

# DAMYSUS: Streamlined BFT Consensus Leveraging Trusted Components

Submission #157

## Abstract

Recently, *streamlined* Byzantine Fault Tolerant (BFT) consensus protocols, such as HotStuff, have been proposed as a means to circumvent the inefficient view-changes of traditional BFT protocols, such as PBFT. Several works have detailed trusted components, and BFT protocols that leverage them to tolerate a minority of faulty nodes and use a reduced number of communication rounds.

Inspired by these works we identify two basic trusted services, respectively called the Checker and Accumulator services, which can be leveraged by streamlined protocols. Based on these services, we design Damysus, a streamlined protocol that improves upon HotStuff’s resilience and uses less communication rounds. In addition, we show how the Checker and Accumulator services can be adapted to develop Chained-Damysus, a *chained* version of Damysus where operations are pipelined for efficiency. We prove the correctness of Damysus and Chained-Damysus, and evaluate their performance showcasing their superiority compared to previous protocols.

## 1 Introduction

Consensus is a crucial building block of many distributed systems including state machine replication and blockchains. Such systems comprise multiple software and hardware components that are bound to eventually fail, potentially causing system malfunction or unavailability. To this end, a distributed system relies on a consensus protocol to agree on actions critical to the system’s correct operation.

In particular, Byzantine Fault Tolerant (BFT) consensus protocols allow systems to withstand arbitrary failures [1, 2]. However, being able to tolerate arbitrary (Byzantine) component failures, induces a non-negligible system overhead, namely in the required number of nodes that should be present in the system as well as the communication complexity. For example, traditional BFT protocols, such as PBFT [3], typically require  $3f+1$  nodes to tolerate up to  $f$  Byzantine faulty nodes. Moreover, albeit being efficient in the normal case, i.e., when the current leader is correct, traditional BFT protocols use complex “view-change” operations to replace possibly malicious (compromised) leaders and often require the nodes to “synchronize” and transfer their states.

In order to circumvent such complex view changes, recent “streamlined” BFT protocols such as HotStuff [4] and Streamlet [5] rotate the leader on each command. Instead of performing a state transfer when switching to another leader, they utilize a simpler locking mechanism combined with an

additional communication phase (compared to PBFT-like protocols) to guarantee both safety and liveness. Hence, streamlined protocols are conceptually simpler. Another distinguishing feature of streamlined protocols is their linear message complexity as opposed to the quadratic all-to-all communication of traditional solutions. Nevertheless, streamlined solutions still require at least  $3f+1$  nodes to tolerate  $f$  Byzantine nodes, and they introduce an additional communication phase that inflicts higher system latency.

For some applications being able to tolerate more than one-third of Byzantine faults is crucial [6]. Traditional BFT protocols have typically relied on *hybrid* solutions to reduce the number of nodes required to ensure system safety as well as to decrease system latency. Those solutions rely on secure elements that execute key functionalities in trusted execution environments such as inside Intel SGX enclaves [7]. For example, MinBFT [8] relies on a trusted message counter to increase the resilience of PBFT, permitting a system of size  $N$  to tolerate  $\lfloor \frac{N-1}{2} \rfloor$  faulty nodes while simultaneously reducing the number of communication rounds in normal case operations from 3 in PBFT to 2 in MinBFT.

While hybrid solutions have been largely explored in the context of traditional BFT protocols focusing on minimizing the functionalities dedicated to trusted components, they have not yet been explored to the same extent in the context of streamlined protocols. To this end, this paper investigates the design of hybrid streamlined solutions focusing on HotStuff-like protocols. Interestingly, we show that the use of trusted message counters (simplest known trusted components used in traditional BFT protocols) [8, 9, 10] is not sufficient for simultaneously making HotStuff-like protocols more resilient and faster (see §4).

Still adhering to a minimalist approach when designing trusted components, we introduce two simple trusted services called CHECKER and ACCUMULATOR, which contribute respectively to increasing a streamlined protocol’s resilience and to decreasing its latency. Roughly speaking, our CHECKER service helps ensure that nodes cannot vote for conflicting blocks, or lie about the blocks they have already voted for, while our ACCUMULATOR service guarantees that leaders can only propose blocks that are consistent with previous votes. We then present Damysus,<sup>1</sup> a streamlined protocol, based on *basic* HotStuff, which makes use of the CHECKER and ACCUMULATOR services. Damysus requires  $2f+1$  nodes to tolerate  $f$  Byzantine failures and can terminate in 2 communication phases as opposed to 3 in basic HotStuff.

<sup>1</sup>Damysus was the fastest of all the Giants in the Greek mythology.

In this sense, Damysus simultaneously tolerates a minority of Byzantine nodes and reduces the number of communication phases to its known minimum. The introduced trusted services are simple and do not require changes to HotStuff's overall behavior, preserving linear view change and optimistic responsiveness [4] where nodes reach consensus at the pace of the actual (vs. maximum) network delay.

We also generalize our results showcasing how we can seamlessly transfer Damysus to Chained-Damysus using the exact same trusted services to also support a *chained* version of HotStuff, where operations are pipelined so that multiple blocks are concurrently processed.

In short, our contribution can be summarized as follows.

1. We introduce two simple trusted services, which when applied to HotStuff (without changing its overall behavior), result in Damysus, a streamlined consensus protocol with improved resilience and latency. Damysus requires  $2f+1$  nodes and terminates in 2 phases.
2. We also introduce Chained-Damysus, a chained version of Damysus for higher throughput.
3. We provide theoretical proofs that all protocols are safe and live.
4. We show the performance superiority of our protocols compared to baseline HotStuff-like protocols. For example, using 256B payload messages, Damysus has an average throughput increase of 87.5% and an average latency decrease of 45% compared to basic HotStuff. Chained-Damysus, compared to chained HotStuff, has a throughput increase of 50.5% and latency decrease of 32.1%.

The rest of the paper is organized as follows. §2 reviews the related work. §3 provides background on HotStuff. §4 presents the functionality of our CHECKER and ACCUMULATOR trusted services. §5 describes our system model. §6 introduces Damysus and gives the intuition behind its design principles. §7 presents Chained-Damysus. §8 provides a performance evaluation. §9 concludes this paper.

## 2 Related Work

We summarize our analysis of the related work in Table 1 and elaborate further in the following.<sup>2</sup>

Hybrid BFT protocols, which leverage the use of trusted hardware components to provide better fault-tolerance and performance, have been studied in the context of leader-based BFT protocols [13, 14, 8, 9, 10, 15] as well as leaderless ones [16]. In this section, we focus on the former.

Previous literature explored techniques such as attested append-only memory (A2M) [13], and trusted incrementer (TrInc) [14] where multiple trusted monotonic counters are used to prevent equivocation and tolerate more faults. Later,

MinBFT [8] presented an approach that provides the same properties using a single trusted monotonic counter while additionally reducing the number of phases in PBFT from 3 to 2 in the normal case operation. CheapBFT [9], as well as its generalization ReBFT [17], further improved this result by enabling the system to “optimistically” operate with as little as  $f+1$  active replicas, while  $f$  other replicas stay passive during the normal case operation. Hybster [10] also explored the possibility to parallelize instances with TrIncX, a TrInc-like trusted monotonic counter. FastBFT [11] is similar to CheapBFT in the sense that it is both optimistic (making use of  $f+1$  active nodes in the normal case operation) and relies on a trusted monotonic counter. In addition, FastBFT relies on a TEE-based secret sharing mechanism, which compounded with trusted monotonic counters, allows reducing the message complexity of the protocol. Recently, TBFT [15] was proposed to improve performance by requiring backups to communicate only with the leader (with complexity  $O(N)$  as in HotStuff) rather than an all-to-all communication pattern (with complexity  $O(N^2)$  as in PBFT). In addition to relying on trusted monotonic counters, TBFT uses a trusted message sharing mechanism to generate quorum certificates representing messages collected from  $f+1$  replicas.

In the context of streamlined BFT protocols, hybrid solutions were first mentioned for LibraBFT [18], a protocol built on HotStuff, to possibly reduce the attack surface. Fast-HotStuff [19] is a version of HotStuff with one less communication phase, which does not rely on trusted components. Instead, leaders send proofs that the blocks they extend are the highest received blocks. This requires larger messages (containing an aggregated vector of  $2f+1$  quorum certificates) but improves latency and prevents from creating forks.

Later, HotStuff-M and VABA-M [12] proposed the use of trusted components in the form of trusted logs to improve fault tolerance without worsening the communication complexity of the underlying BFT protocol: one trusted log is used for each protocol phase with an additional log for tracking views. This approach also requires changes to the underlying BFT protocol: the protocol would need to use and maintain expander graphs to diffuse messages. Such graphs introduce more network traffic, extra storage overhead and their impact on throughput/latency remains open.

Simultaneously achieving improved communication complexity with optimal resilience of streamlined BFT protocols remains an open challenge. This drives us to propose the first hybrid streamlined BFT system, which tolerates a minority of Byzantine faults while reducing the number of communication phases to improve the latency and throughput of the underlying algorithm. In addition, unlike HotStuff-M, our approach does not change the overall behavior of the underlying protocol, retaining its conceptual simplicity and does not introduce extra network traffic in any phase. It also arguably requires minimal storage (including only a counter

<sup>2</sup>In Table 1, the communication steps exclude the interaction with clients, which are common to all protocols: requests sent by clients to replicas, and replies sent by replicas to clients.

**Table 1.** Comparative analysis of Damysus and Chained-Damysus with the related work. View-change communication steps are indicated between parentheses when applicable. Message counts include self-messages.

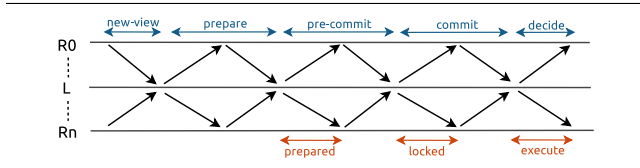
	# replicas	Commun. steps	#Msgs normal case	#Msgs view change	Optimistic exec.?	Trusted component & Storage complexity
PBFT [3]	$3f+1$	3 (+2)	$18f^2+15f+3$	$9f^2+6f+1$	No	-
FastBFT [11]	$f+1$ act. & $f$ pass.	5 (+3)	$6f+5$	$8f^2+8f+2$	Yes	Secret generation - Constant
MinBFT [8]	$2f+1$	2 (+3)	$4f^2+6f+2$	$8f^2+6f+1$	No	Trusted counter - Constant
CheapBFT [9]	$f+1$ act. & $f$ pass.	3 (+3)	$2f^2+4f+2$	$8f^2+6f+1$	Yes	Trusted counter - Constant
HotStuff [4]	$3f+1$	8	$24f+8$	-	No	-
HotStuff-M [12]	$2f+1$	11	$(24+9d)f+(8+3d)$	-	No	Append-only logs - Linear with # msgs
Damysus	$2f+1$	6	$12f+6$	-	No	Checker & Accumulator - Constant
Chained-Damysus	$2f+1$	6	$12f+6$	-	No	Checker & Accumulator - Constant

representing the current view/phase and some information regarding the locally prepared and locked blocks).

### 3 HotStuff in a Nutshell

HotStuff [4] is a BFT protocol whose communication complexity grows linearly with the total number  $N$  of nodes, making it remarkably efficient. HotStuff requires  $N \geq 3f+1$  nodes to tolerate  $f$  Byzantine faults. Nodes build a chain of blocks by voting for extensions, which are proposed by the leaders of *views* (successive rounds).

HotStuff comes in two versions: (1) *Basic* HotStuff that adds a single block per view, and (2) *chained* (or *pipelined*) HotStuff that allows pipelined votes to simultaneous progress on several blocks per view.

**Figure 1** Communication phases in HotStuff.

**Communication phases** To execute a block, HotStuff employs several *phases* per view  $v$ , which Fig. 1 illustrates. In that figure, time progresses from left to right, and  $L$  is the leader of the view  $v$  among the nodes  $R_0, \dots, R_n$ . HotStuff is referred to as a 3-phase protocol, as it has 3 *core* phases: prepare, pre-commit, and commit to agree on blocks, complemented by 2 additional half phases: new-view to submit latest prepared blocks, and decide to execute blocks once it has been judged safe to do so.

(1) In the **new-view** phase<sup>3</sup>, each node sends its latest prepared block and view number to the leader;

(2) In the **prepare** phase, after having received  $2f+1$  pairs of a block and a view number, the leader creates a new block by extending the block with the largest view number among all received pairs. The leader proposes a block by sending it to all other nodes who will vote for it if it is correct.

<sup>3</sup>Different from [4], we group new-view messages in a new-view phase because they are necessary when a view is interrupted.

(3) In the **pre-commit** phase, the leader gathers  $2f+1$  votes for its proposal and marks the proposed block as being *pre-prepared* (marked in orange in Fig. 1). The collection of  $2f+1$  votes forms a quorum certificate indicating that the proposed block is prepared. The leader forwards this certificate to the other nodes, which mark the proposed block as being prepared after verifying the quorum certificate. Once a node has prepared a block, it votes for it in this phase.

(4) In the **commit** phase, the leader gathers a quorum certificate on a prepared block, *locks* it (marked in orange in Fig. 1) and forwards the quorum certificate to the other nodes so that they lock the same block. Once a node has locked a block, it votes for it in this phase.

(5) In the **decide** phase, the leader gathers a quorum certificate of  $2f+1$  votes on a locked block, *executes* that block (marked in orange in Fig. 1) and forwards the certificate to the other nodes so that they also execute the locked block.

Basic HotStuff uses a different leader in each view, and each view is in charge of proposing a single extension to the current blockchain.

**Locking scheme** A new-view message only contains the certificate of what a node believes to be the latest prepared block, rather than all certificates for all prepare messages since the last checkpoint as in traditional protocols such as PBFT [3] and BFT-SMaRt [20]. While this makes view-changes more efficient in terms of communication complexity and message size, as explained in [4, §4.4], HotStuff requires a “locking” phase between the prepare and execution phase to guarantee safety and liveness, where blocks are locked only after they are prepared (for liveness) and before they are executed (for safety). In addition, as the messages sent to extend a block in the prepare phase do not prove the non-existence of a conflicting extension on the same block, it is possible that different subsets of nodes have prepared and locked conflicting extensions of the same block. In that case, only one of those extensions will be accepted (safety), and nodes that locked on a conflicting block have to be allowed to replace it (liveness). To achieve this, a node accepts a block proposed by a leader if it extends the latest block it locked (necessary for safety), or if it extends a block that was prepared at a view higher than the view of the latest block



it locked (necessary for liveness)—the conjunction of those two predicates is referred to as the `SAFENODE` predicate [4].

**Liveness** HotStuff guarantees liveness after GST, the Global Stabilization Time (see §5). This is done through a mechanism that allows correct nodes to eventually be in the same view for long enough to make progress. This mechanism can, for example, be an exponential backoff mechanism where nodes set a timer when starting a new view; exponentially increase their timeout when a time out occurs; and cancel the timer and linearly decrease the timeout when views succeed. In addition, a rotating leader election scheme ensures that there is a single leader per view, and that a correct leader is always eventually elected. HotStuff also guarantees responsiveness – after GST, with a correct leader, nodes make progress by waiting for  $N-f$  messages.

## 4 CHECKER and ACCUMULATOR Services

Previous efforts in the literature presented BFT protocols that increased the resilience of PBFT using increasingly simpler trust assumptions starting from tamperproof distributed components connected through dedicated channels separated from the payload network [21], moving to trusted append-only logs [13], and finally reaching the minimal trusted counters [14, 8]. In this sense, a reader might wonder whether the simplest established trusted components, i.e., a single trusted counter, can be used to increase the resilience of streamlined protocols. We answer this question by the negative in the following.

First let us recall how these simple trusted counters are used. On each client request the leader calls its trusted counter, which generates a unique number along with a signed certificate that proves that the number is assigned to that request. The other nodes go through a similar process when replying to a message from the leader. This way, all correct replicas can safely consider that all messages with a given identifier are identical. Modifying PBFT to allow each node to leverage a trusted counter increases the resilience of PBFT from  $\lfloor \frac{N-1}{3} \rfloor$  to  $\lfloor \frac{N-1}{2} \rfloor$ , where  $N$  is the total number of nodes, and reduces the number of phases from 3 to 2. This is achieved as new-view messages contain view-change certificates that require even Byzantine nodes to include the history of messages they have prepared. However, contrary to PBFT, HotStuff does not need or use view-change certificates. It is up to the new leader to chose which prepared block to extend. This prepared block will have a certificate, but a Byzantine leader could produce an old certificate, and the backups would not have a way to verify whether the leader correctly picked the latest prepared block. A simple trusted message counter is therefore not sufficient to lower HotStuff's number of replicas from  $3f+1$  to  $2f+1$ , i.e., to increase its resilience, let alone to reduce its number of phases.

Let us consider a version of HotStuff with  $2f+1$  replicas only and the same number of phases as basic HotStuff, where

each node is augmented with a trusted message counter. We demonstrate the following non-safe scenario, where  $f=1$ , with nodes  $i, j, k$ . A single Byzantine node  $i$  can convince a leader  $j$  to execute a block  $b$  and then become leader and take advantage of slow communication to convince the third node  $k$ , which has not heard about  $b$ , to execute a conflicting block. The problem is that a HotStuff node  $k$  cannot be expected to receive all messages from  $i$ , as these messages have to go through the leaders, and therefore might not know that  $i$  voted for  $b$ .

To this end, our approach relies on augmenting trusted counters with additional secure storage to guarantee liveness and safety, which later leads to our `CHECKER` service. In order to increase the resilience, the additional secure storage would need to persist both prepared and locked blocks. The latter are needed as otherwise a Byzantine node could lead  $f$  correct but late nodes to lock and then further execute conflicting blocks. However, using only this simple yet additional storage is not enough to reduce the number of communication rounds but merely to increase resilience.

In order to also reduce the number of communication rounds (for better latency) a leader should certify that it has selected the block with the highest view from the new-view messages it has received. This functionality is captured by our `ACCUMULATOR` service, which forces a leader to chose the latest prepared block among the received ones. This is valid with quorums of size  $f+1$ , as using the `CHECKER` nodes are forced to relay the latest prepared block (as they store it). Interestingly, when using this `ACCUMULATOR` function, the original suggested trusted counter with augmented storage (to store prepared and locked blocks) becomes slightly simpler: there is no need to store the locked blocks anymore, leading to our current `CHECKER`.

Therefore, our approach relies on `CHECKER` and `ACCUMULATOR` as two local trusted components, which must be available at all nodes, to achieve respectively two goals: (1) increased Byzantine resilience, and (2) reduced latency. Each component is identified by a unique identifier stored with the component. Both components perform key operations, and produce *certificates* to guarantee that those operations have been performed. These certificates are signed messages, which are produced and verified using *private/public keys* stored within the services (§5 provides further details on signatures and §6 shows the certificates used by Damysus). Private keys need to be kept confidential to prevent hosts of trusted components from forging commitments.

### 4.1 CHECKER

The `CHECKER` trusted service is used by Damysus for increased Byzantine resilience, allowing to tolerate more Byzantine nodes. In each execution of Damysus, all replicas interact with their local `CHECKER`. It has two purposes: (1) it assigns to each message a unique identifier thanks to a monotonic counter, which is incremented each time a message is signed,

to prevent nodes from equivocating; and (2) it stores information about relevant blocks, here about prepared and locked blocks, to guarantee that nodes cannot lie about the blocks they have last prepared and locked. To this end, this service maintains the following state, which needs to be securely stored in memory, and can only be modified by the service:

- a *monotonically increasing counter* keeping track of the current view and phase
- some information regarding the *latest prepared and locked blocks*. As explained in §6.3, for the purpose of Damysus, we store a view number and a hash value per prepared and locked block.

Furthermore, it provides the following interface:

- **TEEprepare**(*block, certificate*): this function takes a block  $b$  (we see below that blocks can often be replaced by their hash values) and a certificate, which includes a block  $b'$ , and a proof that  $b'$  is one of the latest prepared blocks. It then returns a new certificate that contains the current value of the monotonic counter, certifying that  $b$  is a safe (w.r.t. **SAFE**NODE) proposal extending  $b'$ , after which this counter is increased.
- **TEEstore**(*certificate*): this function takes a certificate that some block  $b$  extends some other block  $b'$  as generated by the **CHECKER** service, verifies it, stores  $b$  as the latest prepared or locked block in the state, depending on the phase, and generates a new certificate that this operation has been executed.
- **TEEsig**( $\epsilon$ ): this function generates a certificate for the stored latest prepared block.

## 4.2 ACCUMULATOR

The **ACCUMULATOR** trusted service is used by Damysus to reduce latency by reducing the number of communication phases. In each execution of Damysus, only the leader interacts with its **ACCUMULATOR**. Its goal is to certify that some block has the highest view among a given set of blocks. This is used in Damysus to ensure that a leader at anytime proposes the latest prepared block in the system (if it exists). As for **CHECKER**, it also generates certificates, however, the ones it generates are not tied to a monotonic counter, but instead indicate how many messages have been accumulated so far. We therefore call such a certificate an *accumulator* here.

This service does not require further storage, and provides the following interface:

- **TEEstart**(*certificate*): this function takes a certificate that some block  $b$  is the latest prepared block (see above), and generates an initial *accumulator*, indicating that only  $b$  was considered so far.
- **TEEaccum**(*accumulator, certificate*): this function takes an *accumulator* certifying that some block  $b$  has the highest view  $v$  among  $k$  blocks signed by  $k$  different nodes denoted  $K$  here, as well as a certificate signed by  $j$  that some block  $b'$  is the latest prepared block,

prepared in view  $v'$ , and generates a new accumulator certifying that  $b$  has the highest view  $v$  among  $k+1$  blocks if  $v \geq v'$  and  $j \notin K$ .

- **TEEfinalize**(*accumulator*). The accumulators considered so far have to keep track of the nodes that have signed the certificates used to generate those accumulators. Once a given *threshold* has been reached, **TEEfinalize** replaces the set of node ids  $\{i_1, \dots, i_k\}$  stored in such an accumulator by their numbers  $k$

This service is intended to be used by leaders as follows.

- (1) Before calling the service, collect  $C$ , a list of certificates signed by different nodes, and select from  $C$  the certificate  $c$  with the highest view (here the view at which the block certified by that certificate was prepared).
- (2) Call **TEEstart**( $c$ ) to create an initial accumulator.
- (3) Iterate **TEEaccum** on the remaining certificates in  $C$ , expanding the accumulator with the ids of the nodes that signed those certificates.
- (4) Call **TEEfinalize** on that accumulator.

## 4.3 On the Use of CHECKER and ACCUMULATOR

For space limitations, we use in this paper the **CHECKER** and **ACCUMULATOR** combined, since their combination leads to the most efficient hybrid protocol. However, we point the reader to the extended version of this paper [22] for the individual hybrid protocols that are based solely on either the **CHECKER** (to increase resilience) or the **ACCUMULATOR** service (to decrease latency). Remark that when used in combination with the **ACCUMULATOR** service (§4.2), the **CHECKER** service does not require to store locked blocks (leading to simpler **TEEprepare** and **TEEstore** functions) because leaders are constrained to extend the highest latest prepared block out of a quorum of such messages received in new-view messages.

## 5 System Model

**Replicas and leader** We consider a static system consisting of  $N$  replicas out of which at most  $f$  can be faulty (i.e., Byzantine). Depending on the protocol,  $N$  will either be  $3f+1$  or  $2f+1$ . For simplicity, we refer to a replica by using a unique id. We assume that each view has a unique leader, which is chosen deterministically and known to all nodes.

**Trusted components** Each replica executes trusted components that provide the **CHECKER** and **ACCUMULATOR** services, resulting in a *hybrid* fault model, where at each faulty node all components can be tampered with except the ones providing these trusted services.

**Communications** We assume that replicas communicate by exchanging messages over a fully connected communication network. Communications are reliable (i.e., messages are not lost). We adopt the partial synchrony model, where there is a known bound  $\Delta$  and an unknown Global Stabilization Time (GST), such that after GST, all communications arrive within time  $\Delta$  [23].

**Signatures** Replicas and trusted components rely on an asymmetric signature scheme. A digital signature  $\sigma$  is generated via the **SIGN** function, and verified using the **VERIFY** function. We assume that a digital signature contains the identity of the signing replica or component, which is obtained using  $\sigma.\text{id}$ .

**Blocks** A block  $b$  contains transactions submitted by clients. Our protocols work at the block level, and we therefore leave abstract the internal details of transactions, which are mostly application specific. We assume a cryptographic secure hash function **H** that is used to hash blocks. We write  $h$  for the hash value of a block. We write  $b > h$  when  $b$  is a direct extension of a block with hash value  $h$ . We also write  $b_1 > b_2$  for  $b_1 > \mathbf{H}(b_2)$ . The relation  $>$  can easily be checked if, e.g., blocks store the hash values of the blocks they extend. We write  $>^+$  for its transitive closure, and  $>^*$  for its reflexive and transitive closure. When  $b_2 >^* b_1$ , we say that  $b_2$  is a *descendant* of  $b_1$  and that  $b_1$  is an *ancestor* of  $b_2$ . We say that a block  $b_1$  *conflicts with* a block  $b_2$  if  $\neg b_2 >^+ b_1$  and  $\neg b_1 >^+ b_2$ . We also assume a **createLeaf** function that creates a new block extending a parent block (or simply its hash value) with client transactions. Let **G** stand for the genesis block.

## 6 Damysus

This section presents Damysus, a 2-phase protocol that tolerates  $\lfloor \frac{N-1}{2} \rfloor$  Byzantine replicas ( $f$  out of  $2f+1$  replicas). Damysus requires each node to be equipped with an instance of each of our trusted services described in §4 to achieve respectively two goals: (1) an instance of the **CHECKER** service for increased resilience against Byzantine nodes; and (2) an instance of the **ACCUMULATOR** service for reduced latency.

We next present an overview of Damysus and discuss how the **CHECKER** and **ACCUMULATOR** components are used.

### 6.1 Overview

At the beginning of every view its leader awaits to receive  $f+1$  new-view messages from other processes. A new-view message from process  $p$  contains a quorum certificate of the latest prepared block at  $p$ . A leader then selects from these new-view messages, the certificate  $c$  for the prepared block with the highest view (this view being the view at which the block was prepared), and submits it to its **ACCUMULATOR** component, along with the other  $f$  new-view messages. If  $c$  was indeed a certificate for the highest prepared blocks among the  $f+1$  submitted certificates, and those  $f+1$  certificates have been produced by distinct nodes, then **ACCUMULATOR** generates a certificate of its own (a signed message), called an *accumulator* here, guaranteeing that the above has been checked. A leader then uses this *accumulator* to build its proposal. A leader's proposal is first submitted to its local **CHECKER** component, which validates a single proposal per view, and then it is sent out to the other replicas in a prepare

message. All other nodes only accept proposals validated by the current leader's **CHECKER**.

Upon receiving the leader's proposal  $b$ , every replica uses a modified **SAFE**NODE predicate to check whether it should accept  $b$ . This **SAFE**NODE checks that  $b$  extends the prepared block certified by the **ACCUMULATOR**. If the leader's proposal is accepted, the replicas also use their **CHECKER** components to verify the *accumulator*, and in turn sign the proposal. Those signatures are partially signed prepare votes, which the replicas send to the leader.

When the leader receives  $f+1$  prepare votes for its proposal, it combines them into a certificate, which it broadcasts in the pre-commit phase. Upon receipt of a prepare certificate for  $b$  from the leader, the nodes use their **CHECKER** components to store the view at which  $b$  is prepared, i.e., the current view, and for code simplicity  $b$ 's hash value too. Storing only views would require trusted components to perform further computations to check the “correctness” of those blocks, which can be omitted by simply storing the hash values of those blocks. Prepared blocks are stored in **CHECKER** to ensure that nodes (even Byzantine ones) relay their prepared messages in new-view messages. Once a **CHECKER** has stored a prepared block, it emits a signature guaranteeing that this operation has been done. Such a signature is a partially signed pre-commit vote, which is sent to the leader.

When the leader receives  $f+1$  pre-commit votes, it combines them into a certificate and sends it in a decide message to all other replicas. Upon receiving a decide message, a replica executes the transactions, increments its view number, and starts the next view.

### 6.2 Definitions

In what follows, we introduce some definitions that we use later to formally describe Damysus.

**Phases** Messages generated by trusted components include a tag  $ph \in \{\text{nv\_p}, \text{prep\_p}, \text{pcom\_p}\}$  indicating the phase in which they were generated: Those are shorthands for new-view, prepare, and pre-commit respectively.

**Steps** A step  $(v, ph)$  is a pair of a view  $v$  and a phase  $ph$ . To progress through views and phases, nodes increment steps as follows:  $(v, \text{nv\_p})++ = (v, \text{prep\_p})$ ,  $(v, \text{prep\_p})++ = (v, \text{pcom\_p})$ , and  $(v, \text{pcom\_p})++ = (v+1, \text{nv\_p})$ .

**Commitments** A commitment  $\phi$  (generated by **CHECKER** components) is of the form  $\langle h, v, h', v', ph, \vec{\sigma}^n \rangle$ , where  $h$  and  $h'$  are hash values of blocks,  $v$  and  $v'$  are view numbers,  $ph$  is a phase, and  $\vec{\sigma}^n = [\sigma_1, \dots, \sigma_n]$  is a list of signatures on data  $(h, v, h', v', ph)$ . We simply write  $\vec{\sigma}$  when the number of signatures in the list is irrelevant.

We refer to  $h, v$  as the “proposed (hashed) block/view” of the commitment, and  $h', v'$  as the “justification (hashed) block/view” of the commitment.

We sometimes refer to a commitment as an  $n$ -commitment if it contains  $n$  signatures. In the case of a 1-commitment,



we write  $\langle h, v, h', v', ph, \sigma \rangle$  for  $\langle h, v, h', v', ph, [\sigma] \rangle$ . We define **C-combine**( $[\langle h, v, h', v', ph, \sigma_1 \rangle, \dots, \langle h, v, h', v', ph, \sigma_n \rangle]$ ) to be  $\langle h, v, h', v', ph, [\sigma_1, \dots, \sigma_n] \rangle$ . This is used to create quorum certificates out of partial votes. Given a list  $\vec{\phi}$  of the form  $[\langle h_1, v_1, h'_1, v'_1, ph_1, \sigma_1 \rangle, \dots, \langle h_n, v_n, h'_n, v'_n, ph_n, \sigma_n \rangle]$ , we define the following operation to check whether enough messages of a certain kind have been received: let **C-match**( $\vec{\phi}, k, h, v, ph$ ) be true iff: (1)  $n = k$ ; (2) all  $n$  signatures have been created by different nodes; and (3)  $\forall i \in \{1, \dots, n\}. h = h_i \wedge v = v_i \wedge ph = ph_i$ . Because not all fields are always necessary in commitments, we use  $\perp$  to indicate that a field (more precisely a hash value or view field) is not filled out. The symbol  $\perp$  is considered syntactically different from all hash values and view numbers. Given a commitment  $\phi$  of the form  $\langle h, v, h', v', ph, \vec{\sigma}^n \rangle$ , let  $\phi \cdot \text{Hprep}$  be  $h$ ;  $\phi \cdot \text{Vprep}$  be  $v$ ;  $\phi \cdot \text{Hjust}$  be  $h'$ ;  $\phi \cdot \text{Vjust}$  be  $v'$ ;  $\phi \cdot \text{phase}$  be  $ph$ ; and  $\phi \cdot \text{sign}$  be  $\vec{\sigma}^n$ .

**Accumulator** An accumulator *acc* (generated by the ACCUMULATOR service) is built out of CHECKER-generated commitments and has one of the following forms: (1)  $\langle v, v', h, \vec{i}, \sigma \rangle$ , where  $h$  is the hash value of a block prepared at view  $v'$ , and  $\vec{i}$  is a list of node ids that have prepared blocks at most as high as  $h$ , including a node that has prepared  $h$  at view  $v'$ ; (2)  $\langle v, v', h, n, \sigma \rangle$ , indicating that  $h$  is the hash value of the highest prepared block among  $n$  (a natural number) checked certificates used to generate the accumulator. Let  $\vec{i}_1 @ \vec{i}_2$  be the list  $\vec{i}_2$  appended to the list  $\vec{i}_1$ . Let *acc* · *hash* be  $h$ . Let the size of an accumulator be defined as follows:  $|\langle v, v', h, \vec{i}, \sigma \rangle| = |\vec{i}|$ , and  $|\langle v, v', h, n, \sigma \rangle| = n$ .

**Message** Message exchanged by nodes are either of the form  $\phi$  or  $\langle b, acc, \sigma \rangle$ .

### 6.3 Instances of the Trusted Services

We now describe the instances of our CHECKER and ACCUMULATOR trusted services (presented in §4) used by Damysus.

**CHECKER instance** As mentioned above Damysus's instance of the CHECKER trusted service stores both the view and hash value of the latest prepared block. Also, for convenience, the monotonic counter is split into a view and a phase. More precisely, the CHECKER component maintains a state that comprises the following:

<i>prepv</i>	the view of the last prepared block
<i>preph</i>	the hash value of the last prepared block
<i>view</i>	the current view
<i>phase</i>	the current phase
<i>priv<sub>c</sub></i>	the trusted component's private key
<i>pubs</i>	public keys

The *certificates* produced by the CHECKER service are all 1-commitments, and the ones they act upon are accumulators for **TEEprepare**, and  $(f+1)$ -commitments for **TEEstore**.

**TEEsing** creates a 1-commitment using the prepared block stored in the trusted component. Replicas are required to use this function to generate new-view certificates so that they cannot lie about the last block they have prepared. It

can also be used by late nodes to increase their view and phase number without going through all the steps. As no phase number is checked when calling **TEEsing**, note that crucially, this function generates a commitment  $\phi$  where the proposed hash value  $\phi \cdot \text{Hprep}$  is  $\perp$  so that  $\phi$  can only be used as a commitment for a **nv\_p** phase.

**TEEprepare** takes the pair of the hash value  $h$  of a block proposed by a leader and an accumulator (generated by the leader) corresponding to the justification of that block. If the accumulator was generated by the current leader, i.e., for the current, view then **TEEprepare** generates a 1-commitment, which stands for a partially signed prepare vote.

**TEEstore** takes a  $(f+1)$ -commitment, and checks whether that commitment is for a block that has been prepared in the current view by a quorum of nodes, in which case a 1-commitment is generated, which stands for a partially signed pre-commit vote.

**ACCUMULATOR instance** The ACCUMULATOR component maintains a private key, as well as public keys, which are shared with those of the CHECKER component.

The **TEEstart** function takes a 1-commitment  $\phi$ , and turns it into an accumulator *acc* if  $\phi$ 's signature can be verified, in particular registering in *acc* the id of the node that signed  $\phi$ . This initial commitment is meant to be the one for the prepared block with the highest view number among a collection of  $f+1$  such 1-commitments received by the leader at the beginning of the prepare phase. The leader then iterates over the remaining  $f$  commitments, each time calling **TEEaccum** to check that those commitments are for lower view numbers and signed by different nodes, and updating *acc* accordingly, i.e., recording the ids of the  $f$  signers. Finally, **TEEfinalize** is called to turn *acc* of the form  $\langle v, v', h, [i_1, \dots, i_{f+1}], \sigma \rangle$  into an *acc'* of the form  $\langle v, v', h, f+1, \sigma \rangle$ .

### 6.4 The Damysus Algorithm

Fig. 2 presents Damysus's pseudocode. In the prepare phase, the leader accumulates  $f+1$  new-view messages by calling **accumList** (see Fig. 2a, ln. 7), i.e., the leader selects the new-view message for the block prepared at the highest view among all  $f+1$  received messages, and certifies that it is indeed the highest using its ACCUMULATOR. The leader then extends that highest prepared block with a new block  $b$  (Fig. 2a, ln. 8), and prepares  $b$  using **TEEprepare** (Fig. 2a, ln. 9) of the CHECKER. This guarantees that the algorithm advances to the next step, and generates as well a signature of the newly proposed block (its hash value), which stands for a partially signed prepare vote.

The leader then sends the new proposal  $b$ , along with *acc*, the accumulator it has generated and used to create  $b$ , and the signature of its prepared commitment, to all other nodes. Then, the other nodes check that  $b$  extends the block contained in *acc* (Fig. 2a, ln. 16), and also prepare the block proposed by the leader using their CHECKER (Fig. 2a, ln. 17).

**Figure 2** Damysus ( $2f+1$  replicas and 2 phases)**(a) Non-trusted code at replica  $i$** 

```

1: view = 0 // current view (duplicates the one in the TEE)
2: pubs // public keys (duplicate the ones in the TEE)
3:
4: // prepare phase
5: as a leader
6:   wait for  $\vec{\phi}$  s.t. C-match( $\vec{\phi}$ ,  $f+1$ ,  $\perp$ , view, nv_p)
7:   acc := accumList( $\vec{\phi}$ )
8:   b := createLeaf(acc.hash, client transactions)
9:    $\phi_{prep}$  := TEEprepare(H(b), acc) // also sent to itself
10:  send  $\langle b, acc, \phi_{prep} \cdot \text{sign} \rangle$  to backups
11:
12: as a backup
13:  wait for  $\langle b, \langle \text{view}, v', h', f+1, \sigma \rangle, \sigma' \rangle$  from the leader
14:  acc :=  $\langle \text{view}, v', h', f+1, \sigma \rangle$ 
15:   $\phi_{prep}$  :=  $\langle \text{H}(b), \text{view}, h', v', \text{prep\_p}, \sigma' \rangle$ 
16:  abort if  $\neg(\text{VERIFY}(\phi_{prep})_{pubs} \wedge b > h')$ 
17:  send  $\phi' := \text{TEEprepare}(\text{H}(b), acc)$  to leader
18:
19: // pre-commit phase
20: as a leader
21:  wait for  $\vec{\phi}$  s.t. C-match( $\vec{\phi}$ ,  $f+1$ , h, view, prep_p)
22:  send  $\phi := \text{C-combine}(\vec{\phi})$  to all replicas
23:
24: all replicas
25:  wait for  $\langle h, \text{view}, h', v', \text{prep\_p}, \vec{\sigma}^{f+1} \rangle$  from leader
26:  send  $\phi := \text{TEEstore}(\langle h, \text{view}, h', v', \text{prep\_p}, \vec{\sigma}^{f+1} \rangle)$  to leader
27:
28: // decide phase
29: as a leader
30:  wait for  $\vec{\phi}$  s.t. C-match( $\vec{\phi}$ ,  $f+1$ , h, view, com_p)
31:  send  $\phi := \text{C-combine}(\vec{\phi})$  to all replicas
32:
33: all replicas
34:  wait for  $\langle h, \text{view}, \perp, \perp, \text{pcom\_p}, \vec{\sigma}^{f+1} \rangle$  from leader
35:  abort if  $\neg \text{VERIFY}(\langle h, \text{view}, \perp, \perp, \text{pcom\_p}, \vec{\sigma}^{f+1} \rangle)_{pubs}$ 
36:  execute  $b$  corresponding to  $h$  & reply to clients
37:
38: // new-view phase
39: all replicas
40:  when executed or timeout
41:  // call TEEsign until the node can generate
42:  // a new-view-phase commitment for the next view
43:   $(v, ph) := (\text{view}, \text{prep\_p}); \text{view}++$ 
44:  while  $(v, ph) \neq (\text{view}, \text{nv\_p})$  do
45:     $\phi := \text{TEEsign}(); (v, ph) := (\phi \cdot \text{Vprep}, \phi \cdot \text{phase})$ 
46:  end while
47:  send  $\phi$  to view's leader
48:

```

**(b) TEE code**

```

49: function accumList( $\vec{\phi}$ )
50:    $\phi_0 := \text{commitment } \phi \in \vec{\phi} \text{ with highest } \phi \cdot \text{vjust}$ 
51:   acc := TEEstart( $\phi_0$ )
52:   for  $\phi \in \vec{\phi} \setminus \{\phi_0\}$  do acc := TEEaccum(acc,  $\phi$ )
53:   return TEEfinalize(acc)

(b) TEE code
6: (prepv, preph) = (0, H(G)) // latest prepared block
7: (view, phase) = (0, nv_p) // current step
8: privc, pubs // private (confidential) and public keys
9:
10: function createUniqueSign(h, h', v')
11:    $\phi := \langle h, \text{view}, h', v', \text{phase}, \sigma \rangle$  where  $\sigma := \text{SIGN}(h, \text{view}, h', v', \text{phase})_{privc}$ 
12:   (view, phase)++ // increased to avoid equivocation
13:   return  $\phi$ 
14:
15: function TEEsign()
16:   return  $\phi := \text{createUniqueSign}(\perp, \text{preph}, \text{prepv})$ 
17:
18: function TEEprepare(h, acc) where acc is  $\langle v, v', h', f+1, \sigma \rangle$ 
19:   if  $\text{VERIFY}(acc)_{pubs} \wedge \text{view} = v \wedge h \neq \perp$  then
20:     return  $\phi := \text{createUniqueSign}(h, h', v')$ 
21:   end if
22:
23: function TEEstore( $\phi$ ) where  $\phi$  is  $\langle h, v, h', v', ph, \vec{\sigma}^{f+1} \rangle$ 
24:   if  $\text{VERIFY}(\phi)_{pubs} \wedge \text{view} = v \wedge ph = \text{prep\_p}$  then
25:     preph := h; prepv := v
26:     return  $\phi' := \text{createUniqueSign}(h, \perp, \perp)$ 
27:   end if
28:
29: function TEEstart( $\phi$ ) where  $\phi$  is  $\langle \perp, v, h', v', \text{nv\_p}, \sigma \rangle$ 
30:   if  $\text{VERIFY}(\phi)_{pubs}$  then
31:      $\sigma' := \text{SIGN}(v, v', h', [\sigma \cdot \text{id}])_{privc}$ 
32:     return acc :=  $\langle v, v', h', [\sigma \cdot \text{id}], \sigma' \rangle$ 
33:   end if
34:
35: function TEEaccum(acc,  $\phi$ ) where acc is  $\langle v_1, v'_1, h_1, \vec{i}, \sigma \rangle$ 
   and  $\phi$  is  $\langle \perp, v_2, h_2, v'_2, \text{nv\_p}, \sigma_2 \rangle$ 
36:   if  $\left( \begin{array}{l} v_1 = v_2 \wedge v'_1 \geq v'_2 \wedge \sigma_2 \cdot \text{id} \notin \vec{i} \\ \wedge \text{VERIFY}(acc)_{pubs} \wedge \text{VERIFY}(\phi)_{pubs} \end{array} \right)$  then
37:      $\sigma' := \text{SIGN}(v_1, v'_1, h_1, \vec{i} @ [\sigma_2 \cdot \text{id}])_{privc}$ 
38:     return acc' :=  $\langle v_1, v'_1, h_1, \vec{i} @ [\sigma_2 \cdot \text{id}], \sigma' \rangle$ 
39:   end if
40:
41: function TEEfinalize(acc) where acc is  $\langle v, v', h, \vec{i}, \sigma \rangle$ 
42:   if  $\text{VERIFY}(acc)_{pubs}$  then
43:     return acc' :=  $\langle v, v', h, |\vec{i}|, \sigma \rangle$  where  $\sigma := \text{SIGN}(v, v', h, |\vec{i}|)_{privc}$ 
44:   end if

```

In the next pre-commit phase, correct nodes call **TEEstore** to store  $b$  in their CHECKER component. Nodes are forced to store  $b$  during this phase to generate a partially signed pre-commit vote. This guarantees that if a block is executed,  $f+1$  nodes (possibly Byzantine) have stored it in their CHECKER components, and will relay it in their new-view messages.

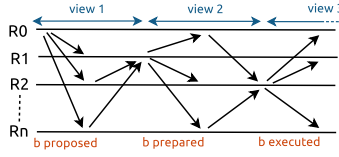
Thanks to the ACCUMULATOR, nodes do not need to lock blocks as done in HotStuff: no commit phase is needed and in the next decide phase, the leader gathers  $f+1$  pre-commit

commitments, and sends a pre-commit ( $f+1$ )-commitment to all nodes, which use it to execute  $b$ . No locking is needed because the ACCUMULATOR service guarantees that leaders only propose blocks that extend the highest prepared block.

## 6.5 Proof of Correctness - Overview

The intuition behind Damysus's correctness stems from the following two points (see [22] for a detailed proof).



**Figure 3** Communication pattern in Chained-Damysus.

(1) The proof relies on demonstrating that a Byzantine leader cannot propose various valid proposals for a single view, or even alter what has been proposed over the course of the different phases within a single view. This is guaranteed by the **CHECKER** that only certifies one proposed value as valid per view (all messages are assigned a unique counter/step using **createUniqueSign**). All other nodes only accept valid messages certified by the **CHECKER**. Hence any vote in any phase of a given view is on the same proposal.

In addition, prepared blocks are also stored in **CHECKER** for both liveness and safety, otherwise leaders might only propose blocks based on  $f+1$  new-view messages from late or Byzantine nodes. **TEESign** in particular enforces that nodes, even Byzantine ones, send their latest prepared blocks.

(2) In basic HotStuff, every node locks a block  $b$  locally once it knows that it has been prepared by a quorum of nodes so that it will later accept another block  $b'$  only if it extends  $b$  or if  $b'$  extends a block that was prepared for a higher view than  $b$ 's (possibly because  $b$  did not get executed, but the system still managed to make some progress on a conflicting block). Without locking, nodes would be allowed to accept any new block, even blocks that are in direct conflict (at the next view, but not extending) with the ones they have executed. Thanks to the **ACCUMULATOR**, locking is unnecessary because when receiving a proposal from the leader with a valid accumulator, nodes are then guaranteed, that this proposal must extend the highest prepared block, and therefore, for example, the block that they have last executed.

## 7 Chained-Damysus

This section presents Chained-Damysus, which concurrently processes blocks to further improve performance.

### 7.1 Algorithm

Chained-Damysus improves on Damysus by allowing nodes to pre-commit the predecessor block  $b_0$  of a block  $b$  when preparing  $b$ , while at the same time deciding the predecessor  $b_1$  of  $b_0$  (Damysus does not have a commit phase). Fig. 3 illustrates Chained-Damysus's communication pattern, while Fig. 5 provides its pseudocode. The leader  $R_0$  of view 1 extends the genesis block with a new block  $b$ , which is proposed to all nodes in this view 1. All nodes vote for  $b$  by sending their votes to the leader  $R_1$  of view 2. In case  $R_1$  receives  $f+1$  votes (prepare messages) for  $b$ , those votes form a new certificate for block  $b$  from view 1. The leader  $R_1$  of view 2 can then use this certificate when creating a new

**Figure 4** Example of a chain of blocks with a blank block.

block  $b'$ . In that case,  $R_1$  does not need to call its trusted accumulator (the condition ln. 8 in Fig. 5a is then false) as it has a certificate from the previous view, and does not need to justify that choice of certificate. However, in case  $R_1$  does not receive  $f+1$  votes for  $b$ , it needs to obtain the latest prepare certificate. To achieve this, nodes always send new-view messages along with their prepare messages (ll. 18, 19, 32, and 33 in Fig. 5a).<sup>4</sup> Once  $R_1$  has collected  $f+1$  new-view messages, it selects the latest prepare certificate, and calls its trusted accumulator to certify that the latest prepare certificate was indeed selected. Once again,  $R_1$  makes use of this certificate when creating a new block  $b'$ . In both cases,  $R_1$  sends its new block to the other nodes. This process is then reiterated in the subsequent views.

In this chained (or pipelined) version, while some other block  $b'$  is being proposed in view 2,  $b$  is being prepared, and while some other block is being proposed in view 3,  $b$  is executed and  $b'$  is being prepared. Chained-Damysus, like Damysus, requires one phase less than chained HotStuff and has a better resilience: it tolerates less than half of the system being corrupt (compared to one-third).

We now detail how we adapt some of the concepts introduced previously to handle a chained version.

**Blocks** As in chained HotStuff [4, §5], we now replace **createLeaf** by **createChain**, which generates blank blocks to fill gaps in the chain. Precisely, the leader of view  $v$  that builds a proposal extending a block proposed in view  $v'$  invokes **createChain** to create blank blocks for all views in  $[v'+1, v-1]$ . We assume that a block stores: (1) as before, a pointer to its parent block (a hash value), accessible using  $b.\text{parent}$ ; (2) its justification accessible using  $b.\text{just}$ , which is a **certificate** as defined below.

Fig. 4 depicts a chain with one blank block between  $b_2$  and  $b_3$ , where the top blue arrows capture  $b.\text{just}$ , while the bottom orange arrows capture  $b.\text{parent}$ .

In chained HotStuff, the oldest block in a chain of four consecutive blocks without blanks gets executed. In the above example,  $b_3$  gets executed in view 7, when  $b_6$  is proposed. As Chained-Damysus requires one phase less, only three consecutive blocks are necessary. Therefore,  $b_3$  gets executed in view 6, when  $b_5$  is proposed.

**Certificates** Certificates (stored within blocks and accessible via  $b.\text{just}$ ) are required to be of size  $f+1$ , and nodes reject blocks such that  $|b.\text{just}| \neq f+1$ . A certificate is now either a *quorum certificate* or an accumulator. For simplicity, we write  $qc$  for both. A *quorum certificate*  $qc$  is of the form

<sup>4</sup>In the pseudo-code, nodes send two messages for simplicity, while in practice those two messages can be combined.

**Figure 5** Chained-Damysus ( $2f+1$  replicas and 3 phases).**(a) Non-trusted code at replica  $i$** 

```

1: view = 1 // current view (duplicates the one in the TEE)
2: pubs // public keys (duplicate the ones in the TEE)
3: qcprep :=  $\perp$  // latest prepared certificate
4: blocks // mapping from views to proposed blocks
5:
6: // prepare phase
7: as a leader
8:   if qcprep.cview  $\neq$  view-1 then
9:     // we don't have the latest certificate
10:    wait for  $\vec{\phi}$  s.t. C-match( $\vec{\phi}$ , f+1,  $\perp$ , view-1, nv_p)
11:    qcprep := accumList( $\vec{\phi}$ )
12:   end if
13:   b := createChain(qcprep, client transactions)
14:   blocks[view] := b
15:   b0 := blocks[b.just.view] if it is defined, else abort
16:   abort if H(blocks[b.just.view])  $\neq$  b.just.hash
17:    $\phi_{prep}$  := TEEprepare(b, b0)
18:   send (b,  $\phi_{prep}$ .sign) to backups
19:   send  $\phi_{nv}$  := TEEsign() to leader of view view+1
20:
21: all replicas
22:   wait for (b,  $\sigma'$ ) from the leader
23:   abort if view  $\neq$  b.just.cview+1
24:   b0 := blocks[b.just.view] if it is defined, else abort
25:   abort if H(blocks[b.just.view])  $\neq$  b.just.hash
26:   b1 := blocks[b0.just.view] if it is defined, else abort
27:   abort if H(blocks[b0.just.view])  $\neq$  b0.just.hash
28:   if  $\neg$  leader of view then
29:      $\phi_{prep}$  := (H(b), view,  $\perp$ ,  $\perp$ , prep_p,  $\sigma'$ )
30:     abort if  $\neg$ (VERIFY( $\phi_{prep}$ )pubs  $\wedge$  b  $\succ^+$  b.just.hash)
31:     blocks[view] := b
32:     send  $\phi'$  := TEEprepare(b, b0) to leader of view view+1
33:     send  $\phi_{nv}$  := TEEsign() to leader of view view+1
34:   end if
35:   if b.parent = H(b0)  $\wedge$  b0.parent = H(b1) then
36:     execute b1 (and previous blocks) & reply to clients
37:   end if
38:   if  $\neg$  leader of view+1 then view++
39:

```

```

40: as a leader of next view view+1
41:   wait for  $\vec{\phi}$  s.t. C-match( $\vec{\phi}$ , f+1, h, view, prep_p)
42:   ( $\langle h, view, \perp, \perp, prep_p, \vec{\sigma}^{f+1} \rangle$ ) := C-combine( $\vec{\phi}$ )
43:   qcprep := (view, h,  $\vec{\sigma}^{f+1}$ ); view++
44:
45: // new-view phase
46: upon timeout// executed by all replicas
47:   ( $v, ph$ ) := (0, prep_p); view++
48:   while ( $v, ph$ )  $\neq$  (view, nv_p) do
49:      $\phi$  := TEEsign(); ( $v, ph$ ) := ( $\phi$ .view,  $\phi$ .phase)
50:   end while
51:   send  $\phi$  to view's leader
52:
53: function accumList( $\vec{\phi}$ )
54:    $\phi_0$  := message  $\phi \in \vec{\phi}$  with highest  $\phi$ .vcomm
55:   acc := TEEstart( $\phi_0$ )
56:   for  $\phi \in \vec{\phi} \setminus \{\phi_0\}$  do acc := TEEaccum(acc,  $\phi$ )
57:   return TEEfinalize(acc)

```

**(b) TEE code**

```

1: (prepv, preph) = (0, H(G)) // latest prepared block
2: (view, phase) = (0, nv_p) // current step
3: privc, pubs // private (confidential) and public keys
4:
5: functions createUniqueSign(h, h', v'), TEEsign()
6:   // Same as in Fig. 2b
7:
8: function TEEprepare(b, b0)
9:   qc := b.just
10:   if ( VERIFY(qc)pubs  $\wedge$  view=qc.cview+1
11:      $\wedge$  qc.hash = H(b0) ) then
12:     if b.parent = H(b0) then
13:       preph := qc.hash; prepv := qc.view
14:     end if
15:     return  $\phi'$  := createUniqueSign(H(b),  $\perp$ ,  $\perp$ )
16:   end if
17: functions TEEstart( $\phi$ ), TEEaccum(acc,  $\phi$ ), TEEfinalize(acc)
18:   // Same as in Fig. 2b

```

$\langle v, h, \vec{\sigma} \rangle$ , where  $h$  is the hash value of a block prepared in view  $v$ , and  $\vec{\sigma}$  is a list of partial signatures on the prepared block at view  $v$ . If a certificate  $qc$  is of the form  $\langle view, h, \vec{\sigma}^n \rangle$ , then let  $qc.cview = qc.view = view$ ,  $qc.hash = h$ , and  $|qc| = n$ . If  $qc$  is an accumulator of the form  $\langle view, v, h, n, \sigma \rangle$ , then let  $qc.cview = view$ ,  $qc.view = v$ ,  $qc.hash = h$ , and  $|qc| = n$ . Given a certificate  $qc$ ,  $qc.view$  refers to the view at which the certificate was created, while  $qc.cview$  refers to the view at which  $qc.hash$  was certified. For example, an accumulator of the form  $\langle v+1, v, h, n, \sigma \rangle$  certifies in view  $v+1$  that  $h$  was the hash value of the highest block elected in view  $v$ .

**Steps** Our chained version has only two types of phases: prepare and view-change identified by the tag  $nv\_p$  and  $prep\_p$ , respectively. Consequently, the increment function is now defined by  $(v, prep\_p)++ = (v, nv\_p)$  and  $(v, nv\_p)++ = (v+1, prep\_p)$ . Steps are incremented in a way that guarantees that a commitment generated at a view  $v$  is tagged with

$v-1$ , regardless of the origin of the commitment, i.e., whether it came from prepare votes or from new-view messages. Due to this step reordering, nodes now start at view 1.

**Commitments** A commitment  $\phi$  is now either a new-view commitment of the form  $\langle \perp, view, h, v, nv\_p, \sigma \rangle$  used by nodes to relay their latest prepared block (its hash value), where  $view$  is the current view, and  $h$  is the hash value of the latest prepared block at view  $v$ ; or it is a prepare commitment of the form  $\langle h, view, \perp, \perp, prep\_p, \vec{\sigma}^n \rangle$  used by nodes to vote on blocks in the prepare phase, where  $h$  is the hash value of the block currently being prepared in view  $view$ . In both cases, we write  $\phi.view$  for  $view$ . In the first case we write  $\phi.hcomm$  for  $h$  and  $\phi.vcomm$  for  $v$ ; and in the second case we write  $\phi.hcomm$  for  $h$  and  $\phi.vcomm$  for  $view$ . Note that  $\phi.hcomm$  is the hash value of a block that received  $f+1$  votes in view  $\phi.vcomm$ , while  $\phi.view$  is the view in which this information is signed and sent.

## 7.2 Proof of Correctness - Overview

We prove Chained-Damysus's safety and liveness in [22], and provide here a high-level view of the safety proof, i.e., that correct nodes do not execute conflicting blocks.

A block is executed if it is the oldest in a chain of three consecutive blocks, i.e., if it is certified in three consecutive views. For safety, we must prove that the executed oldest block  $b_2$  of any such chain  $C_2$  extends the executed oldest block  $b_1$  of any other such chain  $C_1$  (assuming  $b_2$ 's view is higher than  $b_1$ 's). If the two chains overlap, then they cannot be inconsistent because at most one block can be certified per view, as nodes can only vote once per view. If the two chains do not overlap, then  $b_2$ 's view is higher than the views of  $C_1$ 's blocks. We can then derive that  $b_2$  is a descendant of the middle block  $b'_1$  of  $C_1$ , and therefore that  $b_2 \succ^* b'_1 \succ b_1$ . We differentiate two cases: (1) the leaders of the views between  $C_1$  and  $C_2$  have created certificates out of the prepare messages they received, and therefore did not need to call their trusted accumulator, in which case  $b_2 \succ^* b'_1$ ; and (2) a leader  $L$  of a view between  $C_1$  and  $C_2$  extended a different certificate than the one created at the end of the previous view. In that case, because leaders (even Byzantine ones) have to justify that they have selected the latest prepare certificate using their trusted accumulators, and because nodes (even Byzantine ones) must relay their latest prepare certificates in new-view messages, we are again guaranteed that  $L$  will only be able to accumulate  $f+1$  votes for its proposal if that proposal extends  $b'_1$ .

## 8 Evaluation

**Implemented protocols** This section evaluates the performance of Damysus (§6) and Chained-Damysus (§7) and respectively compares them to those of basic and chained HotStuff. In order to better quantify the impact of the CHECKER and ACCUMULATOR trusted services on performance we also implement two non-chained additional versions of Damysus. One where nodes are only equipped with a CHECKER, but not an ACCUMULATOR, called Damysus-C, and one where nodes are only equipped with an ACCUMULATOR, but not a CHECKER, called Damysus-A. We recall from §4.3 that by using a single component, Damysus-C merely improves the resilience while Damysus-A only reduces the communication phases. We summarize in what follows the protocols we evaluate (C. stands for CHECKER and A. for ACCUMULATOR).

name	nodes	phases	trusted comp.
basic HotStuff	$3f + 1$	3	none
Damysus-C	$2f + 1$	3	C.
Damysus-A	$3f + 1$	2	A.
Damysus	$2f + 1$	2	C. & A.
chained HotStuff	$3f + 1$	3	none
Chained-Damysus	$2f + 1$	2	C. & A.

**Development environment** All protocols are implemented in C++<sup>5</sup>, and whenever applicable use Intel SGX enclaves [7]

to run trusted services. We use SGX because of its Linux SDK that provides a convenient development environment. Although our trusted services are generic enough to be potentially implemented in any trusted execution environment, we acknowledge that SGX security vulnerabilities have been described in the literature [24]. Replicas use ECDSA signatures with prime256v1 elliptic curves (available in OpenSSL [25]), and are interconnected using Salticidae [26].

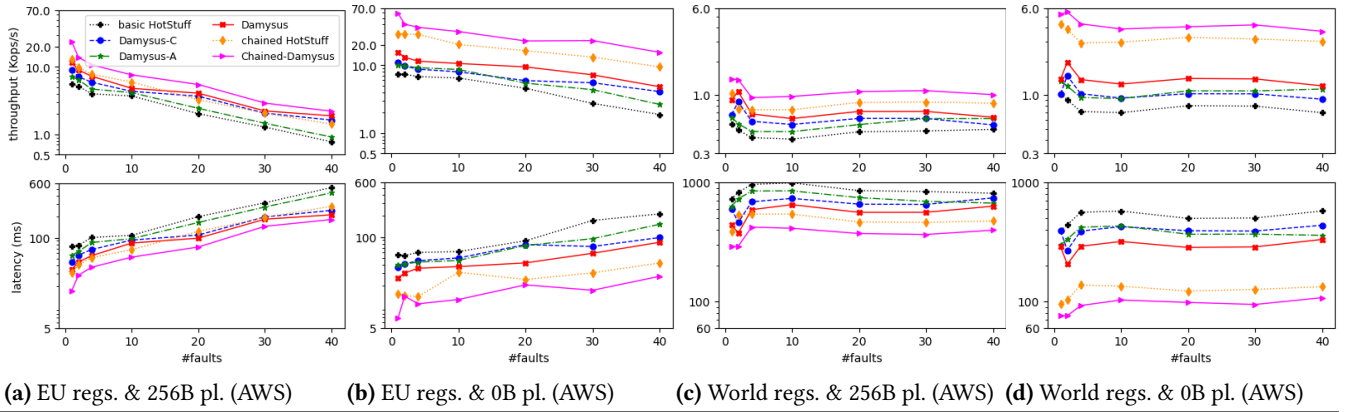
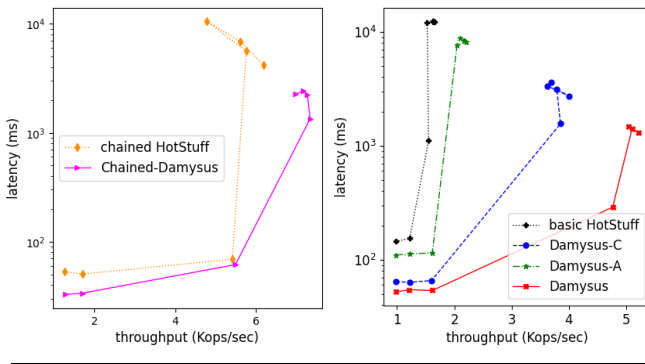
Note that basic and chained HotStuff do not rely on trusted components, while all the other protocols do and make use of SGX to implement them.

**Deployment settings** We deploy our protocols on AWS EC2 machines with one t2.micro instance per node instance to reproduce a real world environment. Each figure presents the average of 100 repetitions, where each repetition involves 30 views. We vary the fault threshold  $f \in \{1, 2, 4, 10, 20, 30, 40\}$  (i.e. up to 121 nodes) with blocks of 400 transactions, and use payloads of 0B and 256B. In addition to the payload, a transaction contains  $2 \times 4B$  for metadata: a client id, and a transaction id. Therefore, the experiments with payloads of 0B for example involve blocks of size  $400 \times 2 \times 4B$ .

**Regional deployment** Fig. 6a shows an experiment with nodes deployed across 4 regions in Europe (Ireland, London, Paris, and Frankfurt), and payloads of size 256B. As we can see there, compared to basic HotStuff, Damysus-C has an average throughput increase of 59.7%, Damysus-A of 19.3%, and Damysus of 87.5%. Moreover, Damysus-C has an average latency decrease of 35.9%, Damysus-A of 16.6%, and Damysus of 45%. In addition, compared to chained HotStuff, Chained-Damysus has a throughput increase of 50.5% and latency decrease of 32.1%. Fig. 6b uses a similar setting, but now payloads are of size 0B. As we can see there, compared to basic HotStuff, Damysus-C has an average throughput increase of 54.6%, Damysus-A of 36.7%, and Damysus of 107.1%. Moreover, Damysus-C has an average latency decrease of 31.8%, Damysus-A of 27.4%, and Damysus of 50.6%. In addition, compared to chained HotStuff, Chained-Damysus has a throughput increase of 57.4% and latency decrease of 33.1%. We can observe from those experiments that each of the trusted components has a substantial effect on increasing the throughput and decreasing the latency.

**Cross continent deployment** Fig. 6c shows an experiment with nodes deployed across 11 regions across the world (4 in the US in North Virginia, Ohio, North California, and Oregon; 4 in Europe in Ireland, London, Paris, and Frankfurt; 2 in Asia in Singapore, Sydney; and 1 in Canada Central), and payloads of size 256B. We observe that compared to basic HotStuff Damysus-C has an average throughput increase of 35.1%, Damysus-A of 18.4%, and Damysus of 61.6%. Note that Damysus performs better than a simple combination of CHECKER (as in Damysus-C) and ACCUMULATOR (as in Damysus-A), as it eliminates the need of CHECKER storing locked blocks, as explained in §4. Moreover, Damysus-C has an average latency decrease of 24.2%, Damysus-A of 14.0%,

<sup>5</sup>Link to repository removed for double blind review.

**Figure 6** Throughput (top Figures, in Kops/s with log scale) and latency (bottom Figures, ms) varying the regions and payloads.**Figure 7** Max. throughput (Kops/s) vs. latency (ms).

and Damysus of 36.6%. In addition, compared to chained HotStuff, Chained-Damysus has a throughput increase of 35.2% and latency decrease of 24.8%. Fig. 6d uses a similar setting, but now payloads are of size 0B. As we can see there, compared to basic HotStuff, Damysus-C has an average throughput increase of 33.1%, Damysus-A of 38.2%, and Damysus of 78.6%. Moreover, Damysus-C has an average latency decrease of 23.3%, Damysus-A of 27.0%, and Damysus of 43.0%. In addition, compared to chained HotStuff, Chained-Damysus has a throughput increase of 32.2% and latency decrease of 23.7%.

In all cases, the hybrid protocols have a higher throughput and lower latency than the corresponding non-hybrid protocols. The most efficient protocol is always Damysus, which combines the advantages of Damysus-C and Damysus-A. The same conclusions apply to the chained versions.

**Comparison with a given  $N$**  In addition to comparing the protocols for a given number of faults  $f$ , we also compare them depending on the number of nodes  $N$ . As one can observe in the above experiments, for a system of a given size, such as  $N = 61$ , Damysus and Chained-Damysus tolerate more faults than basic and chained HotStuff ( $\lfloor \frac{N-1}{2} \rfloor = 30$  compared to  $\lfloor \frac{N-1}{3} \rfloor = 20$ ), and also perform better. For example, for  $N = 61$ , Chained-Damysus and chained HotStuff have comparable throughput and latency in Fig. 6a, while

Chained-Damysus has a 23Kops/sec throughput and a 17ms latency compared to chained HotStuff's 16Kops/sec throughput and 24ms latency in Fig. 6b.

**Throughput versus latency** Fig. 7 evaluates all protocols while increasing the throughput until system saturation. For this experiment, we configure the protocols using  $f = 1$ , blocks of 400 transactions, 0B payloads, and we use 4 regions in Europe. We increase the rate at which clients submit transactions, until we saturate the systems. In the left figure, we evaluate the chained protocols using 10 clients, which send a total of 250000 transactions each; while in the right figure we evaluate the basic protocols using 6 clients, which send a total of 50000 transactions each. We observe that Chained-Damysus has a lower latency than HotStuff, and a higher maximum throughput. Similarly, we observe that all basic hybrid versions have a lower latency and higher maximum throughput than HotStuff. Damysus performs better than Damysus-C, which performs better than Damysus-A.

## 9 Conclusion

This paper explored the design space of hybrid solutions that leverage arguably minimalist trusted components to improve the performance and resilience of HotStuff-like protocols. We showcased that the simplest known trusted components in hybrid solutions for traditional BFT protocols are not sufficient to deliver the same guarantees for streamlined BFT protocols. To this end, we introduced two simple trusted components whose services are locally required by every node in the system. Using these trusted components, which we called CHECKER and ACCUMULATOR, we introduced Damysus a hybrid variant of basic HotStuff. Damysus simultaneously improved resilience, reduced communication complexity and latency, and over-performed HotStuff in terms of throughput. Damysus preserves HotStuff's overall behavior to retain linear view change and optimistic responsiveness. We finally extended Damysus to Chained-Damysus, which supports chained operations. our results to the chained version of HotStuff presenting Chained-Damysus.



## References

- [1] E. Androulaki et al. 2018. Hyperledger fabric: a distributed operating system for permissioned blockchains. In *EUROSYS*. ACM.
- [2] C. Cachin and M. Vukolic. 2017. Blockchain consensus protocols in the wild (keynote talk). In *DISC (LIPIcs)*. Volume 91. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.
- [3] M. Castro and B. Liskov. 1999. Practical byzantine fault tolerance. In *OSDI*. USENIX Association.
- [4] M. Yin, D. Malkhi, M. K. Reiter, G. Golan-Gueta, and I. Abraham. 2019. Hotstuff: BFT consensus with linearity and responsiveness. In *PODC*. ACM.
- [5] B. Y. Chan and E. Shi. 2020. Streamlet: textbook streamlined blockchains. In *AFT*. ACM.
- [6] M. Correia, N. F. Neves, and P. Verissimo. 2004. How to tolerate half less one byzantine nodes in practical distributed systems. In *SRDS*. IEEE Computer Society.
- [7] [n. d.] SGX. <https://software.intel.com/en-us/sgx>.
- [8] G. S. Veronese, M. Correia, A. N. Bessani, L. C. Lung, and P. Verissimo. 2013. Efficient byzantine fault-tolerance. *IEEE Trans. Computers*, 62, 1, 16–30.
- [9] R. Kapitza, J. Behl, C. Cachin, T. Distler, S. Kuhnle, S. V. Mohammadi, W. Schröder-Preikschat, and K. Stengel. 2012. Cheapbft: resource-efficient byzantine fault tolerance. In *EuroSys*. ACM.
- [10] J. Behl, T. Distler, and R. Kapitza. 2017. Hybrids on steroids: sgx-based high performance BFT. In *EuroSys*. ACM.
- [11] J. Liu, W. Li, G. O. Karame, and N. Asokan. 2019. Scalable byzantine consensus via hardware-assisted secret sharing. *IEEE Trans. Computers*, 68, 1, 139–151.
- [12] S. Yandamuri, I. Abraham, K. Nayak, and M. K. Reiter. 2021. Communication-efficient bft protocols using small trusted hardware to tolerate minority corruption. *Cryptology ePrint Archive*, Report 2021/184. (2021).
- [13] B.-G. Chun, P. Maniatis, S. Shenker, and J. Kubiatowicz. 2007. Attested append-only memory: making adversaries stick to their word. In *SOSP*. ACM.
- [14] D. Levin, J. R. Douceur, J. R. Lorch, and T. Moscibroda. 2009. TrInc: small trusted hardware for large distributed systems. In *NSDI*. USENIX Association.
- [15] J. Zhang, J. Gao, K. Wang, Z. Wu, Y. Lan, Z. Guan, and Z. Chen. Tbf: understandable and efficient byzantine fault tolerance using trusted execution environment. (2021).
- [16] B. Arun. 2019. *Revamping Byzantine Fault Tolerant State Machine Replication with Decentralization, Trusted Execution, and Practical Transformations*. PhD thesis. Blacksburg, Virginia.
- [17] T. Distler, C. Cachin, and R. Kapitza. 2016. Resource-efficient byzantine fault tolerance. *IEEE Trans. Computers*, 65, 9, 2807–2819.
- [18] S. Bano, M. Baudet, A. Ching, A. Chursin, G. Danezis, F. Garillot, Z. Li, D. Malkhi, O. Naor, D. Perelman, and A. Sonnino. State machine replication in the libra blockchain. (2019).
- [19] M. M. Jalalzai, J. Niu, and C. Feng. 2020. Fast-hotstuff: A fast and resilient hotstuff protocol. *CoRR*, abs/2010.11454. arXiv: 2010.11454.
- [20] A. N. Bessani, J. Sousa, and E. A. P. Alchieri. 2014. State machine replication for the masses with BFT-SMART. In *DSN*. IEEE.
- [21] M. Correia, N. F. Neves, and P. Verissimo. 2004. How to tolerate half less one byzantine nodes in practical distributed systems. In *SRDS*. IEEE Computer Society.
- [22] Anonymous. [n. d.] Technical Report: Streamlined BFT Consensus Using Trusted Components. [https://www.dropbox.com/sh/gasx9ndczq8sfj1/AAADuxKem\\_y12VgMYBe38vE1va?dl=0](https://www.dropbox.com/sh/gasx9ndczq8sfj1/AAADuxKem_y12VgMYBe38vE1va?dl=0).
- [23] C. Dwork, N. A. Lynch, and L. J. Stockmeyer. 1988. Consensus in the presence of partial synchrony. *J. ACM*, 35, 2, 288–323.
- [24] K. Murdock, D. F. Oswald, F. D. Garcia, J. V. Bulck, D. Gruss, and F. Piessens. 2020. Plundervolt: software-based fault injection attacks against intel SGX. In *SP*. IEEE.
- [25] Last accessed: Feb 25, 2020. Openssl. <https://www.openssl.org/>.
- [26] [n. d.] Salticidae. <https://github.com/Determinant/salticidae>.