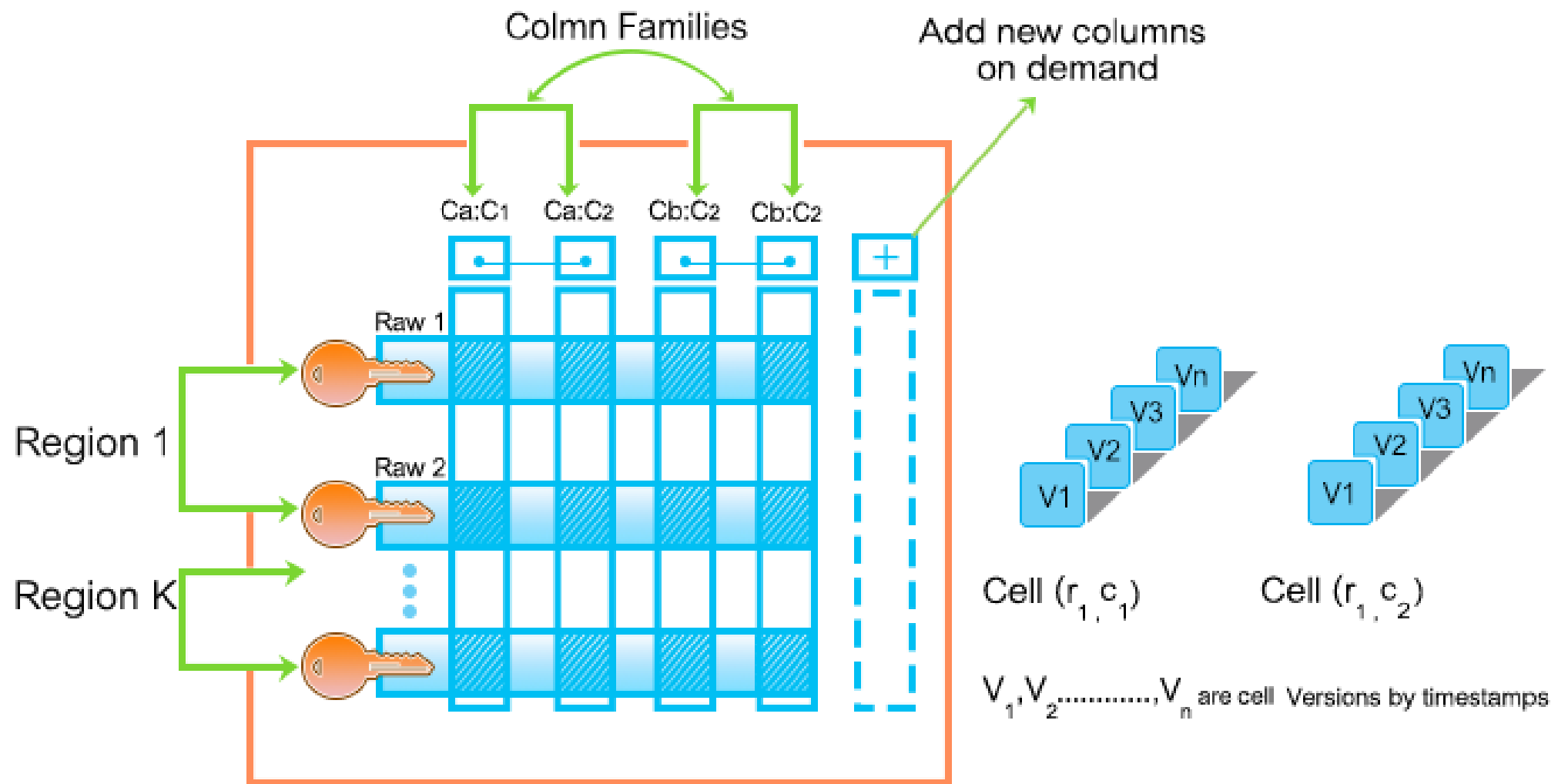


HIVE and HBASE

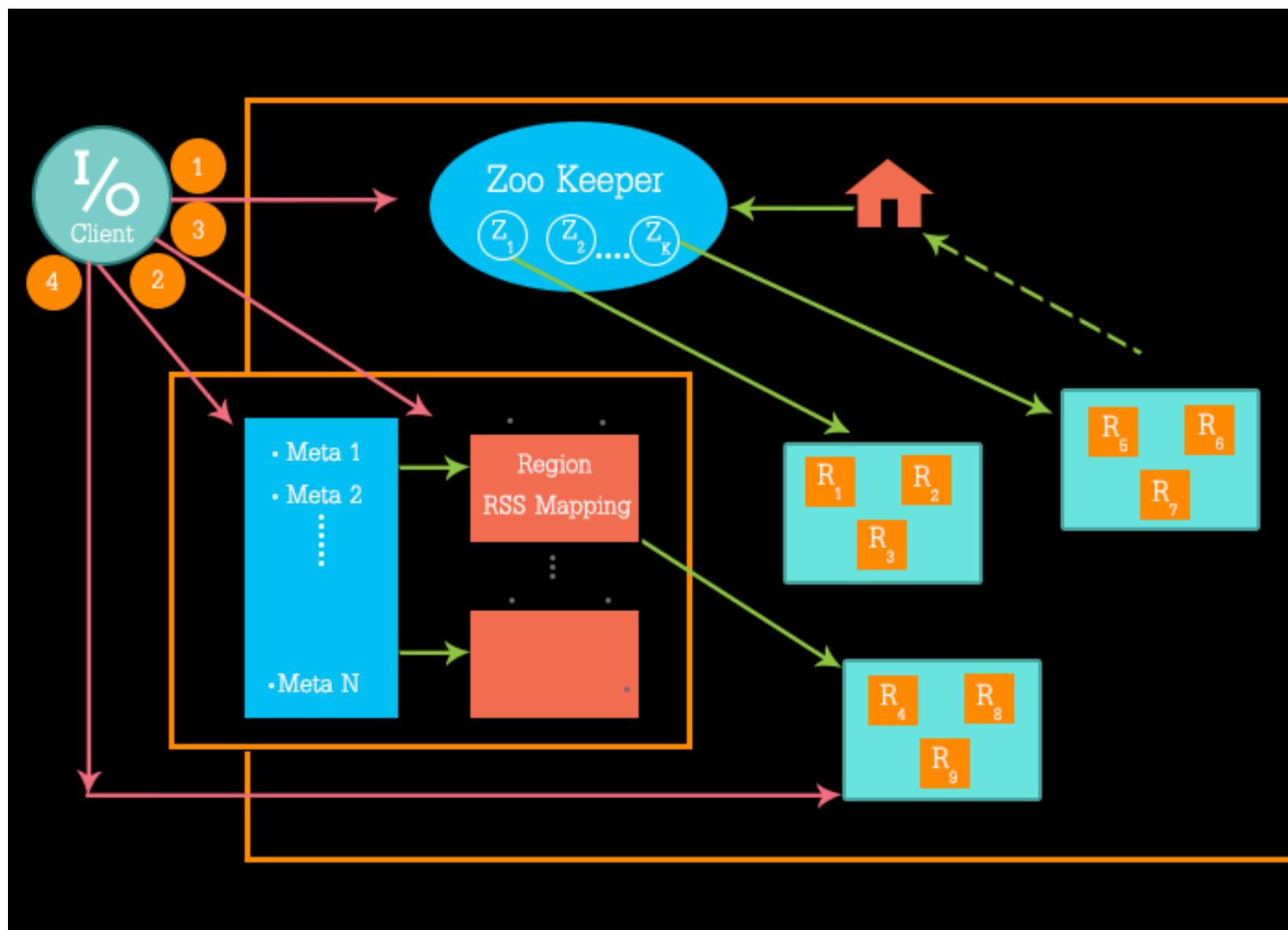
Data Model

- A collection of labelled Tables
- The columns are grouped in column families which can be uniquely identified by column-family prefix
- Every row is uniquely identified by a primary key
- There can be multiple cells with different versions for the same row and column.
- New column families can be added or removed as and when required.



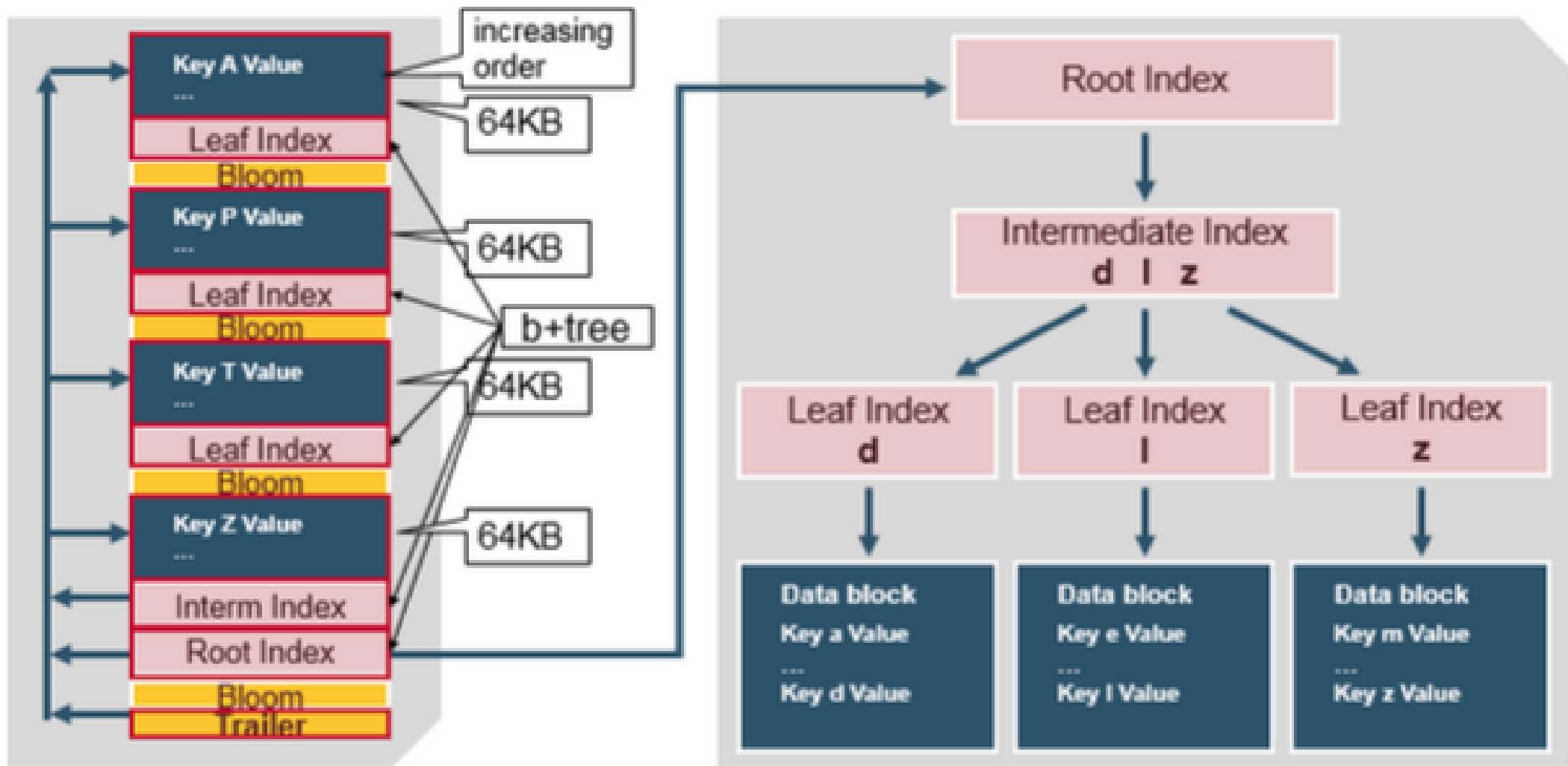
HBase Architecture

- As the number of rows in a table grows beyond a pre-decided limit, the rows are grouped to form regions of the table.
- These regions can be distributed over different nodes (Region Server Slaves) in HBase cluster.
- There is one HBase Master responsible for keeping information and managing regions of table stored on different nodes called Region Server Slaves.
- Zookeeper is appointed by the HBase master to keep active watch over every Region Server Slave by creating a znode entry corresponding to every Region Server Slave.



Hbase file structure

- Data is stored in an HFile which contains sorted key/values.
- An HFile contains a multi-layered index which allows Hbase to seek data without having to read the whole file.
- The multi level index is like a B+ tree.
- Key value pairs are stored in increasing order
- Indexes point by row key to the key value data in 64KB “blocks”
- Each block has its own leaf-index
- The last key of each block is put in the intermediate index
- The root index points to the intermediate index

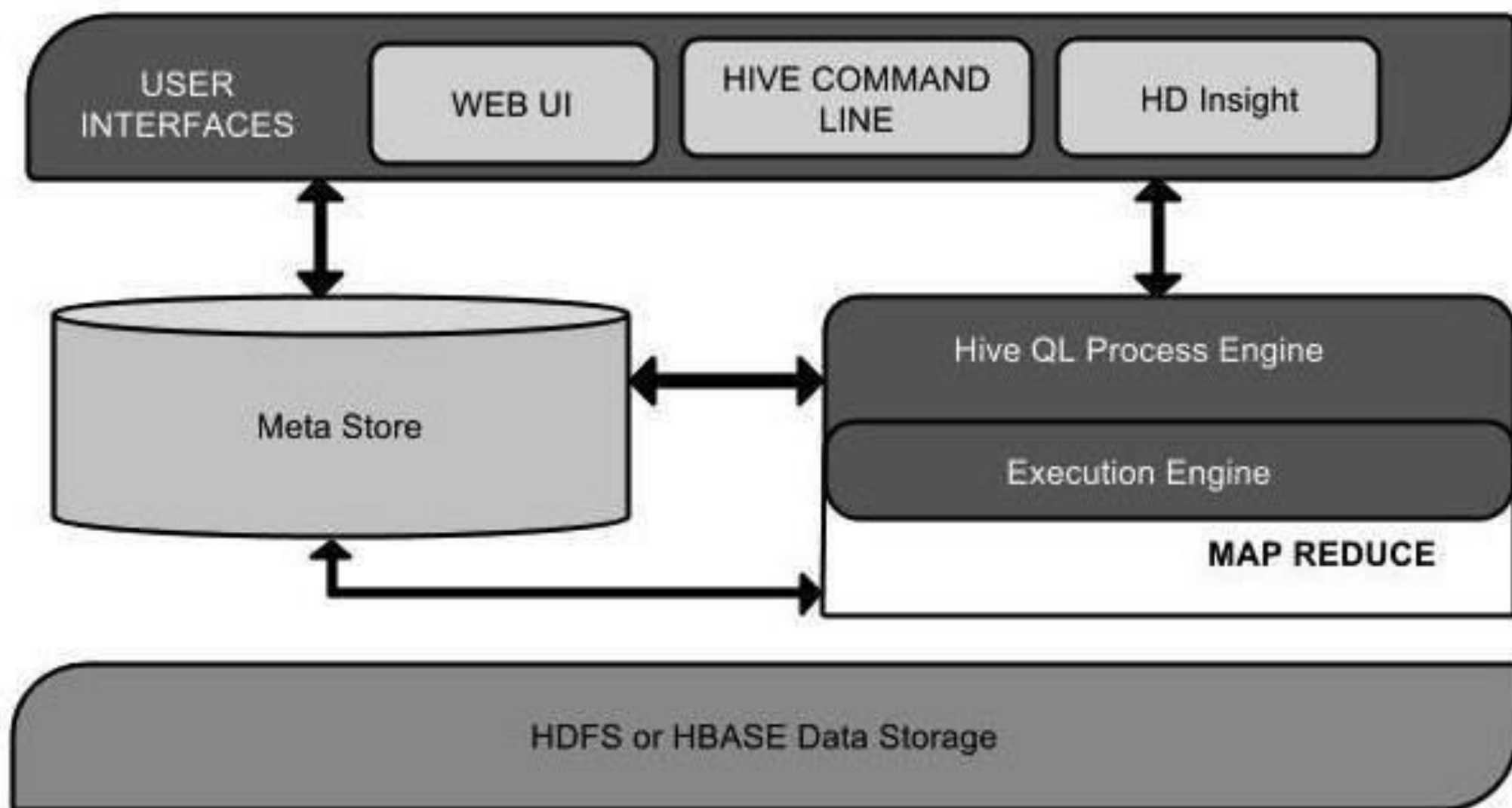


HIVE

- Hive is an open source data warehouse system for querying and analysing large datasets stored in Hadoop files.
- It has three main functions, data summarization, querying and analysis.
- It supports queries expressed in a language called HIVEQL which automatically translates SQL like queries into MapReduce jobs.
- In addition, it also supports custom MapReduce queries.

Architecture

- Hive chooses respective database servers to store the schema or Metadata of tables, databases, columns in a table, their data types, and HDFS mapping.
- HiveQL is similar to SQL for querying on schema info on the Metastore.
- The conjunction part of HiveQL process Engine and MapReduce is Hive Execution Engine. Execution engine processes the query and generates results as same as MapReduce results.
- Hadoop distributed file system or HBASE are the data storage techniques to store data into file system.



Working of HIVE

1 – Execute Query

2 – Get Plan

3 – Get Meta Data

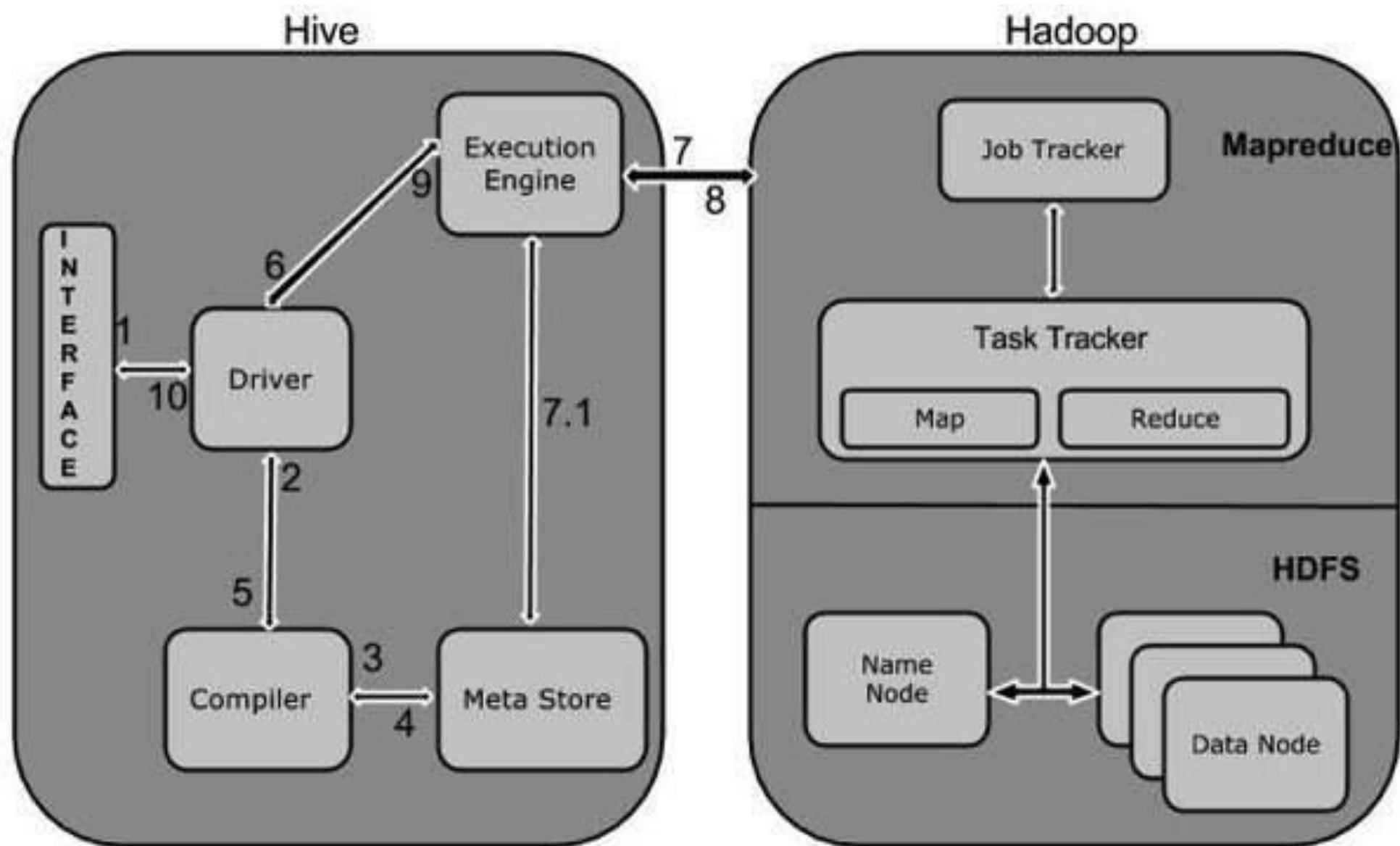
4 - Send Meta Data

5 – Send Plan

6 – Execute Plan

7 – Execute Job (internally MapReduce job)

7.1 – Metadata Ops



Working of HIVE

8 – fetch Result

9 – Send results

10 - Send results

Thank You