

# How to clean up storage in Kylin 4

- Background
  - Directory tree structure under Kylin 4.0 's working dir
  - Summary
- How to use
  - Option Table
  - List help information
  - List directory which to be deleted
  - Deleted them after confirm
  - Only delete stale job\_tmp and unreferenced cuboid files

## Background

Kylin will generate temporary files in HDFS during the cube building; Besides, when purge/drop/merge cubes, some parquet files may be left in HDFS and will no longer be queried; Although Kylin has started to do some automated garbage collection, it might not cover all cases; You can do an offline storage cleanup periodically.

## Directory tree structure under Kylin 4.0 's working dir

### Working Dir(ROOT)

- **{PROJECT\_NAME} [managed by tool]**
  - **parquet**
    - {CUBE\_NAME} [managed by tool]
      - {SEGMENT\_NAME} [managed by tool]
        - {CUBOID\_ID}
          - parquet files
  - **spark\_log**
    - driver
      - {JOB\_ID}
        - drivers' log of cubing job
    - executor
      - {JOB\_ID}
        - executors' log of cubing job
  - **dict/global\_dict [managed by tool]**
    - {CUBE\_NAME}
      - {COLUMN\_NAME}
        - dict files
  - **table\_snapshot [managed by tool]**
    - {SCHEMA\_NAME.TABLE\_NAME}
      - {JOB\_ID}
        - parquet files
  - **job\_tmp [managed by tool]**
    - {JOB\_ID}
      - TBD
- **cube\_statistics**
  - {CUBE\_NAME}
    - {JOB\_ID}
      - seq file of cuboid 's HLL
- **\_sparder\_log**
  - {DATE}
    - executors 's log of query job
- **resources-jdbc**
  - TBD

## Summary

In above directory tree, the directory which end with "**managed by tool**" means **StorageCleanupJob** will try to check and delete useless files under these directory.

For directory **table\_snapshot**, **dict/global\_dict**, **parquet/{CUBE\_NAME}**, **parquet/{CUBE\_NAME}/{SEGMENT\_NAME}** , Kylin will mark files which is unreferenced and stale(by checking last modified time) as garbage.

For directory **job\_tmp**, Kylin will only check last modified time.

## How to use

### Option Table

Option	Data Type	Default Value	Comment
--------	-----------	---------------	---------

delete	Boolean	false	Boolean, whether or not to do real delete operation. Default value is false, means a dry run.
cleanupTableSnapshot	Boolean	true	Boolean, whether or not to delete unreferenced snapshot files. Default value is true .
cleanupGlobalDict	Boolean	true	Boolean, whether or not to delete unreferenced global dict files. Default value is true .
cleanupJobTmp	Boolean	false	Boolean, whether or not to delete job tmp files. Default value is false .
cleanupThreshold	Integer	168	Integer, used to specific delete unreferenced storage that have not been modified before how many hours (recent files are protected). Default value is 168 hours.

## List help information

options
<pre>[root@cdh-master apache-kylin-4.0.0-SNAPSHOT-bin]# bin/kylin.sh org.apache.kylin.tool.StorageCleanupJob -help Retrieving hive dependency... Retrieving hadoop conf dir... Retrieving Spark dependency... ... Running org.apache.kylin.rest.job.StorageCleanupJob -help usage: org.apache.kylin.rest.job.StorageCleanupJob   -cleanupGlobalDict &lt;cleanupGlobalDict&gt;      Boolean, whether or not to   delete unreferenced global   dict files. Default value   is true .   -cleanupJobTmp &lt;cleanupJobTmp&gt;               Boolean, whether or not to   delete job tmp files.   Default value is false .   -cleanupTableSnapshot &lt;cleanupTableSnapshot&gt; Boolean, whether or not to   delete unreferenced   snapshot files. Default   value is true .   -cleanupThreshold &lt;cleanupThreshold&gt;         Integer, used to specific   delete unreferenced   storage that have not been   modified before how many   hours (recent files are   protected). Default value   is 168 hours.   -delete &lt;delete&gt;                             Boolean, whether or not to   do real delete operation.   Default value is false,   means a dry run.</pre>

## List directory which to be deleted

bin/kylin.sh org.apache.kylin.tool.StorageCleanupJob
--

## Deleted them after confirm

bin/kylin.sh org.apache.kylin.tool.StorageCleanupJob --delete true
--

## Only delete stale job\_tmp and unreferenced cuboid files

<pre>bin/kylin.sh org.apache.kylin.tool.StorageCleanupJob --delete true \ --cleanupJobTmp ture -cleanupTableSnapshot false \ -cleanupGlobalDict false --cleanupThreshold 24</pre>
---

