# 口唇動作と音声のタイミング構造に基づく話者検出

堀井 悠† 川嶋 宏彰† 松山 隆司†

† 京都大学大学院情報学研究科 〒 606-8501 京都市左京区吉田本町 E-mail: †horii@vision.kuee.kyoto-u.ac.jp, ††{kawashima,tm}@i.kyoto-u.ac.jp

**あらまし** 人の発話における口唇動作と音声変化の間には, ずれを伴う複雑な時間的構造が存在するため, これらの 共起性をフレーム単位のモデルで表現することはしばしば困難である. そこで本研究では, 話者を撮影して得た口唇 動作と音声の特徴量系列をそれぞれ時間的に分節化して時区間系列の対とし, これら時区間の時間関係に基づいてメ ディア間の系統的時間差(タイミング構造)をモデル化することで, 複数人物の発話シーンに対して話者検出(どの 人物が発話しているかの判別)を行う手法を提案する. 本手法を用いて, 近接した複数人物のうちいずれが発話して いるかを, 単一のカメラとマイクのみでも高精度に検出できることを実験により確認した. **キーワード** 話者検出, タイミング構造, 視聴覚統合, hybrid dynamical system

# Speaker Detection Using the Timing Structure between Lip Motion and Speech Signal

## Yu HORII<sup>†</sup>, Hiroaki KAWASHIMA<sup>†</sup>, and Takashi MATSUYAMA<sup>†</sup>

† Graduate School of Informatics, Kyoto University Yoshida-Honmachi, Sakyo, Kyoto, 606-8501 Japan E-mail: †horii@vision.kuee.kyoto-u.ac.jp, ††{kawashima,tm}@i.kyoto-u.ac.jp

**Abstract** In this paper, we propose a novel approach to speaker detection using the cue of *timing structure* between audio and visual information. We first extract a pair of feature sequences of lip motion and sound, and segment each sequence into temporal intervals. Then, we construct a cross-media timing-structure model of human speech by learning the temporal relations of overlapping intervals. Based on the learned model, we realize speaker detection by evaluating the timing structure between the observed video and audio. Our experimental result shows the effectiveness of using temporal relations of intervals for speaker detection.

Key words speaker detection, timing structure, auditory-visual integration, hybrid dynamical system

## 1. はじめに

我々人間は,話者の顔位置や口唇動作などの視覚情報と,音 源方向や音声変化などの聴覚情報を,その共起性に基づいて相 互補完的に統合することで,他者の発話状態(いつ誰が話した か)の認識を行っている.機械による人の発話状態認識におい ても,視聴覚メディアの統合による実現が期待されており,例 えば,会議や講義といったシーンに対して,その自動撮影のた めに話者位置を検出する方法が広く研究されている[1]~[3].

これらの手法の多くは、背景差分を用いた人物位置推定とマイ クロホンアレイによる音源定位の結果を、Coupled HMM(Hidden Marcov Model)や DBN(Dynamic Baysian Network)を 用いて統合することにより、話者検出を実現している [4], [5]. しかし、音源定位の空間分解能以上に人物が近接している場合 や、検出対象空間でのカメラおよびマイクロホン配置に制約が ある場合など、検出困難となる状況もしばしば存在する.

このような状況でも高精度に発話状態を認識するためには, 人や音源の位置情報だけではなく,人の発話における口唇動作 と音声変化の間の共起性を利用する方法が考えられる.発話認 識で用いられる従来のメディア統合手法は,サンプリングした データや特徴抽出時のフレームを単位として,同一フレームや 隣接フレームでの共起性や特徴量相関をモデル化している[6]. しかし,/a/や/o/などの母音の発声において口唇動作の開始が 音声よりも先行するように,これらは必ずしも完全に同期する ものではなく,従来のフレームベースの統合手法では,このよ うなずれを伴う共起性を十分に表現できない[7].実際に,人の 知覚においても,許容される視覚情報と聴覚情報の時間的ずれ には,ある程度の広がりがあることが知られている[8],[9].

本研究ではこの点に着目し、口唇動作と音声変化の間の系統 的時間差を伴う時間的構造を**タイミング構造**と呼び、これをよ



図1 全体の処理の流れ

り直接的にモデル化し,話者検出に利用する手法を提案する. 具体的には,映像中の複数人物のうち,いずれかの人物が発話 している区間について,発話者の正確な判別を目的とする.

メディア信号間のタイミング構造をモデル化することで, 音 声にあった口唇映像を生成する先行研究として文献 [10] がある. これに対して,本研究ではこのようなタイミング構造を.メ ディア信号の変換・生成ではなく,信号間の時間的構造の整合 性評価に用いる.すなわち,発話というイベントを認識する上 で,タイミング構造がどの程度の有効性を有するかを検討する.

本手法は、人物同士が非常に近接している場合や、カメラと マイクが各1台のみの状況でも実現可能であり、遠隔会議シス テムでの話者追跡による自動撮影だけでなく、アーカイブされ た映像コンテンツの分析などへの応用が期待できる.

### 2. 状況設定とアプローチ

本稿では,話者検出にタイミング構造を利用することの有効 性の評価に焦点をあてるため,以下のように単純化した状況を 設定する.

話者検出の対象とするデータは、カメラ方向を向いている複数人物の発話シーンを、カメラとマイク各1台ずつを用いて同時キャプチャした時系列画像および音声信号である。

- 画像には唇の遮蔽がない正面顔が映っている (video).
- 話者以外の人物の口唇動作も同時に起こる(video).
- 音声に関して一定レベルのノイズは許容する (audio).
- 同時発話は起こらない (audio).

この状況設定における話者検出のアプローチを以下に示す (図1参照).処理は学習フェーズと認識フェーズの2段階から なる.

まず,学習フェーズでは,人物の発話シーンから口唇動作と 音声の特徴量系列を抽出し,各メディア信号のモデル学習と特 徴量系列の分節化を行う.そして,得られた区間系列対を学習 データとして,口唇動作と音声の間に存在する時間的構造を学 習することによって,発話のタイミング構造モデルを獲得する. これは,口唇動作と音声強度の変化パターンに関して,両者の 時間的ずれの許容範囲を確率分布として表現したものである.

認識フェーズでは、新たな観測データに対して、シーンにい



図 2 時区間ハイブリッドダイナミカルシステム (IHDS)

る人物全員の口唇動作と音声の特徴量系列を抽出し,分節化を 行う.次に,あらかじめ学習フェーズで獲得したタイミング構 造モデルを用いて,各人物の口唇動作と音声の間のタイミング 構造を評価し,評価値の最も高い人物を話者として検出する.

なお、本研究では不特定人物を対象とした話者検出を目指す ため、個人差の大きな口唇映像や形状パラメタ、音素変化の特 徴量ではなく、口唇の特徴点の垂直方向の変化や音声パワーと いった、比較的個人差の少ない特徴量を利用する.

続く第3節では特徴量系列を分節化するための単一メディア 信号のモデル化法について述べ,第4節ではメディア間のタイ ミング構造をモデルとして表現する方法を述べる.第5節では 学習したモデルを用いた新たな観測データの評価方法について 述べる.第6節では実験により提案手法の有効性を検証する.

## 3. 単一メディア信号の時区間表現

本研究では、各フレームにおいて口唇動作および音声の特徴 量を抽出し、それぞれのベクトル系列を時系列データとして用 いる.これらの時系列データは、繰り返し出現する要素的な変 化の組み合わせによって表現することができると仮定する.本 稿では、時系列データを複数の線形システムで表現される時区 間に分割して扱うためのモデルとして、文献[10],[11] で提案 されている時区間ハイブリッドダイナミカルシステム(IHDS, Interval-based Hybrid Dynamical System)を用いる.

#### 3.1 時区間ハイブリッドダイナミカルシステム

IHDS は,信号の要素的な変化を表現する力学系(ここでは 線形システム)と,線形システム間の遷移を表す有限状態確率 オートマトンから構成される.これに類する構造を持つモデル として,SLDS(Switching Linear Dynamical System)[12]が 挙げられる.しかし,SLDSは線形システム内の状態変化と離 散的状態遷移を共に物理的時間でモデル化しており,一つの線 形システムが持続する時間が短いほど尤度が高くなるという点 で,分節化結果が観測ノイズの影響を受けやすいと考えられる.

一方,離散状態の持続時間を確率分布として直接モデル化 する枠組みとして Segment Model [7] があり、本稿で用いる IHDS はこの一形態であるといえる(Segment 内の信号変化の モデルに力学系を利用する). 一般に Segment Model などの ように多くのパラメタを有する場合、与えられた観測信号から モデルを推定することは困難であるが、IHDS では、要素的な 信号変化を表現する線形システムを観測信号からボトムアップ に推定する手法を提供しており,さらに,新たな観測信号の分 節化を精度よく行うことが可能である.

以上の理由から、本稿では観測信号を分節化する具体的モデルとして IHDS を用いるが、第4節以降で述べるメディア信号間のタイミング構造のモデル化には、上述の SLDS や Segment Model を含む、他のモデルおよび分節化手法も利用できる.

3.1.1 システムアーキテクチャ

IHDSは2層構造を持ち、第1層は複数の離散状態間の確率的 遷移をモデル化する有限状態確率オートマトンであり、第2層は 各離散状態に付随する複数の線形システム $D = \{D_1, ..., D_N\}$ を持つ内部状態空間(IHDS に属す全線形システムが共有する n次元ベクトル空間)である(図2参照). さらに、これら2 層の統合のために時区間の概念を導入し、各時区間に、属性と してオートマトンの離散状態  $q_i$  とその持続時間  $\tau$  を持たせる. このとき、離散状態  $q_i$  と線形システム  $D_i$  を対応づけることで、 オートマトンによる内部状態のダイナミクスを制御できる.以 下では、離散状態  $q_i$ ,持続時間  $\tau$  の時区間を  $\langle q_i, \tau \rangle$  と表す.

3.1.2 線形システム

線形システム Di の内部状態遷移を,

 $\boldsymbol{x}_{t} = F^{(i)} \boldsymbol{x}_{t-1} + \boldsymbol{g}^{(i)} + \boldsymbol{\omega}_{t}^{(i)}$ (1)

と表す. ただし,  $x_t$  は時刻 t における内部状態を表す n 次元ベ クトル,  $F^{(i)}$  は遷移行列,  $g^{(i)}$  はバイアスである. また,  $\omega_t^{(i)}$ はプロセスノイズであり, ガウス分布でモデル化される. 各線 形システムは  $F^{(i)}$ ,  $g^{(i)}$ ,  $\omega_t^{(i)}$ の分布をパラメタとして持つ.

3.1.3 区間に基づくオートマトンの離散状態遷移

IHDS でのオートマトンは,持続長を持つ離散状態系列(時 区間系列)を生成する.生成される時区間系列 $I = [I_1, ..., I_K]$ に単純マルコフ性を仮定し,さらに,区間同士にはギャップや オーバーラップがないものとする.このとき,時区間  $\langle q_i, \tau \rangle$ が時区間  $\langle q_i, \tau_p \rangle$  の後に続いて起こる確率(区間遷移確率)  $P(I_k = \langle q_j, \tau \rangle | I_{k-1} = \langle q_i, \tau_p \rangle)$ のモデル化を行う.

#### 3.2 IHDS の学習

学習データとして,特徴量ベクトル系列のみが与えられている とする.このとき,線形システムの固有値制約に基づく階層的ク ラスタリングによるシステム個数およびパラメタ概値の推定と, システム個数を固定しての EM (Expectation-Maximization) アルゴリズム[13]によるパラメタ調整という2段階の学習によ り, IHDS のシステム同定が行える.また,同時に学習データ の特徴量ベクトル系列に対する分節化処理も行える.このアル ゴリズムの詳細については,紙面の都合から文献[11]に委ねる.

## 3.3 単一メディア信号の分節化

観測系列が与えられると, IHDS は, 観測した信号系列を最 もよく表現できる時区間系列を, Viterbi アルゴリズムに基づ く尤度計算 [11] によって求める. これにより, 新たに観測され た口唇動作および音声に対して, 学習した IHDS を用いて特徴 量ベクトル系列の分節化を行い, 時区間系列に変換することが 可能である.



図 3 区間対の始点差・終点差によるタイミング構造の表現

## 4. メディア間のタイミング構造のモデル化

各メディア信号をそれぞれ別の IHDS でモデル化することで, それぞれの信号は区間系列として表現される.以下では,2つ のメディア信号 *S*,*S*'の時間的構造について扱うため,それぞ れのメディア信号を「'」の有無によって区別して表記する.

4.1 メディア信号におけるモードおよびモード対時間差分布

メディア信号に含まれる要素的な変化の事象をモードと呼び、ここでは各モード $M_i$ を IHDS に属す線形システム  $D_i$ と 一対一に対応させる.このとき、メディア信号 S におけるモー ド $M_i$  (i = 1, ..., N)の時区間と、メディア信号 S' における モード $M'_p$  (p = 1, ..., N')の時区間とが、どのような時間的 関係で開始・終了するかをモデル化したものを、メディア信号 S と S' のタイミング構造モデルと定義する.

メディア信号 S, S'の区間系列を,  $\mathcal{I} = [I_1, ..., I_K], \mathcal{I}' = [I'_1, ..., I'_{K'}]$ とすると, 区間対の総数は  $K \times K'$ となるが, こ れらの中には時間的に非常に離れた区間対も存在する.そこで, ここでは時間的に離れたところにある区間同士の相互依存性は 小さいと仮定し,オーバーラップする区間対についてのみモデ ル化を行う.すなわち,考慮する区間対  $(I_k, I'_{k'})$ の条件を,

$$[b_k, e_k] \cap [b'_{k'}, e'_{k'}] \neq \emptyset \tag{2}$$

とする. ただし,  $b_k$ ,  $e_k$  は区間  $I_k$  の開始時刻, 終了時刻を表す. 具体的には, 2 つのメディア信号 S, S' のオーバーラップす る区間対  $(I_k, I'_{k'})$  において, それぞれのモードが  $M_i, M'_p$  であ るときに, 両区間の始点差が  $b_k - b'_{k'}$ , 終点差が  $e_k - e'_{k'}$  にな る同時確率をモデルとする. すなわち, 式 (2) の条件を満たす 区間対  $(I_k, I'_{k'})$  に表れるモード対  $(M_i, M'_p)$  について, その開 始時刻の差  $d_b$  および終了時刻の差  $d_e$  の同時分布

$$P(b_{k} - b'_{k'} = d_{\mathbf{b}}, e_{k} - e'_{k'} = d_{\mathbf{e}} | m_{k} = M_{i},$$
  
$$m'_{k'} = M'_{p}, [b_{k}, e_{k}] \cap [b'_{k'}, e'_{k'}] \neq \emptyset)$$
(3)

をタイミング構造のモデルとし、これをモード対時間差分布と 呼ぶ.これは始点差と終点差を軸とする2次元ユークリッド空 間における分布となり(図3参照),例えば、分布が原点付近 に高い山を持つ場合は、その2つのモードは開始と終了が共に 同期する傾向があることを意味する.

#### 4.2 タイミング構造モデルの学習

学習データとして,2つのメディア S,S' のモード区間系列 *I*,*I*' が与えられたとする.このとき,モード対 (*M<sub>i</sub>*,*M'<sub>p</sub>*)の 時間差分布は,学習データの区間系列においてオーバーラップ を含む区間対のうち,モード対 (*M<sub>i</sub>*,*M'<sub>p</sub>*)を持つものを見つけ, 始点差と終点差を軸とする2次元ユークリッド空間に投票する ことによって学習する.ただし,学習サンプル数は有限なので, 実際はガウス分布等の分布関数を各投票点に畳み込むことで全 体の分布を得る.評価実験では,2次元ガウス関数を用いた.

これを全モード対について行うことで、モード対時間差分布 が学習できる。メディア信号 S, S' におけるモード集合の要素 数がそれぞれ N, N' の場合、合計  $N \times N'$  個の分布が得られ る.これによって、次式に示すような区間対  $(I_k, I'_{k'})$  について の関数  $F(I_k, I'_{k'})$  を得たことになる.

$$F(I_k, I'_{k'}) = P(b_k - b'_{k'} = d_{\mathbf{b}}, e_k - e'_{k'} = d_{\mathbf{e}} | m_k = M_i,$$
  
$$m'_{k'} = M'_p, [b_k, e_k] \cap [b'_{k'}, e'_{k'}] \neq \emptyset)$$
(4)

ただし、区間対  $(I_k, I'_{k'})$  が持つモード対を  $(M_i, M'_p)$ , 区間対 の始点差,終点差をそれぞれ  $d_b, d_e$  とした.

## 5. タイミング構造評価に基づく話者検出

2つのメディア信号 S, S' として, 区間系列  $\mathcal{I} = [I_1, \ldots, I_K]$ ,  $\mathcal{I}' = [I'_1, \ldots, I'_{K'}]$ が観測されたとする.このとき,区間系列対  $(\mathcal{I}, \mathcal{I}')$ が持つ時間的構造が,学習によって獲得したタイミング 構造モデルとどの程度一致しているかを評価する.

区間系列 I, I' のうち、オーバーラップを含む区間対  $(I_k, I'_{k'})$ の集合を P とする. このとき、この集合 P が持つタイミング 構造の評価値を、P に含まれる区間対がそれぞれの始点差・終 点差を持つ同時確率により定める. すなわち、区間対  $(I_k, I'_{k'})$ の評価値を式 (4) の  $F(I_k, I'_{k'})$  とし、P の評価値  $\hat{F}(P)$  を集合 P に含まれる全区間対について  $F(I_k, I'_{k'})$  を掛け合わせて、

$$\hat{F}(\mathcal{P}) = \left[\prod_{(I_k, I'_{k'}) \in \mathcal{P}} F(I_k, I'_{k'})\right]^{\frac{1}{n(\mathcal{P})}}$$
(5)

と定める. ただし,  $\mathcal{P}$  の要素数  $n(\mathcal{P})$  による正規化を行っている. さらに, メディア信号対 (S, S') 全体の評価値 E(S, S') を,

$$E(S,S') = \log \hat{F}(\mathcal{P}) = \frac{1}{n(\mathcal{P})} \sum_{(I_k,I'_{k'})\in\mathcal{P}} \log F(I_k,I'_{k'}) \quad (6)$$

と定義する.ここで対数をとるのは、実際の計算を行う際に、 計算機上でのアンダーフローを防ぐためである.

話者検出にあたって、2者X,Yの口唇動作を表すメディア 信号として $S^{(v,X)}$ , $S^{(v,Y)}$ ,音声信号として $S^{(a)}$ が観測された とする.このとき、 $E(S^{(v,X)},S^{(a)})$ と $E(S^{(v,Y)},S^{(a)})$ を比較す ることによって、その評価値の大きい方の口唇動作が、音声と の間に学習したタイミング構造モデルと近い時間的構造を持っ ている、すなわち、発話者のデータであると判断できる.

## 6. 評価実験

提案手法の評価のため、実験を行った.まず、学習データと



図 4 発話シーンをキャプチャした画像の例(学習データ)



図 5 撮影環境のレイアウト

して2者の発話シーンをキャプチャしたデータを用い,発話の タイミング構造モデルを学習した(6.1節).続いて,学習デー タと同一の2者をキャプチャした新たなデータに対して,学習 したモデルを用いたタイミング構造評価を行い,話者検出を 行った(6.2節).さらに,学習データとは別の5者の発話シー ンを用い,学習したモデルの汎化性能を検証した(6.3節).

#### 6.1 発話のタイミング構造モデルの学習

学習データとして2者の発話シーンをキャプチャし,以下に 述べる手順で発話のタイミング構造モデルを学習した.なお, ここで学習するモデルは人物間で共通のものとした.

6.1.1 時系列画像および音声データのキャプチャ

図4に示すような2者(人物X,Yとする)の発話シーンを キャプチャした.キャプチャは,図5に示すようにカメラ<sup>(注1)</sup> とマイク<sup>(注2)</sup>を各1台ずつ配置し,多少の残響とPCファン音 程度の音声ノイズのみが存在する室内環境で行った.撮影画像 は解像度640×480 pixels(口唇領域が約40×20 pixels),フ レームレート60fps,8ビット量子化,グレースケール画像と し,音声は標本化レート48kHz,16ビット量子化とした.

キャプチャしたデータは計 17010 フレーム(≃ 4.7 分)で,2 者は 1~2 センテンス(約 20~30 秒)毎に発話を交代した.

6.1.2 口唇動作と音声の特徴量抽出

本研究は,不特定人物に対する話者検出の実現を目的とする ため,不特定話者に対応したモデルを用いる必要があり,その ためには個人差の小さな特徴量を用いることが望ましい.

口唇動作に関しては,個人差が表れやすい口唇形状そのも のや,口形状の横方向(x方向)の変化は特徴量として適当で

<sup>(</sup>注1): Flea (Point Grey Research 社), 固定焦点 4mm CCTV レンズ (FUJINON, YF4A-2) を使用.

<sup>(</sup>注2):コンデンサマイク(SONY, ECM-23F5), USB Audio Capture (Roland, UA-1000) を使用.



図 6 (a) AAM の学習用画像の例. 特徴点は全 58 点(うち口唇領域 に 8 点), (b) キャプチャした顔画像系列の例, (c) AAM によ る特徴点抽出結果の例.

はない.そこで、口唇動作の特徴量として、下唇に対応する 特徴点の y 座標のフレーム間差分を用いた.具体的には、ま ず、撮影画像の一部に特徴点を手動で与え、顔の AAM (Active Appearance Model) [14] を学習し、それを用いて時系列画像 全体の顔特徴点座標(全 58 点、口唇輪郭に 8 点)を抽出した (図 6 参照).そして、特徴点座標系列に対して、全特徴点が 同一平面上に存在するという仮定をおき、両目端点と鼻下の計 5 点 (図 6(a), 14, 18, 22, 26, 53)を基準として特異値分解に よるスケールの正規化[15]を行ったのち、上唇(図 6(a), 40~ 44)の重心の位置合わせと回転補正を行った.このとき、下唇 の 5 点 (図 6(a), 40, 44~47)の y 座標のフレーム間差分を抽 出した 5 次元ベクトル系列を、口唇動作の特徴量系列とした.

一方, 音声の特徴量には音声信号のパワーを使用した. これ は,本研究では複数人の音声を単一のメディア信号として取り 扱うため,個人差に影響されないパワーベースの特徴量が望ま しいと考えられるためである.音声に関しては,幅 33.3msの 窓を 16.6ms 間隔で掛け,パワーを抽出した.処理には HTK Ver.3.4 [16] を用い,1次元特徴量系列を得た.

6.1.3 特徴量ベクトル系列の分節化

口唇動作および音声の特徴量ベクトル系列のうち,いずれか の人物が発話状態にある区間(17010 フレーム中 13533 フレー ム<sup>(計3)</sup>)を手動で切り出し,各時点での発話者の口唇動作と音 声信号を用いて,各メディア信号を表現する IHDS の学習およ び学習データの分節化を行った(3.2節).なお,各 IHDS の モード数は,モデル化誤差カーブを基準にして,口唇動作・音 声ともに5とした.分節化結果の例を図7,図8に示す.

6.1.4 発話のタイミング構造モデルの学習

口唇動作と音声変化を表す時区間系列をそれぞれメディア信 号 *S*, *S*' とおき,4.2 節に述べた方法で,モード対時間差分布を 学習した.すなわち,オーバーラップを含む区間対の始点差と 終点差を 2 次元平面にプロットし,各投票点に標準偏差 3,共分 散 0 の 2 次元ガウス分布を畳み込む<sup>(注4)</sup>ことによって,全モー



図 9 モード対時間差分布表. 横軸が始点差,縦軸が終点差で,それぞ れ ±50 フレームの範囲の分布を濃淡により可視化したもの.

ド対についてモード対時間差分布を得た.結果を図9に示す.

例えば、video mode 4 と audio mode 0 の時間差分布をみる と、始点差は分散が大きいが、終点は同期する傾向にあること がわかる.video mode 4 は口を開く動作に、audio mode 0 は 無音区間にそれぞれ対応しているので、この学習結果は、人の 発話に「口を開き終わってから、発声が起こる」という傾向が あることに対応していると解釈できる.

## 6.2 タイミング構造評価に基づく話者検出

学習したタイミング構造モデルを用いて,新たな観測データ に対する話者検出を行った.テスト用データとして,学習デー タと同一の2者X,Yの発話シーンを新たにキャプチャし,3.3 節に述べた方法で分節化を行った時区間系列データ12本(そ

場合は約 100ms, 遅れる場合は約 250ms までである. 今回は 2 $\sigma \simeq 100ms$  となるよう,標準偏差を 3 とした (フレームレート 60fps なので, $\sigma \simeq 50ms$ ).

 <sup>(</sup>注3): 人物 X の発話が 6761 フレーム、人物 Y の発話が 6772 フレーム、モデルが一方の人物に特化することを防ぐため、2 者のデータをほぼ同量とした.
 (注4): テレビ視聴の際の映像と音声について、どの程度のずれであれば不自然 さを感じないかを調べた実験 [8] によると、その範囲は映像に対して音声が進む

れぞれ約 2000 フレーム,計 23831 フレーム)を用いた.なお, キャプチャの環境は学習データと同じとした.

以下では、口唇動作、音声の時区間系列をそれぞれ video 信号, audio 信号と呼ぶ. さらに、音声に対する口唇動作の状態 として、次の3つの用語を定義する.

- **発話** (Utterance): 音声に対応した口唇動作をしている.
- 沈黙 (Silence): 口唇動作をしていない.

 ロパク(Fake lip motion):音声と無関係な口唇動作を している.表情変化や,呟きに伴う口唇動作などが含まれる.

#### 6.2.1 観測信号の評価方法

video 信号  $S^{(v)}$  と audio 信号  $S^{(a)}$  のタイミング構造評価関数 を  $E_i(S^{(v)}, S^{(a)})$  とし、 $E_i(S^{(v,X)}, S^{(a)})$  と  $E_i(S^{(v,Y)}, S^{(a)})$  の うち、値が大きい方の video 信号が、発話状態に対応するデー タであると判断する、以下では、 $S^{(v)}$  と  $S^{(a)}$  に幅 T フレーム の窓をかけ、その時間範囲の評価値を求める関数として、 $E_1 \sim E_3$  ( $E_3$  が提案手法) の 3 種類を定義する (図 10 参照).

a) 同一フレームでのモード対共起確率による評価

video 信号, audio 信号の同一フレームにおけるモード対  $(m_t^{(v)}, m_t^{(a)})$ が,  $(v_c, a_c)$ となる確率  $P(m_t^{(v)} = v_c, m_t^{(a)} = a_c)$ (t によらず一定) を, あらかじめ学習データから求めておく. $S^{(v,X)}, S^{(a)}$ の時刻 t における実際のモード対を  $(v_t, a_t)$  とする とき,評価関数  $E_1$  を次式により定義する.

$$E_1(S^{(v)}, S^{(a)}) = \frac{1}{T} \sum_{t=0}^{T-1} \log P\left(m_t^{(v)} = v_t, m_t^{(a)} = a_t\right)$$
(7)

b) 隣接フレーム間のモード遷移確率による評価

Coupled HMM で用いられているような隣接フレーム間の モード遷移確率を用いて評価を行う. a) において,前フレー ムのモード対  $(m_{t-1}^{(v)}, m_{t-1}^{(a)})$  が,  $(v_{p}, a_{p})$  であるときのモード 遷移確率  $P(m_{t}^{(v)} = v_{c}, m_{t}^{(a)} = a_{c} | m_{t-1}^{(v)} = v_{p}, m_{t-1}^{(a)} = a_{p})$  (t に よらず一定)を学習データから得ておく<sup>(注5)</sup>. このとき,評価 関数  $E_{2}$ を次式により定義する.

$$E_2(S^{(v)}, S^{(a)}) = \frac{1}{T-1} \sum_{t=1}^{T-1} \log P\left(m_t^{(v)} = v_t, m_t^{(a)} = a_t \middle| m_{t-1}^{(v)} = v_{t-1}, m_{t-1}^{(a)} = a_{t-1}\right)$$
(8)

#### c) モード対時間差分布による評価

本研究で提案するモード対時間差分布を用いた評価の方法と して,式(6)において,次の式(9)に定義する集合 *P*を用いた ものを評価関数 *E*<sub>3</sub>とする.すなわち,

$$\mathcal{P} = \left\{ \left( I_{k_{v}}^{(v)}, I_{k_{a}}^{(a)} \right) \middle| I_{k_{v}}^{(v)} \in \mathcal{I}^{(v)}, I_{k_{a}}^{(a)} \in \mathcal{I}^{(a)}, \\ \left[ b_{k_{v}}^{(v)}, e_{k_{v}}^{(v)} \right] \cap \left[ b_{k_{a}}^{(a)}, e_{k_{a}}^{(a)} \right] \neq \emptyset \right\},$$
(9)

$$E_{3}(S^{(v)}, S^{(a)}) = \frac{1}{n(\mathcal{P})} \sum_{\substack{(I_{k_{v}}^{(v)}, I_{k_{a}}^{(a)}) \in \mathcal{P}}} \log F(I_{k_{v}}^{(v)}, I_{k_{a}}^{(a)})$$
(10)



図 10 各評価関数で用いられる時間的構造. (a) 同一フレームでのモー ド対共起確率による評価, (b) 隣接フレーム間のモード遷移確 率による評価, (c) モード対時間差分布による評価(提案手法).



図 11 テスト用データに対する評価値の時間的変化の例.上段は各人物の口唇動作状態の変化を示したものである.窓幅 180 フレーム(3 秒),間隔 30 フレームで窓をかけ,その窓範囲内のオーバーラップする区間系列対について,評価関数 E<sub>3</sub> によりタイミング構造評価を行い,各窓の中央の時刻における評価値とした.上段は口パク状態が出現しないデータ(R = 100.0%).下段は口パク状態を十分に含むデータ(R = 94.1%).

とする. ただし,  $F(I_{k_v}^{(v)}, I_{k_a}^{(a)})$ の定義は式 (4) の通りである. 6.2.2 実験結果および考察

評価値の時間変化をみるため、テスト用データに対して幅 T = 180 フレーム(3 秒)の窓を30 フレーム間隔で掛け、各窓 範囲のデータについて、各人物の評価値を求めた.ただし、各 窓範囲の評価値を、その窓の中央の時刻における評価値と定め た.提案手法( $E_3$ )による評価結果の例を図11に示す.

図 11 のグラフから,発話状態の人物の評価値が他方の人物 よりも大きくなっている様子が確認できる.非発話者が沈黙状 態にあるとき(図 11 上段),その評価値の差は特に大きい.

続いて,検出精度の定量的評価のため,いずれかの人物が発 話している区間のうち,その発話者の評価値が他者の評価値よ

<sup>(</sup>注5):学習時に出現しないモード対の出現確率を0としないよう,補正した値 を用いる必要がある(ゼロ頻度問題のディスカウンティング).ここでは、全モー ド対の出現回数に一律に0.5を加えるという単純な加算法を用いた.なお、a) では、出現回数0のモード対がなかったため補正は行っていない.

表1 評価方法による検出正解率の比較.単位は%.括弧内はデータ点の数.  $E_3$ が提案手法である(6.2.1節を参照).  $R_{\rm sil} = 非発話者が沈黙状態のときの検出正解率, <math>R_{\rm fak} = 非発話者がロパク状態のときの検出正解率, <math>R_{\rm X} = 人物 X$ が発話者のときの検出正解率,  $R_{\rm Y} = 人物 Y$ が発話者のときの検出正解率, R = 2体での検出正解率を表す.

使用した評価関数	$R_{\rm sil}$ [%]		$R_{\mathrm{fak}}$ [%]		$R_{\rm X}$ [%]		$R_{\rm Y}$ [%]		R [%]	
$E_1$	55.0	(170)	51.4	(179)	86.1	(309)	13.4	(40)	53.1	(349)
$E_2$	36.2	(112)	49.1	(171)	77.7	(279)	1.3	(4)	43.1	(283)
E <sub>3</sub> (提案手法)	94.2	(291)	81.6	(284)	82.2	(295)	94.0	(280)	87.5	(575)

りも高くなっている区間の割合を求めた.以下では、この割合 を検出正解率と呼び、Rと表記する.また、特に非発話者が 沈黙、ロパク状態にある区間での検出正解率を順に R<sub>sil</sub>, R<sub>fak</sub> とし、発話者が人物 X, Y である区間での検出正解率を順に R<sub>X</sub>, R<sub>Y</sub> とする.テスト用データに対して、各評価関数を用い て R<sub>sil</sub>, R<sub>fak</sub> を求めた結果を図 12 に、R<sub>sil</sub>, R<sub>fak</sub>, R<sub>X</sub>, R<sub>Y</sub>, R を 求めた結果を表 1 にそれぞれ示す.

図 12 から,提案手法では非発話者が沈黙,ロパクのいずれの状態であっても,他の比較手法よりも高い精度で話者検出が行えていることがわかる.非発話者が沈黙状態の場合とロパク状態の場合を比較すると, *R*<sub>sil</sub>の方が*R*<sub>fak</sub>より 10 ポイント以上大きい.加えて,非発話者がロパク状態での正解検出時の評価値の差が平均 2.2 であったのに対し,沈黙状態の場合には平均 34.4 となり,顕著な差が見られた.これは,沈黙状態ではvideo 信号に長い時区間が多くなり,時間差分布の中心から離れるような時間的構造が現れやすいためだと考えられる.

さらに,表1で人物ごとの検出正解率を比較すると,評価関 数 $E_1, E_2$ による評価では $R_Y$  が $R_X$ に対して著しく低く,そ の結果として全体についてのRの値も 50% 程度となっており, 話者検出が正しく行えていないことがわかる.人物ごとの検出 正解率に大きな差が生じた原因として,これらの評価方法が モード対の出現頻度そのものをモデルとしていることが挙げら れる.個人差の小さい特徴量を選ぶことで,両人物の口唇動作 を共通の IHDS で記述できるが,個人ごとの口唇動作の癖の違 いのため,時区間系列のモード遷移やその出現頻度には,ある 程度の個人差が現れる.今回の実験では,(キャプチャ映像から の主観的な判断ではあるが)人物 X よりも人物 Y のほうが口 をはっきり動かして発話する傾向にあった.その結果として人 物 Y の方がモード遷移が激しく起こったため,人物 X の発話 に現れやすいモードを含むモード対の出現頻度が相対的に高く なり,検出正解率に差が生じたと考えられる.

これに対して,提案手法(*E*<sub>3</sub>)では,いずれの人物に対して も 80%以上が正しく検出できている.これは,モード対の出現 頻度を考慮せず,時間的構造の妥当性あるいは崩れに基づいて 評価値を算出することによって,このような口唇動作の個人差 を吸収できたためだと考えられる.

#### 6.3 汎化性能に関する実験

学習したタイミング構造モデルが不特定人物に対して適用で きるかを評価するため、モデル学習に用いたデータとは別の5 者(人物1,...,人物5)の発話シーンをキャプチャした(図13 参照).撮影環境は学習用データと同じく図5の通りとした.



図 12 タイミング構造評価方法の比較.非発話者が沈黙状態であると きの検出正解率 R<sub>sil</sub>(黒),および非発話者がロパク状態であ るときの検出正解率 R<sub>fak</sub>(白)を各評価方法ごとに示す.



図 13 学習データとは異なる人物をキャプチャした画像の例

撮影したデータは8本,計12837フレームで,その大部分で少 なくとも1名が無視できない程度の口パク動作を行った.

これらのデータに対して、6.1節で得たタイミング構造モデル (人物 X, Y から学習したもの)を用いて、タイミング構造の評価を行い、各時点で最大の評価値をもつ人物を話者として検出した.このときの各評価関数による検出正解率の結果を表 2 に示す.ただし、 $R_i$  (i = 1, ..., 5) は人物 i が発話者のときの検出正解率、R は全体の検出正解率を表す.

表2より,提案手法(*E*<sub>3</sub>)が,他の手法(*E*<sub>1</sub>,*E*<sub>2</sub>)よりも高 精度に話者検出を実現していることがわかる.評価関数*E*<sub>1</sub>,*E*<sub>2</sub> を用いた場合は,6.2節と同様に人物ごとの検出正解率に大き なばらつきがあり,結果として全体の検出正解率も低くなった.

一方,提案手法では人物2と人物5の検出正解率がともに 80%程度となった.これは6.2節の実験の*R*<sub>fak</sub>(ロパク状態 に対する検出正解率)と同程度の精度であり,他の人物のデー

表 2 学習データとは異なる人物に対する検出正解率.単位は%.括弧内はデータ点の数. *E*<sub>3</sub> が 提案手法である(6.2.1節を参照). *R<sub>i</sub>* = 人物 *i* が発話者のときの検出正解率を表す.

使用した評価関数	$R_1$ [%]		$R_2$ [%]		$R_3$ [%]		$R_4$ [%]		$R_5 ~[\%]$		R [%]	
$E_1$	61.8	(21)	6.3	(4)	53.1	(34)	6.0	(6)	3.5	(2)	21.0	(67)
$E_2$	47.1	(16)	0.0	(0)	32.8	(21)	0.0	(0)	3.5	(2)	12.2	(39)
E3 (提案手法)	51.5	(17)	78.1	(50)	67.2	(43)	64.0	(64)	84.2	(48)	69.8	(222)

タで学習したモデルが適用できたと判断できる.しかし,他の 人物 1,3,4の検出正解率は 50~60%前後と低い結果となった. 発話者ごとの平均話速を比較したところ,人物 2 および 5 の平 均話速は他の人物よりも速く,学習データ(人物 X,Y)に近 い値であった.このことは,発話者の話し方によって,タイミ ング構造にも違いが現れる可能性を示している.ただし,個人 間の検出正解率の分散は,他の手法と比較して小さくなってお り,話者の個性の影響をある程度吸収できているといえる.

### 7. おわりに

本研究では,発話時の口唇動作と音声の変化パターンの間に 存在する共起性や系統的時間差(タイミング構造)をモデル化 し,このモデルに基づいて観測信号の評価を行うことで,話者 検出を行う手法を提案した.本稿に示した実験は予備的な規模 のものではあるが,口唇動作と音声変化の間で許容されるタイ ミング構造を,モード対時間差分布によって直接モデル化する ことによって,フレーム単位で共起性をモデル化する場合に比 べて,高い精度で話者検出が可能となることを確認した.

一方で、今回提案したタイミング構造モデルは、オーバー ラップする区間対のみから時間差分布を学習するという、比較 的単純なものにとどまっている.非オーバーラップ区間対への 拡張として、オーバーラップはしないが近傍にある区間対も含 めてモデル化する方法が考えられる.そこで実際に、15フレー ム以内の時間差範囲にある区間対も含めて<sup>(注60)</sup>時間差分布を学 習し、6.2節と同様の実験を行ったが、検出精度に明確な差は 見られなかった.今回の評価方法では、考慮する時間範囲を広 げると評価に使う区間対の数も急激に増加するため、特徴的な 情報を持つモード対が逆に埋もれてしまった可能性がある.こ のため、精度の向上のためには、単に時間範囲を広げるだけで はなく、評価に用いるラベル対の選別やタイミング構造モデル 自体の拡張をあわせて検討していく必要がある.

さらに、タイミング構造の有効性評価に焦点をあてるため、 本稿では比較的単純な状況を設定したが、実際の応用を行う上 では、全員が沈黙状態にある場合の考慮や、一時的な唇の遮蔽・ 同時発話への対応、多人数のデータに基づく不特定話者モデル の学習などが必要となる.また、6.3に述べたように、発話者 の話し方がタイミング構造に与える影響については興味深く、 両者の関係を深く調べる必要がある.より大規模な定量評価と あわせ、これらを今後の課題とする.

**謝辞** 本研究の一部は,科学研究費補助金 No.18049046 の 補助を受けて行った.

#### 文 献

- 大西正輝, 影林岳彦, 福永邦雄: "視聴覚情報の統合による会議 映像の自動撮影", 電子情報通信学会論文誌 D-II, 85, 3, pp. 537-542 (2002).
- [2] D. Gatica-Perez, G. Lathoud, I. McCowan, J. Odobez and D. Moore: "Audio-visual speaker tracking with importance particle filters", IEEE International Conference on Image Processing (ICIP) (2003).
- [3] 西口敏司, 東和秀, 亀田能成, 角所考, 美濃導彦: "講義自動撮影 における話者位置推定のための視聴覚情報の統合", 電気学会論 文誌 C, **124**, 3, pp. 729–739 (2004).
- [4] V. Pavlović, A. Garg, J. Rehg and T. Huang: "Multimodal speaker detection using error feedback dynamic Bayesian networks", Proc. Computer Vision and Pattern Recognition, pp. 34–43 (2000).
- [5] 吉村隆, 浅野太, 本村陽一, 麻生英樹, 市村直幸, 山本潔, 中村哲:
  "実環境における発話区間検出のための音響情報と画像情報の統合", 電子情報通信学会技術研究報告. SP, 103, 26, pp. 13–18 (2003).
- [6] A. V. Nefian, L. Liang, X. Pi, X. Liu and K. Murphy: "Dynamic Bayesian networks for audio-visual speech recognition", EURASIP Journal on Applied Signal Processing, 2002, 11, pp. 1–15 (2002).
- [7] M. Ostendorf, V. Digalakis and O. A. Kimball: "From HMMs to segment models: a unified view of stochastic modeling for speech recognition", IEEE Transactions on Acoustics, Speech and Signal Processing, 4, 5, pp. 360–378 (1996).
- [8] 鎧沢勇,滝川啓,大久保栄,渡辺義郎:"衛星通信を利用した画像 会議におけるエコー及び伝搬遅延の影響",電子通信学会論文誌, J64-B, 11, pp. 1281–1288 (1981).
- [9] A. Peregudov, K. Glasman and A. Logunov: "Relative timing of sound and vision: evaluation and correction", Proceedings of the Ninth International Symposium on Consumer Electronics, pp. 198–202 (2005).
- [10] 川嶋宏彰,松山隆司:"時区間ハイブリッドダイナミカルシステムを用いたマルチメディア・タイミング構造のモデル化",情報処理学会論文誌,48,12, pp. 3680–3691 (2007).
- [11] H. Kawashima and T. Matsuyama: "Multiphase learning for an interval-based hybrid dynamical system", IEICE transactions on fundamentals of electronics, communications and computer sciences, 88, 11, pp. 3022–3035 (2005).
- [12] V. Pavlović, J. M. Rehg and J. MacCormick: "Learning switching linear models of human motion", Proc. Neural Information Processing Systems (2000).
- [13] A. P. Dempster, N. M. Laird and D. B. Rubin: "Maximum likelihood from incomplete data via the EM algorithm", J. R. Statist. Soc. B, **39**, pp. 1–38 (1977).
- [14] T. F. Cootes, G. J. Edwards and C. J. Taylor: "Active appearance model", Proc. European Conference on Computer Vision, pp. 484–498 (1998).
- [15] K. S. Arun, T. S. Huang and S. D. Blostein: "Least-squares fitting of two 3-d point sets", IEEE Trans. on Pattern Analysis and Machine Intelligence, 9, 5, pp. 698–700 (1987).
- [16] S. Young, G. Evermann, M. Gales, T. Hain, D. Kershaw, X. Liu, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev and P. Woodland: "The HTK Book (for HTK Version 3.4)", Cambridge University Engineering Department (2006).

<sup>(</sup>注6):式 (2) の区間対  $I_k, I'_{k'}$  の条件を,  $[b_k - \alpha, e_k + \alpha] \cap [b'_{k'}, e'_{k'}] \neq \emptyset$ とした.  $\alpha$  は考慮する時間範囲 ( $\alpha = 15$  フレーム) である.