

A Leader-Free Byzantine Consensus Algorithm^{*}

Fatemeh Borran and André Schiper

Ecole Polytechnique Fédérale de Lausanne (EPFL), Switzerland

Abstract. The paper considers the consensus problem in a partially synchronous system with Byzantine faults. It turns out that, in the partially synchronous system, all deterministic algorithms that solve consensus with Byzantine faults are leader-based. This is not the case of benign faults, which raises the following fundamental question: is it possible to design a deterministic Byzantine consensus algorithm for a partially synchronous system that is not leader-based? The paper gives a positive answer to this question, and presents a leader-free algorithm that is resilient-optimal and signature-free.

1 Introduction

In a distributed system of n processes, where each process has an initial value, Byzantine consensus is the problem of agreeing on a common value, even though some of the processes may fail in arbitrary, even malicious, ways. Consensus is related to the implementation of state machine replication, atomic broadcast, etc. It was first identified by Pease, Shostak and Lamport [1], formalized as the *interactive consistency* problem and solved in a synchronous system. An algorithm achieves interactive consistency if it allows the nonfaulty processes to come to a consistent view of the initial values of all the processes, including the faulty ones. Once interactive consistency has been achieved, the nonfaulty processes can reach consensus by applying a deterministic averaging or filtering function on the values of their view. It is shown in [1] that in a synchronous system $3t + 1$ processes are needed to solve the Byzantine consensus problem without signatures, where t is the maximum number of Byzantine processes.

Later, Fischer, Lynch and Peterson [2] proved that in an asynchronous system, no deterministic asynchronous consensus protocol can tolerate even a single crash failure. The problem can however be solved using randomization even with Byzantine faults, with at least $5t + 1$ processes, as shown by BenOr [3] and Rabin [4]. Later, Bracha [5] increased the resiliency of the randomized algorithm to $3t + 1$ using a “reliable broadcast” primitive.

In 1988, Dwork, Lynch and Stockmeyer [6], considered an asynchronous system that eventually becomes synchronous (called *partially synchronous system*). The consensus algorithms proposed in [6], for benign and for Byzantine faults, achieve safety in all executions, while guaranteeing liveness only if there exists

^{*} Research funded by the Swiss National Science Foundation under grant number 200021-111701.

a period of synchrony. Recently, several papers have considered the partially synchronous system model for Byzantine consensus [7,8,9,10].

However, [10] points out a potential weakness of these Byzantine consensus algorithms, namely that they can suffer from “performance failure”. According to [10], a performance failure occurs when messages are sent slowly by a Byzantine leader, but without triggering protocol timeouts, and the paper points out that the PBFT leader-based algorithm [7] is vulnerable to such an attack. Similar arguments are mentioned in [11] and in [12], where Lamport suggests the use of a virtual leader.

Interestingly, all deterministic Byzantine algorithms for non-synchronous systems are leader-based, e.g., [6,7,8,9]. Even the protocol in [10] is leader-based. However, the authors of [10] managed to make the leader-based protocol less vulnerable to performance failure attacks than PBFT [7] through a complicated mechanism that enables non-leader processes to (i) aggressively monitor the leader’s performance, and (ii) compute a threshold level of acceptable performance. Note that randomized consensus algorithms such as [3,4] are not leader-based. This raises the following fundamental question: is it possible to design a deterministic Byzantine consensus algorithm for a partially synchronous system that is not leader-based? With such an algorithm, performance failure of Byzantine processes might be harmless.

One may imagine that leader-free (non-leader-based) algorithms for benign faults might be extended for Byzantine faults. A leader-free algorithm typically consists of a sequence of rounds, where in each round all processes send messages to all, and a correct process updates its value based on the values received. It is not difficult to design an algorithm based on this all-to-all communication pattern that does not violate the validity and agreement properties of consensus, even with Byzantine faults. However, termination requires that in some round r all correct processes receive exactly the same set of messages (from correct *and* from faulty processes). Let us denote this property for round r by $uniform(r)$. Indeed, if $uniform(r)$ holds and each correct process applies a deterministic function to the received values, the configuration becomes univalent. Can we ensure the existence of a round r in which $uniform(r)$ holds?

For benign faults, it is easy to guarantee that during the synchronous period of the partially synchronous system, in every round r , all correct processes receive messages from the same set of processes. This is not the case for Byzantine faults. In round r , a Byzantine process could send a message to some correct process, and no message to some other correct process. If this happens, $uniform(r)$ does not hold. Therefore one may think that with Byzantine faults the leader is needed to ensure termination, and conclude that no deterministic leader-free Byzantine consensus algorithm could exist in a partially synchronous system. In this paper, we show that this intuition is wrong.

Our new idea is the following. We started from the observation that leader-free consensus algorithms exist for the synchronous system, both for benign faults (e.g., the *FloodSet* algorithm [13]) and for Byzantine faults (e.g., the algorithm based on interactive consistency [1]). However, these algorithms violate

agreement if executed during the asynchronous period of a partially synchronous system. Therefore we tried to combine these algorithms with a second algorithm that never violates agreement in an asynchronous system. This methodology turned out to be successful, and the resulting leader-free Byzantine consensus algorithm, is presented here. The algorithm requires $3t + 1$ processes and does not rely on digital signatures.

The rest of the paper is structured as follows. We define the consensus problem and our system model in Section 2. Our methodology to derive a leader-free consensus algorithm for Byzantine faults is presented in Section 3.¹ Future work is discussed in Section 4. Finally, we conclude the paper in Section 5.

2 Definitions and System Model

2.1 Byzantine Consensus

We consider a set Π of n processes, among which at most t can be Byzantine faulty. Nonfaulty processes are called correct processes. Each process has an initial value. We formally define consensus by the following properties:

- *Strong validity*: If all correct processes have the same initial value, this is the only possible decision value.
- *Agreement*: No two correct processes decide differently.
- *Termination*: All correct processes eventually decide.

2.2 System Model

We consider a partially synchronous system as defined in [6] in which processes communicate through message passing. As in [6], we consider an abstraction on top of the system model, namely a round model, defined next. Using this abstraction, rather than the raw system model, improves the clarity of the algorithms and simplifies the proofs.

There are two fault models considered with Byzantine processes: “authenticated Byzantine” faults, and “Byzantine” faults [6]. In both models a faulty process behaves arbitrarily, but in the authenticated Byzantine model messages can be signed by the sender, and it is assumed that the signature cannot be forged by any other process. No signatures are used with Byzantine faults, but the receiver of a message knows the identity of the sender.

2.3 Basic Round Model

In the round model, processing is divided into rounds of message exchange. Each round r consists of a *sending step* denoted by S_p^r (sending step of p for round r), and of a *state transition step* denoted by T_p^r . In a sending step, each process sends a message to all. A subset of the messages sent is received at the beginning of the state transition step: messages can get lost, and a message sent in round

¹ A simpler algorithm that uses digital signatures is proposed in [14].

r can only be received in round r . We denote by σ_p^r the message sent by p in round r , and by μ_p^r the messages received by process p in round r (μ_p^r is a vector of size n , where $\mu_p^r[q]$ is the message received from q). Based on μ_p^r , process p updates its state in the state transition step.

Let *GSR* (*Global Stabilization Round*) be the smallest round, such that for all rounds $r \geq \text{GSR}$, the message sent in round r by a correct process q to a correct process p is received by p in round r . This is formally expressed by the following predicate (where \mathcal{C} denotes the set of correct processes):

$$\forall r \geq \text{GSR} : \mathcal{P}_{\text{good}}(r), \text{ where } \mathcal{P}_{\text{good}}(r) \equiv \forall p, q \in \mathcal{C} : \mu_p^r[q] = \sigma_q^r.$$

An algorithm that ensures — in a partially synchronous system — the existence of *GSR* such that $\forall r \geq \text{GSR} : \mathcal{P}_{\text{good}}(r)$, is given in [6]. Note that “ $\forall r \geq \text{GSR} : \mathcal{P}_{\text{good}}(r)$ ” is sufficient for the termination of our algorithms, but not necessary. If the system is synchronous, the following stronger property can be ensured: $\forall r : \mathcal{P}_{\text{good}}(r)$.

3 From Synchrony to Partial Synchrony

In this section we explain our methodology to design a leader-free consensus algorithm that tolerates Byzantine faults without signatures. We start with a leader-free consensus algorithm for Byzantine faults in a synchronous system model, and then extend it to a leader-free consensus algorithm in a partially synchronous system.

3.1 Leader-Free Consensus Algorithm for a Synchronous System

One of the first consensus algorithms that tolerates Byzantine faults in synchronous systems was proposed by Pease, Shostak and Lamport [1]. It is based on an algorithm that solves the *interactive consistency* problem, which consists for each correct process p to compute a vector of values, with an element for each of the n processes, such that

- The correct processes compute exactly the same vector;
- The element of the vector corresponding to a given correct process is the initial value of that process.

The algorithm presented in [1] is not leader-based, does not require signatures, tolerates $t < n/3$ Byzantine faults, and consists of $t + 1$ rounds of exchange of messages. We briefly recall the principle of this algorithm (see Algorithm 1).

The information maintained by each process during the algorithm can be represented as a tree (called *Exponential Information Gathering (EIG)* tree in [13,15]), in which each path from the root to a leaf contains $t + 2$ nodes. Thus the height of the tree is $t + 1$. The nodes are labeled with sequences of processes’ identities in the following manner. The root is labeled with the empty sequence λ ($|\lambda| = 0$). Let i be an internal node in the tree with label $\alpha = p_1 p_2 \dots p_r$; for every $q \in \Pi$ such that $q \notin \alpha$, node i has one child labeled αq . Node i with label α will be simply called “node α ”.

Algorithm 1. *EIGByz* with $n > 3t$ (code of process p)

```

1: Initialization:
2:    $W_p := \{\langle \lambda, v_p \rangle\}$  /*  $v_p$  is the initial value of process  $p$ ;  $val_p(\lambda) = v_p$  */
3: Round  $r$ : /*  $1 \leq r \leq t + 1$  */
4:    $S_p^r$ :
5:   send  $\{\langle \alpha, v \rangle \in W_p : |\alpha| = r - 1 \wedge p \notin \alpha \wedge v \neq \perp\}$  to all processes
6:    $T_p^r$ :
7:   for all  $\{q \mid \langle \alpha, v \rangle \in W_p \wedge |\alpha| = r - 1 \wedge q \in \Pi \wedge q \notin \alpha\}$  do
8:     if  $\langle \beta, v \rangle$  is received from process  $q$  then
9:        $W_p := W_p \cup \{\langle \beta q, v \rangle\}$  /*  $val_p(\beta q) = v$  */
10:    else
11:       $W_p := W_p \cup \{\langle \beta q, \perp \rangle\}$  /*  $val_p(\beta q) = \perp$  */
12:    if  $r = t + 1$  then
13:      for all  $\langle \alpha, v \rangle \in W_p$  from  $|\alpha| = t$  to  $|\alpha| = 1$  do
14:         $W_p := W_p \setminus \langle \alpha, v \rangle$  /* replace  $val_p(\alpha)$  ... */
15:        if  $\exists v'$  s.t.  $|\langle \alpha q, v' \rangle \in W_p| \geq n - |\alpha| - t$  then
16:           $W_p := W_p \cup \langle \alpha, v' \rangle$  /* ... with  $newval_p(\alpha)$  */
17:        else
18:           $W_p := W_p \cup \langle \alpha, \perp \rangle$  /* ... with  $newval_p(\alpha)$  */
19:        for all  $q \in \Pi$  do /* level 1 of the tree */
20:           $M_p[q] := v$  s.t.  $\langle q, v \rangle \in W_p$ 

```

Intuitively, $val_p(p_1 p_2 \dots p_r)$ (which denotes the value of node $p_1 p_2 \dots p_r$ in p 's tree) represents the value v that p_r told p at round r that p_{r-1} told p_r at round $r - 1$ that \dots that p_1 told p_2 at round 1 that p_1 's initial value is v . Each correct process p maintains the tree using a set W_p of pairs $\langle \text{node label}, \text{node value} \rangle$. At the beginning of round r , each process p sends the $(r - 1)$ th level of its tree to all processes (line 5). When p receives a message from q in format $\langle p_1 p_2 \dots p_r, v \rangle$, it adds $\langle p_1 p_2 \dots p_r q, v \rangle$ to its set W_p (line 9). If p fails to receive a message it expects from process q , p simply adds $\langle p_1 p_2 \dots p_r q, \perp \rangle$ to its set W_p (line 11).

Information gathering as described above continues for $t + 1$ rounds, until the entire tree has been filled in. At this point the second stage of local computation starts. Every process p applies to each subtree a recursive data reduction function to compute a new value (lines 13 to 18). The value of the reduction function on p 's subtree rooted at a node labeled α is denoted $newval_p(\alpha)$. The reduction function is defined for a node α as follows.

- If α is a leaf, its value does not change ($newval(\alpha) = val(\alpha)$);
- Otherwise, if there exists v such that $n - |\alpha| - t$ children have value v , then $newval(\alpha) = v$, else $newval(\alpha) = \perp$ (lines 16 and 18).

The reason for a quorum of size $n - |\alpha| - t$ can be explained as follows.² Each correct process, at the end of round $t + 1$, has constructed a tree with $t + 2$ levels. Any node in level $0 < k < t + 1$ has $n - k$ children and a label α such that $|\alpha| = k$. If α is a label with only correct processes, then all its children except t (i.e., $n - k - t$ children) have the same value.

At the end of round $t + 1$, every correct process p constructs a vector \mathbf{M}_p of size n (corresponding to the level 1 of its tree), where $\mathbf{M}_p[q]$ is the new value of process q (line 20). *EIGByz* ensures that:

² Since $n > 3t$, this quorum can be replaced by $\frac{n+t}{2} - |\alpha|$ (see [1]).

- The correct processes compute exactly the same vector, i.e., $\forall p, q \in \mathcal{C} : \mathbf{M}_p = \mathbf{M}_q$, and
- The element of the vector corresponding to a given correct process q is the initial value of that process, i.e., $\forall p, q \in \mathcal{C} : \mathbf{M}_p[q] = v_q$.

Therefore, a correct process can decide by applying a deterministic function on its vector \mathbf{M}_p . The *EIGByz* algorithm ensures the following property:

$$(\forall r, 1 \leq r \leq t+1 : \mathcal{P}_{good}(r)) \Rightarrow \forall p, q \in \mathcal{C} : (\mathbf{M}_p = \mathbf{M}_q) \wedge (|\mathbf{M}_p| \geq |\mathcal{C}|) \quad (1)$$

where $|\mathbf{M}_p|$ denotes the number of non- \perp elements in vector \mathbf{M}_p , and $|\mathcal{C}|$ denotes number of correct processes. The premise holds if the system is synchronous.

3.2 Extending *EIGByz* for a Partially Synchronous Model

If Algorithm 1 is executed in a partially synchronous system, it does not ensure $\forall p, q \in \mathcal{C} : (\mathbf{M}_p = \mathbf{M}_q) \wedge (|\mathbf{M}_p| \geq |\mathcal{C}|)$. Therefore, it cannot ensure the agreement property of Byzantine consensus. However, following two properties hold for Algorithm 1 in synchronous as well as in asynchronous periods:

$$\forall p, q \in \mathcal{C} : \mathbf{M}_p[q] \in \{v_q, \perp\} \quad (2)$$

$$\forall q \in \Pi \setminus \mathcal{C}, \exists v \text{ s.t. } \forall p \in \mathcal{C} : \mathbf{M}_p[q] \in \{v, \perp\} \quad (3)$$

where v_q is the initial value of process q . The proofs are in [14].

To ensure agreement in a partially synchronous system, we need to combine Algorithm 1 with another algorithm. We show below two such algorithms: (i) a simple algorithm (Algorithm 2), which requires $n > 5t$, and (ii) a more complex algorithm with optimal resilience $n > 3t$ (Algorithm 3). In both cases, Algorithm 2 and Algorithm 3 ensure agreement, while Algorithm 1 ensures termination.

Consensus Algorithm with $n > 5t$. We start with a simple parameterized consensus algorithm (see Algorithm 2). Parametrization allows us to easily adjust the algorithm to ensure agreement for different fault models. The algorithm was first presented in [16] as *One Third Rule* algorithm ($T = E = 2n/3$) to tolerate $t < n/3$ benign faults. The parameterized version was given in [17] to tolerate “corrupted communication”. Here, since we consider “Byzantine process faults” we need different values for the parameters. Note that in the context of Byzantine faults, Algorithm 2 alone does not ensure termination.

The algorithm consists of a sequence of phases ϕ , where each phase has two rounds $2\phi - 1$ and 2ϕ . Round 2ϕ is a normal round; to ensure termination, round $2\phi - 1$ will have to be simulated by Algorithm 1. Each process p has a single variable x_p , and in every round p sends x_p to all processes. Parameter T (line 7) refers to a “threshold” for updating x_p , and parameter E (line 13) refers to “enough” same values to decide.³

With Byzantine faults, Algorithm 2 ensures agreement with $E \geq (n + t)/2$ and $T \geq 2n - 2E + 2t$. Strong validity requires $T \geq 2t$ and $E \geq t$. Termination,

³ The notation μ_p^r is introduced in Section 2.3.

Algorithm 2. Byzantine algorithm with $n > 5t$ (code of process p)

```

1: Initialization:
2:    $x_p := v_p \in V$  /*  $v_p$  is the initial value of  $p$  */
3: Round  $r = 2\phi - 1$  : /* round simulated by  $t + 1$  micro-rounds of Algorithm 1 */
4:    $S_p^r$ :
5:   send  $\langle x_p \rangle$  to all processes
6:    $T_p^r$ :
7:   if number of non- $\perp$  elements in  $\mu_p^r > T$  then
8:      $x_p :=$  smallest most frequent non- $\perp$  element in  $\mu_p^r$ 
9: Round  $r = 2\phi$  :
10:   $S_p^r$ :
11:  send  $\langle x_p \rangle$  to all processes
12:   $T_p^r$ :
13:  if more than  $E$  elements in  $\mu_p^r$  are equal to  $v \neq \perp$  then
14:    DECIDE( $v$ )

```

together with Algorithm 1, requires $n - t > T$ and $n - t > E$. Putting all together, for the case $E = T$, we get $T = E = 2(n + t)/3$ and $n > 5t$. The proofs of agreement and strong validity are in [14]. We discuss now termination. For termination, it is sufficient for Algorithm 2 to have one round $r = 2\phi - 1$ in which the following holds (where $|\mu_p^r|$ denotes the number of non- \perp elements in vector μ_p^r):

$$\forall p, q \in \mathcal{C} : (\mu_p^r = \mu_q^r) \wedge (|\mu_p^r| > T) \quad (4)$$

and one round $r + 1 = 2\phi$ in which we have:

$$\forall p \in \mathcal{C} : |\mu_p^{r+1}| > E. \quad (5)$$

If (4) holds, all correct processes set x_p to the some common value v_0 in round r (line 8), and if (5) holds all correct processes decide v_0 in round $r + 1$ (line 14).

By comparing (1) with (4) and (5), it is easy to see that Algorithm 1 ensures (4) and (5) if it is executed after *GSR*, and we have $|\mathcal{C}| > T$ and $|\mathcal{C}| > E$ (where $|\mathcal{C}| = n - t$). Therefore, the idea is to replace the send/receive of round $2\phi - 1$ of Algorithm 2 by the $t + 1$ micro-rounds of Algorithm 1. In other words, we simulate round $r = 2\phi - 1$ of Algorithm 2 using the $t + 1$ micro-rounds of Algorithm 1:

- Each instance of Algorithm 1 is started with $W_p = \{\langle p, x_p \rangle\}$, where x_p is defined in Algorithm 2;
- At the end of these $t + 1$ micro-rounds, the vector M_p computed by Algorithm 1 is the vector μ_p of messages received by p in round r ($M_p[q] = \mu_p[q] = \perp$ means that p did not receive any message from q in round r).

Note that, the *One Third Rule* algorithm (Algorithm 2 with $T = E = 2n/3$) cannot be used with Byzantine faults because of the agreement problem. Using *EIGByz*, a Byzantine process cannot send different values to different processes in a single round, however, it can send different values to different processes in different rounds which violates agreement.

Consensus Algorithm with $n > 3t$. As Algorithm 2 requires $n > 5t$, its resilience is not optimal. Here we show a new algorithm, which uses mechanisms

from several consensus algorithms, e.g., Ben-Or [3], and PBFT [7] with strong validity, and requires only $n > 3t$ (see Algorithm 3). Note that, as for Algorithm 2, Algorithm 3 ensures strong validity and agreement, but not termination. As for Algorithm 2, termination is ensured by simulating the first round of each phase ϕ of Algorithm 3 by $t + 1$ micro-round of Algorithm 1.

Algorithm 3 consists of a sequence of phases ϕ , where each phase has three rounds ($3\phi - 2, 3\phi - 1, 3\phi$). Each process p has an estimate x_p , a vote value $vote_p$ (initially $?$), a timestamp ts_p attached to $vote_p$ (initially 0), and a set $pre-vote_p$ of valid pairs $\langle vote, ts \rangle$ (initially \emptyset). The structure of the algorithm is as follows:

- If a correct process p receives the same estimate v in round $3\phi - 2$ from $n - t$ processes, then it accepts v as a valid vote and puts $\langle v, \phi \rangle$ in $pre-vote_p$ set. The pre-vote set is used later to detect an invalid vote.
- If a correct process p receives the same pre-vote $\langle v, \phi \rangle$ in round $3\phi - 1$ from $n - t$ processes, then it votes v (i.e., $vote_p = v$) and updates its timestamp to ϕ (i.e., $ts_p = \phi$).
- If a correct process p receives the same vote v with the same timestamp ϕ in round 3ϕ from $2t + 1$ processes, it decides v .

The algorithm guarantees that (i) two correct processes do not vote for different values in the same phase ϕ ; and (ii) once $t + 1$ correct processes have the same vote v and the same timestamp ϕ , no other value can be voted in the following phases. We discuss now agreement and termination. The full proofs are in [14].

Agreement. A configuration is v -valent if (i) $\exists \phi$ such that at least $t + 1$ correct processes p have $(vote_p, ts_p) = (v, ts)$ with $ts \geq \phi$, and (ii) the other correct processes q have $(vote_q, ts_q) = (v' \neq v, ts')$ with $ts' < \phi$.

Let ϕ_0 be the smallest round in which some correct process decides v (line 26). By line 25 at least $t + 1$ correct processes p have $vote_p = v$, $ts_p = \phi_0$, and $x_p = v$ from line 20; the other correct processes q with $vote_q \neq v$ have $ts_q < \phi_0$ from line 19. Therefore the v -valent definition holds. We denote the former set by $\Pi_{=\phi_0}$, and the latter by $\Pi_{<\phi_0}$. Processes in $\Pi_{=\phi_0}$ keep $x_p = vote_p = v$ from phase ϕ_0 onward, and processes in $\Pi_{<\phi_0}$ can only update $vote_p$ to $?$ or v , as we explain now. This ensures agreement.

First, by lines 10 and 13, it is impossible for a correct process to have two different values with the same timestamp in its pre-vote set. By lines 27-30, in phase ϕ_0 , processes in $\Pi_{<\phi_0}$ can only update $vote_p$ to $?$; processes in $\Pi_{=\phi_0}$ do not update neither $vote_p$, nor x_p to some value $\neq v$. By lines 10-14, in phase $\phi_0 + 1$, correct processes can only update x_p to v and can only add $(v, \phi_0 + 1)$ to $pre-vote_p$. Therefore in round $3(\phi_0 + 1) - 1$, correct processes can only update $vote_p$ to v , i.e., only v can be decided in phase $\phi_0 + 1$. The same reasoning can be repeated for all phases after phase $\phi_0 + 1$.

Termination. We explain intuitively termination by considering the smallest phase ϕ such that $3\phi - 2 \geq GSR$. We distinguish two cases: (i) at the beginning of round $3\phi - 2$, all correct processes have $vote_p = ?$, and (ii) at the beginning of round $3\phi - 2$ at least one correct process has $vote_p \neq ?$.

Algorithm 3. Byzantine algorithm with $n > 3t$ (code of process p)

```

1: Initialization:
2:    $x_p := v_p \in V$  /*  $v_p$  is the initial value of  $p$  */
3:    $pre\text{-}vote_p := \emptyset$ 
4:    $vote_p \in V \cup \{?\}$ , initially ?
5:    $ts_p := 0$ 

6: Round  $r = 3\phi - 2$ : /* round simulated by  $t + 1$  micro-rounds of Algorithm 1 */
7:    $S_p^r$ :
8:   send  $\langle x_p, vote_p \rangle$  to all processes
9:    $T_p^r$ :
10:  if at least  $n - t$  elements in  $\mu_p^r$  are equal to  $\langle -, ? \rangle$  then
11:     $x_p :=$  smallest most frequent element  $\langle x, - \rangle$  in  $\mu_p^r$ 
12:     $pre\text{-}vote_p := pre\text{-}vote_p \cup \{\langle x_p, \phi \rangle\}$ 
13:  if at least  $n - t$  elements in  $\mu_p^r$  are equal to  $\langle v, - \rangle$  then
14:     $pre\text{-}vote_p := pre\text{-}vote_p \cup \{\langle v, \phi \rangle\}$ 

15: Round  $r = 3\phi - 1$ :
16:    $S_p^r$ :
17:   send  $\langle v \mid \langle v, \phi \rangle \in pre\text{-}vote_p \rangle$  to all processes
18:    $T_p^r$ :
19:   if at least  $n - t$  elements in  $\mu_p^r$  are equal to  $\langle v \rangle$  then
20:      $vote_p := v$ ;  $ts_p := \phi$ ;  $x_p := v$ 

21: Round  $r = 3\phi$ :
22:    $S_p^r$ :
23:   send  $\langle vote_p, ts_p, pre\text{-}vote_p \rangle$  to all processes
24:    $T_p^r$ :
25:   if at least  $2t + 1$  elements in  $\mu_p^r$  are equal to  $\langle v \neq ?, \phi, - \rangle$  then
26:     DECIDE( $v$ )
27:   if exists  $\langle v \neq ?, ts, - \rangle$  in  $\mu_p^r$  s.t.  $vote_p \neq v$  and  $ts > ts_p$  then
28:     if exists  $t + 1$  elements  $\langle -, -, pre\text{-}vote \rangle$  in  $\mu_p^r$  s.t.  $\langle v, ts' \rangle \in pre\text{-}vote$  and  $ts' \geq ts$  then
29:        $vote_p := ?$ ;  $ts_p := 0$ ;  $x_p := v$ 
30:   if  $vote_p \neq ?$  then  $x_p := vote_p$ 

```

Case (i): Consider round $3\phi - 2$. Since we are after *GSR*, Algorithm 1 ensures that all correct processes p receive the same set $\mu_p^{3\phi-2}$ of messages with $|\mu_p^{3\phi-2}| \geq |C|$ (see formula (1)), i.e., all correct processes p set x_p to the same common value v (line 11), and add the pair $\langle v, \phi \rangle$ to $pre\text{-}vote_p$ (line 12). It follows that, in round $3\phi - 1$, all correct processes p set $vote_p$ to v (line 20), and all correct processes decide v in round 3ϕ (line 26).

Case (ii): This case is more complex to expose. Consider round 3ϕ , and let q be a correct process with the highest timestamp ts_q and $vote_q = v \neq ?$ at the beginning of round 3ϕ . Line 19 ensures that for any other correct process q' with $ts_{q'} = ts_q$, we have $vote_{q'} = vote_q$. Since $3\phi > GSR$, all correct processes p with $vote_p \neq v$ execute lines 27-29. Therefore, at the end of round 3ϕ all correct processes p have $x_p = v$ and $vote_p \in \{v, ?\}$, i.e., all correct processes p start round $3\phi + 1 = 3(\phi + 1) - 2$ with $x_p = v$. If the condition of line 10 holds, then the most frequent pair received is $\langle v, - \rangle$, i.e., $\langle v, \phi + 1 \rangle$ is added to $pre\text{-}vote_p$ (line 12). The condition of line 13 necessary holds at each correct process, i.e., $\langle v, \phi + 1 \rangle$ is added to $pre\text{-}vote_p$ (line 14). Therefore, at the end of round $3\phi + 1$, all correct processes p only have $\langle y, \phi + 1 \rangle$ with $y = v$ in $pre\text{-}vote_p$. It follows that, in round $3\phi + 2$, all correct processes p set $vote_p$ to v (line 20), and all correct processes decide v in round $3\phi + 3$ (line 26).

Note that in Algorithm 3, the set $pre\text{-}vote_p$ can be bounded, based on the following observation. For instance, if $\langle v, \phi \rangle \in pre\text{-}vote_p$ and p wants to add $\langle v, \phi' \rangle$ into its pre-vote with $\phi' > \phi$, then $\langle v, \phi \rangle$ becomes obsolete.

3.3 Summary

The following table summarizes our results. The second column shows the smallest number of processes needed for each algorithm. The third and fourth columns give an upper bound on number of rounds needed for a single consensus in both best and worst cases. The best case is when the system is synchronous from the beginning, i.e., $GSR = 0$. Both algorithms require n^2 messages per round.

	# processes	# rounds (best case)	# rounds (worst case)
Algorithm 2	$5t + 1$	$t + 2$	$GSR + 2(t + 2) - 1$
Algorithm 3	$3t + 1$	$t + 3$	$GSR + 2(t + 3) - 1$

3.4 Optimizations

We describe two possible optimizations that can be applied to our leader-free Byzantine consensus algorithm.

Early termination. The “early termination” optimization can be applied to Algorithm 1 (*EIGByz*). Algorithm 1 always requires $t + 1$ rounds, even in executions in which no process is faulty. With early termination, the number of rounds can be reduced in such cases.

Let f denote the actual number of faulty processes in a given execution. Moses and Waarts in [18] present an early termination version of the exponential information gathering protocol for Byzantine consensus that requires $n > 4t$ and terminates in $\min\{t + 1, f + 2\}$ rounds. The idea is the following. Consider some node α in p 's tree. Process p may know that a quorum (i.e., $n - |\alpha| - t$) of correct children of node α store the same value. When this happens, process p can already determine the value of $newval_p(\alpha)$, and can stop at the end of the next round. The paper presents another early termination protocol with optimal resiliency ($n > 3t$) that terminates in $\min\{t + 1, f + 3\}$ rounds. These two optimizations can be applied to Algorithm 1.

One round decision. The “one round decision” optimization is relevant to Algorithm 2. One round decision means that if all correct processes start with the same initial value, and the system is synchronous from the beginning, then correct processes decide in one single round. Algorithm 2 does not achieve one round decision, because the simulation of Algorithm 1 (*EIGByz*) appears in each phase, including phase 1. To achieve one round decision, we simply skip round 1, and start Algorithm 2 with round 2. If all correct processes start with the same initial value, and $GSR = 0$, then correct processes decide in one round.

The fact that our one round decision algorithm requires “only” $n > 5t$ is not in contradiction with the result in [19], which establishes the lower bound $n = 7t + 1$ for one-step decision. The reason is that we assume for fast decision a partially synchronous system with $GSR = 0$, i.e., the system is initially synchronous, while [19] considers a system that is initially asynchronous.

4 Discussion and Future Work

In a partially synchronous system the predicate $\mathcal{P}_{good}(r)$ can be ensured using the implementations given in [6]. Actually, [6] distinguishes two variants of partial synchrony: (a) one in which the communication bound Δ and the process bound Φ are *unknown*, and (b) one in which these bounds are known but hold only eventually. The implementation of the round model slightly differs depending on the partial synchrony variant that is considered. We consider here model (a), which is also the model considered in the leader-based Castro-Liskov PBFT protocol [7]. In this model a standard technique, used for example in PBFT, is to have exponentially growing timeouts. For example, in PBFT whenever the leader changes (i.e., the recovery protocol has to be executed), the timeout for the next leader is doubled. Taking this leader-based protocol as a case study, Amir et al. [10] pointed its vulnerability to performance degradation under an attack. Indeed in PBFT, f consecutive Byzantine leaders, say l_1, l_2, \dots, l_f could do construct the following attack. The first leader l_1 is mute, the timeout expires, the recovery protocol is activated, and the algorithm switches to the next leader (rotating coordinator) while doubling the timeout. The same happens for leaders l_2 to l_{f-1} until l_f becomes leader. The last leader l_f sends its message as late as possible, but not too late to remain leader. If l_f remains leader forever, then the time required for any request (instance of consensus) is high.

Although PBFT does not assume a round-based model as we do in this paper, the performance failure attack is possible in the case of a leader-based protocol implemented in the round-based model, in the case the round-based model is constructed on top of a partially synchronous model of type (a). However, we believe that this is not the case for leader-free algorithms, i.e., performance failure attacks are not effective in this case. The intuition is that, once the timeout of a correct process becomes large enough to receive all messages from correct processes, Byzantine processes cannot introduce an attack that forces the correct process to double its timeout. Our future work is to validate this intuition analytically and/or experimentally, and to understand under which conditions leader-free algorithms outperform leader-based algorithms.

5 Conclusion

All previously known deterministic consensus algorithms for partially synchronous systems and Byzantine faults are leader-based. However, leader-based algorithms are vulnerable to performance degradation, which occurs when the Byzantine leader sends messages slowly, but without triggering timeouts. In the paper we have proposed a deterministic (no randomization), leader-free Byzantine consensus algorithm in a partially synchronous system. Our algorithm is resilient-optimal (it requires $3t + 1$ processes) and signature-free (it doesn't rely on digital signatures). To the best of our knowledge this is the first Byzantine algorithm that satisfies all these characteristics. We have also presented optimizations for the Byzantine consensus algorithm, including one-round decision. Finally, a simpler leader-free consensus algorithm that uses digital signatures is proposed.

We have designed our algorithms using a new methodology. It consists of extending a synchronous consensus algorithm to a partially synchronous consensus algorithm using an asynchronous algorithm. The asynchronous protocol ensures safety (i.e., agreement and strong validity), while the synchronous algorithm provides liveness (i.e., termination) during periods of synchrony.

Acknowledgments. We would like to thank Martin Hutle, Nuno Santos and Olivier Rütli for their comments on an earlier version of the paper.

References

1. Pease, M., Shostak, R., Lamport, L.: Reaching agreement in the presence of faults. *J. ACM* 27(2), 228–234 (1980)
2. Fischer, M.J., Lynch, N.A., Paterson, M.S.: Impossibility of distributed consensus with one faulty process. *JACM* 32(2), 374–382 (1985)
3. Ben-Or, M.: Another advantage of free choice (extended abstract): Completely asynchronous agreement protocols. In: *PODC 1983*, pp. 27–30. ACM, New York (1983)
4. Rabin, M.: Randomized Byzantine generals. In: *Proc. Symposium on Foundations of Computer Science*, pp. 403–409 (1983)
5. Bracha, G.: An asynchronous $[(n - 1)/3]$ -resilient consensus protocol. In: *PODC 1984*, pp. 154–162. ACM, New York (1984)
6. Dwork, C., Lynch, N., Stockmeyer, L.: Consensus in the presence of partial synchrony. *JACM* 35(2), 288–323 (1988)
7. Castro, M., Liskov, B.: Practical Byzantine fault tolerance and proactive recovery. *ACM Transactions on Computer Systems (TOCS)* 20(4), 398–461 (2002)
8. Martin, J.-P., Alvisi, L.: Fast Byzantine consensus. *IEEE Transactions on Dependable and Secure Computing* 3(3), 202–215 (2006)
9. Kotla, R., Alvisi, L., Dahlin, M., Clement, A., Wong, E.: Zyzzyva: speculative byzantine fault tolerance. *SIGOPS Oper. Syst. Rev.* 41(6), 45–58 (2007)
10. Amir, Y., Coan, B., Kirsch, J., Lane, J.: Byzantine replication under attack. In: *DSN 2008*, pp. 197–206 (2008)
11. Clement, A., Wong, E., Alvisi, L., Dahlin, M., Marchetti, M.: Making Byzantine fault tolerant systems tolerate Byzantine faults. In: *NSDI 2009*, Berkeley, CA, USA, pp. 153–168. USENIX Association (2009)
12. Lamport, L.: State-Machine Reconfiguration: Past, Present, and the Cloudy Future. In: *DISC Workshop on Theoretical Aspects of Dynamic Distributed Systems* (September 2009)
13. Lynch, N.A.: *Distributed Algorithms*. Morgan Kaufmann, San Francisco (1996)
14. Borran, F., Schiper, A.: A Leader-free Byzantine Consensus Algorithm. Technical report, EPFL (2009)
15. Attiya, H., Welch, J.: *Distributed Computing: fundamentals, simulations, and advanced topics*. John Wiley & Sons, Chichester (2004)
16. Charron-Bost, B., Schiper, A.: The Heard-Of model: computing in distributed systems with benign faults. *Distributed Computing* 22(1), 49–71 (2009)
17. Biely, M., Widder, J., Charron-Bost, B., Gaillard, A., Hutle, M., Schiper, A.: Tolerating corrupted communication. In: *PODC 2007*, pp. 244–253. ACM, New York (2007)
18. Moses, Y., Waarts, O.: Coordinated traversal: $(t + 1)$ -round byzantine agreement in polynomial time. In: *FOCS*, pp. 246–255 (1988)
19. Song, Y.J., van Renesse, R.: Bosco: One-step byzantine asynchronous consensus. In: Taubenfeld, G. (ed.) *DISC 2008*. LNCS, vol. 5218, pp. 438–450. Springer, Heidelberg (2008)