```
$ hbase shell ↵

hbase(main):001:0> list ↵
    TABLE
        0 row(s)

hbase(main):002:0>  create 'tab2', 'cf'

> describe tab2 ↵

    cf is column family

> list 'tablename'

    show particular table information

> put 'tab2' 'row1' , 'cf:a' 100 ↵

> put 'tab2', 'row4'  'cf:b', 1000 ↵

> put 'tab2', 'row1', 'cf:b', 1000 ↵

> put 'tab2', 'row2', 'cf:a', 10 ↵

> put 'tab2', 'row2', 'cf:x' 1000 ↵

> scan 'tab2' ↵
```

| Row | Column_cell |
|-----|-------------|
| row1 | — |
| row1 | — |
| | — |
| row1 | — |
| row2 | — |
| row2 | — |

```
> Create 'hr', 'emp', 'dept', 'mngrinfo'
> describe 'hr'

> put 'hr', 'emp1', 'emp:ecode', 'a101'
> put 'hr', 'emp1', 'emp:esal', 2500
> put 'hr', 'emp1', 'dept:dname', 'it'
> put 'hr', 'emp1', 'dept:dloc', 'hyd'
> scan 'hr'

> get 'hr', 'emp2'
```

(RDBMS)

TABLE HAS PRIMARY KEY:

```
Sqoop import --connect jdbc:mysql://localhost/mydb --username
root --table emp --hbase-table hr1 --column-family emp
```

The primary key of RDBMS table is directed to row key of hbase table, Remaining columns of RDBMS table are directed to column family of hbase table. The above process will be failed if hbase table is not existed.

To create hbase base table dyanmically

```
Sqoop import --connect jd
     -
     -
     -
     --hbase-create-table
```

The option hbase-table will import data if table existed.
If table not existed it will create.

IF TABLE DOESN'T HAVE PRIMARY KEY :-

SQOOP import --connect Jdbc:mysql://localhost/mydb

--username _

--pwd _

--table Staff -m 1 --hbase-table hr3 --hbase-row-key ecode
--column-family emp --hbase-create-table

MYSQL: Tab

primarykey K   a   b   c   d

requirement:

a, b columns to be assigned to cf1 column family

c, d columns to be assigned to cf2 column family

Step1: $ Sqoop import --connect Jdbc:mysql://localhost/mydb

--username root

--pwd --

--table   Tab   --columns a, b, K, --hbase-table Test

--column-family cf1 --hbase-create-table

Step2: $ Sqoop import --connect Jdbc:mysql://localhost:/mydb

--username -root

--table tab --column K, C, d --hbase-table Test --column-
family cf2 --hbase-create-table

HBASE + hive integration :

hive> create table hbaseImage( ecode String, ename String, sex String,
sal int, dname String, dloc string)

> STORED BY 'org.apache.hadoop.hive.hbase.HBaseStorageHandler'

> WITH SERDEPROPERTIES(" hbase.columns.mapping" = ":key emp.name
emp:sex, emp:sal, dept:dname, dept:dloc")

> TBLPROPERTIES (" hbase.table.name"= " htab1");   ←  mandatory

## mapping:

| hive | hbase |
|------|-------|
| ecode | ——————— rowkey |
| name | ——————— emp:name |
| sex | ——————— emp:sex |
| sal | ——————— emp:sal |
| dname | ——————— dept:dname |
| dloc | ——————— dept:dloc |

one to one mapping is done.

STORED BY : SERDE class name : which serializes and deserializes hbase format to hive task and hive format to hbase table and the class name is ~~hbase~~ HBaseStorageHandler which is available in hbase-core.Jar.

SERDE PROPERTIES: Columns mapping will be done

* Columns of hive table and Hbase table mapping will be done
* mapping style is one to one mapping ( 1st column to 1st column, and so on )

Ex 1:

Image ( name, age, ecode, city )
                    ↳ key

they follow following
. hbase.Columns.mapping:
    = emp : name, emp:age, :key, emp:city

order is mandatory.

TBLPROPERTIES (optional):

→ When hive table is created, HBase table will also be created with given name

ex: TBLPROPERTIES (" hbase.table.name = " htab1").

if this option is not specified, with same name of hive table, hbase table will be created.

Integrating hive table with existed hbase table:

→ TBLPROPERTIES (" _ ", "—") → give the name of existed table.

Copying hdfs file data into HBase table:

* files delimiter is 'comma'

ex:-    101, AA, 2000        103, CC, 3000
        102, BB, 3000

1) hive> create table stage (ecode string, ename string, esal int)
   row format delimited fields terminated by ',';

2) hive> load data inpath 'file1.txt' into table stage;

3) hive> Insert overwrite table Image1
   Select ecode, ename, esal from stage;

to find total salary of hbase table

hive> Select sum(sal) from Image1;

   This query is processed by hbase.

Send Hbase table Aggregated results into DB2 table of target Reporting System.

HBASE DATA SCHEMA :

HBASE Table name : mytab
   data schema

| Row | cell | | |
| --- | --- | --- | --- |
| 101 | emp.name = Aminnudth | 103 | emp:sal = 16000 |
| 101 | emp.sal = 10000 | 103 | dept:dname=HR |
| 101 | dept:dname=IT | 103 | dept:dloc = pune |
| 101 | dept:dloc =hyd | | |
| 102 | emp.name = Avinash | | |
| 102 | emp.sal = 20000 | | |
| 102 | dept:dloc = hyd | | |
| 102 | dept:dname = HR | | |
| 103 | emp:name= Ankit | | |

Step1. hive> create table empimage (ecode string, name string, sal int, dname string, dloc string) STORED BY "org.apache.Hadoop.hive. hbase. HBaseStorageHandler" WITH SERDE PROPERTIES ("hbase.column. mapping" = " : key, emp:name, emp:sal, dept:dname, dept:dloc") TBLPROPERTIES ("hbase.table.name" = "mytab"); ↵

Step2: hive> create table Reshbase (dloc string, totalsal int);

Step3: hive> Insert overwrite table Reshbase

   Select dloc, sum(sal) from empimage group by dloc;

Step4:

   $ sqoop export -- connect jdbc:DB2://db.myorg.com/mydb

      -- username Devool  --password ney3123

      -- table LocAggr  --export-dir /usr/hive/warehouse/Resh base/

                                              00000 0

(Qn)

Send Hbase table Aggregated Results into DB2 table of target Reporting System.

Exporting HBase Table data (as it is) into RDBMS using SQOOP

   $ sqoop export -- connect jdbc:mysql://localhost/mydb
   -- username root
   -- table target
   -- hbase-table mytab
   -- column-family cf1 ↵

The row key of hbase table is mapped with primary key of table.
If RDBMS is not have primary key

     ex:- you want to send rowkey to ecode(not primary key)

ifca --hbase-table mytab --hbase-row-key ecode
                                                 ↓
                                   Column of RDBMS.

NOTE: while exporting -m 1 is not required.

limitations of hive and hbase integrations:-
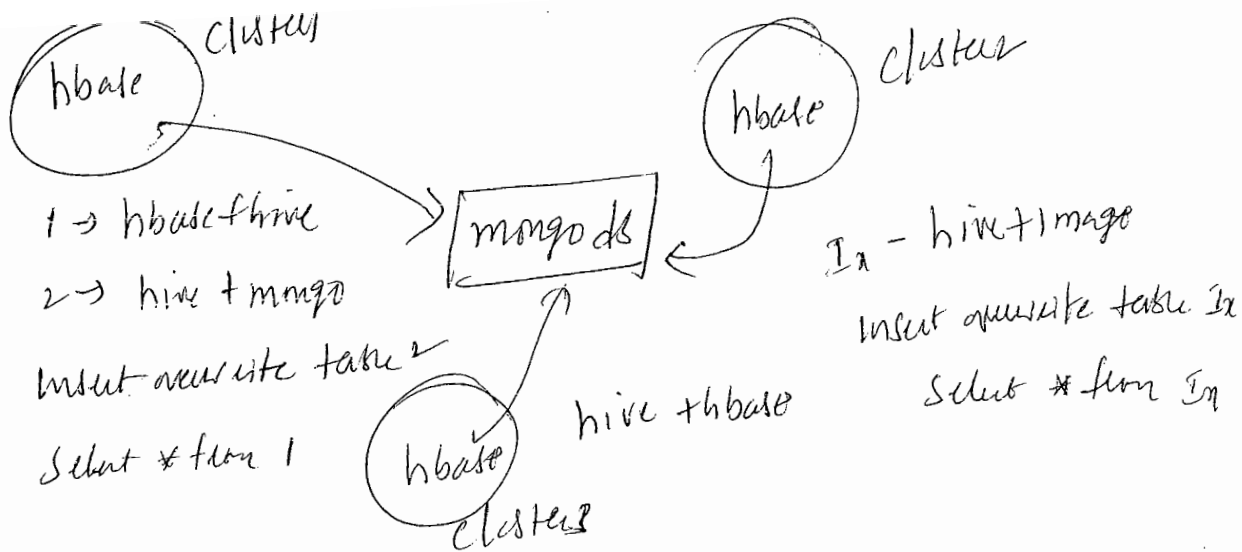
1) not possible for Remote hbase servers

2) We are not interact with Remote hbase (only with local hbase) solution for this is convert into HDFS they interact with Remote server. hive is able to integrate with Remote mongo DB.

ex: create table Imagemongo (-,-,-) STORED BY 'org.apache. hadoop. hive. mongo. mongoloader' with serde properties ("_"::"_") collection (" mongol") IPAddress (" 190: 390: 700: 804)

hbase — cluster

hbase — cluster

hbase — cluster3

mongo db

1 → hbase + hive

2 → hive + mongo

Insert overwrite table 2

Select * from 1

hive + hbase

$I_n$ – hive + Image

Insert overwrite table $I_n$

Select * from $I_n$

currently hive integration is available with

document
- 1) mongodb
- 2) couch db

(key, value)
- 3) Riac
- 4) PNUTS (not open source)

graphs
- 5) green plum
- 6) Neo4J

7) HBASE  8) Cassandra

both local and Remote integration

not Remote integration

# HADOOP CASCADING.

OR

In mapReduce, running job will be failed in following situations.

(1) if output path is already existed

(2) if any bad record is encountered.

## Solutions in CASCADING FOR ABOVE PROBLEMS:

cascading has SinkMode to control the file existence operations, Such as failed, overwrite, Reuse (appending)

SinkMode KEEP

SinkMode REPLACE.

SinkMode UPDATE

KEEP : Job will be failed if output file path is already existed that means it won't stop distrub data of file.

REPLACE : It overwrites into file if file (output path) is already existed.

UPDATE :- it appends new data to the existing file.

In cascading files (I/o files) are specified in the form of Taps. there are three types of Taps.

    1) Source Tap

    2) Sink Tap

    3) Trap Tap

Source Tap : To specify input file path

Sink Tap : to specify output file path

Trap Tap : a file to capture the bad data

## CASCADING TERMINOLOGY :

Tuple : tuple is a data record

Field : piece of data ex: line, word ...etc

Pipe : is an operation, to be performed on data

Tap : indicate a file (I/o)

Scheme : Schema design (field structure)

Input Scheme : input file structure (field list of input file)

Output Scheme : output file structure (field list of output file)

## Types of pipes (operations):

each : To perform functions/filters on each tuple, similar to foreach of pig

every : To perform Aggregation on Grouped data

GROUPBY : to group tuples based on fields

CoGroup : to Join multiple taps similar to SQL Joins such as inner/outer, left/white

Aggregator : to specify Aggregator function.

Pipe Assembly : Assembling All operations togeather in a Sequence.

PipeStream : the data which passes through a pipe is called pipe Stream

FlowConnector :- which connects Sourcetap, Sinketap, assembly

Flow : is a dataflow execution process, which Starts from head of the Assembly to tail of the Assembly.

Structure of cascading program :-

≡} Import of cascading classes.

public class MainClass
{
    public static void main(String args[]) {

        Source Scheme declaration;
        Sink Scheme declaration;
        Source tap;
        Sink tap;
        head of the Assembly

        ———————— operation(pipe1)        Flow instance;
        ——————— operation(pipe2)         execute flow;
        ≡              3                              }
        ≡              4                    }
                       5

    Flow Connector

Assign mainClass of Application Jar file

**Scheme** → to define Source & Sink Schemas.

Pkg : cascade.pipe.Scheme

ex:- Scheme    SourceScheme = new TextLine(new Fields("line"));

Scheme   SinkScheme = new TextLine(new Fields("word", "count"));

**Fields** — to specify field List of file

Pkg :- cascade.pipe.Field

**Tap** : to define input-file (Source tap) and output-file (Sink-tap)

Pkg : cascade.pipe.Tap

Eg :   Tap  Source = new Hfs(SourceScheme, inputpath);

Tap   Sink = new Hfs(SinkScheme, outputpath); SinkMode.REPLACE
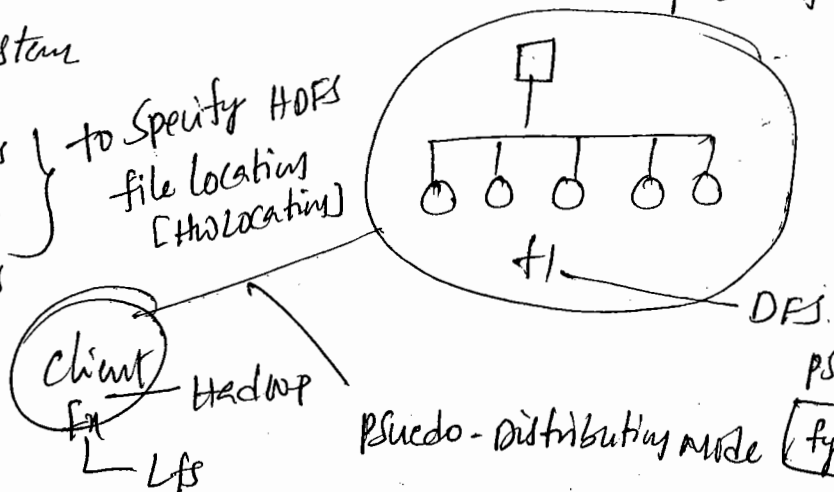
↓

default Sinkmode is

Sinkmode.KEEP

**Hfs** : Hadoop-file System
**Lfs** : Local-file System
**Dfs** : Distributed file System

Hadoop cluster in production

Pkgs :- cascade.pipe.Hfs  } to specify HDFS
        cascade.pipe.Lfs  }  file locations
        cascade.pipe.Dfs  }  [this location]



— DFS

Client — Hadoop
└ Lfs

Psuedo - Distribution mode Hfs

Psuedo Mode

HFS → Refers HDFS file of your psuedo-mode
LFS → Refers HDFS file of Your Client Machine
DFS → Refers HDFS file of Cluster (production)

Defining Assembly : (Header)

Pipe assembly = new Pipe ("wordcount");    ⌐→ Job name

pkg : cascade.pipe. pipe

Pipe Operations :

define a function

String regex = "[^ ]*";

Function f = new RegexGenerator(new Fields ("line"), regex).

(or)

Function f = new StringTokenizer(new Fields ("line")). nextToken().

Applying Functions For each Tuple :

assembly = new Each (new Fields (line), f);

(Grouping)

assembly = new GroupBy(assembly, new Fields ("word"));
                                 assembly, new Fields (line), f);

Defining an Aggregator :-

Aggregator count = new Aggregator (new Fields(wdd), Count);
                                    Count

Applying Aggregation on group :

assembly = new Every (assembly, count);

Flow Connector Instance :-

property p = new property ();

Flow Connector.setApplicationJarClass (Main.class);

... flow connector ... -> flow connector CFS)

## Defining Flow:

Flow f = fc. connect (assembly, source, sink)

( Jobname (wordcount), source Sink, assembly);

## Executing flow:

f. complete ();

## Example program:

import Java. util. regex. RegexGenerator;

import Java. io. IOException;

import cascade. pipe. pipe;

import cascade. pipe. Hfs;

import cascade. pipe. scheme;

import cascade. pipe. Tap;

import cascade. pipe. Fields;

import cascade. pipe. Each;

import    "    , .Function, GroupBy, Aggregator, FlowConnector, property,

Flow;

import og. apache. hadoop. util. GenericOptionparser

```
public class WordCount {

    public static void main (String str[]) {

        String[] files = new GenericOptionparser (args). getRemaining
                                                            Args();
Scheme SourceScheme = new TextLine (new Fields("date" "line") 
                                                 "offset" "line"));

Scheme SinkSource = new TextLine (new Fields ("word", "count"));

Tap Source = new Hfs (SourceScheme fields Cos);
```

```
Tap sink = new Hfs(SinkSource, files[1], SinkMode REPLACE);

Pipe assembly = new Pipe("n wdd-count"); header of assembly

    String regex = "[^ ]*";

Function func = new RegexGenerator(new Fields("line"), regex);

// Adding optional pipes to Assembly

assembly = new Each(assembly, new Fields("line"), func);

assembly = new GroupBy(assembly, new Fields("word"));

Aggregator count = new Count( new Fields("word"));
                         COUNT

assembly  = new Every(assembly, new Fields("count") count);

Properties p = new Properties();                  Tail of Assembly

FlowConnector .setApplicationJarClass (wddCount.class);

FlowConnector fc = new FlowConnector(p);

Flow f = fc. Count Connect("word-count" Source, Sink, assembly
                                              ↓                    );
                                         Source, Sink, assembly
f. complete();                           get together.
}
}
    executing Same as MR code
```