

Chapter 4

Authenticated Agreement

Byzantine nodes are able to lie about their inputs as well as received messages. Can we detect certain lies and limit the power of byzantine nodes? Possibly, the authenticity of messages may be validated using signatures?

4.1 Agreement with Authentication

Definition 4.1 (Signature). *If a node never signs a message, then no correct node ever accepts that message. We denote a message $\text{msg}(x)$ signed by node u with $\text{msg}(x)_u$.*

Remarks:

- Algorithm 4.2 shows an agreement protocol for binary inputs relying on signatures. We assume there is a designated “primary” node p . The goal is to decide on p 's value.

Algorithm 4.2 Byzantine Agreement with Authentication

Code for primary p :

```
1: if input is 1 then  
2:   broadcast  $\text{value}(1)_p$   
3:   decide 1 and terminate  
4: else  
5:   decide 0 and terminate  
6: end if
```

Code for all other nodes v :

```
7: for all rounds  $i \in 1, \dots, f + 1$  do  
8:    $S$  is the set of accepted messages  $\text{value}(1)_u$ .  
9:   if  $|S| \geq i$  and  $\text{value}(1)_p \in S$  then  
10:    broadcast  $S \cup \{\text{value}(1)_v\}$   
11:    decide 1 and terminate  
12:   end if  
13: end for  
14: decide 0 and terminate
```

Theorem 4.3. *Algorithm 4.2 can tolerate $f < n$ byzantine failures while terminating in $f + 1$ rounds.*

Proof. Assuming that the primary p is not byzantine and its input is 1, then p broadcasts $\text{value}(1)_p$ in the first round, which will trigger all correct nodes to decide for 1. If p 's input is 0, there is no signed message $\text{value}(1)_p$, and no node can decide for 1.

If primary p is byzantine, we need all correct nodes to decide for the same value for the algorithm to be correct. Let us assume that p convinces a correct node v that its value is 1 in round i with $i < f + 1$. We know that v received i signed messages for value 1. Then, v will broadcast $i + 1$ signed messages for value 1, which will trigger all correct nodes to also decide for 1. If p tries to convince some node v late (in round $i = f + 1$), v must receive $f + 1$ signed messages. Since at most f nodes are byzantine, at least one correct node u signed a message $\text{value}(1)_u$ in some round $i < f + 1$, which puts us back to the previous case. \square

Remarks:

- The algorithm only takes $f + 1$ rounds, which is optimal as described in Theorem 3.20.
- Using signatures, Algorithm 4.2 solves consensus for any number of failures! Does this contradict Theorem 3.12? Recall that in the proof of Theorem 3.12 we assumed that a byzantine node can distribute contradictory information about its own input. If messages are signed, correct nodes can detect such behavior – a node u signing two contradicting messages proves to all nodes that node u is byzantine.
- Does Algorithm 4.2 satisfy any of the validity conditions introduced in Section 3.1? No! A byzantine primary can dictate the decision value. Can we modify the algorithm such that the correct-input validity condition is satisfied? Yes! We can run the algorithm in parallel for $2f + 1$ primary nodes. Either 0 or 1 will occur at least $f + 1$ times, which means that one correct process had to have this value in the first place. In this case, we can only handle $f < \frac{n}{2}$ byzantine nodes.
- In reality, a primary will usually be correct. If so, Algorithm 4.2 only needs two rounds! Can we make it work with arbitrary inputs? Also, relying on synchrony limits the practicality of the protocol. What if messages can be lost or the system is asynchronous?
- Zyzzyva uses authenticated messages to achieve state replication, as in Definition 1.8. It is designed to run fast when nodes run correctly, and it will slow down to fix failures!

4.2 Zyzzyva

Definition 4.4 (View). *A view V describes the current state of a replicated system, enumerating the $3f + 1$ replicas. The view V also marks one of the replicas as the primary p .*

Definition 4.5 (Command). *If a client wants to update (or read) data, it sends a suitable command c in a **Request** message to the primary p . Apart from the command c itself, the **Request** message also includes a timestamp t . The client signs the message to guarantee authenticity.*

Definition 4.6 (History). *The history h is a sequence of commands c_1, c_2, \dots in the order they are executed by Zyzzyva. We denote the history up to c_k with h_k .*

Remarks:

- In Zyzzyva, the primary p is used to order commands submitted by clients to create a history h .
- Apart from the globally accepted history, node u may also have a local history, which we denote as h^u or h_k^u .

Definition 4.7 (Complete command). *If a command completes, it will remain in its place in the history h even in the presence of failures.*

Remarks:

- As long as clients wait for the completion of their commands, clients can treat Zyzzyva like one single computer even if there are up to f failures.

In the Absence of Failures

Algorithm 4.8 Zyzzyva: No failures

- 1: At time t client u wants to execute command c
 - 2: Client u sends request $\mathbf{R} = \mathbf{Request}(c, t)_u$ to primary p
 - 3: Primary p appends c to its local history, i.e., $h^p = (h^p, c)$
 - 4: Primary p sends $\mathbf{OR} = \mathbf{OrderedRequest}(h^p, c, \mathbf{R})_p$ to all replicas
 - 5: Each replica r appends command c to local history $h^r = (h^r, c)$ and checks whether $h^r = h^p$
 - 6: Each replica r runs command c_k and obtains result a
 - 7: Each replica r sends $\mathbf{Response}(a, \mathbf{OR})_r$ to client u
 - 8: Client u collects the set S of received $\mathbf{Response}(a, \mathbf{OR})_r$ messages
 - 9: Client u checks if all histories h^r are consistent
 - 10: **if** $|S| = 3f + 1$ **then**
 - 11: Client u considers command c to be complete
 - 12: **end if**
-

Remarks:

- Since the client receives $3f + 1$ consistent responses, all correct replicas have to be in the same state.
- Only three communication rounds are required for the command c to complete.

- Note that replicas have no idea which commands are considered complete by clients! How can we make sure that commands that are considered complete by a client are actually executed? We will see in Theorem 4.23.
- Commands received from clients should be ordered according to timestamps to preserve the causal order of commands.
- There is a lot of optimization potential. For example, including the entire command history in most messages introduces prohibitively large overhead. Rather, old parts of the history that are agreed upon can be truncated. Also, sending a hash value of the remainder of the history is enough to check its consistency across replicas.
- What if a client does not receive $3f + 1$ $\text{Response}(a, \text{OR})_r$ messages? A byzantine replica may omit sending anything at all! In practice, clients set a timeout for the collection of Response messages. Does this mean that Zyzyyva only works in the synchronous model? Yes and no. We will discuss this in Lemma 4.26 and Lemma 4.27.

Byzantine Replicas

Algorithm 4.9 Zyzyyva: Byzantine Replicas (append to Algorithm 4.8)

- 1: **if** $2f + 1 \leq |S| < 3f + 1$ **then**
 - 2: Client u sends $\text{Commit}(S)_u$ to all replicas
 - 3: Each replica r replies with a $\text{LocalCommit}(S)_r$ message to u
 - 4: Client u collects at least $2f + 1$ $\text{LocalCommit}(S)_r$ messages and considers c to be complete
 - 5: **end if**
-

Remarks:

- If replicas fail, a client u may receive less than $3f + 1$ consistent responses from the replicas. Client u can only assume command c to be complete if all correct replicas r eventually append command c to their local history h^r .

Definition 4.10 (Commit Certificate). *A commit certificate S contains $2f + 1$ consistent and signed $\text{Response}(a, \text{OR})_r$ messages from $2f + 1$ different replicas r .*

Remarks:

- The set S is a commit certificate which proves the execution of the command on $2f + 1$ replicas, of which at least $f + 1$ are correct. This commit certificate S must be acknowledged by $2f + 1$ replicas before the client considers the command to be complete.
- Why do clients have to distribute this commit certificate to $2f + 1$ replicas? We will discuss this in Theorem 4.21.

- What if $|S| < 2f + 1$, or what if the client receives $2f + 1$ messages but some have inconsistent histories? Since at most f replicas are byzantine, the primary itself must be byzantine! Can we resolve this?

Byzantine Primary

Definition 4.11 (Proof of Misbehavior). *Proof of misbehavior of some node can be established by a set of contradicting signed messages.*

Remarks:

- For example, if a client u receives two $\text{Response}(a, \text{OR})_r$ messages that contain inconsistent OR messages signed by the primary, client u can prove that the primary misbehaved. Client u broadcasts this proof of misbehavior to all replicas r which initiate a view change by broadcasting a IHatePrimary_r message to all replicas.

Algorithm 4.12 Zyzyzyva: Byzantine Primary (append to Algorithm 4.9)

```

1: if  $|S| < 2f + 1$  then
2:   Client  $u$  sends the original  $R = \text{Request}(c, t)_u$  to all replicas
3:   Each replica  $r$  sends a  $\text{ConfirmRequest}(R)_r$  message to  $p$ 
4:   if primary  $p$  replies with OR then
5:     Replica  $r$  forwards OR to all replicas
6:     Continue as in Algorithm 4.8, Line 5
7:   else
8:     Replica  $r$  initiates view change by broadcasting  $\text{IHatePrimary}_r$  to all
       replicas
9:   end if
10: end if

```

Remarks:

- A faulty primary can slow down Zyzyzyva by not sending out the OrderedRequest messages in Algorithm 4.8, repeatedly escalating to Algorithm 4.12.
- Line 5 in the Algorithm is necessary to ensure liveness. We will discuss this in Theorem 4.27.
- Again, there is potential for optimization. For example, a replica might already know about a command that is requested by a client. In that case, it can answer without asking the primary. Furthermore, the primary might already know the message R requested by the replicas. In that case, it sends the old OR message to the requesting replica.

Safety

Definition 4.13 (Safety). *We call a system safe if the following condition holds: If a command with sequence number j and a history h_j completes, then for any command that completed earlier (with a smaller sequence number $i < j$), the history h_i is a prefix of history h_j .*

Remarks:

- In Zyzzyva a command can only complete in two ways, either in Algorithm 4.8 or in Algorithm 4.9.
- If a system is safe, complete commands cannot be reordered or dropped. So is Zyzzyva so far safe?

Lemma 4.14. *Let c_i and c_j be two different complete commands. Then c_i and c_j must have different sequence numbers.*

Proof. If a command c completes in Algorithm 4.8, $3f + 1$ replicas sent a $\text{Response}(a, \text{OR})_r$ to the client. If the command c completed in Algorithm 4.9, at least $2f + 1$ replicas sent a $\text{Response}(a, \text{OR})_r$ message to the client. Hence, a client has to receive at least $2f + 1$ $\text{Response}(a, \text{OR})_r$ messages.

Both c_i and c_j are complete. Therefore there must be at least $2f + 1$ replicas that responded to c_i with a $\text{Response}(a, \text{OR})_r$ message. But there are also at least $2f + 1$ replicas that responded to c_j with a $\text{Response}(a, \text{OR})_r$ message. Because there are only $3f + 1$ replicas, there is at least one correct replica that sent a $\text{Response}(a, \text{OR})_r$ message for both c_i and c_j . A correct replica only sends one $\text{Response}(a, \text{OR})_r$ message for each sequence number, hence the two commands must have different sequence numbers. \square

Lemma 4.15. *Let c_i and c_j be two complete commands with sequence numbers $i < j$. The history h_i is a prefix of h_j .*

Proof. As in the proof of Lemma 4.14, there has to be at least one correct replica that sent a $\text{Response}(a, \text{OR})_r$ message for both c_i and c_j .

A correct replica r that sent a $\text{Response}(a, \text{OR})_r$ message for c_i will only accept c_j if the history for c_j provided by the primary is consistent with the local history of replica r , including c_i . \square

Remarks:

- A byzantine primary can cause the system to never complete any command. Either by never sending any messages or by inconsistently ordering client requests. In this case, replicas have to replace the primary.

View Changes

Definition 4.16 (View Change). *In Zyzzyva, a view change is used to replace a byzantine primary with another (hopefully correct) replica. View changes are initiated by replicas sending IHatePrimary_r to all other replicas. This only happens if a replica obtains a valid proof of misbehavior from a client or after a replica fails to obtain an OR message from the primary in Algorithm 4.12.*

Remarks:

- How can we safely decide to initiate a view change, i.e. demote a byzantine primary? Note that byzantine nodes should not be able to trigger a view change!

Algorithm 4.17 Zyzzyva: View Change Agreement

- 1: All replicas continuously collect the set H of `IHatePrimaryr` messages
 - 2: **if** a replica r received $|H| > f$ messages or a valid `ViewChange` message **then**
 - 3: Replica r broadcasts `ViewChange(H^r, h^r, S_l^r)r`
 - 4: Replica r stops participating in the current view
 - 5: Replica r switches to the next primary " $p = p + 1$ "
 - 6: **end if**
-

Remarks:

- The $f + 1$ `IHatePrimaryr` messages in set H prove that at least one correct replica initiated a view change. This proof is broadcast to all replicas to make sure that once the first correct replica stopped acting in the current view, all other replicas will do so as well.
- S_l^r is the most recent commit certificate that the replica obtained in the ending view as described in Algorithm 4.9. S_l^r will be used to recover the correct history before the new view starts. The local histories h^r are included in the `ViewChange(H^r, h^r, S_l^r)r` message such that commands that completed after a correct client received $3f + 1$ responses from replicas can be recovered as well.
- In Zyzzyva, a byzantine primary starts acting as a normal replica after a view change. In practice, all machines eventually break and rarely fix themselves after that. Instead, one could consider to replace a byzantine primary with a fresh replica that was not in the previous view.

Algorithm 4.18 Zyzzyva: View Change Execution

- 1: The new primary p collects the set C of `ViewChange(H^r, h^r, S_l^r)r` messages
 - 2: **if** new primary p collected $|C| \geq 2f + 1$ messages **then**
 - 3: New primary p sends `NewView(C)p` to all replicas
 - 4: **end if**
 - 5: **if** a replica r received a `NewView(C)p` message **then**
 - 6: Replica r recovers new history h_{new} as shown in Algorithm 4.20
 - 7: Replica r broadcasts `ViewConfirm(h_{new})r` message to all replicas
 - 8: **end if**
 - 9: **if** a replica r received $2f + 1$ `ViewConfirm(h_{new})r` messages **then**
 - 10: Replica r accepts $h^r = h_{\text{new}}$ as the history of the new view
 - 11: Replica r starts participating in the new view
 - 12: **end if**
-

Remarks:

- Analogously to Lemma 4.15, commit certificates are ordered. For two commit certificates S_i and S_j with sequence numbers $i < j$, the history h_i certified by S_i is a prefix of the history h_j certified by S_j .
- Zyzyyva collects the most recent commit certificate and the local history of $2f + 1$ replicas. This information is distributed to all replicas, and used to recover the history for the new view h_{new} .
- If a replica does not receive the $\text{NewView}(C)_p$ or the $\text{ViewConfirm}(h_{new})_r$ message in time, it triggers another view change by broadcasting IHatePrimary_r to all other replicas.
- How is the history recovered exactly? It seems that the set of histories included in C can be messy. How can we be sure that complete commands are not reordered or dropped?

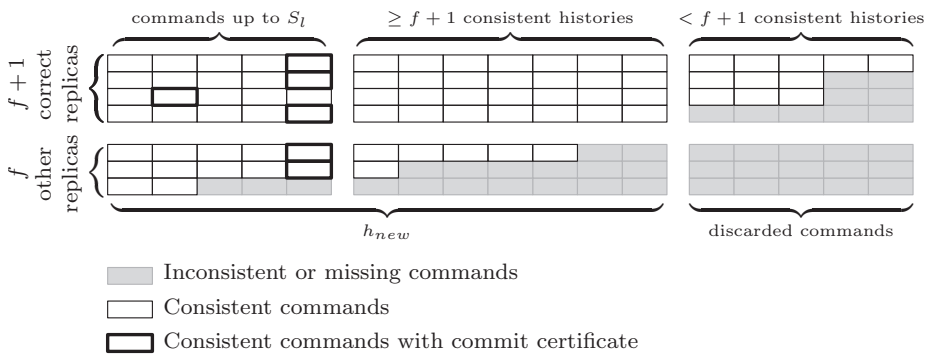


Figure 4.19: The structure of the data reported by different replicas in C . Commands up to the last commit certificate S_l were completed in either Algorithm 4.8 or Algorithm 4.9. After the last commit certificate S_l there may be commands that completed at a correct client in Algorithm 4.8. Algorithm 4.20 shows how the new history h_{new} is recovered such that no complete commands are lost.

Algorithm 4.20 *Zyzyva*: History Recovery

```

1:  $C$  = set of  $2f + 1$   $\text{ViewChange}(H^r, h^r, S^r)_r$  messages in  $\text{NewView}(C)_p$ 
2:  $R$  = set of replicas included in  $C$ 
3:  $S_l$  = most recent commit certificate  $S_l^r$  reported in  $C$ 
4:  $h_{\text{new}}$  = history  $h_l$  contained in  $S_l$ 
5:  $k = l + 1$ , next sequence number
6: while command  $c_k$  exists in  $C$  do
7:   if  $c_k$  is reported by at least  $f + 1$  replicas in  $R$  then
8:     Remove replicas from  $R$  that do not support  $c_k$ 
9:      $h_{\text{new}} = (h_{\text{new}}, c_k)$ 
10:  end if
11:   $k = k + 1$ 
12: end while
13: return  $h_{\text{new}}$ 

```

Remarks:

- Commands up to S_l are included into the new history h_{new} .
- If at least $f + 1$ replicas share a consistent history after the last commit certificate S_l , also the commands after that are included.
- Even if $f + 1$ correct replicas consistently report a command c after the last commit certificate S_l , c may not be considered complete by a client, e.g., because one of the responses to the client was lost. Such a command is included in the new history h_{new} . When the client retries executing c , the replicas will be able to identify the same command c using the timestamp included in the client's request, and avoid duplicate execution of the command.
- Can we be sure that all commands that completed at a correct client are carried over into the new view?

Lemma 4.21. *The globally most recent commit certificate S_l is included in C .*

Proof. Any two sets of $2f + 1$ replicas share at least one correct replica. Hence, at least one correct replica which acknowledged the most recent commit certificate S_l also sent a $\text{LocalCommit}(S_l)_r$ message that is in C . \square

Lemma 4.22. *Any command and its history that completes after S_l has to be reported in C at least $f + 1$ times.*

Proof. A command c can only complete in Algorithm 4.8 after S_l . Hence, $3f + 1$ replicas sent a $\text{Response}(a, \text{OR})_r$ message for c . C includes the local histories of $2f + 1$ replicas of which at most f are byzantine. As a result, c and its history is consistently found in at least $f + 1$ local histories in C . \square

Lemma 4.23. *If a command c is considered complete by a client, command c remains in its place in the history during view changes.*

Proof. We have shown in Lemma 4.21 that the most recent commit certificate is contained in C , and hence any command that terminated in Algorithm 4.9

is included in the new history after a view change. Every command that completed before the last commit certificate S_l is included in the history as a result. Commands that completed in Algorithm 4.8 after the last commit certificate are supported by at least $f + 1$ correct replicas as shown in Lemma 4.22. Such commands are added to the new history as described in Algorithm 4.20. Algorithm 4.20 adds commands sequentially until the histories become inconsistent. Hence, complete commands are not lost or reordered during a view change. \square

Theorem 4.24. *Zyzyva is safe even during view changes.*

Proof. Complete commands are not reordered within a view as described in Lemma 4.15. Also, no complete command is lost or reordered during a view change as shown in Lemma 4.23. Hence, Zyzyva is safe. \square

Remarks:

- So Zyzyva correctly handles complete commands even in the presence of failures. We also want Zyzyva to make progress, i.e., commands issued by correct clients should complete eventually.
- If the network is broken or introduces arbitrarily large delays, commands may never complete.
- Can we be sure commands complete in periods in which delays are bounded?

Definition 4.25 (Liveness). *We call a system **live** if every command eventually completes.*

Lemma 4.26. *Zyzyva is live during periods of synchrony if the primary is correct and a command is requested by a correct client.*

Proof. The client receives a $\text{Response}(a, \text{OR})_r$ message from all correct replicas. If it receives $3f + 1$ messages, the command completes immediately in Algorithm 4.8. If the client receives fewer than $3f + 1$ messages, it will at least receive $2f + 1$, since there are at most f byzantine replicas. All correct replicas will answer the client's $\text{Commit}(S)_u$ message with a correct $\text{LocalCommit}(S)_r$ message after which the command completes in Algorithm 4.9. \square

Lemma 4.27. *If, during a period of synchrony, a request does not complete in Algorithm 4.8 or Algorithm 4.9, a view change occurs.*

Proof. If a command does not complete for a sufficiently long time, the client will resend the $R = \text{Request}(c, t)_u$ message to all replicas. After that, if a replica's $\text{ConfirmRequest}(R)_r$ message is not answered in time by the primary, it broadcasts an IHatePrimary_r message. If a correct replica gathers $f + 1$ IHatePrimary_r messages, the view change is initiated. If no correct replica collects more than f IHatePrimary_r messages, at least one correct replica received a valid $\text{OrderedRequest}(h^p, c, R)_p$ message from the primary which it forwards to all other replicas. In that case, the client is guaranteed to receive at least $2f + 1$ $\text{Response}(a, \text{OR})_r$ messages from the correct replicas and can complete the command by assembling a commit certificate. \square

Remarks:

- If the newly elected primary is byzantine, the view change may never terminate. However, we can detect if the new primary does not assemble C correctly as all contained messages are signed. If the primary refuses to assemble C , replicas initiate another view change after a timeout.

Chapter Notes

Algorithm 4.2 was introduced by Dolev et al. [DFF⁺82] in 1982. Byzantine fault tolerant state machine replication (BFT) is a problem that gave rise to various protocols. Castro and Liskov [MC99] introduced the Practical Byzantine Fault Tolerance (PBFT) protocol in 1999, applications such as Farsite [ABC⁺02] followed. This triggered the development of, e.g., Q/U [AEMGG⁺05] and HQ [CML⁺06]. Zyzzyva [KAD⁺07] improved on performance especially in the case of no failures, while Aardvark [CWA⁺09] improved performance in the presence of failures. Guerraoui et al. [GKQV10] introduced a modular system which allows to more easily develop BFT protocols that match specific applications in terms of robustness or best case performance.

This chapter was written in collaboration with Pascal Bissig.

Bibliography

- [ABC⁺02] Atul Adya, William J. Bolosky, Miguel Castro, Gerald Cermak, Ronnie Chaiken, John R. Douceur, Jon Howell, Jacob R. Lorch, Marvin Theimer, and Roger P. Wattenhofer. Farsite: Federated, available, and reliable storage for an incompletely trusted environment. *SIGOPS Oper. Syst. Rev.*, 36(SI):1–14, December 2002.
- [AEMGG⁺05] Michael Abd-El-Malek, Gregory R Ganger, Garth R Goodson, Michael K Reiter, and Jay J Wylie. Fault-scalable byzantine fault-tolerant services. *ACM SIGOPS Operating Systems Review*, 39(5):59–74, 2005.
- [CML⁺06] James Cowling, Daniel Myers, Barbara Liskov, Rodrigo Rodrigues, and Liuba Shrira. Hq replication: A hybrid quorum protocol for byzantine fault tolerance. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation, OSDI '06*, pages 177–190, Berkeley, CA, USA, 2006. USENIX Association.
- [CWA⁺09] Allen Clement, Edmund L Wong, Lorenzo Alvisi, Michael Dahlin, and Mirco Marchetti. Making byzantine fault tolerant systems tolerate byzantine faults. In *NSDI*, volume 9, pages 153–168, 2009.
- [DFF⁺82] Danny Dolev, Michael J Fischer, Rob Fowler, Nancy A Lynch, and H Raymond Strong. An efficient algorithm for byzantine

- agreement without authentication. *Information and Control*, 52(3):257–274, 1982.
- [GKQV10] Rachid Guerraoui, Nikola Knežević, Vivien Quéma, and Marko Vukolić. The next 700 bft protocols. In *Proceedings of the 5th European conference on Computer systems*, pages 363–376. ACM, 2010.
- [KAD⁺07] Ramakrishna Kotla, Lorenzo Alvisi, Mike Dahlin, Allen Clement, and Edmund Wong. Zyzzyva: speculative byzantine fault tolerance. In *ACM SIGOPS Operating Systems Review*, volume 41, pages 45–58. ACM, 2007.
- [MC99] Barbara Liskov Miguel Castro. Practical byzantine fault tolerance. In *OSDI*, volume 99, pages 173–186, 1999.