

発表原稿を利用した講演の同時翻訳システム*

5 Q - 6

柏岡 秀紀, 松本 賢司, 田中 英輝†

ATR 音声言語通信研究所‡

e-mail: {kashioka,kmatsu,tanakah}@slt.atr.co.jp

1 はじめに

音声翻訳システムの利用形態として、講演やニュースなどの独話を対象とした同時通訳は重要な課題であり、現在、幾つかの研究が進められている[1, 2, 3]。同時翻訳システムを構築するには、同時性を保つための認識や翻訳の処理時間、翻訳の精度など多くの課題がある。

講演の同時通訳では、講演の発表原稿やキーワードリストなどの情報が予め得られることが多い。一般的な翻訳システムでは、この様な情報が利用出来る場面はほとんどない。これらの情報を効率的に利用するシステムの構築が期待される。本稿では、予め発表原稿が得られるという設定でのシステムについて検討する。

実際の講演で発表原稿が得られることは稀であると思われる。しかしながら、身近な状況で我々自身が国際会議で発表するときに、一度、日本語で発表原稿を作成し、英語に翻訳、添削を行い、英語の発表原稿、およびスライドを作成し、発表に臨むという手順を取る人も少くない。日英翻訳を対象に、発表原稿とその英訳を利用した同時通訳システムを考え、実現可能性をさぐるために予備検討を行った。実際に国際会議で発表を行った研究員に協力してもらい、3つの講演で発表の書き起しと発表原稿との比較実験を行った。本稿では、比較実験の概要とその結果に基く予備検討について報告する。

2 発表原稿と発話の照合

本稿で想定しているシステムは、講演の発表原稿が予め与えられる所に特徴がある。基本的な考えは、「講演発表に伴い発表原稿と発話を順次照合し、一致する原稿に対応する訳文の音声を合成する。」というものである。予め準備された英語を出力するため、訳質の高い出力が期待できる。また、照合の単位を工夫することで、同時性を保つことも期待できる。本稿では、発表原稿と発話との照合について検討する。日本語の発表原稿に対する英語原稿の作成については、本稿の議論の対象外とする。

発表原稿と発話との一致を判定するためには、少なくとも以下の課題を考慮する必要がある。

- 照合の単位：文(節)、句、文節、単語など
- 類似表現の照合：同義語、言い換えなど
- 照合結果の出力：照合の成否、出力のタイミングなど

Machine Translation Model Using Manuscript for Monologue

Hideki Kashioka, Kenji Matumoto, Hideki Tanaka
ATR Spoken Language Translation Research Laboratories

- 照合失敗時の処理：照合の結果、適切な出力がない、あるいは過剰に出力がある場合の処理

ここでの照合は、発表原稿を想定したものであるが、発表と照合するデータを、予め得られる情報から構築した対訳データベースに置き換えることができる。これにより、発表原稿が得られるという稀な状況設定の制約を取り扱えるが、前述の課題に加えて、訳質の高さを保つ為の工夫が必要となる。

3 実験

国際会議で発表する 3 名の研究員(以下、A、B、C とする)の協力により、国際会議での発表内容に関して日本語で発表原稿を作成し、日本語の発表を収録し、書き起しを行った。各発表は、10 分から 20 分程度の内容である。

各研究員の日本語の発表原稿は、以下の手順により作成された。A、B の 2 名は、国際会議の後、実際に利用した英語の発表原稿を翻訳し、日本語原稿を作成した。両研究員の発表内容は、これまでに研究会等で数回発表を行っている。C は、国際会議発表準備として、日本語で発表練習したものを書き起し、日本語の発表原稿を作成した。国際会議の英語の発表原稿は、それを翻訳し、修正を加えたものである。また、本発表内容については、これまでに 1 度、日本語での発表をしている。

収録した日本語での発表は、A、B は、より発表に近い形式で会議室でスライドを投影しながら原稿を見ずに発表したものである。C の収録は、収録室でスライドを投影せず原稿を見ながら発表したものである。

実験では、音声認識結果を利用する望ましいが、予備検討ということもあり、簡単のために発表原稿と発表の書き起しを比較した。発表原稿と書き起しの比較に際し、単語表記の揺れ、書き起しの誤りは修正し、言い誤りを削除し、Juman(Ver.3.61)を利用して形態素に分割し比較した。表 1 に発表原稿と書き起しの特徴を示す。

全体としての発表原稿と書き起しの一致の割合を見るために、共通の形態素の述べ出現数の一致率を見ると、80% から 90% が一致していることがわかった。そこで、文を単位として考え、書き起しと一致している原稿の文の対応を取った。表 2 にその結果を示す。対応する文は、語順も考慮した共通の形態素部分列の割合が高い文を一致した文と考えた。ただし、対応する文を取り出すのに、原稿、書き起しの各文に順に番号を付け、書き起しの文番号 m に対して、原稿の m ± n 文の範囲にある文の中から共通の形態素部分列の割合の高い文を取り出している。n の値として、A、B では 5、C で

表 1: 原稿と書き起しの特徴

	A	B	C
原稿文数	71	78	122
書き起し文数	73	82	127
原稿異り形態素数	307	338	618
書き起し異り形態素数	311	357	658
共通異り形態素数	283	307	501
原稿延べ形態素数	1392	1643	2794
書き起し延べ形態素数	1469	1785	3470
共通延べ形態素数	1273	1500	2458

は 10 とした。文単位の認定は、原稿については作成者、書き起しについては転記者の主観による。

表 2: 比較結果

	A	B	C
延べ一致形態素数の割合	0.89	0.88	0.78
文対応の割合の平均	0.79	0.75	0.62

ここで得られた文の対応で、過不足なく内容が一致しているもの、部分的に一致しているもの、全く一致していないものの割合を、表 3 に示す。

表 3: 文対応の内容の一致

	A	B	C
過不足なく一致	65	73	89
部分的に一致	6	8	28
不一致	2	1	10

4 考察

本節では 2 節で示した課題について検討する。まず、照合の単位であるが、前節で示した文の対応では、表 3 より 70% から 90% 近くが過不足なく照合できている。このことから、照合の単位として文は適切であると思われるが、文の切れ目の判断を考慮すると、ポーズ等で区切られる節や句等の単位での照合が期待される。また、部分的に一致している文は、細かくみれば以下の様に分類できる。具体例とともに以下に示す。

- 複数の文により完全に一致

書き起し：

まず最初に直接音に対して指向性を形成します。

その後 L-1 個の反射音に対して指向性を形成し それらの信号を同相化し... できます。

原稿：

最初に直接音に対して指向性を形成し 次に L-1 個の反射音に対して指向性を形成します。

それからそれらの出力信号を同相化して... できます。

- 原稿にない表現が含まれる書き起し：

この時目的とする音がシータ方向から到来する場合各マイクロホンで受音した信号には 到来時間差 タウ が生じます。

そこでこの タウ を補正遅延器によって... できます。

原稿：

そこでこの到来時間差を遅延器により... できます。

- 原稿の表現が削除されている書き起し：

一致数は単純に何語一致したかです。

原稿：

一致数は単純に 詰語集合のうち 何語一致したかです。

これらの例をみると、文より短い単位の照合が望まれる。ただし、対応する単位の照合が成功しても、出力する翻訳の適切な単位とはいえない。

次に類似表現の照合であるが、類似表現にも色々なタイプが考えられる。“線形”と“直線”的な単語同士、“システムの構築を行うため”と“システムを実現するため”の様な句、節の単位での類似表現がある。

さらに出力として、対応のとれた文の番号をみると、発表原稿の文番号が前後することはほとんどなく、全体的に話の流れは一致しているといえる。スライドを利用しているため、原稿とほぼ同じであり、発表の順序が入れ替わったとしても、同じスライド内の入れ替わりであると予想される。また、原稿の内容で書き起しに現れていない内容はほとんどなかった。対応する番号で原稿の文番号の現れていないものは、「後で具体例を示します」などの付加的なものと、発表の最後のまとめで、今後の予定を説明している部分であった。

5 まとめ

本稿では、発表原稿を利用した同時翻訳システムの構築にむけて、発表原稿と実際の発表講演の比較を行うことで、その有効性、実現可能性について検討した。3 件という小数の発表についての比較結果であるが、スライドを利用した発表において、原稿と発表書き起しとの対応は、文を単位にある程度実現できることがわかった。ただし、対応する文の抽出法、照合の単位や類似表現の照合などについては、今後の十分な検討が必要である。

参考文献

- [1] 渡辺, 松原, 外山, 稲垣: “英日同時翻訳のための漸進的日本語生成”, 言語処理学会第 6 回年次大会, pp.272-275, 2000.
- [2] 丸山, 熊野, 柏岡: “日本語における独話の特徴と文分割”, 言語処理学会第 7 回年次大会発表論文集, pp.429-432, 2001.
- [3] 柏岡: “講演の同時通訳データ作成と分析”, 信学技報, 思考と言語, TL2000-33, pp.61-66, 2000.