

密级：公开

机器翻译分词研究

中国科学技术信息研究所

2011 年 6 月

机器翻译分词研究

Research on Chinese Word Segmentation in Machine Translation

课题组长：王惠临

课题成员：石崇德 何彦青 张均胜 刘丹

中国科学技术信息研究所

摘要

本研究对中文分词研究的历史、现状、难点以及主要切分方法进行了广泛调研，并对机器翻译中使用断字、机械分词和字标注分词等不同分词方法进行了对比实验，通过分析和统计不同切分方法对机器翻译的影响，总结和归纳了包括总分词数、分词词表数和分词错误等影响机器翻译的主要因素，以作为机器翻译中分词优化方向的指导。

Abstract

This research focused on how different Chinese Word Segmentation (CWS) influences machine translation. We firstly broadly examined the history, current state, difficulties, and mostly used paradigms of CWS. Then we compared different CWS methods like char-based, dictionary-based and char-labeling-based segmentation and how they worked in statistical machine translation systems. According to different impact of these methods in machine translation, we conclude that three factors in different segmentation play the most important part in machine translation: total number of segmented words, total number of words in vocabulary and segment error. These three factors is our direction of optimize CWS in machine translation.

目 录

一、分词研究现状与不足.....	1
1.1 分词的影响.....	1
1.2 中文分词的困难.....	2
二、主要分词方法.....	3
2.1 机械分词方法.....	3
2.2 基于统计和机器学习的方法.....	4
三、机器翻译中的分词研究现状.....	5
四、翻译对比实验.....	6
五、结果分析.....	8
5.1 总分词数.....	8
5.2 词表词数.....	9
5.3 分词错误.....	10
总结.....	13
参考文献.....	13