ANT202

# What's new with Amazon EMR

Vincent Gromakowski

Principal Solutions Architect, Analytics
Amazon Web Services

aws

# Amazon EMR

## BIG DATA ANALYTICS USING OPEN-SOURCE FRAMEWORKS: APACHE SPARK, PRESTO, TRINO, HADOOP, HIVE, HBASE & FLINK

### Differentiated performance for runtimes

Performance-optimized runtime for popular frameworks like Spark, Hive, Presto, and Flink with 100% open-source API compatibility

### Latest open-source features

New open-source features available within 30 days of release in open source

### Best price-performance for big data analytics

Reduce cost using Amazon EC2 Spot, Amazon EMR managed scaling, and per-second billing

### Self-service data science

New!

Data science IDE with EMR Studio and deep integration with Amazon SageMaker Studio provides ability to use open-source UX and frameworks to build, visualize, and debug applications

### Run workloads on Amazon EC2, Amazon EKS, or on premises

New!

EMR provides flexibility to run big data workloads on EC2, EKS, and on premises with AWS Outposts

### S3 data lake integration

New!

Fine-grained access controls with AWS Lake Formation and Apache Ranger, and integrations with Apache HUDI to enable Amazon S3 data lake use cases

# Best price-performance for big data analytics

aws

# Differentiated Spark runtime performance

**Over 3x faster than standard Apache Spark 3.0 in derived TPC-DS 3 TB benchmark**

**Takes advantage of AWS-native Graviton2 instances to provide the best performance**

**100% compliance with open-source APIs makes moving applications to EMR easy**

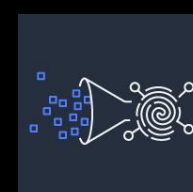**Performance improvements are enabled by default**

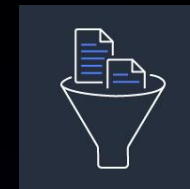Dynamic-sized executors

Adaptive join selection

Dynamic pruning of data columns

Operator optimization

Early worker allocation

Intelligent filtering

Parallel/async initialization

Redundant scan elimination

Data pre-fetch

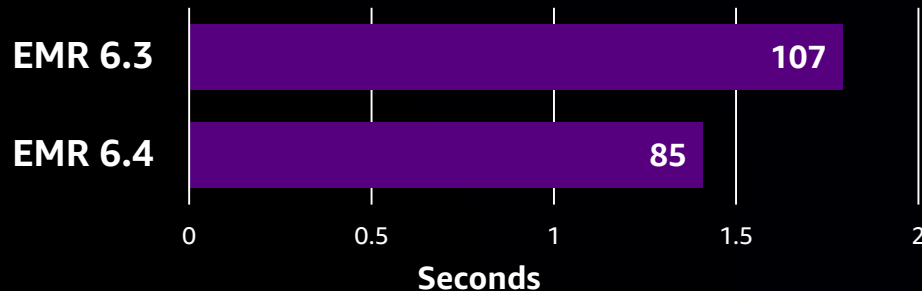Broadcast join w/o statistics

Stats inference

Optimized metadata fetch

aws

# Amazon EMR runtime for Apache Hive

**1.25X FASTER PERFORMANCE WITH APACHE HIVE 3.1.2 ON EMR 6.4**

## Geometric mean of 98 derived query runtimes
**(lower is better)**

| | |
|---|---|
| EMR 6.3 | 107 |
| EMR 6.4 | 85 |

0    0.5    1    1.5    2

**Seconds**

## Apache Hive 3.1.2 on EMR 6.4 vs. EMR 6.3*

*Based on TPC-DS 3 TB benchmarking running 16 node M5.8xlarge cluster*

## EMR's performance-optimized Apache Hive runtime

### Best performance

- **1.25x faster on geometric mean**
- **Up to 2x improvement on individual queries**
- Improves query planning time for AWS Glue Data Catalog
- Improves query execution time for ORC data from Amazon S3

**100% compliant with open-source Apache Hive APIs**

# Additional capabilities to reduce costs

WITH EMR, YOU CAN DO WAY MORE WITH WAY LESS



## Performance optimizations

- Runtime improvements
- Transactions in data lakes

## Compute optimizations

- Graviton instances
- Spot instances
- Instance fleets

## Cluster management

- Managed scaling
- Cluster auto-termination

aws

# Improvements in cluster startup times

STARTING OR SCALING AN EMR ON EC2 CLUSTER IS NOW 35% FASTER

Task nodes provisioned alongside core and primary nodes

Have ready-to-use proxy instances available to improve start time for cluster launched in private subnet

Optimized retry policies for EC2 throttling

aws

# AWS **Graviton2** instances have the best price-performance within their instance families

**REALIZE UP TO 30% BETTER PRICE-PERFORMANCE WITH GRAVITON2 INSTANCES**

**12%–16% improvement in performance** compared to M5 instance types

**20% lower cost** vs. same-sized comparable M5 instances

**Up to 30% better price-performance**

aws

# Amazon EC2 Spot Instances

ACCELERATE COMPUTE FOR LESS

**Low, predictable prices**

Up to 90% discount over On-Demand prices

**Faster results**

Increase throughput up to 10x while staying in budget

**Easy to use**

Launch through AWS services (ex. ECS, EKS, AWS Batch, EMR) or integrated third parties

aws

# Benefits of EMR with Spot Instances

EMR IS THE BEST FIT FOR RUNNING BIG DATA WORKLOADS USING SPOT INSTANCES

### Accelerate compute

Run parallel tasks on a multitude of instance types running Spark, Hive, Flink, or Presto

### Further reduce costs

Access instances at up to a 90% discount vs. On-Demand

### Build for scale

Quickly ramp up short-lived but massive data jobs by scaling compute and storage independently

aws

# Amazon EMR on Spot Instances

**REDUCE TIME-TO-INSIGHT WITH HYPER-PARALLELIZED WORKLOADS**



# parallelized nodes

Time

Job running time: **10 hours**

# parallelized nodes

Time

Job running time: **1 hour**

aws

# Scale up cluster with Spot Instances

10-node cluster running for 14 hours
**Cost = 1.0 * 10 * 14 = $140**

Add 10 more nodes on Spot

**20-node cluster running for 7 hours**
**Cost = 1.0 * 10 * 7 = $70**
**= 0.5 * 10 * 7 = $35**

**Total $105**

**50% less runtime (hours: 14 → 7)**
**25% less cost (dollars: 140 → 105)**

# Use a mix of On-Demand and Spot Instances for different scenarios

| Scenario | Leader node | Core nodes | Task nodes |
|---|---|---|---|
| Long-running clusters and data warehouses | On-Demand | On-Demand or instance-fleet mix | Spot or instance-fleet mix |
| Cost-driven workloads (without a production SLA) | Spot | Spot | Spot |
| Data-critical workloads | On-Demand | On-Demand | Spot or instance-fleet mix |
| Application testing | Spot | Spot | Spot |

# Amazon EMR on Spot Instances

## EMR IS THE BEST PLACE TO RUN BIG DATA WORKLOADS USING SPOT CAPACITY

Save 75%–90% on compute with Spot Instances

Low, predictable prices

Minimal interruptions, <5%

No bidding

Salesforce has observed an increase in iterative job performance and saved 80% vs. On-Demand pricing

*"With AWS, we can manage flexible capacity changes, contain overall costs on daily compute tasks, and manage overall infrastructure growth."*

Roopak Gupta

Vice President, Software Engineering, Salesforce DMP

aws

# Managed scaling feature overview

Constantly improving EMR managed algorithm that gives you a fully managed experience

High-resolution metrics enabled with managed scaling

Only min/max cost constraints configurations required

More data points and faster reaction time than auto scaling

Save 20%–60% of costs

aws

# Managed scaling enhancements

**NEW ENHANCEMENTS ENABLED BY DEFAULT TO FURTHER REDUCE COSTS AND SPOT INTERRUPTIONS**

Capacity awareness in instance groups enabled by default for all supported EMR versions

Integrated with real-time EC2 Spot capacity metrics to scale the right task group based on instance pool depth

Shuffle awareness enabled by default from EMR 6.4

Ensure nodes with active shuffle data are not scaled down

Support for PrestoDB and Trino available from EMR 6.4

Sign up for preview access via aws-support@amazon.com

"Acxiom uses Spark on Amazon EMR on Spot Instances to run 3 trillion inferences in less than 15 hours. By using Amazon EMR, we could utilize Spot compute capacity across the entire AWS Region and speed up the run time of our inference pipeline that typically took 11–15 days every month to under 15 hours."

Varadarajan "Raj" Srinivasan
Sr. Director, ML Engineering and Data Science, Acxiom

# Self-service data science with EMR Studio and Amazon SageMaker Studio

aws

# EMR Studio

**Single sign-on integration with IdP**

**Fully managed Jupyter Notebooks**

**Integrated with Git Repositories**

**Simplified debugging with Spark UI and YARN UI**

**Browse, create, or delete EMR clusters**

**Run Notebooks in workflows using APIs**

**Run interactive data analysis using EMR or EKS clusters**

# Data science and engineering workflows

# Data science and engineering workflows

# Log in to EMR Studio without logging into the AWS Management Console



**Data scientist**

**Data engineer**

Log in with corporate identity using AWS Single Sign-On

# EMR Studio gives you a fully managed notebook



**Workspaces help organize notebooks**

**Workspaces share similar properties**

**Fully managed Jupyter Notebooks**

**Write Python, R, PySpark, Scala**

# Workspace: Single IDE for interactive data analysis

## WITH CURATED LIST OF EXTENSIONS TO ENHANCE THE JUPYTER EXPERIENCE

File browser

**Attach/detach EMR clusters**

Kernel management

Git operations

**Workspace Git link**

**Sample notebooks**

Table of contents

**SQL Explorer**

**Collaboration**

COMING SOON

COMING SOON

---

File    Edit    View    Run    Kernel    Git    Tabs    Settings    Help

**Controls Panel**
Manage features and extensions

**Launcher Tab**
Launch notebooks or extensions

Notebook

| P | P | S | S |
|---|---|---|---|
| Python 3 | PySpark | Spark | SparkR |

Console

| P | P | S | S |
|---|---|---|---|
| Python 3 | PySpark | Spark | SparkR |

Other

| $_ | | {...} | | M | | | |
|---|---|---|---|---|---|---|---|
| Terminal | Collaboration | Notebook Examples | Text File | Markdown File | Python File | Show Contextual Help | SQL Explorer |

# SQL Explorer integrated in Jupyter



**Data catalog browser**

**Multiple SQL editors**

# Data science and engineering workflows

# Simple to load custom libraries and kernels



**Install notebook-scoped libraries with PySpark kernel**

**Install additional Python libraries and kernels on the leader node of the cluster**

# Data science and engineering workflows

# Simple to connect to code repositories



**Connect to AWS CodeCommit, GitHub, and Bitbucket**

**Select existing or add new Git repositories**

# Collaborate in real time

**Enable Workspace collaboration**

**Invite collaborators to your Workspace**

# Collaborate in real time
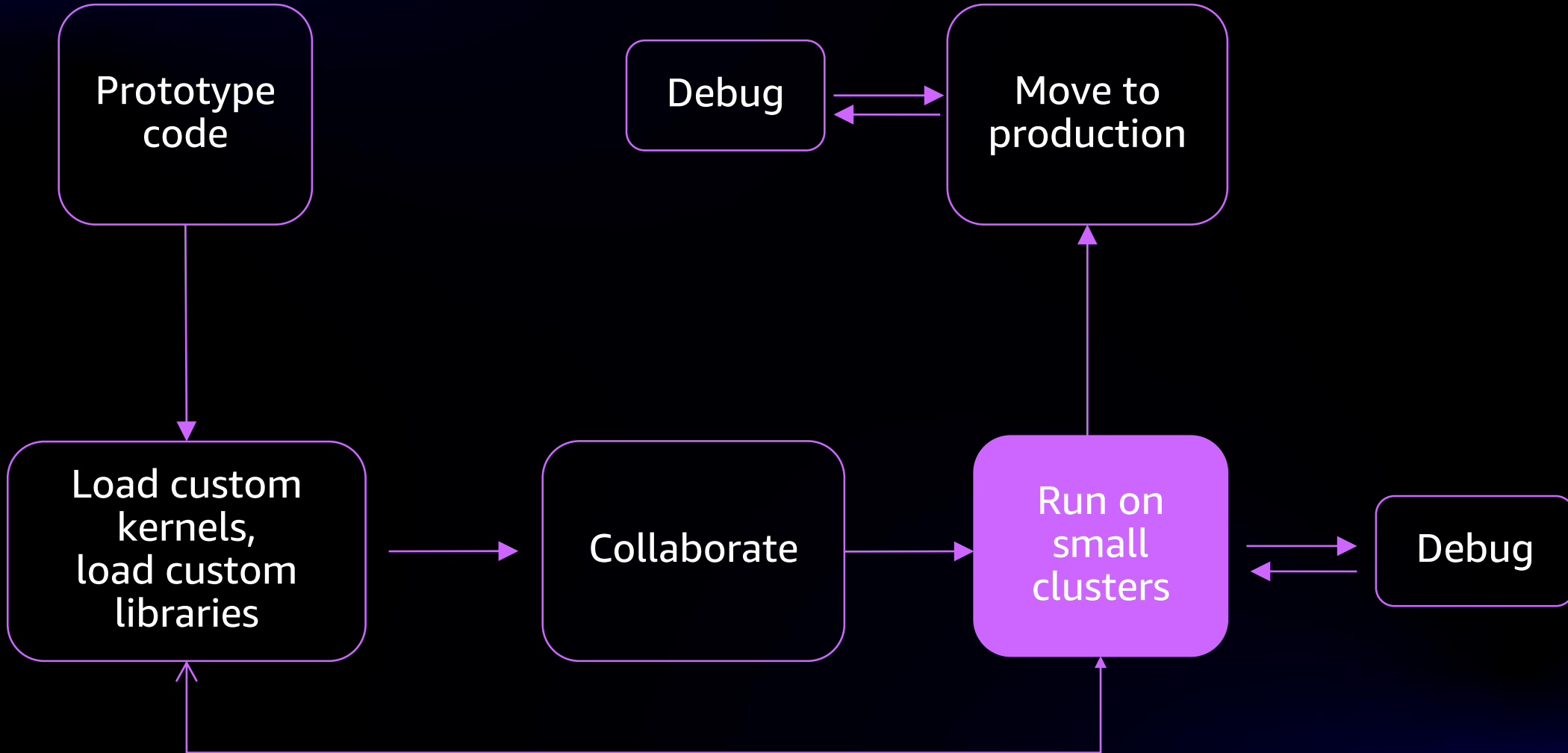
```
[ ]: import random

     NUM_SAMPLES = 100


[ ]: amshuang@amazon.com
     def inside(p):
         x, y = random.random(), random.random()
         return x*x + y*y < 1

     count = sc.parallelize(range(0, NUM_SAMPLES)) \
                 .filter(inside).count()
     print("Pi is roughly %f" % (4.0 * count / NUM_SAMPLES))
```
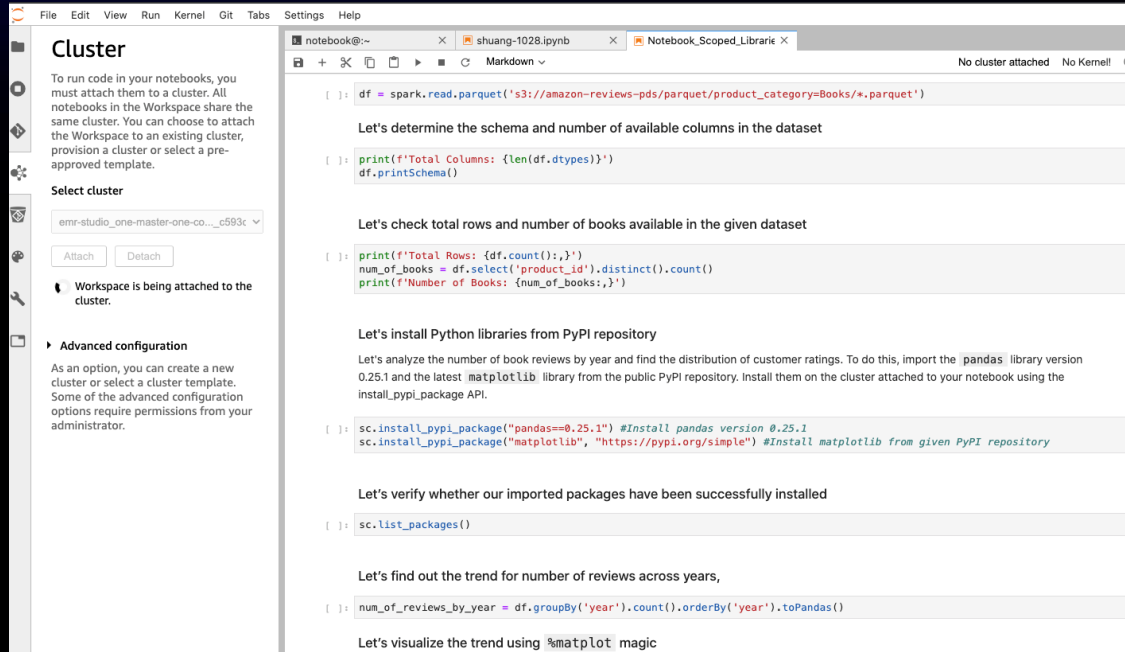
**Collaborate in real time**
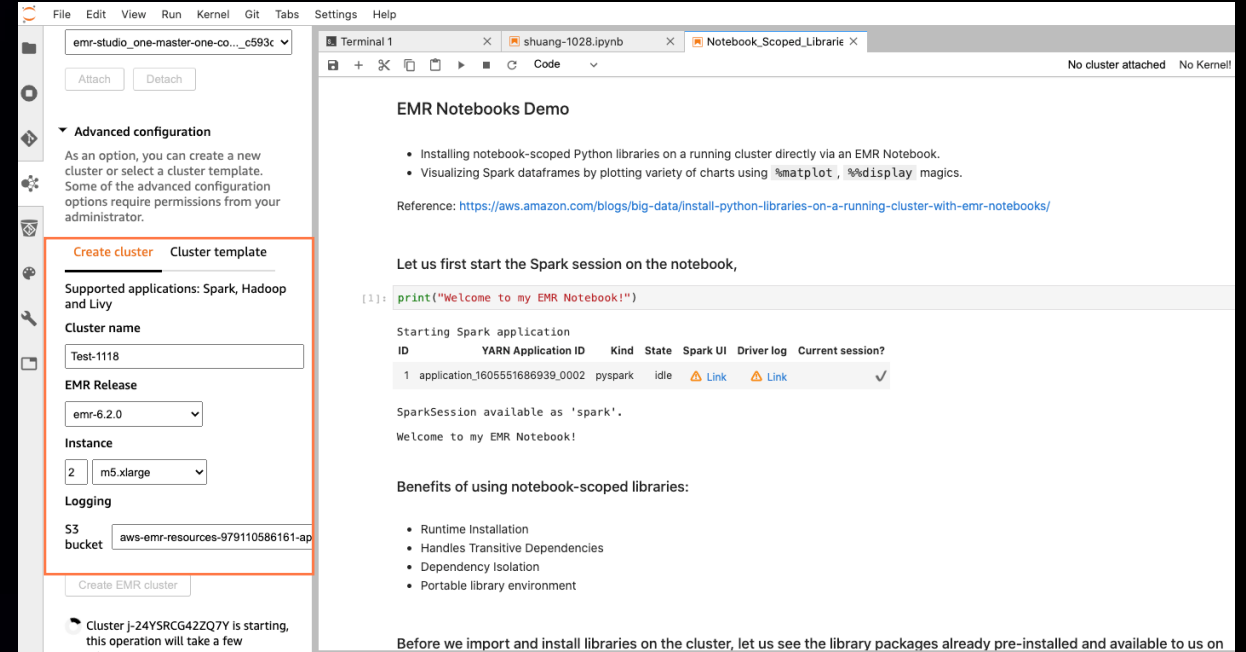
# Data science and engineering workflows

# Single-click attach to clusters to run jobs



**Attach workspace to an existing EMR cluster**

**Provision EMR clusters using simple configurations (you can limit users to either cluster templates or creating their own EMR cluster)**

# Single-click attach to clusters to run jobs



**Provision EMR clusters using preconfigured cluster templates via AWS Service Catalog**

**Connecting to clusters from Amazon EKS**

# Data science and engineering workflows

# Live debugging is simple



**Debug by clicking the "Spark UI" link in notebook to navigate to the live on-cluster Spark UI**

**View debugging information in Spark history server for the application in a separate browser tab**

# Data science and engineering workflows

# Simplify building pipelines from notebooks



Run notebooks as pipelines via Amazon Managed Workflows for Apache Airflow (MWAA)

Parameterize and chain notebooks that can be run as pipelines

# Schedule notebook pipelines



**Schedule notebooks from within Jupyter**

# Data science and engineering workflows

# Debugging production pipelines is easy



**Browse all clusters in one place**

**Narrow down clusters for investigation using filters such as cluster state**

# Debugging production pipelines is easy



**Diagnose jobs on both active and terminated clusters using Spark UI, Tez UI, and Yarn Timeline Service**

**Overlay execution context on jobs, even for terminated clusters and jobs**

"Choosing EMR Studio as our official workflow for Jupyter Notebooks on EMR has enabled us to reduce costs and time spent supporting data users. The built-in Git-based workflow has streamlined our previously cluttered landscape of notebooks. Connecting to an EMR cluster is as simple as selecting it in a dropdown box, avoiding the need to have personal clusters running 24/7."

**Phil Austin,
Director of DevOps
Verana Health**

"EMR Studio allows us to prototype Spark applications and data science models that power large-scale data processing and transformations. The integrated development environment makes it easy for data scientists and engineers to perform ad hoc analysis and debug data processing workloads."

**Saba El-Hilo**
**Head of Data Platform, Mapbox**

# Deep integration between EMR and SageMaker Studio

Process petabyte-scale data easily to train ML models using EMR Spark, Hive, and Presto from SageMaker Studio

Use EMR's integration with EC2 Spot and Graviton instances to run large-scale data processing at lower costs

Discover and connect to an EMR cluster from SageMaker Studio

Run Apache Spark, Hive, and Presto jobs on EMR from SageMaker Studio

Use familiar debugging tools such as Spark UI

Create, scale, and auto-terminate EMR clusters using AWS Service Catalog templates

aws

# Run Spark workloads on Amazon EKS easily

# Amazon EMR on Amazon EKS

Simplify infrastructure management

Consolidate multiple versions of Spark on same EKS cluster

Simplify Spark application upgrades

Add Multi-AZ resiliency by EKS with worker nodes across multiple AZs

Application and dependencies

Amazon EMR builds and runs Spark runtime as containers

Amazon EMR-specific components

Spark

OS

K8s

OS and infrastructure provisioned and managed by Amazon **EKS**

OS

EC2 instances

AWS Fargate pods

aws

# Consolidate workloads with EMR on EKS

**WORKLOAD CONSOLIDATION DRIVES HIGHER RESOURCE UTILIZATION AND LOWER COSTS**

## Amazon EMR on Amazon EC2

Amazon
EMR cluster

EMR 5.X

YARN

Amazon
EMR cluster

EMR 5.Y

YARN

Amazon
EMR cluster

EMR 6.Z

YARN

## Amazon EMR on Amazon EKS

Amazon
EKS cluster

EMR 5.X

EMR 5.Y

EMR 6.Z

Other
workloads

K8s

# Running jobs on EMR on EKS is easy

**NO NEED TO LEARN DIFFERENT TOOLS TO RUN SPARK JOBS USING EMR ON EKS**

AWS CLI/SDK

EMR Studio,
self-managed
notebooks

Apache
Airflow

AWS Step
Functions

```
aws emr-containers start-job-run \
    --virtual-cluster-id cluster_id \
    --name sample-job-name \
    --execution-role-arn execution-role-arn \
    --release-label emr-6.3.0-latest \
    --job-driver '{
        "sparkSubmitJobDriver": {
            "entryPoint": "local:///usr/lib/spark/examples/src/main/python/pi.py",
            "sparkSubmitParameters": "--conf spark.executor.instances=2 --conf spark.executor.memory=2G --conf
spark.executor.cores=2 --conf spark.driver.cores=1"
        }
    }'
```

aws

# Custom container images

- Install and configure packages specific to different workloads in different container images

- Use custom image validation tool to detect errors when creating custom container image

Create a custom container image

Upload image to your registry

Use image in multiple EMR on EKS jobs with near instantaneous job start time

aws

# Pod templates

Optimize price-performance by scheduling Spark executors to run on Amazon EC2 Spot or Graviton instances

Run a separate "sidecar" container next to the Spark driver or executor for logging or additional monitoring purposes

Run an "init" container that prepares the environment, e.g., downloads and installs dependencies

"Migration to EMR on EKS from open-source Spark on Kubernetes helped us to consolidate on two fronts – multiple Spark versions on same EKS cluster and Spark workloads alongside other workloads on same EKS cluster. This consolidation led to significant cost savings and reduced operational overhead."

**Ujjwal Sarin**
**Data Platform Engineering, Stitch Fix**

"The Battery Data Science Team at Rivian redefines how batteries are developed, monitored, and improved with near-real-time processing of massive datasets from our R&D labs and vehicle test fleet. As Rivian begins production and the volume of our battery data grows significantly, EMR on EKS enables our team to seamlessly scale our analytics capabilities. We use a variety of AWS-managed services, including ECR for Docker, EKS with Fargate for Kubernetes, and MSK for Kafka, for the flexibility and development speed of the open-source frameworks we want without the headaches of managed infrastructure."

**Rob Jenks**
**Principal Software Engineer – Cloud, Rivian**

# EMR simplifies building data lakes

aws

# Streaming data ingestion pipelines

**Writing dataset challenges**

- Make atomic changes
- Reader – writer isolation
- High throughput ingestion
- Small file compactions
- Concurrent writers
- Clustering by secondary keys

Ingest streaming telemetry events

**1**

Data lake

Mobile clients

API clients

Amazon Managed Streaming for Kafka

Amazon EMR

Spark streaming ingestion on EMR

**2**

Amazon S3

User

**Query dataset challenges**

- Incremental query
- Time travel query
- Multiple engine support
  - Presto, Hive, Spark SQL
  - Amazon Athena, Amazon Redshift Spectrum

# CDC ingestion pipelines challenges

Change data capture (CDC) from Amazon Aurora database is ingested using AWS DMS into a raw Amazon S3 data lake (ingestion tier)

**1**

Amazon RDS
(Hive metastore)

Amazon
Aurora → AWS DMS → Amazon S3
ingestion tier → Amazon EMR
(Spark) → Amazon S3
analytics tier

**2**

Change log from datasets are ingested using a Spark streaming job

## Challenges

- Make atomic changes
- Reader – writer isolation
- High throughput ingestion
- Small file compactions
- Row-level upserts and deletes
- Clustering by secondary keys

# GDPR (data erasure) pipelines challenges

Amazon RDS
(Hive metastore)



Monthly GDPR
deletion pipelines

1

Monthly deletion
pipelines

Amazon EMR
(Spark)

Amazon S3
analytics tier

2

Spark Batch
deletion jobs

## Challenges

- Row-level upserts and deletes
- Concurrent writers

# Apache Hudi enables transactional data lakes

### Ingestion

- Transactions (ACID) – reader and writer isolation
- Transactions (ACID) – concurrent writer support
- Record-level upserts and deletes
- High throughput streaming ingestion
- Spark, Flink, and Java Writer Support
- Automatic compaction of small files
- Spark SQL DML support (Hudi 0.9.0)  **NEW!**

### Query

- Spark, PrestoDB/Trino, and Hive support
- Efficient queries across partitions and files
- Incremental query support
- Time travel query support

# Apache Hudi enables transactional data lakes

Administration



- Async background compaction of files
- Async background sorting and clustering of keys
- Automatically clean up files beyond retention period
- Metrics for past commits or rollbacks

aws

# Easily operationalize at scale using detailed metrics

**AMAZON CLOUDWATCH INTEGRATION**

| All metrics | Graphed metrics | Graph options | Source |

N. Virginia ⌄    All  ›  Hudi  ›  Hudi Table, Metric Type    🔍 Search for any metric, dimension or resource id    Graph search

| ☐ | Hudi Table (40) | Metric Type | Metric Name |
| --- | --- | --- | --- |
| ☐ | tpcds_store_sales_3TB_08 | gauge | commit.totalScanTime |
| ☐ | tpcds_store_sales_3TB_08 | gauge | commit.totalUpdateRecordsWritten |
| ☐ | tpcds_store_sales_3TB_08 | gauge | commit.totalUpsertTime |
| ☐ | tpcds_store_sales_3TB_08 | gauge | finalize.duration |
| ☐ | tpcds_store_sales_3TB_08 | gauge | finalize.numFilesFinalized |
| ☐ | tpcds_store_sales_3TB_08 ▾ | count ▾ | timer.clean ▾ |
| ☐ | tpcds_store_sales_3TB_08 | count | timer.commit |
| ☐ | tpcds_store_sales_3TB_08 | gauge | TimelineService.TOTAL_CHECK_TIME |

**No. of commits**

**Total records written**

# Apache Hudi is widely supported on AWS

Spark, Hive, Presto, Flink
Support on Amazon EMR

AWS Glue Catalog
and ETL support

AWS Lake Formation
FGAC support

Amazon Athena
Native query support

Amazon Redshift Spectrum
Native query support

AWS Database Migration Service
CDC ingestion support

Amazon CloudWatch
Integration for metrics

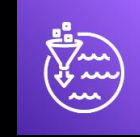# Security controls in EMR

# Comprehensive security features

**Isolation**

**Authentication**

**Authorization**

**Encryption**

**Audit**

VPC

Private subnets

Security groups

LDAP

Kerberos

AWS SSO
(EMR Studio)

AWS IAM
(EMR Studio)

Cluster IAM role

User execution
role (preview)   **NEW!**

FGAC using
Apache Ranger

FGAC using AWS
Lake Formation
(preview)   **NEW!**

Encryption at rest

Encryption in
transit

Key management

Audit using
Ranger via
Amazon
CloudWatch
Logs

Audit using
AWS Lake
Formation
via AWS
CloudTrail   **NEW!**

# One cluster – One role



User1

Amazon EMR

**Cluster Role1**

User2

Amazon EMR

**Cluster Role2**

Amazon Kinesis
Data Streams

Amazon S3

Amazon DynamoDB

# Enable multi-tenant shared clusters

**NEW!**

**User execution role (preview): User1 has access to Stream1, Bucket1, and Table1**

**Execution Role1**

User1

User2

**Single shared cluster**

Amazon EMR

Stream1

Amazon Kinesis
Data Streams

Bucket1

Amazon S3

Table1

Amazon DynamoDB

# Enable multi-tenant shared clusters

**User Execution Role (preview): User2 has access to Stream2, Bucket2, and Table2**
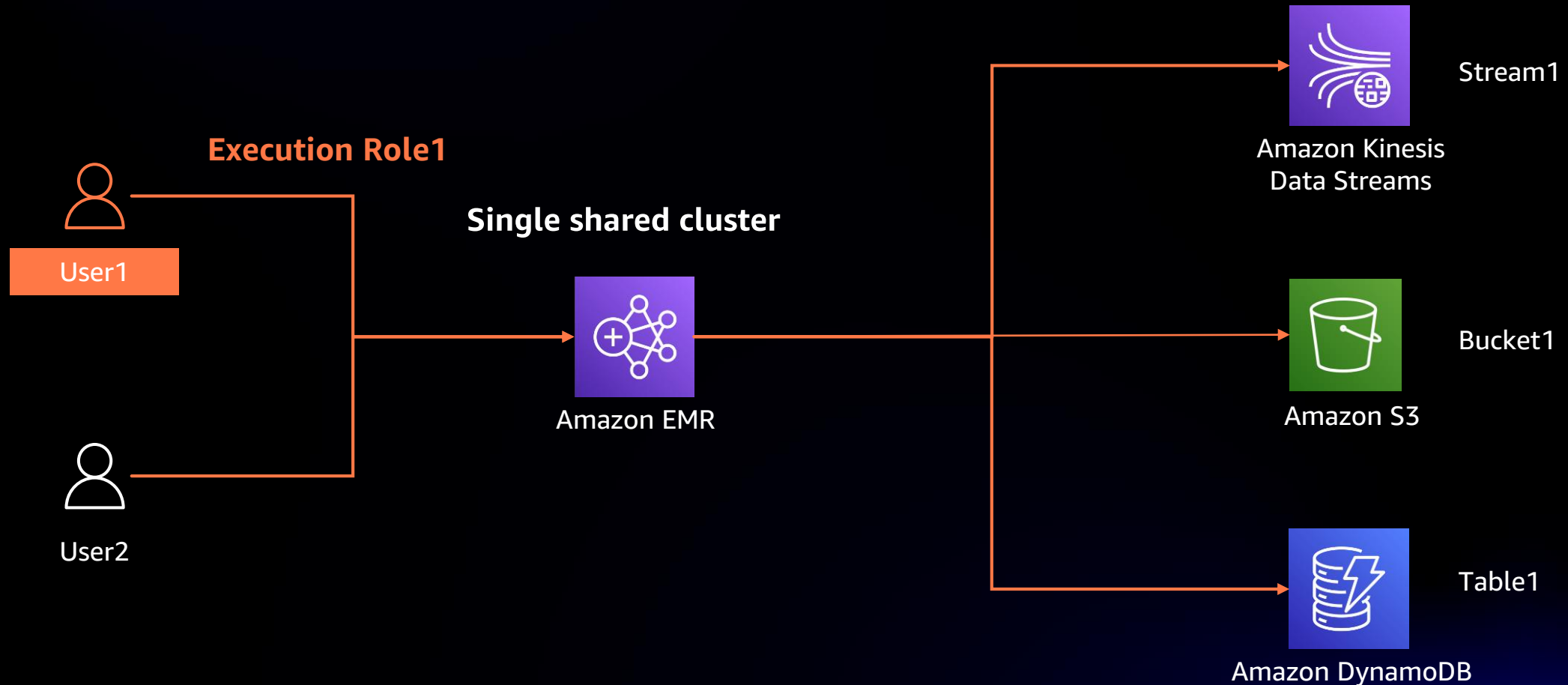
User1

**Multi-tenant shared cluster**

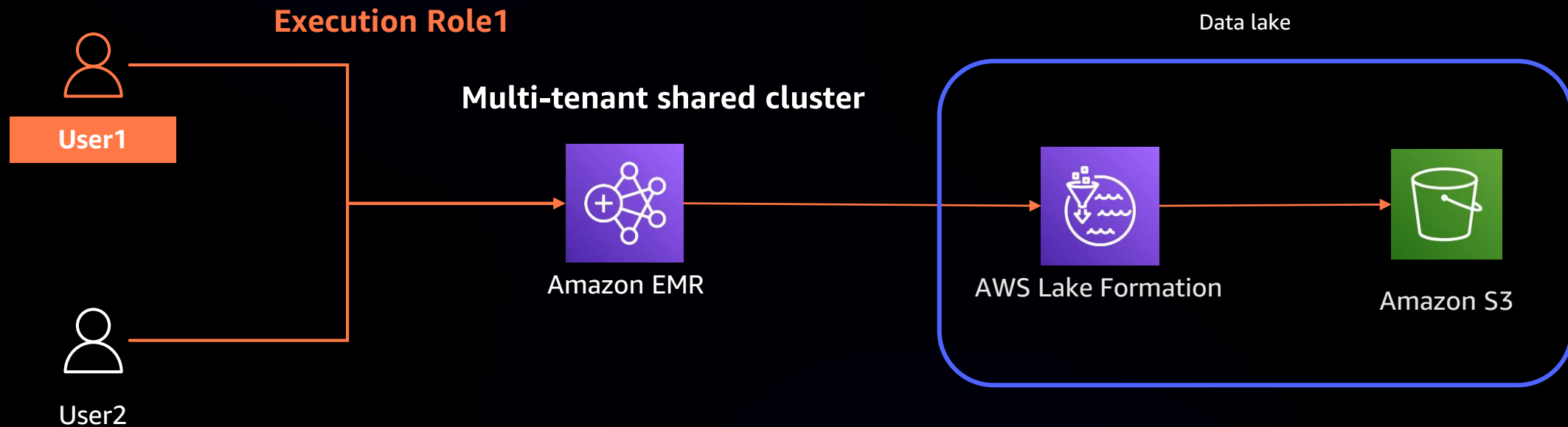Amazon EMR

User2

**Execution Role2**

Stream2

Amazon Kinesis
Data Streams

Bucket2

Amazon S3
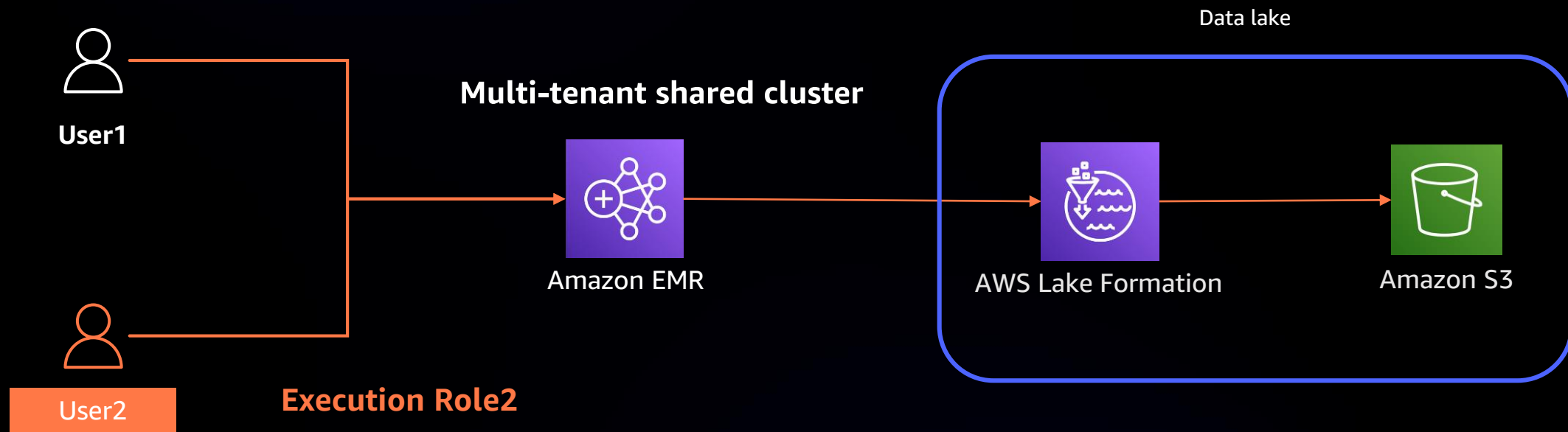
Table2

Amazon DynamoDB

# Enable fine-grained access control

**FGAC using AWS Lake Formation (preview): User1 has access to Table1, Columns 1–10**

# Enable fine-grained access control

**FGAC using AWS Lake Formation (preview): User2 has access to Table1, Columns 5–10**



Data lake

**User1**

**Multi-tenant shared cluster**

Amazon EMR

AWS Lake Formation

Amazon S3

User2

**Execution Role2**

# Thank you!

Vincent Gromakowski

gromav@amazon.com