

# Decentralized Peer to Peer End to End Encrypted File Syncing and Versioning

Lokesh Raj Arora, K.M. Uma Maheswari

Received 14 November 2018 ▪ Revised: 30 November 2018 ▪ Accepted: 06 December 2018

**Abstract:** Creation of an application framework or proof of concept for the decentralized end to end encrypted real time communication between two or more hosts in the network. This gives users more control over who and what to share over the network and also give a bidirectional support built in the very protocol, rather than using some hacks using polling and long polling and preventing loss of data during transfer of data, in protocols like  $\mu$ TP.

**Index Terms:** Decentralized, Peer to Peer.

**Keywords:** Decentralization, End to End Encryption, Merkle Trees.

## INTRODUCTION

Numerous data sets are shared online today utilizing HTTP what's more, FTP, which need worked in help for variant control or substance tending to of information. This outcomes in connect decay and substance float as records are moved, refreshed or on the other hand erased, prompting a disturbing rate of vanishing information references in regions, for example, distributed logical writing. Distributed storage administrations like S3 guarantee accessibility of information, yet they have a concentrated center point and-talked net- working model and are in this manner constrained by their bandwidth, which means famous records can turn out to be very expensive to share. Administrations like Dropbox and Google Drive give rendition control and synchronization over distributed storage administrations which fixes numerous issues with broken connections however depend on restrictive code and administrations expecting clients to store their information on brought together cloud framework which has suggestions on cost, exchange speeds, merchant lock-in and client security. Dispersed record sharing instruments can turn out to be quicker as records turn out to be increasingly mainstream, expelling the data transfer capacity bottleneck and making record dispersion less expensive. They additionally use connect goals and disclosure frameworks which can counteract broken connections meaning if the first source goes disconnected other reinforcement sources can be naturally found. Anyway these document sharing devices today are not upheld by Web programs, don't have great protection ensures, and don't give a component for refreshing documents without redistributing another data set which could mean altogether re-downloading information you as of now have.

## CREATION OF A REPOSITORY

DP2PFS is a data set synchronization convention that does not accept a data set is static or that the whole data set will be downloaded.

The convention is freethinker to the fundamental transport for example you could execute DP2PFS over transporter pigeon. Information is put away in a configuration called SLEEP (Ogden and Buus 2017), portrayed in its very own paper. The key properties of the DP2PFS configuration are clarified in this area.

**3.1 Content Integrity:** Information and distributor trust- worthiness is confirmed through utilization of marked hashes of the substance.

**3.2 Decentralized Mirroring:** Users sharing the same DP2PFS automatically discover each other and exchange data in a swarm.

**3.3 Network Privacy:** DP2PFS provides certain privacy guarantees including end-to-end encryption.

**3.4 Incremental Versioning:** Data sets can be efficiently synced, even in real time, to other peers.

**3.5 Random Access:** Huge file hierarchies can be efficiently traversed remotely.

### A. Content Integrity

Content honesty or integrity implies having the capacity to confirm the information you got is precisely the same form of the information that you anticipated. This is critical in a disseminated framework as this component will get off base information sent by terrible friends. It additionally has suggestions for reproducibility as it gives you a chance to allude to a particular adaptation of a data set. Connection decay, when joins online quit settling, and substance float, when information changes however the connection to the information continues as before, are two normal issues in information examination. For instance, one day a record called data.zip might change, yet a regular HTTP connects to the record does exclude a hash of the substance, or give an approach to get refreshed metadata, so customers that just have the HTTP connect have no real way to check if the record changed without downloading the whole record once more. Alluding to a record by the hash of its substance is called substance address ability, and lets clients not just confirm that the information they get is the rendition of the information they need, yet in addition gives individuals a chance to refer to explicit renditions of the information by alluding to a particular hash.

DP2PFS uses *BLAKE2b* (Aumasson et al. 2013) cryptographically secure hashes to address content. Hashes are arranged in a Merkle tree (Mykletun, Narasimha, and Tsudik 2003), a tree where each non-leaf node is the hash of all child nodes. Leaf nodes contain pieces of the data set. Due to the properties of secure cryptographic hashes the top hash can only be produced if all data below it matches exactly. If two trees have matching top hashes then you know that all other nodes in the tree must match as well, and you can conclude that your data set is synchronized. Trees are selected as the primary data structure because of its certain properties that make it easier to resemble the file tree.

All messages in the DP2PFS convention are scrambled and marked utilizing the open key amid transport. This implies that except if you know the open key (*for example except if the DP2PFS connect was imparted to you*) at that point you won't have the capacity to find or speak with any part of the swarm for that Dat. Anybody with people in general key can check that messages, (for example, sections in a DP2PFS Stream) were made by a holder of the private key. Each archive has a relating private key which is kept in your home envelope and never shared. DP2PFS never uncovered either people in general or private key over the system. Amid the disclosure stage the *BLAKE2b* hash of the open key is utilized as the discovery key. This implies the first key is unimaginable to find (*except if it was shared openly through a separate channel*) since just the hash of the key is uncovered openly.

### B. Decentralized Mirroring

DP2PFS is a distributed convention intended to trade bits of a data set among a swarm of companions. When a companion gains their first bit of information in the data set they can turn into an incomplete mirror for the data set. In the event that another person gets in touch with them and necessities the piece they have, they can share it. This can happen at the same time while the companion is as yet downloading the pieces they need from others.

1) Peer Connections: After the disclosure stage, DP2PFS ought to have a rundown of potential information sources to attempt and contact. DP2PFS employments either TCP, HTTP or UTP (Rossi et al. 2010). UTP utilizes LEDBAT which is intended to not take up all accessible data transfer capacity on a system (for example with the goal that other individuals sharing WiFi can in any case utilize the Internet), and is as yet dependent on UDP so works with NAT traversal procedures like UDP gap punching. HTTP is supported for similarity with static record servers and internet browser customers. Note that these are the conventions we support in the reference DP2PFS usage, yet the convention itself is transport rationalist. On the off chance that an HTTP source is determined repository will favor that one over different sources. Generally when DP2PFS gets the IP and port for a potential TCP or UTP source it endeavors to interface utilizing the two conventions. On the off chance that one interfaces to begin with, DP2PFS prematurely ends the other one. On the off chance that none associate, DP2PFS will attempt again until it chooses that source is disconnected or inaccessible and afterward quits endeavoring to associate with them. Sources DP2PFS can associate with go into a rundown of known great sources, so that if/when the Internet association goes down DP2PFS can utilize that rundown to reconnect to known great sources again rapidly.

On the off chance that DP2PFS gets a great deal of potential sources it picks a bunch at irregular to attempt and associate with and keeps the rest around as extra sources to utilize later on the off chance that it chooses it needs more sources. When a duplex double association with a remote source is open DP2PFS at that point layers on the Hypercore convention, a message-based replication convention that permits two companions to impart over a stateless station to re-journey and trade information. You

open separate imitation channels with numerous companions on the double which permits customers to parallelize information asks for over the whole pool of companions they have built up associations with.

### C. Network Privacy

On the Web today, with SSL, there is a certification that the traffic between your PC and the server is private. For whatever length of time that you believe the server to not spill your logs, assailants who block your system traffic won't almost certainly read the HTTP traffic traded among you and the server. This is a genuinely straight-forward model as customers just need to confide in a solitary server for some space.

There is an inborn trade off in distributed frameworks of source revelation versus client security. The more sources you contact and request a little information, the more sources you trust to keep what you requested private. Our objective is to have Dat be configurable in regard to this trade off to enable application designers to meet their possess security rules.

It is up to customer projects to settle on plan choices around which revelation systems they trust. For precedent if a Dat customer chooses to utilize the BitTorrent DHT to find companions, and they are hunting down openly shared Dat key (for example a key referred to openly in a distributed logical paper) with known substance, at that point as a result of the security structure of the BitTorrent DHT it winds up open information what key that customer is scanning for.

A customer could decide to just utilize disclosure systems with certain security ensures. For instance a customer could just interface with an endorsed rundown of sources that they trust, like SSL. For whatever length of time that they trust each source, the encryption incorporated with the Dat organize convention will keep the Dat key they are searching for from being spilled.

### D. Incremental Versioning

Given a surge of parallel information, we part the stream into pieces, hashes each lump, and masterminds the hashes in a particular kind of Merkle tree that permits for certain replication properties.

DP2PFS is likewise ready to completely or in part synchronize streams in a conveyed setting regardless of whether the stream is being added to. This is practiced by utilizing the informing convention to cross the Merkle tree of remote sources and get a key arrangement of hubs.

Because of the low-level, message-arranged plan of the replication convention, distinctive hub traversal systems can be executed.

There are two types of versioning performed automatically by DP2PFS. Metadata is stored in a folder called .dat in the root folder of a repository, and data is stored as normal files in the root folder.

1. Metadata Versioning: DP2PFS attempts however much as could be expected to go about as a coordinated reflection of the condition of an envelope and every one of its substance.

When bringing in documents, DP2PFS utilizes an arranged, profundity first recursion to list every one of the documents in the tree. For each document it discovers, it gets the file system metadata (filename, Stat object, and so on) and checks if there is as of now a passage for this filename with this precise metadata as of now spoken to in the DP2PFS store metadata. On the off chance that the document with this metadata coordinates precisely the most up to date rendition of the record metadata put away in DP2PFS, at that point this record will be skipped (no change). On the off chance that the metadata contrasts from the current existing one (or then again there are no sections for this filename at all in the history), at that point this new metadata passage will be attached as the new 'most recent' adaptation for this record in the attach just SLEEP metadata content register.

2. Content Versioning: Notwithstanding putting away a verifiable record of file system metadata, the substance of the records themselves are too fit for being put away in a variant controlled way. The default stockpiling framework utilized in DP2PFS stores the documents as records. This has the upside of being straight-forward for clients to see, however the drawback of not putting away old variants of substance of course.

As opposed to other form control frameworks like Git, DP2PFS as a matter of course just stores the present arrangement of checked out documents on plate in the vault organizer, not old adaptations. It stores all past metadata for old forms in .dat. Git for instance stores all past content renditions and all past metadata variants in the .git organizer. Since Dat is intended for bigger data sets, in the event that it put away all past record forms in .dat, at that point the .dat envelope could without much of a stretch

top off the clients hard drive coincidentally. Along these lines DP2PFS has numerous capacity modes dependent on utilization.

Hypercore registers incorporate a discretionary information record that stores all lumps of information. In DP2PFS, just the metadata.data record is utilized, however the content.data record isn't utilized. The default conduct is to store the current records just as expected documents. On the off chance that you need to run a 'recorded' hub that keeps every single past form, you can design DP2PFS to utilize the content.data record. For instance, on a mutual server with heaps of capacity you most likely need to store all variants. However on a workstation machine that is just getting to a subset of one form, the default method of putting away all metadata in addition to the present arrangement of downloaded documents is satisfactory, in light of the fact that you realize the server has the full history.

3. Merkle Trees: Registers in DP2PFS use a specific method of encoding a Merkle tree where hashes are positioned by a scheme called binary in-order interval numbering or just "bin" numbering. This is just a specific, deterministic way of laying out the nodes in a tree. For example a tree with 7 nodes will always be arranged like this: 0 1 2 3 4 5 6 In DP2PFS, the hashes of the chunks of files are always even numbers, at the wide end of the tree. So the above tree had four original values that become the even numbers: chunk0 -> 0 chunk1 -> 2 chunk2 -> 4 chunk3 -> 6

In the resulting Merkle tree, the even and odd nodes store different information:

i) Evens List of data hashes [chunk0, chunk1, chunk2, .

ii) Odds List of Merkle hashes (hashes of child even nodes) [hash0, hash1, hash2, .

These two lists get interleaved into a single register such that the indexes (position) in the register are the same as the bin numbers from the Merkle tree. All odd hashes are derived by hashing the two child nodes, e.g. given hash0 is hash(chunk0) and hash2 is hash(chunk1), hash1 is hash(hash0 + hash2). For example a register with two data entries would look something like this (pseudo code):

i) hash(chunk0)

ii) hash(hash(chunk0) + hash(chunk1))

iii) hash(chunk1)

It is possible for the in-order Merkle tree to have multiple roots at once. A root is defined as a parent node with a full set of child node slots filled below it.

For example, this tree has 2 roots (1 and 4)

0 1 2 4 This tree has one root (3): 0 1 2 3 4 5 6 This one has one

root (1): 0 1 2

### Random Access

DP2PFS pursues the following access capabilities:

- Support large file hierarchies (millions of files in a single repository).
- Support efficient traversal of the hierarchy (listing files in arbitrary folders efficiently).
- Store all changes to all files (metadata and/or content). • List all changes made to any single file.
- View the state of all files relative to any point in time.
- Subscribe live to all changes (any file). • Subscribe live to changes to files under a specific path.
- Efficiently access any byte range of any version of any file.
- Allow all of the above to happen remotely, only syncing the minimum metadata necessary to perform any action.
- Allow efficient comparison of remote and local repository state to request missing pieces during synchronization.
- Allow entire remote archive to be synchronized, or just some subset of files and/or versions.

The way DP2PFS accomplishes these are through a combination of storing all changes in Hypercore feeds, but also using strategic metadata indexing strategies that support certain queries efficiently to be performed by traversing the Hypercore feeds.

#### 1) Scenario: Consuming a File from a Specific Offset

Tom has a data set in Dat, Bobby needs to get to a 200MB CSV called dog\_dna.csv put away in the remote storehouse, yet just needs to get to the 20MB scope of the CSV traversing from 40MB 50MB.

Tom has never spoken with Alice, and is beginning crisp with no learning of this Dp2pfs store other than that he realizes he needs dog\_dna.csv at a particular balance.

To begin with, Tom approaches Bobby through the dp2pfs convention for the metadata he needs to determine dog\_dna.csv to the right metadata feed passage that speaks to the document he needs.

Note: In this situation we accept Tom needs the most recent adaptation of dog\_dna.csv. It is additionally possible to do this for a particular more established form. Tom initially sends a Request message for the most recent section in the metadata feed. Bobby reacts. Tom takes a gander at the esteem, and utilizing the query calculation depicted underneath sends another Request message for the metadata hub that is nearer to the filename he is searching for. This rehashes until Alice send Bob the Coordinating metadata section. This is the goals that utilizes  $\log(n)$  round treks, however there are approaches to upgrade this by having Tom send extra grouping numbers to Bob that assistance him cross in less round excursions.

In the metadata record Tom got for dog\_dna.csv there is the byte balanced to the start of the document in the information feed. Tom adds his +30MB counterbalance to this esteem and begins mentioning bits of information beginning at that byte balance.

This technique attempts to permit any byte scope of any document to be gotten to without the need to synchronize the full metadata for all records in advance.

## EXISTING WORK

### A. Git

Git advanced the possibility of a directed non-cyclic diagram (DAG) joined with a Merkle tree, an approach to speak to changes to information where each change is tended to by the protected hash of the change in addition to all progenitor hashes in a chart.

This gives an approach to confide in information trustworthiness, as the main way a particular hash could be inferred by another companion is on the off chance that they have similar information and change history required to imitate that hash.

This is imperative for reproducibility as it gives you a chance to trust that a particular git submit hash alludes to a particular source code state.

Decentralized form control instruments for source code like Git give a convention to proficiently downloading changes to a lot of records, yet are improved for content records and have issues with huge documents.

Arrangements like GitLFS settle this by utilizing HTTP to download vast documents, instead of the Git convention. GitHub offers Git-LFS facilitating yet charges archive proprietors for data transfer capacity on well-known records. Building a circulated conveyance layer for documents in a Git archive is troublesome due to plan of Git Packfiles which are delta compacted archive expresses that don't effortlessly bolster arbitrary access to byte goes in past record adaptations.

### B. BitTorrent

BitTorrent executes a swarm based document sharing convention for static data sets. Information is part into fixed measured pieces, hashed, and afterward that hash is utilized to find peers that have similar information. Favorable position of utilizing BitTorrent for data set exchanges is that download transfer speed can be completely immersed. Since the record is part into pieces, and companions can productively find which pieces every one of the friends they are associated with have, it implies one companion can download non-covering locales of the from numerous friends at the equivalent time in parallel, amplifying system throughput.

Fixed estimated lumping has downsides for information that changes. BitTorrent accept all metadata will be exchanged in advance which makes it unreasonable for gushing or refreshing substance. Generally BitTorrent customers isolate information into 1024 pieces meaning substantial could have a substantial lump estimate which impacts arbitrary access execution (for example for gushing video).

Another downside of BitTorrent is because of the way customers promote and find different companions in nonappearance of any convention level protection or trust. From a client security point of view, BitTorrent spills what clients are getting to or endeavoring to get to, and does not give a similar perusing protection works as frameworks like SSL.

### C. Kademlia Distributed Hash Table

Kademlia ( and 2002) is a disseminated hash table, an appropriated key/esteem store that can fill a comparative need to DNS servers yet has no hard coded server addresses. All customers in Kademlia are additionally servers. For whatever length of time that you know something like one address of another companion in the system, you can inquire them for the key you are endeavoring to discover, and they will either have it or give you some other individuals to talk to that are bound to have it.

On the off chance that you don't have an underlying friend to converse with you, most customers utilize a bootstrap server that haphazardly gives you a friend in the system to begin with. In the event that the bootstrap server goes down, the system still capacities as long as different strategies can be utilized to bootstrap new peers, (for example, sending them peer addresses through side channels like how .downpour documents incorporate tracker addresses to attempt on the off chance that Kademlia finds no friends).

Kademlia is unmistakable from past DHT structures due to its effortlessness. It utilizes a straightforward XOR activity between two keys as its "remove" metric to choose which peers are nearer to the information being hunt down.

On paper, it appears as though it wouldn't fill in as it doesn't consider things like ping latency or data transfer throughput. Rather, its structure is extremely straightforward deliberately to limit the measure of control/tattle messages and to limit the measure of multifaceted nature required to actualize it. By and by Kademlia has been amazingly fruitful and is broadly conveyed as the "Mainline DHT" for BitTorrent, with help in all famous BitTorrent customers today.

Because of the effortlessness in the first Kademlia plan various assaults, for example, DDOS as well as have been illustrated. There are convention augmentations () which in specific cases relieve the impacts of these assaults, for example, BEP 44 which incorporates a DDOS relief method. In any case anybody utilizes Kademlia ought to know about the confinements.

### D. Peer to Peer Streaming Peer Protocol

PPSPP (IETF RFC 7574, (Bakker, Petrocco, and Grishchenko 2015)) is a convention for live spilling content over a shared system. In it they characterize a particular sort of Merkle Tree that takes into account subsets of the hashes to be mentioned by a friend so as to lessen the time-till-playback for end clients. BitTorrent for instance exchanges all hashes in advance, which isn't reasonable for live spilling.

Their Merkle trees are requested utilizing a plan they call "receptacle numbering", which is a strategy for deterministically organizing an affix just log of leaf hubs into an all together format tree where non-leaf hubs are inferred hashes.

In the event that you need to check a particular hub, you just need to demand its kin's hash and all its uncle hashes. PPSPP is exceptionally worried about lessening round outing time and time-till-playback by taking into account numerous sorts of improvements, for example, to pack the same number of hashes into data grams as conceivable while trading tree data with friends.

Despite the fact that PPSPP was planned with spilling video as a primary concern, the capacity to demand a subset of metadata from a huge as well as gushing data set is entirely alluring for some different sorts of data sets.

### E. WebTorrent

With WebRTC, programs would now be able to make distributed associations legitimately to different programs. BitTorrent utilizes UDP attachments which aren't accessible to program JavaScript, so can't be utilized as-is on the Web.

WebTorrent actualizes the BitTorrent convention in JavaScript utilizing WebRTC as the vehicle. This incorporates the BitTorrent square trade convention as well as the tracker convention actualized as it were that can empower half-and-half hubs, talking all the while to both BitTorrent and WebTorrent swarms (if a customer is fit for making both UDP attachments too as WebRTC attachments, for example, Node.. Trackers are presented to web customers over HTTP or WebSockets.

### F. Inter Planetary File System

IPFS is a group of utilization and system conventions that have distributed record sharing and information lasting heated in.

IPFS abstracts organize conventions also, naming frameworks to give an elective application conveyance stage to the present Web. For instance, rather than utilizing HTTP and DNS straightforwardly, in IPFS you would utilize LibP2P streams and IPNS so as to access the highlights of the IPFS stage.

### G. Certificate Transparency / Secure Registers

The UK Government Digital Service have created the idea of a register which they characterize as a computerized open record you can trust. In the UK government registers are starting to be steered as an approach to uncover basic open informational indexes in a manner where purchasers can confirm the information has not been altered, and enables the information distributors to refresh their informational indexes after some time.

The plan of registers was propelled by the foundation backing the Certificate Transparency (Laurie, Langley, and Kasper 2013) venture, started at Google, which gives an administration over SSL testaments that empowers specialist organizations to compose testaments to a dispersed open record. Any customer or specialist co-op can check if an authentication they got is in the record, which secures against purported "rebel testaments".

## CONCLUSION

This is study that explores the possibility of having a decentralized medium to transfer files and synchronize between them, which is available in every computing device. The main advantages of an implementation is Data Privacy, faster data replication in highly changing scenarios. Using this implementation, we can harness the power of versioning by changing only files that needs to be changed rather than keeping two immutable versions of a file.

Having conducted an expansive survey on the Topic of Decentralized File Sharing and Syncing, I have dropped to a conclusion that the vision of the project is perfectly aligned with what the survey says. There is a need of a software / implementation which makes decentralized file sharing a common thing and is also robust enough for everyday use.

There are various researches conducted in this field, but none seem to address the crucial user privacy related chores of every technology possible, one of the research paper proposes creation of a protocol which is secure and robust by default, but lacks a robust solution for it. Hence, I have come to a conclusion that I shall be moving forward with this project and creating a proof of concept implementation of the so called protocol.

## REFERENCES

- [1] Git [<https://en.wikipedia.org/wiki/Git>]
- [2] BitTorrent [<http://web.cs.ucla.edu/classes/cs217/05BitTorrent.pdf>]
- [3] Inter Planetary File System  
[<https://www.researchgate.net/publication/263930348IPFS-ContentAddressedVersionedP2PFileSystem>]
- [4] Kademlia Distributed Hash Table [<https://pub.tik.ee.ethz.ch/students/2006-So/SA-2006-19.pdf>]
- [5] SLEEP-Syncable Ledger of Exact Events Protocol  
<https://github.com/datproject/docs/blob/master/papers/sleep.pdf>]
- [6] Research on file synchronization backup system based on security policy  
[<https://ieeexplore.ieee.org/document/6014262>]
- [7] File Synchronization and Sharing : User Practices and Challenges  
[<https://www.asis.org/asist2014/proceedings/submissions/papers/161paper.pdf>]