# What to expect

## 01
### Interactive data integration

Why interactive data integration is important and hard to do

## 02
### Key AWS Glue innovations

How AWS Glue Interactive Sessions and AWS Glue Studio Notebooks help give a better interactive experience

## 03
### Build interactive applications

How to unlock new paradigms for data integration and Spark development

# Why customers do data integration

**Meaningful insights**

**Increased collaboration**

**Faster decisions**

**Break data silos**

# AWS Glue

**FAST, OPEN, AND SCALABLE
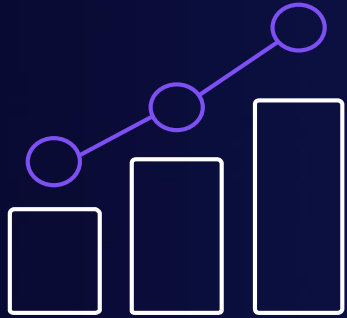DATA INTEGRATION SERVICE**

Integrate data faster

Scale with your growth

No servers or capacity to manage

# Understanding your data assets . . . interactively

**EXPONENTIAL DATA GROWTH**

**DATA FROM NEW SOURCES**

**INCREASINGLY DIVERSE DATA**

# Data interactivity for different personas
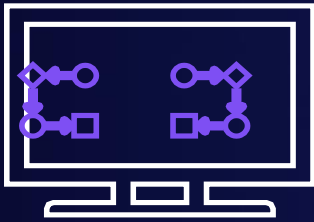
**NO OR LOW CODE**

**DATA ENGINEERS**

**SPARK DEVELOPERS**

# New kind of interactive data integration

**REAL-TIME/ SLA SENSITIVE**

**COMPLEX**

**COST-SENSITIVE**

# Current interactive Spark development

**Complex**
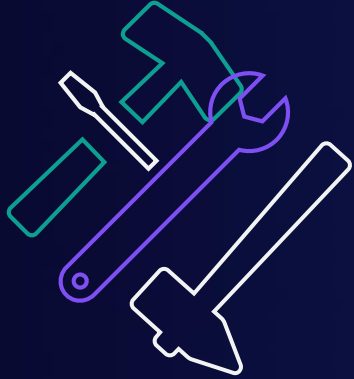
Hard to get started

**Lengthy**

Time to first query
takes 10–30 minutes

**Wasted cost**

Over-provisioning

No cost control
mechanisms

# Current interactive data integration process

**Slow and lengthy**

**Poor feedback mechanism**

**Works only for certain users**

aws Glue

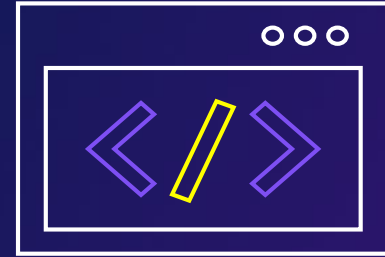[interactive_sessions: ⬇ ]:

# AWS Glue Interactive Sessions
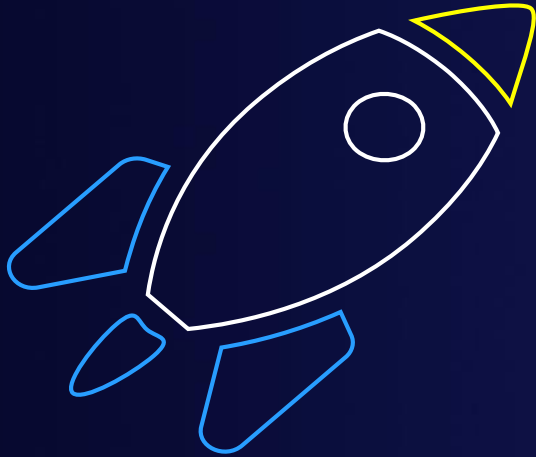
Rapid Spark
development

Pay as you go

Development tool
of your choice

# AWS Glue Interactive Sessions is rapid

Rapid creation of serverless Spark

Execute first code in seconds

Configure and install packages in seconds

Dedicated resources for every session

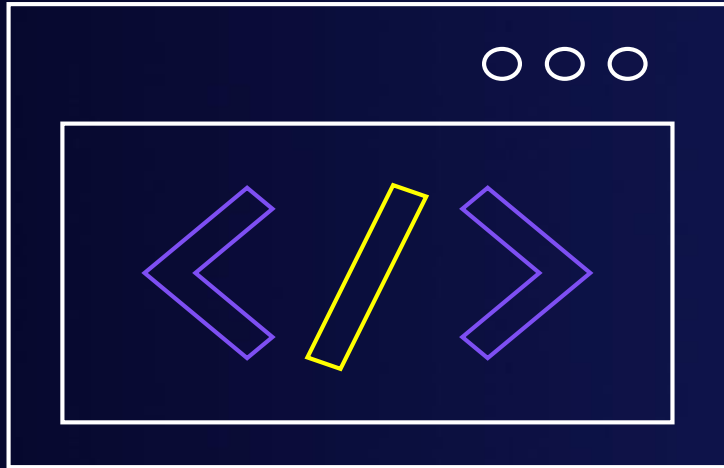# AWS Glue Interactive Sessions is frugal

Pay as you use

Inactive timeout & cleanup

Per second billing with
1-minute minimum

# AWS Glue Interactive Sessions is easy to use

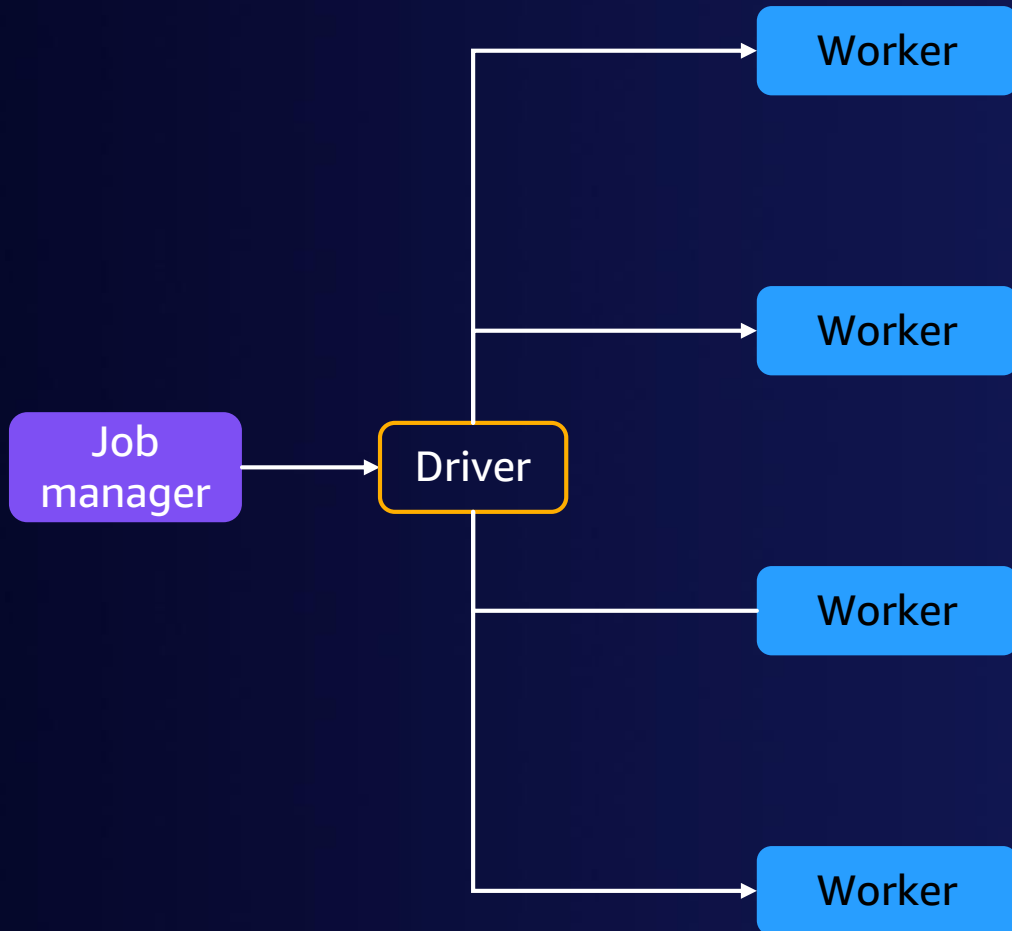No complex setup needed

Live serverless data integration

Works with you favorite IDE

# Same serverless execution engine as Glue jobs

Job manager → Driver → Worker, Worker, Worker, Worker

Jobs are divided into stages

Data is divided into partitions (shards) that are processed concurrently

Job manager schedules task on worker

Glue's execution scales to 1000s of worker

You pay for compute used

# AWS Glue Interactive Sessions is RESTful

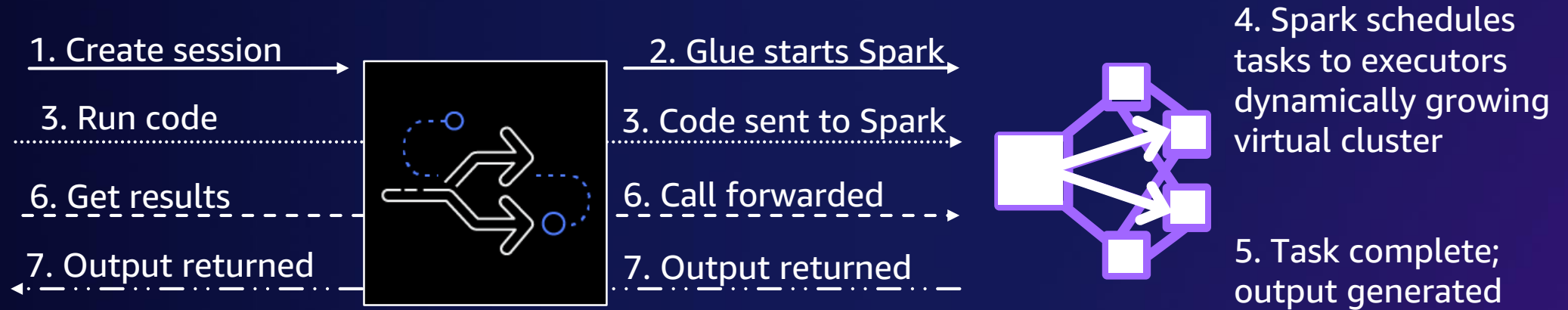Integrate your own applications using APIs

# Built on AWS Glue Execution Engine

1. Create session →

2. Glue starts Spark →

3. Run code ⟶

3. Code sent to Spark ⟶

4. Spark schedules tasks to executors dynamically growing virtual cluster

6. Get results ⟶

6. Call forwarded ⟶

5. Task complete; output generated

7. Output returned ←

7. Output returned ←

In notebook, session starts when first code cell is executed

**a. Warm pool**

AZ1

AZ2

**b. EC2**

# AWS Glue Interactive Sessions API

create_session() ——————— Start a session with your parameters

run_statement() ——————— Submit Spark code to your session

get_statement() ——————— Get returned output of submitted code

stop_session() ——————— Stop the session

# AWS Glue Interactive Sessions for Jupyter

Open-source Jupyter Kernel

Choose your favorite IDE

Jupyter Notebooks & Lab

Visual Studio Code

IntelliJ | PyCharm | DataSpell

# Configure with Magics
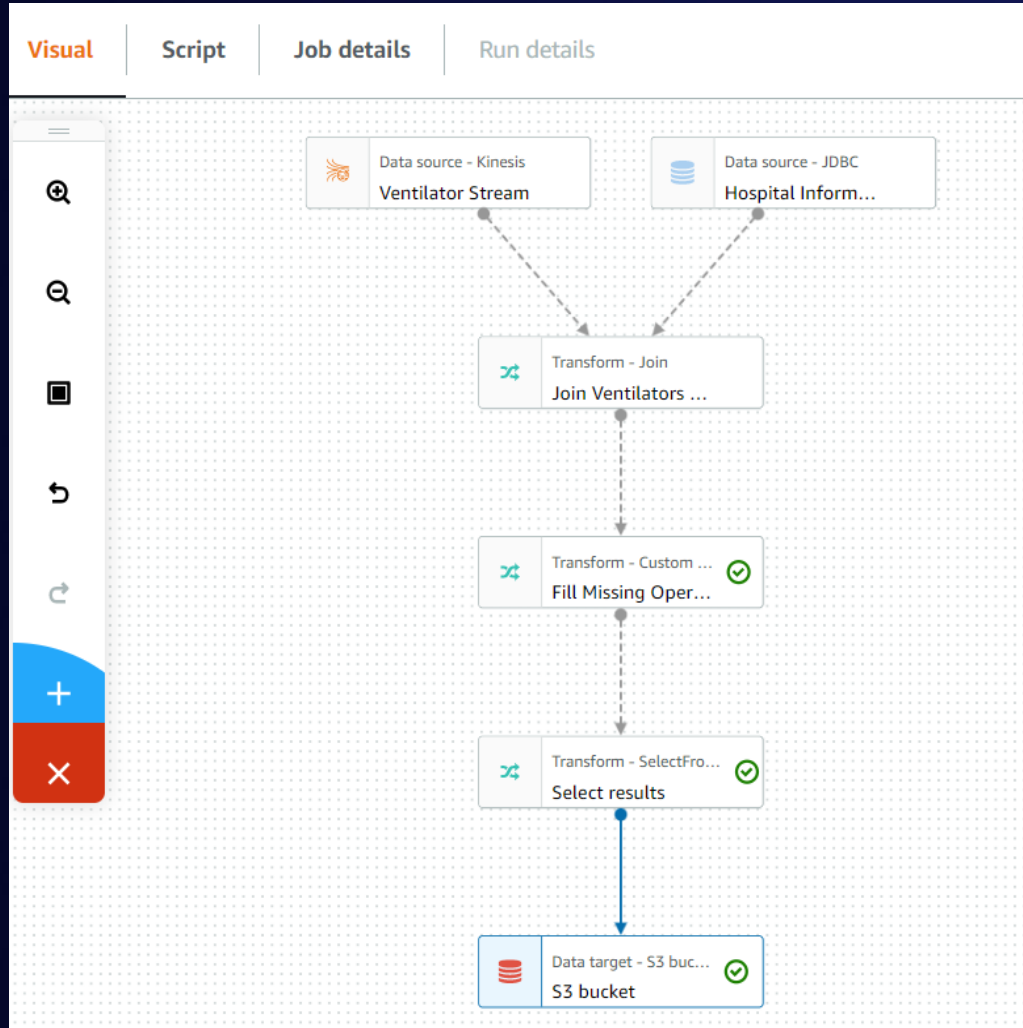
# Config stays with code

```
[ ]:   %streaming
       %worker_type G.2X
       %connection my-vpc
       %%additional_python_modules awswrangler, s3://mybucket/mycustom.whl
       %%configure
         {
           "glue_version": "3.0",
           "number_of_workers": "20",
           "--conf": "spark.serializer=org.apache.spark.serializer.KryoSerializer"
         }
```

```
[ ]:   dyf = glue_context.create_dynamic_frame.from_catalog()
```

# AWS Glue Studio: Visual ETL interface

## MAKES IT EASY TO AUTHOR, RUN, AND MONITOR AWS GLUE ETL JOBS



Author AWS Glue jobs visually without coding

Monitor 1000s of jobs through a single pane of glass

Distributed processing without the learning curve

Advanced transforms though code snippets

# AWS Glue Studio benefits
## Visual ETL in AWS Glue

**Rapid development**
Author ETL jobs faster and with less debugging

**Monitor jobs in one place**
See all your job runs through a single pane of glass

**Streaming and batch ETL**
Use one service to process data from all your sources

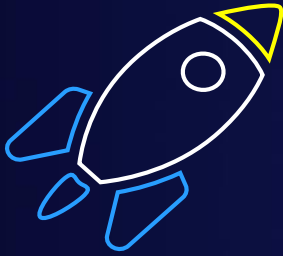**Perform complex transformations**
Use visual transforms or write code snippets to clean and prepare your data for analysis

**Process structured and semi-structured data**
Use Glue studio to handle structured and semi-structured data like IoT logs
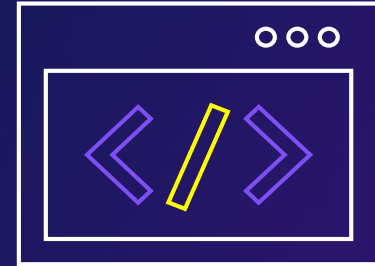
# AWS Glue Studio Notebooks

Serverless

Free

1-click job creation
& scheduling

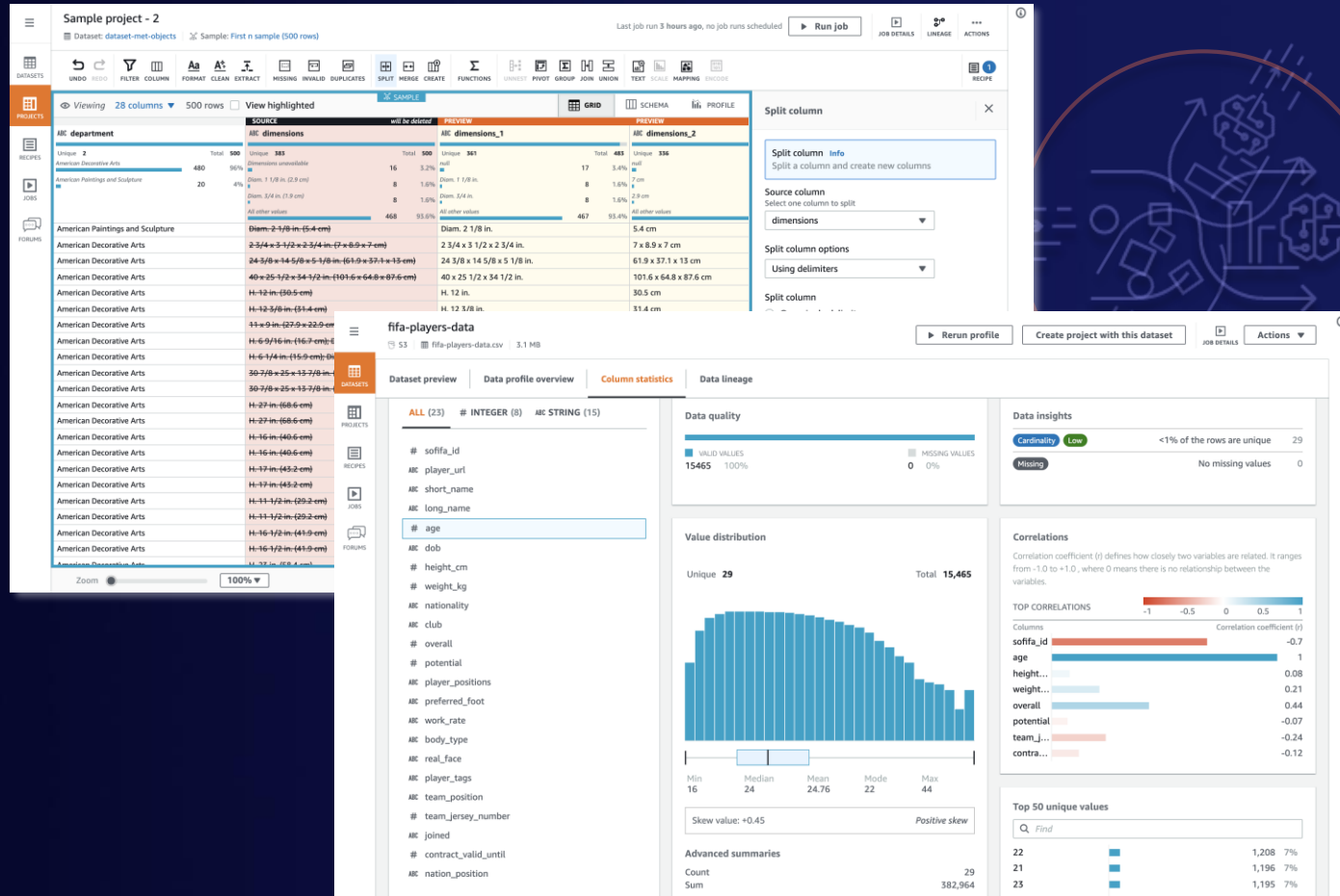# AWS Glue Interactive Sessions applications

### Direct API access to Glue Spark

Including Glue Streaming

### Extend applications into Spark

# AWS Glue DataBrew

## POWERED BY AWS GLUE INTERACTIVE SESSIONS



Chose from 350+ transformations to prepare data

Identify dataset quality with 50+ statistics

Share, reuse, and schedule transforms across users and teams

Powered by Interactive Sessions

# Summary

- Build faster, easier, and more cost-effectively
- End-to-end data integration
- Interact with Apache Spark using your favorite IDE
- Build interactive applications

# What would you build?

# Thank you!

Saroj Yadav

saroyada@amazon.com
https://www.linkedin.com/in/sarojyadav/

Zach Mitchell

mitczach@amazon.com
https://www.linkedin.com/in/zachary-mitchell-ab882853/

aws