

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4864783号  
(P4864783)

(45) 発行日 平成24年2月1日(2012.2.1)

(24) 登録日 平成23年11月18日(2011.11.18)

(51) Int.Cl. F I  
**G 1 O L 15/20 (2006.01)** G 1 O L 15/20 3 6 O Z  
**G 1 O L 15/02 (2006.01)** G 1 O L 15/02 2 0 O Z

請求項の数 7 (全 20 頁)

(21) 出願番号	特願2007-76928 (P2007-76928)	(73) 特許権者	000208891
(22) 出願日	平成19年3月23日 (2007.3.23)		K D D I 株式会社
(65) 公開番号	特開2008-233782 (P2008-233782A)		東京都新宿区西新宿二丁目3番2号
(43) 公開日	平成20年10月2日 (2008.10.2)	(74) 代理人	100106909
審査請求日	平成21年7月10日 (2009.7.10)		弁理士 棚井 澄雄
		(74) 代理人	100064908
			弁理士 志賀 正武
		(74) 代理人	100089037
			弁理士 渡邊 隆
		(72) 発明者	遠藤 俊樹
			埼玉県ふじみ野市大原2丁目1番15号
			株式会社K D D I 研究所内
		(72) 発明者	加藤 恒夫
			埼玉県ふじみ野市大原2丁目1番15号
			株式会社K D D I 研究所内

最終頁に続く

(54) 【発明の名称】 パタンマッチング装置、パタンマッチングプログラム、およびパタンマッチング方法

(57) 【特許請求の範囲】

【請求項1】

外部より入力された音声データまたは画像データの特徴量を算出する分析手段と、  
 前記分析手段で算出された前記特徴量を正規化する正規化手段と、  
 前記正規化手段で正規化された正規化済み特徴量に基づいて、パタンマッチングを行う  
 パタンマッチング手段と、

を備えたパタンマッチング装置において、

前記正規化手段は、

前記音声データの全フレーム数または前記画像データ全体の前記特徴量の平均値である  
 全体平均値を取得する全体平均取得手段と、

前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の平  
 均値である局所平均値を計算する局所平均計算手段と、

前記局所平均値に基づいて、前記音声データの局所のフレーム数または前記画像デー  
 アの局所範囲の前記特徴量の分散値である局所分散値を計算する局所分散計算手段と、

前記全体平均値と複数の前記局所分散値とに基づいて前記特徴量を正規化する正規化処  
 理計算手段と、

を備えたことを特徴とするパタンマッチング装置。

【請求項2】

前記全体平均取得手段は、前記音声データの全フレーム数または前記画像データ全体の  
 前記特徴量から前記全体平均値を計算する

ことを特徴とする、請求項 1 に記載のパターンマッチング装置。

【請求項 3】

前記全体平均取得手段は、予め記憶した所定値を前記全体平均値とすることを特徴とする、請求項 1 に記載のパターンマッチング装置。

【請求項 4】

パターンマッチングの対象とする前記特徴量が含まれる範囲を同定する範囲同定手段を備え、

前記全体平均取得手段は、前記範囲同定手段で同定された範囲に基づく前記音声データの全フレーム数または前記画像データ全体の前記特徴量から前記全体平均値を計算することを特徴とする、請求項 2 に記載のパターンマッチング装置。

10

【請求項 5】

前記局所平均計算手段は、過去に計算した前記局所平均値を重み付けした値に基づいて、前記局所平均値を計算し、

前記局所分散計算手段は、過去に計算した前記局所分散値を重み付けした値に基づいて、前記局所分散値を計算する

ことを特徴とする請求項 1 ~ 4 に記載のパターンマッチング装置。

【請求項 6】

外部より入力された音声データまたは画像データの特徴量を算出する分析手段と、

前記分析手段で算出された前記特徴量を正規化する正規化手段と、

前記正規化手段で正規化された正規化済み特徴量に基づいて、パターンマッチングを行うパターンマッチング手段と、

20

としてコンピュータを機能させるためのパターンマッチングプログラムにおいて、

前記正規化手段は、

前記音声データの全フレーム数または前記画像データ全体の前記特徴量の平均値である全体平均値を取得する全体平均取得手段と、

前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の平均値である局所平均値を計算する局所平均計算手段と、

前記局所平均値に基づいて、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の分散値である局所分散値を計算する局所分散計算手段と、

前記全体平均値と複数の前記局所分散値とに基づいて前記特徴量を正規化する正規化処理計算手段と、

30

としてコンピュータを機能させるためのパターンマッチングプログラム。

【請求項 7】

外部より入力された音声データまたは画像データの特徴量を算出する分析ステップと、

前記分析ステップで算出された前記特徴量を正規化する正規化ステップと、

前記正規化ステップで正規化された正規化済み特徴量に基づいて、パターンマッチングを行うパターンマッチングステップと、

を備えたパターンマッチング方法において、

前記正規化ステップは、

前記音声データの全フレーム数または前記画像データ全体の前記特徴量の平均値である全体平均値を取得する全体平均取得ステップと、

40

前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の平均値である局所平均値を計算する局所平均計算ステップと、

前記局所平均値に基づいて、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の分散値である局所分散値を計算する局所分散計算ステップと、

前記全体平均値と複数の前記局所分散値とに基づいて前記特徴量を正規化する正規化処理計算ステップと、

を備えたことを特徴とするパターンマッチング方法。

【発明の詳細な説明】

【技術分野】

50

## 【 0 0 0 1 】

本発明は、パタンマッチング装置、パタンマッチングプログラム、およびパタンマッチング方法に関する。

## 【背景技術】

## 【 0 0 0 2 】

音声認識装置は、入力音声信号から抽出された時系列の音響特徴量を、母音や子音などの音素を単位として、音響特徴量空間における確率密度分布が予め学習された音響モデルと照合することにより認識結果を得る。確率モデルである音響モデルは、音響特徴量の入力に対して、その音素らしさのスコア（音響尤度）を出力する。音声認識装置は文法と単語辞書の制約に従って音素らしさのスコア（音響尤度）を発声全体に渡って累積し、累積スコアが最も高い単語の並びを認識結果として出力する。

10

## 【 0 0 0 3 】

音響特徴量は多次元ベクトルの時系列データであり、各次元において各音素に該当するデータの頻度分布を集計すると正規分布に近い形状、もしくは複数の正規分布の和に近い形状になる。こうした音響特徴量の分布を表現するために、音響モデルの確率密度分布は多次元正規分布もしくは複数の多次元正規分布によって表現される。しかし、実際の照合においては、マイク特性のばらつき、話者による違い、背景雑音などにより、入力音響特徴量の分布と音響モデルの確率密度分布との間に mismatches が生じ、認識率低下の原因となる。入力音響特徴量と音響モデルの照合において、この mismatches を解消する手法として、ケプストラム平均値正規化（CMN: Cepstral Mean Normalization）という手法が

20

## 【 0 0 0 4 】

## 【数 1】

$$x_c(t) = x(t) - E(x) \quad \dots(1) \quad E(x) = \frac{1}{T} \sum_{t=1}^T x(t) \quad \dots(2)$$

30

## 【 0 0 0 5 】

一方、MVNとは、発声の各時刻の音響特徴量を、その発声全体の平均値と分散で正規化して、基準系の正規分布N（平均0、分散1）に揃えることで、マイク特性などによる入力音響特徴量の分布と音響モデルの確率密度分布との mismatches を低減する手法である。MVN前の各次元の音響特徴量を  $x(t)$ 、MVN後の音響特徴量を  $x_m(t)$  とすると、MVNの操作は（3）～（5）式で表される。

40

## 【 0 0 0 6 】

【数2】

$$x_m(t) = \frac{x(t) - E(x)}{\sqrt{V(x)}} \quad \dots(3) \quad E(x) = \frac{1}{T} \sum_{t=1}^T x(t) \quad \dots(4)$$

$$V(x) = \frac{1}{T-1} \sum_{t=1}^T \{x(t) - E(x)\}^2 \quad \dots(5)$$

10

【0007】

また、音声に限らず、静止画および動画についても、CMNおよびMVNにて正規化が可能である。静止画像の場合、各次元の画像特徴量を $x_{i,j}$ 、CMN後の画像特徴量を $x_{c,i,j}$ とすると、CMNの操作は式(6)、(7)で表される。I、Jは静止画の縦軸、横軸のブロック数を表す。

【0008】

【数3】

$$x_{c,i,j} = x_{i,j} - E(x_{i,j}) \quad \dots(6) \quad E(x_{i,j}) = \frac{1}{I \times J} \sum_{i=1}^I \sum_{j=1}^J x_{i,j} \quad \dots(7)$$

20

【0009】

一方、MVNでは、MVN前の各次元の画像特徴量を $x_{i,j}$ 、MVN後の画像特徴量を $x_{m,i,j}$ とすると、MVNの操作は式(8)~(10)で表される。

【0010】

【数4】

$$x_{m,i,j} = \frac{x_{i,j} - E(x_{i,j})}{\sqrt{V(x_{i,j})}} \quad \dots(8) \quad E(x_{i,j}) = \frac{1}{I \times J} \sum_{i=1}^I \sum_{j=1}^J x_{i,j} \quad \dots(9)$$

$$V(x_{i,j}) = \frac{1}{(I-1) \times (J-1)} \sum_{i=1}^I \sum_{j=1}^J \{x_{i,j} - E(x_{i,j})\}^2 \quad \dots(10)$$

40

【0011】

動画の場合、各次元の動画特徴量を $x_{i,j,t}$ 、CMN後の動画特徴量を $x_{c,i,j,t}$ とすると、CMNの操作は式(11)、(12)で表される。I、Jは動画の縦軸、横軸のブロック数、Tはフレーム数を表す。

【0012】

【数5】

$$x_{c_{i,j,t}} = x_{i,j,t} - E(x_{i,j,t}) \cdots (11) \quad E(x_{i,j,t}) = \frac{1}{I \times J \times T} \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^T x_{i,j,t} \cdots (12)$$

【0013】

一方、MVNでは、MVN前の各次元の動画特徴量を $x_{i,j,t}$ 、MVN後の動画特徴量を $x_{m_{i,j,t}}$ とすると、MVNの操作は式(13)~(15)で表される。

10

【0014】

【数6】

$$x_{m_{i,j,t}} = \frac{x_{i,j,t} - E(x_{i,j,t})}{\sqrt{V(x_{i,j,t})}} \cdots (13)$$

$$E(x_{i,j,t}) = \frac{1}{I \times J \times T} \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^T x_{i,j,t} \cdots (14)$$

20

$$V(x_{i,j,t}) = \frac{1}{(I-1) \times (J-1) \times (T-1)} \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^T \{x_{i,j,t} - E(x_{i,j,t})\}^2 \cdots (15)$$

【0015】

ただし、発声全体の平均値や分散を用いるCMNやMVNは、発声が終わるまで正規化後の音響特徴量が得られないために照合処理の開始が遅れ、発声終了から認識結果出力までの待ち時間を長くしてしまうというデメリットがある。この処理遅れを低減する手法として、発声全体の代わりに数十~数百ミリ秒の局所の区間から平均値や分散を算出して正規化に用いる手法が提案されている。以降、発声全体から計算した平均値を用いて音響特徴量を正規化する手法をバッチCMN、発声の一部区間から計算した平均値を用いて音響特徴量を正規化する手法をセグメンタルCMNとよぶ。同様に、発声全体から計算した平均値と分散値を用いて音響特徴量を正規化する手法をバッチMVN、発声の一部区間から計算した平均値と分散値を用いて音響特徴量を正規化する手法をセグメンタルMVNと呼ぶ。

30

【0016】

また、特徴量の量子化を仮定しない平均値・分散正規化(MVN)において算出した分散の値がゼロもしくはゼロに近い小さな値の場合には分散正規化を行わない手法も知られている(例えば、特許文献1参照)。

40

【特許文献1】特開2002-278586号公報

【発明の開示】

【発明が解決しようとする課題】

【0017】

しかし、CMNでは、バッチCMNの方が、セグメンタルCMNよりも長い音声区間から特徴量の平均値を算出するため、精度が高く認識率の改善効果が高いが、入力音響特徴量の分布のばらつきと、参照する音響モデルの確率密度分布のばらつきまで揃えることはできない。

50

## 【0018】

また、MVNでは、バッチMVNは発声全体の音響特徴量の分布を平均0、分散1に正規化するが、音声認識の単位となる音素ごとの分布に着目すると、分散は正規化されていない。一方、セグメンタルMVNで平均・分散の計算区間を1音素相当の時間長(数十から数百ミリ秒)に設定すれば、音素ごとの分布の分散を正規化するのに近い効果が得られる。ただし、短時間の平均値も0に正規化されるので、すべての音素の分布の平均値が0に近づくため重なりが大きくなり(図2参照)、音素の識別能力の低下を招く。

## 【0019】

また、特許文献1では、発声全体の平均値と分散値を用いて正規化するバッチMVNと、局所の平均値と分散値を用いて音響特徴量を正規化するセグメンタルMVNへの適用についてのみ述べられており、前述の音素の識別能力の低下を招くという問題点を解決することができない。

10

## 【0020】

すなわち、CMNでは分布のばらつき(分散)を正規化することができず、セグメンタルMVNでは、音素ごとの分散の正規化に近い効果があるが、音素間で分布の平均値が近づいてしまい音素の識別能力が低下してしまうという問題がある。

## 【0021】

また、上記の課題は、音声に限らず、外部より入力されたデータの特徴量を算出し、算出した特徴量を正規化し、正規化済み特徴量に基づいてパタンマッチングを行うパタンマッチング装置にも当てはまる。

20

## 【0022】

本発明は、上記の課題を解決するためになされたものであり、特徴量の識別能力を低下させることなく特徴量を正規化することが可能なパタンマッチング装置、パタンマッチングプログラム、およびパタンマッチング方法を提供することを目的とする。

## 【課題を解決するための手段】

## 【0023】

本発明は、外部より入力された音声データまたは画像データの特徴量を算出する分析手段と、前記分析手段で算出された前記特徴量を正規化する正規化手段と、前記正規化手段で正規化された正規化済み特徴量に基づいて、パタンマッチングを行うパタンマッチング手段と、を備えたパタンマッチング装置において、前記正規化手段は、前記音声データの全フレーム数または前記画像データ全体の前記特徴量の平均値である全体平均値を取得する全体平均取得手段と、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の平均値である局所平均値を計算する局所平均計算手段と、前記局所平均値に基づいて、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の分散値である局所分散値を計算する局所分散計算手段と、前記全体平均値と複数の前記局所分散値とに基づいて前記特徴量を正規化する正規化処理計算手段と、を備えたことを特徴とするパタンマッチング装置である。

30

## 【0024】

また、本発明の前記全体平均取得手段は、前記音声データの全フレーム数または前記画像データ全体の前記特徴量から前記全体平均値を計算することを特徴とする。

40

## 【0025】

また、本発明の前記全体平均取得手段は、予め記憶した所定値を前記全体平均値とすることを特徴とする。

## 【0026】

また、本発明は、パタンマッチングの対象とする前記特徴量が含まれる範囲を同定する範囲同定手段を備え、前記全体平均取得手段は、前記範囲同定手段で同定された範囲に基づく前記音声データの全フレーム数または前記画像データ全体の前記特徴量から前記全体平均値を計算することを特徴とする、請求項2に記載のパタンマッチング装置である。

## 【0027】

また、本発明の前記局所平均計算手段は、過去に計算した前記局所平均値により重み付

50

けた値に基づいて、前記局所平均値を計算し、前記局所分散計算手段は、過去に計算した前記局所分散値により重み付けした値に基づいて、前記局所分散値を計算することを特徴とする。

【0028】

また、本発明は、外部より入力された音声データまたは画像データの特徴量を算出する分析手段と、前記分析手段で算出された前記特徴量を正規化する正規化手段と、前記正規化手段で正規化された正規化済み特徴量に基づいて、パタンマッチングを行うパタンマッチング手段と、としてコンピュータを機能させるためのパタンマッチングプログラムにおいて、前記正規化手段は、前記音声データの全フレーム数または前記画像データ全体の前記特徴量の平均値である全体平均値を取得する全体平均取得手段と、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の平均値である局所平均値を計算する局所平均計算手段と、前記局所平均値に基づいて、前記第2の範囲に含まれる前記特徴量の分散値である局所分散値を計算する局所分散計算手段と、前記全体平均値と複数の前記局所分散値とに基づいて前記特徴量を正規化する正規化処理計算手段と、としてコンピュータを機能させるためのパタンマッチングプログラムである。

10

【0029】

また、本発明は、外部より入力された音声データまたは画像データの特徴量を算出する分析ステップと、前記分析ステップで算出された前記特徴量を正規化する正規化ステップと、前記正規化ステップで正規化された正規化済み特徴量に基づいて、パタンマッチングを行うパタンマッチングステップと、を備えたパタンマッチング方法において、前記正規化ステップは、前記音声データの全フレーム数または前記画像データ全体の前記特徴量の平均値である全体平均値を取得する全体平均取得ステップと、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の平均値である局所平均値を計算する局所平均計算ステップと、前記局所平均値に基づいて、前記音声データの局所のフレーム数または前記画像データの局所範囲の前記特徴量の分散値である局所分散値を計算する局所分散計算ステップと、前記全体平均値と複数の前記局所分散値とに基づいて前記特徴量を正規化する正規化処理計算ステップと、を備えたことを特徴とするパタンマッチング方法である。

20

【発明の効果】

【0030】

本発明によれば、特徴量の識別能力を低下させることなく特徴量を正規化することができる。

30

【発明を実施するための最良の形態】

【0031】

以下、図面を参照し、本発明の実施形態を説明する。図1は本発明の一実施形態による音声認識装置の構成を示している構成図である。音響分析部101は、マイク等より入力された音声データに対して音響分析を行い、音響特徴量を計算する。入力は、プッシュ・トゥークで制御することも可能である。また、音響分析部101は、計算した音響特徴量を一時的にバッファに記憶させる。正規化処理部102は、音響分析部101がバッファに記憶させた音響特徴量を、音響特徴量の平均値および分散値を用いて正規化処理を行う。正規化処理については後述する。音響モデル学習部103は、学習用音声データに対して、認識対象の音声データと同一の音響分析を音響分析部101で行い、正規化処理部102で正規化を行って得た、学習用音声データの音響特徴量を音響モデル記憶部104に記憶させる。言語モデル記憶部105は、単語辞書や文法を記憶する。認識処理部106は、認識対象の音声データに対して音響分析部101で音響分析を行い、正規化処理部102で正規化処理を行って得た、認識対象の音声データの音響特徴量と音響モデル記憶部104が記憶している学習用音声データの音響特徴量および言語モデルが記憶している単語辞書や文法を用いてパタンマッチングを行い、認識結果を出力する。

40

【0032】

[第1の実施形態]

50

まず、本発明の第1の実施形態を説明する。図3は、本実施形態による正規化処理部102の構成を示している。マイクなどから音声認識装置に入力された1発声全体の音響特徴量は、音響分析部101によって図示せぬバッファに格納されている。全体平均計算部301は、発声全体に対応したフレーム数T内の音響特徴量をバッファから読み出し、その平均値を計算する。発声全体の長さとして、単語の長さ、音声の切れ目までの長さ、句読点から句読点までの長さ、入力された音声全体の長さなどを用いることが可能である。発声全体の音響特徴量の平均値 $E(x)$ は計算式(16)で求める。

【0033】

【数7】

$$E(x) = \frac{1}{T} \sum_{t=1}^T x(t) \quad \dots(16)$$

10

【0034】

局所平均計算部302は、予め設定した局所のフレーム数 $\tau$ 内の発声に対応した音響特徴量をバッファから読み出し、その平均値を計算する。局所のフレーム数 $\tau$ は、音素の長さとして、例えば数十から数百ミリ秒に対応した数である。音素の長さなので発声する単語や人によって変動するが、本実施形態では固定値を使用する。局所のフレーム数 $\tau$ の音響特徴量の平均値 $E_\tau(x)$ は計算式(17)で求める。

20

【0035】

【数8】

$$E_\tau(x) = \frac{1}{\tau} \sum_{t=1}^{\tau} x(t) \quad \dots(17)$$

30

【0036】

局所分散計算部303は、予め設定した局所のフレーム数 $\tau$ 内の発声に対応した音響特徴量をバッファから読み出し、その分散値を、局所平均計算部302で算出した平均値に基づいて計算する。局所のフレーム数 $\tau$ の音響特徴量の分散値 $V_\tau(x)$ は計算式(18)で求める。

【0037】

【数9】

$$V_\tau(x) = \frac{1}{\tau-1} \sum_{t=1}^{\tau} \{x(t) - E_\tau(x)\}^2 \quad \dots(18)$$

40

【0038】

正規化処理計算部304は、正規化前の音声特徴量から全体平均計算部301で算出した発声全体に対しての音響特徴量の平均値を減算し、局所分散計算部303で算出した局所のフレーム数 $\tau$ の音響特徴量の分散値で割ることで、正規化後の音響特徴量 $x(t)$ を求めることができる(計算式(19)参照)。

50



【 0 0 3 9 】

【 数 1 0 】

$$x_{\tau}(t) = \frac{x(t) - E(x)}{\sqrt{V_{\tau}(x)}} \quad \dots(19)$$

【 0 0 4 0 】

10

上述したとおり、入力音響特徴量に対して、発声全体の平均値による正規化処理を行うことにより、すべての音素の分布の位置を音響モデルの該当音素の分布に揃え、更に局所の分散値による正規化処理によって、全音素の分布の重なりを抑制しつつ正規分布に近づける効果を持つ（図4参照）。その結果、音素間の識別精度を低減することなく、背景雑音や残響などによる音響モデルと入力された音響特徴量のミスマッチ成分を低減することができ、音声認識精度の劣化を低減することができる。

【 0 0 4 1 】

なお、全体平均計算部301で、平均を求めるフレーム数を発声時間に対応する数としたが、代わりにフレーム数 $\tau$ を予め設定してもよい。

【 0 0 4 2 】

20

[第2の実施形態]

次に、本発明の第2の実施形態を説明する。図5は、本実施形態による正規化処理部102の構成を示している。本実施形態では、対象とする局所のフレーム数での局所平均値および局所分散値を算出する際に、1つ前の局所のフレーム数の音響特徴量から計算した局所平均値（以下、1つ前の局所平均値と記す。）および1つ前の局所のフレーム数の音響特徴量から計算した局所分散値（以下、1つ前の局所分散値と記す。）を用いることを特徴とする。突発的な雑音が音声認識装置に入力された場合、局所平均値および局所分散値が大きく変わり、入力された音声データを正しく認識することが困難となるが、1つ前の局所平均値および1つ前の局所分散値を用いることで、突発的に音声認識装置に雑音が入力された場合でも局所平均値および局所分散値が大きく変わらず、音声認識制度の劣化を低減することができる。

30

【 0 0 4 3 】

全体平均計算部501は、第1の実施形態と同様に音声認識装置に入力された発声全体に対応した音響特徴量をバッファから読み出し、その平均値 $E(x)$ を計算する。局所平均計算部502は、予め設定した局所のフレーム数 $\tau$ 内の発声に対応した音響特徴量をバッファから読み出し、その平均値を計算する。その際忘却係数 $\alpha$ を予め設定し、1つ前の局所平均値を重み付け加算する。1つ前の局所平均値を重み付け加算した、局所のフレーム数 $\tau$ の音響特徴量の局所平均値 $E_p(t)$ は計算式(20)で求める。

【 0 0 4 4 】

【 数 1 1 】

40

$$E_p(t) = \frac{\alpha}{\tau} \sum_{i=t-\tau}^t x(i) + (1-\alpha)E_p(t-1) \quad \dots(20)$$

【 0 0 4 5 】

局所分散計算部503は、予め設定した局所のフレーム数 $\tau$ 内の発声に対応した音響特徴量をバッファから読み出し、その分散値を、局所平均計算部502で算出した平均値に

50

基づいて計算する。その際忘却係数  $\alpha$  を予め設定し、1つ前の局所平均値を重み付け加算する。1つ前の局所平均値を重み付け加算した、局所のフレーム数  $\tau$  の音響特徴量の局所分散値  $V_p(t)$  は計算式(21)で求める。

【0046】

【数12】

$$V_p(t) = \frac{\alpha}{\tau - 1} \sum_{i=t-\tau}^t (x(i) - E_p(i))^2 + (1 - \alpha)V_p(t-1) \quad \dots(21)$$

10

【0047】

正規化処理計算部504は、正規化前の音声特徴量から全体平均計算部501で算出した発声全体の音響特徴量の平均値を減算し、局所分散計算部503で算出した局所のフレーム数  $\tau$  の音響特徴量の分散値で割ることで、正規化後の音響特徴量  $x_p(t)$  を求めることができる(計算式(22)参照)。

【0048】

【数13】

$$x_p(t) = \frac{x(t) - E_p(x)}{\sqrt{V_p(x)}} \quad \dots(22)$$

20

【0049】

上述したとおり、入力音響特徴量に対して、発声全体の平均値による正規化処理を行うことにより、すべての音素の分布の位置を音響モデルの該当音素の分布に揃え、更に局所の分散値による正規化処理によって、全音素の分布の重なりを抑制しつつ正規分布に近づける効果を持つ(図4参照)。その結果、音素間の識別精度を低減することなく、背景雑音や残響などによる音響モデルと入力された音響特徴量のミスマッチ成分を低減することができ、音声認識精度の劣化を低減することができる。さらに、突発的な雑音が音声認識装置に入力された場合、入力された音声データを認識することが困難となるが、1つ前の局所平均値および1つ前の局所分散値を用いることで、突発的に音声認識装置に雑音が入力された場合でも平均値が大きく変わらず、音声認識精度の劣化を低減することができる。

30

【0050】

[第3の実施形態]

次に、本発明の第3の実施形態を説明する。図6は、本実施形態による正規化処理部102の構成を示している。本実施形態では、実施形態1での発声全体の音響特徴量の平均値を算出する代わりに、予め計算した固定の平均値を用いることを特徴とする。これにより、発声全体から音響特徴量の平均値を計算する必要がないため、音響特徴量の正規化が完了するまでの待ち時間が、局所分散の計算に必要な時間となり、リアルタイム処理が可能となる。

40

【0051】

固定平均値記憶部601は、予め設定した音響特徴量の平均値  $E_f(x)$  を記憶する。固定値は、前の発声の平均値を用いる、もしくは過去の莫大な音声データから求めることなどが可能である。

【0052】

50

局所平均計算部 602、局所分散計算部 603 は、第 1 の実施形態と同様に局所平均値および局所分散値を算出する。正規化処理計算部 602 は、固定平均値記憶部 601 に記憶された固定平均値を用い、正規化前の音声特徴量から固定平均値を減算し、局所分散計算部 603 で算出した局所のフレーム数 の音響特徴量の分散値で割ることで、正規化後の音響特徴量  $x_f(t)$  を求めることができる（計算式（23）参照）。

【0053】

【数14】

$$x_f(t) = \frac{x(t) - E_f(x)}{\sqrt{V_\tau(x)}} \quad \dots(23)$$

10

【0054】

上述したとおり、発声全体の音響特徴量の平均値を算出する代わりに、予め計算した固定の平均値を用いることで、発声全体から音響特徴量の平均値をリアルタイムに計算する必要がない。これにより、音響特徴量の正規化が完了するまでの待ち時間が局所分散の計算に必要な時間となり、リアルタイム処理が可能となる。また、入力音響特徴量に対して、発声全体の平均値による正規化処理を行うことにより、すべての音素の分布の位置を音響モデルの該当音素の分布に揃え、更に局所の分散値による正規化処理によって、全音素の分布の重なりを抑制しつつ正規分布に近づける効果を持つ（図4参照）。その結果、音素間の識別精度を低減することなく、背景雑音や残響などによる音響モデルと入力された音響特徴量のミスマッチ成分を低減することができ、音声認識精度の劣化を低減することができる。

20

【0055】

[第4の実施形態]

次に、本発明の第4の実施形態を説明する。図7は、本実施形態による正規化処理部102の構成を示している。本実施形態では、全体平均計算部702の前段に音声検出部を設ける事により音声区間を同定し、音声区間とその前後の数十ミリ秒を加えた時間に対応するフレーム数  $\tau'$  での平均値を用いて正規化することを特徴とする。これにより、発声終了後に無音区間が長く続いた場合においても、正規化処理までの待ち時間を短くすることが可能となる。

30

【0056】

音声検出部701は、入力された音響特徴量に音声特有の特徴が含まれていることを検出し、音声区間を同定する。音声特有の特徴としては、音声のパワー、ケプストラム値、周波数などを用いることが可能である。全体平均計算部702は、音声検出部701で同定した音声区間とその前後の数十ミリ秒を加えた時間に対応するフレーム数  $\tau'$  での発声に対応した音響特徴量をバッファから読み出し、その平均値を計算する。 $\tau'$  の音響特徴量の平均値  $E_{\tau'}(x)$  は計算式（24）で求める。

40

【0057】

【数15】

$$E_{\tau'}(x) = \frac{1}{\tau'} \sum_{t=1}^{\tau'} x(t) \quad \dots(24)$$

【0058】

50

局所平均計算部 703 は、予め設定した局所のフレーム数 内の発声に対応した音響特徴量をバッファから読み出し、その平均値を計算する。局所のフレーム数は、音素の長さとして、例えば数十から数百ミリ秒に対応した数である。音素の長さなので発声する単語や人によって変動するが、本実施形態では固定値を使用する。局所のフレーム数 の音響特徴量の平均値  $E_{\tau}(x)$  は計算式 (17) で求める。

【0059】

局所分散計算部 704 は、予め設定した局所のフレーム数 内の発声に対応した音響特徴量をバッファから読み出し、その分散値を、局所平均計算部 703 で算出した平均値に基づいて計算する。局所のフレーム数 の音響特徴量の分散値  $V_{\tau}(x)$  は計算式 (18) で求める。

10

【0060】

正規化処理計算部 705 は、正規化前の音声特徴量から全体平均計算部 702 で算出した発声全体の音響特徴量の平均値を減算し、局所分散計算部 704 で算出した局所のフレーム数 の音響特徴量の分散値で割ることで、正規化後の音響特徴量  $x_{\tau'}(t)$  を求めることができる (計算式 (25) 参照)。

【0061】

【数 16】

$$x_{\tau'}(t) = \frac{x(t) - E_{\tau}(x)}{\sqrt{V_{\tau}(x)}} \quad \dots(25)$$

20

【0062】

上述したとおり、全体平均計算部 702 の前段に音声検出部 701 を設ける事により音声区間を同定し、音声区間とその前後の数十ミリ秒を加えた時間に対応するフレーム数  $\tau'$  での平均値を用いて正規化することにより、発声終了後に無音区間が長く続いた場合においても、正規化処理までの待ち時間を短くすることが可能となる。

【0063】

30

[第 5 の実施形態]

次に、本発明の第 5 の実施形態を説明する。図 8 は本実施形態による画像認識装置の構成を示している構成図である。図 4 において、マイクから入力された音声データの代わりにカメラから入力された画像とし、単語辞書・文法と音響モデルの代わりにオブジェクトモデルとし、音声認識結果の代わりに画像認識結果と置き換えることで、画像認識への適用も可能となる。

【0064】

画像特徴量分析部 801 は、カメラから入力された画像データに対して画像特徴量分析を行い、画像特徴量を計算する。正規化処理部 802 は、画像特徴量分析部 801 で計算した画像特徴量を画像特徴量の平均値および分散値を用いて正規化処理を行う。正規化処理については後述する。オブジェクトモデル学習部 803 は、学習用画像データに対して、認識対象の画像データと同一の画像特徴量分析を画像特徴量分析部 801 でを行い、正規化処理部 802 で正規化を行って得た、学習用画像データの画像特徴量をオブジェクトモデル 804 に記憶させる。認識処理部 805 は、認識対象の画像データに対して画像特徴量分析部 801 で画像特徴量分析を行い、正規化処理部 802 で正規化処理を行って得た、認識対象の画像データの画像特徴量とオブジェクトモデル 804 が記憶している学習用画像データの画像特徴量を用いて認識処理を行い、認識結果を出力する。

40

【0065】

図 9 を参照し本実施形態における画像の正規化処理について説明する。図 9 は、本実施形態による正規化処理部 102 の構成を示している。全体平均計算部 901 は、カメラ等

50

から画像認識装置に入力された画像データ全体の画像特徴量をバッファから読み出し、その平均値を計算する。画像データ全体の画像特徴量の平均値  $E(x_{i,j})$  は計算式(26)で求める。I、Jは静止画の縦軸、横軸のブロック数を表す。

【0066】

【数17】

$$E(x_{i,j}) = \frac{1}{I \times J} \sum_{i=1}^I \sum_{j=1}^J x_{i,j} \quad \dots(26)$$

10

【0067】

局所平均計算部902は、予め設定した画像データの局所範囲における画像特徴量の平均値を計算する。局所範囲としては、正規化対象画像範囲を含む周囲数ブロックなどを用いることが可能である。局所範囲(k,l)の画像特徴量の平均値  $E(x_{k,l})$  は計算式(27)で求める。I、Jは静止画の局所範囲での縦軸、横軸のブロック数を表す。

【0068】

【数18】

$$E(x_{k,l}) = \frac{1}{K \times L} \sum_{k=1}^K \sum_{l=1}^L x_{k,l} \quad \dots(27)$$

20

【0069】

局所分散計算部903は、予め設定した画像データの局所範囲における画像特徴量の分散値を、局所平均計算部902で算出した平均値に基づいて計算する。局所範囲(k,l)の画像特徴量の分散値  $V(x_{k,l})$  は計算式(28)で求める。

【0070】

【数19】

$$V(x_{k,l}) = \frac{1}{(K-1) \times (L-1)} \sum_{k=1}^K \sum_{l=1}^L \{x_{k,l} - E(x_{k,l})\}^2 \quad \dots(28)$$

30

【0071】

正規化処理計算部904は、正規化前の画像特徴量から全体平均計算部901で算出した画像全体の画像特徴量の平均値を減算し、局所分散計算部903で算出した予め設定した画像データの範囲における画像特徴量の分散値で割ることで、正規化後の画像特徴量  $x_{k,l}$  を求めることができる(計算式(29)参照)。

40

【0072】

【数 20】

$$x_{k,l} = \frac{x_{i,j} - E(x_{i,j})}{\sqrt{V(x_{k,l})}} \quad \dots(29)$$

10

【0073】

上述したとおり、画像認識においても、画像特徴量に対して画像全体の平均値による正規化処理を行うことにより、すべての画像特徴量の分布の位置をオブジェクトモデルの該当画像特徴量の分布に揃え、更に局所の分散値による正規化処理によって、全画像特徴量の分布の重なりを抑制しつつ正規分布に近づける効果を持つ。その結果、画像特徴量の識別精度を低減することなく、影や輝度などによるオブジェクトモデルと入力された画像特徴量のミスマッチ成分を低減することができ、画像認識精度の劣化を低減することができる。

【0074】

なお、画像認識については、平面画像だけではなく、3D画像でも可能である。3D画像を作成する際にカメラの位置によって、対象物の陰が変わるが、本発明の正規化を用いることで、画像特徴量のミスマッチ成分を低減することができ、画像認識精度の劣化を低減することができる。

20

【0075】

[第6の実施形態]

また、画像認識に時間要素を取り入れることで、動画についても動画特徴量のミスマッチ成分を低減することができ、動画認識精度の劣化を低減することができる。

【0076】

本発明の第6の実施形態を説明する。図10は本実施形態による動画認識装置の構成を示している構成図である。図4において、マイクから入力された音声データの代わりにカメラから入力された動画とし、単語辞書・文法記憶部と音響モデル記憶部の代わりにオブジェクトモデル記憶部とし、音声認識結果の代わりに動画認識結果と置き換えることで、動画認識への適用も可能となる。

30

【0077】

動画特徴量分析部1001は、カメラから入力された動画データに対して動画特徴量分析を行い、動画特徴量を計算する。正規化処理部1002は、動画特徴量分析部1001で計算した動画特徴量を動画特徴量の平均値および分散値を用いて正規化処理を行う。正規化処理については後述する。オブジェクトモデル学習部1003は、学習用動画データに対して、認識対象の動画データと同一の動画特徴量分析を動画特徴量分析部1001で行い、正規化処理部1002で正規化を行って得た、学習用動画データの動画特徴量をオブジェクトモデル1004に記憶させる。認識処理部1005は、認識対象の動画データに対して動画特徴量分析部1001で動画特徴量分析を行い、正規化処理部1002で正規化処理を行って得た、認識対象の動画データの動画特徴量とオブジェクトモデル1004が記憶している学習用動画データの動画特徴量を用いて認識処理を行い、認識結果を出力する。

40

【0078】

図11を参照し本実施形態における動画の正規化処理について説明する。図11は、本実施形態による正規化処理部102の構成を示している。全体平均計算部1101は、カメラ等から動画認識装置に入力された動画データ全体の動画特徴量をバッファから読み出し、その平均値を計算する。動画データ全体の動画特徴量の平均値  $E(x_{i,j}, \tau)$  は計

50

算式(30)で求める。I、Jは動画の縦軸、横軸のブロック数、Tはフレーム数を表す。

【0079】

【数21】

$$E(x_{i,j,t}) = \frac{1}{I \times J \times T} \sum_{i=1}^I \sum_{j=1}^J \sum_{t=1}^T x_{i,j,t} \quad \dots(30)$$

10

【0080】

局所平均計算部1102は、予め設定した動画データの局所範囲における画像特徴量の平均値を計算する。局所範囲としては、正規化対象動画範囲を含む周囲数ブロックおよび局所のフレーム数を用いることが可能である。局所範囲(k, l)および局所のフレーム数の動画特徴量の平均値E(x<sub>k,l,τ</sub>)は計算式(31)で求める。I、Jは動画の局所範囲での縦軸、横軸の区間のブロック数、Tは局所のフレーム数を表す。

【0081】

【数22】

$$E(x_{k,l,\tau}) = \frac{1}{K \times L \times T} \sum_{k=1}^K \sum_{l=1}^L \sum_{\tau=1}^T x_{k,l,\tau} \quad \dots(31)$$

20

【0082】

局所分散計算部1103は、予め設定した動画データの局所範囲における動画特徴量の分散値を、局所平均計算部1102で算出した平均値に基づいて計算する。局所範囲(k, l)および局所のフレーム数の動画特徴量の分散値V(x<sub>k,l,τ</sub>)は計算式(32)で求める。

30

【0083】

【数23】

$$V(x_{k,l,\tau}) = \frac{1}{(K-1) \times (L-1) \times (T-1)} \sum_{k=1}^K \sum_{l=1}^L \sum_{\tau=1}^T \{x_{k,l,\tau} - E(x_{k,l,\tau})\}^2 \quad \dots(32)$$

40

【0084】

正規化処理計算部1104は、正規化前の動画特徴量から全体平均計算部1101で算出した動画全体の動画特徴量の平均値を減算し、局所分散計算部1103で算出した予め設定した動画データの範囲および局所のフレーム数における動画特徴量の分散値で割ることで、正規化後の動画特徴量x<sub>k,l,τ</sub>を求めることができる(計算式(33)参照)。

【0085】

【数 2 4】

$$x_{k,l,\tau} = \frac{x_{i,j,t} - E(x_{i,j,t})}{\sqrt{V(x_{k,l,\tau})}} \quad \dots(33)$$

【0086】

10

以上、この発明の実施形態について図面を参照して詳述してきたが、具体的な構成はこの実施形態に限られるものではなく、この発明の要旨を逸脱しない範囲の設計等も含まれる。

【0087】

例えば、音声、画像、および動画について詳述してきたが、音声、画像、および動画に限らず、入力されたデータの特徴量に基づいてパターンマッチングを行う認識装置にも本発明が適用可能である。

【0088】

また、第2～第4の実施形態については音声認識について説明したが、画像認識および動画認識についても適用可能である。

20

【0089】

また、図1などに示す正規化処理部102の機能を実現するためのプログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録されたプログラムをコンピュータシステムに読み込ませ、実行することにより、正規化処理を行ってもよい。なお、ここでいう「コンピュータシステム」とは、OSや周辺機器等のハードウェアを含むものであってもよい。

【0090】

また、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境（あるいは表示環境）も含むものとする。また、「コンピュータ読み取り可能な記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM、フラッシュメモリ等の書き込み可能な不揮発性メモリ、CD-ROM等の可搬媒体、コンピュータシステムに内蔵されるハードディスク等の記憶装置のことをいう。

30

【0091】

さらに「コンピュータ読み取り可能な記録媒体」とは、インターネット等のネットワークや電話回線等の通信回線を介してプログラムが送信された場合のサーバやクライアントとなるコンピュータシステム内部の揮発性メモリ（例えばDRAM（Dynamic Random Access Memory））のように、一定時間プログラムを保持しているものも含むものとする。

【0092】

また、上記プログラムは、このプログラムを記憶装置等に格納したコンピュータシステムから、伝送媒体を介して、あるいは、伝送媒体中の伝送波により他のコンピュータシステムに伝送されてもよい。ここで、プログラムを伝送する「伝送媒体」は、インターネット等のネットワーク（通信網）や電話回線等の通信回線（通信線）のように情報を伝送する機能を有する媒体のことをいう。また、上記プログラムは、前述した機能の一部を実現するためのものであってもよい。さらに、前述した機能をコンピュータシステムにすでに記録されているプログラムとの組み合わせで実現できるもの、いわゆる差分ファイル（差分プログラム）であってもよい。

40

【図面の簡単な説明】

【0093】

【図1】本発明の一実施形態による音声認識装置の構成を示した構成図である。

【図2】セグメンタルMVN手法の正規化処理による分布の変化の様子を示した図である

50



- 【図3】本発明の第1の実施形態による正規化処理部の構成を示した図である。
- 【図4】本発明の正規化処理による分布の変化の様子を示した図である。
- 【図5】本発明の第2の実施形態による正規化処理部の構成を示した図である。
- 【図6】本発明の第3の実施形態による正規化処理部の構成を示した図である。
- 【図7】本発明の第4の実施形態による正規化処理部の構成を示した図である。
- 【図8】本発明の第5の実施形態による画像認識装置の構成を示した構成図である。
- 【図9】本発明の第5の実施形態による正規化処理部の構成を示した図である。
- 【図10】本発明の第6の実施形態による画像認識装置の構成を示した構成図である。
- 【図11】本発明の第6の実施形態による正規化処理部の構成を示した図である。

10

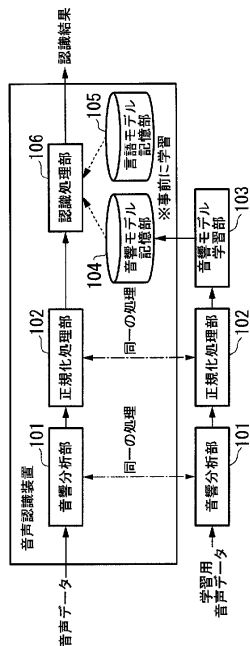
【符号の説明】

【0094】

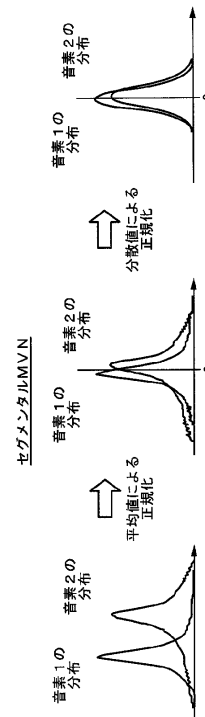
101・・・音声分析部、102,802,1002・・・正規化処理部、103・・・音響モデル学習部、104・・・音響モデル、105・・・言語モデル、106,805・・・認識処理部、301,501,702,901,1101・・・全体平均計算部、302,502,602,703,902,1102・・・局所平均計算部、303,503,603,704,903,1103・・・局所分散計算部、304,504,604,705,904,1104・・・正規化処理計算部、601・・・固定平均値記憶部、701・・・音声検出部、801・・・画像特徴量分析部、803,1003・・・オブジェクトモデル学習部、804,1004・・・オブジェクトモデル、1001・・・動画特徴量分析部

20

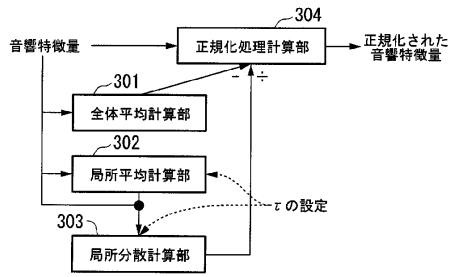
【図1】



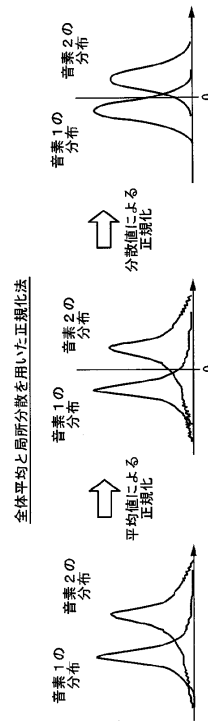
【図2】



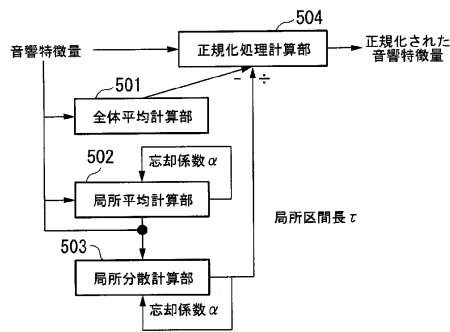
【図3】



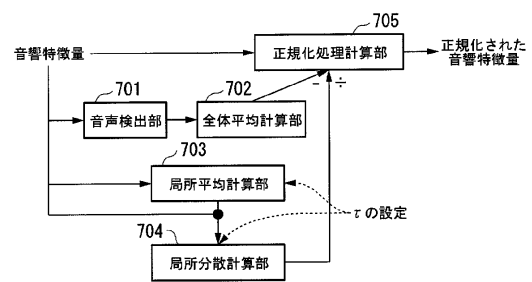
【図4】



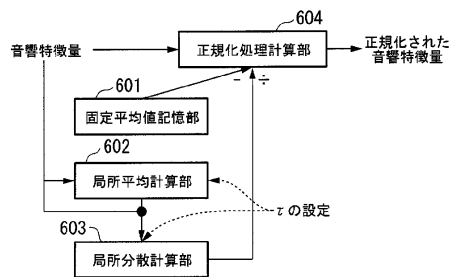
【図5】



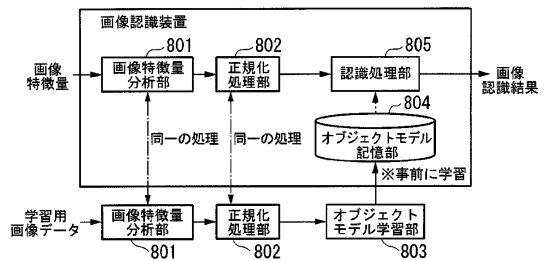
【図7】



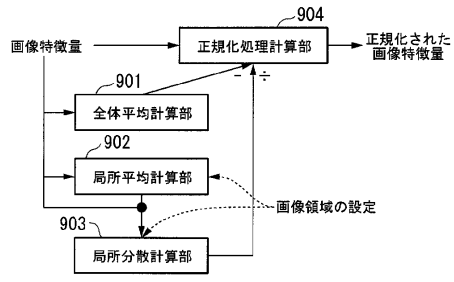
【図6】



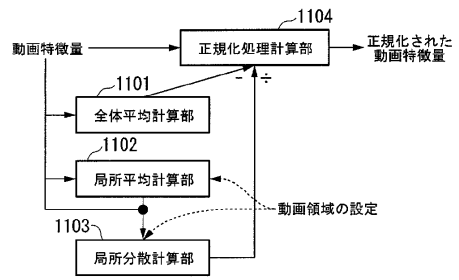
【図8】



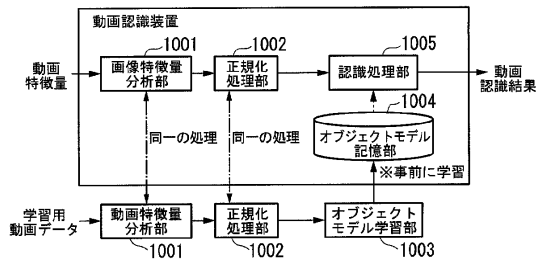
【図9】



【図11】



【図10】



---

フロントページの続き

審査官 前田 祐希

- (56)参考文献 特開2006-084659(JP,A)  
特表2007-536562(JP,A)  
特表2005-521091(JP,A)  
特開平11-085200(JP,A)  
特開2003-167599(JP,A)  
原 一眞, 柘植 覚, 獅子堀 正幹, 北 研二, 黒岩 眞吾, "実時間分散正規化手法の検討 A Study on Real-time Cepstral Variance Normalization", 日本音響学会2003年秋季研究発表会講演論文集 - I - THE 2003 AUTUMN MEETING OF THE ACOUSTICAL SOCIETY OF JAPAN, 日本, 社団法人日本音響学会, 2003年 9月17日, p151~152

(58)調査した分野(Int.Cl., DB名)

G10L 15/00 - 17/00