

学校编码: 10384  
学 号: 20051302328

分类号\_\_\_\_\_密级\_\_\_\_\_  
UDC \_\_\_\_\_

厦门大学  
硕士 学位 论文

基于浅层句法分析的翻译模板  
自动获取研究

Translation Templates' acquisition based on  
**Partial parsing**

林 哲 辉

指导教师 : 史晓东 教授  
专业名称 : 计算机应用  
论文提交日期 : 2008 年 月  
论文答辩日期 : 2008 年 月  
学位授予日期 : 2008 年 月

答辩委员会主席: \_\_\_\_\_

评 阅 人: \_\_\_\_\_

2008 年 月

# 厦门大学学位论文原创性声明

兹呈交的学位论文，是本人在导师指导下独立完成的研究成果。  
本人在论文写作中参考的其他个人或集体的研究成果，均在文中以  
明确方式标明。本人依法享有和承担由此论文产生的权利和责任。

声明人（签名）：

年   月   日

厦门大学博硕士论文摘要库

# 厦门大学学位论文著作权使用声明

本人完全了解厦门大学有关保留、使用学位论文的规定。厦门大学有权保留并向国家主管部门或其指定机构送交论文的纸质版和电子版，有权将学位论文用于非赢利目的的少量复制并允许论文进入学校图书馆被查阅，有权将学位论文的内容编入有关数据库进行检索，有权将学位论文的标题和摘要汇编出版。保密的学位论文在解密后适用本规定。

本学位论文属于

1. 保密（），在 年解密后适用本授权书。
2. 不保密（）

（请在以上相应括号内打“√”）

作者签名： 日期： 年 月 日

导师签名： 日期： 年 月 日

厦门大学博硕士论文摘要库

## 摘要

传统的 EBMT(Example-Based Machine Translation, 基于实例的机器翻译)方法是建立在大规模的实例库基础之上的，存在着精确匹配率不高，模糊匹配时产生译文质量较差等缺点。利用翻译模板可以有效的解决翻译实例的数据稀疏问题、简化实例库的规模并提高实例匹配的精确率。因此，本文提出了基于两级翻译模板的方法，提出了两级翻译模板的体系结构并给出了定义，给出了从句子对齐的双语语料库自动抽取两级翻译模板的方法，并设计了利用模板库进行模板匹配的算法。

文中首先对机器翻译的重要技术尤其是基于EBMT的翻译方法进行了整理和分析，其次，概括了翻译模板提出的意义，提出了本文翻译模板的定义和研究策略。最后，在此基础上实现了两级翻译模板的抽取。

基于模板的翻译方法的核心问题是模板的抽取与匹配算法。本文在抽取翻译模板的过程中采用了浅层句法分析方法，主要运用了组块识别。浅层句法分析的优点是可以识别出确定性高的部分分析结果，减少句法分析中的歧义，从而降低句法分析的难度。具体方法是：在模板抽取过程中首先对翻译实例进行浅层句法分析，得到双语组块信息，并确定实例的框架，根据对齐约束提取短语翻译模板，在此基础上，提取出句子级的翻译模板。在模板匹配过程中首先对输入句进行组块分析，然后在系统模板库中检索匹配模板。检索过程中兼顾模板的结构信息及子块信息。

**关键词：**机器翻译 翻译模板 EBMT

厦门大学博硕士论文摘要库

## Abstract

EBMT ( Example-Based Machine Translation) systems are based on large scale example corpus in traditional, having the defect of low precision of matching. Translation template can solve the problem of data sparsity , large storage space and low matching precision of examples. Therefore, this thesis proposes one method which automatically extracts bi-level translation templates from sentence-aligned bilingual corpus, we propose its definition and matching algorithm. Finally, We proposed the matching algorithm as well.

Firstly, we analyze some important techniques of machine translation, especially the EBMT method. Secondly, we summarize the meaning of translation template, propose our definition and research strategy. Finally, we accomplish the extraction of bi-level translation templates.

The key point of EBMT method is translation templates' extraction and match algorithm. We used partial parsing in the process of extraction, chunking in general. The advantage of partial parsing is the robust result that can decrease the ambiguity of parsing. Concretely, partial parsing on the sentence-aligned bilingual corpus and get the chunking results, then we set the frame of the sentence, after extraction of phrase templates, we can get the sentence template. The matching module searches the most matching template for input sentence in the database, with the information of the chunking result and frame structure. The templates matching algorithm gets the searching result by using the key word in the sentence frame and its chunking template.

**Key Words:** Machine Translation; Translation Template; EBMT

厦门大学博硕士论文摘要库

# 目 录

|                                 |           |
|---------------------------------|-----------|
| <b>第一章 引言 . . . . .</b>         | <b>1</b>  |
| 1.1 自然语言处理与机器翻译综述 .....         | 1         |
| 1.2 机器翻译简介 .....                | 1         |
| 1.3 语料库与翻译模板 .....              | 2         |
| 1.4 本文的工作 .....                 | 3         |
| 1.4.1 研究背景和目标.....              | 3         |
| 1.4.2 正文组织.....                 | 3         |
| <b>第二章 机器翻译技术概述 . . . . .</b>   | <b>5</b>  |
| 2.1 机器翻译问题及其研究 .....            | 5         |
| 2.1.1 问题描述.....                 | 5         |
| 2.1.2 历史综述.....                 | 6         |
| 2.2 基于规则的机器翻译 .....             | 7         |
| 2.3 基于统计的机器翻译 .....             | 9         |
| 2.3.1 概述.....                   | 9         |
| 2.3.2 基于信源信道模型的统计机器翻译方法.....    | 10        |
| 2.3.3 语言模型.....                 | 10        |
| 2.3.4 翻译模型.....                 | 11        |
| 2.3.5 解码算法.....                 | 12        |
| 2.4 基于实例的机器翻译 .....             | 12        |
| 2.5 机器翻译主流评测方法 .....            | 13        |
| 2.5.1 几个简单的机器翻译自动评价指标.....      | 14        |
| 2.5.2 IBM 的 BLEU 评价方法.....      | 14        |
| 2.5.3 NIST 评价.....              | 15        |
| <b>第三章 翻译模板及其研究意义 . . . . .</b> | <b>17</b> |
| 3.1 EBMT .....                  | 17        |

|       |                            |           |
|-------|----------------------------|-----------|
| 3.1.1 | EBMT 的基本原理 .....           | 17        |
| 3.1.2 | EBMT 的基本流程 .....           | 17        |
| 3.1.3 | 泛化的 EBMT .....             | 18        |
| 3.2   | <b>翻译模板及其研究现状 .....</b>    | <b>19</b> |
| 3.2.1 | kaji 的方法.....              | 20        |
| 3.2.2 | CMU 的方法 .....              | 21        |
| 3.2.3 | Bilkent 的方法.....           | 22        |
| 3.2.4 | 东北大学的方法.....               | 22        |
| 3.2.5 | 厦门大学的实验.....               | 23        |
| 3.3   | <b>定义层次 .....</b>          | <b>24</b> |
| 3.4   | <b>本文的定义 .....</b>         | <b>24</b> |
|       | <b>第四章 翻译模板的抽取算法 .....</b> | <b>26</b> |
| 4.1   | <b>浅层句法分析综述 .....</b>      | <b>26</b> |
| 4.1.1 | 浅层句法分析方法.....              | 27        |
| 4.1.2 | 组块识别.....                  | 27        |
| 4.1.3 | 基本名词短语（BaseNP）识别.....      | 28        |
| 4.1.4 | CRF(条件随机场).....            | 28        |
| 4.2   | <b>应用.....</b>             | <b>29</b> |
| 4.3   | <b>短语抽取方法 .....</b>        | <b>31</b> |
| 4.4   | <b>模板抽取的难点 .....</b>       | <b>32</b> |
| 4.4.1 | 关于分析深度.....                | 32        |
| 4.4.2 | 关于框架结构.....                | 32        |
| 4.5   | <b>模板抽取算法 .....</b>        | <b>33</b> |
| 4.5.1 | 基于双语分析算法描述.....            | 33        |
| 4.5.2 | 基于单语分析算法描述.....            | 35        |
| 4.5.3 | 应用 .....                   | 36        |
| 4.6   | <b>模板翻译概率 .....</b>        | <b>37</b> |
| 4.7   | <b>句子级模板的抽取 .....</b>      | <b>38</b> |
| 4.8   | <b>相关实验 .....</b>          | <b>38</b> |

|                                |                   |           |
|--------------------------------|-------------------|-----------|
| 4.8.1                          | 语料库及其预处理.....     | 38        |
| 4.8.2                          | 部分实验结果.....       | 42        |
| 4.8.3                          | 例子.....           | 42        |
| 4.9                            | 结果分析 .....        | 45        |
| <b>第五章 翻译模板的存储与匹配 .....</b>    |                   | <b>47</b> |
| 5.1                            | 模板库的存储组织 .....    | 47        |
| 5.1.1                          | 句子级模板的存储和索引.....  | 47        |
| 5.1.2                          | 短语级模板的存储.....     | 48        |
| 5.2                            | 模板的匹配 .....       | 48        |
| 5.2.1                          | 句子级模板的匹配.....     | 49        |
| 5.2.2                          | 短语级模板的匹配.....     | 49        |
| 5.3                            | 译文生成 .....        | 49        |
| <b>第六章 基于模板的 EBMT 系统 .....</b> |                   | <b>51</b> |
| 6.1                            | EBMT 系统的类型.....   | 51        |
| 6.2                            | EBMT 系统的关键问题..... | 51        |
| 6.3                            | 系统流程 .....        | 52        |
| 6.4                            | 相关实验 .....        | 54        |
| 6.5                            | 小结.....           | 55        |
| <b>第七章 总结与展望 .....</b>         |                   | <b>56</b> |
| 7.1                            | 结论.....           | 56        |
| 7.2                            | 今后的工作方向 .....     | 56        |
| <b>参考文献 .....</b>              |                   | <b>59</b> |
| <b>攻读硕士学位期间发表论文 .....</b>      |                   | <b>63</b> |
| <b>致谢 .....</b>                |                   | <b>65</b> |
| <b>附录 1 .....</b>              |                   | <b>67</b> |

|            |    |
|------------|----|
| 附录 2 ..... | 69 |
|------------|----|

厦门大学博硕士论文摘要库

## Contents

|  |           |
|--|-----------|
| <b>Chapter1 Introduction.....</b>                                | <b>1</b>  |
| 1.1    A Survey of NLP and MT .....                              | 1         |
| 1.2    Introduction to Machine Translation .....                 | 1         |
| 1.3    Corpus and Translation Template .....                     | 2         |
| 1.4    Overview of Our Work .....                                | 3         |
| 1.4.1    Research Background and Target.....                     | 3         |
| 1.4.2    Structure of The Content.....                           | 3         |
| <b>Chapter2 A survey on Machine Translation Techniques .....</b> | <b>5</b>  |
| 2.1    Machine Translation and Related Research.....             | 5         |
| 2.1.1    Problem Description .....                               | 5         |
| 2.1.2    Overview of History .....                               | 6         |
| 2.2    RBMT.....   | 7         |
| 2.3    SMT.....  | 9         |
| 2.3.1    Introduction.....                                       | 9         |
| 2.3.2    Noisy Channel-Based Model SMT .....                     | 10        |
| 2.3.3    Language Model .....                                    | 10        |
| 2.3.4    Transltion Model.....                                   | 11        |
| 2.3.5    Decoder Algorithm.....                                  | 12        |
| 2.4    EBMT.....   | 12        |
| 2.5    Evaluating Measure of MT .....                            | 13        |
| 2.5.1    Some Simple Evaluate Guideline .....                    | 14        |
| 2.5.2    BLEU .....  | 14        |
| 2.5.3    NIST.....   | 15        |
| <b>Chapter 3 Translation Template and Its Significance .....</b> | <b>17</b> |
| 3.1    EBMT.....   | 17        |
| 3.1.1    Keystone of EBMT .....                                  | 17        |

|            |   |           |
|------------|---|-----------|
| 3.1.2      | Flow of EBMT .....                              | 17        |
| 3.1.3      | Extensive EBMT .....                            | 18        |
| <b>3.2</b> | <b>Translation Template and Its Status.....</b> | <b>19</b> |
| 3.2.1      | kaji's Method .....                             | 20        |
| 3.2.2      | CMU's Method .....                              | 21        |
| 3.2.3      | Bilkent's Method.....                           | 22        |
| 3.2.4      | NEU's Method .....                              | 22        |
| 3.2.5      | XMU's Method .....                              | 23        |
| <b>3.3</b> | <b>Definition Arrangement.....</b>              | <b>24</b> |
| <b>3.4</b> | <b>Our Definition .....</b>                     | <b>24</b> |

## **Chapter 4 The Algorithm of Translation Templates' Acqusition .....26**

|            |  |           |
|------------|--|-----------|
| <b>4.1</b> | <b>Introduction to Partial Parsing.....</b>                    | <b>26</b> |
| 4.1.1      | Means of Partial Parsing .....                                 | 27        |
| 4.1.2      | Chunk Parsing .....  | 27        |
| 4.1.3      | BaseNP.....  | 28        |
| 4.1.4      | CRF .....  | 28        |
| <b>4.2</b> | <b>Application.....</b>  | <b>29</b> |
| <b>4.3</b> | <b>Method of Phrase Extraction.....</b>                        | <b>31</b> |
| <b>4.4</b> | <b>Difficulty of Translation Templates' Acqusition.....</b>    | <b>32</b> |
| 4.4.1      | Analysis Depth.....  | 32        |
| 4.4.2      | Frame Structure.....   | 32        |
| <b>4.5</b> | <b>The Algorithm of Translation Templates' Acqusition.....</b> | <b>33</b> |
| 4.5.1      | Billigual-Based Algorithm .....                                | 33        |
| 4.5.2      | Unllingual-Based Algorithm .....                               | 35        |
| 4.5.3      | Application.....   | 36        |
| <b>4.6</b> | <b>Translation Probability of Translation Templates .....</b>  | <b>37</b> |
| <b>4.7</b> | <b>Extraction of Sentence Template .....</b>                   | <b>38</b> |
| <b>4.8</b> | <b>Related Experiments.....</b>                                | <b>38</b> |
| 4.8.1      | Corpus and Preprocess .....                                    | 38        |

|  |   |           |
|--|---|-----------|
| 4.8.2  | Experiment Results .....                          | 42        |
| 4.8.3  | Example .....                                     | 42        |
| <b>4.9</b>   | <b>Analyse of Results .....</b>                   | <b>45</b> |
| <b>CHapter 5 Storage and Matching of Translation Templates .....</b> |   | <b>47</b> |
| <b>5.1</b>   | <b>Storage form of Translation Template .....</b> | <b>47</b> |
| 5.1.1  | Storage and Index of SentenceTemplate .....       | 47        |
| 5.1.2  | Storage of Chunk Template.....                    | 48        |
| <b>5.2</b>   | <b>Matching of templates .....</b>                | <b>48</b> |
| 5.2.1  | Matching of Sentence Template.....                | 49        |
| 5.2.2  | Matching of Chunk Template .....                  | 49        |
| <b>5.3</b>   | <b>Translation Generate .....</b>                 | <b>49</b> |
| <b>Chapter 6 An EBMT System Based on Templates .....</b>             |   | <b>51</b> |
| <b>6.1</b>   | <b>Types of EBMT.....</b>                         | <b>51</b> |
| <b>6.2</b>   | <b>Key Problem of EBMT .....</b>                  | <b>51</b> |
| <b>6.3</b>   | <b>System Flow.....</b>                           | <b>52</b> |
| <b>6.4</b>   | <b>Related Experiment .....</b>                   | <b>54</b> |
| <b>6.5</b>   | <b>Brief Summary .....</b>                        | <b>55</b> |
| <b>Chapter 7 Conclusion and Prospect.....</b>                        |   | <b>56</b> |
| <b>7.1</b>   | <b>Summary.....</b>                               | <b>56</b> |
| <b>7.2</b>   | <b>Future works .....</b>                         | <b>56</b> |
| <b>Acknowledge.....</b>  |   | <b>65</b> |
| <b>Reference.....</b>  |   | <b>59</b> |
| <b>Appendix 1 .....</b>  |   | <b>67</b> |
| <b>Appendix 2 .....</b>  |   | <b>69</b> |

Degree papers are in the "[Xiamen University Electronic Theses and Dissertations Database](#)". Full texts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to [etd@xmu.edu.cn](mailto:etd@xmu.edu.cn) for delivery details.

厦门大学博硕士论文摘要库