

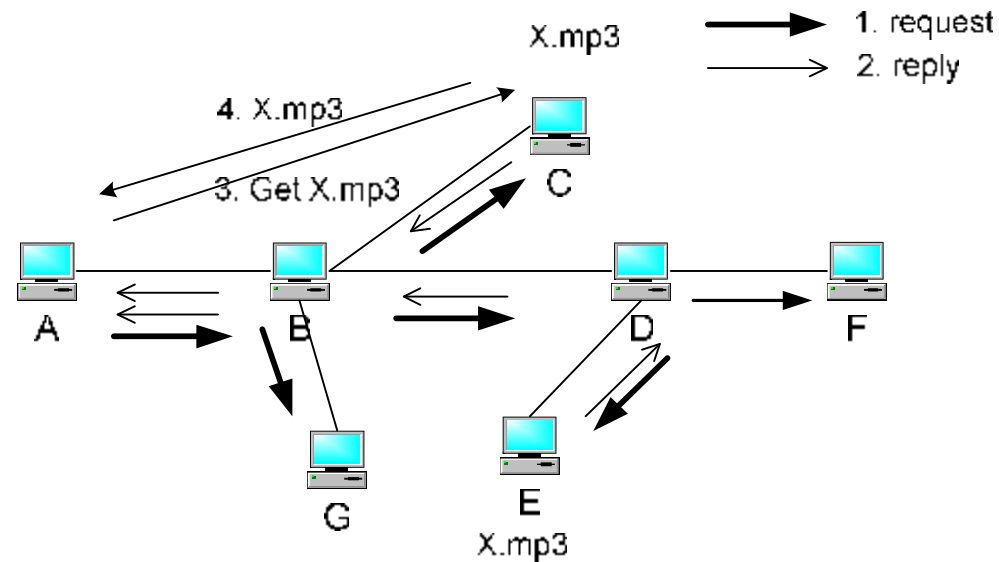


Using of Semantic Friends for Search in P2P Network

Evgeny Linsky

Unstructured P2P Network

- n "Search": performed using flooding
- n Unstructured: neighbors are chosen in arbitrary manner



n Advantages

- n easy leave/join procedure
- n support complex search queries (a*ccc?ff.mp3)

n Disadvantages

- n scales badly, i.e. many nodes --- huge control traffic
- n search: $\sim n$, n --- number of nodes



Problem Statement

- n Goal: optimize flooding using model of requests
- n General ideas
 - n All documents and nodes are divided into thematic classes (by interest)
 - n All documents have different popularity
- n Assumptions
 - n Network is stable: no one joins and leaves the network
 - n Network is fully connected (one-hop network)
 - n Node searches the same query periodically



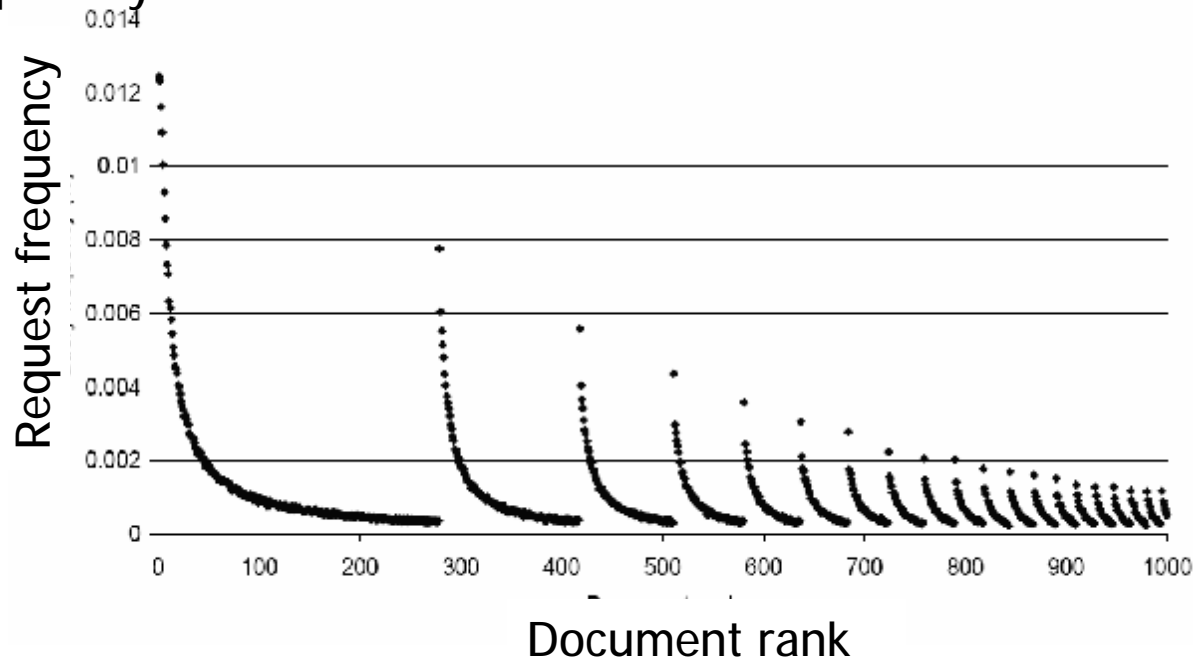
Paper Used as Base

- n “Exploiting semantic proximity in peer-to-peer content searching”
 - n Spyros Voulgaris, Anne-Marie Kermarrec, Laurent Massoulié, Maarten van Steen
- n Contribution
 - n Model of requests
 - n Algorithm for the given model

Model of Requests

n Ideas

- n Documents (d_i) are divided into classes (c_j)
- n Documents and classes have different popularity (Zipf law)
 - n $P(d_i) = 1/I$, $\text{Size}(c_j) \sim 1/j$
- n Every user is interested in one class
 - n $P = \alpha$ (0.8) --- user requests document in his own class
 - n $P = 1 - \alpha$ (0.2) --- user searches for document according to general popularity
- n Popularity is static





Base Algorithm

- n During “warming up” use flooding
 - n If document was found on some node, add this node to semantic friends list
- n When the list of semantic friends is formed
 1. Send request to semantic friends
 2. If this search fails, use flooding
- n List of semantic friends has finite size
- n List updating policy --- modification of LRU



Base Algorithm

n Disadvantages

- n Performance criterion – Hit Ratio of semantic friends list
- n Maximization of Hit Ratio does not have big impact on search delay
 - n Requests to low probability class are rare
 - n Increase of Hit Ratio for this class does not influence on total overheads

n Our goal: minimize average number of transmissions for one served request



Proposed Solution

n FileList Algorithm

- n Node, which successfully served the request, sends FileList in addition to requested file
- n FileLists are stored and used during search in semantic friends

n Advantages

- n Lead to decreasing of average number of transmissions for one served query
- n Adaptable to changes in request model

n Disadvantages

- n Storage and analysis of large data arrays is hard for restricted devices => limitation of the list size



Limited FileList

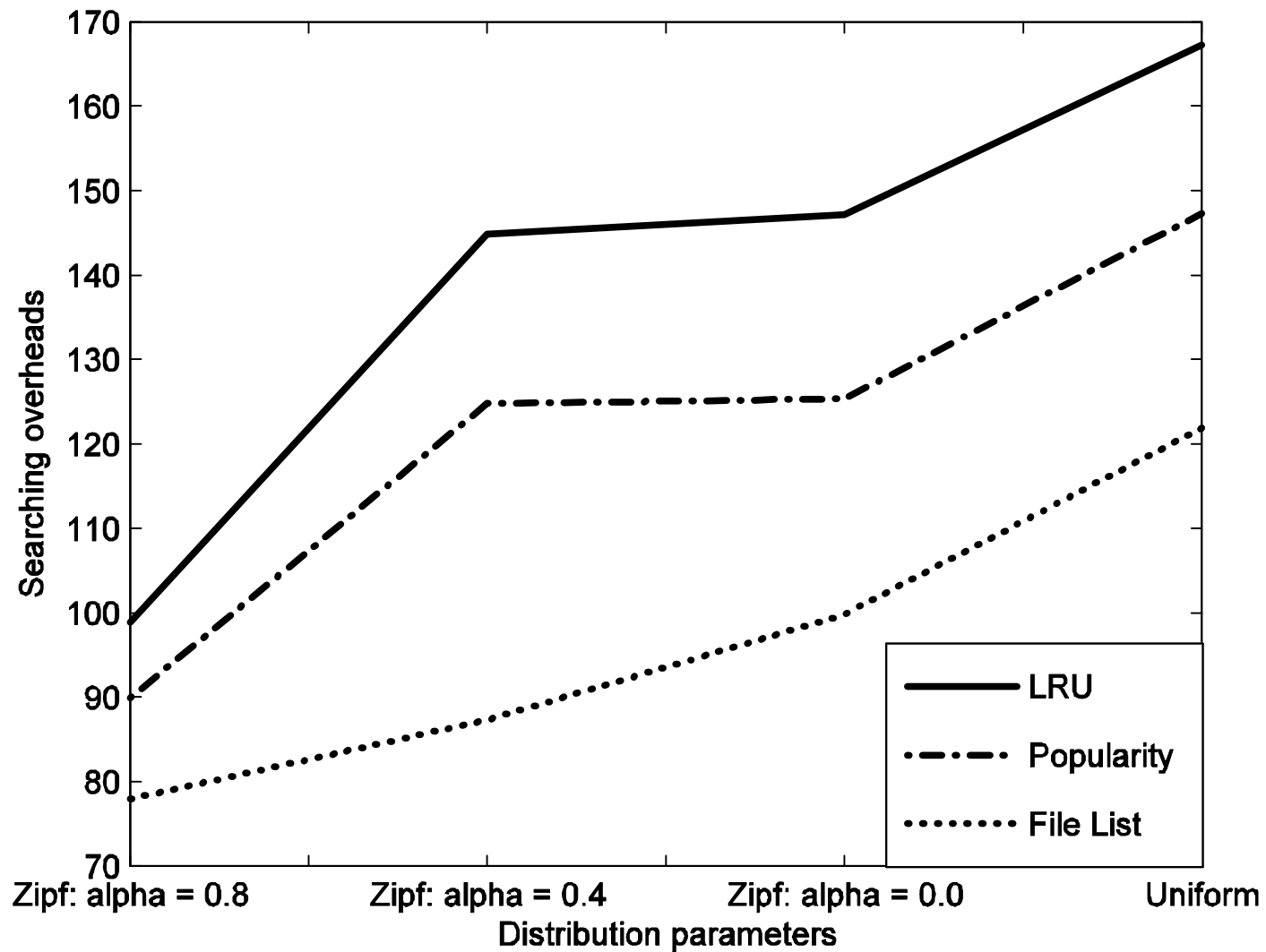
- n Node collects statistics about requests trying to reconstruct request model
 - n Node estimates probability of documents in its own group
- n Semantic friends list management policy
 - n Include node, only if its FileList maximizes the future successful search probability
 - n Probability is calculated using described above estimations

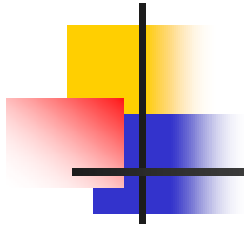


Simulation Results

Algorithm name	Searching overhead
LRU	100
Popularity	90
FileList	78

Simulation Results





Next Step

Service Discovery in Smart Spaces



Ubiquitous Computing

- n Many local wireless networks
 - n including multi-hop
 - n including mobile and infrastructureless
- n Users roam between such networks
- n Goal: find services in local proximity (Resource Discovery)
 - n Closest free printer
 - n Free parking slot
 - n Web-service, which could convert RTF to PDF
 - n Shared collection of mp3 files



Current Solution

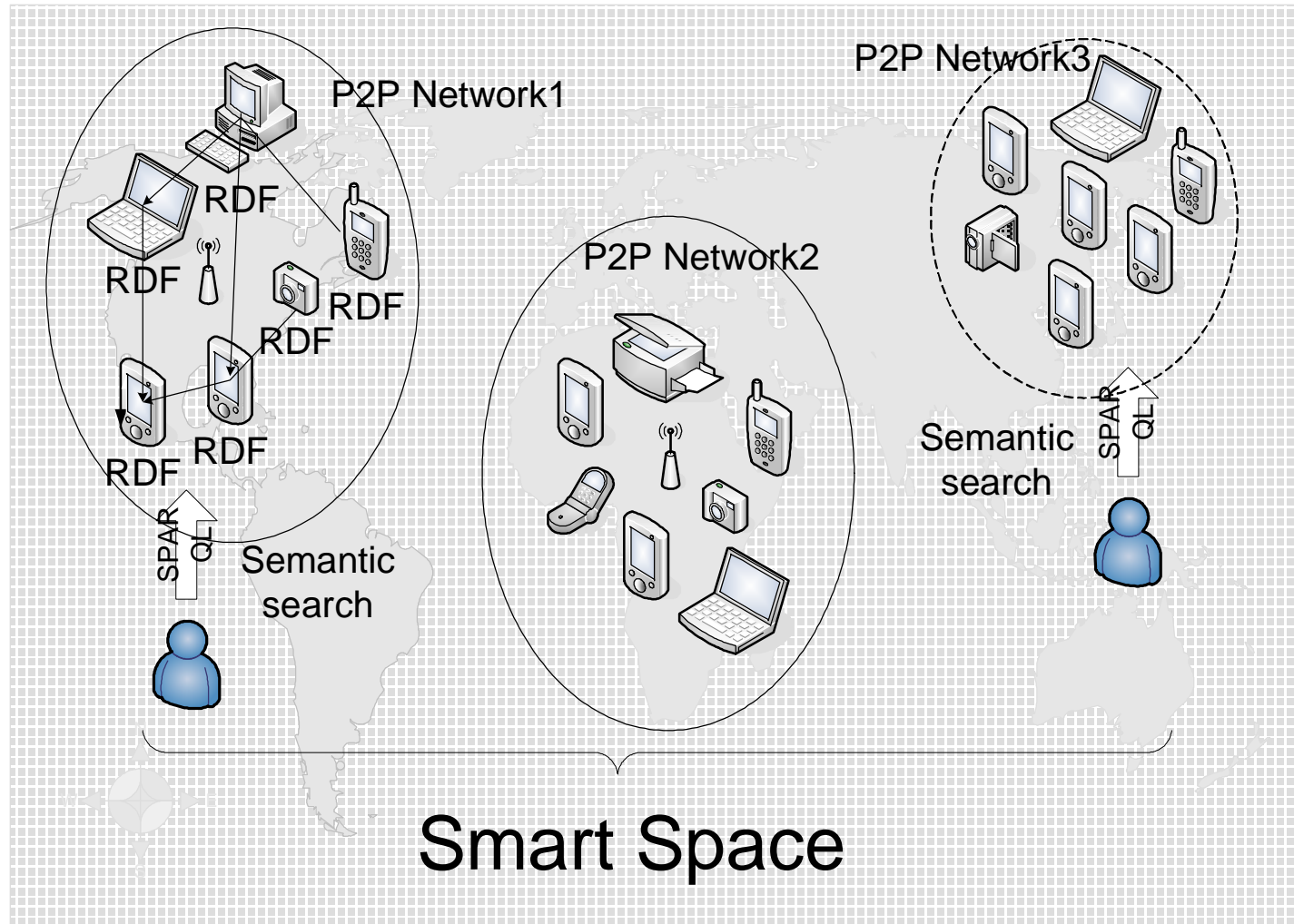
- n Documents and services are described using natural language
- n Dedicated server scans the network indexing the documents
- n Users sends requests to dedicated server
- n Requests are formulated using natural queries
 - n Server uses complex algorithms for finding relevant documents
 - n Results are analyzed manually by human



Approach

- n Documents and services has formal description on RDF language
 - n Simplify relevancy analysis
 - n Results could be interpreted by program
- n No dedicated server
 - n Single point of failure
 - n Additional overheads for maintaining infrastructure
- n Search is implemented as distributed algorithms, e.g. using flooding

Smart Space





Problem Statement

- n Flooding search

- n Large searching overheads
 - n Large delay

- n Optimization idea

- n Directed search

- n Request is rebroadcasted only to selected neighbors
 - n Selection is done using history of previously performed requests

- n Disadvantages

- n Not all possible documents can be find
 - n Tradeoff: search delay vs QoS
 - n $QoS = \frac{\text{Number of Results}}{\text{Number of matching documents}}$