

Catchain Consensus: An Outline

modified by TOSDAO, designed by Telegram

April 24, 2023

Abstract

The aim of this text is to provide an outline of the Catchain Consensus Protocol, a Byzantine Fault Tolerant (BFT) protocol specifically crafted for block generation and validation in the TOS Blockchain [3]. This protocol can potentially be used for purposes other than block generation in a proof-of-stake (PoS) blockchain; however, the current implementation uses some optimizations valid only for this specific problem.

Contents

1	Overview	2
2	Catchain Protocol	5
3	Block Consensus Protocol	16

1 Overview

The Catchain Consensus protocol builds upon the overlay network construction protocol and the overlay network broadcast protocol of TOS Network ([3]). The Catchain Consensus protocol itself can be decomposed into two separate protocols, one more low-level and general-purpose (the *Catchain protocol*¹), and the other the high-level *Block Consensus Protocol (BCP)*, which makes use of the Catchain protocol. Higher levels in the TOS protocol stack are occupied by the block generation and validation levels; however, all of them are executed essentially locally on one (logical) machine, with the problem of achieving consensus on the newly-generated block delegated to the Catchain protocol level.

Here is an approximate diagram of the protocol stack employed by TOS for block generation and distribution, showing the correct place of the Catchain Consensus protocol (or rather its two component protocols):

- *Top-level*: Block generation and block validation software, logically running on a stand-alone logical machine, with all the inputs provided and outputs handled by the lower-level protocols. The job of this software is to either generate a new valid block for a blockchain (a shardchain or the masterchain of the TOS Blockchain; cf. [3] for a discussion of shardchains and the masterchain), or to check the validity of a block generated by somebody else.
- *(TOS) Block consensus protocol*: Achieves (byzantine fault tolerant) consensus on the block to be accepted as the next one in the current validator group for the masterchain or a shardchain. This level makes use of (the abstract interface of) the block generation and validation software, and builds upon the lower-level Catchain protocol. This protocol is explained in more detail in Section 3.
- *Catchain protocol*: Provides secure persistent broadcasts in an overlay network (e.g., the task group of validators for a specific shardchain or the masterchain dedicated to generation, validation, and propagation of new blocks in this shardchain or masterchain), and detects attempts

¹The original name of this protocol used during the initial stage of the research and development phase was *catch-chain* or *catchchain*, because it essentially is a special *blockchain* dedicated to *catching* all events important for the consensus protocol; after saying and writing this name a lot of times it gradually got contracted to “catchain”.

of “cheating” (protocol violation) on the part of some participants. This protocol is explained in more detail in Section 2.

- *(TOS Network) overlay broadcast protocol*: A simple best-effort broadcast protocol for overlay networks in the TOS Network as described in [3]. Simply broadcasts received broadcast messages to all neighbors in the same overlay network that did not receive a copy of these messages before, with minimal effort dedicated to keeping copies of undelivered broadcast messages for a short period of time.
- *(TOS Network) overlay protocol*: Creates overlay networks (cf. [3]) inside the ADNL protocol network, manages neighbor lists for these overlay networks. Each participant of an overlay network tracks several neighbors in the same overlay network and keeps dedicated ADNL connections (called “channels”) to them, so that incoming messages can be efficiently broadcast to all neighbors with minimal overhead.
- *Abstract Datagram Network Layer (ADNL) protocol*: The basic protocol of the TOS Network, that delivers packets (datagrams) between network nodes identified only by 256-bit abstract (ADNL) addresses, which effectively are cryptographic keys (or their hashes).

This text aims to describe only the second and the third protocol in this suite, namely, the (TOS) block consensus protocol and the (TOS) Catchain protocol.

We would like to point out here that the author of this text, while providing the general guidelines of how this protocol should be designed (on the lines of “let’s create a BFT-hardened group broadcast message system, and run a suitably adapted simple two-phase or three-phase commit protocol on top of this system”) and participating in several discussions during the development and implementation of the protocol, is definitely not the only designer of this protocol and especially of its current implementation. This is the work of several people.

A few words on the efficiency of the combined Catchain Consensus protocol. Firstly, it is a true Byzantine Fault Tolerant (BFT) protocol, in the sense that it eventually achieves consensus on a valid next block of the blockchain even if some participants (validators) exhibit arbitrarily malicious behavior, provided these malicious participants are less than one third of the total number of the validators. It is well-known that achieving BFT consensus

is impossible if at least one third of participants are malicious (cf. [5]), so the Catchain Consensus protocol is as good as theoretically possible in this respect. Secondly, when the Catchain Consensus was first implemented (in December 2018) and tested on up to 300 nodes distributed all over the world, it achieved consensus on a new block in 6 seconds for 300 nodes and in 4–5 seconds for 100 nodes (and in 3 seconds for 10 nodes), even if some of these nodes fail to participate or exhibit incorrect behavior.² Since the TOS Blockchain task groups are not expected to consist of more than a hundred validators (even if a total of a thousand or ten thousand validators are running, only a hundred of them with the largest stakes will generate new masterchain blocks, and the others will participate only in the creation of new shardchain blocks, each shardchain block generated and validated by 10–30 validators; of course, all numbers given here are configuration parameters (cf. [3] and [4]) and can be adjusted later by a consensus vote of validators if necessary), this means that the TOS Blockchain is able to generate new blocks once every 4–5 seconds, as originally planned. This promise has been further tested and found out to be fulfilled with the launch of the Test Network of the TOS Blockchain a couple of months later (in March 2019). Therefore, we see that the Catchain Consensus protocol is a new member of the ever-growing family of practical BFT protocols (cf. [2]), even though it is based on slightly different principles.

²When the ratio of the malicious or non-participating or very slow validators grows up to one third, the protocol exhibits graceful degradation, with the block consensus time growing very slowly—say, by at most half a second—until the critical value of one third is almost achieved.

2 Catchain Protocol

We have already explained in the Overview (cf. 1) that the BFT consensus protocol used by the TOS Blockchain for achieving consensus on new blockchain blocks consists of two protocols. We provide here a brief description of the *Catchain protocol*, the lower-level of these two protocols that could be potentially used for purposes other than BFT consensus for blocks. The source code for the Catchain protocol resides in subdirectory `catchain` of the source tree.

2.1. Prerequisites for running the Catchain protocol. The main prerequisite for running (an instance of) the Catchain protocol is the ordered list of all nodes that are participating (or allowed to participate) in this specific instance of the protocol. This list consists of public keys and ADNL addresses of all participating nodes. It has to be provided from the outside when an instance of the Catchain protocol is created.

2.2. Nodes participating in the block consensus protocol. For the specific task of creating new blocks for one of the blockchains (i.e., the masterchain or one of the active shardchains) of the TOS Blockchain, a special task group consisting of several validators is created. The list of members of this task group is used both to create a private overlay network inside ADNL (this means that the only nodes that can join this overlay network are those explicitly listed during its creation) and to run the corresponding instance of the Catchain protocol.

The construction of this list of members is the responsibility of the higher levels of the overall protocol stack (the block creation and validation software) and therefore is not the topic of this text ([4] would be a more appropriate reference). It is sufficient to know at this point that this list is a deterministic function of the current (most recent) masterchain state (and especially of the current value of the configuration parameters, such as the active list of all validators elected for creating new blocks along with their respective weights). Since the list is computed deterministically, all validators compute the same lists and in particular each validator knows in which task groups (i.e., instances of the Catchain protocol) it participates without any further need for network communication or negotiation.³

³If some validators have an outdated masterchain state, they may fail to compute correct task group lists and to participate in the corresponding catchains; in this respect,

2.2.1. Catchains are created in advance. In fact, not only the current values of the lists alluded to above are computed, but also their immediately subsequent (future) values are computed as well, so that the Catchain is usually created in advance. In this way it is already in place when the first block has to be created by the new instance of the validator task group.

2.3. The genesis block and the identifier of a catchain. A *catchain* (i.e., an instance of the Catchain protocol) is characterized by its *genesis block* or *genesis message*. It is a simple data structure containing some magic numbers, the purpose of the catchain (e.g., the identifier of the shardchain, for which the blocks will be generated, and the so-called *catchain sequence number*, also obtained from the masterchain configuration and used to distinguish subsequent instances of the catchain generating “the same” shardchain, but possibly with different participating validators), and, most importantly, the list of all participating nodes (their ADNL addresses and Ed25519 public keys as explained in **2.1**). The Catchain protocol itself uses only this list and the SHA256 hash of the overall data structure; this hash is used as an internal identifier of the catchain, i.e., of this specific instance of the Catchain protocol.

2.3.1. Distribution of the genesis block. Note that the genesis block is not distributed among the participating nodes; it is rather computed independently by each participating node as explained in **2.2**. Since the hash of the genesis block is used as the catchain identifier (i.e., identifier of the specific instance of the Catchain protocol; cf. **2.3**), if a node (accidentally or intentionally) computes a different genesis block, it will be effectively locked out from participating in the “correct” instance of the protocol.

2.3.2. List of nodes participating in a catchain. Note that the (ordered) list of nodes participating in a catchain is fixed in the genesis block and hence it is known to all the participants and it is unambiguously determined by the hash of the genesis block (i.e., the catchain identifier), provided there are no (known) collisions for SHA256. Therefore, we fix the number of participating nodes N in the discussion of one specific catchain below, and assume that the nodes are numbered from 1 to N (their real identities may be looked up

they are treated as if they were malicious or malfunctioning and do not affect the overall validity of the BFT protocol as long as less than one third of all validators fail in this fashion.

in the list of participants using this index in range $1 \dots N$). The set of all participants will be denoted by I ; we assume that $I = \{1 \dots N\}$.

2.4. Messages in a catchain. Catchain as a process group. One perspective is that a catchain is a (*distributed*) *process group* consisting of N known and fixed (*communicating*) *processes* (or *nodes* in the preceding terminology), and these processes generate *broadcast messages*, that are eventually broadcast to all members of the process group. The set of all processes is denoted by I ; we usually assume that $I = \{1 \dots N\}$. The broadcasts generated by each process are numbered starting from one, so the n -th broadcast of process i will receive *sequence number* or *height* n ; each broadcast should be uniquely determined by the identity or the index i of the originating process and its height n , so we can think of the pair (i, n) as the natural identifier of a broadcast message inside a process group.⁴ The broadcasts generated by the same process i are expected to be delivered to every other process in exactly the same order they have been created, i.e., in increasing order of their height. In this respect a catchain is very similar to a process group in the sense of [1] or [7]. The principal difference is that a catchain is a “hardened” version of a process group tolerant to possible Byzantine (arbitrarily malicious) behavior of some participants.

2.4.1. Dependence relation on messages. One can introduce a *dependence relation* on all messages broadcast in a process group. This relation must be a strict partial order \prec , with the property that $m_{i,k} \prec m_{i,k+1}$, where $m_{i,k}$ denotes the k -th message broadcast by group member process with index i . The meaning of $m \prec m'$ is that m' *depends on* m , so that the (broadcast) message m' can be processed (by a member of the process group) only if m has been processed before. For instance, if the message m' represents the reaction of a group member to another message m , then it is natural to set $m \prec m'$. If a member of the process group receives a message m' before all its dependencies, i.e., messages $m \prec m'$, have been processed (or *delivered* to the higher-level protocol), then its processing (or *delivery*) is delayed until all its dependencies are delivered.

We have defined the dependence relation to be a strict partial order, so it must be transitive ($m'' \prec m'$ and $m' \prec m$ imply $m'' \prec m$), antisymmetric (at most one of $m' \prec m$ and $m \prec m'$ can hold for any two messages m and m') and anti-reflexive ($m \prec m$ never holds). If we have a smaller set

⁴In the Byzantine environment of a catchain this is not necessarily true in all situations.

of “basic dependencies” $m' \rightarrow m$, we can construct its transitive closure \rightarrow^+ and put $\prec := \rightarrow^+$. The only other requirement is that every broadcast of a sender depends on all previous broadcasts of the same sender. It is not strictly necessary to assume this; however, this assumption is quite natural and considerably simplifies the design of a messaging system inside a process group, so the Catchain protocol makes this assumption.

2.4.2. Dependence set or cone of a message. Let m be a (broadcast) message inside a process group as above. We say that the set $D_m := \{m' : m' \prec m\}$ is the *dependence set* or *dependence cone* of message m . In other words, D_m is the *principal ideal* generated by m in the partially ordered finite set of all messages. It is precisely the set of all messages that must be delivered before m is delivered.

2.4.3. Extended dependence cone of a message. We also define D_m^+ , the *extended dependence cone* of m , by $D_m^+ := D_m \cup \{m\}$.

2.4.4. Cones, or ideals with respect to \prec . More generally, we say that a subset D of messages is a *cone* if it is an ideal with respect to the dependence relation \prec , i.e., if $m \in D$ and $m' \prec m$ imply $m' \in D$. Of course, the dependence cone D_m and the extended dependence cone D_m^+ of any message m are cones (because any principal ideal in a partially ordered set is an ideal).

2.4.5. Identification of cones with the aid of vector time. Recall that we have assumed that any message depends on all preceding messages of the same sender, i.e. $m_{i,s} \prec m_{i,s+1}$ for any $i \in I$ and any $s > 0$, such that $m_{i,s+1}$ exists. This implies that any cone D is completely characterized by N values $\text{VT}(D)_i$ indexed by $i \in I$:

$$\text{VT}(D)_i := \sup\{s \in \mathbb{N} : m_{i,s} \in D\} = \inf\{s \in \mathbb{N}_0 : m_{i,s+1} \notin D\} \quad (1)$$

(if no message $m_{i,s}$ is in D , we set $\text{VT}(D)_i := 0$). Indeed, it is clear that

$$m_{i,s} \in D \Leftrightarrow s \leq \text{VT}(D)_i \quad (2)$$

We say that the vector $\text{VT}(D) = (\text{VT}(D)_i)_{i \in I} \in \mathbb{N}_0^I$ with non-negative components $\text{VT}(D)_i$ is the *vector time* or *vector timestamp* corresponding to cone D (cf. [1] or [7] for a more detailed discussion of vector time).

2.4.6. Partial order on vector timestamps. We introduce a partial order \leq on the set of all possible vector times \mathbb{N}_0^I , which is the product of the usual orders on \mathbb{N}_0 :

$$\mathbf{x} = (x_i)_{i \in I} \leq \mathbf{y} = (y_i)_{i \in I} \quad \text{iff} \quad x_i \leq y_i \quad \text{for all } i \in I \quad (3)$$

It is immediate that $D \subset D'$ iff $\text{VT}(D) \leq \text{VT}(D')$; therefore, VT is a strict order-preserving embedding of the set of all cones contained in the set of all messages into \mathbb{N}_0^I .

2.4.7. Vector timestamp $\text{VT}(m)$ of a message m . Given any message m , we define its *vector timestamp* $\text{VT}(m)$ as $\text{VT}(D_m)$. In other words, message m can be delivered only after the first $\text{VT}(m)_j$ messages generated by process j are delivered, and this is true for all $j \in I$.

If i is the sender of message m , and s is the height of message m , so that $m = m_{i,s}$, then $\text{VT}(m)_i = s - 1$. We can define the *adjusted vector timestamp* $\text{VT}^+(m)$ of message m by setting $\text{VT}^+(m)_j = \text{VT}(m)_j$ for $j \neq i$, $\text{VT}^+(m)_i = \text{VT}(m)_i + 1 = s$. Alternatively, $\text{VT}^+(m) = \text{VT}(D_m^+)$, where $D_m^+ := D_m \cup \{m\}$ is the *extended dependence cone* of m (cf. **2.4.3**).

Note that $m' \preceq m$ iff $D_{m'}^+ \subset D_m^+$ iff $\text{VT}^+(m') \leq \text{VT}^+(m)$ in \mathbb{N}_0^I , where $m' \preceq m$ means “ $m' \prec m$ or $m' = m$ ”. Similarly, $m' \prec m$ iff $D_{m'}^+ \subset D_m$ iff $\text{VT}^+(m') \leq \text{VT}(m)$. In other words, *the dependence relation \prec on (some or all) messages is completely determined by the adjusted vector timestamps of these messages.*

2.4.8. Using vector timestamps to correctly deliver broadcast messages. Vector timestamps can be used (in non-byzantine settings) to correctly deliver messages broadcast in a process group.⁵ Namely, suppose that every broadcast message $m = m_{i,s}$ contains the index of its sender i and the vector timestamp of this message $\text{VT}(m)$. Then each receiver j knows whether the message can be delivered or not. For this, j keeps track of the cone C_j of all messages delivered so far, for example by maintaining a *current timestamp* $\text{VT}(j)$ equal to $\text{VT}(C_j)$. In other words, $\text{VT}(j)_k$ is the count of messages of sender k processed by j so far. If $\text{VT}(m) \leq \text{VT}(j)$, then the message m is delivered immediately and $\text{VT}(j)$ is updated to $\text{sup}(\text{VT}(j), \text{VT}^+(m))$ afterwards; this is equivalent to increasing $\text{VT}(j)_i$ by

⁵We assume that all broadcast messages in the process group are “causal broadcasts” or “cbcast” in the terminology of [1], because we only need cbcasts for the implementation of Catchain protocol and Catchain consensus.

one, where i is the original sender of message m . If this condition is not met, then m may be put into a waiting queue until $\text{VT}(j)$ becomes large enough. Instead of passively waiting for the required broadcasts, j can construct the list of message indices (i', s') that are implicitly mentioned in $\text{VT}(m)$ of some received but not delivered message m , and request messages with these indices from the neighbors from which j learned about m and $\text{VT}(m)$; an alternative strategy (actually employed by the current implementation of the Catchain protocol) is to request these messages from randomly chosen neighbors from time to time. The latter strategy is simpler because it does not require remembering the immediate sources of all received messages (which may become unavailable anyway).

2.5. Message structure in a catchain. Catchain as a multi-blockchain.

The message structure in a catchain is a bit more complicated than described above because of the necessity to support a BFT protocol. In particular, vector timestamps are not sufficient in a Byzantine setting. They have to be complemented by descriptions based on maximal elements of a dependence cone (such descriptions are typically used in non-byzantine settings only when the process group is very large, so that vector timestamp sizes become prohibitive).

2.5.1. Describing cones by means of their maximal elements. An alternative way (to using a vector timestamp) of describing a message cone D is by listing all its *maximal elements* $\text{Max}(D)$, i.e. elements $m \in D$, such that $m \prec m'$ does not hold for any $m' \in D$. Of course, one needs a suitable way of referring to messages without including them completely in order for this representation to be practical.

2.5.2. Message identifiers inside a catchain. Catchain protocol uses SHA256 hashes of (suitably serialized) messages as their unique identifiers. If we assume that there are no collisions for SHA256 (computable in reasonable, e.g., polynomial time), then a message m is completely identified within the process group by its hash $\text{SHA256}(m)$.

2.5.3. Message headers. The header of a message $m = m_{i,s}$ inside a catchain (i.e., an instance of the Catchain protocol) always contains the index i of its sender, the height s , the catchain identifier (i.e., the hash of the genesis message, cf. 2.3) and the set of hashes of maximal elements of the dependence cone of m , i.e., the set $\{\text{SHA256}(m') : m' \in \text{Max}(D_m)\}$. In particular, the

hash $\text{SHA256}(m_{i,s-1})$ of the previous message of the same sender is always included since $m_{i,s-1} \in \text{Max}(D_m)$ if $s > 1$; for performance reasons, there is a separate field in the message header containing $\text{SHA256}(m_{i,s-1})$. If $s = 1$, then there is no previous message, so the hash of the genesis message (i.e., the catchain identifier, cf. **2.3**) is used instead.

The vector timestamp $\text{VT}(m)$ is not included in the message header; however, the header implicitly determines $\text{VT}(m)$ since

$$\text{VT}(m) = \sup_{m' \in D_m} \text{VT}^+(m') = \sup_{m' \in \text{Max}(D_m)} \text{VT}^+(m') \quad (4)$$

Note that the message header is a part of the message, and in particular the hash of a message (i.e., the message identifier) depends on all data listed in the header. Therefore, we assume that the message identifier implicitly determines all the dependencies of the corresponding message (if there are no known collisions for SHA256).

2.5.4. Message signatures. Apart from that, every message in a catchain is signed by its creator. Since the list of participating nodes (processes) in a catchain is known in advance, and this list includes the public keys of all processes, these message signatures can be checked by a receiving process immediately after a message is received. If the signature is invalid, the message is discarded without any further processing.

2.5.5. Message encryption. All messages in a catchain are also encrypted before being transferred from a node to its neighbor in the private overlay network underlying the catchain. However, this encryption is performed by lower-level network protocols (such as ADNL) and is not relevant to the discussion here. We would like to mention that correct encryption is possible here only because the list of participating processes includes not only the public keys of all processes, but also their ADNL addresses (which effectively are public encryption keys for network transmission).

Notice that even if the encryption had been absent, this would not violate the BFT properties of the protocol, because faking a message from another sender would not be possible because of the signatures. However, this might lead to a leak of information to outside observers, which is often undesirable.

2.5.6. Alternative perspective: a catchain as a multi-blockchain. Note that all messages created by the same sender i in a catchain turn out to have a simple “blockchain structure”, because the header of $m_{i,s+1}$ contains

the hash $\text{SHA256}(m_{i,s})$ (among other hashes of messages from $\text{Max}(D_{m_{i,s+1}})$) of the previous message of sender i . In this way each process i generates a simple blockchain consisting of its messages, with each “block” of this blockchain corresponding to one message and referring to the previous block by its hash, and sometimes includes references to blocks (i.e., messages) of other processes by mentioning the hashes of these blocks in its blocks. Each block is signed by its creator. The resulting structure is very similar to that of an “asynchronous payment channel” considered in [3, 5], but with N participants instead of 2.

2.6. Message propagation in a catchain. Now we are ready to describe message propagation in a catchain. Namely:

- The (lower-level) overlay network protocol maintains a list of neighbors in the private overlay network underlying the catchain and provides ADNL channels to each of these neighbors. This private overlay network has the same list of members (processes, nodes) as the catchain, and the neighbors of each node form an (oriented) subgraph on the set of all participating nodes. This (essentially random) subgraph is strongly connected with probability very close to one.
- Each process generates some new messages from time to time (as needed by the higher-level protocol). These messages are augmented by catchain message headers as outlined in **2.5.3**, signed, and propagated to all known neighbors using the ADNL channels established by the overlay protocol.
- In contrast with the usual simple overlay broadcast protocol, the messages received from neighbors are not immediately rebroadcast to all other neighbors that are not known yet to have a copy of them. Instead, the signature is checked first, and invalid messages are discarded. Then the message is either delivered (if all its dependent messages have already been delivered), or put into a waiting queue. In the latter case, all the required messages mentioned in its header (i.e., the set $\text{Max}(D_m)$) are pulled from the neighbor that sent this message (apart from that, attempts to download these missing messages from random neighbors are performed from time to time). If necessary, this process is repeated recursively until some messages can be delivered. Once a message is ready for local delivery (i.e., all its dependencies are already present), it is also rebroadcast to all neighbors in the overlay network.

- Apart from the recursive “pull” mechanism described above, a faster vector timestamp-based mechanism is also used, so that messages can be queried from neighbors by their senders and heights (learned from the vector timestamps of received messages). Namely, each process sends a special query containing the current vector timestamp to a randomly chosen neighbor from time to time. This peer-to-peer query leads to its receiver sending back all or some messages unknown to the sender (judging by their vector timestamps).
- This faster vector timestamp-based mechanism can be disabled for messages originating from certain senders as soon as a “fork” is detected, i.e., a second message with the same sender i and height s , but with a different hash, is learned from a neighbor, for example, during the fast or slow “pull” process. Once a fork created by i is detected, the corresponding component V_{T_i} of all subsequent vector timestamps is set to a special value ∞ to indicate that comparing the values of these components does not make sense anymore.
- When a message is delivered (to the higher-level protocol), this message is added into the cone C of processed messages of the current process (and the current vector timestamp is updated accordingly), and all subsequent messages generated by the current process will be assumed to depend on all the messages delivered so far (even if this is not logically necessary from the perspective of the higher-level protocol).
- If the set $\text{Max}(C)$ of the maximal elements of the cone of processed messages becomes too large (contains more elements than a certain amount fixed in advance by the genesis message of the catchain), then the Catchain protocol asks the higher-level protocol to generate a new message (empty if no useful payload is available). After this new message is generated (and immediately delivered to the current process), C is updated and $\text{Max}(C)$ consists of only one element (the new message). In this way the size of $\text{Max}(C)$ and therefore the size of the message header always remain bounded.
- Once a message m is delivered and the set C is modified to include this message, a timer is set, and after some small delay the higher-level protocol is asked to create a new message (empty if necessary), so that this new message m^* would refer to the new C , similarly to

the procedure described in the previous item. This new message m^* is pushed to all neighbors; since its header contains $\text{Max}(C)$ for the new C , and $m \in C$, the neighbors learn not only about the newly-generated message m^* , but also about the original received message m . If some neighbors do not have a copy of m yet, they would require one (from the current process or not).

- All (broadcast) messages received and created in a catchain are stored into a special local database. This is especially important for newly-created messages (cf. **3.3.2**): if a message is created and sent to neighbors, but not saved into the database (and flushed to disk) before the creating process crashes and is restarted, then another message with the same sender and height can be created after restart, thus effectively leading to an involuntary “fork”.

2.7. Forks and their prevention. One can see that the multi-blockchain structure of a catchain outlined above (with references to other blocks by their hashes and with signatures) leaves very little possibility for “cheating” in a consensus protocol built upon a catchain (i.e., using the catchain as a means for broadcasting messages inside a process group). The only possibility that is not detected immediately consists of creating two (or more) different versions of the same message $m_{i,s}$ (say, $m'_{i,s}$ and $m''_{i,s}$), and sending one version of this message $m'_{i,s}$ to some peers and a different version $m''_{i,s}$ to others. If s is minimal (for a fixed i), then this corresponds to a *fork* in blockchain terminology: two different next blocks $m'_{i,s}$ and $m''_{i,s}$ for the same previous block $m_{i,s-1}$.

Therefore, the Catchain protocol takes care to detect forks as soon as possible and prevent their propagation.

2.7.1. Detection of forks. The detection of forks is simple: if there are two different blocks $m'_{i,s}$ and $m''_{i,s}$ with the same creator $i \in I$ and the same height $s \geq 1$, and with valid signatures of i , then this is a fork.

2.7.2. Fork proofs. Block signatures in the Catchain protocol are created in such a way that creating *fork proofs* (i.e., the proof that a process i has intentionally created a fork) is especially simple since it is the hash of a very small structure (containing a magic number, the values of i and s , and the hash of the remainder of the message) that is actually signed. Therefore, only two such small structures and two signatures are required in a fork proof.

2.7.3. External punishment for creating forks. Notice that an external punishment for creating catchain forks may be used in the proof-of-stake blockchain generation context. Namely, the fork proofs may be submitted to a special smart contract (such as the elector smart contract of the TOS Blockchain), checked automatically, and some part or all of the stake of the offending party may be confiscated.

2.7.4. Internal processing of forks. Once a fork (created by i) is detected (by another process j), i.e. j learns about two different messages $m_{i,s}$ and $m'_{i,s}$ created by i and having same height s (usually this happens while recursively downloading dependencies of some other messages), j starts ignoring i and all of its subsequent messages. They are not accepted and not broadcast further. However, messages created by i prior to the fork detection may be still downloaded if they are referred to in messages (blocks) created by processes that did not see this fork before referring to such messages created by i .

2.7.5. Accepting messages from a “bad” process is bad. Furthermore, if process i learns about a fork created by process j , then i shows this to its neighbors by creating a new service broadcast message that contains the corresponding fork proof (cf. **2.7.2**). Afterwards, this and all subsequent messages of j cannot directly depend on any messages by the known “bad” producer i (but they still can refer to messages from another party k that directly or indirectly refer to messages of i if no fork by i was known to k at the time when the referring message was created). If j violates this restriction and creates messages with such invalid references, these messages will be discarded by all honest processes in the group.

2.7.6. The set of “bad” group members is a part of the intrinsic state. Each process i keeps its own copy of the set of known “bad” processes in the group, i.e., those processes that have created at least one fork or have violated **2.7.5**. This set is updated by adding j into it as soon as i learns about a fork created by j (or about a violation of **2.7.5** by j); after that, a callback provided by the higher-level protocol is invoked. This set is used when a new broadcast message arrives: if the sender is bad, then the message is ignored and discarded.

3 Block Consensus Protocol

We explain in this section the basic workings of the TOS Block Consensus Protocol (cf. **1**), which builds upon the generic Catchain protocol (cf. **2**) to provide the BFT protocol employed for generating and validating new blocks of the TOS Blockchain. The source code for the TOS Block Consensus protocol resides in subdirectory `validator-session` of the source tree.

3.1. Internal state of the Block Consensus Protocol. The higher-level Block Consensus Protocol introduces a new notion to the catchain: that of an *internal state* of the Block Consensus Protocol (BCP), sometimes also (not quite correctly) called “the internal state of the catchain” or simply *catchain state*. Namely, each process $i \in I$ has a well-determined internal state σ_{C_i} after a subset of messages (actually always a dependence cone) C_i is delivered by the Catchain protocol to the higher-level protocol (i.e., to the Block Consensus Protocol in this case). Furthermore, this state $\sigma_{C_i} = \sigma(C_i)$ depends only on cone C_i , but not on the identity of the process $i \in I$, and can be defined for any dependence cone S (not necessarily a cone C_i of delivered messages for some process i at some point).

3.1.1. Abstract structure of the internal state. We start with an abstract structure of the internal state employed by BCP; more specific details will be provided later.

3.1.2. Updating the internal state. The Catchain protocol knows nothing about the internal state; it simply invokes appropriate callbacks supplied by the higher-level protocol (i.e., the BCP) whenever a message m is delivered. It is the job of the higher-level protocol to compute the new state $\sigma_{S'}$ starting from the previously computed state σ_S and the message m , where $S' = S \cup \{m\}$ (and necessarily $S \supset D_m$, otherwise m could not have been delivered at this point).

3.1.3. Recursive formula for updating the internal state. The abstract setup for computing σ_S for all cones S consists of three components:

- A value σ_\emptyset for the initial state (this value actually depends on the genesis block of the catchain; we ignore this dependence here because we consider only one catchain at this point).

- A function f that computes the state $\sigma_{D_m^+}$ from the previous state σ_{D_m} and the newly-delivered message m :

$$\sigma_{D_m^+} = f(\sigma_{D_m}, m) \quad (5)$$

where D_m is the dependence cone of message m and $D_m^+ = D_m \cup \{m\}$ its extended dependence cone (cf. **2.4.3**). In most cases, f will actually satisfy the stronger condition

$$\sigma_{S \cup \{m\}} = f(\sigma_S, m) \quad \text{if } S \text{ and } S \cup \{m\} \text{ are cones and } m \notin S \quad (6)$$

However, this stronger condition is not required by the update algorithm.

- A “merge function” g that computes $\sigma_{S \cup T}$ from σ_S and σ_T :

$$\sigma_{S \cup T} = g(\sigma_S, \sigma_T) \quad \text{for any cones } S \text{ and } T \quad (7)$$

(the union of two cones always is a cone). This function g is applied by the update algorithm only in the specific case $T = D_m^+$ and $m \notin S$.

3.1.4. Commutativity and associativity of g . Note that (7) (for arbitrary cones S and T) implies associativity and commutativity of g , at least when g is applied to possible states (values of form σ_S for some cone S). In this respect g defines a commutative monoid structure on the set $\Sigma = \{\sigma_S : S \text{ is a cone}\}$. Usually g is defined or partially defined on a larger set $\tilde{\Sigma}$ of state-like values, and it may be commutative and associative on this larger set $\tilde{\Sigma}$, i.e., $g(x, y) = g(y, x)$ and $g(x, g(y, z)) = g(g(x, y), z)$ for $x, y, z \in \tilde{\Sigma}$ (whenever both sides of the equality are defined), with σ_\emptyset as an unit, i.e., $g(x, \sigma_\emptyset) = x = g(\sigma_\emptyset, x)$ for $x \in \tilde{\Sigma}$ (under the same condition). However, this property, useful for the formal analysis of the consensus algorithm, is not strictly required by the state update algorithm, because this algorithm uses g in a deterministic fashion to compute σ_S .

3.1.5. Commutativity of f . Note that f , if it satisfies the stronger condition (6), must also exhibit a commutativity property

$$f(f(\sigma_S, m), m') = f(f(\sigma_S, m'), m) \quad (8)$$

whenever S is a cone and m and m' are two messages with $D_m \subset S$, $D_{m'} \subset S$, $m \notin S$ and $m' \notin S$, because in this case $S \cup \{m\}$, $S \cup \{m'\}$ and $S \cup \{m, m'\}$

are also cones, and (6) implies that both sides of (8) are equal to $\sigma_{S \cup \{m, m'\}}$. Similarly to **3.1.4**, f is usually defined or partially defined on the product of a larger set $\tilde{\Sigma}$ of state-like values and of a set of message-like values; it may exhibit the “commutativity” property (8) or not on this larger set. If it does, this might be useful for formal analysis of the algorithms relying on σ_S , but this property is not strictly necessary.

3.1.6. The state update algorithm. The state update algorithm (independently executed by each process i) employed by the catchain (actually by the higher-level BCP) uses σ_\emptyset , f and g as follows:

- The algorithm keeps track of all $\sigma_{D_m^+}$ for all messages m delivered so far.
- The algorithm keeps track of σ_{C_i} , where C_i is the current dependence cone, i.e., the set of all messages m delivered (to the current process i). The initial value of σ_{C_i} is σ_\emptyset .
- When a new message m is delivered, the value of σ_{D_m} is computed by a repeated application of g since $D_m = \bigcup_{m' \in D_m} D_{m'}^+ = \bigcup_{m' \in \text{Max}(D_m)} D_{m'}^+$; therefore, if $\text{Max}(D_m) = \{m'_1, \dots, m'_k\}$, then

$$\sigma_{D_m} = g\left(\dots g\left(g\left(\sigma_{D_{m'_1}^+}, \sigma_{D_{m'_2}^+}\right), \sigma_{D_{m'_3}^+}\right), \dots \sigma_{D_{m'_k}^+}\right) \quad . \quad (9)$$

The set $\text{Max}(D_m)$ is explicitly listed in the header of message m in some fixed order m'_1, \dots, m'_k ; the above formula is applied with respect to this order (so the computation of D_m is deterministic). The first element in this list always is the previous message of the sender of m , i.e., if $m = m_{i,s+1}$, then $m'_1 = m_{i,s}$.

- After that, the value of $\sigma_{D_m^+}$ is computed by an application of f : $\sigma_{D_m^+} = f(\sigma_{D_m}, m)$. This value is memorized for future use.
- Finally, when a new message m is delivered to the current process i , thus updating C_i to $C'_i := C_i \cup \{m\}$, the algorithm uses the computed value $\sigma_{D_m^+}$ to update the current state

$$\sigma_{C'_i} = g(\sigma_{C_i}, \sigma_{D_m^+}) \quad (10)$$

This state, however, is “virtual” in the sense that it can be slightly changed later (especially if g is not commutative). Nevertheless, it is

used to make some important decisions by the higher-level algorithm (BCP).

- Once a new message m is generated and locally delivered, so that C_i becomes equal to D_m^+ , the previously computed value of σ_{C_i} is discarded and replaced with $\sigma_{D_m^+}$ computed according to the general algorithm described above. If g is not commutative or not associative (for example, it may happen that $g(x, y)$ and $g(y, x)$ are different but equivalent representations of the same state), then this might lead to a slight change of the current “virtual” state of process i .
- If the lower-level (catchain) protocol reports to the higher-level protocol that a certain process $j \notin i$ is “bad” (i.e., j is found out to have created a fork, cf. 2.7.6, or to have knowingly endorsed a fork by another process, cf. 2.7.5), then the current (virtual) state σ_{C_i} is recomputed from scratch using the new set $C'_i = \bigcup_{m \in C_i, m \text{ was created by "good" process } k} D_m^+$ and the “merge” function g applied to the set of $\sigma_{D_m^+}$ where m runs through the set of last messages of the processes known to be good (or through the set of maximal elements of this set). The next created outbound message will depend only on the messages from C'_i .

3.1.7. Necessity to know the internal state of the other processes.

Formula (9) implies that process i must also keep track of $\sigma_{D_m^+}$ for all messages m , created by this process or not. However, this is possible since these internal states are also computed by appropriate applications of the update algorithm. Therefore, BCP computes and remembers all $\sigma_{D_m^+}$ as well.

3.1.8. Function f would suffice. Notice that the update algorithm applies g only to compute $\sigma_{S \cup D_m^+} = g(\sigma_S, \sigma_{D_m^+})$ when S is a cone containing D_m , but not containing m . Therefore, every actual application of g could have been replaced by an application of f satisfying the extended property (6):

$$\sigma_{S \cup D_m^+} = g(\sigma_S, \sigma_{D_m^+}) = f(\sigma_S, m) \tag{11}$$

However, the update algorithm does not use this “optimization”, because it would disable the more important optimizations described below in 3.2.4 and 3.2.5.

3.2. The structure of the internal state. The structure of the internal state is optimized to make the *transition function* f of (5) and the *merge*

function g of (7) as efficiently computable as possible, preferably without the need of potentially unbounded recursion (just some loops). This motivates the inclusion of additional components into the internal state (even if these components are computable from the remainder of the internal state), which have to be stored and updated as well. This process of including additional components is similar to that employed while solving problems using dynamic programming, or to that used while proving statements by mathematical (or structural) induction.

3.2.1. The internal state is a representation of a value of an abstract algebraic data type. The internal representation of the internal state is essentially a (directed) tree (or rather a directed acyclic graph) or a collection of nodes; each node contains some immediate (usually integer) values and several pointers to other (previously constructed) nodes. If necessary, an extra *constructor tag* (a small integer) is added at the beginning of a node to distinguish between several possibilities. This structure is very similar to that used to represent values of abstract algebraic data types in functional programming languages such as Haskell.

3.2.2. The internal state is persistent. The internal state is *persistent*, in the sense that the memory used to allocate the nodes which are part of the internal state is never freed up while the catchain is active. Furthermore, the internal state of a catchain is actually allocated inside a huge contiguous memory buffer, and new nodes are always allocated at the end of the used portion of this buffer by advancing a pointer. In this way the references to other nodes from a node inside this buffer may be represented by an integer offset from the start of the buffer. Every internal state is represented by a pointer to its root node inside this buffer; this pointer can be also represented by an integer offset from the start of the buffer.

3.2.3. The internal state of a catchain is flushed to an append-only file. The consequence of the structure of the buffer used to store the internal states of a catchain explained above is that it is updated only by appending some new data at its end. This means that the internal state (or rather the buffer containing all the required internal states) of a catchain can be flushed to an append-only file, and easily recovered after a restart. The only other data that needs to be stored before restarts is the offset (from the start of the buffer, i.e., of this file) of the current state of the catchain. A simple key-value database can be used for this purpose.

3.2.4. Sharing data between different states. It turns out that the tree (or rather the dag) representing the new state $\sigma_{S \cup \{m\}} = f(\sigma_S, m)$ shares large subtrees with the previous state σ_S , and, similarly, $\sigma_{S \cup T} = g(\sigma_S, \sigma_T)$ shares large subtrees with σ_S and σ_T . The persistent structure used for representing the states in BCP makes it possible to reuse the same pointers inside the buffer for representing such shared data structures instead of duplicating them.

3.2.5. Memoizing nodes. Another technique employed while computing new states (i.e., the values of function f) is that of *memoizing new nodes*, also borrowed from functional programming languages. Namely, whenever a new node is constructed (inside the huge buffer containing all states for a specific catchain), its hash is computed, and a simple hash table is used to look up the latest node with the same hash. If a node with this hash is found, and it has the same contents, then the newly-constructed node is discarded and a reference to the old node with the same contents is returned instead. On the other hand, if no copy of the new node is found, then the hash table is updated, the end-of-buffer (allocation) pointer is advanced, and the pointer to the new node is returned to the caller.

In this way if different processes end up making similar computations and having similar states, large portions of these states will be shared even if they are not directly related by application of function f as explained in **3.2.4**.

3.2.6. Importance of optimization techniques. The optimization techniques **3.2.4** and **3.2.5** used for sharing parts of different internal states inside the same catchain are drastically important for improving the memory profile and the performance of BCM in a large process group. The improvement is several orders of magnitude in groups of $N \approx 100$ processes. Without these optimizations BCM would not be fit for its intended purpose (BFT consensus on new blocks generated by validators in the TOS Blockchain).

3.2.7. Message m contains a hash of state $\sigma_{D_m^+}$. Every message m contains a (Merkle) hash of (the abstract representation of) the corresponding state $\sigma_{D_m^+}$. Very roughly, this hash is computed recursively using the tree of nodes representation of **3.2.1**: all node references inside a node are replaced with (recursively computed) hashes of the referred nodes, and a simple 64-bit hash of the resulting byte sequence is computed. This hash is also used for memoization as described in **3.2.5**.

The purpose of this field in messages is to provide a sanity check for the

computations of $\sigma_{D_m^+}$ performed by different processes (and possibly by different implementations of the state update algorithm): once $\sigma_{D_m^+}$ is computed for a newly-delivered message m , the hash of computed $\sigma_{D_m^+}$ is compared to the value stored in the header of m . If these values are not equal, an error message is output into an error log (and no further actions are taken by the software). These error logs can be examined to detect bugs or incompatibilities between different versions of BCP.

3.3. State recovery after restart or crashes. A catchain is typically used by the BCP for several minutes; during this period, the program (the validator software) running the Catchain protocol may be terminated and restarted, either deliberately (e.g., because of a scheduled software update) or unintentionally (the program might crash because of a bug in this or some other subsystem, and be restarted afterwards). One way of dealing with this situation would be to ignore all catchains not created after the last restart. However, this would lead to some validators not participating in creating any blocks for several minutes (until the next catchain instances are created), which is undesirable. Therefore, a catchain state recovery protocol is run instead after every restart, so that the validator can continue participating in the same catchain.

3.3.1. Database of all delivered messages. To this end, a special database is created for each active catchain. This database contains all known and delivered messages, indexed by their identifiers (hashes). A simple key-value database suffices for this purpose. The hash of the most recent outbound message $m = m_{i,s}$ generated by the current process i is also stored in this database. After restart, all messages up to m are recursively delivered in proper order (in the same way as if all these messages had been just received from the network in an arbitrary order) and processed by the higher-level protocol, until m finally is delivered, thus recovering the current state.

3.3.2. Flushing new messages to disk. We have already explained in 2.6 that newly-created messages are stored in the database of all delivered messages (cf. 3.3.1) and the database is flushed to disk before the new message is sent to all network neighbors. In this way we can be sure that the message cannot be lost if the system crashes and is restarted, thus avoiding the creation of involuntary forks.

3.3.3. Avoiding the recomputation of states $\sigma_{D_m^+}$. An implementation might use an append-only file containing all previously computed states as described in 3.2.3 to avoid recomputing all states after restart, trading off disk space for computational power. However, the current implementation does not use this optimization.

3.4. High-level description of Block Consensus Protocol. Now we are ready to present a high-level description of the Block Consensus Protocol employed by TOS Blockchain validators to generate and achieve consensus on new blockchain blocks. Essentially, it is a three-phase commit protocol that runs over a catchain (an instance of the Catchain protocol), which is used as a “hardened” message broadcast system in a process group.

3.4.1. Creation of new catchain messages. Recall that the lower-level Catchain protocol does not create broadcast messages on its own (with the only exception being service broadcasts with fork proofs, cf. 2.7.5). Instead, when a new message needs to be created, the higher-level protocol (BCP) is asked to do this by invoking a callback. Apart from that, the creation of new messages may be triggered by changes in the current virtual state and by timer alarms.

3.4.2. Payload of catchain messages. In this way the payload of catchain messages is always determined by the higher level protocol, such as BCP. For BCP, this payload consists of

- Current Unix time. It must be non-decreasing on subsequent messages of the same process. (If this restriction is violated, all processes processing this message will tacitly replace this Unix time by the maximum Unix time seen in previous messages of the same sender.)
- Several (zero or more) *BCP events* of one of the admissible types listed below.

3.4.3. BCP events. We have just explained that the payload of a catchain message contains several (possibly zero) BCP events. Now we list all admissible BCP event types.

- `SUBMIT(round, candidate)` — suggest a new block candidate
- `APPROVE(round, candidate, signature)` — a block candidate has passed local validation

- $\text{REJECT}(\text{round}, \text{candidate})$ — a block candidate has failed local validation
- $\text{COMMITSIGN}(\text{round}, \text{candidate}, \text{signature})$ — a block candidate has been accepted and signed
- $\text{VOTE}(\text{round}, \text{candidate})$ — a vote for a block candidate
- $\text{VOTEFOR}(\text{round}, \text{candidate})$ — this block candidate must be voted for in this round (even if the current process has another opinion)
- $\text{PRECOMMIT}(\text{round}, \text{candidate})$ — a preliminary commitment to a block candidate (used in three-phase commit scheme)

3.4.4. Protocol parameters. Several parameters of BCP must be fixed in advance (in the genesis message of the catchain, where they are initialized from the values of the configuration parameters extracted from the current masterchain state):

- K — duration of one attempt (in seconds). It is an integer amount of seconds in the current implementation; however, this is an implementation detail, not a restriction of the protocol
- Y — number of *fast* attempts to accept a candidate
- C — block candidates suggested during one round
- Δ_i for $1 \leq i \leq C$ — delay before suggesting the block candidate with priority i
- Δ_∞ — delay before approving the null candidate

Possible values for these parameters are $K = 8$, $Y = 3$, $C = 2$, $\Delta_i = 2(i - 1)$, $\Delta_\infty = 2C$.

3.4.5. Protocol overview. The BCP consists of several *rounds* that are executed inside the same catchain. More than one round may be active at one point of time, because some phases of a round may overlap with other phases of other rounds. Therefore, all BCP events contain an explicit round identifier *round* (a small integer starting from zero). Every round is terminated either by (collectively) accepting a *block candidate* suggested by one of the participating processes, or by accepting a special *null candidate*—a

dummy value indicating that no real block candidate was accepted, for example because no block candidates were suggested at all. After a round is terminated (from the perspective of a participating process), i.e., once a block candidate collects COMMITSIGN signatures of more than $2/3$ of all validators, only COMMITSIGN events may be added to that round; the process automatically starts participating in the next round (with the next identifier) and ignores all BCP events with different values of *round*.⁶

Each round is subdivided into several *attempts*. Each attempt lasts a predetermined time period of K seconds (BCP uses clocks to measure time and time intervals and assumes that clocks of “good” processes are more or less in agreement with each other; therefore, BCP is not an asynchronous BFT protocol). Each attempt starts at Unixtime exactly divisible by K and lasts for K seconds. The attempt identifier *attempt* is the Unixtime of its start divided by K . Therefore, the attempts are numbered more or less consecutively by 32-bit integers, but not starting from zero. The first Y attempts of a round are *fast*; the remaining attempts are *slow*.

3.4.6. Attempt identification. Fast and slow attempts. In contrast with rounds, BCP events do not have a parameter to indicate the attempt they belong to. Instead, this attempt is implicitly determined by the Unix time indicated in the payload of the catchain message containing the BCP event (cf. 3.4.2). Furthermore, the attempts are subdivided into *fast* (the first Y attempts of a round in which a process takes part) and *slow* (the subsequent attempts of the same round). This subdivision is also implicit: the first BCP event sent by a process in a round belongs to a certain attempt, and Y attempts starting from this one are considered fast by this process.

3.4.7. Block producers and block candidates. There are C designated block producers (member processes) in each round. The (ordered) list of these block producers is computed by a deterministic algorithm (in the simplest case, processes $i, i+1, \dots, i+C-1$ are used in the i -th round, with the indices taken modulo N , the total number of processes in the catchain) and is known to all participants without any extra communication or negotiation. The processes are ordered in this list by decreasing priority, so the first member of the list has the highest priority (i.e., if it suggests a block candidate in

⁶This also means that each process implicitly determines the Unixtime of the start of the next round, and computes all delays, e.g., the block candidate submission delays, starting from this time.

time, this block candidate has a very high chance to be accepted by the protocol).

The first block producer may suggest a block candidate immediately after the round starts. Other block producers can suggest block candidates only after some delay Δ_i , where i is the index of the producer in the list of designated block producers, with $0 = \Delta_1 \leq \Delta_2 \leq \dots$. After some predetermined period of time Δ_∞ elapses from the round start, a special *null candidate* is assumed automatically suggested (even if there are no explicit BCP events to indicate this). Therefore, at most $C + 1$ block candidates (including the null candidate) are suggested in a round.

3.4.8. Suggesting a block candidate. A block candidate for the TOS Blockchain consists of two large “files” — the block and the collated data, along with a small header containing the description of the block being generated (most importantly, the complete *block identifier* for the block candidate, containing the workchain and the shard identifier, the block sequence number, its file hash and its root hash) and the SHA256 hashes of the two large files. Only a part of this small header (including the hashes of the two files and other important data) is used as *candidate* in BCP events such as SUBMIT or COMMITSIGN to refer to a specific block candidate. The bulk of the data (most importantly, the two large files) is propagated in the overlay network associated with the catchain by the streaming broadcast protocol implemented over ADNL for this purpose (cf. [3, 5]). This bulk data propagation mechanism is unimportant for the validity of the consensus protocol (the only important point is that the hashes of the large files are part of BCP events and hence of the catchain messages, where they are signed by the sender, and these hashes are checked after the large files are received by any participating nodes; therefore, nobody can replace or corrupt these files). A SUBMIT(*round, candidate*) BCP event is created in the catchain by the block producer in parallel with the propagation of the block candidate, indicating the submission of this specific block candidate by this block producer.

3.4.9. Processing block candidates. Once a process observes a SUBMIT BCP event in a delivered catchain message, it checks the validity of this event (for instance, its originating process must be in the list of designated producers, and current Unixtime must be at least the start of the round plus the minimum delay Δ_i , where i is the index of this producer in the list of designated producers), and if it is valid, remembers it in the current

catchain state (cf. 3.1). After that, when a streaming broadcast containing the files associated with this block candidates (with correct hash values) is received (or immediately, if these files are already present), the process invokes a validator instance to validate the new block candidate (even if this block candidate was suggested by this process itself!). Depending on the result of this validation, either an `APPROVE(round, candidate, signature)` or a `REJECT(round, candidate)` BCP event is created (and embedded into a new catchain message). Note that the *signature* used in `APPROVE` events uses the same private key that will ultimately be used to sign the accepted block, but the signature itself is different from that used in `COMMITSIGN` (the hash of a structure with different magic number is actually signed). Therefore, this interim signature cannot be used to fake the acceptance of this block by this particular validator process to an outside observer.

3.4.10. Overview of one round. Each round of BCP proceeds as follows:

- At the beginning of a round, several processes (from the predetermined list of designated producers) submit their block candidates (with certain delays depending on their producer priority) and reflect this fact by means of `SUBMIT` events (incorporated into catchain messages).
- Once a process receives a submitted block candidate (i.e., observes a `SUBMIT` event and receives all necessary files by means external to the consensus protocol), it starts the validation of this candidate and eventually creates either an `APPROVE` or a `REJECT` event for this block candidate.
- During each *fast attempt* (i.e., one of the first Y attempts) every process votes either for a block candidate that has collected the votes of more than $2/3$ of all processes, or, if there are no such candidates yet, for the valid (i.e., `APPROVED` by more than $2/3$ of all processes) block candidate with the highest priority. The voting is performed by means of creating `VOTE` events (embedded into new catchain messages).
- During each *slow attempt* (i.e., any attempt except the first Y) every process votes either for a candidate that was `PRECOMMITTED` before (by the same process), or for a candidate that was suggested by `VOTEFOR`.

- If a block candidate has received votes from more than $2/3$ of all processes during the current attempt, and the current process observes these votes (which are collected in the catchain state), a PRECOMMIT event is created, indicating that the process will vote only for this candidate in future.
- If a block candidate collects PRECOMMITs from more than $2/3$ of all processes inside an attempt, then it is assumed to be accepted (by the group), and each process that observes these PRECOMMITs creates a COMMITSIGN event with a valid block signature. These block signatures are registered in the catchain, and are ultimately collected to create a “block proof” (containing signatures of more than $2/3$ of the validators for this block). This block proof is the external output of the consensus protocol (along with the block itself, but without its collated data); it is ultimately propagated in the overlay network of all full nodes that have subscribed to new blocks of this shard (or of the masterchain).
- Once a block candidate collects COMMITSIGN signatures from more than $2/3$ of all validators, the round is considered finished (at least from the perspective of a process that observes all these signatures). After that, only a COMMITSIGN can be added to that round by this process, and the process automatically starts participating in the next round (and ignores all events related to other rounds).

Note that the above protocol may lead to a validator signing (in a COMMITSIGN event) a block candidate that was REJECTED by the same validator before (this is a kind of “submitting to the will of majority”).

3.4.11. VOTE and PRECOMMIT messages are created deterministically. Note that each process can create at most one VOTE and at most one PRECOMMIT event in each attempt. Furthermore, these events are completely determined by the state σ_{D_m} of the sender of catchain message m containing such an event. Therefore, the receiver can detect invalid VOTE or PRECOMMIT events and ignore them (thus mitigating byzantine behavior of other participants). On the other hand, a message m that should contain a VOTE or a PRECOMMIT event according to the corresponding state σ_{D_m} but does not contain one can be received. In this case, the current implementation automatically creates missing events and proceeds as if m had contained

them from the very beginning. However, such instances of byzantine behavior are either corrected or ignored (and a message is output into the error log), but the offending processes are not otherwise punished (because this would require very large misbehavior proofs for outside observers that do not have access to the internal state of the catchain).

3.4.12. Multiple VOTES and PRECOMMITs of the same process. Note that a process usually ignores subsequent VOTES and PRECOMMITs generated by the same originating process inside the same attempt, so normally a process can vote for at most one block candidate. However, it may happen that a “good” process indirectly observes a fork created by a byzantine process, with VOTES for different block candidates in different branches of this fork (this can happen if the “good” process learns about these two branches from two other “good” processes that did not see this fork before). In this case, both VOTES (for different candidates) are taken into account (added into the merged state of the current process). A similar logic applies to PRECOMMITs.

3.4.13. Approving or rejecting block candidates. Notice that a block candidate cannot be APPROVED or REJECTED before it has been SUBMITTED (i.e., an APPROVE event that was not preceded by a corresponding SUBMIT event will be ignored), and that a candidate cannot be approved before the minimum time of its submission (the round start time plus the priority-dependent delay Δ_i) is reached, i.e., any “good” process will postpone the creation of its APPROVE until this time. Furthermore, one cannot APPROVE more than one candidate of the same producer in the same round (i.e., even if a process SUBMITS several candidates, only one of them—presumably the first one—will be APPROVED by other “good” processes; as usual, this means that subsequent APPROVE events will be ignored by “good” processes on receipt).

3.4.14. Approving the null block candidate. The implicit null block candidate is also explicitly approved (by creating an APPROVE event) by all (good) processes, once the delay Δ_∞ from the start of the round expires.

3.4.15. Choosing a block candidate for voting. Each process chooses one of the available block candidates (including the implicit null candidate) and votes for this candidate (by creating a VOTE event) by applying the following rules (in the order they are presented):

- If the current process created a PRECOMMIT event for a candidate during one of the previous attempts, and no other candidate has collected votes from more than $2/3$ of all processes since (i.e., inside one of the subsequent attempts, including the current one so far; we say that the PRECOMMIT event is still *active* in this case), then the current process votes for this candidate again.
- If the current attempt is fast (i.e., one of the first Y attempts of a round from the perspective of the current process), and a candidate has collected votes from more than $2/3$ of all processes during the current or one of the previous attempts, the current process votes for this candidate. In the case of a tie, the candidate from the latest of all such attempts is chosen.
- If the current attempt is fast, and the previous rules do not apply, then the process votes for the candidate with the highest priority among all *eligible candidates*, i.e., candidates that have collected APPROVES (observable by the current process) from more than $2/3$ of all processes.
- If the current attempt is slow, then the process votes only after it receives a valid VOTEFOR event in the same attempt. If the first rule is applicable, the process votes according to it (i.e., for the previously PRECOMMITTED candidate). Otherwise it votes for the block candidate that is mentioned in the VOTEFOR event. If there are several such valid events (during the current attempt), the candidate with the smallest hash is selected (this may happen in rare situations related to different VOTEFOR events created in different branches of a fork, cf. **3.4.12**).

The “null candidate” is considered to have the least priority. It also requires an explicit APPROVE before being voted for (with the exception of the first two rules).

3.4.16. Creating VOTEFOR events during slow attempts. A VOTEFOR event is created at the beginning of a slow attempt by the *coordinator* — the process with index $attempt \bmod N$ in the ordered list of all processes participating in the catchain (as usual, this means that a VOTEFOR created by another process will be ignored by all “good” processes). This VOTEFOR event refers to one of the block candidates (including the null candidate) that have collected APPROVES from more than $2/3$ of all processes, usually randomly chosen among all such candidates. Essentially, this is a suggestion

to vote for this block candidate directed to all other processes that do not have an active PRECOMMIT.

3.5. Validity of BCP. Now we present a sketch of the proof of validity of TOS Block Consensus Protocol (BCP) described above in **3.4**, assuming that less than one third of all processes exhibit byzantine (arbitrarily malicious, possibly protocol-violating) behavior, as it is customary for Byzantine Fault Tolerant protocols. During this subsection, we consider only one round of BCP, subdivided into several attempts.

3.5.1. Fundamental assumption. Let us emphasize once again that we assume that *less than one third of all processes are byzantine*. All other processes are assumed to be *good*, i.e., they follow the protocol.

3.5.2. Weighted BCP. The reasoning in this subsection is valid for the *weighted variant of BCP* as well. In this variant, each process $i \in I$ is pre-assigned a positive weight $w_i > 0$ (fixed in the genesis message of the catchain), and statements about “more than 2/3 of all processes” and “less than one third of all processes” are understood as “more than 2/3 of all processes by weight”, i.e., “a subset $J \subset I$ of processes with total weight $\sum_{j \in J} w_j > \frac{2}{3} \sum_{i \in I} w_i$ ”, and similarly for the second property. In particular, our “fundamental assumption” **3.5.1** is to be understood in the sense that “the total weight of all byzantine processes is less than one third of the total weight of all processes”.

3.5.3. Useful invariants. We collect here some useful invariants obeyed by all BCP events during one round of BCP (inside a catchain). These invariants are enforced in two ways. Firstly, any “good” (non-byzantine) process will not create events violating these invariants. Secondly, even if a “bad” process creates an event violating these invariants, all “good” processes will detect this when a catchain message containing this event is delivered to BCP and ignore such events. Some possible issues related to forks (cf. **3.4.12**) remain even after these precautions; we indicate how these issues are resolved separately, and ignore them in this list. So:

- There is at most one SUBMIT event by each process (inside one round of BCP).
- There is at most one APPROVE or REJECT event by each process related to one candidate (more precisely, even if there are multiple candidates created by the same designated block producer, only one of them

can be APPROVED by another process).⁷ This is achieved by requiring all “good” processes to ignore (i.e., not to create APPROVES or REJECTS for) all candidates suggested by the same producer but the very first one they have learned about.

- There is at most one VOTE and at most one PRECOMMIT event by each process during each attempt.
- There is at most one VOTEFOR event during each (slow) attempt.
- There is at most one COMMITSIGN event by each process.
- During a slow attempt, each process votes either for its previously PRECOMMITTED candidate, or for the candidate indicated in the VOTEFOR event of this attempt.

One might somewhat improve the above statements by adding the word “valid” where appropriate (e.g., there is at most one *valid* SUBMIT event...).

3.5.4. More invariants.

- There is at most one eligible candidate (i.e., candidate that has received APPROVES from more than $2/3$ of all processes) from each designated producer, and no eligible candidates from other producers.
- There are at most $C + 1$ eligible candidates in total (at most C candidates from C designated producers, plus the null candidate).
- A candidate may be accepted only if it has collected more than $2/3$ PRECOMMITs during the same attempt (more precisely, a candidate is accepted only if there are PRECOMMIT events created by more than $2/3$ of all processes for this candidate and belonging to the same attempt).

⁷In fact, REJECTS appear only in this restriction, and do not affect anything else. Therefore, any process can abstain from sending REJECTS without violating the protocol, and REJECT events could have been removed from the protocol altogether. Instead, the current implementation of the protocol still generates REJECTS, but does not check anything on their receipt and does not remember them in the catchain state. Only a message is output into the error log, and the offending candidate is stored into a special directory for future study, because REJECTS usually indicate either the presence of a byzantine adversary, or a bug in the collator (block generation) or validator (block verification) software either on the node that suggested the block or on the node that created the REJECT event.

- A candidate may be VOTED for, PRECOMMITTED, or mentioned in a VOTEFOR only if it is an *eligible candidate*, meaning that it has previously collected APPROVES from more than $2/3$ of all validators (i.e., a valid VOTE event may be created for a candidate only if APPROVE events for this candidate have been previously created by more than $2/3$ of all processes and registered in catchain messages observable from the message containing the VOTE event, and similarly for PRECOMMIT and VOTEFOR events).

3.5.5. At most one block candidate is accepted. Now we claim that *at most one block candidate can be accepted (in a round of BCP)*. Indeed, a candidate can be accepted only if it collects PRECOMMITs from more than $2/3$ of all processes inside the same attempt. Therefore, two different candidates cannot achieve this during the same attempt (otherwise more than one third of all validators must have created PRECOMMITs for two different candidates inside an attempt, thus violating the above invariants; but we have assumed that less than one third of all validators exhibit byzantine behavior). Now suppose that two different candidates c_1 and c_2 have collected PRECOMMITs from more than $2/3$ of all processes in two different attempts a_1 and a_2 . We may assume that $a_1 < a_2$. According to the first rule of **3.4.15**, each process that has created a PRECOMMIT for c_1 during attempt a_1 must continue voting for c_1 in all subsequent attempts $a' > a_1$, or at least cannot vote for any other candidate, unless another candidate c' collects VOTES of more than $2/3$ of all processes during a subsequent attempt (and this invariant is enforced even if some processes attempt not to create these new VOTE events for c_1 , cf. **3.4.11**). Therefore, if $c_2 \neq c_1$ has collected the necessary amount of PRECOMMITs during attempt $a_2 > a_1$, there is at least one attempt a' , $a_1 < a' \leq a_2$, such that some $c' \neq c_1$ (not necessarily equal to c_2) has collected VOTES of more than $2/3$ of all processes during attempt a' . Let us fix the smallest such a' , and the corresponding $c' \neq c_1$ that has collected many votes during attempt a' . More than $2/3$ of all validators have voted for c' during attempt a' , and more than $2/3$ of all validators have PRECOMMITTED for c_1 during attempt a_1 , and by the minimality of a' there was no attempt a'' with $a_1 < a'' < a'$, such that a candidate distinct from c_1 collected more than $2/3$ of all votes during attempt a'' . Therefore, all validators that PRECOMMITTED for c_1 could vote only for c_1 during attempt a' , and at the same time we supposed that c' has collected votes from more than $2/3$ of all validators during the same attempt a' . This implies that

more than $1/3$ of all validators have somehow voted both for c_1 and c' during this attempt (or voted for c' while they could have voted only for c_1), i.e., more than $1/3$ of all validators have exhibited byzantine behavior. This is impossible by our fundamental assumption **3.5.1**.

3.5.6. At most one block candidate may be PRECOMMITTED during one attempt. Note that all valid PRECOMMIT events (if any) created inside the same attempt must refer to the same block candidate, by the same reasoning as in the first part of **3.5.5**: since a valid PRECOMMIT event for a candidate c may be created only after votes from more than $2/3$ of all processes are observed for this candidate inside the same attempt (and invalid PRECOMMITs are ignored by all good processes), the existence of valid PRECOMMIT events for different candidates c_1 and c_2 inside the same attempt would imply that more than one third of all processes have voted both for c_1 and c_2 inside this attempt, i.e., they have exhibited byzantine behavior. This is impossible in view of our fundamental assumption **3.5.1**.

3.5.7. A previous PRECOMMIT is deactivated by the observation of a newer one. We claim that *whenever a process with an active PRECOMMIT observes a valid PRECOMMIT created by any process in a later attempt for a different candidate, its previously active PRECOMMIT is deactivated*. Recall that we say that a process has an *active* PRECOMMIT if it has created a PRECOMMIT for a certain candidate c during a certain attempt a , did not create any PRECOMMIT during any attempts $a' > a$, and did not observe votes of more than $2/3$ of all validators for any candidate $\neq c$ during any attempts $a' > a$. Any process has at most one active PRECOMMIT, and if it has one, it must vote only for the precommitted candidate.

Now we see that if a process with an active PRECOMMIT for a candidate c since attempt a observes a valid PRECOMMIT (usually by another process) for a candidate c' created during some later attempt $a' > a$, then the first process must also observe all dependencies of the message that contains the newer PRECOMMIT; these dependencies necessarily include valid VOTES from more than $2/3$ of all validators for the same candidate $c' \neq c$ created during the same attempt $a' > a$ (because otherwise the newer PRECOMMIT would not be valid, and would be ignored by the other process); by definition, the observation of all these VOTES deactivates the original PRECOMMIT.

3.5.8. Assumptions for proving the convergence of the protocol. Now we are going to prove that the protocol described above *converges* (i.e.,

terminates after accepting a block candidate) with probability one under some assumptions, which essentially tell us that there are enough “good” processes (i.e., processes that diligently follow the protocol and do not introduce arbitrary delays before sending their new messages), and that these good processes enjoy good network connectivity at least from time to time. More precisely, our assumptions are as follows:

- There is a subset $I^+ \subset I$ consisting of “good” processes and containing more than $2/3$ of all processes.
- All processes from I^+ have well-synchronized clocks (differing by at most τ , where τ is a bound for network latency described below).
- If there are infinitely many attempts, then infinitely many attempts are “good” with respect to network connectivity between processes from I^+ , meaning that all messages created by a process from I^+ during this attempt or earlier are delivered to any other process from I^+ within at most $\tau > 0$ seconds after being created with probability at least $q > 0$, where $\tau > 0$ and $0 < q < 1$ are some fixed parameters, such that $5\tau < K$, where K is the duration of one attempt.
- Furthermore, if the protocol runs for infinitely many attempts, then any arithmetic progression of attempts contains infinitely many “good” attempts in the sense described above.
- A process from I^+ creates a VOTEFOR during a slow attempt after some fixed or random delay after the start of the slow attempt, in such a way that this delay belongs to the interval $(\tau, K - 3\tau)$ with probability at least q' , where $q' > 0$ is a fixed parameter.
- A process from I^+ , when it is its turn to be the coordinator of a slow attempt, chooses a candidate for VOTEFOR uniformly at random among all eligible candidates (i.e., those candidates that have collected APPROVES from more than $2/3$ of all validators).

3.5.9. The protocol terminates under these assumptions. Now we claim that *(each round of) the BCP protocol as described above terminates with probability one under the assumptions listed in 3.5.8*. The proof proceeds as follows.

- Let us assume that the protocol does not converge. Then it continues running forever. We are going to ignore the first several attempts, and consider only attempts $a_0, a_0 + 1, a_0 + 2, \dots$ starting from some a_0 , to be chosen later.
- Since all processes from I^+ continue participating in the protocol, they will create at least one message not much later than the start of the round (which may be perceived slightly differently by each process). For instance, they will create an APPROVE for the null candidate no later than Δ_∞ seconds from the start of the round. Therefore, they will consider all attempts slow at most KY seconds afterwards. By choosing a_0 appropriately, we can assume that all attempts we consider are slow from the perspective of all processes from I^+ .
- After a “good” attempt $a \geq a_0$ all processes from I^+ will see the APPROVES for the null candidate created by all other processes from I^+ , and will deem the null candidate eligible henceforth. Since there are infinitely many “good” attempts, this will happen sooner or later with probability one. Therefore, we can assume (increasing a_0 if necessary) that there is at least one eligible candidate from the perspective of all processes from I^+ , namely, the null candidate.
- Furthermore, there will be infinitely many attempts $a \geq a_0$ that are perceived slow by all processes from I^+ , that have a coordinator from I^+ , and that are “good” (with respect to the network connectivity) as defined in 3.5.8. Let us call such attempts “very good”.
- Consider one “very good” slow attempt a . With probability $q' > 0$, its coordinator (which belongs to I^+) will wait for $\tau' \in (\tau, K - 3\tau)$ seconds before creating its VOTEFOR event. Consider the most recent PRECOMMIT event created by any process from I^+ ; let us suppose it was created during attempt $a' < a$ for some candidate c' . With probability $qq' > 0$, the catchain message carrying this PRECOMMIT will be already delivered to the coordinator at the time of generation of its VOTEFOR event. In that case, the catchain message carrying this VOTEFOR will depend on this PRECOMMIT(c') event, and all “good” processes that observe this VOTEFOR will also observe its dependencies, including this PRECOMMIT(c'). We see that *with probability at least qq' , all processes from I^+ that receive the VOTEFOR event during*

a “very good” slow attempt receive also the most recent PRECOMMIT (if any).

- Next, consider any process from I^+ that receives this VOTEFOR, for a randomly chosen eligible candidate c , and suppose that there are already some PRECOMMITs, and that the previous statement holds. Since there are at most $C + 1$ eligible candidates (cf. **3.5.4**), with probability at least $1/(C + 1) > 0$ we’ll have $c = c'$, where c' is the most recently PRECOMMITTED candidate (there is at most one such candidate by **3.5.6**). In this case, all processes from I^+ will vote for $c = c'$ during this attempt immediately after they receive this VOTEFOR (which will be delivered to any process $j \in I^+$ less than $K - 2\tau$ seconds after the beginning of the attempt with probability qq'). Indeed, if a process j from I^+ did not have an active PRECOMMIT, it will vote for the value indicated in VOTEFOR, which is c . If j had an active PRECOMMIT, and it is as recent as possible, i.e., also created during attempt a' , then it must have been a PRECOMMIT for the same value $c' = c$ (because we know about at least one valid PRECOMMIT for c' during attempt a' , and all other valid PRECOMMITs during attempt a' must be for the same c' by **3.5.6**). Finally, if j had an active PRECOMMIT from an attempt $< a'$, then it will become inactive once the VOTEFOR with all its dependencies (including the newer PRECOMMIT(c')) has been delivered to this process j (cf. **3.5.7**), and the process will again vote for the value c indicated in VOTEFOR. Therefore, all processes from I^+ will vote for the same $c = c'$ during this attempt, less than $K - 2\tau$ seconds after the beginning of the attempt (with some probability bounded away from zero).
- If there are no PRECOMMITs yet, then the above reasoning simplifies further: all processes from I^+ that receive this VOTEFOR will immediately vote for the candidate c suggested by this VOTEFOR.
- In both cases, all processes from I^+ will create a VOTE for the same candidate c less than $K - 2\tau$ seconds from the beginning of the attempt, and this will happen with a positive probability bounded away from zero.
- Finally, all processes from I^+ will receive these VOTES for c from all processes from I^+ , again less than $(K - 2\tau) + \tau = K - \tau$ seconds

after the beginning of this attempt, i.e., still during the same attempt (even after taking into account the imperfect clock synchronization between processes from I^+). This means that they will all create a valid PRECOMMIT for c , i.e., the protocol will accept c during this attempt with probability bounded away from zero.

- Since there are infinitely many “very good” attempts, and the probability of successful termination during each such attempt is $\geq p > 0$ for some fixed value of p , the protocol will terminate successfully with probability one.

References

- [1] K. BIRMAN, *Reliable Distributed Systems: Technologies, Web Services and Applications*, Springer, 2005.
- [2] M. CASTRO, B. LISKOV, ET AL., *Practical byzantine fault tolerance*, *Proceedings of the Third Symposium on Operating Systems Design and Implementation* (1999), p. 173–186, available at <http://pmg.csail.mit.edu/papers/osdi99.pdf>.
- [3] N. DUROV, *TOS Network*, 2017.
- [4] N. DUROV, *TOS Network Blockchain*, 2018.
- [5] L. LAMPORT, R. SHOSTAK, M. PEASE, *The byzantine generals problem*, *ACM Transactions on Programming Languages and Systems*, **4/3** (1982), p. 382–401.
- [6] A. MILLER, YU XIA, ET AL., *The honey badger of BFT protocols*, Cryptology e-print archive 2016/99, <https://eprint.iacr.org/2016/199.pdf>, 2016.
- [7] M. VAN STEEN, A. TANENBAUM, *Distributed Systems*, 3rd ed., 2017.