ABD213

# AWS re:INVENT

## How to Build a Data Lake with AWS Glue Data Catalog

Prajakta Damle
Senior Product Manager, AWS Glue

AWS re:Invent

aws

## What to expect from this session

Data challenge today

What is a data lake?

What is AWS Glue Data Catalog?

How does AWS Glue catalogue my data?

My data is catalogued, what's next?

Q&A

# Your data today

*Multiple sources and formats... and growing everyday*

## Documents and files



Clickstream data

Mobile app data

Spreadsheets

Infrastructure logs

Social media data

## Records



Amazon RDS

ERP

Amazon DynamoDB

On Premises databases

Amazon Redshift

## Streams



AWS IoT

Amazon Kinesis Streams

Device data

Amazon Kinesis Firehose

Sensor data

---

# Why is this a new problem?

Structured data



Unstructured and Semi-structured data



Web and mobile data

Logs

Social Media data    Spreadsheets

Streaming data    IOT data

Dark data

# Dark data challenge

**Data Volume**

- Generated Data
- Available for Analysis

1990　2000　2010　2020

# Multiple consumers and requirements

Data Scientists

Business Users

Analysts

Applications

Agile　Real time

Flexible　Scale

**Data duplication**

# What is a data lake?

Collect and store all data, at any scale, and low cost

Help locate, curate, and secure your data

Provide democratized access to data within your organization

Quickly and easily perform new types of data analysis

# Benefits of a data lake

Quickly ingest and store any type of data, at any scale, and at low cost

Have a single source of truth and quickly search and find the relevant data

Easily query the data through a unified set of tools

# Layers of a data lake



**AI**
- Lex
- Amazon Polly
- Amazon Rekognition
- Amazon ML

**Analyze**
- Athena
- AMAZON QUICKSIGHT
- EMR
- AMAZON REDSHIFT

**Secure**
- Identity & Access Management
- Security Token Service
- CloudWatch
- CloudTrail
- Key Management Service

**Store**
- Kinesis
- Direct Connect
- AMAZON S3
- Snowball
- Database Migration Service

# The missing piece

> A unified view into your data no matter where it is stored

> Integration with your analytics tools

> A way to automatically build your metadata and keep it in sync with your data as it evolves

# What is AWS Glue?

**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. Easily customize this code or bring your own.

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

# Select AWS Glue customers



NTT **docomo**

*News Corp*

**AUTODESK.**

**MediaMath**

**MERCK**

my**T**omorrows

**OLX** Group

**ost**

**Expedia®**

**amazon** *Prime Air*

**21ST CENTURY FOX**

aws

# AWS Glue Components

### Data Catalog

**Discover**

Apache Hive Metastore compatible

Integrated with AWS services

Automatic crawling

### Job Authoring

**Develop**

Auto-generates ETL code

Python and Apache

Spark

Edit, debug, and share

### Job Execution

**Deploy**

Serverless execution

Flexible scheduling

Monitoring and alerting

---

# What is the AWS Glue Data Catalog?

**Unified metadata repository** across relational databases, Amazon RDS, Amazon Redshift, and Amazon S3...with support for more coming soon!

- Get a **single view** into your data, no matter where it is stored
- Automatically **classify** your data in one central list that is **searchable**
- Track data evolution using **schema versioning**
- **Query** your data using Amazon Athena or Amazon Redshift Spectrum
- **Apache Hive metastore compatible**; can be used as an external metastore for applications running on Amazon EMR

Data lake on Amazon S3 with AWS Glue



Logical data lake with AWS Glue

AWS Glue supports using Zepellin Notebook that connects to the Data Catalog. We can spin up a Zepellin Notebook and connect it to the Glue Serverless environment, then read and write data from various data sources by leveraging the metadata that is in the Data Catalog.
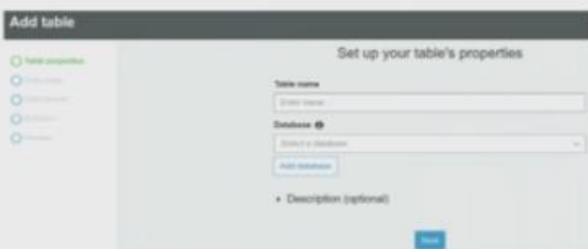
This is good for quick testing if you know the schema



This works great if you have a handful of tables to use.

## Easier way to build the Glue Data Catalog

1. Tell us where your data is
2. Tell us how often you want to check for updates
   And you are done! **Your Data Catalog** is ready for search and querying

You can simply tell Glue where your data is coming from and it will automatically generate a table for you using that source, then it will keep the table data up to date for you to run jobs on. We can have tables that are mapped to data sources in S3, from Redshift, data warehouse, PostGres DBs.

## What are crawlers?

Crawlers automatically build your Data Catalog and keep it in sync
- Scan your data stored in various data stores, extract metadata and data statistics, and add table definitions to your Data Catalog
  - Classify data using built-in and custom classifiers
  - You can write your own using Grok expressions

- Discover new data, extracts schema definitions
  - Detect schema changes and version tables
  - Detect Hive style partitions on Amazon S3

- Run on demand or on a schedule; serverless – only pay when crawler runs

# GitHub timeline data



githubarchive.org

20+ event types

unique payload
per event type

# A table in the Glue Data Catalog



Table properties

Data statistics

Table schema

Nested fields

When Glue crawls the dataset above, it generates the above table entries in the Data Catalog as above.

# How is my data classified?

Crawlers apply a set of classifiers to the data as they scan it and add the metadata as Tables to the Data Catalog.

A **classifier** recognizes the format of your data and generates a schema.
It returns a certainty number between 0.0 and 1.0, which helps crawlers determine if there is a match.

Glue has a list of in-build classifiers that are applied with every crawl. But you can **write your own!**

You can set up your crawler with an ordered set of classifiers. Crawlers invoke classifiers in the order they were provided until a match is found.

---

# Crawlers: automatic schema inference

# Detecting schema similarity

**Schema A root**
- name: str
- id: num
- addr
  - street: str
  - city: str
  - zip: num

**Schema B root**
- name: str
- id: num
- addr: str

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

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

# What can crawlers classify

Create additional Custom Classifiers with Grok!

**IAM Role**

**Glue Crawler**
- JDBC Connection
- Object Connection

**Databases**
Amazon RDS

**Data Warehouse**
Amazon Redshift

**Data Lakes**
Amazon S3

**Built-In Classifiers**

MySQL
MariaDB
PostreSQL
Oracle
Microsoft SQL Server
Amazon Aurora

Amazon Redshift

Avro
Parquet
ORC
XML
JSON & BSON
Logs
(Apache (Grok), Linux(Grok), MS(Grok), Ruby, Redis, and many others)
Delimited
(comma, pipe, tab, semicolon)
Compressions
(ZIP, BZIP, GZIP, LZ4, Snappy)

# How can I write my own classifiers?

You can write a custom classifier by providing a Grok pattern and a classification string for the matched schema.

A Grok pattern is a named set of regular expressions (regex) that are used to match data one line at a time.

Example:
%{TIMESTAMP_ISO8601:timestamp}
\[%{MESSAGEPREFIX:message_prefix}\]
%{CRAWLERLOGLEVEL:loglevel} :
%{GREEDYDATA:message}

**Classifier name**

Id Crawler logs

**Classification**

crawlerlogs

Describes the format of the data classified or a custom label.

**Grok pattern**

%{TIMESTAMP_ISO8601:timestamp} \[%{MESSAGEPREFIX:messi

Built-in and custom named patterns used to parse your data into a structured schema. For more information, see the list of built-in patterns.

**Custom patterns**

```
1  CRAWLERLOGLEVEL (BENCHMARK|ERROR|WARN|INFO|TRACE)
2  MESSAGEPREFIX .*-.*-.*-.*-.*
3
```

Optional custom building blocks for the grok pattern.

---

# Custom classifiers

### 1. Write a custom classifier

**Classifier name**

Id Crawler logs

**Classifier type**

◉ Grok   ○ XML

**Classification**

crawlerlogs

Describes the format of the data classified or a custom label.

**Grok pattern**

%{TIMESTAMP_ISO8601:timestamp} \[%{MESSAGEPREFIX:mes

### 2. Add it to your crawler

Classifiers infer the schema of your data. The first classifier in the list of custom classifiers to recognize your data is used. Subsequent classifiers are skipped. Built-in classifiers are used if you do not supply a classifier that matches.

| Custom classifiers | | Selected classifiers | |
|---|---|---|---|
| Showing: 1 - 2 | | Showing: 1 - 1 | |
| Classifier | Classificatic | Classifier | Classificatic |
| Id Crawle... | crawlerlogs  Add | Id Crawler logs | crawlerlogs  ✕ |
| MyCusto... | MyLogFo...  Add | | |

**Tables**  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.

| ☐ Name | Database | Location | | Classification | Last updated | Deprecated | |
|---|---|---|---|---|---|---|---|
| ☐ exportedlogs | mys3 | s3 | /exportedlogs/ | crawlerlogs | 27 October 2017 12:22 PM U... | | |

# Automatically detected partitions



Available partitions

# How are partitions detected?

**S3 bucket hierarchy**



sim=.93  month=Nov

sim=.99  date=10  ...  sim=.95  date=15

file 1  ...  file N    file 1  ...  file N

**Table definition**

| Column | Type |
|--------|------|
| month  | str  |
| date   | str  |
| col 1  | int  |
|        | float |
| ⋮      | ⋮    |

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

# Automatic schema versioning

Automatically update table version as data evolves

---

# Import/Export your metadata

Import from an external metastore          Export to an external metastore



Find the *import/export ETL script* on Glue's GitHub repository

---

# Your data is catalogued…what's next?

## Quickly find your data

Search on key terms | Save results and come back to it later

Tables  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.

Query data in Amazon Athena

You can do text based search and also filter on key attributes from the Glue console to find relevant datasets. You can also start using this in queries with Athena.



## Analyze same data with different engines

AMAZON S3 → AWS GLUE CRAWLERS → AWS GLUE DATA CATALOG → AMAZON ATHENA / AMAZON EMR → AMAZON QUICKSIGHT / AMAZON REDSHIFT SPECTRUM

1. Crawlers scan your data sets and populate the Glue Data Catalog

2. The Glue Data Catalog serves as a central metadata repository

3. Once catalogued in Glue, your data is immediately available for analytics

# What is Amazon Athena?

Interactive query service to analyze data in Amazon S3 using standard SQL

No infrastructure to set up or manage and no data to load

| Query Instantly | Pay per query | Open | Easy |
|---|---|---|---|
| Zero setup cost; just point to Amazon S3 and start querying | Pay only for queries run; save 30-90% on per-query costs through compression | ANSI SQL interface, JDBC/ODBC drivers, multiple formats, compression types, and complex joins and data types | Serverless: zero infrastructure, zero administration Integrated with Amazon QuickSight |

# What is Amazon EMR?

Analytics and ML at scale with 19 open-source projects

Integration with AWS Glue Data Catalog for Apache Spark, Apache Hive, and Presto

Enterprise-grade security

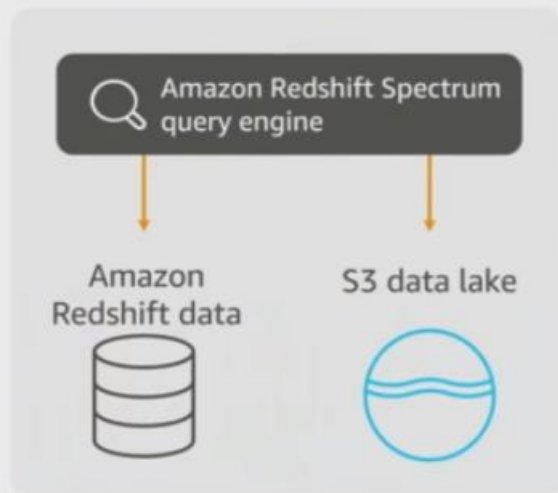| Latest versions | Low cost | Use S3 storage | Easy |
|---|---|---|---|
| Updated with the latest open source frameworks within 30 days of release | Flexible billing with per-second billing, EC2 spot, reserved instances and auto-scaling to reduce costs 50-80% | Process data directly in the Amazon S3 data lake securely with high performance using the EMRFS connector | Launch fully managed Apache Hadoop & Apache Spark in minutes; no cluster setup, node provisioning, cluster tuning |

# What is Amazon Redshift Spectrum?
## Extend the data warehouse to your S3 data lake



Amazon Redshift Spectrum query engine

Amazon Redshift data

S3 data lake

Exabyte Amazon Redshift SQL queries against S3

Join data across Amazon Redshift and S3

Scale compute and storage separately

Stable query performance and unlimited concurrency

Parquet, ORC, Grok, Avro, & CSV data formats
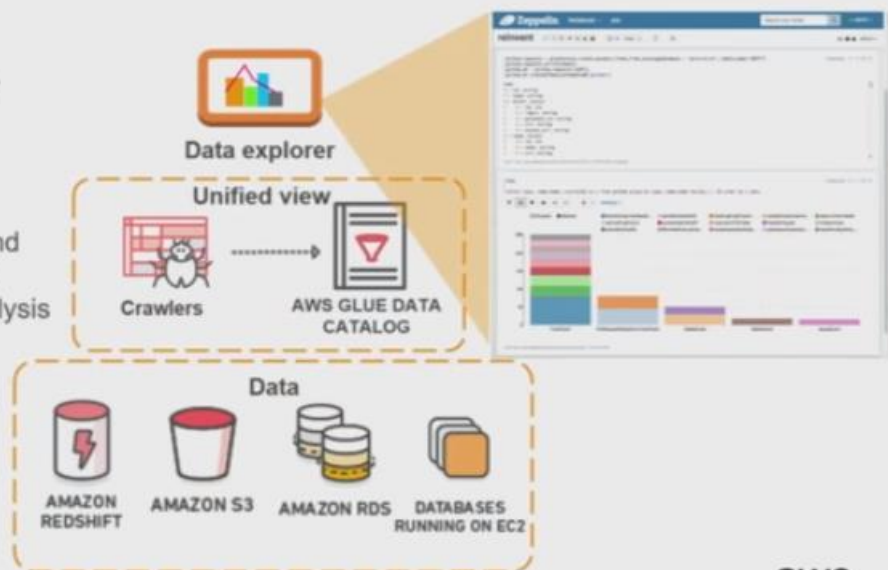
Pay only for the amount of data scanned

---

# Serverless data exploration

> Data scientists want fast access to disparate datasets for data exploration

> Glue automatically catalogues heterogeneous data sources, and offers serverless Apache Spark infrastructure for interactive analysis

> Gain insight in minutes without the need to configure and operationalize infrastructure



Data explorer

Unified view

Crawlers

AWS GLUE DATA CATALOG

Data

AMAZON REDSHIFT

AMAZON S3

AMAZON RDS

DATABASES RUNNING ON EC2

# Move data across storage systems

Data Stores

AMAZON REDSHIFT

AMAZON S3

AMAZON RDS

DATABASES RUNNING ON EC2

AWS GLUE CRAWLERS

**Unified view**

AWS GLUE DATA CATALOG

AWS GLUE ETL

AMAZON REDSHIFT

AMAZON S3

AMAZON RDS

DATABASES RUNNING ON EC2

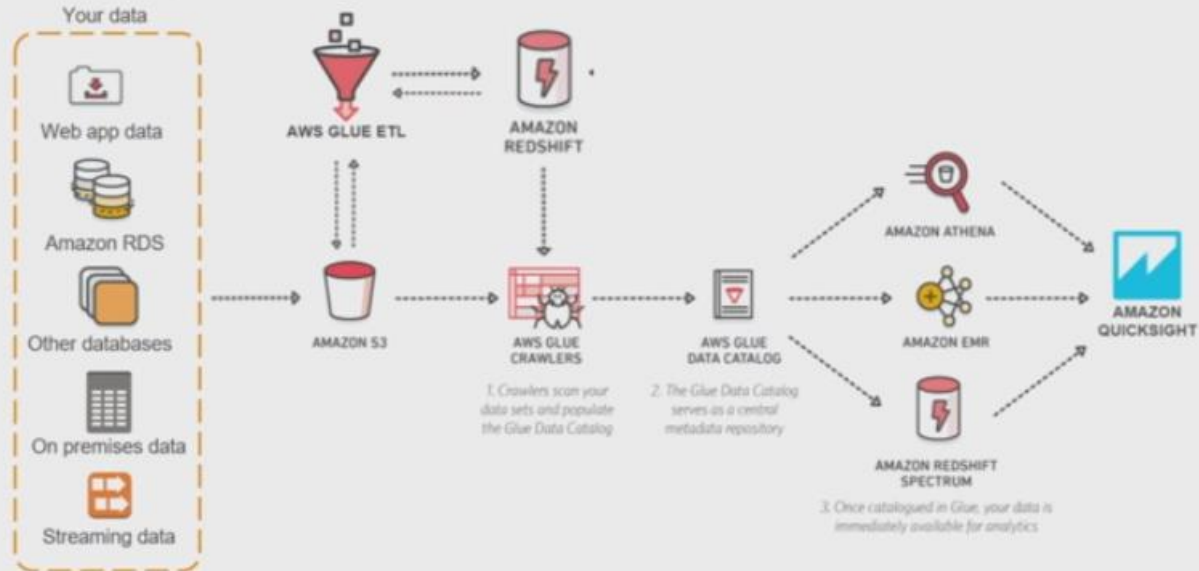# Data lake vs. data warehouse

| Data lake | Data warehouse |
|---|---|
| Semi-structured /Unstructured /structured data | Structured data |
| Schema on read | Schema on write |
| Data science, predictive analysis, BI use cases | SQL based BI use cases |
| Great for storing granular data; raw as well as processed data | Great for storing frequently accessed data as well as data aggregates and summary |
| Separation of compute and storage | Tightly coupled compute and storage |

# Interoperate data lake and data warehouse

# Key announcements (coming soon)

> Write Glue ETL jobs in **Scala**, in addition to PySpark

> Glue available in eu-west-1 (Ireland)

> Glue available in ap-northeast-1(Tokyo)