



# Primeros Pasos con AWS Glue

👤 Created by	👤 Juan Alejandro Carrillo Jaimes
🕒 Created time	@August 9, 2025 11:07 AM
☰ Tags	<span>AWS</span> <span>AWS Glue</span> <span>Cloud</span> <span>ETL</span> <span>Fundamentals</span>

Configuraciones y Ejercicios

[Configuración de permisos de IAM para AWS Glue](#)

[Simplifique la administración de canalizaciones ETL](#)

[Explore, experimente y procese datos de forma interactiva](#)

[Descubra datos de manera eficiente](#)

[Soporta varios marcos de procesamiento y cargas de trabajo](#)

[Beneficios de AWS Glue](#)

[Data Integration Options](#)

[Opciones de Integración de Datos en AWS Glue](#)

[Interfaces para crear y ejecutar estos jobs:](#)

[Orígenes de datos:](#)

[Destinos típicos:](#)

[Data Preparation](#)

[¿Qué es AWS Glue DataBrew?](#)

[Funcionalidades Clave:](#)

Características Principales de AWS Glue

[Descubrir y Organizar datos](#)

[Transformar, preparar y limpiar los dato para su análisis](#)

[Crear y supervisar pipelines de datos](#)

[Diseño de Pipelines Complejos](#)

[Ejecución y Automatización](#)

[Monitoreo y Auditoría](#)

[Trazabilidad y Mantenimiento](#)

Conceptos de AWS Glue

[Flujo de AWS Glue explicado paso a paso](#)

1. [Data Stores](#)

2. [Crawler](#)

3. [AWS Glue Data Catalog](#)

4. [AWS Management Console](#)

5. [Job](#)

6. **Script (Transformación)**
7. **Trigger / Schedule or Event**
8. **Transform → Load**

Términos clave conectados al flujo:  
[Sources](#)

## Configuraciones y Ejercicios



### Configuración de permisos de IAM para AWS Glue



La siguiente configuración de AWS IAM, está diseñada para otorgar acceso a los recursos de AWS Glue y ejecutar tareas en AWS Glue Data Quality.

1. Ir a la AWS Glue Console Manager, y dar clic en "Prepare your account for AWS Glue"

The screenshot shows the AWS Glue Getting Started page. It features three main sections:

- Prepare your account for AWS Glue**: Includes a shield icon, a description about granting access to AWS Glue and setting a default IAM role, and a button to "Set up roles and users".
- Catalog and search for datasets**: Includes a magnifying glass icon, a description about viewing databases & tables and cataloging data using Crawlers, and a button to "Go to the Data Catalog".
- Move and transform data**: Includes a circular icon with arrows, a description about using Zero-ETL integrations to replicate data, and buttons to "Go to Zero-ETL integrations" and "Go to ETL jobs".

2. Seleccionamos el rol por defecto de Glue, su configuración de permisos IAM depende directamente del caso de uso que se necesite.

The screenshot shows the "Choose IAM users and roles for AWS Glue" section. It includes a "How it works" box explaining that IAM identities (roles or users) need access to AWS Glue, and a "Selected roles" list.

**Selected roles (1)**

Role name	Creation date (UTC)
AWSGlueServiceRole	October 17, 2024 at 15:43:36

- Elegimos un bucket, o la opcion de todos los buckets, se recomienda manejar un permiso sobre un bucket, para tener mayor control de dependencias.

**Selected S3 locations (1)**

Choose the list of S3 locations to be accessible by the Glue users and roles.

Name	Creation date (UTC)
gft-sources-us-east-1	August 9, 2025 at 16:33:15

**Data access permissions**

Set the type of data access for the Glue users and roles.

**Data access permissions**

- Read only (recommended)
- Read and write

Cancel Previous Next

- Elegimos por defecto el AWS Glue Service Role, para realizar la tarea.

Step 1 Choose IAM users and roles for AWS Glue

Step 2 Grant Amazon S3 access

Step 3 **Choose a default service role**

Step 4 Review and confirm

**Choose a default service role**

**How it works**

AWS Glue uses an IAM service role to run jobs, access data, and run Data Quality tasks. We recommend that you start with the standard AWSGlueServiceRole as the default. [Learn more](#)

**Choose a default AWS Glue service role**

**IAM role for AWS Glue**

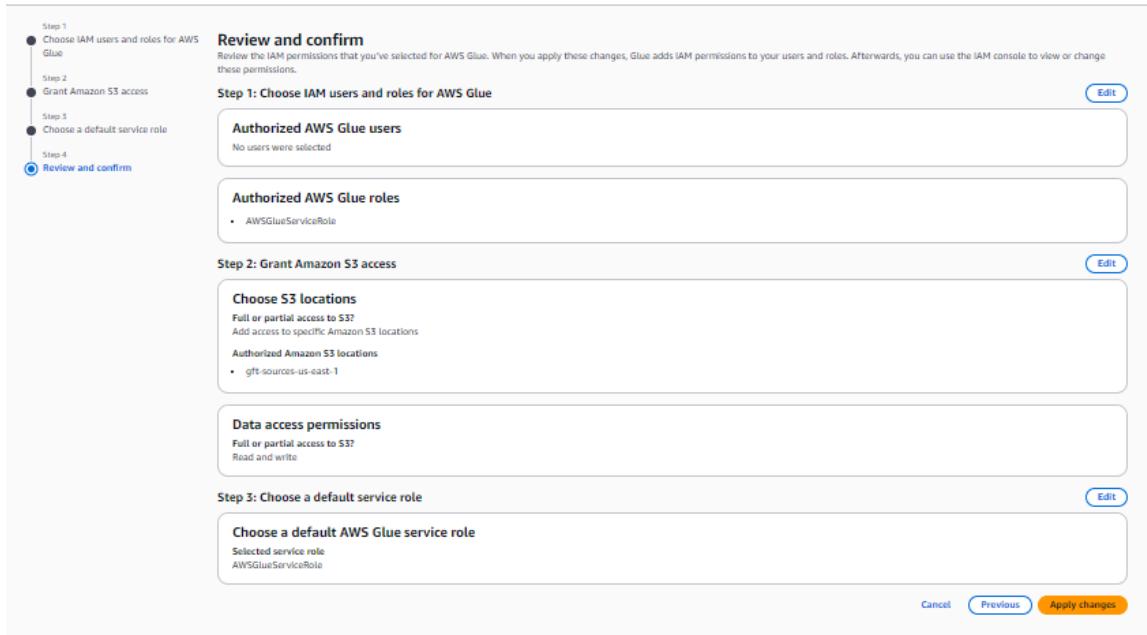
- Update the standard AWS Glue service role and set it as the default (recommended)  
AWS will update the role with the IAM policies needed to run AWS Glue jobs, then set it as the default.
- Set an existing IAM role as the default  
Select an IAM role that you've configured to use as an AWS Glue service role. Glue will set this role as the default, but won't add any permissions to it. [Learn more](#)

The following IAM role will be created and automatically configured for you:

- AWSGlueServiceRole

Cancel Previous Next

- Revisamos y confirmamos.



## Simplifique la administración de canalizaciones ETL

Elimine la administración de infraestructura con aprovisionamiento automático y administración de trabajadores, y consolide todas sus necesidades de integración de datos en un único servicio.

## Explore, experimente y procese datos de forma interactiva

Mediante las sesiones interactivas de AWS Glue, los ingenieros de datos pueden explorar y preparar datos de forma interactiva utilizando el entorno de desarrollo integrado (IDE) o el bloc de notas de su elección.

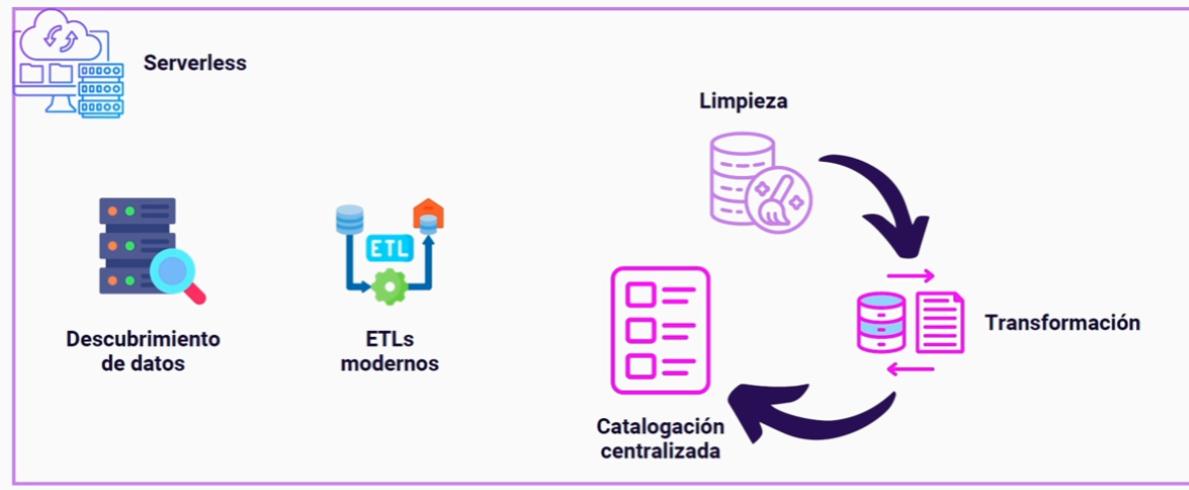
## Descubra datos de manera eficiente

Identifique datos rápidamente en AWS, en las instalaciones y en otras nubes y, a continuación, póngalos a disposición al instante para realizar consultas y transformaciones.

## Soporta varios marcos de procesamiento y cargas de trabajo

Soporta más fácilmente varios marcos de procesamiento de datos, como ETL y ELT, y varias cargas de trabajo, incluidos lotes, microlotes y streaming.

## Beneficios de AWS Glue

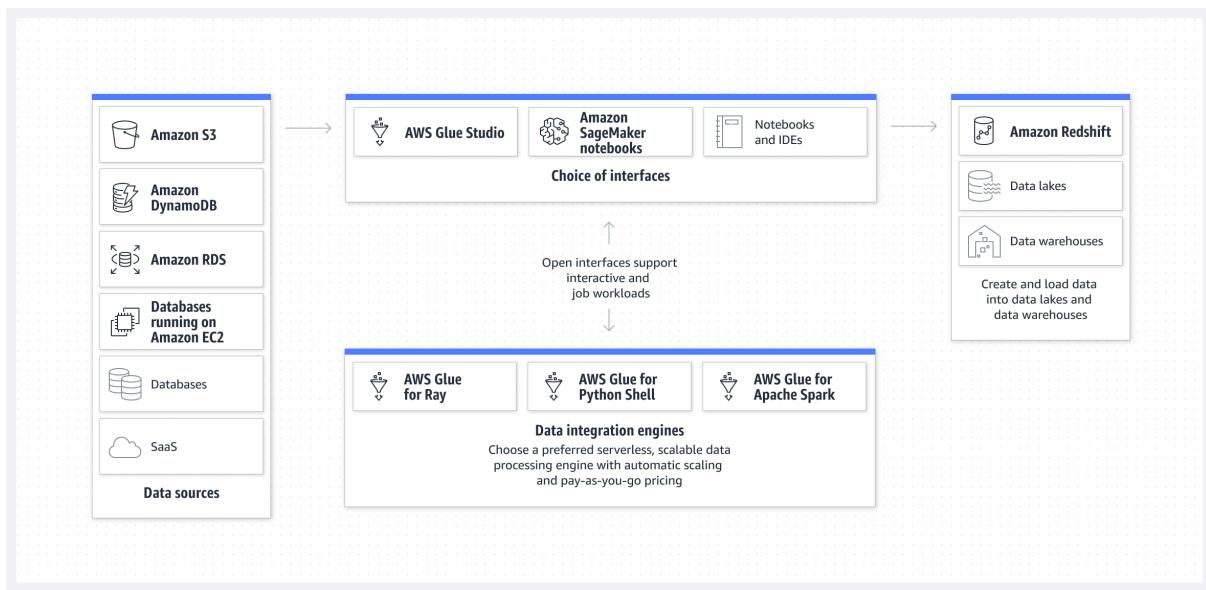


- Compatible con todas las cargas de trabajo:** Soporte flexible para ETL, ELT, batch, streaming y más, sin lock-in.
- Escala bajo demanda:** Escala de petabytes, facturación por uso y tamaño de datos.
- Herramientas a medida:** Soporte a todos los usuarios de datos, desde desarrolladores a usuarios de negocio.
- Asistente de IA generativa:** Obtenga ayuda impulsada por IA a lo largo de su viaje de integración de datos, desde la generación automática de código ETL hasta la modernización de sus trabajos de Spark. AWS Glue proporciona generación de código inteligente, actualizaciones de Spark asistidas por IA (vista previa) y solución de problemas de Spark integrada.
- Todo en uno:** Capacidades completas de integración de datos en un servicio sin servidor.



**Serverless:** No hay infraestructura para administrar

## Data Integration Options



## 🔧 Opciones de Integración de Datos en AWS Glue

### 1. AWS Glue for Apache Spark

- Ideal para cargas de trabajo grandes y transformaciones complejas.
- Procesamiento distribuido y escalable.
- Soporta `DynamicFrame`, PySpark y particionamiento automático.

### 2. AWS Glue for Python Shell

- Usado para tareas ligeras o scripts rápidos.
- Ideal para validación de archivos, llamadas a APIs, procesamiento de metadatos.
- Usa solo Python estándar (no Spark).

### 3. AWS Glue for Ray

- Motor más nuevo orientado a procesamiento paralelo con Python puro.
- Bueno para tareas que no requieren Spark pero sí paralelismo.
- Compatible con librerías como Pandas, NumPy y Dask.



## Interfaces para crear y ejecutar estos jobs:

- **AWS Glue Studio:** interfaz gráfica para diseñar y correr ETLs.
- **Amazon SageMaker Notebooks:** desarrollo en notebooks Jupyter.
- **Notebooks e IDEs:** integración con VS Code, Jupyter local, etc.



## Orígenes de datos:

- Amazon S3, DynamoDB, RDS, EC2, SaaS, etc.



## Destinos típicos:

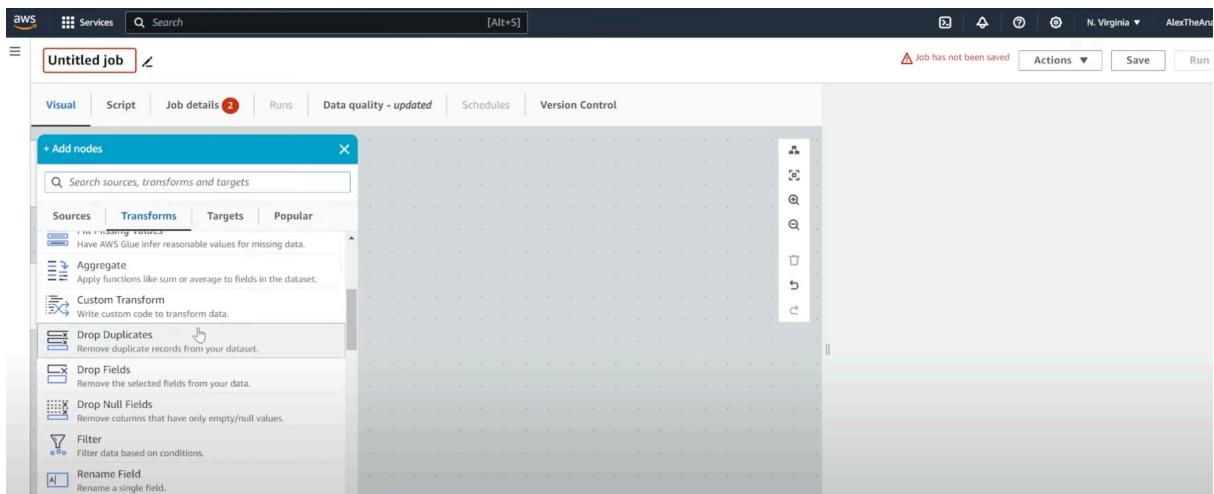
- Amazon Redshift (data warehouse), Data Lakes.

# Data Preparation



Con [AWS Glue DataBrew](#), puede explorar y experimentar con datos directamente desde su lago de datos, almacenes de datos y bases de datos, incluidos Amazon S3, Amazon Redshift, AWS Lake Formation, Amazon Aurora y Amazon Relational Database Service (RDS). Puede elegir entre más de 250 transformaciones predefinidas en DataBrew para automatizar tareas de preparación de datos como filtrar anomalías, estandarizar formatos y corregir valores no válidos.

## ¿Qué es AWS Glue DataBrew?



**AWS Glue DataBrew** es una herramienta visual de preparación de datos que permite a analistas y científicos de datos limpiar y transformar datos **sin escribir código**. Está diseñada para trabajar directamente con fuentes de datos como Amazon S3, Redshift, RDS, Aurora y más.

## Funcionalidades Clave:

### 1. Conectividad con fuentes de datos

- Se conecta directamente a data lakes, data warehouses y bases de datos.

### 2. +250 transformaciones predefinidas

- Puedes limpiar, normalizar, transformar y enriquecer datos sin necesidad de programar.

### 3. Evaluación de calidad de datos

- Perfilado de datos para detectar patrones, valores nulos, outliers y anomalías.

### 4. Automatización escalable

- Reutiliza pasos de transformación automáticamente cuando llegan nuevos datos.

### 5. Publicación en Amazon S3

- Guarda los datos transformados directamente en S3, listos para usarse en análisis o ML.

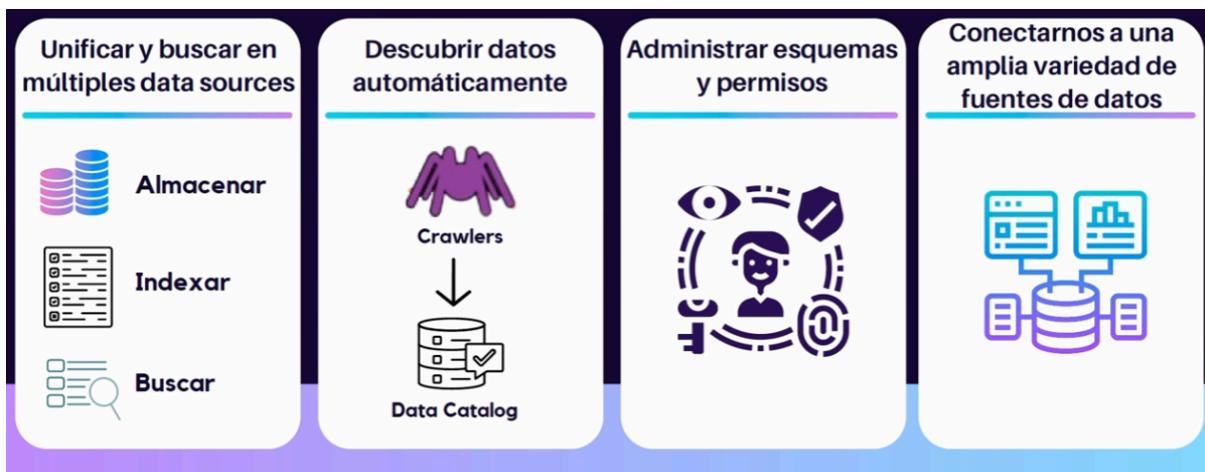
### 6. Obtención rápida de insights

- Datos listos para alimentar dashboards, sistemas de BI o modelos de Machine Learning.

## Características Principales de AWS Glue



### Descubrir y Organizar datos



AWS Glue facilita la **unificación y organización de datos dispersos** provenientes de diversas fuentes tanto en la nube como on-premise. Este proceso se apoya en varias capacidades clave:

- **🔍 Descubrimiento Automático con Crawlers:** Los *crawlers* examinan automáticamente tus datos, infieren los esquemas y tipos de datos, y los registran en el **AWS Glue Data Catalog**. Esto elimina la necesidad de definir manualmente estructuras complejas.
- **🗃 Catálogo Centralizado de Datos (Glue Data Catalog):** Funciona como un metastore que organiza tus datasets en bases de datos y tablas accesibles para herramientas como Amazon Athena, Redshift Spectrum y EMR.
- **🔒 Gestión de Accesos y Permisos:** Se integra con **AWS Lake Formation** y **IAM** para establecer permisos granulares a nivel de base de datos, tabla, columna o fila.
- **🔗 Conectividad Multifuente:** Glue puede conectarse a datos almacenados en Amazon S3, Amazon RDS, Redshift, DynamoDB, y también a fuentes

externas mediante conectores JDBC y conectores personalizados para SaaS o entornos on-premise.

-  **Soporte para Esquemas Evolutivos:** Permite gestionar cambios en los esquemas a lo largo del tiempo, versión tras versión, sin romper los pipelines existentes.

## Transformar, preparar y limpiar los dato para su análisis



AWS GLUE

### Transformar, preparar y limpiar los datos para su análisis



- 01 Transforme visualmente los datos con una interfaz de arrastrar y soltar (drag-and-drop)
- 02 Cree piplines de ETLs complejos con una programación de jobs simple
- 03 Limpie y transforme los datos en streaming
- 04 Deduplicar y limpiar datos con Machine Learning
- 05 Crear notebooks de trabajo integrados
- 06 Editar, depurar y probar el código de los ETL
- 07 Definir, detectar y corregir datos confidenciales

AWS Glue permite crear flujos de transformación y pipelines ETL de manera intuitiva y potente, con las siguientes capacidades:

-  **Interfaz Visual Drag-and-Drop (AWS Glue Studio):**

Diseña flujos ETL complejos sin necesidad de escribir código. Ideal para usuarios técnicos y no técnicos que desean transformar datos desde diversas fuentes como Amazon S3, RDS, Redshift, DynamoDB y más.

-  **Orquestación Flexible de Pipelines de Datos:**

Crea pipelines escalables con ejecuciones programadas (basadas en calendario), bajo demanda o disparadas por eventos (por ejemplo, carga de archivos en S3).

-  **Procesamiento en Streaming:**

Transforma datos en tiempo real a medida que llegan, lo que permite casos de uso como monitoreo de eventos, análisis en vivo y limpieza sobre datos en tránsito.

-  **Limpieza de Datos con IA y Machine Learning (ML):**

Utiliza funcionalidades avanzadas como **FindMatches ML Transform** para realizar deduplicación inteligente y limpieza automática con algoritmos de aprendizaje automático.

-  **Notebooks Serverless Integrados:**

Experimenta, depura y ajusta tus scripts ETL en entornos interactivos como Jupyter Notebooks sin tener que aprovisionar servidores, gracias a la integración con Amazon SageMaker y Glue Notebooks.

-  **Edición y Depuración en IDEs o Notebooks:**

Desarrolla tus ETLs directamente en entornos conocidos como Visual Studio Code, notebooks o Glue Studio, con ejecución de secciones interactivas que facilitan la depuración paso a paso.

-  **Detección y Protección de Datos Sensibles:**

Identifica y remedia automáticamente datos confidenciales (como PII) en tus pipelines y lagos de datos, ayudando a cumplir normativas como GDPR, HIPAA o CCPA.

## Crear y supervisar pipelines de datos

# Crear y supervisar pipelines de datos

Escalar automáticamente en función de la carga de trabajo



Ejecutar y monitorear jobs



Automatizar jobs con activadores basados en eventos



Definir flujos de jobs para ETLs y actividades de integración



AWS Glue permite **construir y orquestar pipelines de datos escalables y totalmente administrados**, diseñados para manejar desde simples transformaciones hasta flujos de datos complejos y dependientes. A continuación se detallan las capacidades clave:

## Diseño de Pipelines Complejos

- Puedes **encadenar múltiples jobs, crawlers y triggers**, creando flujos de trabajo ETL bien definidos y modulares.
- Los *triggers* permiten iniciar procesos basados en eventos específicos (por ejemplo, carga de archivos en S3) o en horarios programados (*cron*).
- Los flujos de trabajo (*Workflows*) permiten visualizar la ejecución completa del pipeline en forma de grafo.

## Ejecución y Automatización

- AWS Glue **escala automáticamente los recursos subyacentes** en función del tamaño del dataset y complejidad del job, eliminando la necesidad de configurar clústeres manualmente.
- Puedes automatizar la ejecución de jobs mediante **triggers condicionales**, activando procesos de manera inteligente según el éxito, fallo o finalización de tareas previas.

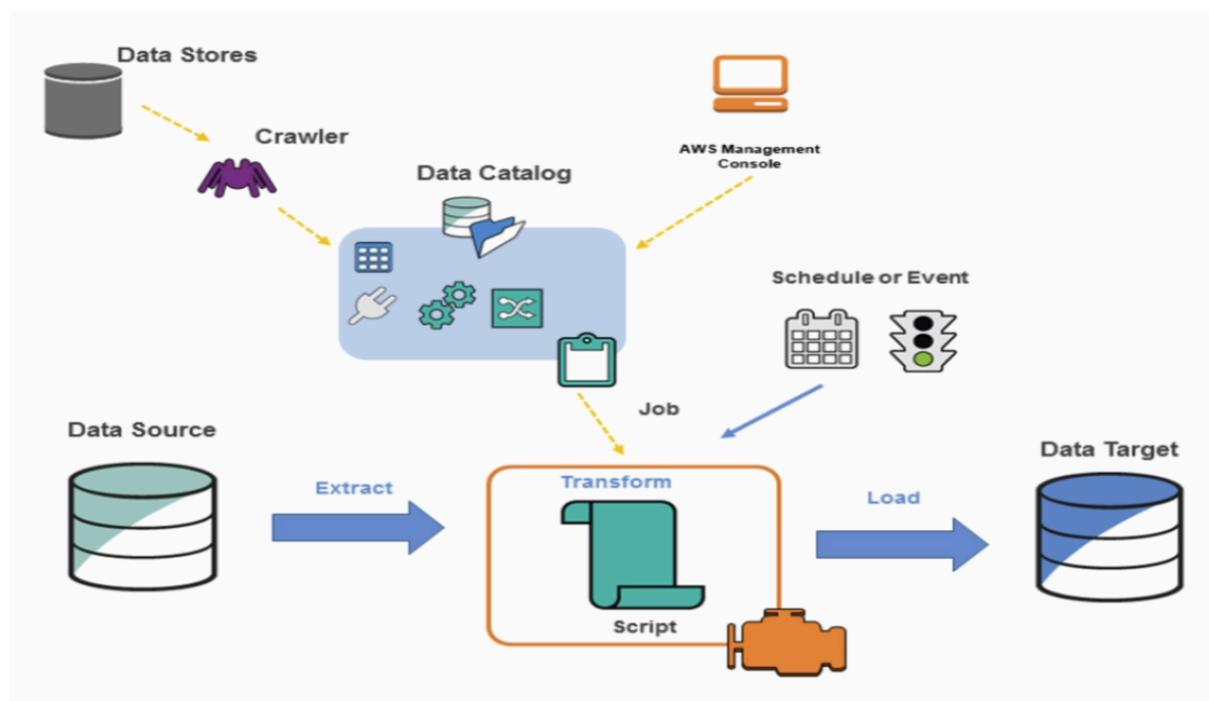
## Monitoreo y Auditoría

- La ejecución de jobs puede **monitorearse en tiempo real** desde AWS Glue Studio o mediante APIs.
- Integración con **AWS CloudTrail** para auditar todas las acciones realizadas en Glue: ejecución de jobs, activación de triggers, uso de crawlers, etc.
- Glue registra métricas clave en **Amazon CloudWatch**, permitiendo configurar alertas o dashboards para detectar cuellos de botella, errores o tareas fallidas.

## Trazabilidad y Mantenimiento

- Cada componente del pipeline (jobs, crawlers, triggers) está versionado y puede gestionarse de forma independiente.
- Las dependencias y relaciones entre tareas se gestionan desde **Glue Workflows**, lo que facilita el mantenimiento y la evolución del sistema ETL.

# Conceptos de AWS Glue



## Flujo de AWS Glue explicado paso a paso

## 1. Data Stores

- Representan las **fuentes de datos** y **destinos** como buckets S3, bases de datos relacionales (RDS, Aurora, etc.), NoSQL, etc.
  - Estos datos pueden estar **estructurados, semi-estructurados o no estructurados**.
  - Se conectan mediante **AWS Glue Connections**, que son objetos en el **Glue Data Catalog** usados para acceder de forma segura a estas fuentes o destinos.
- 

## 2. Crawler

- El **Crawler** escanea un Data Store (por ejemplo, un bucket S3 con archivos CSV).
  - Usa **clasificadores** para identificar el formato (CSV, JSON, Parquet, etc.).
  - Automáticamente infiere el esquema (nombres de columnas, tipos de datos) y lo **registra como tablas** en el **Glue Data Catalog**.
  - Se puede ejecutar de forma manual o automática (por eventos o triggers).
- 

## 3. AWS Glue Data Catalog

- Es un almacén o **repositorio centralizado de metadatos**.
  - Contiene:
    - **Bases de datos lógicas** (agrupan tablas).
    - **Tablas** (esquemas de archivos o bases de datos reales).
    - **Conexiones** (credenciales + ubicación).
  - Cada vez que un crawler corre, actualiza el catálogo.
- 

## 4. AWS Management Console

- Es la **interfaz para configurar** crawlers, jobs, triggers, conexiones, etc.
  - También puedes lanzar trabajos manualmente desde aquí.
  - Alternativamente, puedes usar **AWS CLI** o **SDKs**.
- 

## 5. Job

- Un **Job de Glue** es una unidad de procesamiento que ejecuta un **script de transformación** de datos.
  - Puede ser de tipo:
    - **Spark** (PySpark o Scala): procesamiento distribuido.
    - **Python Shell**: scripts ligeros como validación, pequeñas tareas.
  - Puede usar:
    - **DynamicFrame**: estructura optimizada de Glue (más flexible que DataFrame).
    - **DataFrame** de Spark para transformaciones más avanzadas.
- 

## 6. **Script (Transformación)**

- Aquí defines tu lógica ETL:
    - Limpieza, filtrado, mapeo, enriquecimiento.
    - En tu caso: validaciones, flags por FICO, ratios, etc.
  - Se puede desarrollar usando:
    - **Glue Studio (visual)**
    - **Development Endpoints + Jupyter Notebook Server** (interactivo)
    - **VS Code + CLI Toolkit** (como lo haces tú)
  - Se ejecuta en una infraestructura administrada por AWS usando **Workers (DPU)**.
- 

## 7. **Trigger / Schedule or Event**

- Un Job puede ejecutarse:
    - Por programación (**Schedule**).
    - Por evento (**por ejemplo, cuando se sube un archivo a S3**).
    - En cadena (cuando termina otro job, **Trigger**).
  - Los **Triggers** permiten definir workflows complejos.
- 

## 8. **Transform → Load**

- El script toma datos del Data Source, aplica transformaciones y los carga en el **Data Target**.
  - El destino puede ser:
    - S3 (en formatos optimizados como Parquet).
    - Amazon Redshift (usando JDBC Sink).
    - Bases de datos relacionales.
  - Este paso es típicamente el resultado final de un Job exitoso.
- 



## Términos clave conectados al flujo:

Término	Descripción	Rol en el flujo
<b>Glue Data Catalog</b>	Catálogo de metadatos	Centraliza info de fuentes/destinos
<b>Clasificador</b>	Identifica formato de archivo	Usado por el Crawler
<b>Conexión</b>	Objeto del Catálogo con credenciales	Para acceder a RDS, Redshift, etc.
<b>Crawler</b>	Escanea fuentes y crea tablas	Detecta y registra datos automáticamente
<b>Base de Datos (Glue)</b>	Agrupador lógico de tablas	Organiza tus datasets
<b>Development Endpoint</b>	Entorno para desarrollo de scripts	Útil para Jupyter Notebooks
<b>Dynamic Frame</b>	Estructura de Glue optimizada	Más robusta que Spark DF
<b>Job</b>	Proceso ETL que ejecuta scripts	Core de la lógica empresarial
<b>Notebook Server</b>	Entorno interactivo	Desarrollo exploratorio
<b>Script</b>	Código PySpark/Scala/Python	Contiene la lógica de transformación
<b>Tabla</b>	Representación del esquema	Se almacena en el Data Catalog
<b>Transformación</b>	Proceso de cambio en los datos	Parte del script de ETL
<b>Trigger</b>	Desencadenador de Jobs	Automatización
<b>Worker (DPU)</b>	Unidad de procesamiento	Se configura por job

## Sources

<https://aws.amazon.com/es/glue/features/databrew/>

<https://docs.aws.amazon.com/glue/latest/dg/machine-learning-teaching.html>