

Components in FoundationDB

The FoundationDB architecture chooses a decoupled design, where processes are assigned different heterogeneous roles (e.g., Coordinators, Storage Servers, Master). Scaling the database is achieved by horizontally expanding the number of processes for separate roles:

Coordinators

All clients and servers connect to a FoundationDB cluster with a cluster file, which contains the IP:PORT of the coordinators. Both the clients and servers use the coordinators to connect with the cluster controller. The servers will attempt to become the cluster controller if one does not exist, and register with the cluster controller once one has been elected. Clients use the cluster controller to keep an up-to-date list of proxies.

Cluster Controller

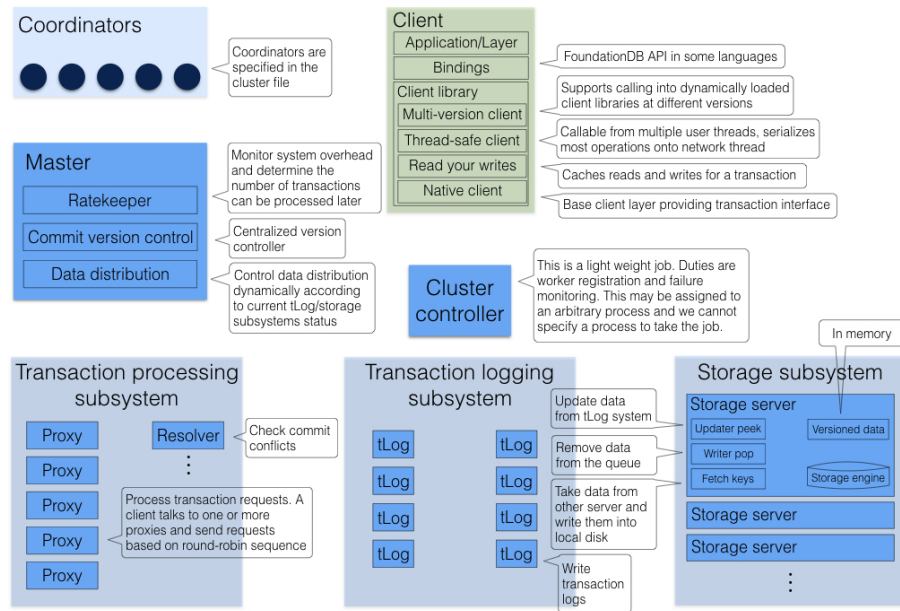
The cluster controller is a singleton elected by a majority of coordinators. It is the entry point for all processes in the cluster. It is responsible for determining when a process has failed, telling processes which roles they should become, and passing system information between all of the processes.

Master

The master is responsible for coordinating the transition of the write sub-system from one generation to the next. The write sub-system includes the master, proxies, resolvers, and transaction logs. The three roles are treated as a unit, and if any of them fail, we will recruit a replacement for all three roles. The master provides the commit versions for batches of the mutations to the proxies.

Historically, Ratekeeper and Data Distributor are coupled with Master on the same process. Since 6.2, both have become a singleton in the cluster. The life time is no longer tied with Master.

Components



Proxies

The proxies are responsible for providing read versions, committing transactions, and tracking the storage servers responsible for each range of keys. To provide a read version, a proxy will ask all other proxies to see the largest committed version at this point in time, while simultaneously checking that the transaction logs have not been stopped. Ratekeeper will artificially slow down the rate at which the proxy provides read versions.

Commits are accomplished by:

- Get a commit version from the master.
- Use the resolvers to determine if the transaction conflicts with previously committed transactions.
- Make the transaction durable on the transaction logs.

The key space starting with the `\xff` byte is reserved for system metadata. All mutations committed into this key space are distributed to all of the proxies through the resolvers. This metadata includes a mapping between key ranges and the storage servers which have the data for that range of keys. The proxies provides this information to clients on-demand. The clients cache this mapping; if they ask a storage server for a key it does not have, they will clear their cache and get a more up-to-date list of servers from the proxies.

Transaction Logs

The transaction logs make mutations durable to disk for fast commit latencies. The logs receive commits from the proxy in version order, and only respond to the proxy once the data has been written and fsync'd to an append only mutation log on disk. Before the data is even written to disk we forward it to the storage servers responsible for that mutation. Once the storage servers have made the mutation durable, they pop it from the log. This generally happens roughly 6 seconds after the mutation was originally committed to the log. We only read from the log's disk when the process has been rebooted. If a storage server has failed, mutations bound for that storage server will build up on the logs. Once data distribution makes a different storage server responsible for all of the missing storage server's data we will discard the log data bound for the failed server.

Resolvers

The resolvers are responsible determining conflicts between transactions. A transaction conflicts if it reads a key that has been written between the transaction's read version and commit version. The resolver does this by holding the last 5 seconds of committed writes in memory, and comparing a new transaction's reads against this set of commits.

Storage Servers

The vast majority of processes in a cluster are storage servers. Storage servers are assigned ranges of key, and are responsible to storing all of the data for that range. They keep 5 seconds of mutations in memory, and an on disk copy of the data as of 5 second ago. Clients must read at a version within the last 5 seconds, or they will get a `transaction_too_old` error. The SSD storage engine stores the data in a B-tree based on SQLite. The memory storage engine store the data in memory with an append only log that is only read from disk if the process is rebooted. In the upcoming FoundationDB 7.0 release, the B-tree storage engine will be replaced with a brand new Redwood engine.

Data Distributor

Data distributor manages the lifetime of storage servers, decides which storage server is responsible for which data range, and ensures data is evenly distributed across all storage servers (SS). Data distributor as a singleton in the cluster is recruited and monitored by Cluster Controller. See [internal documentation](#).

Ratekeeper

Ratekeeper monitors system load and slows down client transaction rate when the cluster is close to saturation by lowering the rate at which the proxy provides read versions. Ratekeeper as a singleton in the cluster is recruited and monitored by Cluster Controller.

Clients

A client links with specific language bindings (i.e., client libraries) in order to communicate with a FoundationDB cluster. The language bindings support loading multiple versions of C libraries, allowing the client communicates with older version of the FoundationDB clusters. Currently, C, Go, Python, Java, Ruby bindings are officially supported.

Transaction Processing

A database transaction in FoundationDB starts by a client contacting one of the Proxies to obtain a read version, which is guaranteed to be larger than any of commit version that client may know about (even through side channels outside the FoundationDB cluster). This is needed so that a client will see the result of previous commits that have happened.

Then the client may issue multiple reads to storage servers and obtain values at that specific read version. Client writes are kept in local memory without contacting the cluster. By default, reading a key that was written in the same transaction will return the newly written value.

At commit time, the client sends the transaction data (all reads and writes) to one of the Proxies and waits for commit or abort response from the proxy. If the transaction conflicts with another one and cannot commit, the client may choose to retry the transaction from the beginning again. If the transaction commits, the proxy also returns the commit version back to the client. Note this commit version is larger than the read version and is chosen by the master.

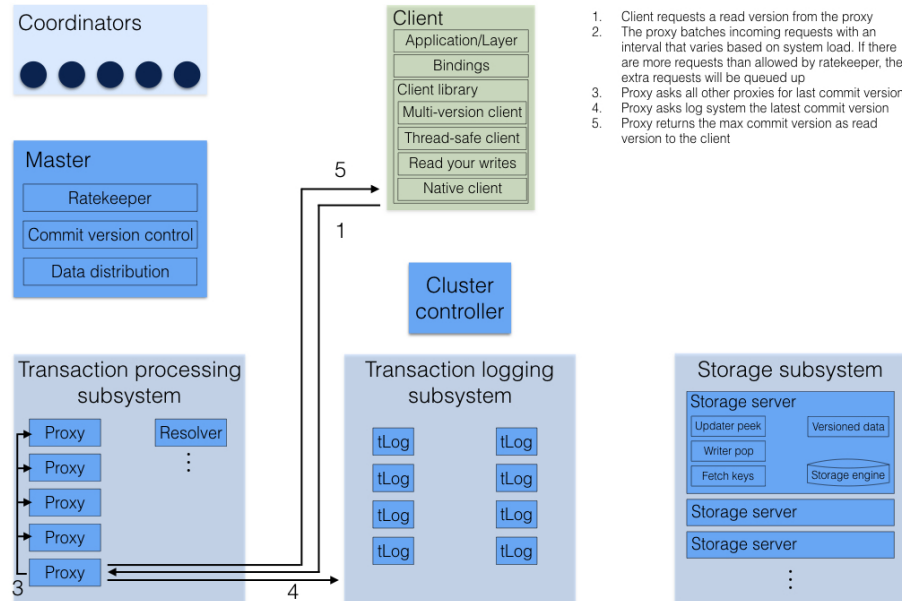
The FoundationDB architecture separates the scaling of client reads and writes (i.e., transaction commits). Because clients directly issue reads to sharded storage servers, reads scale linearly to the number of storage servers. Similarly, writes are scaled by adding more processes to Proxies, Resolvers, and Log Servers in the transaction system.

Determine Read Version

When a client requests a read version from a proxy, the proxy asks all other proxies for their last commit versions, and checks a set of transaction logs

satisfying replication policy are live. Then the proxy returns the maximum commit version as the read version to the client.

Start a transaction



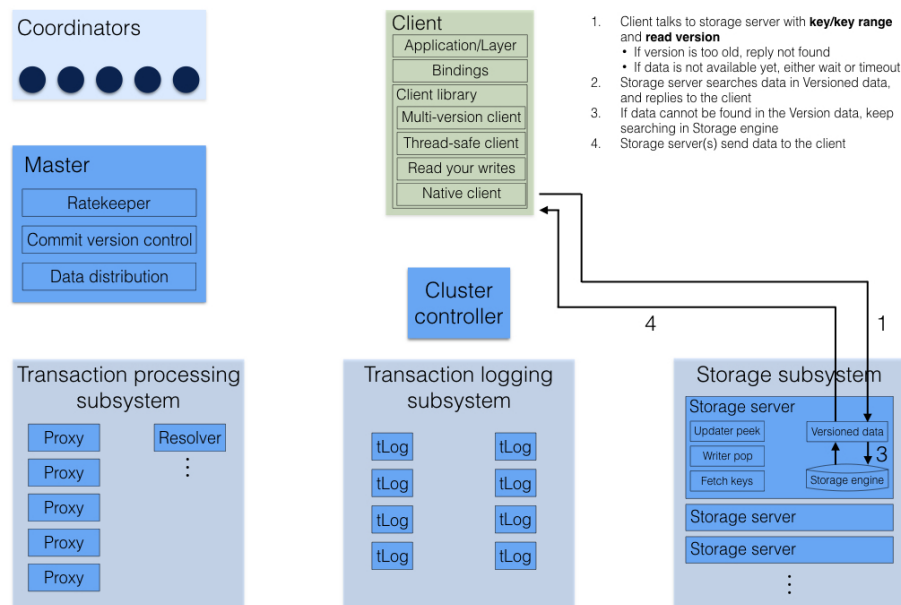
The reason for the proxy to contact all other proxies for commit versions is to ensure the read version is larger than any previously committed version. Consider that if proxy A commits a transaction, and then the client asks proxy B for a read version. The read version from proxy B must be larger than the version committed by proxy A. The only way to get this information is by asking proxy A for its largest committed version.

The reason for checking a set of transaction logs satisfying replication policy are live is to ensure the proxy is not replaced with newer generation of proxies. This is because proxy is a stateless role recruited in each generation. If a recovery has happened and the old proxy is still live, this old proxy could still give out read versions. As a result, a read-only transaction may see stale results (a read-write transaction will be aborted). By checking a set of transaction logs satisfying replication policy are live, the proxy makes sure no recovery has happened, thus the read-only transaction sees the latest data.

Note that the client cannot simply ask the master for read versions. The master gives out versions to proxies to be committed, but the master does not know when the versions it gives out are durable on the transaction logs. Therefore it is not safe to do reads at the largest version the master has provided because that version might be rolled back in the event of a failure, so the client could end up reading data that was never committed.

In order for the client to use versions from the master, the client needs to wait until all in-flight transaction-batches (a write version is used for a batch of transactions) to commit. This can take a long time and thus is inefficient. Another drawback of this approach is putting more work towards the master, because the master role can't be scaled. Even though giving out read-versions isn't very expensive, it still requires the master to get a transaction budget from the Ratekeeper, batches requests, and potentially maintains thousands of network connections from clients.

Read transaction flow



Transaction Commit

A client transaction commits in the following steps:

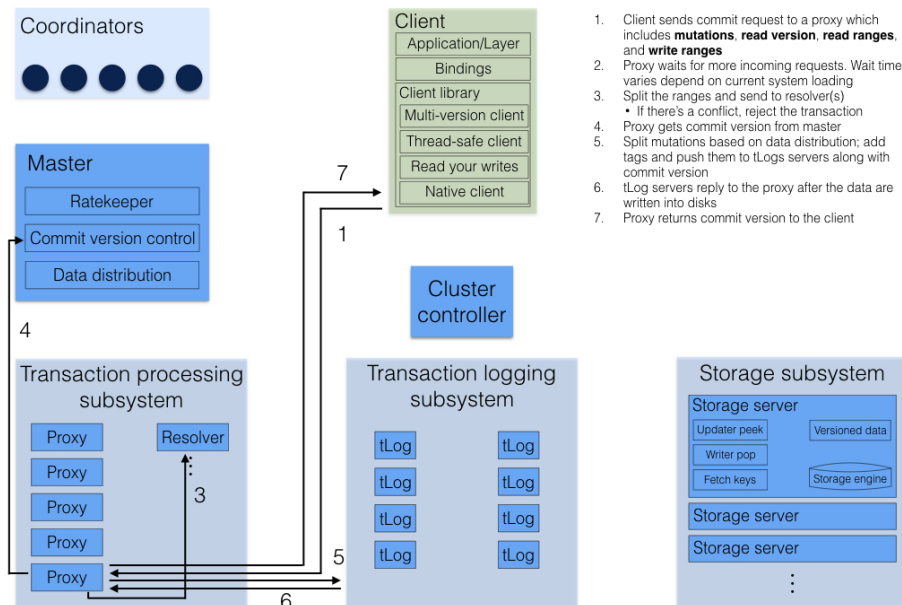
1. A client sends a transaction to a proxy.
2. The proxy asks the master for a commit version.
3. The master sends back a commit version that is higher than any commit version seen before.
4. The proxy sends the read and write conflict ranges to the resolver(s) with the commit version included.
5. The resolver responds back with whether the transaction has any conflicts with previous transactions by sorting transactions according to their commit versions and computing if such a serial execution order is conflict-free.
 - If there are conflicts, the proxy responds back to the client with a

not_committed error.

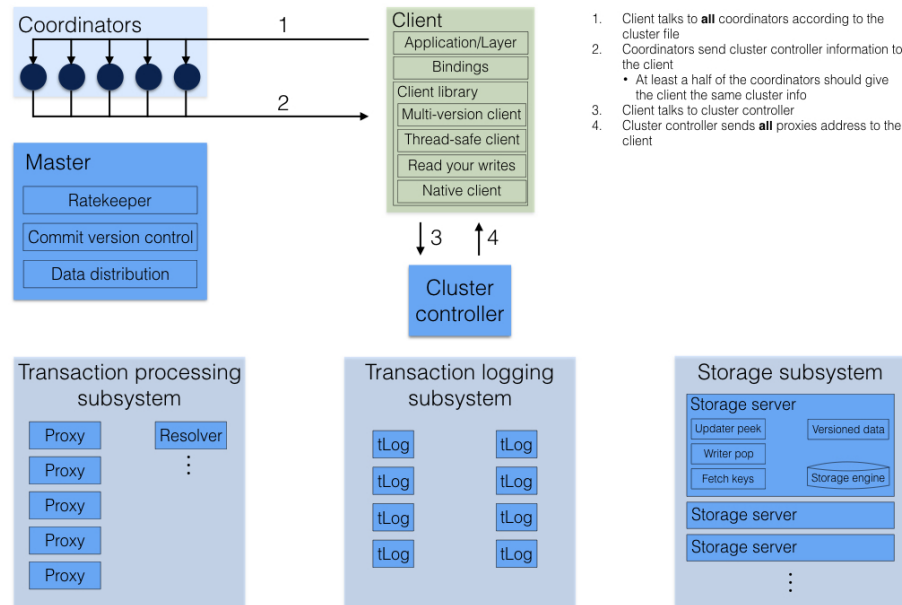
- If there are no conflicts, the proxy sends the mutations and commit version of this transaction to the transaction logs.
6. Once the mutations are durable on the logs, the proxy responds back success to the user.

Note the proxy sends each resolver their respective key ranges, if any one of the resolvers detects a conflict then the transaction is not committed. This has the flaw that if only one of the resolvers detects a conflict, the other resolver will still think the transaction has succeeded and may fail future transactions with overlapping write conflict ranges, even though these future transaction can commit. In practice, a well designed workload will only have a very small percentage of conflicts, so this amplification will not affect performance. Additionally, each transaction has a five seconds window. After five seconds, resolvers will remove the conflict ranges of old transactions, which also limits the chance of this type of false conflict.

Write transaction flow



Start and open a database

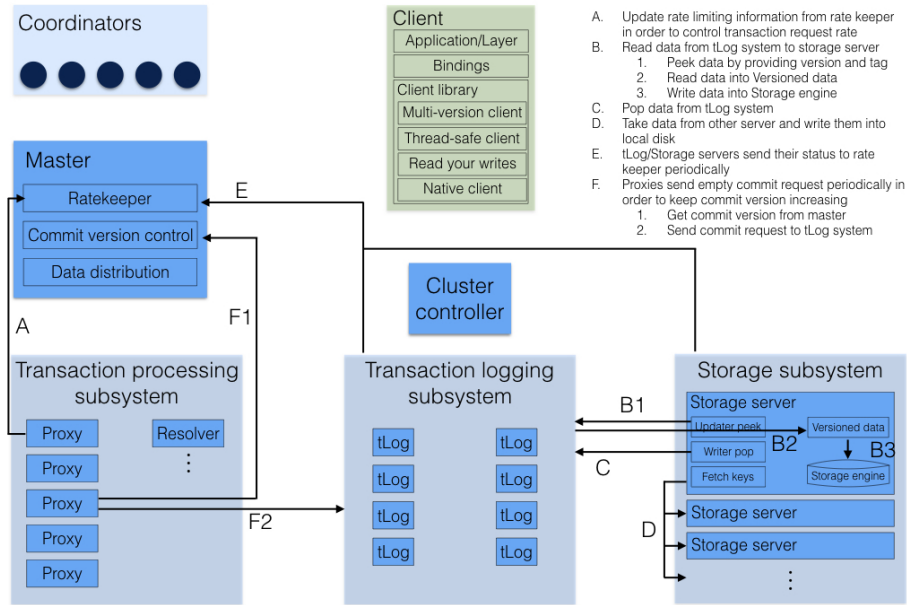


Background Work

There are a number of background work happening besides the transaction processing:

- **Ratekeeper** collects statistic information from proxies, transaction logs, and storage servers and compute the target transaction rate for the cluster.
- **Data distribution** monitors all storage servers and perform load balancing operations to evenly distribute data among all storage servers.
- **Storage servers** pull mutations from transaction logs, write them into storage engine to persist on disks.
- **Proxies** periodically send empty commits to transaction logs to keep commit versions increasing, in case there is no client generated transactions.

Background works



Transaction System Recovery

The transaction system implements the write pipeline of the FoundationDB cluster and its performance is critical to the transaction commit latency. A typical recovery takes about a few hundred milliseconds, but longer recovery time (usually a few seconds) can happen. Whenever there is a failure in the transaction system, a recovery process is performed to restore the transaction system to a new configuration, i.e., a clean state. Specifically, the Master process monitors the health of Proxies, Resolvers, and Transaction Logs. If any one of the monitored process failed, the Master process terminates. The Cluster Controller will detect this event, and then recruits a new Master, which coordinates the recovery and recruits a new transaction system instance. In this way, the transaction processing is divided into a number of epochs, where each epoch represents a generation of the transaction system with its unique Master process.

For each epoch, the Master initiates recovery in several steps. First, the Master reads the previous transaction system states from Coordinators and lock the coordinated states to prevent another Master process from recovering at the same time. Then the Master recovers previous transaction system states, including all Log Servers' Information, stops these Log Servers from accepting transactions, and recruits a new set of Proxies, Resolvers, and Transaction Logs. After previous Log Servers are stopped and new transaction system

is recruited, the Master writes the coordinated states with current transaction system information. Finally, the Master accepts new transaction commits. See details in this [documentation](#).

Because Proxies and Resolvers are stateless, their recoveries have no extra work. In contrast, Transaction Logs save the logs of committed transactions, and we need to ensure all previously committed transactions are durable and retrievable by storage servers. That is, for any transactions that the Proxies may have sent back commit response, their logs are persisted in multiple Log Servers (e.g, three servers if replication degree is 3).

Finally, a recovery will fast forward time by 90 seconds, which would abort any in-progress client transactions with `transaction_too_old` error. During retry, these client transactions will find the new generation of transaction system and commit.

commit_result_unknown error: If a recovery happened while a transaction is committing (i.e., a proxy has sent mutations to transaction logs). A client would have received `commit_result_unknown`, and then retried the transaction. It's completely permissible for FDB to commit both the first attempt, and the second retry, as `commit_result_unknown` means the transaction may or may not have committed. This is why it's strongly recommended that transactions should be idempotent, so that they handle `commit_result_unknown` correctly.

Resources

[Forum Post](#)

[Existing Architecture Documentation](#)

[Summit Presentation](#)

[Data Distribution Documentation](#)

[Recovery Documentation](#)