AWS
re:Invent

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

**ANT205-R**

# Achieving your modern data architecture

Santosh Chandrachood

General Manager
AWS Glue

# Agenda

Modern data architecture on AWS

End-to-end data life cycle on the modern data architecture

Data governance and data mesh in action

Journey towards modern data architecture

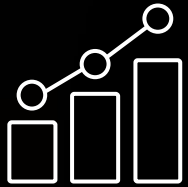# Deriving insights from data is hard

AWS CAN HELP



**Data silos**

**People silos**

**Business silos**

# Challenges in data silos

Data is growing exponentially

From new sources

Increasingly diverse

More users needing secure access

Cost and performance

# Challenges in people silos

**UNIQUE USER SKILLSETS**

**TOOL PREFERENCES**

**REQUIRED PROCESSES**

# Challenges in business silos

**COST**

**LEGACY INFRASTRUCTURE**

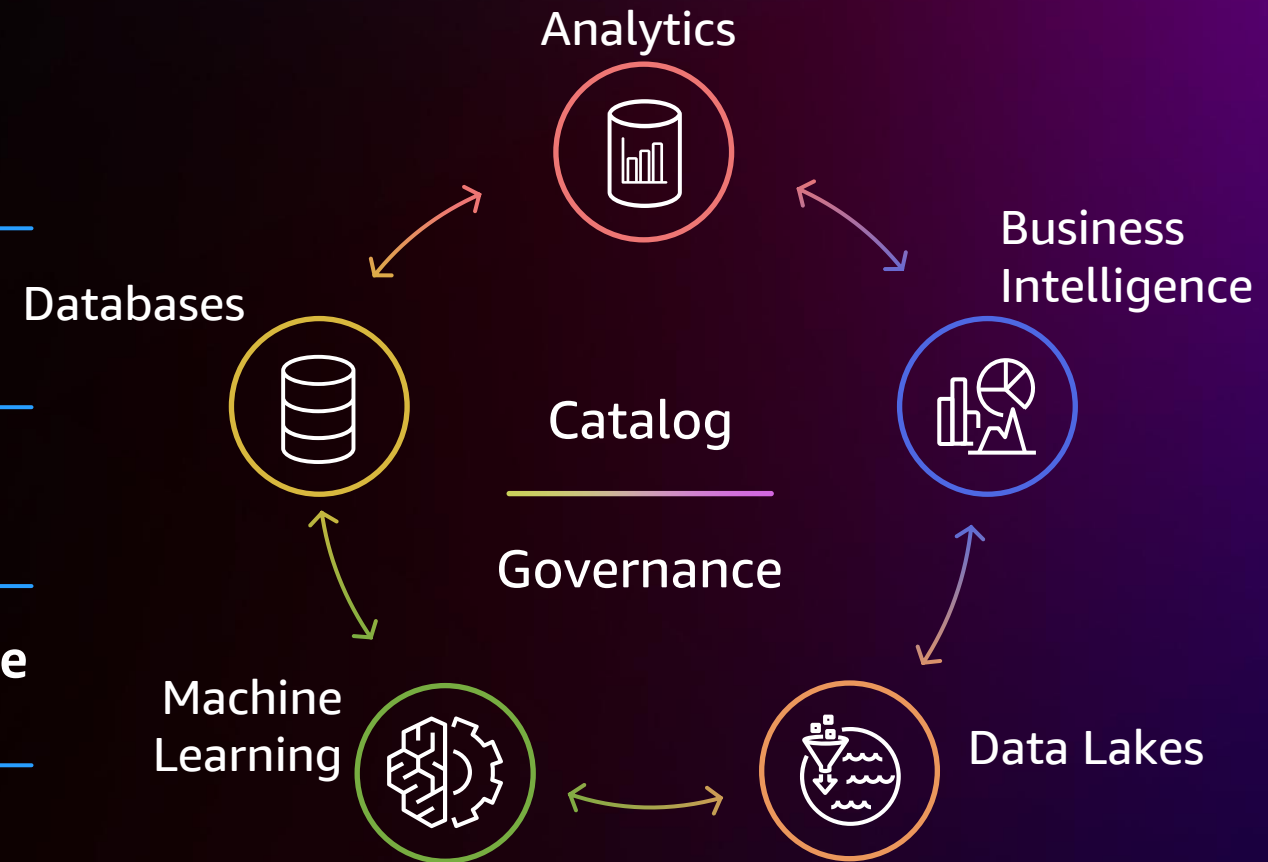**AGILITY**

# The five pillars of a modern data architecture

**Unified analytics**

**Highest performance at the lowest cost**
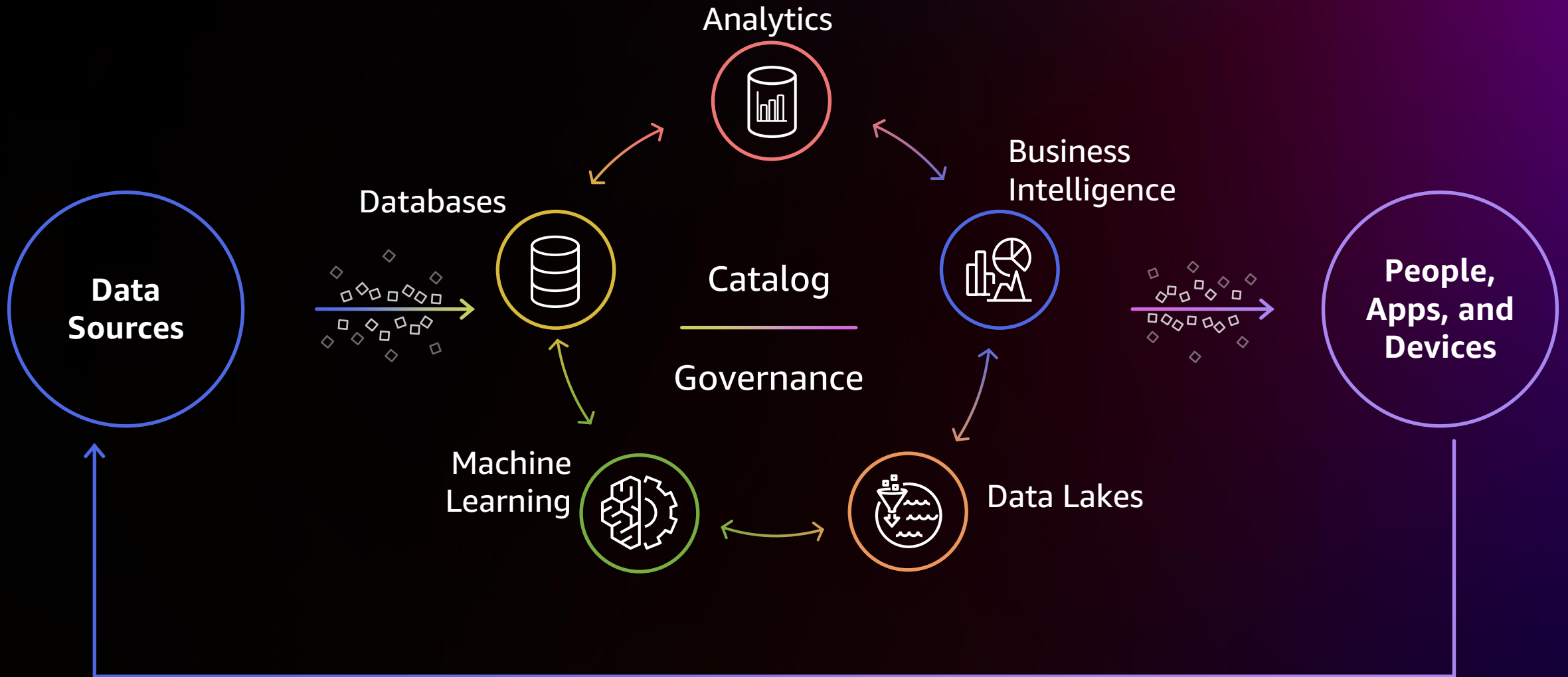
**Machine learning integration**

**Unified data access, security and governance**

**Insights for everyone**

Analytics

Business Intelligence

Databases

Catalog

Governance

Machine Learning

Data Lakes

# Modern data architecture

# Warner Bros. Games scales data analytics

## Challenge:

Measure the business and provide meaningful feedback without getting in the way of the creative process. The business isn't static. They have data velocity, volume, variety, and voracity. They need a scalable infrastructure, data federation and democracy, and a way to act on the data.

## Solution:

A modern data architecture with Amazon Redshift, Amazon EMR, and AWS Glue.

## Result:

One version of the data for all stakeholders to access with increased scalability, lowered overhead costs, and more compute and memory for the same cost.

Amazon Redshift     AWS Glue     Amazon EMR

aws

# Modern data architecture's five benefits

**Unified analytics**

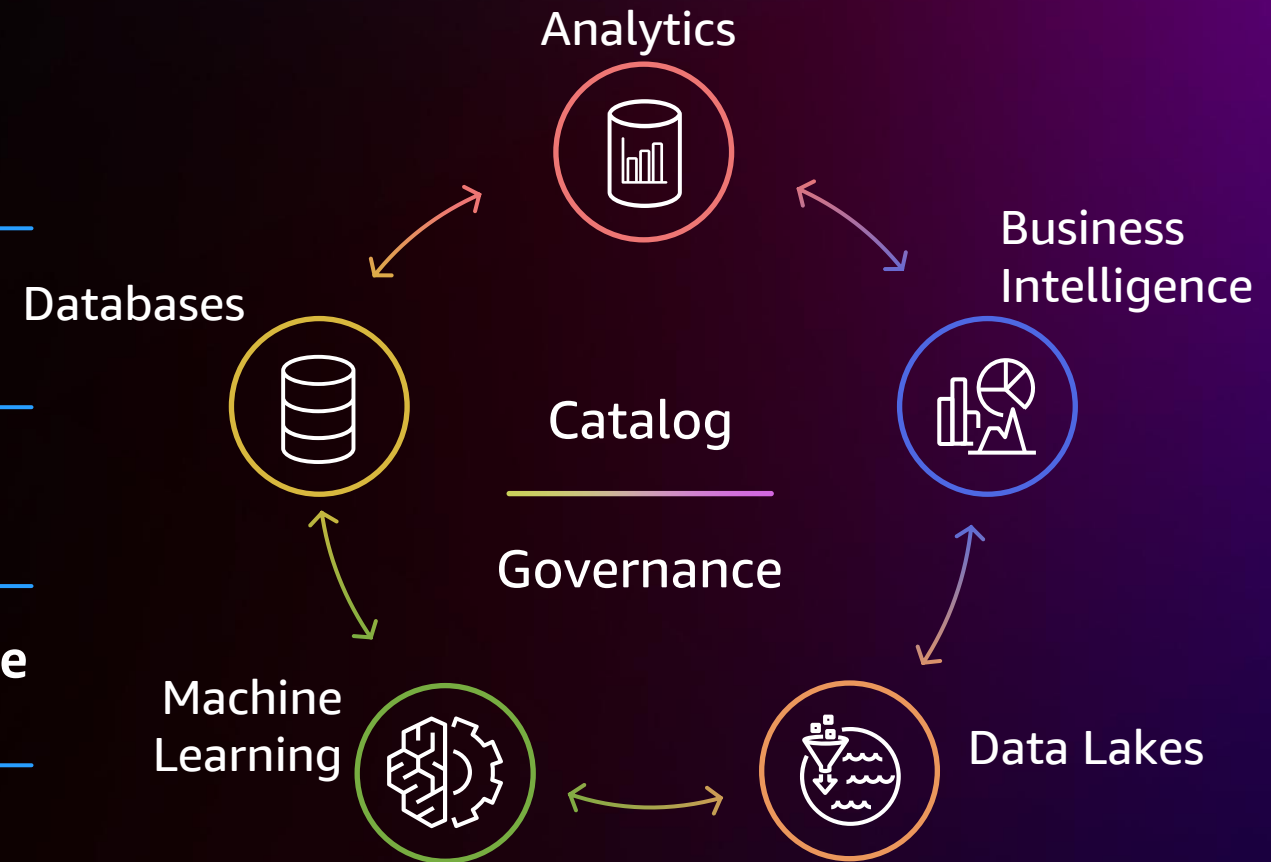**Highest performance at the lowest cost**

**Machine learning integration**

**Unified data access, security and governance**
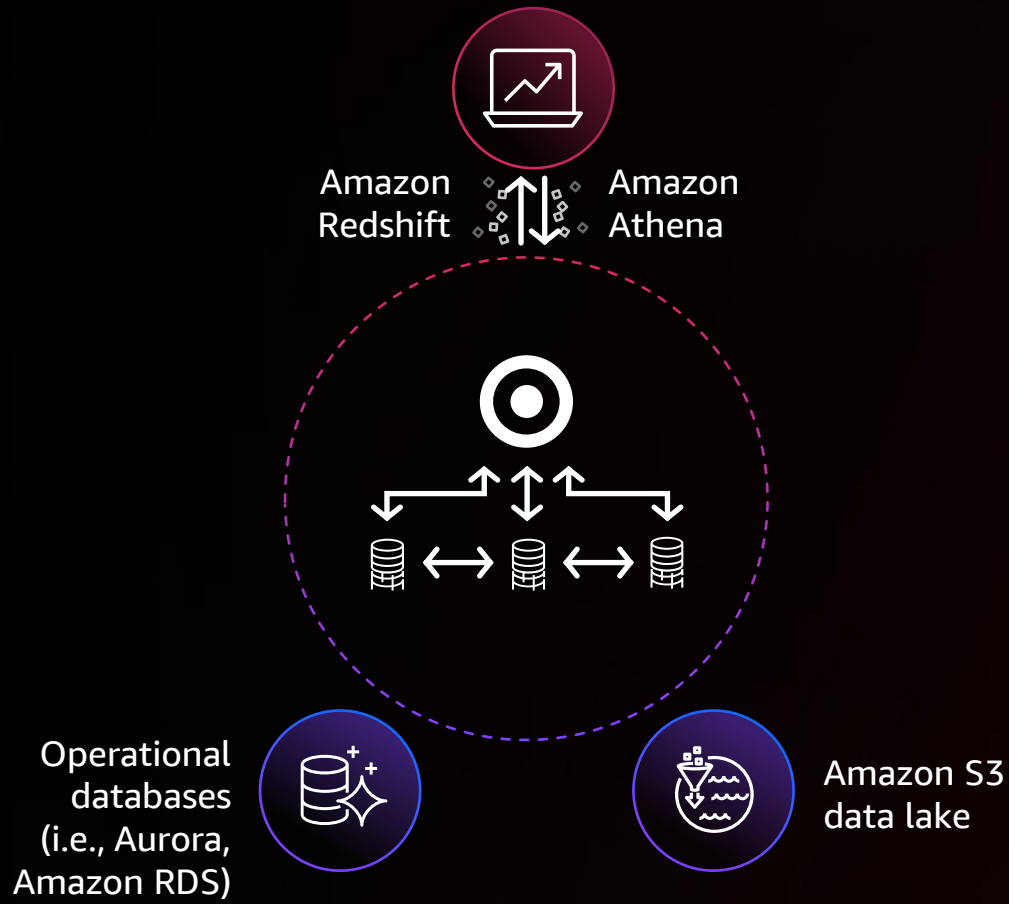
**Insights for everyone**

Analytics

Business Intelligence

Databases

Catalog

Governance

Machine Learning

Data Lakes

# Unified analytics

**Data access anywhere**

**Bringing your data in**

**Right tools for analytics**

# Federated query in Amazon Redshift and Amazon Athena

Amazon Redshift    Amazon Athena

Operational databases (i.e., Aurora, Amazon RDS)

Amazon S3 data lake

**Integrate varied data stores** with data warehouse and Amazon S3 data lake

Analytics on combined data **without data movement** and ETL delays

**Flexible and easy** way to ingest data, avoiding complex ETL pipelines

**NEW** **[PREVIEW]**

# Amazon Aurora Zero-ETL to Amazon Redshift

Eliminates the need to build and maintain complex ETL pipelines

Run near-real-time analytics and machine learning using Amazon Redshift on petabytes of transactional data from Amazon Aurora

Derive insights using advanced analytics in Amazon Redshift from data consolidated from multiple Amazon Aurora database clusters

**NEW** **[GA]**

# Amazon Redshift integration for Apache Spark

## SIMPLIFY AND SPEED UP APACHE SPARK APPLICATIONS ACCESSING AMAZON REDSHIFT DATA FROM AWS ANALYTICS SERVICES

**Author Apache Spark applications using Java, Python, Scala, with access to rich, curated data in your data warehouse**

**No manual setup and maintenance of uncertified versions of Spark-Amazon Redshift open-source connectors**

**Advanced pushdown optimizations in the Apache Spark-Amazon Redshift connector accelerate 3 TB out-of-the-box TPC-DS queries by 10x**

**Improved security with IAM-based credentials**

AWS Glue

Amazon EMR

Amazon SageMaker

*Pre-packaged Amazon Redshift connector for Spark*

# Broadest and most cost-effective set of analytics services

**Interactive query**
Amazon Athena

**Big data processing**
Amazon EMR

**Operational and log analytics**
Amazon OpenSearch Service

**Real-time analytics**
Amazon Kinesis and Amazon MSK

**Business intelligence**
Amazon QuickSight

**Data warehouse**
Amazon Redshift

**Governance & data lakes**
Amazon S3, AWS Lake Formation, AWS Glue Data Catalog

**Data integration & 3P sources**
AWS Glue, ADX, and Amazon AppFlow

# Modern data architecture's five benefits

**Unified analytics**

---

**Highest performance at the lowest cost**

---

**Machine learning integration**

---

**Unified data access, security and governance**

---

**Insights for everyone**

Analytics

Business
Intelligence

Databases

Catalog

Governance

Machine
Learning

Data Lakes

# Highest performance at the lowest cost

**Scale linearly with predictable high performance**

**Maximize your cost savings**

**Self-learning, self-tuning system to enhance performance**

# Price performance innovations in 2022

**Amazon EMR**

EMR Serverless GA

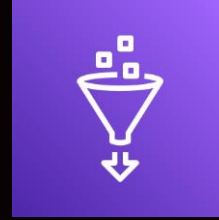11-16% performance improvement with Graviton2 at 20%+ reduced cost

68% performance increase for Spark

19% performance boost on Graviton3

**Amazon Athena**

Applications start up in under a second

**AWS Glue**

Auto Scaling

Flex Jobs

2.4x performance increase for Spark

**Amazon OpenSearch Service**

40% query boost on Graviton2

38% indexing boost on Graviton2

**Amazon Redshift**

Up to 5x better price performance vs. other cloud data warehouses

Up to 7x better price performance vs. other cloud data warehouses on high concurrency

Low latency workloads like dashboarding applications

# Amazon Redshift performance improvements

| Compute | System | Autonomics |
|---|---|---|
| Vectorized scans for Amazon Redshift tables | String-encoding for in-memory perf | Performance mode |
| Write/Commit performance | CaaS cache pre-warming | Auto WLM enhancements |
| Snapshot isolation | CaaS region expansion | ATO enhancements |
| Concurrency scaling writes (GA) | Incremental updates of MVs on datashares | Advisor enhancements |

# Differentiated performance

## ON SPARK, PRESTO, AND HIVE

**3.9x**

faster than standard Apache Spark 3.0 in TPC-DS 3TB benchmark

**4.2x**

faster than standard OSS Trino 388 in TPC-DS 3TB benchmarks

**11–16%**

performance improvement with Graviton2 at **20%+ reduced cost**

**100%**

open-source API compliant
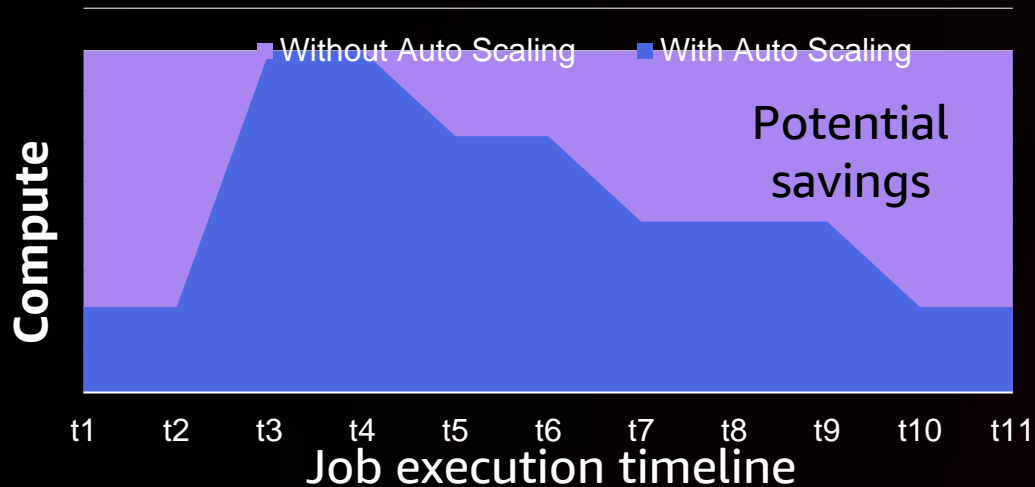
**30%**

Better price-performance with Graviton2

# AWS Glue cost optimization

## Auto scaling

Automatically resize compute for lower cost



Without Auto Scaling    With Auto Scaling

Potential savings

Compute

t1   t2   t3   t4   t5   t6   t7   t8   t9   t10   t11

Job execution timeline

Reduce cost by 20-40%

Simplify capacity and performance planning

## Flex

spare capacity execution



Up to
**34% cost savings**

Cost effective for **one-time data-load** workloads
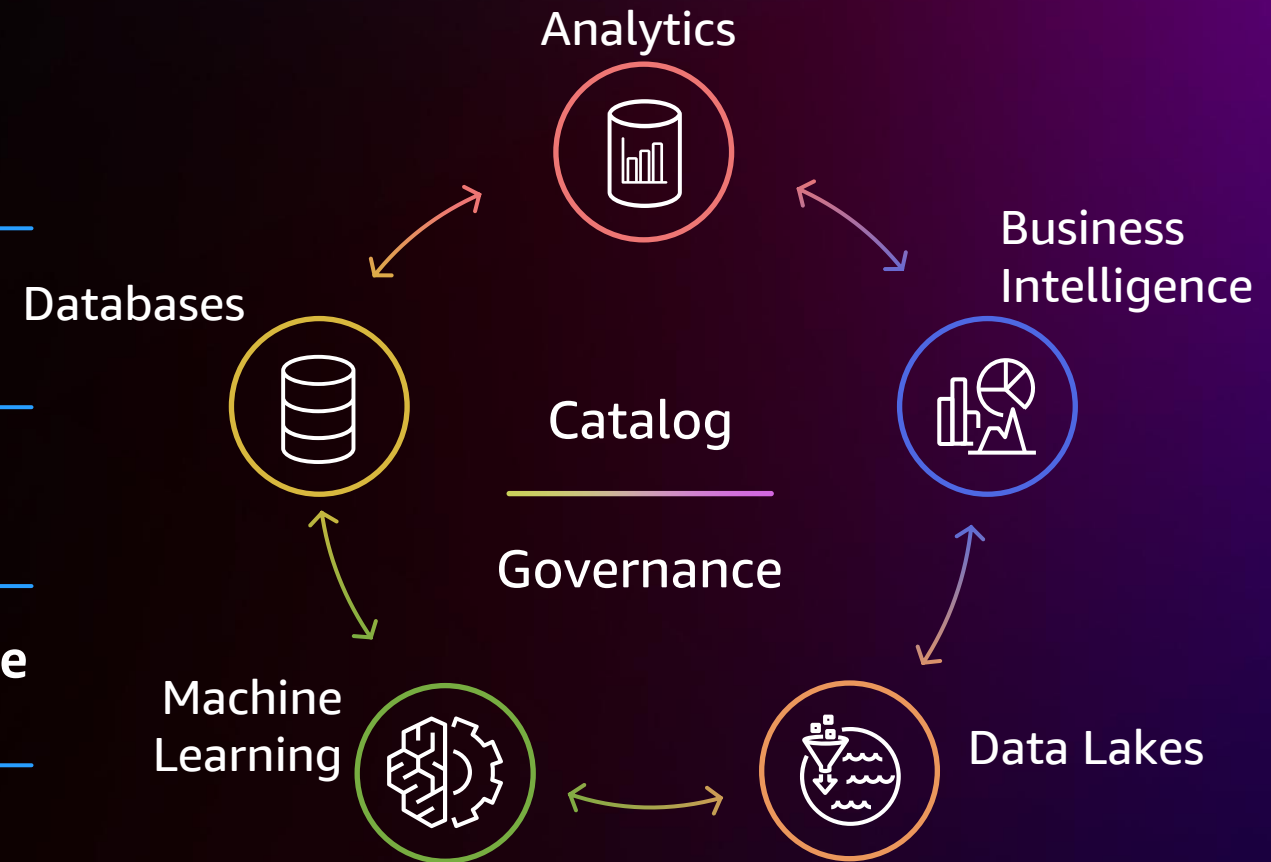
# Modern data architecture's five benefits

**Unified analytics**

**Highest performance at the lowest cost**

**Machine learning integration**

**Unified data access, security and governance**

**Insights for everyone**

Analytics

Business Intelligence

Databases

Catalog

Governance

Machine Learning

Data Lakes

# Connecting data services and ML to drive more value

Databases + Data warehouses and data lakes + Business intelligence tools

AMAZON
AURORA ML

AMAZON
NEPTUNE ML

AMAZON
REDSHIFT ML

AMAZON
ATHENA ML

AMAZON
QUICKSIGHT Q

# Amazon Redshift ML

## EASILY CREATE AND TRAIN ML MODELS USING SQL QUERIES WITH AMAZON SAGEMAKER
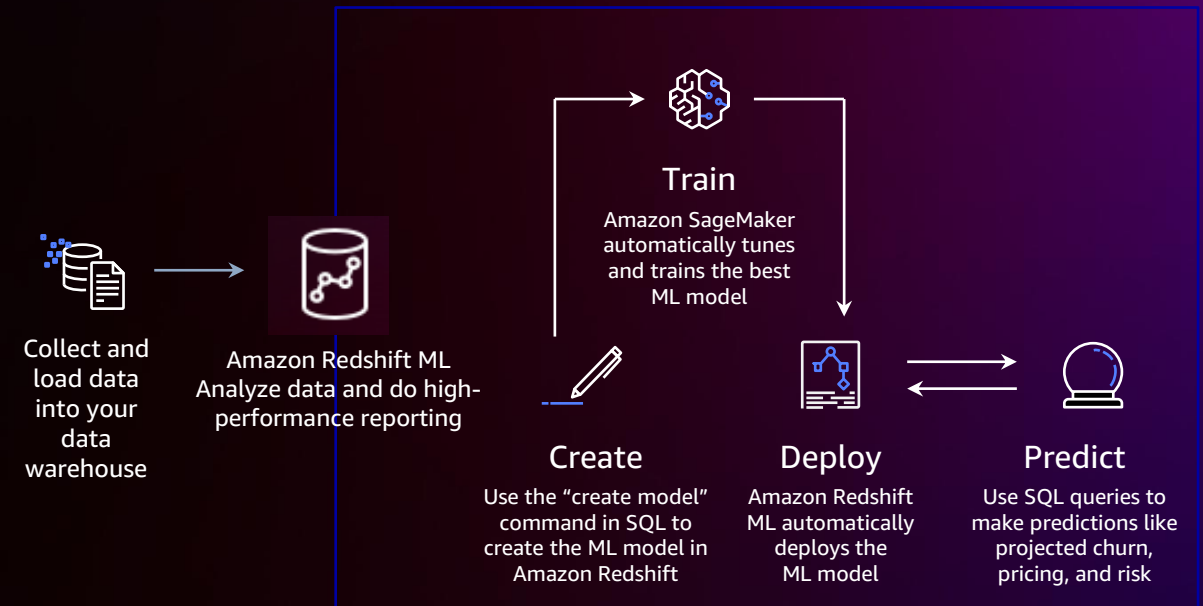
# 80+ billion
## predictions per week

Train and create ML models using SQL

Automatic pre-processing, creation, training, deployment, and inferencing of models
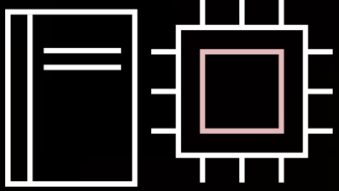
SageMaker models for in-database or remote inference

Supervised and unsupervised trainings

Collect and load data into your data warehouse

Amazon Redshift ML
Analyze data and do high-performance reporting

**Train**
Amazon SageMaker automatically tunes and trains the best ML model

**Create**
Use the "create model" command in SQL to create the ML model in Amazon Redshift

**Deploy**
Amazon Redshift ML automatically deploys the ML model

**Predict**
Use SQL queries to make predictions like projected churn, pricing, and risk

jobcase

# Amazon SageMaker Studio Universal Notebooks

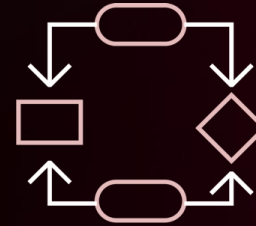## PERFORM DATA ENGINEERING, ANALYTICS AND ML IN ONE NOTEBOOK

**Built-in data and analytics integration**

Connect with Amazon EMR, AWS Glue and data lakes on Amazon S3

**Interactive data preparation**

Interactively query, analyze and transform wide range of data

**Inline debugging and monitoring**

Visually debug and monitor Spark jobs inline in same notebook

**Build ML workflows**

Build end to end ML workflows without leaving the notebook

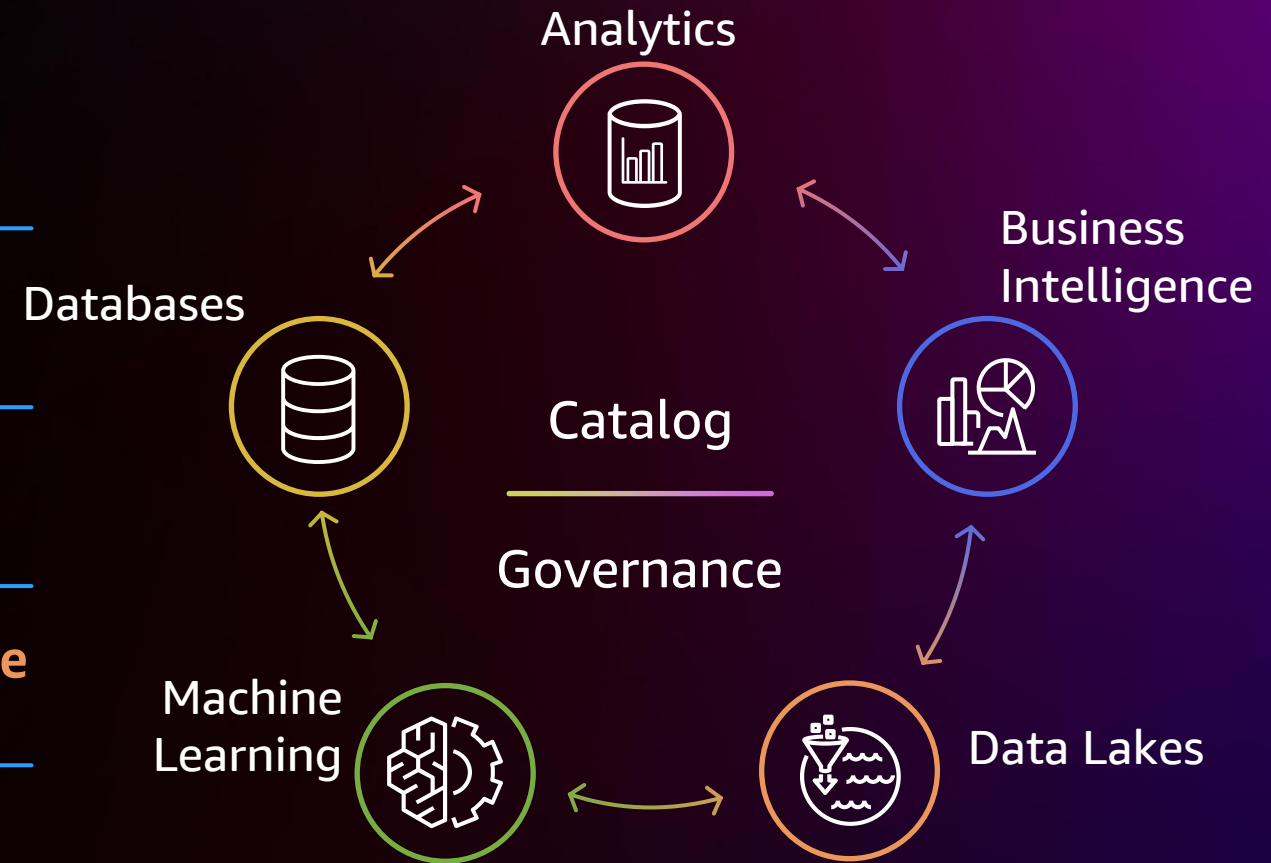# Modern data architecture's five benefits

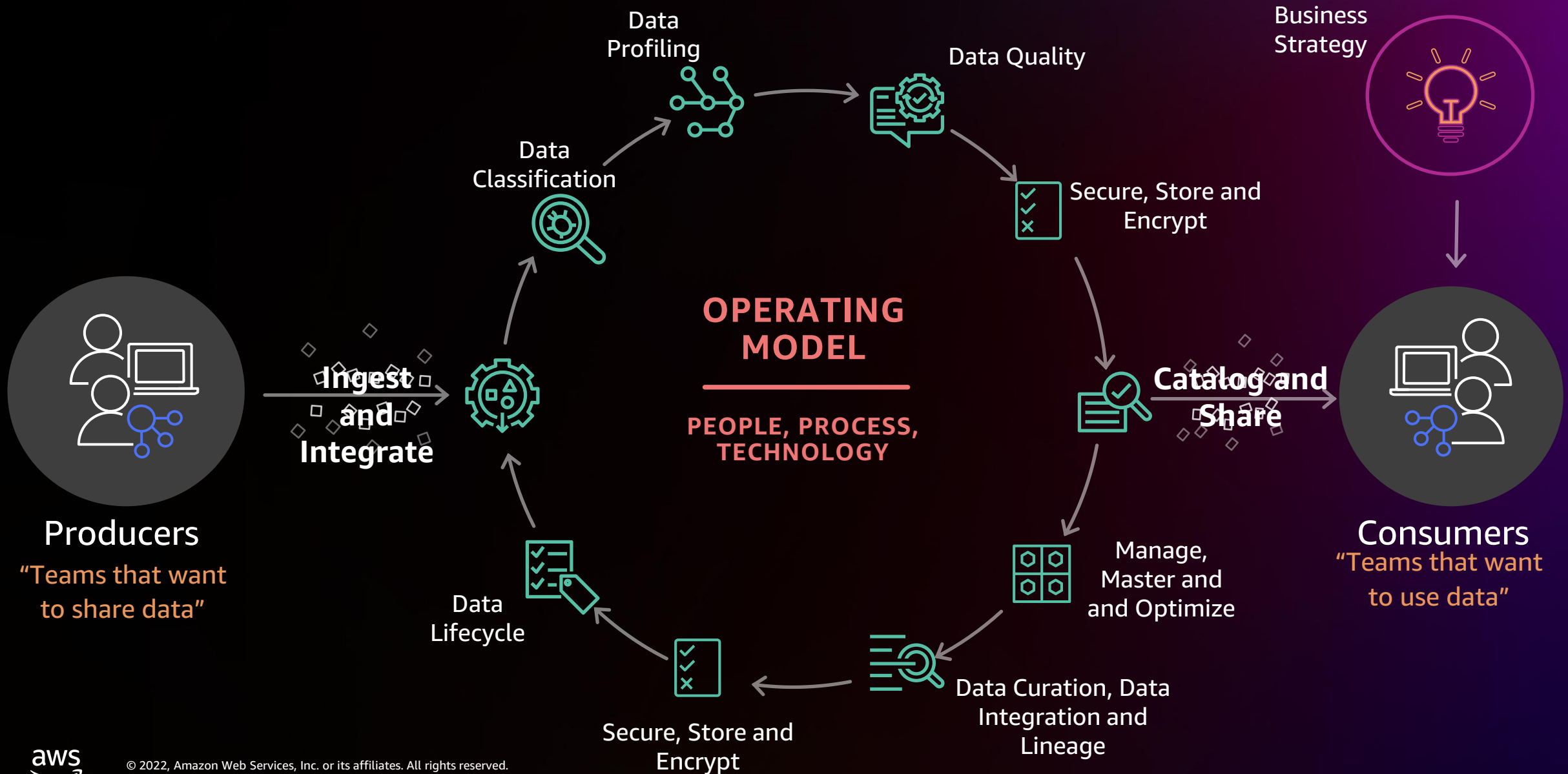**Unified analytics**

**Highest performance at the lowest cost**

**Machine learning integration**

**Unified data access, security and governance**

**Insights for everyone**
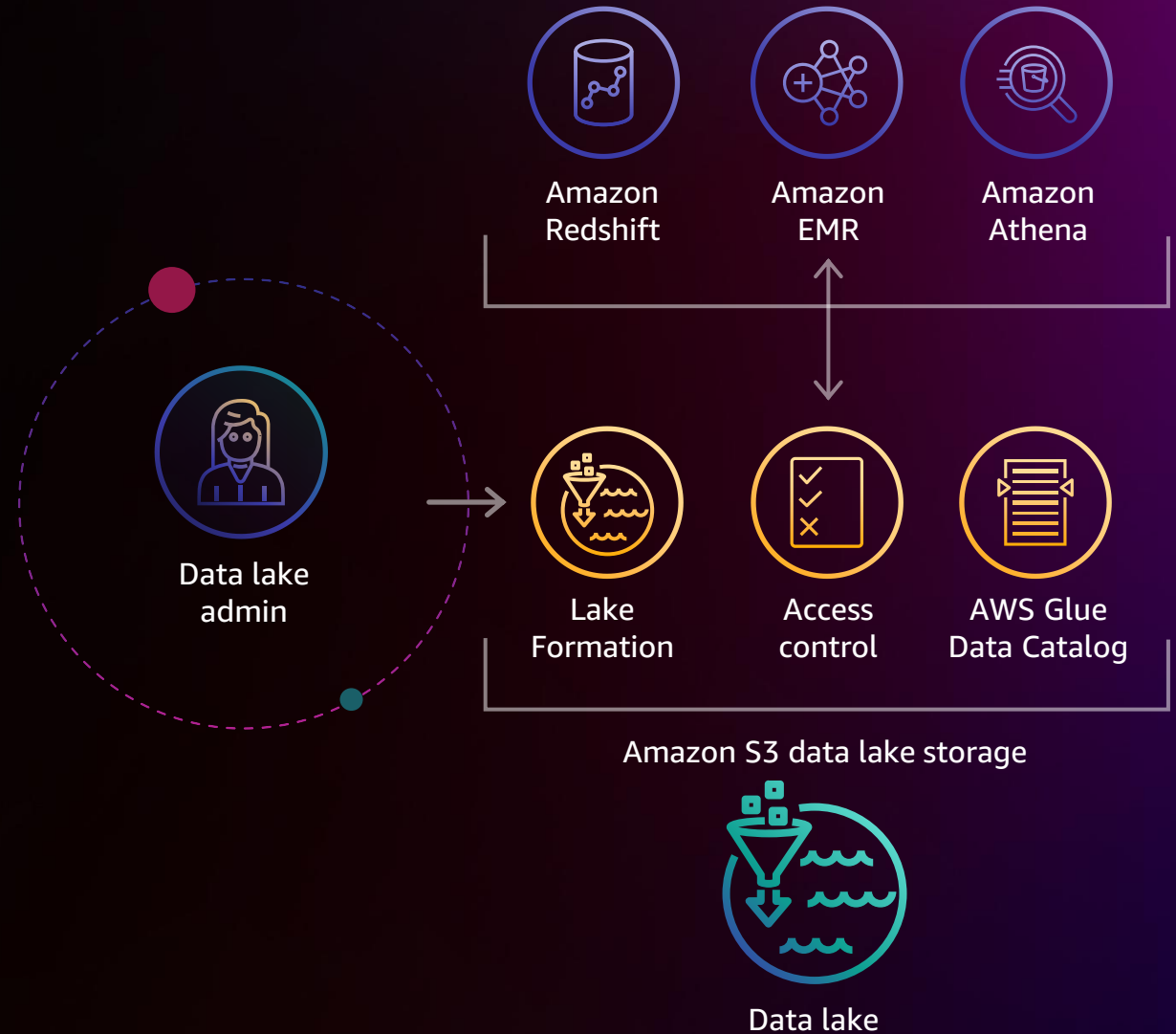


Analytics

Business Intelligence

Databases

Catalog

Governance

Machine Learning

Data Lakes

# Data Governance starts with business



Data Profiling

Data Quality

Business Strategy

Data Classification

Secure, Store and Encrypt

**OPERATING MODEL**
—
**PEOPLE, PROCESS, TECHNOLOGY**

**Ingest and Integrate**

**Catalog and Share**

Producers
"Teams that want to share data"

Consumers
"Teams that want to use data"

Data Lifecycle

Secure, Store and Encrypt

Manage, Master and and Optimize

Data Curation, Data Integration and Lineage

# AWS Lake Formation unifies data governance

Simplify security management with **Lake Formation**

Data lake admin

Amazon Redshift

Amazon EMR

Amazon Athena

Lake Formation

Access control

AWS Glue Data Catalog

Amazon S3 data lake storage

Data lake

# Core components of Amazon DataZone

**Data producers**

Bring data from different sources, and across accounts/Regions

**Amazon DataZone**

Data portal

APIs

**Organizational domains**

Business data catalog

Data projects

Governance and access control

**Data consumers**

Simplify access, collaboration, and consumption using different tools

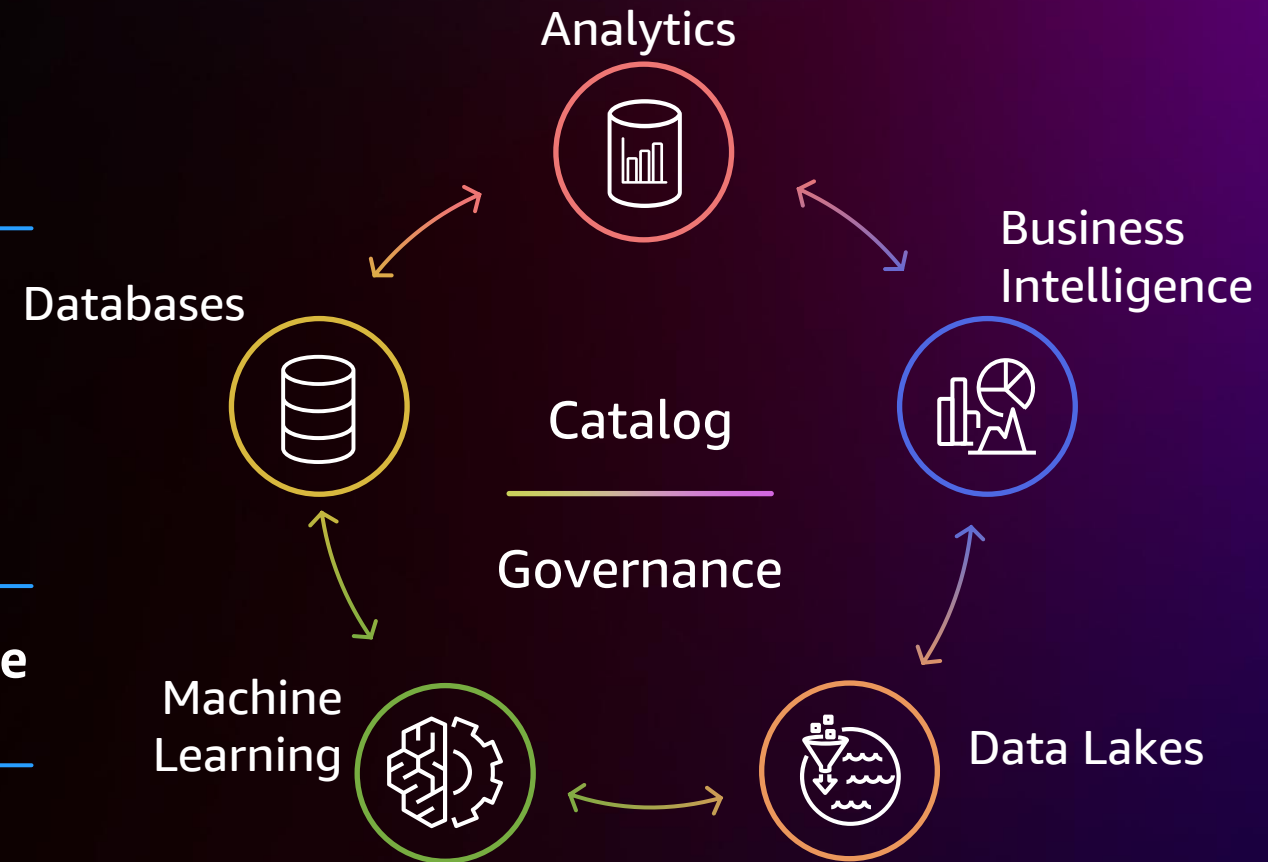# Modern data architecture's five benefits

**Unified analytics**

**Highest performance at the lowest cost**

**Machine learning integration**

**Unified data access, security and governance**

**Insights for everyone**

Analytics

Business Intelligence

Databases

Catalog

Governance

Machine Learning

Data Lakes

# Insights for everyone

Focus on data
without managing
infrastructure

Choose your tools based
on your skillset

# Serverless is a key for your data infrastructure

## The benefits of serverless

Faster time to market

Zero infrastructure management

Pay for what you use

Automatic scaling

Compute provisioning

Automated patching

Automatic failover

Advanced monitoring

Backup and recovery

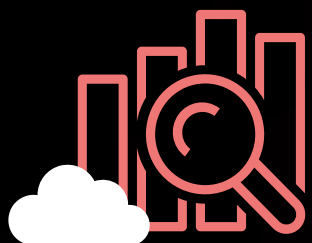Routine maintenance

Security and industry compliance

YOU

**focus on insights**

aws

**takes care of the rest**

# A full stack of serverless options for data analytics in the cloud

**AWS Glue**
Data integration, ETL, and Catalog

**Amazon AppFlow**
SaaS integration

**Amazon OpenSearch Service**
Log and search analytics

**Amazon Redshift**
Data warehousing

**Amazon MSK**
Real-time analytics

**AWS Analytics**

**Amazon EMR**
Big data processing

**Amazon QuickSight**
Visualization

**Amazon Athena**
Interactive analytics

**AWS Lake Formation**
Data lake setup management and governance

**Amazon Kinesis**
Real-time analytics

# Amazon OpenSearch Serverless

Amazon OpenSearch Service securely unlocks real-time search, monitoring, and analysis of operational data

## Easy to administer
No sizing, scaling, and tuning of clusters, and no shard and index lifecycle management

## Fast
Automatically scale resources to maintain consistently fast data ingestion rates and query response times

## Ecosystem
Get started in seconds using the same OpenSearch Service clients, pipelines, and APIs

## Cost-effective
Pay only for the resources consumed
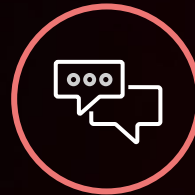
# Amazon SageMaker Canvas

Build ML models and generate accurate predictions — no code required

- Quickly access and prepare data for Machine Learning
- Built-in AutoML to build models and generate accurate predictions
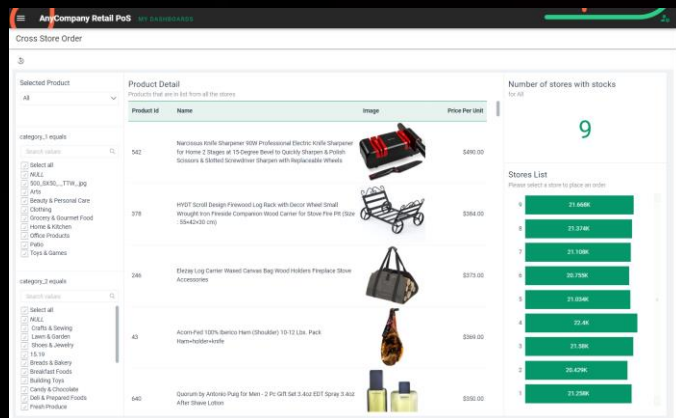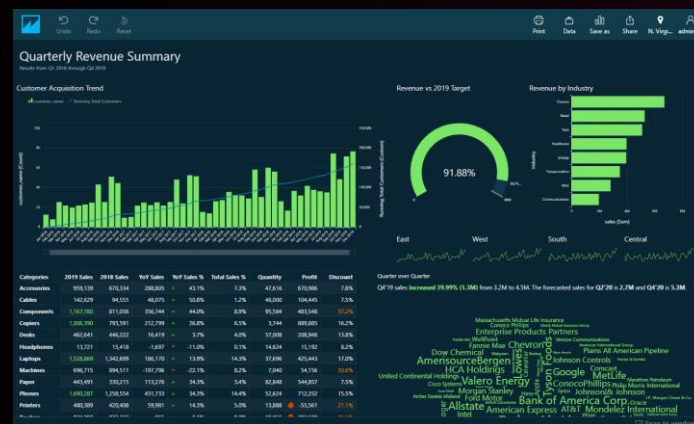- Share ML models and collaborate with data science teams
- Usage-based pricing to avoid licensing fees and reduce TCO

# Insights through Amazon QuickSight



## Embedded insights

Enhance customer-facing products and monetize data assets

## Interactive dashboards, meaningful insights

Dashboards, visualizations, and ad-hoc analysis primarily for internal audiences

## Enterprise reporting

Static, highly formatted, email-based reporting distributed to large internal or external audiences

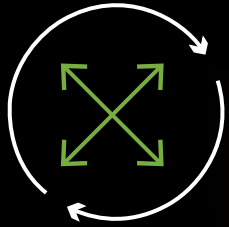# Amazon QuickSight Q

Ask natural language questions about your data
and get answers in seconds

Type your question and get instant answer

# Modern Data Architecture is making it easier to unlock the value of data across the end-to-end data journey

### Scalable
Performance at scale

### Unified
Connect to all your data

### Comprehensive
Tools for all your workloads

### Governed
End-to-end governance

# Agenda

Modern data architecture on AWS

End-to-end data life cycle on the modern data architecture

Data governance and data mesh in action

Journey towards modern data architecture

# End to end data life cycle

| Ingest | Store | Transform & Catalog | Analyze & Visualize | Predict | Share |
|--------|-------|---------------------|---------------------|---------|-------|
| From any source including on-premises, real-time | Any amount of data at Exabyte scale | Data preparation, transformation, and providing seamless data access | End-to-end analytics & visualization for any use case | Most comprehensive set of ML and AI Services | With just a few lines of code |

**Unified Security, Governance, and Data Access**

# Typical architecture

## #1: Ingest

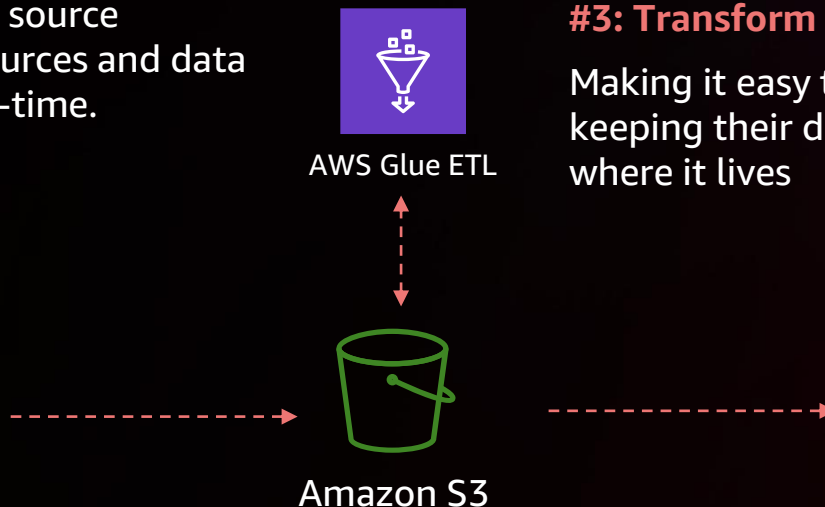Ingesting data from any source including on-premise sources and data that is generated in real-time.

On-premises

Streams

Databases

Logs

AWS Glue ETL

Amazon S3

## #2: Store

Storing both transactional data in databases and analytical data in data warehouses and data lakes at any scale.

## #3: Transform & Catalog
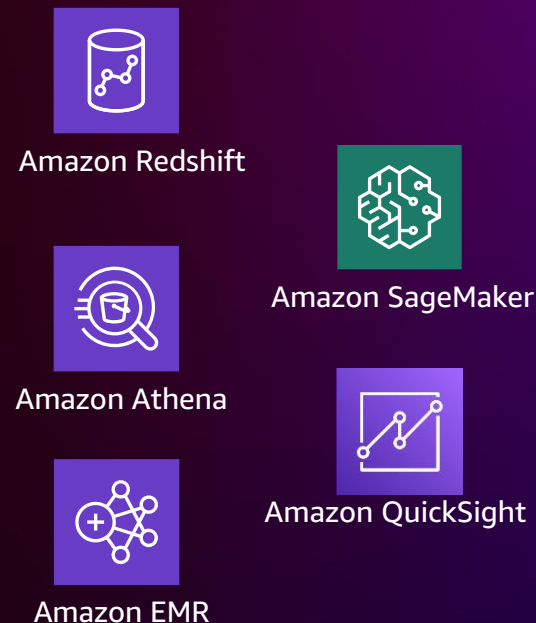
Making it easy to access their data and keeping their data in sync regardless of where it lives

AWS Glue Data Catalog

## #4: Analyze & Visualize

Analyzing data using any of ad hoc queries, distributed frameworks and search engines, and visualize the data on dashboards

Amazon Redshift

Amazon Athena

Amazon EMR

Amazon SageMaker

Amazon QuickSight

## #5: Predict

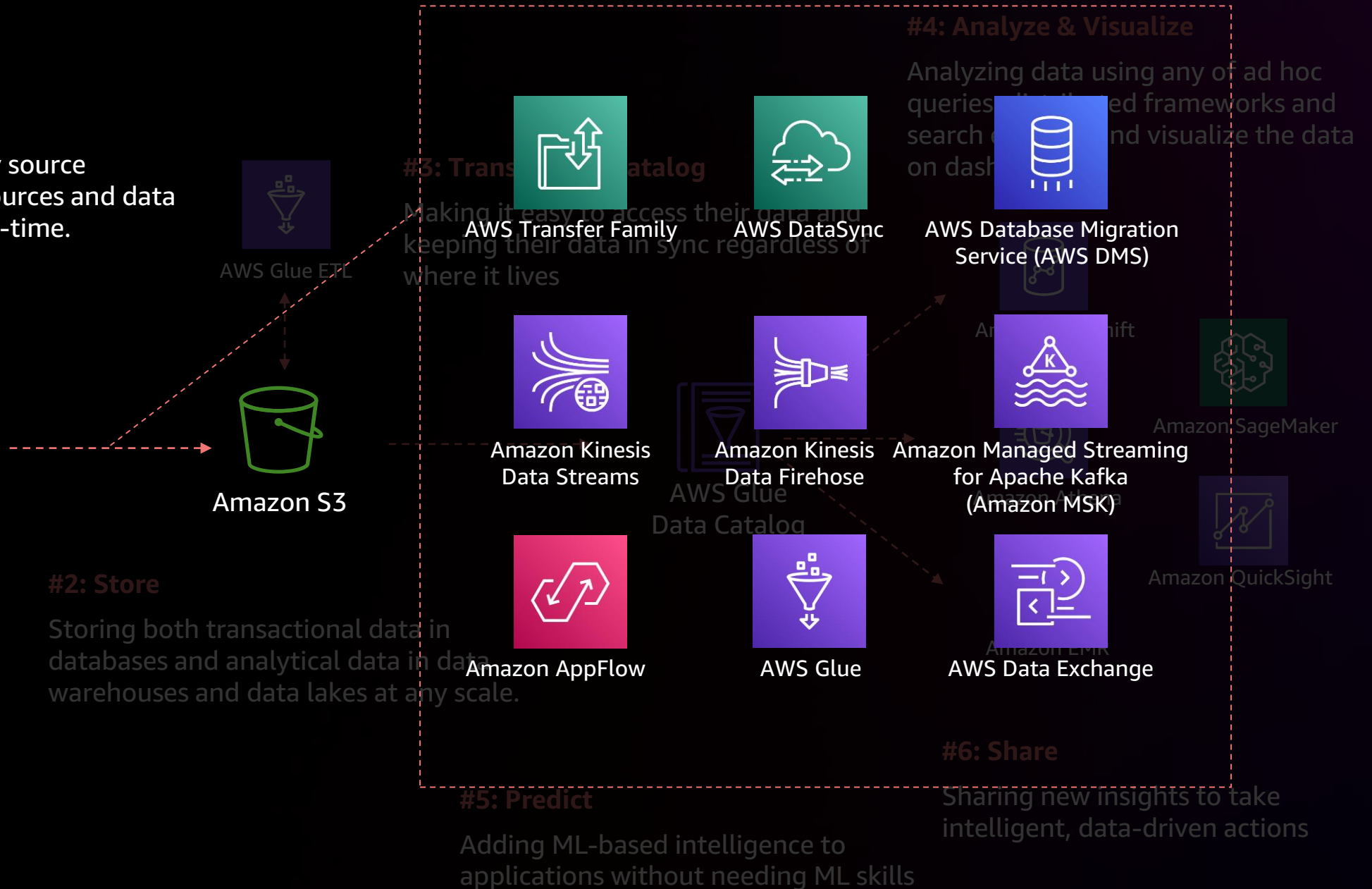Adding ML-based intelligence to applications without needing ML skills

## #6: Share

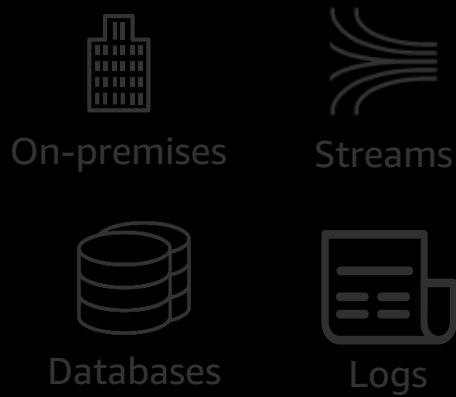Sharing new insights to take intelligent, data-driven actions

aws

# Ingest

## #1: Ingest

Ingesting data from any source including on-premise sources and data that is generated in real-time.

**On-premises**

**Streams**

**Databases**

**Logs**

**Amazon S3**

AWS Glue ETL

## #3: Transform & Catalog

Making it easy to access their data and keeping their data in sync regardless of where it lives

**AWS Transfer Family**

**AWS DataSync**

## #4: Analyze & Visualize

Analyzing data using any of ad hoc queries distributed frameworks and search engines and visualize the data on dashboards

**AWS Database Migration Service (AWS DMS)**

Amazon Redshift

AWS Glue Data Catalog

**Amazon Kinesis Data Streams**

**Amazon Kinesis Data Firehose**

**Amazon Managed Streaming for Apache Kafka (Amazon MSK)**

Amazon SageMaker

Amazon Athena

## #2: Store

Storing both transactional data in databases and analytical data in data warehouses and data lakes at any scale.

**Amazon AppFlow**

**AWS Glue**

**AWS Data Exchange**

Amazon EMR

Amazon QuickSight

## #6: Share

Sharing new insights to take intelligent, data-driven actions

## #5: Predict

Adding ML-based intelligence to applications without needing ML skills

# Store



## #1: Ingest

Ingesting data from any source including on-premise sources and data that is generated in real-time.
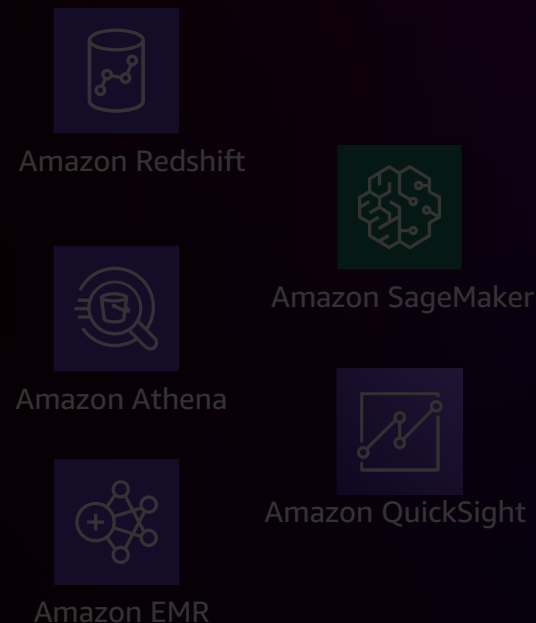
## #3: Transform & Catalog

Making it easy to access their data and keeping their data in sync regardless of where it lives

## #4: Analyze & Visualize

Analyzing data using any of ad hoc queries, distributed frameworks and search engines, and visualize the data on dashboards

AWS Glue ETL

On-premises    Streams

Databases    Logs

Amazon S3    Amazon Redshift    AWS Glue Data Catalog

Amazon Redshift

Amazon SageMaker

Amazon Athena

Amazon QuickSight

Amazon EMR

## #2: Store

Storing both transactional data in databases and analytical data in data warehouses and data lakes at any scale.

## #6: Share

Sharing new insights to take intelligent, data-driven actions

## #5: Predict

Adding ML-based intelligence to applications without needing ML skills

# Amazon S3 for data lakes

**AN OBJECT STORAGE SERVICE OFFERING INDUSTRY-LEADING SCALABILITY, DATA AVAILABILITY, SECURITY, AND PERFORMANCE**

**Durability, availability,** and **scalability**

**Easy to use** with **cost optimization:** Intelligent tiering

Most ways to **get data in**

Amazon S3

Most **object-level controls**

**Broad portfolio of analytics tools**

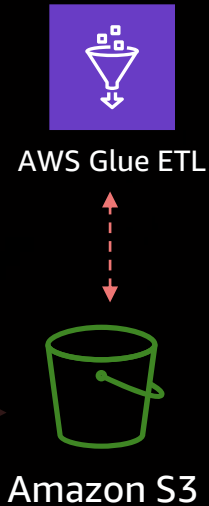**Security, compliance,** and **audit** capabilities

**Cold storage and archive** capabilities

# Transform & Catalog

**#1: Ingest**

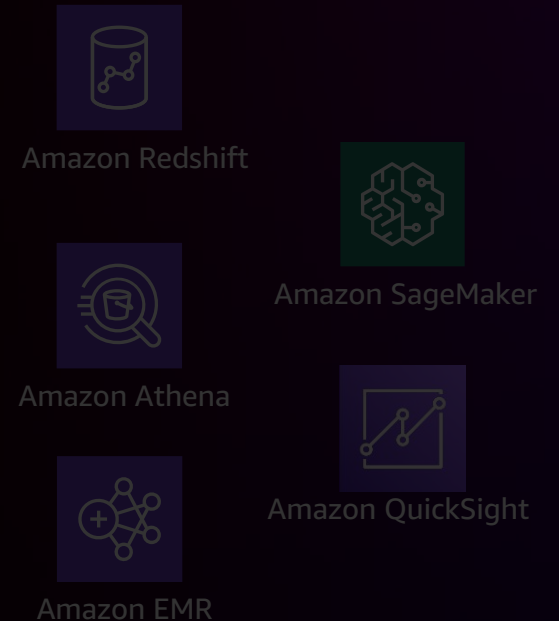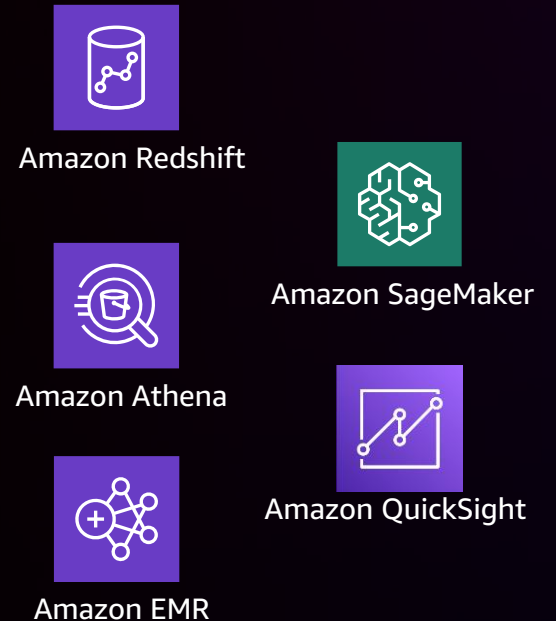Ingesting data from any source including on-premise sources and data that is generated in real-time.

**#3: Transform & Catalog**

Making it easy to access their data and keeping their data in sync regardless of where it lives

**#4: Analyze & Visualize**

Analyzing data using any of ad hoc queries, distributed frameworks and search engines, and visualize the data on dashboards

AWS Glue ETL

On-premises

Streams

Databases

Logs

Amazon S3

AWS Glue Data Catalog

Amazon Redshift

Amazon SageMaker

Amazon Athena

Amazon QuickSight

Amazon EMR

**#2: Store**

Storing both transactional data in databases and analytical data in data warehouses and data lakes at any scale.

**#6: Share**

Sharing new insights to take intelligent, data-driven actions

**#5: Predict**

Adding ML-based intelligence to applications without needing ML skills

# Analyze & Visualize

**#4: Analyze & Visualize**

Analyzing data using any of ad hoc queries, distributed frameworks and search engines, and visualize the data on dashboards

**#1: Ingest**

Ingesting data from any source including on-premise sources and data that is generated in real-time.

**#3: Transform & Catalog**

Making it easy to access their data and keeping their data in sync regardless of where it lives

AWS Glue ETL

On-premises

Streams

Databases

Logs

Amazon S3

AWS Glue Data Catalog

Amazon Redshift

Amazon SageMaker

Amazon Athena

Amazon QuickSight

Amazon EMR

**#2: Store**

Storing both transactional data in databases and analytical data in data warehouses and data lakes at any scale.

**#6: Share**

Sharing new insights to take intelligent, data-driven actions

**#5: Predict**

Adding ML-based intelligence to applications without needing ML skills

# Predict and share

**#1: Ingest**

Ingesting data from any source
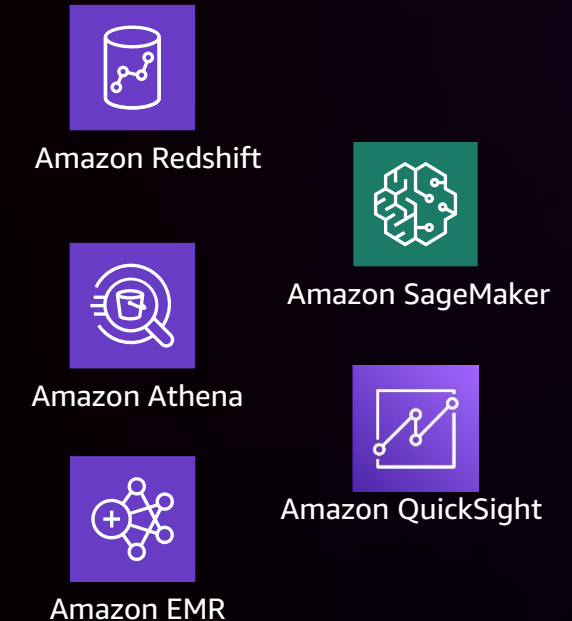including on-premise sources and data
that is generated in real-time.

**#3: Transform & Catalog**

Making it easy to access their data and
keeping their data in sync regardless of
where it lives

**#4: Analyze & Visualize**

Analyzing data using any of ad hoc
queries, distributed frameworks and
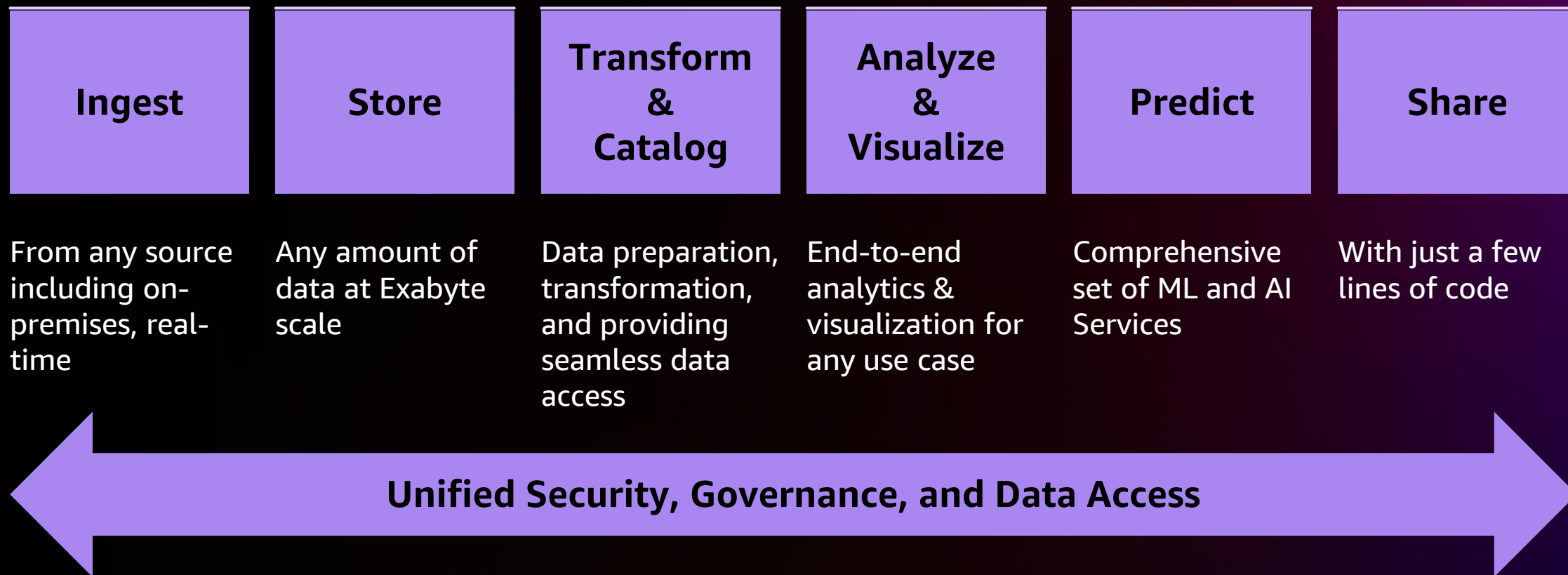search engines, and visualize the data
on dashboards

AWS Glue ETL

On-premises

Streams

Databases

Logs

Amazon S3

AWS Glue
Data Catalog

Amazon Redshift

Amazon Athena

Amazon EMR

Amazon SageMaker

Amazon QuickSight

**#2: Store**

Storing both transactional data in
databases and analytical data in data
warehouses and data lakes at any scale.

**#5: Predict**

Adding ML-based intelligence to
applications without needing ML skills

**#6: Share**

Sharing new insights to take
intelligent, data-driven actions

# Unified Security, Governance, and Data Access

| Ingest | Store | Transform & Catalog | Analyze & Visualize | Predict | Share |
|--------|-------|---------------------|---------------------|---------|-------|
| From any source including on-premises, real-time | Any amount of data at Exabyte scale | Data preparation, transformation, and providing seamless data access | End-to-end analytics & visualization for any use case | Comprehensive set of ML and AI Services | With just a few lines of code |

**Unified Security, Governance, and Data Access**

# Agenda

Modern data architecture on AWS

End-to-end data life cycle on the modern data architecture

Data governance and data mesh in action

Journey towards modern data architecture

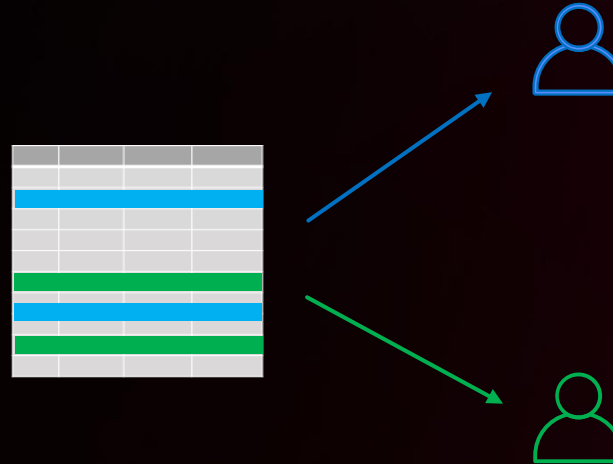# Challenges securing and sharing data

On-premises

Streams

Databases

Logs

Amazon S3

AWS Glue
Data Catalog

AWS Glue ETL

Amazon Athena

Amazon EMR

Amazon Redshift
Spectrum

Amazon
SageMaker

**Challenge #2: Data sharing**

Sharing across accounts and
organizations is cumbersome

**Challenge #1: Security and governance**

Managing fine grained permissions at
scale is difficult and error-prone

# Why is managing data lake permissions hard?

**Unifying permissions across the data lake stack**

Analytics engine

Data Catalog

Storage

Split storage, metadata, and compute

Each system has different permissions

Syncing permissions is error-prone

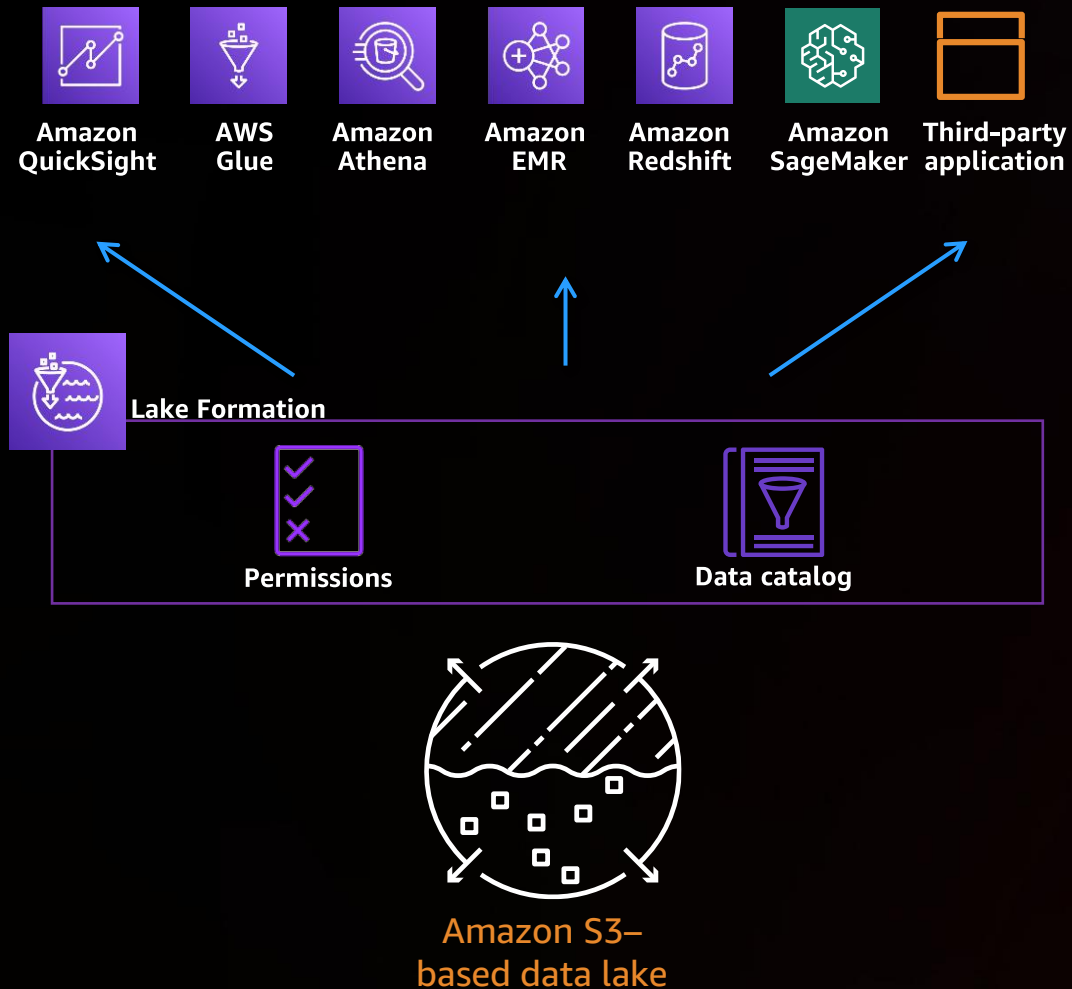**Enforcing fine-grained permissions to restrict access**

Data lakes contain a lot of data

Users should only access portions

Thousands of resources and tens of thousands of users

**Ensuring that data access complies with regulations**

Democratize data access

Regulations and governance

Monitor and audit data access

# Lake Formation permissions model

Amazon QuickSight · AWS Glue · Amazon Athena · Amazon EMR · Amazon Redshift · Amazon SageMaker · Third-party application

**Lake Formation**

Permissions

Data catalog

Amazon S3–based data lake

DB-style fine-grained permissions on resources

Scale permissions management Lake Formation Tag-Based Access Control (LF-TBAC)

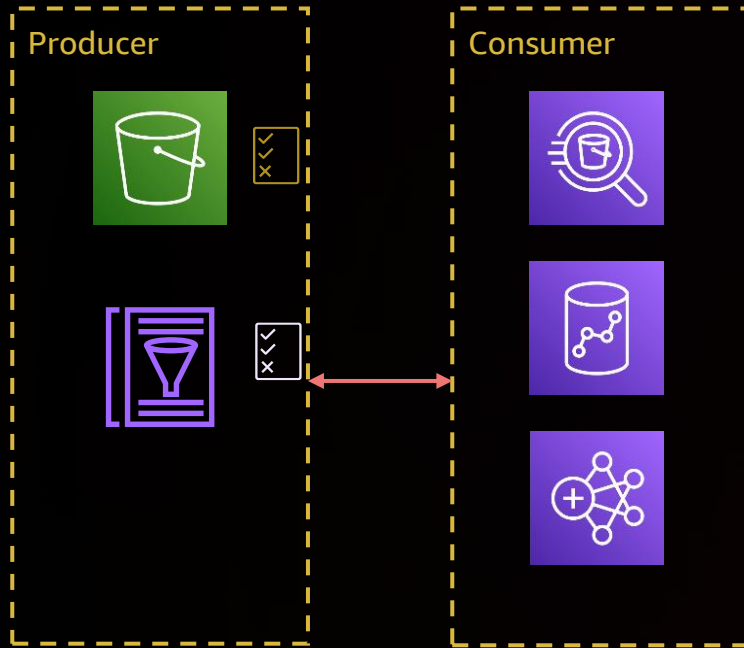Unified Amazon S3 permissions

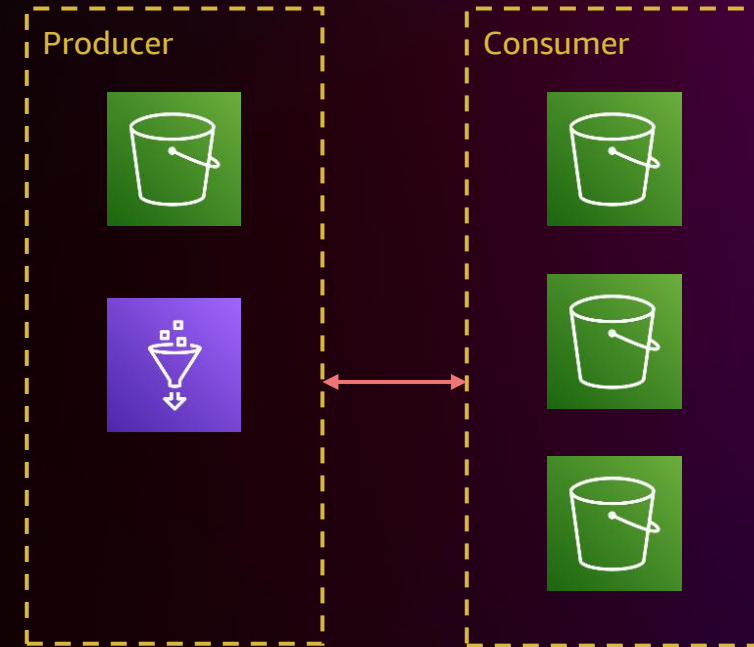Integrated with services and tools

Easy to audit permissions and access

# Challenge: Data sharing

Sharing across accounts and organizations is cumbersome

On-premises

Streams

Databases

Logs

Amazon S3

AWS Glue
Data Catalog

AWS Glue ETL

Amazon Athena

Amazon EMR

Amazon Redshift Spectrum

Amazon SageMaker

# Why is sharing data across accounts hard?

## To share data

**Producer**

**Consumer**

Manges multiple Amazon S3 and IAM policies

Lacks discoverability

Policy size limits (coarse grained)

## Duplicating data

**Producer**

**Consumer**

ETL pipelines

Multiple redacted copies

Expensive, brittle, and error-prone

# Common data sharing patterns

**Single account**

**Centralized solution
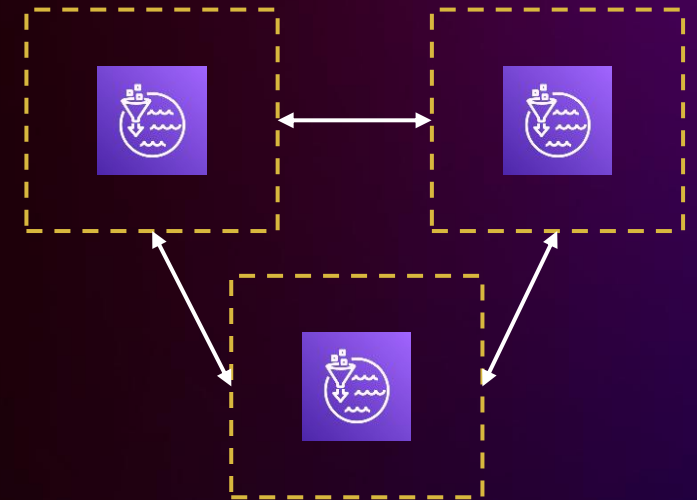(Hub and spoke)**

**Data mesh**

Producer

Consumer

Consumer

Consumer

**Centralized
single account**

**Hub and spoke
multi-account**

**Data mesh
central governance**

**Simple to get started**

**Cross-organization**
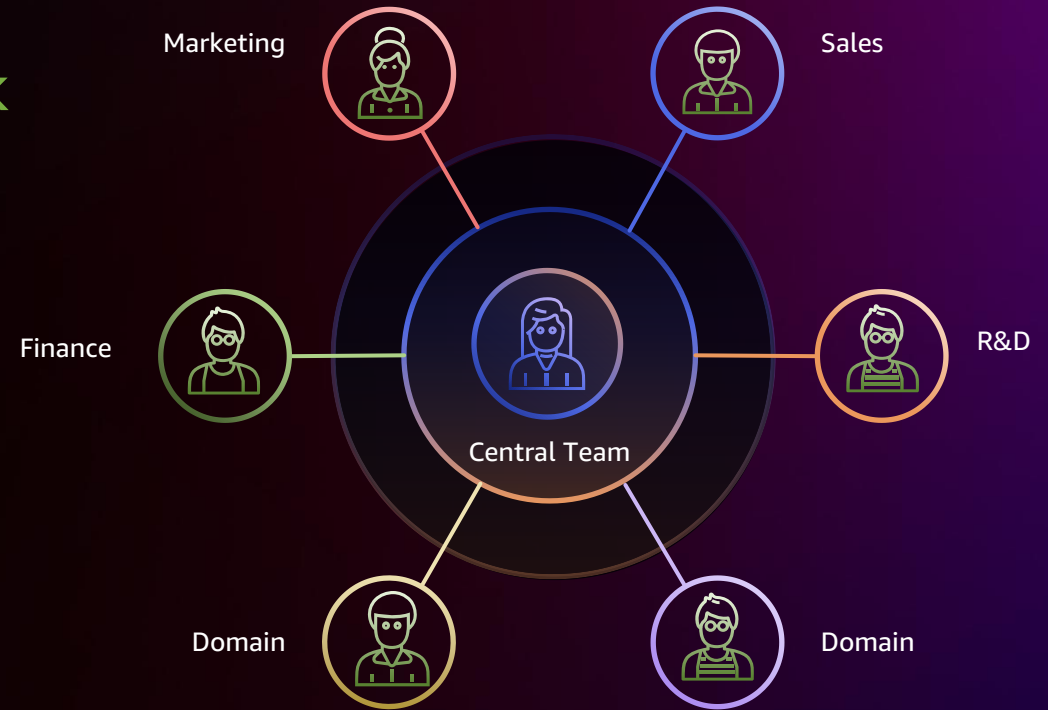
**Organizational autonomy**

# Centralized solution challenges

Central data team becomes the bottleneck

Fail to scale data consumers

Unable to discover and consume the data

No central Data Governance
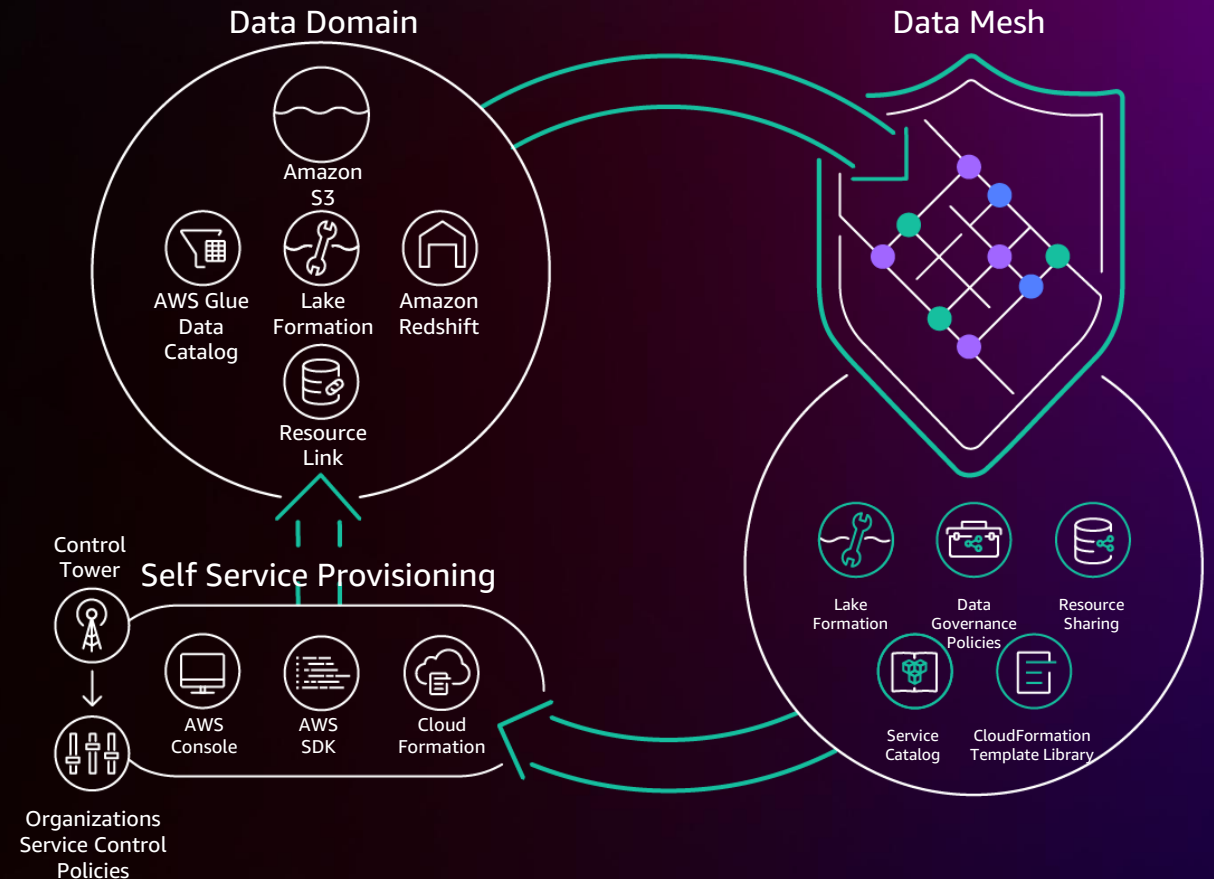
Lack of data auditability

# Why data mesh?

Treats existing data platforms as **independent** domains

Improves **data governance** by pushing access policy to data domains

Establishes a **central** mechanism for **data discovery**
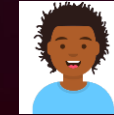
Provides **self-service data sharing** features



Data Domain

Amazon S3

AWS Glue Data Catalog

Lake Formation

Amazon Redshift

Resource Link

Control Tower

Self Service Provisioning

AWS Console

AWS SDK

Cloud Formation

Organizations Service Control Policies

Data Mesh

Lake Formation

Data Governance Policies

Resource Sharing

Service Catalog

CloudFormation Template Library

# Data mesh – Four core principles

**Data Owner**

Data Domain Ownership

**Data Steward**
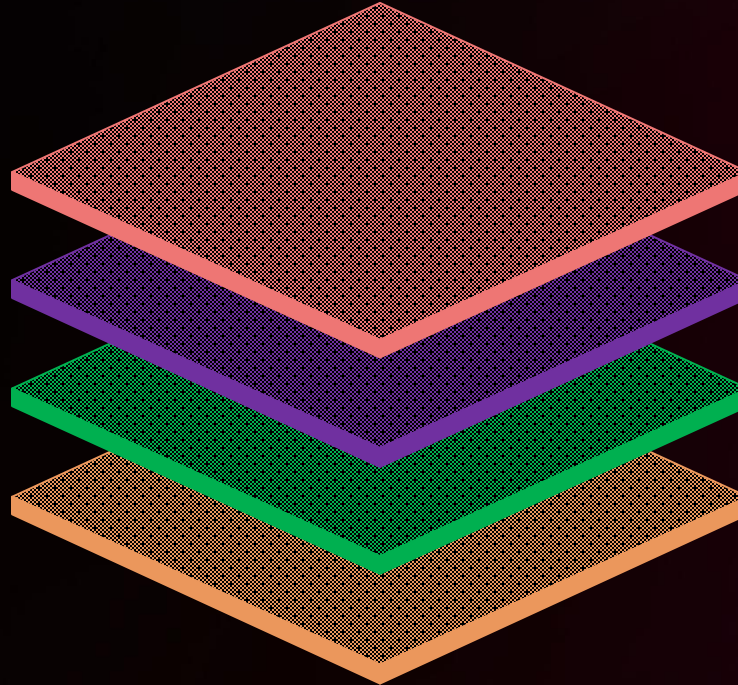
Federated Computational Governance

**Data Engineer**
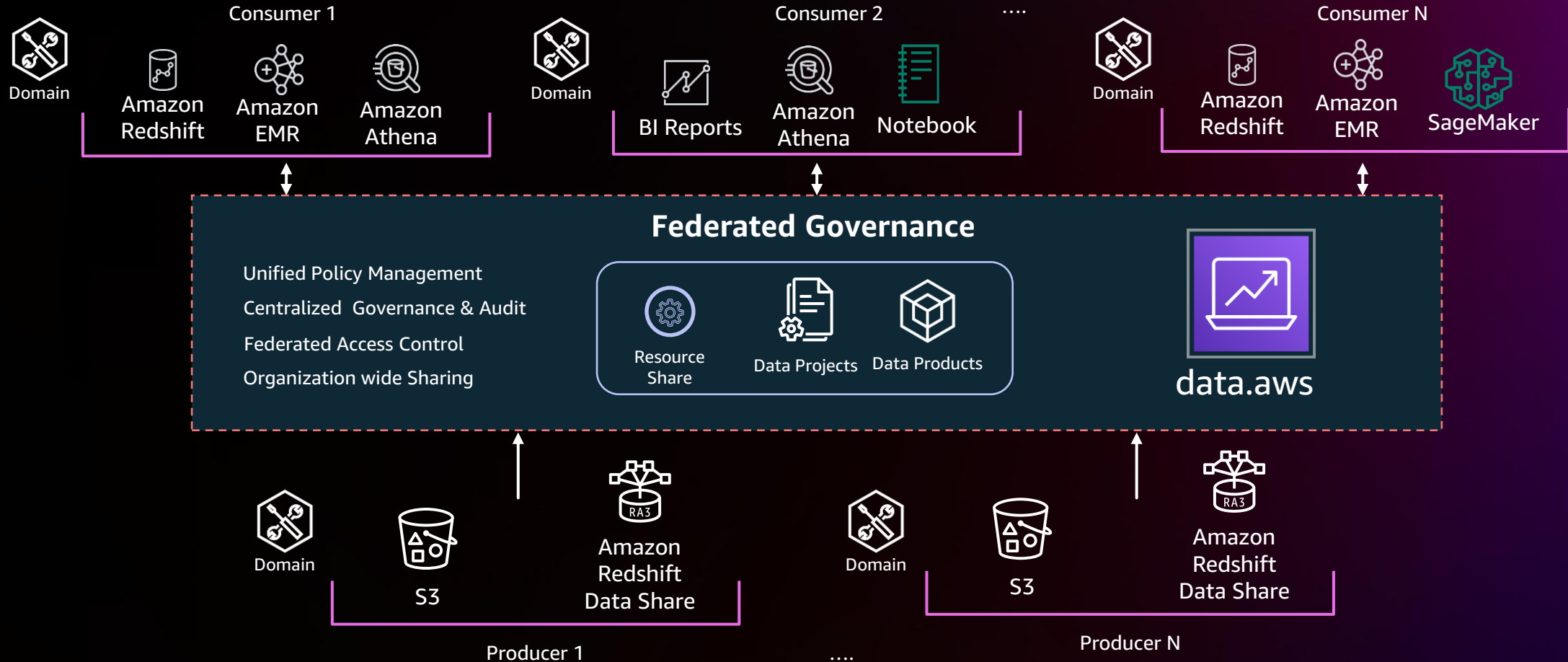
Data as a Product

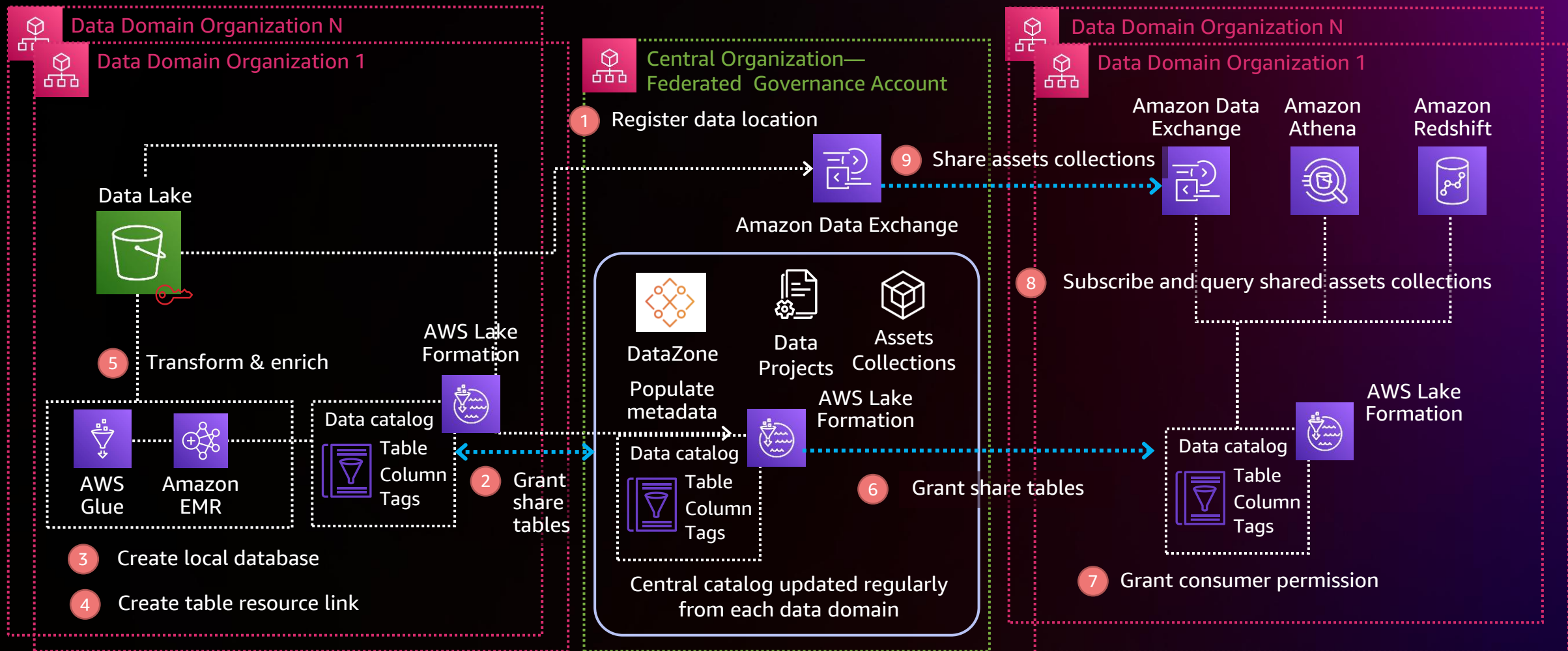**Data Consumer**

Self-Serve Sharing

# Data mesh architecture

DECENTRALIZED, LIGHTWEIGHT FEDERATED GOVERNANCE ACROSS DOMAIN-ORIENTED DATA SYSTEMS TO DRIVE GOVERNED SHARING
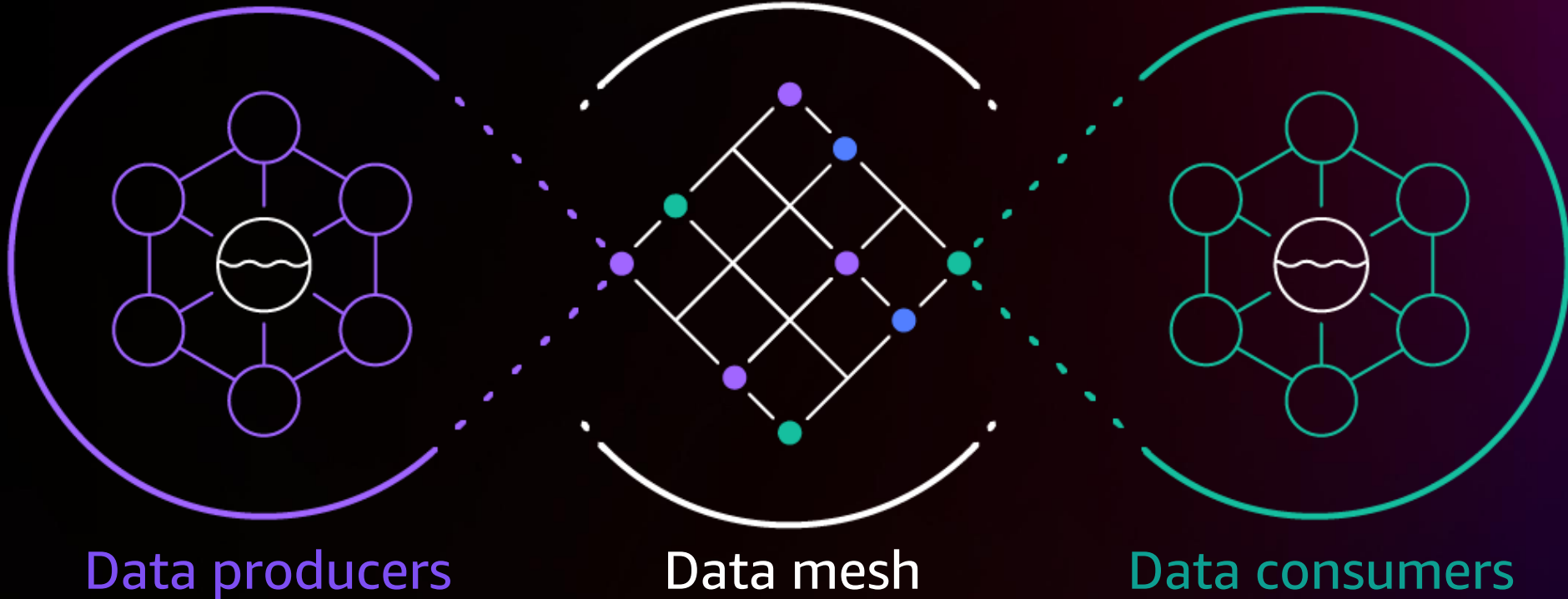
# Data mesh architecture pattern: data lake assets collections sharing

Data Domain Organization N

Data Domain Organization 1

Central Organization—Federated Governance Account

Data Domain Organization N

Data Domain Organization 1

Amazon Data Exchange

Amazon Athena

Amazon Redshift

**1** Register data location

**9** Share assets collections

Amazon Data Exchange

Data Lake

DataZone

Data Projects

Assets Collections

**8** Subscribe and query shared assets collections

**5** Transform & enrich

Populate metadata

AWS Lake Formation

AWS Lake Formation

AWS Lake Formation

Data catalog

Table Column Tags

Data catalog

Table Column Tags

Data catalog

Table Column Tags

AWS Glue

Amazon EMR

**2** Grant share tables

**6** Grant share tables

**3** Create local database

**7** Grant consumer permission

**4** Create table resource link

Central catalog updated regularly from each data domain

aws

# Security, Governance and Data Access with Data Mesh

Data mesh **unifies** Security, Governance, and Data Access of modern data architecture



Data producers          Data mesh          Data consumers

# Agenda

Modern data architecture on AWS

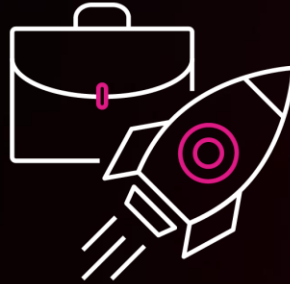End-to-end data life cycle on the modern data architecture

Data governance and data mesh in action

Journey towards modern data architecture

# Journey towards modern data architecture



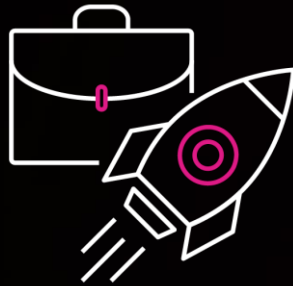Start small



Scale fast



Think big

# Journey towards modern data architecture

**Start small**

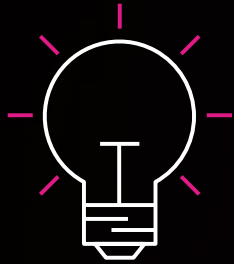Start from **small subset** in your data platform

# Journey towards modern data architecture

Scale fast

Scale **fast** to achieve your business goal based on data

# Journey towards modern data architecture

Think big

Expand your data platform
to support **more workload**
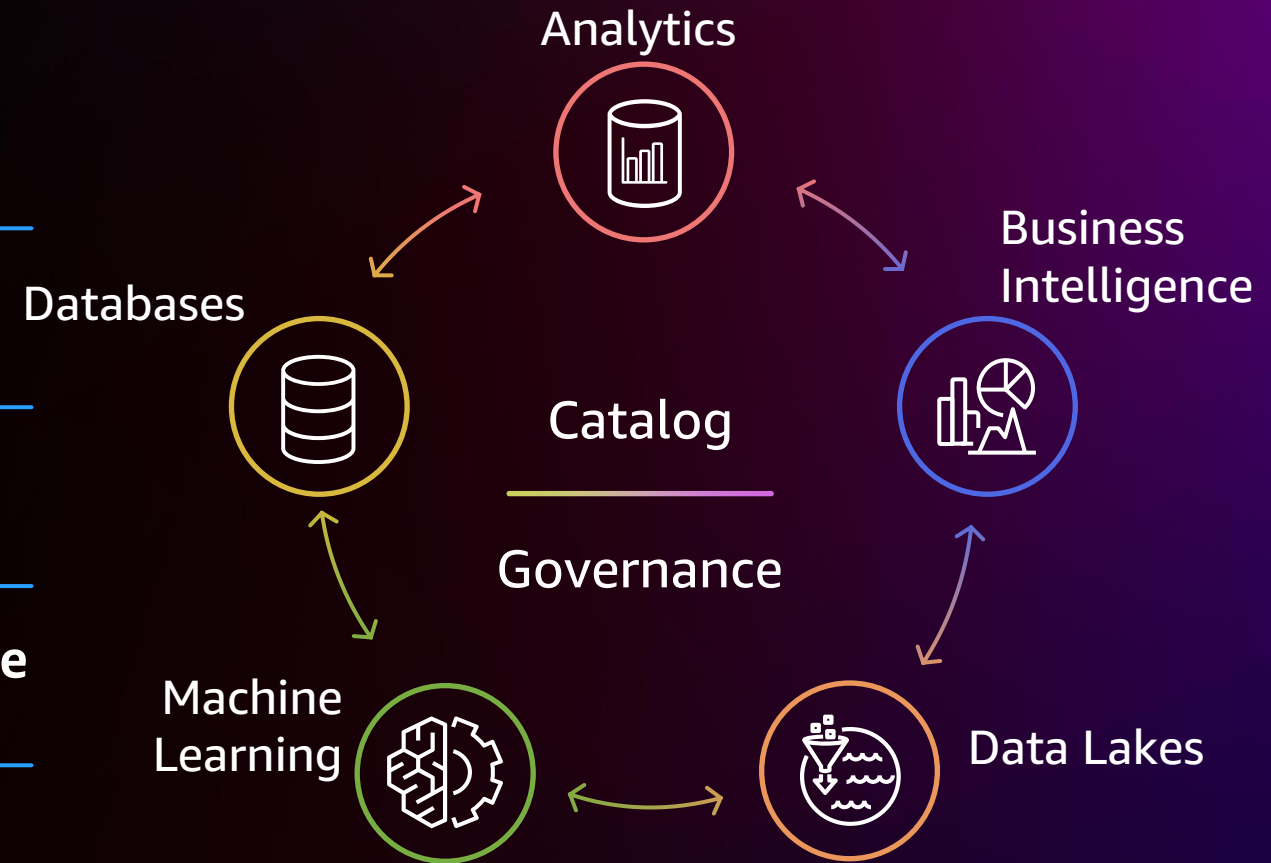and **advanced goals**

# Conclusion

**Unified analytics**

**Highest performance at the lowest cost**

**Machine learning integration**

**Unified data access, security and governance**

**Insights for everyone**



Analytics

Business Intelligence

Databases

Catalog

Governance

Machine Learning

Data Lakes

# Learn more at re:Invent 2022

**Swami Sivasubramanian, Vice President of AWS Data and Machine Learning – Keynote**
Wednesday November 30 | 8:30 AM – 10:30 AM PST | The Venetian

**ANT203-L (LVL 200) Unlock the value of your data with AWS analytics**
Wednesday November 30 | 2:30 PM – 3:30 PM PST | The Venetian

**ANT223 (LVL 200) Simplify and accelerate data integration & ETL modernization with AWS Glue**
Wednesday November 30 | 12:15 PM – 1:15 PM PST | MGM Chairmans 368

**ANT310 (LVL 300) Build a data mesh with AWS Lake Formation and AWS Glue**
Wednesday November 30 | 05:30 PM – 07:30 PM PST | MGM Grand

**ANT344 (LVL 300) Democratize data with governance – Connect people, data, and tools with Amazon DataZone**
Wednesday November 30 | 02:30 PM – 03:30 PM PST | MGM Grand

# Thank you!

Santosh Chandrachood

sanchas@amazon.com

Please complete the session survey in the **mobile app**