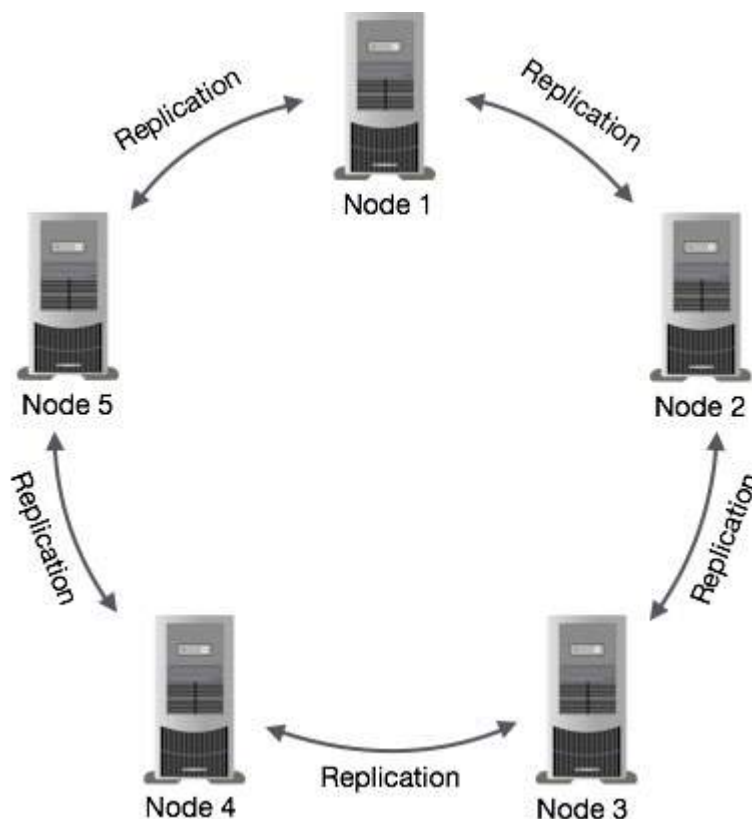# CASSANDRA - ARCHITECTURE

The design goal of Cassandra is to handle big data workloads across multiple nodes without any single point of failure. Cassandra has peer-to-peer distributed system across its nodes, and data is distributed among all the nodes in a cluster.

- All the nodes in a cluster play the same role. Each node is independent and at the same time interconnected to other nodes.

- Each node in a cluster can accept read and write requests, regardless of where the data is actually located in the cluster.

- When a node goes down, read/write requests can be served from other nodes in the network.

## Data Replication in Cassandra

In Cassandra, one or more of the nodes in a cluster act as replicas for a given piece of data. If it is detected that some of the nodes responded with an out-of-date value, Cassandra will return the most recent value to the client. After returning the most recent value, Cassandra performs a **read repair** in the background to update the stale values.

The following figure shows a schematic view of how Cassandra uses data replication among the nodes in a cluster to ensure no single point of failure.



**Note** – Cassandra uses the **Gossip Protocol** in the background to allow the nodes to communicate with each other and detect any faulty nodes in the cluster.

## Components of Cassandra

The key components of Cassandra are as follows −

- **Node** − It is the place where data is stored.

- **Data center** − It is a collection of related nodes.

- **Cluster** − A cluster is a component that contains one or more data centers.

- **Commit log** − The commit log is a crash-recovery mechanism in Cassandra. Every write operation is written to the commit log.

- **Mem-table** − A mem-table is a memory-resident data structure. After commit log, the data will be written to the mem-table. Sometimes, for a single-column family, there will be multiple mem-tables.

- **SSTable** − It is a disk file to which the data is flushed from the mem-table when its contents reach a threshold value.

- **Bloom filter** − These are nothing but quick, nondeterministic, algorithms for testing whether an element is a member of a set. It is a special kind of cache. Bloom filters are accessed after every query.

## Cassandra Query Language

Users can access Cassandra through its nodes using Cassandra Query Language *CQL*. CQL treats the database *Keyspace* as a container of tables. Programmers use **cqlsh:** a prompt to work with CQL or separate application language drivers.

Clients approach any of the nodes for their read-write operations. That node *coordinator* plays a proxy between the client and the nodes holding the data.

## Write Operations

Every write activity of nodes is captured by the **commit logs** written in the nodes. Later the data will be captured and stored in the **mem-table.** Whenever the mem-table is full, data will be written into the **SStable** data file. All writes are automatically partitioned and replicated throughout the cluster. Cassandra periodically consolidates the SSTables, discarding unnecessary data.

## Read Operations

During read operations, Cassandra gets values from the mem-table and checks the bloom filter to find the appropriate SSTable that holds the required data.

Loading [MathJax]/jax/output/HTML-CSS/jax.js