

構音障害者の音声認識の検討

松政 宏典[†] 滝口 哲也[†] 有木 康雄[†] 李 義昭^{††} 中林 稔堯^{†††}

[†] 神戸大学工学部 〒657-8501 兵庫県神戸市灘区六甲台町 1-1
^{††} 追手門学院大学経済学部 〒567-8502 大阪府茨木市西安威 2-1-15
^{†††} 神戸大学発達科学部 〒657-8501 兵庫県神戸市鶴甲 3-11

E-mail: †mattu28@me.cs.scitec.kobe-u.ac.jp, ††{takigu,ariki,nakaba}@kobe-u.ac.jp,
†††chao55@res.otemon.ac.jp

あらまし 本稿ではアテトーゼ型脳性マヒによる構音障害者の音声認識の検討を行う。近年、音声認識技術は向上し汎用性の高いモデル作りも可能となっている。しかし、構音に障害のある人など発話スタイルが健常者と異なる人の音声認識は、健常者用の音響モデルを用いた場合、認識精度が劣化する。そこで、本稿では構音障害者の音響モデルの構築を行い、精度の改善を試みる。アテトーゼ型脳性マヒによる構音障害者は、筋肉の緊張のため1回目の発話が不安定となることがある。従来、対数スペクトルに対し離散コサイン変換を適用した MFCC (Mel Frequency Cepstral Coefficients) を特徴量として用いるが、本稿では離散コサイン変換の代わりに、PCA (Principal Component Analysis) を用いた発話スタイル正規化手法を提案し、その有効性を示す。
キーワード 構音障害, 言語障害, 脳性マヒ

A Study on Speech Recognition of a Person with Articulation Disorders

Hironori MATSUMASA[†], Tetsuya TAKIGUCHI^{††}, Yasuo ARIKI^{†††}, Ichao LI^{†††}, and Toshitaka

NAKABAYASHI^{†††}

[†] Faculty of Engineering, Kobe University Rokkodaicho 1-1, Nada-ku, Kobe, Hyogo, 657-8501 Japan
^{††} Faculty of Economics, Otemon Gakuin University Nishiai 2-1-15, Ibaraki, Osaka, 567-8502 Japan
^{†††} Faculty of Human Development, Kobe University Tsurukabuto 3-11, Nada-ku, Kobe, Hyogo, 657-8501 Japan

E-mail: †mattu28@me.cs.scitec.kobe-u.ac.jp, ††{takigu,ariki,nakaba}@kobe-u.ac.jp,
†††chao55@res.otemon.ac.jp

Abstract We investigate speech recognition of a person with articulation disorders by athetoid type of cerebral palsy. Recently, the accuracy of speaker-independent speech recognition has been remarkably improved by use of stochastic modeling of speech. However, the use of those acoustic models causes degradation of speech recognition for a person with different speech style (e.g., articulation disorders). In this paper, we have tried to build the acoustic model for a person with articulation disorders. The articulation of the first utterance tends to become unstable due to strain of a muscle and that causes degradation of speech recognition, where MFCC (Mel Frequency Cepstral Coefficients) is used as speech features. Therefore we propose a robust feature extraction method based on PCA (Principal Component Analysis) instead of MFCC. Its effectiveness is confirmed by word recognition experiments.

Key words articulation disorders, cerebral paralysis

1. はじめに

情報技術が向上し、近年、福祉分野への情報技術の適用が行われている。例えば、画像認識技術を用いた手話認識 [1] や、文書内の文字の音声化などが行われている [2]。また、音声合成を用いて、発話障害者支援のための音声合成器の作成なども行われている [3]。

音声認識技術は現在、車内でのカーナビの操作、会議においての書き起こしなど様々な環境下や場面において使用される機会が増加している。対象者が子供である場合などには精度が低下することがわかっている [4]。文献 [5] では、構音障害者音声を対象とした音響モデル適応の検証を行っているが、言語障害者などの障害者を対象としたものは非常に少ない。現在、日本だけでも構音障害者も含まれる言語障害者が3万4000人もいることから十分なニーズがあり、研究の必要性があるといえる [6] [7]。

言語障害の原因の一つとして、脳性マヒが考えられる。脳性マヒとは新生児1000人あたりおよそ2人の割合で発生する。脳性マヒの定義として厚生省は「受胎から生後4週以内の新生児までの間に生じた、脳の非進行性病変に基づく、永続的な、しかし変化しうる運動および姿勢の異常である。その症状は満2歳までに発現する。」(1968年)と定義している。

脳性マヒの原因は、中枢神経系の損傷によるものであり、それに伴う運動障害であると考えられている。発症時期として出生前(胎内感染、母体の中毒、栄養欠損など)、分娩時(脳外傷、脳出血、脳の無酸素状態、胎児黄疸、仮死状態、未熟での出産など)、出生後(脳内出血、脳炎など)の3つの要因が考えられている。

脳性マヒは次のように分類される。1) 痙直型 2) アテトーゼ型 3) 失調型 4) 緊張低下型 5) 固縮型、それぞれの症状が混合して現れる混合型もある。

本稿ではアテトーゼ型の脳性マヒによる構音障害者を対象としている。アテトーゼ型とは脳性マヒ患者の約10~15%に発生する症状であり、大脳基底核と呼ばれる視床下部、脳幹、小脳と関連を持ち随意運動、姿勢、筋緊張を調節する働きをしている部位に損傷をうけたためアテトーゼと呼ばれる不随意運動が伴う型である。アテトーゼは緊張時や意図的な動作を行う際に出現しやすい。症状は軽度から重度まで様々であり、知能障害を合併していないケースや比較的知能障害の程度が軽いケースも多いのが特徴である [8] [9]。そこで本

稿では、まず知能障害を合併していないアテトーゼ型に着目した。アテトーゼ型の構音障害者の場合、最初の動作において緊張状態により、通常よりも不安定になる場合がある。そこで本稿では、複数回連続発話の収録を行ない、アテトーゼによる発話スタイルの変動の影響を調べる。

従来の音声認識では、対数スペクトルに対し離散コサイン変換を適用したMFCCを特徴量として用いるが、本稿では離散コサイン変換ではなく2回目以降のより安定したデータを利用した、PCA (Principal Component Analysis) による発話変動にロバストな手法を提案し、その有効性を示す。

2. PCA を用いた特徴量抽出

2.1 提案手法

音声認識システムにおいて従来はMFCC (Mel Frequency Cepstral Coefficient) が音声特徴量として用いられている。MFCCではメル尺度フィルタバンクの短時間対数エネルギー出力系列に対して、離散コサイン変換 (Discrete Cosine Transformation: DCT) を適用し、ケプストラムが得られる。そして音声のスペクトル包絡成分に対応する低次ケプストラムのみを抽出し、特徴量として音声認識に用いられる。本稿では、発話スタイルの変動にロバストな特徴量抽出法として、離散コサイン変換の代わりにPCAを用いた手法を検討する(図1)。文献 [10] では残響下でのロバストな特徴量抽出として対数スペクトル上で、DCTを行わずPCAを行う事によって有効なスペクトル情報の抽出を可能としている。本稿では2回目以降の安定した発話を用いて主軸を求め、1回目の発話(調音不安定音声)に対して対数スペクトル上でPCAを適用する。本手法においては、安定した音声成分は低次に、調音不安定成分は高次に集まる。この結果PCAにより調音不安定成分の抑圧が行われると期待でき、より有効なスペクトル情報の抽出が可能となる。

2.2 発話スタイル変動成分の抑圧

短時間分析によって得られたフレーム n 、周波数 ω の観測音声を $X_n(\omega)$ 、クリーン音声を $S_n(\omega)$ とする。ここで1回目の音声を以下の式で表現する。

$$X_n(\omega) = S_n(\omega)H(\omega) \quad (1)$$

1回目発話には発話スタイル変動成分 H が畳み込まれていると考える。次に対数変換を行い観測信号を S と H の加算で表す。

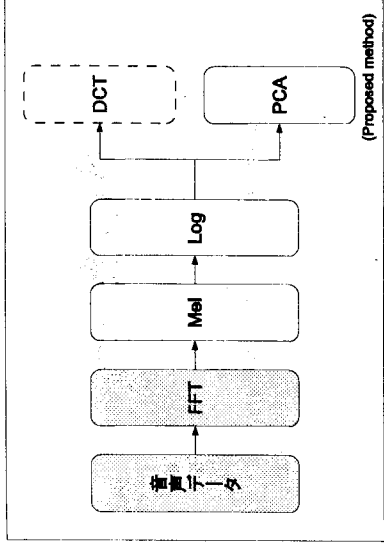


図1 PCAを用いた特徴量抽出

$$\log X_n(\omega) = \log S_n(\omega) + \log H(\omega) \quad (2)$$

ここで観測信号 X に対して PCA を適用すると、

- クリーン音声の S の主なエネルギーは D 個の主な固有値に集中する。
 - それ以外の固有値に対応する主な成分は、付加成分である。
- と期待できる。ここで、主な D 個の固有値に対応する固有ベクトルを $\mathbf{V} = [\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(D)}]$ とすると、この \mathbf{V} を用いて以下のようなフィルタを考える。

$$\hat{S} = \mathbf{V}\mathbf{X} \quad (3)$$

このフィルタリングによって付加成分 $H(\omega)$ を抑圧することが出来る。また主軸 \mathbf{V} の計算をクリーン音声のみから行えば、クリーン音声の構造のみを考慮したフィルタを作成することが出来る。本稿では、主軸 \mathbf{V} の計算には2回目以降の発話音声を用いて、上記フィルタリングを1回目の調音不安定音声に適用する。

2.3 PCA [11]

PCA (Principal Component Analysis) の目的は、データの本質的な構造を残しながら次元数を削減することにある。元の次元空間から削減した次元空間への有効的な選択方法は、分散の最大となる空間を選択することである。ここで、元の d 次元を \tilde{d} 次元に削減すると考えると、元の特徴空間から部分空間への変換行列 \mathbf{A} は、

$$\mathbf{A} = (\mathbf{u}_1, \dots, \mathbf{u}_{\tilde{d}}) \quad (4)$$

与えられる。 \mathbf{u} は d 次元からなる正規直交基底とする。特徴ベクトル \mathbf{x} は

$$\mathbf{y} = \mathbf{A}^t \mathbf{x} \quad (5)$$

に変換される。また基底の正規直交性から

$$\mathbf{A}^t \mathbf{A} = \mathbf{I} \quad (6)$$

が成立する。ただし \mathbf{I} は \tilde{d} 次元単位行列である。このときパターン数を n 、原特徴空間でのパターン平均を \mathbf{m} 、部分空間でのパターン平均を $\tilde{\mathbf{m}}$ とすると

$$\begin{aligned} \tilde{\mathbf{m}} &= \frac{1}{n} \sum_{\mathbf{y} \in Y} \mathbf{y} \\ &= \mathbf{A}^t \mathbf{m} \end{aligned} \quad (7)$$

であるから、部分空間でのパターン分散 $\tilde{\sigma}^2(\mathbf{A})$ は

$$\begin{aligned} \tilde{\sigma}^2(\mathbf{A}) &= \frac{1}{n} \sum_{\mathbf{y} \in Y} (\mathbf{y} - \tilde{\mathbf{m}})^t (\mathbf{y} - \tilde{\mathbf{m}}) \\ &= \frac{1}{n} \sum_{\mathbf{x} \in X} \text{tr}(\mathbf{A}^t (\mathbf{x} - \mathbf{m})(\mathbf{A}^t (\mathbf{x} - \mathbf{m}))^t) \\ &= \text{tr}(\mathbf{A}^t \frac{1}{n} \sum_{\mathbf{x} \in X} ((\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^t) \mathbf{A}) \\ &= \text{tr}(\mathbf{A}^t \sum \mathbf{A}) \end{aligned} \quad (8)$$

となる。ここで \sum はパターン集合の原特徴空間における共分散行列 (covariance matrix) を表す。また、 \sum はパターン \mathbf{x} の集合を、 \sum は \mathbf{x} が式 (5) で変換されたパターン \mathbf{y} の集合を表す。

式 (8) より分散を最大にする \mathbf{A} を求めることは、式 (6) の制約条件の基で $\text{tr}(\mathbf{A}^t \sum \mathbf{A})$ を最大にする \mathbf{A} を求める変分問題になる。 \mathbf{A} を \tilde{d} 次元対角行列とすると、

$$\mathbf{A}^t \sum \mathbf{A} = \mathbf{\Lambda} \quad (9)$$

となり、 \mathbf{A} は \sum を対角化する行列である。 \sum の d 個の固有値を λ_i ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$) とすれば、式 (8)、(9) より

$$\begin{aligned} \max(\tilde{\sigma}^2(\mathbf{A})) &= \max(\text{tr}(\mathbf{A}^t \sum \mathbf{A})) \\ &= \max(\text{tr} \mathbf{\Lambda}) \\ &= \sum_{i=1}^{\tilde{d}} \lambda_i \end{aligned} \quad (10)$$

となり、 $\tilde{\sigma}^2(\mathbf{A})$ を最大にする変換行列 \mathbf{A} は \sum の上位 \tilde{d} 個の固有値 $\lambda_1, \dots, \lambda_{\tilde{d}}$ に対応する \tilde{d} 個の正規直交ベクトルを列とする行列として求まる。図2にPCAの計算手順の概要を示す。

3. 認識実験

3.1 実験条件

実験用データとして構音障害者、健常者それぞれ1

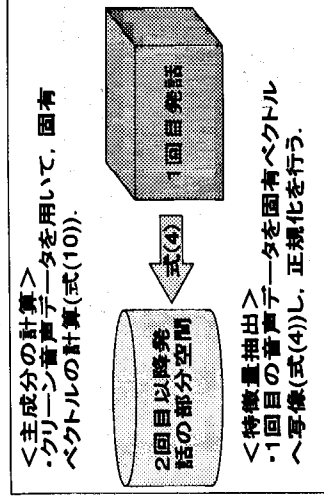


図2 PCAの計算手順

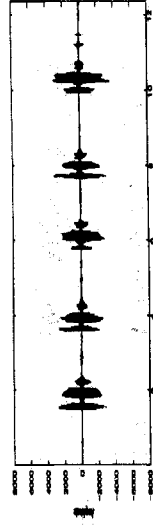


図3 収録データ

サンプリング周波数	16 kHz
ハミング窓長	25 msec
フレーム周期	10 msec
音響モデル	monophone (3 状態 43 音素)
特徴パラメータ	12 次 MFCC+Δ 12 次 MFCC
混合分布数	16
テストデータ	1050 (210 単語 × 5 回)
辞書	210 単語

名のデータを収録した。発話内容として ATR 音素パラ
ンス単語 (216 単語) から 210 単語を無作為に選択し
た。収録は各単語を 5 回連続発声し (図 3)、その後、
各発話を手動で切り出した。図 4 に構音障害者、図 5
に健常者の音声例を示す。

3.2 汎用モデルでの認識実験

初めに、汎用モデルでの認識実験を行った。認識実
験には「Julian 記述文法音声認識キット v3.1」[12] を
使用し、音響モデルは認識キットに含まれている不特
定話者音響モデルを用いた。追加の学習はせず、孤立
単語での認識を行った。実験条件を表 1 に示し、認識
結果を図 6 に示す。健常者では 88.4% の精度が得られ
るが、構音障害者では 2.4% しか認識できず、従来の汎
用モデルでは発話スタイルが健常者と異なるため認識
が困難である事がわかる。

3.3 構音障害者音響モデルでの認識実験

汎用モデルでの認識が困難であることから、構音障
害者の音響モデルを作成し認識実験を行った。1 回目
の発話の認識を行う場合は 2~5 回目の発話を用いて音

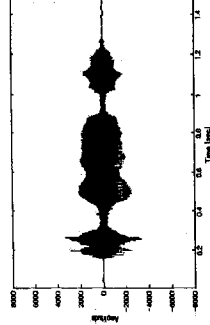


図4 構音障害者の音声例//akegata

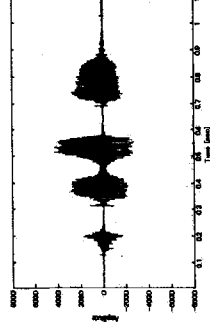


図5 健常者の音声例//akegata

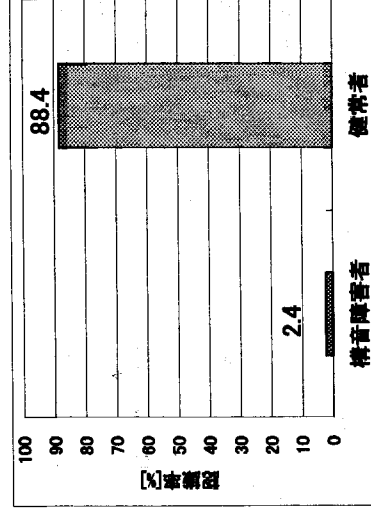


図6 汎用モデルでの認識実験結果

響モデルを作成した。これを各発話に対して行う。初
期モデルの作成、学習、認識には HTK [13] を用いた。
実験条件を表 2 に示し認識結果を図 7 に示す。

サンプリング周波数	16 kHz
ハミング窓長	25 msec
フレーム周期	10 msec
音響モデル	monophone (3 状態 54 音素)
特徴パラメータ	12 次 MFCC+Δ 12 次 MFCC (正規化)
混合分布数	6
テストデータ	1050 (210 単語 × 5 回)
辞書	210 単語

特定話者モデルを用いる事で構音障害者において、
87.2% の認識が得られた。健常者における発話回数ご
との認識率を図 8 に示し、構音障害者における発話回
数ごとの認識率を図 9 に示す。構音障害者において、1

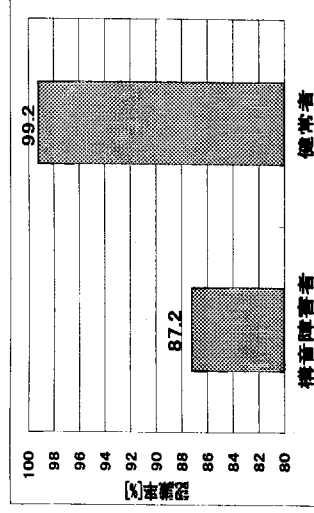


図7 特定話者モデルでの認識実験結果

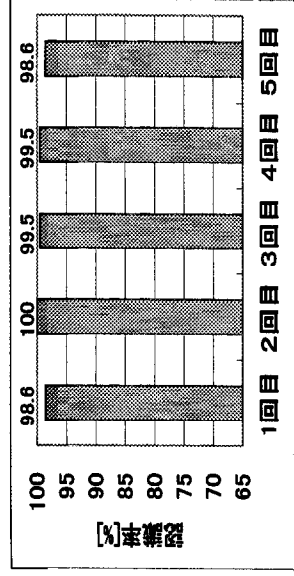


図8 健常者における発話回数ごとの認識率

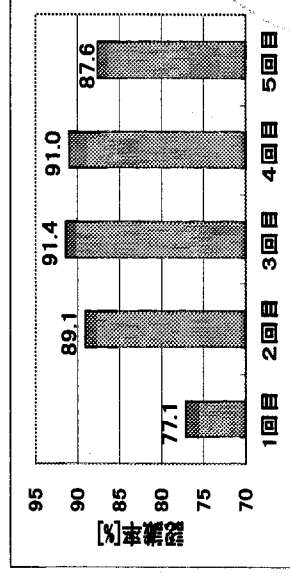


図9 構音障害者における発話回数ごとの認識率

回目発話の認識率が77.1%と他の発話に比べると著しく低下している。これは1回目の発話は最初の意図的な動作であり、他の発話よりも緊張状態に陥っていると考えられる。そのためアテトリーゼが生じて調音が困難になり、発話スタイルが不安定となることから認識精度が低下したと考えられる。

3.4 提案手法による認識実験

メルフイルタバンク出力24次元に対しPCAを適用した結果を示す。得られた値を基本係数とし、基本係数+ Δ 係数を音声認識の特徴量とした。今回は主成分の個数は11,13,15,17,19個として実験を行った。図10に1回目発話における結果を示し、図11に主成分が17個の場合の結果を示し、図12に結果を示す。

DCTの代わりにPCAを適用することにより、1回

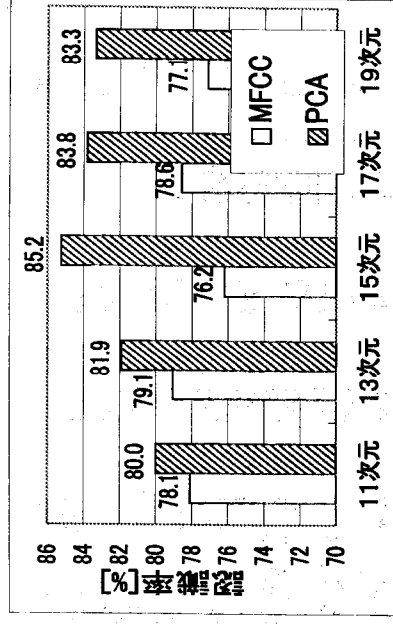


図10 提案手法による1回目発話の認識率

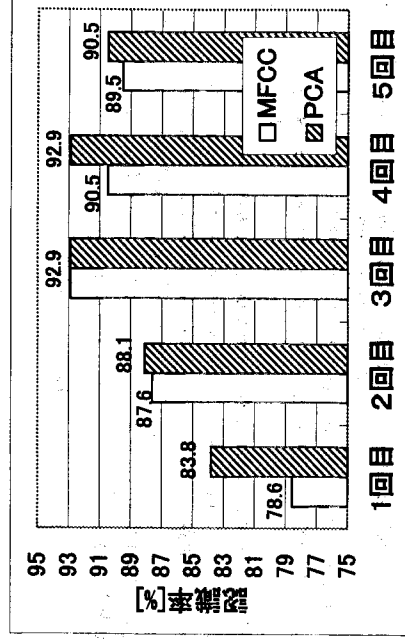


図11 提案手法による発話回数ごとの認識率 (17次元)

目発話において、主成分15個で85.2%まで認識率が改善された。DCTを用いる場合では最も良い場合でも79.1%(13次元)の認識精度なので、DCTよりもPCAの方が1回目発話の調音不安定音声において、より有効的な特徴量抽出法であるといえる。また、全体の認識精度も平均1.3%改善された。

4. おわりに

構音障害者の音声認識の検討を行った。構音障害者の発話スタイルが健常者と異なることから汎用モデルでは認識が困難であるが、構音障害者用の音響モデルを用いることで音声認識技術の適用が可能であることがわかった。1回目の調音不安定発話に対する特徴抽出方法としてMFCCにおけるDCTの代わりに、PCAを用いる手法を検討した。提案手法によって6.1%の改善が得られた。今後は、構音障害者特有の特徴量の検討や、様々な音声認識手法の適用を行い認識率の改善に取り組んでいく。また更に対象者を増やしていく予

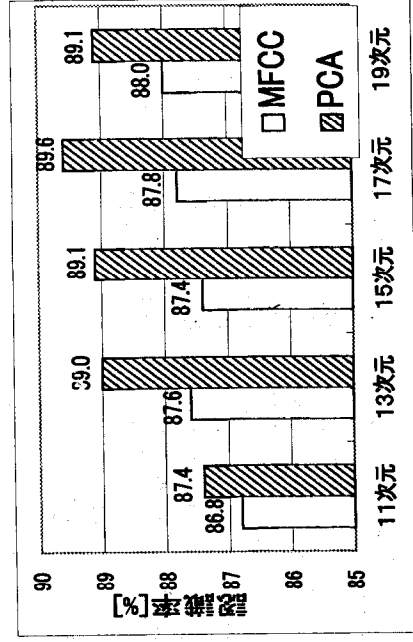


図 12 提案手法による認識実験結果

定である。

文 献

- [1] 佐川浩彦, 酒匂裕, 大平栄二, 崎山朝子, 阿部正博, “圧縮連続 DP 照合を用いた手話認識方式,” 電子情報通信学会論文誌, Vol. J77-D2, No.4, pp. 753-763, 1994.
- [2] 鈴木悠司, 平岩裕康, 竹内義則, 松本哲也, 工藤博章, 大西昇, “視覚障害者めための環境内の文字情報抽出システム,” 電子情報通信学会技術研究報告, WIT2003-314, pp. 13-18, 2003.
- [3] 藤謙一郎, 伊福部達, 青村茂, “発話障害者支援のための音声合成器の基礎的設計,” 電子情報通信学会技術研究報告, SP2006-321, pp. 59-64, 2006.
- [4] 鮫島充, 李晃伸, 猿渡洋, 鹿野清宏, “子供音声認識のための音響モデルの構築および適応手法の評価,” 電子情報通信学会技術研究報告, SP2004-114, pp. 109-114, 2004.
- [5] 中村圭吾, 田村直良, 鹿野清宏, “発話障害者音声を対象にした健常者音響モデルの適応と検証,” 日本音響学会講演論文集, 3-7-4, pp. 109-110, 2005.
- [6] 内閣府 “平成 18 年 障害者白書,” <http://www8.cao.go.jp/shougai/>
- [7] 厚生労働省 “平成 13 年 身体障害児・者実態調査結果,” <http://www.mhlw.go.jp/houdou/2002/08/h0808-2.html>
- [8] S. Terry Canale, 落合直之, 藤井克之, “キャンベル整形外科手術書 第 4 巻 小児の神経障害/小児の骨折・脱臼,” エルゼビア・ジャパン, 2004.
- [9] “脳性マヒについて,” <http://www.geocities.co.jp/SweetHome/6954/nouseimahi.html>
- [10] 滝口哲也, 有木康雄, “Kernel PCA を用いた残響下におけるロバスト特徴量抽出の検討,” 情報処理学会論文誌, Vol.47, No.6, pp. 1767-1773, 2006.
- [11] 石井健一郎, 上田修功, 前田英作, 村瀬洋, “わかりやすいバーターン認識,” オーム社, 1998.
- [12] “大語彙連続音声認識システム Julius,” <http://julius.sourceforge.jp/index.htm>
- [13] S. Young et. al., “The HTK Book,” Entropic Labs and Cambridge University, 1995-2002.