

IPFS - Content-Addressed, Versioned, P2P File System (DRAFT 3)

Juan Benet
juan@benet.ai

ABSTRACT

The InterPlanetary File System (IPFS) is a peer-to-peer distributed file system that seeks to connect all computing devices with the same system of files. In some ways, IPFS is similar to the Web, but IPFS could be seen as a single BitTorrent swarm, exchanging objects within one Git repository. In other words, IPFS provides a high throughput content-addressed block storage model, with content-addressed hyperlinks. This forms a generalized Merkle DAG, a data structure upon which one can build versioned file systems, blockchains, and even a Permanent Web. IPFS combines a distributed hashtable, an incentivized block exchange, and a self-certifying namespace. IPFS has no single point of failure, and nodes do not need to trust each other.

1. INTRODUCTION

There have been many attempts at constructing a global distributed file system. Some systems have seen significant success, and others failed completely. Among the academic attempts, AFS [8] has succeeded widely and is still in use today. Others [9, 5] have not attained the same success. Outside of academia, the most successful systems have been peer-to-peer file-sharing applications primarily geared toward large media (audio and video). Most notably, Napster, KaZaA, and BitTorrent [2] deployed large file distribution systems supporting over 100 million simultaneous users. Even today, BitTorrent maintains a massive deployment where tens of millions of nodes churn daily [27]. These applications saw greater numbers of users and files distributed than their academic file system counterparts. However, the applications were not designed as infrastructure to be built upon. While there have been successful repurposings¹, no general file-system has emerged that offers global, low-latency, and decentralized distribution.

Perhaps this is because a “good enough” system for most use cases already exists: HTTP. By far, HTTP is the most successful “distributed system of files” ever deployed. Coupled with the browser, HTTP has had enormous technical and social impact. It has become the de facto way to transmit files across the Internet. Yet, it fails to take advantage of dozens of brilliant file distribution techniques invented in the last fifteen years. From one perspective, evolving Web infrastructure is nearly impossible, given the number of backwards compatibility constraints and the number of

strong parties who have invested in the current model. But from another perspective, new protocols have emerged and gained wide use since the emergence of HTTP. What is lacking is upgrading design: enhancing the current HTTP web, and introducing new functionality without degrading user experience.

Industry has gotten away with using HTTP this long because moving small files around is relatively cheap, even for small organizations with lots of traffic. But we are entering a new era of data distribution with new challenges: (a) hosting and distributing petabyte datasets, (b) computing on large data across organizations, (c) high-volume high-definition on-demand or real-time media streams, (d) versioning and linking of massive datasets, (e) preventing accidental disappearance of important files, and more. Many of these can be boiled down to “lots of data, accessible everywhere.” Pressed by critical features and bandwidth concerns, we have already given up HTTP for different data distribution protocols. The next step is making them part of the Web itself.

Orthogonal to efficient data distribution, version control systems have managed to develop important data collaboration workflows. Git, the distributed source code version control system, developed many useful ways to model and implement distributed data operations. The Git toolchain offers versatile versioning functionality that large file distribution systems severely lack. New solutions inspired by Git are emerging, such as Camlistore [3], a personal file storage system, and Dat [17] a data collaboration toolchain and dataset package manager. Git has already influenced distributed file system design [11], as its content-addressed Merkle DAG data model enables powerful file distribution strategies. What remains to be explored is how this data structure can influence the design of high-throughput oriented file systems, and how it might upgrade the Web itself.

This paper introduces IPFS, a novel peer-to-peer version-controlled file system seeking to reconcile these issues. IPFS synthesizes learnings from many past successful systems. Careful interface-focused integration yields a system greater than the sum of its parts. The central IPFS principle is modeling *all data* as part of the same Merkle DAG.

2. BACKGROUND

This section reviews important properties of successful peer-to-peer systems, which IPFS combines.

2.1 Distributed Hash Tables

Distributed Hash Tables (DHTs) are widely used to coor-

¹For example, Linux distributions use BitTorrent to transmit disk images, and Blizzard, Inc. uses it to distribute video game content.

dinate and maintain metadata about peer-to-peer systems. For example, the BitTorrent Mainline DHT tracks sets of peers part of a torrent swarm.

2.1.1 Kademia DHT

Kademlia [12] is a popular DHT that provides:

1. Efficient lookup through massive networks: queries on average contact $\lceil \log_2(n) \rceil$ nodes (e.g. 20 hops for a network of 1,000,000 nodes).
2. Low coordination overhead: it optimizes the number of control messages it sends to other nodes.
3. Resistance to various attacks by preferring long-lived nodes.
4. Wide usage in peer-to-peer applications, including Gnutella and BitTorrent, forming networks of over 20 million nodes [27].

2.1.2 Coral DSHT

While some peer-to-peer file systems store data blocks directly in DHTs, this “wastes storage and bandwidth, as data must be stored at nodes where it is not needed” [7]. The Coral DSHT extends Kademia in three particularly important ways:

1. Kademia stores values in nodes whose identifiers are “nearest” (using XOR-distance) to the key. This does not take into account application data locality, ignores “far” nodes that may already have the data, and forces “nearest” nodes to store it, whether they need it or not. This wastes significant storage and bandwidth. Instead, Coral stores addresses to peers who can provide the data blocks.
2. Coral relaxes the DHT API from `get_value(key)` to `get_any_values(key)` (the “sloppy” in DSHT). This still works since Coral users only need a single (working) peer, not the complete list. In return, Coral can distribute only subsets of the values to the “nearest” nodes, avoiding hot-spots (overloading *all the nearest nodes* when a key becomes popular).
3. Additionally, Coral organizes a hierarchy of separate DSHTs called *clusters* depending on region and size. This enables nodes to query peers in their region first, “finding nearby data without querying distant nodes” [7] and greatly reducing the latency of lookups.

2.1.3 S/Kademia DHT

S/Kademia [1] extends Kademia to protect against malicious attacks in two particularly important ways:

1. S/Kademia provides schemes to secure `NodeId` generation, and prevent Sybil attacks. It requires nodes to create a PKI key pair, derive their identity from it, and sign their messages to each other. One scheme includes a proof-of-work crypto puzzle to make generating Sybils expensive.
2. S/Kademia nodes lookup values over disjoint paths, in order to ensure honest nodes can connect to each other in the presence of a large fraction of adversaries in the network. S/Kademia achieves a success rate of 0.85 even with an adversarial fraction as large as half of the nodes.

2.2 Block Exchanges - BitTorrent

BitTorrent [4] is a widely successful peer-to-peer filesharing system, which succeeds in coordinating networks of untrusting peers (swarms) to cooperate in distributing pieces of files to each other. Key features from BitTorrent and its ecosystem that inform IPFS design include:

1. BitTorrent’s data exchange protocol uses a quasi tit-for-tat strategy that rewards nodes who contribute to each other, and punishes nodes who only leech others’ resources.
2. BitTorrent peers track the availability of file pieces, prioritizing sending rarest pieces first. This takes load off seeds, making non-seed peers capable of trading with each other.
3. BitTorrent’s standard tit-for-tat is vulnerable to some exploitative bandwidth sharing strategies. PropShare [10] is a different peer bandwidth allocation strategy that better resists exploitative strategies, and improves the performance of swarms.

2.3 Version Control Systems - Git

Version Control Systems provide facilities to model files changing over time and distribute different versions efficiently. The popular version control system Git provides a powerful Merkle DAG² object model that captures changes to a file system tree in a distributed-friendly way:

1. Immutable objects represent Files (`blob`), Directories (`tree`), and Changes (`commit`).
2. Objects are content-addressed, by the cryptographic hash of their contents.
3. Links to other objects are embedded, forming a Merkle DAG. This provides many useful integrity and workflow properties.
4. Most versioning metadata (branches, tags, etc.) are simply pointer references, and thus inexpensive to create and update.
5. Version changes only update references or add objects.
6. Distributing version changes to other users is simply transferring objects and updating remote references.

2.4 Self-Certified File Systems - SFS

SFS [14, 13] proposed compelling implementations of both (a) distributed trust chains, and (b) egalitarian shared global namespaces. SFS introduced a technique for building *Self-Certified File Systems*, addressing remote file systems using the following scheme:

`/sfs/<Location>:<HostID>`

where `Location` is the server network address, and:

`HostID = hash(public_key || Location)`

Thus the *name* of an SFS file system certifies its server. The user can verify the public key offered by the server, negotiate a shared secret, and secure all traffic. All SFS instances share a global namespace where name allocation is cryptographic, not gated by any centralized body.

²Merkle Directed Acyclic Graph – similar but more general construction than a Merkle Tree. Deduplicated, does not need to be balanced, and non-leaf nodes contain data.

3. IPFS DESIGN

IPFS is a distributed file system which synthesizes successful ideas from previous peer-to-peer systems, including DHTs, BitTorrent, Git, and SFS. The contribution of IPFS is simplifying, evolving, and connecting proven techniques into a single cohesive system, greater than the sum of its parts. IPFS presents a new platform for writing and deploying applications, and a new system for distributing and versioning large data. IPFS could even evolve the Web itself.

IPFS is peer-to-peer; no nodes are privileged. IPFS nodes store IPFS objects in local storage. Nodes connect to each other and transfer objects. These objects represent files and other data structures. The IPFS Protocol is divided into a stack of sub-protocols responsible for different functionality:

1. **Identities** - manage node identity generation and verification. Described in Section 3.1.
2. **Network** - manages connections to other peers, uses various underlying network protocols. Configurable. Described in Section 3.2.
3. **Routing** - maintains information to locate specific peers and objects. Responds to both local and remote queries. Defaults to a DHT, but is swappable. Described in Section 3.3.
4. **Exchange** - a novel block exchange protocol (BitSwap) that governs efficient block distribution. Modelled as a market, weakly incentivizes data replication. Trade strategies are swappable. Described in Section 3.4.
5. **Objects** - a Merkle DAG of content-addressed immutable objects with links. Used to represent arbitrary data structures, e.g. file hierarchies and communication systems. Described in Section 3.5.
6. **Files** - versioned file system hierarchy inspired by Git. Described in Section 3.6.
7. **Naming** - A self-certifying mutable name system. Described in Section 3.7.

These subsystems are not independent; they are integrated and leverage blended properties. However, it is useful to describe them separately, building the protocol stack from the bottom up.

Notation: data structures and functions below are specified in Go syntax.

3.1 Identities

Nodes are identified by a `NodeId`, the cryptographic hash³ of a public key, created with S/Kademlia's static crypto puzzle [1]. Nodes store their public and private keys (encrypted with a passphrase). Users are free to instantiate a "new" node identity on every launch, though that loses accrued network benefits. Nodes are incentivized to remain the same.

```
type NodeId Multihash
type Multihash []byte
// self-describing cryptographic hash digest

type PublicKey []byte
```

³Throughout this document, *hash* and *checksum* refer specifically to cryptographic hashes of data.

```
type PrivateKey []byte
// self-describing keys
```

```
type Node struct {
    NodeId NodeId
    PubKey PublicKey
    PrivKey PrivateKey
}
```

S/Kademlia-based IPFS identity generation:

```
difficulty = <integer parameter>
n = Node{}
do {
    n.PubKey, n.PrivKey = PKI.genKeyPair()
    n.NodeId = hash(n.PubKey)
    p = count_preceding_zero_bits(hash(n.NodeId))
} while (p < difficulty)
```

Upon first connecting, peers exchange public keys, and check that `hash(other.PublicKey)` equals `other.NodeId`. If not, the connection is terminated.

Note on Cryptographic Functions.

Rather than locking the system to a particular set of function choices, IPFS favors self-describing values. Hash digest values are stored in `multihash` format, which includes a short header specifying the hash function used, and the digest length in bytes:

```
<function code><digest length><digest bytes>
```

This allows the system to (a) choose the best function for the use case (e.g. stronger security vs faster performance), and (b) evolve as function choices change. Self-describing values allow using different parameter choices compatibly.

3.2 Network

IPFS nodes communicate regularly with hundreds of other nodes in the network, potentially across the wide Internet. The IPFS network stack features:

- **Transport:** IPFS can use any transport protocol, and is best suited for WebRTC data channels [16] (for browser connectivity) or uTP (LEDBAT [23]).
- **Reliability:** IPFS can provide reliability if underlying networks do not provide it, using uTP (LEDBAT [23]) or SCTP [24].
- **Connectivity:** IPFS also uses the ICE NAT traversal techniques [22].
- **Integrity:** IPFS optionally checks integrity of messages using a hash checksum.
- **Authenticity:** IPFS optionally checks authenticity of messages by digitally signing them with the sender's private key.

Note on Peer Addressing.

IPFS can use any network; it does not rely on or assume access to IP. This allows IPFS to be used in overlay networks. IPFS stores addresses as `multiaddr`-formatted byte strings for the underlying network to use. `multiaddr` provides a way to express addresses and their protocols, including support for encapsulation. For example:

```
# an SCTP/IPv4 connection
/ip4/10.20.30.40/sctp/1234/

# an SCTP/IPv4 connection proxied over TCP/IPv4
/ip4/5.6.7.8/tcp/5678/ip4/1.2.3.4/sctp/1234/
```

3.3 Routing

IPFS nodes require a routing system that can find (a) other peers' network addresses and (b) peers who can serve particular objects. IPFS achieves this using a DSHT based on S/Kademlia and Coral, using the properties discussed in Section 2.1. The size of objects and use patterns of IPFS are similar to Coral [7] and Mainline [27], so the IPFS DHT makes a distinction for values stored based on their size. Small values (equal to or less than 1 KiB) are stored directly on the DHT. For larger values, the DHT stores references, which are the `NodeIds` of peers who can serve the block.

The interface of this DSHT is the following:

```
type IPFSRouting interface {

    FindPeer(node NodeId)
    // gets a particular peer's network address

    SetValue(key []bytes, value []bytes)
    // stores a small metadata value in DHT

    GetValue(key []bytes)
    // retrieves small metadata value from DHT

    ProvideValue(key Multihash)
    // announces this node can serve a large value

    FindValuePeers(key Multihash, min int)
    // gets a number of peers serving a large value
}
```

Note that different use cases will call for substantially different routing systems (e.g. DHT in wide networks, static HT in local networks). Thus the IPFS routing system can be swapped for one that fits users' needs. As long as the interface above is met, the rest of the system will continue to function.

3.4 Block Exchange - BitSwap Protocol

In IPFS, data distribution happens by exchanging blocks with peers using a BitTorrent inspired protocol: BitSwap. Like BitTorrent, BitSwap peers are looking to acquire a set of blocks (`want_list`), and have another set of blocks to offer in exchange (`have_list`). Unlike BitTorrent, BitSwap is not limited to the blocks in one torrent. BitSwap operates as a persistent marketplace where nodes can acquire the blocks they need, regardless of what files those blocks are part of. The blocks could come from completely unrelated files in the file system. Nodes come together to barter in the marketplace.

While the notion of a barter system implies a virtual currency could be created, this would require a global ledger to track ownership and transfer of the currency. This can be implemented as a BitSwap strategy, and will be explored in a future paper.

In the base case, BitSwap nodes have to provide direct value to each other in the form of blocks. This works fine

when the distribution of blocks across nodes is complementary, meaning they have what the other wants. Often, this will not be the case. In some cases, nodes must *work* for their blocks. In the case that a node has nothing that its peers want (or nothing at all), it seeks the pieces its peers want, with lower priority than what the node wants itself. This incentivizes nodes to cache and disseminate rare pieces, even if they are not interested in them directly.

3.4.1 BitSwap Credit

The protocol must also incentivize nodes to seed when they do not need anything in particular, as they might have the blocks that others want. Thus, BitSwap nodes send blocks to their peers optimistically, expecting the debt to be repaid. But leeches (free-loading nodes that never share) must be protected against. A simple credit-like system solves the problem:

1. Peers track their balance (in bytes verified) with other nodes.
2. Peers send blocks to debtor peers probabilistically, according to a function that falls as debt increases.

Note that if a node decides not to send to a peer, the node subsequently ignores the peer for an `ignore_cooldown` timeout (the default for BitSwap is 10 seconds). This prevents senders from trying to game the probability by just causing more dice-rolls.

3.4.2 BitSwap Strategy

The differing strategies that BitSwap peers might employ have wildly different effects on the performance of the exchange as a whole. In BitTorrent, while a standard strategy is specified (tit-for-tat), a variety of others have been implemented, ranging from BitTyrant [10] (sharing the least-possible), to BitThief [10] (exploiting a vulnerability and never share), to PropShare [10] (sharing proportionally). A range of strategies (good and malicious) could similarly be implemented by BitSwap peers. The choice of function, then, should aim to:

1. maximize the trade performance for the node, and the whole exchange
2. prevent freeloaders from exploiting and degrading the exchange
3. be effective with and resistant to other unknown strategies
4. be lenient to trusted peers

The exploration of the space of such strategies is future work. One choice of function that works in practice is a sigmoid, scaled by a *debt ratio*:

Let the *debt ratio* r between a node and its peer be:

$$r = \frac{\text{bytes_sent}}{\text{bytes_recv} + 1}$$

Given r , let the probability of sending to a debtor be:

$$P(\text{send} | r) = 1 - \frac{1}{1 + \exp(6 - 3r)}$$

As you can see in Figure 1, this function drops off quickly as the nodes' *debt ratio* surpasses twice the established credit.

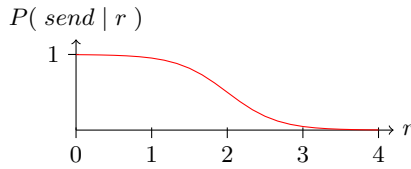


Figure 1: Probability of Sending as r increases

The *debt ratio* is a measure of trust: lenient to debts between nodes that have previously exchanged lots of data successfully, and merciless to unknown, untrusted nodes. This (a) provides resistance to attackers who would create lots of new nodes (Sybil attacks), (b) protects previously successful trade relationships, even if one of the nodes is temporarily unable to provide value, and (c) eventually chokes relationships that have deteriorated until they improve.

3.4.3 BitSwap Ledger

BitSwap nodes keep ledgers accounting the transfers with other nodes. This allows nodes to keep track of history and avoid tampering. When activating a connection, BitSwap nodes exchange their ledger information. If it does not match exactly, the ledger is reinitialized from scratch, losing the accrued credit or debt. It is possible for malicious nodes to purposefully “lose” the ledger, hoping to erase debts. It is unlikely that nodes will have accrued enough debt to warrant also losing the accrued trust; however the partner node is free to count it as misconduct, and refuse to trade.

```
type Ledger struct {
  owner      NodeId
  partner    NodeId
  bytes_sent int
  bytes_rcv  int
  timestamp  Timestamp
}
```

Nodes are free to keep the ledger history, though it is not necessary for correct operation. Only the current ledger entries are useful. Nodes are also free to garbage collect ledgers as necessary, starting with the less useful ledgers: the old (peers may not exist anymore) and small.

3.4.4 BitSwap Specification

BitSwap nodes follow a simple protocol:

```
// Additional state kept
type BitSwap struct {
  ledgers map[NodeId]Ledger
  // Ledgers known to this node, inc inactive

  active map[NodeId]Peer
  // currently open connections to other nodes

  want_list []Multihash
  // checksums of blocks this node wants

  have_list []Multihash
  // checksums of blocks this node has
}

type Peer struct {
```

```
  nodeid NodeId
  ledger Ledger
  // Ledger between the node and this peer

  last_seen Timestamp
  // timestamp of last received message

  want_list []Multihash
  // checksums of all blocks wanted by peer
  // includes blocks wanted by peer's peers
}

// Protocol interface:
interface Peer {
  open (nodeid :NodeId, ledger :Ledger);
  send_want_list (want_list :WantList);
  send_block (block :Block) -> (complete :Bool);
  close (final :Bool);
}
```

Sketch of the lifetime of a peer connection:

1. Open: peers send ledgers until they agree.
2. Sending: peers exchange want lists and blocks.
3. Close: peers deactivate a connection.
4. Ignored: (special) a peer is ignored (for the duration of a timeout) if a node's strategy avoids sending.

Peer.open(NodeId, Ledger).

When connecting, a node initializes a connection with a **Ledger**, either stored from a connection in the past or a new one zeroed out. Then, sends an **open** message with the **Ledger** to the peer.

Upon receiving an **open** message, a peer chooses whether to activate the connection. If – according to the receiver's **Ledger** – the sender is not a trusted agent (transmission below zero, or large outstanding debt) the receiver may opt to ignore the request. This should be done probabilistically with an **ignore_cooldown** timeout, as to allow errors to be corrected and attackers to be thwarted.

If activating the connection, the receiver initializes a **Peer** object with the local version of the **Ledger** and sets the **last_seen** timestamp. Then, it compares the received **Ledger** with its own. If they match exactly, the connections have opened. If they do not match, the peer creates a new zeroed out **Ledger** and sends it.

Peer.send_want_list(WantList).

While the connection is open, nodes advertise their **want_list** to all connected peers. This is done (a) upon opening the connection, (b) after a randomized periodic timeout, (c) after a change in the **want_list** and (d) after receiving a new block.

Upon receiving a **want_list**, a node stores it. Then, it checks whether it has any of the wanted blocks. If so, it sends them according to the *BitSwap Strategy* above.

Peer.send_block(Block).

Sending a block is straightforward. The node simply transmits the block of data. Upon receiving all the data, the receiver computes the Multihash checksum to verify it matches the expected one, and returns confirmation.

Upon finalizing the correct transmission of a block, the receiver moves the block from `want_list` to `have_list`, and both the receiver and sender update their ledgers to reflect the additional bytes transmitted.

If a transmission verification fails, the sender is either malfunctioning or attacking the receiver. The receiver is free to refuse further trades. Note that BitSwap expects to operate on a reliable transmission channel, so transmission errors – which could lead to incorrect penalization of an honest sender – are expected to be caught before the data is given to BitSwap.

Peer.close(Bool).

The `final` parameter to `close` signals whether the intention to tear down the connection is the sender's or not. If false, the receiver may opt to re-open the connection immediately. This avoids premature closes.

A peer connection should be closed under two conditions:

- a `silence_wait` timeout has expired without receiving any messages from the peer (default BitSwap uses 30 seconds). The node issues `Peer.close(false)`.
- the node is exiting and BitSwap is being shut down. In this case, the node issues `Peer.close(true)`.

After a `close` message, both receiver and sender tear down the connection, clearing any state stored. The `Ledger` may be stored for the future, if it is useful to do so.

Notes.

- Non-open messages on an inactive connection should be ignored. In case of a `send_block` message, the receiver may check the block to see if it is needed and correct, and if so, use it. Regardless, all such out-of-order messages trigger a `close(false)` message from the receiver to force re-initialization of the connection.

3.5 Object Merkle DAG

The DHT and BitSwap allow IPFS to form a massive peer-to-peer system for storing and distributing blocks quickly and robustly. On top of these, IPFS builds a Merkle DAG, a directed acyclic graph where links between objects are cryptographic hashes of the targets embedded in the sources. This is a generalization of the Git data structure. Merkle DAGs provide IPFS many useful properties, including:

1. **Content addressing:** all content is uniquely identified by its multihash checksum, **including links**.
2. **Tamper resistance:** all content is verified with its checksum. If data is tampered with or corrupted, IPFS detects it.
3. **Deduplication:** all objects that hold the exact same content are equal, and only stored once. This is particularly useful with index objects, such as git **trees** and **commits**, or common portions of data.

The IPFS Object format is:

```
type IPFSLink struct {
    Name string
    // name or alias of this link
```

```
    Hash Multihash
    // cryptographic hash of target

    Size int
    // total size of target
}

type IPFSObject struct {
    links []IPFSLink
    // array of links

    data []byte
    // opaque content data
}
```

The IPFS Merkle DAG is an extremely flexible way to store data. The only requirements are that object references be (a) content-addressed, and (b) encoded in the format above. IPFS grants applications complete control over the data field; applications can use any custom data format they chose, which IPFS may not understand. The separate in-object link table allows IPFS to:

- List all object references in an object. For example:

```
> ipfs ls /XLZ1625Jjn7SubMDgEyeaynFuR84ginqvzb
XLKkgq61DYaQ8NhkcqyU7rLcnSa7dSHQ16x 189458 less
XLHBNmRQ5sJJrdMPuu48pzeyTtRo39tNDR5 19441 script
XLf4hwVHsVuZ78FZK6fozf8Jj9WEURMbCX4 5286 template

<object multihash> <object size> <link name>
```

- Resolve string path lookups, such as `foo/bar/baz`. Given an object, IPFS resolves the first path component to a hash in the object's link table, fetches that second object, and repeats with the next component. Thus, string paths can walk the Merkle DAG no matter what the data formats are.
- Resolve all objects referenced recursively:

```
> ipfs refs --recursive \
    /XLZ1625Jjn7SubMDgEyeaynFuR84ginqvzb
XLLxhdgJcXzLbtsLRL1twCHA2NrURp4H38s
XLKkgq61DYaQ8NhkcqyU7rLcnSa7dSHQ16x
XLHBNmRQ5sJJrdMPuu48pzeyTtRo39tNDR5
XLWVQDqxo9Km9zLyquoC9gAP8CL1gWnHZ7z
...
```

A raw data field and a common link structure are the necessary components for constructing arbitrary data structures on top of IPFS. While it is easy to see how the Git object model fits on top of this DAG, consider these other potential data structures: (a) key-value stores (b) traditional relational databases (c) linked data triple stores (d) linked document publishing systems (e) linked communications platforms (f) cryptocurrency blockchains. These can all be modeled on top of the IPFS Merkle DAG, which allows any of these systems to use IPFS as a transport protocol for more complex applications.

3.5.1 Paths

IPFS objects can be traversed with a string path API. Paths work as they do in traditional UNIX file systems and the Web. The Merkle DAG links make traversing it easy. Note that full paths in IPFS are of the form:

```
# format
/ipfs/<hash-of-object>/<name-path-to-object>

# example
/ipfs/XLYkgq61DYaQ8NhkcqyU7rLcnSa7dSHQ16x/foo.txt
```

The `/ipfs` prefix allows mounting into existing systems at a standard mount point without conflict (mount point names are of course configurable). The second path component (first within IPFS) is the hash of an object. This is always the case, as there is no global root. A root object would have the impossible task of handling consistency of millions of objects in a distributed (and possibly disconnected) environment. Instead, we simulate the root with content addressing. All objects are always accessible via their hash. Note that this means that given three objects in path `<foo>/bar/baz`, the last object is accessible by all:

```
/ipfs/<hash-of-foo>/bar/baz
/ipfs/<hash-of-bar>/baz
/ipfs/<hash-of-baz>
```

3.5.2 Local Objects

IPFS clients require some *local storage*, an external system on which to store and retrieve local raw data for the objects IPFS manages. The type of storage depends on the node's use case. In most cases, this is simply a portion of disk space (either managed by the native file system, by a key-value store such as `leveldb` [6], or directly by the IPFS client). In others, for example non-persistent caches, this storage is just a portion of RAM.

Ultimately, all blocks available in IPFS are in some node's *local storage*. When users request objects, they are found, downloaded, and stored locally, at least temporarily. This provides fast lookup for some configurable amount of time thereafter.

3.5.3 Object Pinning

Nodes who wish to ensure the survival of particular objects can do so by **pinning** the objects. This ensures the objects are kept in the node's *local storage*. Pinning can be done recursively, to pin down all linked descendent objects as well. All objects pointed to are then stored locally. This is particularly useful to persist files, including references. This also makes IPFS a Web where links are *permanent*, and Objects can ensure the survival of others they point to.

3.5.4 Publishing Objects

IPFS is globally distributed. It is designed to allow the files of millions of users to coexist together. The DHT, with content-hash addressing, allows publishing objects in a fair, secure, and entirely distributed way. Anyone can publish an object by simply adding its key to the DHT, adding themselves as a peer, and giving other users the object's path. Note that Objects are essentially immutable, just like in Git. New versions hash differently, and thus are new objects. Tracking versions is the job of additional versioning objects.

3.5.5 Object-level Cryptography

IPFS is equipped to handle object-level cryptographic operations. An encrypted or signed object is wrapped in a special frame that allows encryption or verification of the raw bytes.

```
type EncryptedObject struct {
    Object []bytes
    // raw object data encrypted

    Tag []bytes
    // optional tag for encryption groups
}

type SignedObject struct {
    Object []bytes
    // raw object data signed

    Signature []bytes
    // hmac signature

    PublicKey []multihash
    // multihash identifying key
}
```

Cryptographic operations change the object's hash, defining a different object. IPFS automatically verifies signatures, and can decrypt data with user-specified keychains. Links of encrypted objects are protected as well, making traversal impossible without a decryption key. It is possible to have a parent object encrypted under one key, and a child under another or not at all. This secures links to shared objects.

3.6 Files

IPFS also defines a set of objects for modeling a versioned file system on top of the Merkle DAG. This object model is similar to Git's:

1. **block**: a variable-size block of data.
2. **list**: a collection of blocks or other lists.
3. **tree**: a collection of blocks, lists, or other trees.
4. **commit**: a snapshot in the version history of a tree.

I hoped to use the Git object formats exactly, but had to depart to introduce certain features useful in a distributed file system, namely (a) fast size lookups (aggregate byte sizes have been added to objects), (b) large file deduplication (adding a **list** object), and (c) embedding of **commits** into **trees**. However, IPFS File objects are close enough to Git that conversion between the two is possible. Also, a set of Git objects can be introduced to convert without losing any information (UNIX file permissions, etc).

Notation: File object formats below use JSON. Note that this structure is actually binary encoded using Protocol Buffers, though IPFS includes import/export to JSON.

3.6.1 File Object: blob

The **blob** object contains an addressable unit of data, and represents a file. IPFS Blobs are like Git blobs or file system data blocks. They store the users' data. Note that IPFS files can be represented by both **lists** and **blobs**. Blobs have no links.

```

{
  "data": "some data here"
  // blobs have no links
}

```

3.6.2 File Object: list

The **list** object represents a large or deduplicated file made up of several IPFS **blobs** concatenated together. **lists** contain an ordered sequence of **blob** or **list** objects. In a sense, the IPFS **list** functions like a file system file with indirect blocks. Since **lists** can contain other **lists**, topologies including linked lists and balanced trees are possible. Directed graphs where the same node appears in multiple places allow in-file deduplication. Of course, cycles are not possible, as enforced by hash addressing.

```

{
  "data": ["blob", "list", "blob"],
  // lists have an array of object types as data
  "links": [
    { "hash": "XLYkgq61DYaQ8NhkcqyU7rLcnSa7dSHQ16x",
      "size": 189458 },
    { "hash": "XLHBNmRQ5sJJrdMPuu48pzeyTtRo39tNDR5",
      "size": 19441 },
    { "hash": "XLWVQDqxo9Km9zLyquoC9gAP8CL1gWnHZ7z",
      "size": 5286 }
  ]
  // lists have no names in links
}

```

3.6.3 File Object: tree

The **tree** object in IPFS is similar to Git's: it represents a directory, a map of names to hashes. The hashes reference blobs, lists, other trees, or commits. Note that traditional path naming is already implemented by the Merkle DAG.

```

{
  "data": nil,
  // trees have no data, only links
  "links": [
    { "hash": "XLYkgq61DYaQ8NhkcqyU7rLcnSa7dSHQ16x",
      "name": "less", "size": 189458 },
    { "hash": "XLHBNmRQ5sJJrdMPuu48pzeyTtRo39tNDR5",
      "name": "script", "size": 19441 },
    { "hash": "XLWVQDqxo9Km9zLyquoC9gAP8CL1gWnHZ7z",
      "name": "template", "size": 5286 }
  ]
  // trees do have names
}

```

3.6.4 File Object: commit

The **commit** object in IPFS represents a snapshot in the version history of any object. It is similar to Git's, but can reference any type of object. It also links to author objects.

```

{
  "data": {
    "type": "tree",
    "date": "2014-09-20 12:44:06Z",
    "message": "This is a commit message."
  },
  "links": [

```

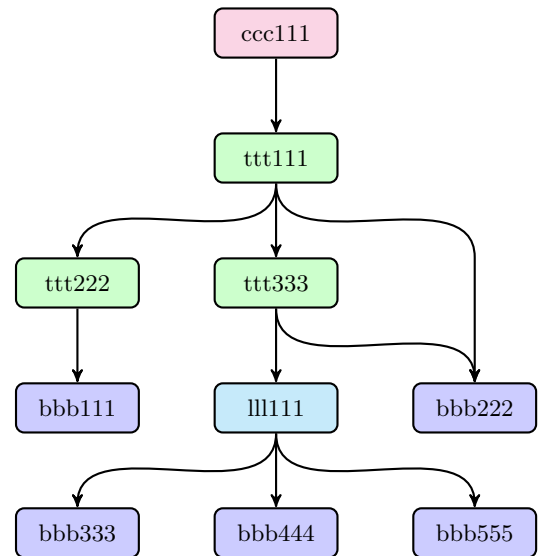


Figure 2: Sample Object Graph

```

ipfs file-cat <ccc111-hash> --json
{
  "data": {
    "type": "tree",
    "date": "2014-09-20 12:44:06Z",
    "message": "This is a commit message."
  },
  "links": [
    { "hash": "<ccc000-hash>",
      "name": "parent", "size": 25309 },
    { "hash": "<ttt111-hash>",
      "name": "object", "size": 5198 },
    { "hash": "<aaa111-hash>",
      "name": "author", "size": 109 }
  ]
}

ipfs file-cat <ttt111-hash> --json
{
  "data": nil,
  "links": [
    { "hash": "<ttt222-hash>",
      "name": "ttt222-name", "size": 1234 },
    { "hash": "<ttt333-hash>",
      "name": "ttt333-name", "size": 3456 },
    { "hash": "<bbb222-hash>",
      "name": "bbb222-name", "size": 22 }
  ]
}

ipfs file-cat <bbb222-hash> --json
{
  "data": "blob222 data",
  "links": []
}

```

Figure 3: Sample Objects


```

{ "hash": "XLa1qMBKiSEEDhojb9FFZ4tEvLf7FEQdhdU",
  "name": "parent", "size": 25309 },
{ "hash": "XLGw74KAy9junbh28x7ccWov9inu1Vo7pnX",
  "name": "object", "size": 5198 },
{ "hash": "XLF2ipQ4jD3UdeX5xp1KBgeHRhemUtaA8Vm",
  "name": "author", "size": 109 }
]
}

```

3.6.5 Version control

The `commit` object represents a particular snapshot in the version history of an object. Comparing the objects (and children) of two different commits reveals the differences between two versions of the file system. As long as a single `commit` and all the children objects it references are accessible, all preceding versions are retrievable and the full history of the file system changes can be accessed. This falls out of the Merkle DAG object model.

The full power of the Git version control tools is available to IPFS users. The object model is compatible, though not the same. It is possible to (a) build a version of the Git tools modified to use the IPFS object graph, (b) build a mounted FUSE file system that mounts an IPFS `tree` as a Git repo, translating Git file system read/writes to the IPFS formats.

3.6.6 File System Paths

As we saw in the Merkle DAG section, IPFS objects can be traversed with a string path API. The IPFS File Objects are designed to make mounting IPFS onto a UNIX file system simpler. They restrict `trees` to have no data, in order to represent them as directories. The `commits` can either be represented as directories or hidden from the file system entirely.

3.6.7 Splitting Files into Lists and Blob

One of the main challenges with versioning and distributing large files is finding the right way to split them into independent blocks. Rather than assume it can make the right decision for every type of file, IPFS offers the following alternatives:

- (a) Use Rabin fingerprints [21] as in LBFS [15] to pick suitable block boundaries.
- (b) Use the rsync [26] rolling-checksum algorithm, to detect blocks that have changed between versions.
- (c) Allow users to specify block-splitting functions highly tuned for specific files.

3.6.8 Path Lookup Performance

Path-based access traverses the object graph. Retrieving each object requires looking up its key in the DHT, connecting to peers, and retrieving its blocks. This is considerable overhead, particularly when looking up paths with many components. This is mitigated by:

- **tree caching**: since all objects are hash-addressed, they can be cached indefinitely. Additionally, `trees` tend to be small in size so IPFS prioritizes caching them over `blobs`.
- **flattened trees**: for any given `tree`, a special `flattened tree` can be constructed to list all objects reachable

from the `tree`. Names in the `flattened tree` would really be paths parting from the original tree, with slashes.

For example, the `flattened tree` for `ttt111` in Figure 2:

```

{
  "data":
    ["tree", "blob", "tree", "list", "blob" "blob"],
  "links": [
    { "hash": "<ttt222-hash>", "size": 1234
      "name": "ttt222-name" },
    { "hash": "<bbb111-hash>", "size": 123,
      "name": "ttt222-name/bbb111-name" },
    { "hash": "<ttt333-hash>", "size": 3456,
      "name": "ttt333-name" },
    { "hash": "<lll111-hash>", "size": 587,
      "name": "ttt333-name/lll111-name"},
    { "hash": "<bbb222-hash>", "size": 22,
      "name": "ttt333-name/lll111-name/bbb222-name" },
    { "hash": "<bbb222-hash>", "size": 22
      "name": "bbb222-name" }
  ]
}

```

3.7 IPNS: Naming and Mutable State

So far, the IPFS stack forms a peer-to-peer block exchange constructing a content-addressed DAG of objects. It serves to publish and retrieve immutable objects. It can even track the version history of these objects. However, there is a critical component missing: mutable naming. Without it, all communication of new content must happen off-band, sending IPFS links. What is required is some way to retrieve mutable state at *the same path*.

It is worth stating why – if mutable data is necessary in the end – we worked hard to build up an *immutable* Merkle DAG. Consider the properties of IPFS that fall out of the Merkle DAG: objects can be (a) retrieved via their hash, (b) integrity checked, (c) linked to others, and (d) cached indefinitely. In a sense:

Objects are **permanent**.

These are the critical properties of a high-performance distributed system, where data is expensive to move across network links. Object content addressing constructs a Web with (a) significant bandwidth optimizations, (b) untrusted content serving, (c) permanent links, and (d) the ability to make full permanent backups of any object and its references.

The Merkle DAG (immutable content-addressed objects) and Naming (mutable pointers to the Merkle DAG) instantiate a dichotomy present in many successful distributed systems. These include the Git Version Control System, with its immutable objects and mutable references; and Plan9 [18], the distributed successor to UNIX, with its mutable Fossil [20] and immutable Venti [19] file systems. LBFS [15] also uses mutable indices and immutable chunks.

3.7.1 Self-Certified Names

Using the naming scheme from SFS [14, 13] gives us a way to construct self-certified names (in a cryptographically assigned global namespace) that are mutable. The IPFS scheme is as follows.

1. Recall that in IPFS:

```
NodeId = hash(node.PubKey)
```

2. We assign every user a mutable namespace at:

```
/ipns/<NodeId>
```

3. A user can publish an Object to this path **Signed** by her private key, say at:

```
/ipns/XLF2ipQ4jD3UdeX5xp1KBgeHRhemUtaA8Vm/
```

4. When other users retrieve the object, they can check the signature matches the public key and NodeId. This verifies the authenticity of the Object published by the user, achieving mutable state retrieval.

Note the following details:

- The **ipns** (InterPlanetary Name Space) separate prefix is to establish an easily recognizable distinction between *mutable* and *immutable* paths, for both programs and human readers.
- Because this is *not* a content-addressed object, publishing it relies on the only mutable state distribution system in IPFS, the Routing system. The process is (1) publish the object as a regular immutable IPFS object, (2) publish its hash on the Routing system as a metadata value:

```
routing.setValue(NodeId, <ns-object-hash>)
```

- Any links in the Object published act as sub-names in the namespace:

```
/ipns/XLF2ipQ4jD3UdeX5xp1KBgeHRhemUtaA8Vm/  
/ipns/XLF2ipQ4jD3UdeX5xp1KBgeHRhemUtaA8Vm/docs  
/ipns/XLF2ipQ4jD3UdeX5xp1KBgeHRhemUtaA8Vm/docs/ipfs
```

- It is advised to publish a **commit** object, or some other object with a version history, so that clients may be able to find old names. This is left as a user option, as it is not always desired.

3.7.2 Human Friendly Names

While IPNS is indeed a way of assigning and reassigning names, it is not very user friendly, as it exposes long hash values as names, which are notoriously hard to remember. These work for URLs, but not for many kinds of offline transmission. Thus, IPFS increases the user-friendliness of IPNS with the following techniques.

Peer Links.

As encouraged by SFS, users can link other users' Objects directly into their own Objects (namespace, home, etc). This has the benefit of also creating a web of trust (and supports the old Certificate Authority model):

```
# Alice links to bob Bob  
ipfs link /<alice-pk-hash>/friends/bob /<bob-pk-hash>  
  
# Eve links to Alice  
ipfs link /<eve-pk-hash>/friends/alice /<alice-pk-hash>
```

```
# Eve also has access to Bob  
/<eve-pk-hash>/friends/alice/friends/bob
```

```
# access Verisign certified domains  
/<verisign-pk-hash>/foo.com
```

DNS TXT IPNS Records.

If `/ipns/<domain>` is a valid domain name, IPFS looks up key **ipns** in its DNS TXT records. IPFS interprets the value as either an object hash or another IPNS path:

```
# this DNS TXT record  
ipfs.benet.ai. TXT "ipfs=XLF2ipQ4jD3U ..."  
  
# behaves as symlink  
ln -s /ipns/XLF2ipQ4jD3U /ipns/fs.benet.ai
```

Proquint Pronounceable Identifiers.

There have always been schemes to encode binary into pronounceable words. IPNS supports Proquint [28]. Thus:

```
# this proquint phrase  
/ipns/dahih-dolij-sozuk-vosah-luvar-fuluh  
  
# will resolve to corresponding  
/ipns/KhAwNprxYVxKqpDZ
```

Name Shortening Services.

Services are bound to spring up that will provide name shortening as a service, offering up their namespaces to users. This is similar to what we see today with DNS and Web URLs:

```
# User can get a link from  
/ipns/shorten.er/foobar  
  
# To her own namespace  
/ipns/XLF2ipQ4jD3UdeX5xp1KBgeHRhemUtaA8Vm
```

3.8 Using IPFS

IPFS is designed to be used in a number of different ways. Here are just some of the usecases I will be pursuing:

1. As a mounted global file system, under `/ipfs` and `/ipns`.
2. As a mounted personal sync folder that automatically versions, publishes, and backs up any writes.
3. As an encrypted file or data sharing system.
4. As a versioned package manager for *all* software.
5. As the root file system of a Virtual Machine.
6. As the boot file system of a VM (under a hypervisor).
7. As a database: applications can write directly to the Merkle DAG data model and get all the versioning, caching, and distribution IPFS provides.
8. As a linked (and encrypted) communications platform.

9. As an integrity checked CDN for large files (without SSL).
10. As an encrypted CDN.
11. On webpages, as a Web CDN.
12. As a new Permanent Web where links do not die.

The IPFS implementations target:

- (a) an IPFS library to import in your own applications.
- (b) command line tools to manipulate objects directly.
- (c) mounted file systems, using FUSE [25] or as kernel modules.

4. THE FUTURE

The ideas behind IPFS are the product of decades of successful distributed systems research in academia and open source. IPFS synthesizes many of the best ideas from the most successful systems to date. Aside from BitSwap, which is a novel protocol, the main contribution of IPFS is this coupling of systems and synthesis of designs.

IPFS is an ambitious vision of new decentralized Internet infrastructure, upon which many different kinds of applications can be built. At the bare minimum, it can be used as a global, mounted, versioned file system and namespace, or as the next generation file sharing system. At its best, it could push the Web to new horizons, where publishing valuable information does not impose hosting it on the publisher but upon those interested, where users can trust the content they receive without trusting the peers they receive it from, and where old but important files do not go missing. IPFS looks forward to bringing us toward the Permanent Web.

5. ACKNOWLEDGMENTS

IPFS is the synthesis of many great ideas and systems. It would be impossible to dare such ambitious goals without standing on the shoulders of such giants. Personal thanks to David Dalrymple, Joe Zimmerman, and Ali Yahya for long discussions on many of these ideas, in particular: exposing the general Merkle DAG (David, Joe), rolling hash blocking (David), and S/Kademlia Sybil protection (David, Ali). And special thanks to David Mazieres, for his ever brilliant ideas.

6. REFERENCES TODO

References

- [1] I. Baumgart and S. Mies. S/Kademlia: a practicable approach towards secure key-based routing. In *Parallel and distributed systems, 2007 international conference on*. Vol. 2. IEEE, 2007, pp. 1–8.
- [2] I. BitTorrent. BitTorrent and μ Torrent software surpass 150 million user milestone. Jan. 2012. URL: http://www.bittorrent.com/intl/es/company/about/ces_2012_150m_users.
- [3] Camlistore. URL: <http://camlistore.org/>.
- [4] B. Cohen. Incentives build robustness in BitTorrent. In *Workshop on economics of peer-to-peer systems*. Vol. 6, 2003, pp. 68–72.
- [5] F. Dabek, M. F. Kaashoek, D. Karger, R. Morris, and I. Stoica. Wide-area cooperative storage with CFS. In *Proceedings of the eighteenth ACM symposium on operating systems principles*. ACM, 2001, pp. 202–215.
- [6] J. Dean and S. Ghemawat. LevelDB—a fast and lightweight key/value database library by Google. 2011.
- [7] M. J. Freedman, E. Freudenthal, and D. Mazieres. Democratizing content publication with Coral. In *NSDI*. Vol. 4, 2004, pp. 18–18.
- [8] J. H. Howard, M. L. Kazar, S. G. Menees, D. A. Nichols, M. Satyanarayanan, R. N. Sidebotham, and M. J. West. Scale and performance in a distributed file system. *ACM transactions on computer systems (TOCS)*, 6(1):51–81, 1988.
- [9] J. Kubiawicz, D. Bindel, Y. Chen, S. Czerwinski, P. Eaton, D. Geels, R. Gummadi, S. Rhea, H. Weather- spoon, W. Weimer, et al. Oceanstore: an architecture for global-scale persistent storage. *ACM sigplan notices*, 35(11):190–201, 2000.
- [10] D. Levin, K. LaCurts, N. Spring, and B. Bhattacharjee. BitTorrent is an auction: analyzing and improving BitTorrent’s incentives. In *ACM SIGCOMM computer communication review*. Vol. 38. (4). ACM, 2008, pp. 243–254.
- [11] A. J. Mashtizadeh, A. Bittau, Y. F. Huang, and D. Mazieres. Replication, history, and grafting in the Ori file system. In *Proceedings of the twenty-fourth ACM symposium on operating systems principles*. ACM, 2013, pp. 151–166.
- [12] P. Maymounkov and D. Mazieres. Kademlia: a peer-to-peer information system based on the xor metric. In, *Peer-to-peer systems*, pp. 53–65. Springer, 2002.
- [13] D. Mazieres and F. Kaashoek. Self-certifying file system, 2000.
- [14] D. Mazieres and M. F. Kaashoek. Escaping the evils of centralized control with self-certifying pathnames. In *Proceedings of the 8th ACM SIGOPS european workshop on support for composing distributed applications*. ACM, 1998, pp. 118–125.
- [15] A. Muthitacharoen, B. Chen, and D. Mazieres. A low-bandwidth network file system. In *ACM SIGOPS operating systems review*. Vol. 35. (5). ACM, 2001, pp. 174–187.
- [16] A. Narayanan, C. Jennings, A. Bergkvist, and D. Burnett. WebRTC 1.0: real-time communication between browsers. W3C Working Draft. W3C, Sept. 2013. URL: <http://www.w3.org/TR/2013/WD-webrtc-20130910/>.
- [17] M. Ogden. Dat: version, fork, and sync data. URL: <https://github.com/maxogden/dat>.
- [18] R. Pike, D. Presotto, K. Thompson, H. Trickey, et al. Plan 9 from bell labs. In *Proceedings of the summer 1990 UKUUG conference*. London, UK, 1990, pp. 1–9.
- [19] S. Quinlan and S. Dorward. Venti: a new approach to archival storage. In *Fast*. Vol. 2, 2002, pp. 89–101.
- [20] S. Quinlan, J. McKie, and R. Cox. Fossil, an archival file server. 2003.
- [21] M. O. Rabin et al. *Fingerprinting by random polynomials*. Center for Research in Computing Techn., Aiken Computation Laboratory, Univ., 1981.

- [22] J. Rosenberg and A. Keranen. Interactive Connectivity Establishment (ICE): a protocol for Network Address Translator (NAT) traversal for offer/answer protocols, 2013.
- [23] S. Shalunov, G. Hazel, J. Iyengar, and M. Kuehlewind. Low extra delay background transport (LEDBAT). *Draft-ietf-ledbat-congestion-04. txt*, 2010.
- [24] R. R. Stewart and Q. Xie. *Stream control transmission protocol (SCTP): a reference guide*. Addison-Wesley Longman Publishing Co., Inc., 2001.
- [25] M. Szeredi et al. FUSE: filesystem in userspace. URL: <http://fuse.sourceforge.net/>.
- [26] A. Tridgell, P. Mackerras, et al. The rsync algorithm, 1996.
- [27] L. Wang and J. Kangasharju. Measuring large-scale distributed systems: case of BitTorrent Mainline DHT. In *Peer-to-peer computing (P2P), 2013 IEEE thirteenth international conference on*. IEEE, 2013, pp. 1–10.
- [28] D. S. Wilkerson. A proposal for proquints: identifiers that are readable, spellable, and pronounceable. *arXiv preprint arXiv:0901.4016*, 2009.