

Handheld Machine Translation System Based on Constraint Synchronous Grammar

Fai Wong, Francisco Oliveira, Sam Chao, Chi-Wai Tang

Faculty of Science and Technology

University of Macau

{derekfw, olifran, lidiasc, kevintang}@umac.mo

Abstract

The advancements of mobile technology permit handheld devices to be smaller, versatile, and have more processing power. On the other hand, the development of complex applications which require more processing capabilities are being developed rapidly nowadays. The implementation of Machine Translation (MT) systems with high translation quality is always considered difficult in desktop devices. In order to understand the languages deeply, large amounts of knowledge and processing capabilities are always required to guarantee the translation quality. This turns out that it is even more challenging for handheld devices. As a result, this paper introduces the application of MT based on Constraint Synchronous Grammar (CSG) in devices with limited resources. Since CSG describes syntactic structures of two languages simultaneously based on feature constraints, the analysis and the generation of the translation can be done at one stage to lower the process complexity.

1 Introduction

With the rapid development of mobile technology, many useful translation applications are being ported to handheld devices, including pocket bilingual dictionaries, and Machine Translation (MT) systems. Pocket bilingual dictionaries provide not only the translation of words and phrases to the other languages but also the pronunciation. Although the size of the dictionaries is small, users

can access easily the meaning of a word at any time. On the other hand, porting MT applications from powerful desktop computers to handheld devices is always difficult. Since they require a large amount of computation and linguistic knowledge to guarantee the quality of the translation, when the central processing unit, memory, and storage requirements cannot fulfill the processing capability of the handheld devices, it is impossible to guarantee the quality of the translation. As a result, a more careful design has to be accomplished, including the methodologies applied in MT and the way that these systems should be deployed in the handheld devices.

In the literature, there are many different approaches applied in the MT field. Those techniques can be generally classified into three main categories: Rule based, Example based, and Statistic based.

Rule based MT (Bennett and Slocum, 1985) approach is based on a set of linguistic grammar rules for handling the translation. Recently, researchers considered relationships between syntactical structures of two languages simultaneously based on synchronous grammar rules in the parsing process. Wu (1995) proposed Inversion Transduction Grammar for defining a single parsing structure based on a set of brackets to account for both languages simultaneously. Multiple Context Free Grammar (Seki et al., 1991) was used by defining a set of functions for non-terminal symbols in the productions for interpreting the symbols during the generation phase. Deneefe and Knight (2009) proposed a practical way in developing a MT system

based on Synchronous Tree Adjoining Grammar. Although the accuracy can be guaranteed for small and close domains, the construction of rules can be expensive and time consuming.

Example based MT (Brown, 1996) analyzes different pieces of bilingual examples stored in parallel corpora for generating the translation. However, it often depends on the quality of the examples and the similarity function applied.

With the large available of digitized parallel language resources nowadays, Statistic based MT (Lopez, 2008) has become a new research trend. Probabilities are estimated between the translation of words and the ordering of the sentences from the text corpora. The accuracy often depends with the information of the digitized resources.

In the design of MT systems in the handheld devices, one common approach is to build up a thin client with fat server architecture. The client is the handheld device, and it provides an interface to input the source sentence and show the translation results. The translation task is done in the server side and their communication is through different communication channels. John Hutchins (2005) mentioned some software sellers that provide short message services (SMS) to translate one language to another by sending SMS requests to a specific number. Since the size of the text is limited, it may need to send and receive several messages when the text gets longer.

Some considered Wireless connections such as Wi-Fi or 3G enabled devices. Yamabana et al. (2003) introduced a client-server speech translation system. The translation modules reside in the server side, and the core translation is based on Lexicalized Tree AutoMata-based Grammar formalism. Michael Paul et al. (2008) presented two translation services for mobile phones. Translation queries are sent to the Server side, and they are translated by a phrase-based statistical MT system. Frago et al. (2011) presented an augmented reality translation system TranslatAR for a specific mobile device. The system first recognizes text captured from the video camera of the mobile device, then it is sent to Google translate for translation, and finally the result is replaced on top of the original text in the captured freeze video frame. The above approaches suffer from the location dependency and the communication cost. First, not all locations have wireless connections available, and users cannot do the translation at their desired

places. Second, the cost of using wireless services is very high when the translation involves large amounts of texts.

Another approach is to design and implement all the required translation modules in the handheld devices. Waibel et al. (2003) demonstrated an application in pocket digital assistants. They investigated two interlingua based approaches, including knowledge based and statistical based method. Zhang and Vogel (2007) proposed a phrase based SMT system called Pandora that runs directly in these devices. Different solutions are proposed to overcome the processing and memory limitations, including: the transformation of words into integer symbols to reduce the size, the cross-indexing of the phrase translation models to reduce the redundancy in the representation, and the serialization of the model files to directly access the external storage card instead of the dynamic random access memory. Gao et al. (2008) introduced a Mastor speech-to-speech MT system optimized for handheld devices. There are two components in the translation module: statistical natural language understanding (NLU) and statistical natural language generation (NLG) module. NLU module is based on a statistical parser that analyses statistical decision-tree models in the identification of the meaning and structure of the input sentence. NLG component is responsible for the ordering and generation of the target language.

This paper presents the application of Constraint Synchronous Grammar (CSG) (Wong et al., 2005) formalism to MT for handheld devices. The whole translation system is integrated in the handheld devices so that users can perform the translation at any time and location without considering the communication network cost. Since the system is designed for mobile environment, a simple to use interface is developed, as shown in Figure 1.

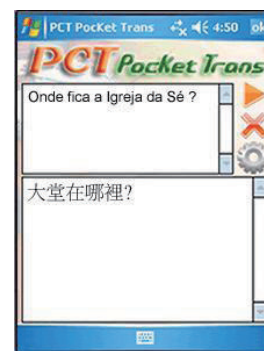


Figure 1. System's Interface in Mobile Device

The application runs in Windows Mobile environment, the user only needs to type in the source sentence in the upper area and press the button on the device to trigger the translation. After processing, the bottom area displays the translation results.

CSG is a variation of synchronous grammar that is used to express syntactic relationships between the source sentential pattern with one or more target patterns. The selection of the most suitable target is based on the defined feature constraints for each grammar rule. Since CSG uses feature structures in describing detailed information for each constituent (including Part-of-Speech, gender, number agreement, sense, etc), it effectively removes ambiguities in the analysis and parsing. In the design of the system, only the most important translation modules are implemented in the handheld devices, including the preprocessing module for the morphological analysis, and the CSG translation module for the analysis and generation of the translation. As long as suitable bilingual CSG rules are added in the knowledge base, different language pairs can be considered, and in this paper, a Portuguese to Chinese translation based on CSG is used as an example to demonstrate its feasibility in the handheld devices.

This paper is organized as follows. An introduction of CSG is given in section 2. The design and implementation of the translation system in the traditional and in the mobile environment are detailed in section 3. The semi-automatic acquisition of CSG rules is presented in section 4. Evaluations are given in section 5, followed by a conclusion.

2 Constraint Synchronous Grammar

Constraint Synchronous Grammar is based on the formalism of Context Free Grammar (CFG) to the case of synchronous. In CSG formalism, it consists of a set of production rules that describes the sentential patterns of the source text and target translation patterns.

In CSG, every production rule is in the form of the example shown in (1). The source pattern $NP_1 VP NP_2 PP NP_3$ is associated with two target patterns, including $NP_1 VP^1 NP_3 VP^2 NP_2$, and $NP_1 VP NP_3 NP_2$ respectively.

The determination of the suitable generative rule is based on the control conditions defined by rule.

The one satisfying all the conditions determines the relationship between the source and target sentential pattern. For example, if the category of the verb is $vb1$, and the sense of the subject, indirect, and direct objects governed by the verb, VP , corresponds to the first, second, and the third nouns (NP), then the source pattern $NP_1 VP NP_2 PP NP_3$ is associated with the target pattern $NP_1 VP^1 NP_3 VP^2 NP_2$.

$$S = NP_1 VP^* NP_2 PP NP_3 \{ \begin{array}{l} [NP_1 VP^1 NP_3 VP^2 NP_2; C_1] \\ [NP_1 VP NP_3 NP_2; C_2] \end{array} \}$$

$$C_1 = \{ VP_{\text{category}} = vb1, \quad (1)$$

$$\begin{array}{l} VP_{\text{sense subject}} = NP_1_{\text{sense}}, \\ VP_{\text{sense indirect object}} = NP_2_{\text{sense}}, \\ VP_{\text{sense object}} = NP_3_{\text{sense}} \end{array}$$

$$C_2 = \{ VP_{\text{sense subject}} = NP_1_{\text{sense}}, \\ VP_{\text{sense indirect object}} = NP_2_{\text{sense}} \}$$

Their relationship is established by the given subscripts and the sequence is based on the target sentential pattern. As an example, in the first generative rule, although the first NP in the source pattern corresponds to the first NP in the target one, the sequence for the second and third NP in the source are changed in the target sentential pattern. The asterisk “*” indicates the head element, and its usage is to propagate all the related features/linguistic information of the head symbol to the reduced non-terminal symbol in the left hand side. The use of the “*” is to achieve the property of features inheritance in CSG formalism.

CSG is especially effective for modeling non-standard linguistic phenomena for languages which are structurally different. The ordering of the constituents is modeled easily by using the subscripts and the sequence defined in CSG production rule. The discontinuity between words in different languages is solved by defining non-terminal symbols that appear in the source but not the target pattern or vice-versa. As an example, in the first target sentential pattern in production (1), two verbs (VP^1 and VP^2) are associated with the source one. Similarly, the consideration of constituents that are disappeared or shown in the target syntactical pattern is handled in a similar way.

CSG rules are parsed by a modified version of generalized LR algorithm (Tomita, 1987), a shift-

reduce approach based on an extended LR parsing table. Besides having the actions to be accomplished (shift, reduce, accept), and the state of the parser at different stages of parsing, the table is extended by taking feature's constraints and target rules into consideration. In other words, as the parser identifies CSG productions through the normal shift actions, it checks the associated constraints to determine if the current reduce action is valid or not.

3 System's Design and Architecture

The architecture and design are different in the desktop computers and handheld devices. Figure 2 shows the translation process in the desktop computers based on CSG formalism.

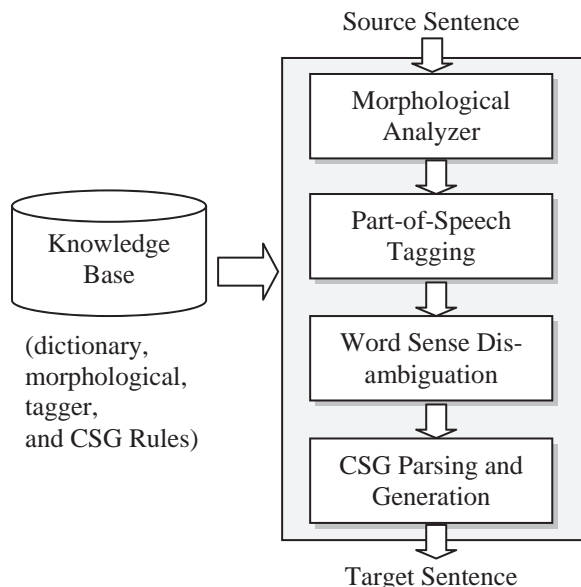


Figure 2. MT System's Design in Desktop Computers

The input sentence is first morphologically analyzed in order to restore the word to its original form. The next step is to determine the grammatical categories for each word based on a probabilistic tagger. A word sense disambiguation module based on context collocation to select the most probable target translation is considered afterwards. Finally, in the parsing stage, the source sentence is parsed based on the constraints defined while at the same time, it generates the corresponding target translation. In this design, the MT system requires a large amount of data to success, including bilin-

gual lexicon; morphological, tagger, and disambiguation rules; and CSG rules.

In the handheld devices environment, in order to reduce the processing complexity and the amount of resources in the handheld devices, the architecture of the MT system is simplified to meet the requirements, as shown in Figure 3.

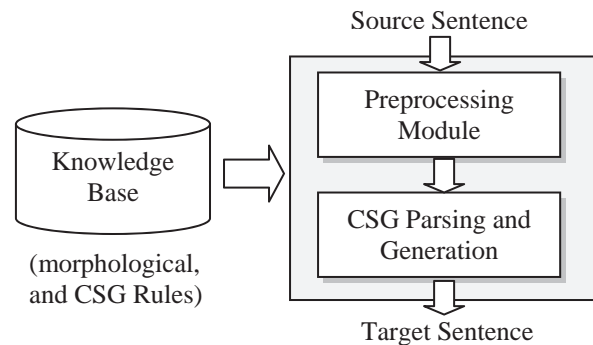


Figure 3. MT System's Design in Handheld Devices

The input Portuguese sentence is first preprocessed by the morphological analyzer. However, if the input sentence is a different language, for example the Chinese language, a segmentation module for the identification of word boundaries should be incorporated. Although Part-of-Speech Tagging module is not considered, lexicalization is applied to merge lexicon entries with Constraint Synchronous Grammars. The main idea is to associate syntactic contexts to each lexicon, and let the whole translation process dependent only on the defined rules. The sentence is then parsed based on the constraints defined afterwards. Unlike traditional transfer-based MT architectures where the translation process is carried out in pipeline using different sets of representation rules for structure analysis and transformation for the analyzer, transfer, and generation modules, the proposed MT system only requires a one stage analysis based on CSG rules. Since the structures of parallel languages are synchronized in the formalism, their structural differences are also captured and described by the grammar. Hence, translation of an input text essentially involves three steps in the parsing and generation module. First, the structure of an input sentence is analyzed using the source component's rules from the CSG productions. Second, the reduction process is based on the feature constraints defined. Finally, the selected target sequence is used to generate the corresponding

The restoration of the words into its original format and the generation of CFG rules are done automatically by the annotation system developed in the creation of CSG rules. Moreover, this annotation system provides an interface to show the syntactic structure converted from CSG rules of the sentence and let users to extend it with proper translations and target syntactic sequential patterns. Instead of working in the grammar rules directly, the annotation system provides a semi-automatic construction of lexicalized grammars with the help of the users.

5 Evaluation and Discussion

The effectiveness of the proposed MT system based on CSG formalism running in mobile devices is investigated for handling the translation between Portuguese and Chinese. A prototyping system is built on a Windows Mobile handheld device with a 400 MHz Central Processing Unit (CPU), 256 MB Read Only Memory (ROM), and 128 MB Random Access Memory (RAM).

The translation quality is measured by three automatic evaluation metrics, including NIST (Doddington, 2002), BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and by human assessments. The average value of three human evaluated results is considered in human assessments, and the quality of the translated sentences are classified as Good, Acceptable, and Bad. The size of the MT's knowledge used for the experiments includes 200 entries for morphological rules and 718 entries for CSG rules.

| Evaluation Methods | MT in Handheld Device |
|--------------------|-----------------------|
| BLEU | 0.8026 |
| NIST | 9.5181 |
| METEOR | 0.8671 |
| Good | 68.6% |
| Acceptable | 17.5% |
| Bad | 13.9% |

Table 1. Evaluation Results in Close Domain

In the first experiment, a close domain and a text file with 8k size are considered to evaluate the quality of the MT system running in the handheld device. Evaluation results are shown in Table 1.

Another experiment is conducted by comparing the translation quality and efficiency between the system running in the MT system and the one in the desktop computer. Moreover, another test suite of 6k size is considered to evaluate the system when the knowledge is out of domain. Evaluation results are shown in Table 2.

| Evaluation Methods | Handheld Device | Desktop Computer |
|--------------------|-----------------|------------------|
| BLEU | 0.2963 | 0.3508 |
| NIST | 4.1111 | 5.4384 |
| METEOR | 0.4326 | 0.5430 |
| Good | 27.8% | 20.6% |
| Acceptable | 10% | 25.3% |
| Bad | 62.2% | 54.1% |

Table 2. Evaluation Results in open domain

The scores are directly affected by several factors. In a close domain, since the knowledge acquired covers most of the sentences, the system achieves a BLEU score of 0.8026 and the accuracy is 86.1% for Good and Acceptable cases. However, in the open domain, the accuracy drops drastically, due to the following reasons. Since some automatic evaluation methods rely on n-gram co-occurrence precision, it may generate low scores even if the translated sentence is correct when compared with the reference translations that use different synonyms. Moreover, Callison-Burch et al. (2006) argued that BLEU evaluation metric is not adequately suitable for Rule based MT systems. Second, as all the words are lexicalized and stored in the format of CSG rules, if sentences contain words that do not exist in our knowledge, an incorrect translation will be generated. The same case happens when the rules acquired do not cover the syntactical structure of the source sentences.

In some sense, there is a trade-off between efficiency and accuracy. In order to reduce the execution time, memory and disk space in the mobile environment, in our approach, the translation task is mainly accomplished by the CSG parsing and generation module and one knowledge base. This approach not only reduces the complexity of the MT system effectively but also increases the translation efficiency in the mobile environment. However, the side effect is that the accuracy decreases when the knowledge is out of domain. This can be

reflected by comparing the results performed in the handheld device and the desktop computer based on the second test suite. By using the same set of grammar rules, although many syntactic structures are not covered for both devices, in desktop computer, it achieves much better results in all the automatic evaluation methods compared with the one running in the handheld device. Similar improvements are also obtained in human assessments.

In terms of translation efficiency, a traditional MT system installed in the desktop computer is compared with the proposed MT approach in the mobile device. Hardware specifications, and the average response time are concluded in Table 3.

| Device | Hardware/Time | Specification |
|------------------|---------------------------------------|---------------|
| Mobile Device | CPU | 400 MHz |
| | RAM | 128 MB |
| | Data Size | 2 MB |
| | Average Response Time for translation | 2.5 seconds |
| Desktop Computer | CPU | 1.6 GHz |
| | RAM | 128 MB |
| | Data Size | > 100 MB |
| | Average Response Time for translation | 0.8 seconds |

Table 3. Translation Efficiency and Hardware Specifications

In terms of speed, the results give an estimate on how fast the system runs in translating a sentence rather than comparing them directly. The main reason is because the CPU and memory are not comparable with each other. On the other hand, in the desktop environment, it requires much more memory space than the one running in the handheld device for such improvement.

6 Conclusion

In this paper, the application of Machine Translation based on Constraint Synchronous Grammar formalism in the handheld devices is presented. Due to the limitations of the handheld devices, in our design, the whole system only consists of a CSG Knowledge Base and two modules: preprocessing and CSG parsing and generation module. In Portuguese-Chinese translation, the first module

restores the original format of the words in the source sentence, and the second module is responsible for the analysis of the source and the generation of the target sentence based on the constraints defined. Since the knowledge is highly dependent in the CSG rules, semi-automatic extraction methodologies are proposed, including the acquisition of skeletal syntactic structures, and lexicalized synchronous grammar based on statistical tools. In order to ensure the translation quality and to remove unnecessary disambiguation, the help of linguistics is considered by adding suitable constraints and analyzing the correctness of the extracted rules. Finally, experiments are conducted and a prototyping system is built to evaluate the effectiveness of the proposed system.

Acknowledgments

This work was partially supported by the Research Committee of University of Macau under grant UL019/09-Y2/EEE/LYP01/FST, and also supported by Science and Technology Development Fund of Macau under grant 057/2009/A2.

References

- Adam Lopez. 2008. *Statistical Machine Translation*. ACM Computing Surveys 40(3).
- Alex Waibel, Ahmed Badran, Alan W Black, Robert Frederking, Donna Gates, Alon Lavie, Lori Levin, Kevin Lenzo, Laura Mayfield Tomokiyo, Juergen Reichert, Tanja Schultz, Dorcas Wallace, Monika Woszczyna, and Jing Zhang. 2003. *Speechalator: Two-way Speech-to-Speech Translation In Your Hand*. In Proceedings of HLT-NAACL 2003 Demonstrations, Edmonton, pp. 29-30.
- Anne Abeillé. 1991. *Une grammaire lexicalisée d'arbres adjoints pour le français*. Ph.D. thesis, Université Paris.
- Chris Callison-Burch, Miles Osborne and Philipp Koehn. 2006. *Re-evaluating the Role of Bleu in Machine Translation Research*. In Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics, pp. 249-256.
- Dekai Wu. 1995. *Grammarless extraction of phrasal translation examples from parallel texts*. In Proceedings of the 6th International Conference on Theoretical and Methodological Issues in Machine Translation, Leuven, Belgium, pp. 354-372.
- Fai Wong, Dong-Cheng Hu, Yu-Hang Mao, Ming-Chui Dong, and Yi-Ping Li. 2005. *Machine Translation Based on Constraint-Based Synchronous Grammar*. In Proceedings of the 2nd International Joint Confe-

- rence on Natural Language, Jeju Island, Republic of Korea, pp. 612-623.
- Franz Josef Och, and Hermann Ney. 2003. *A systematic comparison of various statistical alignment models*. Computational Linguistics, 29(1):19-51.
- George Doddington. 2002. *Automatic evaluation of machine translation quality using n-gram co-occurrence statistics*. In Proceedings of the Second International Conference on Human Language Technology Research, San Diego, California, pp. 138-145.
- Hiroyuki Seki, Takashi Matsumura, Mamoru Fujii, Tadao Kasami. 1991. *On multiple context-free grammars*. Theoretical Computer Science 88(2):191-229.
- John Hutchins. 2005. *Current commercial machine translation systems and computer-based translation tools: system types and their uses*. International Journal of Translation 17(1-2):5-38.
- Kishore Papineni, Salim Roukos, Todd Ward, Wei-Jing Zhu. 2002. *BLEU: a Method for Automatic Evaluation of Machine Translation*. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), Philadelphia, Pennsylvania, pp. 311-318.
- Kiyoshi Yamabana, Seiya Osada, Ken Hanazawa, Akiyoshi Okumura, Ryosuke Isotani, Takao Watanabe. 2003. *A Speech Translation System with Mobile Wireless Clients*. In Proceedings of the 41st Annual Meeting on Association for Computational Linguistics, Sapporo, Japan, pp. 133-136.
- Masaru Tomita. 1987. *An efficient augmented-context-free parsing algorithm*. Computational Linguistics 13(1-2):31-46.
- Michael Paul, Hideo Okuma, Hirofumi Yamamoto, Eiichiro Sumita, Shigeki Matsuda, Tohru Shimizu, and Satoshi Nakamura. 2008. *Multilingual Mobile-Phone Translation Services for World Travelers*. In Proceedings of the 22nd International Conference on Computational Linguistics: Demonstration Papers, Manchester, United Kingdom, pp. 165-168.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. *Moses: Open Source Toolkit for Statistical Machine Translation*. In Proceedings of the ACL 2007 Demo and Poster Sessions, Prague, pp. 177-180.
- Ralf D. Brown. 1996. *Example-Based machine translation in the Pangloss system*. In Proceedings of the 16th International Conference on Computational Linguistics (COLING-96), Copenhagen, Denmark, pp. 169-174.
- Satanjeev Banerjee, and Alon Lavie. 2005. *METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments*. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization, Ann Arbor, Michigan, pp. 65-72.
- Steve DeNeefe, and Kevin Knight. 2009. *Synchronous Tree Adjoining Machine Translation*. In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, Singapore, pp. 727-736.
- Victor Fragoso, Steffen Gauglitz, Shane Zamora, Jim Kleban, Matthew Turk. 2011. *TranslatAR: A Mobile Augmented Reality Translator*. In Proceedings of the IEEE Workshop on Applications of Computer Vision, Kona, Hawaii, pp. 497-502.
- Winfield S. Bennett, Jonathan Slocum. 1985. *The LRC Machine Translation System*. Computational Linguistics 11(2-3):111-121.
- Ying Zhang, Stephan Vogel. 2007. *PanDoRA: A Large-scale Two-way Statistical Machine Translation System for Hand-held Devices*. In Proceedings of MT Summit XI, Copenhagen, Denmark, pp. 10-14.
- Yuqing Gao, Bowen Zhou, Weizhong Zhu, and Wei Zhang. 2008. *Handheld Speech to Speech Translation System*. Automatic Speech Recognition on Mobile Devices and over Communication Networks, Springer London.