# Unicode Based Bilingual Sindhi-English Pictorial Dictionary for Children

**Zeeshan Bhatti[1,*], Imdad Ali Ismaili[1], Dil Nawaz Hakro[1], Ahmad Waqas[2]**

[1]Institute of Information and Communication Technology, University of Sindh, Jamshoro, Paksitan
[2]Department of Computer Science, Sukkur Institute of Business Administration, Pakistan
*Corresponding author: zeeshan.bhatti@usindh.edu.pk

**Abstract**  Dictionaries are the fundamental building block of any language facilitating the language learners, scholars, teachers, students and native speakers with vital understanding of linguistics resources. In this paper we discuss the implementation of Unicode based bilingual dictionary in Sindhi to English and English to Sindhi; targeted specifically for children. There doesn't exists any such type of software system which is specifically for the young school going children of Sindh Province, that would facilitate and aid them in learning English language. This application thus is the first ever attempt to cultivate English to Sindhi and Sindhi to English Pictorial dictionary for Sindhi speaking Children. The system is developed using Java programming language with database managed through Hash-table. The system uses java beans to create and store the dictionary word information as objects and serialized to save it in a file using java's object serialization technique. Initially a two set of user interfaces are created. First is the administrator section that is used to enter and save words and meanings along with their corresponding picture. The second section is the end user section that is used by the students to search and view the Sindhi-English word with picture. The words included are specifically targeted towards the children of Sindh and contain approximately thirty Three thousand words.

**Cite This Article:** Zeeshan Bhatti, Imdad Ali Ismaili, Dil Nawaz Hakro, and Ahmad Waqas, "Unicode based Bilingual Sindhi-English Pictorial Dictionary for Children." *American Journal of Software Engineering* 2, no. 1 (2014): 1-7. doi: 10.12691/ajse-2-1-1.

## 1. Introduction

There are various simple definitions of the term dictionary given by the various authors according to their writing experience and their intellectual approach. In brief, Dictionary is just like a reference book which provided the meaning and additional information such as pronunciations, origin etc. [1]. With advent of technology the use of dictionaries has seen a paradigm shift from hard copy based book, to digital applications available on computers, online or in cell phones. A pictorial or illustration based dictionary uses pictures or drawings to illustrate and visually elaborate what the desired word means. Generally pictorial dictionaries are organized by the topic instead of being an alphabetic list of words [2]. Pictorial dictionaries include a small corpus of words, minimizing the dictionaries size and complexity. The reason is that these dictionaries hold only those words which are represented and elaborated with their pictures [2].

Sindhi is the official language of Sindh province in Pakistan and currently there are more than 34.4 million Sindhi speaking people in Pakistan [3,4] and is regarded as one of the oldest spoken languages of the world [5]. Sindhi is also spoken widely in various regions of India by

approximately 2.8 million people [6]. Unfortunately very little work has been done with respect to Sindhi language for developing a digital Sindhi Dictionaries and this small amount of work can be only seen in [7,8]. However Sindhi Language is influencing and making its way into digital age by means of research conducted in Sindhi NLP and linguistics area; from developing Sindhi Typing Tutor [9], Sindhi Academic portal [10] to developing GUI's in Sindhi [7]. The work done by Mahar et al. in Speech Tagging for Sindhi Language with rule based technique [11] and using Word Net [12] address key problems in Sindhi computational processing in NLP areas. The Sindhi Word Segmentation problem was addressed using Lexicon-driven approach [13] and Sindhi Text was segmented into word tokens in [14] by Mahar et al. the Sindhi Corpus was constructed by Rahman M.U., [15] and techniques for developing Sindhi text to speech system was discussed by [16,17,18]. Nevertheless, there have been selected efforts for the computerization of standard Sindhi dictionary, based on third-party plugins and system dependent features discussed in [8], whereas a generic framework was proposed by Ismaili et al. to develop Unicode based Sindhi Dictionaries [7]. However to best of our knowledge no other work has been done in the area of Sindhi Digital Dictionary development especially with pictorial architecture. English to Sindhi and Sindhi to English Children's Pictorial Dictionary is the first ever

attempt of developing such an application in Sindhi Language and only application which provides a great help to school children and in rural areas because it contains not only Sindhi meanings but also contains pictures.

This paper shows the implementation of Java built framework that is not dependent on any third-party software or plugin. This system can be run on any computer platform with any underlying operating system, as its completely built using Java with its platform independent features. The development of such dictionaries and their provision through Pakistan's National Digital Library would highly facilitate the research work on Sindhi language, literature and linguistics [7].

This paper is further divided into four main sections. In the next subsection we discuss the research and source of various reference materials obtained for this project and later the main aims and objectives of this paper. The section two elaborates the design and methodology of this research project. The graphical user interface and project workflow is discussed in section three whereas section four presents the project evaluation. Finally conclusion and future work is given in section five with acknowledgment in the end.

## 1.1. Aims and Objective of the Project

This dictionary is specially designed for those school children who are not aware of the English words and their Sindhi meanings. Even some of the teenage Sindhi students are not aware and don't understand English with Sindhi meanings. Therefore, this application is also included with images or pictures for every word along with its meaning, so that the young students of Sindh can understand and learn English language by the help of images/pictures.

There are various aims and goals behind developing this application. Some of the specific aims of creating this application are given as under.

### 1.1.1. Sindhi Unicode Based Application

One of the major aims of creating this type of application is to provide a digital dictionary which is based on Sindhi Unicode. This will greatly help to provide convenience in using this application and displaying Sindhi text on old system that are not updated and are available in under developed Sindhi Schools. The Unicode provides a standard platform for displaying Sindhi text Glyphs and typing Sindhi text without additional external plugin support.

### 1.1.2. No Regional Language Support

The main target was to develop the application with all built-in features that does not require any regional language support. Thus we used Java with Unicode character system in such a way so that the user does not need to install any other external plugin, nor any regional language setting in their system. There is also no need to install any external keyboard layout for Sindhi, as the system uses internal keyboard mapping mechanism discussed in "Sindhi Typing tutor Application" [9].

### 1.1.3. No 3rd Party Plugin

This application is very much efficient in retrieval of lexis (words) because we have not involved any third party plugin (any database) that adds its own delay for the retrieval of words, because the control will be passed to it and it will execute the query for the retrieval.

### 1.1.4. No Additional Sindhi Fonts Required

This application does not require any additional Sindhi fonts. This is one of the specific criteria which will be appreciated by everyone. As it uses Java technology with Unicode character system, the need of specific fonts to display character glyphs is no longer required and Sindhi text have been made compatible with standard default fonts available in all computer systems.

### 1.1.5. Platform Independent

This application is platform independent which means that this can be used everywhere, on every system. Its use is not restricted on any specific platform as it uses java virtual machine as its underlying architecture.

### 1.1.6. User oriented Application

This application is designed in such a way that it is user oriented in nature and can easily be operated by a novice user of any age.

### 1.1.7. No Delay in Retrieval

This application is very much efficient in retrieval of words. It does not introduce any delay in retrieval, as it has no third party plug- in and it uses java's Hash table structure to store and retrieve data.

### 1.1.8. Desktop Based Application

This application is desktop based application and is very convenient to use. It guides users especially the young school children to learn their native language.

## 1.2. Research and Reference Material

We have searched and looked through various areas and places in order to design, develop and collect the words with relevant pictures for this Unicode based bilingual Sindhi-English dictionary. The various areas which are searched out for the collection of data and to design this Sindhi Unicode application include numerous websites, books, dictionaries along with Sindhi Language Authority, Hyderabad and Sindhi Adabi Board was also contacted. To accumulate English words, various English dictionaries were studied and used such as such as "Webster's new world college dictionary" [1]. For the attaining Sindhi meanings, "Sindhica Pictorial Dictionary" [19] was used. Sindhi meanings of the English words are obtained from the study of various Sindhi dictionaries. A lot of Sindhi meanings are collected from general observation and also trough manual translations of the English text.

To collect the images/pictures of the words and meanings, some of the small pictorial dictionaries are studied [21] along with MSN-Encarta [20]. We also used Adobe Photoshop tool to design images at our own. Most of the images were generally obtained from google images.

## 2. Design and Methodology

This project of Unicode based Sindhi – English pictorial dictionary consists of two key sections. First section entails the administrator control interface that allows the developer to create, append, modify and update the Sindhi-to-English dictionary along with specifying and uploading the picture illustration of each word. This section of application is developed on the principles discussed in [7]. The architecture of this Bilingual pictorial dictionary is based on Java Technology. The use of Java technology provides platform independence and supports the Unicode character encoding scheme very efficiently. The data structure framework implements Hash table structure of Java for the proficient storage and retrieval of the dictionary data, as opposed to conventional approaches, which make use of third-party databases [7].

### 2.1. Project Activity

When the user runs the children pictorial dictionary, first a welcome screen of the dictionary will be shown. After that pictorial dictionary class will be loaded in memory. Next loading the dictionary class initializes the GUI components which are based on java's swing framework. Then the Dictionary files are loaded into the system using object Input Stream class that reads the file as on java Object. Once the file has been read and the java bean object loaded into the memory, then that object is put inside a hash table structure. As Sindhi characters are not organized in its appropriate sequence in the Unicode character schemes [22], thus adhoc techniques of sorting Sindhi Words are performed at runtime and finally the words are loaded into the GUI. The activity diagram in Figure 1, illustrates the actual working of the pictorial dictionary.
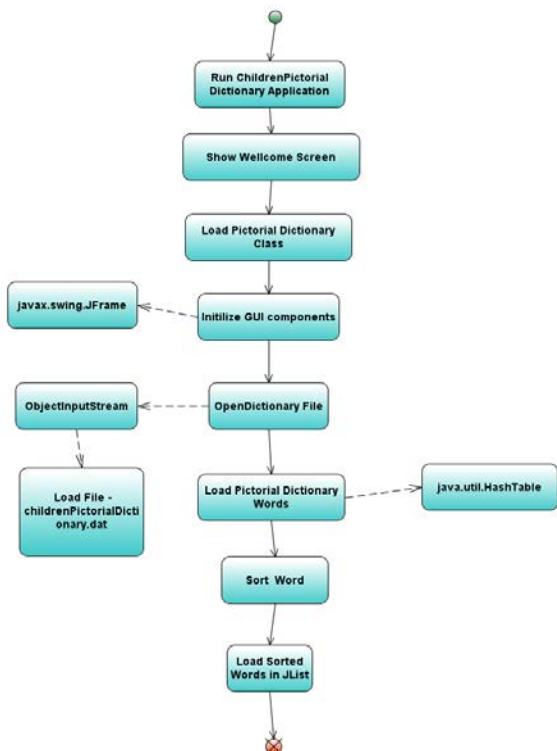


**Figure 1.** Activity Diagram of the system

### 2.2. Class Structure of the Project

The Figure 2 presents the class diagram of children's pictorial dictonary project. The main design of the application is divided into three classes, welcome screen, main class and pictorial_com class. Pictorial_om class is divided further into two classes, Sindhi Unicode and pictorial dictionary class. Pictorial dictionary is further distributed into Sindhi keyboard, Pictorial dictionary database and English dictionary. Sindhi Unicode class is divided into Sindhi Unicode encoding.
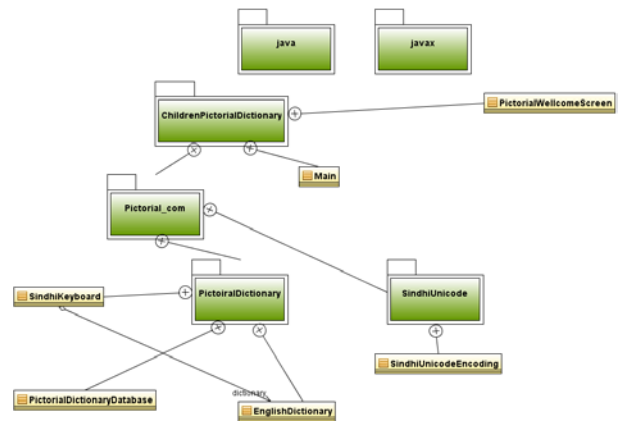


**Figure 2.** Class Structure of the System

The Figure 3 shows the details of various classes along with attributes and Operators.
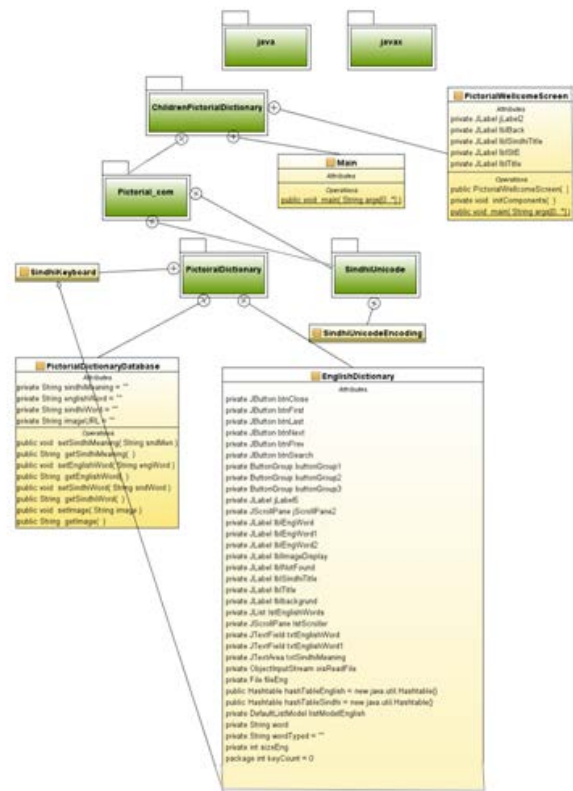


**Figure 3.** Class Diagram of the system

### 2.3. Hash Table Structure

Our system implements Hash table to save the data instead of any conventional database. A hash table or hash map is a data structure that uses a hash function to

efficiently map certain identifiers or keys (e.g., person names) to associated values (e.g., their telephone numbers). The hash function is used to transform the key into the index (the hash) of an array element (the slot or bucket) where the corresponding value is to be sought [23].

The main purpose of using hash tables over other data structures is speed. This advantage is more apparent when the number of entries is large (thousands or more) as aimed in our project to incorporate at least more than Fifty Thousand words in the dictionary. The other reason is that the Hash tables are particularly efficient when the maximum number of entries can be predicted in advance, so that the bucket array can be allocated once with the optimum size.
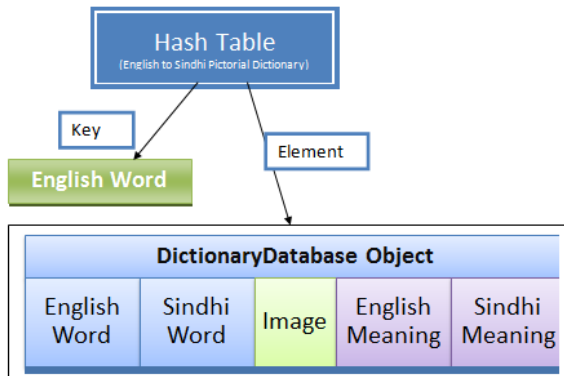


**Figure 4.** Hash table structure for the Dictionary Database object

The Hash table structure used in our dictionary is shown in Figure 4, where we had employed English Word as the Key identifier for English to Sindhi Dictionary and Sindhi Word in Sindhi to English Dictionary. For the associated element, we created a custom java bean object to save various information's such as Sindhi Word, English Word, Image File, English meaning and Sindhi Meaning.

## 2.4. Pictorial Dictionary Database Class



**Figure 5.** Pictorial Dictionary Database Class

A Java Class "Pictorial Dictionary Database" is created where the actual information of the dictionary is saved. The pictorial dictionary database class shown in Figure 5 contains the attributes and the operands for saving the dictionary data in the customized database. As discussed earlier this database consists of hash table and java bean

object. The java beans contain the data which are saved in the hash table, which then is saved as an object in a file. There are four attributes used in the Pictorial dictionary database as shown in Figure 5. The attributes used in the database are private and eight customized operational functions are used to store and retrieve the data from the bean object. These methods and attributes are saved as a single java object, and then this object is later put inside a hash table with corresponding word as key.

The methods declared are for saving and retrieval of English Word, Sindhi Word, English Meaning or Sindhi Meaning and finally fetching image or picture for the corresponding word. This same class would be used for both categories of dictionary that is Sindhi to English and English to Sindhi Dictionary. With record of each Sindhi Word, its relevant information is saved in the Dictionary Database bean as object and then the record is put into the hash table. Finally the entire hash table object is saved into a file for permanent storage usi/ng Object Serialization.

## 3. GUI and System Workflow

As discussed the project consists of two sections; Administrator and User Interface. The administrator section is consists of two internal sections, English to Sindhi dictionary and Sindhi to English Dictionary. The admin of the application selects which dictionary type he wants to work on as shown in Figure 6.



**Figure 6.** Administrator Page for selecting dictionary type



**Figure 7.** Locate the Database file

As the admin selects the dictionary type, he is asked to locate the database file from the computer as shown in Figure 7. One major technique that's been articulated in

our project is also based according to [7]. According to this approach we provide a single input mechanism for both categories of dictionaries. By this we mean that, when the user chooses for English to Sindhi dictionary as its primary input source, the Sindhi meaning entered for each English word is tokenized and separated into multiple Sindhi words as discussed in [24] and the information is saved in Sindhi to English dictionary with its relevant object. Through this approach user only needs to enter a single record for English to Sindhi dictionary, and the Sindhi to English dictionary is automatically created by the system.

The Figure 8 shows the main GUI of the administrator page. The UI has been designed so that the data can be entered and saved quickly and easily. The user is able to insert dictionary records with word, meaning and also load an image file. The administrator can easily traverse the previously entered records, along with editing and deleting the invalid entries.



**Figure 8.** Main GUI of the Administrator Page

## 3.1. GUI for the User

The second section of the project develops the GUI for the end User. The administrator page will only be available to the developers for them to quickly and easily enter the database records. However the User GUI will be distributed publically for the general use. For the General public a simplified welcome screen is designed as shown in Figure 9.

In this figure it is shown that when we select any English word, it will show the Sindhi meaning and its picture also.
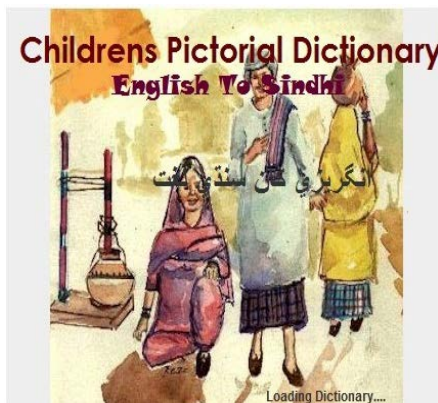


**Figure 9.** Welcome Screen for the general user

As discussed the project involves two separate types of dictionary; English to Sindhi and Sindhi to English. For each type of dictionary, separate GUI has been designed and developed to incorporate the usability aspects, as shown in Figure 10 (a–e).



**Figure 10 (a).** GUI of User section for English to Sindhi Dictionary



**Figure 10 (b).** GUI of User section for English to Sindhi Dictionary
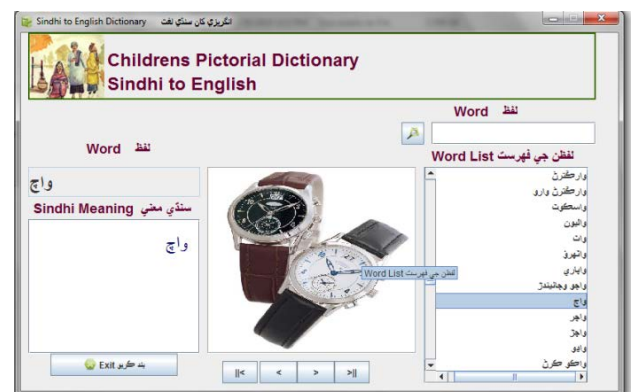


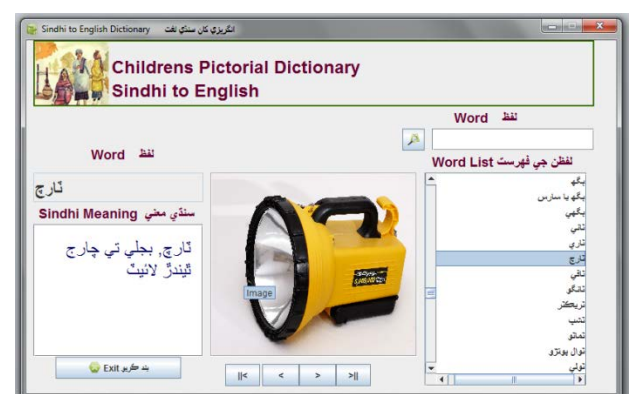**Figure 10 (c).** GUI of User section for English to Sindhi Dictionary



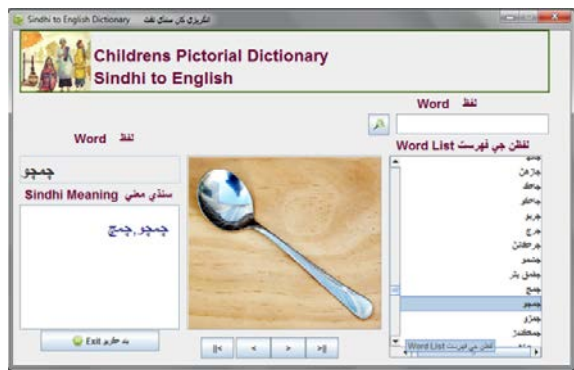**Figure 10 (d).** GUI of User section for English to Sindhi Dictionary

**Figure 10 (e).** GUI of User section for English to Sindhi Dictionary

Figure 10 (a) and Figure 10 (b) shows the interface of English to Sindhi Pictorial Dictionary, whereas Figure 10 (c) to Figure 10 (e) present the interface of Sindhi to English Pictorial Dictionary.

It is essential here to mention that the project does not use and copyrighted images and sources. We have ensured that all the data and pictures are obtained from open source and free with non-copyrighted mediums.

# 4. Project Evaluation

It is observed that in the projects of dictionaries, whether English, Urdu or Sindhi, it is very much necessary that such type of application projects should be user oriented. To achieve this goal we firstly settled meetings with reputed educational personalities, and the various editors, got their notion on how should our Sindhi dictionary be implemented. The second major criteria which we have considered in developing this application is how we can make this dictionary more and more easy for new users especially for children because our application is mainly developed for school children, young Childs and also for people living in rural areas. We also met with some project developers for the suggestions and guidelines in order to work in right directions. After that we studied various projects of the same criteria that have been made for other languages, studied their design methodology and mechanism as discussed in [25] for Pictorial dictionary of Ancient Athens, for printed Farsi subwords [26], Arabic [27] and finally we look into document image database [28].

Finally we conducted usability test and survey amongst students by letting the students use the dictionary and latter asking them to fill in simple questioners regarding the dictionary. Figure 11 shows the results of the usability survey. Amongst the 100 students that were asked to take part in the evaluation, 55 showed complete satisfaction with only 10 being unsatisfied by the system with certain reservations.
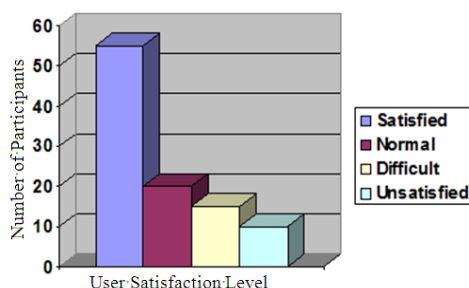


**Figure 11.** Graph showing the satisfaction level of the dictionary users

# 5. Conclusion and Future Work

In this paper we discussed the development of a Unicode based bilingual Pictorial dictionary for children in English to Sindhi and Sindhi to English. There is a great need of such an application, which is user oriented, easy to use and more convenient for children belonging to rural areas of Sindh, especially for young school going children. Our application is not only beneficial for those children who are aware of the Sindhi meaning but also for those who don't know Sindhi meaning. This application does not need any regional language support and also doesn't rely on any external fonts, neither any additional keyboard layout for Sindhi Language is required. This application implements an internal keyboard mapping mechanism. It's based on Unicode character standards and is independent on and 3rd party plugins. The system uses hash table for its internal database where we use java bean class, saved as an object of hash table element. One of the main features is that this dictionary includes the pictures along with every word. The dictionary currently contains approximately thirty three thousand words in its repository along with pictures. In future we intent to increase the number of words and reduce the image size and dictionary size for faster loading and processing.

# Acknowledgement

# References

[1]    Webster's New World College Dictionary, Fourth Edition, 2002.

[2]    Wikipedia, the free encyclopedia "Picture dictionary" URL: http://en.wikipedia.org/wiki/Picture_dictionary. Retrieved on 25/11/2013.

[3]    Allana, G.A., "The Origin and Growth of Sindhi Language", Institute of Sindhology, University of Sindh, Jamshoro, Pakistan, 2002.

[4]    I.A. Ismaili, Zeeshan Bhatti, A. A. Shah "Design and Development of Graphical User Interface for Sindhi Language (GUISL)". Mehran University Research Journal of Engineering & Technology, Volume 30, No. 4, October 2011.

[5]    Allana, G.A., "An Introduction to Sindhi Literature", Sindhi Adabi Board, Jamshoro, Pakistan, 1991.

[6]    Sindhi Language, From Wikipedia, the free encyclopedia, Web URL: http://en.wikipedia.org/wiki/Sindhi_language, Retrieved on January 2013.

[7]    I.A. Ismaili, Zeeshan. Bhatti, A. A. Shah, "Towards a Generic Framework for the Development of Unicode Based Digital Sindhi Dictionaries". Mehran University Research Journal of Engineering & Technology Volume 31, No. 1, January 2012.

[8]    Soomro, H.K., Shah, A.A., and Shaikh, A.A., "Development of Computerized Sindhi to English and English to Sindh Dictionary", Mehran University Research Journal of Engineering & Technology, Volume 23, No. 4, pp. 289-296, Jamshoro, Pakistan, October, 2004.

[9]    Bhatti Z., Ismaili, I.A., Khan, W.I., Nizamani, A.S., (2013) "Development of Unicode based Sindhi Typing System", Journal of Emerging Trends in Computing and Information Sciences, Volume 4, Issue 3, pp-309-314.

[10] Bhatti, Z. , Hakro, D. N., & Jarwar, A. A. (2013). "Sindhi Academic Informatic Portal". American Journal of Information Systems, 1(1), 21-25.

[11] Mahar, J. A., & Memon, G. Q. (2010, February). Rule Based Part of Speech Tagging of Sindhi Language. In *Signal Acquisition and Processing, 2010. ICSAP'10. International Conference on* (pp. 101-106). IEEE.

[12] Mahar, J. A., & Memon, G. Q. (2010). Sindhi Part of Speech Tagging System using WordNet. *International Journal of Computer Theory and Engineering*,*2*(4), 538-545.

[13] Mahar, J. A., Memon, G. Q., & Danwar, S. H. (2011). Algorithms for Sindhi word segmentation using Lexicon-driven approach. *International journal of academic research*, *3*(3).

[14] Mahar, J.A., Shaikh, H., & Memon, G. Q (2012) Model for Sindhi Text Segmentation into Word Tokens. *Sindh University ResearchJjournal (Science Series)* , Vol. 44(1) pp.4*3*-48.

[15] Rahman M U (2010). Towards Sindhi Corpus Construction, Conference on Language and Technology, Lahore, Pakistan.

[16] Mahar J A (2012). Statistical approaches to diacritics res-toration in Sindhi text to speech synthesis system, Ph. D Thesis, Hamdard University, Karachi, Pakistan.

[17] Shaikh, H., Mahar, J. A., & Malah, G. A. (2013). Digital Investigation of Accent Variation in Sindhi Dialects. *Indian Journal of Science and Technology*, *6*(10), 5429-5433.

[18] Shah, A. A., Ansari, A. W., & Das, L. (2004). Bi-Lingual Text to Speech Synthesis System for Urdu and Sindhi. In *National Conf. on Emerging Technologies* (pp. 20126-130).

[19] Gopang, H.B.A., Memon, F.R., (2007) "Sindhica Pictorial Dictionary: English to Sindhi", Sindhica, 2007. Pp. 143.

[20] Encarta, M. S. N. (2003). Dictionary.

[21] Arora, B. (2010). Pictorial Dictionary. Har-Anand Publications.

[22] Hussain, S., & Durrani, N. A Study on Collation of Languages from Developing Asia. PAN Localization, Center for Research in Urdu Language Processing National University of Computer and Emerging Sciences, Lahore, Pakistan.

[23] Wikipedia, The Free Encylcopedia "Hash Table", http://en.wikipedia.org/wiki/Hash_Table (retrieved on August 203).

[24] Bhatti, Z. Ismaili, I.A., Soomro, W.J., and Hakro, D., (2013) "Word Segmentation Model for Sindhi Text" *merican Journal of Computing Research Repository, 1(1).*

[25] Travlos, J. (1971). Pictorial Dictionary of Ancient Athens (pp. 42-42). London: Thames and Hudson.

[26] Ebrahimi, A., & Kabir, E. (2008). A pictorial dictionary for printed Farsi subwords. Pattern recognition letters, 29(5), 656-663.

[27] Al-Nafjan, A. (2010). The Design and Development of Labib: An Arabic Sign Language Educational Tool. In Proceedings of World Conference on Educational Multimedia, Hypermedia and Telecommunications 2010 (pp. 3231-3234). Chesapeake, VA: AACE. Retrieved November 27, 2013 from http://www.editlib.org/p/35103.

[28] Akbari, M., & Azimi, R. (2010, February). Document image database indexing with pictorial dictionary. In Second International Conference on Digital Image Processing (pp. 75462R-75462R). International Society for Optics and Photonics.