

Peer to Peer - Appunti

Francesco Lorenzoni

February 2024

Contents

1	Introduction	7
1.1	Blockchain concepts	7
1.1.1	TriLemma	8
1.2	P2P Systems	8
1.2.1	Semi-Decentralized systems	8
1.2.2	Fully decentralized systems	8
1.3	P2P Overlay network	8
1.3.1	Unstructured overlay	8
1.3.2	Structured overlays	10
1.3.3	Hierarchical overlays	10
1.3.4	Summary	12
2	Distributed Hash Tables	13
2.1	Building DHT	14
2.1.1	Peers joining and leaving	14
2.2	Data Lookup	15
2.2.1	Addressing data	15
2.2.2	API, Lookup and Various Properties	16
3	Kademlia	17
3.1	Structure	17
3.1.1	Assigning keys to leaves	18
3.2	Distance - XOR Metric	18
3.3	Routing Table	18
3.4	Key Lookup	18
3.5	Protocol Messages	18
4	BitTorrent	19
4.1	Deeper into BitTorrent	19
4.1.1	Glossary	19
4.1.2	Protocol Overview	20
4.2	Pieces selection	20
4.2.1	Free Riders	21
4.3	DHT and BitTorrent	21
5	Blockchain	23

Course info

...

Chapter 1

Introduction

Opposed to Client-server architectures where there are end-hosts and dedicated-hosts (servers), in P2P Systems there are only end-nodes which directly communicate with each other; they have an “on/off” behaviour, and they handle **churn**¹. However, in P2P systems servers are still needed, but only as *bootstrap servers*, typically allowing for new nodes to easily join the P2P network.

Peers’ connection in P2P is called *transient*, meaning that connections and disconnections to the network are very frequent.

Notice that since each time a peers connects to the P2P network it may have a different IP address, resources cannot be located using IP, but a different method at application layer must be used.

Definition 1.1 (P2P System) *A peer to peer system is a set of autonomous entities (peers) able to auto-organize and sharing a set of distributed resources in a computer network.*

The system exploits such resources to give a service in a complete or partial decentralized way

Definition 1.2 (P2P System - Alternative definition) *A P2P system is a distributed system defined by a set of nodes interconnected able to auto-organize and to build different topologies with the goal of sharing resources like CPU cycles, memory, bandwidth. The system is able to adapt to a continuous churn of the nodes maintaining connectivity and reasonable performances without a centralized entity (like a server)*

1.1 Blockchain concepts

Definition 1.3 (Blockchain) ◊ *a write-only, decentralized, state machine that is maintained by untrusted actors, secured by economic incentive*
◊ *cannot delete data*
◊ *cannot be shut down or censored*
◊ *supports defined operations agreed upon by participants*
◊ *participants may not know each other (public)*
◊ *in actors best interest is to play by the rules*

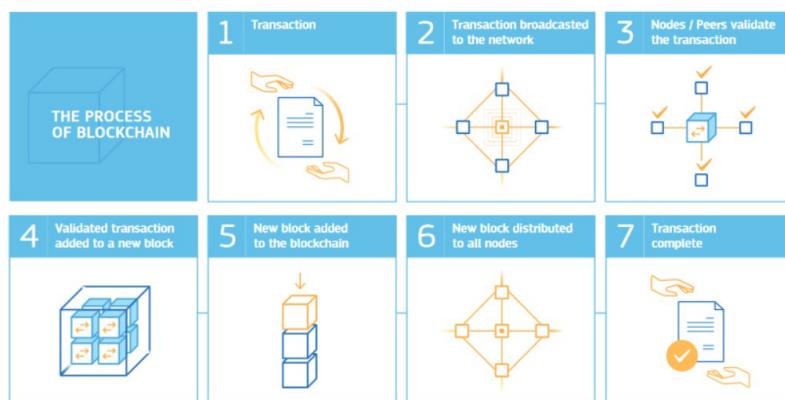


Figure 1.1: Blockchain process

¹ “churn” will be a recurring term. In Italian it means “rimescolare”

Bitcoin were developed as an alternative way to exchange money which wouldn't need intermediaries such as banks. Today, *Ethereum* is becoming more and more popular. NFT² allow to establish the owner of a digital artwork, by generating a token using a blockchain.

1.1.1 TriLemma

The Blockchain **trilemma** states that a blockchain **cannot** simultaneously provide *Decentralization*, *Security* and *Scalability*.

1.2 P2P Systems

1.2.1 Semi-Decentralized systems

An example is **Napster**, released in 2001. Napster used servers only to allow users to locate peers which could provide the desired file, delegating the actual file exchange to peers, allowing for a very few server needed.

For the first time users are called *peers*, and the systems implemented in this way *peer-to-peer systems*

Napster had many strengths common to many P2P systems, from whose emerges the ability of peers to act both as server and a client, but also suffered from weaknesses derived from its centralization, at least for “node discovery”. Napster centralized server represents a design bottleneck, and also made it target of legal attacks.

1.2.2 Fully decentralized systems

Gnutella is similar to Napster, but here no centralized server exists. Peers establish *non-transient* direct connections to search files, not to actually transfer them.

- Cons*
1. High network traffic
 2. No structured search
 3. Free-riding

1.3 P2P Overlay network

In P2P systems there is an overlay network at application level operating on top of the underlying (IP) network.

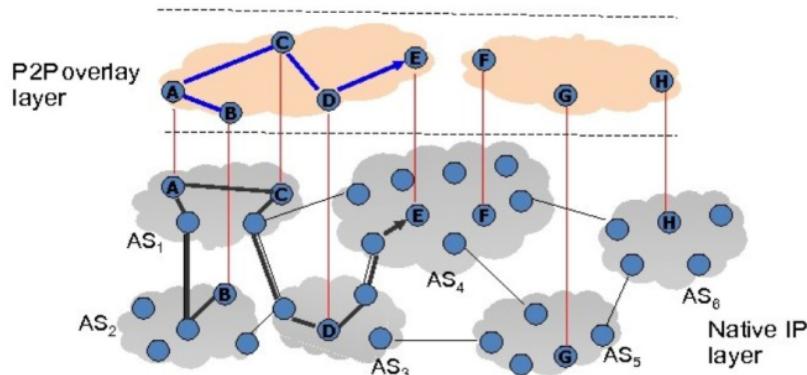


Figure 1.2: P2P Overlay networks

A P2P **protocol**—defined over the P2P overlay— defines the set of messages that the peers exchange.

1.3.1 Unstructured overlay

²Non Fungible Tokens

The two key issues here are:

- ◊ how to **bootstrap** on the network?
- ◊ how to **find content** without a central index?

Possible lookup algorithms are the following, but they all are not very scalable, and are costful in terms of performance:

- ◊ **Flooding**
- ◊ **Expanding ring**
- ◊ **Random walk**

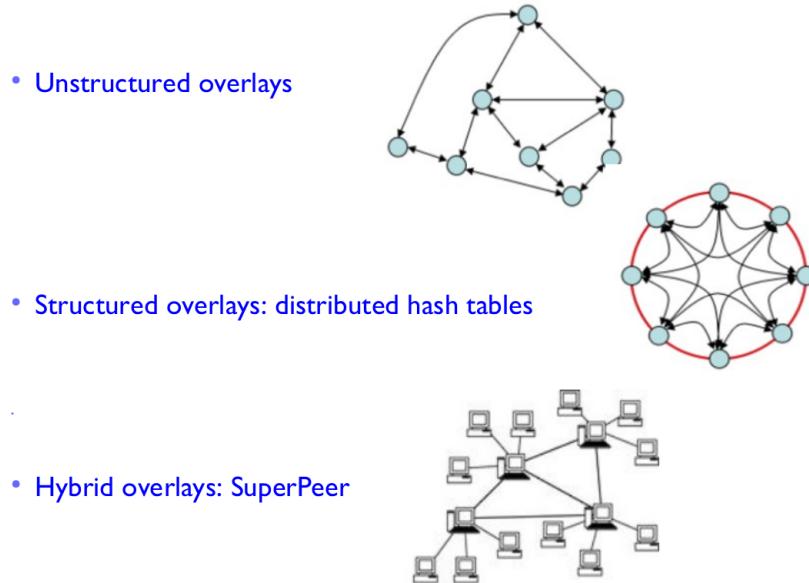


Figure 1.3: P2P Overlay Network classification

Flooding

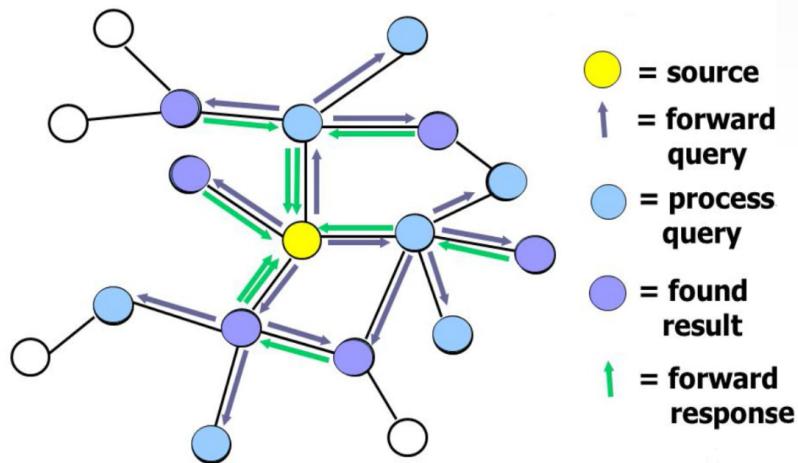


Figure 1.4: Flooding search in unstructured overlay

Messages have a *TTL* to limit the number of hops when propagating, but also a *unique identifier* to detect cycles.
 Flooding is not only for searching, but also to propagate transactions in the P2P network underlying a **blockchain**

Expanding Ring

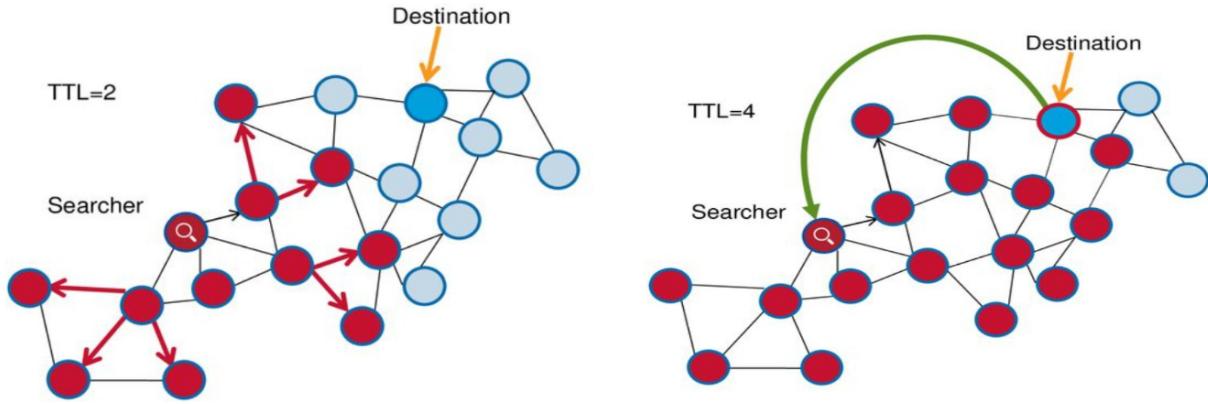


Figure 1.5: Expanding ring/Iterative Deepening

This technique consists in repeated flooding with an increasing TTL, implementing a BFS search.

Random walk

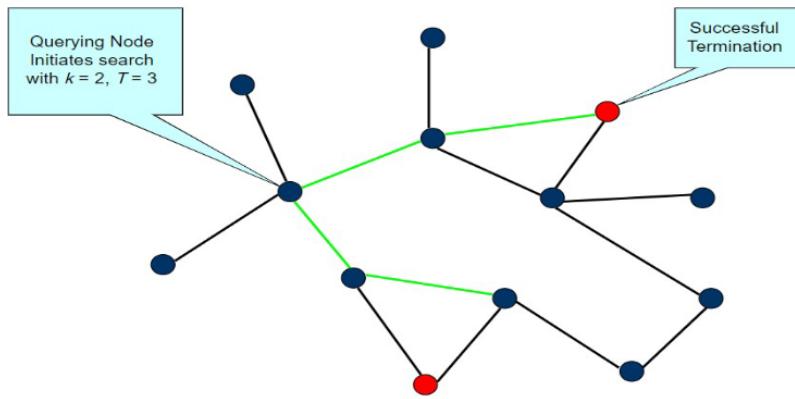


Figure 1.6: Random Walk

k indicates the number of “walkers” to be generated by the querying node. Each green path corresponds to a “walker”.

A path is constructed by taking successive (single) steps in random directions defined by a *Markov Chain*, which is “memory-less”. Note that only one successor node is chosen at each step.

The path is bounded by a TTL.

Random walk avoids an exponential increase in the amount of messages, which becomes an issue for vastly populated networks.

Paths can be stopped by a TTL but also by checking periodically with the destination whether the stop condition has been met.

A querying node can also bias its walks towards high-degree nodes: higher probability to choose the highest degree neighbor.

1.3.2 Structured overlays

The choice of the neighbours is defined according to a given criteria, resulting in a **structured** overlay network. The goal is to guarantee scalability by providing:

- ◊ *key-based lookup*
- ◊ information lookup has a given *complexity* e.g. $\mathcal{O}(\log N)$

1.3.3 Hierarchical overlays

Peers connect to **Super-Peers** which know (“*index*”) Peer resources. The flooding is restricted to Super-Peers, but still allowing for resources to be directly exchanged between the peers.

Lookup complexity is less and the scalability is improved, but there is a lower resistance to super-peers churn.

In some cases —such as Gnutella— peers are “*self-promoted*” to super-peers, while in others they are statically defined.

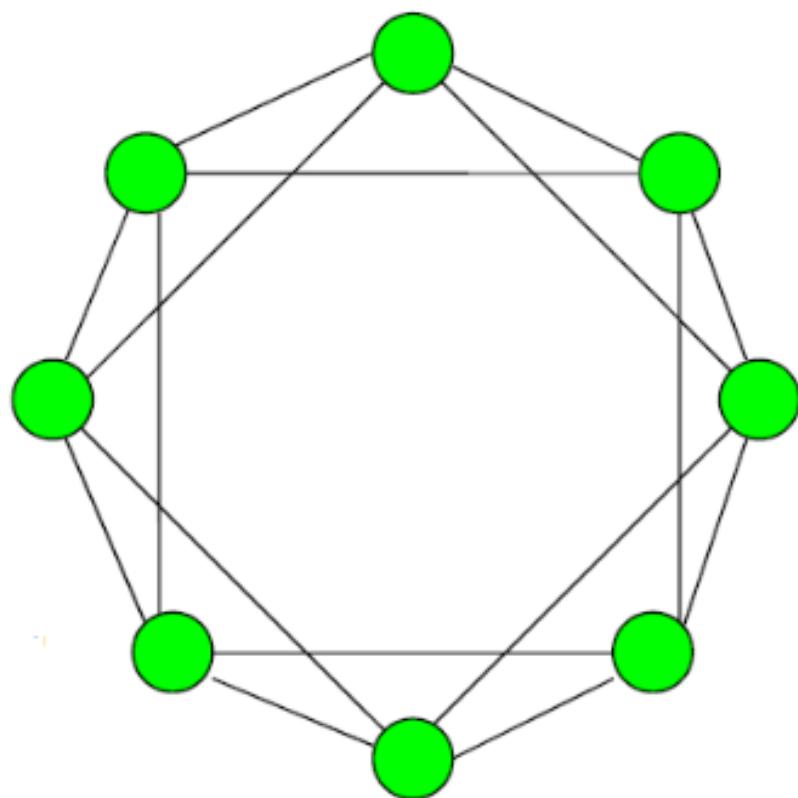


Figure 1.7: Structured overlay

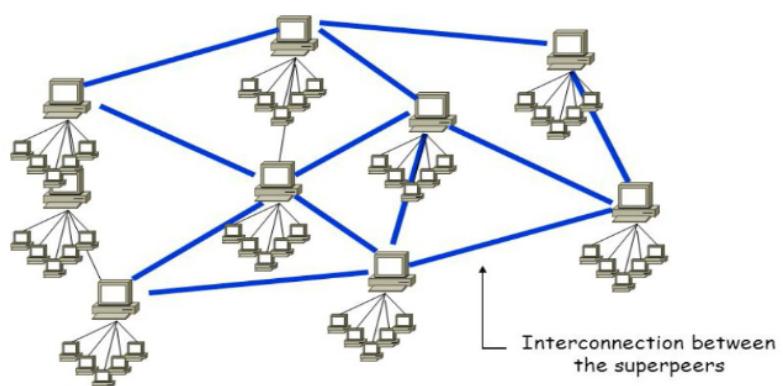


Figure 1.8: Hierarchical overlay

1.3.4 Summary

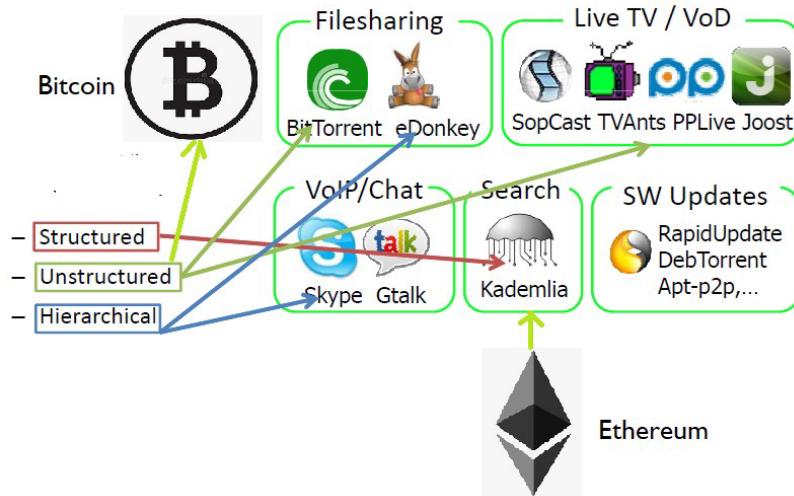


Figure 1.9: Overlay structure for known applications

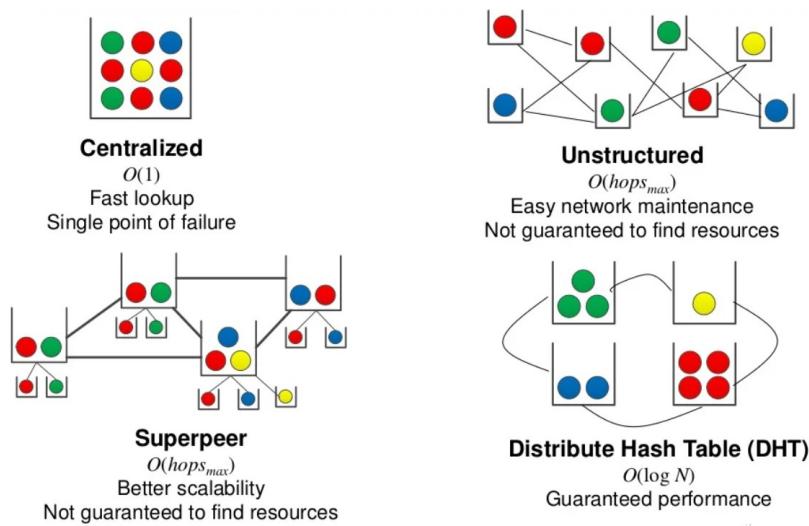


Figure 1.10: Overlays summary
DHT will be discussed in the next chapter

Chapter 2

Distributed Hash Tables

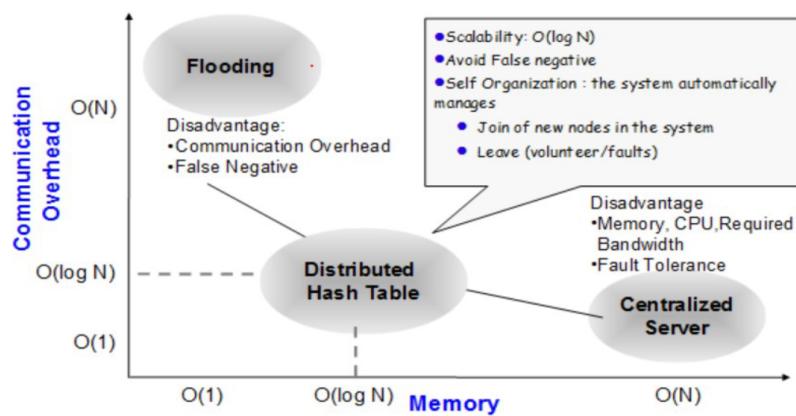


Figure 2.1: DHT Motivations

The key idea is to split the hash tables into several parts and distribute them to several servers, and to use hash of resources (or of the URLs of resources) as a key to map them to a dynamically changing set of web caches, but with each key mapped to single server; so that each machine (user) can locally compute which web cache should contain the required resource, referenced by an URL.

This technique is extended to DHT for P2P systems.

However, rehashing is a problem in dynamic scenarios if the hashing scheme depends directly on the number of servers: 99% of keys have to be remapped, resulting in a lot of messages exchange.



Figure 2.2: Rehashing problem

Consistent hashing is a set of hash techniques which guarantees that adding more nodes/remove nodes implies moving only a minority of data items. Each node manages—instead of a set of sparse keys—an interval of consecutive hash keys, and intervals are joined/splitted when nodes join/leave the network and keys redistributed between adjacent peers.

2.1 Building DHT

- ◊ Use a logical name space, called *identifier space* consisting of identifiers $\{0, 1, 2, \dots, N - 1\}$
- ◊ define identifier space as a *logical ring* modulo N
- ◊ every node picks a random identifier through Hash H .

```
space N=16 {0,...,15}
• five nodes a, b, c, d, e
• H(a) = 6
• H(b) = 5
• H(c) = 0
• H(d) = 11
• H(e) = 2
```

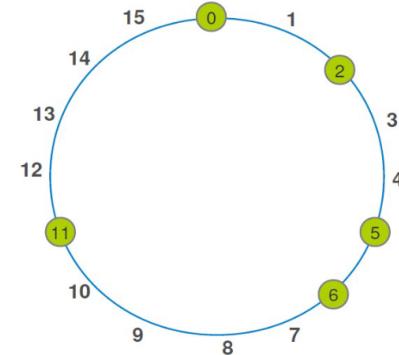


Figure 2.3: Identifier space

2.1.1 Peers joining and leaving

When a new node is **added**, we map the keys between the new node and the previous node in the hash ring to point to the new node; those the keys will no longer be associated with their old nodes.

When a node is **removed** from the hash ring, only the keys associated with that node are rehashed and remapped rather than remapping all the keys.

In case a node suddenly disconnects from the network, all data stored on it are lost if they are not stored on other nodes; to avoid such a problem:

- ◊ introduce some redundancy (data replication)
- ◊ information loss: periodical information refresh

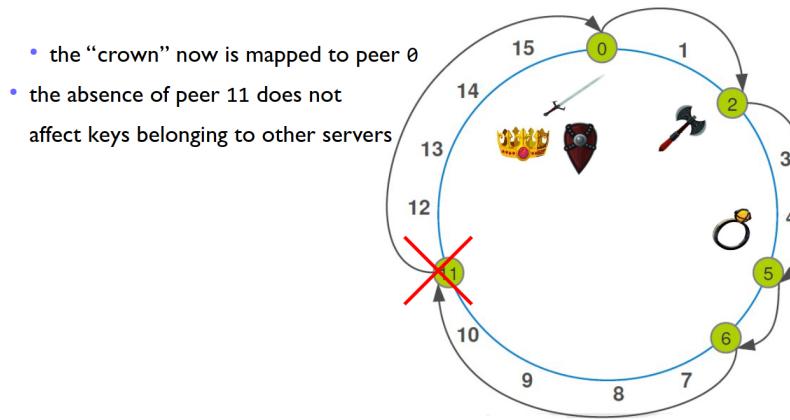


Figure 2.3: Peer 11 leaves the Network
In case a peer leaves, its keys can easily be remapped to its successor

When the hash table is **resized**, on the average, only $\frac{k}{n}$ keys need to be remapped on average, where k is the number of keys and n is the number of servers.

2.2 Data Lookup

- finger/routing table:
 - point to `succ(n+1)`
 - point to `succ(n+2)`
 - point to `succ(n+4)`
 - point to `succ(n+8)`
 - ...
 - point to `succ(n+2M-1)` (M number of bits for the identifiers)
- distance always halved to the destination.

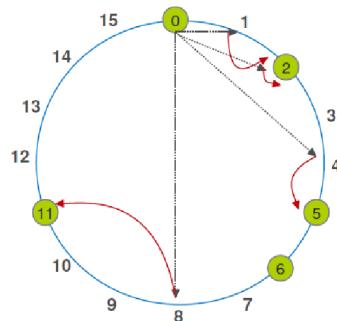


Figure 2.4: Exponential Search for DHT

The data lookup can be implemented by using exponential search, rather than performing a walk by asking each peer for its successor

Data Lookup can be sped up even more, by computing the hash $h(x)$ of the searched object, and propagating the query to farthest node¹ which has an identifier smaller than $h(x)$, which then recursively applies the same algorithm, until the object is found.

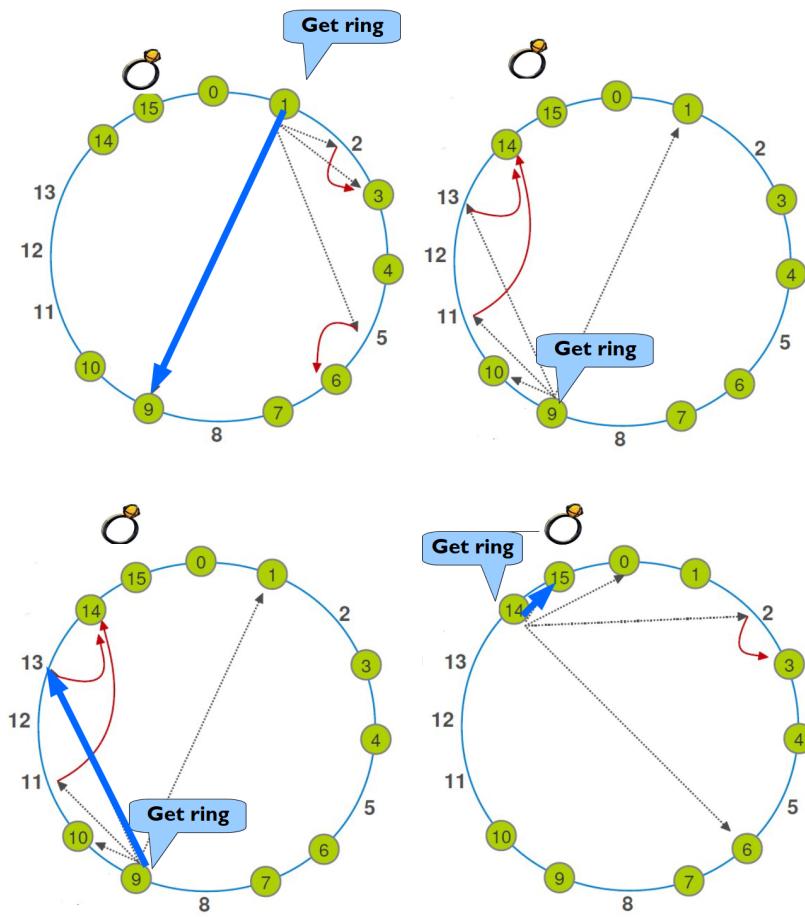


Figure 2.5: Lookup performed in the CHORD DHT

2.2.1 Addressing data

Data was usually addressed by **location**, a `http://` link to locate resources; Such link is an identifier that points to a particular location on the web.

This approach forces us all to pretend that the data are in only one location.

¹Which is found using exponential search

IPFS instead uses **content addressing**, which exploits the cryptographic hash of the content to identify it.

2.2.2 API, Lookup and Various Properties

Most DHT provide a simple interface **PUT, GET, Value**, usually without the possibility to move keys.

Approach	Memory for each node	Communication Overhead	Complex Queries	False Negatives	Robustness
Central Server	$O(N)$	$O(1)$	✓	✓	✗
Pure P2P (flooding)	$O(1)$	$O(N^2)$	✓	✗	✓
DHT	$O(\log N)$	$O(\log N)$	✗	✓	✓

Figure 2.6: Lookup time complexity comparison

DHT

- ◊ Routing is based on key (unique identifier)
- ◊ Key are uniformly distributed to the DHT nodes
 1. Bottleneck avoidance
 2. Incremental insertion of the keys
 3. Fault tolerance
- ◊ Auto organizing system
- ◊ Simplex and efficient organization
- ◊ The terms “Structured Peer-to-Peer“ and “DHT“ are often used as synonyms

Chapter 3

Kademlia

Kademlia is a protocol used by some of the largest public DHTs

- ◊ BitTorrent Mainline DHT
- ◊ Ethereum P2P network
- ◊ IPFS

It has three key characteristics which are not offered by other DHTs

1. routing information spreads automatically as a side-effect of lookups
2. flexibility to send multiple requests in parallel to speed up lookups by avoiding timeout delays (parallel routing)
3. iterative routing

At each routing step of the query, the queried node sends a report to the starting querying node, even if it could not answer the query.

3.1 Structure

Kademlia exploits the leaves of a **Trie**¹ to define the logical identifier space;

Note that not all leaves correspond to nodes (peers)

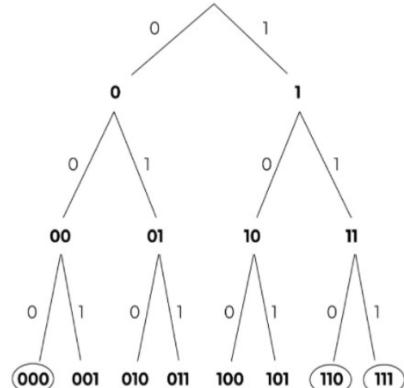


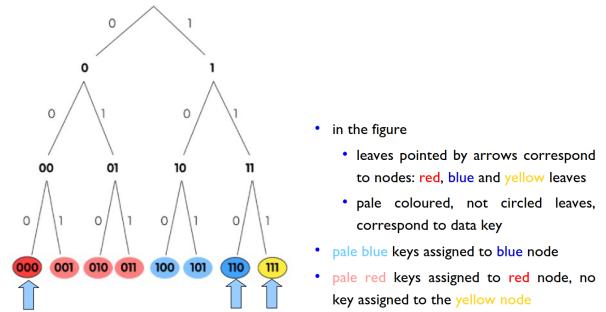
Figure 3.1: Trie

¹k-ary search tree and prefix tree

3.1.1 Assigning keys to leaves

The rule to partition the keys (content) among the nodes must respect the rules of *consistent hashing*.

Definition 3.1 (Partitioning rule) A key is assigned to the node with the “lowest common ancestor”:
Find the longest prefix between the key and the node identifier, and then assign the key to such node.



3.2 Distance - XOR Metric

How to compute the *distance* between two nodes?

3.3 Routing Table

In order to look for data, Kademlia’s key idea is to store a logarithmic number of node IDs and their corresponding IP addresses and some contact taken from the identifier trie.

TODO

3.4 Key Lookup

TODO

3.5 Protocol Messages

TODO reformat this

1. **FIND_NODE** $v \rightarrow w$ (T) (v,w nodes, T target of the look up)

the recipient of the message (w) returns k (IP address, UDP port, Node ID) triples for the k nodes it knows about closest to the target T.

these triples can come from a single k-bucket, or they may come from multiple k-buckets if the closest k-bucket is not full.

in any case, the recipient must return k items, unless there are fewer than k nodes in all its k-buckets combined, in which case it returns every node it knows about

2. **FIND_VALUE** $v \rightarrow w(T)$ (v, w nodes, T value looked up)

in: T, 160-bit ID representing a value

out: if a value corresponding to T is present in the queried node (w), the associated data is returned otherwise it is equivalent to **FIND_NODE** and w returns a set of k triples

If **FIND_VALUE** returns a list of other peers, it is up to the requester to continue searching for the desired value from that list

3. **PING** $v \rightarrow w$

probe node w to see if its online

4. **STORE** $v \rightarrow w$ (Key, Value)

instructs node w to store a $\langle key, value \rangle$ pair
the node has been retrieved through a

The actual **Lookup algorithm** is based on **FIND_NODE**. many **FIND_NODE** can be executed in parallel, according to α that is a system-wide concurrency parameter.

With $\alpha = 1$, the lookup algorithm is similar to *Chord*, one step progress each time

Lookup procedure is the same for **FIND_VALUE** and **FIND_NODE**

Chapter 4

BitTorrent

The goal of *Content Distribution Networks* is to distribute web contents to hundreds of thousands or millions of simultaneous users, exploiting data and/or service replication on different **mirror servers**.

In **P2P CDN** the initial file request are served by a centralized server, and further requests served by peers which have already received and replicated the files (**seeders**), without involving the initial server.

BitTorrent in a nutshell

- ◊ Basically a *Content Distribution Network* (CDN)
- ◊ A distributed set of hosts cooperating to distribute large data set to end users.
- ◊ Efficient content distribution systems using *file swarming*
- ◊ Does *not* perform all the functions of a typical P2P system, like searching
- ◊ Rather than providing a search protocol itself, was designed to integrate seamlessly with the Web and made file descriptors available via Web, which could be searched with standard Web search
- ◊ *File swarming*: a peer makes whatever portion of the file that is downloaded immediately available for sharing

4.1 Deeper into BitTorrent

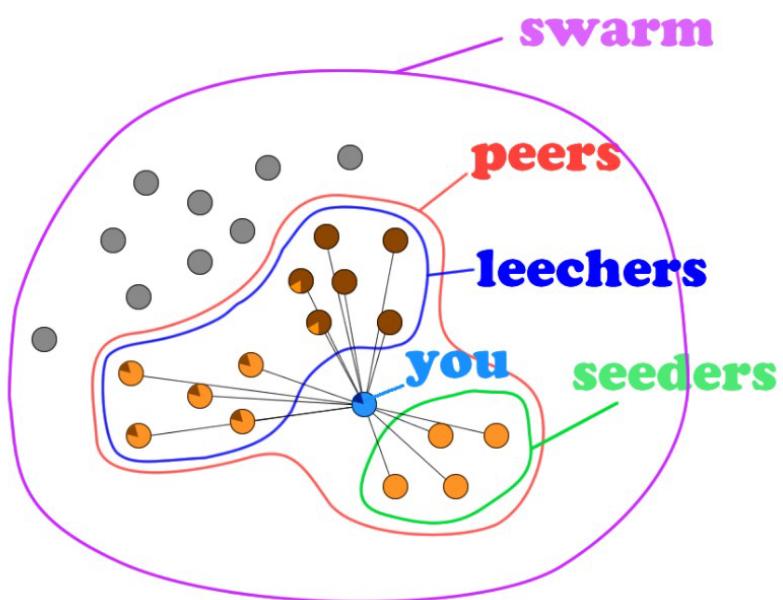


Figure 4.1: Swarm schema

4.1.1 Glossary

- ◊ **tracker**: active entity which coordinates the peers sharing the file, taking trace of who is currently providing the content

- Joe connects to the tracker announcing the content
- the tracker now knows Joe is providing the file
- ◊ **.torrent** a descriptor of the file to be published on a server, which includes a reference to a tracker
- ◊ **swarm** set of peers collaborating to the distribution of the same file coordinated by the same tracker
- ◊ **seeder** peer which owns all the parts of the file
- ◊ **leecher** peer which has some part or no part of the file and downloads the file from the seeders and/or from other lechers.

4.1.2 Protocol Overview

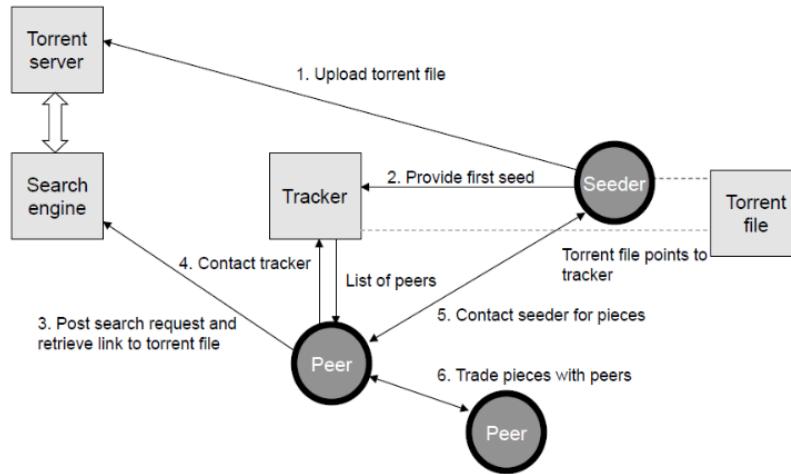


Figure 4.2: BitTorrent protocol overview
BitTorrent protocol is built on top of HTTP

- Seeder*
1. Upload the .torrent on a Torrent Server
 2. Opens a connection to the Tracker and informs it of its own existence: for the moment, it is the only peer which owns the file
 3. Retrieves the file descriptor (.torrent) and opens it through the BitTorrent client
 4. Opens a connection to the tracker and informs it of its own existence and receives from the tracker a list of peers of the swarm
 5. Opens a set of connections with other peers of the swarm.

Objects are serialized in **Bencode**, which is —not popular as JSON— used only in torrent; provides 4 data types: String, Integer, Lists and Dictionaries.

Content is split into chunks called pieces (256KB - 2MB): when a peer receives a piece, it becomes the seeder of that piece.

There is a SHA-1 hash per piece stored in the .torrent file, used to check the piece once it is fully downloaded, allowing to require retransmission in case the check fails.

Pieces size got adapted to have a reasonably small .torrent file

Pieces are then split in **subpieces (blocks)** of 16KB, with each one downloadable from a different peer, optimizing the bandwidth and allowing *pipelining*, decreasing the overall download time.

Trackers keep a database of swarms identified by torrent hash, and knows also the state of each peer in each swarm. In the last versions, **trackerless** BitTorrent uses *Kademlia DHT* to avoid the centralization point of the tracker.

4.2 Pieces selection

The order in which pieces are selected by different peers is critical for good performance, to avoid making peers end up stuck with the same pieces.

- Peers*
- ◊ **Strict Priority**
Complete the “assembling” of a piece before asking for another piece
 - ◊ **Rarest First**
Download the rarest pieces first
 - ◊ **Random First Piece**
Choose a random piece —only— in the bootstrap phase
 - ◊ **Endgame**
When the file download is almost terminated, the remaining pieces are required in *parallel* to all peers who own them. This policy is executed for a small period of time

4.2.1 Free Riders

Free riders in BitTorrent are peers that do not put their bandwidth at disposal of the community. Several non official BitTorrent clients enable the user to limit the upload bandwidth as they like.

However, an approach to solve this problem is based on **reciprocity**, allowing a client to obtain a good service if and only if it gives a good service to the community, by exploiting a dynamic strategy based on connection monitoring called “Tit for Tat”, implemented using **choking**:

choking means *temporarily* refusing to upload to another peer, but still downloading from them; the principle is to upload to peers who have uploaded to us.

Choking

The local peer can receive data from a remote peer if

- ◊ The local peer is *interested* in the remote peer
- ◊ The remote peer *unchoke* the local peer

Choking only peers that upload the most to the local peers would lead to ignoring peers that recently join the network and to the lack of discovery of connections actually better than the used ones.

To avoid this, BitTorrent uses **optimistic unchoking**, i.e. one random peer is being unchoked.

Then, every 30s an interested and choked peer is selected at random **planned optimistic unchoke** (POU), and if this new connection turns out to be better than one of the existing unchoked connections, it will replace it.

In case a peer is choked by everyone, it follows an **anti-snubbing** policy, by increasing the number of simultaneous optimistic unchoke to more than one.

For *seeders* this schema does clearly not apply, since they do not have to download anything; hence they use a different choking algorithm: unchoke peers with the highest upload rate, ensuring that pieces get uploaded and replicated faster.

4.3 DHT and BitTorrent

Kademlia is the protocol used by the largest public DHTs. BitTorrent Inc. introduces its own DHT, called *Mainline DHT*. With respect to Kademlia there are some improvements concerning

- ◊ Routing table management
- ◊ Look-up

The main purpose of Mainline DHT is to provide a “trackerless” peer discovery mechanism to locate peers belonging to a swarm.

Chapter 5

Blockchain

The basic concepts concerning Blockchains are

- ◊ *Ledger*
- ◊ *Consensus* in a distributed environment
- ◊ Tamper freeness
- ◊ Proof of ownership
- ◊ Permissioned and permissionless blockchains

Each **block** is made up of *Data*, *Hash* and the *Hash of the previous block*.

Tamper freeness refers to changing one hash causes changing the hash of the following blocks, implying not only to recompute some hashes, but also to find a value that combined with the new hash solves the *Proof of Work*.

A **ledger** acts like a notary, and is replicated on each node of a P2P network, it is immutable and benefits of the tamper freeness property.

The ledger is like a bullettin storing operations and their order. It must be an **append-only** list of events, and also **tamper-proof**.

If a ledger is organized as a list of blocks, we call it a **blockchain**.

Consensus is the mechanism which defines who decides which operation will be added to the blockchain, and which operation among those to be confirmed will be added.