re:Invent

NOV. 28 - DEC. 2, 2022 | LAS VEGAS, NV

ANT304

Run serverless Spark workloads with AWS analytics

Noritaka Sekiyama
Principal Big Data Architect, AWS Glue
AWS

Vincent Gromakowski
Principal SA, Analytics
AWS



Agenda

Why Apache Spark on AWS

Use cases across AWS services for running Spark

Focus on AWS Glue and Amazon EMR



Apache Spark is popular





Support multiple workloads in one application



Accelerate the speed of analytics and processing



Highly fault tolerant



Developer friendly, supporting multiple programming languages



Why Spark on AWS



on



Key enabler across multiple AWS services

Provide Spark offering optimized for AWS

Eliminate Spark deployment and update challenges

Deployed on the most comprehensive cloud platform

Supported by the largest partner community

Have serverless deployment options

Unique AWS features to improve resiliency, scale, and reduce costs



Spark on AWS comparison



Amazon EMR Serverless

Use case General purpose big data

When to use Run big data apps with open source frameworks like Spark and Hive

Data Analysts Personas

Data Engineers **Data Scientists**



AWS Glue for Apache Spark

Data integration

Build data lakes, DW, and data pipelines for different workloads (ETL/ELT/Streaming)

> **Data Engineers ETL Developers Data Architects** Data Stewards



Cost model



Amazon EMR Serverless



AWS Glue for Apache Spark

vCPU

Per vCPU per hour

RAM

Per GB per hour

Per DPU per hour (4 vCPU, 16 GB RAM)

Storage

Per storage GB per hour

Amazon EMR Serverless



Amazon EMR Serverless

SERVERLESS GENERAL-PURPOSE BIG DATA PLATFORM



Run big data applications on Spark and Hive



60-day Spark release cadence



"Code first" approach



Petabyte scale



Bring your own customizations



Transactional support



Use cases







Big data analytics

- Large-scale data processing and what-if analysis using statistical algorithms and predictive models
- Specific or latest Spark versions
- Custom environment setup

Interactive workloads

 Start immediately on the workers without any delays (pre-init capacity)

Lift and shift

- Migrate existing Hadoop applications from onpremise
- Migrate from EMR on EC2 or EMR on EKS to serverless



Benefits



Advanced configuration control and flexibility

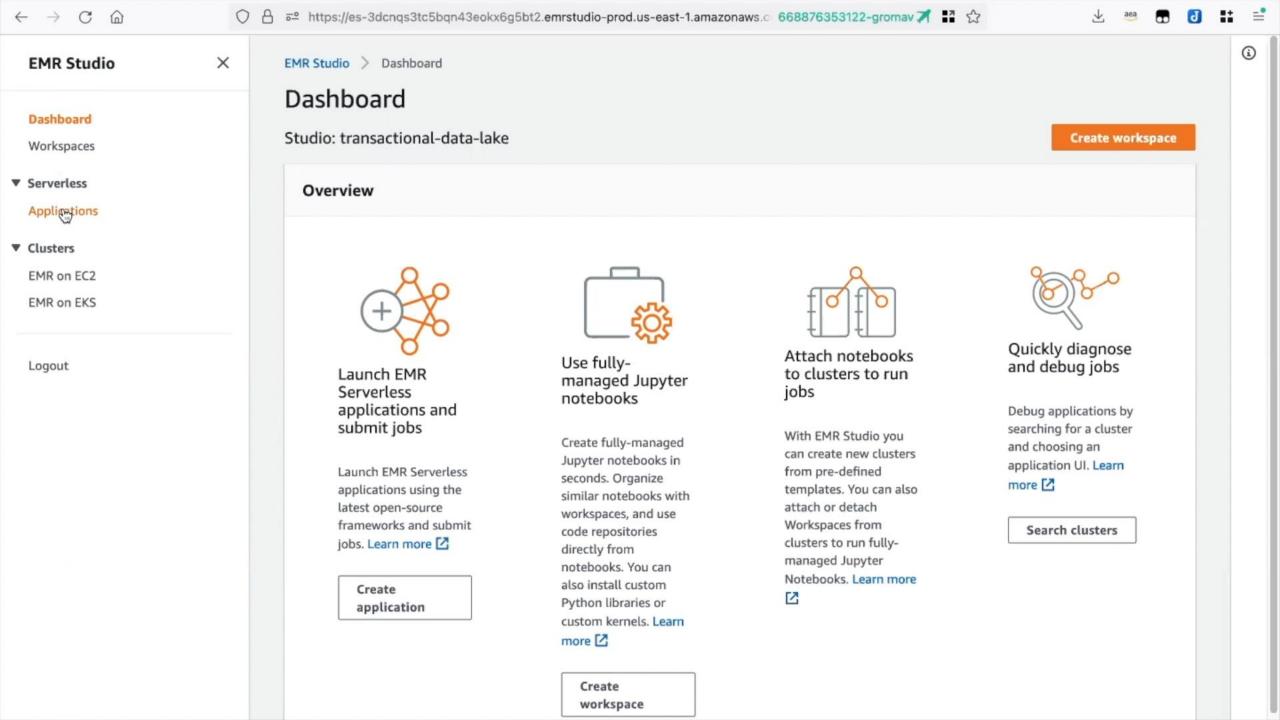


Latest Spark versions



70–90 second start time





Advanced Spark tuning

Override

- Spark properties defined in the EMR Serverless application (driver and executor conf)
- Any advanced Spark configuration including JVM tuning, offheap, etc. for performance



Python versions and dependencies

- Package a Python virtual env from a Docker image
- Optionally, use a different Python version
- Upload the archive to Amazon S3
- Load the package in the EMR Serverless job

```
aws emr-serverless start-job-run \
--application-id $APPLICATION_ID \
--execution-role-arn $JOB_ROLE_ARN \
--job-driver '{

"sparkSubmit": {

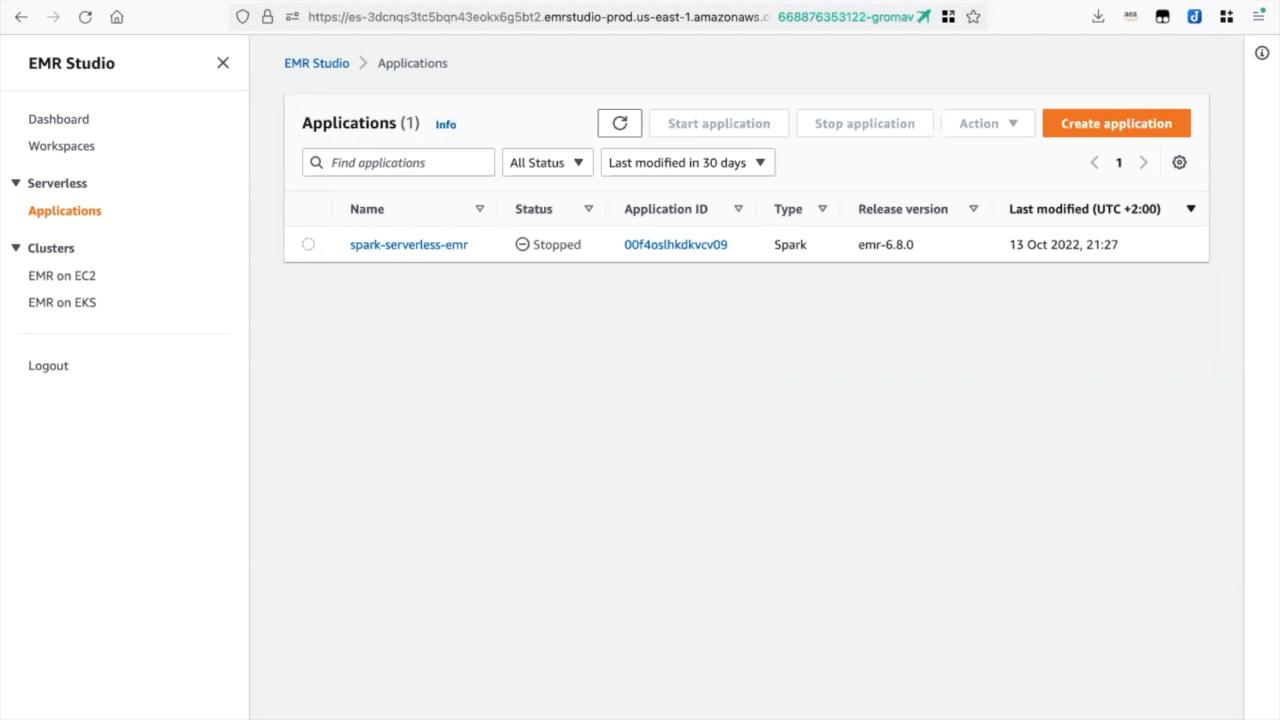
"entryPoint": "s3://'${$3_BUCKET}'/code/pyspark/ge_profile.py",

"sparkSubmitParameters": "--conf spark.archives=s3://'${$3_BUCKET}'/artifacts/pyspark/pyspark_ge.tar.gz#environment

--conf spark.emr-serverless.driverEnv.PYSPARK_PYTHON=./environment/bin/python

--conf spark.emr-serverless.executorEnv.PYSPARK_PYTHON=./environment/bin/python"
```

```
RG PYTHON VERSION=3.10.6
FROM --platform=linux/amd64 amazonlinux:2 AS base
ARG PYTHON_VERSION
# Install Python 3.10.6 - Note that python 3.10 requires OpenSSL >= 1.1.1
RUN yum install -y qcc openssl11-devel bzip2-devel libffi-devel tar gzip wget make && \
    wget https://www.python.org/ftp/python/${PYTHON_VERSION}/Python-${PYTHON_VERSION}.taz && \
   tar xzf Python-${PYTHON_VERSION}.tgz && \
   cd Python-${PYTHON_VERSION} && \
   ./configure --enable-optimizations && \
   make install
ENV VIRTUAL_ENV=/opt/venv
RUN python3 -m venv $VIRTUAL_ENV --copies
RUN cp -r /usr/local/lib/python3.10/* $VIRTUAL_ENV/lib/python3.10/
# Ensure our python3 executable references the virtual environment
ENV PATH="$VIRTUAL_ENV/bin:$PATH"
# Upgrade pip (good practice), install dependencies and venv-pack
RUN python3 -m pip install --upgrade pip && \
   python3 -m pip install venv-pack==0.2.0 great_expectations==0.15.6
   venv-pack -o /output/pyspark_${PYTHON_VERSION}.tar.gz --python-prefix /home/hadoop/environment
COPY --from=base /output/pyspark_${PYTHON_VERSION}.tar.gz /
```



AWS Glue



AWS Glue





Focus on data

Low maintenance serverless solution



Powerful open source engine

No lock in, support from wide innovative eco system



Scale on demand

Avoid licensing cost and infrastructure idle time



All in one

Support all your users personas and workloads



AWS Glue overview

SERVERLESS DATA INTEGRATION SERVICE

Connectors



Data warehouse



Data lakes



Marketplace



NoSQL



Streams

Author



Visual



Notebook



Built-in transformations

Serverless infrastructure



IDE

Operationalize



Workflow



Monitoring



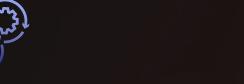
Data Catalog

Data Management

Sensitive data detection

Engines

Schedule



Choice of data integration engines



Compliance and security

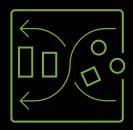


Data quality

Stream data processing



Use cases







Various data sources

- Built-in connectors (S3, Amazon Redshift, Amazon RDS/JDBC, Amazon DynamoDB, Amazon Kinesis/Kafka, MongoDB, etc.)
- Connector marketplace

Data management

- Incremental processing
- 250+ transformations
- Sensitive data detection
- Data lake format support (Hudi, Delta Lake, Iceberg)

Low-code job authoring

- AWS Glue Studio visual editor with visual transformations and data preview
- Automatic code generation



Benefits



Moderate level of configurations and flexibility

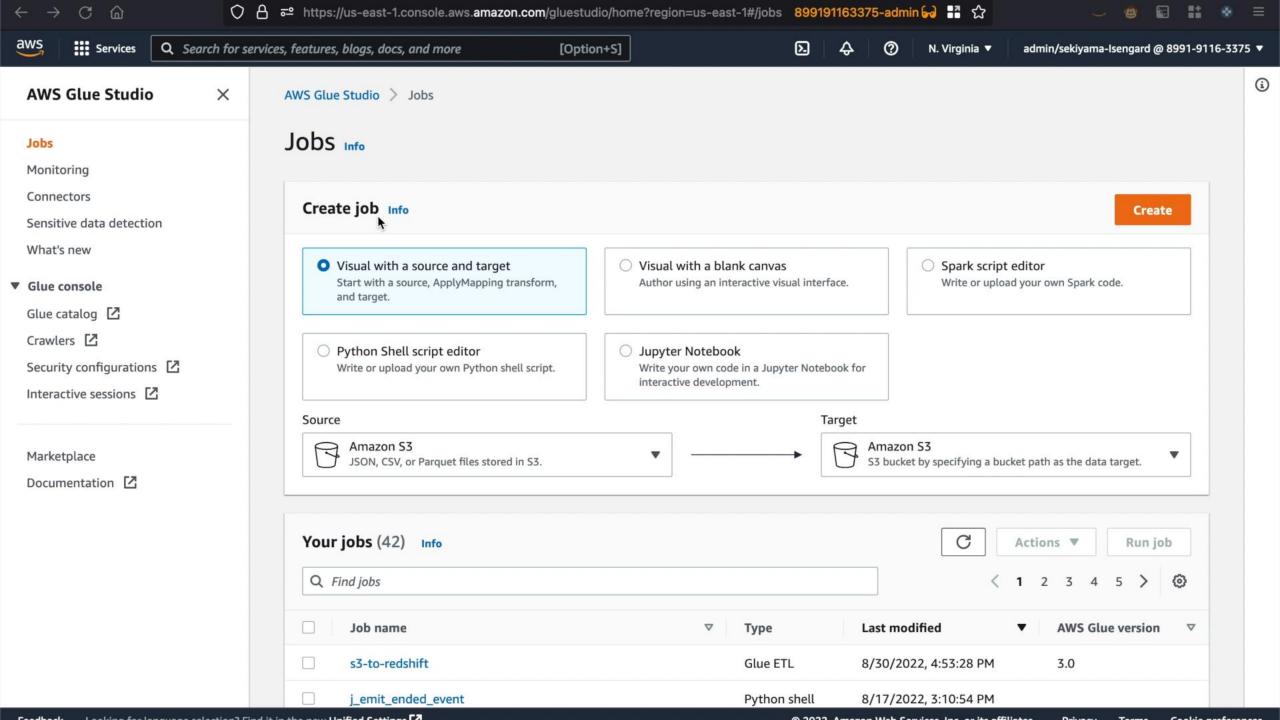


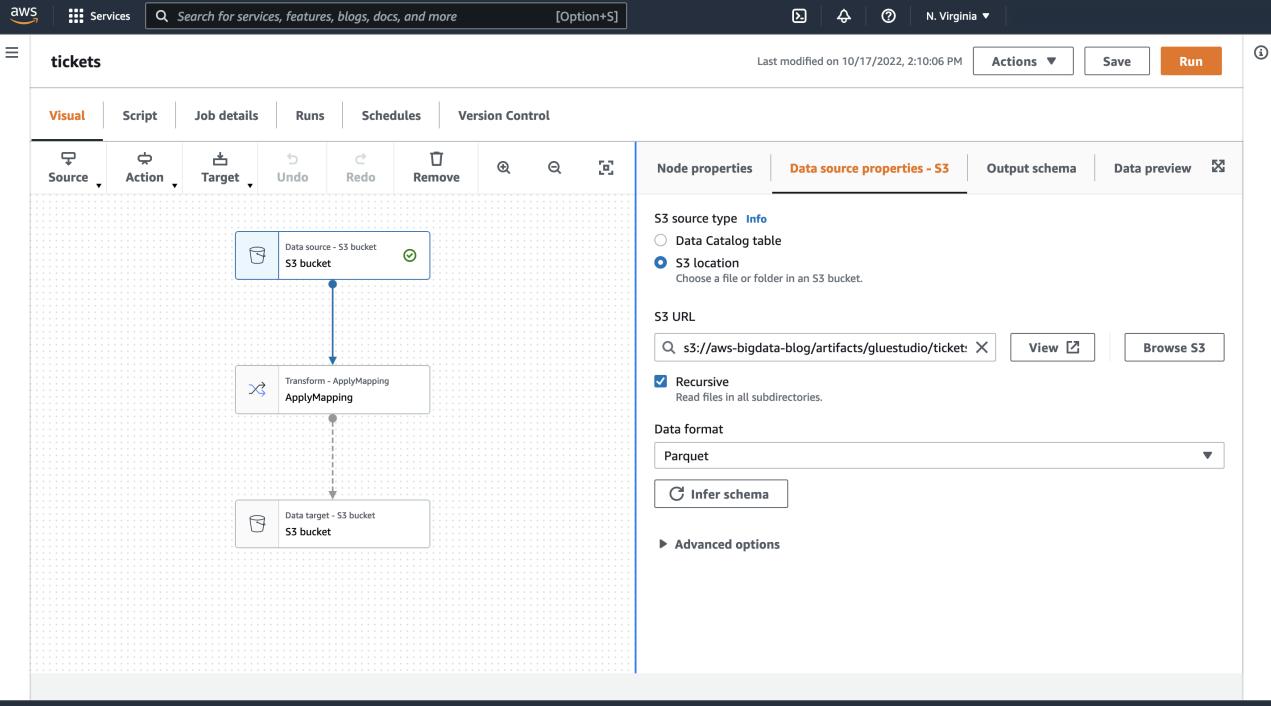
Many options to write and run code; IDE, notebook, GUI

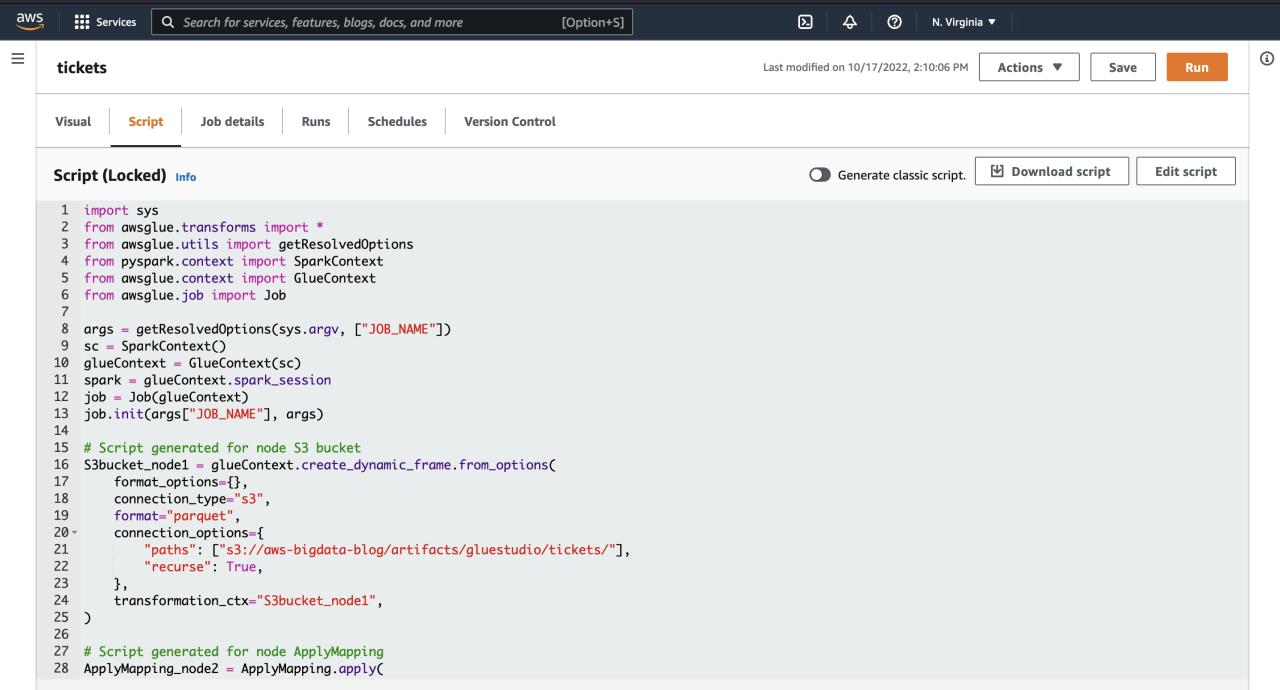


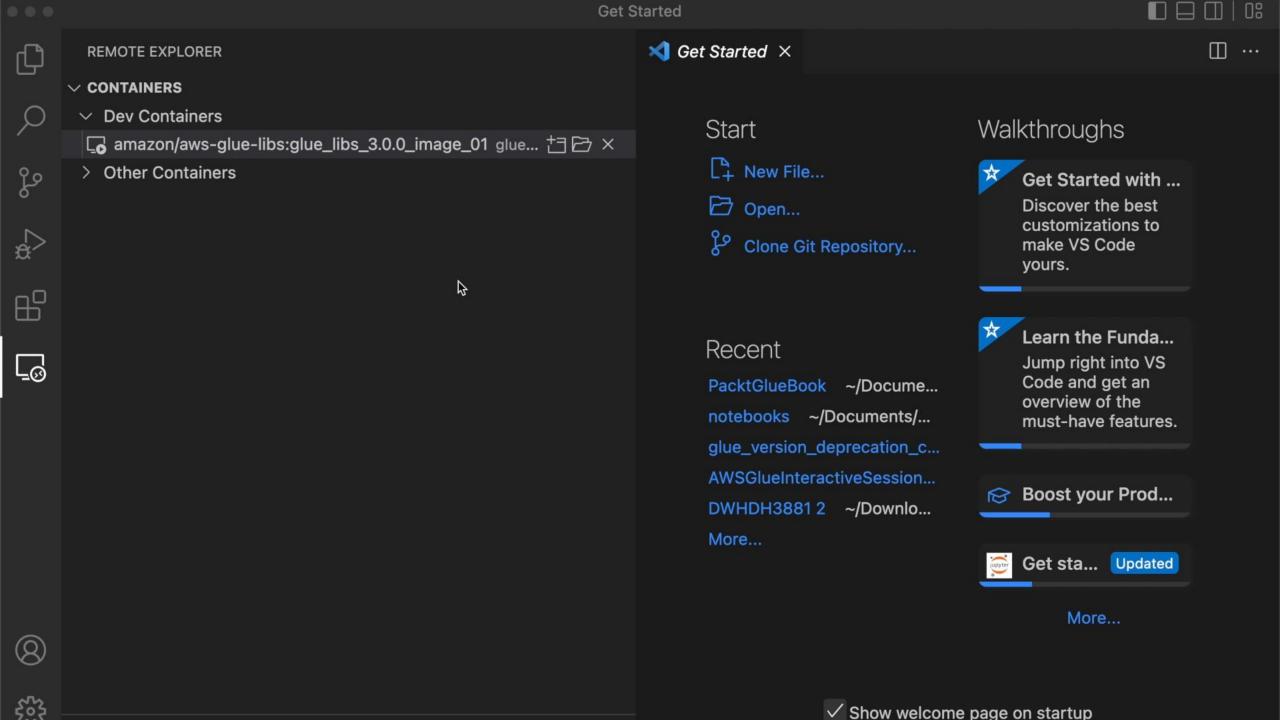
5–10 seconds of start time

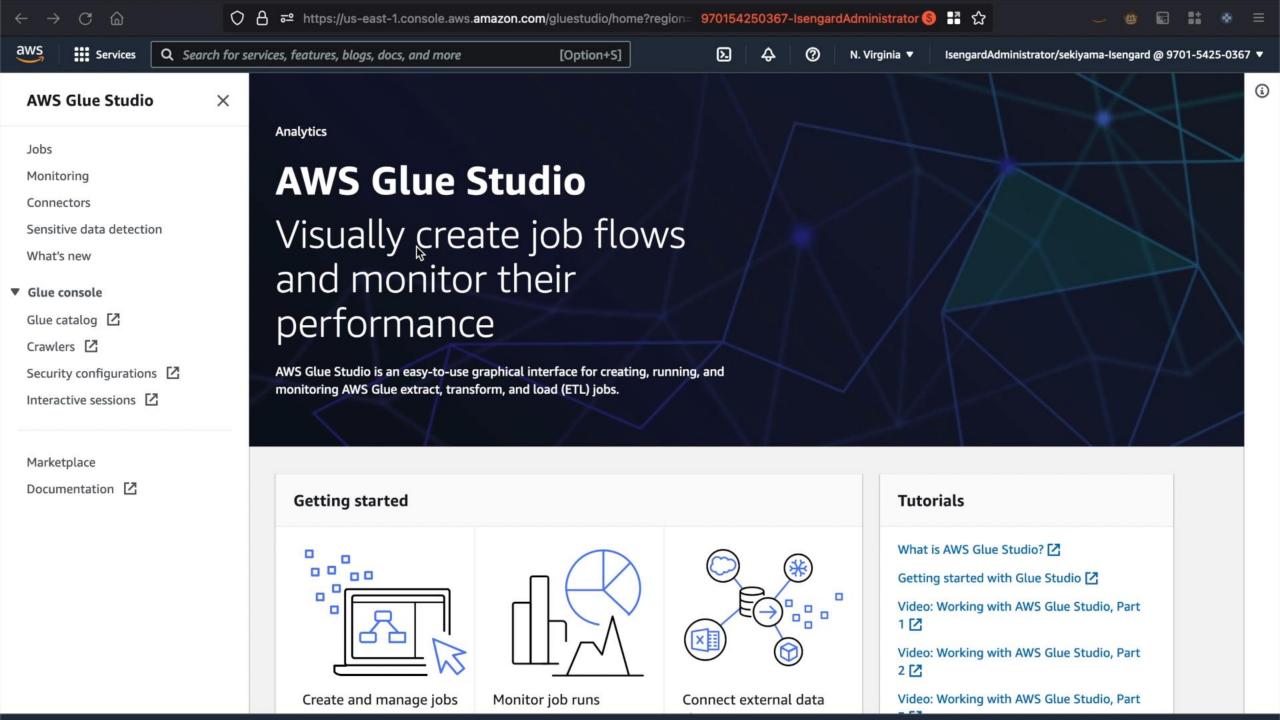




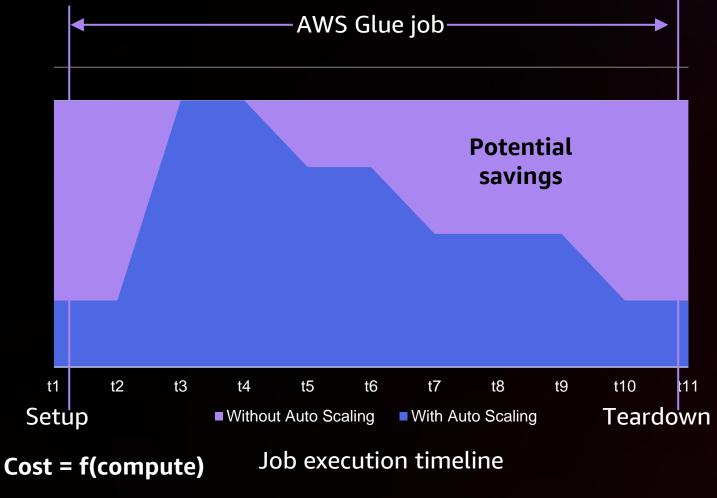








AWS Glue Auto Scaling





Reduce cost



Automate



Simplify capacity planning



AWS Glue Flex jobs



Standard

class

~1 min job start times Predictable job latencies

Enables micro-batching Latency-sensitive workloads



Flex

class

Up to 34% cost savings

Cost effective for one-time data-load workloads



Extra library management in AWS Glue

Python

packages

pip integration to install additional libraries

Key: --additional-python-modules

Value: boto3==1.24.89

or

Key: --additional-python-modules
Value: s3://path_to_your_whl_file

Java/Scala

JAR files

Extra JAR support to install additional libraries

Key: --extra-jars

Value: s3://path_to_your_jar_file



Q&A



Thank you!

Noritaka Sekiyama



sekiyama@amazon.com

Vincent Gromakowski gromav@amazon.com



Please complete the session survey in the mobile app

