

Stellar Consensus by Instantiation

Giuliano Losa* Eli Gafni†
Galois, Inc. UCLA
giuliano@galois.com eli@ucla.edu

David Mazières‡
Stanford
<http://www.scs.stanford.edu/~dm/addr/>

August 12, 2019

Abstract

Stellar introduced a new type of quorum system called a Federated Byzantine Agreement System. A major difference between this novel type of quorum system and a threshold quorum system is that each participant has its own, personal notion of a quorum. Thus, unlike in a traditional BFT system, designed for a uniform notion of quorum, even in a time of synchrony one well-behaved participant may observe a quorum of well-behaved participants, while others may not.

To tackle this new problem in a more general setting, we abstract the Stellar Network as an instance of what we call *Personal Byzantine Quorum Systems*. Using this notion, we streamline the theory behind the Stellar Network, removing the clutter of unnecessary details, and refute the conjecture that Stellar’s notion of intact set is optimally fault-tolerant. Most importantly, we develop a new consensus algorithm for the new setting.

*Funding: BSF Grant 2014226, NSF Grant 1655166, a gift from the Stellar Development Foundation, and Galois, Inc.

†Funding: BSF Grant 2014226, NSF Grant 1655166, and a gift from the Stellar Development Foundation

‡Funding: Stanford Center for Blockchain Research

1 Introduction

We study the consensus problem in a new type of quorum system that we call a Personal Byzantine Quorum System (abbreviated PBQS). In a PBQS, each node has its own, personal notion of what a quorum is, subject to the requirement that if Q_p is a quorum of p and $p' \in Q_p$ then there is a quorum $Q_{p'}$ of p' included in Q_p . Justifying this rather strong requirement on the intuitive level, p' being in a personal quorum of p has the connotation that p ‘trusts’ p' . Hence, its quorum should contain at least one quorum of p' . The reverse is not required since trust might not be reflexive, i.e. p' might not trust p even though p trusts p' .

In contrast to PBQSSs, traditional Byzantine quorum systems are uniform, in the sense that every participant has the same notion of a quorum. Under the assumptions of quorum intersection (i.e., that every two quorums intersect at a well-behaved participant) and quorum availability (i.e., that at least one quorum is exclusively well-behaved), one can implement consensus under eventual synchrony [6]. However, traditionally, the ability to implement consensus using quorums is all or nothing; as soon as two quorums fail to intersect at a well-behaved participant, or if no quorum is available, no subset of the participants can solve consensus.

In a PBQS, it is possible that a subset S of the participants has intersecting quorums, in which case we say S is intertwined, while the system as a whole does not. Relying on quorum intersection to ensure safety to S is straightforward. However, suppose S_1 and S_2 are each intertwined but $S_1 \cup S_2$ is not. In this case there is no way to keep $S_1 \cup S_2$ in agreement, but we can still keep each set internally in agreement. It is also possible that $S_1 \cup S_2$ is not intertwined even though $S_1 \cap S_2 \neq \emptyset$. In this case, can a consensus algorithm also ensure liveness to S_1 and S_2 ? This seems impossible since, if S_1 and S_2 diverge, a participant belonging to both S_1 and S_2 has to pick a side and violate safety on the other side in order to make progress. Those observations raise the problem of determining, given an instance of PBQS and a set B of malicious participants, for which family of sets both safety and liveness are achievable, and whether there is an optimal such family. In Section 2, we give necessary conditions for a family of sets to enjoy consensus and we define the notion of a *consensus cluster*, for which we show how to solve consensus in Section 3.

Another crucial technical difference between PBQSSs and traditional Byzantine quorum systems is that since members do not know what constitutes a quorum for another participant, even in a synchronous period, we face the asynchronous phenomenon that one well-behaved member observes a quorum of well-behaved participants, while others do not. BFT algorithms for eventually synchronous systems can experience this phenomenon only during the period of asynchrony.

Why is it important to study PBQSSs? Beyond theoretical curiosity, PBQSSs successfully abstract a deployed, real-world system: the Stellar Network. We found designing a BFT consensus algorithm which is both safe and live under these conditions to be challenging. Indeed, the Stellar Consensus Protocol [14]

(SCP) has only been proved non-blocking when there are Byzantine failures. Here, we propose an algorithm which is safe and live, albeit impractical. Nevertheless, it serves our purpose of showing that while the Stellar network is optimally fault-tolerant for safety, the family of sets that enjoy both safety and liveness is not optimal as previously conjectured. Furthermore, our algorithm guarantees termination in the eventually synchronous model. Whether a practical protocol can achieve these properties is still an open question.

In addition to introducing the PBQS model, we make the following contributions:

- We design an unauthenticated BFT consensus algorithm using idea from Dwork et al.[6] to solve consensus for the Stellar Network’s consensus clusters.
- We refute the conjecture made in the Stellar Whitepaper [14] that intact sets are optimal for consensus. Indeed we suspect that our generalization of intact sets called consensus clusters are optimal.
- We show that the Stellar Network may harbor several disjoint consensus clusters which can nevertheless remain internally in agreement and live. Past work on federated Byzantine agreement systems [14, 7] (FBAS) assumes global quorum intersection and leaves the reader pondering whether all guarantees collapse should this assumption be violated.

Finally, we formalize the static properties of PBQSs and Stellar’s federated Byzantine agreement systems in Isabelle/HOL; the formal theory is available in the Archive of Formal Proofs [12].

2 Personal Byzantine Quorum Systems

In this section we formalize the Personal Byzantine Quorum System Model (the PBQS Model), we define what it means to solve consensus in this model, we observe that global consensus is impossible even without faults, we give lower bounds on what subsets of participants can possibly enjoy consensus, and we define the notion of a consensus cluster. In a consensus algorithm, different consensus clusters may diverge, but, as we show in the next section, consensus is solvable under eventual synchrony within a consensus cluster. The main technical result of this section is that maximal consensus clusters are disjoint, as it is an obvious requirement for consensus.

Definition 1. A PBQS consists of a set of participants P , a set $B \subseteq P$ of Byzantine participants, a set $W = P \setminus B$ of well-behaved participants, and a function mapping a participant p to its non-empty set of quorums, which are subsets of P . The participants’ quorums must be such that:

Property 1 (Quorum sharing). *If Q_p is a quorum of p and $p' \in Q_p$ then there exists a quorum $Q_{p'}$ of p' such that $Q_{p'} \subseteq Q_p$.*

In other words, property 1 states that a quorum Q of some participant p must contain a quorum of every one of its members. As we show in Lemma 4, this remarkably simple property is sufficient to give a mathematically pleasing structure, obviously required if each consensus cluster is to be internally consistent, to PBQs: Maximal consensus clusters are disjoint.

2.1 Consensus Algorithms in PBQs

We assume that the participants communicate via a fully-connected point-to-point message-passing network. (In the Stellar network this is accomplished using an overlay network and signatures.) This means that a participant always knows the identity of the well-behaved sender of a message that it receives. However, message content is not authenticated (in keeping with the current Stellar Modus Operandi of not forwarding signatures) and therefore a participant cannot trust what a sender p says it heard from sender q . Well-behaved participants take steps according to the algorithm they are given, while Byzantine participants may take arbitrary steps. Each well-behaved participant is scheduled infinitely often and a message sent from a well-behaved participant to a well-behaved participant is eventually delivered.

A consensus algorithm consists of a non-terminating sequential program run by each participant in the system. The program can send and receive messages as well as take local computation steps. Initially, a participant starts with a unique identifier, a set of quorums, the set of all participants (used for round-robin leader election, which is replaced by a probabilistic election algorithm in the Stellar Network), and an input value, all of which are accessible to its program. Crucially, a participant does not know a priori the quorums of other participants (it only knows its own set of quorums). In the Stellar Network, a participant learns one of its own quorums only when it receives messages from all members of that quorum, but this difference is not of consequence. A participant also does not know which participants are Byzantine. At any point, a participant's program may produce a unique, irrevocable decision value.

Definition 2 (Intertwined). We say that a set S of well-behaved participants is intertwined when for every two sets Q and Q' which are both quorums of some (possibly different) members of S , we have $Q \cap Q' \cap W \neq \emptyset$.

Note that, by definition, two intertwined participants cannot have empty quorums.

Definition 3 (Quorum-based algorithm). We say that a consensus algorithm is quorum-based when:

1. If a well-behaved participant p decides, then there must be a quorum Q of p such that p received at least one message from each member of Q .
2. If Q is a quorum of a participant p , $p \in W$, and v is a possible input value, then there exists an execution in which only p and members of Q take steps, and p eventually outputs v .

As we have already noted, a PBQS may, for example, harbor two intertwined sets S_1 and S_2 such that $S_1 \cup S_2$ is not intertwined. As implied by the following lemma, in this case no quorum-based algorithm can solve consensus for $S_1 \cup S_2$.

Lemma 1. *Consider two participants p and p' , $p \neq p'$, and two quorums Q, Q' such that Q is a quorum of p and Q' is a quorum of p' and $(Q \cap Q') \setminus B = \emptyset$. Then no quorum-based algorithm can guarantee agreement between p and p' .*

Proof. By definition of quorum-based algorithm, there are two executions e and e' such that (a) only p and members of Q take steps in e and p decides value v in e , and (b) only p' and members of Q' take steps in e' and p' decides value $v' \neq v$ in e' . Because Q and Q' are disjoint, the execution $e \cdot e'$ consisting of the concatenation of e and e' is also an execution. Moreover, agreement is violated in $e \cdot e'$. \square

Lemma 1 shows that, in general, consensus in a PBQS is not solvable globally. Instead, we reformulate the consensus problem such that, given a PBQS \mathcal{U} and a family of sets of participants depending on \mathcal{U} (and thus on the quorum slices and on W), the traditional properties of consensus have to be guaranteed only to each set in the family.

Definition 4 (The PBQS Consensus Problem). In the PBQS consensus problem for a PBQS \mathcal{U} and a family of sets of participants $\{S_i\}$ (depending on \mathcal{U}), we require that for every set S_i in the family:

- **Agreement:** no two members of S_i decide different values.
- **Liveness:** every member of S_i eventually decides some value.
- **Non-triviality:** if only well-behaved participants take steps and a member of S_i decides, then it decides the input value of some well-behaved participant.

Note that the definition above does not preclude any participant from taking steps in the algorithm; instead, the definition gives guarantees only to sets in the family.

In Section 2.3, we define the family of consensus clusters, and we show in Section 3 that PBQS consensus is solvable for consensus clusters. Another, more restrictive, family for which PBQS consensus is solvable is the family of intact sets, as defined in the Stellar Whitepaper. In Section 5, we show that every intact set is a consensus cluster but that the reverse is not true. In this sense, it shows that intact sets cannot be optimal for PBQS consensus. Definition 4 also raises the question of whether there exists an optimal family (in the sense of inclusion) for which PBQS consensus is solvable. We leave this question open, although we conjecture that the consensus clusters family is optimal.

2.2 A Necessary Condition for Liveness

Next we observe that if every quorum Q of a participant p contains a Byzantine node, then it is impossible to guarantee liveness for p because malicious participants can always remain silent. This is formalized using the notion of blocking set:

Definition 5 (Blocking). If R is a set of participants, we say that p is blocked by R , or equivalently that R blocks p , when every quorum of p intersects R . We denote the set of participants blocked by R by $BlockedBy(R)$, and the set of sets that each blocks p , called p 's blocking sets, by $Blocking(p)$.

Lemma 2. *If p is blocked by B then no quorum-based algorithm can ensure liveness to p .*

Proof. If all malicious participants remain silent, then there is no quorum Q such that p eventually receives a message from every member of Q . Therefore, by requirement 1, p never decides. \square

An interesting question is whether q who is blocked by $BlockedBy(B)$ shares the same fate as p who is blocked by B . The answer is positive and a consequence of the quorum sharing property, as implied by the following lemma.

Lemma 3. *In a personal quorum system, for every set of participants R , we have*

$$BlockedBy(BlockedBy(R)) = BlockedBy(R).$$

Proof. Suppose that $p \in BlockedBy(BlockedBy(R))$ but $p \notin BlockedBy(R)$. Hence, there is a quorum Q of p that does not intersect R . However, since $p \in BlockedBy(BlockedBy(R))$, Q must contain p' which is $BlockedBy(R)$. By the quorum sharing property, Q contains a quorum Q' of p' , and by the virtue of p' being blocked by R , Q' contains a member of R . Since $Q' \subseteq Q$, we conclude that Q contains a member of R , and this is a contradiction. \square

Corollary 1. *If p is well-behaved and is not blocked by B , then p has a quorum consisting exclusively of well-behaved participants that are not blocked by B .*

2.3 Consensus Clusters

In this section we define consensus clusters and we show that maximal consensus clusters are disjoint. Consensus clusters can be thought of as disjoint islands which can be kept internally consistent and live by a consensus algorithm, but which may diverge from each other.

Definition 6 (Consensus cluster). A subset $S \subseteq W$ of the well-behaved participants is a consensus cluster when:

- Quorum Intersection: S is intertwined.

- **Quorum Availability:** If $p \in S$ then there is a quorum Q_p of p such that $Q_p \subseteq S$.

Note that, by quorum availability, a member of a consensus cluster must have a quorum, and, by quorum intersection, all its quorums must be non-empty.

We now show that maximal consensus clusters are disjoint. This depends crucially on the quorum sharing property (Property 1) of quorums (i.e., if Q is a quorum of p and $p' \in Q$ then Q contains a quorum of p').

Definition 7. A consensus cluster C is maximal when no strict superset of C is a consensus cluster.

Lemma 4. Consider a personal quorum system. If C_1 and C_2 are two consensus clusters and $C_1 \cap C_2 \neq \emptyset$, then $C_1 \cup C_2$ is a consensus cluster.

Proof. Consider $p \in C_1$ and $q \in C_2$. It suffices to show that p and q are intertwined (quorum availability is immediate). Consider two quorums Q_p and Q_q of p and q , and a quorum Q_m of a participant $m \in C_1 \cap C_2$ such that $Q_m \subseteq C_1$. Since m and q are intertwined by virtue of belonging to C_2 , it follows that Q_q and Q_m have non-empty intersection in C_1 . Let $n \in C_1$ be a member of this intersection. By the quorum sharing property, Q_q contains a quorum Q_n of n . Since both n and p belong to C_1 they are intertwined. Consequently Q_p and Q_n intersect at a well-behaved participant. Since $Q_n \subseteq Q_q$, we get that Q_p and Q_q intersect at a well-behaved participant, and we are done. □

Corollary 2. Maximal consensus clusters are disjoint.

Finally, we present the two properties, Properties 2 and 3, that, as shown in the next section, are sufficient to solve PBQS consensus for consensus clusters.

Property 2 (quorum of member of C , blocks all members of C). *If C is a consensus cluster and Q is a quorum of a member of C , then $Q \cap W$ blocks every member of C .*

Proof of Property 2. Consider $p \in C$. By the quorum-intersection property of consensus clusters, all quorums of p intersect Q at a well-behaved participant. Thus $Q \cap W$ intersects all quorums of p , and we conclude that $Q \cap W$ blocks p . □

Property 3 (blocking set of member of C contains a member of C). *If C is a consensus cluster, $p \in C$, and R blocks p , then $R \cap C \neq \emptyset$.*

Proof of Property 3. By definition of blocking sets, R intersects all quorums of p . Moreover, by the quorum-availability property of consensus clusters, p has a quorum $Q_p \subseteq C$. Thus, R intersects C . □

3 Solving Consensus under Eventual Synchrony in a PBQS

3.1 The Key Insight

Most eventually-synchronous BFT consensus algorithms [6, 3, 11, 5, 1, 8], whether they use authenticated messages or not, rely for liveness on the fact that if two participants p, p' receive the same messages then p observes a quorum (or blocking set) if and only if p' does. For example, this is used by PBFT's leader to convince other participants to prepare its value by attaching signed messages that prove that the value cannot contradict a past decision. In the unauthenticated BFT algorithm of Dwork et al. [6] (Algorithm 3), liveness is ensured by the fact that, during synchrony, a participant that locks a value at the highest round causes all other locks to be released because, thanks to reliable broadcast, the corresponding quorum is observed by all in a timely manner.

Unfortunately, those techniques fail in a PBQS because the notion of quorum is not shared by the participants: even if all participants receive the same messages, one may observe a quorum while the other does not.

The key observation that we make to solve this problem is the following. Consider a consensus cluster C . If, instead of just observing a quorum, a member p of C observes a quorum Q that unanimously states having observed a quorum making statement s , then all members of C that receive the same messages as p can derive that there is a unanimous quorum of some member of C making statement s . This is because, by Property 2, $Q \cap W$ blocks all members of C and, by Property 3, a blocking set contains a member of C , which can be trusted when it reports that a quorum of C unanimously makes statement s .

3.2 The Consensus Algorithm

We assume eventual synchrony, i.e., that there is a time GST after which (a) the messages between well-behaved participants are reliably delivered within a time bound Δ and (b) the relative rate of the clocks of any two well-behaved participants is bounded by a constant ρ . GST, Δ , and ρ are fixed but unknown to the participants.

The consensus algorithm is described in pseudocode in Figure 1. It consists of an unbounded sequence of rounds, where each participant progresses from round to round as instructed by a clock-synchronization protocol described in Section 3.3. The clock-synchronization protocol guarantees that there is a round GSR happening after GST such that for the round GSR and every round after GSR, members of a consensus cluster proceed from round to round synchronously, always receiving each other's messages.

Each four consecutive rounds form an *epoch*. Each epoch has a unique leader chosen round-robin. We refer to the individual rounds within an epoch as *phases*. Nodes broadcast their state at each phase. The algorithm uses a few key concepts:

- A participant *locks* a value v with an associated epoch e when it suspects that v might become decided at epoch e ; if it later observes that the value was in fact not decided, then it unlocks it. Locks ensure that, within a consensus cluster, a value locked by a quorum can never be unlocked.
- A participant p considers a value *final* when it observes that no member of its consensus cluster, should p belong to a consensus cluster, can decide something different.
- A participant p *decides* a value when it observes that no participant that is intertwined with p may make a conflicting decision.
- Participants maintain a *candidate value* and keep track of the *progress-round* of their candidate; a participant assigns progress-round r to its candidate when it adopts it from the leader in round r or when it observes a unanimous quorum with the same candidate and progress round $r - 1$ (we sometimes refer to progress phase when the epoch is clear from the context).

With those concepts in mind, the phases proceed as follows:

- In phase 1, a leader proposes a candidate value on which to try to agree. A node adopts the leader's value unless it suspects that a different value was decided. A node that adopts the leader's value updates its progress round to the current round.
- In phases 2 to 4, a participant p sends a candidate value to all and expects a quorum that agrees with that candidate. At each of those phases, if the expected quorum materializes, the participant updates the progress round of the candidate to the current round. The crucial property of the scheme is that, after GST, if the candidate of a member of consensus cluster C successfully progresses to phase $i > 1$, then all members of C will infer that v has progressed to phase $i - 1$; this is because if Q is a quorum of a member of C , then $Q \cap W$ is a blocking set for all members of C , and a blocking sets contains a member of C and thus can be trusted.
- A processor unlocks its candidate if it gets a “proof” that, after its locking epoch, there was a quorum for another candidate. This is accomplished by observing a unanimous blocking set for a value that progressed to at least phase 2 in a higher epoch.
- A participant whose candidate progresses to phase 3 considers its candidate value *locked* (because it suspects that it may become final in phase 4), and a participant whose candidate progresses to phase 4 considers its candidate final. A final value is not revocable; in contrast “locking” is. At any time, if a participant observes a quorum unanimously declaring the same value v as final, then it decides v . While intertwined participants that are not part of a consensus cluster may disagree on final values (because those

participants may be convinced to unlock arbitrarily by Byzantine participants), they cannot disagree on decisions because final is irrevocable; this is the purpose of the concept of final value.

- A participant that unlocks a value keeps a record that this value was previously locked and at what epoch, and it includes all those records in its messages. Other participants can then check that the unlock steps are valid, by making sure they can derive first-hand that a quorum justifies each unlock step; this prevents a well-behaved participants outside a consensus cluster from “contaminating” the consensus cluster because of a bogus unlock step. This is essential because members of a consensus cluster might depend on a well-behaved outsider for quorum intersection.

The fact that final values are irrevocable guarantees that two intertwined participants never disagree. The crux of the algorithm’s liveness is that a quorum with progress phase 2 suffices to unlock a value, while it takes a quorum with progress phase 3 to lock a value; this ensures that, after GST, the highest lock causes all other locks among a consensus cluster to become unlocked and the leader to adopt the corresponding value. A decision is then necessarily reached in the next epoch.

3.3 Clock Synchronization

We now describe a clock-synchronization algorithm adapted from the Stellar Consensus Protocol [14], which is simpler than the algorithm of Dwork et al. A participant p running the clock-synchronization protocol continuously advertises its current round $r[p]$ to all other participants, and it updates its round according to the following rules:

1. If p hears from a quorum whose members all advertise a round greater or equal to $r[p]$, then p arms a timer of duration $r[p] \cdot T_0$, where T_0 is some base timeout (e.g., 1 second).
2. If p ’s timer fires, p increments its current round.
3. If there is a round $r' > r[p]$ such that p hears from a blocking set whose members all advertise a round greater or equal to r' , then p cancels any pending timeout and advances $r[p]$ to r' .

Now consider a consensus cluster C . By Property 2, rule 3 ensures that, after GST, any members of C that straggle in lower rounds catch up in constant time d_1 to the highest round that is advertised unanimously by the well-behaved portion of a quorum Q of C (because $Q \cap W$ is a blocking set for members of C). Since a blocking set must contain a member of C , rule 3 cannot be used by Byzantine participants to bring well-behaved participants to a round that was not already started by a member of C . Finally, rules 1 and 2 ensure that, despite Byzantine behavior, the first member of C to enter round r stays in round r for a duration proportional to r . Thus, round progression slows down

```

struct NODESTATE {
  round round, progress ▷ initially 1, 0
  value val ▷ initial input
  bool locked, final ▷ initially false, false
  epoch lockEpoch
  set⟨messages⟩ received ▷ all received messages with valid unlockHistory
  set⟨pair⟨epoch, value⟩⟩ unlockHistory ▷ all values ever unlocked
}
define epoch( $r$ ) =  $\lceil r/4 \rceil$ 
define leader( $e$ ) = participant ( $e \bmod N$ ) ▷  $N$  is the number of participants
define phase( $r$ ) =  $r - 4 \cdot \text{epoch}(r) + 1$ 
define valid( $h$ ) =  $\forall (e, v) \in h$ , a quorum sent messages w. epoch(progress) >  $e$  and
val  $\neq v$ .

method NODESTATE::BEGINROUND()
  broadcast(this) ▷ send (round, progress, val, final, lockHistory) to all nodes

method NODESTATE::ENDROUND()
  let  $F \leftarrow \{m \mid m \in \text{received} \text{ and } m.\text{val} = \text{val} \text{ and } m.\text{final}\}$ 
  if final and  $\{m.\text{sender} \mid m \in F\}$  is a quorum then
    decide(val)
  end if
  if phase(round) = 1 then
    let leaderState  $\leftarrow m \mid m \in \text{received} \text{ and } m.\text{round} = \text{round}$ 
    and  $m.\text{sender} = \text{leader}(\text{epoch}(\text{round}))$ 
    if !locked or leaderState.val = val then
      val  $\leftarrow$  leaderState.val
      progress  $\leftarrow$  round
    end if
  else ▷ phases 2–4
    let  $R \leftarrow \{m \mid m \in \text{received} \text{ and } m.\text{val} = \text{val} \text{ and } m.\text{progress} = \text{round} - 1\}$ 
    if  $\{m.\text{sender} \mid m \in R\}$  contains a quorum then
      progress  $\leftarrow$  round
      if phase(round) = 3 then
        locked  $\leftarrow$  true
        lockEpoch  $\leftarrow$  epoch(round)
      else if phase(round) = 4 then
        locked  $\leftarrow$  true
        final  $\leftarrow$  true ▷ can no longer unlock
      end if
    end if
  if phase(round)=4 then
    let  $B \leftarrow$  highest-round unanimous-val blocking set w. phase(progress) = 2
     $w \leftarrow$  unanimous value in  $B$ 
     $e \leftarrow$  epoch(round) in  $B$ 
    if !locked or (locked and !final and  $e > \text{lockEpoch}$  and  $w \neq \text{val}$ ) then
      locked  $\leftarrow$  false
      unlockHistory.insert((epoch(round),val))
      val  $\leftarrow$   $w$  ▷ only matters if we are next leader
    end if
    progress  $\leftarrow$  round ▷ sets phase(progress)=4, not checked in phase 1
  end if
end if

```

Figure 1: Algorithm pseudocode

linearly with time, and there eventually comes a round GSR after which rounds are long enough for all members of C to receive each other’s messages. Note the timer duration in Rule 1 can be change, e.g., to obtain an exponential increase in round duration.

4 Consensus in Federated Byzantine Agreement Systems

In this section we show that, despite their seemingly unrealistic features, PBQSs are a useful model of Stellar’s federated Byzantine agreement systems (FBASs). More precisely:

- We instantiate the consensus algorithm of Section 3 to FBASs, providing effective ways to implement its steps.
- Given a FBAS, we define a corresponding PBQS and we show that, under eventual synchrony, the instantiated consensus algorithm behaves similarly to its counterpart in the PBQS model.

The results of this section show that consensus clusters can be kept safe and live in a federated Byzantine agreement system that does not enjoy system-wide quorum intersection, whereas previous work on the subject made the assumption of system-wide quorum intersection. This is important because, in practice, misconfigured participants, rival factions, or compromised participants could, in violating quorum intersection, yield several disjoint consensus clusters.

4.1 Federated Byzantine Agreement Systems

In a FBAS, each participant chooses a set of slices, which are sets of participants. A participant p considers a set Q to be a quorum when (a) p has at least one slice that is a subset of Q and (b) every member of Q has a slice that is a subset of Q . Practical aspects of FBASs are beyond the scope of this paper, and we refer the reader to Mazières [14] for such matters. What we will say is that it is intended that a participant will trust any information unanimously agreed upon by any of its slices, and thus a quorum is, intuitively, a set that trusts itself.

Slice-based quorums have the advantage that any new participant can join or leave the system without coordination (to join, all it needs to do is join the communication substrate; in practice, this is an overlay network emulating a point-to-point network using public-key cryptography). Moreover, any participant can also reconfigure its slices unilaterally, without coordination, e.g., to remove participants it deems unreliable or to add newcomers. On the flip side, without further assumptions, there is no guarantee that quorums will intersect, and the set of participants at a given time is generally unknown. For the analysis that follows, we assume that the set of participants is unknown but fixed and that the participants’ slices do not change throughout an execution.

Three key aspects of federated Byzantine agreement systems prevent a straightforward analogy with PBQs for the purpose of solving consensus:

1. Since each participants self-declares its set of slices (e.g., by broadcasting it), participants discover their quorums as they receive the slices of other participants. Byzantine participants have the opportunity to declare arbitrary slices and shape the quorums of well-behaved participants.
2. The algorithms of Section 3 require checking whether a set of participants is a blocking set. Doing this check by enumerating quorums is not practical even if all slices are known because the number of quorums of a participant may be exponential in the size of the system.
3. The set of participants is unknown, and thus round-robin leader-election is impossible.

4.2 Abstracting Federated Byzantine Agreement Systems

In a FBAS, participants discover quorums as they learn about the slices of other participants. Therefore, for a participant, the notion of quorum is not fixed; instead, it is augmented with new quorums as the participant learns about the slices of other participants. We call the quorums of a participant p at time t the observed quorums of p at time t . We now define a fixed notion of abstract quorums, which form a PBQS, and relate them to observed quorums.

Definition 8 (Abstract Quorums). A set Q is an abstract quorum of participant p when $p \in B$ or p has a slice contained in Q , and every well-behaved member of Q has a slice contained in Q .

Note that the definition of abstract quorum places requirements only on well-behaved nodes. Hence it is not computable by the participants, who do not know which participants are well-behaved. The following three lemmas are direct consequences of the definition of abstract quorum.

Lemma 5. *Abstract quorums form a PBQS.*

Proof. From the definition of abstract quorum we immediately get that if Q is an abstract quorum of p and $p' \in Q$, then Q is an abstract quorum of p' . \square

Lemma 6. *If Q is an observed quorum of a well-behaved participant p at some time t , then Q is an abstract quorum.*

Lemma 7. *Assume that Q is an abstract quorum of $p \in W$ consisting exclusively of well-behaved participants. Then, under eventual synchrony and assuming that participants do advertise their slices: shortly after GST, Q is an observed quorum of p .*

Proof. Since Q is exclusively well-behaved, shortly after GST, all well-behaved participants receive the slices of the members of Q and can check whether Q is a quorum of theirs. \square

Lemma 6 shows that the set of abstract quorums is an over-approximation of the observed quorums. Because all the algorithms presented so far use the notion of quorum only positively (i.e. adding quorums can only enable more behaviors), Lemma 6 implies that abstract quorums are a safe abstraction of the Stellar Network when considering those algorithms, and substituting the notion of observed quorum for quorum in those algorithms does not compromise their safety properties. Lemma 7 shows that, after GST, a well-behaved participant has observed all its abstract quorums. Since the liveness of the consensus algorithm depends only on the behavior of its maximal consensus cluster, we conclude that the instantiation of the algorithm to the FBAS model preserves liveness.

4.3 Checking Whether a Set is Blocking

The algorithms of Section 3 depend on the ability for a participant p to compute whether a given set R is one of its blocking sets. Even if all slices were known, doing so by enumerating p 's quorums is not practical because, by virtue of how quorums are defined in a FBAS, p may have a number of quorums that is exponential in the size of the system. Instead, we now show that there is a recursive algorithm to check whether a set is a blocking set (a) without enumerating quorums and (b) relying only on the knowledge of the slices of well-behaved participants. This algorithm can be run locally or as a distributed algorithm, e.g., as in Stellar's Federated Voting algorithm [14]. It relies on the notion of slice-blocking.

Definition 9 (Slice-Blocking). We say that the set of participants R slice-blocks p when R intersects each slice of p .

Definition 10 (Inductively Blocked). If R is a set of participants, the set of participants inductively blocked by R , denoted R^* , is defined computationally as follows. Start with $R^* = \emptyset$. While a fixpoint is not reached, repeat the following step: add to R^* all the participants that are slice-blocked by $R^* \cup R$.

A participant can compute locally whether some set R is blocking based on its knowledge of other's slices. However, if its knowledge of slices is incomplete, it might wrongly believe that R is not blocking. This can only remove behaviors in the algorithms of Section 3, because blocking set is used only positively, and thus, with Lemma 8, the substitution of inductively blocking for blocking does not impact safety.

Finally, Lemma 9 shows that, after GST, well-behaved blocking sets are reliably identified by well-behaved participants using the notion of inductively blocking. Thus, liveness is also preserved when substituting inductively blocking for blocking.

Lemma 8. *At any time, if R inductively blocks $p \in W$ then R blocks p in the abstract quorum system.*

Proof. Assume by induction that if p' is in a slice of p and p' is inductively blocked by R , then all quorums of p' intersect R .

Now suppose by contradiction that R does not block p in the abstract system, i.e. that Q is an abstract quorum of p and $R \cap Q = \emptyset$. Since Q is an abstract quorum of p , there must be a slice s_p of p such that $s_p \subseteq Q$. Moreover, since R inductively blocks p , then s_p must have a member p' that is inductively blocked by R . By the quorum-sharing property, Q is an abstract quorum of p' . Thus $Q \cap R \neq \emptyset$, which is a contradiction. \square

Lemma 9. *If $p \in W$ and $R \subseteq W$ blocks p in the abstract quorum system, then, shortly after GST, R inductively blocks p .*

Proof. First, observe that, shortly after GST, p knows all the slices of the well-behaved participants. Thus, suppose that p knows all the slices of the well-behaved participants.

Suppose that R does not inductively block p according to p . Then, by definition, there is a slice s_p of p whose members are not inductively blocked by R and such that $s_p \cap R = \emptyset$. Since the members of s_p are not inductively blocked by R , then, for every $p' \in s_p \setminus B$, we also have that there is a slice s'_p of p' whose members are not inductively blocked by R and such that $s'_p \cap R = \emptyset$ (we have to exclude B from s_p since p might not know the slices of Byzantine participants; in the worst case, none of those are observed inductively blocked). Continuing inductively in this fashion, we obtain an abstract quorum Q of p which does not intersect R , and we have only used the slices of well-behaved participants. This contradicts the fact that R blocks p in the abstract quorum system. \square

4.4 Leader Election

As noted before, round-robin leader-election is impossible in a FBAS because the set of participants is in general unknown. In this section we show how to probabilistically elect a leader. However, we give no bound on the probability of success, except that it is non-zero. Devising an efficient leader-election mechanism, or, more generally, a conciliator[2] mechanism, is left open.

To agree on a common leader taken among C with non-zero probability, every participant p selects at random a participant p' belonging to one of its slices or itself. If $p = p'$, then p elects itself as leader and broadcasts (leader, p). Otherwise, it waits to receive a broadcast of the form (leader, p'') from p' , and then elects the participant p'' as leader and broadcasts (leader, p'').

We now show that, through this process, members of C agree on a common leader taken among C with non-zero probability.

Definition 11. Graph $D(S)$ If S is a set of participants, the directed graph $D(S)$ is defined as the graph whose set of vertices is S , and where there is an edge from n_1 to n_2 when $n_2 \neq n_1$ and n_2 is in a slice of n_1 .

Lemma 10. *If C is a consensus cluster, $p \in C$, and Q is a quorum of a member of C , then Q is reachable from p in $D(C)$.*

Proof. Since $p \in C$ and C is a consensus cluster, there is a quorum Q' of p such that $Q' \subseteq C$. Now suppose that Q is not reachable from p in $D(C)$. Then, with $Q' \subseteq C$, we get that $Q' \cap Q = \emptyset$. This contradicts the assumption that C is a consensus cluster. \square

Definition 12. Elementary quorum An elementary quorum is a quorum Q such that no strict subset of Q is a quorum.

Note that, by definition, every quorum contains an elementary quorum.

Lemma 11. *If n_1 and n_2 are members of an elementary quorum q consisting exclusively of well-behaved participants, then there is a path in $D(q)$ from n_1 to n_2 .*

Proof. Suppose q is an elementary quorum and that $n_1, n_2 \in q$ and n_2 is not reachable from n_1 in $D(q)$. Then consider the set S of participants that are reachable from n_1 in $D(q)$. By our assumption above, n_2 does not belong to S . Thus S is a strict subset of q . Moreover, every member n of S has a slice $s_n \subseteq q$. Additionally, consider that we must have that $s_n \subseteq S$, as otherwise a participant outside S would be reachable from n_1 . Thus every member of S has a slice in S , and therefore S is a quorum. Since S is a strict subset of q , this contradicts the fact that q is an elementary quorum. \square

Lemma 12. *If C is a consensus cluster, then there exists a member of C that is reachable in $D(C)$ from every other participant in C .*

Proof. Since C is a quorum, C contains an elementary quorum Q . By Lemma 10, Q is reachable from every member n of C in $D(C)$. Moreover, by Lemma 11, every member of Q is reachable in $D(Q)$ from every other member of Q . Thus, because $D(Q) \subseteq D(C)$, every member of Q is reachable in $D(C)$ from every member of C . \square

Lemma 13. *If C is a consensus cluster, then, with non-zero probability, every member of C elects the same leader $l \in C$.*

Proof. Note that the leader-election algorithm can be seen as randomly selection edges in $D(P)$ (where P is the set of participants). Because there is a member n of C reachable in D from all other members of C in $D(C)$ (and because well-behaved participant have a finite number of outgoing edges), then with non-zero probability the edges selected by the leader-election algorithm will form a sink tree rooted at n , who will be elected unique leader by all members of C . \square

5 Related Work

Federated Byzantine quorum systems were first introduced in the Stellar Whitepaper by Mazières [14], who also proposes the notion of intact set and a consensus algorithm for intact sets, the Stellar Consensus Protocol (SCP). The epidemic propagation mechanism and the clock-synchronization protocol presented in the present paper are taken from the Stellar Whitepaper. Mazières also discusses more practical aspects of the Stellar Network.

One important contribution of the present paper is that Stellar’s intact sets, conjectured in the Stellar Whitepaper to be optimal for consensus, are in fact not the biggest sets for which an algorithm can solve consensus. An intact set is a subset S of W such that every member of S is well-behaved and: (a) if Q and Q' are quorums of S , then $Q \cap Q' \cap S \neq \emptyset$; (b) S is a quorum. Comparing the definitions of consensus cluster and intact set, it is easy to see that any intact set is also a consensus cluster. However, as shown by the following lemma, there are some consensus clusters that are strictly bigger than any intact set.

Lemma 14. *There are some configurations in which a set S is a consensus cluster but S is not intact and S has no intact superset.*

Proof. Consider a system of three well-behaved participants p_1 , p_2 , and p_3 (note that there are no malicious participants) where p_1 has a single slice $\{p_1\}$, p_2 has two slices $\{p_1, p_2\}$ and $\{p_2, p_3\}$, and p_3 has two slices $\{p_1, p_3\}$ and $\{p_2, p_3\}$. According to those slices, the quorums are $\{p_1\}$, $\{p_1, p_2, p_3\}$, $\{p_2, p_3\}$, $\{p_1, p_2\}$, and $\{p_1, p_3\}$. In this system, $C = \{p_2, p_3\}$ is a consensus cluster but is not intact, because $Q_1 = \{p_1, p_2\}$ and $Q_2 = \{p_1, p_3\}$ intersect outside C . Moreover, the only strict superset of C , $\{p_1, p_2, p_3\}$, is not intact because the quorums $\{p_1\}$ and $\{p_2, p_3\}$ do not intersect. \square

Another novel aspect of the present paper compared to the Stellar Whitepaper is that we do not assume global quorum intersection; nevertheless, we show that consensus clusters enjoy safe and live consensus. This is important because it shows that safety and liveness guarantees do not collapse system-wide in the face of misconfigurations or attacks.

We have studied federated quorum system under the assumption that well-behaved participants do not change their slices. However, in practice, well-behaved participants might change their slices to eliminate unreliable participants or add newcomers. The Stellar Whitepaper also analyzes this situation.

García-Pérez and Gotsman [7] conduct a detailed study of Stellar’s federated Byzantine quorum systems and the implementation of broadcast abstractions therein. They also propose the notion of subjective dissemination quorum system (subjective DQS) in which, like in a PBQS, each participant has its own set of quorums. However, subjective DQSs have two crucial differences compared to PBQSs: subjective DQSs have system-wide quorum intersection and they do not have Property 1 (which says that a quorum is a quorum for all its members). In the absence of system-wide quorum intersection, Property 1 of PBQSs ensures that maximal consensus clusters are disjoint (Lemma 4). Without it,

maximal consensus clusters may intersect, which implies that consensus is not solvable even for consensus clusters (a participant in the intersection may have to violate safety on one side in order to make progress).

Ripple [15] introduced the first permissionless quorum-based consensus protocol. In the XRP Ledger Consensus Protocol, each participant p is responsible for configuring its own UNL, which is a list of other participants that p will accept messages from. Moreover, p will accept as a quorum any set of participants consisting of more than a fixed fraction (defined system-wide by the protocol, e.g. 80%) of its UNL. Maintaining agreement in Ripple’s protocol rests on the assumption that participants will provide sufficiently overlapping UNLs (roughly 90% for every pair of participants, in the most adversarial model of Chase and MacBrough [4]).

Traditional Byzantine quorum systems are uniform, in the sense that every participant has the same notion of quorum. Uniform Byzantine quorum systems are studied in details by Malkhi and Reiter [13]. More complex types of uniform quorum systems are studied by Guerraoui and Vukolić [9]. General Byzantine adversaries [10] do not give rise to a PBQS because participants have global knowledge of the adversary in this model.

Acknowledgment: The authors are in debt to an anonymous reviewer who suspected that our algorithm had a flaw. Indeed, that suspicion was correct.

References

- [1] Ittai Abraham, Guy Gueta, Dahlia Malkhi, and Jean-Philippe Martin. Revisiting fast practical byzantine fault tolerance: Thelma, Velma, and Zelma. *arXiv preprint arXiv:1801.10022*, 2018.
- [2] James Aspnes. A modular approach to shared-memory consensus, with applications to the probabilistic-write model. *Distributed Computing*, 25(2):179–188, 2012.
- [3] Miguel Castro and Barbara Liskov. Practical byzantine fault tolerance and proactive recovery. *ACM Transactions on Computer Systems (TOCS)*, 20(4):398–461, 2002.
- [4] Brad Chase and Ethan MacBrough. Analysis of the XRP ledger consensus protocol. *arXiv preprint arXiv:1802.07242*, 2018.
- [5] 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.
- [6] Cynthia Dwork, Nancy Lynch, and Larry Stockmeyer. Consensus in the presence of partial synchrony. *Journal of the ACM (JACM)*, 35(2):288–323, 1988.

- [7] Álvaro García-Pérez and Alexey Gotsman. Federated byzantine quorum systems. In *22nd International Conference on Principles of Distributed Systems (OPODIS 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- [8] 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.
- [9] Rachid Guerraoui and Marko Vukolić. Refined quorum systems. *Distributed Computing*, 23(1):1–42, 2010.
- [10] Martin Hirt and Ueli Maurer. Complete characterization of adversaries tolerable in secure multi-party computation. In *PODC*, volume 97, pages 25–34, 1997.
- [11] Ramakrishna Kotla, Lorenzo Alvisi, Mike Dahlin, Allen Clement, and Edmund Wong. Zyzzyva: Speculative byzantine fault tolerance. In *Proceedings of Twenty-first ACM SIGOPS Symposium on Operating Systems Principles*, SOSP '07, pages 45–58, New York, NY, USA, 2007. ACM.
- [12] Giuliano Losa. Stellar quorum systems. *Archive of Formal Proofs*, August 2019. http://isa-afp.org/entries/Stellar_Quorums.html, Formal proof development.
- [13] Dahlia Malkhi and Michael Reiter. Byzantine quorum systems. *Distributed computing*, 11(4):203–213, 1998.
- [14] David Mazieres. The Stellar Consensus Protocol: A federated model for internet-level consensus. *Stellar Development Foundation*, page 32, 2015.
- [15] David Schwartz, Noah Youngs, and Arthur Britto. The Ripple protocol consensus algorithm, 2014. https://ripple.com/files/ripple_consensus_whitepaper.pdf.