Università della Svizzera italiana Faculty of Informatics

### Query Driven Mining of Citation Networks for Patent Citation Retrieval and Recommendation

Parvaz Mahdabi †‡ Fabio Crestani †

University of Lugano, Switzerland † Idiap Research Institute, Martigny, Switzerland ‡

CIKM 2014

### Patent Document

- Patent classifications
- Inventor information
- Title
- Abstract
- Description
- Claims



### Hierarchical Structure of IPC classes H ELECTRICITY

#### H03 BASIC ELECTRONIC CIRCUITRY

H03B GENERATION OF OSCILLATIONS, DIRECTLY OR BY FREQUENCY-CHANGING, BY CIRCUITS EMPLOYING ACTIVE ELEMENTS WHICH OPERATE IN A NON-SWITCHING MANNER; GENERATION OF NOISE BY SUCH CIRCUITS ...

> Modifications of generator to compensate for variations in physical values, e.g. power supply, load, temperature

H03B 5/04

### Patent Retrieval Versus Standard Information Retrieval



### Web Search



### Prior-art Search



# Challenges of Prior-art Search

- A full patent application instead of a keyword query
- Term mismatch
  - Non-standardized acronyms: invented by authors
  - Homonyms: bus (I- motor vehicle, 2- within a computer system)
  - Synonyms: signal and wave
- Incorporating different relevance evidences that come together with a patent application

# Research Questions

- How can we distinguish relevant docs from non-relevant docs
  - using the contextual similarity information diffused over the patent citation network
- Can this information be leveraged for formulating a better query?

# Hypothesis

 Term distribution of influential documents in the citation graph might help to mitigate the term mismatch between the query and the relevant documents



### Related Work

- Textual and categorical similarity
  - Lopez and Romary (CLEF 2010)
  - Magdy and Jones (CLEF 2010)
  - D'hondt et. al (CLEF 2011)
  - Bashir and Rauber (ECIR 2010)
  - Mahdabi et. al (SIGIR 2012, SIGIR 2013)

- Citation Network Analysis
  - Fujii (SIGIR 2007)
- Network analysis derived from interacting companies and inventors
  - Yang et. al (CIKM 2011)
  - Tang et. al (SIGKDD 2012)
  - Wu et. al (WSDM 2013)
- Temporal Citation Analysis
  - Wang et. al (CIKM 2014)
  - Zhang et. al (CIKM 2014)











# Building the Network

- Given a query patent, we retrieve an initial rank list using lexical similarity (root set)
- All documents that cite or are cited by a document in the root set are collected (base set)
- Used EPO web service to extract citation information



### Problem Definition

- Directed weighted graph G = (V, E)
- V is a set of |V| = N patent documents
- A set of citation relationships between patent documents  $E \in V \times V$
- A set of attributes associated to a patent (applicant, inventor, conceptual tags)

$$X = \{x_1, x_2, \dots, x_N\}$$

• Goal: suggest a ranked list of citations for a patent based on its content and attributes

# Initial Query Model

 Estimate a query by quantifying the difference between the language model of the query document and the language model of the collection (the cross entropy)

### Outline

- Introduction and challenges
- Related work and definitions
- Building the network
- Time aware network analysis
- Reranking and query expansion
- Experimental results
- Conclusions

 The initial probability is exponentially discounted according to the age of the node

$$\rho_i = e^{\frac{-\text{age}}{\tau_d}}$$

• Compensating for the bias of the page rank algorithm against recent documents

# Weights for Citation Relationships

- Number of common classification codes
- Number of common inventors
- Number of common applicants
- Lexical similarity (cosine similarity)
- Difference of the publication date
- Combination of different weights

## Weighted PageRank

$$PR - W(u) = \lambda \cdot \frac{O_u}{\sum_{p \in L_v} O_p} + (1 - \lambda) \cdot \sum_{v \in B_u} \frac{PR(v)}{L_v}$$

### Re-ranking and Query Expansion

- Identify and weigh terms in influential documents belonging to the citation graph
- Use weighted page rank values as document prior to guide term selection

$$P(t|Q_{cit}) = Z_t \sum_{D \in \mathcal{O}} P(t|D)P(D)$$

Interpolate citation query model with initial query model

## Experimental Settings

- Dataset
  - two patent collections: CLEF-IP 2010 and CLEF-IP 2011 (relevance judgement: citations)
  - size of each collection: 1.3 million patents
- Evaluation metric
  - PRES: designed for recall-oriented application
- Baseline
  - second best participant of the CLEF-IP 2010 (no bias favoring applicant's initial citation list\*\*)

\* Magdy and Jones: SIGIR 2010 \*\* Magdy, Lopez and Jones: ECIR 2010

### Results on CLEF-IP 2010



AQE-TPR: after QE - weighted temporal PR

### Results on CLEF-IP 2010



AQE-TPR: after QE - weighted temporal PR

### Results on CLEF-IP 2010



# Results on Different Technological fields

MAP @ 1000



- A: Human necessities
- B: Performing operations and transporting
- C: Chemistry and Metallurgy
- **D:** Textiles and Papers
- E: Fixed Constructions
- F: Mechanical Engineering
- G: Physics
- H: Electricity



# Results on Different Technological fields

MAP @ 1000



- A: Human necessities
- B: Performing operations and transporting
- C: Chemistry and Metallurgy
- **D:** Textiles and Papers
- E: Fixed Constructions
- F: Mechanical Engineering
- G: Physics
- H: Electricity



AQE-PR

# Results on Different Technological fields

MAP @ 1000



- A: Human necessities
- B: Performing operations and transporting
- C: Chemistry and Metallurgy
- **D:** Textiles and Papers
- E: Fixed Constructions
- F: Mechanical Engineering
- G: Physics
- H: Electricity



AQE-PR

### Conclusions

- Built a directed weighted graph of citations, encoding contextual similarities as edge weights: common IPC classes, applicants and inventors and temporal distance
- Found influential documents using time aware weighted page rank on the citation network
- Candidate documents are used to draw terms for query expansion
- Lexical and categorical similarities work best for recall, temporal information improves the MAP

Università della Svizzera italiana Faculty of Informatics

### Query Driven Mining of Citation Networks for Patent Citation Retrieval and Recommendation

Parvaz Mahdabi +‡ Fabio Crestani +

University of Lugano, Switzerland † Idiap Research Institute, Martigny, Switzerland ‡

CIKM 2014