

CS 378 – Big Data Programming

Lecture 18

Hadoop Ecosystem

Hadoop Ecosystem

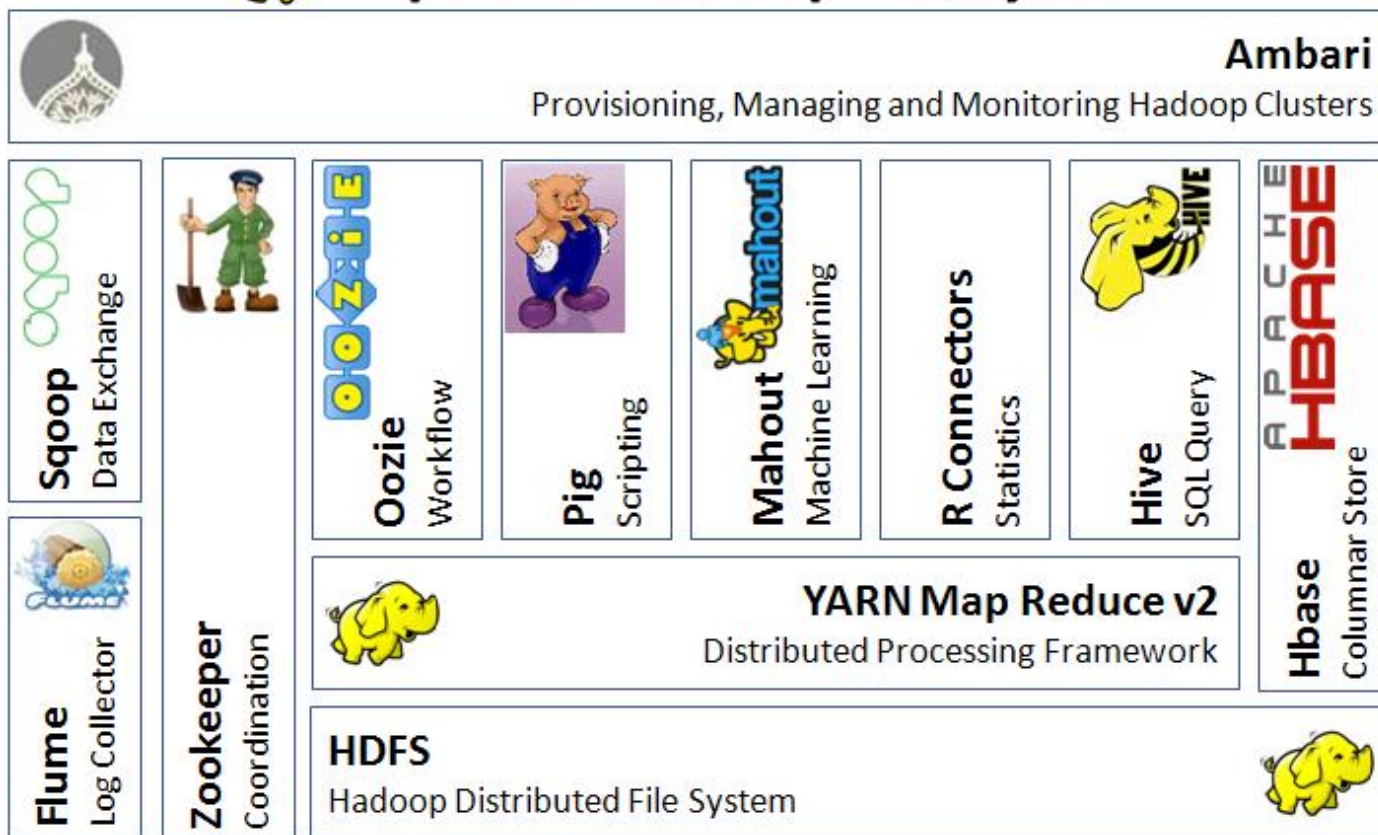
- Many other tools have been implemented on
 - Hadoop
 - HDFS (Hadoop Distributed File System)
- We'll briefly discuss a few
 - HBase
 - ZooKeeper
 - Pig, Impala
 - Hive

Hadoop Ecosystem

thebigdatablog.weebly.com



Apache Hadoop Ecosystem



HBase

- Column-orient database
 - Implemented on top of HDFS
 - Distributed
- Goal is to scale to very large datasets
 - With real-time read/write access

Column-oriented Database

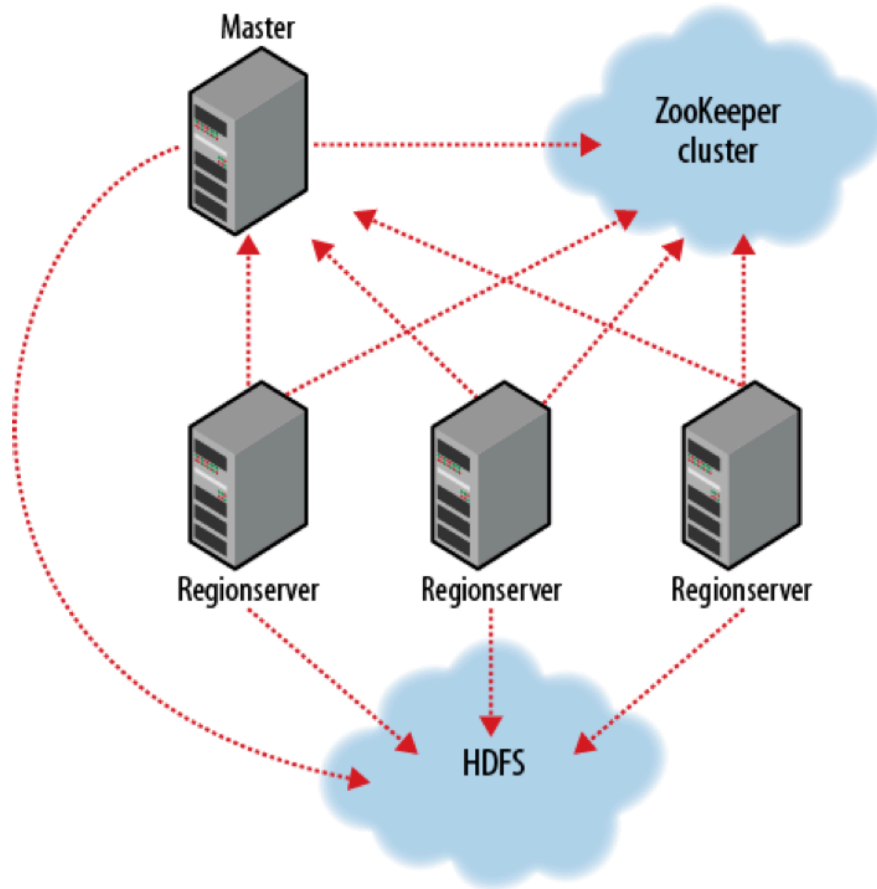
- Table cells are the unit of access
 - Content is uninterpreted array of bytes
 - A cell is versioned (can have multiple versions)
- Table cell is accessed by
 - Row, column, and version (often a timestamp)
- Columns are grouped into families

Column-oriented Database

- New column family members can be added
- Column family members are stored together
- For best performance, family members should be accessed together
- Rows can be subset into regions

Hbase

Figure 13-1 from Hadoop The Definitive Guide 3rd Edition



ZooKeeper

- Messaging and synchronization in a distributed environment
 - Distributed queues, locks
 - Leader election among a group of peers
- High availability (tolerates failures)
- Loosely coupled interactions
 - Rendezvous mechanism

Pig

- Higher level data structures and operations
 - Higher level than Java code for map-reduce job
- Language: Pig Latin
 - Operations and transformations on data
 - Pig converts these to map-reduce jobs for you
- Think of it as a query language for data in HDFS

Pig Examples

Summarization

SQL

The Numerical Aggregation pattern is analogous to using aggregates after a GROUP BY in SQL:

```
SELECT MIN(numericalcol1), MAX(numericalcol1),  
       COUNT(*) FROM table GROUP BY groupcol2;
```

Pig

The GROUP ... BY expression, followed by a FOREACH ... GENERATE:

```
b = GROUP a BY groupcol2;  
c = FOREACH b GENERATE group, MIN(a.numericalcol1),  
                      MAX(a.numericalcol1), COUNT_STAR(a);
```

Pig Examples

Binning

Pig

The SPLIT operation in Pig implements this pattern.

```
SPLIT data INTO  
  eights IF col1 == 8,  
  bigs IF col1 > 8,  
  smalls IF (col1 < 8 AND col1 > 0);
```

Pig Examples

Join

Pig

Pig has native support for a replicated join through a simple modification to the standard join operation syntax. Only inner and left outer joins are supported for replicated joins, for the same reasons we couldn't do it above. The order of the data sets in the line of code matters because all but the first data sets listed are stored in-memory.

```
huge = LOAD 'huge_data' AS (h1,h2);
smallest = LOAD 'smallest_data' AS (ss1,ss2);
small = LOAD 'small_data' AS (s1,s2);
A = JOIN huge BY h1, small BY s1, smallest BY ss1 USING 'replicated';
```

Impala

- Interactive SQL for data in HDFS, HBase
- SQL processing engine
 - Parallel execution
 - Horizontal scaling
- Runs on each data node
 - Direct access to HDFS, HBase (no map-reduce)

Hive

- Data warehouse on top of Hadoop
- SQL for access
 - Hive converts a query into a series of map-reduce steps
- Various Hive clients are available
 - JDBC, ODBC, Thrift, ...