

Identification, Tracking and Impact: Understanding the trade secret of catchphrases

Jagriti Jalal
IIT Kharagpur, India
jagritijalal@iitkgp.ac.in

Mayank Singh
IIT Gandhinagar, India
singh.mayank@iitgn.ac.in

Arindam Pal
Data61, CSIRO & Cyber Security CRC
Sydney, New South Wales, Australia
arindamp@gmail.com

Lipika Dey
TCS Innovation Labs, India
lipika.dey@tcs.com

Animesh Mukherjee
IIT Kharagpur, India
animeshm@cse.iitkgp.ac.in

ABSTRACT

Understanding the topical evolution in industrial innovation is a challenging problem. With the advancement in the digital repositories in the form of patent documents, it is becoming increasingly more feasible to understand the innovation secrets – ‘catchphrases’ – of organizations. However, searching and understanding this enormous textual information is a natural bottleneck. In this paper, we propose an unsupervised method for the extraction of catchphrases from the abstracts of patents granted by the U.S. Patent and Trademark Office over the years. Our proposed system achieves substantial improvement, both in terms of precision and recall, against state-of-the-art techniques. As a second objective, we conduct an extensive empirical study to understand the temporal evolution of the catchphrases across various organizations. We also show how the overall innovation evolution in the form of introduction of newer catchphrases in an organization’s patents correlates with the future citations received by the patents filed by that organization. Our code and data sets will be placed in the public domain.

CCS CONCEPTS

• **Social and professional topics** → **Patents**; • **Information systems** → *Data mining*; Information extraction.

KEYWORDS

Patents; digital library;

1 INTRODUCTION

As software and other products are becoming more complex, the number and size of patent documents are increasing gradually. Automated patent document processing systems are essential to extract information and gain insights from this ever-increasing collection of patent databases. Catchphrases provide a concise representation of the content of a document. A catchphrase is a well-known

word or phrase encapsulating the particular concept or subject of a document. They contain all the important legal and technical aspects, instead of just summarizing the document. They have numerous applications such as document categorization, clustering, summarization, indexing, topic search, quantifying semantic similarity with other documents, and conceptualizing particular knowledge domain of the document [12, 16]. However, since only a small minority of documents have author-assigned catchphrases, and manual assignment of catchphrases to existing documents is time-consuming, the automation of the catchphrase extraction process is highly desirable. In the current study, catchphrases represent innovation topics. Figure 1 presents example catchphrases from two different patent abstracts.

In this paper, we propose an unsupervised method for the extraction of catchphrases from the abstracts of patents granted by the U.S. Patent and Trademark Office over the years. The key contributions of this paper are as follows.

- We propose an unsupervised technique for catchphrase identification and ranking in patent documents.
- We conduct robust evaluations and comparison against several state-of-the-art baselines.
- As a secondary objective, we study the evolution of catchphrases present in the patents filed by various organizations over time.
- We bring forth some of the unique temporal characteristics of these catchphrases and show how these are correlated to the overall future citation count of the patents filed by an organization.
- The catchphrase evolution study further unfolds that companies get polarized based on whether the patent documents keep re-using the same catchphrases over time or they introduce newer catchphrases as time progresses.

2 RELATED WORK

A variety of techniques have been applied for automated keyword extraction like locating important phrases by analyzing markups like capitalization, section headings and emphasized texts [17]; building phrase dictionary by parts-of-speech (POS) tagging of word sequences [18]; thesaurus-based keyphrase indexing [29]; domain-specific keyphrase extraction [11, 23] and several other supervised methods such as KEA [30], MAUI [22], back-of-the-book

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

JCDL '20, August 1–5, 2020, Virtual Event, China

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7585-6/20/06...\$15.00

<https://doi.org/10.1145/3383583.3398512>

ID: US06681004

Abstract: The telephone memory aid provides a database to a primary party for storing and retrieving **personal information** about a secondary party, including **summary information** related to **communication exchanges** between the primary and secondary parties. The summary information includes, for example, the date and time of prior telephone calls and the topics discussed. This secondary party information, including the summaries of prior telephone calls, is available for review by the primary party during future phone calls with the secondary party. The telephone memory aid also facilitates entry of information into the **database** through **speech recognition** algorithms and through **question** and answer **sessions** with the primary and secondary parties.

ID: US06680003

Abstract: The present invention concerns **chiral doping agents** allowing a modification to be induced in the spiral pitch of a **cholesteric liquid crystal**, said doping agents including a **biactivated chiral unit** at least one of whose functions allows a chemical link to be established with an isomerisable group, for example by **radiation**, said group possibly having a polymerisable or co-polymerisable end chain. These new chiral **doping agents** find application in particular in a **color display**.

Figure 1: Example abstracts from USPTO patents US06681004 and US06680003. The highlighted set of words are identified as catchphrases from IPC (described in Section 5).

indexing using catchphrase extraction [10], MAUI with text denoising [26], CSSeer [8] etc. In recent years artificial neural networks (ANNs) are being used to build predictive models to rank words in a document [5] and then select keywords based on these ranks. It has been widely recognized that the innovative capability of a firm is a critical determinant of its performance and competitive edge [3, 13, 15]. Since patents are a direct outcome of the inventive process and are broken down by technical fields, they are considered indicators of not only the rate of the innovative activities of a firm but also its direction [1, 2, 4]. Many previous studies have examined the relationships between the patenting activities of a company and its market value [14, 24]. Bornmann and Daniel [6] precisely reviews the citing behavior of scientists and shows the role of citations as a reliable measure of impact. Cheng et al. [9] shows that some indicators of patent quality are statistically significant to return on assets. Lee et al. [19] assesses future technological impacts by employing the future citation count as a proxy while Lee et al. [20] employs various patent indicators such as novelty and scope, as features of an ANN for early identification of emerging technologies.

3 DATASETS AND PREPROCESSING

The current study requires a rich time-stamped dataset. We, therefore, leverage two independent data sources. These are:

- (1) **The patent dataset:** We compile the first dataset by crawling the full-text patent articles, available at the United States Patent and Trademark Office (USPTO¹). It comprises patents granted weekly (Tuesdays) from January 1, 2003, to May 18, 2018 (excluding images/drawings). The patents are available as XML encoded files with English as the primary language. Out of all the curated documents, in this study, we only consider those patents for which the abstract information is present (see Table 1 for statistics).
- (2) **The newsgroup corpus:** We also use another data source, the *20 Newsgroups Dataset*² donated by T. Mitchell in 1999. It includes one thousand Usenet articles each from 20 newsgroups like 'alt.atheism', 'comp.graphics', 'talk.politics.guns', etc. Approximately 4% of the articles are crossposted. This serves as a non-patent corpus to estimate the importance of a word specifically in the domain of the patents concerning a non-patent domain (see Table 1 for statistics).

Patent	Year range	2003–2018
	Number of patents	3,915,639
	Number of patents with abstract	3,486,866
Newsgroup	Year range	1993–2017
	Number of articles	19,997
	Number of words	–
	Language	English

Table 1: General statistics about the patent dataset and the newsgroup corpus. A large fraction (89%) of patents have abstract information.

Pre-processing: For both of the above, we performed several pre-processing tasks such as a sentence to lowercase conversion, removal of special characters except apostrophe and periods, lemmatization, and multiple white-spaces removals.

4 CATCHPHRASE EXTRACTION

Catchphrase extraction is a challenging problem mainly due to the diversity and unavailability of large-text annotated datasets. We, therefore, present an unsupervised method for catchphrase extraction. We propose a two-stage extraction strategy that identifies relevant candidate catchphrases in a given patent article. In the first stage, we select the candidate catchphrases. This is followed by candidate catchphrase ranking in the second stage. Next, we describe the two stages in detail.

4.1 Stage-1: Candidate selection

In the first stage, we select candidate catchphrases from each patent’s abstract. Empirically, we observe that all catchphrases are n-gram noun phrases, for example, unigrams (e.g. *communication, dielectrometry, etc.*), bigrams (e.g. *consecutive bit, voice synthesizer, etc.*),

¹<https://bulkdata.uspto.gov/>

²<https://archive.ics.uci.edu/ml/datasets/Twenty+Newsgroups>

trigrams (e.g. *integrated circuit device, hydrogen chloride gas, etc.*) or quadrigrams (e.g. *commercially available synthesis tool, electric signal processing board, etc.*). We, therefore, perform part-of-speech-tagging (POS) of each abstract text to identify noun phrases. Currently, we leverage python’s state-of-the-art NLP library SpaCy³. Note that, we experimented with two text processing approaches before noun phrase identification: (i) with stopwords (WS), and (ii) without stopwords (WOS). WS represents that no stopwords were removed from the abstracts, whereas, WOS represents that all stopwords in the abstract text were removed beforehand. Abstracts with stopwords (WS) led to better quality extraction results due to the existence of stop-words in noun phrases. We discuss the results in detail in Section 5. Table 2 presents statistics of extracted candidate phrases from the dataset.

Word n-grams	Count
Unigrams	208,105
Bigrams	2,616,762
Trigrams	4,432,251
Quadrigrams	2,138,696
Total	9,395,814

Table 2: Count of n-gram noun phrases generated from patent dataset.

4.2 Stage-2: Candidate ranking

Candidate phrases obtained in the first stage are ranked in this stage. The ranking algorithm is based upon two empirical findings: (i) *how well the phrase describes the document’s topic*, and (ii) *how specific is the phrase to the patent literature*. Our proposed method unifies both of these findings by combining a frequency-based measure with an information-theoretic measure. Given a patent document d and a set of candidate phrases c_d obtained in the previous stage, we compute the phrase score $PS(c, d)$ for each phrase $c \subseteq c_d$.

$$PS(c, d) = \sum_{i=1}^{|c|} \{\log(\text{score}(t_i))\} \cdot KLI(c, d) \quad (1)$$

where, t_i denotes the i^{th} term in an n-gram candidate phrase c , $\text{score}(t_i)$ denotes the score of the i^{th} term by estimating the importance of the term specifically in the patent domain relative to a non-patent domain and $KLI(c, d)$ represents the Kullback-Leibler divergence informativeness specifying how well a candidate phrase c represents a document d . The term $\text{score}(t)$ in the above equation is computed as

$$\text{score}(t) = \frac{\text{Importance}(t, C_p)}{\text{Importance}(t, C_n) + 1} \quad (2)$$

Again, here, C_p and C_n represents the patent collection and non-patent (in our case, the newsgroup) collection. The importance(t, C) of a term t in a given collection $C \in \{C_p, C_n\}$ is measured in terms of the collection frequency CF and the document frequency DF . $CF(t, C)$ represents how many times the term t appeared in the

entire collection C . $DF(t, C)$ represents the count of documents where the term t appeared. It is computed as

$$\text{Importance}(t, C) = \frac{CF(t, C)}{DF(t, C) + 1} \quad (3)$$

$KLI(c, d)$ denotes an information theoretic measure to compute how informative the phrase is in the given document d . It is computed as:

$$KLI(c, d) = \frac{TF(c, d)}{|d|} \cdot \log \frac{\frac{TF(c, d)}{|d|}}{\frac{CF(c)}{n}} \quad (4)$$

where, $TF(c, d)$ represents how many times c appeared in document d . $CF(c)$ denotes how many times c appeared in the entire patent collection C_p . $|d|$ and $|n|$ represents total number of n-grams in document d and C_p respectively.

The above scoring method results in a ranking of candidate phrases. We select top-ranked candidates such as top-5, top-10, top-20, etc., and evaluate our unsupervised method in the next section.

5 EXPERIMENTS

In this section, we describe the experimental settings, baselines and the evaluation metrics. We construct a collection of possible catchphrases from the International Patent Classification (IPC) list. This list is maintained by the *World Intellectual Property Organization* (WIPO)⁴. The IPC provides a hierarchical system of language independent symbols for the classification of patents and utility models according to the different areas of technology to which they pertain. The hierarchy comprises eight high-level categories:

- (1) **Cat-1:** Human necessities
- (2) **Cat-2:** Performing operations; Transporting
- (3) **Cat-3:** Chemistry; Metallurgy
- (4) **Cat-4:** Textiles; Paper
- (5) **Cat-5:** Fixed constructions
- (6) **Cat-6:** Mechanical engineering; Lighting; Heating; Weapons; Blasting
- (7) **Cat-7:** Physics
- (8) **Cat-8:** Electricity

In each of these high-level categories, several sub-categories exist. An n-gram phrase represents each category. We term these phrases as *ground truth catchphrases* (GTC). Overall, we obtained 22,855 GTC such as "actuators", "cleaning fabrics", "feedback arrangements in control systems", etc. We use GTC to evaluate our proposed catchphrase extraction method. Table 3 presents examples of GTC for each high-level category. Next, we present three state-of-the-art baselines.

5.1 Baselines

- (1) **Keyphrase extraction algorithm (KEA):** KEA [30] is a supervised machine learning toolkit that extracts keyphrases and ranks them. The original algorithm was trained on scientific documents and uses a trained Naive Bayes model. We trained KEA for patent documents leveraging a similar training procedure.

³<https://spacy.io/usage/linguistic-features>

⁴<http://www.wipo.int/classifications/ipc/ipcpub/>

Category	Unigrams	Bigrams	Trigrams	Quadgrams
Cat-1	rhinoscopes	dental surgery	table service equipment	foodstuffs containing gelling agents
Cat-2	thwarts	rivet hearths	making plough shares	making plastics bushes bearings
Cat-3	riboflavin	septic tanks	acetone carboxylic acid	chromising of metallic material surfaces
Cat-4	carding	carbon filaments	opening fiber bales	drying wet webs in paper-making
Cat-5	collieries	suspension bridges	setting anchoring bolts	freezing for sinking mine shafts
Cat-6	thermal	diesel engines	portable accumulator lamps	treating internal-combustion engine exhaust
Cat-7	ozotypy	investigating abrasion	measuring electric supply	incineration of solid radioactive waste
Cat-8	rheostats	electric accumulator	thermo magnetic devices	electric amplifiers for amplifying pulse

Table 3: Examples of ground truth catchphrases for each high-level category available in the International Patent Classification (IPC) list.

- (2) **Legal:** Mandal et al. [21] also follow an unsupervised approach for identification of catchphrases from legal court cases. The scoring is done as:

$$PS(c, C_p, C_{np}) = \log \left[\sum_{i=1}^{|c|} Score(t_i, C_p, C_{np}) \right] \cdot KLI(c, d) \quad (5)$$

where $score(t_i, C_p, C_{np})$ and $KLI(c, d)$ can be calculated using equations 2 and 4 respectively. Note the change in the formula in equation 5 compared to equation 1. This modification as we shall see almost doubles our performance.

- (3) **KLIP:** Tomokiyo and Hurst [27], Verberne et al. [28] proposed a Kullback-Leibler (KL) divergence based phrase assignment score which is a linear combination of two different scores:

- (a) *KL informativeness (KLI):* KLI measures how well a candidate phrase represents a document. It is computed using equation 4.

- (b) *KL phraseness (KLP):* KLP score is computed specifically for multi-word phrases. It compensates for low frequency of multi-word phrases by assigning higher weights to longer phrases:

$$KLP(c, d) = \frac{TF(c, d)}{|d|} \cdot \log \frac{\frac{TF(c, d)}{|d|}}{\prod_{i=1}^{|c|} \frac{freq(t_i, d)}{|d|}} \quad (6)$$

where, t_i is the i^{th} term of the phrase c , and $freq(t_i, d)$ is the frequency of the term t_i in document d .

- (4) **BM25:** BM25 [25] is a well-known measure for scoring documents with respect to a given query. We use this function for assigning score to an extracted candidate phrase c in a given document d . The scoring function is:

$$score(c, d) = IDF(c) \cdot \frac{TF(c, d) \cdot (k_1 + 1)}{TF(c, d) + k_1 \cdot \left(1 + b + b \cdot \frac{|d|}{avgdl}\right)} \quad (7)$$

where $TF(c, d)$ is the term frequency of phrase c in the document d . k_1 and b are free parameters. We choose $k_1 \in [1.2, 2.0]$ and $b = 0.75^5$. $IDF(c)$ is the inverse document frequency of the candidate phrase c , calculated as

$$IDF(c) = \log \frac{n - DF(c) + 0.5}{DF(c) + 0.5} \quad (8)$$

where $DF(c)$ is the document frequency of the phrase c in the collection.

Note that KEA is a supervised machine learning model whereas Legal, KLIP and BM25 are unsupervised methods.

5.2 Evaluation measures

We evaluate our proposed method against the three baselines. We use two standard evaluation measures: (i) **Macro precision**, and (ii) **Macro recall**. These metrics are computed by macro-averaging the precision/recall values computed for every patent.

$$\text{Macro precision} = \frac{\sum_{i=1}^{|T|} precision_i}{|T|} \quad (9)$$

$$\text{Macro recall} = \frac{\sum_{i=1}^{|T|} recall_i}{|T|} \quad (10)$$

where $precision_i$ and $recall_i$ are the precision and recall values computed for i^{th} patent in our test dataset T . The precision and recall values for the i^{th} patent are computed as follows

$$precision_i = \frac{DCP_i}{DC_i} \quad (11)$$

$$recall_i = \frac{DCP_i}{CPG_i} \quad (12)$$

where DC_i , DCP_i , and CPG_i represents the number of catchphrases in the i^{th} patent that are detected, detected and present in GTC, and present in GTC respectively.

As KEA requires training, we partition our dataset into two classes: (i) train and (ii) test. Train split consist of 2,055,588 (65%) patent documents. Test split consist of 1,106,883 (35%) patent documents. For a fair comparison, we evaluate our proposed method against baselines (described in Section 5.1) using only the test split.

5.3 Results and discussion

Table 4 compares our proposed catchphrase extraction approach against state-of-the-art baselines. We outperform all baselines by a substantially high margin. The second best system in terms of precision is KEA, whereas the second-best system in terms of recall is a mix between KLIP and KEA. The baseline Legal performed worst among all the baselines, which is possible because of the fact that the authors take a logarithm of the sum of all the scores rather than the sum of the logarithms of the scores. The former measure

⁵We select these values as per previous literature [25].

undermines the contribution of the scores from each term and is therefore ineffective and is rather unintuitive.

6 TEMPORAL STUDY

In this section, we intend to show the usability of catchphrase extraction. We claim that catchphrase evolution presents a fair understanding of the changing innovation trends of companies. We conduct several interesting temporal studies to understand the emergence of new research topics in the industry. In this study, we select top-10 companies from three industrial segments: (1) Software⁶, (2) Hardware⁷, and (3) Mobile Phones⁸. Table 5 presents the list of top-10 companies in each of the above three segments.

In subsequent sections, we analyze patents filed by these companies over the years. In our patent dataset, each company can have several variations in name due to multiple research groups, geographical locations, subsidiaries, headquarters, etc. For example, IBM is present as ‘International Business Machines Corporation Armonk’, ‘International Business Machines Laboratory Inc.’, etc. We overcome these inconsistencies by manually annotating name variations. However, we claim that basic string matching techniques can easily automate this normalization. Besides, we eliminate frequently occurring catchphrases like, ‘method’, ‘present invention’, etc., to ignore noisy/redundant signals. This filtering process was automated by removing catchphrases with top-10 document frequencies. We next, present how catchphrases can be leveraged in understanding the topical evolution of companies.

6.1 Topic evolution

In this section, we study the topical evolution of companies. We leverage the *Jaccard Similarity* (JS) between the catchphrases to compute the topical overlap between patents filed in consecutive years by a specific company. We conduct this experiment for 11 years between 2006–2016. Figure 2 shows temporal profiles of a three-year moving average over JS for each of the three segments. We observe that Baidu in *Hardware* segment while Oppo, Vivo, and OnePlus in *Mobile Phones* segments exhibit relatively low similarity between catchphrases over the years. However, most of the companies have similarity curves with multiple peaks with an overall increase in the JS values over the years. For this analysis, we only considered 2-gram catchphrases. However, we found similar observations for higher n-gram catchphrases. If an organization is filing patents on the same topics over the years, the JS value will only increase; on the other hand, if an organization is continuously filing patents on newer topics, the JS value is expected to decline.

6.2 Categorization

Further, we conduct a nuanced study to understand this temporal behavior. We classify each company’s similarity profile into five categories [7] based on the number and location of peaks. A peak in the similarity profile of a company represents a high topical similarity between consecutive years followed by a topical drifting off period. We leverage the peak identification method proposed

by Chakraborty et al. [7]. Note that peaks occurring in consecutive years are considered as a single peak. The categories are:

- (1) **MONINC**: Similarity profile that monotonically increases. The peak occurs in the last year.
- (2) **MONDEC**: Similarity profile that monotonically decreases. The peak occurs in the first year.
- (3) **PEAKINIT**: Similarity profile that consists single peak within the first three years but not the first year.
- (4) **PEAKLATE**: Similarity profile that consists single peak after the initial three years but not the last year.
- (5) **PEAKMULT**: Similarity profile consisting of multiple peaks.
- (6) **OTHERS**: Similarity profiles that do not qualify into the above categories are kept in this category. They mainly consist of profiles with extremely low JS values for each year.

Table 6 shows categorization results. We find no company in **MONDEC** and **PEAKINIT** categories. Majority of the companies are present in the **PEAKMULT** category followed by **PEAKLATE** category. Companies in **OTHERS** category have very less number of filed patents. Three out of four companies in **OTHERS** category are recently launched mobile companies.

Even though, **PEAKMULT** category consists multiple peaks, we observe two distinct fluctuation patterns. We term these patterns as (i) **STABLE** and (ii) **UNSTABLE**. In **STABLE**, the profile looks considerably less fluctuating. The profile highly fluctuates in **UNSTABLE** category. We quantify the above fluctuating patterns by leveraging the average value of JS. Given, $JS(c)$ is the similarity profile for a company c , average value of JS ($avg_{JS}(c)$) is computed as:

$$avg_{JS}(c) = \frac{\min(JS(c)) + \max(JS(c))}{2} \quad (13)$$

Empirically, we observe that companies with $avg_{JS} > 0.1$ can be classified as **UNSTABLE**, while the rest can be classified as **STABLE**. Table 7 shows companies in the **PEAKMULT** category that are further categorized into **STABLE** and **UNSTABLE**. Among, **STABLE** and **UNSTABLE** sub-categories, the former contains more (=7) companies than the latter (=5).

6.3 Citation count

Citations, in the scholarly world, determine the popularity of research papers/authors/organizations. Here, we adopt a similar analogy for patent articles. A patent citation is a document cited by an applicant, third party, or a patent office examiner because its content relates to a patent application. We compute the citation count of a patent p by summing the citations received by p . For the current study, we construct citer-cited pairs by extracting references present in patent texts and use these pairs to compute patent citation counts.

Next, we create multiple citation zones based on the citation count of a patent. We define four distinctive zones: (i) very low, (ii) low, (iii) medium, and (iv) high, to study the influence of the JS profile of a company on the number of citations received by its patents. Table 8 presents zoning statistics of the complete dataset. Out of 3,829,153 patent articles, 1,499,175 have zero citation count.

Next, we relate similarity profiles and citation count zones. For each company, we measure the fraction of patents in different citation zones. We leverage histograms as a visualization tool to

⁶https://en.wikipedia.org/wiki/List_of_the_largest_software_companies

⁷<https://www.investopedia.com/articles/investing/012716/worlds-top-10-hardware-companies-aaplibm.asp>

⁸<https://www.researchsnipers.com/top-10-largest-mobile-companies-in-the-world-2018/>

z	PRECISION					RECALL				
	Our Model	KEA	Legal	BM25	KLIP	Our Model	KEA	Legal	BM25	KLIP
10	0.253	0.192	0.075	0.080	0.120	0.773	0.557	0.255	0.265	0.386
15	0.231	0.148	0.128	0.131	0.146	0.910	0.566	0.559	0.567	0.623
20	0.217	0.133	0.156	0.156	0.164	0.945	0.568	0.750	0.749	0.772
15%	0.260	0.200	0.056	0.060	0.108	0.750	0.555	0.172	0.185	0.323
20%	0.240	0.156	0.109	0.111	0.132	0.886	0.563	0.448	0.457	0.528

Table 4: Comparison of our proposed method against the baselines: Precision and recall values at different top-ranks ($z \in \{10, 15, 20, 15\%, 20\%\}$) of extracted catchphrases.

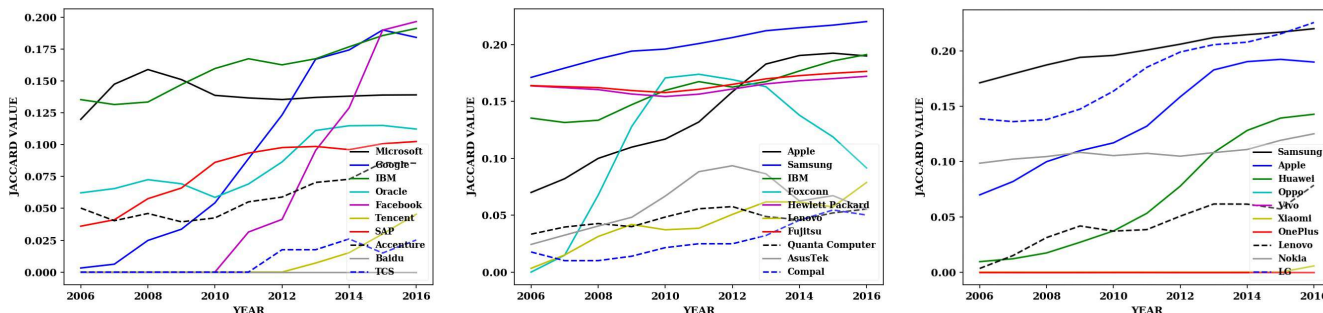


Figure 2: Moving average of catchphrases similarity between consecutive years for – *Software* (left), *Hardware* (center), and *Mobile Phone* (right) companies.

Software	Hardware	Mobile Phone
Microsoft	Apple	Samsung
Google	Samsung	Apple
IBM	IBM	Huawei
Oracle	Foxconn	Oppo
Facebook	Hewlett Packard	Vivo
Tencent	Lenovo	Xiaomi
SAP	Fujitsu	OnePlus
Accenture	Quanta Computer	Lenovo
TCS	AsusTek	Nokia
Baidu	Compal	LG

Table 5: Top-10 *Software*, *Hardware*, and *Mobile Phone* companies selected from three publicly available lists.

Category	Count	Names
MONINC	4	Tencent, Samsung, Xiaomi, Lenovo
MONDEC	0	
PEAKINIT	0	
PEAKLATE	6	Facebook, TCS, Huawei, AsusTek, Foxconn, Compal
PEAKMULT	12	HP, SAP, Accenture, Nokia, Fujitsu, Quanta Computer, Microsoft, IBM, Oracle, Google, Apple, LG
OTHERS	4	Baidu, Oppo, Vivo, OnePlus

Table 6: Categorization of top-10 *Software*, *Hardware*, and *Mobile Phone* companies based on temporal catchphrase similarity profile. No company was classified in *MONDEC* and *PEAKINIT* category.

conduct this study. In Figure 3, we observe that the fraction of patents in *Medium* and *High* citation zones in *PEAKLATE* category are relatively higher than in *MONINC* category. This indicates that the introduction of diversity in topics over time helps in enhancing the future citations of the patents filed by a company.

Figure 4 compares two subcategories of *PEAKMULT*. We observe that the fraction of patent falling under the *Medium*, and *High* citation zones in *UNSTABLE* category is relatively higher than *STABLE* categories implying that the companies with high fluctuations in similarity profiles perform better in terms of receiving citation counts. A possible explanation is that the companies with relatively specialized research domain file patents which attract lesser citations than the companies with diversified research domain.

Lastly, we study *OTHERS* category in Figure 5. Quite surprisingly, we observe that the fraction of patents in *Medium* and *High* citation zones in *OTHERS* category is relatively higher than the rest of the categories described above in Figures 3 and 4.

6.4 Catchphrases in the *STABLE* and *UNSTABLE* groups

In this section, we analyze the extent of usage of certain catchphrases (bigrams and trigrams) by a company. We rank the catchphrases based on document frequency, i.e. the number of patent documents a catchphrase is present in. Tables 9 and 10 show the top-10 bigrams for companies present in the *STABLE* and *UNSTABLE*

Company	avgJS	Category
Nokia	0.040	STABLE
Fujitsu	0.085	STABLE
Quanta Computer	0.069	STABLE
Microsoft	0.105	STABLE
Accenture	0.040	STABLE
SAP	0.048	STABLE
Hewlett Packard	0.084	STABLE
LG	0.223	UNSTABLE
Oracle	0.117	UNSTABLE
Google	0.121	UNSTABLE
Apple	0.197	UNSTABLE
IBM	0.124	UNSTABLE

Table 7: List of companies in PEAKMULT that are classified into STABLE and UNSTABLE sub-categories along with the average value of Jaccard Similarity ($avgJS$) used for categorization.

Category	Citation Count	Patent Count
Very Low	0	1,499,175
Low	$0 < x < 5$	1,274,029
Medium	$5 \leq x < 25$	840,461
High	$x \geq 25$	215,488

Table 8: Patent citation zones with distinct citation count ranges.

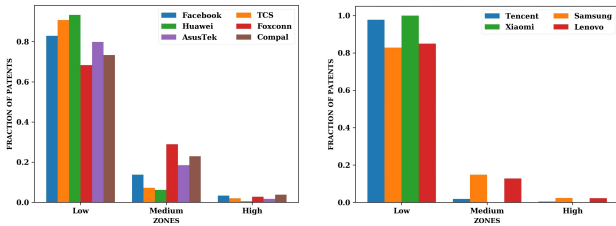


Figure 3: Citation count zones vs similarity profiles: Fraction of patents in PEAKLATE (left) and MONINC (right) category companies in each citation count zone.

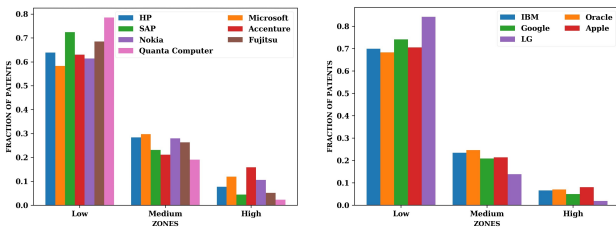


Figure 4: Citation count zones vs similarity profiles: Fraction of patents in STABLE (left) and UNSTABLE (right) category companies in each citation count zone.

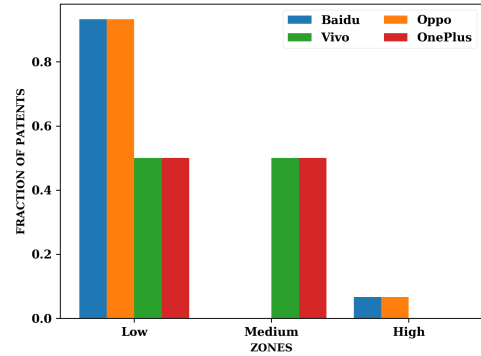


Figure 5: Citation count zones vs similarity profiles: Fraction of patents in OTHERS in each citation count zone.

groups respectively. Table 11 and 12 show top-10 trigrams for the same companies. Last, Table 13 notes the top-10 bigrams and trigrams from the entire stable and unstable categories taken together. While the STABLE group is concerned more about computer systems, the UNSTABLE group is more about electronic device parts.

6.5 Temporal visualizations

In this section, we study the catchphrase evolution of companies. As a popular visualization tool, we leverage *word clouds*. We create word clouds for each company between the years 2003–2016. Due to space constraints, in Figure 6, we only consider word clouds for one representative company from STABLE, UNSTABLE, PEAKLATE and MONINC categories at three representative years. We claim that catchphrase evolution presents a fair understanding of the changing innovation trends of companies. Note that we consider only bigram catchphrases in this study. We can conduct a similar study for any company in different years⁹.

In Figure 6a, we study catchphrase evolution for Microsoft (a representative company in the STABLE group). We observe a shift from traditional topics such as client-server models, databases, basic Web development, etc. (in 2003), toward full-fledged Web search and Internet technologies (in 2010). In 2016, the focus shifted to mobile devices and gesture identification. The above trends coincide with several product releases such as BING (a search engine released in 2009)¹⁰ and Lumia (mobile phones released in 2015)¹¹.

In Figure 6b, we study catchphrase evolution for Oracle (a representative company in the UNSTABLE category). Oracle seems to have shifted its focus from traditional database topics like relational databases, query, etc. (in 2003), toward the development of software as a service (SAAS) in 2010. In 2016, it continued to focus on services with a major emphasis on reliable authentication mechanisms in the cloud. These innovation trends resulted in several products like *Oracle cloud* (cloud computing service launched in 2016), Primavera (an enterprise project portfolio management software acquired by Oracle in 2008), etc.

⁹The detailed word clouds for all companies in our dataset are available at <http://tinyurl.com/y5ynhj9n>

¹⁰[https://en.wikipedia.org/wiki/Bing_\(search_engine\)](https://en.wikipedia.org/wiki/Bing_(search_engine))

¹¹https://en.wikipedia.org/wiki/Microsoft_Lumia_435

HP	MS	SAP	Accenture	Nokia	Fujitsu	Quanta Computer
print job	client device	business process	processing device	user interface	closed position	circuit board
one aspect	search result	application server	third party	communication device	inner surface	display panel
printing system	application program	application program	real-world environment	one embodiment	opposite side	second image
second set	user input	software application	mobile device	computer program	upper surface	one side
operating system	computing system	business application	invention concern	telecommunication system	longitudinal axis	second end
second side	search engine	system method	educational material	telecommunication network	another embodiment	battery module
second position	data store	data structure	computer-implemented method	access point	opposite end	second position
second portion	least portion	system software	communication network	data transmission	open position	one end
display device	subject matter	user input	synchronized video	least part	bottom surface	portable computer
present disclosure	client computer	business object	solution information	second device	side wall	power supply

Table 9: Bi-grams with the top 10 document frequency values in STABLE category.

IBM	Oracle	Google	Apple	LG
top surface	application server	user interface	integrated circuit	second electrode
operating system	operating system	present disclosure	second set	one side
storage device	data structure	one example	user input	common electrode
computer program	second set	system method	one example	light source
second set	software application	example method	first set	lcd device
computing system	one technique	content item	operating system	control information
another embodiment	source code	user input	another embodiment	display device
user interface	computer-implemented method	one processor	least portion	array substrate
drain region	another aspect	user device	host device	drain electrode
data structure	database object	subject matter	client device	washing machine

Table 10: Bi-grams with the top 10 document frequency values in UNSTABLE category.

HP	MS	SAP	Accenture	Nokia	Fujitsu	Quanta Computer
storage area network	host operating system	first data object	dual information system	first base station	user's head	portable electronic apparatus
least one component	client computing device	one general aspect	telecommunication industry taxonomy	packet data network	first second portion	mobile communication device
fluid ejection assembly	mobile communication device	business process model	contact center representative	least one parameter	user's foot	second frequency band
first second set	user's interaction	second user input	contact center system	first network element	least one opening	third conductor arm
least one component	least one implementation	least one service	context-appropriate enforcing completion	user equipment due	thinning spraying irrigation	portable computer system
disclosed embodiment relate	distributed computing system	core software platform	location-based service system	wireless communication device	patient's body	second radiating element
least one surface	application program interface	least one attribute	cognitive educational experience	least one cell	least one side	blade server system
inkjet ink composition	one computing device	second data object	individualized learning experience	wireless communication system	least one aperture	service agent server
central processing unit	client computer system	one exemplary embodiment	user's comprehension	wireless communication device	storied index rating	printed circuit board
graphical user interface	wireless access point	related method system	object recognition analysis	second base station	usda hardness zone	wireless communication device

Table 11: Tri-grams with the top 10 document frequency values in STABLE category.

IBM	Oracle	Google	Apple	LG
first conductivity type	current result list	one search result	electronic device housing	light guide plate
field effect transistor	flexible extensible architecture	first search result	scrolling 3d manipulation	digital broadcasting system
second dielectric layer	computer program product	image sensor interface	intuitive hand configuration	second semiconductor layer
gate dielectric layer	distributed computing environment	disclosed subject matter	hand approach touch	liquid crystal cell
data communication network	graphical user interface	image search result	proximity-sensing multi-touch surface	light emitting diode
integrated circuit device	data storage system	client computing device	wireless communication circuitry	main service data
direct physical contact	data processing system	client computing device	antenna resonating element	image display device
second conductivity type	application programming interface	second computing device	computer readable medium	serving base station
database management system	database management system	mobile communication device	wireless communication system	light emitting diode
buried insulator layer	contention management mechanism	distributed storage system	wireless electronic device	first second electrode

Table 12: Tri-grams with the top 10 document frequency values in UNSTABLE category.

Similarly, in Figure 6c, we study catchphrase evolution for Facebook (a representative company in PEAKLATE category). As Facebook started its operations from 2004, we present visualizations for three years, 2010, 2013 and 2016. The initial focus was to develop technical features like news feeds, membership, etc. In the year 2013, these trends shift toward instant messaging aspects. In the year 2016, the catchphrases show a distinct innovation pattern

of restricting and disclosing data availability. Facebook Messenger (introduced in 2011) is one of the products developed between 2011–2013¹².

We study Samsung as a representative company in MONINC category (see Figure 6d). Primarily Samsung's major focus lies in traditional electronics innovation. Recent trends suggest an increased

¹²<https://en.wikipedia.org/wiki/Facebook>

STABLE		UNSTABLE	
Bi-grams	Tri-grams	Bi-grams	Tri-grams
closed position	first second portion	top surface	second semiconductor layer
another embodiment	user's head	user interface	printed circuit board
opposite side	least one opening	second set	light guide plate
inner surface	least one side	operating system	light emitting diode
upper surface	user's foot	present disclosure	digital broadcasting system
one aspect	least one aperture	least portion	first semiconductor layer
longitudinal axis	patient's body	another embodiment	light emitting diode
open position	thinning spraying irrigation	system method	liquid crystal cell
second position	central processing unit	data structure	first second electrode
opposite end	storie index rating	computing system	serving base station

Table 13: Bi-grams and Tri-grams with the top 10 document frequency values.

focus on mobile technologies such as user interfaces, display units, etc.

7 CONCLUSION AND FUTURE WORK

In this paper, we propose an unsupervised catchphrase identification and ranking system. Our proposed system achieves a substantial improvement, both in terms of precision and recall, against state-of-the-art techniques. We demonstrate the usability of this extraction by analyzing how topics evolve in patent documents and how these evolution patterns shape the future citation count of the patents filed by a company.

In the future, we plan to extend the current work by developing an online interface for automatic catchphrase identification. We also plan to understand the influence of catchphrase evolution on the company's revenue.

REFERENCES

- [1] Daniele Archibugi and Mario Planta. 1996. Measuring technological change through patents and innovation surveys. *Technovation* 16, 9 (1996), 451–519.
- [2] Kendall W Artz, Patricia M Norman, Donald E Hatfield, and Laura B Cardinal. 2010. A longitudinal study of the impact of R&D, patents, and product innovation on firm performance. *Journal of product innovation management* 27, 5 (2010), 725–740.
- [3] Richard A Bettis and Michael A Hitt. 1995. The new competitive landscape. *Strategic management journal* 16, S1 (1995), 7–19.
- [4] Nicholas Bloom and John Van Reenen. 2002. Patents, real options and firm performance. *The Economic Journal* 112, 478 (2002), C97–C116.
- [5] Zvi Boger, Tsvi Kuflik, Peretz Shoval, and Bracha Shapira. 2001. Automatic keyword identification by artificial neural networks compared to manual identification by users of filtering systems. *Information Processing & Management* 37, 2 (2001), 187–198.
- [6] Lutz Bornmann and Hans-Dieter Daniel. 2008. What do citation counts measure? A review of studies on citing behavior. *Journal of documentation* 64, 1 (2008), 45–80.
- [7] Tanmoy Chakraborty, Suhansanu Kumar, Pawan Goyal, Niloy Ganguly, and Animesh Mukherjee. 2015. On the categorization of scientific citation profiles in computer science. *Commun. ACM* 58, 9 (2015), 82–90.
- [8] Hung-Hsuan Chen, Puktada Treeratpituk, Prasenjit Mitra, and C Lee Giles. 2013. CSSeer: an expert recommendation system based on CiteseerX. In *Proceedings of the 13th ACM/IEEE-CS joint conference on Digital Libraries*. ACM, 381–382.
- [9] Yin-Hui Cheng, Fu-Yung Kuan, Shih-Chieh Chuang, and Yun Ken. 2009. Profitability decided by patent quality? An empirical study of the US semiconductor industry. *Scientometrics* 82, 1 (2009), 175–183.
- [10] Andras Csomai and Rada Mihalcea. 2008. Linguistically motivated features for enhanced back-of-the-book indexing. *Proceedings of ACL-08: HLT* (2008), 932–940.
- [11] Eibe Frank, Gordon W Paynter, Ian H Witten, Carl Gutwin, and Craig G Nevill-Manning. 1999. Domain-specific keyphrase extraction. In *16th International joint conference on artificial intelligence (IJCAI 99)*, Vol. 2. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 668–673.
- [12] Parthasarathy Gopavarapu, Line C Pouchard, and Santiago Pujol. 2016. Increasing datasets discoverability in an engineering data platform using keyword extraction. In *Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries*. ACM, 225–226.
- [13] Christine Greenhalgh and Mark Longland. 2005. Running to stand still?—the value of R&D, patents and trade marks in innovating manufacturing firms. *International Journal of the Economics of Business* 12, 3 (2005), 307–328.
- [14] Bronwyn H Hall, Adam Jaffe, and Manuel Trajtenberg. 2005. Market value and patent citations. *RAND Journal of economics* (2005), 16–38.
- [15] Constance E Helfat and Margaret A Peteraf. 2003. The dynamic resource-based view: Capability lifecycles. *Strategic management journal* 24, 10 (2003), 997–1010.
- [16] Steve Jones and Gordon Paynter. 1999. Topic-based browsing within a digital library using keyphrases. In *Proceedings of the fourth ACM conference on Digital libraries*. ACM, 114–121.
- [17] Bruce Krulwich and Chad Burkey. 1997. The InfoFinder agent: Learning user interests through heuristic phrase extraction. *IEEE Expert* 12, 5 (1997), 22–27.
- [18] Leah S Larkey. 1999. A patent search and classification system. In *Proceedings of the fourth ACM conference on Digital Libraries*. ACM, 179–187.
- [19] Changyong Lee, Yangrae Cho, Hyeonju Seol, and Yongtae Park. 2012. A stochastic patent citation analysis approach to assessing future technological impacts. *Technological Forecasting and Social Change* 79, 1 (2012), 16–29.
- [20] Changyong Lee, Ohjin Kwon, Myeongjung Kim, and Daeil Kwon. 2018. Early identification of emerging technologies: A machine learning approach using multiple patent indicators. *Technological Forecasting and Social Change* 127 (2018), 291–303.
- [21] Arpan Mandal, Kripabandhu Ghosh, Arindam Pal, and Saptarshi Ghosh. 2017. Automatic Catchphrase Identification from Legal Court Case Documents. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. ACM, 2187–2190.
- [22] Olena Medelyan, Eibe Frank, and Ian H Witten. 2009. Human-competitive tagging using automatic keyphrase extraction. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 3-Volume 3*. Association for Computational Linguistics, 1318–1327.
- [23] Chau Q Nguyen and Tuoi T Phan. 2009. An ontology-based approach for key phrase extraction. In *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*. Association for Computational Linguistics, 181–184.
- [24] Sooyoung Oh, Zhen Lei, Prasenjit Mitra, and John Yen. 2012. Evaluating and ranking patents using weighted citations. In *Proceedings of the 12th ACM/IEEE-CS joint conference on Digital Libraries*. ACM, 281–284.
- [25] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
- [26] Rushdi Shams and Robert E Mercer. 2012. Investigating keyphrase indexing with text denoising. In *Proceedings of the 12th ACM/IEEE-CS joint conference on Digital Libraries*. ACM, 263–266.
- [27] Takashi Tomokiyo and Matthew Hurst. 2003. A language model approach to keyphrase extraction. In *Proceedings of the ACL 2003 workshop on Multiword expressions: analysis, acquisition and treatment-Volume 18*. Association for Computational Linguistics, 33–40.
- [28] Suzan Verberne, Maya Sappelli, Djoerd Hiemstra, and Wessel Kraaij. 2016. Evaluation and analysis of term scoring methods for term extraction. *Information Retrieval Journal* 19, 5 (2016), 510–545.
- [29] Ian H Witten and Olena Medelyan. 2006. Thesaurus based automatic keyphrase indexing. In *Proceedings of the 6th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL'06)*. IEEE, 296–297.
- [30] Ian H Witten, Gordon W Paynter, Eibe Frank, Carl Gutwin, and Craig G Nevill-Manning. 2005. KEA: Practical Automated Keyphrase Extraction. In *Design and Usability of Digital Libraries: Case Studies in the Asia Pacific*. IGI Global, 129–152.