

Cloudera DataFlow for Data Hub 7.2.17

Streaming Analytics

Date published: 2019-12-16

Date modified: 2023-06-27

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

Streaming Analytics in Cloudera.....	4
What is Apache Flink?.....	5
Core features of Flink.....	6
What is SQL Stream Builder?.....	9
Key features of SSB.....	10
SQL Stream Builder architecture.....	11

Streaming Analytics in Cloudera

Cloudera Streaming Analytics (CSA) offers real-time stream processing and streaming analytics powered by Apache Flink. Flink implemented on CDP provides a flexible streaming solution with low latency that can scale to large throughput and state. In addition to Flink, CSA includes SQL Stream Builder (SSB) to offer data analytical experience using SQL queries on your data streams.

Key features of Cloudera Streaming Analytics

Apache Flink

CSA is powered by Apache Flink that offers a framework for real-time stream processing and streaming analytics. CSA offers the features and functionalities of the upstream Apache Flink integrated on CDP Public Cloud.

SQL Stream Builder

SQL Stream Builder is a job management interface to compose and run Continuous Streaming SQL on streams using Apache Flink as an engine, as well as to create REST APIs for the results.

Cloudera Platform

Implementing Flink on the Cloudera Platform allows you to easily integrate with Runtime components, and have all the advantages of cluster and service management with Cloudera Manager.

Streaming Platform

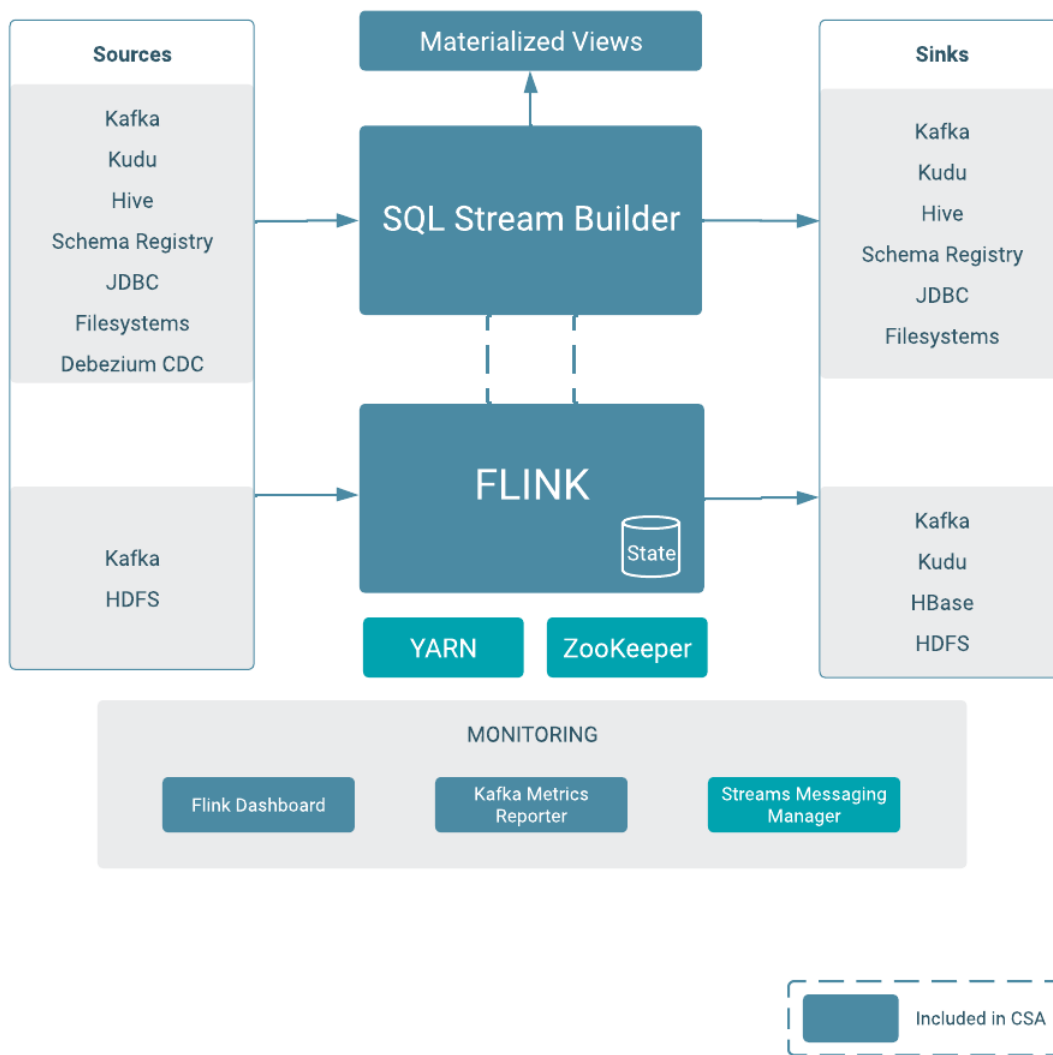
For streaming analytics, CSA fits into a complete streaming platform augmented by Apache Kafka, Schema Registry, Streams Messaging Manager in the Cloudera Runtime stack.

Supported Connectors

CSA offers a set of connectors for Flink and SSB from which you can choose from based on your requirements. Kafka, HBase, HDFS, Kudu and Hive connectors are available for Flink. Kafka, HDFS/S3, JDBC and a set of the Debezium CDC connectors are available for SSB. Other than the connectors, SSB also supports Schema Registry, Hive and Kudu as catalogs.

Monitoring Solutions

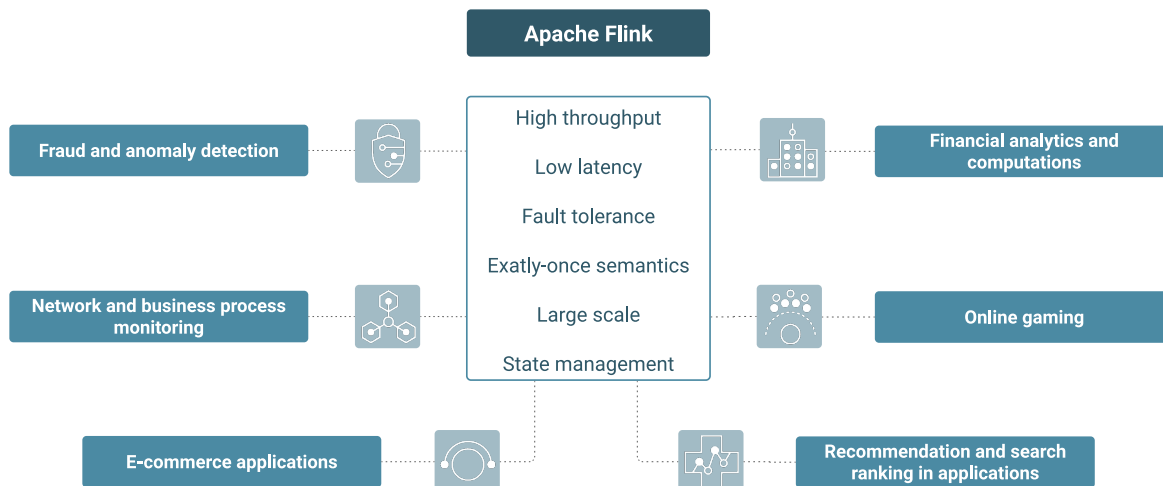
In CDP Public Cloud, Streams Messaging components provide tools that support the operational needs of CSA. For example, when Kafka is used as a connector, you can use Kafka Metrics Reporter and Streams Messaging Manager (SMM) for Kafka management and alerting of Kafka actions. Beside SMM, you can use the Flink Dashboard to monitor your Flink and SSB jobs.



What is Apache Flink?

Flink is a distributed processing engine and a scalable data analytics framework. You can use Flink to process data streams at a large scale and to deliver real-time analytical insights about your processed data with your streaming application.

Flink is designed to run in all common cluster environments, perform computations at in-memory speed and at any scale. Furthermore, Flink provides communication, fault tolerance, and data distribution for distributed computations over data streams. A large variety of enterprises choose Flink as a stream processing platform due to its ability to handle scale, stateful stream processing, and event time.

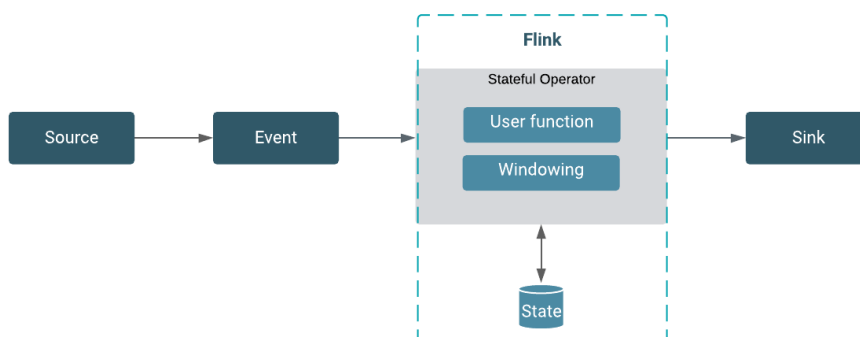


Core features of Flink

Learn more about the specific details of Flink architecture, the DataStream API, how Flink handles event time and watermarks, and how state management works in Flink.

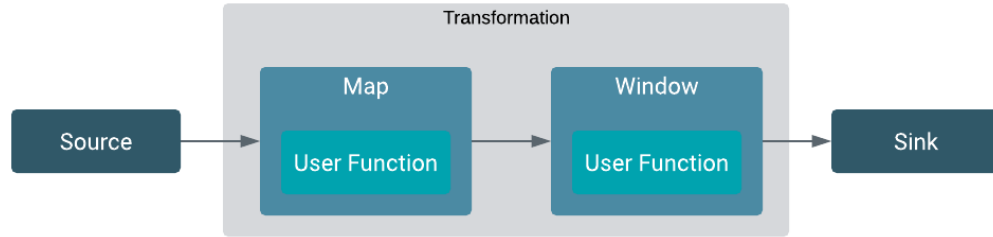
Architecture

The two main components for the task execution process are the Job Manager and Task Manager. The Job Manager on a master node starts a worker node. On a worker node the Task Managers are responsible for running tasks and the Task Manager can also run more than one task at the same time. The resource management for the tasks are completed by the Job manager in Flink. In a Flink cluster, Flink jobs are executed as YARN applications. HDFS is used to store recovery and log data, while ZooKeeper is used for high availability coordination for jobs.



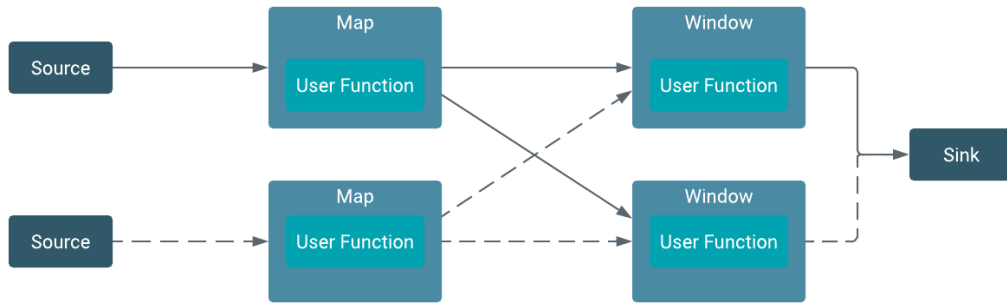
DataStream API

The DataStream API is used as the core API to develop Flink streaming applications using Java or Scala programming languages. The DataStream API provides the core building blocks of the Flink streaming application: the datastream and the transformation on it. In a Flink program, the incoming data streams from a source are transformed by a defined operation which results in one or more output streams to the sink.



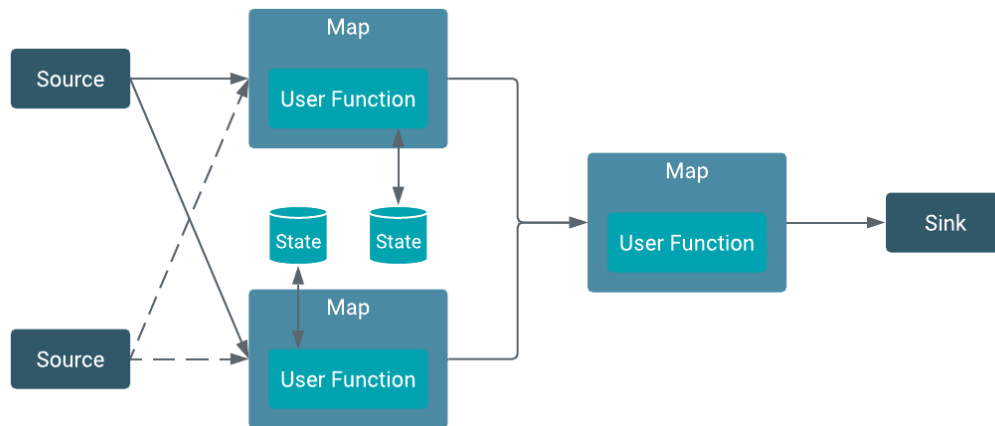
Operators

Operators transform one or more DataStreams into a new DataStream. Programs can combine multiple transformations into sophisticated data flow topologies. Other than the standard transformations like map, filter, aggregation, you can also create windows and join windows within the Flink operators. On a dataflow one or more operations can be defined which can be processed in parallel and independently to each other. With windowing functions, different computations can be applied to different streams in the defined time window to further maintain the processing of events. The following image illustrates the parallel structure of dataflows.



State and state backend

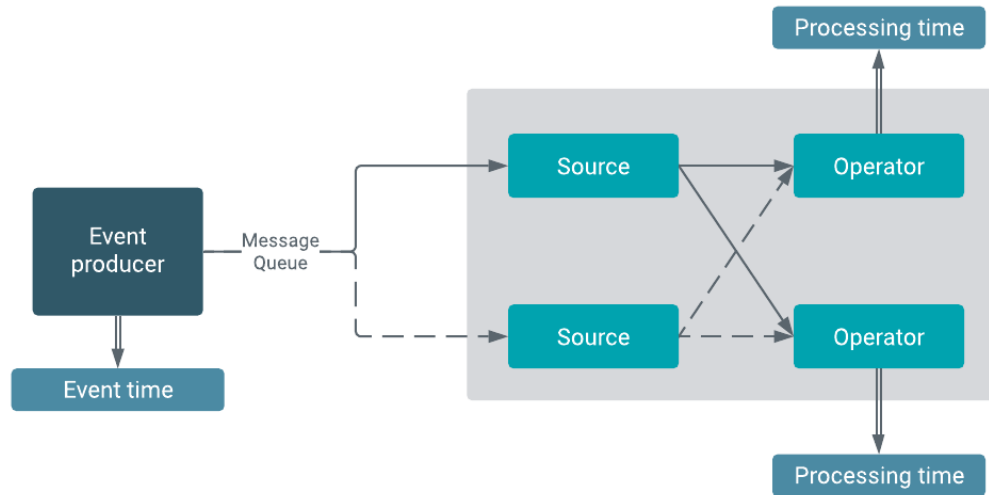
Stateful applications process dataflows with operations that store and access information across multiple events. You can use Flink to store the state of your application locally in state backends that guarantee lower latency when accessing your processed data. You can also create checkpoints and savepoints to have a fault-tolerant backup of your streaming application on a durable storage.



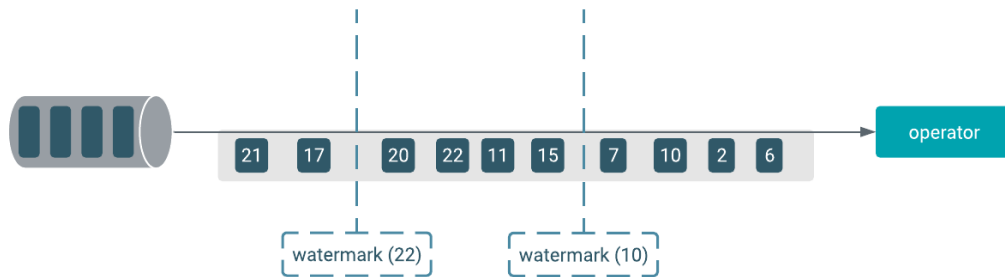
Event time and watermark

In time-sensitive cases where the application uses alerting or triggering functions, it is important to distinguish between event time and processing time. To make the designing of applications easier,

you can create your Flink application either based on the time when the event is created or when it is processed by the operator.

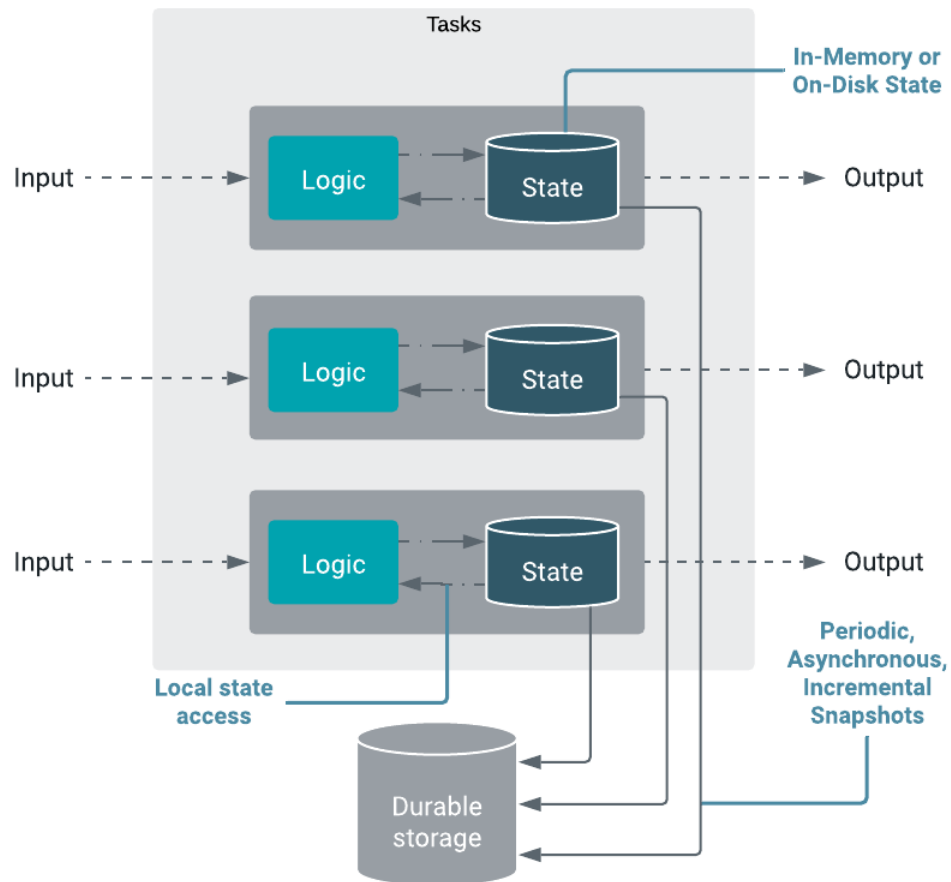


With only the event time, it is not clear when the events are processed in the application. To track the time for an event time based application, watermark can be used.



Checkpoints and savepoints

Checkpoints and savepoints can be created to make the Flink application fault tolerant throughout the whole pipeline. Flink contains a fault tolerance mechanism that creates snapshots of the data stream continuously. The snapshot includes not only the dataflow, but the state attached to it. In case of failure, the latest snapshot is chosen and the system recovers from that checkpoint. This guarantees that the result of the computation can always be consistently restored. While checkpoints are created and managed by Flink, savepoints are controlled by the user. A savepoint can be described as a backup from the executed process.



Related Information

[Flink application structure](#)

[Configuring RocksDB state backend](#)

[Enabling checkpoints for Flink applications](#)

[Enabling savepoints for Flink applications](#)

What is SQL Stream Builder?

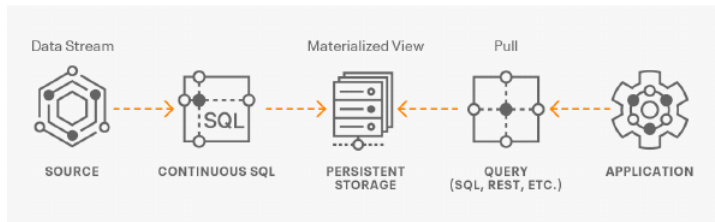
Cloudera Streaming Analytics offers an easy to use and interactive SQL Stream Builder as a service to create queries on streams of data through SQL.

The SQL Stream Builder (SSB) is a comprehensive interactive user interface for creating stateful stream processing jobs using SQL. By using SQL, you can simply and easily declare expressions that filter, aggregate, route, and otherwise mutate streams of data. SSB is a job management interface that you can use to compose and run SQL on streams, as well as to create durable data APIs for the results.

What is Continuous SQL?

SSB runs Structured Query Language (SQL) statements continuously, this is called Continuous SQL or Streaming SQL. Continuous SQL can run against both bounded and unbounded streams of data. The results are sent to a sink of some type, and can be connected to other applications through a Materialized View interface. Compared to traditional SQL, in Continuous SQL the data has a start, but no end. This means that queries continuously process results. When

you define your job in SQL, the SQL statement is interpreted and validated against a schema. After the statement is executed, the results that match the criteria are continuously returned.



Integration with Flink

SSB runs in an interactive fashion where you can quickly see the results of your query and iterate on your SQL syntax. The executed SQL queries run as jobs on the Flink cluster, operating on boundless streams of data until canceled. This allows you to author, launch, and monitor stream processing jobs within SSB as every SQL query is a Flink job. You can use Flink and submit Flink jobs without using Java, as SSB automatically builds and runs the Flink job in the background.

As a result of Flink integration, you are able to use the basic functionalities offered by Flink. You can choose exactly once processing, process your data stream using event time, save your jobs with savepoints, and use Flink SQL to create tables and use connectors based on your requirements. As a result of the various connectors, you are able to enrich your streaming data with data from slowly changing connectors.

Key features of SSB

SQL Stream Builder (SSB) within Cloudera supports out-of-box integration with Flink. Using Flink SQL, you can create tables directly from the Streaming SQL Console window or built-in templates. For integration with Business Intelligence tools you can create Materialized Views.

Flink SQL

SQL Stream Builder allows you to use Data Definition Language (DDL), Data Manipulation Language (DML) and Query Language directly from the Streaming SQL Console.

For more information about the supported Flink SQL statements, see the [Flink SQL statements](#) section.

Change Data Capture

SQL Stream Builder supports PostgreSQL, Oracle, MySQL, Db2 and SQL Server as Debezium connectors using Flink SQL. With Change Data Capture (CDC), you can capture changes in your databases and update your applications with the newly added data.

For more information about Debezium, see the [official documentation](#).

Session Cluster

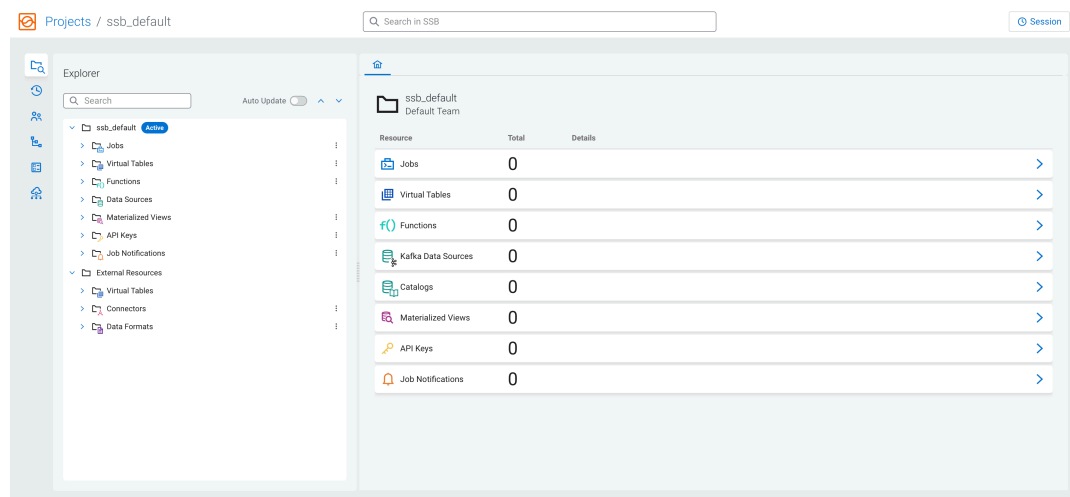
When you start SSB, a session cluster is deployed to share and maintain resources. The submitted SQL jobs are executed as Flink jobs in the same session cluster that share a Job Manager. The properties for the session cluster can be viewed using the Streaming SQL Console, and can be modified with SET statements in the SQL window.

Built-in Templates

The Built-in Templates in SSB allows you to quickly and simply create tables for your SQL queries. You only need to provide the connection and job specific information to the template to use it in SSB.

Streaming SQL Console

SSB comes with an interactive user interface that allows you to easily create, and manage your SQL jobs in one place. It allows you to create and iterate on SQL statements with robust tooling and capabilities. Query parsing is logged to the console, and results are sampled back to the interface to help with iterating on the SQL statement as required.



Materialized Views

SSB has the capability to materialize results from a Streaming SQL query to a persistent view of the data that can be read through REST and over the PG wire protocol. Applications can use this mechanism to query streams of data in a way of high performance without deploying additional database systems. Materialized Views are built into the SQL Stream Builder service, and require no configuration or maintenance. The Materialized Views act like a special kind of sink, and can even be used in place of a sink. They require no indexing, storage allocation, or specific management.

Input Transform

In case you are not aware of the incoming data structure or raw data is being collected from for example sensors, you can use the Input Transform to clean up and organize the incoming data before querying. Input transforms also allow access to Kafka header metadata directly in the query itself. Input transforms are written in Javascript and compiled to Java bytecode deployed with the Flink jar.

User-defined Functions

You can create customized and complex SQL queries by using User-defined Functions to enrich your data, apply computations or a business logic on it. User defined functions are written in Javascript or Java language.

Related Information

[Integration with Flink](#)

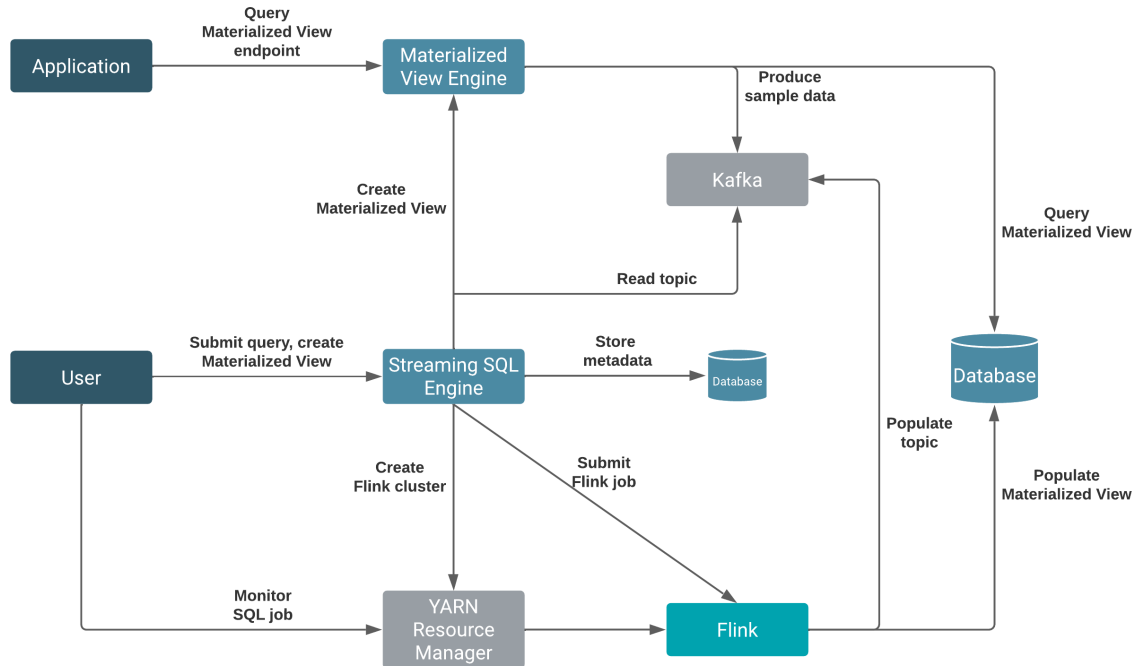
[Flink SQL statements](#)

SQL Stream Builder architecture

The SQL Stream Builder (SSB) service is integrated on Cloudera Data Platform (CDP) and connected to Flink and its services. The SSB architecture includes SQL Stream Engine and Materialized View Engine. These main components within SSB are responsible for executing jobs, populating topics, creating metadata and querying data that happens in the background.

SSB consists of the following main components:

- SQL Stream Engine
- Materialized View Engine



The primary point of user interaction for SQL Stream Builder is the Streaming SQL Console User Interface running on the Streaming SQL Engine. When a query is submitted on the Streaming SQL Console, the Streaming SQL Engine automatically creates a Flink job in the background on the cluster. SSB also requires a Kafka service on the same cluster. This mandatory Kafka service is used to automatically populate topics for the websocket output. The websocket output is needed for sampling data to the Console, and when no table is added to output the results of the SQL query.

When a Materialized View query is submitted, Flink generates the data to the Materialized View database from which the Materialized View Engine queries the required data.

Database management in SSB

SSB uses databases in the following cases:

- To store metadata of SQL jobs
- To store data for creating Materialized Views
- As a connector for Flink SQL

In CDP Public Cloud, PostgreSQL is supported as a default database for SQL Stream Builder and the Materialized View database.