# DELOS Task 2.8:
# Personalized Query Routing in
# Peer-to-Peer Federations of Digital Libraries

**Christian Zimmer** (czimmer@mpi-inf.mpg.de)

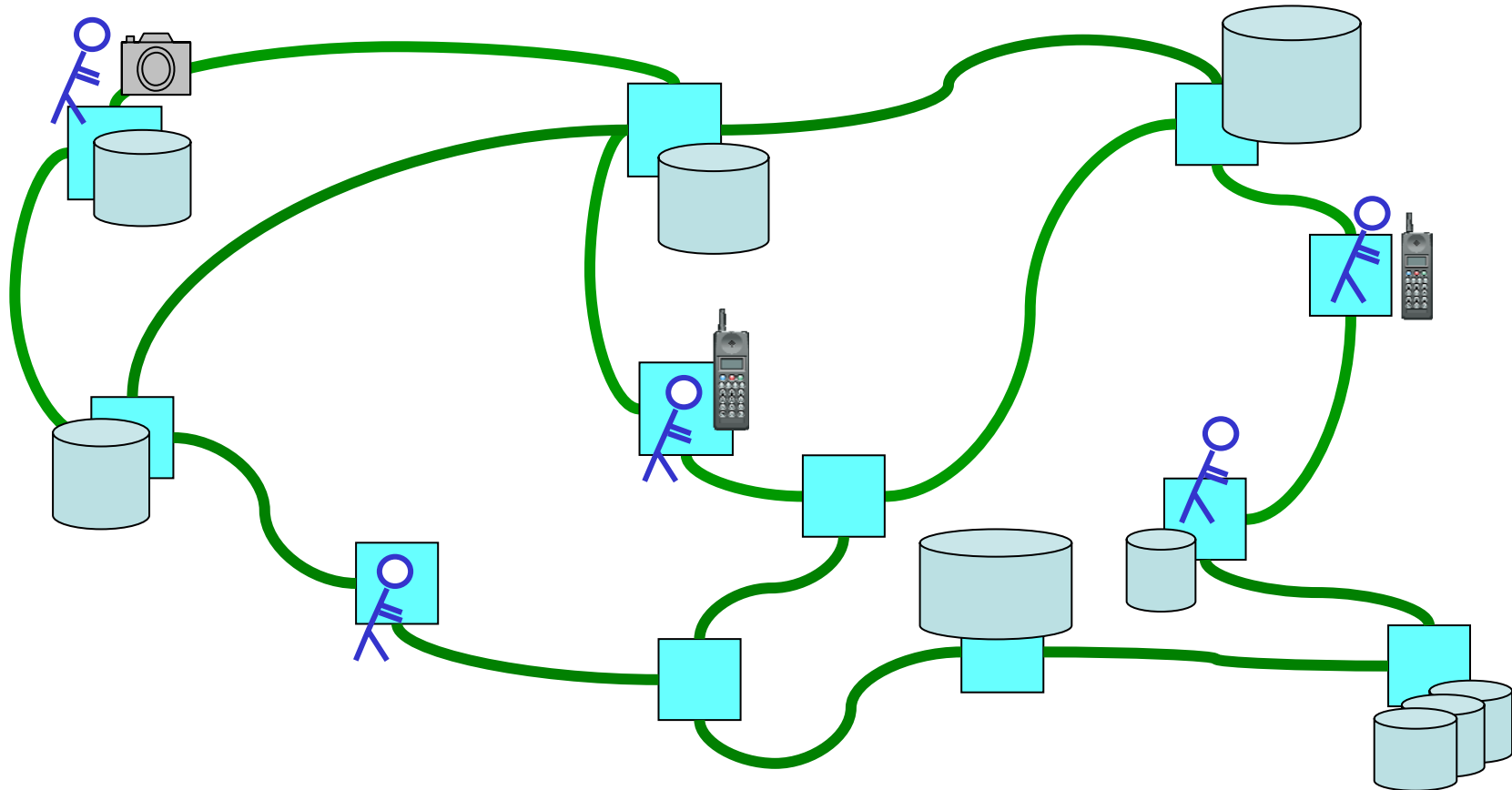**Gerhard Weikum** (weikum@mpi-inf.mpg.de)

max planck institut informatik

DELOS

NETWORK OF EXCELLENCE ON DIGITAL LIBRARIES
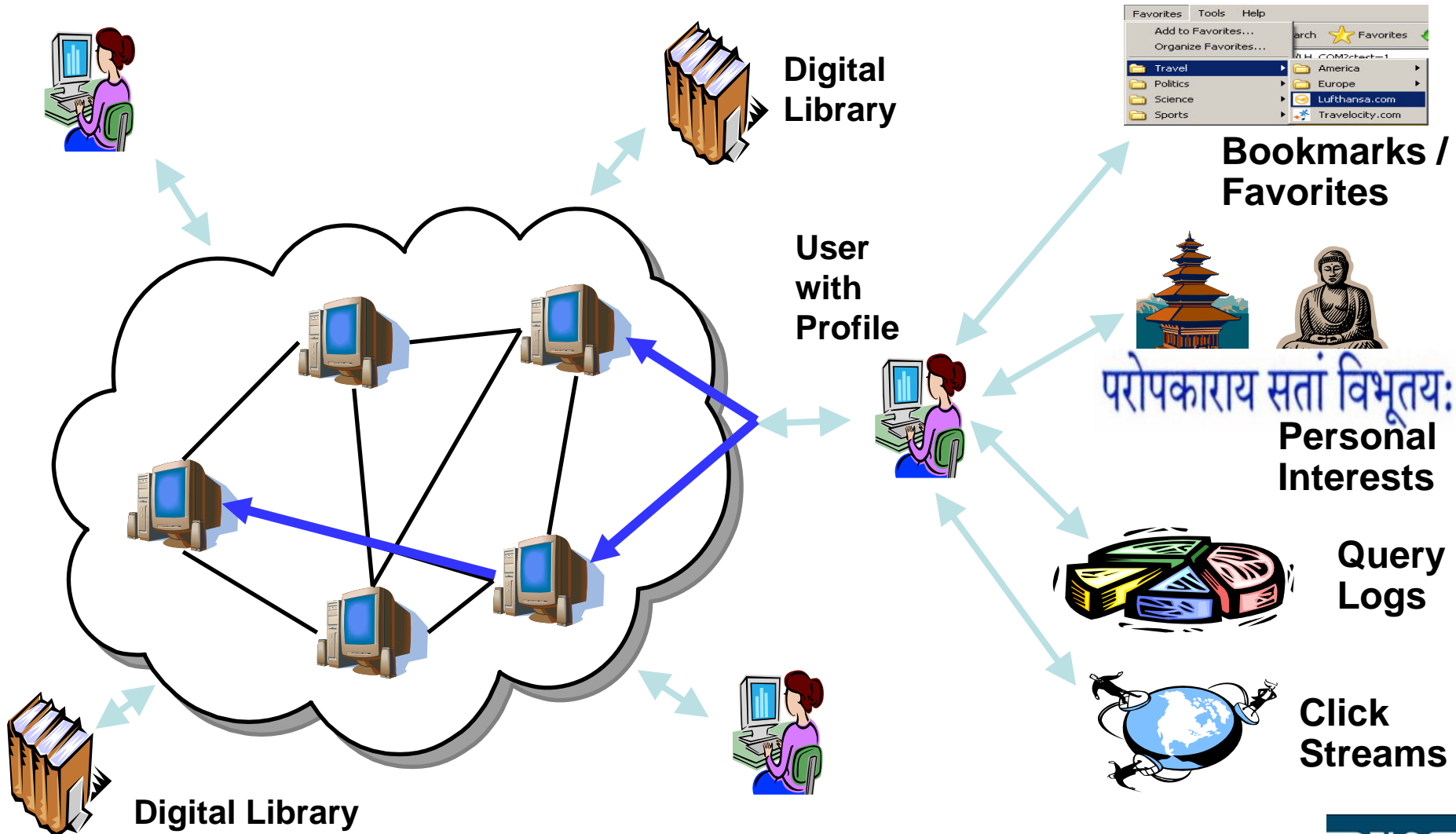
# P2P Architecture for DLs and DL Users

Self-organizing overlay networks for info sharing, PubSub, recommendations, search, routing (e.g. BitTorrent, Skype, etc.)



Peers:
- **DLs, Citation Servers, Annotation Servers, Image Repositories, Public Databases, Web Archives, News Feeds, Blogs, etc.**
- **Users, Mobile Devices, etc.**

# Opportunities and Challenges of Personalized P2P Search



Digital Library

Bookmarks / Favorites

User with Profile

परोपकाराय सतां विभूतयः

Personal Interests

Query Logs

Click Streams

Digital Library

# Task 2.8: Goal and Partners

**Goal:**

**models and strategies for personalized query routing (selecting peers based on user profile & history)**

**Partners and their Expertise:**

- **Max-Planck Institute for Informatics Saarbrücken (Gerhard Weikum):** P2P Web search
- **National University of Athens (Yannis Ioannidis):** user profiles, preference queries
- **University for Health Sciences Innsbruck (Hans-Jörg Schek):** relevance feedback, e-health apps
- **University of Duisburg-Essen (Norbert Fuhr):** P2P IR, DL agents
- **Masaryk University Brno (Pavel Zezula):** distributed similarity search
- **ETH Zurich (Donald Kossmann):** scalable, personalized PubSub, desktop search

# Outline

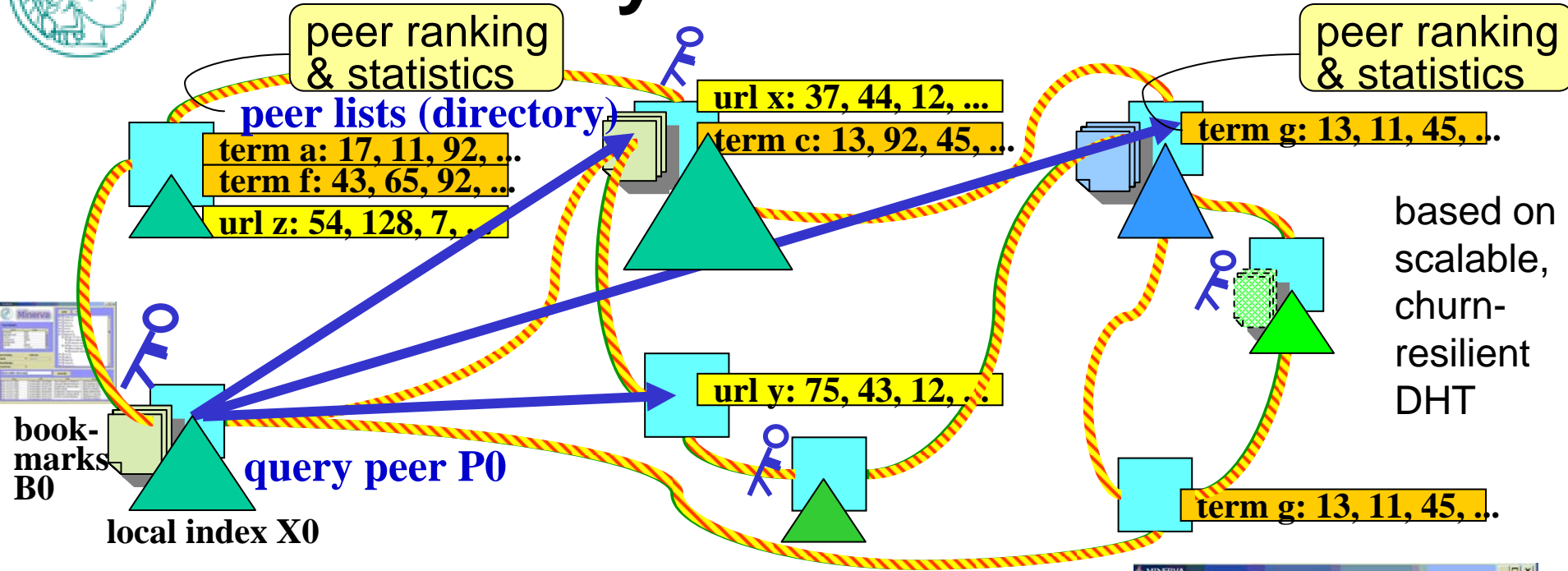✓ Motivation and Research Direction

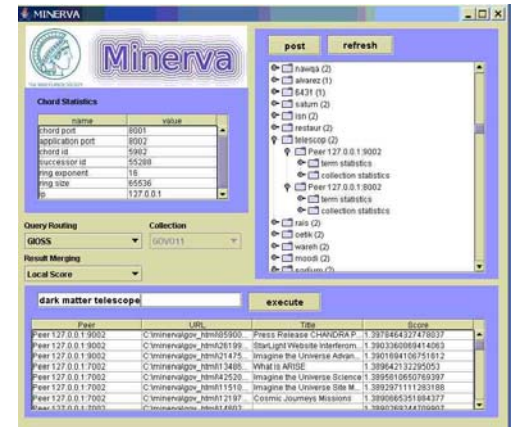- P2P Search Engine

- Query Routing

- Conclusion

# Minerva System Architecture



**Query routing** aims to optimize benefit/cost
driven by distributed statistics on
peers' content similarity, content overlap,
freshness, authority, trust, performability etc.

Dynamically precompute „good peers"
to maintain a **Semantic Overlay Network**
using random but biased graphs

# Minerva at Work

✓**Peers Registering with MINERVA**

   ✓Join DHT-style directory and inspect system status

   ✓Post statistical metadata about local index

   ✓Inspect metadata of other peers

✓**Query Routing and Processing with MINERVA**

   ✓Enter keyword query

   ✓Gather metadata from distributed directory to perform Query Routing

   ✓Execute query at selected peers using top-k query execution strategies

✓**Query Result Merging and Display**

   ✓Merge results into single result list at querying peer

   ✓Click on query results to view (cached copies of) web pages

# Outline

✓ Motivation and Research Direction

✓ P2P Search Engine

• Query Routing

• Conclusion

# Quality&Overlap-Aware Query Routing [SIGIR'05]

Select peers with highest benefit/cost ratio where
- benefit(Pi) ~ sim (X0, Xi) and ~ 1/overlap(X0, Xi)
  or using bookmarks B0, Bi for personalization & efficiency
- cost(Pi) ~ estimated response time or communication costs

precompute sim: $KL(X0, Xi) := \sum_{terms\ x} freq(x, X0)\ log \dfrac{freq(x, X0)}{freq(x, Xi)}$

estimate overlap by Bloom filters, hash sketches, or MIPs

Experiments:

based on 100 .Gov partitions (1.25 Mio. docs), assigned to 50 peers,
with each peer holding 10 partitions and 80% overlap for $P_i$, $P_{i+1}$
with 50 TREC-2003 Web queries, e.g.: „juvenile delinquency"



recall vs. # queried peers — overlap-aware Minerva, CORI baseline

# Considering Term Correlations [IPTPS'06]

<u>Problem:</u> DHT-based Per-Term Directory loses term correlations such as „Michael Jordan" or „Native American Music"

<u>Solution:</u>
- peers perform frequent-itemset mining on local query log
- correlated termsets posted to all single-term directory peers
- directory peers collect postings for termsets from all peers
- query routed to single-term peers, evaluated over max. termsets
- all communication piggybacked on normal traffic, no extra cost

experiments based on 750 peers
with .Gov partitions,
running expanded queries
from TREC-2003 Web track;
<u>examples:</u>
„marijuana legalization drug abuse ...",
„wireless communication broadcasting"

# Distributed Similarity Search in Metric Spaces

Problem:

Scalable distributed indexing of

data objects for kNN queries with metric distances

satisfying triangle inequality dist(x,z) $\leq$ dist(x,y) + dist(y,z)

Approach: [Delos 2005]

• embed data objects into distance-preserving vector space

• map kNN queries into range queries

• index by dynamic partitioning across peers of DHT

*Example: Edit Distance*

query q: *Mex Plank Institute*　　　　　and then submitted

　should be corrected into　　　　　to P2P search

query q': *Max Planck Institut*　　　　　(joint work MPII & MUNI)

　based on P2P directory

# Continuous Queries in P2P Publish-Subscribe

IR (Information Retrieval):

best results for                    vs.

one-time query

IF (Information Filtering):

alerting about new docs

that match standing query

State-of-the-art IF considers only exact matches
and has only coarse-grained topics for personalization

Challenge (work in progress):

Approximate IF

should alert the user about vague matches

and may miss some docs with low probability

for better P2P scalability and churn-resilience,

and can support fine-grained personalization

# Outline

✓ Motivation and Research Direction

✓ P2P Search Engine

✓ Query Routing

- Conclusion

# Conclusion

**P2P search engines have great potential:**

• harness local resources for power search engine

• rich models for content extraction, annotation, summarization, and indexing of text, images, speech, audio&video, feeds, portals

• customization and personalization

• collaboration & recommendation networks with other peers

• naturally fits with mobile clients and context awareness

• naturally gears for rich cognitive model of user behavior

• no monopoly, no central profiling or bias

**Query routing is the key issue in P2P search**

**Task 2.8: 6 partners (MPII, NUA, UMIT, UniDU, MUNI, ETHZ)**

• complementary expertise and potential for synergies

• collaboration started (dedicated 2-day workshop, bilateral visits)

max planck institut informatik

DELOS
NETWORK OF
EXCELLENCE ON
DIGITAL
LIBRARIES