

# MYSTICETI: Low-Latency DAG Consensus with Fast Commit Path

Kushal Babel  
Cornell Tech, USA  
IC3

Andrey Chursin  
Mysten Labs, USA

George Danezis  
Mysten Labs, UK  
University College London

Lefteris Kokoris-Kogias  
Mysten Labs, Greece  
IST Austria

Alberto Sonnino  
Mysten Labs, UK  
University College London

## Abstract

We introduce MYSTICETI-C a byzantine consensus protocol with low-latency and high resource efficiency. It leverages a DAG based on Threshold Clocks and incorporates innovations in pipelining and multiple leaders to reduce latency in the steady state and under crash failures. MYSTICETI-FPC incorporates a fast commit path that has even lower latency. We prove the safety and liveness of the protocols in a byzantine context. We evaluate MYSTICETI and compare it with state-of-the-art consensus and fast path protocols to demonstrate its low latency and resource efficiency, as well as more graceful degradation under crash failures. MYSTICETI is the first byzantine protocol to achieve WAN latency of 0.5s for consensus commit, at a throughput of over 50k TPS that matches the state-of-the-art.

## 1 Introduction

Several recent blockchains, such as Sui [7, 38], have adopted consensus protocols based on certified directed acyclic graphs (DAG) of blocks, such as Narwhal-Tusk [17], Bullshark [30], and recent proposals such as Shoal [31]. By design, these consensus protocols scale well in terms of throughput with a performance of 100k TPS of raw transactions. Notably, using Bullshark coupled with the MoveVM, Sui has processed a peak of 65m programmable transaction blocks on 27 July 2023, sustained throughput of over 700 TPS over the day, using a single consensus worker [7].

However, certified DAG consensus suffers from two disadvantages: (1) the certified DAG requires multiple round-trips: to broadcast each block between validators, get signatures, and re-broadcast certificates. This leads to higher latency than traditional consensus protocols [20, 35]; (2) since all blocks need to be signed by all, signature generation and verification consume a large amount of CPU on each validator which grows with the number of validators. This burden is particularly heavy when a crash-recovered validator syncs to the DAG and is overwhelmed by signature verification.

A separate line of work explores the processing of transactions without or before reaching consensus, such as in FastPay [5], Zef [6], and Astro [15]. These systems use reliable broadcast instead of consensus to commit transactions

that only access state controlled by a single party, and refer to this mechanism as a *fast path*. The Sui Lutris [7] mechanism underlying the Sui blockchain combines a fast path with a black-box certified DAG consensus. This composition is generic and leads to very low latencies for fast path transactions (the vast majority of transactions on the peak 27 July 2023 date). But it also leads to (1) increased latencies for other transactions requiring the consensus path and overall increased sync latency due to a separate post-consensus checkpoint mechanism, and (2) additional signature generation and verification since each transaction needs to be certified separately. The latter means that the validator’s CPU is largely devoted to performing cryptographic operations rather than processing transactions. The need to sequence valid certificates into the consensus also reduces its capacity and imposes the CPU-heavy certificate verification burden on the critical path of consensus.

In this work, we present MYSTICETI a family of protocols to safely commit distributed transactions in a Byzantine setting that focuses on low-latency and low-CPU operation. MYSTICETI-C is a consensus protocol based on a threshold logical clock [18] DAG of blocks, that commits without the need to explicitly certify each block. We extend it to include a fast path with MYSTICETI-FPC, leading to very low-latency commits, without the need to generate an explicit certificate for each transaction. Both designs yield lower latency as well as lower CPU utilization. Both protocols could be augmented with multiple workers to support higher throughput, but our experiments show that current peaks of activity for all blockchains can comfortably be served by a single worker<sup>1</sup>.

**Contributions.** We make the following contributions:

- We present MYSTICETI-C, a DAG-based byzantine consensus algorithm and its proofs of safety and liveness. Notably, it implements a generic pipelined and multi-leader commit rule where every single block can be directly committed significantly reducing latency even when failures occur. We show it has a low commit latency, and exceeds the throughput of one worker Narwhal-based consensus.

<sup>1</sup><https://app.artemis.xyz/comparables>

- We also present MYSTICETI-FPC that offers feature parity with Sui Lutris [7], i.e. both a fast path and a consensus path, as well as safe checkpointing and epoch close mechanisms. We show that MYSTICETI-FPC has a fast path latency comparable with Zef [6] and Fastpay [5] but higher throughput due to lower CPU utilization and batching.
- We fully implement both protocols and perform an experimental evaluation on a Wide Area Network. We show their latency-throughput characteristics are superior to certified DAG-based designs on the consensus mode; and competitive in the fast path while their throughput is far superior due to lower CPU overheads.

## 2 Background

We consider a message-passing system in which each epoch  $n = 3f + 1$  validators process transactions using the MYSTICETI protocols. In every epoch, a computationally bound adversary that controls the network can statically corrupt an unknown set of up to  $f$  validators. We call these validators *byzantine* and they can deviate from the protocol arbitrarily. The rest of the validators (at least  $2f + 1$ ) are *correct* or *honest* and follow the protocol faithfully.

For the description of the protocol, we assume that links between honest parties are reliable and authenticated. That is, all messages among honest parties eventually arrive and a receiver can verify the sender’s identity. The adversary is computationally bound and the usual security properties of cryptographic hash functions, digital signatures, and other cryptographic primitives hold. Under these assumptions, Section 5 shows that the MYSTICETI protocols are safe, in that, no two correct validators commit inconsistent transactions.

Validators are communicating over a partially synchronous network. There exists a time called Global Stabilization Time (GST) and a known finite time bound  $\Delta$ , such that any message sent by a party at time  $x$  is guaranteed to arrive by time  $\Delta + \max\{\text{GST}, x\}$ . Within periods of synchrony (after GST) the MYSTICETI protocols are also live in that they are guaranteed to commit transactions from correct validators.

Following prior work [17, 23, 30] we focus on byzantine atomic broadcast for MYSTICETI. Additionally for MYSTICETI-FPC, we show that the fast-path transactions subprotocol satisfies reliable broadcast within an epoch [7], but allows for recovery of equivocating objects across epochs without losing safety at the epoch boundaries. More formally:

**Reliable broadcast.** Each validator  $v_k$  broadcasts messages by calling  $r\_broadcast_k(m, q)$ , where  $m$  is a message and  $q \in \mathbb{N}$  is a sequence number. Every validator  $v_i$  has an output  $r\_deliver_i(m, q, v_k)$ , where  $m$  is a message,  $q$  is a sequence number, and  $v_k$  is the identity of the validator that called the corresponding  $r\_broadcast_k(m, q)$ . The reliable broadcast abstraction guarantees the following properties:

**Agreement:** If an honest validator  $v_i$  outputs  $r\_deliver_i(m, q, v_k)$ , then every other honest validator  $v_j$  eventually outputs  $r\_deliver_j(m, q, v_k)$ .

**Integrity:** For each sequence number  $q \in \mathbb{N}$  and validator  $v_k$ , an honest validator  $v_i$  outputs  $r\_deliver_i(m, q, v_k)$  at most once regardless of  $m$ .

**Validity:** If an honest validator  $v_k$  calls  $r\_broadcast_k(m, q)$ , then every honest validator  $v_i$  eventually outputs  $r\_deliver_i(m, q, v_k)$ .

Additionally, for byzantine atomic broadcast, each honest validator  $v_i$  can call  $a\_broadcast_i(m, q)$  and output  $a\_deliver_i(m, q, v_k)$ . A byzantine atomic broadcast protocol satisfies reliable broadcast (agreement, integrity, and validity) as well as:

**Total order:** If an honest validator  $v_i$  outputs  $a\_deliver_i(m, q, v_k)$  before  $a\_deliver_i(m', q', v'_k)$ , then no honest party  $v_j$  outputs  $a\_deliver_j(m', q', v'_k)$  before  $a\_deliver_j(m, q, v_k)$ .

Finally, most prior work defines properties as if the protocol runs in a single epoch. This, however, is unrealistic as there is validator churn. To this end, we extend all the protocols to also take as a parameter the epoch number and all properties should hold inside a single epoch. Fortunately, the definition of reliable broadcast allows the recovery of liveness for blocked sequence numbers that are equivocated inside an epoch. More specifically we define equivocation tolerance as follows:

**Equivocation tolerance** If a byzantine validator  $v_k$  concurrently called  $r\_broadcast_k(m, q, e)$  and  $r\_broadcast_k(m', q, e)$  with  $m \neq m'$  then the rest of the validators either  $r\_deliver_i(m, q, v_k, e)$ , or  $r\_deliver_i(m', q, v_k, e)$ , or there is a subsequent epoch  $e' > e$  where  $v_k$  is honest, calls  $r\_broadcast_k(m'', q, e')$  and all honest validators  $r\_deliver_i(m'', q, v_k, e')$ ,

## 3 The MYSTICETI-C Protocol

We describe the MYSTICETI-C consensus protocol. Section 4 describes MYSTICETI-FPC, a variant incorporating a fast path.

### 3.1 MYSTICETI-C overview

MYSTICETI-C allows a committee of validators to open a consensus channel for an *epoch*, sequence several messages within it, and then eventually close the channel at the end of the epoch. The MYSTICETI-C protocol proceeds in a sequence of *rounds*. At the end of every round, each honest validator broadcasts a *unique signed block* for the round. During a round validators receive transactions from users, as well as blocks from other validators. They construct their block to contain both *references to blocks* from past rounds, always starting from their own latest block; as well as *fresh transactions* not already included indirectly in the past blocks. Once a block contains references to at least  $2f + 1$  validator blocks from the previous round, and after a delay, it can be signed by the validator and disseminated to the other validators.

To commit transactions, the basic variant of MYSTICETI-C relies on committing a *common sequence of blocks from leaders* at specific *leader rounds*. Rounds are structured to be a

sequence of a leader round, followed by one or more support rounds, and finally a decision round. For each leader round all correct validators determine a leader using a deterministic method based on the round’s number. Within the subsequent support rounds, blocks support the first leader block (in case of equivocation) indirectly included in their block. When a block indirectly includes  $2f + 1$  validator blocks that support a leader block we say the block certifies the leader block. If a leader block is certified by  $2f + 1$  blocks in the decision round, we initiate extending the leader commit sequence. First, any previous uncommitted leader block that is certified by at least one block in the causal history of the leader is committed, before the final leader block is committed.

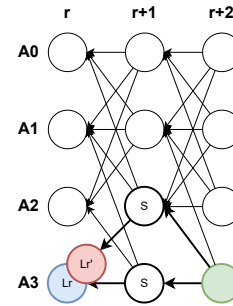
The sequence of leaders committed is consistent across all correct validators. Each leader block commits the full causal history of blocks and contained transactions, that are not already part of a previous commit. The algorithm to transform the sequence of leader commits and expand it to transaction commits can be arbitrary as long as all new transactions in blocks are included in the sequence in a deterministic manner. The basic variant is extended to what we call the universal commit rule (Section 3.4). The universal commit rule runs concurrently multiple virtual MYSTICETI-C effectively multiplexing commit sequences stemming from every block virtually acting as a potential leader. This allows for a reduction in commit latency in the common case as well as under failures.

### 3.2 Transactions submission, support, and certificates on consensus blocks

A client submits transactions to a validator who includes it inside their next block. If the transaction does not appear in the consensus output within some time, the client picks another validator and retries. Since MYSTICETI-C implements a BAB the transactions at this stage are treated as a payload of bytes which will be forwarded to the execution engine [22, 28] after the transaction first appears in the total ordering.

The main unit of communication in MYSTICETI-C is the block, which includes (1) the author  $A$  of the block along with their signature on the full block contents, (2) a round number  $r$ , (3) a sequence containing at least  $2f + 1$  distinct hashes of blocks from previous rounds, and (4) a list of transactions. By convention, the first hash must be to the latest past block from  $A$ . A block is valid if it is signed by a valid validator; all hashes point to distinct valid blocks from previous rounds; the first block links to a block from  $A$ ; and within the sequence of past blocks, there are  $2f + 1$  blocks from the previous round  $r - 1$ . We index each block by the triplet  $B \equiv (A, r, h)$ , comprised of the author  $A$ , the round  $r$ , and the hash  $h$  of the block contents. A correct validator produces at most one unique block per round.

A block  $B'$  supports a past block  $B \equiv (A, r, h)$ , if in the depth-first search performed starting at  $B'$  and recursively following all blocks in the sequence of blocks hashed, block



**Figure 1.** Block  $(A_3, r + 2, \_)$  (green) contains support originating from different validators for both  $(A_3, r, L_r)$  (blue) and  $(A_3, r, L'_r)$  (red) equivocating blocks. If any of the blocks will gather  $2f + 1$  support it will be certified, and we show that at most one may do so.

$B$  is the first block encountered for validator  $A$  at round  $r$ . Once a block  $B$  is supported by  $2f + 1$  distinct validator blocks, we say that the block  $B$  is *certified*, and any block that links to over  $2f + 1$  blocks that support  $B$  represents a *certificate* for the block  $B$ . As we will see in the next section this concept of support and certificates is critical to achieve safety in MYSTICETI-C.

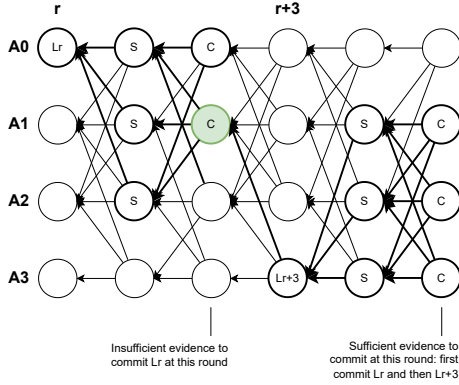
We note that support and certification are properties of blocks and are monotonic properties: if one correct validator ever observes a block supporting another, or certifying another this will remain true forever, and also will be considered true by any other honest validator once they process the block. As Figure 1 illustrates, a block  $(A_3, r + 2, \_)$  (green) that supports block  $(A_3, r, L_r)$  (blue) may contain also support for an equivocating block  $(A_3, r, L'_r)$  (red).

### 3.3 Leaders and basic consensus and commit rule

MYSTICETI-C rounds are composed of an ever-repeating sequence of phases which we call a wave: a single leader round, followed by a fixed number of one or more support rounds, and finally a single commit round (there must be at least three phases in one wave). This pattern of rounds continues until the consensus channel is closed on a final commit round. Without losing generality we present the protocol with a single support round, and 3 rounds in total. This version requires two timeouts, one for the leader round and one for the support round<sup>2</sup>.

At each leader round a single validator is deterministically chosen to act as the *leader*: the block from this validator in this round is used – if committed – to extend the sequence of commits, and as a result, the sequence of transactions. The exact algorithm by which leaders are chosen does not matter as long as all correct validators agree on the leader for all

<sup>2</sup>An alternative design would be with two support rounds (4 rounds in total) which requires a timeout only for the leader round, but in our initial experiments the 3-round design was faster.



**Figure 2.** An illustration of the commit rule.

leader rounds, and eventually, a correct validator acts as a leader (Section 5).

During the support round, validators generate blocks with sufficient delay to allow a block from the leader to be included in their subDAGs and gain support. Finally, in the commit round, validators wait enough time to witness sufficient support. If  $2f + 1$  blocks certify a leader block from the leader round, we start the algorithm to extend the sequence of leader blocks committed. To extend the sequence of committed leader blocks we apply the following algorithm: the latest leader block  $L$  will be committed last (since this triggered the commit). Then we consider the previous leaders: we pick the previous uncommitted leader  $L'$  with the largest round *that has a certificate in the sub-graph* from block  $L$ . And apply the commit algorithm recursively with  $L \leftarrow L'$ . The algorithm stops when the last committed leader is discovered. Then leaders are committed from the oldest to the newest by unwinding the recursion. This guarantees that even if validators commit leaders at different points in (logical threshold) time they do not diverge. Algorithm 3 describes the base commit algorithm.

Given the sequence of leader blocks committed, a deterministic algorithm can be applied to extend the sequence with all the blocks linked by the leaders, and then expand the blocks to the transactions contained in blocks. For example, each leader block can commit the full dag of blocks it links to that are not yet committed, and all transactions not ordered previously their transactions, ordered by ascending rounds.

**Illustration of commit rule.** Figure 2 illustrates an example of the commit rule in action. Leader  $L_r$  does not have enough certificates in the decision round  $r + 2$  to commit since only  $2 (< 2f + 1)$  validators issue blocks that certify it. However, leader block  $L_{r+3}$  has enough (i.e.,  $3 = 2f + 1$  for  $f = 1$ ) certificates in the decision round  $r + 5$ , to initiate a commit. Before committing block  $L_{r+3}$ , the leader block  $L_r$  is committed, since  $L_{r+3}$  contains a certificate for  $L_r$  in its

---

### Algorithm 1 Offset Logic

---

```

1: procedure WAVELENGTH( $r$ )
2:   return  $(r - \text{roundOffset}) / \text{waveLength}$ 

3: procedure LEADERROUND( $w$ )
4:   return  $w * \text{waveLength} + \text{roundOffset}$ 

5: procedure DECISIONROUND( $w$ )
6:   return  $w * \text{waveLength} + \text{waveLength} - 1 + \text{roundOffset}$ 

7: procedure GETPREDEFINEDLEADER( $w$ )
8:    $r_{\text{leader}} \leftarrow \text{LEADERROUND}(w) + \text{LeaderOffset}$ 
9:   return PREDEFINEDLEADER( $r_{\text{leader}}$ )

```

---

causal history (green block). So the sequence is extended by  $L_r$  then  $L_{r+3}$ .

### 3.4 The universal commit rule

The basic MYSTICETI-C is a good primer to understand why the protocol works. However, it suffers from two shortcomings that increase latency. First, it only commits every 3 rounds resulting in 5 rounds to commit for some transactions. Additionally, a crashed leader increases latency by at least 3 more rounds. We resolve these challenges through the universal committer abstraction. Unlike the basic MYSTICETI-C that tries to commit one block per wave, the universal MYSTICETI-C tries to commit every block of the wave. To achieve this we need to make sure that there are no gaps in the commit sequence. For this reason, we introduce the notion of slots. A slot is a tuple (validator, round) and is either empty or has the proposal of the validator for the round. Additionally, a slot can have three states ‘to-commit’, ‘to-skip’, and ‘undecided’. The origins of this approach can be traced back to Multi-Paxos [25]. We can decide how many slots per round we want to instantiate, but it is a configuration parameter fixed for the epoch. The basic committer has one slot every 3 rounds, the leader slot. The benefit of more than one slot per round is to mask crash-faults but if the validator in the slot is Byzantine they can manipulate their slot to be undecided and have the same effect in latency as if it was an unmasked crash fault. Hence, we decided to have 2 slots per round as a good tradeoff for our experiments.

The new commit rule is as follows. First, we get a total ordering between all pending slots, this is deterministic and follows the ordering of rounds. Within one round the ordering can be fixed or changed per round (e.g. round robin). For example, if we have 4 validators ( $v_1, v_2, v_3, v_4$ ) and two rounds ( $r_1, r_2$ ) a potential total ordering of slots with two slots per round is:  $[(v_1, r_1), (v_2, r_1), (v_2, r_2), (v_3, r_2)]$ . Where this ordering represents a FIFO queue ( $(v_1, r_1)$  at the head). Once we have this slot ordering we run the classic MYSTICETI-C protocol described above, but once we enter a new round we virtually treat every block as if it is the leader of a wave. So for a slot  $S$  at round  $r$ , the moment we enter its decision round ( $r + 2$ ) we run the leader-commit checks. More specifically, if at the



---

**Algorithm 2** Consensus helper functions

---

```

1: procedure GETLEADERBLOCK( $w$ )
2:    $r_{leader} \leftarrow \text{LEADERROUND}(w)$ 
3:    $id \leftarrow \text{GETPREDEFINEDLEADER}(r_{leader})$ 
4:   if  $\exists b \in \text{DAG}[r_{leader}]$  s.t.  $b.author = id$  then return  $b$ 
5:   return  $\perp$ 

6: procedure GETFIRSTVOTINGBLOCKS( $w$ )
7:    $r_{voting} \leftarrow \text{LEADERROUND}(w) + 1$ 
8:   return  $\text{DAG}[r_{voting}]$ 

9: procedure GETDECISIONBLOCKS( $w$ )
10:   $r_{decision} \leftarrow \text{DECISIONROUND}(w)$ 
11:  return  $\text{DAG}[r_{decision}]$ 

12: procedure LINK( $b_{old}, b_{new}$ )
13:  return exists a sequence of  $k \in \mathbb{N}$  blocks  $b_1, \dots, b_k$  s.t.  $b_1 = b_{old}, b_k = b_{new}$  and  $\forall j \in [2, k] : b_j \in \bigcup_{r \geq 1} \text{DAG}[r] \wedge b_{j-1} \in b_{j-1}.parents$ 

14: procedure ISVOTE( $b_{vote}, b_{leader}$ )
15:  function SUPPORTEDBLOCK( $b, id, r$ )
16:    if  $r \geq b.round$  then return  $\perp$ 
17:    for  $b' \in b.parents$  do
18:      if  $(b'.author, b'.round) = (id, r)$  then return  $b'$ 
19:       $res \leftarrow \text{SUPPORTEDBLOCK}(b', id, r)$ 
20:      if  $res \neq \perp$  then return  $res$ 
21:    return  $\perp$ 
22:     $(id, r) \leftarrow (b_{leader}.author, b_{leader}.round)$ 
23:    return  $\text{SUPPORTEDBLOCK}(b_{vote}, id, r) = b_{leader}$ 

24: procedure ISCERT( $b_{cert}, b_{leader}$ )
25:   $res \leftarrow |\{b \in b_{cert}.parents : \text{ISVOTE}(b, b_{leader})\}|$ 
26:  return  $res \geq 2f + 1$ 

27: procedure SKIPPEDLEADER( $w$ )
28:   $r_{leader} \leftarrow \text{LEADERROUND}(w)$ 
29:   $id \leftarrow \text{GETPREDEFINEDLEADER}(r_{leader})$ 
30:   $B \leftarrow \text{GETFIRSTVOTINGBLOCKS}(w)$ 
31:   $res \leftarrow |\{b \in B \text{ s.t. } \forall b' \in b.parents : b'.author \neq id\}|$ 
32:  return  $res \geq 2f + 1$ 

33: procedure SUPPORTEDLEADER( $w$ )
34:   $b_{leader} \leftarrow \text{GETLEADERBLOCK}(w)$ 
35:   $B \leftarrow \text{GETDECISIONBLOCKS}(w)$ 
36:  if  $|\{b' \in B : \text{ISCERT}(b', b_{leader})\}| \geq 2f + 1$  then return  $b_{leader}$ 
37:  return  $\perp$ 

38: procedure CERTIFIEDLINK( $b_{anchor}, b_{leader}$ )
39:   $w \leftarrow \text{WAVENUMBER}(b_{leader}.round)$ 
40:   $B \leftarrow \text{GETDECISIONBLOCKS}(w)$ 
41:  return  $\exists b \in B$  s.t.  $\text{ISCERT}(b, b_{leader}) \ \& \ \text{LINK}(b, b_{anchor})$ 

```

---

decision round ( $r + 2$ ) there are  $2f + 1$  certificates of a block  $B$  proposed for  $S$ , then the slot is tagged as ‘to-commit’.

Additionally, we also run a ‘skip’ check. This helps reduce the undecided slots (which was a non-issue in the basic MYSTICETI-C). Our check ensures that there are at least  $2f + 1$  blocks at  $r + 2$  that do not support each of the (potentially equivocated) blocks proposed for  $S$ . If this is true then we tag

---

**Algorithm 3** Baseline MYSTICETI-C

---

```

1: procedure TRYCOMMIT( $r_{committed}, r_{highest}$ )
2:   $w_{committed} \leftarrow \text{WAVENUMBER}(r_{committed})$ 
3:   $w_{highest} \leftarrow \text{WAVENUMBER}(r_{highest})$ 
4:   $sequence \leftarrow []$ 
5:  for  $w \in [w_{committed} + 1$  up to  $w_{highest}]$  do
6:     $status \leftarrow \text{TRYDIRECTDECIDE}(w)$ 
7:    if  $status = \text{Commit}(b_{leader})$  then
8:       $res \leftarrow \text{INDIRECTDECIDE}(w, b_{leader})$ 
9:       $sequence \leftarrow sequence || res$ 
10:    $sequence \leftarrow sequence || status$ 
11:  return  $sequence$ 

12: procedure TRYDIRECTDECIDE( $w$ )
13:  if SKIPPEDLEADER( $w$ ) then return Skip( $w$ )
14:   $b_{leader} \leftarrow \text{SUPPORTEDLEADER}(w)$ 
15:  if  $b_{leader} \neq \perp$  then return Commit( $b_{leader}$ )
16:  return  $\perp$ 

17: procedure INDIRECTDECIDE( $w, b_{anchor}$ )
18:   $sequence \leftarrow []$ 
19:   $b \leftarrow b_{anchor}$ 
20:   $w_{committed} \leftarrow \text{WAVENUMBER}(b_{anchor}.round)$ 
21:  for  $w' \in [w - 1$  down to  $w_{committed} + 1]$  do
22:     $b_{leader} \leftarrow \text{GETLEADERBLOCK}(w')$ 
23:    if CERTIFIEDLINK( $b, b_{leader}$ ) then
24:       $sequence \leftarrow sequence || \text{Commit}(b_{leader})$ 
25:       $b \leftarrow b_{leader}$ 
26:    else
27:       $sequence \leftarrow sequence || \text{Skip}(w')$ 
28:  return  $sequence$ 

```

---

$S$  as ‘to skip’, this means that we know that even if a future leader decides to commit it will not find a certificate during the indirect commit rule for  $S$  so we can already decide to skip (see Section 5). Finally, if neither of the two holds, we keep the slot as ‘undecided’.

Once we have run all these checks for all slots in round  $r$  we check the head of the slot-queue. If the slot is ‘to-commit’ we commit the block and its causal history; if the slot is ‘to-skip’ then we skip the slot. We continue emptying the queue until it is empty (i.e. if we are at  $r + 2$  we have decided every slot in or before  $r$ ) or the head of the queue is ‘undecided’. One final optimization that we do is that if all slots of a round are decided to skip then, since this round was the decision round of some previous block, we automatically assign the next round to become the decision round. Similarly to the baseline protocol, this is safe because all the in-between slots are skips and hence do not affect the state. Algorithm 4 describes the universal commit rule.

### 3.5 Block and commit timestamp

One final functionality we want to have in MYSTICETI is that of exposing timestamps. MYSTICETI includes a timestamp in each block and for each commit. Validators include the current time in each block they create. When a block is received its timestamp is validated by checking that the time included

**Algorithm 4** Universal MYSTICETI-C

---

```

1: procedure TRYCOMMIT( $r_{committed}, r_{highest}$ )
2:    $sequence \leftarrow []$ 
3:   for  $r \in [r_{highest} \text{ down to } r_{committed} + 1]$  do
4:     for  $l \in [k - 1 \text{ down to } 0]$  do
5:        $i \leftarrow r \% \text{wave\_length}$ 
6:        $c \leftarrow \text{BaseCommitter}(i, l)$ 
7:        $w \leftarrow c.\text{WAVELENGTH}(r)$ 
8:       if  $c.\text{LEADERROUND}(w) \neq r$  then continue
9:        $status \leftarrow c.\text{TRYDIRECTDECIDE}(w)$ 
10:      if  $status = \perp$  then
11:         $status \leftarrow \text{TRYINDIRECTDECIDE}(c, w, sequence)$ 
12:       $sequence \leftarrow status || sequence$ 
13:     $decided \leftarrow []$ 
14:    for  $status \in sequence$  do
15:      if  $status = \perp$  then break
16:       $decided \leftarrow decided || status$ 
17:    return  $decided$ 

18: procedure TRYINDIRECTDECIDE( $c, w, sequence$ )
19:    $r_{decision} \leftarrow c.\text{DECISIONROUND}(w)$ 
20:    $anchors \leftarrow [s \in sequence \text{ s.t. } r_{decision} < s.\text{round}]$ 
21:   for  $a \in anchors$  do
22:     if  $a = \perp$  then return  $\perp$ 
23:     if  $a = \text{Commit}(b_{anchor})$  then
24:        $b_{leader} \leftarrow c.\text{GETLEADERBLOCK}(w)$ 
25:       if  $c.\text{CERTIFIEDLINK}(b_{anchor}, b_{leader})$  then
26:         return  $\text{Commit}(b_{leader})$ 
27:       else
28:         return  $\text{Skip}(w)$ 
29:   return  $\perp$ 

```

---

is greater or equal to the timestamps of included blocks, otherwise, reject the block as invalid. Honest validators will only include blocks into their blocks with past timestamps, and if a block is received with a future timestamp a validator must wait before including (or rejecting) it.

As a result, if a Byzantine validator introduces a block too far in the future, such a block will be rejected. The small variation in the local clocks of validators is mitigated by the implementation, by suspending the block in memory for a short duration if that block’s timestamp is only slightly ahead of the current local time.

When MYSTICETI consensus emits a commit, it also associates a timestamp with this commit, known as a commit timestamp. Commit timestamp is denoted as a maximum of the timestamp(s) of leader block(s) of such commit and the timestamp of the previous commit. As such, MYSTICETI commit timestamps are guaranteed to be monotonically increasing. We have to include the commit timestamp of the previous commit in the maximum, since when pipelining the consecutive commit leader blocks are not necessarily linked by parent-child relationship, and thus cannot guarantee monotonicity.

**4 The MYSTICETI-FPC Protocol**

In MYSTICETI-FPC validators include transactions in their blocks as in MYSTICETI-C, but also include explicit votes for past transactions within their blocks. A correct validator votes for a transaction if it does not conflict with any other transactions they voted for. A block that causally contains  $2f + 1$  votes for a transaction certifies the transaction, and once a transaction is certified validators may execute it unless it requires consensus. When a commit occurs only transactions certified by committed blocks are included in the common sequence of transactions to be executed.

**4.1 Fast-path with MYSTICETI-FPC**

Section 3 discussed how to achieve consensus using MYSTICETI-C. Nevertheless, the benefits of MYSTICETI-C can also be extended to blockchains that have a consensusless path such as Sui Lutriss [7]. These hybrid blockchains use the observation that some objects, such as coins, that only touch state controlled by a single party need not undergo consensus – they can be safely finalized through a fast path utilizing reliable broadcast. Such objects are said to have an ‘owned-object’ type. We call transactions that have all their inputs as owned objects as fast-path transactions. Unlike prior works [7], the fast path in MYSTICETI-FPC is embedded inside the block DAG structure itself. This removes the need for a validator to sign each fast-path transaction individually. Instead, a validator’s fast path votes are piggybacked on its signed blocks being produced already as part of the consensus protocol. This simple optimization has three benefits:

1. The number of signature generation and verification operations is significantly reduced, thus the compute bottleneck is alleviated.
2. A separate post-consensus checkpointing mechanism is no longer required, thus reducing the sync latency. The consensus commits themselves serve as checkpoints.
3. The epoch close mechanism (Section 4.2) is simplified.

Although these are key benefits compared to prior work, there is an interesting tradeoff when we compare MYSTICETI-FPC to MYSTICETI-C. We now first describe the MYSTICETI-FPC protocol and then discuss the latency tradeoff with MYSTICETI-C protocol.

In addition to the block contents of MYSTICETI-C, blocks in MYSTICETI-FPC also contain explicit votes for transactions that have at least one owned-object input. Transactions are received from users, and a validator includes a transaction in its block if it does not conflict with any other transaction that it has voted for in the past. Note that a validator implicitly votes for the transactions included in its blocks. A validator, in its block  $B$ , includes an explicit vote for a transaction  $t$ , if (1)  $t$  appears in the causal history of  $B$ ; and (2)  $t$  does not conflict with any other transaction that it has voted for in the past. In our implementation, we represent the vote for a transaction  $t$  appearing in block  $b$  at position  $i$  as the tuple

$(b, i)$ . Votes for transactions appearing consecutively from the positions  $s$  to  $e$  are succinctly represented as the tuple  $(b, s, e)$ . Of course, votes for a transaction that appears at multiple positions or in multiple blocks are added up against the same transaction.

**Fast-path execution.** A validator may safely execute a fast path transaction as soon as it observes blocks from  $2f + 1$  validators that contain a vote for the transaction. Thanks to the quorum intersection, no correct validator will ever execute two conflicting fast-path transactions.

**Fast-path finality.** A block is said to certify a transaction if it links to blocks from at least  $2f + 1$  validators that contain a vote for the transaction. When blocks from  $2f + 1$  validators certify a fast path transaction, it is considered final (see Theorem 5.18), i.e. its effects are guaranteed to persist across epoch boundaries and reconfiguration of validators.

**Consensus path.** Consensus path transactions are executed after and in the order they are finalized by the commit rule. The commit rule of MYSTICETI-FPC is identical to that of MYSTICETI-C except for one important detail of the algorithm to convert the committed sequence of blocks to the committed sequence of transactions contained in those blocks. While MYSTICETI-C allows this algorithm to commit all the contained transactions, MYSTICETI-FPC places one filter on the algorithm — transactions that have at least one owned-object input are committed only if the blocks contain  $2f + 1$  votes for the transaction. Thanks to the quorum intersection argument, this filter ensures, in Theorem 5.19, that the transactions executed by the fast path and the consensus path do not conflict with each other. Notice that a fast path transaction that is committed by the consensus path will have enough votes to be executed by the fast path already.

## 4.2 Epoch close mechanism

Blockchain protocols usually operate in epochs, in order to allow validators to leave and new validators to join the system at epoch boundaries. Additionally, the epoch boundary also serves as a natural point, for protocols that have a consensusless path, to *unlock* those transactions that have lost liveness due to equivocation from the client [7]. This committee *reconfiguration* must ensure one key safety property — transactions finalized in an epoch should continue to persist across subsequent epochs. In other words, no transactions that may be committed in future epochs should conflict with transactions finalized in the current epoch.

The safety of reconfiguration is ensured by including all the finalized transactions from the current epoch in the causal history of the epoch’s final commit, which also serves as the starting state for the next epoch. Reconfiguration safety is trivial to guarantee in systems that require all transactions to undergo consensus, such as MYSTICETI-C, due to

the total ordering property of consensus. A deterministic consensus commit  $c$  serves as the epoch boundary between epoch  $e$  and  $e + 1$ , such that all transactions finalized in epoch  $e$  appear in and before commit  $c$ .

Reconfiguration solutions are however non-trivial for systems that have a consensusless fast path, such as MYSTICETI-FPC. There is a race between finalized transactions being included in the consensus commits and new transactions being finalized by the fast path. If the epoch is closed trivially, the final commit of the epoch may fail to include all the transactions finalized by the fast path, violating the key safety property of reconfiguration. In what follows, we describe the mechanism for closing an epoch safely in MYSTICETI-FPC.

Recall that a block  $B$  is said to certify a transaction  $tx$  if the causal history of  $B$  contains  $2f + 1$  votes for  $tx$ . However, we now introduce an overriding bit, called *epoch-change bit* in MYSTICETI-FPC blocks, which when set to 1 (default set to 0), indicates that the block does not certify any transaction, regardless of the causal history. This epoch-change bit, in effect, allows for pausing the consensusless fast path of MYSTICETI-FPC near the closing of the epoch, thus averting the race condition highlighted above.

Epoch change begins at a pre-determined commit, for example, when a smart contract indicates that the new committee is ready to take over. As soon as an honest validator determines that the epoch change has begun, it stops including transactions or casting votes for any fast path transactions and sets the epoch-change bit to 1 in all its future blocks for this epoch. Additionally, although the validator continues to advance its rounds and participate in consensus, it stops contributing to the processing and finalization of fast-path transactions. Once blocks from  $2f + 1$  validators with the epoch-change bit set are committed by the consensus path, the epoch is considered closed. As we prove in Theorem 5.18, this epoch close mechanism guarantees the key safety property that transactions finalized in an epoch continue to persist in all subsequent epochs. The liveness of the algorithm is trivial as it simply piggybacks on the liveness of MYSTICETI-C. Once the epoch is considered closed, any continuing validator may unlock the fast path transactions, which could not get finalized due to equivocation by the client, for fresh votes in the subsequent epochs allowing for equivocation tolerance.

## 5 Security Arguments

In this Section, we show the correctness of MYSTICETI. A validator  $v_k$  broadcasts messages calling  $a\_bcast_k(b, r)$ , where  $b$  is a block signed by validator  $v_k$  and  $r$  is the block’s round number, i.e.  $r = b.round$ . Every validator  $v_i$  has an output  $a\_deliver_i(b, b.round, v_k)$ , where  $v_k$  is the author of  $b$  and the validator that called the corresponding  $a\_bcast_k(b, b.round)$ .

**Lemma 5.1.** *If at a round  $x$ ,  $2f + 1$  blocks from distinct authorities certify a block  $B$ , then all blocks at future rounds ( $> x$ ) will link to a certificate for  $B$  from round  $x$ .*

*Proof.* Each block links to  $2f + 1$  blocks from the previous round. For the sake of contradiction, assume that a block in round  $r (> x)$  does not link to a certificate from round  $x$ . If  $r = x + 1$ , by the standard quorum intersection argument, a correct validator equivocated in round  $x$ , which is a contradiction. Similarly, if  $r > x + 1$ , by the standard quorum intersection argument, a correct validator's block in round  $r - 1$  does not link to its own block in round  $x$ , which is also a contradiction.  $\square$

**Lemma 5.2.** *If a correct validator commits some block in a slot  $s$ , then no other correct validator decides to directly skip the slot  $s$ .*

*Proof.* A validator  $X$  decides to directly skip a slot  $s$  if there is no support during the support rounds for any block corresponding to  $s$ . If another validator committed some block  $b$  for slot  $s$ , at least  $f + 1$  correct validators supported  $b$ . By the quorum intersection argument,  $X$  must have observed at least one validator supporting  $B$ , which is a contradiction.  $\square$

**Lemma 5.3.** *If a correct validator directly commits some block in a slot  $s$ , then no other correct validator decides to skip the slot  $s$ .*

*Proof.* For the sake of contradiction, assume that a correct validator  $X$  directly commits block  $b$  in slot  $s$  while another correct validator  $Y$  decides to skip the slot.  $Y$  can decide to skip the slot  $s$  in one of two ways: (a)  $Y$  directly skipped  $s$  because there was no support during the support rounds for any block corresponding to  $s$ , or (b)  $Y$  skipped  $s$  during the recursive commits triggered by a direct commit of a later slot.

Case (a). Direct contradiction of Lemma 5.2.

Case (b). Let block  $b'$  denote the leader block, committed during the recursive indirect commits, that allowed  $Y$  to decide  $s$  as skipped. Due to the commit rule, the round number of  $b'$  is greater than the decision round of  $s$ , and  $b'$  does not link to a certificate for  $b$ . Since  $X$  committed  $b$ , there are  $2f + 1$  certificates for  $b$  in its decision round, leading to a contradiction due to Lemma 5.1.  $\square$

**Lemma 5.4.** *For any slot  $s \equiv (v, r)$ , a correct validator never supports two distinct block proposals from validator  $v$  in round  $r$  across all of its blocks.*

*Proof.* By definition, a block can only support at most a single proposal for a particular slot  $s$ . Block support is calculated through a depth-first traversal of the referenced blocks, such that the first block corresponding to  $s$  encountered during the

traversal is supported. Since a correct validator first includes a reference to its own block from the previous round, once a correct validator supports a certain block for  $s$ , it continues to support the same block in all of its future blocks.  $\square$

**Lemma 5.5.** *For any slot, at most a single block will ever be certified, i.e. gather a quorum ( $2f + 1$ ) of support.*

*Proof.* For contradiction's sake, assume that two distinct block proposals for a slot gather a quorum of support. By the standard quorum intersection argument, a correct validator supports two distinct blocks for the same slot, which is a contradiction of the proved Lemma 5.4.  $\square$

As a result of Lemma 5.5, we get the following corollary:

**Corollary 5.6.** *No two correct validators commit distinct blocks for the same slot.*

**Lemma 5.7.** *All correct validators have a consistent state for each slot, i.e. if two validators have decided the state of a slot, then both either commit the same block or skip the slot.*

*Proof.* Let  $[x_i]_{i=0}^n$  and  $[y_i]_{i=0}^m$  denote the state of the slots for two correct validators  $X$  and  $Y$ , such that  $n$  and  $m$  are respectively the indices of the highest committed slot. WLOG  $n \leq m$ . Any slot decided by  $X$  higher than  $n$  are direct skips and are therefore consistent with  $Y$  due to Lemma 5.2. We now prove, by induction, statement  $P(i)$  for  $0 \leq i \leq n$ : if  $X$  and  $Y$  both decide the slot  $i$ , then both either commit the same block or skip the slot.

Base Case:  $i = n$ .  $X$  directly commits the slot  $i$  as it is the highest committed slot for  $X$ . Due to Lemma 5.3, if  $Y$  decides the slot  $i$ , then it must also commit the slot  $i$ . By Corollary 5.6,  $Y$  commits the same block.

Assuming  $P(i)$  is true for  $k + 1 \leq i \leq n$ , we now prove  $P(k)$ . Similar to the base case, if one validator decides to directly commit a block in slot  $k$ , then the other validator, if it also decides slot  $k$ , decides to commit the same block. If one validator decides to directly skip slot  $k$ , then the other validator, if it also decides slot  $k$ , decides to skip due to Lemma 5.2. We now analyze the only remaining case where  $X$  and  $Y$  indirectly decide the slot  $k$ . Let  $k'$  denote the first slot  $> k$  with a round number higher than the decision round of  $k$ . There exist slots  $k_x (\geq k')$  and  $k_y (\geq k')$  such that  $X$  commits block  $b_x$  in  $k_x$  while skipping all slots in  $[k', k_x)$  and  $Y$  commits block  $b_y$  in  $k_y$  while deciding to skip all slots in  $[k', k_y)$ . As  $k_x \leq n$ , it follows from the induction hypothesis that  $k_x = k_y$  and  $b_x = b_y = b$ . Since the indirect decision of  $X$  and  $Y$  for slot  $k$  depends entirely on the causal history of the same block  $b$ , both validators decide the slot  $k$  identically.  $\square$

**Lemma 5.8.** *All correct validators commit a consistent sequence of leader blocks (i.e., the committed leader sequence of one correct validator is a prefix of another's).*



*Proof.* The committed sequence of leader blocks is nothing but the sequence of committed blocks before the first undecided slot. The statement is then a direct implication of Lemma 5.7.  $\square$

**Theorem 5.9** (Total Order). *MYSTICETI-C satisfies the total order property of Byzantine Atomic Broadcast.*

*Proof.* Correct validators deliver blocks by using an identical deterministic algorithm to order the causal history of committed leader blocks. Since a correct validator has all of the causal histories of a block when the block is added to its DAG, and the sequence of committed leader blocks of one validator is a prefix of another’s (Lemma 5.8), all correct validators deliver a consistent sequence of blocks, i.e. the sequence of blocks delivered from one validator is a prefix of another. The total order property of BAB immediately follows.  $\square$

**Theorem 5.10** (Integrity). *MYSTICETI-C satisfies the integrity property of Byzantine Atomic Broadcast.*

*Proof.* The algorithm to linearize the causal history of a committed leader block removes any block with duplicate sequence numbers before delivering the sequence of blocks.  $\square$

**Lemma 5.11** (Round-Synchronization). *After GST all honest parties will enter the same round within  $\Delta$ .*

*Proof.* After GST all messages sent before GST deliver within  $\Delta$ . This means that if  $r$  is the highest round any honest validator proposed a block for before GST, then every honest validator will receive the block proposal of the honest validator at  $GST + \Delta$  and also enter  $r$ .  $\square$

**Lemma 5.12** (Leader-Proposal). *After GST an honest leader’s proposal will get votes from every honest validator.*

*Proof.* After GST if an honest validator enters wave  $w$ , then it has to broadcast the last block of wave  $w - 1$ . Within  $\Delta$  the honest leader (and every other honest party) will receive the block and adopt the parents, being able to also enter wave  $w$  as they are all synchronized (Lemma 5.11). Then the honest leader will directly propose its block. Since the timeout is set to  $2 * \Delta$  the leader’s proposal of wave  $w$  will arrive before the first honest validator times out hence, every honest validator will vote for the leader.  $\square$

**Lemma 5.13** (Sufficient Votes). *After GST all honest validators will create a certificate for the honest leader.*

*Proof.* By Lemma 5.12 all honest validators will vote for an honest leader after GST. For an honest validator to propose a block at the decision round it needs to (a) get the proposal of the leader and (b) have  $2f + 1$  parents. All honest validators receive the leader proposal within  $\Delta$  since the leader is honest. Additionally from the moment one honest validator advances to the decision round all honest validators will

receive its block proposal and adopt the parents within  $\Delta$ . As a result, by the code, all honest validators wait for  $2 * \Delta$  before giving up the certificate creation and will receive the votes from all honest validators witnessing a certificate  $\square$

**Lemma 5.14.** *The round-robin schedule of leaders in MYSTICETI ensures that in any window of  $3f + 3$  rounds, there are three consecutive rounds with honest primary leaders. A primary leader is the leader of the first slot of a round.*

*Proof.* There are  $3f + 1$  groups of three consecutive rounds. Due to the round-robin schedule, each of the honest validators must be the primary leader in exactly 3 of these groups. As there are  $2f + 1$  honest validators, due to the pigeon-hole principle, one group must contain  $\lceil \frac{3*(2f+1)}{3f+1} \rceil = 3$  honest leaders.  $\square$

**Lemma 5.15.** *After GST any undecided slot eventually gets decided.*

*Proof.* Let there be an undecided slot  $s$  in round  $r$ . After GST, due to Lemma 5.14, there will eventually be an honest leader for the first slots  $s_0, s_1$  and  $s_2$  of rounds  $k, k + 1$  and  $k + 2$  respectively, where  $k > r$ . By Lemma 5.13, the honest leader’s blocks will have  $2f + 1$  certificates and be scheduled for a commit. We now prove that by induction, all slots in round  $\leq k - 1$  get decided. In the base case, any undecided slots in rounds  $k - 3, k - 2$  or  $k - 1$  get decided by the commits in slots  $s_0, s_1$  and  $s_2$  respectively, as they are the first slots higher than the respective decision rounds. For the induction step, any undecided slot  $s$  in round  $x \leq k - 4$  also gets decided since  $s_0$  is higher than the decision round of  $x$  and there are no undecided slots between  $s$  and  $s_0$  due to the induction hypothesis.  $\square$

**Theorem 5.16** (Consensus Liveness). *After GST an honest leader’s proposal will commit.*

*Proof.* By Lemma 5.13 there will be  $2f + 1$  certificates for the leader, one per honest party. By the code an honest validator tries to commit the leader for every block they get so eventually they will get the  $2f + 1$  certificates. The validator schedules the block to be committed. By Lemma 5.15, all prior undecided blocks will eventually be decided, and the validator will deliver the honest leader’s block.  $\square$

**Theorem 5.17** (Agreement). *MYSTICETI-C satisfies the agreement property of Byzantine Atomic Broadcast.*

*Proof.* If a correct validator outputs  $a\_deliver_i(b, r, v_k)$ , then it must have committed a sequence of leader blocks  $L = l_0, l_1 \dots l_n$  such that the deterministic algorithm to deliver blocks from the sequence  $L$  delivers block  $b$ . Another correct validator  $Y$  that has not delivered  $b$  will eventually see a proposal  $b'$  from an honest leader in round  $r' > r$  as per the leader schedule of MYSTICETI-C. Due to Theorem 5.16, after GST,  $Y$  will commit the leader’s block  $b'$ . Due to Lemma 5.8,  $Y$  will also commit the leader sequence  $L$  before committing

$b'$ . Since  $Y$  follows an identical deterministic algorithm as  $X$  to deliver blocks from the committed sequence of leader blocks, it also delivers  $b'$  eventually.  $\square$

### 5.1 Security Arguments of MYSTICETI-FPC

**Theorem 5.18** (Epoch close safety). *Transactions finalized by MYSTICETI-FPC in an epoch continue to persist in all subsequent epochs.*

*Proof.* It is sufficient to prove that all fast-path transactions that are considered final have one certifying block committed in the current epoch. For contradiction's sake, assume that the epoch closed before any certifying block for a finalized transaction  $tx$  could be committed. For the epoch to close, blocks from  $2f + 1$  validators with the epoch-change bit set must be committed. Since  $tx$  is finalized,  $2f + 1$  validators, by definition, publish a block that certifies the transaction. By quorum intersection, one honest validator  $v$  published a block  $B_1$  in round  $r_1$  certifying transaction  $tx$ , whereas a block  $B_2$  in round  $r_2$  from  $v$  with epoch-change bit set must have been committed. All blocks published by  $v$  in rounds  $\geq r_2$  also have the epoch-change bit set. Because blocks with the epoch-change bit set, by definition, do not certify any transaction,  $B_1$  is necessarily published in an earlier round than that of  $B_2$  (i.e.  $r_1 < r_2$ ).  $B_1$  is therefore contained in the causal history of  $B_2$ , and must also have been committed, which is a contradiction.  $\square$

**Theorem 5.19** (MYSTICETI-FPC Safety). *An honest validator in MYSTICETI-FPC never finalizes two conflicting transactions.*

*Proof.* Transactions that have an owned object as input require votes from  $2f + 1$  validators to be finalized. If two conflicting fast paths are finalized, an honest validator must have voted for both transactions (by quorum intersection), hence a contradiction. Using a similar argument, a fast path transaction does not conflict with a consensus path transaction, as the consensus path in MYSTICETI-FPC finalizes a transaction with an owned object input only if it has votes from  $2f + 1$  validators.  $\square$

**Theorem 5.20** (Fast-Path Liveness). *An honest fast-path transaction will commit after GST.*

The proof is the same as consistent broadcast. We do it after GST assuming the epoch does not end. If the epoch has infinite length then we can convert all references to  $\Delta$  with “eventually” and the proof will work in asynchrony.

*Proof.* An honest validator will submit a fast-path transaction that does not have equivocation. As a result, all honest validators will receive it after  $\Delta$  and vote. These votes will appear in the DAG after at most  $4 * \Delta$  since any round has at most duration of  $\text{timeout} + \Delta = 3 * \Delta$ . In the next round, every honest validator will reference the  $2f + 1$  votes in their DAG and execute.  $\square$

**Theorem 5.21** (Equivocation-Tolerance). *If a faulty validator  $v_k$  concurrently called  $r\_bcast_k(m, q, e)$  and  $r\_bcast_k(m', q, e)$  with  $m \neq m'$  then the rest of the validators either  $r\_deliver_i(m, q, v_k, e)$ , or  $r\_deliver_i(m', q, v_k, e)$ , or there is a subsequent epoch  $e' > e$  where  $v_k$  is honest, calls  $r\_bcast_k(m'', q, e')$  and all honest validators  $r\_deliver_i(m'', q, v_k, e')$ ,*

*Proof.* For the case that validators  $r\_deliver_i(m', q, v_k, e)$  it is a direct result of Theorem 5.20. Otherwise, from the code of the epoch change when the epoch ends all validators forget the locks they have taken on messages without certificates. As a result in a future epoch  $e'$  where  $v_k$  is honest and does not equivocate it will be able to commit  $m$  again from Theorem 5.20.  $\square$

## 6 Implementation

We implement a networked multi-core MYSTICETI validator in Rust. It uses `tokio`<sup>3</sup> for asynchronous networking, utilizing TCP sockets for communication without relying on any RPC frameworks. For cryptographic operations, we rely on the efficient `ed25519-consensus`<sup>4</sup> for asymmetric cryptography and `blake2`<sup>5</sup> for cryptographic hashing. To ensure data persistence and crash recovery, we've integrated a Write-Ahead Log (WAL), seamlessly tailored to our specific requirements. We've intentionally avoided key-value stores like `RocksDB`<sup>6</sup> to eliminate associated overhead and periodic compaction penalties. Our implementation optimizes I/O operations by employing vectored writes<sup>7</sup> for efficient multi-buffer writes in a single syscall. For reading the WAL, we make use of memory-mapped files while carefully minimizing redundant data copying and serialization. We use the `minibytes`<sup>8</sup> crates to efficiently work with memory-mapped file buffers without unsafe code.

While all network communications in our implementation are asynchronous, the core consensus code runs synchronously in a dedicated thread. This approach facilitates rigorous testing, mitigates race conditions, and allows for targeted profiling of this critical code path.

In addition to regular unit tests, we have two supplementary testing utilities. First, we've developed a simulation layer that replicates the functionality of the `tokio` runtime and TCP networking. This simulated networking layer accurately simulates real-world WAN latencies, while our `tokio` runtime simulator employs a discrete event simulation approach to mimic the passage of time. Utilizing this simulator, we can test a wide range of scenarios on a single machine and accurately estimate resulting latencies. It's worth noting that we've found these simulated latencies, such as commit latency, to closely mirror those observed in real-world

<sup>3</sup><https://tokio.rs>

<sup>4</sup><https://github.com/penumbra-zone/ed25519-consensus>

<sup>5</sup><https://github.com/RustCrypto/hashes>

<sup>6</sup><https://rocksdb.org>

<sup>7</sup><https://linux.die.net/man/3/writev>

<sup>8</sup><https://github.com/facebook/sapling/tree/main/eden/scm/lib/minibytes>

cluster testing, provided that the cross-validator latency distribution in the simulated network is correctly configured. Second, we've created an "orchestrator" - a command-line utility designed to set up real-world clusters of MYSTICETI with machines distributed across the globe. The simulator has proven indispensable in identifying correctness defects, while the orchestrator has been instrumental in pinpointing and addressing efficiency bottlenecks.

We are open-sourcing our MYSTICETI implementation, along with its simulator and orchestration utilities<sup>9</sup>.

## 7 Evaluation

We evaluate the throughput and latency of MYSTICETI through experiments on Amazon Web Services (AWS). We then show its performance improvements over several state-of-the-art protocols. Despite the large number of BFT consensus protocols [11, 12, 14, 19, 26, 32, 40], we opt to compare MYSTICETI-C with vanilla HotStuff [41], HotStuff-over-Narwhal (called *Narwhal-HotStuff*) [17], and Bullshark [30]. We select these protocols for the availability of open-source implementations and detailed benchmarking scripts, their similarity to MYSTICETI, and their adoption in real-world deployments. We specifically select the Jolteon [20] variant of HotStuff as it has been adopted by Flow [36], Diem [4], Aptos [34], and Monad [27]. We also select the Narwhal-HotStuff variant as it operates on a structured DAG as MYSTICETI and is the most performant variant of HotStuff. We finally select Bullshark as it is a performant DAG-based protocol adopted by the Sui blockchain [7, 38] and in the roadmap for integration within Aptos. We evaluate the 1-worker variants of the Narwhal-based systems (that is, Narwhal-HotStuff and Bullshark). We also evaluate the fast path MYSTICETI-FPC against Zef [6] (in its default configuration, with 10 shards), which is the state-of-the-art fast path protocol that serves as the foundation for the Linera blockchain [37].

Throughout our evaluation, we particularly aim to demonstrate the following claims. **C1:** MYSTICETI-C has higher throughput and drastically lower latency than the baseline state-of-the-art protocols. **C2:** MYSTICETI-C has a similar throughput to the baseline protocols but maintains sub-second latencies when operating in the presence of crash faults. **C3:** MYSTICETI-FPC maintains the same latency as the baseline state-of-the-art consensus-less protocol but with drastically higher throughput. Note that evaluating BFT protocols in the presence of Byzantine faults is an open research question [3], and state-of-the-art evidence relies on formal proofs of safety and liveness (see Section 5).

### 7.1 Experimental setup

In the following graphs, each data point is the average latency and the error bars represent one standard deviation (error bars are sometimes too small to be visible on the graph). We instantiate several geo-distributed benchmark clients within

<sup>9</sup><https://github.com/MystenLabs/mysticeti/tree/paper> (commit aee594d)

each validator submitting transactions at a fixed rate for a duration of 10 minutes. We experimentally increase the load of transactions sent to the systems, and record the throughput and latency of commits. As a result, all plots illustrate the 'steady state' latency of all systems under low load, as well as the maximal throughput they can serve after which latency grows quickly. Transactions in the benchmarks are random and contain 512 bytes, and MYSTICETI is instantiated with two leaders per round (see Section 3.4).

When referring to *latency*, we mean the time elapsed from when the client submits the transaction to when it receives confirmation of the transaction's finality. When referring to *throughput*, we mean the number of committed transactions over the entire duration of the run. Appendix A provides a tutorial to reproduce our experiments.

### 7.2 Benchmark in ideal conditions

Figure 3 illustrates the Latency (seconds) - Throughput (Transactions per second, TPS) relationship for MYSTICETI-C compared with other consensus protocols, for a small deployment of 10 validators (up to 3 tolerable failures) and a larger deployment of 50 validators (up to 16 tolerable failures). The systems run in ideal conditions, without faults.

At a steady state of 50k to 100k TPS for both network sizes MYSTICETI-C exhibits sub-second latency, a factor 2x-3x lower than the fastest protocols, namely HotStuff, and Narwhal-HotStuff. Bullshark uses a certified DAG and worker architecture and is over 3x slower in terms of latency compared with MYSTICETI-C for low system loads. In terms of throughput, the smaller MYSTICETI-C network scales extremely well and achieves a throughput of over 400k TPS before latency reaches 1.5s, that is, comparable to the latency of state-of-the-art systems. The larger deployment scales to 120k TPS before latency goes over 1.5s, which is comparable to the single worker variant of Narwhal-based designs, and HotStuff variants. This illustrates that the single-host throughput efficiency of MYSTICETI-C is higher than for previous designs. Note that current real-world blockchains combined<sup>10</sup> process fewer than 100M transactions per day, equivalent to about 1.2k TPS, well within the steady state low-latency parameter space for MYSTICETI-C, without any further scaling strategies (which we discuss later).

These observations validate our claim **C1** showing that MYSTICETI-C has higher throughput and drastically lower latency than the baseline state-of-the-art protocols.

### 7.3 Benchmark with faults

Figure 4 illustrates the performance of HotStuff, Narwhal-HotStuff, Bullshark, and MYSTICETI-C when a committee of 10 parties suffers 0 to 3 crash faults (the maximum that can be tolerated in this setting). HotStuff suffers a massive degradation in both throughput and latency. With 3 faults, the throughput of HotStuff drops to a few hundred TPS and

<sup>10</sup>Estimates from <https://app.artemis.xyz/comparables>

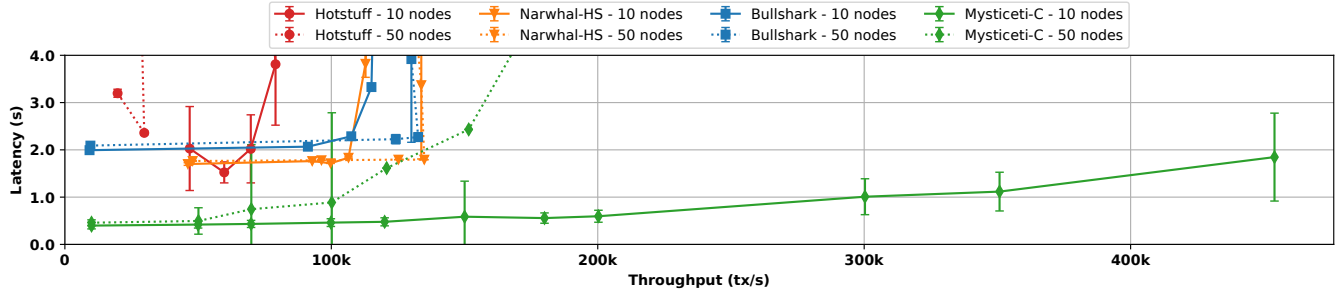


Figure 3. Throughput-Latency graph comparing MYSTICETI-C performance with state-of-the-art consensus protocols.

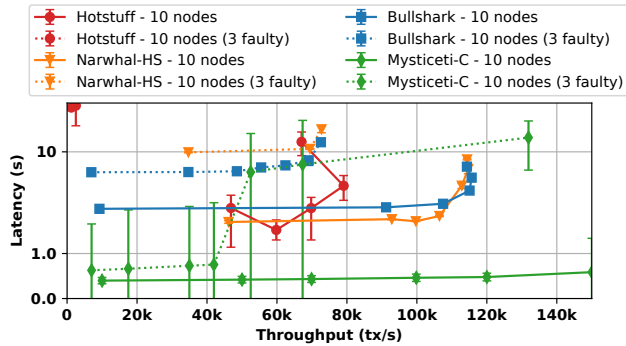


Figure 4. Throughput - Latency under crash faults. Note the log y-axis, for latency.

its latency exceeds 15s. Narwhal-HotStuff, Bullshark, and MYSTICETI-C maintain a good level of throughput: the underlying DAG continues collecting and disseminating transactions despite the faults. Narwhal-HotStuff and Bullshark can process about 60-80k TPS in about 8-10 seconds. In contrast, MYSTICETI-C can process up to 50k TPS while maintaining sub-second latency and up to 80k TPS with comparable latency to Narwhal-HotStuff and Bullshark. As a result, MYSTICETI-C demonstrates a 15-20x latency improvement compared to the baseline state-of-the-art protocols.

These observations validate our claim C2 showing that MYSTICETI-C can handle a similar throughput to state-of-the-art consensus protocols but with sub-second latency despite the presence of crash faults.

### 7.4 Benchmark of the fast path

Figure 5 illustrates the Latency - Throughput of fast path commits for MYSTICETI-FPC, compared with Zef [6] when deployed without privacy protections<sup>11</sup>. Both systems run in ideal conditions, without faults. We observe that for low loads both protocols have a comparable latency of around 0.25s. However, as the load increases a Zef host has to verify and produce an increasing number of signatures, proportional to the throughput times the number of validators. As a

<sup>11</sup>Zef can also be instantiated to leverage the Coconut threshold credentials system [29] to provide privacy guarantees at the cost of performance.

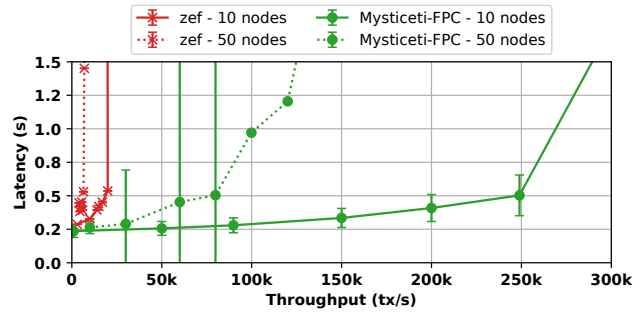


Figure 5. Throughput - Latency comparison for fast path commits between MYSTICETI-FPC and Zef

result throughput tops at 20k TPS for a small Zef network and 7K TPS for a larger network, at a latency of 0.5s. MYSTICETI-FPC avoids the need for individual signature verification for each transaction. At a low load, its latency is similar to Zef at 0.25s. However, as the load increases MYSTICETI-FPC can process many more messages on a single host, namely 175k TPS for a small network and 80K for a larger network, at a latency of less than 0.5s. This is a single host throughput improvement of 8x-10x compared with Zef. We acknowledge that the Zef design can scale by adding additional hosts per validator, and sharding. However, this leads to additional hardware cost meaning that MYSTICETI-FPC is an order of magnitude more resource efficient for the same latency.

We thus validate our claim C3 showing that MYSTICETI-FPC offers the same latency as state-of-the-art consensus-less protocols but with significantly higher throughput.

## 8 Related Work

MYSTICETI is a family of protocols designed to support next-generation distributed ledgers. To this end, its goal is to capture as wide a range of distributed ledgers as possible whether consensus-based or consensus-less. The pioneer on hybrid distributed ledgers is the Sui Lutriss blockchains [7] which has been productionized by Sui [38]. However, the design of Sui Lutriss focuses on providing a glue between the two distinct use-cases of consensus-based and consensus-less distributed ledgers, or in the production code a glue



of FastPay [5] and Bullshark [30]. This design process of starting with the to-be-glued components and ending in a final system has led to significant inefficiencies such as multiple rebroadcasting of the same data as well as signature verification costs. Unlike Sui Lutris, MYSTICETI is designed from first principles and as a result shows a potential halving of the latency, matching the lower bounds of PBFT [9] for consensus and Reliable Broadcast [8] for consensus-less distributed ledgers with equivocation tolerance.

We now focus on the different variants of MYSTICETI, namely MYSTICETI-C and MYSTICETI-FPC. We already discussed the core benefits of MYSTICETI-FPC in terms of much lower CPU cost. In addition, it inherits the ability to change epochs, reconfigure the validator set, and tolerate equivocations from Sui Lutris. These benefits can also be used to embed other broadcast-based protocols like FastPay [5], Zef [6], and Astro [15] in order to get better privacy guarantees.

In terms of consensus, the most recent DagRider [23], Narwhal-Tusk [17], Bullshark [30], DisperseLedger [40] were the inspiration for using a structured DAG and defining a safe commit rule on it. However, they all use a DAG of certified blocks which increases both latency and implementation complexity. MYSTICETI uses instead a DAG of signed but not certified blocks, reducing latency significantly. Notably, Narwhal-based designs use a worker-primary architecture to increase throughput. MYSTICETI-C can be adapted to this architecture, by acting as a primary for any number of workers in case additional throughput is needed. The concurrent work of Shoal [31] defines a pipelined variant of Bullshark and other Narwhal-based protocols, this can be seen as a similar contribution to the universal committer of MYSTICETI if set with one leader per round. We plan in our future work to adapt our universal committer to Bullshark in order to have a fair comparison over certified DAGs. Additionally, Shoal and HammerHead [39] propose leader reputation protocols inspired by Carousel [13]. MYSTICETI-C could adopt these designs to select more reliable leaders, but for liveness, it would need to adopt a leader rotation schedule where leaders remain static for 3 rounds.

Previous consensus protocols such as Hashgraph [2] and Blockmania [16] also use a DAG of signed but not certified blocks: however they use DAGs that are not structured as threshold clocks [18] making the proofs of safety for them very complex. Notably, MYSTICETI-C works in only 3 message communication rounds which match PBFT, and is optimal latency [1, 10] without the use of optimistic methods like Zyzzyva [24]. This is lower than the state-of-the-art Jolteon [21] currently deployed in multiple blockchains [27, 34–36]. The reason is that these protocols focus on linear communication complexity whereas MYSTICETI-C embraces its cubic cost and amortizes it using the DAG structure as first proposed by DagRider and Narwhal. Finally, non-byzantine variants of consensus have also been defined on threshold

clocks, such as Que-Paxa [33] but cannot be deployed as part of distributed ledgers.

## 9 Conclusion

We introduced MYSTICETI, a threshold clock-based byzantine consensus protocol with the lowest WAN latency of 0.5s and the ability to process over 50k TPS at this latency for single-host nodes, far exceeding the needs of blockchains today (which consume in total about 1.2k TPS). Its fast path achieves even lower latency at 0.25s but is over 8x more resource efficient compared with protocols with explicit certificates. It is also more crash tolerant using multiple leaders per round, implemented through a universal commit rule.

We leave several explorations for the future. For use cases requiring higher throughput, we note that MYSTICETI-C can be augmented with workers, in a similar way to Tusk and Bullshark. This would allow it to scale without known bounds, at the cost of additional latency (a round trip) to coordinate workers and primaries. An alternative approach would be to run multiple MYSTICETI-C instances in parallel, something we feel is under-explored but inspired us to have explicit votes in MYSTICETI-FPC. The structure of MYSTICETI-FPC has all nodes timestamping transactions through their votes and may be useful for implementing MEV protections.

Finally, we note that as the latency of consensus under low load shrinks (now 0.5s) the latency advantages of the fast path diminish. It is an open industrial question whether use cases that require low latency justify the complexity of dual path systems going forward, as the latency gap closes.

## Acknowledgements

We would like to thank Dmitry Perelman, Mingwei Tian, Tasos Kichidis, and Xun Li from the MystenLabs engineering team for the great discussions that improved this work. This work was conducted while Kushal Babel was interning with Mysten Labs.

## References

- [1] ABRAHAM, I., NAYAK, K., REN, L., AND XIANG, Z. Good-case latency of byzantine broadcast: A complete categorization. In *Proceedings of the 2021 ACM Symposium on Principles of Distributed Computing* (New York, NY, USA, 2021), PODC’21, Association for Computing Machinery, p. 331–341.
- [2] BAIRD, L. The swirls hashgraph consensus algorithm: Fair, fast, byzantine fault tolerance. *Swirls Tech Reports SWIRLDS-TR-2016-01, Tech. Rep 34* (2016), 9–11.
- [3] BANO, S., SONNINO, A., CHURSIN, A., PERELMAN, D., LI, Z., CHING, A., AND MALKHI, D. Twins: Bft systems made robust. In *25th International Conference on Principles of Distributed Systems* (2022).
- [4] BAUDET, M., CHING, A., CHURSIN, A., DANEZIS, G., GARILLOT, F., LI, Z., MALKHI, D., NAOR, O., PERELMAN, D., AND SONNINO, A. State machine replication in the libra blockchain. *The Libra Assn., Tech. Rep 1*, 1 (2019).
- [5] BAUDET, M., DANEZIS, G., AND SONNINO, A. Fastpay: High-performance byzantine fault tolerant settlement. In *AFT ’20: 2nd ACM Conference on Advances in Financial Technologies*, New York, NY, USA, October 21–23, 2020 (2020), ACM, pp. 163–177.

- [6] BAUDET, M., SONNINO, A., KELKAR, M., AND DANEZIS, G. Zef: Low-latency, scalable, private payments. *CoRR abs/2201.05671* (2022).
- [7] BLACKSHEAR, S., CHURSIN, A., DANEZIS, G., KICHIDIS, A., KOKORIS-KOGIAS, L., LI, X., LOGAN, M., MENON, A., NOWACKI, T., SONNINO, A., ET AL. Sui lutris: A blockchain combining broadcast and consensus. Tech. rep., Technical Report. Mysten Labs. <https://sonnino.com/papers/sui-lutris.pdf>, 2023.
- [8] BRACHA, G., AND TOUEG, S. Asynchronous consensus and broadcast protocols. *J. ACM* 32, 4 (1985), 824–840.
- [9] CASTRO, M., AND LISKOV, B. Practical byzantine fault tolerance. In *Proceedings of the Third USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, New Orleans, Louisiana, USA, February 22–25, 1999 (1999), M. I. Seltzer and P. J. Leach, Eds., USENIX Association, pp. 173–186.
- [10] CHAN, B. Y., AND PASS, R. Simplex consensus: A simple and fast consensus protocol. Cryptology ePrint Archive, Paper 2023/463, 2023. <https://eprint.iacr.org/2023/463>.
- [11] CHAN, B. Y., AND SHI, E. Streamlet: Textbook streamlined blockchains. In *Proceedings of the 2nd ACM Conference on Advances in Financial Technologies* (2020), pp. 1–11.
- [12] CHEN, J., GUPTA, S., SONNINO, A., KOKORIS-KOGIAS, L., AND SADOOGHI, M. Resilient consensus sustained collaboratively. *arXiv preprint arXiv:2302.02325* (2023).
- [13] COHEN, S., GELASHVILI, R., KOGIAS, L. K., LI, Z., MALKHI, D., SONNINO, A., AND SPIEGELMAN, A. Be aware of your leaders. In *International Conference on Financial Cryptography and Data Security* (2022), Springer, pp. 279–295.
- [14] COHEN, S., GOREN, G., KOKORIS-KOGIAS, L., SONNINO, A., AND SPIEGELMAN, A. Proof of availability & retrieval in a modular blockchain architecture. *Cryptology ePrint Archive* (2022).
- [15] COLLINS, D., GUERRAOU, R., KOMATOVIC, J., MONTI, M., XYGKIS, A., PAVLOVIC, M., KUZNETSOV, P., PIGNOLET, Y.-A., SEREDINSCHI, D.-A., AND TONKIKH, A. Online payments by merely broadcasting messages (extended version). *arXiv preprint arXiv:2004.13184* (2020).
- [16] DANEZIS, G., AND HRYSZYN, D. Blockmania: from block dags to consensus. *arXiv preprint arXiv:1809.01620* (2018).
- [17] DANEZIS, G., KOKORIS-KOGIAS, L., SONNINO, A., AND SPIEGELMAN, A. Narwhal and tusk: a dag-based mempool and efficient BFT consensus. In *EuroSys '22: Seventeenth European Conference on Computer Systems, Rennes, France, April 5 - 8, 2022* (2022), Y. Bromberg, A. Kermerrec, and C. Kozyrakis, Eds., ACM, pp. 34–50.
- [18] FORD, B. Threshold logical clocks for asynchronous distributed coordination and consensus. *CoRR abs/1907.07010* (2019).
- [19] GAO, Y., LU, Y., LU, Z., TANG, Q., XU, J., AND ZHANG, Z. Dumbo-ng: Fast asynchronous bft consensus with throughput-oblivious latency. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security* (2022), pp. 1187–1201.
- [20] GELASHVILI, R., KOKORIS-KOGIAS, L., SONNINO, A., SPIEGELMAN, A., AND XIANG, Z. Jolteon and ditto: Network-adaptive efficient consensus with asynchronous fallback. In *International Conference on Financial Cryptography and Data Security* (2022), Springer, pp. 296–315.
- [21] GELASHVILI, R., KOKORIS-KOGIAS, L., SONNINO, A., SPIEGELMAN, A., AND XIANG, Z. Jolteon and ditto: Network-adaptive efficient consensus with asynchronous fallback. In *Financial Cryptography and Data Security - 26th International Conference, FC 2022, Grenada, May 2-6, 2022, Revised Selected Papers* (2022), I. Eyal and J. A. Garay, Eds., vol. 13411 of *Lecture Notes in Computer Science*, Springer, pp. 296–315.
- [22] GELASHVILI, R., SPIEGELMAN, A., XIANG, Z., DANEZIS, G., LI, Z., MALKHI, D., XIA, Y., AND ZHOU, R. Block-stm: Scaling blockchain execution by turning ordering curse to a performance blessing. In *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming, PPOPP 2023, Montreal, QC, Canada, 25 February 2023 - 1 March 2023* (2023), M. M. Dehnavi, M. Kulkarni, and S. Krishnamoorthy, Eds., ACM, pp. 232–244.
- [23] KEIDAR, I., KOKORIS-KOGIAS, E., NAO, O., AND SPIEGELMAN, A. All you need is DAG. In *PODC '21: ACM Symposium on Principles of Distributed Computing, Virtual Event, Italy, July 26-30, 2021* (2021), A. Miller, K. Censor-Hillel, and J. H. Korhonen, Eds., ACM, pp. 165–175.
- [24] KOTLA, R., ALVISI, L., DAHLIN, M., CLEMENT, A., AND WONG, E. L. Zyzzyva: speculative byzantine fault tolerance. In *Proceedings of the 21st ACM Symposium on Operating Systems Principles 2007, SOSP 2007, Stevenson, Washington, USA, October 14-17, 2007* (2007), T. C. Bressoud and M. F. Kaashoek, Eds., ACM, pp. 45–58.
- [25] LAMPORT, L. Paxos made simple. *ACM SIGACT News (Distributed Computing Column)* 32, 4 (Whole Number 121, December 2001) (2001), 51–58.
- [26] MALKHI, D., AND SZALACHOWSKI, P. Maximal extractable value (mev) protection on a dag. *arXiv preprint arXiv:2208.00940* (2022).
- [27] MONAD. Monadbft: pipelined two-phase hotstuff consensus. <https://docs.monad.xyz/technical-discussion/consensus/monadbft> (2023).
- [28] PATRIGNANI, M., AND BLACKSHEAR, S. Robust safety for move. In *36th IEEE Computer Security Foundations Symposium, CSF 2023, Dubrovnik, Croatia, July 10-14, 2023* (2023), IEEE, pp. 308–323.
- [29] SONNINO, A., AL-BASSAM, M., BANO, S., MEIKLEJOHN, S., AND DANEZIS, G. Coconut: Threshold issuance selective disclosure credentials with applications to distributed ledgers. In *26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019* (2019), The Internet Society.
- [30] SPIEGELMAN, A., AND ALBERTO SONNINO, N. G., AND KOKORIS-KOGIAS, L. Bullshark: DAG BFT protocols made practical. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security, CCS 2022, Los Angeles, CA, USA, November 7-11, 2022* (2022), H. Yin, A. Stavrou, C. Cremers, and E. Shi, Eds., ACM, pp. 2705–2718.
- [31] SPIEGELMAN, A., AURN, B., GELASHVILI, R., AND LI, Z. Shoal: Improving DAG-BFT latency and robustness. *CoRR abs/2306.03058* (2023).
- [32] SPIEGELMAN, A., AURN, B., GELASHVILI, R., AND LI, Z. Shoal: Improving dag-bft latency and robustness. *arXiv preprint arXiv:2306.03058* (2023).
- [33] TENNAGE, P., BASESCU, C., KOKORIS-KOGIAS, L., SYTA, E., JOVANOVIĆ, P., ESTRADA-GALIÑANES, V., AND FORD, B. Quepaxa: Escaping the tyranny of timeouts in consensus. In *Proceedings of the 29th Symposium on Operating Systems Principles, SOSP 2023, Koblenz, Germany, October 23-26, 2023* (2023), J. Flinn, M. I. Seltzer, P. Druschel, A. Kaufmann, and J. Mace, Eds., ACM, pp. 281–297.
- [34] THE APTOS TEAM. Aptos. <https://aptoslabs.com>, 2023.
- [35] THE DIEM TEAM. Diembft v4. <https://developers.diem.com/papers/diem-consensus-state-machine-replication-in-the-diem-blockchain/2021-08-17.pdf>, 2021.
- [36] THE FLOW TEAM. The flow blockchain. <https://flow.com>, 2023.
- [37] THE LINERA TEAM. Linera. <https://linera.io>, 2023.
- [38] THE SUI TEAM. The sui blockchain. <http://sui.io>, 2023.
- [39] TSIMOS, G., KICHIDIS, A., SONNINO, A., AND KOKORIS-KOGIAS, L. Hammerhead: Leader reputation for dynamic scheduling. *arXiv preprint arXiv:2309.12713* (2023).
- [40] YANG, L., PARK, S. J., ALIZADEH, M., KANNAN, S., AND TSE, D. {DispersedLedger}:: {High-Throughput} byzantine consensus on variable bandwidth networks. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)* (2022), pp. 493–512.
- [41] YIN, M., MALKHI, D., REITER, M. K., GOLAN-GUETA, G., AND ABRAHAM, I. Hotstuff: BFT consensus with linearity and responsiveness. In *Proceedings of the 2019 ACM Symposium on Principles of Distributed Computing, PODC 2019, Toronto, ON, Canada, July 29 - August 2, 2019* (2019), P. Robinson and F. Ellen, Eds., ACM, pp. 347–356.

## A Reproducing Experiments

We provide the orchestration scripts<sup>12</sup> used to benchmark the codebase evaluated in this paper on AWS.

**Deploying a testbed.** The file ‘~/aws/credentials’ should have the following content:

```
[default]
aws_access_key_id = YOUR_ACCESS_KEY_ID
aws_secret_access_key = YOUR_SECRET_ACCESS_KEY
```

configured with account-specific AWS *access key id* and *secret access key*. It is advise to not specify any AWS region as the orchestration scripts need to handle multiple regions programmatically.

A file ‘settings.json’ contains all the configuration parameters for the testbed deployment. We run the experiments of Section 7 with the following settings:

```
{
  "testbed_id": "${USER}-hammerhead",
  "cloud_provider": "aws",
  "token_file": "/Users/${USER}/.aws/credentials",
  "ssh_private_key_file": "/Users/${USER}/.ssh/aws",
  "regions": [
    "us-east-1",
    "us-west-2",
    "ca-central-1",
    "eu-central-1",
    "ap-northeast-1",
    "ap-northeast-2",
    "eu-west-1",
    "eu-west-2",
    "eu-west-3",
    "eu-north-1",
    "ap-south-1",
    "ap-southeast-1",
    "ap-southeast-2"
  ],
  "specs": "m5d.8xlarge",
  "repository": {
    "url": "https://github.com/AUTHOR/REPO.git",
    "commit": "main"
  }
}
```

where the file ‘/Users/\$USER/.ssh/aws’ holds the ssh private key used to access the AWS instances, and ‘AUTHOR’ and ‘REPO’ are respectively the GitHub username and repository name of the codebase to benchmark.

The orchestrator binary provides various functionalities for creating, starting, stopping, and destroying instances. For instance, the following command to boots 2 instances per region (if the settings file specifies 13 regions, as shown in the example above, a total of 26 instances will be created):

```
cargo run --bin orchestrator -- testbed deploy --instances 2
```

<sup>12</sup><https://github.com/MystenLabs/mysticeti/tree/paper> (commit aee594d)

The following command displays the current status of the testbed instances

```
cargo run --bin orchestrator testbed status
```

Instances listed with a green number are available and ready for use and instances listed with a red number are stopped. It is necessary to boot at least one instance per load generator, one instance per validator, and one additional instance for monitoring purposes (see below). The following commands respectively start and stop instances:

```
cargo run --bin orchestrator -- testbed start
cargo run --bin orchestrator -- testbed stop
```

It is advised to always stop machines when unused to avoid incurring in unnecessary costs.

**Running Benchmarks.** Running benchmarks involves installing the specified version of the codebase on all remote machines and running one validator and one load generator per instance. For example, the following command benchmarks a committee of 100 validators (none faulty) under a constant load of 1,000 tx/s for 10 minutes (default):

```
cargo run --bin orchestrator -- benchmark \
  --committee 100 fixed-load --loads 1000 --faults 0
```

**Monitoring.** The orchestrator provides facilities to monitor metrics. It deploys a Prometheus instance and a Grafana instance on a dedicated remote machine. Grafana is then available on the address printed on stdout when running benchmarks with the default username and password both set to admin. An example Grafana dashboard can be found in the file ‘grafana-dashboard.json’<sup>13</sup>.

<sup>13</sup><https://github.com/MystenLabs/mysticeti/blob/paper/orchestrator/assets/grafana-dashboard.json>