

“© 2015 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.”

Modeling Technological Topic Changes in Patent Claims

Hongshu Chen^{1,2}, Yi Zhang^{1,2}, Guangquan Zhang¹, Donghua Zhu², Jie Lu¹

¹Decision Systems & e-Service Intelligence Lab, Centre for Quantum Computation & Intelligent Systems
Faculty of Engineering and Information Technology, University of Technology Sydney, Sydney, Australia

²School of Management and Economics, Beijing Institute of Technology, Beijing, China

Abstract—Patent claims usually embody the most essential terms and the core technological scope to define the protection of an invention, which makes them the ideal resource for patent content and topic change analysis. However, manually conducting content analysis on massive technical terms is very time consuming and laborious. Even with the help of traditional text mining techniques, it is still difficult to model topic changes over time, because single keywords alone are usually too general or ambiguous to represent a concept. Moreover, term frequency which used to define a topic cannot separate polysemous words that are actually describing a different theme. To address this issue, this research proposes a topic change identification approach based on Latent Dirichlet Allocation to model and analyze topic changes with minimal human intervention. After textual data cleaning, underlying semantic topics hidden in large archives of patent claims are revealed automatically. Concepts are defined by probability distributions over words instead of term frequency, so that polysemy is allowed. A case study using patents published in the United States Patent and Trademark Office (USPTO) from 2009 to 2013 with Australia as their assignee country is presented to demonstrate the validity of the proposed topic change identification approach. The experimental result shows that the proposed approach can be used as an automatic tool to provide machine-identified topic changes for more efficient and effective R&D management assistance.

I. INTRODUCTION

Patent claims, as an important part of unstructured segments of a patent document, hold explicit information and implicit knowledge revealing technological concepts, topics and related R&D activities with concise, but precise language [1, 2]. It is often argued as a valuable source for the detection of technological changes and to gain technological insight [3-5]. Since manually conducting content analysis on massive patent documents is very time-consuming and laborious, in recent years, one of the fundamental changes to research in R&D management is the access to extremely powerful information techniques and a vast amount of digital and textual data [6]. In particular, for efficient patent analysis, automatic approaches to assist domain experts and decision makers to mine and understand large volumes of patent documents have drawn increasing attention and still are in great demand [7].

Much effort has been devoted to reveal latent knowledge from the textual data of patent documents. Watts and Porter [8] suggested an approach to investigate terminological trends by tracking the historical change of keywords. Yoon and Park [9] presented a keyword-based morphology study to identify the detailed configurations of promising technology.

Zhang and his colleague [10] introduced a term clumping approach based on Principal Components Analysis to explore keywords and main phrases in abstract from scientific literature. In addition, text analytics have already been applied to Technology Intelligence application *TrendPerceptor* [11], *Techpioneer* [12], *VantagePoint* [13] and *Aureka* [14] to determine hidden concepts and relationships, where clustering, classification and mapping techniques were used to support further content analysis of technological documents. However, before most of these applications are applied, several sets of keywords need to be defined in advance, which still derive from the opinion and knowledge of domain experts. Moreover, the outcomes of traditional text mining techniques are single keywords with ranking, yet these words alone are usually too general or misleading for indicating a concept, especially when there are polysemous words actually describing different topics [7].

To overcome these limitations, this research proposes a topic change identification approach using Latent Dirichlet Allocation. Unsupervised topic modeling is applied to vast amounts of target patent claims, providing a corpus structure with minimal human intervention. There is no pre-set classification or keywords list for this approach and the results are discovered in a completely unsupervised way. In addition, instead of using single terms, topics are represented by probability distributions over words. The actual semantic meaning of a topic is able to be delivered in this way, and at the same time the polysemous words which are actually depicting different concepts, can also be separated. After revealing topics from patent sub-collections of different years, a topic change model is utilized to identify topic changes over time. Finally, to demonstrate the performance of our proposed approach, patents published during years 2009 to year 2013 in the United States Patent and Trademark Office (USPTO) with Australia as their assignee country, are selected to present a case study. The experimental result demonstrates that the proposed approach is able to provide machine-identified topic changes automatically without any pre-setting of keywords. The outcomes of our approach will be used to serve R&D management assistance.

This paper is organized as follows: Section II reviews related research developments by introducing patent data in technological research and Latent Dirichlet Allocation. Section III describes the proposed topic change identification approach step by step. Section IV carries out experiments using USPTO patents to demonstrate the proposed approach in a real patent analysis context. The conclusions and future study are addressed in Section V.

II. LITERATURE REVIEW

A. Patent Data in Technological Research

Patent documents are composed of structured information and unstructured descriptions of inventions. Analytical approaches based on structured data of patents, such as issue date, inventor, assignees or International Patent Classification, have played the major role in both theoretical and practical research [15-17]. However, the unstructured data in patent documents, such as abstracts, claims, and descriptions usually contain much more abundant information than the structured sections, since they contain significant characteristics, detailed functionalities, or major contributions of technologies. Therefore, in recent years, there has been a lot of interest in applying text mining techniques to unstructured patent data to set domain analysts free from studying and understanding massive amounts of technological content [7, 18].

Patent claims, as an important part of unstructured segments of a patent document, embody all the important technical features of an invention with the most essential technological terms to define the protection [19]. A patent claim usually consists of three parts: a Preamble that serves as an introductory section to recite the primary purpose, function or properties; a transition phrase, such as comprising, having including, consisting of, etc.; a 'body' that contains the elements or steps that together describe the invention [20-22]. Claims, on one hand, reveal the core inventive topics and the major technological scope of a patent; on the other hand, they are written in concise, but precise language, which make them the best resource for identifying technological topics and facilitating patent document analysis [1, 2, 20, 23]. This research utilizes patent claims as the main source of topic change analysis.

Among patent databases from different countries, the United States Patent and Trademark Office (USPTO) database is mostly used because patents submitted in other countries are often also simultaneously submitted in the United States [24].

B. Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) [25] is a probabilistic model that aims to estimate the properties of multinomial observations by unsupervised learning. It gives an estimation of the latent semantic topics hidden in large archives of documents, and indicates the probabilities of how various documents belong to different topics. LDA has been used as a very efficient tool to assist topic discovery and analysis, in practice. For example, Griffiths and Steyvers [26] applied LDA-based topic modeling to discover the hot topics covered by papers in Proceedings of the National Academy of Sciences of the United States of America (PNAS); Yang et al. [27] proposed a Topic Expertise Model (TEM) based on LDA to jointly model topics and expertise for Community Question Answering (CQA) with Stack Overflow data; Kim and Oh [28] proposed a framework based on LDA to identify

important topics and their meaningful structure within the news archives on the Web.

The graphical model of LDA is presented in Fig. 1, showing three rectangular plates where: D denotes the overall documents in a corpus; K indicates the topic numbers for D ; and N_d stands for the term number of d^{th} document in document collection D . Each node in Fig. 1 stands for a random variable in the generative process of LDA. All the plates in the figure indicate replication. On the left of the figure, $\vec{\vartheta}_d$ stands for the topic proportions for the d^{th} document. For document d , the topic assignments are Z_d , where $Z_{d,n}$ indicates the topic assignment of the n^{th} word in the d^{th} document. On the right of the figure, the topics themselves are illustrated by $\vec{\varphi}_{1:K}$, where each $\vec{\varphi}_k$ is a distribution over vocabularies. All of the unshaded circles indicate hidden nodes. The shaded circles are observable nodes, where $W_{d,n}$ stands for the n^{th} word in document d . Finally, α and β are two hyperparameters that determine the amount of smoothing applied to the topic distributions for each document and the word distributions for each topic [25, 29-31]. In summary, the generative process of LDA can be denoted by the joint distribution of the random variables in Fig. 1.

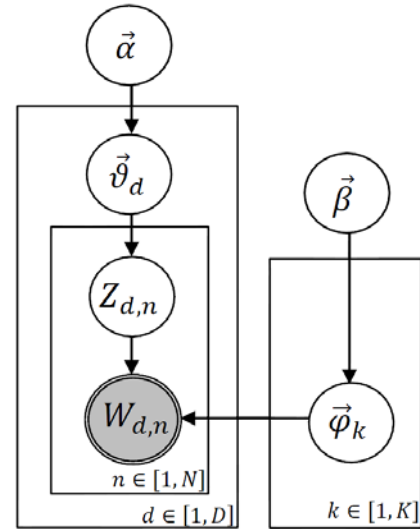


Fig 2. The graphical model of Latent Dirichlet Allocation

The parameters of LDA need to be estimated by an iterative approach. Among existing approaches, Gibbs sampling, which is one of the most commonly used methods, is an approximate inference algorithm based on the Markov Chain Monte Carlo (MCMC) and widely used to estimate the assignment of words to topics by observed data [26, 32, 33]. The randomness introduced by the initiation of the sampling affects the estimation of probabilities in LDA, so that the result is slightly different even with exactly the same setting of input and parameters; yet on the whole, the results of different experiments won't change much.

III. METHODOLOGY

This section explains the details of our proposed topic change identification approach. The framework is given first, then each detailed step is illustrated.

A. Framework

The overall framework of our proposed topic change identification approach is shown in Fig. 2. Users first initiate search statements to declare their domain analytic needs. Patent ID, title, claims, issue time, assignees, United States Patent Classification (USPC) and other information of target patents are then crawled into a database waiting for further analysis. To identify topic changes over time, the whole patent collection is first divided into several sub-collections and labeled with their corresponding issue year. Subsequently, for each sub-collection, patent claims and titles embodying essential technical terms, and USPC providing a general understanding of the domain classification are extracted from the target patents database. The two plates in the figure indicate replication.

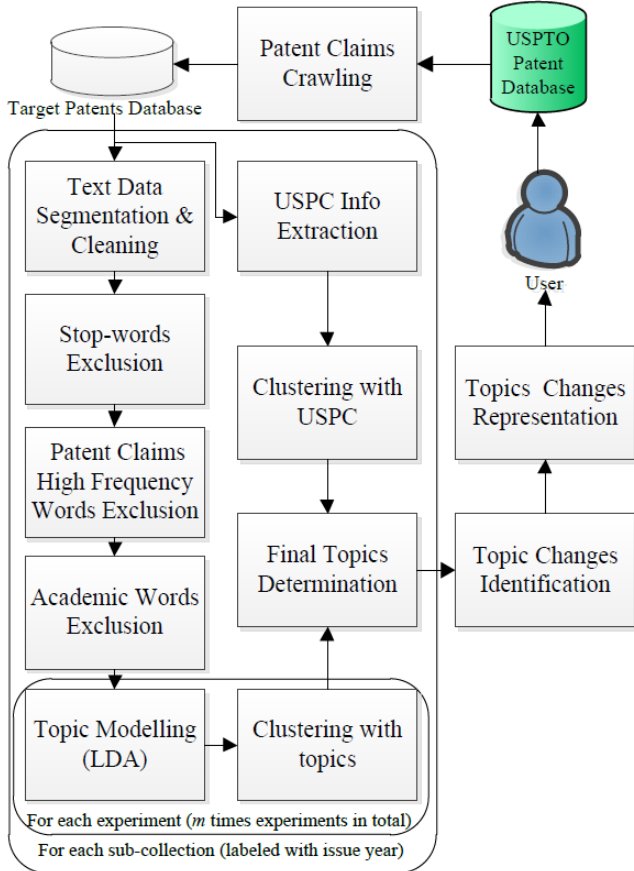


Fig. 2 The framework of the proposed topic change identification approach

Textual data composed by claims and titles, after data segmentation and cleaning, are then placed into a series of words exclusion modules to filter out the most common function words, high frequency words that commonly appeared in patent claims, and academic words with vague

and general meanings. Then, the prepared text will be passed to the topic modeling module. Meanwhile, the USPC information of the corresponding patents is extracted to assist final topic determination. As mentioned, the randomness introduced by the initiation of the sampling will affect the result of LDA. To acquire the most reliable topics of the corpus, we utilize USPC as a measurement to evaluate results from m times experiments. Patents are clustered with both their USPC and topic assignments. The final topic modeling result is the one that provides the most similar clusters to the USPC clustering outcome. Finally, with all the topics estimated from patent sub-collections of different years, topic changes over time can be identified and presented to users.

B. Patent Claim Text Cleaning

Patent claims are a special kind of textual data that contain plenty of technical terms, specific words serving as transition phrases and numerous academic words that describe invention outcomes. Among all the terms that one claim may contain, only technical terms provide most meaningful information reflecting technological topics. Therefore, for patent collections of each year, as shown in Fig. 2, before modeling topics with LDA, we utilize three modules to remove general words from the corpus of patents as follows:

- Stop words such as *the, that, these*
- High frequency words in patent claims such as *claimed, comprising, invention*
- General academic words such as *research, approach, data*.

The stop words list we applied is from an information retrieval Resources link from Stanford University [34]; the patent claim commonly used phrases are summarized from a Transitional Phrase page on Wikipedia [35]; the general academic words list is provided by the University of Nottingham, we select the top 100 most frequent academic words and remove them from our final corpus [10,36].

C. Topic Modeling

LDA utilizes a probability distribution over words, instead of a single term, to define a concept, delivering the semantic meaning of the topic and, at the same time, allowing polysemy. Thus it is very suitable for “understanding” the content of large corpuses such as emails, news, scientific papers and our main data source here, patent claims. After removing all commonly used words from the corpus, we utilize LDA to generate several groups of topics for a number of patent sub-collections, which are labeled by their corresponding issue year. In a sub-collection of the corpus, the claims and title of each patent constitute one document, and the number of documents equals the number of patents; the USPC and other structural information are stored alone in a single file to assist further topic determination. All the textual documents in the corpus are seen as mixtures of a number of topics; each topic is seen as a distribution over various vocabularies. Here we present the global topics

as $\vec{P}_{1:t} = (\vec{P}_1, \vec{P}_2, \dots, \vec{P}_i, \dots, \vec{P}_t)$, where \vec{P}_i stand for the topics of the i^{th} sub-collection of the corpus. The relationship between sub-collections and topics is illustrated in Fig. 3.

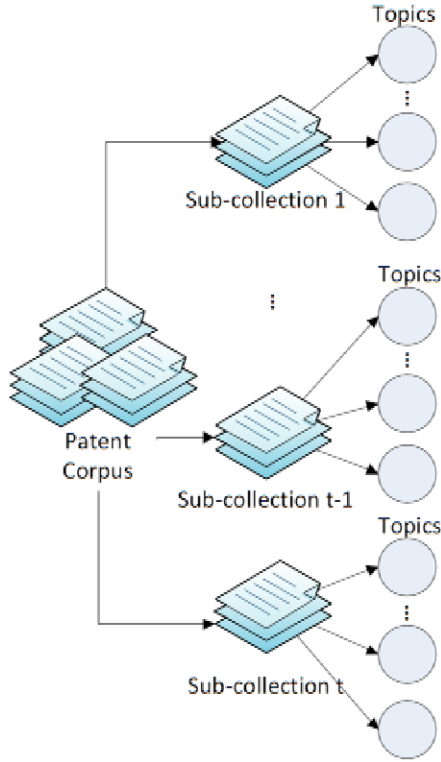


Fig. 3 Relationships between sub-collections and topics

Since we know nothing about the word distributions composing the topics and the topic distributions composing the documents, before topic modeling, assumptions need to be first drawn to determine the parameters k, α, β of LDA. According to previous research, hyper-parameters α, β of the Dirichlet distribution in LDA have a smoothing effect on multinomial parameters; that is, the lower the values of α and β are, the more decisive topic associations there will be [30]. This research sets $\alpha = 0.5$ and $\beta = 0.1$, which are commonly used in LDA applications [37]. For the setting of K , higher K will reduce the topical granularity but increase the processing time. Therefore, during the implementation, K needs to be decided case by case, balancing user requirement and time consumption. Different parameter settings may improve modeling performance, yet optimizing these parameters is beyond the scope of this paper. We then apply Gibbs sampling to infer the needed distributions in LDA. Since the initial values of variables are determined randomly in Gibbs sampling, the outputs of LDA in multiple experiments with a same corpus are slightly different. To ensure the final topic modeling estimation as reliable as possible, evaluation criteria will be needed for the topics finalization.

D. Final Topics Determination

As a predefined classification hierarchy built on domain expert judgments, USPC provides a general understanding of the technical domain of concern to one patent, but most of the time, provides only a general understanding. Because patents covering similar topics are usually assigned to a same main USPC, this research uses the main USPC to judge which estimation is closer to the actual topic distribution.

For a sub-collection of corpus, multiple LDA experiments will produce a number of topic assignment matrixes, each indicating the topic distribution proportions of patent documents in the corresponding trial. As shown in Fig. 2, there will be m times experiments for every sub-collection; and after performing each time run, patents in the sub-collection are clustered with their calculated topic assignments using the hierarchical clustering approach [38]. Meanwhile, the same group of patents will be also clustered with USPC information. The closer the two clustering results are, the more reliable the topic modeling result is.

Specifically, the values of indexes Jaccard, Folkes & Mallows and F1 of m times experiments are used to measure the similarity of the clustering results of two groups, one by topics and the other by USPC. The three indices are listed as follows [39]:

$$J = a / (a + b + c) \quad (1),$$

$$FM = a / \sqrt{r_1 \cdot r_2} \quad (2),$$

$$F_\beta = \frac{(\beta^2 + 1)r_1 r_2}{\beta^2 r_1 + r_2} \quad (3),$$

where J stands for Jaccard coefficient, FM indicates Folkes & Mallows index, F_β presents the F1 indice. In addition, $r_1 = a / (a + b)$, $r_2 = a / (a + c)$, where a represents the number of patents that belong to the same cluster of topics and to the same USPC in our case, b is the number of patents that belong to the same cluster of topics but to different USPC, and c is the number of patents that belong to different clusters of topics but to the same USPC. The topic modeling result that provides the highest index values is the optimal one.

E. Topic Change Identification

After locating the final topics and words underlying the sub-collections of our corpus, we are able to identify the topic change over time. As show in Fig. 4, we compare two groups of topics deriving from different corpus sub-collections, calculating words with a similarity between each topic in \vec{P}_i and all the topics in \vec{P}_{i+1} . If two topics under different sub-collections contain approximately the same group of words, then we believe that these two topics are actually one topic evolving from year to year. However, if there is no similar topic that can be matched in the previous topic collection, then the un-matched topic in the later year is the new one.

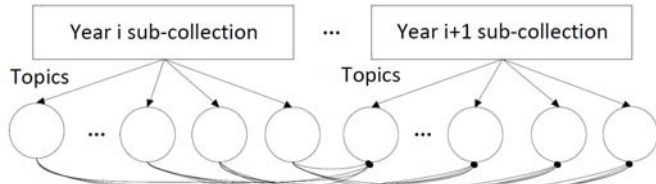


Fig. 4 Topic change identification model

IV. CASE STUDY

A. Data Collection

To demonstrate the performance of our proposed approach, patents published during years 2009 to year 2013 in USPTO (<http://www.uspto.gov/>) with Australia as their assignee country are selected to present a case study. There are 7071 target patents covering 343 different main USPC¹². Their patent ID, titles, issue time, inventors, Assignees, United States Patent Classification (USPC), International Patent Classification (IPC) and most importantly, their claims, are clawed from USPTO and placed in a patents tool for further processing. The claims and title for each patent constitute one document in our corpus, which totals 7071 documents on the whole. Then the whole document collection was divided into five sub-collections to present technological feature and essential terms of inventions by Australia assignees in the past five years. The detailed documents number published every year from 2009 to 2010, the term number and USPC number in each corresponding sub-collection are shown in table 1. Although the documents number declined from year 2011, the term number kept rising, which implies that the average complexity of patent claims description is increasing in the recent three years. We also observe that the number of USPC in 2010 had a visible growth, suggesting that there may be a group of new topics appearing in year 2010.

TABLE 1. THE NUMBER OF DOCUMENTS, TERMS AND USPC OF PATENTS PUBLISHED EACH YEAR

Year	Doc NO.	Term NO.	USPC NO.
2009	1174	19796	199
2010	1613	24726	233
2011	1746	23757	228
2012	1256	25102	233
2013	1282	29714	227

B. Topics Determination

Before topic modeling, as mentioned, a number of parameters need to be set first, including the number of topics, α, β of Dirichlet distribution and the number of iterations for Gibbs sampling. In this case study, we applied $K = 10$ with model hyper-parameters $\alpha = 0.5, \beta = 0.1$ to our target documents, to balance the topical granularity, convenience of understanding, and the speed of processing.

Observation for each year were performed 5 ($m = 5$) runs with 2000 iterations of Gibbs sampling. Indices Folkes & Mallows, Jaccard, and F1 are calculated after we clustered the patents using both topic assignment and USPC information. The detailed index values of five times experiments are listed in Table 2, where we observe directly that the 3rd experiment (E3) of documents sub-collection in 2009, the 5th experiment of documents sub-collection in 2010 (E5), the 4th experiment of documents sub-collection in 2011 (E4), the 2nd experiment of documents sub-collection in 2012 and the 3rd experiment (E3) of documents sub-collection in 2013 have the largest value of all three indexes. Thus the topics and parameters provided by these five trials are the final topic modeling result. There are 10 topics describing the essential technological content and feature for each year; and every topic is presented with 10 words given highest probability by this topic.

The topic modeling results are discovered in a completely unsupervised way, with no pre-set classification or domain knowledge assistance. In the past five years, patents owned by Australia assignees cover several important technological topics, such as printhead and nozzle, alkyl compound, pressure apparatus and antibody sequence. The more the topic words are taken into consideration to describe a topic, the more clear and specific the topical semantic meaning will be. Specifically, the topics for each year are presented as follows. The order of the topics is random, and the numbers behind words are the probability values of corresponding topic words. Details of all the topics, the top 10 ranked words and their corresponding probabilities, are shown in the table 1 in the Appendix.

- The topics of year 2009 include printhead (0.0418) cartridge (0.0353), image (0.0217) device (0.0244), ink (0.0442) nozzle (0.0334), composition (0.0095) material (0.0065), portion (0.0246) assembly (0.0132), roller (0.0142) device (0.0122), alkyl (0.0109) compound (0.0183) formula (0.0111), computer (0.0079) gaming (0.0088), signal (0.0278) sensor (0.0108) and antibody (0.0379) sequence (0.0220).
- The topics of year 2010 contain portion (0.0217) assembly (0.0090), light (0.0131)/optical (0.0104) device (0.0104), ink (0.0518) printhead (0.0476), layer (0.0101) material (0.0144), computer (0.0191) memory (0.0253) plurality (0.0161), coded (0.0252) device (0.0269), antibody (0.0117) sequence (0.0172), pressure (0.0164) apparatus (0.0370), alkyl (0.0096) compound (0.0184) and electrode (0.0146) system (0.0175).
- The topics of year 2011 include layer (0.0166) material (0.0188), portion (0.0260) assembly (0.0202), ink (0.0579) printhead (0.0457), acid (0.0201) sequence (0.0234), alkyl (0.0142) compound (0.0159), pressure (0.0161) apparatus (0.0226), light (0.0133) device (0.0114), image (0.0170) print (0.0449), coded (0.0211) device (0.0207) and plurality (0.0084) apparatus (0.0096).

¹ Data accessed in March 2014

² All plant patents are seen as having one same USPC for calculation convenience.

TABLE 2. INDEXES INFORMATION FOR THE FINAL CHOSEN EXPERIMENT RESULT

Year	Index	E 1	E 2	E 3	E 4	E 5
2009	FM	0.2376	0.2803	0.2845	0.2739	0.1948
2009	DJC	0.1217	0.1500	0.1505	0.1436	0.0962
2009	F1	0.2169	0.2608	0.2616	0.2511	0.1755
2010	FM	0.2668	0.2152	0.2253	0.3125	0.3688
2010	DJC	0.1357	0.1037	0.1077	0.1634	0.2017
2010	F1	0.2389	0.1880	0.1944	0.2809	0.3356
2011	FM	0.2521	0.2484	0.2334	0.2604	0.2541
2011	DJC	0.1334	0.1300	0.1166	0.1342	0.1294
2011	F1	0.2354	0.2301	0.2089	0.2366	0.2292
2012	FM	0.3060	0.3202	0.2773	0.2820	0.2686
2012	DJC	0.1756	0.1853	0.1539	0.1632	0.1521
2012	F1	0.2987	0.3127	0.2667	0.2806	0.2640
2013	FM	0.2984	0.2989	0.3356	0.3177	0.3086
2013	DJC	0.1753	0.1749	0.1986	0.1876	0.1794
2013	F1	0.2983	0.2977	0.3313	0.3159	0.3042

- The topics of year 2012 cover configured (0.0165) signal (0.0325), fluid (0.0209) chamber (0.0145), portion (0.0240) assembly (0.0213), gaming (0.0513) system (0.0205), light (0.0145) lens (0.0067), signal (0.0104) sensor (0.0093), layer (0.0119) material (0.0196), portion (0.0164) apparatus (0.0101), computer (0.0202) memory (0.0150) and acid (0.0151) sequence (0.0162).
- The topics of year 2013 comprise portion (0.0200) assembly (0.0122), gaming (0.0451) controller (0.0226), configured (0.0181) signal (0.0206), cushion (0.0345) mask (0.0287), acid (0.0167) sequence (0.0158), wireless (0.0132) signal (0.0092) sensor (0.0109), layer (0.0120) material (0.0135), optical (0.0095) lens (0.0098), message (0.0103) system (0.0272) and alkyl (0.0132) compound (0.0160).

C. Topic Change Identification

After discovering main topics hidden in patent claims of each year, we then use the topic change model to identify the topic variation from years 2009 to 2013. For different groups of topics associated with two consecutive years, we conduct traversal comparison between the topics that belong to the later year with the topics related to the previous year. Topics that contain very similar words are considered as the same topic experiencing innovation; while topics that cannot match any existing ones count as new topics. Fig. 5 illustrates the new topics that arose each year after 2009, by presenting the top 10 words for each topic using Pajek [40].

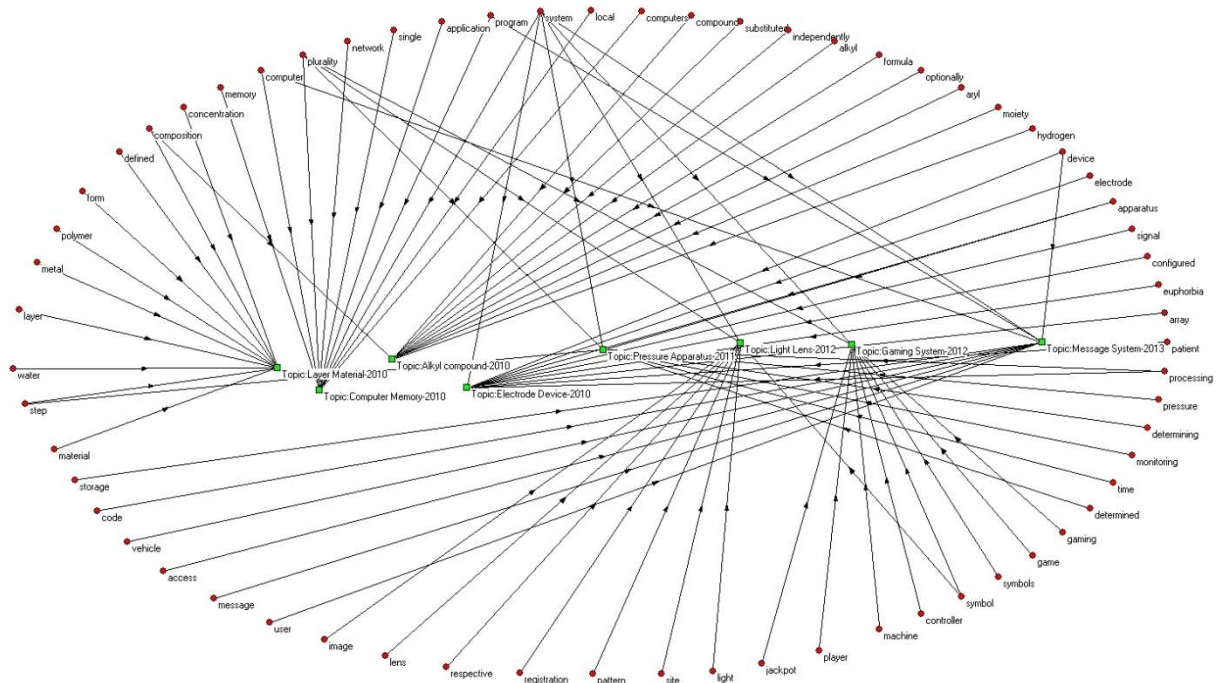


Fig.5 New topics and topmost frequent words of each topic from 2010 to 2013

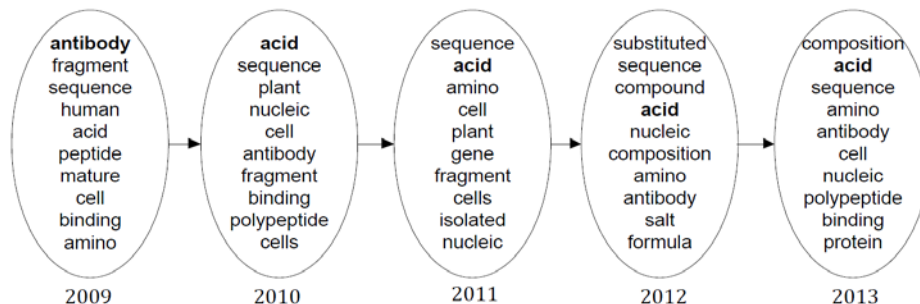


Fig.6 An example of the topic “antibody” evolving over time

In year 2010, four new topics appeared, including layer material that related to metal and polymer composition, electrode device, computer memory and alkyl compound. In year 2011, one new topic appeared, pressure apparatus. Then year 2012 introduced two new topics including light lens and gaming system/controller. Finally, for year 2013, computer system related to vehicle and message appeared as a new topic. All the topics above were identified without assistance of pre-set domain knowledge, which demonstrates the validity of our proposed topic change identification approach. The detailed words and their corresponding probabilities of these new topics are highlighted in boldface in the table 1 of the Appendix.

Moreover, we can also use the proposed approach to discover how the detailed content of a certain topic evolves from year to year. In the case study, topic antibody fragment/sequence is chosen as an example. As shown in Fig. 6, we observe directly that the word distribution composing the topic develops over time. In year 2009, human and peptide were in the top words list, yet after this, the stress of the topic itself moved to plant, amino acid, nucleic acid and polypeptide. The word ‘acid’, instead of ‘antibody’, ranked higher from year 2010 to 2013. The variation of the content of this topic may suggest that, in this area, the key point of technological research and development has shifted to amino/nucleic acid sequence.

V. CONCLUSION AND FUTURE WORK

This paper proposed an unsupervised topic change identification approach using Latent Dirichlet Allocation. Patent claims that embody the most significant technological feature and terms are chosen as the main textual data source of our research. To improve the usage of LDA in patent topic extraction, we utilize USPC as a measurement of different estimations, to reduce the randomness effect on the topic modeling. Machine-identified topics are then placed into a topic change model to locate topic variation over time. Since there is no need to define any keywords in advance and all topics are automatically identified in an unsupervised way, this approach is able to set domain experts and analysts free from reading, understanding and summarizing massive technical documents and records. Finally, a case study, using

USPTO patents published during the years 2009 to 2013 with Australia as their assignee country, is presented. The experimental results demonstrate that the proposed approach can be used as an automatic tool to extract topics and identify topic changes from a large volume of patent documents. From the application perspective, the discovered topic variations can be utilized to assist further decision making in R&D management, especially for newly created innovative enterprises, for example, to provide a full understanding of the topic structure of a certain industry, seek technological opportunities and so on.

As patents and other technological indicators are generating and accumulating in an increasing rate, approaches for automatically identifying topic changes using data mining and machine learning methods will continue to be emphasized. In future work, we will keep focusing on locating topic changes that associate with more meaningful temporal segmentation, like trend turning intervals [41], to identify and analyze the context that contributes to trend changing of patenting activities.

ACKNOWLEDGEMENTS

The work presented in this paper is partly supported by the Australian Research Council (ARC) under Discovery Project DP140101366 and the National High Technology Research and Development Program of China (Grant No.2014AA015105).

REFERENCES

- [1] Xie, Z. and Miyazaki, K., “Evaluating the effectiveness of keyword search strategy for patent identification,” *World Patent Information*, vol. 35(1), pp. 20-30, 2013.
- [2] WIPO. “Patent Cooperation Treaty (PCT) Article 6. Claims,” Retrieved 2014, World Intellectual Property Organization, <http://www.wipo.int/pct/en/texts/articles/a6.htm>.
- [3] Campbell, R.S.; “Patent trends as a technological forecasting tool”, *World Patent Information*, vol. 5(3), pp. 137-143, 1983.
- [4] Ernst, H.; “The use of patent data for technological forecasting: the diffusion of CNC-technology in the machine tool industry”, *Small Business Economics*, vol.9(4), pp.361-381, 1997.
- [5] WIPO, “WIPO Intellectual Property Handbook: Policy, Law and Use. 2 ed”, *Technological and Legal Developments in Intellectual Property*, Vol. 489, pp.17-40, 2004.
- [6] Daim, T.U., Kocaoglu, D.F., Anderson, T.R., “Using technological intelligence for strategic decision making in high technology

- environments”, *Technological Forecasting and Social Change*, vol.78(2), pp. 197-198, 2011.
- [7] Tseng, Y.H., Lin, C.J., Lin, Y.I., “Text mining techniques for patent analysis”, *Information Processing & Management*, vol. 43(5), pp. 1216-1247, 2007.
- [8] Watts, R.J., Porter, A.L., “Innovation Forecasting”, *Technological Forecasting and Social Change*, vol. 56(1), pp.25-47, 1997.
- [9] Yoon, B. and Park, Y., “A systematic approach for identifying technology opportunities: Keyword-based morphology analysis”, *Technological Forecasting and Social Change*, vol. 72(2), pp. 145-160, 2005.
- [10] Zhang, Y., Porter, A., Hu, Z., Guo, Y., Newman, N.C., “Term clumping” for technical intelligence: A case study on dye-sensitized solar cells, *Technological Forecasting and Social Change*, vol. 85(0), pp. 26-39, 2014.
- [11] Yoon, J. and Kim, K., “TrendPerceptor: A property-function based technology intelligence system for identifying technology trends from patents”, *Expert Systems with Applications*, vol.39(3), pp. 2927-2938, 2012.
- [12] Yoon, B., “On the development of a technology intelligence tool for identifying technology opportunity”, *Expert Systems with Applications*, vol. 35(1), pp. 124-135, 2008.
- [13] Zhu, D. and Porter, A.L., “Automated extraction and visualization of information for technological intelligence and forecasting”, *Technological Forecasting and Social Change*, vol. 69(5),pp. 495-506, 2002.
- [14] Trippe, A.J.; “Patinformatics: Tasks to tools”, *World Patent Information*, vol. 25(3), pp. 211-221, 2003.
- [15] Lai, K.K. and Wu, S.J., “Using the patent co-citation approach to establish a new patent classification system”, *Information Processing & Management*, vol.41(2), pp. 313-330, 2005.
- [16] Sheikh, N., Gomez, F. A., Cho, Y.,Siddappa, J., “Forecasting of advanced electronic packaging technologies using bibliometric analysis and Fisher-Pry diffusion model”, in *Technology Management in the Energy Smart World (PICMET)*, 2011 Proceedings of PICMET '11, 2011.
- [17] Nishijima, Y., Anzai, T., Sengoku, S., Application of bibliometric analysis to market analysis. in *Technology Management in the IT-Driven Services (PICMET)*, 2013 Proceedings of PICMET '13, 2013.
- [18] Camus, C. and Brancalion, R., “Intellectual assets management: from patents to knowledge”, *World Patent Information*, vol. 25(2),pp. 155-159, 2003.
- [19] Tong, X. and Frame, J.D., “Measuring national technological performance with patent claims data”, *Research Policy*, vol. 23(2), pp. 133-141, 1994.
- [20] Yang, S. and Soo, V., “Extract conceptual graphs from plain texts in patent claims”, *Engineering Applications of Artificial Intelligence*, vol. 25(4), pp. 874-887, 2012.
- [21] USPTO. “Manual of Patent Examining Procedure: Claim Interpretation. Patent Laws, Regulations, Policies & Procedures, Chapter 2100, Section 2111”, published 2012, Retrieved 2014, <http://www.uspto.gov/web/offices/pac/mpep/s2111.html>.
- [22] Sheldon, J.G.; How to Write a Patent Application, *Practising Law Institute*, 1995.
- [23] Novelli, E.; “An examination of the antecedents and implications of patent scope”, *Research Policy*, 2014(In press).
- [24] USPTO. United States Patent and Trademark Office <http://www.uspto.gov/patents/index.jsp>.
- [25] Blei, D.M., Ng, A.Y., Jordan, M.I., “Latent dirichlet allocation”, *the Journal of machine Learning research*, vol.3, pp. 993-1022, 2003.
- [26] Griffiths, T.L. and Steyvers, M., “Finding scientific topics”, *Proceedings of the National academy of Sciences of the United States of America*, vol.101,Suppl 1, pp. 5228-5235, 2004.
- [27] Yang, L., Qiu, M.,Gottipati, S.,Zhu, F.,Jiang, J., Sun, H.,Chen, Z., “CQArank: jointly model topics and expertise in community question answering”, in Proceedings of *the 22nd ACM international conference on Conference on information & knowledge management*, ACM, 2013.
- [28] Kim, D. and Oh, A., “Topic Chains for Understanding a News Corpus”, in *Computational Linguistics and Intelligent Text Processing*, A. Gelbukh, Editor, Springer Berlin Heidelberg. pp. 163-176, 2011
- [29] Steyvers, M. and Griffiths, T., “Probabilistic topic models”, in *Latent Semantic Analysis: A road to meaning*, D.M. T. Landauer, S. Dennis, and W. Kintsch, Editor, Laurence Erlbaum, 2007.
- [30] Heinrich, G.; “Parameter estimation for text analysis”, Fraunhofer, IGD: Darmstadt, Germany,2005.
- [31] Blei, D.M.;“Probabilistic topic models”,*Communications of the ACM*,vol.55(4),pp.77-84, 2012.
- [32] Noel, G.E. and Peterson, G.L., “Applicability of Latent Dirichlet Allocation to multi-disk search”, *Digital Investigation*, 2014.
- [33] Lukins, S.K., Kraft, N.A., and Etkorn, L.H., “Bug localization using latent Dirichlet allocation”, *Information and Software Technology*, vol. 52(9), pp. 972-990, 2010.
- [34] Lewis,D.D., Yang,Y., Rose,G.T., Li, F., “SMART stopword list”, *Journal of Machine Learning Research*, published 2004, <http://jmlr.csail.mit.edu/papers/volume5/lewis04a/a11-smart-stop-list/english.stop>.
- [35] Wikipedia. Transitional phrase. 2014, http://en.wikipedia.org/wiki/Transitional_phrase.
- [36] Haywood,S. AcademicVocabulary, published 2003, Academic Vocabulary - University of Nottingham <http://www.nottingham.ac.uk/alzsh3/acvocab/wordlists.htm>
- [37] Kolteov, S., Koltsova, O., and Nikolenko, S., “Latent dirichlet allocation: stability and applications to studies of user-generated content”, in Proceedings of *the 2014 ACM conference on Web science*, ACM: Bloomington, Indiana, USA, pp. 161-165, 2014.
- [38] Steinbach, M., Karypis, G., and Kumar, V., “A comparison of document clustering techniques”, in *KDD workshop on text mining*, Boston. 2000.
- [39] Halkidi, M., Batistakis, Y., and Vazirgiannis, M., “On clustering validation techniques”, *Journal of Intelligent Information Systems*, vol. 17(2-3), pp. 107-145, 2001.
- [40] Batagelj, V., and Mrvar, A. Pajek—analysis and visualization of large networks Springer, 2004.
- [41] Chen, H., Zhang, G., Lu, J., "A patent time series processing component for technology intelligence by trend identification functionality," *Neural Computing and Applications*, vol.(26:2),pp345-353,2015.

APPENDIX

TABLE 1. THE TOP 10 RANKED WORDS OF TOPICS FOR YEARS FROM 2009 TO 2013 AND THEIR CORRESPONDING PROBABILITIES

Year 2009									
Topic 1		Topic 2		Topic 3		Topic 4		Topic 5	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
printhead	0.0418	device	0.0244	ink	0.0442	step	0.0116	portion	0.0246
ink	0.0353	image	0.0217	ejection	0.0336	composition	0.0095	body	0.0150
print	0.0333	coded	0.0209	nozzle	0.0334	gas	0.0088	assembly	0.0132
printer	0.0252	system	0.0195	inkjet	0.0307	leach	0.0081	surface	0.0110
media	0.0229	sensing	0.0181	printhead	0.0245	material	0.0065	extending	0.0092
cartridge	0.0138	digital	0.0132	drop	0.0229	acid	0.0064	wall	0.0091
module	0.0137	computer	0.0105	apparatus	0.0224	fuel	0.0063	mask	0.0081
printing	0.0135	camera	0.0101	actuator	0.0220	water	0.0059	adapted	0.0076
assembly	0.0132	identity	0.0092	element	0.0191	polymer	0.0058	substantially	0.0072
configured	0.0124	position	0.0086	chamber	0.0189	ph	0.0055	support	0.0069
Topic 6		Topic 7		Topic 8		Topic 9		Topic 10	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
support	0.0152	compound	0.0183	system	0.0116	signal	0.0278	antibody	0.0379
roller	0.0142	formula	0.0111	material	0.0090	sensor	0.0108	fragment	0.0246
device	0.0122	alkyl	0.0109	game	0.0088	signals	0.0107	sequence	0.0220
drive	0.0109	independently	0.0102	plurality	0.0087	frequency	0.0089	human	0.0219
assembly	0.0101	layer	0.0098	computer	0.0079	device	0.0087	acid	0.0177
mechanism	0.0082	optionally	0.0095	gaming	0.0073	input	0.0084	peptide	0.0175
surface	0.0080	base	0.0088	entry	0.0072	output	0.0081	mature	0.0164
frame	0.0075	detector	0.0087	torque	0.0063	apparatus	0.0081	cell	0.0157
position	0.0071	substituted	0.0087	object	0.0058	processing	0.0071	binding	0.0138
mounted	0.0067	reflector	0.0087	service	0.0054	power	0.0067	amino	0.0133
Year 2010									
Topic 1		Topic 2		Topic 3		Topic 4		Topic 5	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
portion	0.0217	signal	0.0240	ink	0.0518	material	0.0144	memory	0.0253
surface	0.0126	light	0.0131	printhead	0.0476	step	0.0136	computer	0.0191
outer	0.0095	system	0.0121	nozzle	0.0214	water	0.0101	plurality	0.0161
assembly	0.0090	optical	0.0104	inkjet	0.0183	layer	0.0101	network	0.0155
body	0.0088	device	0.0104	print	0.0176	metal	0.0088	single	0.0143
extending	0.0086	image	0.0083	assembly	0.0172	polymer	0.0081	application	0.0141
wall	0.0080	power	0.0076	printer	0.0156	form	0.0070	program	0.0133
support	0.0076	frequency	0.0076	media	0.0127	defined	0.0067	system	0.0117
upper	0.0073	output	0.0069	ejection	0.0126	composition	0.0066	local	0.0103
frame	0.0071	sensor	0.0067	configured	0.0110	concentration	0.0063	computers	0.0097
Topic 6		Topic 7		Topic 8		Topic 9		Topic 10	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
device	0.0269	acid	0.0199	apparatus	0.0370	compound	0.0184	system	0.0175
coded	0.0252	sequence	0.0172	air	0.0214	substituted	0.0183	device	0.0154
system	0.0245	plant	0.0159	pressure	0.0164	independently	0.0140	electrode	0.0146
print	0.0190	nucleic	0.0152	fluid	0.0148	alkyl	0.0096	apparatus	0.0107

computer	0.0168	seq	0.0146	valve	0.0144	formula	0.0094	signal	0.0105
sensing	0.0161	cell	0.0136	flow	0.0140	optionally	0.0092	configured	0.0095
user	0.0149	antibody	0.0117	chamber	0.0131	aryl	0.0065	euphoria	0.0095
media	0.0115	fragment	0.0088	system	0.0129	moiety	0.0051	array	0.0079
mobile	0.0109	binding	0.0086	inlet	0.0083	composition	0.0049	patient	0.0074
indicative	0.0101	polypeptide	0.0086	outlet	0.0071	hydrogen	0.0046	processing	0.0071

Year 2011

Topic 1		Topic 2		Topic 3		Topic 4		Topic 5	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
material	0.0188	portion	0.0260	ink	0.0579	sequence	0.0234	optionally	0.0228
layer	0.0166	assembly	0.0202	printhead	0.0457	acid	0.0201	substituted	0.0224
step	0.0130	mask	0.0113	nozzle	0.0282	seq	0.0179	compound	0.0159
composition	0.0083	support	0.0110	inkjet	0.0170	amino	0.0138	alkyl	0.0142
range	0.0070	frame	0.0105	assembly	0.0163	cell	0.0130	lens	0.0102
polymer	0.0064	surface	0.0095	chamber	0.0118	plant	0.0120	independently	0.0089
coating	0.0060	outer	0.0087	integrated	0.0116	gene	0.0113	optical	0.0079
metal	0.0058	wall	0.0084	printer	0.0113	fragment	0.0096	aryl	0.0074
solution	0.0057	extending	0.0071	fluid	0.0107	cells	0.0085	zone	0.0070
forming	0.0056	body	0.0069	plurality	0.0103	isolated	0.0084	lower	0.0067
Topic 6		Topic 7		Topic 8		Topic 9		Topic 10	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
apparatus	0.0226	signal	0.0203	print	0.0449	system	0.0289	system	0.0108
flow	0.0191	light	0.0133	media	0.0296	coded	0.0211	step	0.0099
air	0.0180	power	0.0120	printer	0.0177	device	0.0207	apparatus	0.0096
gas	0.0180	device	0.0114	image	0.0170	computer	0.0186	plurality	0.0084
water	0.0178	wireless	0.0103	controller	0.0148	memory	0.0140	pressure	0.0078
pressure	0.0161	apparatus	0.0090	module	0.0141	sensing	0.0130	determining	0.0076
valve	0.0158	source	0.0090	game	0.0131	plurality	0.0114	processing	0.0066
device	0.0129	plurality	0.0078	gaming	0.0129	identity	0.0109	monitoring	0.0058
fluid	0.0124	electrical	0.0078	configured	0.0127	indicative	0.0101	time	0.0057
humidifier	0.0110	optical	0.0074	printing	0.0120	position	0.0086	determined	0.0055

Year 2012

Topic 1		Topic 2		Topic 3		Topic 4		Topic 5	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
signal	0.0325	fluid	0.0209	portion	0.0240	gaming	0.0513	light	0.0145
configured	0.0165	gas	0.0172	assembly	0.0213	game	0.0504	plurality	0.0114
frequency	0.0132	flow	0.0151	support	0.0126	system	0.0205	system	0.0107
optical	0.0116	chamber	0.0145	mask	0.0106	symbols	0.0190	site	0.0075
sound	0.0116	system	0.0132	system	0.0087	symbol	0.0186	pattern	0.0070
system	0.0103	valve	0.0129	element	0.0080	plurality	0.0185	registration	0.0070
power	0.0092	water	0.0121	nasal	0.0073	controller	0.0172	respective	0.0068
control	0.0090	inlet	0.0099	adapted	0.0072	machine	0.0166	lens	0.0067
electrical	0.0088	pressure	0.0097	frame	0.0071	player	0.0157	symbol	0.0063
device	0.0087	liquid	0.0078	extending	0.0066	jackpot	0.0127	image	0.0063
Topic 6		Topic 7		Topic 8		Topic 9		Topic 10	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
time	0.0112	material	0.0196	portion	0.0164	system	0.0202	substituted	0.0204
determining	0.0107	layer	0.0119	apparatus	0.0101	computer	0.0202	optionally	0.0190

signal	0.0104	polymer	0.0100	surface	0.0101	memory	0.0150	sequence	0.0162
test	0.0093	metal	0.0093	device	0.0098	device	0.0139	compound	0.0157
sensor	0.0093	surface	0.0092	body	0.0088	user	0.0128	acid	0.0151
flow	0.0089	electrically	0.0074	upper	0.0088	plurality	0.0081	seq	0.0095
waveform	0.0085	step	0.0067	extending	0.0087	coded	0.0078	nucleic	0.0084
pressure	0.0085	conductive	0.0064	lower	0.0081	content	0.0078	composition	0.0079
predetermined	0.0070	cell	0.0057	container	0.0081	printed	0.0071	amino	0.0072
plant	0.0068	component	0.0056	assembly	0.0073	image	0.0069	antibody	0.0069

Year 2013

Topic 1		Topic 2		Topic 3		Topic 4		Topic 5	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
portion	0.0200	game	0.0555	signal	0.0206	cushion	0.0345	composition	0.0234
assembly	0.0122	gaming	0.0451	configured	0.0181	mask	0.0287	seq	0.0184
body	0.0107	symbol	0.0322	apparatus	0.0145	portion	0.0285	acid	0.0167
surface	0.0091	plurality	0.0274	device	0.0139	assembly	0.0191	sequence	0.0158
extending	0.0079	symbols	0.0238	stimulation	0.0105	frame	0.0186	amino	0.0102
wall	0.0073	controller	0.0226	signals	0.0097	support	0.0168	antibody	0.0091
housing	0.0072	player	0.0189	system	0.0096	structure	0.0154	cell	0.0076
position	0.0070	system	0.0177	power	0.0096	full-face	0.0124	nucleic	0.0071
relative	0.0063	arranged	0.0152	flow	0.0091	nasal	0.0122	polypeptide	0.0068
outer	0.0062	machine	0.0128	electrical	0.0086	underlying	0.0121	binding	0.0066
Topic 6		Topic 7		Topic 8		Topic 9		Topic 10	
Word	Probability	Word	Probability	Word	Probability	Word	Probability	Word	Probability
device	0.0286	material	0.0135	image	0.0236	system	0.0272	substituted	0.0583
wireless	0.0132	layer	0.0120	oligonucleotide	0.0120	computer	0.0260	optionally	0.0513
system	0.0115	fluid	0.0102	lens	0.0098	user	0.0154	compound	0.0160
plurality	0.0112	gas	0.0094	optical	0.0095	program	0.0112	alkyl	0.0132
sensor	0.0109	flow	0.0084	antisense	0.0086	message	0.0103	independently	0.0129
signal	0.0092	water	0.0083	light	0.0085	access	0.0088	formula	0.0084
processing	0.0088	liquid	0.0081	plurality	0.0077	vehicle	0.0071	alkenyl	0.0084
control	0.0088	surface	0.0075	system	0.0070	code	0.0061	salt	0.0076
devices	0.0087	step	0.0067	laser	0.0063	storage	0.0060	alkynyl	0.0066
component	0.0082	electrode	0.0066	step	0.0062	device	0.0059	acceptable	0.0065