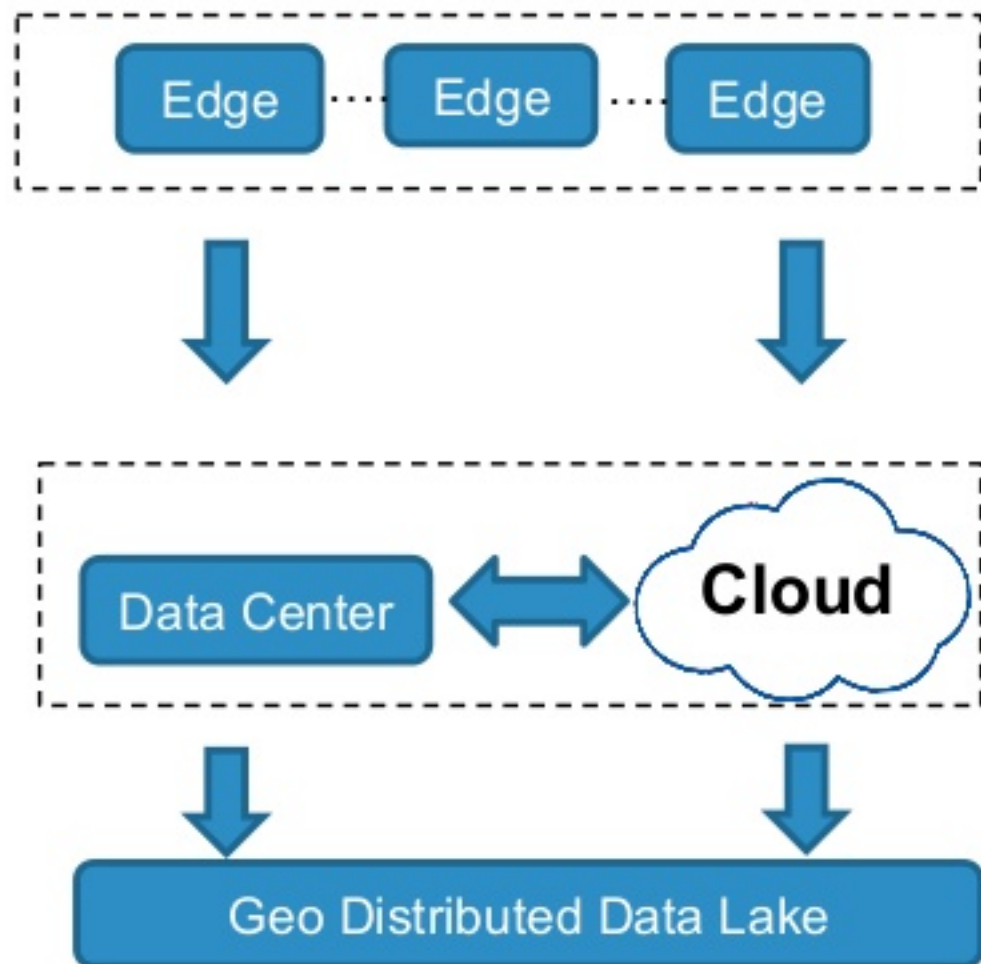# Agenda

- IoT Data Management Challenges

- NetApp Data Fabric Architecture for Big Data

- IoT Customer Use Cases on Data Fabric

- Q&A

IoT Data Management Challenges

# IoT Data Flow



**EDGE**

1. Data is created
2. Data is analyzed in realtime
3. Data is aggregated and sent to Core

**CORE**

1. Data is stored
2. Data is analyzed
3. Data is protected

# IoT Data Management Challenges

| Collect | Transport | Store | Analyze | Protect |
|---------|-----------|-------|---------|---------|

101000010101
101010101000

101000101011
101010101000

101000101011
101010101000

101000101011
101010101000

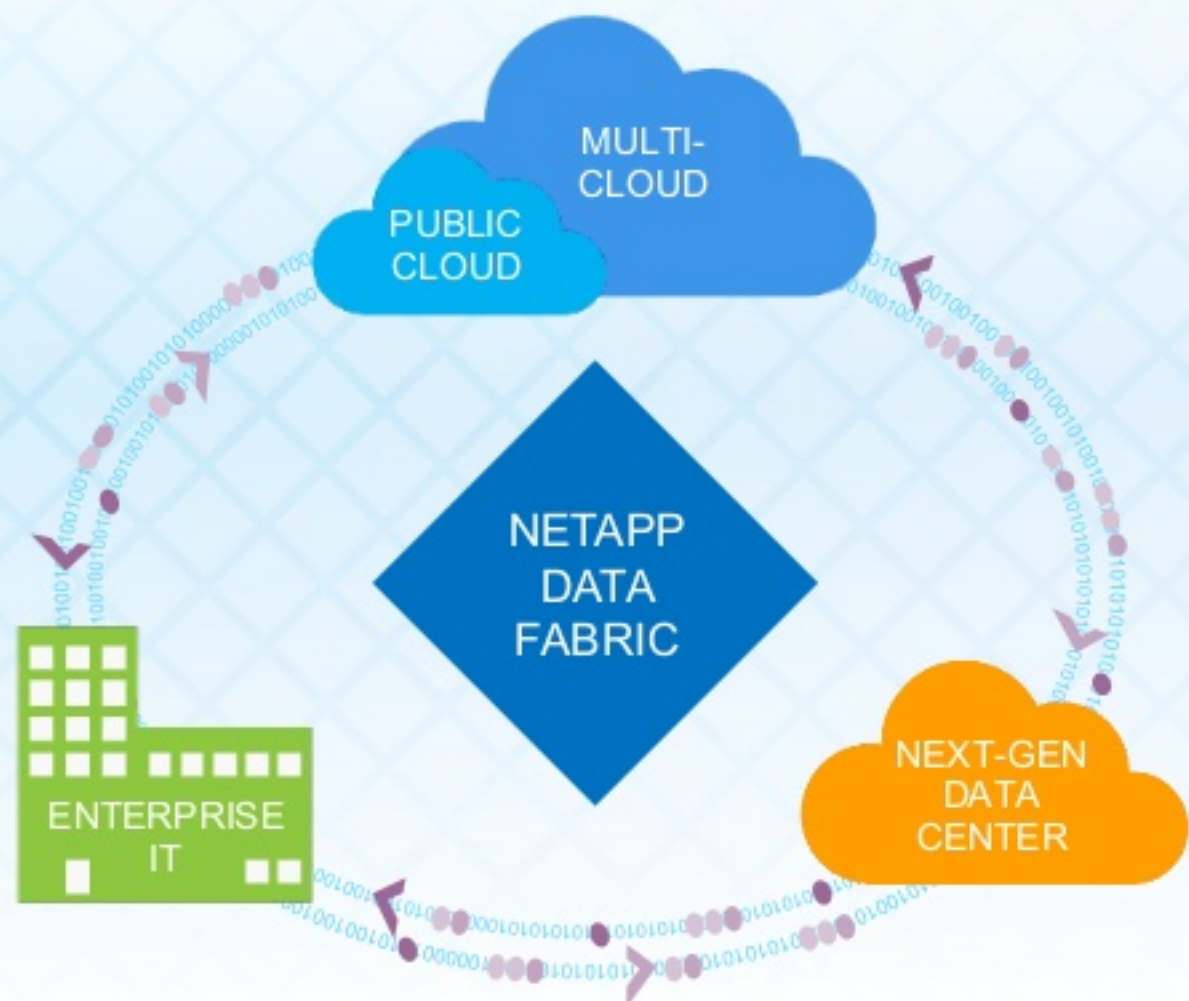- **Flexibility and Agility**
- **Cost**
- **Data Protection**

# NetApp Data Fabric Architecture for Big Data

# The NetApp Data Fabric

Helping customers unleash data to address their business imperatives

MULTI-CLOUD

PUBLIC CLOUD

NETAPP DATA FABRIC

ENTERPRISE IT

NEXT-GEN DATA CENTER
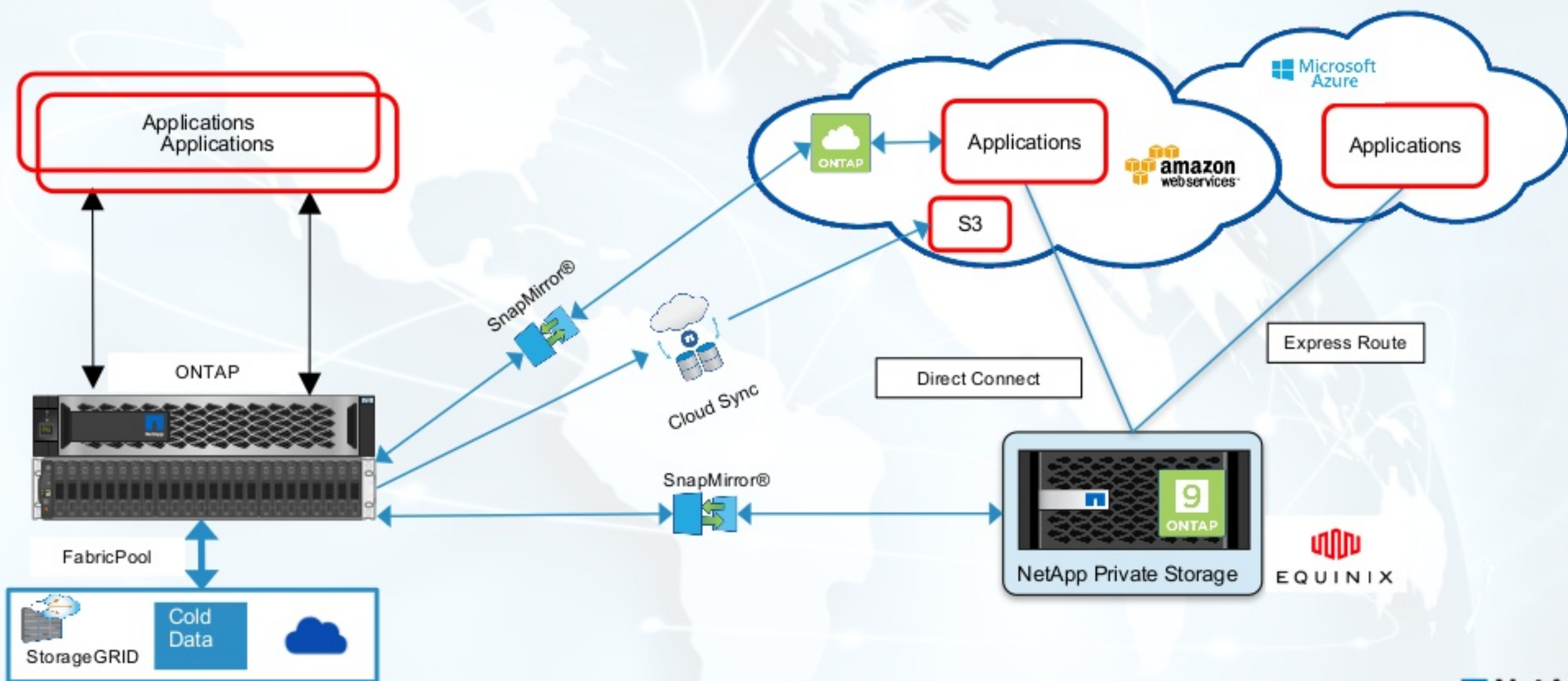
**HARNESS**
the power of the hybrid cloud

**BUILD**
a next-generation data center

**MODERNIZE**
storage through data management

**∏ NetApp**

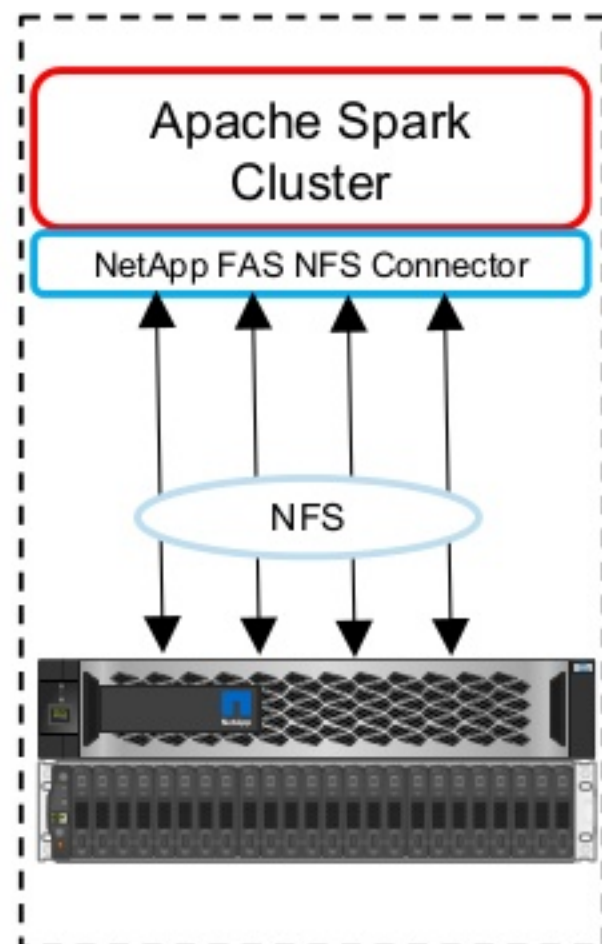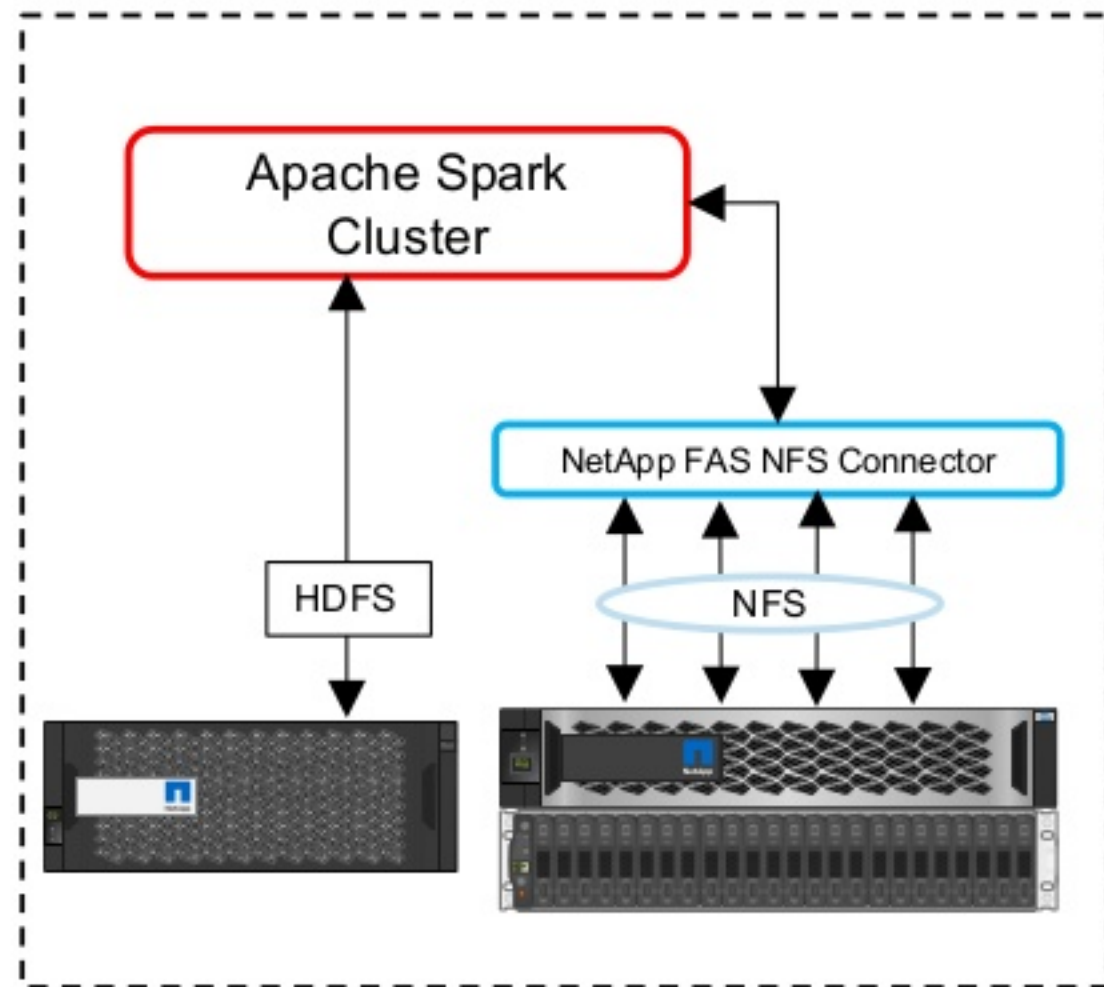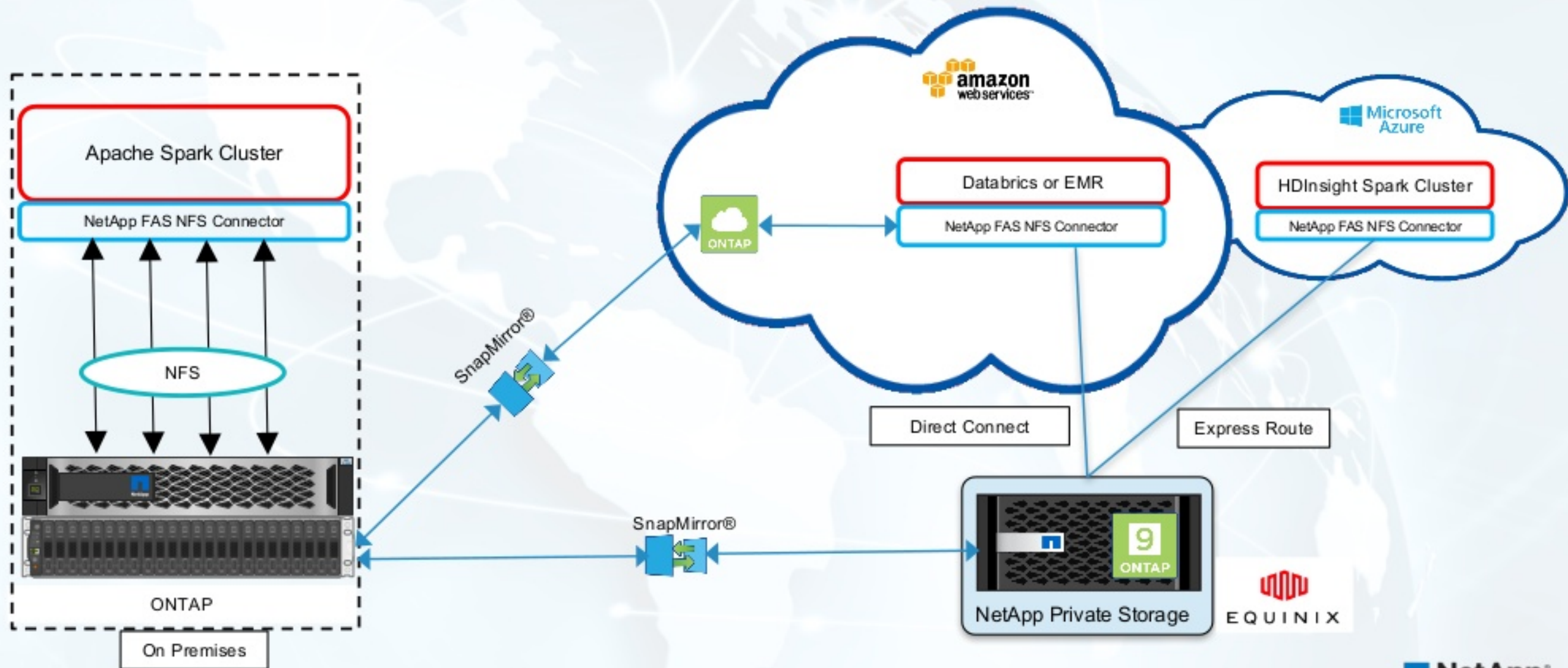# Introducing Data Fabric Building Blocks for Analytics

# In Place Analytics:

**Key Benefits**

- Avoid data move to HDFS. Reduce replicas

- Scale compute and storage independently

- Enterprise data protection

- Hybrid cloud deployment

- Hortonworks Certified

Apache Spark Cluster

NetApp FAS NFS Connector

NFS

Confit 1 : NFS as a Storage

Apache Spark Cluster

NetApp FAS NFS Connector

HDFS

NFS

Confit 2 : HDFS and NFS in Single Spark Cluster

**NetApp**

# Analytics with Data Fabric
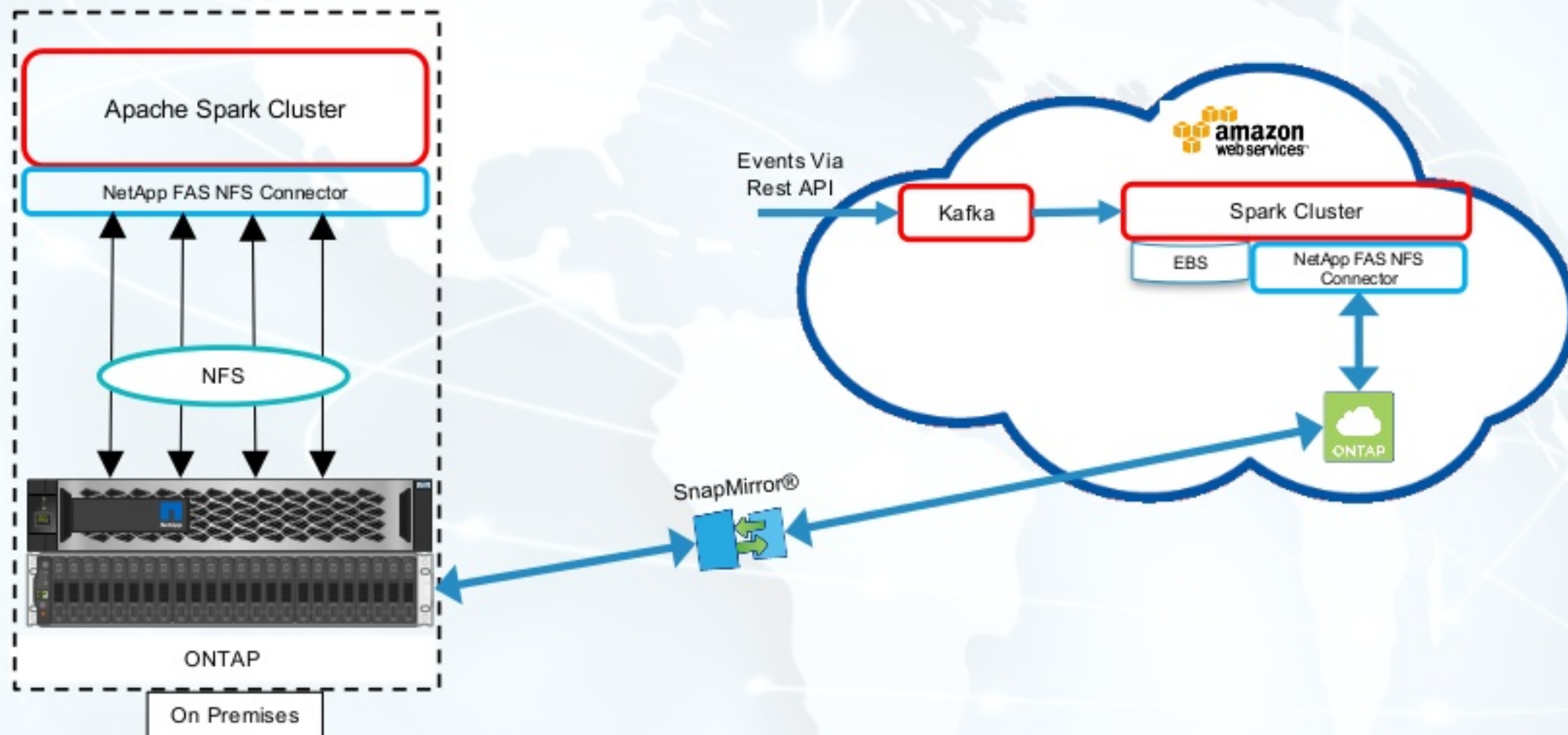
# Customer Scenario

Broadcasting Provider

- IoT data received in AWS and analyzed using **Apache Spark cluster** in AWS

- Data Management Challenge:

  - How to Backup 10 TB data without load on cluster?

  - How to protect the data to on-premises?

# An Architecture for Processing IoT-Data Ingested in Cloud
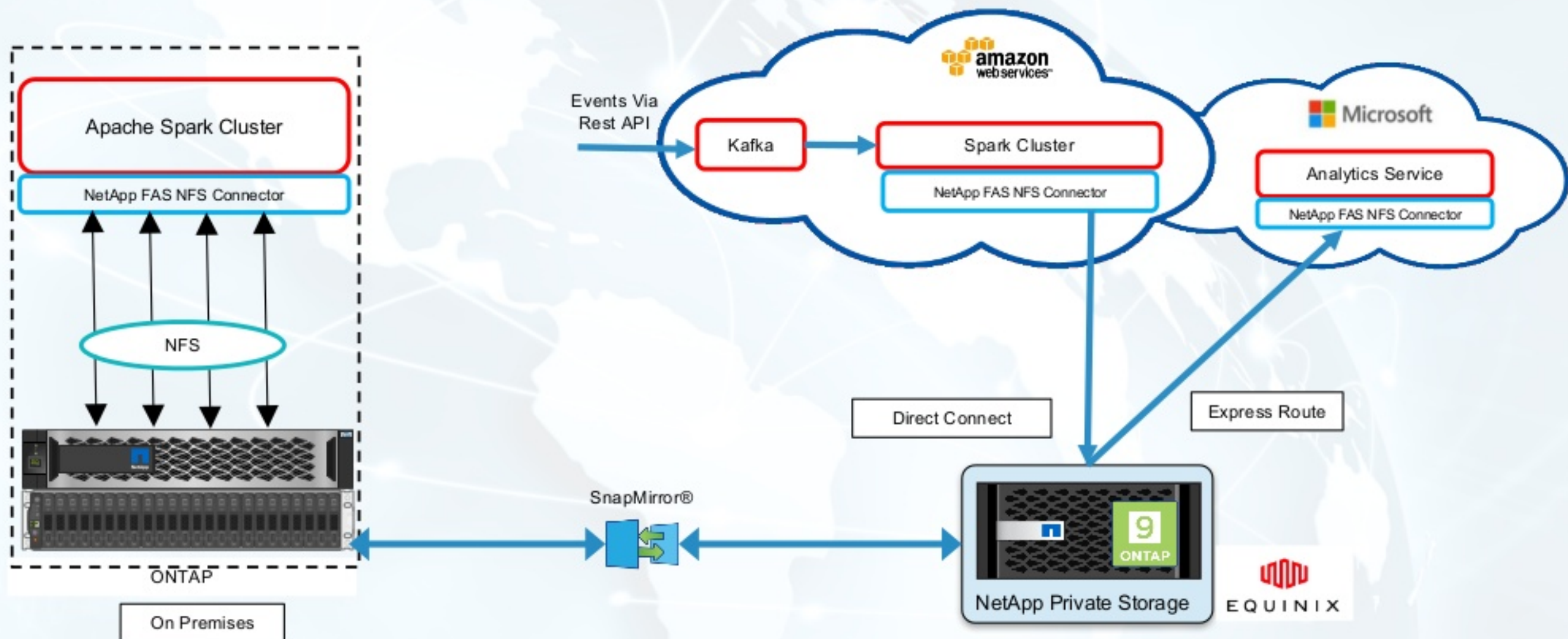
Backup and DR to On Premises

# Customer Scenario

- IoT data is received in AWS and analyzed using **Apache Spark Cluster** in AWS

- Data Management Challenge:
  - How to reduce the solution cost?
  - How to consume analytics services in data center and multiple clouds?

SPARK
SUMMIT
2017

# An Architecture for Processing IoT-Data Ingested in Cloud
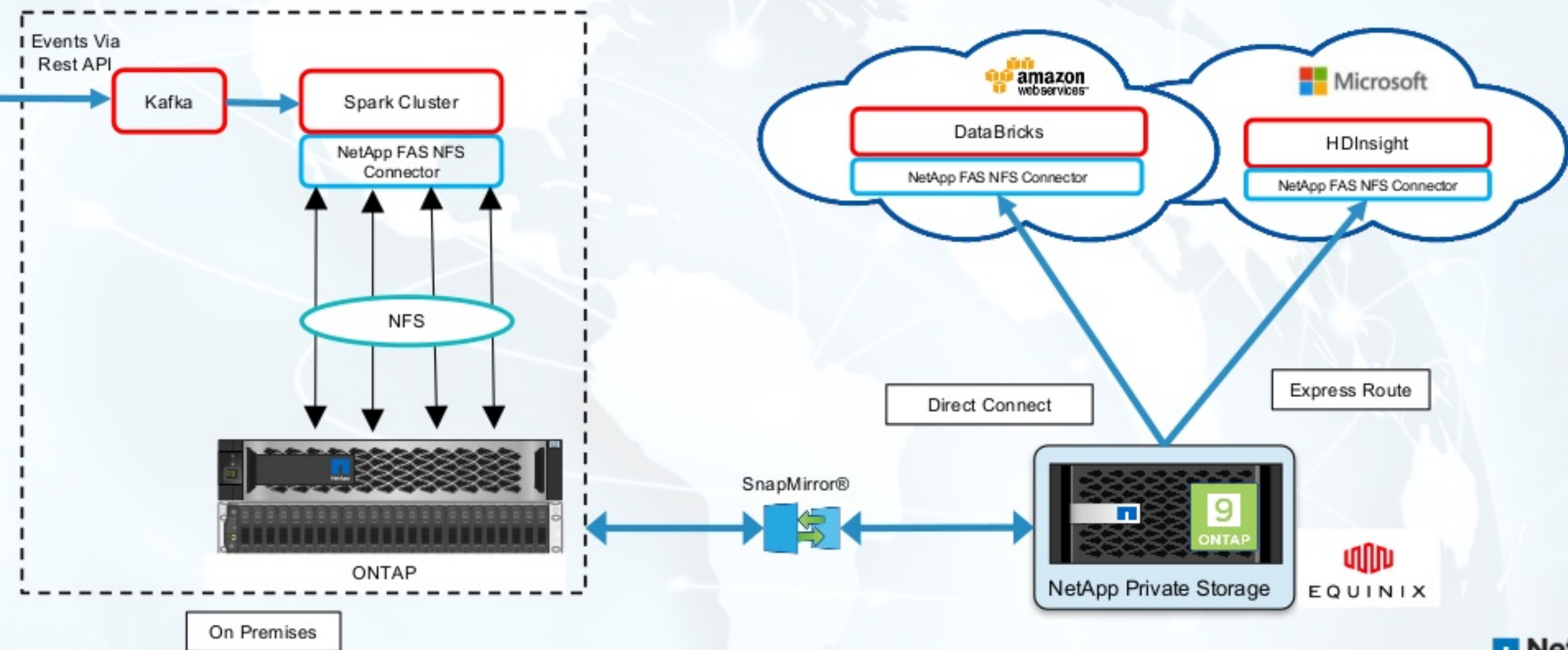
## Multi Cloud Connectivity

# Customer Scenario

Insurance Company

- IoT data received on-premises and analyzed using **Cloudera Spark Cluster** across data center and cloud

- Data Management Challenge:

  - How to leverage cloud computation for analytics ?

  - How to consume legacy data (7PB) for analytics?

# An Architecture for Processing IoT-Data Ingested on premises
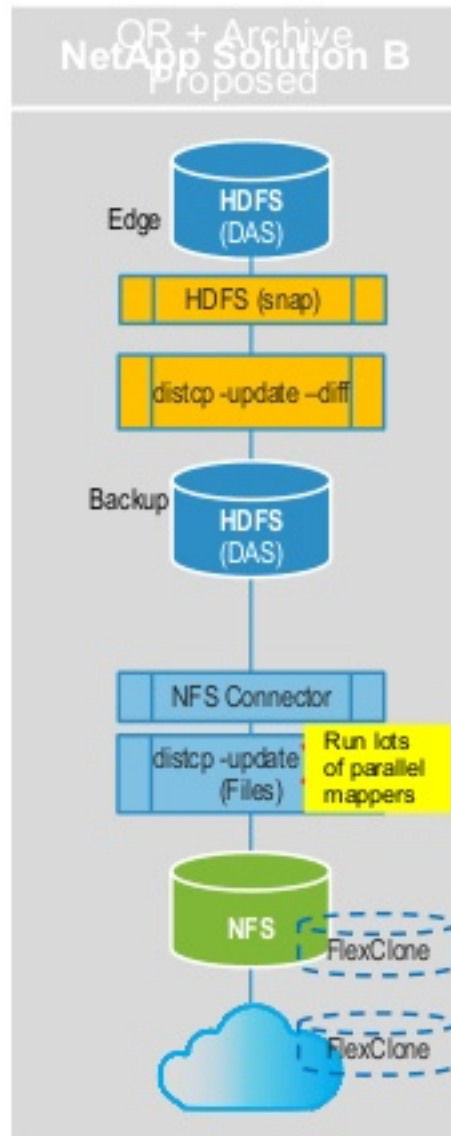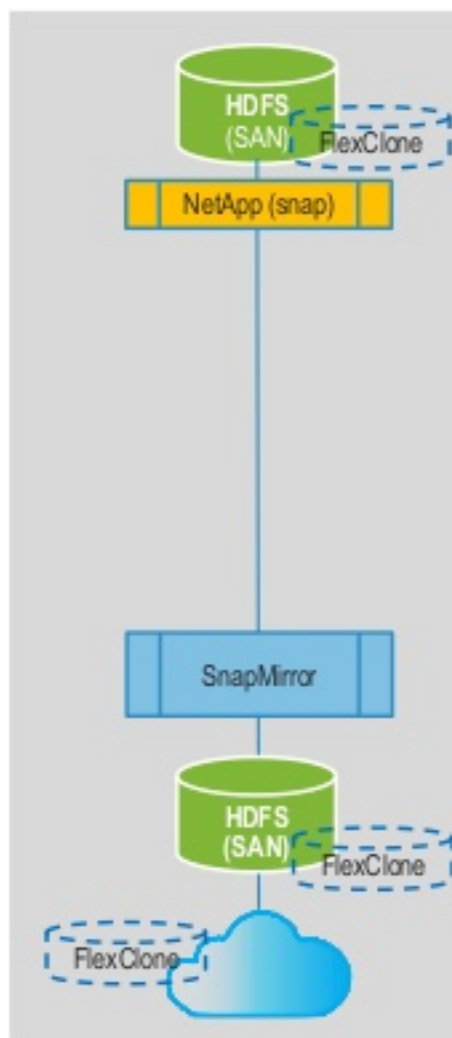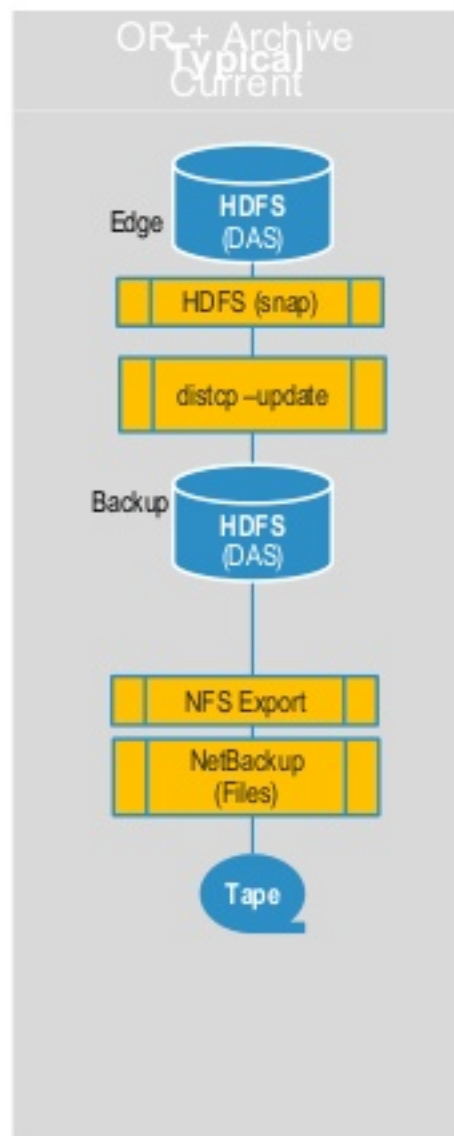
## DR in Cloud; Analytic across data centers

# Customer Scenario

Large Bank

- IoT data received on-premises and stored in Hadoop Data Lake. Data needs to be backed up for compliance

- Data Management Challenge:
    – How to reduce the backup window and optimize solution cost?
    – How to minimize impact on analytics performance during backup?

# Use Case: Backup for IoT Data



**NetApp Backup Solution A**

- NetApp Snapshot Backup
- Backup Archival
- Cloud Compatible

**NetApp Backup Solution B**

- Hadoop Native Support
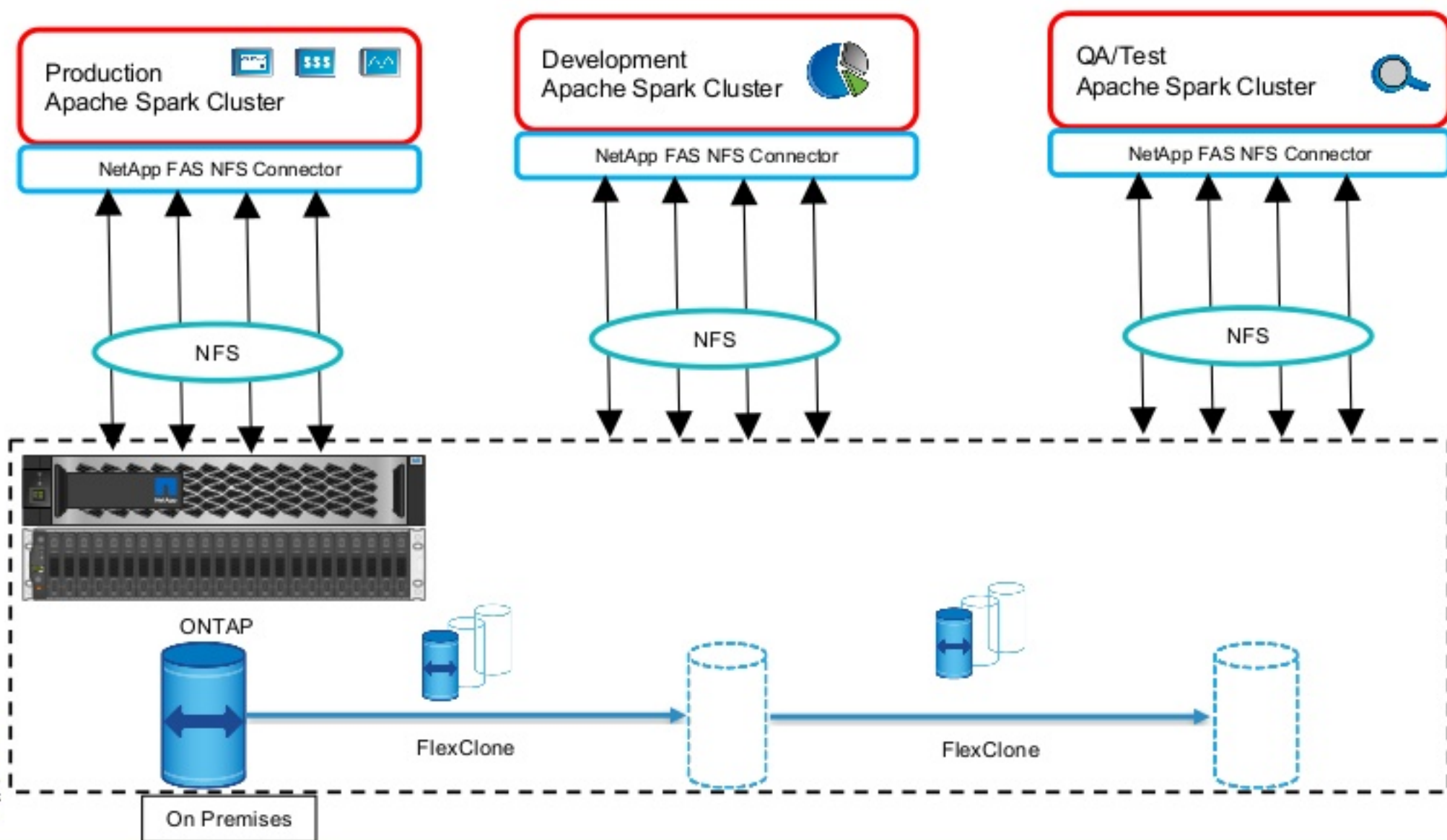- Offload Backup Operation
- Enterprise Management

# Customer Scenario

Online Music Distribution

- Large Hadoop Data Lake implementation on premises with Multiple Spark clusters

- Data Management Challenge:

    - How to make data available for dev/test teams?

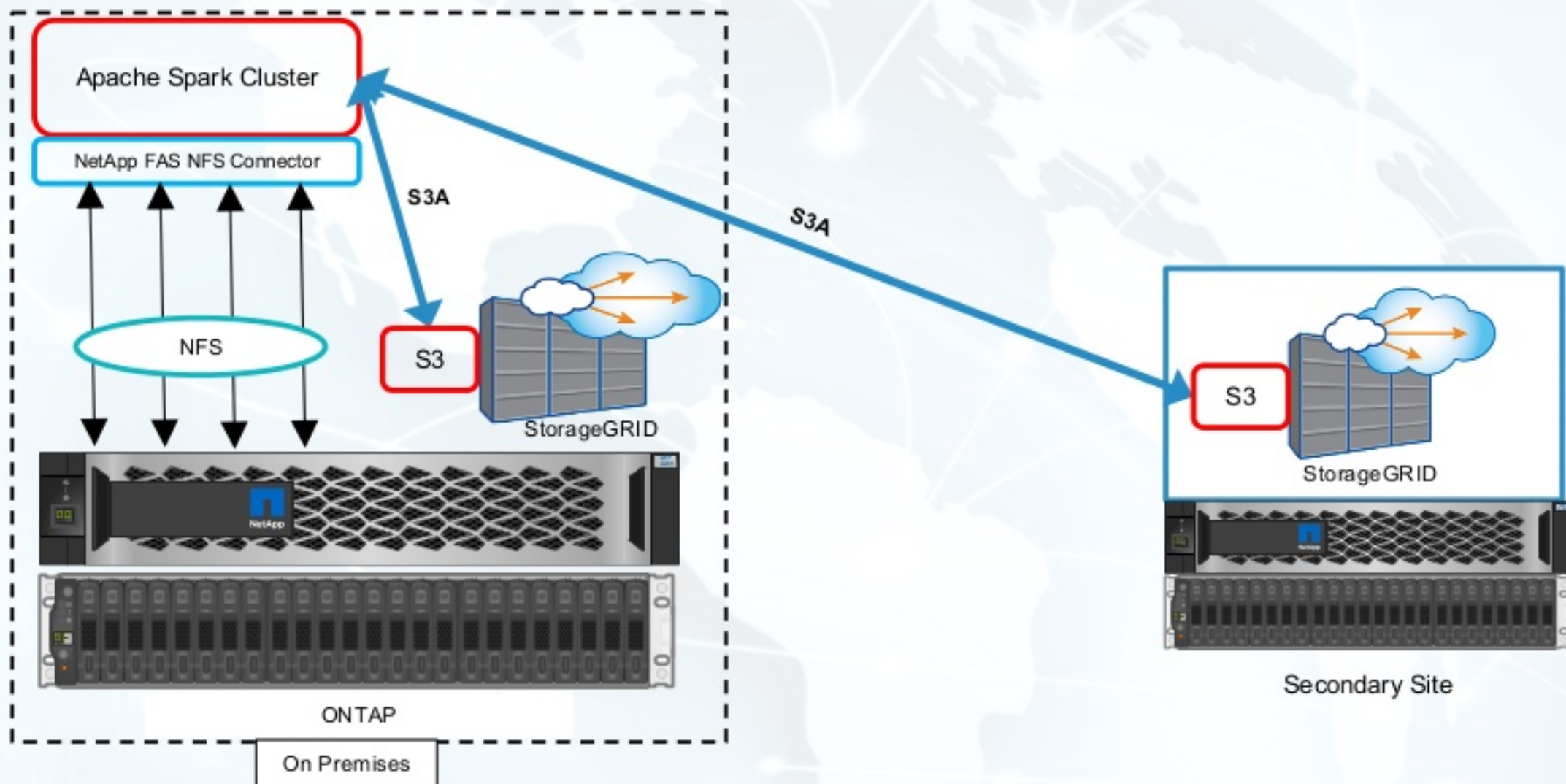    - How to build the new cluster in minutes from an existing cluster?

# Customer Scenario

Online Marketplace

- Run analytics on archival data in object store

- Data Management Challenge:

  – How to run Hadoop jobs in object store

  – Archive the Hadoop data on primary or a secondary site

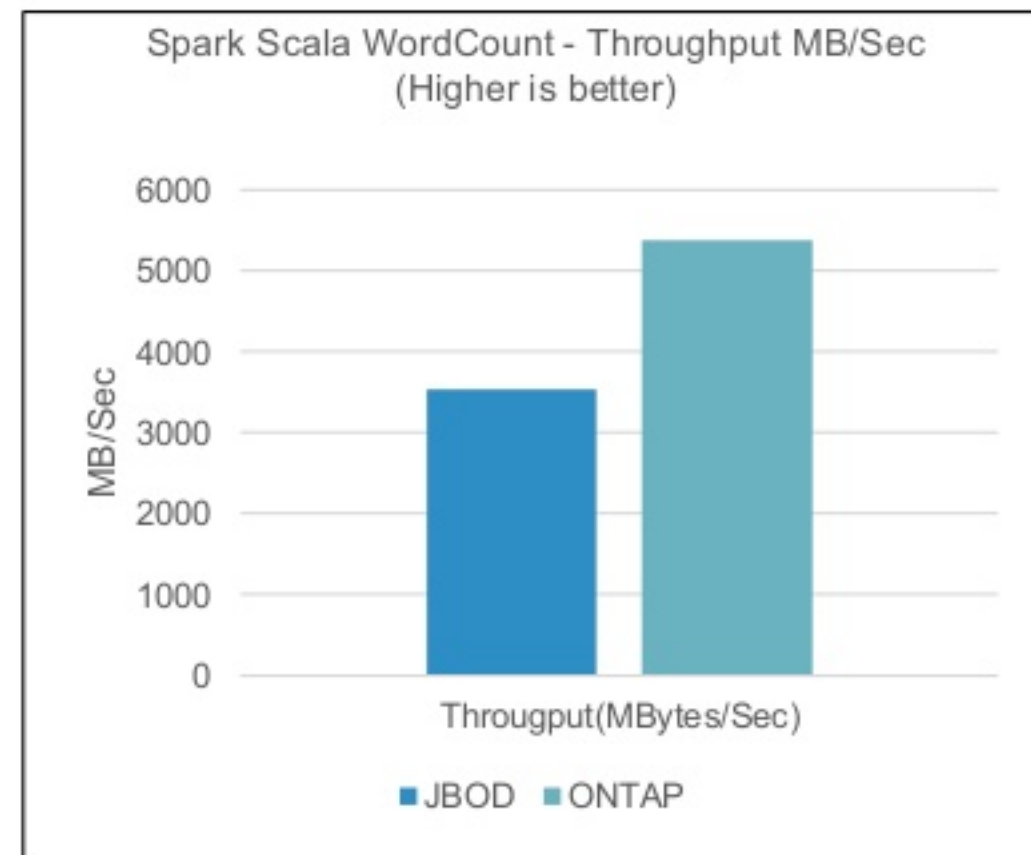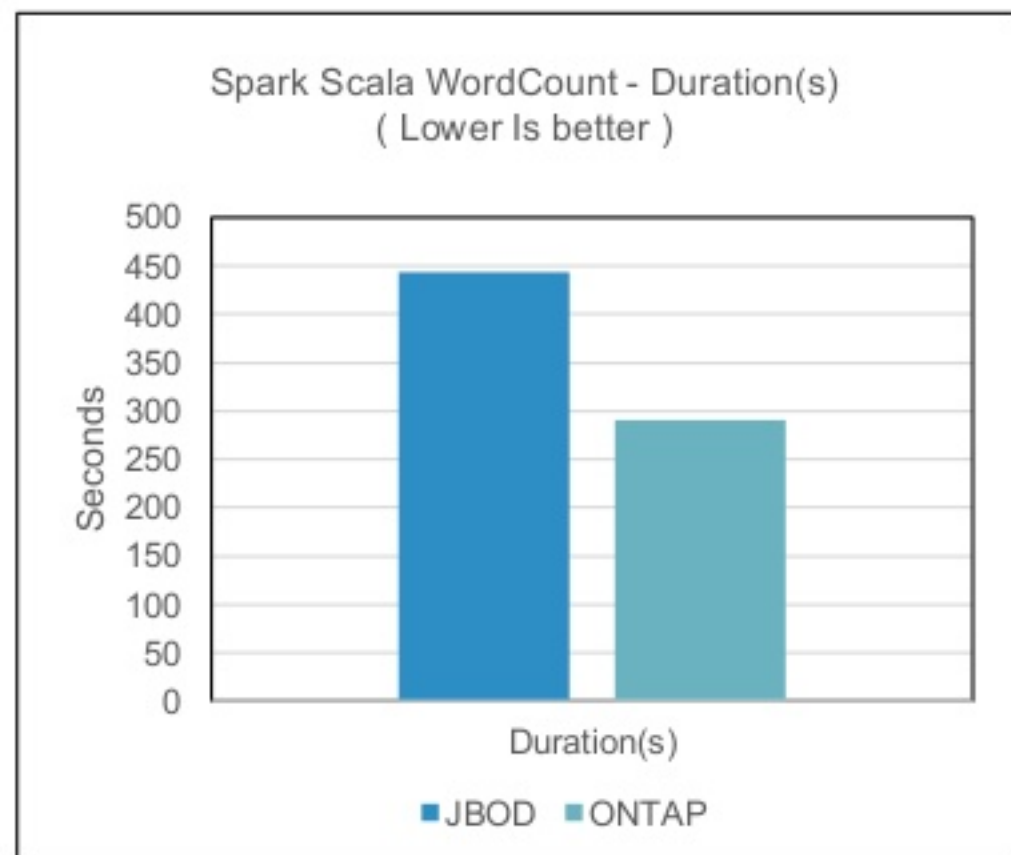# Analytics with NetApp StorageGRID

In place analytics with StorageGRID

Apache Spark Cluster

NetApp FAS NFS Connector

S3A

S3A

NFS

S3

StorageGRID

S3

StorageGRID

ONTAP

On Premises

Secondary Site

**NetApp**

# Spark Performance
## HiBench – Spark Scala Word Count



Spark Scala WordCount - Duration(s) ( Lower Is better )



Spark Scala WordCount - Throughput MB/Sec (Higher is better)

- Input dataset size – 1.5TB
- ONTAP– 52% better than JBOD

| Type | Hadoop Worker Nodes | Drives per Node | Number of Storage Arrays |
|------|--------|-------|---------|
| JBOD | 6 | 12 | NA |
| ONTAP | 6 | 6 | 1 |

# Key Takeaways

**Flexibility and Agility**

- On Demand analytics with Hybrid Cloud/Multi Cloud deployments
- Rapid provisioning of clusters for test/dev environments

**Lower Cost**

- Add storage capacity without adding compute nodes
- One copy vs 3 copies of data for HDFS
- Data Tiering with FabricPool

**Enterprise Data Protection**

- Efficient backup, DR and Archival

**n NetApp**

# Further Resources

- Please visit us at: Booth #512

- Visit our Big Data Website: www.netapp.com/bigdata

**■ NetApp**

Q & A

# Thank You.

Nilesh Bagad: nileshb@netapp.com

Karthikeyan Nagalingam: nkarthik@netapp.com

SPARK SUMMIT 2017