

構音障害者のための話者性を維持した HMM 音声合成システムの提案 *

上田怜奈 (神戸大), 滝口哲也 (神戸大/JST さきがけ), 有木康雄 (神戸大)

1 はじめに

本研究では脳性麻痺から起こる構音障害を持つ人々を支援するための音声合成法を提案する。構音障害者にとって健常者との会話は困難を伴うものである。障害者支援のための研究において, *Veaux et al.* [1] は筋萎縮性側索硬化症 (ALS) 患者のための話者性を維持した音声再構築を試みた。また, 山岸ら [2] は様々な人々の音声を集めデータベースとし, それを用いた ALS 患者のための TTS システムを構築した。構音障害者のコミュニケーションの障害となりうる要因としてはピッチ, スペクトルなどの問題が挙げられる。このような問題に対して, 本研究では構音障害者の話者性を維持しつつ聞き取りやすさを向上させる HMM 音声合成システムを提案する。学習データに構音障害者と健常者音声を使用し, 聞き取りにくさの原因となる構音障害者音声の成分を健常者音声によって補充し, より聞き取りやすい音声を実現する。評価実験を通して本研究の提案法が障害者の話者性は維持しつつより聞き取りやすい合成音声を実現出来ていることを示す。

2 構音障害者のための HMM 音声合成

構音障害者の音声は収録した段階で不安定な音声となっているため, 構音障害者の音声から得られた音声特徴でパラメータ学習をすると得られる合成音は聞き取りづらいものになってしまう。そこで本研究では, 話者性の近い健常者と構音障害者の両方の音声を学習データとして, 話者性は維持しつつより聞き取りやすい合成音を作成した。Fig. 1 は提案手法の概要である。提案法において, 構音障害者と健常者の両方を学習データとして使用する。初めに, STRAIGHT[3] を用いて二人の話者から 3 つの音声パラメータ (F0 概形, スペクトラム包絡, 非周期成分 (AP)) を抽出する。特徴量を抽出したのち, 健常者の F0 系列を修正する (2.1 節)。音素継続長モデルについては構音障害者, 健常者それぞれのコンテキスト依存ラベルからそれぞれのモデルを作成したのち修正を行い, 修正後音素継続

長モデルを得る (2.2 節)。その後, 修正後音素継続長モデルから生成したコンテキスト依存ラベル系列と学習した HMM に基づいて, スペクトラム, F0, AP パラメータが生成される。F0 パラメータは修正した F0 モデルから生成される。AP パラメータは構音障害者のモデルから生成する。スペクトルパラメータは構音障害者と健常者両方のモデルから生成され, その後スペクトル修正を行う (2.3 節)。最後に, パラメータ系列を合成フィルタにかけることによって合成音が生産される。2.1 節, 2.2 節, 2.3 節では F0・音素継続長・スペクトルに対する処理の詳細を記述する。

2.1 F0 系列の修正

構音障害者の F0 系列はしばしば不安定なものであるため, 本研究の F0 の修正法では, 健常者の F0 系列を基本として F0 モデルを学習する。F0 系列に構音障害者の話者性を付与するため, F0 系列を構音障害者の特徴へと変換する。F0 モデルはこの変換後の F0 系列を用いて学習するので, 構音障害者の話者性が含まれていることになる, F0 系列の変換には Eq. (1) のような線形変換を利用する。

$$\hat{w}_t = \frac{\sigma_x}{\sigma_w} (w_t - \mu_w) + \mu_x \quad (1)$$

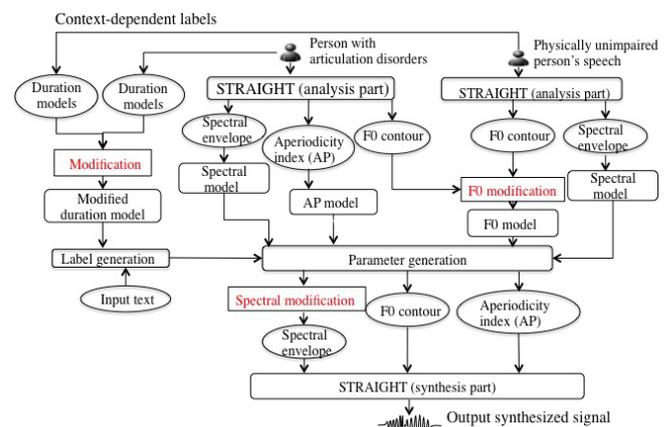


Fig. 1: 構音障害者のための HMM 音声合成手法の概要

*Individuality-Preserving Voice Reconstruction for Articulation Disorders Using HMM Text-to-Speech Synthesis, by Reina Ueda (Kobe University), Tetsuya Takiguchi (Kobe University/JST PRESTO), Yasuo Arika (Kobe University)

Type	母音	子音	無音区間
平均	4.54	4.91	6.25
分散	48.73	60.76	137.67

(a) 構音障害者

Type	母音	子音	無音区間
平均	3.29	3.65	4.89
分散	11.26	14.75	58.26

(b) 健常者

Table 1: 学習後音素継続長モデルの平均, 分散

Eq. (1) において, w_t は健常者の t フレーム目の対数 F0, μ_w, σ_w は健常者の F0 系列の平均・分散, μ_x, σ_x は構音障害者の対数 F0 系列の平均・分散をそれぞれ表している.

2.2 音素継続長の修正

構音障害者の話速は健常者のものと比べて間延びしたのとなっており, このことが聞き取りにくさの原因となっている. Table 1 は構音障害者, 健常者の学習データからそれぞれ音素継続長モデルを作成し, 母音・子音・無音区間ごとに分布の平均, 分散を算出したものである. 各話者の数値を見比べると, 平均, 分散ともどの場合でも構音障害者の数値が高くなっていることが分かる. 特に, 子音の平均値や母音, 子音, 無音区間の分散値が高くなっていることが合成音の聞き取りにくさに影響していると考えられる. そこで, 本研究では音素継続長モデルを修正し, 話者性は維持しつつより聞き取りやすくなる音声を作成する. 提案法では, 健常者の音素継続長モデルをベースとして修正を行う (Fig. 2). そして, 健常者のモデル中の母音の平均値に対して修正を行う. 修正はノードごとに以下のように行う.

$$\hat{y}_i = y_i - \mu_y + \mu_z \quad (2)$$

$$\mu_y = \frac{\sum_{i=1}^I \mu_{yi}}{I} \quad (3)$$

$$\mu_z = \frac{\sum_{i=1}^I \mu_{zi}}{I} \quad (4)$$

Eq. (2) において, y_i は健常者音素継続長モデル中の i 番目のノードの平均値, μ_y, μ_z は Eq. (3), Eq. (4) のようにして求められる. Eq. (3), Eq. (4) において, I はモデル内の母音の全ノード数, μ_{yi} は健常者モデルの i 番目の母音ノードの平均値, μ_{zi} は構音障害者モデルの i 番目の母音ノードの平均値をそれぞれ表して

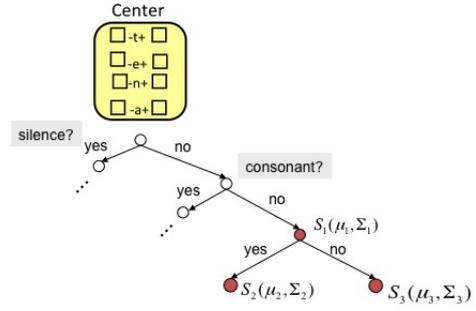


Fig. 2: 健常者音素継続長モデルの修正

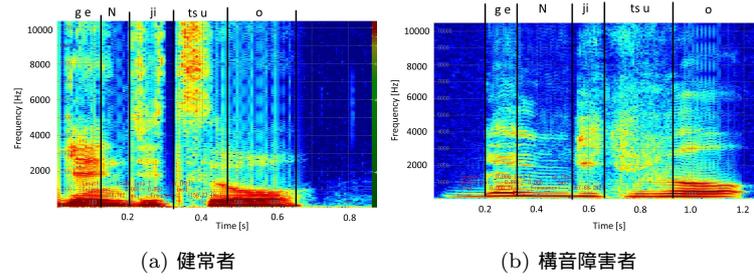


Fig. 3: 元音声スペクトルの一例 // geNjitsu o

いる.

2.3 スペクトル系列の修正

Fig. 3 は健常者と構音障害者の元音声の“現実を”と発声しているスペクトログラムである. Fig. 3 にあるように, 構音障害者のスペクトルの高周波成分は健常者のものと比べて弱くなっている. これは構音障害者の発声の子音成分が弱くなっておりそのことが聞き取りにくさの原因となっていることを示している. そこで Fig. 1 のようにテキストが入力された後, それぞれの話者のスペクトルモデルからスペクトルパラメータを生成する. そして高周波成分を健常者のスペクトルパラメータで補完し, 低周波域は構音障害者のスペクトルパラメータを使用し話者性を維持しつつより聞き取り易くなるように修正を行う. このような修正はすべての音素に対して行うのではなく, 摩擦音, 破擦音等高周波成分にパワーを持つ音素 (sh/s/z/ch/ts/j) に対してのみ行い, その他の音素に対しては修正を行わず構音障害者のスペクトルパラメータを使用する. この修正は以下の式で実現される.

$$\hat{S}^{(ij)} = f_{PU}^{(j)} S_{PU}^{(ij)} + f_{AD}^{(j)} S_{AD}^{(ij)} \quad (5)$$

このとき, $S_{PU}, S_{AD}, \hat{S}, i, j$ はそれぞれ健常者スペクトル (Physically Unimpaired), 構音障害者スペクトル (Articulation Disorder), 修正後スペクトル, フレームのインデックス, 周波数次元のインデックスを

示している．重み関数 f_{PU} , f_{AD} は以下のように定義される．

$$f_{PU}^{(j)} = \frac{1}{1 + e^{(-j+c)}} \quad (6)$$

$$f_{AD}^{(j)} = \frac{1}{1 + e^{(j-c)}} \quad (7)$$

このとき, f_{PU} は健常者スペクトルに対する重み関数, f_{AD} は構音障害者に対する重み関数, c は制御変数をそれぞれ表している．Eq. (5) を用いることにより, 高周波領域では健常者のスペクトル成分によって補完され, より子音部分が明瞭に聞こえるようにし, 低周波領域では構音障害者のスペクトル成分を保持することにより話者性を保つということを実現する．周波数の閾値を制御する変数 c は Eq. (6) によって閾値が 4000Hz になるように設定する．

$$c = \frac{4000}{f_s} \times D \quad (8)$$

Eq. (8) において, f_s はサンプリング周波数, D はスペクトルの次元数を表している．

3 評価実験

3.1 実験条件

学習データには構音障害者の男性 1 名, 健常者男性 1 名を使用した．健常者選定の際, 障害者と話者性が大きく異なる健常者音声を採用すると生成される合成音の話者性も低下する恐れがある．そこで本研究では ATR データベースセット B の話者 10 名と構音障害者間の話者性の類似度を求めるためにスペクトラムの話者間の距離をそれぞれ算出し, その値が最も小さい話者を採用した．その結果, 今回は健常者音声として男性話者 MTK の音声を採用した．健常者音声は ATR データベース 503 文, 障害者音声は収録した同じデータベース中の 429 文を使用した．特徴量についてはサンプリング周波数は 16 kHz, フレームシフト 5 ms で音声特徴量は STRAIGHT を用いて抽出し, スペクトルパラメータ系列は, 25 次元のメルケプストラムとその Δ , $\Delta\Delta$ を使用, 学習・合成には 5 状態のコンテキスト依存 HMM を使用した．提案法の有効性を示すため本研究では話者性と聞き取りやすさの 2 つの観点から実験を試みた．10 人の日本人に対してヘッドホンで聴取実験を行った．話者性に関する実験には本研究では DMOS (Degradation Mean Opinion Score) テストを実施した．このテストではリファレンス音声と評価対象音声を聞き比べ, 評価対象音声と比べて劣化しているかを 5 段階 (5:劣化が全く認めら

Table 2: 実験で比較した合成音の生成条件

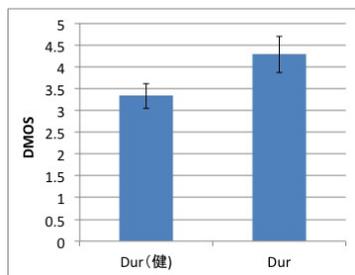
Type \ Model	Duration	F0	AP	Spectral
ADM	AD	AD	AD	AD
Dur(健)	PU	AD	AD	AD
Dur	Mod	AD	AD	AD
Dur_F0	Mod	Mod	AD	AD
Dur_Spe	Mod	AD	AD	Mod
Dur_F0_Spe	Mod	Mod	AD	Mod

(AD: 構音障害者, PU: 健常者, Mod:修正有り)

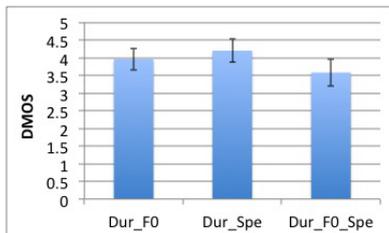
れない 4:劣化が認められるが気にならない 3:劣化がわずかに気になる 2:劣化が気になる 1:劣化が非常に気になる) で評価した．聞き取りやすさの実験は一対比較法で行った．聴取実験は二回に分けて行い, 一回目の実験では音素継続長修正の効果を, 二回目の実験では F0・スペクトル修正の効果をそれぞれ検証した．

3.2 実験結果

一回目の実験では音素継続長修正の効果について検証した．実験にあたり, Table 2 のうち ADM, Dur(健), Dur の 3 種類の合成音を用意した．Fig. 4a は話者性に関する実験結果である．Fig. 4a より, Durの方が Dur(健) よりも優位な結果となった．このことから, 健常者の音素継続長モデルをそのまま使うよりも修正後音素継続長モデルを使う方が話者性が保たれることが分かった．聞き取りやすさに関する実験では ADM と Dur, Dur と Dur(健) をそれぞれ比較した．Fig. 5a より, ADM よりも Dur のほうが優位であるとわかる．ここから, 修正後音素継続長モデルを利用するほうが, 構音障害者の音素継続長モデルを利用するよりも聞き取りやすさは向上することがわかった．Fig. 5b より, Dur(修) と Dur(健) がほぼ同じ評価となった．ここから, 健常者の音素継続長モデルを利用した時と, 構音障害者の音素継続長モデルを利用したときでは聞き取りやすさはほぼ同じで遜色はないことが分かる．以上の結果より音素継続長モデルの修正が話者性と聞き取りやすさの両面から見て有効であることがわかった．二回目の実験では F0 修正, スペクトル修正の効果について検証した．Fig. 4b は話者性に関する実験結果である．ここでは Dur をリファレンスとして 5 段階 DMOS で Dur_F0, Dur_Spe, Dur_F0_Spe を評価した．Fig. 4b より, Dur_Spe, Dur_F0, Dur_F0_Spe の順で話者性を保持出来ていることが分かる．Dur_F0, Dur_Spe の間にはそれほど大きな差は見られなかつ



(a) リファレンス音声:ADM



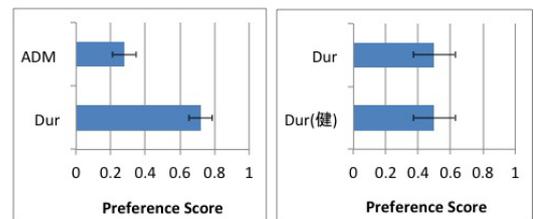
(b) リファレンス音声:Dur

Fig. 4: 話者性に関する比較

たが, Dur_F0_Spe のスコアは他の2条件と比べて低い値となっている. このことから F0 修正, スペクトル修正の両方を適用すると被験者は話者性が下がったと感ずることが分かる. 聞き取りやすさに関する実験では, Dur と Dur_F0, Dur と Dur_Spe, Dur と Dur_F0_Spe をそれぞれ聞き比べてよりどちらがより聞き取りやすいかを評価した. Fig. 5c より, Dur と Dur_F0 を比較したとき Dur_F0 のほうが大幅に優位な結果となった. これは F0 修正が聞き取りやすさに対して大きな効果があることを示している. Fig. 5d は Dur と Dur_Spe を比較したときの実験結果である. Fig. 5d において, Dur と Dur_Spe 間で有意差は確認できなかった. これは, スペクトル修正の対象となる音素が非常に少ないことが一因していると考えられる. Fig. 5e は Dur と Dur_F0_Spe を比較したときの実験結果である. Fig. 5e より, Dur と Dur_F0_Spe を比較したとき, Dur_F0_Spe のほうが優位な結果となった. このことから, F0 修正, スペクトル修正の両方を適用すると被験者は聞き取りやすさが向上したと感ずることが分かる.

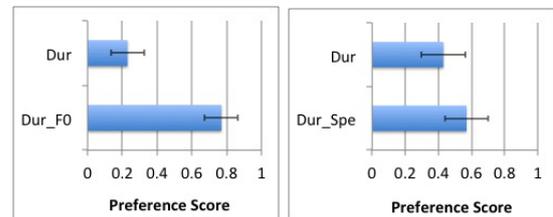
4 おわりに

本研究では構音障害者のための話者性を維持した HMM 音声合成手法を提案した. 実験を通して F0・音素継続長修正法が補正前の合成音と比較して話者性を維持し聞き取りやすい音声を実現出来ることが示された. 今後はより広い周波数域に対応したスペクトル修正法を検討していきたい.



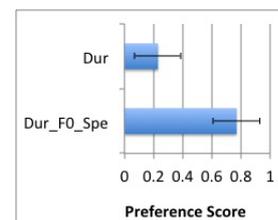
(a)

(b)



(c)

(d)



(e)

Fig. 5: 聞き取り易さに関する比較

謝辞

本研究の一部は, JST さきがけの支援を受けたものである.

参考文献

- [1] C. Veaux *et al.*, "Using HMM-based speech synthesis to reconstruct the voice of individuals with degenerative speech disorders," in *Proc. of Interspeech*, 2012.
- [2] J. Yamagishi *et al.*, "Speech synthesis technologies for individuals with vocal disabilities: Voice banking and reconstruction," *Acoustical Science and Technology*, vol. 33, no. 1, pp. 1–5, 2012.
- [3] H. Kawahara *et al.*, "Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based F0 extraction: Possible role of a repetitive structure in sounds," *Speech communication*, vol. 27, pp. 187–207, 1999.