

Theoretical Overview of Machine translation

Mohamed Amine Chéragui¹

¹ African University, Adrar, Algeria,
m_cheragui@univadrar.org

Abstract. The demand for language translation has greatly increased in recent times due to increasing cross-regional communication and the need for information exchange. Most material needs to be translated, including scientific and technical documentation, instruction manuals, legal documents, textbooks, publicity leaflets, newspaper reports etc. Some of this work is challenging and difficult but mostly it is tedious and repetitive and requires consistency and accuracy. It is becoming difficult for professional translators to meet the increasing demands of translation. In such a situation the machine translation can be used as a substitute.

This paper offers a brief but condensed overview of Machine Translation (MT). Through the following points: History of MT, Architectures of MT, Types of MT, and evaluation of M T.

Keywords: History of MT, Architecture of MT, Types of MT, evaluation of MT.

1 Introduction

After 65 years, this field is one of the oldest applications of computers. Over the years, Machine Translation has been a focus of investigations by linguists, psychologists, philosophers, computer scientists and engineers. It will not be an exaggeration to state that early work on MT contributed very significantly to the development of such fields as computational linguistics, artificial intelligence and application-oriented natural language processing.

Machine translation, commonly known as MT, can be defined as “translation from one natural language (source language (SL)) to another language (target language (TL)) using computerized systems and, with or without human assistance”[1] [2].

We try to give in this paper a coherent, if necessarily brief and incomplete, the development has been the field of machine translation through four points which are: first of all surveys the chronological development of machine translation, the different approaches developed (linguistic and computational), the types of machine translation and finally, we try to answer an important question which is how to evaluate a machine translation?

2 History of Machine Translation

Although we may trace the origins of machine translation (MT) back to seventeenth century ideas of universal (and philosophical) languages and of 'mechanical' dictionaries, it was not until the twentieth century that the first practical suggestions could be made. The history of machine translation can be divided into five (05) periods [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12] :

2.1 First period (1948-1960): The beginning.

- 1949 : Warren Weaver in his Memorandum of 1949 proposed the first ideas on the use of computers in translation, by adopting the term computer translation.
- 1952 : The first symposium of machine translation, entitled Conference on Machine Translation, held in July 1952 at MIT under leadership of Yehoshua Bar-Hillel.
- 1954 : The development of the first automatic translator (very basic) by a group of researchers from Georgetown University in collaboration with IBM, which translates into more than sixty (60) Russian sentences into English. The authors claimed that within three to five years, machine translation would not be a problem.
- 1954 : Victor Yngve published the first journal on MT, entitled « Mechanical translation devoted to the translation of languages by the aid of machines ».

2.2 Second Period (1960-1966) Parsing and disillusionment

- Early 1960s This parsing is put forward as the only possible avenue of research to advance the machine translation. Thus there are already many parsers developed from different types of grammars, such as grammar and dependency grammar Tesnière stratificationnelle Lamb
- 1961 : In February of this year that computational linguistics is born, thanks to weekly lectures organized by David G. Hays at the Rand Corporation in Los Angeles. These conferences will be included as papers at the First International Conference on Machine Translation of Languages and Applied Language Analysis of Teddington in September 1961 with the participation of linguists and computer scientists involved in the translation as: Paul Garvin, Sydney M. Lamb, Kenneth E. Harper, Charles Hockett, Martin Kay and Bernard Vauquois.
- 1964 : the creation of committee ALPAC(Automatic Language Processing Advisory Committee) with American government to studies the perspectives and the chances of machine translation
- 1966 : ALPAC published his famous rapport in which it concluded that its works on machine translation is just wasting of time and money ; the conclusion of this rapport is it had a negative impact on their search (MT) for a number of years

2.3 Third period (1966-1980): New birth and hope

- 1970 : Start of the project REVERSO by a group of Russian researchers.
- 1970 : Development of System SYSTRAN1 (Russian-English) by Peter Toma, who was at that time a member of a group search for Georgetown.
- 1976 : Creation of system WEATHER in the project TAUM (machine translation in the university of Montreal) under the direction of Alai Colmerauer for the machine translation weather forecasts for the general public, this system was created by group of researchers
- 1978 : Creation of system ATLAS2 by the Japanese firm FUJITSU, this translator was based on rules also he is able to translate from Korean to Japanese and vice versa

2.4 Fourth Period (1980-1990): Japanese invaders

- 1982 : The Japanese firm SHARP markets its Automatic translator DUET (English - Japanese), this translator was based on rules an approach to translation transfer
- 1983: as computer giant, NEC develops it's own system of translation based on algorithm called PIVOT. Marketed under the name of Honyaku Adaptor II, the version public the system of translation of NEC is also based on the method of pivot, by using Interlingua.
- 1986: Development of system PENSEE by OKI3, which is a translator (Japanese-English) based on rules.
- 1986: The group Hitachi developed his own translation system based on rules (which is an approach taken by transfer), christened on HICATS (Hitachi Computer Aided Translation System / Japanese- English).

2.5 Fifth Period (since 1990): the Web and the new vague of translators

- 1993: The project C-STAR (Consortium for Speech Translation Advanced Research) is an international cooperation. The theme of project is the machine translation of the parole in the field of tourism (dialogue client travel agent), by videoconference. these project birth the system C-STAR I which dealt three (03) languages (English, German et Japanese) and made the first demonstrations transatlantic trilingual in January 1993
- 1998: Marketing the translator REVERSO by the company Softissimo.
- 2000: the Development of system ALPH by Japanese laboratory ATR, this translator (Japanese-English and Chinese - English) takes an approach based on examples.

¹ The same translator was adopted by the European commission 1976 for the translation (Japanese-English)

² Currently we are in version 14 of the translator.

³ OKI : founded in 1881 Oki Electric Industry Co, is a Japanese manufacturer of telecommunications

- 2005: The appearance of the first web site for automatic translation ,like Google (<http://translate.google.fr/>).
- 2007: METIS-II is a hybrid machine translation system, in which insights from Statistical, Example based, and Rule-based Machine Translation (SMT, EBMT, and RBMT respectively) are used.
- 2008 : 23% of internet users, have used the machine translation and 40 % considering doing so
- 2009: 30% the professionals have used the machine translation and 18% perform a proofreading.
- 2010: 28% of internet users, have used the machine translation and 50% planning to do.

3 Architectures of machine translation systems

Different strategies have been adopted by different researchers at different times in the history of machine translation. The choice of strategy reflects one side of the depth and linguistic diversity but also the grandeur of ambition on the other side. There are generally two types of architecture for machine translation, which are:

3.1 Linguistic Architecture

In the linguistic architecture there are three basic approaches being used for developing MT systems that differ in their complexity and sophistication. These approaches are:

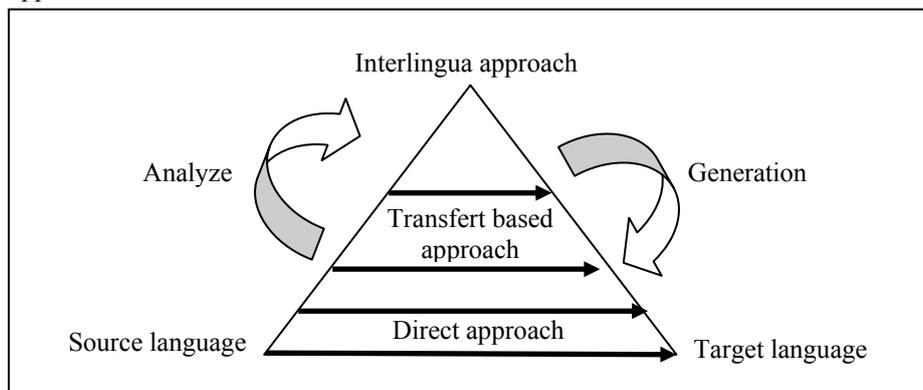


Fig1. The Vauquois triangle

- **Direct approach:** In direct translation, translation is direct from the source text to the target text. The vocabularies of SL texts are analyzed as needed for the resolution of SL ambiguities, for the correct identification of TL expressions as well as for the specification of word order in TL. This approach involves taking a string of words from the source language, removing the morphological

inflection from words to obtain the base forms, and looking them up in a bilingual dictionary between the source and the target languages. Components of this system are a large bilingual dictionary and a program for lexically and morphologically analyzing and generating texts [13].

- **Transfer-based approach:** In the Transfer approach, translation is completed through three stages: the first stage consists in converting SL texts into an intermediate representation, usually parse trees; the second stage converting these representations into equivalent ones in the target language; and the third one is the generation of the final target text [13].

In the transfer approach, the source text is analyzed into an abstract representation that still has many of the characteristics of the source, but not the target, language. This representation can range from purely syntactic to highly semantic. In the syntactic transfer, some type of tree manipulation into a target language tree converts the parse tree of the source input. This can be guided by associating feature structures with the tree. Whatever representation is used, transfer to the target language is done using rules that map the source language structures into their target language equivalents. Then in the generation stage, the mapped target structure is altered as required by the constraints of the target language and the final translation is produced.

- **Interlingua approach:** The Interlingua approach is the most suitable approach for multilingual systems. It has two stages: Analysis (from SL to the Interlingua) and Generation (from the Interlingua to the TL). In the analysis phase, a sentence in the source language is analyzed and then its semantic content is extracted and represented in the Interlingua form representation, where an Interlingua is an entirely new language that is independent of any source or target language and is designed to be used as an intermediary internal representation of the source text. The analysis phase is followed by the generation of the target sentences from the Interlingua representation. An analysis program for a specific SL can be used for more than one TL since it is SL-specific and not oriented to any particular TL. Furthermore, the generation program for a particular TL can be used again for translation from every SL to this particular TL since it is TL-specific and not designed for input from a particular SL [13].

3.2 Computational Architecture

- **Rule Based approach:** rule-based MT has two approaches: Interlingua and transfer. Rule-Based MT Systems rely on different levels of linguistic rules for translation. This MT research paradigm has been named rule-based MT due to the use of linguistic rules of diverse natures. For instance, rules are used for lexical transfer, morphology, syntactic analysis, syntactic generation, etc. In RBMT the translation process consists of:
 - Analyzing input text morphologically, syntactically and semantically.
 - Generating text via structural conversions based on internal structures.

The steps mentioned above make use of a dictionary and a grammar, which must be developed by linguists. This requirement is the main problem of RBMT as it is a time-consuming process to collect and spell out this knowledge, frequently referred as knowledge acquisition problem. It is not just very hard to develop and maintain the rules in this type of system, but one is not guaranteed to get the system to operate as well as before the addition of a new rule. RBMT systems are large-scale rule based systems; whereas their computational cost is high, since they must implement all aspects whether syntactic, semantic, structural transfer etc. as rules [14].

- **Corpus-based approach:** Corpus-Based Machine Translation, also referred as data driven machine translation, is an alternative approach for machine translation to overcome the knowledge acquisition problem of rule-based machine translation. There are two types of CBMT Statistical Machine Translation (SMT) and Example-Based Machine Translation (EBMT). Corpus-based MT automatically acquires the translation knowledge or models from bilingual corpora. Since this approach has been designed to work on large sizes of data, it has been named Corpus-Based MT ([17], [18], [16] and [15]).
- **Hybride approach:** Some recent work has focused on hybrid approaches that combine the transfer approach with one of the corpus-based approaches. This was designed to work with fewer amounts of resources and depend on the learning and training of transfer rules. The main idea in this approach is to automatically learn syntactic transfer rules from limited amounts of word-aligned data. This data contains all the needed information for parsing, transfer, and generation of the sentences ([19] and [20]). The following section covers part of the MT literature that gives details of specific systems for deriving the appropriate translation using different approaches.

4 Types of Machine Translation

4.1 Machine Translation for Watcher (MT-W)

This is intended for readers who wanted to gain access to some information written in foreign language who are also prepared to accept possible bad translation rather than nothing. This was the type of MT envisaged by the pioneers. This came in with the need to translate military technological documents. This was almost the dictionary-based translation far away from linguistic based machine translation [25].

4.2 Machine Translation for Revisers (MT-R)

This type aims at producing raw translation automatically with a quality comparable to that of the first drafts produced by human. The translation output can be considered only as brush-up so that the professional translator freed from that very boring and time consuming task can be promoted to revisers [25].

4.3 Machine Translation for Translators (MT-T)

This aims at helping human translators do their job by providing on-line dictionaries, thesaurus and translation memory. This type of machine translation system is usually incorporated into the translation work stations and the PC based translation tools. “Tools for individual translators have been available since the beginning of office automation.” And those systems running on standard platforms and integrated with several text processors are the ones that attained operational and commercial success [25].

4.4 Machine Translation for Authors (MT-A)

This aims at authors wanting to have their texts translated into one or several languages and accepting to write under control of the system or to help the system disambiguate the utterance so that satisfactory translation can be obtained without any revision. This is an “interactive MT, The interaction was however done both during analysis and during transfer, and not by authors, but by specialists of the system and language(s).” In short, there have been no operational successes yet in MT-A, but the designs are becoming increasingly user oriented and geared towards the right kind of potential users, people users, people needing to produce translations, preferably into several languages [25].

5 Evaluation of Machine Translation Systems

Evaluating Machine translation system is important not only for its potential users and buyers, also to researchers and developers. Various types of evaluation have been developed, such as :

5.1 BLEU (BiLingual Evaluation Understudy)

The BLEU metric, proposed by Papineni in 2001 was the first automatic measurement accepted as a reference for the evaluation of translations. The principle of this method is to calculate the degree of similarity between candidate (machine) translation and one or more reference translations based on the particular n-gram precision. The BLEU score is defined by the following formula [21]:

$$\text{BLEU} = \text{BP} \times e^{\left(\sum_{n=1}^N w_n \log p_n\right)} \quad (1)$$

Where:

- “ p_n ”: the number of n-grams of machine translation is also present in one or more reference translation, divided by the number of total n-grams of machine translation.
- “ w_n ”: positive weights.

- “BP”: Brevity Penalty, which penalizes translations for being “too short”. The brevity penalty is computed over the entire corpus and was chosen to be a decaying exponential in “r/c”, where “c” is the length of the candidate translation and “r” is the effective length of the reference translation.

$$BP = \begin{cases} 1 & \text{Si } c > r \\ e^{1-\frac{r}{c}} & \text{Si } c \leq r \end{cases} \quad (2)$$

5.2 WER (Word Error Rate)

The WER metric, Proposed by Popovic and Ney in 2007. Originally used in Automatic Speech Recognition, compares a sentence hypothesis refers to a sentence based on the Levenshtein distance. It is also used in machine translation to evaluate the quality of a translation hypothesis in relation to a reference translation. For this, the idea is to calculate the minimum number of edits (insertion, deletion or substitution of the word) to be performed on hypothesis translation to make it identical to the reference translation. The number of editss to be performed, noted “ $d_L(\text{ref}, \text{hyp})$ ” is then divided by the size of the reference translation, denoted “ N_{ref} ” as shown in the following formula [22]:

$$WER = \frac{1}{N_{ref}} \times d_L(\text{ref}, \text{hyp}). \quad (3)$$

Where:

- $d_L(\text{ref}, \text{hyp})$: is the Levenshtein distance between the reference translation “ref” and the hypothesis tanslation “hyp”.

A shortcoming of the WER is the fact that it does not allow reordering of words, whereas the word order of the hypothesis can be different from word order of the reference even though it is correct translation.

5.3 PER (Position-independent word Error Rate)

The PER metric, proposed by Tillman in 1997. compare the words of machine translation with those of the reference regardless of their sequence in the sentence. The PER score is defined by the following formula [23]:

$$PER = \frac{1}{N_{ref}} \times d_{per}(\text{ref}, \text{hyp}). \quad (4)$$

Where:

- d_{per} : calculates the difference between the occurrences of words in machine translation and the translation of reference.

A shortcoming of the PER is the fact that the word order can be important in some cases.

5.4 TER (Translation Error Rate)

The TER metric, proposed by Snover in 2006. Is defined as the minimum number of edits needed to change a hypothesis so that it exactly matches one of the references. The possible edits in TER include insertion, deletion, and substitution of single words, and an edit which moves sequences of contiguous words. Normalized by the average length of the references. Since we are concerned with the minimum number of edits needed to modify the hypothesis, we only measure the number of edits to the closest reference. The TER score is defined by the following formula [24]:

$$TER = \frac{Nb (op)}{Avreg N_{Ref}} \quad (5)$$

Where:

- Nb (op) : is the minimum number of edits;
- Avreg N_{ref}: the average size in words references.

6 Conclusion

In conclusion, we can say that the field of machine translation has been and remains a key focus of research on natural language processing and that led to the development of many positive results. However, perfection is still far away. If the translators have today reached a level of reliability and efficiency in a technical text, perfection is still a long way in the literary text, overwhelmed by the intricacies, the puns and colorful expressions. We think it must look to the construction of a translator hybrid (combining statistical and rules) at the end to increase the performance of the translation system.

References

1. Hutchins, W. J. and Somers, H. L., An introduction to machine translation, Academic Press, London. (1992)
2. Baumgartner-Bovier, “ La traduction automatique, quel avenir ? Un exemple basé sur les mots composés ”, Cahiers de Linguistique Française N°25, (2003).
3. J. Chandioux, “Histoire de la traduction automatique au Canada”, journal des traducteurs, vol. 22, n° 1, p. 54-56, (1977).
4. H. Kaji, “HICATS/JE : A Japanese-to-English Machine Translation System Based on Semantics ”, Machine Translation Summit, (1987).
5. Y. Lepage, E. Denoual, “ALEPH: an EBMT system based on the preservation of proportional analogies between sentences across languages ”, (2005).
6. Y. Fukumochi, “A Way of Using a Small MT System in Industry ”, the 5th Machine Translation Summit, July 10-13, (1995).

7. M. Cori et J. Léon , “ La constitution du TAL Étude historique des dénominations et des concepts ”, TAL. Volume 43 – n° 3/(2002).
8. C. Granell, “La Traduction automatique, Pour qui ? Pour Quoi ? ”, Support de cours, Novembre (2010).
9. P. P. Monty, “Traduction statistique par recherche locale”, , Mémoire de Maitre des sciences en informatique, Université de Montréal, (2010).
10. F. Yvon, “Une petite introduction au traitement Automatique du langage naturel, support de cours ”, Ecole Nationale Supérieure des télécommunications, Avril (2007).
11. C. Fuchs, B. Habert, “ Introduction le traitement automatique des langues : des modèles aux ressources ”, Article paru dans Le Français Moderne LXXII Volume1, (2004).
12. P. Bouillon, “Traitement automatique des langues naturelles ”, édition Duculot, (1998).
13. Hutchins J., Machine Translation: A Brief History, Concise History of the Language Sciences: From the Sumerians to the Cognitivists. Koerner E. F. K. and Asher R. E. (ed.). Oxford: Pergamon Press, pp. 431- 445, (1995).
14. Sumita E., Iida H., and Kohyama H., Translating with Examples: A New Approach to Machine Translation, the Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Language, pp. 203–212, (1990).
15. Lavie L., Vogel S., Peterson E., Probst K., Font-Llitió A., Reynolds R., Carbonell J., and Cohen R., Experiments with a Hindi-to-English Transfer-Based MT System under a Miserly Data Scenario, ACM Transactions on Asian Language Information Processing TALIP, Papineni, K., Roukos, S., Ward, and T.: Maximum Likelihood and Discriminative, pp.143 – 163, (2004).
16. Imamura K., Okuma H., Watanabe T., and Sumita E., Example-based Machine Translation Based on Syntactic Transfer with Statistical Models, Proceedings of the 20th International Conference on Computational Linguistics, Vol. 1, University of Geneva, Switzerland, pp. 99-105, August (2004).
17. Imamura K., Doctor's Thesis Automatic Construction of Translation Knowledge for Corpus-based Machine Translation, May 10, (2004).
18. Lavie L., Vogel S., Peterson E., Probst K., Wintner S., and Eytani Y., Rapid Prototyping of A Transfer-Based Hebrew-to-English Machine Translation System, Proceedings of the 10th International Conference on Theoretical and Methodological Issues in Machine Translation TMI-04. Baltimore, MD USA, pp.1-10, October (2004).
19. Probst K., Peterson E., Carbonell J and Levin L., MT for Minority Language Using Elicitation-based Learning of Syntactic Transfer Rules. Machine Translation 17: 245-270, Kluwer Academic Publishers, pp. 245 – 270, (2002).
20. Zantout R., and Guessoum A., Arabic Machine Translation: A Strategic Choice for the Arab World, journal of King Saud University, Volume 12, pp. 299-335, (2000).
21. K. Papineni, S. Roukos, T. Ward, and W. Zhu, “Bleu: a Method for Automatic Evaluation of Machine Translation”, Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), Philadelphia, pp. 311-318, (2002).
22. M. Popovic and H. Ney.” Word error rates: Decomposition over POS classes and applications for error analysis”. In Proceedings of ACL Workshop on Machine Translation.
23. C. Tillman, , S. Vogel, H. Ney, H. Sawaf, and A. Zubiaga. .”Accelerated DP-based search for statistical translation”. In Proceedings of the 5th European Conference on Speech Communication and Technology, pp- 2667.2670. Rhodes, Greece. (1997).
24. M. Snover, B. Dorr, , R. Schwartz, L. Micciulla, J. Makhoul.: “A Study of Translation Edit Rate with Targeted Human Annotation”. In Proceedings of AMTA, Boston, (2006).
25. Abdullah H. Homiedan, “Machine translation”, Journal of King Saud University, Language & Translation Vol 10, pp.1.21, (1998).