



Democratizing your organization's data analytics experience

ANURAJ PAHUJA

Sr. Solutions Architect
WWPS – State and Local
Government

Amazon Web Services

anurajpa@amazon.com

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Andrew Henderson

Sr. Solutions Architect
WWPS – State and Local Government
Amazon Web Services

andrewhn@amazon.com

Agenda

Cloud strategies and data gravity

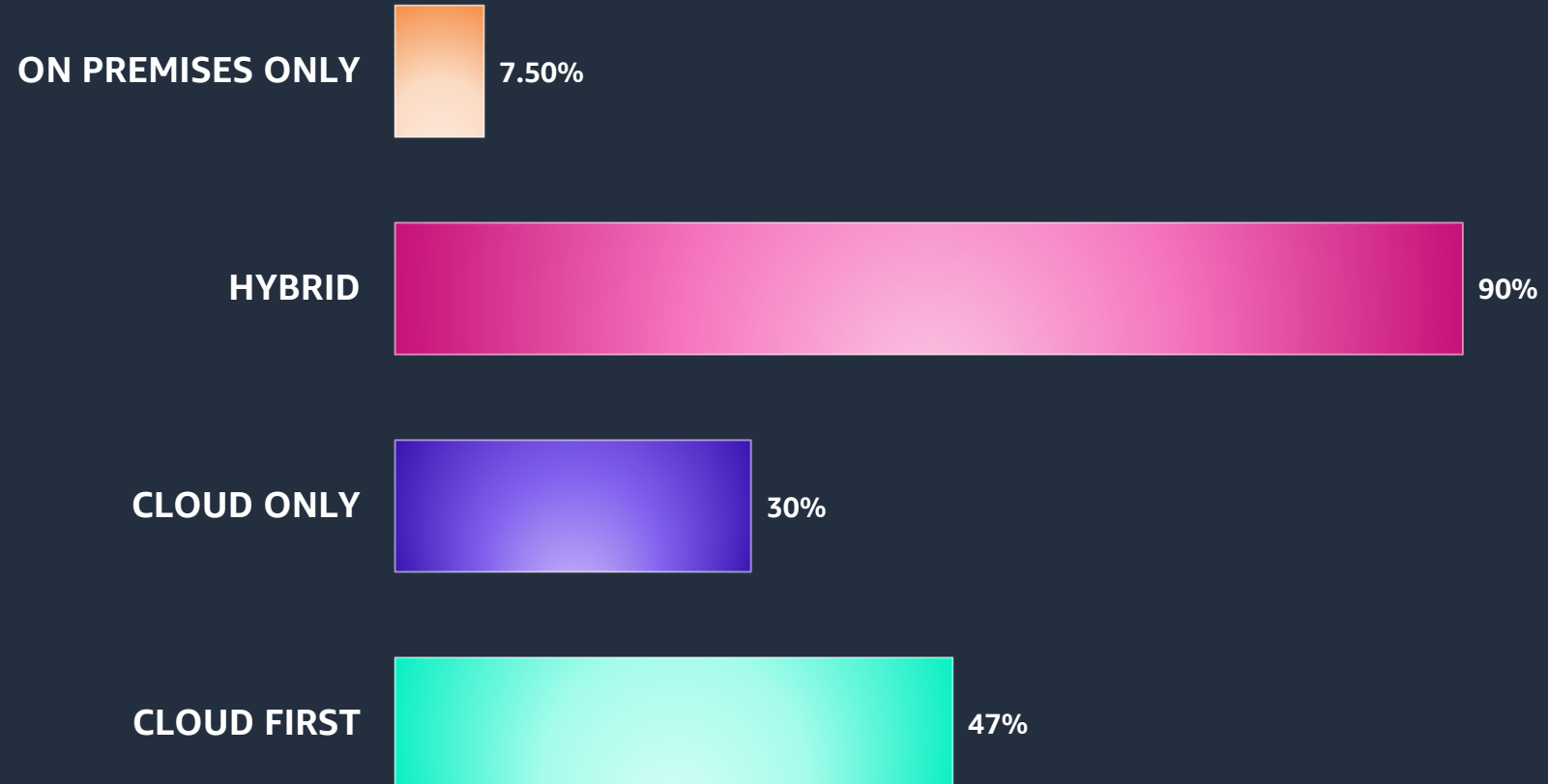
Democratizing analytics

Ease of use

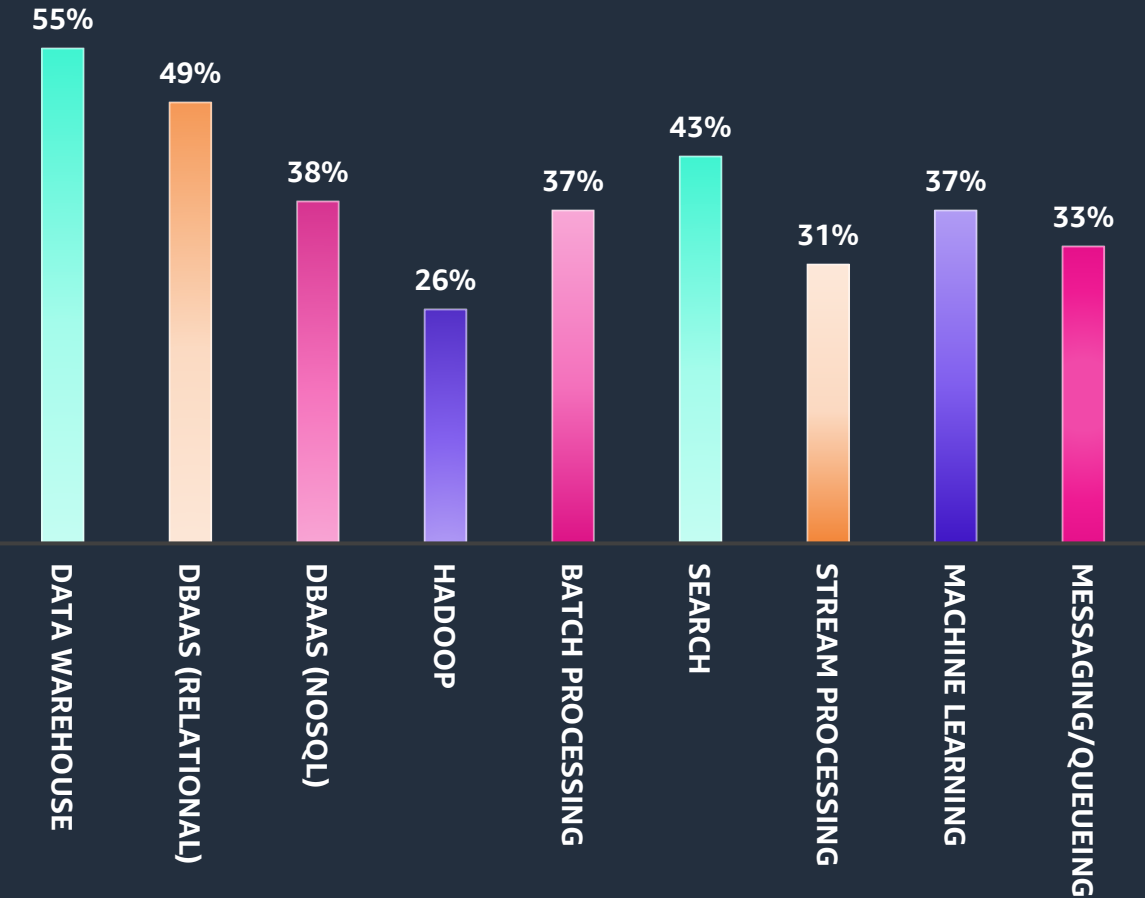
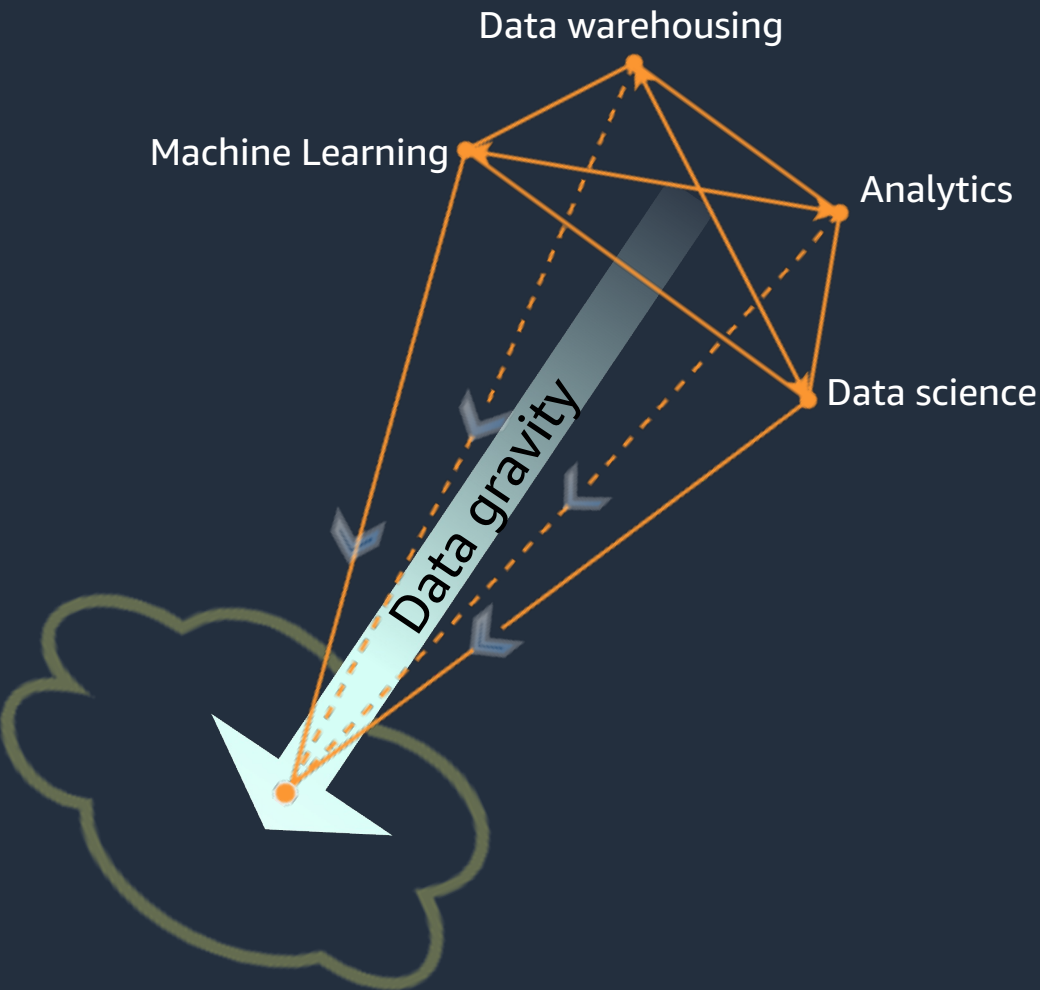
Price performance



Cloud strategies



Data gravity



Source: [Flexera cloud computing trends](#)



Continued Analytics Cloud Growth

The Big Data and Analytics software and cloud services has reached \$90.4B spend in 2021, with 44% deployed in the cloud and the remaining 56% on-premises.

-IDC

Organizations will move more than 70% of their advanced analytics (enriched with AI/ML) to the cloud by 2024.

-Gartner



Data challenges

Cost of data management

Interoperability

Operational freedom

Scale-at-speed

Data driven



A word cloud of data-related terms on a blue sky background. The words are arranged in a scattered, non-uniform pattern. The terms include: Search, Messaging, Interactive analytics, Batch Processing, Blockchain, Streaming data, SaaS, Columnar, Structured data, Data warehouse, Observational data, Data lake, PaaS, IoT data, IaaS, Relational data, Key-value data, Graph data, Machine learning, Transactional data, and Hadoop.

Search

Messaging

Interactive analytics

Batch Processing

Blockchain

Streaming data

SaaS

Columnar

Structured data

Data warehouse

Observational data

Data lake

PaaS

IoT data

IaaS

Relational data

Key-value data

Graph data

Machine learning

Transactional data

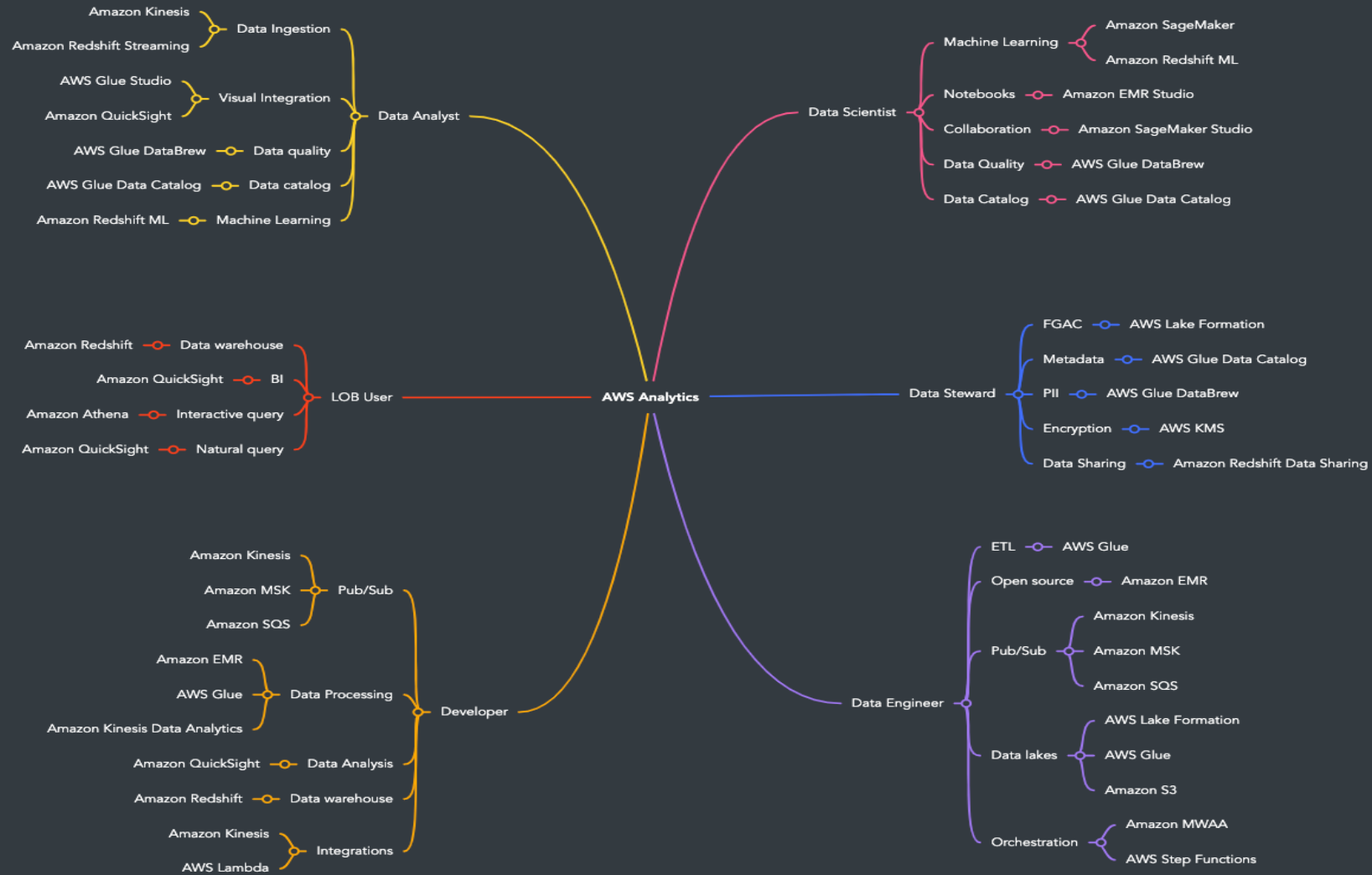
Hadoop

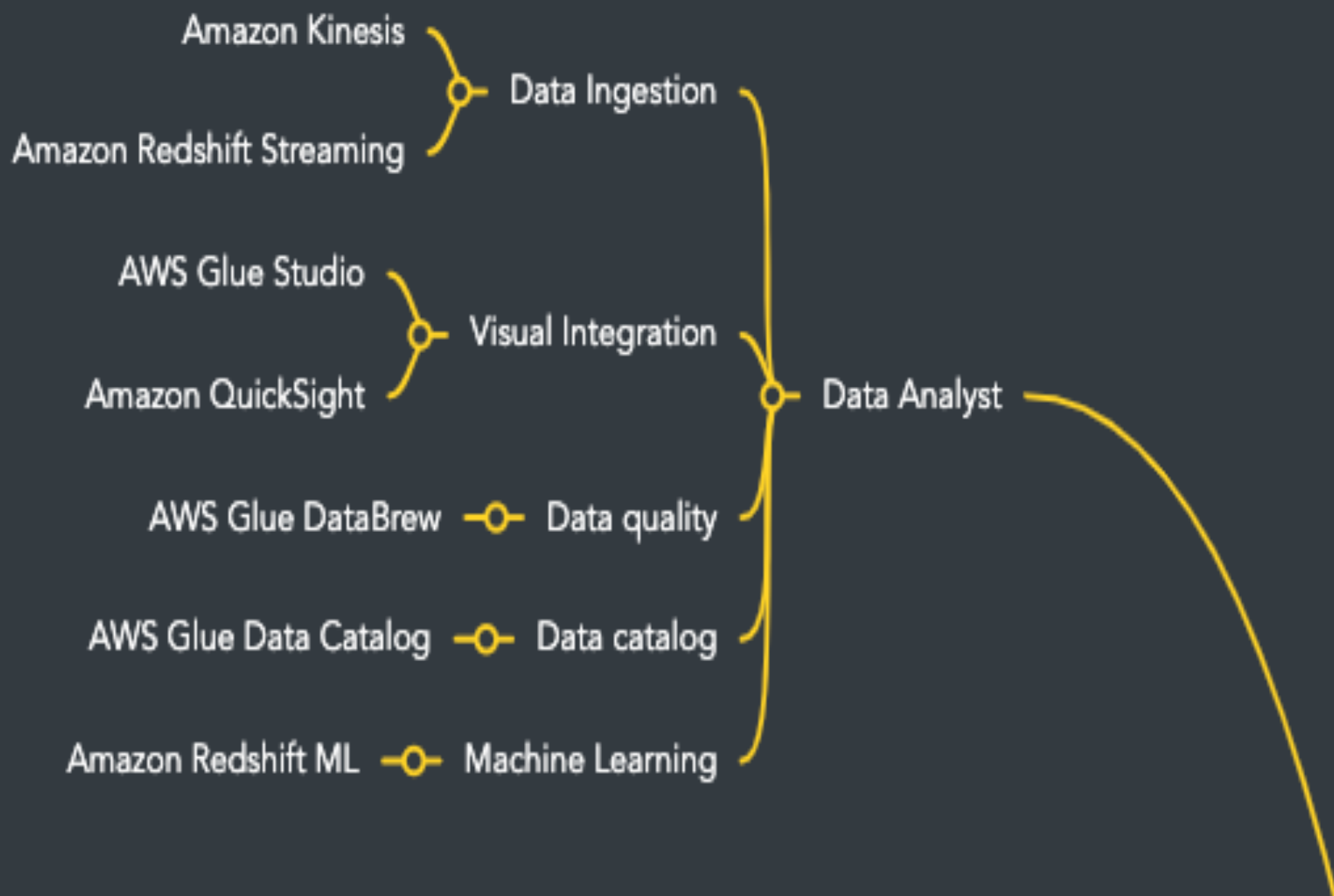
Democratizing analytics

**Make analytics
available,
accessible and
affordable**



AWS analytics mind-map





Data Scientist

Machine Learning

Amazon SageMaker

Amazon Redshift ML

Notebooks

Amazon EMR Studio

Collaboration

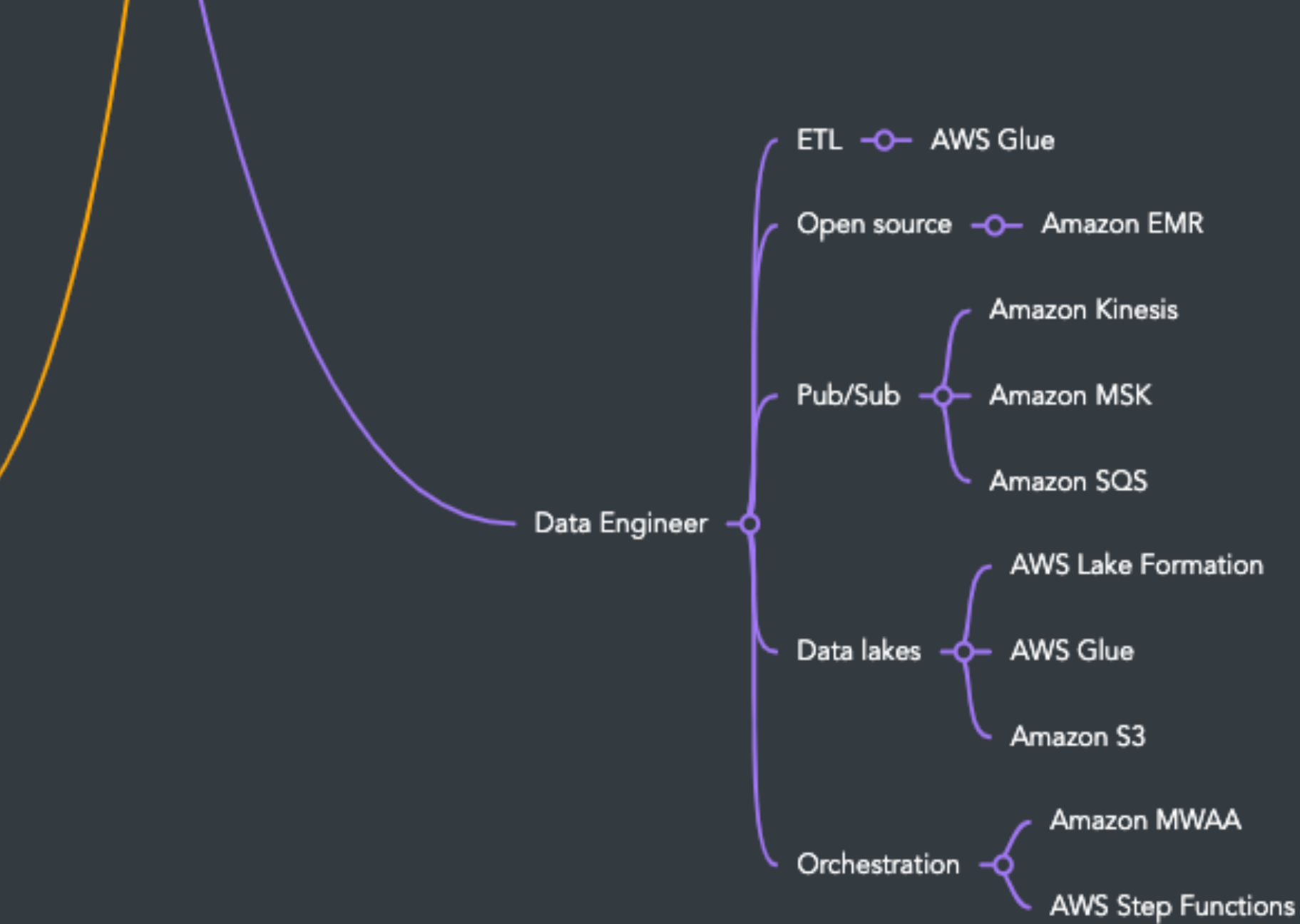
Amazon SageMaker Studio

Data Quality

AWS Glue DataBrew

Data Catalog

AWS Glue Data Catalog





Democratizing analytics

**Make analytics
available,
accessible and
affordable**



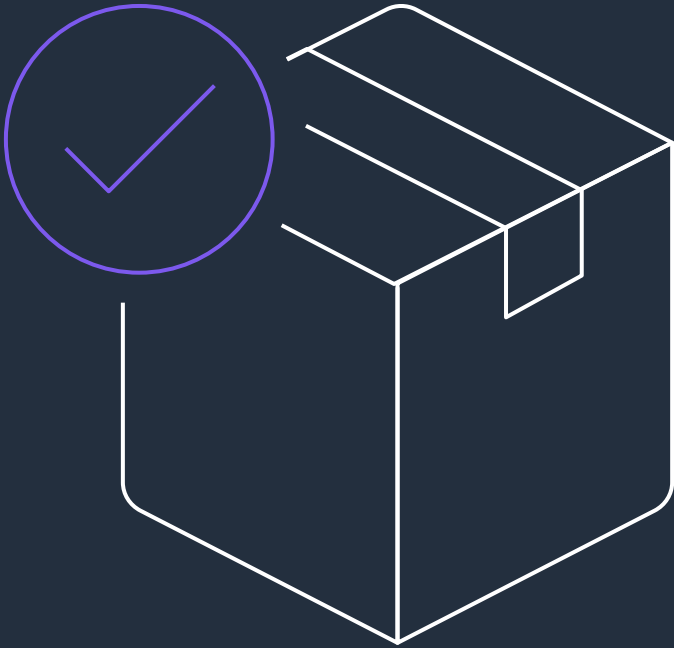
AWS differentiators



Ease of use



Ease of use



Low barrier to entry



Reduced operational burden



Low code / No code experience

Ease of use by AWS

Low barrier
to entry



Intuitive

Start quick / Fail fast

Open to a wider audience

Ease of use by AWS

Low barrier
to entry



Intuitive

Start quick / Fail fast

Open to a wider audience

Reduced
operational
burden



Automation

Monitoring

Operations

Ease of use by AWS

Low barrier
to entry



Intuitive
Start quick / Fail fast
Open to a wider audience

Reduced
operational
burden



Automation
Monitoring
Operations

Low code / No
code experience



Increased business agility
Rapid development / higher
productivity
Reduced OpEx

Price performance



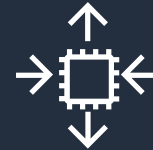
Price performance



Performance pricing



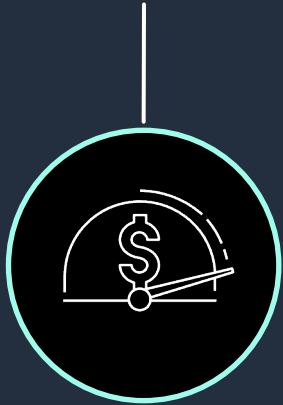
Do more with less



Best fit

Price performance by AWS

Performance
pricing



Consumption based
pricing models

Continuous
performance
improvements



Price performance by AWS

Performance
pricing



Consumption based
pricing models

Continuous
performance
improvements

Do more
with less

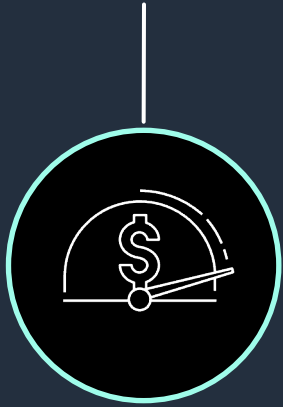


Iterative feature
development

3P and native
integration support

Price performance by AWS

Performance
pricing



Consumption based
pricing models

Continuous
performance
improvements

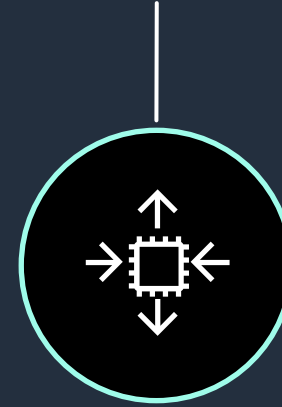
Do more
with less



Iterative feature
development

3P and native
integration support

Best fit



Deployment choices

Amazon EMR

BIG DATA ANALYTICS USING OPEN-SOURCE FRAMEWORKS: APACHE SPARK, PRESTO, TRINO, HIVE, HBASE, HUDI AND 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 EC2 Spot, EMR Managed Scaling and per-second billing



Self service data science

Data Science IDE with EMR Studio and Deep integration with Sagemaker Studio provides ability to use open source UX and frameworks to build, visualize and debug applications



Run workloads on EC2, EKS or on-premises

EMR provides flexibility to run big data workloads on EC2, EKS, and on-premises with Outpost



S3 Data Lake Integration

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



Amazon EMR

3.9x

Faster than
standard Apache
Spark 3.0 in
TPC-DS 3 TB
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



Amazon EMR deployment options

Feature	
Multi-AZ Availability	
OSS frameworks	
Ability to choose OSS version	
Automatic resource scaling	
Ability to choose instance type	
Ability to use EC2 Spot	
Pricing	
Ability to allocate costs	



Amazon EMR deployment options

Feature	Amazon EMR on EC2	
Multi-AZ Availability	No (clusters run in a single AZ)	
OSS frameworks	Spark, Hive, Presto, Trino, Flink	
Ability to choose OSS version	Yes	
Automatic resource scaling	Yes	
Ability to choose instance type	Yes	
Ability to use EC2 Spot	Yes	
Pricing	By instance type used	
Ability to allocate costs	Per cluster	



Amazon EMR deployment options

Feature	Amazon EMR on EC2	Amazon EMR on EKS
Multi-AZ Availability	No (clusters run in a single AZ)	Yes (with multi-AZ EKS clusters)
OSS frameworks	Spark, Hive, Presto, Trino, Flink	Spark
Ability to choose OSS version	Yes	Yes
Automatic resource scaling	Yes	Yes
Ability to choose instance type	Yes	Optional (use EC2 instances or AWS Fargate)
Ability to use EC2 Spot	Yes	Yes
Pricing	By instance type used	By vCPU and memory used
Ability to allocate costs	Per cluster	Per application



Amazon EMR deployment options

Feature	Amazon EMR on EC2	Amazon EMR on EKS	Amazon EMR Serverless
Multi-AZ Availability	No (clusters run in a single AZ)	Yes (with multi-AZ EKS clusters)	Yes (automated job redirection)
OSS frameworks	Spark, Hive, Presto, Trino, Flink	Spark	Spark, Hive
Ability to choose OSS version	Yes	Yes	Yes
Automatic resource scaling	Yes	Yes	Yes
Ability to choose instance type	Yes	Optional (use EC2 instances or AWS Fargate)	No
Ability to use EC2 Spot	Yes	Yes	No
Pricing	By instance type used	By vCPU and memory used	By vCPU and memory used
Ability to allocate costs	Per cluster	Per application	Per application or per job



Amazon Athena



SERVERLESS

ZERO setup cost

Serverless: zero infrastructure, zero administration



PAY PER QUERY

Pay only for queries run

\$5/TB

Save **30%–90%** on per-query costs through compression



OPEN AND FLEXIBLE

ANSI SQL

JDBC/ODBC drivers

Multiple formats, compression types, and complex joins and data types



EASY TO USE

Point to S3 and start querying

DDL operations

Query concurrency

Integrated data connectors



Amazon Athena



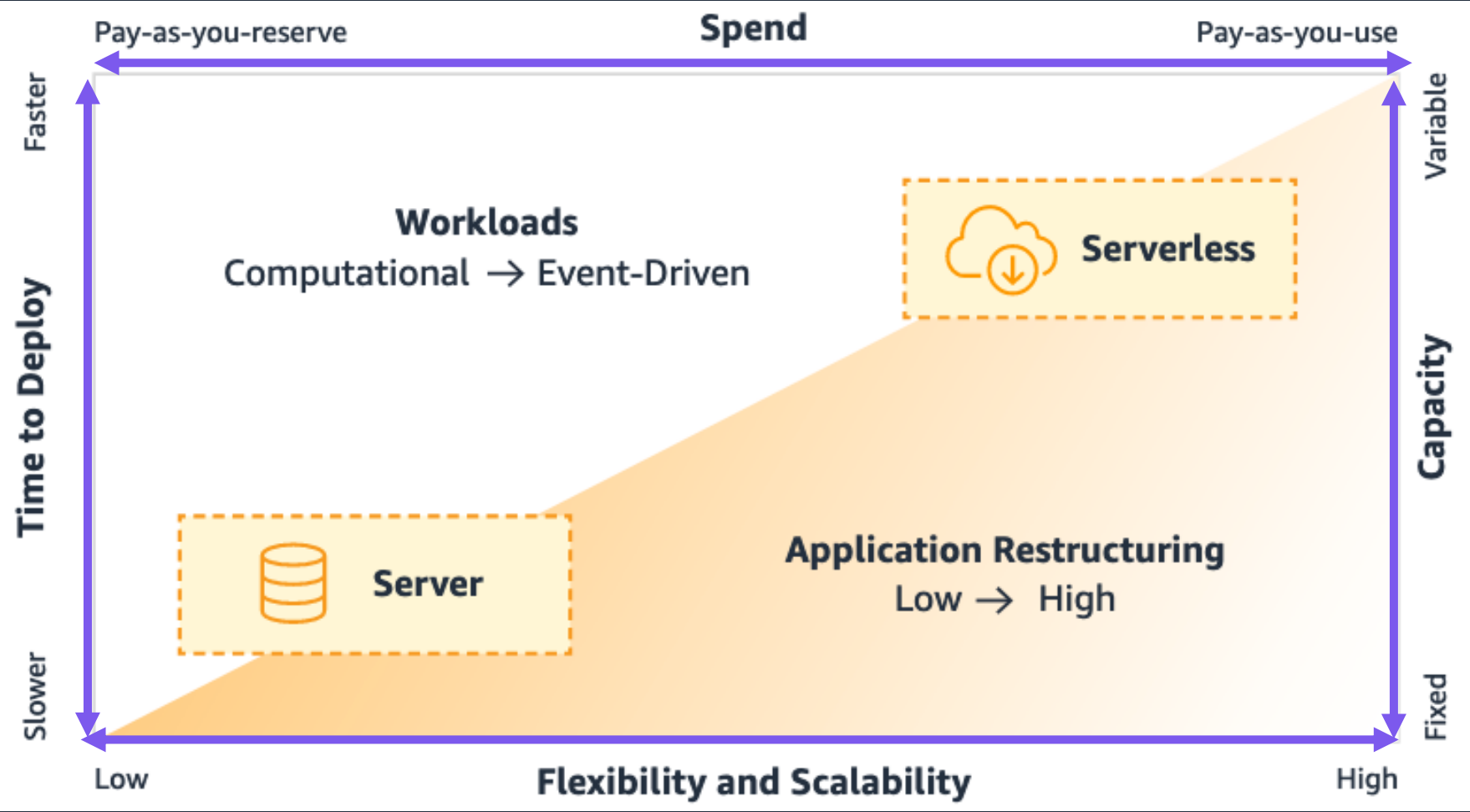
What is a Data Lake? (aka layered & flexible storage)

A data lake is a **centralized repository** that allows you to store all your **structured and unstructured** data at **any scale**

You can store your data as-is, without having to first structure the data, and then easily run **different types of analytics or transformations**



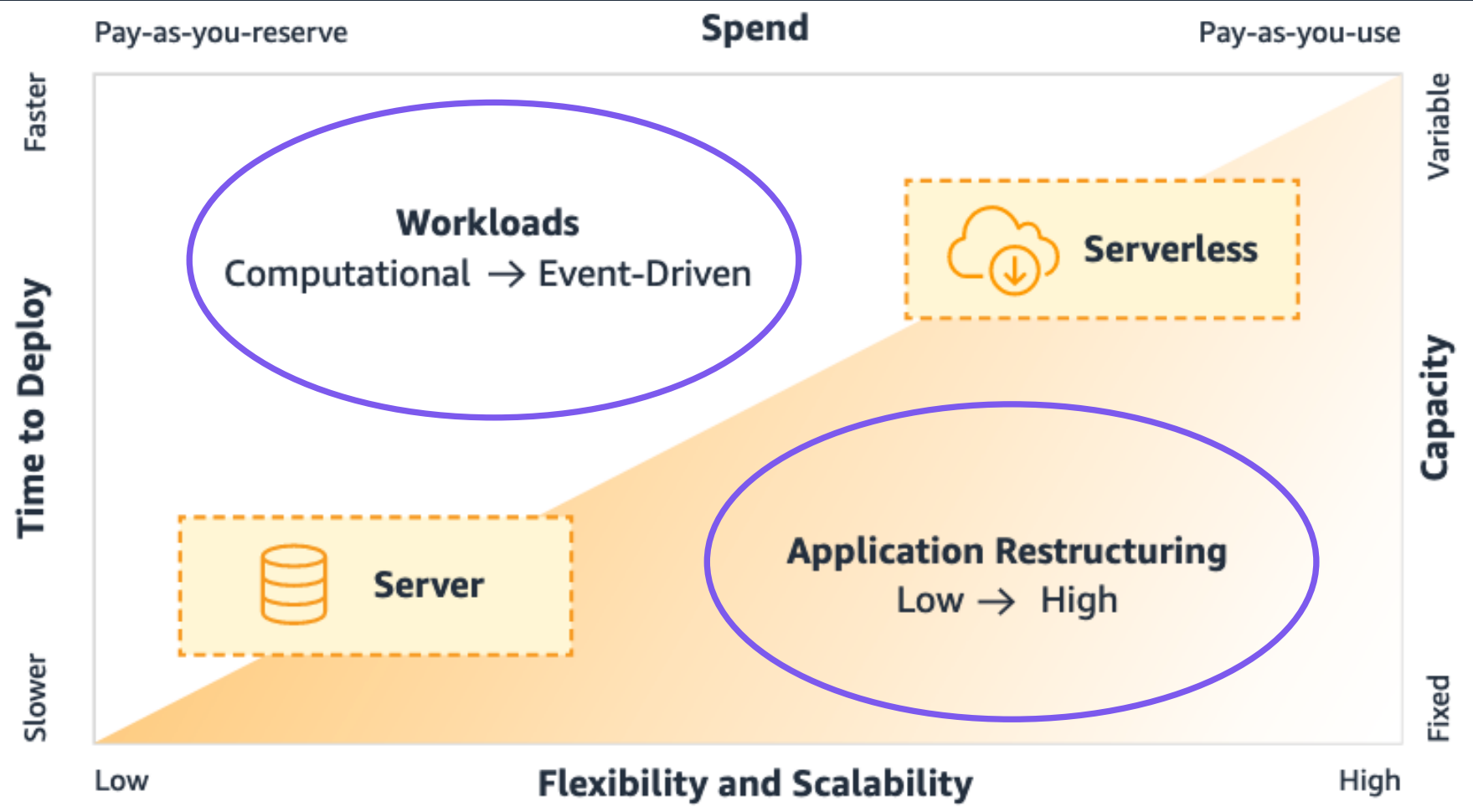
Serverless TCO



Source: [Deloitte](#)



Serverless TCO



Source: [Deloitte](#)



Serverless data analytics on AWS

AWS has the **most serverless options**
for data analytics in the cloud

INTERACTIVE
QUERY



AMAZON
ATHENA

BIG DATA
PROCESSING



AMAZON
EMR

REAL-TIME
ANALYTICS



AMAZON
MSK

REAL-TIME
ANALYTICS



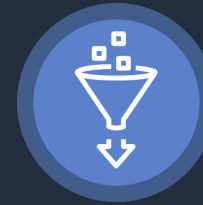
AMAZON
KINESIS

DATA
WAREHOUSING



AMAZON
REDSHIFT

DATA
INTEGRATION



AWS
GLUE

DATA
VISUALIZATION



AMAZON
QUICKSIGHT

DATA LAKE SETUP
MANAGEMENT AND
GOVERNANCE



AWS LAKE
FORMATION



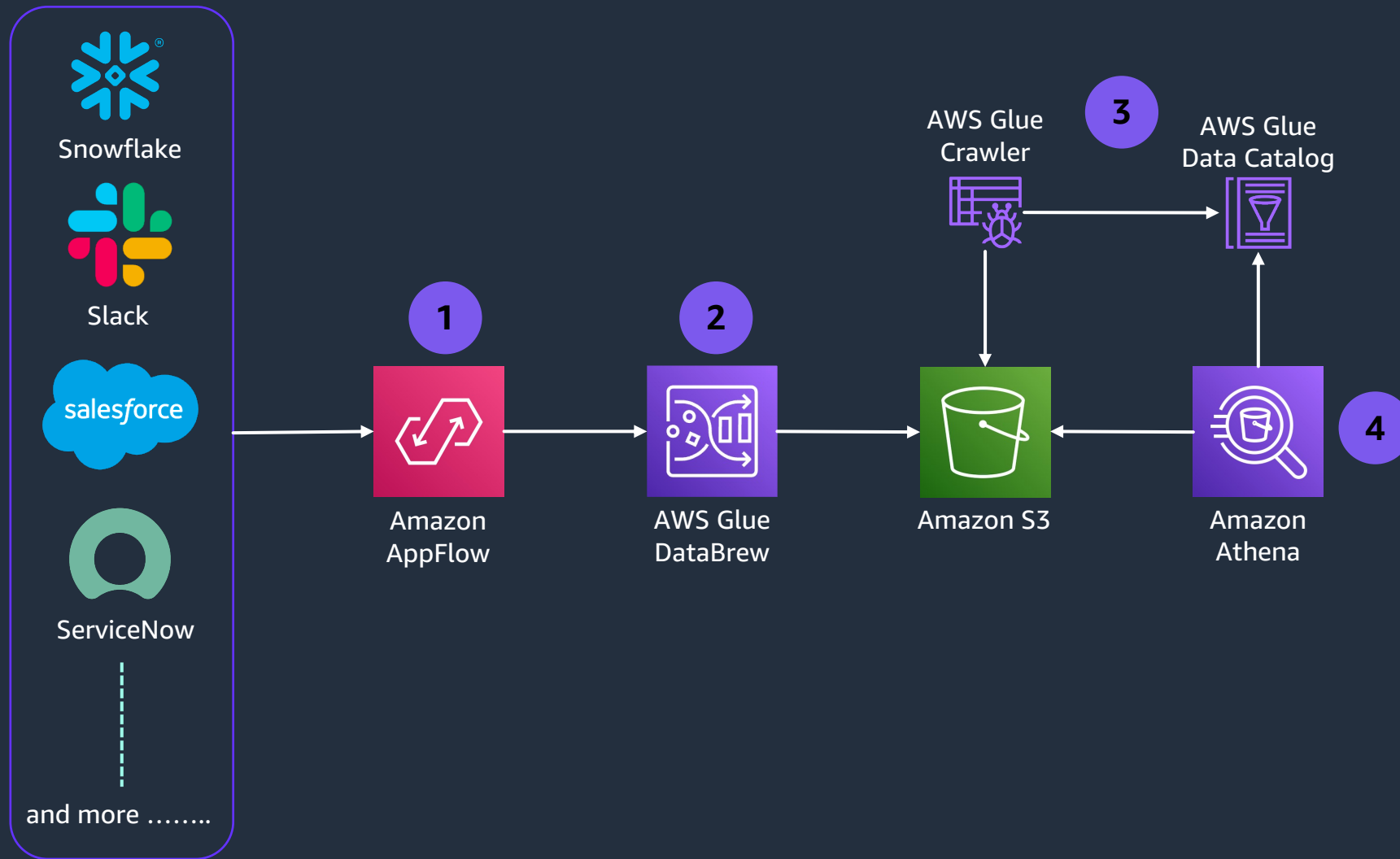
AWS differentiators



Demo



Demo architecture



Useful links
[YouTube Demo](#)
[AWS Blog Post](#)

Request for Survey



Track: Data and Analysis Track

Topic: Democratizing your organization's data analytics experience