

# バイリンガル旅行会話コーパスに見られる 話し言葉の特徴分析

竹沢 寿幸 † 白井 諭 † 大山 芳史 ‡

† ATR 音声言語通信研究所

‡ NTT コミュニケーション科学基礎研究所

あらまし 大語彙かつ多様な表現の扱える音声翻訳を目指して、新聞記事など書き言葉の機械翻訳と会話文など話し言葉の機械翻訳の扱う内容の違いに関する実際的かつ定量的な調査分析を実施した。まず、語彙的な特性に関する数値情報を報告する。次に、用言と格要素相当語句から構成される基本構文表現の日英対訳に関する数値情報を報告する。さらに、技術の現状を把握するための事例研究として、ATR 音声翻訳通信研究所で構築された日英音声翻訳システム ATR-MATRIX が扱っている範囲を論じる。

## Characteristics of colloquial expressions in a bilingual travel conversation corpus

Toshiyuki TAKEZAWA † Satoshi SHIRAI † Yoshifumi OYAMA ‡

† ATR Spoken Language Translation Research Laboratories

‡ NTT Communication Science Laboratories

**Abstract:** In order to develop a speech translation system that has a large vocabulary and accepts various expressions, we have carried out a practical and quantitative investigation of similarities and differences among the tasks of machine translation for written language such as newspaper articles, and those for spoken language such as conversations in daily life. First, we mention the characteristics of vocabularies for colloquial expressions in conversations. Next, we report the characteristics of basic sentence patterns which consist of one predicate and essential case phrases from the viewpoint of translation from Japanese to English. Finally, we discuss state-of-the-art technology based on a case study of a Japanese-to-English speech translation system ATR-MATRIX built by ATR Interpreting Telecommunications Research Laboratories.

### 1 まえがき

現在の音声翻訳は、限定されたタスクへの適用を想定し、比較的小規模の語彙(約13,000語)と、表現の種類をある程度限定して翻訳実験を行なっている[1, 2, 3]。音声翻訳の適用範囲の拡大を目指すには、大語彙、多様な表現、分野適応性の問題を解決する必要がある。新聞記事等

の書き言葉向きの機械翻訳システムは音声翻訳システムに比較して大語彙を扱っているが、テキストの翻訳と会話の翻訳は違った難しさがあるという指摘[4]がある。会話文など話し言葉向きの機械翻訳システムの研究[5]はあるものの、書き言葉向きの機械翻訳システムとは別に研究がなされており、会話文など話し言葉の機械翻訳とテキストなど書き言葉の機械翻訳の扱う内

容の違いに関する分析調査はこれまで十分に  
なされていなかった。

近年 NTT が構築した産業経済記事など記述  
文を対象とする日英機械翻訳システム ALT-J/E  
[6] の意味辞書の一部が「日本語語彙大系」とし  
て出版され [7]、その CD-ROM 版 [8] も利用で  
きるようになった。一方、ATR 音声翻訳通信  
研究所が収集した「バイリンガル旅行会話コー  
パス」 [9, 10] も広く公開されている。

実際的かつ定量的な調査分析データは書き言  
葉向きの機械翻訳システムで会話調の話し言葉  
を扱う際の指針や有益な知見を与えるものと期  
待できる。そこで、機械処理可能な大語彙辞書  
である日本語語彙大系を用いて電子化されたバ  
イリンガル旅行会話コーパスの特徴分析を試み  
た。まず、コーパスに含まれる単語を抽出し、  
それらが日本語語彙大系の単語体系に含まれて  
いるか調査した [11]。次に、用言と格要素相当  
語句から構成される基本構文表現を抽出し、そ  
れらが日本語語彙大系の構文体系に含まれてい  
るか、含まれているならばさらに英訳の妥当性  
を調査した。書き言葉向きの大語彙辞書を対照  
データとすることで、会話特有の語彙、単語の  
連結により構成される表現の特徴を検討する。  
本稿では、それらの数値情報を報告するととも  
に、音声翻訳システムの現状について議論す  
る。

2 で分析調査対象であるバイリンガル旅行会  
話コーパスと日本語語彙大系の概要について述  
べる。3 で旅行会話に現れた会話特有の語彙に  
ついて述べる。4 で旅行会話に現れた会話特有  
の基本構文表現について述べる。5 で音声翻訳  
システムの技術の現状について議論する。最後  
に 6 で全体をまとめる。

## 2 分析調査対象

### 2.1 バイリンガル旅行会話コーパス

ATR 音声翻訳通信研究所では音声翻訳研究の  
ためにバイリンガル旅行会話コーパス [9, 10] を  
構築した。良質で大量の基礎資料を得るため、  
通訳者を介したバイリンガル会話を収集した。  
また実用性を考えて、タスクとしては多くの入

表 1: バイリンガル旅行会話コーパスの概要

収集会話数	618
異なり話者数	71
異なり通訳者数	23
発話総数	16,107
日本語形態素延べ数	301,961

表 2: 日本語語彙大系の概要 (部分)

収録データ	件数
単語体系	30 万語
構文体系	6,000 用言 14,000 文型

に利用可能なホテル予約を中心とした旅行会話  
を選んだ。旅行会話は日本語話者、英語話者と  
日英方向、英日方向の 2 名の通訳者の合計 4 名  
によってなされる。話者の一方はホテルのフロ  
ント係であり、他方は外国人の旅行者である。  
2 名の通訳者は音声翻訳システムの代りとして  
振る舞う。具体的には、話者の 1 回の発話は 10  
秒以内とし、相手が話している間に割り込むこ  
とは禁止した。そして、通訳者はそのような 1  
回の発話毎に逐次的に通訳を行なう。このよう  
な制約を設けることにより、極端に長い発話や  
発話の重なりを避けることができるため、「近  
未来の音声翻訳システム」研究の基礎資料とし  
て適していると判断できた。表 1 にバイリン  
ガル旅行会話コーパスの概要を示す。関連情報  
は <http://results.atr.co.jp/products/> をご覧い  
ただきたい。

### 2.2 日本語語彙大系

NTT では、産業経済記事など記述文を対象と  
する日英機械翻訳システム ALT-J/E [6] を実現  
し、その意味辞書の一部が日本語語彙大系とし  
て出版されている [7]。日本語の語彙 30 万語を  
3,000 種類の意味属性で分類し、さらに 6,000  
語の用言には日英の文型 14,000 パターンが付与

表記形 | 読み | 標準形 | 品詞 | (出現頻度)

の | ノ | の | 格助詞 | (69)  
 の | ノ | の | 終助詞 | (4)  
 の | ノ | の | 準体助詞 | (337)  
 の | ノ | の | 連体助詞 | (4135)

図 1: 助詞「の」の例

されている。表 2 に日本語語彙大系のうち本稿で調査対象とした項目の概要を示す。分析調査にはその CD-ROM 版 [8] および必要に応じて適宜 ALT-J/E [6] の辞書情報を活用した。関連情報は <http://www.kecl.ntt.co.jp/icl/mtg/resources/GoiTaikei/> をご覧いただきたい。

### 3 旅行会話に現れた会話特有の語彙

バイリンガル旅行会話コーパスのうち日英方向の発話を対象に、単語の抽出と頻度のカウントを行なった [11]。助詞「の」の例を図 1 に示す。記号 | で区切られた項目は左から順に「表記形」「読み」「標準形」「品詞」を表し、() 内の数値が出現頻度である。

表 3 に助詞の例を示すように、旅行会話コーパスの品詞体系 [12] は日本語語彙大系や ALT-J/E の辞書情報とは異なる。そこで、その違いを考慮して、表 3 の () 内に記すように、旅行会話コーパスの準体助詞は日本語語彙大系や ALT-J/E の形式名詞に、旅行会話コーパスの係助詞は日本語語彙大系や ALT-J/E の副助詞に、旅行会話コーパスの並立助詞、連体助詞、引用助詞は日本語語彙大系や ALT-J/E の格助詞に対応させて照合を行ない、日本語語彙大系の単語体系に基づく被覆率を調査した。

図 1 の例では、格助詞「の」、準体助詞「の」、連体助詞「の」は日本語語彙大系や ALT-J/E に同等とみなせる語彙項目が含まれていると判断できた。しかし、終助詞「の」は日本語語彙大系や ALT-J/E には含まれていなかった。したがって、異なり語数を基準とする被覆率は 3/4 となり、75.0% である。延べ語数を基準とする被覆率は 4541/4545 となり、

表 3: 品詞体系の違い (助詞の例)

旅行会話コーパス	日本語語彙大系や ALT-J/E
格助詞	格助詞
準体助詞	なし (形式名詞)
係助詞	なし (副助詞)
副助詞	副助詞
並立助詞	なし (格助詞)
接続助詞	接続助詞
終助詞	終助詞
連体助詞	なし (格助詞)
引用助詞	なし (格助詞)

表 4: 会話特有の語彙の割合

	延べ語数	異なり語数
含まれるもの	80,685	2,493
含まれないもの	18,828	1,269
被覆率	81.1%	66.3%

99.9% である。

バイリンガル旅行会話コーパスには言い直し等の単語の断片が含まれており、それらは形態素の品詞に「その他」と付与されている。また、バイリンガル旅行会話コーパスでは活用する語は語幹と語尾に分割している。そこで、品詞「その他」「語尾」等を除き、被覆率を求めた。結果を表 4 に示す。異なり語基準による品詞毎の被覆率を図 2 に示す。

日本語語彙大系の単語体系が含まない単語は大きく二つに分類できる。一つは会話調の話し言葉特有の表現であり、もう一つは旅行という分野に依存する語彙である。それぞれの例を次に示す。

- 会話調の話し言葉特有の表現: 主に感動詞、副詞、接続詞、助詞、助動詞等

「ありがとうございました(感動詞)」、  
 「いらっしゃいませ(感動詞)」、「すみませんが(副詞)」、「そうしますと(接続詞)」、「じゃ(格助詞)」、「の(終助詞)」、「ちゃ(助動詞)」等

- 旅行という分野に依存する語彙: 主に普通

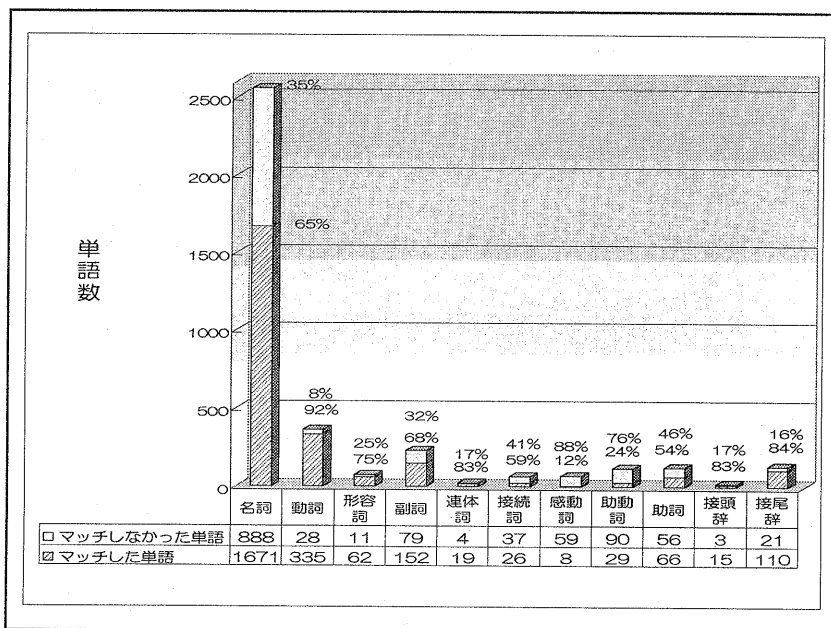


図 2: 異なり語基準による品詞毎の被覆率

名詞、固有名詞等

「さば寿司(普通名詞)」、「カナディアンロッキー(固有名詞)」、「嵐山線(固有名詞)」等

日本語語彙大系の単語体系に含まれない単語のうち、名詞類が異なり語数で 888 語、延べ語数で 7,261 語を占める。大まかに名詞類がすべて旅行という分野に依存する語彙であると近似すれば、日本語語彙大系の単語体系に含まれない単語のうち旅行という分野に依存する語彙の占める割合は、異なり語数基準で約 70%、延べ語数基準で約 40% である。逆に、その残りが会話調の話し言葉特有の表現の占める割合である。

#### 4 旅行会話に現れた会話特有の基本構文表現

バイリンガル旅行会話コーパスのうち日英方向の発話を対象に、用言と格要素相当語句から

構成される基本構文表現の抽出を行なった。いわゆる「ダ文」は会話調の言い回しに多用される。その事例の一部を図 3 に示す。内容項目は左から順に出現頻度、日本語構文、英語構文である。出現頻度は日本語構文と英語構文のペアを単位として求めた。日本語構文はダ文表現と格要素相当語句から構成される。ダ文表現の前後を / で囲んで示し、ダ文の「だ」の直前に # を付与する。「です」等の表層表現となっているものもすべて「#だ」の形で抽出した。主題を表す係助詞「は」で格助詞「が」に置き換え可能な場合は「が((は))」と記述した。助詞を伴わない格要素相当語句において格要素を示す助詞を補うことができる場合は \* で囲んだ形式でそれを記述した。英語構文全体は " で囲んで示し、さらに日本語のダ文表現に相当する部分を / で囲んで記述する。英語に人称代名詞が含まれる場合はそれを one に置き換えた。英語構文として連結する必要はないが、日本語構文の対訳として必要な要素は、| で区切って記述した。

- 1 / 雨 #だ / "it /rain/"
- 3 / 初めて #だ / "/be/ one's first time"
- 1 / 初めて #だ / "/be/ one's first visit"
- 2 / 初めて #だ / "/be/ the first time"
- 1 奈良が ((は)) / 初めて #だ / "this /be/ one's first trip | to Nara"
- 1 着物が ((は)) / 初めて #だ / "/be/ one's first time | to put kimono on"
- 1 着物が ((は)) / 初めて #だ / "/be/ one's first time | wearing a kimono"
- 1 それは / お困り #だ / "that /be/ a problem"
- 1 バスを / ご利用 #だ / "/take/ the bus"
- 1 こちらが / チケット #だ / "here /be/ one's tickets"
- 1 サービス料 \* が \* / 込み #だ / "/include/ service charges"
- 1 東京成田空港を / 出発 #だ / "/leave/ Tokyo Narita Airport"
- 1 一つ目が / 東寺駅 #だ / "the first stop /be/ Toji Station"
- 1 茶室が / 入り用 #だ / "/need/ a tearoom"
- 3 / ご存じ #だ / "/know/"
- 6 / お泊まり #だ / "/stay/"

図 3: いわゆる「ダ文」の例 (一部)

表 5: 会話特有の基本構文表現の割合

	一般		ダ文	
	個数	割合	個数	割合
英訳妥当	327	9.1%	1	0.2%
英訳ほぼ決まる	1413	39.5%	173	26.1%
訳し分け必要等	1840	51.4%	488	73.7%
合計	3580	100.0%	662	100.0%

同様にして動詞等の一般的な用言とその格要素相当語句から構成される基本構文表現も抽出した。得られた数値情報を表 5 に示す。延べ基準でも異なり基準でも割合は大きく変わらないため、異なり基準の数字のみ示す。

動詞等の一般的な用言とその格要素相当語句から構成される基本構文表現に対して、日本語語彙大系の構文体系で妥当な英訳が得られるものは全体の約 1 割であった。例えば「地下鉄烏丸線に乗る」という表現で「地下鉄烏丸線」が鉄道であることがわかれば妥当な英訳 “/take/ the Karasuma subway line” が得られる。残り

の約 9 割は日本語語彙大系の構文体系では妥当な英訳が得られない。しかしながら、約 4 割のものは英語の用言がほぼ決まるものであった。代表的な事例は次の通りである。

- 外来語: / ファックスする / → /fax/  
「ファックスする (fax)」「チェックインする (check-in)」という言い回しが会話調の話し言葉ではしばしば使われる。このような事例は英語の用言がほぼ決まる。
- 敬語: 予約を / 確認いたす / → /confirm/ one's reservation  
接客業務のような会話では特に敬語が多用される。敬語は日本語語彙大系の構文体系には含まれていない。「予約を / 確認いたす /」のように、英語の用言がほぼ決まるものが多数ある。
- 敬語表現: マイクロバスを / ご利用に / なる / → /use/ the shuttle service  
「なる」そのものは敬語ではないが、「お(ご)・・・になる」という言い回しは

尊敬の表現となる。「お(ご)」と「になる」の間の内容語に関する情報を活用することで、英語の用言がほぼ決まる。

- 名前の伝達: わたくしが((は))|...と申す  
→ my name /be/ ...

接客業務のような会話では「申す」「いう」は名前を伝える形式で使われることがほとんどである。「申す」という用言は日本語語彙大系の構文体系に含まれているが、このような英訳は含まれていなかった。

さらに、英語の用言が一つに絞りきれないものは約5割あった。代表的な事例は次の通りである。

- 多義:

例えば「いらっしゃる」という表現は「居る(exist)」と「来る(come)」という二つの意味がある。また「文楽が/いい/」という表現は「文楽を勧める(/recommend/ bunraku)」と「文楽が素晴らしい(bunraku /be fine/)」という二つの意味がある。

- 婉曲表現:

会話調の言い回しでは「朝食が((は))|別に/なる/ → /there be/ a separate charge for breakfast」や「昼の部が|[数]時に/なる/ → the matinee /start/ at [num] a.m.」、「|[数]円と/なる/ → /be/ [num] yen」等の「になる」「となる」という言い回しが多用される。これらはいわゆる「ダ文」に相当する婉曲表現と言えるものである。

一方、いわゆる「ダ文」については、日本語語彙大系の構文体系で妥当な英訳が得られるものは「雨だ(it rains)」の1例のみであった。つまり、ほとんどすべてのものが日本語語彙大系の構文体系では妥当な英訳が得られない。約4分の1のものは「日本食が大好きだ(/like/ Japanese food)」 「茶室が入り用だ(/need/

a tearoom)」 「乗り換えが必要だ(/need/ to transfer)」のように英語の用言がほぼ決まるものであった。残りの約4分の3は「シビックを二台だ(two Civic)」に対して英語の用言として rent や reserve を訳出することが適切であるとすれば、その発話だけで適当な英語用言を訳出するのが難しい事例である。

## 5 議論および今後の課題

### 5.1 音声翻訳システムの現状

ATR 音声翻訳通信研究所で構築された日英音声翻訳システム ATR-MATRIX [1] を事例として取り上げ、音声翻訳システムの現状について考察する。会話文の特徴と考えられる「ダ文」に対して、システムから得られる事例を次にいくつか列挙する。

- 茶室が入り用です (The tea ceremony room is necessary)  
→ The tea ceremony room is necessary
- 支払いはカードです (I will pay for it by card)  
→ The payment is by the card
- シビックが二台です (Two Civic)  
→ Civic is two

「茶室が入り用です」のように英語の用言がほぼ決まると分類できるものについては妥当な英訳が得られている。それ以外の場合でも「支払いはカードです」のように頻度が多いものについては用例が整備されているものもある。しかしながら、頻度の少ない「シビックが二台です」のようなものについては日本語の構文を反映した英語の構文に対して内容語を置き換えたものが出力される。

これまでの音声翻訳システムの研究では窓口業務のような場面が想定されていた。日英双方の音声翻訳システムを用いた対話実験では、音声認識や翻訳に若干の誤りがあっても、発話者の再発話や相手からの確認質問によって、与

えられた課題は達成できることが確認されている [3]。発話全体が翻訳できない場合でも翻訳可能な部分を翻訳する部分翻訳機能 [13] が有効であるのは、窓口業務のような場面で協調的な応答が期待できる場合にはキーワードがわかるだけでも助かる状況があるためである。

しかしながら、大語彙かつ多様な表現の扱える音声翻訳を目指すならば、キーワードが伝わるだけでは不十分な状況が考えられる。会話文を対象とするとしても、構内アナウンスのように相手からの確認が期待できない場合には品質の良い翻訳が必須となる。したがって、一般的な用言の約5割、「ダ文」の約4分の3を占める訳し分けが必要な基本構文表現については、翻訳の品質を高める必要がある。

## 5.2 今後の課題

基本構文表現を抽象化することにより汎用性の向上が期待できるため、会話文の日英翻訳のための文型辞書として整備する計画である。

また、表4において、異なり語数に対する被覆率より延べ語数に対する被覆率が高いことから、限定した領域を想定すれば、そこに現れる語彙は絞れる可能性があるといえる。しかし、表5の基本構文表現は、異なり基準でも延べ基準でも割合が大きく変わらなかった。つまり、語彙は絞れたとしても、その組合わせである基本構文表現は絞りきれない、あるいは、まだコーパスの量が十分でないという解釈ができる。そのため、対訳コーパスの稠密度を把握する手段を検討する予定である。

## 6 むすび

NTT が書き言葉を対象として構築した大語彙辞書である日本語語彙大系を用いて ATR 音声翻訳通信研究所が収集したバイリンガル旅行会話コーパスの特徴分析を行なった。テキストの翻訳と会話の翻訳は違った難しさがあるという定性的ないし思索的な指摘はあったが、会話文など話し言葉の機械翻訳とテキストなど書き言葉の機械翻訳の扱う内容の違いに関する分析調査はこれまで十分になされていなかった。そ

こで、機械処理可能な大語彙辞書や電子化されたバイリンガル会話コーパスが公開されたことを背景に、実際のかつ定量的な調査分析を試みた。このような調査分析は書き言葉向きの機械翻訳システムで会話調の話し言葉を扱う際の指針や有益な知見を与えるものと期待できる。

## 謝辞

分析作業に協力いただいた NTT アドバンステクノロジー株式会社 サービスシステム事業部 言語処理システム部の皆様に感謝する。

## 参考文献

- [1] T. Takezawa, T. Morimoto, Y. Sagisaka, N. Campbell, H. Iida, F. Sugaya, A. Yokoo, and S. Yamamoto, "A Japanese-to-English speech translation system: ATR-MATRIX," *Proc. International Conference on Spoken Language Processing*, pp. 2779-2782, 1998.
- [2] E. Sumita, S. Yamada, K. Yamamoto, M. Paul, H. Kashioka, K. Ishikawa, and S. Shirai, "Solutions to problems inherent in spoken-language translation: the approach of ATR-MATRIX," *Proc. Machine Translation Summit*, pp. 229-235, 1999.
- [3] F. Sugaya, T. Takezawa, A. Yokoo, and S. Yamamoto, "End-to-end evaluation in ATR-MATRIX: speech translation system between English and Japanese," *Proc. EUROSPEECH*, pp. 2431-2434, 1999.
- [4] 長尾真, "情報技術の新時代に向けて," 情報処理, Vol. 41, No. 1, pp. 48-49, 2000.
- [5] 古瀬蔵, 山本和英, 山田節夫, "構成素境界解析を用いた多言語話し言葉翻訳," 自然言語処理, Vol. 6, No. 5, pp. 63-91, 1999.

- [6] 八巻俊文, 大山芳史, 白井諭, 横尾昭男, “機械翻訳特集: 日英機械翻訳システム ALT-J/E の研究開発,” *NTT R&D*, Vol. 46, No. 12, pp. 1391-1398, 1997.
- [7] 池原悟, 宮崎正弘, 白井諭, 横尾昭男, 中岩浩巳, 小倉健太郎, 大山芳史, 林良彦 (編), “日本語語彙大系,” 岩波書店, 1997.
- [8] 池原悟, 宮崎正弘, 白井諭, 横尾昭男, 中岩浩巳, 小倉健太郎, 大山芳史, 林良彦 (編), “日本語語彙大系 CD-ROM 版,” 岩波書店, 1999.
- [9] T. Morimoto, N. Uratani, T. Takezawa, O. Furuse, Y. Sobashima, H. Iida, A. Nakamura, Y. Sagisaka, N. Higuchi, and Y. Yamazaki, “A speech and language database for speech translation research,” *Proc. International Conference on Spoken Language Processing*, pp. 1791-1794, 1994.
- [10] 竹沢寿幸, 中村篤, 隅田英一郎, “ATR の会話音声翻訳研究用データベース,” *音声研究*, Vol. 4, No. 2, pp. 16-23, 2000.
- [11] 竹沢寿幸, 大山芳史, “書き言葉向き大語彙辞書を用いたバイリンガル旅行会話コースの特徴分析,” *言語処理学会第6回年次大会発表論文集*, pp. 75-78, 2000.
- [12] 竹沢寿幸, “音声言語データベースの日本語形態素情報マニュアル — 最終版 —,” *ATR テクニカルレポート*, TR-IT-0315, 1999.
- [13] 脇田由実, 河井淳, 飯田仁, “意味的類似性を用いた音声認識正解部分の特定法と正解部分のみ翻訳する音声翻訳手法,” *自然言語処理*, Vol. 5, No. 4, pp. 111-125, 1998.