



Urdu Component Development Project

Urdu SpellChecker Application v1.0

October 31, 2007

**CENTER FOR RESEARCH IN URDU LANGUAGE PROCESSING
NATIONAL UNIVERSITY OF COMPUTER AND EMERGING SCIENCES, LAHORE
PAKISTAN**

Table of Contents

1	Introduction.....	4
2	Application Usage	4
3	File Formats	5
3.1	<i>Config File</i>	5
3.2	<i>Urdu Alphabet File</i>	5
3.3	<i>Wordlist File</i>	5
3.4	<i>Bi-Gram Posting Lists File</i>	5
3.5	<i>Bi-Gram PL Index Matrix</i>	5
3.6	<i>Uni-Gram Posting Lists File</i>	5
3.7	<i>Uni-Gram PL Index List</i>	6
3.8	<i>Hash Table File</i>	6

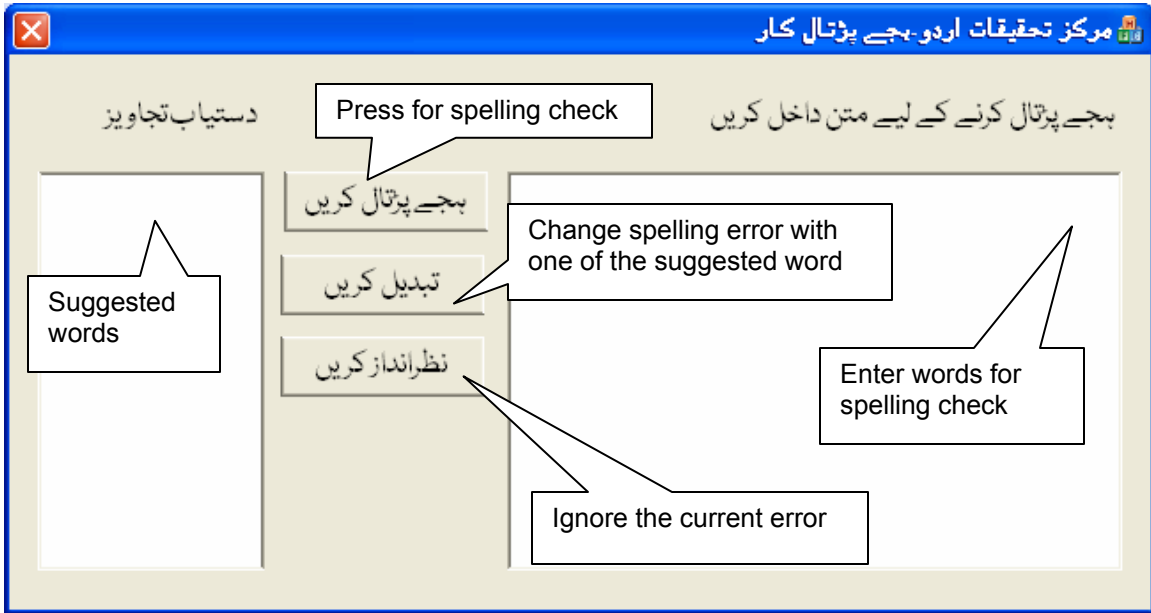
Revision History

Name	Change Date	Version	Description of Changes
Atif Gulzar	31-10-2007	1.0.0.0	Initial document

1 Introduction

SpellChecker is a sample application developed using SpellChecker Utility (released by CRULP). SpellChecker application provides support for spelling check for Urdu language. This document enlists SpellChecker application specifications and its usage.

2 Application Usage



3 File Formats

3.1 Config File

This file contains the file paths to be used in the API. This file also contains initialization status of the application. If the initialize status is set to “YES” in config.ini, the *initialize* function loads the wordlist into hash table and generates the bigram and unigram files. Then it sets the initialize status to “NO”. If the initialize status is set to “NO”, the *initialize* function loads the precompiled hash table, unigram and bigram. This file should exist at the root path of application.

3.2 Urdu Alphabet File

This file (in Unicode format) contains a list of letters in Urdu alphabet, in order. Format of the file is simple; each line has only one letter. The order of letters is used to index into ‘posting list index matrices or lists’ (i.e. Bigram).

3.3 Wordlist File

This file (in Unicode format) contains a list of words with frequencies. The first line of this file contains the total number of words in the file. The remaining each line contains a word and its frequency.

3.4 Bi-Gram Posting Lists File

This file contains a number of posting lists. Each posting list contains indices to word forms in the hash table. These posting list indices are 4-byte numbers. A particular posting list contains indices to all words in which a particular bi-gram exists at any word position. These posting lists are arranged according to bi-grams and are referred from the offsets and counts in the Bi-Gram PL Index Matrix.

3.5 Bi-Gram PL Index Matrix

The 2D matrix contains indices to Bi-gram posting lists. A Bi-Gram is pair of characters that can occur in a word at any position. The matrix is a square matrix i.e. it has equal number of rows and columns, each of size equal to number of characters defined at start of the file. The matrix contains elements for all possible bi-gram combinations of allowed character set defined at start of file. Each matrix element (which is an index to one of the bi-gram posting lists) would take 6 bytes: 4 bytes for posting list offset and 2 bytes for number of following records in the posting list. Please note that, if a particular bi-gram does not exist i.e. if a particular element in the Bi-Gram PL Index Matrix does not point any word in the posting lists file, then its offset field would contain -1 and count field would contain 0.

3.6 Uni-Gram Posting Lists File

The file contains a number of posting lists. Each posting list contains indices to word forms in the hash table. These posting list indices are 4-byte numbers. A particular posting list contains indices to all words having word length less than or equal to 3 characters, and in which a particular uni-gram exists at any word position. These posting lists are arranged according to uni-grams and are referred from the offsets and counts in the Uni-Gram PL Index List.

3.7 Uni-Gram PL Index List

Uni-Gram PL Index List is a list for uni-grams, which have word length less than or equal to 3 characters, and which can occur at any word position. Each list element contains offsets and count of posting list. If a particular uni-gram does not exist, its offset field would contain -1 and count field would contain 0. The size of the list is the same as the size of listed character set in Urdu alphabet file.

3.8 Hash Table File

This file is in binary format and contains all searchable lexemes (without diacritics). These word forms are given positions (or bucket) in table based of hashing. Universal hash function will be used to generate hash key. Linear chaining within the same array will be implemented for probing using "Next" field. Please note that entry can be single-word words or multi-words compounds. Each bucket is of fixed length and contains following fields:

Field	Word without Diacritics and spaces	Frequency	Next
Size (in Bytes)	40	4	4

Total Size (in Bytes) = 48

Detail of Lexicon Fields

Word without Diacritics

A null terminated string of 20 Unicode characters (total 40 bytes required) that represents word form lexeme without spaces and marked diacritics. Please note that all lexemes in the table may not be unique and duplication is allowed.

Frequency

A 4-byte integer that contains the frequency of a word obtained form a large corpus.

Next

A 4-byte integer that points to another bucket of the same Hash Table by storing the index of next bucket in the collision list. In result, a list of collided words is maintained in the same hash table. The last element of the list has 0 in its "Next" field. In case, a collision occurs, the collided word is added at the end of the list.