

Big Data Processing

Anusuriya Devaraju

01.03.2017

What is Big Data?



Source: <https://thingsthatresonate.wordpress.com/tag/knowledge/>

The Big Data Pipeline



Outline

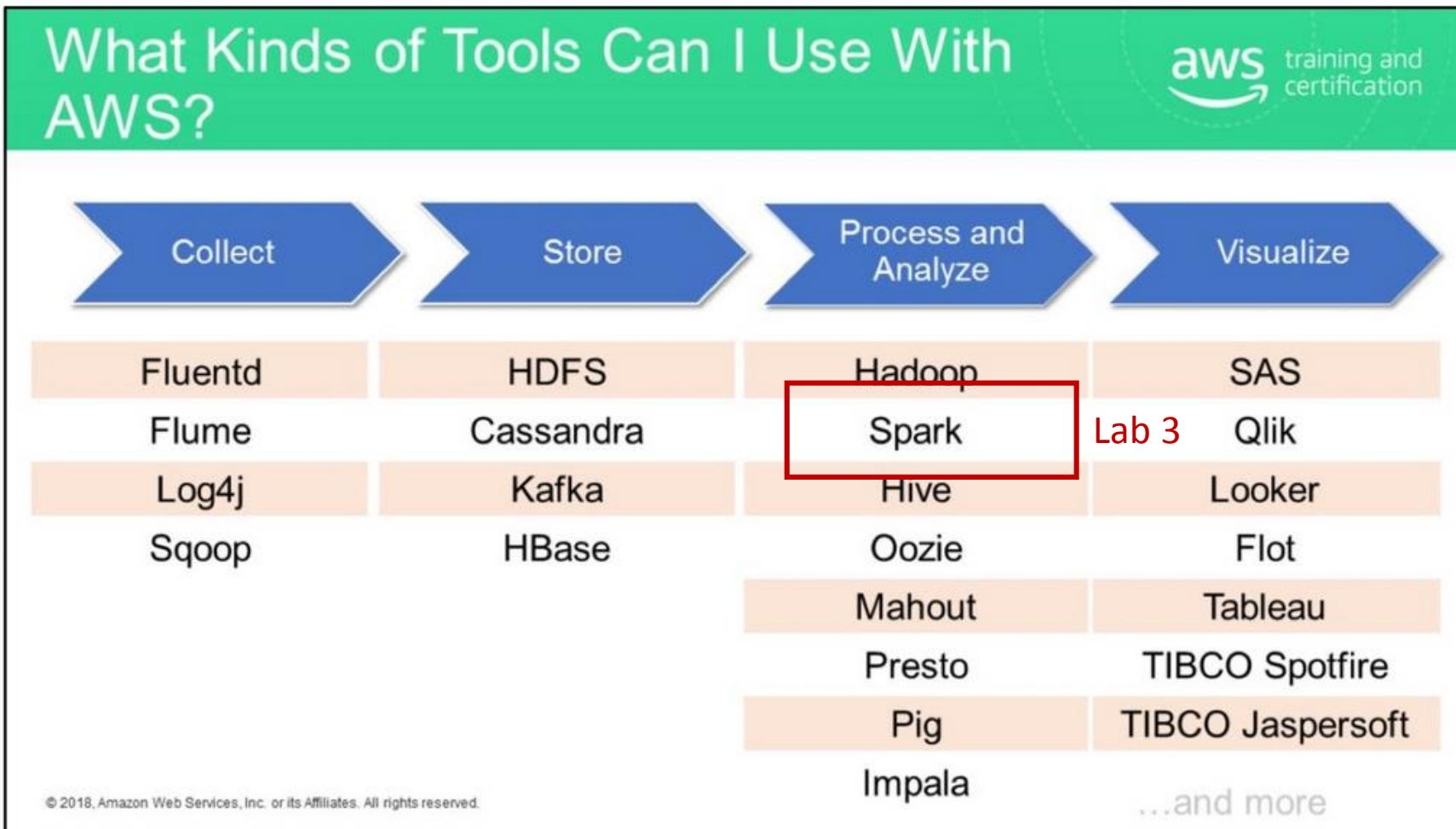
- Basic Concepts
- Hands-on & Review Questions
 - Lab 1 - Using Amazon Athena to Analyze Log Data
 - Lab 2 - Processing Server Logs with Hive on Amazon EMR
 - Lab 3 - Apache Spark on Amazon Web Services EMR

Reference: Some of the following slides were modified from AWS. Big Data on AWS 3.2 (EN): Student Guide. AWS/Gilmore. VitalBook file.

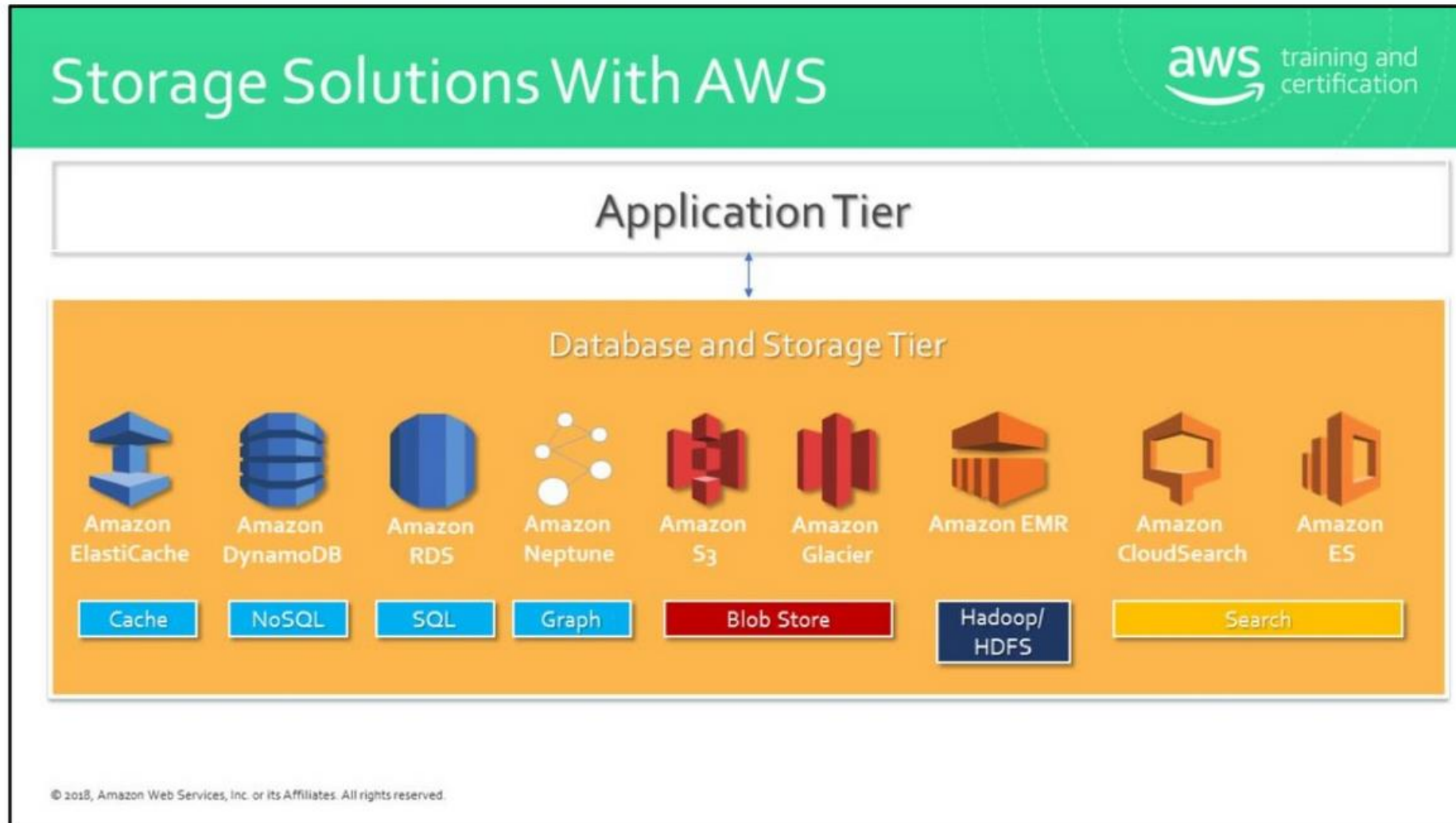
From AWS Solutions to the Big Data Pipeline



Other Tools That Can be Used With AWS



Storage Solutions With AWS



Outline

- Basic Concepts
- Hands-on & Review Questions
 - **Lab 1 - Using Amazon Athena to Analyze Log Data**
 - Lab 2 - Processing Server Logs with Hive on Amazon EMR
 - Lab 3 - Apache Spark on Amazon Web Services EMR

Lab 1 - Using Amazon Athena to Analyze Log Data

- The following log data is provided in an Amazon S3 bucket stored in the US East (North Virginia) Region.

```
2015-01-01T00:00:00.022719Z elb_demo_005 244.218.91.244:2255 172.36.231.239:443
0.000878 0.000803 0.000891 200 200 0 1886 "GET https://www.example.com/jobs/376
HTTP/1.1" "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko)
Chrome/47.0.2526.111 Safari/537.36" DHE-RSA-AES128-SHA TLSv1.2
```

- Raw data stored in plain text
 - a) Raw data stored in plain text
 - b) Compressed data in gzip format
 - c) Compressed and Partitioned data split into sub-directories
 - d) Columnar data, stored in Parquet format (compressed using Snappy compression)

Lab 1 - Results

TASKS	RUNTIME & DATA SCANNED
Task 1: Raw data stored in plain text	(Run time: 19.57 seconds, Data scanned: 90.72GB)
Task 2: Compressed data in gzip format	(Run time: 24.5 seconds, Data scanned: 13.08GB)
Task 3: Compressed and Partitioned data split into sub-directories	(Run time: 19.54 seconds, Data scanned: 432.46MB)
Task 4: Columnar data (Parquet with Snappy compression)	(Run time: 2.57 seconds, Data scanned: 18.44MB)

Lab 1 – Conclusions and Q&A

Q1. Name few benefits of Amazon Athena

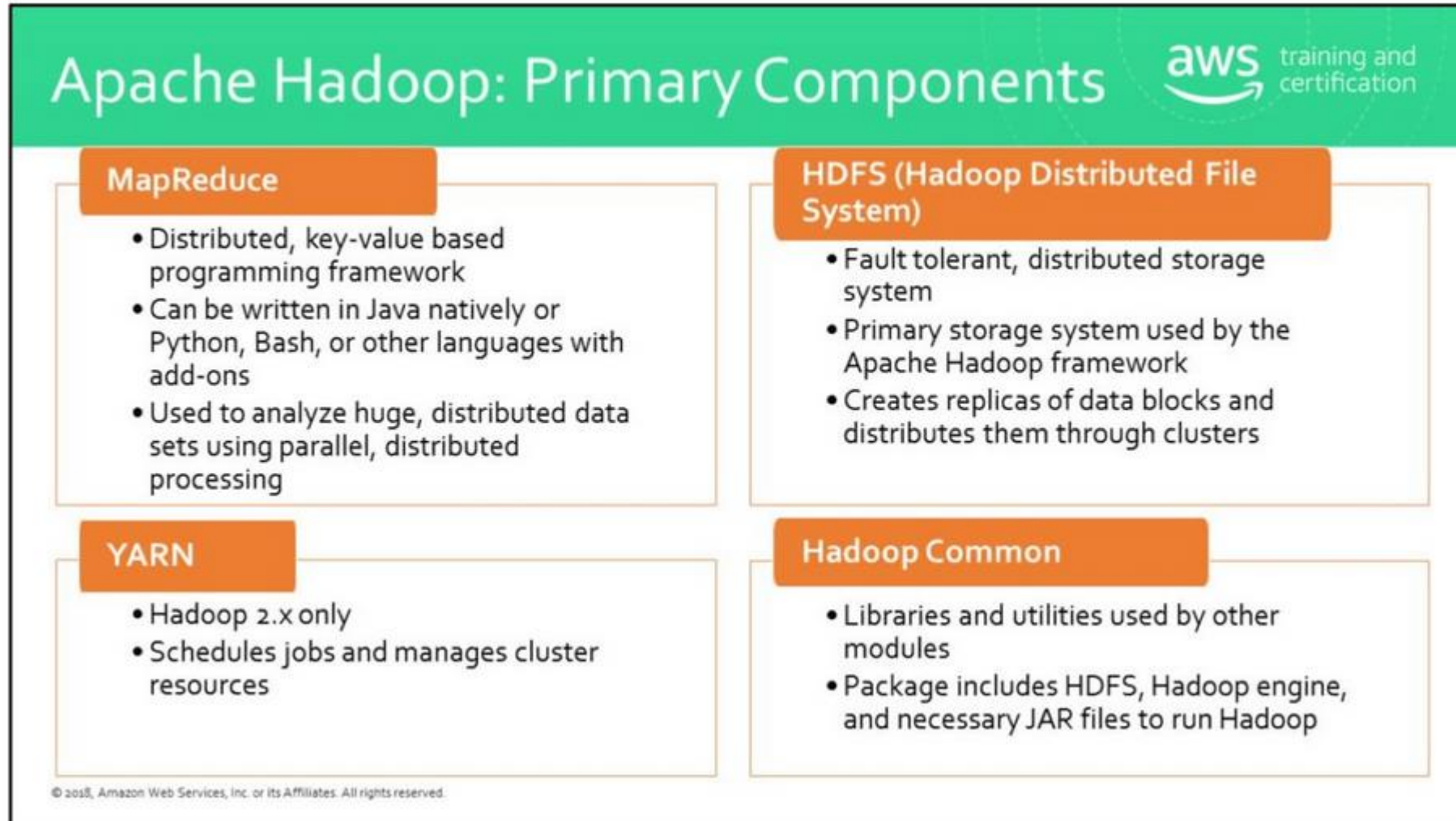
Q2. In your Lab 1, which data format improve the performance and save cost in Amazon Athena, and why?

Q3. Where does Amazon Athena store information and schemas about database?

Outline

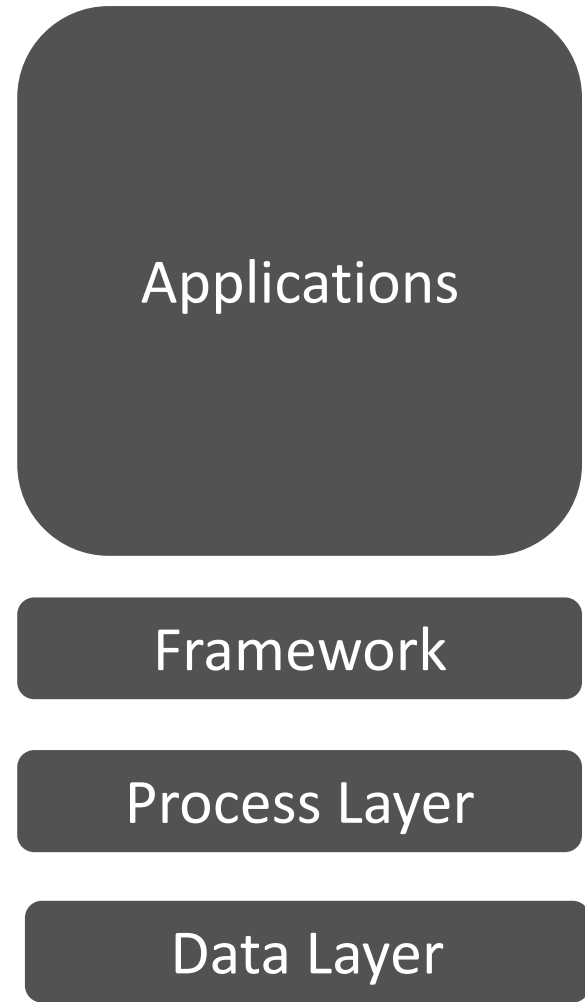
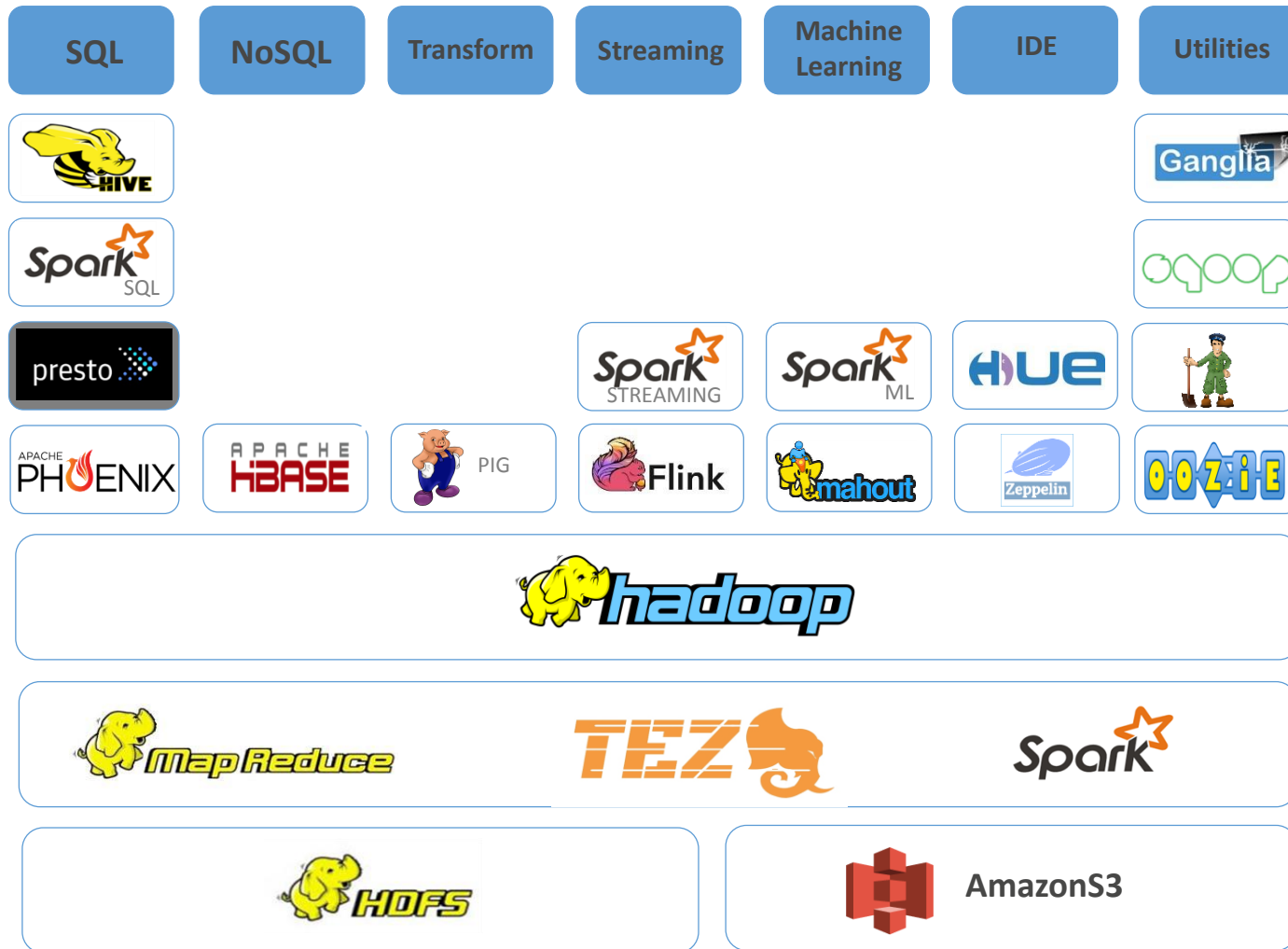
- Basic Concepts
- Hands-on & Review Questions
 - Lab 1 - Using Amazon Athena to Analyze Log Data
 - **Lab 2 - Processing Server Logs with Hive on Amazon EMR**
 - Lab 3 - Apache Spark on Amazon Web Services EMR

Apache Hadoop

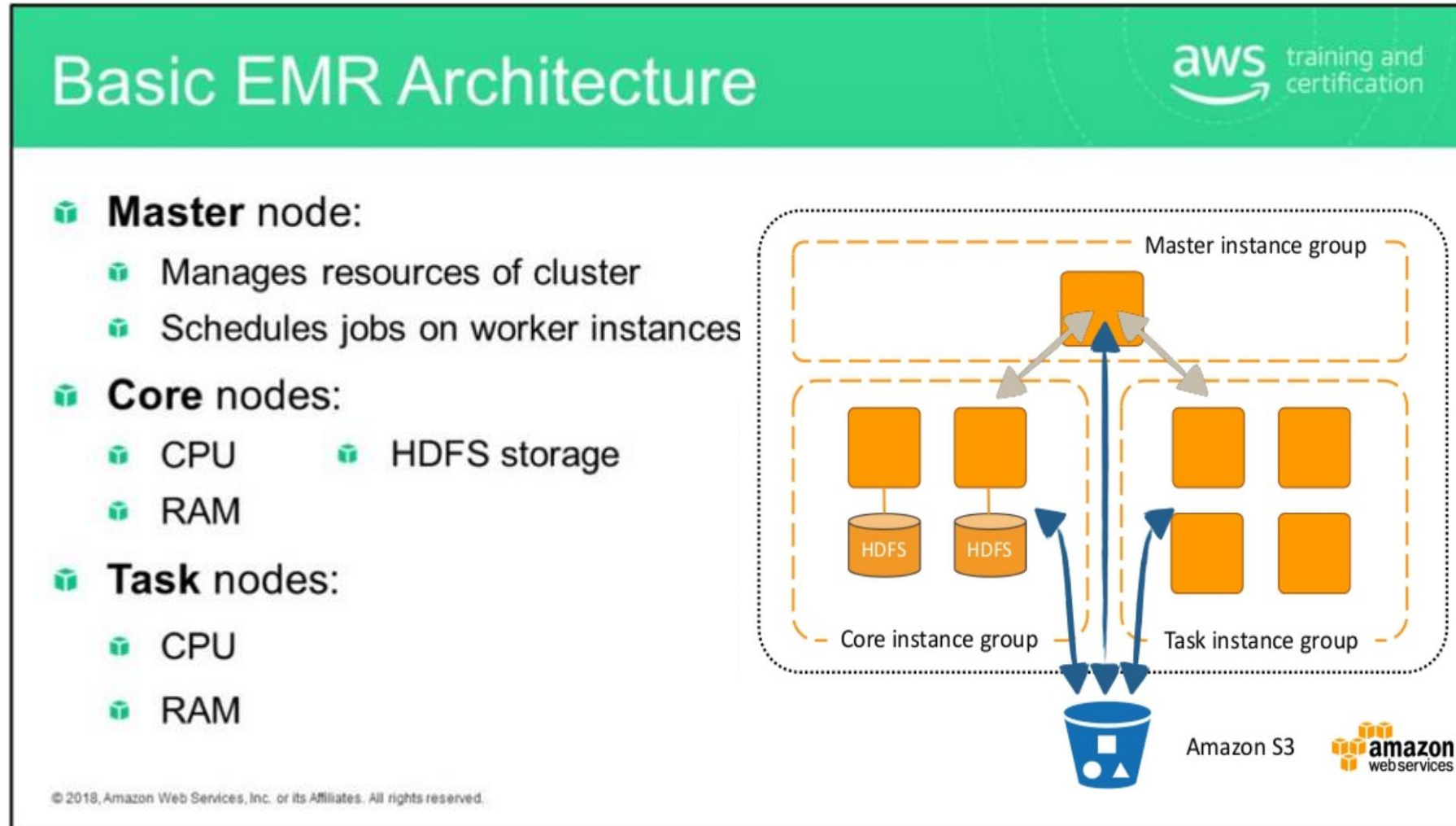


Amazon EMR

- Managed cluster platform
- Run Hadoop, Spark and other applications
- Easy – launch cluster in minutes!
- Use HDFS and S3 file systems
- Resize a running cluster
- Low cost –only pay for the resources you actually use!

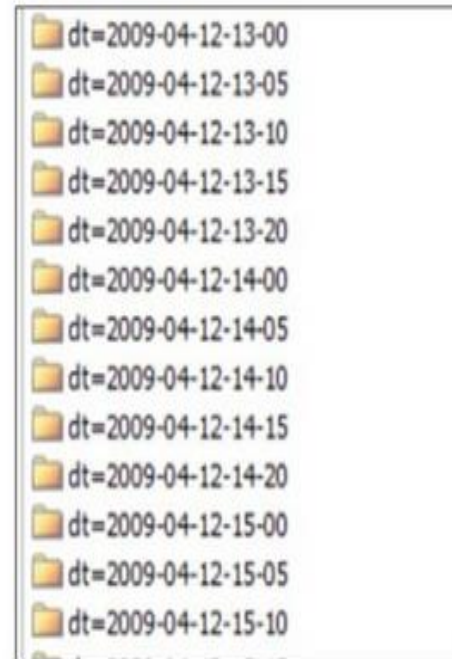


Amazon EMR Architecture



Lab 2 - Processing Server Logs with Hive on Amazon EMR

- Launch an Amazon EMR cluster with Apache Hive installed
- Use Hive commands to create external relation database tables from log data stored in Amazon S3
- Use Hive to query the tables you create and persist in Amazon S3.

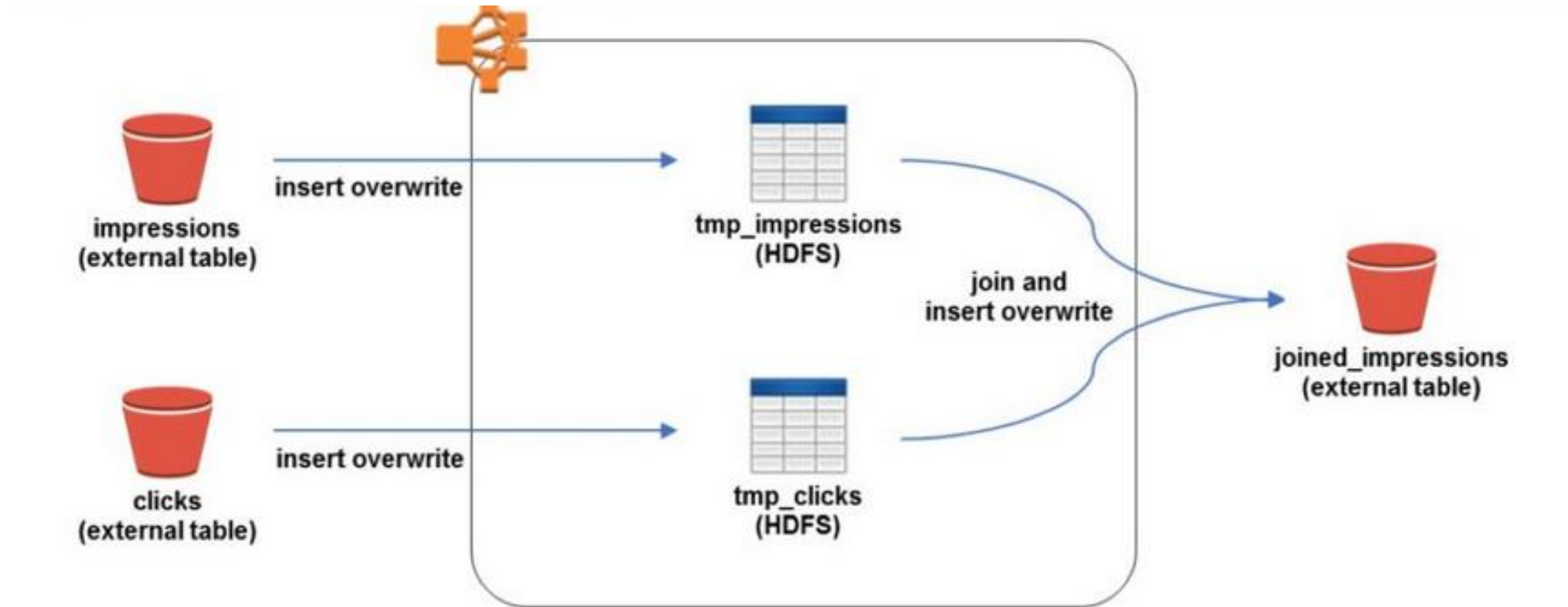


Stored log structure in Amazon S3

```
{
  "number": "95144", "referrer": "pdaipal.com", "processId": "1762", "adId": "94hLr4fOMx6u3HqK9xHdoXKqkatee", "browserCookie": "xmbatcmid", "userCookie": "xmbatcmid", "requestEndTime": "1239541871000", "impressionId": "2xP7V9YgDdGdcmFqPCwHSL9cV9Ln", "userAgent": "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; Trident/4.0)", "timers": {
    "modelLookup": "0.277", "requestTime": "0.8404", "threadId": "13", "ip": "20.143.145.26", "modelId": "hxxiuxduah", "hostname": "ec2-0-51-75-39.amazonaws.com", "sessionId": "UphjCmFIBt8KLaXZMnlp9lbXKPS", "requestBeginTime": "1239541870000"
  },
  "number": "97742", "referrer": "pdaipal.com", "processId": "1601", "adId": "nU35OND2:NEj6N4prfu9GK78v7dk12", "browserCookie": "hivvkrachd", "userCookie": "dhQFwSeAPdQdaFN95219aPv48N3BK", "requestEndTime": "1239541797000", "impressionId": "DLF6TfBpgDmGfRaxGCj89dtcE0nV", "userAgent": "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322; .NET CLR 2.0.50727; .NET CLR 3.0.04506.30; InfoPat)", "timers": {
    "modelLookup": "0.2844", "requestTime": "0.6756", "threadId": "51", "ip": "45.130.207.41", "modelId": "hxxiuxduah", "hostname": "ec2-0-51-75-39.amazonaws.com", "sessionId": "CKr06c8BpgLUMjg197W09689X6Jris", "requestBeginTime": "1239541794000"
  },
  "number": "102584", "referrer": "pandamio.com", "processId": "1646", "adId": "C4V3aHfPacH8pdc8R72H52wGFtJr2", "browserCookie": "hioibhniks", "userCookie": "wFioHtrVQki3nEJkMLLjoqC9oa", "requestEndTime": "1239541516000", "impressionId": "oadbpggHfijLkmdQmankCCX4nifwe", "userAgent": "Mozilla/5.0 (Windows; U; Windows NT 5.1; en-US) AppleWebKit/530.5 (KHTML, like Gecko) Chrome/2.0.172.37 Safari/530.5", "timers": {
    "modelLookup": "0.2726", "requestTime": "0.784", "threadId": "70", "ip": "51.140.241.78", "modelId": "hxxiuxduah", "hostname": "ec2-0-51-75-39.amazonaws.com", "sessionId": "HTCojSVeH6daXq8GaxjVG22H0osbc", "requestBeginTime": "1239541514000"
  }
}
```

Log data (sample)

Lab 2 - Processing Server Logs with Hive on Amazon EMR



Lab 2 - Review

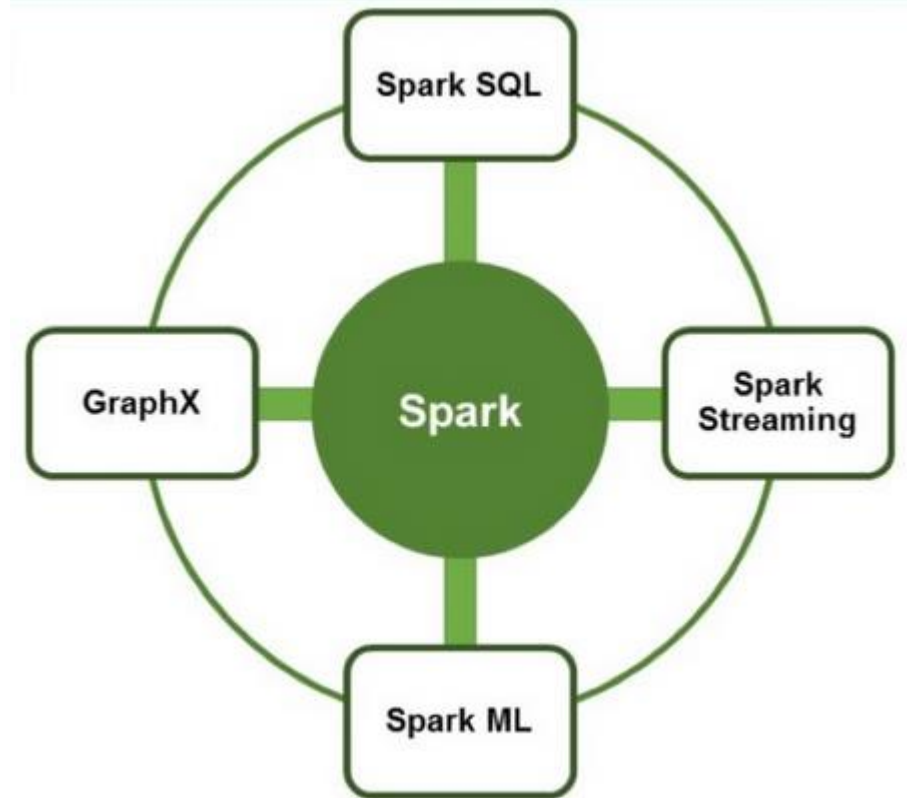
- What are the primary Apache Hadoop components?
- List applications for which Apache Hadoop is not suitable.
- What is the difference between a transient and a long-running cluster?
- List advantages of running Hive and Amazon EMR

Outline

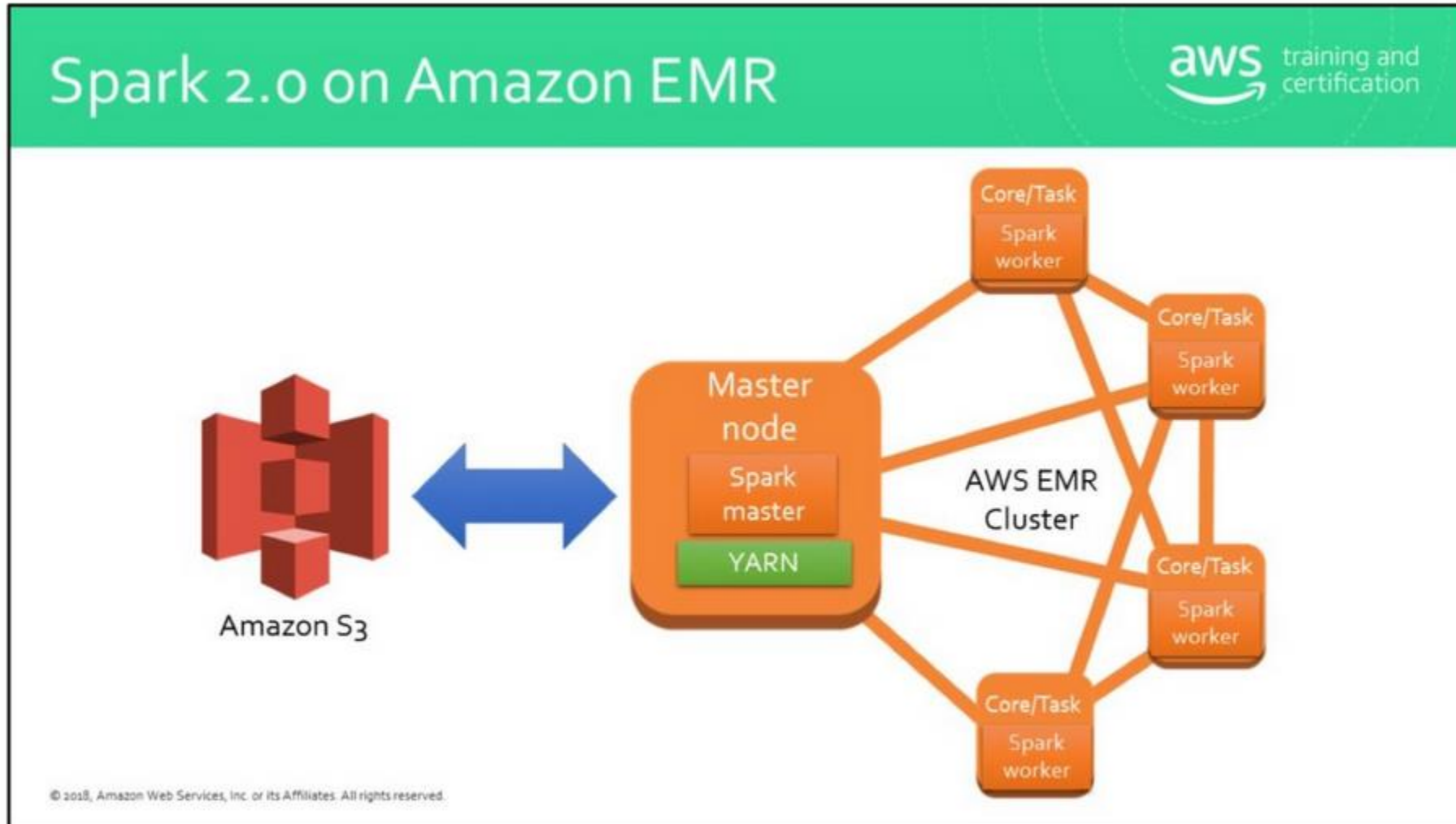
- Basic Concepts
- Hands-on & Review Questions
 - Lab 1 - Using Amazon Athena to Analyze Log Data
 - Lab 2 - Processing Server Logs with Hive on Amazon EMR
 - **Lab 3 - Apache Spark on Amazon Web Services EMR**

Apache Sparks

- Allows in-memory data mining and querying big datasets at fast speeds
- Resilient distributed datasets (RDD)
- Faster than MapReduce!
- Supports batch, interactive and streaming applications
- Can read and store data in HDFS, Amazon S3, and other databases.
- JDBC/ODBC connector compatible



Spark 2.0 on Amazon EMR



Lab 3 - Apache Spark on Amazon Web Services EMR

1. Create an AWS S3 storage bucket to hold data inputs, data outputs and logs
2. Create, Configure and Launch an Amazon EMR Cluster
3. Log in to the cluster master node
4. Run Spark-SQL Command Line Interface and Issue commands to create tables and run queries
5. Shut down cluster and remove any temporary resources

Lab 3 - Review

- SparkQL Exercise: Find the total trading volume (trade_size) for each stock before 12:00 noon
- What are the alternatives to load data to S3 from your local PC?

Data61 Research Big Data Cluster

