

RCILTS-Kannada Indian Institute of Science, Bangalore

Introduction

The Resource Centre for Indian Language Technology Solutions-Kannada at Indian Institute of Science chose a very broad spectrum of activities, and undertook several activities related languages other than Kannada. Participation by a large number of faculty of the Institute enabled the Resource Centre to create knowledge bases, web sites, OCR for Tamil and Kannada, speech synthesis for Tamil and Kannada, resource information on European history in Hindi for high school children, a bilingual interactive German-Hindi course, CDs in support of learning Kannada at primary and high school level, a variety of tools for language processing, word net in Kannada, and investigate into language identification. It was also recognized that language technology segment of the IT industry remained very small in India for a variety of reasons. Two important approaches were taken towards improving this. One is to get some jobs done by the industry instead of getting them done through project assistants, and the other one is to make available many of the basic tools developed in the public domain. The web site KANNUDI should meet the long felt needs of Kannadigas to know about their language and state. A large number of scholars in various disciplines actively contributed in making this web site very rich. An organization like Kannada Sahitya Parishat felt enthused enough to permit this resource centre to create a web site on their activities and have it co-located with KANNUDI. In an attempt to make a difference to an important segment of population, namely primary and high school students, Bodhana Bharati series were created, which provide interactive learning of Kannada language as per the state syllabus. Even the process of creation sensitized a large number of teachers to the usefulness of IT tools in education, and enabled them to actively participate in creating the CDs. It was realized that manpower availability in the area of language technology is very poor. The web site LT-IISC is created that

would enthuse both faculty and students at engineering colleges to do their mini projects and main projects in language technology related activities.

This phase of Resource Centre created adequate momentum in creating IT tools for Kannada and made a beginning that would ensure that Kannadigas need not feel alien in their own state.

1. Web Sites And Support To Instruction

1.1 Kannudi

Kannudi is bilingual web site for the benefit of Kannadigas to learn about their language and heritage. Non-Kannadigas will also have an opportunity to know about Karnataka. It is mainly aimed at non-specialists, and provides a first level overview of a large number of topics.



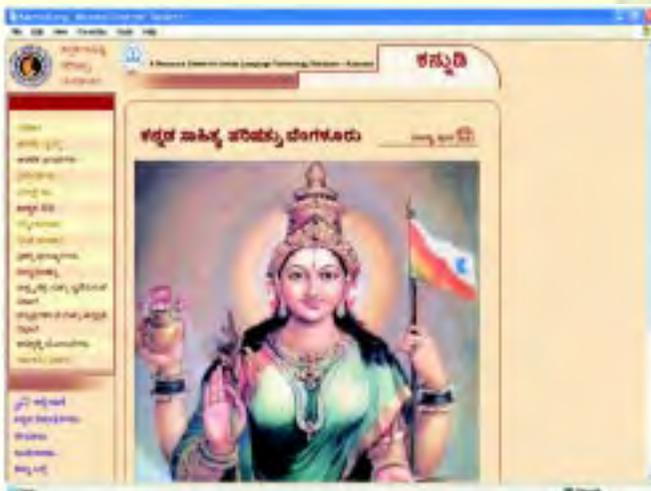
The contents of this site are organized under 11 topics: Language, Epigraphs, Literature, Folklore, Geography, History, Arts, Classics, Personalities, Temples and Festivals, and Cultural Societies and their Activities. Some of the important features of Kannudi are:

- The site has more than 2000 pages of content.
- Most of the articles are written by experts in the respective areas.
- Images have been incorporated at all possible places.

- Brief life sketches of a large number of personalities of Karnataka are prepared.
- Several classics of Kannada literature are made available on the web site.



The web site also hosts the web site of Kannada Sahitya Parishad.



The contents of the site are continuously enhanced. The web site address is www.kannudi.org.

1.2 LT-IISC

This is a web site to meet the needs of students, teachers, and developers interested and concerned with language technologies in general and in particular with emphasis on Kannada language. The main sections of the web site are Standards, language resources, office automation, web

technologies, optical character recognition, speech technologies, machine translation, multi-lingual issues, applications, and open source software. The website will have a download section common to all topics which would also have all the products generated at IISc. The website also features a Discussion Forum on each topic with which interested people can exchange their views regarding technologies, solve problems (if any), etc. Information on events, seminars, workshops in the area of language technologies is also provided. Annotated links to other websites on language technologies, language products, and standards is also given.

The address of the site is www.ltiisc.org.

Publications:

1. Narmadamba, K.: e-learning in Kannada (Kannadadhalli e Kalike) (Book), 2002, Kannada Sahitya Parishat.
2. Narmadamba, K.: Kannada and computers in Samyukta Karnataka, April, 2003.
3. Narmadamba, K.: Speech variation in South India languages with respect to dialect, emotion and style, Sadhane, Bangalore University, May, 2003.

1.3 Bodhana Bharathi: Multimedia Educational CDs for 7th, 8th and 10th Standards

It has been observed that the quality of teaching Kannada language at primary and high school levels needs considerable improvement. It was, therefore, decided to create interactive multimedia learning material, useful to both teachers and students, in the form of CDs. This CD series are named as “Bodhana Bharathi”. These CDs are meant to supplement the text books and classroom teaching. The framework within which these CDs are developed is,

1. To improve listening, and writing skills.
2. To enable the students perform better in the examinations.

3. To appreciate the background in which the lessons were created by the original authors.
4. To make them understand and feel proud about their nationality and culture, literature, people and language.
5. To enjoy poems and songs.

Some of the specific aspects of the instructional material include,

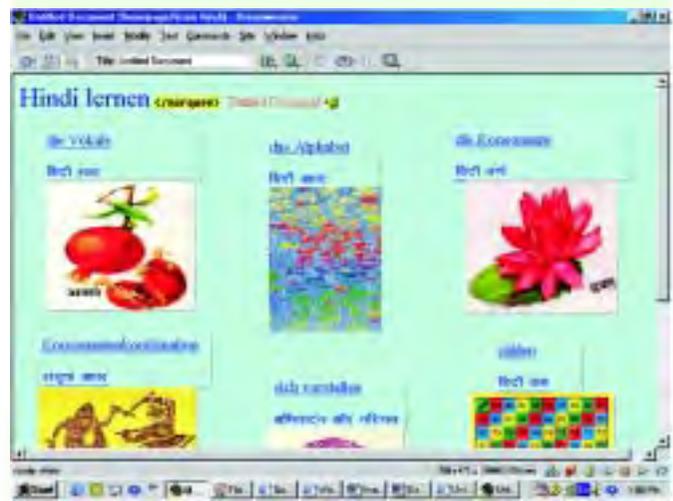
1. Introducing the student to the writers of the lessons and their works.
2. To make them understand the important points and the moral of the lessons.
3. Exposing them to new words that they come across in the lessons.
4. To help them recite poems appropriately.
5. Testing their knowledge in the lessons by giving them unit tests.



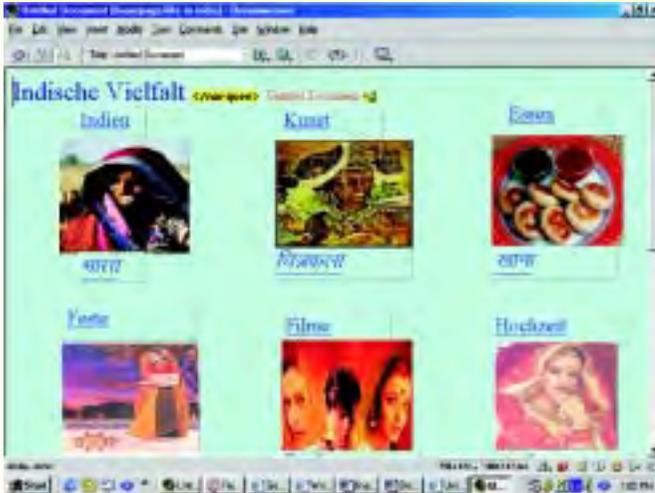
The material is divided into three sections Teacher's section, Student's section and a common section. The teachers section has information such as lesson plan, objectives and others. The students section contains extra questions, new words, etc. The common section comprises of a preamble, summary of the lesson, values in the lesson, key points of the lesson supported by visuals, model question paper with partial interaction and answers, and meaning of words with relations,.

1.4 Bilingual Instructional Aid for learning German through Hindi

Both Hindi and German belong to the Indo-European family of languages. There are a lot of similarities between the two in terms of grammar, vocabulary etc. These two are semantically quite close to each other. As a consequence it is easier to explain an unknown German word, phrase or idiom in Hindi and vice versa than using English. Moreover with the growing number of joint ventures and offshore projects there has been increasing interaction between Germans and Indians. If the visitors from Germany need to stay in India for a longer period, then some knowledge of Hindi would be of great advantage in making professional, social and cultural contacts. Similar would be the case of Indian visitors to Germany. This was the motivation for developing a bilingual instructional aid for Hindi-German in the form of a web-based material. This web site has four sections



Learn Hindi section teaches Hindi alphabets, vowels, consonants, counting and methods of introducing oneself.



Traveling to India section gives information on how to speak at Airport, Restaurant, Shopping complex, Bank, Post Office, Hospital, Railway Station, and for Renting a House.

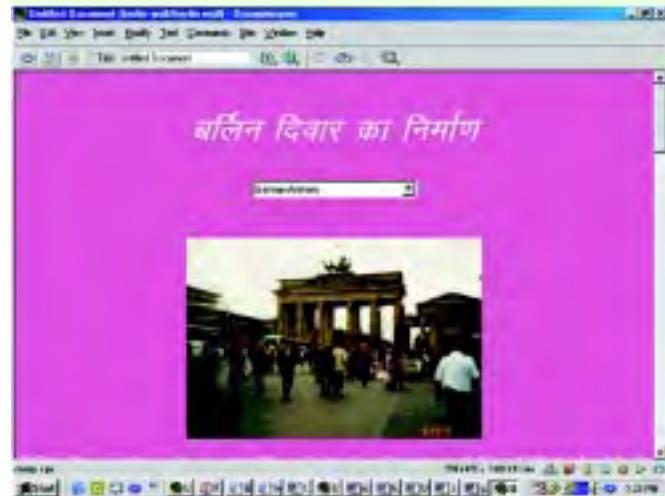
Life in India section has information on culture, lifestyles, tradition, food etc.

Exercises section provides practicing exercises.

This information can be accessed at www.mgmt.iisc.ernet.in/~fls/projects and www.mgmt.iisc.ernet.in/~fls/german

1.5 Information Base in Hindi Pertaining to German History

The Information Base primarily concentrates on bringing information about Germany in perspective, which appeals to Indian school children. The aim of this project is also to put all this material into Hindi so as to make the information accessible to as many school children as possible. NCERT has now undertaken to provide all school children access to information materials both in Hindi and English. This project made use of the knowledge of the foreign language experts at the Institute, and evolved materials in Hindi pertaining to German history. This was done in conjunction with the prescribed syllabus.



The Web-site contains information on German Anthem, Germany in Europe, German Weather, German Population, World War II, Construction and Fall of Berlin wall, German Automobiles, German Beer and Wine, German Culture, German Economy, German Unification, German Literature, Political Parties, Social security system, European Union, Indo-European Languages, and German Cities-Berlin, Brandenburg, Cologne, Hamburg, Frankfurt, Muenich, Trier.

2. Knolwedge Bases

2.1 Sudarshana : A Web Knowledge Base on Darshana Shastras

The main objective was to create a “knowledge base” containing basic texts, commentaries and free translations of texts on the Systems of Indian Philosophy, popularly known as Darshanas, in multimedia form, and have them web-enabled, and make them available in the form of CDs. The target audience is classified into Naïveté, Students, and Scholars.

All the resource material and the experts on all Darshanas are identified. All the web sites on Sanskrit and Darshana-s in particular are surveyed comprehensively. A white paper on Shad-darshana-s is created. A series of lectures by one of the experts on Ishvara in all darshana-s and “Arthasamgraha” is recorded. The audio files and transcribed version of these lectures is made available on the site. Partial translation in English is available. Nyayakusumanjali and other Sanskrit

texts referred for the above discourses are also made available on-line. Glossary of technical terms in Nyaya, Vaisheshika and Mimamsa darshana-s in English and Sanskrit is created. A relevant search engine for both Sanskrit and English is also made available



For any clarifications Dr. N.R.Srinvasa Raghavan can be contacted at raghavan@mgmt.iisc.ernet.in

2.2 Indian Logic Systems

The website“PRAMITI : PRAMana – IT with Indian logic” is being developed to provide comprehensive information on various aspects of Indian Logic and its applicability to Computer Science and Information Technology. The website is designed to be user-friendly by using a concise and objective style. The web address is www.pramiti.org.



2.3 Indian Aesthetics

A website on Indian Aesthetics in Product Design has been created with following features:

- Identifies the body of knowledge in Indian Aesthetics from Literature, Culture and Philosophy.
- A Visual Database of Indian Products, Visuals and Examples.
- Methods for Designers to enable incorporation of “Indianness” into product design.
- Learning materials for students in design programs regarding Indian Aesthetics, and methodology of incorporation into product designs.

The contents are broadly classified as “philosophical”, “cultural”, and “pragmatic”

The philosophical section contains extracts of ancient texts, book reviews, review and research papers. It has a rich visual database of images from crafts, artifacts and general visuals. It has a comprehensive glossary and rich bibliography.



The cultural section has articles on Indian philosophy and Western Philosophy. Case studies containing research in the area of product and visual semantics are also included.



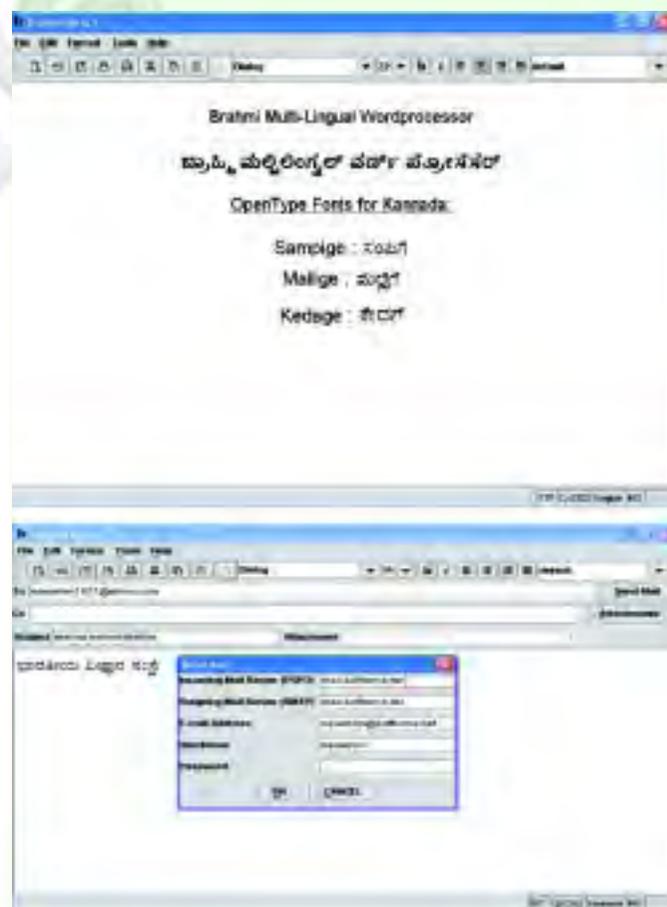
The pragmatic section discusses the learning material for design students

3. Technologies and Language Resources

3.1. BRAHMI: Kannada Indic Input method, Word Processor

Some scripts have hundreds of individual characters, and it is not easy to put all these characters on a standard keyboard that is designed for simpler scripts. Input Methods are developed to provide a better approach for inputting text for these languages. The Brahmi Kannada Input Method (BKIM) is one such method which allows user to give input in Kannada. The BKIM is developed using Java which provides an input method framework that enables the user to enter text directly into the text component. Application developers who wish to have direct Kannada input in their applications can make use of the BKIM. As it is developed in Java, it is also platform independent. The BKIM uses the KGP keyboard layout (standardized by Govt. of Karnataka) and each keystroke is mapped for their corresponding Unicode characters. It uses Open type fonts.

The Brahmi - Multilingual Word processor is a demonstration of Brahmi Kannada Input Method. The word processor has the option of choosing the input method from among nine Indian Languages (at present only Kannada is



enabled) and English. Apart from normal features like file open, print, cut, copy, paste, font features, paragraph settings, color changes it also has some unique features such as saving file in different encodings (UTF-8, UTF-16, Unicode Big Endian, etc.), search and replace Unicode text and email facility. It has an easy access graphical toolbar and comprehensive help.

3.2 OPENTYPE FONTS: Sampige, Mallige, Kedage

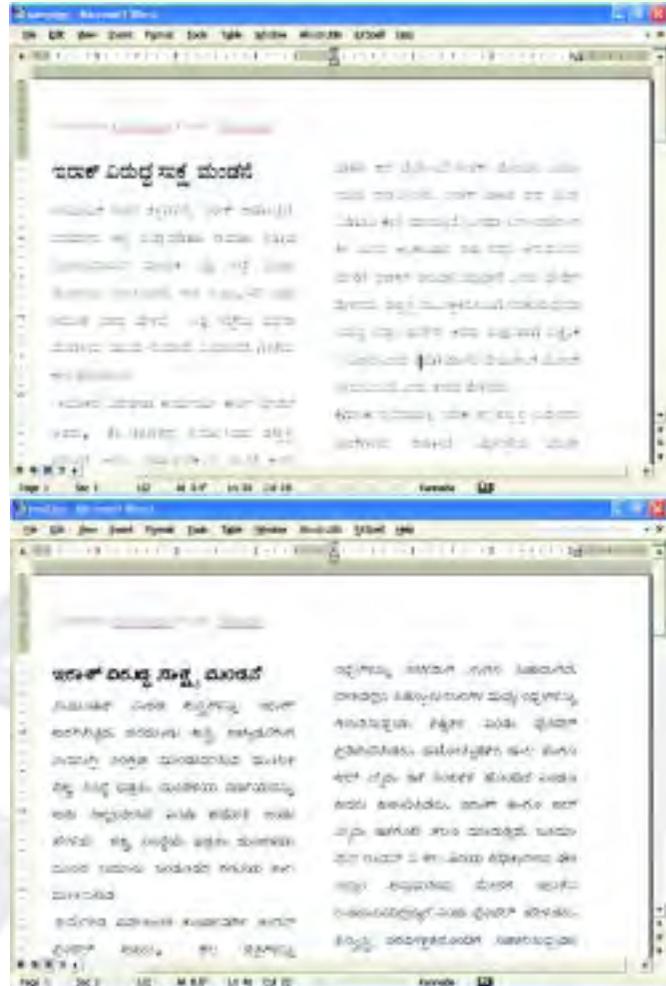
OpenType Font is a new cross platform font file format developed jointly by Adobe Systems Incorporated and Microsoft. Based on Unicode standard the OpenType format is an extension of the TrueType SFNT format that can now support PostScript font data and new Typographic features.

Open Type offers several compelling advantages:

- A single, cross-platform font file that can be used on both Macintosh and Windows platforms
- An expanded character set based on the international Unicode encoding standard for rich linguistic support
- Advanced typographic capabilities related to glyph positioning and glyph substitution that allow for the inclusion of numerous alternate glyphs — such as old-style figures, small capitals and swashes — in one font file
- A compact font outline data structure for smaller font file sizes
- OpenType fonts are best suited for meeting the demand for complex script handling and high quality typography in today's global publishing and communication environment.

Three OpenType fonts called “Sampige”, “Kedage” and “Mallige” are developed. Sampige and Kedage are text based fonts whereas Mallige is a handwritten type of font. The fonts have been

tested and are distributed along with Brahmi – Multilingual Word processor and also individually. People are free to use these fonts in their applications and/or for documentation purposes.



3.3 Kannada Wordnet : WordNet is an on-line lexical reference system whose design is inspired by current psycholinguistic theories of human lexical memory. The nouns, verbs, and adjectives of a language are organized into synonym sets, each representing one underlying lexical concept. Different relations link the synonym sets. Synchronic organization of lexical knowledge and structured organization of words is helpful for natural language processing.

Although the design of Kannada WordNet has been inspired by the famous English WordNet, the unique features of Kannada WordNet are ‘Graded antonyms and meronymy relationships’ and ‘efficient underlying database design’.

Nominal as well as verbal compounding, complex verb constructions also play a vital role.

There are different organizing principles for different syntactic categories. Two kinds of relations are recognized: lexical and semantic in WordNet. Lexical relations hold between word forms; semantic relations hold between word meanings. The basic categories such as nouns, verbs, adjectives and adverbs are organized into synonym sets, each representing one underlying lexical concept. Synsets or the synonymy sets are the basic building blocks. The synonym set serves as an identifying definition of lexical concepts.

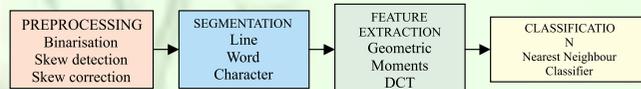
The semantic relations discussed in the Kannada WordNet are Synonymy, Antonymy, Hypernymy –hyponymy, Meronymy – holonymy, Entailment, troponymy. This WordNet will initially be built for 2000 words.



3.4 OCR for Tamil

A, Optical Character Recognition (OCR) system capable of converting multi-lingual manuscripts to machine readable codes is one of the key steps in working towards the goal of machine translation. There are numerous applications that OCR systems have to offer which are of help in day-to-day activities of life. These include mass conversion of existing documents and literature into electronic format, reading aid for the blind as a part of text-to-speech converter, automatic sorting of mails in postal department and processing of bank documents, machine transliteration/translation of documents and literature of other scripts, etc.

A OCR for Tamil script is developed that works in a multi-font and multi-size scenario. The input to the system is a scanned or digitized document and the output is in TAB code. The documents are expected to contain text only. The process sequence is given in the following block diagram.



Preprocessing is the first step in OCR, which involves binarisation, skew detection and correction. Binarisation is the process of converting the input gray scale image scanned with a resolution of 300 dpi into a binary image with foreground as white and back ground as black. Suitable techniques have been used to take care of contrast in the images. The skew angle of the document is estimated using a combination of Hough transform and Principal Component Analysis. Segmentation is done by first detecting lines, and then detecting the words in each line followed by detection of the individual characters in each word. Horizontal and vertical projection profiles are employed for line and word detection, respectively. Connected component analysis is performed to extract the individual characters. The segmented characters are normalized to a

predefined size and thinned before recognition phase.

The segmented symbols are fed to the classifier for recognition. Tamil alphabet set contains 154 different symbols. The characters are divided into some clusters based on domain knowledge, to reduce the recognition time and a smaller probability of confusion. This is accomplished by designing a three level, tree structured classifier to classify Tamil script symbols.

A line in any Tamil text has three different segments; upper, middle and lower. Depending upon the occupancy of these segments, these symbols are divided into one of the four different classes. The number of dots present in each segment contributes a lot for the classification, which is dependent on the scanning resolution. This constitutes the first level clustering

The second level of classification based on Matras/ extensions is applied only to symbols which have upward matras and downward extensions. The classes are further divided into groups, depending on the type of ascenders and descenders present in the character. This classification is feature based i.e. the feature vectors of the test symbol are compared with the feature vectors of the normalized training set. The features used in this level are second order geometric moments and the classifier employed is the nearest neighbor classifier.

Feature-based recognition is performed at the third level. For each of the groups, the symbol normalization scheme is different. The dimensions of the feature vector are different for different groups, as their normalization sizes are different. Truncated Discrete Cosine Transform (DCT) coefficients are used as features at this level of classification. Nearest neighbor classifier is used for the classification of the symbols.

The system has been tested on files taken from Tamil magazines and novels which are scanned at

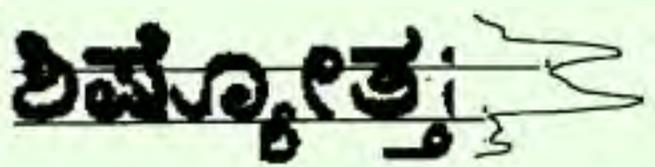
300 dpi (dots per inch). It has been ensured that the test files contain almost all the symbols present in the script. An accuracy of over 99% has been achieved on the training set and 98% on other samples. The sample size was 100.

[1] K G Aparna and A G Ramakrishnan, "A complete Tamil Optical Character Recognition System", *Document Analysis Systems V*, 5th International Workshop, DAS 2002, Princeton, NJ, USA, August 19-21, 2002, pp. 53-57.

[2] K G Aparna and A G Ramakrishnan, "Tamil Gnani – an OCR on Windows", *Proc. Tamil Internet 2001*, Kuala Lumpur, August 26-28, 2001, pp. 60-63.

3.5 OCR of Printed Text Documents in Kannada

The input to the system is the image of the printed page obtained by scanning on a flat bed scanner at 300 DPI resolution and converting the image into binary by making use of a global threshold selected automatically for each page. This image is processed to remove any skew (so that the text lines are aligned horizontally) using a Hough transform based technique. Next, the individual lines in the image and the words in each line are separated using Projection Profile based methods. Due to the special characteristics of Kannada script, separating the individual characters in the word is not an attractive choice. Hence a novel segmentation algorithm is developed which segments words into sub-character level so that each akshara may be composed of many segments. A pattern classification method is used which is based on the Support Vector Machines to assign a classification label to each segment so obtained. After labeling individual segments, rules are used on how aksharas are composed to finally effect

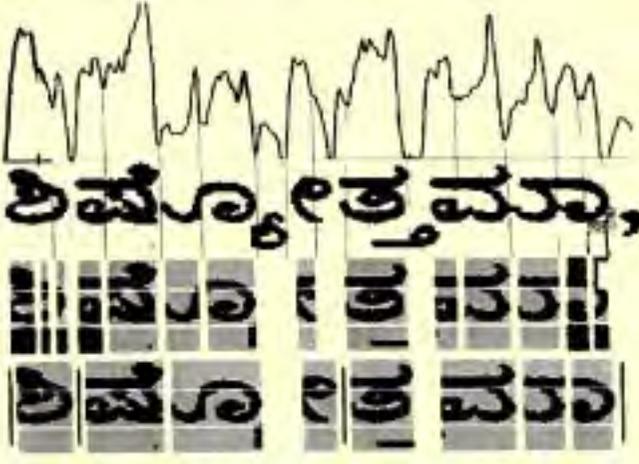


recognition of individual aksharas. The final output of the system is an ASCII file compatible with kantex typesetting package for Kannada which is built around the standard Latex system.

The words are first vertically segmented into three zones as shown in the figure.

The next task is to segment the three zones horizontally. The middle zone is the most critical since it contains a major portion of the letter.

After segmenting the document each segment has to be recognized to effect final recognition.



Features (set of numbers that capture the salient characteristics of the segment image) from each segment are extracted. The characters in Kannada have a rounded appearance. Therefore features which can extract the distribution of the ON pixels in the radial and the angular directions will be effective in capturing the shapes of the characters.



The features extracted have to be classified using a classifier. Here a Support Vector Machine (SVM) classifier is used.

The output after classification is transformed into a format which can be loaded into a Kannada editing package. The input is usually ASCII text in which aksharas are encoded ASCII strings.

ಅದು ಸಿದ್ಧನ ಕಾಯಕದ ಹೊತ್ತು. ನೂರಾರು ದೇವರನ್ನು ನೆನೆ
ಆದ್ರೆ ಸಿದ್ಧನಿಗೆ ಬಾಯಲ್ಲಿ ಯಾವ ದೇವರೂ ಬರಲಿಲ್ಲ. ರುದ್ರಕ
ತಿವಿತ್ತಿದ್ದವು. ರುದ್ರವನ ಗುಡಿಸಿಲಿಂದ ಕೇರಿ ದಾಟಿ, ಊರದ
ಹಿಂದಾಗಿಲಿಗೆ ಬಂದ. ಮಾಯದ ಬೆಳಕು ಮೂಡುತಿತ್ತು. "ತಾ
ಮರುದಿನ ಶಾಲೆಗೆ ತಂದೆ ಬಂದ
ಅವನು ಮನೆಯ ಹೊರಗೆ ಏನೂ
ಕಾಸೂ ಕೊಡುವುದಿಲ್ಲ" ಎಂದರು.
ಸ್ಥಾಪಿಸಿದ್ದರು. ಮೈಸೂರು ಅರಸರು
ತರಿಸಿ ಇಲ್ಲಿ ಇರಿಸಿದ್ದಾರೆ. ಈಗ
ದೊಡ್ಡದಾಗಿ ಮಾಡಿದ್ದಾರೆ. ಹೆಚ್ಚು, ಹೆಚ್ಚು

ಅದು ಸಿದ್ಧನ ಕಾಯಕದ *ಫು ನ*ರಾರ ದೇವರನ್ನು ನೆನೆ
ರು ಸಿದ್ಧನಿಗೆ ಬಾಯಲ್ಲಿ ಯಾವ ದೇವರೂ ಬರಲಿಲ್ಲ ರುದ್ರಕ
ತಿವಿತ್ತಿದ್ದವು ರುದ್ರವನ ಗುಡಿಸಿಲಿಂದ ಕೇರಿ ದಾಟಿ ಊರದ
ಹಿಂದಾಗಿಲಿಗೆ ಬಂದ ಮಾಯದ ಬೆಳಕು ಮೂಡುತಿತ್ತು ಫ
ಮರುದಿನ ಶಾಲೆಗೆ ತಂದೆ ಬಂದ
ಅವನು ಮನೆಯ ಹೊರಗೆ ಏನೂ ಕಾ
ಕಾಸೂ ಕೊಡುವುದಿಲ್ಲ ಎಂದರು
ಸ್ಥಾಪಿಸಿದ್ದರು *ಮೈಸೂರು ಅರಸರು ರ
ತರಿಸಿ ಇಲ್ಲಿ ಇರಿಸಿದ್ದಾರೆ ಈಗ ನ*
ದೊ ಡ್ಡದಾಗಿ ಮಾಡಿದ್ದಾರೆ ಹೆಚ್ಚು ಹೆಚ್ಚು

The system was tested on pages scanned from Kannada magazines, etc. Currently the system recognizes about 85% of the aksharas correctly.

4. Research

4.1 Automatic Classification of Languages Using Speech Signals

Automatic language identification is the problem of identifying the language being spoken from a sample of speech from an unknown speaker. Among various approaches for language identification, the phone recognition offers considerable promise, as it incorporates sufficient knowledge of the phonology of the language to be identified, without incurring the significantly higher cost of word based approaches. An approach based on sub-words that does not require manually labeled data in any of the languages recognized is used. Particularly, the

focus was on the specific architecture termed Parallel Phone Recognition (PPR), and the system is referred to as parallel Sub-word Recognition System.

Research in automatic language identification requires a large corpus of multi-lingual speech data to capture many sources of variability within and across the languages. Among various Indian languages six languages Hindi, Kannada, Malayalam, Marathi, Tamil and Telugu were selected. English as spoken by the same Indian speakers is the seventh language. For each language, twenty adult speakers of different age and gender are selected. Care was taken to ensure that a speaker was chosen for a particular language, only if he/she had that language as a native language particularly during childhood. Speech was collected using a Sennheiser HMD224 noise-canceling microphone and low pass filtered at 7.6 kHz. The recording protocol was designed to obtain Digits and days of the week, Numbers in English, English alphabets, commonly used words in native languages, railway reservation words in native language and in English, Banking words in native and English language. In the elicited free speech, personal details with name, age, native language, profession, family, passage reading in Hindi, English and native languages and etc.

Here three approaches are studied, namely, Parallel Sub-Word Recognition.(PSWR), SWRLM, and Parallel – SWRLM. The following tables show the performance of PWSR system on OGI-TS and ILDB data base with different scores.

Comparison of LID performance of PSWR system on OGITS and ILDB database
MLC with Acoustic some PA
Training : 50 utterances /Language Test : 20 utterances/ Language (OGI-TS) Training : 40 utterances /language Test : 20 utterances/language (ILDB)

	OGI-TS		ILDB	
	Training	Test	Training	Test
No-BR	91.6	67.5	63.3	71.6
BR-Z	90.2	62.3	80.8	76.6

LID accuracy of PSWR system
(6 languages) Averaged over 6 languages
(HI, KA, MA, MR, TA, TL)
For different classifiers (MLC and GC) and different scores (PA, PAL, PL)
Training : 40 utterances / language : Test : 20 utterances / language

		Training			Test		
		P _A	P _{AL}	P _L	P _A	P _{AL}	P _L
MLC	No-BR	63.3	64.6	27.5	71.6	73.5	28.3
	BR-Z	80.8	81.6	27.8	76.6	76.6	30.0
GC	ED	35.0	36.6	22.5	31.6	32.5	23.5
	WED	35.0	36.6	24.16	23.3	24.16	20.8
	MD	91.6	91.2	39.5	84.2	84.0	35.0

This SWRLM approach uses a single front-end sub word recognizer followed by N back-end LMs for an N language LID task. The front-end SWR can be language dependent or language independent. The SWRLM performance will improve if we use language independent single-word unit inventory which is obtained from all languages in LID task.

Results: the table III shows the LID performance of SWRLM on OGI-TS database for each of the front ends for both the cases of MLC. Table IV shows the performance of SWRLM on ILDB for all six front ends and both cases of MLC.

From the above tables we note that LID performance on ILDB is much better than on OGI-TS database for training data. Whereas for test data, OGI-TS database gives better results as compared to ILDB database.

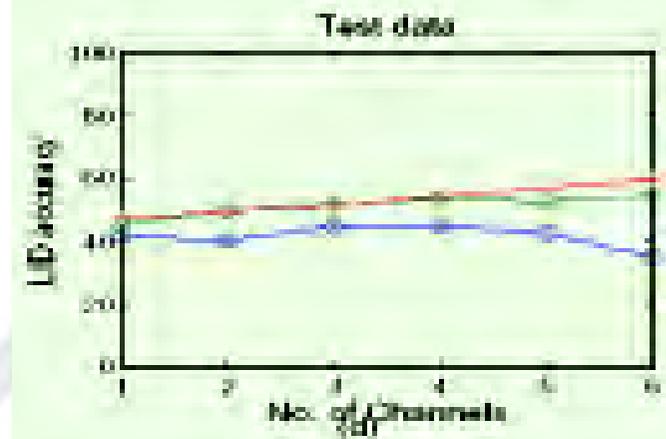
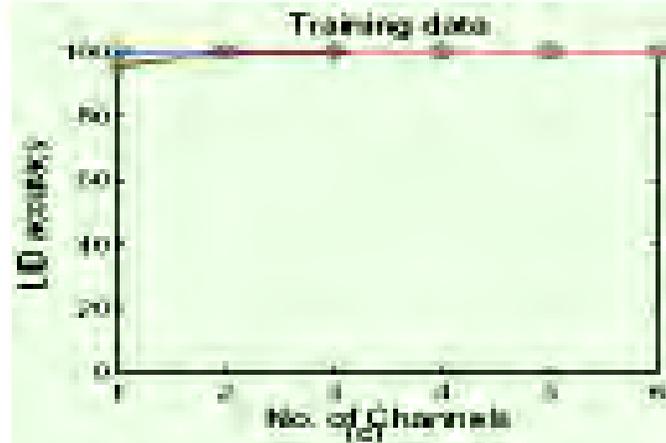
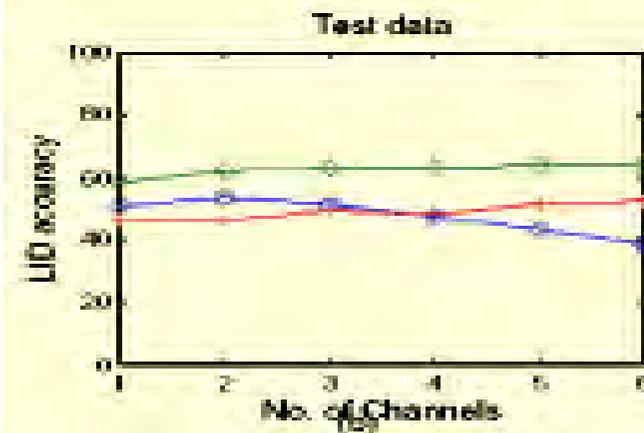
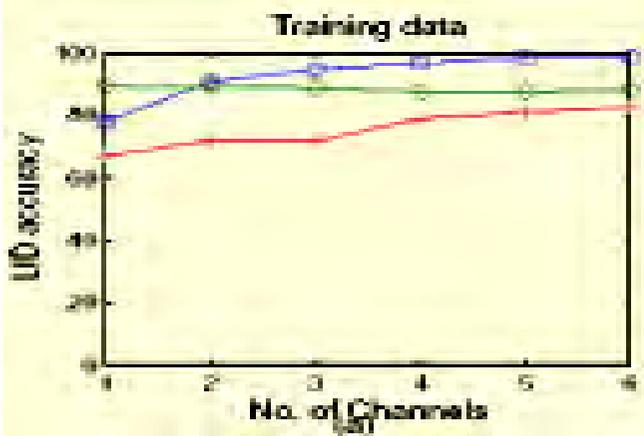
As we know that the sounds in the language to be identified do not always occur in the one language used to train the front-end sub-word recognizer. Thus it seems natural to look for a way to incorporate phones from more than one language into a SWRL like systems. Alternatively another approach is simply to run multiple SWRLM systems in parallel with the single language SWR each trained in different languages. Therefore P-SWRLM uses multiple front-end SWRs with each

No-BR and BR-Z

Front-end	No-BR		BR-Z	
	Training	Test	Training	Test
EN	91.15	54.16	76.84	51.66
GE	89.38	50.83	65.43	46.67
HI	81.83	61.67	36.02	22.5
JA	92.12	56.67	51.76	29.16
MA	91.80	68.33	86.01	63.3
SP	90.83	60.0	86.65	65.0

SWR followed by N number of backend LMs for a N language LID task.

The database and preprocessing with parameterization is identical to that in PSWR approach. The bias removal was done by Zissman method on PI. Results of LID accuracy for six languages Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu PSWRLM system, with two types of classifiers Maximum Likelihood Classifier and Gaussian Classifier are given in the following plots that compare the performance of LID accuracy



on PSWRLM system for two classifiers i.e., MLC and GC-MD, and also across OGI databases and for ILDB database.

The following conclusions may drawn:

- The sub-word approach to LID, developed in our lab holds good promise for LID among a small set of languages.
- The LID performance is quite good among south Indian languages, although they share a lot of phonetic structure and even vocabulary.
- PSWR approach is more promising than the SWRLM and PSWRLM approaches to LID.

4.2 Algorithms for Kannada Speech Synthesis

The basic units, namely, CV, VC, VCV and VCCV have identified and recorded. A framework has been standardized for creation and handling of database of spoken basic units. A synthesis scheme, based on waveform concatenation of basic units has been attempted.

New techniques for pitch detection and modification, and speech synthesis with emotion are proposed.

5. Publications

1. R. Muralishankar, K. Suresh and A. G. Ramakrishnan, "DCT based approaches to Pitch Estimation", submitted to Signal Processing.
2. K. Suresh and A. G. Ramakrishnan, "A DCT based approach to Estimation of Pitch", Proc. Intern. Conf. on Multimedia Processing and Systems, Chennai, Aug. 13-15, 2000, pp. 54-57.
3. R. Murali Shankar and A. G. Ramakrishnan, "Robust Pitch detection using DCT based Spectral Autocorrelation", Proc. Intern. Conf. on Multimedia Processing and Systems, Chennai, Aug. 13-15, 2000, pp. 129-132.
4. R. Murali Shankar and A. G. Ramakrishnan, "Synthesis of Speech with Emotions", Proc. Intern. Conf. on Commn., Computers and Devices, Vol. II, Kharagpur, Dec. 14-16, 2000, pp. 767-770.
5. Anoop Mohan, Harish T, A. G. Ramakrishnan and R. Muralishankar, "Pitch modification sans pitch marking", submitted to Intern. Conf. Acoustics, Speech and Signal Processing, May 7-11, Salt Lake City, Utah, USA, 2001.

*Courtesy: Prof. N.J. Rao
Indian Institute of Science
Centre for Electronics Design and Technology
(CEDT)
Bangalore – 560 012
(RCILTS for Kannada)
Tel. : 00-91-80-3466022, 3942378, 3410764
E-Mail: njrao@mgmt.iisc.ernet.in,
chairman@mgmt.iisc.ernet.in*