

N-grams: A Tool for Repairing Word Order Errors in ill-formed Texts

Theologos Athanaselis, Stelios Bakamidis, Ioannis Dologlou and Konstantinos Mamouras

Abstract—This paper presents an approach for repairing word order errors in English text by reordering words in a sentence and choosing the version that maximizes the number of trigram hits according to a language model. A possible way for reordering the words is to use all the permutations. The problem is that for a sentence with length N words the number of all permutations is $N!$. The novelty of this method concerns the use of an efficient confusion matrix technique for reordering the words. The confusion matrix technique has been designed in order to reduce the search space among permuted sentences. The limitation of search space is succeeded using the statistical inference of N-grams. The results of this technique are very interesting and prove that the number of permuted sentences can be reduced by 98,16%. For experimental purposes a test set of TOEFL sentences was used and the results show that more than 95% can be repaired using the proposed method.

Keywords—Permutations filtering, Statistical language model N-grams, Word order errors, TOEFL

I. INTRODUCTION

SYNTAX is the word used to describe relationships of words in sentences of a language. What appears to be given in all languages is that words can not be randomly ordered in sentences, but that they must be arranged in certain ways, both globally and locally. For example, in English the normal way of ordering elements is subject, verb, object (Boy meets girl) [1]. Subjects and objects are composed of noun phrases, and within each noun phrase are elements such as articles, adjectives, and relative clauses associated with the nouns that head the phrase (the tall woman who is wearing a hat). On the other hand, there are languages that appear a word order freedom like Modern Greek. It is a highly flexible language when it comes to word order. The functions of the nouns are very clear due to the morphological forms. In English, the position of the nouns tells the listener what role the nouns play. Hence the strict rule of SVO (subject-verb-object) does not apply to Greek. Native speakers of a language seem to have a sense about the order of constituents of a

Manuscript received February 9, 2006. Theologos Athanaselis is with the Institute for Language and Speech Processing, Greece, P.C 15125 (corresponding author to provide phone: +302106875416; fax:+302106875300; e-mail: tathana@ilsp.gr).

Stelios Bakamidis, is with the Institute for Language and Speech Processing, Athens, Greece, P.C 15125 (e-mail: bakam@ilsp.gr).

Ioannis Dologlou is with the Institute for Language and Speech Processing, Athens, Greece, P.C 15125 (e-mail: ydol@ilsp.gr).

Konstantinos Mamouras, is with the Institute for Language and Speech Processing, Athens, Greece, P.C 15125 (e-mail: kmam@ilsp.gr).

phrase, and such knowledge appears to be outside of what one learns in school [2].

Automatic grammar checking is traditionally done by manually written rules, constructed by computer linguists. Methods for detecting grammatical errors without manually constructed rules have been presented before. Atwell [3] uses the probabilities in a statistical part-of the speech tagger, detecting errors as low probability part of speech sequences. Golding [4] showed how methods used for decision lists and Bayesian classifiers could be adapted to detect errors resulting from common spelling confusions among sets such as “there”, “their” and “they’re”. He extracted contexts from correct usage of each confusable word in a training corpus and then identified a new occurrence as an error when it matched the wrong context. Chodorow and Leacock [5] suggested an unsupervised method for detecting grammatical errors by inferring negative evidence from edited textual corpora. Heift [6],[7] released the German Tutor, an intelligent language tutoring system where word order errors are diagnosed by string comparison of base lexical forms. Bigert and Knutsson [8] presented how a new text is compared to known correct text and deviations from the norm are flagged as suspected errors. Sjobergh [9] introduced a method of grammar errors recognition by adding errors to a lot of (mostly error free) unannotated text and by using a machine learning algorithm.

Unlike most of the approaches, the proposed method is applicable to any language (language models can be computed in any language) and does not work only with a specific set of words. The use of parser and/or tagger is not necessary. Also, it does not need a manual collection of written rules since they are outlined by the statistical language model. A comparative advantage of this method is that avoids the laborious and costly process of collecting word order errors for creating error patterns. Finally, the performance of the method does not depend on the word order patterns which vary from language to language and for that reason it can be applied to any other language with less fixed word order.

The paper is structured as follows: the architecture of the entire system and a description of each component follow in section 2. The language model is described in section 3. The 4th section shows how permutations are filtered by the proposed method. The 5th section specifies the method that is used for searching valid trigrams in a sentence. The results of using TOEFL’s experimental scheme are discussed in section 6. Finally, the concluding remarks are made in section 7.

II. SYSTEM ARCHITECTURE

Writers sometimes make errors that violate language's grammar e.g (sentences with wrong word order). This paper presents a new method for repairing sentences with word order errors that is based on the conjunction of a new confusion technique with a statistical language model. It is straight forward that the best way for reconstructing a sentence with word order errors is to reorder the words. However, the question is how it can be achieved without knowing the attribute of each word. Many techniques have been developed in the past to cope with this problem using a grammar parser and rules. However, the success rates reported in the literature are in fact low. A way for reordering the words is to use all the possible permutations. The crucial drawback of this approach is that given a sentence with length N words the number of all permutations is $N!$. This number is very large and seems to be restrictive for further processing. The novelty of the proposed method concerns the use of a technique for filtering the initial number of permutations. The process of repairing sentences with word-order errors incorporates the followings tools:

1. a simple, and efficient confusion matrix technique
2. and language model's trigrams and bigrams.

Consequently, the correctness of each sentence depends on the number of valid trigrams. Therefore, this method evaluates the correctness of each sentence after filtering, and provides as a result, a sentence with the same words but in correct order.

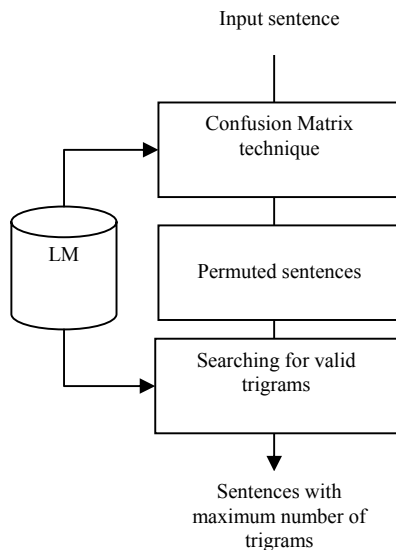


Fig. 1 The architecture of the proposed system

III. LANGUAGE MODEL

The language model (LM) that is used subsequently is the standard statistical N-grams [10]. The N-grams provide an estimate of $P(W)$, the probability of observed word sequence W . Assuming that the probability of a given word in an utterance depends on the finite number of preceding words,

the probability of N-word string can be written as:

$$P(W) = \prod_{i=1}^N P(w_i | w_{i-1}, w_{i-2}, \dots, w_{i-(n-1)}) \quad (1)$$

One major problem with standard N-gram models is that they must be trained from some corpus, and because any particular training corpus is finite, some perfectly acceptable N-grams are bound to be missing from it. That is, the N-gram matrix for any given training corpus is sparse; it is bound to have a very large number of cases of putative "zero probability N-grams" that should have some non zero probability. Some part of this problem is endemic to N-grams; since they can not use long distance context, they always tend to underestimate the probability of strings that happen no tot have occurred nearby in their training corpus. There are some techniques that can be used in order to assign a non zero probability to these zero probability N-grams. In this work, the language model has been trained using BNC and consists of trigrams with Good-Turing discounting [11] and Katz back off [12] for smoothing. BNC contains about 6.25M sentences and 100 million words.

TABLE I
 THE NUMBER OF DIFFERENT ELEMENTS OF LANGUAGE MODEL

<i>Elements of language model</i>	<i>number</i>
<i>unigrams</i>	<i>126062</i>
<i>bigrams</i>	<i>8166674</i>
<i>trigrams</i>	<i>8033315</i>

The next figure depicts the number of trigrams for different spaces of logarithmic probabilities. Note that the minimum logarithmic probability of trigrams is -5,84 while the maximum equivalent is very close to zero. The log scale have been split into 100 equal spaces and the figure shows that the 80% of trigrams have log probabilities greater than -3,33.

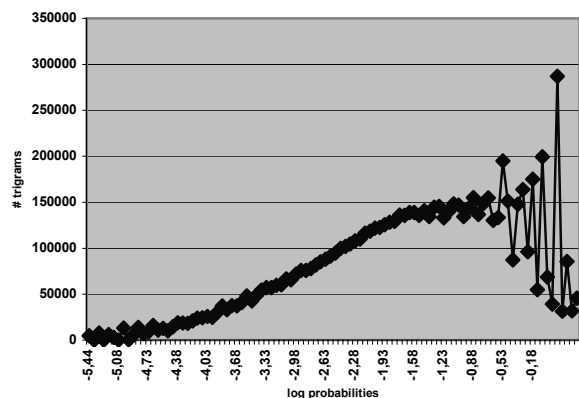


Fig. 2 The distribution of trigrams according to their probabilities

IV. FILTERING PERMUTATIONS

Considering that an ungrammatical sentence includes the correct words but in wrong order, it is plausible that generating all the permuted sentences (words reordering) one of them will be the correct sentence (words in correct order). The question here is how feasible is to deal with all the permutations for sentences with large number of words. Therefore, a filtering process of all possible permutations is necessary. The filtering involves the construction of a confusion matrix $N \times N$ in order to extract possible permuted sentences.

Given a sentence $a = [w[0], w[1], \dots, w[n-1], w[n]]$ with N words, a confusion matrix $A \in R^{N \times N}$ can be constructed,

TABLE II
 CONSTRUCTION OF THE CONFUSION MATRIX $N \times N$, FOR A GIVEN SENTENCE
 $a = [w[0], w[1], \dots, w[n-1], w[n]]$

WORD	w[0]	w[1]	w[n]
w[0]	P[0,0]	P[1,0]	P[n,0]
w[1]	P[0,1]	P[1,1]	P[n,1]
.
w[n]	P[0,n]	P[1,n]	P[n,n]

The size of the matrix depends on the length of the sentence. The objective of this confusion matrix is to extract the valid bigrams according to the language model. The element $P[i, j]$ indicates the validness of each pair of words $(w[i]w[j])$ according to the list of language model's bigrams. If a pair of two words $(w[i]w[j])$ cannot be found in the list of language model bigrams then the corresponding $P[i, j]$ is taken equal to 0 otherwise it is equal to one. Hereafter, the pair of words with $P[i, j]$ equals to 1 is called as valid bigram. Note that, the number of valid bigrams is M lower than the size of the confusion matrix which is $N(N-1)$, since all possible pairs of words are not valid according to the language model. In order to generate permuted sentences using the valid bigrams all the possible words' sequence must be found. This is the search problem and its solution is the domain of this filtering process.

As with all the search problems there are many approaches. In this paper a left to right approach is used. To understand how it works the permutation filtering process, imagine a network of N layers with N states. The factor N concerns the number of sentence's words. Each layer corresponds to a position in the sentence. Each state is a possible word. All the states on layer 1 are then connected to all possible states on the second layer and so on according to the language model.

The connection between two states (i, j) of neighboring layers $(N-1, N)$ exists when the bigram $(w[i]w[j])$ is valid. This network effectively visualizes the algorithm to obtain the permutations. Starting from any state in layer 1 and moving forward through all the available connections to the N -th layer of the network, all the possible permutations can be obtained. No state should be "visited" twice in this movement.

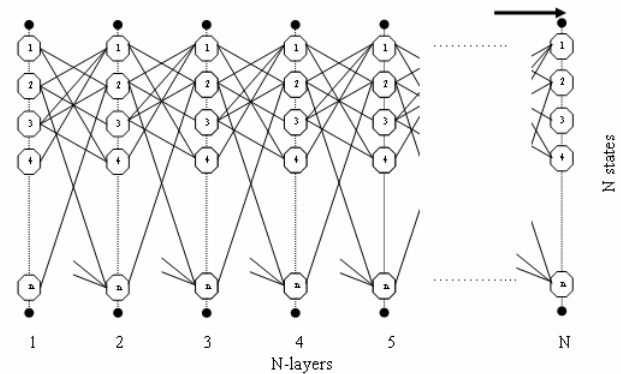


Fig. 3 Illustration of the lattice with N -layers and N states

V. SEARCHING VALID TRIGRAMS

The prime function of this approach is to decompose any input sentence into a set of trigrams. To do so, a block of words is selected. In order to extract the trigrams of the input sentence, the size of each block is typically set to 3 words, and blocks are normally overlapped by two words. Therefore, an input sentence of length N , includes $N-2$ trigrams.

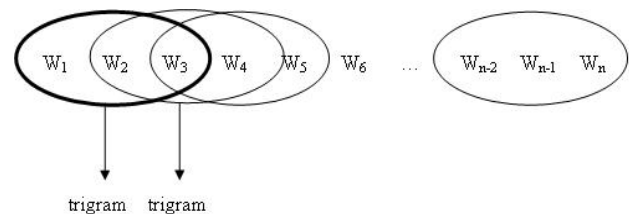


Fig. 4 It illustrates the way of decomposing the sentence into a set of bigrams and trigrams. The input sentence has the following words order $W_1 W_2 \dots W_{n-2} W_{n-1} W_n$

The second step of this method involves the search for valid trigrams for each sentence. In the third step of this method the number of valid trigrams per each permuted sentence is calculated. Considering that the sentence with no word-order errors has the maximum number of valid trigrams, it is expected that any other permuted sentence will have less valid trigrams. Although some of the sentence's trigrams may be typically correct, it is possible not to be included into the list of LM's trigrams.

The plethora of LM's trigrams relies on the quality of corpus. The lack of these valid trigrams does not affect the

performance of the method since the corresponding trigrams of the permuted sentence will not be included into LM as well. The criterion for ranking all the permuted sentences is the number of valid trigrams. The system provides as an output, a sentence with the maximum number of valid trigrams. In case where two or more sentences have the same number of valid trigrams a new distance metric should be defined. This distance metric is based on the total logarithmic probability of the trigrams. The total logarithmic probability is computed by adding the logarithmic probability of each trigram, whereas the probability of non valid trigrams is assigned to -100. Therefore the sentence with the maximum probability is what the system responses.

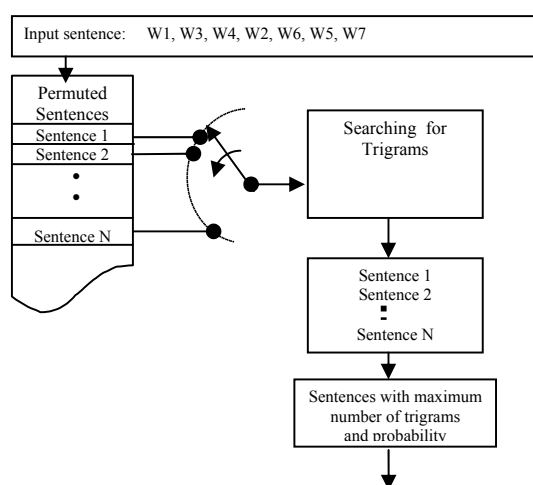


Fig. 5 The architecture of subsystem for repairing sentences with word-order errors. It is based on the algorithm of searching valid trigrams according to the LM

VI. EXPERIMENTATION

A. Experimental Scheme

The experimentation involves a test set of 550 sentences. These sentences have been selected randomly from the section “Structure” of TOEFL past exams [13],[14]. The TOEFL test refers to the Test of English as a Foreign Language. The TOEFL program is designed to measure the ability of non-native speakers to read, write and understand English as used at college and university in North America. The Structure section focuses on recognizing vocabulary, grammar and proper usage of standard written English. There are two types of questions in the Structure section of the TOEFL test. One question type presents candidates with a sentence containing a blank line. Test-takers must choose a word or phrase that appropriately fills in the blank.

The other question type consists of complete sentences with four separate underlined words. Candidates must choose which of the four underlined answer choices contains an error in grammar or usage. For experimental purposes our test set consists of sentences for TOEFL’s word order practice. These sentences are selected from the list of the answer choices but

are not the correct ones. Note that the test sentences are not included into the training set of the statistical language model that is used as tool for the proposed method and 90% of the test words belong to the BNC vocabulary (training data). The goal of the experimental scheme is to confirm that the outcome of the method (sentence with best score) is the TOEFL’s correct answer. It is shown that the corpus contains sentences of length between 4 and 12 words.

B. Error’s Profile

A report of gathered data of this study is presented in the current section. It discusses a categorization of sentences found in the test set according to the length and the type of each sentence; and also it describes the distribution of errors in the whole test data and in different types of sentences [15]. The table below depicts the number of corpus’ sentences as a function of their length.

TABLE III
 NUMBER OF SENTENCES WITH RESPECT TO THEIR LENGTH

# of words per sentences	number of TOEFL sentences
7	72
8	67
9	78
10	93
11	105
12	135

The following figure depicts the percentage of different type sentences in the test set. As it is appeared the test set contains 319 positive sentences which constitute the 58% of the total sentences. Next most frequent type of sentences is the questions with 31%. The negative sentences are 6 times less frequent than positive sentences, with 10% in total. Finally, the imperative sentences constitute the 1% of the total sentences.

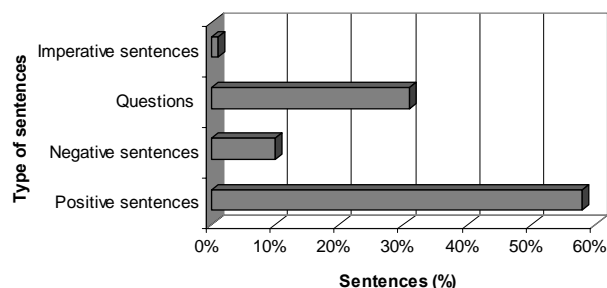


Fig. 6 The percentage of different type sentences in the corpus

The test sentences display 5 different word order errors

[16],[17]. The word order errors concern the transposition of Verbs, Nouns, Adjectives, Adverbs, and Pronouns, thus violating the sentences' word order constraints [18]. The most common errors are the Verb transposition with 35.0% and the adverb transpositions with 30.5% in total. The errors with adjectives transpositions present a lower percentage (19.9%). Noun transpositions are less frequent with 11.4%. The errors with Pronouns are least frequent with 3.4%.

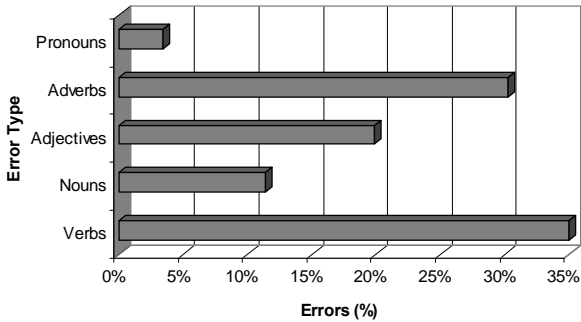


Fig. 7 Word Order errors distribution in TOEFL test set

The next figure shows the distribution of word order errors for each type of sentence. According to the above figure, the most frequent errors in the whole test set are the verb transpositions with 35.0% in total; this holds for all different types of sentences except from the category of the questions where the most frequent word order errors are the adverb transpositions. Regarding the imperative sentences it can be observed that there are no pronoun, noun and adverb transpositions.

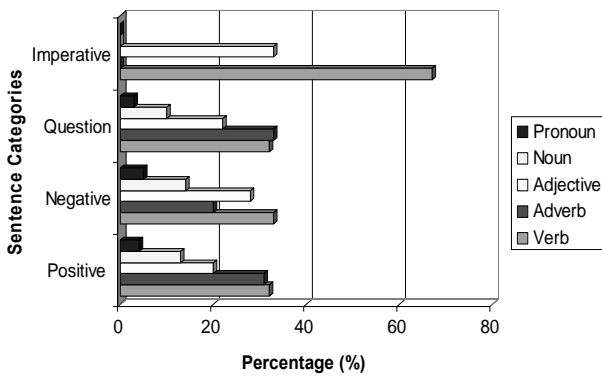


Fig. 8 Word Order errors distribution in different types of sentences

C. Experimental Results

Figure 7 shows the repairing results using the test sentences. This figure depicts the capability of the system to give as output the correct sentences in the 10-best list. The x-axis corresponds to the place of the correct sentence into this list. The last position (11) indicates that the correct sentence is out of this list.

It is obvious that the system's performance for detecting and repairing method of ill-formed sentences with word order

errors depends mainly on the quality of the corpus. The high success rate of the system is achieved using the grammatically and syntactically correct sentences of BNC.

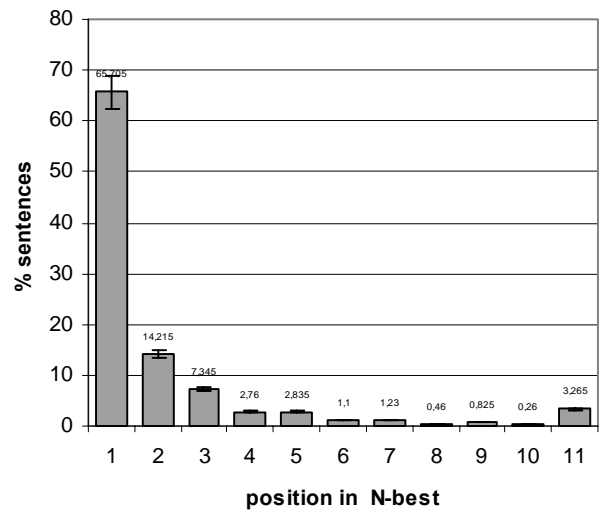


Fig. 9 The percentage of test sentences in different places into the N-best list (N=10)

The findings from the experimentation show that 96,735% of the test sentences have been repaired using the proposed method (True Corrections) (included into 10-best sentences). On the other hand, the result for 3,265% of the test sentences was false (False Corrections). In case of "False Corrections" the system's response does not include the correct sentence into the 10-best. The incorrect output of the system can be explained considering that some TOEFL words are not included into the BNC vocabulary, hence some of the sentences' trigrams are considered as invalid.

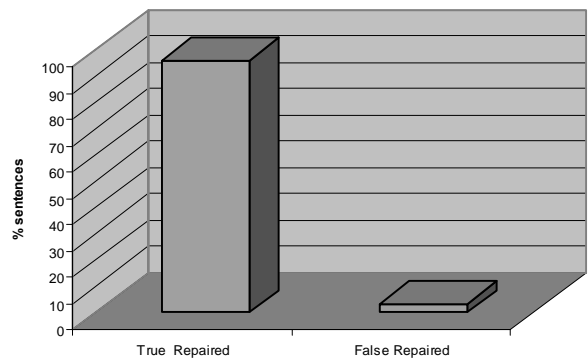


Fig. 10 The percentage of sentences with True and False corrections

D. Results using the confusion matrix technique

The number of permutations that are extracted with the filtering process is significantly lower than the corresponding value without filtering, especially for large sentences. For sentences with length up to 8 words, the number of permutations is slightly lower when the filtering process is

used, while for sentences with length greater than 8 words the filtering process provides a drastical reduction of permutations. It is obvious that the performance of filtering process depends mainly on the number of valid bigrams. This implies that the language model's reliability affects the outcome of the system and especially of the filtering process.

TABLE IV
 THE MEAN VALUE OF PERMUTATIONS FOR TOEFL SENTENCES

words	No filtering	With filtering
7	5040	768
8	40320	4627
9	362880	32451
10	3628800	246987
11	39916800	2167890
12	479001600	8790541

The next figure shows the impact of the filtering process on the permutations. In case of sentences with 12 words, the filtering enhances the performance of the proposed system and reduces the computational load.

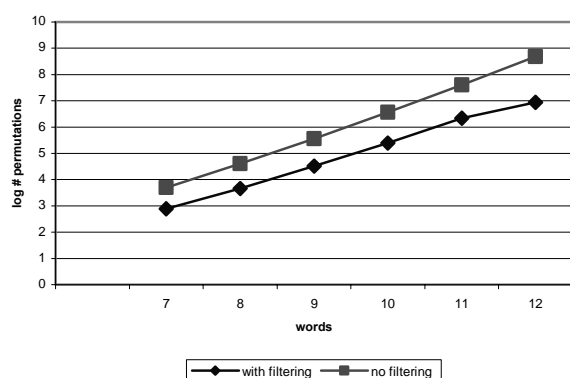


Fig. 11 The number of permutations with and without filtering for sentences with length from 5 to 12. The symbol (■) denotes the log number of the sentence's permutations without filtering while the symbol (◆) presents the log number of the permutations extracted from the filtering method

VII. CONCLUSIONS

Recognising and repairing sentences with word order errors is a challenge ready to be addressed. The proposed method is effective in repairing erroneous sentences. Therefore the method can be adopted by a grammar checker as a word order repairing tool. The necessity of the grammar checkers in educational purposes and e-learning is more than evident. Another aspect of the method's effects is the ability of using different text corpora to distinguish different writing styles. It is interesting that the system does not only detect errors as other approaches do but also repairs the ill-formed sentences.

The findings show that most of the sentences can be

repaired by this method independently from the sentence's length and the type of word order errors. By the permutation's filtering process, the system takes advantage of better performance, rapid response and smaller computational space.

One of the key questions is whether the use of other kinds of statistical language models (skipping, clustering) can improve the performance of the proposed system. The issue certainly invites research. Another issue that should be investigated is whether the language model in conjunction with the attributes of each word can give better results.

REFERENCES

- [1] J. A., Hawkins, A Performance Theory of Order and Constituency. Cambridge, Cambridge University Press, 1994.
- [2] D., Schneider, K.F., McCoy, Recognizing syntactic errors in the writing of second language learners, Proceedings of the 17th international conference on Computational linguistics, 1198-1204, 1998.
- [3] E.S., Atwell, How to detect grammatical errors in a text without parsing it. In Proceedings of the 3rd EAACL, 38-45, 1987.
- [4] A., Golding, A Bayesian hybrid for context-sensitive spelling correction. Proceedings of the 3rd Workshop on Very Large Corpora, 39-53. 1995
- [5] M.,Chodorow, C., Leacock. An unsupervised method for detecting grammatical errors. In Proceedings of NAACL'00, 140-147. 2000.
- [6] T. Heift, Designed Intelligence: A Language Teacher Model, Unpublished Ph.D. Dissertation, Simon Fraser University, 1998
- [7] T. Heift, Intelligent Language Tutoring Systems for Grammar Practice. Zeitschrift für Interkulturellen Fremdsprachenunterricht (Online), 6 (2), 15 pp. 2001
- [8] J., Bigert, O., Knutsson. Robust error detection: A hybrid approach combining unsupervised error detection and linguistic knowledge. In Proceedings of Robust Methods in Analysis of Natural language Data, (ROMAND 2002), 10-19, 2002.
- [9] J., Sjöbergh, Chunking: an unsupervised method to find errors in text, Proceedings of the 15th Nordic Conference of Computational Linguistics, NODALIDA 2005, 2005
- [10] S.J., Young., Large Vocabulary Continuous Speech Recognition, IEEE Signal Processing Magazine 13, (5), 45-57, 1996.
- [11] I.J., Good, The population frequencies of species and the estimation of population parameters. Biometrika, 40(3 and 4):237-264, 1953.
- [12] S.M., Katz, Estimation of probabilities from sparse data for the language model component of a speech recogniser. IEEE Transactions on Acoustics, Speech and Signal Processing, 35(3):400-401, 1987.
- [13] C. M., Feyton, Teaching ESL/EFL with the internet. Merrill Prentice-Hall, 2002.
- [14] K.S., Folse, Intermediate TOEFL Test Practices (rev. ed.). Ann Arbor, MI: The University of Michigan Press, 1997.
- [15] J. C., Park, M., Palmer, and G., Washburn, An English grammar checker as a writing aid for students of English as a second language, In Proceedings of Conference on Applied Natural Language Process, New Brunswick, NJ, 1997.
- [16] R., Murphy, Order of several describing words together (adjectives), English Grammar in Use Cambridge University Press, Cambridge, Unit 95, 1990.
- [17] J., Eastwood, Order of place, time and frequency words (never, often), Oxford Practice Grammar Oxford University Press, Oxford. Unit 89, 1997.
- [18] E., Izumi, K., Uchimoto, T., Saiga, T., Supnithi, H. Isahara, Automatic error detection in the Japanese learners English spoken data. In Companion Volume to the Proceedings of ACL '03, 145-148, 2003