



## Automatic extraction of significant terms from the title and abstract of scientific papers using the machine learning algorithm: A multiple module approach

Bhaskar Mukherjee<sup>a</sup> and Debasis Majhi<sup>b</sup>

<sup>a</sup>Professor, Department of Library & Information Science, Banaras Hindu University, Varanasi,  
Email: mukherjee.bhaskar@gmail.com

<sup>b</sup>Junior Research Fellow, Department of Library & Information Science, Banaras Hindu University, Varanasi,  
Email: debasismajhidlis@gmail.com

*Received: 09 February 2023; accepted: 27 March 2023*

Keyword extraction is the task of identifying important terms or phrase that are most representative of the source document. Although the process of automatic extraction of keywords from title is an old method, it was mainly for extraction from a single web document. Our approach differs from previous research works on keyword extraction in several aspects. For those who are non-expert of the scientific fields, understating scientific research trends is difficult. The purpose of this study is to develop an automatic method of obtaining overviews of a scientific field for non-experts by capturing research trends. This empirical study excavates significant term extraction using Natural Language Processing (NLP) tools. More than 15000 titles saved in a .csv file was our dataset and scripts written in Python were our process to compare how far significant terms of scientific title corpus are similar or different to the terms available in the abstract of that same scientific article corpus. A light-weight unsupervised title extractor, Yet Another Keyword Extractor (YAKE) was used to extract the results. Based on our analysis, it can be concluded that these algorithms can be used for other fields too by the non-experts of that subject field to perform automatic extraction of significant words and understanding trends. Our algorithm could be a solution to reduce the labour-intensive manual indexing process.

**Keywords:** Data mining, Title extraction, Natural Language Processing, YAKE, NLTK, Keyword Extraction-NLP

### Introduction

This paper is an attempt to extract significant keywords from the title and abstract of a scientific paper by using machine learning process to understand the trends in research on Bioinformatics. While the title of a research paper reveals the thought contents of the overall research, the abstract represents the gist findings of the scientific research. Often, authors choose title carefully to express the research correctly, however, number of times during reviews, the reviewers or the editors suggest alternative title and ask author to change the title without changing the intention of the paper. The reason behind such suggestions is mainly because of compliance with the journal's guidelines. Other reasons may be to keep those terms in the title that searchers usually approach on the online platforms so that the research paper can be accessed by the searcher. It may also be to avoid any such words that are not part of the research or confusing in nature.

On web, search engines play an important role in searching precise results. The search engine bot or spider crawls hundreds of billion pages using their

web crawler. In getting search engine optimization (SEO) success, link building between search query and metatags of a research paper is crucial. Among various metatags, title tag is the most important anchor as it is meant to provide a clear and comprehensive idea of what the research paper is all about. Over the past few years, user search behaviour factors were being discussed a lot as logical proof of choosing right title, by contrast authors sometime chose interesting, catchy title to gather attention. More than the author names, the title affects how a paper is perceived by others. Whimsical, amusing or hilarious titles though appealing, may be wrongly categorized in automatic indexing process and may be cited less often<sup>1</sup>.

Term extraction from the titles of the scientific research papers is not a complex task in present days. There are so many simple algorithms available through which we can extract single-word terms from the title. However, extraction of multiple terms through machine learning algorithm with correct semantic value is a complex task. Because a single word when combined with other words reveals

different meaning. For example, 'Information Storage and Retrieval' represent a different context than the words like 'information', 'storage' or 'retrieval'. 'Citation metrics' have different connotation than the word 'citation' and 'metrics'. For creating a meaningful network diagram, significant words with correct semantics are more important than creating cluster of significant terms.

Like title, abstract is yet another important element of a research paper to include important words related to methodology, results and findings. For appropriate indexing purpose and for retrieval, abstract plays an important role in actual data representation and interpretation<sup>2</sup>. The title and abstract of a research paper are often freely available to the readers and easily accessible through journal webpages, indexing databases or search engine result page (SERP). In this paper, we take a machine learning approach to address the problem and compare the output by developing frameworks using two modules of python language. The key issues are to define a better specification on machine extraction by comparing the results of title and abstract extraction of research papers.

With the advancement of science, more and more interdisciplinary research is being carried out throughout the world. Biological research nowadays generates large volume of biomolecular data. The open sharing of these data is crucial towards groundbreaking scientific achievements<sup>3</sup>. The reason behind choosing bioinformatics as a field to test our frameworks was that this recently emerging interdisciplinary field has wide application in human endeavour which collects, stores, analyses and disseminates data on DNA and amino acid sequence using computer technology and statistics.

### Review of literature

Xue<sup>4</sup> conducted a study for HTML text extraction of web pages by using supervised machine learning approach. They found that the use of both extracted titles and title fields is always better than use of title field alone. And the use of extracted titles is particularly helpful in the task of name page findings. In contrast, Gali<sup>5</sup> observed that HTML text extraction from service-based web pages fails. Because advertisements in the service-based web pages often give more visual emphasis than main headings. Therefore, they proposed a novel method that combines statistical features, linguistic knowledge,

and text segmentation. In another study Uzun<sup>6</sup> presented a hybrid approach that contains two steps. The first step discovers informative contents using Decision Tree Learning of machine learning method and creates rules from the results. The second step extracts informative contents using rules obtained from the first step.

There are few other studies which deals with title extraction process. For example, Giuffrida et al.<sup>7</sup> extracted contents from PostScript files using postotext tool, another rule-based system PDFX is described by Constatin<sup>8</sup>. Tkaczyk<sup>9</sup> developed a comprehensive open-source system for extracting metadata and parsed bibliographic references from scientific articles in born digital format.

Rinartha and Kartika<sup>10</sup> suggested a model for keyword extraction by Rapid Automatic Keyword Extraction (RAKE). They suggested that word frequency and RAKE can be combined to get more accurate results in extracting keywords from title. Wand<sup>11</sup> developed an algorithm that can automatically extract keywords from the meta-information of each article and generate the basic data for review articles by implementing Rapid Automatic Keyword Extraction algorithm on the title and abstract of each article. Finally, they classified all extracted keywords into class by calculating Levenshtein distance between each of them.

Nakajima<sup>12</sup> performed an analysis using 12 years of articles on high-temperature superconductors and developed a method by examining research article, review literature and co-citation among research articles. In another study, Gunawan<sup>13</sup> applied TextRank algorithm to extract the keywords and then generate keyphrase based on the rank on the vertices. Their results explained that the number of assigned keywords play a crucial role in determining the recall value. The recall value of TextRank algorithm raises from 38.4% (5 keywords) to 55.26% (10 keywords) and 61.54% (15 keywords). Jiang<sup>14</sup> suggested a model of automatic extraction keywords in the field third-generation semiconductor materials based on entity recognition (ER) and relation extraction (RE) techniques. Based on this ER and RE, a neural network using domain knowledge (DKNet) is proposed to improve ER performance. Adjusting keyword sequence of each entity type as prior knowledge, adds a dedicated embedding to encode entity categories, then combines prior knowledge and encoded vectors with the context through a gated

information fusion module to assist recognition. By applying this process, they have suggested a multi-aspect attention-based network model (MANet) to enhance the attention to relation-indicative words.

### Objectives of the study

- To understand the text-characteristics in scientific titles and their abstracts for a multidisciplinary field like bioinformatics; and
- To compare the output of N-gram text extraction results by developing algorithm through YAKE module.

### Methodology

This research is based on Natural Language Processing (NLP) applications for automatic text extraction from scientific articles. For comparing the results of different modules of python compatible modules for text extraction, we needed raw data. We downloaded 15000 records from Web of Science (WoS) database on bioinformatics for the period 2021 to 2022. We used Pandas, a library in python that has functions like analysing, cleaning, exploring, and processing data. Along with pandas, we have used different modules of Natural Language Toolkit (NLTK) like PorterStemmer, WordNetLemmatizer, Yake, CounterVector modules for title extraction. The WordNetLemmatizer was used to group together the different inflected forms of a word so that they can be considered as a single word (eg. Rocks->Lemmatization -> rock; corpora-> Lemmatization -> corpus; Caring -> Lemmatization -> Care etc.). On the other hand, PorterStemmer was used to stem the word so that not only noting but chopping the word into its root form could be achieved (eg. "python", "pythoner", "pythoning", "pythoned" ->PorterStemmer -> python).

Stop words or bag of insignificant words are in-built feature of NLTK which can eliminate commonly used English stop words from title. But on analysing those words, we observed that the lists of such words are insufficient. We therefore developed a dictionary of stop words and run the same in python library. Our newly developed stop wordlist also includes those words which have no significance in scientific titles (eg. 'named', 'formerly', 'accordingly', 'across', 'actually' etc.).

By using above mentioned process, we converted the text into small segments, called tokens, consisting of words. Most of these words are unigram or a single word. As most of the word thus segmented do not serve our needs, we attempted to use tools for bigrams and trigrams. By fixing the lower and upper boundary of n-values, different n-grams were generated using YAKE and CounterVector modules. Yake is a light-weight unsupervised automatic keyword extraction of NLTK of Pandas which relies on statistical text features extracted from a single document to select most relevant keywords of a text<sup>15</sup>. On the other hand, CounterVector is another model of text classification where text data have been vectorized and total occurrence of word can be calculated. It is needless to mention here that NLP models do not understand the textual data, therefore, it is essential to convert such data to vector/numbers.

### Results

Figure 1 displays the results of our program and Table 1 analyses the result. It is seen that research article titles on an average consists of 14 words (Sd 4.2), out of which almost 70% words are significant (Sd 2.9) and remaining 30% are stop words. The average number of words per abstract is 220 and the average number of significant words is 133.

```

IDLE Shell 3.9.6
File Edit Shell Debug Options Window Help
Python 3.9.6 (tags/v3.9.6:db3ff76, Jun 28 2021, 15:26:21) [MSC v.1929 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:\Users\HF\AppData\Local\Programs\Python\Python39\Final Codes\term-count and seperate-csv to py and csv--important.py
      Title                                     Des
0      Improving bioinformatics software quality thro... [improving, bioinformatics, software, quality,...
1      Principles and Validation of Bioinformatics Pi... [principles, validation, bioinformatics, pipel...
2      General considerations for online teaching pra... [general, considerations, online, teaching, pr...
3      Community development, implementation, and ass... [community, development,, implementation,, ass...
4      Visualizing the knowledge structure and evolut... [visualizing, knowledge, structure, evolution,...
...
14995 AMPK upregulates K(Ca)2.3 channels and amelior... [ampk, upregulates, k(ca)2.3, channels, amelio...
14996 Molecular Cloning and Expression Analysis of A... [molecular, cloning, expression, analysis, aux...
14997 Identification of an epithelial-mesenchymal tr... [identification, epithelial-mesenchymal, trans...
14998 Long intergenic non-protein coding RNA 1094 (L... [long, intergenic, non-protein, coding, rna, l...
14999 Construction and Analysis of mRNA and lncRNA R... [construction, analysis, mrna, lncrna, regulat...

[15000 rows x 2 columns]
>>>

```

Fig. 1 — Text extraction of significant words excluding stop words

Although a script can be written for identifying the common and uncommon words of a title but the output results we got by applying those script did not fulfil our requirements, and therefore was discarded.

As seen in Fig. 1, although titles were broken into single significant words but the same does not reveal the actual context of the subject. For example, in the line 0 we got a result for the title “Proteomics analysis of lipid droplets in mouse neuroblastoma cells” as ['proteomics', 'analysis', 'lipid', 'droplets', 'mouse', 'neuroblastoma', 'cells']. The term lipid and droplets as single term have no significance until we got our results as ‘lipid droplets’. Similarly for the title “The expression and role of tubulin polymerization-promoting protein 3 in oral squamous cell carcinoma” extraction of each significant word does not serve our purpose if we do not get the results like ‘squamous cell carcinoma’ as a single significant term. Furthermore, term extraction in this process does not provide the overall frequency of significant terms to understand the trends in research.

To solve the problem, we have deployed N-gram technique. In natural language processing, N-grams are continuous sequence of n-number of neighbouring words or tokens in a document. There are several ways N-gram can be deployed. For this study, we have tested N-gram using two different modules of python: one for counting the frequency and another for measuring the most and least important terms. CounterVector is a module which helps in transforming a given text into vector based on its frequency. Scikit-learn’s CounterVector with N-gram features of panda module have been deployed here to extract significant keywords from a corpus of text and count its frequency. In Figure 2, the extracted results of our dataset are shown. Here, results show bi-grams like ‘hepatocellular carcinoma’, ‘breast cancer’, ‘cell carcinoma’ along with their frequency in the whole corpus of titles. Along with that, tri-grams like ‘squamous cell carcinoma’ was also extracted from the text corpus. It is seen that ‘bioinformatics analysis’ (1084 times), ‘hepatocellular carcinoma’ (460) and ‘breast cancer’ (445 times), are

Table 1—Frequency of significant words in titles and abstracts of scientific articles

	Title Characters	Title words	Significant words (Title)	Abstract words	Significant words (Abstract)
Average Frequency	113.23	14.39	10.34	219.94	133.39
Sd	32.19394	4.203348	2.942147	50.88	32.71

```

IDLE Shell 3.9.6
File Edit Shell Debug Options Window Help
Python 3.9.6 (tags/v3.9.6:db3ff76, Jun 28 2021, 15:26:21) [MSC v.1929 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:\Users\HP\AppData\Local\Programs\Python\Python39\Final Codes\counting significant terms - n grams-bio-info.py
frequency          bigram/trigram
0      1084      bioinformatics analysis
1      460      hepatocellular carcinoma
2      445      breast cancer
3      436      cell carcinoma
4      350      gastric cancer
5      300      colorectal cancer
6      286      integrated bioinformatics
7      283      lung adenocarcinoma
8      273      squamous cell
9      268      genome wide
10     261      squamous cell carcinoma
11     247      cancer cells
12     246      non coding
13     245      sars cov
14     245      lung cancer
15     234      signaling pathway
16     218      cell proliferation
17     211      long non
18     210      long non coding
19     206      immune infiltration
20     205      proliferation migration
21     204      gene expression
22     198      integrated bioinformatics analysis
23     195      hub genes
24     181      coding rna
>>>

```

Fig. 2 — Text extraction of significant words from Title using CounterVector

the top three significant words in our dataset. To make a cross-check, we have also tallied these frequencies with the actual dataset available in .csv format and found that the extracted terms and frequency remain same in both.csv files and extracted results of our program.

Further, we have attempted to apply the same code for extraction of frequently used terms in abstracts of the titles. Figure 3 shows the result. Comparing the rank of abstract terms with the title terms, it was seen that there is difference in their rank. The term ‘bioinformatics analysis’ ranks first in title but it ranks second in abstract. Similarly the top ranked term ‘gene expression’ in the abstract is the 13<sup>th</sup> ranked term in title.

Only 36% (9 out of top 29) terms were similar in both title and abstract extracted terms and remaining 64% terms are unique in both title and abstract. And some of the terms in abstract were also repeatedly expressed, like ‘protein-protein’, ‘protein interaction’ and ‘protein-protein interaction’, however such repetition is not seen in title extracted terms.

*Numbers in brackets are the frequency*

Next, we have attempted to examine how far the algorithm generated keywords are like author generated keywords. From the results as shown in the table 2, algorithm generated keywords are quite accurate and occurrence of such keywords are in

```

IDLE Shell 3.9.6
File Edit Shell Debug Options Window Help
Python 3.9.6 (tags/v3.9.6:db3ff76, Jun 28 2021, 15:26:21) [MSC v.1929 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:\Users\HP\AppData\Local\Programs\Python\Python39\Final Codes\counting abstract terms - n grams-bio-info.py
frequency          bigram/trigram
0          5123      bioinformatics analysis
1          3978      differentially expressed
2          3785      gene expression
3          2888      hub genes
4          2126      expressed genes
5          2093      signaling pathway
6          1908      differentially expressed genes
7          1880      enrichment analysis
8          1791      cell proliferation
9          1777      breast cancer
10         1681      expression levels
11         1553      protein protein
12         1472      qrt pcr
13         1457      protein interaction
14         1424      migration invasion
15         1385      cell lines
16         1364      genes degs
17         1347      proliferation migration
18         1334      gene ontology
19         1330      protein protein interaction
20         1275      covid 19
21         1253      expressed genes degs
22         1199      expression omnibus
23         1194      gene expression omnibus
24         1160      luciferase reporter
>>>
    
```

Fig. 3 — Keywords extraction from Abstract through CounterVector (top 25 keywords)

Table 2—Comparison between algorithm-generated title & abstract keywords with author-generated keywords

Author produced keywords	Title extracted keywords	Abstract extracted keywords
Bioinformatics (2596)	bioinformatics analysis (1084)	bioinformatics analysis (5123)
bioinformatics analysis (910)	hepatocellular carcinoma (460)	differentially expressed (3978)
Prognosis (805)	breast cancer (445)	gene expression (3785)
Biomarker (428)	cell carcinoma (436)	hub genes (2888)
Proliferation (329)	gastric cancer (350)	expressed genes (2126)
Proteomics (313)	colorectal cancer (300)	signaling pathway (2093)
hepatocellular carcinoma (284)	integrated bioinformatics (286)	differentially expressed genes (1908)
machine learning (272)	lung adenocarcinoma (283)	enrichment analysis (1880)
breast cancer (270)	squamous cell (273)	cell proliferation (1791)
Immune infiltration (265)	genome wide (268)	breast cancer (1777)

accordance with the author generated keywords. This justifies the validity of our algorithm.

Although CounterVectorize clarifies the most frequent words in the text corpus but unable to identify the more important or less important (significant) keywords of scientific research titles. To overcome the drawbacks, we have written a python script using YAKE (Yet Another Keyword Extraction) module of panda that uses statistical features of a single document to get better results. Figure 4 shows our output.

On analyzing the results of N-gram as generated through YAKE, we observed that N-grams created by the YAKE are more appropriate than others, like SpaCy, RAKE, Gensim. In our result the term 'cell lung cancer' scored  $2.344861670606526e-06$  followed by 'integrated bioinformatics analysis' with score  $2.3901992076133613e-06$ . In our dataset of 15000 titles, these two terms came 126 and 74 times respectively. In YAKE, lower the score, the more relevant keyword is.

To explore to what extent the title extracted keywords are similar with the abstract extracted keywords, we have deployed the abstracts of the titles in YAKE module. As shown in Figure 5, abstracts produced a different set of words than title. Only 24% terms have similarity in both title and abstract.

Remaining 76% terms are different in these two sets of results. It is interesting to note here that the term 'cell lung cancer' received highest attention in terms from title extraction process but does not come under top 25 abstract extracted words. Similarly, variation in results have been observed for the term 'gene expression'.

## Discussion

Proper indexing of research paper titles is essential for make it searchable through search engines. However, for humans, it is quite difficult to handle such a huge volume of data and index it accurately. This research is an attempt to understand whether it is possible by machine to extract correct keywords that represent the actual content of the scientific titles.

In this work, we present an algorithm that generate data to be further used for identifying significant areas of a research field of the titles by leveraging automatic keyword extraction and similarity calculation. We used three different tools of natural language processing viz, simple, moderate and sophisticated, for this purpose and compared the results of automatic extraction from title and abstract of research articles. The reason behind choosing YAKE for keyword extraction was that it followed unsupervised algorithm, and there was no need to train the corpus. It does not need any pre-defined

```

IDLE Shell 3.9.6
File Edit Shell Debug Options Window Help
Python 3.9.6 (tags/v3.9.6:db3fff76, Jun 28 2021, 15:26:21) [MSC v.1929 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:\Users\HP\AppData\Local\Programs\Python\Python39\Final Codes\yake-term frequency.py
Keyphrase: cell lung cancer : score: 2.344861670606526e-06
Keyphrase: integrated bioinformatics analysis : score: 2.3901992076133613e-06
Keyphrase: bioinformatics analysis : score: 2.884403078221545e-06
Keyphrase: squamous cell carcinoma : score: 3.2623199137346124e-06
Keyphrase: cancer bioinformatics analysis : score: 4.0794758246379234e-06
Keyphrase: cell renal cell : score: 6.685372542626656e-06
Keyphrase: identification expression analysis : score: 7.177257149219912e-06
Keyphrase: breast cancer cells : score: 7.205535966323828e-06
Keyphrase: cancer cell proliferation : score: 7.33475243723005e-06
Keyphrase: renal cell carcinoma : score: 7.599641059439793e-06
Keyphrase: cell carcinoma : score: 9.21351692751987e-06
Keyphrase: cancer cells : score: 9.842377942530743e-06
Keyphrase: comprehensive bioinformatics analysis : score: 1.042616390081492e-05
Keyphrase: bioinformatics analysis reveals : score: 1.1078026579082515e-05
Keyphrase: cancer integrated bioinformatics : score: 1.2126987566445583e-05
Keyphrase: identification key genes : score: 1.2449271523862322e-05
Keyphrase: lung cancer cells : score: 1.2947714442044728e-05
Keyphrase: bioinformatics analysis identification : score: 1.3310232829966296e-05
Keyphrase: colorectal cancer cells : score: 1.538596364707445e-05
Keyphrase: cell carcinoma cells : score: 1.558561058152372e-05
Keyphrase: analysis : score: 1.6334905361925996e-05
Keyphrase: breast cancer : score: 1.6713321420877236e-05
Keyphrase: cell carcinoma bioinformatics : score: 1.676827502436704e-05
Keyphrase: immune cell infiltration : score: 1.6967156233554673e-05
Keyphrase: cancer cells targeting : score: 1.702486124254899e-05
>>> |

```

Fig. 4 — Keywords extraction from the Title using YAKE (top 25 keywords)

```

IDLE Shell 3.9.6
File Edit Shell Debug Options Window Help
Python 3.9.6 (tags/v3.9.6:db3ff76, Jun 28 2021, 15:26:21) [MSC v.1929 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:\Users\HP\AppData\Local\Programs\Python\Python39\Final Codes\yake-abstract-term frequency.py
Keyphrase: differentially expressed genes : score: 2.2805602220349346e-07
Keyphrase: gene expression : score: 3.133170846807409e-07
Keyphrase: gene expression omnibus : score: 3.225214714900365e-07
Keyphrase: gene expression analysis : score: 4.3678232951949966e-07
Keyphrase: bioinformatics analysis : score: 4.381256553355884e-07
Keyphrase: gene expression data : score: 5.871991477971958e-07
Keyphrase: analysis gene expression : score: 6.213382434009784e-07
Keyphrase: gene expression profiles : score: 9.422045908367214e-07
Keyphrase: data gene expression : score: 1.1324554993231634e-06
Keyphrase: bioinformatics analysis revealed : score: 1.143191565365895e-06
Keyphrase: genes : score: 1.373552429352326e-06
Keyphrase: gene expression levels : score: 1.4926376263698394e-06
Keyphrase: expressed genes : score: 1.4970997652178773e-06
Keyphrase: hub genes : score: 1.5652857290846443e-06
Keyphrase: encyclopedia genes genomes : score: 1.6556494308785462e-06
Keyphrase: genes gene expression : score: 1.8052643190615083e-06
Keyphrase: expression hub genes : score: 1.8077721029299807e-06
Keyphrase: expression : score: 1.8152617893290628e-06
Keyphrase: expression levels genes : score: 1.8468228258474285e-06
Keyphrase: pathway enrichment analysis : score: 1.8619425812665415e-06
Keyphrase: cancer cells : score: 1.8846374498298157e-06
Keyphrase: protein expression levels : score: 1.8920837965894199e-06
Keyphrase: immune cell infiltration : score: 1.9002653142879986e-06
Keyphrase: analysis : score: 1.9796023487035342e-06
Keyphrase: gene set enrichment : score: 2.1077817640931045e-06
>>>

```

Fig. 5 — Keywords extraction from the Abstract using YAKE (top 25 keywords)

rules, dictionary or thesaurus. It uses statistical features from the text and can extract keywords from a large dataset without re-training. Other procedure of NLP techniques such as part-of-speech tagging (PoS), name-entry recognition (NER), normalization, linguistic parsing or stemming requires specialized tools<sup>15</sup>. For counting the score, YAKE judges the 5 features of a word: casing, word position, word frequency, word relatedness to context, and word different sentence. These scores gained from the 5 features combined into a single score for each keyword. This process was assumed to be better than the Term Frequency-Inverse Document Frequency (TF-IDF) algorithm.

Natural Language Processing although allows machine to interpret human language, however for retrieval, common words (like articles, propositions, pronouns, conjunctions etc.) never serve any purpose. Therefore, text pre-processing is essential for developing effective machine extraction. We observed that scientific research articles usually contain  $\pm 4$  such common words. While the readily available stop words list that are derived from non-technical resources are available in NLP domain, rigorously identified subject specific insignificant, uninformative stop words list is needed for effective text extraction process.

For the present research, although we have developed an additional list of almost 535 such stop

words which are insignificant in scientific field, however we observed there is no universal rule for some terms like ‘analysis’, ‘differential’ ‘invasion’, ‘network’ ‘computation’ etc. These words can seriously harm topic formation because they create spurious co-occurrence of un-related words. They may probably occupy the top positions of the multiple topics, consequently resulting in inscrutable topic of a domain. Therefore, sophisticated statistical metrics from term frequency is needed for better precision than recall.

Our analysis also shows that machine extracted keywords are quite similar with the author-provided keywords and sometimes machine extracted keywords are more accurate with the actual words of scientific titles. How far these techniques are useful for subjects like humanities or social science by maintaining the thematic structure of the titles needs to be excavated.

We have used the YAKE modules of NLTK to extract the significant words from the titles and abstracts of a corpus of scientific titles. Although our results shows a fair picture to understand the trend of a subject but this technique is inappropriate for extraction of keyword at micro level of granularity (line-by-line) for a multiple row based file. Although the keyword-based approach has the advantage of identifying the uses of technologies for large corpus of data without expert knowledge of the domain, but

it cannot easily represent how technologies are used in their area. New techniques like semantic relations between technologies through a combination of linguistic patterns and statistical methods and Subject-Action-Object (SAO) network analysis have been evolved, considerable number of research are yet to be conducted with large scale real dataset.

### Conclusion

With rapidly growing research literature in every field, indexing keywords manually is an impossible task. The CounterVector algorithm used in this paper can be used for generating keywords in other fields. The algorithm could be a solution to reduce the labour-intensive manual indexing process. Our algorithm has certain limitations which requires human experts to illustrate significance of some keywords.

### References

- 1 Bavdekar S B, Formulating the right title for a research article, *Journal of the Association of Physicians of India*, 64 (2016) 53–6.
- 2 Alexandrov AV and Hennerici MG, Writing good abstracts, *Cerebrovascular Discovery*, 23 (2007) 256–59.
- 3 Vamathevan J, Apweiler R and Birney E, Biomolecular Data Resources: Bioinformatics Infrastructure for Biomedical Data Science, *Annual Review of Biomedical Data Science*, 2 (2019) 199–222.
- 4 Xue Y, Hu Y, Xin G, Song R, Shi S, Cao Y, Lin CY and Li H, Web page title extraction and its application, *Information Processing & Management*, 43 (2007) 1332–47. <https://doi.org/10.1016/j.ipm.2006.11.007>.
- 5 Gali N, Content-Based Title Extraction from Web Page. *12th International Conference on Web Information Systems and Technologies*, 2016, pp. 204–10. <https://doi.org/10.5220/0005794102040210>.
- 6 Uzun E, Agun HV and Yerlikaya T, A hybrid approach for extracting informative content from web pages, *Information Processing & Management*, 49(4) (2013) 928–44. <https://doi.org/10.1016/j.ipm.2013.02.005>.
- 7 Giuffrida G, Shek E C and Yang J, Knowledge-based metadata extraction from postscript files. *Proceedings of the fifth ACM conference on Digital Libraries*, (2000) 77–84. <https://doi.org/10.1145/336597.336639>
- 8 Constantin A, Pettifer S and Voronkov A, PDFX: fully-automated pdf-to-xml conversion of scientific literature. *ACM Symposium on Document Engineering*, (2013) 177–180.
- 9 Tkaczyk P, Szostek P, Dendek J, Fedoryszak M and Bolikowski L, CERMINE -- Automatic Extraction of Metadata and References from Scientific Literature, *11th IAPR International Workshop on Document Analysis Systems, Tours, France* (2014) 217–221. doi: 10.1109/DAS.2014.63.
- 10 Rinatha, Komang and Kartika L G S, Rapid Automatic Keyword Extraction and Word Frequency in Scientific Article Keywords Extraction. *3rd International Conference on Cybernetics and Intelligent System (ICORIS)*, (2021) 1–4. *IEEE Xplore*, <https://doi.org/10.1109/ICORIS52787.2021.9649458>.
- 11 Wang J, Su G, Wan C, Huang X and Sun L, A keyword-based literature review data generating algorithm—analyzing a field from scientific publications, *Symmetry*, 12(6) (2020) 903, <https://doi.org/10.3390/sym12060903>.
- 12 Nakajima R and Nobuyuki M, Topic Extraction to Provide an Overview of Research Activities: The Case of the High-Temperature Superconductor and Simulation and Modelling, *Journal of Information Science*, 47(5) (2021) 590–608. <https://doi.org/10.1177/0165551520920794>.
- 13 Gunawan D, Purnamasari F, Ramadhiana R and Rahmat RF, Keyword extraction from scientific articles in Bahasa Indonesia using TextRank algorithm, In *4<sup>th</sup> International Conference on Electrical, Telecommunication and Computer Engineering (ELTICOM)*, (2020) 260–64.
- 14 Jiang X, He K and Yang B, Automatic information extraction in the third-generation semiconductor materials domain based on DKNet and MANet, *IEEE Access*, 10 (2022) 29367–76.
- 15 Campos R, Mangaravite V, Pasquali A, Jorge A, Nunes C and Jatowt A, YAKE! keyword extraction from single documents using multiple local features, *Information Sciences*, 509 (2020) 257–89.