# Bulk Loading Your Big Data into Apache HBase, a Full Walkthrough

Jean-Daniel Cryans

Strata + Hadoop World NYC 2014

# About me

- Software Engineer at Cloudera, Storage team.
- Apache HBase committer since 2008, PMC member.

**cloudera**
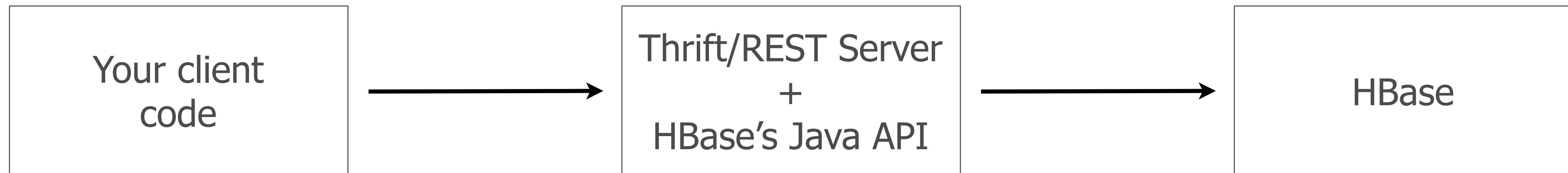Ask Bigger Questions

# Agenda

**1.HBase's write path**

2.Bulk loading concepts

3.ETL example

4.Issues and gotchas

**cloudera**
Ask Bigger Questions

# Getting your BIG data in HBase
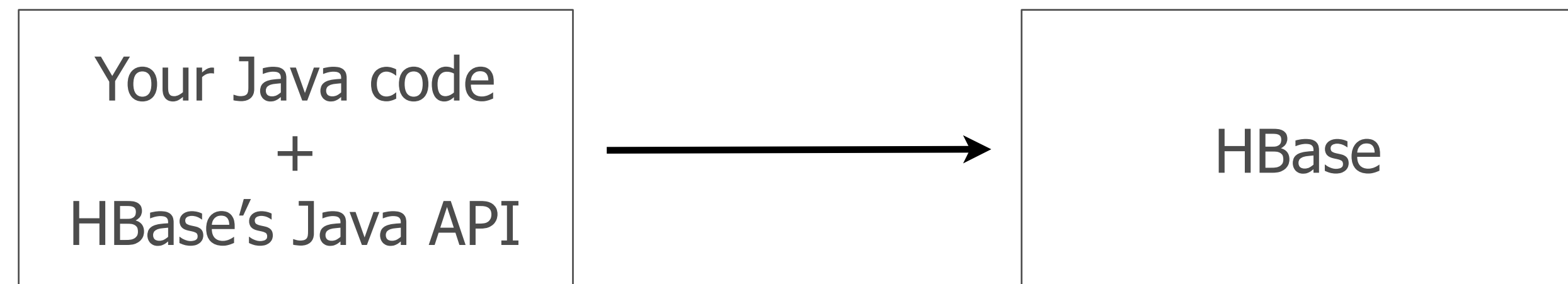
- Thrift/REST
- Java API
- MapReduce

**cloudera**®
Ask Bigger Questions

# Getting your BIG data in HBase

- Thrift/REST
  - Low throughput due to indirection.
  - Need a way to have many clients.

```
┌─────────────┐      ┌─────────────────┐      ┌──────────┐
│ Your client │ ───▶ │ Thrift/REST Server │ ───▶ │  HBase   │
│    code     │      │        +         │      │          │
│             │      │  HBase's Java API │      │          │
└─────────────┘      └─────────────────┘      └──────────┘
```
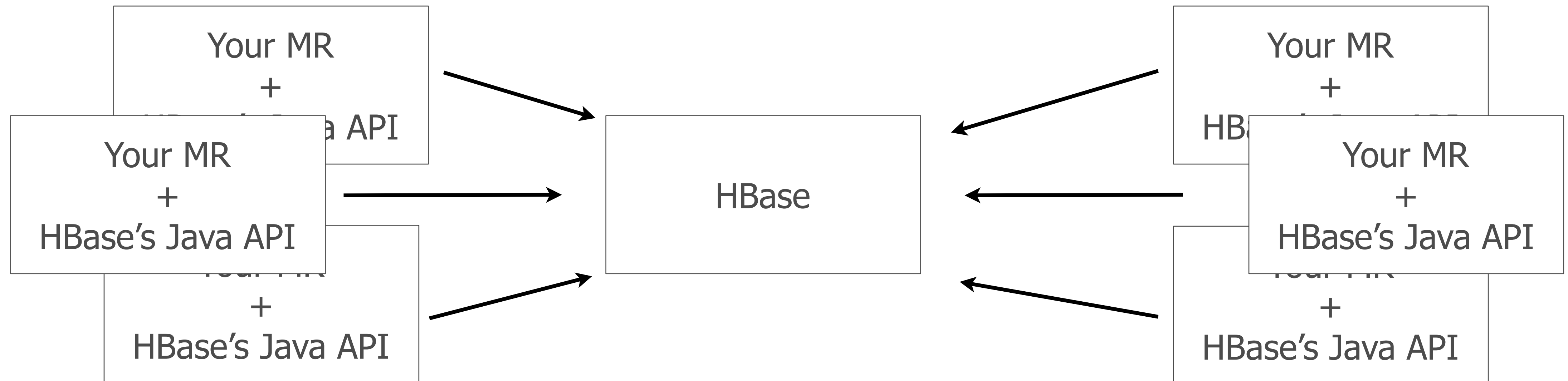
# Getting your BIG data in HBase

- Java API
  - Indirection problem is solved.
  - Still need a way to have many clients.

```
┌─────────────────────┐                    ┌─────────────────────┐
│   Your Java code     │                    │                     │
│         +            │ ────────────────▶  │       HBase         │
│   HBase's Java API   │                    │                     │
└─────────────────────┘                    └─────────────────────┘
```
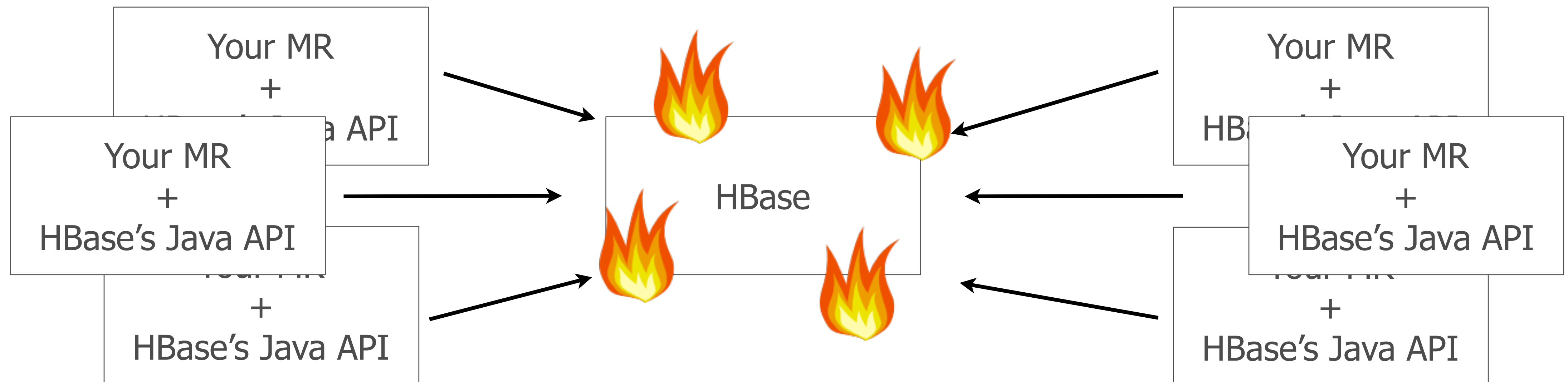
# Getting your BIG data in HBase

- ## MapReduce
  - ### No indirection.
  - ### No distribution problem, but...

**cloudera®**
Ask Bigger Questions

# Getting your BIG data in HBase

- MapReduce
  - No indirection.
  - No distribution problem, but...

**cloudera**®
Ask Bigger Questions

hbase-user@hadoop.apache.org

My region servers are always dying???

Hey list,

I'm using HBase 0.94 and trying to import a few TBs of data. Originally it was slow when sending the data from Python, I estimated it would take over a month in the best case, but now I wrote this MR job that's super fast for a few hours but then everything crashes!

When my region servers die I see a lot of HDFS stack traces and eventually there's a spooky YouAreDeadException.

Can someone help please?

Thx,

J-D

**cloudera**
Ask Bigger Questions

cloudera®
Ask Bigger Questions

# Log-structured merge-trees



$C_1$ tree       $C_0$ tree

Disk       Memory

http://www.cs.umb.edu/~poneil/lsmtree.pdf

cloudera®
Ask Bigger Questions
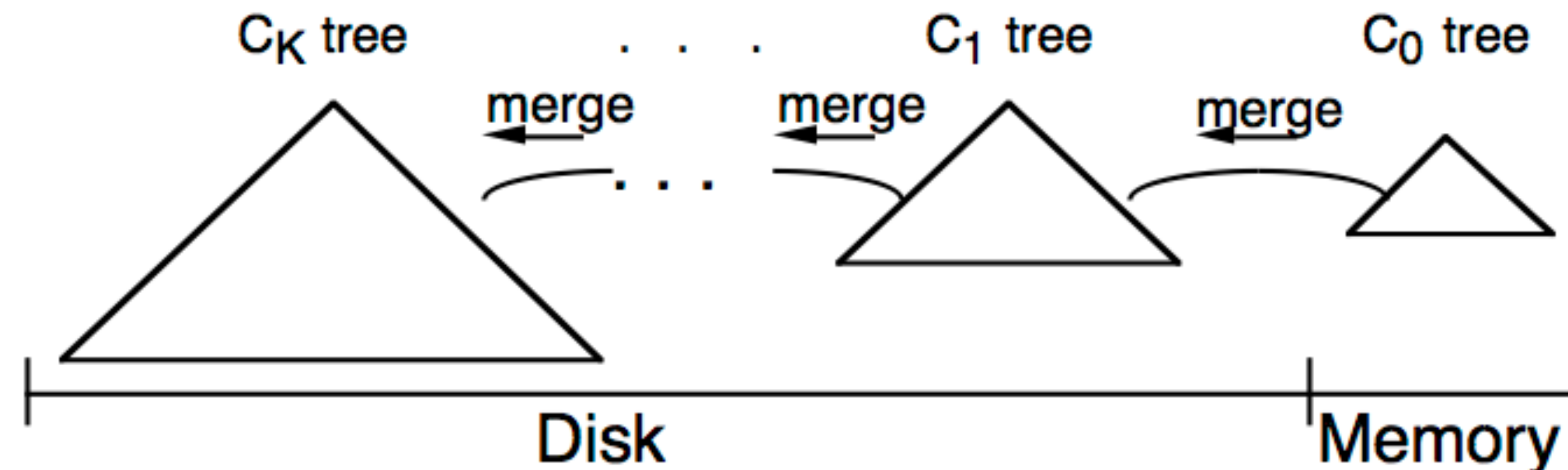
# A quick intro to LSM trees

- Data is written in memory to C0.
- C0 flushes upon reaching a certain threshold.
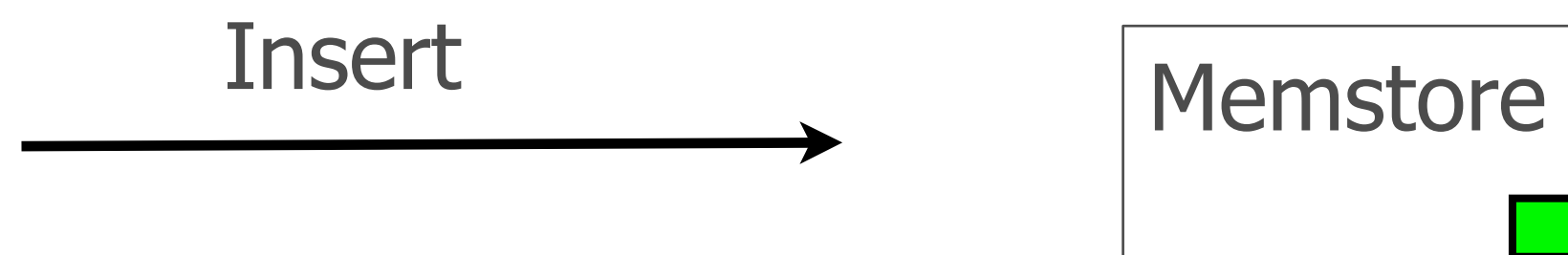- On-disk components C1-Ck are merged in the background.



Figure 3.1. An LSM-tree of K+1 components

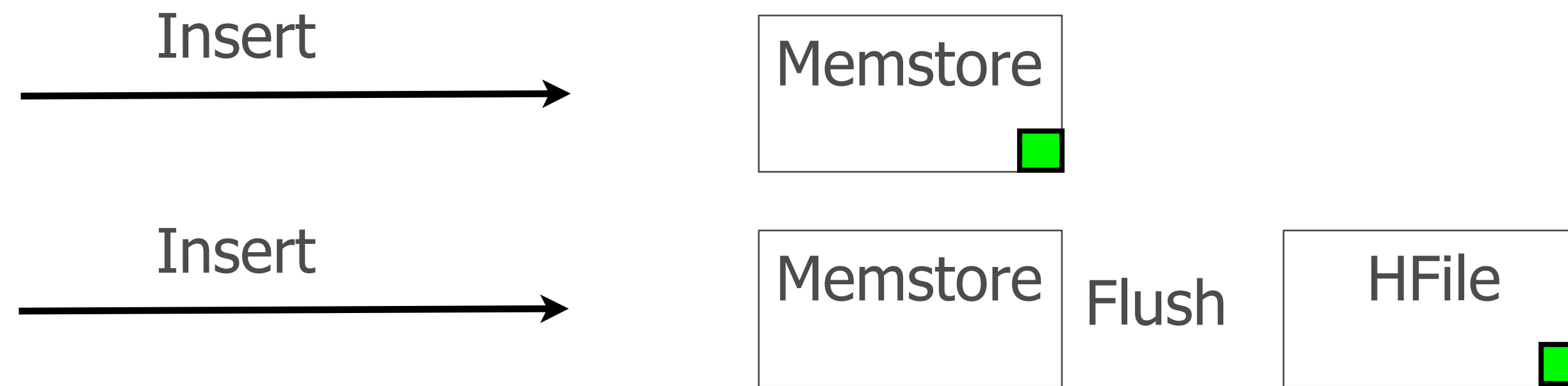http://www.cs.umb.edu/~poneil/lsmtree.pdf

# LSM trees in HBase

**cloudera®**
Ask Bigger Questions

# LSM trees in HBase

Insert →

Memstore

# LSM trees in HBase

Insert →

Memstore ▪

Insert →

Memstore    Flush    HFile ▪

**cloudera**®
Ask Bigger Questions

# LSM trees in HBase

Insert → Memstore
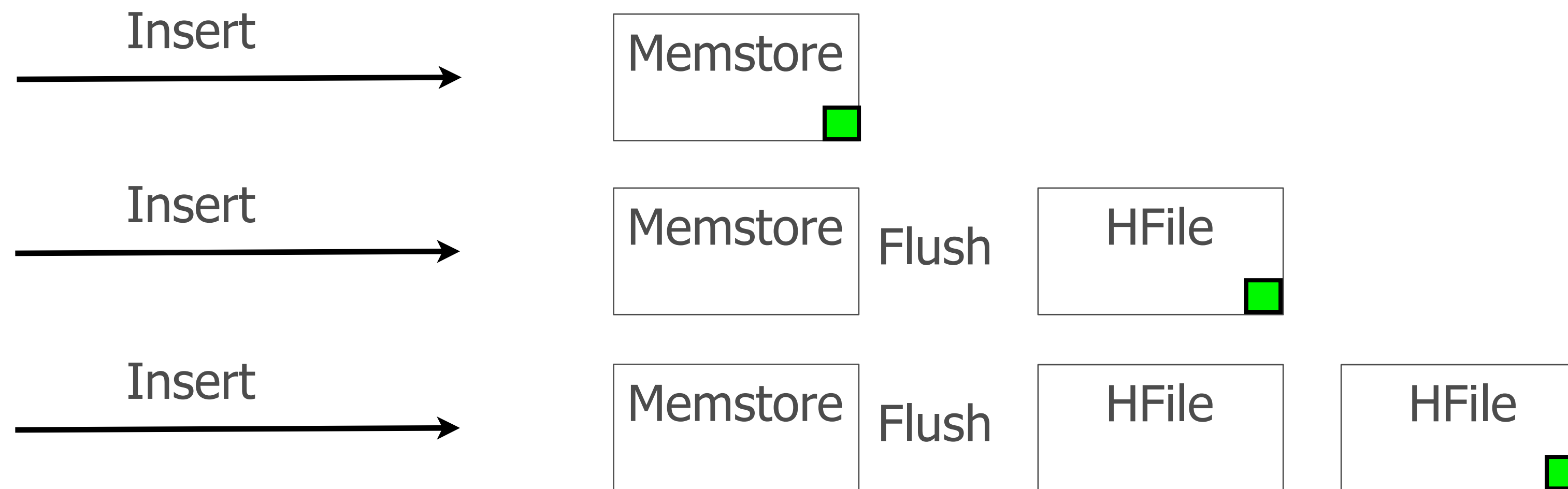
Insert → Memstore Flush HFile

Insert → Memstore HFile

**cloudera®**
Ask Bigger Questions

# LSM trees in HBase

# LSM trees in HBase

Insert →   Memstore 🟩

Insert →   Memstore   Flush   HFile 🟩

Insert →   Memstore   Flush   HFile   HFile 🟩

Insert →   Memstore   HFile   HFile 🟩

**cloudera®**
Ask Bigger Questions

# LSM trees in HBase

Insert → Memstore

Insert → Memstore  Flush  HFile

Insert → Memstore  Flush  HFile  HFile

Insert → Memstore  Flush  HFile  HFile  HFile

cloudera®
Ask Bigger Questions

# LSM trees in HBase

Insert →  | Memstore 🟩 |

Insert →  | Memstore | Flush | HFile 🟩 |

Insert →  | Memstore | Flush | HFile | HFile 🟩 |

Insert →  | Memstore | Flush | HFile | HFile | HFile 🟩 |

| HFile 🟩 |  Compaction!

cloudera®
Ask Bigger Questions

# LSM trees in HBase

Inserts

Inserts

**cloudera®**
Ask Bigger Questions

# LSM trees in HBase

Inserts

→
→
→

Inserts

→
→
→

Memstore

**cloudera**
Ask Bigger Questions

# LSM trees in HBase

Inserts

Memstore

Flush!
Flush!
Flush!

Inserts

**cloudera®**
Ask Bigger Questions

# LSM trees in HBase

Inserts

Inserts

Memstore

Flush!
Flush!
Flush!

HFile

HFile

HFile

# LSM trees in HBase

Inserts

Memstore

Inserts

Memstore

Flush!
Flush!
Flush!

HFile

HFile

HFile

**cloudera**®
Ask Bigger Questions

# LSM trees in HBase

Inserts

Inserts

Memstore

Memstore

Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!

HFile

HFile

HFile

cloudera®
Ask Bigger Questions

# LSM trees in HBase

Inserts

Memstore

Inserts

Memstore

Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!

| HFile | | HFile | | HFile | |
|---|---|---|---|---|---|
| HFile | HFile | HFile | HFile | HFile | |
| HFile | HFile | HFile | | | |

# LSM trees in HBase

Inserts

Memstore

Inserts

Memstore

Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!

| HFile | HFile | HFile |
|-------|-------|-------|

| HFile | HFile | HFile | HFile | HFile |
|-------|-------|-------|-------|-------|
| HFile | HFile | HFile | | |

Compaction!

# LSM trees in HBase

Inserts

Memstore

Inserts

Memstore

Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!
Flush!

| HFile | HFile | HFile |
|-------|-------|-------|

| HFile | HFile | HFile | HFile | HFile |
|-------|-------|-------|-------|-------|
| HFile | HFile | HFile | | |

Compaction!

- How many times will data be rewritten?
- What kind of tuning could make this better?
- What about splitting those regions?

cloudera®
Ask Bigger Questions

# LSM trees in HBase

HFile

Or is there a way to just get the final result directly in HBase?

**cloudera**®
Ask Bigger Questions

# Agenda

1. HBase's write path
2. **Bulk loading concepts**
3. ETL example
4. Issues and gotchas

**cloudera**®
Ask Bigger Questions

# Bulk loading overview

- Goal: generate data files in HBase's own format, respecting the region boundaries, and give them to the region servers.

- Use cases:

**cloudera**®
Ask Bigger Questions

# Bulk loading overview

- Goal: generate data files in HBase's own format, respecting the region boundaries, and give them to the region servers.
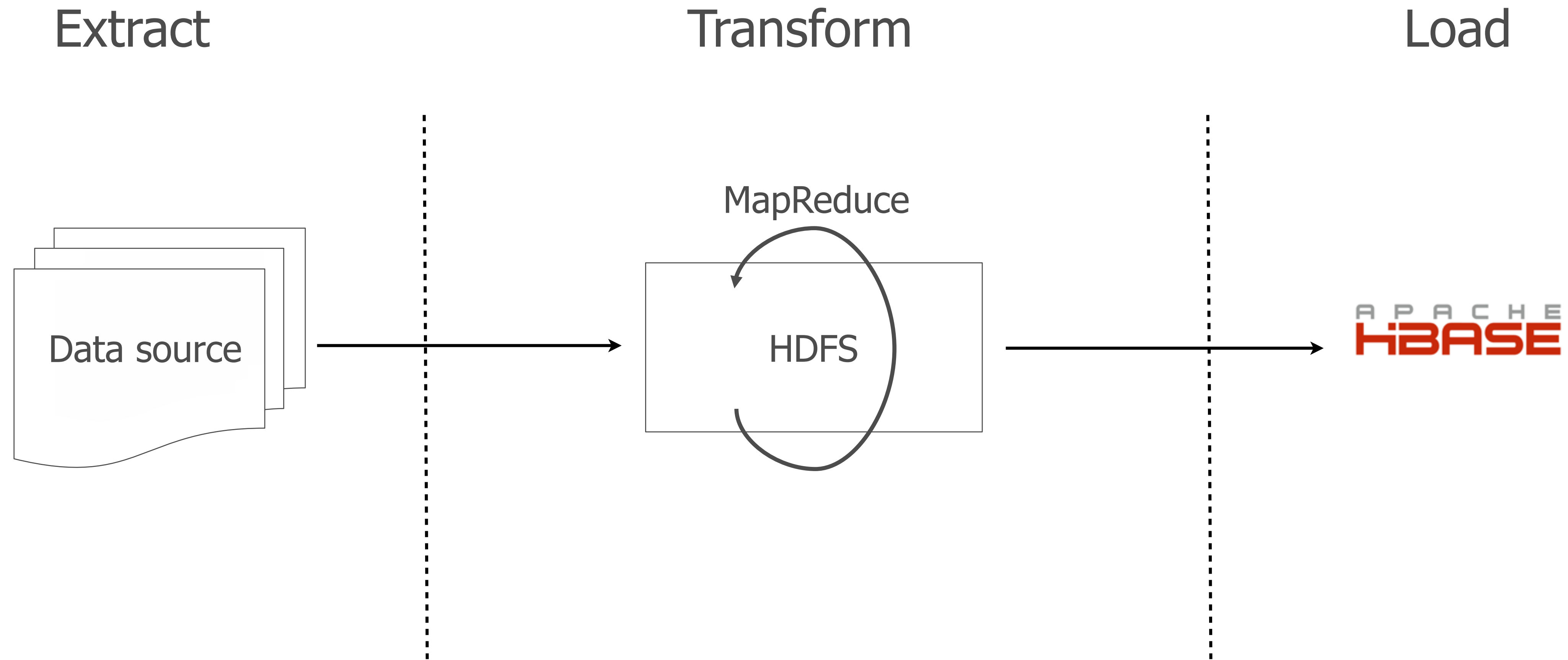
- Use cases:

  Initial Data Import
  Example:

**cloudera**®
Ask Bigger Questions

# Bulk loading overview

- Goal: generate data files in HBase's own format, respecting the region boundaries, and give them to the region servers.
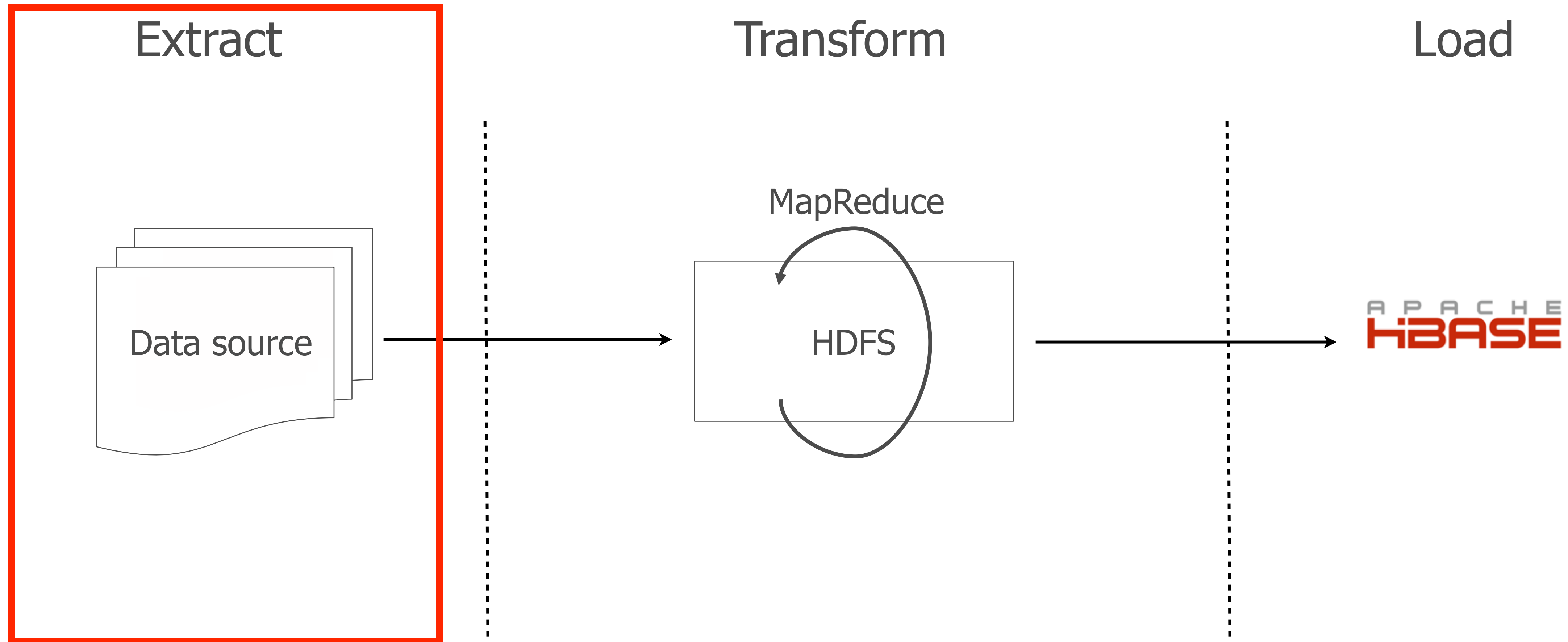
- Use cases:

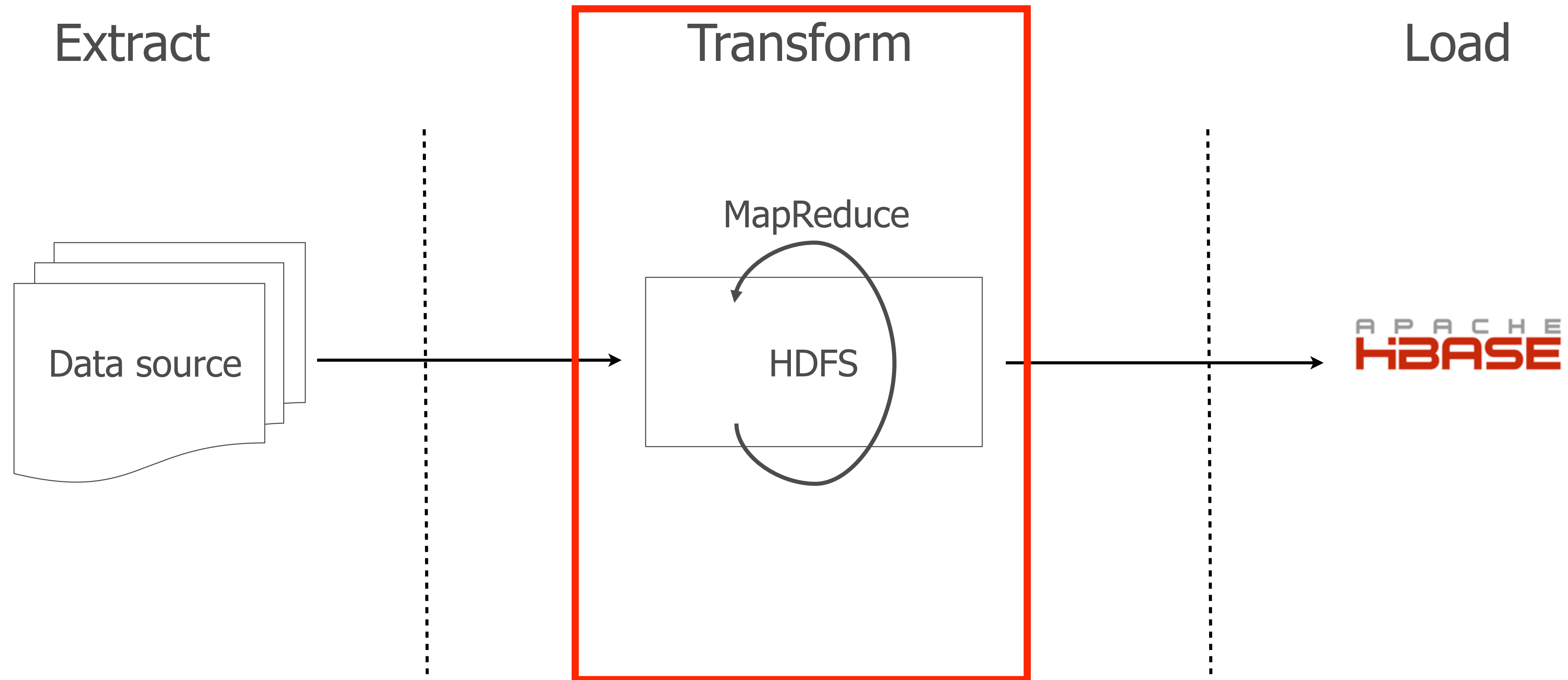Initial Data Import
Example:



Regular Imports
Example:

17

cloudera®
Ask Bigger Questions

# Bulk loading data flow

Extract                    Transform                    Load

MapReduce

Data source            HDFS

APACHE
HBASE

cloudera®
Ask Bigger Questions

# Bulk loading data flow

| Extract | Transform | Load |
|---------|-----------|------|
| Data source | MapReduce HDFS | APACHE HBASE |

19

**cloudera**
Ask Bigger Questions

# Bulk loading data flow

Extract

Transform

Load

MapReduce

Data source

HDFS

APACHE
**HBASE**

**cloudera**
Ask Bigger Questions

# Transforming data into HFiles

```java
HTable table = new HTable(conf, tableName);
job.setReducerClass(PutSortReducer.class);
Path outputDir = new Path(hfileOutPath);
FileOutputFormat.setOutputPath(job, outputDir);
job.setMapOutputKeyClass(ImmutableBytesWritable.class);
job.setMapOutputValueClass(Put.class);
HFileOutputFormat.configureIncrementalLoad(job, table);
```

**cloudera®**
Ask Bigger Questions

# Transforming data into HFiles

```
HFileOutputFormat.configureIncrementalLoad(job, table);
```

**cloudera**®
Ask Bigger Questions

# Transforming data into HFiles

```
HFileOutputFormat.configureIncrementalLoad(job, table);
```

Mapper 1

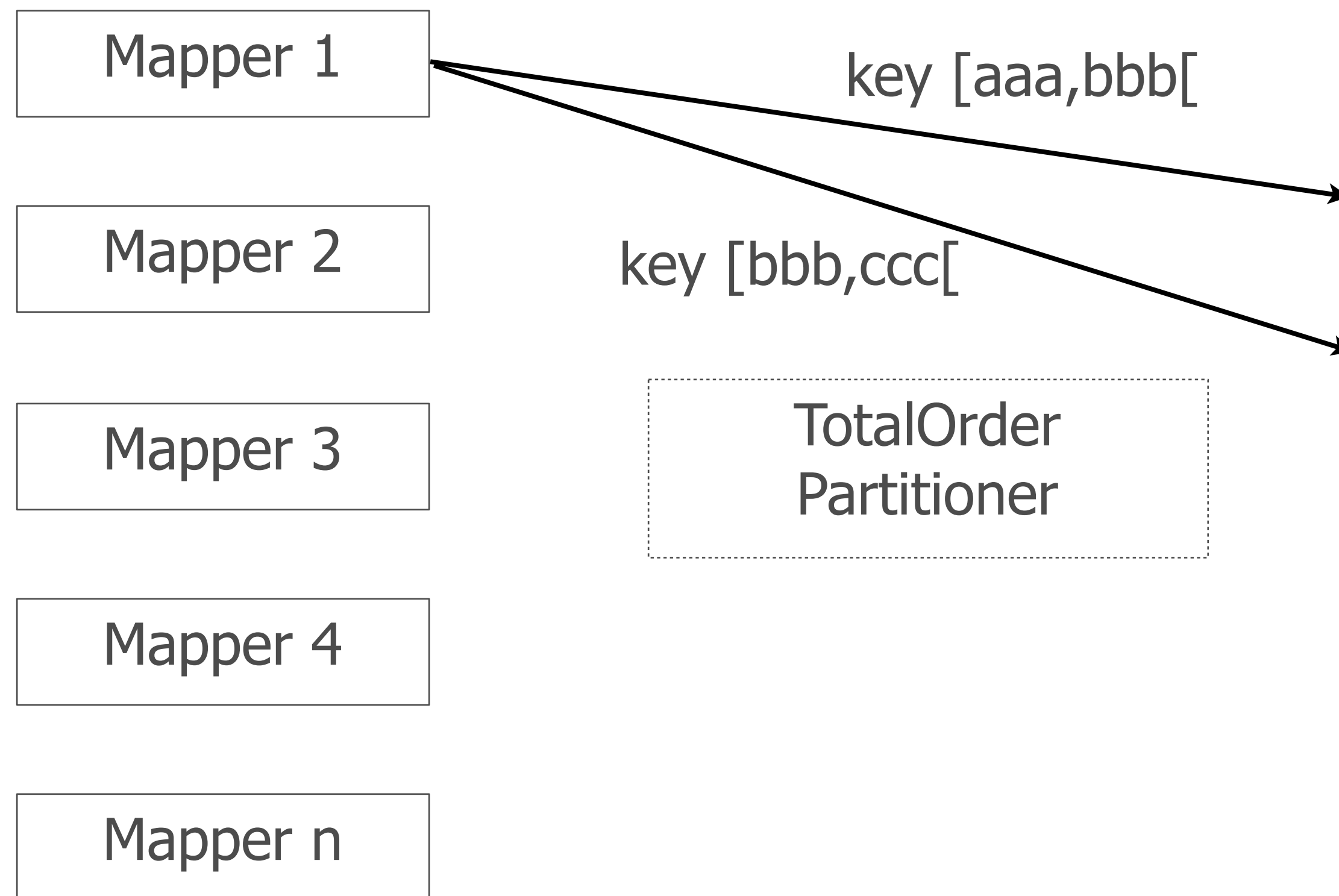Mapper 2

Mapper 3

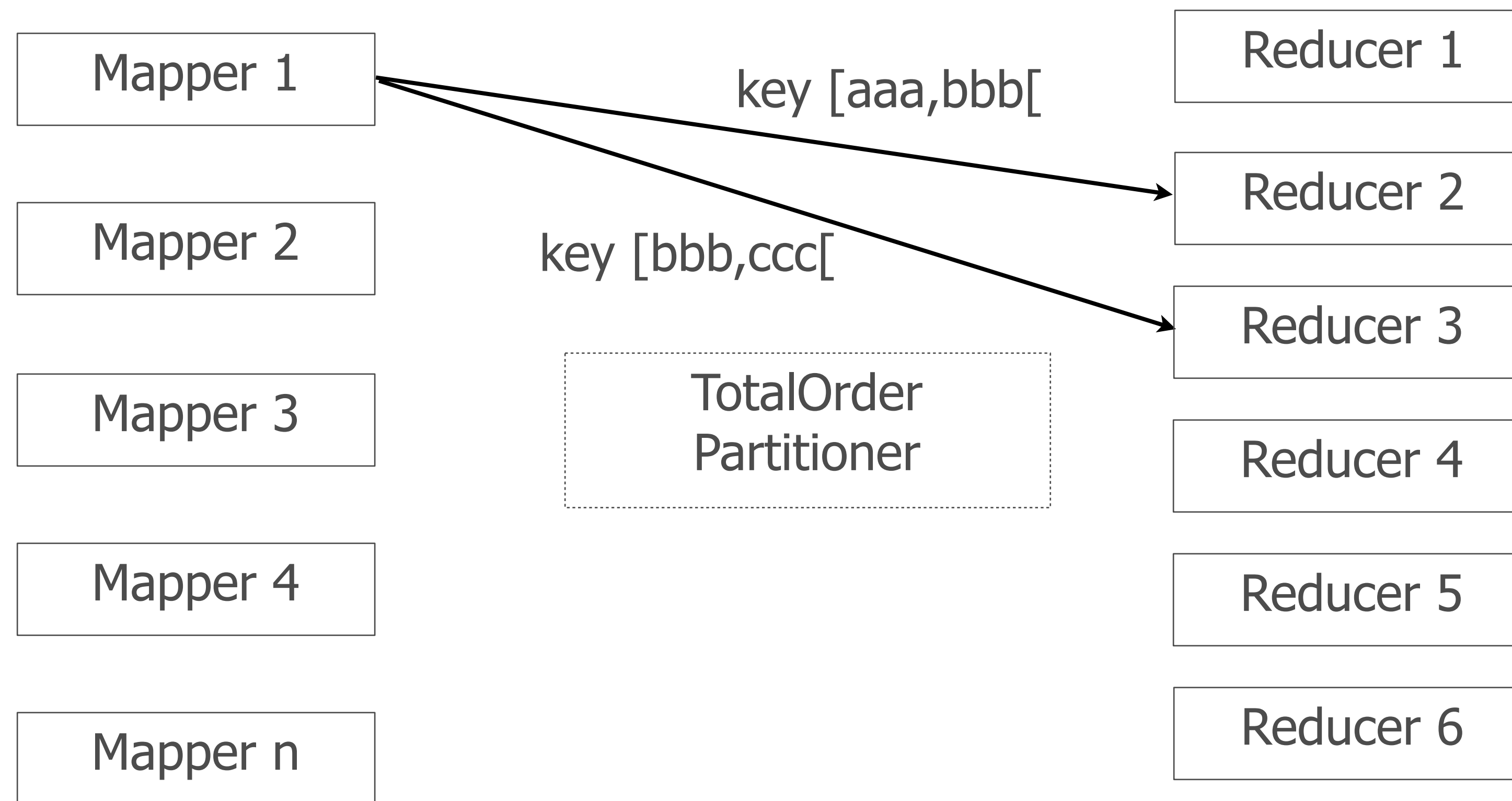Mapper 4

Mapper n

**cloudera®**
Ask Bigger Questions

# Transforming data into HFiles

```
HFileOutputFormat.configureIncrementalLoad(job, table);
```

# Transforming data into HFiles

```
HFileOutputFormat.configureIncrementalLoad(job, table);
```



Mapper 1

Mapper 2

Mapper 3

Mapper 4

Mapper n

key [aaa,bbb[

key [bbb,ccc[

TotalOrder
Partitioner

Reducer 1

Reducer 2

Reducer 3

Reducer 4

Reducer 5

Reducer 6

Each reducer outputs
one file per region.

**cloudera**®
Ask Bigger Questions

# Bulk loading data flow

Extract            Transform            Load

Data source   →   MapReduce / HDFS   →   APACHE HBASE

cloudera
Ask Bigger Questions
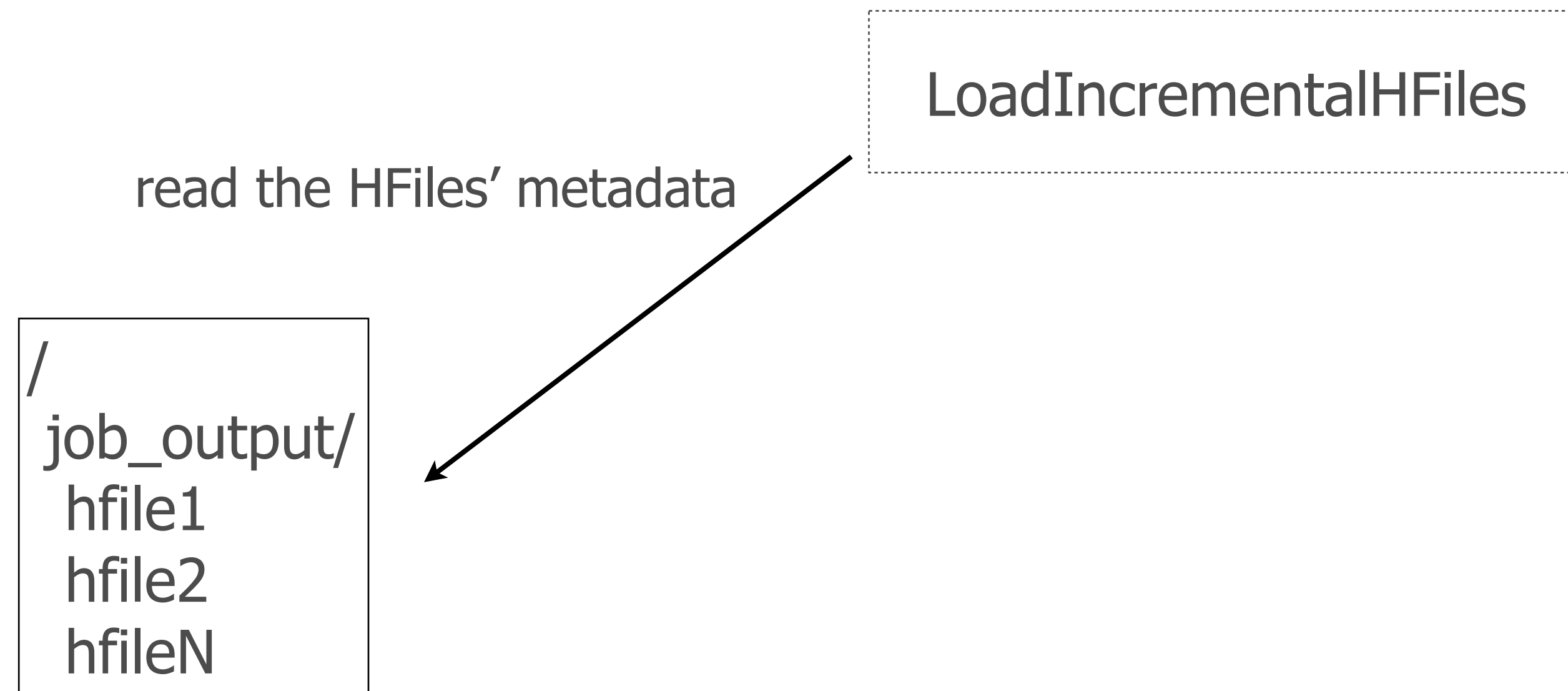
# Loading HFiles

```
$ hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles
<files_location> <table_name>
```
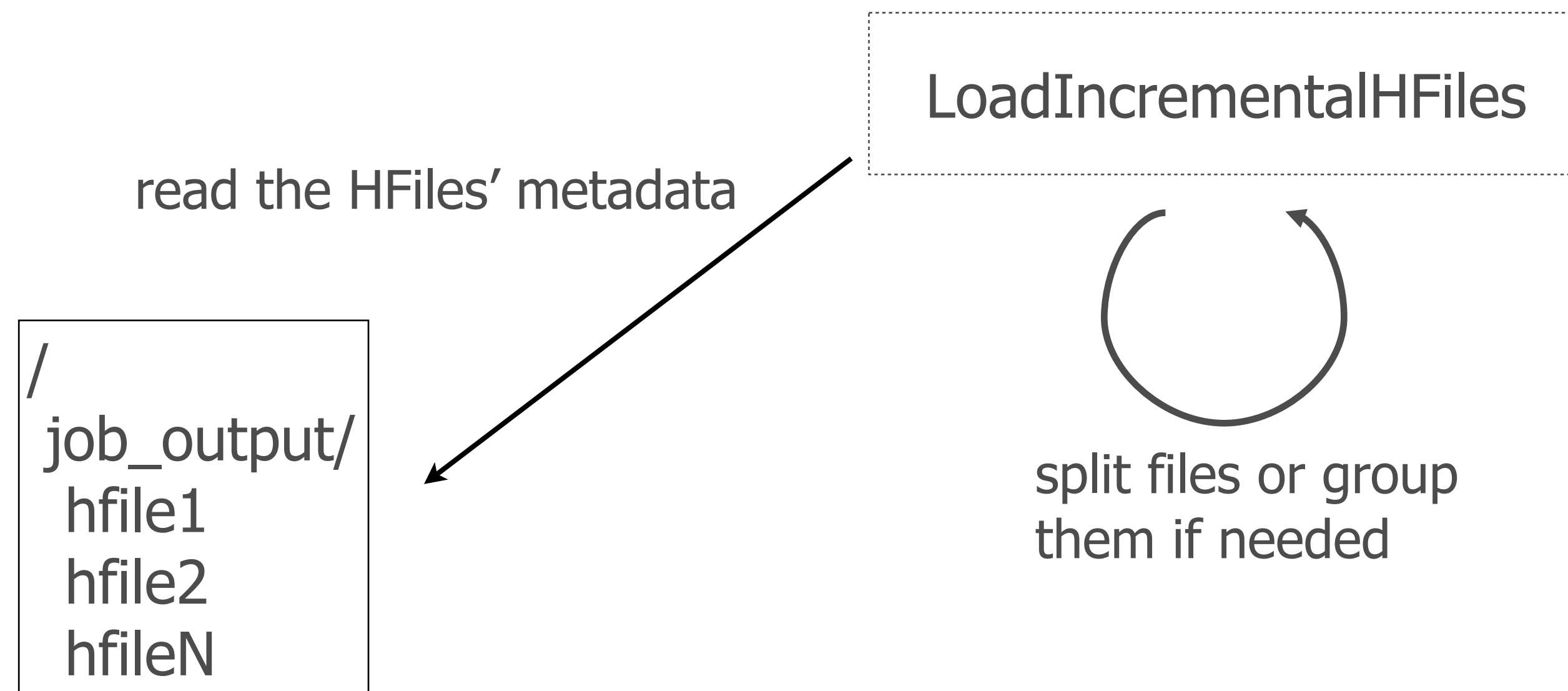
```
/
 job_output/
  hfile1
  hfile2
  hfileN
```

# Loading HFiles

```
$ hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles
<files_location> <table_name>
```
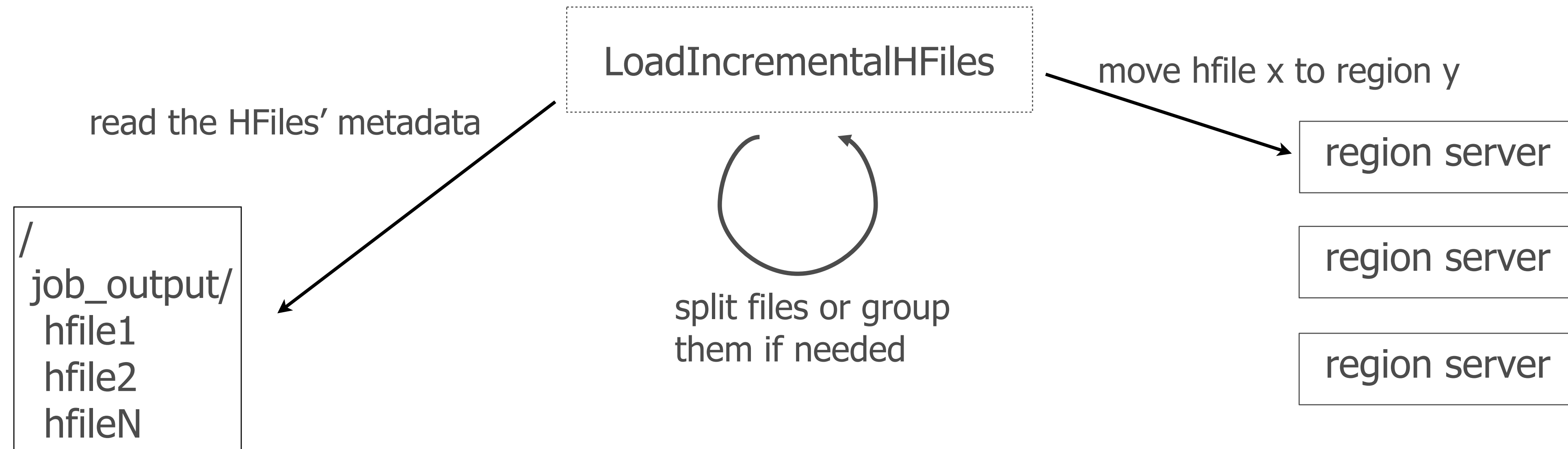
LoadIncrementalHFiles

read the HFiles' metadata

```
/
 job_output/
  hfile1
  hfile2
  hfileN
```

# Loading HFiles

```
$ hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles
<files_location> <table_name>
```

LoadIncrementalHFiles

read the HFiles' metadata

/
 job_output/
  hfile1
  hfile2
  hfileN

split files or group
them if needed

# Loading HFiles

```
$ hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles
<files_location> <table_name>
```

LoadIncrementalHFiles

move hfile x to region y

read the HFiles' metadata

/
 job_output/
  hfile1
  hfile2
  hfileN

split files or group
them if needed

region server

region server

region server

**cloudera**®
Ask Bigger Questions

# Agenda

1. HBase's write path
2. Bulk loading concepts
3. **ETL example**
4. Issues and gotchas

cloudera®
Ask Bigger Questions

# MySQL Import

- Extract
  - CSV dump into file.
- Transform
  - Map columns, create HFiles.
- Load
  - Use LoadIncrementalHFiles.
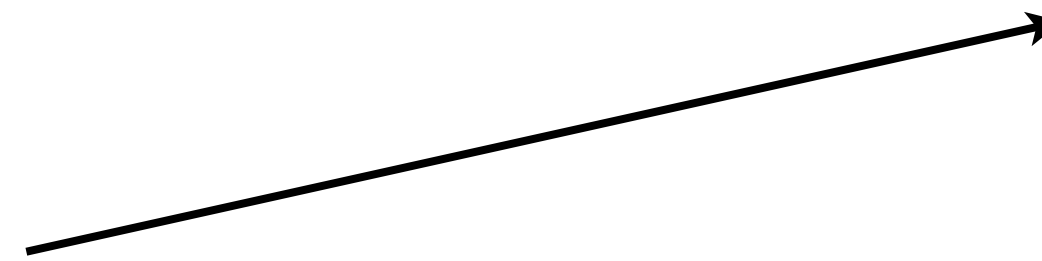
**cloudera**
Ask Bigger Questions

# Extract



```
SELECT * INTO OUTFILE 'dump.csv'
    FIELDS TERMINATED BY ',' OPTIONALLY ENCLOSED BY '"'
    LINES TERMINATED BY '\n'
FROM table
```

hdfs dfs -put dump.csv

dump.csv

# Transform



Map dump.csv

Reduce to output/

```
hadoop jar /usr/lib/hbase/hbase-0.98.6-cdh5.2.0-security.jar importtsv
-Dimporttsv.separator=,
-Dimporttsv.bulk.output=output
-Dimporttsv.columns=HBASE_ROW_KEY,f:col1,f:col2 table-name dump.csv
```
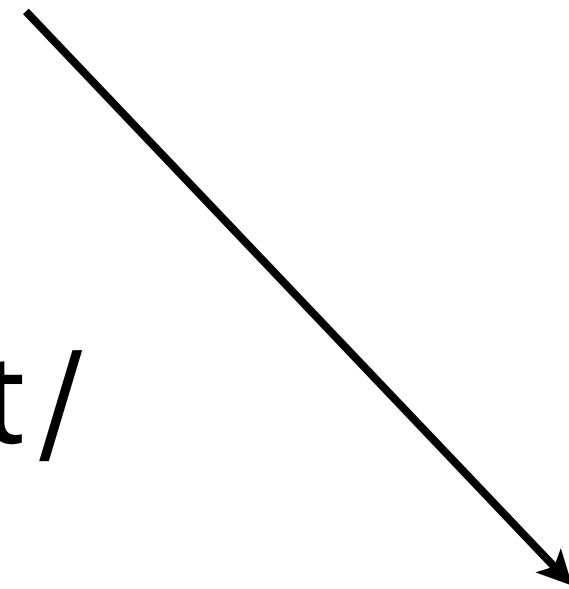
# Transform

Map dump.csv

Reduce to output/

```
hadoop jar /usr/lib/hbase/hbase-0.98.6-cdh5.2.0-security.jar importtsv
-Dimporttsv.separator=,
-Dimporttsv.bulk.output=output
-Dimporttsv.columns=HBASE_ROW_KEY,f:col1,f:col2 table-name dump.csv
```
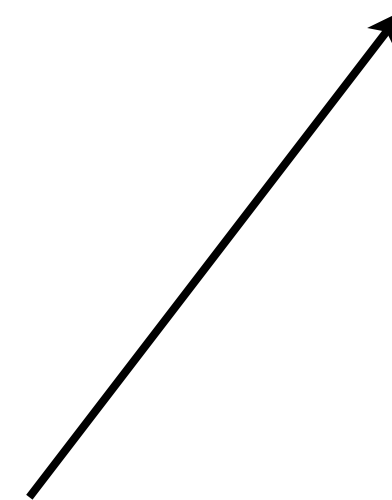
# Transform



Map dump.csv

Reduce to output/

```
hadoop jar /usr/lib/hbase/hbase-0.98.6-cdh5.2.0-security.jar importtsv
-Dimporttsv.separator=,
-Dimporttsv.bulk.output=output
-Dimporttsv.columns=HBASE_ROW_KEY,f:col1,f:col2 table-name dump.csv
```

# Load



List the files
under output/

Tell each RS
to move them.

`hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles output table-name`

# Agenda

1. HBase's write path
2. Bulk loading concepts
3. ETL example
4. **Issues and gotchas**

# Planning the bulk load; gotchas

- Initial import
  - Tables must still be created, pre-split.

```
create 'table-name', {NAME => 'f'},   {SPLITS => ['a', 'b', 'c', 'd']}
```

  - Plan for the files to fit in the regions else it will split.

```
alter 'table-name', {MAX_FILESIZE => 10737418240}
```

**cloudera**®
Ask Bigger Questions

# Planning the bulk load; gotchas

- ## Regular import
  - ### Loading data on HDFS still not free, IO-wise.
    - Especially the Transform phase.

  - ### Data won't be in the block cache once Loaded.

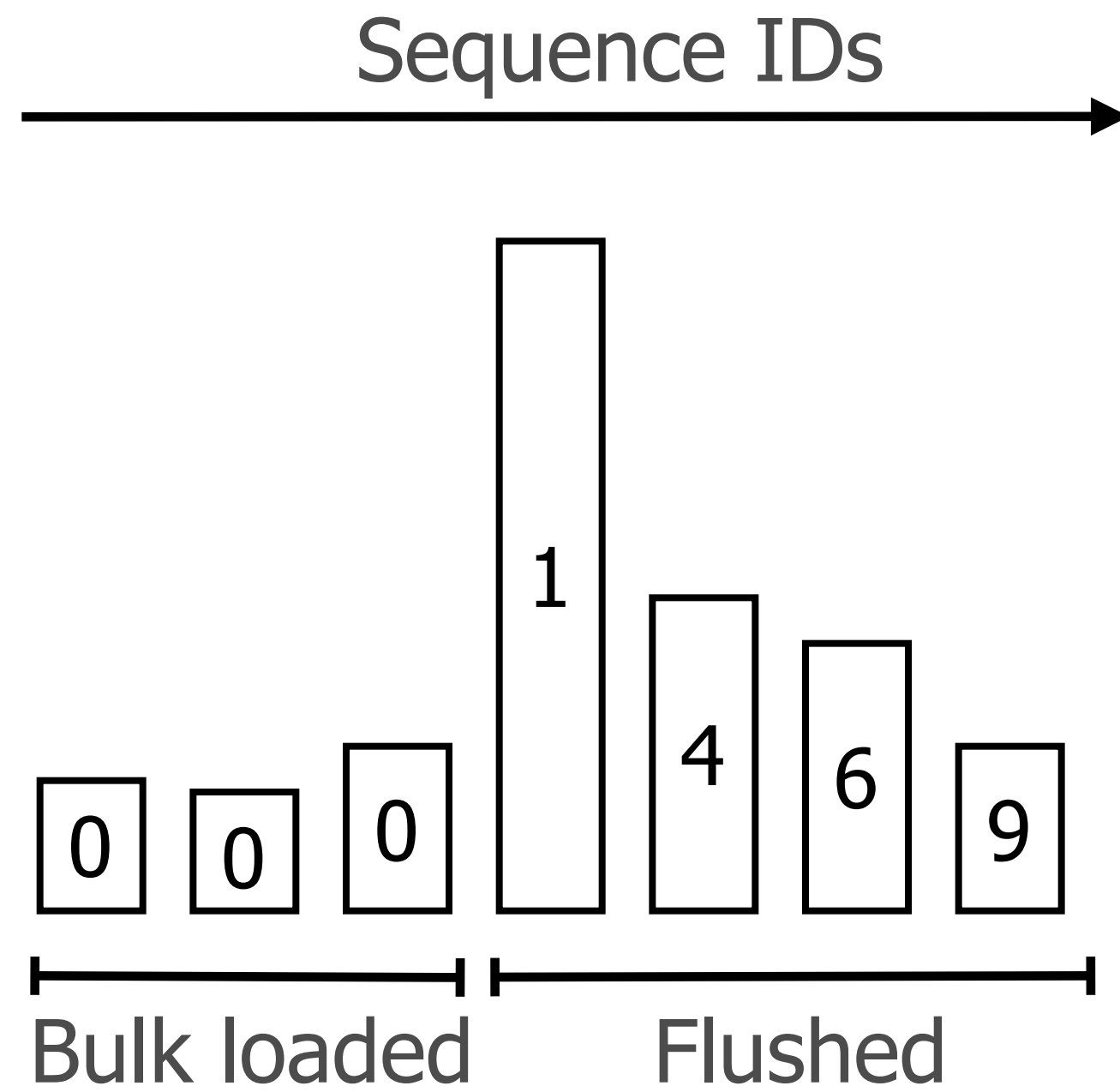| Cache Hit Ratio |
|---|
| 70% |

  - ### Block locality isn't guaranteed.

| Block locality |
|---|
| 0 |

**cloudera**
Ask Bigger Questions

# Gotchas: Security

- Problem:
  - The user "hbase" must move files it doesn't have access to.
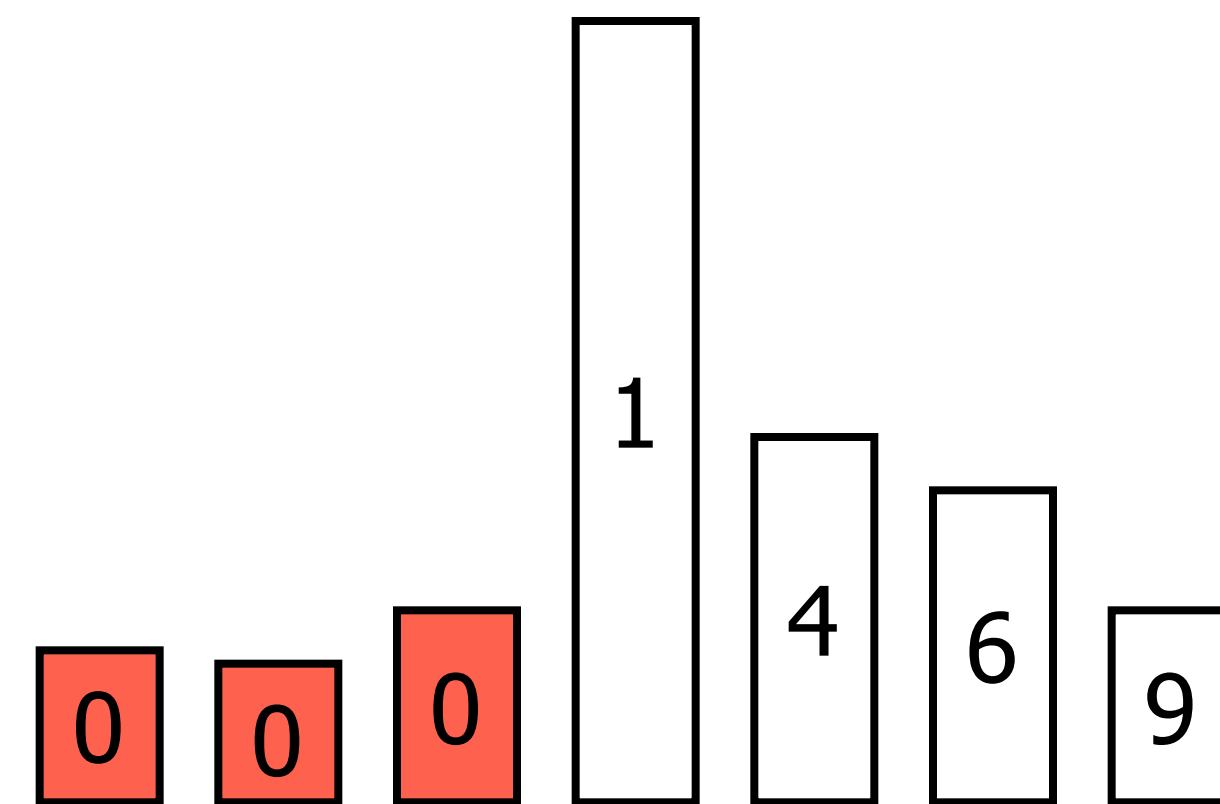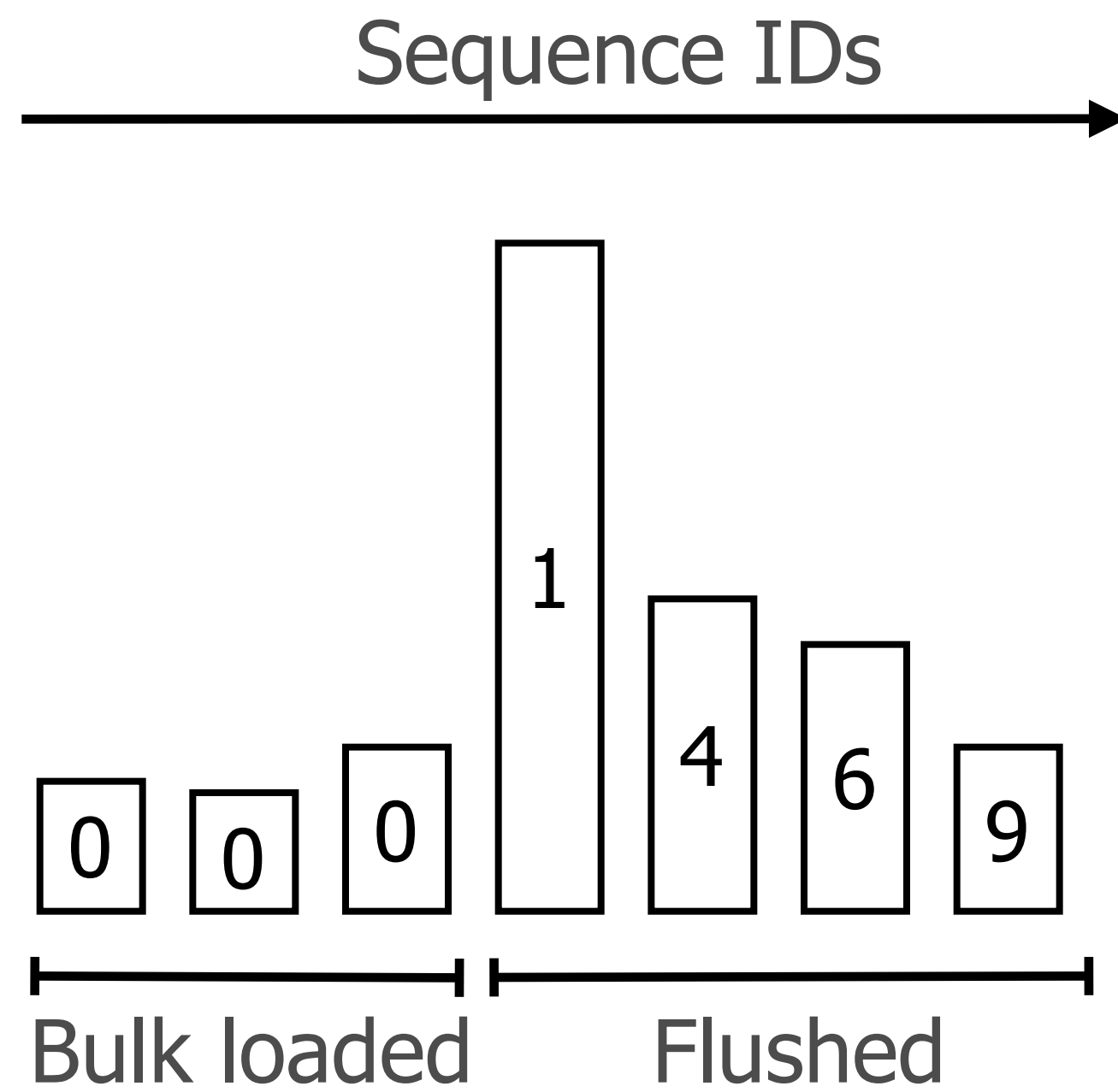
**cloudera**®
Ask Bigger Questions

# Gotchas: Security

- ## Problem:
  - The user "hbase" must move files it doesn't have access to.

- ## SecureBulkLoadEndpoint
  - Must be installed as part of enabling security.
  - A secret staging directory with 777 perms is used.
  - LoadIncrementalHFiles moves files there and then the RS moves it into its regions' directories.
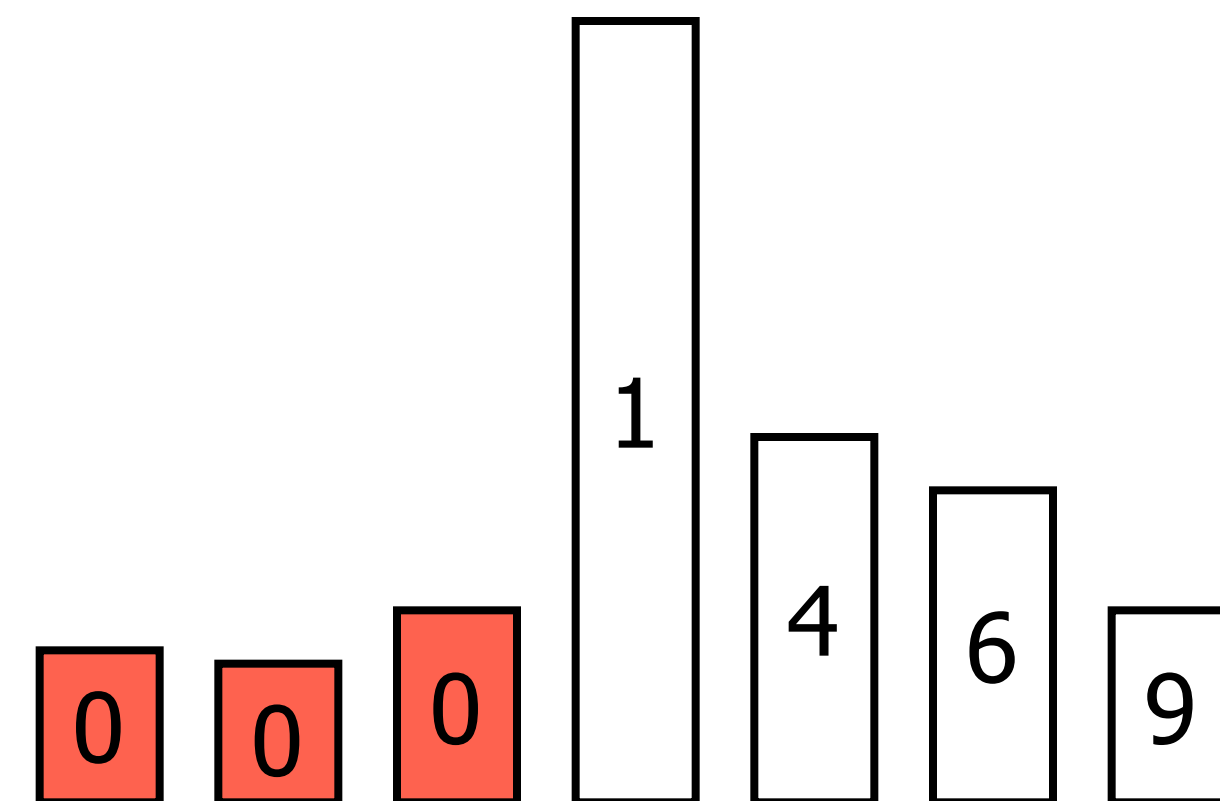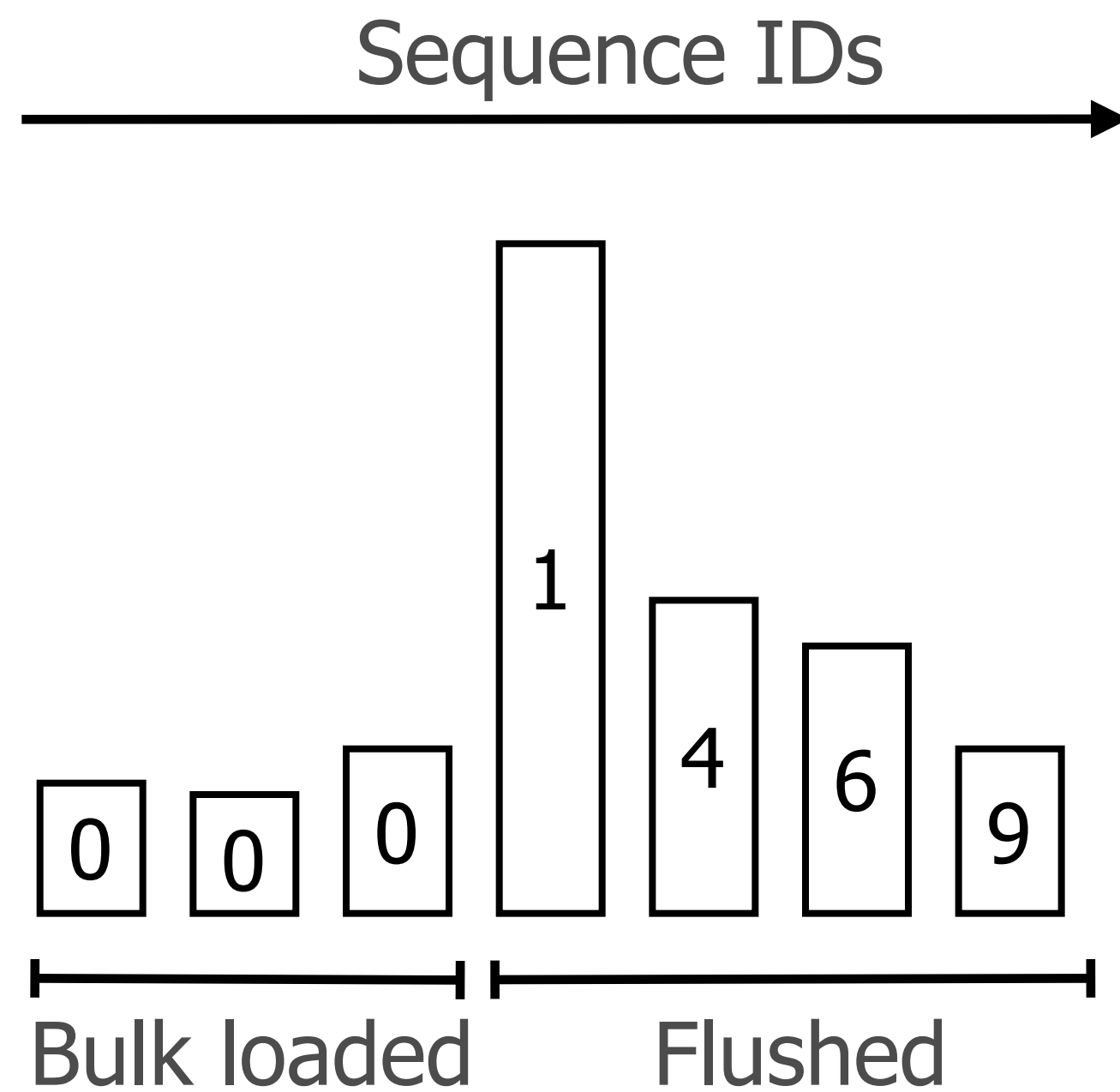
# Gotchas: HBASE–8521 (fixed 0.94.13+)

Sequence IDs



Bulk loaded     Flushed

**cloudera**®
Ask Bigger Questions

# Gotchas: HBASE-8521 (fixed 0.94.13+)



Sequence IDs

Bulk loaded    Flushed

1 4 6 9

1 4 6 9

3 HFiles fit the compaction
selection criteria...

# Gotchas: HBASE-8521 (fixed 0.94.13+)



Sequence IDs

Bulk loaded    Flushed

3 HFiles fit the compaction selection criteria...

Major compaction

It means we also select all the files that come after!

cloudera
Ask Bigger Questions

# Gotchas: HBASE-8521 (fixed 0.94.13+)

Sequence IDs →

Minor compaction

1

4

6

9

0 0 0

Bulk loaded    Flushed

1

4

6

9

0 0 0

3 HFiles fit the compaction
selection criteria...

1

4

6

9    11    14    17

Solution: assign sequence IDs to the bulk loaded files. Optional in 0.94, on by default in 0.96+.

**cloudera**
Ask Bigger Questions

Gotchas: HBASE-8283 (fixed 0.94.9+)

Sequence IDs

1

0  0  0      4    6    9

Bulk loaded      Flushed

# Gotchas: HBASE–8283 (fixed 0.94.9+)

Sequence IDs →

1

0 0 0 | 4 6 9

Bulk loaded | Flushed

Option 1 | Option 2

1

0 0 0 | 4 6 9

New selection algorithm considers multiple alternatives and doesn't work in only one direction
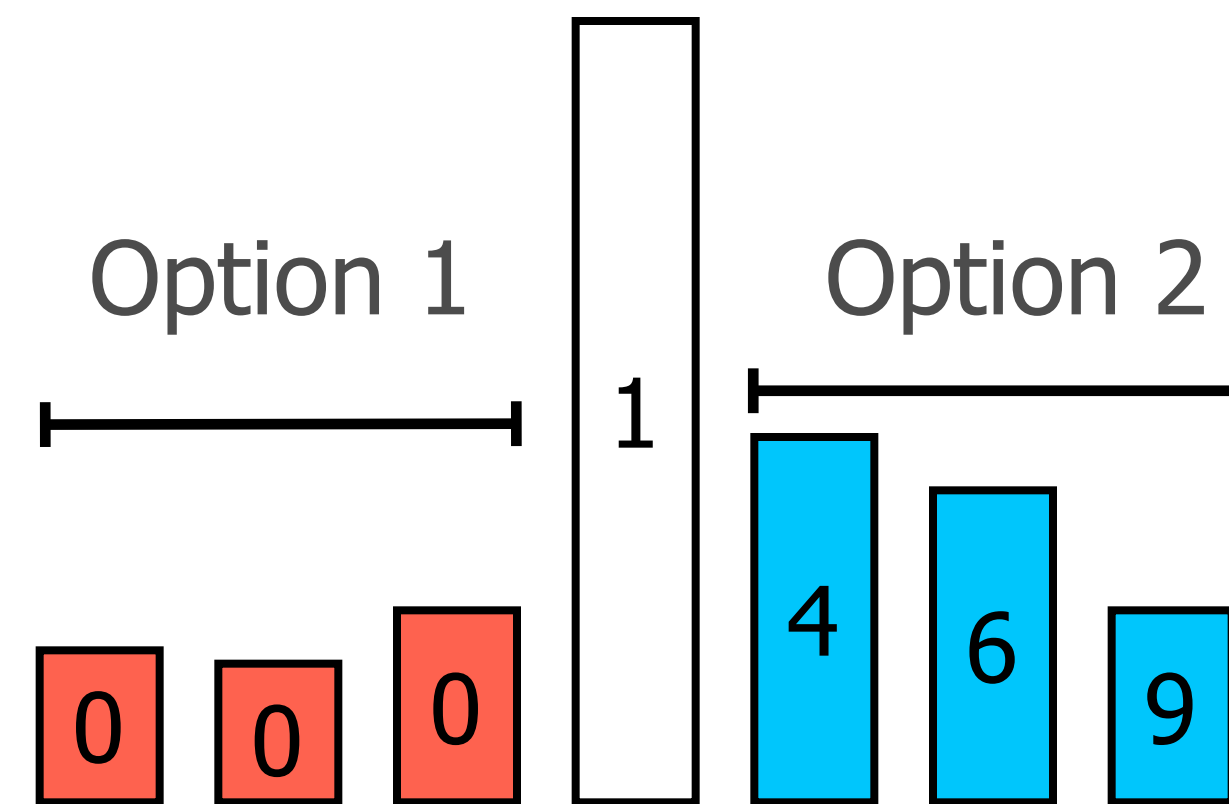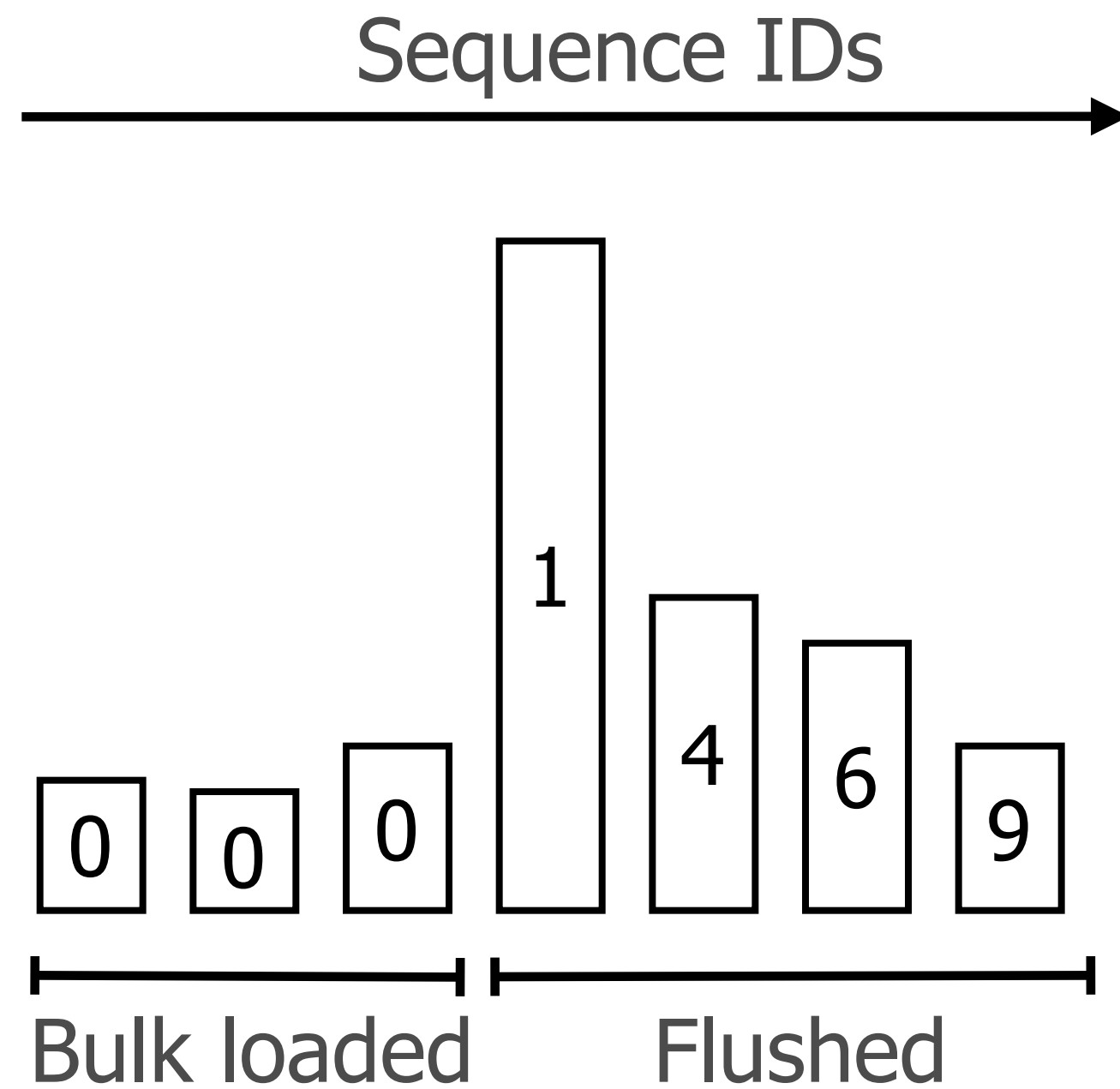
cloudera
Ask Bigger Questions

# Gotchas: HBASE–8283 (fixed 0.94.9+)

Sequence IDs
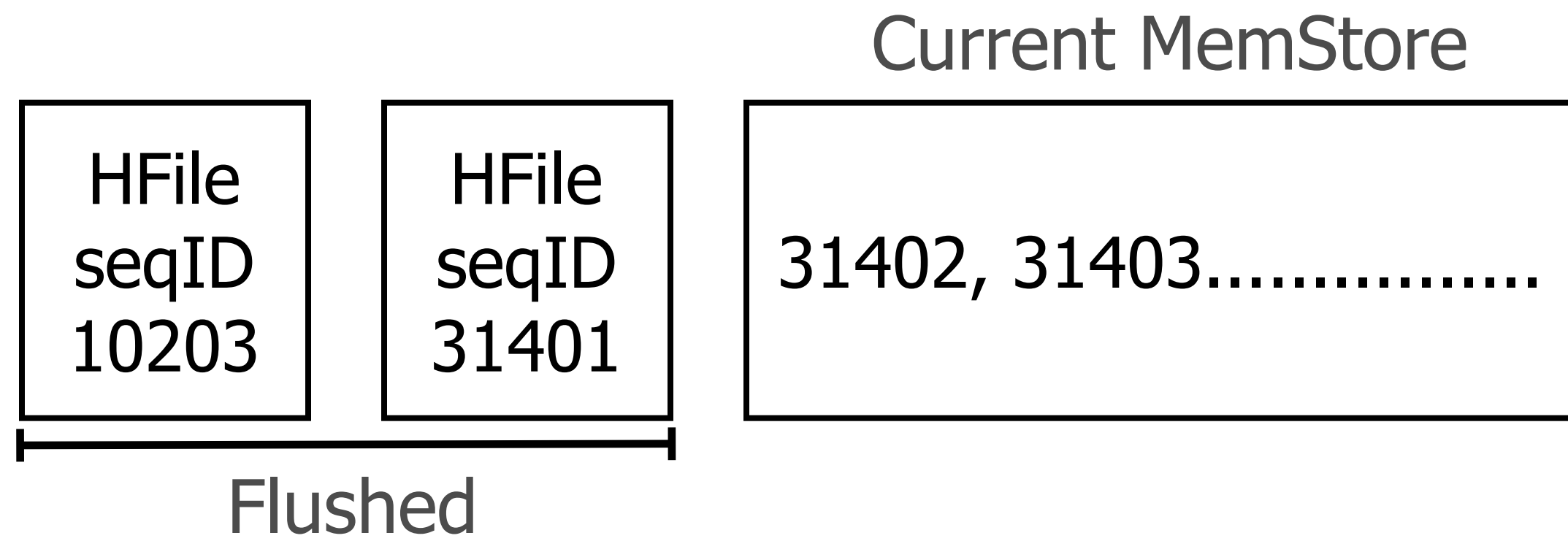
New selection algorithm considers multiple alternatives and doesn't work in only one direction

The selection that lowers the amount of seeks while compacting the least amount of data is chosen

# HBASE–10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 | 31402, 31403................ |

Flushed

# HBASE-10958 AKA Blindspot

Current MemStore

HFile
seqID
10203

HFile
seqID
31401

31402, 31403................

Flushed

The bulk loaded file is assigned
a sequence ID that ends up
somewhere in the MemStore's

HFile
seqID
31634

**cloudera**
Ask Bigger Questions

# HBASE-10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 | 31402, 31403................ |

Flushed

When loading edits from a failed RS's log, replay only the edits coming after the HFile's highest sequence ID.

The bulk loaded file is assigned a sequence ID that ends up somewhere in the MemStore's

HFile seqID 31634

cloudera®
Ask Bigger Questions

# HBASE–10958 AKA Blindspot

Current MemStore

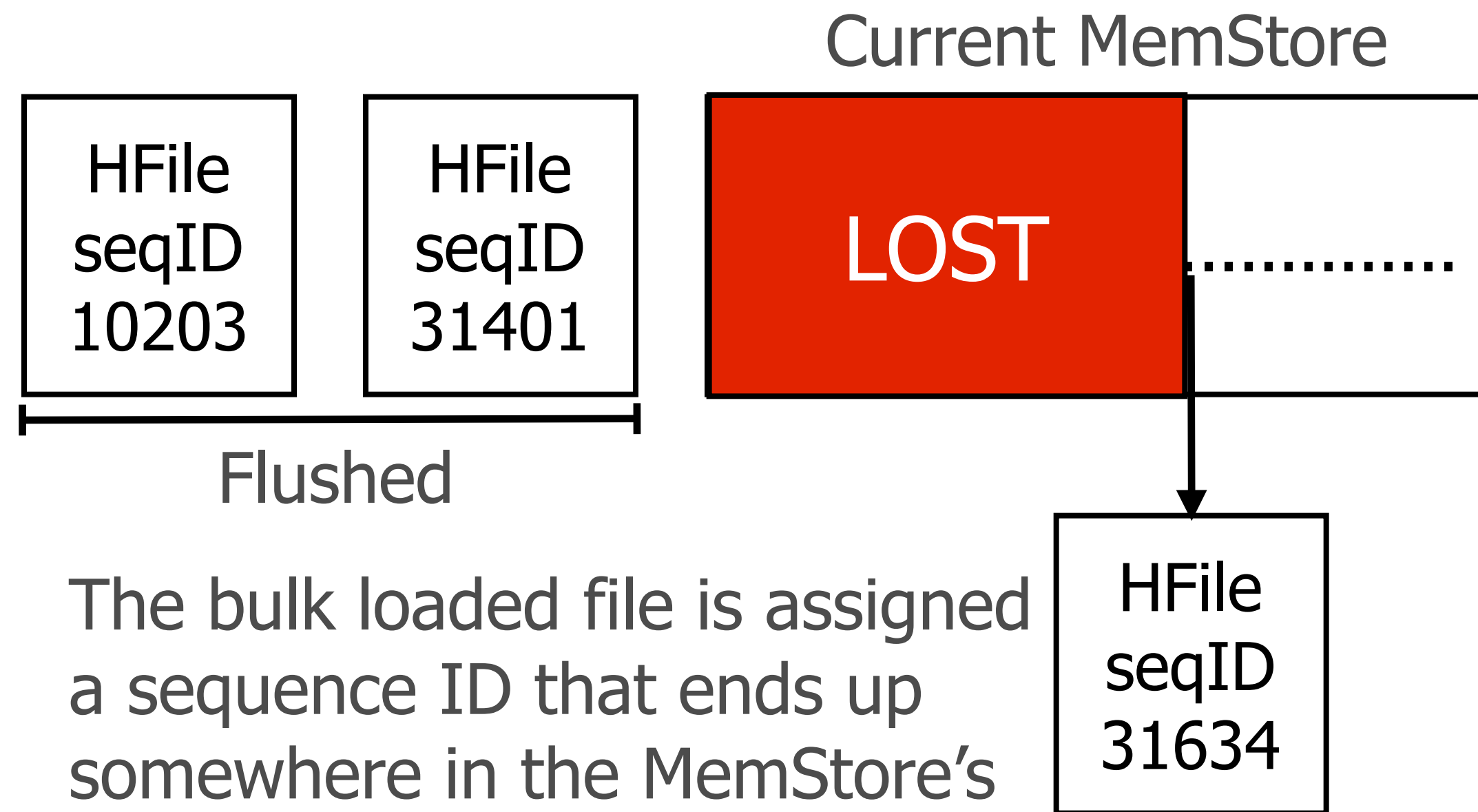| HFile<br>seqID<br>10203 | HFile<br>seqID<br>31401 | LOST | ............. |
|---|---|---|---|

Flushed

The bulk loaded file is assigned a sequence ID that ends up somewhere in the MemStore's

HFile
seqID
31634

When loading edits from a failed RS's log, replay only the edits coming after the HFile's highest sequence ID.

**cloudera**®
Ask Bigger Questions

# HBASE–10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 |
|---|---|

**Flushed**

LOST ··············

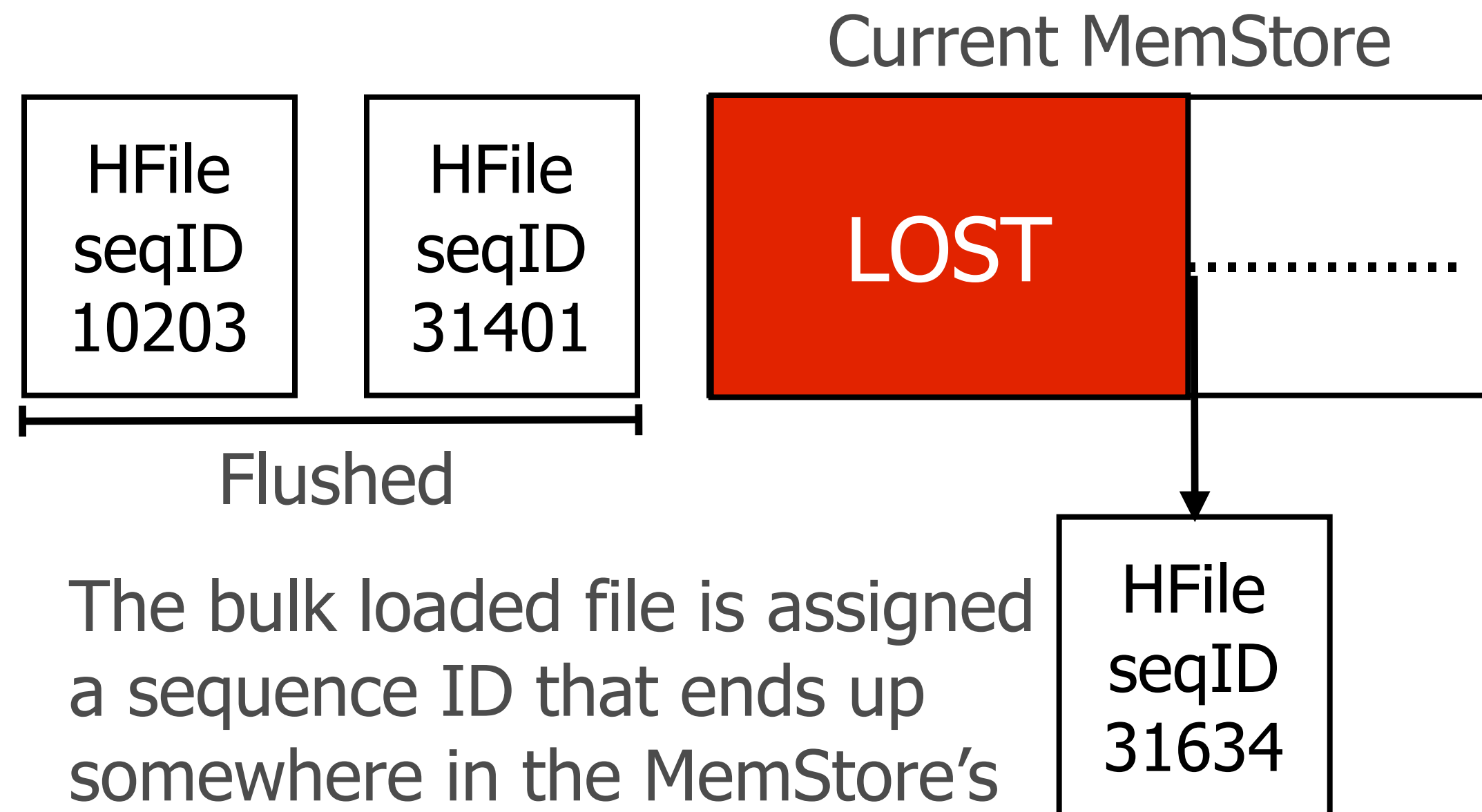The bulk loaded file is assigned a sequence ID that ends up somewhere in the MemStore's

HFile seqID 31634
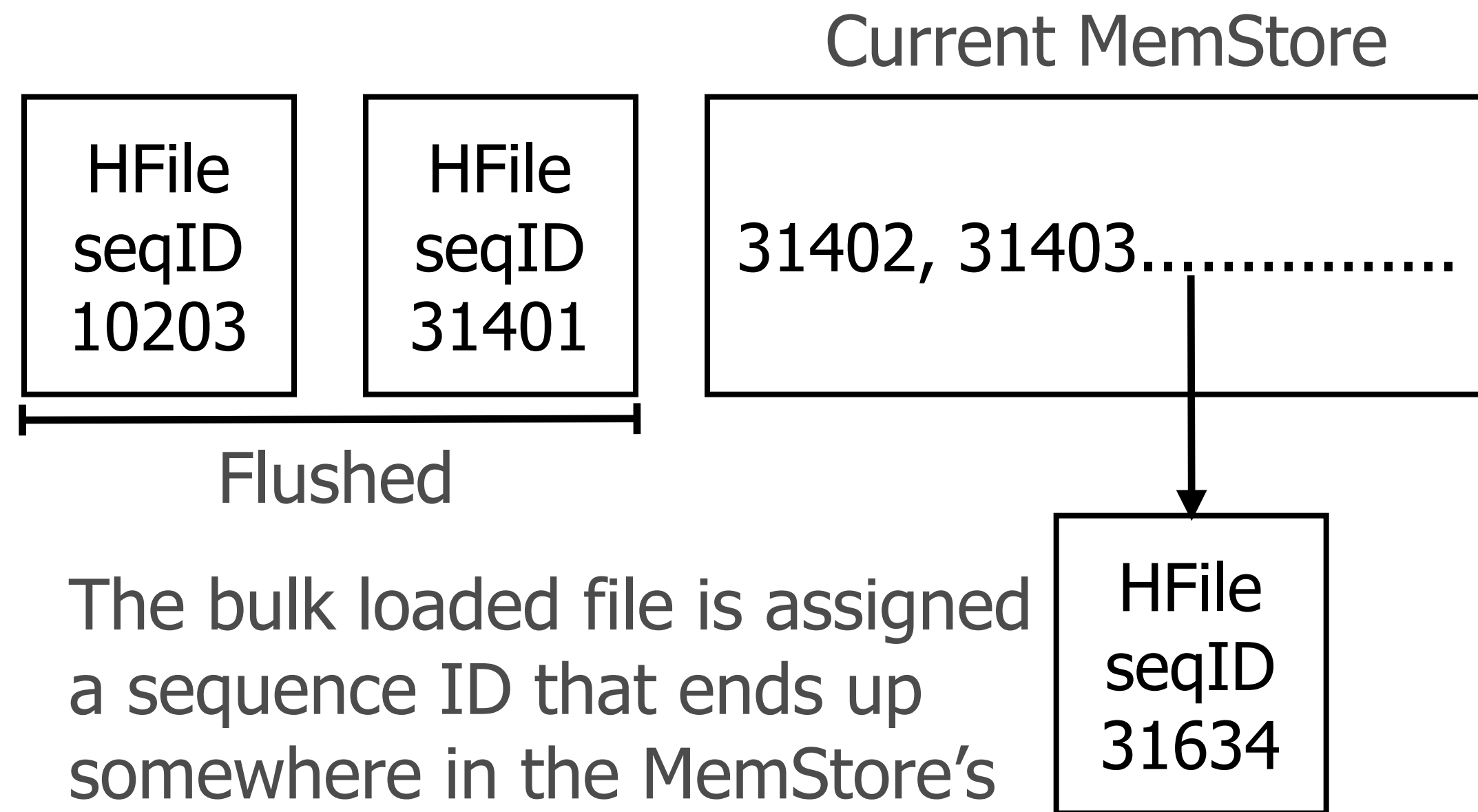
When loading edits from a failed RS's log, replay only the edits coming after the HFile's highest sequence ID.

Thankfully, we can recognize when a file is a bulk loaded one…

**cloudera**
Ask Bigger Questions

# HBASE–10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 | 31402, 31403................ |

Flushed

The bulk loaded file is assigned a sequence ID that ends up somewhere in the MemStore's

HFile seqID 31634

When loading edits from a failed RS's log, replay only the edits coming after the HFile's highest sequence ID.

Thankfully, we can recognize when a file is a bulk loaded one...
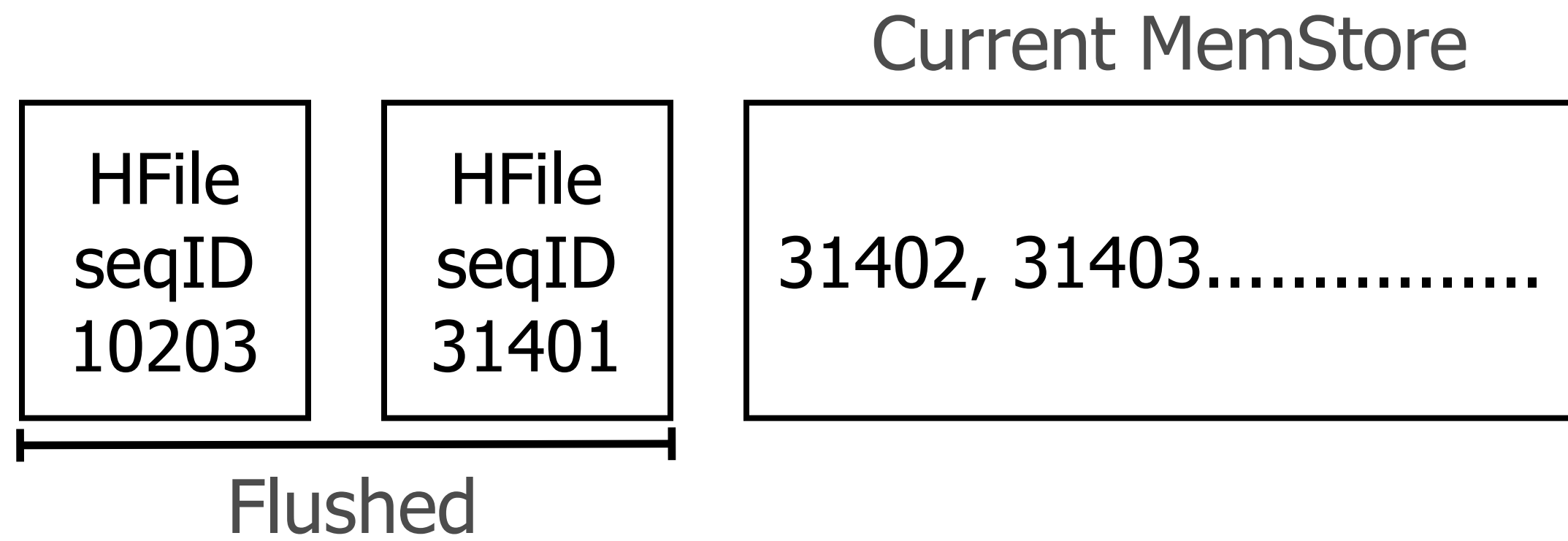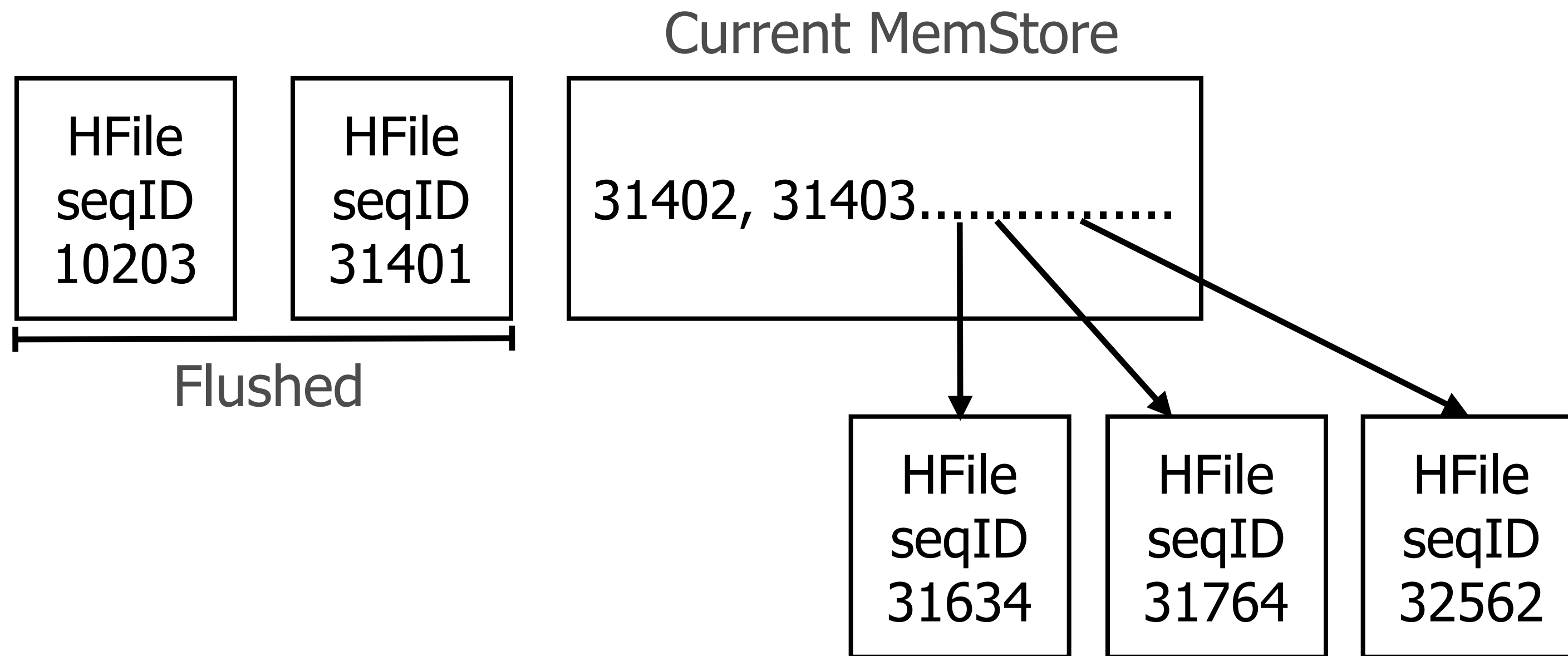
# HBASE–10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 | 31402, 31403............... |

Flushed

# HBASE–10958 AKA Blindspot

# HBASE–10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 | 31402, 31403................ |

Flushed

**cloudera**®
Ask Bigger Questions

# HBASE–10958 AKA Blindspot

Current MemStore

| HFile<br>seqID<br>10203 | HFile<br>seqID<br>31401 |

31402, 31403................

Flushed

The bulk loaded status is lost through compaction, the resulting HFile looks like any other!

HFile
seqID
32562

cloudera
Ask Bigger Questions

# HBASE–10958 AKA Blindspot

Current MemStore

| HFile seqID 10203 | HFile seqID 31401 | Blindspot ...... |

Flushed

HFile seqID 32562

The bulk loaded status is lost through compaction, the resulting HFile looks like any other!

The MemStore's data, starting from the beginning and up to 32562, will be lost even if it was logged.

**cloudera**
Ask Bigger Questions

# HBASE–10958 AKA Blindspot

Current MemStore
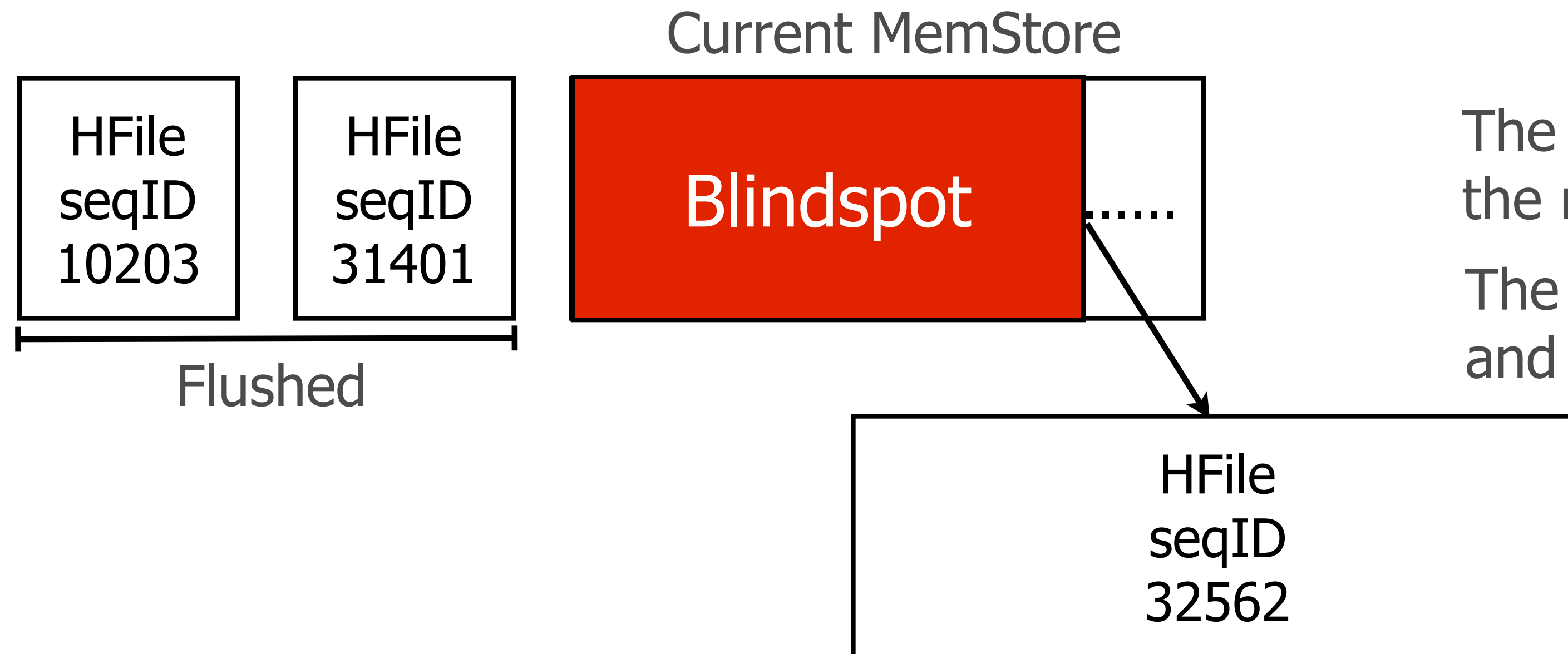
HFile
seqID
10203

HFile
seqID
31401

Blindspot ......

Flushed

HFile
seqID
32562

The bulk loaded status is lost through compaction, the resulting HFile looks like any other!

The MemStore's data, starting from the beginning and up to 32562, will be lost even if it was logged.

**The proposed solution is to force flush when bulk loading with sequence IDs, since the way it currently works goes against log replay's assumptions.**

cloudera®
Ask Bigger Questions

# In conclusion

- Loading via the "normal" APIs can be slow and/or disrupt the cluster.

# In conclusion

- Loading via the "normal" APIs can be slow and/or disrupt the cluster.
- Bulk loading can create files that HBase can directly use.

**cloudera**
Ask Bigger Questions

# In conclusion

- Loading via the "normal" APIs can be slow and/or disrupt the cluster.

- Bulk loading can create files that HBase can directly use.

- Useful for your original data import or incremental ones.

**cloudera**
Ask Bigger Questions

# In conclusion

- Loading via the "normal" APIs can be slow and/or disrupt the cluster.

- Bulk loading can create files that HBase can directly use.

- Useful for your original data import or incremental ones.

- Recommended to use HBase versions released during the past year.

cloudera®
Ask Bigger Questions

@jdcryans