

Cassandra is a distributed storage system that manages structured data across many commodity servers. It does so by providing high availability, and no single point of failure. It runs on the infrastructure of nodes spread across different data centers. Thus, the reliability and scalability of the system is handled by Cassandra.

The tables are defined as multidimensional maps, which are key indexed. The keys for the rows and columns are generally strings of 16 to 36 bytes length. There are two kinds of column families: Simple and Super, which are column families within another column family. Columns can be sorted on the basis of the time or name. Different operations that can be performed on the data are as follows, where columnName can be a column, column family, super column:

- insert(table, key, rowMutation)
- get(table, key, columnName)
- delete(table, key, columnName)

Different modules of the distributed systems of the Cassandra are:

Partitioning, Membership, Replication, Bootstrapping, Scaling, Local Persistence.

Each of these modules are implemented in Java. These modules follow a message driven architecture, where the tasks and the messages are divided into different stages, on the basis of the SEDA architecture. All these modules synchronously handle read/write requests. These requests are routed to any of the nodes in the cluster. On receiving a write request, the node determines the replicas, and routes the request accordingly and waits for a quorum of replicas to acknowledge the write completion. On receiving a read request, the system routes the request to the closest replica or all the replicas, based on the consistency guaranteed, and waits for the response.

Strong Points:

1. Cassandra uses Zookeeper for electing a leader amongst all the nodes. Thus, each of the nodes asks the leader for its respective range of replicas.
2. Data persistence on the local file system. The writes performed are sequential to the disks and indexes are generated for efficient lookup. As time passes, a process runs in background for all such files, to merge into one.
3. Column indices are used to retrieve the columns which are away from the multidimensional keys, thereby providing faster lookups.

Weak Points:

1. Cassandra being a completely decentralized system , does need Zookeeper for some co-ordination among the nodes.
2. Detecting failure takes much time, as the size of the clusters grows in the system.
3. The paper provides very little information for Scuttlebutt, an entropy-based gossip algorithm.

Question:

What is the mechanism of the Scuttlebutt algorithm?

