

Hybrides zusammenfassungs-basiertes P2P Information Retrieval mit partiellem Transfer der Indexdaten

Semantic Routing by Histogramms

Martin Eisenhardt

Lst. für Medieninformatik
Fakultät für Wirtschaftsinformatik und Angewandte Informatik
Otto-Friedrich-Universität Bamberg

Workshop der DBAG
Schloß Dagstuhl, 30. Juni - 02. Juli 2005



Outline

- 1 Einführung und Motivation**
 - P2P Computing
 - Inhaltsbasierte Suche in P2P Netzwerken
 - Herausforderungen
 - Basistechnologie
- 2 Bamberger Ansatz**
 - Basic Idea
 - Generierung der Cluster-Histogramme
 - Retrieval Workflow
 - Index Swapping
- 3 Experimentelle Ergebnisse**
 - Hypothese: Indexdaten liegen auf den Peers geclustert vor
 - Hypothese 2: Indexdaten sind gleichverteilt
- 4 Fazit**
 - Zusammenfassung & Ausblick



Was ist ein Peer?

Dictionary Lookup ...

Peer

- ① A person who has equal standing with another or others, as in rank, class, or age: *children who are easily influenced by their peers.*
- ② Other meanings:
 - ① A nobleman.
 - ② A man who holds a peerage by descent or appointment.
- ③ *Archaic:* A companion; a fellow: "To stray away into these forests drear,/Alone, without a peer." (John Keats).

Webster's Dictionary



Was ist ein Peer?

Genauer: Was versteht man in Netzwerken unter einem Peer?

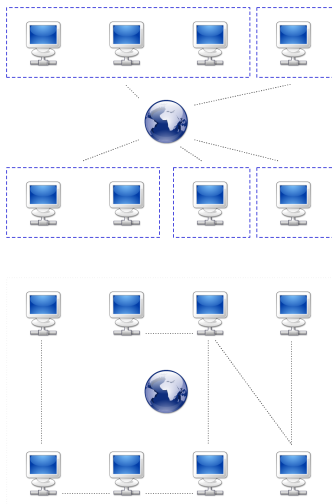
Peer-to-Peer (P2P) computing benötigt **keinerlei zentrale Koordinierung oder Steuerung**.

- Jeder Peer **erbringt Dienste** für das Netzwerk.
- Jeder Peer **nutzt Dienste** anderer Peers.
- Jeder Peer nimmt an der **verteilten Administration** des Netzes teil.

Da also alle beteiligten Entitäten als Gleiche zusammenarbeiten, werden sie **Peers** genannt.



Anatomie eines P2P-Netzwerks



- Die Peers sind über ein **physisches** Netzwerk verbunden.
 - Campus-Netzwerk, Unternehmensnetz, Internet, ...
- Die Peers legen ein logisches **Overlay-Netzwerk** über das physische Netzwerk.
- Die Struktur/Topologie des Overlay-Netzwerks ist für das jeweils gewählte P2P-Netzwerk spezifisch.

Ziel: Inhaltsbasierte Suche in P2P-Netzwerken

- Die ersten P2P-Systeme boten Schlüsselwortsuche auf Dateinamen an.
 - Anwendung in bekannt-berüchtigten P2P-Filesharing-Börsen.
 - *“Ich will **alle** Dateien, bei denen Madonna im Dateinamen enthalten ist.”*
 - Unflexibel, nicht sehr mächtig.
- Heute wird mehr erwartet: **Inhaltsbasierte Ähnlichkeitssuche**.
 - Nutzer sollen Inhalte auf fremden Peers durchsuchen können.
 - Analogon: verteiltes Google.
 - Nicht beschränkt auf Textdokumente, auch Multimedia-Dokumente sollen zugreifbar sein.



Basis für inhaltsbasierte Suche

- Modell: Wie definiert sich der Inhalt eines Dokuments?
 - Welche Teile des Dokuments sind wichtig?
 - Wie beschreibt man ein Dokument?
- Index des Inhalts der Dokumente:
 - Effiziente Datenstrukturen nutzen.
 - Speichern von Informationen über die Dokumente.
- Retrieval:
 - Nutzer stellt Anfrage, die sein Informationsbedürfnis repräsentiert.
 - Zugriff auf Indexdaten, best-matchende Dokumente werden zurückgeliefert.



Herausforderungen der inhaltsbasierten Suche in P2P-Netzen

- Große Mengen an Indexdaten liegen verteilt vor, müssen für die Beantwortung einer Anfrage ausgewertet werden.
- Es gibt keine zentrale Instanz hierfür.
- Wie geben wir Peers Zugriff auf entfernte Indexdaten?
- Zu betrachtende Aspekte:
 - Effizienz
 - Verfügbarkeit
 - Robustheit
 - Skalierbarkeit
 - Sicherheit
 - Anonymität



Zugriff auf entfernte Indexdaten

Verschiedene Ansätze, um Zugriff auf Indexdaten in P2P-Netzen zu gewährleisten:

- Indexdaten am Entstehungsort speichern, Anfragen per Broadcast verteilen.
 - System: sehr frühe File-Sharing-Systeme, z.B. Gnutella
- Jeder Peer ist für einen Teil des Schlüsselraums zuständig und verwaltet Indexdaten, die in diesen Bereich fallen.
 - Indexdaten gehen bei einem *disgraceful leave* verloren.
 - Dadurch können Dokumente nicht mehr gefunden werden, obwohl sie noch im Netz vorhanden sind.
 - Effizient sind nur schlüsselbasierte und Exact-Match-Anfragen möglich.
 - System: z.B. Chord, CAN.



Zugriff auf entfernte Indexdaten — cont'd

- Jeder Peer verwaltet Indexdaten für eigene Dokumente und verteilt kompakte Zusammenfassungen.
 - Verlässt der Peer das Netzwerk, gehen Dokumente *und* zugehörige Dokumente verloren.
 - Nur die **kompakten Zusammenfassungen** werden im Netz verteilt.
 - **Semantic routing** anhand der kompakten Zusammenfassungen wird möglich.
 - System: PlanetP.



PlanetP

PlanetP ist ein an der Rutgers University entwickeltes P2P-Netz zum verteilten Text Retrieval.

- Ein Peer in einem PlanetP-Netzwerk erstellt invertierte Listen der auf ihm vorliegenden Dokumente.
- Zusätzlich erstellt er als kompakte Repräsentation einen Bloom-Filter.
 - Bloom Filter können in äußerst kompakter Form kodieren, ob eine Entität (ein Term) zu einer Menge gehört oder nicht.
- Die Bloom-Filter werden per Rumor Spreading (Gossiping) verteilt.



Rumor Spreading / Gossiping

- Rumor Spreading ist eine Technik zum Verteilen von Informationen in Netzwerken.
- Ein Peer überträgt von Zeit zu Zeit neue Informationen an seine Nachbarn im Netzwerk.
- Zwei Ansätze:
 - Beim **Push-based Rumour Spreading** überträgt ein Peer pro Interval die neuen Informationen an eine von ihm gewählte Untermenge seiner Nachbarn.
 - Beim **Pull-based Rumour Spreading** fragt ein Peer von Zeit zu Zeit einen seiner Nachbarn nach neuen Informationen.
- Push-based Rumour Spreading erzeugt viel Overhead, z.B. werden Informationen an Peers weitergegeben, die diese Information bereits haben.
- Pull-based Rumour Spreading lässt sich weiter optimieren.

PlanetP

Retrieval

Retrieval in einem PlanetP-Netzwerk folgt diesem Workflow:

- 1 Nutzer stellt an seinem Peer p_q die Anfrage.
- 2 p_q kann anhand der lokal vorliegenden Bloom-Filter ermitteln, welche Peers p_i vermutlich relevante Dokumente besitzen.
- 3 Die p_i werden kontaktiert und die Anfrage übertragen, die p_i geben relevante Dokumente zurück.
- 4 Eine Heuristik sorgt dafür, dass nach einer bestimmten Anzahl von Kontakten die Anfrage endet.

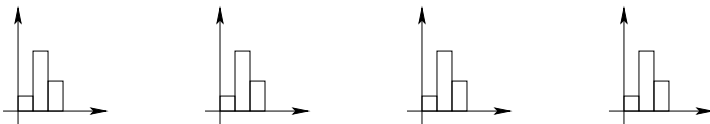
Dieses Verfahren garantiert **nicht**, dass alle relevanten Dokumente im Netzwerk gefunden werden.



Cluster-Histogramme als kompakte Zusammenfassungen

Um multimediales Retrieval in einem PlanetP-artigen P2P-Netzwerk zu ermöglichen, wird eine geeignete kompakte Zusammenfassung benötigt. Unser Ansatz:

- Alle Dokumente im P2P-Netz werden global in Cluster/Klassen eingeteilt.
- Jeder Peer bestimmt dann, wie viele seiner Dokumente in die einzelnen Klassen fallen.
- So können **(Cluster-) Histogramme** der Peers erstellt werden.



Cluster-Histogramme als kompakte Zusammenfassungen

Vorteile von Cluster-Histogrammen:

- Äußerst kompakt, Größe hängt nur von der Anzahl der Klassen ab.
- Anzahl der Dokumente im Netz/pro Peer hat **keinen Einfluss** auf die Histogramm-Größe.
- Einfach zu berechnen, sobald globale Klassifikation vorliegt.

Problem:

Wie berechnet man verteilt in einem P2P-Netz Cluster-Histogramme?

Workflow für die Generierung der Cluster-Histogramme

Schritte für die Cluster-Histogramm-Generation:

- 1 Globale Klassifikation aller Dokumente im Netzwerk.
- 2 Verteilen der Klassifikation an alle Peers.
- 3 Erstellen der Cluster-Histogramme auf den Peers.
- 4 Verteilen der Cluster-Histogramme per Rumour Spreading.



Verteilte Klassifikation von Dokumenten

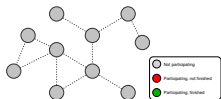
Um eine disjunkte Klassifikation aller Dokumente im Netz zu erhalten, kombinieren wir zwei Verfahren:

- k -means clustering für die globale Dokument-Klassifikation.
 - Verfahren gut geeignet für die Klassifikation von Dokumenten.
 - Inhärent **datunenabhängig**, gut verteilbar.
 - **Kein direkter Zugriff** auf alle Objekte/Dokumente notwendig.
- Probe-Echo-Mechanismus verteilt/sammelt Informationen.
 - Generiert *kostengünstigsten* Spannbaum eines Graphen/Netzwerks.
 - Zwei Phasen: Expansion (PROBE) und Kontraktion (ECHO).

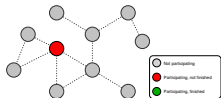


PROBE/ECHO-Mechanismus

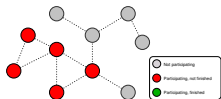
1



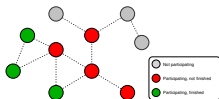
2



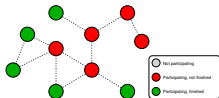
3



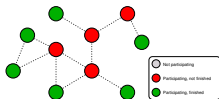
4



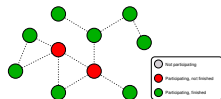
5



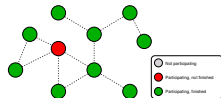
6



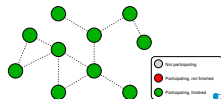
7



8



9



Retrieval Workflow

Im beschriebenen Ansatz kann wie nach Dokumenten gesucht werden:

- Nutzer sucht von seinem Peer p_q aus nach Dokumenten ähnlich dem Anfragedokument d_q .
- d_q wird einem Cluster c_q zugeteilt.
- p_q kann nun anhand der Cluster-Histogramme die anderen Peers p_i ranken.
- Die Peers p_i werden nun anhand des Rankings kontaktiert.
- Heuristik begrenzt Anfragedauer.
- Anfrageergebnisse der einzelnen Peers müssen in ein gemeinsames Ergebnis überführt werden.



Optimierung durch partiellen Transfer von Indexdaten

Überlegung

Für eine Anfrage werden zunächst diejenigen Peers kontaktiert, die im Anfrage-Cluster c_q besonders viele Indexdaten haben.

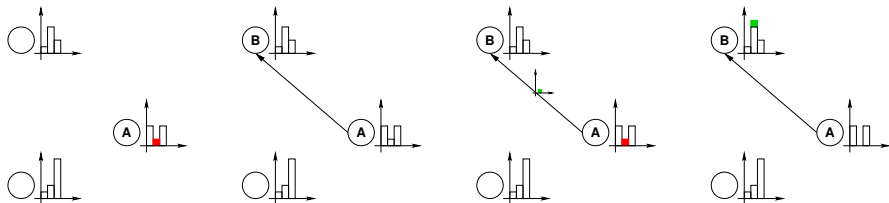
Wie kann man nun die Effizienz des Retrievals steigern?

Index Swapping

Ein Peer p_a soll versuchen, ein Indexdatum idx an einen anderen Peer p_b zu verschieben, wenn

- p_a für den Cluster dieses Indexdatums nur wenige andere Indexdaten hat und
- p_b für diesen Cluster vielen Indexsätze hat.

Index Swapping Illustrated



Schritt 1

Peer A identifiziert einen Cluster, in dem er wenig Expertise besitzt.

Schritt 2

Peer B wird als kompetentester Peer für diesen Cluster identifiziert.

Schritt 3

Peer A überträgt Indexdaten an Peer B. A entfernt diese aus seinem Index, B fügt sie zu seinem Index hinzu.

Resultat

Bei künftigen Anfragen kann Peer B mehr Ergebnisse liefern, es müssen weniger Peers kontaktiert werden.

Hypothese 1: Geclusterte Indexdaten

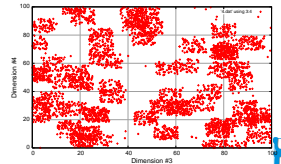
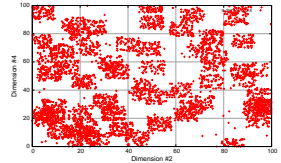
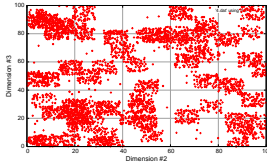
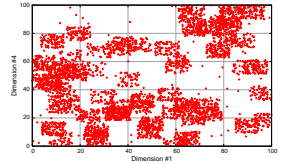
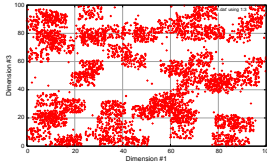
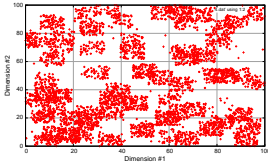
Hypothese

Der Nutzer eines Peers hat gewisse Interessen (z.B. Datenbanken, Gartenbau, ...). Daher liegen auf seinem Peer semantisch eng verwandte Dokumente vor.

- Text Retrieval: Dokumente mit ähnlichen Begriffen.
- Multimedia Retrieval: Foto-Serien ähnlicher Motive.
- Daher werden auch die Indexdaten **innerhalb des Peers homogen** und **zwischen Peers heterogen** sein.
- Somit erhöht Index Swapping die Performance nicht viel, da sich die Peers nicht viel weiter spezialisieren können.



Synthetische geclusterte Indexdaten



Synthetische geclusterte Indexdaten

Mehrere Testkollektionen:

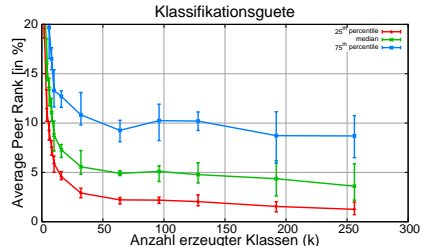
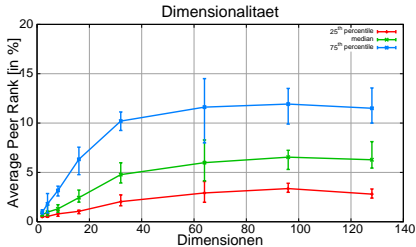
- 100,000 Objekte,
- variiert wird: Dimensionalität, Anzahl von Clustern in den Indexdaten, und Noise.
- Round-Robin-Verteilung auf 1000 Peers:
 - Peers haben ein *center of interest*.
 - Peers ziehen abwechselnd das jeweils relevanteste Dokument aus der Kollektion.
 - Mit und ohne Replikation.

Maß für Performance

Average Peerrank — wieviele Peers müssen kontaktiert werden, um mindestens n der Top- N relevanten Dokumente für eine Anfrage zu finden.



Dimensionality & k



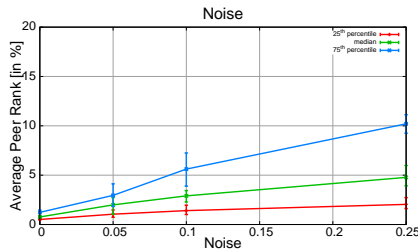
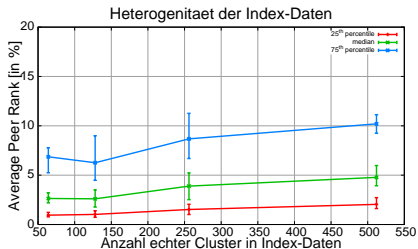
Dimensionalität

Der Einfluss der Dimensionalität auf die Performance ist begrenzt. Für höhere Dimensionalität wächst der durchschnittliche Peerrank nur noch wenig an.

k – Klassifikationsgüte

Größere Genauigkeit bei der globalen Klassifikation der Dokumente bringt bis zu einem gewissen k Vorteile, danach rechtfertigt der Performanzzuwachs den zusätzlichen Aufwand nicht mehr.

Heterogenität der Index-Daten & Noise



Heterogenität der Indexdaten

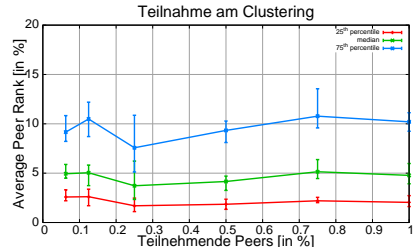
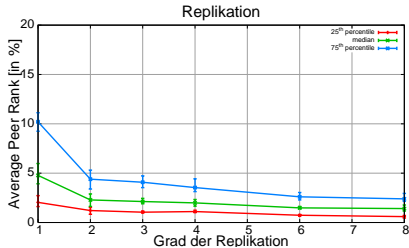
Das Verfahren arbeitet am besten mit möglichst homogenen Daten, höhere Heterogenität der Indexdaten lässt den Aufwand etwas ansteigen, er bleibt aber vertretbar.

Noise

Zwar erschwert Noise das Retrieval, aber auch bei 25% Noise müssen nur 10% der Peers kontaktiert werden, um 15 der Top-20 Dokumente zu finden.



Replikation & Teilnehmer am Clustering



Replikation

Ein geringes Maß an Replikation führt zu einem deutlich verringerten Average Peerrank, höhere Replikationsstufen ergeben keinen deutlichen Performanzzuwachs.

Teilnehmer am Clustering

Nehmen nicht alle Peers an der globalen Klassifikation des Dokumente teil, ändern sich die Ergebnisse nur unwesentlich.

Ergebnis

- Das Verfahren skaliert gut in Abhängigkeit von
 - Dimensionalität
 - Klassifikationsgüte
- Die Heterogenität der Index-Daten und Noise in der globalen Kollektion haben nur begrenzten Einfluss auf die Retrieval-Performance.
- Replikation in Maßen hilft.
- **Wichtig:** Selbst wenn nur wenige Peers an der globalen Klassifikation teilnehmen, bleibt die Retrieval-Performance hoch.
 - Churn im Netzwerk gut beherrschbar.
 - Robust gegen Änderung der Kollektion(en).
 - Neues Clustering nur in großen Intervallen notwendig.
 - Last gut verteilbar, Teilnahme nur jedes n -te Mal.



Hypothese 2: Index-Daten sind nicht-geclustert

Hypothese

Unabhängig von den Nutzerinteressen liegen zumindest die Index-Daten auf den Peers **nicht-geclustert** vor.

- In diesem Fall kann Index Swapping zu Vorteilen beim Retrieval führen, da die Peers schrittweise für eine *künstliche* Clusterung der Index-Daten sorgen.
- Peers spezialisieren sich auf einzelne Cluster.
- Für eine gegebene Anfrage gibt es dann nur noch wenige kompetente Peers.

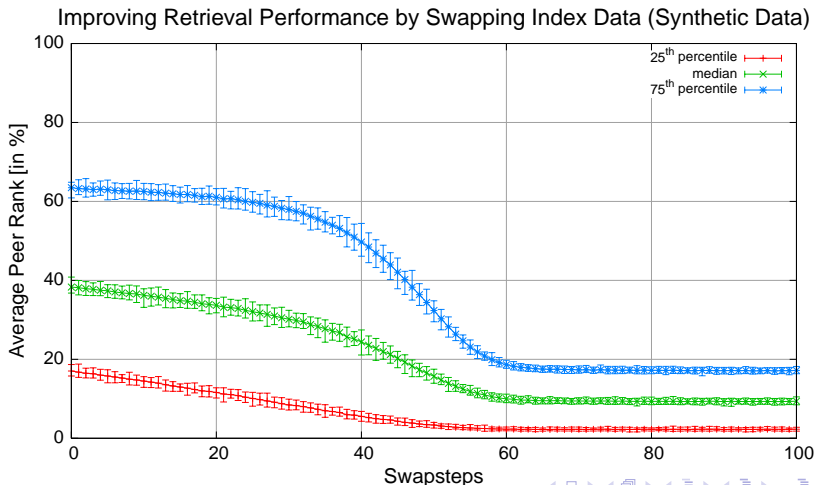


Verwendete Kollektionen

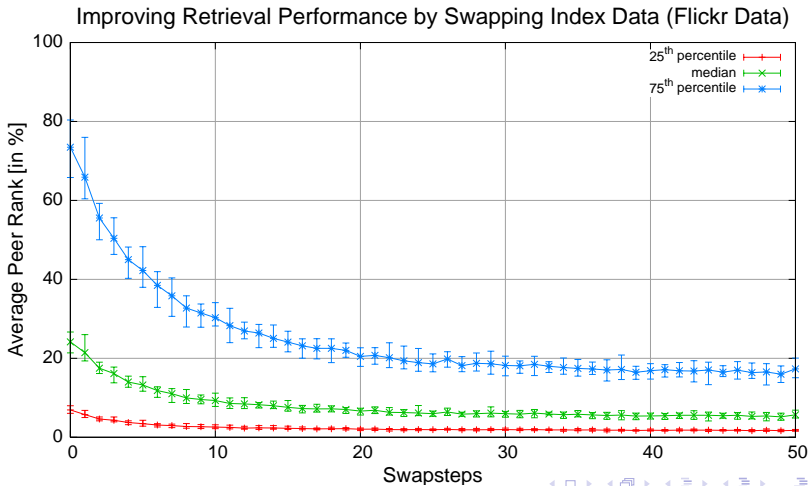
- Synthetische Daten
 - 50.000 Datensätze
 - 32 Dimensionen
 - Gleichverteilung der Indexdaten

- Real-World-Daten
 - 50.000 Bilder von Flickr
 - Flickr ist ein freier Online-Bild-Speicher
 - Tauschen der Bilder erlaubt
 - Zuordnung: ein Flickr-Nutzer entspricht einem Peer

Index Swapping: gleichverteilte, synthetische Daten

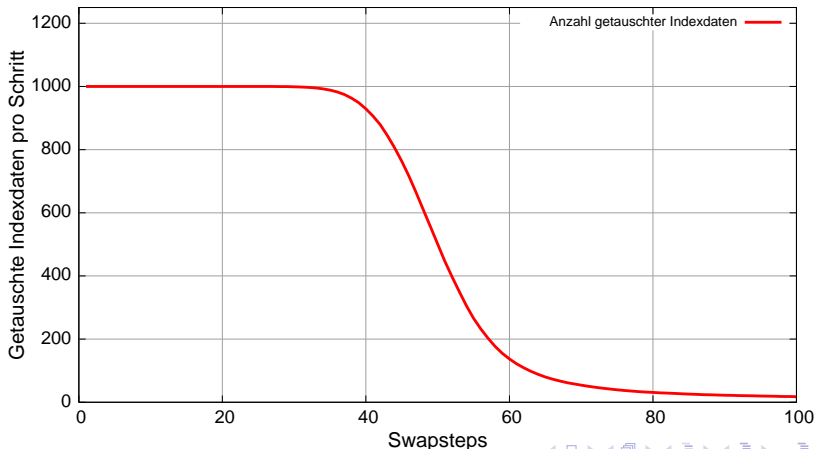


Index Swapping bei Real-World-Daten



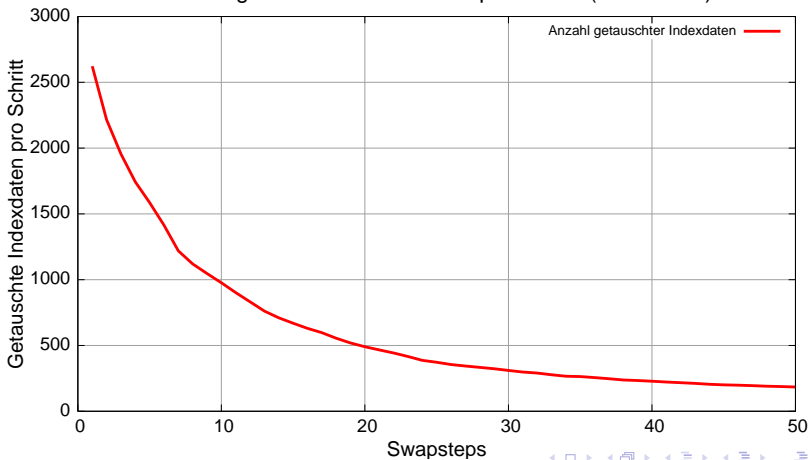
Kosten des Index Swapping (synthetische Daten)

Anzahl getauschter Indexdaten pro Schritt (Synthetische Daten)



Kosten des Index Swapping (Real-World-Daten)

Anzahl getauschter Indexdaten pro Schritt (Flickr Data)



Zusammenfassung & Ausblick

- **Zusammenfassungsbasiertes Semantic Routing** leitet Anfragen effizient zu relevanten Dokumenten.
- **Cluster-Histogramme** führen zu guten Ergebnissen.
- Einfluss von hoher Dimensionalität und Churn ist begrenzt.
- Größere k haben keinen signifikanten Effekt.
- **Index Swapping** führt zu nochmals deutlich erhöhten Performance.

Ausblick:

- Der hier beschriebene Ansatz skaliert nur bis zu einigen hundert/tausend Peers.
- Eine Hierarchisierung des Netzes löst dieses Problem → **Rumorama (Talk von Dr. Wolfgang Müller)**.