

## Kubernetes Architecture.

It consists of multiple worker nodes; it could be physical or virtual. Each Physical node has its own disks (SSD or HDD), network, DRAM. Each node may host multiple pods & containers.

## FoundationDB Deployment.

It consists of hundreds of nodes per datacenter; with cross-Datacenter (DC) latency of around 20-30 milliseconds & inner-DC latency of 0.2 milliseconds. If one DC fails, DCs from other regions are active. Each DC can tolerate 2 storage node failures.

It is able to manage transactions with multi version concurrency control & optimistic concurrent control, layers can be built on top to support other data models;

It follows 'Divide-and-conquer' approach, so each component is like a model that handles a specific task.

It consists of:

## 1. Control plane

- Co-ordinators that form a PAXOS group to reach consensus on some decision.
- Singletons meaning that there is only one process across all of the MDB nodes.

## 2. Data Plane

- Transaction System that assigns read & commit versions to transactions.
- Log Systems that store transaction logs.
- Storage System: Read the transaction logs.

The system is designed to handle failure, whenever encountered with a failure, a reconfiguration process brings the transaction management system to a new configuration. It has a Sequence process monitoring the health of proxies, Resolvers & log servers. If any one of monitored processes fails or the database configuration changes the sequence process terminates.

## FDB Transactions

The client sends a request to a ProxyServer asking for a read version, here the proxy server gets the read version from the sequencer. The client issues reads directly to the StorageServers. During write there is Buffer locally at the client; The client sends read & write sets to a proxy server. The proxy asks the sequencer for a commit version. The proxy will ask Resolvers to check for read-write conflicts. The proxy sends committed mutations to log servers for persistence. The proxy reports the committed version to the sequencer before return success to client.

## FDB Resolvers Conflict Checking.

When a transaction commits, it includes read & write sets which track the keys it has read from & written to. The resolvers job is to detect conflicts between transactions by comparing their read & write sets. If 2 transactions attempt to modify overlapping keys, or if one reads a key that another transaction writes, the resolver flags a conflict. The conflicting transaction will be aborted to maintain serializability.

## FDB Logging Protocol

Once the conflict has been resolved, the next step is that the proxy would try to shift all of the mutation log to the designated log server. The proxy uses an in-memory shard map to determine the responsible storage servers for the key range & attach storage server tags to the mutation. Log message headers contain sequence numbers & metadata for tracking. After all replica log servers confirms durability, the proxy updates the <sup>KCV</sup> Known Committed Version with the latest log sequence number, guaranteeing that the data is consistent & durable across the systems.

## FDB Storage Server.

- Maintains versions of data in the last 5 seconds in memory.
- An older version is stored persistently on KV engine (SQLite).

## How FDB Storage Server Serves a Read Request?

FDB storage server handles a read request by checking the KCV & last applied LSN for the key requested. If the request version is less than the known commit version, the server returns the value. However, if the request version is greater than the transaction may not have fully committed yet. The server will attempt to serve the most recent committed version yet, it either waits for the logs to be applied or directs the client to try another server storing the key.

## 1. What does commit do in a read-only transaction (no write)?

A. It will skip the commit protocol because read can be serialized before any write that happens in concurrent.

## 2. Is it possible to have multiple resolvers resolve conflicts for a transaction?

A. When the transaction provides with a set of range, it partitions into different resolvers to resolve conflict; ∴ Yes, it is possible.

## 3. What if a storage server is lagging behind on applying the redo logs &amp; a client requests a version of a key pair it does not have?

A. It will wait or the client will try another storage server before giving up.

## Data Model built on FDB.

## - Graph

It's elements include vertices (nodes) & edges. Each vertex / edge has one label & zero, one or more properties. A edge connects 2 vertices.

## Janus Graph

It is designed to manage large number of nodes. It is open source & uses Tinkerpop, a graph computing framework with Gremlin query language.

## Graph Operations.

## - Add a vertex.

```
g.addV("airport")
    .property("city", "Los Angeles")
    .property("code", "LAX")
```

```
g.addV("airport")
    .property("city", "Sydney")
    .property("code", "SYD")
```

## - Add an edge

```
g.addEdge("route")
    .property("dist", "11235")
    .from(V(1)).to(V(2))
```

Add 2 keys; one represents source vertex & other destination vertex

## - Get all outgoing edges of a vertex.

```
g.v(1).outE("route").valueMap()
```

## - Delete a vertex.

```
g.v(1).drop()
```



## Janus Graph Indices

## Composite Index Example.

It constructs a key based on the label & the property of the value is the vertex id, so in order to get vertex id for any vertex; it just needs one look up.

## Vertex-Centric Index Example.

This is to deal with the super node issue, where we can have any vertex with thousands of vertices that connect to it.

## Mixed indices.

Are external indices, to support global range queries & text search.

FoundationDB is a distributed key-value store that is horizontally scalable & strong consistency.