

2022-02-24

Title to be decided

David Kleingeld

Email

Contents

1 Background	2
1.1 Distributed Computing	2
1.2 Faults and Delays	2
1.3 Consensus Algorithms	3
1.4 File System	10
1.5 Distributed file systems	11

1 Background

1.1 Distributed Computing

When state of the art hardware is no longer fast enough to run a system the only option is scaling out. Then there is a choice, do you buy expensive, reliable high performance supercomputer or commodity servers connected by Ip and ethernet? This is the choice between High Performance (HPC) and Distributed Computing. With HPC faults in the hardware are rare and can be handled by restarting, easing development. In a distributed context faults are the norm, restarting the entire system is not an option or you would be down all the time. Resilience against faults comes at an, often significant, cost to performance. It may also limit scalability. As the scale of a system increases so does the frequency with which one of the parts fails. Even the most robust part will fail and given enough of them the system will fail frequently. Therefore at very large scales HPC is not even an option.

1.2 Faults and Delays

Before we can build a fault resistant system we need to know what we need to keep in mind. While hardware failures are, the norm in distributed computing, faults are not the only issue to keep in mind.

It is entirely normal for the clock of a computer to run slightly to fast or to slow. The drift will be tens of milliseconds [1] unless special measures are taken¹. Worse a process can be paused and then resumed at any time. Such a pause could be because the process thread is pre-empted, because its virtual machine is paused or because the process was stopped and resumed after a while².

In a distributed system the computers that form the system or *the nodes*, are connected by IP over ethernet. Ethernet gives no guarantee a packet is delivered on time or at all. A node can be unreachable before suddenly working fine again.

¹One could synchronize the time within a datacenter or provide nodes with more accurate clocks

²On linux by sending SigStop then SigCont

A system model is an abstraction defining what an algorithm can assume. Regarding timing there are three models.

1. The Synchronous model allows an algorithm to assume that clocks are synchronized within some bound and network traffic will arrive within a fixed time.
2. The Partially synchronous model is a more realistic model. Most of the time clocks will be correct within a bound and network traffic will arrive within a fixed bound. However sometimes clocks will drift unbounded and some traffic might be delayed forever.
3. The Asynchronous model has no clock, it is very restrictive.

For most distributed systems we work with the Partially Synchronous model. We assume hardware faults cause a crash from which the node can be recovered later. Either automatically as it restarts or after maintenance.

1.3 Consensus Algorithms

In this world where the network can not be trusted, time lies to us and servers will randomly crash and burn how can we get anything done at all? Let's discuss how we can build a system we can trust, a system that behaves *consistently*. To build such a system we need the parts that make up the system to agree with each other, they must have *Consensus*. Here I discuss three well known solutions. Before we get to that let's look at the principle that underlies them all: *The truth is defined by the majority*.

Quorums

Imagine a node hard at work processing requests from its siblings, suddenly it gets pre-empted. The other nodes notice it is no longer responding and declare it dead, they don't know its threads got paused. A few seconds later the node responds again as if nothing had happened, and in truth, unless it checks the system clock, from its perspective no time has passed. Alternatively a network error might partition the system, each group of servers can reach each other but not the others. The nodes in the group will declare those in the other group dead and continue their work. Usually this results in data loss if the work progresses at all.

We can solve this by voting over each decision. It will be a strange vote, no node cares about the decision itself. In most implementations the nodes only check if it regards the sender as trustworthy or alive and then vote yes. To prove liveness the vote proposal could include a number. Voters only vote yes if the number is correct. For example if the number is the highest they have seen. If a majority votes yes the node that requested the vote can be sure it's not dead or disconnected. This is the idea behind "Quorums," majorities of nodes that vote.

Paxos

The PAXOS algorithm[5] uses a quorum to provide consensus. It does this by choosing a value among proposals such that only that value can be read as the accepted value. Usually it is used to build a fault tolerant distributed state machine.

In PAXOS there are three roles: proposer, acceptor and learner. It is possible for nodes to fulfill only one or two of these roles. For the rest of this explanation assume each node fulfills all three. To reach consensus on a new value we go through two phases: prepare and accept. Once the majority of the nodes has accepted a proposal the value of that proposal has been chosen. In PAXOS nodes keep track of the highest proposal number n they have seen. Let's go through a PAXOS iteration from the perspective of a node trying to share something, *a value*.

In the first phase a new *value* is proposed by our node. It sends a *prepare* request to a majority of acceptors. The request contains a proposal number n higher than the highest number our node has seen up till now. The number is unique to our node. Each acceptor only responds if our number n is the highest it has seen. If an acceptor had already accepted one or more requests it includes the accepted proposal with the highest n in its response.

In phase two our node checks if it got a response from the majority. Our node is going to send an accept request back to those nodes. The content of the accept request depends on what our node received in response to its prepare request:

1. Our node received a response with number n_p . This means an acceptor has already accepted a value. If we continued with our own value the system would have two different accepted values. Therefore the content of our accept request will be the value from proposal n_p .

2. It received only acknowledging replies and none contained a previously accepted value. The system has not yet decided on a value. The content of our accept request will be the value or node wants to propose but with our number n .

The acceptors accept the request if they did not yet receive a prepare request numbered greater than n . On accepting a request an acceptor sends a message to all learners³. This way the learners learn a new value as soon as its ready.

Let's get a feeling why this works by looking at what happens during node failure. Imagine a case where a minimal majority m accept value v_a . A single node in m freezes after the first learners learned of the now chosen value v_a . After freezing $m - 1$ of the nodes will reply v_a as value to learners. The learners will conclude no value has been chosen given $m - 1$ is not a majority⁴. As seen above acceptors change their value if they receive a higher numbered accept request. If a single node changes its value to v_b consensus will break since v_a has already been seen as the chosen value by a learner. A new proposal that can result into higher numbered accept requests needs a majority response. A majority response will include a node from $m - 1$. That node will include v_a as the accepted value. The value for the accept request changes to v_a . No accept request with another value than v_a can thus be issued. Another value v_b will therefore never be accepted. The new accept by issued to a majority changing at least one node to have accepted v_a . Now at least $m + 1$ nodes have v_a as accepted value.

To build a distributed state machine you run multiple instances of PAXOS. This is often referred to as MULTI-PAXOS. The value for each instance is a command to change the shared state. Unfortunately multi paxos is not specified in literature and has never been verified.

Raft

The PAXOS algorithm allows us to reach consensus on a single value. The RAFT algorithm allows us to share a log between nodes. We can only append to and reading from the log. That is the log will always be the same on all nodes. As long as a majority of the nodes still function the log will be readable and appendable.

³remember usually every node is a learner

⁴this is not yet inconsistent, PAXOS does not guarantee consistency over whether a value has been chosen

RAFT maintains a single leader. The leader is decided on by a *quorum*. Appending to the log is sequential because only the leader is allowed to append. There are two parts to RAFT, *electing leaders* and *log replication*.

Leader election A RAFT[6] cluster starts without a leader and when it has a leader it can fail at any time. The cluster must be able to reliably decide on a new leader. Nodes in RAFT start as followers, monitor the leader by waiting for heartbeats. If a follower does not receive a heartbeat from a *valid* leader on time it will try to become the leader, it becomes a candidate. In a fresh cluster without a leader one or more nodes become candidates.

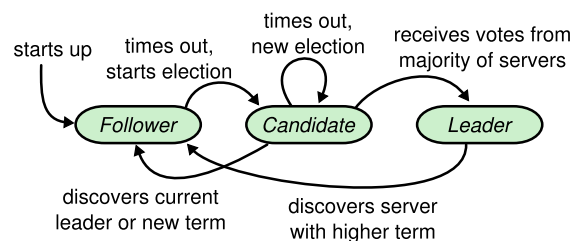


Figure 1: A RAFT node states. Most of the time all nodes except one are followers. One node is a leader. As failures are detected by time outs the nodes change state. Adjusted from [6]

A candidate tries to get itself elected. For that it needs the votes of a majority of the cluster. It asks all nodes for their vote. Note that servers vote only once and only if the candidate would become a *valid* leader. If a majority of the cluster responds with a vote the candidate it becomes the leader. If it takes too long to receive a majority of the votes a candidate starts a fresh election. When there are multiple candidates asking for votes the votes might split⁵, no candidate then reaches a majority. A candidate immediately loses the election if it receives a heartbeat from a *valid* leader. These state changes are illustrated in fig. 1.

In RAFT time can be divided in terms. A term is the period from the election of a node to leader to its failure or a failed election where no node won, illustrated in fig. 2. Terms are used to determine the cluster leader, the leader with the highest terms. A heartbeat is *valid* if, as far as the receiver knows, it originates from the current cluster leader. A message can only be from the *current* leader if the message *term* is equal or higher than the receiving nodes

⁵Election timeouts are randomised so this does not repeat infinitely.

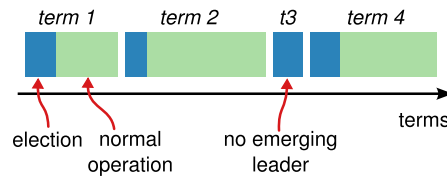


Figure 2: An example of how time is divided in terms in a RAFT cluster

term. If a node receives a message with a higher term it replaces its own with that of the message.

When a node starts its *term* is zero. If a node becomes a candidate it increments its *term* by one. Imagine a candidate with a *term* equal or higher than that of the majority of the cluster. When receiving a vote request the majority will determine this candidate could become a *valid* leader. This candidate will get the majority vote in the absence of another candidate and become *the* leader.

Log replication To append an entry to the log a leader sends an append request to all nodes. Messages from invalid leaders (section 1.3) are rejected. The leader knows an entry is committed after a majority of nodes acknowledged. For example entry 5 in fig. 3 is committed. The leader includes the index up to which entries are committed in all its messages. This means entries will become committed on all followers at the latest with the next heartbeat. If the leader approaches the timeout and no entry needs to be added it sends an empty append. There is no need for a special heartbeat message.

There are a few edge cases that require careful followers to be careful when appending. A follower may have an incomplete log if it did not receive a previous append, it may have messages the leader does not and finally we must prevent a candidate missing committed entries from becoming the leader.

1. To detect missing logs entries are indexed incrementally. The leader includes the index of the previous log entry in the request and the term when it was appended. If a follower's last log entry does not match the included entry and term it responds with an error. The leader will send

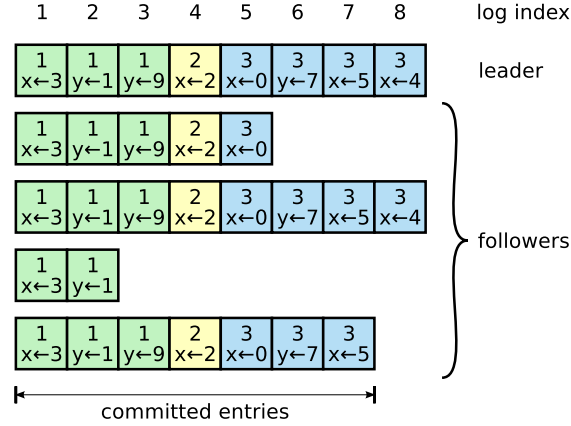


Figure 3: Logs for multiple nodes. Each row is a different node. The log entries are commands to change shared variables x and y to different values. The shade boxes and the number at their top indicate the term.

back its complete log for the follower to duplicate⁶.

2. When a follower has entries the leader does not these will occupy indices the leader will use in the future. This happens when a previous leader crashed having pushed logs to some but not majority of followers. When the leader pushes a new entry with index k the follower will notice it already has an entry with index k . At that point it simply overwrites what it had at k .
3. A new leader missing committed entries will push a wrong log to followers missing entries. For an entry to be committed it must be stored on the majority of the cluster. To win the election a node has to have the votes from the majority. Thus restricting followers to only vote for candidates that are as up-to-date as they are is enough.

Log compaction Keeping the entire log since the cluster was first started is rather inefficient. Especially as nodes are added and need to get send the log. Usually raft is used to build a state machine, in that case sending over the state is faster than the log. The state machine send is a snapshot of the system state, only valid for the moment in time it was taken. All nodes take snapshots independently of the leader.

⁶This is rather inefficient, in the next paragraph we will come back to this

To take a snapshot a node writes the state machine to file together with the current *term* and *index*. It then discards all committed entries up to the snapshotted *index* from its log.

Followers that lag far behind now send the snapshot and all logs that follow. The follower wipes its log before accepting the snapshot.

Consensus as a service

So the problem of consensus has been solved, but the solutions are non trivial to implement. We need to shape our application to fit the models the solutions use. If we are using `RAFT` that means our systems must be built around a log. Then we need to build the application praying we implement the solution correctly. A popular alternative is to use a coordination service. Here we look at `ZOOKEEPER`[3] an example of a wait free coordination service. I first focus on the implementation, drawing parallels to `RAFT` before we look at the *Api* `ZOOKEEPER` exposes.

A ZooKeeper cluster has a designated leader that replicates a database on all nodes. Only the leader can modify the database, therefore only the leader handles *write* requests. The nodes handle *read* requests themselves, *write* requests are forwarded to the leader. To handle a *write* request `ZOOKEEPER` relies on a consensus algorithm called `ZAB`[4]. It is an atomic broadcast capable of crash recovery. It uses a strong leader quite similar to `RAFT` and guarantees that changes broadcast by the leader are delivered in the order they were sent and after changes from previous (crashed) leaders. `ZAB` can deliver messages twice during recovery. ZooKeeper turns change requests into *idempotent* transactions. These can be applied multiple times with the same result bypassing the double delivery problem of `ZAB`.

ZooKeeper exposes its strongly consistent database as a hierarchical name space of *znodes*. Each *znode* contains a small amount of data and is identified by its *path*. Using the *Api* clients can operate on the *znodes*. They can:

1. create new *znodes*
2. change the data in a *znode*
3. delete existing *znodes*
4. sync with the leader
5. query if a *znode* exists
6. read the data in a *znode*
7. get list of the children of the *znode*

The last three operations support watching, the client then gets a single

notification when the result of the operation changed. The notification does not contain any information regarding the change and the value is no longer watched after the notification. Clients can use this to ensure locally cached data is up-to-date.

Clients communicating outside of ZooKeeper need to take special measures to ensure consistency. For example let client A update some znode from value v_1 to value v_2 then communicates to client B, through another medium then ZooKeeper. In this example client A and B are connected to different ZooKeeper nodes. If the communication from A causes B to read the znode it will get the previous value v_1 from ZooKeeper if its ZooKeeper node is lagging behind. To avoid this race condition client B can call *sync* first which will make the zookeeper node process all outstanding requests from the leader before returning.

RAFT has this same race condition. Intuitively it seems as we can just ensure a heartbeat has passed, by then all (functioning) node will be updated. However this is not enough, as a faulty node could be suspended and not notice it is outdated at all. Instead the solution is to include the last log index client A saw in its communication to client B. PAXOS does not *need* to suffer from this problem but it can. In PAXOS reading means asks a *learner* for the value, the *learner* can then ask the majority of the system if they have accepted a value and if so what it is. Usually PAXOS is optimized up by making acceptors inform learners of a change. In this case a leader that missed a message from an acceptor that there is a value will incorrectly return to client B there is no value.

1.4 File System

A file system is split into two parts, the files and the directory structure. File properties, or metadata, such as its name, identifier, size etc are stored in the directory. Typically the directory entry itself only contains the file name and its unique identifier. Using the identifier the other metadata for the file can be fetched. The content of the file is split into blocks these blocks are stored on stable storage such as an hard drive or ssd. The file system defines an API to allow modifying the files system providing ways to *create*, *read*, *write*, *seek* and *truncate* files.

Usually the system adds a distinction between open and closed files. The apis *read* *write* and *seek* are then only allowed on open files. This makes it possible to provide some consistency guarantees in a concurrent environment. For example allowing a file to be opened only if it was not already open. This can

prevent a user from corrupting data by writing from multiple processes at the same place in the file. There is no risk to reading the same file from multiple process, even while appending to it from other processes⁷. To allow such use a file systems can define opening a file in read-only, append-only or read-write mode. On Linux this is opt in⁸. Even more semantics exist for example allowing opening multiple non overlapping ranges of a file for writing.

1.5 Distributed file systems

Here I will discuss the two most widely used distributed file systems. We will look at how they work and the implementation. Before I get to that I will use a very basic file sharing system, Network File System, to illustrate why these distributed systems need their complexity.

Network File System

A basic way to share files is to expose a filesystems via a network to share. For this you use a Network File System. These integrate in the interface of the client. A widely supported system is NFS. In NFS a part of a local directory is exported/shared by a local NFS-server. Other machines can then connect and overlay part of their directory with the exported one. The NFS protocol forwards file operations from the client to the host over the network. When an operation has been applied on the host the result is traced back to the client. To increase performance the client (almost always) caches file blocks and metadata.

In a shared environment it is commonplace for multiple users to simultaneously access the same files. In NFS this can be problematic, as meta data is cached new files can appear to other users after 30 seconds. Further more simultaneous writes can become interleaved as each write gets split into multiple network packets [8, p. 527], writing corrupt data. Version 4 improves the semantics respecting unix advisory file locks [7]. Most applications do not take advisory locks into account still risking data corruption.

⁷The OS can ensure append writes are serialized, this is useful for writing to a log file where each write call appends an entire log line to a file opened in append mode

⁸see flock or fcntl or mandatory locking

Google file system

The Google File System [2] was developed in 2003 in a response to Google's rapidly growing search index. The key to the system is the separation of the control plane from the data plane. That is, the file data is stored on many *chunk servers* while a single server coordinates where that data is stored, how it is replicated and access to it. The single coordinating server dramatically simplifies the design, GFS needs no consensus algorithm as all decisions are made on a single machine. When a GFS client wants to use a file it contacts the main node for metadata. The client then uses the metadata to determine on which chunk servers the file content is located. Finally it streams the bytes directly of the chunk servers. GFS is meant to serve huge files as many small files could easily bottleneck the single coordinating server. With large files the load will be dominated by streaming bytes. File metadata a single machine can handle those requests, simplifying consistency. The much higher load of streaming file data is spread over many *chunk servers*.

Hadoop FS - implementation of GFS that has expanded over the years. drives Hadoop (map reduce) - high availability module uses quorum - opt in consistency

Ceph

Subtree partitioning

References

- [1] M Caporali and R Ambrosini. “How closely can a personal computer clock track the UTC timescale via the internet?” In: *European Journal of Physics* 23.4 (June 2002), pp. L17–L21. doi: 10.1088/0143-0807/23/4/103. url: <https://doi.org/10.1088/0143-0807/23/4/103>.
- [2] Sanjay Ghemawat, Howard Gobioff and Shun-Tak Leung. “The Google File System”. In: *Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles*. SOSP ’03. Bolton Landing, NY, USA: Association for Computing Machinery, 2003, pp. 29–43. isbn: 1581137575. doi: 10.1145/945445.945450. url: <https://doi.org/10.1145/945445.945450>.
- [3] Patrick Hunt et al. “{ZooKeeper}: Wait-free Coordination for Internet-scale Systems”. In: *2010 USENIX Annual Technical Conference (USENIX ATC 10)*. 2010.
- [4] Flavio P Junqueira, Benjamin C Reed and Marco Serafini. “Zab: High-performance broadcast for primary-backup systems”. In: *2011 IEEE/IFIP 41st International Conference on Dependable Systems & Networks (DSN)*. IEEE. 2011, pp. 245–256.
- [5] Leslie Lamport. “Paxos Made Simple”. In: *ACM SIGACT News (Distributed Computing Column)* 32, 4 (Whole Number 121, December 2001) (Dec. 2001), pp. 51–58. url: <https://www.microsoft.com/en-us/research/publication/paxos-made-simple/>.
- [6] Diego Ongaro and Ousterhout John. *In Search of an Understandable Consensus Algorithm (Extended Version)*. <https://raft.github.io/>. accessed 15-Feb-2022. 2014.
- [7] S Shepler et al. *Network File System (NFS) version 4 Protocol*. RFC 3530. IETF, Apr. 2003. url: <https://www.ietf.org/rfc/rfc3530.txt>.
- [8] Abraham Silberschatz, Peter Baer Galvin and Greg Gagne. *Operating system concepts*. John Wiley & Sons, 2014.