

Cloudera Runtime 7.3.1

Upgrading Apache Spark

Date published: 2020-07-28

Date modified: 2024-12-10

CLOUDERA

<https://docs.cloudera.com/>

Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Upgrading Spark 2 to Spark 3 for Cloudera Private Cloud 7.3.1.....	4
Upgrading from 7.1.9 SP1.....	5
Upgrading Apache Spark 2.4.8 on Cloudera Private Cloud 7.1.9 SP1 to Spark 3.4.1 on 7.3.1.....	5
Upgrading Apache Spark 2.4.8 (with CDS 3.3.2) on 7.1.9 SP1 to Spark 3 on 7.3.1.....	6
Upgrade from Spark 3.3.2 (CDS).....	8
Upgrading from 7.1.8.....	9
Upgrading Apache Spark 2.4.8 on 7.1.8 to Spark 3 on 7.3.1.....	10
Upgrading Apache Spark 2.4.8 on 7.1.8 to Spark 3 on 7.3.1.....	11
Upgrading Apache Spark 2.4.8 (with CDS 3.3.x) on 7.1.8 to Spark 3 on 7.3.1.....	12
Upgrading Apache Spark 2.4.8 (with CDS 3.3.x) on 7.1.8 to Spark 3 on 7.3.1.....	13
Upgrading Apache Spark 3.3.x (CDS) on 7.1.8 to Spark 3 on 7.3.1.....	14
Upgrading from 7.1.7.....	15
Upgrading Apache Spark 2.4.7 on 7.1.7 SP3 to Spark 3 on 7.3.1.....	15
Upgrading Apache Spark 2.4.7 (with connectors) on 7.1.7 SP3 to Spark 3 on 7.3.1.....	17
Upgrading Apache Spark 2.4.7 (with CDS 3.2.3) on 7.1.7 SP3 to Spark 3 on 7.3.1.....	17
Upgrading Apache Spark 2.4.7 (with CDS 3.2.3) on 7.1.7 SP3 to Spark 3 on 7.3.1.....	19
Upgrading Apache Spark 3.2.3 on 7.1.7 SP3 to Spark 3 on 7.3.1.....	20
Migrating Spark applications.....	20
Java versions.....	21
Scala versions.....	21
Python versions.....	21
Spark commands.....	22
Spark connectors.....	22
Logging.....	22
Third-party libraries.....	23
Spark behavior changes.....	23
Apache Spark Migration guides.....	23
Spark 2 to Spark 3 workload refactoring.....	23
Unsupported features.....	26
Post-migration checklist.....	26
Benchmark testing.....	27
Troubleshooting.....	27

Upgrading Spark 2 to Spark 3 for Cloudera Private Cloud 7.3.1

Upgrading Apache Spark from version 2 to version 3 in Cloudera Private Cloud is a process that involves:

1. Intermediate in-place cluster upgrade tasks, due to different support for Spark 3 versions of connectors in Cloudera versions.
2. Intermediate Spark application migration tasks, due to minor or maintenance Spark version changes.
3. Pre-application migration tasks.
4. Application migration from Spark 2 to Spark 3, due to major Spark version changes.
5. Post-application migration tasks.
6. In-place cluster upgrade tasks.



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

7. Spark application migration tasks, due to minor or maintenance Spark version changes.



Note: The application migration tasks (see *Migrating Spark Applications*) are done on your application codebase.



Note: Depending on your environment and Spark applications, some of these steps are not necessary.

This table below provides links to the upgrade guides based on the version of Cloudera Private Cloud you're using, and the source versions of Spark in your environment.



Important:

If you're using Spark 2 with any of the following connectors, follow the upgrade steps in the relevant *Upgrade guide (with connectors)* pages:

- Oozie
- Solr
- Phoenix
- Hive Warehouse Connector
- Spark Schema Registry

Each upgrade guide contain all steps needed to upgrade Spark 2 to Spark 3 and perform the upgrade to Cloudera Private Cloud version 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Documentation
7.1.9 SP1	2.4.8	None	Upgrade guide
	2.4.8	3.3.2 (CDS parcel)	Upgrade guide
	None	3.3.2 (CDS parcel)	Upgrade guide
7.1.8	2.4.8	None	Upgrade guide Upgrade (with connectors) guide
	2.4.8	3.3.0 or 3.3.2 (CDS parcel)	Upgrade guide Upgrade guide (with connectors)
	None	3.3.0 or 3.3.2 (CDS parcel)	Upgrade guide
7.1.7 SP1	2.4.7	None	Upgrade guide
			Upgrade guide (with connectors)

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Documentation
	2.4.7	3.2.3 (CDS parcel)	Upgrade guide
	None	3.2.3 (CDS parcel)	Upgrade guide (with connectors)
	None	3.2.3 (CDS parcel)	Upgrade guide

Related Information

[Migrating Spark applications](#)

Upgrading from 7.1.9 SP1

Upgrading Apache Spark 2.4.8 on Cloudera Private Cloud 7.1.9 SP1 to Spark 3.4.1 on 7.3.1

The following steps will help you upgrade from Apache Spark 2.4.8 on Cloudera Private Cloud version 7.1.9 SP1 to Spark 3.4.1 on version 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version
7.1.9 SP1	2.4.8	none	7.3.1	3.4.1

Pre-application migration tasks

Procedure

Install the CDS parcel and Spark 3 services, as described in the [CDS parcel documentation](#). A short overview of the process is as follows:

- Check that all the software prerequisites are satisfied.
- In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
- Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.
- Add the Spark 3 service to your cluster.
- Return to the Home page.
- Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

Spark application migration (from Spark 2 to Spark 3)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 2.4.8 to 3.3.2.

- Check the supported Java versions.
- Check the supported Scala version.
- Check the supported Python versions.
- Account for changed or versioned Spark commands in your code. (spark-submit, pyspark, etc.)
- Check supported versions for Spark connectors.
- Check the logging library used in your code.
- Check the compatibility of 3rd-party libraries used in your code.
- Check Spark behavior changes and refactor your code.

Post-application migration tasks

Procedure

1. Stop the Livy (Livy for Spark 2) and Spark 2 (SPARK_ON_YARN) services.
2. Delete the Spark 2 and Livy for Spark 2 services.
3. Move Spark 2 event logs to the Spark 3 event logs directory.

In-place cluster upgrade

Before you begin



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).
2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Spark application migration (from Spark 3.x to Spark 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.2 to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading Apache Spark 2.4.8 (with CDS 3.3.2) on 7.1.9 SP1 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.8 on Cloudera Private Cloud 7.1.9 SP1 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ¹
7.1.9 SP1	2.4.8	3.3.2 (CDS)	7.3.1	3.4.1	no

Application migration tasks from Spark 2 to Spark 3

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 2.4.8 to 3.3.2.

- a) Check the supported Java versions.
- b) Check the supported Scala version.
- c) Check the supported Python versions.
- d) Account for changed or versioned Spark commands in your code. (spark-submit, pyspark, etc.)
- e) Check supported versions for Spark connectors.
- f) Check the logging library used in your code.
- g) Check the compatibility of 3rd-party libraries used in your code.
- h) Check Spark behavior changes and refactor your code.

Post-application migration tasks

Procedure

1. Stop the Livy (Livy for Spark 2) and Spark 2 (SPARK_ON_YARN) services.
2. Delete the Spark 2 and Livy for Spark 2 services.
3. Move Spark 2 event logs to the Spark 3 event logs directory.

In-place cluster upgrade

Before you begin



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Spark application migration (from Spark 3.x to Spark 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.2 to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrade from Spark 3.3.2 (CDS)

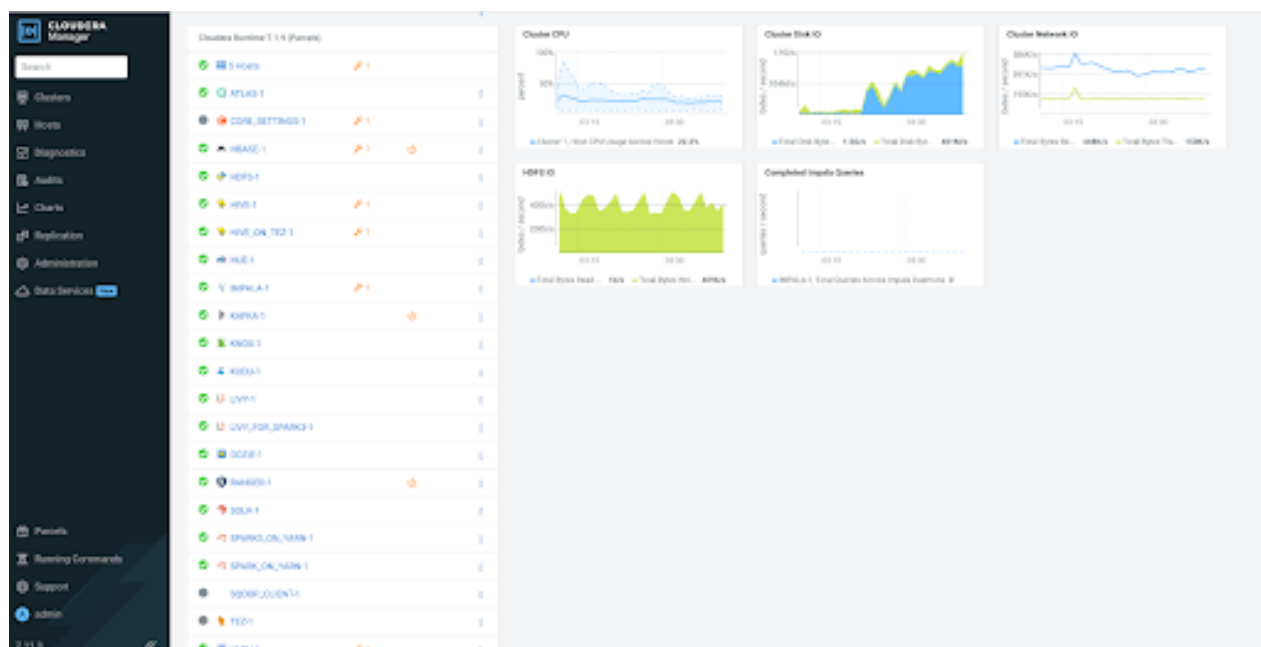
About this task

The following steps will help you upgrading from Apache Spark 2.4.8 on Cloudera Private Cloud 7.1.9 SP1 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ²
7.1.9 SP1	none	3.3.2 (CDS)	7.3.1	3.4.1	no

Before you begin

Initial state:



Procedure

In-place cluster upgrade

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Spark application migration

3. Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.2 to 3.4.1
 - a) Refactor your Spark application code.

In-place cluster upgrade

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Spark application migration (from Spark 3.x to Spark 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.2 to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading from 7.1.8

Upgrading Apache Spark 2.4.8 on 7.1.8 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.8 on Cloudera Private Cloud 7.1.8 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ³
7.1.8	2.4.8	none	7.3.1	3.4.1	no

Pre-application migration tasks

Procedure

Install the CDS parcel and Spark 3 services, as described in the [CDS parcel documentation](#). A short overview of the process is as follows:

- Check that all the software prerequisites are satisfied.
- In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
- Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.
- Add the Spark 3 service to your cluster.
- Return to the Home page.
- Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

Application migration tasks (Spark 2 to Spark 3)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 2.4.8 to 3.3.0

- Check the supported Java versions.
- Check the supported Scala version.
- Check the supported Python versions.
- Account for changed or versioned Spark commands in your code. (spark-submit, pyspark, etc.)
- Check supported versions for Spark connectors.
- Check the logging library used in your code.
- Check the compatibility of 3rd-party libraries used in your code.
- Check Spark behavior changes and refactor your code.

Post-application migration tasks

Procedure

- Stop the Livy (Livy for Spark 2) and Spark 2 (SPARK_ON_YARN) services.
- Delete the Spark 2 and Livy for Spark 2 services.
- Move Spark 2 event logs to the Spark 3 event logs directory.

In-place cluster upgrade

Before you begin



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Application migration tasks (Spark 3.x to Spark 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.x to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading Apache Spark 2.4.8 on 7.1.8 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.8 on Cloudera Private Cloud 7.1.8 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ⁴
7.1.8	2.4.8	none	7.3.1	3.4.1	yes

Intermediate in-place cluster upgrade

Procedure

1. Upgrade your cluster to Cloudera Private Cloud 7.1.9 SP1 with CDS 3.3.2.
 - a) Review the Supported Upgrade Paths for your upgrade.
 - b) Review the Requirements and Supported Versions for your upgrade
 - c) Review the Release Notes for the version of Cloudera Private Cloud Base you are upgrading to.
 - d) Review the Cloudera Upgrade/Migrate Troubleshooting Articles.
 - e) Gather information on your deployment.
 - f) Plan how and when to begin your upgrade.
 - g) If necessary, Upgrade the JDK.
 - h) If necessary, Upgrade the Operating System.
 - i) Upgrade Cloudera Manager.
 - j) Use Cloudera Manager to upgrade from lower version of Cloudera Runtime to a higher version of Cloudera Runtime.
 - k) Check that all the software prerequisites are satisfied for CDS 3.3.2.
 - l) In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
 - m) Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.
 - n) Add the Spark 3 service to your cluster.
 - o) Return to the Home page.
 - p) Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

For more information on in-place cluster upgrades and installing the CDS parcel, see:

- [Upgrading Cloudera Private Cloud Base to a higher version](#)
- [CDS parcel documentation](#)

2. Once complete, continue the upgrade with the steps described in [Upgrading Apache Spark 2 \(with 3.3.2 CDS\) on 7.1.9 SP1 to Spark 3 on 7.3.1](#).

Upgrading Apache Spark 2.4.8 (with CDS 3.3.x) on 7.1.8 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.8 (with CDS 3.3.0) on Cloudera Private Cloud 7.1.8 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ⁵
7.1.8	2.4.8	3.3.0 or 3.3.2 (CDS)	7.3.1	3.4.1	no

Application migration tasks (from Spark 2 to 3)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 2.4.8 to 3.3.0

- a) Check the supported Java versions.
- b) Check the supported Scala version.
- c) Check the supported Python versions.
- d) Account for changed or versioned Spark commands in your code. (spark-submit, pyspark, etc.)
- e) Check supported versions for Spark connectors.
- f) Check the logging library used in your code.
- g) Check the compatibility of 3rd-party libraries used in your code.
- h) Check Spark behavior changes and refactor your code.

Post-application migration tasks

Procedure

1. Stop the Livy (Livy for Spark 2) and Spark 2 (SPARK_ON_YARN) services.
2. Delete the Spark 2 and Livy for Spark 2 services.
3. Move Spark 2 event logs to the Spark 3 event logs directory.

In-place cluster upgrade

Before you begin



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Application migration tasks (from Spark 3.x to 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.x to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading Apache Spark 2.4.8 (with CDS 3.3.x) on 7.1.8 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.8 (with CDS 3.3.x) on Cloudera Private Cloud 7.1.8 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ⁶
7.1.8	2.4.8	3.3.0 or 3.3.2 (CDS)	7.3.1	3.4.1	yes

Intermediate in-place cluster upgrade

Procedure

1. Upgrade your cluster to Cloudera Private Cloud 7.1.9 SP1 with CDS 3.3.2.
 - a) Review the Supported Upgrade Paths for your upgrade.
 - b) Review the Requirements and Supported Versions for your upgrade
 - c) Review the Release Notes for the version of Cloudera Private Cloud Base you are upgrading to.
 - d) Review the Cloudera Upgrade/Migrate Troubleshooting Articles.
 - e) Gather information on your deployment.
 - f) Plan how and when to begin your upgrade.
 - g) If necessary, Upgrade the JDK.
 - h) If necessary, Upgrade the Operating System.
 - i) Upgrade Cloudera Manager.
 - j) Use Cloudera Manager to upgrade from lower version of Cloudera Runtime to a higher version of Cloudera Runtime.
 - k) Check that all the software prerequisites are satisfied for CDS 3.3.2.
 - l) In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
 - m) Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.
 - n) Add the Spark 3 service to your cluster.
 - o) Return to the Home page.
 - p) Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

For more information on in-place cluster upgrades and installing the CDS parcel, see:

- [Upgrading Cloudera Private Cloud Base to a higher version](#)
- [CDS parcel documentation](#)

2. Once complete, continue the upgrade with the steps described in [Upgrading Apache Spark 2 \(with 3.3.2 CDS\) on 7.1.9 SP1 to Spark 3 on 7.3.1](#).

Upgrading Apache Spark 3.3.x (CDS) on 7.1.8 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 3.3.0 or 3.3.2 (CDS) on Cloudera Private Cloud 7.1.8 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ⁷
7.1.8	none	3.3.0 or 3.3.2 (CDS)	7.3.1	3.4.1	no

In-place cluster upgrade

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).
2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Application migration tasks (Spark 3.x to Spark 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.3.x to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading from 7.1.7

Upgrading Apache Spark 2.4.7 on 7.1.7 SP3 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.7 on Cloudera Private Cloud 7.1.7 SP3 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ⁸
7.1.7 SP3	2.4.7	none	7.3.1	3.4.1	no

Pre-application migration tasks

Procedure

Install the CDS parcel and Spark 3 services, as described in the [CDS parcel documentation](#). A short overview of the process is as follows:

- a) Check that all the software prerequisites are satisfied.
- b) In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
- c) Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.

- d) Add the Spark 3 service to your cluster.
- e) Return to the Home page.
- f) Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

Application migration tasks (Spark 2 to 3)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 2.4.7 to 3.2.3

- a) Check the supported Java versions.
- b) Check the supported Scala version.
- c) Check the supported Python versions.
- d) Account for changed or versioned Spark commands in your code. (spark-submit, pyspark, etc.)
- e) Check supported versions for Spark connectors.
- f) Check the logging library used in your code.
- g) Check the compatibility of 3rd-party libraries used in your code.
- h) Check Spark behavior changes and refactor your code.

Post-application migration tasks

Procedure

1. Stop the Livy (Livy for Spark 2) and Spark 2 (SPARK_ON_YARN) services.
2. Delete the Spark 2 and Livy for Spark 2 services.
3. Move Spark 2 event logs to the Spark 3 event logs directory.

In-place cluster upgrade

Before you begin



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

Procedure

1. Upgrade the Cloudera Manager Server to 7.13.1.0
 - a) Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - b) Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - c) Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

2. Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Application migration tasks (Spark 3.x to 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.2.3 to 3.4.1

- a) Refactor your Spark application code.

Fianl steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading Apache Spark 2.4.7 (with connectors) on 7.1.7 SP3 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.7 on Cloudera Private Cloud 7.1.7 SP3 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ⁹
7.1.7 SP3	2.4.7	none	7.3.1	3.4.1	yes

Intermediate in-place cluster upgrade

Procedure

1. Upgrade your cluster to Cloudera Private Cloud 7.1.9 SP1 with CDS 3.3.2.
 - a) Review the Supported Upgrade Paths for your upgrade.
 - b) Review the Requirements and Supported Versions for your upgrade
 - c) Review the Release Notes for the version of Cloudera Private Cloud Base you are upgrading to.
 - d) Review the Cloudera Upgrade/Migrate Troubleshooting Articles.
 - e) Gather information on your deployment.
 - f) Plan how and when to begin your upgrade.
 - g) If necessary, Upgrade the JDK.
 - h) If necessary, Upgrade the Operating System.
 - i) Upgrade Cloudera Manager.
 - j) Use Cloudera Manager to upgrade from lower version of Cloudera Runtime to a higher version of Cloudera Runtime.
 - k) Check that all the software prerequisites are satisfied for CDS 3.3.2.
 - l) In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
 - m) Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.
 - n) Add the Spark 3 service to your cluster.
 - o) Return to the Home page.
 - p) Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

For more information on in-place cluster upgrades and installing the CDS parcel, see:

- [Upgrading Cloudera Private Cloud Base to a higher version](#)
- [CDS parcel documentation](#)

2. Once complete, continue the upgrade with the steps described in [Upgrading Apache Spark 2 \(with 3.3.2 CDS\) on 7.1.9 SP1 to Spark 3 on 7.3.1](#).

Upgrading Apache Spark 2.4.7 (with CDS 3.2.3) on 7.1.7 SP3 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.7 (with CDS 3.2.3) on Cloudera Private Cloud 7.1.7 SP3 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ¹⁰
7.1.7 SP3	2.4.7	3.2.3 (CDS)	7.3.1	3.4.1	no

Application migration tasks (Spark 2 to 3)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 2.4.7 to 3.2.3

- Check the supported Java versions.
- Check the supported Scala version.
- Check the supported Python versions.
- Account for changed or versioned Spark commands in your code. (spark-submit, pyspark, etc.)
- Check supported versions for Spark connectors.
- Check the logging library used in your code.
- Check the compatibility of 3rd-party libraries used in your code.
- Check Spark behavior changes and refactor your code.

Post-application migration tasks

Procedure

- Stop the Livy (Livy for Spark 2) and Spark 2 (SPARK_ON_YARN) services.
- Delete the Spark 2 and Livy for Spark 2 services.
- Move Spark 2 event logs to the Spark 3 event logs directory.

In-place cluster upgrade

Before you begin



Important: Cluster upgrade to version 7.3.1 is only allowed if all Spark 2 applications have been migrated to Spark 3 and tested in the **source cluster**.

Procedure

- Upgrade the Cloudera Manager Server to 7.13.1.0
 - Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

- Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Application migration tasks (Spark 3.x to 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.2.3 to 3.4.1

- a) Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

1. Check the status of your clusters.
2. Perform benchmark testing on your applications. See [Spark Application Migration](#).

Upgrading Apache Spark 2.4.7 (with CDS 3.2.3) on 7.1.7 SP3 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 2.4.7 (with CDS 3.2.3) on Cloudera Private Cloud 7.1.7 SP3 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ¹¹
7.1.7 SP3	2.4.7	3.2.3 (CDS)	7.3.1	3.4.1	yes

Intermediate in-place cluster upgrade

Procedure

1. Upgrade your cluster to Cloudera Private Cloud 7.1.9 SP1 with CDS 3.3.2.
 - a) Review the Supported Upgrade Paths for your upgrade.
 - b) Review the Requirements and Supported Versions for your upgrade
 - c) Review the Release Notes for the version of Cloudera Private Cloud Base you are upgrading to.
 - d) Review the Cloudera Upgrade/Migrate Troubleshooting Articles.
 - e) Gather information on your deployment.
 - f) Plan how and when to begin your upgrade.
 - g) If necessary, Upgrade the JDK.
 - h) If necessary, Upgrade the Operating System.
 - i) Upgrade Cloudera Manager.
 - j) Use Cloudera Manager to upgrade from lower version of Cloudera Runtime to a higher version of Cloudera Runtime.
 - k) Check that all the software prerequisites are satisfied for CDS 3.3.2.
 - l) In the **Cloudera Manager Admin Console**, add the CDS parcel repository to the Remote Parcel Repository URLs in Parcel Settings.
 - m) Download the CDS parcel, distribute it to the hosts in your cluster, and activate it.
 - n) Add the Spark 3 service to your cluster.
 - o) Return to the Home page.
 - p) Click the stale configuration icon to launch the **Stale Configuration wizard** and restart the necessary services.

For more information on in-place cluster upgrades and installing the CDS parcel, see:

 - [Upgrading Cloudera Private Cloud Base to a higher version](#)
 - [CDS parcel documentation](#)
2. Once complete, continue the upgrade with the steps described in [Upgrading Apache Spark 2 \(with 3.3.2 CDS\) on 7.1.9 SP1 to Spark 3 on 7.3.1](#).

Upgrading Apache Spark 3.2.3 on 7.1.7 SP3 to Spark 3 on 7.3.1

The following steps will help you upgrading from Apache Spark 3.2.3 (CDS) on Cloudera Private Cloud 7.1.7 SP3 to Spark 3.4.1 on 7.3.1.

Source cluster version	Source cluster Spark 2 version	Source cluster Spark 3 version	Target cluster version	Target cluster Spark 3 version	Spark 2 used with connectors ¹²
7.1.7 SP3	none	3.2.3 (CDS)	7.3.1	3.4.1	no

In-place cluster upgrade

Procedure

- Upgrade the Cloudera Manager Server to 7.13.1.0
 - Back up the Cloudera Manager server databases, working directories, and several other entities. These backups can be used to restore your Cloudera Manager deployment if there are problems during the upgrade.
 - Upgrade the Cloudera Manager server software on the Cloudera Manager host using package commands from the command line (for example, yum on RHEL systems). Cloudera Manager automates much of this process and is recommend for upgrading and managing your CDH/Cloudera Runtime clusters.
 - Upgrade the Cloudera Manager agent software on all cluster hosts. The Cloudera Manager upgrade wizard can upgrade the agent software (and, optionally, the JDK), or you can install the agent and JDK software manually. The CDH or Cloudera Runtime software is not upgraded during this process.

For more information, see [Upgrading Cloudera Manager 7](#).

- Cloudera Manager handles deactivation of the SPARK3 parcel. The parcel itself is not removed but deactivated.

Application migration tasks (Spark 3.x to 3.4.1)

Procedure

Follow the [Spark application migration documentation](#) to migrate your Apache Spark Applications from version 3.2.3 to 3.4.1

- Refactor your Spark application code.

Final steps

Procedure

After the upgrade and application migration are complete:

- Check the status of your clusters.
- Perform benchmark testing on your applications. See [Spark Application Migration](#).

Migrating Spark applications

How to refactor Spark 2 workloads to Spark 3 during the upgrade/migration process due to the removal of Spark 2 in Cloudera Private Cloud.

Introduction

The purpose of this document is to gather all information required to carry out a Spark application migration between different versions.

The necessary set of steps largely depends on the source and target Spark versions, while major version changes require considerable effort, minor and maintenance version changes mostly require only small config or no adjustments.

Major version migration

Migration between major versions requires considerable effort and taking into account many factors.

This documentation focuses on migrating applications from Spark 2 to Spark 3, the two major versions currently supported by Cloudera on different versions of Cloudera Private Cloud.

Java versions

Cloudera currently supports 3 major JDK versions in general:

- 8
- 11
- 17

Refer to *Support Matrix* for a list of supported versions of Java.



Important: For Java 11, setting `-Dio.netty.tryReflectionSetAccessible=true` is required for the Apache Arrow library. This prevents the `java.lang.UnsupportedOperationException: sun.misc.Unsafe or java.nio.DirectByteBuffer.(long, int) not available` error when Apache Arrow uses Netty internally.

Related Information

[Cloudera Support Matrix](#)

Scala versions

As Cloudera only supports only Spark 2 applications compiled with Scala **2.11**, and Spark 3 applications with Scala **2.12**, a major version change always require:

1. Spark Scala applications to be recompiled with Scala 2.12,
2. adjusting the dependencies to use Spark 3 version binaries provided by Cloudera in the *public maven repository* and the Scala 2.12 version of third-party libraries.

Scala version changes can also require source code changes, for which see the *Scala documentation*.

Related Information

[Cloudera public maven repositories](#)

[Cloudera Private Cloud | Using the Cloudera Runtime Maven repository](#)

[Scala documentation | Scala](#)

Python versions

The supported versions of Python can change between Spark versions. Refer to the table below for details on the supported Python versions for each Spark version, and follow the *Python documentation* to adjust your application.

Spark version	Minimum supported Python version	Maximum supported Python version
3.5.0	3.8	3.11
3.4.0	3.7	3.11
3.3.2	3.7	3.10
3.3.0	3.7	3.10

Spark version	Minimum supported Python version	Maximum supported Python version
3.2.3	3.6	3.9
2.4.8	2.7/3.4	3.7

Related Information

[Python documentation | Python](#)

Spark commands

Cloudera supports multiple versions of Spark, depending on the version of Cloudera Private Cloud clusters. The general (unversioned) Apache Spark commands (`spark-submit`, `pyspark`, etc.) can point to different versions based on the cluster version.

1. The original commands always point to the earliest available version of Spark in the distribution.

For example, the `spark-submit` command points to **Spark 2 in version 7.1.9 SP1**, but points to **Spark 3 in version 7.3.1**.

2. Other available Spark 3 versions can be used via versioned commands.

For example, the `spark3-submit` command points to Spark 3 in all versions.

Spark connectors

Spark 3 supports certain Spark connectors from certain versions.

If Spark 2 connectors are used, please take the connectors into account when choosing the minimum Cloudera Private Cloud version you need to upgrade to when migrating a Spark application to a higher version.

- *Hive Warehouse Connector* for Spark 3 is supported from:
 - Cloudera Private Cloud version 7.1.8
 - Cloudera Distribution of Spark (CDS) version 3.3.0
- *HBase connector* for Spark 3 is supported from:
 - Cloudera Private Cloud version 7.1.7
 - Cloudera Distribution of Spark (CDS) version 3.2
- *Phoenix* connector for Spark 3 is supported from:
 - Cloudera Private Cloud version 7.1.8
 - Cloudera Distribution of Spark (CDS) version 3.3.0
- *Oozie* for Spark 3 is supported from:
 - Cloudera Private Cloud version 7.1.9
 - Cloudera Distribution of Spark (CDS) version 3.3.2
- *Solr* for Spark 3 is supported from:
 - Cloudera Private Cloud version 7.1.9
 - Cloudera Distribution of Spark (CDS) version 3.3.2
- *Spark Schema Registry* is supported from:
 - Cloudera Private Cloud version 7.1.9 SP1
 - Cloudera Distribution of Spark (CDS) version 3.3.2

Logging

Since Apache Spark 3 has transitioned from log4j to log4j2, you need to adjust the logging library and/or logging configuration used in your application.



Important: The log4j and reload4j runtime libraries are no longer included in the classpath by default. They have been replaced with the analogous log4j2 libraries.

Third-party libraries

When migrating between versions, ensure that your 3rd-party runtime dependencies align with the Spark versions.

Spark behavior changes

As Apache Spark evolves, its behavior can change between major and minor versions, but many times legacy configurations are added to support the old behavior.

As configurations can be defined at multiple levels, restoring the old behavior might require changing the application itself, the application starting commands/scripts, and/or the default Spark configurations defined on the cluster.

Apache Spark Migration guides

The comprehensive guide of behavior changes between versions are available in the Apache Spark Migration Guides.

Always refer to the following documents to ensure that your Spark application remains compatible with newer Spark versions:

- [Migration Guide: Spark Core](#)
- [Migration Guide: SQL, Datasets and DataFrame](#)
- [Structured Streaming Programming Guide](#)
- [Machine Learning Library \(MLlib\) Guide](#)
- [Upgrading PySpark](#)

Spark 2 to Spark 3 workload refactoring

The following list summarizes the most important behavior changes from Spark 2 to Spark 3, and gives examples on how to refactor the Spark 2 application to become Spark 3 compatible.

The list is not exhaustive, refer to the *Apache Spark Migration guides* for the complete list.

Spark Core

Spark Core language/syntactic-level changes

Spark 2	Spark 3	Refactor action
TaskContext.isRunningLocally	Deprecated method, removed.	Remove TaskContext.isRunningLocally if used in code.
ShuffleBytesWritten and shuffleRecordsWritten (ShuffleWriteMetrics class)	bytesWritten and recordsWritten (org.apache.spark.status.api.v1.OutputMetrics class)	Use bytesWritten and recordsWritten, available in class org.apache.spark.status.api.v1.OutputMetrics.
org.apache.spark Class Accumulator	org.apache.spark.util.AccumulatorV2	Replace org.apache.spark.Accumulator with org.apache.spark.util.AccumulatorV2.
For non-struct types, (e.g. int, string, array, Dataset.groupByKey) results in a grouped dataset with key attribute is wrongly named as value.	For non-struct types (e.g. int, string, array, Dataset.groupByKey) results to a grouped dataset with key attribute is named as key.	Refactor the value attribute used in logic to key. To preserve the old behavior, set spark.sql.legacy.dataset.nameNonStructGroupingKeyAsValue to false.


Spark SQL

Spark SQL language/syntactic-level changes

Spark 2	Spark 3	Refactor action
Path option is overwritten if one path parameter is passed to <code>DataFrameReader.load()</code> , <code>DataFrameWriter.save()</code> , <code>DataStreamReader.load()</code> , or <code>DataStreamWriter.start()</code> .	Path option cannot coexist when the following methods are called with path parameter(s): <code>DataFrameReader.load()</code> , <code>DataFrameWriter.save()</code> , <code>DataStreamReader.load()</code> , or <code>DataStreamWriter.start()</code> .	Remove the path option if it's the same as the path parameter, or add it to the <code>load()</code> parameter if you do want to read multiple paths. To ignore this check, set <code>spark.sql.legacy.pathOptionBehavior.enabled</code> to <code>true</code> .
<code>count(tblName.*)</code> works.	An exception is thrown if <code>count(tblName.*)</code> is used for getting the number of records in the table.	Refactor the code to use <code>count(*)</code> , or expand the columns manually. (Example: <code>count(col1, col2)</code> .) To restore the old behavior, set <code>spark.sql.legacy.allowStarWithSingleTableIdentifierInCount</code> to <code>true</code> .

Spark SQL configuration-level changes

Spark 2	Spark 3	Refactor action
SET command works for SparkConf entries.	<code>AnalysisException</code> error is thrown if SET command is used to modify the SparkConf entries.	Remove SET commands for SparkConf entries from your code. You can enter SparkConf values at the cluster level by entering them in the cluster's Spark configuration and restarting the cluster. To disable the check, set <code>spark.sql.legacy.setCommandRejectsSparkCoreConfs</code> to <code>false</code> .
The second argument of <code>date_add</code> function (<code>num_days</code>) can be a fraction, as it gets casted to <code>Int</code> internally.	The second argument of <code>date_add</code> function (<code>num_days</code>). If an integer is not provided, an <code>AnalysisException</code> is thrown.	Make sure that in code always integer is passed as the second argument to <code>date_add</code> and <code>date_subtract</code> function.
Fractional and string types are allowed in <code>percentile_approx</code> third argument i.e. accuracy, as it gets casted to <code>Int</code> internally.	<code>percentile_approx</code> third argument accuracy can only be integer. If an integer is not provided, an <code>AnalysisException</code> is thrown.	Make sure that in code always integer is passed as the third argument to <code>percentile_approx</code> function.
Hash expressions can be applied on <code>MapType</code> elements.	Hash expressions are prohibited on <code>MapType</code> elements.	If hash expression is applied on map type, refactor the code to remove it, OR set <code>spark.sql.legacy.allowHashOnMapType</code> to <code>true</code> .
a map can be created with duplicate keys via built-in functions like <code>CreateMap</code> , <code>StringToMap</code> , <code>map_from_arrays</code> etc.	Spark throws <code>RuntimeException</code> when duplicated keys are found in <code>Map</code> . Users may still read map values with duplicate keys from data sources which do not enforce it (for example, Parquet).	If duplicate keys are passed into built in functions to create a map then try to remove duplicate keys OR set <code>spark.sql.mapKeyDedupPolicy</code> to <code>LAST_WIN</code> , the map keys are deduplicated.
the resulting date is adjusted in <code>add_months</code> , when the original date is a last day of months. For example, adding a month to 2019-02-28 results in 2019-03-31.	the <code>add_months</code> function does not adjust the resulting date to a last day of month if the original date is a last day of months. For example, adding a month to 2019-02-28 results in 2019-03-28.	Adjust the code according to logic if required.
multiple from-to units is allowed in <code>Interval</code> literal.	multiple from-to units <code>Interval</code> literal is not allowed.	Remove multiple from-to units is allowed in <code>Interval</code> literal. Adjust the code according to logic if required.
Dataset query success if it contains ambiguous column reference that is caused by self join.	Dataset query fails if it contains ambiguous column reference that is caused by self join. This is because Spark cannot resolve Dataset column references that point to tables being self joined, and <code>df1("a")</code> is exactly the same as <code>df2("a")</code> in Spark.	Use aliases. For example: <code>df2.as("purchases").join(df1.as("devices"), col("devices.key1") === col("purchases.key2")).show()</code>
invalid time zone ids are silently ignored and replaced by GMT timezone.	invalid time zone ids are rejected, and Spark throws <code>java.time.DateTimeException</code> .	rectify to correct Zone ID.

Spark 2	Spark 3	Refactor action
for Parsing and formatting of timestamp and date strings, <code>java.text.SimpleDateFormat</code> is used for timestamp/date string conversions, and the supported patterns are described in <code>SimpleDateFormat</code> .	<code>DateTimeFormatter</code> under the hood or Parsing and formatting of timestamp and date strings. Strict checking of Input is performed.	Refactor the code to correct pattern matching for Input OR set <code>spark.sql.legacy.timeParserPolicy</code> to <code>LEGACY</code> to restore the behavior, OR set it to <code>CORRECTED</code> and treat it as an invalid datetime string.
datetime pattern letter F is aligned to week of month that represents the concept of the count of weeks within the month where weeks start on a fixed day-of-week.	datetime pattern letter F is aligned to day of week in month that represents the concept of the count of days within the period of a week where the weeks are aligned to the start of the month.	Refactor the code to accommodate new behavior of pattern F.
<code>SparkContext</code> can be created in executors.	an exception will be thrown when creating <code>SparkContext</code> in executor.	Refactor the code to remove the creation of Spark context OR allow it by setting the configuration <code>spark.executor.allowSparkContext</code> when creating <code>SparkContext</code> in executors.
<code>TRANSFORM</code> operator can support alias in inputs.	<code>TRANSFORM</code> operator can't support alias in inputs.	Refactor the code to remove aliases from Inputs.
Loading and saving of timestamps from and to Parquet files does not fail if the timestamps are before 1900-01-01 00:00:00Z.	Loading and saving of timestamps from and to Parquet files fails if the timestamps are before 1900-01-01 00:00:00Z.	Ensure that Input reads do not contain timestamps before 1900-01-01 00:00:00Z. Alternatively, set <code>spark.sql.parquet.int96RebaseModeInWrite</code> to <code>CORRECTED</code> to write the datetime values as it is.
The <code>Char(n)</code> type handled inconsistently, depending on whether the table is partitioned or not.	In upstream Spark 3 the <code>spark.sql.legacy.charVarcharAsString</code> config was introduced, but does not solve all incompatibilities.	A new configuration <code>spark.cloudera.legacy.charVarcharLegacyPadding</code> is introduced in Cloudera to keep the full compatibility.  Note: Although this workaround is available, using <code>CHAR</code> is discouraged.
The Row field names are sorted alphabetically when constructing with named arguments for Python versions 3.6 and above.	The Row field names are no longer sorted alphabetically.	To enable sorted fields by default as in Spark 2.4, set the environment variable <code>PYSPARK_ROW_FIELD_SORTING_ENABLED</code> to true for both the executors and the driver.

Spark SQL property-level changes

Spark 2	Spark 3	Refactor action
When there is nested CTE with a conflicting name, outer CTE definitions take precedence.	When there is nested CTE with conflicting name, Spark throws an <code>AnalysisException</code> by default, and forces users to choose the specific substitution order they wanted. If the value of <code>spark.sql.legacy.ctePrecedencePolicy</code> is set to <code>CORRECTED</code> (recommended), inner CTE definitions take precedence over outer definitions. If the value of <code>spark.sql.legacy.ctePrecedencePolicy</code> is set to <code>LEGACY</code> , outer CTE definitions take precedence over inner definitions.	Set <code>spark.sql.legacy.ctePrecedencePolicy</code> to <code>CORRECTED</code> .
Type conversions during table insertion are allowed as long as they are valid Cast.	The type coercion is performed as per the ANSI SQL standard.	Ensure the type coercion is performed as per the ANSI SQL standard. Alternatively, set <code>spark.sql.storeAssignmentPolicy</code> to <code>Legacy</code> to restore previous behavior.

Spark storage location configuration changes

To execute workloads in Cloudera Public Cloud, modify the local data storage locations to cloud storage (for example, from HDFS to S3 bucket).

The following example shows a sample workload, with the modified data location **highlighted in bold**.

Spark 3.2 (HDFS)	Spark 3.2 (S3)
spark-shell	spark-shell
<pre>scala> spark.sql("CREATE TABLE IF NOT EXISTS default.sales_spark_2(Region string, Country string,Item_Type string,Sales_Channel string,Order_Priority string,Order_Date date,Order_ID int,Ship_Date date,Units_sold string,Unit_Price string,Unit_cost string,Total_revenue string,Total_Cost string,Total_Profit string) row format delimited fields terminated by ','") scala> spark.sql("load data local inpath '/tmp/sales.csv' into table default.sales_spark_3")</pre>	<pre>scala> spark.sql("CREATE TABLE IF NOT EXISTS default.sales_spark_2(Region string, Country string,Item_Type string,Sales_Channel string,Order_Priority string,Order_Date date,Order_ID int,Ship_Date date,Units_sold string,Unit_Price string,Unit_cost string,Total_revenue string,Total_Cost string,Total_Profit string) row format delimited fields terminated by ','") scala> spark.sql("load data inpath 's3://[*** BUCKET ***/sales.csv' into table default.sales_spark_3")</pre>
<pre>scala> spark.sql("select count(*) from default.sales_spark_3").show()</pre>	<pre>scala> spark.sql("select count(*) from default.sales_spark_3").show()</pre>

Unsupported features

Unsupported Spark 3 features in Cloudera Private Cloud.

ZSTD compression in ORC data source

spark.sql.orc.compression.codec config doesn't accept zstd value.

Apache Jira:[SPARK-33978](#)

spark.hadoopRDD.ignoreEmptySplits

Causes issues in HBase TableInputFormat.

Apache Jira:[SPARK-34809](#)

LDAP authentication for livy-server

Open CVEs in Apache Directory Server dependency, LDAP based authentication is not supported in Livy Server.

Apache Jira:[LIVY-356](#)

Thrift ldap authentication, based on ldapurl, basedn, domain

Open CVEs in Apache Directory Server dependency, LDAP based authentication is not supported in Livy Thrift Server.

Apache Jira:[LIVY-678](#)

For more information, see *Unsupported Apache Spark Features*.

Related Information

[Unsupported Apache Spark Features](#)

Post-migration checklist

Benchmark testing

After all post-migration configurations are performed, perform benchmark testing on the new Apache Spark version.

Troubleshooting

Troubleshoot failed or slow performing workloads by analyzing the application event and driver logs, and fine tune the workloads for better performance.