

StreamSets Data Collector 1.5.1.2 Release Notes

August 4, 2016

New Features and Enhancements

We're happy to announce a new version of StreamSets Data Collector.

This version includes several bug fixes, described below.

Feel free to check out the [Documentation](#) for this release.

Upgrade

You can upgrade a previous version of Data Collector to version 1.5.1.2. For instructions on upgrading, see the [Upgrade Documentation](#).

Fixed Issues

The following table lists the known issues that are fixed with this release.

JIRA	Description
SDC-3616	The Hadoop FS destination incorrectly updates the idle timeout for files that are idle when the destination writes data to other files.
SDC-3602	The Hive Metadata processor and Hive Metastore destination do not correctly route to error when Hive query failures occur.
SDC-3600	The Hive Metadata processor and the Hive Metastore destination use the incorrect field types.
SDC-3588	The Hive Metastore destination uses the incorrect date format for Avro schema serialization.

Known Issues

Please note the following known issues with this release.

For a full list of known issues, check out [our JIRA](#).

JIRA	Description
SDC-3445	When multiple pipelines write JSON data to MapR FS, the pipelines might encounter an exception.

StreamSets Data Collector 1.5.1.2 Release Notes

SDC-3357	<p>If you run Data Collector from Docker, you cannot shut down Data Collector by running <code>docker stop</code> or pressing Ctrl+C from the Docker Quickstart Terminal.</p> <p>Workaround: In the Data Collector console, click Administration > Shut Down.</p>
SDC-3356	<p>Using the following commands to shut down or restart Data Collector does not properly complete the shutdown:</p> <ul style="list-style-type: none"> • <code>service sdc stop</code> • <code>service sdc restart</code> <p>Workaround: In the Data Collector console, click Administration > Shut Down or Administration > Restart.</p>
SDC-3234	<p>Cluster streaming pipelines that run on YARN use the YARN user instead of the Data Collector user to run executors.</p>
SDC-3133	<p>When you upgrade Data Collector from the RPM package, the environment configuration file <code>\$SDC_DIST/libexec/sdcd-env.sh</code> is overwritten.</p> <p>Workaround: Back up the <code>sdcd-env.sh</code> file before you upgrade.</p>
SDC-3017	<p>Hadoop FS does not close and rename open files when the pipeline stops. Open files are named <code>prefix_tmp</code>. When the pipeline restarts, Hadoop FS continues processing the files, closing and renaming them as expected.</p> <p>Workaround: If your workflow immediately processes output files from Hadoop FS and your workflow requires synchronized data, pause downstream processing before you stop the pipeline.</p>
SDC-2998	<p>The JDBC Consumer origin does not correctly parse Oracle timestamp fields.</p> <p>Workaround: Set the JVM system property <code>oracle.jdbc.J2EE13Compliant</code> to <code>true</code> in the <code>SDC_JAVA_OPTS</code> environment variable in the <code>sdc-env.sh</code> or <code>sdcd-env.sh</code> file. For example:</p> <pre>export SDC_JAVA_OPTS="-Xmx1024m -Xms1024m -XX:PermSize=128m -XX:MaxPermSize=256m -Doracle.jdbc.J2EE13Compliant=true -server \${SDC_JAVA_OPTS}"</pre> <p>Setting the property to <code>true</code> returns <code>java.sql.Timestamp</code> instead of <code>oracle.sql.TIMESTAMP</code> for the timestamp SQL type.</p>
SDC-2950	<p>When a pipeline writes error records to Elasticsearch, the record header information - error code, error message, and error stage - is not preserved.</p>
SDC-2822	<p>If you configure a Kafka Producer destination to write one message per batch, and then use a cluster pipeline to process that data from the Kafka cluster, the cluster pipeline might encounter an out of memory error.</p>

StreamSets Data Collector 1.5.1.2 Release Notes

SDC-2586	<p>To process records larger than 1 MB, you must configure the <code>DataFactoryBuilder.OverRunLimit</code> property. However, this property is not configurable in the Data Collector configuration file, <code>\$SDC_CONF/sdc.properties</code></p> <p>Workaround: Set the value of <code>DataFactoryBuilder.OverRunLimit</code> property in the <code>SDC_JAVA_OPTS</code> environment variable in the Data Collector environment file, <code>\$SDC_DIST/libexec/sdc-env.sh</code> or <code>\$SDC_DIST/libexec/sdcd-env.sh</code></p> <p>Set the property greater than the largest record you want to process. For example, to process records up to 2 MB, set the property to 2097152 as follows:</p> <pre>SDC_JAVA_OPTS="-DDataFactoryBuilder.OverRunLimit=2097152"</pre>
SDC-2552	<p>When an invalid topic is specified for a MapR Streams Consumer or a MapR Streams Producer, the following incorrect message displays:</p> <pre>CONTAINER_0701 - Stage 'com_streamsets_pipeline_stage_origin_maprstreams_MapRStreamsDS ource_1' initialization error: org.apache.kafka.common.config.ConfigException: No bootstrap urls given in bootstrap.servers</pre> <p>This message is misleading because MapR Streams does not support the <code>bootstrap.servers</code> option.</p>
SDC-2374	<p>A cluster mode pipeline can hang with a <code>CONNECT_ERROR</code> status. This can be a temporary connection problem that resolves, returning the pipeline to the <code>RUNNING</code> status.</p> <p>If the problem is not temporary, you might need to manually edit the pipeline state file to set the pipeline to <code>STOPPED</code>. Edit the file only after you confirm that the pipeline is no longer running on the cluster or that the cluster has been decommissioned.</p> <p>To manually change the pipeline state, edit the following file: <code>\$SDC_DATA/runInfo/<cluster pipeline name>/<revision>/pipelineState.json</code></p> <p>In the file, change <code>CONNECT_ERROR</code> to <code>STOPPED</code> and save the file.</p>
SDC-2359	<p>Due to a Kafka issue, a pipeline with a Kafka Consumer or Kafka Producer can hang during validation or display initialization errors when unable to connect to a Kafka 0.9.0.0 broker. The Data Collector log indicates that the broker is unavailable.</p> <p>For more information about the Kafka JIRA, see https://issues.apache.org/jira/browse/KAFKA-2880.</p>
SDC-2133	<p>You cannot stop a pipeline in the middle of a long-running processor, such as a Jython Evaluator performing a complex script.</p>

StreamSets Data Collector 1.5.1.2 Release Notes

SDC-1731	<p>When using the Kafka Consumer or Kafka Producer on HDP 2.3 with Kerberos enabled, set the Kafka broker configuration property <code>security.inter.broker.protocol</code> to PLAINTEXT.</p> <p>When enabling Kerberos, HDP 2.3 sets the <code>security.inter.broker.protocol</code> property to PLAINTEXTSASL, which is not supported.</p> <p>If the property is not set to PLAINTEXT, when the pipeline starts, validation errors indicate a problem connecting to Kafka.</p>
SDC-1567	<p>You cannot use cluster mode pipelines to read from HDP 2.3 due to a HDP integration issue with Kafka and Spark Streaming.</p> <p>For more information about the HDP issue, see http://hortonworks.com/community/forums/topic/kafka-and-spark-streaming-nosuchmethoderror-kafka-consumer-simpleconsumer/.</p>
SDC-891	<p>At this time, writing to error records to file is not supported for cluster mode pipelines.</p> <p>Workaround: Write error records to Kafka or to an SDC RPC pipeline.</p>
SDC-890	<p>For cluster mode pipelines configured to stop on error or to stop upon reaching a memory limit, the Data Collector cannot stop all worker pipelines as expected.</p> <p>Workaround: To stop all pipelines, use the Stop icon in the Data Collector console.</p>

Contact Information

For more information about StreamSets, visit our website: <http://streamsets.com/>.

To report an issue, ask for help, or find out about our next meetup, check out our Community page: <http://streamsets.com/community/>.

For general inquiries, email us at info@streamsets.com.