# BIG DATA (HTTPS://TIMEPASSTECHIES.COM/)

## Hive 6 tutorial 9 - Hive performance tuning using join optimization with common, map, bucket and skew join

Map Reduce (https://timepasstechies.com/category/programming/data-analytics/mapreduce/)

🕓 August, 2017   👤 adarsh (https://timepasstechies.com/author/adarshgorur/)   💬 Leave a

comment (https://timepasstechies.com/hive-tutorial-9-hive-performance-tuning-using-join-optimization-common-map-bucket-skew-join/#respond)

Hive (https://timepasstechies.com/category/programming/data-analytics/hive/)

Hdfs & Yarn (https://timepasstechies.com/category/programming/data-analytics/hdfs/)

Introduction to PyS
DataCamp Online (

Pig (https://timepasstechies.com/category/programming/data-analytics/pig/)

Online coding and data
science at your own pa
datacamp.com

Oozie (https://timepasstechies.com/category/programming/data-analytics/oozie/)

**Learn more**

Hbase (https://timepasstechies.com/category/programming/data-analytics/hbase/)

Design Patterns (https://timepasstechies.com/category/programming/design-patterns/)

### Common join

Streaming (https://timepasstechies.com/category/stream-processing/) and works for most of the time. For common joins, we need to make sure the big table is on the right-most side or specified by hit, as follows.

## Posts

```
/*+ STREAMTABLE(stream_table_name) */.
```

## Map join

Map join is used when one of the join tables is small enough to fit in the memory, so it is very fast but limited.Hive can convert map join automatically with the following settings.

```
1.  SET hive.auto.convert.join=true; --default false
2.
3.  SET hive.mapjoin.smalltable.filesize=600000000; --default 25M
4.
5.  SET hive.auto.convert.join.noconditionaltask=true; --default false. Set to true so that
6.
7.  SET hive.auto.convert.join.noconditionaltask.size=10000000; --The default value control
```

Once autoconvert is enabled, Hive will automatically check if the smaller table file size is bigger than the value specified by hive.mapjoin.smalltable.filesize, and then Hive will convert the join to a common join. If the file size is smaller than this threshold, it will try to convert the common join into a map join. Once autoconvert join is enabled, there is no need to provide the map join hints in the query.

## Bucket map join

Bucket map join is a special type of map join applied on the bucket tables. To enable bucket map join, we need to enable the following settings.

```
1.   SET hive.auto.convert.join=true; --default false
2.
3.   SET hive.optimize.bucketmapjoin=true; --default false
```

In bucket map join, all the join tables must be bucket tables and join on buckets columns. In addition, the buckets number in bigger tables must be a multiple of the bucket number in the small tables.

## Sort merge bucket (SMB) join

SMB is the join performed on the bucket tables that have the same sorted, bucket, and join condition columns. It reads data from both bucket tables and performs common joins (map and reduce triggered) on the bucket tables. We need to enable the following properties to use SMB.

```
1.   SET hive.input.format=org.apache.hadoop.hive.ql.io.BucketizedHiveInputFormat
2.   SET hive.auto.convert.sortmerge.join=true;
3.   SET hive.optimize.bucketmapjoin=true;
4.   SET hive.optimize.bucketmapjoin.sortedmerge=true;
5.   SET hive.auto.convert.sortmerge.join.noconditionaltask=true;
```

## Sort merge bucket map (SMBM) join

SMBM join is a special bucket join but triggers map-side join only. It can avoid caching all rows in the memory like map join does. To perform SMBM joins, the join tables must have the same bucket, sort, and join condition columns. To enable such joins, we need to enable the following settings.

```
1.   SET hive.auto.convert.join=true;
2.   SET hive.auto.convert.sortmerge.join=true
3.   SET hive.optimize.bucketmapjoin=true;
4.   SET hive.optimize.bucketmapjoin.sortedmerge=true;
5.   SET hive.auto.convert.sortmerge.join.noconditionaltask=true;
```

```
   6.    SET hive.auto.convert.sortmerge.join.bigtable.selection.policy=org.apache.hadoop.hive.q
```

## Skew join

When working with data that has a highly uneven distribution, the data skew could happen in such a way that a small number of compute nodes must handle the bulk of the computation. The following setting informs Hive to optimize properly if data skew happens.

```
   1.    SET hive.optimize.skewjoin=true; --If there is data skew in join, set it to true. Defaul
   2.    SET hive.skewjoin.key=100000; --This is the default value. If the number of key is bigg
```

> Note : Skew data could happen on the GROUP BY data too. To optimize it, we need to do the following settings to enable skew data optimization in the GROUP BY result.
>
> SET hive.groupby.skewindata=true;
>
> Once configured, Hive will first trigger an additional MapReduce job whose map output will randomly distribute to the reducer to avoid data skew.

**RELATED**

Hive Tutorial 2 - hive dml, hive inner join, hive outer join, hive cross join, hive map join, hive left semi join, hive union all , hive union, hive intercept and hive minus (https://timepasstechies.com/hive-tutorial-2-hive-dml-hive-inner-join-hive-outer-join-hive-cross-join-hive-map-join-hive-left-semi-join-hive-union-hive-union-hive-intercept-hive-minus/)

spark dataframe and dataset loading and saving data, spark sql performance tuning - tutorial 19 (https://timepasstechies.com/spark-dataframe-dataset-loading-saving-data-spark-sql-performance-tuning/)
November, 2017
In "Data Analytics"

Hive tutorial 8 - Hive performance tuning using data file optimization using file format, compression and storage optimization (https://timepasstechies.com/hive-tutorial-8-hive-performance-tuning-using-data-file-optimization-using-file-format-compression-storage-optimization/)
August, 2017
In "Data Analytics"

August, 2017
In "Data Analytics"

Posted in: Data Analytics (https://timepasstechies.com/category/programming/data-analytics/), Hive
(https://timepasstechies.com/category/programming/data-analytics/hive/), performance tuning (https://timepasstechies.com/category/performance-
tuning/)
Filed under: hive (https://timepasstechies.com/tag/hive/), hive performance tuning (https://timepasstechies.com/tag/hive-performance-tuning/)

← Hive tutorial 8 – Hive performance tuning
using Job and query optimization with local
mode, jvm reuse and parallel execution
(https://timepasstechies.com/hive-tutorial-8-
hive-performance-tuning-using-job-query-
optimization-local-mode-jvm-reuse-parallel-
execution/)

Hive tutorial 10 – Hive example for writing
custom user defined function →
(https://timepasstechies.com/hive-tutorial-9-
hive-example-writing-custom-user-defined-
function/)

## LEAVE A REPLY

Your email address will not be published. Required fields are marked *

COMMENT

NAME *

EMAIL *

WEBSITE

☐ NOTIFY ME OF FOLLOW-UP COMMENTS BY EMAIL.
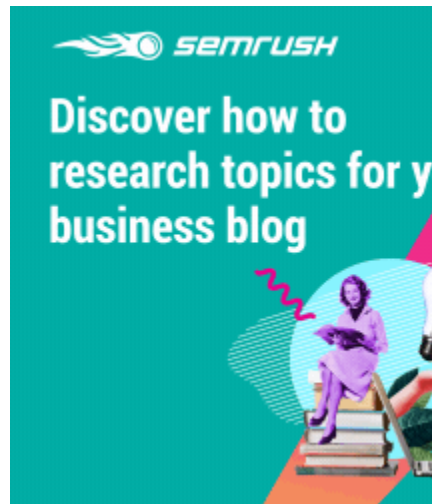
☐ NOTIFY ME OF NEW POSTS BY EMAIL.

POST COMMENT

Search …

**RECENT POSTS**

Scala code to get a secret stored in Azure key vault from databricks (https://timepasstechies.com/scala-code-to-get-a-secret-stored-in-azure-key-vault-from-databricks/)

How to read write data from Azure Blob Storage with Apache Spark (https://timepasstechies.com/how-to-read-write-data-from-azure-blob-storage-with-apache-spark/)

scala code to copy azure blob from one container to another (https://timepasstechies.com/scala-code-to-copy-azure-blob-from-one-container-to-another/)

HOME (HTTPS://TIMEPASSTECHIES.COM)    CONTACT ME (HTTPS://TIMEPASSTECHIES.COM/CONTACT/)    ABOUT ME (HTTPS://TIMEPASSTECHIES.COM/ABOUT/)