

Towards Next-Generation Peer-to-Peer Systems



Magnus Kolweyh
aka risq
University of Bremen

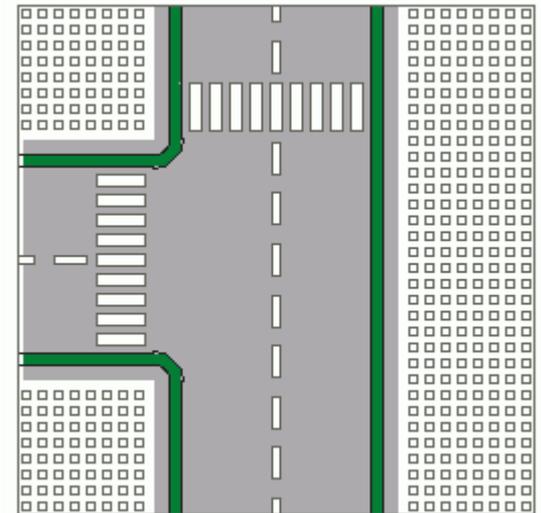
mag@tzi.de

Motivation

- Present an overview of interesting P2P systems
- Offer some knowledge out of P2P science
- Discuss novel implemented approaches and concepts
- Show some current measurement studies
- Inspire developers
- **Ask the usual P2P suspects**

Roadmap

- Science vs. Peer-to-Peer
- P2P generations
- Challenges and concepts
- Current trends in file sharing
- Services for Peer-to-Peer systems
- Example: Data Mining
- Prospects



Science and Peer-to-Peer

...or how to survive as an evil file sharing PhD student...



Edutella - P2P for the Semantic Web

Anthill



Project
JXTA

Don't touch too much file sharing..illegal content



BitTorrent



g²dn

apple_juice

Suprnova dead ? ..ok..next one please

P2P Generations ?!

Old school

Phones
Arpanet
Bulletin Boards

File sharing

Napster
Gnutella
Edonkey

Messenger

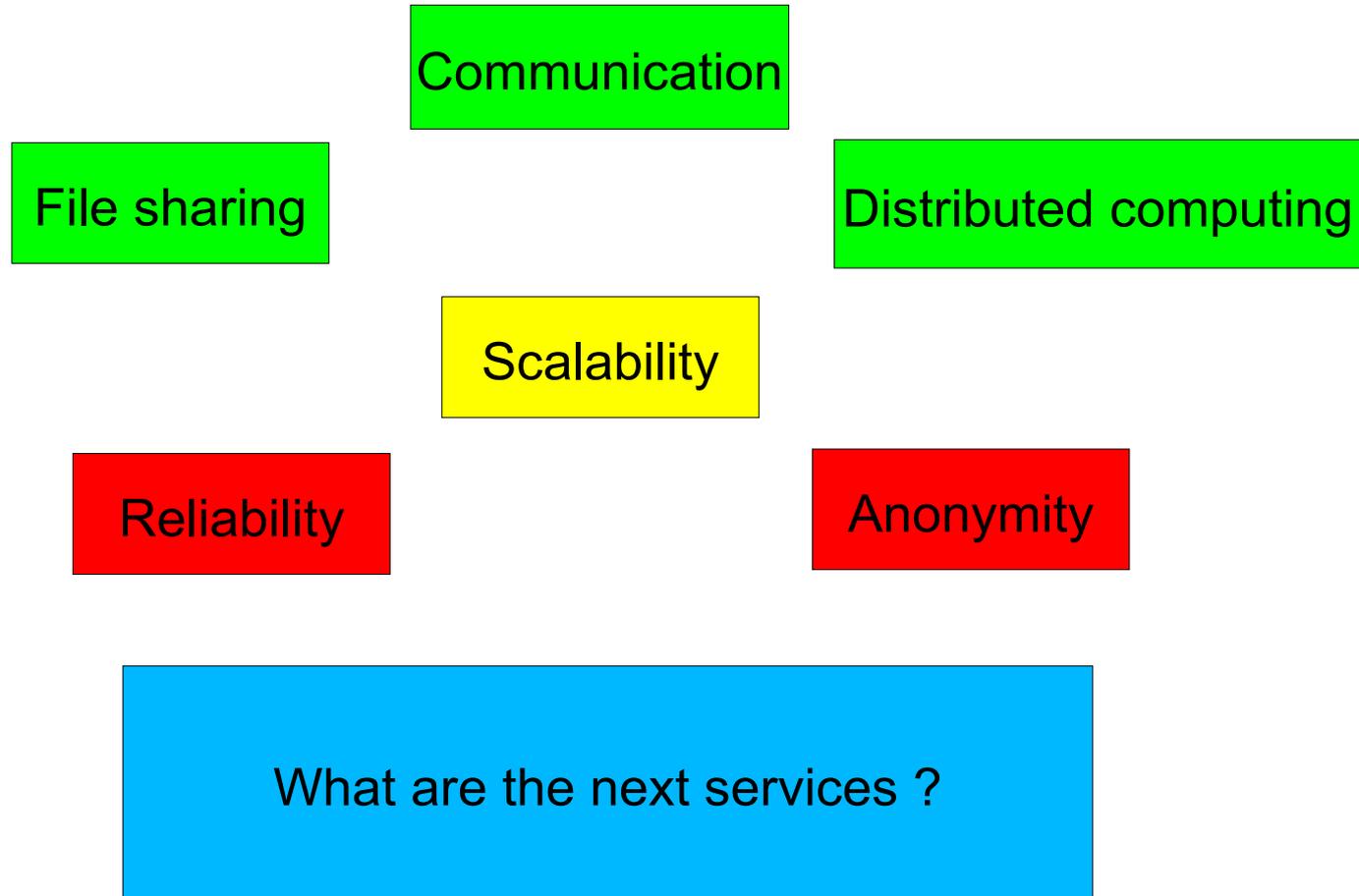
ICQ
M\$ and AOL crap

Distributed computing

Zetagrid
SETI@HOME

..awkward !! ..any ideas ?

Let's think of services



Peer-to-Peer Communities

	<p>News & Announcements Read this before submitting your first post to any forum Moderator Global Moderators</p>	472	495	Fri Oct 01, 2004 1:51 am tomk →
	<p>Frequently Asked Questions Some of the most commonly heard questions in the Gentoo Community, along with answers. Moderator Global Moderators</p>	67	82	Wed Sep 08, 2004 4:01 pm pip →
	<p>Installing Gentoo Having non-GUI problems with the Installation Guide? If you're still working your way through it, or just need some info before you start your install, this is the place. All other questions go elsewhere. Moderator Global Moderators</p>	23984	140332	Fri Oct 01, 2004 7:00 am Elamite →
	<p>Multimedia Help with creation, editing, or playback of sounds, images, or video. XMMS, mplayer, grip, cdparanoia and anything else that makes a sound or plays a video. Moderator Global Moderators</p>	14094	74655	Fri Oct 01, 2004 7:11 am sear →
	<p>Desktop Environments Problems with GUI applications? Questions about X, KDE, Gnome, Fluxbox, etc.? Come on in. NOTE: For multimedia, go up one forum Moderator Global Moderators</p>	31234	188905	Fri Oct 01, 2004 7:11 am eeknay →
	<p>Networking & Security Having problems getting connected to the internet or running a server? Wondering about securing your box? Ask here. Moderator Global Moderators</p>	21844	109526	Fri Oct 01, 2004 7:11 am BoBB →
	<p>Kernel & Hardware Kernel not recognizing your hardware? Problems with power management or PCMCIA? What hardware is compatible with Gentoo? See here.</p>	18237	107848	Fri Oct 01, 2004 7:07 am lamekain →



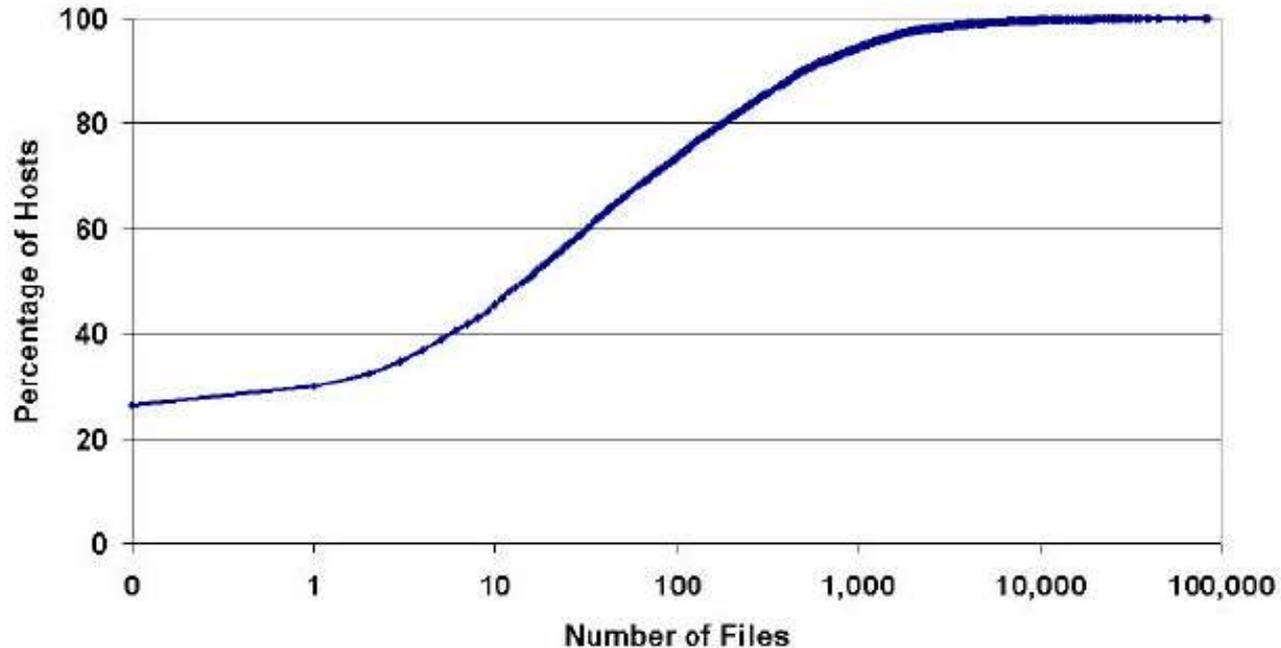
Boards

Community basis

Files

Free riding

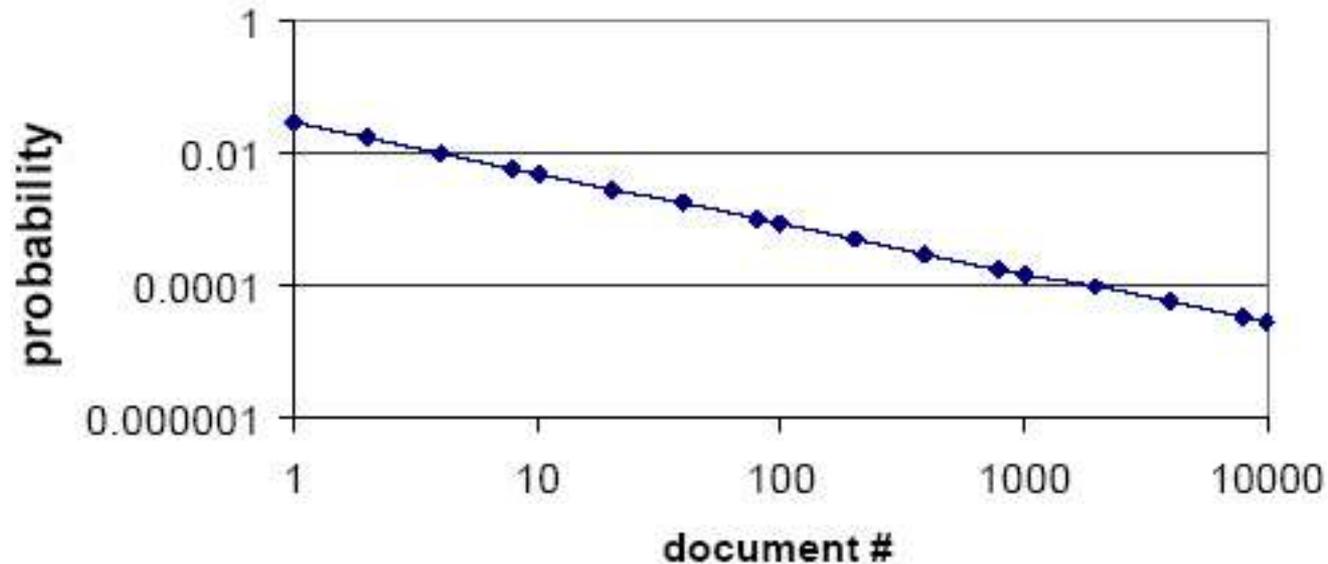
CDF of Number of Shared Files (Gnutella)



Adar/Huberman

20% host 98% of all content
60-70% peers don't share

Data distribution



$y = r^{-b}$ -> Zipf Distribution of data, not random !
(Examples: wealth among firms / words in human languages)

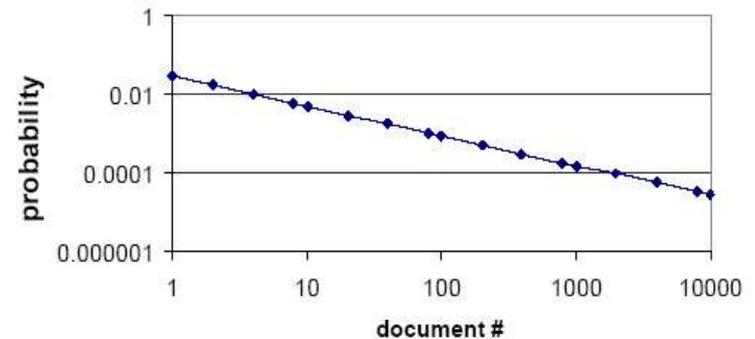
Performance issues

For instance Gnutella I

- Simple Protocol
- Plain keyword searching
- Random selection of the peers
- No central point of failure
- No caching of the peers or the data
- No redundancy of the data
- Vulnerable to DoS attacks
- Heavy messaging
- Scales bad

Time for improvements

- Data structures
 - Distributed Hash tables
 - Routing Indices
- Graphs
 - Small world effect
 - Power law distribution
- Semantic overlays
 - Description of content repository
- Agent based systems
 - Swarm-Intelligence, Learning
- Consider data distribution

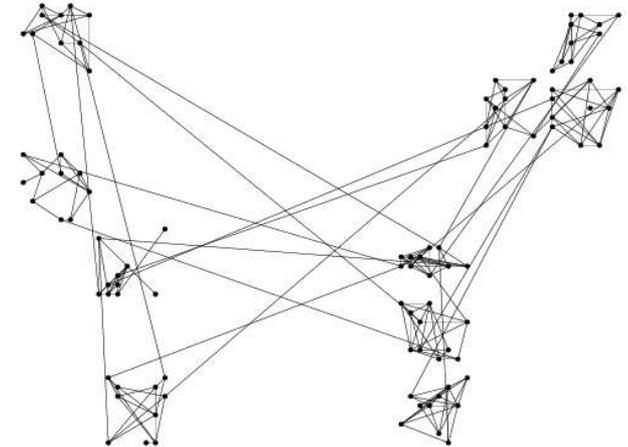


It's a small world after all...



The Kevin Bacon Game

Small World Graph



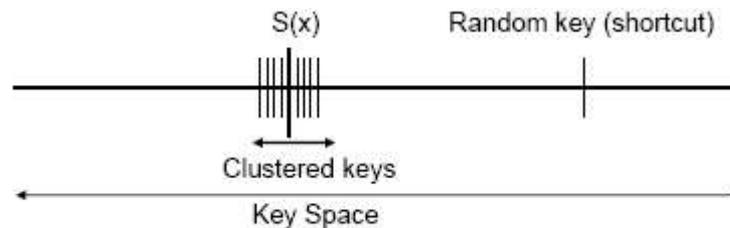
- Neuronal networks of worms
- Baseball players
- Power grids
- Web graphs
- Gnutella

Small worlds and P2P

2 approaches

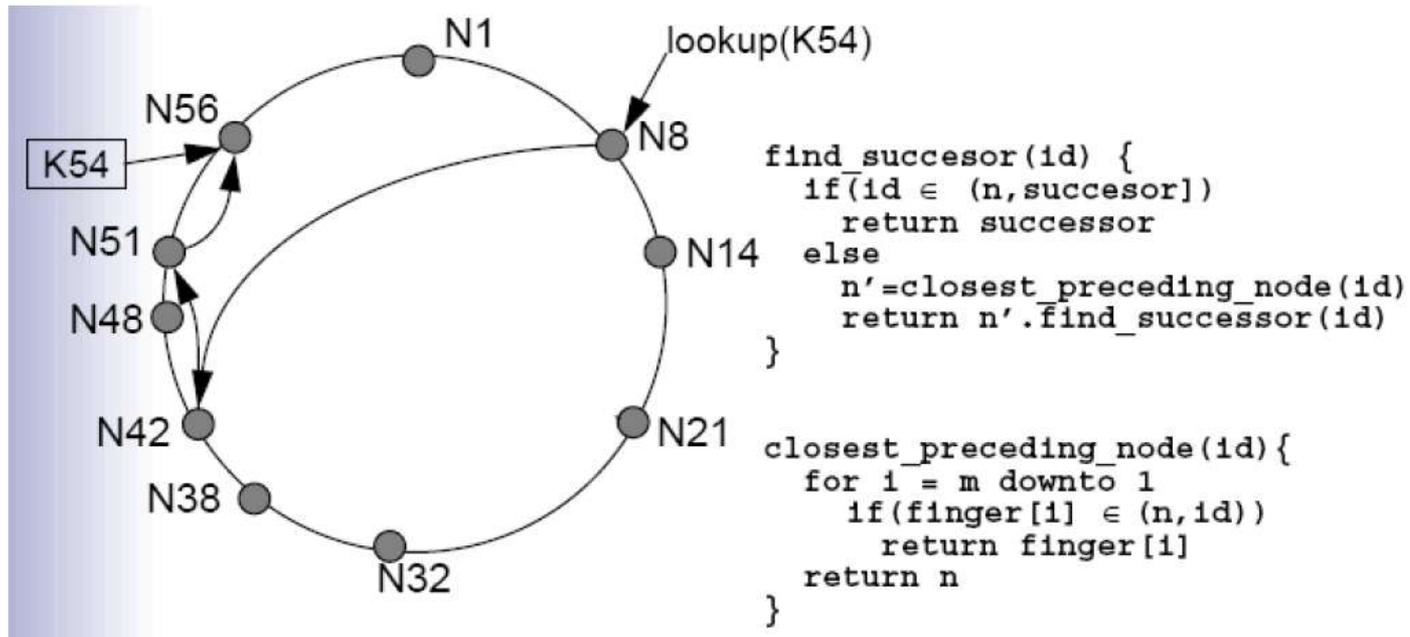
- Adapt the distributed protocol/algorithm to the six degrees of separation (the small world) world
- Build systems with small world attributes

Example:



CHORD – Scalable P2P

Load Balance, Decentralization, Scalability,
Availability, Flexible naming space



Scalable systems

Concepts

- CAN
- PASTRY
- TAPESTRY
- JXTA

Applications

 The Free Network Project



Edutella

Back to file sharing

Assumptions made by popular media

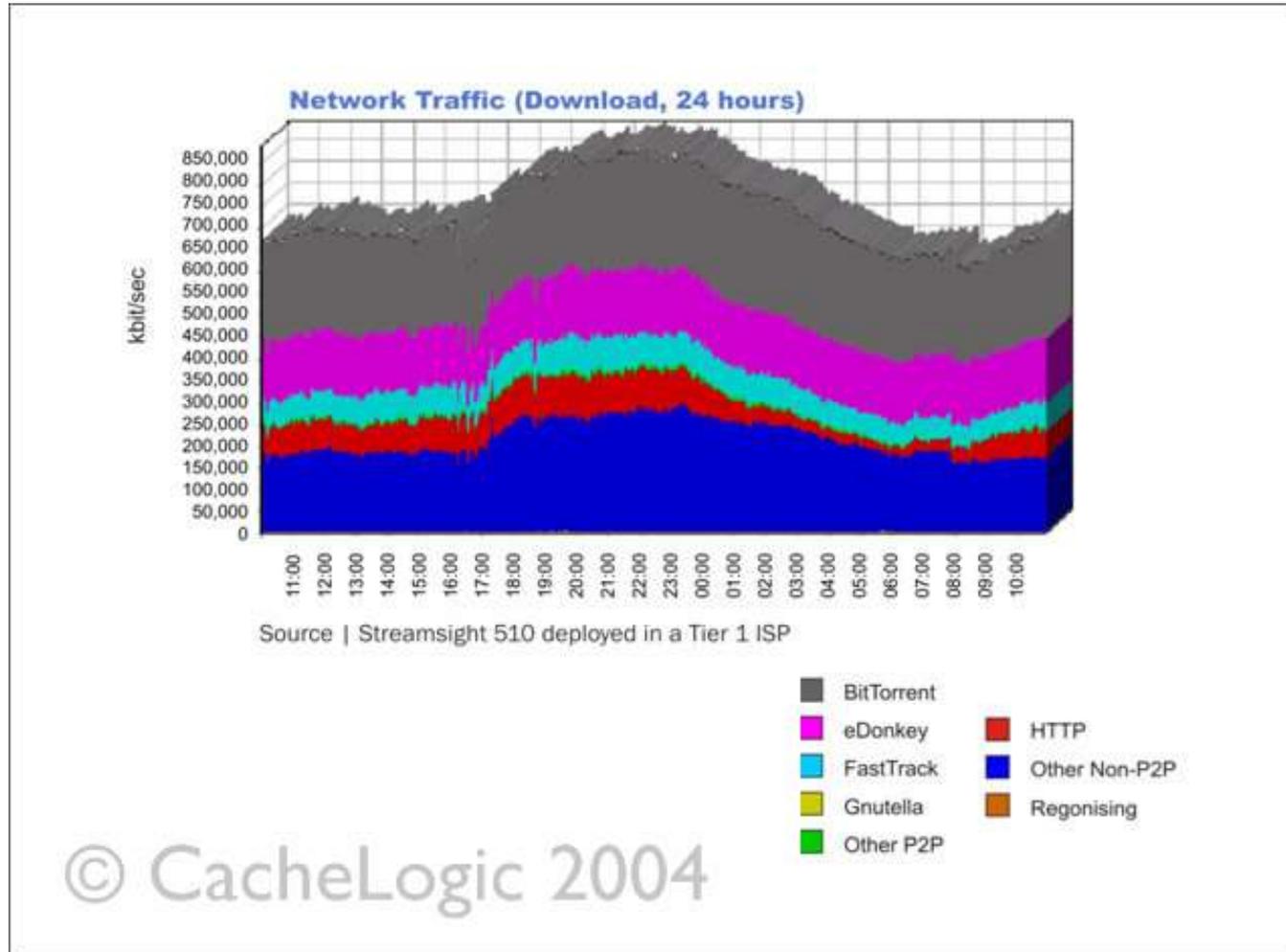
- File sharing is on the decline
- Those nets are all about music and video
- Edonkey is the new leader, ahead of Kazaa
- P2P = illegal sharing of files



What we will do here

- System analysis (traffic, content, distribution)
- Services for real world P2P systems

File sharing traffic



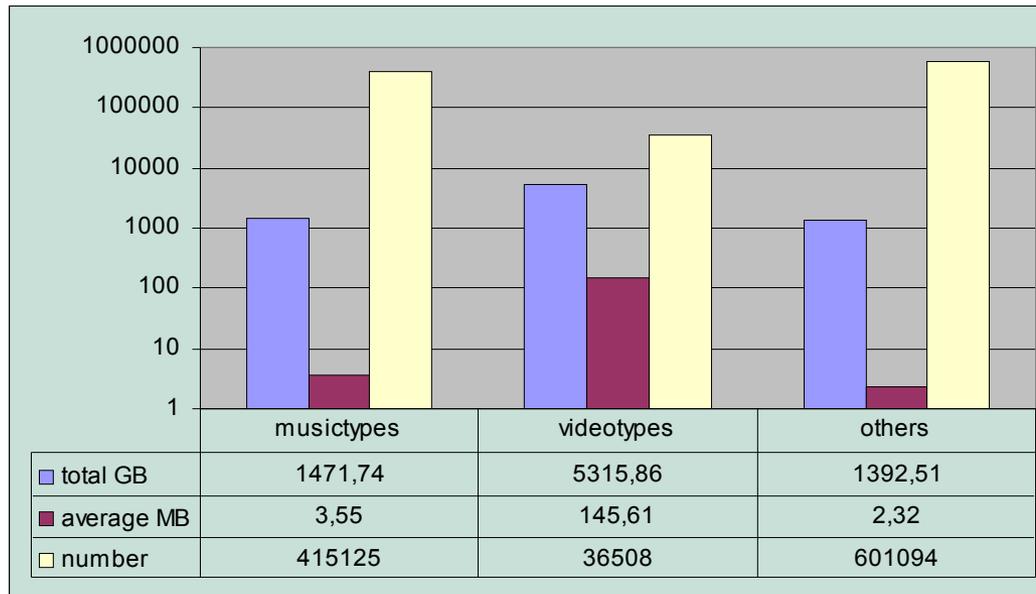
Where have all the flowers gone..

- Traditional traffic measurements don't work for P2P
- P2P applications use various ports dynamically
- P2P protocols use common ports (80) to jump over firewall/NAT barriers (e.g. jxta uses http)
- Users and projects switched to Bittorrent recently



---> **P2P is not on decline but is hiding !**

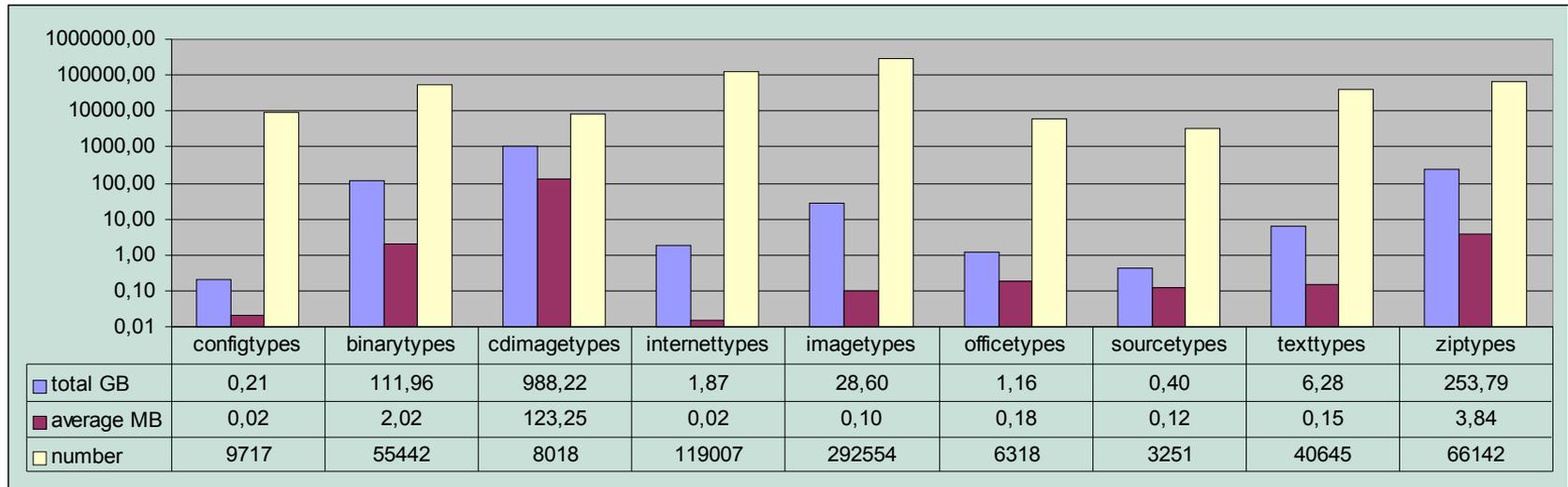
File sharing content



Direct Connect Hubs
November 04

Just sound and video content ?
What other content is there ?

What other content types ?



- High diversity of the files
- Vast amount of Internet data (html-files..)
- Too much data to ignore => NONE movie/music places

Services for huge databases

Goals Personalization, collaboration

Example How much customers from Schöneberg
buy ice cream on Fridays ?

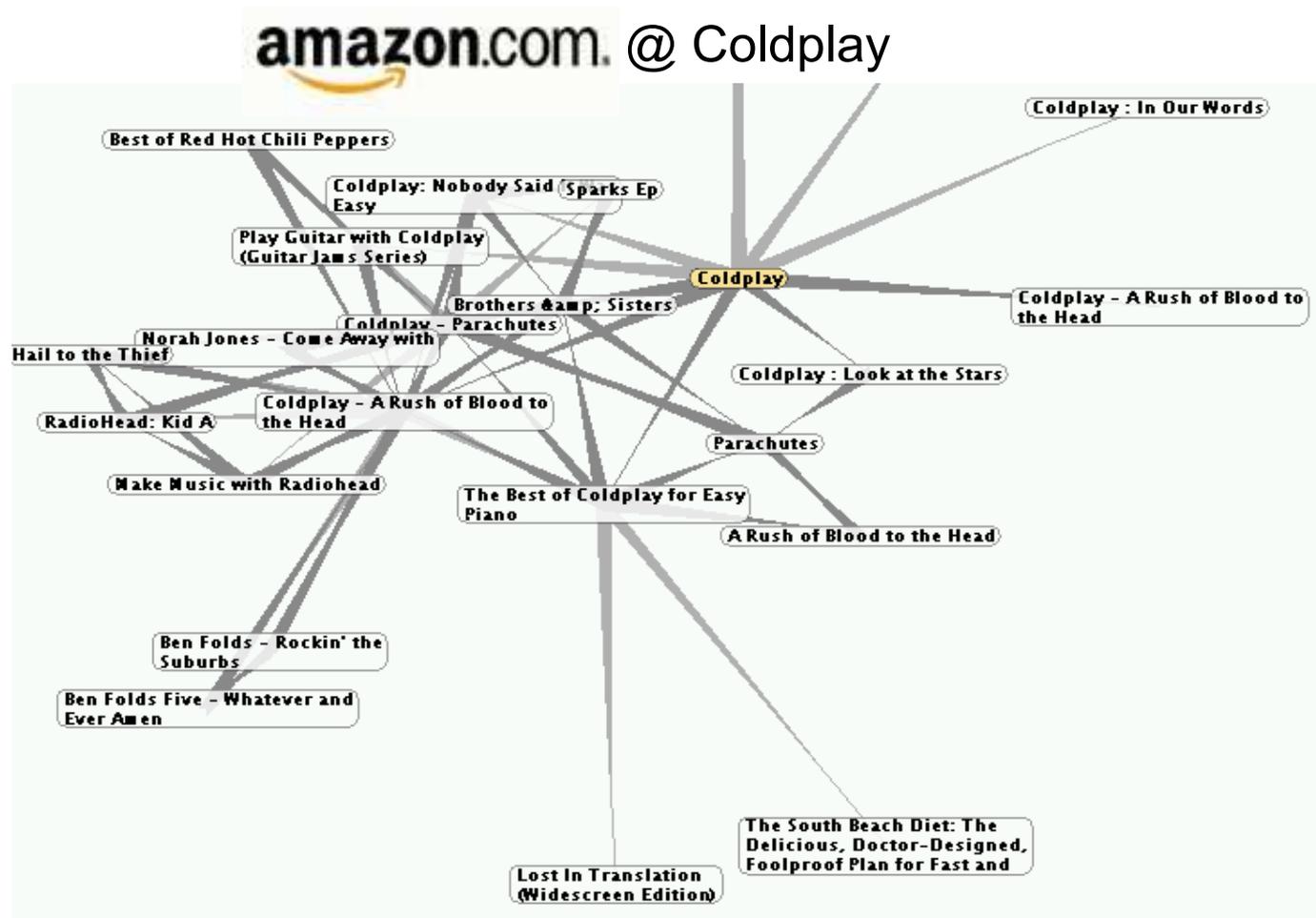
Users who bought this book also bought...

Services Collaborative filtering, recommender systems

How to get the recommendations ?

- User Input Community-based, but users are lazy !
- Automatic

amazon recommendation service



Do we really need this ?

- **Recommendation service**
- **Knowledge discovery**



Types of RS

Bayesian networks

Model-based, Decision trees
complex, not very fast

Clustering

Clutch users into groups with
similar interest
**recommended product is
the average of a group**

Association rules

Relations between users or data
**sounds easy but needs to be
adapted to P2P** →

Challenges for P2P RS

Mining strategy

Data extraction and conversion

Validate and represent the rules

Data Mining algorithm

scalable, dynamic, ad-hoc,
asynchronous, suitable

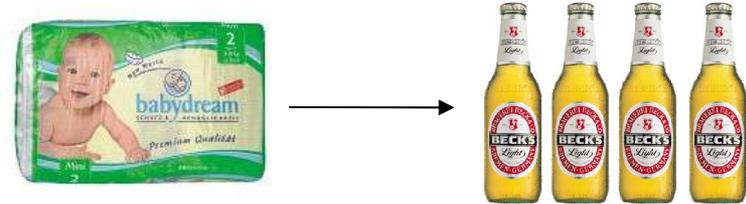
Association Rule Mining

Set of items: $I = \{I_1, I_2, \dots, I_m\}$

Transactions: $D = \{t_1, t_2, \dots, t_n\}, t_j \in I$

Item set X: $\{I_{i1}, I_{i2}, \dots, I_{ik}\} \in I$

Support $supp(X) = |X(t)| / |D|$



Transactions	Items			
T1	Beer	Linux	Honey	
T2	Diapers	Wine	Linux	
T3	Beer	Linux		
T4	Honey	Beer	Diapers	Bread
T5	Wine	Bread	Linux	Beer

- Association Rule $X \Rightarrow Y$ Implication $X \Rightarrow Y \quad X, Y \subseteq I \quad X \cap Y = \emptyset$
- Support $X \Rightarrow Y$ Transactions $X \cup Y$
- Confidence $X \Rightarrow Y$ $supp(X \cup Y) / supp(X)$
- Support-Confidence Framework: $X \Rightarrow Y$
 - $supp(X \cup Y) \geq minsupp$
 - $conf(X \Rightarrow Y) \geq minconf$

Task: Search frequent items, prune rules that don't interest

Distributed ARM in P2P

- Traditional ARM: *Sequentially*
Methods: Central Processing, Broadcasting, Global Synchronisation
- Peer to Peer: *Distributed*
 - Synchronisation: Autonomous behaviour of every node
 - Information never is up-to-date, highly dynamic
- Global Synchronisation: Computing costs are too high
Here: Local Algorithm: *Local majority voting*
(Wolff/Schuster)

Direct Connect network

Infrastructure

Structured by Hubs

User oriented

Content oriented

Data

Peers 341089

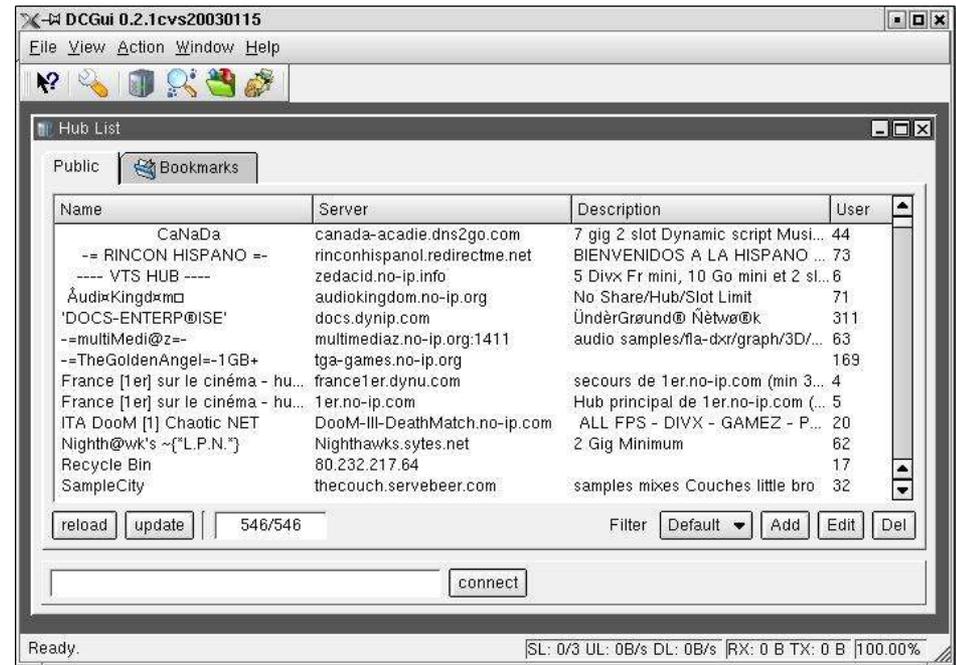
Data size **14486.50 Terra bytes!**

Clients

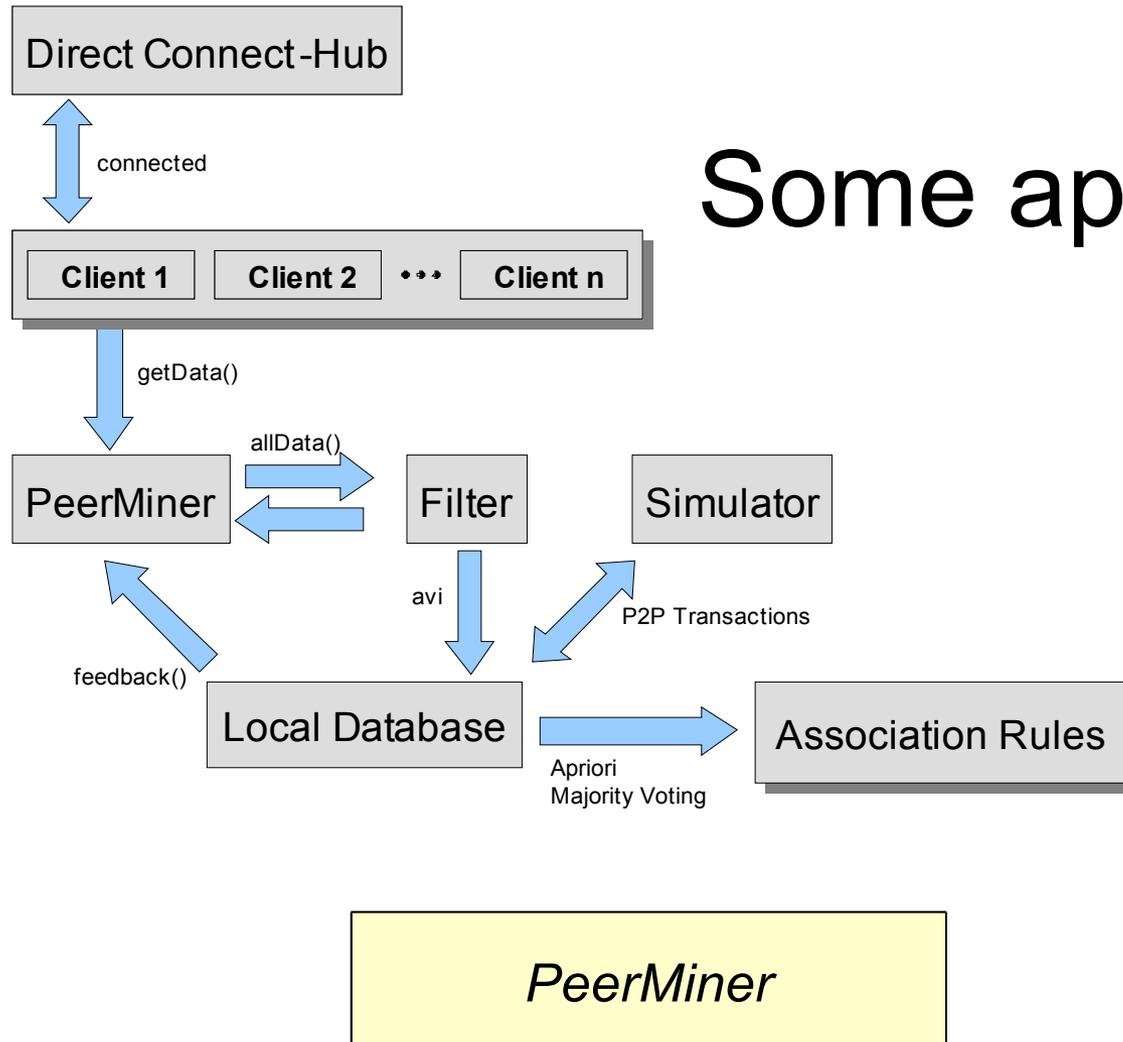
Original NeoModus DC

DC++

Valknut (!) ..aka DC-GUI



Some approach



Let's get some ~~movies~~

..ooops..sorry..rules

Content avi-video

Hubs Chosen by regular expressions ("movies",...)

Attributes 10

Instances 151

Minimum support 0.10

Minimum confidence 0.5

transaction database

Mystic River
Spartan
Kill Bill
Love actually
...

some popular rules

Love actually=no => Mystic River=yes && Kill Bill=no
Mystic River=yes && Love actually=no => Kill Bill=no
Love actually=no && Kill Bill=no => Mystic River=yes
Spartan=yes => Kill Bill=no
Love actually=no => Mystic River=yes



Current work

- Implementation of Distributed ARM Algorithms on P2P networks
- Social structures of P2P communities
- Distributed Hashing in JXTA
- Develop more accurate P2P simulators
- Measurement studies of other P2P systems (Bittorrent)

Summary

- Scientific P2P
- Current trends in file sharing
- Scalability and routing
- Coding resources
- Novel services for P2P systems
- Association Rule Mining



mag@tzi.de

*.... Questions and
discussion please...*