

# 自動翻訳から翻訳支援へ，そして…

(株) 富士通研究所 知能システム研究部

潮田 明, 富士 秀, 大倉 清司, 山下 達雄

ushioda@jp.fujitsu.com / fuji.masaru@jp.fujitsu.com / okura.seiji@jp.fujitsu.com / yto@jp.fujitsu.com

「フルコースを堪能した男は、飲み干したコーヒーカップを満足げな表情で静かに置いた。男はアルコールのせい或少し上気している様子だった。しかし給仕がそっと置いていった伝票を手を取った瞬間、赤みを帯びた男の顔は青ざめていった。」

下手な小説の一節のようであるが、何が起きたかは賢明なる読者には明らかであろう。一般に言葉を他の言語に翻訳するにはまずその内容を理解しなければならないといわれている。確かにその通りなのであるが、翻訳を計算機にやらせようとした場合、言葉の真の意味を理解させてからというのは残念ながら現実的なアプローチではない。上の文章は現状レベルの機械翻訳にはやや手ごわい相手であるが、以下の程度の訳を出すことは現在の機械翻訳でも可能である。

(機械翻訳による翻訳文)

*The man who was satisfied of a full course quietly put the coffee cup that had been drunk up in the pleased expression. The man seemed to have flushed a little probably because of alcohol. However, the face of the man who had a tinge of red was pale momentarily at the time of having taken the slip that the waiter had quietly put.*

この翻訳でどのくらい意味が通じるかを調べるために英語のネイティブスピーカー 5 人に翻訳文を読んでもらった後に 2 つの質問をしたところ全員が以下のように回答した。

- ① Did the man like the taste of the dishes ? ==> Yes
- ② Did the man like the price of the dishes ? ==> No

これは当然我々の直感に合った回答だが、このような答えを出すためには、言葉の意味（内容）を正しく理解することが必要である。しかし果たして計算機に同じ質問をして同じように答えさせることができるだろうか？ ①には答えられても②に答えられる計算システムはまずないといってよいだろう。つまり、ある意味で機械翻訳システムは日本語の内容を真に理解せず正しく内容を英語で伝えたことになる。

こういった現象が人間のしゃべる言葉のどのくらいの範囲まで起き得るかは定かでないが、少なくともこれまでの機械翻訳の研究は言葉の意味の真の理解を迂回して翻訳を行う道を求め続けてきたし、今後も急速に変わることはないであろう。むしろ逆に意味の理解からさらに遠ざかって、表面的な現象をたくさん寄せ集めて字面上の統計的な特徴を基に翻訳を行おうという動きが活発になってきている。そのようなアプローチで最終目標である完璧な翻訳品質に到達できるかは専門家の間でも意見

の分かれるところであるが、一方視点を変えると新しい展開がみえてくる。

つまり、目標を完璧な自動翻訳に置くのではなく、翻訳するのはあくまで人間であり、人間の翻訳スピードを計算機の助力により飛躍的に速めてやるという発想である。真の意味理解は人間にまかせ、機械は字面上の翻訳知識を大量に蓄積して、必要なときに必要に応じて人間に提示するという、人間と機械の協調作業が、完璧な自動翻訳に到る過程でのマイルストーンとして浮上してきた。

本稿では、これまでの機械翻訳のたどってきた道りを簡単に振り返りながら、翻訳支援という別の流れとどのように融合して新しい方向に向かっていくのかを概観する。

## 【機械翻訳の進化と現状】

機械翻訳は、これまで主に「自動的に翻訳するための機械」として開発されてきた。インターネットの出現によって日本人が外国語文書に触れる機会が増大したため、大量の英語文書を高速に和訳して概略をつかみたいという需要が増えてきたが、そのためのツールとしては威力を発揮するようになってきた。外国語が苦手な一般ユーザの情報収集ツールとしては市民権を得ているといえよう。しかし、やはりそこには訳質の問題が壁として立ちはだかり、これ以上の利用には結びついていないのが現状である。

なお、ここで1つ気をつけたいのが用語である。文書を自動的に翻訳するためのソフトウェアの技術は、英語で“machine translation”(MT)と呼ばれるが、これを「機械翻訳」と訳して使うことが多い。しかし、「自動翻訳」も同様の意味で用いられる。また、製品として商用となっている機械翻訳のソフトウェアは「翻訳ソフト」等と呼ばれることが多く、一般の人はこの名称の方が馴染みが深いであろう。

### ◆機械翻訳がたどってきた道のり◆

コンピュータ言語と違い、人間が使う言語（自然言語）にはあいまい性が存在し、自動処理は非常に難しい。しかしそんな中で、Noam Chomsky が1957年に提案した変形生成文法は、人間の言語をコンピュータで扱うための具体的な道筋を示す記念碑的な理論となった。

言語の処理の中でも、言語間の翻訳はかなり初期の頃から注目を集めていた。機械翻訳の研究は1950年代にヨーロッパで始まり、次いで1960年代初頭にはアメリカでの研究も活発化する。日本でも、1959年に通産省電気試験所（現在の産業総合研究所）において初歩的な英文和訳機 YAMATO が発表され、やや遅れて、九州大学の超高周波研究室でも独文和訳機の試作が行われた。

ところが1976年、機械翻訳システム評価のために設けられたアメリカの ALPAC 委員会で、機械翻訳システムは人間の言語の複雑性を処理するための十分な能力を持ち得ない、と報告された。この報告によって、アメリカ政府の研究開発費の大半が打ち切れ、またそれ以外の国の機械翻訳プロジェクトにも大きな負の影響を与えた。

その一方、日本では1970年代末から大手メーカーによる機械翻訳ソフトの研究開発が始まった。1982年には科学技術庁機械翻訳プロジェクト（Mu プロジェクト）が開始され本格的な機械翻訳研究の幕開けとなった。

1980年代初頭には機械翻訳システムの商品化も進められた。初期の機械翻訳といえば、ハードとともに販売され、ソフト自体の価格も数百万円と高価なものであった。この頃は主にルールベースとって解析規則（文法）を書いて機械翻訳を制御していた。特に構造が違う言語間の翻訳（たとえば英日翻訳）では、文の構造解析や、さらに進んで文の意味解析まで対処する文法が開発された。

やがて1990年代、インターネットの立ち上がりとともに多くのベンダが機械翻訳市場に参入し、安い値段（数万円）で翻訳ソフトを販売、競争が激化した。この頃になると、ハードの進歩、国家プロジェクトによるコーパス（電子化された大量テキスト）作成などを背景として、コーパスを使った機械翻訳用辞書の学習アルゴリズムや、翻訳文法の習得などの研究が盛んに行われた。また「例による翻訳」(Example-Based MT<sup>1</sup>)も同様に実現味をおびてきた。さらに近年では、ルールベースと例による翻訳を組み合わせた「パターン翻訳」という枠組みを採用するシステムも多く現れた。

### ◆機械翻訳の仕組みと精度向上の方向性◆

#### 基礎編

翻訳を行うシステムの大半で用いられている最も基本的な処理方式がルールベース翻訳である。これは、前述の Chomsky の提唱したアイデアが基になっており、文法規則を使って翻訳処理を行う。文法規則は言語ごとに異なるので、言語ごとに用意する必要がある。たとえば非常に単純な日本語文は以下の規則で生成できる。

S → NP VP #S (文) は名詞句 (NP) と動詞句 (VP) からなる。  
VP → NP V #動詞句 (VP) は名詞句 (NP) と動詞 (V) からなる。  
NP → N JO #名詞句 (NP) は名詞 (N) と助詞 (JO) からなる。  
N → 私 #名詞 (N) には「私」がある。  
N → 彼女 #名詞 (N) には「彼女」がある。  
V → 見る #動詞 (V) には「見る」がある。  
JO → が #助詞 (JO) には「が」がある。  
JO → を #助詞 (JO) には「を」がある。

これらの規則から、「私が彼女を見る」「彼女が私を見る」「彼女を私が見る」「私を彼女が見る」が生成できる。翻訳を行うにはさらに、「が」は主語を表す、という情報や、「見る」は他動詞で主語と目的語をとる、という情報などが必要になってくる。「見る ⇔ see」など、対訳辞書も必要になる。

## 深層情報を充実させる

言語の翻訳では、単純な文法規則だけでは不十分で、原言語の深い意味を解釈して意味的な構造にマッピングし、そこから訳文を生成するアプローチが必要な場合もある。原文と、出力すべき訳文の構造が極端に違う場合にこのアプローチは不可欠である。実際には言語非依存の構造にマッピングするのは難しいため、たとえば日本語寄りの意味構造、英語寄りの意味構造へのマッピングを行うことも多い。たとえば、「これで復旧が容易になる」という文を日英翻訳する場合、日本語の意味構造は

<BECOME> 手段 → <THIS>  
<BECOME> 対象 → <RECOVER>  
<BECOME> 目標 → <EASY>

という構造で表すことができるが、ここからそのまま英語を出力すると、Recovery becomes easy by this. となってしまう。しかし、上の日本語寄りの意味構造から、次のように英語寄りの意味構造にマッピングを行うと ...

<MAKE> 動作主 → <THIS>  
<MAKE> 対象 → <RECOVER>  
<MAKE> 目標 → <EASY>

この意味構造から、This makes recovery easy. という英語が生成できる。

ただし、この深層処理も万能というわけではなく、本当は表層的な情報を入れておけば単純に処理できるような場合でも、必要以上に深い処理が行われて誤った結果になることも起きる。極端な例では、「おはようございます。」という文の深い意味を解釈して変換するよりも、「おはようございます⇔ Good morning」というように文をまるまる入れて変換したほうがシンプルで確実である。

## 表層情報を充実させる

深層処理の対極にあるのが表層処理である。

いくら複雑な規則を作っても解析が失敗する場合は必ずあるわけで、それならいっそ頻度の高い表現や用語は大量のコーパスから拾ってきてそのまま入れたほうがうまくいく場合も出てくる。このように、表層に近い情報を大量に収集するようなアプローチとしてコーパスからの翻訳知識獲得の流れがある。

対訳コーパス（電子化された大量の対訳文）からは、まずは、文そのものやフレーズ等の文の断片を取り出して利用することができる。深い処理を通ることなく、フ

レーズ等を組み合わせることによって翻訳を進めるアプローチも試みられている。

さらには、対訳コーパスからは未知語学習や文法規則学習を行うこともできる。未知語とは辞書に入っていない単語のことである。日本語の複合語（2単語以上で1つの意味を表す語句）は1年で数十万語の新語が出てくるといわれている。対訳コーパスから、文の対応づけ→単語対応づけ→未登録単語の対応づけという方法で、未知語の対訳候補が抽出できる。また、対訳コーパスは「この文はこう訳す」という例の集合とみなすことができ、文の言語間のマッピングを統計的に学習させることができる。

## 表層処理の改良版アプローチ

表層処理で大量の翻訳知識を取り込むことができて、そのままでは柔軟性が低い場合も多い。たとえば、対訳コーパスの文やフレーズ等は、構成単語が少し違っていただけで入力にヒットしない場合もある。これを克服するための技術がパタン翻訳であり、表層情報に単純なルールを取り込んだものと考えられることができる。別の見方をすると、パタン翻訳とはルールベース翻訳を拡張したもので、「パタン」によって翻訳規則を記述する。通常の翻訳ルールに比べ、単純であるのが特徴といえる。

たとえば以下のようなパタンをつくることができる。<..> は変数を表し、N1, N2 は変数のラベルを表し、原文と訳文で対応している必要がある。ここでは変数には名詞がくる、という制約がある。

<N1:They> gave <N2:him> the fact.  
<N1: 彼ら> は <N2: 彼> に事実を知らせた。

このパタンを使うと、たとえば Tom gave the man Mary saw yesterday the fact. という文を、「トムは昨日メアリが会った男に事実を知らせた」と翻訳できる。

パタン翻訳はこのようなアプローチの一例だが、表層情報と深層情報の両方を充実させるための試みは今後さまざまなかたちで続けられるだろう。

## ◆誰にとって役に立つのか？◆

現状の機械翻訳システムはまだまだ「完璧な」翻訳を行うことはできない。そのため、たとえば英語力のある人は、英日機械翻訳の出力よりは英語を直接読んだほうが良いと考えているだろう。しかし、同時に英語に自信がなく、完璧でなくても何かに頼りたいと思う多くの人々がいることも確かであろう。「どの程度の英語力の

TOEIC 得点	読解得点向上	印象向上
～ 490	○	○
495 ～ 690	○	×
690 ～	×	×

表-1 英語力と機械翻訳の読解・印象

人にとって機械翻訳が役に立つか」ということを実験で求めた研究<sup>4)</sup>があるので、この研究を紹介して読解における機械翻訳の現状を説明してみたい。

研究では、機械翻訳ユーザの英語能力の尺度として、多くの受験者を持つ TOEIC を用いた。ここでの考え方を一言でいえば、TOEIC で 500 点を取った人が機械翻訳の出力を利用することによって 600 点の成績を修めることができれば、この人にとってはこの機械翻訳システムは有効であるといえる、ということである。

手順としては、英語で書かれた TOEIC の読解用文書を機械翻訳で和訳し、その訳文を読んで質問に答える。このときの正解率が、英語原文を読んだときの正解率と比べて向上したかどうかを測定する。この実験を、幅広い TOEIC 得点層の被験者に対して行い、統計的な有意性を検証するわけである。この実験によって、どの程度の英語力があれば読解に有意差があるかが分かる。また、この実験を行った後に、「英語原文と機械翻訳和文のどちらのほうが分かりやすいか」という、被験者の直感に関するヒアリングも行って数値で表現し、統計処理を行った。

なお、機械翻訳文の提示の仕方としては、被験者に機械翻訳文のみを単独で提示するという方法と、原文と機械翻訳文を並べて提示するという方法とがあるが、より現実の利用形態に近い後者の結果について述べる。

表-1 は、被験者の英語能力 (TOEIC 得点) と、読解得点の向上および分かりやすさの印象の向上の関係を表したものである。英語が苦手と思われる TOEIC 低得点取得者は、読解得点と分かりやすさの印象の両方が改善しており、英語を得意とする TOEIC 高得点取得者は機械翻訳の利用は読解も印象もよくなっていない。これはある程度予想された通りだが、その中間層として、読解は向上しているが、印象は向上していない被験者が存在する。

この層の人たちにとっては、機械翻訳の訳文に問題があって印象はよくないが、実際問題としては原文そのままよりも理解の助けになっているということだろう。

図-1 は、さらに、TOEIC の実施団体が公開している TOEIC スコアと人数分布の統計グラフに今回の実験結果を重ね合わせたものである。読解という用途に限ると、

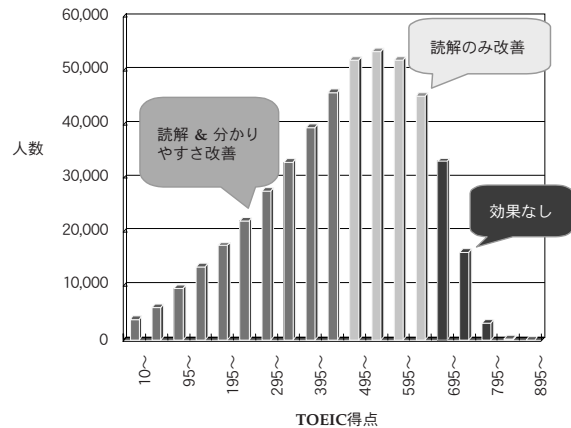


図-1 英語力分布と機械翻訳効果

不完全な機械翻訳でも多くの人がメリットを感じられるということが分かるのである。

#### ◆現状の機械翻訳の限界◆

さてここまで開発の歴史から始まり、機械翻訳について外観してきた。自動で翻訳を行う機械を作るにあたって、表層のアプローチと深層のアプローチの両側面から技術を極める努力が続けられているが、当面は英語の苦手な人にとっての読解支援という範囲の利用にとどまっているのが現状である。

#### 【翻訳支援の起源と現状】

これまでお話ししてきた機械翻訳による自動翻訳技術の進化とは別の世界で、発想の原点は違うが、やはり翻訳の自動化につながるような動きがある。翻訳業界における「翻訳支援」がそれである。

ここでは、主体はあくまでも翻訳を行う翻訳専門家であり、その人たちが支援システムを使うことによっていかに効率を上げられるかが重要なのである。つまり、システム自体の精度や速度を測るのではなく、翻訳支援を導入したときの人手作業の効率向上が最大の関心事となる世界である。

#### ◆欧米で始まったローカライゼーション◆

産業翻訳において翻訳支援がいち早く注目され実用化されたのは、ローカライゼーション業界である。ローカライゼーションとは耳慣れない言葉だが、その典型的な例としては、たとえば、ある欧米のソフトウェアの日本語

版を開発する際に発生する一切の翻訳業務がある。ソフトウェアのインタフェースで使用される対話のテキストの翻訳や、ソフトウェアの取り扱いマニュアルの日本語化等である。

一般的に、ソフトウェア製品の改版頻度は非常に短く、そのたびに大量のローカライズ作業が発生する。しかし、ある版の翻訳は、その前の版の翻訳結果をそのまま参考にできる場合が多い。このことに着目した支援ツールメーカーでは、作成された対訳例文をデータベースに蓄積しておき、これを参照しながら効率よく作業を行うような技術が発達した。業界では、このようにして蓄積した対訳例文を「翻訳メモリ」と呼ぶことが多い。

このような翻訳業界のニーズから発展した翻訳支援技術は、実用レベルで使われるようになってきた。現状では、翻訳業界において唯一効果が認められた支援システムが、このローカライズ業界における対訳例文再利用のシステムだといってもよいだろう。

## ◆翻訳業界の市場◆

翻訳業界の市場規模は、統一的な集計が難しいことから、正確な数値は出されていないが、全世界的にかなり大規模であることは確かである。

翻訳というと、一般的には文学作品等のいわゆる文芸翻訳を想像することが多いが、実際に分量が大量に発生して、大きな市場を形成しているのは産業翻訳である。この中で日⇄英間の翻訳は国際的にも大きな比重を占めている。

## ◆「翻訳メモリ」とは？◆

翻訳メモリとは、過去の訳例（原文と訳文が1組になった翻訳事例）をデータベースに登録し、一致や類似検索により再利用する技術である。コンピュータがすべて翻訳する機械翻訳とは異なり、ユーザが訳例を検索するといった、人間が翻訳するときの支援ツールとして位置づけられる。

マニュアルの改版やアニュアルレポート、変更箇所が少ない場合や、同じ文言を使いまわすことの多い翻訳に用いることが多い。また、ソフトウェアのローカライゼーション（メニューやヘルプなどの翻訳）にも多く用いられている。データベースに登録されている類似した訳を使いまわすことにより、翻訳効率の向上や用語統一がはかれる。短時間で質の高い翻訳が行えるという効果があり、コストの削減につながる。

## 例：翻訳メモリの使い方

簡単な例で翻訳メモリの使い方を説明する。あらかじめ訳例データベースには、大量の英日対訳文が登録されているものとする。ここで、翻訳者が「I buy an apple.」という英文を翻訳したいとする。翻訳者はこの英文をキーに訳例データベース（翻訳メモリ）を検索する。すると、データベースから「I buy an apple.」に似た英文「I eat an apple.」を持つ訳例データが得られる。翻訳者は、その訳例データ中の日本語訳「私はリンゴを食べる。」を編集して翻訳文を完成させる。翻訳者は、英文の差分「buy-eat」により、「食べる」をbuyの訳語である「買う」に置き換える。こうして、「私はリンゴを買う。」という訳文ができ上がる。

原文：I buy an apple.

↓

訳例：I eat an apple. / 私はリンゴを食べる。

↓

訳文：私はリンゴを買う。

## 翻訳メモリの業界標準フォーマット

XMLによる翻訳メモリデータを記述する標準的なフォーマットとして、TMXがある。LISA<sup>2)</sup>というローカライゼーションに関する非営利団体により提供されている。LISAのサイトでDTDなどの必要なものが入手できる。TMXにより、各ベンダが開発した翻訳メモリソフトの翻訳メモリデータを相互に利用できるようにすることができる。

近年、このTMXのサポートをサポートするソフト、つまり、TMXによるエクスポートやTMXで書かれた翻訳メモリデータをインポート可能な翻訳メモリソフトが増えてきている。TMX自体は交換用フォーマットであるため、対訳文ペアの格納などの最低限の簡単な情報のみを対象としている。対訳文間での単語やフレーズの対応などのより複雑な情報は扱うことはできない。

## 翻訳メモリの限界

類似度の高い訳例は翻訳対象文との差異が少ないため、容易に再利用できる。

たとえば、マニュアルの改版における以下のような例である。この場合は、ほとんど一致する訳例が得られ、バージョン番号などの差異を変更するだけでよい。

原文：Thank you for purchasing ej/je translation software  
'atlas V9.0'.

↓

訳例：Thank you for purchasing ej translation software 'atlas V8.0'. / このたびは、英日翻訳ソフト『ATLAS V8.0』をお買い上げいただきまして、誠にありがとうございます。

↓

訳文：このたびは、英日・日英翻訳ソフト『ATLAS V9.0』をお買い上げいただきまして、誠にありがとうございます。

しかし、上記のようなテンプレート的な文ではない場合、このように満足のいく訳例が検索されることはまれである。

たとえば、データベースに登録されている訳例が少ない場合や他の分野・用途の訳例のみが登録されている場合などは、「若干の変更でそのまま再利用できる訳例」はほとんど検索されない。

この場合、検索結果には類似度の低い訳例しか含まれないことになる。実際には、類似度が低いといっても、文全体での類似度が低いのであって、文の一部分に着目するとほとんど一致するようなものもある。

これらをいかに効率的に再利用できるかが翻訳支援において重要となる。

#### ◆そして、従来型翻訳支援の限界…◆

ここで述べた従来型の翻訳支援は、過去の例文をそのままに近いかたちで利用しようという、いわば表層的なアプローチをとってきた。これは、ローカライゼーションという、ある程度対象分野を絞って作業を行う環境であるがゆえに出てきた発想である。とはいえ、このアプローチの問題は「類似度の高い場合のみ有効」ということであり、現実に存在する文書を見渡すと従来型の翻訳支援の適用範囲はあまりにも狭すぎるという問題は解決されていない。

### 【翻訳支援に機械翻訳を取り入れる】

従来型の翻訳支援は翻訳メモリという表層のアプローチに依存したものであり、この限界を克服するべく機械翻訳の深層的要素を取り入れようという動きが一部で活発化している。翻訳支援に機械翻訳技術を導入することによって得られると期待されるメリットは大きく分けて2つある。1つは、機械翻訳の品質は全体としてはまだ十分ではないとはいえ、訳の中には実用上十分使える部分もあるという点である。もう1つは、機械翻訳技術を導入することにより、より付加価値の高い翻訳メモリの構築が可能になり、過去に蓄積した翻訳知識からより柔

軟に必要な知識を取り出すことができるであろうという点である。たとえば、機械翻訳の文構造解析機能を使うことにより、文単位でのみ対応のとれた翻訳メモリの中のさらに細かい部分表現同士のマッチングがとれ、それによってきめの細かい訳例の再利用が可能になると期待される。これが実現できれば、蓄積された例文と入力 of 翻訳対象文の類似度がそれほど高くなくても必要な知識が得られるようになり、ひいては、対象文書の幅も格段に広がると期待されるわけである。

なお、翻訳支援はあくまでも人手作業用の環境であるから、単に翻訳メモリと機械翻訳を合体させるだけではまったく不十分である。翻訳作業のフローを最大限に意識しながら翻訳メモリと機械翻訳が有機的に融合され、自然な作業の流れを作り出すように設計されていなければならないのである。

#### ◆システム開発の動向◆

産業翻訳用の翻訳支援システムに機械翻訳を取り入れる試みはかなり前から行われていたが、最近になってようやく現実的な場面での翻訳効率化を意識したシステムが出てきている。

1つは、従来から機械翻訳の開発を行ってきたメーカーが翻訳支援用の機能を追加してきたものである。このようなメーカーは、言語を扱うための幅広い技術力を持っている場合も多いが、問題は、翻訳業界とのつながりが少なく、現実に即した開発ができてこなかった面がある。2つ目は、翻訳業界主導の翻訳メモリシステムの強化で、でき合いの機械翻訳を取り込むような動きもみられる。しかし、翻訳業界自体には、言語処理のためのソフトウェア開発技術を持ったところが少なく、開発上のネックとなっている。

#### ◆機械翻訳と翻訳メモリの融合◆

ここでは、翻訳メモリと機械翻訳を融合した翻訳支援システムの一例として、富士通研究所でプロトタイプを開発した翻訳支援システム Cliché ついて述べる<sup>3)</sup>。このシステムは、従来の機械翻訳と翻訳メモリを有機的に統合し、翻訳メモリの機能を高度化することによって、産業翻訳での本格的な利用を現実しようとしている。

Cliché は、翻訳メモリ+機械翻訳のこれまでにない統合型翻訳支援システムである(図-2)。クライアントシステムは翻訳対象文から原文を1文ずつ切り出して翻訳エディタに提示する。翻訳者はエディタ上で訳文を作成するが、その際 GUI を通して、機械翻訳サーバおよび

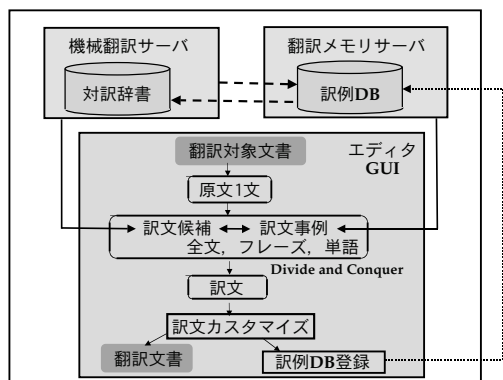


図-2 統合型翻訳支援システム

翻訳メモリサーバにアクセスし、原文やその一部に対する自動翻訳結果や訳例検索結果から有用な表現を選択し作成中の訳文に簡単な操作で挿入することができる。訳文作成の手順は翻訳者の自由であるが、多数の翻訳者による試験運用から分かった効率的な手順は以下の通りである。

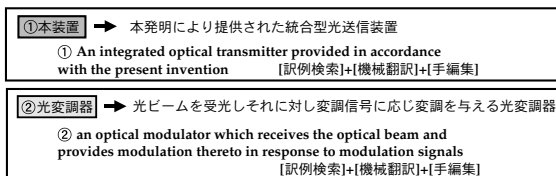
まず基本的な方針は“divide-and-conquer”である(図-3)。これは Example-Based MT における機械処理の理想的な工程に似ている。まず文の大きな構造をつかみ訳文の骨格を決める。そして骨格の主要要素をそれぞれ翻訳し、最後に各要素の翻訳を骨格にはめ込んで訳文を完成させる。要素自体が長い、あるいは複雑な場合はその要素をさらに骨格と要素に分割して翻訳を進める。このように各要素を分割していくと、ある時点で同一の表現が訳例中に見つかるか、あるいは機械翻訳(部分翻訳)で正しく訳されるようになる。作成された訳文は最後に訳文スタイルのユーザカスタマイゼーションが施されて翻訳文書に挿入される。このような翻訳の進め方は、特許などのように一文が長い文書に特に有効である。

翻訳メモリ中の各訳例は作成・登録時に機械翻訳にかけられ、解析情報付の訳例として構造化されて保存される。訳例の原文と訳文の単語などの対応情報も同時に格納される(図-4)。たとえば「This is a pen.」(原文)「これはペンである。」(訳文)という簡単な訳例で考えてみると、従来技術ではこれを単に文字列情報としてしか持たないので、たとえば「pen」と「ペン」が対応していることが即座に分からなかった。本システムでは、原文、訳文ともに機械翻訳で解析を行い、単語の対応づけをする。その結果、「This」と「これ」、「pen」と「ペン」の対応がつけられる。検索時にこの解析済訳例が表示されるが、対応情報も表示される。これにより、結果表示画面において従来の翻訳メモリより効果的なインタ

原文：  
 [本発明により提供された統合型光送信装置は、] [光ビームを受光しそれに対し変調信号に応じ変調を与える光変調器] によって特徴付けられる。

骨格文：①本装置 は、②光変調器 によって特徴付けられる。

→ ① This device is characterized by ② an optical modulator. [機械翻訳]



最終訳文

[An integrated optical transmitter provided in accordance with the present invention] is characterized by [an optical modulator which receives the optical beam and provides modulation thereto in response to modulation signals.]

(a) 骨格文 - 主要素の分割と翻訳結果の結合

② 光ビームを受光しそれに対し変調信号に応じ変調を与える光変調器

- ②-1 光ビームを受光し receives the optical beam [訳例検索]
- ②-2 変調信号 modulation signals [機械翻訳]
- ②-3 変調を与える provides modulation [訳例検索]
- ②-4 光変調器 an optical modulator [機械翻訳]

⇒ ② an optical modulator which receives the optical beam and provides modulation thereto in response to modulation signals [手編集]

(b) [訳例検索] + [機械翻訳] + [手編集] の例

図-3 Cliché における翻訳作業の流れ

フェースが可能になった。訳例の解析は機械翻訳により行われるが、機械翻訳の精度が上がれば対応づけの精度も向上する。一方、自動語句対応づけシステムにより翻訳メモリから機械翻訳の辞書も構築することが可能になった。翻訳メモリの質が上がり、量が増えればそれだけ抽出できる語句の質・量も多くなる。つまり、翻訳メモリが賢くなれば機械翻訳も賢くなり、機械翻訳が賢くなればさらに翻訳メモリも賢くなる。

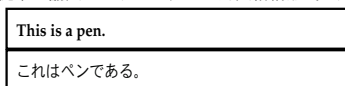
ソフトウェアではいかによい機能があってもインタフェースによっては使い勝手が非常に悪いものが出てしまうため、インタフェースの作りこみは重要な技術要素となる。本システムは企画段階から翻訳支援という観点で画面が設計され、ユーザビリティテストなどを通じて使いやすさの向上が図られた。具体的にはウィンドウ内のエディタや結果画面の配置、よく使われるボタンの配置などである。

#### ◆翻訳メモリの高機能化◆

翻訳メモリの検索の仕組みについて、ここでは例として Cliché<sup>3), 5)</sup> での具体的な方法を用いて説明する。

大量の訳例データに対しては、シーケンシャルな類似

従来の翻訳メモリシステムの対訳格納形式：文字列



開発した翻訳ワークベンチの対訳格納形式：構造的

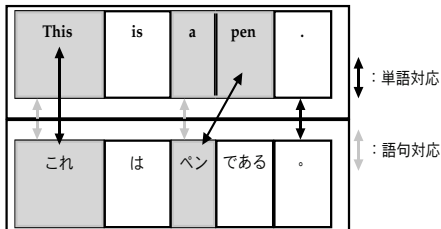


図-4 翻訳メモリの訳例格納形式

検索は速度から見て実用的ではない。そこで Cliché では、検索はインデックス検索による「絞り込み」と「マッチング」の2段階の処理を行っている。

事前に大量の訳例データに高速類似検索用のインデックスを付与し、類似度の高い順に指定された数だけ絞り込む。Cliché ではインデックスには suffix array を採用している。

この絞り込まれた検索結果（訳例）に対し、マッチング、つまり、検索キー文と訳例原文の一致個所の認識を行う。一般に、ダイナミックプログラミング（DP）による手法を用いる。

### 検索結果表示

検索の結果、類似した訳例が得られるわけだが、それをユーザに提示する際には、検索キー文と訳例原文で一致した個所、または、相違個所をハイライトする方法がある。これらの対応情報は DP の結果から得られ、扱いやすいため、TRADOS など多くのシステムで採用されている。

さらに、訳例原文と訳例訳文の対応をハイライト表示すると、ユーザにとって訳例の使える個所が判別しやすくなる。数字や日付などの要素については、簡単な文字列処理で認識できるため、原文・訳文間で対応づけしやすい。

Cliché では、機械翻訳システムの辞書を用いた形態素解析処理により訳例の両言語の文中の単語対応を得て、それらの対応をハイライト表示している。また、これに検索キー文と訳例原文の対応情報を統合して、「3つ組」ハイライト表示を行っている。

図-5 に3つ組表示の例を示す。検索結果の各訳例は、上から、検索キー文、訳例原文、訳例訳文の3文を単位としてボックス表示している。「I」-「I」-「私」、「pen」

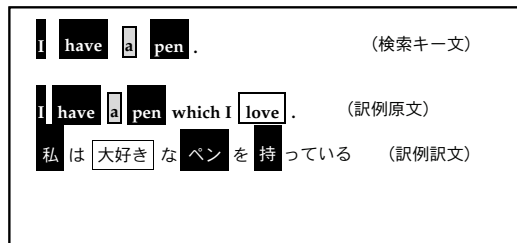


図-5 3つ組表示例

- 「pen」-「ペン」などが3つ組に当たる。たとえば、図中の検索キー文の「pen」にマウスを合わせると、訳例原文の「pen」と訳例訳文の「ペン」もハイライトされ、ユーザは対応個所を容易に認識できるといった動的ハイライトを行うことができる。

3つ組表示により、検索キー文中の単語について、この訳例中での訳が一目で分かり、ユーザによる効率的な訳語選択が可能となる。また、3つ組単語対応をヒントに使うことで、ユーザによるフレーズの把握が容易になり、訳例の部分利用が促進される。

### ◆どのくらいの効率化につながるのか？◆

翻訳支援システムの効率評価の基本的なアイデアは、システムによる支援ありの場合と支援なしの場合の効率を比較することによって効率化を見積もることである。とはいえ、現実の実験では効率化を定量的に測定するのは難しい。これは、翻訳者、支援システム、データ類等の一切を含めた要素が実験環境に含まれており、それらすべてを考慮した評価手法を設計しなければならないからである。以下では、Cliché を対象にした効率化測定実験を例に評価手法について説明する。

留意すべき点の1つとして翻訳速度と訳質の関係がある。翻訳作業効率の測定において、単に作業速度だけを測定しても有意な差が出ない場合が多い。これは、実際の人間の作業では、訳質を犠牲にしてまで高速に作業したり、また逆に訳質ばかりに拘って必要以上に作業が遅くなるような場合が頻繁に発生するからである。このことから、作業速度と訳質の両方を測定することによって初めて有効な測定ができる。

また、支援あり翻訳と支援なし翻訳のような2条件間の測定を比較する際には、条件の適用順が結果に影響を与える。これは、被験者が対象文書を扱うにつれて対象文書に慣れて処理速度が上がっていくためである。この問題を解決するために、条件1の文書と条件2の文書を交互に評価対象とした。



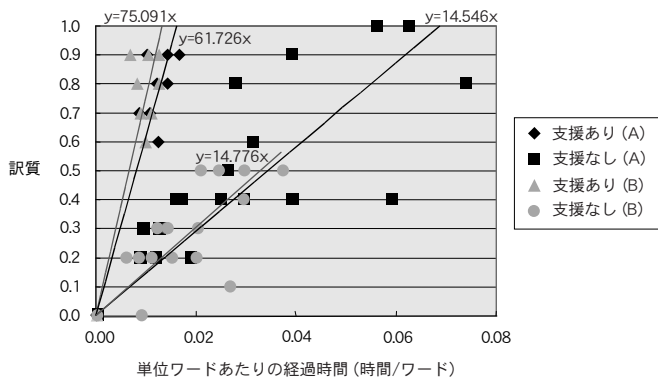


図-6 翻訳時間と訳質

図-6では、Clichéにおける翻訳時間と訳質の推移を、2人の翻訳者A、Bが支援あり条件と支援なし条件で行ったときのデータを重ねて表示している。この線の傾きが効率化の度合いに相当するわけだが、この例では4倍程度の効率化が測定されたところを示している。

## 【究極の翻訳自動化とは？】

機械翻訳はその歴史の中で、意味処理などの深層的な情報の利用と、コーパス処理等の表層的な情報の両方を充実させる方向で研究が進められてきた。その時々で、いずれかが一時的に脚光を浴びることはあっても、結局はその両者を充実させることが高精度自動翻訳の実現のためには必要であろう。

究極の翻訳自動化とはどのようなものかまだ眺望できる域には技術は達していないが、言語処理において人間と同等のパフォーマンスを求めるとすれば、やはり人間なみの言葉の理解が必要であると考えられるのは自然であろう。しかし冒頭でも述べた通り、これまでの機械翻訳は言葉の真の意味の理解を経ずに翻訳を行おうとしてきた。これは、そもそも「意味」とは何で、どのように計算機上で意味を表現し計算処理したらよいのか、という人工知能の根源的問題にまだ明確な解が与えられていないためである。

一方で、チェスや囲碁、将棋などのゲームにおいて人間並みかそれ以上の力量を誇るプログラムが出現しているが、そのメカニズムが人間の思考と同じかといえれば答えは否であることを考えると、翻訳においても人間とは異なる「思考」に基づくアプローチも可能かもしれない。これらのプログラムは人間のように振る舞わなくても、定石などの人間の知識の断片は解探索のためのヒューリスティックスとして大量に組み込まれている。また、

人の声を計算機で生成する音声合成においても実際の人間の肉声の断片をたくさん計算機内に蓄積しておき、状況に応じて断片を繋ぎ合わせて人間の発声に近い音を生成する手法が近年多く使われるようになってきている。このようにみると、機械翻訳においても人間の翻訳知識の断片を限りなく蓄積しておき状況により繋ぎ合わせる手法が有望かもしれない。Example-Based MTはまさにそういった観点に基づくものであるが、ただ単に用例を集めるというだけでは、言葉のバリエーションが無限に近く存在するためすべての可能な表現を蓄積するのは不可能に近い。そこで文法的に同等の表現は同じ規則に従うと仮定することにより解析や生成の規則を用いて「断片知識」の数を減らすルールベースの手法が必要になる。またExample-Based MTとルールベースMTをつなぎ合わせて使おうという試みも近年行われているが、どのように繋ぎ合わせたら最も効果的かという問題はまだまだほとんど手つかずの状態にある。

今回紹介した機械翻訳と翻訳メモリの融合による翻訳支援システムの構築は、直接の目的は人間の翻訳作業を効率化させようというものだが、その先にExample-Based MTとルールベースMTの最適な結合方法を探るという意図も込められている。すなわちまずは人間に、断片再利用ができる高機能翻訳メモリといったExample-Based MTの素材と、ルールベース機能を備えた機械翻訳による翻訳結果を提供し、人間による翻訳作業を効率化させる。その過程で、人間がどのように翻訳メモリの一部と機械翻訳結果の一部を繋ぎ合わせて訳文を完成させていくかを分析しモデル化することができれば、「継ぎはぎ」の作業をさらに効率化するためのツールの作成が期待できる。このようにツールを用いた人間の翻訳作業の各工程を分析し、各工程の効率化のためのまた別のツールを用いた半自動化を行い、ついには自動化が果たされるならば、はじめから全自動を前提とした従来の機械翻訳とは動作原理がまったく異なる機械翻訳システムの構築が可能かもしれない。

### 参考文献

- 1) Sato, S. and Nagao, M.: Toward Memory-based Translation, Proceedings of the 13th International Conference on Computational Linguistics, pp.247-252 (1990).
- 2) <http://www.lisa.org/>
- 3) 潮田 明, 富士 秀, 大倉清司, 山下達雄: 機械翻訳と訳例検索を統合した翻訳支援システム, 言語処理学会第9回年次大会予稿集 (2003).
- 4) 富士 秀, 畠中伸敏, 伊藤悦雄, 亀井真一郎, 隈井裕之, 介弘達也, 吉見毅彦, 井佐原均: 機械翻訳システムの有効性の評価~どのような人にとってMTは役立つか~, 言語処理学会第8回年次大会予稿集 (2002).
- 5) 山下達雄, 富士 秀, 大倉清司, 潮田 明: 翻訳支援に有効な訳例検索の類似度計算方式と検索結果提示方式, 言語処理学会第9回年次大会予稿集 (2003).
- 6) 野口正一 監修, 牧野武則 著: 図解 自然言語処理, オーム社, pp.2-3 (1991).

(平成15年6月2日受付)

