

Paxos made moderately working

May 2020

Contents

1	Introduction	2
2	Consistency and progress condition	3
3	Assumptions	3
4	Design concerns	4
5	The single-decree synod	5
5.1	Mathematical results	5
5.2	The preliminary protocol	6
5.3	The basic protocol	8
5.3.1	Step 1	8
5.3.2	Step 2	9
5.3.3	Step 3	10
5.3.4	Step 4	11
5.3.5	Step 5	12
5.3.6	Step 6	13
5.4	The complete protocol	13
6	Conclusion and opening on the multi-decree parliament	15

1 Introduction

Distributed systems theory is really a troublesome field.

The main idea is that we want a set of multiple computers to behave like there was only one.

For example, we want our database to be split among multiple machines to avoid things like data losses, service interruption, server overload, storage issues and many other fun problems. Basically we want to improve the performances of a system by getting it to run on multiple machines.

This implies that all the machines composing our system needs to stay synchronized.

There are 3 main components that we want to maintain:

- Consistency: every machines in the system as to agree with each other, you would not want a database that doesn't respond the same thing depending on the actual machine that received the request.
- Availability: I want my system to be always able to respond, I don't want service interruptions
- Partition tolerance: This is a tricky one but basically I want my system to work as intended even if some machines are separated from each other (when a partition occurs). For example if I have 3 machines (A, B and C) in my system, and it turns out that A and B can access each other but can no longer access C (but C is still up, it is just not accessible), then a partition occured and I don't want my system to diverge on the two sides of the partition because they are not able to communicate.

This looks really interesting but one of the first theorem you will learn while working on distributed systems is that it is impossible to maintain these 3 components at once.

Many algorithms have addressed this issue by making concessions. They are also many paper that explains the problems in further details as well as adding other ones (Byzantine generals, FLP ...) but we will not talk about these.

We will talk about one legendary algorithm called Paxos, described in *The Part-Time Parliament*.

This algorithm is known for being really hard to understand but also because its original research paper tells a story to describe the algorithm.

I really wanted to understand this algorithm, so I reimplemented it in C++.

You can find the code [here](#).

The present article acts as a guide to navigate between the implementation and the original paper.

2 Consistency and progress condition

We will first talk of the key parts of the algorithm and how they are presented in the paper.

The author created a fictional civilisation on the island of Paxos to tell his story. The culture of the paxons (residents of Paxos) is really particular because they have the bad habit of leaving whatever they are doing to do their personal activities.

Even the legislators in the parliament could leave in the middle of a vote. Yet paxons had laws, and those laws had to be written somewhere after being voted. How could the legislators agree with each other on the current state of the law, and how could they create a new law?

These two problems are respectively called the "Consistency" and "Progress" conditions in the paper.

Of course, this is an image for a distributed system. The system is the parliament and each legislator is a machine. The fact that a legislator might leave the parliament for his personal activities represents the fact that the machine either crashed or is not accessible anymore. Yet, each legislator had to know all the laws in place at the current time and must be able to vote for a new one.

The authors then describes how paxons represented the law. The law is defined by a sequence of numbered decrees, for example:

- 132: Lamps must use only olive oil
- 155: The olive tax is 3 drachmas per ton

To record the law, one might want to hire a secretary to register every decree, but in Paxos no one was willing to act as a secretary and stay in parliament for the entire session.

Instead, each legislator kept a ledger containing the decrees that were passed. Each legislator had its own version of the law.

We can now better rephrase the consistency condition stated earlier: No two ledgers can contain contradictory information. It means that if a legislator has the decree 132 in the above example, no other legislator can have a different entry for decree 132.

However, it was possible for a legislator to not have any entry for decree 132, for example if he didn't learn about it yet.

Of course, the consistency condition could be trivially fulfilled by leaving every ledger blank so we needed the progress condition.

In Paxos, an atmosphere of mutual trust prevailed, it means that legislators were willing to accept any decree proposed by any other legislator.

This might cause some consistency issues. For example, a group of legislators pass the decree **37: Painting on temple walls is forbidden**. And then leave the parliament before another group of legislators enters and pass (without knowing anything about the previously passed decree) the decree **37: Freedom of artistic expression is guaranteed**.

To solve this, we need at least one legislator in the first group to also be in the second, so he can warn them that this decree number already exists.

To make sure that any two groups of legislators have at least one legislator in common, we simply need to make sure that both of them are a majority. This means that a decree can only be passed if a majority of legislators are in the parliament. This guarantee allows to define the progress condition as follows: If a majority of the legislators were in the parliament, then any decree proposed by a legislator in the parliament would be passed and every decree that had been passed would appear in the ledger of every legislator in the parliament.

We now need to maintain those 2 conditions.

3 Assumptions

The article then states a few assumptions for the organisation of the parliament.

First, every legislator has a personal ledger, in which they write decrees.

This ledger represents any kind of persistent data storage, in our implementation, it will simply be a file, but it could be an actual database.

It is stated that legislator could take notes at the back of their ledgers to not forget what they were doing after leaving the parliament. We will also use files for this matter, we will use them to keep track of data relative to the algorithm that MUST NOT be lost upon crashes.

Legislators would write other notes on a slip of paper which he might (or not) lose upon leaving the parliament. This represents simple memory storage.

Legislators had an hourglass to measure time, it means that our servers must have internal clocks, in our implementation we will not use them as the complete protocol requires because it isn't useful in the case where we are initiating decrees manually (which we do in our implementation).

It is then specified that the parliament had bad acoustics, the only way for legislators to communicate was to send messengers. Messengers could leave the parliament for personal activities the same way legislators did. It means that our machines cannot have direct inter-process communication, we will here use TCP connections. The application layer of our message will be further detailed when we will talk of the protocol.

4 Design concerns

Let's put the basis of our implementation, we will only detail here what is relevant for the protocol, all the socket, events or configuration isn't interesting here so we will only focus on the protocol in itself by making abstraction of everything else.

Let's build our first class, the Legislator, as stated in the assumptions it contains a ledger and some volatile data for the algorithm but we do not know them yet.

```
1  class Legislator
2  {
3  public:
4      Legislator(const LegislatorConfig& config);
5      LegislatorConfig config_;
6
7  private:
8      Ledger ledger;
9  };
10 using shared_legislator = std::shared_ptr<Legislator>;
11
12 extern shared_legislator self;
13 extern std::unordered_map<std::string, shared_legislator> legislators;
```

The LegislatorConfig simply hold the ip and port of the server, as well as a name to identify it. Legislators are created with the factory design pattern using a configuration file, it simply contains configurations of all the legislators in the parliament.

The name of the legislator currently running is given via the command line.

So if you want to simulate a parliament with 2 legislators, you need to run the program twice with the same configuration file, but by putting a different name (that must be both present in the configuration file). There is nothing really worth noting in the Ledger class, it simply contains a bunch of methods to retrieve data (for example decrees).

You will also notice 2 global variables, a map between the name of the legislator and the legislator itself, and a pointer to the legislator that corresponds to that particular machine.

We will add a lot of methods in this class, it is really our central point.

We then have the Decree class.

It is simply a value object containing an int. The decree will only be an int here.

5 The single-decree synod

5.1 Mathematical results

We now have everything we need to start the algorithm, the paper presents it in 2 steps.

First there is the single-decree synod protocol and then the multi-decree parliament protocol which directly derives from the first. The single-decree protocol is itself a derivation of multiple sub-protocols. There is first the preliminary protocol, that is the mathematical basis of the algorithm, then the basic protocol that derives directly from it that respects the consistency condition and then the complete protocol that adds the progress condition.

The goal of the single-decree protocol is to synchronize a single decree between all our legislators. It is worth to note that the proposed implementation is actually only the basic single-decree protocol, the multi-decree part is yet to be done and the complete single-decree protocol is irrelevant in our implementation.

The decree is chosen through a series of numbered ballots.

A ballot is simply a referendum on a single decree, in each ballot each legislator has the choice to either vote or to not vote for the decree.

To each ballot was associate a non-empty set of legislators called the quorum, the ballot succeeded (and the decree was passed) if and only if every legislator in the quorum voted for it. A ballot B is said to be earlier than a ballot C if B's ballot number is inferior to C's.

Let's now consider a set β of ballots, we will see that by adding a few conditions, once a ballot has succeeded, every later ballots will be fore the same decree, thus satisfying the consistency condition.

The consistency condition is satisfied if:

- $B1(\beta)$: Each ballot in β has a unique ballot number
- $B2(\beta)$: The quorum of any two ballots in β have at least one legislator in common
- $B3(\beta)$: For every ballot B in β , if any legislator in B's quorum voted in an earlier ballot in β , then the decree of B equals the decree of the latest of those earlier ballots

Rather than going into the maths of why this works I will try to build a strong intuition of this.

So let's consider ballot B with number 4 for decree D. Given the success of B, every legislators in B's quorum voted for it by definition of a successful ballot.

Let's now consider ballot C with number 5 whose decree is yet to be determined. C's quorum contains at least one legislator in common with B's quorum (because of $B2(\beta)$).

Since every legislator in B's quorum voted in B for decree D, we can immediatly deduct that at least one legislator in C's quorum previously voted in B for decree D. B being the earliest ballot before C, the decree of C MUST be D to satisfy $B3(\beta)$.

Take a moment to think about it, it means that once a ballot succeeded, every later ballot will be for the same decree, if a legislator was not in the parliament at the time the decree was originally passed, he would still learn its existence thanks to the new ballots, holding the same decree.

5.2 The preliminary protocol

We can derive a protocol from the requirement that $B1(\beta) - B3(\beta)$ remain true where β is the set of all ballots that were or are being conducted.

Each legislator could initiate a ballot by choosing its number, decree and quorum. The choice derive directly from the need to maintain $B1(\beta) - B3(\beta)$.

To maintain $B1(\beta)$, each ballot needs to receive a unique number. A legislator can take notes in his ledger to remember of which ballot he already initiated to avoid picking the same number twice. To avoid two different legislators from initiating ballots with the same number, the set of possible ballot number is partitionned among the legislators. There are multiple way of doing this, in our implementation we give ballot numbers to legislators following a round-robin basis.

If we have 3 legislators (A, B and C), legislator A will have ballot numbers 0, 3, 6 ...

legislator B will have 1, 4, 7 ... and legislator C will have 2, 5, 8 ...

To maintain $B2(\beta)$ we can simply choose the quorum as any majority set among the legislators.

$B3(\beta)$ requires that the legislator initiating the ballot knows the decree every legislators in the quorum voted for. To do this we will have to exchange messages with the legislators in the quorum.

Let's now address what messages exactly must be sent to conduct a ballot.

- 1: Legislator l chooses to initiate a new ballot, he do so by choosing a new ballot number b among those still available to him and send a NextBallot(b) message to some set of legislators (in practice we will send it to all the legislators)
- 2: Upon receiving a NextBallot(b), legislator m sends a LastVote(b, v) to l where v is the vote with the largest ballot number less than b that m has cast, it is the latest vote from m before ballot b. A vote is composed of a ballot number, a decree and a voter. If m never voted before b we will consider both the ballot number and the decree of the vote to be -1.
Here, m has to take notes in his ledger to remember the votes he had cast. One thing that we must be careful of is that the set of ballot can change, since legislator l is going to use v as the latest vote from m until b. This means that upon sending a LastVote(b, v), legislator m MUST NOT cast any vote for ballots between v's ballot and b. Otherwise the value of v used by l will be outdated and we will break the consistency.
- 3: After receiving a LastVote(b, v) message from a majority set Q of legislators, legislator creates a new ballot with number b, quorum Q. The decree d is chosen to satisfy $B3(\beta)$, if no legislator in the quorum casted a vote before this (i.e. there are no successful ballots before b), we will make d being equal to b. It is simply an implementation choice, it could be anything but it is the thing we want to synchronize, so here, we will actually synchronize the number of the first successful ballot. l then record the new ballot in the notes of his ledger and send a BeginBallot(b, d) to every legislator in Q.
- 4: Upon receiving a BeginBallot(b, v), legislator m decides wether or not to cast his vote in ballot number b. He may refuse to do so if this violate a promise implied by a previous LastVote message. If he decides to vote, he sends a Voted(b, m) message to l and record his vote in his ledger. Note that legislator m has the option not to vote if he wants to. Every steps in this protocol are actually facultative, even if not responding to a message will surely prevent progress, it cannot break consistency as it cannot make $B1(\beta) - B3(\beta)$ false. This way, even if a legislator leaves the parliament in the middle of a ballot, it cannot break consistency.
- 5: if legislator l has received a Voted(b, m) message from every legislator m in Q, he writes d in his ledger and sends a Success(d) message to every legislators.
- 6: Upon receiving a Success(d) message, a legislator enters decree d in his ledger.

This is a lot to digest, take a deep breath, we will not dive into the implmentation of this right now, they are some restrictions we can apply in order to ease the implementation of this and I also think it is time to address some design again.

The multiple restrictions will result in the basic protocol so that's where we are going to talk about this.

Concerning design, we can start to think about what the messages we will send will look like. I decided to use a format really similar to HTTP, first because it allows me to reuse code from another project, but also because it is really easily scalable. A message will be composed of a method and a set of headers. To put it simply the method will be the type of the message and headers will be the parameters. A complete message will look like this:

```
BeginBallot
receiver: m
sender: 1
ballot: 4
decree: 2
```

So when I said headers were a set, you can obviously see that I lied, it is in fact a map, it will allow us to easily recover key information in the message. I will not dive into the parsing and the serialization of such messages because we do not care, just know that the end of the message is marked by an empty line.

Also note the receiver header, I put it in every message automatically upon sending it and its value is the name of the legislator who will receive the message, this is solely for debugging and logging purposes, to see where the message is actually going.

5.3 The basic protocol

The preliminary protocol require every legislator to record the number of every ballot he initiated, every vote he has cast and every LastVote message he has sent. That's a lot of information to keep track of without even considering how hard it might be to implement.

To answer this problem, the basic protocol is a restriction of the preliminary protocol.

This means that every action in the basic protocol is allowed by the preliminary protocol (the reverse is wrong), since the preliminary protocol maintains the consistency condition, so does the basic one.

For the basic protocol, every legislator l has to keep track in his ledger of:

- `lastTried[l]`: The number of the last ballot that l tried to initiate (we will take the liberty to put -1 if no such ballot exists as every ballot has a positive number in our implementation)
- `prevVote[l]`: The vote cast by p in the highest-numbered ballot in which he voted
- `nextBal[l]`: The largest value of b for which l has sent a `LastVote(b, v)`

The preliminary protocol allows legislator l to conduct any number of ballots concurrently, the basic protocol only allow one, he will ignore every message concerning a ballot different from the last he initiated.

In the preliminary protocol a `LastVote(b, v)` sent by m was the promise to not vote for any ballot between v 's ballot and b . In the basic protocol, it represents the stronger promise to not vote for any ballot numbered lower than b . This might prevent m from casting a vote in step 4 that he would have been allowed to cast in the preliminary protocol but as we saw, the preliminary protocol still allows m to not cast a vote even if he is allowed to so the basic protocol is does require m to perform an action forbidden by the preliminary protocol.

It would be pointless to rewrite the previous protocol and simply change some words so I will dive into the code right now, I recomand to take a look at the preliminary algorithm so you can get a better grasp of its progression.

I purposely removed some details in the code (logs, definitions of certain abstractions ...) to make the code lighter, certain abstractions have rather confusing implementations because I'm a terrible developer, but also because they contain optimizations that will be addressed in the complete protocol, yet I think the name of the fonctions will be good enough for you to understand their goal.

5.3.1 Step 1

```
1  void Legislator::initiate_ballot()
2  {
3      //this attribute keeps tracks of wether or not the ballot has begun
4      has_started = false;
5      //This attribute is a map of the previous votes of the quorum,
6      //it is updated after receiving a LastVote
7      quorum_previous_votes.clear();
8      int new_ballot_number = get_next_ballot_id();
9      //we set lastTried[l] to the newly initiated ballot
10     ledger.set_last_tried(new_ballot_number);
11     send_next_ballot(new_ballot_number);
12 }
13 void Legislator::send_next_ballot(int ballot)
14 {
15     std::string ballot_string = std::to_string(ballot);
16     Message message;
17     message.set_method("NextBallot");
18     message.add_header("ballot", ballot_string);
19     message.add_header("sender", self->config.name);
20     for (auto legislator : legislators)
21         message.send(legislator.second);
22 }
```


5.3.2 Step 2

```
1 //Step 2
2 void Legislator::receive_next_ballot(Message message)
3 {
4     std::string ballot_str = *message.get_header("ballot");
5     int ballot = std::stoi(ballot_str);
6     receive_next_ballot(ballot, *message.get_header("sender"));
7 }
8
9 void Legislator::receive_next_ballot(int ballot, std::string sender)
10 {
11     int next_ballot = ledger.next_bal(); //get nextBal[m]
12     if (ballot <= next_ballot) //ignore the message if the ballot is outdated
13         return;
14     ledger.set_next_bal(ballot); //Set nextBal[m] to b
15     Vote previous_vote = ledger.prev_vote(); //Get the last vote before b
16     send_last_vote(ballot, previous_vote, sender);
17 }
18
19 void Legislator::send_last_vote(int ballot, Vote previous_vote,
20     std::string sender)
21 {
22     std::string ballot_string = std::to_string(ballot);
23     std::string vote_ballot_id_string
24         = std::to_string(previous_vote.ballot_id);
25     std::string decree_string = std::to_string(previous_vote.decree.decree);
26     Message message;
27     message.set_method("LastVote");
28     message.add_header("ballot", ballot_string);
29     message.add_header("vote_ballot_id", vote_ballot_id_string);
30     message.add_header("decree", decree_string);
31     message.add_header("sender", self->config_.name);
32     message.send(legislators[sender]);
33 }
34 }
```

So here, legislator *m* receives a NextBallot, if he already responded to a NextBallot from a later ballot (whose ballot number is higher), *m* discards the message because the ballot is considered to be outdated, as a recall responding to a NextBallot(*b*) is a promise to never vote for a ballot earlier than *b*, it only makes sense to discard it now before even be able to vote.

If *b* is indeed a newer ballot, then *m* retrieves its last vote and sends it to legislator *l*.

5.3.3 Step 3

```
1 void Legislator::receive_last_vote(Message message)
2 {
3     std::string ballot_str = *message.get_header("ballot");
4     std::string vote_ballot_id_str = *message.get_header("vote_ballot_id");
5     std::string sender = *message.get_header("sender");
6     int ballot = std::stoi(ballot_str);
7     int vote_ballot_id = std::stoi(vote_ballot_id_str);
8     int vote_decree = std::stoi(*message.get_header("decree"));
9
10    if (ballot != ledger.last_tried() || has_started)
11        return;
12
13    Decree decree;
14    decree.decree = vote_decree;
15    Vote vote;
16    vote.decree = decree;
17    vote.legislator = sender;
18    vote.ballot_id = vote_ballot_id;
19
20    quorum_previous_votes.insert(std::pair<std::string, Vote>
21        (sender, vote));
22    unsigned int quorum_size = quorum_previous_votes.size();
23    unsigned int nb_legislators = legislators.size();
24    if (quorum_size > nb_legislators / 2)
25        receive_enough_last_vote();
26 }
27
28 void Legislator::receive_enough_last_vote()
29 {
30     has_started = true;
31     int ballot = ledger.last_tried();
32     std::string ballot_str = std::to_string(ballot);
33
34     Vote max_vote;
35     max_vote.ballot_id = -1;
36     for (auto legislator_vote_pair : quorum_previous_votes)
37     {
38         Vote current_vote = legislator_vote_pair.second;
39         if (current_vote.ballot_id > max_vote.ballot_id)
40             max_vote = current_vote;
41     }
42     Decree decree;
43     if (max_vote.ballot_id != -1)
44         decree = max_vote.decree;
45     else
46         decree.decree = ballot;
47     send_begin_ballot(ballot, decree);
48 }
```

There is actually a lot of things going on here, first I will not show the definition of `send_begin_ballot(int, Decree)`; as it is really straightforward, it is basically the same thing than `NextBallot` but we send it only to the quorum and that we also send a decree.

The first thing we check when we receive a `LastVote` is whether or not it is for the ballot we are currently conducting (this ballot is marked by `lastTried[]`).

We also need to check that we did not already start the ballot, we do not want to send a BeginBallot every time we receive a LastVote when you actually had a large enough quorum to send a BeginBallot beforehand, even though it doesn't impact consistency, it still add too much information in the logs.

After checking the ballot is valid, we register the vote in a map of the legislator names and their respective last vote. We then check the size of the map, if the number of elements is greater than half of the total number of legislators, it means that a majority of legislators responded with a LastVote, we can now enter the actual step 3: find the decree and send a BeginBallot.

We first mark the ballot has started, any incoming LastVote for this ballot will be discarded as the quorum just has been established.

We then simply check which vote was the latest among the last votes of the quorum and choose the corresponding decree, if no legislator previously voted we arbitrarily choose the ballot number as the new decree.

5.3.4 Step 4

```
1  void Legislator::receive_begin_ballot(Message message)
2  {
3      std::string ballot_string = *message.get_header("ballot");
4      std::string decree_string = *message.get_header("decree");
5      std::string sender = *message.get_header("sender");
6
7      int ballot = std::stoi(ballot_string);
8      int decree = std::stoi(decree_string);
9      receive_begin_ballot(ballot, decree, sender);
10 }
11
12 void Legislator::receive_begin_ballot(int ballot, int decree_id,
13     std::string sender)
14 {
15     int next_ballot = ledger.next_bal();
16     if (ballot != next_ballot)
17         return;
18     Vote vote;
19     vote.ballot_id = ballot;
20     Decree decree;
21     decree.decree = decree_id;
22     vote.decree = decree;
23     ledger.set_prev_vote(vote);
24
25     send_voted(ballot, decree, sender);
26 }
27
28 void Legislator::send_voted(int ballot, Decree decree, std::string receiver)
29 {
30     Message message;
31     message.set_method("Voted");
32     message.add_header("ballot", std::to_string(ballot));
33     message.add_header("sender", self->config_.name);
34     message.add_header("decree", std::to_string(decree.decree));
35     message.send(legislators[receiver]);
36 }
```

When receiving a BeginBallot(b, d), we first check if we should cast our vote in the ballot, indeed casting a vote might break a promise from an earlier LastVote, so we need to make sure that b is greater or equal to nextBal[m]. You might notice that we also ignore the message if nextBal[m] is greater than b, the reason is simple: this case is not supposed to ever happen, indeed m can receive a BeginBallot(b,

v) only if m is part of the quorum of b, which mean it sent a LastVote(b) previously which mean that nextBallot[m] was previously set to b (but can be bigger if an other ballot was initiated right after). Even though this case is supposed to never happen, we still handle the possibility by ignoring the message.

Once we checked that we could indeed cast the vote, we cast it and send Voted.

5.3.5 Step 5

```
1  void Legislator::receive_voted(Message message)
2  {
3      std::string ballot_str = *message.get_header("ballot");
4      std::string sender = *message.get_header("sender");
5      std::string decree_str = *message.get_header("decree");
6      Decree decree;
7      decree.decree = std::stoi(decree_str);
8
9      receive_voted(std::stoi(ballot_str), decree, sender);
10 }
11
12 void Legislator::receive_voted(int ballot, Decree decree, std::string voter)
13 {
14     if (ballot != ledger.last_tried())
15         return;
16     quorum_previous_votes.erase(voter);
17     if (quorum_previous_votes.size() == 0)
18         receive_enough_voted(ballot, decree);
19 }
20
21 void Legislator::receive_enough_voted(int ballot, Decree decree)
22 {
23     ledger.set_decree(decree);
24     send_success(decree);
25 }
26
27 void Legislator::send_success(Decree decree)
28 {
29     Message message;
30     message.set_method("Success");
31     message.add_header("decree", std::to_string(decree.decree));
32
33     for (auto legislator : legislators)
34         message.send(legislator.second);
35 }
```

When we receive a Voted message for the good ballot, we erase the voter from the quorum, once the quorum is empty, it means that all the quorum has voted, and we can send a Success.

5.3.6 Step 6

```
1 void Legislator::receive_success(Message message)
2 {
3     std::string decree_str = *message.get_header("decree");
4     Decree decree;
5     decree.decree = std::stoi(decree_str);
6     receive_success(decree);
7 }
8
9 void Legislator::receive_success(Decree decree)
10 {
11     ledger.set_decree(decree);
12 }
```

Nothing surprising here.

The basic protocol derive directly from the preliminary protocol and maintains the consistency condition, however we have no guarantee of the progress condition as no legislator is forced to ever initiate a ballot.

The complete protocol, derived from the basic protocol, addresses this problem.

5.4 The complete protocol

The progress condition requires that a legislator eventually performs step 1 but we still need to think of the frequency of the ballots. Indeed, initiating too many ballots can prevent progress. If legislator 1 initiates a ballot number b larger than any other ballot, the $\text{LastVote}(b, v)$ from legislator m in step 2 will prevent m from voting in any previously initiated ballot. If new ballots continually initiated with increasing ballot numbers before the previous ballots have a chance to succeed, then no progress can be made.

The main idea introduced by the complete protocol is that only one legislator can initiate ballot, and that only he can decide of the appropriate time to do so.

This legislator is called the president. The election process is irrelevant but it must occur frequently to make sure the president didn't leave the parliament, we obviously need to change the president in such cases.

So as I mentionned earlier the complete protocol is in fact irrelevant in our case, the reason being that for visualisation purposes, I decided to let the ballot initiation process between the hands of the user, to initiate a ballot, you must send a SIGTSTP to the process of the server.

The problem of president and ballot frequency are now completely in the hands of the user who is totally able to simulate the complete protocol by hand.

However, the complete protocol still addresses an interesting optimization regarding presidency change. (in our case, when the user will send ballots from multiple legislators).

This is linked to the way we choose the next ballot number, the basic protocol specifies that legislator 1 must initiate a ballot with a number higher than $\text{lastTried}[1]$, however, by doing so we have no guarantee that this will be higher than an already passed successful ballot that was initiated by someone else.

To handle the latter possibility, legislator 1 had to learn about any ballot number higher than $\text{lastTried}[1]$ that was conducted before he became president.

At the start of the basic protocol I promised an optimization and this what this is all about. Instead of simply taking a ballot number higher than $\text{lastTried}[1]$, we actually look if an higher numbered ballot as been initiated by someone else, this ballot number will be stored in $\text{nextBal}[1]$.

```

1  int Legislator::get_next_ballot_id()
2  {
3      int previous_ballot_id = ledger.last_tried();
4      if (previous_ballot_id == -1)
5          previous_ballot_id = config.ballot_partition_id - legislators.size();
6      int last_voted_ballot = ledger.next_bal();
7      while (last_voted_ballot > (int)legislators.size() + previous_ballot_id)
8          previous_ballot_id += legislators.size();
9      return previous_ballot_id + legislators.size();
10 }

```

This is only a simple optimization not presented in the original algorithm, the reason it is not presented is because it actually is useless because one problem remains et that response to this problem also handle the above case, but yet, this optimization allow to save a few message sending.

So the problem that remains is that legislator m was conduction ballots while legislator l was outside of the parliament, then legislator l would not have any way to know these ballots in his personal ledger. If its first action after entering the parliament (before receiving any message) was to initiate decree, he will have a lot of troubles to find a suitable ballot number.

To solve this, the complete protocol proposes to send a message in response to a NextBallot(b) or BeginBallot(b, d) instead of ignoring it when b is inferior or equal to nextBal[m].

This message contains the value of nextBal[m], legislator l, can then update its value of lastTried[l] before re-initiating a ballot.

```

1  //Is called in receive_next_ballot and receive_begin_ballot
2  //instead of ignoring the message
3  void Legislator::send_higher_ballot(int ballot, std::string receiver)
4  {
5      Message message;
6      message.set_method("HigherBallot");
7      message.add_header("ballot", std::to_string(ballot));
8      message.send(legislators[receiver]);
9  }
10
11 void Legislator::receive_higher_ballot(Message message)
12 {
13     std::string ballot_str = *message.get_header("ballot");
14     receive_higher_ballot(std::stoi(ballot_str));
15 }
16
17 void Legislator::receive_higher_ballot(int ballot)
18 {
19     int last_tried = ledger.last_tried();
20     if (last_tried > ballot)
21         return;
22
23     while (last_tried + (int)legislators.size() < ballot)
24         last_tried += legislators.size();
25
26     ledger.set_last_tried(last_tried);
27     initiate_ballot();
28 }

```

6 Conclusion and opening on the multi-decree parliament

As I said the implementation is incomplete, I presented here the complete single-decree synod protocol. However the original paper presents the multi-decree parliament protocol which directly derives from the single-decree. Its goal is of course to synchronize not one but multiple decrees between legislators.

This have obviously much more applications than the single-decree protocol but I did not implement it yet, I do not want to risk myself continuing this article on a subject I did not experience myself, but maybe you will be able to understand it by yourself by reading the original paper.

I hope this article gave you the keys you needed to understand the Paxos algorithm with a more pragmatic and intuitive approach.

I will probably implement the multi-decree parliament some day on the same repository, but for now I will have to let you search for it, but do not be scared, it is in fact fairly simple as, I quote,

"The Paxons never bothered writing a precise description of the parliamentary protocol because it was so easily derived from the Synod protocol."