



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 5, Issue 6, June 2017

## A Speech-Based Just-In-Time Retrieval System

Devyani Kolhal<sup>1</sup>, Prof. Ritesh Thakur<sup>2</sup>

M.E, Department of Computer Engineering, IOK College of Engineering, Savitribai Phule Pune University, Pimple  
Jagtap Pune, Mahatashtra, India.<sup>1</sup>

HOD, Department of Computer Engineering, IOK College of Engineering, Savitribai Phule Pune University, Pimple  
Jagtap Pune, Mahatashtra, India<sup>2</sup>

**ABSTRACT:** This paper addresses the problem of key-word extraction from conversations, with the objective of utilizing those watchwords to get better, for each quick discussion piece, a bit number of conceivably pertinent reviews, which may be prescribed to members. Anyways, even a brief piece contains a mixed bag of phrases, which might be conceivably diagnosed with a few issues; also, using a programmed discourse acknowledgment (ASR) framework offers slips amongst them. along those lines, it is hard to surmise efficaciously the statistics wishes of the discussion members. We first recommend a calculation to take away decisive phrases from the yield of an ASR framework (or a guide transcript for testing), which makes usage of topic demonstrating techniques and of a sub modular prize capability which supports differing qualities inside the magic phrase set, to coordinate the potential differing characteristics of topics and reduce ASR commotion. At that point, we recommend a method to deduce numerous topically remote inquiries from this decisive word set, preserving in mind the stop goal to expand the possibilities of creating at any fee one pertinent proposal while making use of those inquiries for over the English Wikipedia. In this if any query has not found document relevant to it, then we apply that single query to search engine dataset to retrieve the document related to that. The proposed systems are assessed as a long way as significance as for dialogue pieces from the Fisher, AMI, and ELEA conversational corpora, appraised by using some human judges. The scores display that our proposition movements forward over past systems that remember simply word recurrence or subject closeness, and speaks to a promising answer for a file recommender framework to be applied as a part of discussions.

**KEYWORDS:** Document recommendation, information retrieval, keyword extraction, meeting analysis, topic modeling.

### I. INTRODUCTION

Human beings are encompassed with the aid of an uncommon abundance of records, accessible as information, databases, or combined media belongings. access to this records is tailored by using the accessibility of appropriate web indexes, however notwithstanding while these are handy, clients often do not start a pursuit, in light of the truth that their modern-day movement does no longer allow them to do as such, or in mild of the fact that they're no longer mindful that relevant information is offered. We acquire on this paper the point of view of within the nick of time restoration, which replies this inadequacy via all of sudden suggesting data that are identified with clients' present sporting events. at the point when these physical activities are commonly conversational, for occurrence while clients take part in a meeting, their data needs can be validated as understood inquiries that are constructed out of sight from the professed phrases, obtained through continuous programmed discourse acknowledgment (ASR). those positive questions are applied to get better and endorse reports from the web or a neighborhood storehouse, which customers can determine to analyze in extra element in the event that they find out them exciting.

The center of this paper is on figuring verifiable questions to a without a second to spare recovery framework for usage in assembly rooms. Conversely to unequivocal talked inquiries that can be made in commercial enterprise internet crawlers, our within the nick of time healing framework should develop certain questions from conversational facts, which includes a much bigger number of words than a question. for example, within the instance examined in



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijirccce.com](http://www.ijirccce.com)

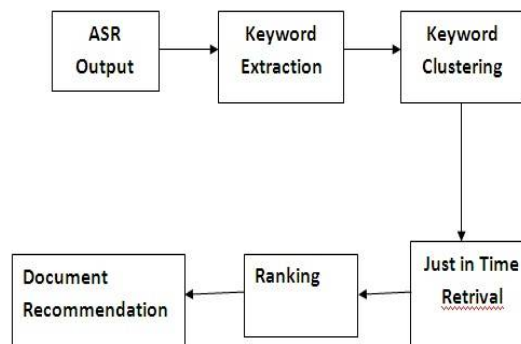
Vol. 5, Issue 6, June 2017

section V-B beneath, in which four people installation collectively a rundown of factors to help them get by using in the mountains, a quick piece of one hundred twenty seconds contains round 250 words, regarding a blended bag of regions, for example, 'chocolate', 'gun', or 'lighter'. What would possibly then be the maximum supportive 3–five Wikipedia pages to prescribe, and the way may a framework consciousness them?

Given the ability kind of subject matters, bolstered via ability ASR slips or discourse disfluencies, (for instance, "rush" in this illustration), our objective is to keep up unique speculations approximately clients' facts wishes, and to provide a touch example of proposals in view of the absolute confidence ones. in this way, we factor at keeping apart a pertinent and various arrangement of catchphrases, organization them into theme specific questions placed by means of importance, and gift clients an example of effects from those questions. The factor based bunching abatements the possibilities of inclusive of ASR errors into the questions,

## II. PROPOSED SYSTEM

The proposed system calculate split decisive words from the yield of an ASR framework (or a guide transcript for testing), which makes utilization of subject matter demonstrating techniques and of a sub modular prize capability which supports differing characteristics in the catchphrase set, to coordinate the ability diverse characteristics of points and lessen ASR commotion. At that factor, we advise a technique to infer special topically remote questions from this magic phrase set, with a selected end goal to extend the pictures of creating no much less than one crucial concept when using these inquiries to pursuit over the English Wikipedia. In this if any query has not found document relevant to it, then we apply that single query to search engine dataset to retrieve the document related to that.



In our system input is ASR output mean audio file. ASR output convert into text by using keyword extraction technique. Clustering perform on keyword using clustering algorithm. Applying ranking algorithm to rank the text file and document will be recommend to the user.

## III. PROPOSED ALGORITHM

### 1. Diverse Keyword Extraction

The benefit of numerous keyword extractions is that the insurance of the principle subjects of the communiqué fragment is maximized. Moreover, a good way to cover greater subjects, the proposed set of rules will pick out a smaller quantity of key phrases from every subject matter. That is acceptable for two reasons. This can result in extra numerous implicit queries, for this reason increasing the kind of retrieved files. And, if words which are in truth ASR noise can create a primary subject matter in the fragment, then the set of rules will pick out a smaller quantity of those noisy keywords compared to algorithms which forget about diversity.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 5, Issue 6, June 2017

## 2. Keyword Clustering

Clusters of keywords are built by keywords for each main topic of the fragment. One cluster contains similar keywords related to one topic. In our system we used K-means clustering algorithm.

## 3. Ranking Algorithm

The ranking function should rank more specific results higher than less specific results. In our system we used page rank algorithm for rank the document.

## IV. RESULT ANALYSIS

TABLE II. NO OF KEYWORD VS KEYWORD EXTRACTION METHODS

No of Keywords	D(.75)	TS	WF	D(.5)
1	0.91	0.925	0.7	0.825
2	0.95	0.825	0.825	0.95
3	0.825	0.7	0.775	0.91
4	0.9	0.7	0.79	0.91

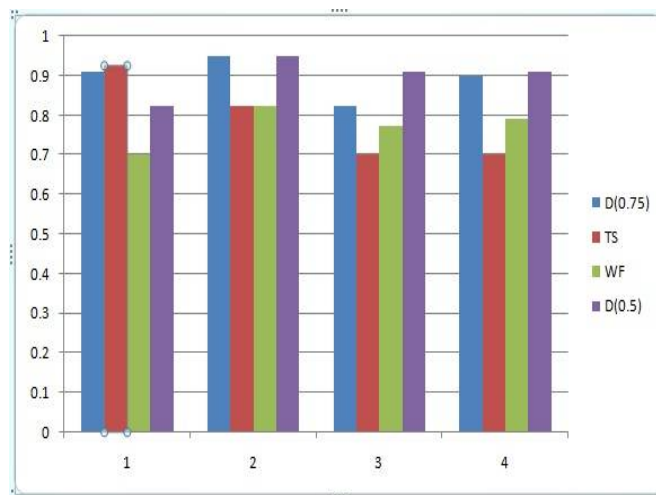


Fig. 2. Graph of No of keyword vs Keyword Extraction Methods

## V. CONCLUSION AND FUTURE WORK

We have taken into consideration a selected form of just-in-time retrieval systems supposed for conversational environments, in which they advise to customers documents which are relevant to their data wishes. We focused on modeling the customer's information wishes by means of deriving implicit queries from short conversation fragments. Those queries are based on sets of keywords extracted from the communication. We have proposed novel numerous keyword extraction approaches which covers the maximal range of essential topics in a fraction. Then, to reduce the noisy impact on queries of the mixture of subjects in keyword set, we proposed a clustering technique to divide the set of keywords into smaller topically-independent subsets constituting implicit queries.



ISSN(Online): 2320-9801  
ISSN(Print): 2320-9798

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 5, Issue 6, June 2017

## REFERENCES

- [1] M. Habibi and A. Popescu-Belis, "Enforcing topic diversity in a document recommender for conversations," in Proc. 25th Int. Conf. Comput.Linguist. (Coling), 2014, pp. 588–599.
- [2] H. P. Luhn, "A statistical approach to mechanized encoding and searching of literary information," IBM J. Res. Develop., vol. 1, no. 4, pp. 309–317, 1957.
- [3] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," Inf. Process. Manage. J., vol. 24, no. 5, pp. 513–523, 1988.
- [4] S. Ye, T.-S. Chua, M.-Y. Kan, and L. Qiu, "Document concept lattice for text understanding and summarization," Inf. Process. Manage., vol.43, no. 6, pp. 1643–1662, 2007.
- [5] A. Csomai and R. Mihalcea, "Linking educational materials to encyclopedic knowledge," in Proc. Conf. Artif. Intell.Educat.: BuildingTechnol. Rich Learn. Contexts That Work, 2007, pp. 557–559.
- [6] D. Harwath and T. J. Hazen, "Topic identification based extrinsic evaluation of summarization techniques applied to conversational speech," in Proc. Int. Conf. Acoust., Speech, Signal Process. (ICASSP), 2012, pp. 5073–5076.
- [7] A. Popescu-Belis, E. Boertjes, J. Kilgour, P. Poller, S. Castronovo, T. Wilson, A. Jaimes, and J. Carletta, "The AMIDA automatic content linking device: Just-in-time document retrieval in meetings," in Proc.5th Workshop Mach. Learn. Multimodal Interact. (MLMI), 2008, pp.272–283.
- [8] A. Popescu-Belis, M. Yazdani, A. Nanchen, and P. N. Garner, "Aspeech-based just-in-time retrieval system using semantic search," in Proc. Annu. Conf. North Amer. Chap. ACL (HLT-NAACL), 2011, pp.80–85.
- [9] P. E. Hart and J. Graham, "Query-free information retrieval," Int. J.Intell. Syst. Technol. Applicat., vol. 12, no. 5, pp. 32–37, 1997.
- [10] B. Rhodes and T. Starner, "Remembrance Agent: A continuously running automated information retrieval system," in Proc. 1st Int. Conf.Pract. Applicat.Intell. Agents Multi Agent Technol., London, U.K.,1996, pp. 487–495.