

Hive



走向分布式？

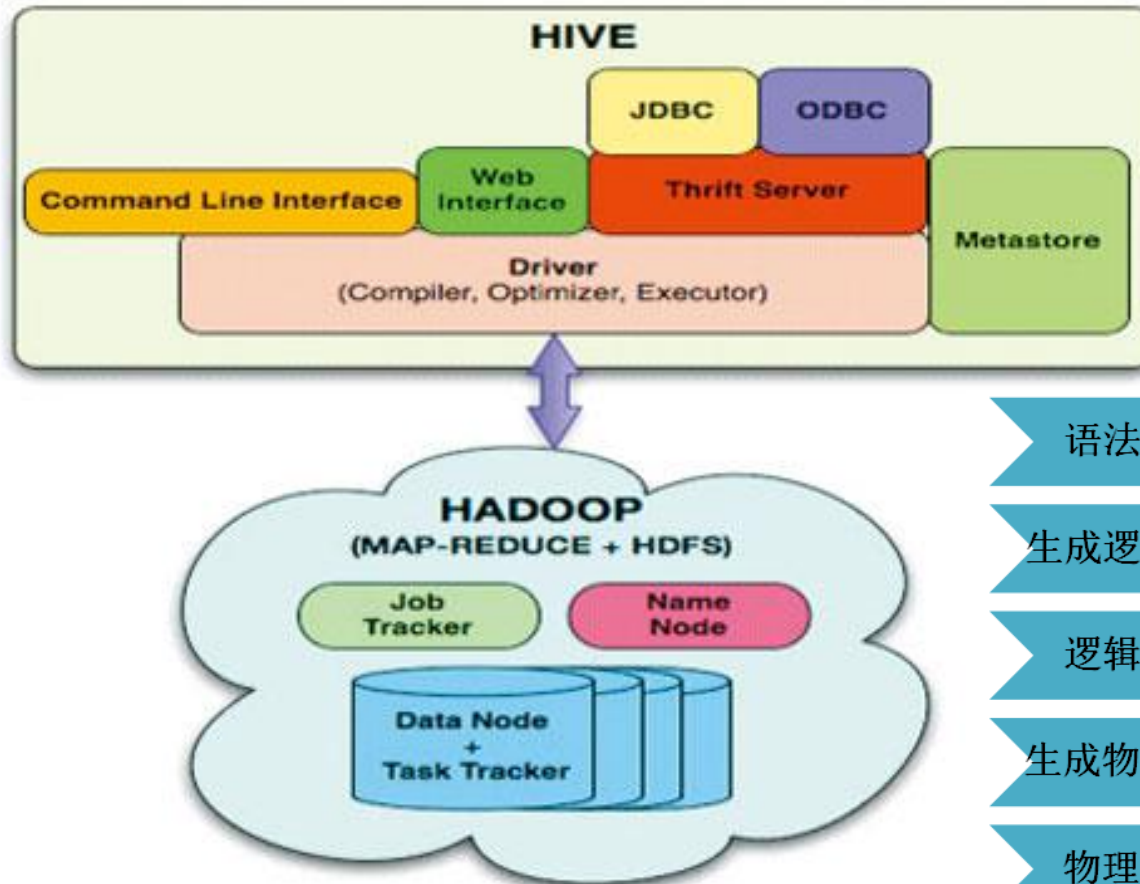
一个系统走向分布式，一定有其不得不为的理由。可扩展性是最常见的理由之一。我先简单的将“可伸缩”的需求分成两种：

- **Data Scalability:** 单台机器的容量不足以 (经济的) 承载所有资料，所以需要分散。如：NoSQL
- **Computing Scalability:** 单台机器的运算能力不足以 (经济的) 及时完成运算所以需要分散。如：科学运算。不管是哪一种需求，在决定采用分布式架构时，就几乎注定要接受一些牺牲：
 1. 牺牲效率：网路延迟与节点间的协调，都会降低执行效率。
 2. 牺牲 AP 弹性：有些在单机上能执行的运算，无法轻易在分布式环境中完成。
 3. 牺牲维护维运能力：分散式架构的问题常常很难重现，也很难追踪。另外，跟单机系统一样，也有一些系统设计上的 tradeoffs(权衡)
 4. CPU 使用效率优化或是 IO 效率优化
 5. 读取优化或是写入优化
 6. 吞吐率优化或是 网络延迟优化
 7. 资料一致性或是资料可得性,选择了不同的 tradeoff，就会有不同的系统架构。

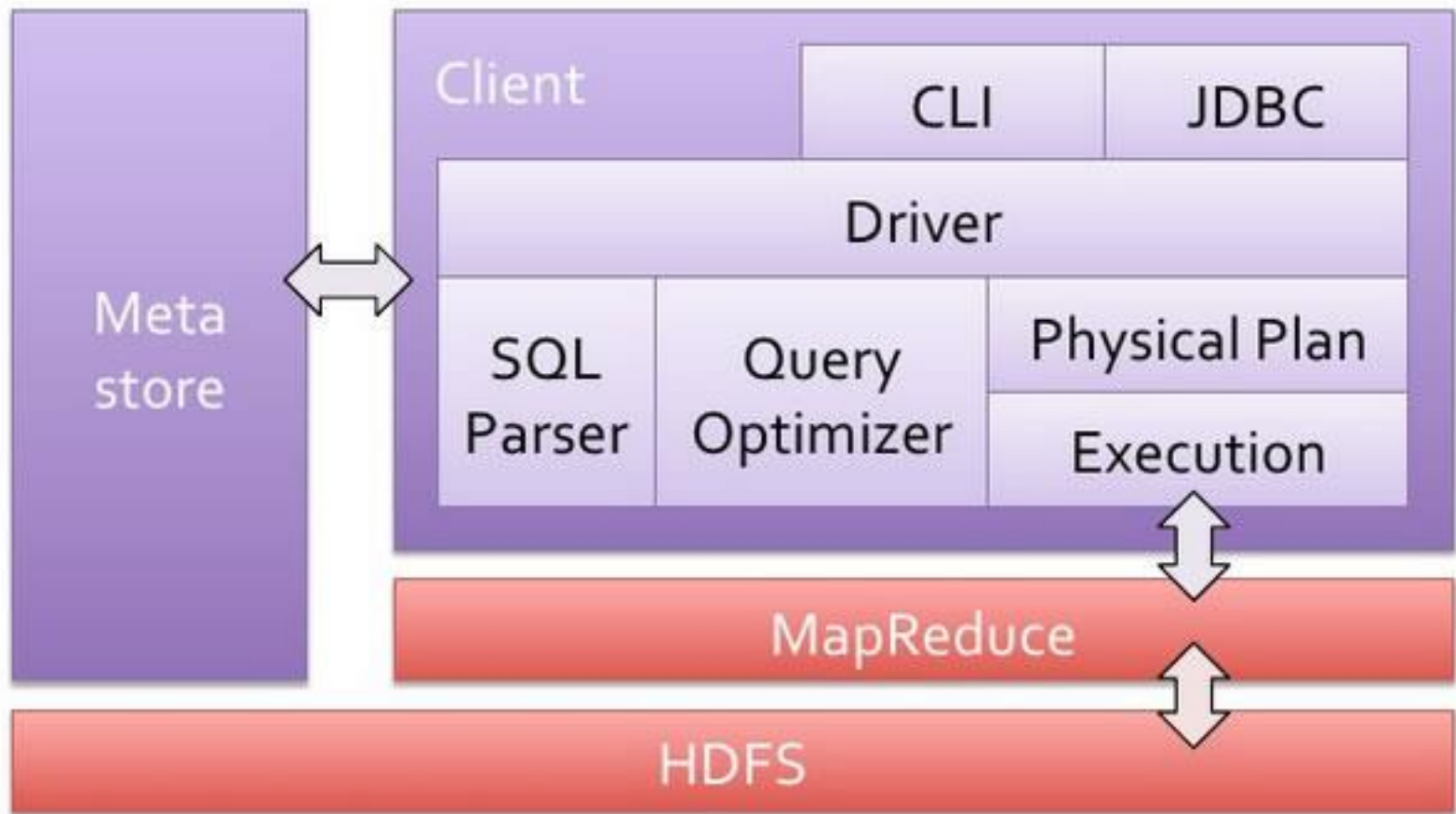
Hive特点

- **ETL**（ **Extraction-Transformation-Loading**）工具
- 构建在hadoop之上的数据仓库
- Hive定义了一种hql语句类似sql查询语句
- 常用于离线数据处理,
- HQL底层解析为MR程序执行

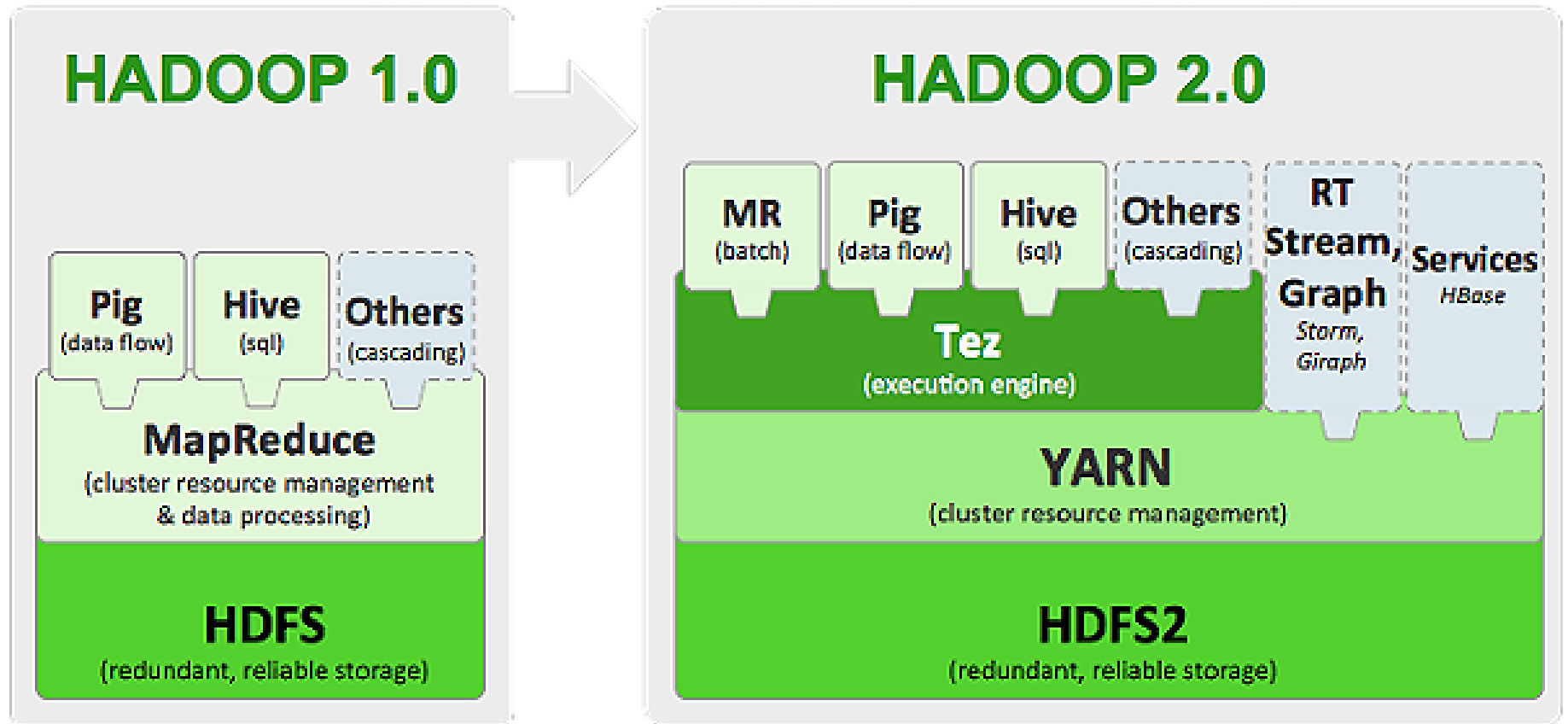
Hive



Hive 构成



Hive on Tez



性能

➤ Impala官网

参考: <http://blog.cloudera.com/blog/2014/09/new-benchmarks-for-sql-on-hadoop-impala-1-4-widens-the-performance-gap>
<http://blog.cloudera.com/blog/2014/05/new-sql-choices-in-the-apache-hadoop-ecosystem-why-impala-continues-to-lead/>

➤ HDP官网

参考: http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.0.0.2/bk_installing_manually_book/content/rpm-chap-tez-2.html
<http://zh.hortonworks.com/blog/evaluating-hive-with-tez-as-a-fast-query-engine/>

LanguageManual

- 1、Hive数据类型
- 2、Hive支持的文件格式
- 3、Hive CLI And ThriftServer/JDBC
- 4、DML/DDDL

参考:

Hive Type: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Types>

Hive Fileformats: <https://cwiki.apache.org/confluence/display/Hive/FileFormats>

Hive Cli: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Cli>

Commands line: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Commands>

LanguageManual: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual>

LanguageManual | DateType

	void	boolean	tinyint	smallint	int	bigint	float	double	decimal	string	varchar	timestamp	date	binary
void to	true	true	true	true	true	true	true	true	true	true	true	true	true	true
boolean to	false	true	false	false	false	false	false	false	false	false	false	false	false	false
tinyint to	false	false	true	true	true	true	true	true	true	true	true	false	false	false
smallint to	false	false	false	true	true	true	true	true	true	true	true	false	false	false
int to	false	false	false	false	true	true	true	true	true	true	true	false	false	false
bigint to	false	false	false	false	false	true	true	true	true	true	true	false	false	false
float to	false	false	false	false	false	false	true	true	true	true	true	false	false	false
double to	false	false	false	false	false	false	false	true	true	true	true	false	false	false
decimal to	false	false	false	false	false	false	false	false	true	true	true	false	false	false
string to	false	false	false	false	false	false	false	true	true	true	true	false	false	false
varchar to	false	false	false	false	false	false	false	true	true	true	true	false	false	false
timestamp to	false	false	false	false	false	false	false	false	false	true	true	true	false	false
date to	false	false	false	false	false	false	false	false	false	true	true	false	true	false
binary to	false	false	false	false	false	false	false	false	false	false	false	false	false	true

隐式转换支持

LanguageManual | CLI/hiveserver/jdbc

- Commands

- 参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Commands>

- CLI

- 参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Cli>

- Hiveserver/jdbc

- 参考: <https://cwiki.apache.org/confluence/display/Hive/Setting+up+HiveServer2>
<https://cwiki.apache.org/confluence/display/Hive/HiveServer>
<https://cwiki.apache.org/confluence/display/Hive/HiveJDBCInterface>

Hive support file formats

- Hive supports several file formats:
- Text File
- SequenceFile
- [RCFile](#)
- [Avro Files](#)
- [ORC Files](#)
- [Parquet](#)
- Custom INPUTFORMAT and OUTPUTFORMAT

LanguageManual | DDL

- HiveQL DDL statements are documented here, including:
- CREATE DATABASE/SCHEMA, TABLE, VIEW, FUNCTION, INDEX
- DROP DATABASE/SCHEMA, TABLE, VIEW, INDEX
- TRUNCATE TABLE
- ALTER DATABASE/SCHEMA, TABLE, VIEW
- MSCK REPAIR TABLE (or ALTER TABLE RECOVER PARTITIONS)
- SHOW DATABASES/SCHEMAS, TABLES, TBLPROPERTIES, PARTITIONS, FUNCTIONS, INDEX[ES], COLUMNS, CREATE TABLE
- DESCRIBE DATABASE/SCHEMA, table_name, view_name

参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+DDL>

LanguageManual | DDL

```
CREATE (DATABASE | SCHEMA) [IF NOT EXISTS] database_name  
  [COMMENT database_comment]  
  [LOCATION hdfs_path]  
  [WITH DBPROPERTIES (property_name=property_value, ...)];
```

```
DROP (DATABASE | SCHEMA) [IF EXISTS] database_name [RESTRICT | CASCADE];
```

```
ALTER (DATABASE | SCHEMA) database_name SET DBPROPERTIES  
(property_name=property_value, ...); -- (Note: SCHEMA added in Hive 0.14.0)
```

```
ALTER (DATABASE | SCHEMA) database_name SET OWNER [USER | ROLE]  
user_or_role; -- (Note: Hive 0.13.0 and later; SCHEMA added in Hive 0.14.0)
```

```
USE database_name;  
USE DEFAULT;
```

LanguageManual | DDL

```
CREATE [TEMPORARY] [EXTERNAL] TABLE [IF NOT EXISTS] [db_name.]table_name -- (Note:
TEMPORARY available in Hive 0.14.0 and later)
[(col_name data_type [COMMENT col_comment], ...)]
[COMMENT table_comment]
[PARTITIONED BY (col_name data_type [COMMENT col_comment], ...)]
[CLUSTERED BY (col_name, col_name, ...) [SORTED BY (col_name [ASC|DESC], ...)] INTO num_buckets
BUCKETS]
[SKEWED BY (col_name, col_name, ...) -- (Note: Available in Hive 0.10.0 and later)]
ON ((col_value, col_value, ...), (col_value, col_value, ...), ...)
[STORED AS DIRECTORIES]
[
[ROW FORMAT row_format]
[STORED AS file_format]
| STORED BY 'storage.handler.class.name' [WITH SERDEPROPERTIES (...)] -- (Note: Available in Hive
0.6.0 and later)
]
[LOCATION hdfs_path]
[TBLPROPERTIES (property_name=property_value, ...)] -- (Note: Available in Hive 0.6.0 and later)
[AS select_statement]; -- (Note: Available in Hive 0.5.0 and later; not supported for external tables)
```

LanguageManual | DDL

file_format:

- : SEQUENCEFILE

- | TEXTFILE -- (Default, depending on
hive.default.fileformat configuration)

- | RCFILE -- (Note: Available in Hive 0.6.0 and later)

- | ORC -- (Note: Available in Hive 0.11.0 and later)

- | PARQUET -- (Note: Available in Hive 0.13.0 and later)

- | AVRO -- (Note: Available in Hive 0.14.0 and later)

- | INPUTFORMAT input_format_classname

OUTPUTFORMAT output_format_classname

LanguageManual | External Tables

```
create external table if not exists test(  
    month_no STRING,  
    day_id  STRING,  
    data_no STRING,  
) row format delimited fields terminated by ','  
lines terminated by '\n' stored as textfile  
location '/user/hadoop/data'
```


LanguageManual | Partitioned Tables

```
CREATE [EXTERNAL] TABLE [IF NOT EXISTS] table_name [(col_name data_type  
[COMMENT col_comment], ...)] [COMMENT table_comment] [PARTITIONED  
BY (col_name data_type [COMMENT col_comment], ...)] [CLUSTERED BY  
(col_name, col_name, ...) [SORTED BY (col_name [ASC|DESC], ...)] INTO  
num_buckets BUCKETS] [ROW FORMAT row_format] [STORED AS file_format]  
[LOCATION hdfs_path]
```

```
CREATE EXTERNAL TABLE bigdata.test_partitioned (  
  username string,  
  age string  
)  
PARTITIONED BY (dt string, hr string)  
ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' LOCATION  
'/user/hadoop/input';  
--STORED AS parquet; 针对impala可以存储为parquet表
```

LanguageManual | as select *

Create Table As Select (CTAS)

Create table test as select * from a;

参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Select>

LanguageManual | DML

Loading files into tables:

```
LOAD DATA [LOCAL] INPATH 'filepath' [OVERWRITE] INTO TABLE tablename [PARTITION (partcol1=val1, partcol2=val2 ...)]
```

Inserting data into Hive Tables from queries:

```
INSERT OVERWRITE TABLE tablename1 [PARTITION (partcol1=val1, partcol2=val2 ...) [IF NOT EXISTS]]  
select_statement1 FROM from_statement;
```

```
INSERT INTO TABLE tablename1 [PARTITION (partcol1=val1, partcol2=val2 ...)] select_statement1 FROM  
from_statement;
```

```
INSERT INTO TABLE pageviews PARTITION (datestamp = '2014-09-23')  
VALUES ('jsmith', 'mail.com', 'sports.com'), ('jdoe', 'mail.com', null);
```

Hive0.14+ :

```
UPDATE tablename SET column = value [, column = value ...] [WHERE expression]
```

```
DELETE FROM tablename [WHERE expression]
```

参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+DML#LanguageManualDML>

LanguageManual | DML

EXPLAIN Syntax:

EXPLAIN [EXTENDED | DEPENDENCY | AUTHORIZATION] query

EXPLAIN select sum(hash(key)), sum(hash(value)) from src_orc_merge_test_part where ds='2012-01-03' and ts='2012-01-03+14:46:31'

ImportExport:

EXPORT TABLE tablename [PARTITION (part_column="value"[, ...])]
TO 'export_target_path'

IMPORT [[EXTERNAL] TABLE new_or_original_tablename [PARTITION (part_column="value"[, ...])]
FROM 'source_path'
[LOCATION 'import_target_path']

参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Explain>
<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+ImportExport>

LanguageManual | DML

Examples

Simple export and import:

```
export table department to 'hdfs_exports_location/department';  
import from 'hdfs_exports_location/department';
```

Rename table on import:

```
export table department to 'hdfs_exports_location/department';  
import table imported_dept from 'hdfs_exports_location/department';
```

Export partition and import:

```
export table employee partition (emp_country="in", emp_state="ka") to 'hdfs_exports_location/employee';  
import from 'hdfs_exports_location/employee';
```

Export table and import partition:

```
export table employee to 'hdfs_exports_location/employee';  
import table employee partition (emp_country="us", emp_state="tn") from 'hdfs_exports_location/employee';
```

Specify the import location:

```
export table department to 'hdfs_exports_location/department';  
import table department from 'hdfs_exports_location/department'  
    location 'import_target_location/department';
```

Import as an external table:

```
export table department to 'hdfs_exports_location/department';  
import external table department from 'hdfs_exports_location/department';
```

参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Explain>
<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+ImportExport>

LanguageManual | DML

Select Syntax

[WITH CommonTableExpression (, CommonTableExpression)*] (Note: Only available starting with Hive 0.13.0)

SELECT [ALL | DISTINCT] select_expr, select_expr, ...

FROM table_reference

[WHERE where_condition]

[GROUP BY col_list]

[**CLUSTER BY** col_list

| [**DISTRIBUTE BY** col_list] [SORT BY col_list]

]

[LIMIT number]

参考: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Select>

LanguageManual | more

UDF:

<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+UDF>

Locking:

https://cwiki.apache.org/confluence/display/Hive/Locking?from=_YKhAQ

Authorization:

<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Authorization>

HiveCLI:

<https://cwiki.apache.org/confluence/display/Hive/HiveClient>

HiveServer2 Clients:

<https://cwiki.apache.org/confluence/display/Hive/HiveServer2+Clients>

Hive HBase Integration:

<https://cwiki.apache.org/confluence/display/Hive/HBaseIntegration>

Beeline Hive Commands:

<https://cwiki.apache.org/confluence/display/Hive/HiveServer2+Clients>

演示

大数据



Thank you

提问时间?

Blog: <http://www.itweet.cn>

PPT: <https://github.com/itweet/course>

Video: <http://www.tudou.com/home/sparkjvm/>

