

cloudera[®]

CDS 2 Powered by Apache Spark

Important Notice

© 2010-2018 Cloudera, Inc. All rights reserved.

Cloudera, the Cloudera logo, and any other product or service names or slogans contained in this document are trademarks of Cloudera and its suppliers or licensors, and may not be copied, imitated or used, in whole or in part, without the prior written permission of Cloudera or the applicable trademark holder. If this documentation includes code, including but not limited to, code examples, Cloudera makes this available to you under the terms of the Apache License, Version 2.0, including any required notices. A copy of the Apache License Version 2.0, including any notices, is included herein. A copy of the Apache License Version 2.0 can also be found here: <https://opensource.org/licenses/Apache-2.0>

Hadoop and the Hadoop elephant logo are trademarks of the Apache Software Foundation. All other trademarks, registered trademarks, product names and company names or logos mentioned in this document are the property of their respective owners. Reference to any products, services, processes or other information, by trade name, trademark, manufacturer, supplier or otherwise does not constitute or imply endorsement, sponsorship or recommendation thereof by us.

Complying with all applicable copyright laws is the responsibility of the user. Without limiting the rights under copyright, no part of this document may be reproduced, stored in or introduced into a retrieval system, or transmitted in any form or by any means (electronic, mechanical, photocopying, recording, or otherwise), or for any purpose, without the express written permission of Cloudera.

Cloudera may have patents, patent applications, trademarks, copyrights, or other intellectual property rights covering subject matter in this document. Except as expressly provided in any written license agreement from Cloudera, the furnishing of this document does not give you any license to these patents, trademarks copyrights, or other intellectual property. For information about patents covering Cloudera products, see <http://tiny.cloudera.com/patents>.

The information in this document is subject to change without notice. Cloudera shall not be liable for any damages resulting from technical errors or omissions which may be present in this document, or from use of this document.

Cloudera, Inc.

**395 Page Mill Road
Palo Alto, CA 94306
info@cloudera.com
US: 1-888-789-1488
Intl: 1-650-362-0488
www.cloudera.com**

Release Information

Version: CDS 2 Powered by Apache Spark
Date: November 2, 2018

Table of Contents

CDS 2 Powered by Apache Spark Overview.....	5
CDS 2 Powered by Apache Spark Release Notes.....	6
CDS Powered by Apache Spark Requirements.....	6
<i>CDH Versions.....</i>	6
<i>Cloudera Manager Versions.....</i>	6
<i>Scala 2.11 Requirement.....</i>	7
CDS Powered by Apache Spark New Features and Changes.....	7
<i>New Features.....</i>	7
CDS Powered by Apache Spark Known Issues.....	7
<i>Spark 2 Version Requirement for Clusters Managed by Cloudera Manager.....</i>	8
<i>Spark Standalone.....</i>	8
<i>HiveOnSpark is not Supported with Spark 2.....</i>	8
<i>SparkOnHBase is not Supported with Spark 2.....</i>	8
<i>Using the JDBC Datasource API to access Hive or Impala is not supported.....</i>	8
<i>Structured Streaming is not supported.....</i>	8
<i>Spark Streaming Direct Connector is not Supported.....</i>	8
<i>Oozie Spark2 Action is not Supported.....</i>	8
<i>SparkR is not Supported.....</i>	8
<i>GraphX is not Supported.....</i>	8
<i>Thrift Server.....</i>	9
<i>Spark SQL CLI is not Supported.....</i>	9
<i>Kudu is not Supported.....</i>	9
<i>Rolling Upgrades are not Supported.....</i>	9
<i>Package Install is not Supported.....</i>	9
<i>Spark Avro is not Supported.....</i>	9
<i>Accessing Multiple Clusters Simultaneously Not Supported.....</i>	9
<i>Hardware Acceleration for MLlib is not Supported.....</i>	9
<i>Long-running apps on a secure cluster might fail if driver is restarted.....</i>	9
Spark 2 Incompatible Changes.....	9
CDS Powered by Apache Spark Fixed Issues.....	10
<i>Issues Fixed in CDS 2.0 - Release 2.....</i>	10
<i>Issues Fixed in CDS 2.0 - Release 1.....</i>	10
CDS Powered by Apache Spark Version, Packaging, and Download Information.....	19
<i>CDS Versions Available for Download.....</i>	20
<i>CDS Maven Artifacts.....</i>	20
Using the CDS Powered by Apache Spark Maven Repository.....	20

<i>CDS 2.0 Powered by Apache Spark Maven Artifacts.....</i>	<i>21</i>
Installing CDS 2 Powered by Apache Spark.....	24
Administering CDS 2 Powered by Apache Spark.....	25
Configuring Spark 2 Tools as the Default.....	25
Running Applications with CDS Powered by Apache Spark.....	26
The Spark 2 Job Commands.....	26
Canary Test for pyspark2 Command.....	26
Fetching Spark 2 Maven Dependencies.....	26
Adapting the Spark WordCount App for Spark 2.....	27
Accessing the Spark 2 History Server.....	27
Troubleshooting for Spark 2.....	28
Frequently Asked Questions about CDS Powered by Apache Spark.....	29
Appendix: Apache License, Version 2.0.....	30

CDS 2 Powered by Apache Spark Overview

**Note:**

This Spark 2.0 documentation refers to the second release of CDS 2 Powered by Apache Spark. This component is generally available and is now supported on CDH 5.7 through CDH 5.10.

The latest release (release 2) addresses a Hive compatibility issue that affects CDH 5.10.1 and higher, CDH 5.9.2 and higher, CDH 5.8.5 and higher, and CDH 5.7.6 and higher. If you are using one of these CDH versions, you must upgrade to the Spark 2.0 release 2 parcel to avoid Spark 2 job failures when using Hive functionality.

[Apache Spark](#) is a general framework for distributed computing that offers high performance for both batch and interactive processing. It exposes APIs for Java, Python, and Scala.

For detailed API information, see the [Apache Spark project site](#).



Note: Although this document makes some references to the external Spark site, not all the features, components, recommendations, and so on are applicable to Spark when used on CDH. Always cross-check the Cloudera documentation before building a reliance on some aspect of Spark that might not be supported or recommended by Cloudera. In particular, see [CDS Powered by Apache Spark Known Issues](#) on page 7 for components and features to avoid.

CDS 2 Powered by Apache Spark consists of Spark core and several related projects:

[Spark SQL](#)

Module for working with structured data. Allows you to seamlessly mix SQL queries with Spark programs.

[Spark Streaming](#)

API that allows you to build scalable fault-tolerant streaming applications.

[MLlib](#)

API that implements common machine learning algorithms.

Cloudera distributes two versions of Apache Spark: 1.6 and 2.0.

Spark 1.6 is distributed as part of Cloudera Enterprise 5.7.x and higher, whose documentation is available at [Cloudera Enterprise 5.7.x Documentation](#).

This document describes the separately released CDS 2.0.

A Spark 1.6 service can coexist with a Spark 2.0 service. The configurations of the two services do not conflict and both services use the same YARN service. The port of the Spark History Server is 18088 for Spark 1.6 and 18089 for Spark 2.0.

Unsupported Features

Consult [CDS Powered by Apache Spark Known Issues](#) on page 7 for a comprehensive list of features that are not supported with CDS 2 Powered by Apache Spark.

Related Information

- [Cloudera Spark forum](#)
- [Apache Spark documentation](#)

CDS 2 Powered by Apache Spark Release Notes

The release notes provide information on requirements, new and changed features, known issues, and fixed issues, and version and packaging information for CDS 2 Powered by Apache Spark.

CDS Powered by Apache Spark Requirements

The following sections describe software requirements for CDS Powered by Apache Spark.

CDH Versions



Important: CDS Powered by Apache Spark is available in parcel format only, and not packages. Because Cloudera Manager does not support using parcels and packages in the same cluster, you cannot use CDS if you are using a package-based installation of CDH.

The CDS parcel version displayed in Cloudera Manager, which is also part of the [parcel file name](#), is structured as follows:

`<CDS_version>-1.<cdh_build_version>.p<patch_version>.<build_number>`

The `<cdh_build_version>` portion is the version of CDH upon which the release was built. It is *not* the minimum supported CDH version.

Supported versions of CDH are described below.

The latest release (release 2) addresses a Hive compatibility issue that affects CDH 5.10.1 and higher, CDH 5.9.2 and higher, CDH 5.8.5 and higher, and CDH 5.7.6 and higher. If you are using one of these CDH versions, you must upgrade to the Spark 2.0 release 2 parcel to avoid Spark 2 job failures when using Hive functionality.

CDS 2 Powered by Apache Spark Version	CDH Version
2.0 Release 2	CDH 5.7, CDH 5.8, CDH 5.9, CDH 5.10
2.0 Release 1	CDH 5.7 up to 5.7.5, CDH 5.8 up to 5.8.4, CDH 5.9 up to 5.9.1, CDH 5.10.0. Spark 2.0 Release 2 is required for any higher maintenance releases in any of these CDH versions.

A Spark 1.6 service can co-exist on the same cluster as Spark 2. The two services are configured to not conflict and run on the same YARN cluster. Spark 2 uses the [external shuffle service](#) from the CDH installation if Spark 1 is already installed, or installs the shuffle service itself if necessary. Only the external shuffle service classes from the CDH installation can be used.

Cloudera Manager Versions



Important:

Because CDS Powered by Apache Spark is only installable using the parcel mechanism, it can only be used on clusters managed by Cloudera Manager. Additionally, because Cloudera Manager does not support using parcels and packages in the same cluster, you cannot use CDS if you are using a package-based installation of CDH.

Applicable versions of Cloudera Manager for Spark 2 are described below.

CDS 2 Powered by Apache Spark Version	Cloudera Manager Version
2.0 Release 2	Cloudera Manager 5.8.3, 5.9 and higher
2.0 Release 1	Cloudera Manager 5.8.3, 5.9 and higher

Scala 2.11 Requirement

Spark 2 does not work with Scala 2.10. Use Scala 2.11 only.

CDS Powered by Apache Spark New Features and Changes

The following sections describe what's new and changed in each CDS Powered by Apache Spark release.

New Features

The following sections describe what's new in each CDS Powered by Apache Spark release.

What's New in CDS 2.0 Release 2

- The latest release (release 2) addresses a Hive compatibility issue that affects CDH 5.10.1 and higher, CDH 5.9.2 and higher, CDH 5.8.5 and higher, and CDH 5.7.6 and higher. If you are using one of these CDH versions, you must upgrade to the Spark 2.0 release 2 parcel to avoid Spark 2 job failures when using Hive functionality.

What's New in CDS 2.0 Release 1

- New `SparkSession` object replaces `HiveContext` and `SQLContext`.
 - Most of the Hive logic has been reimplemented in Spark.
 - Some Hive dependencies still exist:
 - SerDe support.
 - UDF support.
- Added support for the unified Dataset API.
- Faster Spark SQL achieved with whole stage code generation.
- More complete SQL syntax now supports subqueries.
- Adds the `spark-csv` library.
- Backport of SPARK-5847. The root for metrics is now the app name (`spark.app.name`) instead of the app ID. The app ID requires investigation to match to the app name, and changes when streaming jobs are stopped and restarted.

CDS Powered by Apache Spark Known Issues

The following sections describe the current known issues and limitations in CDS Powered by Apache Spark. In some cases, a feature from the upstream Apache Spark project is currently not considered reliable enough to be supported by Cloudera. For a number of integration features in CDH that rely on Spark, the feature does not work with CDS 2 Powered by Apache Spark because CDH components are not introducing dependencies on Spark 2.

Spark 2 Version Requirement for Clusters Managed by Cloudera Manager

**Important:**

Because CDS Powered by Apache Spark is only installable using the parcel mechanism, it can only be used on clusters managed by Cloudera Manager. Additionally, because Cloudera Manager does not support using parcels and packages in the same cluster, you cannot use CDS if you are using a package-based installation of CDH.

All CDH clusters managed by a single Cloudera Manager instance must use exactly the same version of CDS Powered By Apache Spark. Make sure to install or upgrade the CSDs and parcels across all machines of all clusters at the same time.

Spark Standalone

Spark Standalone is not supported for Spark 2.

HiveOnSpark is not Supported with Spark 2

The HiveOnSpark module is a CDH 5 component that has a dependency on Apache Spark 1.6. Because CDH 5 components do not have any dependencies on Spark 2, the HiveOnSpark module does not work with CDS Powered by Apache Spark. You can still use Spark 2 with Hive using other methods.

SparkOnHBase is not Supported with Spark 2

The SparkOnHBase module is a CDH 5 component that has a dependency on Apache Spark 1.6. Because CDH 5 components do not have any dependencies on Spark 2, the SparkOnHBase module does not work with CDS Powered by Apache Spark. You can still use Spark 2 with HBase using other methods.

Using the JDBC Datasource API to access Hive or Impala is not supported

Structured Streaming is not supported

Cloudera does not support the Structured Streaming API because it is an experimental API.

Spark Streaming Direct Connector is not Supported

Spark 2 does not support the Spark Streaming direct connector that uses the new Kafka consumer API, available starting Apache Kafka 0.9 (Cloudera Kafka 2.0) for secure clusters. Therefore, you cannot use Spark 2 to read data from Kafka using the new direct connector. Consequently, you cannot read data from a secure cluster that uses Kerberos with Spark Streaming. You can still use the older Spark Streaming direct connector, which uses the old Kafka consumer API, to read data from Kafka in a non-secure cluster.

Oozie Spark2 Action is not Supported

The Oozie Spark action is a CDH component that has a dependency on Spark 1.6. Because CDH components do not have any dependencies on Spark 2, the Oozie Spark action does not work with Spark 2.

SparkR is not Supported

SparkR is not supported for Spark 2. (SparkR is also not supported in CDH with Spark 1.6.)

GraphX is not Supported

GraphX is not supported for Spark 2. (GraphX is also not supported in CDH with Spark 1.6.)

Thrift Server

The Thrift JDBC/ODBC server is not supported for Spark 2. (The Thrift server is also not supported in CDH with Spark 1.6.)

Spark SQL CLI is not Supported

The Spark SQL CLI is not supported for Spark 2. (The Spark SQL CLI is also not supported in CDH with Spark 1.6.)

Kudu is not Supported

The Kudu integration for Spark only works with Spark 1.6.

Rolling Upgrades are not Supported

Rolling upgrades are not possible from Spark 1.6 bundled with CDH, to CDS 2 Powered by Apache Spark.

Package Install is not Supported

CDS 2 Powered by Apache Spark is only installable as a parcel.

Spark Avro is not Supported

The `spark-avro` library is not integrated into the Spark 2.0 parcel.

Accessing Multiple Clusters Simultaneously Not Supported

Spark does not support accessing multiple clusters in the same application.

Hardware Acceleration for MLlib is not Supported

This feature, part of the GPL Extras package for CDH, is not supported with the CDS Powered By Apache Spark 2. This feature is supported for Spark 1.6.

Long-running apps on a secure cluster might fail if driver is restarted

If you submit a long-running app on a secure cluster using the `--principal` and `--keytab` options in cluster mode, and a failure causes the driver to restart after 7 days (the default maximum HDFS delegation token lifetime), the new driver fails with an error similar to the following:

```
Exception in thread "main"
org.apache.hadoop.ipc.RemoteException(org.apache.hadoop.security.token.SecretManager$InvalidToken):
token <token_info> can't be found in cache
```

Workaround: None

Affected Versions: All CDS 2.0, 2.1, and 2.2 releases

Fixed Versions: CDS 2.3 Release 2

Apache Issue: [SPARK-23361](#)

Cloudera Issue: CDH-64865

Spark 2 Incompatible Changes

The following sections describe changes in CDS 2 Powered by Apache Spark that might require special handling during upgrades, or code changes within existing applications.

Incompatible Changes in CDS 2.0

- `HiveContext` and `SQLContext` have been removed.
- `DataFrames` have been removed from the Scala API. `DataFrame` is now a special case of `Dataset`.
- Spark 2.0 and higher do not use an assembly JAR for standalone applications.

CDS Powered by Apache Spark Fixed Issues

The following sections describe the issues fixed in each CDS Powered by Apache Spark release.

Issues Fixed in CDS 2.0 - Release 2

- [\[SPARK-4563\]](#)[CORE] Allow driver to advertise a different network address.
- [\[SPARK-18993\]](#) Unable to build/compile Spark in IntelliJ due to missing Scala deps in spark-tags
- [\[SPARK-19314\]](#) Do not allow sort before aggregation in Structured Streaming plan
- [\[SPARK-18762\]](#) Web UI should be http:4040 instead of https:4040
- [\[SPARK-18745\]](#) java.lang.IndexOutOfBoundsException running query 68 Spark SQL on (100TB)
- [\[SPARK-18703\]](#) Insertion/CTAS against Hive Tables: Staging Directories and Data Files Not Dropped Until Normal Termination of JVM
- [\[SPARK-18091\]](#) Deep if expressions cause Generated SpecificUnsafeProjection code to exceed JVM code size limit

Issues Fixed in CDS 2.0 - Release 1

- [\[SPARK-4563\]](#)[CORE] Allow driver to advertise a different network address.
- [\[SPARK-18685\]](#)[TESTS] Fix URI and release resources after opening in tests at ExecutorClassLoaderSuite
- [\[SPARK-18677\]](#) Fix parsing ['key'] in JSON path expressions.
- [\[SPARK-18617\]](#)[SPARK-18560][TESTS] Fix flaky test: StreamingContextSuite. Receiver data should be deserialized properly
- [\[SPARK-18617\]](#)[SPARK-18560][TEST] Fix flaky test: StreamingContextSuite. Receiver data should be deserialized properly
- [\[SPARK-18274\]](#)[ML][PYSPARK] Memory leak in PySpark JavaWrapper
- [\[SPARK-18674\]](#)[SQL] improve the error message of using join
- [\[SPARK-18617\]](#)[CORE][STREAMING] Close "kryo auto pick" feature for Spark Streaming
- [\[SPARK-17843\]](#)[WEB UI] Indicate event logs pending for processing on h...
- [\[SPARK-17783\]](#)[SQL][BACKPORT-2.0] Hide Credentials in CREATE and DESC FORMATTED/EXTENDED a PERSISTENT/TEMP Table for JDBC
- [\[SPARK-18640\]](#) Add synchronization to TaskScheduler.runningTasksByExecutors
- [\[SPARK-18553\]](#)[CORE] Fix leak of TaskSetManager following executor loss
- [\[SPARK-18597\]](#)[SQL] Do not push-down join conditions to the left side of a Left Anti join [BRANCH-2.0]
- [\[SPARK-18118\]](#)[SQL] fix a compilation error due to nested JavaBeans
- [\[SPARK-17251\]](#)[SQL] Improve `OuterReference` to be `NamedExpression`
- [\[SPARK-18436\]](#)[SQL] isin causing SQL syntax error with JDBC
- [\[SPARK-18519\]](#)[SQL][BRANCH-2.0] map type can not be used in EqualTo
- [\[SPARK-18053\]](#)[SQL] compare unsafe and safe complex-type values correctly
- [\[SPARK-18504\]](#)[SQL] Scalar subquery with extra group by columns returning incorrect result
- [\[SPARK-18477\]](#)[SS] Enable interrupts for HDFS in HDFSMetadataLog
- [\[SPARK-18546\]](#)[CORE] Fix merging shuffle spills when using encryption.
- [\[SPARK-18547\]](#)[CORE] Propagate I/O encryption key when executors register.
- [\[SPARK-16625\]](#)[SQL] General data types to be mapped to Oracle
- [\[SPARK-18462\]](#) Fix ClassCastException in SparkListenerDriverAccumUpdates event

- [\[SPARK-18459\]](#)[SPARK-18460][STRUCTUREDSTREAMING] Rename triggerId to batchId and add triggerDetails to json in StreamingQueryStatus (for branch-2.0)
- [\[SPARK-18430\]](#)[SQL][BACKPORT-2.0] Fixed Exception Messages when Hitting an Invocation Exception of Function Lookup
- [\[SPARK-18400\]](#)[STREAMING] NPE when resharding Kinesis Stream
- [\[SPARK-18300\]](#)[SQL] Do not apply foldable propagation with expand as a child [BRANCH-2.0]
- [\[SPARK-18337\]](#) Complete mode memory sinks should be able to recover from checkpoints
- [\[SPARK-16808\]](#)[CORE] History Server main page does not honor APPLICATION_WEB_PROXY_BASE
- [\[SPARK-17348\]](#)[SQL] Incorrect results from subquery transformation
- [\[SPARK-18416\]](#)[STRUCTURED STREAMING] Fixed temp file leak in state store
- [\[SPARK-18432\]](#)[DOC] Changed HDFS default block size from 64MB to 128MB
- [\[SPARK-18010\]](#)[CORE] Reduce work performed for building up the application list for the History Server app list UI page
- [\[SPARK-18382\]](#)[WEBUI] "run at null:-1" in UI when no file/line info in call site info
- [\[SPARK-18426\]](#)[STRUCTURED STREAMING] Python Documentation Fix for Structured Streaming Programming Guide
- [\[SPARK-17982\]](#)[SQL][BACKPORT-2.0] SQLBuilder should wrap the generated SQL with parenthesis for LIMIT
- [\[SPARK-18387\]](#)[SQL] Add serialization to checkEvaluation.
- [\[SPARK-18368\]](#)[SQL] Fix regexp replace when serialized
- [\[SPARK-18342\]](#) Make rename failures fatal in HDFSBackedStateStore
- [\[SPARK-18280\]](#)[CORE] Fix potential deadlock in `StandaloneSchedulerBackend.dead`
- [\[SPARK-17703\]](#)[SQL][BACKPORT-2.0] Add unnamed version of addReferenceObj for minor objects.
- [\[SPARK-18137\]](#)[SQL] Fix RewriteDistinctAggregates UnresolvedException when a UDAF has a foldable TypeCheck
- [\[SPARK-18283\]](#)[STRUCTURED STREAMING][KAFKA] Added test to check whether default starting offset in latest
- [\[SPARK-18125\]](#)[SQL][BRANCH-2.0] Fix a compilation error in codegen due to splitExpression
- [\[SPARK-17849\]](#)[SQL] Fix NPE problem when using grouping sets
- [\[SPARK-17693\]](#)[SQL][BACKPORT-2.0] Fixed Insert Failure To Data Source Tables when the Schema has the Comment Field
- [\[SPARK-17981\]](#)[SPARK-17957][SQL][BACKPORT-2.0] Fix Incorrect Nullability Setting to False in FilterExec
- [\[SPARK-18189\]](#)[SQL][FOLLOWUP] Move test from RepITSuite to prevent java.lang.ClassCircularityError
- [\[SPARK-17337\]](#)[SPARK-16804][SQL][BRANCH-2.0] Backport subquery related PRs
- [\[SPARK-18200\]](#)[GRAPHX][FOLLOW-UP] Support zero as an initial capacity in OpenHashSet
- [\[SPARK-18200\]](#)[GRAPHX] Support zero as an initial capacity in OpenHashSet
- [\[SPARK-18111\]](#)[SQL] Wrong approximate quantile answer when multiple records have the minimum value(for branch 2.0)
- [\[SPARK-18160\]](#)[CORE][YARN] spark.files & spark.jars should not be passed to driver in yarn mode
- [\[SPARK-16796\]](#)[WEB UI] Mask spark.authenticate.secret on Spark environ...
- [\[SPARK-18133\]](#)[BRANCH-2.0][EXAMPLES][ML] Python ML Pipeline Examl...
- [\[SPARK-18144\]](#)[SQL] logging StreamingQueryListener\$QueryStartedEvent
- [\[SPARK-18114\]](#)[HOTFIX] Fix line-too-long style error from backport of SPARK-18114
- [\[SPARK-18148\]](#)[SQL] Misleading Error Message for Aggregation Without Window/GroupBy
- [\[SPARK-18189\]](#)[SQL] Fix serialization issue in KeyValueGroupedDataset
- [\[SPARK-18114\]](#)[MESOS] Fix mesos cluster scheduler generage command option error
- [\[SPARK-18030\]](#)[TESTS] Fix flaky FileStreamSourceSuite by not deleting the files
- [\[SPARK-18143\]](#)[SQL] Ignore Structured Streaming event logs to avoid breaking history server (branch 2.0)
- [\[SPARK-16312\]](#)[FOLLOW-UP][STREAMING][KAFKA][DOC] Add java code snippet for Kafka 0.10 integration doc
- [\[SPARK-18164\]](#)[SQL] ForeachSink should fail the Spark job if `process` throws exception
- [\[SPARK-16963\]](#)[SQL] Fix test "StreamExecution metadata garbage collection"
- [\[SPARK-17813\]](#)[SQL][KAFKA] Maximum data per trigger
- [\[SPARK-18132\]](#) Fix checkstyle

- [\[SPARK-18009\]](#)[SQL] Fix ClassCastException while calling toLocalIterator() on dataframe produced by RunnableCommand
- [\[SPARK-16963\]](#)[STREAMING][SQL] Changes to Source trait and related implementation classes
- [\[SPARK-13747\]](#)[SQL] Fix concurrent executions in ForkJoinPool for SQL (branch 2.0)
- [\[SPARK-18063\]](#)[SQL] Failed to infer constraints over multiple aliases
- [\[SPARK-16304\]](#) LinkageError should not crash Spark executor
- [\[SPARK-17733\]](#)[SQL] InferFiltersFromConstraints rule never terminates for query
- [\[SPARK-18022\]](#)[SQL] java.lang.NullPointerException instead of real exception when saving DF to MySQL
- [\[SPARK-16988\]](#)[SPARK SHELL] spark history server log needs to be fixed to show https url when ssl is enabled
- [\[SPARK-18070\]](#)[SQL] binary operator should not consider nullability when comparing input types
- [\[SPARK-17624\]](#)[SQL][STREAMING][TEST] Fixed flaky StateStoreSuite.maintenance
- [\[SPARK-18044\]](#)[STREAMING] FileStreamSource should not infer partitions in every batch
- [\[SPARK-17153\]](#)[SQL] Should read partition data when reading new files in filestream without globbing
- [\[SPARK-18093\]](#)[SQL] Fix default value test in SQLConfSuite to work rega...
- [\[SPARK-17810\]](#)[SQL] Default spark.sql.warehouse.dir is relative to local FS but can resolve as HDFS path
- [\[SPARK-18058\]](#)[SQL] [BRANCH-2.0]Comparing column types ignoring Nullability in Union and SetOperation
- [\[SPARK-17123\]](#)[SQL][BRANCH-2.0] Use type-widened encoder for DataFrame for set operations
- [\[SPARK-17698\]](#)[SQL] Join predicates should not contain filter clauses
- [\[SPARK-17986\]](#)[ML] SQLTransformer should remove temporary tables
- [\[SPARK-16606\]](#)[MINOR] Tiny follow-up to , to correct more instances of the same log message typo
- [\[SPARK-17853\]](#)[STREAMING][KAFKA][DOC] make it clear that reusing group.id is bad
- [\[SPARK-16312\]](#)[STREAMING][KAFKA][DOC] Doc for Kafka 0.10 integration
- [\[SPARK-17812\]](#)[SQL][KAFKA] Assign and specific startingOffsets for structured stream
- [\[SPARK-17929\]](#)[CORE] Fix deadlock when CoarseGrainedSchedulerBackend reset
- [\[SPARK-17926\]](#)[SQL][STREAMING] Added json for statuses
- [\[SPARK-17811\]](#) SparkR cannot parallelize data.frame with NA or NULL in Date columns
- [\[SPARK-18034\]](#) Upgrade to MiMa 0.1.11 to fix flakiness
- [\[SPARK-17999\]](#)[KAFKA][SQL] Add getPreferredLocations for KafkaSourceRDD
- [\[SPARK-18003\]](#)[SPARK CORE] Fix bug of RDD zipWithIndex & zipWithUniqueId index value overflowing
- [\[SPARK-17989\]](#)[SQL] Check ascendingOrder type in sort_array function rather than throwing ClassCastException
- [\[SPARK-17675\]](#)[CORE] Expand Blacklist for TaskSets
- [\[SPARK-17623\]](#)[CORE] Clarify type of TaskEndReason with a failed task.
- [\[SPARK-17304\]](#) Fix perf. issue caused by TaskSetManager.abortIfCompletelyBlacklisted
- [\[SPARK-15865\]](#)[CORE] Blacklist should not result in job hanging with less than 4 executors
- [\[SPARK-15783\]](#)[CORE] Fix Flakiness in BlacklistIntegrationSuite
- [\[SPARK-15783\]](#)[CORE] still some flakiness in these blacklist tests so ignore for now
- [\[SPARK-15714\]](#)[CORE] Fix flaky o.a.s.scheduler.BlacklistIntegrationSuite
- [\[SPARK-10372\]](#) [CORE] basic test framework for entire spark scheduler
- [\[SPARK-16106\]](#)[CORE] TaskSchedulerImpl should properly track executors added to existing hosts
- [\[SPARK-18001\]](#)[DOCUMENT] fix broke link to SparkDataFrame
- [\[SPARK-17711\]](#)[TEST-HADOOP2.2] Fix hadoop2.2 compilation error
- [\[SPARK-17731\]](#)[SQL][STREAMING][FOLLOWUP] Refactored StreamingQueryListener APIs for branch-2.0
- [\[SPARK-17841\]](#)[STREAMING][KAFKA] drain commitQueue
- [\[SPARK-17711\]](#) Compress rolled executor log
- [\[SPARK-17751\]](#)[SQL][BACKPORT-2.0] Remove spark.sql.eagerAnalysis and Output the Plan if Existed in AnalysisException
- [\[SPARK-17731\]](#)[SQL][STREAMING] Metrics for structured streaming for branch-2.0
- [\[SPARK-17892\]](#)[SQL][2.0] Do Not Optimize Query in CTAS More Than Once #15048
- [\[SPARK-17819\]](#)[SQL][BRANCH-2.0] Support default database in connection URIs for Spark Thrift Server
- [\[SPARK-17953\]](#)[DOCUMENTATION] Fix typo in SparkSession scaladoc

- [\[SPARK-17863\]](#)[SQL] should not add column into Distinct
- [\[SPARK-17387\]](#)[PYSPARK] Creating SparkContext() from python without spark-submit ignores user conf
- [\[SPARK-17834\]](#)[SQL] Fetch the earliest offsets manually in KafkaSource instead of counting on KafkaConsumer
- [\[SPARK-17876\]](#) Write StructuredStreaming WAL to a stream instead of materializing all at once
- [\[SPARK-16827\]](#)[BRANCH-2.0] Avoid reporting spill metrics as shuffle metrics
- [\[SPARK-17782\]](#)[STREAMING][KAFKA] alternative eliminate race condition of poll twice
- [\[SPARK-17790\]](#)[SPARKR] Support for parallelizing R data.frame larger than 2GB
- [\[SPARK-17884\]](#)[SQL] To resolve Null pointer exception when casting from empty string to interval type.
- [\[SPARK-17808\]](#)[PYSPARK] Upgraded version of Pyrolite to 4.13
- [\[SPARK-17880\]](#)[DOC] The url linking to `AccumulatorV2` in the document is incorrect.
- [\[SPARK-17816\]](#)[CORE][BRANCH-2.0] Fix ConcurrentModificationException issue in BlockStatusesAccumulator
- [\[SPARK-17346\]](#)[SQL][TESTS] Fix the flaky topic deletion in KafkaSourceStressSuite
- [\[SPARK-17738\]](#)[TEST] Fix flaky test in ColumnTypeSuite
- [\[SPARK-17417\]](#)[CORE] Fix # of partitions for Reliable RDD checkpointing
- [\[SPARK-17832\]](#)[SQL] TableIdentifier.quotedString creates un-parseable names when name contains a backtick
- [\[SPARK-17806\]](#) [SQL] fix bug in join key rewritten in HashJoin
- [\[SPARK-17782\]](#)[STREAMING][BUILD] Add Kafka 0.10 project to build modules
- [\[SPARK-17346\]](#)[SQL][TEST-MAVEN] Add Kafka source for Structured Streaming (branch 2.0)
- [\[SPARK-17707\]](#)[WEBUI] Web UI prevents spark-submit application to be finished
- [\[SPARK-17805\]](#)[PYSPARK] Fix in sqlContext.read.text when pass in list of paths
- [\[SPARK-17612\]](#)[SQL][BRANCH-2.0] Support `DESCRIBE table PARTITION` SQL syntax
- [\[SPARK-17792\]](#)[ML] L-BFGS solver for linear regression does not accept general numeric label column types
- [\[SPARK-17750\]](#)[SQL][BACKPORT-2.0] Fix CREATE VIEW with INTERVAL arithmetic
- [\[SPARK-17803\]](#)[TESTS] Upgrade docker-client dependency
- [\[SPARK-17780\]](#)[SQL] Report Throwable to user in StreamExecution
- [\[SPARK-17798\]](#)[SQL] Remove redundant Experimental annotations in sql.streaming
- [\[SPARK-17643\]](#) Remove comparable requirement from Offset (backport for branch-2.0)
- [\[SPARK-17758\]](#)[SQL] Last returns wrong result in case of empty partition
- [\[SPARK-17778\]](#)[TESTS] Mock SparkContext to reduce memory usage of BlockManagerSuite
- [\[SPARK-17773\]](#)[BRANCH-2.0] Input/Output] Add VoidObjectInspector
- [\[SPARK-17549\]](#)[SQL] Only collect table size stat in driver for cached relation.
- [\[SPARK-17559\]](#)[MLLIB] persist edges if their storage level is non in PeriodicGraphCheckpointer
- [\[SPARK-17112\]](#)[SQL] "select null" via JDBC triggers IllegalArgumentException in Thriftserver
- [\[SPARK-17753\]](#)[SQL] Allow a complex expression as the input a value based case statement
- [\[SPARK-17587\]](#)[PYTHON][MLLIB] SparseVector __getitem__ should follow __getitem__ contract
- [\[SPARK-17736\]](#)[DOCUMENTATION][SPARKR] Update R README for rmarkdown,...
- [\[SPARK-17721\]](#)[MLLIB][ML] Fix for multiplying transposed SparseMatrix with SparseVector
- [\[SPARK-17672\]](#) Spark 2.0 history server web Ui takes too long for a single application
- [\[SPARK-17712\]](#)[SQL] Fix invalid pushdown of data-independent filters beneath aggregates
- [\[SPARK-16343\]](#)[SQL] Improve the PushDownPredicate rule to pushdown predicates correctly in non-deterministic condition.
- [\[SPARK-17641\]](#)[SQL] Collect_list/Collect_set should not collect null values.
- [\[SPARK-17673\]](#)[SQL] Incorrect exchange reuse with RowDataSourceScan (backport)
- [\[SPARK-17644\]](#)[CORE] Do not add failedStages when abortStage for fetch failure
- [\[SPARK-17666\]](#) Ensure that RecordReaders are closed by data source file scans (backport)
- [\[SPARK-17056\]](#)[CORE] Fix a wrong assert regarding unroll memory in MemoryStore
- [\[SPARK-17618\]](#) Guard against invalid comparisons between UnsafeRow and other formats
- [\[SPARK-17652\]](#) Fix confusing exception message while reserving capacity
- [\[SPARK-17649\]](#)[CORE] Log how many Spark events got dropped in LiveListenerBus
- [\[SPARK-17650\]](#) malformed url's throw exceptions before bricking Executors

- [\[SPARK-10835\]](#)[ML] Word2Vec should accept non-null string array, in addition to existing null string array
- [\[SPARK-15703\]](#)[SCHEDULER][CORE][WEBUI] Make ListenerBus event queue size configurable (branch 2.0)
- [\[SPARK-4563\]](#)[CORE] Allow driver to advertise a different network address.
- [\[SPARK-17577\]](#)[CORE][2.0 BACKPORT] Update SparkContext.addFile to make it work well on Windows
- [\[SPARK-17210\]](#)[SPARKR] sparkr.zip is not distributed to executors when running sparkr in RStudio
- [\[SPARK-17640\]](#)[SQL] Avoid using -1 as the default batchSize for FileStreamSource.FileEntry
- [\[SPARK-16240\]](#)[ML] ML persistence backward compatibility for LDA - 2.0 backport
- [\[SPARK-17502\]](#)[17609][SQL][BACKPORT][2.0] Fix Multiple Bugs in DDL Statements on Temporary Views
- [\[SPARK-17599\]](#)[\[SPARK-17569\]](#) Backport and to Spark 2.0 branch
- [\[SPARK-17616\]](#)[SQL] Support a single distinct aggregate combined with a non-partial aggregate
- [\[SPARK-17638\]](#)[STREAMING] Stop JVM StreamingContext when the Python process is dead
- [\[SPARK-17613\]](#) S3A base paths with no '/' at the end return empty DataFrames
- [\[SPARK-17421\]](#)[DOCS] Documenting the current treatment of MAVEN_OPTS.
- [\[SPARK-17494\]](#)[SQL] changePrecision() on compact decimal should respect rounding mode
- [\[SPARK-17627\]](#) Mark Streaming Providers Experimental
- [\[SPARK-17512\]](#)[CORE] Avoid formatting to python path for yarn and mesos cluster mode
- [\[SPARK-17418\]](#) Prevent kinesis-asl-assembly artifacts from being published
- [\[SPARK-17617\]](#)[SQL] Remainder(%) expression.eval returns incorrect result on double value
- [\[SPARK-15698\]](#)[SQL][STREAMING] Add the ability to remove the old MetadataLog in FileStreamSource (branch-2.0)
- [\[SPARK-17051\]](#)[SQL] we should use hadoopConf in InsertIntoHiveTable
- [\[SPARK-17513\]](#)[SQL] Make StreamExecution garbage-collect its metadata
- [\[SPARK-17160\]](#) Properly escape field names in code-generated error messages
- [\[SPARK-17100\]](#) [SQL] fix Python udf in filter on top of outer join
- [\[SPARK-16439\]](#) [SQL] bring back the separator in SQL UI
- [\[SPARK-17611\]](#)[yarn][test] Make shuffle service test really test auth.
- [\[SPARK-17433\]](#) YarnShuffleService doesn't handle moving credentials levelDb
- [\[SPARK-17438\]](#)[WEBUI] Show Application.executorLimit in the application page
- [\[SPARK-17473\]](#)[SQL] fixing docker integration tests error due to different versions of jars.
- [\[SPARK-17589\]](#)[TEST][2.0] Fix test case `create external table` in MetastoreDataSourcesSuite
- [\[SPARK-17297\]](#)[DOCS] Clarify window/slide duration as absolute time, not relative to a calendar
- [\[SPARK-17571\]](#)[SQL] AssertOnQuery.condition should always return Boolean value
- [\[SPARK-16462\]](#)[\[SPARK-16460\]](#)[\[SPARK-15144\]](#)[SQL] Make CSV cast null values properly
- [\[SPARK-17586\]](#)[BUILD] Do not call static member via instance reference
- [\[SPARK-17546\]](#)[DEPLOY] start-* scripts should use hostname -f
- [\[SPARK-17541\]](#)[SQL] fix some DDL bugs about table management when same-name temp view exists
- [\[SPARK-17480\]](#)[SQL][FOLLOWUP] Fix more instances which calls List.length/size which is O(n)
- [\[SPARK-17491\]](#) Close serialization stream to fix wrong answer bug in putIteratorAsBytes()
- [\[SPARK-17575\]](#)[DOCS] Remove extra table tags in configuration document
- [\[SPARK-17548\]](#)[MLLIB] Word2VecModel.findSynonyms no longer spuriously rejects the best match when invoked with a vector
- [\[SPARK-17561\]](#)[DOCS] DataFrameWriter documentation formatting problems
- [\[SPARK-17567\]](#)[DOCS] Use valid url to Spark RDD paper
- [\[SPARK-17558\]](#) Bump Hadoop 2.7 version from 2.7.2 to 2.7.3
- [\[SPARK-17484\]](#) Prevent invalid block locations from being reported after put() exceptions
- [\[SPARK-17364\]](#)[SQL] Antlr lexer wrongly treats full qualified identifier as a decimal number token when parsing SQL string
- [\[SPARK-17483\]](#) Refactoring in BlockManager status reporting and block removal
- [\[SPARK-17114\]](#)[SQL] Fix aggregates grouped by literals with empty input
- [\[SPARK-17547\]](#) Ensure temp shuffle data file is cleaned up after error
- [\[SPARK-17521\]](#) Error when I use sparkContext.makeRDD(Seq())

- [\[SPARK-17465\]](#)[SPARK CORE] Inappropriate memory management in `org.apache.spark.storage.MemoryStore`` may lead to memory leak
- [\[SPARK-17463\]](#)[CORE] Make `CollectionAccumulator` and `SetAccumulator`'s value can be read thread-safely
- [\[SPARK-17511\]](#) Yarn Dynamic Allocation: Avoid marking released container as Failed
- [\[SPARK-17514\]](#) `df.take(1)` and `df.limit(1).collect()` should perform the same in Python
- [\[SPARK-17445\]](#)[DOCS] Reference an ASF page as the main place to find third-party packages
- [\[SPARK-16711\]](#) `YarnShuffleService` doesn't re-init properly on YARN rolling upgrade
- [\[SPARK-15074\]](#)[SHUFFLE] Cache shuffle index file to speedup shuffle fetch
- [\[SPARK-17480\]](#)[SQL] Improve performance by removing or caching `List.length` which is $O(n)$
- [\[SPARK-17525\]](#)[PYTHON] Remove `SparkContext.clearFiles()` from the PySpark API as it was removed from the Scala API prior to Spark 2.0.0
- [\[SPARK-17531\]](#) Don't initialize Hive Listeners for the Execution Client
- [\[SPARK-17515\]](#) `CollectLimit.execute()` should perform per-partition limits
- [\[SPARK-17474\]](#) [SQL] fix python udf in `TakeOrderedAndProjectExec`
- [\[SPARK-17485\]](#) Prevent failed remote reads of cached blocks from failing entire job
- [\[SPARK-14818\]](#) Post-2.0 MiMa exclusion and build changes
- [\[SPARK-17503\]](#)[CORE] Fix memory leak in Memory store when unable to cache the whole RDD in memory
- [\[SPARK-17486\]](#) Remove unused `TaskMetricsUIData.updatedBlockStatuses` field
- [\[SPARK-17336\]](#)[PYSARK] Fix appending multiple times to `PYTHONPATH` from `spark-config.sh`
- [\[SPARK-17439\]](#)[SQL] Fixing compression issues with approximate quantiles and adding more tests
- [\[SPARK-17396\]](#)[CORE] Share the task support between `UnionRDD` instances.
- [\[SPARK-17354\]](#) [SQL] Partitioning by dates/timestamps should work with Parquet vectorized reader
- [\[SPARK-17456\]](#)[CORE] Utility for parsing Spark versions
- [\[SPARK-17339\]](#)[CORE][BRANCH-2.0] Do not use path to get a filesystem in `hadoopFile` and `newHadoopFile` APIs
- [\[SPARK-16533\]](#)[CORE] - backport driver deadlock fix to 2.0
- [\[SPARK-17370\]](#) Shuffle service files not invalidated when a slave is lost
- [\[SPARK-17296\]](#)[SQL] Simplify parser join processing [BACKPORT 2.0]
- [\[SPARK-17372\]](#)[SQL][STREAMING] Avoid serialization issues by using Arrays to save file names in `FileStreamSource`
- [\[SPARK-17279\]](#)[SQL] better error message for exceptions during `ScalaUDF` execution
- [\[SPARK-17316\]](#)[CORE] Fix the 'ask' type parameter in 'removeExecutor'
- [\[SPARK-17110\]](#) Fix `StreamCorruptionException` in `BlockManager.getRemoteValues()`
- [\[SPARK-17299\]](#) TRIM/LTRIM/RTRIM should not strips characters other than spaces
- [\[SPARK-16334\]](#) [BACKPORT] Reusing same dictionary column for decoding consecutive row groups shouldn't throw an error
- [\[SPARK-16922\]](#) [\[SPARK-17211\]](#) [SQL] make the address of values portable in `LongToUnsafeRowMap`
- [\[SPARK-17356\]](#)[SQL] Fix out of memory issue when generating JSON for `TreeNode`
- [\[SPARK-17369\]](#)[SQL][2.0] `MetastoreRelation` toJSON throws `AssertException` due to missing `otherCopyArgs`
- [\[SPARK-17358\]](#)[SQL] Cached table(parquet/orc) should be shard between beelines
- [\[SPARK-17353\]](#)[\[SPARK-16943\]](#)[\[SPARK-16942\]](#)[BACKPORT-2.0][SQL] Fix multiple bugs in CREATE TABLE LIKE command
- [\[SPARK-17391\]](#)[TEST][2.0] Fix Two Test Failures After Backport
- [\[SPARK-17335\]](#)[SQL] Fix `ArrayType` and `MapType` `CatalogString`.
- [\[SPARK-16663\]](#)[SQL] desc table should be consistent between data source and hive serde tables
- [\[SPARK-16959\]](#)[SQL] Rebuild Table Comment when Retrieving Metadata from Hive Metastore
- [\[SPARK-17347\]](#)[SQL][EXAMPLES] Encoder in Dataset example has incorrect type
- [\[SPARK-17230\]](#) [SQL] Should not pass optimized query into `QueryExecution` in `DataFrameWriter`
- [\[SPARK-17261\]](#) [PYSARK] Using `HiveContext` after re-creating `SparkContext` in Spark 2.0 throws `"Java.lang.IllegalStateException: Cannot call methods on a stopped sparkContext"`
- [\[SPARK-16935\]](#)[SQL] Verification of Function-related ExternalCatalog APIs
- [\[SPARK-17352\]](#)[WEBUI] Executor computing time can be negative-number because of calculation error
- [\[SPARK-17342\]](#)[WEBUI] Style of event timeline is broken

- [\[SPARK-17355\]](#) Workaround for HIVE-14684 / HiveResultSetMetaData.isSigned exception
- [\[SPARK-16926\]](#) [SQL] Remove partition columns from partition metadata.
- [\[SPARK-17271\]](#)[SQL] Planner adds un-necessary Sort even if child orde...
- [\[SPARK-17318\]](#)[TESTS] Fix RepSuite replicating blocks of object with class defined in repl again
- [\[SPARK-17180\]](#)[\[SPARK-17309\]](#)[\[SPARK-17323\]](#)[SQL][2.0] create AlterViewAsCommand to handle ALTER VIEW AS
- [\[SPARK-17316\]](#)[TESTS] Fix MesosCoarseGrainedSchedulerBackendSuite
- [\[SPARK-17316\]](#)[CORE] Make CoarseGrainedSchedulerBackend.removeExecutor non-blocking
- [\[SPARK-17243\]](#)[WEB UI] Spark 2.0 History Server won't load with very large application history
- [\[SPARK-17318\]](#)[TESTS] Fix RepSuite replicating blocks of object with class defined in repl
- [\[SPARK-17264\]](#)[SQL] DataStreamWriter should document that it only supports Parquet for now
- [\[SPARK-17301\]](#)[SQL] Remove unused classTag field from AtomicType base class
- [\[SPARK-17063\]](#) [SQL] Improve performance of MSCK REPAIR TABLE with Hive metastore
- [\[SPARK-16216\]](#)[SQL][FOLLOWUP][BRANCH-2.0] Bacoport enabling timestamp type tests for JSON and verify all unsupported types in CSV
- [\[SPARK-17216\]](#)[UI] fix event timeline bars length
- [ML][MLLIB] The require condition and message doesn't match in SparseMatrix.
- [\[SPARK-15382\]](#)[SQL] Fix a bug in sampling with replacement
- [\[SPARK-17274\]](#)[SQL] Move join optimizer rules into a separate file
- [\[SPARK-17270\]](#)[SQL] Move object optimization rules into its own file (branch-2.0)
- [\[SPARK-17269\]](#)[SQL] Move finish analysis optimization stage into its own file
- [\[SPARK-17244\]](#) Catalyst should not pushdown non-deterministic join conditions
- [\[SPARK-17235\]](#)[SQL] Support purging of old logs in MetadataLog
- [\[SPARK-17246\]](#)[SQL] Add BigDecimal literal
- [\[SPARK-17165\]](#)[SQL] FileStreamSource should not track the list of seen files indefinitely
- [\[SPARK-17242\]](#)[DOCUMENT] Update links of external dstream projects
- [\[SPARK-17231\]](#)[CORE] Avoid building debug or trace log messages unless the respective log level is enabled
- [\[SPARK-17205\]](#) Literal.sql should handle Infinity and NaN
- [\[SPARK-15083\]](#)[WEB UI] History Server can OOM due to unlimited TaskUIData
- [\[SPARK-16700\]](#)[PYSPARK][SQL] create DataFrame from dict/Row with schema
- [\[SPARK-17167\]](#)[2.0][SQL] Issue Exceptions when Analyze Table on In-Memory Cataloged Tables
- [\[SPARK-16991\]](#)[\[SPARK-17099\]](#)[\[SPARK-17120\]](#)[SQL] Fix Outer Join Elimination when Filter's isNotNull Constraints Unable to Filter Out All Null-supplying Rows
- [\[SPARK-17061\]](#)[\[SPARK-17093\]](#)[SQL][BACKPORT] MapObjects should make copies of unsafe-backed data
- [\[SPARK-17193\]](#)[CORE] HadoopRDD NPE at DEBUG log level when getLocationInfo == null
- [\[SPARK-17228\]](#)[SQL] Not infer/propagate non-deterministic constraints
- [\[SPARK-16216\]](#)[SQL][BRANCH-2.0] Backport Read/write dateFormat/timestampFormat options for CSV and JSON
- [\[SPARK-16781\]](#)[PYSPARK] java launched by PySpark as gateway may not be the same java used in the spark environment
- [\[SPARK-17086\]](#)[ML] Fix InvalidArgumentException issue in QuantileDiscretizer when some quantiles are duplicated
- [\[SPARK-17186\]](#)[SQL] remove catalog table type INDEX
- [\[SPARK-17194\]](#) Use single quotes when generating SQL for string literals
- [\[SPARK-13286\]](#) [SQL] add the next expression of SQLException as cause
- [\[SPARK-17182\]](#)[SQL] Mark Collect as non-deterministic
- [\[SPARK-16550\]](#)[\[SPARK-17042\]](#)[CORE] Certain classes fail to deserialize in block manager replication
- [\[SPARK-17162\]](#) Range does not support SQL generation
- [\[SPARK-16320\]](#)[DOC] Document G1 heap region's effect on spark 2.0 vs 1.6
- [\[SPARK-17085\]](#)[STREAMING][DOCUMENTATION AND ACTUAL CODE DIFFERS - UNSUPPORTED OPERATIONS]
- [\[SPARK-17115\]](#)[SQL] decrease the threshold when split expressions
- [\[SPARK-17098\]](#)[SQL] Fix `NullPropagation` optimizer to handle `COUNT(NULL) OVER` correctly
- [\[SPARK-12666\]](#)[CORE] SparkSubmit packages fix for when 'default' conf doesn't exist in dependent module

- [\[SPARK-17124\]](#)[SQL] RelationalGroupedDataset.agg should preserve order and allow multiple aggregates per column
- [\[SPARK-17104\]](#)[SQL] LogicalRelation.newInstance should follow the semantics of MultiInstanceRelation
- [\[SPARK-17150\]](#)[SQL] Support SQL generation for inline tables
- [\[SPARK-17158\]](#)[SQL] Change error message for out of range numeric literals
- [\[SPARK-17149\]](#)[SQL] array.sql for testing array related functions
- [\[SPARK-17113\]](#) [SHUFFLE] Job failure due to Executor OOM in offheap mode
- [\[SPARK-16686\]](#)[SQL] Remove PushProjectThroughSample since it is handled by ColumnPruning
- [\[SPARK-11227\]](#)[CORE] UnknownHostException can be thrown when NameNode HA is enabled.
- [\[SPARK-16994\]](#)[SQL] Whitelist operators for predicate pushdown
- [\[SPARK-16961\]](#)[CORE] Fixed off-by-one error that biased randomizeInPlace
- [\[SPARK-16947\]](#)[SQL] Support type coercion and foldable expression for inline tables
- [\[SPARK-17069\]](#) Expose spark.range() as table-valued function in SQL
- [\[SPARK-17117\]](#)[SQL] 1 / NULL should not fail analysis
- [\[SPARK-16391\]](#)[SQL] Support partial aggregation for reduceGroups
- [\[SPARK-16995\]](#)[SQL] TreeNodeException when flat mapping RelationalGroupedDataset created from DataFrame containing a column created with lit/expr
- [\[SPARK-17096\]](#)[SQL][STREAMING] Improve exception string reported through the StreamingQueryListener
- [\[SPARK-17102\]](#)[SQL] bypass UserDefinedGenerator for json format check
- [\[SPARK-15285\]](#)[SQL] Generated SpecificSafeProjection.apply method grows beyond 64 KB
- [\[SPARK-17084\]](#)[SQL] Rename ParserUtils.assert to validate
- [\[SPARK-17089\]](#)[DOCS] Remove api doc link for mapReduceTriplets operator
- [\[SPARK-16964\]](#)[SQL] Remove private[sql] and private[spark] from sql.execution package [Backport]
- [\[SPARK-17065\]](#)[SQL] Improve the error message when encountering an incompatible DataSourceRegister
- [\[SPARK-16508\]](#)[SPARKR] Split docs for arrange and orderBy methods
- [\[SPARK-17027\]](#)[ML] Avoid integer overflow in PolynomialExpansion.getPolySize
- [\[SPARK-16966\]](#)[SQL][CORE] App Name is a randomUUID even when "spark.app.name" exists
- [\[SPARK-17013\]](#)[SQL] Parse negative numeric literals
- [\[SPARK-16975\]](#)[SQL] Column-partition path starting '_' should be handled correctly
- [\[SPARK-17022\]](#)[YARN] Handle potential deadlock in driver handling messages
- [\[SPARK-17018\]](#)[SQL] literals.sql for testing literal parsing
- [\[SPARK-17015\]](#)[SQL] group-by/order-by ordinal and arithmetic tests
- [\[SPARK-15899\]](#)[SQL] Fix the construction of the file path with hadoop Path for Spark 2.0
- [\[SPARK-17011\]](#)[SQL] Support testing exceptions in SQLQueryTestSuite
- [\[SPARK-17007\]](#)[SQL] Move test data files into a test-data folder
- [\[SPARK-17008\]](#)[\[SPARK-17009\]](#)[SQL] Normalization and isolation in SQLQueryTestSuite.
- [\[SPARK-16866\]](#)[SQL] Infrastructure for file-based SQL end-to-end tests
- [\[SPARK-17010\]](#)[MINOR][DOC] Wrong description in memory management document
- [\[SPARK-15639\]](#) [\[SPARK-16321\]](#) [SQL] Push down filter at RowGroups level for parquet reader
- [\[SPARK-16324\]](#)[SQL] regexp_extract should doc that it returns empty string when match fails
- [\[SPARK-16522\]](#)[MESOS] Spark application throws exception on exit.
- [\[SPARK-16905\]](#) SQL DDL: MSCK REPAIR TABLE
- [\[SPARK-16956\]](#) Make ApplicationState.MAX_NUM_RETRY configurable
- [\[SPARK-16950\]](#) [PYSPARK] fromOffsets parameter support in KafkaUtils.createDirectStream for python3
- [\[SPARK-16610\]](#)[SQL] Add `orc.compress` as an alias for `compression` option.
- [\[SPARK-16563\]](#)[SQL] fix spark sql thrift server FetchResults bug
- [\[SPARK-16953\]](#) Make requestTotalExecutors public Developer API to be consistent with requestExecutors/killExecutors
- [\[SPARK-16586\]](#)[CORE] Handle JVM errors printed to stdout.
- [\[SPARK-16936\]](#)[SQL] Case Sensitivity Support for Refresh Temp Table

- [\[SPARK-16457\]](#)[SQL] Fix Wrong Messages when CTAS with a Partition By Clause
- [\[SPARK-16939\]](#)[SQL] Fix build error by using `Tuple1` explicitly in StringFunctionsSuite
- [\[SPARK-16409\]](#)[SQL] regexp_extract with optional groups causes NPE
- [\[SPARK-16911\]](#) Fix the links in the programming guide
- [\[SPARK-16870\]](#)[DOCS] Summary:add "spark.sql.broadcastTimeout" into docs/sql-programming-gu...
- [\[SPARK-16932\]](#)[DOCS] Changed programming guide to not reference old accumulator API in Scala
- [\[SPARK-16925\]](#) Master should call schedule() after all executor exit events, not only failures
- [\[SPARK-16772\]](#)[PYTHON][DOCS] Fix API doc references to UDFRegistration + Update "important classes"
- [\[SPARK-16750\]](#)[FOLLOW-UP][ML] Add transformSchema for StringIndexer/VectorAssembler and fix failed tests.
- [\[SPARK-16907\]](#)[SQL] Fix performance regression for parquet table when vectorized parquet record reader is not being used
- [\[SPARK-16863\]](#)[ML] ProbabilisticClassifier.fit check thresholds' length
- [\[SPARK-16877\]](#)[BUILD] Add rules for preventing to use Java annotations (Deprecated and Override)
- [\[SPARK-16880\]](#)[ML][MLLIB] make ann training data persisted if needed
- [\[SPARK-16875\]](#)[SQL] Add args checking for DataSet randomSplit and sample
- [\[SPARK-16802\]](#) [SQL] fix overflow in LongToUnsafeRowMap
- [\[SPARK-16873\]](#)[CORE] Fix SpillReader NPE when spillFile has no data
- [\[SPARK-14204\]](#)[SQL] register driverClass rather than user-specified class
- [\[SPARK-16714\]](#)[\[SPARK-16735\]](#)[\[SPARK-16646\]](#) array, map, greatest, least's type coercion should handle decimal type
- [\[SPARK-16796\]](#)[WEB UI] Visible passwords on Spark environment page
- [\[SPARK-16831\]](#)[PYTHON] Fixed bug in CrossValidator.avgMetrics
- [\[SPARK-16787\]](#) SparkContext.addFile() should not throw if called twice with the same file
- [\[SPARK-16850\]](#)[SQL] Improve type checking error message for greatest/least
- [\[SPARK-16836\]](#)[SQL] Add support for CURRENT_DATE/CURRENT_TIMESTAMP literals
- [\[SPARK-16062\]](#) [\[SPARK-15989\]](#) [SQL] Fix two bugs of Python-only UDTs
- [\[SPARK-16837\]](#)[SQL] TimeWindow incorrectly drops slideDuration in constructors
- [\[SPARK-15541\]](#) Casting ConcurrentHashMap to ConcurrentMap (master branch)
- [\[SPARK-16558\]](#)[EXAMPLES][MLLIB] examples/mllib/LDAExample should use MLVector instead of MLlib Vector
- [\[SPARK-16734\]](#)[EXAMPLES][SQL] Revise examples of all language bindings
- [\[SPARK-16818\]](#) Exchange reuse incorrectly reuses scans over different sets of partitions
- [\[SPARK-15869\]](#)[STREAMING] Fix a potential NPE in StreamingJobProgressListener.getBatchUIData
- [\[SPARK-16774\]](#)[SQL] Fix use of deprecated timestamp constructor & improve timezone handling
- [\[SPARK-16791\]](#)[SQL] cast struct with timestamp field fails
- [\[SPARK-16778\]](#)[SQL][TRIVIAL] Fix deprecation warning with SQLContext
- [\[SPARK-16805\]](#)[SQL] Log timezone when query result does not match
- [\[SPARK-16813\]](#)[SQL] Remove private[sql] and private[spark] from catalyst package
- [\[SPARK-16812\]](#) Open up SparkILoop.getAddedJars
- [\[SPARK-16800\]](#)[EXAMPLES][ML] Fix Java examples that fail to run due to exception
- [\[SPARK-16748\]](#)[SQL] SparkExceptions during planning should not wrapped in TreeNodeException
- [\[SPARK-16761\]](#)[DOC][ML] Fix doc link in docs/ml-guide.md
- [\[SPARK-16750\]](#)[ML] Fix GaussianMixture training failed due to feature column type mistake
- [\[SPARK-16664\]](#)[SQL] Fix persist call on Data frames with more than 200...
- [\[SPARK-16772\]](#) Correct API doc references to PySpark classes + formatting fixes
- [\[SPARK-16764\]](#)[SQL] Recommend disabling vectorized parquet reader on OutOfMemoryError
- [\[SPARK-16740\]](#)[SQL] Fix Long overflow in LongToUnsafeRowMap
- [\[SPARK-16639\]](#)[SQL] The query with having condition that contains grouping by column should work
- [\[SPARK-15232\]](#)[SQL] Add subquery SQL building tests to LogicalPlanToSQLSuite
- [\[SPARK-16730\]](#)[SQL] Implement function aliases for type casts
- [\[SPARK-16729\]](#)[SQL] Throw analysis exception for invalid date casts

- [\[SPARK-16621\]](#)[SQL] Generate stable SQLs in SQLBuilder
- [\[SPARK-16633\]](#)[\[SPARK-16642\]](#)[\[SPARK-16721\]](#)[SQL] Fixes three issues related to lead and lag functions
- [\[SPARK-16724\]](#) Expose DefinedByConstructorParams
- [\[SPARK-16672\]](#)[SQL] SQLBuilder should not raise exceptions on EXISTS queries
- [\[SPARK-16722\]](#)[TESTS] Fix a StreamingContext leak in StreamingContextSuite when eventually fails
- [\[SPARK-14131\]](#)[STREAMING] SQL Improved fix for avoiding potential deadlocks in HDFSMetadataLog
- [\[SPARK-16715\]](#)[TESTS] Fix a potential ExprId conflict for SubexpressionEliminationSuite. "Semantic equals and hash"
- [\[SPARK-16485\]](#)[DOC][ML] Fixed several inline formatting in ml features doc
- [\[SPARK-16703\]](#)[SQL] Remove extra whitespace in SQL generation for window functions
- [\[SPARK-16698\]](#)[SQL] Field names having dots should be allowed for datasources based on FileFormat
- [\[SPARK-16648\]](#)[SQL] Make ignoreNullsExpr a child expression of First and Last
- [\[SPARK-16699\]](#)[SQL] Fix performance bug in hash aggregate on long string keys
- [\[SPARK-16515\]](#)[SQL][FOLLOW-UP] Fix test `script` on OS X/Windows...
- [\[SPARK-16690\]](#)[TEST] rename SQLTestUtils.withTempTable to withTempView
- [\[SPARK-16380\]](#)[EXAMPLES] Update SQL examples and programming guide for Python language binding
- [\[SPARK-16651\]](#)[PYSARK][DOC] Make `withColumnRenamed/drop` description more consistent with Scala API
- [\[SPARK-16650\]](#) Improve documentation of spark.task.maxFailures
- [\[SPARK-16287\]](#)[HOTFIX][BUILD][SQL] Fix annotation argument needs to be a constant
- [\[SPARK-16287\]](#)[SQL] Implement str_to_map SQL function
- [\[SPARK-16334\]](#) Maintain single dictionary per row-batch in vectorized parquet reader
- [\[SPARK-16656\]](#)[SQL] Try to make CreateTableAsSelectSuite more stable
- [\[SPARK-16644\]](#)[SQL] Aggregate should not propagate constraints containing aggregate expressions
- [\[SPARK-16440\]](#)[MLLIB] Destroy broadcasted variables even on driver
- [\[SPARK-5682\]](#)[CORE] Add encrypted shuffle in spark
- [\[SPARK-16901\]](#) Hive settings in hive-site.xml may be overridden by Hive's default values
- [\[SPARK-16272\]](#)[CORE] Allow config values to reference conf, env, system props.
- [\[SPARK-16632\]](#)[SQL] Use Spark requested schema to guide vectorized Parquet reader initialization
- [\[SPARK-16634\]](#)[SQL] Workaround JVM bug by moving some code out of ctor.
- [\[SPARK-16505\]](#)[YARN] Optionally propagate error during shuffle service startup.
- [\[SPARK-14963\]](#)[MINOR][YARN] Fix typo in YarnShuffleService recovery file name
- [\[SPARK-14963\]](#)[YARN] Using recoveryPath if NM recovery is enabled
- [\[SPARK-16349\]](#)[SQL] Fall back to isolated class loader when classes not found.
- [\[SPARK-16119\]](#)[sql] Support PURGE option to drop table / partition.

CDS Powered by Apache Spark Version, Packaging, and Download Information

The following sections provide links to the parcel and service descriptor files for the different CDS versions, as well as information about using CDS with Maven.

CDS Versions Available for Download



Note: The parcel version displayed in Cloudera Manager, which is also part of the parcel file name, is structured as follows:

`<CDS_version>-1.<cdh_build_version>.p<patch_version>.<build_number>`

For example:

`2.0.0.cloudera2-1.cdh5.7.0.p0.118100`

The `<cdh_build_version>` portion is the version of CDH upon which the release was built. It is *not* the minimum supported CDH version.

To view the supported CDH versions and other requirements, see [CDS Powered by Apache Spark Requirements](#) on page 6.

Table 1: Available CDS Versions

Version	CSD	Parcel
2.0 Release 2	SPARK2_ON_YARN-2.0.0.cloudera2.jar	The exact parcel name is dependent on the OS. You can find all the parcels at http://archive.cloudera.com/spark2/parcels/2.0.0.cloudera2/ .
2.0 Release 1	SPARK2_ON_YARN-2.0.0.cloudera1.jar	The exact parcel name is dependent on the OS. You can find all the parcels at http://archive.cloudera.com/spark2/parcels/2.0.0.cloudera2/ .

CDS Maven Artifacts

For information about using CDS Maven artifacts, see [Using the CDS Powered by Apache Spark Maven Repository](#) on page 20.

Using the CDS Powered by Apache Spark Maven Repository



Important: CDS 2 no longer includes an assembly JAR. When you build an application JAR, *do not* include CDH or CDS JARs, because they are already provided. If you do, upgrading CDH or CDS can break your application. To avoid this situation, set the Maven dependency `scope` to `provided`. If you have already built applications which include the CDH or CDS JARs, update the dependency to set `scope` to `provided` and recompile.

For information on how to use CDH artifacts, see [Using the CDH 5 Maven Repository](#).

If you want to build applications or tools for use with CDS Powered by Apache Spark, and you are using Maven or Ivy for dependency management, you can pull the CDS artifacts from the Cloudera Maven repository. The repository is available at <https://repository.cloudera.com/artifactory/cloudera-repos/>.

The following is a sample POM (`pom.xml`) file:

```
<project xmlns="http://maven.apache.org/POM/4.0.0"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
http://maven.apache.org/maven-v4_0_0.xsd">
```

```

<repositories>
  <repository>
    <id>cloudera</id>
    <url>https://repository.cloudera.com/artifactory/cloudera-repos/</url>
  </repository>
</repositories>
</project>

```

For more information about the Maven artifacts for each CDS release, see the following topics:

CDS 2.0 Powered by Apache Spark Maven Artifacts

The following tables lists the `groupId`, `artifactId`, and `version` required to access the artifacts for each CDS 2.0 Powered by Apache Spark release:

CDS 2.0 Release 2 Maven Artifacts

The following pom fragment shows how to access a CDS 2.0 Release 2 artifact from a Maven POM.

```

<dependency>
  <groupId>org.apache.spark</groupId>
  <artifactId>spark-core_2.11</artifactId>
  <version>2.0.0.cloudera2</version>
  <scope>provided</scope>
</dependency>

```

The complete artifact list for this release follows.

Project	groupId	artifactId	version
Apache Spark	org.apache.spark	spark-catalyst_2.11	2.0.0.cloudera2
	org.apache.spark	spark-core_2.11	2.0.0.cloudera2
	org.apache.spark	spark-graphx_2.11	2.0.0.cloudera2
	org.apache.spark	spark-hive-exec_2.11	2.0.0.cloudera2
	org.apache.spark	spark-hive_2.11	2.0.0.cloudera2
	org.apache.spark	spark-launcher_2.11	2.0.0.cloudera2
	org.apache.spark	spark-mllib-local_2.11	2.0.0.cloudera2
	org.apache.spark	spark-mllib_2.11	2.0.0.cloudera2
	org.apache.spark	spark-network-common_2.11	2.0.0.cloudera2
	org.apache.spark	spark-network-shuffle_2.11	2.0.0.cloudera2
	org.apache.spark	spark-network-yarn_2.11	2.0.0.cloudera2
	org.apache.spark	spark-repl_2.11	2.0.0.cloudera2

Project	groupId	artifactId	version
	org.apache.spark	spark-sketch_2.11	2.0.0.cloudera2
	org.apache.spark	spark-sql_2.11	2.0.0.cloudera2
	org.apache.spark	spark-streaming-flume-sink_2.11	2.0.0.cloudera2
	org.apache.spark	spark-streaming-flume_2.11	2.0.0.cloudera2
	org.apache.spark	spark-streaming-kafka-0-8_2.11	2.0.0.cloudera2
	org.apache.spark	spark-streaming_2.11	2.0.0.cloudera2
	org.apache.spark	spark-tags_2.11	2.0.0.cloudera2
	org.apache.spark	spark-unsafe_2.11	2.0.0.cloudera2
	org.apache.spark	spark-yarn_2.11	2.0.0.cloudera2

CDS 2.0 Release 1 Maven Artifacts

The following pom fragment shows how to access a CDS 2.0 Release 1 artifact from a Maven POM.

```
<dependency>
  <groupId>org.apache.spark</groupId>
  <artifactId>spark-core_2.11</artifactId>
  <version>2.0.0.cloudera1</version>
  <scope>provided</scope>
</dependency>
```

The complete artifact list for this release follows.

Project	groupId	artifactId	version
Apache Spark	org.apache.spark	spark-catalyst_2.11	2.0.0.cloudera1
	org.apache.spark	spark-core_2.11	2.0.0.cloudera1
	org.apache.spark	spark-graphx_2.11	2.0.0.cloudera1
	org.apache.spark	spark-hive_2.11	2.0.0.cloudera1
	org.apache.spark	spark-launcher_2.11	2.0.0.cloudera1
	org.apache.spark	spark-mllib-local_2.11	2.0.0.cloudera1
	org.apache.spark	spark-mllib_2.11	2.0.0.cloudera1

Project	groupId	artifactId	version
	org.apache.spark	spark-network-common_2.11	2.0.0-cloudera1
	org.apache.spark	spark-network-shuffle_2.11	2.0.0-cloudera1
	org.apache.spark	spark-network-yarn_2.11	2.0.0-cloudera1
	org.apache.spark	spark-repl_2.11	2.0.0-cloudera1
	org.apache.spark	spark-sketch_2.11	2.0.0-cloudera1
	org.apache.spark	spark-sql_2.11	2.0.0-cloudera1
	org.apache.spark	spark-streaming-flume-sink_2.11	2.0.0-cloudera1
	org.apache.spark	spark-streaming-flume_2.11	2.0.0-cloudera1
	org.apache.spark	spark-streaming-kafka-0-8_2.11	2.0.0-cloudera1
	org.apache.spark	spark-streaming_2.11	2.0.0-cloudera1
	org.apache.spark	spark-tags_2.11	2.0.0-cloudera1
	org.apache.spark	spark-unsafe_2.11	2.0.0-cloudera1
	org.apache.spark	spark-yarn_2.11	2.0.0-cloudera1

Installing CDS 2 Powered by Apache Spark

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

CDS 2 Powered by Apache Spark is distributed as two files: a CSD file and a parcel, both of which need to be installed on the cluster. You can install Spark2 using the following instructions.



Important:


Because CDS Powered by Apache Spark is only installable using the parcel mechanism, it can only be used on clusters managed by Cloudera Manager. Additionally, because Cloudera Manager does not support using parcels and packages in the same cluster, you cannot use CDS if you are using a package-based installation of CDH.

If your Cloudera Manager Server does not have Internet access, you can use the CSD and proceed with the following instructions:

1. Download the [Spark2 CSD](#).
2. Install the Spark2 CSD into Cloudera Manager as described in [Installing an Add-on Service](#).
3. In the Cloudera Manager Admin Console, add the [Spark2 parcel repository](#) to the [remote repository URLs](#).



Note: If your Cloudera Manager Server does not have Internet access, you can use the Spark 2 parcel files, put them into a [new parcel repository](#), and then configure the Cloudera Manager Server to target this newly-created repository.

4. Download the Spark2 parcel, distribute the parcel to the hosts in your cluster, and activate the parcel. See [Managing Parcels](#).
5. [Add the Spark 2 service](#) to your cluster. When configuring the assignment of role instances to hosts, add a [gateway role](#) to every host. The History Server port is 18089 instead of the usual 18088.
6. Return to the Home page by clicking the Cloudera Manager logo.
7. Click  to invoke the cluster restart wizard.
8. Click **Restart Stale Services**.
9. Click **Restart Now**.
10. Click **Finish**.

Administering CDS 2 Powered by Apache Spark

Most administration tasks are the same whether you are using Spark 1 or Spark 2. To configure and manage Spark, follow the procedures in the Cloudera Enterprise [Spark Guide](#).

In addition, follow these procedures that are specific to Spark 2:

Configuring Spark 2 Tools as the Default



Important: If you configure Spark 2 as the default, modules such as SparkOnHBase and HiveOnSpark no longer work due to dependencies on Spark 1.6 in CDH. For more information, see [CDS Powered by Apache Spark Known Issues](#) on page 7.

When you start trying out Spark 2, you can do most of your testing by running the standard Spark 1 commands such as `pyspark` and `spark-shell` alongside their Spark 2 equivalents such as `pyspark2` and `spark2-shell`. All of these commands are represented as symbolic links in `/usr/bin`.

If you are testing a workflow that has the original command names hardcoded in other scripts, you might configure the system so that issuing the `pyspark` command really runs the `pyspark2` script, and so on for other Spark-related binaries. This change is done using the Linux “alternatives” mechanism, which keeps track of the appropriate target for each of the `/usr/bin` symlinks.

To use Spark 2 tools as the default, run the following script *on all hosts in the cluster*:

```
for binary in pyspark spark-shell spark-submit; do
  # Generate the name of the new binary e.g. pyspark2, spark2-shell, etc.
  new_binary=$(echo $binary | sed -e 's/spark/spark2/')
  # Update the old alternative to the client binary to the new client binary
  # Use priority 11 because the default priority with which these alternatives are
  # created is 10
  update-alternatives --install /usr/bin/${binary} ${binary} /usr/bin/${new_binary} 11
done
# For configuration, we need to have a separate command
# because the destination is under /etc/ instead of /usr/bin like for binaries.
# The priority is different - 52 because Cloudera Manager sets up configuration symlinks
# with priority 51.
update-alternatives --install /etc/spark/conf spark-conf /etc/spark2/conf 52
```

To remove this setting and return to using the Spark contained in CDH, run the following script *on all hosts in the cluster*. It removes the Spark 2 targets of the symlinks and points those symlinks back to the original Spark-related scripts:

```
for binary in pyspark spark-shell spark-submit; do
  new_binary=$(echo $binary | sed -e 's/spark/spark2/')
  update-alternatives --remove ${binary} /usr/bin/${new_binary}
done
update-alternatives --remove spark-conf /etc/spark2/conf
```

Running Applications with CDS Powered by Apache Spark

With CDS Powered by Apache Spark, you can run Apache Spark 2 applications locally or distributed across a cluster, either by using an interactive shell or by submitting an application. Running Spark applications interactively is commonly performed during the data-exploration phase and for ad hoc analysis.

The Spark 2 Job Commands

With Spark 2, you use slightly different command names than in Spark 1, so that you can run both versions of Spark side-by-side without conflicts:

- `spark2-submit` instead of `spark-submit`.
- `spark2-shell` instead of `spark-shell`.
- `pyspark2` instead of `pyspark`.

For development and test purposes, you can also configure each host so that invoking the Spark 1 command name runs the corresponding Spark 2 executable. See [Configuring Spark 2 Tools as the Default](#) on page 25 for details.

Canary Test for pyspark2 Command

The following example shows a simple `pyspark2` session that refers to the `SparkContext`, calls the `collect()` function which runs a Spark 2 job, and writes data to HDFS. This sequence of operations helps to check if there are obvious configuration issues that prevent Spark 2 jobs from working at all. For the HDFS path for the output directory, substitute a path that exists on your own system.

```
$ hdfs dfs -mkdir /user/jdoe/spark
$ pyspark2
...
SparkSession available as 'spark'.
>>> strings = ["one","two","three"]
>>> s2 = sc.parallelize(strings)
>>> s3 = s2.map(lambda word: word.upper())
>>> s3.collect()
['ONE', 'TWO', 'THREE']
>>> s3.saveAsTextFile('hdfs:///user/jdoe/spark/canary_test')
>>> quit()
$ hdfs dfs -ls /user/jdoe/spark
Found 1 items
drwxr-xr-x - jdoe spark-users 0 2016-08-26 14:41 /user/jdoe/spark/canary_test
$ hdfs dfs -ls /user/jdoe/spark/canary_test
Found 3 items
-rw-r--r-- 3 jdoe spark-users 0 2016-08-26 14:41 /user/jdoe/spark/canary_test/_SUCCESS
-rw-r--r-- 3 jdoe spark-users 4 2016-08-26 14:41
/user/jdoe/spark/canary_test/part-00000
-rw-r--r-- 3 jdoe spark-users 10 2016-08-26 14:41
/user/jdoe/spark/canary_test/part-00001
$ hdfs dfs -cat /user/jdoe/spark/canary_test/part-00000
ONE
$ hdfs dfs -cat /user/jdoe/spark/canary_test/part-00001
TWO
THREE
```

Fetching Spark 2 Maven Dependencies

The Maven coordinates are a combination of `groupId`, `artifactId` and `version`. The `groupId` and `artifactId` are the same as for the upstream Apache Spark project. For example, for `spark-core`, `groupId` is `org.apache.spark`, and `artifactId`

is `spark-core_2.11`, both the same as the upstream project. The version is different for the Cloudera packaging: see [Using the CDS Powered by Apache Spark Maven Repository](#) on page 20 for the exact name depending on which release you are using.

Adapting the Spark WordCount App for Spark 2

The following pom fragment shows how to access a CDS Powered by Apache Spark artifact from a Maven POM.

```
<dependency>
  <groupId>org.apache.spark</groupId>
  <artifactId>spark-core_2.11</artifactId>
  <version>2.0.0.cloudera2</version>
  <scope>provided</scope>
</dependency>
```

Use this dependency definition to update `pom.xml` for the example described in [Developing and Running a Spark WordCount Application](#). If you are using a different CDS version, see [Using the CDS Powered by Apache Spark Maven Repository](#) on page 20.

To account for changes in the Spark 2 API, before building the example, make the following updates to `com.cloudera.sparkwordcount.JavaWordCount`:

- Add `import java.util.Iterator;`
- Replace all instances of `Iterable` with `Iterator`.
- Perform the following replacements:
 - `return Arrays.asList(s.split(" "));` to `return Arrays.asList(s.split(" ")).iterator();`
 - `return chars;` to `return chars.iterator();`

Accessing the Spark 2 History Server

The Spark 2 history server is available on port 18089, rather than port 18088 as with the Spark 1 history server.

Troubleshooting for Spark 2

Troubleshooting for Spark mainly involves checking configuration settings and application code to diagnose performance and scalability issues.

Commercial support for GA version

Cloudera customers with commercial support can now use normal support channels for the CDS 2.0 component.

Error instantiating Hive metastore class

The latest release (release 2) addresses a Hive compatibility issue that affects CDH 5.10.1 and higher, CDH 5.9.2 and higher, CDH 5.8.5 and higher, and CDH 5.7.6 and higher. If you are using one of these CDH versions, you must upgrade to the Spark 2.0 release 2 parcel to avoid Spark 2 job failures when using Hive functionality.

When you encounter a problem due to the Hive compatibility issue, the error stack starts like this:

```
java.lang.RuntimeException: Unable to instantiate
  org.apache.hadoop.hive.q1.metadata.SessionHiveMetaStoreClient
    at org.apache.hadoop.hive.metastore.MetaStoreUtils.newInstance
      (MetaStoreUtils.java:1545)
    at org.apache.hadoop.hive.metastore
      .RetryingMetaStoreClient.<init>(RetryingMetaStoreClient.
```

The solution is to upgrade to 2.0 Release 2 or higher.

Wrong version of Python

Spark 2 requires Python 2.7 or higher. You might need to install a new version of Python on all hosts in the cluster, because some Linux distributions come with Python 2.6 by default. If the right level of Python is not picked up by default, set the `PYSPARK_PYTHON` and `PYSPARK_DRIVER_PYTHON` environment variables to point to the correct Python executable before running the `pyspark2` command.

API changes that are not backward-compatible

Between Spark 1.6 and Spark 2.0, some APIs have changed in ways that are not backward compatible. Recompile all applications to take advantage of Spark 2 capabilities. For any compilation errors, check if the corresponding function has changed in Spark 2, and if so, change your code to use the latest function name, parameters, and return type.

A Spark component does not work or is unstable

Certain components from the Spark ecosystem are explicitly not supported with CDS 2 Powered by Apache Spark. Check against the compatibility matrix for Spark to make sure the components you are using are all intended to work with CDS 2 Powered by Apache Spark and CDH.

Frequently Asked Questions about CDS Powered by Apache Spark

**Note:**

This Spark 2.0 documentation refers to the second release of CDS 2 Powered by Apache Spark. This component is generally available and is now supported on CDH 5.7 through CDH 5.10.

The latest release (release 2) addresses a Hive compatibility issue that affects CDH 5.10.1 and higher, CDH 5.9.2 and higher, CDH 5.8.5 and higher, and CDH 5.7.6 and higher. If you are using one of these CDH versions, you must upgrade to the Spark 2.0 release 2 parcel to avoid Spark 2 job failures when using Hive functionality.

This Frequently Asked Questions (FAQ) page covers general information about CDS Powered by Apache Spark, coexistence with Spark 1, and other questions that are relevant for early adopters of the latest Spark 2 features.

Running Spark 1 and Spark 2 Side-by-Side

The Spark 2 service does not conflict with Spark 1 if it is installed. The history server uses a different port. Spark 2 shares the Spark 1 shuffle service if already available, or installs the shuffle service if not.

Why doesn't feature or library XYZ work?

A number of features, components, libraries, and integration points from Spark 1.6 are not supported with CDS 2 Powered by Apache Spark. See [CDS Powered by Apache Spark Known Issues](#) on page 7 for details.

Appendix: Apache License, Version 2.0

SPDX short identifier: Apache-2.0

Apache License
Version 2.0, January 2004
<http://www.apache.org/licenses/>

TERMS AND CONDITIONS FOR USE, REPRODUCTION, AND DISTRIBUTION

1. Definitions.

"License" shall mean the terms and conditions for use, reproduction, and distribution as defined by Sections 1 through 9 of this document.

"Licensor" shall mean the copyright owner or entity authorized by the copyright owner that is granting the License.

"Legal Entity" shall mean the union of the acting entity and all other entities that control, are controlled by, or are under common control with that entity. For the purposes of this definition, "control" means (i) the power, direct or indirect, to cause the direction or management of such entity, whether by contract or otherwise, or (ii) ownership of fifty percent (50%) or more of the outstanding shares, or (iii) beneficial ownership of such entity.

"You" (or "Your") shall mean an individual or Legal Entity exercising permissions granted by this License.

"Source" form shall mean the preferred form for making modifications, including but not limited to software source code, documentation source, and configuration files.

"Object" form shall mean any form resulting from mechanical transformation or translation of a Source form, including but not limited to compiled object code, generated documentation, and conversions to other media types.

"Work" shall mean the work of authorship, whether in Source or Object form, made available under the License, as indicated by a copyright notice that is included in or attached to the work (an example is provided in the Appendix below).

"Derivative Works" shall mean any work, whether in Source or Object form, that is based on (or derived from) the Work and for which the editorial revisions, annotations, elaborations, or other modifications represent, as a whole, an original work of authorship. For the purposes of this License, Derivative Works shall not include works that remain separable from, or merely link (or bind by name) to the interfaces of, the Work and Derivative Works thereof.

"Contribution" shall mean any work of authorship, including the original version of the Work and any modifications or additions to that Work or Derivative Works thereof, that is intentionally submitted to Licensor for inclusion in the Work by the copyright owner or by an individual or Legal Entity authorized to submit on behalf of the copyright owner. For the purposes of this definition, "submitted" means any form of electronic, verbal, or written communication sent to the Licensor or its representatives, including but not limited to communication on electronic mailing lists, source code control systems, and issue tracking systems that are managed by, or on behalf of, the Licensor for the purpose of discussing and improving the Work, but excluding communication that is conspicuously marked or otherwise designated in writing by the copyright owner as "Not a Contribution."

"Contributor" shall mean Licensor and any individual or Legal Entity on behalf of whom a Contribution has been received by Licensor and subsequently incorporated within the Work.

2. Grant of Copyright License.

Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable copyright license to reproduce, prepare Derivative Works of, publicly display, publicly perform, sublicense, and distribute the Work and such Derivative Works in Source or Object form.

3. Grant of Patent License.

Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable (except as stated in this section) patent license to make, have made, use, offer to sell, sell, import, and otherwise transfer the Work, where such license applies only to those patent claims

licensable by such Contributor that are necessarily infringed by their Contribution(s) alone or by combination of their Contribution(s) with the Work to which such Contribution(s) was submitted. If You institute patent litigation against any entity (including a cross-claim or counterclaim in a lawsuit) alleging that the Work or a Contribution incorporated within the Work constitutes direct or contributory patent infringement, then any patent licenses granted to You under this License for that Work shall terminate as of the date such litigation is filed.

4. Redistribution.

You may reproduce and distribute copies of the Work or Derivative Works thereof in any medium, with or without modifications, and in Source or Object form, provided that You meet the following conditions:

1. You must give any other recipients of the Work or Derivative Works a copy of this License; and
2. You must cause any modified files to carry prominent notices stating that You changed the files; and
3. You must retain, in the Source form of any Derivative Works that You distribute, all copyright, patent, trademark, and attribution notices from the Source form of the Work, excluding those notices that do not pertain to any part of the Derivative Works; and
4. If the Work includes a "NOTICE" text file as part of its distribution, then any Derivative Works that You distribute must include a readable copy of the attribution notices contained within such NOTICE file, excluding those notices that do not pertain to any part of the Derivative Works, in at least one of the following places: within a NOTICE text file distributed as part of the Derivative Works; within the Source form or documentation, if provided along with the Derivative Works; or, within a display generated by the Derivative Works, if and wherever such third-party notices normally appear. The contents of the NOTICE file are for informational purposes only and do not modify the License. You may add Your own attribution notices within Derivative Works that You distribute, alongside or as an addendum to the NOTICE text from the Work, provided that such additional attribution notices cannot be construed as modifying the License.

You may add Your own copyright statement to Your modifications and may provide additional or different license terms and conditions for use, reproduction, or distribution of Your modifications, or for any such Derivative Works as a whole, provided Your use, reproduction, and distribution of the Work otherwise complies with the conditions stated in this License.

5. Submission of Contributions.

Unless You explicitly state otherwise, any Contribution intentionally submitted for inclusion in the Work by You to the Licensor shall be under the terms and conditions of this License, without any additional terms or conditions. Notwithstanding the above, nothing herein shall supersede or modify the terms of any separate license agreement you may have executed with Licensor regarding such Contributions.

6. Trademarks.

This License does not grant permission to use the trade names, trademarks, service marks, or product names of the Licensor, except as required for reasonable and customary use in describing the origin of the Work and reproducing the content of the NOTICE file.

7. Disclaimer of Warranty.

Unless required by applicable law or agreed to in writing, Licensor provides the Work (and each Contributor provides its Contributions) on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied, including, without limitation, any warranties or conditions of TITLE, NON-INFRINGEMENT, MERCHANTABILITY, or FITNESS FOR A PARTICULAR PURPOSE. You are solely responsible for determining the appropriateness of using or redistributing the Work and assume any risks associated with Your exercise of permissions under this License.

8. Limitation of Liability.

In no event and under no legal theory, whether in tort (including negligence), contract, or otherwise, unless required by applicable law (such as deliberate and grossly negligent acts) or agreed to in writing, shall any Contributor be liable to You for damages, including any direct, indirect, special, incidental, or consequential damages of any character arising as a result of this License or out of the use or inability to use the Work (including but not limited to damages for loss of goodwill, work stoppage, computer failure or malfunction, or any and all other commercial damages or losses), even if such Contributor has been advised of the possibility of such damages.

9. Accepting Warranty or Additional Liability.

Appendix: Apache License, Version 2.0

While redistributing the Work or Derivative Works thereof, You may choose to offer, and charge a fee for, acceptance of support, warranty, indemnity, or other liability obligations and/or rights consistent with this License. However, in accepting such obligations, You may act only on Your own behalf and on Your sole responsibility, not on behalf of any other Contributor, and only if You agree to indemnify, defend, and hold each Contributor harmless for any liability incurred by, or claims asserted against, such Contributor by reason of your accepting any such warranty or additional liability.

END OF TERMS AND CONDITIONS

APPENDIX: How to apply the Apache License to your work

To apply the Apache License to your work, attach the following boilerplate notice, with the fields enclosed by brackets "[]" replaced with your own identifying information. (Don't include the brackets!) The text should be enclosed in the appropriate comment syntax for the file format. We also recommend that a file or class name and description of purpose be included on the same "printed page" as the copyright notice for easier identification within third-party archives.

```
Copyright [yyyy] [name of copyright owner]

Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License.
```