

構音障害者の少量学習データによる音声合成の検討*

☆南阪竜翔, 滝口哲也, 有木康雄(神戸大)

1 はじめに

本研究では脳性麻痺などにより健常者と異なった発話となる構音障害者のコミュニケーション支援としてのテキスト音声合成(Text-To-Speech)及び声質変換を提案する。

構音障害者は健常者と比べて多くの発話データを収録することが困難である。しかしDNNによるテキスト音声合成では多くのデータ量が必要となる。そこでまず健常者で音声合成モデルを生成し、得られた合成音を声質変換することで構音障害者の音声合成システムを作成する。

テキスト音声合成とは、任意に与えられたテキストから対応する音声を合成する技術である。これまでテキスト音声合成を実現するための多くの手法が提案されており、従来手法として隠れマルコフモデル(Hidden Markov Model: HMM)を用いたもの[1]が最も代表的であった。

近年ではDNN(Deep Neural Network)を用いた音声合成[2]が、従来の隠れマルコフモデルを用いた場合と比べ高音質の合成音が作成できるためDNNを用いた音声合成が主となっている。音声合成は障害者支援にも用いられ、山岸ら[3]は様々な人々の音声を集めデータベースを構成し、ALS患者の為のTTSシステムを構築した。

声質変換とは、入力音声に対して、発話内容は保持しながら、話者性や感情といった情報を変換する技術である。統計的声質変換の代表例としてGMM(Gaussian Mixture Model)[4]がある。近年では深層学習に対する注目の高まりからDNNを用いたもの[5]もある。これらの声質変換には同一内容、同一発話長の発話が必要であるが、それを必要としない適応型制限ボルツマンマシン(Adaptive Restricted Boltzmann Machine)を用いた研究[6]も行われている。

本稿では、Bidirectional LSTM[7]によるテキスト音声合成とGMM(Gaussian Mixture Model)による声質変換を用いた少量学習データによる音声合成について述べる。

2 提案手法

テキスト音声合成ではBidirectional LSTMを用いる。声質変換にはGMMを用いる。Fig. 1に本研究の概要を示す。

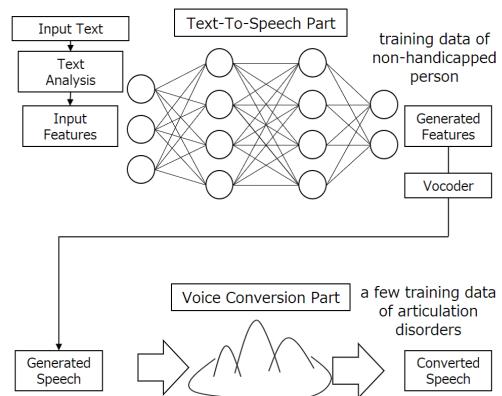


Fig. 1: A flow of proposed speech synthesis system

2.1 DNN テキスト音声合成

まずテキストから当該音素、アクセント型、フレーム位置などの情報を抽出する。その情報をバイナリや連続値で表現し、フレームで連結したものとを言語特徴量としてDNNの入力データとする。教師データとして、音声からボコーダーを用いて抽出された音声パラメータと動的特徴量をフレームで連結したものを用いる。

音声合成時は任意のテキストから言語特徴量を抽出してモデルに入力する。そして得られた出力から最尤推定を用いて音響特徴量を推定し、ボコーダーを用いて音声を合成する。

2.1.1 Bidirectional LSTM 音声合成

Bidirectional LSTMはRecurrent Neural Networksに基いて構成される。隠れ層を $\mathbf{h} = (h_1, \dots, h_T)$ 、出力を $\mathbf{y} = (y_1, \dots, y_T)$ 、入力を $\mathbf{x} = (x_1, \dots, x_T)$ 、時系列を $t = 1 \dots T$ とすると Recurrent Neural Network (RNN)における、隠れ層の状態と出力の関係は以下のように

* Speech synthesis system using small amounts of data for articulation disorders, by Ryuka N, Tetsuya Takiguchi, Yasuo Ariki (Kobe univ.)

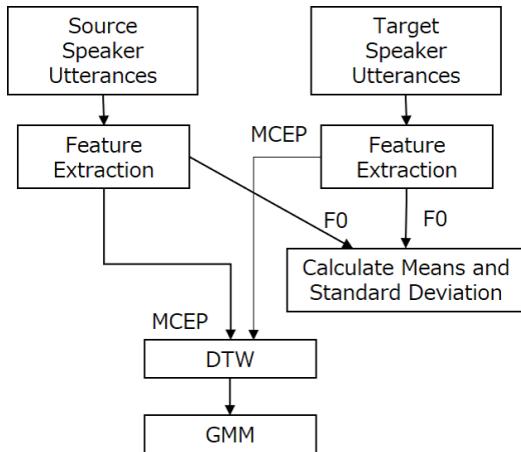


Fig. 3: GMM Voice Conversion system

3 評価実験

3.1 実験条件

実験データには構音障害者の男性1名、健常者の男性1名を使用した。健常者の選択にはヘッドホンを用いた聴取実験により構音障害者に一番声質が近い話者を選択し、ATR 音素バランス450文を学習、43文を評価、10文をテストに用いて実験を行った。サンプリング周波数は16kHz、フレームシフトは5msとした。

言語特徴量にはコンテキストラベルに対してHTS[8]形式のQuestionを適用して抽出した979次元を使用する。音響特徴量にはWORLD[9]を用いて抽出したスペクトル包絡にメルフィルタバンクを使用し、低次60次元メルケプストラム係数、対数基本周波数F0、帯域非周期性指標1次元とそれらの2次までの動的特徴量、1次元の有声無声パラメータを用いた。

言語特徴量は最小値0、最大値1となるように、次元ごとに正規化を行った。音響特徴量は平均0分散1となるように、正規化を行った。

声質変換においてはテキスト音声合成において得られた音声と構音障害者のATR音素バランス50文を学習に、10文をテストに用いた。GMM学習にはスペクトル包絡にメルフィルタバンクを使用し、低次60次元メルケプストラム係数を用いた。基本周波数は式(15)で変換を行った。帯域非周期性指標は入力話者のものを用いた。発話長はDPマッチングにより同一発話長とした。

評価基準として、メルケプストラム歪み(Melcepstrum Distortion:MelCD)を用いた。比較対象として、構音障害者の少量のデータのみを用いた合成音を用いた。学習データはそれぞれ50

文または、100文を用いて比較した。

3.2 実験結果・考察

声質変換によって得られた音声のスペクトログラムをFig. 4に示す。

構音障害者の音声は健常者に比べて高周波成分が弱い傾向がある。変換後のスペクトログラムにおいても2000Hz以上のパワーが弱くなっている。また低周波数域のフォルマントにおいても変換前に比べ構音障害者に近づいていることが分かる。

示した音声は「青い青い海は女性の美しさを持っている。」という文であるが、特に「持っている。」の部分において変換前のスペクトログラムから構音障害者のスペクトログラムへと近づきが見られた。今回入力話者のdurationをそのまま用いたが、構音障害者のdurationを用いることで更に話者性を考慮した音声合成が可能であると考える。

MelCDによる評価結果をFig. 5に示す。評価結果はテスト10文の平均を取っている。構音障害者の50文、100文で音声合成システムを構築した場合と比べ、提案手法の(健常者合成+)声質変換を用いることにより、学習データ50文でもいい結果が得られた。

4 おわりに

本稿では、Bidirectional LSTMによるTTSとGMMによる声質変換を組み合わせた構音障害者の少量学習データによる音声合成についての提案を行った。

Fig. 4より構音障害者のスペクトログラムにおける高周波成分特徴や低周波におけるフォルマントが上手く学習されていることを確認した。今後は少量の構音障害者データのみ用いたテキスト音声合成と提案手法の音声合成の主観評価による比較を行う。

また、パラレルデータを得るためにDPマッチングを行ったが、アライメント誤りが音質に影響している可能性がある。アライメントを必要としない非パラレルデータ声質変換についても検討を行い、比較する。

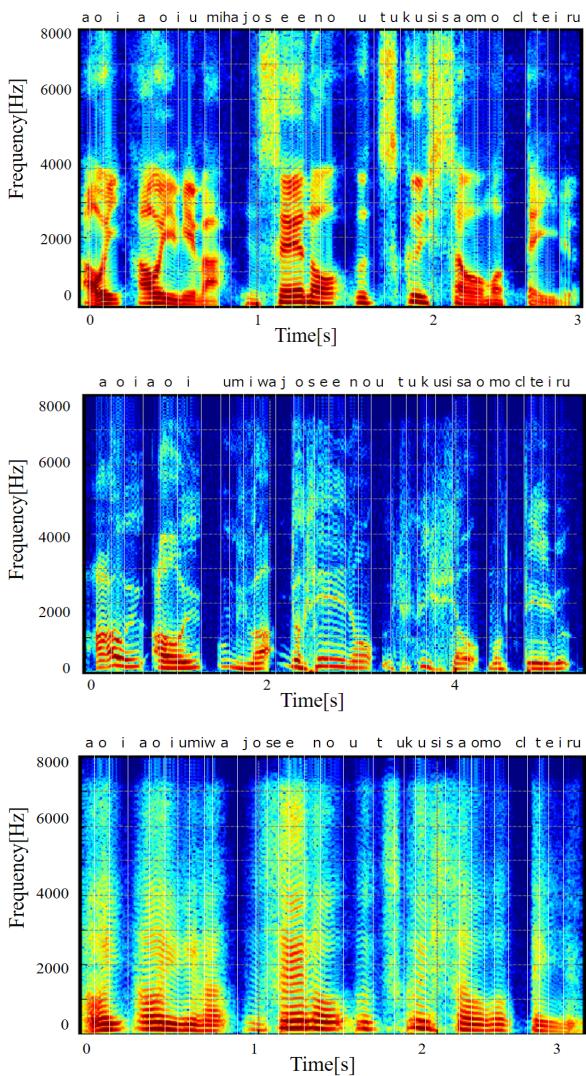


Fig. 4: Comparison of spectrogram
(Top: Source, Middle: Target, Bottom: Converted)

参考文献

- [1] K. Tokuda *et al.*, “Speech parameter generation algorithms for HMM-based speech synthesis,” in *ICASSP*, 2000, pp. 1315-1318.
- [2] H. Ze *et al.*, “Statistical parametric speech synthesis using deep neural networks,” in *ICASSP*, 2013, pp. 7962-7966.
- [3] J. Yamagishi *et al.*, “Speech synthesis technologies for individuals with vocal disabilities: Voice banking and reconstruction,” in *Acoustical Science and Technology*, vol. 33, no.1, pp. 1-5, 2012.
- [4] T. Toda *et al.*, “Voice Conversion Based on Maximum-Likelihood Estimation of Spectral Parameter Trajectory,” in *IEEE*

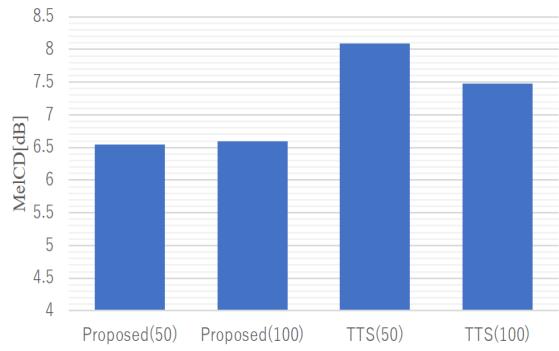


Fig. 5: Mel-cepstrum Distortion

Transactions on Audio, Speech, and Language Processing, Vol. 15, No. 8, 2007.

- [5] S. Desai *et al.*, “Voice conversion using Artificial Neural Networks,” in *ICASSP*, 2009.
- [6] T. Nakashika *et al.*, “Non-Parallel Training in Voice Conversion Using an Adaptive Restricted Boltzmann Machine,” in *IEEE Transactions on Audio, Speech, and Language Processing*, Vol. 24, No. 11, 2016
- [7] Y. Fan *et al.*, “TTS Synthesis with Bidirectional LSTM based Recurrent Neural Networks ,” in *INTERSPEECH 2014*.
- [8] Z. HeiGa *et al.*, “HMM 音声合成システム (HTS) の開発 ,” in *IPSJ SIG Technical Reports*, 2007.
- [9] M. Morise *et al.*, “WORLD: a vocoder-based high-quality speech synthesis system for real-time applications ,” in *IEICE TRANSACTIONS on Information and Systems*, vol. 99, no. 7, pp. 1877-1884, 2016.