

技術論文表題の英和自動翻訳の試み

長尾 真, 辻井潤一, 建部周二(京大・工)

1.はじめに

機械翻訳システムのように、非常に多くの未解決の問題を持っているシステムを作成していく場合には、まず現在の技術の範囲内で比較的簡単に作成できるプロトタイプのシステムを作り、そのシステムでの経験から、残された問題と、もしそれが解決したとしたら、どの程度のパフォーマンス向上が得られるかを整理し、次に解決すべき問題点を明確にしておくべきであると思われる。特に、工学的立場からは、あまりに詳細な言語現象についての解説を行なっても、それが全体に及ぼすパフォーマンス向上につながらない場合もあることを考えると、プロトタイプシステムの作成と、その問題点の明確化は重要な作業になる。

我々はきのうのようなプロトタイプシステムとして、対象を技術論文の表題にかぎった英和自動翻訳システムを作成した。近年のデータベースシステムの発展と情報検索への要求の高まりによって、論文表題の翻訳への要求は非常に高まつており、日本では日本科学情報センター、外国では米、独、その他の国で具体的にこの種の翻訳への要求がある。

本報告は完全な形で稼動している機械翻訳システムについての報告というよりも、現在我々の作ったプロトタイプシステムで明らかになつた将来の問題点、次に我々が行なうべき課題について述べることにする。

2.機械翻訳システムの概要

自動翻訳は一般に次のようないくつかのステップをへる。

- | | | |
|----------------|--------------------|---------------|
| (i) 原文入力 | (ii) 原文の形態素解析 | (iii) 原文の構造解析 |
| (iv) 原文の意味解析 | (v) 原文情報から訳文情報への変換 | (vi) 訳文の構造合成 |
| (vii) 訳文の形態素合成 | (viii) 出力プリント | |

これらのうち(ii)の形態素解析 (vii)の訳文の形態素合成の部分は、過去にしつかりして研究もあり、やれればほぼ出来る範囲のものである。ここでは研究の焦点を(viii)と(v)にあて、(iv)についてある程度の処理を行なうという立場をとつた。その理由は次の通りである。

(i) 意味の問題にはいくつかのアプローチがあるが、せまく限られに明確な世界について適用はできても、一般的の分野に適用できるような確立された手法はいまだにない。

(ii) 意味を十分にあつかうためには、現実世界に関する知識も採用しなければならない。

(iii) 技術論文の表題の翻訳というような実際的な分野の場合に、(i)(ii)のような意味を積極的にあつかうのではなく、逆にどこまで意味のことを考えずに実用的なところまでゆけるかを試みる方が実際的である。

(iv) 計算機はぼう大なデータの処理にむけていまから、複雑でこつに方法を考えるよりも、種々の場合をなるべく多く集めて、それによって表題の解析と翻訳を考える方がよいのではないか。

このように、論文表題の翻訳を单纯なプログラムとできるだけ多くの場合をお

語では、複数個の単語で1つの概念を表わすことが多い。この場合、誤出の過程で個々の単語の誤語から全体の誤語を合成していくにのでは、ほとんど理解不可能な誤文が生じる。人間の専門家が専門用語の誤語を決定する場合に、個々の英語の単語の誤から合成するのではなく、その句全体が表現している概念を自然な日本語で表現することによって行なっていることを考えると、すなはち、専門用語の日本語訳を決定する際に人間の創意が入っている以上、この過程を計算機に行なわせることはほとんど不可能であり、これらの句も1つの単語と考える必要がある。そこで我々はこれらの句に相当する専門用語も1つの独立した辞書項目として設定することにした。実際には、その句中に出て来る単語の辞書記述中にそれらの句を登録している。

この他、各分野でよく使われる固有の表現 (Idiom 的な表現) — time varying (時間とともに変化する), settling time (整定時間) — あるいは、英語全体にみられる idiom 的表現 (based on —, perpendicular to — 等の辞前置詞) もまた、辞書中に記述され、辞書引きの段階で単語と同じように取扱われる。これを(2)の例に於いて Rigid Idiom (以下 RI) と呼ぶ。

一般に、単語の単位を大きくとればとるほど、辞書の規模は大きくなるが、遂に後の処理が簡単になり、誤文の質も向上する。特に、ある特定の専門分野に対象を固定した場合には、この手法は有効であろう。この Rigid Idiom を選定する過程で、学会発行の学術用語集等を参考にしたが、用語数が少なく、ほとんど後に並んでいた。現実的な機械翻訳システム作成のためにには、現在使われている用語を網羅的に収集して Terminology Data Bank 的なものを作成する必要がある。

(4) 表題文においては、連續した名詞、形容詞によって名詞句を構成する電文調の表現が多くあらわれる。

通常の文体では、前置詞等を使って結合される名詞が、表題文においては、單に連続するだけで結合され名詞句を構成することが多い (例 — High Resolution Bragg Reflection Method, Pulse Compression Distance Measuring Equipment System etc.)。

これらの名詞間、形容詞一名詞間の結合関係を同定し、誤文合成の際に適切な語順変換と、助詞の插入を行なうことほもろん理想的ではあるが、これを行なうためには、かなり深い意味的な解析が必要とし、現在の時点ではほとんど不可能である。

一方、日本語においても、表題文においては、かなり長い漢字連続で英語と同様な表現をとることが可能であり (上記例の場合には、「高分解グラッゲ反射手法」、「パルス圧縮距離測定装置システム」となる)、しかもその中の英語と日本語の語順はほとんど変化しない。

したがって、現在の我々のシステムにおいては、このような長い名詞連続内の単語間の詳細な係り受けの関係は同定せず、单にそのような長い名詞句があつたということだけを簡単な Transition Network (以下 TN) 文法で認識するにした。現在のシステムでは日本語と英語の語順が変化する時は、前置詞、—ed形、—ing形等の単語があらわれた場合だけであると考え、形容詞を含む名詞連続では語順がかわらないと考えている。したがって、ここで使われる TN はこのような語順の切れ目を示す語があらわれるまでを、読みとばして I

おったにデータ、法則とで行なってみようというところに研究の焦点をあてた。したがって入力文、合成文とともに形態素の処理は行なっていはず、構文と変換の部分のみのテスト結果がわかるようになっている。

また、上のような翻訳のメインの流れの他に、我々が今回作成したように実験的なシステムにおいては、辞書記述やルールの変更を含む様々な試みを、できるだけ容易に行なえることが重要になる。したがって、表題文の KWIC、辞書・ルール記述の Editor 等の機械翻訳システムのためのサポート・システムも同時に作成した。全体のシステム構成を図 1 に示す。なお、翻訳のための論文表題は、日本科学技術センターの文献速報磁気テープからとり、1 万件の英文表題と、それに対して人間によってつけられた和文表題とがつづかれている。辞書としては、1 万件の英文表題にあらわれた 9,800 語（形態素処理は行なっていないので、morphological variant もすべて独立して 1 語として登録されている）について、品詞、日本語訳がつけられている。

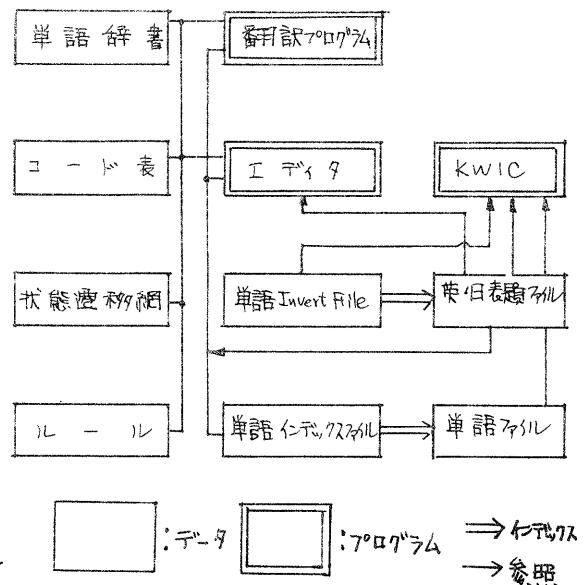


図 1. システム構成

3. 機械翻訳システムの構成

3.1. システム構成の考え方

論文表題は、機械翻訳システムの対象としてみると、次のような特徴を持っている。

(1) 表題文中の単語は専門用語がほとんどで、他の分野に比べて多義語の問題を排けることができる。

もちろん、専門用語であっても各専門分野ごとに別個の誤語がつけられてい場合もあるが、各誤語ごとに分野コードを指定する等、簡単な措置で多義語を排けることができる。ただし、誤語と分野コードの指定は現在のシステムで不行なっていい。

(2) 表題文固有の文体構造がある。

特に、Experiment on —, Application of — to — 等のように、表題文全体の構造を決める英語における特有の言い回し方は、それに応じる日本語訳のパターンを用意することにより、前置詞の誤出方法など、従来困難と考えられていた問題をある程度排除するのに有効である。我々は、この種の特定な単語に固有な文体パターンをFlexible Idiom(以下 FI)と呼び、各単語の辞書に登録している。

(3) 専門用語は、必ずしも 1 語単位ではなく、句に対応することが多い。

例えば、Large Scale Integrated Circuit(大規模集積回路)のように、専門用

つの名詞句としてまとめてゆく。(にだし、 $\{\text{冠詞}\} + \{-\text{ing形}\}$ のように明らかに形容詞的に使われていることがわかる $-\text{ing形}$, $-\text{ed形}$ の場合には、語順を変換する必要がないので、TNによってそのまま読みとばされる)。

(5)表題文にあらわれる文構造は、書き換え規則で処理する必要があるほど多様ではない。

言語分析において、書き換え規則に代表される規則の体系を使用する大きな理由は、言語表現の持つ「無限」の可能性を、有限個の規則で処理したいためである。しかしながら、表題文においては、表題文全体の長さが制限されてしまうこと、埋込み文等の複雑な文体がほとんどあらわれないこと等の理由によって、すくなくとも品詞並びに還元してみる限り、規則で処理しなければならないほど多くの無限の可能性があるとは考えられまい。しかも、前述の(3)の項で述べたように、句で表現される専門用語を1つの辞書項目と考え、また、(4)で述べたように長い名詞連続をTNによって1つの名詞に縮退する操作を行なうと、可能な品詞の並びはさらに局限される。

このような考え方から、我々は表題文の品詞列を処理する際に書き換え規則等の規則による処理は行なわずに、表題文にあらわれる可能な品詞並びをそのまま列挙しておき、これと入力表題文の品詞並びとのパターンマッチングを行なうことにより、表題文全体の訳出語順を決定することにした。入力表題文から、種々の操作を経て得られた縮退しに品詞並びを、その入力表題文の骨格パターン、システムによってあらかじめ用意されている品詞並びを文体パターンと呼ぶ。したがって、入力の骨格パターンとシステムに登録されている文体パターンのパターン・マッチングによって、訳出語順が決定される。

一般に規則の数を少なくするためにには、局所的な言語の構造を処理するための規則を用意し、その組み合せで文の大域的な構造を把握することになる。同じように、解析の結果を訳出する場合にも、逆方向の規則によって、局所的な構造を一部分ずつ線形表現に変換していくことになり、全体として訳のバランスや、より大域的な観点から決めるべき語順の決定において困難がある。

以上のような理由からも、表題文のように文體に強い制限がある場合には、我々のとった文型パターンによる訳出がすぐれでいると思われる。もちろん、後述する文型パターンに付随した意味処理の指定等が、本来的には1つの統語規則について宣言されるべきものが、同じ部分構造を持ついくつもの文型パターンについて宣言されなければならぬという欠点がある。

3-2.システムの構成

表題文の翻訳は次の各ステップ^oによって行なわれる(図2に処理の流れと実例を示す(図2は末尾))。

Step 1: 辞書引きおよびRIの処理

Step 2: LOCAL AND の処理

Step 3: 名詞連続の処理

Step 4: 骨格パターンと文型パターンとのパターン・マッチング処理

Step 5: 簡単な意味処理とFIを使って訳出語順の決定

Step 6: 日本語文の合成

以下の節では表題文の各ステップ^oの処理の概要について述べる。

3-2-1. Step1 — 辞書引きおよび Rigid Idiom の処理

現在のシステムでは、英語の形態素処理を全く行なわないために、morphological variant もすべて辞書項目として含まれている。もちろん、実用化の際に規則化可能な部分を規則化することによって、辞書項目数を減少させる必要がある。

ただし、同じ 'V + ing' の形式で作られた語であっても、accounting, bonding, scattering, engineering 等はほとんどいつも名詞的に、superconducting は形容詞的に、using はもっぱら目的語としてその直前の名詞を修飾して前置詞的に、determining は他動詞的に使われる。一般に 'V + ing' 形はこれらのすべての使われ方をする可能性があるが、これらの語については、分野を固定するとほとんどある一定の使われ方をする。したがって、辞書記述の中にはこの固定的な使われ方をするということを記述しておくことにより、後の段階での複雑な処理を排除することができます。また、time varying, based on のように、RI に組込まれているものもある。英語の形態素処理を一般的なルールで行なうことが意味をもつために、それに応える日本語の形態素発生がうまくルール化でき、しかも、英語の解析、語順決定等がその形態素で変形を行うに元の単語とは無関係に実行できることが必要であるが、上述のようにもとの動詞が何であったかによって、'V + ing' 形の果す役割がわかる場合には、結局単語ごとに ing 形の付いた場合の記述が必要となり、個々の V + ing 形を辞書項目として立てるなどと、ほとんど大差がないことになり、いたずらに処理を複雑にするだけで、利点がないことになる。

3-2-2. Step2 — LOCAL AND の処理

ここでは、and だけではなく、and・or・but のような等位接続詞によって作られる並列句一般の処理を行なう。實際には、対象とした表題文中には、等位接続詞としては and だけしかあらわれなかつて（一万表題中、and を含む表題 722 件、内 27 件以上の and を含むもの 51 件）。

等位接続詞を含む句の解釈には、よく知られているように、scope ambiguity がある。例えば、「Controlability and Stability of the system」を「システムの可制御性と安定性」と誤すべきか、「可制御性とシステムの安定性」と誤すべきかは、結局 Controlability, Stability, System の 3 語間の意味的関係によって決定せざるを得ず、意味処理の問題と密接に関連する。これとシステムの安定性と可制御性」のように、並列されている句の順序を入れ替えて、日本語においても同じ曖昧さを持つ表現に置き換えることもできるが、本質的な解決にはならない。現在の我々のシステムでは、このような曖昧さを持つ等位接続詞の解釈は未解決な問題として残されている。

我々のシステムでは、文全体の大域的な構造を参照しながらも、局部的な品詞列だけから決定できる「形容詞 and 形容詞」等の連續だけを処理する（図3）。

マッチすべきパターン → 結果	
(top)	
det	adj and adj n → det adj n
adv	
ing	
ed	
v	
prep	
ing and ing	→ ing

図3. 局所的に処理可能な AND

この種の局所的な修飾語の and 結合として処理できること例は、722件中約10%の75件であった。それ以外の and の用例はすべて骨格パターン中に残されることになる。

実際には、骨格パターン中に and を残すという二つの方法は、次の二つの理由から望ましいことではない。

(1)システムが用意する文型パターンの数が組み合せ的に増大する。

例えば、文型パターンとして ' $n + prep + n$ ' があれば、これに対応して、' $n + and + n + prep + n$ '...等の各名詞句が and を含む場合の文型パターンをすべて用意する必要があり、冗長度の高いパターンの数が飛躍的に増大してしまう。

(2)次に述べる Step 3 では、長い名詞句を 1 つの名詞に縮約する操作を行なうが、次のようないくつかの例のためにこの操作を and の周辺では全く実行できないことになり、骨格パターンに残ってしまう単語が増大する。

例: Combat vehicle and aircraft stabilization system

(この例では、Combat が vehicle, aircraft の両者を形容し、また動詞stabilize の名詞化 stabilization が vehicle, aircraft の両者を目的語としてとているため、'combat vehicle → vehicle' 'aircraft stabilization system → system' というように and で結合されている 2 つの部分を単独に縮約することができる。)

以上のようなことから、and のような等位接続詞を含む句の処理は、次の Step 3, 4 へ行くまでに処理しなければならないが、この部分の処理はかなり heuristic な手法で行なわれるを得ないと思われる。

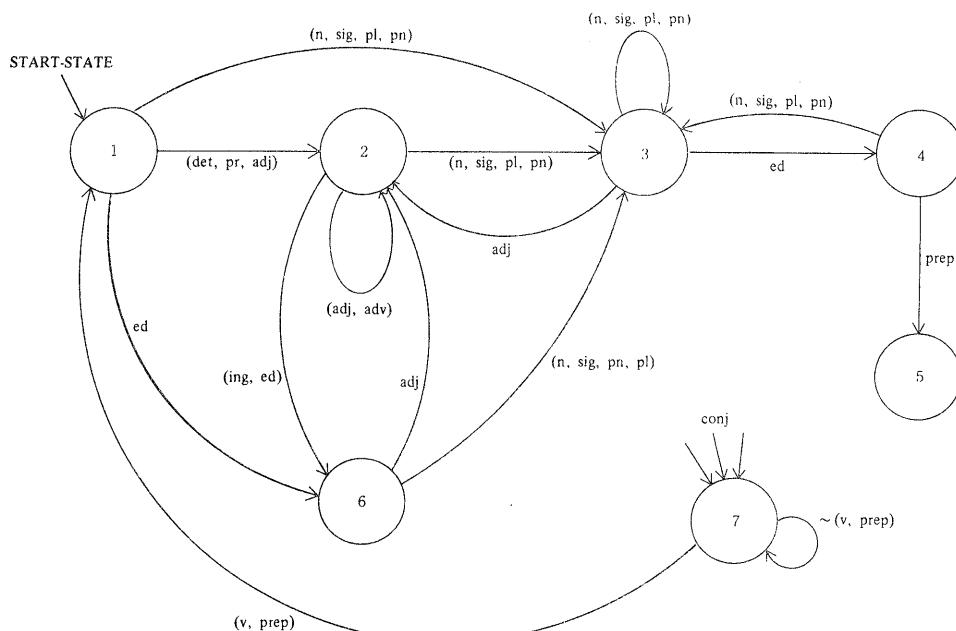


図4. 遞続名詞、形容詞を処理するTN

3-2-3. Step 3 — 名詞連続のTNによる処理

ここまで Step でつくれた品詞列は、RI と簡単な構造の並列句を処理してだけで、入力表題の品詞列とはほぼ同じものである。従って、このままで入力表題の数だけ累ねて品詞列があり得る。この Step 3 では、語順の変換を必要とする。 $'n + n'$ やそれに形容詞がもう入されたような品詞列を、1つの名詞句に縮退する。これを行なうための TN を図 4 に示す。前に注意したように、この TN はあくまで Acceptor 型の TN で、入力品詞列に対する構造を作り出可 transducer ではない。

この Step によって、次のような句が1つの單語に縮退される。

- (例) This film waveguides → waveguides
 High current arc discharge → discharge
 Pressure control system → system
 High quality phosphor screens → screens
 Transmission Line Power Flows → Flows
 An Automated General Purpose Test system → system
 Eddy current method → method
 Time varying strong magnetic fields → fields
 The refractive index profile → profile
 stochastically disturbed industrial Plants → Plants

3-2-4. Step 4 — 骨格パターンと文型パターンとのパターン・マッチング処理

以上の Step で、入力表題文に対応して、図 5 に示すような骨格パターンが抽出される。Step 4 では、これに対してシステムに用意された文型パターンとのマッチングが行なわれる。(1) Experimental Method of Measuring the Anisotropy Field in Small Samples

Skelton Pattern: Method of Measuring Field in Samples

品詞列 : n prep ing n prep n

- (2) A Laser Doppler Technique for Measuring Flow Velocities in High Current Arc Discharge

Skelton Pattern: Technique for Measuring Velocities in Discharge

品詞列 : n prep ing n prep n

- (3) Voltage Time Integral Method for Measuring Machine Inductance

Skelton Pattern: Method for Measuring Inductance

品詞列 : n prep ing n

- (4) Diffuse Field Technique for Measuring Directivity Index

Skelton Pattern: Technique for Measuring Index

品詞列 : n prep ing n

- (5) Apparatus for Thermopower Measurements on Organic Conductors

Skelton Pattern: Apparatus for Measurements on Conductors

品詞列 : n prep n prep n

- (6) The Design of a Pressure Control System for the Pepsios Spectrometers

Skelton Pattern: Design of System for Spectrometers

品詞列 : n prep n prep n

- (7) Evaluation of High Quality Phosphor Screens for Image Tubes

Skelton Pattern: Evaluation of Screens for Tubes

品詞列 : n prep n prep n

図5 入力表題文とその骨格パターン

3-2-5. Step 5
— 読出語順の決定
 Step 4 で入力文の骨格パターンが見つけられるが、一般に1つの文型パターンには複数個の読み出語順があるのでも、この Step では單語の意味記述と FI を用いて読み出語順が決定される。

FI は、ある特定の單語に固有な表題文の文体を列挙しているので、これに合致する場合には、その読み出語順及び前置詞句等の読み出

それが優先される。もし、 F_1 を持つような単語がばかりにり、もしあつたとしても、入力表題文の骨格パターンがそれに合致しない場合には、システムの用意する一般の文型パターンに対する誤が採用される。図6にその模式図を示す。

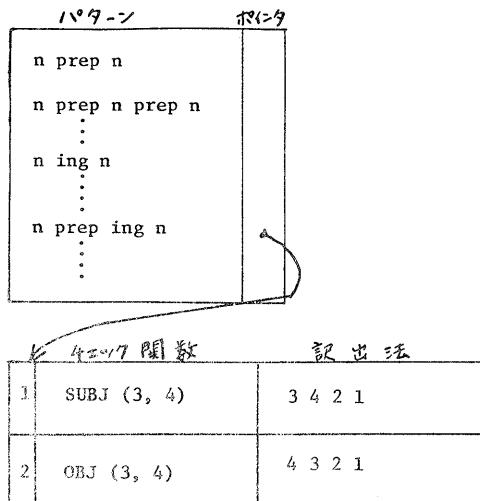


図6. 文型パターンと誤出語順の決定

すなはち、—ing形がある場合には、peratureのように、日本語訳において他の単語との結びつき方によって語順が変化するためである。これを処理するために、我々は各動詞に関する、その動詞が主語、目的語としてどのような意味カテゴリの名詞をとるかを記述し、これに対応して、名詞の意味カテゴリに分類した。その1例を図7に示す。

名詞の意味カテゴリ：動詞の意味記述は、このような簡単なものでは可もないが、現在どの程度の意味記述が実際に複数の良い訳文を得るために必要な検討を行なっている。

ここでいう誤出語順の決定は、実際には骨格パターン中の係受け関係同定によるが、この過程で問題になるのは、前置詞句の修飾先である。英語では、前置詞句はそれが修飾する語よりも後に置かれ、日本語では全くその逆になるので、英語の語順を完全に逆転させておけば、英語において前置詞句の修飾先が曖昧であれば、同じように日本語においても曖昧な表現が得られ、プロトタイプの翻訳システムとしては十分である。しかしながら、考えてみると英語と前置詞句が動詞（—ing形、—ed形を含む）を修飾しているか、名詞を修飾しているかによって、対応する日本語への修飾句の語尾変化がかかる（連用修飾句か連体修飾句か）ため、このStep 5

図に示すように、各誤出語順には、その誤出語順を意味的に支持するエック関数が用意されている。すなはち、我々のシステムにおいては、必ずしも意味的な記述が完全に与えられることを仮定しているので、Wilks の preference semantics の手法と同様に、あくまで意味処理は補助手段になっている。意味処理関数をはたらかせた結果、意味的に整合がとれていることが判った場合には、その誤出語順の優先度が上がるに至るだけで、意味的な整合が確認されなかったとしても、必ずしもその誤出語順が棄却されたことはならない。

誤出語順の決定においては—ing形の存在がもっとも大きな影響を与える。measuring device, measuring tem-perature, measuring technique, principle, method approach

道具	probe, instrument, equipment system, set, machine unit
理論	technique, principle, method approach
観点	velocity, displacement acceleration, position, energy inductance, resistance capacitance, conductance temperature mobility, permittivity intensity level parameter coefficient, time, angle
物体	laser, beam, light solenoid, coil film, tape, disc metal MOS, device, waveguide diode, transistor oil, water, liquid, solid air, earth star, stone

図7. 名詞のカテゴリ分類

における係り受け関係の決定は、より質の高い訳文を得るために重要である。現在のシステムでは、次にく3 Step 6 の日本語の形態素発生の処理を行なっていなければ、このStep 5で動詞にかかるとそれに前置詞句もすべて連体修飾語として生成しており、これが訳出結果の質を落す原因にもなっている。このStep 5の訳出語順の決定をうまく形態素発生に結びつけるStep 6を開発する必要がある。(現在のシステムでは、なるべく動詞的概念も名詞化して表現しているため、これによる不自然さはいくぶん緩和している。 —(例) A Method of measuring velocity by doppler effect

→ 「ドッペラー効果による速度測定方法」(現システム)
 → 「ドッペラー効果にて速度を測定する方法」)

4. おわりに

現在、本システムは京大大型機センターでPL/Iを使って実現されている。本システムに翻訳された実行例を図8に示す。現時点での問題点は、本文中に述べたが、これ以上により本格的なシステムを作成するためには、文型パターン、RI、FIの網羅的な収集と、単語数の増加が必須である。現在、表題数を増加させて時に、骨格パターンがどのように増えるか、JICST発行の文献データの中にある表題をかなりの安定度で訳出するためには、どの程度の文型パターン、RI、FIを用意すれば良いか等の評価を進めている。これらの評価の他に、訳文の質を人間の翻訳と比較して、定量的に検討していくことも重要なよう。

今後に残されたシステム構成上の問題点は①ANDの処理、②名詞、動詞の意味記述、③語順が確定した時点での日本語の形態素発生 等が考えられる。

```
► GENERATION OF STANDARD EM FIELDS USING TEM TRANSMISSION CELLS.

< 001 N      > < 002 PREP > < 003 N      > < 004 SIG   > < 005 PL    > < 006 PREP >
< 007 SIG   > < 008 N      > < 009 PL    >

< 001 N      > < 002 PREP > < 003       > < 004       > < 005 N      > < 006 PREP >
< 007       > < 008       > < 009 N     >

SKELTON PATTERN: GENERATION OF FIELDS USING CELLS

BUNBANGO=      22**DAI      1 BANME**

GENERATION OF STANDARD EM FIELDS USING TEM TRANSMISSION CELLS.
TEM伝送セルを用いた標準EM場の生成
BUNBANGO=      23

► TRANSIENT ELECTROMAGNETIC PROPERTIES OF TWO, INFINITE, PARALLEL WIRES.

< 001 ADJ   > < 002 ADJ   > < 003 PL    > < 004 PREP > < 005 N      > < 006 ADJ   >
< 007 ADJ   > < 008 PL    >

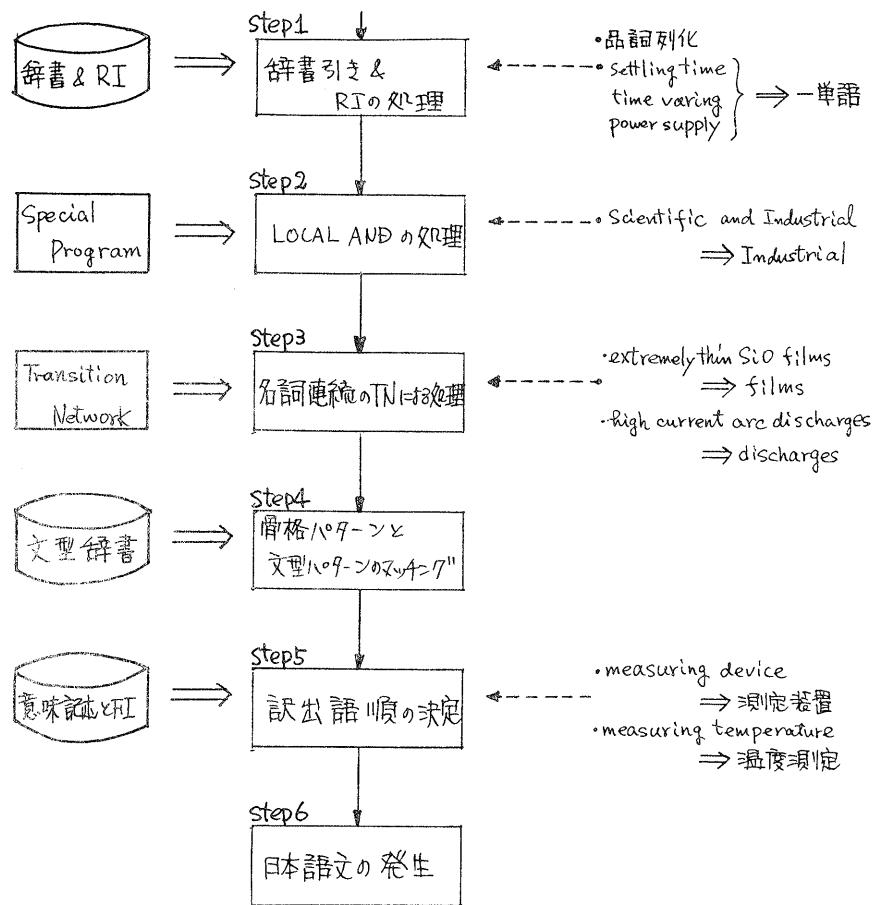
< 001       > < 002       > < 003 N      > < 004 PREP > < 005       > < 006       >
< 007       > < 008 N      >

SKELTON PATTERN: PROPERTIES OF WIRES

BUNBANGO=      23**DAI      1 BANME**

TRANSIENT ELECTROMAGNETIC PROPERTIES OF TWO, INFINITE, PARALLEL WIRES.
二無限平行ワイヤの過度電磁性質
```

図8 翻訳システムの実行例



翻訳ステップの例

入力 : Industrial and Scientific Techniques for Measuring field Effect Mobility

Step 1 : Industrial and Scientific Techniques for Measuring Effect Mobility

Step 2 : Scientific Techniques for Measuring Effect Mobility

Step 3 : Techniques for Measuring Mobility

Step 4 : <文型ハセターン> m prep ing m

読み出結果 : 電界効果移動度測定のための工業的及び科学的手法

図2. 翻訳システムの流れと実行例