

Hudi cleaner:

Before:

Objects (7) Info

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 Inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

| <input type="checkbox"/> | Name | Type | Last modified | Size | Storage class |
|--------------------------|---|---------------------------|--|----------|---------------|
| <input type="checkbox"/> | <div> <div>📁</div> <div> hoodie_partition_metadata </div> </div> | hoodie_partition_metadata | September 30, 2024, 16:48:42 (UTC+05:30) | 96.0 B | Standard |
| <input type="checkbox"/> | <div> <div>📁</div> <div> hoodie/ </div> </div> | Folder | - | - | - |
| <input type="checkbox"/> | <div> <div>📄</div> <div> 21f40bb7-f082-40a1-b98f-b141ad66628b-O-129-590_20240930112344952.parquet </div> </div> | parquet | September 30, 2024, 16:33:52 (UTC+05:30) | 426.6 KB | Standard |
| <input type="checkbox"/> | <div> <div>📄</div> <div> 21f40bb7-f082-40a1-b98f-b141ad66628b-O-164-737_2024093011281868.parquet </div> </div> | parquet | September 30, 2024, 16:58:19 (UTC+05:30) | 426.6 KB | Standard |
| <input type="checkbox"/> | <div> <div>📄</div> <div> 21f40bb7-f082-40a1-b98f-b141ad66628b-O-30-157_20240930111812735.parquet </div> </div> | parquet | September 30, 2024, 16:48:44 (UTC+05:30) | 426.1 KB | Standard |
| <input type="checkbox"/> | <div> <div>📄</div> <div> 21f40bb7-f082-40a1-b98f-b141ad66628b-O-62-300_20240930112120711.parquet </div> </div> | parquet | September 30, 2024, 16:51:34 (UTC+05:30) | 426.3 KB | Standard |
| <input type="checkbox"/> | <div> <div>📄</div> <div> 21f40bb7-f082-40a1-b98f-b141ad66628b-O-94-443_20240930112247543.parquet </div> </div> | parquet | September 30, 2024, 16:52:55 (UTC+05:30) | 426.5 KB | Standard |

After adding cleaner policy:

hudi_table/

Copy S3 URI

Objects

Properties

Objects (5) Info

🔄

Copy S3 URI

Copy URL

Download

Open

Delete

Actions ▼

Create folder

Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 Inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

🔍 Find objects by prefix

<

1

>

⊙

| <input type="checkbox"/> | Name ▲ | Type ▼ | Last modified ▼ | Size ▼ | Storage class ▼ |
|--------------------------|---|---------------------------|--|----------|-----------------|
| <input type="checkbox"/> | .hoodie_partition_metadata | hoodie_partition_metadata | September 30, 2024, 17:03:13 (UTC+05:30) | 96.0 B | Standard |
| <input type="checkbox"/> | .hoodie/ | Folder | - | - | - |
| <input type="checkbox"/> | e76d4482-7360-4c0e-ae46-e5b6579bec32-0_0-212-910_20240930113259879.parquet | parquet | September 30, 2024, 17:03:14 (UTC+05:30) | 426.1 KB | Standard |
| <input type="checkbox"/> | e76d4482-7360-4c0e-ae46-e5b6579bec32-0_0-244-1053_20240930113402833.parquet | parquet | September 30, 2024, 17:04:09 (UTC+05:30) | 426.3 KB | Standard |
| <input type="checkbox"/> | e76d4482-7360-4c0e-ae46-e5b6579bec32-0_0-276-1196_20240930113420641.parquet | parquet | September 30, 2024, 17:04:26 (UTC+05:30) | 426.5 KB | Standard |

Cleaner configurations

'hoodie.cleaner.policy': 'KEEP_LATEST_COMMITS', # Use KEEP_LATEST_COMMITS policy

'hoodie.cleaner.max.commits': '3', # Keep the latest 3 commits

'hoodie.cleaner.parallelism': '4', # Number of parallel cleaner threads

Based on the Hudi options you provided, the cleaning policy is set to KEEP_LATEST_COMMITS, which means Hudi will automatically clean up older versions of data, keeping only the latest 3 commits .

Since the automatic cleaning is governed by this setting, Hudi will automatically clean the table during the commit phase (upsert or write operations) by keeping only the latest 3 commits and deleting older files.

Here's an explanation of the three cleaner configuration options:

1. 'hoodie.cleaner.policy': 'KEEP_LATEST_COMMITS'

- **Purpose:** This configuration sets the cleaning policy for the Hudi table.
- **Details:**
 - **Policy:** KEEP_LATEST_COMMITS ensures that only the latest commits are retained, and the older commits are cleaned up. This helps in managing storage by removing outdated or unnecessary data versions.
 - **Use case:** It's useful when you want to limit the number of data versions in the table, keeping only the most recent versions to optimize space.

2. 'hoodie.cleaner.max.commits': '3'

- **Purpose:** This option defines how many of the most recent commits should be retained during the cleaning process.
- **Details:**
 - **Setting:** 3 means Hudi will keep the latest 3 commits and delete data files related to older commits during the cleaning process.

- **Use case:** This ensures that even if multiple versions of the data exist, the cleaner will keep the most recent 3 versions, thus providing a rollback buffer, but not overwhelming the storage with too many versions.

3. 'hoodie.cleaner.parallelism': '4'

- **Purpose:** This setting controls the level of parallelism during the cleaning process.
- **Details:**
 - **Setting:** 4 means that Hudi will run 4 cleaner threads in parallel.
 - **Impact:** Higher parallelism increases the speed of cleaning, especially for large datasets, by cleaning multiple partitions or files at the same time. However, setting parallelism too high may cause more resource contention.
 - **Use case:** This is useful in environments with sufficient resources (e.g., a multi-core machine or distributed cluster) to speed up the cleaning process without overwhelming the system.

In Summary:

- **KEEP_LATEST_COMMITS:** Defines the strategy to retain only the latest commits.
- **max.commits = 3:** Retains only the last 3 commits, helping manage storage efficiently.
- **parallelism = 4:** Speeds up the cleaning process by using 4 concurrent threads.

These settings together help balance version retention and efficient cleaning in a resource-optimized manner.