

日英機械翻訳用スケルトン-フレッシュ型構文意味辞書の構成

6Q-8

横尾昭男 中岩浩巳 白井 諭 池原 悟

NTT情報通信網研究所

1. はじめに

筆者らは日英機械翻訳の研究を進めており、日本語の意味構造を解析して英語の意味構造に変換するために構文意味辞書を構築してきた。従来の辞書記述方式では、1つの意味表現構造として1つのパターンですべてを記述していたが、ここでは、そのパターン構造を骨格部分(スケルトン)と肉付部分(フレッシュ)に分離して辞書に格納する方式を提案する。

本方式により、既存の約13,000パターン*1を再構築した結果、構造の持つ意味の共通化を行うことができ、少数の骨格構造で大多数のパターンを記述できることを示す。

*1:当初15,000パターン作成したが、その後縮退などにより、現状では約2,000パターン減少している。

2. 構文意味辞書の構成

NTTで研究を進めている日英機械翻訳システムALT-J/Eでは、「構造は意味の一部である」という考えに基づき、単位文を変換するときには、「多段変換方式」に従った構文意味辞書を構築してきた(図1) [1], [2]。

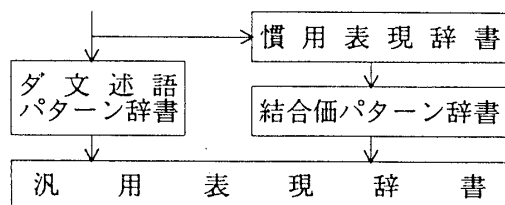


図1 多段変換用構文意味辞書

例として、「～が～を～に提案する」という文を変換する結合価パターンを図2に示す

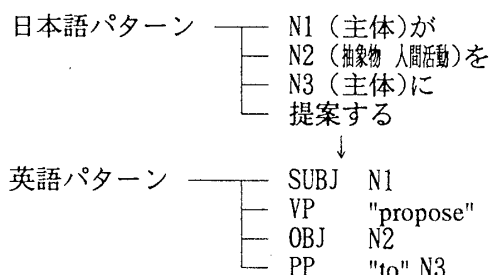


図2 「提案する」の結合価パターン

この英語パターンは、全体を1つの構造とするので、全体を扱いやすいという利点はあるが、複雑な構造になると全体を捉えにくくなり、共通部分について重複記述する必要があるので、全体の記述量が増えるという問題がある。

3. スケルトン-フレッシュ型構文意味辞書

3.1 構文意味辞書の構成法

英語パターンの重複を考慮し、全体を効率よく記述するため、図2における英語パターンを骨格部分と肉付部分に分離して記述する方式を提案する。骨格部分と肉付部分は、それぞれ以下の特徴を持つ。

骨格部分(スケルトン) … 構造の意味を表す

- ・ 述語、格要素(前置詞句も含む)、副詞句、修飾句節の構造を表す
- ・ 単語の品詞も記述する

肉付部分(フレッシュ) … 要素の意味を表す

- ・ 格要素の変数(Nx)へのポインタ、単語の字面等を記述する

骨格部分は、図1における各辞書共通に1本作成し、それを新英語パターン構造辞書と呼ぶ。肉付部分は、各辞書ごとに対応づけて作成し、全体をまとめて新英語パターン辞書と呼ぶ。

図2における「提案する」に対応する骨格部分と肉付部分を図3に示す。

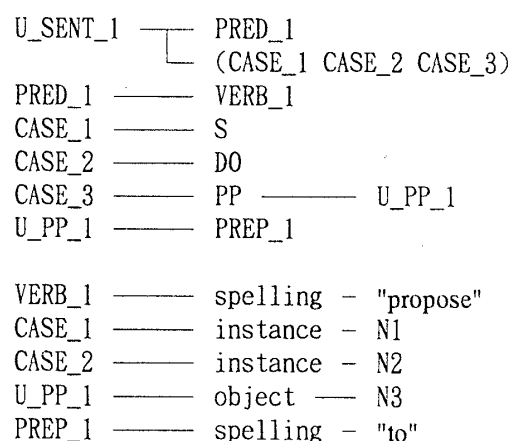


図3 スケルトン-フレッシュ型構文意味辞書の骨格部分(上)と肉付部分(下)

3. 2 骨格構造のカバー率

3. 1 で示した考え方にに基づき、既存の英語パターン辞書（結合価パターン、慣用表現辞書、ダ文述語パターン、計 13,225）を再構築した。その結果、全体を 446 種類の骨格構造(スケルトン)で記述することができた。骨格構造の出現頻度の高い方から見た順位とカバー率（そこまでの骨格構造の累積で、全パターンのどれだけを表現できるかの百分率）の関係を表 1 に示す。

表 1 骨格構造の出現順位とカバー率の関係

出現順位	9	17	40	95	214	314	446
カバー率	70.6	80.7	90.0	95.0	98.0	99.0	100

骨格構造のうち出現頻度の高い10種類について、その骨格構造、出現頻度、頻度累計、カバー率を表 2 に示す。また、頻度順上位 100 までの骨格構造のカバー率を図 4 に示す。

表 2 骨格構造(スケルトン)の上位 10 構造
(全パターン数: 13225)

順位	骨格構造(スケルトン)	頻度	累計	カバー率
1	Nx Vt Nx.	2622	2622	19.83
2	Nx Adj.	1893	4515	34.14
3	Nx Vt Nx prep Nx.	1436	5951	45.00
4	Nx Vi prep Nx.	1083	7034	53.19
5	Nx Adj prep Nx.	981	8015	60.60
6	Nx Vt noun prep Nx.	406	8421	63.67
7	Nx Vi.	402	8823	66.71
8	Nx Vt noun.	272	9095	68.77
9	Nx Vi prep Nx prep Nx.	240	9335	70.59
10	Nx be Vpp prep Nx.	231	9566	72.33

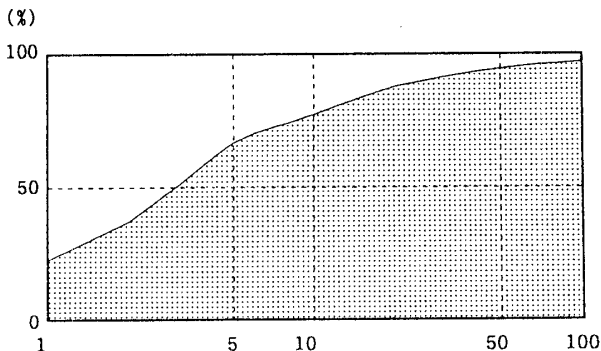


図 4 上位 100 の骨格構造のカバー率

以上から、少数の骨格構造により効率よく英語パターンを記述できるということがいえる。辞書メンテナンス・ツールでは、比較的出現頻度の高い 52 種類の骨格構造を「高頻度骨格構造」としており、これらで 12,084 パターン、91.37% をカバーしている。

4. 利用者用構文意味辞書への応用

ALT-J/E では、結合価パターン辞書は、一般的な用言を格納した一般パターン辞書と、分野に依存した用言を格納した専門パターン辞書で構成される。専門パターン辞書は、ある意味では利用者が自分で作成する利用者パターン辞書と捉えることができる。そこで、専門パターン辞書（全体で 206 パターン）について、骨格構造の出現頻度とカバー率を調査した。その結果を図 5 に示す。

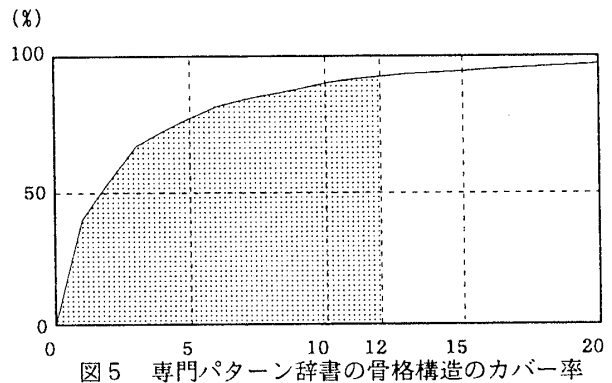


図 5 専門パターン辞書の骨格構造のカバー率

上図のうち、頻度順 12 位までの骨格構造(カバー率: 92.7%, 網掛け部分)は、高頻度骨格構造に含まれており、利用者が個人用のパターンを作成するときには、既存の骨格構造を参照しながら、格要素の変数名や単語の字面だけを指定すればよいので、新規登録の効率化を図ることができた^[3]。

5. おわりに

本方式により、既存の 13,225 のパターンを 446 の骨格部分に分類することができ、出現頻度が 40 位までの骨格構造で全体の 90% をカバーすることが分かった。また、利用者が構文意味辞書を登録する場合もこの 40 の構造に入ることが多く、利用者辞書を構築する上での負担を軽減できることを示した。

また、日本語の単一用言文が英語の 500 文型程度で記述可能なことが分かったので、今後は、本記述形式による辞書の英文解析への応用についても検討する予定である。

〈参考文献〉

- [1] 白井, 池原, 横尾, 中岩: 「前編集不要の日英機械翻訳の実現に向けて」, NTT R&D, Vol.40, No.7, pp.897-904, 1991
- [2] 池原, 宮崎, 横尾: 「日英機械翻訳のための意味解析用の知識とその分解能」, 情処論, Vol.34, No.8, pp.1692-1704, 1993
- [3] 白井, 横尾, 池原, 井上: 「日英機械翻訳用構文意味辞書の記述精度の向上と作成支援」, 48 情処全大 6Q-9, 1994