

# **A Comprehensive Analysis of Bio-Inspired Speech indexing using Latent Semantic Analysis**

<sup>1\*</sup> **R. Eswin Pria Angel**

Assistant Professor, DoS in Computer Science,  
SBRR Mahajana First Grade College (A), PG Wing, KRS Road, Metagalli, Mysuru, Karnataka, India

<sup>2</sup> **R.C. Evangeline**

Assistant Professor, Dept. of ISE, Nitte Meenakshi Institute of Technology, Bangalore, Karnataka, India

<sup>3</sup> **Dr.P.Raviraj**

Professor & Head, Dept. of CSE,  
GSSS Institute of Engineering and Technology for Women, Affiliated to VTU-Belagavi, Mysuru, Karnataka, India.

## **Abstract:**

A comprehensive analysis of bio-inspired speech indexing using Latent Semantic Analysis (LSA) explores the application of LSA in extracting semantic information from speech data. By leveraging LSA's ability to capture latent semantic relationships between words and concepts, this approach enhances the effectiveness of indexing and retrieving meaningful content from large speech datasets. The integration of bio-inspired techniques, such as neural networks and self-organizing maps, further refines the indexing process, making it more efficient and adaptable to complex audio environments. The analysis examines how LSA can overcome challenges related to noise, data sparsity, and the dynamic nature of spoken language, ultimately improving the accuracy and scalability of speech indexing systems.

**Keywords:** Latent Semantic Analysis, cognitive science, robotics, and cybernetics

## **1. Introduction**

Modern enterprises face two key challenges: the exponential growth of information and the rising importance of extracting value from it. Handling terabytes of text, including emails, requires rapid access and meaningful analysis, even for small and medium-sized businesses. As research becomes increasingly interdisciplinary, solving complex problems demands crossing traditional boundaries, integrating diverse disciplines, and applying multiple methods. This approach fosters innovation by examining subjects from various perspectives. This volume exemplifies the essential interdisciplinary nature of 21st-century research.

## 2. Challenges in Speech Processing

Speech processing is a fascinating field with immense potential to benefit society and offers numerous research opportunities at a fundamental level. Natural Language Processing (NLP), though inherently complex, has attracted countless researchers due to its unique challenges. Linguistics continues to grapple with understanding how natural language is acquired, produced, and processed in multilingual contexts, requiring tools for interaction, translation, and human-computer interfaces. NLP extends beyond audio processing to include pre-processing, noise reduction, segmentation, classification, information retrieval, and meaning analysis. It also powers speech automation systems like Automatic Voice Response (AVR) and speaking robots. This work focuses on semantic analysis, addressing the challenge of identifying meaning, mapping context, and creating truly meaningful speech automation systems.

## 3. Audio or Speech Indexing

Indexing organizes data types like text, images, audio, video, and multimedia for efficient storage and retrieval based on input queries. It involves categorization, summarization, content-based matching, and retrieval. Audio indexing requires annotation tools to insert text and process audio based on predefined indices. Matching is text-based, involving word frequency analysis and statistical comparisons for retrieval. Challenges arise with large audio datasets, as raw data lacks utility for matching, speech-derived text omits prosodic features, and feature representation demands significant effort and precision.

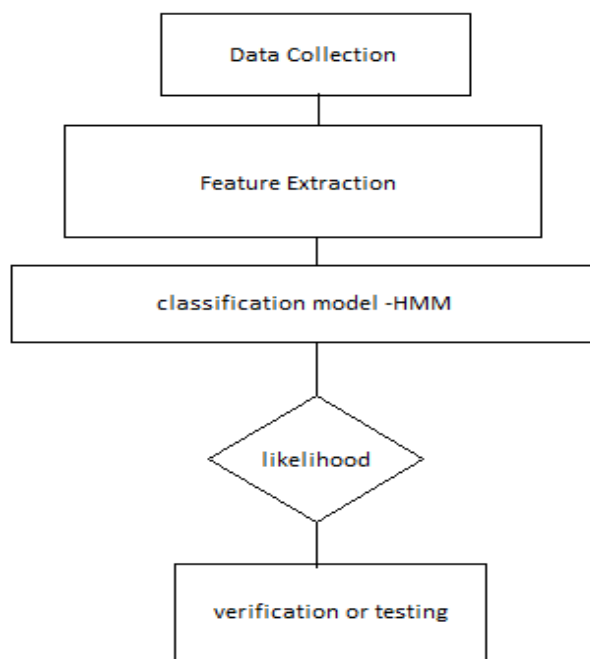


Fig1. Semantic concept extraction process

Figure 1 illustrates the process of audio clip processing, beginning with data collection and feature extraction. Key features include volume standard deviation, volume undulation, zero-

crossing rate (silence vs. non-silence), pitch contour, pitch variation, and frequency metrics like centroid, bandwidth, and sub-bands. These features support audio indexing tasks such as speaker recognition using prosodic elements and language identification through speech rate, syllable timing, pitch contour, and differential pitch. Additional metrics include syllable distance and duration for detailed analysis.

Audio indexing tasks

- ❖ Categorization of audio clips
- ❖ Categorization of musical instruments
- ❖ Categorization of speakers
- ❖ Categorization of languages
- ❖ Music / speech discrimination
- ❖ Spoken document retrieval

Task based features such as audio clip indexing using spectral centroid, musical classification such as loudness, pitch, bandwidth, harmonicity speech discrimination with formant frequencies or peaks and their time of occurrence can be taken for research works using simple mechanisms.

#### **4. Approaches to Semantic Processing**

Efforts to incorporate semantic information into speech or audio processing date back nearly half a century but with help of text which is annotated side by side. Over the years, designers have followed various approaches to integrating some degree of semantic processing into their information retrieval systems:

- ❖ Auxiliary Structures
- ❖ Local Co-Occurrence Statistics
- ❖ Latent Semantic Indexing

##### **4.1 Auxiliary Structures**

Controlled vocabularies, like dictionaries and thesauri, enhance queries by incorporating broader, narrower, and related terms, addressing issues of synonymy (similar meanings) and polysemy (multiple meanings). Tools like WordNet and domain-specific ontologies have expanded semantic constructs by representing concepts and their relationships. While controlled vocabularies improve efficiency in information retrieval for well-defined topics with standardized terminology, they struggle with unstructured data and the diverse, expansive needs of modern enterprises.

Some other drawbacks of using auxiliary structures:

- ❖ Establishing useful controlled vocabularies requires lots of human input and oversight.
- ❖ Language rapidly evolves, requiring the constant updating of controlled vocabularies.

- ❖ Controlled vocabularies can often represent the world view of their creators, introducing a potential source for conceptual mismatches.
- ❖ Controlled vocabularies capture a world view at a particular point in time. They can be difficult to modify as concepts change in a specific topic area.

## **4.2 Local Co-Occurrence Statistics**

Statistical co-occurrence, explored since the 1950s and widely used in the 1990s for synonym mining and word-to-word translation, involves counting how often term pairs appear together within a sliding window in a document. While simple, this method captures only a fraction of the semantic information in a text and relies on prior knowledge of the content, which is challenging for large, unstructured collections. Studies show that only about 25% of text information is local, limiting the effectiveness of co-occurrence-based approaches in most applications.

## **4.3 Latent Semantic Indexing**

Latent Semantic Indexing (LSI) is a statistical information retrieval method that retrieves text based on concepts rather than specific keywords. Developed at Bell Labs in the 1980s, LSI uses a term-document matrix, applies term weighting, and performs Singular Value Decomposition (SVD) to uncover patterns in term-concept relationships. By reducing the dimensions of the term space, LSI establishes associations between terms in similar contexts, enabling conceptually relevant query results even without shared words.

LSI excels at conceptual matching and has been shown to capture causal, goal-oriented, and taxonomic relationships. Experiments reveal surprising similarities between LSI and human text categorization. Its superiority in extracting semantic information has been validated through tasks like document categorization. Notably, LSI achieved the best-ever results in the Reuters 21578 benchmark, a global standard for automated document categorization, demonstrating its efficacy in processing and categorizing large text collections.

## **5. Semantic Analysis**

While syntactic structure has been explored in language modeling, semantic analysis has received less attention. This study introduces three novel language modeling techniques that incorporate semantic analysis for spoken dialog systems. Semantic analysis involves assigning semantic tags and constructing hierarchical groupings, similar to syntactic parsing, with shallow parsing (single-level hierarchy) and full parsing (multi-level hierarchy). Tags represent the meanings or concepts of words, with shallow parsers (chunkers or classers) capturing basic structures, while full parses handle deeper semantic hierarchies.

This chapter uses bracket notation to represent semantic parses, where words are paired with tags and grouped with labeled tokens (e.g., [LABEL and LABEL]). If tags are not used, only the label tokens are included. To enable semantic parsing, speech data must first be

converted into text using a Speech-to-Text system or Part-of-Speech (POS) tagging. Figures 2 and 3 provide examples of shallow and full semantic parses.

**5.1 Shallow semantic analysis using weighted finite state transducers**

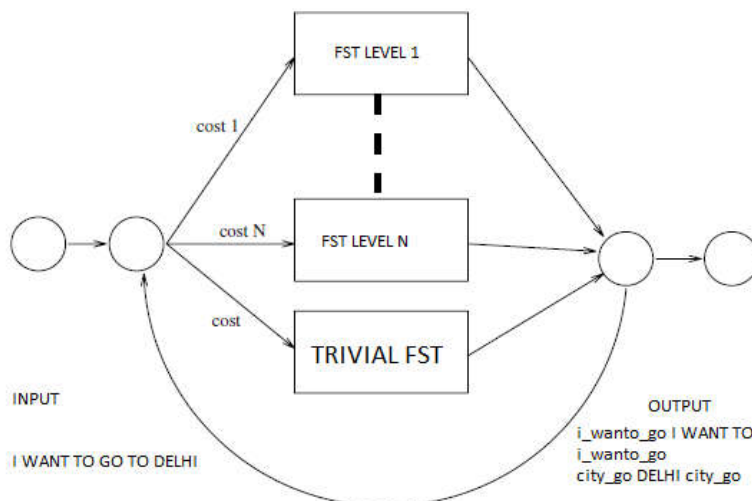


Fig2. Semantic concept extraction process

It is much easier to write phrasal concept grammars than grammars which accept whole sentences. For example, one can come up with a grammar for all expressions that refer to a date. Given a list of locations that we are interested in, we can write a CFG that accepts all of these location expressions. Similarly, many other concepts like these can be represented by phrasal CFGs.

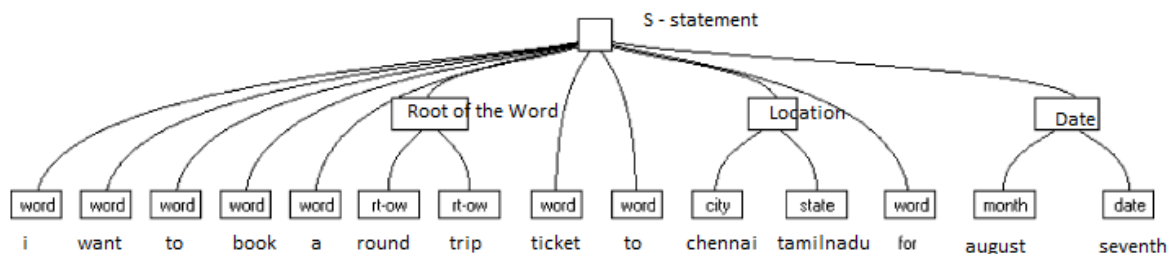


Fig3. Semantic classifier with output

A sentence can have multiple parses due to overlapping concept context-free grammars (CFGs); for example, "two twenty" might represent a time, apartment number, or flight number. To handle this, probabilities (or costs) are assigned to each phrasal finite state transducer (FST) to prefer certain parses. Costs are initialized uniformly, with a higher cost for trivial FSTs, favoring simpler parses. Training is performed using the Viterbi Expectation Maximization algorithm, with performance enhanced by rescoring parses using N-gram statistics for concept sequences. Sentences accepted by the weighted FST (WFST) can be represented by a stochastic context-free grammar (SCFG), which resembles a stochastic recursive transition network (SRTN) for handling recursive structures.

## **6. Latent Semantic Analysis**

Latent Semantic Analysis (LSA) is a method for extracting and representing the contextual meaning of words through statistical computations applied to large text corpora. It relies on the idea that the contexts in which words appear help determine their semantic similarity. LSA's adequacy is validated by its ability to replicate human performance on vocabulary tests, word sorting, and semantic priming tasks. By analyzing large language datasets, LSA represents words and passages as points in a high-dimensional semantic space using Singular Value Decomposition (SVD), a technique akin to factor analysis.

LSA provides a model for approximating human-like semantic judgments, offering insights into word-word, word-passage, and passage-passage relationships. Its key advantage lies in its ability to uncover deeper statistical relationships beyond simple co-occurrence, improving its predictive accuracy over traditional methods. Although LSA is based purely on textual analysis and lacks direct sensory or experiential input, it can approximate knowledge about the world as reflected in language.

LSA does not account for word order, syntax, or morphology, limiting its precision in some cases. It differs from other statistical approaches by using detailed word-context patterns across large contexts, like sentences and paragraphs, rather than relying solely on word pair co-occurrence. Additionally, LSA emphasizes reducing the dimensionality of data to improve its cognitive modeling accuracy, akin to postulated semantic features in psycholinguistics. Despite its limitations, LSA's ability to process text and model human cognition has proven useful in various research applications.

## **7. Semantics of LSA-based language models**

Landauer and Dumais (1997) show that Latent Semantic Analysis (LSA) can capture higher-order semantic similarities, meaning that if A and B are similar, and B and C are similar, A and C are also similar, even if they don't appear together. They also highlight LSA's inductive power, where 75% of its knowledge about words comes from documents that do not directly contain those words. Their experiments demonstrated that LSA's ability to identify synonyms improves as the training corpus size increases, even without direct synonym inclusion.

### **7.1 Lexical meaning and sentence meaning**

We distinguish between lexical meaning, sentence meaning, and utterance meaning. This discussion focuses on specific concepts from each category: "synonymy" for lexical meaning, and "predication" and "negation" for sentence meaning. Utterance meaning, which involves speech acts (e.g., questions, promises) and discourse relations like conversational implicatures, is not addressed here. Sentence meaning, central to logical semantics, deals with truth and related concepts such as negation, predication, and quantification. Synonymy involves interchangeable words with different meanings in specific contexts, predication uses first-order logic to analyze logical consequences, and negation helps derive word meanings as vectors.

## 8. Conclusion

Self-Organizing Maps (SOM) have been effectively used to organize large text archives by converting them into smoothed histograms of word categories. SOM can cluster documents based on document vectors, which are weighted averages of word vectors decoded from speech. The goal is to associate documents with index terms that capture their latent semantics, aiding in effective query processing. Audio indexing involves multiple stages: recording, classifying, decoding, indexing, and query processing.

SOM and LSA-based audio indexing are ideal for large datasets, such as broadcast news or flight data recorders. Text data derived from speech is converted into sparse matrices, and Singular Value Decomposition (SVD) reduces computational complexity. Random mapping can generate orthogonal vectors, allowing efficient approximation with lower complexity than SVD, even with large vocabularies.

Clustering, particularly using SOM, mitigates the effects of incorrect or missing decoded words by mapping documents based on their content. Unlike traditional methods like K-means, SOM adapts to the fine structure of data, offering smoother representation and better accuracy in dense areas. SOM's 2D grid visualization aids in understanding document clusters, with average perplexities ranging between 1.5 and 2.3. The best performance comes from LSA-based SVDSOM, which can retrieve a large amount of correlated text.

## References:

- [1]. Davies, K.: 1999, 'The IBM Conversational Telephony System for Financial Applications', in Proceedings of EuroSpeech'1999, Budapest, pp. 275–278
- [2]. Ward, Todd, 2000, How long until a high school student can build a language understanding system. In: ICSLP.
- [3]. Erdogan, Hakan, Sarikaya, Ruhi, Gao, Yuqing, Picheny, Michael, 2002. Semantic structured language models, In:ICSL
- [4]. Y. Deng and S. Khudanpur, 2003, Latent semantic information in maximum entropy language models for conversational speech recognition. In Proceedings of HLTNAACL, pages 56–63, Edmonton.
- [5]. J.R. Bellegarda. 2000a, Exploiting latent semantic information in statistical language modeling. Proceedings of the IEEE, 88(8):1279–1296.
- [6]. J.R. Bellegarda. 2000b, Large vocabulary speech recognition with multispans statistical language models. IEEE Transactions on Speech and Audio Processing, 8(1):76– 84.
- [7]. Jerome Bellegarda, 2004, Latent Semantic Language Modeling for Speech Recognition, Mathematical Foundations of Speech and Language Processing.
- [8]. D.Jurafsky and Martin, 2000, Speech and Language Processing, Prentice Hall, pp.223-231.

[9]. Chen, Stanley F. and Joshua Goodman,1998, An Empirical Study of Smoothing Techniques for Language Modeling. Harvard Computer Science Group Technical Report TR-10-98.

[10]. Grice (1981). "Presupposition and Conversational Implicature", in P. Cole (ed.), Radical Pragmatics, Academic Press, New York, pp. 183–198.