

John Richardson (黒橋教授)

「Improving Statistical Machine Translation with Target-Side Dependency Syntax」
 (目的言語側の依存構文による統計的機械翻訳の改善)

平成 28 年 9 月 23 日授与

Machine Translation (MT) is the application of Natural Language Processing that focuses on the automatic translation between languages. Translation is particularly challenging for language pairs with widely different grammatical structures, such as English and Japanese. Syntax-based MT is a translation paradigm based on the principle of generalizing language with grammatical rules. This additional layer of abstraction enables the design of more robust and flexible translation rules. However, the majority of previous approaches to syntax-based MT have employed only source-side grammar (known as ‘tree-to-string MT’). This is mainly because syntactic analysis is difficult, prone to error and resulting systems can become overly complicated.

While there have been previous studies on exploiting target-side syntax (‘tree-to-tree MT’), results have not been promising. Our aim is to analyze in detail the effectiveness of target-side syntax in the modern world of machine translation. We ask whether the potential improvement in translation quality is able to outweigh the increased complexity of employing a structured target-side representation (in particular, dependency parses).

This thesis begins with an overview of machine translation, outlining the major paradigms and methods of evaluation. We continue by outlining the case study of a state-of-the-art dependency tree-to-tree system, KyotoEBMT (see Figure 1), which we have been developing as a core component of our research on syntax-based MT. The design and extraction of dependency tree-to-tree translation rules are discussed. Analysis of the system gives empirical evidence of the advantages and disadvantages of syntax-based approaches and provides a starting point for our investigation.

We proceed to analyze two major aspects of translation where target-side syntax can be effective: word order and translation fluency. We discuss our approaches to each of these areas, describing experiments assessing the effectiveness of our proposed approaches and discussing the potential impact of each method.

While this thesis concentrates on statistical syntax-based approaches, the field has recently seen a surge in interest in translation methods based on neural networks. The final chapter presents an overview of future work that could incorporate ideas from this paradigm. We conclude by discussing the potential impact and future directions of our work.

The screenshot shows the KyotoEBMT interface. At the top, there are flags for Japan, China, and the UK. Below the flags, there are two text boxes. The left box contains the Japanese sentence: 「本稿では依存構造に基づく用例ベース機械翻訳システムを紹介する。」. The right box contains the English translation: 「Example based machine translation system based on dependency structure are introduced in this paper .」. Below the text boxes is a button with the text '>>'. Underneath the button is a table with two columns: the source sentence and the target sentence. The table is titled '*** Input and Output Dependency Trees ***'. The source sentence is: 「本稿では依存構造に基づく用例ベース機械翻訳システムを紹介する。」. The target sentence is: 「Example based machine translation system based on dependency structure are introduced in this paper .」. The table shows the dependency trees for both sentences. The source tree has nodes: r[0] 本稿, r[0] で, r[0] は, r[6] 依存, r[5] 構造, r[5] に, r[5] 基づく, r[4] 用例, r[3] ベース, r[2] 機械, r[2] 翻訳, r[1] システム, r[1] を, r[0] 紹介, r[0] する, r[7] 。. The target tree has nodes: r[4] an*, r[4] example, r[3] based, r[2] machine, r[2] translation, r[1] system, r[5] based, r[5] on, r[6] dependency, r[6] structure, r[5] .*, r[0] are*, r[0] introduced, r[0] in, r[0] this, r[0] paper, r[7] .*. Below the table is a section titled '*** List of Used Translation Examples ***'. It contains two examples. The first example is: 「[0] NICT JE SP-train-G-0654753」. The source sentence is: 「本稿では依存構造に基づく用例ベース機械翻訳システムを紹介する。」. The target sentence is: 「Example based machine translation system based on dependency structure are introduced in this paper .」. The second example is: 「[1] NICT JE SP-train-R-0064303」. The source sentence is: 「CADは、CAMを説明する。」. The target sentence is: 「CAD explains CAM.」.

Figure 1: An example translation with KyotoEBMT showing translation examples with dependency trees.