

Byzantine Set-Union Consensus using Efficient Set Reconciliation

Florian Dold

Inria

Email: florian.dold@inria.fr

Christian Grothoff

Inria

Email: christian.grothoff@inria.fr

Abstract—Applications of secure multiparty computation such as certain electronic voting or auction protocols require Byzantine agreement on large sets of elements. Implementations proposed in the literature so far have relied on state machine replication, and reach agreement on each individual set element in sequence.

We introduce *set-union consensus*, a specialization of Byzantine consensus that reaches agreement over whole sets. This primitive admits an efficient and simple implementation by the composition of Eppstein’s set reconciliation protocol with Ben-Or’s ByzConsensus protocol.

A free software implementation of this construction is available in GNUet. Experimental results indicate that our approach results in an efficient protocol for very large sets, especially in the absence of Byzantine faults. We show the versatility of set-union consensus by using it to implement distributed key generation, ballot collection and cooperative decryption for an electronic voting protocol implemented in GNUet.

I. INTRODUCTION

Byzantine consensus is a fundamental building block for fault-tolerant distributed systems. It allows a group of peers to reach agreement on some value, even if a fraction of the peers are controlled by an active adversary. Earlier theory-oriented work on Byzantine consensus has focused on finding a single agreement on a binary flag or bit string. More recent approaches for practical applications are mainly based on state machine replication (SMR), wherein peers agree on a sequence of state machine transitions. State machine replication makes it relatively easy to lift existing, non-fault-tolerant services to a Byzantine fault-tolerant implementation [1]. Each request from a client triggers a state transition in the replicated state machine that provides the service.

A major shortcoming of the SMR is that all requests to the service need to be individually agreed upon in sequence by the replica peers of the state machine. This is undesirable since in unoptimized SMR protocols such as PBFT [1], a single transition requires $O(n^2)$ messages to be exchanged for n replicas. Some implementations try to address this inefficiency by optimistically processing requests and falling back to individual Byzantine agreements only when Byzantine behavior is detected. In practice this leads to very complex implementations whose correctness is hard to verify and that have weak progress guarantees [2].

The canonical example for a service where this inefficiency becomes apparent is the aggregation of values submitted by clients into a set. This scenario is relevant for the implementation of secure multiparty computation protocols such

electronic voting [3], where ballots must be collected, and auctions [4], where bids must be collected. A direct implementation that reaches agreement on a set of m elements with SMR requires m sequential agreements, each consisting of $O(n^2)$ messages.

We introduce Byzantine set-union consensus (BSC) as an alternative communication primitive that allows this aggregation to be implemented more efficiently. In order to implement the set aggregation service described above, the peers first reconcile their sets using an efficient set reconciliation protocol that is not fault-tolerant but where the complexity is bounded even in the case of failures. Then, they use a variant of ByzConsensus [5] to reach Byzantine agreement on the union.

We assume a partially synchronous communication model, where non-faulty peers are guaranteed to successfully receive values transmitted by other non-faulty peers within an existing but unknown finite time bound [6]. Peers communicate over pairwise channels that are authenticated. Message delivery is reliable (i.e. messages arrive uncorrupted and in the right order) but the receipt of messages may be delayed. We make the same assumption as Castro and Liskov [1], [7] about this delay, namely that it does not grow faster than some (usually exponential) function of wall clock time. We assume a computationally unbounded adversary that can corrupt at most $t = \lceil n/3 \rceil - 1$ peers creating Byzantine faults. The adversary is static, that is the set of corrupted peers is fixed before the protocol starts, but this set is not available to the correct peers. The actual number of faulty peers is denoted by f , with $f \leq t$.

The BSC protocol has message complexity $O(mn + n^2)$ when no peers show Byzantine behavior. When f peers show Byzantine behavior, the message complexity is $O(mnf + kfn^2)$, where k is the number of valid set elements exclusively available to the adversary. We will show how k can be bounded for common practical applications, since in the general case k is only bounded by the bandwidth available to the adversary. In practice, we expect kf to be significantly smaller than m . Thus, $O(mnf + kfn^2)$ is an improvement over using SMR-PBFT which would have complexity $O(mn^2)$.

We have created an implementation of the BSC protocol by combining Ben-Or’s protocol for Byzantine consensus [5] with Eppstein’s protocol for efficient set reconciliation [8]. We demonstrate the practical applicability of our resulting abstraction by using BSC to implement distributed key generation,

ballot collection and cooperative decryption from the Cramer-Gennaro-Schoenmakers remote electronic voting scheme [3] in the GNUet framework.

II. BACKGROUND

A. Byzantine consensus

The Byzantine consensus problem [9] is a generalization of the consensus problem where peers might also exhibit Byzantine faults. A fundamental result is that no Byzantine consensus protocol with n peers can support $\lceil n/3 \rceil$ or more Byzantine faults in the partially synchronous model [6].

Many specific variants of the agreement problem (such as interactive consistency [10], k -set consensus [11], or leader election [12] and many others [13]) exist. We will focus on the consensus problem, wherein each peer in a set of peers $\{P_1, \dots, P_n\}$ starts with an initial value $v_i \in M$ for an arbitrary fixed set M . At some point during the execution of the consensus protocol, each peer irrevocably decides on some output value $\hat{v}_i \in M$. Informally, a protocol that solves the consensus problem must fulfill the following properties:¹

- *Agreement*: If peers P_i, P_j are correct, then $v_i = v_j$.
- *Termination*: The protocol terminates in a finite number of steps.
- *Validity*: If all correct peers have the same input value \tilde{v} , then all correct peers decide on \tilde{v} .

Some definitions of the consensus problem also include *strong validity*, which requires the value that is agreed upon to be the initial value of some correct peer [14]. The consensus protocol presented in this paper does not offer strong validity. Early attempts at implementing Byzantine consensus with state machine replication are SecureRing [15] and Rampart [16]; they suffered from sacrificing correctness for progress guarantees in the presence of asynchrony [1].

Castro and Liskov’s Practical Byzantine Fault Tolerance (PBFT) [1], [7] does not suffer from this problem. PBFT guarantees progress as long as the message delay does not grow indefinitely for some fixed growth function². PBFT uses a leader to coordinate peers (called *replicas* in BPFT terminology). When replicas detect that the leader is faulty, they run a leader-election protocol to appoint a new leader.

The approach taken by BPFT (and several derived protocols) has several problems [2]: In practice, malicious participants are able to slow down the system significantly or even reduce the throughput to practically zero. Correctness proofs for the respective protocols and the implementation of state machine replication are notoriously difficult [17].

Some more recent Byzantine state machine replication protocols such as Q/U [18] or Zyzzyva [19] have less overhead per request since they optimize for the non-Byzantine case. This comes, however, often at the expense of robustness in the presence of Byzantine faults [2].

¹Different variations and names can be found in the literature. We have chosen a definition that extends to our generalization to sets later on.

²In practice, exponential back-off is used.

B. Gradecast

A key building block for our protocol is Feldman’s Gradecast protocol [20]. In contrast to an unreliable broadcast, Gradecast provides correctness properties to the receivers, even if the leader is exhibiting Byzantine faults.

In a Gradecast, a leader P_L broadcasts a message m among a fixed set $\mathcal{P} = \{P_1, \dots, P_n\}$ of peers. For notational convenience, we assume that $P_L \in \mathcal{P}$. These are the communication steps for peer P_i :

- 1) LEAD: If $i = L$, send the input value v_L to \mathcal{P}
- 2) ECHO: Send the value received in LEAD to \mathcal{P} .
- 3) CONFIRM: If a common value \bar{v} was received at least $n-t$ times in round ECHO, send \bar{v} to \mathcal{P} . Otherwise, send nothing.

Afterwards, each peer assigns a confidence value $c_i \in \{0, 1, 2\}$ that “grades” the correctness of the broadcast. The result is a graded result tuple $\langle \hat{v}_i, c_i \rangle$ containing the output value \hat{v}_i and the confidence c_i . The grading is done with the following rules:

- If some \hat{v} was received at least $n-t$ times in CONFIRM, output $\langle \hat{v}, 2 \rangle$.
- Otherwise, if some \hat{v} was received at least $t+1$ times in CONFIRM, output $\langle \hat{v}, 1 \rangle$.
- Otherwise, output $\langle \perp, 0 \rangle$. Here, \perp denotes a special value that indicates the absence of a meaningful value.

For the c_i , the following correctness properties must hold:

- 1) If $c_i \geq 1$ then $\hat{v}_i = \hat{v}_j$ for correct P_i and P_j
- 2) If P_L is correct, then $c_i = 2$ and $\hat{v}_i = v_L$ for correct P_i .
- 3) $|c_i - c_j| \leq 1$ for correct P_i and P_j .

When a correct peer P_i receives a Gradecast with confidence 2, it can deduce that all other peers received the same message, but some other peers might have only received it with a confidence of 1. Receiving a Gradecast with confidence 1 also guarantees that all other correct peers received the same message. However, it indicates that the leader behaved incorrectly. No assumption can be made about the confidence of other peers. Receiving a Gradecast with confidence 0 indicates that the leader behaved incorrectly and, crucially, that all other correct peers *know* that the leader behaved incorrectly.

A simple counting argument proves that the above protocol satisfies the three Gradecast properties. [20]

C. ByzConsensus

ByzConsensus [5] uses Gradecast to implement a consensus protocol for simple values. Each peer begins with a starting value $s_i^{(1)}$ and the list of all n participants \mathcal{P} . Each peer also starts with an empty blacklist of corrupted peers. If a peer is ever blacklisted, it is henceforth excluded from the protocol. In ByzConsensus, Gradecast is used to force corrupt peers to either expose themselves as faulty—and consequently be excluded—by gradecasting a value with low confidence, or to follow the protocol and allow all peers to reach agreement.

ByzConsensus consists of at most $f+1$ sequentially executed super-rounds $r \in 1 \dots f+1$ where $f \leq t$. In each super-round, each peer leads a Gradecast using their candidate value

$s_i^{(r)}$; these n Gradecasts can be executed in parallel. Leaders where the Gradecast results in a confidence of less than 2 are put on the blacklist. Recall that different correct peers might receive a Gradecast with different confidence; thus peers do not necessarily agree on the blacklist.

At the end of each super-round, each peer computes a new candidate value $s_i^{(r+1)}$ using the value that was received most often from the Gradecasts with a confidence of at least 1. If $s_i^{(r)}$ was received more than $n - t$ times, then $r = f$ and the next round is the last round.

If the final candidate value does not receive a majority of at least $2t + 1$ among the n Gradecasts, or if the blacklist has more than t entries, then the protocol failed: either more than t faults happened or, in the partially synchronous model, correct peers did not receive a message within the designated round due to the delayed delivery.

ByzConsensus has message complexity $O(fn^3)$. While the asymptotic message complexity is obviously worse than the $O(n^2)$ of PBFT, there is a way to use set reconciliation to benefit from the parallelism of the Gradecast rounds and thereby reduce the complexity to $O(fn^2)$.

D. Set reconciliation

The goal of set reconciliation is to identify the differences between two large sets, say S_a and S_b , that are stored on two different machines in a network. A simple but inefficient solution would be to transmit the smaller of the two sets, and let then receiver compute and announce the difference. Research has thus focused on protocols that are more efficient than this naive approach with respect to the amount of data that needs to be communicated when the sets S_a and S_b are large, but their symmetric difference $S_a \oplus S_b$ is small.

A practical protocol was first proposed by Eppstein et al. in 2011. [8] It is based on invertible Bloom filters (IBFs), a data structure that is related to Bloom filters [21]. An attractive property of this approach is that IBFs are used both to construct an estimator for the size of the symmetric difference between two sets, as well as for the the reconciliation itself, which requires this estimate.

An existing generalization of IBFs to multi-party set reconciliation [22] based on network coding requires trusted intermediaries, and is thus not applicable in the presence of Byzantine faults.

III. OUR APPROACH

We now describe how to combine the previous approaches into a protocol for Byzantine fault-tolerant set consensus. The goal of the adversary is to sabotage timely consensus among correct peers, e.g. by increasing message complexity or forcing timeouts.

A major difficulty with agreeing on a set of elements as a whole is that malicious peers can initially withhold elements from the correct peers and later send them only to a subset of the correct peers. This could possibly happen at a time when it is too late to reconcile the remaining difference caused by distributing these elements. We assume that the number of

these elements that are initially known to the adversary but not to all correct peers is bounded by k , where k exists but is not necessarily known to the correct participants.

A. Definition

We now give a definition of set-union consensus that is motivated by practical applications to secure multiparty computation protocols such as electronic voting, which are discussed in more detail in Section VII.

Consider a set of n peers $\mathcal{P} = \{P_1, \dots, P_n\}$. Fix some (possibly infinite) universe M of elements that can be represented by a bit string. Each peer P_i has an initial set $S_i^{(0)} \subseteq M$.

Let $R : 2^M \rightarrow 2^M$ be an idempotent function that canonicalizes subsets of M by replacing multiple conflicting elements with the lexically smallest element in the conflict set and removes invalid elements. What is considered conflicting or invalid is application-specific. During the execution of the set-union consensus protocol, after finite time each peer P_i irrevocably commits to a set S_i such that:

- 1) For any pair of correct peers P_i, P_j it holds that $S_i = S_j$.
- 2) If P_i is correct and $e \in S_i^0$ then $e \in S_i$.
- 3) The set S_i is canonical, that is $S_i = R(S_i)$.

The canonicalization function allows us to set an upper bound on the number of elements that can simultaneously be in a set. For example in electronic voting, canonicalization would remove malformed ballots and combine multiple different (encrypted) ballots submitted by the same voter into a single “invalid” ballot for that voter.

B. Byzantine set-union consensus (BSC) protocol

Recall that every peer P_i , $0 < i \leq n$ starts with a set $S_i^{(0)}$. The BSC protocol incorporates two subprotocols, bounded Eppstein set reconciliation and lower bound agreement and uses those to realize an efficient variant of ByzConsensus.

The basic problem solved by the two subprotocols is bounding the cost of Eppstein’s set reconciliation. Given a set size difference between two peers of k , the expected cost of Eppstein’s set reconciliation is $O(k)$ if both participants are honest. However, we need to ensure that malicious peers cannot raise the complexity to $O(m)$ where m is the size of the union. There are two attacks, discussed in the next section, that could degrade the performance of the unmodified Eppstein protocol, causing it to either send or receive significantly more than $O(k)$ elements.

1) *Bounded Eppstein set reconciliation:* To send more than $O(k)$ elements, a malicious peer can send IBFs that fail to decode, and cause the reconciliation algorithm to fall back to sending the complete set, creating $O(m)$ traffic instead of $O(k)$. This can be thwarted by forcing senders prove that their sets are large enough to justify the IBF decoding failures.

In our bounded Eppstein set reconciliation protocol, a peer that observes a modest (i.e. logarithmic in m) number of IBF decoding failures requires a probabilistic proof of the size of the sender’s set using sampling. Here, a receiver experiencing IBF decoding failures transmits a *challenge*. The sender must then respond with sample elements close to the

challenge's value(s). The receiver can then estimate the size of the sender's set from the proximity of the results to the challenge using [23]. If the size estimate is too low to justify the IBF decoding failures, the receiver aborts the reconciliation and blacklists the sender.

2) *Lower bound agreement*: To provide a lower bound on the permissible set size for set reconciliation, BSC first executes a protocol for *lower bound agreement* (LBA). In this first step, every correct peer P_i learns a superset $S_i^{(1)}$ of the union of all correct peers' initial sets, as well as a lower bound ℓ_i for the minimum number of elements shared by all correct peers where $n - \ell_i \leq k$. Note that neither $S_i^{(1)} = S_j^{(1)}$ nor $\ell_i = \ell_j$ necessarily hold even for correct peers P_i and P_j . Our LBA protocol proceeds in three steps:

- (i) All peers reconcile their initial set with each other, using pairwise (bounded) Eppstein set reconciliation.
- (ii) All peers send their current set size to each other, and each peer P_i sets ℓ_i to the $(t + 1)$ -smallest set size that P_i received.
- (iii) All peers again reconcile their sets with each other, using pairwise (bounded) Eppstein set reconciliation.

The third step is necessary to ensure that every correct P_i has at least ℓ_i elements, since malicious peers could use the k elements initially withheld to force an honest peer's set size below the $(t + 1)$ -smallest set size. Thanks to the repetition even if ℓ_i is different for each peer, it is guaranteed that P_i has at least ℓ_i elements in common with every other good peer.

In subsequent set reconciliations, ℓ_i can be used to bound the traffic that malicious peers are able to cause by falsely claiming to have a large number of elements missing. LBA itself has complexity $O(nmf)$: initially all malicious peers can *once* claim to have empty sets with all other peers. LBA ensures that for the remainder of the protocol, a correct peer with m_i elements can stop sending elements to malicious peer P_M after P_M requested $m_i - \ell_i \leq k$ elements.

3) *Exact set agreement*: After LBA, an *exact set agreement* is executed, where all peers reach Byzantine agreement for a super-set of the set reached in LBA. The exact set agreement is implemented by executing a variant of ByzConsensus which instead of sending values reconciles sets.

The Gradecast is adapted as follows:

- (i) LEAD: If $i = L$, reconcile the input set V_L with \mathcal{P} .
- (ii) ECHO: Reconcile the set received in LEAD with \mathcal{P} .
- (iii) CONFIRM: Let \mathcal{U}_E be the union of all sets received in the ECHO round, and $N_E(e)$ the number of times a single set element e was received.

If $\bigvee_{e \in \mathcal{U}_E} t < N_E(e) < n - t$, send \perp (where $\perp \neq \emptyset$). Otherwise send $\mathcal{U}_E - \{e \mid N_E(e) \leq t\}$ to \mathcal{P} .

The grading rules are also adapted to sets. Let \mathcal{U}_C be the union of sets received in CONFIRM, $N_C^+(e)$ the number of times a single element $e \in \mathcal{U}_C$ was received, and $N_C^-(e)$ the number of sets (not \perp) received in CONFIRM that excluded e .

- If $\bigwedge_{e \in \mathcal{U}} N_C^+(e) \geq n - t \vee N_C^-(e) \geq n - t$,
output $\{\{e \mid N_C^+(e) \geq n - t\}, \perp\}$.

- Otherwise if $\bigwedge_{e \in \mathcal{U}_C} N_C^+(e) > t \wedge N_C^+(e) \geq N_C^-(e)$
or $\bigwedge_{e \in \mathcal{U}_C} N_C^-(e) > t \wedge N_C^-(e) > N_C^+(e)$,
output $\{\{e \mid N_C^+(e) > t \wedge N_C^+(e) \geq N_C^-(e)\}, \perp\}$.
- Otherwise, output $\langle \perp, 0 \rangle$.

Similar to ByzConsensus, the BSC consists of at most $f + 1$ super-rounds, where $f \leq t$. Each peer P_i starts with $S_i^{(1)}$ as its current set. In sequential super-rounds, all peers lead a Gradecast for their candidate set. Like in ByzConsensus, if P_i receives a Gradecast with a confidence value that is not 2, then P_i puts the leader of the Gradecast on its blacklist, and correct peers stop all communication with peers on their blacklist.

At the end of each super-round, peers update their candidate set as follows. Let n' be the number of leaders that gradecasted a set with a non-zero confidence. The new candidate set contains all set elements that were included in at least $\lceil n'/2 \rceil$ sets that were gradecasted with a non-zero confidence value. If all elements occur with a $(n - t)$ -majority, then the next round is the last round. The output of the consensus protocol is the candidate set after the last round—or failure if $f > t$.

We give a correctness proof that generalizes Feldman's proof for Gradecast of single values [24, Section 4.1].

Lemma 1. *If two correct peers send sets $A \neq \perp$ and $B \neq \perp$ respectively in CONFIRM, then $A = B$.*

Proof. Proof by contradiction and counting argument. Assume w.l.o.g. that $e \in A$ and $e \notin B$. At least $n - t$ peers must have echoed a set that includes e to the first peer. Suppose f of these peers were faulty, then at least $n - t - f > t$ good peers included e in the ECHO transmission to the second peer. If $e \notin B$, then $t < N_E(e) < n - t$. In this case, an honest second peer must output $B = \perp$. Contradiction. \square

Theorem 2. *The generalization of Gradecast to sets satisfies the three Gradecast properties.*

Proof. We show that each property holds:

- Property 1 (If $c_i, c_j \geq 1$ then $\hat{V}_i = \hat{V}_j$ for correct P_i and P_j): Assume w.l.o.g. that $e \in \hat{V}_i \setminus \hat{V}_j$. For $e \in \hat{V}_j$, P_i must have received e at least $N_C^+(e) > t$ times in CONFIRM. Given $f \leq t$ failures, at least one honest peer must thus have included e in CONFIRM. According to Lemma 1, then all $n - f$ honest peers must either include e in CONFIRM or send \perp . Because \perp is not a set, this leaves at most all $f \leq t$ faulty peers that can send a set without e . But for $e \notin \hat{V}_j$ we need $N_C^-(e) \geq t + 1$. Contradiction.
- Property 2 (If P_L is correct, then $c_i = 2$ and $\hat{V}_i = \hat{V}_L$ for correct P_i): All $n - f \geq n - t$ good peers ECHO and CONFIRM the same set. By the grading rules, they must output a confidence of 2.
- Property 3 ($|c_i - c_j| \leq 1$ for correct P_i and P_j): Proof by contradiction. Assume w.l.o.g. $c_i = 2$ and $c_j = 0$. $c_i = 2$ implies that for each $x \in \hat{V}_i$ at least $n - t$ peers (and thus $(n - t) - f \geq t + 1$ correct peers) must have sent a set in CONFIRM that includes x . For any $y \notin \hat{V}_i$, $n - t$ peers

(and thus $(n - t) - f \geq t + 1$ correct peers) must have sent a non- \perp set in CONFIRM that excludes y .

Given $c_j = 0$, there must have been an element e such that, $N_C^+(e) \leq t$ and $N_C^-(e) \leq t$ for P_j . However, we just derived that for all elements either $N_C^+(e) > t$ or $N_C^-(e) > t$. Contradiction. \square

Given the Gradecast properties for sets, the correctness argument given by Ben-Or [5] for the Byzantine consensus applies to BSC’s generalization to sets.

As described, the protocol has complexity $O(mnf + fkn^3)$. However, the n parallel set reconciliation rounds in each super-round can be combined by tagging the set elements that are being reconciled in the LEAD, ECHO and CONFIRM rounds with the respective leader L . Because LBA (via $n - \ell_i \leq k$) and bounded Eppstein set reconciliation limit mischief for the combined super-round, each malicious peer can, as leader, *once* cause bounded set reconciliation during the ECHO round to all-to-all transmit at most k extra elements, resulting in a total of $O(fkn^2)$ extra traffic over all $f + 1$ rounds. Before exposing themselves this way, non-leading malicious peers can only cause $O(f^2kn)$ additional traffic during all ECHO rounds. Finally, malicious peers can also cause at most $O(fkn^2)$ traffic in the CONFIRM round. Thus, BSC has overall message complexity of $O(mnf + fkn^2)$.

IV. IMPLEMENTATION

We implemented the BSC protocol in the SET and CONSENSUS services of GUNet (<https://gnunet.org/>).

A. The GUNet framework

GUNet is composed of various components that run in separate operating system processes and communicate via message passing. Components that expose an interface to other components are called *services* in GUNet. The main service used by our implementation is the CADET service, which offers pairwise authenticated end-to-end encryption between all participants. CADET uses a variation of the Axolotl public key ratcheting scheme and double-encrypts using both TwoFish and AES. [25] The resulting encryption is relatively expensive compared to the other operations, and thus dominates in terms of CPU consumption for the experiments.

B. Set reconciliation

Set reconciliation is implemented in the SET service. The SET service provides a generic interface for set operations between two peers; the operations currently implemented are the IBF-based set reconciliation and set intersection [26].

In addition to the operation-specific protocols, the following aspects are handled generically (i.e. independent of the specific remote set operation) in the SET service:

Local set operations

Applications need to create sets and perform actions (iteration, insertion, deletion) on them locally.

Concurrent modifications

While a local set is in use in a network operation, the application may still continue to mutate that set.

To allow this without interfering with concurrent the network operations, changes are versioned. A network operation only sees the state of a set at the time the operation was started.

Lazy copying

Some applications building on the SET service—especially the CONSENSUS service described in the next section—manage many local sets that are large but only differ in a few elements. We optimize for this case by providing a lazy copy operation which returns a logical copy of the set without duplicating the sets in memory.

Negotiating remote operations

In a network operation, the involved peers have one of two roles: The acceptor, which waits for remote operation requests and accepts or rejects them, as well as the initiator, which sends the request.

Our implementation estimates the initial difference between sets only using *strata estimators* as described by Eppstein [8]. However, we compress the strata estimator—which is 60KB uncompressed—using `gzip`. The compression is highly effective at reducing bandwidth consumption due to the high probability of long runs of zeros or ones in the most sparse or most dense strata respectively.

We also use a *salt* when deriving the bucket indices from the element keys. When the decoding of an IBF fails, the IBF size is doubled and the salt is changed. This prevents decoding failures in scenarios where keys map to the same bucket indices even modulo a power of two, where doubling the size of the IBF does not remove the collision.

C. Set-Union consensus

To keep the description of the set-union consensus protocol in the previous section succinct, we merely stated that peers efficiently transmit sets using the reconciliation protocol. However, given that the receiving peer has usually many sets to reconcile against, an implementation needs to be careful to ensure that it scales to large sets as intended.

The key goal is to avoid duplicating full sets and to instead focus on the differences. New sets usually differ in only a few elements, thus our implementation avoids copying entire sets. Instead, in the leader round we just store the set of differences with a reference to the original set. In the ECHO and CONFIRM round, we also reconcile with respect to the set we received from the leader, and not a peer’s current set. In the ECHO round, we only store one set and annotate each element to indicate which peer included or excluded that element. This also allows for a rather efficient computation of the set to determine the \perp -result in the CONFIRM round.

D. Evaluating malicious behavior

For the evaluation, our CONSENSUS service can be configured to exhibit the following types of adversarial behavior:

- *SpamAlways*: A malicious peer adds a constant number of additional elements in every reconciliation.

- *SpamLeader*: A malicious peer adds a constant number of additional elements in reconciliations where the peer is the leader.
- *SpamEcho*: A malicious peer adds a constant number of additional elements in echo rounds.
- *Idle*: Malicious peers do not participate actively in the protocol, which amounts to a crash fault from the start of the protocol. This type of behavior is not interesting for the evaluation, but used to test the implementation with regards to timeouts and majority counting.

For the *Spam-** behaviors, two different variations are implemented. One of them (“*-replace”) always generates new elements for every reconciliation. This is not typical for real applications where the number of stufferable elements ought to be limited by set canonicalization. However, this shows the performance impact in the worst case. The other variation (“*-noreplace”) reuses the same set of additional elements for all reconciliations, which is more realistic for most cases. We did not implement adversarial behaviour where elements are elided, since the resulting traffic is the same as for additional elements, and memory usage would only be reduced.

V. EXPERIMENTAL RESULTS

All of the experiments were run on a single machine with a 24-core 2.30GHz Intel Xeon E5-2630 CPU, and GUNet SVN revision 36765. We used the `gnunet-consensus-profiler` tool, which is based on GUNet’s TESTBED service [27], to configure and launch multiple peers on the target system. We configured the profiler to emulate a network of peers connected in a clique topology (via loopback, without artificial latency). Elements for the set operations are randomly generated and always 64 bytes large.

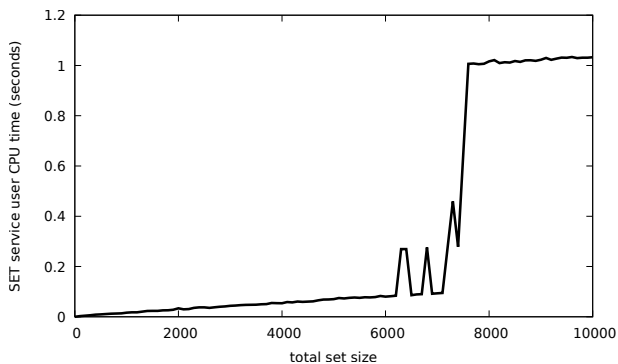
Bandwidth consumption was measured using the statistics that GUNet’s CADET service [25] provides. Processor time was measured using GUNet’s resource reporting functionality, which uses the `wait3` system call for that purpose.

A. Set reconciliation

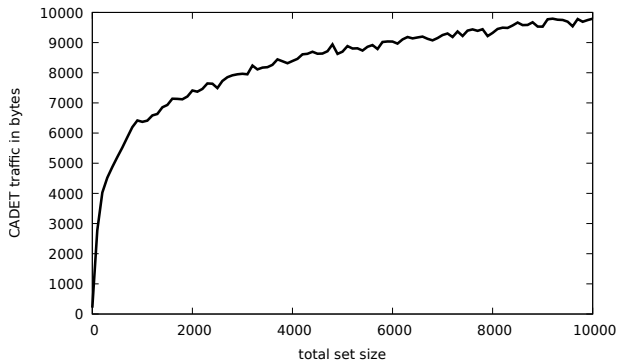
Figure 1 summarizes experimental result for the set reconciliation protocol between two peers. We first measured the behavior of the set reconciliation if identical sets were given to both peers (Figure 1a and 1b). Figure 1a shows that total CPU utilization generally grows slowly as the set size increases. The sudden jump in processing time that is visible at around 7,000 elements can most likely be explained by cache effects. The effect could not be observed when we ran the experiment under profiling tools.

Figure 1b shows that bandwidth consumption does not grow linearly with the total set size, as long as the set size difference between the two peers is small. The logarithmic increase of the traffic with larger sets can be explained by the compression of strata estimators: The k -th strata samples the set with probability 2^{-k} , and for small input sets the strata tends to contain long runs of zeros that are more easily compressed.

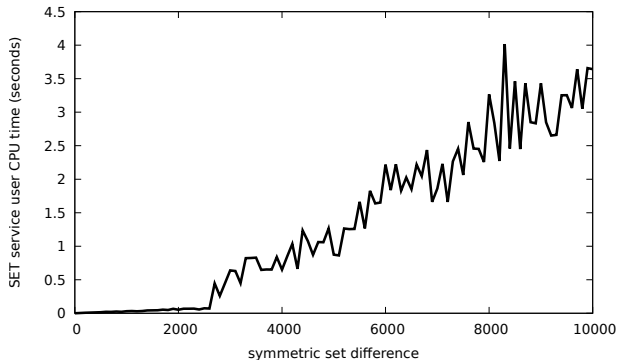
We also measured the behavior of the set reconciliation implementation if the sets differed. Figure 1c and 1d show



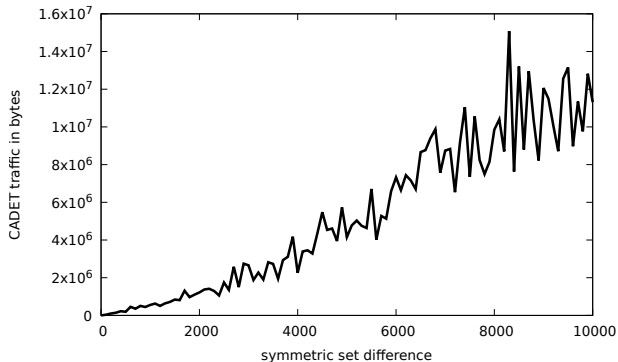
(a) CPU system time for the SET service in relation to total set size.



(b) CADET traffic for the SET service in relation to total set size.



(c) CPU system time for the SET service in relation to symmetric set difference.



(d) CADET traffic for the SET service in relation to symmetric difference.

Fig. 1: SET reconciliation benchmarks. No common elements. Average over five executions.

that—as expected—CPU time and bandwidth do grow linearly with the symmetric difference between the two sets.

Closer analysis of the data (not shown here due to space limitations) suggests that our difference estimator tends to underestimate the difference for larger symmetric differences. The estimation could be improved according to Eppstein et al. by combining strata estimation with MinWise difference estimators [28], which are more accurate for larger differences but less accurate for smaller ones.

B. Set consensus

For our experiments with the BSC implementation, all ordinary peers start with the same set of elements; different sets would only affect the all-to-all union phase of the protocol which does pairwise set reconciliation, resulting in increased bandwidth and CPU consumption proportional to the set difference as shown in the previous section.

As expected, traffic increases cubically with the number of peers when no malicious peers are present (Figure 2a). Most of the CPU time (Figure 2b) is taken up by CADET, which uses expensive cryptographic operations [25]. Since we ran the experiments on a multicore machine, the total runtime follows the same pattern as the traffic (Figure 2c).

We now consider the performance implications from the presence of malicious peers. Figures 3 and Figure 4 show that bandwidth and runtime increase linearly with the additional elements malicious peers can exclusively supply, in contrast to the sub-linear growth for the non-Byzantine case (Figure 1b).

Figure 4 highlights how the different attack strategies impact the number of additional elements that were received during set reconciliations: The number of stuffed elements for the “SpamEcho” behavior is significantly larger than for “SpamLead”, since multiple ECHO rounds are executed for one LEAD round, and the number of stuffed elements is fixed per reconciliation. When malicious peers add extra elements during the LEAD round, the effect of that is amplified, since all correct receivers have to re-distribute the additional elements in the ECHO/CONFIRM round. Even though adding elements in the LEAD round requires the least bandwidth from the leader the effect on traffic and latency is the largest (see Figures 2d and 3).

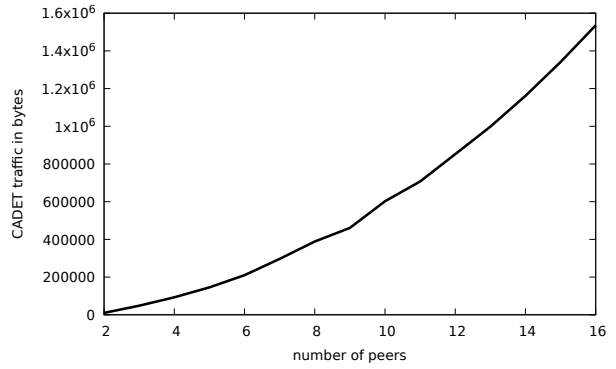
As expected, when the number of stuffed elements is limited to a fixed set, the effect on the performance is limited (“SpamAll-noreplace” in Figures 2d, 3, 4).

VI. OPPORTUNITIES FOR FURTHER IMPROVING BSC

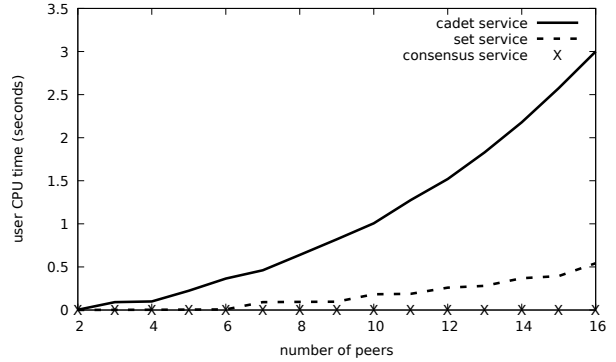
We now discuss some of the key limitations of the current implementation and, how it could be optimised further.

A. Extension to partial synchrony

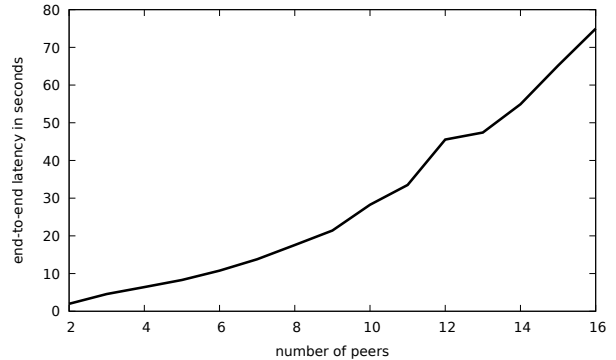
The prototype used in the evaluation only works in the synchronous model. It would be trivial to extend it to the partially synchronous model with synchronous clocks by using the same construction as BPFT [1], namely retrying the protocol with larger round timeouts (usually doubled on each retry) when it did not succeed.



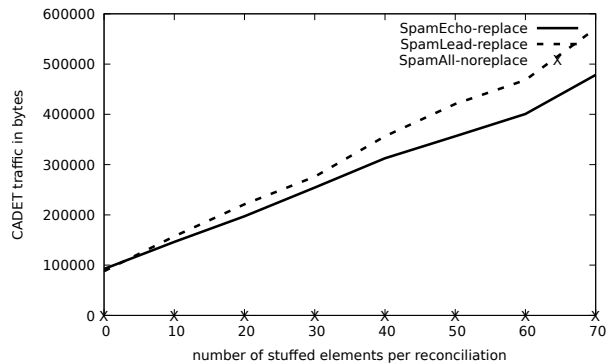
(a) CADET traffic per peer for 100 elements and only correct peers.



(b) CPU of consensus for 100 elements of 64 bytes and only correct peers.



(c) Runtime of consensus for 100 elements of 64 bytes and only correct peers.



(d) CADET traffic for consensus on 100 elements of 64 bytes and one malicious peer with the indicated mode.

Fig. 2: Consensus benchmarks. Average over five executions.

It might be worthwhile to further investigate the Byzantine round synchronization protocols discovered independently by Attya and Dolev [29] as well as Dwork, Lynch and Stockmeyer [6]. Running a Byzantine clock synchronization protocol interleaved with consensus protocol might lead to a protocol with lower latency, since the timeouts are dynamically adjusted instead of being increased for each failed iteration.

B. Persistent data structures

Both the SET and CONSENSUS service have to store many variations of the same set when faulty peers elide or add elements. While the SET service API already supports lazy copying, the underlying implementation is inefficient and based on a log of changes per element with an associated version number. It might be possible to reduce memory usage and increase performance of the element storage by using data structures that are more well suited, such as the persistent data structures described by Okasaki [30].

C. Fast dissemination

Recall that in order to be included in the final set, an element must be sent to at least $t + 1$ peers, so that at least one correct peer will receive the element. In applications of set-union

consensus such as electronic voting, the effort to the client should be minimized, and thus in practice elements might be sent only to $t + 1$ peers, which would lead to large initial symmetric differences between peers.

A possible optimization would be to add another dissemination round that only requires $n \log_2 n$ reconciliations to achieve perfect element distribution when only correct peers are present. The n^2 reconciliations that follow will consequently be more efficient, since no difference has to be reconciled when all peers are correct. In the presence of faulty peers, the optimization adds more overhead due to the additional dissemination round.

More concretely, in the additional dissemination round the peers reconcile with their 2^ℓ -th neighbour (for some arbitrary, fixed order on the peers) in the ℓ -th subround of the dissemination round. After $\lceil \log_2 \rceil$ of these subrounds, the elements are perfectly distributed as long as every peer passed along their current set correctly.

VII. APPLICATION TO SMC

Secure multiparty computation (SMC) is an area of cryptography that is concerned with protocols that allow a group of peers $\mathcal{P} = P_1, \dots, P_n$ to jointly compute a function $y = f(x_1, \dots, x_n)$ over private input values x_1, \dots, x_n without using a trusted third party [31]. Each peer P_i contributes its own input value x_i , and during the course of the SMC protocol, P_i ideally only learns the output y , but no additional information about the other peers' input values. Applications of SMC include electronic voting, secure auctions and privacy-preserving data mining.

SMC protocols often assume a threshold $t < n$ on the amount of peers controlled by an adversary, which is typically either *honest-but-curious* (i.e. tries to learn as much information as possible but follows the protocol) or *actively malicious*. The actively malicious case mandates the availability of Byzantine consensus as a building block [32].³

In practical applications, the inputs typically consist of sets of values that were given to the peers \mathcal{P} by external clients: In electronic voting protocols the peers need to agree on the set of votes; with secure auctions, the peers need to agree on bids, and so on.

In this section, we focus on one practical problem, namely electronic voting. We show how BSC is useful at multiple stages of the protocol, and discuss how our approach differs from existing solutions found in the literature.

A. Bulletin board for electronic voting

The *bulletin board* is communication abstraction commonly used for electronic voting [33], [34]. It is a stateful, append-only channel that participants of the election can post messages to. Participants of the election identify themselves with a

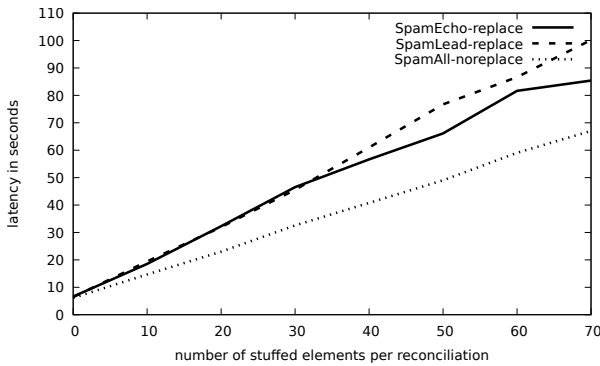


Fig. 3: Latency for consensus with 4 peers on 100 elements of 64 bytes and one malicious peer with the indicated mode. Average over five executions.

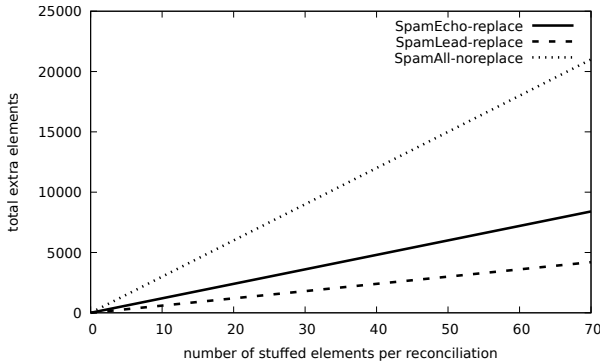


Fig. 4: Total number of extra elements received by each peer for consensus on 100 elements of 64 bytes and one malicious peer with the indicated mode. Average over five executions.

³An attempt has been made to relax the definition of SMC to alleviate this requirement, resulting in a weaker definition that includes non-unanimous *aborts* as a possible result [31]. This definition is mainly useful in scenarios without an non-faulty $2/3$ majority, where Byzantine consensus is not possible in the asynchronous model [6].

public signing key and must sign all messages that they post to the bulletin board. The posted messages are publicly available to facilitate independent auditing of elections.

Existing work on electronic voting either does not provide a Byzantine fault-tolerant bulletin board in the first place [35] and instead relies on trusted third parties, or suggests the use of state machine replication [3].

Some of the bulletin board protocols surveyed by Peters [34] use threshold signatures to certify to the voter that the vote was accepted by a sufficiently large fraction of the peers that jointly provide the bulletin board service. With a naive approach, the message that certifies acceptance by t peers is the concatenation of the peers' individual signatures and thus $O(t)$ bits large. Threshold signature schemes allow smaller signatures, but at the expense of a more complex protocol. Since the number of peers is typically not very large, a linear growth in t is acceptable, which makes the simple scheme sufficient for practical implementations.

It is easy to implement a variant of the bulletin board with set-union consensus. In contrast to traditional bulletin boards, this variant has *phases*, where posted messages are only visible after the group of peers have agreed that a phase is concluded. The concept of phases maps well to the requirements of existing voting protocols. Every phase is implemented with one set-union consensus execution. To guarantee that a message is posted to the bulletin board, it must be sent to at least one correct peer from the group of peers that jointly implements the bulletin board.

B. Distributed threshold key generation and cooperative decryption

Voting schemes as well as other secure multiparty computation protocols often rely on threshold cryptography [36]. The basic intuition behind threshold cryptography is that some operations—such as signing a message or decrypting a ciphertext—should only succeed if a large enough fraction

of some group of peers cooperate. Typically the public key of the threshold cryptosystem is publicly known, while the private key is not known by any entity but reconstructible from the shares that are distributed among the participants, for example with Shamir's secret sharing scheme [37].

Generating this shared secret key either requires a trusted third party, or a protocol for distributed key generation [38], [39]. The former is undesirable for most practical applications since it creates a single point of failure.

In a distributed key generation protocol, each peer contributes a number of *pre-shares*. The peers agree on the set of pre-shares and each peer re-combines them in a different way, yielding the shares of the private threshold key.

In the key generation protocol used for the Cramer et al. voting scheme, the number of pre-shares that need to be agreed upon is quadratic in the number of peers. Every peer needs to know every pre-share, even if it is not required by the individual peer for reconstructing the share, since the pre-shares are accompanied by non-interactive proofs of correctness. Thus the number of values that need to be agreed upon is quadratic in the number of peers, which makes the use of set-union consensus attractive compared to individual agreement.

Even though the pre-shares can be checked for correctness, Byzantine consensus on the set of shares is still necessary for the case when a malicious peer submits an incorrect share to only some peers. Without Byzantine consensus, different correct recipients might exclude different peers, resulting in inconsistent shares.

Similarly, when a message that was encrypted with the threshold public key shall be decrypted, every peer contributes a *partial decryption* with a proof of correctness. While the set of partial decryptions is typically linear in the number of peers, set-union consensus is still a reasonable choice here, this way the whole system only needs one agreement primitive.

C. Electronic voting with homomorphic encryption

Various conceptually different voting schemes use homomorphic encryption; we look at the scheme by Cramer et al. [3] as a modern and practical representative. A fundamental mechanism of the voting scheme is that a set of voting authorities A_1, \dots, A_n establish a threshold key pair that allows any entity that knows the public part of the key to encrypt a message that can only be decrypted when a threshold of the voting authorities cooperate. The homomorphism in the cryptosystem enables the computation of an encrypted tally with only the ciphertext of the submitted ballots. Ballots represent a choice of one candidate from a list of candidate options. The validity of encrypted ballot is ensured by equipping them with a non-interactive zero-knowledge proof of their validity.

It is assumed that the adversary is not able to corrupt more than $1/3$ of the authorities. The voting process itself is then facilitated by all voters encrypting their vote and submitting it to the authorities. The encrypted tally is computed by every authority and then cooperatively decrypted by the authorities and published. Since correct authorities will only agree to

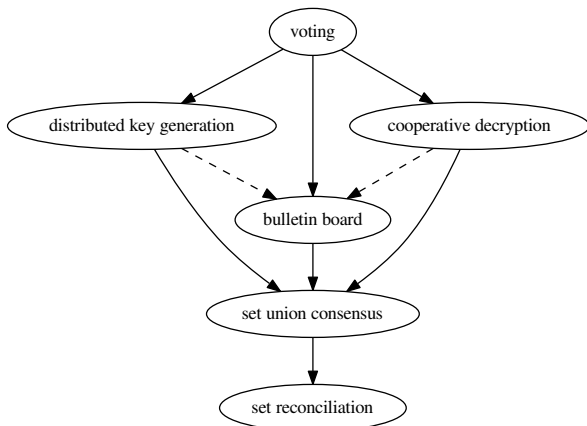


Fig. 5: Relation of different SMC protocols and communication primitives in GNUet. Dashed arrows indicate optional dependencies.

decrypt the final tally and not individual ballots, the anonymity of the voter is preserved. For the voting scheme to work correctly, all correct peers must agree on exactly the same set of ballots before the cooperative decryption process starts, otherwise the decryption of the tally will fail.

Using BSC for this final step to agree on a set of ballots again makes sense, as the number of ballots is typically much larger than the number of authorities. Figure 5 summarizes the various ways how BSC and is used in our implementation [40] of Cramer-style [3] electronic voting.

VIII. CONCLUSION

Given m ballots, n authorities, f Byzantine faults and k ballots exclusively available to the adversary, voting with BSC achieves a complexity of $O(mn + (f+k)n^3)$, which in practice is better than the $O(mn^2)$ complexity of using SMR as m is usually significantly larger than n . Equivalent arguments hold for other applications requiring consensus over large sets. Furthermore, BSC remains advantageous in the absence of Byzantine failures, and the set reconciliation makes it particularly efficient at handling various non-Byzantine faults.

ACKNOWLEDGMENT

This work benefits from the financial support of the Brittany Region (ARED 9174) and the Renewable Freedom Foundation. We thank Jeffrey Burdges and the anonymous reviewers for comments on an earlier draft of this paper.

REFERENCES

- [1] M. Castro, B. Liskov *et al.*, "Practical byzantine fault tolerance," in *OSDI*, vol. 99, 1999, pp. 173–186.
- [2] A. Clement, E. L. Wong, L. Alvisi, M. Dahlin, and M. Marchetti, "Making byzantine fault tolerant systems tolerate byzantine faults," in *NSDI*, vol. 9, 2009, pp. 153–168.
- [3] R. Cramer, R. Gennaro, and B. Schoenmakers, "A secure and optimally efficient multi-authority election scheme," *European transactions on Telecommunications*, vol. 8, no. 5, pp. 481–490, 1997.
- [4] P. Bogetoft, D. L. Christensen, I. Damgård, M. Geisler, T. Jakobsen, M. Krøigaard, J. D. Nielsen, J. B. Nielsen, K. Nielsen, J. Pagter *et al.*, "Secure multiparty computation goes live," in *Financial Cryptography and Data Security*. Springer, 2009, pp. 325–343.
- [5] M. Ben-Or, D. Dolev, and E. N. Hoch, "Simple gradedcast based algorithms," *arXiv preprint arXiv:1007.1049*, 2010.
- [6] C. Dwork, N. Lynch, and L. Stockmeyer, "Consensus in the presence of partial synchrony," *Journal of the ACM (JACM)*, vol. 35, no. 2, pp. 288–323, 1988.
- [7] M. Castro and B. Liskov, "Practical byzantine fault tolerance and proactive recovery," *ACM Transactions on Computer Systems (TOCS)*, vol. 20, no. 4, pp. 398–461, 2002.
- [8] D. Eppstein, M. T. Goodrich, F. Uyeda, and G. Varghese, "What's the difference?: efficient set reconciliation without prior context," in *ACM SIGCOMM Computer Communication Review*, vol. 41, no. 4. ACM, 2011, pp. 218–229.
- [9] L. Lamport, R. Shostak, and M. Pease, "The byzantine generals problem," *ACM Transactions on Programming Languages and Systems (TOPLAS)*, vol. 4, no. 3, pp. 382–401, 1982.
- [10] M. J. Fischer and N. A. Lynch, "A lower bound for the time to assure interactive consistency," DTIC Document, Tech. Rep., 1981.
- [11] R. De Prisco, D. Malkhi, and M. Reiter, "On k -set consensus problems in asynchronous systems," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 12, no. 1, pp. 7–21, 2001.
- [12] N. Malpani, J. L. Welch, and N. Vaidya, "Leader election algorithms for mobile ad hoc networks," in *Proceedings of the 4th international workshop on Discrete algorithms and methods for mobile computing and communications*. ACM, 2000, pp. 96–103.
- [13] M. J. Fischer, N. A. Lynch, and M. Merritt, "Easy impossibility proofs for distributed consensus problems," *Distributed Computing*, vol. 1, no. 1, pp. 26–39, 1986.
- [14] G. Neiger, "Distributed consensus revisited," *Information Processing Letters*, vol. 49, no. 4, pp. 195–201, 1994.
- [15] K. P. Kihlstrom, L. E. Moser, and P. M. Melliar-Smith, "The securering protocols for securing group communication," in *System Sciences, 1998., Proceedings of the Thirty-First Hawaii International Conference on*, vol. 3. IEEE, 1998, pp. 317–326.
- [16] M. K. Reiter, "The rampart toolkit for building high-integrity services," in *Theory and Practice in Distributed Systems*. Springer, 1995, pp. 99–110.
- [17] P.-L. Aublin, R. Guerraoui, N. Knežević, V. Quéma, and M. Vukolić, "The next 700 bft protocols," *ACM Trans. Comput. Syst.*, vol. 32, no. 4, pp. 12:1–12:45, Jan. 2015. [Online]. Available: <http://doi.acm.org/10.1145/2658994>
- [18] M. Abd-El-Malek, G. R. Ganger, G. R. Goodson, M. K. Reiter, and J. J. Wylie, "Fault-scalable byzantine fault-tolerant services," *ACM SIGOPS Operating Systems Review*, vol. 39, no. 5, pp. 59–74, 2005.
- [19] R. Kotla, L. Alvisi, M. Dahlin, A. Clement, and E. Wong, "Zyzyva: speculative byzantine fault tolerance," in *ACM SIGOPS Operating Systems Review*, vol. 41, no. 6. ACM, 2007, pp. 45–58.
- [20] P. Feldman and S. Micali, "Optimal algorithms for byzantine agreement," in *Proceedings of the twentieth annual ACM symposium on Theory of computing*. ACM, 1988, pp. 148–161.
- [21] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors," *Communications of the ACM*, vol. 13, no. 7, pp. 422–426, 1970.
- [22] M. Mitzenmacher and R. Pagh, "Simple multi-party set reconciliation," *arXiv preprint arXiv:1311.2037*, 2013.
- [23] N. Evans, B. Polot, and C. Grothoff, "Efficient and secure decentralized network size estimation," in *Proceedings of the 11th international IFIP TC 6 conference on Networking-Volume Part I*. Springer-Verlag, 2012, pp. 304–317.
- [24] P. N. Feldman, "Optimal algorithms for byzantine agreement," Ph.D. dissertation, Massachusetts Institute of Technology, 1988.
- [25] B. Polot and C. Grothoff, "Cadet: Confidential ad-hoc decentralized end-to-end transport," in *Ad Hoc Networking Workshop (MED-HOC-NET), 2014 13th Annual Mediterranean*. IEEE, 2014, pp. 71–78.
- [26] S. Tarkoma, C. E. Rothenberg, and E. Lagerspetz, "Theory and practice of bloom filters for distributed systems," *Communications Surveys & Tutorials, IEEE*, vol. 14, no. 1, pp. 131–155, 2012.
- [27] S. H. Totakura, "Large scale distributed evaluation of peer-to-peer protocols," Master's Thesis, Technische Universität München, Garching bei München, 06/2013 2013.
- [28] P. Li and A. C. König, "Theory and applications of b-bit minwise hashing," *Communications of the ACM*, vol. 54, no. 8, pp. 101–109, 2011.
- [29] C. Atiya, D. Dolev, and J. Gil, "Asynchronous byzantine consensus," in *Proceedings of the third annual ACM symposium on Principles of distributed computing*. ACM, 1984, pp. 119–133.
- [30] C. Okasaki, *Purely functional data structures*. Cambridge University Press, 1999.
- [31] S. Goldwasser and Y. Lindell, "Secure multi-party computation without agreement," *Journal of Cryptology*, vol. 18, no. 3, pp. 247–287, 2005.
- [32] J. Saia and M. Zamani, "Recent results in scalable multi-party computation," in *SOFSEM 2015: Theory and Practice of Computer Science*. Springer, 2015, pp. 24–44.
- [33] J. D. C. Benaloh, *Verifiable secret-ballot elections*. Yale University, Department of Computer Science, 1987.
- [34] R. Peters, "A secure bulletin board," Master's Thesis, Technische Universiteit Eindhoven, 2005.
- [35] B. Adida, "Helios: Web-based open-audit voting," in *USENIX Security Symposium*, vol. 17, 2008, pp. 335–348.
- [36] Y. G. Desmedt, "Threshold cryptography," *European Transactions on Telecommunications*, vol. 5, no. 4, pp. 449–458, 1994.
- [37] A. Shamir, "How to share a secret," *Communications of the ACM*, vol. 22, no. 11, pp. 612–613, 1979.
- [38] P.-A. Fouque and J. Stern, "One round threshold discrete-log key generation without private channels," in *Public Key Cryptography*. Springer, 2001, pp. 300–316.
- [39] T. P. Pedersen, "A threshold cryptosystem without a trusted party," in *Advances in CryptologyEUROCRYPT91*. Springer, 1991, pp. 522–526.
- [40] F. Dold, "Cryptographically secure, distributed electronic voting," Bachelor's Thesis, Technische Universität München, 2014.