# Building next-generation data lakes and analytics on AWS

Ninad Phatak
Big Data Solutions Architect, AISPL
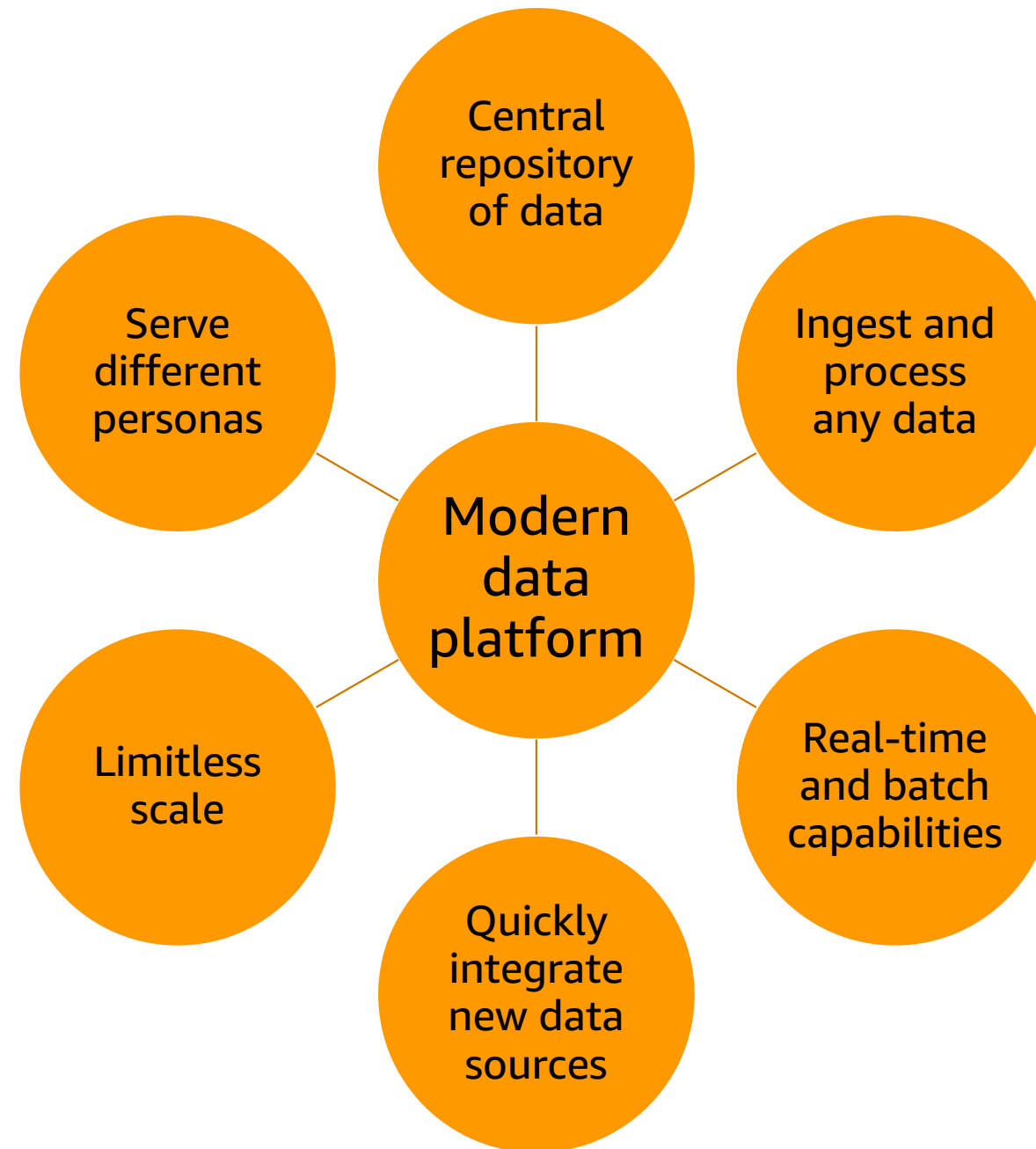
# Agenda

Characteristics of a modern data platform
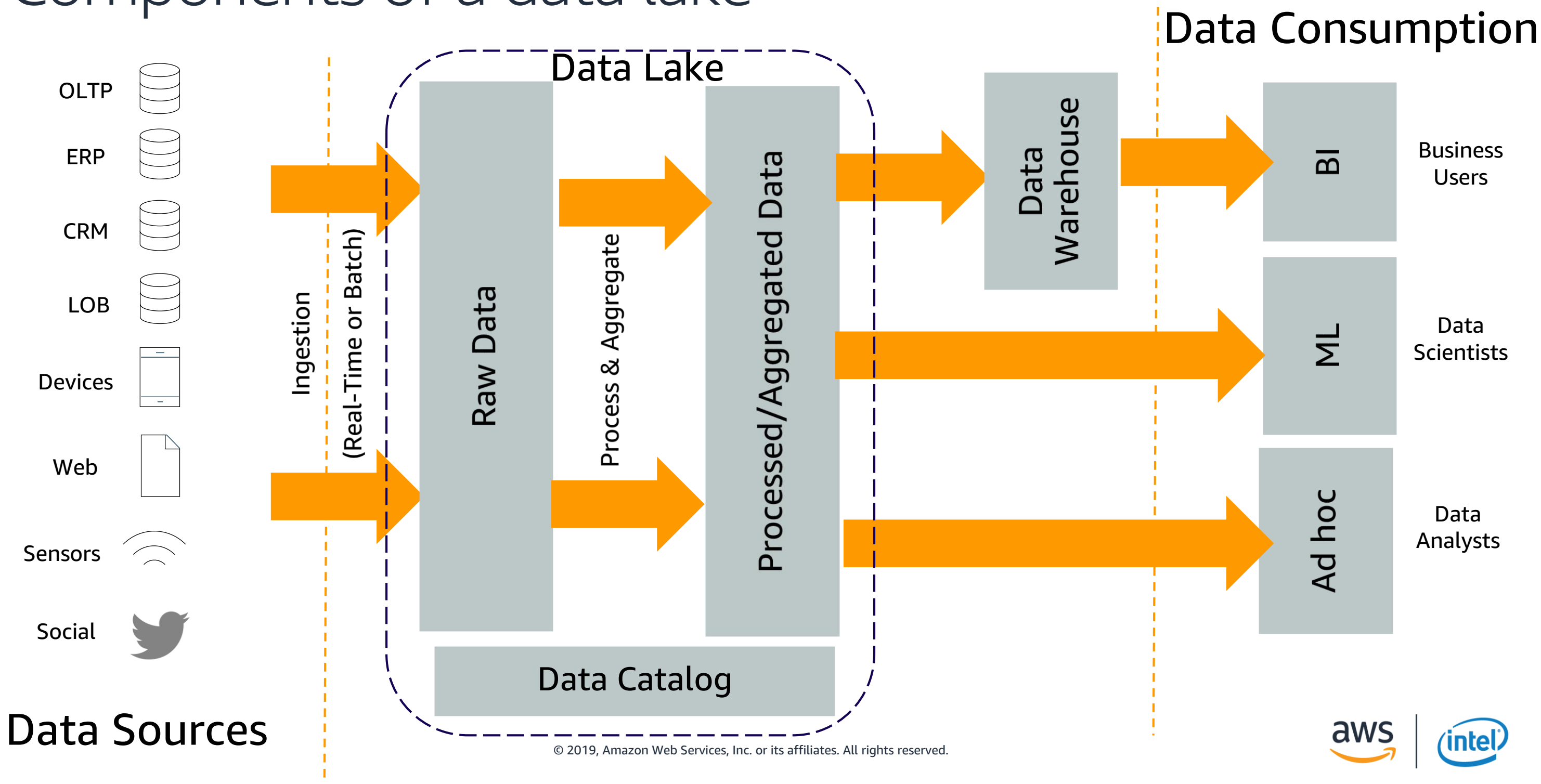
Components of a data lake on AWS

- Core data lake

- Data ingestion

- Data processing

- Consumption of data and insights

AWS lake formation

# A modern data platform/data lake
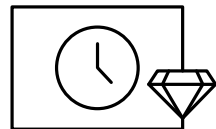
# Components of a data lake

# The core data lake

# Amazon S3: Object storage

Secure, highly scalable, durable object storage with millisecond latency for data access

Store any type of data — web sites, mobile apps, corporate applications, and IoT sensors

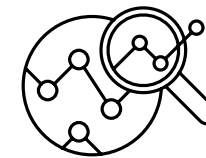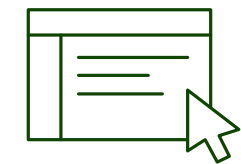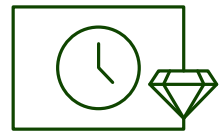| Durability, availability, and scalability | Security and compliance | Query in place | Flexible management |
|---|---|---|---|
| Built for eleven 9s of durability; data distributed across 3 physical facilities in an AWS region; automatically replicated to any other AWS region | Three different forms of encryption; encrypts data in transit when replicating across regions; log and monitor with AWS CloudTrail, use ML to discover and protect sensitive data with Amazon Macie | Run analytics and ML on data lake without data movement; Amazon S3 Select can retrieve subset of data, improving analytics performance by 400% | Classify, report, and visualize data usage trends; objects can be tagged to see storage consumption, cost, and security; build lifecycle policies to automate tiering and retention |

aws | intel

# Amazon S3 Glacier: Backup and archive

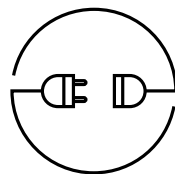Secure, durable, and extremely low-cost storage for data archiving and long-term backup

Store data at $0.004/GB/month

**Durability, availability, and scalability**

Built for eleven 9s of durability; data distributed across 3 physical facilities in an AWS region; automatically replicated to any other AWS region

**Retrieves data in minutes**

Three retrieval options to fit your use case; expedited retrievals with Amazon S3 Glacier Select can return data in minutes

**Secure**

Log and monitor with AWS CloudTrail; Amazon S3 Glacier Vault Lock enables WORM storage capabilities, helping satisfy compliance requirements
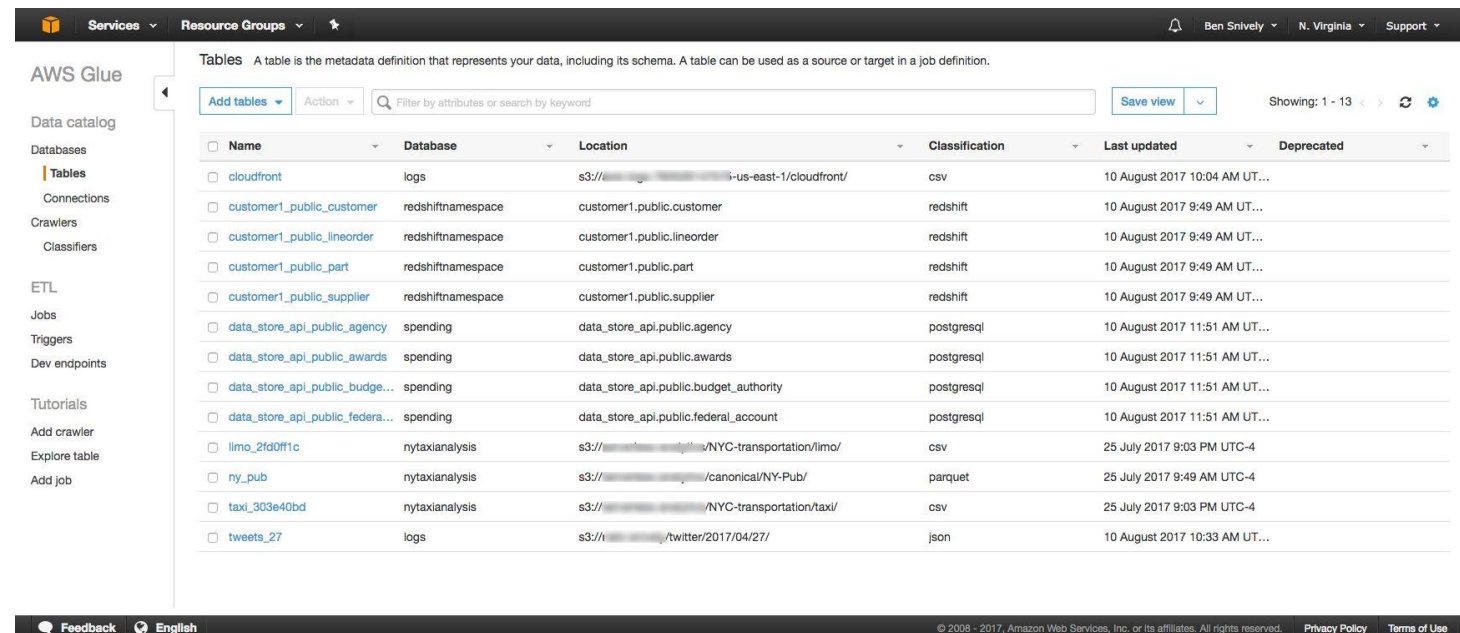
**Inexpensive**

Lowest-cost AWS object storage class, allowing you to archive large amounts of data at a very low cost

aws | intel

# AWS Glue Data Catalog

Bring in metadata from a variety of data sources into a single categorized list that is searchable



- Hive metastore-compatible metadata repository

- Used across multiple services

- Automatically discover data using crawlers

- Track data evolution using schema versioning

# Data ingestion

# Batch data ingestion



AWS Glue

sFTP/PUT

Source

ETL Tools on Amazon EC2

Amazon S3 (Raw Data)

# Real-time data ingestion using Amazon Kinesis services

## Collect, process, and analyze data streams in real time

Amazon EMR/Spark/
Amazon SageMaker

Custom code
on Amazon EC2

AWS Lambda

**Kinesis Data Streams**

Ingest store data streams

**Kinesis Data Analytics**

Aggregate, filter, enrich data

**Kinesis Data Firehose**

Egress data streams

Amazon S3

Amazon Redshift

Amazon Elasticsearch Service

Splunk

- Real-time
- Fully managed
- Scalable
- Secure
- Cost-effective

# Amazon Managed Streaming for Kafka (Amazon MSK)



- Fully compatible with Apache Kafka v1.1.1 and 2.1.0

- Clusters are set up automatically

- Provision Apache Kafka brokers and storage

- Create and tear down clusters on demand

- Apache ZooKeeper is under the hood, highly available, and included with each cluster at no additional cost

# Data processing

# Amazon EMR: Enterprise-grade Hadoop & Spark

## Scale to any size

**Amazon EMR**

**Data lake on AWS**

- Scale compute (Amazon EMR) & storage (Amazon S3) independently

- Store and process any amount of data—PB to EBs

- Provision one, hundreds, or thousands of nodes

- Automatic scaling

- Integrated with Spot Instances

# AWS Glue: ETL service
## Make ETL scripting and deployment easy



- Run your Spark code in a serverless manner

- Automatically generate ETL code

- Code is customizable with Python and Scala

- Endpoints provided to edit, debug, and test code

- Jobs are scheduled or event-based

# Consumption of data and insights

# Amazon Redshift: Data warehousing

Fast, powerful, simple, and fully managed data warehouse at 1/10th the cost

Massively parallel, scale from gigabytes to petabytes

**Fast at scale**

Columnar storage technology to improve I/O efficiency and scale query performance

**Open file formats**

Analyze optimized data formats on the latest SSD, and all open data formats in Amazon S3

**Secure**

Audit everything; encrypt data end to end; extensive certification and compliance

**Inexpensive**

As low as $1,000 per terabyte per year, 1/10th the cost of traditional data warehouse solutions; start at $0.25 per hour

aws | intel

# Extend Amazon Redshift queries to Amazon S3

**Amazon Redshift**

**Amazon S3**

**Redshift Query Engine**

**Redshift Data**

**Data Lake**

Query S3 directly or join data across Amazon Redshift and S3

Scale Amazon Redshift compute and storage separately

Support for CSV, Parquet, ORC, Grok, Avro, and more formats

aws | (intel)

# Amazon QuickSight

## No servers to manage

Amazon QuickSight is a fully managed cloud application, meaning there's no upfront cost, software to deploy, capacity planning, maintenance, upgrades, or migrations.

## Pay for what you use

Provide read-only access to interactive dashboards and pay only when your users access them with pay-per-session pricing. There are no upfront costs, no annual commitments, and no charges for inactive users.

## Scalable

Amazon QuickSight automatically scales with your usage and activity, with no need for additional infrastructure. From 10 users to 10,000, Amazon QuickSight seamlessly grows with you.

## Fully integrated

Amazon QuickSight integrates with your data sources and other AWS services like Amazon Redshift, S3, Amazon Athena, Amazon Aurora, Amazon RDS, IAM, AWS CloudTrail, Amazon Cloud Directory, and more—providing you with everything you need to build an end-to-end BI solution.

# Amazon Athena

- Query data from Amazon S3 directly

- Use ANSI SQL with support for complex joins, nested queries, and window functions

- Leverages Presto under the covers

- Pay $5/TB scanned by your query

- Workgroups for cost control

# Amazon SageMaker: Build, train, and deploy ML models at scale

Pre-built notebooks for common problems

Built-in, high-performance algorithms

One-click training on the highest-performing infrastructure

Model optimization

One-click deployment

Fully managed with automatic scaling for 75% less

# Going back to the data lake schematic

# Data lake on AWS



OLTP

ERP

Batch

CRM

LOB

sFTP/PUT

Devices

Web

Real-Time

Sensors

Social

AWS Glue

ETL on EC2

Amazon Kinesis
Data Firehose

**Data Lake**

Amazon S3 (Raw)

Amazon EMR

Amazon S3
(Processed/Aggregate)

AWS Glue Data
Catalog

Amazon Redshift

Amazon QuickSight

Business
Users

Amazon Athena

Data
Analysts

Amazon SageMaker

Data
Scientists

aws | intel

# AWS Lake Formation

# AWS Lake Formation
## Build a secure data lake in days



| Stage 1 | Stage 2 | Stage 3 | Stage 4 |
|---------|---------|---------|---------|
| **Ingest and Register** | **Security and Control** | **Collaborate and Use** | **Monitor and Audit** |
| Automatically ingest, clean, organize, and catalog structured and semi-structured data, including logs and data from relational databases into your data lake. | Define access controls that provide the right data to the right users, groups, and roles. Database, table, and column permissions offer granular security. | Search and discover using catalog metadata. All access is checked against policy, so your data is protected even if tools change or new data arrives. | Review data access history across multiple use cases from a single place and download data access audit logs. |
| **Create data lake** | **Table permissions** | **Search data catalog** | **Monitor activity** |
| **Import data** | **User permissions** | **Add metadata** | **Download logs** |

# AWS Lake Formation relieves data lake heavy lifting

## Loading data

- **Blueprints / Data Importers** – templates for data ingestion, ETL, metadata (schema) and partition management

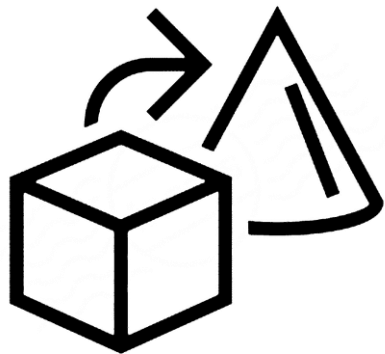- **ML Transformations** – ML algorithms that customers can use to create their own ML Transforms (e.g., record de-duplication, match finding)

# AWS Lake Formation relieves data lake heavy lifting

## Discoverability and governance

- **Enhanced governance layer** – security and governance layer at the Data Catalog level

- **Enhanced Data Catalog** – enable users to record more metadata and tag Data Catalog objects (i.e., databases, tables, columns)

# Learn from AWS experts. Advance your skills and knowledge. Build your future in the AWS Cloud.

**Digital Training**

Free, self-paced online courses built by AWS experts

**Classroom Training**

Classes taught by accredited AWS instructors

**AWS Certification**

Exams to validate expertise with an industry-recognized credential

**Ready to begin building your cloud skills?**
**Get started at: https://www.aws.training/**

aws | intel

# Why work with an APN Partner?

**APN Partners** are uniquely positioned to help your organization at any stage of your cloud adoption journey, and they:

- Share your goals—focused on your success

- Help you take full advantage of all the business benefits that AWS has to offer

- Provide services and solutions to support any AWS use case across your full customer life cycle

**APN Partners with deep expertise in AWS services:**

**AWS Managed Service Provider (MSP) Partners**

APN Partners with cloud infrastructure and application migration expertise

**AWS Competency Partners**

APN Partners with verified, vetted, and validated specialized offerings

**AWS Service Delivery Partners**

APN Partners with a track record of delivering specific AWS services to customers

**Find the right APN Partner for your needs:** https://aws.amazon.com/partners/find/

# Thank you for attending AWS Innovate

We hope you found it interesting! A kind reminder to **complete the survey.**
Let us know what you thought of today's event and how we can improve the event experience for you in the future.

- aws-apac-marketing@amazon.com
- twitter.com/AWSCloud
- facebook.com/AmazonWebServices
- youtube.com/user/AmazonWebServices
- slideshare.net/AmazonWebServices
- twitch.tv/aws