**World Scientific**
www.worldscientific.com

# Machine Translation of English Noun Phrases into Arabic

KHALED SHAALAN

Computer Science Department,
Faculty of Computers and Information, Cairo Univ.,
5 Tharwat St., Orman, Giza, Egypt
*shaalan@mail.claes.sci.eg*


AHMED RAFEA

Computer Science Department, American University in Cairo,
113, Sharia Kasr El-Aini, P.O. Box 2511, 11511, Cairo, Egypt
*rafea@aucegypt.edu*


AZZA ABDEL MONEIM

Central Lab. for Agricultural Expert Systems (CLAES),
El-Nour St., Dokki, Giza, Egypt
*azza@mail.claes.sci.eg*


HODA BARAKA

Computer Engineering Department,
Faculty of Engineering, Cairo University, Dokki, Giza, Egypt
*Hbarka@idsc1.gov.eg*

The present work reports our attempt in automating the translation of English noun phrase (NP) into Arabic. Translating NP is a very important task toward sentence translation since NPs form the majority of textual content of the scientific and technical documents. The system is implemented in Prolog and the parser is written in DCG formalism. The paper also describes our experience with the developed MT system and reports results of its application on real titles of theses and journals from the computer science domain.

*Keywords*: Machine Translation; Transfer Approach; Noun Phrases; Arabic Language Processing.

# 1.    Introduction

Machine translation (MT) is the area of information technology and applied linguistics dealing with the translation of human languages such as English and Arabic. With globalization and expanding trade, demand for translation is set to grow. Computer technology has been applied to technical translation in order to improve one or both of the following factors (Trujillo, 1999): 1) **Speed**: Translation by or with the aid of machines can be faster than manual translation, and 2) **Cost**: Computer aids in translation can reduce the cost per word of a translation. In addition, the use of MT can result in improvements in quality, particularly in the use of consistent terminology within a text or for a particular kind of client

As English is a universal language, most of the researches in Arabic MT are mainly concentrated on the translation between English and Arabic. This will help in simplifying the Arab communication with other countries. These systems are based mainly on the transfer model. Ibrahim (1991) discussed the problem of the English-Arabic translation of the embedded idioms and proverb expressions in the English sentences. Rafea et al. (1992) developed an English-Arabic MT system, which translates a sentence from the domain of the political news of the Middle East. Maalej (1994) discussed the MT of English nominal compounds into Arabic. It has been motivated by their frequent occurrence in referring and naming in all text-types. Pease et al. (1996) developed a system, which translates medical texts from English to Arabic. El-Desouki et al. (1996) discussed the necessity of modular programming for English-Arabic MT. A translation of an English subset of a knowledge base, written in KROL (Shaalan et al., 1998), to the corresponding Arabic phrases is described in (El-Saka et al., 1999). Mokhtar (2000) developed an English-Arabic MT system, which is applied to abstracts from the field of Artificial Intelligence.

On the contrary, little work has been done in developing Arabic-English MT systems. Al Barhamtoshy (1995) proposes a translation method for compound verbs. Shaalan (2000) described a tool for translating the Arabic interrogative sentence into English. Chalabi (2001) presented an Arabic-English MT engine that allows any Arabic user to search and navigate through the Internet using the Arabic language. Othman et al. (2003) developed an efficient chart parser that will be used for translating Arabic sentence.

The present work addresses the translation of a fairly complex English NP into Arabic, which is an important task for automating the translation between English and Arabic sentences. Justeson et al. (1995) reported that English NPs form the majority of textual content of the scientific and technical documents. These have motivated us to start by translating English NP as a first step toward

sentences translation. We also found that the translation of the title of scientific texts can be done using this system.

The next section outlines the overall architecture of the proposed Arabic-English MT system. The following sections describe the main components of the system. We also describe how we evaluated the correctness of our MT system. In a concluding section, we discuss its application on real titles of theses from the computer science domain and present some final remarks.

## 2.  Overall Structure of the System

There are three basic approaches being used for developing MT systems that differ in their complexity and sophistication. These approaches are: direct approach, transfer-based approach, and interlingua approach. The current work follows the transfer-based MT approach. There are many factors which make transfer an attractive design for MT (Trujillo, 1999): 1) Many systems are bilingual, or their principal use is for translation in one direction between a limited number of languages; 2) Where full multilinguality is required, it is possible to have a hub language into and out of which all translation is done; and 3) Portions of transfer modules can be shared when closely related languages are involved. For example, an English-Portuguese module may share several transformations with an English-Spanish module. The architecture of the transfer-based English to Arabic MT system is given in Fig. 1, with three main components: an analyzer component, a transfer component, and a generation component.

Prolog is one of the most widely used programming languages in computational linguistics. We have chosen Prolog as the implementation language of our translation tool. Among the features that make it attractive are its efficient unification, its declarative nature and its backtracking regime. The tool is implemented in SICStus Prolog and the parser is written in DCG formalism. DCG translate grammar rules directly to Prolog, producing a simple top-down parser.

## 3.  Syntax Analysis

The development of the parser is a two-step process. In the first step, we acquire the rules that constitute a grammar for the English NP that gives a precise account of what it is for a NP to be grammatically correct. The grammar covers the titles of theses from the computer science domain. The grammar of NPs is acquired from the analysis of 50 titles. Our analysis indicates that NPs either occur in a simple form or in a complex form. The complex form of a NP is two
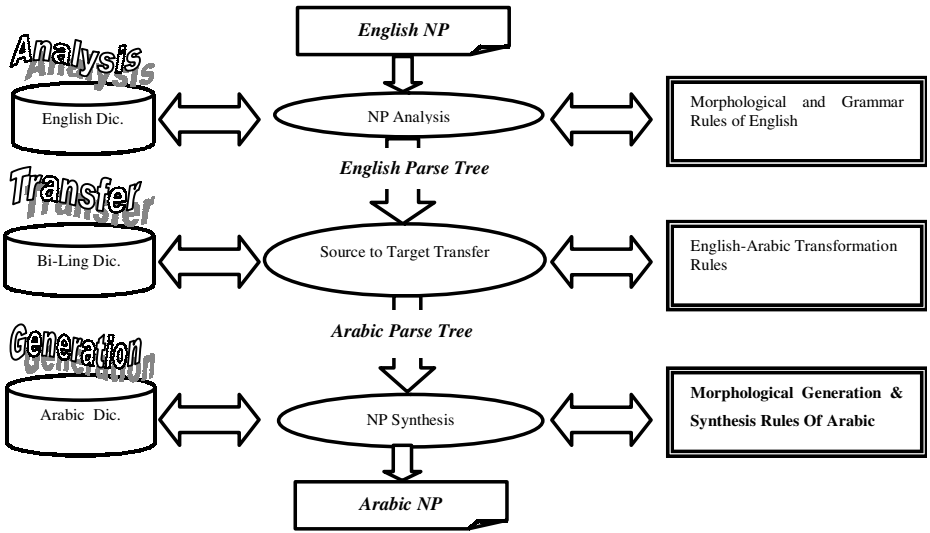
Figure 1.  Overall structure of English-Arabic noun phrase translator.

or more simple NPs separated by: a connector, a preposition, or a special word — collectively called *separator*. Special words are the special symbol ":", colon, or a word that can be used to recognize the beginning of a new NP. The second step is to implement the parser that assigns grammatical structure onto input NP.

In order to implement the parser, it was needed to perform morphological analysis on the inflected English words. An English monolingual dictionary was also needed to successfully implement the morphological analyzer. The morphological analyzer returns to the parser the words in its primitive form with some additional information such as the number of a noun. Entries of the English dictionary can be stems of nouns, adjectives, quantifiers and separators. Considering nouns, they are stored in the dictionary in their singular form. In addition, the dictionary also includes entries for irregular plurals.

## 4.  Syntactic Transfer MT

Syntactic transfer systems rely on mappings between the surface structure of sentences: a collection of tree-to-tree transformations is applied recursively to the analysis tree of the source language in order to construct a target language analysis tree (Arnold, 1993 & Arnold et al., 1994). The tree-to-tree transformation algorithm is a recursive, non-deterministic, top-down process in which one side of the tree-to-tree transfer rules is matched against the input structure, resulting in the structure on the right-hand-side. The transfer component has a bilingual dictionary to perform English to Arabic translation.

In our NP translator, the actual translation occurs in the transfer phase. The following paragraphs describe the problems encountered while designing this phase. These problems are regarded as peculiar to translation, since they arise from the divergences and mismatches between source and target NPs.

## 4.1. Lexical transfer

The lexical transfer rules relate every word in the source sentence representation to some corresponding target language representation, e.g.

$$\text{networks} \Leftrightarrow \text{شبكة}$$
$$\text{performance} \Leftrightarrow \text{اداء}$$
$$\text{evaluation} \Leftrightarrow \text{تقييم}$$

The features associated with the English word node must also be translated in some way. These features are pairs that consist of an attribute, such as a number, and a value, such as singular. The rules relevant to the features are straightforward, indicating that the given values are simply carried over from source language representation to the target language representation, e.g.

$$\{\text{number} = \text{sg}\} \Leftrightarrow \{\text{number} = \text{sg}\}$$

These dictionary rules can be seen as relating leaves (the word nodes) on the English parse tree to leaves on the target Arabic tree.

## 4.2. Structural transfer

The structural rules relate other parts and nodes of the two trees to each other. There is a relationship between adjacent lexical units in the English NP. This concerns the preserved order, the relative positioning (precedence), of lexical units in the NP. For example, adjective in English precedes its noun, like in (good man), while in Arabic, noun precedes the adjective. Consequently, restructuring of the English parse tree is needed to conform to the target Arabic grammar. The transfer rules described here deal with the restructuring of the parse tree and reordering of words.

In the structural rules that follow, the LHS describes an English structure and the RHS describes the Arabic one, and $1, $2, …$k are variables interpreted as standing for pieces of English structure on one side, and for their translations on the other side. Arnold (1994) introduced the method of defining the structure transfer rules described here.

## SIMPLE TRANSFER

The main difference between the English and Arabic parse tree representation is that words are in reverse order. A relatively simple straightforward example where a more complex example is called for involves the translation of the NP "networks performance evaluation" into "تقييم أداء شبكة", which shows the switching of words. Such a rule might look as follows:

(1)     $[w_i:\$1, w_{i+1}:\$2, \ldots, w_k:\$k] \Leftrightarrow$        $(1 \leq i \leq k)$
        $[w_k:\$k, w_{k-1}:\$k-1, \ldots, w_i:\$i]$       $(1 \leq i \leq k)$

This rule says that the translation of the word at level $i$ is switched with the word at level $k - i + 1$, where $k$ is the number of NPs equivalent to maximum (sub)tree level, see Fig. 2. This is done in Prolog as follows:
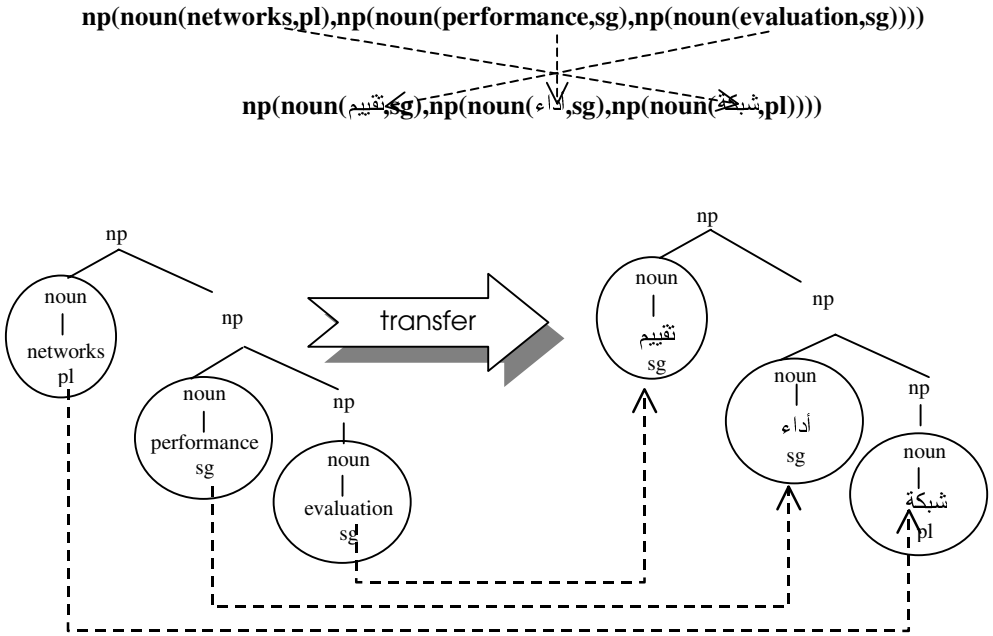
np(noun(networks,pl),np(noun(performance,sg),np(noun(evaluation,sg))))

np(noun(تقييم,sg),np(noun(أداء,sg),np(noun(شبكة,pl))))



Figure 2. Simple transfer.

## COMPOUND TRANSFER

The translation rule of a compound NP is as follows.

(2)     $[NP:\$1, Sept:\$2, NP:\$3] \Leftrightarrow$
        $[NP:\$1, Sept:\$2, NP:\$3]$

As an example, consider the translation of the NP "intelligent search system for bibliographic services" into "بحث ذكي لخدمة بيبلوجرافي". The rule associates variable $1 with the subtree of "intelligent search system", $2 with the node for "for", and $3 with the subtree for "bibliographic services". Translating each of these then becomes a separate task for transfer, which operates on these subtrees in the same way as on the original tree — attempting to find rules which deal with these sorts of structure, see Fig. 3.

WORD OMISSION

In the translation process, lexical mapping is inevitable. In any transfer process, the one-to-zero mapping is undesirable. One-to-zero mapping means the lack of equivalent lexical word in the target language (Trujillo, 1999). In this work we found this gap when we translate a NP that contains an "of" separator. This translation is described in two steps: restructuring of the English parse tree and reordering of words. The translation rule of a NP that contains the separator "of" is as follows.

(3)    [NP:$1, Sept:of, NP:$2] ⇨ [NP:$3[NP:$2, NP:$1]] ⇔
       [NP:$3]

As an example, consider the translation of the NP "performance evaluation of routing algorithms" into "تقييم أداء خوارزمية مسارات". The rule associates variable $1 with the subtree of "performance evaluation" and $2 with the subtree for "routing algorithms". In the first step, the original tree is restructured by dropping the "of" node and switching the arguments of the "of". This yields a tree representation of "routing algorithms performance evaluation" that is to be associated with $3. In the second step, the normal transfer rules are applied to the tree representation associated with $3 to get the translation of the original tree, see Fig. 4.

ORDER OF THE HEAD WORD

The translation of a NP starting in a quantifier is described by a special rule that is as follows.

(4)    [wi:$1[quantifier form], wi+1:$2, …, wk:$k] ⇔          $(1 \le i \le k)$
       [wi:$1, wk:$k, wk−1:$k−1, …, wi:$i]                    $(1 \le i \le k)$

This rule says that the quantifier that fills the head of the NP is translated to an equivalent word that fills the head, i.e. at the same position, of target NP. The rest of the NP is translated in the normal transfer rules. As an example, consider the
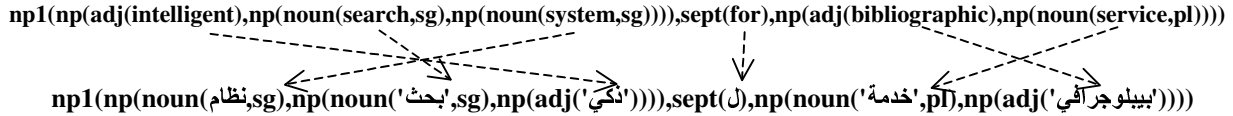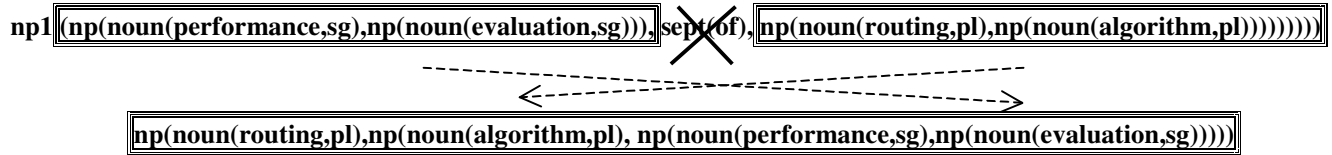
np1(np(adj(intelligent),np(noun(search,sg),np(noun(system,sg)))),sept(for),np(adj(bibliographic),np(noun(service,pl))))

np1(np(noun(نظام,sg),np(noun('بحث',sg),np(adj('ذكي')))),sept(ل),np(noun('خدمة',pl,np(adj('بيبلوجرافي'))))

Figure 3.  Compound transfer.

**Step 1**

np1 (np(noun(performance,sg),np(noun(evaluation,sg))), sept(of), np(noun(routing,pl),np(noun(algorithm,pl)))))))))

np(noun(routing,pl),np(noun(algorithm,pl), np(noun(performance,sg),np(noun(evaluation,sg)))))

**Step 2**

np(noun(مسارات,pl),np(noun(خوارزمية,pl), np(noun(أداء,sg),np(noun(تقييم,sg)))))

np(noun(تقييم,sg),np(noun(أداء,sg),np(noun(خوارزمية,pl),np(noun(مسارات,pl))))))
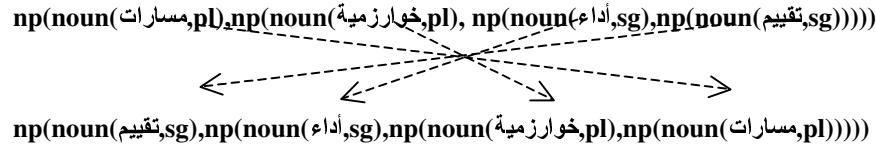
Figure 4.  "of" transfer.

translation of the NP "comparative study for some software engineering tools" into "دراسة مقارنة لبعض أداة هندسه برمجيات", see Fig. 5.

**np1(np(noun(comparative,sg),np(noun(study,sg))),sept(for),**
**np(quantifier(some),np(noun(software,pl),np(noun(engineering,sg),np(noun(tool,pl)))))**
**np1(np(noun(دراسة,sg),np(noun(مقارنة,sg))),sept(ل),**
**np(quantifier(بعض),np(noun(أداة,pl),np(noun(هندسه,sg),np(noun(برمجيات,pl)))))**

Figure 5  Quantifier transfer.

The translation of a NP starting in a gerund form is described by a special rule that is as follows.

| | | |
|---|---|---|
| (5) | [wi:$1[-ing form], wi+1:$2, …, wk:$k] $\Leftrightarrow$ | $(1 \leq i \leq k)$ & $1 \notin$ termbase |
| | [wi:$1, wk:$k, wk−1:$k−1, …, wi:$i] | $(1 \leq i \leq k)$ |

This rule says that the gerund form that fills the head of the NP is translated to an equivalent word that fills the head, i.e. at the same word order, of target NP. The rest of the NP is translated in the normal transfer rules. As an example, consider the translation of "automating software testing" into "أتمتة اختبار برمجيات", see Fig. 6.

**np(noun(automating,sg),np(noun(software,pl),np(noun(testing,sg))))**

**np(noun(أتمتة,sg),np(noun(اختبار,sg),np(noun(برمجيات,pl))))**

Figure 6.  Gerund transfer.

One further condition is added in order to apply this rule. The gerund word that fills the head of the NP must not be a computer term. As an example, consider the translation of the NP "switching techniques for multimedia information transmission over a wide-area network" into "الوسائط المتعددة علي شبكة تسعة تقنيات تحويل لنقل معلومات". Since "switching" is a gerund computer term, the above rule cannot be applied. Special words like the one given have to be kept in a termbase in order to prevent this rule from being applied.

The translation of a NP starting in a word ending in (able) form is described by a special rule that is as follows.

| | | |
|---|---|---|
| (6) | [wi:$1[-able form], wi+1:$2, …, wk:$k] $\Leftrightarrow$ | $(1 \leq i \leq k)$ |
| | [wi:$1, wk:$k, wk−1:$k−1, …, wi:$i] | $(1 \leq i \leq k)$ |

As an example, consider the translation of "reusable problem solving components" into "اعادة أستخدام مكون حل مشكلة", see Fig. 7.

np(adj(reusable,sg),np(noun(problem,sg),np(noun(solving,sg),np(noun(component,pl)))))

np(noun(اعادة أستخدام,sg),np(noun(مكون,pl),np(noun(حل,sg),np(noun(مشكلة,sg)))))
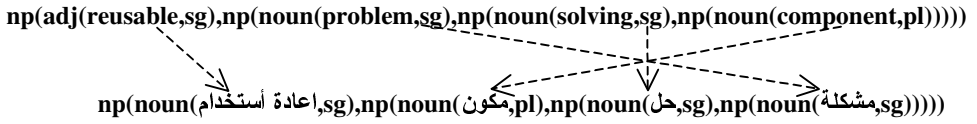
Fig. 7 "-able" transfer.

## 5. Syntactic Generation

The generation component comprises Arabic dictionary, Arabic morphological synthesizer, and Arabic NP synthesizer. The Arabic morphological synthesizer generates the inflected Arabic word according to the Arabic agreement relationship. This relationship is between words in certain context such that a word in one position follows the word in a corresponding position in some aspects: such as number (single, plural), sex (male, female), and definition (definite, indefinite). The role of this module is to synthesize: defined noun, defined adjective, feminine adjective, and plural noun (sound masculine plural, sound feminine plural, and irregular plural).

The NP synthesis consists of three steps. First, the right nouns are built according to the number and definition features. Second, we ensure the agreement relationship between descriptive nouns and adjectives with regard to the sex and definition features. Third, the transformed tree is traversed in a depth-first manner to produce the surface Arabic NP.

## 6. Evaluation of the MT System

We developed an evaluation methodology to compare the original manual translation of 66 new real thesis titles from the computer science domain, different from those 50 titles used to develop the MT system, and the MT system output of these 66 titles. These title contain 156 simple NP. The following steps describe this methodology:

1. Compare the original translation with the system output taking into consideration that one title may contain more than one simple NP. (A human translator was involved in this activity.)
2. Classify the problems that arise from the mismatches between the two translations.

3  Assign a suitable score (0 – 10) for each simple NP in the title translation based on the identified problem and its deviation from the human translation.

4. If one NP in the translation contains more than one problem the score is calculated as the product of the score due to each individual problem divided by $10^{n-1}$; where $n$ is the number of problems.

5. Determine the overall correctness score by summing up the scores of all NPs and compute the percentage.

Applying this evaluation methodology has resulted in giving 47 matched NP translations. The remaining 109 NPs have problems. The classification of these problems sorted according to their seriousness takes into consideration the number of occurrences and deviation from the correct translation as follows:

1. The synonym of a word is used, e.g the English noun "acquisition" may be translated into Arabic to either "استخلاص" or "اكتساب". The translated title is not affected due to this problem and therefore a NP having this problem is given a score of 10. The total number of occurrences of this problem is 67.

2. In the English annexation "Arabic text retrieval and classification", the postfixed NP "Arabic text" comes before the annexed structure "retrieval and classification". The MT system considers the word "retrieval" as the annexed noun without the noun "classification" which is considered as part of the NP coming after the connector "and". This is a difficult problem and for this reason we give score 7 to NPs having this problem. The total number of occurrences of this problem is 1.

3. In a NP that is formed from annexation structure, the human translators prefer sometimes to change the position of the adjective to follow the annexed noun directly and add an extra preposition to connect this described noun with the rest of the annexation. For example, "Investigating a new constraint solving approach" is translated into "بحث طريقة جديدة لـحل القيود". Although the automatic translation is acceptable as a possible translation we give score 8 to NPs having this problem just because it does not exactly match the human translation. The total number of occurrences of this problem is 4.

4. In a NP that is formed from annexation structure the annexed noun and postfixed noun are sometimes different in gender and number. When an adjective is given to this annexation, the automatic translation gives to the adjective the gender of the last annexed noun in the annexation. This is not necessarily the gender of the annexed noun. The correct translation is to let the adjective gender be the same as the postfixed noun. For example, "intelligent tutoring systems" should be translated to "نظم التعليم الذكية" instead of "نظم التعليم الذكي" because "intelligent" refers to "systems" which

is feminine and not to "tutoring" which is masculine. We give score 8 to NPs having this problem. The total number of occurrences of this problem is 4.

5. A scientific term or an (ing) form of a noun is not usually defined in English, e.g. "constraint logic programming" or "controlling". We noticed that human translators tend to add the definition article to these terms. Although the automatic translation is acceptable as a possible translation, we give score 8 to NPs having this problem just because it is not exactly matching the human translation. The total number of occurrences of this problem is 8.

6. A NP representing a scientific term such as "data mining", decision support system", etc., has special translation to Arabic because literal translation is not acceptable. We give score 8 to NPs having this problem since it is not a serious error in the MT system. The total number of occurrences of this problem is 9.

7. In most of the cases there is no need to translate the preposition "of" to an equivalent Arabic preposition. Sometimes, this is not correct, e.g. when translating "automatic generation of explanation", we have to add "ل" to our translation to be correct. We noticed that in Arabic if the adjective follows directly the annexed noun in the annexation NP, we have to add "ل" to connect this described annexed noun to the reset of the annexation. We give score 6 to this problem because we cannot generalize any finding at the moment and more research is needed to decide on the translation of the preposition "of". The total number of occurrences of this problem is 6.

8. The translation of a preposition is different. We give score 7 to NPs having this problem since it is not an easy problem to solve. The total number of occurrences of this problem is 17.

The evaluation results show that the 156 NPs take 1439 total score, which means that about 92% of the translations were correct. Table 1 describes the distribution of the 156 NP's based on their scores.

Table 1.  Evaluation results.

| Score range | No. of NPs | % |
|---|---|---|
| 10 | 114 | 73 |
| [8–10) | 21 | 13 |
| [6–8) | 18 | 12 |
| [4–6) | 3 | 2 |
| Total | 156 | 100 |

## 7.  Conclusion

This paper has concentrated on the issues in the design and implementation of a transfer-based MT system, which translates a fairly complex English NP into Arabic. We showed that the MT approach is promising and can be used to automate the translation of the titles of the thesis in the computer science domain.

We have collected 116 real titles of refereed thesis from the computer science domain. These were available to us each of which has its Arabic translation. We took 50 titles out of the 116 titles as a training sample. This helped us in acquiring the rules for analysis, transfer, and generation components of the MT system. The remaining 66 thesis titles are used for evaluating the approach and the correctness of the MT system. The evaluation is based on comparing the system output with the human translations. The 66 titles include 156 simple NP. In the case where there was a difference between the human translation and the machine translation it concerned only a fragment of the entire NP. These means that the automatic translation was partially correct.

The problems found are classified, explained and assigned a score. The overall evaluation results, according to the presented methodology, were satisfactory. It showed that the 156 NPs took 1439 total score, which means that about 92% of the translations were correct. It is also worthy to mention that 86% of the 66 thesis titles got a score in the range 8 to 10

In future works, the problems mentioned can be improved by storing in the lexicon synonyms of words, implementing the transfer rules that deals with the special cases for mapping the English parse tree into the Arabic parse tree, and implementing the synthesis rules that introduces new words in the output translations for some special cases. These improvements will raise the correctness of the translations from 92% into 100%.

## References

[1] A. Al Barhamtoshy, Arabic to English Translator of Compound Verbs, *in Proceedings of the Annual Conference on Computer Science & Statistical ISSR*, December 1995.

[2] D.J. Arnold, *Sur la Conception du Transfert, La Traductique*, P. Bouillon and A. Clas, (Eds.), Les Presses de l'universite de Montreal, Montreal, 1993, pp 64–76. A version of this-in English "*Transfer*" — can be found in http://clwww.essex.ac.uk/~doug/transfer.ps.gz.

[3] D.J. Arnold, Lorna Balkan, Siety Meijer, R. Lee Humphreys and Louisa Sadler, *Machine Translation: An Introductory Guide*, Blackwells-NCC, London, 1994, available at http://clwww.essex.ac.uk/~doug/book/book.html.

[4]  A.S. Chalabi, Web-based Arabic-English MT engine, *in Proceedings of the Arabic NLP Workshop at ACL/EACL,* 2001.

[5]  A. El-Desouki, A. Abd Elgawwad and M. Saleh, A Proposed Algorithm for English-Arabic Machine Translation System, *in Proceedings of the 1st KFUPM Workshop on Information and Computer Sciences (WICS): Machine Translation*, Dhahran, Saudi Arabic, 1996.

[6]  T. El-Saka, A. Rafea, M. Rafea and M. Madkour, English to Arabic Knowledge Base Translation Tool, *in Procceding of the 7th International Conference on Artificial Intelligence Applications*, Cairo, Feb., 1999.

[7]  M. Ibrahim, *A Fast and Expert Machine Translation System Involving Arabic Language*, Ph. D. Thesis, Cranfield Institute of Technology, U.K., 1991.

[8]  J. Justeson and S. Katz, Technical Terminology: Some Linguistic Properties and an Algorithm for Identification in Text, in *Natural Language Engineering*, Vol. 1, No. 1, 9–27, 1995.

[9]  Z. Maalej, English-Arabic Machine Translation of Nominal Compounds, *in Proceedings of the Workshop on Compound Nouns: Multilingual Aspects of Nominal Composition.* Geneva: ISSCO, pp. 135–146, 1994.

[10] H. Mokhtar, *An Automatic System for English-Arabic Translation of Scientific Text (SEATS)*, Master thesis, Computer Engineering Department Faculty of Engineering, Cairo University, 2000.

[11] E. Othman, K. Shaalan and A. Rafea, A Chart Parser for Analyzing Modern Standard Arabic Sentence, *in Proceedings of the MT Summit IX Workshop on Machine Translation for Semitic Languages: Issues and Approaches*, New Orleans, Louisiana, U.S.A., 2003.

[12] C. Pease and A. Boushaba, Towards an Automatic Translation of Medical Terminology and Texts into Arabic, *in Proceedings of the Translation in the Arab World*, King Fahd Advanced School of Translation, Nov. 27–-30, 1996.

[13] A. Rafea, M. Sabry, R. El-Ansary, S. Samir, Al-Mutargem, A Machine Translator for Middle East News, *in Proceedings of the 3rd International Conference and Exhibition on Multi-lingual Computing*, December 1992.

[14] K. Shaalan, M. Rafea and A. Rafea, KROL: A Knowledge Representation Object Language on Top of Prolog, *Expert Systems with Applications: An International Journal*, Vol. 15, 33–46, 1998.

[15] K. Shaalan, Machine Translation of Arabic Interrogative Sentence into English, *in Proceedings of the 8th International Conference on Artificial Intelligence Applications*, Egyptian Computer Society (EGS), Egypt, 2000, pp. 473-483

[16] A. Trujillo, *Translation Engines: Techniques for Machine Translation*, Springer Verlag, 1999.