

Generating Links for Patent Documents: an Automatic Approach using Computational Intelligence

C. M. Souza, M. E. Santos, M. R. G. Meireles

Pontifical Catholic University of Minas Gerais, Brazil

`cinthia.mikaela@sga.pucminas.br`

`matheus.santos.1004060@sga.pucminas.br`

`magali@pucminas.br`

Abstract. Patents are organized into classification systems, which assist offices and users in the process of seeking and retrieving such documents. A wide variety of users use the patent systems and the information contained in these documents. However, patents are complex legal documents with a significant number of technical and descriptive details, which makes it difficult to identify and analyze the information contained in these documents. An automatic link system associated with some of the terms found in the patents would provide quick access to the concepts contained in specific knowledge bases. This work presents results of a project in which the objective is the automatic generation of links in patent documents. The experiments were conducted with four subgroups of the United States Patent and Trademark Office (USPTO), which uses the Cooperative Patent Classification (CPC) system. As the patent documents did not have keywords, the meaningful terms were selected using the algorithm χ^2 , for which the contents of the entire patent document were used. Some keywords with more than one meaning were disambiguated using a specific algorithm, generating a file with useful information used in the experiments. The links were generated based on Wikipedia articles and the USPTO patent database. The use of the patent database as a possible destination for the link is intended to cover cases in which Wikipedia has no articles on certain terms and also to provide an alternative source that may assist readers in understanding those documents. It is expected, with the creation of automated links, to make it easier to access concepts related to the terms presented by the documents and to understand the information disclosed by the inventors.

Categories and Subject Descriptors: H.2.8 [**Database Management**]: Database Applications; I.2.7 [**Artificial Intelligence**]: Natural Language Processing

Keywords: Word Sense Disambiguation, Keywords Extraction, Link creation, Patents

1. INTRODUCTION

Patents are an important knowledge source and, therefore, their analysis has been considered a useful tool for research and for management development. Ouellette [2017] conducted a survey of 832 researchers to assess the importance of patent study. Most researchers in different fields of knowledge have stated that they have found useful information in the documents, but acknowledged that there is room for improvement, particularly regarding the accessibility and understanding of information contained in patents. Many of the interviewees stated that it is possible to find information unavailable in the scientific literature and that patents are an underutilized complement to the dissemination of scientific knowledge.

In general, patent documents are divided into sections, defined by patent offices. The main sections of these documents are, title, abstract, claims and descriptions. Each of these sections have par-

The authors thank the financial support of the Pontifical Catholic University of Minas Gerais, the National Council for Scientific and Technological Development (CNPq, grant 429144/2016-4) and the Foundation for Research Support of the State of Minas Gerais (FAPEMIG, grant APQ 01454-17).

Copyright©2019 Permission to copy without fee all or part of the material printed in JIDM is granted provided that the copies are not made or distributed for commercial advantage, and that notice is given that copying is by permission of the Sociedade Brasileira de Computação.

ticular characteristics. Abstracts are characterized by their complex syntactic structure and generic vocabulary. The claims have more specific information of the invention. Usually, the language used is more formal, and sometimes more obscure, than everyday language. The description section has more distinctive information describing the features of the invention [Codina-Filbà et al. 2017; Durham 2018]. These characteristics make it difficult, in many cases, to consult and understand the information, thus generating an underutilization of patent knowledge [Seneviratne 2018].

The increase in the number of patenting processes along with the need to access the information contained in the patent databases motivate researchers to develop efficient techniques for extracting knowledge from such documents [Khode and Jambhorkar 2017]. One of the factors hindering the process of analyzing and extracting knowledge from the information contained in patents is related to the fact that patents are complex documents with technological details, legal language and exhaustive descriptions [Meireles et al. 2016]. In this context, conventional approaches for information retrieval are difficult to apply. The in-depth study of patents has yet to become more accessible, by identifying potential areas of research for the scientific community and generating useful information in the processes of decision-making in the area of competitive intelligence. The proposal of creating an automatic link generation system, which allows simplified access to the concepts related to the terms presented by patents, solves this problem, beckoning with the possibility of quickly accessing the knowledge base related to the theme proposed by the patent.

The automatic determination of a link includes the identification of possible text fragments that should be associated with knowledge bases. In most cases, keywords are selected and their extraction can be done by supervised or unsupervised methods [Mihalcea and Csomai 2007]. Before selecting the texts that will be associated with certain words, it is necessary to define the context in which the word is being used. This is because several language units have different meanings, which is a common feature in many languages and a problem that needs to be addressed in depth.

This article presents results of a project that pursues the automatic generation of links in patent documents. Reginaldo et al. [2017] state, in their work, that the algorithm χ^2 achieved the best results for keyword extraction in this context. To achieve a more satisfactory performance of the keyword extraction process, all sections of the patent were explored. This destination of the links are generated using two separate databases, the Wikipedia article database and the USPTO patent database.

This work was divided into 5 sections. Section 2 presents the main concepts used, as well as a description of the algorithms implemented in the processes of keyword extraction of links and Word Sense Disambiguation (WSD). Section 3 presents the proposed methodology, the used database and the methodological steps of the work. Sections 4 and 5 show the results, analysis and final considerations.

2. AUTOMATIC LINK GENERATION

The automatic link generation process can be divided into three distinct steps: identifying the source of a link, determining the appropriate document to be associated with the terms selected as the source, and generating the link in the document. In the second step, we must solve the problem of WSD. Given this, this section will be divided into four subsections. The first three sections address the main concepts of the work and the last one presents some related works.

2.1 Keywords Extraction

Keywords are a set of relevant terms that sufficiently describe a given text document. Because it is considered a compact form of document representation, keywords are often used in the task of locating information in large databases. Due to these characteristics, keywords are considered an important resource in text mining tasks, natural language processing and information retrieval [Onan et al. 2016; Duari and Bhatnagar 2019]. In some cases, as in patent documents, the keywords are not defined by

the authors of the document and therefore it is necessary to develop a method or select an algorithm that extracts the words considered significant for the document and can represent it in a system of information retrieval.

The χ^2 algorithm, defined by Manning and Schütze [1999], evaluates the independence between two variables and compares observed and expected values, evaluating how far apart they are. According to Manning and Schütze [1999] the chi-squared statistic sums the difference between observed and expected values in all squares of the table, scaled by the magnitude of the expected values. This algorithm is used to order the words according to their dependence on the patent, so that the greater the score given χ^2 to a word, the greater its dependence on the document. Even if the algorithm accepts that a word is independent, the note given by it to the word is simply added to the ordering in lower positions [Reginaldo et al. 2017]. The χ^2 is defined by Equation 1:

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}, \tag{1}$$

where, O_{ij} is the observed value, which represents the number of occurrences of the word in the document, i is the line and j is the column, and E_{ij} is the expected value calculated using Equation 2 and Table I:

$$E_{ij} = \frac{\sum(column_j) \times \sum(line_i)}{\sum(total)}. \tag{2}$$

Table I. Contingency Table

Words/Documents	P ₁	...	P _m	
Word ₁	Occurrences of the Word ₁ in Document P ₁	...	Occurrences of the Word ₁ in Document P _m	Sum (line)
.
.
Word _n	Occurrences of the Word _n in Document P ₁	...	Occurrences of the Word _n in Document P _m	Sum (line)
	Sum (column)	...	Sum (column)	Sum (Internal cells)

Source: Adapted from Mihalcea and Csomai [2007]

2.2 Word Sense Disambiguation

The techniques of WSD aim to computationally identify the meaning of a word, taking into account the context in which it is inserted. Therefore, given a document with a sequence of words $T = w_1, w_2, \dots, w_n$, this technique aims to give meaning to all or some words of that document. This task can be performed with only one lexical sample or with all the words in the document. Generally, the lexical sample is more used, because a wide coverage of domains is necessary to carry out the sense disambiguation of all the words [Corrêa Jr et al. 2018].

JobimText is a WSD algorithm proposed by Panchenko et al. [2017]. This algorithm receives, as input, the word that will be disambiguated and the context to which it belongs. In addition, it is necessary to define some parameters, such as the model used and the output format. Among the

available models, the ensemble model is the most complete, since it searches for the meaning of the ambiguous word in the inventory of word meanings. If a word is outside this vocabulary, then it is disambiguated by using the super meaning inventory. This template was created from a text corpus which is a combination of Wikipedia, ukWaC, corpus LCC News and Gigaword. For the realization of WSD, this algorithm performs the following steps:

- Extraction of context features computing word and feature similarities;
- Word meaning induction;
- Labeling of clusters with hypernyms and images (hypernym is a word with a broad meaning constituting a category in which words with more specific meanings fall ²);
- Sense disambiguation of words in context based on the induced inventory.

In the end, the algorithm returns the meaning of the word, its hypernyms, the set of related words taken from the dictionary, and a set of phrases to exemplify the meaning of the disambiguated word. The context words are those that co-occur with the word ambiguous destination in the given meaning, and they are also returned with the words related to the disambiguated word, taken from the text itself, and the level of trust of the WSD. The confidence level is a metric that evaluates the WSD result. It is calculated from the extraction of hypernyms. For this, the algorithm ranks the hypernyms using functions that relate the word to the set of words of the cluster and to the hypernym. Figure 1 shows a diagram with the steps of the WSD algorithm.

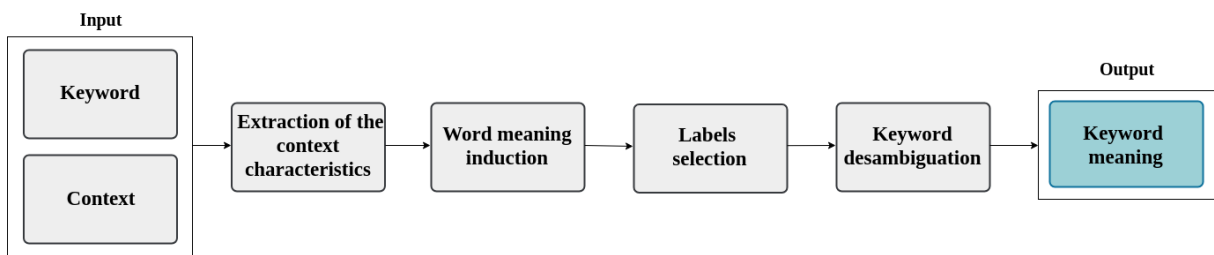


Fig. 1. Steps of the WSD process

2.3 Identification of the link's destination

Identification of the link's destination consists in finding and selecting appropriate documents to be used as the destination for the link [Erbs et al. 2011]. There are different approaches that intend to solve this problem, for example by using the text content, the title of the document and also the links present in the possible destination documents [Seneviratne 2018].

Many papers use the Wikipedia article database as the destination for the link, due to its large number of articles and also the quality of these articles in technical domains. The task of associating terms with Wikipedia articles is called Wikification [Tsunakawa and Kaji 2015]. However, in some cases, it may occur that the Wikipedia database was not sufficient, and it becomes necessary to use other alternative databases.

2.4 Related Works

Linking web data to relevant knowledge base articles has become popular. Some studies on automatic linking of text to important database has captured the interest of the research community [Gardner

²en.oxforddictionaries.com

and Xiong 2009]. The discovery of automatic links has been widely discussed in the literature and has comprehensive solutions in Wikipedia texts and scientific articles [Seneviratne 2018]. The majority of the work is focused in linking Wikipedia texts to their referent Wikipedia pages [Mihalcea and Csomai 2007; Cucerzan 2007; Jana et al. 2017]. Han et al. [2011] and Ratinov et al. [2011] focus on generating links in documents that already have the keywords defined. However, in some documents, keywords are not previously defined, such as in patents. Thus, it is necessary to use techniques capable of extracting these words. Works such as of Mihalcea and Csomai [2007], Erbs, Zesch and Gurevych [2011] and Jana et al. [2017] are described in this subsection, and provide a good perspective on how the automatic link generation task is explored.

Mihalcea and Csomai [2007] used Wikipedia for the automatic extraction of keywords and for the WSD process. The system developed by the authors automatically extracts the keywords, makes the WSD process, and generates the link with the Wikipedia page. For the extraction of the keyword, three methods were tested, tf-idf, χ^2 , and Keyphraseness. For the keyword sense disambiguation, the authors tested two methods. The first one, a knowledge-based approach, and the other one, based on data-driven. In the end, the Wikify! system presented superior results compared to the competitive baselines.

Erbs, Zesch and Gurevych [2011] evaluated, in their work, the performance of the link discovery task, for Wikipedia database, using the content of the text, the title of the document and also the links present in the possible destination documents. In the end, the authors concluded that, in documents that have a large amount of links added manually, the approach based on links has a good result. However, most documents do not have a large amounts of links. In these cases, the approach using text content performs better.

Jana et al. [2017] presented a project to generate links in abstracts of scientific documents with Wikipedia articles. They performed the extraction of the important mentions of the scientific text using tf-idf, together with a set of intelligent filters. Afterwards, for each mention, they extracted a list of candidate entities (Wikipedia links). These entities were classified and punctuated according to their similarity, and finally, based on this score, the entity for link generation was selected. The results show that the methodology used helps to improve the performance of the wikification task in scientific articles.

Although the task of link generation is widely discussed, it is still little explored in the domain of patent knowledge. There are currently several patent databases and search systems that aim to facilitate the search for information in these documents. However, none of them offer features such as Wikipedia's to facilitate understanding of these documents. Moreover, the fact that patents do not have previously defined keywords is a major problem in this task, since there is no great discussion in the literature about the performance of keyword extraction methods in the patent domain. There is a gap between existing technology solutions and desired solutions to make patent knowledge accessible [Seneviratne 2018].

3. METHODOLOGY

This section presents the database created for the experiments and the description of the methodological steps.

3.1 Database

The database used in the experiment is provided by the United States Patent and Trademark Office (USPTO), whose classification system is the Cooperative Patent Classification (CPC). CPC classifies patents into sections, classes, subclasses, groups, and subgroups. For this work, four subgroups, G06K 7/1443, G06K 7/1447, G06K 7/1452 and G06K 7/1456 of the G06K subclass, named "recognition of

data, presentation of data, record carriers, handling record carriers", were randomly selected. Figure 2 illustrates the hierarchical organization of CPC system. The subgroups used in this work are represented by their suffixes, 43, 47, 52, 56, and were highlighted in Figure 2. These subgroups are at the lowest level of the CPC hierarchy.

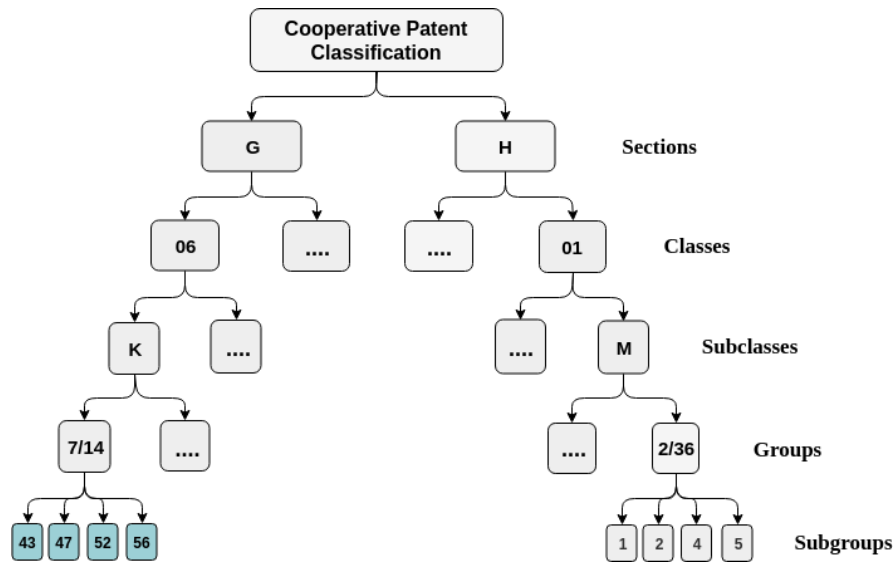


Fig. 2. Hierarchical organization of CPC system

The used database is composed of 999 patents and was updated on April 30, 2019. To validate the methodology in this work, 10 patents of each subgroup were selected. Table II shows the name of the subgroup in the CPC classification system and the distribution of the patents in each of them.

CPC Codes	Number of patents
G06K 7/1443	490
G06K 7/1447	297
G06K 7/1452	82
G06K 7/1456	130

3.2 Methodological Steps

The proposed methodology is presented in Figure 3, showing the main steps of link generation. The first step is the preprocessing of the patent document. The used algorithm performs the removal of stopwords, special characters and the stemization of words. In addition, the algorithm uses a vocabulary based on Wikipedia titles to generate significant n-grams. Thus, only the words which are titles of Wikipedia articles are identified as n-grams, unlike those generated by algorithms that consider only the number of occurrences of two or more words together. This algorithm receives as input a patent document and returns a matrix of occurrences, which contains information on how many times a given word occurs in the document. All sections of the document were used to generate this matrix.

After the preprocessing has been executed, there is an array of documents by words where the occurrence of the words in the document is computed. The second step consists of extracting keywords

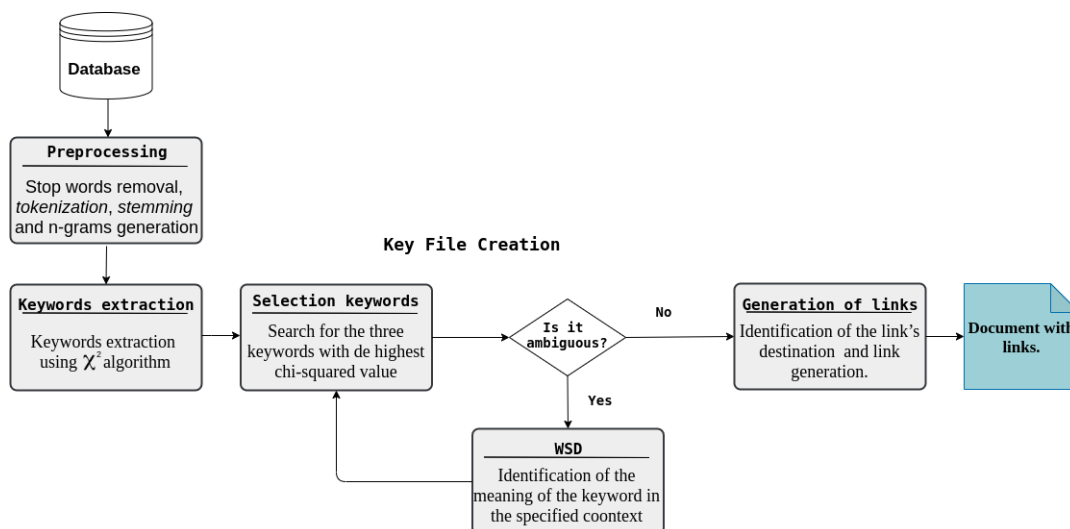


Fig. 3. Proposed methodology

using the algorithm χ^2 . This algorithm receives as input the matrix generated by the preprocessing and returns a list of keywords.

The third step creates the key files for each one of subgroups used in the experiments. Each key file consists of the patent identification code, three keywords with the highest χ^2 index, along with their respective index of χ^2 and the meaning of the keyword if it is ambiguous. For keywords that have ambiguous meaning, the first four hypernyms were ordered and stored according to their reliability index. To identify if the keyword has ambiguous meaning, the Wikipedia API was used. Key files are generated using unprocessed keywords as they are found in the original text. So, after finding the keywords for each preprocessed text, the keywords are selected from the unprocessed text.

The WSD is performed, if necessary, in the fourth step. The used WSD algorithm receives the keyword and a paragraph from the patent text that owns the keyword. Both inputs are not pre-processed. In this case, this paragraph is the context that the algorithm used to find the meaning of the word. For this work, the ensemble model was chosen. This algorithm returns the meaning of the word and four hypernyms. The reliability of the WSD process is then verified with the use of the metric provided by the algorithm. The trust value (0 to 100%) indicates whether the meaning of the keyword and the content of the page selected for the link's destination are associated with the same context.

In the fifth step, the link generation was performed. Two distinct databases were used in this task, the Wikipedia article database, and the USPTO patent database presented in Table II. From the list of keywords generated in the second methodological step, the three keywords with the highest value of χ^2 were extracted for each document. To identify the destination of the link, initially, the implemented algorithm receives a keyword and searches for a Wikipedia page with the corresponding content and returns the page link. This algorithm uses a Python library called "wikipedia", which encapsulates the MediaWiki API¹ for this purpose. This library allows access to Wikipedia data and metadata via API. In order to access this data, the user provides an input data and the algorithm provides a Wikipedia page. However, in some cases, ambiguity in the meaning of the keyword may occur. In this case, the algorithm makes an exception, notifying that there is more than one content-related page. For the treatment of the exception, we use the keyword meaning in the document, present in the key file created in the third methodological step. Using the meaning, the search for the keyword at the base of Wikipedia is carried out again. If a page that has the keyword in the specified meaning is still

¹Application Programming Interface

not found, the algorithm uses the other meaning returned by the WSD algorithm. At the end, if no page is found, the algorithm returns, stating that it was not possible to find a Wikipedia page.

To identify the patent document used as the destination for the link, the implemented algorithm searches, in the key file, the patents that contain the source keyword of the link with the same meaning. In some cases, more than one patent meets this criterion. In this case, the algorithm chooses the patent, in which the selected keyword has the highest value of χ^2 . The algorithm starts from the principle that the best destination for the link is the patent that has the keyword with the same meaning and the highest value of χ^2 . Keywords in one of the four subgroups only links to those in the other three subgroups. As the proposed work is focused on the lowest level of the CPC hierarchy, when creating a link with a document of the same subgroup, the generated link can target a very similar document, not adding value to the explanation of the keyword.

After defining the link destination patent, a code in Python language was implemented using Web Scraping techniques to access the USPTO database and to choose for the link of the selected patent. For the effective creation of the link in the patent document, the keyword is replaced by the URL in HTML code of the page. Figure 4 presents a diagram of the algorithm implemented to link the patent with the two selected database.

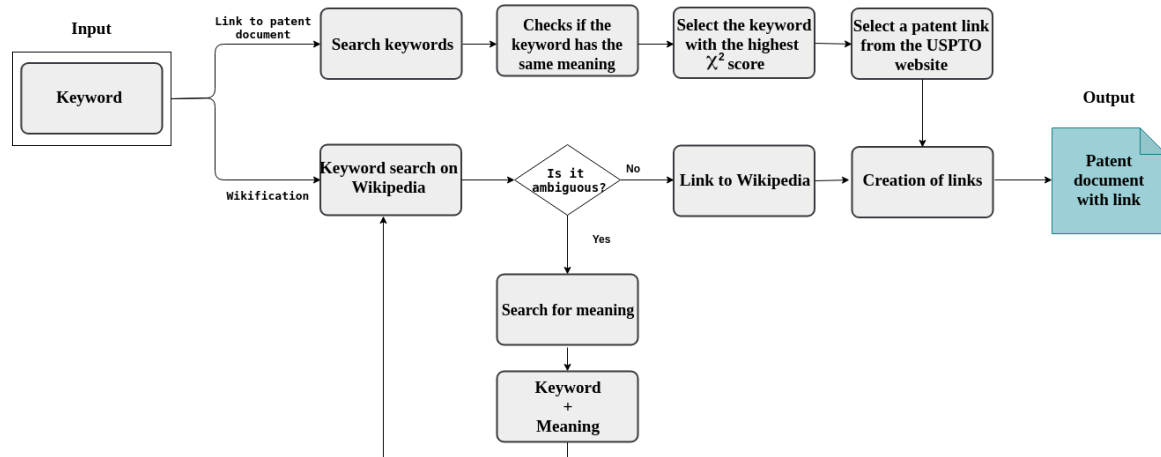


Fig. 4. Identification of the link's destination algorithm steps

4. RESULTS AND DISCUSSION

In the majority of the cases, the proposed methodology was able to correctly identify the Wikipedia pages associated with the keyword origin of the link, when it had no ambiguous meaning. As the vocabulary used for preprocessing is from Wikipedia titles, all extracted keywords represent a page. This makes it easier to correctly identify the destination of the link. However, when the keyword has ambiguous meaning, it is necessary to identify a page that is able to explain the term according to the specified meaning. In this case, the identification of the correct destination for the link is directly associated with the reliability of the WSD algorithm. Therefore, the better the performance of the WSD algorithm, the better is the result obtained by the link destination identification algorithm. It is worth mentioning that identifying the correct link for the keyword is not only related to identifying a page that has the title equal to keyword, but, rather to identifying a page that somehow addresses keyword-related content, giving readers the possibility to understand more clearly the subject matter covered in the patent.

Overall, out of the 40 patents selected for the test database, a total of 120 keywords were extracted. Of these, 80 were identified as unambiguous. Analyzing this group of keywords, it was observed that it is made up of keywords that really have no ambiguity, such as **microcontroller**, **digital asset management**, **greyscale**, **geometric transformation**, **color image sensor**, **image acquisition**, **homography**, **input device**, among others. In addition, 30 keywords were identified as ambiguous. They were **pipe**, **product**, **tag**, **relationship**, **cell**, **filter**, **feature**, among others. The average reliability of sense disambiguation for the 30 keywords was of 76.19%, with values between 41.82% and 100%. Figure 5 shows the amount of keywords distributed in each reliability interval. Of the keywords extracted, 10 were not found in the wikipedia database, showing the need to use another database. In this group of words, we found words such as **dynamic range**, **center line**, **audible**, among others.

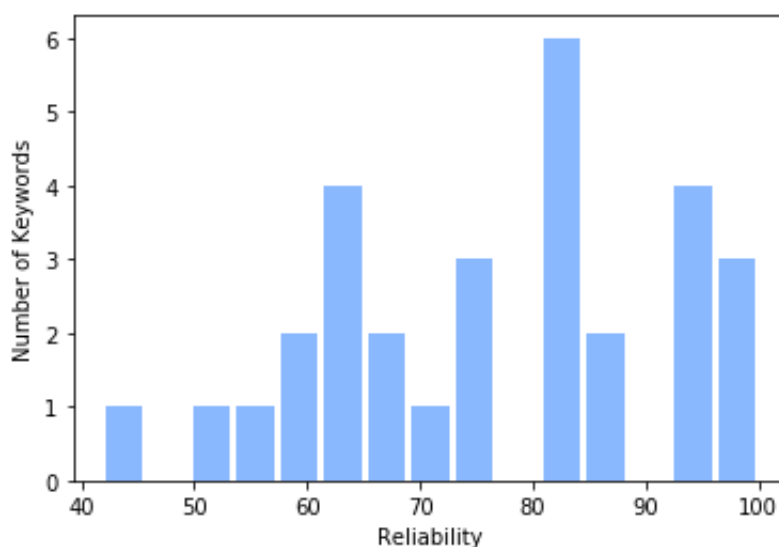


Fig. 5. Reliability of the WSD process

To exemplify the results obtained, we selected four patents that have at least one ambiguous keyword. Tables III, IV, V and VI present some of the results obtained in the WSD step. The first column presents the extracted keywords that have ambiguous meaning, the second one, the hypernyms, the third one, the meaning of the word in the specific context of the patent and the fourth presents the reliability calculated by WSD algorithm. Before each table, the three extracted keywords and the context for sense disambiguation of the selected keywords are presented. The selected patents were identified as P_1 , P_2 , P_3 and P_4 .

Patent P_1

Keywords: compression, **filters**, input buffer.

Context for sense disambiguation of the keyword **filters**:

"An apparatus, and corresponding method for finding an area of interest in an input image, **filters** the input image with at least one compression filter to generate a final compressed image. The apparatus determines the location of an artifact in the final compressed image and then determines the location of the area of interest in the input image according to the location of the artifact in the final compressed image. An apparatus, and corresponding method for finding one or more artifacts in a two-dimensional image, receives a first row of the image and generates a list of regions in accordance with the first row. The apparatus receives a next row of the image and updates the list of regions

in accordance with the next row. The apparatus determines whether a region in the list of regions corresponds to an artifact in the image and, if so, selects the region as the artifact."

Table III. Results of the WSD process for P_1

Keyword	Hypernyms	Meaning	Reliability
filters	device equipment technique application	device	87,35%

Patent P_2

Keyword: **magnitude**, candidate, fft.

Context for sense disambiguation of the keyword **magnitude**:

" 11. The method of claim 1, wherein determining whether the candidate start region includes at least a portion of a candidate oriented encoded signal includes: determining a number of pixels in the candidate start region that have an edge **magnitude** exceeding an edge **magnitude** threshold; and comparing the number of pixels to a pixel count threshold."

Table IV. Results of the WSD process for P_2

Keyword	Hypernyms	Meaning	Reliability
magnitude	factor parameter property condition	factor	74,48%

Patent P_3

Keyword: qr, **registration**, dynamic image.

Context for sense disambiguation of the keyword **registration**:

"The present disclosure relates to systems and methods for decoding intrinsic matrixed bar codes, such as Quick Response ("QR") codes. In one implementation, system for decoding an intrinsic matrixed bar code may include an image-receiving device, a processor configured to execute instructions, and a memory storing the instructions. The instructions may include instructions to: receive an image using the image-receiving device; determine whether the received image contains at least one **registration** mark; when the received image is determined to have at least one **registration** mark: determine, from the at least one **registration** mark, one or more coordinates; when the received image is determined not to have at least one **registration** mark: determine one or more coordinates based on a detected color shifting of a portion of the image; and extract one or more matrixed bar codes overlaid on the received image and located at the one or more coordinates."

Table V. Results of the WSD process for P_3

Keyword	Hypernyms	Meaning	Reliability
registration	service information issue document	service	75,71%

Patent P_4

Keyword: image capturing device, **product**, linear slide.

Context for sense disambiguation of the keyword **product**:

"A system for acquiring multi-angle images of a **product** includes a workstation having a working surface for placing a **product**, a camera supporting member having a vertical axis, and an image capturing device movably attached to the camera supporting member so that it may move along the vertical axis of the camera supporting member. The system captures and analyzes a digital image of a **product** to detect the vertical center of the **product**, and adjusts the position of the image capturing device along the vertical axis so that the vertical center of the **product** is proximate to the vertical center of the image. The system may also have a turntable and additionally rotate the turntable at multiple capturing angles and capture one or more additional digital images of the **product** at various capturing angles and store the one or more additional images in a **product** database."

Table VI. Results of the WSD process for P_4

Keyword	Hypernyms	Meaning	Reliability
product	item		
	industry	item	92,40%
	service		
	material		

By analyzing the presented results, it can be seen that the reliability values found by the WSD algorithm were satisfactory. In general, ambiguous keywords are not necessarily complex terms, but they are of great relevance to the document. When analyzing the keyword and its meaning obtained through the WSD algorithm, it was possible to verify that, in many cases, this word does not have a common meaning. Therefore, it becomes necessary to have a great understanding of the content of the patent to correctly identify the meaning and to evaluate if it is coherent. The possibility of generating links to other databases besides Wikipedia comes with the need to supply the cases in which it is not possible to find a page in Wikipedia that addresses the subject of the keyword in a specific meaning. The experiments showed that the algorithm used to identify the link to the patent database presents a satisfactory result, since it was possible to select a patent capable of explaining the content of the keyword. However, using the key file with only three of the best keywords limited the number of possible patents to be used as the link target. Figure 6 shows the word **metadata** identified, in the patent of subgroup G06K 7/1443 as a keyword. This keyword was linked to another patent from subgroup G06K 7/1447 and to a Wikipedia article.

5. FINAL CONSIDERATIONS

Patents are legal documents with a significant number of technical and descriptive details, which makes their analysis very complex. Access to information in such documents is often laborious, due to the difficulty imposed by technical language and poorly designed writing styles, contrary to the main objective of a patent system of sharing knowledge.

The proposal to create an automatic link generation system is an alternative to provide a simpler access to the knowledge bases related to the content of the patent document. With this proposal, we expect to contribute to the study of links in patents, facilitating the understanding of the information contained in these documents and promoting the dissemination of scientific knowledge associated with the technological advances proposed by these inventions. As a continuation of this work, an user experiment will be conducted to determine which terms in the patents used in the experiments require reference to external knowledge. As users with different types of background may indicate different words to link, these selected words can be considered gold standards to evaluate the performance of our experiments.

United States Patent
Vorabbi, et al.

10,078,774
September 18, 2018

SIMD-based system for multiple decode of captured images

Abstract

A decoding device may include a processor that may include a core component configured to: analyze watermark metadata to identify a watermark ROI from among multiple candidate ROIs in response to generation of the watermark metadata, determine whether rectification is to be performed within the watermark ROI, perform watermark decoding with the rectified watermark ROI data to decode data encoded within a digital watermark within the rectified watermark ROI, and output the decoded data to a server via a network in response to the decoding. The core component may be configured to perform at least one of the operations to determine whether the rectification is to be performed within the watermark ROI, to determine the watermark ROI, or to determine the watermark ROI data to be decoded within the watermark ROI.

Inventors: V...
Applicant: D...
Assignee: D...
Family ID: 63...
Appl. No.: 15...
Filed: A...

United States Patent
Jones, et al.
7,095,871
August 22, 2006

Digital asset management and linking media signals with related data using watermarks

Abstract

A method of performing digital asset management of media content. In this method, a watermark reader device reads a watermark embedded into media content. The watermark conveys watermark information, such as a content identifier and creator identifier. The reader forwards the watermark information to a router. The router then uses the watermark information to find a metadata database identifier. It then sends a request for metadata along with the watermark information to the metadata database identified by the metadata database identifier. The metadata database uses the watermark information to find the metadata for the content. The metadata database uses the watermark information to find the metadata for the content.

Current U.S. Patent No. 10,078,774
Current CPC Class: G06F 16/00
Current International Class: H04N 1/00
Field of Search: 707/100, 707/101, 707/102, 707/103, 707/104, 707/105, 707/106, 707/107, 707/108, 707/109, 707/110, 707/111, 707/112, 707/113, 707/114, 707/115, 707/116, 707/117, 707/118, 707/119, 707/120, 707/121, 707/122, 707/123, 707/124, 707/125, 707/126, 707/127, 707/128, 707/129, 707/130, 707/131, 707/132, 707/133, 707/134, 707/135, 707/136, 707/137, 707/138, 707/139, 707/140, 707/141, 707/142, 707/143, 707/144, 707/145, 707/146, 707/147, 707/148, 707/149, 707/150, 707/151, 707/152, 707/153, 707/154, 707/155, 707/156, 707/157, 707/158, 707/159, 707/160, 707/161, 707/162, 707/163, 707/164, 707/165, 707/166, 707/167, 707/168, 707/169, 707/170, 707/171, 707/172, 707/173, 707/174, 707/175, 707/176, 707/177, 707/178, 707/179, 707/180, 707/181, 707/182, 707/183, 707/184, 707/185, 707/186, 707/187, 707/188, 707/189, 707/190, 707/191, 707/192, 707/193, 707/194, 707/195, 707/196, 707/197, 707/198, 707/199, 707/200, 707/201, 707/202, 707/203, 707/204, 707/205, 707/206, 707/207, 707/208, 707/209, 707/210, 707/211, 707/212, 707/213, 707/214, 707/215, 707/216, 707/217, 707/218, 707/219, 707/220, 707/221, 707/222, 707/223, 707/224, 707/225, 707/226, 707/227, 707/228, 707/229, 707/230, 707/231, 707/232, 707/233, 707/234, 707/235, 707/236, 707/237, 707/238, 707/239, 707/240, 707/241, 707/242, 707/243, 707/244, 707/245, 707/246, 707/247, 707/248, 707/249, 707/250, 707/251, 707/252, 707/253, 707/254, 707/255, 707/256, 707/257, 707/258, 707/259, 707/260, 707/261, 707/262, 707/263, 707/264, 707/265, 707/266, 707/267, 707/268, 707/269, 707/270, 707/271, 707/272, 707/273, 707/274, 707/275, 707/276, 707/277, 707/278, 707/279, 707/280, 707/281, 707/282, 707/283, 707/284, 707/285, 707/286, 707/287, 707/288, 707/289, 707/290, 707/291, 707/292, 707/293, 707/294, 707/295, 707/296, 707/297, 707/298, 707/299, 707/300, 707/301, 707/302, 707/303, 707/304, 707/305, 707/306, 707/307, 707/308, 707/309, 707/310, 707/311, 707/312, 707/313, 707/314, 707/315, 707/316, 707/317, 707/318, 707/319, 707/320, 707/321, 707/322, 707/323, 707/324, 707/325, 707/326, 707/327, 707/328, 707/329, 707/330, 707/331, 707/332, 707/333, 707/334, 707/335, 707/336, 707/337, 707/338, 707/339, 707/340, 707/341, 707/342, 707/343, 707/344, 707/345, 707/346, 707/347, 707/348, 707/349, 707/350, 707/351, 707/352, 707/353, 707/354, 707/355, 707/356, 707/357, 707/358, 707/359, 707/360, 707/361, 707/362, 707/363, 707/364, 707/365, 707/366, 707/367, 707/368, 707/369, 707/370, 707/371, 707/372, 707/373, 707/374, 707/375, 707/376, 707/377, 707/378, 707/379, 707/380, 707/381, 707/382, 707/383, 707/384, 707/385, 707/386, 707/387, 707/388, 707/389, 707/390, 707/391, 707/392, 707/393, 707/394, 707/395, 707/396, 707/397, 707/398, 707/399, 707/400, 707/401, 707/402, 707/403, 707/404, 707/405, 707/406, 707/407, 707/408, 707/409, 707/410, 707/411, 707/412, 707/413, 707/414, 707/415, 707/416, 707/417, 707/418, 707/419, 707/420, 707/421, 707/422, 707/423, 707/424, 707/425, 707/426, 707/427, 707/428, 707/429, 707/430, 707/431, 707/432, 707/433, 707/434, 707/435, 707/436, 707/437, 707/438, 707/439, 707/440, 707/441, 707/442, 707/443, 707/444, 707/445, 707/446, 707/447, 707/448, 707/449, 707/450, 707/451, 707/452, 707/453, 707/454, 707/455, 707/456, 707/457, 707/458, 707/459, 707/460, 707/461, 707/462, 707/463, 707/464, 707/465, 707/466, 707/467, 707/468, 707/469, 707/470, 707/471, 707/472, 707/473, 707/474, 707/475, 707/476, 707/477, 707/478, 707/479, 707/480, 707/481, 707/482, 707/483, 707/484, 707/485, 707/486, 707/487, 707/488, 707/489, 707/490, 707/491, 707/492, 707/493, 707/494, 707/495, 707/496, 707/497, 707/498, 707/499, 707/500, 707/501, 707/502, 707/503, 707/504, 707/505, 707/506, 707/507, 707/508, 707/509, 707/510, 707/511, 707/512, 707/513, 707/514, 707/515, 707/516, 707/517, 707/518, 707/519, 707/520, 707/521, 707/522, 707/523, 707/524, 707/525, 707/526, 707/527, 707/528, 707/529, 707/530, 707/531, 707/532, 707/533, 707/534, 707/535, 707/536, 707/537, 707/538, 707/539, 707/540, 707/541, 707/542, 707/543, 707/544, 707/545, 707/546, 707/547, 707/548, 707/549, 707/550, 707/551, 707/552, 707/553, 707/554, 707/555, 707/556, 707/557, 707/558, 707/559, 707/560, 707/561, 707/562, 707/563, 707/564, 707/565, 707/566, 707/567, 707/568, 707/569, 707/570, 707/571, 707/572, 707/573, 707/574, 707/575, 707/576, 707/577, 707/578, 707/579, 707/580, 707/581, 707/582, 707/583, 707/584, 707/585, 707/586, 707/587, 707/588, 707/589, 707/590, 707/591, 707/592, 707/593, 707/594, 707/595, 707/596, 707/597, 707/598, 707/599, 707/600, 707/601, 707/602, 707/603, 707/604, 707/605, 707/606, 707/607, 707/608, 707/609, 707/610, 707/611, 707/612, 707/613, 707/614, 707/615, 707/616, 707/617, 707/618, 707/619, 707/620, 707/621, 707/622, 707/623, 707/624, 707/625, 707/626, 707/627, 707/628, 707/629, 707/630, 707/631, 707/632, 707/633, 707/634, 707/635, 707/636, 707/637, 707/638, 707/639, 707/640, 707/641, 707/642, 707/643, 707/644, 707/645, 707/646, 707/647, 707/648, 707/649, 707/650, 707/651, 707/652, 707/653, 707/654, 707/655, 707/656, 707/657, 707/658, 707/659, 707/660, 707/661, 707/662, 707/663, 707/664, 707/665, 707/666, 707/667, 707/668, 707/669, 707/670, 707/671, 707/672, 707/673, 707/674, 707/675, 707/676, 707/677, 707/678, 707/679, 707/680, 707/681, 707/682, 707/683, 707/684, 707/685, 707/686, 707/687, 707/688, 707/689, 707/690, 707/691, 707/692, 707/693, 707/694, 707/695, 707/696, 707/697, 707/698, 707/699, 707/700, 707/701, 707/702, 707/703, 707/704, 707/705, 707/706, 707/707, 707/708, 707/709, 707/710, 707/711, 707/712, 707/713, 707/714, 707/715, 707/716, 707/717, 707/718, 707/719, 707/720, 707/721, 707/722, 707/723, 707/724, 707/725, 707/726, 707/727, 707/728, 707/729, 707/730, 707/731, 707/732, 707/733, 707/734, 707/735, 707/736, 707/737, 707/738, 707/739, 707/740, 707/741, 707/742, 707/743, 707/744, 707/745, 707/746, 707/747, 707/748, 707/749, 707/750, 707/751, 707/752, 707/753, 707/754, 707/755, 707/756, 707/757, 707/758, 707/759, 707/760, 707/761, 707/762, 707/763, 707/764, 707/765, 707/766, 707/767, 707/768, 707/769, 707/770, 707/771, 707/772, 707/773, 707/774, 707/775, 707/776, 707/777, 707/778, 707/779, 707/780, 707/781, 707/782, 707/783, 707/784, 707/785, 707/786, 707/787, 707/788, 707/789, 707/790, 707/791, 707/792, 707/793, 707/794, 707/795, 707/796, 707/797, 707/798, 707/799, 707/800, 707/801, 707/802, 707/803, 707/804, 707/805, 707/806, 707/807, 707/808, 707/809, 707/810, 707/811, 707/812, 707/813, 707/814, 707/815, 707/816, 707/817, 707/818, 707/819, 707/820, 707/821, 707/822, 707/823, 707/824, 707/825, 707/826, 707/827, 707/828, 707/829, 707/830, 707/831, 707/832, 707/833, 707/834, 707/835, 707/836, 707/837, 707/838, 707/839, 707/840, 707/841, 707/842, 707/843, 707/844, 707/845, 707/846, 707/847, 707/848, 707/849, 707/850, 707/851, 707/852, 707/853, 707/854, 707/855, 707/856, 707/857, 707/858, 707/859, 707/860, 707/861, 707/862, 707/863, 707/864, 707/865, 707/866, 707/867, 707/868, 707/869, 707/870, 707/871, 707/872, 707/873, 707/874, 707/875, 707/876, 707/877, 707/878, 707/879, 707/880, 707/881, 707/882, 707/883, 707/884, 707/885, 707/886, 707/887, 707/888, 707/889, 707/890, 707/891, 707/892, 707/893, 707/894, 707/895, 707/896, 707/897, 707/898, 707/899, 707/900, 707/901, 707/902, 707/903, 707/904, 707/905, 707/906, 707/907, 707/908, 707/909, 707/910, 707/911, 707/912, 707/913, 707/914, 707/915, 707/916, 707/917, 707/918, 707/919, 707/920, 707/921, 707/922, 707/923, 707/924, 707/925, 707/926, 707/927, 707/928, 707/929, 707/930, 707/931, 707/932, 707/933, 707/934, 707/935, 707/936, 707/937, 707/938, 707/939, 707/940, 707/941, 707/942, 707/943, 707/944, 707/945, 707/946, 707/947, 707/948, 707/949, 707/950, 707/951, 707/952, 707/953, 707/954, 707/955, 707/956, 707/957, 707/958, 707/959, 707/960, 707/961, 707/962, 707/963, 707/964, 707/965, 707/966, 707/967, 707/968, 707/969, 707/970, 707/971, 707/972, 707/973, 707/974, 707/975, 707/976, 707/977, 707/978, 707/979, 707/980, 707/981, 707/982, 707/983, 707/984, 707/985, 707/986, 707/987, 707/988, 707/989, 707/990, 707/991, 707/992, 707/993, 707/994, 707/995, 707/996, 707/997, 707/998, 707/999, 707/1000.

9189670
9542732
2002/0122564

The screenshot shows the Wikipedia article for 'Metadata'. The article text states: 'Metadata is "data [information] that provides information about other data".^[1] Many distinct types of metadata exist, among these **descriptive metadata**, **structural metadata**, **administrative metadata**,^[2] **reference metadata** and **statistical metadata**.^[3]' It also includes a list of bullet points: 'Descriptive metadata describes a resource for purposes such as discovery and identification. It can include elements such as title, abstract, author, and keywords.' and 'Structural metadata is metadata about containers of data and indicates how compound objects are put together.' There is an image of a card catalog and a caption: 'In the 2010s, metadata typically refers to digital forms, but traditional card catalogues contain metadata, with cards holding...'

Fig. 6. Example of a patent with the link to two databases

REFERENCES

CODINA-FILBÀ, J., BOUAYAD-AGHA, N., BURGA, A., CASAMAYOR, G., MILLE, S., MÜLLER, A., SAGGION, H., AND WANNER, L. Using genre-specific features for patent summaries. *Information Processing & Management* 53 (1): 151–174, 2017.

CORRÊA JR, E. A., LOPES, A. A., AND AMANCIO, D. R. Word sense disambiguation: A complex network approach. *Information Sciences* vol. 442-443, pp. 103–113, 2018.

CUCERZAN, S. Large-scale named entity disambiguation based on wikipedia data. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*. pp. 708–716, 2007.

DUARI, S. AND BHATNAGAR, V. scake: Semantic connectivity aware keyword extraction. *Information Sciences* vol. 477, pp. 100–117, 2019.

DURHAM, A. L. *Patent law essentials: A concise guide*. ABC-CLIO, 2018.

- ERBS, N., ZESCH, T., AND GUREVYCH, I. Link discovery: A comprehensive analysis. In *2011 IEEE Fifth International Conference on Semantic Computing*. IEEE, pp. 83–86, 2011.
- GARDNER, J. J. AND XIONG, L. Automatic link detection: A sequence labeling approach. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management*. CIKM '09. ACM, New York, NY, USA, pp. 1701–1704, 2009.
- HAN, X., SUN, L., AND ZHAO, J. Collective entity linking in web text: a graph-based method. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. ACM, pp. 765–774, 2011.
- JANA, A., MOORIYATH, S., MUKHERJEE, A., AND GOYAL, P. Wikim: metapaths based wikification of scientific abstracts. In *2017 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*. IEEE, pp. 1–10, 2017.
- KHODE, A. AND JAMBHORKAR, S. A literature review on patent information retrieval techniques. *Indian Journal of Science & Technology* 10 (37), 2017.
- MANNING, C. D. AND SCHÜTZE, H. *Foundations of statistical natural language processing*. MIT press, 1999.
- MEIRELES, M. R. G., FERRARO, G., AND GEVA, S. Classification and information management for patent collections: a literature review and some research questions. *Information Research* 21 (1), 2016.
- MIHALCEA, R. AND CSOMAI, A. Wikify!: linking documents to encyclopedic knowledge. In *Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management*. CIKM '07. ACM, pp. 233–242, 2007.
- ONAN, A., KORUKOĞLU, S., AND BULUT, H. Ensemble of keyword extraction methods and classifiers in text classification. *Expert Systems with Applications* vol. 57, pp. 232–247, 2016.
- OUELLETTE, L. L. Who reads patents? *Nature biotechnology* 35 (5): 421–424, 2017.
- PANCHENKO, A., RUPPERT, E., FARALLI, S., PONZETTO, S. P., AND BIEMANN, C. Unsupervised does not mean uninterpretable: The case for word sense induction and disambiguation. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*. Vol. 1. pp. 86–98, 2017.
- RATINOV, L., ROTH, D., DOWNEY, D., AND ANDERSON, M. Local and global algorithms for disambiguation to wikipedia. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. Association for Computational Linguistics, pp. 1375–1384, 2011.
- REGINALDO, T. V., LUCINDO, D. L. B., MEIRELES, M. R. G., PATROCÍNIO JÚNIOR, Z. K. G., AND ALMEIDA, P. E. M. A comparison of algorithms for the extraction of keywords in a patent database. *Proceedings of the XXXVIII Iberian Latin-American Congress on Computational Methods in Engineering*, 2017.
- SENEVIRATNE, D. *Patent Link Discovery*. Ph.D. thesis, Queensland University of Technology, 2018.
- TSUNAKAWA, T. AND KAJI, H. Towards cross-lingual patent wikification. *Proceedings of 6th Workshop on Patent and Scientific Literature Translation (PSLT6)* vol. 6, pp. 89, 2015.