CLEF-IP 2009: retrieval experiments in the Intellectual Property domain

Giovanna Roda Matrixware Vienna, Austria

CLEF 2009 / 30 September - 2 October, 2009

(中) (종) (종) (종) (종) (종)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

¹http://www.clef-campaign.org

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

Previous work on patent retrieval:

¹http://www.clef-campaign.org

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

Previous work on patent retrieval:

 $\bullet~{\rm ACM}~{\rm SIGIR}$ 2000 Workshop

¹http://www.clef-campaign.org

Previous work on patent retrieval:

- $\bullet~{\rm Acm}~{\rm Sigir}$ 2000 Workshop
- N_{TCIR} workshop series since 2001

¹http://www.clef-campaign.org

Previous work on patent retrieval:

- $\bullet~{\rm Acm}~{\rm Sigir}$ 2000 Workshop
- NTCIR workshop series since 2001 Primarily targeting Japanese patents.

¹http://www.clef-campaign.org

Previous work on patent retrieval:

- $\bullet~{\rm Acm}~{\rm Sigir}$ 2000 Workshop
- NTCIR workshop series since 2001 Primarily targeting Japanese patents.
 - ad-hoc task (goal: find patents on a given topic)
 - invalidity search (goal: find patents invalidating a given claim)
 - patent classification according to the F-term system

¹http://www.clef-campaign.org

Legal and economic implications of patent search.

- patents are legal documents
- patent portfolios are assets for enterprises
- a single patent search can be worth several days of work

High recall searches

Missing even a single relevant document can have severe financial and economic impact. For example, when a granted patent becomes invalidated because of a document omitted at application time.

The main task in the $\rm CLEF-IP$ track was to find prior art for a given patent.

▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ ― 臣 … 釣��

The main task in the $\rm CLEF-IP$ track was to find $\it prior \ art$ for a given patent.

Prior art search

Prior art search consists in identifying all information (including non-patent literature) that might be relevant to a patent's claim of novelty.

 $\langle \Box \rangle$

(中) (문) (문) (문) 문

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• before filing a patent application (*novelty search* or *patentability search* to determine whether the invention fulfills the requirements of

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• before filing a patent application (*novelty search* or *patentability search* to determine whether the invention fulfills the requirements of

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

novelty

• before filing a patent application (*novelty search* or *patentability search* to determine whether the invention fulfills the requirements of

(日) (四) (문) (문) (문)

- novelty
- inventive step

- before filing a patent application (*novelty search* or *patentability search* to determine whether the invention fulfills the requirements of
 - novelty
 - inventive step
- before grant results of search constitute the search report attached to patent document

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

- before filing a patent application (*novelty search* or *patentability search* to determine whether the invention fulfills the requirements of
 - novelty
 - inventive step
- before grant results of search constitute the search report attached to patent document
- *invalidity search*: post-grant search used to unveil prior art that invalidates a patent's claims of originality

< 🗆 🕨



• patentese: language used in patents is not natural

▲ロト ▲御ト ▲ヨト ▲ヨト 三ヨ つくぐ

- patentese: language used in patents is not natural
- patents are linked (by citations, applicants, inventors, priorities, ...)

▲ロト ▲御ト ▲ヨト ▲ヨト 三ヨ つくぐ

- patentese: language used in patents is not natural
- patents are linked (by citations, applicants, inventors, priorities, ...)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• available classification information (IPC, ECLA)

- patentese: language used in patents is not natural
- patents are linked (by citations, applicants, inventors, priorities, ...)

< D >

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• available classification information (IPC, ECLA)

Outline

1 Introduction

- Previous work on patent retrieval
- The patent search problem
- CLEF-IP the task

2 The CLEF-IP Patent Test Collection

- Target data
- Topics
- Relevance assessments

3 Participants

- 4 Results
- 5 Lessons Learned and Plans for 2010

(日) (四) (王) (王) (王)

12

6 Epilogue

Outline

Introduction

- Previous work on patent retrieval
- The patent search problem
- CLEF-IP the task

2 The CLEF-IP Patent Test Collection

- Target data
- Topics
- Relevance assessments

3 Participants

4 Results

5 Lessons Learned and Plans for 2010

6 Epilogue

The CLEF-IP Patent Test Collection

The CLEF-IP collection comprises

• target data: 1.9 million patent documents pertaining to 1 million patents (75GB)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

The CLEF-IP Patent Test Collection

The CLEF-IP collection comprises

• target data: 1.9 million patent documents pertaining to 1 million patents (75GB)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• 10,000 topics

The $\operatorname{CLEF-IP}$ collection comprises

- target data: 1.9 million patent documents pertaining to 1 million patents (75GB)
- 10,000 topics
- relevance assessments (with an average of 6.23 relevant documents per topic)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

The $\operatorname{CLEF-IP}$ collection comprises

- target data: 1.9 million patent documents pertaining to 1 million patents (75GB)
- 10,000 topics
- relevance assessments (with an average of 6.23 relevant documents per topic)

Target data and topics are multi-lingual: they contain fields in English, German, and French.

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

The data was provided by Matrixware in a standardized XML format for patent data (the Alexandria XML scheme).

<?xml version="1.0" encoding="UTF-8" ?>

- <patent-document status="new" lang="EN" ucid="EP-0481532-B1" country="EP" doc-number="0481532" kind="B1" date="19951227">

(日) (四) (문) (문) (문)

- + <bibliographic-data>
- + <description lang="EN" status="new">
- + <claims lang="DE" status="new">
- + <claims lang="EN" status="new">
- + <claims lang="FR" status="new">
- + <legal-status status="new">
- </patent-document>

<?xml version="1.0" encoding="UTF-8" ?>

- <patent-document status="new" lang="EN" ucid="EP-0481532-B1" country="EP" doc-number="0481532" kind="B1" date="19951227">
 - + <bibliographic-data>
 - <description lang="EN" status="new">
 - The present invention relates to a semiconductor memory device, more particularly it relates to a nonvolatile memory device constituted by combining a volatile memory cell and a nonvolatile memory cell including a floating gate circuit element.
 - Recently, in a static random access memory device (RAM), a volatile memory cell is combined with a floating gate circuit element to obtain a nonvolatile memory cell which is used to constitute a nonvolatile memory device. In a nonvolatile memory device of this type, the circuit configuration of each memory cell tends to be complex, and so the size of each memory cell tends to be large. However, this tendency leads to degradation in the reliability and integration of the memory device. In view of this problem,

Field: description Language: German English French

- + <bibliographic-data>
- + <description lang="EN" status="new">
- <claims lang="DE" status="new">
 - + <claim num="xx">
 - <claim num="xx">
 - <claim-text>Halbleiter-Speichervorrichtung nach Anspruch 1, bei der die Schreib-Schaltungsanordnung eine erste Schreibschaltung umfaßt, die mit dem anderen Anschluß des Tunnelkondensators (TC*1, TC*2, TC*1, TC*1, TC*1, TC*), TCa) verbunden ist, und eine zweite Schreibschaltung aufweist, die mit dem einen Anschluß des Tunnelkondensators kapazitiv gekoppelt ist.</claim-text>
 - </claim>
 - <claim num="xx">
 - <claim-text>Halbleiter-Speichervorrichtung nach Anspruch 2, bei der die flüchtige Speicherzelle (1) zwei Transistoren (Q₁, Q₂) aufweist, die zueinander kreuzgekoppelt sind, wobei die Spannungen der Gateanschlüsse der zwei Transistoren der nicht flüchtigen Speicherzelle (21) als

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Field: claims

Language: German English French

- + <bibliographic-data>
- + <description lang="EN" status="new">
- + <claims lang="DE" status="new">
- <claims lang="EN" status="new">
 - <claim num="xx">

<claim-text>A semiconductor memory device comprising a volatile memory cell (1, 11, 11') and a nonvolatile memory cell (16, 17, 20, 21, 22, 23) corresponding to the volatile memory cell, said nonvolatile memory cell comprising a transistor (Q₅₁, Q₇₃, Q₇₁, Q₇₃, TM) which has a floating gate and is operable to turn on or off in response to memorised data, a tunnel capacitor (TC₆₁, TC₁₂, TC₁₁, TC₁₁, TC₁₁, TC₃₁, TCa), one electrode of which is connected to the floating gate, and write circuitry coupled to the tunnel capacitor, for supplying a high voltage (VH1, VH2; VH1, VH) as a write voltage to said tunnel capacitor in response to the memorised data of said volatile memory cell: characterised in that: said write circuitry includes first and second high voltage sources (VH1, VH2;

Field: **claims** Language: German **English** French

- + <bibliographic-data>
- + <description lang="EN" status="new">
- + <claims lang="DE" status="new">
- + <claims lang="EN" status="new">
- <claims lang="FR" status="new">
 - <claim num="xx">

<claim-text>Dispositif de mémoire à semiconducteur comprenant une cellule de mémoire volatile (1, 11, 11') et une cellule de mémoire non volatile (16, 17, 20, 21, 22, 23) correspondant à la cellule de mémoire volatile, ladite cellule de mémoire non volatile comprenant un transistor (Q₆₁, Q₇₃, Q₆₁, Q₉₃, TM) qui comporte une grille flottante et qui sert à activer ou désactiver, en réponse à des données mémorisées, un condensateur tunnel (TC₆₁, TC₆₂, TC₇₁, TC₇₁, TC₇₁, TC₃) dont une électrode est connectée à la grille flottante, et un circuit d'écriture couplé au condensateur tunnel, pour appliquer une tension élevée (VH1, VH2; VH1, VH) en tant que tension d'écriture audit condensateur tunnel, en réponse

Field: **claims** Language: German English **French** The task for the $\rm CLEF{-}IP$ track was to find prior art for a given patent.

▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ ― 臣 … 釣��

The task for the $\rm CLEF-IP$ track was to find prior art for a given patent. But:

▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ ― 臣 … 釣��

The task for the ${\rm CLEF-IP}$ track was to find prior art for a given patent. Dut

But:

• patents come in several versions corresponding to the different stages of the patent's life-cycle

The task for the ${\rm CLEF-IP}$ track was to find prior art for a given patent.

But:

- patents come in several versions corresponding to the different stages of the patent's life-cycle
- not all versions of a patent contain all fields

▲□▶ ▲圖▶ ▲필▶ ▲필▶ - 필
How to represent a patent topic?



We assembled a "virtual patent topic" file by



We assembled a "virtual patent topic" file by

• taking the B1 document (granted patent)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

We assembled a "virtual patent topic" file by

- taking the B1 document (granted patent)
- adding missing fields from the most current document where they appeared

▲ロト ▲御ト ▲ヨト ▲ヨト 三ヨ つくぐ

Patents to be used as topics were selected according to the following criteria:

- availability of granted patent
- 2 full text description available
- at least three citations
- at least one highly relevant citation

Introduction

- Previous work on patent retrieval
- The patent search problem
- CLEF-IP the task

2 The CLEF-IP Patent Test Collection

- Target data
- Topics
- Relevance assessments

3 Participants

4 Results

5 Lessons Learned and Plans for 2010

6 Epilogue

Sources of citations:



Sources of citations:

• applicant's disclosure: the USPTO requires applicants to disclose all known relevant publications

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

Sources of citations:

- applicant's disclosure: the USPTO requires applicants to disclose all known relevant publications
- Patent office search report: each patent office will do a search for prior art to judge the novelty of a patent

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Sources of citations:

- applicant's disclosure: the USPTO requires applicants to disclose all known relevant publications
- Patent office search report: each patent office will do a search for prior art to judge the novelty of a patent
- opposition procedures: patents cited to prove that a granted patent is not novel

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Extended citations as relevance assessments



direct citations and their families

Extended citations as relevance assessments



(日) (四) (王) (王) (王)

12

direct citations of family members ...

Extended citations as relevance assessments



... and their families

A patent family consists of patents granted by different patent authorities but related to the same invention.

A patent family consists of patents granted by different patent authorities but related to the same invention.

simple family all family members share the same priority number



A patent family consists of patents granted by different patent authorities but related to the same invention.

simple family all family members share the same *priority number* extended family there are several definitions, in the INPADOC database all documents which are directly or indirectly linked via a priority number belong to the same family

Patent families



Patent documents are linked by priorities

Patent families



Patent documents are linked by priorities

INPADOC family.

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Patent families



Patent documents are linked by priorities

CLEF-IP uses simple families.

Outline

Introduction

- Previous work on patent retrieval
- The patent search problem
- CLEF-IP the task

2 The CLEF-IP Patent Test Collection

- Target data
- Topics
- Relevance assessments

3 Participants

- 4 Results
- 5 Lessons Learned and Plans for 2010

6 Epilogue

Participants



- 15 participants
- 48 runs for the main task
- 10 runs for the language tasks

(日) (四) (王) (王) (王)

臣

Participants

- 1 Tech. Univ. Darmstadt, Dept. of CS, Ubiquitous Knowledge Processing Lab (**DE**)
- 2 Univ. Neuchatel Computer Science (CH)
- 3 Santiago de Compostela Univ. Dept. Electronica y Computacion (**ES**)
- 4 University of Tampere Info Studies (FI)
- 5 Interactive Media and Swedish Institute of Computer Science (SE)
- 6 Geneva Univ. Centre Universitaire d'Informatique (**CH**)
- 7 Glasgow Univ. IR Group Keith (UK)
- 8 Centrum Wiskunde & Informatica Interactive Information Access (**NL**)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- 9 Geneva Univ. Hospitals Service of Medical Informatics (CH)
- 10 Humboldt Univ. Dept. of German Language and Linguistics (DE)
- 11 Dublin City Univ. School of Computing (IE)
- 12 Radboud Univ. Nijmegen Centre for Language Studies & Speech Technologies (NL)
- 13 Hildesheim Univ. Information Systems & Machine Learning Lab (**DE**)
- 14 Technical Univ. Valencia Natural Language Engineering (**ES**)
- 15 Al. I. Cuza University of Iasi Natural Language Processing (**RO**)

A system based on Alfresco² together with a Docasu³ web interface was developed. Main features of this system are:

(□) (@) (E) (E) E

²http://www.alfresco.com/

³http://docasu.sourceforge.net/

A system based on Alfresco² together with a $Docasu^3$ web interface was developed.

◆□>
◆□>
●□>

Main features of this system are:

• user authentication

²http://www.alfresco.com/

³http://docasu.sourceforge.net/

A system based on $\mathsf{Alfresco}^2$ together with a Docasu^3 web interface was developed.

◆□> <@> < ≥> < ≥> < ≥</p>

Main features of this system are:

- user authentication
- run files format checks

²http://www.alfresco.com/

³http://docasu.sourceforge.net/

A system based on $\mathsf{Alfresco}^2$ together with a Docasu^3 web interface was developed.

Main features of this system are:

- user authentication
- run files format checks
- revision control

◆□> <@> < ≥> < ≥> < ≥</p>

< □ >

²http://www.alfresco.com/

³http://docasu.sourceforge.net/

Who contributed

These are the people who contributed to the $\ensuremath{\mathrm{CLEF}}\xspace{-}\ensuremath{\mathrm{IP}}\xspace$ track:

▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ ― 臣 … 釣��

 \bullet the $\rm CLEF{-}IP$ steering committee:

• the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

• the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

 \bullet Helmut Berger who invented the name $\rm CLEF{-}IP$

- the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker
- \bullet Helmut Berger who invented the name $\rm CLEF{-}IP$
- Florina Piroi and Veronika Zenz who walked the walk

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker
- \bullet Helmut Berger who invented the name $\rm CLEF{-}IP$
- Florina Piroi and Veronika Zenz who walked the walk
- the patent experts who helped with advice and with assessment of results

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker
- \bullet Helmut Berger who invented the name $\rm CLEF{-}IP$
- Florina Piroi and Veronika Zenz who walked the walk
- the patent experts who helped with advice and with assessment of results

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• the Soire team

- the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker
- \bullet Helmut Berger who invented the name $\rm CLEF{-}IP$
- Florina Piroi and Veronika Zenz who walked the walk
- the patent experts who helped with advice and with assessment of results
- the Soire team
- Evangelos Kanoulas and Emine Yilmaz for their advice on statistics

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで
These are the people who contributed to the $\ensuremath{\mathrm{CLEF}}\xspace{-}\ensuremath{\mathrm{IP}}\xspace$ track:

- the CLEF-IP steering committee: Gianni Amati, Kalervo Järvelin, Noriko Kando, Mark Sanderson, Henk Thomas, Christa Womser-Hacker
- \bullet Helmut Berger who invented the name $\rm CLEF{-}IP$
- Florina Piroi and Veronika Zenz who walked the walk
- the patent experts who helped with advice and with assessment of results
- the Soire team
- Evangelos Kanoulas and Emine Yilmaz for their advice on statistics

< 🗆 🕨

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• John Tait

Outline

Introduction

- Previous work on patent retrieval
- The patent search problem
- CLEF-IP the task

2 The CLEF-IP Patent Test Collection

- Target data
- Topics
- Relevance assessments

3 Participants

4 Results

5 Lessons Learned and Plans for 2010

6 Epilogue

▲ロト ▲御ト ▲ヨト ▲ヨト 三ヨ つくぐ

Precision, Precision@5, Precision@10, Precision@100

▲ロト ▲御ト ▲ヨト ▲ヨト 三ヨ つくぐ

Precision, Precision@5, Precision@10, Precision@100

• Recall, Recall@5, Recall@10, Recall@100

• Precision, Precision@5, Precision@10, Precision@100

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- Recall, Recall@5, Recall@10, Recall@100
- MAP

- Precision, Precision@5, Precision@10, Precision@100
- Recall, Recall@5, Recall@10, Recall@100
- MAP
- nDCG (with reduction factor given by a logarithm in base 10)

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

Some participants were disappointed by their poor evaluation results as compared to other tracks



How to interpret the results

MAP = 0.02 ?

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 … 釣��

There are two main reasons why evaluation at $\rm CLEF-IP$ yields lower values than other tracks:



There are two main reasons why evaluation at $\rm CLEF-IP$ yields lower values than other tracks:

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

• citations are incomplete sets of relevance assessments

There are two main reasons why evaluation at $\rm CLEF-IP$ yields lower values than other tracks:

- O citations are incomplete sets of relevance assessments
- e target data set is fragmentary, some patents are represented by one single document containing just title and bibliographic references (thus making it practically unfindable)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Still, one can sensibly use evaluation results for comparing runs assuming that

Still, one can sensibly use evaluation results for comparing runs assuming that $% \left({{{\boldsymbol{x}}_{i}}} \right)$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

Incompleteness of citations is distributed uniformly

Still, one can sensibly use evaluation results for comparing runs assuming that

- Incompleteness of citations is distributed uniformly
- Same assumption for unfindable documents in the collection

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

Still, one can sensibly use evaluation results for comparing runs assuming that

- Incompleteness of citations is distributed uniformly
- Same assumption for unfindable documents in the collection

Incompleteness of citations is difficult to check not having a large enough gold standard to refer to.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Still, one can sensibly use evaluation results for comparing runs assuming that

- **1** incompleteness of citations is distributed uniformly
- Same assumption for unfindable documents in the collection

Incompleteness of citations is difficult to check not having a large enough gold standard to refer to.

Second issue: we are thinking about re-evaluating all runs after removing unfindable patents from the collection.

MAP: best run per participant



MAP: best run per participant

Group-ID	Run-ID	MAP	R@100	P@100
humb	1	0.27	0.58	0.03
hcuge	BiTeM	0.11	0.40	0.02
uscom	BM25bt	0.11	0.36	0.02
UTASICS	all-ratf-ipcr	0.11	0.37	0.02
UniNE	strat3	0.10	0.34	0.02
TUD	800noTitle	0.11	0.42	0.02
clefip-dcu	Filtered2	0.09	0.35	0.02
clefip-unige	RUN3	0.09	0.30	0.02
clefip-ug	infdocfreqCosEnglishTerms	0.07	0.24	0.01
cwi	categorybm25	0.07	0.29	0.02
clefip-run	ClaimsBOW	0.05	0.22	0.01
NLEL	MethodA	0.03	0.12	0.01
UAIC	MethodAnew	0.01	0.03	0.00
Hildesheim	MethodAnew	0.00	0.02	0.00

Table: MAP, P@100, R@100 of best run/participant (S)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• 7 patent search professionals



- 7 patent search professionals
- judged in average 264 documents per topics



- 7 patent search professionals
- judged in average 264 documents per topics
- not surprisingly, rankings of systems obtained with this small collection do not agree with rankings obtained with large collection

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- 7 patent search professionals
- judged in average 264 documents per topics
- not surprisingly, rankings of systems obtained with this small collection do not agree with rankings obtained with large collection

Investigations on this smaller collection are ongoing.

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

The rankings of runs obtained with the three sets of topics (S=500, M=1000, XL=10,000) are highly correlated (Kendall's $\tau > 0.9$) suggesting that the three collections are equivalent.



◆□▶ <□▶ <□▶ <□▶ <□▶</p>

As expected, correlation drops when comparing the ranking obtained with the 12 manually assessed topics and the one obtained with the \geq 500 topics sets.



I didn't have time to read the working notes ...



... so I collected all the notes and generated a Wordle



▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

... so I collected all the notes and generated a Wordle



They're about patent retrieval.

イロト イヨト イヨト イ

Refining the Wordle



I ran an Information Extraction algorithm in order to get a more meaningful picture

Refining the Wordle



Refining the Wordle










Humboldt's University working notes





• a layered evaluation model is needed in order to measure the impact of each single factor to retrieval effectiveness

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

• a layered evaluation model is needed in order to measure the impact of each single factor to retrieval effectiveness

• provide images (they are essential elements in chemical or mechanical patents, for instance)

- a layered evaluation model is needed in order to measure the impact of each single factor to retrieval effectiveness
- provide images (they are essential elements in chemical or mechanical patents, for instance)
- investigate query reformulations rather than one query-result set

- a layered evaluation model is needed in order to measure the impact of each single factor to retrieval effectiveness
- provide images (they are essential elements in chemical or mechanical patents, for instance)
- investigate query reformulations rather than one query-result set

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

• extend collection to include other languages

- a layered evaluation model is needed in order to measure the impact of each single factor to retrieval effectiveness
- provide images (they are essential elements in chemical or mechanical patents, for instance)
- investigate query reformulations rather than one query-result set

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

- extend collection to include other languages
- include an annotation task

- a layered evaluation model is needed in order to measure the impact of each single factor to retrieval effectiveness
- provide images (they are essential elements in chemical or mechanical patents, for instance)
- investigate query reformulations rather than one query-result set
- extend collection to include other languages
- include an annotation task
- include a categorization task

< 🗆 🕨

▲ロト ▲圖ト ▲国ト ▲国ト 三国 - のへで

• we have created a large integrated test collection for experimentations in patent retrieval

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

- we have created a large integrated test collection for experimentations in patent retrieval
- \bullet the ${\rm CLEF-IP}$ track had a more than satisfactory participation rate for its first year

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

- we have created a large integrated test collection for experimentations in patent retrieval
- the CLEF-IP track had a more than satisfactory participation rate for its first year
- the right combination of techniques and the exploitation of patent-specific know-how yields best results

▲ロト ▲御ト ▲ヨト ▲ヨト 三ヨ つくぐ

Thank you for your attention.

▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ ― 臣 … 釣��