# Generating Text Description from Spoken Content using LDA

**Sayali S. Chavan[1], Ramesh M. Kagalkar[2]**

M.E Student, Computer Engg., Dr. DY Patil School of Engg And Technology, Pune, India[1]

Research Scholar and Asst. Professor, Computer Engg. , Dr. DY Patil School of Engg And Technology, Pune, India[2]

**Abstract**: Spoken content retrieval refers to directly indexing and retrieving spoken content based on the audio rather than text descriptions. This potentially eliminates the requirement of producing text descriptions for multimedia content for indexing and retrieval purposes, and is able to precisely locate the exact time the desired information appears in the multimedia. Spoken content retrieval has been very successfully achieved with the basic approach of cascading automatic speech recognition with text information retrieval: after the spoken content is transcribed into text or lattice format, a text retrieval engine searches over the ASR output to find desired information. Latent Dirichlet allocation algorithm is used for clustering. It is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar. Depends on automatic speech recognition output the algorithm will perform clustering on that content and return the expected result. After converting audio to text we apply LDA algorithm for clustering the input data.

**Keywords**: Spoken content retrieval, spoken term detection, query by example, semantic retrieval, joint optimization, pseudo-relevance feedback, graph-based random walk, unsupervised acoustic pattern discovery, query expansion, interactive retrieval, summarization, key term extraction.

## I. INTRODUCTION

Today the Internet has become an everyday part of human life. Internet content is indexed, retrieved, searched, and browsed primarily based on text, and the success of these capabilities has not only changed our lives, but generated a very successful global industry in Internet-based content and services. Although multimedia Internet content is growing rapidly, with shared videos, social media, broadcasts, etc., as of today, it still tends to be processed primarily based on the textual descriptions of the content offered by the multimedia providers. As automatic speech recognition (ASR) technologies continue to advance, it is reasonable to believe that speech and text offerings will eventually be symmetric, since they are alternative representations of human language, in the spoken and written form, respectively, and the transformation between the two should be direct and straightforward. With this perspective, spoken content retrieval, or indexing and retrieving multimedia content from its spoken part, is an important key to easier browsing and retrieving of multimedia content in the future. In cases where the essence of the multimedia content is captured by its audio, especially for broadcast programs, lectures, meetings, etc., indexing and retrieving the content based on the spoken part not only eliminates the extra requirements of producing the text description for indexing purposes, but can precisely locate the exact time when the desired information appears in the multimedia. The basic scenario for spoken content retrieval is therefore the following: when the user enters a *query*, which can be either in textual or spoken form, the system is expected to search over the spoken content and return relevant hits, possibly including the corresponding multimedia (e.g., video). In recent years, spoken content retrieval has achieved

significant advances by primarily cascading ASR output with text information retrieval techniques [6-8]. With this approach, the spoken content is first converted into word sequences or lattices via ASR. In order to cope with ASR errors, lattices have been used to represent the spoken content instead of a single word sequence [9]–[13], and sub word-based techniques have been used to some extent to address the out-of-vocabulary (OOV) problem [14–18]. For a subsequent user query (represented by lattices if spoken [8]), the text retrieval engine searches over the ASR output, and returns the relevant spoken content. After converting audio to text we apply LDA algorithm for clustering the input data. It is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar. For example, if observations are words collected into documents, it posits that each document is a mixture of a small number of topics and that each word's creation is attributable to one of the document's topics. LDA is an example of a topic model and was first presented as a graphical model for topic discovery [19-20].

## II. RELATED WORK

In this area, initially specify the notations utilized in this paper, examine some safe primitives utilized in our secure de duplication. Lin-shan Lee [1] that Spoken content retrieval refers to directly indexing and retrieving spoken content based on the audio rather than text descriptions. This potentially eliminates the requirement of producing text descriptions for multimedia content for indexing and retrieval purposes, and is able to precisely locate the exact time the desired information appears in the multimedia. Spoken content retrieval has been very successfully achieved with the basic approach of cascading automatic speech recognition (ASR) with text information retrieval:

after the spoken content is transcribed into text or lattice format, a text retrieval engine searches over the ASR output to find desired information. This framework works well when the ASR accuracy is relatively high, but becomes less adequate when more challenging real-world scenarios are considered, since retrieval performance depends heavily on ASR accuracy. The basic Framework of cascading ASR with text retrieval in order to have retrieval performances that is less dependent on ASR accuracy. The emergence of another approach to spoken content retrieval: to go beyond the basic framework of cascading ASR with text retrieval in order to have retrieval performances that are less dependent on ASR accuracy.

C. Chelba, T. Hazen, and M. Saraclar in [2] has proposed that Ever-increasing computing power and connectivity bandwidth, together with falling storage costs, are resulting in an overwhelming amount of data of various types being produced, exchanged, and stored. Consequently, information search and retrieval has emerged as a key application area. Text-based search is the most active area, with applications that range from Web and local network search to searching for personal information residing on one's own hard-drive. Speech search has received less attention perhaps because large collections of spoken material have previously not been available. However, with cheaper storage and increased broadband access, there has been a subsequent increase in the availability of online spoken audio content such as news broadcasts, podcasts, and academic lectures. A variety of personal and commercial uses also exist. As data availability increases, the lack of adequate technology for processing spoken documents becomes the limiting factor to large-scale access to spoken content. In this article, we strive to discuss the technical issues involved in the development of information retrieval systems for spoken audio documents, concentrating on the issue of handling the error-ful or incomplete output provided by ASR systems. We focus on the usage case where a user enters search terms into a search engine and is returned a collection of spoken document hits. This makes automatic approaches for indexing and searching spoken document collections very desirable. An ideal system would simply concatenate an automatic speech recognition (ASR) system with a standard text indexing and retrieval system. The speech recognition systems are not yet robust enough to produce high quality transcriptions for unconstrained speech audio in uncontrolled recording environments.

M. Larson and G. J. F. Jones put forth the concept of Speech media [3], that is, digital audio and video containing spoken content has blossomed in recent years. Large collections are accruing on the Internet as well as in private and enterprise settings. This growth has motivated extensive research on techniques and technologies that facilitate reliable indexing and retrieval. Spoken content retrieval (SCR) requires the combination of audio and speech processing technologies with methods from information retrieval (IR). SCR research initially investigated planned speech structured in document-like units, but has subsequently shifted focus to more informal spoken content produced spontaneously, outside of the studio and in conversational settings. This survey provides an overview of the field of SCR encompassing component technologies, the relationship of SCR to text IR and automatic speech recognition and user interaction issues. It is aimed at researchers with backgrounds in speech technology or IR who are seeking deeper insight on how these fields are integrated to support research and development, thus addressing the core challenges of SCR. It research initially investigated planned speech structured in document-like units, but has subsequently shifted focus to more informal spoken content produced spontaneously, outside of the studio and in conversational settings. Spoken content retrieval (SCR) requires the combination of audio and speech processing technologies with methods from information retrieval.

L.-s. Lee and B. Chen explained the concept in [4] that the retrieval can be performed based on the full content, the summaries/titles/topic labels, or both. The input spoken document length is too short. Maximum size of file can be uploaded. Spoken documents (or associated multimedia content) are in fact better understood and reorganized in a way that retrieval/browsing can be performed easily. For example, they are now in the form of short paragraphs, properly organized in some hierarchical visual presentation with titles/summaries/topic labels as references for retrieval and browsing. The retrieval can be performed based on the full content, the summaries/titles/topic labels, or both. In this article, this is referred to as spoken document understanding and organization for efficient retrieval/browsing applications. The purpose of this article is to present a concise, comprehensive, and integrated overview of related areas in a unified context of spoken document understanding and organization for efficient retrieval/browsing applications. In addition, we present an initial prototype system we developed at National Taiwan University as a new example of integrating the various technologies and functionalities.

M. Saraclar made a study on the recent work on spoken document retrieval in [5] has suggested that it is adequate to take the single best output of ASR, and perform text retrieval on this output. This is reasonable enough for the task of retrieving broadcast news stories, where word error rates are relatively low, and the stories are long enough to contain much redundancy. But it is patently not reasonable if one's task is to retrieve a short snippet of speech in a domain where WER's can be as high as 50%; such would be the situation with teleconference speech, where one's task is to find if and when a participant uttered a certain phrase. In this paper we propose an indexing procedure for spoken utterance retrieval that works on lattices rather than just single-best text. We demonstrate that this procedure can improve F scores by over five points compared to single best retrieval on tasks with poor WER and low redundancy. The representation is flexible so that we can represent both

word lattices, as well as phone lattices, the latter being important for improving performance when searching for phrases containing OOV words. Recent work on spoken document retrieval has suggested that it is adequate to take the single best output of ASR, and perform text retrieval on this output. It is patently not reasonable if one's task is to retrieve a short snippet of speech in a domain where WER's can be as high. An easy way to comply with the conference paper formatting requirements is to use this document as a template and simply type your text into it.

## III. PROPOSED APPROACH FRAMEWORK AND DESIGN

We have proposed a system in which the output retrieved from the ASR is taken as a input by the LDA algorithm and applying the algorithm the output is derived in the text format. The output derived is the analysis output extracted from the audio file [21-25]. The architecture of the proposed system is shown in the diagram below.

### A. Architecture

Retrieval-oriented Acoustic Modeling under Relevance Feedback Scenario: Relevance feedback well used in text retrieval is useful to integrate the ASR and retrieval modules as a whole and optimize the overall retrieval performance, rather than considering them as two cascaded independent components. When a query is entered by the user, the system offers a ranked list of retrieved objects to the user. Because the scores used for ranking the objects depend on the acoustic models, the objects below item 3 not yet viewed by the user can thus be re-ranked[26,27,28,29]. In this way, the acoustic models can be "adapted locally" considering the specific query and the corresponding feedback entered by the individual user, resulting in "query-specific" acoustic models, to be used for the unlimited number of acoustic conditions for the spoken content. An interactive retrieval process incorporating user actions may produce better retrieval results and user experiences. The LDA algorithm is used for clustering. Depends on ASR output the LDA will perform clustering on that content and return the expected result. After converting audio to text we apply LDA algorithm for clustering the input data [30,31].

### B. Proposed Work

#### Retrieval-Oriented Language Modeling

In keyword spotting, it was found that boosting the probabilities of n-grams including query terms by repeating the sentences including the query terms in the language model. Robust Automatic Transcription of Speech (RATS) program and the NIST OpenKWS13 Evaluation. Similar concept was also used in neural network based language models (NNLM), who's input is a history word sequence represented by a feature vector, and the output is the probability distribution over the words [32].

#### Retrieval-Oriented Decoding.

It has been proposed that the search with OOV queries can be achieved in two steps. In this framework, each utterance has a word-based and a sub word-based lattices. When an OOV query is entered, in the first step, a set of utterances which possibly contain the OOV query is obtained by searching over the sub word-based lattices. *Retrieval-Oriented Confusion Models* Some effort has been made to model the occurrence of the recognition errors in a systematic way, referred to as confusion models here, and to try to optimize such models to have better retrieval performance [33]. There can be at least three ways to achieve this goal: *Query transformation* (to transform the word or sub word sequence of each query into the sequences that the query tends to be mis-recognized to, and the new set of sequences are used to retrieve the lattices), *Spoken Content transformation* (to transform the recognition output for the spoken content instead of the query), and *Fuzzy match* (defining a distance between different word or sub word sequences, and the lattices containing word or sub word sequences sufficiently close to the query being retrieved). An interactive retrieval process incorporating user actions may produce better retrieval results and user experiences. The LDA algorithm is used for clustering. Depends on ASR output the LDA will perform clustering on that content and return the expected result. After converting audio to text we apply LDA algorithm for clustering. [34-35].
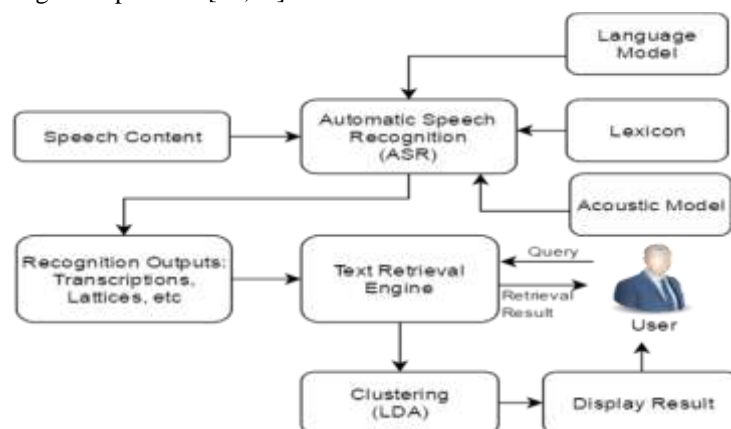


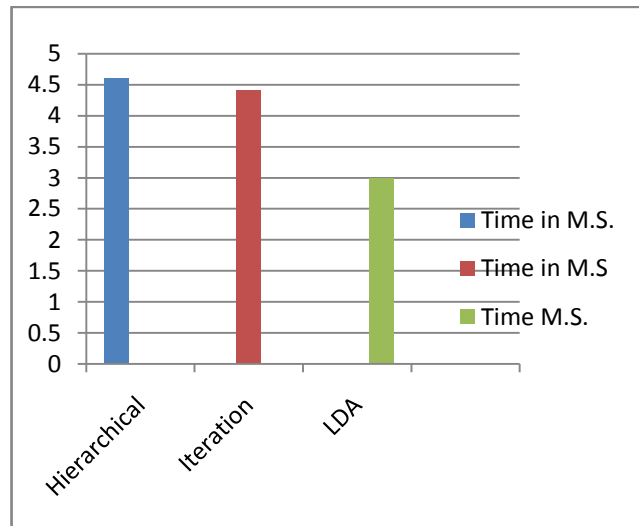Fig. 1 Architecture of proposed system

## IV. METHODOLOGY

An overview of the system is presented in Fig.1. For an input audio sample, the following steps are followed:

Input:     A – Audio

Steps:

1. User will give voice (Audio) to system. Apply signal processing on the audio file for noise elimination.
2. Apply language model, lattices and perform the lexicon analysis. These are implemented by ASR..
3. Recognize the output provided by the ASR. The text Retrieval Engine retrieves the text.
4. Apply LDA algorithm. The algorithm runs in following separate steps
   1) Calculate TF and IDF:
   2) TF =Term Frequency
   3) IDF= Inverse Document Frequency
   4) $TF(t)$ = (Number of times term t appears in a document) / (Total number of terms in the document)
   5) $IDF(t)$ = $\log_e$(Total number of documents / Number of documents with term t in it).
   6) Value = TF * IDF
5. The expected result will be displayed to end user as per user input query.

derived output is "Java simple dot net". This is improper statement. Hence future scope is to get the proper output. In S2 the input is related to the market conversation. And the derived output is quiet similar to the expected output but is grammatically not corrected.  In S3 the conversation is taken from the garment shop where the lady is buying the cloths and she is comparing the cloths and the output is driven by making the comparison between the types of the cloths. Finally in S4 conversation is between the Owner and the Employee where the owner is asking the employee to do his job quickly. And the employee is ready to submit the job by today.



| Sr.No | Audio Samples | Audio input | SIZE (KB) | Actual output | Expected Output |
|-------|--------------|-------------|-----------|---------------|-----------------|
| 1. | S1 | Student 1: My project topic is on java. Student 2: My project topic is on dot net language. Student 1: Is java simple? Student 2: Yes it is simple. Dot net is tough language. | 10 | Java simple dot net | Java is simple than dot net. |
| 2. | S2 | Man: How much is the cost of potatoes. Women: Sir, it is 40 Rs. Per kg. Man: It seems costlier than yesterday. Women: It was 45 Rs. yesterday. | 10 | Potatoes costlier | The potatoes are costlier today. |
| 3. | S3 | Girl: What type of garment is this? Man: This is cotton cloth. Girl: I want a synthetic cloth. The cotton cloths are costlier. Man: Here it is synthetic. Girl: It is nice than the cotton. Man: Bye this madam as is it beautiful. | 12 | Cotton costlier synthetic beautiful | Synthetic cloths are beautiful and cotton cloths are costlier. |
| 4. | S4 | Owner: Did u do your job? Employee: No sir, it is incomplete. Owner: When it will be completed? Employee: It will be completed today. Owner: Do it fast as our customers are waiting for reply. Employee: Yes sir it will be done today. | 14 | Incomplete job complete today | Employee will complete his job today. |

Fig. 2  Result table

## V. Graph

In this section, time require to process audio depend on size of audio is shown. Time requires for Hierarchical, Iteration and LDA audio extraction is given in the graph. As compare to Hierarchical and Iteration the LDA algorithm takes less time. Hence LDA algorithm is better among hierarchical and iteration algorithm.

## VI. Conclusion

Many advanced application tasks of spoken language processing were solved by cascading a set of modules in early stages of developments. Take the spoken dialogue system as an example, which was actually built in early years by cascading ASR, natural language understanding, dialogue management, natural language generation and TTS. Today the spoken dialogue is already a full-fledged independent area far beyond the above cascading framework. Good examples include the dialogue managers based on Partially Observable Markov Decision Process taking the uncertainty of ASR and spoken language understanding into considerations, and learning the policy of dialogue manager and natural language generator jointly. These novel techniques beyond the cascading framework have turned the pages of the research and development of spoken dialogues. Another example is speech translation, in which jointly optimizing the ASR module and its downstream text processing module is also considered as a major trend. We believe similar developments may be experienced in spoken content retrieval in the future. Cascading ASR with text retrieval has been very successful in this area, but inevitably becomes less adequate for more challenging real-world tasks.

## References

[1] G. Tur and R. DeMori, Spoken Language Understanding: Systems for Extracting Semantic Information from Speech. New York, NY, USA: Wiley, 2011, ch. 15, pp. 417–446.

[2] C. Chelba, T. Hazen, and M. Saraclar, "Retrieval and browsing of spoken content," IEEE Signal Process. Mag., vol. 25, no. 3, pp. 39–49, May 2008.

[3] M. Larson and G. J. F. Jones, "Spoken content retrieval: A survey of techniques and technologies," Found. Trends Inf. Retr., vol. 5, pp. 235–422, 2012.

[4] L.-s. Lee and Y.-C. Pan, "Voice-based information retrieval–how far are we from the text-based information retrieval?," in Proc. ASRU, 2009.

[5] L.-s. Lee and B. Chen, "Spoken document understanding and organization,"IEEE Signal Process. Mag., vol. 22, no. 5, pp. 42–60, Sep. 2005.

[6] K. Koumpis and S. Renals, "Content-based access to spoken audio," IEEE Signal Process. Mag., vol. 22, no. 5, pp. 61–69, Sep. 2005.

[7] A. Mandal, K. Prasanna Kumar, and P. Mitra, "Recent developments in spoken term detection: A survey," Int. J. Speech Technol., pp. 1–16, 2013.

[8] T. K. Chia, K. C. Sim, H. Li, and H. T. Ng, "A lattice-based approach to query-by-example spoken document retrieval," in Proc. SIGIR, 2008.

[9] M. Saraclar, "Lattice-based search for spoken utterance retrieval," in Proc. HLT-NAACL, 2004, pp. 129–136.

[10] C. Chelba and A. Acero, "Position specific posterior lattices for indexing speech," in Proc. 43rd Annu. Meeting Assoc. Comput. Linguist., 2005, pp. 443–450.

[11] D. Vergyri, I. Shafran, A. Stolcke, R. R. Gadde, M. Akbacak, B. Roark, and W. Wang, "The SRI/OGI 2006spoken term detection system," in Proc. Interspeech, 2007.

[12] J. Mamou, B. Ramabhadran, and O. Siohan, "Vocabulary independent spoken term detection," in Proc. SIGIR, 2007.

[13] D. R. H. Miller, M. Kleber, C. lin Kao, O. Kimball, T. Colthurst, S. A. Lowe, R. M. Schwartz, and H. Gish, "Rapid and accurate spoken term detection," in Proc. Interspeech, 2007.

[14] K. Ng, "Subword-based approaches for spoken document retrieval," Ph.D. dissertation, Mass. Inst. of Technol., Cambridge, MA, USA, 2000.

[15] J. S. Garofolo, C. G. P. Auzanne, and E. M. Voorhees, The TREC Spoken Document Retrieval Track: A Success Story, 2000.

[16] [Online]. Available: http://speechfind.utdallas.edu/ [17] J. Ogata and M. Goto, "Podcastle: Collaborative training of acoustic models on the basis of wisdom of crowds for podcast transcription," in Proc. Interspeech, 2009.

[17] C. Alberti, M. Bacchiani, A. Bezman, C. Chelba, A. Drofa, H. Liao, P. Moreno, T. Power, A. Sahuguet, M. Shugrina, and O. Siohan, "An audio indexing system for election video material," in Proc. ICASSP, 2009, pp. 4873–4876.

[18] J. Glass, T. J. Hazen, S. Cyphers, I. Malioutov, D. Huynh, and R. Barzilay, "Recent progress in the MIT spoken lecture processing project," in Proc. Interspeech, 2007.

[19] S.-Y. Kong, M.-R. Wu, C.-K. Lin, Y.-S. Fu, and L.-s. Lee, "Learning on demand–course lecture distillation by information extraction and semantic structuring for spoken documents," in Proc. ICASSP, 2009, pp. 4709–4712.

[20] Shivaji J. Chaudhari and Ramesh M. Kagalkar, "A Review of Automatic Speaker Age Classification, Recognition and Identifying Speaker Emotion Using Voice Signal", International Journal of Science and Research (IJSR 2014), ISSN(Online): 2319-7064,Volume 3 Issue 11, November 2014.

[21] Kaveri Kamble and Ramesh Kagalkar," A Review: Translation of Text to Speech Conversion for Hindi Language", International Journal of Science and Research (IJSR) ISSN (Online): 2319-7064, Vol. 3 Issue 11, November 2014.

[22] Ajay R. Kadam and Ramesh M. Kagalkar," Predictive Sound Recognition System", International Journal of Advance Research in Computer Science and Management Studies, Volume 2, Issue 11, November 2014.

[23] Shivaji J. Chaudhari and Ramesh M. Kagalkar," Automatic Speaker Age Estimation and Gender Dependent Emotion Recognition", International Journal of Computer Applications( IJCA) (0975 - 8887),Volume 117 No. 17, May 2015.

[24] Shivaji J. Chaudhari and Ramesh M. Kagalkar , "A Methodology for Efficient Gender Dependent Speaker Age and Emotion Identification System", International Journal of Advanced Research in Computer and Communication Engineering(IJARCCE) ISSN 2319- 5940,Volume 4, Issue 7, July 2015.

[25] Kaveri Kamble and Ramesh Kagalkar, "Audio Visual Speech Synthesis and Speech Recognition for Hindi Language" , International Journal of Computer Science and Information Technologies(IJCSIT) ISSN (Online): 0975-9646, Vol. 6 Issue 2, April 2015.

[26] Kaveri Kamble and Ramesh Kagalkar , "A Novel Approach for Hindi Text Description to Speech and Expressive Speech Synthesis", International Journal of Applied Information Systems (IJAIS) ISSN 2249-0868, Vol. 8 Issue 7, May 2015.

[27] Ajay R. Kadam and Ramesh M. Kagalkar," Audio Scenarios Detection Technique", International Journal of Computer Applications (IJCA), Volume 120- Number 16, June 2015.

[28] Ajay R. Kadam and Ramesh M. Kagalkar, " A Review Paper on Predictive Sound Recognition System", CiiT International Journal of Software Engineering and Technology, June Issue 2015(Online), June 2015.

[29] Amitkumar Shinde and Ramesh M. Kagalkar, "Sign Language Recognition for Deaf Sign User", International Journal for Research in Applied Science & Engineering Technology

(IJRASET) ©IJRASET, Volume 2, Issue XII, December, ISSN: 2321-9653, 2014.

[30]   Amit kumar and Ramesh Kagalkar, "Methodology for Translation of Sign Language into Textual Version in Marathi", CiiT, International Journal of Digital Image Processing, Volume 07, No.08, Aug 2015.

[31]   Amitkumar Shinde and Ramesh M. Kagalkar," Advanced Marathi Sign Language Recognition using Computer Vision", International Journal of Computer Applications, (ISSN:0975 – 8887) , Volume 118,  No. 13, May 2015.

[32]   Rashmi. B. Hiremath and  Ramesh. M. Kagalkar, "Methodology for Sign Language Video Interpretation in Hindi Text Language ", International Journal of Innovative Research in Computer and Communication Engineering, Volume. 4, Issue 5, May 2016.

[33]   Rashmi. B. Hiremath and Ramesh. M. Kagalkar, "Sign Language Video Processing for Text Detection in Hindi Language", International Journal of Recent Contributions from Engineering, Science and IT,  Volume  4, No 3, 2016.

[34]   Rashmi. B. Hiremath and Ramesh. M. Kagalkar" A Methodology for Sign Language Video Analysis and Translation into Text in Hindi Language", CiiT International Journal of Fuzzy Systems, Volume  8, No 5, 2016.

## BIOGRAPHY

Authors have the option to publish a biography together with the paper, with the academic qualification, past and present positions, research interests, awards, etc. This increases the profile of the authors and is well received by international reader.

**Ramesh M. Kagalkar** He was born on Jun 1st, 1979 in Karnataka, India and presently working as an Assistant. Professor, Department of Computer Engineering, Dr. D Y Patil School of Engineering and Technology, Charoli, B.K.Via Lohegaon, Pune, Maharashtra, India. He has 14.5 years of teaching experience at various institutions. He is a Research scholar in Visveswaraiah Technological University, Belgaum, He had obtained M.Tech (CSE) Degree in 2006 from VTU Belgaum and He received BE (CSE) Degree in 2001 from Gulbarga University, Gulbarga Karnataka, India. He is the author of two text book; 1.Advance Computer Architecture, 2. The Swift Practical Approach of Learning C-Programming in LAP LAMBERT Academic Publishing, Germany (Available in online) and  One of his research article A Novel Approach for Privacy Preserving has been consider as text in LAP LAMBERT Academic Publishing, Germany (Available in online). He is waiting for submission of two research articles for patent right. He has published more than 35 research papers in International Journals and presented few of there in international conferences. His main research interest includes Image processing, Gesture recognition, Speech processing, Voice to sign language and CBIR. Under his guidance Ten ME students awarded degree in SPPU, Pune, three students at the edge of completion their ME final dissertation reports and four students started new research work and they have publish their research papers on International Journals and International conference. He can be contacted by email rameshvtu10@gmail.com.

**Sayali S. Chavan** is M.E 2nd year student of Computer Engg. Department, Dr. D. Y. Patil School of Engg. And Technology, Lohegaon, Pune. And her research interest includes computer engineering & networking. E-mail-sayali.chavan002@gmail.com