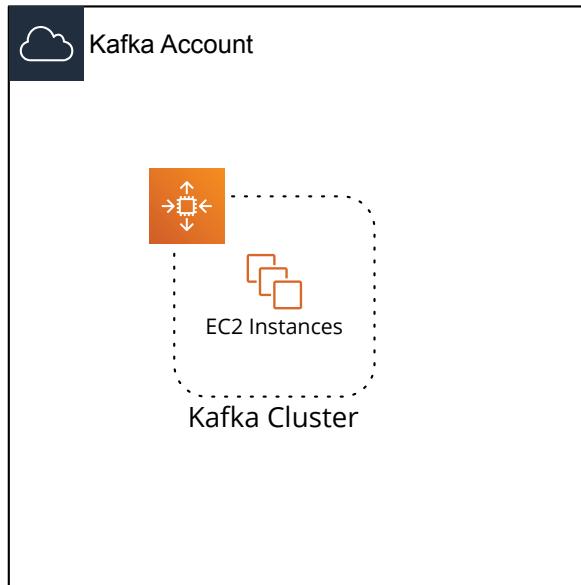
# Kafka Connect Deployment



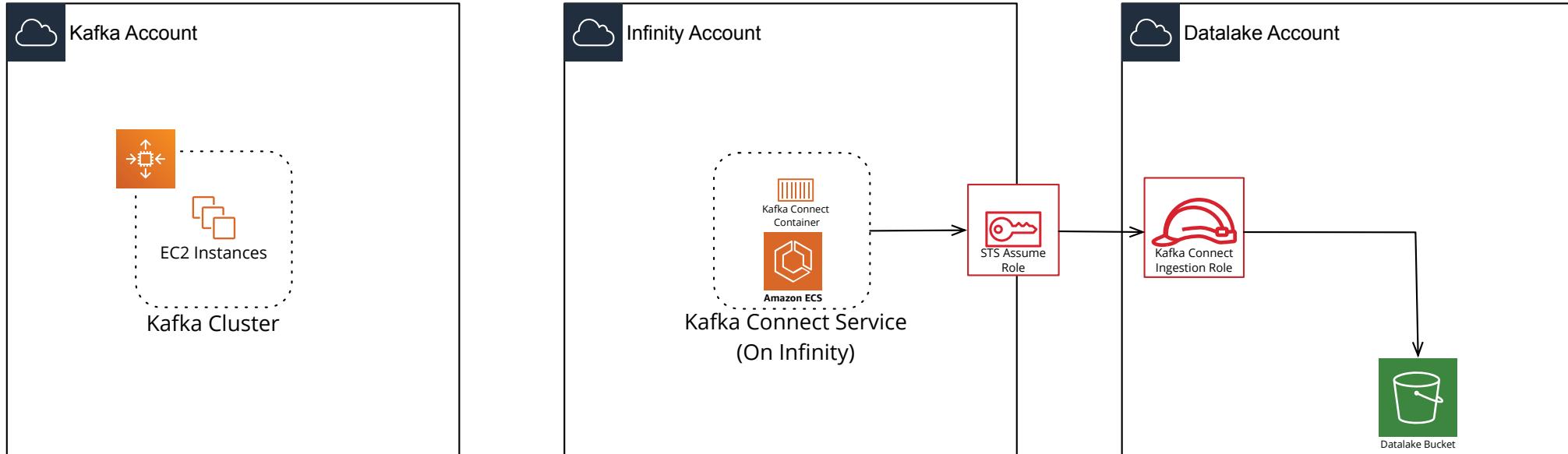
# Kafka Connect Deployment



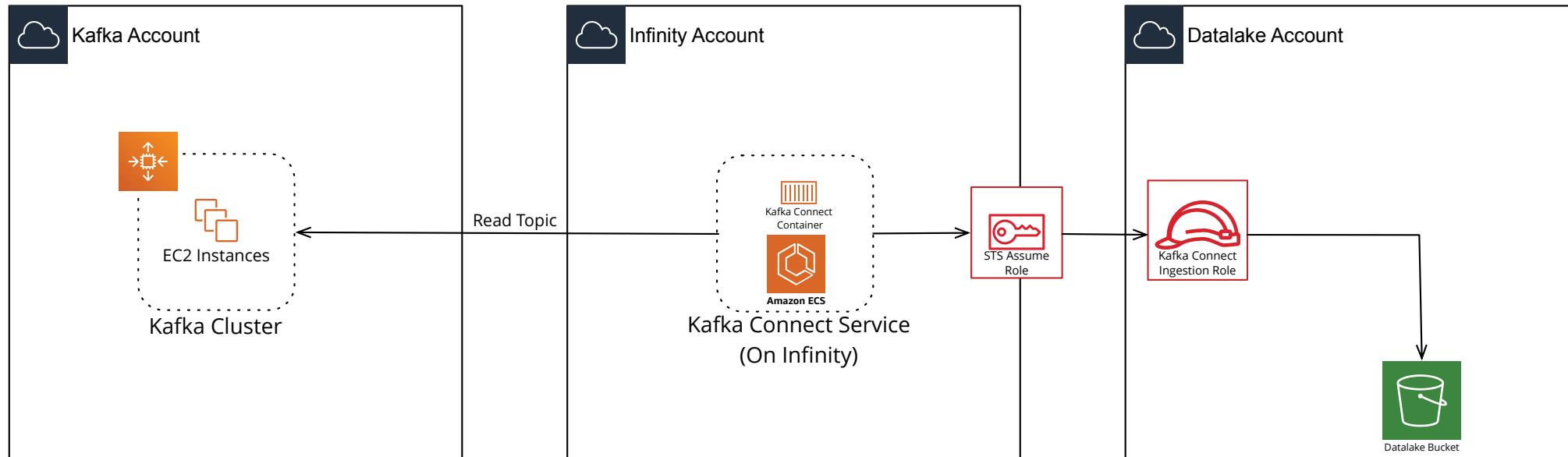
# Kafka Connect Deployment



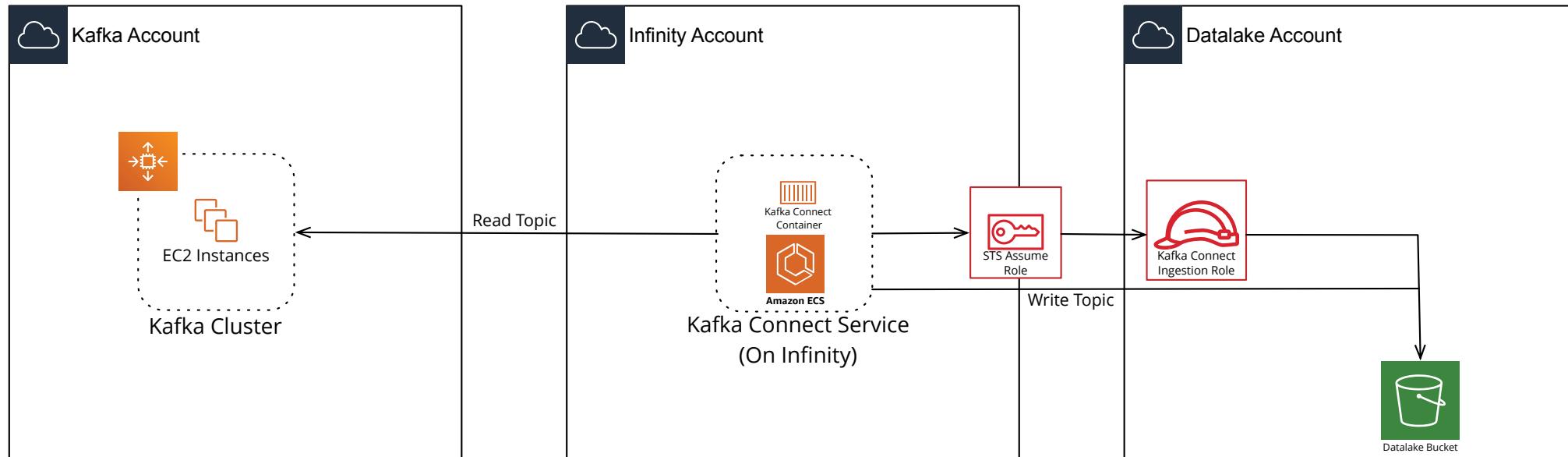
# Kafka Connect Deployment



# Kafka Connect Deployment



# Kafka Connect Deployment



# Hadoop Rest API Ingestion Mechanism

---

# Hadoop Rest API – A Motivation



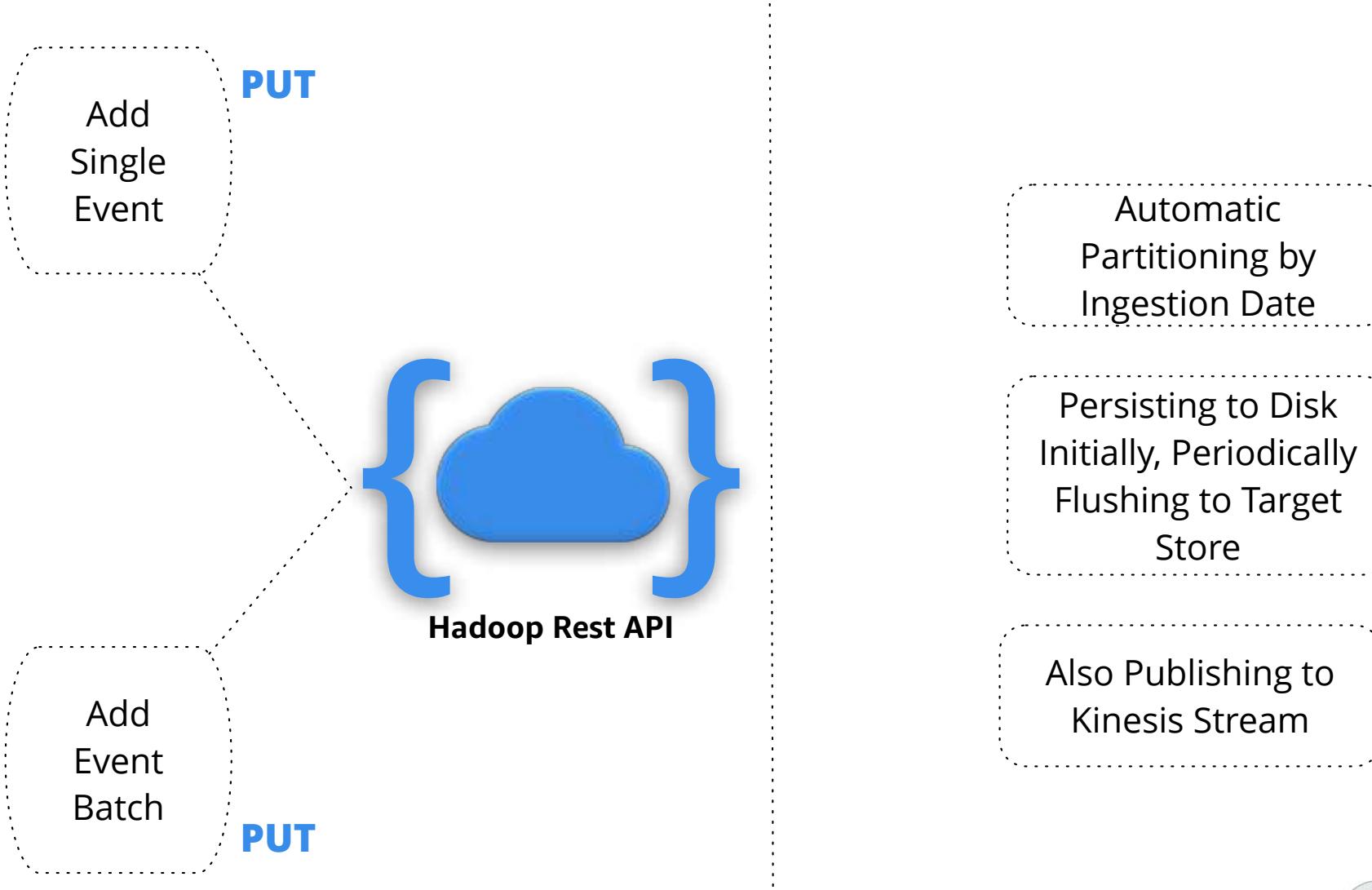
Unified Set of  
Ingestion Operations

Stable User-Facing  
Endpoints, Varying  
Backend

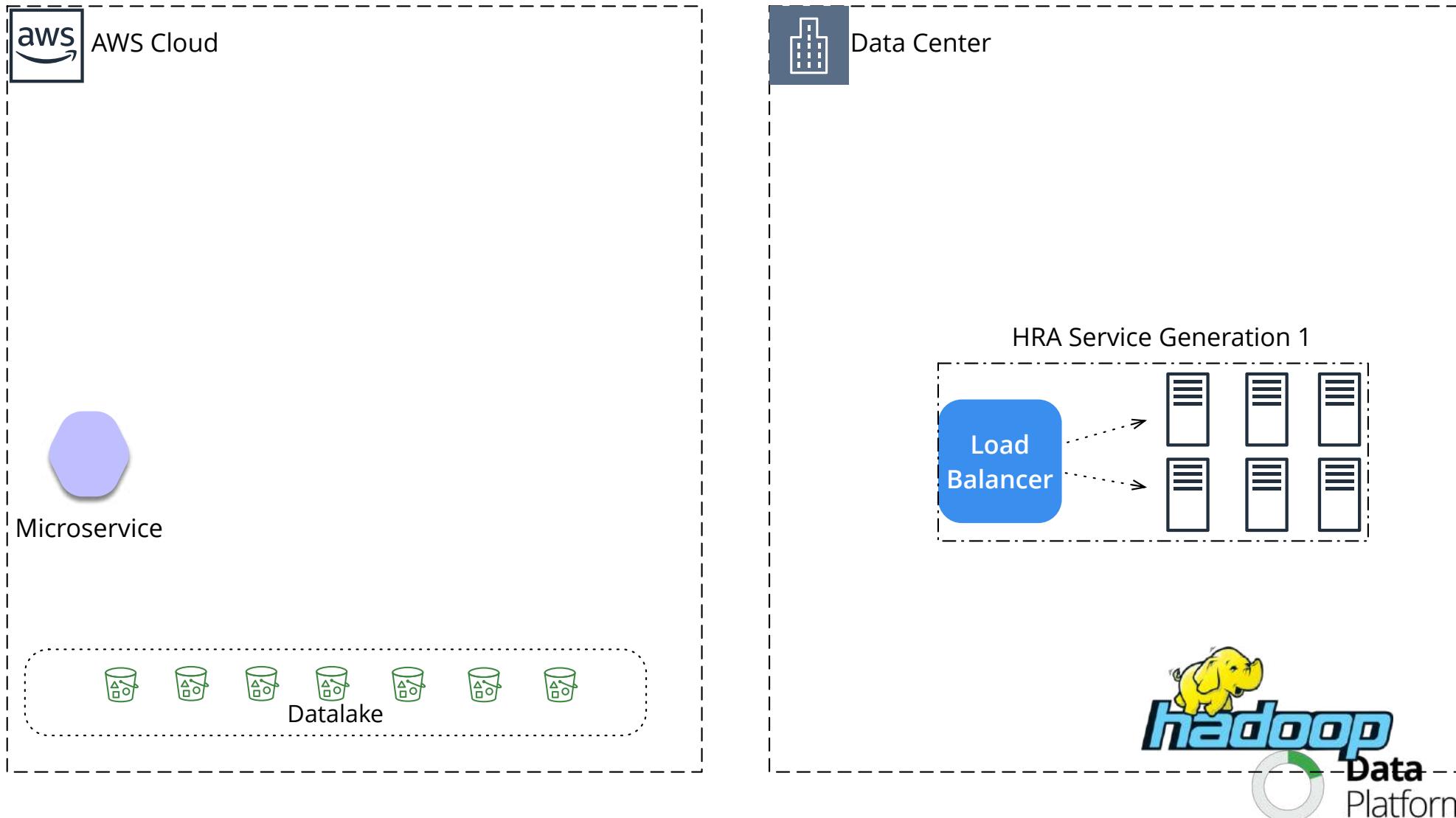
Adds Reception  
Timestamp Field

Performs Event  
Validation

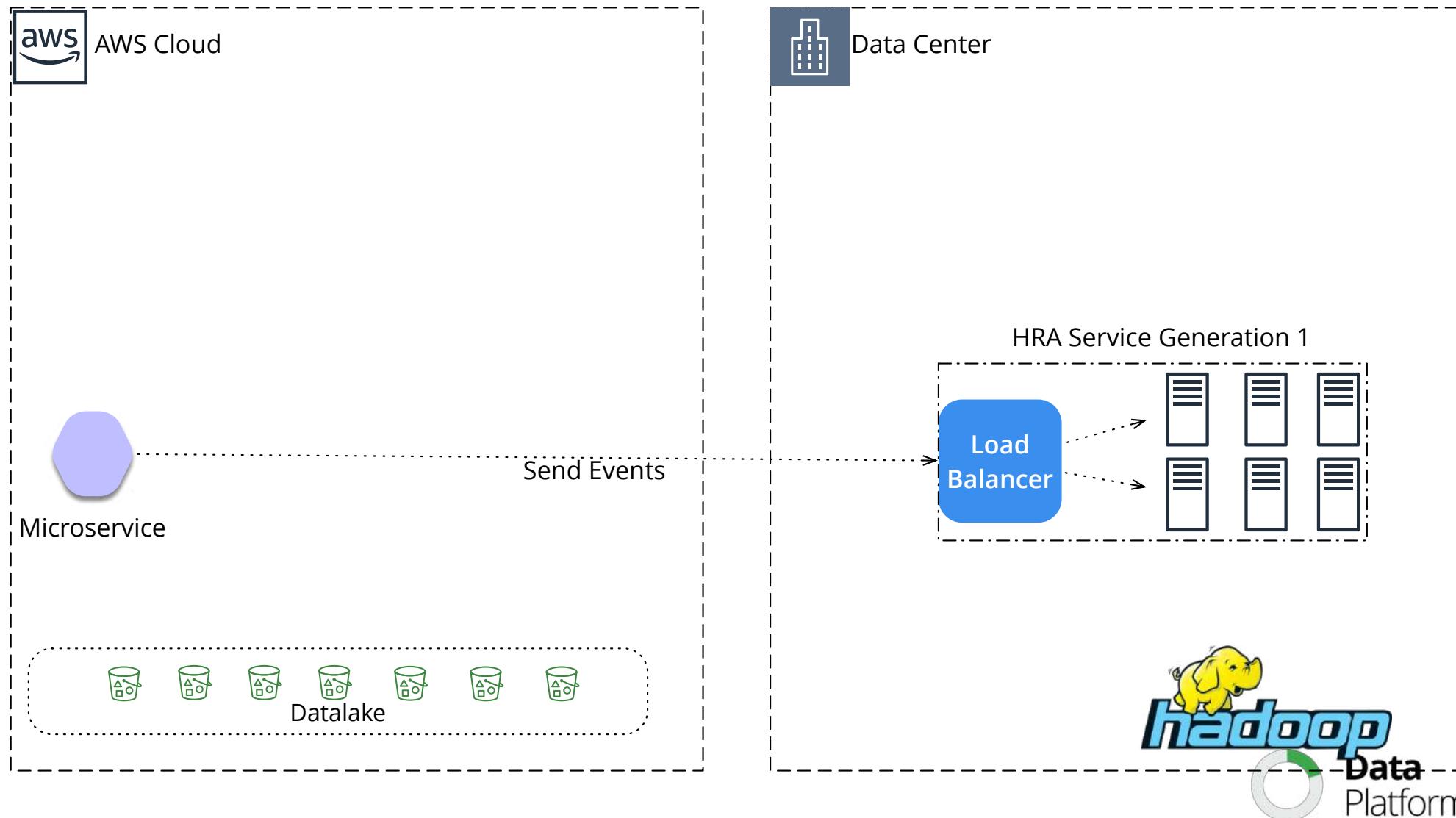
# Hadoop Rest API – Operation Overview



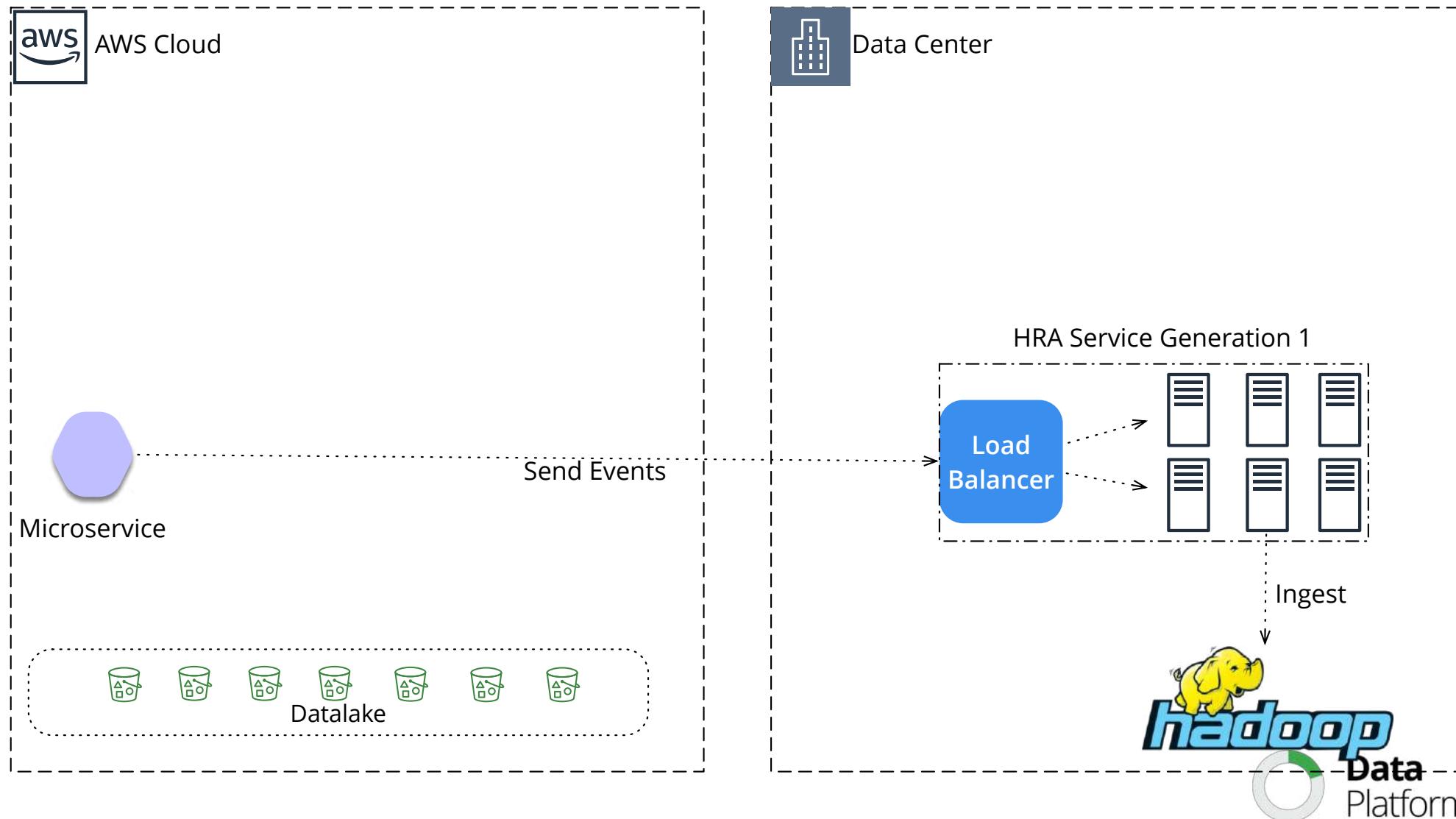
# Hadoop Rest API – Architecture Evolution (Phase 1)



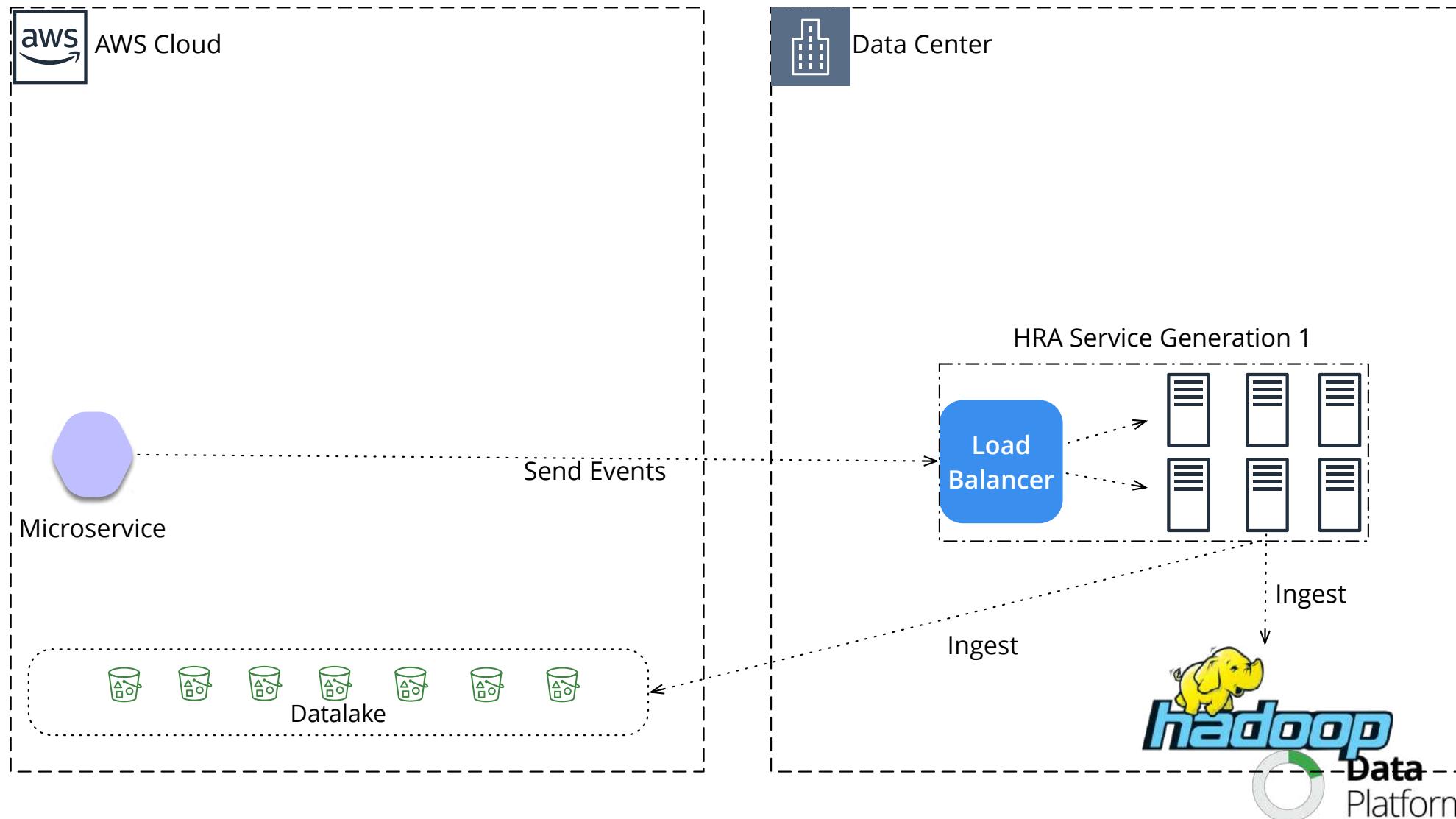
# Hadoop Rest API – Architecture Evolution (Phase 1)



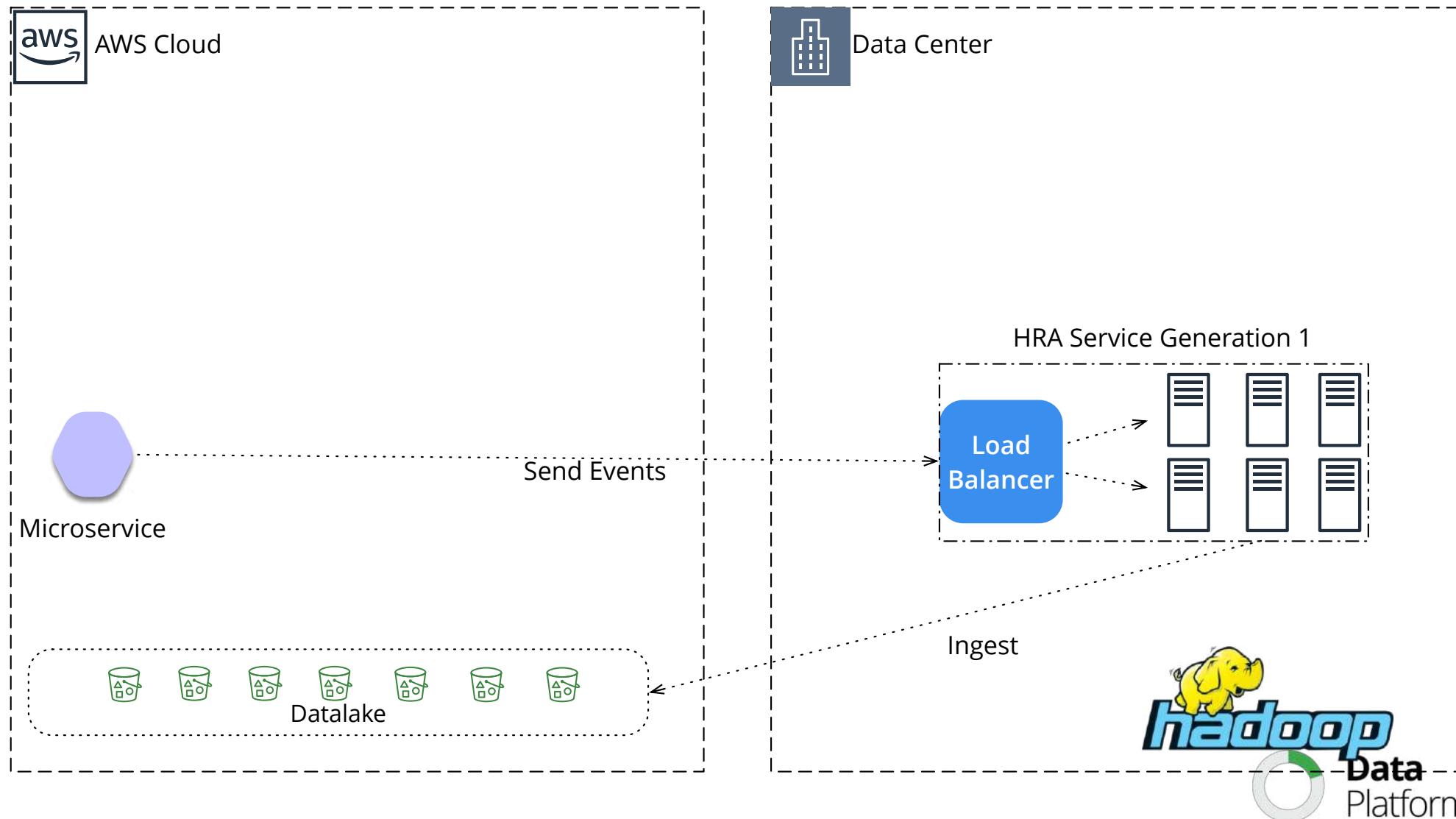
# Hadoop Rest API – Architecture Evolution (Phase 1)



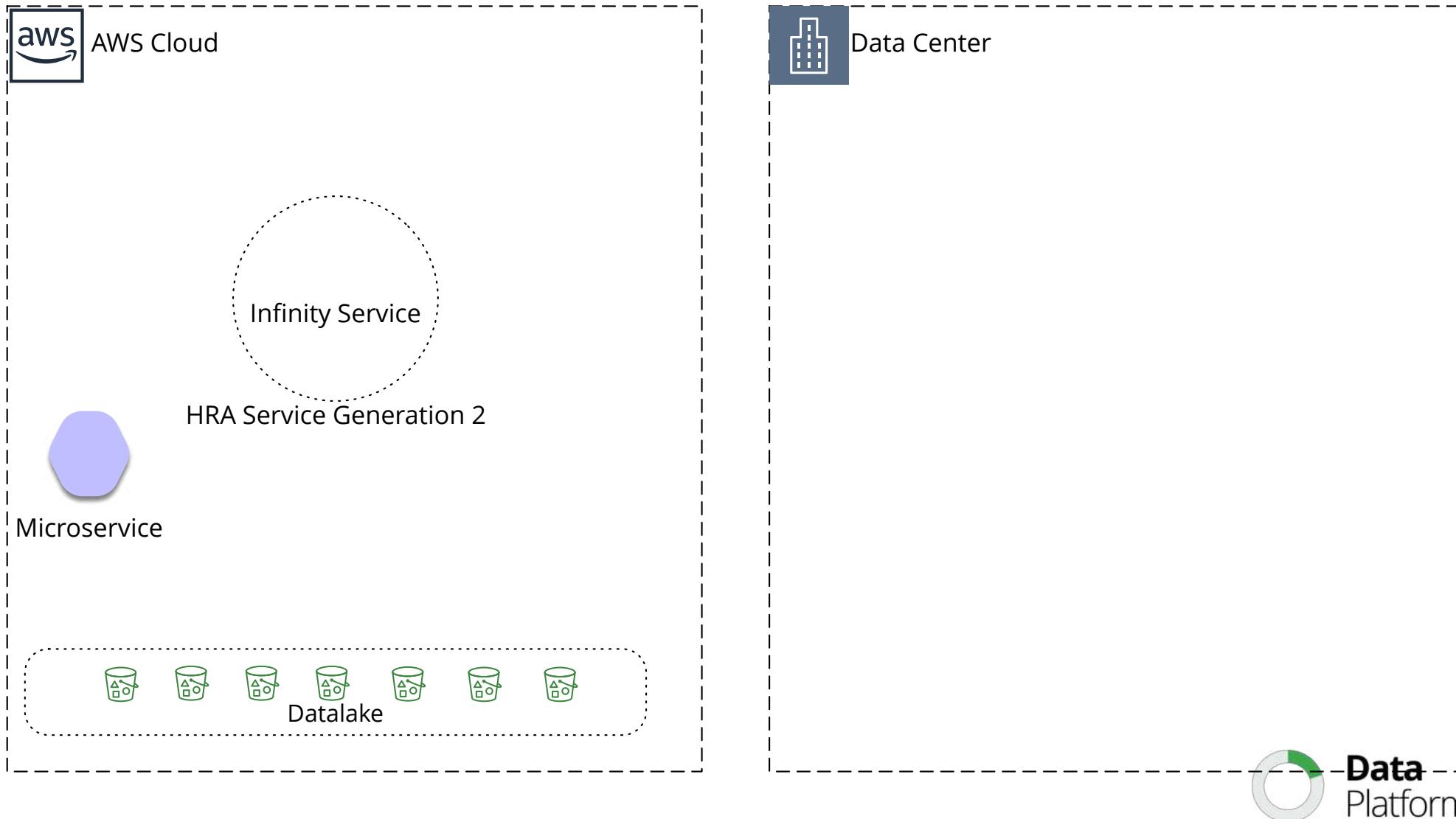
# Hadoop Rest API – Architecture Evolution (Phase 2)



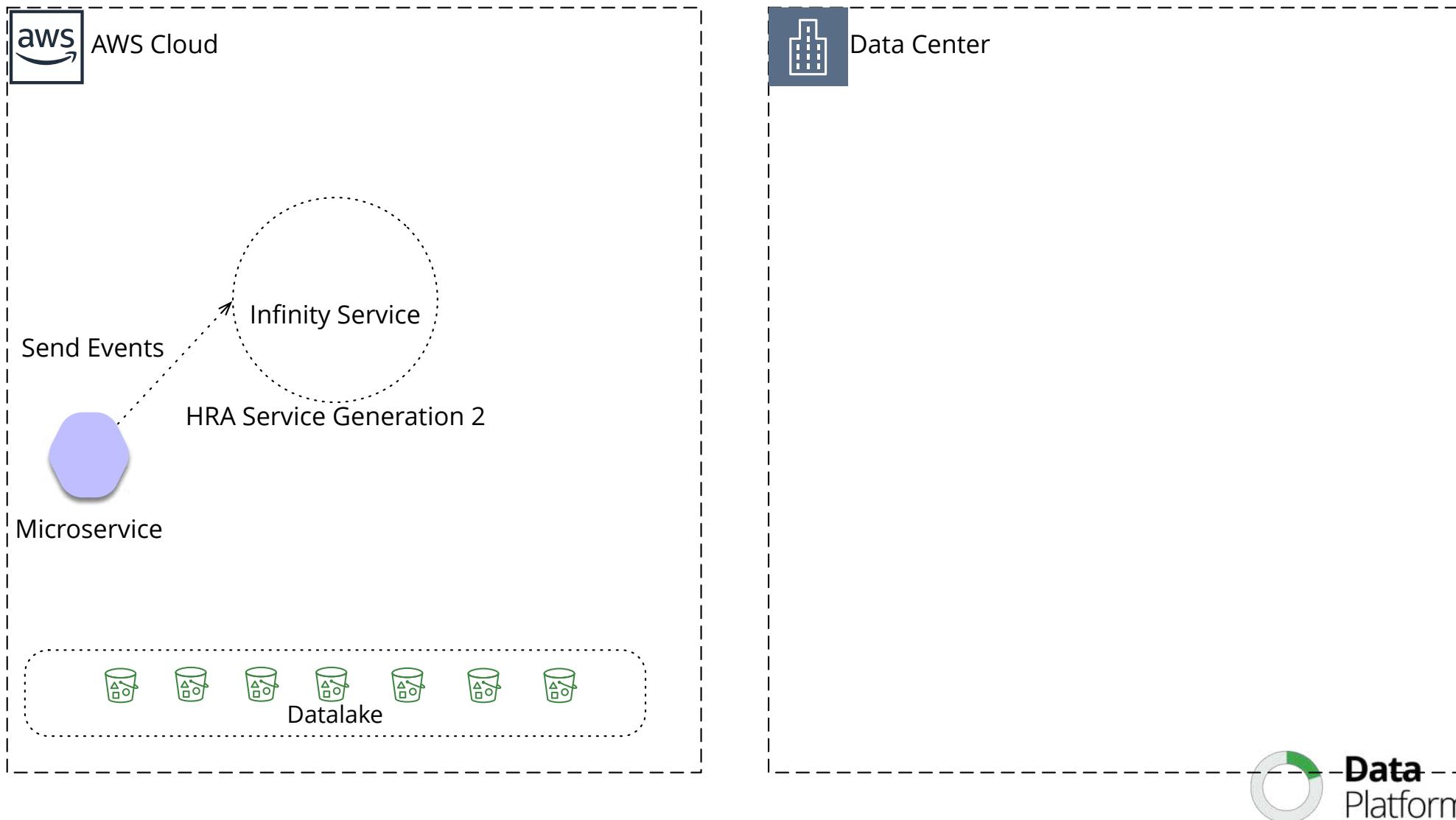
# Hadoop Rest API – Architecture Evolution (Phase 2)



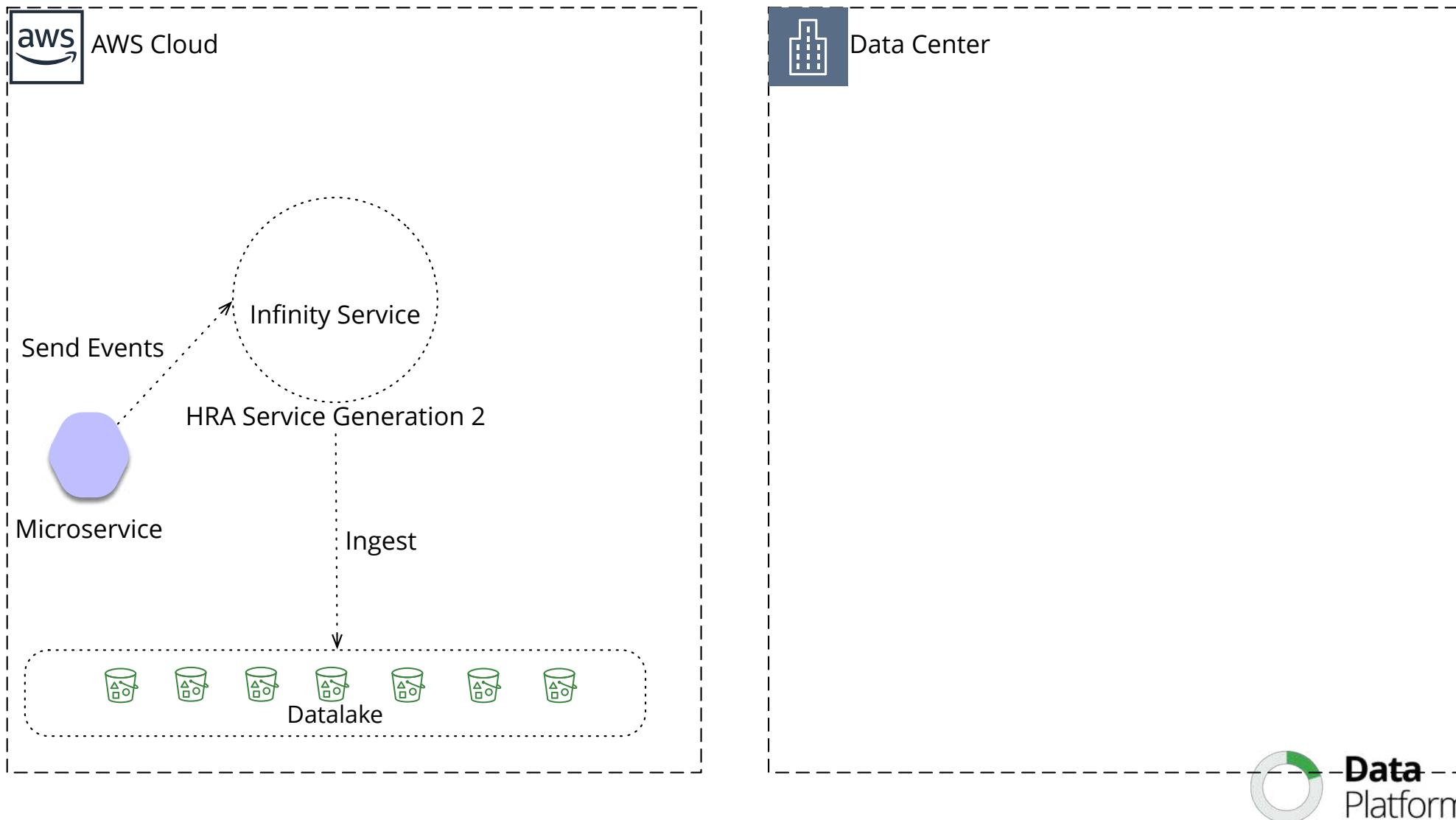
# Hadoop Rest API – Architecture Evolution (Phase 3)



# Hadoop Rest API – Architecture Evolution (Phase 3)



# Hadoop Rest API – Architecture Evolution (Phase 3)



# Self-Service ETL

# Self Service ETL Requirements

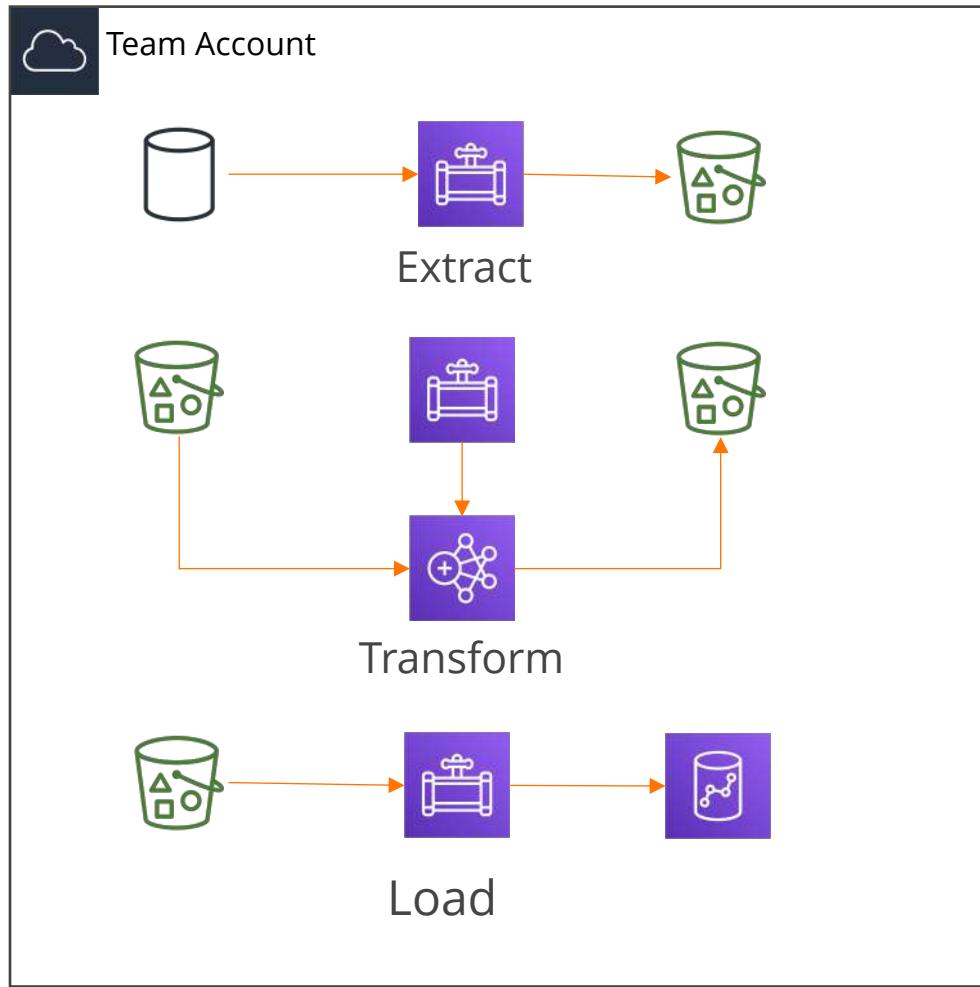
**Data Landscape Manifesto Principle #5:** Data consumers are responsible for the definition and visualization of metrics and for driving the implementation and maintenance of these metrics.

- ➔ Data Consumers get a lot responsibility
- ➔ The Data Platform must empower users to fulfill this responsibility
- ➔ It must be easy to develop, test and maintain data processes

**Multi-Account setup:** Data processes must run in data consumer's account

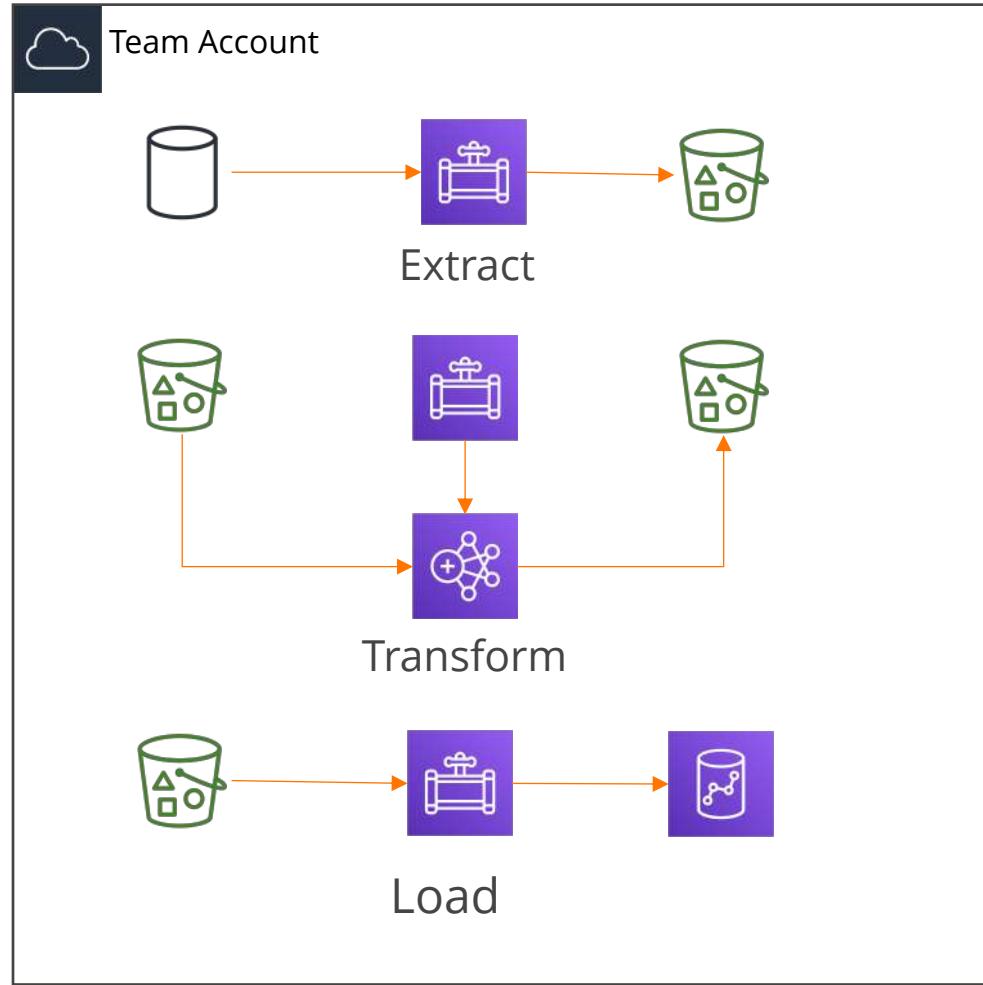
- ➔ No centrally run processing tool
- ➔ AWS Managed service is preferred

# AWS Data Pipeline



ETL runs in team account

# AWS Data Pipeline



Teams run ETL in their accounts



It's easy .... Or is it?



Data  
Platform

SCOUT 24

```

1 {
2   "objects": [
3     {
4       "schedule": {
5         "ref": "DefaultSchedule"
6       },
7       "directoryPath": "s3://sourceBucket",
8       "name": "SourceS3DataNode",
9       "id": "DataNodeId_yIW0K",
10      "type": "S3DataNode"
11    },
12    {
13      "period": "1 days",
14      "name": "Every 1 day",
15      "id": "DefaultSchedule",
16      "type": "Schedule",
17      "startAt": "FIRST_ACTIVATION_DATE_TIME"
18    },
19    {
20      "output": {
21        "ref": "RedshiftDataNodeId_sAlcV"
22      },
23      "input": {
24        "ref": "S3DataNodeId_VJmwh"
25      },
26      "schedule": {
27        "ref": "DefaultSchedule"
28      },
29      "name": "Load",
30      "id": "CopyActivityId_jPabN",
31      "runsOn": {
32        "ref": "ResourceId_YzpFO"
33      },
34      "type": "RedshiftCopyActivity",
35      "insertMode": "OVERWRITE_EXISTING"
36    },
37    {
38      "schedule": {
39        "ref": "DefaultSchedule"
40      },
41      "resourceRole": "DataPipelineDefaultResourceRole",
42      "role": "DataPipelineDefaultRole",
43      "name": "Ec2Instance",
44      "id": "ResourceId_YzpFO",
45      "type": "Ec2Resource"
46    },
47    {
48      "schedule": {
49        "ref": "DefaultSchedule"
50      },
51      "database": {
52        "ref": "DatabaseId_nwLUX"
53      },
54      "name": "SourceSqlDataNode",
55      "id": "SqlDataNodeId_nhC9I",
56      "type": "SqlDataNode",
57      "table": "sourceTable"
58    },
59    {
60      "output": {
61        "ref": "DataNodeId_yIW0K"
62      },
63      "input": {
64        "ref": "SqlDataNodeId_nhC9I"
65      },
66      "schedule": {
67        "ref": "DefaultSchedule"
68      },
69      "name": "Extract",
70      "workerGroup": "local-worker-1",
71      "id": "CopyActivityId_0yKN4",
72      "type": "CopyActivity"
73    },
74    {
75      "schedule": {
76        "ref": "DefaultSchedule"
77      },
78      "name": "EmrResource",
79      "releaseLabel": "emr-5.21.0",
80      "id": "ResourceId_tBzz1",
81      "type": "EmrCluster"
82    },
83    {
84      "schedule": {
85        "ref": "DefaultSchedule"
86      },
87      "name": "Transform",
88      "runsOn": {
89        "ref": "ResourceId_tBzz1"
90      },
91      "id": "ShellCommandActivityId_NM2Ga",
92      "type": "ShellCommandActivity",
93      "command": "echo \\\"Transform\\\""
94    },
95    {
96      "failureAndRerunMode": "CASCADE",
97      "schedule": {
98        "ref": "DefaultSchedule"
99      },
100     "resourceRole": "DataPipelineDefaultResourceRole",
101     "role": "DataPipelineDefaultRole",
102     "scheduleType": "cron",
103     "name": "Default",
104     "id": "Default"
105   },
106   {
107     "schedule": {
108       "ref": "DefaultSchedule"
109     },
110     "database": {
111       "ref": "RedshiftDatabaseId_1IMYz"
112     },
113     "name": "TargetRedshiftDataNode",
114     "id": "RedshiftDataNodeId_sAlcV",
115     "type": "RedshiftDataNode",
116     "tableName": "targettable"
117   },
118   {
119     "schedule": {
120       "ref": "DefaultSchedule"
121     },
122     "directoryPath": "s3://target-bucket",
123     "name": "TargetS3DataNode",
124     "id": "S3DataNodeId_VJmwh",
125     "type": "S3DataNode"
126   },
127   {
128     "connectionString": "jdbc:oracle:thin:/host/db",
129     "*password": "geheim",
130     "name": "SourceDB",
131     "id": "DatabaseId_nwLUX",
132     "type": "JdbcDatabase",
133     "jdbcDriverClass": "com.oracle.Driver",
134     "username": "user"
135   },
136   {
137     "connectionString": "jdbc:redshift://cluster//",
138     "*password": "passwo",
139     "name": "TargetDB",
140     "id": "RedshiftData",
141     "type": "RedshiftDat",
142     "username": "user"
143   },
144   ],
145   "parameters": []
146 }

```



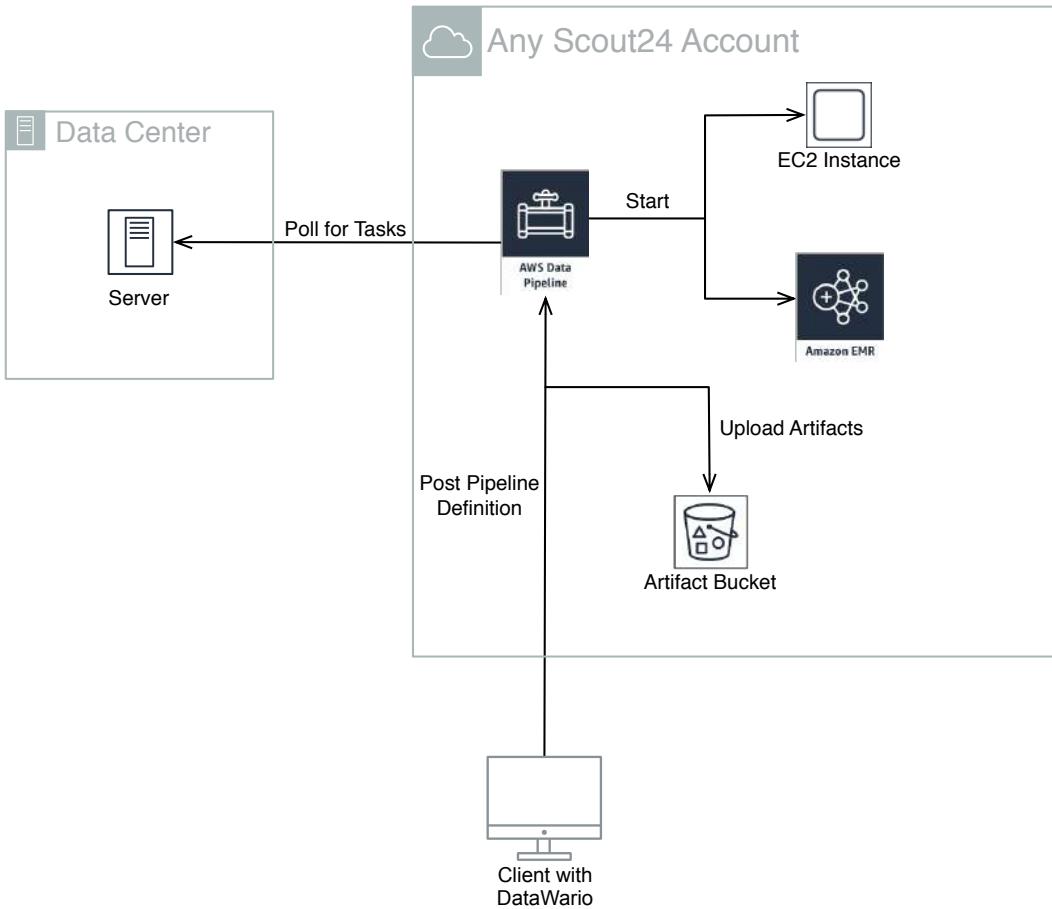
# DataWario to the rescue

```
1 emr:
2 ec2:
3 databases:
4   sourceDb:
5     driverClass: oracle.jdbc.OracleDriver
6     driverUri:ojdbc6.jar
7     user: aws_export
8     password: wonttell
9     url: jdbc:oracle:thin:@//myDB
10    targetDb:
11      type: redshift
12      clusterId: redshift-cluster-id
13      user: redshift-user
14      password: redshift-password
15      dbName: cloud_dwh
16    steps:
17      - database_to_s3:
18        dbRef: sourceDb
19        format: csv
20        selectQuery: select 1 from dual
21        s3Url: s3://someBucket/some.csv
22      - shell:
23        command: "aws s3 cp --recursive s3://someBucket/some.csv s3://someOtherBucket/some.csv"
24        runsOn: emr
25      - s3_to_database:
26        runsOn: ec2
27        dbRef: targetDb
28        insertMode: OVERWRITE_EXISTING
29        commandOptions:
30          - delimiter '\t'
31          - gzip
32          - truncatecolumns
33          - maxerror 1000
34          - blanksasnull
35        table: thetable
36        schema: myschema
37        s3Url: s3://someOtherBucket/some.csv
```

# DataWario

1. Client for AWS Data Pipeline
2. Transforms YAML into JSON format understood by Data Pipeline
3. Integrates with the Scout24 account setup and sets sensible configuration defaults
4. Augments Data Pipeline with additional features
5. Makes testing and debugging faster

# DataWario – Architecture

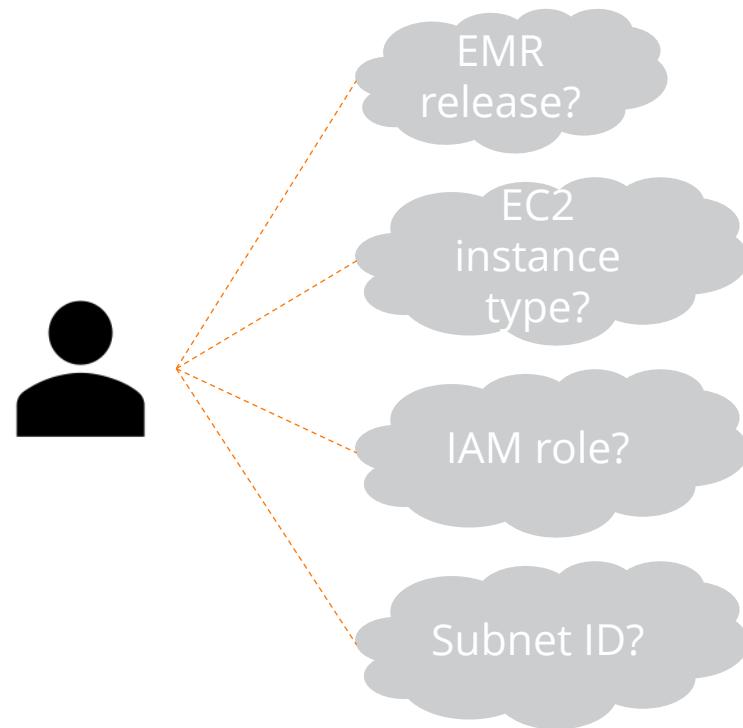


- Java Application, distributed as a JAR
- Shell wrapper (`dw`) for easy usage from the command line

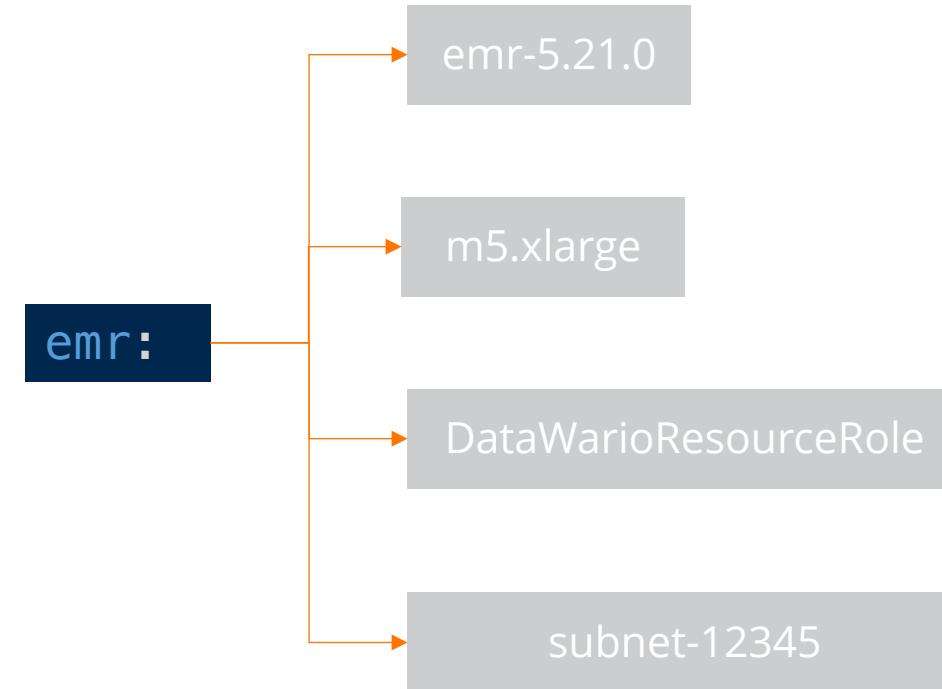
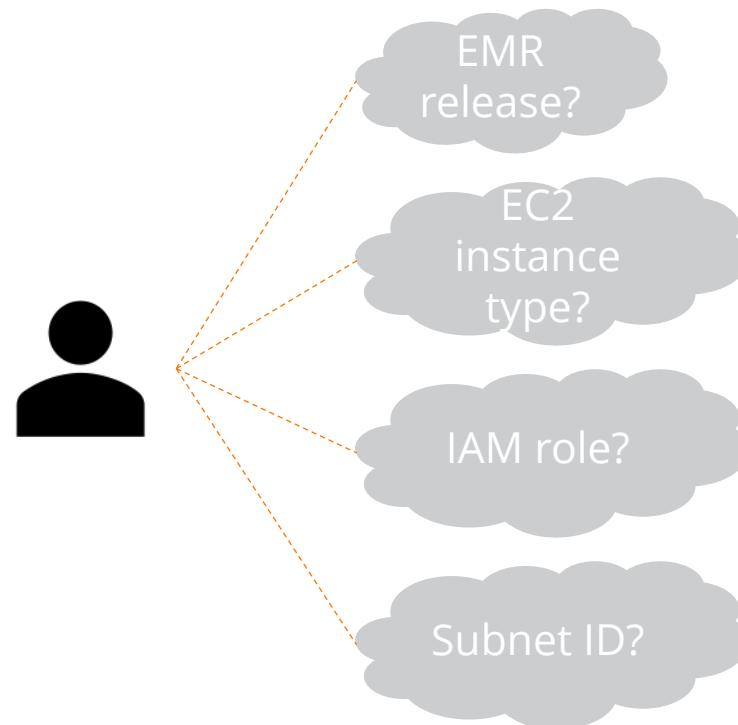
## Deployment of a pipeline

1. User issues `dw deploy`
2. DataWario
  1. uploads artifacts from user's machine to S3
  2. generates JSON definition from YAML
  3. creates data pipeline
3. AWS Data Pipeline runs the defined steps

# Sensible defaults



# Sensible defaults



# Augmenting data pipeline



ID	Make	Model	Price	Mileage
1	VW	Golf	22000	15000
2	BMW	320d	35000	40000

„Built-In“ Copy Activity

1, VW, Golf, 22000, 15000

2, BMW, 320d, 35000, 40000

- CSV (or TSV)
- No schema

# Augmenting data pipeline



ID	Make	Model	Price	Mileage
1	VW	Golf	22000	15000
2	BMW	320d	35000	40000

## „Built-In“ Copy Activity

1, VW, Golf, 22000, 15000

2, BMW, 320d, 35000, 40000

- CSV (or TSV)
- No schema

## „Custom“ Copy Activity

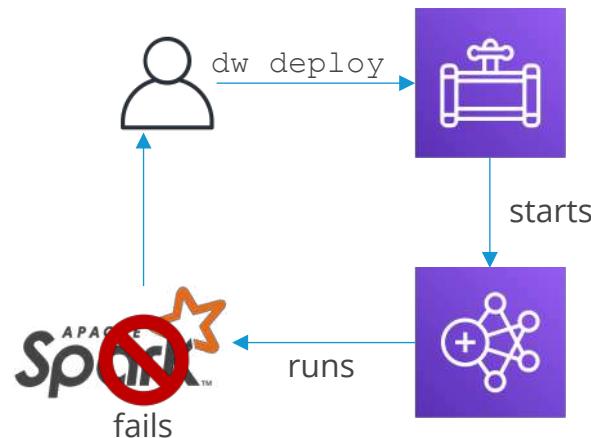


+

```
{
  "fields": [
    {
      "name": "id",
      "type": "long"
    },
    {
      "name": "make",
      "type": "string"
    },
    {
      "name": "model",
      "type": "string"
    },
    {
      "name": "price",
      "type": "double"
    },
    {
      "name": "mileage",
      "type": "double"
    }
  ]
}
```

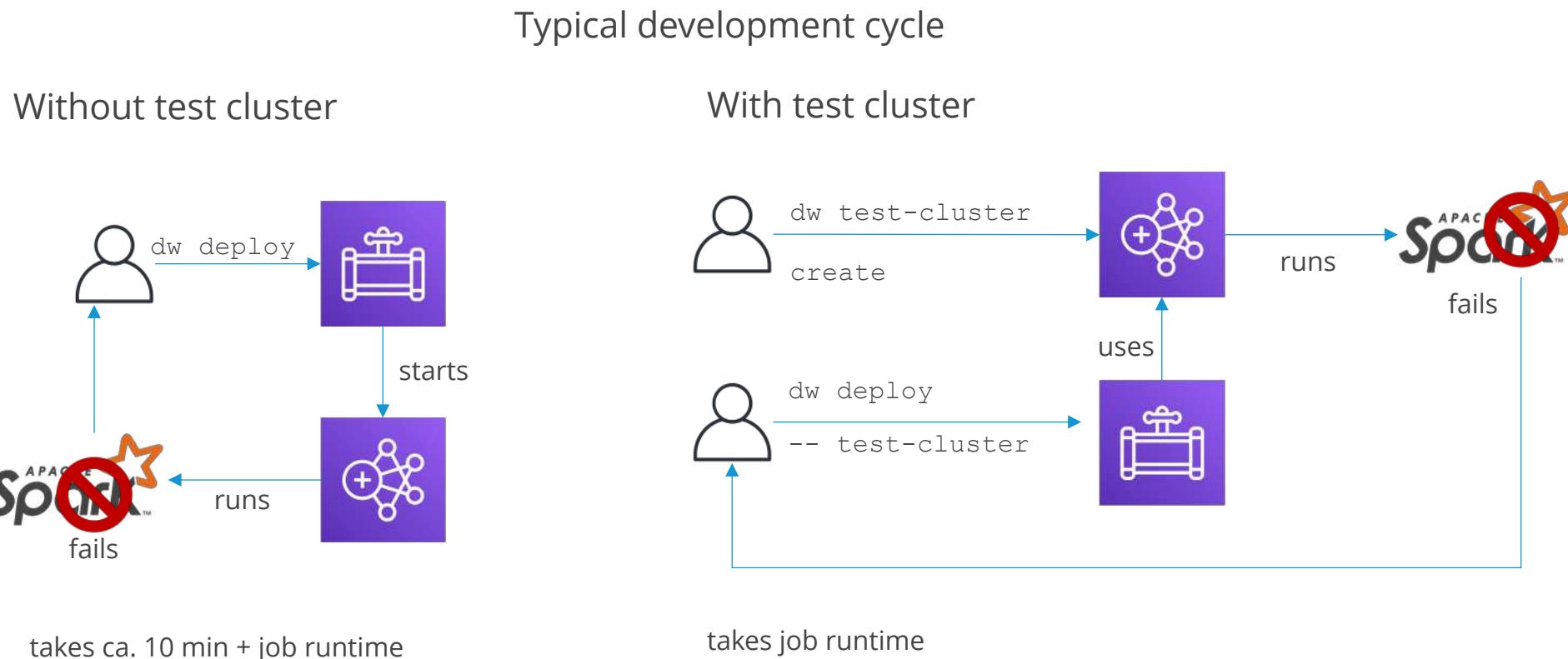
# Speed up pipeline development

Typical development cycle



takes ca. 10 min + job runtime

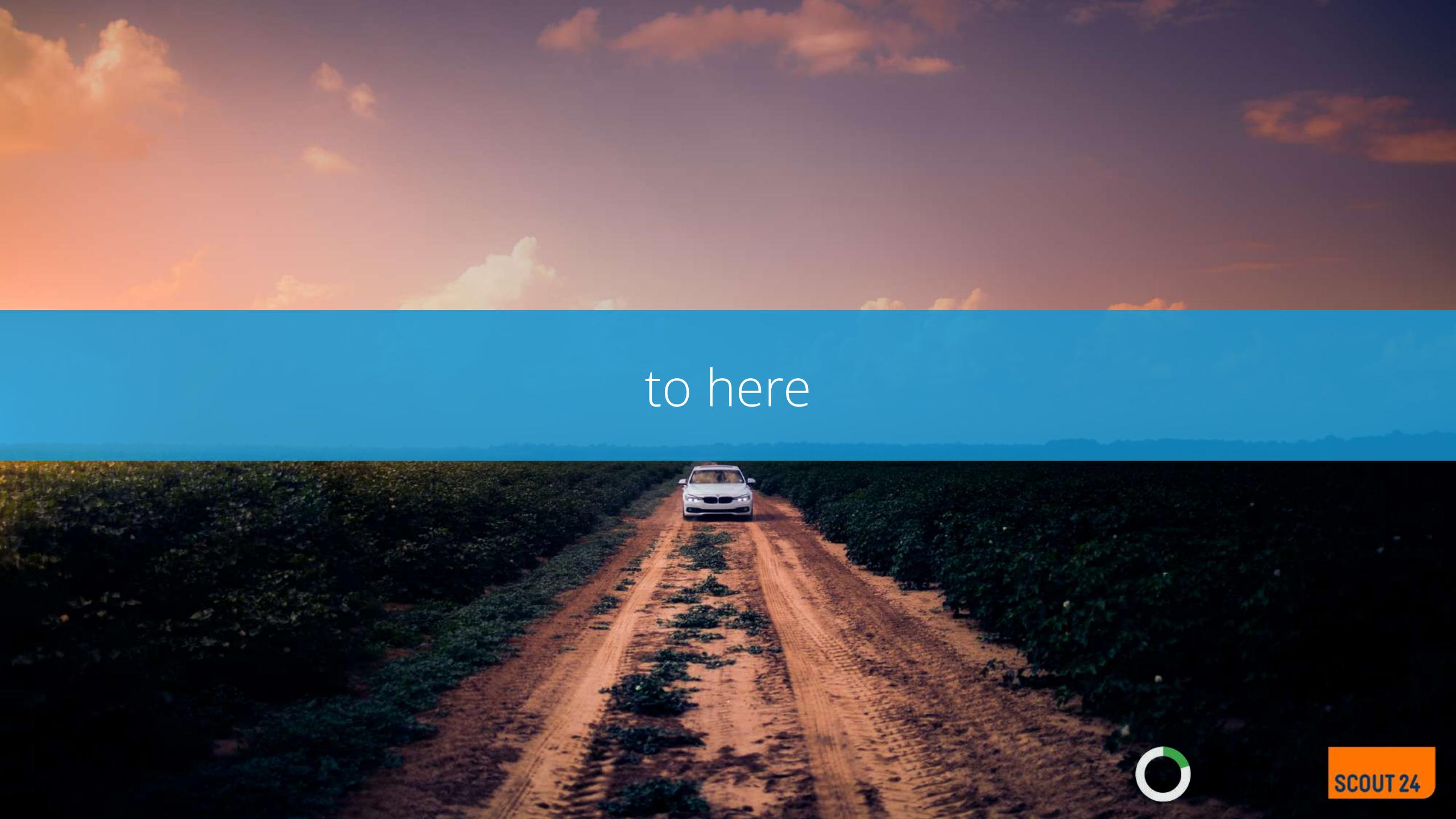
# Speed up test and debug cycle





DataWario brought us from here ...

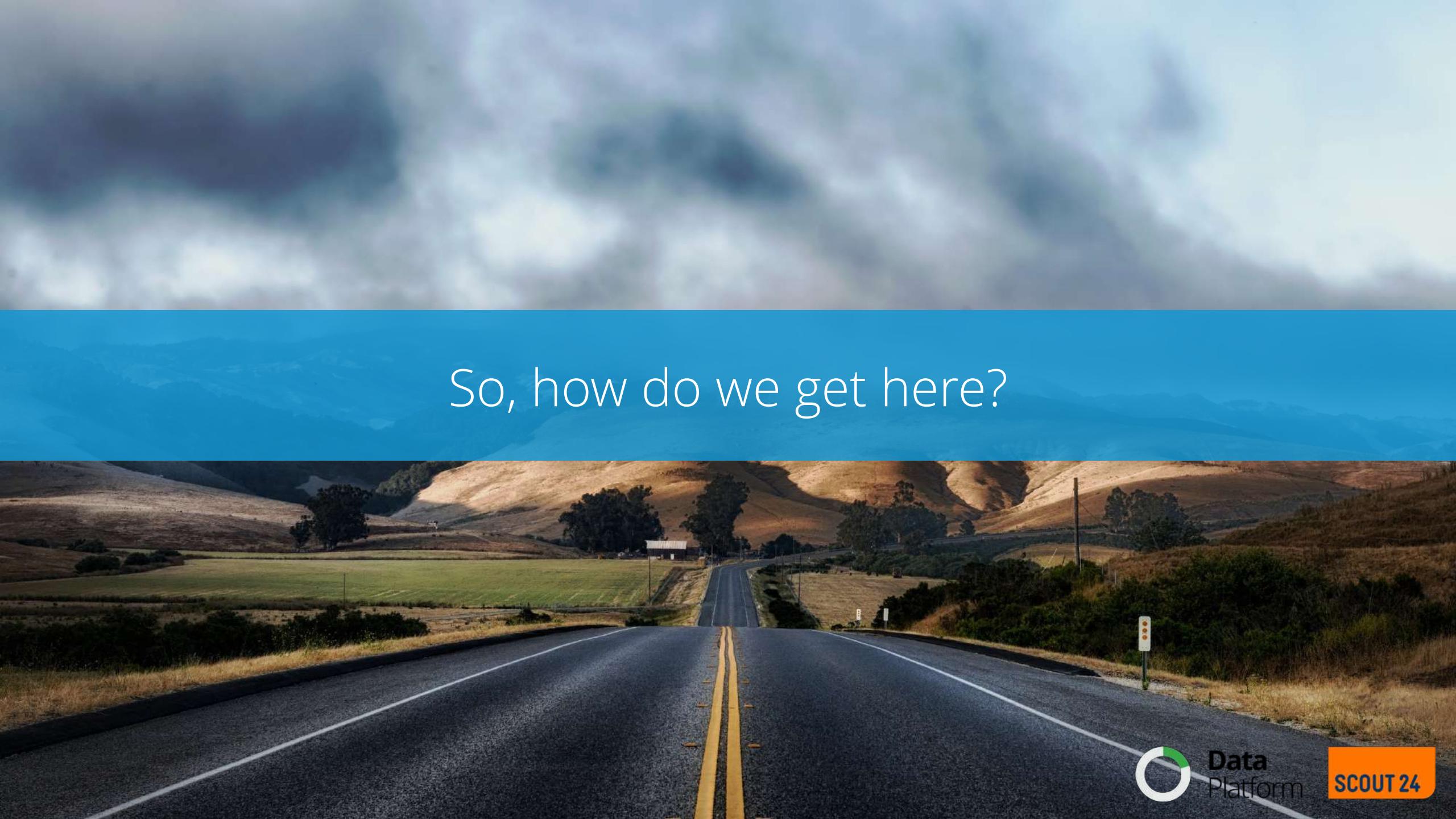
photo by Dirk van der Made ([https://commons.wikimedia.org/wiki/File:DirkvdM\\_cloudforest-jungle.jpg](https://commons.wikimedia.org/wiki/File:DirkvdM_cloudforest-jungle.jpg)), „DirkvdM cloudforest-jungle“, <https://creativecommons.org/licenses/by-sa/3.0/legalcode>



to here



SCOUT24



So, how do we get here?

# Paving the path further

- We built a lot of pipelines, so did our users
- We see repetitive tasks emerge in these pipelines
- Reimplementing solutions over and over again is waste
- We don't like waste

➔ Let's build reusable tools for common tasks

Much of the power of the UNIX operating system comes from a style of program design that makes programs easy to use and, more important, **easy to combine with other programs**. This style has been called the use of *software tools*, and depends more on how the programs fit into the programming environment and how they can be used with other programs than on how they are designed internally. [...] This style was based on the use of **tools**: using programs separately or in combination to get a job done, rather than doing it by hand, by monolithic self-sufficient subsystems, or by special-purpose, one-time programs.

---

*Rob Pike; Brian W. Kernighan (October 1984). "[Program Design in the UNIX Environment](#)"*

# Tools that do one thing, and do that well

- Advantages:
  - Easy to configure, only a few settings
  - Easy to test, because there are less branches to cover
  - Easy to debug, output of each tool can be checked independently
- Drawbacks:
  - More intermittent data written
  - longer runtimes

# Event snapshotter

- Typical questions: What was the headline of a listing on Feb 22, 2016? Was it published? How about March 1?
- Input data: Change events for an entity (e.g. listing)
  - User changed headline of a listing, number of pictures
  - Listing got published/unpublished
- Snapshotter
  1. takes snapshot from the day before
  2. merges all change events from the day, keeps last
  3. writes new snapshot for day
- Simple config (input and output path, name of ID and timestamp column, output path, partition schema)
- No business logic, agnostic of type of entity

# Event Snapshotter Example

```
{"id": 1, "changed": "2019-01-01T15:31:00.000+01:00", "headline": "nice appartment", "published": false}  
{"id": 1, "changed": "2019-01-01T18:31:00.000+01:00", "headline": "A very nice appartment", "published": false}  
{"id": 1, "changed": "2019-01-02T06:31:00.000+01:00", "headline": "A very nice appartment", "published": true}  
{"id": 2, "changed": "2019-01-02T07:28:00.000+01:00", "headline": "Another nice appartment", "published": true}
```

January 1, 2019 (s3://snapshotBucket/snapshot\_day=2019-01-01):

```
{"id": 1, "changed": "2019-01-01T18:31:00.000+01:00", "headline": "A very nice appartment", "published": false}
```

January 2, 2019 (s3://snapshotBucket/snapshot\_day=2019-01-02):

```
{"id": 1, "changed": "2019-01-02T06:31:00.000+01:00", "headline": "A very nice appartment", "published": true}
```

```
{"id": 2, "changed": "2019-01-02T07:28:00.000+01:00", "headline": "Another nice appartment", "published": true}
```

# Aggregator

- Typical Questions:
  - How many published listings did we have on Feb 1 2018?
  - What's the daily average number of pageviews per car make?
- Input data:
  - Interaction events of an entity ( e.g. page viewed, email sent, number called)
  - output of Event Snapshotter
- Simple config (input and output paths, timestamp column, grouping column, aggregation method and column)

# Spark Utils

- Collection of simple tools (usually one class per tool)
- Can be used
  - standalone in a pipeline step
  - as a dependency of your code
- Examples:
  - Create a Hive table on-top of existing data
  - Convert data from one Format to another (e.g. from JSON to Parquet for faster querying)
  - Flatten nested structures
  - publish data quality metrics to CloudWatch
  - Materializing views

# Current State / Outlook

- Data Pipeline with DataWario widely used throughout the company
- Challenges:
  - Data Pipeline not getting much love from AWS (not deprecated, but not getting many new features)
  - Coordination between data pipelines only rudimentary
- Possible solutions:
  - AWS Step Functions in combination with AWS Glue
  - Introduction of a Workflow tool (e.g. AirFlow)

# Self-Service Analytics

# Query Challenges



# What's Ahead

Unlock the Datalake for Scout24's Toolset and Users with Different Skillsets

Data Analysis for Various User Groups

Provide a Timely and Accurate Update of the Metadata Layer

# What's Ahead

Unlock the Datalake for Scout24's Toolset and Users with Different Skillsets

Data Analysis for Various User Groups

Provide a Timely and Accurate Update of the Metadata Layer



OneScout Hive Metastore

# What's Ahead

Unlock the Datalake for Scout24's Toolset and Users with Different Skillsets

Data Analysis for Various User Groups

Provide a Timely and Accurate Update of the Metadata Layer



OneScout Hive Metastore

Personal Analytics Cluster

# What's Ahead

Unlock the Datalake for Scout24's Toolset and Users with Different Skillsets



Data Analysis for Various User Groups



Provide a Timely and Accurate Update of the Metadata Layer

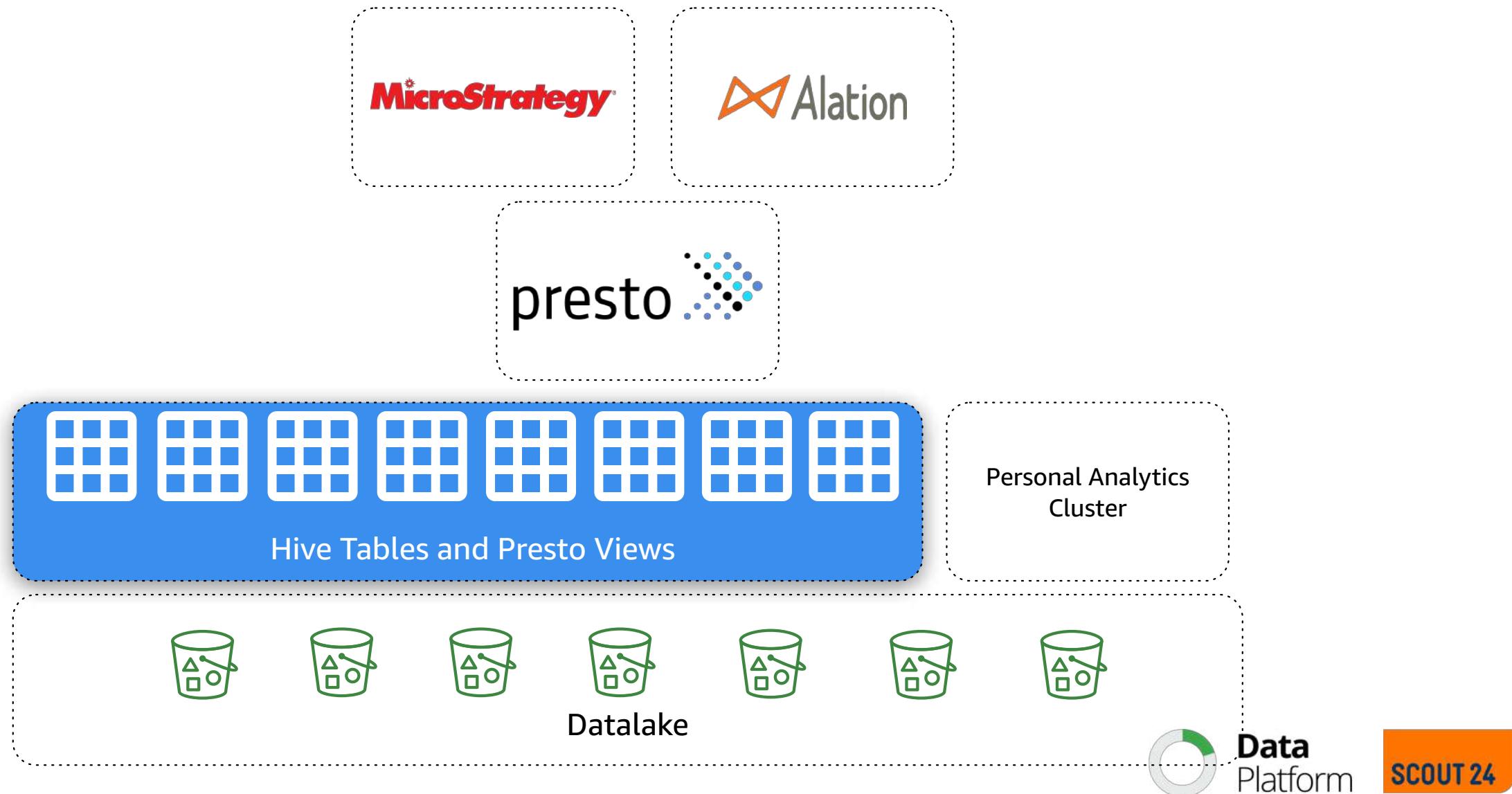


OneScout Hive Metastore

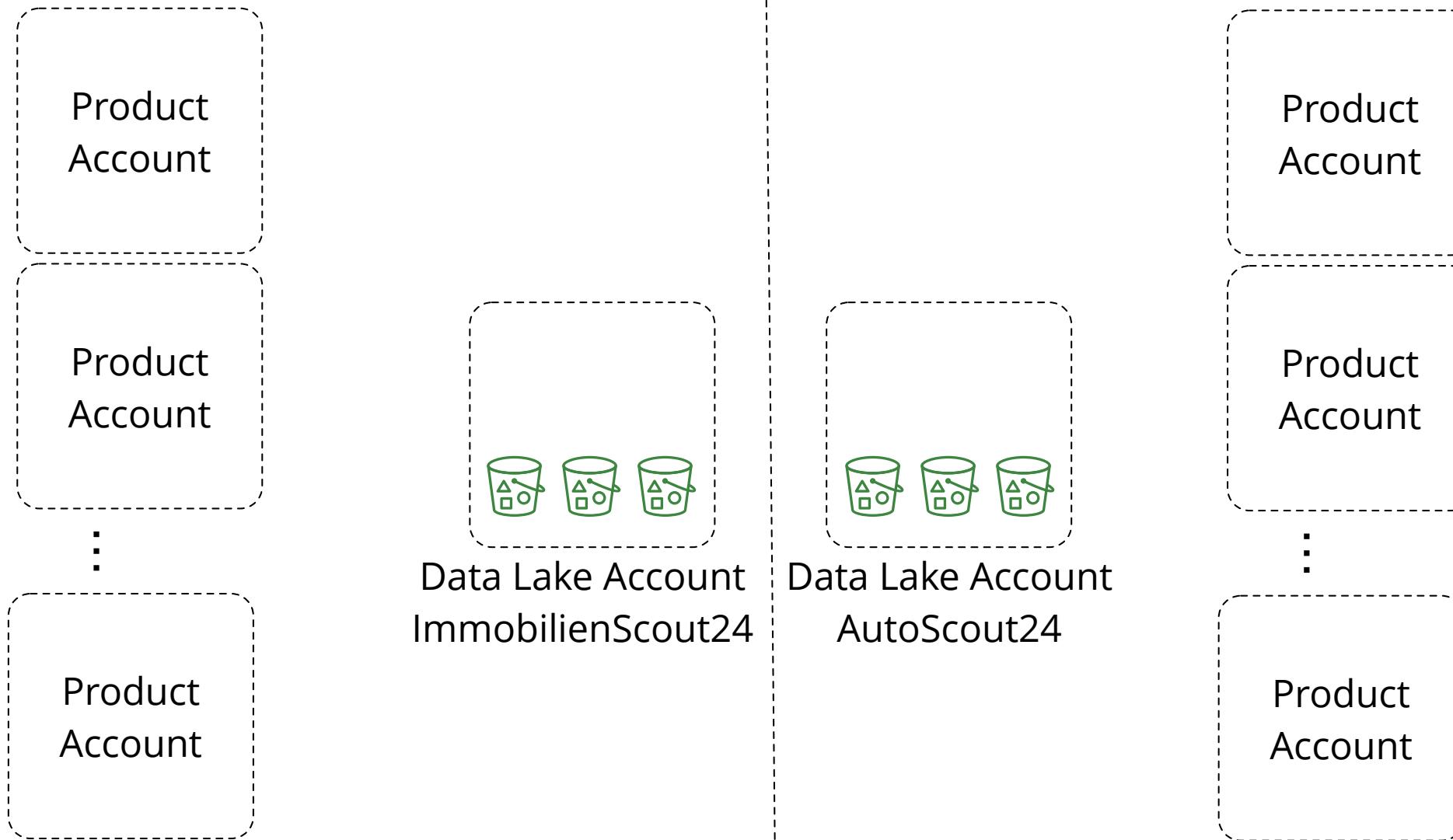
Personal Analytics Cluster

Automatic Hive Partition Detection

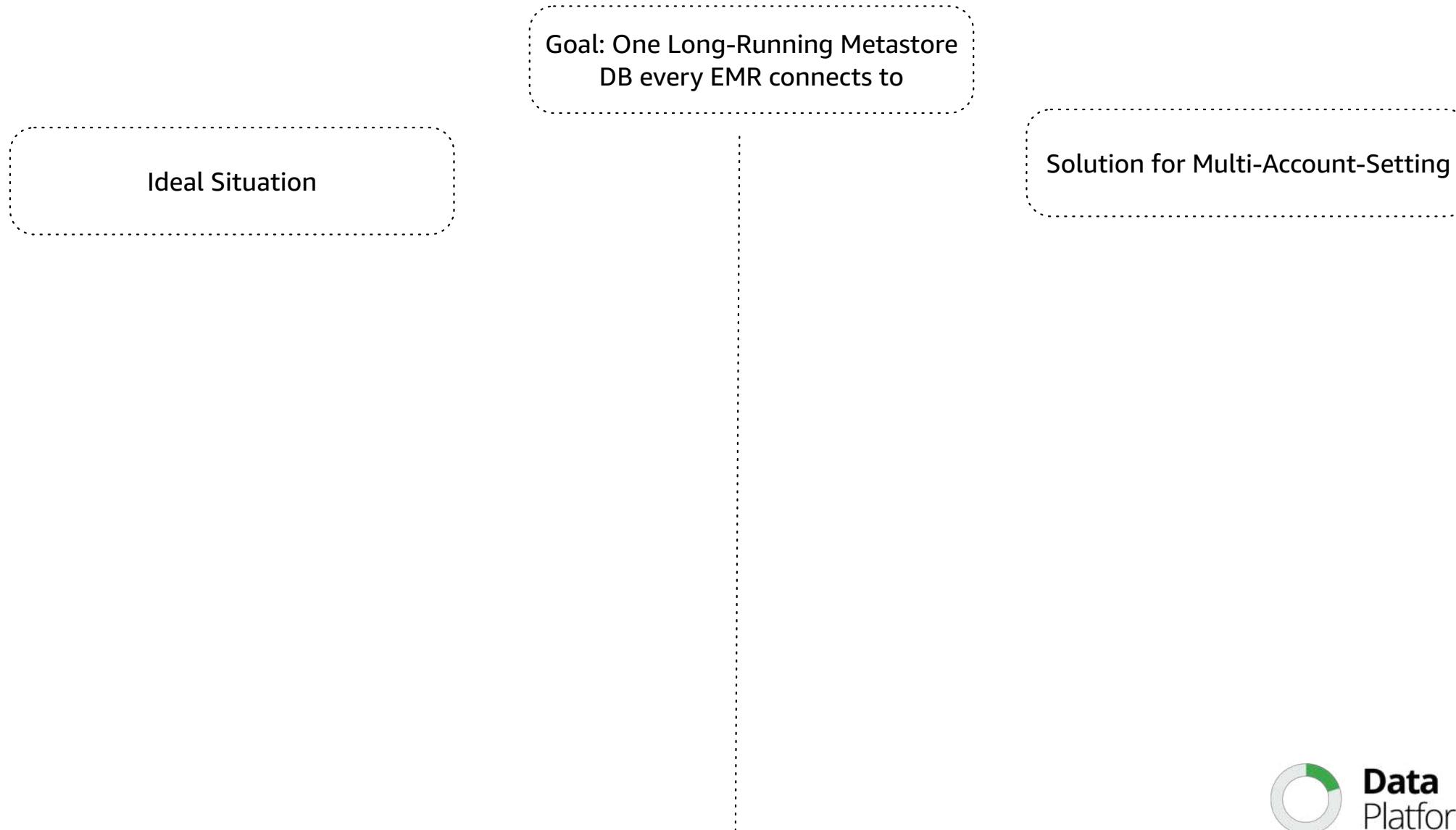
# OneScout Hive Metastore – A Schematic View



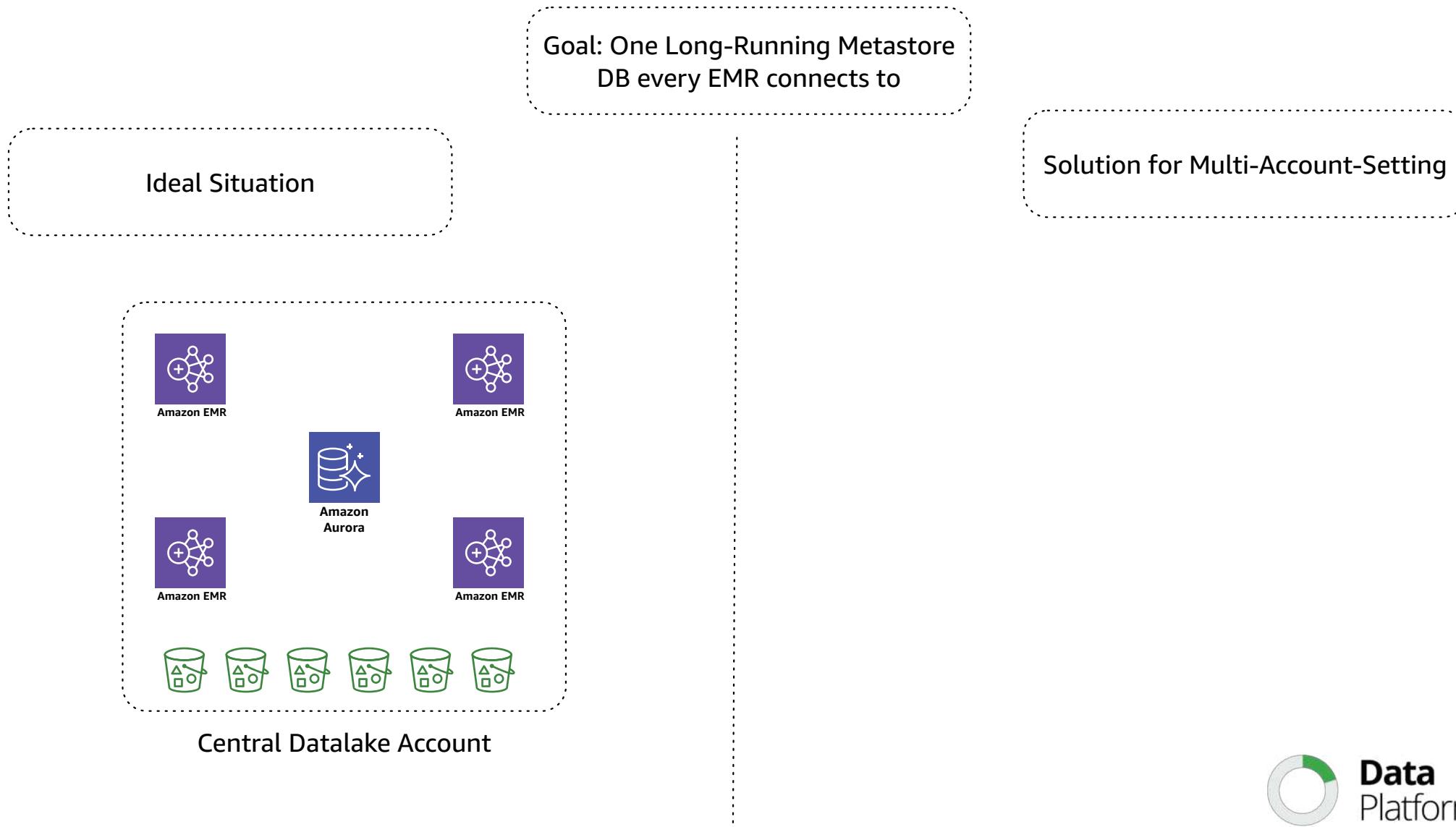
# OneScout Hive Metastore – Recap of Ecosystem



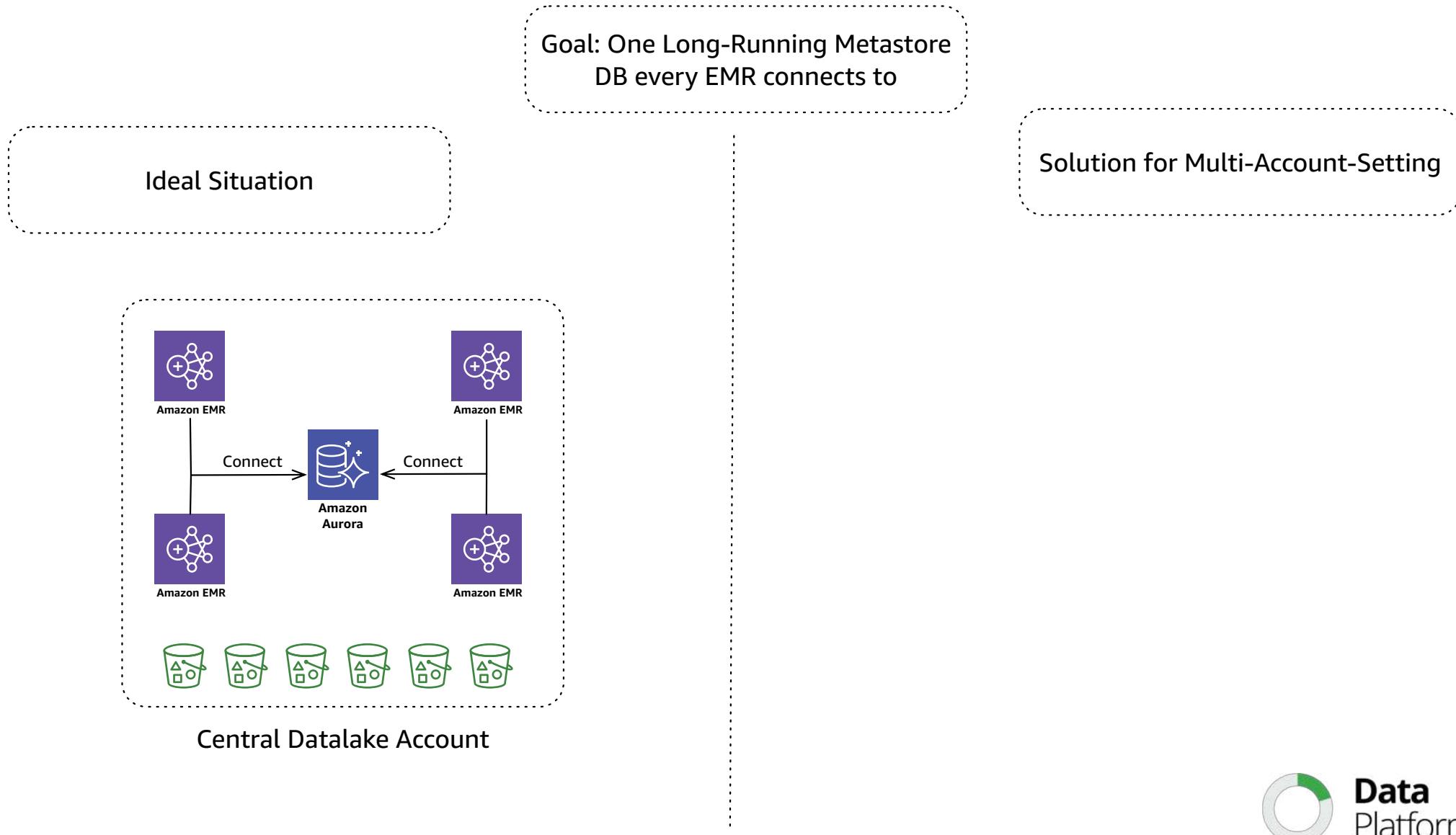
# The Scout24 Hive Metastore Proxy – A Motivation



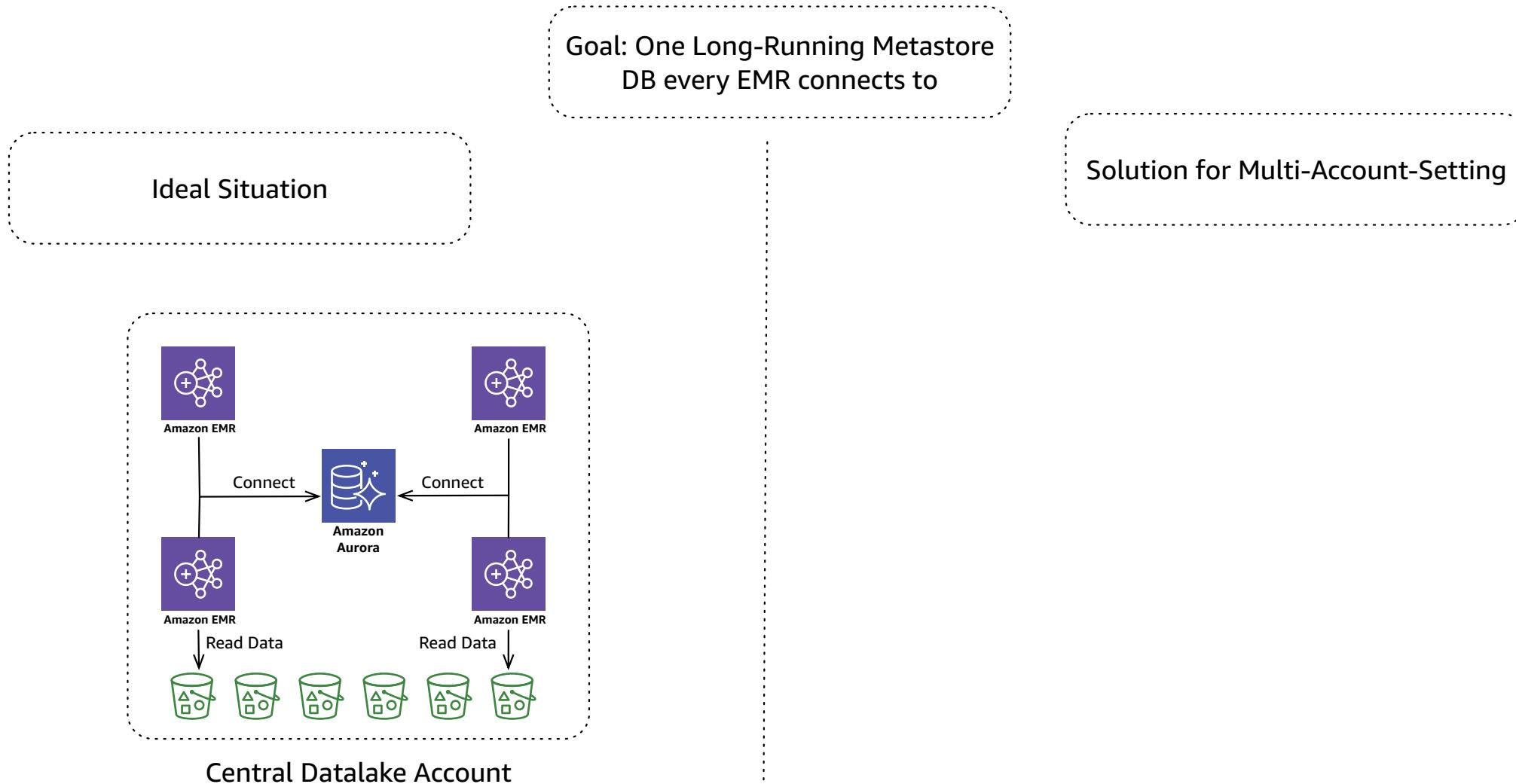
# The Scout24 Hive Metastore Proxy – A Motivation



# The Scout24 Hive Metastore Proxy – A Motivation



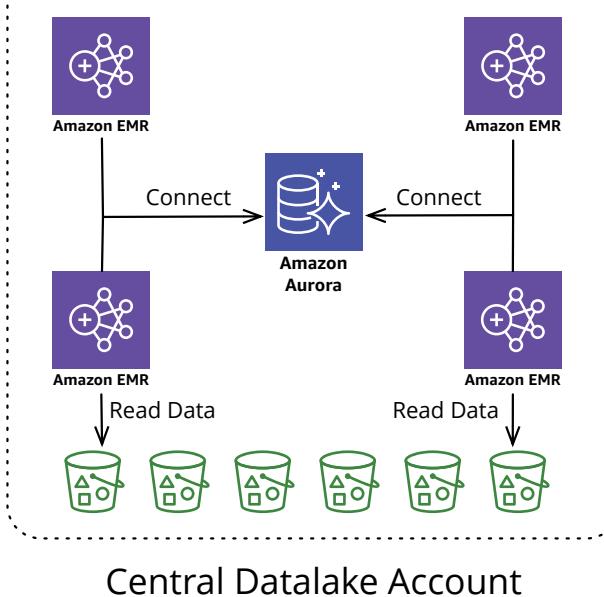
# The Scout24 Hive Metastore Proxy – A Motivation



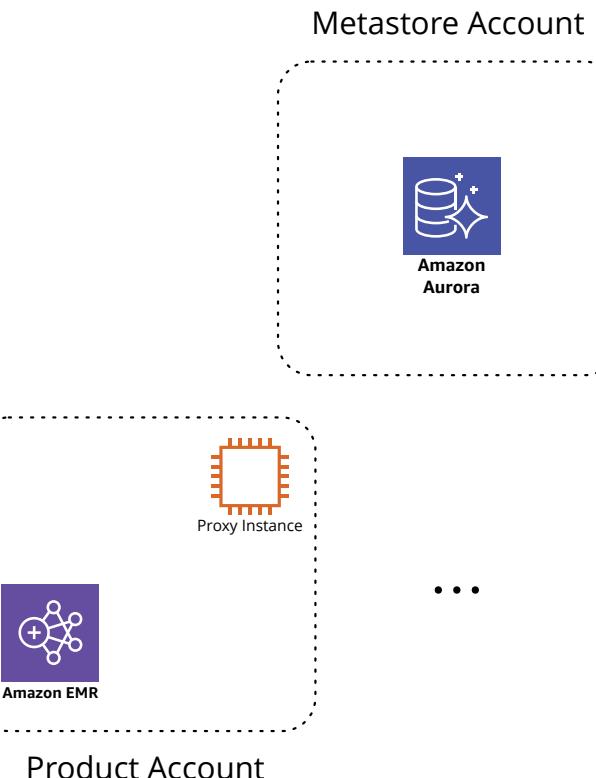
# The Scout24 Hive Metastore Proxy – A Motivation

Goal: One Long-Running Metastore DB every EMR connects to

Ideal Situation



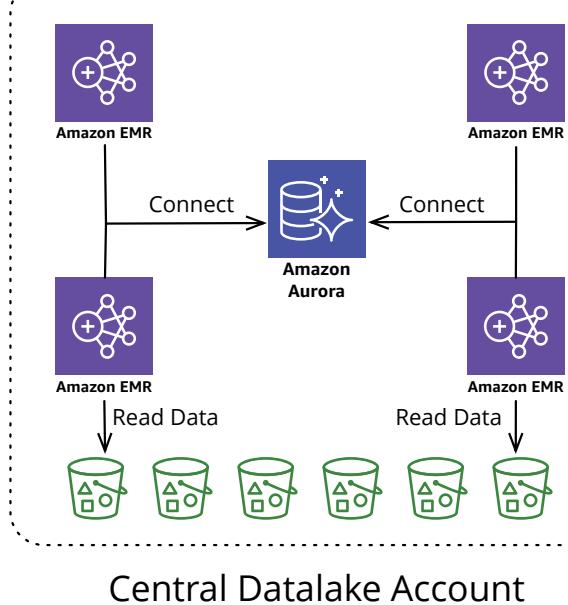
Solution for Multi-Account-Setting



# The Scout24 Hive Metastore Proxy – A Motivation

Goal: One Long-Running Metastore DB every EMR connects to

Ideal Situation



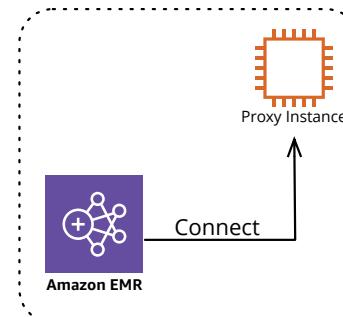
Central Datalake Account

Solution for Multi-Account-Setting

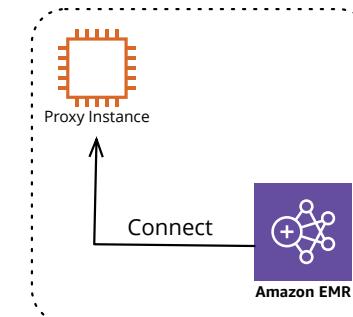
Metastore Account



Amazon Aurora

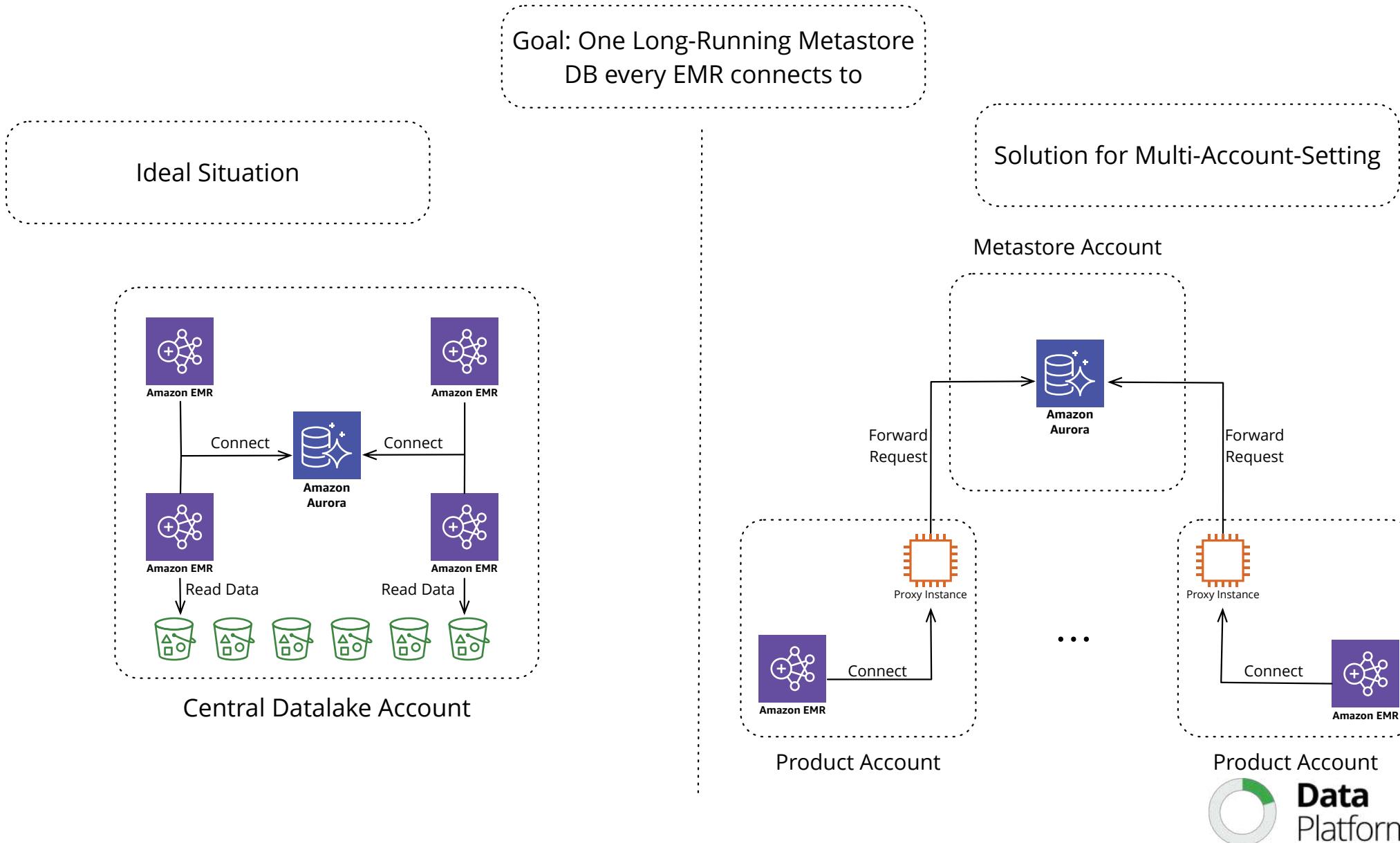


Product Account

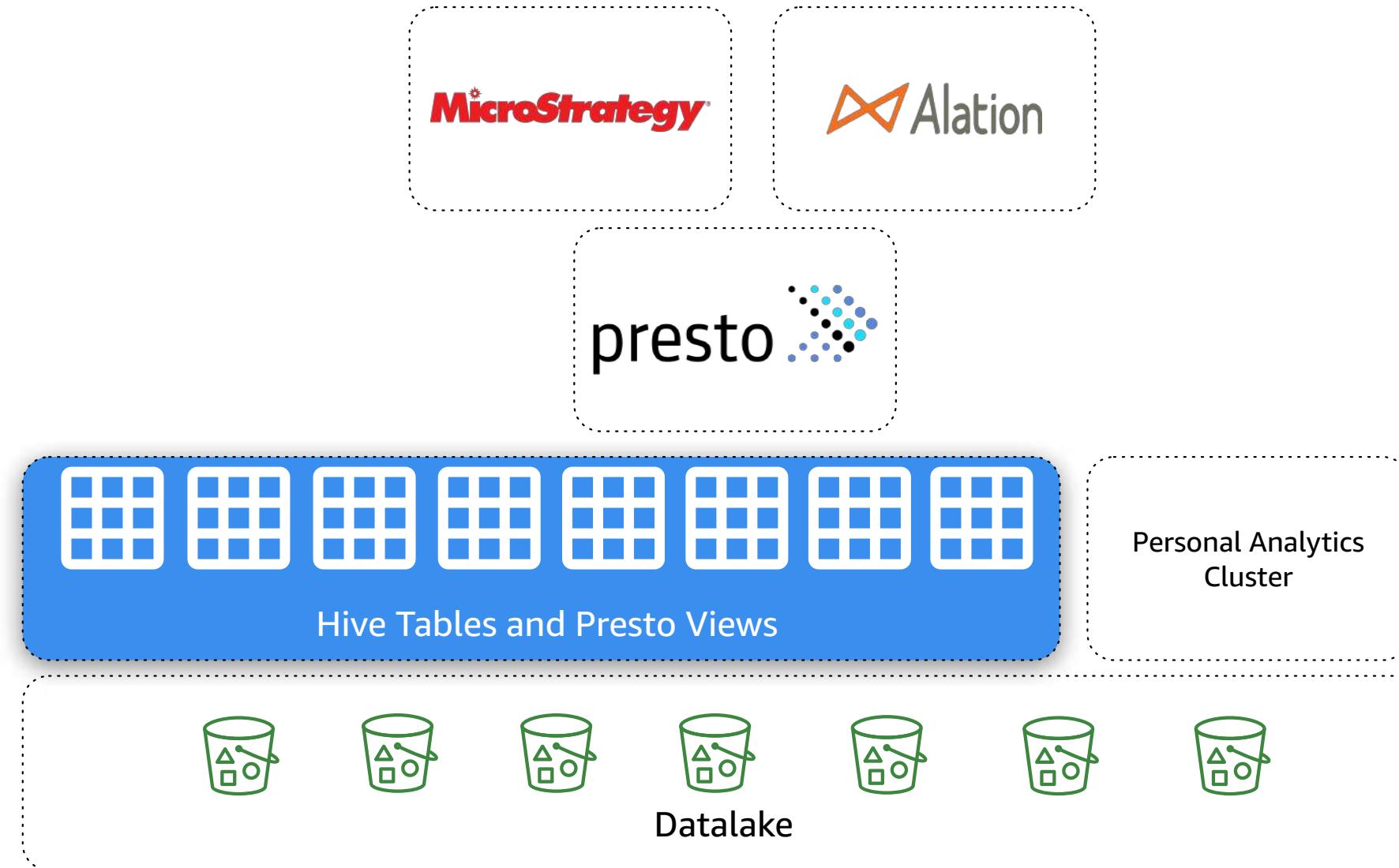


Product Account

# The Scout24 Hive Metastore Proxy – A Motivation



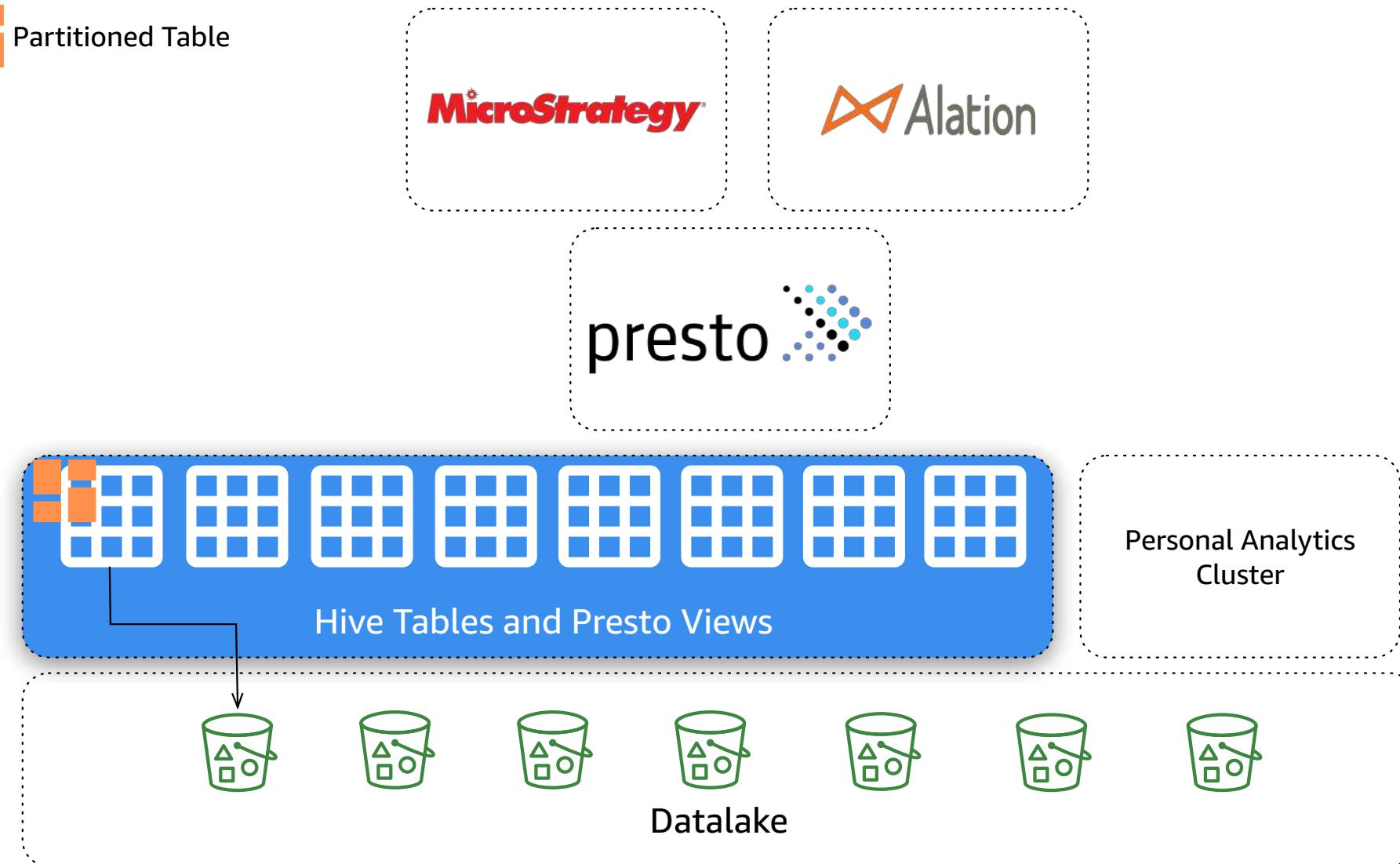
# Automated Partition Detection – A Motivation



# Automated Partition Detection – A Motivation

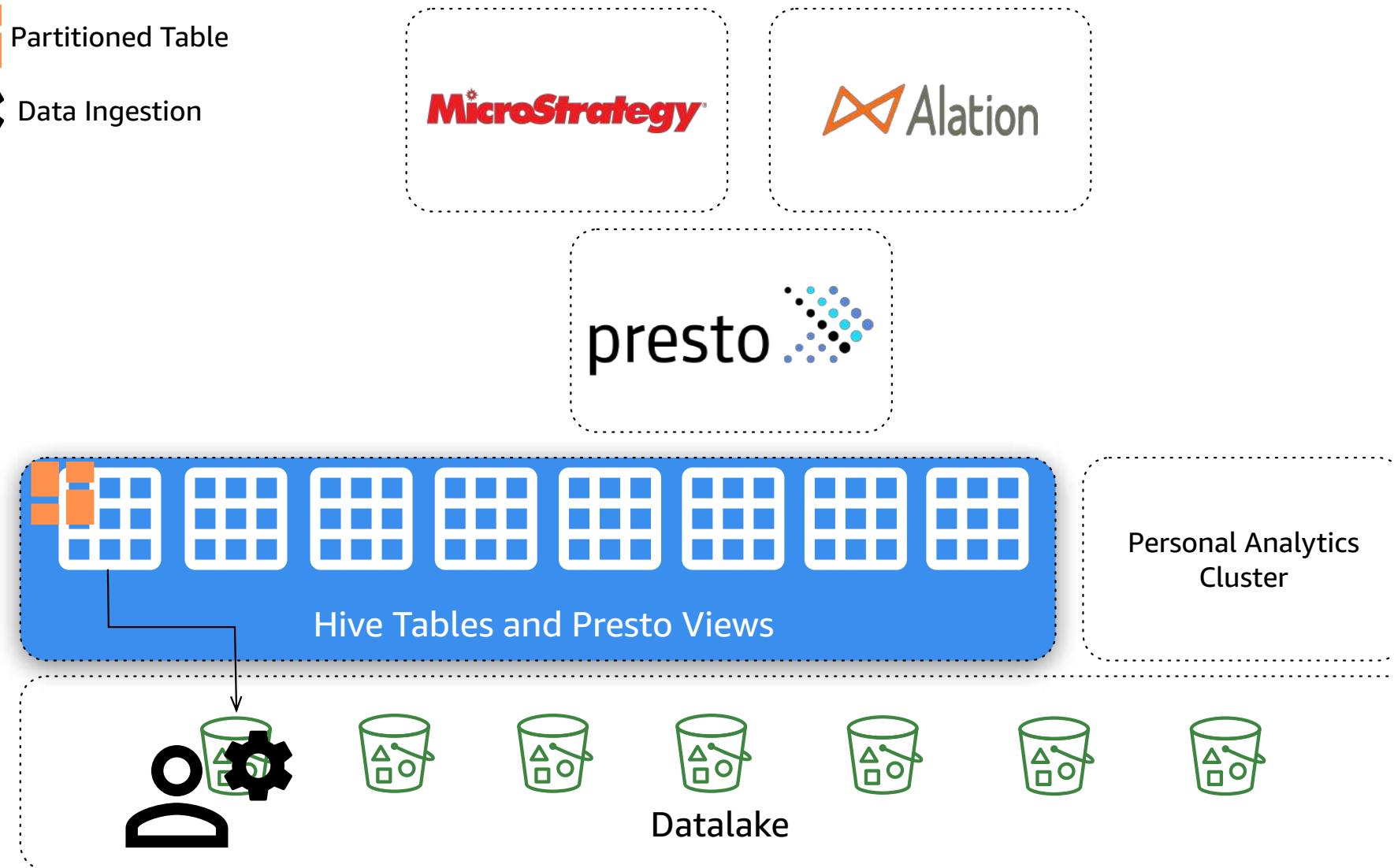


Partitioned Table



# Automated Partition Detection – A Motivation

-  Partitioned Table
-  Data Ingestion



# Automated Partition Detection – A Motivation



Partitioned Table



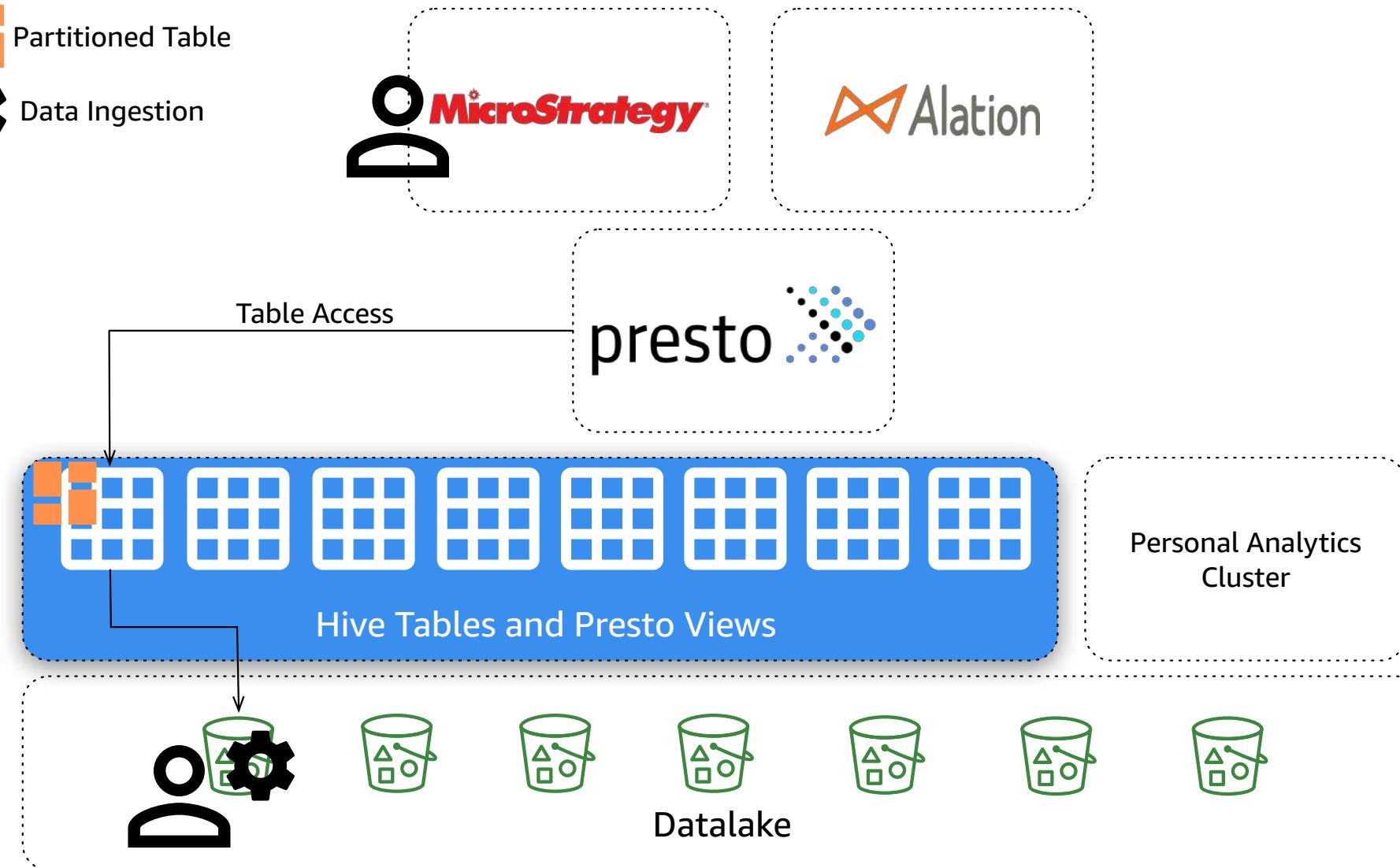
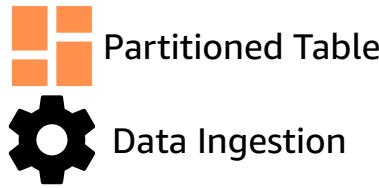
Data Ingestion



Personal Analytics  
Cluster



# Automated Partition Detection – A Motivation



# Automated Partition Detection – A Motivation



Partitioned Table

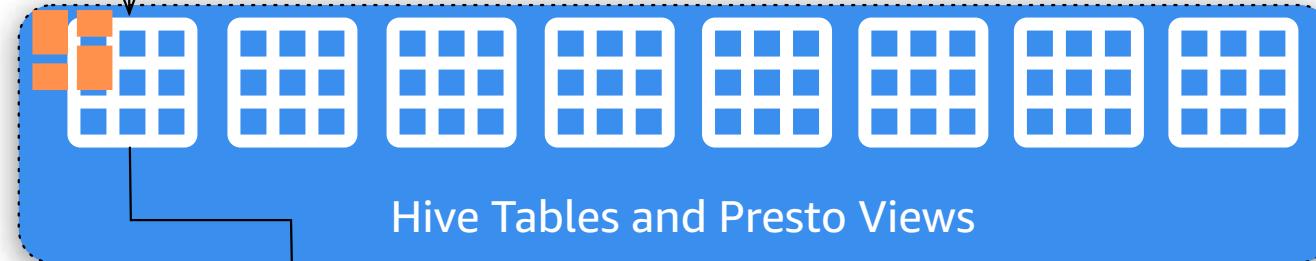


Data Ingestion



Table Access

presto



Personal Analytics  
Cluster



# Automated Partition Detection – A Motivation



Partitioned Table

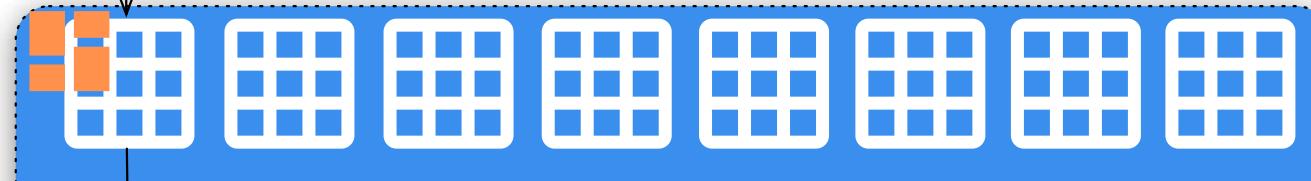


Data Ingestion



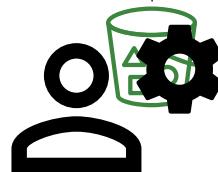
Table Access

presto



Personal Analytics Cluster

Automatic Partition Detection



Datalake

# Automated Partition Detection – A Motivation



Partitioned Table



Data Ingestion



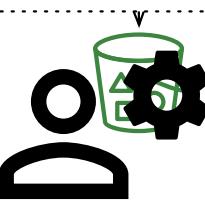
Table Access

presto



Personal Analytics Cluster

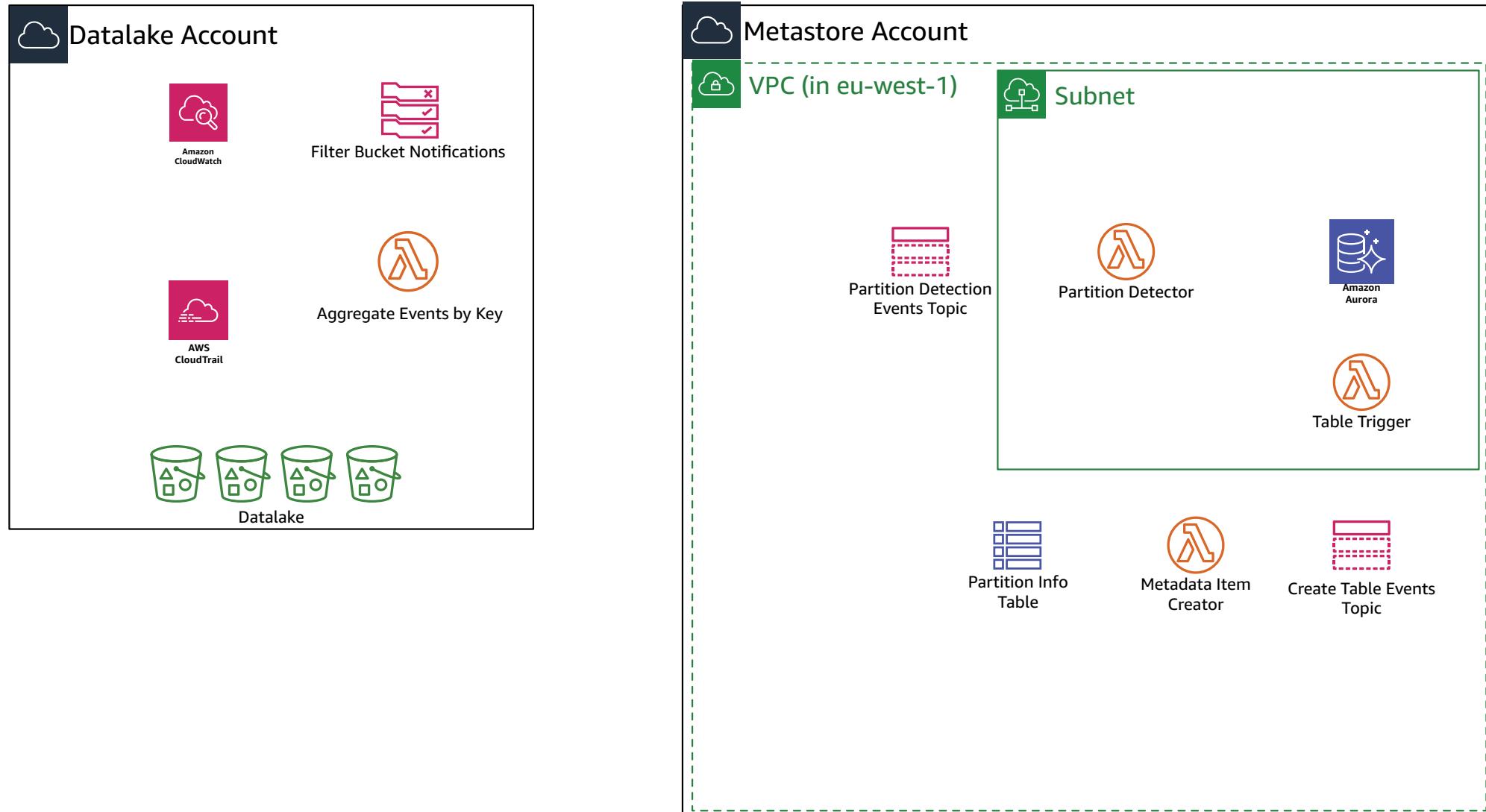
Automatic Partition Detection



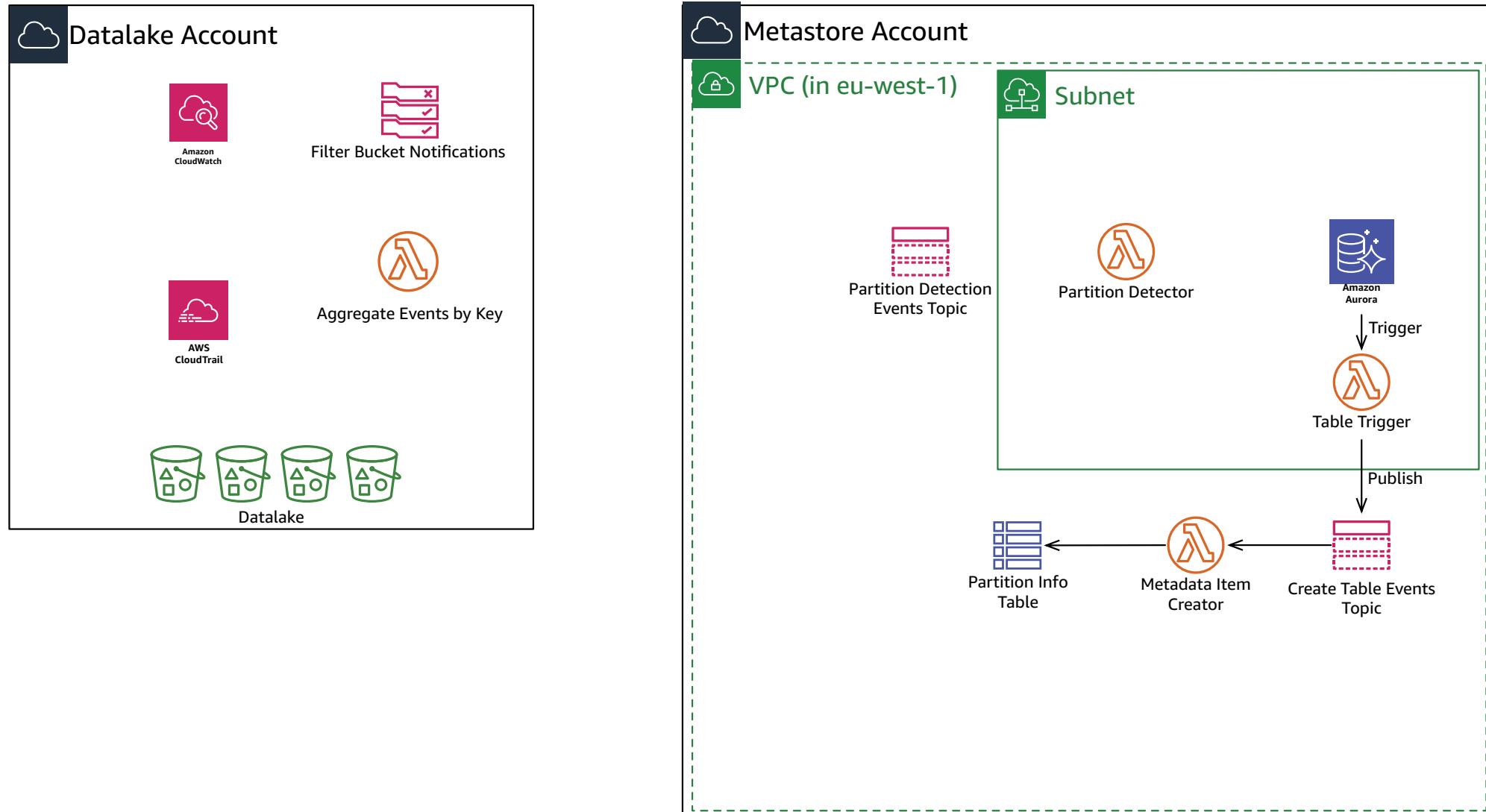
Datalake



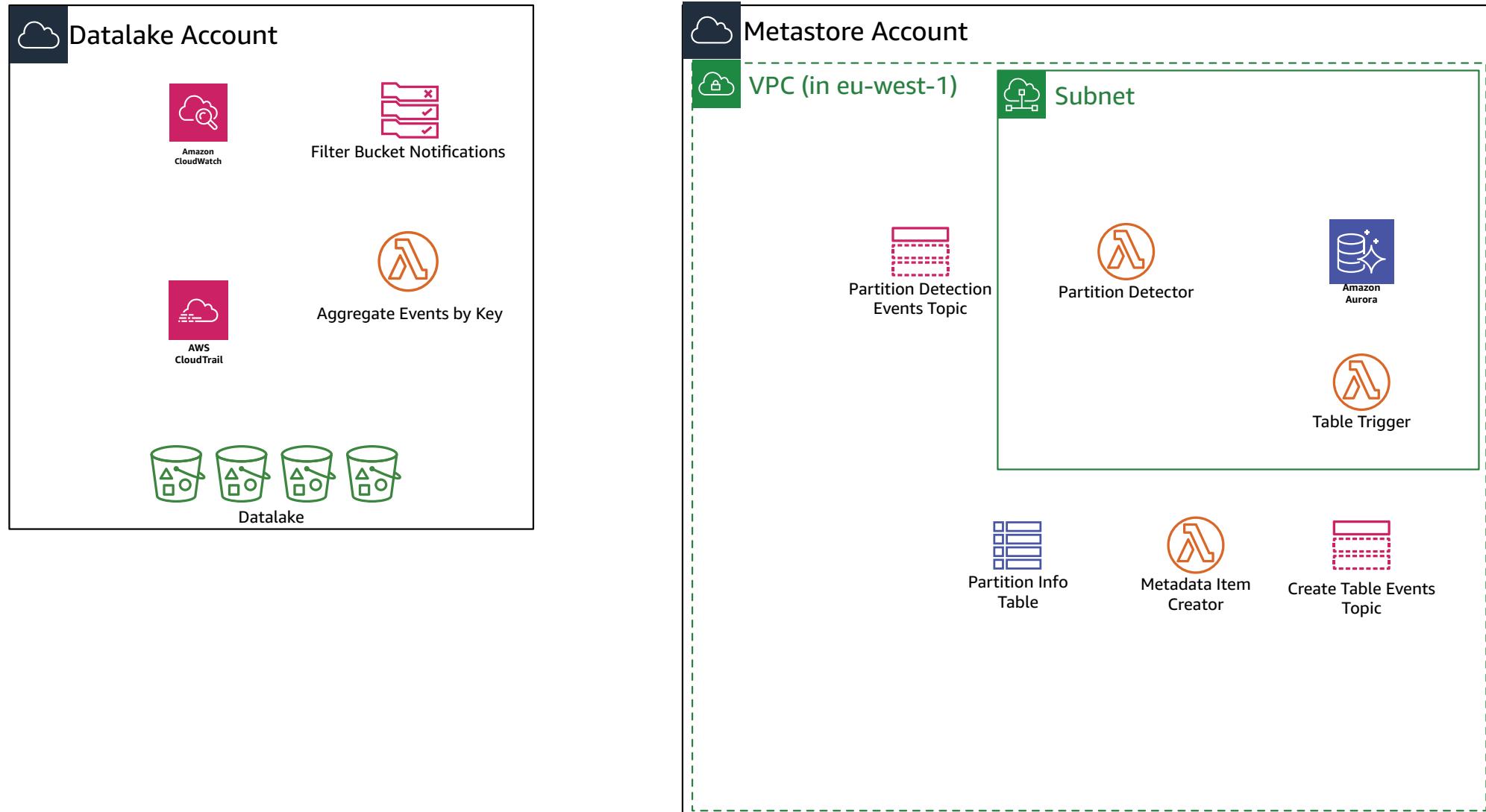
# Partition Detection Architecture



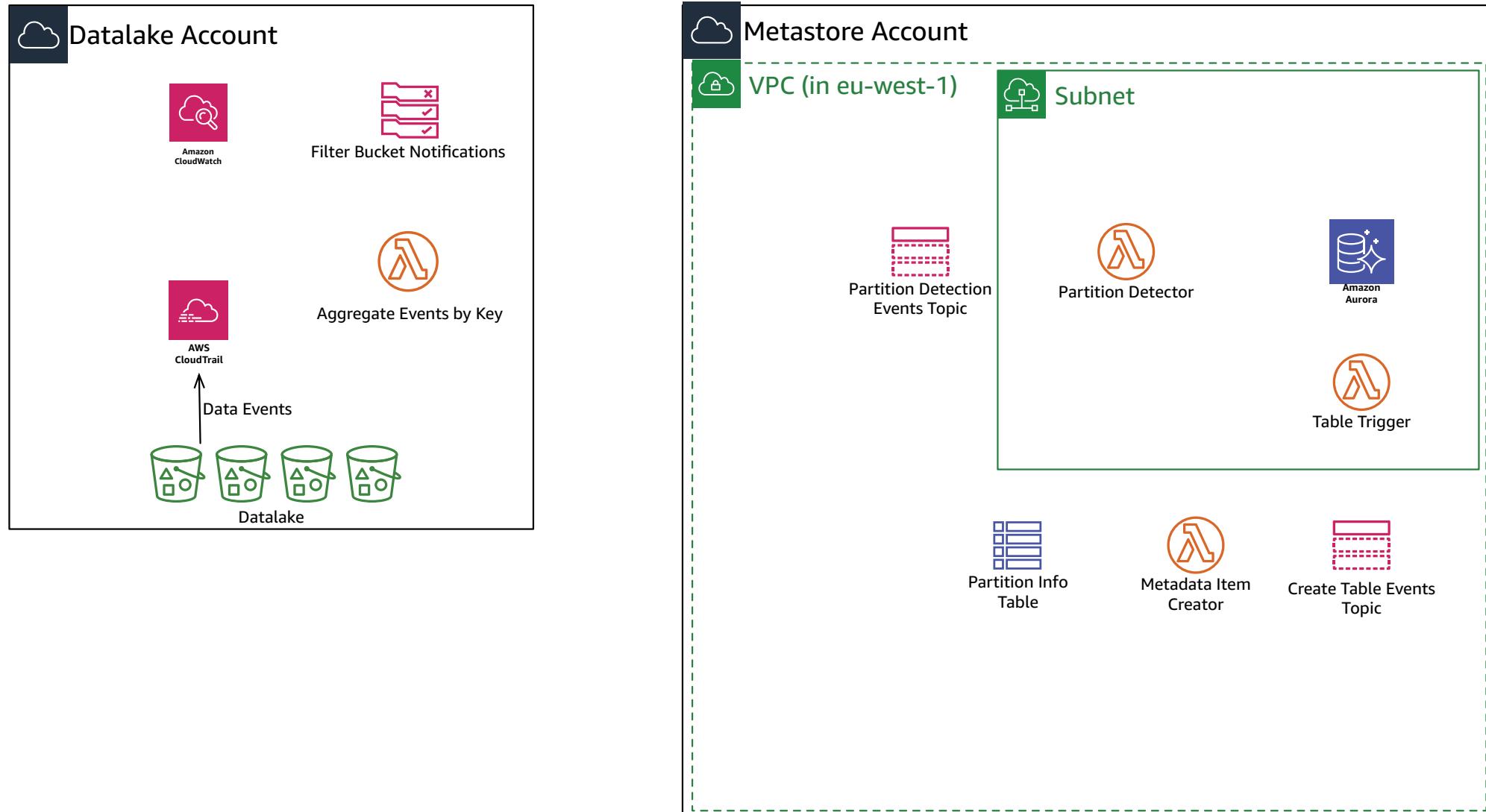
# Partition Detection Architecture



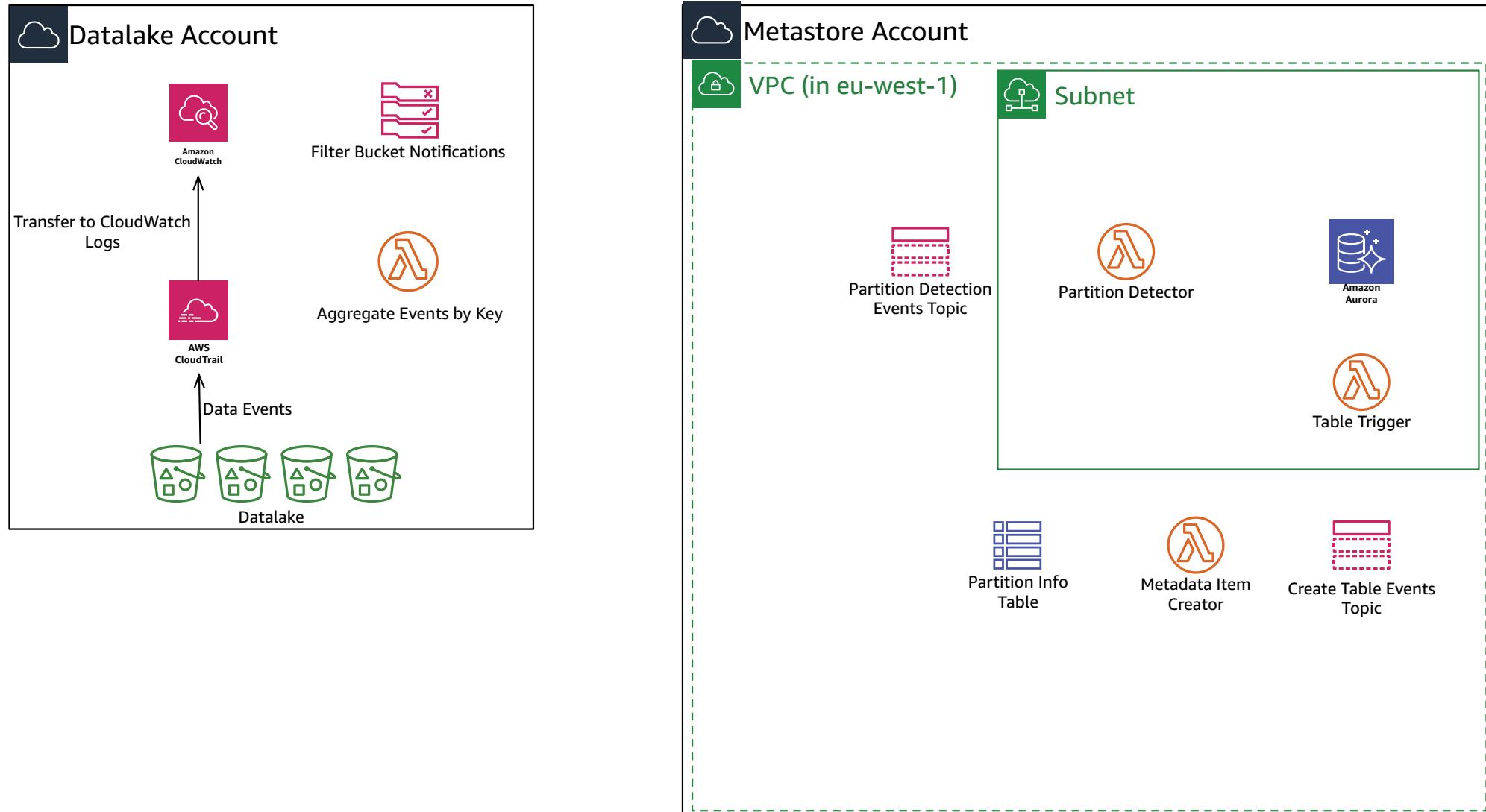
# Partition Detection Architecture



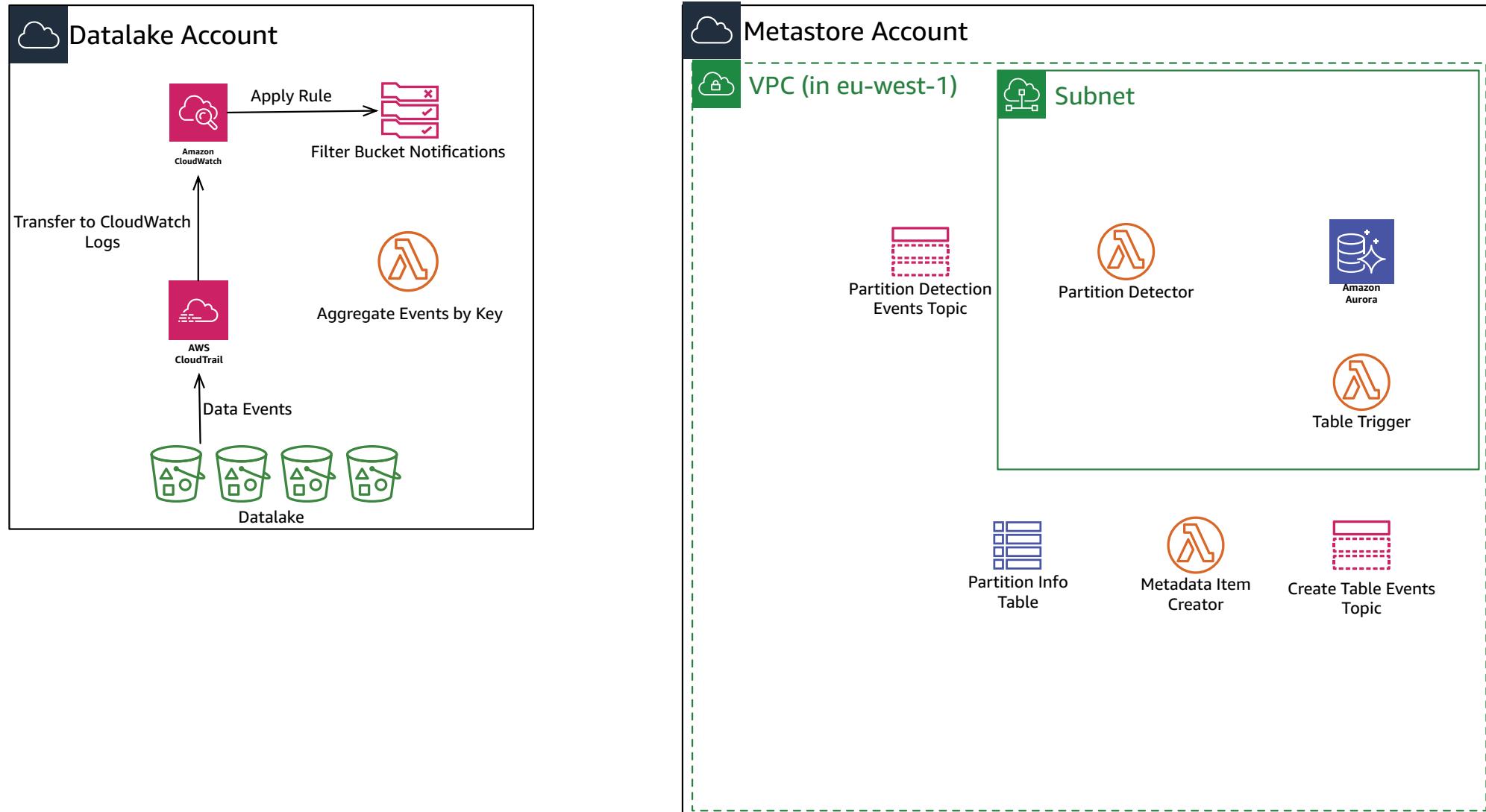
# Partition Detection Architecture



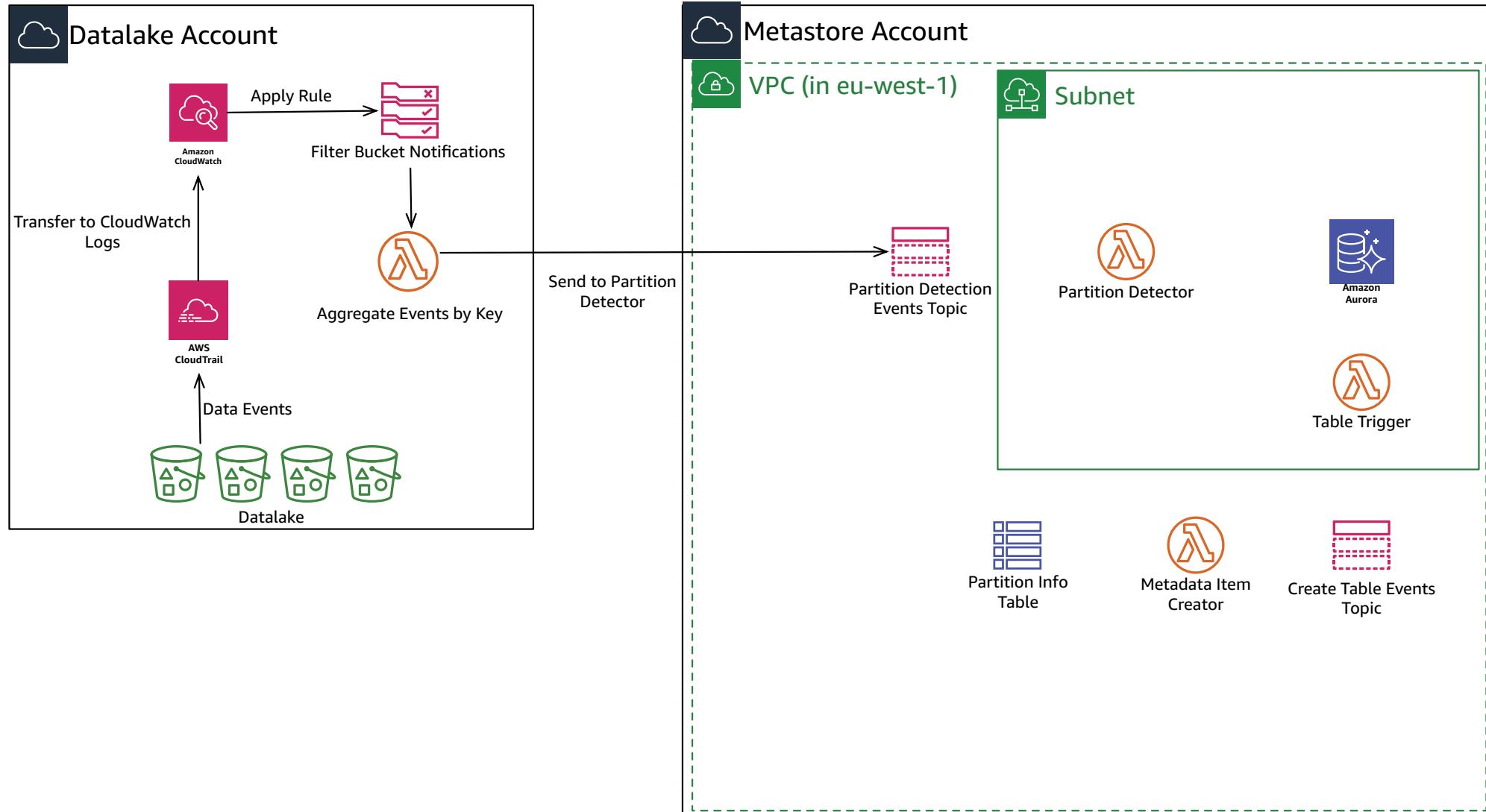
# Partition Detection Architecture



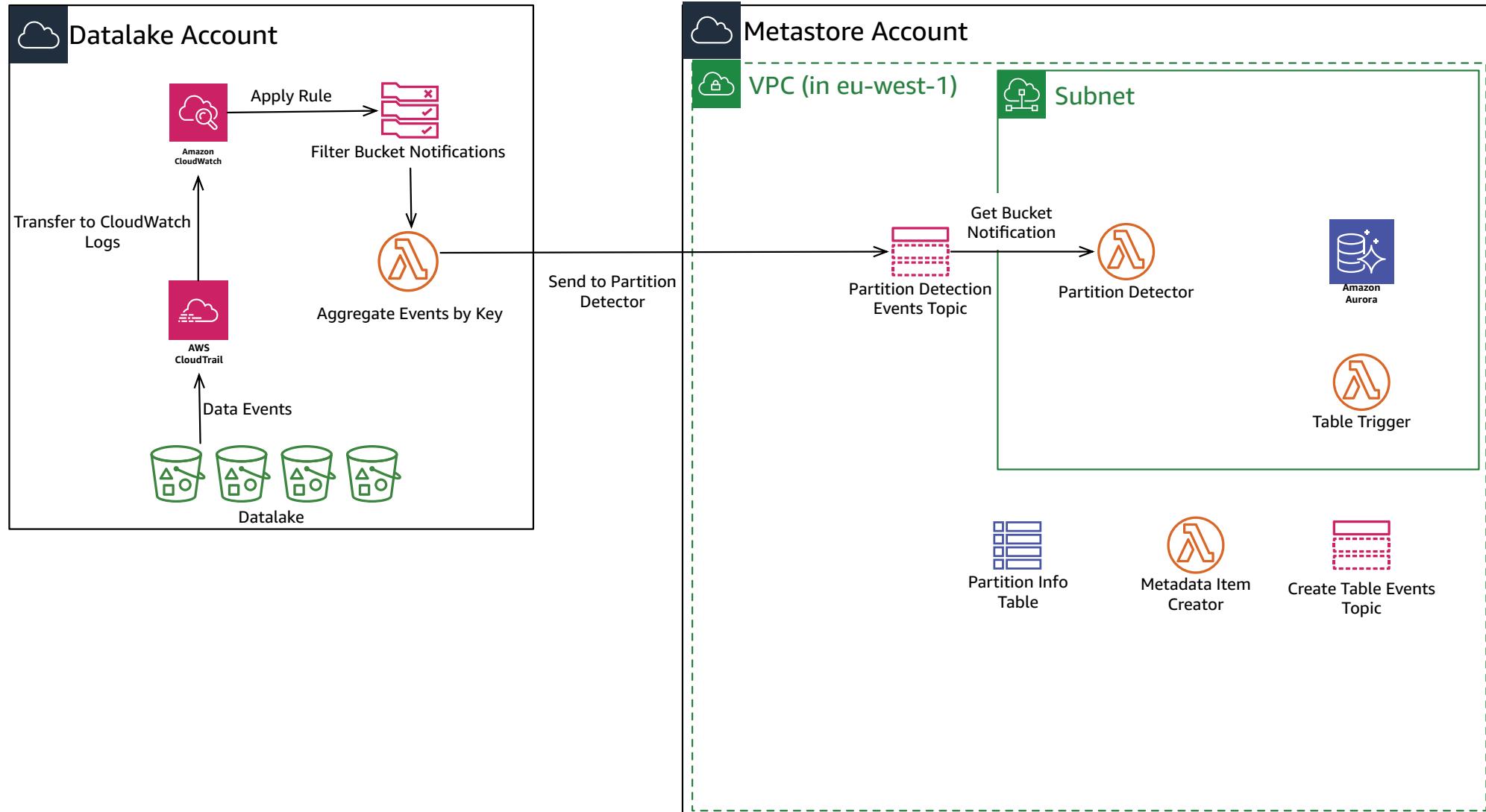
# Partition Detection Architecture



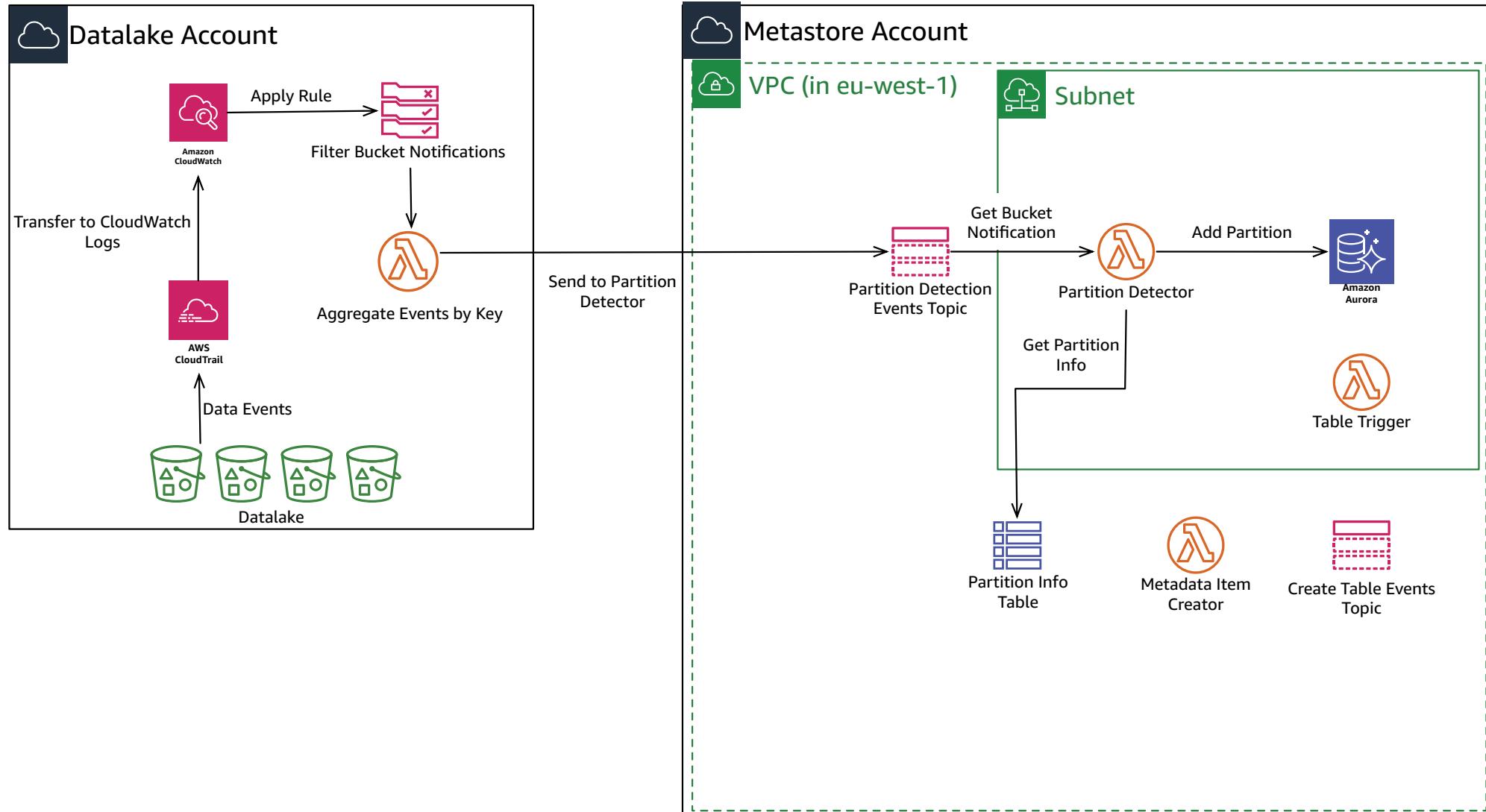
# Partition Detection Architecture



# Partition Detection Architecture



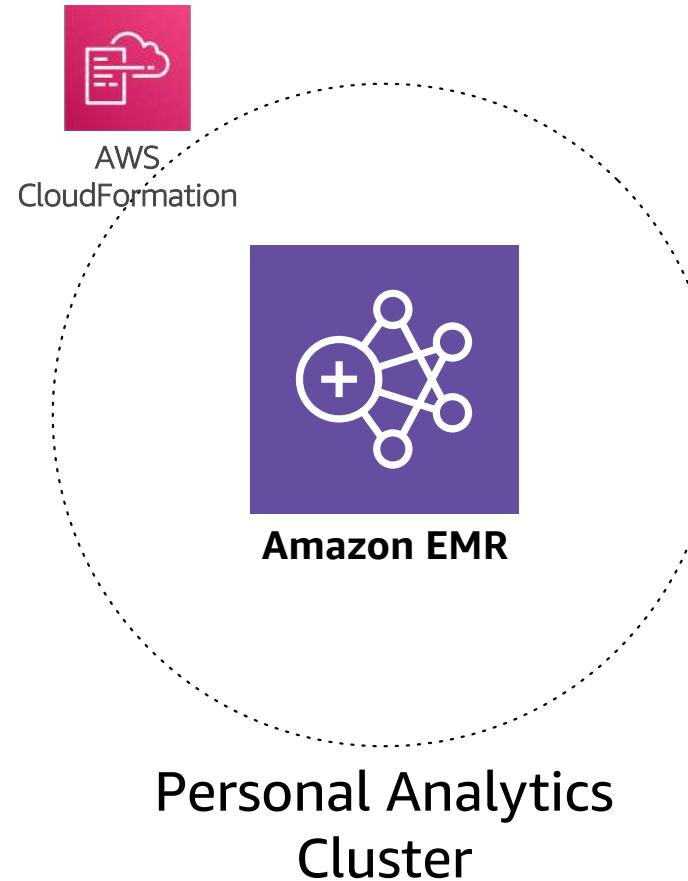
# Partition Detection Architecture



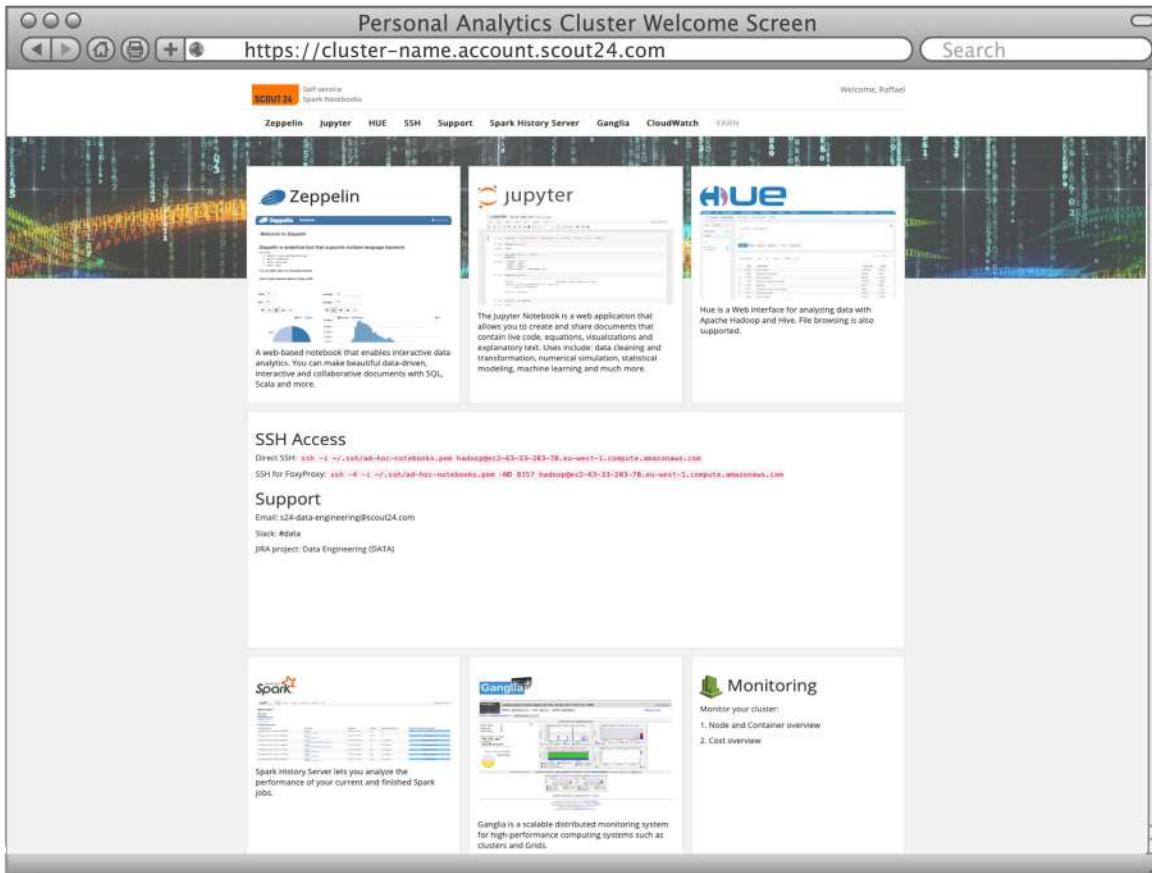
# Personal Analytics Cluster

---

# The Personal Analytics Cluster – An Overview



# The Personal Analytics Cluster – An Overview



Personal Analytics  
Cluster

Easy Access via Web  
Interface

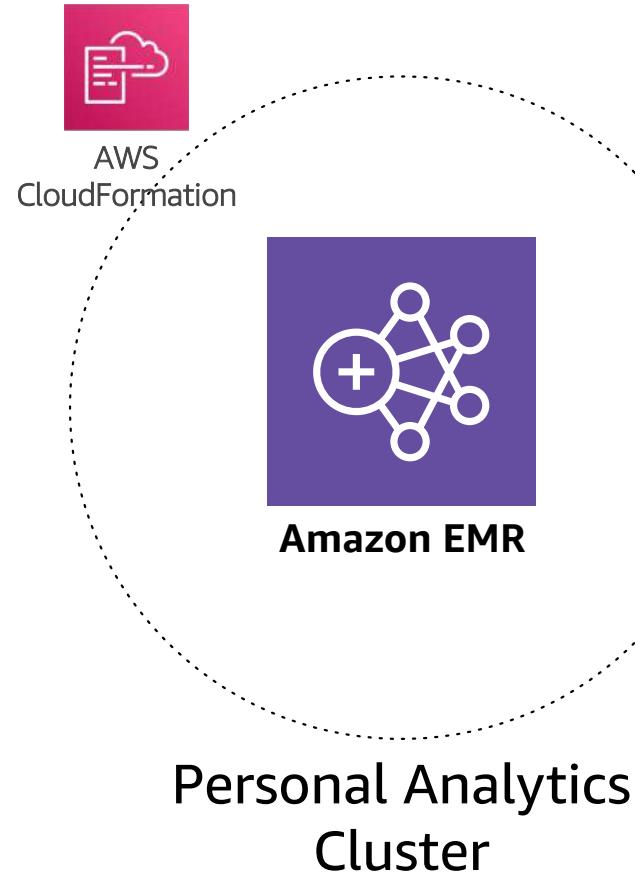
Zeppelin and Jupyter  
Notebook Restore

OneClick Deployment

Managed Scaling and  
Shutdown

Support for Pre-baked  
AMIs and Configs

# The Personal Analytics Cluster – An Overview



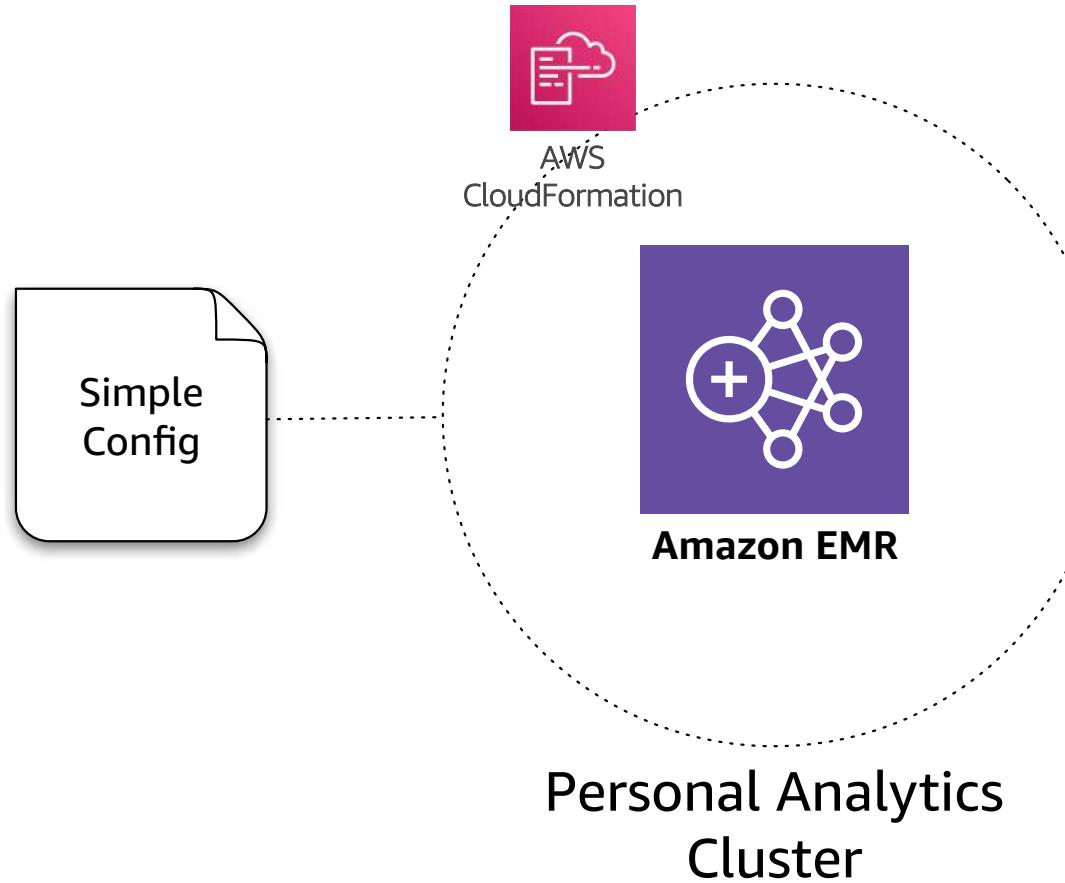
Provisioned Infrastructure

Daily: 20

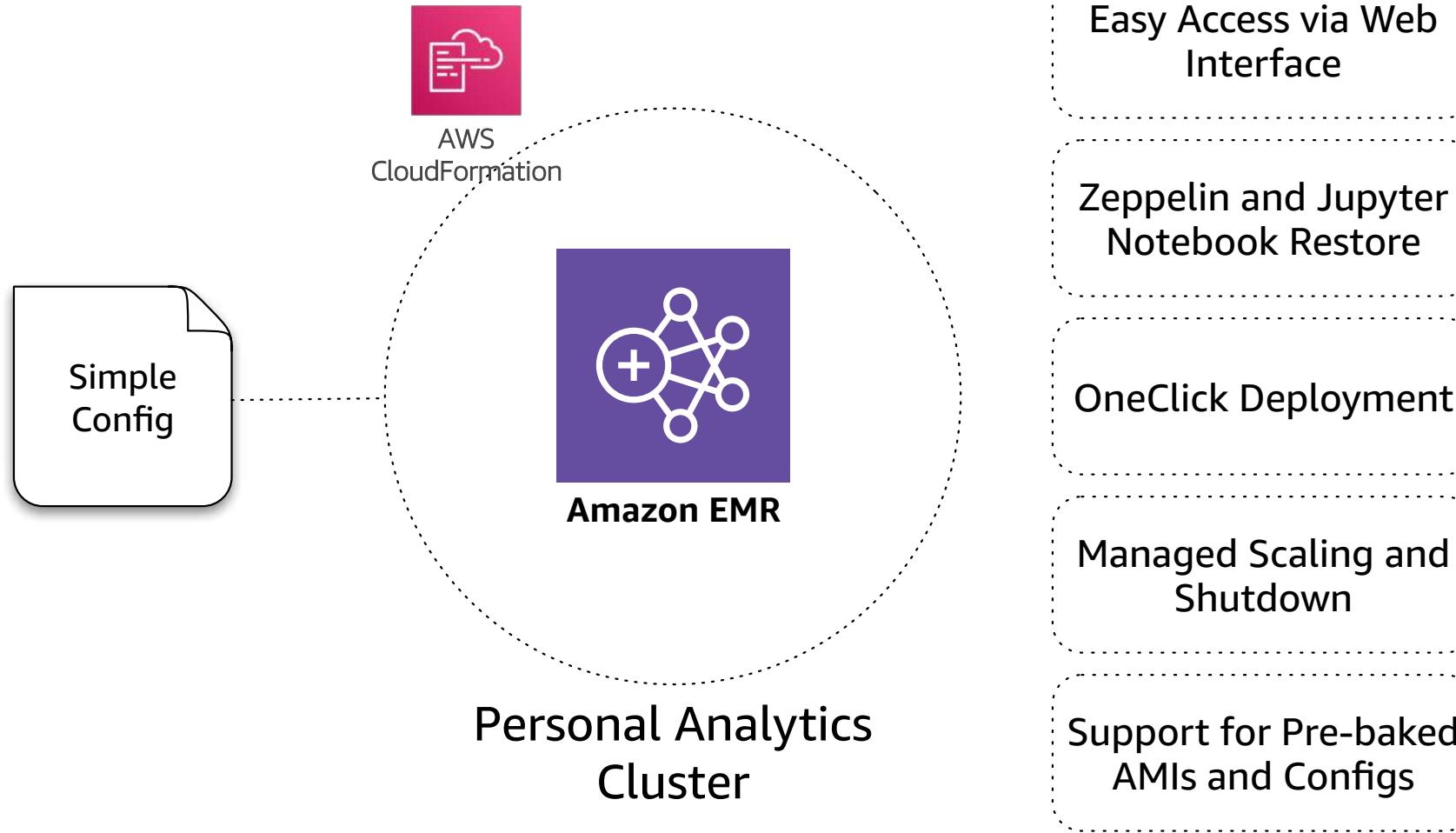
Weekly: 128

Monthly: 389

# The Personal Analytics Cluster – An Overview



# The Personal Analytics Cluster – An Overview



# The Personal Analytics Cluster – Configuration Setup

- Cluster Name (default is usually username)
- AWS Account Name
- Market segment
- Use Case
- Shutdown time
- Instance Type
- Bid Price
- Key Pair
- EBS Volume Size
- Tear Down Existing Cluster (optional)

# The Personal Analytics Cluster – Configuration Setup

```
$ ./deploy-cluster --cluster-name <> \
    --account-name <> \
    --market-segment <> \
    --usecase <> \
    --shutdown-time <> \
    --instance-type <> \
    --bid-price <> \
    --key-pair <> \
    --ebs-volume-size <> \
    --teardown-existing <>
```

Run Custom Build AWS Data / One Click to Data Analysis (Spark) Notebooks / - ☆☆ Start your own EM... ×

General Dependencies Changes Parameters \* Comment and Tags

The below parameters are marked as necessary for review

**Configuration parameters**

Cluster Name  Reset  
(default=your login) \*  
Browser URL will be [clustername]-cluster.  
[accountname].wolke.is

AWS Account name \*  Reset  
The cluster will be launched in this account. Note: teamcity  
must be trusted by this account.

Market/Platform Segment \*  Reset  
Please provide your market or platform segment for cost  
tagging

Use Case \*  Reset  
e.g. Jira Story or Epic, for specific cost tagging and  
monitoring

Shutdown time in Berlin time zone \*  Reset  
At this time your cluster will be automatically shut down  
(full hour only)

EC2 Instance Type \*  Reset  
Your cluster starts with 3 of those and scales to your  
needs (max=20)

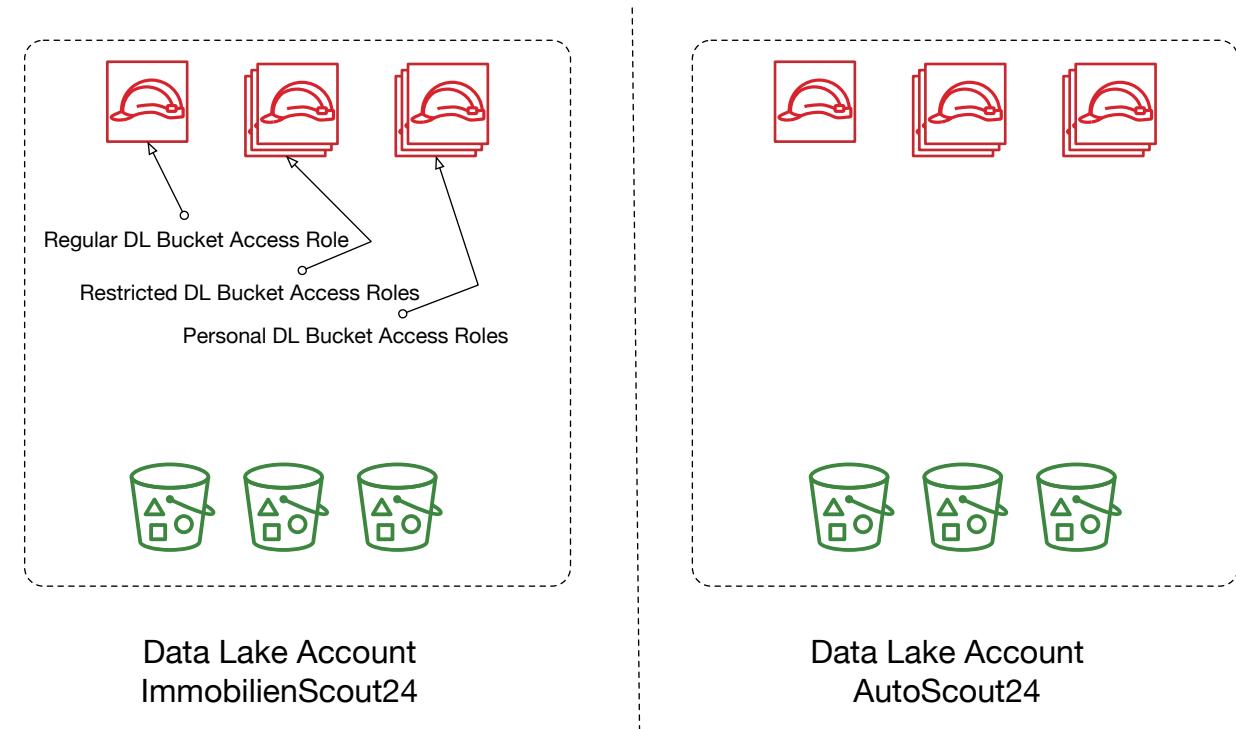
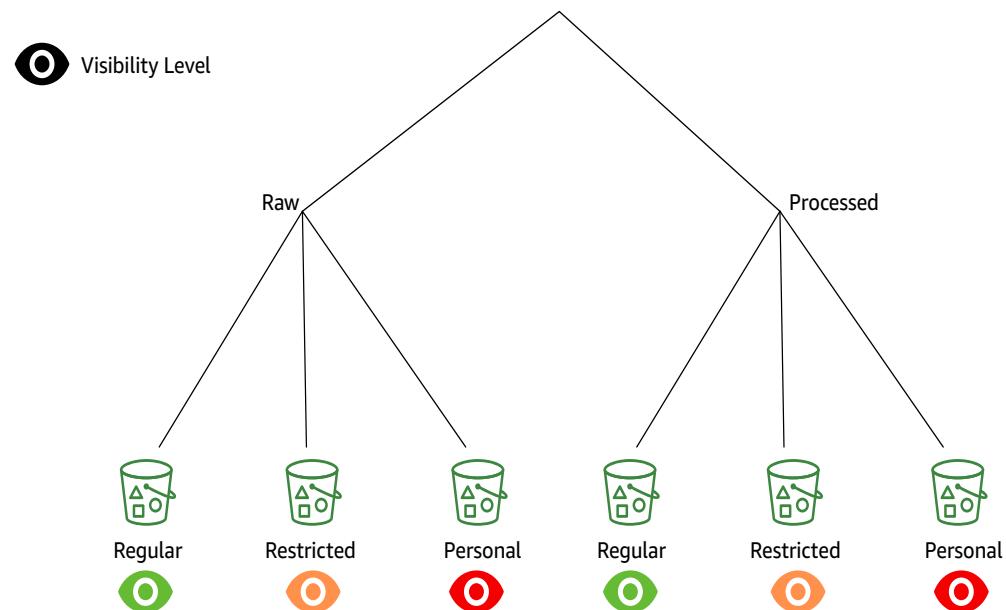
Bid Price \*  Reset  
Max \$/hour for EC2 instances, use default or ask  
DataEngineering

EC2 Key Pair name \*  Reset  
For ssh refer to a valid Key Pair in your account, otherwise  
leave blank

EBS Volume size \*  Reset  
The size of the EBS volumes that will be attached to the  
EMR instances

tearDown Existing PAC \*  Reset

# The Personal Analytics Cluster – Data Lake Access



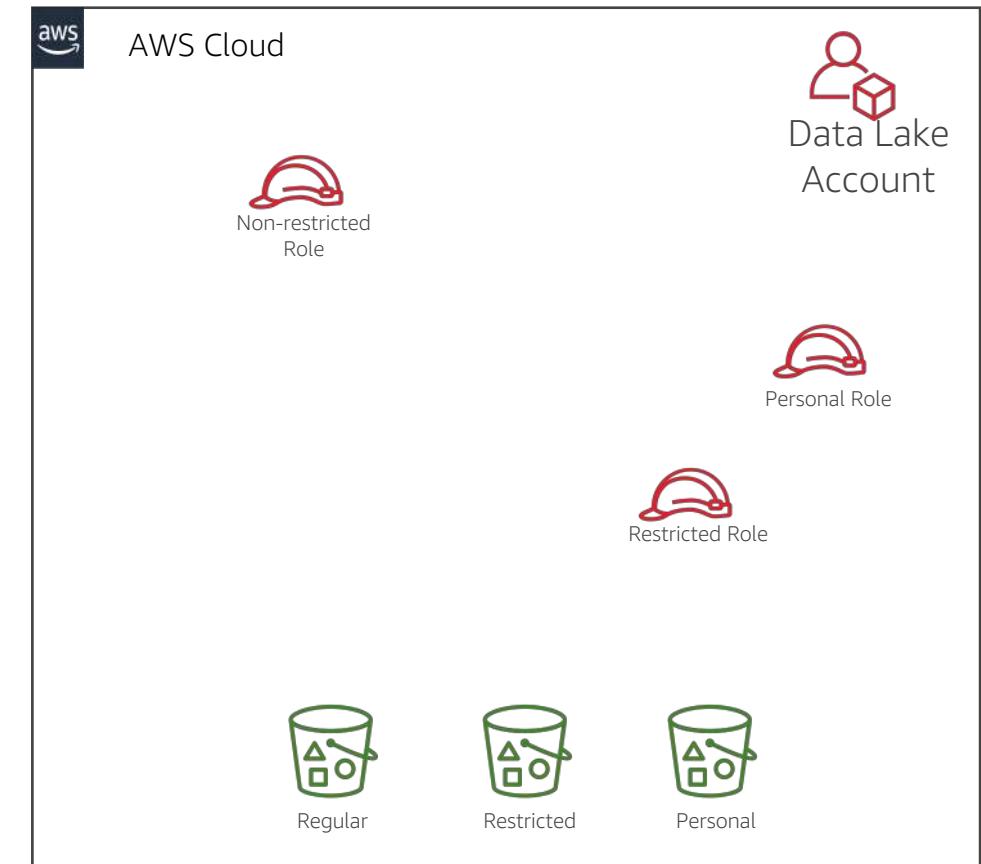
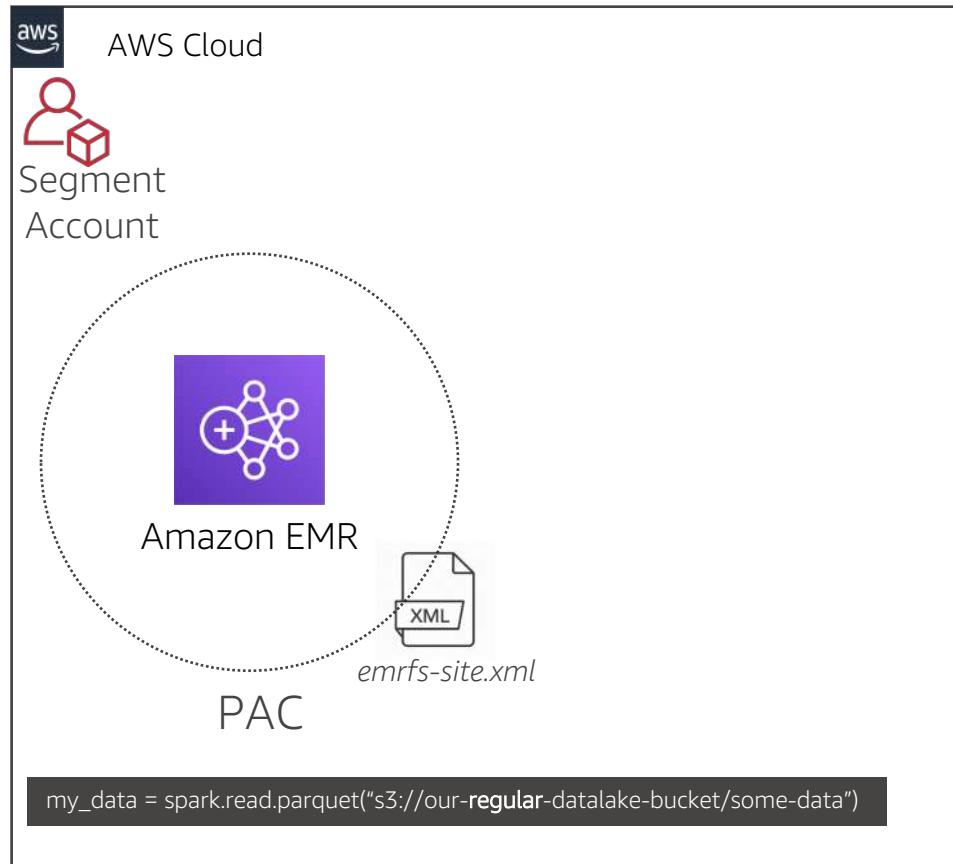
# The Personal Analytics Cluster – Data Lake Access

- Managed with AWS EMR FS Security Configuration
- Integrates with EMR Applications internally; Spark, Presto, etc
- Easy to define
- Maps a defined IAM role to S3 bucket being accessed

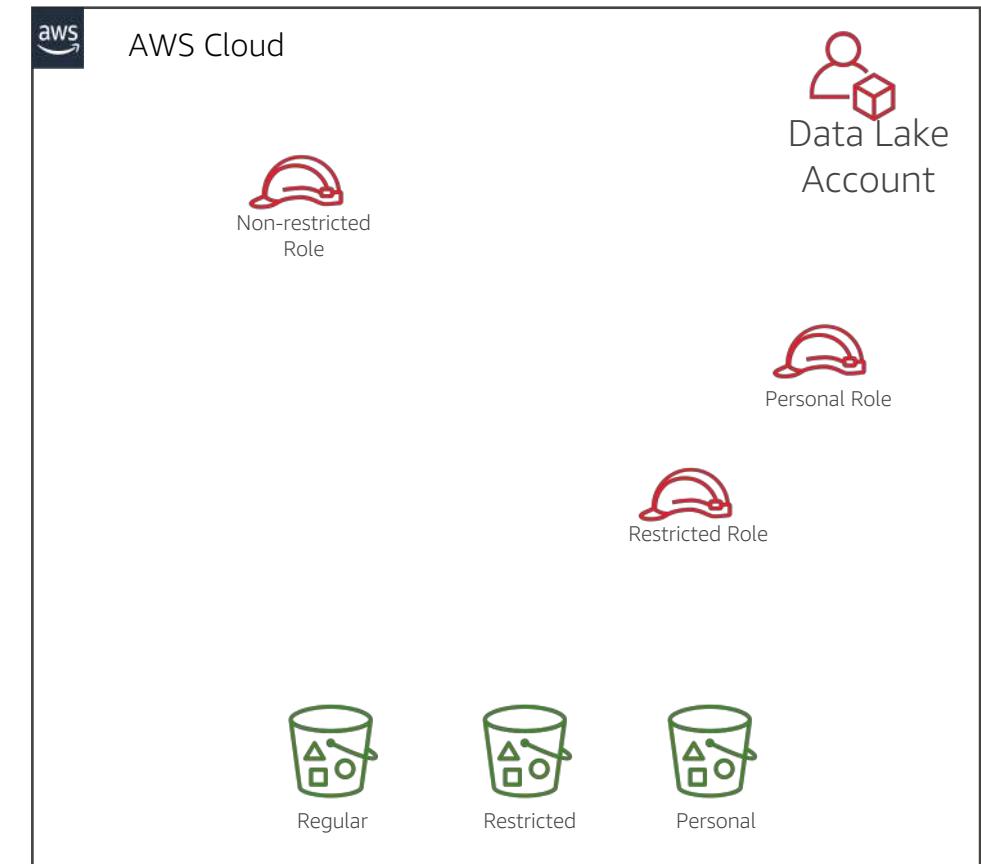
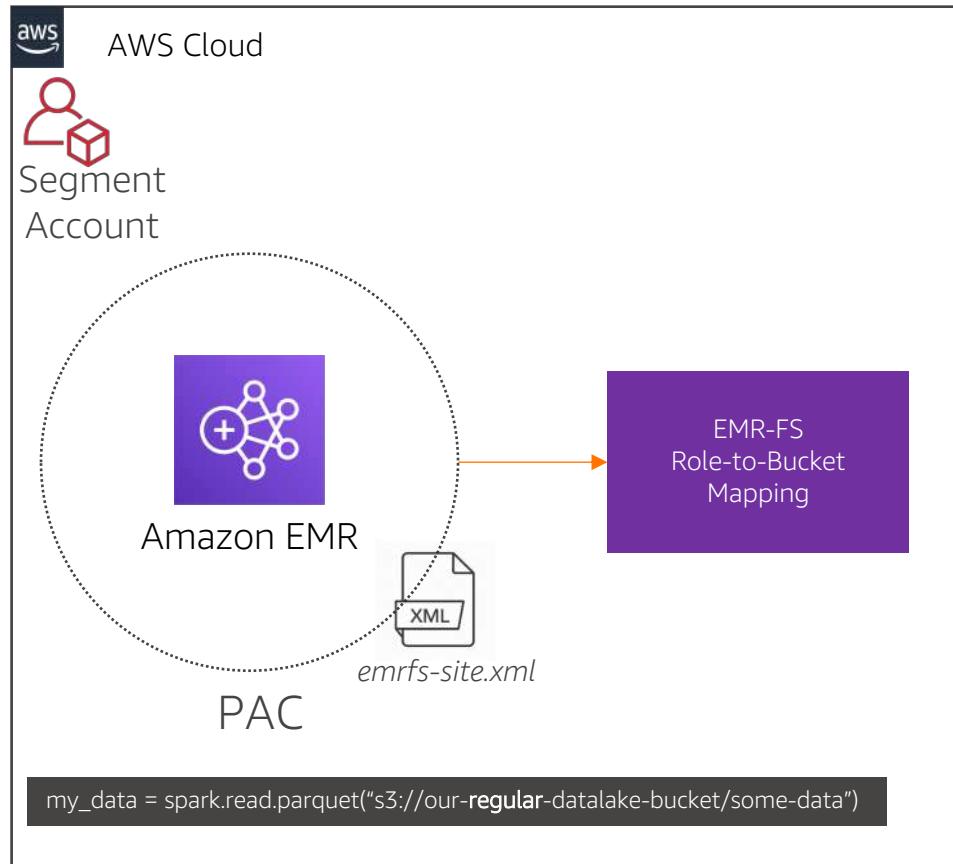
# The Personal Analytics Cluster – Data Lake Access

Scenario One:  
Non-restricted Data Lake Access

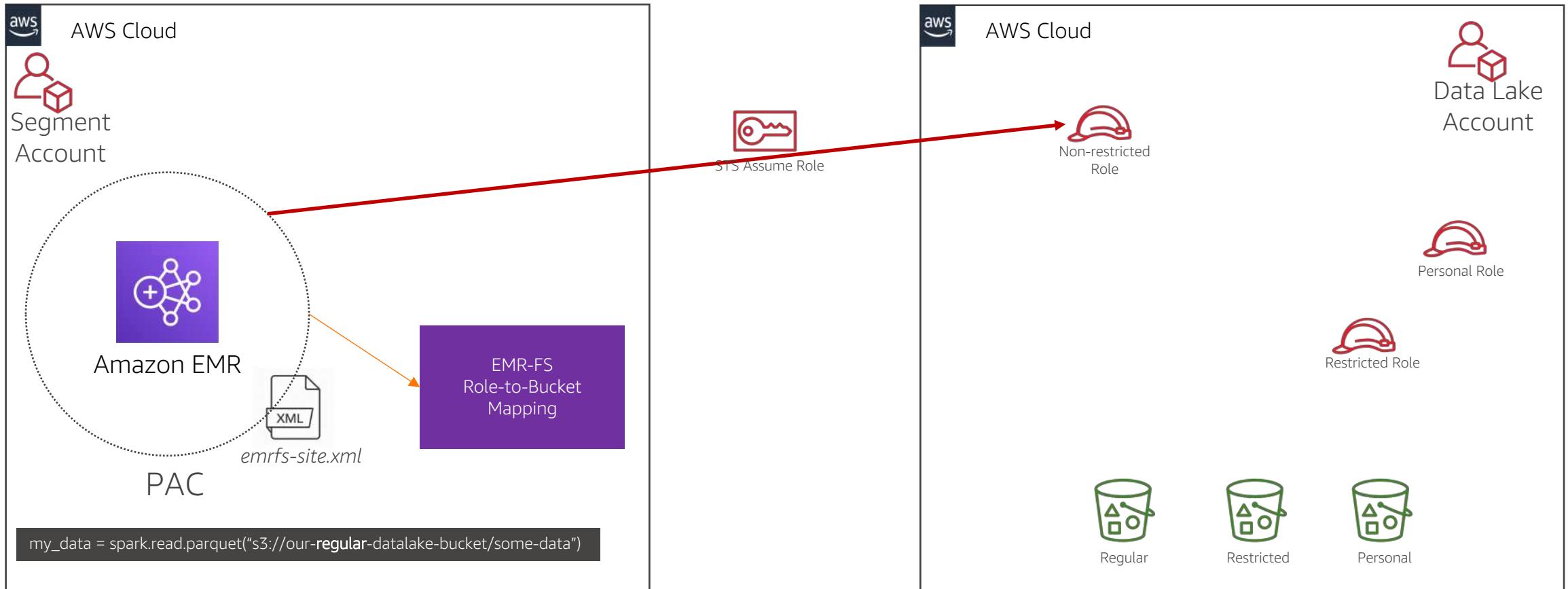
# The Personal Analytics Cluster – Data Lake Access



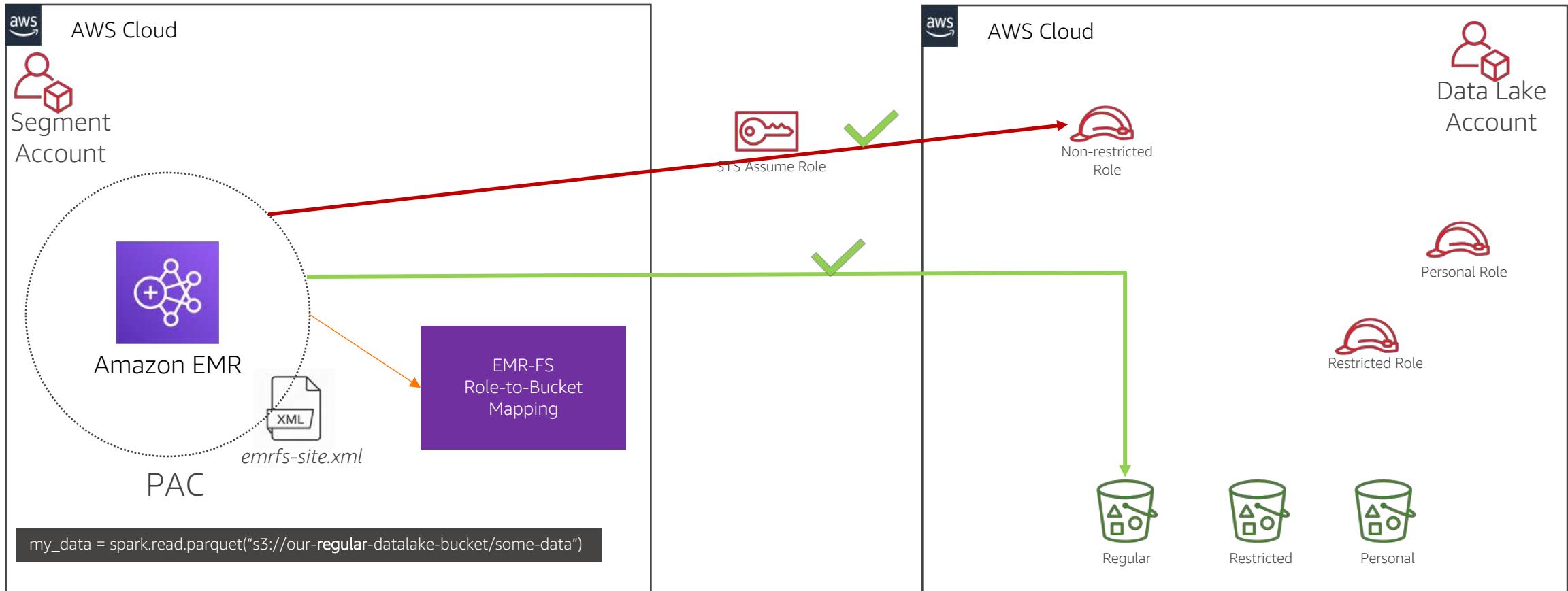
# The Personal Analytics Cluster – Data Lake Access



# The Personal Analytics Cluster – Data Lake Access



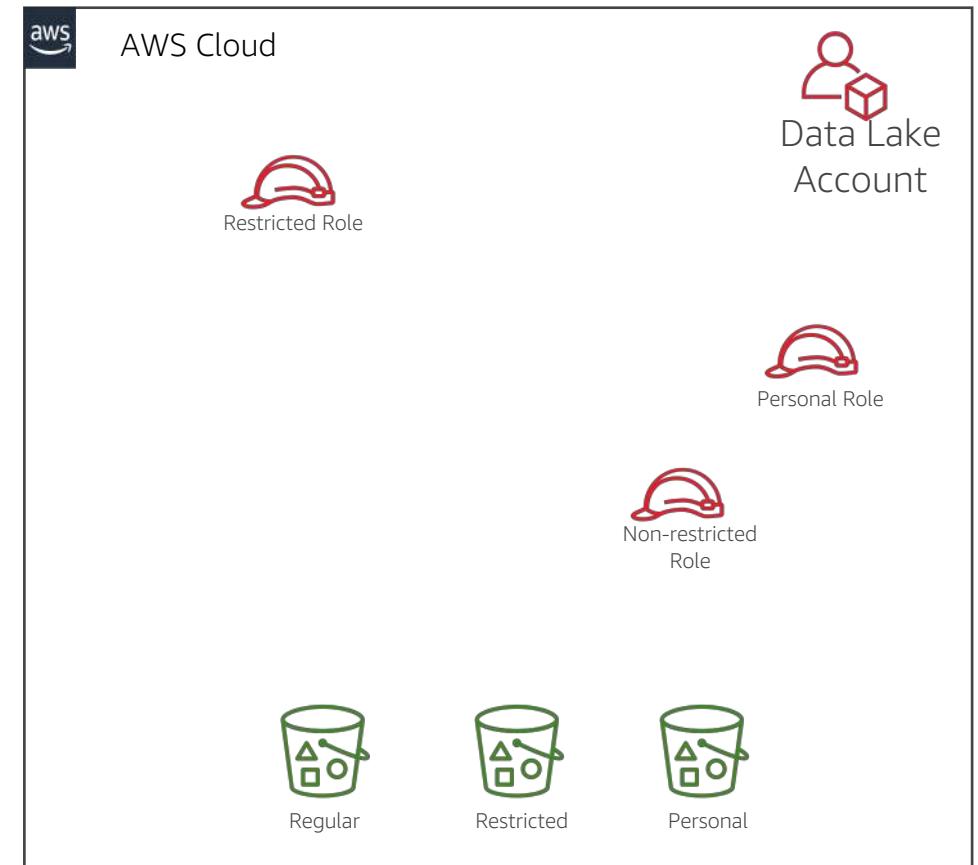
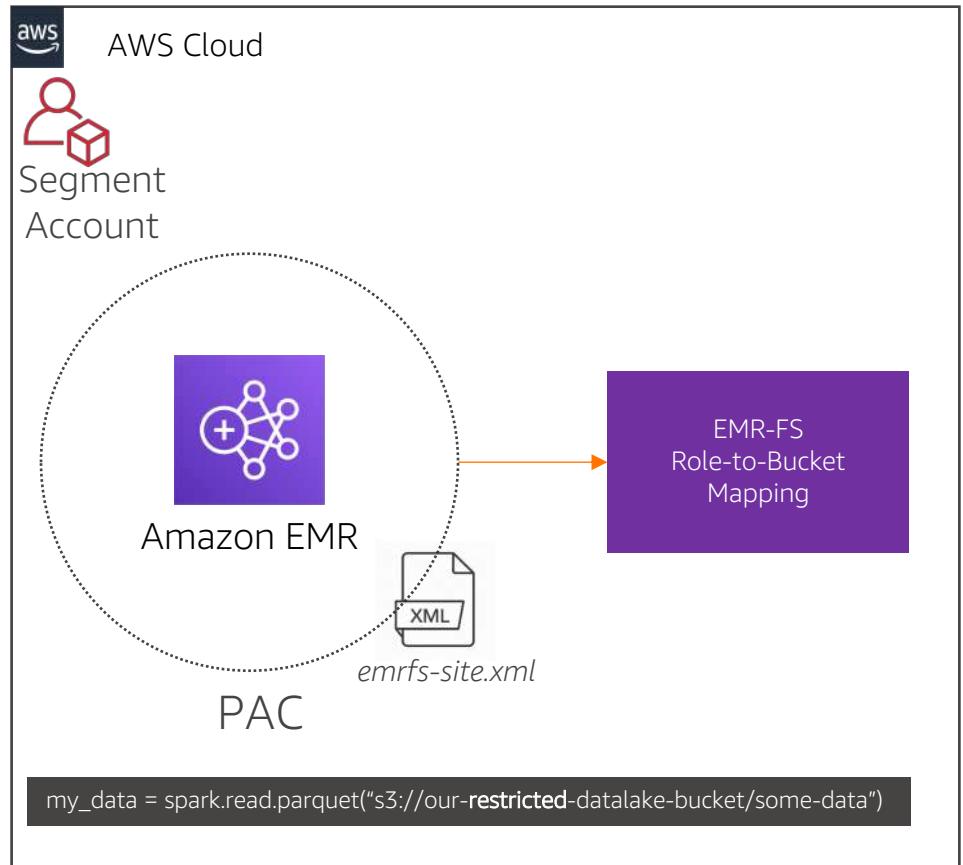
# The Personal Analytics Cluster – Data Lake Access



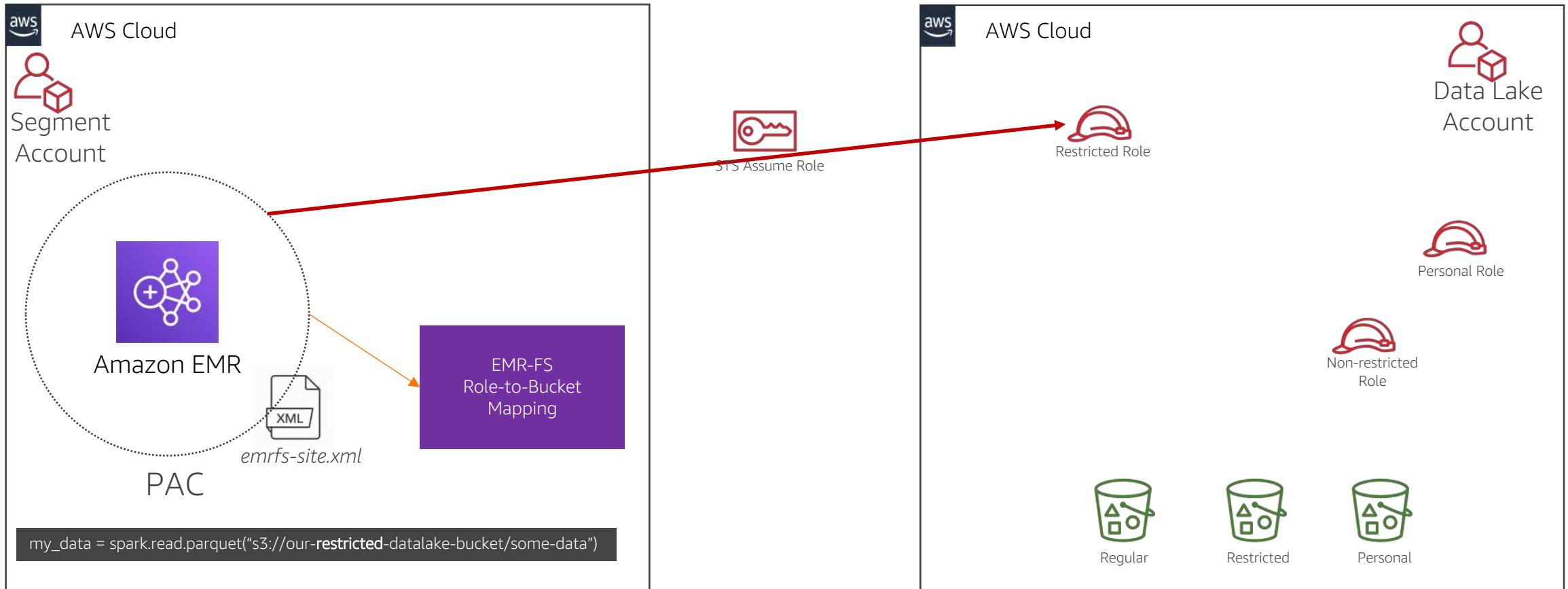
# The Personal Analytics Cluster – Data Lake Access

Scenario Two:  
Restricted/Personal Data Lake Access

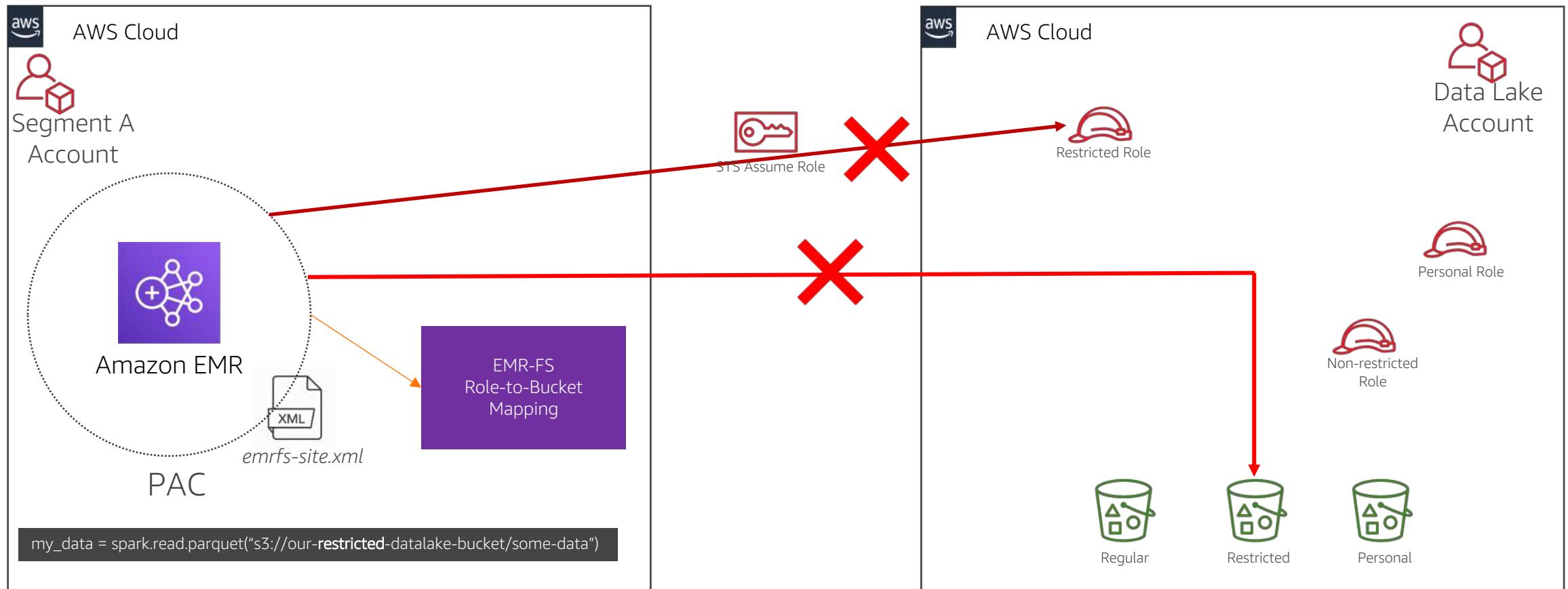
# The Personal Analytics Cluster – Data Lake Access



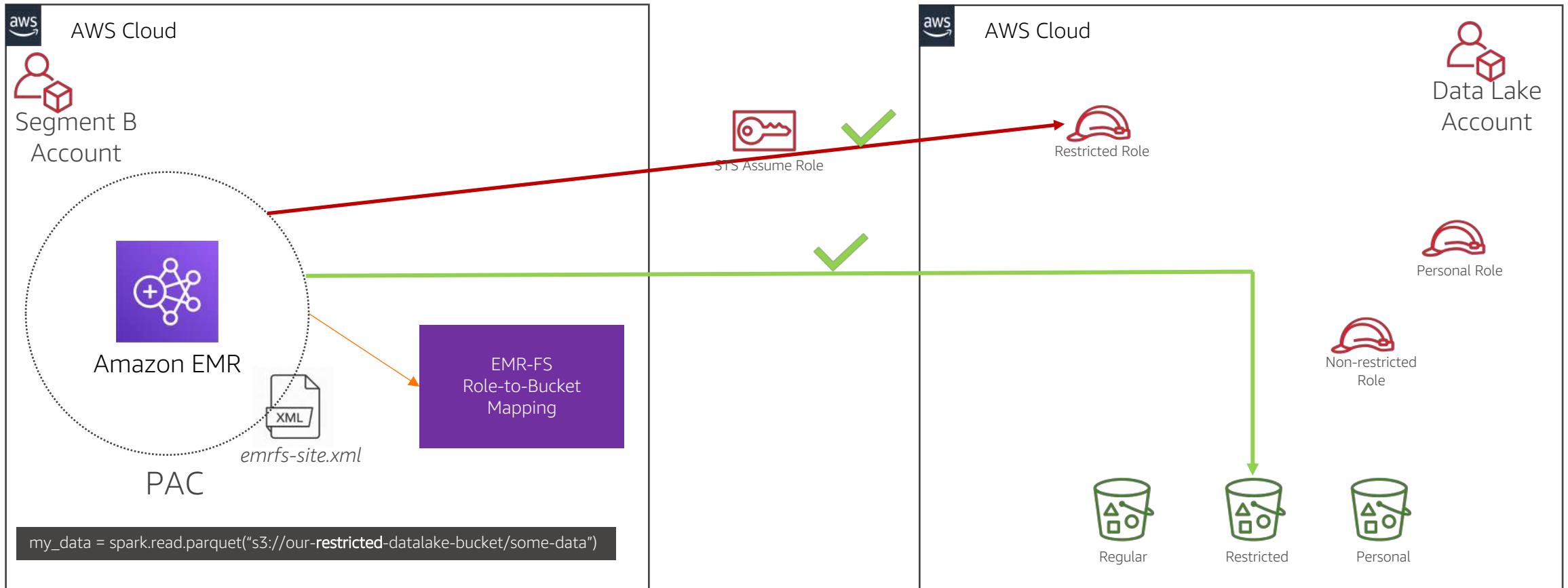
# The Personal Analytics Cluster – Data Lake Access



# The Personal Analytics Cluster – Data Lake Access



# The Personal Analytics Cluster – Data Lake Access



# Data Platform SQL Access

---

# Our Journey to Presto



# Our Journey to Presto



On EMR

Cross Account Support  
(OneScout Hive  
Metastore)

Leverages Datalake  
Access Roles (EMRFS)

Scheduled Scaling  
Configurations

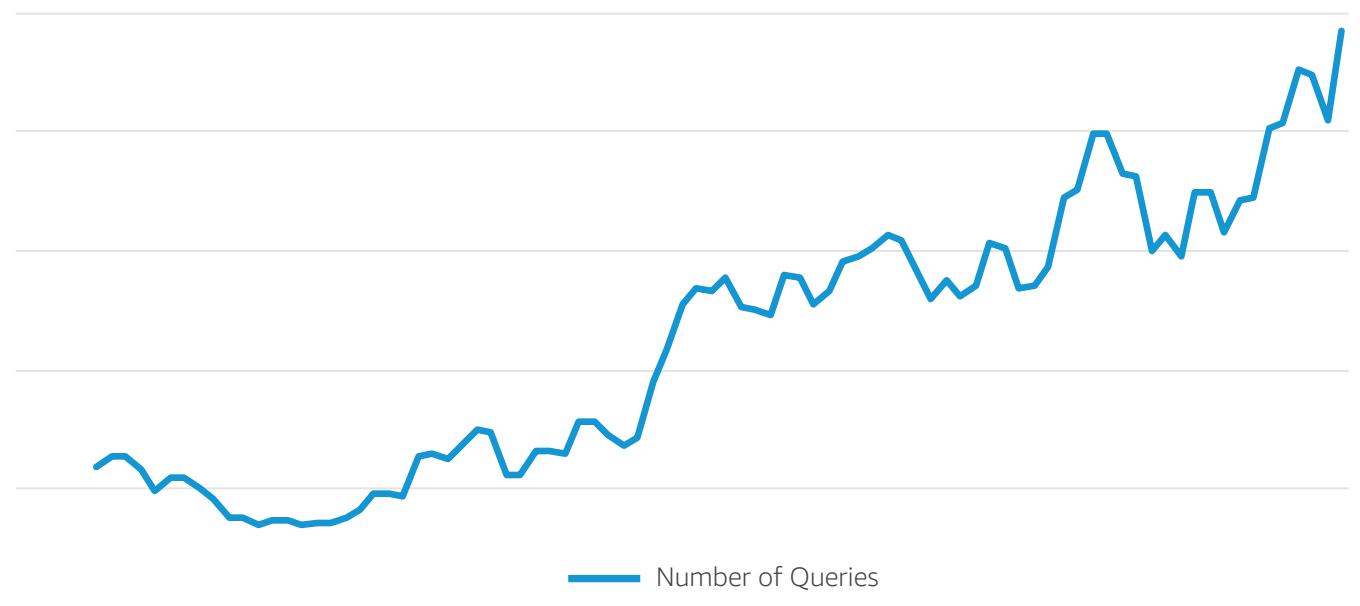
Fits our GDPR Concept  
(multiple isolated  
Clusters)

# Our Journey to Presto

360

Queries per Day

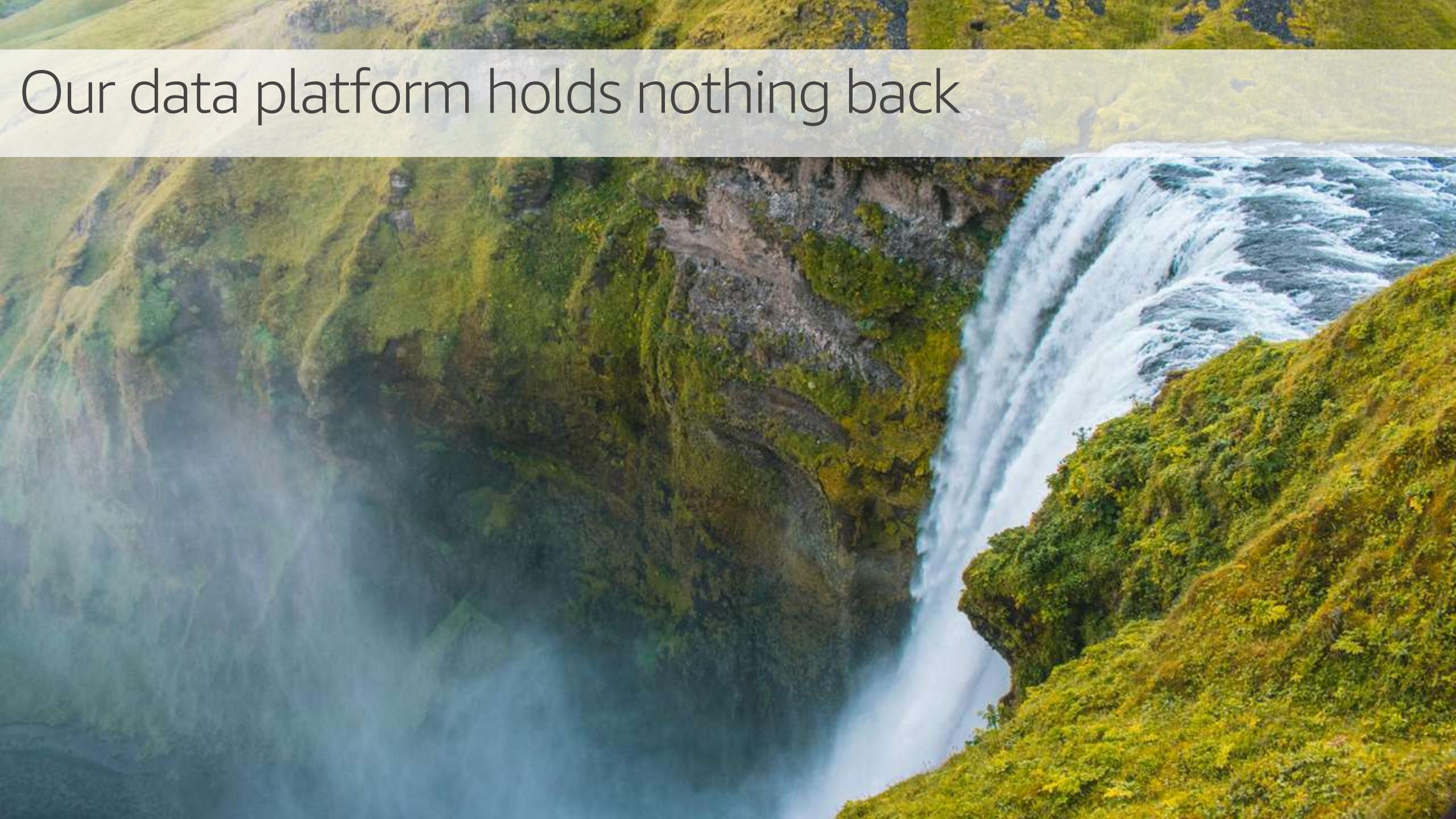
Average Presto Queries per Day



We hope to throw out  
**most** of the custom components  
we build.



Our data warehouse was a bottle neck

A wide-angle photograph of a waterfall. The water falls from a high, steep cliff face covered in lush green moss and vegetation. The waterfall creates a misty spray at the base. The water is a vibrant blue-green color. The background shows more of the same cliff face and vegetation.

Our data platform holds nothing back



# Thank you for your attention!

Christian Dietze

Senior Data Engineer

Scout24 Group

Olalekan Elesin

Data Engineer

Scout24 Group

Raffael Dzikowski

Senior Data Engineer

Scout24 Group

