



## Peer-to-Peer Systems and Security IN2194

### Chapter 1 Peer-to-Peer Systems 1.4 Maintenance, Optimization, and Technical Issues

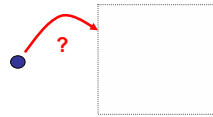
Prof. Dr.-Ing. Georg Carle  
Dipl.-Inform. Heiko Niedermayer



- Bootstrapping
  - How to find and join the Peer-to-Peer overlay?
- Maintenance
  - Churn / Dynamics of Peer-to-Peer networks
  - Measurements
  - Mobility
- Optimization
  - Locality
  - Load
- Technical Issues
  - NAT / Firewall Problems
  - Legacy Applications
- Conclusions



## How to find the Peer-to-Peer overlay?



### Problem

- How do I contact a node in the network?
- How can this be done without a server or a single server?

### Necessary

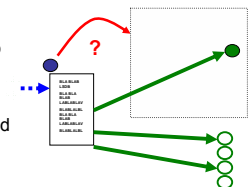
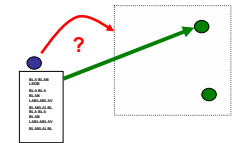
- Some knowledge
  - direct knowledge („You already know it“)
  - indirection via some rendezvous point („Ask someone“)
- Important
  - *We cannot find a service we do not know anything of!*
  - However, more semantically-oriented rendezvous mechanisms might exist in the future (e.g. „give me a news-exchange network for Munich“), but they are not common today.



## How to find the Peer-to-Peer overlay?

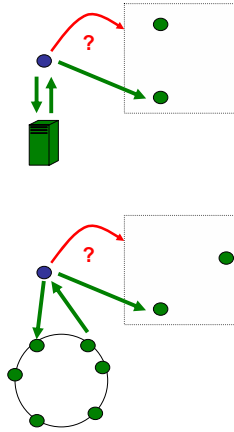
### Solutions

- List with fixed (rendezvous) nodes
  - Method
    - Contact one of the nodes on the list
  - Assumption:
    - There exists a fixed set of nodes that are always accessible, now and in the future.
  - The approach is efficient and simple.
  - List or parts of the list may be updatable
    - e.g. by the user via the website
- Cache with nodes
  - Method
    - Store nodes you have seen. Contact some of them in order to join.
  - Assumption
    - Many nodes stay long in the network and can still be contacted the next time.
  - Cannot solve the initial contact problem (→ no cache yet).
  - Many nodes are stored anyway during runtime for robustness. Problem, what are the long-lived nodes?



## How to find the Peer-to-Peer overlay?

- Acquire node list via some client/server service like HTTP
  - Mechanism
    - Ask a server for a list with current rendezvous nodes in the network.
  - Assumption
    - There is a server that has a knowledge of currently active (rendezvous) nodes in the system. It needs to get either updated regularly or scan the P2P network.
  - For large networks no need to scan all nodes.
  - External nodes get to know network status information
    - Privacy issues, useful information for attacks, etc.
- Acquire node list via some `_other_` Peer-to-Peer mechanism
  - Mechanism
    - Like above, but using a Peer-to-Peer system.
- Ask user
  
- While no individual solution is perfect, a combination of the mechanisms presented should work for most cases.



## Special case: local Peer-to-Peer networks

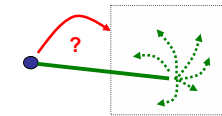
### Motivation

- Example: Friends sit in a park or classroom. They want to form a Peer-to-Peer network.
- In case of wireless networks such networks are called mobile ad-hoc networks (MANETs). Bluetooth is an example.



### Solution

- Use Broadcast or Multicast mechanism of lower layers to find other peers.



## Bootstrapping examples in existing P2P systems

### Edonkey / eMule

- Via node lists (called „server list“, „server.met“ files)
  - Downloaded from somewhere or stored on hard disk
- Via user input / Kad also uses a node cache.

### BitTorrent

- Problem is to find a torrent / tracker with the file of interest.
  - No automatic solution.
- Tracker knows active nodes.
  - Tracker known from .torrent file.

### Skype

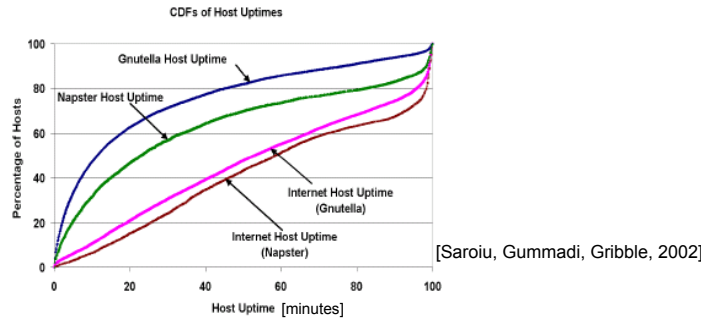
- Host cache: list of cached super peers
- Fixed set of bootstrap super peers (7 in 2004) – in the host cache after installation

## Maintenance

### Maintenance

- Problem
  - Joins and leaves (= „Churn“)
  - Changes in network connectivity
- Goals
  - Keep the network alive
  - Keep the data stored
  - Provide a good quality of service
- Solutions
  - Know more nodes
    - Successor/neighbor lists to recover.
    - Buckets instead of single links.
  - Be up-to-date
    - Frequently check if nodes still exist and routing table is correct.
  - Store data on nodes of a replica set (= group of nodes)
  - Use networks with low join/leave overhead if churn is high.

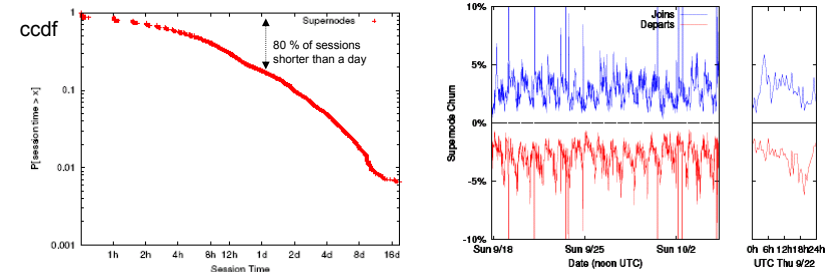
## Churn – Filesharing



### Filesharing

- The Internet Host Uptime figures show how long the host is online. The Gnutella and Napster uptime are how long the application is running on the host.
  - Obvious: Gnutella/Napster Uptime  $\leq$  Host Uptime
- Most nodes are only available for a short time. However, some are online for several days.
- A measurement by Chu, Labonte and Levine reports that short short sessions are even more likely. They speculate about the session time being an overlap of Zipf distributions (heavy-tail / power law distribution).

## Churn – Skype



[Guha, Daswani, Jain, IPTPS 2006]

### Skype

- Super nodes
  - Relatively stable network (session time less than a day does not imply that the node won't be back).
  - Possible approx: Poisson arrival, session time heavy-tail
  - Variations over daytime and weekday.
  - 95 % of super nodes will be there 30 min later

## Mobility

### Mobility

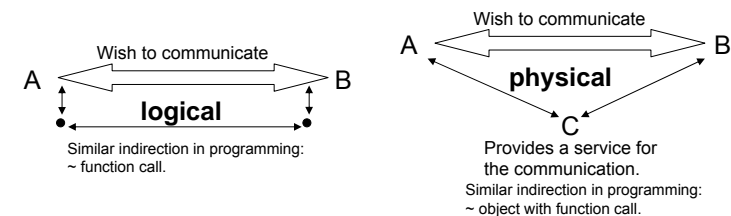
- Users are mobile.
  - Geographic location.
  - Network access changes.
  - Used computing device changes.
  - Consequences
    - Network address changes (IP:Port).
    - Disruptions in physical connection.
- Problem
  - Maintain connectivity / avoid interruption in higher layer communications.
  - Find someone / some service.
  - User contact information has to be updated.
  - Node contact information has to be updated.

### Solution

- All mobility solutions use some kind of indirection.

## Indirection

„Any problem in computer science can be solved with an additional layer of indirection.“  
(David Wheeler) - famous quote - but Wheeler continues this sentence with  
„But that usually will create another problem.“(David Wheeler)



### Indirection

- A and B want to communicate, they use C as intermediate point.
- Indirection usually adds a protocol and a layer in the protocol stack.
- Points of Indirection
  - Logical
    - Use primitives of more abstract protocol (heiko@net instead of EF:40:45:2A:44:55)
  - Physical
    - Use another physical entity that provides the service for the additional abstraction.

## Mobility Support



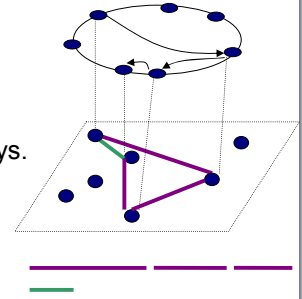
### Mobility Support

- P2P network stores current location of users under their ID
  - Update locator (IP:Port) when locator changes (→ newIP:Port).
  - Inform currently active communication partners.
- Mobility for peers in the network
  - Graceful leave and rejoin.
  - Alternative: If all neighbors can be contacted, update their information.
- Softening the handover
  - A change of locator usually causes a moment of being unreachable.
  - Softening
    - If communication is via indirection point, then buffer some packets.
    - Try to get new locator and update before old locator is down.

## Optimization – Locality

### Locality

- Basic Peer-to-Peer network construction ignores geographic or underlay relationships.
  - Long path with hops all over the world may be used to contact a local neighbor. *Inefficient!*



### Stretch

- Metric for locality in and performance of overlays.
- Idea: Compare distance over overlay with distance in the underlay.

$$\text{Stretch} = \frac{\text{DISTANCE}_{\text{via\_overlay}}}{\text{DISTANCE}_{\text{in\_underlay}}}$$

- Any distance metric can be used, latency (default) or IP hops are common.

## Locality – Approaches

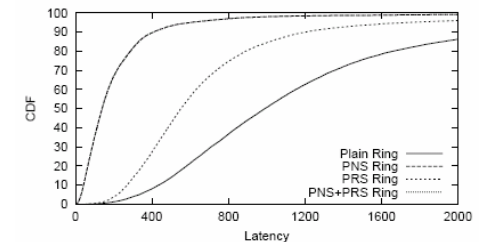
### Basic approaches to improve locality

- Proximity Neighbor Selection (PNS)
  - Given multiple potential nodes for a routing table entry, select the best from the candidates according to proximity metric (e.g. latency, RTT).
- Proximity Route Selection (PRS)
  - Use routing table as candidate set and select next hop as trade-off between latency and reduction of distance in ID space.
  - Alternative understanding of PRS (almost identical to PNS): Given multiple routing entries into one direction, chose the next hop with best proximity metric.
- Proximity Identifier Selection (PIS)
  - Select identifier in a way to minimize proximity to neighbors in overlay.
    - Internet Coordinate Systems are examples.
- Special solutions
  - Pastry: Neighbor Set with local nodes, join via local node
  - Topplus: Use IP addresses as overlay IDs.

## Locality – Comparison & Discussion

### General Results

- PNS better than PRS
  - Reason: PNS uses #candidates \* size(routing table) nodes to optimize. PRS only uses size(routing table) nodes.
- PNS+PRS only slightly better than PNS, in Figure on the right the two curves are indistinguishable.



Gummadi et al, Sigcomm 2003

### Discussion

- Optimizing locality may conflict with the idea of diversity, i.e. to improve robustness by relying not only on geographically local nodes, but on internationally distributed ones.

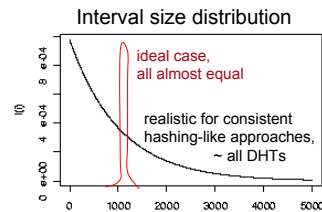
## Optimization – Load Balancing

### What is load?

- Items stored
- Computational effort
- Traffic and work imposed by requests.
- Maintenance traffic

### Problem?

- Peers may be assigned more work due to larger intervals.
  - Low probability of bad luck in randomized scenarios.
- Non-random approaches for work and ID assignment
- Requests may not be uniformly distributed among items. → Zipf's Law



## Zipf's Law

**Zipf's law:** "The popularity of  $i$ th-most popular object is proportional to  $i^{-\alpha}$ ,  $\alpha$ : Zipf coefficient."

- Zipf distribution is Power-Law.
- Zipf-like popularity can be found for websites, words in natural languages, movies.
- However, with static content that is only requested once, there is usually a saturation for popular files and the behaviour not fully Zipf (measurements of Kazaa, Gnutella,...).
- Popularity of items is application-dependent. Not all applications create Zipf-like distribution.

→ Extreme differences in popularity are one argument for the need for load balancing.

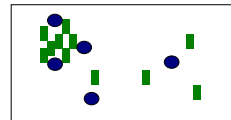
## Necessity of Load Balancing and basic ideas

### Random item IDs

- Load of each node follows binomial distribution with low variance.
  - Variation of the size of ID ranges is dominating.
- Node Balancing is important

### Non-random item IDs

- Need to balance load as items may cluster in certain areas and build hotspots.



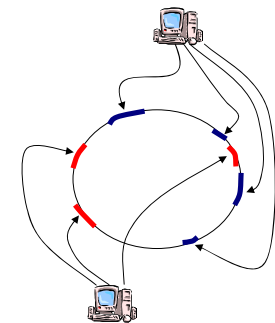
### Basic ideas for load balancing

- Methods to reduce variation in number of items and interval size.
- Share load at hotspots by using nodes in parallel.

## The Virtual Server approach

### Virtual Server

- Each node is represented by  $O(\log n)$  virtual nodes called virtual servers.
  - Having  $O(\log n)$  intervals already averages out some imbalance. Theory of consistent hashing suggests this number (fair with high probability).
- Virtual servers can be transferred from overloaded (heavy) to light nodes
  - The light node does not become heavy.
  - The transferred virtual server is the lightest to make the heavy node light.
  - If no virtual server is heavy enough to make the heavy node light, transfer the heaviest virtual server.
- How do they find each other?
  - Randomly contacting other nodes.
  - Directories with information about heavy and light nodes.



## Simple Efficient Load Balancing

### Simple Efficient Load Balancing (Karger/Ruhl)

- Address-Space Balancing
  - Instead of having  $O(\log n)$  virtual servers, each node has a fixed set of  $O(\log n)$  possible positions, only one of them is active.
  - The active virtual server is selected according to interval size.
    - To avoid the selection of mini-intervals, the „smallest“ address range is selected according to the order  $1 < \frac{1}{2} < \frac{1}{4} < \frac{3}{4} < \frac{1}{8} < \frac{3}{8} < \frac{5}{8} < \dots$
- Item Balancing
  - Nodes randomly connect each other.
  - IF one of them has  $\epsilon$  times less items than the other THEN
    - It moves to the interval of the other node. Both nodes fairly divide the interval so that they share the same number of items.
- Results
  - Approach has provable bounds for address-space and item balancing.
  - Item Balancing needs  $\Omega(\log n)$  random connections for each node per half-life and when the number of items on a node doubles.

## Dealing with Hotspots

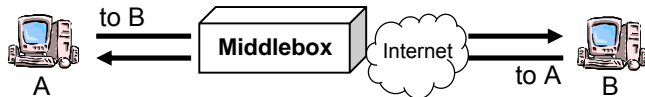
### Hotspots

- Zipf's Law indicates that there may be items with extreme weight, say „60 % of the users want this file“. Such cases cannot be solved with item balancing alone.

### Dealing with Hotspots

- Replica sets
  - A set of nodes is storing an item and any node can answer to queries for reading the item.
  - Basically, replica sets are proposed for fault tolerance, but they can be used to relief hot spots. This only works, however, if nodes in the replica set are on paths towards the item ID and are only hit by a fair share of the queries.
    - Leaf set in Pastry, Predecessors on paths in Kademlia are good choices.
    - Successor list in Chord is a bad choice as they are not hit by queries.
- Structural Replication
  - Allow multiple nodes per key-space. Similar to replica sets, but directly built into the system design, e.g. in the DHT P-Grid (Aberer et. al 2001) or CAN.

## Technical Issues – Middleboxes

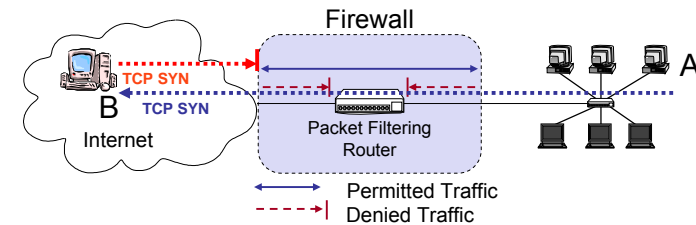


„A middlebox is defined as any intermediary device performing functions other than the normal, standard functions of an IP router on the datagram path between a source host and destination host.“ (RFC 3234, Middleboxes: Taxonomy and Issues, February 2002, [www.ietf.org](http://www.ietf.org))

### Middleboxes

- may prevent hosts on the Internet from connecting to each other
- Firewall
- Network-Address-Translators (NAT)
- ...
- Results for the Kad network [Brunner, 2006]
  - 44 % of peers firewalled or NATed.

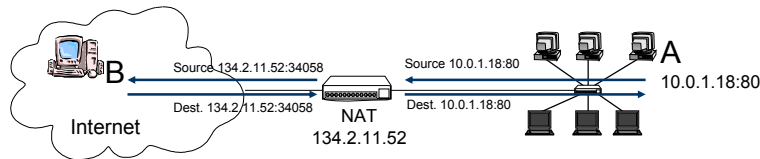
## Firewall



### Firewalls

- Usually at the edges of a local network firewalls are used to filter traffic according to a set of rules.
  - Rules are usually based on IP:Port combinations, e.g. only allow Port 80 (http) traffic to the webserver (IP)
- Typical problem for Peer-to-Peer
  - Connection establishment may only be allowed from inside to the Internet and not from the Internet to the hosts in the local network.

## Network Address Translator (NAT)



### Network Address Translation

- Typical
  - Local computers have private IP addresses (10/8,192.168/24) which are reserved for usage in local networks and which are not routed by routers.
    - computers cannot be addressed with their private IP address
- Address Translation used to translate local private addresses to global public addresses.
  - usually: Local\_IP:Local\_Port <-> Public\_IP:Public\_Port
- Problems
  - Public address different from machine address and application port.
  - NATs usually work dynamically, so application or computer are unreachable before they connect and get a public IP:Port pair.

## Simple solutions

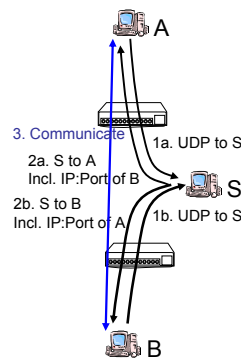
### Simple solutions

- PUSH in Gnutella
  - An existing connection is used to signal the request to the firewalled peer, the firewalled peer then initiates the connection.
- In super peer networks
  - Super peer approach
    - Super peers may not be firewalled or behind NAT.
  - Use super peer as relay
- Use reachable peer as relay
  - Peers need to know peers with public IP addresses.
- Hide behind common ports
  - Many firewalls may allow connections to some common ports (HTTP, SMTP, etc.) and block new connections to unknown ports (typically used by P2P).
  - e.g. Skype tends to listen at port 80 if available.

## Hole Punching (for NAT Traversal)

### UDP Hole Punching

- A and B use a third party S to directly connect to each other.
- Algorithm
  - Let A and B be the two hosts, each in its own private network; N1 and N2 are the two NAT devices; S is a public server with a well-known globally reachable IP address.
  - A and B each begin a UDP conversation with S; the NAT devices N1 and N2 create UDP translation states and assign temporary external port numbers
  - S relays these port numbers back to A and B
  - A and B contact each others' NAT devices directly on the translated ports; the NAT devices use the previously created translation states and send the packets to A and B
- Does not work with all types of NATs, but with the most common ones.
- There is also TCP Hole Punching.



## Standards & Protocols (for NAT Traversal)

### STUN, TURN, ICE

- IETF standards for NAT Traversal
- Based on Hole Punching and/or Relaying
- Transparent approaches – no interaction with the NAT

### UPnP (Universal Plug and Play)

- Developed by Microsoft, very common protocol in home network and consumer devices.
- Interaction with NAT to determine public IP address, configure port mappings, etc.
  - Many current filesharing network use it to open ports if necessary.
- Insecure, not suitable for larger/company networks

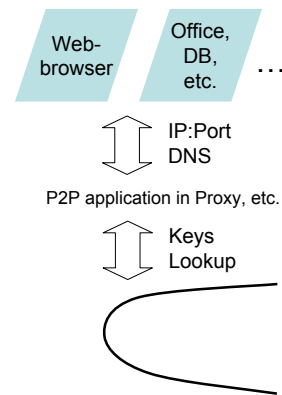


### Legacy Applications

- ❑ Applications that are commonly used and that are not developed for the P2P network.
- ❑ TCP / UDP sockets used to communicate.
- ❑ Resolution of IP addresses via DNS.

### Peer-to-Peer for legacy applications

- ❑ Accept TCP/UDP traffic with IP:Port addressing
- ❑ Intercept DNS to resolve destinations and map IP:Port pairs to keys in the P2P network
  - E.g. based on special non-existing top-level domains, e.g. „schoenerdienst.p2p“
- ❑ Technical solutions
  - Proxy
    - The application sends the packets to the configured proxy. The proxy then modifies and inserts them into the Peer-to-Peer network.
  - Virtual Network Devices (tun/tap)
    - Similar, but Layer 2 approach.



### Conclusion

- ❑ Bootstrapping is not possible without some pre-knowledge or central element.
- ❑ Churn and mobility impose the need to
  - Frequently update neighbor and routing information.
  - Use replication and other mechanisms for redundancy.
- ❑ Locality can be improved. However, it can conflict with diversity.
  - PNS, PRS, PIS
- ❑ Load balancing can be necessary depending on the use case.
  - Virtual nodes / servers.
  - (Structural) Replication
- ❑ Problem of Middleboxes (NAT/FW)
- ❑ Integration of legacy applications.