

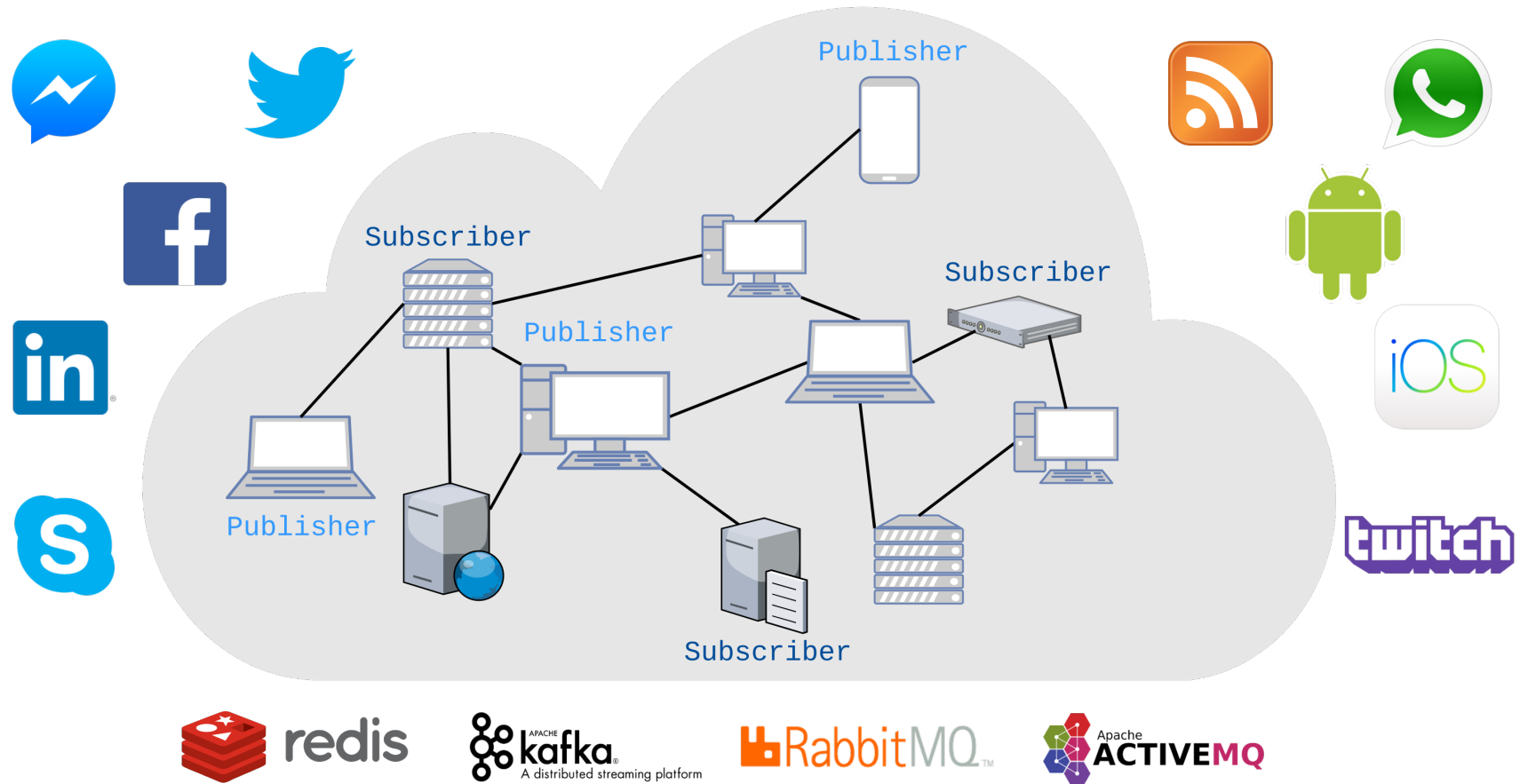
PulsarCast

Scaling Pub-Sub over the distributed web

João Gonalo da Silva Antunes
Instituto Superior Tcnico
INESC-ID

Motivation

World Wide Web



Pub-Sub Paradigm

- Communication paradigm providing **full decoupling** in:
 - Time
 - Space
 - Synchronisation

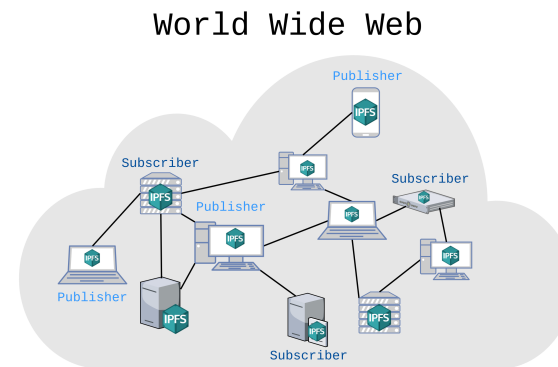
Problems

Lack of a pub-sub system that:

- Scales **for the web**
- Provides reliability:
 - Delivery guarantees
 - Data persistence

Objectives

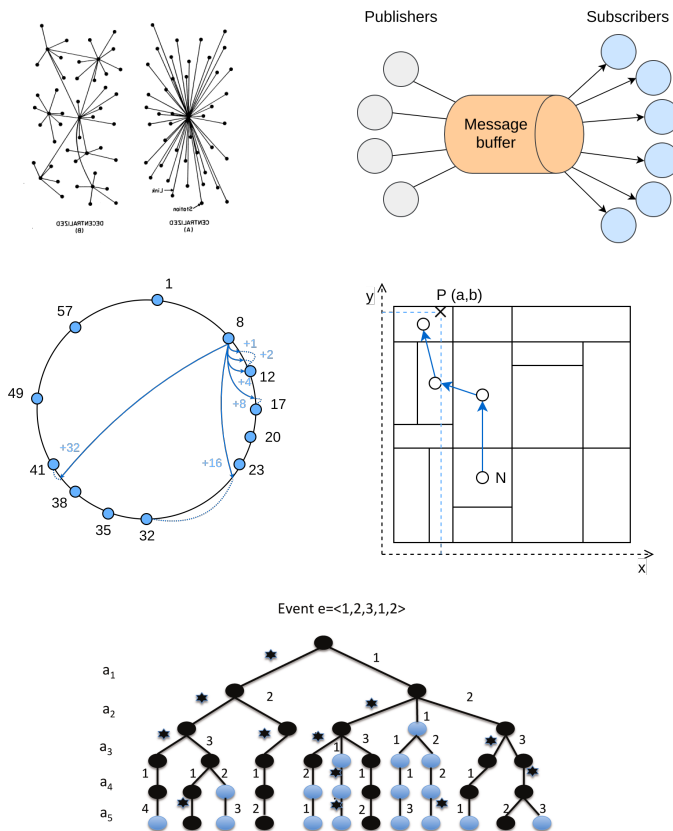
- Decentralised pub-sub architecture using IPFS¹:
 - Highly scalable
 - Reliable
 - Assures persistence
- Develop a system that meets these requirements
- Evaluate its performance on multiple environments



[1] - <https://ipfs.io/>

Related Work

Pub-Sub Systems



Web Technologies



Design Dimensions

- Subscription model
- Network architecture
- Overlay structure
- Subscription management
- Event dissemination

Subscription Model

- Topic based
 - E.g. Redis², Scribe, Tera, Poldercast
- Type based
 - E.g. Hermes
- Content based
 - E.g. Gryphon, Siena, Meghdoot

[2] - <https://redis.io/>

Topic Based

- ✓ A simple notion of group
- ✓ Members of the group receive every message
- ✓ Can build a complex hierarchy
- ✗ Lacks expressiveness

- **Note:** Type based can be seen as a special case of the topic based model fit for type based languages

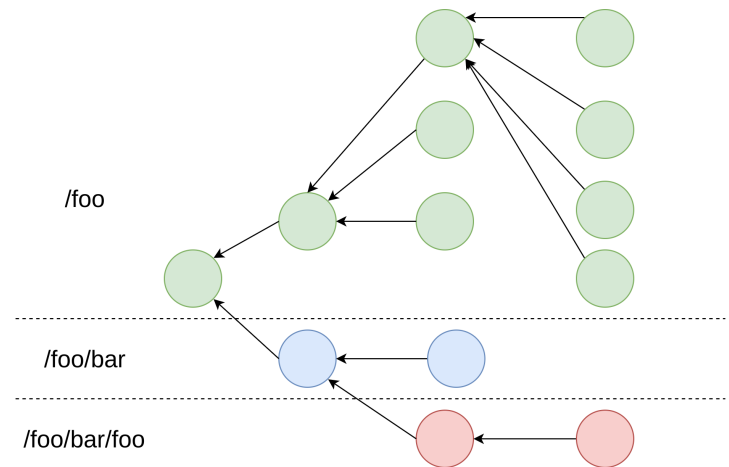


Fig.1: Events on a topic hierarchy

Content Based

- ✓ Filter events based on their content
- ✓ Support for rich and flexible subscriptions
- ✗ Requires complex filtering and message forwarding

Message

```
{  
  exchange: "Euronext Lisboa",  
  company: "CTT",  
  order: "buy",  
  number: "100",  
  price: "5.55",  
}
```

Subscription

```
{  
  exchange: "Euronext Lisboa",  
  order: "buy",  
  number: ">50",  
  price: "<10",  
}
```

Fig.2: JSON subscription example

Network Architecture

- **Centralised:**
 - E.g. Redis, Kafka³, Gryphon
- **Decentralised:**
 - E.g. Scribe, Meghdoot, Poldercast

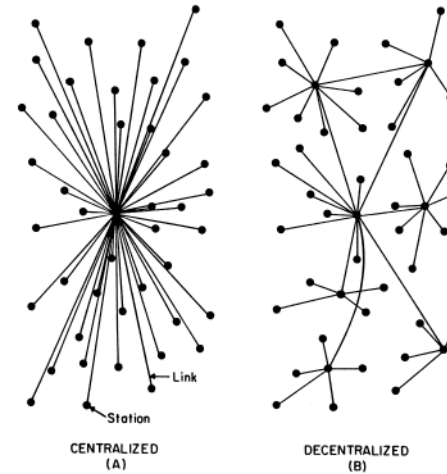


Fig.3: Network architecture overview

[3] - <http://kafka.apache.org/documentation/#introduction>

Centralised

- Focus:
 - ✓ Reliability (with replication)
 - ✓ Consistency
- Lacks:
 - ✗ Scalability
 - ✗ Data Throughput

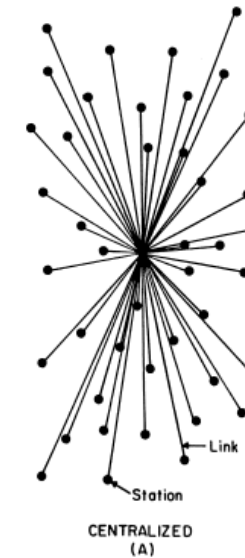


Fig.4: Centralised architecture example

Decentralised

- Focus:
 - ✓ Scalability
 - ✓ Data Throughput
- Lacks:
 - ✗ Reliability
 - ✗ Consistency
 - ✗ Persistence

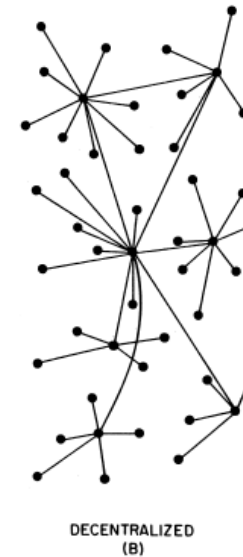
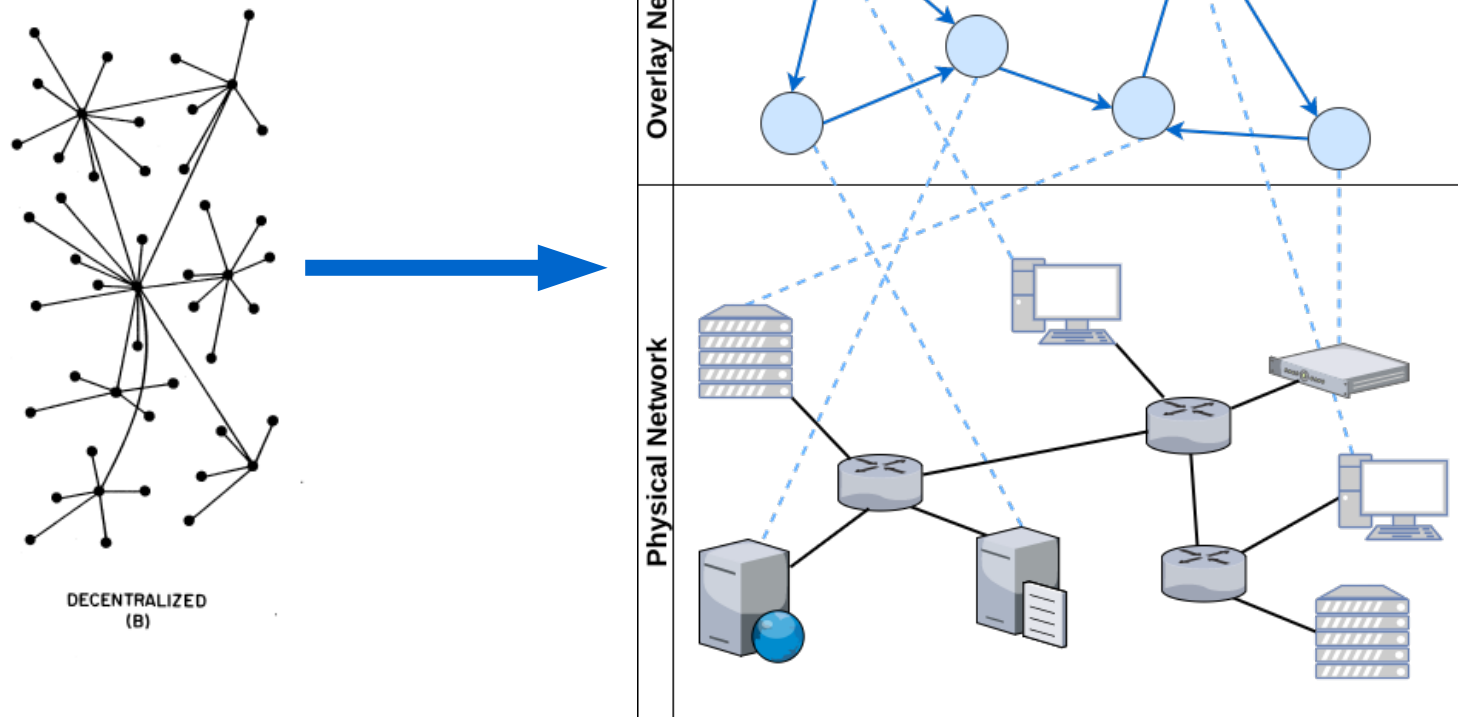


Fig.5: Decentralised architecture example

Network Overlay



Unstructured Overlay

- Each peer connects to a subset of nodes
- No clear structure or hierarchy
- Usage of gossip based membership protocols (e.g. Cyclon[6]) to help preserve:
 - Network diameter
 - Average degree

Structured Overlay

- Peers organise according to a specific structure (e.g. ring, tree, multi-dimensional space).
- Common approach is to use a Distributed Hash Table
 - E.g. Chord, Kademlia, CAN

Structured Overlay - DHT

- Peer identifiers evenly spread across key space
- Ensures the content is evenly distributed
- Queries usually solved in logarithmic time

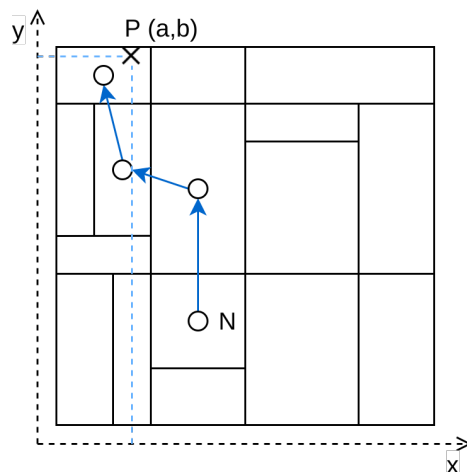


Fig.6: 2 dimensional CAN routing

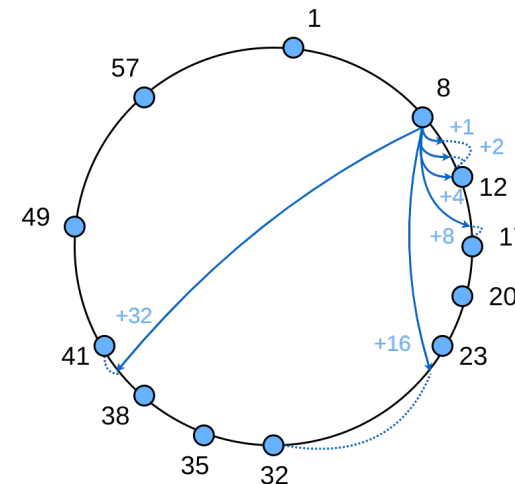


Fig.7: Chord ring

Subscription Management

Event Dissemination

Matching?

Distributed!

Ideally, evenly across the network

System Analysis

	Subscription Model	Architecture	Overlay Structure	Subscription Management	Event Dissemination
Gryphon	Content	(C) broker hierarchy	N/A	Each broker responsible for a subscription scheme	Tree hierarchy
Siena	Content	(C) broker mesh	N/A	Keep state at each node for subscription routes	Flood with cached state
Jedi	Content	(C) broker hierarchy	N/A	Keep state at each node for subscription routes	Tree hierarchy
Bayeux	Topic	Decentralised	Tapestry DHT	Rendezvous node	Multicast tree
Scribe	Topic	Decentralised	Pastry DHT	Rendezvous node	Multicast tree
Meghdoot	Content	Decentralised	CAN DHT	Point in CAN DHT	CAN routing
Hermes	Type	Decentralised	Pastry DHT	Rendezvous node	Multicast tree
Tera	Topic	Decentralised	Gossip based overlay	Unstructured overlay per topic	Random walks & flooding
Mercury	Content	Decentralised	Ring based DHTs	Overlay per attribute in schema	Route through ring overlays
Sub-2-Sub	Content	Decentralised	Ring based DHT & gossip overlay	Clustering of similar subscriptions	Gossip & ring overlay routing
Poldercast	Topic	Decentralised	Ring based DHT, Vicinity & Cyclon	Ring overlay per topic	Ring overlay routing

(C) – Centralised

System Analysis

	Relay Free Routing	Delivery Guarantees	Fault Tolerance
Gryphon	N/A	yes	Best effort
Siena	N/A	yes	Best effort
Jedi	N/A	yes	Best effort
Bayeux	No	yes	Best effort, no subscription persistence
Scribe	No	yes	Best effort, no subscription persistence
Meghdoot	No	yes	Replicated subscriptions
Hermes	No	yes	Best effort
Tera	No	no	Best effort
Mercury	No	yes	Best effort
Sub-2-Sub	No	no	Best effort
Poldercast	Yes	Yes (every publisher is a subscriber)	High resilience to churn, no subscription persistence

Delivery Guarantees – Event delivery guarantees under normal network conditions

Fault Tolerance – Mechanisms to guarantee successful event delivery and subscription matching under heavy churn

Web Technologies

- The Javascript ecosystem
- New network protocols that facilitate P2P communication
- New P2P applications that leverage all of this
 - E.g. IPFS

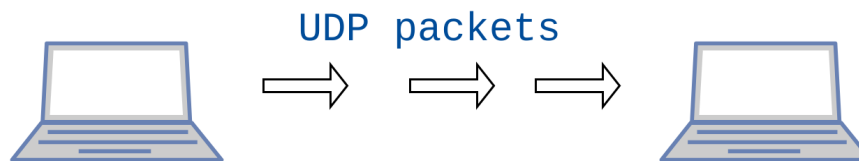
The Javascript Ecosystem

- **Javascript** is one of the main programming languages for the **web**
- Thanks to **NodeJS**, now **runs in servers** also
- Its package manager (**NPM**), is one of the largest package registries in the world
- NPM powers a UNIX-like culture of small modules that work well together



New Protocols

- **WebRTC** made possible full-duplex communication between browsers without an intermediary.
- **QUIC** and **uTP** provided alternatives to TCP, bringing **reliability** and **order delivery** over **UDP**



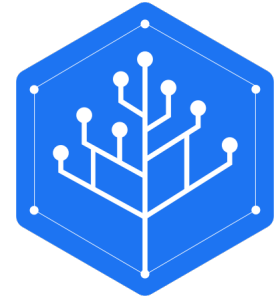
IPFS

- A **P2P hypermedia** protocol designed to create a **persistent, content-addressable** network for the web
- Uses a modular approach through **libp2p** to solve common challenges of P2P applications
- At its core, IPFS uses a Merkle DAG

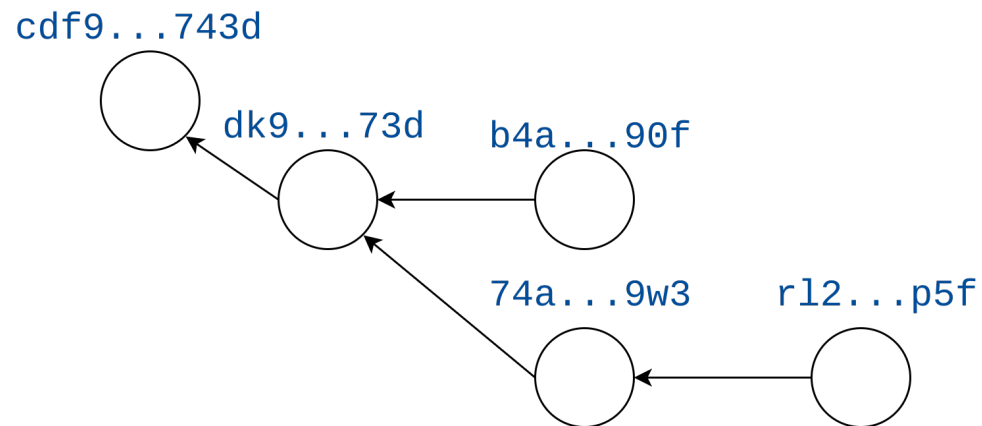


Merkle DAG

- A graph structure used to store and represent **linked data**.
- Each node can be linked to based on the **hash of its content**.
- Referred to as **IPLD** in the IPFS ecosystem.
- Offers **immutability**.

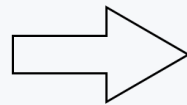


Merkle DAG

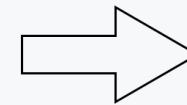


```

{
  "content": {
    "key": "value"
  },
  "merkle-link": "cdf9...743d"
}
  
```



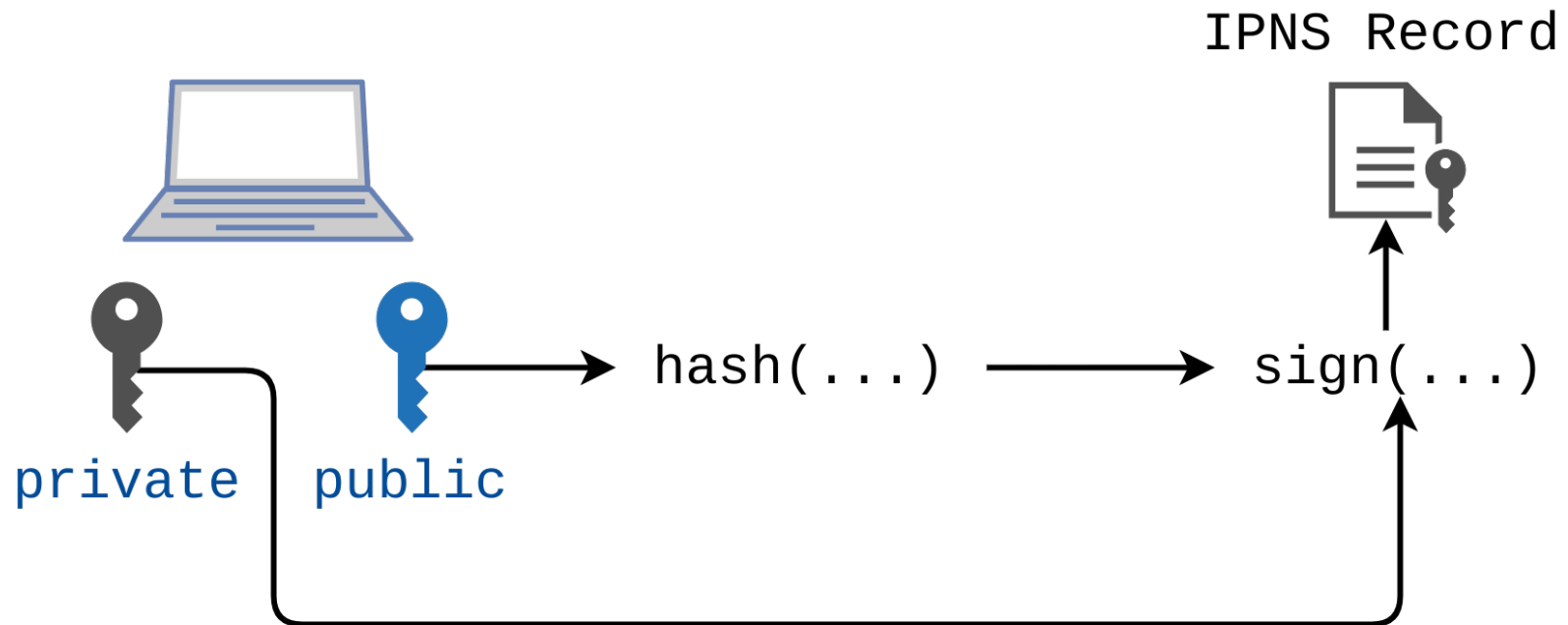
Hash(...)



dk9...73d

IPNS

- A way to offer **mutability**.



Proposed Solution

- A topic based, pub-sub module for the IPFS ecosystem that provides:
 - Data persistence
 - Delivery guarantees
 - High scalability

Architecture

- Take advantage of the already existent Kademlia DHT used by IPFS.

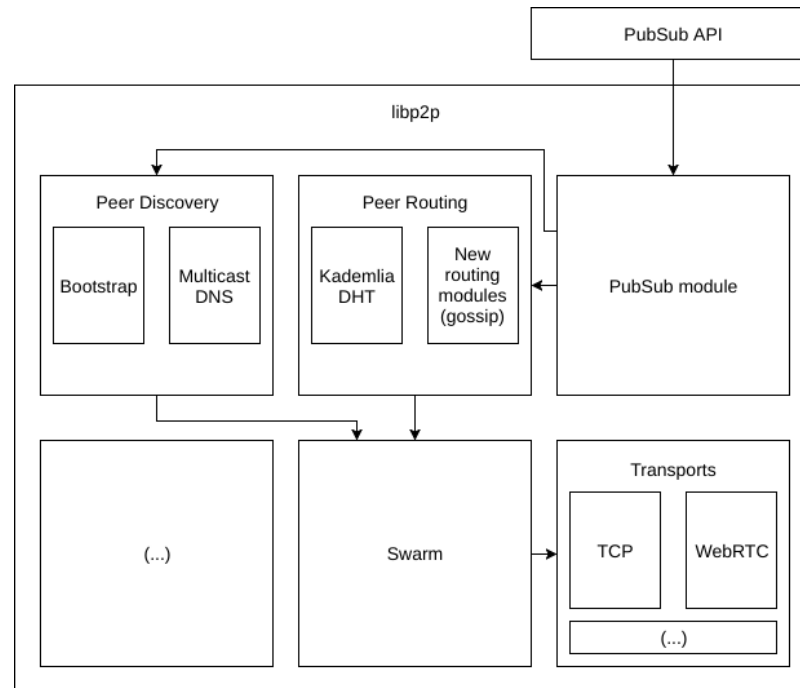


Fig.8: libp2p architecture

Libp2p Kademlia DHT

- put: insert a value with a given key.
- get: get a value of a given key.
- findPeer: find the peer with the given peer ID.
- **getClosestPeers**: find the k closest peers to a given key.
- **provide**: let the network know that this peer can also distribute a given key.
- **findProviders**: find providers for a given key.

Distributed Data Structures

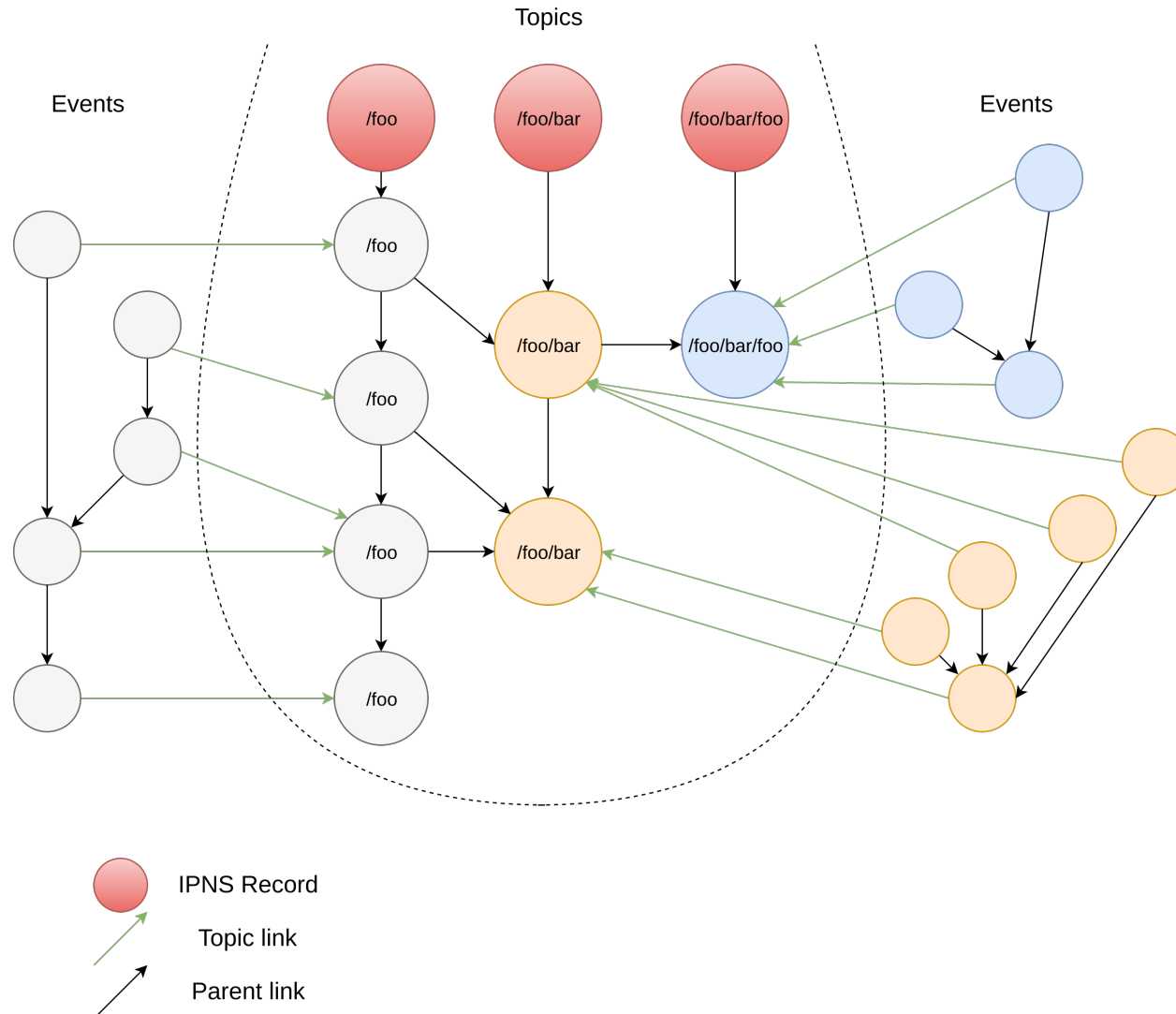
Topic Descriptor

```
{
  name: <topic-name>,
  metadata: <json-object>, // creation date, protocol version, etc.
  #: {
    <sub-topic-name>: <merkle-link to topic descriptor>,
    ...
  },
  parent: <merkle-link to topic descriptor>
}
```

Event Descriptor

```
{
  topic: {
    name: <topic-name>, // Name of the topic
    link: <merkle-link> // Link to the topic of this message
  },
  publisher: <publisher node ID>
  parent: <merkle-link to previous event>,
  metadata: <json-object>, // Timestamp and other relevant info
  payload: <json-object>, // The actual message content
}
```

Distributed Data Structures



Subscription Management

- Using the Kademlia DHT, build a tree structure when creating a new subscription.
-

Event Dissemination

- (mention matching)

Quality of Service

Evaluation

Conclusion

Thank you for your presence!

Questions?

References

Topic Descriptor

```
{  
  name: <topic-name>,  
  metadata: <json-object>, // creation date, protocol version, etc.  
  #: {  
    <sub-topic-name>: <merkle-link to topic descriptor>,  
    ...  
  },  
  parent: <merkle-link to topic descriptor>  
}
```

Event Descriptor

```
{
  topic: {
    name: <topic-name>, // Name of the topic
    link: <merkle-link> // Link to the topic of this message
  },
  publisher: <publisher node ID>
  parent: <merkle-link to previous event>,
  metadata: <json-object>, // Timestamp and other relevant info
  payload: <json-object>, // The actual message content
}
```