

---

# A protocol & event-sourced database for decentralized user-siloed data

PICK<sup>1</sup>, FARMER<sup>1</sup>, SUTULA<sup>1</sup>, ?GOZALISHIVILI OTHERS?<sup>2</sup> AND HILL<sup>1</sup>

<sup>1</sup>*www.textile.io*

<sup>2</sup>*other?*

*contact@textile.io*

*revised October 3, 2019*

---

## 1. INTRODUCTION

Compared to their predecessors, modern cloud-based apps and services provide an extremely high level of convenience. Users can completely forget about data management, and enjoy seamless access to their apps across multiple devices. This convenience is now expected, but has come at the cost of additional consequences for users. One such consequence is that many of the same development patterns that bring convenience (e.g., single sign-on, minimal encryption, centralized servers and databases) also enable, or even require data hoarding by apps. While collecting large amounts of users' data can create value for the companies building apps (e.g., app telemetry, predictive tools, or even new revenue streams), that value flows mostly in one direction: apps collect, process and benefit from a user's private data, but users rarely have access to the original data or the new insights that come from it. Additionally, the data is generally not accessible to other apps and services. This is *app-siloed* data.

While the fact that companies collect user data may not itself be a significant problem, problems may arise if over time a company's incentives shift from *providing* value to users, to *extracting* value from users [1]. When this incentive shift happens, companies that have been creating data-silos may treat that data as new source of value or revenue for the company. It may be possible to stop this trend through extreme privacy measures, government intervention/legislation, or by boycotting companies that collect any form of data. Ideally, there is an alternative approach that allows individuals to capture the value of their data and still allow developers to build new interfaces and experience on top of this data. This approach is called *user-siloed* data and it fundamentally separates apps from their users' data.

One of the most exciting aspects of user-siloed data is the ability to build data-driven apps and services while

users remain in control of their own data. Other projects have identified app-siloed data as a problem [2], [3], and some have identified user-siloed data as a solution [4]. However, none so far have addressed the problem of *how data should be collected to make it extensible*, nor have they provided a sufficiently interoperable protocol for scalable storage, transmission, and use of user-siloed application data.

In this paper, we study existing technologies that could be used to build a network for user-siloed data. We outline six challenge areas: flexible data formats, efficient synchronization, conflict resolution, access-control, scalable storage, and network communication. Based on our findings, we propose a novel architecture for event sourcing (ES) with Interplanetary Linked Data (IPLD), that is designed to store, share, and host user-siloed datasets at scale. Our proposed design leverages new and existing protocols to solve major challenges with building a secure and distributed network for user data while at the same time providing the flexibility and scalability required by today's apps.

## 2. BACKGROUND

We now describe some of the technologies and concepts motivating the design of our novel decentralized database system. We highlight some of the advantages and lessons learned from event sourcing (ES) and discuss drawbacks to using these approaches in decentralized systems. We provide an overview of some important technologies related to the Interplanetary File System (IPFS) that make it possible to rethink ES in a decentralized network. Finally, we cover challenges to security and access control on an open and decentralized network, and discuss how they are designed in popular database management (DBMS) systems outside the IPFS context.

## 2.1. Data Synchronization

To model realistic systems, apps often need to map data between domain models and database tables, where the same data model is used to both query and update a database. To solve *synchronization* it is often helpful to handle updates on just the database and then provide the query and interfaces only later in a DBMS. One powerful approach to synchronization is to use a set of append-only logs to model the state of an object simply by applying its change sequence in the correct order. This concept can be expressed succinctly by the state machine approach [5]: if two identical, deterministic processes begin in the same state and get the same inputs in the same order, they will produce the same output and end in the same state. This is a powerful concept baked into a simple structure, and is at the heart of many distributed database systems [6].

**Logs or Append-only Log** A log is a registry of database transactions that is read sequentially (normally ordered by time) from beginning to end. In distributed systems, logs are often treated as append-only, where changes or updates can only be added to the set and never removed.

### 2.1.1. CQRS, Event Sourcing, and Logs

For most apps, it is critical to have reliable mechanisms for publishing updates and events (i.e., to support event-driven architectures), scalability (optimized write and read operations), forward-compatible application updates (e.g., code changes, retroactive events), auditing systems, etc. To support such requirements, developers have begun to utilize event sourcing and command query responsibility segregation (CQRS) patterns [7], relying on append-only logs to support immutable state histories. Indeed, a number of commercial and open source software projects have emerged in recent years that facilitate ES and CQRS-based designs, including Event Store [8], Apache Kafka [9] and Samza [10], among others [11].

**CQRS** Command query responsibility segregation or CQRS is a design pattern whereby reads and writes are separated into different models, using commands to write data, and queries to read data [12].

**ES** Event sourcing or ES is a design pattern for persisting the state of an application as an append-only log of state-changing events.

A key principal of ES and append-only logs is that all changes to application state are stored as a sequence of events. Because any given state is simply the result of a series of atomic updates, the log can be used to reconstruct past states or process retroactive updates [13]. The same principal means a log can be viewed as a mechanism to support an infinite number of valid state interpretations (see sec. 2.1.2). In other words, with minimal conformity, a single log can model multiple

application states [14].

### 2.1.2. Views & Projections

In CQRS and ES, the separation of write operations from read operations is a powerful concept. It allows developers to define views into the underlying data that are best suited for the use case or user interface they are building. Multiple (potentially very different) views can be built from the same underlying event log.

**View** A (typically highly de-normalized) read-only model of event data. Views are tailored to the requirements of the application, which helps to maximize display and query performance. Views that are backed by a database or filesystem-optimized access are referred to as materialized views.

**Projection** An event handler and corresponding reducer/fold function used to build and maintain a view from a set of (filtered) events. While projections may lead to the generation of new events, their reducer should be a pure function.

Views themselves are enabled by projections<sup>3</sup>, which can be thought of as transformations that are applied to each event in a stream. They update the data backing the views, be this in memory or persisted to a database. In a distributed setting, it may be necessary for projections to define and operate as eventually consistent data structures, to ensure all peers operating on the same stream of events have a consistent representation of the data.

### 2.1.3. Eventual Consistency

The CAP theorem [15], [16] states that a distributed database can guarantee only two of the following three promises at the same time: consistency (i.e., that every read receives the most recent write or an error), availability (i.e., that every request receives a [possibly out-of-date] non-error response), and partition tolerance (i.e., that the system continues to operate despite an arbitrary number of messages being dropped [or delayed] by the network). As such, many distributed systems are now designed to provide availability and partition tolerance by trading consistency for *eventual* consistency. Eventual consistency allows state replicas to diverge temporarily, but eventually arrive back to the same state. While an active area of research, designing systems with provable eventual consistency guarantees remains challenging [17], [18].

**CRDT** A conflict-free replicated data type (CRDT) assures eventual consistency through optimistic replication (i.e. all new updates are allowed) and eventual merging. CRDTs rely on data structures

<sup>3</sup>Terminology in this section may differ from some other examples of ES and CQRS patterns, but reflects the underlying architecture and designs that the Textile team will elaborate on in Section 3

that are mathematically guaranteed to resolve concurrent updates the same way regardless of the order in which those events were received.

How a system provides eventual consistency is often decided based on the intended use of the system. Two well-documented categories of solutions include logs (def. 1) and CRDTs (def. 2).

A common minimum requirement for log synchronization across multiple peers is that the essential order of events is respected and/or can be determined [19], [20]. For these cases, logical clocks are a useful tool for eventual consistency and total ordering [21]. However, some scenarios (e.g., temporarily missing events or ambiguous order) can force a replica into a state that cannot be later resolved without costly recalculation. In specific cases, CRDTs can provide an alternative to log-based consensus (see sec. 2.1.5).

#### 2.1.4. Logical Clocks

In a distributed system with multiple peers (each with an independent clock) creating events, local timestamps can't be used to determine "global" event causality. Machine clocks are never perfectly synchronized [22], meaning that one peer's concept of "now" is not necessarily the same as another. Machine speed, network speed, and other factors compound the issue. For this reason, simple wall-clock time does not provide a sufficient notion of order in a distributed system. Alternatives to wall-clock time exist to help achieve eventual consistency. Examples include various logical clocks (Lamport [22] Schwartz [19], Bloom [23], and Hybrid variants [21], etc.), which use "counter"-based time-stamps to provide partial ordering.

Cryptographically linked events can also represent a clock (see sec. 2.2.3). One such example is called the Merkle-Clock [24], which relies on properties of a Merkle-DAG to provide strict partial ordering between events. This approach does have its limitations however [24, Sec. 4.3]:

Merkle-Clocks represent a strict partial order of events. Not all events in the system can be compared and ordered. For example, when having multiple heads, the Merkle-Clock cannot say which of the events happened before.

#### 2.1.5. Conflict-Free Replicated Data Types

CRDTs (def. 2) are one way to achieve strong eventual consistency, where once all replicas have received the same events, they will arrive at the same final state, *regardless of ordering*<sup>4</sup>. A review of the types of possible CRDTs is beyond the scope of this paper, however, it is important to note their role in eventually

consistent systems and how they relate to clock-based event ordering. See for example [24], [25] for informative reviews of these types of data structures.

Whether a system (e.g., an app) uses a CRDT or a clock-based sequence of events is entirely dependent on the use-case and final data model. While CRDTs may seem superior (and are currently a popular choice among decentralized systems), it is not possible to model every system as a CRDT. Additionally, the simplicity of clock-based sequencing often makes it easier to leverage in distributed systems where data conflicts will only rarely arise. Lastly, logs and CRDTs are not mutually exclusive and can be used together or as different stages of a larger system.

## 2.2. Content-based Addressing

Internet application architecture today is often designed as a system of clients (users) communicating to endpoints (hosts or servers). Communication between clients and endpoints usually happens via the TCP/IP protocol stack and depends on *location-based addressing*. Location-based addressing, where the client makes a request that is routed to a specific endpoint based on prior knowledge (e.g., the domain name or IP address), works relatively well for many use-cases. However, there are many reasons why addressing content by location is problematic, such as duplication of storage, inefficient use of bandwidth, invalid/dead links (link rot), centralized control, and authentication issues.

An alternative to location addressing, called *content-based addressing*, may provide a solution to many of the problems associated with location-based addressing. Content-based addressing is where the content itself is used to create an address which is then used to retrieve said content from the network [26].

### 2.2.1. IPFS & Content-based Addressing

There are a number of systems that utilize content-based addressing to access content and information (e.g. Git, IPFS, Perkeep, Tahoe-LAFS, etc), and there is an active body of literature covering its design and implementation [27]–[29]. The Interplanetary File System (IPFS) — which is a set of protocols to create a content-addressed, peer-to-peer (P2P) filesystem — is one such system [27]. In IPFS, the address for any piece of content is determined based on a cryptographic hash of the content itself. In practice, an IPFS Content Identifier (CID) is a multihash, which is a self-describing "protocol for differentiating outputs from various well-established cryptographic hash functions, addressing size [and] encoding considerations" [30]. That addressing system confers several benefits to the network, including tamper resistance (i.e., a given piece of content has not been modified en route if its hash matches what we were expecting/requested) and de-duplication (i.e., the same content from different peers will produce the same hash

<sup>4</sup>Though non-commutative CRDTs may require a specific ordering of events in certain cases [24]

address). Additionally, IPFS content-based addresses are immutable and universally unique.

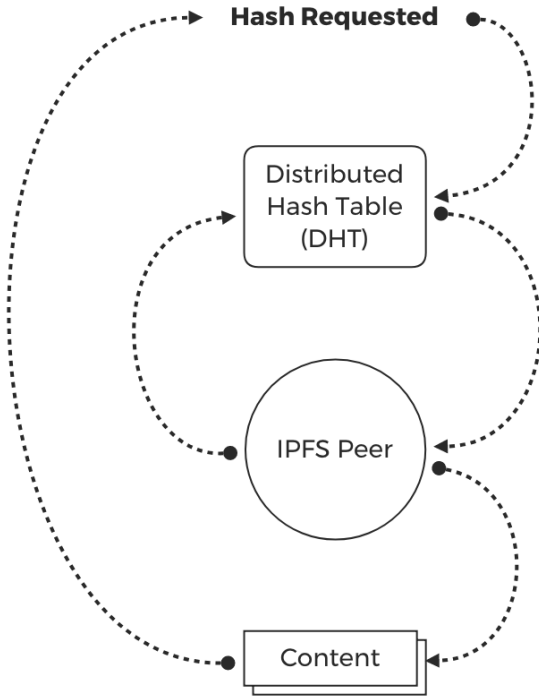


FIGURE 1: The cryptographic hash of content is used to make a request to the network of IPFS peers. Using the built-in routing mechanisms and a distributed hash table (DHT), peers hosting the requested content are identified and content is returned.

While content-based addressing doesn't dictate to a peer *how* to get a piece of content, IPFS (via libp2p<sup>5</sup>) does provide a system for moving content across the network. On the IPFS network, a client who wants specific content requests the CID from the network of IPFS hosts. The client's request is routed to the first host capable of fulfilling the request (i.e., the first host that is actively storing the content behind the given CID). The IPFS network can be seen as a distributed file system, with many of the benefits that come with this type of system design.

### 2.2.2. IPLD

As discussed previously, IPFS uses the cryptographic hash of a given piece of content to define its content-based address (see [27] for details on this process). However, in order to provide standards for accessing content-addressable data (on the web or elsewhere), it is necessary to define a common format or specification. In IPFS and other systems [31], this common data format is called Interplanetary Linked Data (IPLD)<sup>6</sup>. As the name suggests, IPLD is based on principals of linked data [32], [33] with the added capabilities of a content-based addressing storage network.

<sup>5</sup><https://libp2p.io/>

<sup>6</sup><https://ipld.io/>

IPLD is used to represent linked data that is spread across different "hosts", such that everything (e.g., entities, predicates, data sources) [34] uses content-based addresses as unique identifiers. To form its structure, IPLD implements a Merkle DAG, or directed acyclic graph<sup>7</sup>. This allows all hash-linked data structures to be treated using a unified data model, analogous to linked data in the Semantic Web sense [35]. In practice, IPLD is represented as objects, each with **Data** and **Links** fields, where **Data** can be a small blob of unstructured, arbitrary binary data, and **Links** is an array of links to other IPLD objects.

### 2.2.3. Merkle-Clocks

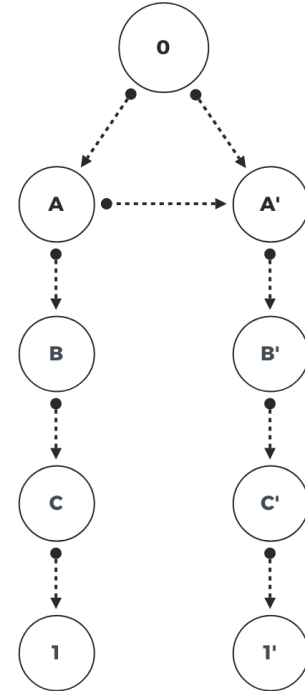


FIGURE 2: Divergent heads in a multi-writer Merkle-DAG

A Merkle-Clock is a Merkle-DAG that represents a sequence of events, or a log [24]. When implemented on IPFS (or an equivalent network where content can be cryptographically addressed and fetched), Merkle-Clocks provide a number of benefits for data synchronization between replicas [24, Sec. 4.3]:

1. Sharing the Merkle-Clock can be done using only the *root* CID. The whole Clock is unambiguously identified by the CID of its root, and its full structure can be traversed as needed.
2. The immutable nature of a Merkle-DAG allows every replica to perform quick comparisons, and fetch only those nodes (leaves) that it is missing.
3. Merkle-DAG nodes are self-verified and immune to corruption and tampering. They can be fetched

<sup>7</sup>Other examples of DAGs include the Bitcoin blockchain or a Git version history.

from any source willing to provide them, trusted or not.

4. Identical nodes are de-duplicated by design: there can only be one unique representation for every event.

On the downside (see also sec. 2.1.4), a Merkle-Clock cannot order divergent heads (or roots). For example, in fig. 2, two replicas (left and right columns) are attempting to write (top to bottom) events to the same Merkle-Clock. After the first replica writes event A, the second writes event A' and properly links to A. At that point, perhaps the two replicas stop receiving events from one another. To a third replica (not pictured) that does continue to receive events, there would now be two independent heads, 1 and 1'. For the third replica, resolving these two logs of events may be costly (many updates happened since the last common node) or impossible (parts of the chain may not be available on the network).

In order to reduce the likelihood of divergent heads, all replicas should be perfectly connected and be able to fetch all events and linkages in the Merkle-Clock. On real-world networks with many often replicas that are often offline (mobile and Internet of things (IoT) devices, laptops, etc.), these conditions are rarely met, making the use of a single Merkle-Clock to synchronize replicas problematic.

## 2.3. Networking

So far we have primarily discussed the mechanics of creating or linking content in a series of updates. Now we will overview some common networking tools for connecting distributed peers who aim to maintain replicas of a shared state. This could be any decentralized network of interacting entities (e.g., cloud servers, IoT devices, botnets, sensor networks, mobile apps, etc) collectively updating a shared state. IPFS contains a collection of protocols and systems to help address the networking needs of different use-cases and devices — be it a phone, desktop computer, browser, or Internet-enabled appliance.

### 2.3.1. Libp2p

The libp2p project provides a robust protocol communication stack. IPFS and a growing list of other projects (e.g., Polkadot, Ethereum 2.0, Substrate, FileCoin, OpenBazaar, Keep, etc) are building on top of libp2p. Libp2p solves a number of challenges that are distinct to P2P networks. A comprehensive coverage of networking issues in P2P systems is out of the scope for this paper, however, some core challenges that libp2p helps to address include network address translator [36] (NAT) traversal, peer discovery and handshake protocols, and even encryption and transport security — libp2p supports both un-encrypted (e.g. TCP, UDP) and encrypted (e.g. TLS, Noise) protocols — among others.

Libp2p uses the concept of a multiaddress to address peers on a network, which essentially models network addresses as arbitrary encapsulations of protocols [37]. In addition to “transport layer” modules, libp2p provides several tools for sharing and/or disseminating data over a P2P network.

### 2.3.2. Pubsub

One of the most commonly used P2P distribution layers built on libp2p, is its Pubsub (or publish-subscribe) system. Pubsub is a standard messaging pattern where the publishers don't know who, if anyone, will subscribe to a given topic. *Publishers* send messages on a given topic or category, and *Subscribers* receive only messages on a give topic to which they are subscribed. Libp2p's Pubsub module can be configured to utilize a *floodsub* protocol — which floods the network with messages, and peers are required to ignore messages in which they are not interested — or *gossipsub* — which is a proximity-aware epidemic Pubsub system, where peers communicate with proximal peers, and messages can be routed more efficiently. In those implementations, there is a benefit to using Pubsub in that no direct connection between publishers and subscribers is required.

Another benefit to using Pubsub is the ability to publish topical sequences of updates to multiple recipients. Like libp2p, encryption is a separate concern and often added in steps prior to data transmission. However, also like libp2p, Pubsub doesn't offer any simple solutions for transferring encryption keys (beyond public keys), synchronizing datasets across peers (i.e. they aren't databases), or enforcing any measures for access control (e.g. anyone subscribed to a topic can also author updates on that topic). To solve some of these challenges, some systems introduce message echoing and other partial solutions. However, it makes more sense to use Pubsub and libp2p as *building blocks* in systems that can effectively solve these issues, such as using multi-modal communication strategies or leveraging tools such as deferred routing (e.g. inboxing) for greater tolerance of missed messages.

### 2.3.3. IPNS

Our discussion of Pubsub and libp2p has so far only dealt with *push-based* transfer of data; but IPFS also offers a useful technology for hosting *pull-/request-based* data endpoints based on an Interplanetary Name System (IPNS). IPNS aims to address the challenge of mutable data within IPFS. It relies on a global namespace (shared by all participating IPFS peers) based on Public Key Infrastructure (PKI). By using IPNS, a content creator generates a new address in the global namespace and points that address to an endpoint (e.g. a CID). Using their private key, a content creator can update the static route to which the IPNS address refers. But IPNS isn't only useful for creating static addresses pointing to content-based addresses; it is also compatible with

TABLE 1: Example Access Control List.

	Create	Delete	Edit	Read
Jane	-	-	-	✓
John	✓	✓	✓	✓
Mary	-	-	✓	✓

external naming systems such as DNS, Onion, or bit addresses.

Unfortunately, many use-cases that require highly mutable data, also require rapid availability of updates, or need flexible multi-party access control, which is not currently viable using IPNS alone. However, taken together, libp2p, pubsub, IPNS, and IPFS more generally provide a useful toolkit for building robust abstractions to deliver fast, scalable, data synchronization on a decentralized network.

## 2.4. Data Access & Control

IPFS is an implementation of PKI, where every node on the network has a key-pair. In addition to using the key-pair for secure communication between nodes, IPFS also uses the key-pair as the basis for identity. Specifically, when a new IPFS node is created, a new key-pair is generated, and this public key is used to generate the node’s Peer IDentity (Peer ID).

### 2.4.1. Agent-centric Security

Agent-centric security refers to the maintenance of data integrity without leveraging a central or blockchain-based consensus. The general approach is to let the reader enforce permissions and perform validations, not the writer or some central authority. Agent-centric security is possible if the reader can reference local-only, tamper-free code or if the local system state can be used to determine whether a given operation (e.g., delete operation) is permitted. Many decentralized networks, such as Secure Scuttlebutt [38] and Holochain [39], make use of agent-centric security. Each of these systems leverage cryptographic signatures to validate peer identities and messages.

### 2.4.2. Access control

All file-systems and databases have some notion of “access control”. Many make use of an access-control list (ACL), which is a list of permissions attached to an object or group of objects [40]. An ACL determines which users or processes can access an object and whether a particular user or process with access can modify or delete an object (see tbl. 1).

Using ACLs in systems where identity is derived from various configurations of PKI has been around for some time [41]. Still, many existing database and communication protocols built on IPFS to date lack

support for an ACL or only have primitive ACL support. Where ACLs are missing, many systems use cryptographic primitives like signature schemes or enable encryption without any role-based configuration. Even more, many systems deploy an all-or-none security model, where those with access to a database have complete access, including write capabilities. Ideally, ACLs are mutable over time, and permission to modify an ACL should also be recorded in an ACL.

Event-driven systems (e.g., event sourcing) often make use of ACLs with some distinct properties. The ACL of an ES-based system is usually a list of access rules built from a series of events. For example, the two events, “grant Bob write access” and “revoke read access from Alice” would together result in a final ACL state where, Bob has read and write access, but Alice does not.

## 3. THE THREADS PROTOCOL

We propose Threads, a protocol and decentralized database that runs on IPFS meant to help decouple apps from user-data. Inspired by event sourcing and object-based database abstractions, Threads is a protocol for creating and synchronizing state across collaborating peers on a network. Threads offer a multi-layered encryption and data access architecture that enables datasets with independent roles for writing, reading, and following changes. By extending on the multiaddress addressing scheme, Threads differs from previous solutions by allowing *pull*-based replica synchronization in addition to *push*-based synchronization that is common in distributed protocols. The flexible event-based structure enables client applications to derive advanced applications states, including queriable materialized views, and custom CRDTs.

In essence, Threads are topic-based collections of single-writer logs. Taken together, these logs represent the current “state” of an object or dataset. The basic units of Threads — Logs and Events — provide a framework for developers to create, store, and transmit data in a P2P distributed network. By structuring the underlying architecture in specific ways, this framework can be deployed to solve many of the problems discussed above.

### 3.1. Event Logs

In multi-writer systems, conflicts arise as disparate peers end up producing disconnected state changes, or changes that end up out of sync. In order to proceed, there must be some way to deal with these conflicts. In some cases (e.g. `ipfs-log` [42]), a Merkle-Clock can be used to induce ordering. This approach cannot achieve a total order of events without implementing a data-layer conflict resolution strategy [24]:

A total order can be useful ...and could be obtained, for example, by considering concurrent events to be equal. Similarly, a

strict total order could be built by sorting concurrent events by the CID or their nodes or by any other arbitrary user-defined strategy based on additional information attached to the clock nodes (data-layer conflict resolution).

Solutions such as **ipfs-log** include a built-in CRDT to manage conflicts not resolved by the Merkle-Clock. This approach works for many cases, but the use of a deterministic resolution strategy can be insufficient in cases with complicated data structures or complicated network topologies. A git merge highlights one such example, in which a predetermined merge strategy could be used, but is not often the best choice in practice. Furthermore, a multi-writer log using a linked data format (e.g., a Merkle-DAG) in imperfect networking or storage environments can lead to states where it is prohibitively difficult (e.g., due to networking, storage, and/or computational costs) to regain consistency.

A promising approach to dealing with this is to leverage the benefits of both a Merkle-Clock for events from any one peer, and a less constrained ordering mechanism to combine events from all peers. In this case, developers can more freely institute their own CRDTs or domain-specific conflict resolution strategies. Additionally, it naturally supports use-cases where all peers contributing to a dataset may not be interested in following or replicating the events of all other peers (e.g., in a Pubsub-based system).

### 3.1.1. Single-writer Event Logs

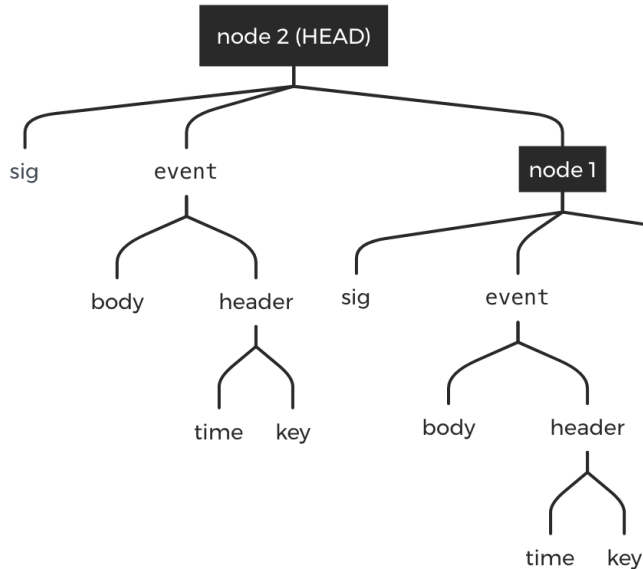


FIGURE 3: A single-writer Merkle-Clock (Event Log).

Our solution to dealing with log conflicts (i.e., divergent Merkle-Clocks) is to institute a *single-writer rule*: A

Log can only be updated by a single replica or *identity*. An Event Log is then a *single-writer Merkle-Clock* that can be totally ordered (fig. 3), and *separate* Event Logs can be composed into advanced structures, including CRDTs [25].

This presents a novel take on multi-writer systems [25], such that conflict resolution is *deferred* to a point at which a decision is actually required. This means that imperfect information may be supplemented along the way without causing conflicts in the mean time. It also means that apps can choose conflict resolution strategies specific to the task at hand. For example, if using a downstream CRDT, ordering is actually irrelevant and can be ignored completely. Alternatively, an additional clock may be required to ensure consistent ordering, such as a vector or Bloom clock (see sec. 2.1.4). Finally, even manual merge-type strategies are possible if this is the desired conflict resolution strategy.

**Writer** The single IPFS Peer capable of writing to an Event Log.

**Reader** Any Peer capable of reading a Log. Practically speaking, this means any Peer with the Log's Read Key (sec. 3.1.3).

**Event** A single node in a Merkle-Clock, stored on IPFS.

For any given Log, Events are authored by a single IPFS Peer, or *Writer*. This single-writer setup is a core feature of Logs, and provides properties unique to the Threads protocol. For clarity, we can similarly define a *Reader* as any other Peer capable of reading a Log. Related, a Merkle-Clock (see sec. 2.2.3) is simply a Merkle-DAG of *Events*.

### Listing 1 The Log multiaddress

```
// Log multiaddress
/ipel/12D3KooWC2zyCVron7AA34N6oKNtaXaZB51feG9rBkr7QbCcw

// Encapsulated multiaddress
/ip4/127.0.0.1/tcp/1234/p2p/12D..dwaA6Qe/ipel/12D..bCcW8ab

// Address book
[p2p/12D..dwaA6Qe/ipel/12D..bCcW8ab
/p2p/12D..dJT6nXY/ipel/12D..bCcW8ab // Follower]
```

### 3.1.2. Multi-addressed Event Logs

Together with a cryptographic signature, an Event is written to a log with an additional Node (see fig. 3) enabling Log verification by Readers (sec. 3.1.3). At a minimum, a Node must link to its most immediate ancestor. However, links to older ancestors are often included as well to improve concurrency during traversal and verification [43].

As shown in sec. 7.1, an Event's actual content (or body), is contained in a separate Node. This allows Events







**Listing 2** The AddrBook interface for storing Log addresses and the KeyBook interface for storing Log keys.

```

type AddrBook interface {
    AddAddr(thread.ID, peer.ID, ma.Multiaddr, time.Duration)
    AddAddrs(thread.ID, peer.ID, []ma.Multiaddr, time.Duration)
    SetAddr(thread.ID, peer.ID, ma.Multiaddr, time.Duration)
    SetAddrs(thread.ID, peer.ID, []ma.Multiaddr, time.Duration)
    UpdateAddrs(thread.ID, peer.ID, oldTTL time.Duration, newTTL time.Duration)
    Addrs(thread.ID, peer.ID) []ma.Multiaddr
    ClearAddrs(thread.ID, peer.ID)
}

type KeyBook interface {
    PubKey(thread.ID, peer.ID) ic.PubKey
    AddPubKey(thread.ID, peer.ID, ic.PubKey) error
    PrivKey(thread.ID, peer.ID) ic.PrivKey
    AddPrivKey(thread.ID, peer.ID, ic.PrivKey) error
    ReadKey(thread.ID, peer.ID) []byte
    AddReadKey(thread.ID, peer.ID, []byte) error
    FollowKey(thread.ID, peer.ID) []byte
    AddFollowKey(thread.ID, peer.ID, []byte) error
}

```

it is necessary to distribute each new key to users by including it in the header of the Event Block. Therefore, the Event Block itself is further encrypted using a *Read* key. The Read Key is not distributed within the Log itself but via a separate (secure) channel to all Peers who require access to the content of the Log.

**Read Key** The Read Key is a symmetric key created by the Log owner and used to encrypt the Content Key in each event.

Finally, the encrypted Event Block, its signature, and the IPLD linkage(s) from an Event to its antecedents are encrypted together using a Follow Key. Follow Keys allow Logs to be *followed* by peers on the network who do not have access to any content within the event. Followers can only see signatures and linkage(s) between Events.

**Follow Key** The Follow Key is a symmetric key created by the Log owner and used to encrypt the entire Event payload before adding the Event to the Log.

### 3.2. Threads

In Textile, the *interface* to Logs is managed as a Thread, which is a collection of Logs on a given topic. Threads are an event sourced, distributed database, and can be used to maintain a single, collaboratively edited, followed, or hosted dataset across multiple Peers. Threads provide the mechanism to combine multiple Logs from individual Writers into singular shared states through the use of either cross-Log sequencing (e.g. using a Bloom Clock, Merkle-Clock, or Hybrid Logical Clock [21]) or a CRDT (sec. 2.1.5).

#### 3.2.1. Identity

A unique Thread IDentity (TID) is used to group together Logs which compose a single dataset and as a topic identifier within Pubsub-based synchronization. The components of a TID are given in eq. 1.

$$\text{Thread ID} = \underbrace{0x02}_{\text{Multibase}} \underbrace{0x01}_{\text{Variant}} \underbrace{0x55}_{\text{Version}} \underbrace{0x39bc1dc03ee8cb5d478e41c}_{\text{Random Number}} \quad (1)$$

TIDs share some similarities with UUIDs [44] (version and variant) and IPFS-based CIDs and are multibase encoded<sup>10</sup> for maximum forward-compatibility. Base32 encoding is used by default, but any multibase-supported string encoding may be used.

**Multibase Prefix** The encoding type used by the multibase encoder. 1 byte.

**Version** ID format version. 8 bytes max. This allows future version to be backwards-compatible.

**Variant** Used to specify thread-level expectations, like access-control. 8 bytes max. See sec. 3.2.2 for more about variants.

**Random Number** A random number of a user-specified length. 16 bytes or more (see lst. 3).

#### 3.2.2. Variants

Certain TID *variants* may be more appropriate than others in specific contexts. For example, Textile provides an *access-controlled* Thread variant, which supports various collaborative structures — e.g., social media feeds, shared documents, blogs, photo albums, etc.

**Raw** This variant declares that consumers are not expected to make additional assumptions. This is the default variant (see lst. 3(a)).

**Access-Controlled** This variant declares that consumers should assume an access control list is composable from Log Events. The ACL represents a permissions rule set that must be applied when reading data (sec. 5.5.2 and lst. 3(b)).

**Listing 3** Identity variants.

```

# (a) Raw identity. V1, 128 bit
bafkxd5bjgi6k4zivuoxyxo4ua4mzyy

# (b) ACL enabled identity. V1, 256 bit.
bafyoioibghzefwldfrwqmqzz2ka66zgmdmgeobw2mimktr5jivsavv

```

#### 3.2.3. Log Synchronization

Log Writers, Readers, and Followers synchronize the state of their Logs by sending and receiving Events. Inspired by Git<sup>11</sup>, a reference to the latest Event in a Log is referred to as the *Head* (or sometimes the *root*).

<sup>10</sup><https://github.com/multiformats/multibase>

<sup>11</sup><https://git-scm.com/>

When a new Event is received, Readers and Followers simply advance their Head reference for the given Log. This is similar to how a system such as OrbitDB [42] works, except we are tracking *multiple* Heads (one per Log), rather than a single Head.

Regardless of the network protocol, Events are transported between Peers in a standardized *Event Envelope*. A new Thread is created by generating a TID and Log. The Log’s creator is the Writer, meaning it has possession of the Log’s Identity, Read, and Follow Keys. All of these keys are needed to compose Events. At this point, the Thread only exists on the Writer’s machine. Whether for collaboration, reading, or following, the process of sharing a Thread with other Peers starts by authoring a special Event called an *Invite*, which contains a set of keys from all of the Thread’s Logs, called a *Key Set*.

**Event Envelope** An over-the-wire message containing an Event and the sender’s signature of the Event.

**Invite** An Event containing a mapping of Log IDs to Key Sets, which can be used to join a Thread. Threads backed by an ACL (sec. 5.5.2) will also include the current ACL for the Thread in an Invite. This enables Peers to invite others to only read or follow a Thread, instead of becoming a full-on collaborator (i.e., a new Log Writer).

**Key Set** A set of keys for a Log. Depending on the context, a Key Set may contain the Follow and Read Key, or just the Follow Key. A Key Set is encrypted with the recipient’s public key.

The Invite is authored in the sender’s Log. Because the recipient does not yet have this Log’s Key Set, the Event is encrypted with the recipient’s public key. If the recipient accepts the Invite, they will author another special Event called a *Join* in a new Log of their own.

**Join** An Event containing an invitee’s new Log ID and Key Set, encrypted with the Key Set of the *inviting Peer’s Log*.

For a Join to be successful, all Log Writers must receive a copy of the new Key Set so they can properly handle future Events in the new Log. Instead of encrypting a Join with the public key of each existing Writer, we can encrypt a single Join with the Key Set of the inviting Peer’s Log, which the other Writers also have. Once a Peer has accepted an Invite, it will receive new Events from Log Writers. In cases where the invitee becomes a collaborator (i.e., a Writer) it is also responsible for sending its own Events out to the network.

### 3.2.3.1. Sending

Sending is performed in multiple phases because, invariably, some Thread participants will be offline or unresponsive:

1. New Events are pushed<sup>12</sup> directly to the Thread’s other Log Writers.
2. New Events are pushed directly to the target Log’s Follower(s), who may not maintain their own Log.
3. New Events are published over gossip-based Pubsub using TID as a topic, which provides potentially unknown Readers or Followers with an opportunity to consume Events in real-time.

Step 2 above allows for *additional* push mechanisms, as followers with public IP addresses become relays:

1. New Events may be pushed directly to web-based participants over a WebSocket.
2. New Events may be pushed to the Thread’s other Log Writers via federated notification services like Apple Push Notification Service (APNS), Google Cloud Messaging (GCM), Firebase Cloud Messaging (FCM), and/or Windows Notification Service (WNS).
3. New Events may trigger web-hooks, which could enable many complex (e.g., IFTTT<sup>13</sup>) workflows.

### 3.2.3.2. Receiving

There are multiple paths to receiving new Events, that together maximize connectivity between Peers who are often offline or unreachable.

1. Log Writers can receive Events directly from the Writer.
2. Events can be pulled from Followers via HTTP, libp2p, RSS, Atom, etc.
  1. In conjunction with push over WebSockets (seen in Step 2 of the additional push mechanisms above), this method provides web-based Readers and Followers with a reliable mechanism for receiving Log Events (fig. 5).
3. Writers and readers can receive new Events via a Pub/Sub subscription at the TID.

### 3.2.4. Log Replication

The notion of the Follow Key (sec. 3.1.3) makes duplicating all Log Events trivial. This allows any Peer on the network to be granted the responsibility of replicating data from another Peer without having read access to the raw Log entries. This type of Log replication can act as a data backup mechanism. It can also be used to build services that react to Log Events, potentially pushing data to disparate, non-Textile systems, especially if the replication service *is* granted read access to the Log Events (sec. 3.2.3).

## 4. THREADS INTERNALS

Previous sections have discussed the core features of the Textile Threads protocol. However, we have not

<sup>12</sup>Here push means “send to multiaddress(es)”, which may designate different protocols, e.g., P2P, HTTP, Bluetooth, etc.

<sup>13</sup><https://ifttt.com>

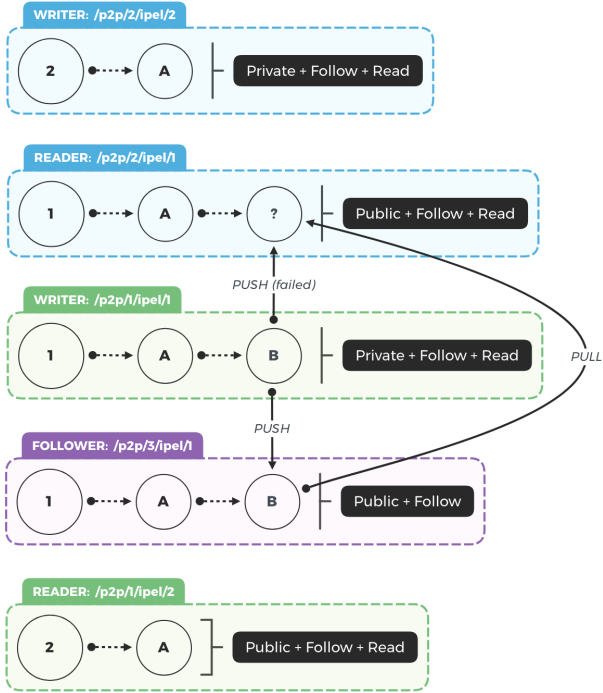


FIGURE 5: A pull-based request from a Follower.

yet discussed dealing with Log Events in practice. In this section, we provide a description of a Threads-compatible Event Store implementation. The Event Store outlined here takes advantage of ideas from several existing CQRS and ES systems (e.g., [45]), as well as concepts and designs from Flux [46], Redux<sup>14</sup> [47] and domain driven design [48] (DDD)<sup>15</sup>. Following this discussion of Threads *internals*, in sec. 5 we outline how it can be used to build intuitive developer-facing application programming interfaces (APIs) to make adopting and using Threads “the right choice” for a wide range of developers.

#### 4.1. Overview

Application state management tools such as Redux provide opinionated frameworks for making state mutations *predictable* by imposing certain restrictions on how and when updates can happen. The architecture of such an application shares several core features with CQRS and ES systems (see [46] and/or [49]). In both cases, the focus is on unidirectional data flows that build downstream views from atomic updates in the form of events (or actions).

We adopt a similar flow in Threads (see fig. 6). Like any CQRS/ES-based system, Threads are built on *Events*. Events are similar to actions in a Flux-based system, and are used to produce *predictable* updates to downstream state. Similarly to a DDD-based pattern, to add an

<sup>14</sup>Redux builds on concepts from CQRS and ES itself, and is arguably an implementation of the Flux application architecture.

<sup>15</sup><https://dddcommunity.org>

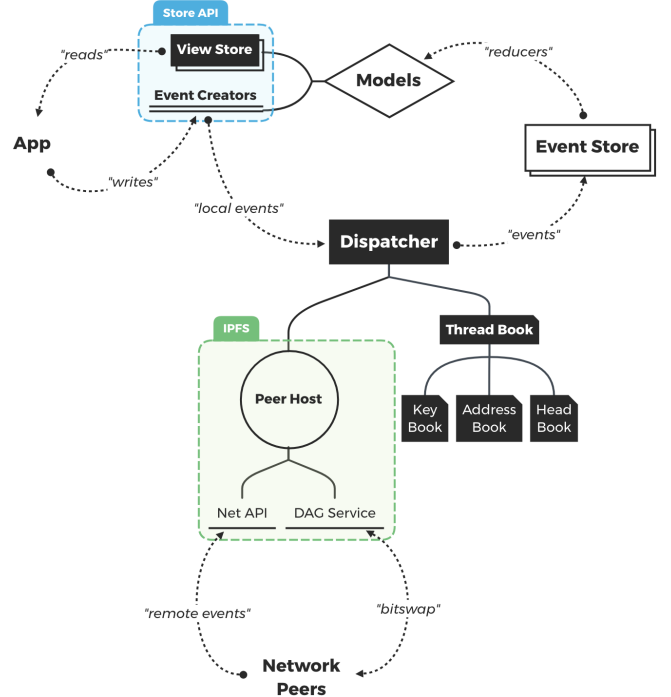


FIGURE 6: Architectural diagram for internal Event Store implementation.

Event to the internal Threads system, we use *Event Creators*, which dispatch Events to the system via a singleton *Dispatcher*. The Dispatcher then stores the derived Event in an *Event Store*, and calls a set of registered *Reducer* functions that mutate a (set of) downstream view *Stores* defined by a corresponding set of view *Models*, all within a single *Transaction*. This unidirectional, transaction-based system provides a flexible framework for building complex event-driven application logic.

##### 4.1.1. Events & Creators

*Events* are at the heart of Threads — every update to local and shared (i.e., across peers) state happens via Events (see also sec. 2.1.1). Events are used to describe “changes to an application state” [13] (e.g., a photo was added, an item was added to a shopping cart, etc). Related, *Event Creators* are used to send Events from a *local* application to an internal Event Store. If built around a specific *domain*, Creators provide bounded context that can be roughly compared to an aggregate root in DDD [48].

##### 4.1.2. Dispatcher

In order to persist and dispatch Events to downstream view Models, a *Dispatcher* is used. As shown in fig. 6, the Dispatcher is at the center of the ES system. All Events must go through the singleton Dispatcher, whether these be from local Event Creators or remote Peers. The Dispatcher is responsible for ensuring that incoming Events are persisted to the Event Store (the “source

of truth” for the system), as well as dispatched to downstream view Models by way of a set of Reducer functions.

#### 4.1.3. Transactions

All Reducer function calls (and “side effects”) due to a given Event happen within a single *Transaction* [50], in order to ensure consistency of both storage (Event Store) and Models. Transactions are similar to Redux *Sagas*<sup>16</sup> in terms of outcome, but with stronger consistency guarantees. Once an Event has been persisted in the internal Event Store, the Dispatcher is responsible for running a view Model’s Reducer callback. In order to make this possible, Models must *register* their Reducer with the Dispatcher.

#### 4.1.4. View Models/Stores

As in many CQRS-based systems, Events are dispatched on the write side, and are *reacted to* on the read side. Reactions happen via Reducer functions, which cause updates to view Models, which in turn provide interfaces for queries and accessing persisted view *Stores* (state). View Stores are then responsible for notifying downstream consumers (application logic) of changes to their state via a *Broadcaster* (i.e., event emitter). They are similar in some respects to a Redux Store, though it is possible to have *multiple* view Models/Stores as in the more general Flux pattern.

A view Model is generally defined by custom update logic (i.e., a Reducer), a (possibly ORM<sup>17</sup>-based) view Store for persistence, an Event Bus (or Broadcaster) for notifying downstream consumers, and a set of query/*Resolver* [45] functions. In practice, a view Model may *also* wrap the Event Creator application logic that is used to generate upstream Events. This provides an intuitive, singular access point to the *local* system, while also leaving room for updates via *external* Peer Events. In practice, the Store is a lightweight interface that can be implemented by one of many database management systems (see sec. 7.2) providing end users/developers the greatest flexibility.

#### 4.1.5. Remote Events

In addition to locally derived Events (i.e., from application logic and user interactions), Threads are designed so that Peers may collaborate in a given Thread via Events. Events generated by other, network peers are called *Remote Events*, and they enter the system via a Peer Host. The Peer Host is responsible for dealing with incoming Events (be they push or pull). These Events are no different from *Local Events*, though in practice the Peer Host is required to validate Remote Events before they are dispatched to the internal system.

## 5. THREAD INTERFACES

To make Threads as easy to adopt and use as possible, Textile has designed a developer facing API on top of the Threads internals that simplifies dealing with events and data, while still maintaining the power and flexibility of CQRS and ES. Developers should not have to learn a whole new set of terms and tools to take advantage of Threads’ capabilities. These simple, public-facing APIs will be comfortable to application developers looking to leverage distributed systems that leverage user-siloed data, with minimal configuration and maximum interoperability. Inspired by tools such as MondoDB<sup>18</sup>, Parse<sup>19</sup>, and Realm<sup>20</sup>, as well as the idea of bounded context and aggregate roots from DDD, we provide simple abstractions for performing distributed database operations as if they were local.

Indeed, the components of a Thread already provide several features that you would expect when operating on a dataset or table within a database: each Thread is defined by a unique ID, and provides facilities for access control and permissions, networking, and more. To illustrate how these underlying components can be combined to produce a simple API with minimal configuration and boilerplate, consider the following example in which we provide pseudo-code for a hypothetical Photos app.

### 5.1. Illustrative Example

To create a useful application, developers start with view **Models**, as in lst. 4. A Model is essentially the public API for the underlying view Models from sec. 4.1.4. Building on this, developers might create a new Thread for a user to store **Contact** information, as well as their mobile phone’s camera roll photos, as in lst. 5. This would create a new view Store “under-the-hood” (with corresponding indexes, etc), to be mutated by incoming Events.

The next step is to actually *create* and add a **Message** object to a shared album. In lst. 6, a message instance is created via the custom **Message** class, and then added to the shared **Dogs** Thread (which represents an album here). Behind the scenes, the Model (which is providing an Event Creator interface) is internally responsible for dispatching the Event through the local Dispatcher.

An example Event is given in lst. 7, and is the result of a new Message Event. Now, any updates to an existing Message instance will automatically generate the required underlying update Events. For example, lst. 8 shows the body text of the previous example being updated, and saved (committed) to the Thread. Behind the scenes, this event will be added to the User’s local Log, and pushed to any peers identified in the Thread’s

<sup>16</sup><https://redux-saga.js.org/>

<sup>17</sup>Object-relational mapping

<sup>18</sup><http://www.mongodb.com>

<sup>19</sup><https://parseplatform.org>

<sup>20</sup><https://realm.io>

**Listing 4** Create a photo entity and some way to represent a photo's author.

```
Photo = NewModel({
  _id: UUID,
  thumbnail: Buffer,
  original: Buffer,
});

Contact = NewModel({
  _id: UUID,
  name: {type: String, index: true},
  avatar: Photo,
});

// Photos may be grouped into messages.
Message = NewModel({
  _id: UUID,
  author: Contact,
  body: String,
  photos: [Photo],
})
```

**Listing 5** Create address book and camera roll stores.

```
AddressBook = NewDocumentStore("AddressBook")
// This store should take Contacts.
AddressBook.AddModel("Contact", Contact)
// If needed, additional models can be added...

// Create another store for camera roll photos.
CameraRoll = NewDocumentStore("CameraRoll")
// This store should take photos.
CameraRoll.AddModel("Photo", Photo)

// Create another store for a shared album.
MyDogsAlbum = NewDocumentStore("Dogs")
MyDogsAlbum.AddModel("Message", Message)

// Messages can also be nested.
Message.AddModel("Message", Message)
```

ACL document (see secs. 3.2.2, 5.4). In practice, the Model generates another Event that carries the diff and a document identifier.

All instances and models in Threads have several special methods and properties specific to the Threads API, several of which we will explore here. Lst. 9 demonstrates several features common in a Threads-based workflow, including queries, creating invites and changing permissions/access control (see also sec. 5.5.2), as well as subscribing to updates and changes at various levels of the Threads API. These subscriptions would enable downstream consumers (views, front-end stores, etc.) to receive updates as changes to the Thread are made via underlying Events.

**Listing 6** Adding data to a shared thread.

```
// Create a message with a photo.
MyMessage = Message.Create({
  author: <author_id>,
  body: "This is Lucas.",
  photos: [{
    thumbnail: <buffer>,
    original: <buffer>
  }]
})
// Now it can be added to Dogs "album".
MyDogsAlbum.Add(MyMessage)
```

**Listing 7** A new Message Event.

```
{
  "body": {
    "data": {
      "author_id": <author_id>,
      "body": "This is Lucas.",
      "photos": [{
        "thumbnail": <buffer>,
        "original": <buffer>
      }]
    }
  },
  "header": {
    "time": 1569434034737,
    "key": "215bs...1DXJ"
  }
}
```

## 5.2. Modules

One of Textile's stated goals is to allow individuals to better capture the value of their data while still enabling developers to build new interfaces and experiences on top of said data. A key piece of this goal is to provide *inter-application* data access to the *same underlying user-siloed data*. In order to meet this goal, it is necessary for developers to be using the same data structure and conventions when building their apps. In conjunction with community developers, Textile will provide a number of *Modules* designed to wrap a given domain (e.g., Photos) into a singular software package to facilitate this. This way, developers need only agree on the given data Module in order to provide seamless inter-application experiences. For example, any developer looking to provide a view on top of a user's Photos (perhaps their phone's camera roll) may utilize the Photos Module (which may be designed as in the example above). They may also extend this Module, to provide additional functionality.

In building on top of an existing Module, developers ensure other application developers are also able to interact with the data produced by their app. This enables tighter coupling between applications, and it



**Listing 8** Message updates are persisted and transmitted automatically.

```
MyMessage.body = "Actually, this is Fido."
MyMessage.Save()

// New Event with diff and document id
{
  "body": {
    "data": {
      "doc_id": MyMessage._id,
      "body": "Actually, this is Fido."
    }
  },
  "header": {
    "time": 1569434035737,
    "key": "iJMfqWy...1qfJyc29RS"
  }
}
```

**Listing 9** Additional Threads-based API functionality.

```
// Query for message, select only thumbnail.
Dogs.FindOne({ "_id": MyMessage._id }, "thumbnail")

// Every (Event increments version tag
MyMessage.Version()

// All doc changes in Dogs Thread
Dogs.subscribe()
// All changes to all Messages in all Threads
Message.Subscribe()
// Changes specific to this document
MyMessage.Subscribe()

// Create invite Event IFF User has permission
Dogs.Grant(<peer_id>, <role>)

// Alter ACL OR create invite if needed/allowed
MyMessage.Grant(<peer_id>, <role>)
```

allows for smaller apps that can focus on a very specific user experience (say, filters on Photos). Furthermore, it provides a *logically centralized*, platform-like developer experience, without the actual centralized infrastructure. APIs for Photos, Messages/Chat, Music, Video, Storage, etc are all possible, extensible, and available to all developers. This is a powerful concept, but it is also flexible. For application developers working on very specific or niche apps with less need for inter-application usability, Modules are not needed, and they can instead focus on custom Models. However, it is likely that developers who build on openly available *standard* Modules will provide a more useful experience for their users, and will benefit from the *network effects* [51] produced by many interoperable apps.

### 5.3. Databases

With these interface simplifications, it is not difficult to imagine even higher-level APIs in which Threads are exposed via interfaces compatible with *existing* datastores or DBMS. Here we draw inspiration from similar projects (e.g., OrbitDB [42]) which have made it much easier for developers familiar with centralized database systems to make the move to decentralized systems such as Threads. For example, a key-value store built on Threads would “map” key-value operations, such as Put, Get, and Del to an internal (i.e., private) Model as in the previous section, with similarly defined methods. The generated Events would then mutate the internal map-like view Model effectively encapsulating the entire Event Store in a database structure that satisfies a given interface (see lst. 10 for example). These too would be distributed as Modules, making it easy for developers to swap in/substitute existing backend infrastructure.

**Listing 10** A proposed key-value store interface

```
type TextileKVStore interface {
  Put(key string, value Node) error
  Get(key string) (Node, error)
  Del(key string) error
}
```

Other database abstractions include a no-sql style document store for storing and indexing arbitrary structs and/or JSON documents. The interface for such as store, again built using a “wrapped” view Model, might look like lst. 11, where Indexable could be satisfied by any structure with a Key field and Query might be taken from the go-datastore interface library<sup>21</sup> or similar.

**Listing 11** A proposed document store interface

```
type TextileDocStore interface {
  Put(doc Indexable) error
  Get(key string) (Indexable, error)
  Del(key string) error
  Query(query Query) ([]Indexable, error)
}
```

Similar abstractions will be used to implement additional database types and functions. Tables, feeds, counters, and other simple stores can also be built on Threads. Each database style would be implemented as a standalone wrapper/software library, allowing application developers to pick and choose the solution most useful to the application at hand. Similarly, more advanced applications could be implemented using a combination of database types, or by examining the source code of these *reference* libraries.

<sup>21</sup><https://github.com/ipfs/go-datastore>

## 5.4. CRDTs

Eventually consistent, CRDT-based structures can also be implemented on top of Threads' Event-driven architecture. CRDT-based Stores are particularly useful for managing views of a document in a multi-peer collaborative editing environment (like Google Docs or similar). For example to support offline-first, potentially concurrent edits on a shared JSON document, one could implement a JSON CRDT datatype [52] that merges updates to a JSON document in a view Model's Reducer function. Libraries such as Automerge<sup>22</sup> provide useful examples of reducer functions that make working with JSON CRDTs relatively straightforward, and implementations in other programming languages are also available [...]. A practical example of using a JSON CRDT in Threads is given in sec. 5.5.2, where it is used to represent updates to an ACL document defined as a default view Model, with interfaces defined for an access-controlled Threads implementation.

## 5.5. Thread Extensions

The Textile protocol provides a distributed framework for building shared, offline first, Stores that are fault tolerant<sup>23</sup>, eventually consistent, and scalable. Any internal implementation details of a compliant Threads *client* may use any number of well-established design patterns from the CQRS and ES (and related) literature to *extend* the Threads protocol with additional features and controls. Indeed, by designing our system around Events, a Dispatcher, and generic Stores, we make it easy to extend Threads in many different ways. Some extensions included by Textile's Threads implementation are outlined in this section to provide some understanding of the extensibility this design affords.

### 5.5.1. Snapshots and Compaction

Snapshots<sup>24</sup> are simply the current state of a Store at a given point in time. They can be used to rebuild the state of a view Store without having to query and replay all previous Events. When a Snapshot is available, a Thread Peer can rebuild the state of a given view Store/Model by replaying only Events generated since the latest Snapshot using the Model's Reducer function. Multiple Peers processing the same Log could create a Snapshot every 1000 Events and be guaranteed to create the exact same Snapshot because each Peer's Event counts are identical<sup>25</sup>.

In practice, Snapshots are written to their own internal Event Store and stored locally. They can potentially

be synced (sec. 3.2.3) to other Peers as a form of data backup or to optimize state initialization when a new Peer starts participating in a shared Thread (saving disk space, bandwidth, and time). They can similarly be used for initializing a local view Store during recovery.

Compaction is a local-only operation (i.e., other Peers do not need to be aware that Compaction was performed) performed on an Event Store to free up local disk space. As a result, it can speed up re-hydration of a downstream Stores's state by reducing the number of Events that need to be processed. Compaction is useful when only the latest Event of a given type is required.

### 5.5.2. Access Control

One of the most important properties of a shared data model is the ability to apply access control rules. There are two forms of access control possible in Threads, Entity-level ACLs and Thread-level ACLs. Thread-level access control lists (ACLs) allow creators to specify who can *follow*, *read*, *write*, and *delete* Thread data. Similarly, Entity-level ACLs provide more granular control to Thread-writers on a per-Entity (see def. 4) basis. Both types of ACLs are implemented as JSON CRDTs (see sec. 5.4) wrapped in a custom view Model (see sec. 5). ACLs implemented as JSON Models provide two advantages over static or external ACL rules (although static and external ACLs are also possible). First, ACLs are fully mutable, allowing developers to create advanced rules for collaboration with any combination of readers, writers, and followers. Second, because ACLs are essentially mutable JSON documents, they can specify their *own editing rules* (i.e. allowing multiple Thread participants to modify the ACL) in a self-referencing way.

**Entity** An Entity is made of of a series of ordered Events referring to a specific entity or object. For example, an ACL JSON document is a single entity made up of a sequence of Thread Events that describe a JSON document. An Entity might have a unique UUID (see lst. 12) which can be referenced across/within Event updates.

---

### Listing 12 Entity Id.

---

```
// UUID
bafykrq5i25vd64ghamtgus6lue74k
```

---

Textile's Threads includes ACL management tooling based on a *Role-based access control* [53] pattern, wherein individuals or groups are assigned roles which carry specific permissions. Roles can be added and removed as needed. Textile ACLs can make use of five distinct roles<sup>26</sup>: *No-access*, *Follow*, *Read*, *Write*, and *Delete*.

---

<sup>26</sup>By default, Threads without access control operate similar to Secure Scuttlebutt (SSB; where every peer consumes what they want and writes what they want).

<sup>22</sup><https://github.com/automerge/automerge>

<sup>23</sup>When using an ACID compliant backing store for example.

<sup>24</sup>The literature around snapshots and other CQRS and ES terms is somewhat confusing, we attempt to use the most common definitions here.

<sup>25</sup>Assuming any network partitions are only short-lived (i.e., that peers are able to share events consistently).



**No-access** No access is permitted. This is the default role.

**Follow** Access to Log Follow Keys is permitted. Members of this role are able to verify Events and follow linkages. The Follow role is used to designate a “follower” peer for offline replication and/or backup.

**Read** Access to Log Read Keys is permitted in addition to Follow Keys. Members of this role are able to read Log Event payloads.

**Write** Members of this role are able to author new Events, which also implies access to Log Follow and Read Keys. At the Thread-level, this means authoring a Log. At the document-level, the Write role means that Events in this Log are able to target a particular document.

**Delete** Members of this role are able to delete Events, which implies access to Log Follow Keys. In practice, this means marking an older Event as “deleted”.

A typical Thread-level ACL JSON (see lst. 13) can be persisted to a local Event Store as part of the flow described in sec. 4. See also sec. 5, and in particular lst. 9 for the public API for editing ACL definitions.

---

**Listing 13** ACL JSON document.

---

```
{
  "_id": "bafykrq5i25vd64ghamtgus6lue74k",
  "default": "no-access",
  "peers": {
    "12D..dwaA6Qe": ["write", "delete"],
    "12D..dJT6nXY": ["follow"],
    "12D..P2c6ifo": ["read"],
  }
}
```

---

The `default` key states the default role for all network peers. The `peers` map is where roles are delegated to specific peers. Here, `12D..dwaA6Qe` is likely the owner, `12D..dJT6nXY` is a designated follower, and `12D..P2c6ifo` has been given read access. A Thread-level ACL has its own document ACL, which also applies to all other document ACLs (see lst. 14). This means that only `12D..dwaA6Qe` is able to alter the access-control list.

---

**Listing 14** Thread and Entity ACL

---

```
{
  "_id": "bafykrq5i25vd64ghamtgus6lue74k-acl",
  "default": "no-access",
  "peers": {
    "12D..dwaA6Qe": ["write", "delete"],
  }
}
```

---

## 6. CONCLUSION

In this paper, we described the challenges and considerations when creating a protocol suitable for large-scale data storage, synchronization, and use in a distributed system. We identified six requirements for enabling *user-siloed* data: flexible data formats, efficient synchronization, conflict resolution, access-control, scalable storage, and network communication. We have introduced a solution to these requirements that extends on IPFS and prior research done by Textile and others, which we term Threads. Threads are a novel data architecture that builds upon a collection of protocols to deliver a scalable and robust storage system for end-user data.

We show that the flexible core structure of single-writer append-only logs can be used to compose higher-order structures such as Threads, Views, and/or CRDTs. In particular, we show that through the design of specific default view Models, we can support important features such as access control lists and common, specialized, or complex data models. The Threads protocol described here is flexible enough to derive numerous specific database types (e.g. key/value stores, document stores, relational stores, etc) and model an unlimited number of applications states. The cryptography used throughout Threads will help shift the data ownership model from apps to users.

### 6.1. Future Work

The research and development of Textile Threads has highlighted several additional areas of work that would lead to increased benefits for users and developers. In particular, we have highlighted network services and security enhancements as core future work. In the following two sections, we briefly outline planned future work in these critical areas.

#### 6.1.1. Enhanced Log Security

The use of a single Read and Follow Key for an entire Log means that, should either of these keys be leaked via malicious (or other/accidental) means, there is no way to prevent a Peer with the leaked keys from listening to Events or traversing the Log history. Potential solutions currently being explored by Textile developers include key rotation at specific Event offsets [54], and/or incorporating the Double Ratchet Algorithm [55] for forward secrecy [56].

#### 6.1.2. Tighter Coupling with Front End Models

Implementing Threads internals (see sec. 4) using similar patterns to common frontend workflows (e.g., Redux) presents opportunities for tighter coupling between “backend” logic and frontend views. This is a major advantage of tools such as `reSolve`<sup>27</sup>, where “system

---

<sup>27</sup><https://reimagined.github.io/resolve/>

changes can be reflected immediately [on the frontend], without the need to re-query the backend” [45]. Textile developers will create frameworks to more directly expose the internals of Threads to frontend SDKs (or DBMS, see sec. 7.2), making it possible to sync application state across and between apps and services on the IPFS network.

### 6.1.3. Textile: The Thread & Bot Network

Threads change the relationship between a user, their data, and the services they connect with that data. The nested, or multi-layered, encryption combined with powerful ACL capabilities create new opportunities to build distributed services, or Bots, in the network of IPFS peers. Based on the Follow Key now available in Threads, Bots can relay, replicate, or store data that is synchronized via real-time updates in a *trustless*, partially trusted, or fully-trusted way. Bots can additionally enhance the IPFS network by providing a framework to build and deploy many new kinds of services available over HTTP or P2P. Services could include simple data archival, caching and republishing, translation, data conversion, and more. Advanced examples could include payment, re-encryption, or bridges to Web 2.0 services to offer decentralized access to Web 2.0.

## 7. APPENDIX

### 7.1. Event node interface

```
// Node is the most basic component of a log.
// Note: In practice, this is encrypted with the Follow Key.
type Node interface {
    ipld.Node

    // Event is the Log update.
    Event() Event

    // Refs are node linkages.
    Refs() []cid.Cid

    // Sig is a cryptographic signature of Event and Refs
    // created with the Log's private key.
    Sig() []byte
}

// Event represents the content of an update.
// Note: In practice, this is encrypted with the Recipient Key.
type Event interface {
    ipld.Node

    // Header provides a means to store a timestamp
    // and a key needed for decryption.
    Header() EventHeader

    // Body contains the content of an update.
    // In practice, this is encrypted with the Header Key
    // or the recipient's public key.
    Body() ipld.Node

    // Decrypt is a helper function that decrypts Body
    // with a key in Header.
    Decrypt() (ipld.Node, error)
}
```

*// EventHeader contains Event metadata.*

```
type EventHeader interface {
    ipld.Node

    // Time is the wall-clock time at which the Event
    // was created.
    Time() int

    // Key is an optional single-use symmetric key
    // used to encrypt Body.
    Key() []byte
}
```

### 7.2. Database Management Systems (DBMS)

1. Relational
  1. MariaDB <https://mariadb.org/>
  2. PostgreSQL <https://www.postgresql.org/>
  3. SQLite <https://www.sqlite.org>
2. Key-value stores
  1. Dynamo <https://aws.amazon.com/dynamodb/>
  2. LevelDB <https://github.com/google/leveldb>
  3. Redis <https://redis.io/>
3. Document stores
  1. CouchDB <http://couchdb.apache.org/>
  2. MongoDB <https://www.mongodb.com/>
  3. RethinkDB <https://rethinkdb.com/>

## REFERENCES

- [1] C. Dixon, “Why Decentralization Matters,” *Medium*, 26-Oct-2018. [Online]. Available: <https://medium.com/s/story/why-decentralization-matters-5e3f79f7638e>. [Accessed: 19-Sep-2019].
- [2] T. Berners-Lee and K. O’Hara, “The read–write Linked Data Web,” *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 371, no. 1987, p. 20120513, Mar. 2013.
- [3] Y. K. de Montjoye, S. S. Wang, A. Pentland, D. T. T. Anh, and A. Datta, “On the Trusted Use of Large-Scale Personal Data.” *IEEE Data Eng. Bull.*, vol. 35, no. 4, pp. 5–8, 2012.
- [4] A. V. Sambra et al., *Solid: A Platform for Decentralized Social Applications Based on Linked Data*. 2016.
- [5] F. B. Schneider, “Implementing fault-tolerant services using the state machine approach: A tutorial,” *ACM*

- Computing Surveys (CSUR)*, vol. 22, no. 4, pp. 299–319, 1990.
- [6] Jay Kreps, “The Log: What every software engineer should know about real-time data’s unifying abstraction,” 16-Dec-2013. [Online]. Available: <https://engineering.linkedin.com/distributed-systems/log-what-every-software-engineer-should-know-about-real-time-datas-unifying> [Accessed: 19-Sep-2019].
- [7] D. Betts, J. Dominguez, G. Melnik, F. Simonazzi, and M. Subramanian, *Exploring CQRS and Event Sourcing: A Journey into High Scalability, Availability, and Maintainability with Windows Azure*, 1st ed. Microsoft patterns & practices, 2013.
- [8] “Event Store.” [Online]. Available: <https://eventstore.org/>. [Accessed: 19-Sep-2019].
- [9] “Apache Kafka.” [Online]. Available: <https://kafka.apache.org/>. [Accessed: 19-Sep-2019].
- [10] “Apache Samza.” [Online]. Available: <http://samza.apache.org/>. [Accessed: 19-Sep-2019].
- [11] M. Kleppmann, *Designing data-intensive applications: The big ideas behind reliable, scalable, and maintainable systems*. " O'Reilly Media, Inc.", 2017.
- [12] Martin Fowler, “CQRS,” *martinfowler.com*, 14-Jul-2011. [Online]. Available: <https://martinfowler.com/bliki/CQRS.html>. [Accessed: 20-Sep-2019].
- [13] M. Fowler, “Event Sourcing,” *martinfowler.com*. [Online]. Available: <https://martinfowler.com/eaDev/EventSourcing.html>. [Accessed: 19-Sep-2019].
- [14] Microsoft Corporation, “Azure Application Architecture Guide,” *Microsoft Docs*. [Online]. Available: <https://docs.microsoft.com/en-us/azure/architecture/guide/>. [Accessed: 19-Sep-2019].
- [15] E. Brewer, “Towards Robust Distributed Systems,” in *19th ACM Symposium on Principles of Distributed Computing (PODC)*, 2000.
- [16] S. Gilbert and N. Lynch, “Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services,” *Acm Sigact News*, vol. 33, no. 2, pp. 51–59, 2002.
- [17] M. Shapiro, N. Preguiça, C. Baquero, and M. Zawirski, “A comprehensive study of Convergent and Commutative Replicated Data Types,” report, Jan. 2011.
- [18] P. S. Almeida, A. Shoker, and C. Baquero, “Delta state replicated data types,” *Journal of Parallel and Distributed Computing*, vol. 111, pp. 162–173, Jan. 2018.
- [19] R. Schwarz and F. Mattern, “Detecting causal relationships in distributed computations: In search of the holy grail,” *Distrib Comput*, vol. 7, no. 3, pp. 149–174, Mar. 1994.
- [20] S. Katz and D. Peled, “Interleaving set temporal logic,” *Theoretical Computer Science*, vol. 75, no. 3, pp. 263–287, 1990.
- [21] S. S. Kulkarni, M. Demirbas, D. Madappa, B. Avva, and M. Leone, “Logical Physical Clocks,” in *Principles of Distributed Systems*, vol. 8878, M. K. Aguilera, E. Quera, and M. Shapiro, Eds. Cham: Springer International Publishing, 2014, pp. 17–32.
- [22] L. Lamport, “Time, clocks, and the ordering of events in a distributed system,” *Commun. ACM*, vol. 21, no. 7, pp. 558–565, Jul. 1978.
- [23] L. Ramabaja, “The Bloom Clock,” May 2019.
- [24] H. Sanjuan, S. Poyhtari, and P. Teixeira, “Merkle-CRDTs,” May-2019.
- [25] V. Enes, P. S. Almeida, and C. Baquero, “The Single-Writer Principle in CRDT Composition,” in *Proceedings of the Programming Models and Languages for Distributed Computing on - PMLDC ’17*, 2017, pp. 1–3.
- [26] R. Mört, *Content Based Addressing : The case for multiple Internet service providers*. 2012.
- [27] J. Benet, “IPFS: Content addressed, versioned, p2p file system,” *arXiv preprint arXiv:1407.3561*, vol. (Draft 3), 2014.
- [28] M. Selimi and F. Freitag, “Tahoe-LAFS Distributed Storage Service in Community Network Clouds,” in *2014 IEEE Fourth International Conference on Big Data and Cloud Computing*, 2014, pp. 17–24.
- [29] S. C. Rhea, R. Cox, and A. Pesterev, “Fast, Inexpensive Content-Addressed Storage in Foundation.” in *USENIX Annual Technical Conference*, 2008, pp. 143–156.
- [30] Protocol Labs, “Multihash,” *Multiformats*. [Online]. Available: <https://multiformats.io/>. [Accessed: 27-Sep-2019].
- [31] Protocol Labs, “Filecoin: A Decentralized Storage Network,” 19-Jul-2017.
- [32] T. Berners-Lee, “Linked Data,” *Design Issues*, 18-Jun-2009. [Online]. Available: <https://www.w3.org/DesignIssues/LinkedData.html>. [Accessed: 20-Sep-2019].
- [33] C. Bizer, T. Heath, and T. Berners-Lee, “Linked data: The story so far,” in *Semantic services, interoperability and web applications: Emerging concepts*, IGI Global, 2011, pp. 205–227.
- [34] T. Heath and C. Bizer, “Linked Data: Evolving the Web into a Global Data Space,” *Synthesis Lectures on the Semantic Web: Theory and Technology*, vol. 1, no. 1, pp. 1–136, Feb. 2011.

- [35] Brendan O’Brien and Michael Hucka, “Deterministic Querying for the Distributed Web,” Nov-2017.
- [36] P. Srisuresh and M. Holdrege, “IP Network Address Translator (NAT) Terminology and Considerations.” [Online]. Available: <https://tools.ietf.org/html/rfc2663>. [Accessed: 20-Sep-2019].
- [37] Protocol Labs, “Multiaddr,” *Multiformats*. [Online]. Available: <https://multiformats.io/>. [Accessed: 20-Sep-2019].
- [38] Secure Scuttlebutt, “Scuttlebutt Protocol Guide.” [Online]. Available: <https://ssbc.github.io/scuttlebutt-protocol-guide/>. [Accessed: 11-Sep-2019].
- [39] Eric Harris-Braun, Nicolas Luck, and Arthur Brock, “Holochain: Scalable agent-centric distributed computing,” Ceptur LLC, 15-Feb-2018.
- [40] R. W. Shirey, “Internet Security Glossary, Version 2,” *Network Working Group*, Aug-2007. [Online]. Available: <https://tools.ietf.org/html/rfc4949>. [Accessed: 20-Sep-2019].
- [41] A. Herzberg, Y. Mass, J. Mihaeli, D. Naor, and Y. Ravid, “Access control meets public key infrastructure, or: Assigning roles to strangers,” in *Proceeding 2000 IEEE Symposium on Security and Privacy. S P 2000*, 2000, pp. 2–14.
- [42] Mark Robert Henderson, Samuli Pöyhtäri, Vesa-Ville Piironen, Juuso Räsänen, Shams Methnani, and Richard Littauer, “The OrbitDB Field Manual,” Haja Networks Oy, 26-Sep-2019.
- [43] A. Meyer, “Bamboo,” 16-Sep-2019. [Online]. Available: <https://github.com/AljoschaMeyer/bamboo>. [Accessed: 20-Sep-2019].
- [44] P. J. Leach, M. Mealling, and R. Salz, “A Universally Unique IDentifier (UUID) URN Namespace,” *Network Working Group*, Jul-2005. [Online]. Available: <https://tools.ietf.org/html/rfc4122>. [Accessed: 20-Sep-2019].
- [45] R. Eremin, “A Redux-Inspired Backend,” *Medium*, 14-Jan-2019. [Online]. Available: <https://medium.com/resolvejs/resolve-redux-backend-ebcfc79bbbea>. [Accessed: 24-Sep-2019].
- [46] Facebook, “Flux: In-Depth Overview,” 2019. [Online]. Available: <http://facebook.github.io/flux/docs/in-depth-overview>. [Accessed: 23-Sep-2019].
- [47] Redux, “Motivation,” *Redux*. [Online]. Available: <https://redux.js.org/introduction/motivation>. [Accessed: 20-Sep-2019].
- [48] E. Evans, *Domain-driven design: Tackling complexity in the heart of software*. Addison-Wesley Professional, 2004.
- [49] D. Abramov, “The Case for Flux,” *Medium*, 03-Nov-2015. [Online]. Available: <https://medium.com/swlh/the-case-for-flux-379b7d1982c6>. [Accessed: 23-Sep-2019].
- [50] T. Haerder and A. Reuter, “Principles of Transaction-oriented Database Recovery,” *ACM Comput. Surv.*, vol. 15, no. 4, pp. 287–317, Dec. 1983.
- [51] C. Shapiro, S. Carl, and H. R. Varian, *Information rules: A strategic guide to the network economy*. Harvard Business Press, 1998.
- [52] M. Kleppmann and A. R. Beresford, “A Conflict-Free Replicated JSON Datatype,” *IEEE Trans. Parallel Distrib. Syst.*, vol. 28, no. 10, pp. 2733–2746, Oct. 2017.
- [53] R. S. Sandhu, E. J. Coyne, H. L. Feinstein, and C. E. Youman, “Role-based access control models,” *Computer*, vol. 29, no. 2, pp. 38–47, Feb. 1996.
- [54] HashiCorp, “Key Rotation.” [Online]. Available: <https://www.vaultproject.io/docs/internals/rotation.html>. [Accessed: 20-Sep-2019].
- [55] M. Marlinspike, “The Double Ratchet Algorithm,” vol. Revision 1, p. 35, Nov. 2016.
- [56] N. Unger *et al.*, “SoK: Secure Messaging,” in *2015 IEEE Symposium on Security and Privacy*, 2015, pp. 232–249.