

AWS GLUE

Simple, flexible, cost-effective ETL

- AWS Glue is a fully managed ETL (extract, transform, and load) service
- Categorize your data, clean it, enrich it and move it reliably between various data stores
- Once catalogued, your data is immediately searchable and queryable across your data silos
- Simple and cost-effective
- Serverless; runs on a fully managed, auto-scaling Spark environment

Why would AWS get into the ETL space?

We have lots of ETL partners

Amazon Redshift Partner Page for Data Integration



but...Customers are still hand-coding ETL...

70% of ETL jobs are hand-coded

With no use of ETL tools

Actually...

It's over 90% in the cloud

Why do we see so much hand-coding?

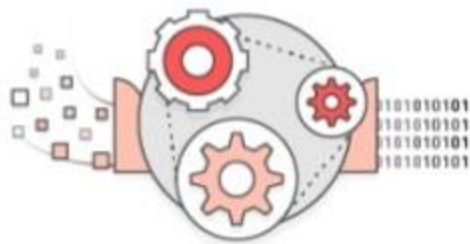
Code is flexible | Code is powerful

You can unit test | You can deploy with other code | You know your dev tools

Hand-coding involves a lot of undifferentiated heavy lifting...

Brittle | Error-prone | Laborious

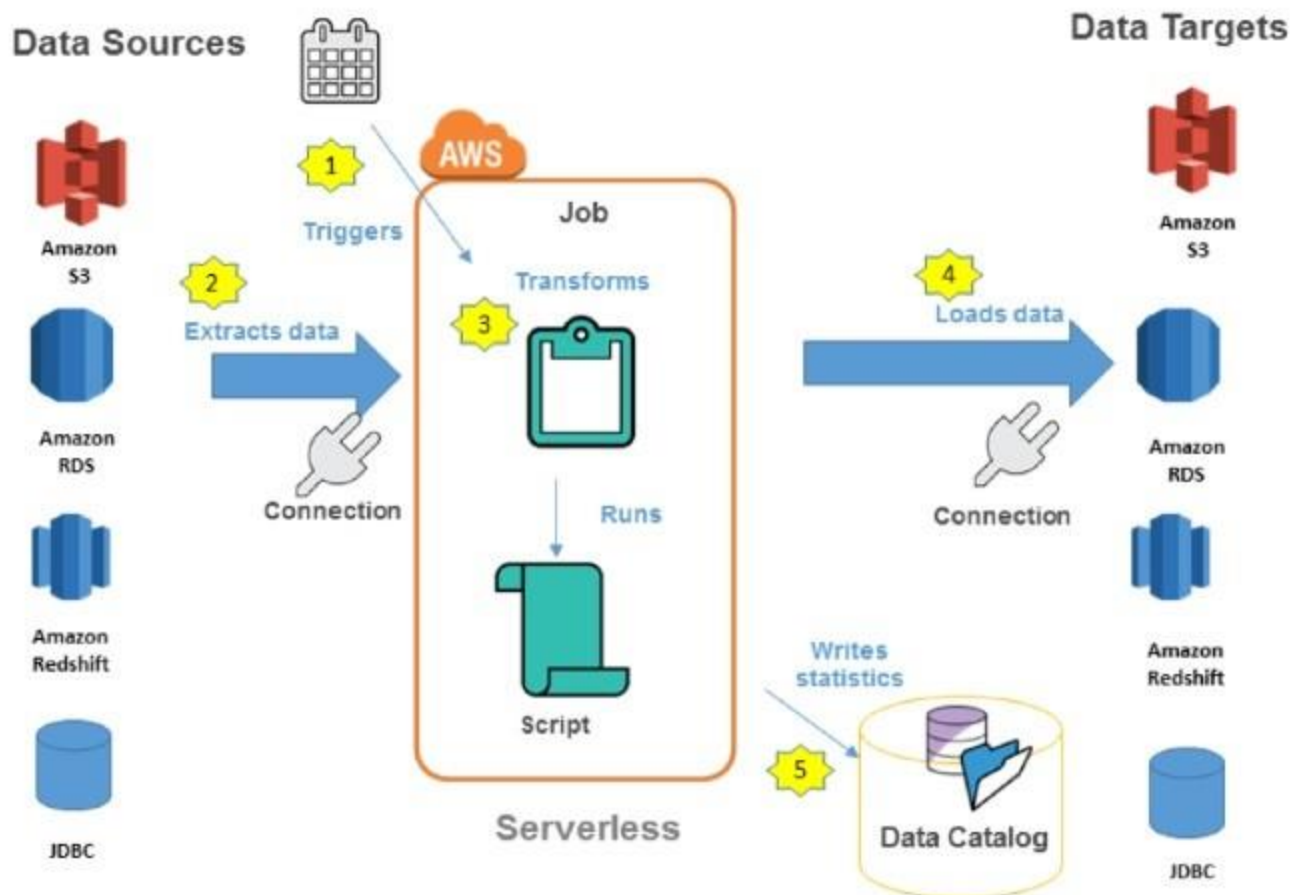
- ▶ As data formats change
- ▶ As target schemas change
- ▶ As you add sources
- ▶ As data volume grows



Glue automates the undifferentiated heavy-lifting of ETL

- ✓ Discover and organize data, regardless of where it lives
- ✓ Focus on writing transformations, not handling undifferentiating heavy lifting
- ✓ ETL jobs run under a Serverless execution model

AWS Glue: big picture



AWS Glue: components



Data Catalog

- Discover and Organize your data in various databases, data warehouses and data lakes



Job Authoring

- Focus on the writing transformations
 - Generate code through a wizard
 - Write your own code



Job Execution

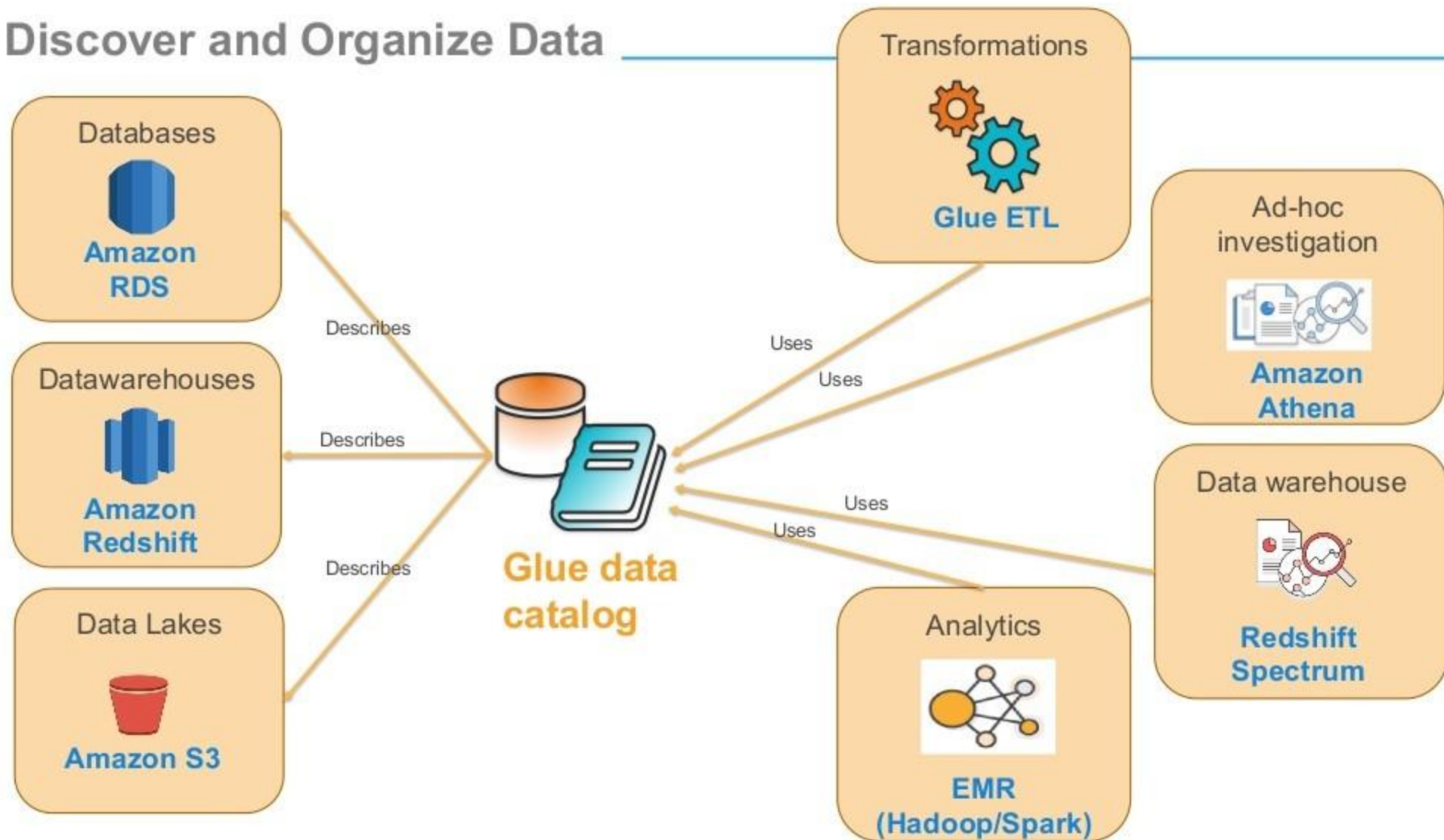
- Runs jobs in Spark containers – automatic scaling based on SLA
- Glue is serverless – only pay for the resources you consume



Glue data catalog

Discover and organize your data sets

Discover and Organize Data

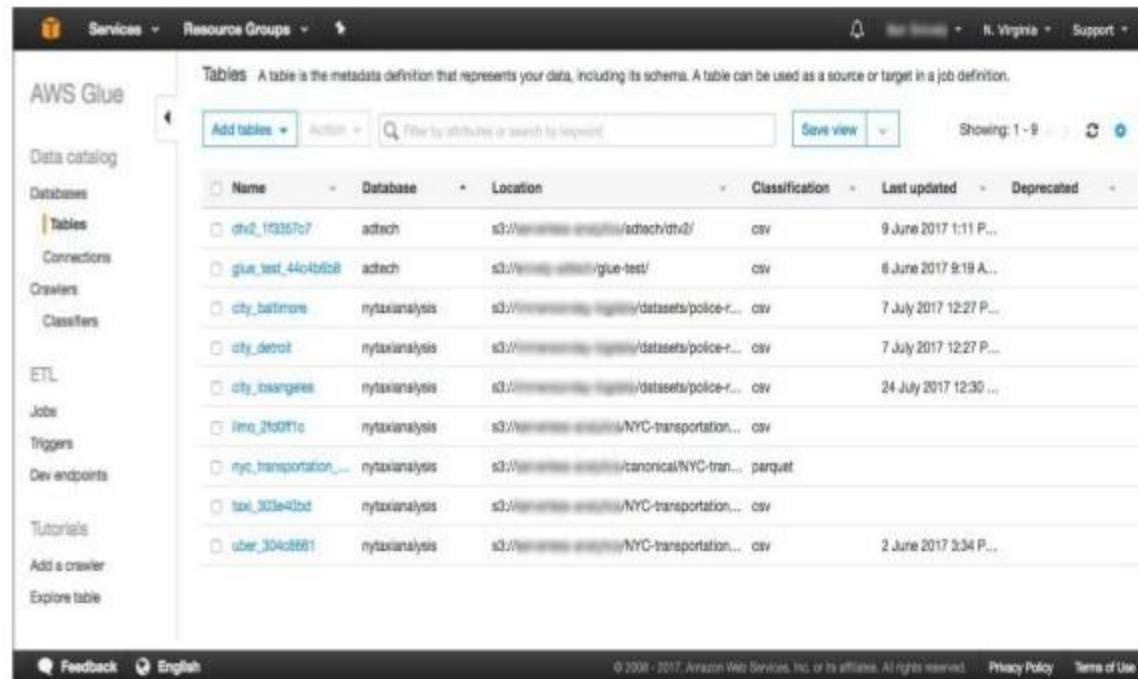


Glue data catalog

Manage table metadata through a Web Interface, Hive metastore API, Hive SQL, or automated through crawlers

Listening to our customers, we've added:

- **Search** metadata for data discovery
- **Connections** to RDS, Redshift and JDBC
- **Classification** for identifying and parsing files
- **Versioning** of table metadata as schemas



The screenshot shows the AWS Glue Data Catalog console. On the left is a navigation menu with options: Data catalog, Databases, Tables (selected), Connections, Crawlers, and Classifiers. The main area displays a list of tables with the following columns: Name, Database, Location, Classification, Last updated, and Deprecated. The table list includes entries like 'dtv2_1f3357c7', 'glue_test_44c4b6c8', and 'city_baltimore'. Above the table list, there is a header section with 'Add tables', a search bar, and a 'Save view' button. Below the table list, there is a footer section with 'Feedback', 'English', and copyright information.

| Name | Database | Location | Classification | Last updated | Deprecated |
|-----------------------|----------------|---|----------------|------------------------|------------|
| dtv2_1f3357c7 | adtech | s3://aws-logs-us-east-1-adtech/dtv2/ | csv | 9 June 2017 1:11 P... | |
| glue_test_44c4b6c8 | adtech | s3://aws-logs-us-east-1-adtech/glue-test/ | csv | 6 June 2017 9:19 A... | |
| city_baltimore | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/datasets/policy-r... | csv | 7 July 2017 12:27 P... | |
| city_detroit | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/datasets/policy-r... | csv | 7 July 2017 12:27 P... | |
| city_houston | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/datasets/policy-r... | csv | 24 July 2017 12:30 ... | |
| lmo_2b50f1c | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/nytaxi-transportation... | csv | | |
| nyc_transportation... | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/canonical/nytaxi-transportation... | parquet | | |
| taxi_303e41bd | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/nytaxi-transportation... | csv | | |
| uber_304c8881 | nytaxianalysis | s3://aws-logs-us-east-1-nytaxi/nytaxi-transportation... | csv | 2 June 2017 3:34 P... | |

Crawlers: auto-populate data catalog

Run crawlers on-demand and on a schedule to **discover** new data and **schema changes**.

Serverless – only pay when crawls run.

Automatic schema inference:

- Built-in classifiers
 - **Detect file type**
 - **Extract schema**
 - **Identify partitions**
- Add your own classifiers
 - Grok for each of use

Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

[Add crawler](#) [Run crawler](#) [Actions](#)

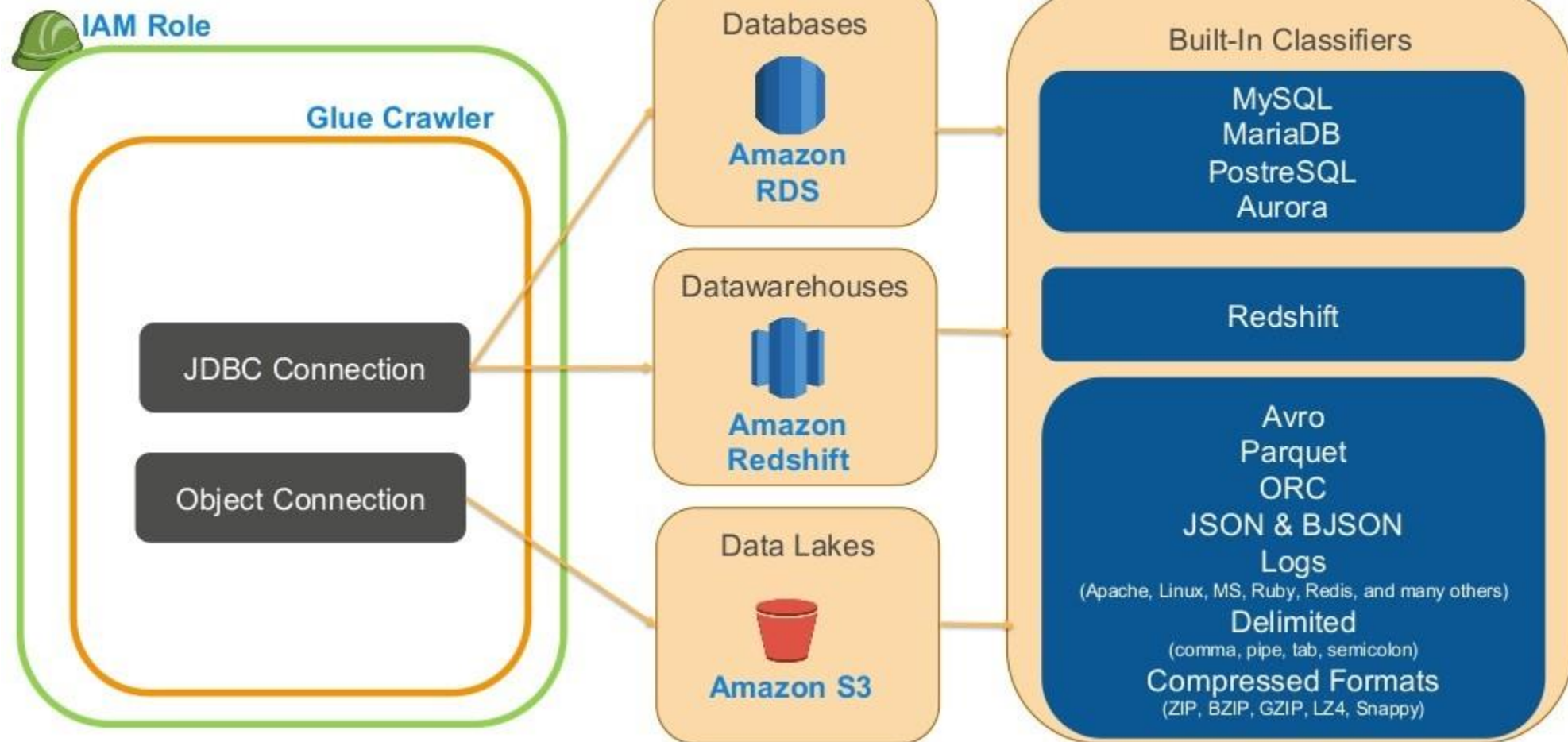
Showing: 1 - 6

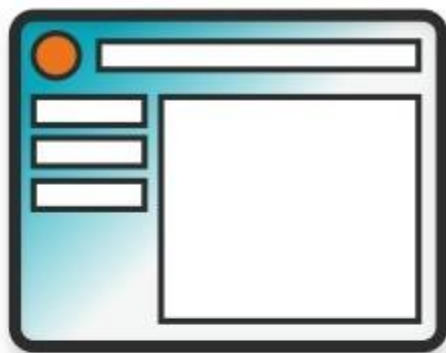
| <input type="checkbox"/> | Name | Schedule | Status | Logs | Last runtime | Median runtime | Tables updated | Tables added |
|--------------------------|-------------------|----------|--------|----------------------|--------------|----------------|----------------|--------------|
| <input type="checkbox"/> | DoubleClickCra... | | Ready | Logs | 15 secs | 15 secs | 0 | 1 |
| <input type="checkbox"/> | FlightCrawling | | Ready | Logs | 20 secs | 20 secs | 0 | 1 |
| <input type="checkbox"/> | NYC-Taxi-Craw... | | Ready | Logs | 13 secs | 13 secs | 0 | 4 |
| <input type="checkbox"/> | Uber | | Ready | Logs | 16 secs | 16 secs | 0 | 1 |
| <input type="checkbox"/> | adTechCrawler | | Ready | Logs | 20 secs | 20 secs | 0 | 1 |
| <input type="checkbox"/> | policeCrawler | | Ready | Logs | 15 secs | 15 secs | 1 | 0 |

© 2008 - 2017, Amazon Web Services, Inc. or its affiliates. All rights reserved. [Privacy Policy](#) [Terms of Use](#)

Crawlers: Classifiers

Create additional Custom Classifiers with Grok!



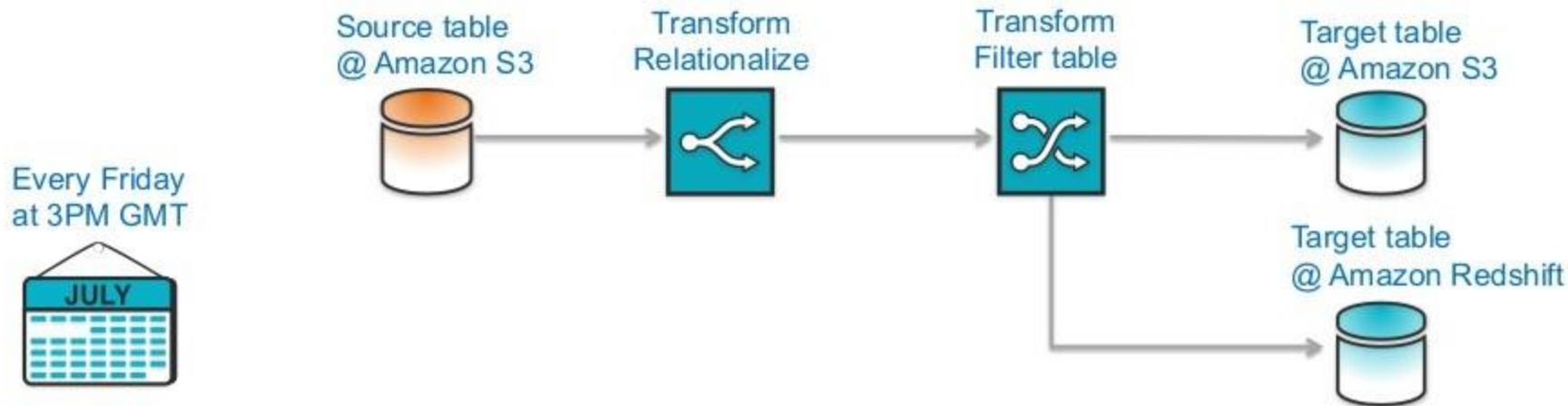


Job authoring in Glue

**You have
choices on how
to get started...**

Script generated by AWS Glue
Existing script brought into AWS Glue
Blank script authored by you

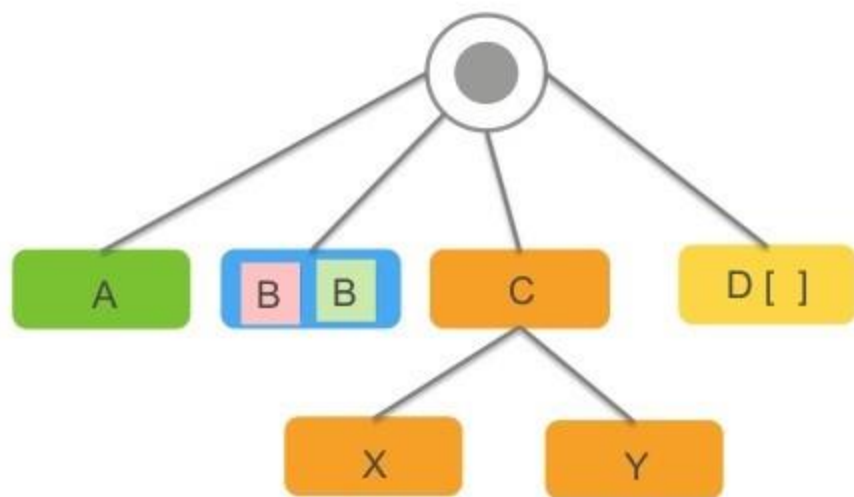
Automatic Code Generation



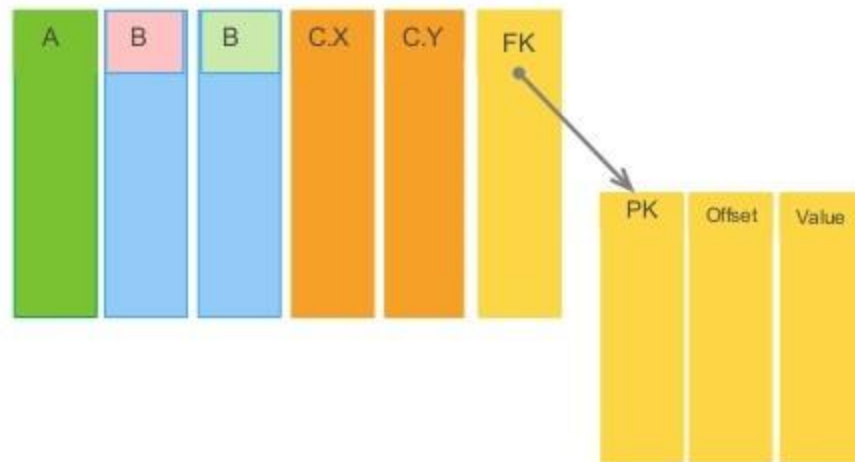
1. Pick sources and targets from the data catalog
2. Glue generates transformation graph and *Python code*
3. Specify trigger condition and execute job

Glue transformations are flexible and adaptive

Semi-structured schema



Relational schema



- Flatten semi-structured objects with arbitrary complexity into relational tables, on-the-fly.
- Pivot arrays and other collection types into a separate table, generating key-foreign key values.
- Modify mapping as the source schema changes, and modify the target schemas as needed.

Glue ETL scripts are forgiving and flexible

The screenshot displays the AWS Glue console interface. On the left, a visual ETL workflow is shown: a source database 'nytaxianalysis' with table 'city_baltimore' is transformed by 'ApplyMapping' and 'SelectFields' into a target database 'incidents' with table 'canonical'. The main area shows the PySpark script for the job, which includes imports for Glue utilities, SparkContext, and GlueContext, followed by job configuration and execution logic. The script is as follows:

```
1 | import sys
2 | from awsglue.transforms import *
3 | from awsglue.utils import getResolvedOptions
4 | from pyspark.context import SparkContext
5 | from awsglue.context import GlueContext
6 | from awsglue.job import Job
7 |
8 | ## @param: [JOB_NAME]
9 | args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10 |
11 | sc = SparkContext()
12 | glueContext = GlueContext(sc)
13 | job = Job(glueContext)
14 | job.init(args['JOB_NAME'], args)
15 |
16 | ## @type: DataSource
17 | ## @args: [database = 'nytaxianalysis', table_name = 'city_baltimore', transformation_ctx = 'datasource0']
18 | ## @return: DataSource
19 | ## @inputs: []
```

Below the script, there are tabs for 'Logs', 'Schema', and 'Statistics'. The bottom of the console shows a footer with 'Feedback', 'English', and copyright information: '© 2009 - 2017, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use'.

- **Human-readable** code run on a scalable platform, PySpark
- **Forgiving** in the face of failures – handles bad data and crashes
- **Flexible**: handles complex semi-structured data, and adapts to source schema changes

Add Custom Modules and Files

- Add external Python modules
- Java JARs required by the script
- Additional files such as configuration, etc.

Parameters (optional)

▸ Advanced properties

▼ Script libraries and job parameters (optional)

☐ Server-side encryption

Python library path

s3://bucket-name/folder-name/file-name

Dependent jars path

s3://bucket-name/folder-name/file-name

Referenced files path

s3://bucket-name/folder-name/file-name



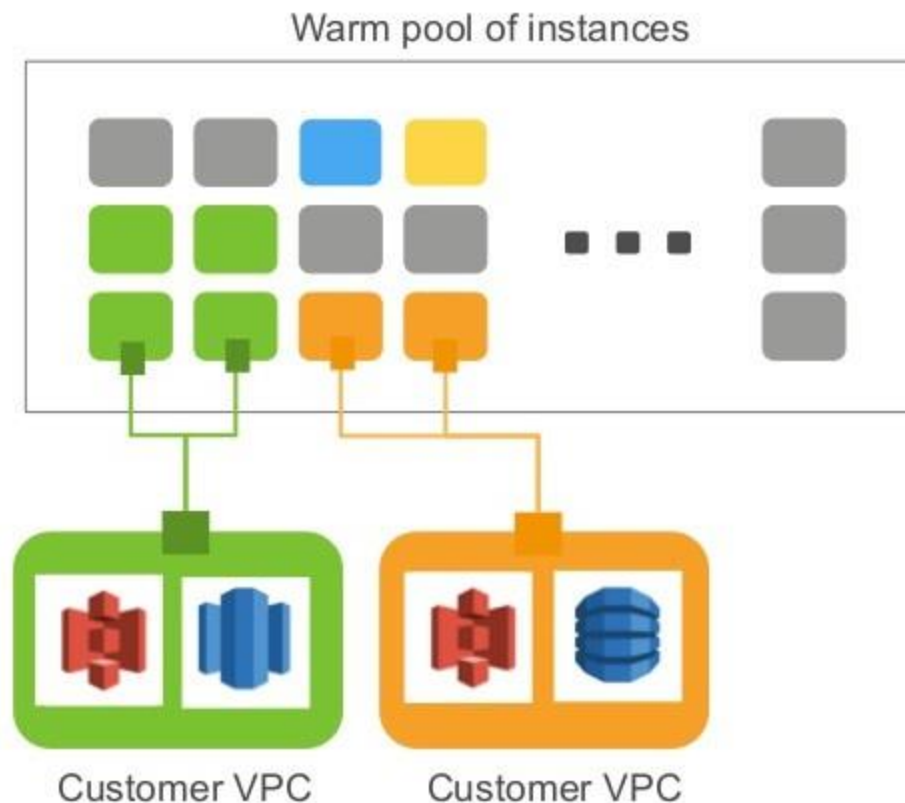
Orchestration & resource management

Fully managed, serverless job execution

Serverless job execution

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

- Warm pools: pre-configured fleets of instances to reduce job startup time
- Auto-configure VPC and role-based access
- Automatically scale resources to meet SLA and cost objectives
- You pay only for the resources you consume while consuming them.



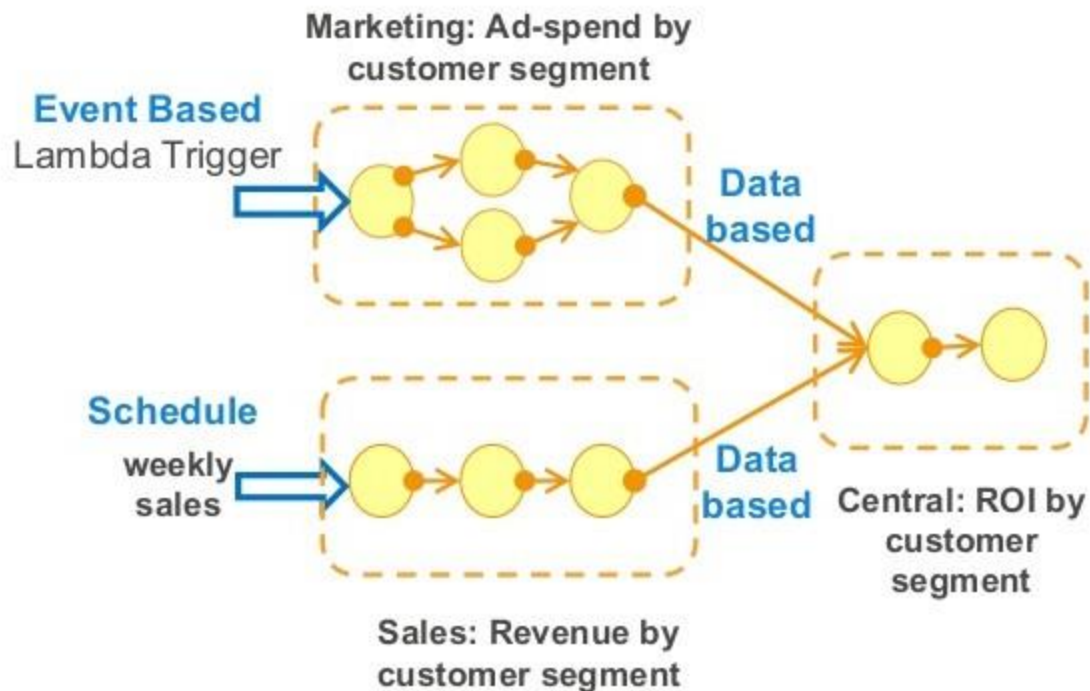
Job composition and triggers

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., data availability, job completion
- **External sources:** e.g., AWS Lambda

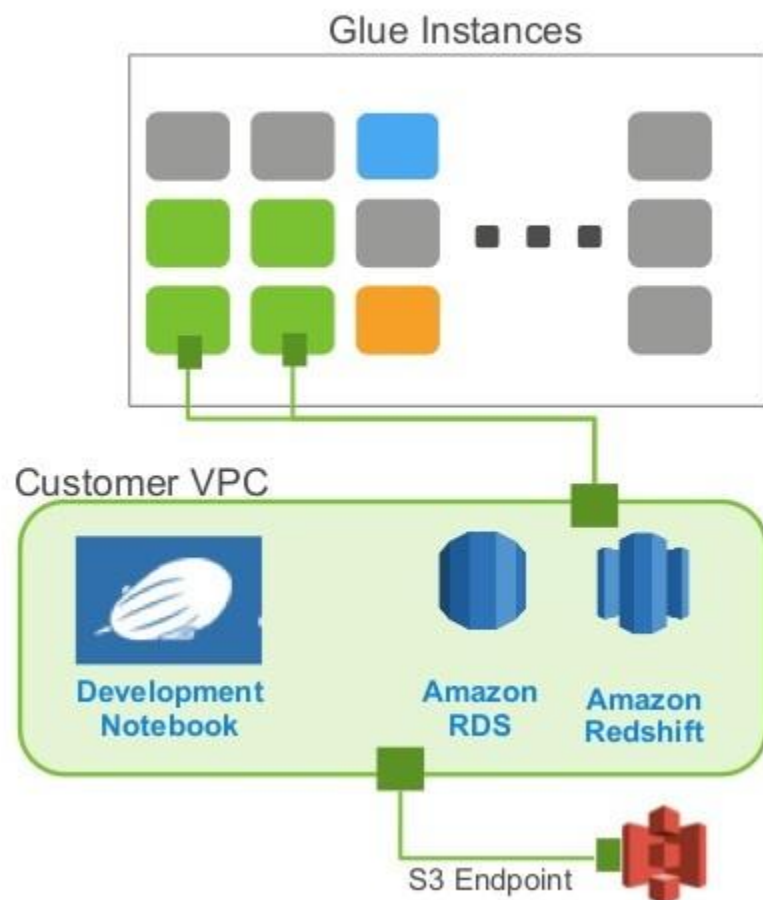


Developer Endpoints

Environment to iteratively develop and test ETL scripts.

Develop your script in a notebook and point to an AWS Glue endpoint to test it.

When you are satisfied with the results of your development process you can create an ETL job that runs your script.



Glue ETL Security Model

- ETL Jobs that do not require special handling can be launched without additional configuration
- For jobs requiring restricted access to data Glue utilizes a VPC endpoint launching dedicated ENIs into the customer's VPC which are assigned private IP address from the customer's subnet.
- For such jobs, Glue also enforces the use of an S3 VPC endpoint

Usage and Pricing

Glue ETL Pricing

- With Glue, you only pay for the time your ETL job takes to run.
 - There are no resources to manage and no upfront costs, and you are not charged for startup or shutdown time.
- You are charged an hourly rate based on the number of Data Processing Units (or DPUs) used to run your ETL job.
 - A single Data Processing Unit (DPU) provides 4 vCPU and 16 GB of memory and corresponding networking capabilities.
- You are billed per hour in increments of 1 minute, rounded up to the nearest minute, with a 10-minute minimum duration for each job.

Glue Data Catalog and Crawler Pricing

Data catalog:

- With the AWS Glue data catalog, you can store up to a million objects per month for free. If you store more than a million objects, you will be charged per 100,000 objects over a million.
 - An object in the AWS Glue data catalog is a table, a partition, or a database.
- The first million access requests per month to the AWS Glue data catalog are free. If you exceed a million requests in a month, you will be charged per million requests over the first million.
 - Some common requests are CreateTable, CreatePartition, GetTable and GetPartitions.

Crawlers:

- You will pay an hourly rate for AWS Glue crawler runtime to populate the Glue data catalog, based on the number of Data Processing Units (or DPUs) used to run your crawler.
 - A single Data Processing Unit (DPU) provides 4 vCPU and 16 GB of memory and corresponding networking capabilities.
 - You are billed in increments of 1 minute, rounded up to the nearest minute.
- Use of AWS Glue crawlers is optional, and you can populate the Glue data catalog directly through the API.