



# DA 440 – Apache Hive Essentials Slide Guide

---

Fall 2015 – Version 5.0.0

This Guide is protected under U.S. and international copyright laws, and is the exclusive property of MapR Technologies, Inc.

© 2016, MapR Technologies, Inc. All rights reserved. All other trademarks cited here are the property of their respective owners.





## DA 440: Apache Hive Essentials

Lesson 1: Apache Hive in the Apache Hadoop Ecosystem

© 2015 MapR Technologies 

1

Welcome to DA 440: Apache Hive Essentials – Lesson 1: Apache Hive in the Hadoop Ecosystem. In this lesson, you will learn what Apache Hive is and how it fits into the Hadoop ecosystem.



## Learning Goals

- ▶ Describe what Hive is and when to use it
- ▶ Describe the data pipeline and use cases
- ▶ Describe how Hive fits in the Hadoop ecosystem
- ▶ Describe data types in Hive

© 2015 MapR Technologies MAPR

2

When you have finished this lesson, you will be describe what Hive is as well as when to use it.

You will also be able to describe the data pipeline, how Hive fits into the data pipeline, and what the pros and cons of using Hive are.

You will learn how Hive fits into the broader Hadoop ecosystem, including how Hive uses MapReduce, how Hive interfaces with other Hadoop tools, and some of the differences between Hive versions.

Finally, we will discuss data types and casting in Hive.



## Learning Goals

- ▶ Describe what Hive is and when to use it
- ▶ Describe the data pipeline and use cases
- ▶ Describe how Hive fits in the Hadoop ecosystem
- ▶ Describe data types in Hive

Let's start with a brief overview of Hive, as well as other SQL-on-Hadoop tools.



## Too Much Data?



You work in a large climate research center. You

- already know SQL
- have accumulated petabytes of weather measurements
- migrated data to HDFS

How do you process this data?

© 2015 MapR Technologies 

4

We are working at an international Climate Research Institute, and you are a data scientist. This research center manages thousands of weather stations around the world. Each weather station provides a number of data points: temperature, wind speed, or a number of other measurements. This data is sent in aggregated batches every month.

This data is organized in tables. Each row contains columns representing the station name, the date the measurements were taken, the temperature, and the wind speed.

You, as a data scientist at the Climate Research Institute who is already familiar with SQL, might want to query this data with SQL.

However, the research institute has accumulated hundreds of petabytes of data. These massive tables are too large for a traditional Relational Database Management System to handle. The Climate Research Institute has decided to migrate all its data to a Hadoop Distributed File System. The HDFS can store, access, and backup petabytes of data easily.

However, traditional SQL does not work on the HDFS. What are your options as a data scientist? One option is Apache Hive.



## What are Hive & HiveQL?



Hive is

- a data warehouse application
- part of the Hadoop ecosystem

Hive Query Language is

- called HiveQL (HQL)
- A SQL-like language used to explore data in HDFS

© 2015 MapR Technologies 

5

What is Hive? Apache Hive was initially developed at Facebook for the purpose of summarizing, querying, and analyzing large amounts of data stored on a Hadoop Distributed File System.

Hive is often used to analyze large databases, such as customer data for online retailers, social media user accounts, or research data like that used at our Climate Research Institute.

Hive is a data warehouse application built on top of Hadoop. Hive can interact with many parts of the Hadoop ecosystem using a query language called HiveQL.

HQL is an SQL-like language that can query data in an HDFS. Hive combines the power of lower level MapReduce code with the ease of a higher level language like SQL.

Since HQL is similar to SQL, it's easy for analysts who are already familiar with SQL to learn HQL. However, keep in mind that HQL is not ANSI-SQL, and this course will discuss some of the differences between HQL and SQL.

These are some of the reasons why Hive is ideal for a situation where your data scientists already know SQL, but your database has become too large for traditional relational database management systems.

The infographic is titled "How Do You Use Hive?" and features four yellow hexagonal boxes arranged in a 2x2 grid. Each box contains a bulleted list of methods:

- Top-left box: • Command line interface (CLI)
  - Hive Shell, Beeline
- Top-right box: • Graphical user interface (GUI)
  - Hue, Beeswax
- Bottom-left box: • Other Hadoop applications
  - HCatalog, Tez, Spark, HBase
- Bottom-right box: • Other languages or applications
  - JDBC
  - ODBC

© 2015 MapR Technologies 6

Now that we have decided Hive is the right tool for us, how can we use Hive?

Early versions of Hive, this command line interface is sometimes referred to as the Hive Shell. However, in version 0.11 a newer command line interface, called Beeline, was introduced. Depending on the version of Hive, HiveServer, and Hadoop you are running, the CLI may be called the Hive Shell or Beeline. There are only minor differences between the two; consult the Hive documentation for more details about these.

For users not comfortable with the command line, there are graphical user interface options like Hue and Beeswax. We will be using the Hive Shell in the labs for this course.

Hive can also be used in conjunction with other Hadoop applications such as Hcatalog, Tez, Spark, and Hbase. Hive is often used in conjunction with Pig as part of an Extract, Transform, Load data pipeline.

Finally, Hive can be used with other programming languages and other applications outside of the Hadoop ecosystem. JDBC allows users to use languages like Java, Python, and Ruby with Hive. ODBC allows users to connect Hive to applications like Microsoft Excel or Tableau.

The screenshot shows the HUE Query Editor interface. At the top, there's a red header bar with the HUE logo and the text "CLI vs GUI". Below the header, there are two separate code input boxes. The left box contains the Beeline command: "beeline> SELECT \* FROM rothera;". The right box contains the Hive command: "hive> SELECT \* FROM rothera;". To the right of these boxes is the main HUE Query Editor window. It has tabs for "Hive Editor" and "Query Editor", with "Query Editor" currently selected. The "Query Editor" tab has sub-options like "Assist" and "Settings". A dropdown menu labeled "DATABASE" is set to "db1". Below the database dropdown is a "Table name..." input field and a list of tables: "chai", "chai1", "cutoversimple", "deltaonecolmorev...", and "temp\_currentview". At the bottom of the editor window are buttons for "Execute", "Save as...", "Explain", and "New query". Below the editor are tabs for "Recent queries", "Query", "Log", "Columns", and "Results". In the bottom right corner of the main window, there's a copyright notice: "© 2015 MapR Technologies" and the MAPR logo.

Here we see the same query, `SELECT * FROM rothera;` entered into the Hive Command Line Interface, the Beeline Shell, and the HUE Graphical User Interface.

As you can see, the syntax of the queries themselves is the same in all three interfaces. There may be slight differences in how each interface handles errors or accesses the file system as well. For the most part, however, the Hive QL you would use is the same regardless of which interface you choose.

INSTRUCTOR: Now is a good time to demonstrate each of these three interfaces. You should have them open and running off-screen, with some sort of data already loaded, so you can easily access them at this point of the lecture.



You should have already installed and connected to the MapR Sandbox during Lesson 0. Try connecting to the Hive CLI, which is how we will primarily be using Hive in this course. You can also try connecting to the Beeline Shell.

Once you have connected to the Hive CLI, you can continue with lesson 1.



## When should you use Hive?

- Easy to learn if you already know SQL
- Widely used in the Hadoop ecosystem
- Good version compatibility
- Query
  - large amounts of data
  - structured data
  - in batches



© 2015 MapR Technologies 

9

Let's review: How do you know if Hive is the right tool for you?

First, it's easy to learn if you already know SQL. Since many data scientists and business analysts are already familiar with SQL, Hive is a great introduction to Hadoop, allowing users to run MapReduce jobs without learning Java code. Hive is also learned by many developers and administrators who work with data scientists.

Another benefit of Hive is that it is widely used in the Hadoop ecosystem. There is a large support network available for learners of Hive: many books, forums, and other groups can provide support as you dive deeper into Hive.

Next, Hive has great compatibility with other tools in the Hadoop ecosystem. Hive has been around long enough that it has been tested thoroughly. Hive works, and each new version has fewer bugs and compatibility issues. While newer tools may have additional features, Hive has stood the test of time, much like SQL has.

Unlike SQL, however, Hive is built to handle very large data sets. While SQL may fail on a large dataset, Hive will reliably complete.

Hive is also good for structured data, like tables. Unstructured data, such as JSON, may require some manipulation before it is usable in Hive. If your data is structured, or it can become structured through an ETL process, then Hive might be a good tool for you.

Finally, although you can do exploratory analysis in Hive, it is commonly used for batch queries. If you already know what you want to query, and can write batch scripts for your dataset, Hive might be a good choice for you.



## When to use SQL or HQL?

SQL	HQL
Works on RDBMS	Works on HDFS
Interactive; queries provide results in near real time	Batch; queries have overhead due to MapReduce
Works best on small or medium datasets (megabytes or gigabytes)	Works best on large datasets (terabytes or petabytes)

© 2015 MapR Technologies

10

This chart summarizes some common use cases for SQL or Hive.

The type of database you are working on, whether a relational database or a Hadoop database, is an important factor when deciding when to use Hive or SQL.

You should be aware that SQL is interactive, while Hive is batch. SQL queries produce results in real-time, while complex Hive queries may take several hours or even days to run.

Finally, the size of the database is also important: SQL works best on small or medium datasets, while Hive is optimized for very large datasets.

In the case of our weather data set, our Climate Research Institute is facing an ever-growing database which we recently moved to an HDFS. We are more concerned with trends over time rather than real-time updates. These long term trends can be queried in batches. Therefore, we will choose Hive over SQL.



## Knowledge Check

Who should be familiar with Hive? Check all that apply.

- a) Data scientists working with data on an HDFS
- b) Business analysts working with data on an RDBMS
- c) Data analysts who want interactive, real time queries
- d) Hadoop developers who work with data scientists

Answers: A and D

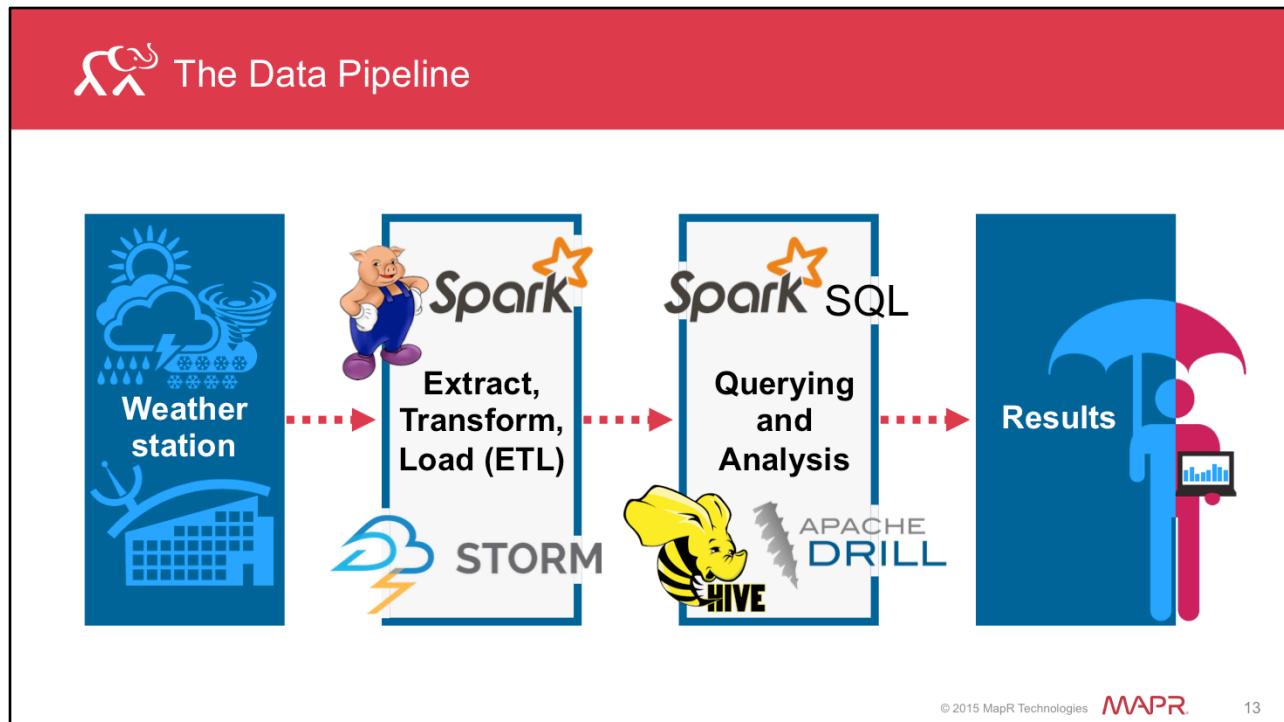
INSTRUCTOR: Take some time to have your students introduce themselves. What role(s) do they see themselves performing with Hive? Are they admins or developers who support data scientists? Are they data scientists or business analysts interested in big data? Are they curious lifelong learners who are taking the course to expand their knowledge of big data technologies?



## Learning Goals

- ▶ Describe what Hive is and when to use it
- ▶ Describe the data pipeline and use cases
- ▶ Describe how Hive fits in the Hadoop ecosystem
- ▶ Describe data types in Hive

In this next section, we'll look at other tools in the data pipeline, and where Hive fits in.



© 2015 MapR Technologies 

13

Let's consider the data pipeline from the perspective of a data scientist at our Climate Research Institute.

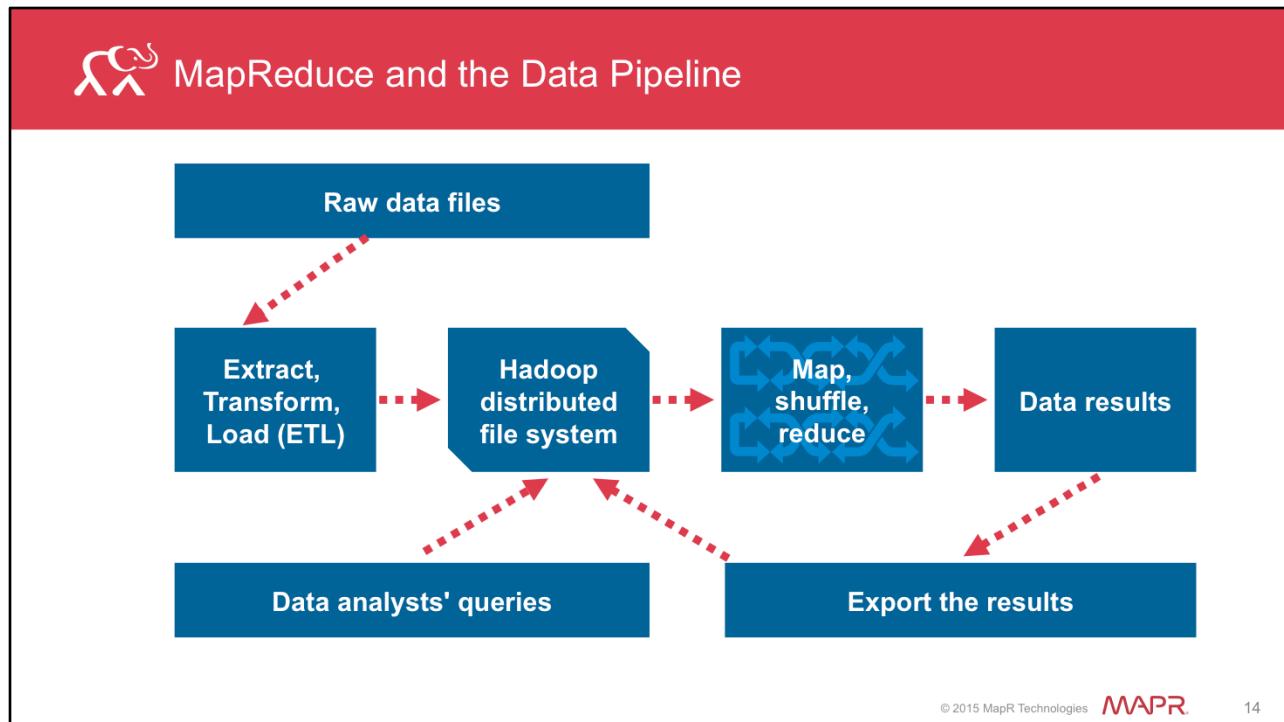
The source of raw data will vary from one project to the next. One example would be the Internet of Things. The Internet of Things refers to networked sensors, such as those found in smart cars or weather stations, that automatically send their data to some central location, such as the Climate Research Institute servers.

However, this data may not be in a structured table, or it may need to be normalized through rounding or some other algorithm. Tools such as Pig, Storm, and Spark would be good choices for this. These tools can automate the process of data ingestion, normalization, structuring, and exporting into the Hadoop Distributed File System with one or more scripts.

While some parts of the early ETL process can be performed by Hive, other tools such as Pig are better suited for this step of the data pipeline.

Once the data is loaded, a data scientist can query and analyze the transformed data. For example, you might want to find out the average temperature in August at weather station across our entire dataset; this is the type of query you can easily perform in Hive. Depending on the nature of your data and your project, there are again a number of tools to choose from, including other SQL-on-Hadoop tools like Drill or SparkSQL.

Finally, you can export the results of your query to use in a report, which you might create with other programs to generate graphs, write up your analysis, or make slide decks for presentations.



This diagram represents the data pipeline from a more technical perspective.

The data pipeline starts with raw data files. This might be click stream data from a website, log files from a server, JSON files, or any other type of data, including the raw sensor data from our weather stations.

In most cases, a raw data file must be extracted, transformed, and loaded into a Hadoop Distributed File System. This includes any data normalization processes you must perform, such as splitting a string, defining delimiters, or rounding, truncating, or otherwise aggregating data.

Once the data is structured into a table, the data analyst can connect to the HDFS through a Hive client, such as the Hive CLI, the Beeline Shell, a GUI like Hue, or another program such as Spark, Tez, or Thrift.

Once the analyst submits their query, Hive works by converting the SQL-like HiveQL into MapReduce code. Therefore, once the Hive code is run, a MapReduce job is queued. This MapReduce process can take several minutes, hours, or even days for larger datasets, so it is important to be patient. Hive is very reliable, and the process will eventually finish in the majority of cases.

Once the MapReduce process is finished, results are output into the command line or other user interface. The analyst can also choose to store these results as a another table or as .csv file for later use.

These exported results can be loaded back into the HDFS, and the cycle can continue; analysts can continue querying manipulated data in this way.



## Discussion – How do you interact with data?

Think about the data you work with, as a developer or data analyst. What part of the data pipeline are you usually involved with? How do you think you might use Hive in your work?

© 2015 MapR Technologies MAPR

15

**INSTRUCTORS:** Take some time to get to know your students, what their skill levels are, and what their goals for this course are.



## Hive & other SQL-on-Hadoop tools

- First SQL-on-Hadoop tool
- New tools like Drill:
  - often faster than Hive
  - may lack key functions
  - may suffer compatibility issues
- Familiarity with SQL & Hive helps you learn new tools



© 2015 MapR Technologies

16

Now that we know where in the data pipeline Hive fits, what about the other SQL-on-Hadoop tools that fit in the same place?

Hive is not the only SQL-on-Hadoop tool, but it was the first, and is compatible with most versions of Hadoop. This also means many data analysts already are familiar with Hive, and may have already written scripts in HQL.

Other SQL-on-Hadoop tools like Impala, Drill, and SparkSQL exist as well, but may not be compatible with all versions of Hadoop. Many data scientists first learn SQL, then Hive, then move on to one of these new tools when they find they need speed or a certain functionality which Hive lacks.

If you would like to learn more about Drill, refer to DA 410: Drill Essentials or DEV 360: Spark Essentials in the MapR Academy course catalog.



## When to use Hive or Pig?

Hive (and HiveQL)	Pig (and Pig Latin)
Declarative language	Procedural language
Used mostly by data analysts	Used by both data analysts and developers
Used to run batch queries on structured data	Used to automate ETL for unstructured data



© 2015 MapR Technologies 

17

Since Pig and Hive are often used together, and there is some overlap in their capabilities, many users want to know when to use one or the other. It is important to note that both of these tools can perform many of the same tasks, but they are often more suited for one task or another.

Hive is a non-procedural, declarative language while Pig is a procedural language. In other words, Hive will process your query in whatever order is most optimized, while Pig will process your query in the order you specify.

Hive is mostly used by data analysts to explore data sets, while Pig is used by both data analysts who want to clean their own data, and by developers who want to process data.

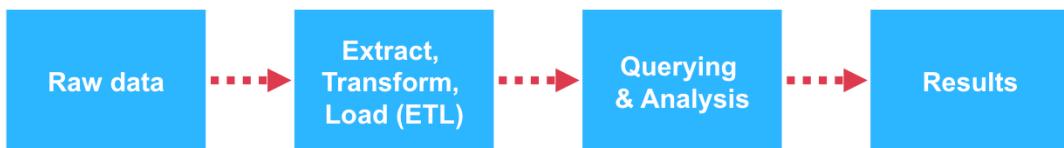
Hive runs batch queries on structured data while Pig can automate ETLs for unstructured data.

In the case of our Climate Research Institute, we will ultimately be using both Hive and Pig. Pig will be used for initial data processing and normalization, while Hive will be used for querying and data analysis.



## Knowledge Check

Give an example of tools or sources used at each of these stages of the data pipeline.



© 2015 MapR Technologies

18

Some possible responses:

- Box 1: Internet of Things
- Box 2: Pig, Spark
- Box 3: Hive, Drill
- Box 4: graph/report



## Learning Goals

- ▶ Describe what Hive is and when to use it
- ▶ Describe the data pipeline and use cases
- ▶ Describe how Hive fits in the Hadoop ecosystem
- ▶ Describe data types in Hive

In this next section, we'll learn a little bit about the technical specifications in Hive, including features found in different versions of Hive and how Hive interacts with other tools in the Hadoop ecosystem

## HiveQL & MapReduce

Hive queries

- work on HDFS
- sent in batches
- Queries sent from Beeline or Hive CL
- Most queries trigger MapReduce

© 2015 MapR Technologies **MAPR** 20

It is important to keep in mind that Hive queries are sent in batches.

Because they are sent in batches, and because some queries trigger MapReduce jobs, which requires a lot of overhead and can use a lot of computational resources. Each step in a Hive query, including joins, selects, or any other computational step, may queue up its own MapReduce job, adding additional computational overhead to the process.

An dataset may be several petabytes in size, and a Hive query may take hours or even days to process due to the computational overhead, but it is reliable and will finish eventually.

For users who are accustomed to the near-instantaneous results of SQL, it may appear that the query has failed. This is often not the case, and it is important to be patient when using Hive.



## Options for optimization of Hive

- Skip reads
  - Partition, ORC
- Minimize shuffle
  - Bucket, Sort
- Hive on Spark
- Hive on Tez



© 2015 MapR Technologies 

21

Speed is a major concern for many Hive users. There are several options for optimizing the performance of Hive.

One way to increase the speed of your hive queries is to skip reads, which minimizes the amount of data Hive has to read in order to perform your query. One way to do this is partitioning tables. Telling Hive to only read a part of a table, rather than all your data, can significantly improve performance.

Another way to skip reads is to take advantage of Ordered Record Columnar files. ORC files, when properly sorted, can also minimize the number of rows Hive must read to find the data which matches your query.

The second common way to optimize Hive queries is to minimize the time spent in the shuffle phase of the MapReduce process. Proper use of bucketing and sorting, especially when joining tables, can help reduce time in the shuffle phase.

Using Hive in conjunction with other Hadoop tools may also provide a performance boost.

Hive on Spark gives users the option to use Hive queries on data, but instead of triggering a MapReduce job, Hive on Spark uses the Apache Spark engine. This means your queries may be faster.

Likewise, Hive on Tez provides another alternative for Hive users to bypass MapReduce jobs.

While we will not cover using Hive on Spark or Hive on Tez in this course, it is important to be aware that these are options you can consider in the future.



## Differences between versions of Hive

```
$ hive --version
```

<http://doc.mapr.com/display/MapR/Ecosystem+Support+Matrix>

Feature	0.12	0.13	1.0	1.1
DATE, TIMESTAMP, VARCHAR data types	✓	✓	✓	✓
Support for Apache Tez	✗	✓	✓	✓
Implicit JOIN; IN, NOT in subqueries; CHAR	✗	✓	✓	✓
Transactions with ACID semantics	✗	✗	✓	✓
XML data type	✗	✗	✓	✓
Support for Apache Spark	✗	✗	✗	✓

© 2015 MapR Technologies MAPR

22

When using Hive, it is important to know what version you are using. You can check this from a bash shell on your MapR Sandbox by typing `hive --version`. Once you know the version of Hive you are using, you can refer to the MapR Ecosystem Support Matrix to make sure it is compatible with other tools you might wish to use in your data pipeline.

It is also important to know that Hive is not ANSI-SQL compliant. While each new version of Hive adds more features, certain features of SQL have not been implemented. This table demonstrates some of the different features available.

For example the use of IN or NOT in subqueries, was not available until version 0.13. Likewise, transactions with ACID semantics and advanced datatypes such as XML were not available until version 1.0.

Version 1.1 features additional support for interacting with many other tools in the Hadoop ecosystem, including Apache Spark.

**INSTRUCTORS:** discuss some of these features with your students; which are necessary for their work?



## Knowledge Check

How does using Hive in conjunction with another tool like Spark increase performance speed?

- a) Spark is a lower-level language than Hive
- b) Spark does not run on an HDFS
- c) Spark does not use resource-heavy MapReduce jobs
- d) Spark only runs with newer versions of Hive

© 2015 MapR Technologies MAPR

23

Answer: C

INSTRUCTORS: Discuss other tools your students use, or have heard of. What are the pros and cons of each?



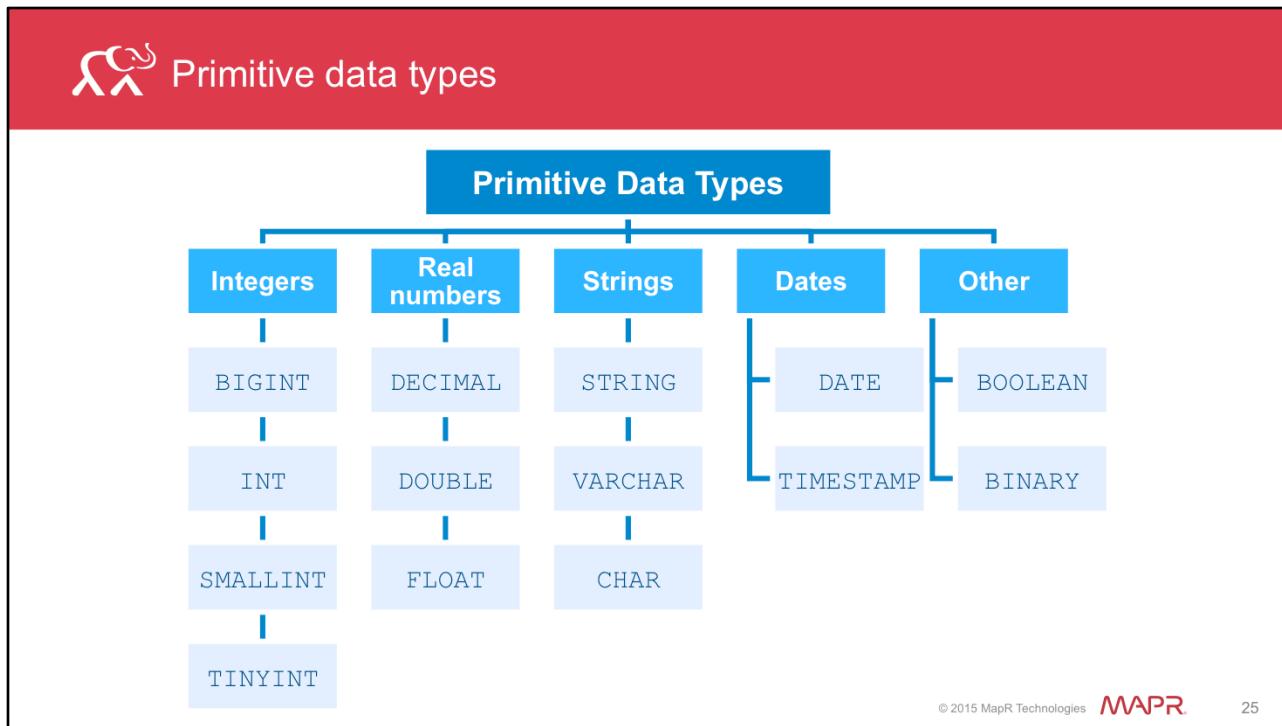
## Learning Goals

- ▶ Describe what Hive is and when to use it
- ▶ Describe the data pipeline and use cases
- ▶ Describe how Hive fits in the Hadoop ecosystem
- ▶ Describe data types in Hive

© 2015 MapR Technologies MAPR

24

Next, we'll learn about data types in Hive.



Hive tables can handle a variety of data types, but each column can only have one data type. You must specify whether each column contains an integer, a float, a string, a time stamp, an array, or some other data type whenever you create a table.

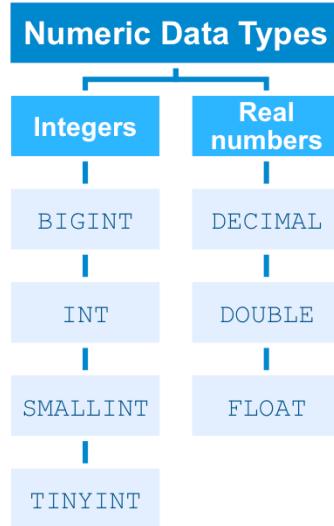
In most cases, Hive is able to cast up the data hierarchy, but not down. For example, a tinyint can be recast as an bigint, but a bigint might not be able to be recast as a tiny int. Likewise, a char can be cast as a string, but not vice-versa.

This diagram shows most of the primitive data types currently available in Hive. We will briefly discuss each, and then turn to complex data types.



## Numeric data

- Integers are signed whole numbers
  - TINYINT -2<sup>7</sup> to 2<sup>7</sup>-1
  - SMALLINT -2<sup>15</sup> to 2<sup>15</sup>-1
  - INT -2<sup>31</sup> to 2<sup>31</sup>-1
  - BIGINT -2<sup>63</sup> to 2<sup>63</sup>-1
- Real numbers include decimals
  - FLOAT 4 bytes
  - DOUBLE 8 bytes
  - DECIMAL 32 bytes
- **NUMERIC, MONEY** datatypes from SQL are not available in HiveQL



© 2015 MapR Technologies 

26

All integers are whole numbers. Integers in Hive are signed, meaning they can be either positive or negative.

The integer data type includes tinyint, smallint, int, and bigint, which correspond to 1, 2, 4, and 8 bytes of information, respectively.

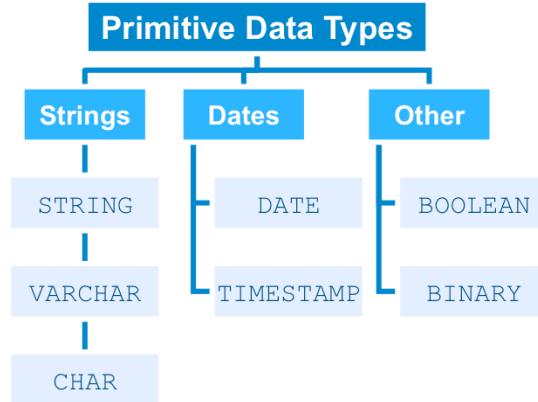
Real numbers includes floats, doubles, and decimals. These can contain decimals of varying precision, with the DECIMAL datatype containing up to 32 digits.

Some commonly used data types in SQL, such as NUMERIC or MONEY, are not available in Hive.



## Other primitive data types

- Strings
  - 'Hello world!'
- Dates
  - DATE 1970-01-24
  - TIMESTAMP  
1970-01-24 08:52:48.123
- Other
  - TRUE or FALSE
  - Array of bytes



© 2015 MapR Technologies 

27

Other common primitive data types include strings, dates, and others.

Strings may be any length, but varchar and char are limited in length. Varchars may be specified as any length between 1 and 65355 characters long. Chars have a maximum of 255 characters.

Dates are useful for ordered data. The DATE data type takes the format YYYY-MM-DD, so that January 24<sup>th</sup>, 1970 would be expressed as 1970 dash 01 dash 24.

Timestamp on the other hand is the number of seconds since Unix time began. A timestamp includes the date in year-month-day format, followed by the time, in hour-minute-second format, followed by 3 decimal places, measuring thousands of a second.

Other primitive data types include Boolean, which is true/false. The binary data type is an array of bytes that is not interpreted by Hive.



## Complex data types

- **ARRAY<datatype>**
  - stationNames ARRAY<STRING>
  - {"rothera", "airport", "mountain"}
- **MAP<primitive,datatype>**
  - stationIDs MAP<INT, STRING>
  - {201:"airport", 403:"mountain"}
- **STRUCT<colname:datatype,...>**
  - stationLocation STRUCT<name:STRING, longitude:INT, latitude:INT>
  - {name:"rothera", longitude:67, latitude:68}

© 2015 MapR Technologies MAPR

28

Hive can also handle more complex data types, such as indexed arrays, key-value maps, ordered key-value maps called structs, and uniontypes.

Arrays must contain objects of the same type. You could have an array of strings or an array of ints, but you cannot have an array which contains both strings and ints. For example, you might have an array of strings describing each of your weather stations

A map pairs one primitive datatype with one other datatype. Maps are unordered by default. For example, if each of your weather stations was also assigned a numeric ID number, you might map this ID, as an INT, to the station name, as a STRING.

A struct is a more complex datatype, which may contain several columns of several datatypes. For example, you might have a struct called stationIDs, which maps your station names, which are strings, and their latitude and longitude, which are INTs. In this example, our struct maps the Rothera station with its longitude and latitude, 67 degrees south and 68 degrees west.



In this lab activity, we will load populate a small table with some sample data to learn how Hive deals with different data types. We will discuss the details of creating and loading tables in the next lesson.

Open your Lab Guide and complete Lab 1.2. When you are finished, you are ready to move on with the lecture.



## Discussion – Data Types

What data types do you use in your work? Do you think it is more appropriate to cast numbers as tinyints, bigints, or floats in your database? Why?

© 2015 MapR Technologies MAPR

30

**INSTRUCTORS:** Have an open-ended discussion with your students about the datatypes found in their different industries. You may want to include a more in-depth discussion of complex datatypes, strings, or numerical data depending on the audience. You might also want to discuss the importance of thinking about the data types you are working with before you create data definitions.



**Next Steps**

**Lesson 2**  
Creating and Loading  
Data in Apache Hive

© 2015 MapR Technologies **MAPR** 31

In the next lesson, we will explore how to create tables and load data into Hive.



## **DA440: Apache Hive Essentials**

Lesson 2: Create and Load Data in Apache Hive

© 2015 MapR Technologies 

1

Welcome to DA440 Lesson 2: Create and Load Data in Apache Hive!

In this lesson, we will learn about Hive's data definition language to create tables and load them with data.



## Learning Goals

- ▶ Create databases and tables
- ▶ Partition and bucket data
- ▶ Load tables with data
- ▶ Alter and drop tables

By the end of this lesson, you will be able to:

- Create databases and simple tables
- Create external and partitioned tables
- Load tables with data
- Alter and drop tables



## Learning Goals

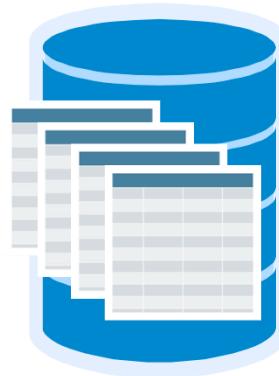
- ▶ Create databases and tables
- ▶ Partition and bucket data
- ▶ Load tables with data
- ▶ Alter and drop tables

Let's start by creating databases and simple tables.



## Databases in Hive

- Databases contain tables
  - default is used if you do not specify a database
- Tables contain data
  - HiveQL data definition creates metadata
- Hive itself does not contain data
  - All data is stored on HDFS
  - Hive provides metadata structure for accessing data on the HDFS



© 2015 MapR Technologies

4

Databases in Hive contain tables, and tables contain data. Hive databases and tables serve as a structure that can be queried with HiveQL; however, Hive itself does not store any data.

Whenever you create a table in Hive, you must specify the database you create it in. If you do not specify which database your table is in, the “default” database is used.

If there are a large number of users on a cluster, grouping tables together into different databases can help organize your data. In the case of the Climate Research Institute, we might want to group data by year or by region, for example.



## Create a new database

- CREATE DATABASE [options] <database name> [options];

```
hive> CREATE DATABASE europe;
OK
Time taken: 0.198 seconds
hive> CREATE DATABASE africa COMMENT
'Contains data from Africa';
OK
Time taken: 0.317 seconds
hive> CREATE DATABASE antarctica
LOCATION '/user/user01/hive/antarctica.db';
OK
Time taken: 0.034 seconds
```

© 2015 MapR Technologies

5

This example shows the syntax for creating new databases from the Hive Command Line Interface. For our Climate Research Center, we want to organize our data logically into several distinct databases. We have decided to make one database for each continent. Here, you can see the creation of the Africa, Antarctica, and Europe databases.

To create a database, at minimum, you must use the CREATE DATABASE command, provide a unique name for your database, and end your statement in a semicolon. We did this to create the Europe database.

There are several options you can add to the CREATE DATABASE command.

For example, you may add a comment about your database. This comment must be inside single quotes, and occurs after the keyword COMMENT following your database name. In this example, we added a comment to the "Africa" database.

You may also specify a location for your database to be stored. To do this, use LOCATION followed by the file path . Here we specified a location for the Antarctica database.



## SHOW and USE databases

- SHOW DATABASES;
- USE <database>;

```
hive> SHOW DATABASES;
OK
africa
antarctica
default
europe
Time taken: 0.024 seconds, Fetched: 4 row(s)
hive> USE antarctica;
OK
Time taken: 0.025 seconds
```

© 2015 MapR Technologies 

6

Once you create databases in your cluster, use SHOW DATABASES to list them. Your databases will be listed in alphabetical order. As you can see, we now have the africa, antarctica, and europe databases created in our cluster, in addition to the default database.

To change the working database, use the USE keyword. In this example, we set Antarctica as our working database. For the remainder of the course, examples will be created inside the Antarctica database, unless otherwise specified.



Now you can create your first Hive database. Remember, each database on your cluster needs a unique name. If you are in a classroom environment where several students are sharing the same cluster, it is important to create a unique database name. Each student should create a database associated with their own username; if your assigned username is user02, for example, then your database will be named user02.

Open your Lab Guide, and complete Lab 2.1. When you are finished, you can explore the Hive command line interface by trying familiar SQL commands, or you can continue with the lecture.



## Create a simple table

- CREATE TABLE <table name> (columnName TYPE) [options];
- USE and dot notation

```
hive> CREATE TABLE south_pole (date DATE, celsius FLOAT);
OK
Time taken: 0.381 seconds
hive> CREATE TABLE europe.london (date DATE, celsius FLOAT);
OK
Time taken: 0.114 seconds
```

© 2015 MapR Technologies 

8

To create a table, use the CREATE TABLE keywords, followed by the name of the table. Like creating a database, Hive will throw an error if the table name already exists.

After the table name, list each column name and its data type in parentheses. We'll discuss data types in Hive later.

As with all Hive commands, end the data definition with a semicolon. By convention, keywords like CREATE and FLOAT are written in allcaps; however, like SQL, Hive is case-insensitive.

Since we already specified we're using the Antarctica database with the USE command, the "south pole" table in this example will be created in the Antarctica database. If we want to add a table called "London" to the "Europe" database, we can either switch databases with the USE command again, or we can use dot notation. By typing "europe.london", we specify the London table is in the Europe database.

To ensure your tables and queries are always made in the correct database, it is a good practice to always use dot notation.



SHOW TABLES

- SHOW TABLES [IN <database>];

```
hive> SHOW TABLES;
OK
south_pole
Time taken: 0.079 seconds, Fetched: 1 row(s)
hive> SHOW TABLES IN europe;
OK
london
Time taken: 0.028 seconds, Fetched: 1 row(s)
```

We can see which tables we have already created by using SHOW TABLES.

Since we set our working database to antarctica, using SHOW TABLES right now will only show the south\_pole table we just created.

If we want to see the london table, we can either set the working database to europe, or we can use SHOW TABLES IN europe; as shown in this example.

SHOW TABLES is a great way to verify you are working in the correct database, and that your tables were created correctly.



## Delimiters & file types

- [FIELDS | ROWS] TERMINATED BY <delimiter>
  - Commas, spaces
  - Tabs '\t'
  - New lines '\n'
- STORED AS <fileformat>
  - TEXTFILE
  - SEQUENCEFILE
  - RCFILE
  - ORC
  - PARQUET

© 2015 MapR Technologies 

10

Table delimiters might be commas, spaces, tabs, new lines, or a variety of other characters such as hashes, dollar signs, or semicolons.

You can also specify how your file is stored. The most common file types are textfile and parquet files.

One common file format is a textfile, which includes .TXT and .CSV.

A sequence file is the native file format of the Hadoop Distributed File System. Sequence files are often split into chunks across mappers.

Record columnar files, or RC files, are a flipped matrix. Columns are stored as rows, and rows are stored as columns. This makes tall narrow data into flat wide data. Depending on the nature of your data, using this file format may help speed processing times.

Ordered record columnar files, abbreviated as ORC, have the same features as RC files, with the addition of headers and footers for metadata and indexing. The ordered nature of ORC files means Hive can sometimes skip reads, again leading to an increase in speed.

Finally, parquet files offer a compressed, efficient columnar file format which works with a variety of tools in the Hadoop Ecosystem, including Hive, Pig, Drill, and Spark. You can read the Apache Parquet documentation for more information about this file format.



## DESCRIBE and COMMENT

- Add the COMMENT option
- DESCRIBE <table>;

```
hive> CREATE TABLE africa.cairo
      > (date DATE COMMENT 'Date recorded',
      > celsius FLOAT COMMENT'Degrees Celsius');
OK
Time taken: 0.056 seconds
hive> DESCRIBE africa.cairo;
OK
date                  date          Date recorded
celsius               float         Degrees Celsius
Time taken: 0.261 seconds, Fetched: 2 row(s)
```

© 2015 MapR Technologies

11

Let's create another table, Cairo, in the Africa database. This time, we'll add some comments to our columns so we know what they are.

For ease of readability, we have divided this table creation into several lines. Remember, Hive won't run your query or command until you enter a semicolon to denote your query or command is complete.

Now that we've created this new table, let's examine it with Describe. Describe will list all the column names, the data type of the column, and any comments we added about the column. In this example we see that the column called "date" contains a data type of "date" and the note that it is the date recorded.

Like SHOW TABLES, Describe is another great way to verify your tables were created successfully.



```
SHOW CREATE <table>;
```

```
hive> SHOW CREATE TABLE europe.london;
OK
CREATE TABLE `europe.london`(
  `date` date,
  `celsuis` float)
ROW FORMAT SERDE
  'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe'
STORED AS INPUTFORMAT
  'org.apache.hadoop.mapred.TextInputFormat'
OUTPUTFORMAT
  'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat'
LOCATION
  'maprfs:/user/hive/warehouse/europe.db/london'
TBLPROPERTIES (
  'transient_lastDdlTime'='1444171998')
```

© 2015 MapR Technologies

12

In addition to DESCRIBE and SHOW TABLES, you can also use SHOW CREATE <table> to verify your data was created correctly.

SHOW CREATE TABLE provides details from the data definition of a table.



Now is a great time to create your first Hive table. Open your Lab Guide, and complete Lab 2.2. When you are finished, you can continue to explore the Hive command line by trying familiar SQL commands on your own, or you can continue with the lectures and labs.



## Knowledge Check

When can you use dot notation instead of the `USE` command?

- a) Whenever you create a new table, regardless of the current database setting
- b) Whenever you create a table outside the current database
- c) Whenever you create a table outside the default database

Answer: A (technically all of the above)

To ensure you always create your tables in the correct databases, you can always use the dot notation and not worry about using the `USE` command.



## Learning Goals

- ▶ Create databases and tables
- ▶ Partition and bucket data
- ▶ Load tables with data
- ▶ Alter and drop tables

Next, we'll create external tables and partitioned tables.



## Create an external table

- Link table directly to a data file
- Specify features of the file, including LOCATION

```
hive> CREATE EXTERNAL TABLE europe.madrid
    > (date DATE, celsius FLOAT)
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > LINES TERMINATED BY '\n'
    > STORED AS TEXTFILE
    > LOCATION '/user/user01/madrid.csv';
OK
Time taken: 0.071 seconds
```

© 2015 MapR Technologies

16

External tables are linked directly to an external data file. This has several advantages. First, you load the data and create the table simultaneously in the same data definition. Second, external tables are read-only. Read-only is a good option if you do not want your users to accidentally edit or delete data. Users can still create and save views from an external table, but the external table will remain unchanged, as it is linked to an external file.

When creating a table, you can specify many features such as the data type of each column, the row format, delimiters, how the file is stored, and its location. In this example, we have specified a row format which is delimited by commas, with lines terminated by a new line character, and is stored as a text file. These settings are common for CSV files.

When you are creating an external table, you must specify a location on your database. This will link your table to data, such as a csv file. In this example, we are linking this table to a file called "madrid.csv".



## Knowledge Check

Why is it important to know how your data is stored?

- a) Hive requires you to specify the file format
- b) Hive requires you to specify the data type
- c) Some file formats are optimized for faster querying in Hive
- d) Some files format or data types cannot be read by Hive
- e) All of the above

Answer: all of the above

INSTRUCTORS: Take some time now to discuss with your students which file formats they plan to use.



## Create a partitioned table

- Partitioning divides data into logically organized chunks

```
hive> CREATE TABLE antarctica.windspeed
  > (year INT, month STRING, knots FLOAT)
  > PARTITIONED BY (station STRING)
  > ROW FORMAT DELIMITED
  > FIELDS TERMINATED BY ','
  > LINES TERMINATED BY '\n'
  > STORED AS TEXTFILE;
OK
Time taken: 0.5 seconds
```

© 2015 MapR Technologies

18

In partitioned tables, the data is divided into smaller chunks. Each chunk of rows will correspond to a partition, which we define logically. These chunks can reduce the number rows Hive must read, which may result in faster processing times.

In this example, we have a large dataset which includes wind speed data, in knots, across several different weather stations. We can partition the data by the station, so that each station can be queried separately without querying the whole table.

When we load this data later, we will load data corresponding to each station separately.



## Create a bucketed table

- Bucketing divides data into evenly-sized chunks by hashing

```
hive> CREATE TABLE antarctica.faraday
      > (year INT, month STRING, knots FLOAT)
      > CLUSTERED BY (year) INTO 10 BUCKETS;
OK
```

© 2015 MapR Technologies

19

Like partitioning, bucketing also divides data into smaller chunks. These chunks also reduce the number of reads Hive must make, which may result in faster processing times. Unlike partitions, buckets are not defined by us. Instead, buckets are sampled from the data by a hashing function.

In this example, we bucket the data from the Faraday weather station into 10 buckets. The data will be bucketed by the year. The hash function will create 10 evenly sized samples of the data, and organize each bucket by the year column.

When we load this data later, the data is loaded all at once. The bucketing occurs automatically.



## Discussion – Partitioning Data

Think about your own data. How might you logically partition it? By date? By categories? How many buckets would be appropriate for various sizes of data?

© 2015 MapR Technologies MAPR

20

INSTRUCTORS: Take some time to discuss this with your students.



Now you will create external and partitioned tables Hive. Open your Lab Guide, and Lab 2.3. When you are finished, you can continue to explore the Hive command line by trying familiar SQL commands on your own, or you can continue with the lectures and labs.



## Learning Goals

- ▶ Create databases and tables
- ▶ Partition and bucket data
- ▶ Load tables with data
- ▶ Alter and drop tables

© 2015 MapR Technologies MAPR

22

In this next section, we'll learn how to load data into internal and partitioned tables.



## Load data into tables

- Use LOAD DATA to populate a table with data

```
hive> LOAD DATA LOCAL INPATH
      > '/user/user01/southpole.csv'
      > INTO TABLE south_pole;
Loading data to table antarctica.south_pole
Table antarctica.south_pole stats: [numFiles=1, totalSize=9229]
OK
Time taken: 0.374 seconds
```

© 2015 MapR Technologies 

23

The LOAD DATA command inserts data from a file into a table. LOAD DATA only works if the table is already defined using the data definition language.

The LOAD DATA command can load data from a local file path. Specify the file path with LOCAL INPATH, then specify the name of the table with INTO TABLE.

This table will be created in the current working database, which is Antarctica. It is good practice to use dot notation to specify the database.



## Load data into partitioned tables

- Use PARTITION to load parts of data into a partitioned table

```
hive> LOAD DATA LOCAL INPATH '/user/user01/wind_Adelaide.csv'  
    > INTO TABLE antarctica.windspeed PARTITION(station = 'Adelaide');  
Loading data to table antarctica.windspeed partition (station=Adelaide)  
OK  
Time taken: 1.124 seconds  
hive> LOAD DATA LOCAL INPATH '/user/user01/wind_Rothera.csv'  
    > INTO TABLE antarctica.windspeed PARTITION(station = 'Rothera');  
Loading data to table antarctica.windspeed partition (station=Rothera)  
OK  
Time taken: 0.86 seconds
```

© 2015 MapR Technologies

24

Since our windspeed table is partitioned by the station name, we need to load each partition separately. We specify which station corresponds to each partition when we load it.

The data in wind\_adelaide.csv does not have a station column, but when we load the data like this, Hive populates the table with a date column and inserts “Adelaide” for all these data points. This example demonstrates loading two such partitions into the windspeed table.



Now you can try to load data into the tables you have created in Hive. Open your Lab Guide, and complete Lab 2.4. When you are finished, you can continue to explore the Hive command line by trying familiar SQL commands on your own, or you can continue with the lectures and labs.



## Learning Goals

- ▶ Create databases and simple tables
- ▶ Create external and partitioned tables
- ▶ Load tables with data
- ▶ Alter and drop tables

© 2015 MapR Technologies MAPR

26

Finally, we'll learn about altering and dropping the tables that we have created.



## Drop tables and databases

- use `DROP TABLE;` to delete a table
- use `DROP DATABASE;` to delete an empty database

```
hive> DROP TABLE south_pole;
OK
Time taken: 0.198 seconds
hive> DROP TABLE africa.cairo;
OK
Time taken: 0.212 seconds
hive> DROP DATABASE africa;
OK
Time taken: 0.199 seconds
```

© 2015 MapR Technologies

27

If you have created a table by mistake, or you just no longer need it, you can delete it using `DROP`.

Note that dropping a table does not necessarily delete your data files. Your data is typically stored on your HDFS, and will remain there unaltered. This is true even if you delete an external table.

You can also delete a database using `DROP DATABASE`, but only if there are no tables inside that database.



Re-create the `south_pole` table

```
hive> CREATE TABLE south_pole (
    > station STRING, year INT, month STRING,
    > celsius FLOAT, fahrenheit FLOAT)
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > LINES TERMINATED BY '\n'
    > STORED AS TEXTFILE;
OK
Time taken: 0.085 seconds
```

© 2015 MapR Technologies

28

Now that we have dropped the `south_pole` table, we will need to recreate it if we want to use it again. Let's say we dropped it because it was created with the wrong data definition. Let's now re-create it, this time with different column names and data types.



## Alter tables

- ALTER TABLE can edit tables and columns

```
hive> ALTER TABLE europe.london RENAME TO europe.dublin;
OK
Time taken: 0.228 seconds
hive> ALTER TABLE europe.dublin CHANGE COLUMN celsius fahrenheit FLOAT;
OK
Time taken: 0.264 seconds
```

© 2015 MapR Technologies

29

What if you have created a table, but then realize you need to make a change to it? Rather than dropping the table and recreating it entirely, you can use ALTER TABLE. This is good for smaller changes, such as renaming a column or recasting a data type.

For example, let's say you realized that the london table actually contains data from dublin. You can rename this table.

You then realize that the celsius column should actually be fahrenheit. Remember to specify the data type when altering columns!



## Discussion – Altering and Dropping Tables

When would you alter a table and when would you drop a table?

© 2015 MapR Technologies MAPR

30

INSTRUCTORS: Discuss some of the different use cases of ALTER and DROP.



So far in these labs, you have created a database, created some tables, and loaded them with some data.

Use the remainder of this lab time to briefly explore these tables by following the directions Lab 2.5.

If you are already familiar with SQL, you might try querying this data using common SQL commands.

When doing the labs, it is recommended that you have a text editor open in order to write queries, rather than writing queries directly in the command line.

When you have finished Lesson 2 in the Lab Guide, you are ready to move on to Lesson 3.



**Next Steps**

**Lesson 3**  
Query Data in  
Apache Hive

© 2015 MapR Technologies **MAPR** 32

In the next lesson, we'll learn about how to query and manipulate data in Hive.



## **DA 440: Apache Hive Essentials**

Lesson 3: Query Data in Apache Hive

© 2015 MapR Technologies **MAPR**

1

Welcome DA 440 Lesson 3 – Query Data in Apache Hive.



## Learning Goals

- ▶ Query tables
- ▶ Manipulate tables
- ▶ Combine & store tables

© 2015 MapR Technologies MAPR

2

By the end of this lesson, you will know how to query data. You will also learn to manipulate data with user defined functions. Finally, we will combine and save our data.



## Learning Goals

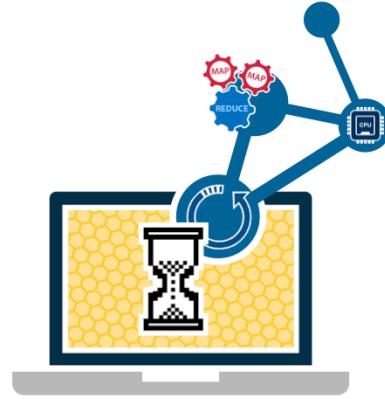
- ▶ Query tables
- ▶ Manipulate tables
- ▶ Combine & store tables

Let's start with some basic queries.



## Hive and MapReduce

- Many Hive queries trigger MapReduce
- MapReduce
  - allows Hive to query data on HDFS
  - uses more resources than SQL or Drill



© 2015 MapR Technologies

4

Before we begin querying data in Hive, recall how Hive differs from SQL.

Unlike SQL, Hive works on a Hadoop Distributed File System by converting SQL-like code into MapReduce code. This triggers a MapReduce job.

Since MapReduce jobs may be resource-intensive, it is important to try to optimize your queries using some of the techniques discussed in Lesson 1.

Users who are used to the speed of SQL sometimes worry their Hive jobs will not complete. Depending on the speed of your system, you may need to be patient as you wait for the MapReduce job to complete.



## The SELECT statement

- SELECT, FROM, and WHERE are similar in Hive and SQL

```
hive> SELECT * FROM south_pole WHERE celsius < -65;
OK
Clean_Air      1995    sep     -65.5   -85.9
Clean_Air      1997    jul     -66.3   -87.340004
Clean_Air      2004    jul     -65.8   -86.44
Clean_Air      1993    aug     -66.6   -87.88
Time taken: 0.137 seconds, Fetched: 4 row(s)
```

© 2015 MapR Technologies

5

The SELECT statement in Hive works similarly to the SELECT statement in SQL, which you may already be familiar with.

For example, let's say we wanted to see all the data from the south pole table where the temperature is less than negative 65 degrees celsius. We can use the asterisk to select all the rows, then specify we are selecting FROM south\_pole, then restrict the statement with the WHERE clause.



count (\*)

```
hive> SELECT count(*) FROM south_pole WHERE celsius < -65;
Query ID = user01_20150828154444_dd93bd89-528c-40de-b77b-a312d81964d8
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1440796388972_0002, Tracking URL = http://maprdemo:8088/
proxy/application_1440796388972_0002/
Kill Command = /opt/mapr/hadoop/hadoop-2.7.0/bin/hadoop job -kill job_1440
796388972_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducer
s: 1
2015-08-28 15:45:11,567 Stage-1 map = 0%,  reduce = 0%
2015-08-28 15:45:23,807 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 2.
79 sec
2015-08-28 15:45:34,775 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU
4.84 sec
MapReduce Total cumulative CPU time: 4 seconds 840 msec
Ended Job = job_1440796388972_0002
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 4.84 sec MAPRFS Read:
0 MAPRFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 4 seconds 840 msec
OK
4
```

- **SELECT count(\*) runs a MapReduce job**
  - **SELECT \*** does not use MapReduce
  - **count(\*)** is a UDF

© 2015 MapR Technologies

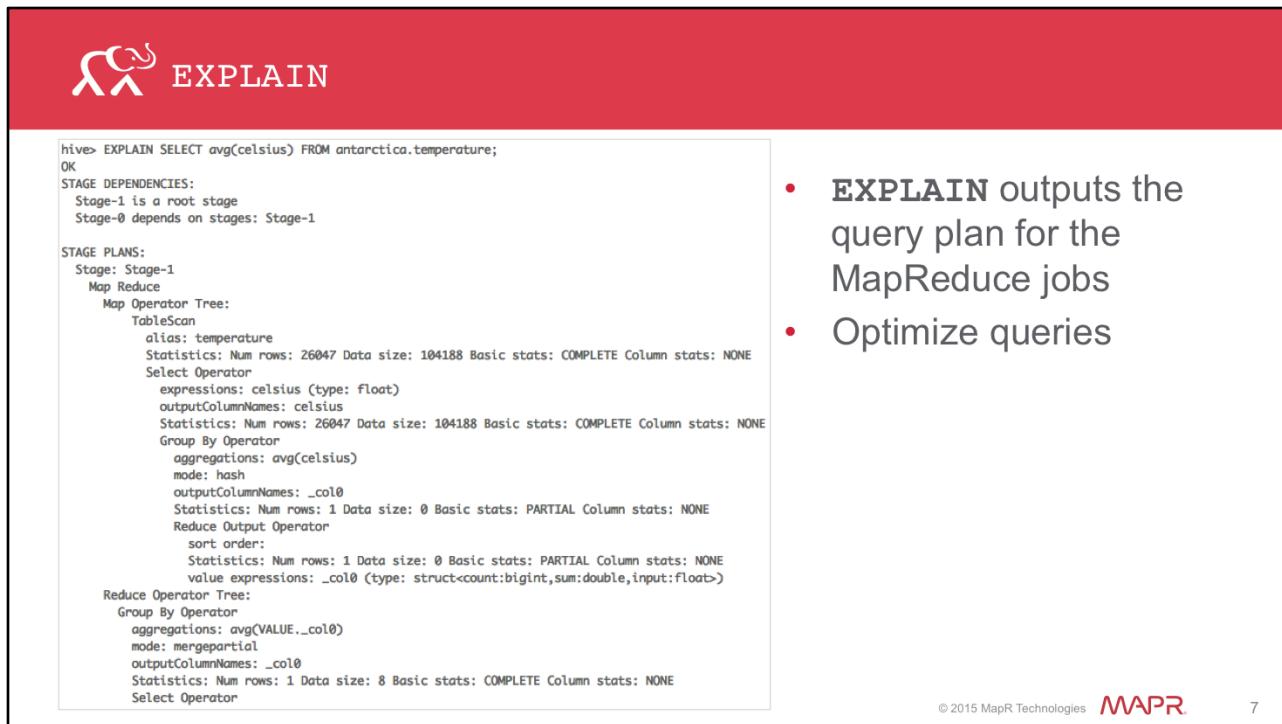
6

What if you wanted to know how many rows of data match your query, rather than printing all the rows? You might be familiar with **SELECT count(\*)** from SQL

Unlike **select \*** from the previous example, **select count(\*)** triggers a MapReduce job. Here you can see an example of what that looks like in the Hive Command Line Interface.

The Hive CLI will start the job, display options for modifying, tracking, and killing the job, and update you on the progress of the map and reduce phases.

This example shows the output of a MapReduce job. In this example, the MapReduce job took 580 milliseconds, then returned a result of 172. So, there are 172 rows of data where the Celsius column contains a value that is less than negative 65 in our South Pole dataset.



The screenshot shows a terminal window with the title "EXPLAIN" in red. The terminal output is as follows:

```
hive> EXPLAIN SELECT avg(celsius) FROM antarctica.temperature;
OK
STAGE DEPENDENCIES:
 Stage-1 is a root stage
 Stage-0 depends on stages: Stage-1

STAGE PLANS:
 Stage: Stage-1
 Map Reduce
  Map Operator Tree:
   TableScan
    alias: temperature
    Statistics: Num rows: 26047 Data size: 104188 Basic stats: COMPLETE Column stats: NONE
   Select Operator
    expressions: celsius (type: float)
    outputColumnNames: celsius
    Statistics: Num rows: 26047 Data size: 104188 Basic stats: COMPLETE Column stats: NONE
   Group By Operator
    aggregations: avg(celsius)
    mode: hash
    outputColumnNames: _col0
    Statistics: Num rows: 1 Data size: 0 Basic stats: PARTIAL Column stats: NONE
   Reduce Output Operator
    sort order:
     Statistics: Num rows: 1 Data size: 0 Basic stats: PARTIAL Column stats: NONE
    value expressions: _col0 (type: struct<count:bigint,sum:double,input:float>)
  Reduce Operator Tree:
   Group By Operator
    aggregations: avg(VALUE._col0)
    mode: mergepartial
    outputColumnNames: _col0
    Statistics: Num rows: 1 Data size: 8 Basic stats: COMPLETE Column stats: NONE
   Select Operator
```

© 2015 MapR Technologies  7

Using EXPLAIN on a Hive query will output the plan for converting HiveQL into a MapReduce job.

Advanced users of Hive may wish to see how their queries get converted into MapReduce code. This can help them identify inefficiencies in the query, and optimize performance.

This screenshot shows part of the output of EXPLAIN on a query.



Now you can try your first queries. Complete Lab 3.1 in your Lab Guide before moving on.

 LIMIT

- LIMIT only displays the first N rows

```
hive> SELECT * FROM south_pole LIMIT 3;
OK
Clean_Air      1986    sep     NULL     NULL
Clean_Air      1987    sep     NULL     NULL
Clean_Air      1988    sep     -58.2   -72.76
Time taken: 0.054 seconds, Fetched: 3 row(s)
```

© 2015 MapR Technologies 

9

Tables in Hive often have a lot of data, since HDFS is designed to handle large quantities of data. If you use SELECT \*, Hive will return the entire dataset, which may be millions or billions of rows. You can use LIMIT to return a smaller subset of rows. You can choose how many rows to display. This is especially useful if you are testing to see if a dataset loaded properly; the first 5 or 10 rows is usually sufficient for this task. In this example, we limited the display to 3 rows.



## DISTINCT

- DISTINCT displays unique values

```
hive> SELECT DISTINCT month FROM south_pole;  
OK  
apr  
aug  
...  
sep  
Time taken: 43.742 seconds, Fetched: 12 row(s)
```

You can also use SELECT DISTINCT to see a specific subset of the data. For example, if you wanted to know which months you have loaded into your south\_pole table, you could use SELECT DISTINCT month FROM south\_pole;



AND, OR, NOT

- Other operators: AND, OR, IS NOT NULL

```
hive> SELECT * FROM south_pole WHERE celsius < -65 AND year > 2000;  
Clean_Air    2004    jul    -65.8    -86.44  
hive> SELECT * FROM south_pole WHERE celsius < -65 OR celsius > -25;  
Clean_Air    1995    sep    -65.5    -85.9  
Clean_Air    1997    jul    -66.3    -87.340004  
Clean_Air    2004    jul    -65.8    -86.44  
Clean_Air    1999    jan    -24.9    -12.82  
Clean_Air    1993    aug    -66.6    -87.88  
Clean_Air    2004    dec    -21.9    -7.419999  
hive> SELECT count(*) FROM south_pole WHERE celsius IS NOT NULL;  
202
```

© 2015 MapR Technologies 

11

Other querying operators might be familiar from SQL.

For example, you might want to select all the data where the temperature is greater than negative 20 degrees and the year is more recent than 2000.

Alternatively, perhaps you only want extreme weather, regardless of whether it's hot or cold. You might select all the data where celsius is greater than thirty OR celsius is less than zero.

You can also use IS NOT NULL in a where clause. For example, if your dataset has many missing values, you could query all the data where the entry IS NOT NULL



## Sort data

- ORDER BY, SORT BY
- ASC, DESC

```
hive> SELECT * FROM south_pole WHERE celsius < -65 SORT BY celsius DESC;
OK
Clean_Air      1995    sep     -65.5   -85.9
Clean_Air      2004    jul     -65.8   -86.44
Clean_Air      1997    jul     -66.3   -87.340004
Clean_Air      1993    aug     -66.6   -87.88
Time taken: 42.777 seconds, Fetched: 4 row(s)
```

© 2015 MapR Technologies

12

ORDER BY performs a total ordering of the data. You can optionally add the parameters ASC or DESC to sort in ascending or descending order. By default, all sorts are in ascending order.

ORDER BY sorts all of the data. This can result in extremely long run times.  
SORT BY works on smaller chunks of the data, and therefore may be better for faster results if you do not need the entire table sorted.

In this example, we selected everything from the south pole table where the celsius column was less than negative 65 degrees, and then sorted it in descending order.



## Knowledge Check

Which query would return the number of rows with a null value in the knots column in the South Pole table?

- a) SELECT \* FROM south\_pole IF NULL;
- b) SELECT \* FROM south\_pole WHERE knots IS NULL;
- c) SELECT count(\*) FROM south\_pole IF NULL;
- d) SELECT count(\*) FROM south\_pole WHERE knots IS NULL;

© 2015 MapR Technologies

13

Answer: D

**INSTRUCTOR:** Ask your students what sort of results each of these queries would return. Alternatively, rather than provide these options, consider collaboratively writing this (and/or other queries) with your students. Discuss a variety of queries you might want to write and try to write them together as a class.

**DISCUSS:** When have you used DISTINCT, NOT, OR, AND, etc.? What are the use cases of these operators in your own experience as a data scientist? How are these similar or different from what you can do in SQL, Hive, or other languages or programs like Drill?

**DISCUSS:** Recall the version chart from Lesson 1. Are there issues with older versions and queries you might want to use?



## Learning Goals

- ▶ Query data
- ▶ Manipulate data
- ▶ Combine & store data

© 2015 MapR Technologies MAPR

14

In this next section, we'll learn about a few common user defined functions.



## Common mathematical operations

- Use UDFs and UDAFs to preform mathematical operations
- UDFs: `abs()`, `round()`, `floor()`, `ceil()`,  
`sin()`, `cos()`, `tan()`, `log()`
- UDAFs: `min()`, `max()`, `avg()`

```
hive> SELECT round(celsius*1.8+32) FROM south_pole;
OK
-73.0
-75.0
-66.0
-75.0
```

© 2015 MapR Technologies 

15

You can use arithmetic operators when querying data. For example, the south pole table only has data in the Celsius column. If you'd like to see this data in degrees Fahrenheit, you can convert it using multiplication and addition, as shown in this example here. We also use the `round()` UDF.

In addition to arithmetic operators, you can use user defined functions and user defined aggregate functions to preform many other common mathematical operations on numerical data.

INSTRUCTORS: Now is a great time to demo some of these functions. Try these queries:

```
SELECT min(celsius), max(celsius) FROM south_pole;
SELECT avg(celsius) FROM south_pole WHERE year > 2000;
SELECT round(fahrenheit) FROM south_pole LIMIT 10;
```



## Common string operations

- Manipulate string data with UDFs
- `lower(s)`, `upper(s)`, `length(s)`, `reverse(s)`,  
`to_date(timestamp)`, `cast(s AS <type>)`,  
`substr(s, pattern)`

```
hive> SELECT upper(station) FROM south_pole LIMIT 3;
OK
CLEAN_AIR
CLEAN_AIR
CLEAN_AIR
Time taken: 0.06 seconds, Fetched: 3 row(s)
```

© 2015 MapR Technologies 

16

In addition to numerical operations, Hive supports many string operations. The `upper()` function is demonstrated in the screenshot shown. This function takes the station name string, which was saved in mixed-case, and converts everything to all caps. The `lower()` function would do the opposite.

A short list of some common string UDFs are shown here. However, this list is not exhaustive, and you should consult your Hive documentation for a complete list.



## UDAFs in subqueries

- Most versions of Hive do not support UDAFs in subqueries



```
hive> SELECT * FROM south_pole WHERE celsius > avg(celsius) AND year > 2000;  
FAILED: SemanticException [Error 10128]: Line 1:41 Not yet supported place for UDAF 'avg'
```



```
hive> SELECT avg(celsius) FROM south_pole WHERE year > 2000;  
OK  
-48.65000004238553  
Time taken: 47.028 seconds, Fetched: 1 row(s)  
hive>  
hive> SELECT * FROM south_pole WHERE celsius > -48 AND year > 2000;  
OK  
Clean_Air    2001    feb    -38.3    -36.94  
Clean_Air    2002    feb    -40.6    -41.079998  
Clean_Air    2003    feb    -37.5    -35.5  
Clean_Air    2004    feb    -35.9    -32.620003
```

© 2015 MapR Technologies

17

It is important to remember that Hive is not ANSI-SQL compliant. Although the HQL queries you have seen so far are similar to SQL, there are some important differences.

For example, currently Hive does not support UDFs inside sub queries, although it is possible that future versions of Hive may support this feature.

This means that if you wanted to select all the data from your south pole weather station where the temperature was greater than average, as in the first query shown, you could not.

If you wanted to perform a similar query, you would first have to determine what the average temperature in the dataset was, and then write a query selecting data which is greater than that number. In this example, we find that the average is about -48. We can then query the data again, based on this number.



Now you can try more queries. Complete Lab 3.2 in your Lab Guide before moving on. Feel free to explore some of the other user defined functions we have discussed.



## Knowledge Check

Write a query to find the dates when record low temperatures occurred in the London table.

- a) `SELECT date FROM europe.london WHERE celsius = min(celsius);`
- b) `SELECT min(celsius) FROM europe.london;` -13.6  
`SELECT date FROM europe.london WHERE celsius = -13.6;`
- c) `X = SELECT min(celsius) FROM europe.london;`  
`SELECT date FROM europe.london WHERE celsius = X;`

© 2015 MapR Technologies 

19

Answer: B

INSTRUCTOR: Rather than provide these options, collaboratively write this (and/or other queries) with your students. Discuss a variety of queries you might want to write and try to write them together as a class.



## Learning Goals

- ▶ Query data
- ▶ Manipulate data
- ▶ Combine & store data

© 2015 MapR Technologies MAPR

20

In this last section, we'll learn how to combine datasets and save them as new tables or files



UNION ALL

- UNION ALL concatenates two or more tables
- CREATE TABLE <table> AS SELECT saves a query as a new table

```
hive> CREATE TABLE windspeed AS SELECT
    > year, month, knots, 'Rothera' AS station FROM wind_rothera
    > UNION ALL SELECT
    > year, month, knots, 'Adelaide' AS STATION FROM wind_adelaide;
hive> SELECT * FROM windspeed LIMIT 3;
OK
1962    jul      18.8    Adelaide
1962    may      14.4    Adelaide
1962    oct      18.7    Adelaide
```

© 2015 MapR Technologies 

21

Union all allows you to concatenate two datasets. For example, we can create a new table, windspeed, using AS SELECT.

In this table, we will combine some data from several different weather stations.

In order to keep the data in this new table organized, we will add a new column called "station", which we will populate with a STRING so we know which station the data comes from. This new column did not exist in either original table.

the UNION ALL command then concatenates the data from Rothera with the data from Adelaide. In this example, a new column, station name, is created with a string we define.



- JOIN combines columns from two or more tables
- AS creates table aliases

```
hive> CREATE TABLE weather AS SELECT  
> t.station, t.year, t.month, t.celsius, w.knots  
> FROM temperature AS t  
> JOIN windspeed AS w  
> ON t.month = w.month  
> AND t.year = w.year  
> AND t.station = w.station;
```

© 2015 MapR Technologies

22

JOIN add columns of data from two or more tables by matching values which are common to each table.

The temperature and windspeed tables each contain data from multiple weather stations.

Let's create a new table called weather. We'll join these two tables on the date and station name. This way, both windspeed and temperature from the same station on the same month will be in the same row.

We'll abbreviate the names of our temperature table AS t and our windspeed table AS w. This abbreviation will make using dot notation easier.

We want to select the station, date, and celsius column from the temperature table, and add the knots column from the windspeed table.

When we join the tables, we have to specify which columns of data should match. We join ON the date column from both temperature and windspeed, AND we also join on the station name from each table.

With this new table, we can make some much more interesting queries. For example, we can find out what the windspeed was on the days with the record lows.



## Save your data from the command line

1. QUIT; exits the Hive Shell
2. Type `hive -e` followed by an HQL query in single quotes
3. The `>` operator exports the table into a CSV file

```
hive> QUIT;  
[user01@maprdemo ~]$ hive -e 'SELECT * FROM antarctica.weather' >  
/user/user01/weather.csv
```

© 2015 MapR Technologies 

23

After you have manipulated data, you might want to save it in a file. You can save your data as a file from the bash command line interface.

To use `hive` from the bash command line, exit the `hive` command line interface with `QUIT`. Then, type `hive dash e` and select everything from the `weather` table in the `antarctica` database. Place this `hive` query inside single quotes.

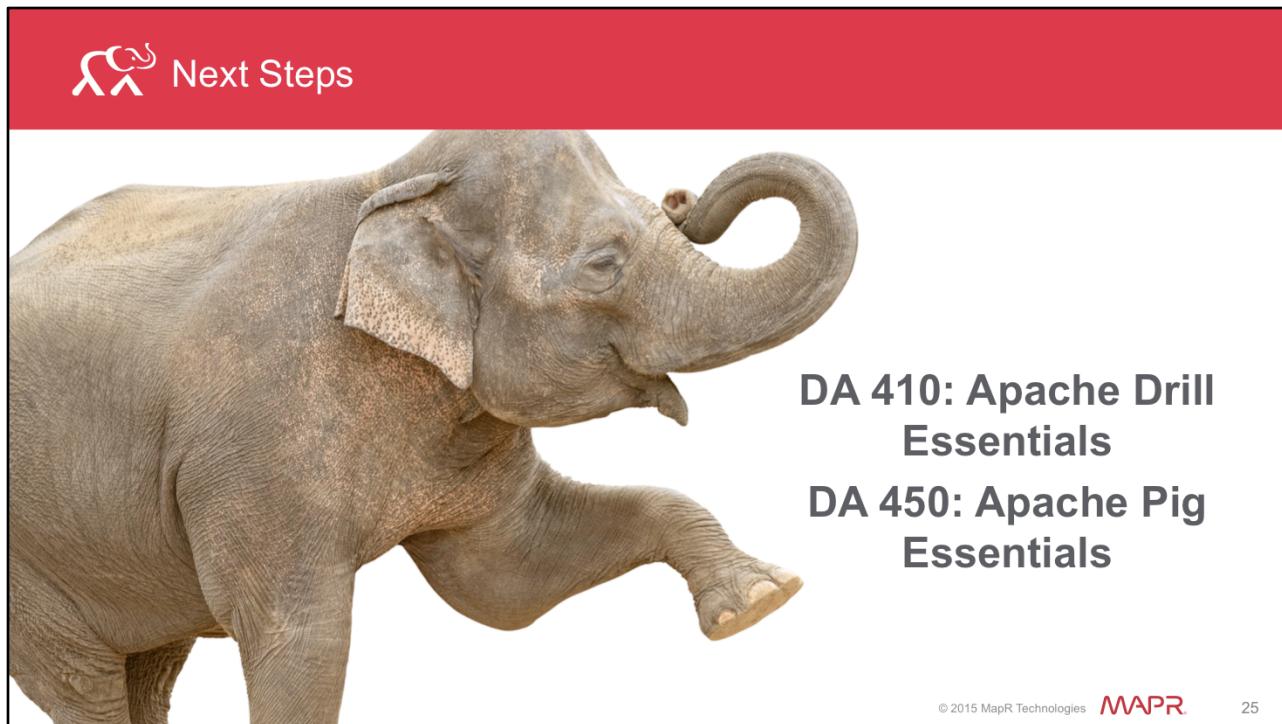
using the arrow operator, we'll save these queries to a new csv file called `weather.csv`. You can then verify your file was created by using `LS`, and view your new file directly from the bash shell using an editor such as `VI`.

You can now use this csv file to load a new table if you would like. This will save you from having to join or create tables from the datasets you loaded previously.



Complete Lab 3.3 in your Lab Guide before moving on. When you are finished, you may explore the weather data set as a data analyst. Be creative and try to come up with your own research questions and queries.

Take a moment to make sure you understand everything from the labs and lectures. Use this time to ask your instructor any clarifying questions.



The slide features a large Indian elephant walking towards the right, positioned on the left side of the frame. The elephant's trunk is raised and curled at the tip. The background is white, and the elephant is rendered in a realistic style.

**Next Steps**

**DA 410: Apache Drill  
Essentials**

**DA 450: Apache Pig  
Essentials**

© 2015 MapR Technologies **MAPR** 25

Thank you! You have now completed DA440: Apache Hive Essentials. Visit the MapR Academy website for more courses, including DA410: Drill Essentials or DA450: Pig Essentials