

機械翻訳用辞書構成

村木 一至
(日本電気(株) CoCoシステム研究所)

1. はじめに

実用的な言語処理システムには、大規模辞書が必要である。そのためには、辞書内容を均質に且つ矛盾なく開発・管理するための方策が必要となる。また、機械翻訳の如く、複数言語を対象とするには、複数言語間の諸規則に関する対応の管理も重要な課題となる。

本稿では、日英機械翻訳システムENVUSの諸規則辞書構成（三層辞書構成・形態・語彙・意味辞書）とその管理について報告する。²⁾

一般的に辞書と言ふとき、様々な内容の、様々な目的でもう辞書を含む。言葉の持つ内容、言葉の言語的性質、特定語を関連する言葉について書かれた多種類の市販辞書がある。

言語処理用機械辞書にしても、ワード用カナ漢文換辞書、言語解釈用文法情報辞書、情報検索用類語辞書など多種類存在する。また、カナ漢文換用辞書にも、その対象によって人名辞書、地名辞書などが存在する。これらは様々な辞書で、その目的と並んで直接扱う機械（通常ソフトウエア）の違いにより、内容、形式が異なる。

機械翻訳用辞書は、原言語解釈・目的言語生成辞書と、原言語と目的言語対応付け辞書（はし機能）を全く含んでない。この辞書内容は、翻訳方式や言語解析・生成方式によつて異なる。また、その構成も処理方式によつて大きく違つて見える。

言語処理や翻訳方式（システム）に依存して大規模辞書を開発しようという研究が報告されている。¹⁾ ところが、この辞書開発方針は、共用利用者辞書と、目的機械辞書をトランシーラーにより対

応付けようとする。トランシーラーを多種類用意することは、同一言語に関する解析・生成用辞書を共用辞書から自動生成する。たゞえ翻訳方式の差違があつても、各翻訳システム用トランシーラーを用意して対処する。

機械翻訳の言語対応付けに於けるイデオム記述の問題を解決しようとする報告がある。³⁾ ところが、單語辞書の1エントリー内が、見出し語に因る多生語、イデオム、形態情報、属性情報、意味情報を記述するためにはどのように構造化すべきかを述べている。

ここに提案する翻訳用辞書構成は、更に以下の問題に対する解決を目的とする。

- 1) 同義語記述
- 2) 多言語間対応記述
- 3) 記述塵覆の解消
- 4) 多様な応用への適合性

その他、5) イデオム記述との処理の簡素化、6) 解析・生成といった処理方向からの独立性を可能とする。

更に、一旦開発した辞書構成に対する辞書管理システム（辞書エディタ）は上記6項目の要請を満足させる辞書を保守・管理するため必須である。

以下の節では、本稿提案による辞書構成が生まれた背景を説明し、辞書構成、管理システムについて言及する。

2. 機械翻訳用辞書構成の要件

本節では、筆者等が開発している機械翻訳システムの説明を通して、辞書構成に關し何が要求されるのかを明らか

大可る。

2.1 VENUS機械翻訳システム

VENUS翻訳システムは、入力日本語文を形態素解析し、単語同定を行つ。この処理には、熟語判定や未定義語処理を含んでゐる。形態素解析は、単語列を生成する。

意味抽出は、入力単語列から文法情報と意味情報を用いて意味構造と呼ぶ木構造を生成する。木構造のノードは概念名（記号）と意味関係名からなる。

概念抽出は、意味関係のうち入力言語の文体を担う特殊な関係を概念関係（事実内容を記述する有意味子関係）と置換し、概念構造を抽出する。

二から三つの処理は、入力言語内に同じくある。之して概念構造を記述する概念記号（概念名と概念関係）は、入力文の意味内容を記述する記述子であり、各言語共通の意味を代表する。例えば、図1の概念構造中の記号“LOCATION”, “TOWER”, “TOKYO”は日英両言語で同一意味内容をもつ。

目的言語内に同じくある処理としては、文体選択、構文生成、形態素生成の3モジュールがある。この生成処理は、概念記号に対して目的言語内の対応する単語表現を辞書により見出し、文法的に説明すれば文表現を得る。

文解析・文生成処理は、概念構造をもつて基本的2分される。この概念構造は、各言語の文体の差違を吸収し、事実内容を充分表現できるよう、全ての概念記号が何らかの表層の単語に対応付けられる。即ち、図1例中の“LOCATION”を代表する関係概念も“ある”, “存在する”, “の(東京のタワー)”, “in”, “exist”, “place”等の表現と対応付けられる。これらを二つに分けて、トランシスファー一方或はトランシスファー処理の対象となる文体变换・構文变换の多くが不要となる。

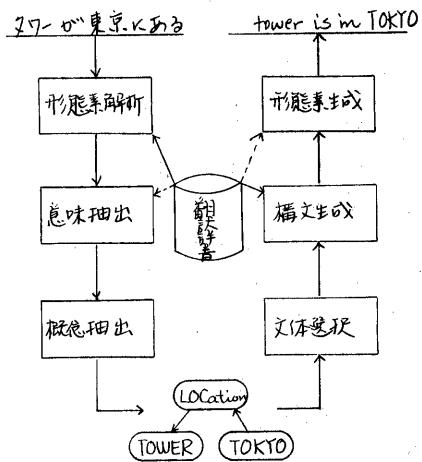


図1. VENUS翻訳手順

2.2 語彙辞書構成の要件

VENUS翻訳システムの処理内容は、翻訳システムやその翻訳方式が異なる場合も何らかの手段で処理可能である事柄である。つまり形態素処理、構文処理、意味・概念処理機能は、多くの翻訳システムの必須機能である。之にて、翻訳用辞書構成に対する要件は以下の通りである。

<要件3 言語の対応付け>

翻訳システム用語彙辞書は、異なる言語間の対応付けを語彙レベルで記述する。VENUSでは、概念記号が併立するなり、日本語の単語と英語単語（より正確には両言語の形態素）を対応付ける。ちなみにも、トランシスファー方式では、日本語単語と英語単語の対応辞書（トランシスファー辞書）がこの対応付けを管理する。この対応は、一般的に対応あり、この対応関係を保持・管理することは、訳語選択の質を向上させるうえでの鍵となる。とくに、辞書開発後翻訳システムに

組み込み、実行評価を行う際に誤つて対応を検出するとか、新たに新しい語を追加する場合には、対応の明示が必要とする。

辞書が2言語以上の対応を保持する必要がある多言語翻訳では、この対応管理はより重要な課題となる。つまり、言語Aの単語αから言語Bの単語βへ対応付けられ、言語Cの単語cから言語Aの単語αへ対応付けられたり、言語Bの単語βは単語αとCへ対して正しい対応を保持していることをどのように検証するか。

<同義語管理>

機械翻訳では、訳文・訳語の均質性（同じくまぎ）で訳文や訳語が存在しないが利用者から期待される複数の同義語が1つの訳語に割り当てられることがある。逆に1つの単語が複数の同義語に対応付けられることがある。とくに單一言語辞書が解析生成の両処理不通用される場合、上記の例で、1対の対応が翻訳辞書を同時に存在し得る。

訳語の統一性・均質性を保持するためには、既対の同義語対応が明示的表現され、その同義語が均質的であるよう保証する手段が必要である。とくに、既に登録のために単語と同義の単語を辞書登録する場合以前に付けた訳語との対応を無視すれば、訳語の不統一や不均質が生じる。

<辞書記述の重複排除>

膨大な辞書情報を記述するのに、全く同一の内容を重複して記述することを回避すべきである。市販のMT化辞書がそのまま利用されることは現状では、負を落として記述量を減らせることが望まれる。例えば、多品詞語も見出し語は1回記述すればよく、同一品詞多義語は、見出し、品詞情報は1回の記述不充分である。また、多品詞語でも意味が同じなら、意味記述は1

回のみでよい。同品詞・同義語ならば、更にし以外は重複して記述してください。

辞書記述重複排除は、單に記述努力を軽減する目的より、記述の誤りの最小化を大にする目的とする。大量の語彙に対し、同一内容を何度も記述すると誤りの発生を避けられない。そのためには、初期辞書記述（登録）段階から何らかの方策を用意しておけばならない。

<不連続語彙登録>

言葉には、不連続な単語（熟語の1部）がある。通常熟語辞書という特別な辞書として開発されており、不連続語彙の中心語基の付属成分として同語基辞書内に記述される。しかし、処理の観点から視ると、このように単語の同定には、通常の分りきし日本語の場合）や派生語認定のための形態素解析（英語の場合）以外の辞書引き動作が必要である。更に、英語不連続単語のように、活用（形態）処理以前と以後両時点の語同定を行う必要が生じることがある。

また、特定語基の付属成分として記述する場合にも、付属成分を入り文中で検索する特別な処理が必要である。

こうして不連続語の辞書記述は、單に記述されるだけではなく、処理手順の簡素化にも何らかの配慮が必要となる。

以上述べた要件の他にも様々な検討課題がある。機械翻訳以外の多様な応用に適し適用可能な内容と、適用可能な構造を持たせる。辞書の管理・維持が容易である。そして、その管理・維持方式がシニアルである。すなはち、辞書記述能力を反える辞書構成方式から、その辞書を機械によつて維持・管理する強力な機能やシニアルを便易・易い辞書管理システムとのインターフェースを満足しなければなら

は、条件は多い。

3. 語彙辞書構成

本節では、語彙辞書の構成法とその特徴について述べる。

3.1 三層辞書構成

三層辞書構成は、形態素辞書、構文辞書、意味(概念)辞書の3つの辞書により構成される。各単語の辞書記述も、3種の情報に区分される。三層辞書の物理構成の概略を図2に示す。

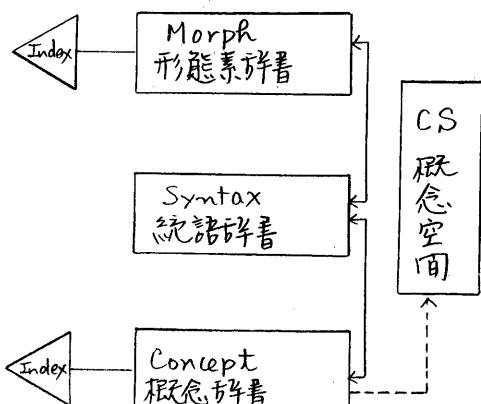


図2. 辞書物理構成

形態素辞書と概念辞書は、各自形態表記と概念記号を見出しつつ検索できることができる。形態素辞書と構文辞書は互いに双向の indexing で結合され、構文辞書と概念辞書間も双向の indexing で連結される。

この三層構成によれば、形態素見出しを与えることにより、形態辞書、構文辞書、概念辞書が検索でき、逆に概念辞書から構文辞書、形態辞書が検索できる。この検索過程は、入力文から概念記号を表現して概念構造を抽出し、概念構造から目的文を生成する ENVS の翻訳過程と論理的に対応する。

併み、1つの形態辞書項目は一般的

複数の構文辞書項目に対応し(多品詞語)、逆に1つの構文辞書項目も多数の形態素辞書項目に対応する(要表記)また、1つの構文辞書項目が複数の概念辞書項目に対応したり(多義語)、1つの概念辞書項目が複数の構文辞書項目に対応したり(同義語)。とくに概念辞書項目から形態素辞書を視ると、複数言語に対する単語の対応が鮮明に理解できる。図3は、これらを模式図で表す。同図は於て、M, S, C, すなはち形態素辞書項目、構文辞書項目、概念辞書項目などを示す。例えば、概念記号「REASON」

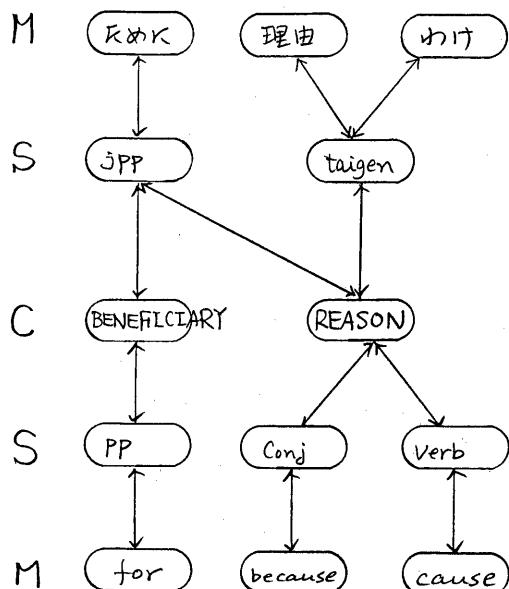


図3. 要言語間の対応模式図

は、構文機能も表記も異なる日本語・英語の表現とどう対応するかが簡単に読み取れる。

3.2 辞書記述

辞書構成が、3種の辞書を如何に物理的に関連付けるかを規定する。その内容と之の構成にどんな意味付けを行う

かは、辞書記述の問題である。

辞書記述に於いてはまず次の制約を設けた。すなはち形態素辞書項目と対応を持つない構文辞書項目は存在しない。つまり、構文機能(構文機能)は、表裏記(形態素辞書見出し)に対し記述するということがある。この制約は、言語処理の直接の対象となる形態素に対し辞書記述を行うという自然なものである。

形態素見出しが、表裏記と発音(読み)からなる。裏表記・同一單語(形態素)に対しては唯1の構文辞書項目を与える。

1形態素に因し、複数の構文機能(例えば品詞)が存在する場合には、複数構文辞書項目に対し唯1の形態素辞書項目を与える。

構文辞書と形態素記述の関係は、構文辞書と概念辞書記述の関係と相似である。つまり、多義語、同義語(言語を越えて同義語を含む)に対する概念辞書項目、構文辞書項目は各自唯一の構文辞書項目、概念辞書項目に対する

二のようだ; 3種の辞書項目にまた1つの語彙を記述する手法に於いては、

- 記述事後の回避と、それによる記述誤りの局所化が図れる。
- 概念辞書を中心とした同義語の明示的管理と、それによる同義語ひうしの記述内容の矛盾排除が図れる。
- 概念辞書を中心とした対訳同義語管理と、それによる誤語品質劣化の防止が図れる。
- 概念辞書を中心とした多数言語にまたがつて語彙対応管理の一元化が図れる。

などの特徴がある。

こうして、

- 各辞書の機能を細分化することにより、簡単な熟語辞書を記述不^可、通常の語彙辞書記述と同じ構造を管理不^可。

通常熟語といういくつかの種類がある。まず、複合語と呼ぶ連続したものの複数の語からなる不連續語(まことば)に区分される。そして、名詞は活用するものと不活用のものに区分される。三層辞書構成に於いて、不連續・活用型熟語を因すべく記述不^可。

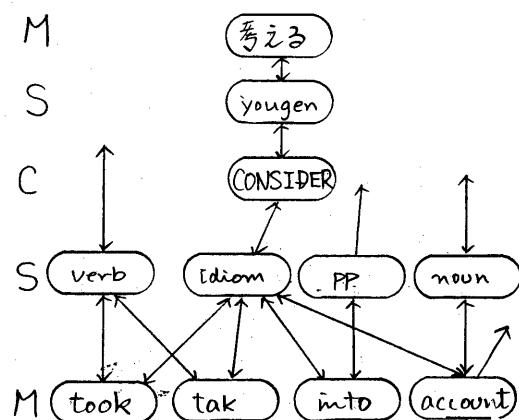


図4. 不連續単語の記述

同図に於いて、構文辞書Sの項目「Idiom」は、構文辞書機能を拡張して熟語保持項目である。この項目には、「考慮す

take/takes/taking/took/taken
-into
-account

る」という意味の熟語が記述されていふ。つまり、以下の特性が言える。
"take", "took", "into", "account" などが熟語要素であると同時に、各自独立して構

文辞書を持つ。故に、二の形態辞書を用いて分かれ書きや形態素解釈時に熟語候補を選びぬくことが出来る。
また熟語認定には、「idiom」項目を中心条件記述することにより対処出来る。
即ち、熟語認定のため再度熟語辞書をアクセスしたり、熟語付属成分を入力文中に探索したりといった無添が省ける。

4. 辞書管理

辞書管理の重要な機能は、辞書検索機能と、辞書項目の登録および削除などの更新機能である。

辞書検索は、形態素見出しと概念記号をB-treeインデックスとして実現され、他の辞書項目へは、関連辞書項目への接続を取ることにより到達出来る。

辞書更新には、辞書管理エディターと、辞書登録ソフトウェアツールを用意している。辞書登録ソフトウェアツールは、書式に従つて書かれた各辞書項目を一括登録する機能を持つ。
逆に対し、辞書管理エディターはオシライン端末により会話的で登録・削除を行うために用意した。図5は、同エディターの概念図を示す。

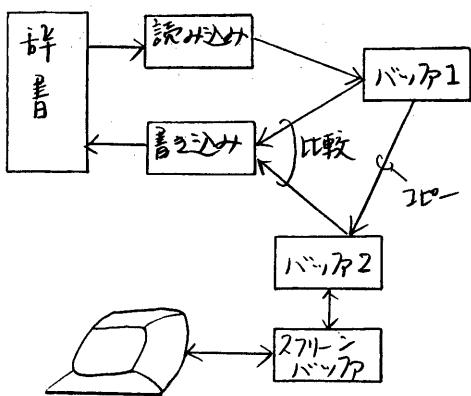


図5. 辞書管理エディター概念図

辞書管理エディターは、基本的形態素見出しと概念見出し（概念記号）の指定による辞書アクセスを提供する。
以後は、どの辞書項目をアクセスするかを、接続された項目のうち画面で表示されるメニューに従い利用者が指定する。

指定された見出しが存在しないときは、登録モードが、見出しが既に存在するときは更新モードが自動的に選択される。
エディター画面のメニューにて従ってアクセスした全ての辞書項目はバックアップされ、エディット結果を保持するバックアップと比較して後更新される。
また、エディター画面はスクリーン機能を持ち、スクリーンドップルクリック管理される。

辞書管理エディターは、之を用いて、接続された辞書項目へもアクセスでき、アクセスした辞書項目が選択可能になる。
このように強力な機能は、一旦開発して辞書内容を比較し記述の誤りを発見訂正するための有力な支援となる。

辞書管理エディターの直感的理解のために、次回6.a～dはエディター端末（PC9800）画面を例示する。
同図は、概念記号「PROCESS」を指定して更新モードで各種辞書項目をアクセスする画面である。

図6.aは、概念記号「PROCESS」を指定してとの画面であり、「PROCESS」に対する形態素表記が同時に表示される。

図6.bは、前画面より形態辞書へのアクセスを指示した結果であり、形態素の選択を得る。

図6.cは、前画面で形態素「刈穂」を選択し、その内容が表示される。

| | |
|---------|-----------------------|
| PROCESS | CCS. (TAG1, OBJ, INS) |
| | |
| | |
| | |
| | |

処理
プロセス
加工
加热処理
PROCESS

IN Conceptual Screen

block number 16289
Index number -1

RUP RDOWN SASK SKIP SAVPTON CPY BF LPT LF

図 6.a 概念辞書項目画面

| 表層 | | 形態 | |
|----|--|----|--|
| A | | F | |
| B | | G | |
| C | | H | |
| D | | I | |
| E | | J | |
| | | | |

IN J. Morph. Screen

KEY INPUT OR NUMBER INPUT (space/figure) =
0 = 処理 1 = プロセス 2 = 加工 3 = 加熱処理

MORPH SYNTA CONCE SKIP SAVE LTRL HALF

図 6.b 表層辞書項目選択画面

| 表層 | | 形態 | |
|-----------|---|----|--|
| A 204 234 | I | F | |
| B | | G | |
| C | | H | |
| D | | I | |
| E | | J | |
| | | | |

IN J. Morph. Screen

block number 16291

CP = PROCESS

図 6.c 表層辞書項目画面

| | |
|----------|----------------------|
| Syntax 1 | Conceptual Primitive |
| | PROCESS |
| | MGRPN |
| | U. 处理 |

IN Syntax Screen

block number 16286

CP = PROCESS -> Mono = 处理 ->

図 6.d 構文辞書項目画面

図 6.d は、前画面の「処理」に対する構文辞書項目の表示である。

このようす画面を介し利用者は三層辞書へのインカラフティアはアクセスを行なうことができる。この場合、もし形態素が複数の構文辞書項目を持つとき、このアクセスメソッドはいかにも頻語候補機能に似た役割も果して得る。

また、この管理エディターは新規登録作業を円滑にするための機能として AKO (A-Kind-Of) 機能をもつ。AKO 機能とは語彙登録時に全ての情報を全く新しい記述する代わりに、多くの情報を共有する既登録辞書項目を指定しその内容を編集する。これら新しい語彙に対する辞書登録を容易にする。

5. おわりに

日本機械翻訳システムの翻訳辞書構成とその管理システムについて述べた。構成の特徴は以下の2点である。

- 形態辞書・構文辞書・概念辞書からなる三層構造を導入した。
 - 概念辞書を要是する多言語間を連結するものとして中心に置いていた。
- この2つの特徴から以下の機能が実現された。
- 類語検索
 - 同義語検索
 - 語彙辞書・熟語辞書の統合
 - 重複記述の削除と共により矛盾の最小化

また、この辞書構成により、概念辞書(概念記号)を中心とした多言語辞

書を統一的に管理することができる。
同様に、同一意味内容をもつ語群を一括管理することにより、誤譲品質の劣化を防止する手立てを提供する。

この構成をもつ辞書データベースは VENUS 翻訳システムで用いられており、いくつもの問題点が明らかにされている。されば、辞書が3種類に分割されており、利用者がそれを明確に意識する必要が生じたこと、そのため、利用者の負担が増加してしまったことを感じられることがある。実際、画面や辞書フォームで従つて辞書登録する際、各々の辞書の対応付けを無茶を利用者が感じることがある。これはあくまでも心理的負担であるが、この辞書データベースの利用をより容易にするための方策を考える必要がある。

また、概念記号は利用者が管理し登録するものである。そのため形態素表記群を代表する道筋であることを保障する必要がある。これは、とくに知的作業に属し、実行にはコストがかかる。もし、概念記号を單に両言語の語彙の一対一対応を保持するための道具と做して、1つの翻訳のインデックスとするならば、概念記号の付加作業は全く機械的に行える。しかし乍ら之を行うと、辞書の多言語対応能力が減少し初期要件を満足することができない。同義語(多言語に対するもの)を管理する概念記号付与作業を支援する方策を検討する必要がある。

今後このような問題を解決するため、辞書管理システムの運用を通して実践的な方策を考えていきたい。

最後に、本研究開発の機会と共に
対する支援を与えられた CoC システム
研究所・メディアテクノロジ研究部十
葉部長ならびに山本課長に深謝します。
同部荻野主任には、管理システム設計
に協力していただきました、ニニに感

謝の意を表します。

参考文献

- 1) 中村順一, 他, 「Mu プロジェクトにおける辞書の運用方式一日更変換辞書と英語生成辞書」、「自然言語処理技術シンポジウム」予稿, 1984, 11
- 2) 荻野, 渋田, 「自然言語処理と辞書管理」, 情報大会第29回, 1984, 9
- 3) OKAJIMA, 他, 「Lexicon Structure for Machine Translation - An Example from English-into-Japanese Translator ATHENE」, 1983 International Conference on Text Processing with a Large Character Set