# Keyword Extraction and Clustering for Document Recommendation in Conversations

**Sucheta Thube[1], Bhavna Visave[2], Aishwarya Walekar[3], Rekha Hiware[4],
Prof. Mrs.Rashmi Bhattad[5]**

[1]*Department of Info.Tech,MMCOE ,Pune,Maharashtra,India*
[2] *Department of Info.Tech,MMCOE ,Pune,Maharashtra,India*
[3] *Department of Info.Tech,MMCOE ,Pune,Maharashtra,India*
[4] *Department of Info.Tech,MMCOE ,Pune,Maharashtra,India*
[5]*Department of Info.Tech,MMCOE ,Pune,Maharashtra,India*

*Abstract — This paper addresses the difficulty of keyword extraction from conversations, with the target of utilizing these watchwords to recover, for each short discussion piece, to a small degree variety of conceivably pertinent reports, which may be prescribed to members. In any case, even a brief piece contains a smorgasbord of words, that square measure conceivably known with many themes; additionally, utilizing a programmed discourse acknowledgment (ASR) framework presents slips among them. on these lines, it's exhausting to surmise properly the information desires of the discussion members. we tend to initial propose a calculation to get rid of decisive words from the yield of an ASR framework (or a manual transcript for testing), that makes utilization of theme demonstrating ways and of a sub standard prize capability that supports differing qualities within the magic word set, to coordinate the potential differing qualities of subjects and reduce ASR commotion. At that time, we tend to propose a method to infer varied locally isolated inquiries from this decisive word set, keeping in mind the top goal to amplify the chances of creating at any rate one pertinent proposal once utilizing these inquiries to get over English Wikipedia. The planned systems square measure assessed as so much as significance as for discussion items from the Fisher, AMI, and ELEA colloquial corpora, appraised by many human judges. The scores demonstrate that our proposition moves forward over past systems that contemplate simply word repeat or theme closeness, and speaks to a promising declare a report recommender framework to be used as an area of discussions.*

*Keywords- Document recommendation, information retrieval, keyword extraction, meeting analysis, topic modeling.*

## I. INTRODUCTION

Humans area unit encompassed by an uncommon abundance of knowledge, accessible as records, databases, or mixed media assets. Access to the current information is tailored by the accessibility of appropriate internet indexes, but still once these area unit accessible, purchasers often do not begin a research, in light-weight of the actual fact that their current action doesn't allow them to try and do per se, or in light-weight of the actual fact that they're not aware that applicable information is accessible. we have a tendency to receive during this paper the purpose of read of within the nick of your time recovery, that replies this inadequacy by suddenly suggesting archives that area unit known with clients' gift exercises. At the purpose once these exercises area unit primarily colloquial, for prevalence once purchasers participate during a meeting, their information desires may be in-contestable  as understood inquiries that area unit engineered out of sight from the professed words, non-heritable through continuous programmed discourse acknowledgment (ASR). These sure queries area unit utilized to recover and counsel reports from the online or a region store, that purchasers will arrange to investigate in additional detail if they discover them intriguing.

The center of this paper is on problem solving verifiable inquiries to a while not an instant to spare recovery framework for utilization in meeting rooms. Conversely to unequivocal talked inquiries that may be created in business internet crawlers, our within the nick of your time recovery framework should develop sure queries from colloquial info, that contains a way larger variety of words than a matter. for instance, within the illustration examined in Section V-B beneath, during which four people originated along a summation of things to assist them get by within the mountains, a brief piece of a hundred and twenty seconds contains around 250 words, about a motley of areas, for instance, 'chocolate', 'gun', or 'lighter'. What may then be the foremost auxiliary 3–5 Wikipedia pages to order, and the way may a framework focus them?

 Given the potential form of themes, strong by potential ASR slips or discourse disfluencies, (for example, "rush" during this illustration), our objective is to stay up completely different speculations concerning clients' information desires, and to gift a little example of proposals visible  of the little doubt ones. During this manner, we have a tendency to purpose at separating a pertinent and numerous arrangement of catchphrases, cluster them into theme specific queries positioned by significance, an gift purchasers an example of results from these queries. the purpose based mostly bunching abatements the chances of as well as ASR blunders into the queries, and also the various qualities of

essential words expands the chances that no but one among the prompt records answers a necessity for information, or will prompt a useful archive whereas taking when its hyperlinks. Case in purpose, whereas a method visible of word repeat would recover the related to Wikipedia pages: 'Light', 'Lighting', and 'Light My Fire' for the said piece, purchasers would lean toward a group, for instance, 'Lighter', "Fleece" and 'Chocolate'. relevance and various qualities may be licensed at 3 stages: at the purpose once removing the magic words; once building one or a couple of sure inquiries; or once re-positioning their outcomes.

## II.     LITERATURE REVIEW

### 1.     Enforcing topic diversity in a document recommender for conversations
**AUTHORS: M. Habibi and A. Popescu-Belis**

This paper addresses the matter of building telegraphic, various and relevant lists of documents, which might be suggested to the participants of a spoken language to meet their data wants while not distracting them. These lists are retrieved sporadically by submitting multiple implicit queries derived from the pronounced words. Every question is expounded to at least one of the topics known within the spoken language fragment preceding the advice, and is submitted to a probe engine over English people Wikipedia. We tend to propose during this paper AN algorithmic rule for various merging of those lists, employing a sub standard reward operate that rewards the topical similarity of documents to the spoken language words in addition as their diversity. We tend to judge the planned technique through crowdsourcing. The results show the prevalence of the varied merging technique over many others that not enforce the variety of topics.

### 2.     A statistical approach to mechanized encoding and searching of literary information
**AUTHORS:** H. P. Luhn

Written communication of concepts is dispensed on the idea of applied math chance therein a author chooses that level of subject specificity which combination of words that he feels can convey the foremost that means. Since this method varies among people and since similar concepts ar thus relayed at totally different levels of specificity and by suggests that of various words, the matter of literature looking by machines still presents major difficulties. A applied math approach to the current downside are printed and therefore the numerous steps of a system supported this approach are delineate. Steps embody the applied math analysis of a set of documents in an exceedingly field of interest, the institution of a group of "notions" and therefore the vocabulary by that they're expressed, the compilation of a thesaurus-type lexicon and index, the automated secret writing of documents by machine with the help of such a lexicon, the secret writing of topological notations (such as branched structures), the recording of the coded info, the institution of a looking pattern for locating pertinent info, and therefore the programming of applicable machines to hold out a probe.

### 3.     Document concept lattice for text understanding and summarization

**AUTHORS**: S. Ye, T.-S. Chua, M.-Y. Kan, and L. Qiu
We argue that the standard of a outline are often evaluated supported what percentage ideas within the original document(s) that may be preserved when account. Here, an idea refers to an abstract or concrete entity or its action typically expressed by numerous terms in text. Outline generation will so be thought of as an optimization downside of choosing a collection of sentences with stripped-down answer loss. During this paper, we have a tendency to propose a document conception lattice that indexes the hierarchy of native topics tied to a collection of frequent ideas and therefore the corresponding sentences containing these topics. The native topics can specify the promising sub-spaces associated with the chosen ideas and sentences. supported this lattice, the outline is Associate in Nursing optimized choice of a collection of distinct and salient native topics that cause greatest coverage of ideas with the given range of sentences. Our summarizer supported the conception lattice has incontestable competitive performance in Document Understanding Conference 2005 and 2006 evaluations moreover as innings tests. 2007 Elsevier Ltd. All rights reserved.

### 4.     Linking educational materials to encyclopedic knowledge

**AUTHORS**: A. Csomai and R. Mihalcea,

This paper describes a system that mechanically links study materials to encyclopedic data, and shows however the supply of such data inside easy reach of the learner will improve each the standard of the data non-inheritable and the time required to get such data.

**5. Remembrance Agent: A continuously running automated information retrieval system.**

**AUTHORS:** B. Rhodes and T. Starner

The Remembrance Agent (RA) could be a program that augments humanm emoryb y displaying an inventory of documents which could be relevant to the user's current context. not like most info retrieval systems, the RA runs unceasingly while not user intervention. Its unnoticeable interface permits a user to pursue or ignore the RA's suggestions as desired.

## III. PROPOSED SYSTEM

The planned strategies area unit evaluated in terms of connexion with relevance language fragments from the Fisher, AMI, and ELEA informal corpora, rated by many human judges. The scores show that our proposal improves over previous strategies that think about solely word frequency or topic similarity, and represents a promising resolution for a document recommender system to be utilized in conversations.
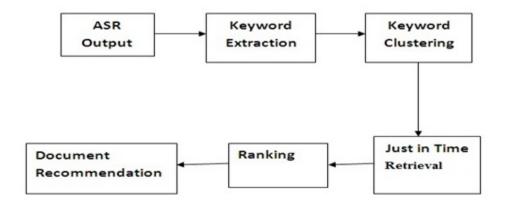


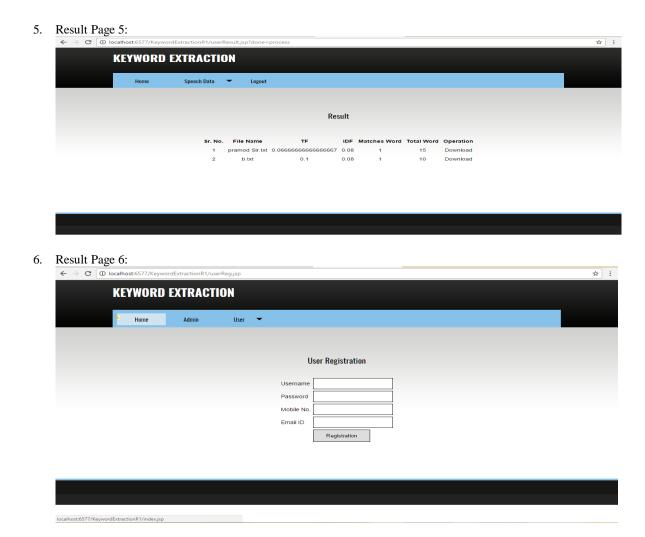Figure: Proposed system architecture

**Advantages of projected System:**

- Reduce Errors.

- Keyword extraction from conversations removes difficulty.

## IV. RESULT AND DISCUSSION

1. Result Page 1:

2.  Result Page 2:



3.  Result Page 3:



4.  Result Page 4:

5. Result Page 5:



6. Result Page 6:



## V.    CONCLUSION

We have thought-about a selected kind of while not an instant to spare recovery frameworks planned for colloquial things, during which they impose to client's archives that are necessary to their knowledge desires. we have a tendency to targeting displaying the client's knowledge desires by obtaining verifiable queries from short discussion items. These queries are in lightweight of sets of important words separated from the discussion. We've planned a unique completely different important word extraction strategy that covers the outside variety of significant themes during a piece. At that time, to minimize the boisterous impact on queries of the mix of themes during a decisive word set, we have a tendency to planned a grouping system to isolate the arrangement of catchphrases into littler topically-autonomous subsets constituting understood inquiries.

We compared the various keyword extraction technique with existing strategies, supported word frequency or topical similarity, in terms of the representativeness of the keywords and also the connectedness of retrieved documents. These were judged by human raters recruited via the Amazon Mechanical Turk crowd sourcing platform. The experiments showed that the various keyword extraction technique provides on the average the foremost representative keyword sets, with the best -NDCG worth, and leading–through multiple topically-separated implicit queries–to the foremost relevant lists of counseled documents. Therefore, imposing each connectedness and variety brings a good improvement to keyword extraction and document retrieval. The keyword extraction technique might be improved by considering n-grams of words additionally to individual words solely, however this needs some adaptation of the complete process chain.

**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.: Building Technol. 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, "A speech-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.

[11] Analytics: An Intelligent Approach in Clinical Trail Management Ankit Lodha* Analytics Operations Lead, Amgen, Thousand Oaks, California, USA.

[12] Agile: Open Innovation to Revolutionize Pharmaceutical Strategy Ankit Lodha University of Redlands, 333 N Glenoaks Blvd #630, Burbank, CA 91502.

[13] Clinical Analytics – Transforming Clinical Development through Big Data Ankit Lodha University of Redlands, 333 N Glenoaks Blvd #630, Burbank, CA 91502.