# Apache Spark and Object Stores —What you need to know
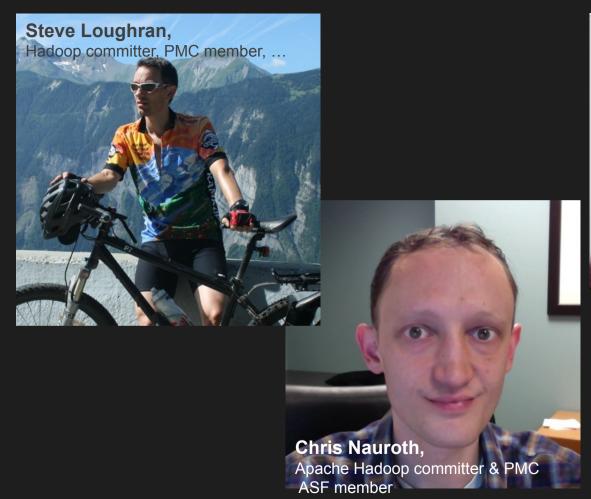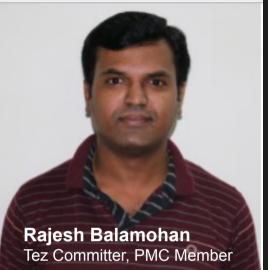
**Steve Loughran**
**stevel@apache.org**

October 2016
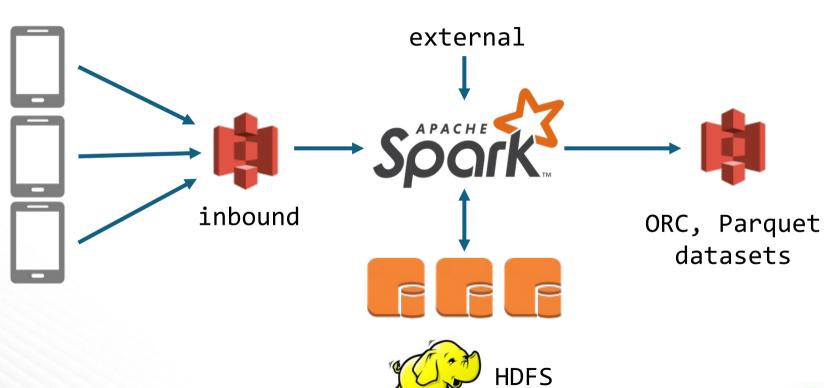
HORTONWORKS®
POWERING THE FUTURE OF DATA™

**Steve Loughran,**
Hadoop committer, PMC member, …

**Rajesh Balamohan**
Tez Committer, PMC Member

**Chris Nauroth,**
Apache Hadoop committer & PMC
ASF member

**Elastic ETL**

external

inbound

APACHE Spark™

ORC, Parquet datasets

HDFS

HORTONWORKS®

**Notebooks**



external

datasets

library

**HORTONWORKS®**

Streaming

# A Filesystem: Directories, Files → Data



`rename("/work/pending/part-01", "/work/complete")`

# Object Store: `hash(name)->blob`

```
hash("/work/pending/part-00")
 ["s01", "s02", "s04"]

hash("/work/pending/part-01")
 ["s02", "s03", "s04"]

copy("/work/pending/part-01",
     "/work/complete/part01")


delete("/work/pending/part-01")
```

# REST APIs

```
PUT /work/pending/part-01
... DATA ...

GET /work/pending/part-01
Content-Length: 1-8192

PUT /work/complete/part01
x-amz-copy-source: /work/pending/part-01

DELETE /work/pending/part-01

HEAD /work/complete/part-01

GET /?prefix=/work&delimiter=/
```

s01

00

s02

01

00

s03

01

s04

00

01

HORTONWORKS®

org.apache.hadoop.fs.FileSystem

hdfs   wasb   s3a   swift   adl   gs

# Four Challenges

1. Classpath
2. Credentials
3. Code
4. Commitment

Let's look At S3 and Azure

**Classpath: fix "No FileSystem for scheme: s3a"**

`hadoop-aws-2.7.x.jar`

```
aws-java-sdk-1.7.4.jar
joda-time-2.9.3.jar
(jackson-*-2.6.5.jar)
```

Get Spark with
Hadoop 2.7+ JARs

See SPARK-7481

**HORTONWORKS®**

**Credentials**

`core-site.xml` or `spark-default.conf`

> **`spark.hadoop.fs.s3a.access.key MY_ACCESS_KEY`**
>
> **`spark.hadoop.fs.s3a.secret.key MY_SECRET_KEY`**

`spark-submit` automatically propagates Environment Variables

> **`export AWS_ACCESS_KEY=MY_ACCESS_KEY`**
>
> **`export AWS_SECRET_KEY=MY_SECRET_KEY`**
>
> NEVER: share, check in to SCM, paste in bug reports…

**HORTONWORKS**®

# Authentication Failure: 403

```
com.amazonaws.services.s3.model.AmazonS3Exception:
 The request signature we calculated does not match
 the signature you provided.
 Check your key and signing method.
```

1. Check `joda-time.jar` & JVM version
2. Credentials wrong
3. Credentials not propagating
4. Local system clock (more likely on VMs)

**HORTONWORKS**®

# Code: just use the URL of the object store

```scala
val csvdata = spark.read.options(Map(
    "header" -> "true",
    "inferSchema" -> "true",
    "mode" -> "FAILFAST"))
  .csv("s3a://landsat-pds/scene_list.gz")
```

...read time O(distance)

HORTONWORKS®

# DataFrames

```
val landsat = "s3a://stevel-demo/landsat"
csvData.write.parquet(landsat)

val landsatOrc = "s3a://stevel-demo/landsatOrc"
csvData.write.orc(landsatOrc)

val df = spark.read.parquet(landsat)
val orcDf = spark.read.parquet(landsatOrc)
```

HORTONWORKS®

# Finding dirty data with Spark SQL

```
val sqlDF = spark.sql(
  "SELECT id, acquisitionDate, cloudCover"
 + s" FROM parquet.`${landsat}`")

val negativeClouds = sqlDF.filter("cloudCover < 0")
negativeClouds.show()
```

\* filter columns and data early
\* whether/when to cache()?
\* copy popular data to HDFS

HORTONWORKS®

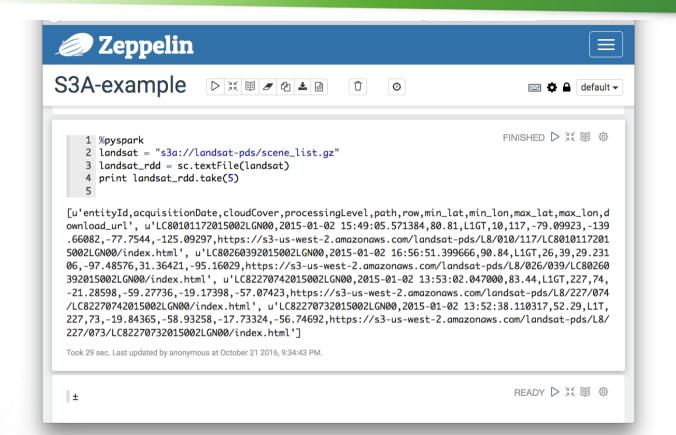# spark-default.conf

```
spark.sql.parquet.filterPushdown true
spark.sql.parquet.mergeSchema false
spark.hadoop.parquet.enable.summary-metadata false

spark.sql.orc.filterPushdown true
spark.sql.orc.splits.include.file.footer true
spark.sql.orc.cache.stripe.details.size 10000

spark.sql.hive.metastorePartitionPruning true
```

HORTONWORKS®

Notebooks? Classpath & Credentials

# The Commitment Problem

- rename() used for atomic commitment transaction

- time to copy() + delete() proportional to data * files

- S3: 6+ MB/s

- Azure: a lot faster —*usually*

```
spark.speculation false
spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version 2
spark.hadoop.mapreduce.fileoutputcommitter.cleanup.skipped true
```

HORTONWORKS®

# What about Direct Output Committers?



© Hortonworks Inc. 2011 – 2016. All Rights Reserved

# Recent S3A Performance (Hadoop 2.8, HDP 2.5, CDH 5.9 (?))

```
// forward seek by skipping stream
spark.hadoop.fs.s3a.readahead.range 157810688

// faster backward seek for ORC and Parquet input
spark.hadoop.fs.s3a.experimental.input.fadvise random

// PUT blocks in separate threads
spark.hadoop.fs.s3a.fast.output.enabled true
```

HORTONWORKS®

# Azure Storage: wasb://

A full substitute for HDFS

HORTONWORKS®
POWERING THE FUTURE OF DATA™

# Classpath: fix "No FileSystem for scheme: wasb"

`wasb://` : Consistent, with very fast rename (hence: commits)

`hadoop-azure-2.7.x.jar`

`azure-storage-2.2.0.jar`

`+ (jackson-core; http-components, hadoop-common)`

# Credentials: `core-site.xml` / `spark-default.conf`

```
<property>

<name>fs.azure.account.key.example.blob.core.windows.net</name>
 <value>0c0d44ac83ad7f94b0997b36e6e9a25b49a1394c</value>
</property>

spark.hadoop.fs.azure.account.key.example.blob.core.windows.net
  0c0d44ac83ad7f94b0997b36e6e9a25b49a1394c


wasb://demo@example.blob.core.windows.net
```

HORTONWORKS®

# Example: Azure Storage and Streaming

```scala
val streaming = new
StreamingContext(sparkConf,Seconds(10))
val azure = "wasb://demo@example.blob.core.windows.net/in"
val lines = streaming.textFileStream(azure)
val matches = lines.map(line => {
    println(line)
    line
  })
matches.print()
streaming.start()
```

\* PUT into the streaming directory

\* keep the dir clean

\* size window for slow scans

**HORTONWORKS**®

# Not Covered

- Partitioning/directory layout

- Infrastructure Throttling

- Optimal path names

- Error handling

- Metrics

**HORTONWORKS**®

# Summary

- Object Stores look just like any other URL

- …but do need classpath and configuration

- Issues: performance, commitment

- *Use Hadoop 2.7+ JARs*

- *Tune to reduce I/O*

- *Keep those credentials secret!*

HORTONWORKS®

APΔCHE:

# BIG_DATA

EUROPE

📅 November 14 - 16, 2016      📍 Melia Sevilla, Seville, Spain

#ApacheBigData

# Backup Slides

# Why run your code in the cloud?

No Upfront
HW Costs

Unlimited
Elastic Scale

Ephemeral &
Long-Running

IT &
Business Agility

**HORTONWORKS®**

History of Object Storage Support

2006 — s3:// –"inode on S3"

2008 — s3n:// "Native" S3 — s3:// Amazon EMR S3

swift:// OpenStack

2013 — s3a:// Replaces s3n

2014 — gs:// Google Cloud

s3a:// Stabilize

2015 — wasb:// Azure WASB

2016 — s3a:// Speed and consistency — adl:// Azure Data Lake — oss:// Aliyun

HORTONWORKS

# Often: Eventually Consistent



GET /work/pending/part-00
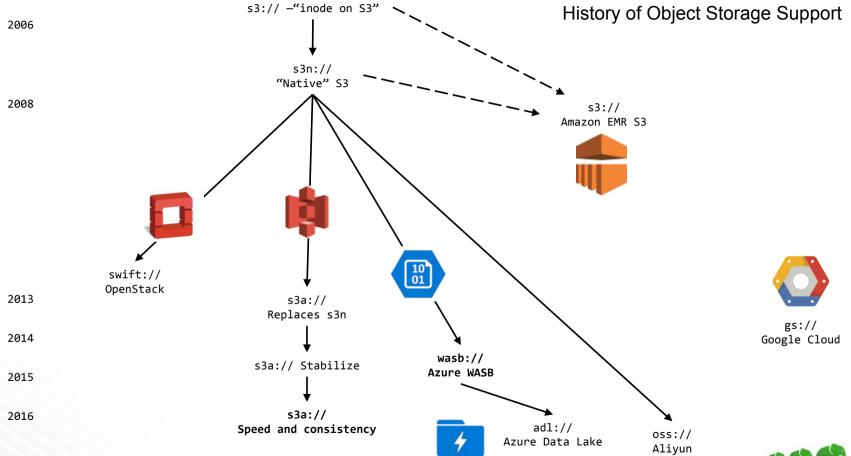
200

DELETE /work/pending/part-00

200

GET /work/pending/part-00

200
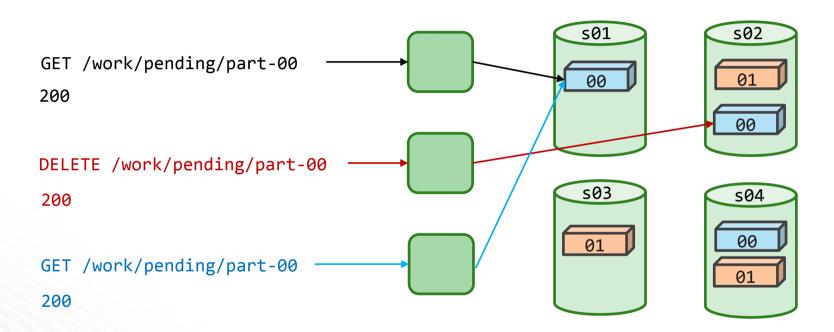
# S3 Server-Side Encryption

- Encryption of data at rest at S3

- Supports the SSE-S3 option: each object encrypted by a unique key using AES-256 cipher

- Now covered in S3A automated test suites

- Support for additional options under development (SSE-KMS and SSE-C)

**HORTONWORKS**®

# Next Steps for all Object Stores

- Output Committers
  - Logical commit operation decoupled from rename (non-atomic and costly in object stores)

- Object Store Abstraction Layer
  - Avoid impedance mismatch with FileSystem API
  - Provide specific APIs for better integration with object stores: saving, listing, copying

- Ongoing Performance Improvement

- Consistency

HORTONWORKS®

# Performance Considerations When Running Queries

- Splits Generation

  – File formats like ORC provides threadpool in split generation

- ORC Footer Cache

  – hive.orc.cache.stripe.details.size > 0

  – Caches footer details; Helps in reducing data reads during split generation

- Reduce S3A reads in Task side

  – hive.orc.splits.include.file.footer=true

  – Sends ORC footer information in splits payload.

  – Helps reducing the amount of data read in task side.

# Cloud Storage Connectors

| Azure | WASB | • Strongly consistent<br>• Good performance<br>• Well-tested on applications (incl. HBase) |
|---|---|---|
| | ADL | • Strongly consistent<br>• Tuned for big data analytics workloads |
| **Amazon Web Services** | S3A | • Eventually consistent - consistency work in progress by Hortonworks<br>• Performance improvements in progress<br>• Active development in Apache |
| | EMRFS | • Proprietary connector used in EMR<br>• Optional strong consistency for a cost |
| **Google Cloud Platform** | GCS | • Multiple configurable consistency policies<br>• Currently Google open source<br>• Good performance<br>• Could improve test coverage |

**HORTONWORKS**

| Scheme | Stable since | Speed | Consistency | Maintenance |
|---|---|---|---|---|
| s3n:// | Hadoop 1.0 | | | (Apache) |
| **s3a://** | **Hadoop 2.7** | **2.8+** | *ongoing* | **Apache** |
| swift:// | Hadoop 2.2 | | | Apache |
| **wasb://** | **Hadoop 2.7** | **Hadoop 2.7** | **strong** | **Apache** |
| adl:// | Hadoop 3 | | | |
| EMR s3:// | AWS EMR | | For a fee | Amazon |
| gs:// | ??? | | | Google @ github |

HORTONWORKS®

# Dependencies in Hadoop 2.8

**`hadoop-aws-2.8.x.jar`**

```
aws-java-sdk-core-1.10.6.jar
aws-java-sdk-kms-1.10.6.jar
aws-java-sdk-s3-1.10.6.jar
joda-time-2.9.3.jar
(jackson-*-2.6.5.jar)
```

**`hadoop-aws-2.8.x.jar`**

```
azure-storage-4.2.0.jar
```