# Best practices for implementing a data lake in Amazon S3

Kumar Nachiketa

Senior Partner Solutions Architect

Amazon Web Services

# Agenda

Data at scale and data lakes

Data lake foundation best practices

Data lake performance & security best practices

AWS Lake Formation demo

# Data at scale and data lakes

# 60 seconds in Viber

**3M** messages sent

**140,000** calls

**1.2M** users log in

**2,000** new users join

**300,000** stickers sent

**10,000** group chat likes

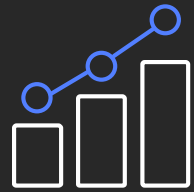**1.5M** photos sent

# **A day** in FINRA

**75 billion** market events
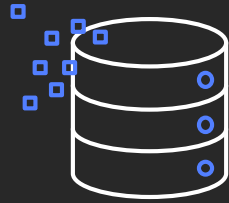
**50,000** files daily

**200** rules to format files

**Half a trillion** validations each day

# Data at scale

**Growing exponentially**
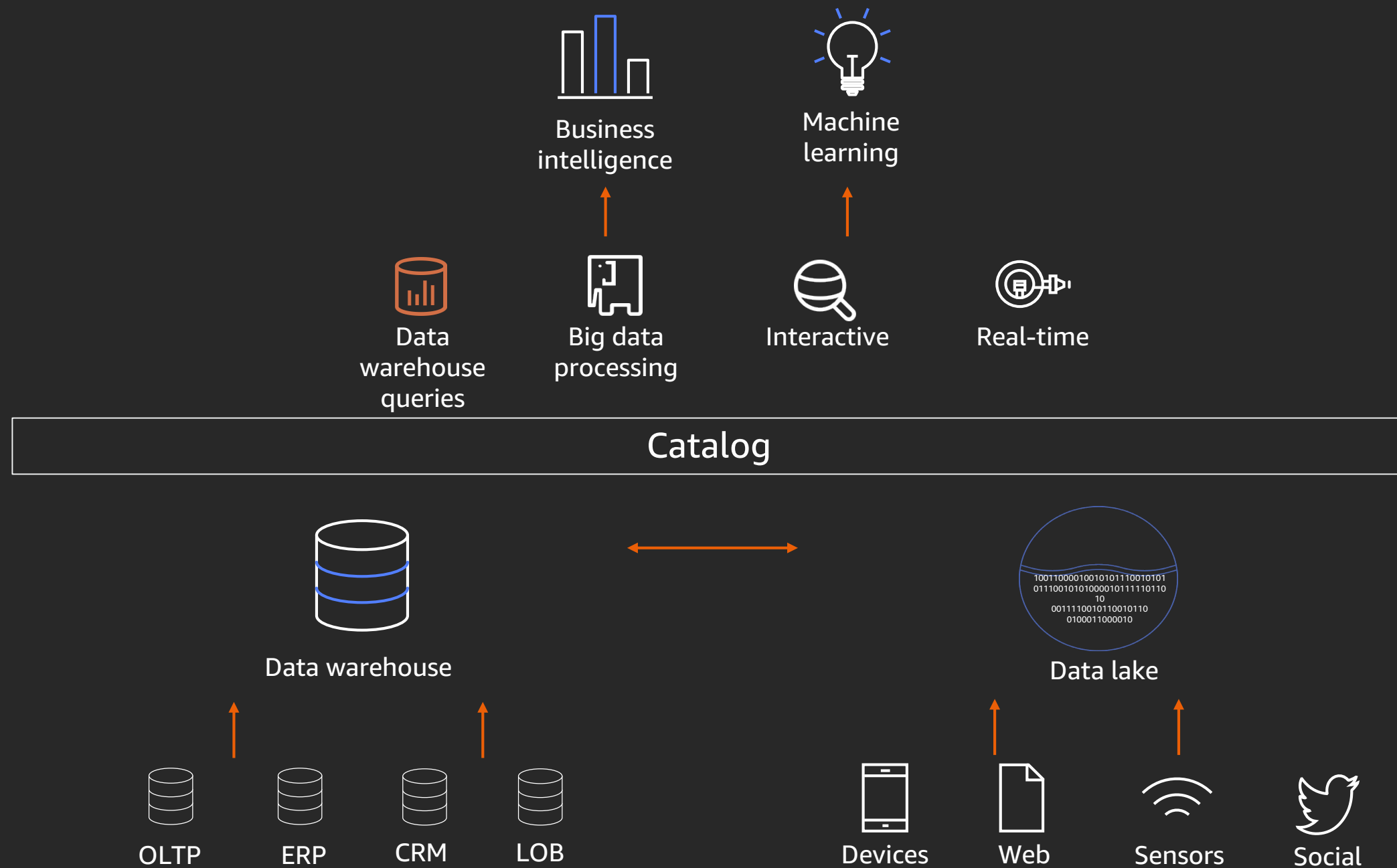
**From new sources**

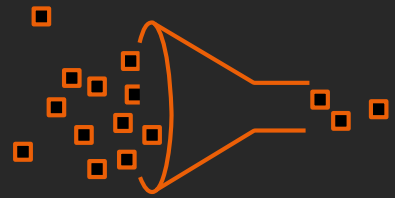**Increasingly diverse**

**Used by many people**

**Analyzed by many applications**

# Defining the data lake

# Amazon S3 is the foundation of any data lake

Multiple data
input sources

Storage scales
on demand

Supports many
unique users
and teams

Analyzed by
many applications

# Amazon S3 as the foundation for data lakes

Amazon Athena

Amazon Elasticsearch Service

Amazon EMR

Amazon Kinesis

Amazon Redshift

AWS Lake Formation & AWS Glue

Amazon Comprehend

Amazon SageMaker

**Amazon S3**

Amazon Rekognition

AWS Snowball

AWS Snowmobile

Amazon Kinesis Data Firehose

Amazon Kinesis Data Streams

Durable, available, exabyte-scalable

Secure, compliant, auditable

High performance

Low-cost storage and analytics

Broad network integration

# Data lake foundation best practices

# Data lake on AWS – S3 at the core

Central storage:
Scalable, secure,
cost-effective

**Amazon
S3**

# Data lake on AWS – access and user interface

AWS
AppSync

Amazon
API Gateway

Amazon
Cognito

Access & user interfaces

Central storage:
Scalable, secure,
cost-effective

**Amazon
S3**

# Data lake on AWS – manage and secure

AWS
AppSync

Amazon
API Gateway

Amazon
Cognito

**Access & user interfaces**

Central storage:
Scalable, secure,
cost-effective

**Amazon
S3**

**Manage & secure**

AWS Key
Management
Service

AWS Identity and
Access Management

AWS CloudTrail

Amazon
CloudWatch

# Data lake on AWS – data ingestion

AWS AppSync

Amazon API Gateway

Amazon Cognito

**Access & user interfaces**

Central storage:
Scalable, secure,
cost-effective

**Amazon S3**

AWS Snowball

Amazon Kinesis Data Firehose

AWS Direct Connect

AWS Database Migration Service

AWS Storage Gateway

Data ingestion

**Manage & secure**

AWS Key Management Service

AWS Identity and Access Management

AWS CloudTrail

Amazon CloudWatch

Data lake on AWS – catalog and search

# Data lake on AWS – analytics, ML, and serving

**Amazon DynamoDB**  
**Amazon Elasticsearch Service**  
**AWS Glue**  
**AWS Lake Formation**

## Catalog & search

**AWS AppSync**  
**Amazon API Gateway**  
**Amazon Cognito**

## Access & user interfaces

### Central storage:
Scalable, secure, cost-effective

**Amazon S3**

## Data ingestion

**AWS Snowball**  
**Amazon Kinesis Data Firehose**  
**AWS Direct Connect**  
**AWS Database Migration Service**  
**AWS Storage Gateway**

## Manage & secure

**AWS Key Management Service**  
**AWS Identity and Access Management**  
**AWS CloudTrail**  
**Amazon CloudWatch**

**Amazon Athena**  
**Amazon EMR**  
**AWS Glue**  
**Amazon Redshift**

**Amazon QuickSight**  
**Amazon Kinesis**  
**Amazon Elasticsearch Service**  
**Amazon DynamoDB**

**Amazon Rekognition**  
**Amazon SageMaker**  
**Amazon RDS**  
**Amazon Neptune**

## Analytics, machine learning & serving

# Data lake design pattern example

aws SUMMIT ONLINE

# Data lake workflow pattern

## Typical steps in building a data lake



1. Set up storage
2. Ingest data
3. Cleanse, prep, and catalog data
4. Configure and enforce security and compliance policies
5. Make data available for analytics

# Data lake ingest and transform patterns

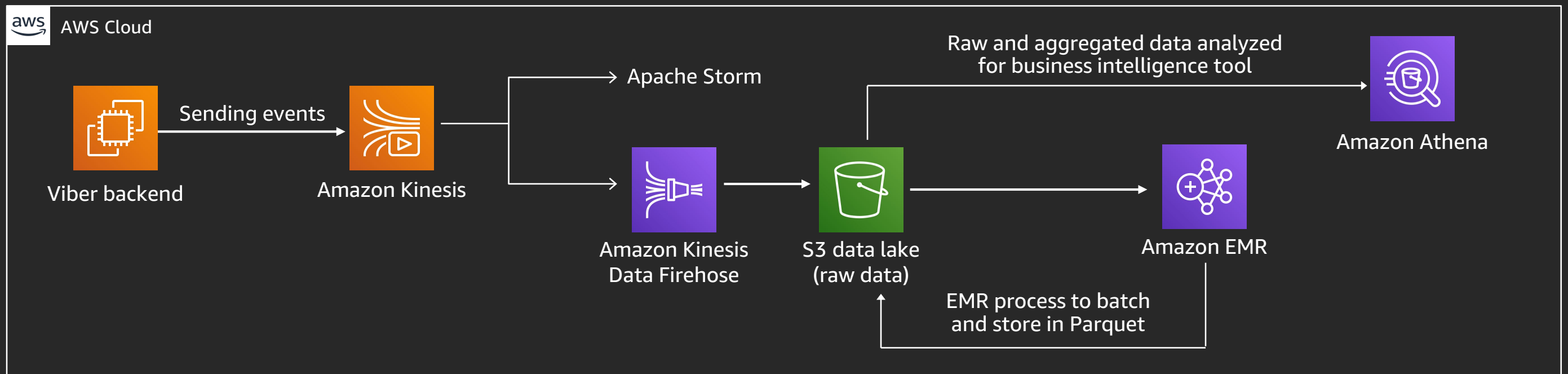Pipelined architectures improve governance, data management, and efficiency

**Raw data**
Amazon S3
Standard

**ETL**
AWS Glue or Amazon EMR

**Production data
(data lake)**
Amazon S3
Intelligent-Tiering

**Data warehouse**
Amazon Redshift

**Triggered code**
AWS Lambda

**ETL and catalog management**
AWS Glue and AWS Lake Formation

**Triggered code**
AWS Lambda

# Viber: Processing events on a data lake

**Workload:** Communications platform serving a billion users worldwide
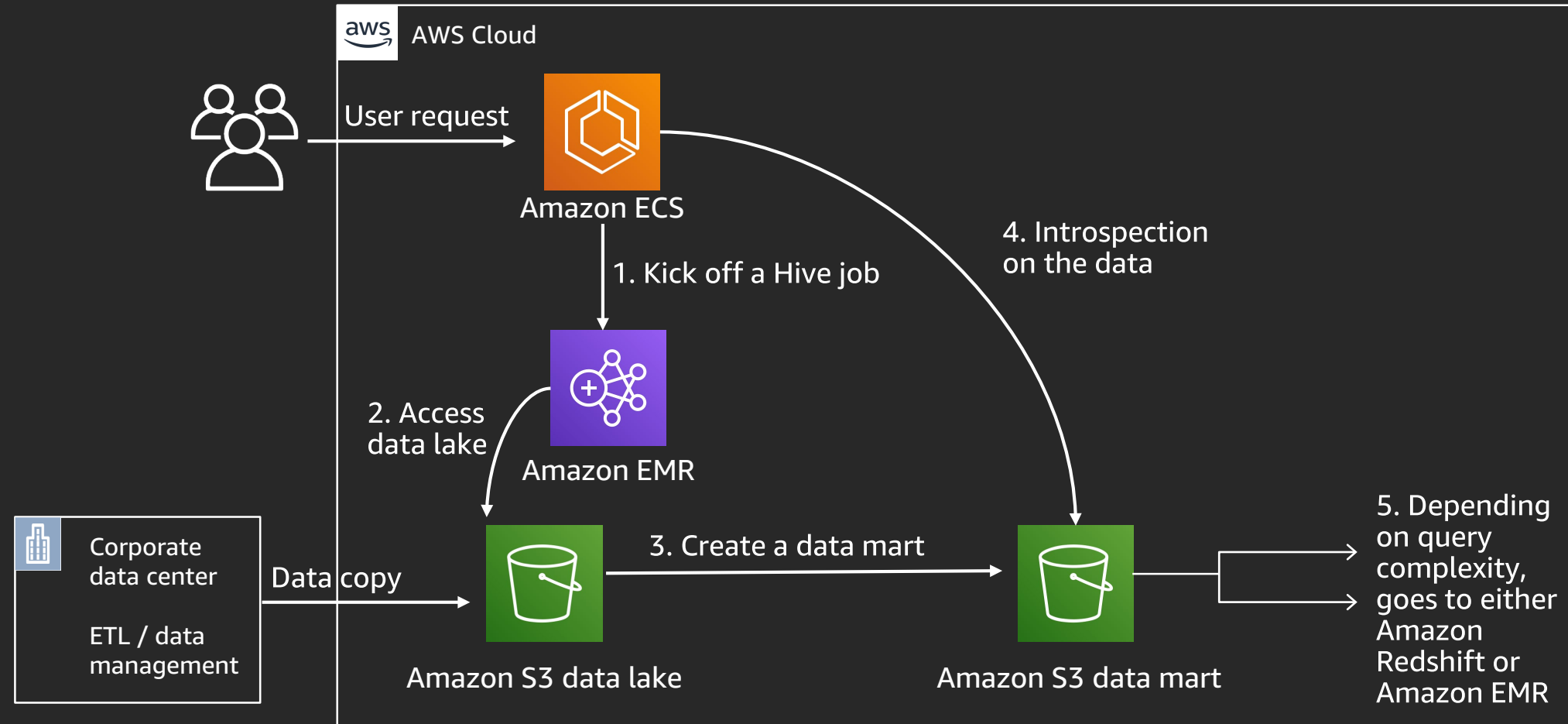
- Processing over 10–15 billion events per day
- Peaking at 300,000 events per second
- Storing many petabytes of data
- Running over 200 events on a single Amazon Kinesis stream

# Viber: Processing events on a data lake

**Workload:** Communications platform serving a billion users worldwide

- Processing over 10–15 billion events per day
- Peaking at 300,000 events per second
- Storing many petabytes of data
- Running over 200 events on a single Amazon Kinesis stream

# FINRA: Petabyte-scale data analysis

**Workload:** Financial regulatory authority providing users with ability to access PBs of data for analytics
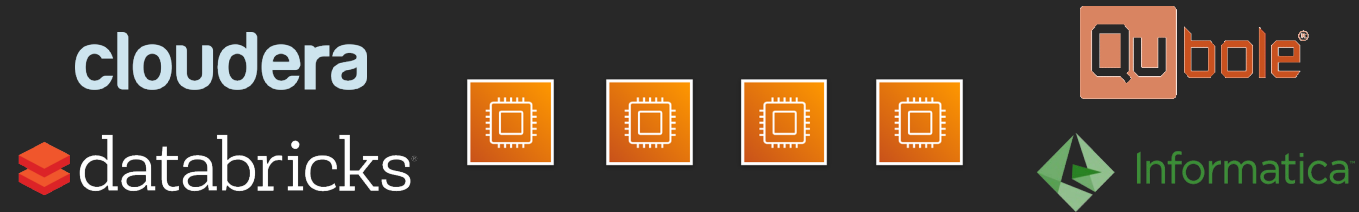
## Architecting on AWS for scale



*"With the new architecture, we uncovered more needs of the end user. This led us to move from one large Amazon Redshift cluster to a blend of querying engines."*

# Running analytics on AWS data lakes

## Lift-and-shift



## AWS Managed Services



Amazon Redshift   AWS Glue   Amazon EMR   Amazon Athena

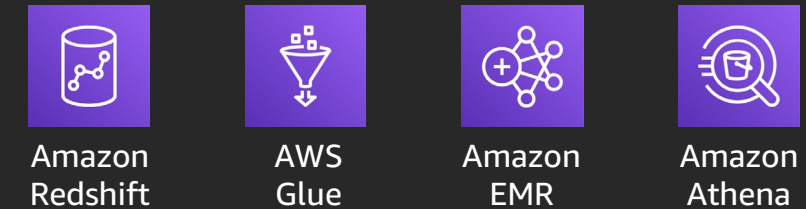| | Lift-and-shift | AWS Managed Services |
|---|---|---|
| **What** | • Run third-party analytics tools on Amazon EC2<br>• Use Amazon EBS and Amazon S3 as data stores<br>• Self-managed environments | • AWS managed and serverless platforms<br>• AWS Glue, Amazon Athena, Amazon EMR, Amazon Redshift<br>• More options to process data in place |
| **Why** | • Simplify on-premises migrations<br>• Use existing tools, code, and customizations<br>• Minimize application changes | • Focus on data outcomes, not infrastructure<br>• Speed adoption of new capabilities<br>• More tightly integrated with AWS security |
| **Consider** | • You provision, manage, and scale<br>• You monitor and manage availability<br>• You own upgrades and versioning | • Utilizing AWS Lake Formation<br>• Flexibility and choice with open data formats<br>• Leverage AWS pace of innovation |

## Amazon S3 is the storage foundation for both approaches
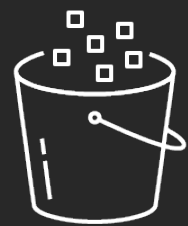
# Amazon S3: Data management, performance, and security

aws SUMMIT ONLINE

# Choosing the right data lake storage class

## Select storage class by data pipeline stage

| Raw data | ETL | Production data lake | Online cool data | Historical data |
|---|---|---|---|---|
| Amazon S3 Standard | Amazon S3 Standard | Amazon S3 Intelligent-Tiering | Amazon S3 Standard Infrequent Access (S3-IA/ZIA) | Amazon S3 Glacier or S3 Glacier Deep Archive |

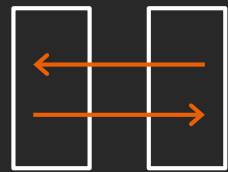| | | | | |
|---|---|---|---|---|
| • Small log files | • Data churn | • Optimized sizes (MBs) | • Replicated DR data | • Historical assets |
| • Overwrites if synced | • Small intermediates | • Many users | • Infrequently accessed | • ML model training |
| • Short-lived | • Multiple transforms | • Unpredictable access | • Infrequent queries | • Compliance/audit |
| • Moved & deleted | • Deletes <30 days | • Long-lived assets | • ML model training | • Data protection |
| • Batched & archived | • Output to data lake | • Hot to cool | | • Planned restores |

## Optimize costs for all stages of data lake workflows
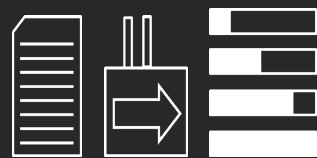
# Data management at scale: Best practices

## Utilize Amazon S3 object tagging
Granularly control access, analyze usage,
manage life cycle policies, and replicate objects

## Implement life cycle policies
Automated, policy-driven archive and data expiration

## Utilize batch operations
Manage millions to billions of objects with a single request

**Plan for rapid growth and automation of management at any scale**

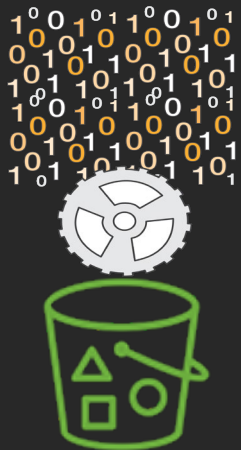# Recommendation: Consider performance design patterns

## Structure key namespace to scale

Most workloads fit in the S3 3,500 PUT / 5,500 GET TPS per key name partition

Amazon S3 automatically creates partitions as data lake use increases

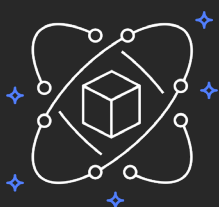Extremely bursty workloads might require customized key name design

## Consider object format and size

Use Parquet or optimized columnar format

Aim for 16–256 MB minimum object size (might require aggregation during ingest)

Perform parallel byte range access (included in AWS SDKs)
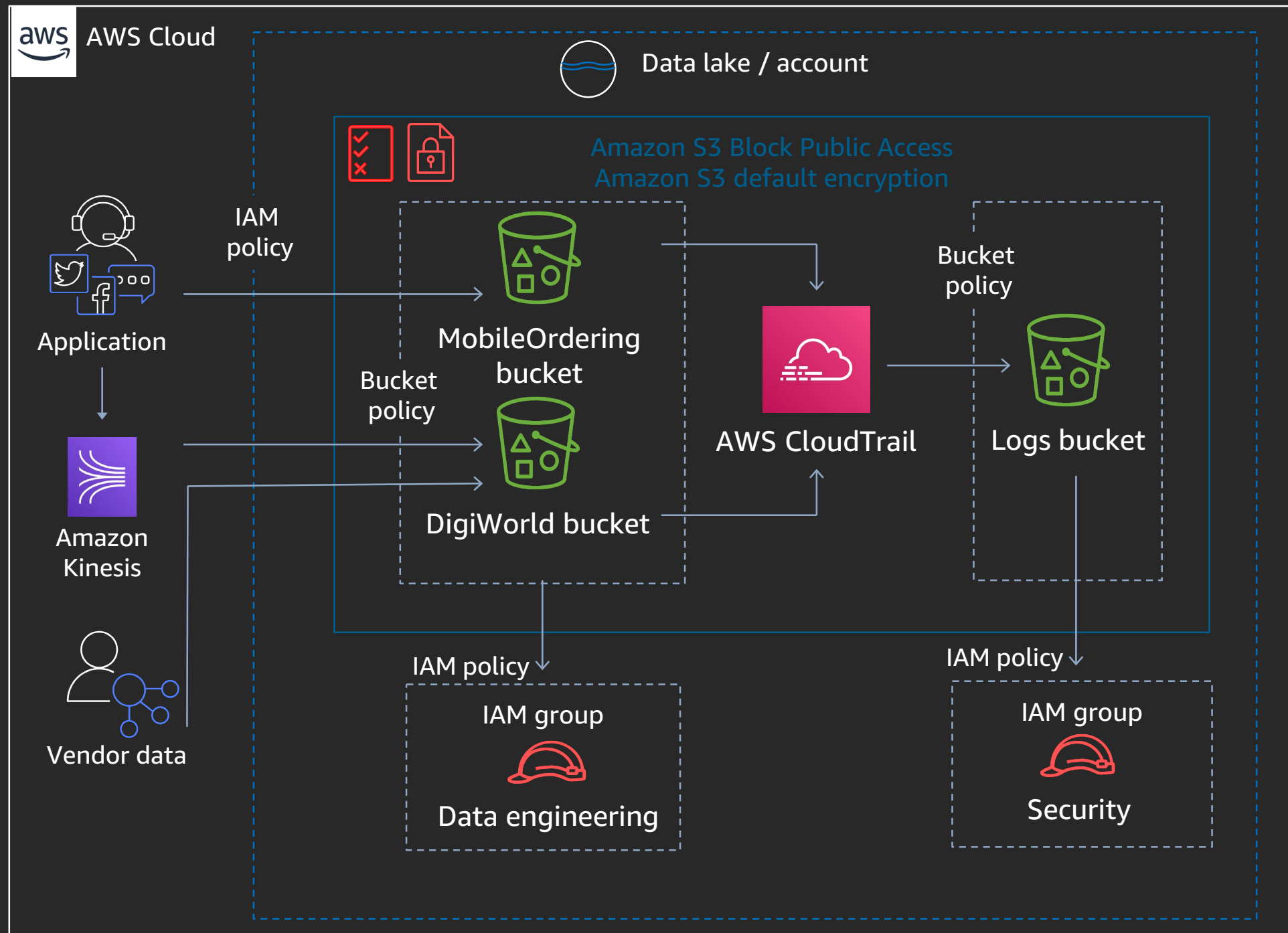
## Use latest SDKs and software versions

AWS SDKs include support for the latest features and optimizations

Amazon EMR 5.18 and above supports S3 Select for Hive, Presto, and Spark

# Secure your data lake



AWS Cloud

Data lake / account

Amazon S3 Block Public Access
Amazon S3 default encryption

Application

IAM policy

MobileOrdering bucket

Bucket policy

DigiWorld bucket

AWS CloudTrail

Bucket policy

Logs bucket

Amazon Kinesis

Vendor data

IAM policy

IAM group

Data engineering

IAM policy

IAM group

Security

Deny access by default
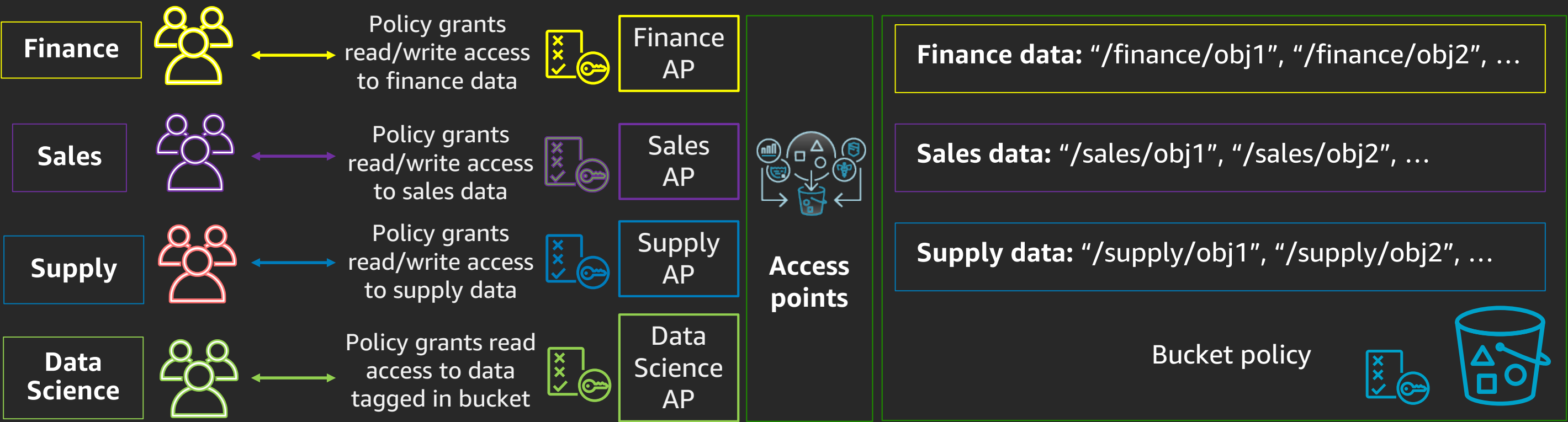
Encrypt your data

Secure multiple data input sources

Provide specific access where appropriate

Support multiple unique users and teams

# Recommendation: Use Amazon S3 access points to manage your application set

New!

Granular control for hundreds of teams accessing your data lake

| Finance | | Policy grants read/write access to finance data | | Finance AP | | Finance data: "/finance/obj1", "/finance/obj2", … |
| Sales | | Policy grants read/write access to sales data | | Sales AP | | Sales data: "/sales/obj1", "/sales/obj2", … |
| Supply | | Policy grants read/write access to supply data | | Supply AP | | Supply data: "/supply/obj1", "/supply/obj2", … |
| Data Science | | Policy grants read access to data tagged in bucket | | Data Science AP | | |

Access points

Bucket policy

# AWS Lake Formation

# AWS Lake Formation

Build a secure data lake in days

## Build data lakes quickly

Move, store, catalog, and clean your data faster

Transform to open formats like Parquet and ORC

ML-based deduplication and record matching

## Simplify security management

Centrally define security, governance, and auditing policies

Enforce policies consistently across multiple services

Integrates with IAM and KMS

## Provide self-service access to data

Build a data catalog that describes your data

Enable analysts and data scientists to easily find relevant data

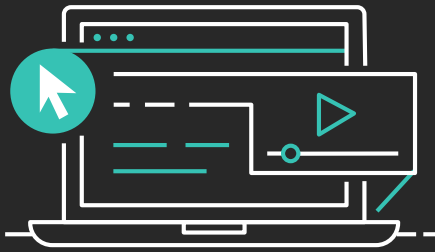Analyze with multiple analytics services without moving data

# Demo

aws SUMMIT

# Overarching takeaways

- Amazon S3 is the foundation for data lakes

- Leverage pipelined architectures to improve governance, data management, and efficiency

- Improve performance by parallelizing access and scaling horizontally

- Privatize your data lake, encrypt everything, and secure specific access to and from that data lake

- Simplify control for shared bucket access by many teams by using S3 access points

# Learn storage with AWS Training and Certification

Resources created by the experts at AWS to help you build cloud storage skills

45+ free digital courses cover topics related to cloud storage, including:

- Amazon S3
- AWS Storage Gateway
- Amazon S3 Glacier

- Amazon Elastic File System (Amazon EFS)
- Amazon Elastic Block Store (Amazon EBS)

Classroom offerings, such as Architecting on AWS, feature AWS expert instructors and hands-on activities

Visit the storage learning path at https://aws.training/storage

# Thank you!

Kumar Nachiketa