

# Design and Implementation of Punjabi Spell Checker

Gurpreet Singh Lehal

Department of Computer Science, Punjabi University, Patiala, India.  
gslehal@yahoo.com

## Abstract :

Spellcheckers are the basic tools needed for word processing and document preparation. Designing a spell checker for Indian languages such as Punjabi poses many new challenges not found in English, which complicates the design of the spell checker. Punjabi language is far different from Western languages in phonetic properties and grammatical rules. Thus the existing algorithms and techniques that are being used to check the spelling and to generate efficient suggestions for mis-spelt words of English and other Western languages are not actually suitable for Punjabi; rather it needs different algorithms and techniques for expected efficiency. This paper presents the complete design and implementation of a Punjabi spell checker.

## Keywords :

Punjabi, Spellchecker, Suggestion List, Typing errors

## Introduction :

Spellcheckers are the basic tools needed for word processing and document preparation. A spell checker is a tool that enables us to check the spellings of the words in a text file, validates them i.e. checks whether they are right or wrongly spelled and in case the spell checker has doubts about the spelling of the word, suggests possible alternatives.

The main steps performed by the spell checker are:

1. Take the word from the file as its input.
2. Pre-process the word
3. Look for the word in dictionary
4. In case, the word exists, pass onto the next one.
5. If the word is not found, then seek for the closest matching patterns and put them up in the form of suggestions.

Even though this appears to be very simple at first glance but designing a spell checker for Indian languages such as Punjabi poses many new challenges not found in English, which complicates the design of the spell checker. Punjabi language is far different from Western languages in phonetic properties and grammatical rules. Thus the existing algorithms and techniques that are being used to check the spelling and to generate efficient suggestions for mis-spelt words of English and other Western languages are not actually suitable for Punjabi; rather it needs different algorithms and techniques for expected efficiency. Some of the typical problems faced during designing the Punjabi are:

1. There is no standardization of Punjabi keyboard layouts. There are more than forty keyboard layouts and fonts commonly being used, which means that the same Punjabi word can be internally stored in forty different ways. As for example, the word ਪੰਜਾਬੀ is internally stored as

- pMjwbl in *Akhar* font

- pMj`bl in *Amrit-Lipi2* font
- pμj;bl in *Anandpur Sahib* font
- gzikph in *Asees* font
- ê³ÛÄiÆ in *Satluj* font

The spell checker has to deal with each of these cases separately and check for the spellings. Even in the same font, a character can be typed and stored in more than one way. As for example, in *Asees* font, the character ਫ਼ can be typed by pressing either a single key *P* or a combination of two keys ;*a*. Similarly the character addak ( ਿ ) can be typed by pressing key *Z* or key ~.

2. Another problem, typical to Indian language scripts was faced. Since Punjabi is not written in linear fashion, so same word could be internally stored in more than one way. The user can type consecutive occurring semi-vowels/upper vowels and half character/lower matra ( \_ or \_ ) in any order and they all look visually similar. For example, the word ਖੁਲ੍ਹੇ can be typed by the typist as ਖ + \_ + ਲ + ਯ + ਿ or as ਖ + \_ + ਲ + ਿ + ਯ and both will be displayed as ਖੁਲ੍ਹੇ. Thus if ਖੁਲ੍ਹੇ has been typed and stored in the database as ਖ + \_ + ਲ + ਯ + ਿ if the user types ਖੁਲ੍ਹੇ as ਖ + \_ + ਲ + ਿ + ਯ then the word will be signalled as wrongly spelled.

3. In some of the Punjabi fonts, the Punjabi characters such as *bindi*, *lava*, *onkar*, *dulainkar* etc. have zero width and so if by mistake a user makes multiple entries of such characters only a single entry is visible. If the spell checker flags such word as misspelled the user will not come to know where the error. Thus for example, consider the word ਪ੍ਰੰਤੂ it has been typed as ਪ ਿ ਿ ਿ ਿ ਿ ਿ ਿ ਿ ਿ and visually the word looks correct but internally it has stored wrongly and the user will not be aware where the error lies.

4. Unlike English, there is no well defined word boundary for Punjabi words written in different Punjabi fonts. As for example, in *Asees* font the following punctuation marks are encoded as Punjabi characters and thus are part of the word ( “ + / : ; ? [ ] \ { } ). But there are many other fonts such as *Akhar*, *Satluj* etc. which do not encode the above punctuation marks as Punjabi characters. So the extraction of word boundary is font dependent in case of Punjabi.

5. There is no standardization of Punjabi spellings. A word may be spelled in more than one way and all the forms may be acceptable.

## Brief Description of Punjabi Language :

Punjabi is the world's 12<sup>th</sup> most widely spoken language. The populace speaking Punjabi is not only confined to North Indian states of India such as Punjab, Haryana and Delhi but is spread over all parts of the world. Punjabi is a phonetic language and commonly written in Gurmukhi script. Some of the major properties of the Gurmukhi alphabet are:

- Gurmukhi script alphabet consists of three (ੳ ਅ ਏ) vowel-carrier letters and nine vowel signs. By using the vowel signs with the three vowel-carrier letters ten vowels are obtained. The vowel-carriers ੳ and ਏ are never used without a vowel sign.
- There are 38 consonants and all the ten vowel-signs are used with all the other consonants (Fig. 1).
- In addition there are two nasal signs ' (bindi) and ˆ (tippi) used for sounds produced through nasal cavity. The symbol *adhak* ˆ, is used to produce the sound of a double consonant.
- There are three half characters in Gurmukhi alphabet. The complete character set of Gurmukhi is depicted in Fig. 1.

<u>Vowels</u>	
ੳ ਆ ਇ ਈ ਉ ਊ ਏ ਐ ਓ ਔ	
<u>Vowel carriers</u>	
ੳ ਅ ਏ	
<u>Consonants</u>	
ੳ ਚ ਛ ਜ ਝ ਵ ਟ ਠ ਡ ਢ ਠ	
ਤ ਥ ਦ ਧ ਨ ਪ ਫ ਬ ਭ ਮ	
ਯ ਰ ਲ ਵ ਝ ਸ ਜ ਖ ਫ ਗ ਲ	
<u>Matras</u>	
ˆ ˆ ˆ ˆ ˆ ˆ ˆ ˆ ˆ ˆ	
<u>Vowel Modifiers or Half Vowels</u>	
ˆ ˆ ˆ	
<u>Half Characters</u>	
ˆ ˆ ˆ	

Fig 1 : Gurmukhi character set

### Lexicon Creation :

The first step in development of the spell checker is the creation of a lexicon of correctly spelled words, which will be used by the spell checker to check the spellings as well as generate the suggestions. Two issues are involved in lexicon development:

- **Size of the lexicon** : There are two approaches followed for storing the lexicon. The first approach stores the root words of a language and the rest of the words are derived from these root words. The other approach is to store all the possible words of the language in the lexicon. We have followed the second approach and stored all the possible forms of words of Punjabi lexicon. Around 1.5 lakh words were identified and stored in the database. The words are arranged according to the word size. During program execution the words are loaded into AVL trees and sixteen different AVL trees used for different word lengths. The number of nodes in each of these AVL tree is shown in Fig. 2.

- **Format of the words** : As there is no standardized Punjabi keyboard, a word in Punjabi may be written in more than forty different ways. It was necessary to formalize a format for storing the lexicon. One option was to store a word in ISCII (Indian Standard Code for Information Interchange) or Unicode. But that option was dropped, since the coding schemes are not closer to how a user actually types. As for example, if a user types ਜ਼ਿੰਸ then it will be stored in ISCII as (ਯ + chr(232) + ਰ + ਿ + ਿ + ਸ), while the user will be typing it as (ਿ + ਯ + ਿ + ਿ + ਸ). It is necessary to store the words in same order as they are typed as this knowledge is helpful in generating the suggestion list. It was also necessary

to assign coding beyond ASCII 127 to the characters, since in most the databases the searching algorithms do not differentiate between upper and lower case characters. Thus two characters assigned codes corresponding to ASCII 'a' and 'A' will be treated to be same for sorting and searching purposes. Also care was taken that if a half character or a lower matra ( ˆ or ˆ) and a semi-vowel occur together, then semi-vowel will be stored after the half character or lower matra. The coding scheme displayed in Table 1 is used to codify the Punjabi letters.

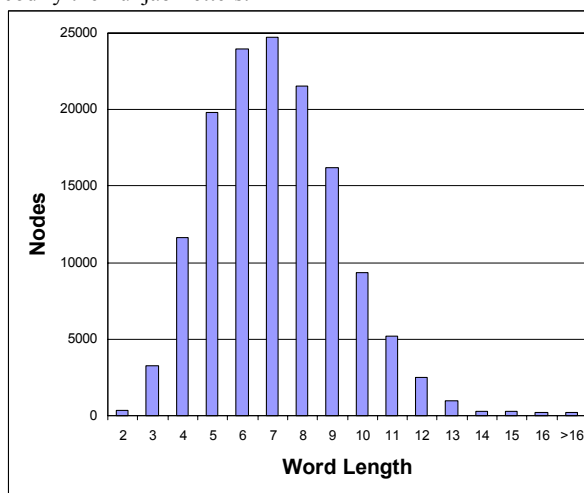


Fig. 2 : Node count in AVL trees of different word lengths

Table 1 : Gurmukhi coding scheme

Code	Punjabi Symbol	Code	Punjabi Symbol
128	ˆ	149	ੜ
129	ˆ	150	ਫ਼
130	ˆ	151	੠
131	ੳ	152	੡
132	ੴ	153	੣
133	ਅ	154	ੵ
134	ੲ	155	੷
135	ਸ	156	ੳ
136	ਸ਼	157	ਬ
137	ਹ	158	ਦ
138	ਕ	159	ਧ
139	ਖ	160	ਨ
140	ਖ਼	161	ਪ
141	ਗ	162	ਫ
142	ਗ਼	163	ਫ਼
143	ਘ	164	ਬ਼
144	ਙ	165	ਭ
145	ਚ	166	ਮ
146	ਛ	167	ਯ
147	ਜ	168	ਰ
148	ਜ਼	169	ਲ
170	ਲ਼	178	ੀ

171	ਐ
172	ਐ
173	ਐ
174	ਐ
175	ਐ
176	ਐ
177	ਐ

179	ਐ
180	ਐ
181	ਐ
182	ਐ
183	ਐ
184	ਐ

**Spell Checker Architecture :**

The major components of the architecture as shown in Fig. 3 are: Tokenisation and normalisation Pre-processing Module, Lexicon Lookup/Error Detection Module and Suggestion Module.

**Tokenisation :**

Tokenisation refers to process of breaking the text into tokens or words using punctuation marks and spaces as delimiters. In case of Punjabi, the punctuation delimiters vary from font to font. As for example, in font *Asees* the punctuation mark ; and : are used to store a Punjabi characters. Similarly, the locations corresponding to the characters ` and ~ are used by majority of the Punjabi fonts. So we have to associate with each font, the valid character set. The spellchecker reads the text character by character along with its font information and uses the information about the character set of that font to tokenise the text. The tokenisation process will check each character in the character set corresponding to that particular font. As for example, the following text ਜੋ ਦਿੱਲੀ; ਵਿੱਚ ਰਹਿੰਦੇ ਹਨ

is encoded internally as follows in *Asees* font:

i' fd~bhl ft~u ojzd/ jB

The tokeniser will break the line into following tokens:

i' fd~bh ft~u ojzd/ jB

We note that in the second token above, the character l has been excluded since it is not part of the

character set for *Asees* font while the characters ' ~ and / have been included in the tokens, since they are part of the character set.

**Normalisation :**

The tokens are then made to pass through a normalisation process to convert them to the format in which the lexicon has been stored. Some of the major steps that are performed in this stage are : removal of redundant zero width characters, changing the order of upper and lower characters and mapping the token to format compatible with the lexicon. Thus for example, the word ਰਹਿੰਦੇ in typed in any Punjabi font will be normalised as {137, 179, 130, 158, 181}.

**Lexicon Lookup/Error Detection :**

In this module the normalised token is searched in the standard dictionary. The standard dictionary has been partitioned into sixteen sub-dictionaries based on the word length and at execution time each of these sub-dictionary is loaded in a height balanced binary search tree (AVL tree). To search for a word of length n, we look for its presence in the AVL tree storing words of length n. Thus ਰਹਿੰਦੇ will be searched in the AVL tree storing lexicon words of length 5. If the word is not found, then it is searched in the AVL tree storing the lexicon of user defined dictionary. If the word is still not found, then it is marked and sent to Suggestion module.

**Suggestion :**

In this module a list of possible correct words is presented to the user. The user selects the correct word, if it is present in the suggestion list, and can give the command to replace a single or all occurrences of the mis-spelled word with the word selected in the suggestion list. The details of this module are discussed in next section.

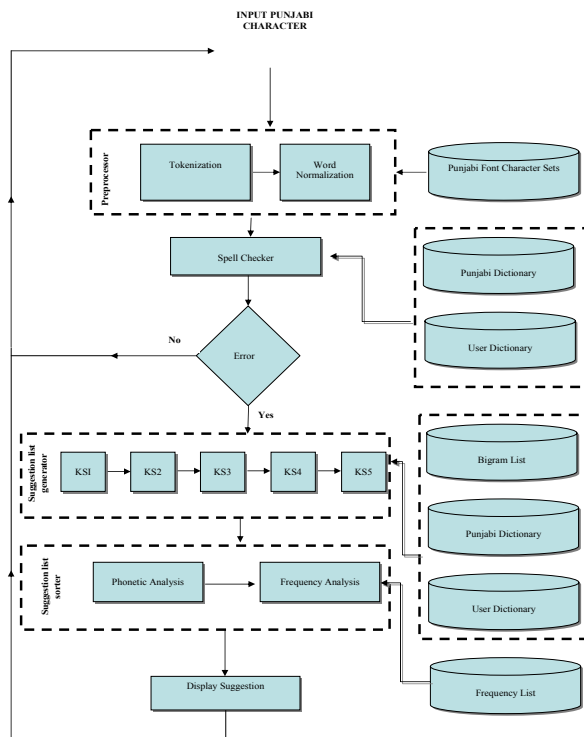
**Suggestion List Generation :**

Once the system has detected an erroneous word, it performs the following steps:

1. Generate a list of candidate corrections
2. Rank the spelling variations
3. Select the highest ranking as the most likely correction

The task of general purpose spelling correction has a long history (e.g. Damerau, 1964; Rieseman and Hanson, 1974; McIlroy, 1982), traditionally focusing on resolving typographical errors such as insertions, deletions, substitutions, and transpositions of letters that result in *unknown words* (i.e. words not found in a trusted lexicon of the language). Typical word processing spell checkers compute for each unknown word a small set of in-lexicon alternatives to be proposed as possible corrections, relying on information about in-lexicon word frequencies and about the most common keyboard mistakes (such as typing *m* instead of *n*) and phonetic/cognitive mistakes, both at word level (e.g. the use of *acceptable* instead of *acceptible*) and at character level (e.g. the misuse of *f* instead of *ph*).

Several approaches based on minimum edit distance, similarity key, rules, N-grams, probability and neural nets are proposed to accomplish the task [1-5]. Of these, minimum edit distance based approaches are the most popular ones. The minimum edit distance is the minimum number of editing operations (insertions, deletions and substitutions) required to transform one text string into another. The distance is also





- ਪੰਜਾਬ ਪੰਜ ਪੰਜਾ ਪੰਜੇਬ ਪੰਜੇ ਪੰਜੀ

**Evaluating the Suggestion List :**

Evaluating the Suggestion List provided by a Spell Checker involves considering three factors:

- Whether the desired word appears on the suggestion list;
- the length of the list of words offered, and;
- the position of the desired word on the list.

The ideal spell checker would offer one word on the list, and that would be correct. If ten words are offered and the correct word is near the bottom, a poor speller must read through the choices and disregard the first few despite their precedence, instead looking for the word which they know, in some way, to be correct. If so many words are offered that they cannot all be seen at once, it is even more difficult to find the correct word.

**Test words :**

We used 225 most commonly mis-spelled words to analyse the performance of the spell checker. The words were drawn from several sources:

- Punjabi corpus prepared by Ministry of Information Technology
- Online Punjabi Newspapers
- Online Punjabi stories
- Punjabi Research Reports

Table 2 : Spellchecker test words, alphabetically listed

ਉਹਨਾ	ਆਣ	ਸਮੇ
ਉਚੀਆਂ	ਆਂਦੀ	ਸਾਹਬ
ਉਚੇ	ਆਵਾਜ਼	ਸਾਬਿਤ
ਉਤਪੱਤੀ	ਐਵੇਂ	ਸਿਘ
ਉਦਾਰਹਣ	ਇਸਤੋਂ	ਸਿਧਾਤਾਂ
ਉਨੱਤੀ	ਇਹਨਾ	ਸਿਰਫ
ਉਨ੍ਹਾ	ਇਹਨੇ	ਸੁਹੱਪਣ
ਉਨਾਂ	ਇੱਕਠ	ਸੁਟ
ਉਨੀ	ਇੱਕਲੇ	ਸੁਟਿਆ
ਉਪਰਲੀ	ਇੰਜਨੀਅਰ	ਸੁਨਣ
ਉਪਲੱਬਧ	ਇਜ਼ਤ	ਸੁਰੂ
ਉਰਫ	ਇੰਨਾ	ਸੂਫੀ
ਊਸ	ਇਨਾਂ	ਸੋਲਾਂ
ਉਹ	ਏਨ੍ਹਾਂ	ਸ਼ਨਾਖਤ
ਅਖਬਾਰ	ਸਖਤ	ਸ਼ਰੀਫ
ਅਖਬਾਰਾਂ	ਸਖਤੀ	ਸ਼ੁਧ
ਅਖਾਂ	ਸਤੱਹੀ	ਹਸ
ਅੰਗਰੇਜ਼ੀ	ਸਬਜੀ	ਹਸਦਿਆਂ
ਅਗੋਂ	ਸਬੰਧਤ	ਹਸਦੀ
ਅਥਰੂ	ਸੰਬੰਧਤ	ਹਸਦੇ
ਅਧਿਅਨ	ਸਬੰਧਾਂ	ਹਕ
ਅਫਗਾਨਿਸਤਾਨ	ਸੱਭ	ਹਥੋਂ
ਅਮ੍ਰਿਤਸਰ	ਸਮਸਿਆ	ਹਦ
ਆਕੇ	ਸਮਾਜਕ	ਹਲ
ਆਜ਼ਾਦੀ	ਸਮੁਚੇ	ਹੂਣ
ਕੱਠੇ	ਦਸਦੀ	ਬੁਲ੍ਹਾਂ
ਕਢ	ਦਸਦੇ	ਭਣੌਈਆ

ਕਢਣ	ਦਾਹੜੀ	ਭਣੌਈਏ
ਕੱਲੀ	ਦਿਤਾ	ਭੁਖਾ
ਕੱਲੇ	ਦੁਨੀਆਂ	ਮਜਬੂਰ
ਕਾਗਜ਼	ਦੇਕੇ	ਮੱਦਦ
ਕਾਨੂਨ	ਦੇਖਕੇ	ਮਧ
ਕਾਲਿਜ਼	ਦੇਂਦੀ	ਮਨੁਖੀ
ਕਿੰਨਾਂ	ਦੇਂਦੀਆਂ	ਮਰਜੀ
ਕਿਵੇ	ਦੇਵੀ	ਮਾਨਣ
ਖਿਚ	ਦੇਵੇ	ਮਾਫ
ਖਿਲਾਫ	ਨਈਂ	ਮਾਂਬਾਪ
ਖੁੱਲਾ	ਨਜਰ	ਮਾਰਕੇ
ਗਲਬਾਤ	ਨਫਰਤ	ਮਿਤਰ
ਗੁਸੇ	ਨਾਹੀਂ	ਮਿਤਰੇ
ਗੁਰੂ	ਨਾਵਾਂ	ਮਿਲਕੇ
ਘਟਨਾਂਵਾ	ਨਿਸਚਾ	ਮੀਤੀ
ਚਕਰ	ਨਿਸਚਿਤ	ਮੁਸਕਾ
ਚਿਠੀ	ਨਿਸਚੇ	ਮੁਸਕਾਉਦੀ
ਚਿੰਨ	ਨਿੱਕਲ	ਮੁਸਕਾਹਟ
ਚੀਜ਼ਾਂ	ਨਿਕੇ	ਮੁੰਹ
ਚੁਕਿਆ	ਨੁਸਖਾ	ਮੁਕ
ਚੁਕੀਆਂ	ਨੁਸਖੇ	ਮੁਢ
ਛਡ	ਨੂ	ਮੁਨਾਫਾ
ਛਡਣ	ਨੋਂ	ਮੁਫਤ
ਛਾਂਵੇ	ਪਛਮੀ	ਮੁੜਕੇ
ਛੁਟ	ਪਰਕਾਸ਼	ਮੈਂਨੂੰ
ਜਹੀ	ਪਰਬੰਧ	ਮੋਤ
ਜਹੀਆਂ	ਪਰਾਪਤ	ਯੂਧ
ਜਹੇ	ਪਰੋਗਰਾਮ	ਯੁਨੀਵਰਸਿਟੀ
ਜਜ	ਪੜ੍ਹਕੇ	ਰਖਕੇ
ਜੱਥੇ	ਪੜ੍ਹਣ	ਰਖੇ
ਜਦੋਂ	ਪ੍ਰਸਿਧ	ਰਫਤਾਰ
ਜਾਕੇ	ਪ੍ਰੰਪਰਾ	ਰਾਜੀ
ਜਾਨਣ	ਪ੍ਰੰਪਰਾਵਾਂ	ਲਗਣ
ਜਿਨਾਂ	ਪ੍ਰੰਤਤਾ	ਲਗਿਆਂ
ਜਿਵੇ	ਪੁਛਦਾ	ਲਗੀਆਂ
ਜ਼ਿਲਾ	ਪੁਜ	ਲਭ
ਜ਼ਿਲੇ	ਪੁਜੇ	ਲਭਦਾ
ਟੁਟ	ਪੁਤਰ	ਲਾਈਕ
ਤਰ੍ਹਾ	ਫੁਲਾਂ	ਲਾਕੇ
ਤੂੰ	ਬਜੇ	ਲਿਜਾਣ
ਤੂਸੀ	ਬਣਕੇ	ਲੇਕੇ
ਤੂੰ	ਬਣਾਕੇ	ਵਖ
ਤੈਂਨੂੰ	ਬਣਣ	ਵਖਰਾ
ਥਾਵਾਂ	ਬਣਾਉਣ	ਵਖਰੀ
ਥੋਡੇ	ਬਣਾਣ	ਵਖਰੇ
ਦਸਣ	ਬਾਹਾਂ	ਵਫਾਦਾਰੀ
ਦਸਣਾ	ਬੁਲ੍ਹ	ਵਰਨਣ
ਵਰਨਣਯੋਗ	ਵਿਚੇ	ਵੇਖਕੇ

These words were fed to the spell checker and statistics such as the position of the word in the suggestion list, the maximum, minimum and average size of the suggestion list were generated (Table 3). It was observed that for 4.39% of the words, the correct word was not displayed in the suggestion list, while for 81.14% of wrongly spelled words, the correct word was displayed on the top of the suggestion list. The minimum and maximum size of the suggestion list is 1 and 52 respectively. The average size of the suggestion list is 15, which is on a higher side, though in 93.4% of cases, the correct word is present in top 10 words of the suggestion list.

**Conclusion :**

This is the first time that a spell checker for Punjabi language has been designed and implemented. The spell checker is part of the commercial Punjabi word processor Akhar. We have only taken care of non-real word errors. Detection and correction of real word errors is a subject of further research.

Table 3 : Average position of the correct word in suggestion list

Position	Percentage of Occurrence
Top	81.14%
Top 3	85.53%
Top 5	89.03%
Top 10	93.42%
In List	95.61%

**References :**

1. E. Brill and R. Moore, " An improved error model for noisy channel spelling correction," *Proceedings of the ACL 2000*, 2000, 286-293.
2. A. R. Golding, " A Bayesian hybrid method for context- sensitive spelling correction ," *Proceedings of the Workshop on Very Large Corpora*, 1995, 39-53.
3. A.R. Golding and D. Roth, " Applying winnow to context-sensitive spelling correction," *Proceedings of ICML*, 1996, 182-190.
4. K. Kukich, " Techniques for automatically correcting words in a text," *Computing Surveys* 24(4), 1992, 377-439.
5. E.J. Yannakoudakis and D. Fawthrop, "An Intelligent spelling corrector," *Information Processing and Management* 19(12), 1983, 101-108
6. F.J. Damerau, " A technique for computer detection and correction of spelling errors," *Communications of ACM* 7(3), 1964, 171-176.
7. V.I. Levenshtein , "Binary codes capable of correcting deletions, insertions and reversals,". *Sov. Phys. Dokl.*,10, 1966, 707-710.
8. G.S. Lehal and M. Bhagat, " Error Pattern in Punjabi Typed Text," *Proceedings of International Symposium on Machine Translation, NLP and TSS*, 2004, 128-141.