

Configuration

Date published: 2019-12-17

Date modified: 2019-12-17

CLOUdera

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

How to configure Apache Flink?.....	4
Setting max parallelism.....	4
Configuring Flink application resources.....	4
RocksDB state backend configuration.....	5
Controlling memory usage.....	5
Configuration parameters for RocksDB.....	6
RocksDB based timers.....	6

How to configure Apache Flink?

Cloudera Streaming Analytics includes Flink with configuration that works out of the box. It is not mandatory to configure Flink to production, but you can use the available configurations to optimize the application behavior in production. Cloudera Manager includes all the necessary configurations for Flink that can also be accessed from the `flink-conf.yaml` file.

Setting max parallelism

The max parallelism is the most essential part of resource configuration for Flink applications as it defines the maximum jobs that are executed at the same time in parallel instances. However, you can optimize max parallelism in case your production goals differ from the default settings.

In a Flink application, the different tasks are split into several parallel instances for execution. The number of parallel instances for a task is called parallelism. Parallelism can be defined at the operator, client, execution environment and system level. Cloudera recommends setting parallelism to a lower value at the initially and increasing it over time, if the job cannot keep up with the input rate.

To configure the max parallelism, `setMaxParallelism` is called as it controls the number of key-groups created by the state backends. A key-group is a partition of an operator state. The number of key-groups determines how data is going to be distributed among the parallel operators. If the key-groups are not distributed evenly, the data distribution is also uneven.

Consider the following aspects when setting the max parallelism:

- The number should be large enough to accommodate expected future load increases as this setting cannot be changed without starting from an empty state.
- If `P` is the selected parallelism for the job, the max parallelism should be divisible by `P` to get even state distribution.
- Please note that larger max parallelism settings have greater cost on the state backend side, for large scale production jobs benchmarking the size of the state based on the maximum parallelism is useful before changing this parameter.

Based on these criteria, Cloudera recommends setting the max parallelism to factorials or other numbers with a large number of divisors (120, 180, 240, 360, 720, 840, 1260), which will make parallelism tuning easier.

Table 1: Reference values

Stateless	In-memory state	RocksDB state
1 million record / sec / core	100 000 records / sec / core	10 000 records / sec / core

Configuring Flink application resources

Generally, Flink automatically identifies the resources for an application. You can configure the Taskmanager, buffer timeout and high availability for a Flink application to maintain resources in a more efficient way.

YARN automatically kills application containers that use more memory than their allocated limit. To avoid Flink TaskManagers getting killed by YARN use the following calculation to set the size: $TM\ total\ size = TM\ Heap + Heap\ -cutoff$. The default heap cutoff is 25% and it is also configurable.

Furthermore, use a YARN queue with preemption disabled to avoid long running jobs being affected when the cluster reaches its capacity limit.

Flink uses network buffers to transfer data from one operator to another. These buffers are filled up with data during the specified time for the timeout. In case of high data rates, the set time is usually never reached. For cases when the data rate is high the throughput can be further increased with setting the buffer timeout to an intentionally higher value due to the characteristics of the TCP channel. However this in turn increases the latency of the pipeline.

Flink on YARN jobs are configured to tolerate a maximum number of failed containers before they terminate. Configure the YARN maximum failed containers setting in proportion to the total parallelism and the expected lifetime of the job.

High Availability is enabled by default in CSA. This eliminates the Job Manager as a single point of failure. You can also tune the application resilience by setting the YARN maximum application attempts, which determines how many times the application will retry in case of failures.

Table 2: Reference values

Configuration	Parameter	Recommended value
TM container memory	-ytm / taskmanager.heap.size	<i>TM Heap + Heap-cutoff</i>
Max parallelism	env.setMaxParallelism(num)	<i>120,720,1260,5040</i>
Container heap cutoff	containerized.heap-cutoff-ratio	<i>0.25-0.75</i>
Buffer timeout	env.setBufferTimeout(millis)	<i>1-100</i>
YARN queue	-yqu	<i>A queue with no preemption</i>
YARN max failed containers	yarn.maximum-failed-containers	<i>3*num_containers</i>
YARN max AM failures	yarn.application-attempts	<i>3-5</i>

RocksDB state backend configuration

You can use RocksDB as a state backend when your Flink streaming application requires a larger state. With `ClouderaConfigurableOptionsFactory`, you are able to configure RocksDB for YARN applications to have a limit on memory usage supported by the Write Buffer Manager.

The default RocksDB state backend configuration available in Flink is not suitable for YARN applications with large state, as it does not effectively limit the native memory usage of the embedded RocksDB database. To effectively enforce a limit on memory usage, support for Write Buffer Manager is introduced in the `ClouderaConfigurableOptionsFactory`.

`ClouderaConfigurableOptionsFactory` can be added to the Flink configuration file to get a good initial configuration of the state backend:

```
state.backend: rocksdb
state.backend.rocksdb.options-factory: org.apache.flink.contrib.streaming
.state.ClouderaConfigurableOptionsFactory
```

For the default RocksDB state backend configuration, see the Apache Flink [documentation](#). See more information about the Write Buffer Manager [here](#).

Related Information

[Stateful Tutorial: Configure the application for production](#)

Controlling memory usage

When you enable `WriteBufferManager` for a Flink application, the block-cache size becomes the single most important option in limiting the memory size. You can set it 256 mb as a starting value, and later increase it based on your production experience.

It is configured with the following option:

- `state.backend.rocksdb.block.cache-size: 256mb`

The block cache serves as the in-memory cache for both reads and writes for the RocksDB instances. Flink creates an independent RocksDB instance for each stateful (keyed state) operator and subtask, so the minimum RocksDB native memory requirement is `num_task_slots * num_stateful_ops * block_cache_size`. The block size should be chosen depending on the size of the state and number of keys, but 256 MB is a good starting value.

If the read performance is insufficient, the block-size can be increased to 512 MB, 1 GB or even 2 GB on machines with large memory sizes.

Once you made the calculation, you need to ensure that the YARN container has enough extra memory on top of the TM heap size to cover the native memory requirement. This can be controlled by the following settings:

- [containerized-heap-cutoff-min](#)
- [containerized-heap-cutoff-ratio](#)

Configuration parameters for RocksDB

Cloudera Streaming Analytics offers you new configurable options factory settings beside the pre-existing one.

In addition to the options factory, the following new options (and their defaults) are introduced:

- `state.backend.rocksdb.disk-type` : *SPINNING*
- `state.backend.rocksdb.writebuffer.manager.size` : *blockcache / 4*
- `state.backend.rocksdb.writebuffer.manager.enabled` : *true*
- `state.backend.rocksdb.cache-index-and-filter` : *true*
- `state.backend.rocksdb.optimize-for-hits` : *true*
- `state.backend.rocksdb.stats.interval-seconds` : *600*

For the list of configuration parameters, see the Apache Flink [documentation](#).

RocksDB based timers

Even if you set RocksDB as a state backend, Flink timers, registered by the window operator or defined by the user, are stored on-heap by default. You can change this by setting the timer service property to RocksDB to reduce the heap consumption.

Timers can be user-defined or registered by the system for window operators. While heap based timers are very quick, they might require too large container sizes.

By setting the [state-backend-rocksdb-timer-service-factory](#) configuration property to *ROCKSDB*, Flink stores all timers in the state backend. This has a severe impact on timer performance, but also greatly reduces the heap consumption.