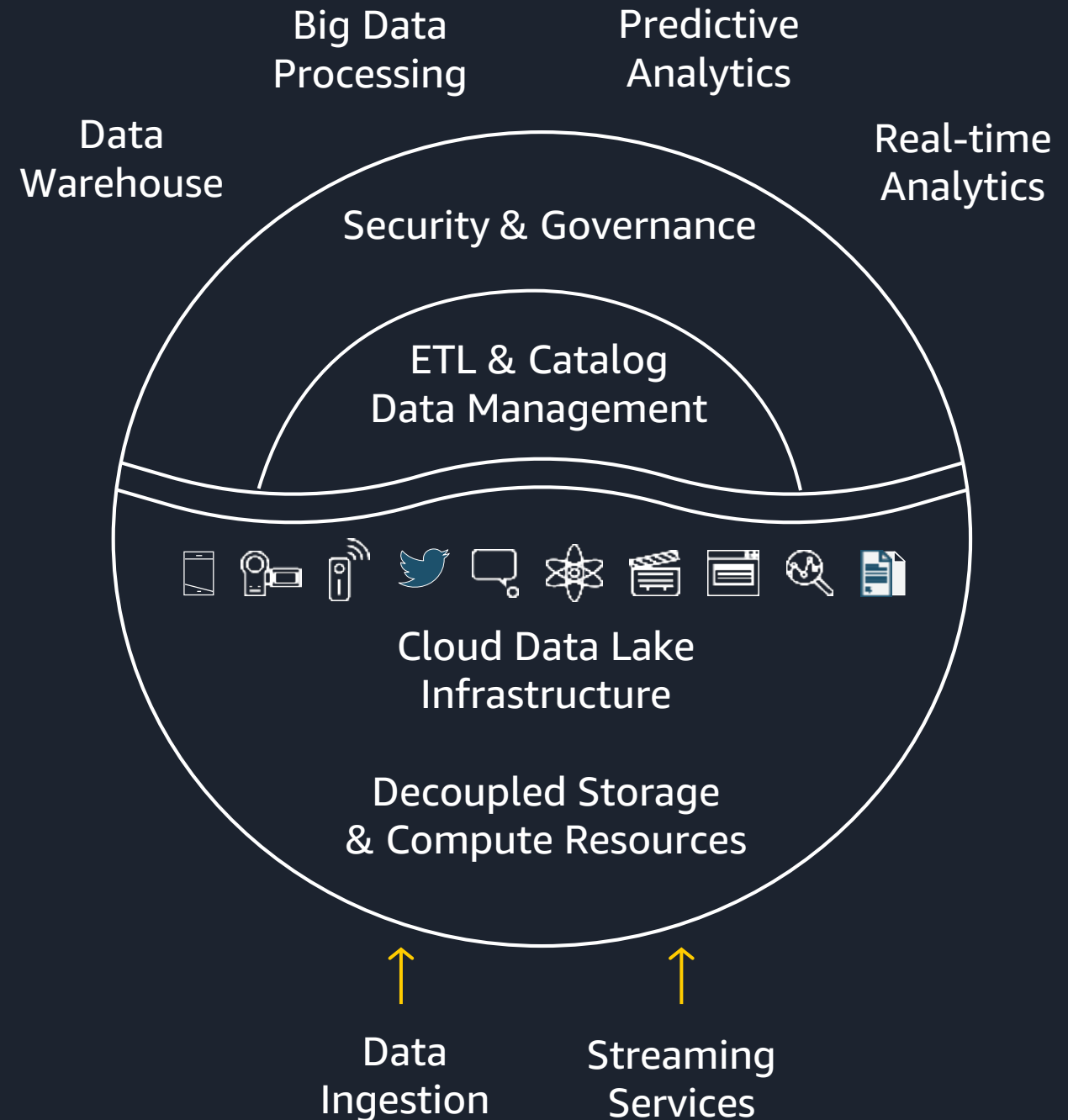**AWS Data, Databases, and Analytics Online Series**

# How to ingest data seamlessly to build your data lake

**Wali Akbari**

Storage Specialist Solutions Architect, AWS

# What is a data lake?

A **data lake** is a centralized

repository that allows you

to ingest, store, and manage

structured and unstructured data

at unlimited scale.

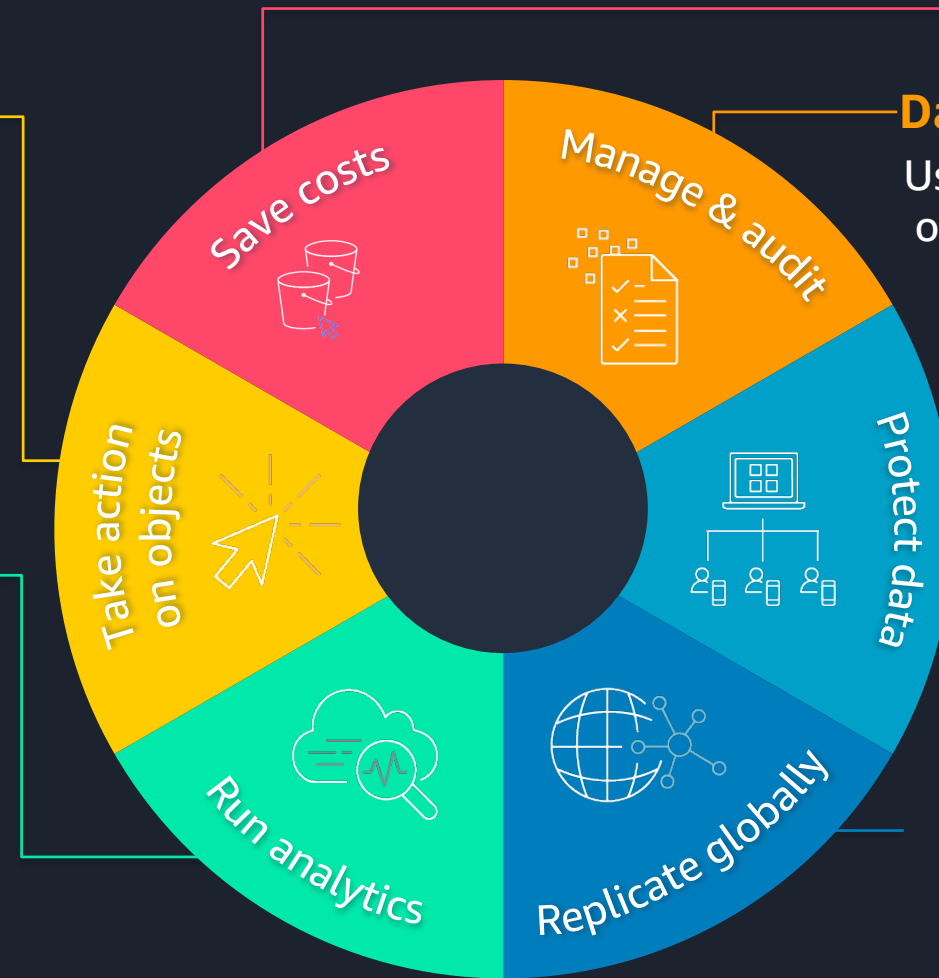Then gain insights through analytics

and machine learning.

Big Data
Processing

Predictive
Analytics

Data
Warehouse

Real-time
Analytics

Security & Governance

ETL & Catalog
Data Management

Cloud Data Lake
Infrastructure

Decoupled Storage
& Compute Resources

Data
Ingestion

Streaming
Services

aws

# Amazon Simple Storage Service (Amazon S3) - Overview



**S3 Batch Operations**
Take actions on objects at scale

**Analytics & file system integration**
S3-integrated analytics applications
AWS Lake Formation to stand up a data lake in days
S3 Select to query data in place
FSx for Lustre for HPC, ML, and media data processing

**Data management tools**
Use tags, buckets, and prefixes to organize data.

**Access management**
Configure access to S3 resources.
Block all public access S3 Block Public Access.

**Cross-Region Replication**
Replicate objects to a different region of your choice.

**S3 Storage Classes**
S3 Standard
S3 Standard-IA
S3 Intelligent-Tiering
S3 One Zone-IA
S3 Glacier
S3 Glacier Deep Archive

Wheel segments: Save costs · Manage & audit · Take action on objects · Protect data · Run analytics · Replicate globally

| Security by design | Compliance programs | 11 9's of durability | Multi-AZ resiliency | Limitless scalability |

aws

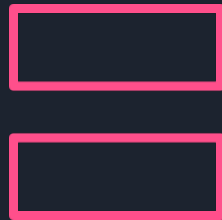# Building a data lake on AWS

**Catalog & Search**
- Amazon DynamoDB
- Amazon Elasticsearch Service
- AWS Glue

**Access & User Interfaces**
- AWS AppSync
- Amazon API Gateway
- Amazon Cognito

**Central Storage**
*Scalable, secure, cost-effective*

Amazon S3

**Data Ingestion**
- ?
- ?
- ?
- ?

**Manage & Secure**
- AWS KMS
- AWS IAM
- AWS CloudTrail
- Amazon CloudWatch

**Analytics & Serving**
- Amazon Athena
- Amazon EMR
- AWS Glue
- Amazon Redshift
- Amazon DynamoDB
- Amazon QuickSight
- Amazon Kinesis
- Amazon Elasticsearch Service
- Amazon Neptune
- Amazon RDS

aws

# Building a data lake on AWS – ingesting data

**Catalog & Search**
- Amazon DynamoDB
- Amazon Elasticsearch Service
- AWS Glue

**Access & User Interfaces**
- AWS AppSync
- Amazon API Gateway
- Amazon Cognito

**Central Storage**
*Scalable, secure, cost-effective*

Amazon S3

**Data Ingestion**
? ? ? ?

**Manage & Secure**
- AWS KMS
- AWS IAM
- AWS CloudTrail
- Amazon CloudWatch

**Analytics & Serving**
- Amazon Athena
- Amazon EMR
- AWS Glue
- Amazon Redshift
- Amazon DynamoDB
- Amazon QuickSight
- Amazon Kinesis
- Amazon Elasticsearch Service
- Amazon Neptune
- Amazon RDS

aws

# Data ingest challenges

What are some challenges customers face when trying to ingest data

Options

Speed

Time and effort

aws

# AWS Storage services portfolio

## Block storage

Amazon EBS

## File storage

Amazon EFS

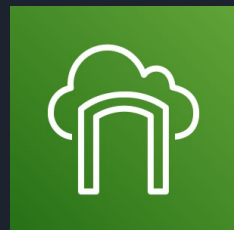Amazon FSx for Windows File Server

Amazon FSx for Lustre
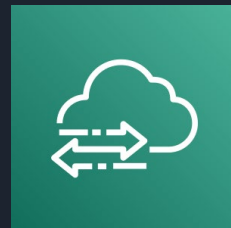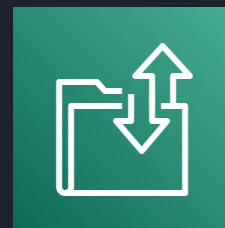
## Object storage

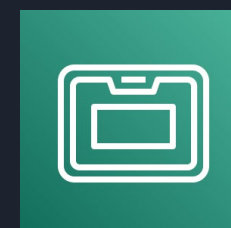Amazon S3

Amazon S3 Glacier

## Hybrid

AWS Storage Gateway
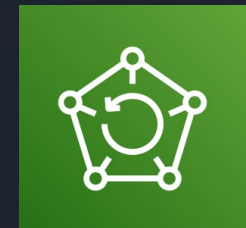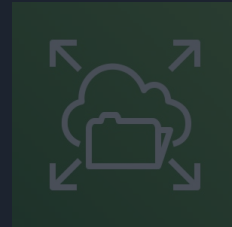
## Transport & edge

AWS DataSync

AWS Transfer Family

AWS Snow* Family

## Backup

AWS Backup

aws

# AWS Storage services portfolio



Block storage

Amazon
EBS

File storage

Amazon
EFS

Amazon FSx for
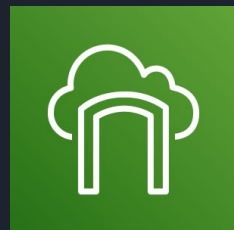Windows File Server

Amazon FSx
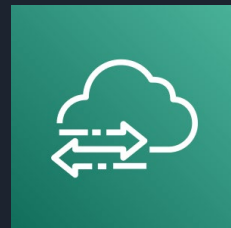for Lustre

Object storage

Amazon
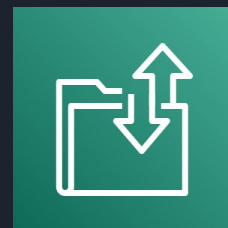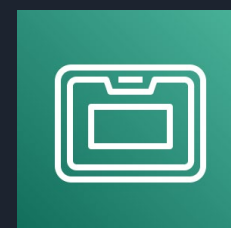S3

Amazon S3
Glacier

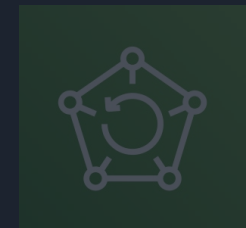Hybrid

AWS Storage
Gateway

Transport & edge
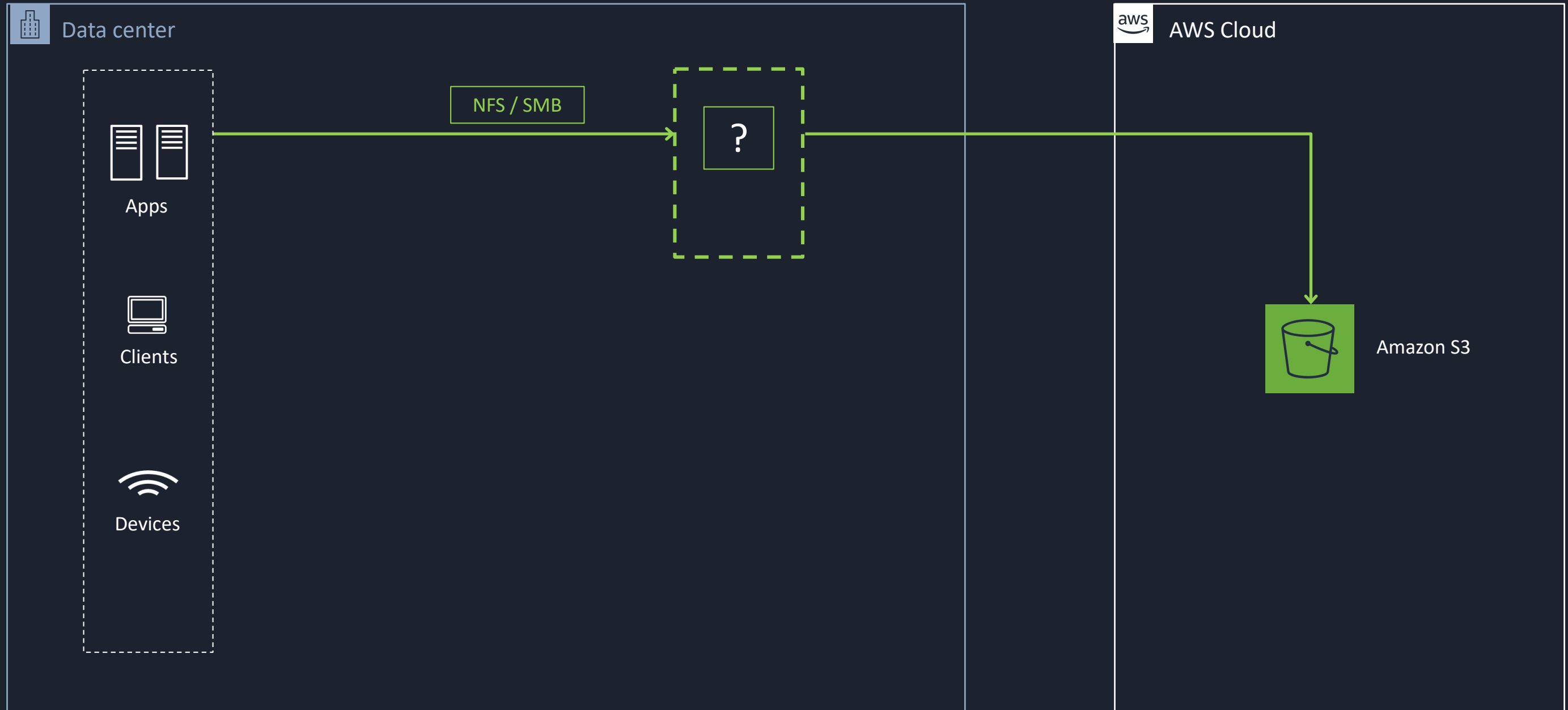
AWS DataSync

AWS Transfer
Family

AWS Snow*
Family

Backup

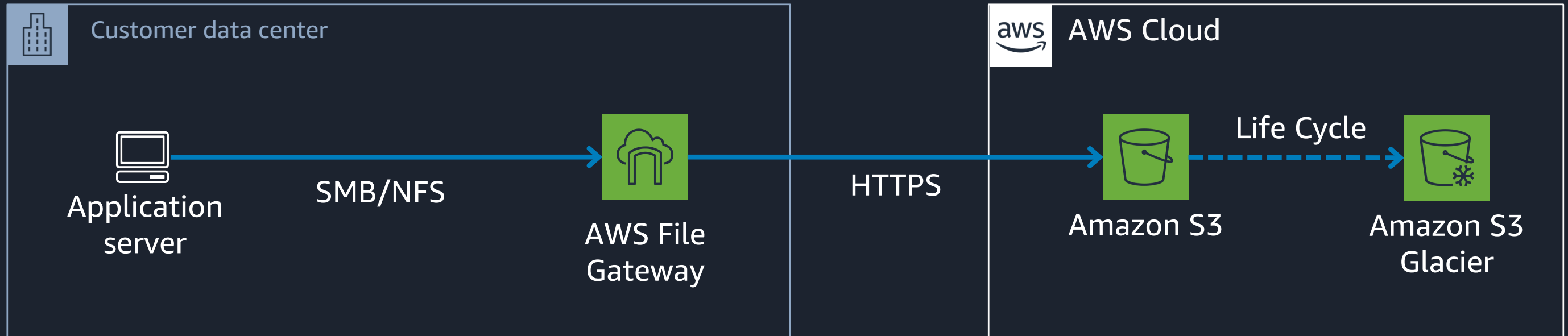AWS Backup

© 2020, Amazon Web Services, Inc. or its Affiliates.

aws

# Data Sources – Online ingest



Data center

Apps

Clients

Devices

NFS / SMB

?

AWS Cloud

Amazon S3

aws

# Data Sources – Online ingest



Data center

Apps

Clients

Devices

NFS / SMB

AWS File
Gateway

AWS Cloud

Amazon S3

aws

# What is AWS File Gateway?



- A virtual or hardware appliance that utilizes a network file system to interface to Amazon S3

- It allows for low latency access to hot data via it's local cache

- Stores file data in its native format as objects in an Amazon S3 bucket

- AWS File Gateway SMB shares can integrate with Microsoft Active Directory

# Using AWS File Gateway - Populate your data lake

o Seamless upload of created data to Amazon S3 using a file share

o Ingest existing data into Amazon S3 via a file share

Sharing your data lake contents

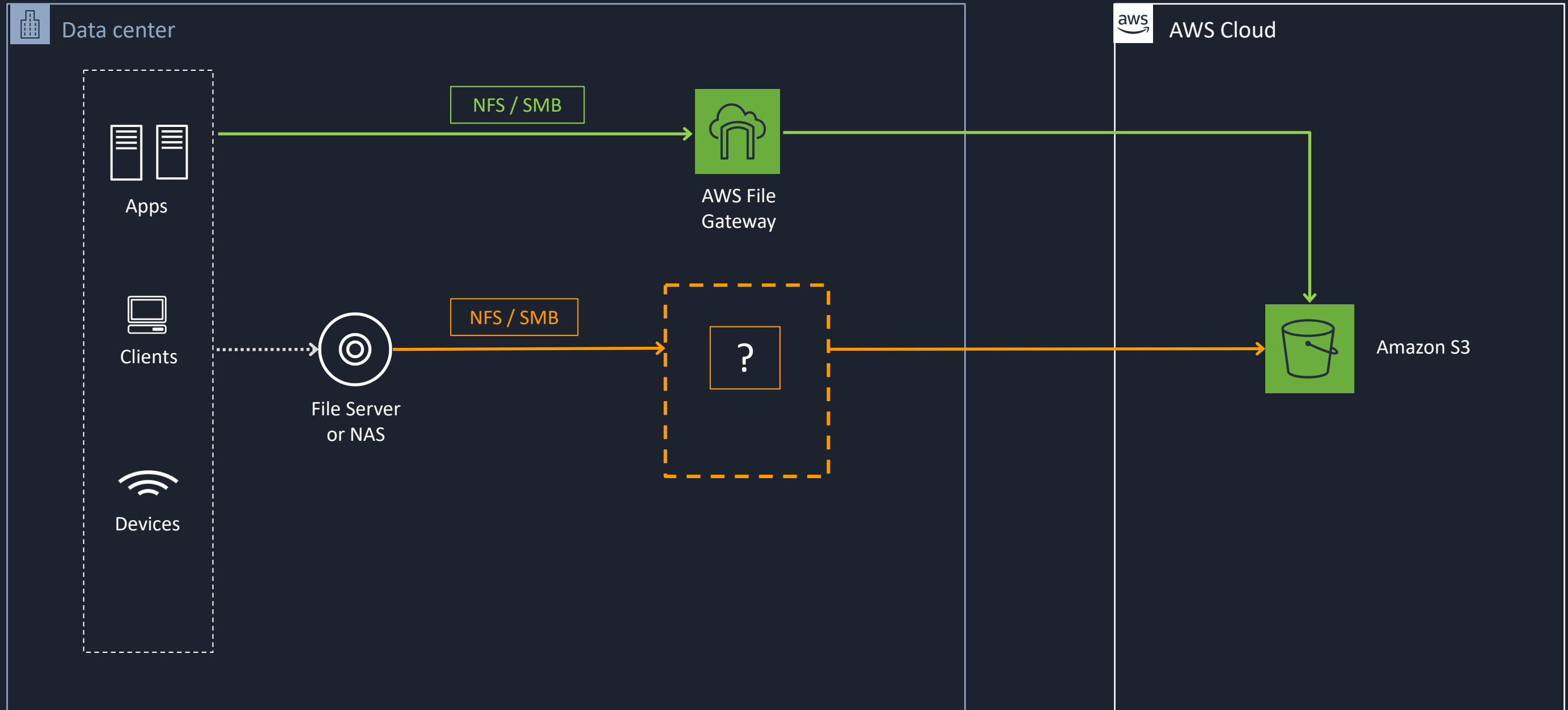o Re-share existing Amazon S3 data via simple file shares

aws

"A big challenge that we face is the integration of the wet labs with the computational aspect of our research.

We use Storage Gateway and DataSync to synchronize our on-premise wet labs with our Amazon storage. By having our scientists **immediately** be able to save their files **directly to the cloud**, they can go on with the next experiment without having to wait for the transfer times and they never run out of space."
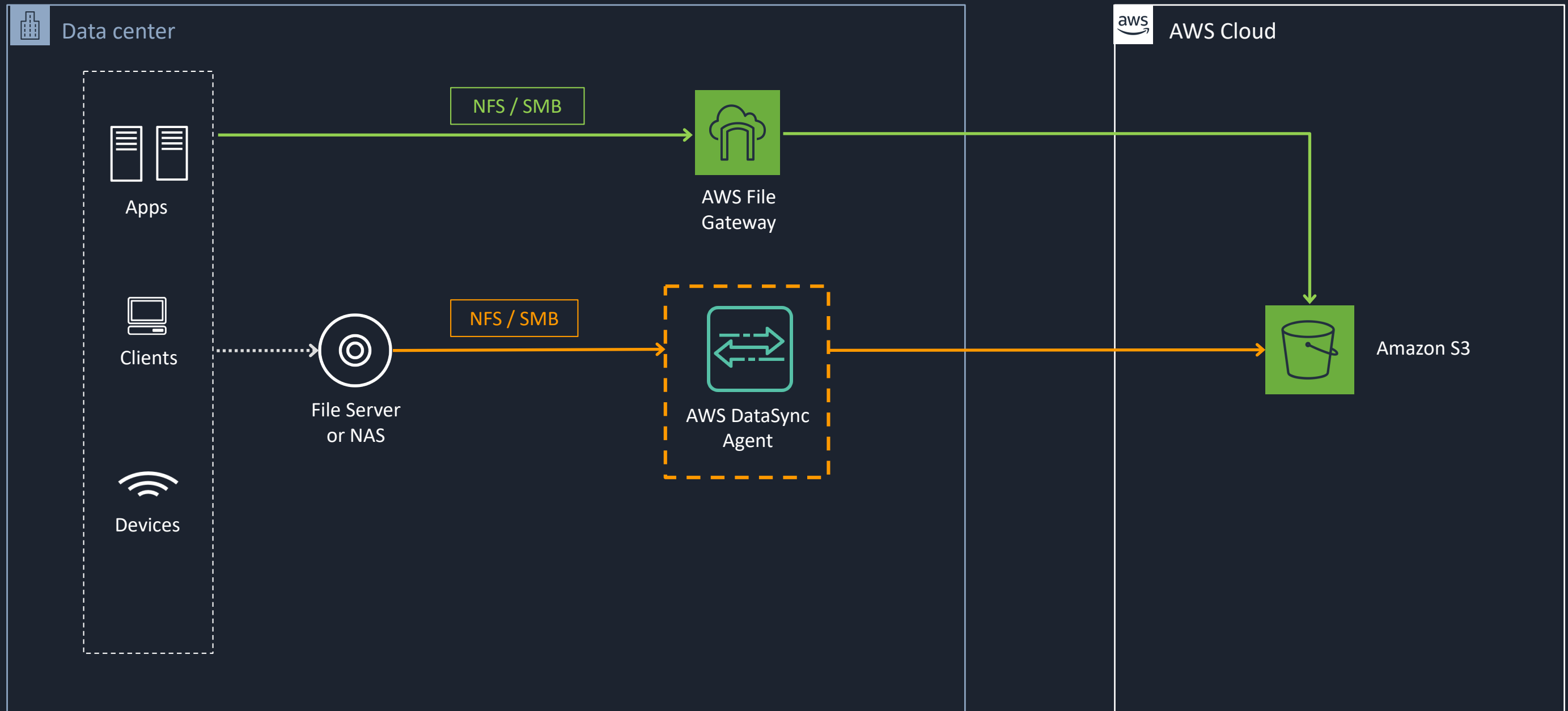
*Lance Smith*
*Associate Director - Celgene*

*https://aws.amazon.com/storagegateway/customers*

aws

# Data Sources – Online ingest



Data center

Apps

Clients

Devices

File Server
or NAS

NFS / SMB

AWS File
Gateway

NFS / SMB

?

AWS Cloud

Amazon S3

# Data Sources – Online ingest



Data center

Apps

Clients

Devices

File Server
or NAS

NFS / SMB

NFS / SMB

AWS File
Gateway

AWS DataSync
Agent

AWS Cloud

Amazon S3

aws

# What is AWS DataSync?

Simplifies, automates, and accelerates your online data transfer

**Migrate** active application data

**Transfer** data for timely processing or Archiving

**Replicate** for data protection and recovery

Transfers up to 10 Gbps per agent

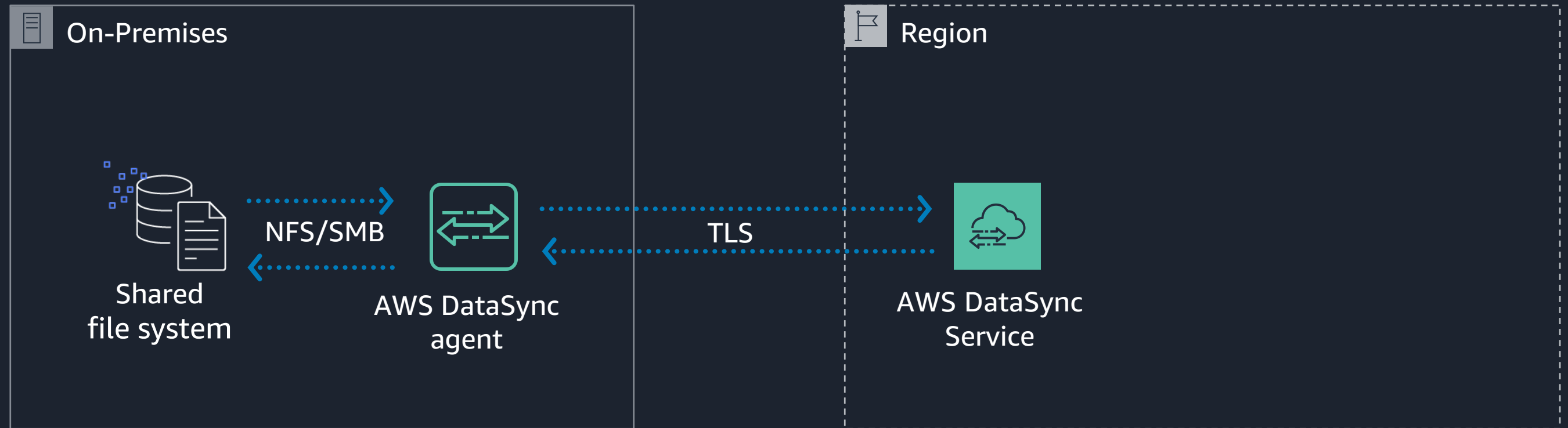Simple data movement to S3, EFS, FSx for Windows

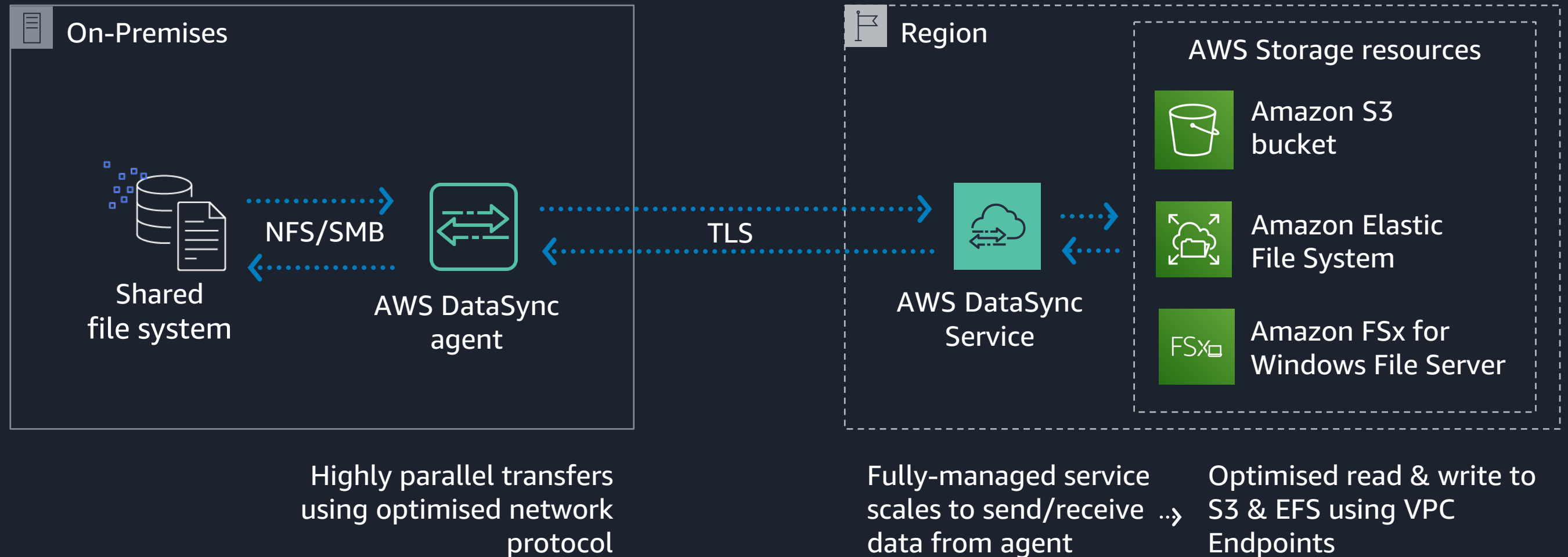Secure and reliable transfers

AWS

AWS integrated

Pay as you go

aws

# How AWS DataSync works



On-Premises

Shared
file system

NFS/SMB

AWS DataSync
agent

TLS

Region

AWS DataSync
Service

aws

# How AWS DataSync works



On-Premises

Shared
file system

NFS/SMB

AWS DataSync
agent

TLS

Region

AWS DataSync
Service

AWS Storage resources

Amazon S3
bucket

Amazon Elastic
File System

Amazon FSx for
Windows File Server

Highly parallel transfers
using optimised network
protocol

Fully-managed service
scales to send/receive
data from agent

Optimised read & write to
S3 & EFS using VPC
Endpoints

The speed and reliability of *network acceleration* software with the
cost-effectiveness of *open source tools*

aws

# Task options

Invoke via schedule, API, or manually

File-level validation

Copy across file metadata

Throttle bandwidth

# AWS DataSync usage scenarios

SMB/NFS data transfer → Amazon S3 storage classes

NFS data transfer → Amazon EFS

SMB data transfer → Amazon FSx for Windows File Server

aws

# Using AWS DataSync - Populate your data lake

o Simplify & accelerate the transfer of data into Amazon S3

o Utilize for end of event batch or scheduled uploads

o A repeatable data transfer mechanism for different use cases

o Think of its benefits with at-scale data ingest, with simplicity and automation in mind

aws

# Problem

Wanted to retire multi-petabyte on-premises Data Domain storage system

Data retention policies required data to be retained for many years

# Solution

Used Amazon S3 for low cost, pay-as-you-go model as well as versioning support

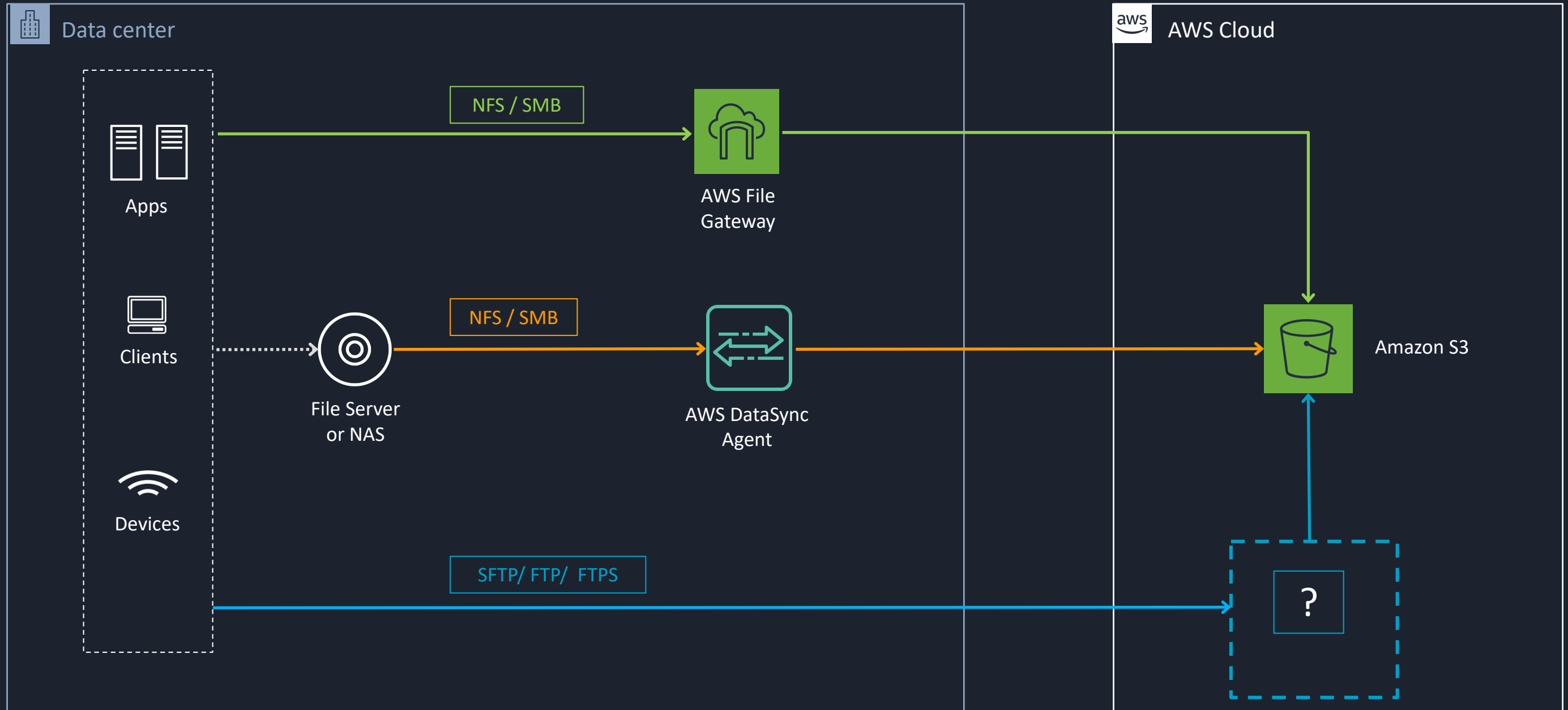Used AWS DataSync to seamlessly move data to Amazon S3

# Outcome

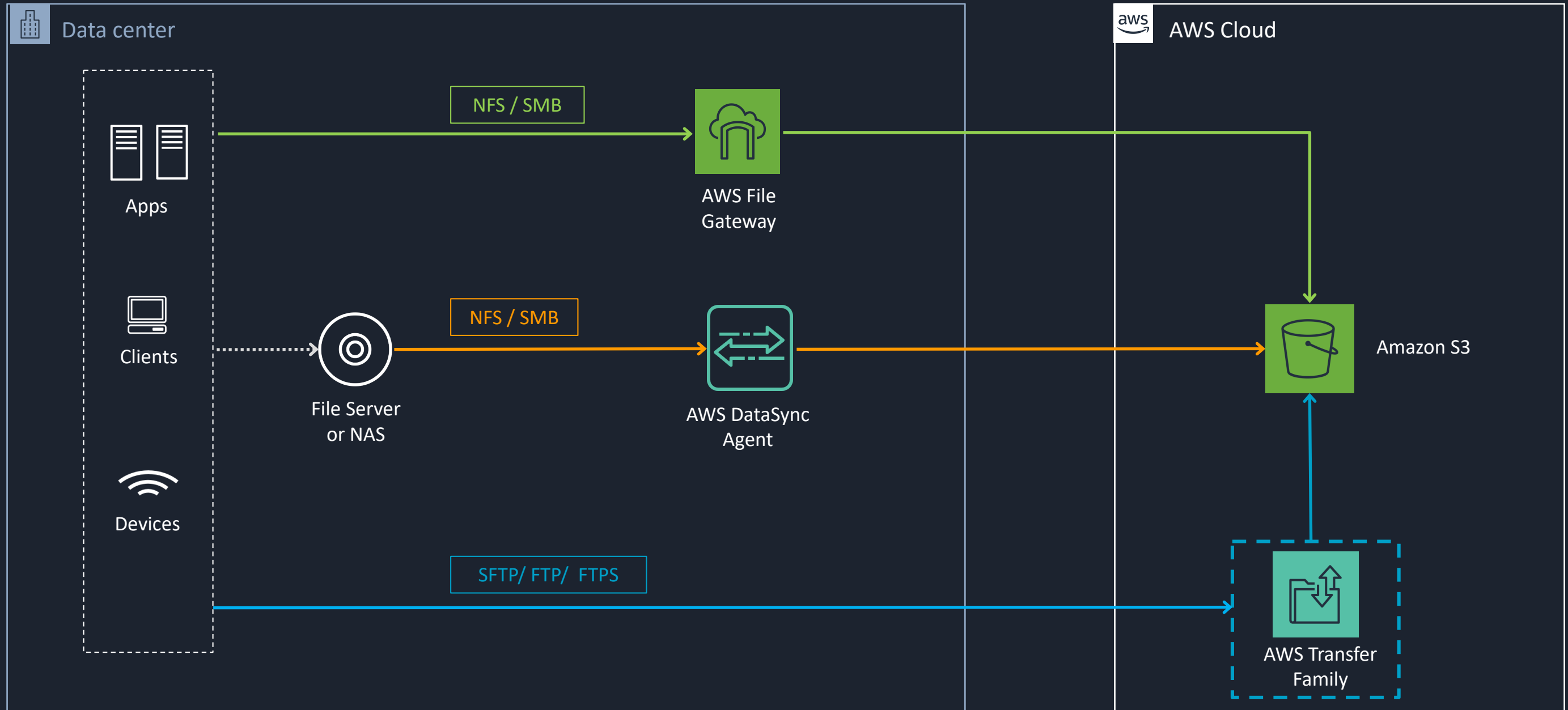Successfully transferred dataset to Amazon S3 with full byte-for-byte verification
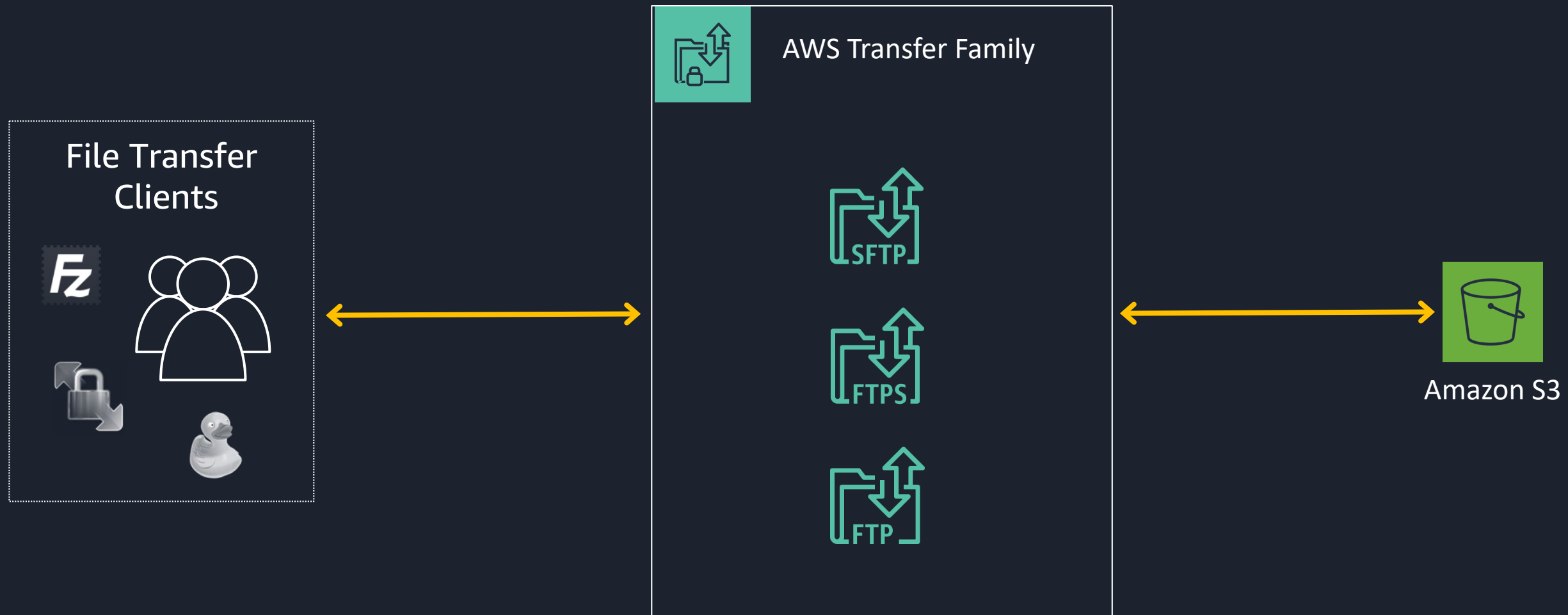
Decommissioned on-premises Data Domain

*"Our petabyte scale data migration journey from on-premises to AWS was accomplished swiftly with minimal effort and was completely self-managed with AWS DataSync. This solution is a game changer!"*

*Satish Kumar*
*Infrastructure Engineer*

https://aws.amazon.com/datasync/customers

# Data Sources – Online ingest



Data center

Apps

Clients

Devices

File Server or NAS

NFS / SMB

AWS File Gateway

NFS / SMB

AWS DataSync Agent

SFTP/ FTP/ FTPS

AWS Cloud

Amazon S3

?

# Data Sources – Online ingest

# What is the AWS Transfer Family?



File Transfer Clients

AWS Transfer Family

SFTP

FTPS

FTP

Amazon S3

aws

# AWS Transfer for SFTP

## Fully managed SFTP service enabling transfer of data into Amazon S3

Seamless migration
of existing SFTP workflows

Fully managed in AWS

Simple

Secure and compliant

Cost-effective

Native integration
with AWS services

aws

# Simple as 1 – 2 – 3

**①** Deploy an SFTP
server endpoint

**②** Select your target
S3 bucket(s)

**③** Set up users

aws

# How AWS SFTP works

Fully managed, highly available service provides secure access to data in Amazon S3



SFTP Users & Clients

No change to existing clients

SSH key-based authentication

Password authentication

AWS Cloud

AWS Transfer for SFTP

Amazon S3

Custom Identity Provider

Amazon API Gateway

AWS Lambda

" Our engineers were able to implement a near real-time customer usage analytics framework **within a week**. This timeframe is profound given the Mobile Call Detail Records ingestion comes in varying sizes and frequencies.

AWS SFTP **scales seamlessly** and makes the files available as S3 Objects, which is just perfect for our needs.

This enabled **event driven ingestion** of data into our **Data Lake**."

*Lambros Kallianiotis*
*Engineering Principal - Belong*

https://aws.amazon.com/aws-transfer-family/customers

# Using AWS Transfer for SFTP - Populate your data lake

o Seamless upload of data to Amazon S3 using SFTP clients

o Think of the service's simplicity at scale

o Utilize for end of event data uploads

aws

# Demo

aws

# Demo Setup – Online Ingest
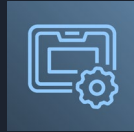
# Offline data ingest methods

aws

# Data Sources – Offline ingest



Data center

Apps

Clients

Devices

GUI / S3 / NFS

?

AWS Cloud

Amazon S3

aws

# Data Sources – Offline ingest

Data center

Apps

Clients

Devices

GUI / S3 / NFS

AWS Snow Family

AWS Cloud

Amazon S3

aws

# AWS Snow Family Portfolio

## AWS Snowball Edge

### Data transfer & edge compute

- 42/100TB storage capacity (S3)
- 10/25/40GE networking
- Data encryption end-to-end
- Rugged 8.5 G impact case
- Chain of Custody, Tamper Detection
- Rain and dust resistant
- EC2/AMI support  for edge computing
- NFSv4 Server

## AWS Snowmobile

### 20+ PB data transfer

- Exabyte-scale storage in a 45ft container (90PB s3/Glacier/EBS)
- 10/25/40GE networking
- Data encryption end-to-end
- S3/Glacier Data import
- Dedicated security personnel
- GPS tracking, alarm monitoring, 24/7 surveillance, and optional additional security

aws

# AWS Snowball Edge import workflow

**Create a job**

**Connect the Snowball Edge**

**Copy & Process Data on device**

**Your data is loaded into your bucket!**

**In transit**

**Delivered to you**

**Sent to AWS**

**Ingest to AWS**

**Job created**

| 1 – 3 Days | 1 – 2 Days | 1-N Days (file transfer speed dependent) | 1 – 2 Days | 1-N Days |

**Job Provisioned**

**At AWS Region**

**Job completed**

aws

# Ingesting data into AWS Snowball Edge

# AWS OpsHub for Snow Family

## GUI for customers to easily manage Snow devices

# Using AWS Snowball Edge - Populating your data lake

o Bulk data transfers - data is not required immediately

o Bulk data transfers -  limited network bandwidth

o Perform data transformation at edge using AWS Snowball Edge compute to pre-process data before ingest into Amazon S3 data lake

aws

# Data transfer at scale with AWS Snowball Edge

PH☆TOBOX

Europe's #1 online photo service

Migrated to AWS from two data centers

- Used AWS Snowball Edge to move 10 PB (5.7 billion) of photos from Dell EMC Isilon and IBM Cleversafe to Amazon S3
- Needed Amazon S3 data durability; higher than colos and other clouds

Shifted investments and focus to innovation and product development for customers, away from IT infrastructure
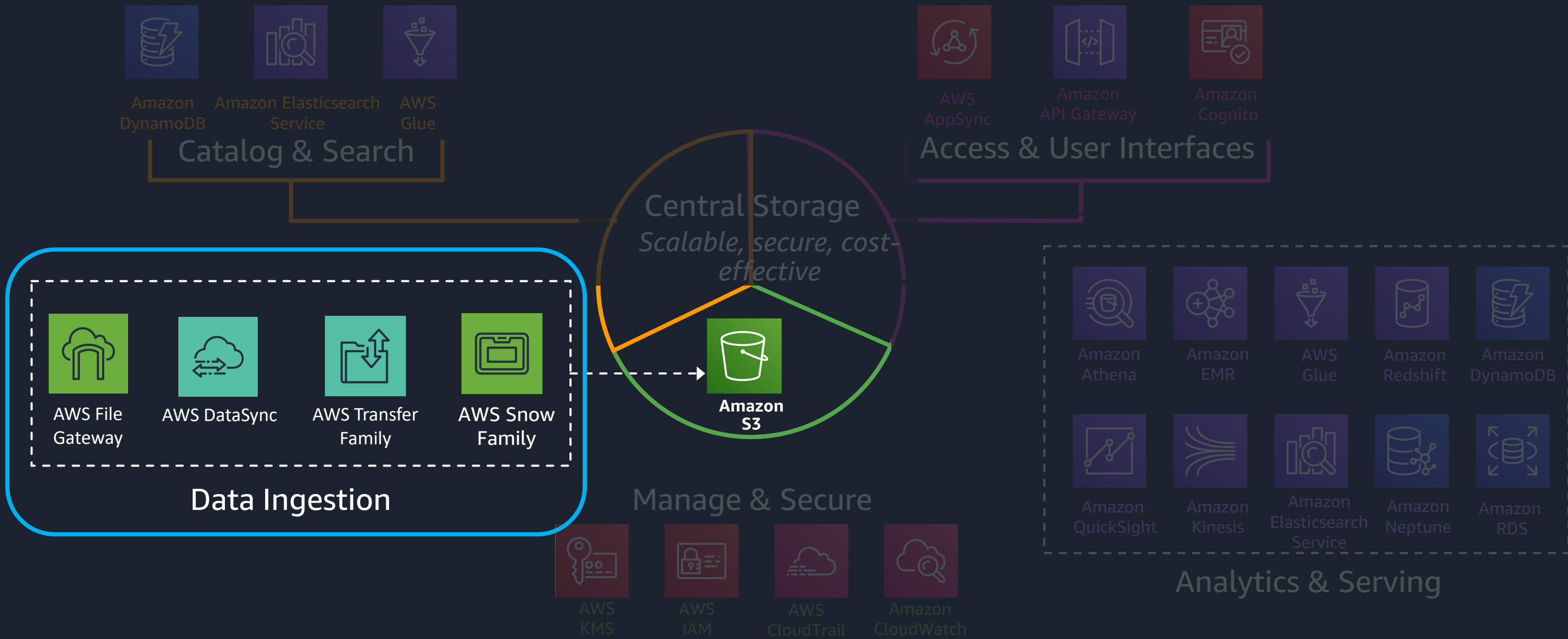
"We've reduced costs. We've improved our customer experience. Generally, we've made our website faster. And that's because AWS manages that infrastructure in a way we could never do internally."

— **Chris Astall, Group Director of Architecture**

https://aws.amazon.com/solutions/case-studies/photobox/

Richard Orme
Group CTO
photoboxgroup

aws

# Recap – Ingesting data into your data lake on AWS

**Catalog & Search**
- Amazon DynamoDB
- Amazon Elasticsearch Service
- AWS Glue

**Access & User Interfaces**
- AWS AppSync
- Amazon API Gateway
- Amazon Cognito

**Central Storage**
*Scalable, secure, cost-effective*

Amazon S3

**Data Ingestion**
- AWS File Gateway
- AWS DataSync
- AWS Transfer Family
- AWS Snow Family

**Manage & Secure**
- AWS KMS
- AWS IAM
- AWS CloudTrail
- Amazon CloudWatch

**Analytics & Serving**
- Amazon Athena
- Amazon EMR
- AWS Glue
- Amazon Redshift
- Amazon DynamoDB
- Amazon QuickSight
- Amazon Kinesis
- Amazon Elasticsearch Service
- Amazon Neptune
- Amazon RDS

aws

# Get up and running with these resources

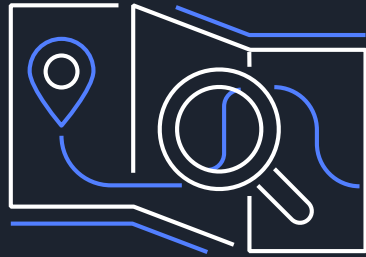Get hands-on experience with the AWS online data migration workshop

https://github.com/aws-samples/aws-online-data-migration-workshop

AWS DataSync : https://aws.amazon.com/datasync

AWS Snow Family :  https://aws.amazon.com/snow

AWS Storage Gateway : https://aws.amazon.com/storagegateway

AWS Transfer Family : https://aws.amazon.com/aws-transfer-family

aws

# AWS Training and Certification



## Training for the Whole Team

Explore tailored Data or Database learning paths for customers and partners

## Flexibility to Learn Your Way

Build cloud skills with free digital Data training courses such as "The elements of Data Science", or dive deep with classroom training

## Validate Skills with AWS Certification

Demonstrate expertise with a Data industry-recognized credential (Data analytics and Database Specialty AWS Certifications)

aws.amazon.com/training/

aws

# Visit the Data, Databases, and Analytics Resource Hub for more resources

Dive deeper with these newly created whitepapers and e-books to help you uncover new insights and value from your data

- An introduction to cloud databases
- Enter the purpose-built database era
- Harness the power of data
- Creating a modern analytics architecture
- The data-driven enterprise
- … and more!



https://tinyurl.com/
aws-data-databases-analytics

**Visit resource hub »**

aws

# Thank you for attending
# AWS Data, Databases, and Analytics Online Series

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

aws

# Thank you!

aws