

# Introduction to AWS Glue

# AWS Glue automates the undifferentiated heavy lifting of ETL

## Discover

Automatically discover and categorize your data making it immediately searchable and queryable across data sources

## Develop

Generate code to clean, enrich, and reliably move data between various data sources; you can also use their favorite tools to build ETL jobs

## Deploy

Run your jobs on a serverless, fully managed, scale-out environment. No compute resources to provision or manage.



Pop-up Loft  
LONDON

# AWS Glue: Components

---



## Data Catalog

- Hive Metastore compatible with enhanced functionality
- Crawlers automatically extracts metadata and creates tables
- Integrated with Amazon Athena, Amazon Redshift Spectrum



## Job Authoring

- Auto-generates ETL code
- Build on open frameworks – Python and Spark
- Developer-centric – editing, debugging, sharing



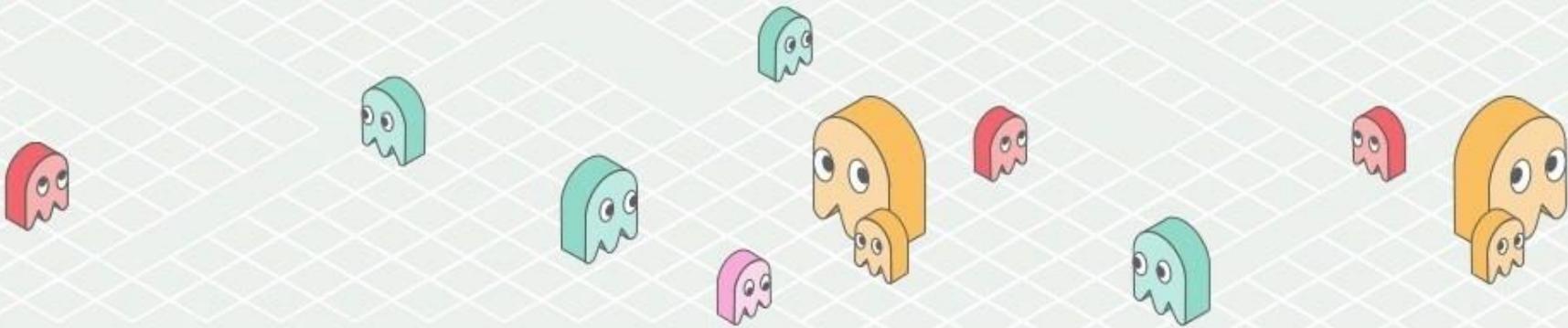
## Job Execution

- Run jobs on a serverless Spark platform
- Provides flexible scheduling
- Handles dependency resolution, monitoring and alerting

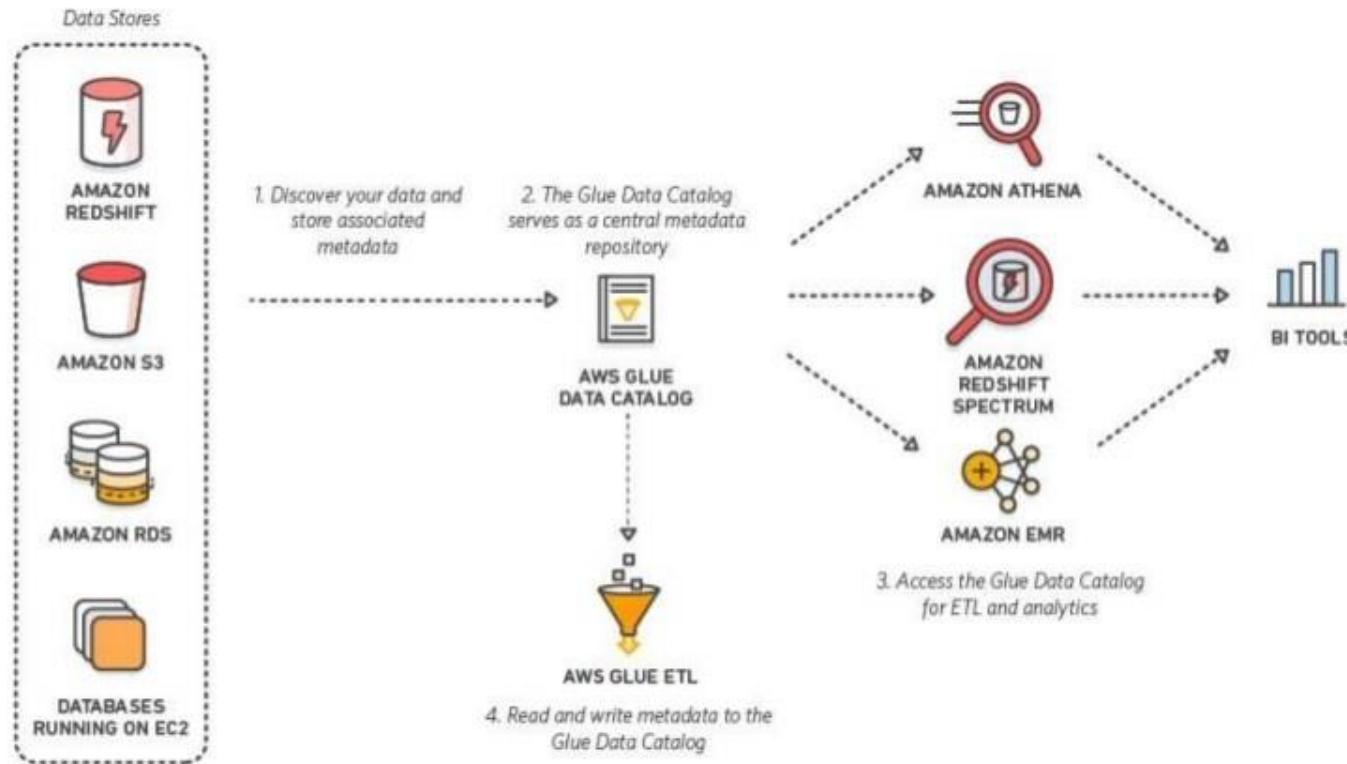


Pop-up Loft  
LONDON

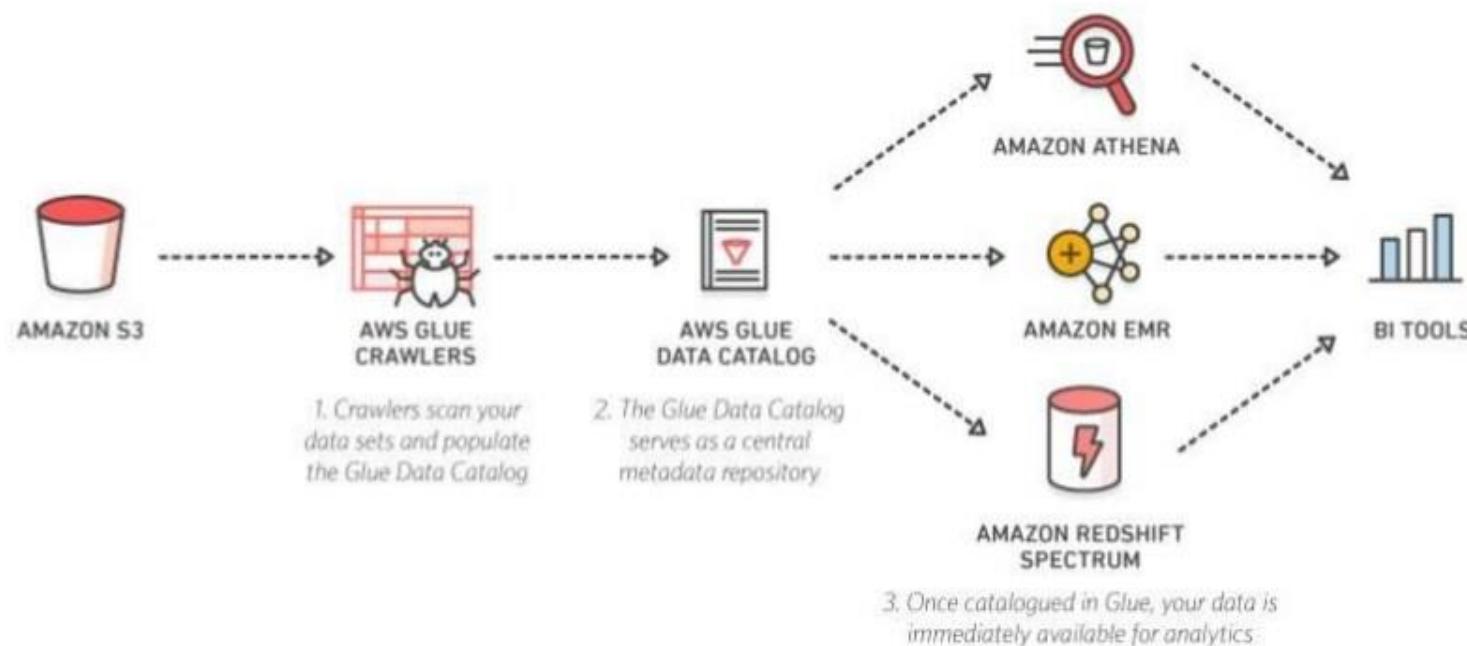
# Common Use Cases for AWS Glue



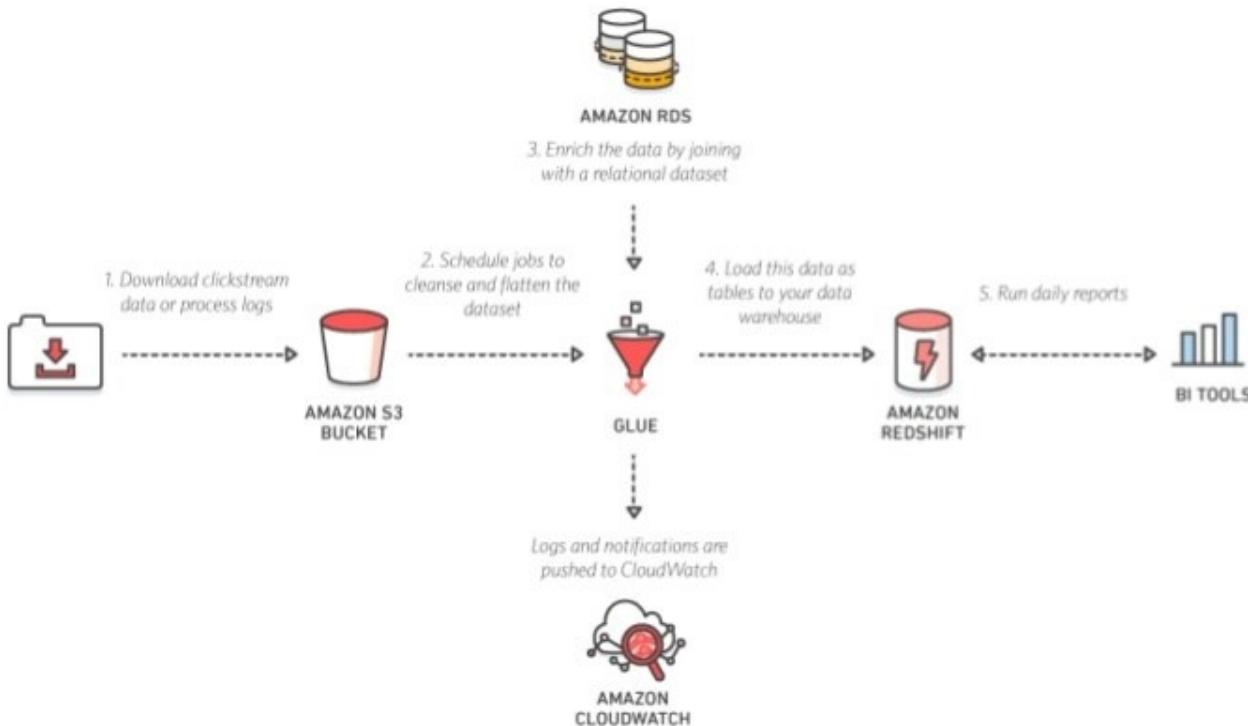
# Understand your data assets



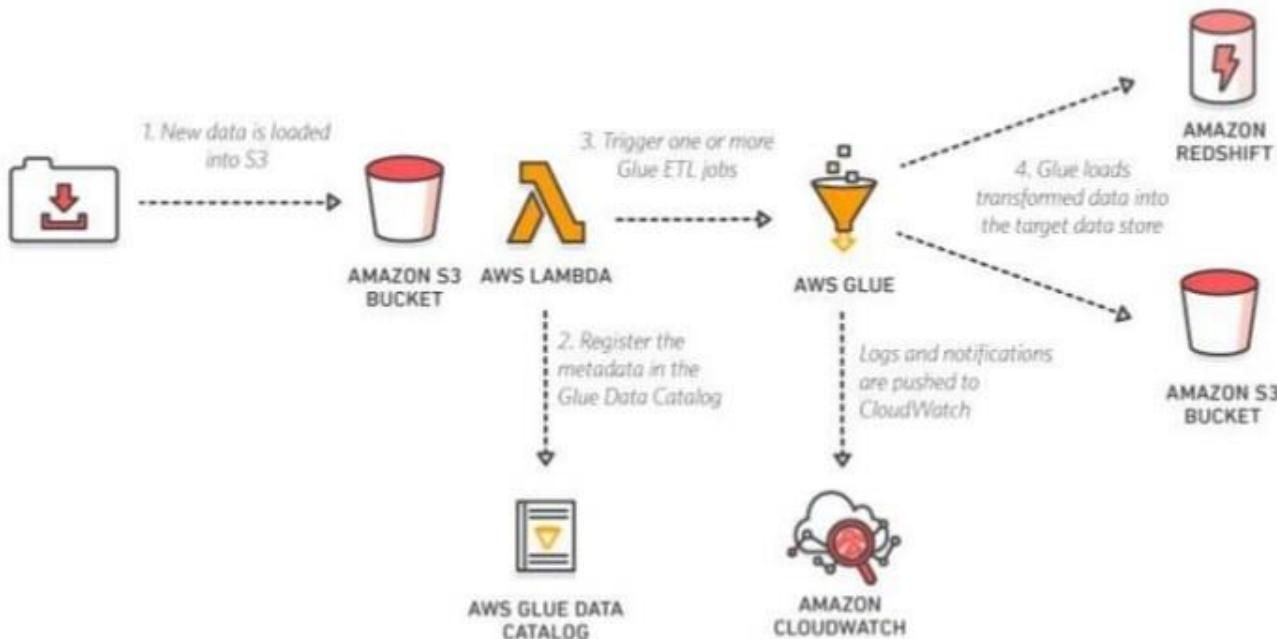
# Instantly query your data lake on Amazon S3



# ETL data into your data warehouse



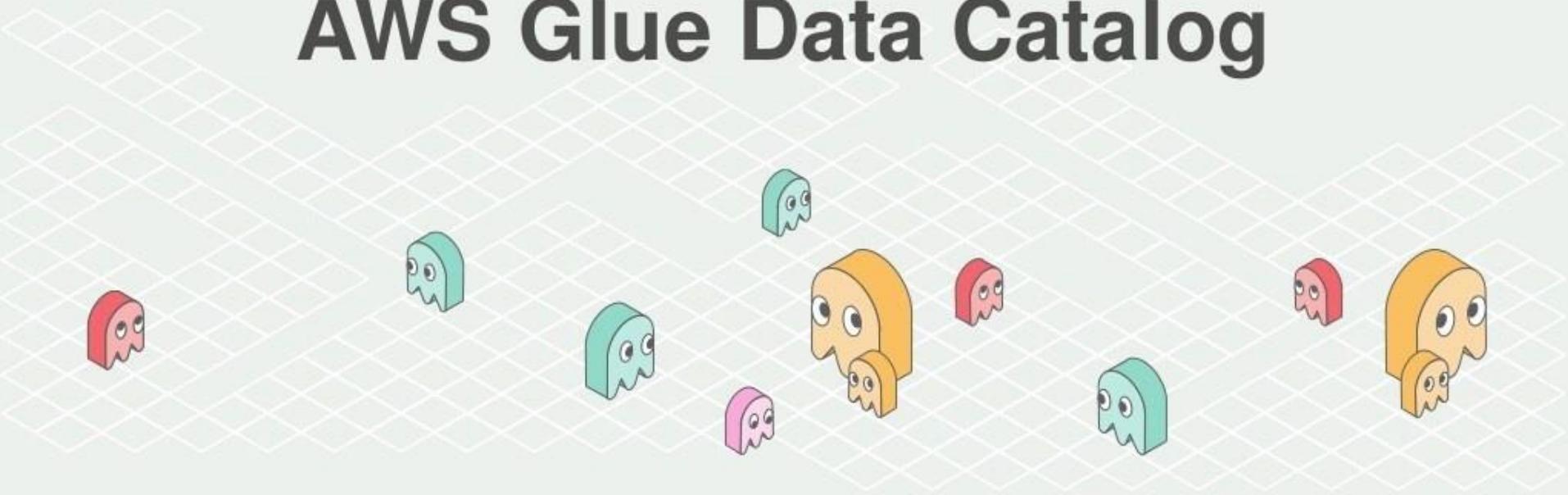
# Build event-driven ETL pipelines





Pop-up Loft  
LONDON

# AWS Glue Data Catalog



# Glue data catalog

---

Manage table metadata through a Hive metastore API or Hive SQL.  
Supported by tools like Hive, Presto, Spark etc.

We added a few extensions:

- **Search** over metadata for data discovery
- **Connection info** – JDBC URLs, credentials
- **Classification** for identifying and parsing files
- **Versioning** of table metadata as schemas evolve and other metadata are updated

Populate using Hive DDL, bulk import, or automatically through **Crawlers**.

## Data Catalog: Crawlers

---

Crawlers automatically build your Data Catalog and keep it in sync

- ▶ Automatically discover new data, extracts schema definitions
  - Detect schema changes and version tables
  - Detect Hive style partitions on Amazon S3
- ▶ Built-in classifiers for popular types; custom classifiers using Grok expressions
- ▶ Run ad hoc or on a schedule; serverless – only pay when crawler runs

# AWS Glue Data Catalog

Bring in metadata from a variety of data sources (Amazon S3, Amazon Redshift, etc.) into a single categorized list that is searchable

The screenshot shows the AWS Glue Data Catalog interface. On the left, there's a sidebar with navigation links for Services, Resource Groups, AWS Glue, Data catalog, Databases, Tables (which is selected), Connections, Crawlers, Classifiers, ETL, Jobs, Triggers, Dev endpoints, Tutorials, Add crawler, Explore table, and Add job. The main area is titled "Tables" with a sub-instruction: "A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition." Below this is a search bar with filters for "Add tables", "Actions", and "Filter by attribute or search by keyword". A table lists 13 entries:

Name	Database	Location	Classification	Last updated	Deprecated
cloudfront	logs	s3://[REDACTED]-us-east-1/cloudfront/	csv	10 August 2017 10:04 AM UTC...	
customer1_public_customer	redshiftnamespace	customer1.public.customer	redshift	10 August 2017 9:49 AM UTC...	
customer1_public_inorder	redshiftnamespace	customer1.public.inorder	redshift	10 August 2017 9:49 AM UTC...	
customer1_public_part	redshiftnamespace	customer1.public.part	redshift	10 August 2017 9:49 AM UTC...	
customer1_public_supplier	redshiftnamespace	customer1.public.supplier	redshift	10 August 2017 9:49 AM UTC...	
data_store_api_public_agency	spending	data_store_api.public.agency	postgresql	10 August 2017 11:51 AM UTC...	
data_store_api_public_awards	spending	data_store_api.public.awards	postgresql	10 August 2017 11:51 AM UTC...	
data_store_api_public_budge...	spending	data_store_api.public.budget_authority	postgresql	10 August 2017 11:51 AM UTC...	
data_store_api_public_feder...	spending	data_store_api.public.federal_account	postgresql	10 August 2017 11:51 AM UTC...	
imo_2010f1e	nytaxanalysis	s3://[REDACTED]-NYC-transportation/imo/	csv	25 July 2017 8:03 PM UTC-4	
ny_pub	nytaxanalysis	s3://[REDACTED]/canonical/NY-Pub/	parquet	25 July 2017 9:49 AM UTC-4	
taxi_303e40bd	nytaxanalysis	s3://[REDACTED]-NYC-transportation/taxi/	csv	25 July 2017 8:03 PM UTC-4	
tweets_2f	logs	s3://[REDACTED]/twitter/2017/04/27/	json	10 August 2017 10:33 AM UTC...	

# Data Catalog: Table details

Table properties

Data statistics

Table schema

Table properties

Data catalog

Databases

Tables

Connections

Crawlers

Classifiers

ETL

Jobs

Triggers

Dev endpoints

Tutorials

Add crawler

Explore table

Add job

AWS Glue

simpletweets\_json

Name: simpletweets\_json  
Description: analytics  
Database: analytics  
Classification: json  
Location: s3://gluesampledata/simpletweets.json  
Connection: [aws-glue](#)  
Deprecated: No  
Last updated: Thu Aug 10 16:25:24 GMT-700 2017

Properties:

sizeKey	456580	objectCount	1	UPDATED_BY_CRAWLER	S3Crawler	CrawlerSchemaSerializationVersion	1.0
recordCount	1001	averageRecordSize	456	CrawlerSchemaDeserializationVersion	1.0	compressionType	none
				typeOfData	file		

Schema:

	Column name	Data type
1	entities	struct
2	id	bigint
3	retweeted	boolean
4	text	string
5	user	struct

Nested fields

user schema details

```
STRUCT
  annotations_analytics BOOLEAN
  description STRING
  features_count INT
  followers_count INT
  friends_count INT
  id INT
  lang STRING
  location STRING
  name STRING
  profile_background_color BOOLEAN
```

# Data Catalog: Version control

## Compare schema versions

Last updated 10 Aug 2017 Table Version 0			
Name	simpletweets_json	Description	
Database	simpletweets_json	Database	simpletweets_json
Classification	json	Classification	json
Location	s3://guvsamplereads/simpletweets.json	Location	s3://guvsamplereads/simpletweets.json
Connection		Connection	
Deprecated	No	Deprecated	No
Last updated	Thu Aug 10 16:25:24 GMT-700 2017	Last updated	Thu Aug 10 16:25:24 GMT-700 2017
sizeKey	46680	objectCount	1
UPDATED_BY_CRAWLER	S3Crawler	Properties	CrawlerSchemaDeserializerVersion: 1.0 recordCount: 1001
Properties	CrawlerSchemaDeserializerVersion: 1.0 recordCount: 1001	Properties	Properties
averageRecordSize	466	CrawlerSchemaDeserializerVersion	1.0
compressionType	none	typeOfData	file
Change Column name Data type Key			
Changed	entities	struct	
	id	bigint	
	retweeted	boolean	
	text	string	
	user	struct	

Last updated 10 Aug 2017 Table Version 1 (Current version)			
Name	simpletweets_json	Description	
Database	simpletweets_json	Database	simpletweets_json
Classification	json	Classification	json
Location	s3://guvsamplereads/simpletweets.json	Location	s3://guvsamplereads/simpletweets.json
Connection		Connection	
Deprecated	No	Deprecated	No
Last updated	Thu Aug 10 16:25:24 GMT-700 2017	Last updated	Thu Aug 10 16:25:24 GMT-700 2017
sizeKey	46680	objectCount	1
UPDATED_BY_CRAWLER	S3Crawler	Properties	CrawlerSchemaDeserializerVersion: 1.0 recordCount: 1001
Properties	CrawlerSchemaDeserializerVersion: 1.0 recordCount: 1001	Properties	Properties
averageRecordSize	466	CrawlerSchemaDeserializerVersion	1.0
compressionType	none	typeOfData	file
Change Column name Data type Key			
Changed	entities	struct	
	id	STRING	
	retweeted	boolean	
	text	string	
	user	struct	

## List of table versions



# Data Catalog: Detecting partitions

S3 bucket hierarchy

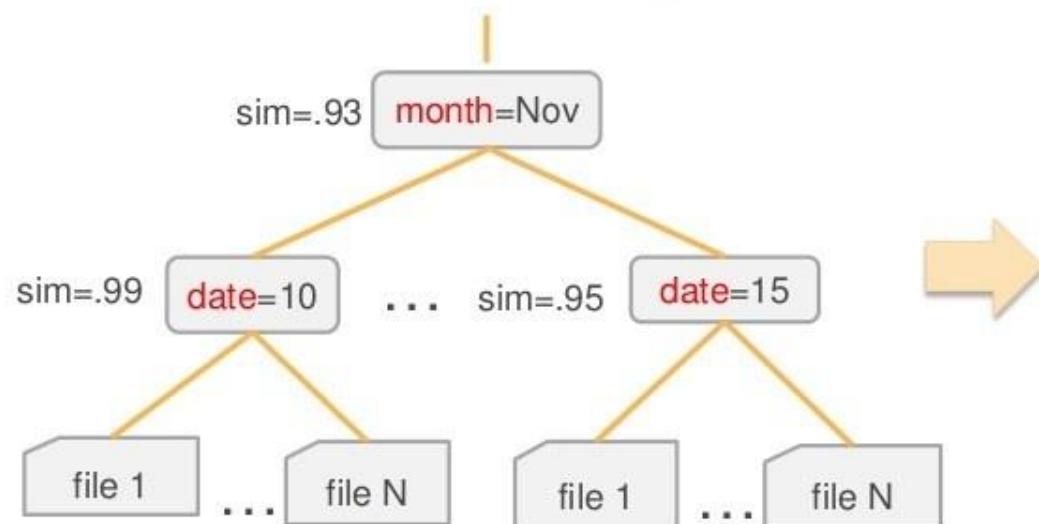


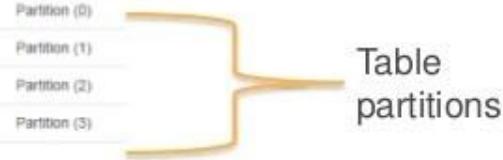
Table definition

Column	Type
month	str
date	str
col 1	int
col 2	float
:	:
:	:
:	:

Estimate schema similarity among files at each level to handle semi-structured logs, schema evolution...

# Data Catalog: Automatic partition detection

12	errorcode	string
13	region	string
14	year	string
15	month	string
16	day	string



Automatically register available partitions

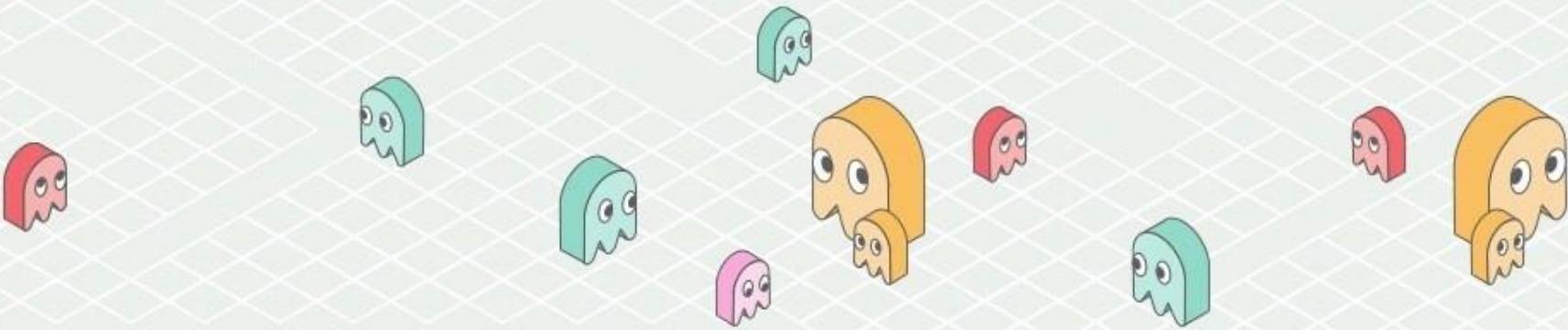
Showing: 1 - 3 < >

region	year	month	day	View files	View properties
us-west-2	2016	02	18	<a href="#">View files</a>	<a href="#">View properties</a>
us-west-2	2016	02	17	<a href="#">View files</a>	<a href="#">View properties</a>
us-west-2	2016	02	16	<a href="#">View files</a>	<a href="#">View properties</a>



Pop-up Loft  
LONDON

# Job Authoring in AWS Glue

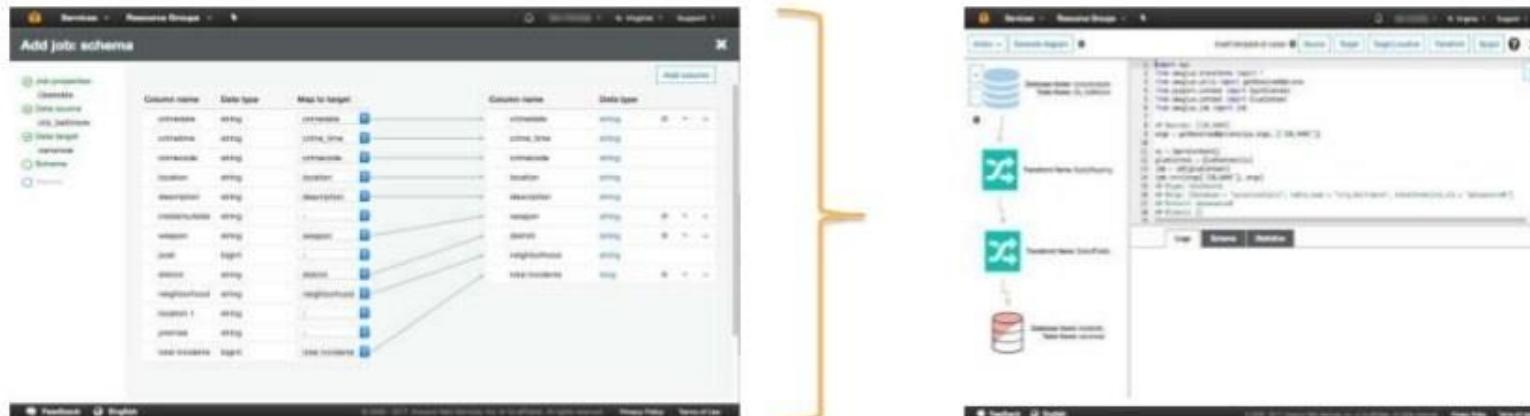


# Job Authoring Choices

- Python code generated by AWS Glue
- Connect a notebook or IDE to AWS Glue
- Existing code brought into AWS Glue



# Job authoring: Automatic code generation



1. Customize the mappings
2. Glue generates transformation graph and **Python** code
3. Connect your **notebook** to development endpoints to customize your code



Pop-up Loft  
LONDON

# Job authoring: ETL code

- **Human-readable**, editable, and portable PySpark code

```
28 sc = SparkContext()
29 glueContext = GlueContext(sc)
30 job = Job(glueContext)
31 job.init(args['JOB_NAME'], args)
32 ## @type: DataSource
33 ## @args: [name_space = "nytaxianalysis", table_name = "taxi_383e48bd", transformation_ctx = "datasource@"]
34 ## @return: datasource0
35 ## @inputs: []
36 datasource0 = glueContext.create_dynamic_frame.from_catalog(name_space = namespace, table_name = tablename, transformation_ctx = "datasource@")
37 RenameField0 = RenameField.apply(frame = datasource0, old_name="lpep_pickup_datetime", new_name="pickup_datetime", transformation_ctx = "RenameField0")
38 RenameField1 = RenameField.apply(frame = RenameField0, old_name="lpep_dropoff_datetime", new_name="dropoff_datetime", transformation_ctx = "RenameField1")
39 RenameField2 = RenameField.apply(frame = RenameField1, old_name="ratecodeid", new_name="ratecode", transformation_ctx = "RenameField2")
```

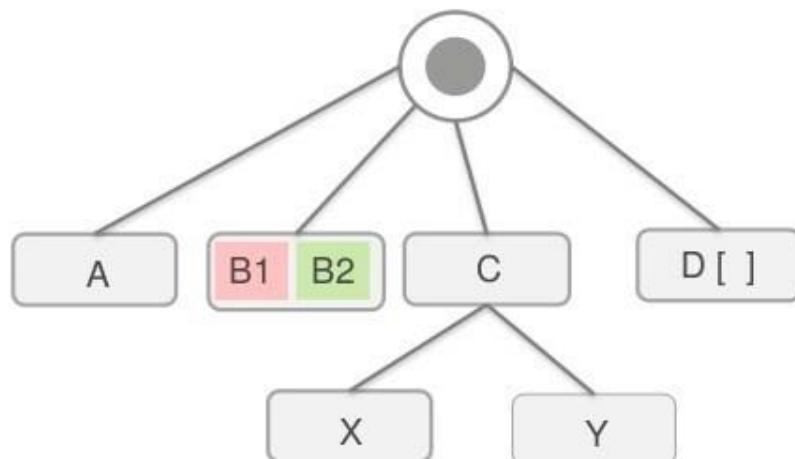
- **Flexible**: Glue's ETL library simplifies manipulating complex, semi-structured data
- **Customizable**: Use native PySpark, import custom libraries, and/or leverage Glue's libraries

```
42 #####
43 ##
44 ## PySpark Logic to do lots of custom stuff...
45 ##
46 #####
47 DataFrame0 = DynamicFrame.toDF(SelectFields0)
48
49 DataFrame0 = DataFrame0.withColumn("pickup_datetime", DataFrame0["pickup_datetime"].cast("timestamp"))
50 DataFrame0 = DataFrame0.withColumn("dropoff_datetime", DataFrame0["dropoff_datetime"].cast("timestamp"))
51 DataFrame0 = DataFrame0.withColumn("type", lit(recordtype))
52
```

- **Collaborative**: share code snippets via GitHub, reuse code across jobs

# Job Authoring: Glue Dynamic Frames

## Dynamic frame schema



Like Spark's Data Frames, but better for:

- Cleaning and (re)-structuring **semi-structured** data sets, e.g. JSON, Avro, Apache logs ...

No upfront schema needed:

- Infers schema on-the-fly, enabling transformations in a **single pass**

Easy to handle the unexpected:

- Tracks new fields, and inconsistent changing data types with **choices**, e.g. integer or string
- Automatically mark and separate error records

# Job Authoring: Glue transforms

Adaptive and flexible

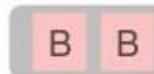
ResolveChoice()



project



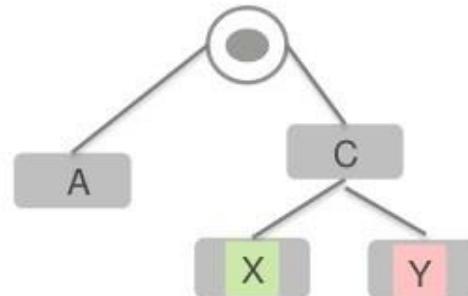
cast



separate into cols

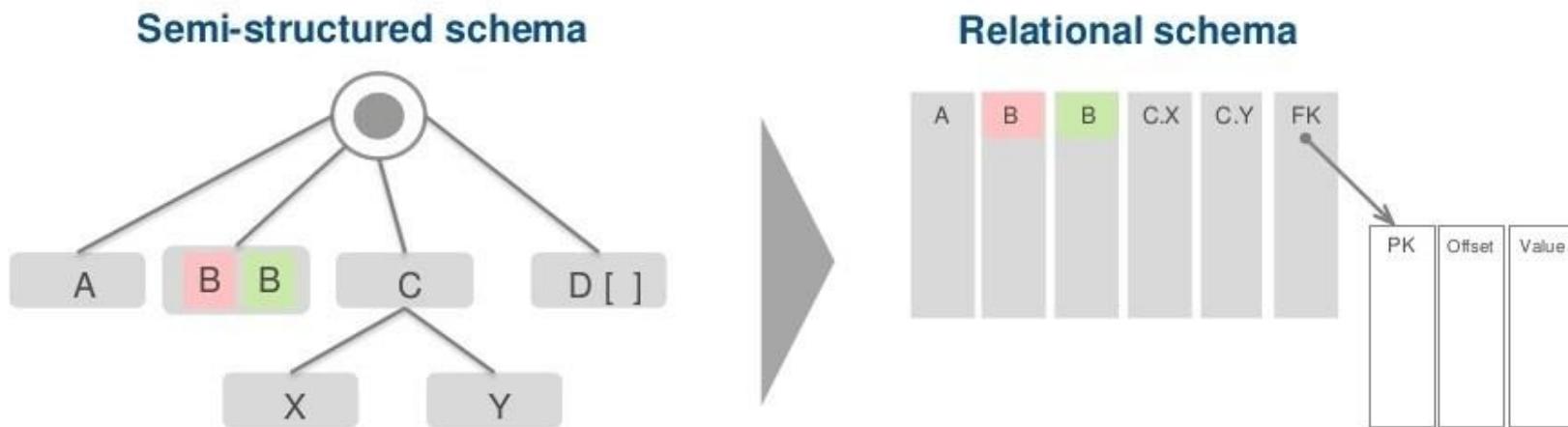


Apply Mapping()



# Job authoring: Relationalize() transform

---



- Transforms and **adds new** columns, types, and tables on-the-fly
- Tracks **keys** and **foreign keys** across runs
- SQL on the relational schema is orders of **magnitude faster** than JSON processing

# Job authoring: Glue transformations

---

## Add transform

Name	Description
DropFields	Drop fields from a DynamicFrame
DropNullFields	DynamicFrame without null fields
Join	Join two DynamicFrames
MapToCollection	Apply a transform to each DynamicFrame in this DynamicFrameCollection
Relationalize	Flatten nested schema and pivot out array columns from the flattened frame
RenameField	Rename a field within a DynamicFrame
SelectFields	Select fields from a DynamicFrame
SelectFromCollection	Select one DynamicFrame from a DynamicFrameCollection
SplitFields	Split fields within a DynamicFrame
SplitRows	Split rows within a DynamicFrame based on comparators
Unbox	Unbox a string field

- **Prebuilt transformation:** Click and add to your job with simple configuration
- **Spigot** writes sample data from DynamicFrame to S3 in JSON format
- **Expanding...** more transformations to come

# Job authoring: Write your own scripts

## fromDF

```
fromDF(dataframe, glue_ctx, name)
```

Converts a DataFrame to a DynamicFrame by converting DynamicRecords to Rows. Returns the new DynamicFrame.

- `dataframe` – The spark SQL dataframe to convert. Required.
- `glue_ctx` – The [GlueContext](#) (p. 164) object that specifies the context for this transform. Required.
- `name` – The name of the resulting DynamicFrame. Required.

## toDF

```
toDF(options)
```

Converts a DynamicFrame to an Apache DataFrame. Returns the new DataFrame.

- `options` – A list of options. Please specify the target type if you choose the Project and Cast action type. Examples are:

```
>>>toDF([ResolveOption("a.b.c", "KeepAsStruct")])  
>>>toDF([ResolveOption("a.b.c", "Project", DoubleType())])
```

Import custom libraries required by your code

Convert to a Spark Data Frame  
*for complex SQL-based ETL*

Convert back to Glue Dynamic Frame  
*for semi-structured processing and AWS Glue connectors*

### Parameters (optional)

► Advanced properties

#### ▪ Script libraries and job parameters (optional)

Server-side encryption

#### Python library path

```
a3://bucket-name/folder-name/file-name
```

#### Dependent jars path

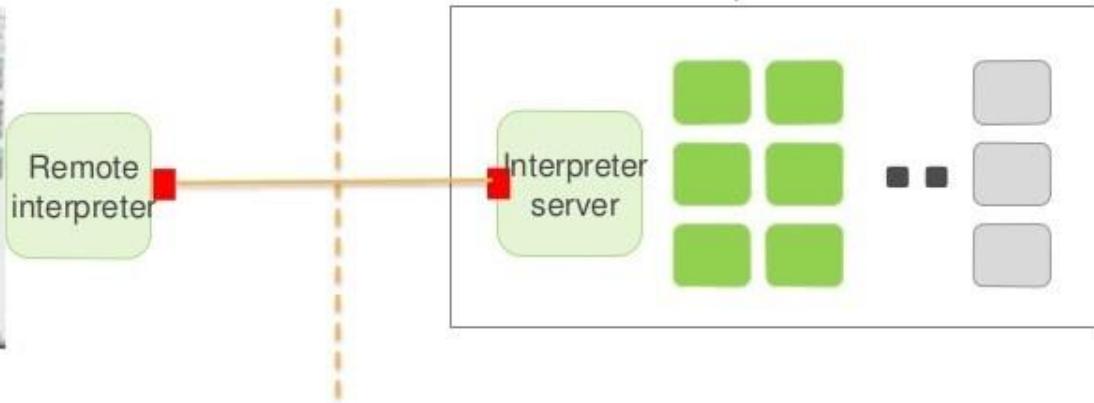
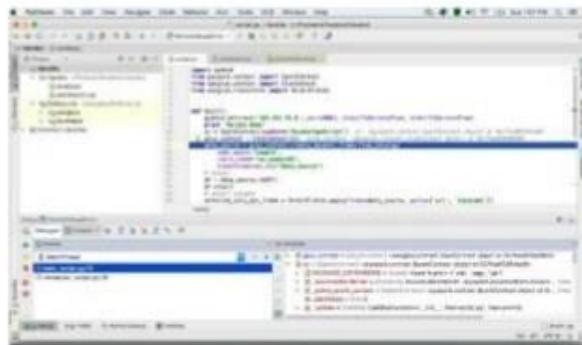
```
a3://bucket-name/folder-name/file-name
```

#### Referenced files path

```
a3://bucket-name/folder-name/file-name
```



# Job authoring: Developer endpoints



- Environment to **iteratively develop** and test ETL code.
- Connect your IDE or notebook (e.g. **Zeppelin**) to a Glue development endpoint.
- When you are satisfied with the results you can create an ETL job that runs your code.

# Job Authoring: Leveraging the community

---

No need to start from scratch.

Use **Glue samples** stored in Github to share, reuse,  
contribute: <https://github.com/awslabs/aws-glue-samples>

- Migration scripts to import existing Hive Metastore data into AWS Glue Data Catalog
- Examples of how to use Dynamic Frames and Relationalize() transform
- Examples of how to use arbitrary PySpark code with Glue's Python ETL library



Download **Glue's Python ETL library** to start developing code in your IDE: <https://github.com/awslabs/aws-glue-lbs>



Pop-up Loft  
LONDON

# Orchestration and Job Scheduling



# Job execution: Scheduling and monitoring

---

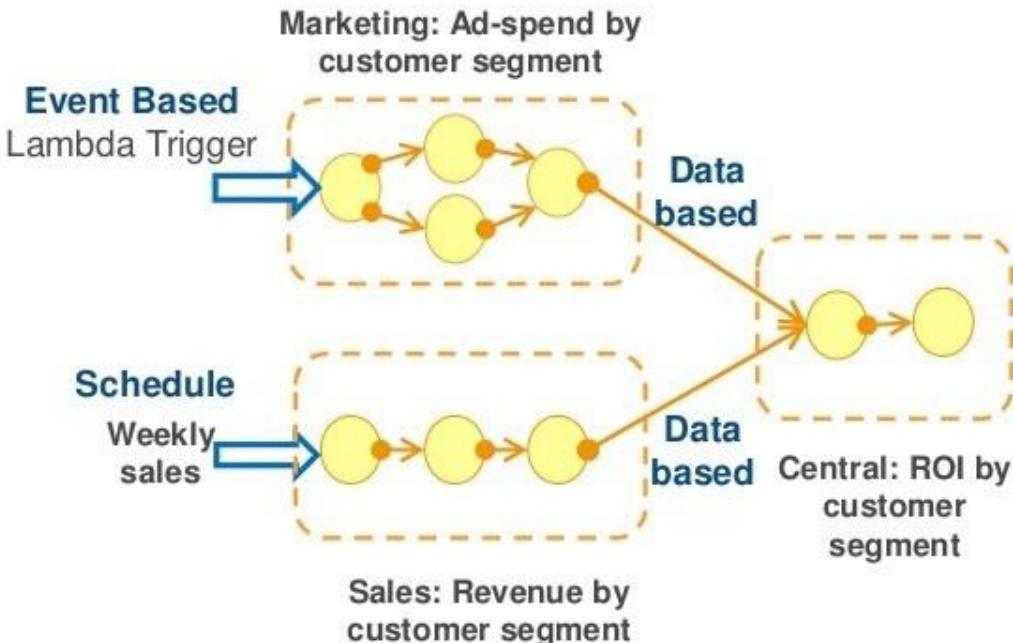
Compose jobs globally with event-based dependencies

- Easy to reuse and leverage work across organization boundaries

Multiple triggering mechanisms

- **Schedule-based:** e.g., time of day
- **Event-based:** e.g., job completion
- **On-demand:** e.g., AWS Lambda
- **More coming soon:** Data Catalog based events, S3 notifications and Amazon CloudWatch events

Logs and alerts are available in Amazon CloudWatch

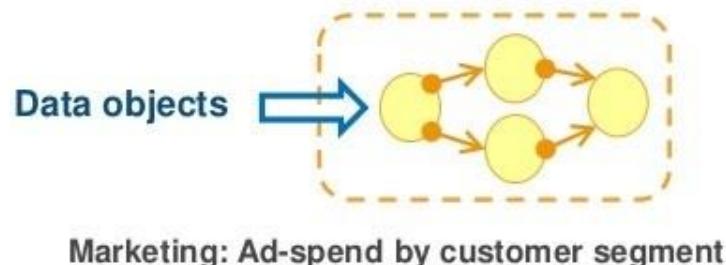


# Job execution: Job bookmarks

Glue keeps track of data that has already been processed by a previous run of an ETL job. This persisted state information is called a **bookmark**.

Option	Behavior
Enable	Pick up from where you left off
Disable	Ignore and process the entire dataset every time
Pause	Temporarily disable advancing the bookmark

**For example**, you get new files everyday in your S3 bucket. By default, AWS Glue keeps track of which files have been successfully processed by the job to prevent data duplication.

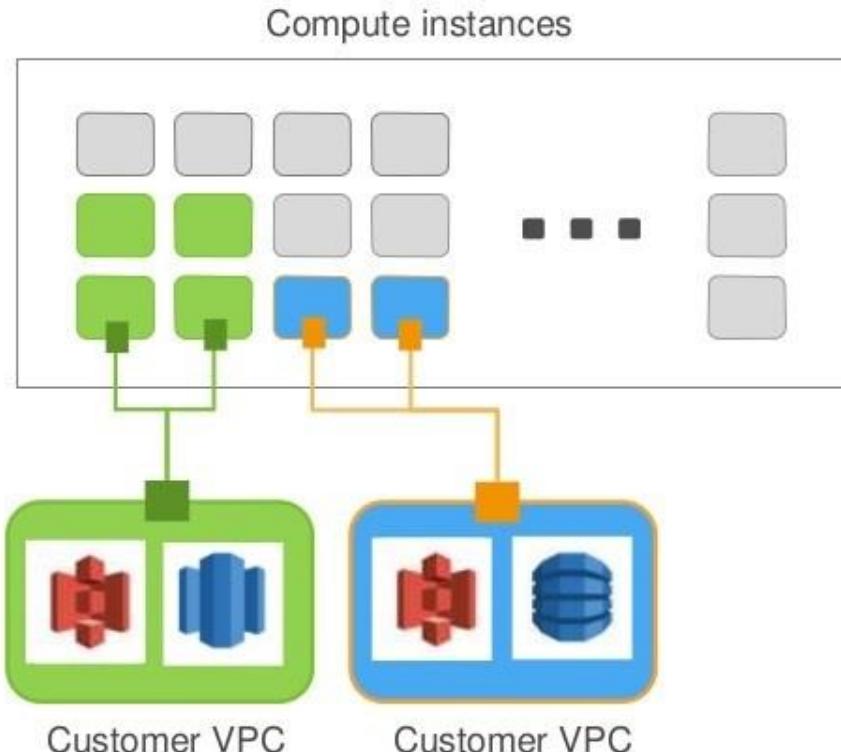


## Job execution: Serverless

---

There is no need to provision, configure, or manage servers

- Auto-configure VPC and role-based access
- Customers can specify the capacity that gets allocated to each job
- Automatically scale resources (on post-GA roadmap)
- You pay only for the resources you consume while consuming them



# Introduction to Amazon Athena

# What to Expect from the Session

Introduction to Serverless

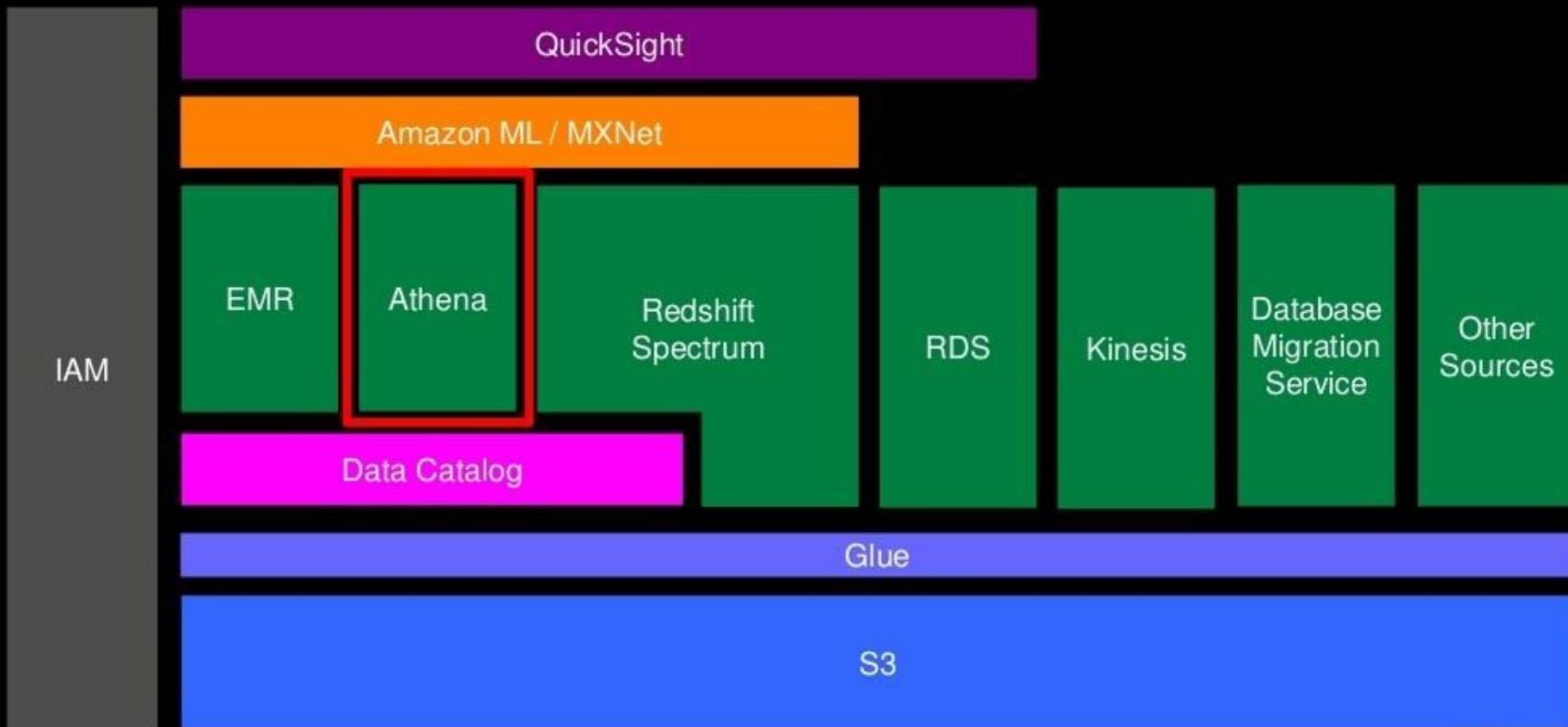
Big data challenges

Introduction to Athena

Demo

Customer Use Cases

# Amazon Analytics End to End Architecture



# Cloud Services Evolution

Virtual  
machines



Managed  
services



Serverless



# Serverless



# Serverless characteristics



No servers to provision  
or manage



Scales with usage



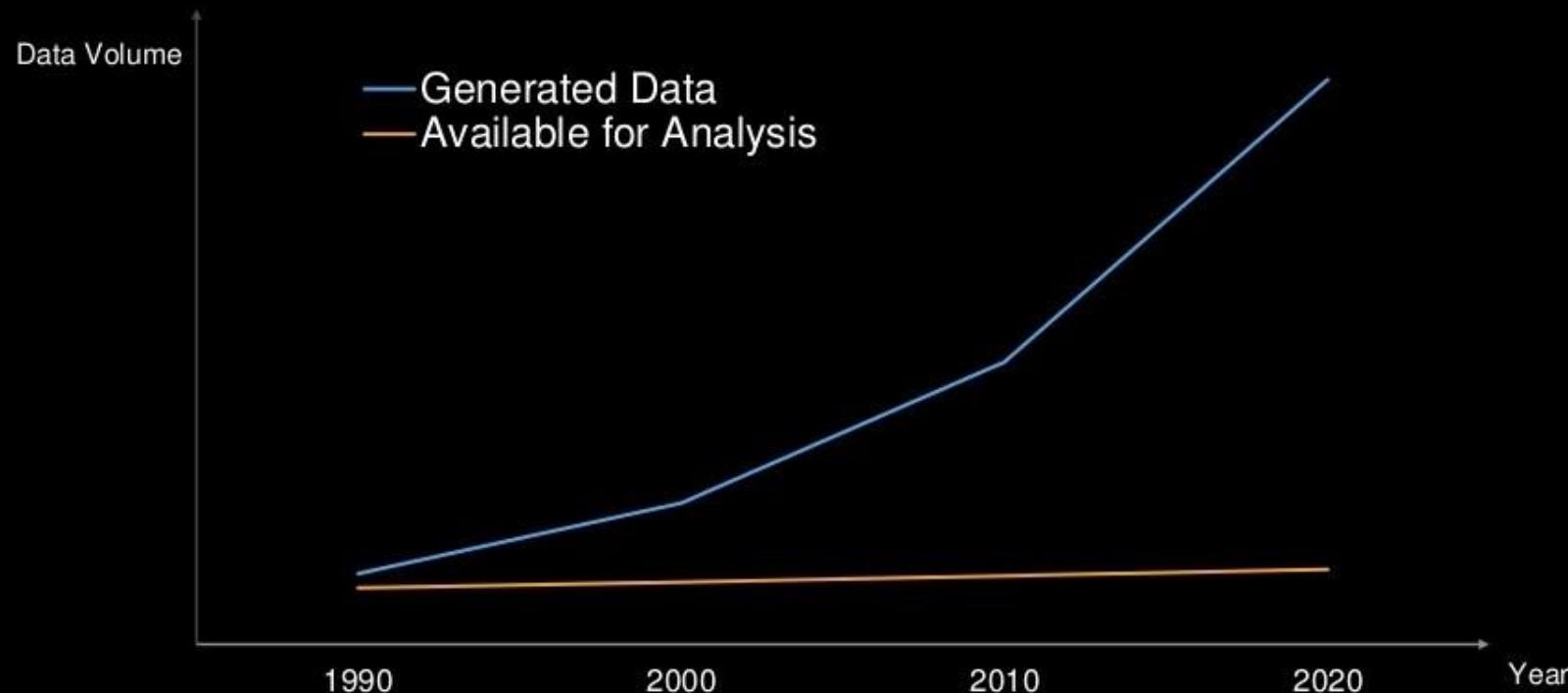
Never pay for idle



Availability and fault  
tolerance built in

# The Dark Data Problem

Most generated data is unavailable for analysis



Sources:

Gartner: User Survey Analysis: Key Trends Shaping the Future of Data Center Infrastructure Through 2011

IDC: Worldwide Business Analytics Software 2012–2016 Forecast and 2011 Vendor Shares





# Amazon Athena

Amazon Athena is an **interactive query service**  
that makes it easy to analyze data directly from  
Amazon S3 using Standard SQL



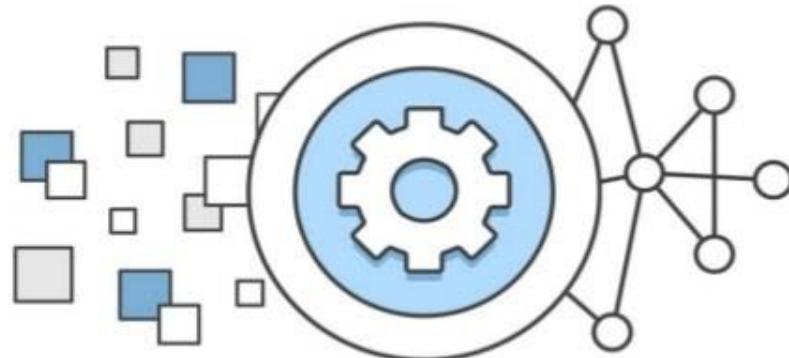
# Value Proposition

- Decouple storage from compute
- Serverless – No infrastructure or resources to manage
- Pay only for data scanned
- Schema on read – Same data, many views
- Secure – IAM for authentication; Encryption at rest
- Standard compliant and open storage formats
- Built on powerful community supported OSS solutions



# Athena is Serverless

- No Infrastructure or administration
- Zero Spin up time
- Transparent upgrades



# Easy To Use

- Log into the Console
- Create a table
  - Type in a Hive DDL Statement
  - Use Glue Data Catalog
  - Use the console Add Table wizard
- Start querying

The screenshot shows the AWS Glue Data Catalog query editor interface. The top bar displays the title "A-Sample-ELB-Query" and the subtitle "Sample Query (Parquet)". Below the title is a code editor containing a Hive DDL statement:

```
1 SELECT elb_name,
2     uptime,
3     downtime,
4     cast(downtime as DOUBLE)/cast(uptime as DOUBLE) uptime_downtime_ratio
5 FROM
6     (SELECT elb_name,
7         sum(case elb_response_code
8             WHEN '200' THEN
9                 1
10            ELSE 0 end) AS uptime, sum(case elb_response_code
11             WHEN '404' THEN
12                 1
13            ELSE 0 end) AS downtime
14     FROM elb_logs_raw_native
15     GROUP BY elb_name)
```

Below the code editor are several buttons: "Cancel" (red), "Save As", "Format Query", and "New Query".

The main area below the buttons is labeled "Results" and shows the status "Running Query...". A note at the bottom states: "You can run another query by clicking on the New Query button. The current query will continue to run in the background."



# Highly Available

- You connect to a service endpoint or log into the console
- Athena uses warm compute pools across multiple Availability Zones
- Your data is in Amazon S3, which is also highly available and designed for 99.99999999% durability

# Query Data Directly from Amazon S3

- No loading of data
- Query data in its raw format
  - Text, CSV, JSON, weblogs, AVRO, AWS service logs
  - Convert to an optimized form like ORC or Parquet for the best performance and lowest cost
- No ETL required (but it can help!)
- Stream data from directly from Amazon S3
- Take advantage of Amazon S3 durability and availability



# Simple Pricing

- DDL operations – FREE
- SQL operations – FREE
- Query concurrency – FREE
- Data scanned - \$5 / TB
- Standard S3 rates for storage, requests, and data transfer apply



# Familiar Technologies Under the Covers



## Used for DDL functionality

- Complex data types
- Multitude of formats
- Supports data partitioning



## Used for SQL Queries

- In-memory distributed query engine
- ANSI-SQL compatible with extensions



# Presto SQL

- ANSI SQL compliant
- Complex joins, nested queries & window functions
- Complex data types (arrays, structs, maps)
- Partitioning of data by any key
  - date, time, custom keys
- Presto built-in functions

```
1 WITH q21_tep1_cached AS
2   (SELECT t_orderkey,
3    count(DISTINCT t_supkey) AS count_supkey,
4    max(t_supkey) AS max_supkey
5    FROM lineitem_parq
6    WHERE t_orderkey IS NOT NULL
7    GROUP BY t_orderkey),
8   q21_tep2_cached AS
9   (SELECT t_orderkey,
10    count(DISTINCT t_supkey) count_supkey,
11    max(t_supkey) AS max_supkey
12    FROM lineitem_parq
13    WHERE t_recupdate = t_commitdate
14    AND t_orderkey IS NOT NULL
15    GROUP BY t_orderkey)
16   SELECT s_name,
17    count() AS numall
18   FROM
19   (SELECT s_name
20    FROM
21    (SELECT t_orderkey,
22     t_supkey,
23     count_supkey,
24     max_supkey
25     FROM q21_tep2_cached t2
26     RIGHT OUTER JOIN
27     (SELECT t_orderkey,
28      t_supkey
29      FROM
30      (SELECT t_orderkey,
31       t_supkey,
32       count_supkey,
33       max_supkey
34       FROM q21_tep1_cached t1
35       JOIN
36       (SELECT t_orderkey,
37        t_supkey
38        FROM orders_parq o
39        JOIN
40        (SELECT s_name,
41         t_orderkey,
42         t_supkey
43         FROM nation_parq n
44         JOIN supplier_s ON s.s_nationkey = n.n_nationkey
45         AND s.s_name = 'SAFEx ANABET'
46         JOIN lineitem_parq l ON l.l_supkey = t1.l_supkey
47         WHERE l.l_recupdate > l.l_commitdate
48         AND l.l_orderkey IS NOT NULL) t3 ON o.o_orderkey = t3.l_orderkey
49         AND o.o_orderstatus = 'F') t2 ON t2.l_orderkey = t1.l_orderkey
50         WHERE (count_supkey > 1)
51         OR (count_supkey = 1
52             AND (t_supkey <= max_supkey))
53         OR ((count_supkey = 1
54             AND (t_supkey > max_supkey))) 13 ON t3.l_orderkey < t2.l_orderkey) b
55         WHERE (count_supkey IS NULL)
56         OR ((count_supkey = 1
57             AND (t_supkey <= max_supkey)))
58         OR ((count_supkey = 1
59             AND (t_supkey > max_supkey))) c
60         GROUP BY s_name
61         ORDER BY numall DESC,
62         s_name LIMIT 300g
```



# Simple Pricing - \$5/TB Scanned

- Pay by the amount of data scanned per query
- Ways to save costs
  - Compress
  - Convert to Columnar format
  - Use partitioning
- Free: DDL Queries, Failed Queries

```
SELECT elb_name,
       uptime,
       downtime,
       cast(downtime as DOUBLE)/cast(uptime as DOUBLE) uptime_downtime_ratio
  FROM
    (SELECT elb_name,
            sum(case elb_response_code
                WHEN '200' THEN
                  1
                ELSE 0 end) AS uptime, sum(case elb_response_code
                WHEN '404' THEN
                  1
                ELSE 0 end) AS downtime
   FROM elb_logs_raw_native
  GROUP BY elb_name)
```

Dataset	Size on Amazon S3	Query Run time	Data Scanned	Cost
Logs stored as Text files	1 TB	237 seconds	1.15TB	\$5.75
Logs stored in Apache Parquet format*	130 GB	5.13 seconds	2.69 GB	\$0.013
Savings	87% less with Parquet	34x faster	99% less data scanned	99.7% cheaper



# CloudTrail

- CloudTrail tracks all Athena APIs
- Can be queried by Athena
- Includes IAM user and account ID of whom executed queries
- Helpful in tracking query usage



# Athena or Redshift Spectrum (or both)?

- **Athena**
  - Separation of storage and compute
  - Good query performance
  - Text and optimized file formats – help reduce cost and improve performance
  - Serverless
- **Redshift Spectrum**
  - Best cost/query – Best query performance
  - Data Warehouse capabilities
  - Join Redshift and S3 data
  - Cluster management overhead



# How to connect



## AWS Glue Data Catalog

Discover and organize your data

## Glue data catalog

Manage table metadata through a Hive metastore API or Hive SQL.  
Supported by tools like Hive, Presto, Spark etc.

We added a few extensions:

- **Search** over metadata for data discovery
- **Connection info** – JDBC URLs, credentials
- **Classification** for identifying and parsing files
- **Versioning** of table metadata as schemas evolve and other metadata are updated

Populate using Hive DDL, bulk import, or automatically through **Crawlers**.



## Data Catalog: Crawlers

Crawlers automatically build your Data Catalog and keep it in sync

- ▶ Automatically discover new data, extracts schema definitions
  - Detect schema changes and version tables
  - Detect Hive style partitions on Amazon S3
- ▶ Built-in classifiers for popular types; custom classifiers using Grok expressions
- ▶ Run ad hoc or on a schedule; serverless – only pay when crawler runs



# Amazon QuickSight



Amazon QuickSight is a Business Analytics Service that lets business users quickly and easily visualize, explore, and share insights from their data.

# Who Is QuickSight For?

QuickSight is designed to give everyday business users the **flexibility** to do easy, self-serve analysis on their data.

QuickSight is also **perfect** for delivering published dashboards throughout the organization.



Discussing Today



Amazon RDS,  
Aurora



Amazon  
Redshift



Amazon  
Athena

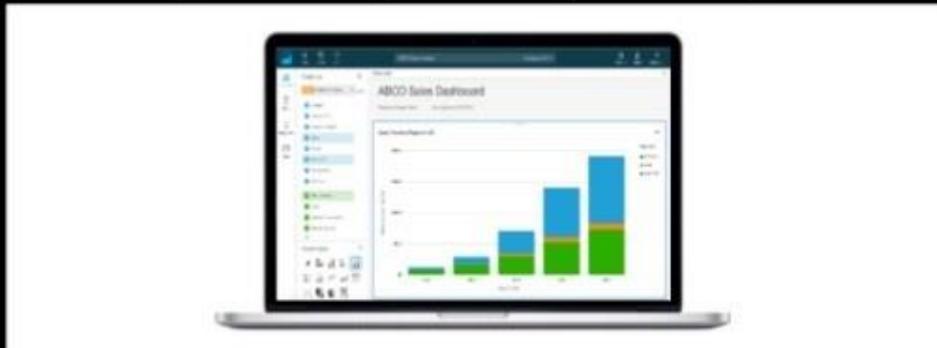


Amazon S3



## Deep Integration with AWS Data Sources

QuickSight is **deeply integrated** with AWS data sources like Redshift, RDS, S3, Athena and others, as well as third-party sources like Excel, Salesforce, as well as on-premises databases.



Flat Files



Microsoft  
Excel



MySQL



PostgreSQL



Microsoft  
SQL Server



# Getting Started with Amazon QuickSight

# AWS Big Data portfolio

## Collect



AWS Direct Connect



AWS Import/Export



Amazon Kinesis

New



Amazon Kinesis Firehose

New



AWS Database  
Migration

## Store



Amazon S3



Amazon RDS/  
Aurora



Amazon  
Glacier



Amazon  
DynamoDB



Amazon  
CloudSearch



Amazon  
Elasticsearch

New



AWS Data  
Pipeline

## Analyze



Amazon  
EMR



Amazon  
Redshift



Amazon  
EC2



Amazon  
Machine  
Learning

New



Amazon  
QuickSight

New



Amazon  
Kinesis  
Analytics

**What are the Big Data challenges  
our customers face?**

How is my marketing campaign performing?

Which devices are showing time for maintenance?

Why is my most profitable region not growing?

Who are my top customers and what are they buying?

## Lots of data

## Lots and lots of questions

## Few insights

How much inventory do I have?

How is my employee satisfaction trending?

What is my product profitability by region?

Has my fraud account expense increased?

# Old-guard BI

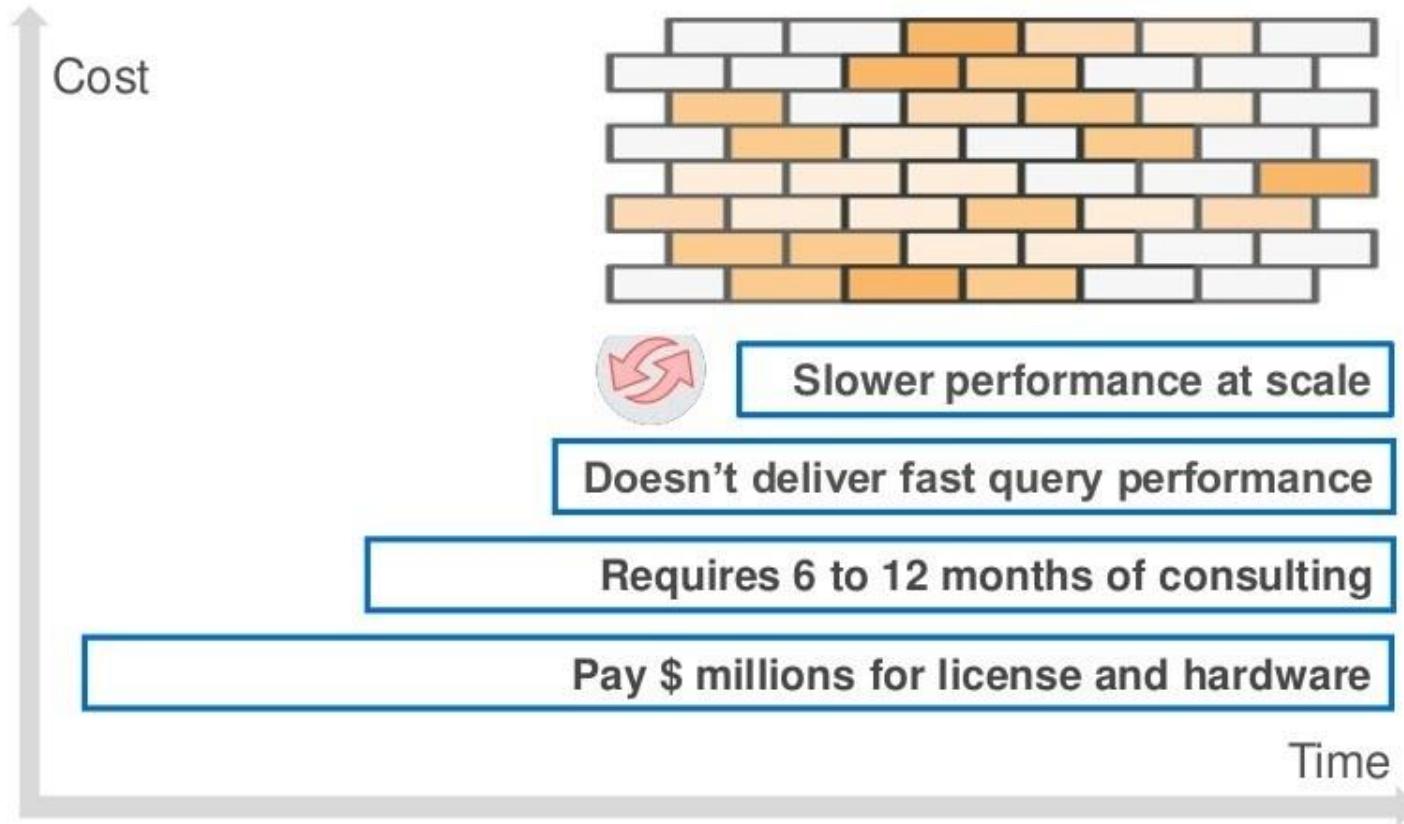
## Costs Too Much

Pay \$ million before seeing first analysis  
3 year TCO \$150 to \$250 per user per month

## Takes Too Long

Spend 6 to 12 months of consulting  
and SW implementation time

# Old-guard BI



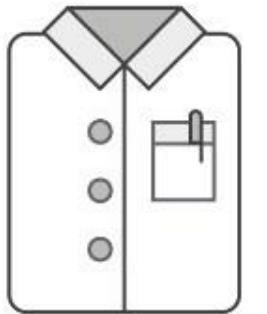
**Freedom to get the real value  
out of the data you have**



# Introducing Amazon QuickSight

**A fast, cloud-powered, BI service for  
1/10<sup>th</sup> the cost of traditional BI software**

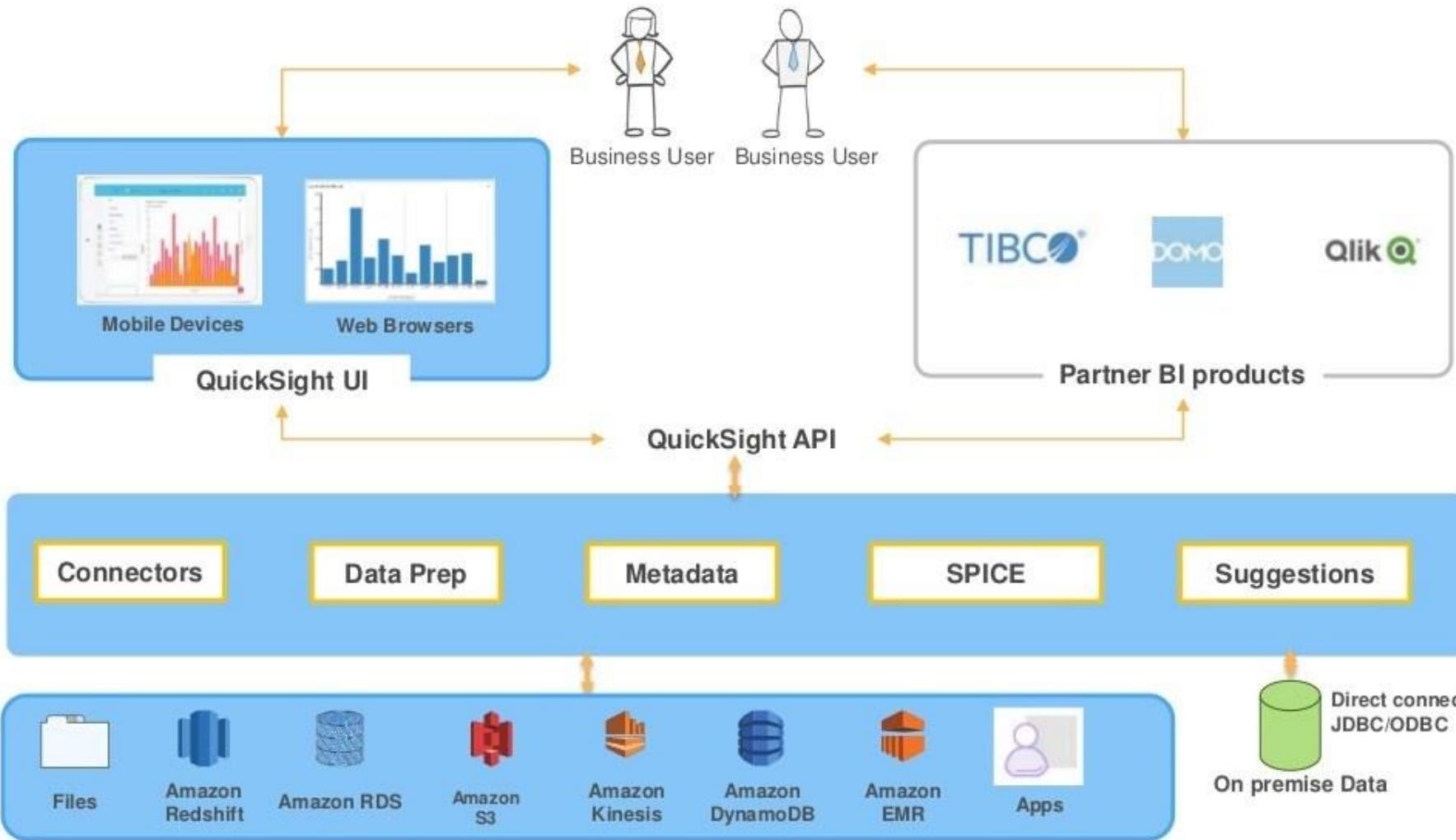
# First analysis in about 60 seconds



→ **Sign-in** →

Business user

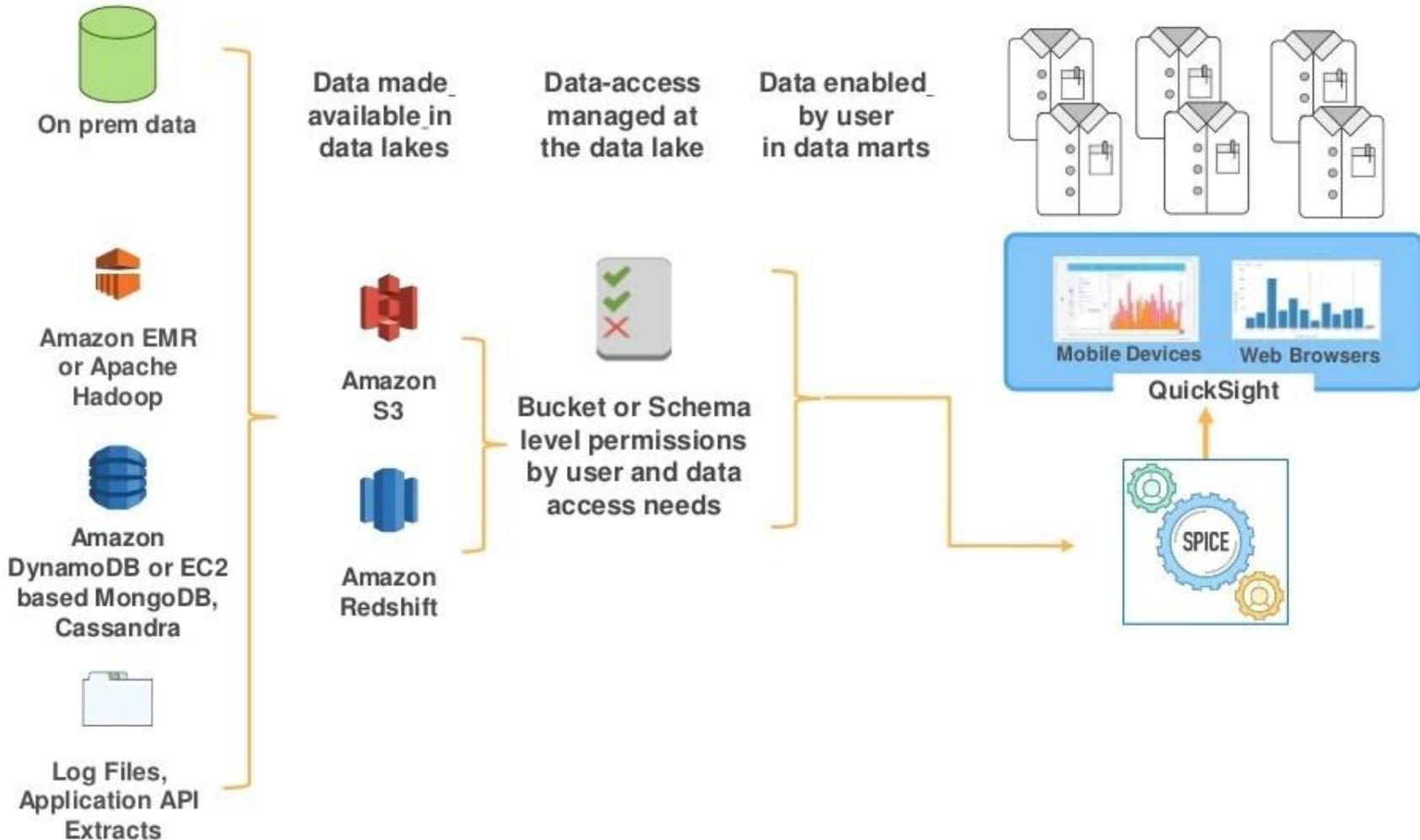




I have multiple data sets both on-premise and on AWS from different sources, and I need to make data available and enable access via Amazon QuickSight.

How do I do this?

1. Data made available in “Data Lakes” using Amazon S3 or Amazon Redshift
2. Data access managed with bucket or schema level policies
3. Data enabled via Amazon QuickSight



# Innovations

- Easy exploration of AWS data
- Fast insights with SPICE (Super-fast, Parallel, In-memory, Calculation Engine)
- Intuitive visualizations and transitions with AutoGraph
- Native mobile experience
- Secure sharing and collaboration through StoryBoards

# Easy exploration of AWS data



Amazon EMR



Amazon RDS



Amazon Kinesis



Amazon S3



Amazon DynamoDB



File Upload



Amazon Redshift



Third Party

Securely discover and connect to AWS data

Quickly explore AWS data sources

- Relational databases
- NoSQL databases
- Amazon EMR, Amazon S3, files
- Streaming data sources

Easily import data from any table or file

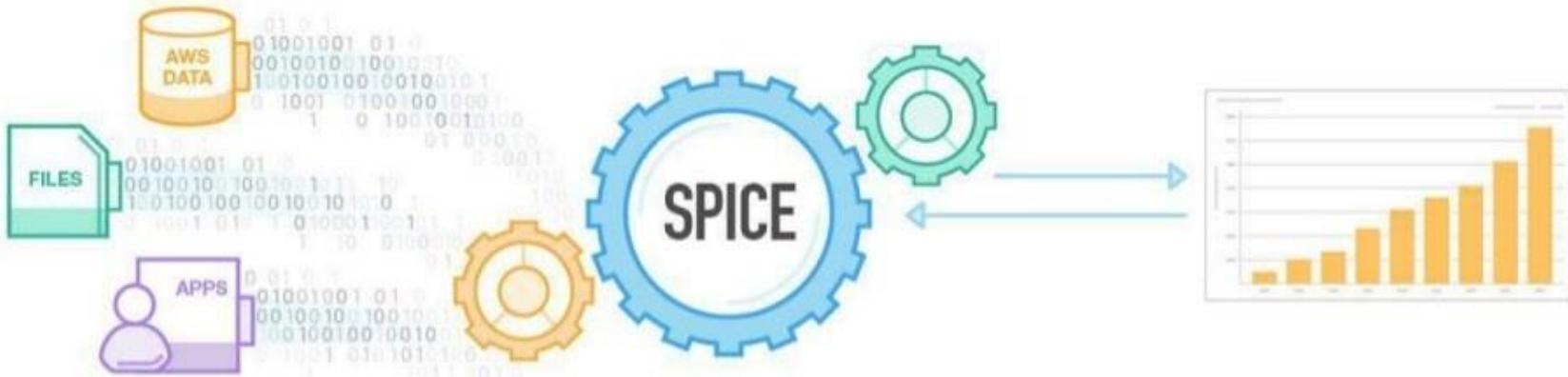
Automatic detection of data types

# **How do I SPICE up my data?**

# Fast insights with SPICE



- Super-fast, Parallel, In-memory optimized, **Calculation Engine**
- 2 to 4x compression columnar data
- Compiled queries with machine code generation
- Rich calculations
- SQL-like syntax
- Very fast response time to queries
- Fully managed – No hardware or software to license



#### DATA SOURCES

Connect to AWS data services;  
upload files; or connect to apps  
such as Salesforce

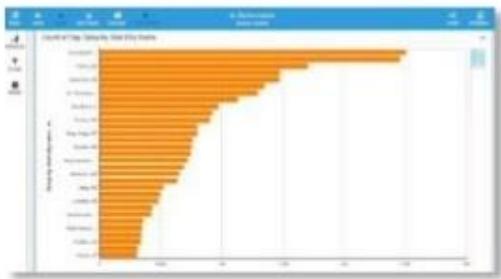
#### IN-MEMORY CALCULATION ENGINE

*The Super-fast, Parallel, In-memory, Calculation Engine ("SPICE") generates answers on large datasets and returns rapid responses*

#### QUICKSIGHT UI

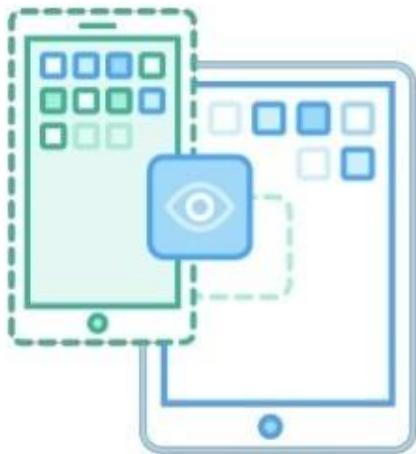
*SPICE allows for very fast analysis and smart visualizations for sharing and collaboration*

# Intuitive visualizations with AutoGraph



- Automatic detection of data types
- Optimal query generation
- Appropriate graph type selection
- Ability to customize the graph type
- Very fast response

# Native mobile experience



- iOS, Android
- Full experience on tablets
- Consumption experience on smart phones
- Very fast response

# Tell a story with your data

- Capture the critical snapshot of analysis
- Build a sequence of analysis
- Share it securely
- Enable interactive exploration
- Very fast response



*Conduct your analysis and capture the visualizations as 'scenes' with annotations and comments*

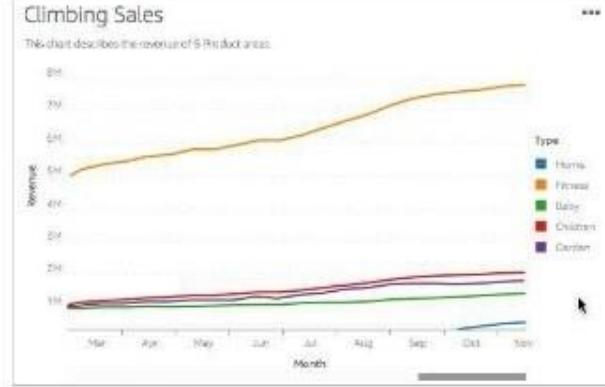
*Share the scenes as a story to be played back anytime for further analysis*

# Dashboard



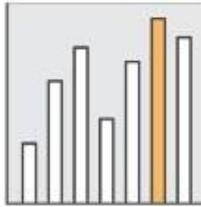
# Embeddable

## Sales Performance





Fast to get started



Easily explore any AWS data



Fast insights with SPICE



Easy to use and share

Effortless scale



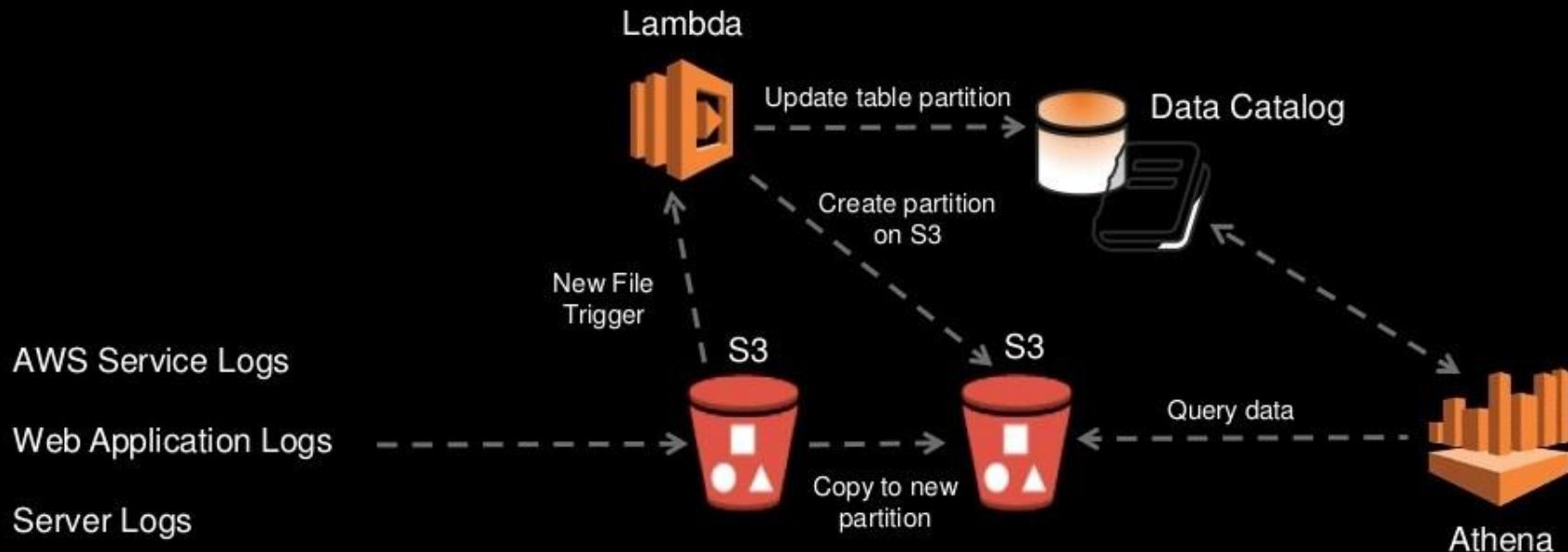
Low cost

# Register for the Preview @ [aws.amazon.com/quicksight](http://aws.amazon.com/quicksight)

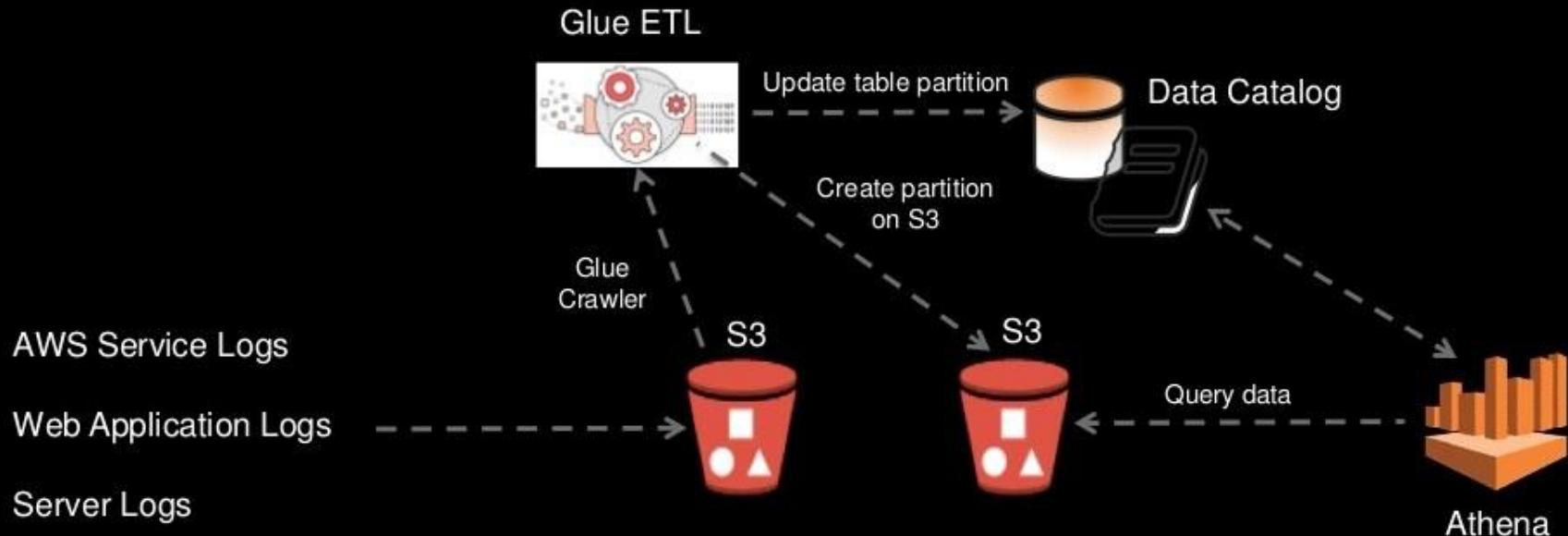


# Popular Customer Use Cases

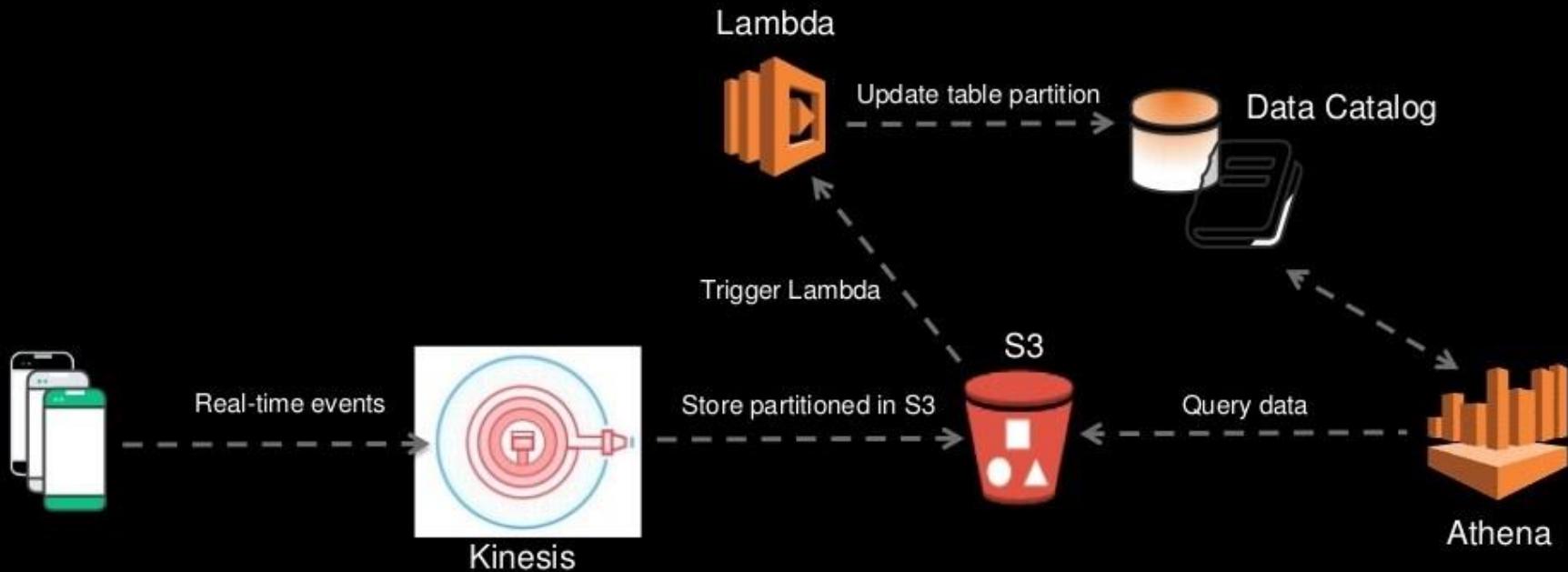
# Log Aggregation



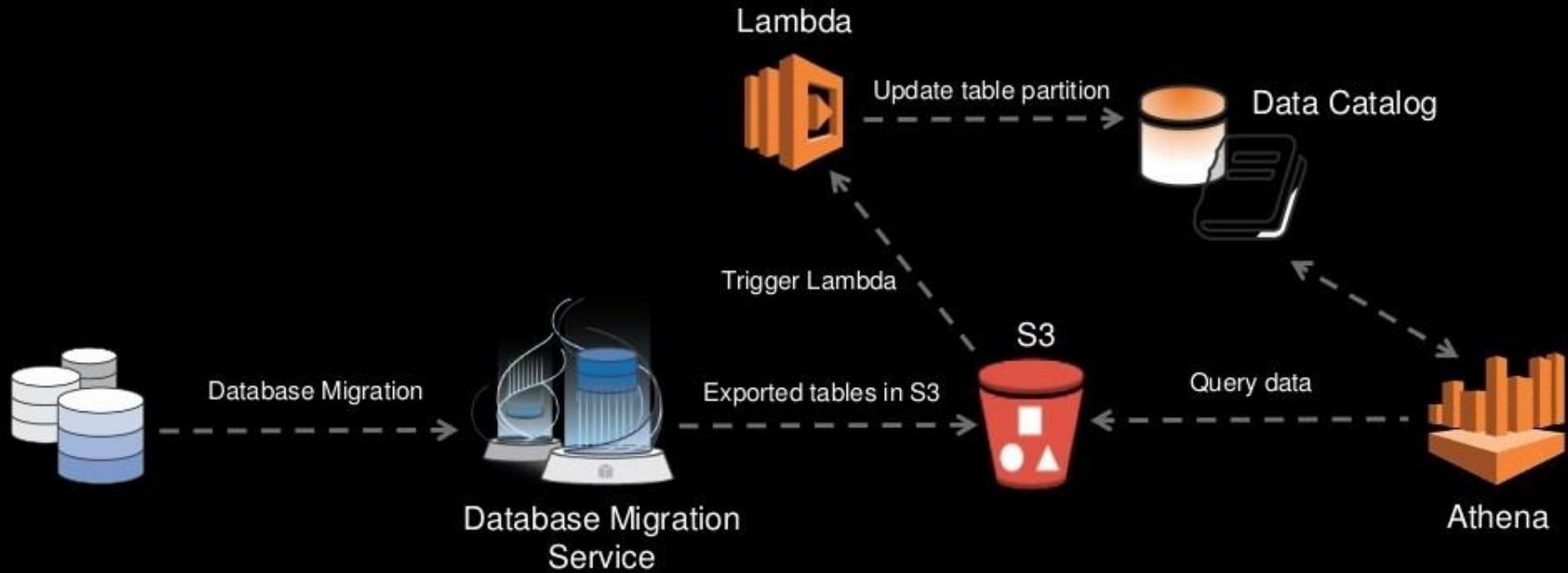
# Log Aggregation with ETL



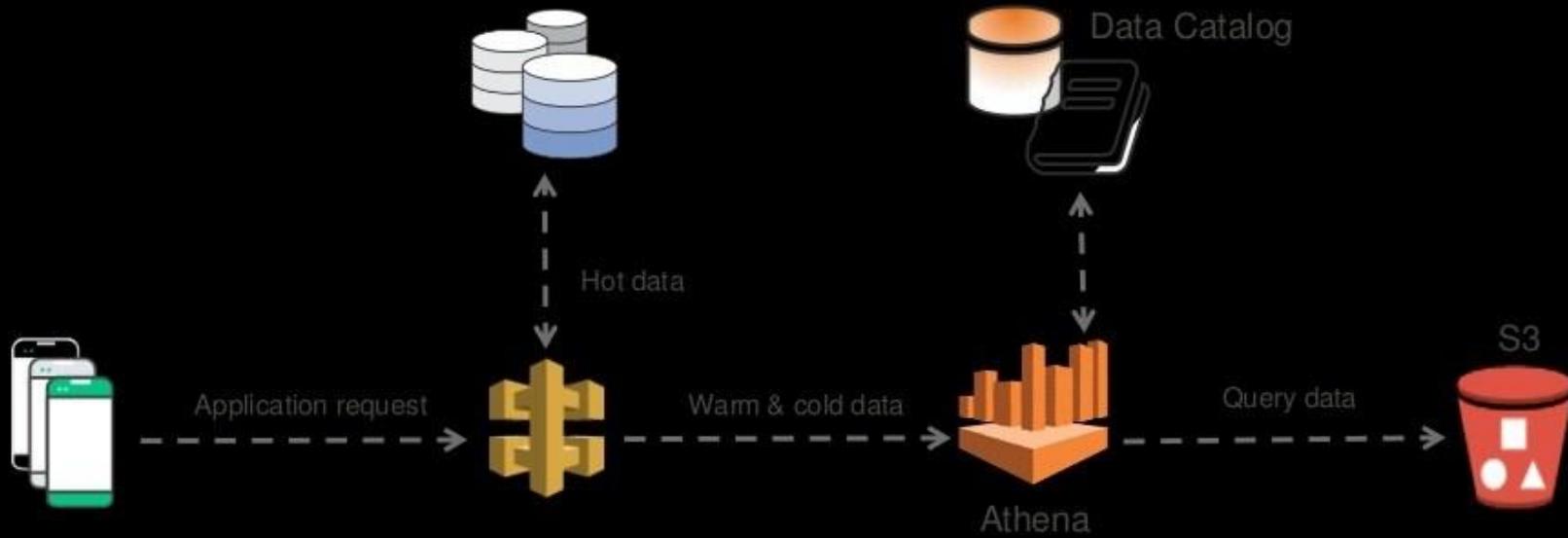
# Real-Time Data Collection



# Data Export



# SaaS Model



# Analytics Reporting\*

