

CellTree: A New Paradigm for Distributed Data Repositories

Anasuya Acharya
Bar-Ilan University
acharya@biu.ac.il

Akash Trehan
Microsoft
aktrehan@microsoft.com

Manoj Prabhakaran
IIT Bombay
mp@cse.iitb.ac.in

March 15, 2021

Abstract

We present CellTree, a new architecture for distributed data repositories. The repository allows data to be stored in largely independent, and highly programmable cells, which are “assimilated” into a tree structure. The data in the cells are allowed to change over time, subject to each cell’s own policies; a cell’s policies also govern how the policies themselves can evolve. A design goal of the architecture is to let a CellTree evolve organically over time, and adapt itself to multiple applications. Different parts of the tree may be maintained by different sets of parties interested in the respective parts, and the core mechanisms used for maintaining the tree can also vary across the tree and over time.

We present provable guarantees of liveness, correctness and consistency (the last one being a generalization of the typical blockchain guarantee of “persistence,” when data is dynamic), when the CellTree architecture is instantiated using a simple set of modules. These properties can be guaranteed for individual cells that satisfy requisite trust assumptions, even if these trust assumptions do not hold for other cells in the tree. This work also includes a concrete design of a working proof-of-concept. Furthermore, we discuss several features of a CellTree that can be exploited by applications. We leave it for future work to develop full-fledged applications on top of this powerful architecture.

1 Introduction

There has been an explosion of interest in the notion of a *distributed ledger*, triggered by the popularity of Bitcoin [16]. The typical distributed ledgers today have the form of a *blockchain*, where each new block points to an earlier block. A variety of ingenious protocols have been developed to add and *immutably* maintain the blocks in such a ledger in an adversarial environment, while incentivizing participation. Blockchain applications today have gone well beyond that envisaged by Bitcoin (namely, to publicly store all the transactions of a cryptocurrency), and have been used (or proposed for use)

as ledgers for supply chains, in gaming, for maintaining public records, and for general purpose contracts and transactions.

However, as blockchains have exploded in popularity and scope, several issues have come to the fore and – despite several proposed alternatives – many of them remain unresolved. Typical blockchain architectures face scalability challenges, as the underlying blockchain structure is ever-growing and all the “full nodes” need to store this entire chain to be able to fully validate *any* block. Further, this data may include illegal content, creating legal complications for the blockchains [13, 17]. Implementation bugs, if exploited, can create irreversible effects, owing to the immutable nature of a blockchain. Another issue of note is that the proof-of-work consensus protocols used in popular blockchains have turned out to be ecologically costly.

Many recent works have proposed fixes to these – using secondary structures [1, 18] to reduce the content stored on a blockchain, using alternate consensus mechanisms [8, 10, 11], and using different graph topologies [2, 19]. However, as these systems grow in scale, or if new attacks emerge due to changing economic incentives or discovery of new bugs, the deployments of these solutions may also need revisions.

In this work, we propose a very different approach to the problem of designing a large scale, long-running system for distributed data storage. A distributed data repository is a complex system, with several constituent components, addressing several disparate sub-problems. A key philosophy of our approach is to *let different solutions coexist* in the system, and to *leave room for the system to organically evolve* over time and across applications. For this, the various sub-problems are delegated to *modules*, and the overall architecture is agnostic to how each module is implemented.

In identifying the modules, we delineate tasks that are typically conflated together in current approaches. For instance, in a typical blockchain, selecting the contents for a new cell, and confirming the cell as part of a blockchain are usually accomplished by a single

mechanism that serves both as a lottery system and a consensus protocol; in our design, these two tasks are assigned to separate modules, and typically separate sets of parties are responsible for the two.

An important consideration in our design is to let different users of the system focus on different parts of the system, unencumbered by the entire system’s data. A related feature is the ability of a part of the whole structure to function as a smaller version of the system, complete with its security guarantees and trust assumptions. In particular, we naturally admit *multi-level confirmation* of new blocks, so that clients trusting lower levels in this hierarchy can get quick confirmation of the addition of a block to the structure, and those who do not trust those levels can wait for a higher level of confirmation.

The CellTree Architecture. We term our new design *CellTree*, contrasting it with a blockchain. A CellTree consists of largely independent *cells*, with their own rules for collecting and updating data. As mentioned above, the data in each cell can evolve subject to policies programmed into the cell, and even these policies can evolve as they permit themselves to. Each cell is “operated” by a crew selected for it (using a selection procedure that is encapsulated into a module that may vary across cells).

The cells act collectively to provide various guarantees. Each cell is addressed by associating it with a node in a tree. Data in the different cells are periodically *assimilated* into the tree. Each cell’s integrity and availability guarantees would depend only on (the crews of) the cells in the path from it to the root of the tree (or to any node whose crew is considered to have an honest-majority). Subtrees could be easily excised from the CellTree, or grafted to multiple locations (effectively making the structure a directed acyclic graph, rather than a tree), with little effect on cells outside the subtree.

The cell structure and the tree structure are complementary, and it is primarily the interface between them that is fixed by the CellTree architecture. The individual cells themselves can be programmed to evolve in customized ways, and the tree’s protocols can also be customized using various modules at the level of nodes, edges and paths.

2 An Organic Design

Before we present the technical details of our construction, we pause to reflect on our key design goals: We seek a flexible distributed repository, where multiple solutions to various sub-problems can co-exist, and even evolve over time, to address different use-cases and attack scenarios. We require the framework to be rich enough to directly support a variety of conventional applications of

distributed ledgers (and more), but should not be over-engineered so that it becomes too rigid for unforeseen future applications. The CellTree framework strives to meet this goal by being *modular*, *cellular* and *evolving*, as sketched below.

Modular We separate out processes (a) for assigning responsibilities to parties, (b) for selecting updates to be applied to the data in the system, and (c) for parties to enforce integrity and availability of the data they are responsible for. Our focus in this work will be on (c), with support for a wide variety of options for the modules for (a) and (b). Indeed, concepts from prior blockchain constructions, like proof-of-work, proof-of-stake and the use of crypto currencies as incentives, as well as relying on permissioned systems, all remain possible means to achieving the goals of (a) and (b) in our system. Even in our solution for (c), various sub-tasks are left as modules, which can be implemented in a variety of ways to achieve different efficiency-robustness trade-offs.

Cellular The system operates in a highly parallel manner, allowing concurrent updates to different *cells* which function largely autonomously and independent of each other, even while collaborating for the purposes of supporting each other. This independence manifests in different ways:

- The parties hosting the distributed system can choose to focus on cells of interest to them, and ignore the remaining.
- Each cell may offer a different functionality.
- As time progresses more cells can be added to the system. Cells could also be excised from the system, with limited impact on the rest of the system.
- Each cell may use its own set of mechanisms for implementing the modules mentioned above.

Evolving We envisage cells as having static addresses, but their data content could change over time. Data updates in a cell must follow an *evolution policy* associated with the cell. Simple examples of evolution policy include one that enforces that the data remains static and one that only allows blocks to be appended to a list of data blocks. The evolution policy in a cell itself can evolve, and this is also dictated by the evolution policy.¹ Consistency and liveness properties will be defined to account for evolving data.

There is one more sense in which a CellTree can evolve: the modules used by various nodes in the tree can evolve, subject to their own evolutionary policies. However, we do not mandate a framework for how modules are chosen and modified, instead leaving this itself to a module.

¹For meaningful guarantees, when a policy rewrites itself, we require the newly resulting policy to validate that the old policy is acceptable to it as a policy to evolve from.

2.1 Design Overview

The basic building block of a CellTree is a *cell*, which carries the data, as well as a (typically smaller) *nucleus*. As detailed in [Section 4](#), the nucleus has code that specifies how a cell can evolve (e.g., cell data can only be appended to).

The cell is hosted in a *node* in a binary tree. Each node is operated by a *crew*. Each node’s crew is also in charge of *monitoring* the nodes in a relatively small subtree rooted at that node. Monitoring involves verifying that the evolution of a nucleus is consistent with its own policies. If so, the monitoring node *assimilates* the updated nucleus into its own hash pointer, which is then propagated to its own ancestors. Further, it also propagates a *proof of assimilation* towards the nodes it is monitoring. These proofs will be verified by a client accessing a cell.

We specify several algorithms executed by the crew members operating a node (or clients accessing a node), to carry out a cell’s evolution, to monitor the evolution of some other cells, and to assimilate updates and propagate assimilation signals up and down the tree. These algorithms are given in [Section 5](#). Many tasks like selecting the next step in a cell’s evolution, selecting crews for new nodes, communicating to other crews, etc. are left to the *modules* listed in [Section 5.2](#).

2.2 CellTree vs. Blockchains

A conventional blockchain, like that of Bitcoin, could be viewed as a single node CellTree, with an appropriately programmed cell, and a very large crew. Conversely, a CellTree has many parallels with a blockchain, with cells in lieu of blocks. But the CellTree architecture differs from blockchains in several important and fundamental ways. Below, we discuss a few such differences.

Multi-Level Confirmation. In a CellTree, a crew operating a node can update its cell autonomously, based on its own local policies, and then have it assimilated into the tree. Unlike in a blockchain architecture, where the whole system has to approve the update, the updates within a cell are wholly determined by the crew of the node hosting the cell (subject to the policies programmed into the cell); a client that trusts the node’s crew has *immediate confirmation* of the update. The purpose of assimilation and higher levels of confirmation is only to protect against (and hence disincentivize) misbehaving crews.

Reversed Hash Pointers. One of the easily spotted difference between the CellTree architecture and blockchains is that the former uses a tree topology, instead of a chain.

At first glance, one may view a blockchain as a *unary* Merkle Tree [\[14, 15\]](#) (i.e., each node has a single child),

while a CellTree uses a binary Merkle tree. However, this comparison is misleading. In a CellTree, the “hash pointers” point from a parent (older node) to its children (newer nodes), whereas in a blockchain, newly added nodes carry hash pointers to existing nodes.² That is, if we view a blockchain as a unary Merkle tree, the “genesis” block would play the role of a leaf, and the most recently added block becomes the root. The reversed direction of hash pointers reflects the fact that in a conventional blockchain, a block gets confirmed when future blocks attach to it, whereas in a CellTree, a cell (or an update to a cell) is confirmed by updating already existing nodes.

Distributed vs. Replicated. All the parties actively involved in maintaining a blockchain replicate the entire blockchain, and their consensus is required while updating the blockchain in any manner. In the CellTree architecture, on the other hand, each individual node is owned and operated by its own crew. Distributing the ownership of nodes to relatively small crews (say, involving a few dozen parties) which can operate in parallel, vastly improves the scalability of the system.

We point out that selecting a crew that can entirely own a node does raise some issues that need to be handled. Firstly, whenever possible, we would like to ensure that a sufficiently large majority of the parties selected for a crew are honest, and secondly, in the event that the crews for some nodes are compromised (say, because a small crew is used), we would still like to retain security guarantees for other nodes. The first issue is left for a module to address: A typical approach would be to randomly sample from a pool that is considered to have a sufficiently large majority of parties (weighted appropriately by the extent of their computational power, stake in the system etc.). Practical instantiations of such schemes appear in protocols like Algorand [\[8\]](#). The second issue is addressed by the CellTree architecture, via *monitoring*: Even if a node’s crew is corrupt, any updates to the node’s cell will be assimilated into the CellTree only if they are validated and agreed upon by the crews monitoring it. (The exact mechanism of arriving at a consensus across crews is left as a module.)

Dynamic Nodes. Crucial to implementing a Merkle tree based data structure is the ability to dynamically update the contents of previously existing nodes in the tree, so that they can assimilate newly created cells (even if we did not plan to support the evolution of individual cells). But allowing cell data itself to evolve, in a programmable manner, brings a whole new dimension to distributed data repositories. We point out a couple of aspects.

²Even architectures like IOTA’s Tangle [\[19\]](#), that do not stick to a chain structure use hash pointers in the same direction as blockchains.

While prior constructions have focused on the immutability or “persistence” guarantee of the blocks in a blockchain, for dynamic data we introduce a notion of *consistency*, to assure that a cell has evolved into its current form in accordance with the policies declared by the cell. These policies govern the modification of the cell’s data as well as the policies themselves.

Though persistence is sometimes desirable, many a time it can be a burden on the system. The CellTree design allows individual cells to forget their history (unless they are programmed to retain it – using a suitably programmed cell, a single node could be used to host an entire blockchain).

Excising Malignant Cells. A blockchain, by design, does not allow removing any blocks already accepted to be part of the chain. This (exacerbated by the need to store the entire chain to verify it) creates practical socio-legal complications when illegal data (say, confidential information) is hosted on a blockchain. A CellTree, in contrast, makes it possible to deactivate “malignant” cells, with little impact on the rest of the tree. The node containing a deactivated cell could be brought back to service (with a fresh cell in it), if all the nodes monitoring it cooperate.

3 Preliminaries

We shall use a *binary tree* structure to hold the data cells. Below we introduce some basic notation and definitions, which we will build upon in [Section 5](#).

- Each node v in the tree is uniquely identified using a binary string that encodes the path from the root to v . The root node corresponds to the empty string, denoted by ϵ , and for any binary string v , the strings $v0$ and $v1$ denote the left and right children, respectively, of the node denoted by v .
- We say $u \preceq v$ if u is a prefix of v . In terms of the tree, u is an ancestor of v (possibly $u = v$). We write $u \prec v$ if $u \preceq v$ and $u \neq v$. We shall also write $u \prec_\ell v$ (resp., $u \preceq_\ell v$) to denote that the distance between u and v is at most ℓ and $u \prec v$ (resp., $u \preceq v$).
- Given a set of nodes L that form an antichain (i.e., no node in L is an ancestor of another), and a node u , we say that $u \preceq L$ if $\forall v \in L, u \preceq v$. We say that $u \preceq L$ if $\exists v \in L$ s.t. $u \preceq v$. Note that if $v \preceq L$, then the set of nodes $\{u : v \preceq u \preceq L\}$ forms a subtree rooted at v with L as its set of “terminal” nodes (i.e., nodes with no outgoing edges).

A (full-domain) *hash function* hash takes an input of arbitrary size and maps it deterministically to a fixed size output. We require a *collision resistant* hash function, in which it is infeasible to find two input strings x and y such that $\text{hash}(x) = \text{hash}(y)$. We shall often apply

the hash function to a list of values (e.g., $\text{hash}(x, y)$), implicitly requiring that the inputs to hash are first *unambiguously* encoded to allow such lists.

A *Merkle Tree* [14, 15] is a technique that allows hashing several blocks of data together into a single block, with the ability to show that a certain data block is part of the hash using a succinct proof [12]. We use a variant of Merkle tree in which the data blocks are associated with all the nodes of a tree (rather than just the leaves). We use the following notation.

- η_v - data block at a node v . (In our construction, η_v itself will be a hash of certain data associated with the node v in a CellTree.) For convenience, we shall require that η_v is defined for all nodes (existing in the tree or not).
- h_v is recursively defined as $h_v := \text{hash}(\eta_v, h_{v0}, h_{v1})$ for nodes that exist in the tree. If a node v does not exist in the tree, $h_v := \text{hash}(\eta_v)$, forming the base of the recursion.
- Given an antichain L , and a node $w \preceq L$, h_w can be computed from the set of values

$$\left(\{\eta_u\}_{w \preceq u \preceq L}, \{h_u\}_{w \preceq \text{parent}(u) \preceq L, u \not\preceq L} \right).$$

We shall also employ a standard digital signature scheme given by a key-generation algorithm and a pair of algorithms (sign, verify). We shall typically use it in the form of a *collective signature* created by a set of parties G (publicly represented by their verification keys in the signature scheme). We write $\text{crewSign.sign}(G, m)$ to denote the collective signature on a message m produced by G , $\{\text{sign}_g(m)\}_{g \in G}$; similarly, we write $\text{crewSign.verify}(G, m, \sigma)$ to denote the verification of a collective signature σ on a message m using the public keys of G . crewSign.verify may be programmed with a parameter that specifies how many missing signatures should be tolerated (additional signatures if any, corresponding to public-keys not in G would typically be ignored).

4 Cell Design

Content of each cell consists of two parts, *cell-data* cdata , and a *nucleus* nuc . The nucleus is an entity that typically encodes a fingerprint of the data, and also constrains the trajectory of the data evolution (in the past and the future). The exact functionality of the nucleus is fully customizable and can be used to design a wide variety of cell behaviors, as we explain below.

Typically, the nucleus contains a summary of the cell-data, in addition to some code. The nuclear code consists of algorithms that dictate cell evolution. For an update of a cell evolve from $(\text{cdata}, \text{nuc})$ to $(\text{cdata}', \text{nuc}')$ to be valid, both the nuclei involved should *agree* to the

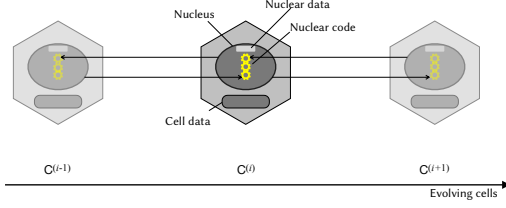


Figure 1 Illustration of cells evolving

The cogwheels indicate the three programs that are part of the nuclear code. The arrows across evolving cells indicate using the next or previous nucleus as inputs to `chkNext` (top cogwheel) and `chkPrev` (bottom cogwheel).

evolution. This is required to enable making inferences about the past and future of a cell from the current nucleus (e.g., to ensure that the data stored in a cell is unmodified and unmodifiable). For efficiency purposes, it will be important that the nuclei should arrive at this agreement without seeing the cell data of the other state (past or future) of the cell, but only the nuclear data. Formally, a nucleus is a pair $(\text{ncode}, \text{ndata})$, where ncode , consists of 3 algorithms (in some well-defined machine model) and ndata holds data for use in these algorithms. These algorithms are for verifying that the cellular data matches the nucleus of the cell (`chkCell`), and for checking if the succeeding and preceding nuclei are compatible with the current nucleus (`chkNext` and `chkPrev`, respectively). Each of these algorithms takes the nucleus it belongs to (code and data), or the entire cell, as its first argument and outputs $b \in \{\text{true}, \text{false}\}$:

$$\begin{aligned} \text{chkCell} &: \text{cell} \mapsto b \\ \text{chkNext} &: (\text{nuc}, \text{nuc}^{(+)}) \mapsto b \\ \text{chkPrev} &: (\text{nuc}, \text{nuc}^{(-)}) \mapsto b \end{aligned}$$

Here, $\text{nuc}^{(+)}$ refers to the new nucleus formed after evolution of nuc , and $\text{nuc}^{(-)}$ refers to the previous nucleus that evolves into nuc . For brevity, we shall write `chkCell(cell)` instead of `cell.nuc.ncode.chkCell(cell)`, `chkNext(nuc, nuc(+))` instead of `cell.nuc.ncode.chkNext(nuc, nuc(+))`, etc.

Initializing an empty cell The nucleus of the empty cell is denoted as nuc_0 . We define $\text{nuc}_0.\text{chkNext}$ to always output `true`, whereas $\text{nuc}_0.\text{chkCell}$ and $\text{nuc}_0.\text{chkPrev}$ always output `false`. This will prevent the acceptance of the empty cell as the result of evolution of any cell.

4.1 Examples of Nucleus Families

We list a few basic examples of nucleus families below. For concreteness, we illustrate their functions `chkPrev`, `chkNext` and `chkCell` in Figure 2.

- **Static data:** `chkPrev` ensures that the previous version of the cell had the empty nucleus. `chkNext` always

returns `false`. This will ensure that a nucleus of this family does not evolve. Then, `chkCell` ensures that the cell-data cannot change either (assuming collision resistance of the hash function).

- **Blockchain Ledger:** A cell can be used to implement an entire blockchain ledger, which is simply a list of data blocks and the only way in which the cell data can be modified is to append blocks.
- **State Machine:** A nucleus can be used to implement a state machine, so that transitions can occur only according to the edges in a directed graph G over a state space. (The state space can be infinite, as in the case of a counter.) The simple version illustrated in Figure 2 does not use any `cdata` and there is no additional information in the nuclear data other than the state.

Private Nucleus. Since the nucleus is available outside of the node containing the cell (to the nodes monitoring it), if the contents of the cell are to be protected, the nucleus should retain the cell’s secrecy. We point out the possibility that the nuclear data can be kept “encrypted” while still making it possible to run the verification algorithms `chkPrev` and `chkNext`. For simple nuclear families this can indeed be achieved efficiently: For instance, for static data or the blockchain ledger above, instead of using the hash of the blocks, a committed hash can be used (e.g., in the random oracle model, appending a random block to the data block before hashing suffices). In the case of the blockchain ledger, we need to also allow virtual datablocks that contain a variable number of blocks (possibly none), so that exactly one virtual block is added during each update, and the actual number of blocks can be hidden. For more complex nuclear functionality, (succinct, non-interactive) zero-knowledge proofs can be used to let `chkPrev` and `chkNext` verify that the updates are valid without revealing anything else.

External References. The machine model for the nuclear code is implemented using a module `exec`. Apart from standard operations, this machine model may allow references to certain external resources, like current time, cells in other nodes of this CellTree, or even data blocks or code in various blockchains. As a simple example, this allows one to use a timestamp included in a nucleus, and let `chkNext` accept a new nucleus only if its timestamp is in the past, but later than the timestamp of the current nucleus. Ability to refer to data in other nodes of the CellTree allows a nucleus, for instance, to prove its relation to cells in other nodes (e.g., the nucleus of a cell containing a ledger can prove that a transaction being added to the cell is also present in a global ledger, whose Merkle hash is stored as a cell in a different node).

Static Data: The nuclear data ndata is simply a hash of the cellular data cdata using a collision-resistant hash function.

$$\begin{aligned}\text{chkCell}(\text{cell}) &= \begin{cases} \text{true} & \text{if } \text{cell.nuc.ndata} = \text{hash}(\text{cell.cdata}) \\ \text{false} & \text{otherwise} \end{cases} \\ \text{chkPrev}(\text{nuc}, \text{nuc}^{(-)}) &= \begin{cases} \text{true} & \text{if } \text{nuc}^{(-)} = \text{nuc}_0 \\ \text{false} & \text{otherwise} \end{cases} \\ \text{chkNext}(\text{nuc}, \text{nuc}^{(+)}) &= \text{false}\end{aligned}$$

Blockchain Ledger: cdata is an array of blocks ($|\text{cdata}|$ refers to the number of blocks in the array and $\text{cdata}[i]$ refers to the i^{th} block). ndata is a triple $(n, \text{blkHash}, \text{chnHash})$ where n is the number of blocks in the ledger and the other two items are hash values (the last block and the root in a Merkle-chain).

$$\begin{aligned}\text{chkCell}(\text{cell}) &= \begin{cases} \text{true} & \text{if } \text{cell.nuc.ndata} = (n, \gamma_n, \text{hash}(\text{hash}(\dots \text{hash}(\text{hash}(\gamma_1, \gamma_2), \gamma_3), \dots), \gamma_{n-1}), \gamma_n))) \\ & \text{where } n = |\text{cell.cdata}| \text{ and } \gamma_i = \text{hash}(\text{cell.cdata}[i]) \\ \text{false} & \text{otherwise} \end{cases} \\ \text{chkNext}(\text{nuc}, \text{nuc}^{(+)}) &= \begin{cases} \text{true} & \text{if } \text{nuc}^{(+).ndata.chnHash} = \text{hash}(\text{nuc.ndata.chnHash}, \\ & \text{nuc}^{(+).ndata.blkHash}), \\ & \text{nuc}^{(+).ncode} = \text{nuc.ncode} \text{ and} \\ & \text{nuc}^{(+).ndata.n} = \text{nuc.ndata.n} + 1 \\ \text{false} & \text{otherwise} \end{cases} \\ \text{chkPrev}(\text{nuc}, \text{nuc}^{(-)}) &= \begin{cases} \text{true} & \text{if } \text{nuc}^{(-)} = \text{nuc}_0, \text{nuc.ndata.n} = 1 \\ & \text{and } \text{nuc.ndata.chnHash} = \\ & \text{nuc.ndata.blkHash} \\ \text{nuc.chkNext}(\text{nuc}^{(-)}, \text{nuc}) & \text{otherwise} \end{cases}\end{aligned}$$

State Machine: A directed graph $G = (V, E)$ over the set of states V is encoded in the algorithms below. ndata is simply a state in V . state_0 refers to the start state of the machine.

$$\begin{aligned}\text{chkCell}(\text{cell}) &= \begin{cases} \text{true} & \text{if } \text{cell.nuc.ndata} \in V \text{ and } \text{cell.cdata} \text{ is empty} \\ \text{false} & \text{otherwise} \end{cases} \\ \text{chkNext}(\text{nuc}, \text{nuc}^{(+)}) &= \begin{cases} \text{true} & \text{if } (\text{nuc.ndata}, \text{nuc}^{(+).ndata}) \in E \text{ and} \\ & \text{nuc}^{(+).ncode} = \text{nuc.ncode} \\ \text{false} & \text{otherwise} \end{cases} \\ \text{chkPrev}(\text{nuc}, \text{nuc}^{(-)}) &= \begin{cases} \text{true} & \text{if } \text{nuc}^{(-)} = \text{nuc}_0 \text{ and} \\ & \text{nuc.ndata} = \text{state}_0 \\ \text{nuc.chkNext}(\text{nuc}^{(-)}, \text{nuc}) & \text{otherwise} \end{cases}\end{aligned}$$

Figure 2 Examples of nucleus families

In each case the three algorithms in the nuclear code are described.

5 Tree Design

In order to maintain the cells in a single ledger, we shall use a tree structure. Each cell will be exclusively associated with a unique node in the tree, and will be implemented by a set of parties – called the node’s *crew*. As time progresses, new cells can be added to the tree and also, each cell may evolve as permitted by its nuclear code.

To make the tree robust to corruption of some cells, we shall have each cell’s evolution monitored by several other cells (more specifically, by several cells that are its ancestors in the tree). Also, a cell’s evolution is not confirmed until it is *assimilated* into the tree. In this section we describe the mechanisms used for cell evolution, assimilation and access. We begin with the structure and contents of a CellTree, before describing the procedures involved.

Nodes of a CellTree. In a CellTree, the location for a cell is a node in a binary tree. The node address, v , is a binary string representation of the path to the node from the root. For each node v of the tree, there shall be a set of parties G_v — the *crew* operating the node v — who implement the protocols associated with that node, and maintain a copy of the contents of the node. Each member of a crew, $g \in G_v$ has a signing/verification key pair $(\text{SK}_g, \text{VK}_g)$ associated with it, and can be uniquely identified by VK_g (or alternately, a hash thereof). For convenience, we shall denote the collection $\{\text{VK}_g | g \in G_v\}$ by $\widehat{\text{VK}}_v$.

Merkle Multi-Trees. We shall use a variant of a Merkle tree in which the underlying graph structure is that of a *multi-tree* rather than a tree. A multi-tree is a directed acyclic graph such that the set of nodes reachable from any node (via the directed edges) is a tree rooted at the node. While there is a unique tree

structure for a CellTree with every node appearing at a fixed location in the tree, the reason we need to use a Merkle Multi-Tree is that the cells in each node can evolve over time, and each version of a cell appears as a separate node in the multi-tree (but has the same address in the binary tree structure).

Recall that in [Section 3](#), we defined a Merkle tree as having data η_u at each node u , and recursively defined a hash pointer for the node $h_u = \text{hash}(\eta_u, h_{u0}, h_{u1})$. But in a Merkle multi-tree, the same binary tree node u is allowed to have different nodes with different values for η_u, h_{u0}, h_{u1} . We denote such a tuple as mNode_u and it is uniquely addressed by h_u . More precisely,

$$\text{mNode}_u = (\text{nuc}_u, \text{nonce}_u, h_{u0}, h_{u1}). \quad (1)$$

where $(\text{nuc}_u, \text{nonce}_u)$ is the actual data associated with mNode_u , and we let $\eta_u = \text{hash}(\text{nuc}_u, \text{nonce}_u)$.

We shall often refer to a “Merkle subtree” that forms a subgraph of a Merkle multi-tree. The Merkle subtree that is reachable from a node u , and terminated by a set of terminals L is recursively defined as follows.

$$\text{mTree}_L^u = \begin{cases} (\text{mNode}_u, \text{mTree}_L^{u0}, \text{mTree}_L^{u1}) & \text{if } u \preceq L \\ \perp & \text{otherwise} \end{cases}$$

Here, mTree_L^u is uniquely specified by the tuple (u, L, h_u) , where h_u is the hash pointer to mNode_u . We let $\text{mTree}_L^u.\text{terminals}$ denote the set of terminals L .

For $u \preceq v \preceq L$, one can *extend* a Merkle subtree rooted at v and terminated at L to one rooted at u and terminated at L , by concatenating it with a Merkle subtree rooted at u and terminated at $\{v\}$. We abbreviate $\text{mTree}_{\{v\}}^u$ as mTree_v^u below. Then, we can denote the above extension as $\text{mTree}_L^u \leftarrow \text{mTree}_v^u \parallel \text{mTree}_L^v$, where the concatenation operator replaces $\text{mTree}_v^v = \perp$ (present in mTree_v^u) with mTree_L^v .

mEval denotes the algorithm that, on input mTree_L^u computes h_u or returns an error, if the input is not a valid Merkle subtree (e.g., if the hash pointers included in a node do not match the values computed for the subtrees of that node).

Proof of Assimilation. A proof of assimilation of a node v , with respect to a *root of assimilation* aroot has the following form:

$$\text{poa}_v^{\text{aroot}} = (\text{mTree}_v^{\text{aroot}}, \sigma^{\text{aroot}}), \quad (2)$$

where $\text{mTree}_v^{\text{aroot}}$ is a Merkle tree rooted at aroot with $\{v\}$ as the terminal set, and σ^{aroot} is a signature of the root hash of this tree, signed by the crew of aroot .

Contents of a node. A node stores the following items (accessed via methods of a module `store`).

- The latest version of its cell. (Updated using method `updateCell` and retrieved using `currentCell`.)
- A list of cells that have been propagated rootwards for assimilation (possibly stored in a compressed fashion), with an associated *hash pointer* (stored using the method `addcell`). Any proofs of assimilation received for each cell are also stored (via `addpoa`).
- An ordered list `UUlist` consisting of nuclei of cell versions since the last rootward propagation. (Items are appended to this list using the method `appendUU`, and the list is retrieved (and emptied) using `flushUU`.)
- A Merkle multi-tree with roots of the form mNode_v , one for each version of the cell that is propagated rootward. (Subtrees of this multi-tree are created and retrieved using methods `addmTree` and `getmTree`).

CellTree Procedures. In [Section 5.1](#) below, we shall specify the following core procedures as part of the CellTree framework:

- **Reading a Cell:** The procedure `READ` will be invoked by a client to read a cell from a node.
- **Cell Evolution:** The procedure `EVOLVE` is invoked by the crew of a node to update its cell.
- **Cell Assimilation:** There are two complementary procedures `ROOTWARDPROPAGATION` and `LEAFWARDPROPAGATION` that are invoked by the crew of a node to communicate with crews of other nodes, during the assimilation process.

In addition to these procedures, the CellTree framework relies on other procedures (e.g., for scheduling the invocation of various procedures, for creating a new node and assigning it a crew, for deciding what a cell should evolve into, etc.), that are all implemented as modules (see [Section 5.2](#) and [Section 6](#)).

Robustness Properties. We define the following desirable properties for a CellTree. These properties are parametrized by a node v , since, even if parts of the tree are malfunctioning, we seek to guarantee these properties to other nodes.

- **consistency(v):** Let C be the set of cells returned by a set of successful invocations of `READ(v)` by honest clients. Then, there exists a sequence of nuclei $\text{nuc}^{(0)}, \dots, \text{nuc}^{(N)}$ such that for all $i \in [1, N]$, both `chkNext($\text{nuc}^{(i-1)}, \text{nuc}^{(i)}$)` and `chkPrev($\text{nuc}^{(i)}, \text{nuc}^{(i-1)}$)` hold and further, for all cell $c \in C$, there exists $i \in [0, N]$, such that `cell.nuc = $\text{nuc}^{(i)}$` and `chkCell(c)` holds.
- **correctness(v):** Any cell returned by `READ(v)` to an honest client is equal to a cell assigned to v using the procedure `EVOLVE` by the crew G_v prior to that (or is the empty cell that every node is initialized with).

- **liveness(v)**: If a cell is assigned to a node v at any point in time (by the crew G_v using the procedure **EVOLVE**), eventually every invocation of **READ(v)** by any honest client will return this cell, or a cell assigned to v subsequently.

Note that the above definitions refer to the crew G_v , but the “correct” G_v is not necessarily well-defined if the crews of the CellTree are malicious. A fully formal definition of the above properties requires a formal definition of the correct crew G_v , which we provide in [Section 6.1](#).

5.1 CellTree Procedures

Now we list the core procedures of the CellTree architecture.

5.1.1 Reading a Cell

To read a cell in a node v , a client executes the procedure **READ** presented below.

```

procedure READ( $v$ ,  $\text{aroot}$ )  $\triangleright \text{aroot} \preceq v$ 
  ( $G'_v, G'_{\text{aroot}}$ )  $\leftarrow$  discover( $v$ ,  $\text{aroot}$ )
  ( $\text{cell}, \text{poa}$ )  $\leftarrow$  fetch.getCell( $G'_v, v$ ,  $\text{aroot}$ )
  VERIFY( $v$ ,  $\text{aroot}$ ,  $G'_{\text{aroot}}$ ,  $\text{cell}$ ,  $\text{poa}$ )
  return  $\text{cell}$ 

procedure VERIFY( $v$ ,  $\text{aroot}$ ,  $G_{\text{aroot}}$ ,  $\text{cell}$ ,  $\text{poa}_v^{\text{aroot}}$ )
  ( $\text{mTree}, \sigma$ )  $\leftarrow$   $\text{poa}_v^{\text{aroot}}$ 
   $h \leftarrow$  MEVAL( $\text{mTree}$ )
  assert  $\text{mTree.terminals} = \{v\} \wedge \text{mTree.root} = \text{aroot}$ 
  assert crewSign.verify( $G_{\text{aroot}}, h, \sigma$ )
  assert  $\text{cell.nuc} = \text{mTree.mNode.nuc}$ 
  assert chkCell( $\text{cell}$ )

```

The client can use any node $\text{aroot} \preceq v$ as the *assimilation root* where it requires that the cell being read has to be assimilated. The client first discovers the crew for the nodes v and aroot , using a module **discover**. (A typical implementation of this module would require the client to already “know” the crew of one or more nodes in the CellTree, and would then carry out a search, relying on the fact that crews of a node are expected to know crews of the nodes they monitor and the crews monitoring their node. An example appears in [Section 6](#).)

The algorithm **READ** also uses a module **fetch** to let a client access the information about a cell from a node’s crew (typically, the latest version of the cell for which the crew possesses a proof of assimilation at the specified assimilation root). The client verifies the proof of assimilation (which certifies a nucleus as being assimilated), ensures that the given cell’s nucleus matches the one in the proof of assimilation and also verifies that the cell contents are consistent with the nucleus.

5.1.2 Cell Evolution

The cell evolution procedure **EVOLVE** is invoked at each crew member of a node by the module **selectCell**, to replace the current cell with a new cell $\text{cell}^{(+)}$ that the module has chosen. We expect that **selectCell** achieves a consensus among the crew members before **EVOLVE** is invoked; on the other hand, the module **store** typically acts locally at each crew member’s storage, without requiring any coordination among the crew members.

```

procedure EVOLVE( $v$ ,  $\text{cell}^{(+)}$ )
   $\text{cell} \leftarrow$  store.currentCell
  assert chkNext( $\text{cell.nuc}$ ,  $\text{cell}^{(+).nuc}$ )
  assert chkPrev( $\text{cell}^{(+).nuc}$ ,  $\text{cell.nuc}$ )
  assert chkCell( $\text{cell}^{(+)}$ )
  store.updateCell( $\text{cell}^{(+)}$ )
  store.appendUU( $\text{cell}^{(+).nuc}$ )

```

5.1.3 Cell Assimilation

Assimilation is carried out by propagating *assimilation signals* from cells towards the root (Rootward Propagation) and *proofs of assimilation* back from the root to all the cells (Leafward Propagation). This process is designed such that the amount of work a node needs to carry out does not grow as the tree grows. However, we do let the amount of storage required at a node to depend on the *depth of that node*, as a version of the cell is retained at least until all proofs of assimilation for a later version are received.

Rootward Propagation: We describe the procedure invoked by the crew of a node to propagate assimilation signals to its monitoring nodes in the tree. The requisite communication between the node’s crew and the crews monitoring it is accomplished through a module **rootward**. This module may use direct communication between the crews of a monitoring and a monitored node, or alternately, use an indirect route (e.g., each node communicating directly only with its parent and children).

During rootward propagation, an *assimilation signal* from each node u is propagated to every node monitoring it. This message contains the following:

$$(\text{UList}_u, \text{mNode}_u)$$

where mNode_u is as defined in [Equation 1](#). We require that nuc_u in mNode_u be the same as the last entry in UList_u (stored together).

In the procedure **ROOTWARDPROPAGATION**, each crew member accepts and verifies a list of nuclear updates from each node it is monitoring. Here, verification involves verifying the consistency between consecutive updates (using **chkNext** and **chkPrev** methods in the two

```

procedure ROOTWARDPROPAGATION( $v$ )
   $L \leftarrow \text{rootward.monitored}$ 
  for  $b \in \{0, 1\}$  do
     $\text{mTree}_L^{vb} \leftarrow \text{ACCEPTVERIFY}(vb, L)$ 
     $h_{vb} \leftarrow \text{MEVAL}(\text{mTree}_L^{vb})$ 
   $\text{UList}_v \leftarrow \text{store.flushUU}$ 
   $\text{cell}_v \leftarrow \text{store.currentCell}$ 
   $\text{nonce}_v \leftarrow \text{generateNonce}$ 
   $\text{mNode}_v \leftarrow (\text{cell}_v.\text{nuc}, \text{nonce}_v, h_{v0}, h_{v1})$ 
  if  $v \neq \epsilon$  then
     $\text{rootward.send}(\text{UList}_v, \text{mNode}_v)$ 
   $\text{mTree}_L^v \leftarrow (\text{mNode}_v, \text{mTree}_L^{v0}, \text{mTree}_L^{v1})$ 
   $h_v \leftarrow \text{hash.eval}(\text{hash.eval}(\text{cell}_v.\text{nuc}, \text{nonce}_v), h_{v0}, h_{v1})$ 
   $\sigma_v \leftarrow \text{crewSign.sign}(h_v)$ 
   $\text{poa}_v^v \leftarrow ((\text{mNode}_v, \perp, \perp), \sigma_v)$ 
   $\text{leafward.send}(\text{poa}_v^v)$  ▷ send to self
   $\text{store.addCell}(\text{cell}_v, h_v)$ 
   $\text{store.addmTree}(h_v, \text{mTree}_L^v)$ 

procedure ACCEPTVERIFY( $u, L$ )
  if  $u \not\lesssim L$  then
    return  $\perp$ 
   $(\text{UList}_u, \text{mNode}_u) \leftarrow \text{rootward.receive}(u)$ 
   $(\text{nuc}_u, \text{nonce}_u, h_{u0}, h_{u1}) \leftarrow \text{mNode}_u$ 
   $\text{lastmTree}_L^u \leftarrow \text{store.getmTree}(L, u)$  ▷ latest version
   $\text{nuc}_u^{(0)} \leftarrow \text{lastmTree}_L^u.\text{mNode}.\text{nuc}$ 
   $\text{nuc}_u^{(1)}, \dots, \text{nuc}_u^{(t)} \leftarrow \text{UList}_u$ 
  if not CONSISTENT( $\text{nuc}_u^{(0)}, \dots, \text{nuc}_u^{(t)}$ ) then
    return  $\text{lastmTree}_L^u$  ▷ reject update
  for  $b \in \{0, 1\}$  do
     $\text{mTree}_L^{ub} \leftarrow \text{ACCEPTVERIFY}(ub, L)$ 
    if  $ub \lesssim L$  then
      if  $\text{MEVAL}(\text{mTree}_L^{ub}) \neq \text{mNode}_u.h_{ub}$  then
        return  $\text{lastmTree}_L^u$  ▷ reject update
    return  $(\text{mNode}_u, \text{mTree}_L^{u0}, \text{mTree}_L^{u1})$ 

procedure CONSISTENT( $\text{nuc}_u^{(0)}, \dots, \text{nuc}_u^{(t)}$ )
  for  $j \in 1, \dots, t$  do
     $x \leftarrow \text{chkNext}(\text{nuc}_u^{(j-1)}, \text{nuc}_u^{(j)})$ 
     $y \leftarrow \text{chkPrev}(\text{nuc}_u^{(j)}, \text{nuc}_u^{(j-1)})$ 
    if  $x = \text{false}$  or  $y = \text{false}$  then
      return  $\text{false}$ 
  return  $\text{true}$ 

```

nuclei), as well as ensuring that the hash pointers included in an updated node points to valid cells that resulted from verified updates. (In case of verification failure, the previous version of the node is retained.) Then, the crew prepares a node in the Merkle multi-tree, with the current version of its cell, with hash pointers to the verified versions of the cells of its two child nodes. This information is propagated to crews monitoring it, along with a the local list of updates since the last rootward propagation.

Note that the recursive subroutine ACCEPTVERIFY could result in the updates of an entire subtree rooted at u being rejected, if the update received from the node u itself fails verification. Further, a faulty update for u could result in updates from its parent (and recursively, all its ancestors till v) being rejected, if the parent of u used a hash pointer that does not match any hash pointer

from u . This may seem to suggest that a malicious crew for a single node can disrupt all updates from its ancestors and descendents. However, such attacks can be avoided by well-designed **rootward** and **leafward** modules. In particular, in the simple instantiation in Section 6, the descendents of a node cannot block the nodes updates from being assimilated (though it allows a node from blocking its descendents' updates).

Apart from the module **rootward**, this procedure uses a few other modules: A module **monitoring** is used to obtain the set (subtree) of nodes v monitors: this set is specified in terms of an antichain L , as $\{u | v \prec u \lesssim L\}$. The module **store** is used to save the Merkle subtree corresponding to the update that was propagated rootward. The module **crewSign** is used to produce a collective signature that is part of the proof of assimilation of the child nodes' cells. This proof is saved using the module **leafward** and recovered later during leafward propagation.

Leafward Propagation Step: We describe the procedure used by the crew of a node to propagate proof of assimilation from an ancestor node **aroot** (possibly itself) to one or more of its subtree nodes (see Equation 2). The procedure uses a module **leafward** to collect the proofs and to propagate them to a subtree (after due verification).

The set of nodes for which a node's crew should prepare the proof is also determined by the module **leafward**: The method **pickleaves** returns a set of terminal nodes L and the proof of assimilation $\text{poaset}_L^{\text{aroot}}$ is prepared for all nodes u s.t. $v \preceq u \lesssim L$. This set of proofs is transmitted to these nodes using the method **leafward.send**.

```

procedure LEAFWARDPROPAGATION( $v$ )
  ▷ Process a proof received
   $(\text{aroot}, G_{\text{aroot}}, \text{poa}_v^{\text{aroot}}) \leftarrow \text{leafward.receive}$ 
   $(\text{mTree}_v^{\text{aroot}}, \sigma^{\text{aroot}}) \leftarrow \text{poa}_v^{\text{aroot}}$ 
   $h_v \leftarrow \text{MEVAL}(\text{mTree}_v^{\text{aroot}})$ 
  assert  $\text{crewSign.verify}(G_{\text{aroot}}, h_v, \sigma^{\text{aroot}})$ 
   $\text{store.addpoa}(\text{poa}_v^{\text{aroot}}, h_v)$ 
   $L \leftarrow \text{leafward.pickleaves}(\text{aroot})$  ▷ propagate till  $L$ 
   $\text{mTree}_L^v \leftarrow \text{store.getmTree}(L, v, h_v)$ 
  if  $\text{aroot} \neq v$  then
     $\text{mTree}_L^{\text{aroot}} \leftarrow \text{mTree}_v^{\text{aroot}} || \text{mTree}_L^v$ 
   $\text{poaset}_L^{\text{aroot}} \leftarrow (\text{mTree}_L^{\text{aroot}}, \sigma^{\text{aroot}})$ 
   $\text{leafward.send}(\text{poaset}_L^{\text{aroot}})$ 

```

5.2 Modules

A CellTree relies on several modules, but is oblivious to their implementation. In this section, we collect all the modules referred to above, as well as some additional modules that are used by a CellTree deployment. While

many of the modules are entirely local to a node’s crew, some modules need to coordinate across multiple nodes.

- **Node Creation and Crew Selection.** Creating a new node and selecting a crew for it are tasks carried out collectively by the nodes which would be monitoring the new node, using a module `selectCrew`. This module could use, e.g., the committee selection protocol used within Algorand [8].

- **Cell Selection.** As mentioned before, cell evolution is triggered by a module `selectCell`, which is responsible for obtaining a consensus among the crew members as to the next version of a cell.

- **Propagation.** As described earlier, the modules `rootward` and `leafward` are used by crews to send and receive assimilation signals and proofs of assimilation. Typically, the `send` and `receive` methods operate in the background and would involve consensus mechanisms and secure communication protocols used within and across crews. The `schedule` invokes the assimilation procedures.

- **Local Storage.** The module `store` is used by each party in a crew to locally store and retrieve various values across separate invocations of the algorithms. Typically, this module requires no communication among crew members (as the consistency guarantees of values being stored are ensured by the other modules).

- **Client Access to a CellTree.** Modules `fetch` and `discover` are used by the procedure READ. These modules may implement access control and denial-of-service protection.

The proofs of assimilation involve a collective signature by the crew (on a hash value). The protocol for creating such signatures and the algorithm for locally verifying them are encapsulated in the module `crewSign`.

- **Code Execution.** The machine model used to execute the nuclear code is specified as a module `exec`. The module may support multiple languages and library functions.

- **Hashing.** The hash algorithm and Merkle tree evaluation are implemented by the module `hash`. A method `generateNonce` included in this module can be used to create nonces that control the hash evaluation (e.g., if a time-stamp in the nonce is in the future, the hash evaluation could return an error).

- **Communication.** The `communication` module is a low level module that manages addresses of crew members and is accessed while running the various protocols in other modules.

- **Selection and Evolution of Modules.** The CellTree framework admits different implementations of the various modules to coexist. But the framework does not specify how modules are chosen, and possibly changed over time, instead delegating it to a module `moduleManager` (which governs its own evolution).

6 Instantiating a CellTree: An Example and Analysis

We sketch an instantiation of a CellTree, using a set of relatively simple modules, and give *provable* robustness guarantees for it (in Section 6.1). This instantiation, which we dub \mathbf{CT}_0 , is only for illustrating the robustness and performance guarantees possible in a CellTree; several improvements are discussed in Section 6.3. A proof-of-concept for the same can be found [here](#).

We sketch the design of \mathbf{CT}_0 below. This design can be implemented entirely using appropriately defined modules in the CellTree architecture above.

- Every node monitors all its descendent nodes at a distance of at most ℓ , which is a fixed parameter for the entire system.
- During the propagation steps, each node’s crew interacts directly only with the crews of its two children and its parent. That is, the crew of a node v has all of the assimilation signals $\{\text{rootward.receive}(u)\}_{vb \leq u}$ (for $b = 0$ or 1) communicated to it by the crew of the node vb ; the method `receive` requires that the assimilation signal it receives for a node is signed by the crews of all the nodes monitoring it; for a signature from a crew, a strict majority of the signatures by individual crew members should verify. When a node v invokes `leafward.pickleaves` it returns $L = \{v0, v1\}$.
- Only the root node ϵ is used as an assimilation root, and hence proofs of assimilation are not created by the other nodes during rootward propagation.
- The module `discover` implements a simple search algorithm (shown below) to iteratively discover the crews of each node in the path from the root to the node to be read. For simplicity, we assume that the crew of the root node G_ϵ is fixed and hardcoded into the `discover` module used by all nodes (at the time of creating each node). The algorithm is designed to prevent errors in discovery, provided that the root crew is “good” and there is a “good” node in every ℓ long path. Once a crew is discovered, a client can cache it.
- For scheduling the propagation algorithms, we use an “epoch period” of Δ (say half an hour) and a “phase shift” $\Phi = \Delta/t$ for some integer t (say, $\Phi = 5$ minutes) as follows: A node at depth d invokes the rootward propagation algorithm periodically, at time T if $T \equiv d\Phi \pmod{\Delta}$.³ This allows the assimilation sig-

³The crew members use roughly synchronized clocks, using a standard protocol like NTP, to determine time. Φ would be chosen to be comfortably larger than the time needed for a crew to carry out the `ROOTWARDPROPAGATION` procedure (up to `rootward.send`) plus possible delays due to clock skews, and Δ would be chosen

nals from nodes at depth d to cascade to the root, with a latency of at most $\Delta + d\Phi$. The root node invokes the procedure `LEAFWARDPROPAGATION` after executing `ROOTWARDPROPAGATION`; in the other nodes, the procedure `LEAFWARDPROPAGATION` is triggered when it receives a proof of assimilation from its parent's `leafward` module.

- We do not fully specify the policy used by the module `selectCrew` to find the crew for a node being created. But we shall require that it uses a final consensus step to certify the selected crew. A party (which considers itself as monitoring the newly created node) accepts the selected crew if it is certified by the signatures of a strict majority among the *multiset union*⁴ of (what the party considers as) the crews of all the nodes that will be monitoring the node being created. This module installs the information regarding the selected crew in all the crews monitoring it, and vice versa.

```

procedure discover( $v$ , aroot)
   $G'_\epsilon \leftarrow G_\epsilon$   $\triangleright G_\epsilon$  is hardcoded
  for all  $u$  s.t.  $\epsilon \prec u \preceq v$  do
    if  $\exists G$  s.t.  $\{G_u^{(w)} : \epsilon \preceq w \prec_\ell u\} = \{G\}$  then
       $G'_u \leftarrow G$ 
    else
      abort
    for all  $w$  s.t.  $u \prec w \preceq_\ell u$  do
       $G_w^{(u)} \leftarrow \text{fetch.getCrew}(G'_u, w)$ 
  return  $(G'_v, G'_{\text{aroot}})$ 

```

6.1 Robustness Analysis

First, we formalize the robustness properties of consistency, correctness and liveness (informally introduced earlier), and present simple assumptions under which these properties are guaranteed by **CT₀**. Recall that the definitions of robustness properties correctness and liveness referred to *the* crew G_v . To formalize this, we shall first define a notion of a node v being *good*, and then we shall define a crew for v .

To define a good node, we shall refer to the set of parties that a party g *considers as the crew* G_u for a node u (such information is installed in g by the module `selectCrew`). Note that a party may not consider any set to be the crew of u (if it is not monitoring or being monitored by u 's crew, or if the node u has not yet been created).

⁴In the resulting multiset a party has weight equal to the number of monitoring crews in which it appears. Majority refers to the weighted majority.

Definition 1 For sets of parties G and \hat{G} and a node u , we say that G considers \hat{G} as G_u if $\exists H \subseteq G$, with $|H| > |G|/2$ and $\forall g \in H$, g is honest and considers \hat{G} as G_u .

Definition 2 We say that a node v is good with respect to a set of parties G_ϵ if there exist a sequence of nodes $\epsilon = v_0 \prec v_1 \prec \dots \prec v_n = v$, and sets of parties $\{\hat{G}_{v_i}\}_{i=0}^n$, such that the following hold:

- $\hat{G}_{v_0} = G_\epsilon$;
- $\forall i \in [0, n]$, \hat{G}_{v_i} considers itself as G_{v_i} ;
- $\forall i \in [0, n]$, \hat{G}_{v_i} considers $\hat{G}_{v_{i+1}}$ as $G_{v_{i+1}}$.

Note that above, each \hat{G}_{v_i} is required to have an honest majority (as in the definition of “considers”). Also, any node being good with respect to G_ϵ , implies that the root node ϵ itself is good.

Recall that the definitions of robustness properties correctness and liveness referred to *the* crew G_v . To formalize these definitions, we shall restrict the definition of these two properties to only good nodes v (with respect to a given G_ϵ), and replace the reference to G_v by (any) crew \hat{G}_v guaranteed by the above definition of v being good.

We shall state our guarantees for the CellTree **CT₀** in terms of the following pre-conditions:

Definition 3 For a node v and a set of parties G_ϵ :

- $A_v(G_\epsilon)$ states that the node v is good with respect to G_ϵ .
- $A_v^i(G_\epsilon)$ states that in any set of i consecutive nodes in the path from the root node ϵ to v (or if no such set exists, then in the entire path) there is at least one node that is good with respect to G_ϵ .

We point out that these assumptions are made *per node*, so that it may hold for some nodes and not others. The assumption of a node being good requires only an honest majority among the crew members, and may arguably be made for *every* node in the CellTree, if `selectCrew` module is implemented is using, say, a scheme like that in Algorand [8], under similar assumptions as used there.

Now we are ready to state the guarantees for **CT₀** (a proof sketch is given in [Appendix C](#)).

Theorem 1 In **CT₀**, let G_ϵ be the root node crew used by the `discover` module. Then,

1. $A_v^\ell(G_\epsilon) \Rightarrow \text{consistency}(v)$,
2. $A_v^\ell(G_\epsilon) \wedge A_v(G_\epsilon) \Rightarrow \text{correctness}(v)$, and
3. $A_v^1(G_\epsilon) \Rightarrow \text{liveness}(v)$.

Our guarantees for a node v only depend on (some of) the nodes in the path from the root to v being good (even

though it may be reasonable to assume that every node is good, by using sufficiently large crews). Not only is this an exponentially small fraction of all the nodes created in the tree (assuming a well-balanced tree), but also it is a set of nodes that are known at the time of creating the node v . Nodes created in the future (even descendants of v) have no effect on any of these guarantees, including liveness.

6.2 Performance Analysis

The CellTree architecture is designed to be scalable, so that a node's communication, computation and storage stay bounded even as the tree grows. The only parameter of the size of the tree that affects a node's complexity is the depth of the node itself (due to the need to store old versions of a cell until newer versions are assimilated at the root), which does not change once the node is created.

Latency: \mathbf{CT}_0 does not exploit the multi-level confirmation feature of the CellTree architecture, as it is the root node which issues all proofs of assimilation. As mentioned earlier, an update at a node at a depth d takes at most $\Delta + d\Phi$ time to reach the root, where Δ is the epoch period and Φ is the phase shift parameter. During the leafward propagation step, at each node LEAFWARDPROPAGATION is triggered right after the parent node has finished their leafward propagation procedure. If the time taken for the procedure to run is at most δ (typically $\delta \ll \Phi$), the proof of assimilation from the root arrives at the node within $\Delta + d(\Phi + \delta)$ time, after an update is made locally at the node.

Storage: Each crew member of a node at depth d stores the following: (1) a single nucleus for each of the nodes monitored by it, (2) the last cell (of its own node) for which a proof of assimilation was received, all the cells whose nuclei have been propagated rootwards but whose proofs of assimilation have not yet arrived, and the current cell, (3) a single proof of assimilation from the root consisting of $O(d)$ hash values and a single crew-signature, and (4) the nuclei of all the local updates since the last rootward propagation. At steady state, the total storage needed can be bounded as follows. Suppose α is a bound on the nucleus size (plus a nonce and the nuclear hash), and β on the cell size, (typically $\alpha \ll \beta$), r is the rate of updates, and L the latency from above. Then the storage required (apart from a proof of assimilation)

$$(2^\ell + rL)\alpha + \left(\frac{L}{\Delta} + 1\right)\beta = d(r\Phi'\alpha + \phi\beta) + c,$$

where $\Phi' = \Phi + \delta \approx \Phi$ and $\phi = \Phi'/\Delta$ is a small constant (e.g., $\phi = 0.2$), and c is independent of the depth. Note that while the storage does depend on the depth, it does so in a moderate way, with the dominant term being $d\phi\beta$.

Communication: A crew member communicates with its peers in the crew, as well as with members of the crews of nodes that monitor it or that it monitors. For concreteness, consider consensus mechanisms which have an initial phase when the actual data is communicated to each party, and a second confirmation stage that uses signatures on hashes of the data; in the optimistic (and typical) case where the confirmation stage accepts the original data, the communication phase need not be repeated. In this case, the total communication incident on each crew member is largely independent of the size of the crews. So we shall focus on the amount of data communicated between *nodes*, rather than the communication incurred by the crew members.

In the assimilation steps, each node communicates directly only with its children and parent nodes. Per epoch, each node receives at most $O(2^\ell r \Delta \alpha)$ bits from the children, and $O(d)$ bits as part of proof of assimilation. The node will also incur communication costs during reads and in other processes it is part of (including `selectCell` for itself and `selectCrew` for nodes it monitors).

The only potential communication bottleneck in the protocol is in the module `discover`. Crew information of nodes monitored by the root (and nodes close to the root) would be needed by almost all the clients in the system. However, the information that different clients recover from the root (and nodes close to it) is essentially the same, and can be periodically *published* by the root. Then, the `discover` module obtains the crew information of the nodes within (say) 2ℓ distance of the root from information published by the root, and if necessary, carries out a tree traversal starting with a node ℓ away from the root.

6.3 Beyond \mathbf{CT}_0

While \mathbf{CT}_0 provides reasonable guarantees, it leaves room for much improvement. We discuss several such possibilities below. Mostly these improvements can be implemented by changing the modules.

Improving Liveness. In \mathbf{CT}_0 , $\text{liveness}(v)$ depends on each ancestor node being good (having an honest majority crew). A single bad ancestor can block v 's updates from propagating to the root. This can be remedied at the expense of using more complex `rootward` and `leafward` modules, that allow parties to directly communicate with the crews of all the nodes they monitor or which monitor them. This can be used to replace the assumption A_v^1 used by \mathbf{CT}_0 for $\text{liveness}(v)$ to merely A_v^ℓ .

Removing the Reliance on a Single Root. Another limitation of \mathbf{CT}_0 is the reliance on G_ϵ being good. This can be remedied by allowing more nodes in the

CellTree to perform the duties of the root: Specifically, many (or all) nodes can issue proofs of assimilation, and the **discover** module can be initialized with the crew information of many such nodes, so that it can choose any of them as its starting point.

Secondly, the assimilation operations carried out by the root’s crew (or other nodes close to it) are publicly verifiable (the only secret information used by them being their signing keys, which have corresponding verification keys already published). Further, a set of subtrees can easily “migrate” to a new CellTree. This makes the CellTree less reliant on the root’s crew, and thereby would also disincentivize the root’s crew from misbehaving.

Modifying the Crews. CT_0 provided no mechanism to modify the crew for any node. In a more realistic instantiation, one would include such a provision (say, by extending the scope of the module **selectCrew**). This feature would be useful for revoking the public keys of crew members, for replacing misbehaving crew members, or for resizing a crew.

Quality of Service. In CT_0 , the propagation schedule, which determines the assimilation rate, is fixed. However, one could support different rates of assimilation to different nodes, thereby providing varying levels of quality of service in terms of efficiency. Further, different nodes could be monitored by different number of nodes, providing different levels of assurance. These changes can be implemented by appropriately modifying the modules **rootward** and **leafward**.

Other Features. CT_0 does not exploit several features offered by the CellTree architecture, like the ability of modules to evolve, and the ability for multiple modules to coexist in a CellTree. We briefly mention several such features, which are discussed in [Appendix B](#).

- **Pruning, Grafting and Mirroring.** Subtrees can be detached (pruned) from a CellTree, or grafted on to a CellTree – possibly in multiple locations, in case of mirroring – with little effect on the other nodes (which are not monitoring any part of the subtree in question).
- **Excising Cells.** A cell’s contents can be deleted, or altered without respecting its program, if *the crews of all the nodes monitoring it* cooperate.
- **Computed Reads.** It is possible to support access to a *function* of a cell’s content (with proof).
- **Secret Cells.** The crew members of a node need not be aware of the contents of the cell, but can still provide authorized clients with access to the cell contents, or functions thereof, via secret-sharing or secure multiparty computation protocols.
- **Computing on Multiple Cells.** Concurrent algorithms can be designed to operate on multiple cells,

working independently on each cell. This enables maintaining multiple views of a database in different nodes (e.g., in a banking application, each customer’s ledger is a partial view of a bank’s central ledger).

- **Saplings.** A sapling is a CellTree whose root is assimilated into another parent CellTree. While the parent’s crews do not monitor a saplings nodes’, they do provide commitment guarantees.
- **Higher Arity Trees.** Higher arity nodes can be easily simulated by allowing the same crew to operate a subtree instead of a single node.
- **Incentivization** The CellTree architecture is agnostic about higher level mechanisms that could be used for incentivizing parties to play the role of crew members. Different parts of the CellTree may employ different incentivization mechanisms.

Finally, we emphasize that a major feature of the CellTree architecture is the ability for different implementations of the same modules to coexist in the tree, offering application specific features in different parts of the tree. Designing such modules and analyzing their effect on the robustness and performance guarantees of a CellTree are left for future research.

7 An Application: Banking

In this section we demonstrate the use of the CellTree architecture in a real application by designing a simple distributed online banking system that uses the CellTree as its backbone. The application is designed to allow transactions across clients and across banks, in a decentralized manner without requiring real-time involvement of the banks.

To support this, the underlying CellTree needs to host two major kinds of cells: the *customer* cells, and the *bank* cells:

A ‘customer’ cell represents an account of a client of a bank and is involved in real-time transactions. Its crew can be selected by the customer, possibly employing professional crew-service providers. A ‘bank’ cell (possibly a branch cell) is operated by a crew selected by the bank, and carries out transactions between banks (once every “epoch” – say, a day). In addition there should be nodes in the CellTree that function as monitors of these nodes; the only requirement on the location and the crew selection process of these nodes is so as to meet the assumptions needed for correctness, consistency and liveness (as in [Section 6.1](#), if CT_0 is used).

We briefly sketch how a transaction is carried out between two accounts. At a high-level, the individual cells would locally update their cell information, with the monitoring cells ensuring that this is done correctly. Since different cells evolve asynchronously, the update would be carried out in a few steps.

Firstly, outside of the CellTree infrastructure, the two account holders agree on the details of the transaction – the account addresses, amount transferred, timestamps and sequence numbers, and random nonces contributed by the two parties. A transaction ID is computed by hashing these transaction details.

To evolve, the two cells play the role of an *initiator* and a *responder*.

1. First, the initiator evolves its cell to incorporate the transaction into an *outstanding transactions* list. It could be an outgoing transfer of money or an incoming transfer.
2. After the above evolution is (sufficiently) assimilated, the responder evolves to incorporate this transaction into its *recently accepted transactions* list. It is required that a transaction with this transaction ID exists in the initiator’s outstanding transactions list (to verify this, the monitors will need to query the initiator’s cell using an appropriate **fetch** module).
3. Finally, the initiator’s cell evolves to close the outstanding transaction (moving it to a *recently closed transactions* list), and update its account balance. For this it is checked that the transaction ID being closed appears in the responder’s recently closed transactions list.

In steps 1 and 2, to be a valid evolution for the sender, the account balance of the sender should permit the current transaction, after subtracting all the outgoing amounts in the outstanding transactions, but not adding any of the outstanding incoming amounts (some flexibility consistent with a bank’s policies could be programmed into the cell).

During every evolution it is also checked that there are no outstanding transactions that have *expired* (i.e., too old to be accepted); all such transactions must be closed or *canceled*, depending on whether the responder has accepted the transaction or not. Also, for efficiency, the “recent” lists must be cleared of sufficiently old entries (e.g., 2 epochs); if an initiator does not close or cancel transaction within this duration (so that it becomes impossible to determine if the transaction was accepted or not), its outstanding transaction will be turned into a *transaction loss*, resulting in the forfeit of the amount involved (outgoing transfers treated as closed, and incoming ones treated as canceled).

Note that the transaction between two accounts does not involve any action from a bank (except if there are agents of the bank involved as crew members in the transacting cells or their monitors). However, if a transaction occurs between accounts in different banks, the banks do need to carry out a transaction between each other. Such transactions are carried out periodically

(once every epoch), using a similar process as above to update their balance sheets.

A customer can also carry out direct transactions with the bank, corresponding to deposits, withdrawals, interest payments and banking charges. The link between physical currency deposited at the bank branches or withdrawn from ATMs is established by the bank. Other instruments like checks, demand drafts and bank transfers using other digital means (other than on the CellTree) will also require the banks to be involved in the transaction.

We rely on the CellTree for the correctness, consistency and liveness conditions. To ensure this, the crews have to be sufficiently trustworthy. These crews could consist of agents of the banking system as well as commercial agents that the account holders can hire. One may also use an instance of the CellTree architecture which provides better liveness guarantees than **CT₀**.

Another important consideration in using a CellTree for bank accounts is that of privacy. This can be addressed using the following measures:

- The **fetch** modules will incorporate access control layers to ensure that only authorized users can read the contents of a cell.
- To facilitate transactions, crews of other cells are allowed to check if a transaction ID is present in the outstanding transactions list of the cell. The **fetch** module will require that the transaction ID is communicated by the reader and the cell responds only with a boolean answer confirming or denying the existence of the transaction ID. Since the transaction ID itself is a hash including random nonces, an untrusted cell’s crew cannot learn anything using such queries.
- In the above scheme, the crews of the cells monitoring a cell learn about all the transactions carried out by a cell. This can be avoided using a nucleus which only involves commitments to the current balance and uses zero-knowledge proofs which are used to check that the amounts involved in the transactions are correctly incorporated into the balance, when a transaction is accepted or closed. (The zero-knowledge proof system used would support showing that the sum of a pair of committed values equals another committed value.)

• Even if using zero-knowledge proofs, the crews of an account’s own cell will need access to the actual balance and current transactions (maintained as contents of the cell). These crew members, chosen by the customers, could be considered trusted by each customer. But for further security, the crew could use secure multi-party computation to protect against a fraction of the crew members being corrupt (we call such cells “secret cells”). Since secret cells involve more computational effort, one may consider a business model that provides crew members for a secret cell for a fee.

Finally, before such a scheme can be practically deployed, it should also include mechanisms for recovery from errors and failures. This can be accomplished within the framework of the CellTree architecture by suitably defining the various modules.

8 Related Work

Merkle trees [14, 15], and succinct proofs using them [12] have been valuable tools in cryptographers’ toolkit for a long time. Cryptographically authenticated blockchains and public ledgers can be traced back to the work of Haber and Scott [9], but became popular only with the advent of Bitcoin [16] and other crypto currencies that used them for recording their transactions. It is outside the scope of this work to survey all the ensuing innovations in this area.

But below we shall discuss a few distributed data repositories which deviate from the blockchain topology, and mention how the CellTree architecture is different from them in its goals and features.

The *Hashgraph* [2] is a distributed ledger with high transaction throughput as compared to blockchains. Its efficient functioning, however, requires that the set of parties involved in the protocol be aware of all the others. (This is comparable to how a single node’s crew operates in a CellTree.) The blocks (or events as they are called) in the Hashgraph form a directed acyclic graph (DAG) with hash pointers as edges (unlike in a blockchain, where the graph is a single path), with new blocks pointing to old blocks (like in a blockchain). As in standard blockchains, the selection and confirmation of new blocks in the ledger are probabilistic, but the ability of all the parties to interact with each other allows the use of an efficient voting protocol (rather than one based on, say, proof-of-work). Also, as in standard blockchains, the desired guarantee is that of immutability or persistence of blocks that are confirmed.

The *Tangle* [19] is a permissionless distributed data structure which also uses a DAG structure to store the transactions, again with the goal of increasing the throughput compared to a blockchain. Tangle allows users to be aware of only parts of the entire data structure. Incidentally, the specific algorithms used for building the DAG structure and considering a node confirmed are known to be susceptible to “*parasitic chain attacks*,” and is the subject of ongoing research [4].

The above two systems store data in a graph with hash-pointers as edges that has the form of a DAG (rather than a path), to increase the throughput of transactions. The CellTree architecture shares this feature, but promises even better performance when multi-level confirmation can be exploited. Also, the other differences that CellTree has with blockchains continue to apply to these systems as well.

The *Inter-Planetary File System* (IPFS) [3] is a peer-to-peer version controlled file system in which data items, with (optional links to other data items) can be stored. While different in its goals from blockchains, IPFS is also a distributed data repository, with parties storing some of the IPFS objects in local storage, and accessing others from a peer. The IPFS uses content-addressed sharing, where the address is a hash of the content (with linked objects replaced by their hashes). To detect and avoid duplication, IPFS uses deterministic hashing (no nonces used) so that the same file is hashed to the same address.

Unlike a data repository that is queried using addresses (e.g., the node address in a CellTree), IPFS does not attempt to provide any form of consensus on the “correct” data. All data items, linking to previously existing data items, are valid, and they have their own content-based address. As such, only liveness (all stored data can be retrieved) is of concern to IPFS.

9 Conclusion

The realization of a robust and immutable global Bulletin Board structure where all parties share the same view has been a long standing dream for the whole community supporting decentralization, and various versions of this have been made possible with the use of robust cryptographic building blocks and by relying on notions such as honest majority and the eventual synchronization of networks. The Blockchain is one such structure that has gained immense popularity as, not only does it have some desirable properties in theory, but it has been shown to work in practice as the backbone of various decentralized protocols, having overcome many game theoretic, social, and political challenges.

In this work we have introduced the CellTree architecture for distributed data repositories. A CellTree is designed to be flexible and heterogenous, allows data and policy evolution, allows parties to focus on only parts of the repository that are of interest to them, and separates out sub-tasks into modules that can be instantiated differently in different parts of the tree. We also illustrated an instantiation of the architecture, **CT₀**, and presented formal security guarantees.

A CellTree is rich in features and is highly programmable. We leave it for future work to exploit this novel architecture for applications.

References

- [1] Adam Back, Matt Corallo, Luke Dashjr, Mark Friedenbach, Gregory Maxwell, Andrew Miller, Andrew Poelstra, Jorge Timón, and Pieter Wuille. Enabling blockchain innovations with

- pegged sidechains. <https://blockstream.com/sidechains.pdf>, 2014.
- [2] Leemon Baird, Mance Harmon, and Paul Madsen. Hedera: A governing council and public hash-graph network. <https://www.hederahashgraph.com/whitepaper>, 2017.
- [3] Juan Benet. IPFS - content addressed, versioned, P2P file system. *CoRR*, abs/1407.3561, 2014.
- [4] Andrew Cullen, Pietro Ferraro, Christopher K. King, and Robert Shorten. Distributed ledger technology for iot: Parasite chain attacks. *CoRR*, abs/1904.00996, 2019.
- [5] Namecoin Foundation. Namecoin. <https://namecoin.org>, 2011.
- [6] VeChain Foundation. Vechain development plan. https://cdn.vechain.com/vechain_ico_ideas_of_development_en.pdf, 2017.
- [7] Juan A. Garay, Aggelos Kiayias, and Nikos Leonardos. The bitcoin backbone protocol: Analysis and applications. <https://eprint.iacr.org/2014/765.pdf>, 2015.
- [8] Yossi Gilad, Rotem Hemo, Silvio Micali, Georgios Vlachos, and Nikolai Zeldovich. Algorand: Scaling byzantine agreements for cryptocurrencies. In *Proceedings of the 26th Symposium on Operating Systems Principles*, pages 51–68. ACM, 2017.
- [9] Stuart Haber and W Scott Stornetta. How to timestamp a digital document. In *CRYPTO*, pages 437–455. Springer, 1990.
- [10] IOHK. Cardano. <https://www.cardano.org>, 2015.
- [11] Aggelos Kiayias, Alexander Russell, Bernardo David, and Roman Oliynykov. Ouroboros: A provably secure proof-of-stake blockchain protocol. In *CRYPTO*, pages 357–388. Springer International Publishing, 2017.
- [12] Joe Kilian. A note on efficient zero-knowledge proofs and arguments. In *Proceedings of the twenty-fourth annual ACM symposium on Theory of computing*, pages 723–732. ACM, 1992.
- [13] Roman Matzutt, Jens Hiller, Martin Henze, Jan Henrik Ziegeldorf, Dirk Müllmann, Oliver Hohlfeld, and Klaus Wehrle. A quantitative analysis of the impact of arbitrary blockchain content on bitcoin. 2018.
- [14] Ralph Merkle. Method of providing digital signatures. US Patent 4309569, 1982.
- [15] Ralph C. Merkle. A digital signature based on a conventional encryption function. In *Advances in Cryptology - CRYPTO '87, A Conference on the Theory and Applications of Cryptographic Techniques, Santa Barbara, California, USA, August 16-20, 1987, Proceedings*, pages 369–378, 1987.
- [16] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system. <http://www.bitcoin.org/bitcoin.pdf>, 2009.
- [17] BBC News. Child abuse images hidden in cryptocurrency blockchain. <https://www.bbc.com/news/technology-47130268>, 2019.
- [18] Joseph Poon and Thaddeus Dryja. The bitcoin lightning network: Scalable off-chain instant payments. <https://lightning.network/lightning-network-paper.pdf>, 2016.
- [19] Serguei Popov. The tangle. http://iotatoken.com/IOTA_Whitepaper.pdf, 2017.
- [20] Gavin Wood. Ethereum: A secure decentralized generalized transaction ledger. <http://gavwood.com/paper.pdf>, 2014.

A Modules: Further Details

In this section, we expand on the modules mentioned in [Section 5.2](#). Here, the modules are listed by name with brief notes on the various methods within them.

selectCrew: A new node is created by the set of nodes who would be monitoring it. They also select a crew for the new node at that time. The exact mechanism used to initiate node creating and carry out crew selection are left to the module. As an example, the crew selection algorithm in this module could be based on the committee selection protocol used within Algorand [8]. We emphasize that, like all modules, multiple implementations of **selectCrew** can operate in various parts of a CellTree simultaneously.

In our implementation for **CT₀** the protocol is initiated by the node whose child node is to be created. The crews of all ℓ ancestors are involved in the selection of the crew and the node is initialized with an empty cell. As part of this module’s functionality, the addresses of the members of the crew selected for a newly created node will be made available to the crews of the monitoring nodes. Also the new node’s crew is given the addresses of the crew members of nodes monitoring it. Further, the module is also expected to set (default) modules for the newly created node.

selectCell: When a node is created, it has an empty cell, which can then evolve over time. As mentioned before, cell evolution is triggered by a module which is responsible for obtaining a consensus among the crew members as to the next version of a cell. As an example, the **selectCell** module could use a payment based model to allow a client to choose the next version of a cell, possibly subject to custom policies. Note that irrespective of the mechanism used by **selectCell**, the CellTree algorithm EVOLVE checks that the selected cell is compatible with the existing cell.

rootward: This module has functions necessary for facilitating the ROOTWARDPROPAGATION procedure. It has the **receive** and **send** methods operating in the background that get assimilation signals from nodes in the monitored subtree and propagate the same to the ancestors monitoring it respectively. When they are invoked by a party as part of a propagation algorithm, they may use local buffers to interface with the background processes. The background processes may involve protocols with multiple rounds of interaction among the parties in a crew or across crews. In particular, a typical instance of the method **receive** will have the parties in a crew use a consensus algorithm to agree on the update signals being received, before storing them in a buffer. The **send** and **receive** methods would typically also involve secure (especially, authenticated) communication protocols for crew-to-crew communication.

Note that as part of rootward propagation, there are two services that a node's crew performs: (1) It checks consistency of nuclear updates at the nodes it is monitoring, and (2) It generates proofs of assimilation. One can configure the system so as to make these operations selective, so that not all crews perform all of these operations at every instance of propagation. For instance, the code execution module **exec** (described below) can be configured to require a payment for carrying out consistency checks, or the module **crewSign** can be configured to produce a signature only once every several updates.

Other functions in **rootward** include the **monitored** function that returns the node ids in the subtree monitored by that node, and **schedule** that determines how often the ROOTWARDPROPAGATION procedure runs. The scheduling can be static (run once every epoch), or dynamic (e.g., triggered by the availability of a new assimilation signal). Scheduling can also react to incentives or negotiations.

leafward: Similar to the **rootward** module, this has all the necessary functions for facilitating the LEAFWARDPROPAGATION procedure. Its **receive** and **send** methods propagate proofs of assimilation. It has its own **schedule** function for triggering the procedure. The **leafward** could be configured to let a crew directly

deliver proofs of assimilation from an assimilation root to any of the nodes it monitors (selected via the method **pickleaves**).

store: This module is used by each party in a crew to locally store and retrieve various values across separate invocations of the algorithms. While the data stored using this module pertains to an entire node, rather than an individual crew member, we envisage this module as operating locally at each crew member location, without any communication. As such, consistency guarantees are derived from the other modules (e.g., **selectCrew**, **selectCell**, **rootward.receive** and **leafward.receive** all would incorporate consensus mechanisms).

The **store** module may also support *garbage collection* to remove old versions of cells and nuclei of monitored nodes.

fetch: This has two major functions:

- **getCell** - Used in the READ procedure, this function allows an external client to query a node's crew for the cell it is hosting. Along with a cell, this also returns a proof of assimilation of the cell, at a specified assimilation root. Note that many versions of a cell may be available at a node, and the module can determine which one to return (possibly after a negotiation between the client and the crew); a natural choice is the latest version which has a proof of assimilation from the specified assimilation root.
- **getCrew** - This is used by crew members and external clients alike as part of the **discover** module to query other crews for either a list of its own crew members or that of a crew it is monitoring.

The **fetch** module may include a layer of customized access control, to restrict the set of clients who are allowed to query a node. It may also incorporate measures to prevent denial-of-service attacks by malicious clients (e.g., by requiring a proof-of-work by the client).

discover: For a client to read from a CellTree node, the client needs to first obtain the addresses of the crew members of the node it is reading from – a functionality encapsulated in the module **discover**. We point out that, within this module, a client will need to know *some* public-key material to bootstrap its access to a CellTree that it is interested in. (This could be in the form of the public-keys of (some of) the crew members for one or more nodes. The module **discover** may itself rely on the module **fetch** to access pieces of information from these nodes.)

crewSign: The proofs of assimilation involve a collective signature by the crew (on a hash value). Such signatures are also widely used in protocols for assimilation, **fetch**, and **selectCrew** where communication

between crews and external clients or other crews need to take place. The protocol for creating such signatures and the algorithm for locally verifying them are encapsulated in the module `crewSign`.

communication: This is a low level module that all other modules access to run protocols involving sending and getting messages to and from other crew members, external clients, or other crews. It internally accesses and maintains address books and uses `crewSign` to authenticate messages.

exec: The machine model used to execute the nuclear code is specified as a module. Making this a module allows different nodes of a CellTree to use different sets of languages. Also, a language can be updated over time, for instance, by introducing new library functions. The `exec` module would use the appropriate interpreter for a nucleus' code (note that we expect that clients reading a node and crews monitoring would use `exec` modules that can correctly interpret its nuclear code).

hash: The module provides methods for evaluating a hash function on a list of inputs (`generate`) and computing the hash-pointer for a Merkle tree (`mEval`). The module could allow multiple hash algorithms (with a header in the output indicating which one is used), and they could be updated over time. A method `generateNonce` included in this module can be used to create nonces with customized content, to control the hash evaluation (e.g., if a time-stamp in the nonce is in the future, the hash evaluation could return an error).

moduleManager: The CellTree framework admits different implementations of the various modules to coexist. But the framework does not specify how modules are chosen, and possibly changed over time. Some modules like `selectCell` are entirely local to a node, and can be freely changed if the node's crew has a consensus on it. But certain other modules like `selectCrew` involve several crews and will need a consensus from all of them. In either case, we propose that module substitutions follow a mechanism similar to cell evolution, in that the current and next modules must agree with each other for a substitution to be considered acceptable. This would allow long term guarantees regarding the functioning of a node and parts of the tree (e.g., a `selectCrew` module may not allow its core mechanism to be changed, but only its parameters).

For the sake of completeness, we define a module `moduleManager` that is used to update all the code and data at each node related to the implementation of various modules, including itself. Among other things, this module could be used to implement bug fixes and patches to protocols.

B Additional Features

A CellTree is highly programmable, thanks to the ability to design nuclei with arbitrary code. Apart from this, the CellTree architecture also provides several possibilities that may not be obvious. Here we mention a few of them. We expect that more possibilities will be identified and exploited as CellTrees are used for more applications.

Pruning, Grafting and Mirroring

To prune or detach a node from the tree, the crew of a node can simply stop propagating the assimilation signals from its subtree to the root, and vice versa. Subsequent updates on this subtree are not assimilated into the original tree. Note that outside of this subtree, the original tree can continue to function unaffected by this change, without requiring the crews to take any special action. Also, the pruned out subtree can continue to function as a live CellTree, with the detaching node taking on the role of the root.

Grafting refers to the reverse operation, in which the root of a CellTree is accepted as a newly created node in another CellTree. We note that the proofs of assimilation of a node v at a node u does not depend on the absolute addresses of u and v , but only in the relative address of v with respect to u (implicit in the order in which hashes are accumulated). This allows the root of a grafted subtree to immediately propagate its assimilation signals towards the new root.

Mirroring refers to when a subtree of a CellTree is grafted on to a different location in the CellTree, without detaching itself from its original location. Mirroring creates multiple paths of assimilation and discovery for a subtree, buying it protection against potential failure of one path. A mirrored subtree receives multiple leafward propagations, one from each location it is mirrored at. (The nodes monitoring a node in the newly mirrored subtree would annotate their proofs of assimilation to indicate when they started monitoring it. Such annotations can be included in the nonces.)

Excising Cells

A cell may be programmed to disallow removing (or even altering) its data. Even so, if the crew of a node wishes to remove (or alter) the contents of its cell from the CellTree, it may do so *in cooperation with the crews of all the nodes monitoring it*. Such a feature, combined with the ability to change the crew of a node, is helpful in removing any illegal content discovered in a node.

Computed Reads

While the method `fetch.getCell` would typically be used to retrieve an entire cell, it could offer a partial view of the cell, provided that the nucleus' `chkCell` algorithm can verify this partial view. For instance, if a

cell contains a blockchain ledger, the partial view may include only the last block, and an accumulated hash of the previous blocks. More generally, `chkCell` could accept a zero-knowledge proof that a given function of the cell is consistent with a commitment (hash) of the cell that is included in the nucleus.

Secret Cells

Using cryptographic techniques, it is possible for a crew to maintain a cell without any crew member (or even a collusion of malicious parties in the crew) learning the contents of the cell. Authorized clients can access this cell data — or a function thereof (as a computed read) — via secure computation protocols. A cell could also be used to pool secret data from a set of clients, and allow authorized clients to query the resulting data.

Computing on Multiple Cells

Evolution of a cell can depend not only on its own data, but also on the cells in other nodes, by using appropriately designed nuclei, using an `exec` module that allows reference to cells in other nodes. Dereferencing would involve reading a cell (subject to access controls). Such computations would be relevant in maintaining multiple views of a database in various nodes (e.g., in a banking application, each customer’s ledger is a partial view of a bank’s central ledger).

We remark that since the nodes are updated individually and asynchronously, the nuclei need to be carefully designed, so that the consistency guarantee of individual nuclei translates to a guarantee that the global updates were carried out correctly.⁵

Saplings

A sapling is a CellTree “supported by” another CellTree for time-stamping purposes. The sapling will be a separate CellTree whose root’s nucleus is included in the nucleus of a node in the parent tree (along with additional information like time-stamps). Saplings can be used to create “enterprise subtrees,” whose contents are entirely hidden from the rest of the tree, but with *commitment guarantees* provided by the larger tree.

⁵As an illustration, suppose we wish to support swapping the cells in nodes u and v (say C_1 and C_2 , respectively). First the cell in u is updated from C_1 to (C_1, C_2) (with its nucleus carrying hashes of both the cells); during the time this updated cell is assimilated by the nodes monitoring u , the cell C_2 in v should not change, so that `chkNext` accepts the update. Later the node v is updated to replace C_2 with C_1 ; to ensure that this update is part of a swap, `chkNext` will confirm that the cell in u has contents (C_1, C_2) , as expected after the above step. Finally, u will be updated again, so that it replaces (C_1, C_2) with C_2 ; `chkNext` will confirm that v has C_1 at this point.

Higher Arity Nodes

Since the depth of a node in the tree affects its performance guarantees, restricting to a binary tree may appear sub-optimal. However, higher arity nodes can be easily simulated by allowing the same crew to operate a subtree instead of a single node. For instance, a crew which operates a node v and its two children $v0$ and $v1$, can simulate a node v with four children ($v00, v01, v10, v11$) but at a smaller depth.

Incentivization

The CellTree architecture is agnostic about higher level mechanisms that could be used for incentivizing parties to play the role of crew members. In particular, clients wishing to load their data into a node’s cell could be required to make a payment (using digital currencies or legal tenders) to the parties in that node’s crew; in turn, the crew members could be required to make payments to the nodes monitoring them, for their continued services.

C Proving Robustness of CT₀

Here we sketch a proof for [Theorem 1](#).

Proof sketch: Firstly, we argue that under the assumption $A_v^\ell(G_\epsilon)$, `discover`(v) may only return a crew that is *considered* to be G_v by a good node u monitoring v (or it will abort). Consider the set of ℓ ancestors of v all of which monitor v . By assumption $A_v^\ell(G_\epsilon)$, there is a node u among them which is good w.r.t. G_ϵ . Then, there exist a sequence of nodes $\epsilon = u_0 \prec u_1 \prec \dots \prec u_n = u$, and sets of parties $\{\hat{G}_{u_i}\}_{i=0}^n$, such that $\hat{G}_{u_0} = G_\epsilon$, and the successive crews consider the next one in the sequence to be the correct crew. Hence if `discover` queries \hat{G}_{u_i} for the crew of u_{i+1} , it will receive $\hat{G}_{u_{i+1}}$ in response. Then, inductively, `discover` will accept only \hat{G}_{u_i} as the crew for u_i (base case being $u_0 = \epsilon$ (or will abort)). Then, it will return a crew only if \hat{G}_{u_n} considers it to be G_v .

Next, we argue that if a proof of assimilation obtained by `READ`(v) verifies, then the cell retrieved along with it must have a nucleus that was verified by a good node u monitoring v . Again, consider the sequence of nodes $\epsilon = u_0 \prec u_1 \prec \dots \prec u_n = u$ from above. The root’s crew signed the proof of assimilation only because all the assimilation signals that were assimilated were signed by all the nodes monitoring the respective nodes. In particular, u_1 which is monitored by $\epsilon = u_0$ must have signed the assimilation signals that the root incorporated in the proof. In turn, u_1 required a signature by u_2 on the assimilation signature that u_1 incorporated into the signal sent to u_0 . Continuing in this manner, we are guaranteed that the proof of assimilation signed by the root corresponds to an assimilation signal sent by $u_n = u$. Thus any cell returned by `READ`(v) will have a nucleus that was assimilated by the node u . u , being a good

node that monitors v , would incorporate a nucleus of v into this signature only after a consistency check.

Now we are ready to prove consistency. For this, we consider the sequence of nuclei that u has maintained for v , $\text{nuc}^{(0)}, \dots, \text{nuc}^{(N)}$ so far. By the above argument, the nuclei of any cell cell accepted by `READ` appears in this sequence (i.e., $\text{cell.nuc} = \text{nuc}^{(i)}$ for some i). Further, u 's consistency checks ensure the conditions $\text{chkNext}(\text{nuc}^{(i-1)}, \text{nuc}^{(i)})$ and $\text{chkPrev}(\text{nuc}^{(i)}, \text{nuc}^{(i-1)})$. Also, the verification in `READ` ensures that $\text{chkCell}(\text{cell})$ holds. These are precisely the conditions required by the consistency guarantee.

The correctness guarantee follows from consistency and the extra condition of $A_v(G_\epsilon)$, as a good node for v will propagate only the updates created by `EVOLVE` for assimilation.

Finally, if $A_v^1(G_\epsilon)$ holds we argue that liveness holds. A non-trivial issue here is to ensure that a node's siblings (or cousins) or descendants should not be able to prevent its rootward propagation. This is ensured in the procedure `ROOTWARDPROPAGATION`, by retaining a valid update from one child, even if the update from its sibling is discarded for failing consistency checks. Also note that if there is a path of all good nodes from ϵ to v , the proof of assimilation is routed back to v 's crew. Finally, we also rely on the fact that `discover`(v) will never abort and the procedure `READ` will also succeed, assuming $A_v^1(G_\epsilon)$. \square