

Kapitel 17: Verteilte Informationssysteme

17.1 Verteilte Systemarchitekturen

17.2 Anfrageausführung in verteilten IR-Systemen

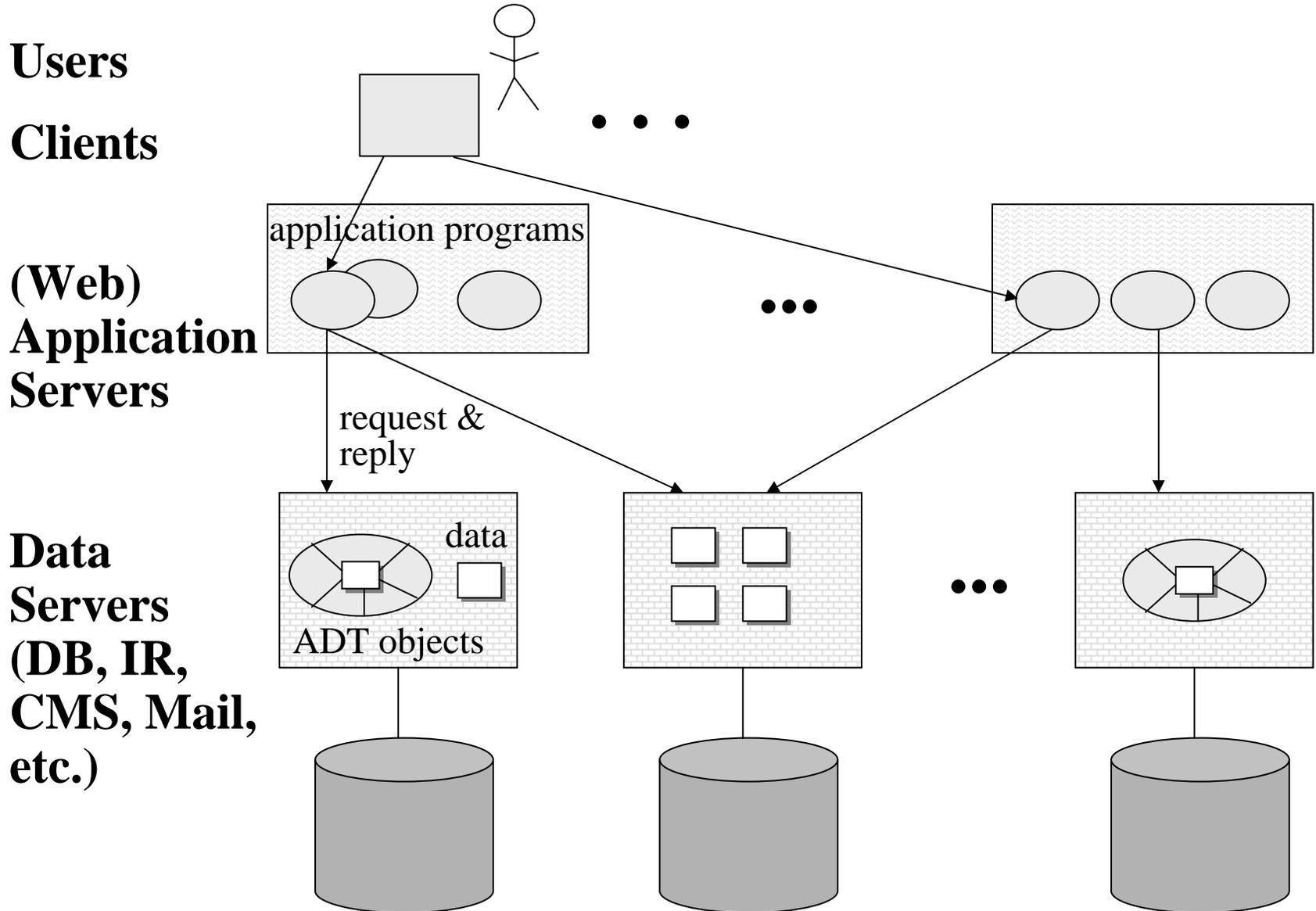
17.3 Skalierbare verteilte Indexierung und Suche

17.4 Anfrageausführung in verteilten DB-Systemen

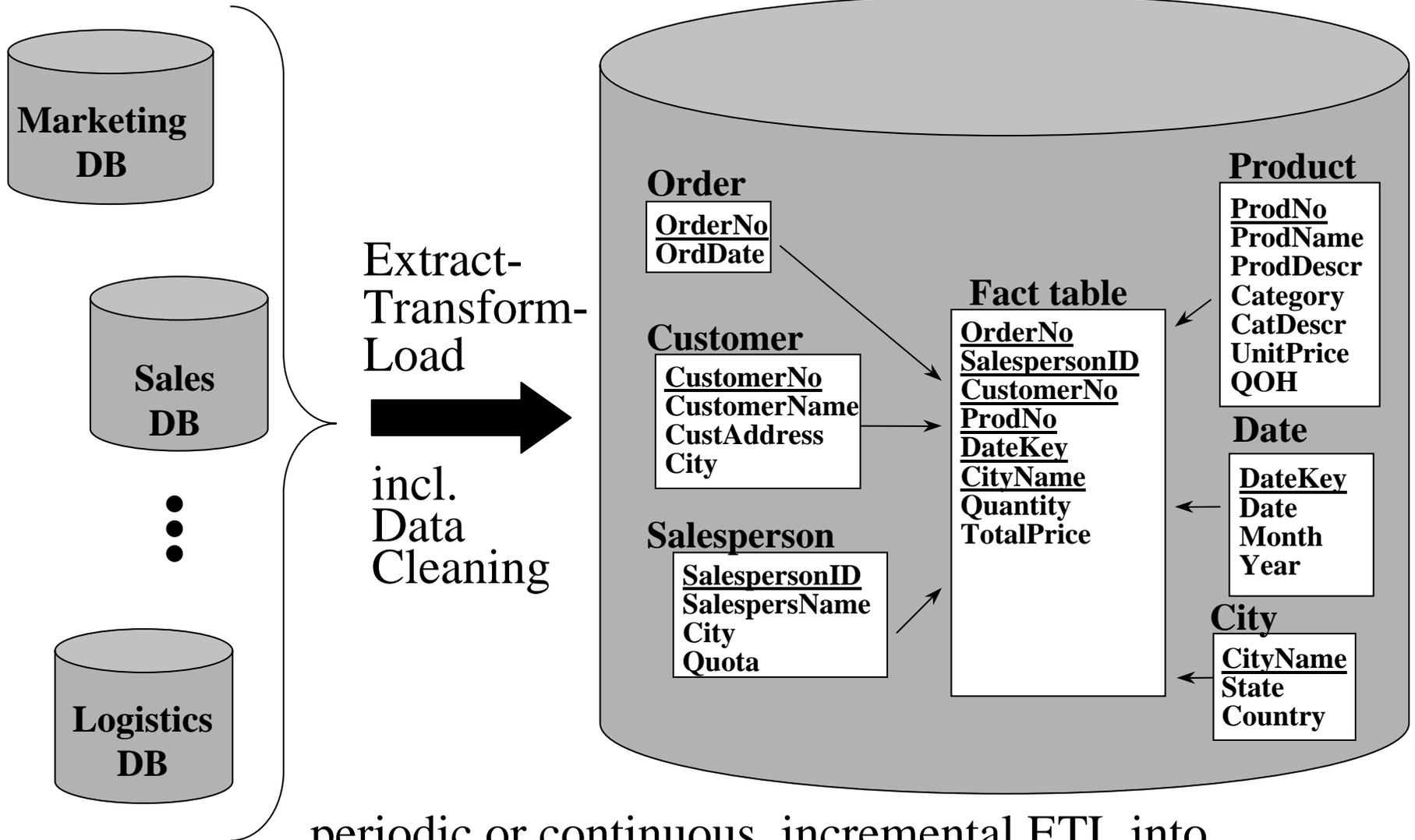
17.1 Verteilte Systemarchitekturen

- **Software-Component-Oriented:**
 - **Client-Server Architecture**
 - **Multi-Tier Architecture**
 - **Federated Systems**
- **Data-Source-Oriented:**
 - **Data Warehouses (and Digital Libraries)**
 - **Wrapper-Mediator Architecture for Information Integration (Example: Internet Portals)**
- **Uncoordinated Decentralization:**
 - **Service-Oriented Architecture for Web Services or Grid**
 - **Peer-to-Peer Systems (P2P)**

Client-Server and (Federated) Multi-Tier Systems

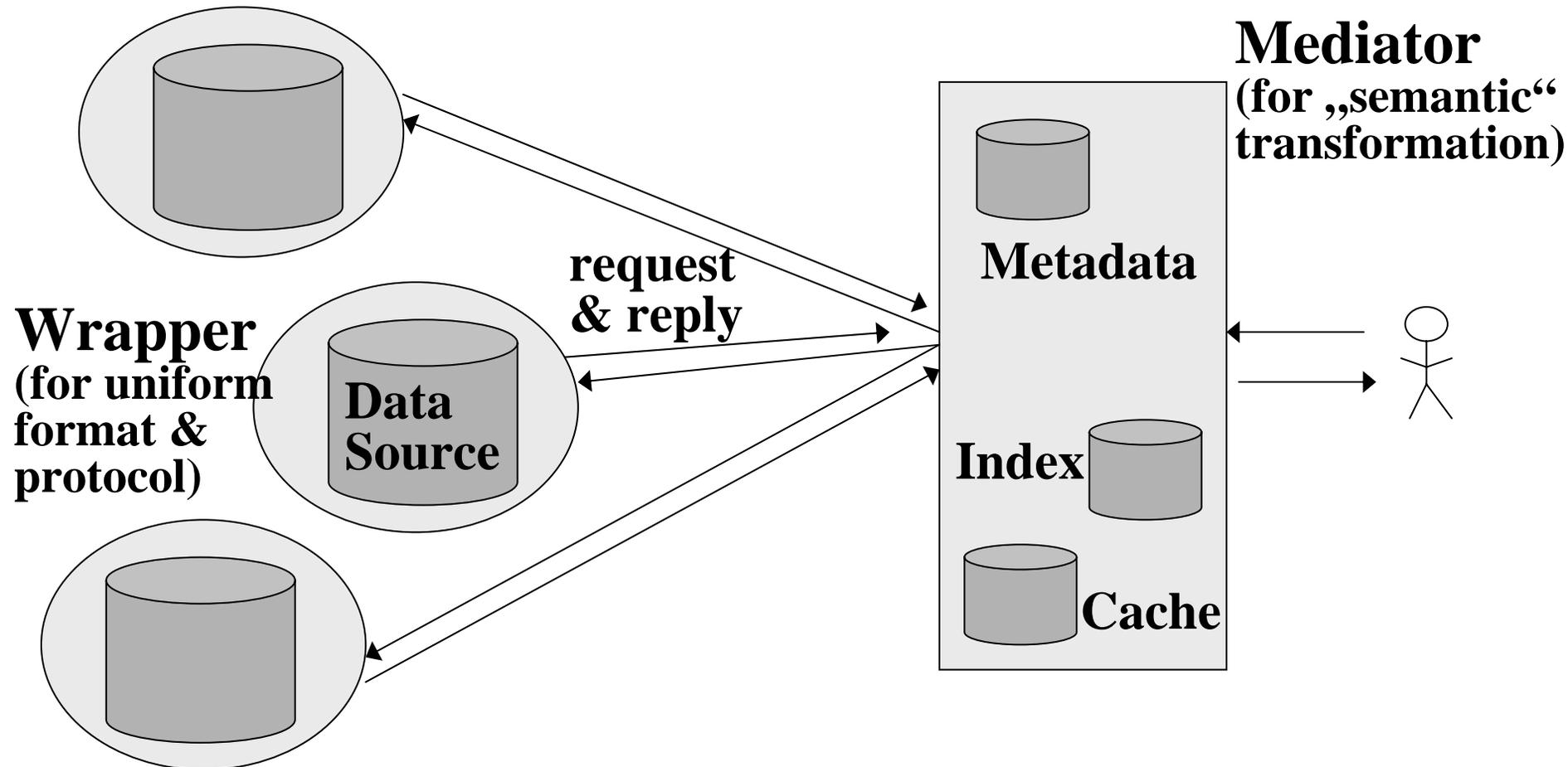


Data Warehouses (and Digital Libraries)



periodic or continuous, incremental ETL into DW with star or snowflake schema (facts, dimensions)

Wrapper-Mediator Architecture for Information Integration Systems



Wrapper
(for uniform
format &
protocol)

Mediator
(for „semantic“
transformation)

request
& reply

**Data
Source**

Metadata

Index

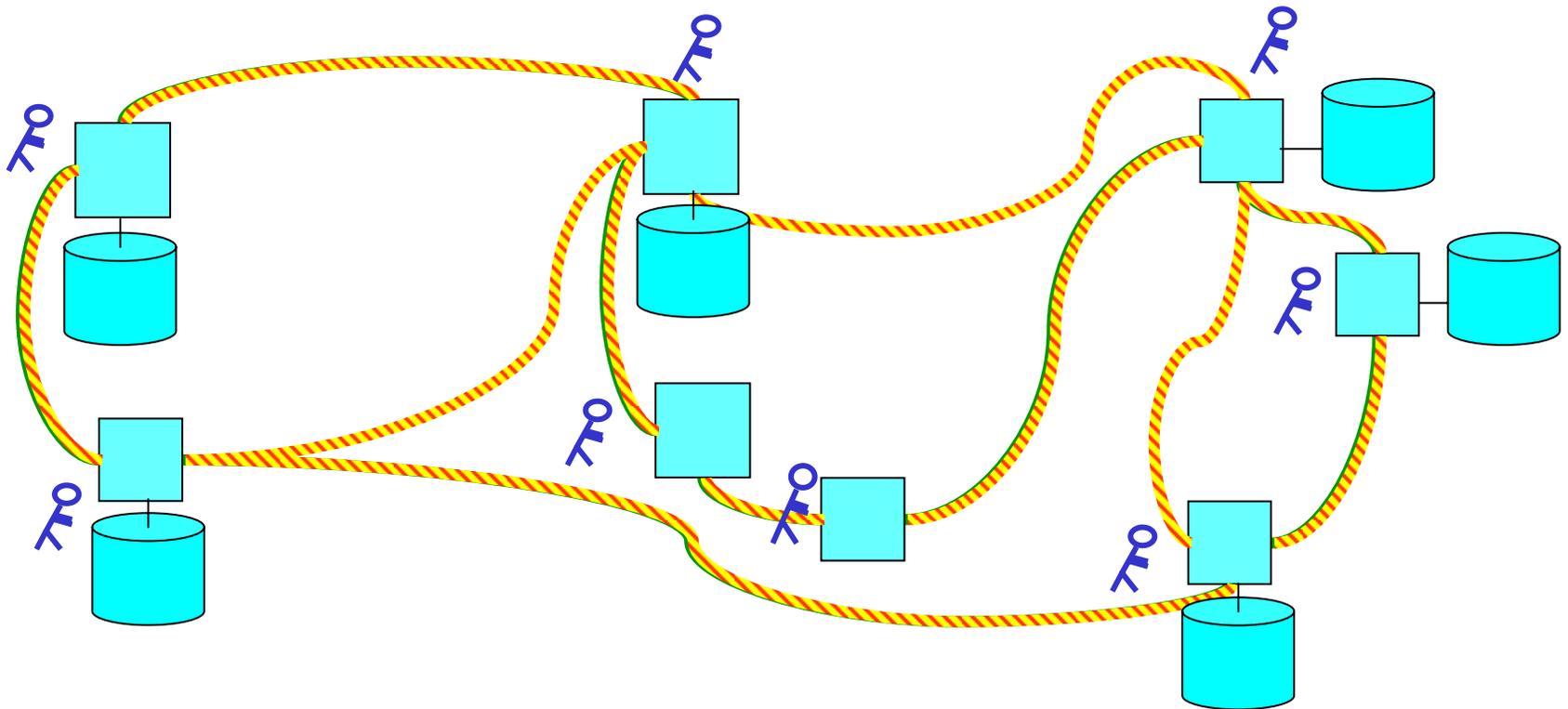
Cache

**Wrappers may be
hand-engineered
or generated**

Examples:

- **Yahoo** for news, business, etc.
- **SRS** for bioinformatics & life sciences
- **intranet portals** for organizations

Peer-to-Peer (P2P) Information Sharing



User post sharable files on autonomous computers (peers)
P2P system maintains (distributed) directory with file names
Users query file names and download files from other peers

Peer-to-Peer (P2P) Architectures

**Decentralized, self-organizing, highly dynamic
loose coupling of many autonomous computers**

Applications:

- **Large-scale distributed computation (SETI, PrimeNumbers, etc.)**
- **File sharing (Napster, Gnutella, KaZaA, etc.)**
- **Publish-Subscribe Information Sharing (Marketplaces, etc.)**
- **Collaborative Work (Games, etc.)**
- **Collaborative Data Mining**
- **(Collaborative) Web Search**

Goals:

- **make systems ultra-scalable and completely self-organizing**
- **make complex systems manageable and less susceptible to attacks**
- **break information monopolies, exploit small-world phenomenon**

1st-Generation P2P

**Napster (1998-2001) and Gnutella (1999-now):
driven by file-sharing for MP3, etc.
very simple, extremely popular**

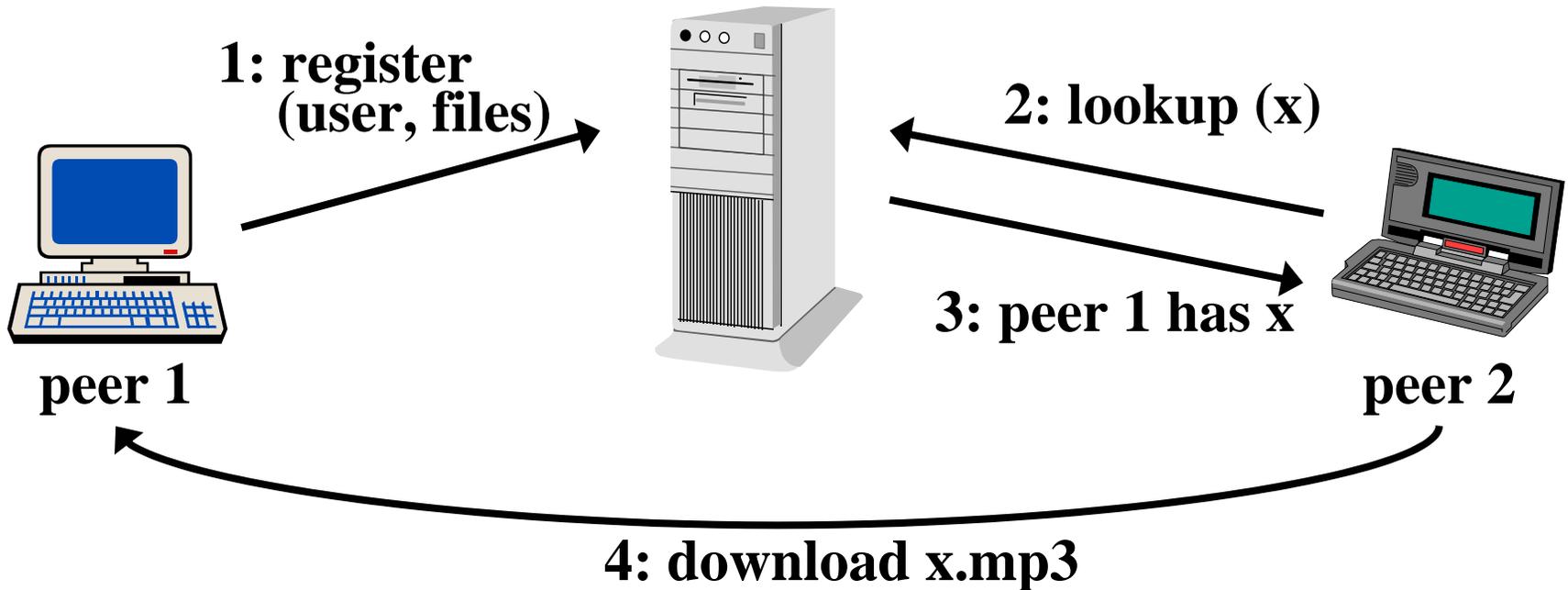
**can be seen as a mega-scale but very simple
publish-subscribe system:**

- **owner of a file makes it available under name x**
- **others can search for x, find copy, download it**

invitation to break the law (piracy, etc.) ?

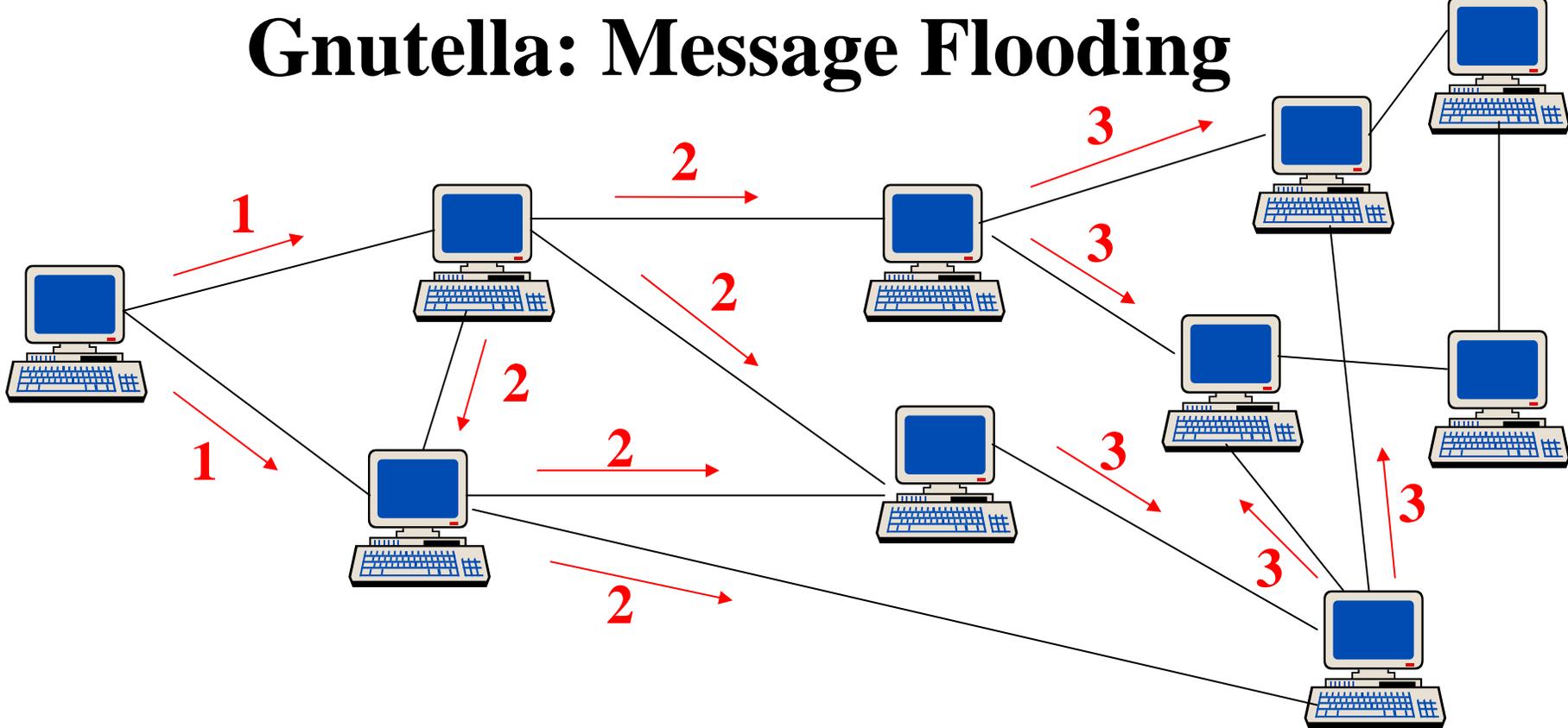
Napster: Centralized Index

Napster server



+ chat room, instant messaging, firewall handling, etc.

Gnutella: Message Flooding



all forward messages carry a TTL tag (time-to-live)

- 1) contact neighborhood and establish virtual topology (on-demand + periodically): *Ping, Pong***
- 2) search file: *Query, QueryHit***
- 3) download file: *Get or Push* (behind firewall)**

2nd-Generation P2P

Freenet

emphasizes anonymity

eDonkey, KaZaA (based on FastTrack), Morpheus, MojoNation, AudioGalaxy, etc. etc.

**commercial, typically no longer open source;
often based on super-peers**

JXTA

(Sun-sponsored) open API

Research prototypes (with much more refined architecture and advanced algorithms):

Chord (MIT), CAN (Berkeley), OceanStore/Tapestry (Berkeley), Farsite (MSR), Spinglass/Pepper (Cornell), Pastry/PAST (Rice, MSR), Viceroy (Hebrew U), P-Grid (EPFL), P2P-Net (Magdeburg), Pier (Berkeley), Peers (Stanford), Kademia (NYU), Bestpeer (Singapore), YouServ (IBM Almaden), Hyperion (Toronto), Piazza (UW Seattle), PlanetP (Rutgers), SkipNet (MSR), Galanx (U Wisconsin), Minerva (MPII), etc. etc.

The Future of P2P: Challenging Requirements

**Unlimited scalability with millions of nodes:
 $O(\log n)$ hops to target, $O(\log n)$ state per node**

**Failure resilience, high availability, self-stabilization
(w.r.t. dynamics: many failures & high churn)**

**Data placement, routing, load management, etc.
in overlay networks**

Robustness to DoS attacks & other traffic anomalies

Trustworthy computing and data sharing

**Incentive mechanisms to reconcile selfish behavior
of individual nodes with strategic global goals**

P2P-Related Technologies

Web Services (SOAP, WSDL, etc.)

for e-business interoperability (supply chains, etc.)

Grid Computing

for scientific data interoperability

Autonomic / Organic / Introspective Computing

for self-organizing, zero-admin operation

Multi-Agent Technology

for interaction of autonomous, mobile agents

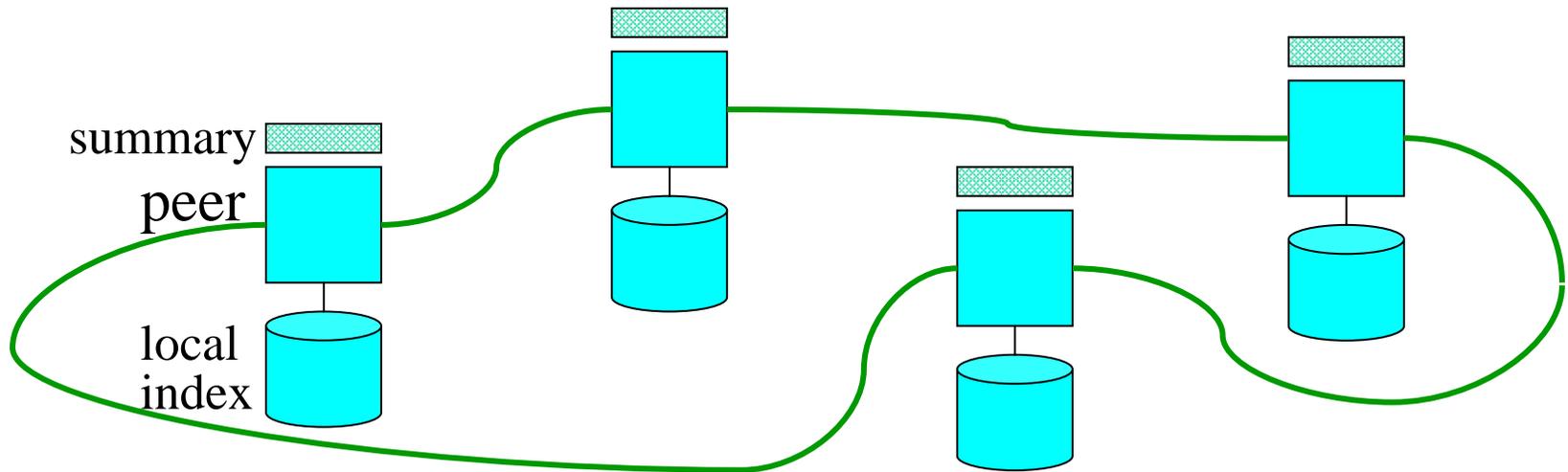
Sensor Networks

for data streams from measurement devices etc.

Content-Delivery Networks (e.g., Akamai)

for large content of popular Web sites

17.2 Anfrageausführung in verteilten IR-Systemen



- every peer is autonomous and has its own **local search engine**
- every peer posts (statistical) **summary info** about its contents (index lists, bookmarks, cached docs, QoS properties, ...)
- **query routing** is driven by similarity to summaries
- summaries are organized into a **(distributed) directory**
 - mapped onto DHT, random-graph overlay network, etc.
 - lazily replicated at additional peers (via „gossiping“)

querying peer needs to

1. determine interesting peers (**query routing**)
2. plan, run, monitor, and adapt **distributed top-k** algorithm
3. **reconcile results** from different peers

Why Peer-to-Peer Web Search?

Goal: Self-organizing P2P Web Search Engine
with Google-or-better functionality

- **Scalable & Self-Organizing** Data Structures and Algorithms
(DHTs, Semantic Overlay Networks, Epidemic Spreading, Distr. Link Analysis, etc.)
- Better Search Result **Quality** (Precision, Recall, etc.)
 - Powerful Search Methods for Each Peer
(Concept-based Search, Query Expansion, Personalization, etc.)
 - Leverage Intellectual Input at Each Peer
(Bookmarks, Feedback, Query Logs, Click Streams, Evolving Web, etc.)
 - Collaboration among Peers
(Query Routing, Incentives, Fairness, Anonymity, etc.)
- Small-World Phenomenon
Breaking Information Monopolies

Differences between Meta and P2P Search Engines

Meta Search Engine

small # sites (e.g., digital libraries)

rich statistics about site contents

static federation of servers

each query fully executed
at each site

interconnection topology
largely irrelevant

P2P Search Engine

huge # sites

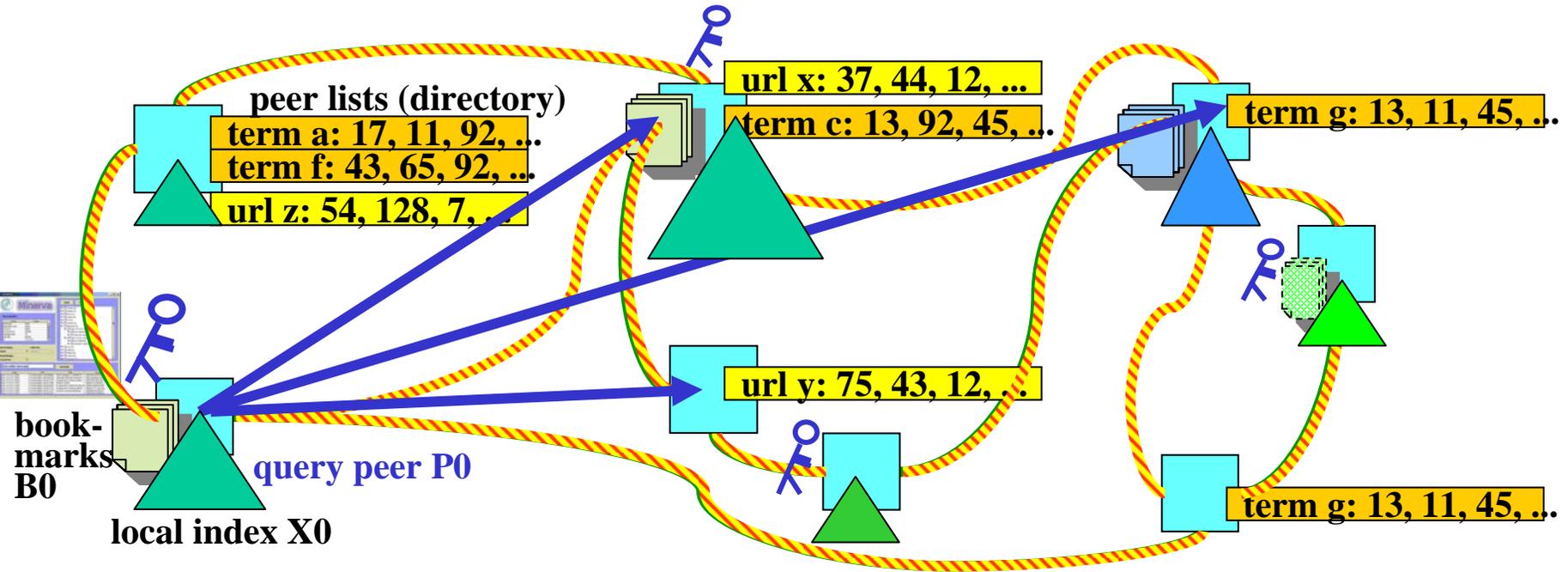
poor/limited/stale summaries

highly dynamic system

single query may need content
from multiple peers

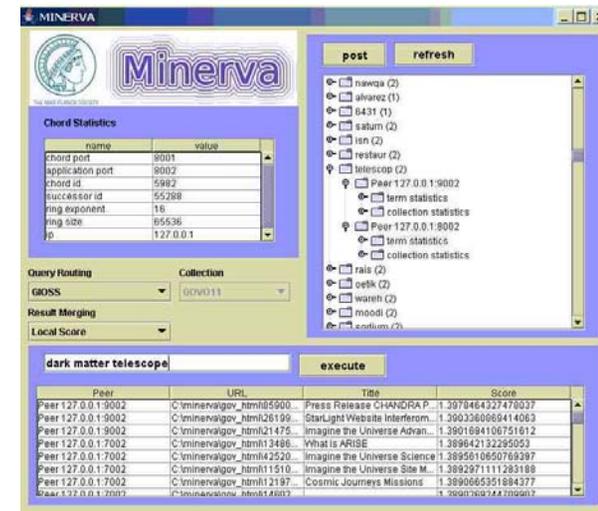
highly dependent on overlay
network structure

P2P Query Routing



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



Framework and Parameters for Query Routing

M terms t_i , N documents d_j , P peers with N_k docs at peer p_k

local measures at peer k:

$tf_i^{(k)}(d)$ – freq. of term i in doc d

$df_i^{(k)}$ - # docs with term i

$idf_i^{(k)}$ - inverse doc freq. of term i

$ttf_i^{(k)}$ – total freq. of term i

$mtf_i^{(k)}$ – max. term freq.

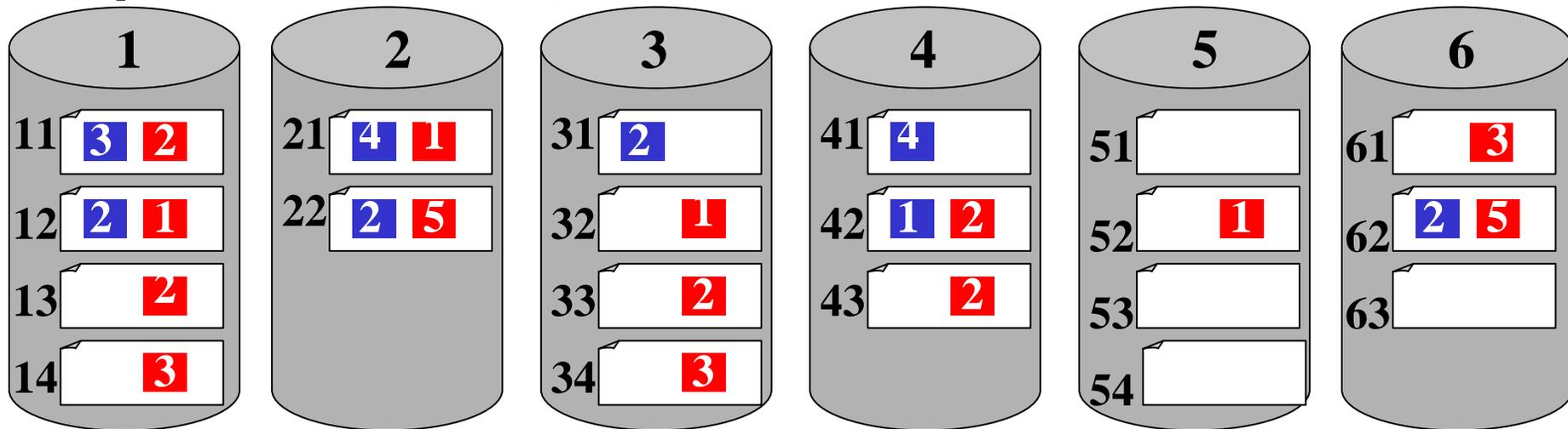
$mdf^{(k)}$ – max. doc freq.

$t^{(k)}$ – # distinct terms

global measures:

$gidf_i$ - inverse doc frequency of term i in P2P system

$gipf_i$ - inverse peer frequency of term i in P2P system



select best peers based on expected **benefit/cost ratio**

Query Routing based on IPF (PlanetP)

Every peer conceptually maintains the global inverse peer frequency (gipf) for each term i :

$$gipf_i = \log \left(1 + \frac{\# \text{ peers}}{\# \text{ peers with term } i} \right)$$

For multi-keyword query q the quality of peer j is:

$$R_j(q) := \sum_{i \in q} gipf_i \cdot \begin{cases} 1 & \text{if peer } j \text{ contains term } i \\ 0 & \text{otherwise} \end{cases}$$

To retrieve top k results for query q :

1. rank peers in descending order of $R_j(q)$
2. contact peers in groups of m in rank order
3. merge results
4. iterate steps 2 and 3 until no peer contributes to top- k result

Example: Query Routing based on IPF

$$\text{gipf}_{\blacksquare} = \log(1 + 6/5)$$

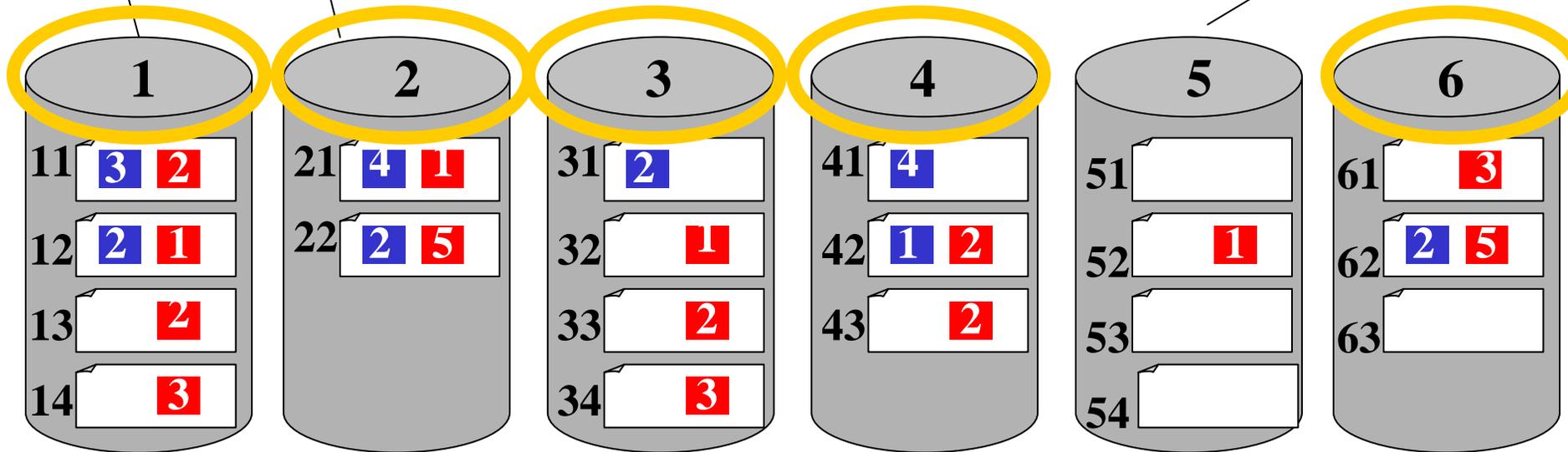
$$\text{gipf}_{\blacksquare} = \log(1+1)$$

$$R_1(\blacksquare \blacksquare) = \log 11/5 * 1 + \log 2 * 1$$

$$R_2(\blacksquare \blacksquare) = \log 11/5 * 1 + \log 2 * 1$$

...

$$R_5(\blacksquare \blacksquare) = \log 2 * 1$$



PlanetP Implementation

Each peer posts its summary in the form of a

Bloom-filter signature:

- bit vector $S[1..s]$ of fixed length s , initially all bits zero
- if peer j has term i it sets bit $h(i)$ to one using a hash function h
- other peers can test if peer j holds term set $\{q_1, \dots, q_k\}$ by looking up $S[h(q_1)], \dots, S[h(q_k)]$ or by computing a bit vector $Q[1..s]$ for $\{q_1, \dots, q_k\}$ and ANDing S with Q , both with the risk of „*false positives*“

Summaries are sent to other peers by asynchronous *gossiping* in a combined push/pull mode:

- *push*: periodically send updates of global registry (small Δs) as „rumors“ to randomly chosen neighbors; stop doing so when n consecutive peers already know the update
- (anti-entropy) *pull*: periodically ask randomly chosen neighbor to send an updated summary of the global registry; alternatively ask push-recipient for recent rumors

Query Routing based on Simple Heuristics

Consider df , tff , mtf , $gipf$ as quality measures of a peer

Choose peers in descending order of

$$\sum_{i \in q} \alpha_1 \log df_i^{(k)} + \alpha_2 \log mtf_i^{(k)} + \alpha_3 \log tff_i^{(k)} \quad \mathbf{or}$$

$$\sum_{i \in q} gipf_i \cdot \left(\alpha_1 \log df_i^{(k)} + \alpha_2 \log mtf_i^{(k)} + \alpha_3 \log tff_i^{(k)} \right)$$

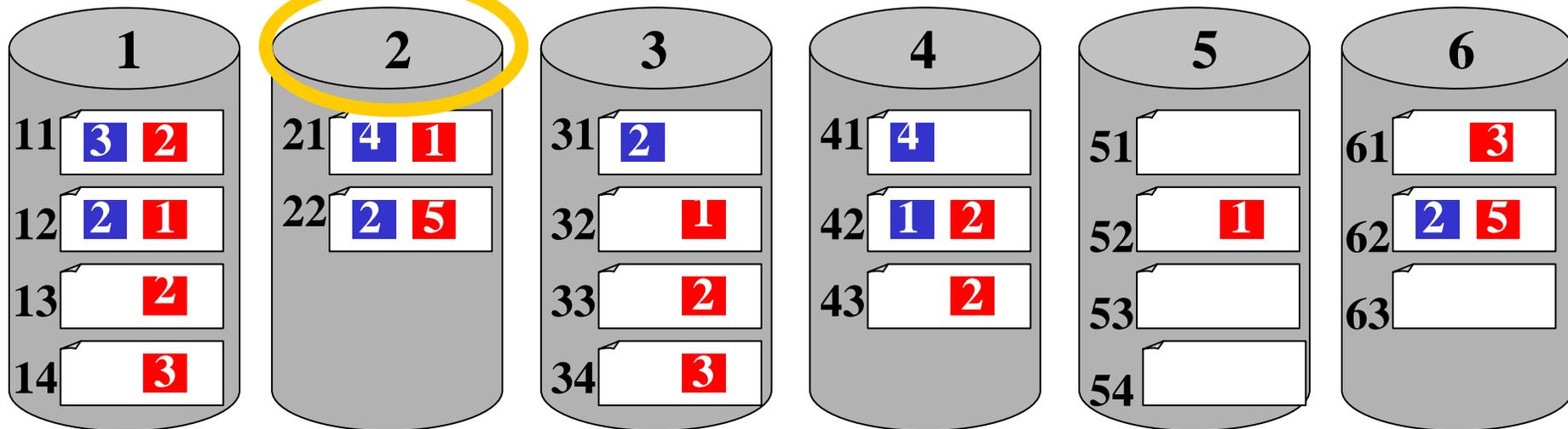
with tunable weights $\alpha_1, \alpha_2, \alpha_3$ such that $\alpha_1 + \alpha_2 + \alpha_3 = 1$

Example: Query Routing Heuristics

$$\text{gipf}_{\text{blue}} = \log(1 + 6/5)$$

$$\text{gipf}_{\text{red}} = \log(1+1)$$

df _{blue} = 2	df _{blue} = 2	df _{blue} = 1	df _{blue} = 2	df _{blue} = 0	df _{blue} = 1
df _{red} = 4	df _{red} = 2	df _{red} = 3	df _{red} = 2	df _{red} = 1	df _{red} = 2
mtf _{blue} = 3	mtf _{blue} = 4	mtf _{blue} = 2	mtf _{blue} = 4	mtf _{blue} = 0	mtf _{blue} = 2
mtf _{red} = 3	mtf _{red} = 5	mtf _{red} = 3	mtf _{red} = 2	mtf _{red} = 1	mtf _{red} = 5
ttf _{blue} = 5	ttf _{blue} = 6	ttf _{blue} = 2	ttf _{blue} = 5	ttf _{blue} = 0	ttf _{blue} = 2
ttf _{red} = 8	ttf _{red} = 6	ttf _{red} = 6	ttf _{red} = 4	ttf _{red} = 1	ttf _{red} = 8



Query Routing based on Statistical Similarity

For query q select peers p with highest value of $\text{sim}(q, p)$, e.g., $\text{cosine}(q, p)$ where p is represented by its centroid

Use **statistical language model** for similarity:

$$KL(q \parallel p) = \sum_{t \in q} P[t|q] \log \frac{P[t|q]}{\lambda P[t|C_p] + (1-\lambda)P[t|G]}$$

where $P[t|q]$, $P[t|C_p]$, $P[t|G]$ are the (estimated) probabilities that term t is generated by the language models for the query q , the corpus C_k of peer k , and the general vocabulary, and λ is a smoothing parameter between 0 and 1

Implementation may estimate $P[t|C_k] \approx \text{tff}_t^{(k)} / \sum_i \text{tff}_i^{(k)}$

The **Kullback-Leibler divergence** (aka. relative entropy) is a measure for the distance between two probability distributions:

$$KL(f \parallel g) := \sum_x f(x) \log \frac{f(x)}{g(x)}$$

CORI Query Routing

Apply probabilistic IR method with heuristic elements (Okapi BM25 term weighting) to peer selection by treating a peer's complete index contents as a „document“

$$P[q|p] \sim$$

$$\sum_{i \in q} \frac{t f_i^{(k)}}{t f_i^{(k)} + 0.5 + 1.5 t^{(k)} / \text{avg}_v(t^{(v)})}$$
$$\cdot \frac{\log(\text{gipf}_t \cdot (0.5 + \# \text{peers}))}{\log(1 + \# \text{peers})}$$

GLOSS Query Routing based on Goodness

Goodness (q, s, l) = $\sum \{\text{sim}(q, d) \mid d \in \text{result}(q, s) \wedge \text{sim}(q, d) > l\}$
for query q, source s, and score threshold l

GLOSS (Glossary Of Servers Server) aims to rank sources by goodness

Approximate goodness by using for source s:

- **$df_i(s)$: number of docs in s that contain term i**
- **$w_i(s)$: $\sum \{\text{tf}_i(d) * \text{idf}_i \mid d \in s\}$ (total weight of term i in s)**

Uniformity assumption:

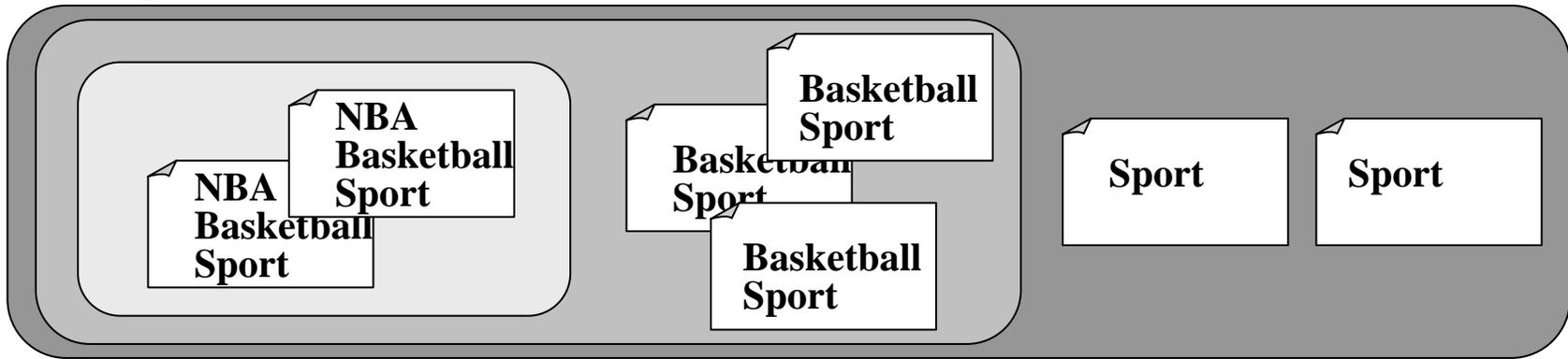
$w_i(s)$ is distributed uniformly over all docs in s that contain i

GLOSS Goodness with High-correlation Assumption

High-correlation assumption:

$df_i(s) \leq df_j(s) \Rightarrow$ every doc in s that contains i also contains j

Example:



For fixed source s and query $q = t_1 \dots t_n$ with $df_i \leq df_{i+1}$ for $i=1..n-1$ consider subqueries $q_p = t_p \dots t_n$ ($p=1..n$).

Every doc d in s that contains only $t_p \dots t_n$ has query similarity

$$sim_p(q, d) = \sum_{i=p..n} t_i \frac{w_i(s)}{df_i(s)}$$

Find p such that $sim_p(q, d) > l$ and $sim_{p+1}(q, d) \leq l$

$$\text{EstGoodness}(q, s, l) = \sum_{i=1..p} (df_i(s) - df_{i-1}(s)) * sim_i$$

GLOSS Goodness with Disjointness Assumption

Disjointness assumption:

$\{d \in s \mid d \text{ contains term } i\} \cap \{d \in s \mid d \text{ contains term } j\} = \emptyset$ for all $i, j \in q$

Uniformity assumption:

$w_i(s)$ is distributed uniformly over all docs in s that contain i

$$sim_i(q, d) = t_i \frac{w_i(s)}{df_i(s)}$$

$$\begin{aligned} \text{EstGoodness}(q, s, l) &= \sum_{i=1..n \wedge sim_i > l} df_i(s) \cdot t_i \frac{w_i(s)}{df_i(s)} \\ &= \sum_{i=1..n \wedge sim_i > l} t_i w_i(s) \end{aligned}$$

Usefulness Estimation Based on MaxSim

Def.: A set S of sources is *optimally ranked* for query q in the order s_1, s_2, \dots, s_m if for every $n > 0$ there exists $k, 0 < k \leq m$, such that s_1, \dots, s_k contain the n best matches to q and each of s_1, \dots, s_k contains at least one of these n matches

Thm.: Let $\text{MaxSim}(q, s) = \max\{\text{sim}(q, d) \mid q \in s\}$.
 s_1, \dots, s_m are optimally ranked for query q if and only if
 $\text{MaxSim}(q, s_1) > \text{MaxSim}(q, s_2) > \dots > \text{MaxSim}(q, s_m)$.

Practical approach („Fast-Similarity method“):

Capture, for each s , $df_i(s)$, $avgw_i(s)$, $maxw_i(s)$ as source summary.
Estimate for query $q = t_1 \dots t_k$

$$\text{MaxSim}(q, s) := \max_{i=1..k} \{t_i * maxw_i(s) + \sum_{v \neq i} t_v * avgw_v(s)\}$$

estimation time linear in query size,
space for statistical summaries linear in #sources * #terms

Overlap Aware Query Routing

First execute q on initiator peer's local index X_0 ,
then select peers P_i with highest **benefit/cost** ratio where

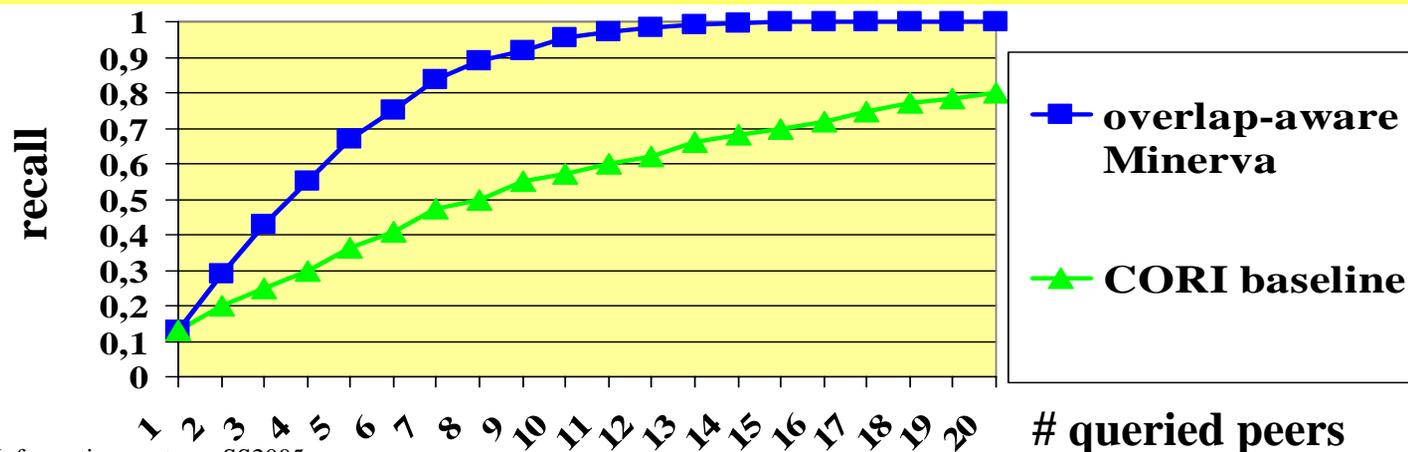
- $\text{benefit}(P_i) \sim \text{sim}(X_0, X_i)$ and $\sim 1/\text{overlap}(X_0, X_i)$
- $\text{cost}(P_i) \sim$ estimated response time or communication costs

precompute sim: $KL(X_0, X_i) := \sum_{\text{terms } x} \text{freq}(x, X_0) \log \frac{\text{freq}(x, X_0)}{\text{freq}(x, X_i)}$

approximate overlap by Bloom filters, hash sketches, etc.

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“



Distributed Query Execution Issues

Algorithm:

- **Determine the number of results to be retrieved from each source a priori based on the source's content quality**

vs.

- **Run distributed version of top-k query processing algorithm**

Dynamic adaptation:

- **Plan query execution only once before initiating it vs.**
- **Dynamic plan adjustment based on sources' result quality and responsiveness (incl. failures)**

Parallelism:

- **Start querying all selected sources in parallel vs.**
- **Consider (initial) results from one source when querying the next sources**

Result Reconciliation

- Case 1: all peers use the same scoring function, e.g. cosine similarities based on tf*idf weights**
- Case 2: peers may use different scoring functions that are publicly known**
- Case 3: peers may use different & unknown scoring functions but provide scored results**
- Case 4: peers provide only result rankings, no scores**

Baseline case:

when **gidf** values are known at the query initiator, we can recompute tf values from the different peers' result docs and compute global scores based on gidf and tf values

Techniques for Result Reconciliation (1)

for case 1:

local sim is

$$lsim(\vec{q}, \vec{d}) = \sum_i \frac{q_i \cdot tf_i(\vec{d}) \cdot lidf_i}{\sqrt{\sum_i q_i^2} \cdot \sqrt{\sum_i tf_i(\vec{d})^2 \cdot lidf_i^2}}$$

global sim is

$$sim(\vec{q}, \vec{d}) = \sum_i \frac{q_i \cdot tf_i(\vec{d}) \cdot gidf_i}{\sqrt{\sum_i q_i^2} \cdot \sqrt{\sum_i tf_i(\vec{d})^2 \cdot gidf_i^2}}$$

either recompute *tf* of result docs and compute *sim*

or submit additional single-term queries (one for each query term) such that each result *d* to the original query *q* is retrieved:

$$lsim(q_i, \vec{d}) = \frac{q_i \cdot tf_i(\vec{d}) \cdot lidf_i}{q_i \cdot \sqrt{\sum_j tf_j(\vec{d})^2 \cdot lidf_j^2}} = \frac{tf_i(\vec{d}) \cdot lidf_i}{\sqrt{\sum_j tf_j(\vec{d})^2 \cdot lidf_j^2}}$$

$$\Rightarrow \frac{lidf_i}{\sqrt{\sum_j tf_j(\vec{d})^2 \cdot lidf_j^2}} = \frac{lsim(q_i, \vec{d})}{tf_i(\vec{d})}$$

**solve equation system
for *tf* and *lidf* values (if possible)
and compute *sim***

Techniques for Result Reconciliation (2)

for case 4:

set global score of doc j retrieved from source i to

$$g(d_j) := 1 - (r_{local}(d_j) - 1) \cdot \frac{r_{min}}{m \cdot r_i} \quad \text{where}$$

- $r_{local}(d_j)$ is the local rank of d_j ,
- r_i is the score of source i among the queried sources,
- r_{min} is the lowest such score, and
- m is the number of desired global results

Intuition:

- initially local ranks are linearly mapped to scores
- the factor $r_{min} / (m r_i)$ is the score difference for consecutive ranks from source i

Wrap-Up: P2P Query Routing & Execution

Research on distributed IR has provided many approaches – principled as well as heuristic ones – that can be carried over to a P2P setting

However, the scale, dynamics, and usage patterns of a P2P search engine entail additional issues, many of which are widely open:

- peer statistics collection and dissemination (frequency, quality, overlap, resource util., response times, etc.)
- precomputation of good peers → „semantic overlay network“
- consideration of strong correlations
- combination with PageRank-style authority etc.
- consideration of P0 query execution and feedback
- coping with tradeoffs in network bandwidth & latency, per-peer resource consumption, search result quality

Literature

- Communications of the ACM, Vol 46, No. 2, Special Section on Peer-to-Peer Computing, Feb. 2003.
- Weiyi Meng, Clement Yu, King-Lup Liu: Building Efficient and Effective Metasearch Engines, ACM Computing Surveys 34(1), 2002
- F.M. Cuenca-Acuna, C. Peery, R.P. Martin, T.D. Nguyen: PlanetP: Using Gossiping to Build Content Addressable Peer-to-Peer Information Sharing Communities, IEEE Symp. on HPDC, 2003
- Jie Lu, Jamie Callan: Content-Based Retrieval in Hybrid Peer-to-Peer Networks, CIKM Conference, 2003.
- Henrik Nottelmann, Norbert Fuhr: Evaluating Different Methods of Estimating Retrieval Quality for Resource Selection, SIGIR 2003
- M. Bender, S. Michel, P. Triantafillou, G. Weikum, C. Zimmer: Improving Collection Selection with Overlap Awareness in P2P Search Engines, SIGIR 2005
- M. Bender, S. Michel, G. Weikum, C. Zimmer: Das MINERVA-Projekt: Datenbankselektion für Peer-to-Peer-Websuche, Informatik Forschung und Entwicklung, 2005

17.3 Skalierbare verteilte Indexierung und Suche

Goals:

Decentralize

data store (MP3 files, Web documents, etc.) or
index (term-doc-score entries) or
directory (statistical info about peers)

across N peers with large N

and provide file-name / keyword lookup

with good properties:

- **scalability:** throughput is proportional to N , response time is independent of N (or grows very slowly with N)
- **efficiency:** overhead should be $O(1)$ or $\leq O(\log N)$
- **load balance:** factor between least loaded and most loaded peer should be $O(\log N)$ or even $O(\log \log N)$
- **robustness to failures:** peers being temporarily down
- **robustness to churn:** peers joining and leaving at high rate
- **self-organization** regarding
data growth, load growth, dynamics of system and load patterns

Structured P2P: Example Chord

Distributed Hash Table (DHT):

map strings (file names, keywords) and numbers (IP addresses) onto very large „cyclic“ key space $0..2^m-1$, the so-called *Chord Ring*

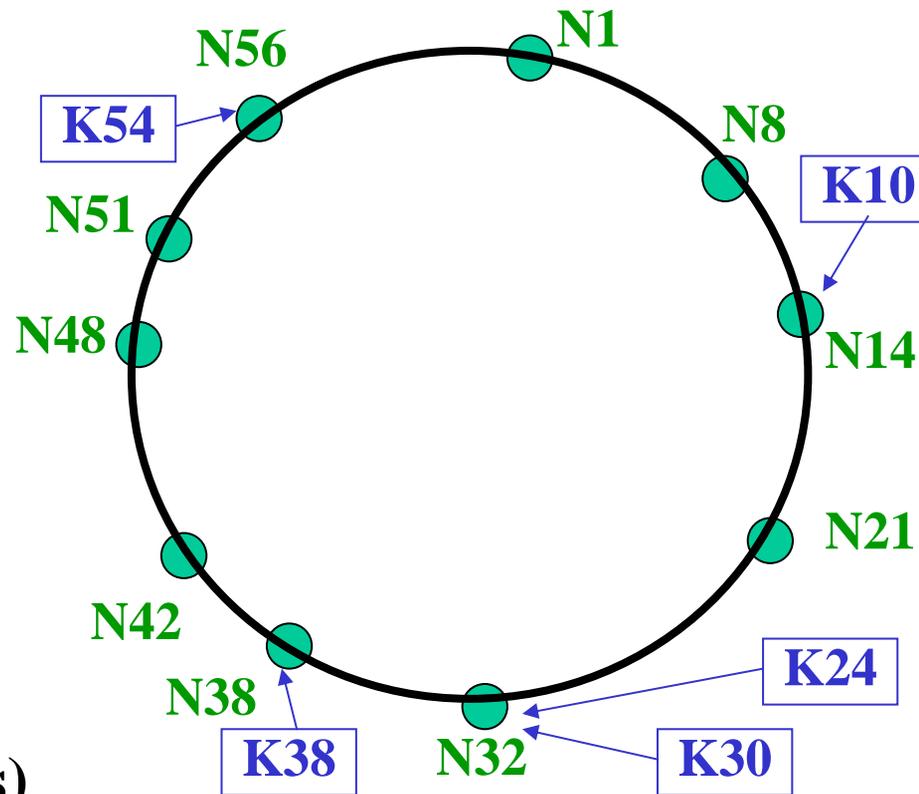
Key k (e.g., *hash(file name)*) is assigned to the node with key n (e.g., *hash(IP address)*) such that $k \leq n$ and there is no node n' with $k \leq n'$ and $n' < n$

Properties:

Unlimited scalability ($> 10^6$ nodes)

$O(\log n)$ hops to target, $O(\log n)$ state per node

Self-stabilization (many failures, high dynamics)



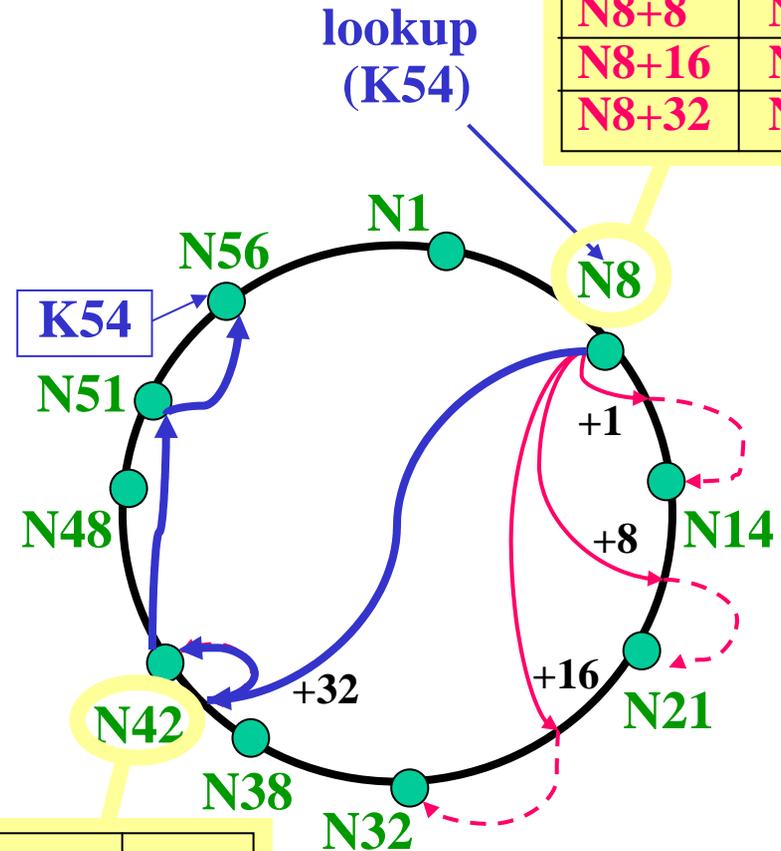
Request Routing in Chord

Every node knows its pred/succ and has a **finger table** with $\log(n)$ pointers: $\text{finger}[i] = \text{successor}(\text{node number} + 2^{i-1})$ for $i=1..m$

Finger table:

N8+1	N14
N8+2	N14
N8+4	N14
N8+8	N21
N8+16	N32
N8+32	N42

For finding key k perform repeatedly:
 determine current node's largest $\text{finger}[i]$ (modulo 2^m) with $\text{finger}[i] \leq k$



N42+1	N48
N42+2	N48
N42+4	N48
N42+8	N51
N42+16	N1
N42+32	N14

pred/succ ring and finger tables require dynamic maintenance
 → stabilization protocol

17.4 Anfrageausführung in verteilten DBS

Beispiel: Netzwerkmonitoring