



# Amazon EMR Studio

Accelerate Apache Spark and Other Big Data  
Application Development with EMR Studio

Damon Cortesi

# Amazon EMR

Easily Run Spark, Hive, Presto, HBase, Flink, and more big data apps on AWS

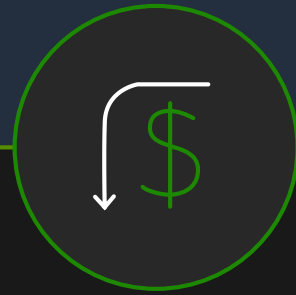
## Latest versions



Updated with latest open source frameworks **within 30 days**

Support for popular OSS like **Flink, Hudi**

## Best Performance at Lowest cost



Spark workloads run **2.4x faster** compared to Open Source

**50–80% reduction** in costs with EC2 Spot and Reserved Instances

**Per-second billing** for flexibility

## Use S3 storage



Process data in S3 **securely** with **high performance** using the EMRFS connector

**Scale Compute and Storage** independent of each other

## Easy & Scalable

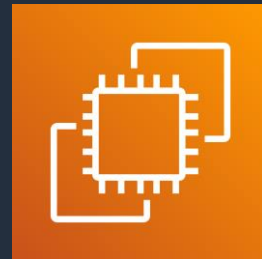


**Fully managed**, no cluster setup, node provisioning or cluster tuning

**Vertical and Horizontal Auto-Scaling** to suit workload demands

# Amazon EMR flexible deployment options

Amazon EMR on Amazon EC2



AWS offers more instance options than any other cloud provider. Choose the instance that gives the best performance or cost for your workload, including Graviton2. Take advantage of On-Demand, Reserved, and Spot Instances to optimize costs.

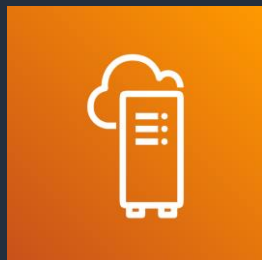
Amazon EMR on Amazon EKS

NEW  
!



Use EMR to automate the provisioning, management, and scaling of Apache Spark jobs on [Amazon Elastic Kubernetes Service](#) (EKS), and take advantage of the optimized EMR runtime.

Amazon EMR on AWS Outposts



Set up, deploy, manage, and scale EMR in your on-premises environments, just as you would in the cloud. AWS Outposts brings AWS services, infrastructure, and operating models to virtually any data center, co-location space, or on-premises facility.

# EMR Studio integrated development environment



Easily build and deploy data science code without logging in to AWS console



Start notebooks in seconds, run jobs later



Build production pipelines simply and flexibly



Save debugging time with native application UIs in one place

# Use existing identity provider with AWS SSO

Easily enable single sign-on in the **AWS SSO** console

- Integrated with AWS organizations
- Choose identity source like Okta, Azure AD, OneLogin, Ping Identity, etc.
- Import users and groups from identity source to AWS SSO

# Abstract infrastructure details using templates

Create cluster templates in **AWS Service Catalog** to simplify running jobs for your data scientists and data engineers –

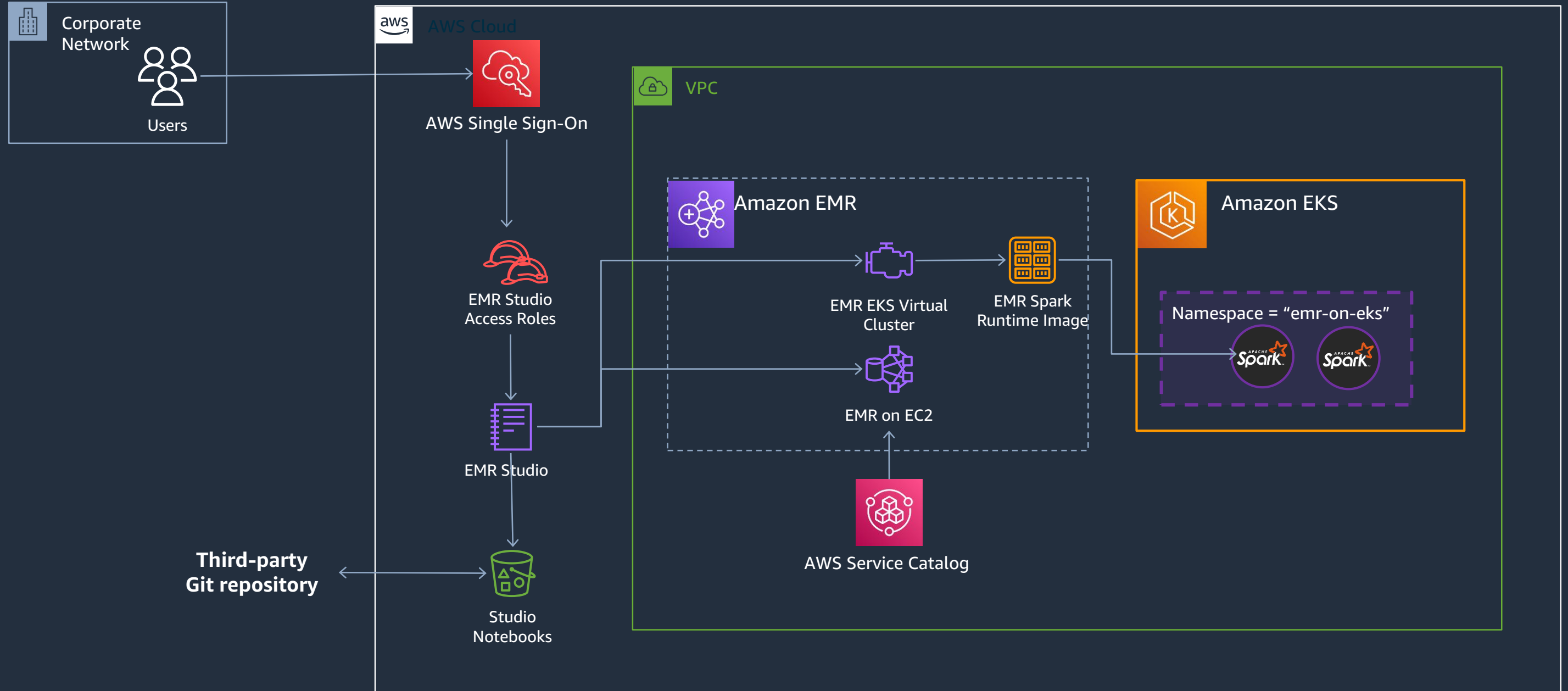
- Control Spark version, Amazon EMR version, etc.
- Restrict number of instances and instance types
- Preconfigure libraries and packages
- Add security settings

# Restrict who can provision a cluster or template

Control who can submit jobs by spinning up clusters on demand instead of keeping long-running clusters for all

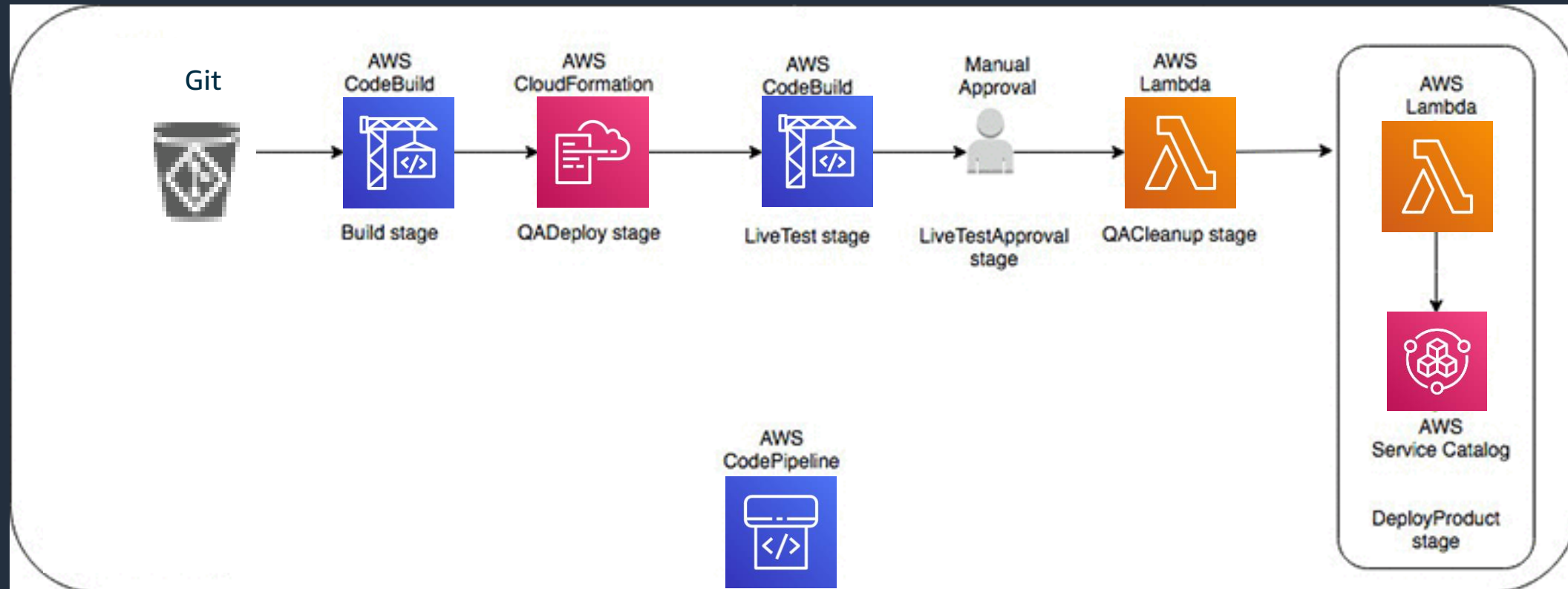
- Allow users and groups assigned to a Studio to access a set of cluster templates
- Authorize users and groups in a Studio to create clusters using configurations
- Limit users and groups to provision permitted templates to run jobs for cost savings

# Example Architecture





# Trigger pipelines from code repositories



Commit code to Git and trigger a continuous integration and deployment pipeline supported by AWS services

# Simplify building pipelines from notebooks

Off DAG: custom\_cluster\_execution\_sensor\_dag Ray testing execution schedule: \*/10 \* \* \* \*

Graph View Tree View Task Duration Task Tries Landing Times Gantt Details Code Trigger DAG Refresh Delete

Base date: 2020-10-26T20:20:00Z Number of runs: 25 Go

MyEmrJobFlowSensor NotebookExecutionSensor PythonOperator

scheduled skipped upstream\_failed up\_for\_reschedule up\_for\_retry failed success running queued no\_status

[DAG]  
create\_cluster\_task  
check\_cluster  
start\_execution\_task  
check\_notebook\_execution

11 AM 12 PM 01 PM

Airflow DAGs Security Browse Admin Docs About 2020-10-26, 20:34:02 UTC assumed-role/Admin/rluamzn-lsengard

Off DAG: custom\_cluster\_execution\_sensor\_dag Ray testing execution schedule: \*/10 \* \* \* \*

Graph View Tree View Task Duration Task Tries Landing Times Gantt Details Code Trigger DAG Refresh Delete

Base date: 2020-10-26T20:20:01Z Number of runs: 25 Run: scheduled\_2020-10-26T20:10:00+00:00 Go

create\_cluster\_task  
check\_cluster  
start\_execution\_task  
check\_notebook\_execution

20:21:00 20:22:00 20:23:00 20:24:00 20:25:00 20:26:00 20:27:00 20:28:00

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

Parameterize and chain notebooks that can be run as pipelines



# Resources

- EMR Studio Workshop: <https://emr-developer-experience.workshop.aws/>
- EMR Studio Templates: <https://github.com/aws-samples/emr-studio-samples>

# Thank you!

