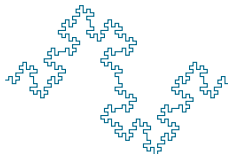


# A Study of Query Reformulation for Patent Prior Art Search with Partial Patent Applications

Mohamed Reda Bouadjenek

**Gabriela Ferraro**

Scott Sanner



# OUTLINE

- Prior art search
- Query reformulation for patents
- Diversification methods for query reformulation
- Experiments
- Results and discussion
- Conclusion and future work

# WHAT PATENTS ARE?

*Patents are legal documents to protect an invention.*

- **Rich meta:** Inventor, Author, Company, Country, Publication year, etc.
- **Predefined document structure:** Title, Abstract, Description and Claims.

**Patent Applications vs. Granted Patents**

# WHAT IS PATENT PRIOR ART SEARCH?

*Finding previously granted patents relevant for a patent application*

- Patent examiners
- Patent authors/inventors

Challenges and data sets:

- ▶ NTCIR (since 2002)
- ▶ TREC-Chem (2007)
- ▶ CLEF-IP (2010/2011)

# PATENT PRIOR ART SEARCH

Why patent prior art search is different to standard Information Retrieval (IR)?

- **Queries** are full patent applications (hundreds of words organized into several sections)
- **Recall-oriented** task (retrieve all relevant documents at early ranks) while text and web search are **precision-oriented** (retrieve a subset of relevant documents)

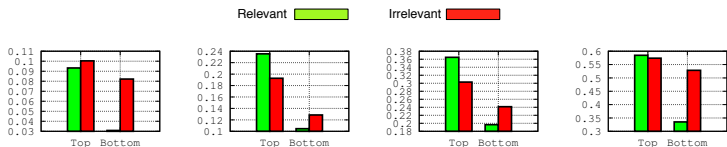
# PATENT PRIOR ART SEARCH WITH PARTIAL APPLICATIONS

*Writing a full patent application is time-consuming and costly*

- We proposed to do patent prior art search with **partial (incomplete) patent applications**

# QUERYING WITH PARTIAL PATENT APPLICATIONS

- ▶ Term overlap (Jaccard Coefficient) of (ir)relevant documents with the result sets for different queries
- ▶ Top 100 and bottom 100 queries
- ▶ Top 10 irrelevant documents ranked by BM25 [?]
- ▶ CLEP-IP 2010



(a) Title query    (b) Abs. query    (c) Claims query    (d) Desc. query

# QUERYING WITH PARTIAL PATENT APPLICATIONS

There are 3 notable trends:

- (i) term overlap increases from *title* to *description* since the query size grows accordingly;
- (ii) the bottom 100 performing queries tend to have much smaller term overlap with the relevant documents than the top 100 queries;
- (iii) the best overlap for any relevant document set for any set of queries is less than one in four terms.

**We investigate Query Reformulation methods**



# QUERY REFORMULATION

*Query reformulation is the process of transforming an initial query  $Q$  to another query  $Q'$ .*

- ▶ **Query Reduction (QR) [?]:** reduces the query such that superfluous information is removed.
- ▶ **Query Expansion (QE) [?]:** enhance the query with additional terms likely to occur in relevant documents.

# QUERY REFORMULATION FOR PATENTS

- ▶ **Query type:** title, abstract, claims, description.

*What part of a partial application an inventor should write to obtain the best search results?*

- ▶ **Relevance model:** BM25, vector space model (VSM):  
TF-IDF [?]

*Which relevance model works best for query reformulation for patent prior art search?*

- ▶ **Query expansion source:** title, abstract, claims, description  
*Are the title words of particularly high value as expansion terms?*

- ▶ **Term selection method:** Rocchio [?], MRRQR

*Which is the best selection method? and with which query type, retrieval model, and term source?*

# QUERY EXPANSION (QE) FRAMEWORKS

*QE aims to alleviate the term mismatch between queries and relevant documents.*

## **Rocchio**

Derives a score for each potential query expansion term and in practice, the top- $k$  scoring terms (often for  $k \ll 200$ ) are used to expand the query and are weighted according to their Rocchio score during the second stage of retrieval.

*What is missed in Rocchio?*

# DIVERSE TERM SELECTION

*Maximal Marginal Relevance* (MMR), a result set diversification algorithm (MMR) [?], usually used for **diverse document selection** (e.g., multi-document summarization)

# NOTATION USED IN MMR QE / QR

		Terms				
		$t_1$	$t_2$	.....	$t_m$	Q
Documents	$d_1$	0.81	0.13	.....	0.28	0.78
	$d_2$	0.11	0.17	.....	0.61	0.51
	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
	$d_n$	0.21	0.1	.....	0.56	0.36

$$t_k^* = \arg \max_{t_k \notin T_{k-1}^*} [\lambda \cos(Q, t_k) - (1 - \lambda) \max_{t_j \in T_{k-1}^*} \cos(t_j, t_k)] \quad (1)$$

# QUERY REDUCTION FRAMEWORKS

*QR aims to short long queries*

We investigate the impact of QR methods when querying with long sections such as *abstract*, *claims* or *description*.

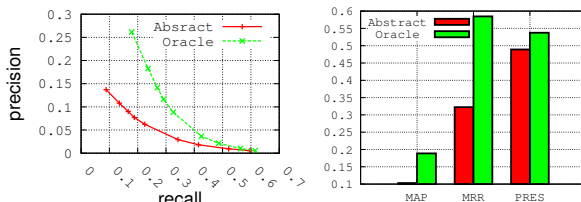


Figure: Sample of terms removed from the abstract section

**MAP:** Mean Average Precision; **MRR:** Mean Reciprocal Rank; **PRES:** Patent Retrieval Evaluation Score

# THE UTILITY OF QUERY REDUCTION FOR 1304 ABSTRACT QUERIES OF THE CLEF-IP 2010 DATASET

## Topic: PAC-1019

**Abstract:** A 5-aminolevulinic acid salt which is useful in fields of microorganisms, fermentation, animals, medicaments, plants and the like; a process for producing the same; a medical composition comprising the same; and a plant activator composition comprising the same.

Term removed	P@5	P@10	R@10	AP	PRES
composit...	<b>0.600</b>	0.300	0.428	<b>0.360</b>	<b>0.829</b>
activ...	0.400	0.300	0.428	0.277	<b>0.809</b>
anim...	<b>0.600</b>	0.300	0.428	<b>0.345</b>	<b>0.798</b>
produc...	0.400	0.300	0.428	<b>0.286</b>	<b>0.797</b>
ferment...	0.200	0.300	0.428	<b>0.283</b>	<b>0.796</b>
microorgan...	<b>0.600</b>	0.300	0.428	<b>0.333</b>	<b>0.793</b>
compris...	0.400	0.300	0.428	0.271	<b>0.790</b>
medica...	0.400	0.300	0.428	<b>0.297</b>	<b>0.789</b>
medic...	0.400	0.300	0.428	<b>0.297</b>	<b>0.787</b>
field...	0.400	0.300	0.428	<b>0.282</b>	<b>0.782</b>
plant...	0.200	0.200	0.285	0.114	0.774
process...	0.400	0.300	0.428	0.279	0.764
acid...	0.400	0.300	0.428	0.252	0.693
salt...	0.200	0.200	0.285	0.216	0.663
aminolevulin...	0.000	0.100	0.142	0.026	0.352
<b>Baseline</b>	0.400	0.300	0.428	0.280	0.777

# EXPERIMENTS SETUP

- ▶ CLEF-IP 2010:
  - ▶ 2.6 million European patent documents
  - ▶ 1303 English topics (queries)
- ▶ CLEF-IP 2011: 3 million patent documents
  - ▶ 2.6 million European patent documents
  - ▶ 1351 English topics (queries)
- ▶ Lucene IR System
- ▶ LucQE: Rocchio method for Lucene
- ▶ Standard stop-words removal
- ▶ Patent-specific stop-words removal [?]
- ▶ Each patent section is indexed in a separate field
- ▶ Queries target all the fields in the index
- ▶ Filtering using the International patent Classification (IPC) of the queries [?, ?]
- ▶ Evaluation on the top 1000 results



# QUERY EXPANSION BASELINES

- ▶ **General QE method**

- ▶ Rocchio []

- ▶ **Patent specific QE methods**

- ▶ **IPC [?]** used the text definitions of the codes assigned to a patent application as a source for expansion.
  - ▶ **WSynSet [?]** used the probability associated with the SynSet entries as a weight for each expanded term in the query.
  - ▶ **USynSet [?]** used uniform weighting for all synonyms of a given term

\* For all methods, their parameters were fixed to their optimal values, which were estimated using the CLEF-IP training queries.

## OTHER QE METHODS

- ▶ Magdy et al. [?] classic techniques of query expansion: WordNet
- ▶ Bashir et al. [?] with SRF set, used a machine learning approach by picking terms that may have a potential positive impact on the retrieval effectiveness.
- ▶ Verma and Varma [?]: used IPC codes as queries, which are expanded using the citation network.

# PSEUDO RELEVANT FEEDBACK (PRF) SIZE

Effect of PRF set with various numbers of feedback documents on the CLEF-IP 2010 dataset.

Query/Source	Metric	Method	5	10	20
Query: Abstract	MAP BL=0.073	Rocchio	0.074	0.072	0.070
		MMRQE	0.074	0.071	0.071
Source: Claims	PRES BL=0.403	Rocchio	0.409	0.409	0.409
		MMRQE	0.411	0.411	0.410
Query: Claims	MAP BL=0.081	Rocchio	0.083	0.080	0.079
		MMRQE	0.082	0.080	0.080
Source: Claims	PRES BL=0.433	Rocchio	0.443	0.445	0.446
		MMRQE	0.445	0.444	0.442

\* 20 terms are used for query expansion

# EXPERIMENTS FOR QE

Experiments options:

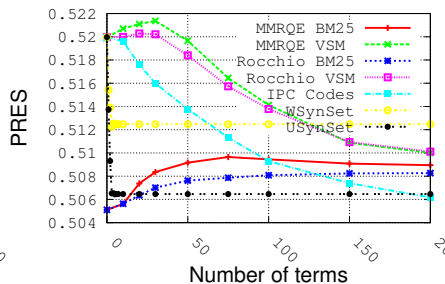
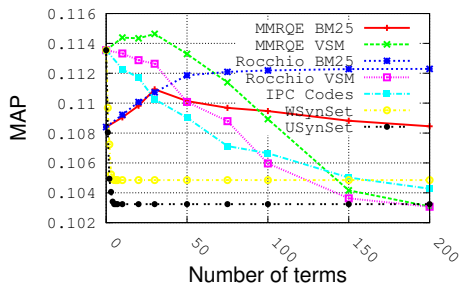
- ▶ **Query type:** {Title, Abstract, Claims, Description}
- ▶ **Query expansion source:** {Title, Abstract, Claims, Description}
- ▶ **Relevance model:** {BM25, Vector-space Model}
- ▶ **Term selection method:** {Rocchio, MMRQE, *etc...*}

## EXPERIEMENTS RESULTS FOR QE

**Query:** Claims

**Date set:** CLEF-IP 2010

**Expansion source:** Abstract



# SAMPLES OF QUERIES (CLEF-IP 2011) WHERE QE IMPROVES THE PERFORMANCE

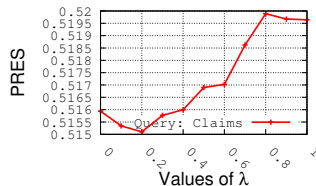
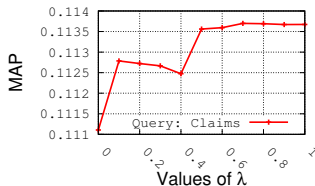
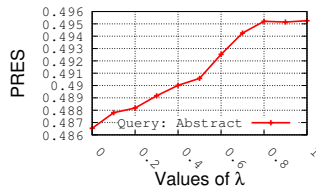
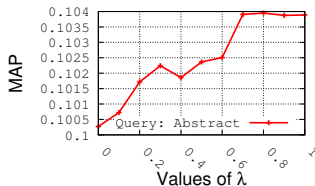
# QUERY REDUCTION BASELINES

- ▶ **General QR method**
  - ▶ **RocchioQR** method for query pruning
- ▶ **Patent specific methods**
  - ▶ **LMQR** [?]: (i) computes Language Modeling similarities by calculating the probability of generating each segment from the top ranked documents; (ii) remove the least similar terms.
  - ▶ **IPC**: (i) rank the terms of the query according to both their frequency in the class code definition, and their frequency in the query. (ii) remove bottom terms of this ranking.

**Other work:** [?] short queries by taking only the first claim of a patent application.

# QR DISCUSSION

Best QR performance results are also obtained when using few documents in the PRF set (top 5). Impact of the diversity parameter  $\lambda$  on the performance of MMRQR on the CLEF-IP 2010 dataset





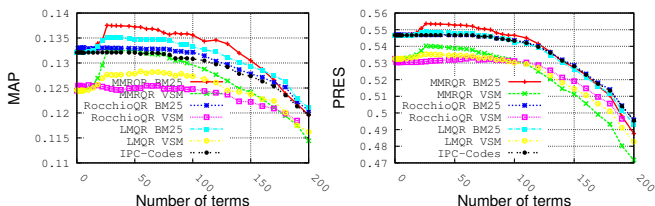
# QR EXPERIMENTS

Experiment options:

- ▶ **Query type:** {Title, Abstract, Claims, Description}
- ▶ **Relevance model:** {BM25, Vector-space Model (VSM)}
- ▶ **Term selection method:** {RocchioQR, MMRQR, *etc...*}

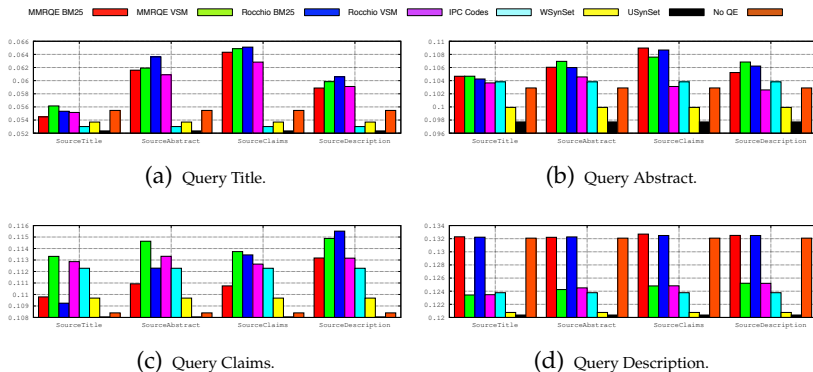
# SAMPLE TABLE!!!

# QR WHILE USING THE DESCRIPTION SECTION FOR QUERYING (CLEF-IP 2010)



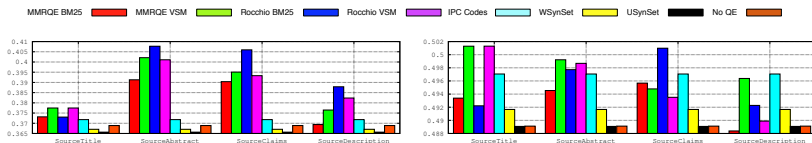
- ▶ (i) for VSM and BM25, MMRQE provides the best performance for MAP and PRES (except for MAP, where Rocchio BM25 provides better performance than MMRQE BM25);
- ▶ (ii) adding more than 50 terms hurts the performance of MMRQE and Rocchio;
- ▶ (iii) exploiting external sources provides poor performance (IPC code definition and SynSets).

# MAP FOR QE METHODS ON CLEF-IP 2010



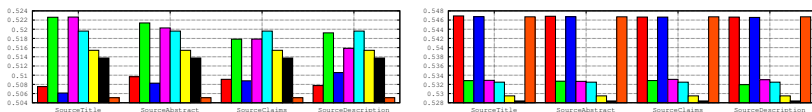
\*for MMRQE  $\lambda = 0.5$

## PRES FOR QE METHODS ON CLEF-IP 2010



(e) Query Title.

(f) Query Abstract.



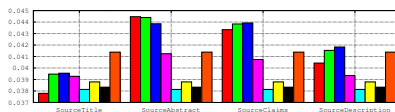
(g) Query Claims.

(h) Query Description.

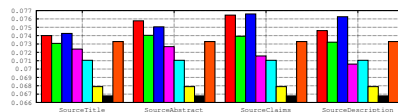
\*for MMRQE  $\lambda = 0.5$

# MAP FOR QE METHODS ON CLEF 2011

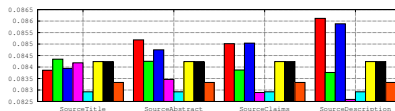
MMRQE BM25 MMRQE VSM Rocchio BM25 Rocchio VSM IPC Codes WSynSet USynSet No QE



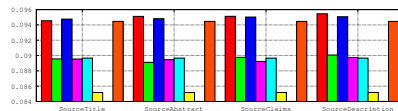
(i) Query Title.



(j) Query Abstract.



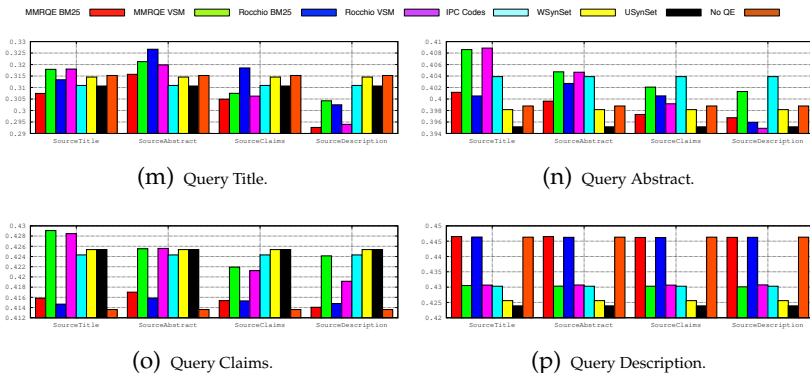
(k) Query Claims.



(l) Query Description.

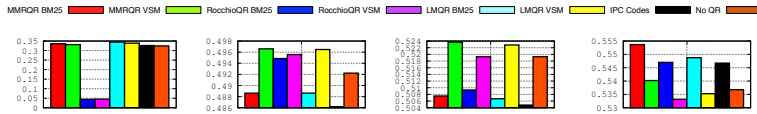
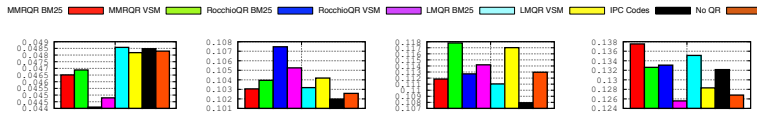
\*for MMRQE  $\lambda = 0.5$

## PRES FOR QE METHODS ON CLEF 2011



\*for MMRQE  $\lambda = 0.5$

# MAP AND PRES FOR QR METHODS ON CLEF 2010



\*for MMRQR  $\lambda = 0.8$



# MAP AND PRES FOR QR METHODS ON CLEF 2011

Contributions are the following:

1. Novel contributions for query expansion and reduction that leverage:
  - ▶ patent structure;
  - ▶ a term diversification technique.
2. A thorough comparative analysis of existing and novel methods for query expansion and reduction in patent prior-art search on standardized datasets of CLEF-IP.

# CONCLUSIONS

- ▶ We analyzed general and specific **QE** and **QR** methods for patent prior art search for partial (incomplete) patent applications (CLEF-IP 2010, 2011);
- ▶ The **claims should be written at early stages of the patent application drafting** (the best section that works with QE/QR (to query with and to use as a source of query expansion/reduction terms)
- ▶ The novel **MMR QE/QR** methods improves results in many cases.

## Future work

Look at more patent-specific methods of and how they can be integrated with methods like MMRQE.