

# 漸進的な音声認識・機械翻訳・テキスト音声合成に基づく音声から音声への同時翻訳\*

○中村 哲, Sashi Novitasari, △帖佐 克己, 柳田 智也  
△二又 航介, 須藤 克仁, Sakriani Sakti (奈良先端大)

## 1 はじめに

音声翻訳は、異なる言語を話す人々のコミュニケーションを支援する技術として研究開発が続けられてきた。音声翻訳は、通常、音声認識、機械翻訳、テキスト音声合成から構成されるが、昨今の深層学習技術の発達によって著しい性能向上が進んでいる。1980年代から始まった種々の研究により、旅行における会話など比較的文が短く、限られたドメインにおいて、発話単位での音声翻訳が実現されつつある。一方で、同時通訳のように発話の終了を待たずに漸進的な翻訳を行う同時音声翻訳の研究が本格化したのはこの10年ほどである [1, 2]。しかし、主として機械翻訳部分の漸進的翻訳手法の要素技術研究に留まっており、同時音声翻訳システム全体の研究はこれからというところである。本稿では、漸進的な音声言語処理に基づく同時音声翻訳システムの実現に向けた研究について紹介する。

## 2 自動同時翻訳の課題

### 2.1 通訳と翻訳

人間による同時通訳 (Simultaneous Interpretation) は、通訳対象の発話を聞き取りながら別の言語への通訳を行い発声するという、非常に高度な専門技能を必要とするタスクである。書き言葉の文書に対する翻訳が静的な入力を対象としたタスクであり、前後の文脈を考慮し、時に外部資料を参照しながら時間をかけて精緻に訳文構成を行うのに対し、話し言葉の発話に対する通訳は動的な入力を対象としたタスクであり、事前の資料と直前までの文脈だけを利用し、情報を要約したりしながら、限られた時間で訳出を行う。本研究では、このような同時通訳にみられるような情報の補完や要約については現時点では扱わないこととし、入力音声に対する漸進的な処理に基づく同時音声翻訳 (Simultaneous Translation) をタスクとして取り組むこととする。

### 2.2 同時通訳における遅延と「順送りの訳」

同時通訳における種々の課題については文献 [3] に詳しいが、その中でも重要な課題の一つとして挙げ

られているのが、言語間の統語構造の違いによって必然的に生じうる訳文構成上の遅延である。文献に挙げられている例を以下に示す。

まず、次の英文を日本語に通訳することを考える (括弧つき数字は説明のために付されたものである)。

(1) The relief workers (2) say (3) they don't have (4) enough food, water, shelter, and medical supplies (5) to deal with (6) the gigantic wave of refugees (7) who are ransacking the countryside (8) in search of the basics (9) to stay alive.

この文を日本語に翻訳すると以下のように訳文となる (括弧つき数字は英文と対応する日本語の節や句を示す)。

(1) 救援担当者は (9) 生きるための (8) 食料を求めて (7) 村を荒らし回っている (6) 大量の難民たちの (5) 世話をするための (4) 十分な食料や水、宿泊施設、医薬品が (3) 無いと (2) 言っています。

一見して分かる通り、(2) の say は日本語では文末に、(9) の to stay alive は日本語では主語の直後に訳出されており、その間は英語と逆順となっている。このような訳文を同時通訳において実現しようとする、通訳者は (2) の動詞を保持したまま (9) までを聞き取り、そこから聞き取った内容を逆順に日本語にして発話する必要がある、こうした処理では、英文を聞き終わるまで通訳発話を開始できないため非常に大きな遅延を生じてしまう。そこで、以下のように訳出を行うことで記憶負荷と遅延の減少を図る「順送りの訳」が用いられる。

(1) 救援担当者たちの (2) 話では (4) 食料、水、宿泊施設、医薬品が (3) 足りず (6) 大量の難民たちの (5) 世話ができないとのこと。 (7) 難民たちは今村々を荒らし回って、 (9) 生きるための (8) 食料を求めているのです。

\*Simultaneous Speech-to-speech Translation based on Neural Incremental ASR, MT, and TTS. by NAKAMURA, Satoshi et al. (Nara Institute of Science and Technology)

この訳では、英文の要素を前から小分けにして訳出し、(7)以降の関係詞節はその手前で一旦文を区切って、関係詞節の内容は「難民」を補足する文として付け加えることで日本語としての自然さを損なわないようにしている。これは英語が主辞前置型 (head-initial) 言語、日本語が主辞後置型 (head-final) 言語であって語順が大きく異なることに起因する。そのため、順送りの訳のような工夫なしでは英語と日本語の間の同時通訳は困難であると言える。

### 2.3 自動同時通訳における遅延

機械による同時音声通訳を構築する場合にも、音声認識、機械通訳、音声合成それぞれの遅延が問題となる。機械通訳においては前節において述べたような言語の構造の違いによって大きな遅延の可能性があるが、順送り通訳を行う仕組みの構築が不可欠である。音声認識では、発話中にデコーディング可能ではあるものの、昨今の双方向 LSTM を用いる方法では発話終了を検出して認識する方法が多い。テキスト音声合成では基本的に発話文が確定してから全体の韻律、音声が生成されるため、漸進的に音声合成をすることはそのままでは難しい。これらが連結されると全体として大きな遅延となり、発話が長いと遅延が致命的となる。本研究では漸進的音声言語処理の技術を統合し、統語構造の差によって生じる遅延の削減を目指す。

## 3 漸進的音声言語処理

本節では同時音声通訳システムを構成する、音声認識・機械通訳・テキスト音声合成における漸進的処理のための手法について簡単に述べる。技術の詳細や個々の性能評価についてはそれぞれ文献 [4, 5, 6] を参照されたい。

### 3.1 漸進的音声認識

音声認識においても注視機構付き系列変換 (attentional sequence-to-sequence) モデルが広く用いられているが、通常注視の対象が文単位の状態系列であることから、漸進的な処理に対応できない。漸進的な処理を実現するために、後方の文脈を参照しないような特殊なモデルや学習方法が提案されている [7]。

我々の研究 [4] では、文全体を入力して注視するモデルを教師 (teacher) とし、漸進的処理のために短いセグメント単位で注視を行うモデルを生徒 (student) として、生徒が教師の注視を再現できるように音声認識の学習を行う手法を提案した。遅延を最小限に留めるために各セグメントの情報のみで音声認識を

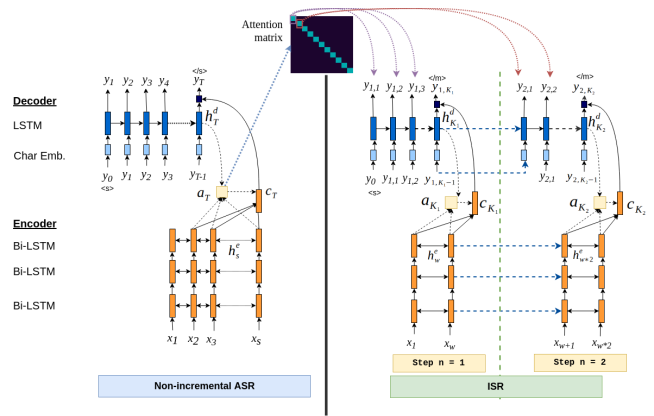


Fig. 1 インクリメンタル音声認識

行うと十分な精度が得られないため、400ms 精度の遅延を許容して対象セグメントの後方の音声特徴量も利用することで文単位の入力を利用した場合からの精度低下を抑えられることが実験的に確認されている。

### 3.2 漸進的機械通訳

機械通訳ではすでに述べたように語順の違いによって低遅延での訳出が難しい場合がある。順送りの訳の実現には依然としてデータ量が不足していることもあるため、現在は多くの研究において (順送りではない) 通常通訳を行った対訳コーパスから通訳の学習を行っているのが実情である。そうした状況で低遅延での同時通訳を実現する方法として提案されたのが *wait-k* [8] と呼ばれる、入力トークン列に対して  $k$  トークンの入力を待ってから通訳出力を開始する方式である。ある時点での訳語選択に必要な情報がそれ以前の入力で得られていない場合は、それ以前の入力から強制的に訳語選択を行うこととなり、ある種の予測として機能する。*wait-k* は非常に単純な方式で実装も容易だが、英語と日本語の間のような語順の差が大きい場合には不十分である。

我々の研究 [5] では後段の入力を適応的に待つ手段として、デコーダの出力記号の一つにトークンを出力せず次の入力を待つことを表す特殊記号を追加し、訳語選択に必要な入力得られていない場合に適応的に入力を待つ方式を提案した。英語から日本語への通訳実験においては、*wait-k* では十分な入力得られず過度な予測を求められるのに対して、提案手法は適応的に入力待機を行い漸進的な通訳による精度低下を小さく抑えられることが確認されている。

### 3.3 漸進的テキスト音声合成

テキスト音声合成における漸進的な処理は、音声認識や機械通訳に比べて従来の取り組みが少ない研究課題と言える。テキスト音声合成でも合成音の予測

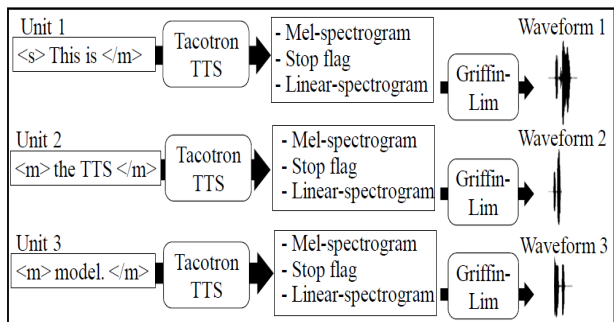


Fig. 2 インクリメンタル音声合成

には周辺の単語から得られる特徴量が不可欠であり、漸進的な処理のために特徴量の予測を HMM に基づくテキスト音声合成に組み込んだ手法が提案されている [9]. ニューラルネットワークに基づく系列変換モデルによる end-to-end 処理はテキスト音声合成でも活用されてきているが、漸進的な処理についてはこれまで試みられていなかった。

我々の研究 [6] では、単語（英語の場合）やアクセント句（日本語の場合）を単位として入力テキストをセグメントに分割し、セグメントごとに音響パラメータ（スペクトログラム）の予測やセグメント終端の予測を行う、漸進的な end-to-end テキスト合成手法を提案した。提案手法を利用した主観評価実験により、1 単語／アクセント句のみの情報に基づく音声合成よりも、多少の遅延を許容して 2-3 単語／アクセント句の情報を利用した音声合成のほうが自然性が高いことが確認されている。

#### 4 音声から音声への同時翻訳システム

前節で述べた漸進的音声認識・機械翻訳・テキスト音声合成技術を利用して、音声から音声への同時翻訳を実現する試作システムを作成した。本試作システムは英語の講演音声日本語の音声に翻訳するものである。

本試作システムは各モジュールを単純にカスケード接続したもので、以下のような手順で処理を行う。

- マイクもしくは音声／動画ファイル入力<sup>1</sup>の英語音声に対する漸進的音声認識を行い、その結果を機械翻訳モジュールに渡す。
- 音声認識モジュールから得られた英語の音声認識結果を日本語に翻訳し、その結果をテキスト音声合成モジュールに渡す。
- 機械翻訳モジュールから得られた日本語への翻

<sup>1</sup>ファイル入力時はファイル読み込みが音声の実時間より高速である点には注意を要する。

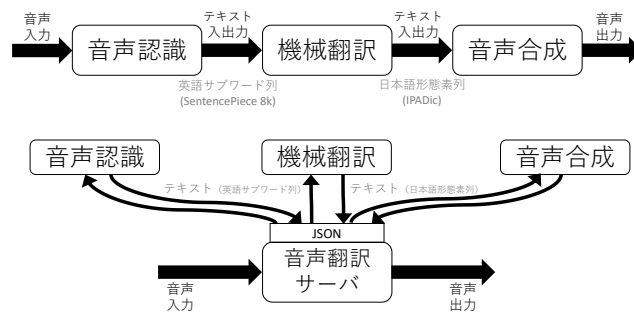


Fig. 3 試作システムの構成例.

訳結果を音声合成し、スピーカーもしくは音声ファイルに出力する。

なお、モジュール間での処理単位の変換を最小限に留めるため、音声認識結果が後段の機械翻訳の入力側と同じサブワードモデルに基づいたサブワード列として得られるように音声認識モデルを学習し、機械翻訳結果が後段のテキスト音声合成で用いる IPADic に基づく日本語形態素の列の形で得られるように機械翻訳モデルを学習した。学習データはそれぞれの論文に記載のものに加え、TED Talks 英語講演と日本語字幕のデータを利用した（テキスト音声合成モデルを除く）。

本試作システムにおける各モジュール間の接続は (1) テキストによる標準入出力（パイプ接続）(2) 処理統括サーバとの相互通信のいずれかで行う設計とした。単一の入力音声に対する処理であれば、各モジュールを別プロセスで駆動した (1) の構成で動作させることが可能である（図 3）。

#### 5 同時通訳データの収集

前節で述べた同時音声翻訳技術・システムの研究に加え、本研究のための講演同時通訳データの収集を継続的に行っている。本稿執筆時点までに、TED Talks を中心に英語から日本語で約 150 時間、日本語から英語で約 110 時間の同時通訳を収集しており、今後も TED 等の講演以外のデータも含め継続していく予定である。

#### 6 課題と今後の展望

本試作システムは音声から音声への同時翻訳を志向した漸進的音声認識・機械翻訳・テキスト音声合成技術の連携によって実現したものである。同時翻訳システム全体としての性能向上のためには当然各モジュール単位での精度および処理効率の向上が必須であるが、それと同時にモジュール間の接続において 1-best の処理結果だけでなく n-best や単語ラティス

のように曖昧性を含んだ結果を渡し、それを考慮した処理が行われることが好ましい。また、近年音声から音声への end-to-end 翻訳 [10] が注目を集めつつあり、同時翻訳においてそうしたアプローチの有用性についての検討が必要であろう。

音声から音声への同時翻訳の今後の展望として、実際の応用において自動同時翻訳がどの程度有用であるかを検証・評価することが考えられる。同時翻訳に関する研究 [2, 8] では BLEU 等の翻訳精度と遅延の大きさのトレードオフとしてその性能を議論してきたが、実際には情報の受け手にとって有用であったか、という観点での議論が必要であろう。さらに、同時通訳のように時間制約の中で情報を要約したり、外部知識や講演資料等に基づいて発話内容を予測したりする等、より通訳に近い処理の実現も将来の目標と考えることができよう。

## 7 おわりに

本稿では、漸進的な音声認識・機械翻訳・テキスト音声合成技術に基づく音声から音声への同時翻訳のアプローチと、その試作システムについて述べた。今後は技術の検討と実際の同時通訳データの蓄積をさらに進める予定である。

謝辞 本研究の一部は JSPS 科研費 JP17H06101 の助成を受けたものである。

## 参考文献

- [1] Srinivas Bangalore, Vivek Kumar Rangarajan Sridhar, Prakash Kolan, Ladan Golipour, and Aura Jimenez. Real-time incremental speech-to-speech translation of dialogs. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 437–445, Montréal, Canada, June 2012. Association for Computational Linguistics.
- [2] Tomiki Fujita, Graham Neubig, Sakriani Sakti, Tomoki Toda, and Satoshi Nakamura. Simple, Lexicalized Choice of Translation Timing for Simultaneous Speech Translation. In *Proceedings of Interspeech*, pp. 3487–3491, 2013.
- [3] 水野的. 同時通訳の理論—認知的制約と訳出方略. 朝日出版社, 2015.

- [4] Sashi Novitasari, Andros Tjandra, Sakriani Sakti, and Satoshi Nakamura. Sequence-to-Sequence Learning via Attention Transfer for Incremental Speech Recognition. In *Proceedings of Interspeech 2019*, pp. 3835–3839, 2019.
- [5] 帖佐克己, 須藤克仁, 中村哲. 英日同時翻訳のための Connectionist Temporal Classification を用いたニューラル機械翻訳. 情報処理学会研究報告 2019-NL-241, 2019.
- [6] Tomoya Yanagita, Sakriani Sakti, and Satoshi Nakamura. Neural iTTS: Toward Synthesizing Speech in Real-time with End-to-end Neural Text-to-Speech Framework. In *Proceedings of the 10th ISCA Speech Synthesis Workshop*, pp. 183–188, 2019.
- [7] Kyuyeon Hwang and Wonyong Sung. Character-level incremental speech recognition with recurrent neural networks. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5335–5339, March 2016.
- [8] Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng, Kaibo Liu, Baigong Zheng, Chuanqiang Zhang, Zhongjun He, Hairong Liu, Xing Li, Hua Wu, and Haifeng Wang. STACL: Simultaneous Translation with Implicit Anticipation and Controllable Latency using Prefix-to-Prefix Framework. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3025–3036, Florence, Italy, July 2019. Association for Computational Linguistics.
- [9] Timo Baumann. Decision tree usage for incremental parametric speech synthesis. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3819–3823, May 2014.
- [10] Takatomo Kano, Sakriani Sakti, and Satoshi Nakamura. Structured-based curriculum learning for end-to-end english-japanese speech translation. In *Proc. Interspeech 2017*, pp. 2630–2634, 2017.