

StreamSets Data Collector 2.1.0.1 Release Notes

October 17, 2016

We're happy to announce a new version of StreamSets Data Collector. This version features a few important bug fixes.

Upgrade

You can upgrade previous versions of Data Collector to version 2.1.0.1. For instructions on upgrading, see the [Upgrade Documentation](#).

Fixed Issues

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

For the full list, click [here](#).

JIRA	Description
SDC-4224	Remove limit on the number of oldest SCNs selected in the Oracle CDC Client origin.
SDC-4087	Data preview fails for pipelines that use the Dev Raw Data Source origin when you refresh the data preview or run data preview with changes.
SDC-3911	<p>The Field Renamer processor does not support quoting regex special characters such as the pipe symbol () in field names.</p> <p>If you followed the workaround in a previous release to manually quote the special character by wrapping it in \Q and \E, you can remove those characters. For example, to specify a field name of /'tag attr' in a previous release, you had to use the workaround of /'tag\Q \Eattr'. Now you can specify that field name as /'tag attr'.</p>

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-4954	<p>The Cassandra destination encounters problems connecting to a Cassandra cluster because the Cassandra stage library directory contains a mixed version of netty JAR files.</p> <p>Workaround:</p>

StreamSets Data Collector 2.1.0.1 Release Notes

	<ol style="list-style-type: none">1. Remove all netty* JAR files from the following directory: <code>\$\$SDC_DIST/streamsets-libs/streamsets-datacollector-cassandra_3-lib/lib</code>2. Download the following netty JAR file: <code>http://central.maven.org/maven2/io/netty/netty-all/4.0.41.Final/netty-all-4.0.41.Final.jar</code>3. Add the netty-all-4.0.41.Final.jar file to the Cassandra stage library directory.
SDC-4212	If you configure a UDP Source or UDP to Kafka origin to enable multithreading after you have already run the pipeline with the option disabled, the following validation error displays: <code>Multithreaded UDP server is not available on your platform.</code> Workaround: Restart Data Collector.
SDC-4172	Data Collector cannot access Vault secrets stored in Hashicorp Vault.
SDC-4128	In cluster mode, Data Collector does not generate log files for worker Data Collectors.
SDC-4090	The MapR FS destination does not support impersonating an HDFS user. Instead, the destination always uses the user account who started the Data Collector to connect to MapR FS.
SDC-4047	The XML Flattener processor fails to parse XML that contains whitespace after the XML prolog. Workaround: Use an Expression Evaluator or scripting processor to remove the whitespace before using the XML Flattener.
SDC-3944	The Hive Streaming destination using the MapR library cannot connect to a MapR cluster that uses Kerberos or username/password login authentication.
SDC-3357	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. Workaround: In the Data Collector console, click Administration > Shut Down .
SDC-3356	Using the following commands to shut down or restart Data Collector does not properly complete the shutdown: <ul style="list-style-type: none">• <code>service sdc stop</code>• <code>service sdc restart</code> Workaround: In the Data Collector console, click Administration > Shut Down or Administration > Restart .
SDC-3234	Cluster streaming pipelines that run on YARN use the YARN user instead of the Data Collector user to run executors.
SDC-3133	When you upgrade Data Collector from the RPM package, the environment configuration file <code>\$\$SDC_DIST/libexec/sdcd-env.sh</code> is overwritten.

StreamSets Data Collector 2.1.0.1 Release Notes

	<p>Workaround: Back up the <code>sdc-d-env.sh</code> file before you upgrade.</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>
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/sdc-d-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-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-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 <code>PLAINTEXT</code>.</p> <p>When enabling Kerberos, HDP 2.3 sets the <code>security.inter.broker.protocol</code> property to <code>PLAINTEXTSASL</code>, which is not supported.</p> <p>If the property is not set to <code>PLAINTEXT</code>, when the pipeline starts, validation errors indicate a problem connecting to Kafka.</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>

StreamSets Data Collector 2.1.0.1 Release Notes

Contact Information

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

To review the latest documentation or try out our tutorials, check out the following links:

- [User Guide](#)
- [User Guide tutorial](#)
- [GitHub tutorials](#)

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.