

# 重度難聴者音声認識のための Deep Canonical Correlation Analysis を用いた 音響特徴量抽出の検討\*

☆高島悠樹 (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

現在, 我が国の障害者手帳を持つ 18 歳以上の人口は 350 万人を超えており, 聴覚・言語障害者の数は 36 万人とされている [1]. 文献 [2] 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ている研究は非常に少ない. 本研究は, コミュニケーション手段として口話を用いる重度難聴者を対象として, 音声と唇形状によるマルチモーダル音声認識を実現し, ユビキタス社会における彼らの生活の支援をすることを目的としている.

人間は発話内容を理解する際, 種々の情報を統合的に利用している. 音声聞き取りが難しい場合, 発話者の顔, 特に唇の動きに注目して発話内容を理解しようとし, 逆に, 唇の動きと音声不一致の場合, 唇の動きに影響されて発話内容を誤って理解してしまうこともある. これは, McGurk effect (マガーク効果) と呼ばれ, 音韻知覚が音声の聴覚情報のみで決まるのではなく, 唇の動きといった視覚情報からも影響を受けることが報告されている [3]. このように人間による発話内容の理解には, 唇の画像と音声情報の統合的利用が極めて重要である.

唇の動きからの発話内容の読み取りは, リップリーディング (読唇) と呼ばれ, 聴覚障害者にとって重要なコミュニケーション手段の一つである. リップリーディングは, 背景雑音に影響されることがないため, 計算機上での実現が期待されている. 例えば, 監視カメラに収録された会話映像のように音声聞き取りにくい場合であっても, リップリーディングであれば発話内容の分析が可能であり, 犯罪の防止や抑止に繋がると考えられる. そのため, 音声の雑音に対して頑健な発話認識を行う手法の一つとして, 音声情報に唇動画像情報を併用して認識を行うマルチモーダル音声認識が注目され, 研究が進められている [4, 5].

重度難聴者は耳で音を聞くことができないため, 正確な発音をすることが難しく, 発話スタイルが健常者と異なる. 彼らのコミュニケーション手段の一つとして口話があり, 訓練により意図した発話の唇の形状を作ることが可能である. そこで, 彼らの音声を認識するために, 唇画像を併用した音声認識シ

テムの構築が望まれる. 重度難聴者を対象としたマルチモーダル音声認識として, CNN (convolutional neural network) を用いた手法 [6] が提案されている. 一般に, 音声をニューラルネットワークを用いてモデリングする際, 教師信号として, 入力特徴量に対応する音素ラベルが使用される. しかし, 重度難聴者音声の場合, 発話スタイルが健常者と異なるため, 強制アライメントにより得られた音素ラベルは誤りを含む. そこで, 文献 [6] では, ネットワークの中間層にボトルネック層を設け, この層のユニットを認識のための特徴量とすることで, 教師信号の誤りに頑健な特徴抽出を行なっている. 本稿では, このアライメントの誤りに対する異なるアプローチとして, 教師信号を使わない, 教師なし学習による検討を行なう.

マルチモーダル学習として, canonical correlation analysis (CCA) を非線形拡張した deep canonical correlation analysis (DCCA) [7] が提案されている. CCA は 2 変数間の相関を最大にするような線形射影行列を学習する手法であり, DCCA は, 2 つのニューラルネットワークにより非線形マッピングされた変数間の相関を最大にするように拡張した枠組みである. CCA と異なり, DCCA はパラメトリックな手法であり, より複雑な変換を学習することができる. DCCA は様々な分類タスク [8, 9] に応用されており, 精度の向上が報告されている. DCCA の目的関数は学習データに対する教師なし学習として設計されており, 上述の誤りを含む音素ラベルを用いる必要がない. 本研究では, DCCA を用いて音声と唇画像からマルチモーダル音声認識のための特徴量抽出法を提案する. マルチモーダルタスクにおいて, モダリティ間の関係性を考慮することは非常に重要なことだと考えられる. 音声認識タスクにおいて, 一般に唇画像は音声と比べて情報量が少ないと考えられる. しかし, 認識に有効な相補的な特徴を抽出することができれば, 音声特徴量の品質が劣化した際に, 唇画像を用いて認識率を補償することができると考えられる. DCCA により互いに高相関な特徴量を抽出することができ, このような効果を持つ特徴量の抽出が期待できる. 以下, 第 2 章で CCA と DCCA について述べ, 第 3 章で提案手法の流れを説明する. 第 4 章で従来の特徴量と比較し, 第 5 章で本稿をまとめる.

\* Audio-Visual Speech Recognition for a Person with Severe Hearing Loss Using Deep Canonical Correlation Analysis, by Yuki Takashima (Kobe University), Tetsuya Takiguchi (Kobe University/JST PRESTO), Yasuo Ariki (Kobe University)

## 2 Canonical Correlation Analysis と非線形拡張

### 2.1 Canonical Correlation Analysis

$X_{audio} \in \mathbb{R}^{d_1 \times N}$ ,  $X_{visual} \in \mathbb{R}^{d_2 \times N}$  をそれぞれ,  $d_1$ ,  $d_2$  次元の音声特徴量及び画像特徴量を  $N$  サンプル並べた行列とする. ここで, これらの行列は平均 0 に正規化されているものとする. CCA において, これらの変数間の相関係数は以下の式で計算される.

$$\rho(\mathbf{a}, \mathbf{b}) = \text{corr}(\mathbf{a}^\top X_{audio}, \mathbf{b}^\top X_{visual}) \quad (1)$$

$$= \frac{\mathbf{a}^\top \Sigma_{av} \mathbf{b}}{\sqrt{\mathbf{a}^\top \Sigma_{aa} \mathbf{a}} \sqrt{\mathbf{b}^\top \Sigma_{vv} \mathbf{b}}} \quad (2)$$

ここで,  $\mathbf{a} \in \mathbb{R}^{d_1}$ ,  $\mathbf{b} \in \mathbb{R}^{d_2}$  は射影ベクトルであり, CCA において推定されるパラメータである.  $\Sigma_{av} \in \mathbb{R}^{d_1 \times d_2}$ ,  $\Sigma_{aa} \in \mathbb{R}^{d_1 \times d_1}$ ,  $\Sigma_{vv} \in \mathbb{R}^{d_2 \times d_2}$  はそれぞれ,  $X_{audio}$  と  $X_{visual}$  の相互共分散行列,  $X_{audio}$  及び  $X_{visual}$  の自己共分散行列を表す.  $\rho(\mathbf{a}, \mathbf{b})$  は,  $\mathbf{a}$ ,  $\mathbf{b}$  のスケールに対して不変であるため, 各標準偏差を 1 と仮定すると, CCA で解くべき問題は以下の最大化問題となる.

$$\max_{\mathbf{a}, \mathbf{b}} \mathbf{a}^\top \Sigma_{av} \mathbf{b} \quad \text{subject to} \quad \mathbf{a}^\top \Sigma_{aa} \mathbf{a} = \mathbf{b}^\top \Sigma_{vv} \mathbf{b} = 1 \quad (3)$$

また,  $L \leq \min(d_1, d_2)$  個の線形射影ベクトルを使用するとき, 音声及び画像特徴量の射影行列はそれぞれ,  $\mathbf{A} \in \mathbb{R}^{d_1 \times L}$  及び  $\mathbf{B} \in \mathbb{R}^{d_2 \times L}$  と書け, 以下のように定式化される.

$$\begin{aligned} & \text{maximize} \quad \text{tr}(\mathbf{A}^\top \Sigma_{av} \mathbf{B}) \\ & \text{subject to} \quad \mathbf{A}^\top \Sigma_{aa} \mathbf{A} = \mathbf{B}^\top \Sigma_{vv} \mathbf{B} = \mathbf{I} \end{aligned} \quad (4)$$

ここで,  $\text{tr}(\cdot)$  と  $\mathbf{I}$  はそれぞれ, 対角成分の和と単位行列を表す.

CCA は,  $\mathbf{T} = \Sigma_{aa}^{-1/2} \Sigma_{av} \Sigma_{vv}^{-1/2}$  の特異値分解により計算される.  $k$  個の射影ベクトルを扱うとき, 射影行列は  $(\mathbf{A}, \mathbf{B}) = (\Sigma_{aa}^{-1/2} \mathbf{U}_k, \Sigma_{vv}^{-1/2} \mathbf{V}_k)$  で与えられる. ここで,  $\mathbf{U}_k \in \mathbb{R}^{d_1 \times k}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d_2 \times k}$  は,  $\mathbf{T}$  の最初から  $k$  個の左及び右特異ベクトルを並べたものである. 実際には, 共分散行列  $\Sigma_{aa}$ ,  $\Sigma_{vv}$  は, 正則行列となるよう正則化を加えて計算される.

### 2.2 Deep Canonical Correlation Analysis

DCCA は, CCA にニューラルネットワークを組み込んだものであり, 複数層による非線形変換を行なう. 音声及び画像特徴量 ( $X_{audio}, X_{visual}$ ) が与えられた時, 音声及び画像のニューラルネットワークの出力を  $f(X_{audio}; \theta_1) \in \mathbb{R}^{o \times N}$  と  $f(X_{visual}; \theta_2) \in \mathbb{R}^{o \times N}$  とする. ここで,  $\theta_1, \theta_2$  はそれぞれ, 音声及び画像の

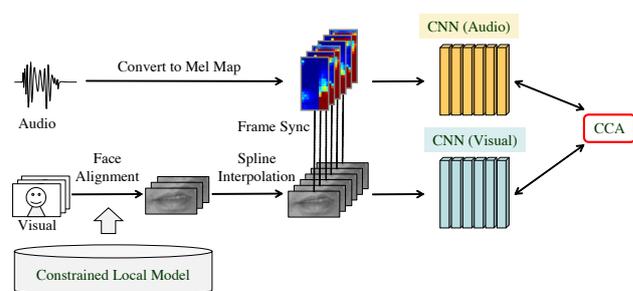


Fig. 1 Deep CCA using CNNs

ネットワークのパラメータを表す. DCCA における相関係数は以下の式で計算される.

$$\text{corr}(\mathbf{a}^\top f(X_{audio}; \theta_1), \mathbf{b}^\top f(X_{visual}; \theta_2)) = \text{tr}(\mathbf{T}^\top \mathbf{T})^{\frac{1}{2}} \quad (5)$$

ここで,  $\mathbf{T} = \hat{\Sigma}_{aa}^{-1/2} \hat{\Sigma}_{av} \hat{\Sigma}_{vv}^{-1/2}$  は 2.1 節と同様である.  $\hat{\Sigma}_{av} = \frac{1}{N-1} X_{audio} X_{visual}^\top$  and  $\hat{\Sigma}_{aa} = \frac{1}{N-1} X_{audio} X_{audio}^\top + r_1 \mathbf{I}$  は, 正則化制約  $r_1 > 0$  を用いて計算された共分散行列であり,  $\hat{\Sigma}_{vv}$  に対しても同様に計算される. DCCA は 2 つのデータ群に対する相関を最大にするように, パラメータ  $\{\theta_1, \theta_2, \mathbf{u}, \mathbf{v}\}$  を同時に学習を行なう. パラメータ  $\{\theta_1, \theta_2\}$  は式 (5) より, 誤差逆伝搬法により学習され, その勾配は以下の式で計算される.

$$\begin{aligned} & \frac{\partial \text{corr}(\mathbf{a}^\top f(X_{audio}; \theta_1), \mathbf{b}^\top f(X_{visual}; \theta_2))}{\partial f(X_{audio}; \theta_1)} \\ &= \frac{1}{N-1} (2 \nabla_{aa} X_{audio} + \nabla_{av} X_{visual}) \end{aligned} \quad (6)$$

ここで,  $\nabla_{ab} = \hat{\Sigma}_{aa}^{-1/2} \mathbf{U} \mathbf{V}^\top \hat{\Sigma}_{vv}^{-1/2}$ ,  $\nabla_{aa} = -\frac{1}{2} \hat{\Sigma}_{aa}^{-1/2} \mathbf{U} \mathbf{D} \mathbf{V}^\top \hat{\Sigma}_{aa}^{-1/2}$  であり,  $f(X_{visual}; \theta_2)$  に対する勾配も同様に計算される.

## 3 DCCA を用いたマルチモーダル特徴量抽出

### 3.1 提案手法の流れ

Fig. 1 に提案する特徴量抽出法の流れを示す. 従来研究 [6] より, 特徴量抽出において CNN が有効であることが示されているため, 本研究では, DCCA における非線形変換に CNN を利用する.

重度難聴者の発話した音声と唇画像からネットワークへの入力特徴量を用意する. 音声特徴量として, 音声信号から計算されたメル周波数スペクトルを数フレーム束ねた 2 次元のメルマップを使用する. 唇領域抽出のためのフェイスアライメントは, 顔モデルを PDM (point distribution model) で表現し, CLM (constrained local model) の枠組みで計算し実現する. 抽出された唇画像は, 音声特徴量のサンプリ

ング周波数に合わせるため、3次スプライン補間を適用する。

音声と唇画像のCNNは、DCCAの目的関数を偏微分し誤差逆伝搬法により学習される。各CNNの学習後、入力音声特徴量と唇画像特徴量をCNNに対して順伝搬させ、得られた出力ユニットを用いて以下の式より線形射影を行なう。

$$\alpha_t = \hat{\Sigma}_{aa}^{-1/2} U_k f(X_t; \theta_1) \quad (7)$$

$$\beta_t = \hat{\Sigma}_{vv}^{-1/2} V_k f(Y_t; \theta_2) \quad (8)$$

ここで、 $(X_t, Y_t)$  はそれぞれ、時刻  $t$  における音声と唇画像の2次元の入力特徴量を表し、 $(\alpha_t \in \mathbb{R}^k, \beta_t \in \mathbb{R}^k)$  は対応する射影後の特徴量を表す。これらの特徴量を連結し、 $[\alpha_t^T \beta_t^T]^T \in \mathbb{R}^{2k}$  をHMMへの入力特徴量として扱い、音声認識を行なう。

### 3.2 重度難聴者音声への応用

重度難聴者の発話スタイルは健常者と異なるため、強制アライメントにより得られた音素ラベルは誤りを含む。Deep neural network (DNN) を用いた手法は教師信号を必要とし、一般に音声を扱う場合、教師信号として入力に対応する音素ラベルが用いられる。しかし、音素ラベルに誤りを含む場合、ネットワークの学習が十分に行なえず、性能向上の妨げになると考えられる。DCCAは教師なし学習の枠組みの1つであり、2変数(モダリティ)間の相関を最大にするように学習するため、音素ラベルを使用しない。DCCAを用いることにより、音声及び唇画像はそれぞれ、互いに高相関となるように変換される。そのため、雑音環境下において、音声信号が劣化しても、唇画像は背景雑音に対して不変であるため、抽出された特徴量はノイズに対する頑健性を持つと期待される。

## 4 評価実験

### 4.1 実験条件

データセットとして、重度難聴者の男性1名の音声及び唇動画を収録し使用した。発話内容はATR音素バランス単語Aセット [10] から選択し、2,620単語を学習、216単語を評価に使用した。音声の標準化周波数は16kHz、語長16bitであり、音響分析にはHamming窓を用いた。STFTにおけるフレーム幅、シフト幅はそれぞれ25ms、5msである。本稿で用いる音響モデルは、54音素のmonophone-HMMで、各HMMの状態数は5、状態あたりの混合分布数は6である。

ケプストラム特徴量であるMFCC+ $\Delta$ + $\Delta\Delta$  (36次元) をベースラインとし、提案手法との比較を行なう。さらに、画像特徴量であるdiscrete cosine trans-

form (DCT) を加えたMFCC+ $\Delta$ + $\Delta\Delta$ +DCT (66次元) をマルチモーダル特徴量として比較する。また、雑音環境下での認識性能を比較するため、音声データに白色雑音 (SNR:20dB, 10dB, 5dB) を加えて評価を行なった。なお、ネットワークの学習にはクリーン音声のみを用いた。

### 4.2 ネットワーク構成

音声CNNの入力層には、39次元のメル周波数スペクトルをフレーム幅13、シフト幅1で分割したメルマップを用いる。画像CNNの入力層には、発話時に顔正面から60fpsで撮影された動画を、(1)画像列に変換し、(2)CLMにより唇領域の輝度画像を抽出、(3)12×24pixelにリサイズを行った上で、(4)スプライン補間によってアップサンプリング(メルマップとの同期)を行った唇画像を用いる。

Table 1に実験に用いたネットワーク構成を示す。ボトルネック特徴量の有効性を確認するため、ボトルネック層を設定している。学習率0.0001、モーメンタム0.99として、確率的勾配法を用いてモデルを学習した。

Table 1 Filter size, number of feature maps and number of MLPs units for each architecture. The value for C indicates the filter size of the convolution layer that has #1 maps. The convolution layer is associated with the pooling layer. The value of S means the pooling factor. The value for M indicates the number of units for each layer in the MLP part.

	Input	C	S	#1	M
Audio CNN	39×13	4×2	3×3	13	108, 30, 108
Visual CNN	12×24	5×5	2×2	13	108, 30, 108

### 4.3 実験結果と考察

DCCAでは、各変数に対する共分散行列を計算する必要がある。より正確な共分散行列の計算のため、ミニバッチサイズは大きく設定する必要がある [8]。そこで、クリーン環境下において、ミニバッチサイズによる影響を調査した。Fig. 2に、ミニバッチサイズを変化させた時の実験結果を示す。ミニバッチサイズを大きくすると認識率が向上することが分かる。以降の実験では、ミニバッチサイズを2,100に設定する。

Fig. 2に、従来手法との比較結果を示す。DCCAとDCCA bottleneckはそれぞれ、最終射影層とボトルネック層から抽出された特徴量を表す。DCCAはDCCA bottleneckと比べて精度が劣化する傾向が見られた。この理由として、音声特徴量が唇画像空

Table 2 Word recognition accuracy for each mini-batch size

# of mini-batches	1,200	1,500	1,800	2,100	2,400
Recognition accuracy [%]	63.89	65.28	66.20	71.76	71.76

間の近くへ射影されたことにより、音声特徴量が持つ情報が失われたためと考えられる。また、DCCA bottleneck は MFCC+DCT と比べて、SNR10dB において、良い精度が得られた。これは、DCCA を用いることで従来の特徴量よりもノイズにロバストな特徴量が得られたためと考えられる。

Fig. 3 に、教師なし学習である従来研究 [6] と比較した認識結果を示す。DCCA は教師なし学習であるため、抽出された特徴量が認識に有効な特徴を示すとは限らない。従来法は誤りを含む音素ラベルを学習に用いているが、明示的に音韻情報をモデルの学習に組み込んでいるため、提案手法と比べて高い認識率が得られたと考えられる。提案手法は教師あり学習である従来手法と比べて、平均 14% の精度の劣化が見られた。

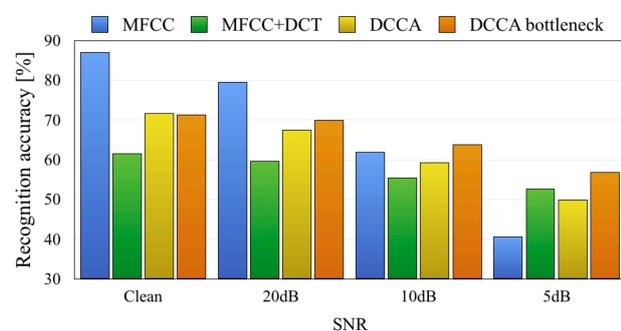


Fig. 2 Word recognition accuracy using HMMs

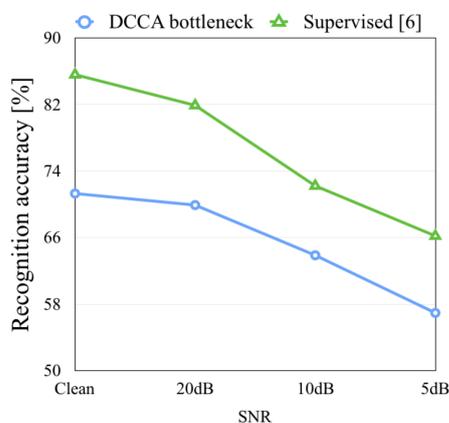


Fig. 3 Word recognition accuracy of unsupervised and supervised training procedure

## 5 おわりに

本稿では、DCCA を用いたマルチモーダル音声認識のための特徴量抽出法を提案した。重度難聴者の音声と唇画像を用いた雑音環境下单語認識実験により、提案手法の評価を行なった。評価実験において、提案手法は従来のケプストラム特徴量と比べて、高い認識率を示した。しかし、教師あり学習による従来研究と比べて、精度が改善されなかった。この原因として、DCCA では音韻情報を使用しないため、認識に有効な特徴量が抽出されなかったためと考えられる。今後は、音韻情報に対しても高い相関を持つような拡張を検討を行なう。

謝辞 本研究の一部は、JST さきがけ JP-MJPR15D2, JSPS 科研費 JP17J04380 の支援を受けたものである。

## 参考文献

- [1] 内閣省, “平成 25 年版障害者白書,” .
- [2] 中村圭吾 *et al.*, “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, pp. 109–110, 2015.
- [3] H. McGurk and J. MacDonald, “Hearing lips and seeing voices,” *Nature*, vol. 264, pp. 746–748, 1976.
- [4] G. Potamianos *et al.*, “Audio-visual automatic speech recognition: An overview,” *Issues in visual and audio-visual speech processing*, vol. 22, pp. 23, 2004.
- [5] Y. Mroueh *et al.*, “Deep multimodal learning for audio-visual speech recognition,” in *ICASSP*, 2015, pp. 2130–2134.
- [6] Y. Takashima *et al.*, “Audio-visual speech recognition using bimodal-trained bottleneck features for a person with severe hearing loss,” in *INTERSPEECH*, 2016, pp. 277–281.
- [7] G. Andrew *et al.*, “Deep canonical correlation analysis,” in *ICML*, 2013, pp. 1247–1255.
- [8] W. Wang *et al.*, “Unsupervised learning of acoustic features via deep canonical correlation analysis,” in *ICASSP*, 2015, pp. 4590–4594.
- [9] N. E.-D. El-Madany *et al.*, “Multiview learning via deep discriminative canonical correlation analysis,” in *ICASSP*, 2016, pp. 2409–2413.
- [10] A. Kurematsu *et al.*, “ATR Japanese speech database as a tool of speech recognition and synthesis,” *Speech Communication*, vol. 9, no. 4, pp. 357–363, 1990.