

衛星放送・英日機械翻訳システムの辞書整備

6B-3

畑田のぶ子* 江原暉将* 山本ゆき**

(*NHK 放送技術研究所 **NHK 情報ネットワーク)

1. はじめに

NHKでは衛星放送に英日機械翻訳システムを試用している。ここでは翻訳者が専門辞書に不足な用語を登録し、翻訳作業をおこなっている。約1年5ヶ月間に翻訳者が登録した3129語について分析した。現在この分析をもとに辞書整備をおこなっている。分析内容と辞書整備に際し、考慮すべき課題について述べる。

2. 衛星放送翻訳システムと辞書整備の概要

衛星放送で機械翻訳の対象となっているニュースは、文化的なトピックスが多く、一日あたり約50文、放送時間は5分である。

現システムでは翻訳者が名詞、動詞、形容詞、副詞を辞書登録でき、登録された語はすでに組み込まれている語より優先的に翻訳に使用される。

辞書整備の概要を図1に示す。今回は翻訳者による用語登録を中心に報告する。

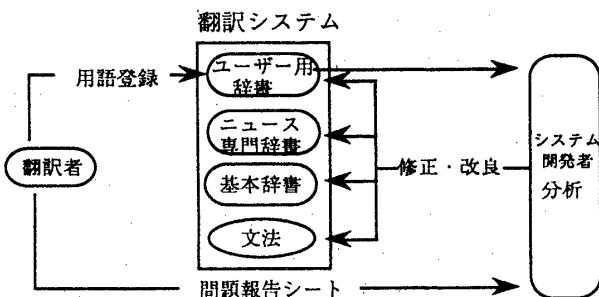


図1. 辞書整備の概要

3. 登録語の分析

1) 登録語の品詞別状況

品詞別状況を図2に示す。

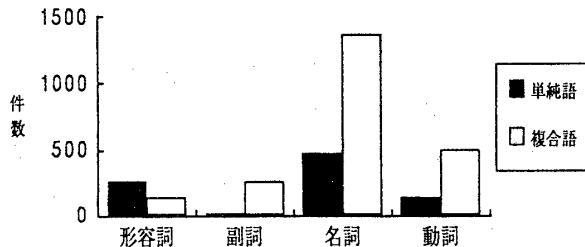


図2. 登録語の品詞別状況

翻訳者の登録語数は全体で3129語、その内58%が名詞、20%が動詞、13%が形容詞、9%が副詞であり、名詞が非常に多い。登録語の内70%は複合語である。

2) 登録語(見出し語)のカバー率

見出し語で現辞書でカバーされていないいなかった語(=未知語)の割合は76%である。複合語にその割合が多く、93%ある。

品詞別割合でみると、未知語であったものの割合は副詞に多く、84%である。これは種々雑多な用語を副詞として登録しているためである。名詞、動詞は77%、形容詞は60%である。

3) 登録語の品詞・訳語のカバー率

見出し語がすでに現辞書にあるにもかかわらず、新たに登録されたものは、現辞書の品詞の不足か訳語の不足が主な原因である。

品詞不足のものは形容詞に特に多く、約9%ある。動詞1.4%、副詞0.7%、名詞0.5%である。訳語不足のものは、形容詞18%、動詞、副詞、名詞12%である。

見出し語がすでに辞書に登録されていて、同一の品詞や訳語があるにもかかわらず、翻訳者が再登録しものについては下記の理由が考えられる。

- ・ 訳語の順位やデフォルトの訳語の設定が翻訳者の希望するものと異なっていた。
- ・ 構文解析により満足する訳語が選択されなかった。
- ・ 辞書データの不備により、意味マーカ等の設定が不足または不適當であった。

4) 登録語の品詞別特徴

形容詞

見出し語が現辞書にあるものの割合が一番多く現辞書の品詞・訳語の不足が原因と思われる。複合語の登録の割合は非常に少なく、他の品詞と比べると割合が逆転している。

- ・ ハイフンを含む語が96件あり、名詞に次いで多い。(ex. adult-to-adult (a):成人間の, family-run (a):家族経営の)

Dictionary Refinement of English-Japanese MT System for Satellite Broadcasting

N. Hatada (1), T. Ehara (1), Y. Yamamoto (2)

(1) NHK Science and Technical Research Laboratories.

(2) NHK Joho Network.

- ・...ed, ...ing の分詞形の見出し語が72件あり、基本辞書にその形での見出し語のないものは59件である。
- ・固有名詞関連の見出し語は58件ある。
- ・前置詞ではじまるものが24件ある。
(これは前置詞句を叙述形容詞として登録しているためである ex. in good health (a):良い健康状態, in question (a):問題になっている)

副詞

- 現辞書にない語の割合が一番が多く、84%を占め、登録語の95%が複合語である。
- ・時を表すもの、場所を表すものが多く、それぞれ133件、19件ある
これらの中には一部分が変化するパターンが多く、月、曜日、季節、人称代名詞、場所などが変化する。(ex. late September (adv):九月末に, late Sunday (adv):日曜日遅く, later this month (adv):今月末に)
 - ・前置詞で始まるものは147件である。

名詞

- 登録語全体の58%を占め、それに伴って未知語の数も多い。
- ・固有名詞、およびそれに準ずるものも多く、304件ある。
 - ・定冠詞で始まるものも多く162件ある。
 - ・ハイフンを含む語は品詞の中では一番多く、107件ある。

動詞

- 名詞に次いで登録語数が多く、登録語の20%を占める。
- ・複合語490件は基本辞書に対してはすべて未知語である。ニュース専門辞書でカバーされているものが29件ある。
 - ・前置詞または副詞小詞で終わるものが166件あり、現辞書に動詞+前置詞または副詞小詞の共起として登録されているものが68件ある。(ex. bring back (vt):持ち帰る, close down (vt):閉鎖する, point out (vt):指摘する)
 - ・上記の他に目的語、その他を伴っていわゆる慣用句として登録しているものが292件あり、基本動詞と関連するものも多い。
(ex. get a haircut (vi):散髪する, get sick (vi):病気になる, get cold feet (vi):おじけづく)
 - ・基本動詞を含むものが121件ある。

4. 辞書整備の課題

1) 辞書の拡充

翻訳者の登録語はニュース関連の未知語を多く含む。これらの用語は順次システム辞書や専門辞書に組み込むことによって、辞書の拡充が可能である。

また、翻訳者が登録した語で翻訳にうまく反映されない語、反映されるが他に悪影響を及ぼすものについては原因を調査し、辞書の登録方法や翻訳文法の改良で対応する必要がある。

2) 用語について

調査した用語の中には、(a) ..ed, ...ing の分詞形をもつ見出し語、(b) 動詞で前置詞または副詞小詞で終わるもの、(c) 基本動詞と関連するもの、(d) 人称代名詞、指示代名詞、月、曜日、などの一部分が変化するパターンをもつものがある。

(a) については動詞の基本形と屈折形の両方が見出し語になっている場合があり、その時は優先順位の高いユーザー用辞書に屈折形を登録するのは問題がある。このような語は基本辞書に移し、ウェイトの調整をする必要がある。

現在は、基本辞書中にある動詞で...ing 形を見出し語にもつ用語をサンプリングし、人手で動詞、形容詞、名詞間でのウェイト調整の実験をおこなっている。この結果を用いて自動的にウェイトを調整する手法が望まれる。

(b) のうち句動詞については動詞と副詞小詞の共起として登録し、目的語が副詞小詞の前にも良いようにすべきである。

(c) については基本動詞の訳し分けに関連する部分と慣用表現に関する部分がある。基本動詞の訳し分けとして広くとらえるべきである。

(d) 一部分が変化するものについては、パターンを利用し、(i) ローカルな文法処理を組み込む方法(ii) 辞書に登録すべき用語を自動生成する方法について検討中である。

5. おわりに

衛星放送・機械翻訳の辞書整備は今回の述べたものと、図1の下側で示した問題シートの分析を用いたものがある。後者については別途報告したい。今後も4.の課題を基に辞書整備、文法の改良を行なっていく予定である。

参考文献

- 相沢、他：放送ニュースへの機械翻訳システムの適用、信学会、NCL91-20, 1991
H.Tanaka: A Method of Translating English Delexical Structures into Japanese, COLING 92, 1992