




# From Data Collection to Actionable Insights in 60 Seconds

**Alex Casalboni**

Technical Evangelist, AWS

 [@alex\\_casalboni](https://twitter.com/alex_casalboni)

# About me

- Software Engineer & Web Developer
- Startupper for 4.5 years
- Serverless Lover & AI Enthusiast
- AWS Customer since 2013



# Agenda

1. Data Challenges
2. Columnar Formats
3. Data Lakes vs. Data Warehouses
4. Serverless Analytics
5. Demo time

# Data Challenges

**Data variety and data volumes** are increasing rapidly



Ingest  
Discover  
Catalog  
Understand  
Curate  
Find insights

**Multiple Consumers and Applications**



UP TO  
75 BILLION  
EVENTS PER  
DAY



Monitors  
99% EQUITIES &  
65% OPTIONS  
in the US



Market  
Reconstruction  
Containing  
TRILLIONS of  
nodes & edges

Over 20 PETABYTES of  
storage



Investor  
PROTECTION



Market  
INTEGRITY

THINK  
BIG





# Customer Needs Come First

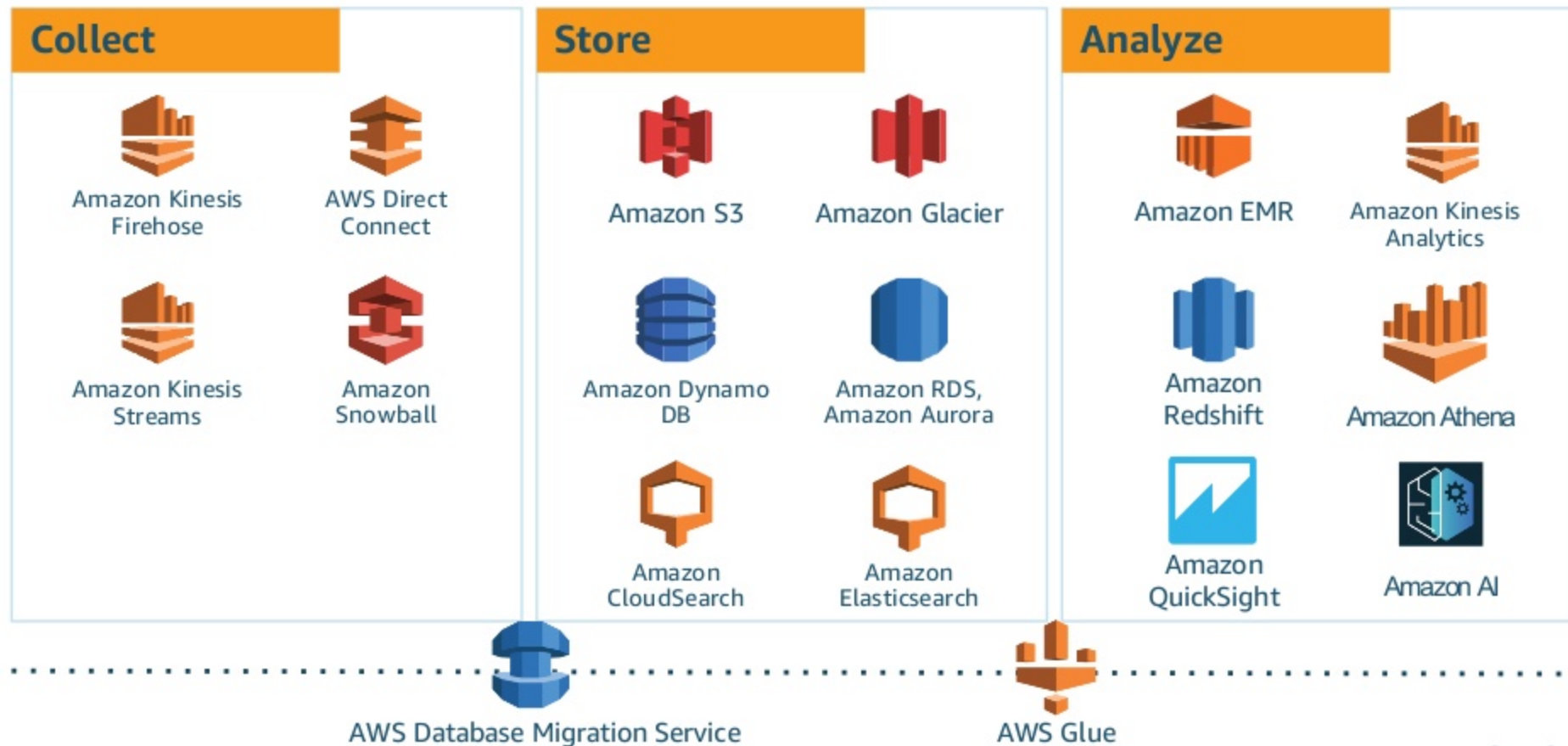


Purpose-built  
engines

Right tool for the job



# Purpose-Built Analytics on AWS



# Columnar data



**Open-source standards (Apache)**

---

**Parquet, ORC, etc.**

---

**Optimize Performance**

---

**Optimize Costs**

---

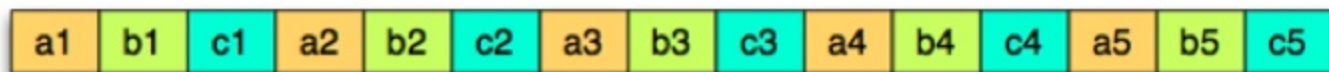
**Analytical queries**



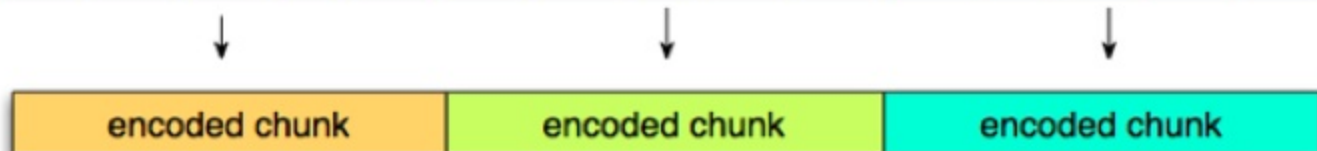
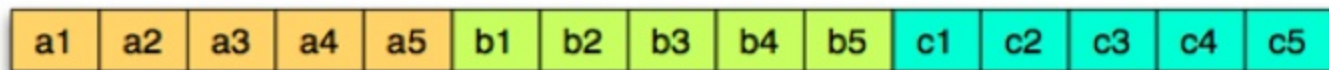
# Under the hood

a	b	c
a1	b1	c1
a2	b2	c2
a3	b3	c3
a4	b4	c4
a5	b5	c5

Row layout



Column layout



# Why it matters

**Big Data Analytics**

**Real-time Analytics**

**Data exploration**



# Traditional Data Warehouse



Relational data

---

Terabytes to Petabytes scale

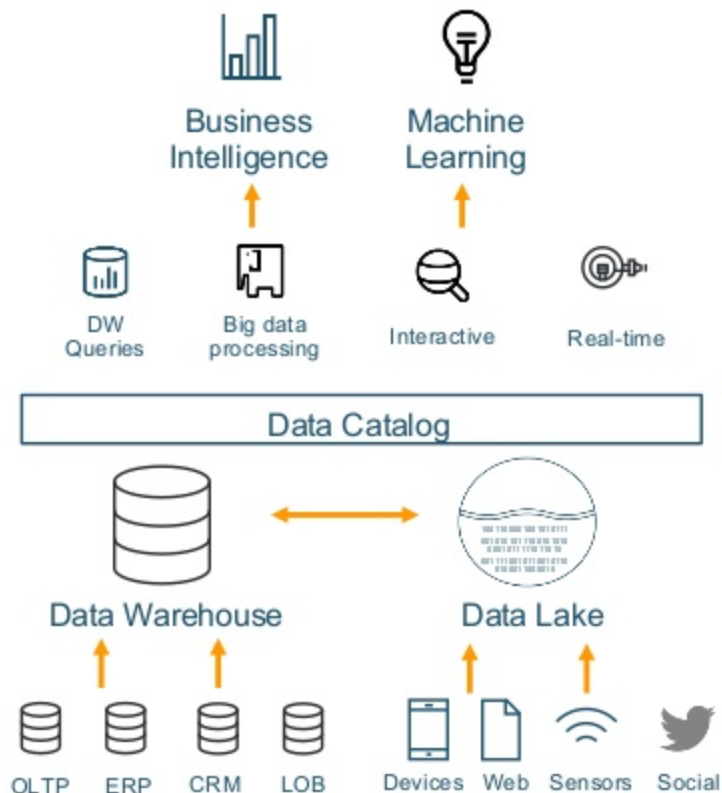
---

Schema defined prior to data load

---

Operational reporting  
and ad-hoc analysis

# Data Lakes extend traditional warehouses



**Relational and non-relational data**

---

**Terabytes to Exabytes scale**

---

**Schema defined during analysis  
(Schema on Read)**

---

**Diverse analytical engines to gain insights**

---

**Designed for low cost storage and analytics**

# Data Lakes on AWS



**Wide variety of ways to bring data in**

---

**Durability and availability at Exabyte scale**

---

**Security, compliance, and audit capabilities**

---

**Run any analytics on the same data without movement**

---

**Scale storage and compute independently**

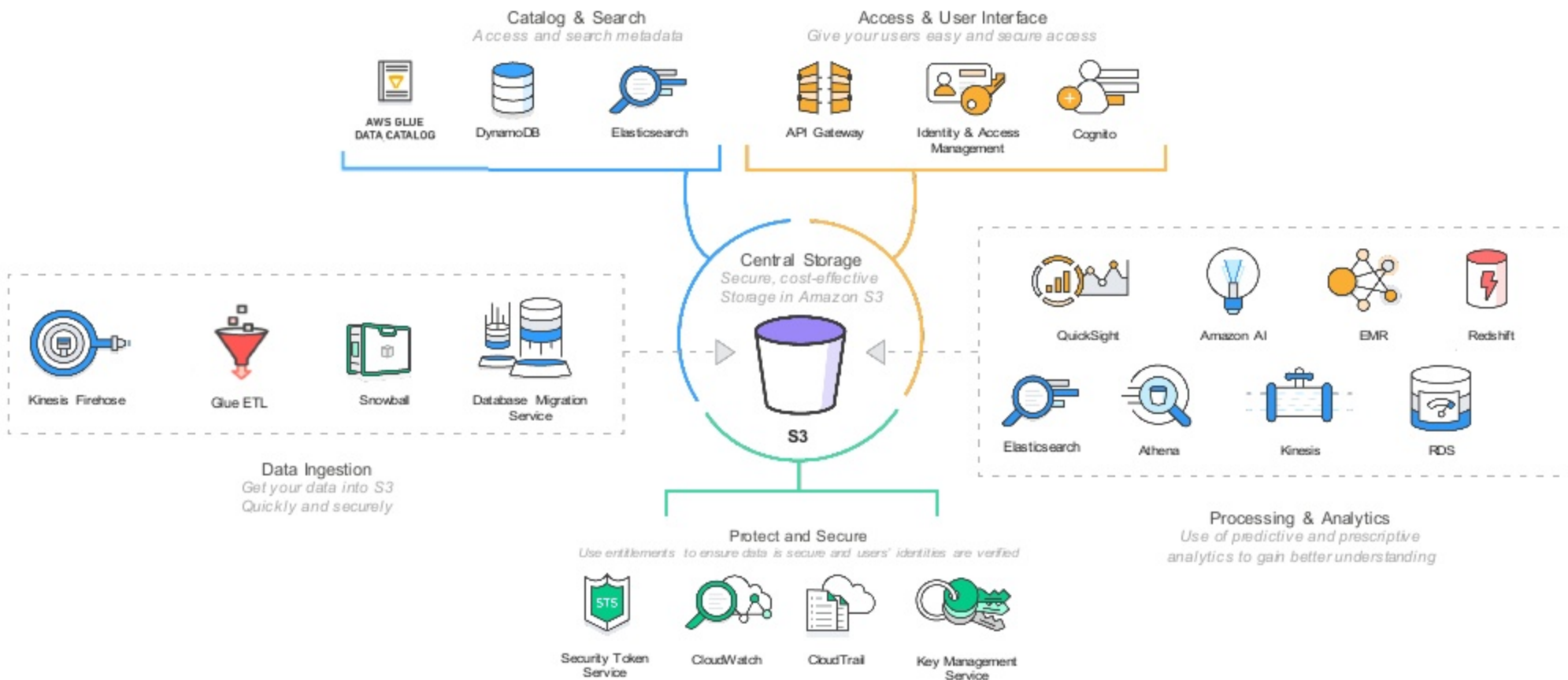
---

**Store at \$0.023 / GB-month**

**Query for \$0.05 / GB scanned**



# Data Lake Components





# Serverless Analytics

Deliver cost-effective analytic solutions faster



Serverless  
Zero infrastructure  
Zero administration



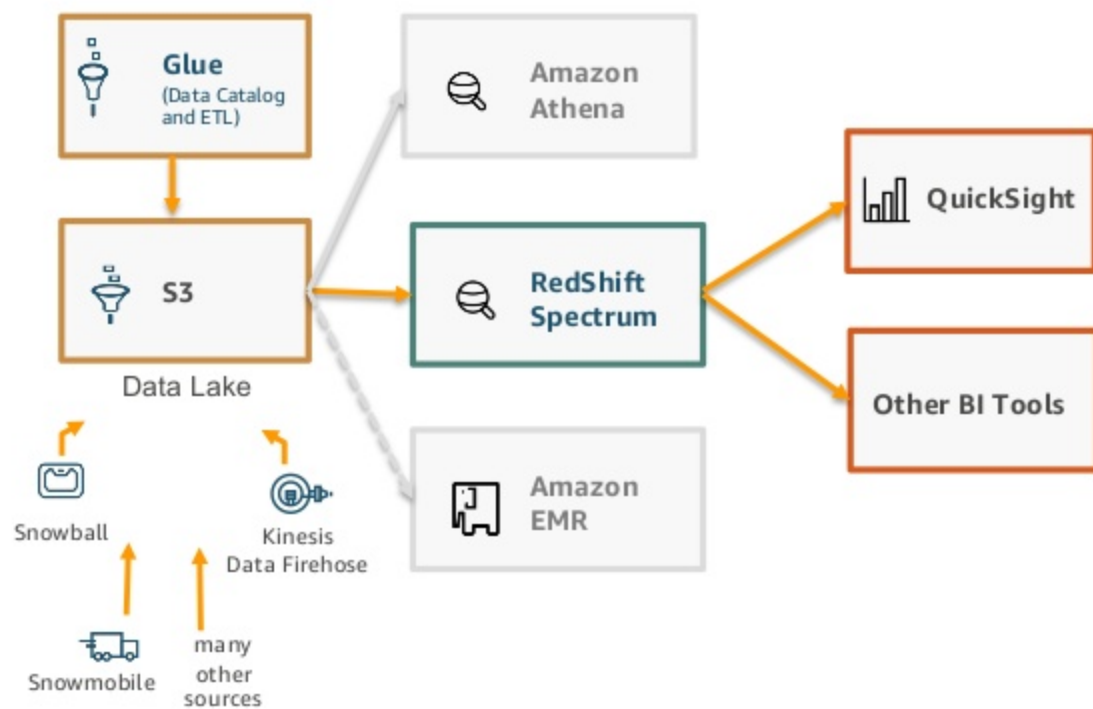
Automatically  
scales resources  
with usage



Pay only for  
what you use,  
not for idle  
resources

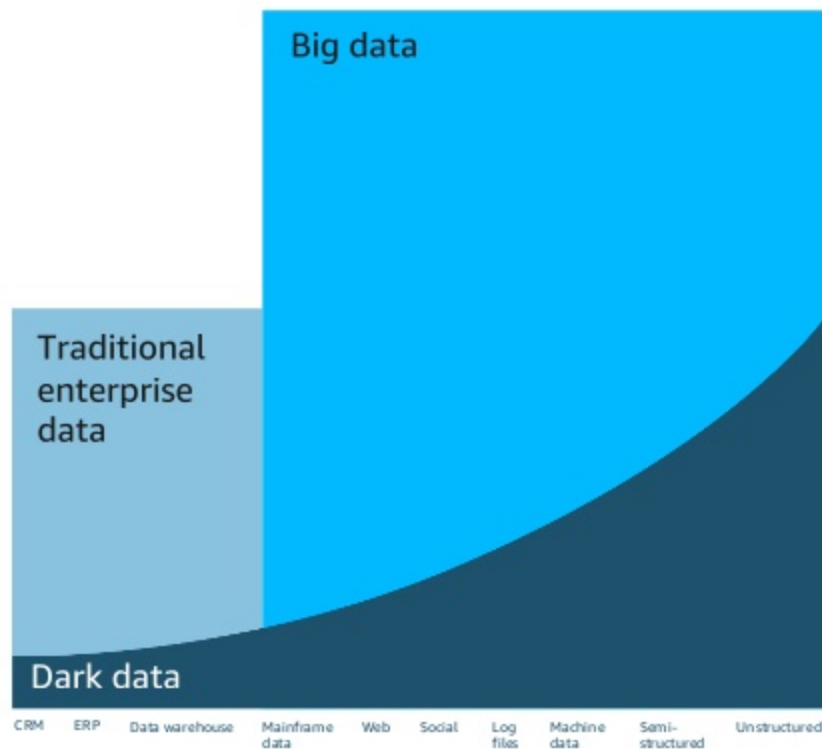


Availability and  
fault tolerance  
built in

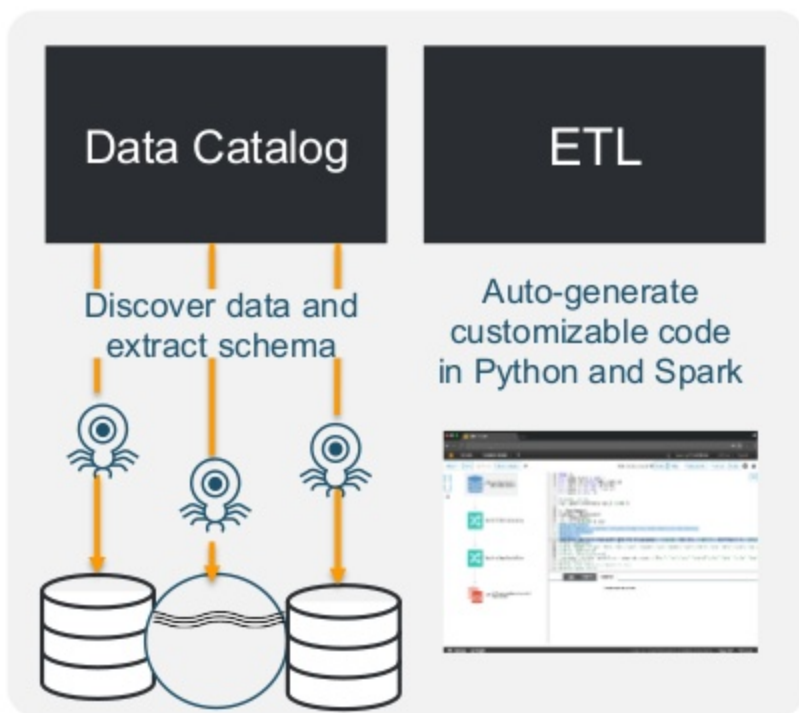


“ Dark data are the information assets organizations collect, process, and store during regular business activities, but generally fail to use for other purposes (for example, analytics, business relationships and direct monetizing). ”

Gartner



# AWS Glue—Serverless Data Catalog & ETL



**Automatically discovers data and stores schema**

---

**Data is immediately searchable  
and available for ETL**

---

**Generates customizable code**

---

**Schedules and runs your ETL jobs**

---

**Serverless Model**

# Crawlers: Automatic Schema Inference

enumerate  
S3 objects

identify file type  
and parse files

semi-structured  
per-file schema

semi-structured  
unified schema

file 1

file 2

...

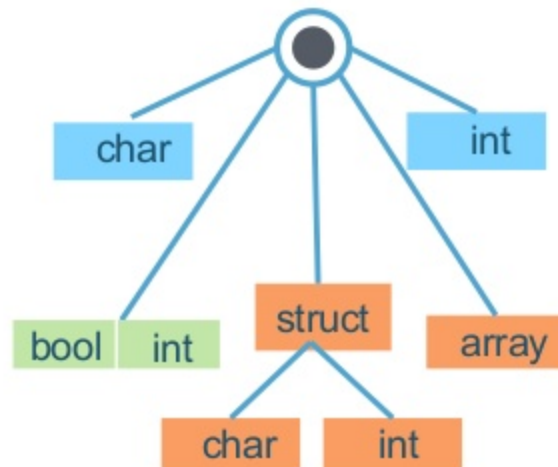
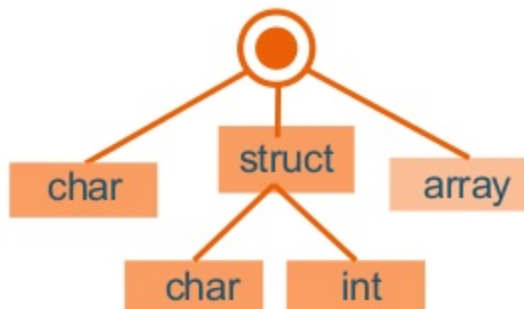
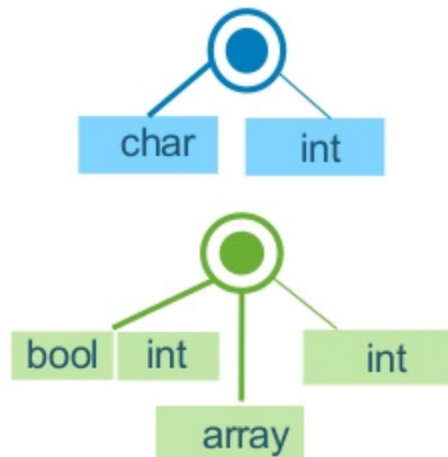
file N

**custom classifiers**

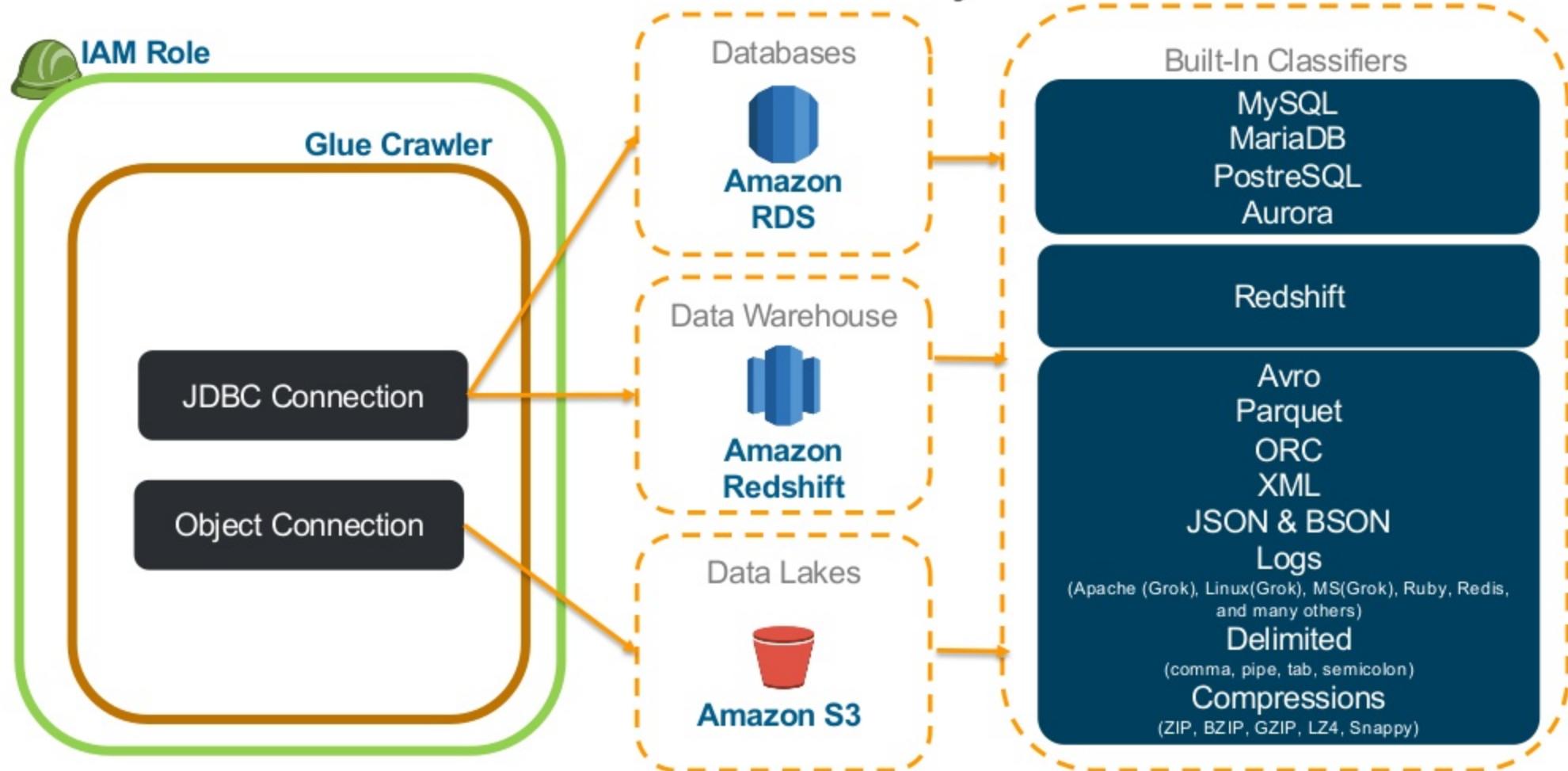
Grok based parser

**built-in classifiers**

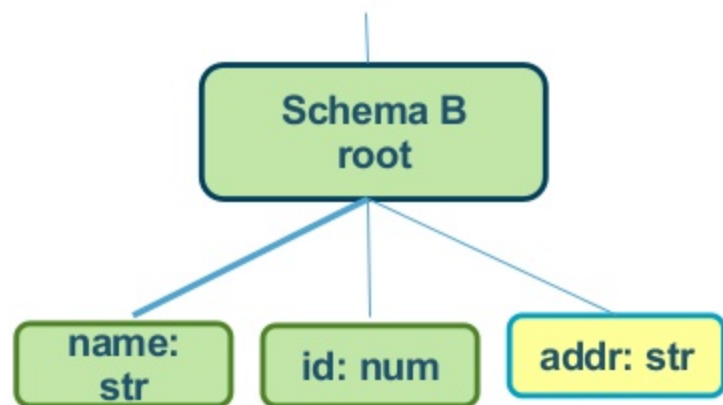
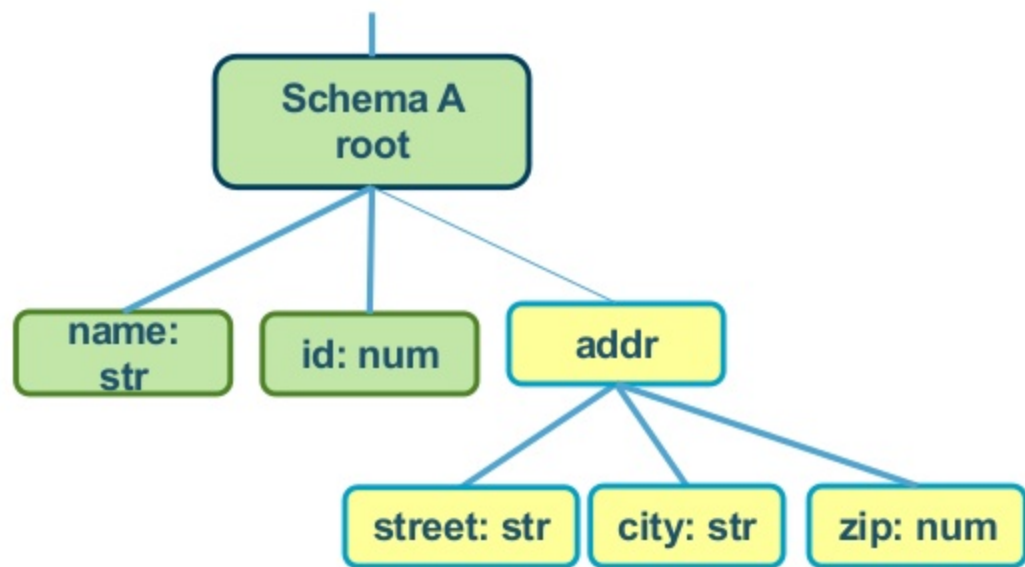
JSON parser  
CSV parser  
Parquet parser  
...



# What can Crawlers Classify?



# Detecting Schema Similarity



## Schema similarity heuristic

- 1 point for matching name
- 1 point for matching data type
- Match when similarity index > 0.7

$$\text{sim} = \frac{\text{intersection}}{\min(A,B)} = \frac{7}{8} = .875$$



## Automatically Detect Partitions

The screenshot shows the AWS Glue console interface. On the left, there is a navigation menu with options like 'Data catalog', 'Tables', 'Connections', 'Crawlers', 'Classifiers', 'ETL', 'Jobs', 'Triggers', 'Dev endpoints', 'Tutorials', 'Add crawler', 'Explore table', 'Add job', and 'Resources'. The main content area displays the details of a crawler named 'glue-crawler'. The crawler is in a 'Stopped' state. The 'Table properties' section shows a table named 'aws-glue-sample-data' with 11 columns: id, type, author, title, year, publisher, crawled\_at, title, year, month, day, and publisher. The crawler is currently in a 'Stopped' state.

### Available partitions

The screenshot displays the AWS Glue console interface for a table named 'githubevents\_data'. The top navigation bar shows the AWS logo, 'Services', 'Resource Groups', and a user profile 'Developer/mashah-Isengard'. Below the navigation bar, the breadcrumb 'Tables > githubevents\_data' is visible. The table's last update date is '22 Nov 2017', and the current version is 'Version (Current version)'. Action buttons include 'Edit table', 'Delete table', 'Close partitions', 'Compare versions', and 'Edit schema'. The table data is shown in a list format with columns 'year', 'month', and 'day'. The data rows are for the year 2017, with months ranging from 02 to 08. Each row has links to 'View files' and 'View properties'. The bottom of the console features a 'Feedback' button, a language selector set to 'English (US)', and links to 'Privacy Policy' and 'Terms of Use'.

year	month	day		
2017	02	15	<a href="#">View files</a>	<a href="#">View properties</a>
2017	03	12	<a href="#">View files</a>	<a href="#">View properties</a>
2017	05	17	<a href="#">View files</a>	<a href="#">View properties</a>
2017	10	12	<a href="#">View files</a>	<a href="#">View properties</a>
2017	12	18	<a href="#">View files</a>	<a href="#">View properties</a>
2017	01	09	<a href="#">View files</a>	<a href="#">View properties</a>
2017	03	07	<a href="#">View files</a>	<a href="#">View properties</a>
2017	06	28	<a href="#">View files</a>	<a href="#">View properties</a>

# Automatic Schema Versioning

Automatically update table version as data evolves

Diagram illustrating Automatic Schema Versioning. An orange bracket connects the 'Version 1' and 'Version 2' table details.

**Version 1** (Last updated 21 Aug 2017 Table)

Name: simpletweets\_json  
Description: simpletweets\_json  
Database: simpletweets\_json  
Classification: json  
Location: s3://gluesampledata/simpletweets\_json  
Connection:   
Deprecated: No  
Last updated: Mon Aug 21 15:23:42 GMT-700 2017  
Input format: org.apache.hadoop.mapred.TextInputFormat  
Output format: org.apache.hadoop.hive.q1.io.HiveIgnoreKeyTextOutputFormat  
Serde serialization lib: org.openx.data.jsonserde.JsonSerDe

Serde parameters: paths entities.id,retweeted,text,user

sizeKey: 456590 objectCount: 1 UPDATED\_BY\_CRAWLER: TestS3Crawler

Table properties: mycustom: abc CrawlerSchemaSerializer/Version: 1.0 recordCount: 1001

averageRecordSize: 456 CrawlerSchemaDeserializer/Version: 1.0

compressionType: none typeOfData: file

Change	Column name	Data type	Key
	id	bigint	
	retweeted	boolean	
	text	string	
	user	string	

**Version 2** (Last updated 25 Nov 2017 Table)

Name: simpletweets\_json  
Description: simpletweets\_json  
Database: simpletweets\_json  
Classification: json  
Location: s3://gluesampledata/simpletweets\_json  
Connection:   
Deprecated: No  
Last updated: Sat Nov 25 12:30:28 GMT-800 2017  
Input format: org.apache.hadoop.mapred.TextInputFormat  
Output format: org.apache.hadoop.hive.q1.io.HiveIgnoreKeyTextOutputFormat  
Serde serialization lib: org.openx.data.jsonserde.JsonSerDe

Serde parameters: paths entities.id,retweeted,text,user

sizeKey: 456590 objectCount: 1 UPDATED\_BY\_CRAWLER: TestS3Crawler

Table properties: mycustom: abc CrawlerSchemaSerializer/Version: 1.0 recordCount: 1001

averageRecordSize: 456 CrawlerSchemaDeserializer/Version: 1.0

compressionType: none typeOfData: file

Change	Column name	Data type	Key
	id	bigint	
	retweeted	boolean	
	text	string	
	user	string	
Added	url	string	

## Other Ways of Creating Tables

## Create table manually

Add table

Table properties

Table alias

Table format

Schemas

Keywords

Set up your table's properties

Table name

Database

Select a database

Add database

Description (optional)

Next

## Run Hive DDL statement

```
1 CREATE EXTERNAL TABLE IF NOT EXISTS elb_logs_raw_native_part (  
2   request_timestamp string,  
3   elb_name string,  
4   request_ip string,  
5   request_port int,  
6   backend_ip string,  
7   backend_port int,  
8   request_processing_time double,  
9   backend_processing_time double,  
10  client_response_time double,  
11  elb_response_code string,  
12  backend_response_code string,  
13  received_bytes bigint,  
14  sent_bytes bigint,  
15  request_verb string,  
16  uri string,  
17  protocol string,  
18  user_agent string,  
19  ssl_cipher string,  
20  ssl_protocol string )  
21 PARTITIONED BY(year string, month string, day string)  
22 NOW FORMAT SERDE 'org.apache.hadoop.hive.serde2.RegexSerDe'  
23 WITH SERDEPROPERTIES (  
24   'serialization.format' = '|', 'input.regex' = '([ ] *) ([ ] *) ([ ] *) :([0-9]*) ([ ] *) :([0-9]*) ([0-9-9]*) ([0-9-9-9]) [^|]*$'  
25 LOCATION 's3://athena-examples/elb/raw/'
```



## Call Glue's CreateTable API

## Import from Apache Hive Metastore



# Amazon Redshift - Data Warehousing

Fast, powerful, simple, and fully managed data warehouse at 1/10 the cost

Massively parallel, scale from gigabytes to petabytes

---

## Fast at any scale



Columnar storage technology to improve I/O efficiency and scale query performance

## Open file formats



Analyze optimized data formats on the latest SSD, and all open data formats in Amazon S3

## Secure



Audit everything; encrypt data end-to-end; extensive certification and compliance

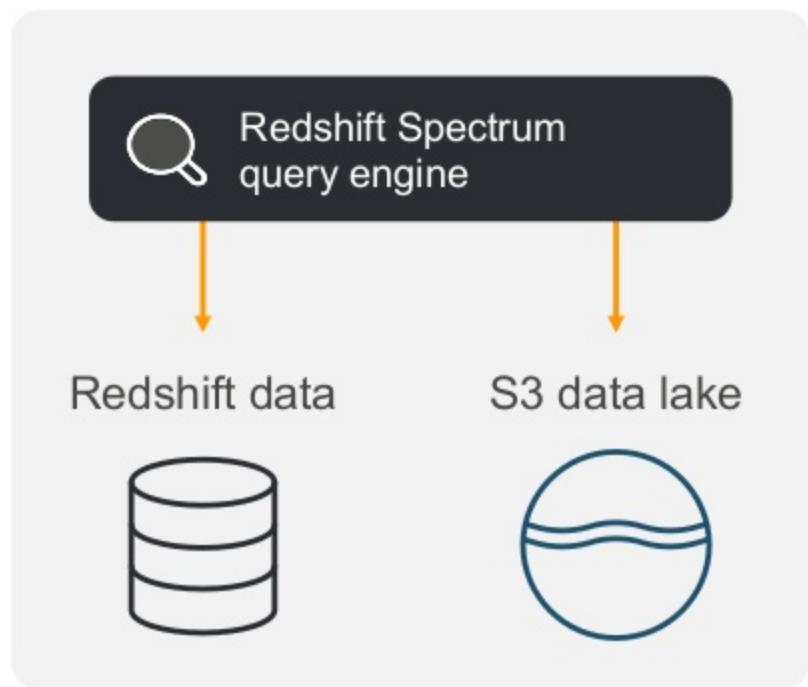
## Inexpensive



As low as \$1,000 per terabyte per year, 1/10th the cost of traditional data warehouse solutions; Start at \$0.25 per hour

# Amazon Redshift Spectrum

Extend the data warehouse to exabytes of data in an S3 data lake



**Exabyte Redshift SQL queries against S3**

---

**Join data across Redshift and S3**

---

**Scale compute and storage separately**

---

**Stable query performance and unlimited concurrency**

---

**CSV, ORC, Grok, Avro, & Parquet data formats**

---

**Pay only for the amount of data scanned**



# Redshift Spectrum

Query your data lake



```
SELECT COUNT(*)  
FROM S3.EXT_TABLE  
GROUP BY ...
```

JDBC/ODBC

Amazon  
Redshift

**Redshift Spectrum**

Scale-out serverless compute



**Amazon S3**

Exabyte-scale object storage



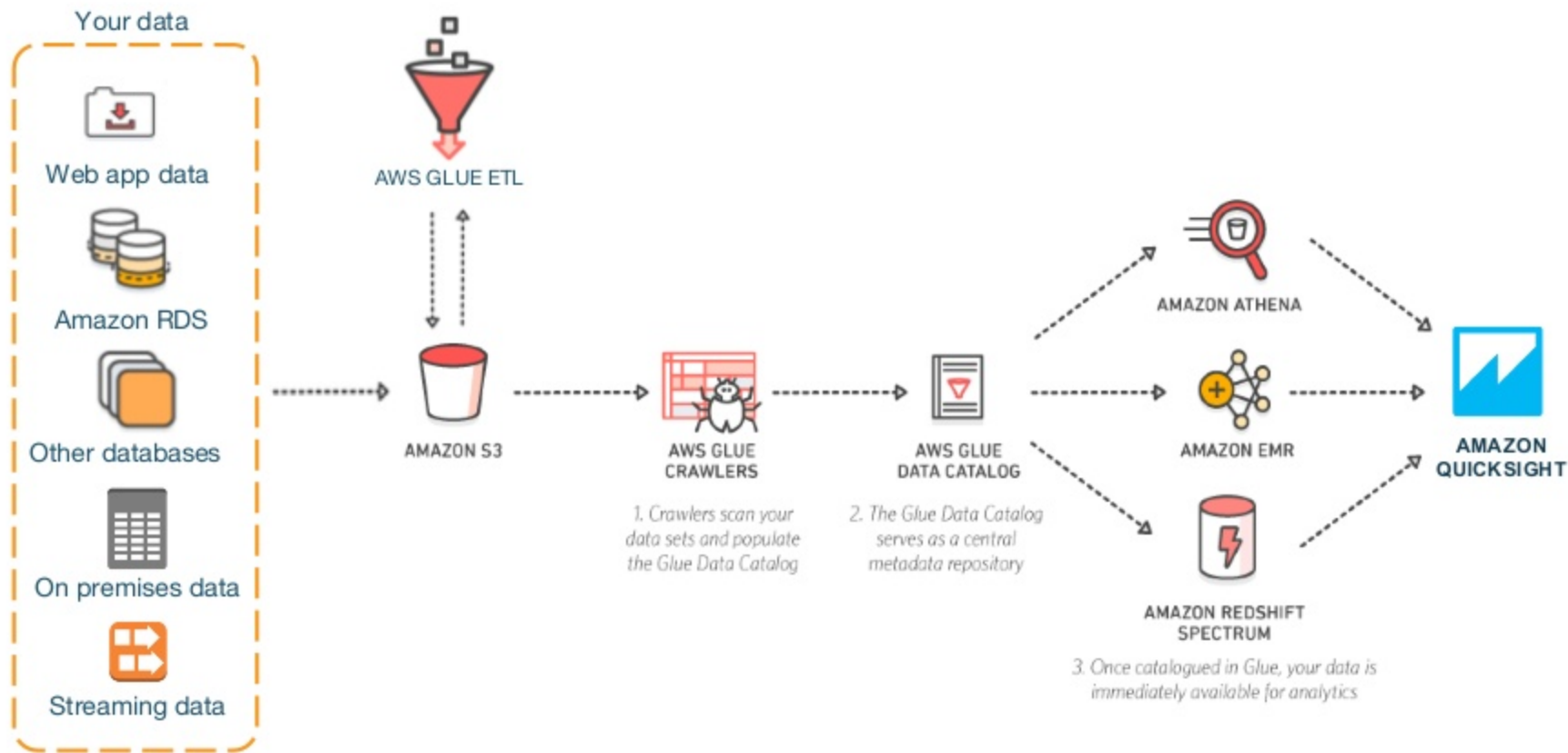
**AWS Glue**

Data Catalog





# Data Lake on Amazon S3 with AWS Glue



# Demo Time

```
SELECT name, avg(value) as average  
FROM table  
WHERE action = 'refill'  
GROUP BY name;
```

**Uncompressed**

**Compressed (-94%)**

**Parquet (-70%)**

**Partitioned (-70%)**

**Overall 99.5% improvement!**

```
SELECT count(*)  
FROM table;
```

**Uncompressed**

**Compressed (-94%)**

**Parquet (-100%)**

**Partitioned (-100%)**

**Overall 100% improvement!**

```
SELECT
  name,
  (
    SELECT count(value)
    FROM table
    WHERE name=t1.name and action='charge'
  ) as charges,
  (
    SELECT count(value)
    FROM table
    WHERE name=t1.name and action='refill'
  ) as refills
FROM table as t1
GROUP BY name;
```

**Uncompressed**

**Compressed (-94%)**

**Parquet (-72%)**

**Partitioned (-100%)**

**Overall 100% improvement!**



**Amazon  
QuickSight**



# Additional Resources

## **Kinesis Data Generator (KDG)**

[github.com/awslabs/amazon-kinesis-data-generator](https://github.com/awslabs/amazon-kinesis-data-generator)

## **Serverless Data Pipeline powered by AWS SAM**

[github.com/alexcasalboni/serverless-data-pipeline-sam](https://github.com/alexcasalboni/serverless-data-pipeline-sam)

## **AWS Big Data Blog**

[aws.amazon.com/blogs/big-data](https://aws.amazon.com/blogs/big-data)



## Did We Scan Your Badge?

---

**Remember to opt-in to AWS communications and you will receive a post-event email with a link to:**

- **AWS Developer Workshop Slides**
- **\$200 in AWS Credits**






# Thank you!

**Alex Casalboni**

Technical Evangelist, AWS

 [@alex\\_casalboni](https://twitter.com/alex_casalboni)