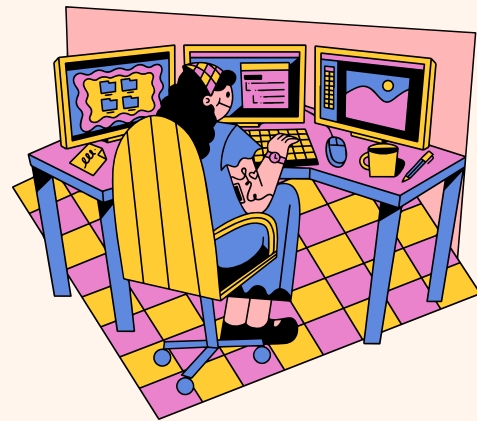


# AWS Cloud Data Pipeline

Serverless Analytics Pipeline Using Amazon Web Services



**Esther Quirós**  
Developer



**Nieves Pérez**  
Data Analyst



**Elena Pavón**  
Data Analyst



**Mayka Durán**  
Data Analyst



**Claudia Cervantes**  
Data Analyst · Scrum Master

*Team Project • AWS re/Start Program*

# ¿Por qué este proyecto?

Como equipo principalmente de Data Analysts, queríamos llevar nuestros conocimientos de Data Analytics a un entorno cloud real utilizando servicios de AWS.

El objetivo fue construir un pipeline automatizado capaz de:

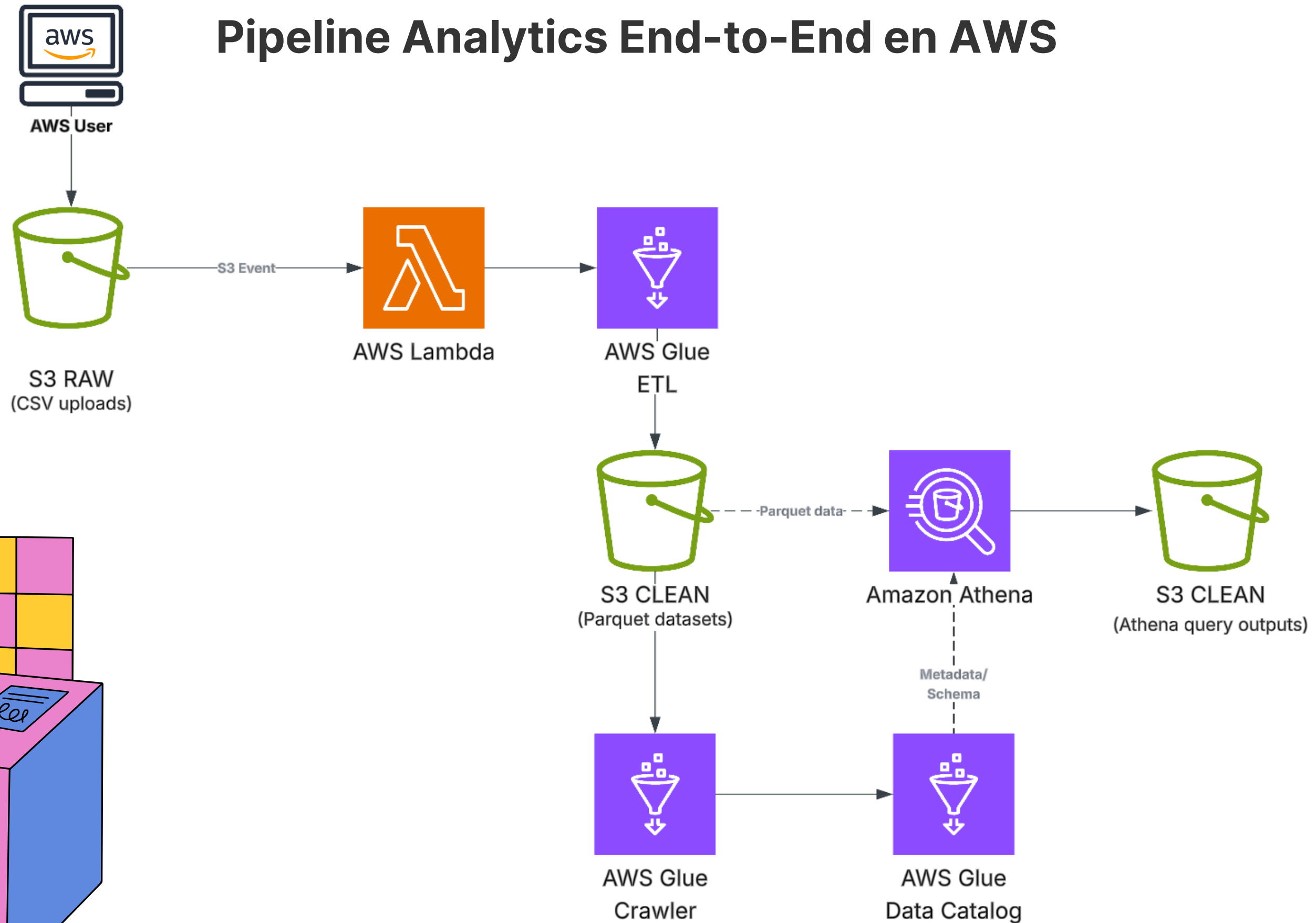
- ingestión de datos
- transformación ETL
- automatización
- almacenamiento en la nube
- consultas SQL analytics



*El dataset fue el caso práctico — el verdadero reto fue construir el pipeline.*

# Arquitectura General

## Pipeline Analytics End-to-End en AWS



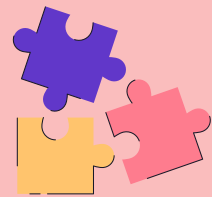
# Amazon S3 como Data Lake



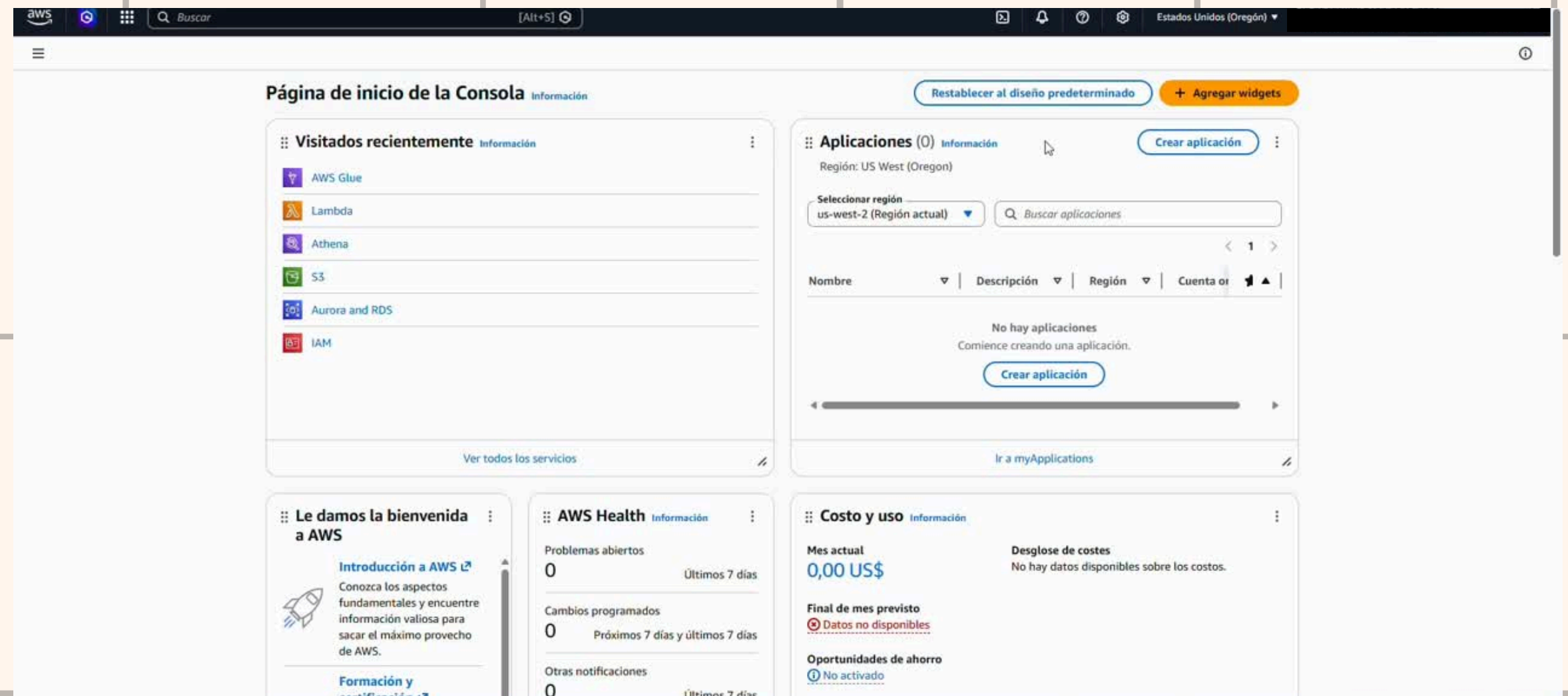
- Amazon S3 = almacenamiento escalable y seguro
- Actúa como Data Lake
- Guarda datos RAW y CLEAN
- Integración con:



AWS Glue  
Amazon Athena  
AWS Lambda



# Organización del almacenamiento



## Bucket 1 — RAW



- Datos originales (CSV)
- Sin modificar
- Fuente del pipeline

¿Por qué lo separamos?



## Bucket 2 — CLEAN



- Datos transformados
- Formato Parquet
- Listos para análisis

- No perder datos originales
- Evitar errores
- Pipeline más claro

# Seguridad y configuración

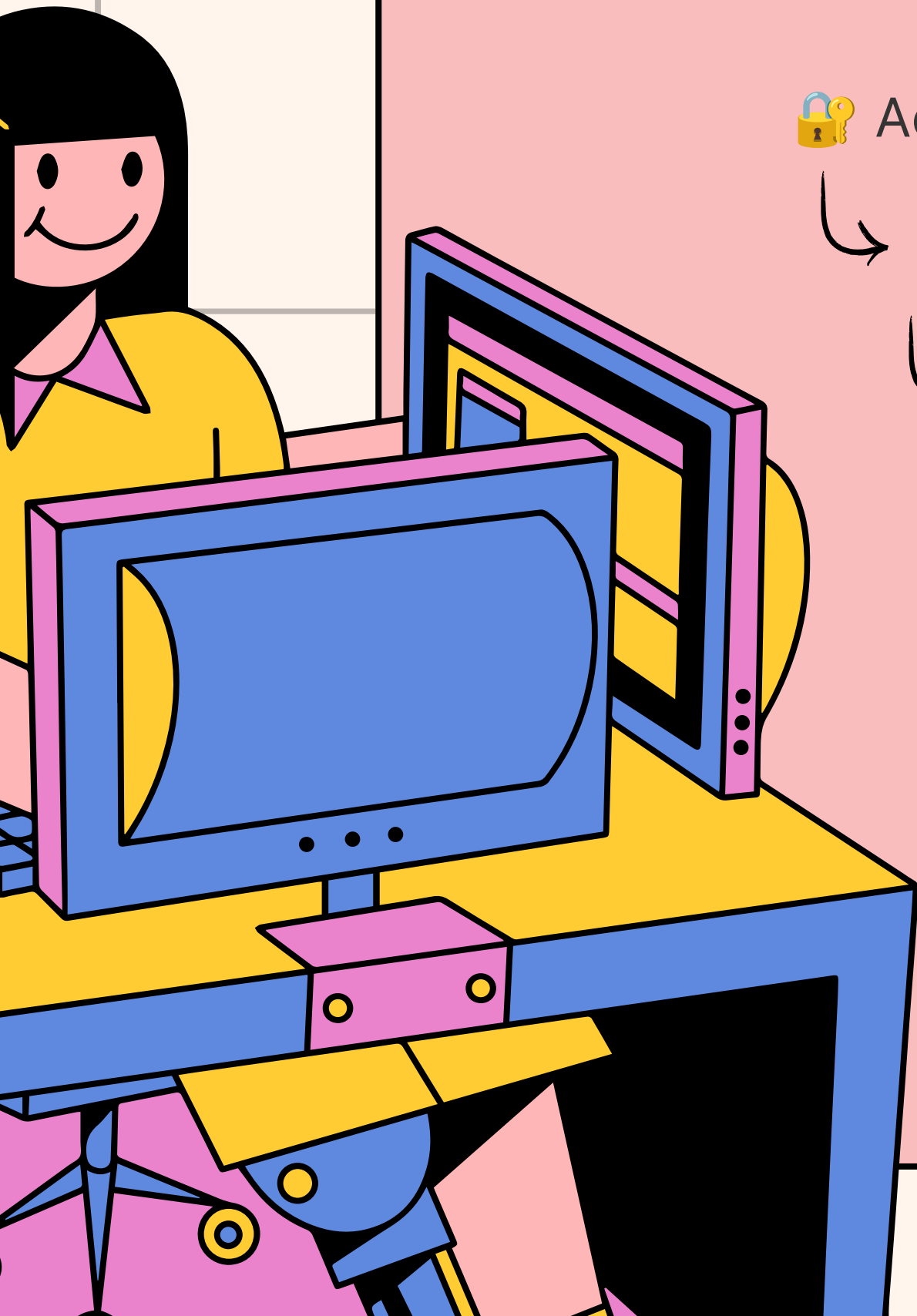
🔑 Acceso mediante roles (IAM)

🔒 Acceso público bloqueado

♻️ Versionado: desactivado  
(laboratorio)

Objetivo

Seguridad + simplicidad + control ✓

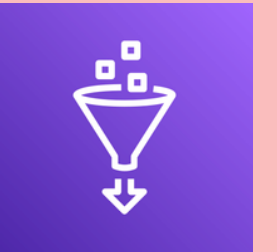


# Configuración de AWS GLUE

¿ cómo preparamos el "Robot" de limpieza?

**Estrategia de Datos:** Optamos por un **\*\*Data Lake con AWS Glue\*\*** en lugar de Aurora SQL.

*¿Qué datos tenemos y que queremos hacer con ellos?*



## CONFIGURACIÓN TÉCNICA

### SCRIPT EDITOR

**Motor:**

Seleccionamos **\*Python Shell\*** por su rapidez de arranque y eficiencia en datasets de tamaño medio (Kaggle).

**Modo "Start Fresh":** Usamos un editor limpio para inyectar nuestro código personalizado de forma directa



**Getting started**  
 ETL jobs  
 Visual ETL  
 Notebooks  
 Job run monitoring  
 Data Catalog tables  
 Data connections  
 Workflows (orchestration)  
 Zero-ETL integrations **New**

**▼ Data Catalog**  
 Catalogs  
 Databases  
 Tables  
 Stream schema registries  
 Schemas  
 Connections  
 Crawlers  
 Classifiers  
 Catalog settings

**Prepare your account for AWS Glue**  
 Admins: Grant access to AWS Glue and set a default IAM role.  
[Set up roles and users](#)

**Catalog and search for datasets**  
 View your databases & tables and catalog data using Crawlers.  
[Go to the Data Catalog](#)

**Move and transform data** **Updated**  
 Use Zero-ETL integrations to replicate data in near real-time, or ETL jobs to transform data in visual, notebook, or code interface.  
[Go to Zero-ETL integrations](#) [Go to ETL jobs](#)

**Resources and tutorials** [↗](#)  
 Getting started with AWS Glue: [Documentation](#) [AWS Training](#)  
 Glue in 5 Minutes Videos: [Authoring](#), [GenAI](#), [Monitoring](#), [Orchestration](#)  
[Using connectors and connections](#)  
[AWS Glue Documentation home](#)

**Data integration and management**  
 Monitor & debug ETL jobs and track usage  
[Go to job run monitoring](#)  
 Connect to your data stores

## 1. Extracción (s3--> Boto3/Pandas)

### 1. Proceso de Transformación:

- Estandarización de nombres
- Eliminación de duplicados
- Manejo de valores nulos
- Homogeneización de Ratings
- Enriquecimiento de fechas
- Parsing de Duración

### 3. Carga (S3 en formato Parquet):

- queries más rápidas
- menor costo Athena
- mejor compresión
- preserva schema
- soporta arrays/listas
- mejor performance

# El Script ETL y la Magia de Parquet

*De CSV "sucio" a un formato de alto rendimiento.*





# DEL ARCHIVO A LA CONSULTA SQL (CRAWLER & CATALOG)

## SEGURIDAD Y PERMISOS:

- USO DEL **LABROLE** (ENTORNO AWS ACADEMY).
- Asociación de políticas clave: ``S3FullAccess`` y ``AWSGlueServiceRole`` para permitir la comunicación entre servicios.

\* Asociación de políticas clave: ``S3FullAccess`` y ``AWSGlueServiceRole`` para permitir la comunicación entre servicios.

## PROCESO DE CATALOGACIÓN:

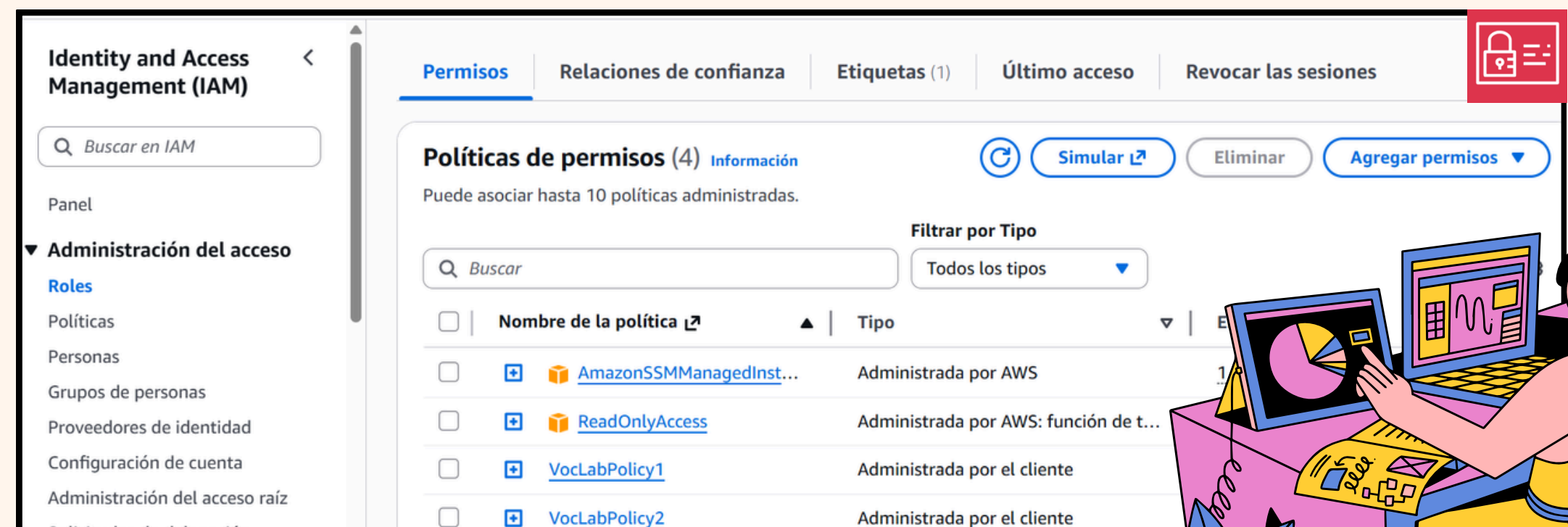
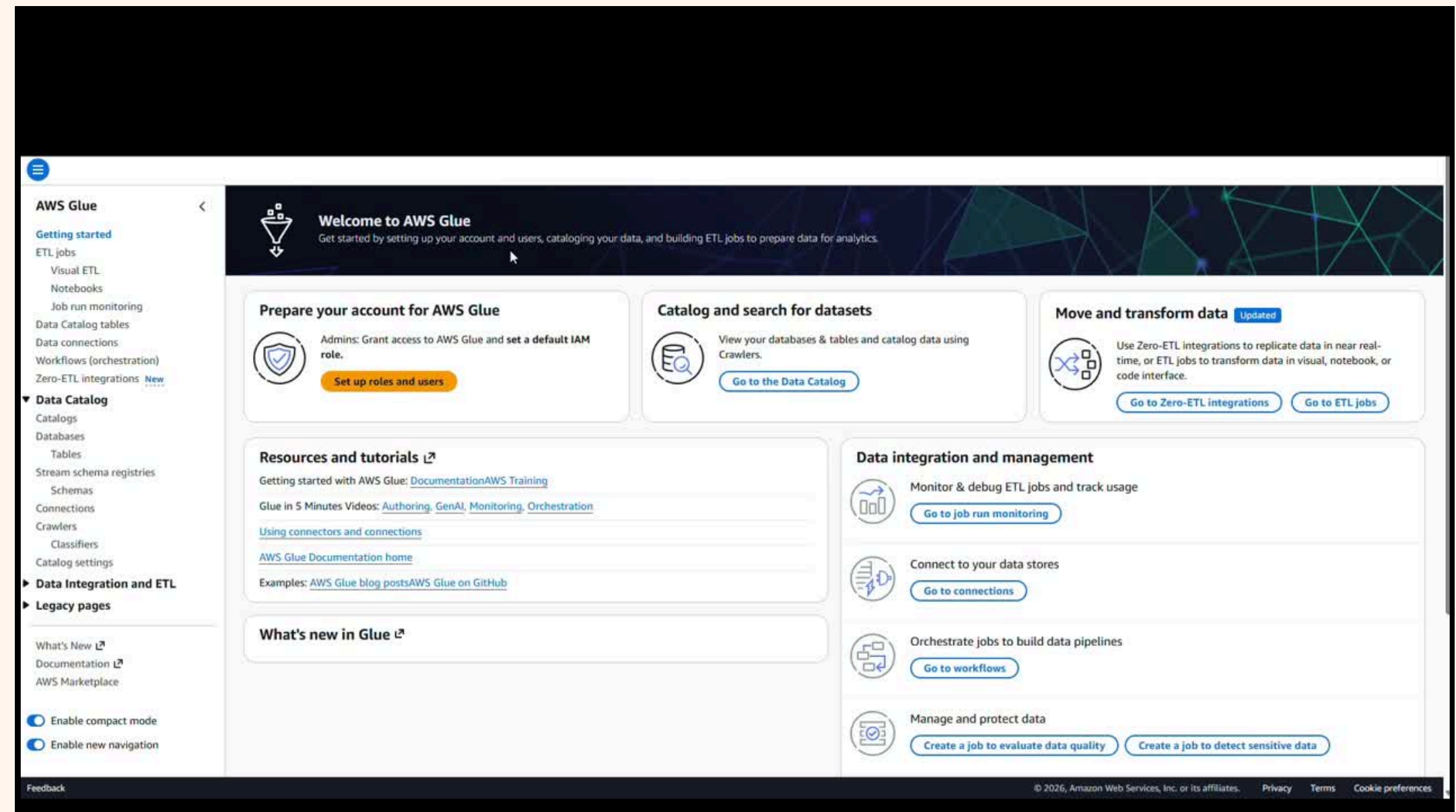
### 1. **Crawler (El Explorador):**

Escanea el *bucket clean*, identifica el esquema de los archivos *Parquet* y detecta cambios.

### 2. **Data Catalog (La Biblioteca):**

Repositorio central de metadatos que guarda el "mapa" de nuestra información.

3. **Database:** Agrupación lógica donde residen las tablas virtuales.

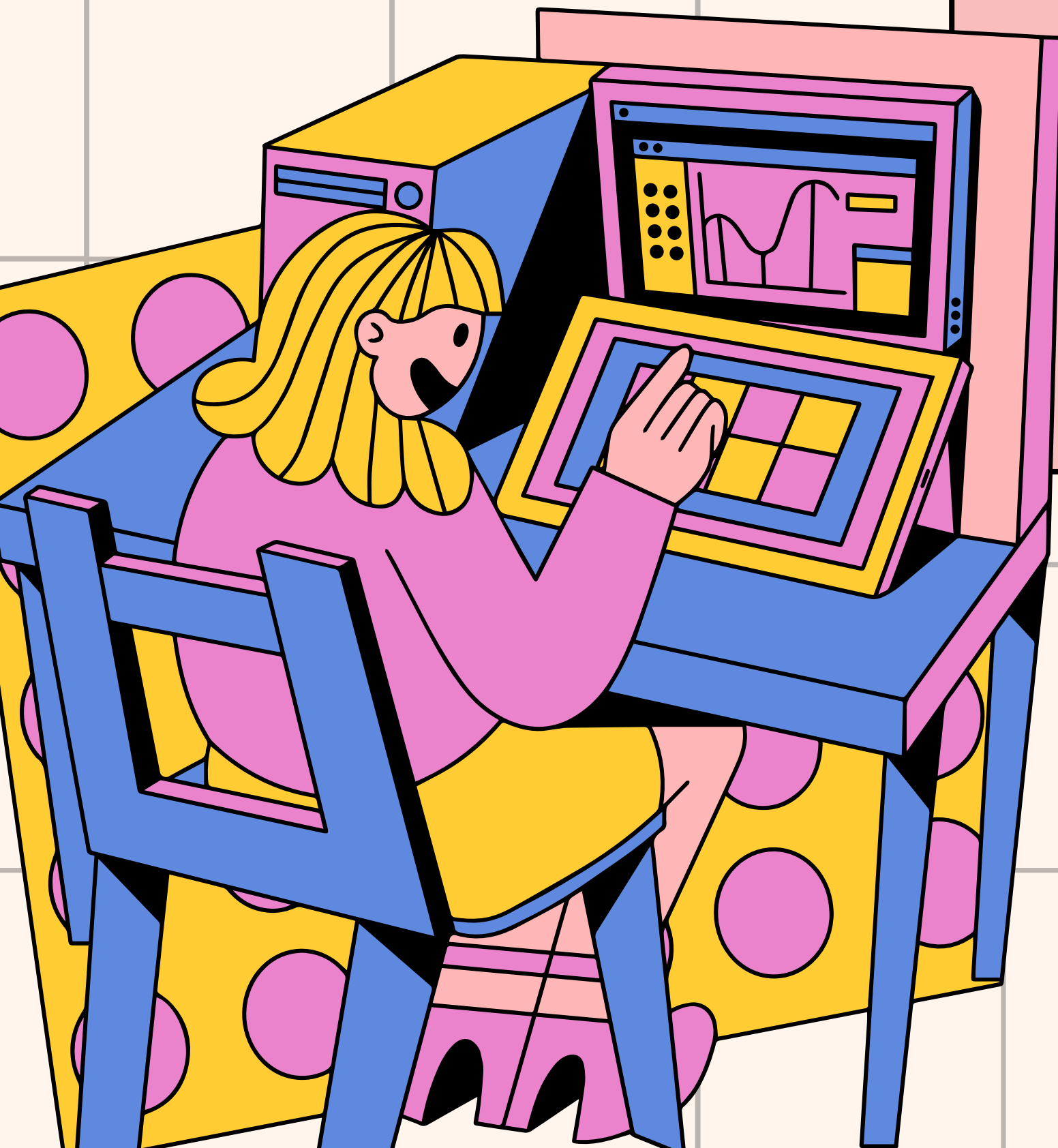




# ¿Qué es Athena?

Servicio de consultas interactivo que permite analizar grandes volúmenes de datos almacenados en Amazon S3 con SQL estándar

- ↳ serverless
- ↳ escala automáticamente
- ↳ funciona directamente sobre Parquet
- ↳ evita cargar datos en MySQL





# ¿Por qué usamos Athena?

Porque nuestro proyecto es un **analytical workload**

✓ Queríamos

- ↳ Analizar datos
- ↳ Agrupar información
- ↳ Generar insights

✗ NO:

- ↳ Manejar transacciones
- ↳ Actualizar usuarios
- ↳ Crear aplicaciones



# Demo de Athena



aws Buscar [Alt+S] United States (Oregon) ▼

**AWS Glue** > Tables

### AWS Glue

- Getting started
- ETL jobs
  - Visual ETL
  - Notebooks
  - Job run monitoring
- Data Catalog tables**
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations New
- ▼ **Data Catalog**
  - Catalogs
  - Databases
    - Tables**
  - Stream schema registries
    - Schemas
  - Connections
  - Crawlers
    - Classifiers
  - Catalog settings
- **Data Integration and ETL**
- **Legacy pages**

---

- What's New [↗](#)
- Documentation [↗](#)
- AWS Marketplace

Enable compact mode  
 Enable new navigation

## 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.

**Tables (1)** Last updated (UTC) May 4, 2026 at 11:38:23 Delete Add tables using crawler Add table

View and manage all available tables.

Choose catalog

Choose database

<input type="checkbox"/>	Name	Database	Location	Classification	Deprecated	View data	Data quality	Column statistics
<input type="checkbox"/>	clean_parquet	prime_analytics	s3://amazon-prime-clean-t	Parquet	-	Table data	View data quality	View statistics

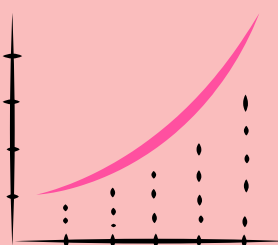
CloudShell Feedback © 2026, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

# ¿Porqué S3 Triggers y AWS Lambda?



## **Eficiencia - Latencia Cero**

Reacción instantánea ante la llegada de datos. Sin esperas manuales ni procesos programados innecesarios.



## **Escalabilidad**

AWS Lambda escala automáticamente para procesar múltiples subidas simultáneas sin cuellos de botella.



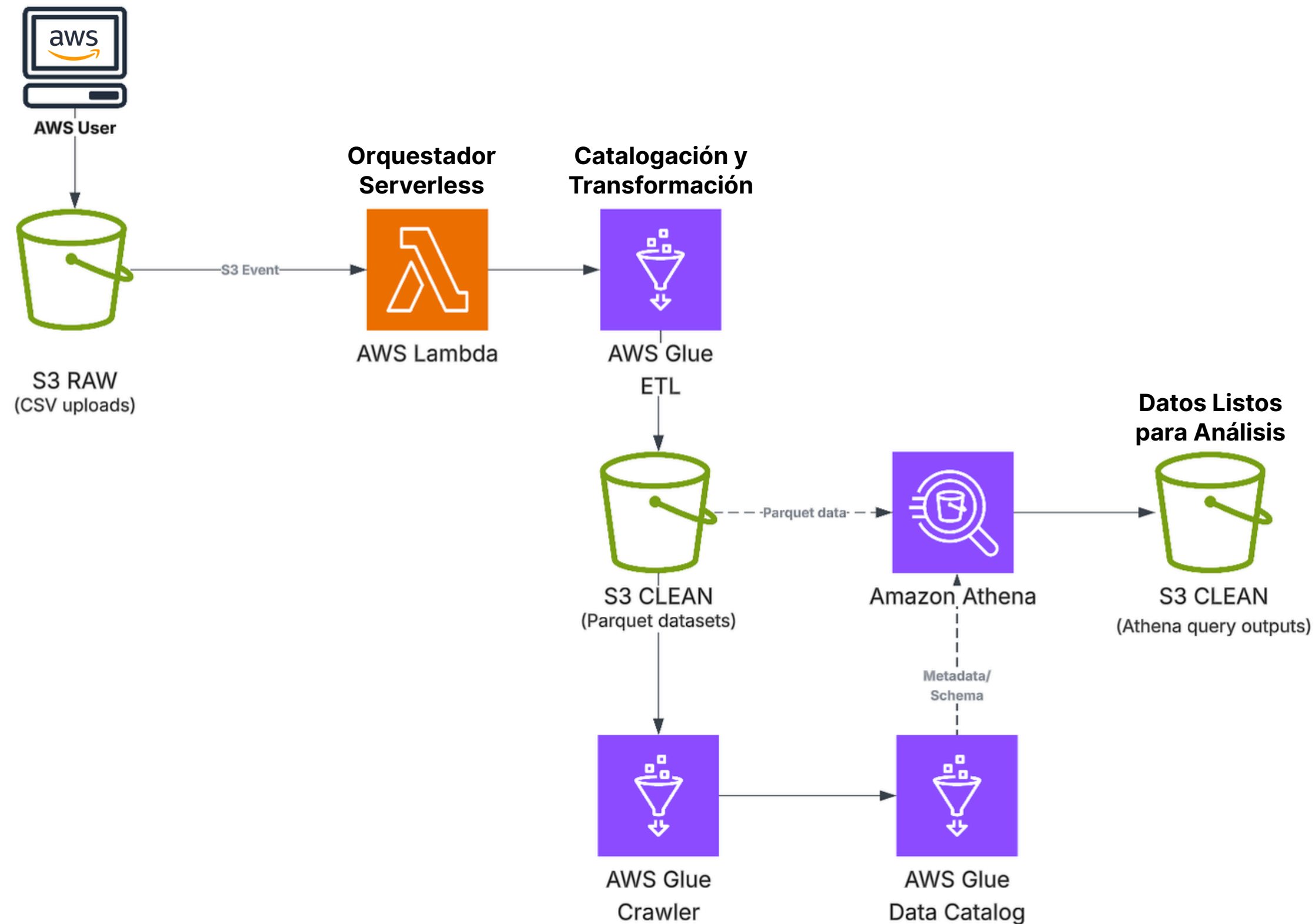
## **Coste \$0**

Solo pagamos por los milisegundos de ejecución. Eficiencia total en el uso de presupuesto cloud.



# Lambda Automation

## El Flujo del Dato



# Demo Lambda



The screenshot shows the AWS console interface. On the left, there's a navigation sidebar with categories like 'Amazon Athena', 'Datos', 'Origen de datos', 'Catálogo', 'Base de datos', 'Tablas y vistas', 'Tablas (1)', and 'Vistas (0)'. The main area is divided into sections: 'Servicios' (127 results), 'Características' (353 results), and 'Resources in us-west-2'. The 'Servicios' section lists 'Lambda' (Ejecuta código sin tener que pensar en los servidores), 'WorkSpaces Secure Browser' (Acceso web seguro nativo en la nube), and 'Amazon Transcribe' (Excelente reconocimiento de voz). The 'Características' section lists 'Lambda Insights' (Función CloudWatch), 'Puntos de acceso de Lambda de objetos' (Función S3), and 'Operaciones por lotes' (Función S3). At the bottom, there's a 'Resources in us-west-2' section with a 'Complete setup' button.

On the right side of the console, a query execution summary is displayed:

- Tiempo en cola: 99 ms
- Tiempo de ejecución: 607 ms
- Datos analizados: 505.45 KB

Below the summary, there are buttons for 'Copiar' and 'Descargar resultados en formato CSV'. A table below shows the results of the query:

total_titulos
56
35

At the bottom of the console, there's a footer with 'CloudShell', 'Comentarios', and copyright information: '© 2026, Amazon Web Services, Inc. o sus filiales. Privacidad Términos Preferencias de cookies'.


# Comprobación

# OK







The screenshot shows the AWS Lambda console for the function 'trigger-glue-prime-etl'. At the top, a green notification bar states: 'El desencadenador amazon-prime-raw-equipos-mdm se agregó correctamente a la función trigger-glue-prime-etl.' Below this, the 'Información general de la función' section is active, showing a diagram with the function icon and a connection to an S3 bucket. The right sidebar displays the function's ARN: 'arn:aws:lambda:us-west-2:613492405651:function:trigger-glue-prime-etl' and its last modification time: 'hace 35 segundos'. The 'Configuración' tab is selected, showing a list of triggers under 'Desencadenadores (1)'. One trigger is listed: 'S3: amazon-prime-raw-equipos-mdm' with the ARN 'arn:aws:s3::amazon-prime-raw-equipos-mdm'. The bottom of the page includes a footer with '© 2026, Amazon Web Services, Inc. o sus filiales.' and links for 'Privacidad', 'Términos', and 'Preferencias de cookies'.

# Arquitectura y decisiones técnicas

Durante el diseño inicial consideramos utilizar servicios como Amazon  RDS y  EC2.

Sin embargo, al analizar el tipo de *workload* y los objetivos analytics del proyecto, evolucionamos hacia una arquitectura serverless basada en S3, Glue y Athena.

## Decisiones clave:

- Data Lake en  S3 con consultas SQL serverless mediante  Athena
- transformación ETL serverless con  AWS Glue
- automatización orientada a eventos mediante  Lambda
- reducción de infraestructura para simplificar costes y operación

*La arquitectura final priorizó escalabilidad, simplicidad operativa y coste eficiente.*



# Conclusión

Construimos un pipeline analytics moderno, automatizado y serverless utilizando servicios reales de AWS.

## El pipeline permitió:

- ingestión y transformación ETL automatizada
- almacenamiento optimizado en formato Parquet
- consultas SQL serverless sobre S3 con Athena
- arquitectura cloud-native escalable
- costes operativos mínimos

## Next Steps:

- alertas automáticas con  SNS +  CloudWatch
- dashboards BI conectados a Athena

**Coste estimado:** ~\$0.77/mes, datasets pequeños y uso académico



# Muchas Gracias!

Pueden consultar el video del despliegue:

