

## Hive

### Transactional and Analytical Processing

Transactional Processing	Analytical Processing
Analyzes individual entries	Analyzes large batches of data
Access to recent or immediate data, from the last few hours or days	Access to older data going back months, or even years
Updates data Involves many write operations	Only reads data
Updates data Fast real-time access	Long running jobs
Usually a single data source	Multiple data sources

### Small data

- ⌘ Both the objects could be achieved using the same database system.
- ⌘ Single machine with backup
- ⌘ Structured, well-defined data
- ⌘ Can access individual records or the entire dataset
- ⌘ No replication, updated data available instantaneously
- ⌘ Different tables store data from different sources


### Big data

- ⌘ Very hard to meet all requirements with the same database system
- ⌘ Data distributed on a cluster with multiple machines
- ⌘ Semi-structured or unstructured data
- ⌘ No random access to data
- ⌘ Data replicated, propagation of updates take time
- ⌘ Different sources may have different unknown formats

## Data warehouse

- ⌘ It typically has long running batch jobs that are used to crunch and much the data.
- ⌘ They are optimized for read operations because there is no real time processing or real time updates on the data.
- ⌘ The data is from multiple sources.
- ⌘ Data stored tends to go back many years in time.
- ⌘ Data may be lagged, not real-time.
- ⌘ Examples of data warehouses – Vertica, Teradata, Oracle, IBM

## Apache Hive

<ul style="list-style-type: none"><li>⌘ It is an open-source data warehouse.</li><li>⌘ Hive is part of the larger Hadoop ecosystem.</li><li>⌘ Hive runs on top of the Hadoop distributed computing framework.</li><li>⌘ It abstracts away the distributed nature of the storage and the processing to present a simple SQL-like interface.</li></ul>	
--	--

## Hive

- ⌘ Hive stores its data in HDFS
- ⌘ Hadoop Distributed File System
  - δ Data is stored as files - text files, binary files
  - δ Partitioned across machines in the cluster
  - δ Replicated for fault tolerance
  - δ Processing tasks parallelized across multiple machines
- ⌘ Hive runs all processes in the form of MapReduce jobs under the hood
- ⌘ MapReduce
  - δ A parallel programming model.
  - δ Defines the logic to process data on multiple machines
  - δ Batch processing operations on files in HDFS

- δ Usually written in Java using the Hadoop MapReduce library

## HIVE QL

- ⊞ Hive Query Language
  - δ A SQL-like interface to the underlying data.
  - δ Modeled on the Structured Query Language (SQL)
  - δ Familiar to analysts and engineers
  - δ Has simple query constructs
    - ⊞ select
    - ⊞ group by
    - ⊞ join
- ⊞ Hive exposes files in HDFS in the form of tables to the user
- ⊞ Write SQL-like query in HiveQL and submit it to Hive
- ⊞ Hive will translate the query to MapReduce tasks and run them on Hadoop
- ⊞ MapReduce will process files on HDFS and return results to Hive



## Hive metastore

- ⊞ A Hive user sees data as if they were stored in tables
- ⊞ Exposes the file-based storage of HDFS in the form of tables
- ⊞ It is the bridge between data stored in files and the tables exposed to users
- ⊞ Stores metadata for all the tables in Hive
- ⊞ Maps the files and directories in Hive to tables
- ⊞ Holds table definitions and the schema for each table
- ⊞ Has information on converting files to table representations
- ⊞ Any database with a JDBC driver can be used as a metastore
- ⊞ Development environments use the built-in Derby database i.e. Embedded metastore



**Hive vs. RDBMS**

HIVE	RDBMS
Large datasets	Small datasets
Parallel computations	Serial computations
High latency	Low latency
Read operations	Read/write operations
Not ACID compliant by default	ACID compliant
HiveQL	SQL

HIVE	RDBMS
<ul style="list-style-type: none"> <li>⊞ Large datasets               <ul style="list-style-type: none"> <li>δ Gigabytes or petabytes</li> <li>δ Calculating Trends</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>⊞ Small datasets               <ul style="list-style-type: none"> <li>δ Megabytes or gigabytes</li> <li>δ Accessing and updating individual records</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>⊞ Parallel Computations               <ul style="list-style-type: none"> <li>δ Distributed system with multiple machines</li> <li>δ Semi-structured data files partitioned across machines</li> <li>δ Disk space cheap, can add space by adding machines</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>⊞ Serial Computations               <ul style="list-style-type: none"> <li>δ Single computer with backup</li> <li>δ Structured data in tables on one machine</li> <li>δ Disk space expensive on a single machine</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>⊞ High Latency               <ul style="list-style-type: none"> <li>δ Records not indexed, cannot be accessed quickly</li> <li>δ Fetching a row will run a MapReduce that might take minutes</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>⊞ Low Latency               <ul style="list-style-type: none"> <li>δ Records indexed, can be accessed and updated fast</li> <li>δ Queries can be answered in milliseconds or microseconds</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>⊞ Read Operations</li> </ul>	<ul style="list-style-type: none"> <li>⊞ Read Operations</li> </ul>

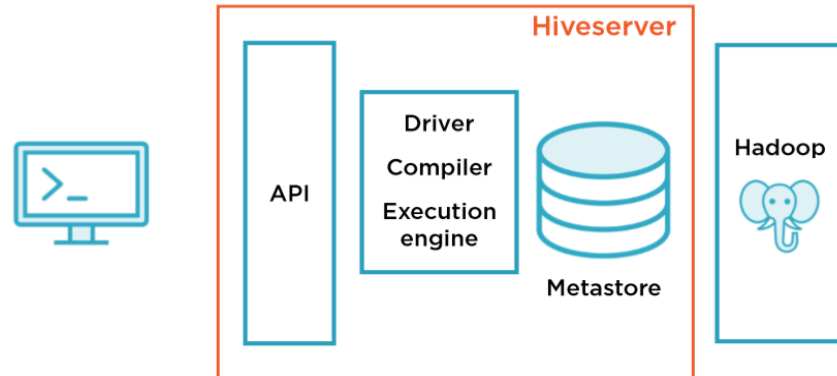
<ul style="list-style-type: none"> <li>δ Not the owner of data. <ul style="list-style-type: none"> <li>▢ Hive stores files in HDFS</li> <li>▢ Hive files can be read and written by many technologies <ul style="list-style-type: none"> <li>⇒ Hadoop, Pig, Spark</li> </ul> </li> <li>▢ Hive database schema cannot be enforced on these files</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>δ Schema-on-read <ul style="list-style-type: none"> <li>▢ Number of columns, column types, constraints specified at table creation</li> <li>▢ Hive tries to impose this schema when data is read</li> <li>▢ It may not succeed, may pad data with nulls</li> </ul> </li> <li>⊞ Read/ Write Operations <ul style="list-style-type: none"> <li>δ Sole gatekeeper for data Schema on write</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>⊞ Not ACID Compliant <ul style="list-style-type: none"> <li>δ Data can be dumped into Hive tables from any source</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>⊞ ACID Compliant <ul style="list-style-type: none"> <li>δ Only data which satisfies constraints are stored in the database</li> </ul> </li> </ul>

HIVEQL	SQL
Schema on read, no constraints enforced	Schema on write keys, not null, unique all enforced
Minimal index support	Indexes allowed
Row level updates, deletes as a special case	Row level operations allowed in general
Many more built-in functions	Basic built-in functions
Only equi-joins allowed	No restriction on joins
Restricted subqueries	Whole range of subqueries

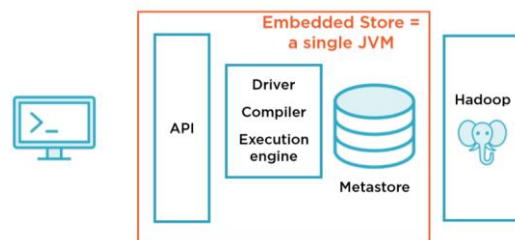
### Basic HIVE BUILDING BLOCKS

- ⊞ There is a terminal window on which the queries are written that look exactly like MySQL.
- ⊞ This talks to the Hiveserver.

- ⌘ The Hiveserver consults the metastore for the Hive table information.
- ⌘ It passes these queries, converts them to MapReduce jobs and submits them to the underlying Hadoop ecosystem where they are run and the results are returned.



- ⌘ Within Hive, the API is the gateway to the external world.
  - δ This allows various clients to connect to Hive and run queries.
- ⌘ The combination of the driver, compiler and execution engine is the heart and soul of Hive.
  - δ It accepts queries through its API, parses it into the right form using metastore information, converts it to individual jobs or tasks that can be executed on Hadoop.
- ⌘ The metastore is also part of Hive and it contains all table metadata information, the table schema, the column data types, how records from Hive can be serialized and deserialized onto HDFS.
- ⌘ All together these components form the Hiveserver.
- ⌘ The Hiveserver passes these MapReduce jobs onto Hadoop in order to be run on a cluster of machines.
- ⌘ Connecting to the API component of the Hiveserver is done by the command line interface.
- ⌘ Its called an embedded store when it is run on a single machine.
- ⌘ The latest versions of Hive run a rewritten server called Hiveserver2
- ⌘ Hiveserver2 supports concurrent clients and better authentication and authorization
- ⌘ Beeline is a command line interface which works with this new server



**Commands**

Command	Explanation	Output
show databases;	Shows the databases existing within Hive.	hive> show databases; OK default 1 row selected (3.61 seconds)
create database <database_name>;	Create a database and store it with the name given.	hive> create database MCA; OK No rows affected (1.483 seconds)  hive> show databases; OK default mca 2 rows selected (0.129 seconds)
use <database_name>;	Use the database	hive> use MCA; OK No rows affected (0.065 seconds)
create table <table_name> (<fieldname> <datatype>, .....);	Create a table	hive> create table customers ( . . > id bigint, . . > name string, . . > address string . . > ); OK No rows affected (4.583 seconds)
create table if not exists <tablename> (<fieldname> <datatype>, .....);	To create a new table only if it does not exist.	hive> create table if not exists orders ( . . > id bigint, . . > product_id string, . . > cust_id bigint, . . > quantity int, . . > amt double . . > ); OK No rows affected (3.022 seconds)
show tables;	Show all tables in the current database.	hive> show tables; OK customers 1 row selected (0.757 seconds)

describe <tablename>;	Shows the schema the table has.	<pre>hive&gt; describe customers; OK id bigint name string address string 3 rows selected (0.703 seconds)</pre>
insert into table <tablename> values (<value1>, <value2>, .....);	Insert data into tables. Rarely used since you load data into tables from some file or some directory in HDFS.	
<pre>hive&gt; insert into table customers values (1, "aaa","bbb");</pre>		
To insert many values at the same time		
<pre>hive&gt; insert into customers values (112, "AAA","ZZZ"), (113,"BBB","YYY"),(114,"CCC","XXX"),(178, "DDD","WWW");</pre>		
select * from <tablename>;	To see all the values	<pre>hive&gt; select * from customers; OK 1 aaa bbb 112 AAA ZZZ 113 BBB YYY 114 CCC XXX 178 DDD WWW 5 rows selected (0.821 seconds)</pre>
select <columnname> from <tablename>	To see a particular column value	<pre>hive&gt; select name from customers; OK aaa AAA BBB CCC DDD 5 rows selected (0.828 seconds)</pre>
select <columnname> from <tablename> where <condition>	To show a particular column's value based on a condition.	<pre>hive&gt; select name from customers where address = "ZZZ"; OK AAA 1 row selected (1.303 seconds)  hive&gt; select * from customers where address = "ZZZ"; OK 112 AAA ZZZ 1 row selected (0.779 seconds)  hive&gt; select name from customers where address = "ZZZ" and id &gt;100; OK AAA DDD EEE 3 rows selected (3.007 seconds)</pre>
select distinct <columnname> from <tablename>;	To show unique values of a particular	<pre>hive&gt; select DISTINCT address from customers;</pre>



	column.	<pre>WWW XXX YYY ZZZ bbb 5 rows selected (73.549 seconds)</pre>
<pre>select &lt;columnname&gt; from &lt;tablename&gt; order by &lt;columnname&gt;;</pre>	To order the result based on a particular column name.	<pre>hive&gt; select name, address from customers order by address; OK DDD WWW CCC XXX BBB YYY EEE ZZZ DDD ZZZ AAA ZZZ aaa bbb 7 rows selected (60.537 seconds)</pre>
<pre>select count(*) from &lt;tablename&gt;;</pre>	To find the total number of records in the table.	<pre>hive&gt; select count(*) from customers; OK 7 1 row selected (57.382 seconds)</pre>
<pre>select sum(&lt;columnname&gt;) from &lt;tablename&gt;;</pre>	To sum a column that contains numbers.	<pre>hive&gt; select sum(id) from customers; OK 750 1 row selected (52.21 seconds)</pre>
<pre>select &lt;aggregate(&lt;columnname&gt;)&gt; from &lt;tablename&gt;;</pre>	To apply any aggregate function on the table.	<pre>hive&gt; select max(id) from customers; OK 178 1 row selected (55.873 seconds)</pre>
<pre>select * from &lt;tablename&gt; limit &lt;n&gt;;</pre>	To display first 'n' rows of a table.	<pre>hive&gt; select * from customers limit 2; OK 1 aaa bbb 112 AAA ZZZ 2 rows selected (0.764 seconds)</pre>
<pre>select &lt;columnname&gt;, &lt;aggregate(&lt;columnname&gt;)&gt; from &lt;tablename&gt; group by &lt;columnname&gt;;</pre> <p>To group the records of a table based on a particular column.</p>		<pre>hive&gt; select address, count(*) from customers group by address; OK WWW 1 XXX 1 YYY 1 ZZZ 3 bbb 1 5 rows selected (45.107 seconds)</pre>

All select commands are MapReduce jobs except select \* since it is a simple read operation on HDFS.

<pre>select * from &lt;table1name&gt; join &lt;table2name&gt; where &lt;table1name&gt;.&lt;columnname&gt; = &lt;table2name&gt;.&lt;columnname&gt; ;</pre>	<p>Join two tables based on equijoin.</p>	<pre>hive&gt; select customers.id, name, product_id from customers join orders ... &gt; where customers.id = orders.cust_id; OK 112 AAA TV 178 DDD phone 2 rows selected (49.577 seconds)</pre>
---	---	---

Hive only supports equi-joins where the 'where' clause on the joined column is an equality clause.

It does not support natural joins.

## HIVE DATA TYPES

### ⌘ Primitive Data types

- ⊘ Boolean
- ⊘ Numeric
- ⊘ String
- ⊘ Timestamp

Boolean		Numeric	
⌘ True or false	⌘ Yes /no questions	⌘ Integers <ul style="list-style-type: none"> <li>⊘ Tinyint: 1 byte, range -128 to 128</li> <li>⊘ Smallint: 2 bytes, range -2<sup>15</sup> to 2<sup>15</sup> - 1</li> <li>⊘ Int: 4 bytes, range -2<sup>31</sup> to 2<sup>31</sup> - 1</li> <li>⊘ Bigint: 8 bytes, range -2<sup>63</sup> to 2<sup>63</sup> - 1</li> </ul>	⌘ Decimals <ul style="list-style-type: none"> <li>⊘ Float: 4 bytes</li> <li>⊘ Double: 8 bytes</li> <li>⊘ Decimal: Arbitrary precision - dec(10, 2) - will occupy 10 digits of space and after the decimal point two spaces</li> </ul>
String		Timestamp	
⌘ String: Unbounded, variable length character string	⊘ Char: Fixed length character string ⊘ Varchar: Bounded,	⌘ Integers: Unix timestamp in nanoseconds ⌘ Floats: Unix timestamp in seconds with decimal precision ⌘ String: JDBC compliant "YYYYMM-DD HH:MM:SS.fffffffff"	

variable length character string	<ul style="list-style-type: none"> <li>⊞ Dates</li> <li>Values of the form "YYYY-MM-DD"</li> </ul>
----------------------------------	--

## Storing Data in HIVE

### ⊞ Data

- δ The records in the table which holds the actual data
- δ Stored in HDFS, the reliable storage for data in Hadoop
- δ Files partitioned across multiple machines in the cluster
- δ Stored in directories under Hive's warehouse directory
- δ hive.metastore.warehouse.dir property in hive-site.xml
- δ Defaults to /user/hive/warehouse

```
C:\hive\apache-hive-2.1.0-bin>hadoop fs -ls /user/hive/warehouse/mca.db
Found 2 items
drwxr-xr-x - raoal supergroup      0 2019-11-13 18:08 /user/hive/warehouse/mca.db/customers
drwxr-xr-x - raoal supergroup      0 2019-11-13 23:04 /user/hive/warehouse/mca.db/orders
```

### ⊞ Metadata

- δ Information about the underlying data in the table
- δ Metastore, acts as a bridge between Hive and files in HDFS
- δ A relational database with information on
  - ▢ databases, tables
  - ▢ columns, owners, storage, serialization/ deserialization information
  - ▢ user supplied metadata

## HIVE TABLES

### ⊞ Managed

- δ Data managed by Hive and stored in the warehouse directory

### ⊞ External

- δ Data not fully managed by Hive and exists outside the warehouse directory

### ⊞ The metadata for both is in the metastore

## Managed Tables

- ⊞ All tables so far have been managed tables
- ⊞ Hive owns the files and directories
- ⊞ These can be modified by other technologies
- ⊞ Deleting a managed table deletes both data and metadata

## External tables

- Share the underlying data across other technologies
- Hadoop, Pig, HBase all of these may access and edit those files
- Deleting an external table deletes only the metadata

## Demo

Use products.csv	It contains id, name and cost of a product.	<table border="1"> <tr> <td>iphone7</td><td>iPhone 7</td><td>950</td></tr> <tr> <td>camera_canon</td><td>Canon 570x</td><td>1000</td></tr> <tr> <td>washingmachine_samsung</td><td>Samsung Swift</td><td>400</td></tr> <tr> <td>tv_vu</td><td>Vu 56 Inch</td><td>600</td></tr> </table>	iphone7	iPhone 7	950	camera_canon	Canon 570x	1000	washingmachine_samsung	Samsung Swift	400	tv_vu	Vu 56 Inch	600
iphone7	iPhone 7	950												
camera_canon	Canon 570x	1000												
washingmachine_samsung	Samsung Swift	400												
tv_vu	Vu 56 Inch	600												
hadoop fs -mkdir /<directoryname>/	Create a directory in Hadoop to store this products.csv	<pre>C:\WINDOWS\system32&gt;hadoop fs -mkdir /data</pre>												
hadoop fs -copyFromLocal <filename> /<directoryname>/	Copy the file products.csv to the hdfs directory created	<pre>C:\hive\apache-hive-2.1.0-bin&gt;hadoop fs -copyFromLocal products.csv /data/ C:\hive\apache-hive-2.1.0-bin&gt;hadoop fs -ls /data Found 1 items -rw-r--r-- 1 raoul supergroup 123 2019-11-15 12:49 /data/products.csv</pre>												
create external table if not exists <tablename> ( <fieldname> <datatype>,..... ) location '/<directoryname>/'	Create the external table	<pre>hive&gt; create external table if not exists products ( . . &gt; id string, . . &gt; title string, . . &gt; cost float . . &gt; ) . . &gt; comment "Table to store information sold in stores" . . &gt; location '/data/'; OK No rows affected (71.257 seconds)  hive&gt; describe products; OK id string title string cost float 3 rows selected (7.26 seconds)</pre>												

		<pre>hive&gt; show tables; OK customers orders products values__tmp__table__2 4 rows selected (2.267 seconds)</pre>	
<ul style="list-style-type: none"> <li>⊞ The external keyword indicates that the data for this table is not managed by Hive.</li> <li>    δ The data for this table exists in an HDFS directory specified by the location keyword.</li> <li>⊞ Make sure that the order in which the columns are specified in the table matches the order in which the data is available in the csv file.</li> <li>⊞ Even though the data has not been inserted when a select command is used all the data that was available in the csv file is available.</li> </ul>			
<p>Hive is schema-on-read.</p> <p>For the schema that has been defined for the table, hive will try to impose the schema on the underlying data which is referenced in HDFS.</p> <p>All the individual records have been squished into one column.</p>		<pre>hive&gt; select * from products; OK iphone7, iPhone 7, 950 camera_canon, Canon 570x, 1000 washingmachine_samsung, Samsung Swift, 400 tv_vu, Vu 56 Inch, 600 6 rows selected (8.542 seconds)</pre>	
hadoop	fs	-ls	<p>Check the contents of the warehouse by using the command</p> <pre>C:\hive\apache-hive-2.1.0-bin&gt;hadoop fs -ls /user/hive/warehouse/mca.db Found 2 items drwxr-xr-x - ra0al supergroup 0 2019-11-13 18:08 /user/hive/warehouse/mca.db/customers drwxr-xr-x - ra0al supergroup 0 2019-11-13 23:04 /user/hive/warehouse/mca.db/orders</pre>
<p>It shows that a directory has been created for the managed tables of customers and orders but not for products.</p> <p>The product table does not live within the Hive warehouse since it is an external table.</p>			
<p>The show tables show the table.</p> <p>The metadata for the table is present in the metastore and the actual data is not in the Hive but in the /data HDFS location</p>			

drop table products;	Drop the table	<pre>hive&gt; drop table products; OK No rows affected (10.403 seconds) hive&gt; show tables; OK customers orders 2 rows selected (0.376 seconds)</pre>
<p>The products.csv is still present in the Hadoop folder.</p> <p>Dropping the table has not deleted the underlying data.</p> <p>This is one of the main differences between managed and external tables in hive.</p>		
row format delimited fields terminated by '<symbol>'	To see that the data is correctly stored in different columns use the command	<pre>hive&gt; create external table if not exists products ( . . &gt; id string, . . &gt; title string, . . &gt; cost float . . &gt; ) . . &gt; comment "Table to store product information stored in stores" . . &gt; row format delimited fields terminated by ',' . . &gt; stored as textfile . . &gt; location '/data/' OK No rows affected (1.017 seconds)</pre>
<p>The format of every row is told to Hive.</p> <p>The row format is delimited fields where each field is terminated by a comma, since it is a comma separated file.</p> <p>Hive considers the default delimiter to be the ASCII value \001.</p>		
stored as <option>	Specify the kind of file being read.	
<p>The default is textfile so the stored as textfile is optional.</p>		
<p><a href="https://cwiki.apache.org/confluence/display/Hive/LanguageManual+DDL#LanguageManualDDL-CreateTable">https://cwiki.apache.org/confluence/display/Hive/LanguageManual+DDL#LanguageManualDDL-CreateTable</a></p>		
select * from <tablename>;	Check the contents	<pre>hive&gt; select * from products; OK iphone7      iPhone 7  950.0 camera_canon Canon 570x 1000.0 washingmachine_samsung Samsung Swift 400.0 tv_vu        Vu 56 Inch 600.0  6 rows selected (1.372 seconds)</pre>

<p>describe &lt;tablename&gt;;</p> <p>formatted</p>	<p>To check whether a table is managed or external use</p>	<pre>hive&gt; describe formatted customers; OK # col_name          data_type          comment id bigint name string address string  # Detailed Table Information Database:           mca Owner:              raoal CreateTime:         Tue Nov 12 09:50:45 IST 2019 LastAccessTime:     UNKNOWN Retention:          0 Location:           hdfs://0.0.0.0:19000/user/hive/warehouse/mca.db/customers Table Type:         MANAGED_TABLE Table Parameters: COLUMN_STATS_ACCURATE {\\"BASIC_STATS\\":\\"true\\"} numFiles            3 numRows             7 rawDataSize         75 totalSize           82 transient_lastDdlTime 1573648726  hive&gt; describe formatted products; OK # col_name          data_type          comment id string title string cost float  # Detailed Table Information Database:           mca Owner:              raoal CreateTime:         Fri Nov 15 20:50:50 IST 2019 LastAccessTime:     UNKNOWN Retention:          0 Location:           hdfs://0.0.0.0:19000/data Table Type:         EXTERNAL_TABLE Table Parameters: EXTERNAL            TRUE comment             Table to store product information stored in stores numFiles            1 totalSize           123 transient_lastDdlTime 1573831250  # Storage Information SerDe Library:      org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe InputFormat:        org.apache.hadoop.mapred.TextInputFormat OutputFormat:       org.apache.hadoop.hive ql.io.HiveIgnoreKeyTextOutputFormat</pre>
<p>like &lt;tablename&gt;</p>	<p>can be used with create command to create a table that is exactly like another table that exists</p>	<pre>hive&gt; create table if not exists fresh_products like products; OK No rows affected (1.591 seconds)  hive&gt; describe fresh_products; OK id string title string cost float 3 rows selected (1.326 seconds)  hive&gt; select * from fresh_products; OK No rows selected (2.001 seconds)</pre>

alter table <tablename> rename to <newname>;	To rename a table use the rename command with alter command.	<pre>hive&gt; alter table fresh_products . . &gt; rename to freshprods; OK No rows affected (3.6 seconds) hive&gt; show tables; OK customers freshprods orders products 4 rows selected (0.358 seconds)</pre>
alter table <tablename> add columns ( <fieldname> <datatype> );	To add columns use the add columns in the alter command.	<pre>hive&gt; alter table freshprods add columns ( . . &gt; expiry_date date . . &gt; comment "Expiry date of fresh produce" . . &gt; ); OK No rows affected (0.565 seconds)  hive&gt; describe freshprods; OK id string title string cost float expiry_date date Expiry date of fresh produce 4 rows selected (0.214 seconds)</pre>
alter table <tablename> change column <oldcolumnname> <old/newcolumnname> <datatype> after <columnname>	To move one column's position use change with after. change column. columns.	<pre>hive&gt; alter table freshprods . . &gt; change column id id string . . &gt; after title; OK No rows affected (0.496 seconds)  hive&gt; describe freshprods; OK title string id string cost float expiry_date date Expiry date of fresh produce 4 rows selected (0.152 seconds)</pre>



Moving columns does not move the underlying data in the

```
hive> alter table products_
. . > change column id id string
. . > after title;
OK
No rows affected (0.184 seconds)
```

```
hive> describe products;
OK
title string
id string
cost float
3 rows selected (0.278 seconds)
```

```
hive> select * from products;
OK
iphone7    iPhone 7  950.0
camera_canon  Canon 570x 1000.0
washingmachine_samsung Samsung Swift 400.0
tv_vu       Vu 56 Inch 600.0
6 rows selected (1.134 seconds)
```

### Temporary tables

- ⌘ Tables created within a Hive session, deleted when the session ends
- ⌘ They are active within the Hive session and as soon as that session ends, the table and all its contents are deleted.
- ⌘ Store temporary data
- ⌘ Tables of the same name can be created by different users
- ⌘ Do not support partitions and indexes
- ⌘ Can have the same name as a permanent table.
  - δ As long as temporary table exists the permanent table cannot be accessed.
  - δ The temporary table hides the permanent table if it has the same name.

**Demo**

Check the tables that exist in the database		<pre>hive&gt; show tables; OK customers freshprods orders products 4 rows selected (14.34 seconds)</pre>
create temporary table <tablename>	Create temporary tables using the key word temporary	<pre>hive&gt; create temporary table test_customers like customers; OK No rows affected (7.389 seconds)  hive&gt; describe test_customers; OK id bigint name string address string 3 rows selected (1.016 seconds)  hive&gt; show tables; OK customers freshprods orders products test_customers 5 rows selected (1.375 seconds)</pre>
insert into <tablename> values <v1, v2, v3>	Insert into the temporary table	<pre>hive&gt; insert into test_customers values (189, "ACAC", "VVV"); WARNINGS: [HIVE-1030107] Hive 0.12.0 does not support INSERT INTO TEMPORARY TABLE</pre>
		<pre>hive&gt; select * from test_customers; OK 189 ACAC VV 1 row selected (2.984 seconds)</pre>
quit; hive	Quit the session and restart hive again.	<pre>hive&gt; show tables; OK customers freshprods orders products 4 rows selected (0.053 seconds)</pre>
The temporary table no longer exists.		
	Create another table name with	<pre>hive&gt; create temporary table customers like orders; OK No rows affected (1.238 seconds)</pre>

	same name as permanent table.	
	The original customers table has been hidden.	<pre>hive&gt; show tables; OK customers freshprods orders products 4 rows selected (0.168 seconds)</pre> <pre>hive&gt; describe customers; OK id bigint product_id string cust_id bigint quantity int amt double 5 rows selected (0.045 seconds)</pre>
	Quit hive and restart it. Use the mca database and see the tables.	<pre>hive&gt; show tables; OK customers freshprods orders products 4 rows selected (1.823 seconds)</pre> <pre>hive&gt; describe customers; OK id bigint name string address string 3 rows selected (1.058 seconds)</pre>

### Inserting data

- ⊞ Standalone
  - δ Insert individual records
  - δ Files
  - δ Other tables

### Demo

- ⊞ Insert using files

	Create a freshproducts.csv and store it in bin directory of hive.	<table border="1"> <tr><td>broccolli</td><td>Broccoli</td><td>5</td></tr> <tr><td>spinach</td><td>Spinach</td><td>7</td></tr> <tr><td>carrot</td><td>Local Carr</td><td>4</td></tr> <tr><td>potato</td><td>Idaho Pot</td><td>4</td></tr> <tr><td>milk</td><td>Skimmed M</td><td>5</td></tr> <tr><td>cheese</td><td>Spicy Garli</td><td>5</td></tr> </table>	broccolli	Broccoli	5	spinach	Spinach	7	carrot	Local Carr	4	potato	Idaho Pot	4	milk	Skimmed M	5	cheese	Spicy Garli	5
broccolli	Broccoli	5																		
spinach	Spinach	7																		
carrot	Local Carr	4																		
potato	Idaho Pot	4																		
milk	Skimmed M	5																		
cheese	Spicy Garli	5																		
alter table <tablename> change column <oldcolumnname> <old/newcolumnname> <datatype> after <columnname>;	Alter the table freshprods	<pre>hive&gt; alter table freshprods change column . . &gt; title title string after id; OK No rows affected (1.58 seconds)</pre>																		
describe <tablename>;	Check the schema of freshprods	<pre>hive&gt; describe freshprods; OK id string title string cost float expiry_date date Expiry date of fresh produce 4 rows selected (0.184 seconds)</pre>																		
load data local inpath <'filename.csv'> into table <tablename>	Load the data from freshproducts.csv into fresprods.	<pre>hive&gt; load data local inpath 'freshproducts.csv' . . &gt; into table freshprods; Loading data to table mca.freshprods OK No rows affected (3.898 seconds)</pre>																		
select * from <tablename>	Check the contents of the table.	<pre>hive&gt; select * from freshprods; OK broccolli Broccoli 5.0 spinach Spinach 7.0 carrot Local Carrots 4.0 potato Idaho Potatoes 4.0 milk Skimmed Milk 5.0 cheese Spicy Garlic 5.0 6 rows selected (4.297 seconds)</pre>																		
<p>Using this command a copy is made into the Hive warehouse directory if this is a managed table. The contents of the csv file has been copied into Hive's warehouse.</p> <p>The keyword local in indicates that the file is present on the local file system and not on HDFS.</p>																				
hadoop fs -ls /user/hive/warehouse/ <database>.db	The path of the file is the hive installation directory.	<pre>c:\BigData\hadoop fs -ls /user/hive/warehouse/mca.db Found 3 items drwxr-xr-x - raoul supergroup 0 2019-11-13 18:08 /user/hive/warehouse/mca.db/customers drwxr-xr-x - raoul supergroup 0 2019-11-19 13:04 /user/hive/warehouse/mca.db/freshprods drwxr-xr-x - raoul supergroup 0 2019-11-13 23:04 /user/hive/warehouse/mca.db/orders</pre>																		

hadoop fs -cat /user/hive/warehouse/<database>.db/<tablename>/*	Check the contents	C:\BigData>hadoop fs -cat /user/hive/warehouse/mca.db/freshprods/* broccoli, Broccoli,5 spinach, Spinach,7 carrot, Local Carrots,4 potato, Idaho Potatoes,4 milk,Skimmed Milk,5 cheese,Spicy Garlic,5
hadoop fs -copyFromLocal <filename>.csv /<hdfs directory>/  hadoop fs -ls /<hdfs directory>/	Copy this file into HDFS and check whether it is copied or not.	C:\hive\apache-hive-2.1.0-bin\bin>hadoop fs -copyFromLocal freshproducts.csv /data/ C:\hive\apache-hive-2.1.0-bin\bin>hadoop fs -ls /data/ Found 2 items -rw-r--r-- 1 raoul supergroup 138 2019-11-20 08:15 /data/freshproducts.csv -rw-r--r-- 1 raoul supergroup 123 2019-11-15 12:49 /data/products.csv
load data inpath '/<hdfs directory>/<filename.csv>' into table <tablename>;	Load the same data again but from the HDFS directory.	hive> load data inpath '/data/freshproducts.csv' into table freshprods; Loading data to table mca.freshprods OK No rows affected (10.98 seconds)
	Check the contents using select command.	hive> select * from freshprods; OK broccoli Broccoli 5.0 spinach Spinach 7.0 carrot Local Carrots 4.0 potato Idaho Potatoes 4.0 milk Skimmed Milk 5.0 cheese Spicy Garlic 5.0 broccoli Broccoli 5.0 spinach Spinach 7.0 carrot Local Carrots 4.0 potato Idaho Potatoes 4.0 milk Skimmed Milk 5.0 cheese Spicy Garlic 5.0 12 rows selected (2.058 seconds)
The new data is appended to the existing table.		
hadoop fs -ls /user/hive/warehouse/<database>.db/<tablename>	Explore the warehouse to see how things have changed.	C:\hive\apache-hive-2.1.0-bin\bin>hadoop fs -ls /user/hive/warehouse/mca.db/freshprods Found 2 items -rwxr-xr-x 1 raoul supergroup 138 2019-11-19 13:04 /user/hive/warehouse/mca.db/freshprods/freshproducts.csv -rwxr-xr-x 1 raoul supergroup 138 2019-11-20 08:15 /user/hive/warehouse/mca.db/freshprods/freshproducts_copy_1.csv
hadoop fs -ls /<hdfs directory>/		C:\hive\apache-hive-2.1.0-bin\bin>hadoop fs -ls /data/ Found 1 items -rw-r--r-- 1 raoul supergroup 123 2019-11-15 12:49 /data/products.csv
The file has been moved from the HDFS directory to the warehouse directory.		

If data is loaded from a file on the local file system a copy will be made in Hive's warehouse.  
If a file from HDFS is loaded into the table it will be moved to Hive's warehouse directory.

load data local inpath '<filename.csv>' overwrite into table <tablename>;	To overwrite existing data in a table	hive> load data local inpath 'freshproducts.csv' . . > overwrite into table freshprods; Loading data to table mca.freshprods OK No rows affected (12.671 seconds)
	Check the contents of the table	hive> select * from freshprods; OK broccoli Broccoli 5.0 spinach Spinach 7.0 carrot Local Carrots 4.0 potato Idaho Potatoes 4.0 milk Skimmed Milk 5.0 cheese Spicy Garlic 5.0 6 rows selected (2.181 seconds)

#### Inserting from other tables

δ Data can be loaded into Hive from other tables too.

insert into table <tablename> select <field1>, <field2>,... from <tablename>	Merge the data of both products and freshprods table.	hive> insert into table products . . > select id, title, cost from freshprods;
The overwrite keyword can also be used with this insert command to completely remove the contents of the existing table.		
	Check the contents	hive> select * from products; OK broccoli Broccoli 5.0 spinach Spinach 7.0 carrot Local Carrots 4.0 potato Idaho Potatoes 4.0 milk Skimmed Milk 5.0 cheese Spicy Garlic 5.0 iphone7 iPhone 7 950.0 camera_canon Canon 570x 1000.0 washingmachine_samsung Samsung Swift 400.0 tv_vu Vu 56 Inch 600.0  12 rows selected (1.286 seconds)

☞ Load multiple tables from single table.

describe <tablename>; alter <tablename> change        column <field>       <field> <datatype> after <field>;	Use the products table as source.  Alter the table such that the id comes before title.	<pre>hive&gt; describe products; OK title string id string cost float 3 rows selected (0.975 seconds) hive&gt; alter table products . . &gt; change column title title string . . &gt; after id; OK No rows affected (1.806 seconds)</pre>
create        table <tablename> (<field1> <datatype>, .....);	Create a new table called prod_name that has id and name as its attributes.  Create another table prod_cost that has id, cost as its attributes.	<pre>hive&gt; create table prod_name (id string, name string); OK No rows affected (10.703 seconds)</pre> <pre>hive&gt; create table prod_cost (id string, cost float); OK No rows affected (0.661 seconds)</pre>
from <sourcetable> insert        overwrite table <table1> select <field1>,... insert into table <table2> select <field1>,...	To populate both the tables with one insert statement	<pre>hive&gt; from products . . &gt; insert overwrite table prod_name . . &gt; select id, title . . &gt; insert into table prod_cost . . &gt; select id, cost;</pre>
	Check the contents	<pre>hive&gt; select * from prod_cost; OK broccoli 5.0 spinach 7.0 carrot 4.0 potato 4.0 milk 5.0 cheese 5.0 iphone7 950.0 camera_canon 1000.0 washingmachine_samsung 400.0 tv_vu 600.0  12 rows selected (1.048 seconds)</pre>

		<pre>hive&gt; select * from prod_name; OK broccoli    Broccoli spinach     Spinach carrot      Local Carrots potato      Idaho Potatoes milk        Skimmed Milk cheese      Spicy Garlic iphone7     iPhone 7 camera_canon Canon 570x washingmachine_samsung Samsung Swift tv_vu       Vu 56 Inch  12 rows selected (0.936 seconds)</pre>
truncate table <tablename>;	Hive uses truncate command to delete all data within a table.	<pre>hive&gt; truncate table freshprods; OK No rows affected (0.811 seconds)</pre>
	Check the contents	<pre>hive&gt; select * from freshprods; OK No rows selected (1.544 seconds)</pre>

### Deleting and updating data in hive

- ⊞ Hive tables do not support row level deletes and updates by default
- ⊞ It is possible to get ACID compliant\* Hive tables by setting up special properties in hive-site.xml

### Partitioning and bucketing of data

- ⊞ Splits data into smaller, manageable parts
- ⊞ Enables performance optimizations
- ⊞ **Partitioning**
  - δ Data may be naturally split into logical units
  - δ Example
    - ▢ Customers portioned based on where they lived.
  - δ Each of these units will be stored in a different directory
  - δ State specific queries will run only on data in one directory
  - δ Splits may not be of the same size



## Bucketing

- δ Size of each split should be the same
- δ Hash of a column value - address, name, timestamp anything
- δ Each bucket is a separate file
- δ Makes sampling and joining data more efficient

## Complex data types

- ⊞ Array
- ⊞ Map
- ⊞ Struct
- ⊞ Union – Rarely uses, incomplete support in Hive

## Array

- ⊞ Collection data type
- ⊞ No fixed size
- ⊞ Entities of the same type
- ⊞ Only arrays of primitive types allowed

## Demo

create table <tablename> ( <field1> array<datatype>.... );	Create tables with arrays.	hive> create table mobilephones ( .. > id string, .. > title string, .. > cost float, .. > colors array<string>, .. > screen_size array<float> .. > );
	Check the schema	hive> describe mobilephones; OK id string title string cost float colors array<string> screen_size array<float> 5 rows selected (1.761 seconds)
insert into table <tablename>  select "stringvalue",  number,  array("stringvalue",.....),	Insert data into table with arrays	hive> insert into table mobilephones .. > select "redminote7", "Redmi Note 7", 300, .. > array("white", "silver", "black"), array(float(4.5)) .. > UNION ALL .. > select "motoGplus", "Moto G Plus", 200, array("black", "gold"), .. > array(float(4.5), float(5.5));

array(<datatype>(number),.....);												
	To select individual arrays	<pre>hive&gt; select id, colors from mobilephones; OK 1       Samsung J7       250.0  ["red","blue","black"] [5.5] 2       One Plus Three   450.0  ["gold","silver"] [4.5,5.5]</pre>										
Select arrayname[<index>];	To select the first element of the colors array	<pre>hive&gt; select id, colors[0] from mobilephones; OK 1       Samsung J7       250.0  red 2       One Plus Three   450.0  gold</pre>										
create table <tablename>( <field1> <datatype>,..... <arrayname> array<datatype>, .....)  row format delimited fields terminated by ``,`  collection items terminated by `#`;	Create a table from a file.	<pre>hive&gt; create table mobilephones ( . . &gt; id string, . . &gt; title string, . . &gt; cost float, . . &gt; colors array&lt;string&gt;, . . &gt; screen_size array&lt;float&gt; . . &gt; ) . . &gt; row format delimited fields terminated by ``,` . . &gt; collection items terminated by `#`; OK No rows affected (2.299 seconds)</pre>										
To identify the collection items, use the keyword collection with the termination value of #.												
	The file contents are	<table><tr><td>samsungj7</td><td>Samsung J7</td><td>250</td><td>red#blue#black</td><td>5.5</td></tr><tr><td>oneplusthree</td><td>One Plus Three</td><td>450</td><td>gold#silver</td><td>4.5#5.5</td></tr></table>	samsungj7	Samsung J7	250	red#blue#black	5.5	oneplusthree	One Plus Three	450	gold#silver	4.5#5.5
samsungj7	Samsung J7	250	red#blue#black	5.5								
oneplusthree	One Plus Three	450	gold#silver	4.5#5.5								
load data local inpath <'filenames.csv'> into table <tablename>;	Load the data using the load command.	<pre>hive&gt; load data local inpath 'mobilephones.csv' . . &gt; into table mobilephones; Loading data to table mca.mobilephones OK No rows affected (2.556 seconds)</pre>										
select * from <tablename>;	To check if it is loaded or not	<pre>hive&gt; select * from mobilephones; OK samsungj7 Samsung J7 250.0 ["red","blue","black"] [5.5] oneplusthree One Plus Three 450.0 ["gold","silver"] [4.5,5.5] 2 rows selected (1.782 seconds)</pre>										

## Map

- ⊞ Unordered collection of pairs
- ⊞ No fixed size
- ⊞ Every entity is a key, value pair
- ⊞ Value is accessed using a unique key
- ⊞ Keys and values have their own data types
  - δ Key can be integer, value can be a struct
  - δ Key can be a string, value can be an integer.

## Demo

<pre>create table &lt;tablename&gt; (   &lt;field1&gt;   &lt;datatype&gt;,   .....   &lt;fieldn&gt;   map&lt;&lt;keydatatype&gt;,   &lt;valuedatatype&gt;&gt; )  row format delimited fields terminated by '\t'  collection items terminated by '#'  map keys terminated by ':';</pre>	<p>Create the table with the map data type.</p> <p>The key is of type string</p> <p>The value is of type boolean</p> <p>The map keys are delimited by :.</p>	<pre>hive&gt; create table mobilephones ( . . &gt; id string, . . &gt; title string, . . &gt; cost float, . . &gt; colors array&lt;string&gt;_ . . &gt; _ . . &gt; screen_size array&lt;float&gt;,_ . . &gt; features map&lt;string, boolean&gt;_ . . &gt; ) . . &gt; row format delimited fields terminated by ',' . . &gt; collection items terminated by '#' . . &gt; map keys terminated by ':'_ OK No rows affected (2.467 seconds)</pre>												
<p>The file contains the following data</p>	<table><tr><td>samsungj7</td><td>Samsung J7</td><td>250</td><td>red#blue#black</td><td>5.5</td><td>camera:true#dualsim:false</td></tr><tr><td>oneplusthree</td><td>One Plus Three</td><td>450</td><td>gold#silver</td><td>4.5#5.5</td><td>autofocus:true</td></tr></table>	samsungj7	Samsung J7	250	red#blue#black	5.5	camera:true#dualsim:false	oneplusthree	One Plus Three	450	gold#silver	4.5#5.5	autofocus:true	
samsungj7	Samsung J7	250	red#blue#black	5.5	camera:true#dualsim:false									
oneplusthree	One Plus Three	450	gold#silver	4.5#5.5	autofocus:true									

load data local inpath <'filename.csv'>  into table <tablename>;	Load the data	<pre>hive&gt; load data local inpath 'mobilephones.csv' . . &gt; into table mobilephones; Loading data to table mca.mobilephones OK No rows affected (1.424 seconds) hive&gt; select * from mobilephones;</pre>
Arrays are identified by square brackets; map is identified by curly braces		<pre>hive&gt; select * from mobilephones; OK samsungj7 Samsung J7 250.0 ["red","blue","black"] [5.5] {"camera":true,"dualsim":false} oneplusthree One Plus Three 450.0 ["gold","silver"] [4.5,5.5] {"autofocus":true} 2 rows selected (0.843 seconds)</pre>
select <mapname> from <tablename>;	Accessing maps	<pre>hive&gt; select id, features from mobilephones; OK samsungj7 {"camera":true,"dualsim":false} oneplusthree {"autofocus":true} 2 rows selected (3.285 seconds)</pre>

### Struct

- ⊞ Logical grouping of data
- ⊞ Can have different data types
- ⊞ Can hold any number of values
- ⊞ Each value referenced by a name

### Demo

<pre>create table &lt;tablename&gt;( &lt;field&gt; &lt;datatype&gt;,..... &lt;field&gt; struct&lt;&lt;field&gt;:&lt;datatype&gt;, &lt;field&gt;: &lt;datatype&gt;,.....&gt; row format delimited fields terminated by `;`;</pre>	<p>Create the table mobilephones</p> <p>Information holds specific info about the phone.</p> <p>The info is battery life and kind of camera.</p>	<pre>hive&gt; create table mobilephones ( . . &gt; id string, . . &gt; title string, . . &gt; cost float, . . &gt; colors array&lt;string&gt;, . . &gt; screen_size array&lt;float&gt;, . . &gt; features map&lt;string, boolean&gt;, . . &gt; information struct&lt;battery:string, camera:string&gt; . . &gt; ) . . &gt; row format delimited fields terminated by `;` . . &gt; collection items terminated by `#` . . &gt; map keys terminated by `:`; OK No rows affected (0.805 seconds)</pre>																		
The file content is as	<table><tr><td>samsungj7</td><td>Samsung J7</td><td>250</td><td>red#blue#black</td><td>5.5</td><td>camera:true#dualsim:false</td><td>24 hours#2MP</td></tr><tr><td>oneplusthree</td><td>One Plus Three</td><td>450</td><td>gold#silver</td><td>4.5#5.5</td><td>autofocus:true</td><td>12 hours</td></tr></table>						samsungj7	Samsung J7	250	red#blue#black	5.5	camera:true#dualsim:false	24 hours#2MP	oneplusthree	One Plus Three	450	gold#silver	4.5#5.5	autofocus:true	12 hours
samsungj7	Samsung J7	250	red#blue#black	5.5	camera:true#dualsim:false	24 hours#2MP														
oneplusthree	One Plus Three	450	gold#silver	4.5#5.5	autofocus:true	12 hours														

shown.  The first data value goes to the battery field and second goes to the camera field.		
describe <tablename>;	Check the schema of the table.	<pre>hive&gt; describe mobilephones; OK id string title string cost float colors array&lt;string&gt; screen_size array&lt;float&gt; features map&lt;string,boolean&gt; information struct&lt;battery:string,camera:string&gt; 7 rows selected (0.217 seconds)</pre>
load data local inpath <'filename.csv'>  into table <tablename>;	Load the data	<pre>hive&gt; load data local inpath 'mobilephones.csv' . . &gt; into table mobilephones; Loading data to table mca.mobilephones OK No rows affected (1.558 seconds)</pre>
Display the data		<pre>hive&gt; select * from mobilephones; OK samsungj7 Samsung J7 250.0 ["red","blue","black"] [5.5] {"camera":true,"dualsim":false} {"battery":"24 hours","camera":"2MP"} oneplusthree One Plus Three 450.0 ["gold","silver"] [4.5,5.5] {"autofocus":true} {"battery":"12 hours","camera":null} 2 rows selected (1.127 seconds)</pre>
Select the map and structure details.		<pre>hive&gt; select features, information from mobilephones; OK {"camera":true,"dualsim":false} {"battery":"24 hours","camera":"2MP"} {"autofocus":true} {"battery":"12 hours","camera":null} 2 rows selected (2.172 seconds)</pre>
Select individual values in a struct.		<pre>hive&gt; select features['camera'], information.battery . . &gt; from mobilephones;  OK true 24 hours 12 hours 2 rows selected (2.035 seconds) hive&gt;</pre>

### Built-in functions

- ⊞ UDF
  - δ User-defined Functions
- ⊞ UDAF
  - δ User-defined Aggregate Functions
- ⊞ UDTF
  - δ User-defined Table-generating Functions

**UDF**

- ⌘ Works on a single row
- ⌘ Outputs a single row
  - δ trim(), concat(), length() round(), floor()
- ⌘ <https://cwiki.apache.org/confluence/display/Hive/languageManual+UDF>

**UDAF**

- ⌘ Works on multiple rows
- ⌘ Outputs a single row
  - δ count(\*), sum(), avg()

**UDTF**

- ⌘ Works on a single row
- ⌘ Outputs multiple rows
  - δ explode(), posexplode()

**Table-generating functions**

- ⌘ explode()
  - δ Flatten the data in arrays and maps

Manager	SubordinateList	Subordinate	Employee	Details	Key	Value
Larry	[Sundar, Eric, Jon]	Sundar	Larry	{"office": "271B", "numReports": 8, "salary": 1}	office	271B
Sergey	[Ruth, Urs]	Eric	Sergey	{"office": "271B", "numReports": 5, "salary": 1}	numReports	8
Sundar	[Susan, Alan, Lazlo]	John	Sundar	{"office": "285", "numReports": 12}	salary	1
		Ruth			office	271B
		Urs			numReports	5
		Susan			salary	1
		Alan			office	285
		Lazlo			numReports	12

**Demo**

select explode(<fieldname>) as <alias name> from <tablename>	Use mobilephones  To explode the values of	hive> select explode(colors) as variants from mobilephones;
---	--	---

	colors array	<pre>OK 0 red 1 blue 2 black 3 gold 4 silver 5 rows selected (4.227 seconds)</pre> <pre>hive&gt; select colors from mobilephones; OK [" red","blue","black"] [" gold","silver"] 2 rows selected (2.621 seconds)</pre>
<pre>select explode(&lt;fieldname&gt;) as (&lt;field1&gt;, &lt;field2&gt;,....) from &lt;tablename&gt;;</pre>	Explode a map data type	<pre>hive&gt; select explode(features) as (feature, present) . . &gt; from mobilephones;</pre> <pre>OK camera true dualsim false autofocus true 3 rows selected (0.813 seconds)</pre>

- ☞ To see that every row of the flattened row is associated with an index value use posexplode.

```
hive> select posexplode(colors) as (index, variants) from mobilephones;
```

```
OK
0 red
1 blue
2 black
3 gold
4 silver
5 rows selected (0.825 seconds)
```