



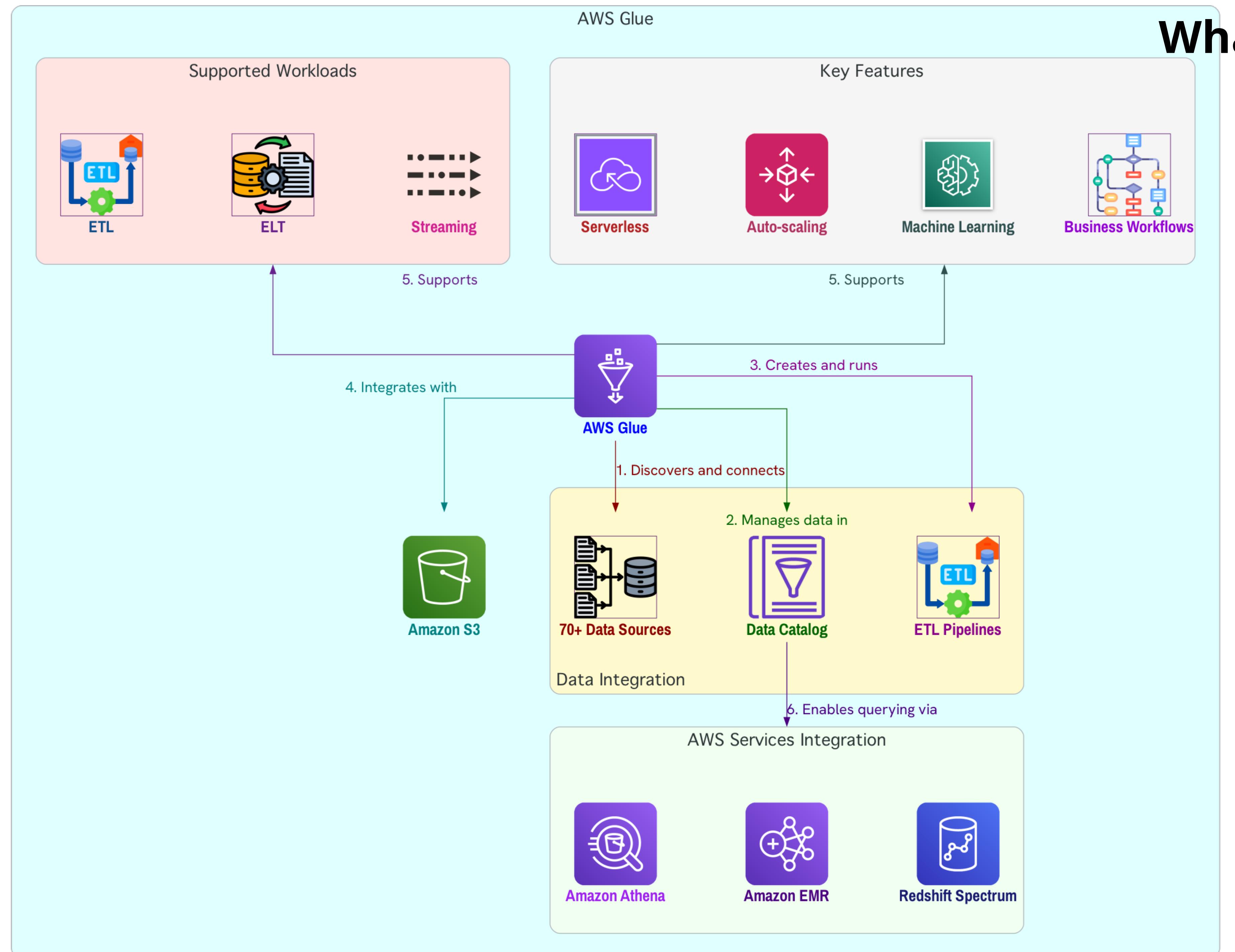
# AWS Glue

# Table of Contents



1. What is AWS Glue?
2. What is AWS Glue Studio?
3. AWS Glue features - 3 categories
4. Discover and organize data
5. Transform, prepare, and clean data for analysis
6. Build and monitor data pipelines
7. How it works
8. AWS Integrations
9. Architecture of an AWS Glue environment

# What is AWS Glue?



1. Discover, prepare, move & integrate data

Serverless data integration service

Analytics, ML, app development

2. 70+ data sources & centralized catalog

Connect diverse sources

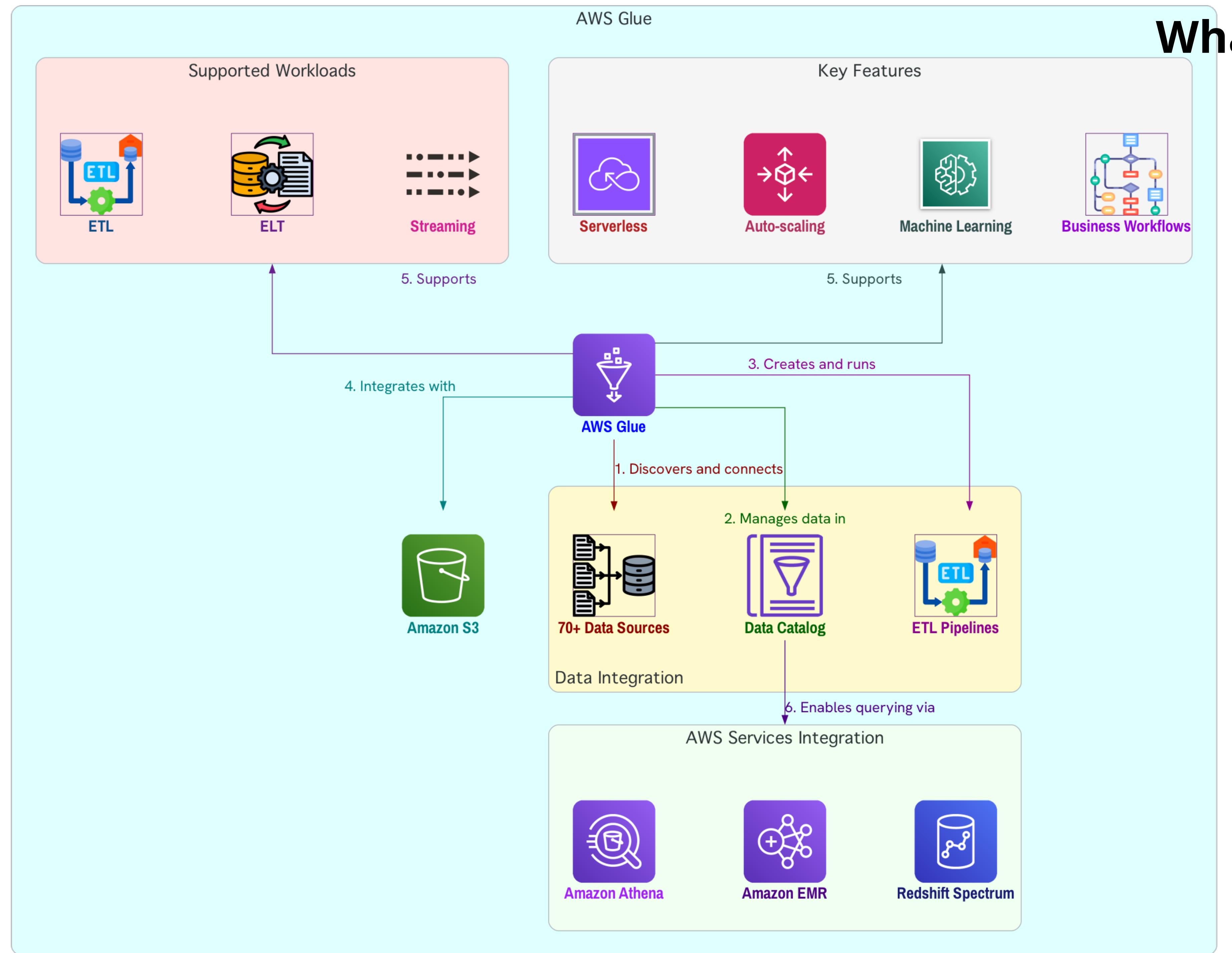
Centralized data management

3. Visual ETL pipeline creation

Visual creation and monitoring

Efficient data lake loading

# What is AWS Glue?



4. 🔎 Instant data querying with AWS services

🔍 Immediate search and query

🤝 Integration with Athena, EMR, Redshift Spectrum

5. ☁ Serverless & scalable architecture

🏗️ No infrastructure management

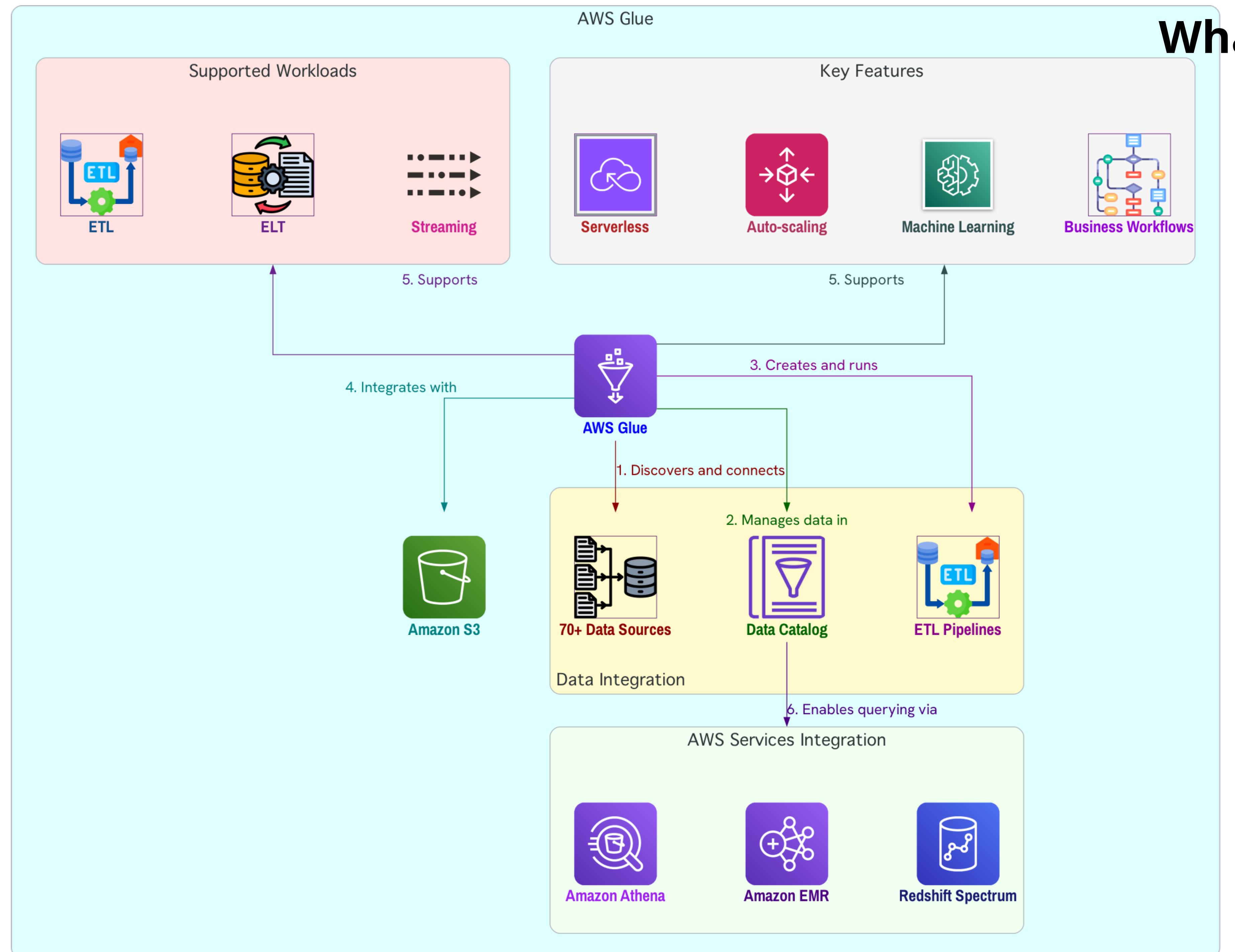
🤖 Consolidated data integration capabilities

6. 👤 Supports various user types

🔧 Flexible workload support

📊 ETL, ELT, streaming in one service

# What is AWS Glue?



7. ■ Easy-to-use interfaces & tools

■ Job-authoring tools

■ Tailored for varied skill sets

8. ■ On-demand scaling for any data size

■ Handles all data types and schemas

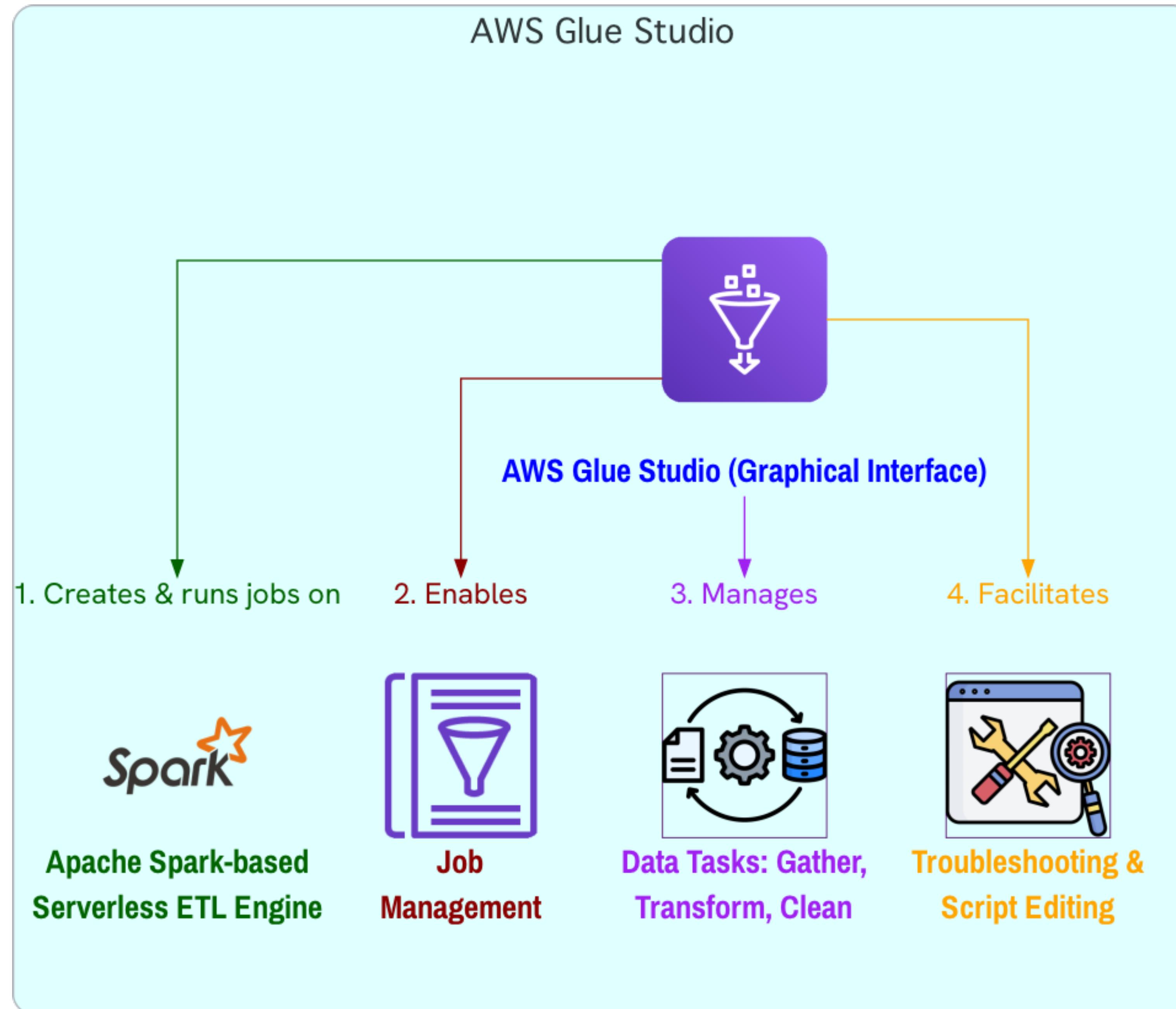
■ Focus on high-value activities

9. ■ \$ Pay-as-you-go billing

■ Cost optimization

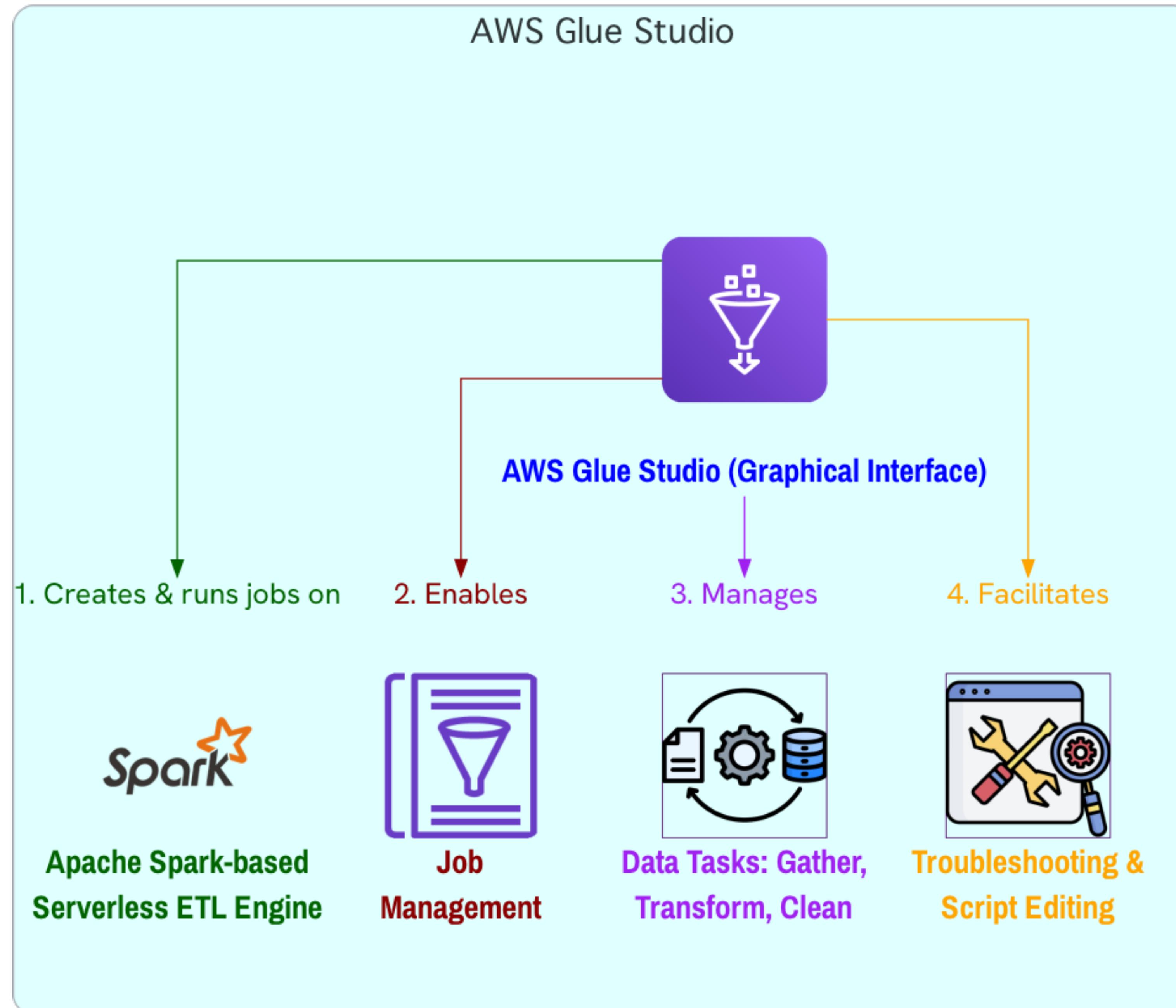
■ Built-in high availability

# What is AWS Glue Studio?



- 1. Graphical interface for AWS Glue
  - User-friendly interface
  - Accessible to various skill levels
- 2. Create, run, and monitor data integration jobs
  - Easy setup
  - Execute processes
  - Track integration workflows
- 3. Visual composition of data transformation workflows
  - Design visually
  - Manage complex data flows
  - Reduce code complexity

# What is AWS Glue Studio?



- 4. ⚡ Apache Spark-based serverless ETL engine
- 🚀 Efficient, scalable processing
- ☁️ No infrastructure management
- 5. 📊 Gather, transform, and clean data
  - Collect from various sources
  - Apply transformations
  - Clean for quality and consistency
- 6. ✖ Troubleshoot and edit job scripts
  - Debug tools
  - Modify scripts
  - Fine-tune processes

# AWS Glue features - 3 categories

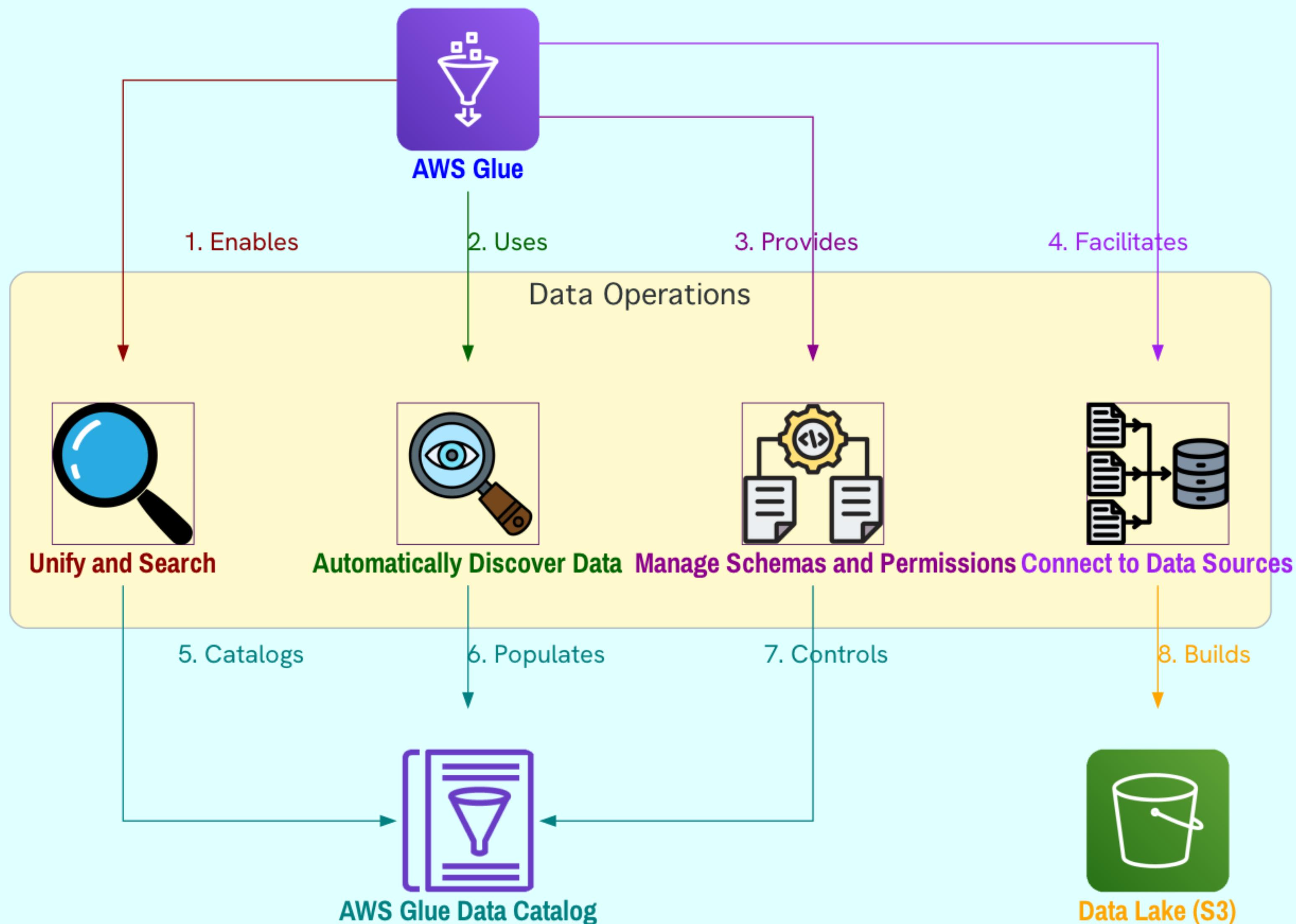
**Discover and organize data**

**Transform, prepare, and clean data for analysis**

**Build and monitor data pipelines**

# Discover and organize data

## Discover and Organize Data



### 1. Discover and organize data

#### Unify and search

Store, index, search data

#### Automatically discover data

Catalog in AWS

#### Manage schemas and permissions

Validate access

Control database/table access

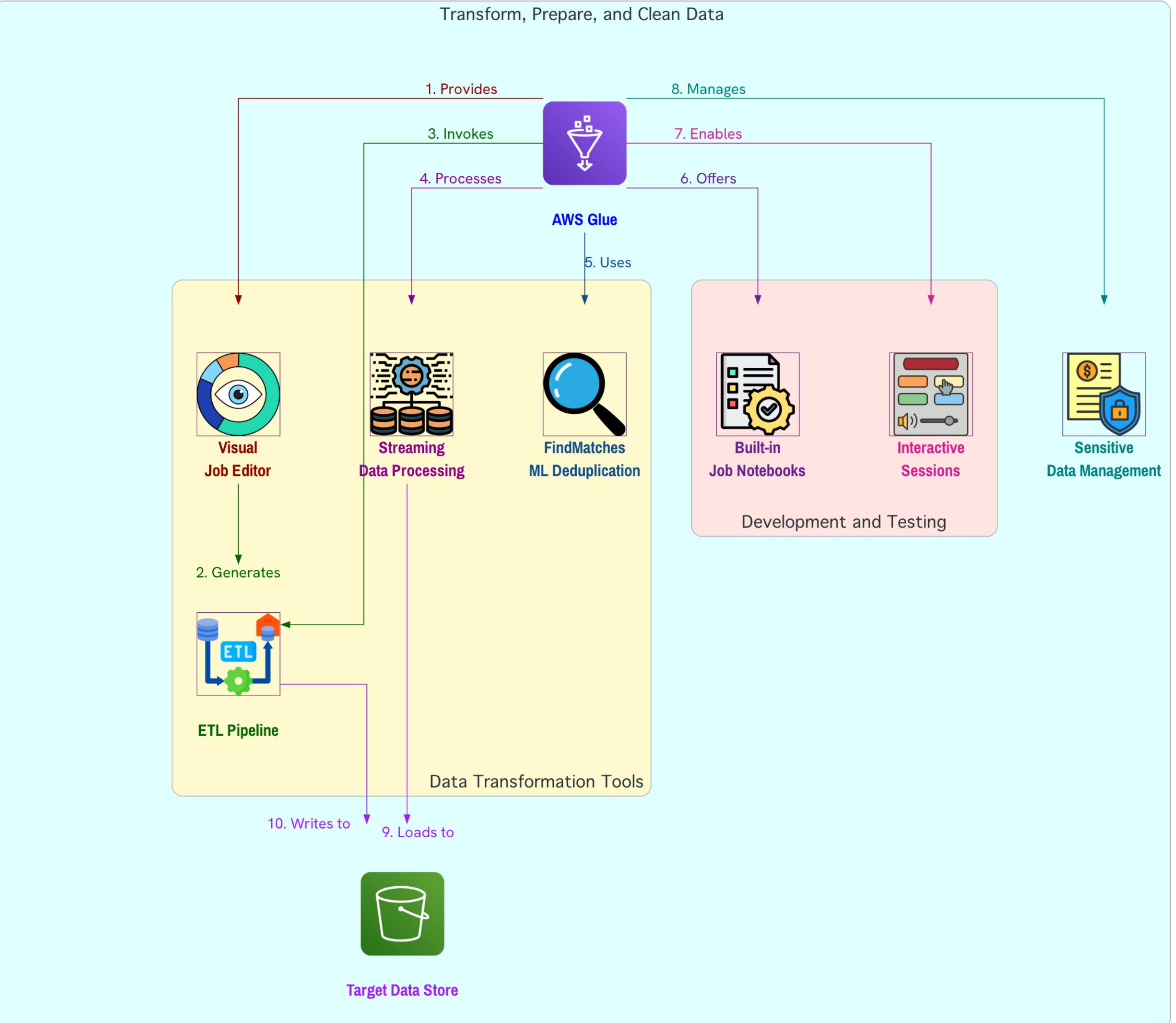
#### Connect to data sources

On-premises sources

AWS sources

Build data lake

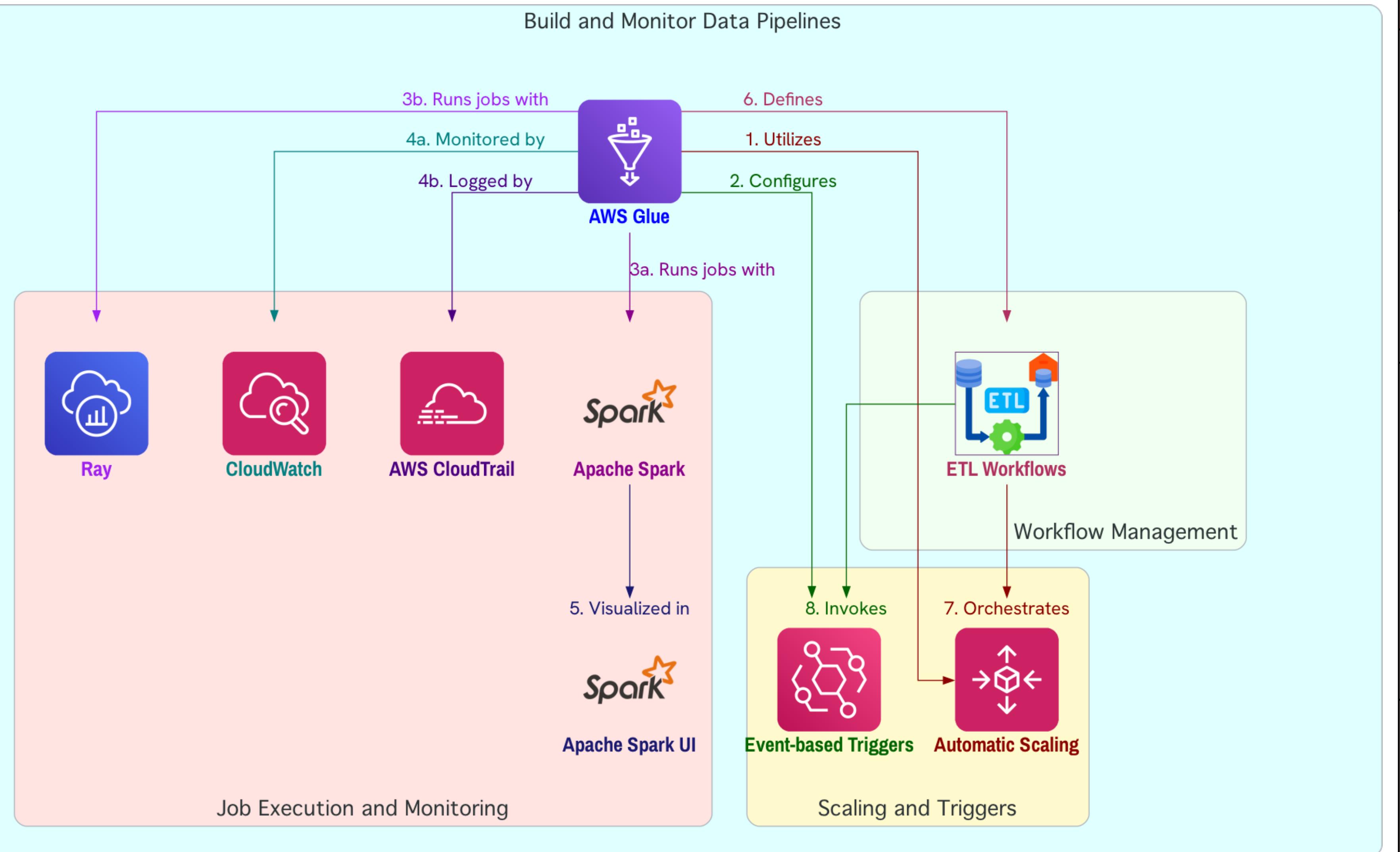
# Transform, prepare, and clean data for analysis



## 2. Transform, prepare, and clean data

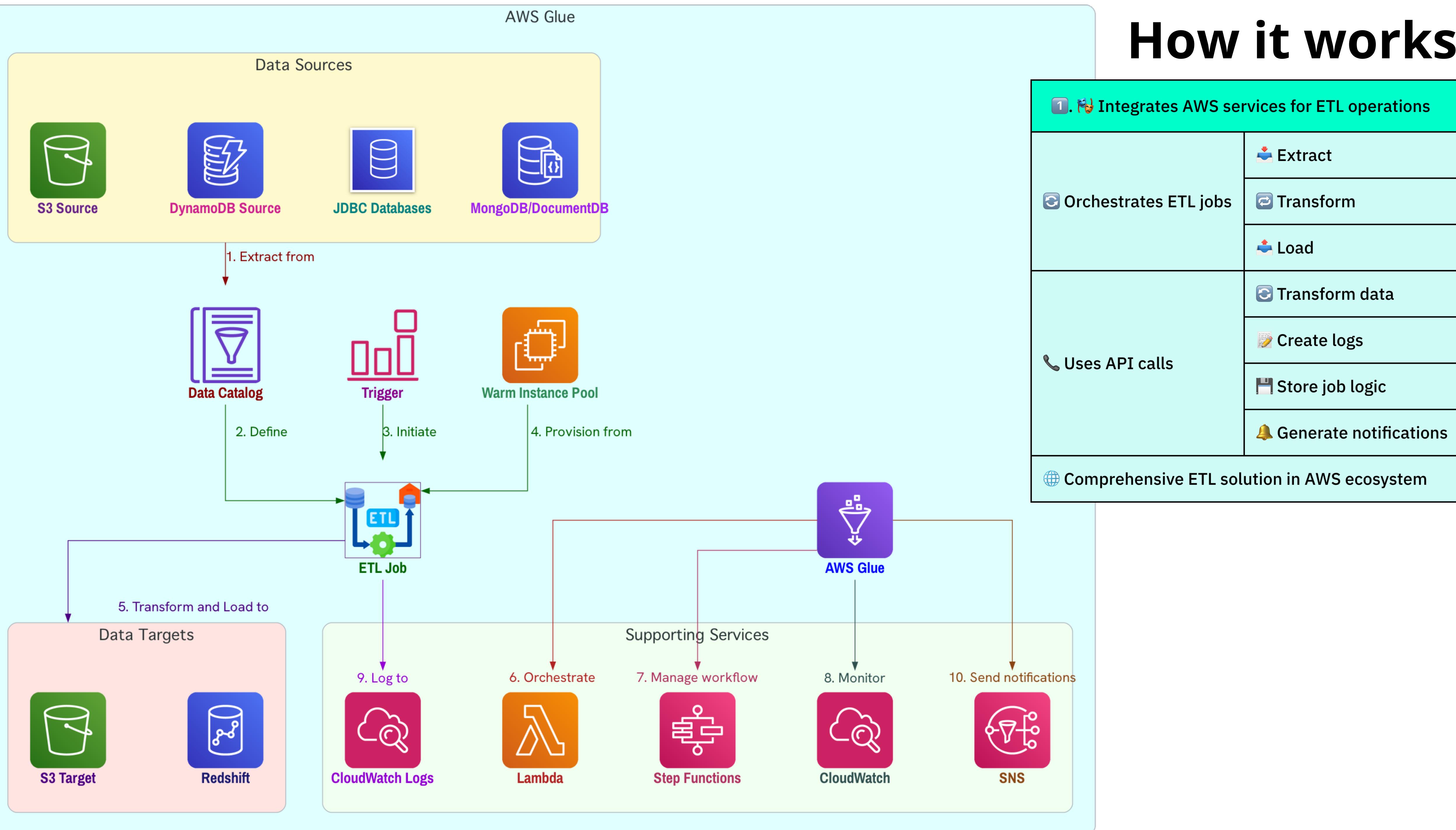
<b>Visual transformation</b>	<ul style="list-style-type: none"><li>Define ETL in visual editor</li><li>Auto-generate ETL code</li></ul>
<b>Complex ETL pipelines</b>	<ul style="list-style-type: none"><li>Schedule-based jobs</li><li>On-demand jobs</li><li>Event-based jobs</li></ul>
<b>Clean streaming data</b>	<ul style="list-style-type: none"><li>Continuous data consumption</li><li>Clean/transform in transit</li><li>Quick analysis in target store</li></ul>
<b>Deduplicate with machine learning</b>	<ul style="list-style-type: none"><li>FindMatches feature</li><li>Clean and prepare data</li><li>Find imperfect matches</li></ul>
<b>Built-in job notebooks</b>	<ul style="list-style-type: none"><li>Serverless notebooks</li><li>Quick setup</li></ul>
<b>Edit, debug, test ETL code</b>	<ul style="list-style-type: none"><li>Interactive sessions</li><li>Explore/prepare interactively</li><li>Use preferred IDE/notebook</li></ul>
<b>Sensitive data management</b>	<ul style="list-style-type: none"><li>Define sensitive data</li><li>Identify sensitive data</li><li>Process in pipeline/data lake</li></ul>

# Build and monitor data pipelines

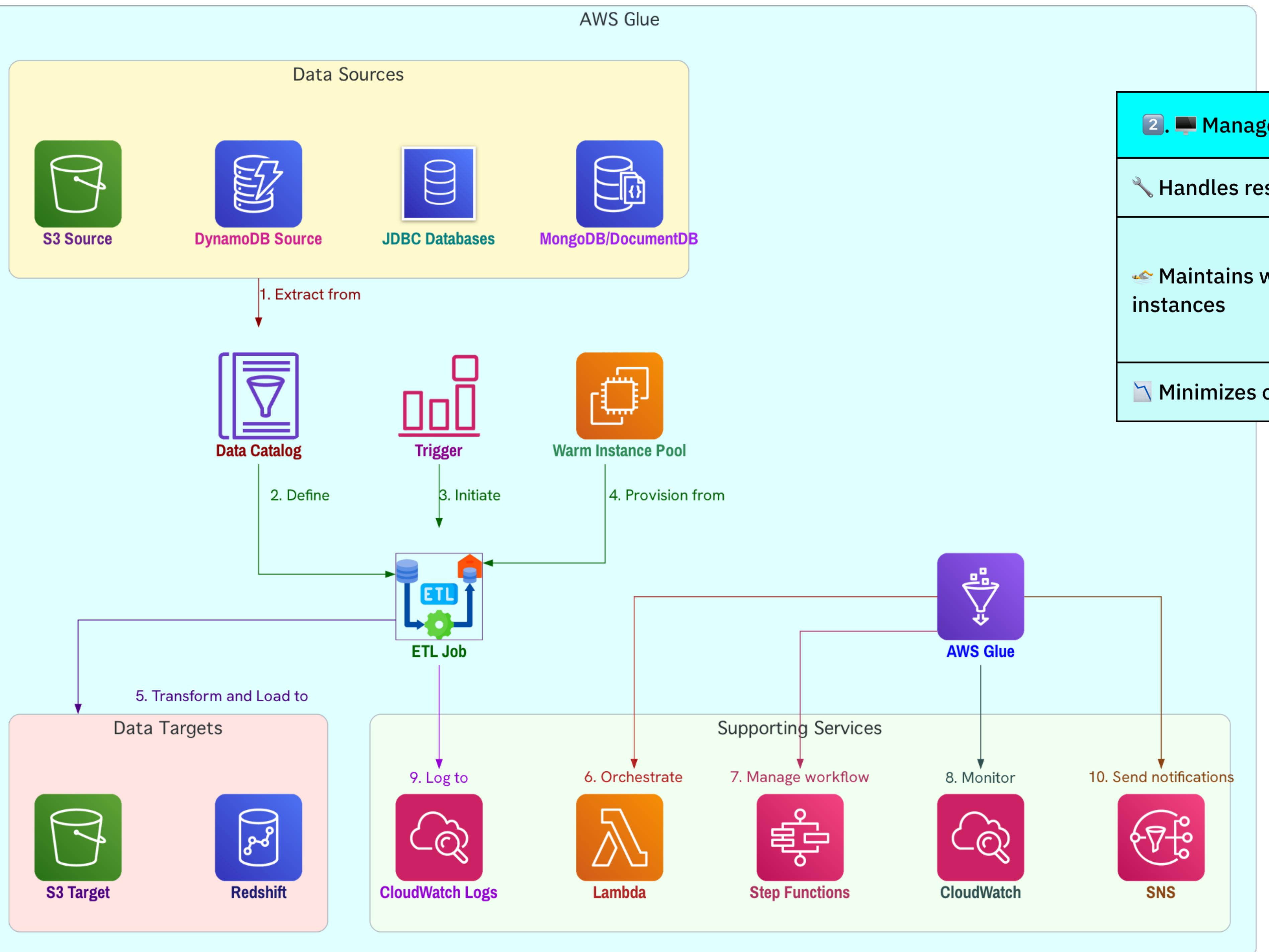


3. Build and monitor data pipelines	
<span>⚖️</span> <b>Automatic scaling</b>	<span>↗️</span> Dynamic resource scaling <span>👷</span> Assign workers as needed
<span>🎬</span> <b>Event-based triggers</b>	<span>▶️</span> Start crawlers/jobs <span>🔗</span> Chain dependent jobs/crawlers
<span>🏃</span> <b>Run and monitor jobs</b>	<span>⚡</span> Spark or Ray engines <span>📊</span> Automated monitoring tools <span>🔍</span> Job run insights <span>🧭</span> AWS CloudTrail <span>🕸️</span> Apache Spark UI
<span>⌚</span> <b>Define workflows</b>	<span>🔄</span> ETL activities <span>🤝</span> Integration activities <span>🕷️</span> Multiple crawlers <span>💼</span> Multiple jobs <span>⌚</span> Multiple triggers

# How it works



# How it works



2. Manages infrastructure automatically

Handles resource provisioning and management

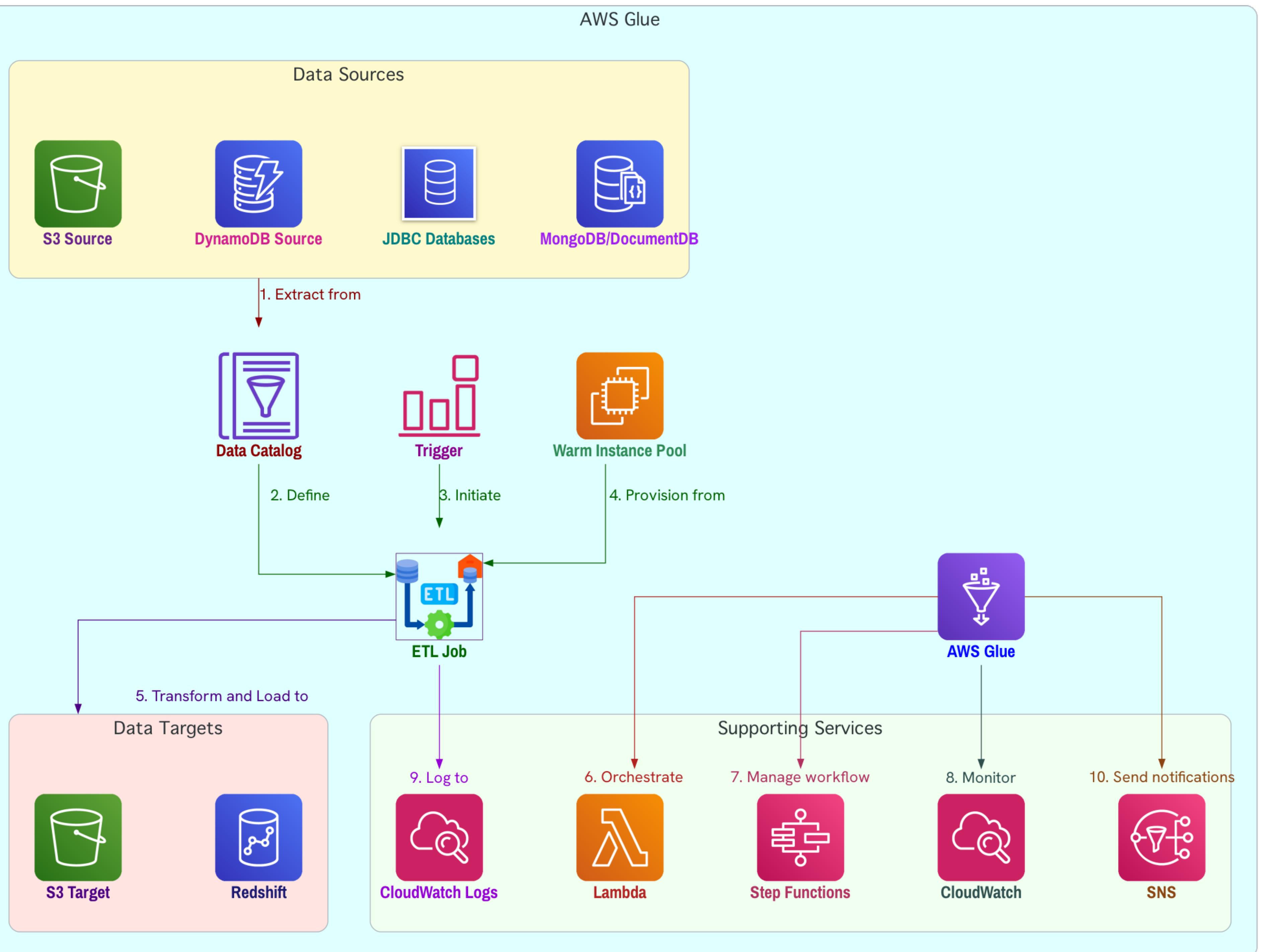
Reduces startup time

Maintains warm pool of instances

Ensures efficient resource allocation

Minimizes operational overhead

# How it works



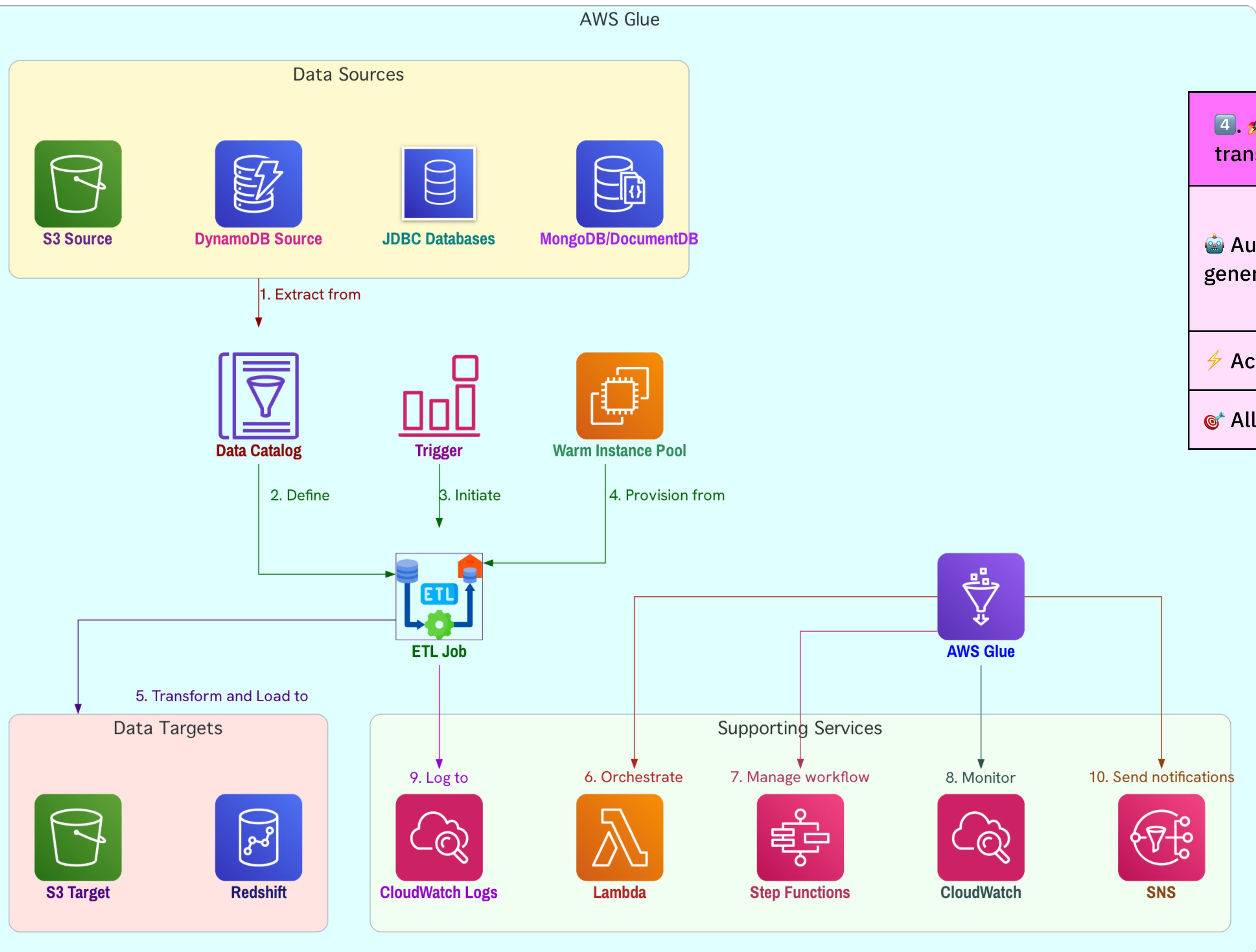
3. Utilizes Data Catalog for job structuring

Central metadata repository

Simplifies job creation and management

Provides clear structure for data assets

# How it works



4. Generates and executes transformation code

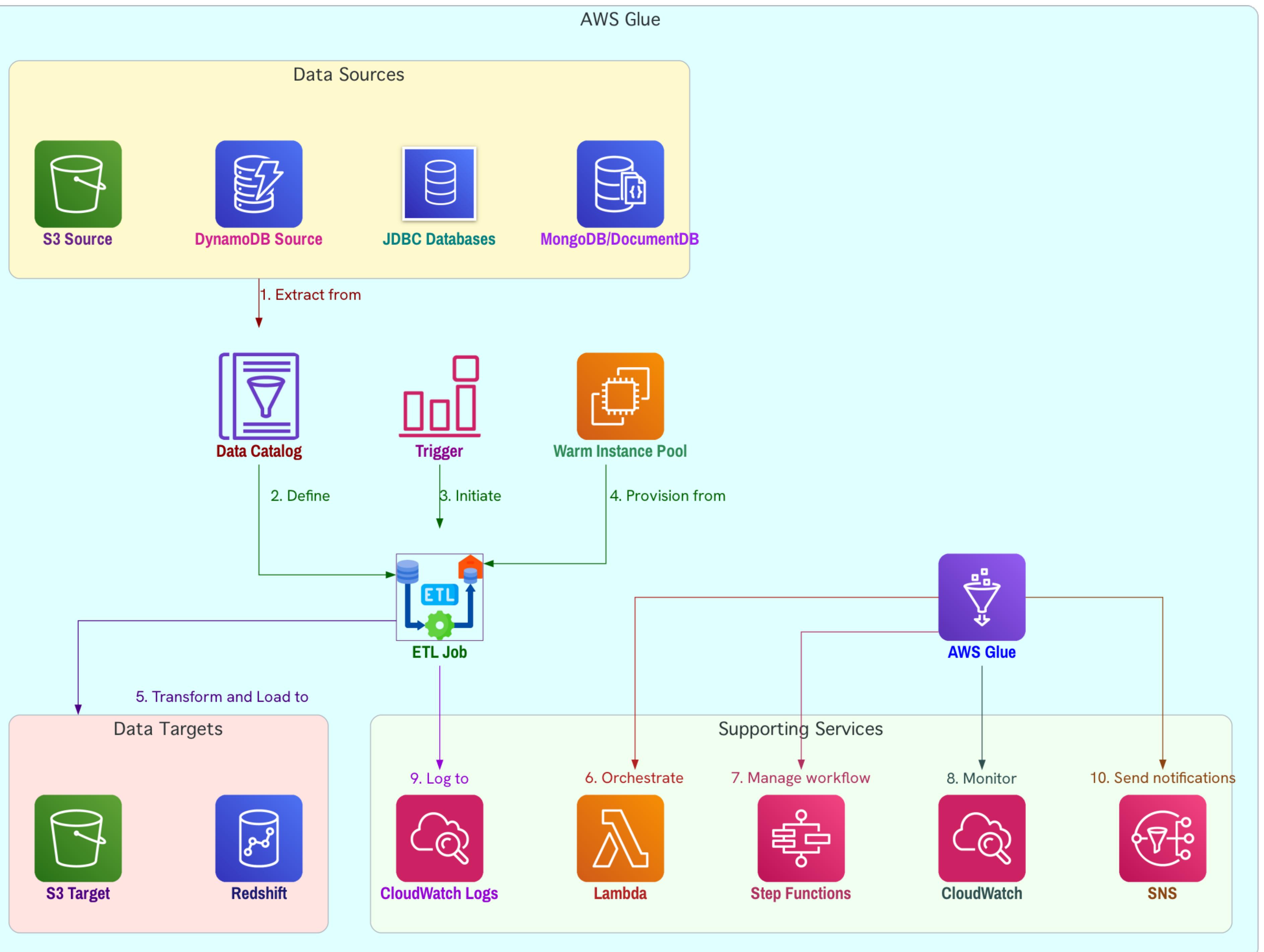
Based on user input  
Automatic code generation

For source and target data

Accelerates development process

Allows focus on data flow definition

# How it works



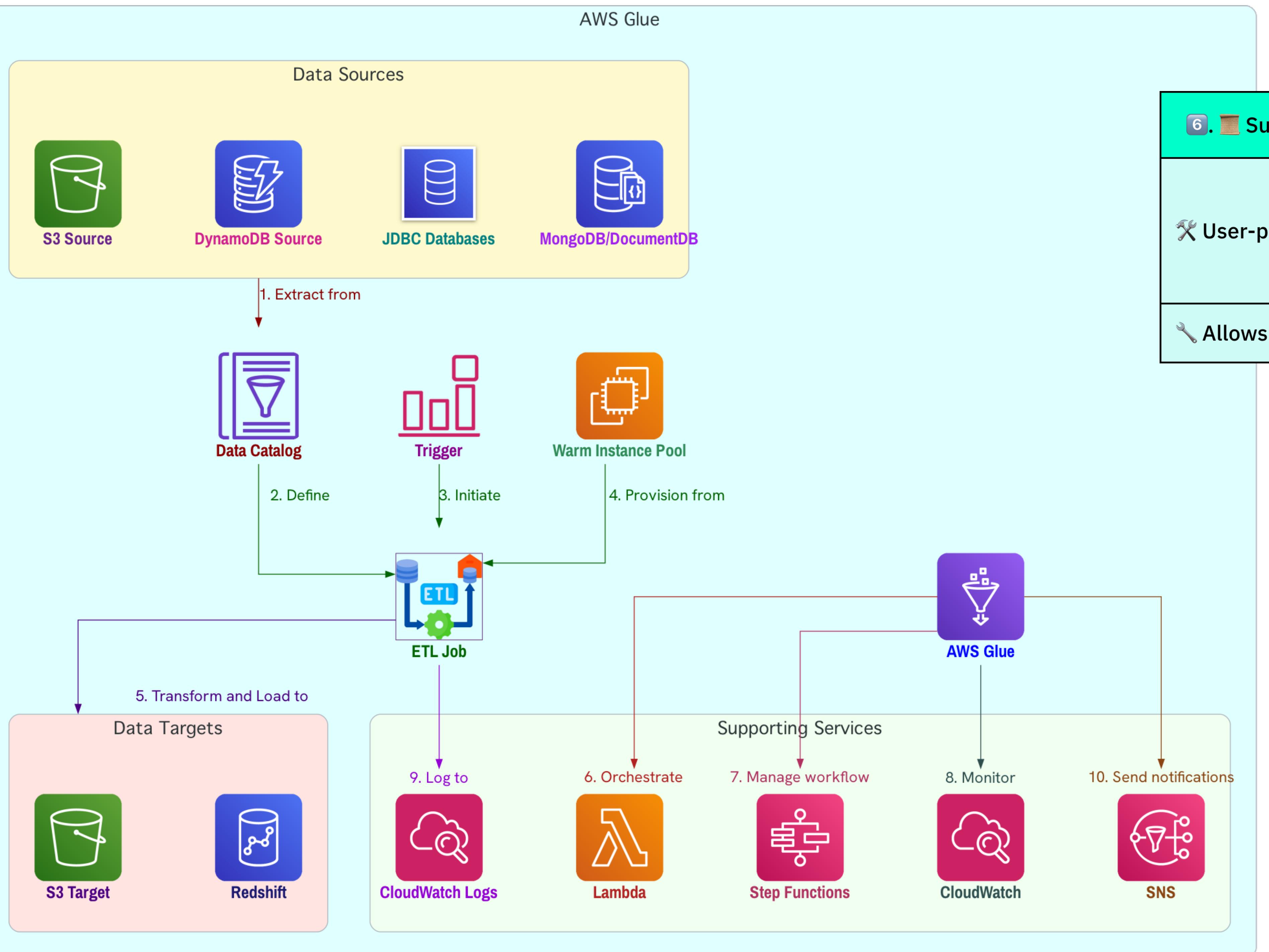
5. ⏰ Schedules and triggers jobs flexibly

Predefined schedules

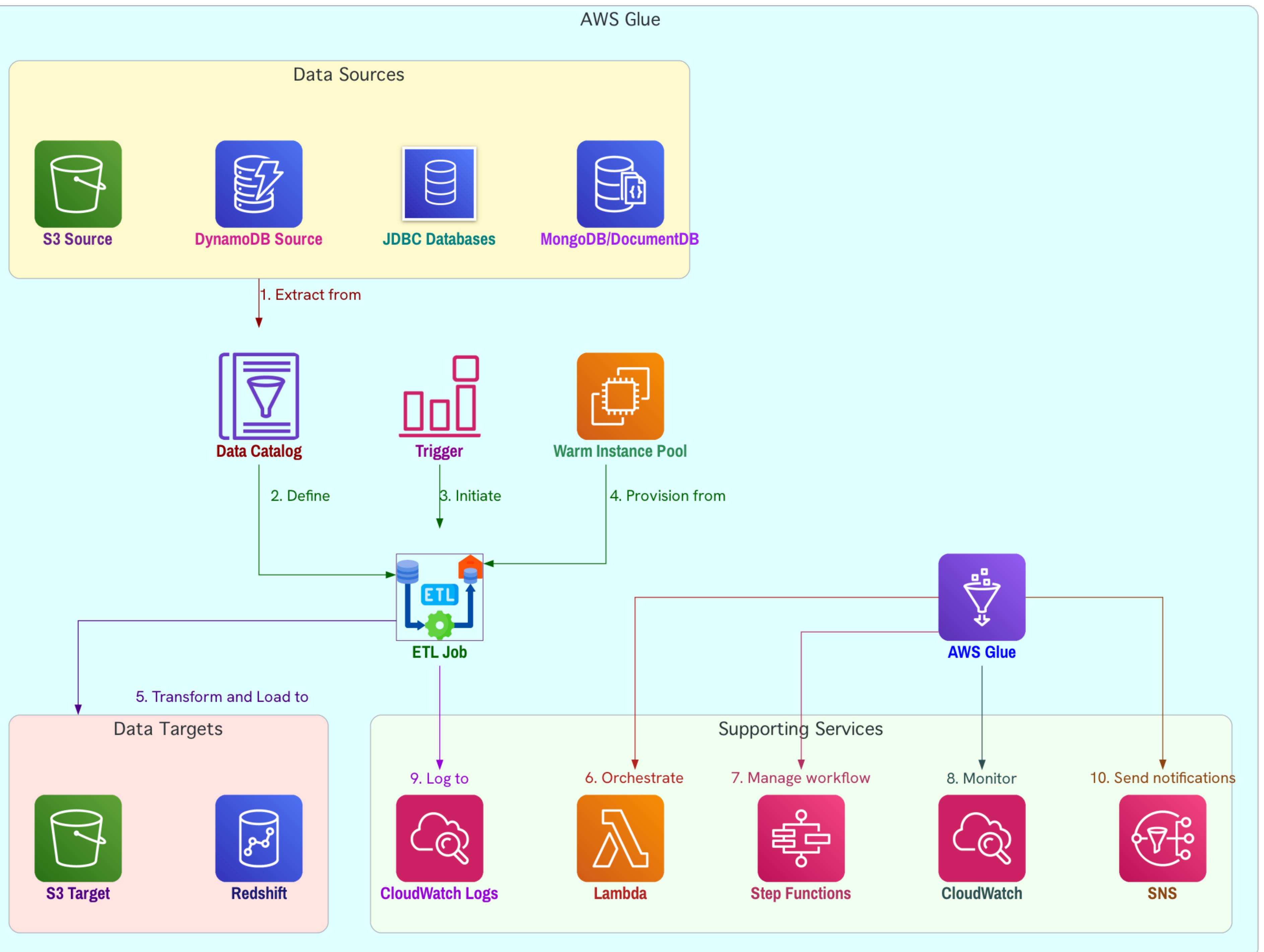
Event-driven triggers

⌚ Supports batch and event-driven architectures

# How it works

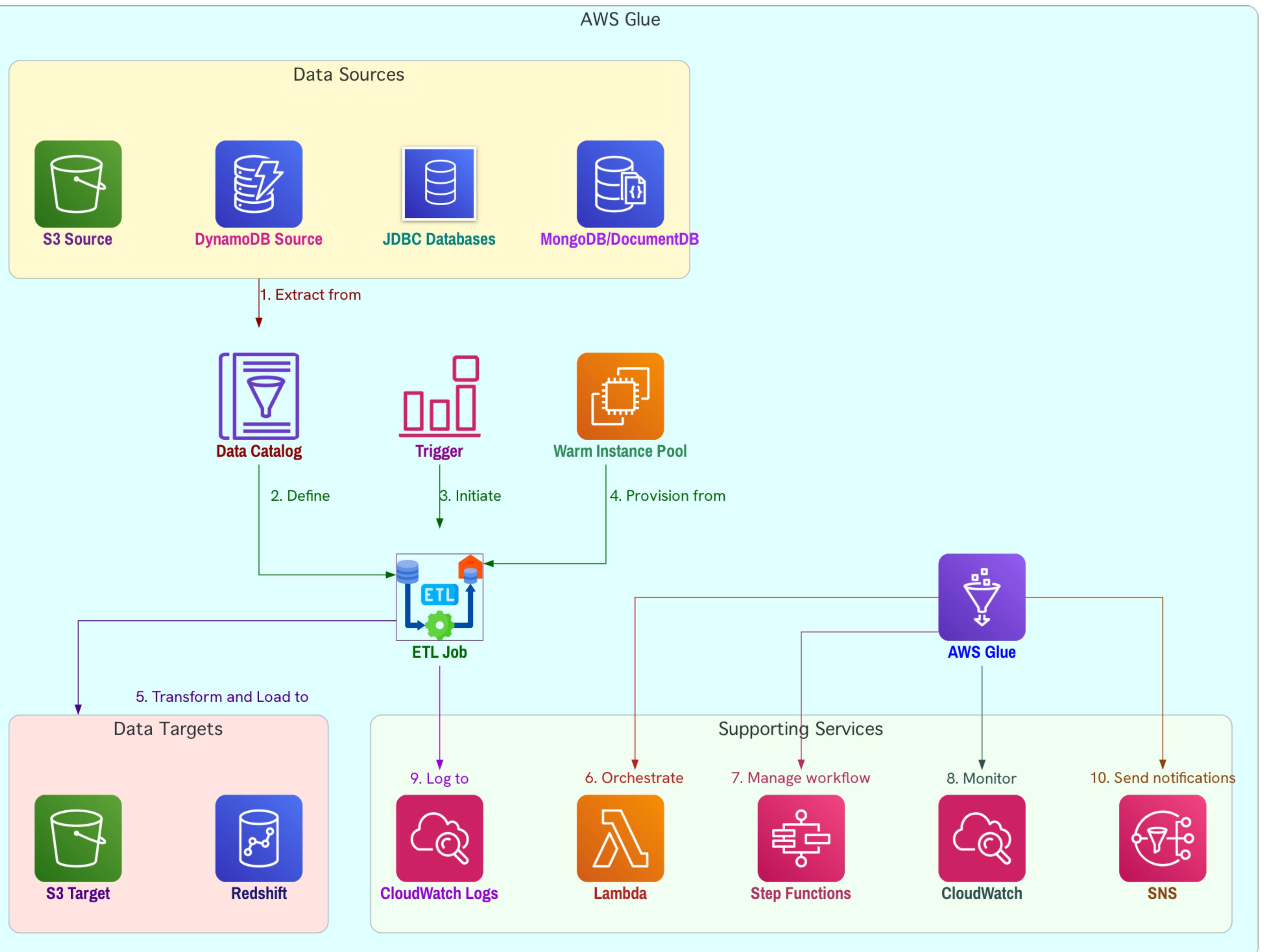


# How it works



- 7. Monitors and logs job performance
- Creates runtime logs
- Generates notifications
- Provides visibility into ETL processes
- Enables quick issue identification
- Facilitates performance optimization

# How it works



8. Streamlines data warehouse and lake creation

Simplifies building process

Supports output stream generation

Enables real-time data processing

Feeds into AWS analytics services

# AWS Integrations

## Data Sources and Destinations



Amazon S3



Amazon DynamoDB



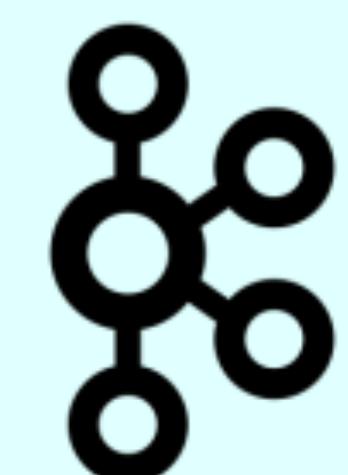
Amazon Redshift



Amazon RDS



Kinesis Data Streams



Apache Kafka



JDBC Databases



MongoDB/DocumentDB

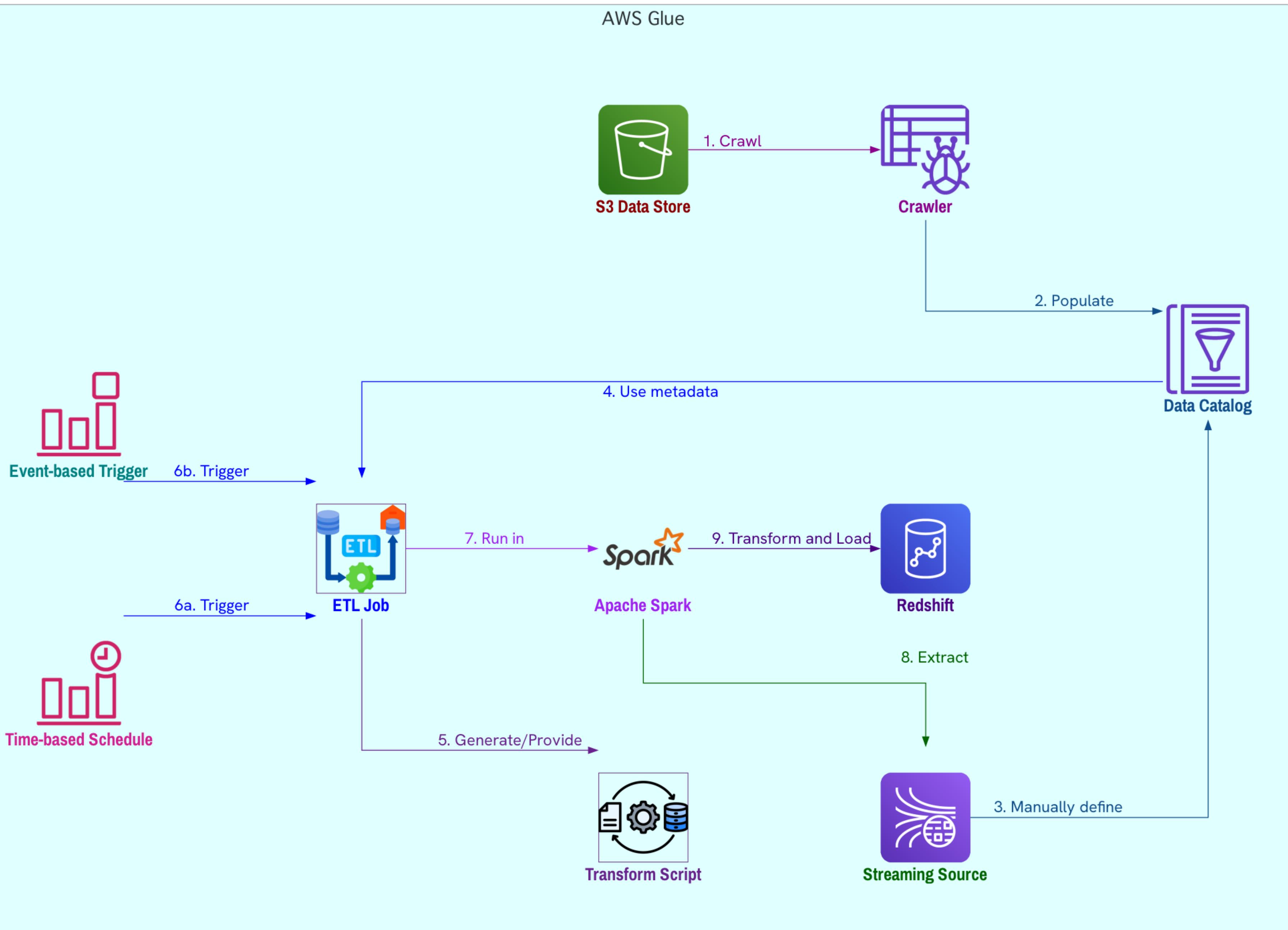


Marketplace Connectors



Apache Spark Plugins

# Architecture of an AWS Glue environment



1. 1 Define ETL jobs for data processing

2. 2 Perform Extract, Transform, Load

3. 3 Move data efficiently

4. 4 Apply necessary transformations

2. 2 Set up crawlers for data store sources

3. 3 Automatically populate Data Catalog

4. 4 Scan data store

5. 5 Create table definitions with metadata

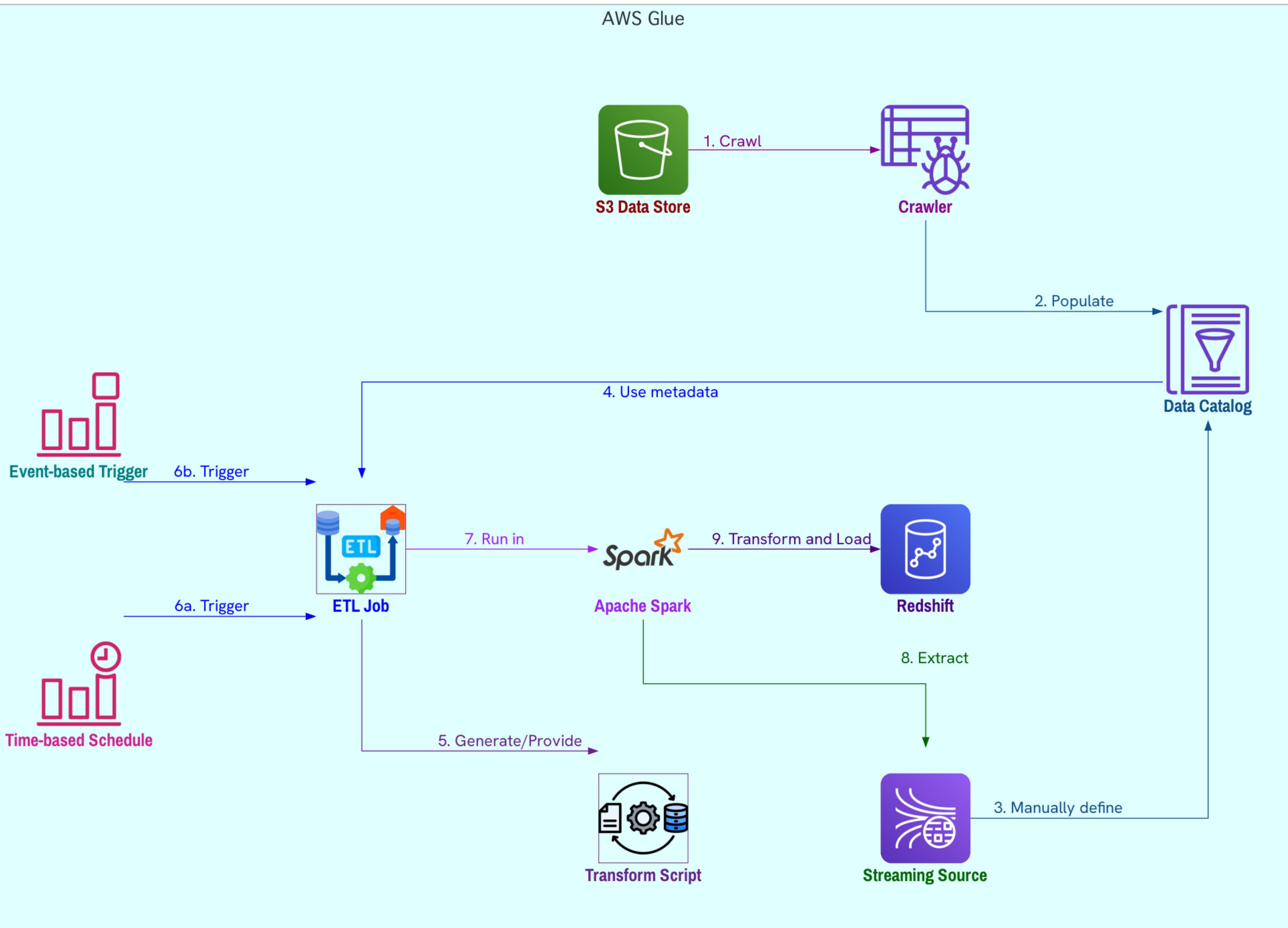
3. 3 Manually define tables for streaming sources

4. 4 Work with streaming data

5. 5 Specify data stream properties

6. 6 Ensure correct data interpretation

# Architecture of an AWS Glue environment



4. 📚 Leverage AWS Glue Data Catalog metadata

🏛️ Central repository for definitions

🔑 Essential for defining ETL jobs

🚩 Provide data structure information

5. 💨 Generate or provide transformation script

🤖 Automatic script generation

👉 Custom script provision

🔧 Complex data transformations

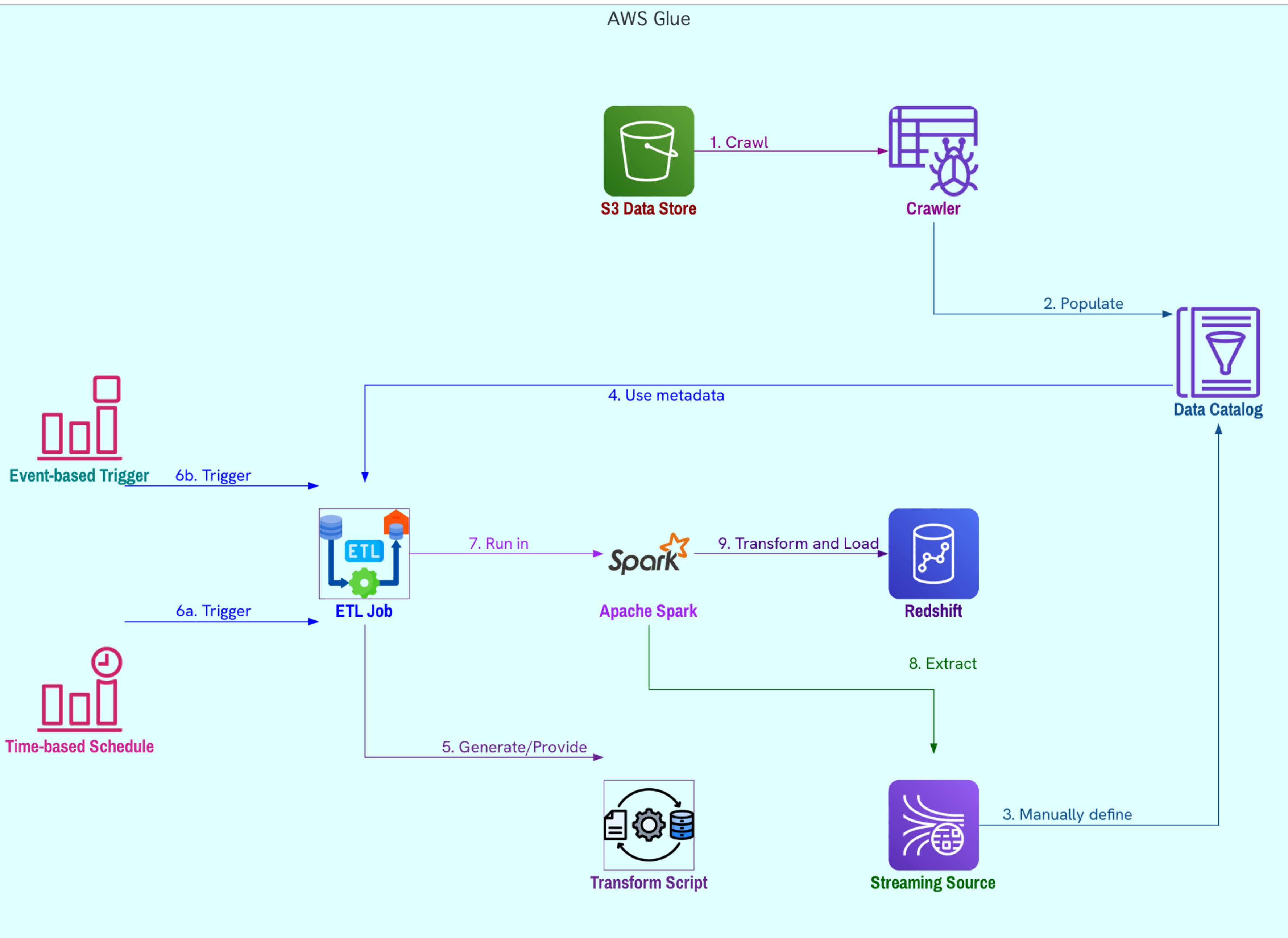
6. ⏳ Configure job triggers

⌚ Time-based triggers

🎭 Event-based triggers

🔧 Flexible scheduling options

# Architecture of an AWS Glue environment



7. Execute job: Extract, Transform, Load

Extract data from source

Apply transformations

Load to target destination

8. Process data in Apache Spark environment

Leverage Spark's capabilities

Efficient, scalable processing

Optimize ETL workloads



**Thanks  
for  
Watching**