

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6234060号
(P6234060)

(45) 発行日 平成29年11月22日 (2017.11.22)

(24) 登録日 平成29年11月2日 (2017.11.2)

(51) Int. Cl.		F I			
G 1 0 L	15/06	(2013.01)	G 1 0 L	15/06	3 0 0 Z
G 1 0 L	15/20	(2006.01)	G 1 0 L	15/20	1 5 3
G 1 0 L	15/14	(2006.01)	G 1 0 L	15/14	2 0 0 Z

請求項の数 12 (全 20 頁)

(21) 出願番号	特願2013-99645 (P2013-99645)	(73) 特許権者	390009531
(22) 出願日	平成25年5月9日 (2013.5.9)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公開番号	特開2014-219605 (P2014-219605A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公開日	平成26年11月20日 (2014.11.20)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
審査請求日	平成27年11月27日 (2015.11.27)		New Orchard Road, Armonk, New York 10504, United States of America
前置審査		(74) 代理人	100108501 弁理士 上野 剛史

最終頁に続く

(54) 【発明の名称】 ターゲットドメインの学習用音声データの生成方法、生成装置、および生成プログラム

(57) 【特許請求の範囲】

【請求項1】

コンピュータの演算処理により、ソースドメインの音声データを利用してターゲットドメインの学習用の音声データを生成する方法であって、

ターゲットドメインのクリーンな音声データを用いて学習された前記ターゲットドメインの混合ガウスモデル(Gaussian mixture model :GMM)を読み出すステップと、

前記ターゲットドメインのGMMを参照して、入力として受け取ったソースドメインの音声データの補正量を求めることにより、前記ソースドメインの音声データを、前記ターゲットドメインの音声データに、該ターゲットドメインの音声データのチャンネル特性に基づきマッピングするステップと、

マッピングした前記ソースドメインの音声データに、前記ターゲットドメインのノイズを加えて、擬似的なターゲットドメインの音声データを出力するステップとを含み、

前記マッピングするステップは、前記ソースドメインの音声データを前記ターゲットドメインの音声データにチャンネルマッピング・パラメータによってマッピングする生成モデル式を決定するステップと、前記生成モデル式を参照して前記ターゲットドメインのGMMから変換したソースドメインのGMMに含まれるチャンネルマッピング・パラメータを、Expectation Maximization(EM)アルゴリズムを用いて推定するステップとを含む、

学習データ生成方法。

【請求項2】

前記チャンネルマッピング・パラメータは、前記擬似的なターゲットドメインの音声デー

タを求めるために前記ソースドメインの音声データから差し引くチャンネルバイアスを含む、請求項1に記載の学習データ生成方法。

【請求項3】

前記チャンネルマッピング・パラメータは、前記擬似的なターゲットドメインの音声データを求めるために前記ソースドメインの音声データに乘算するチャンネル振幅を更に含む、請求項2に記載の学習データ生成方法。

【請求項4】

前記EMアルゴリズムを用いて推定するステップは、変換した前記ソースドメインのGMMにソースドメインの観測値を入力して得られる音響尤度を求めるステップと、求めた音響尤度に基づいた目的関数を最小化する前記チャンネルマッピング・パラメータを求めるステップとを交互に繰り返すステップを含む、請求項1に記載の学習データ生成方法。

10

【請求項5】

前記音響尤度を求めるステップにおいて、現在の推定されている前記チャンネルマッピング・パラメータに基づき算出した雑音成分を参照する、請求項4に記載の学習データ生成方法。

【請求項6】

前記生成モデル式を参照して前記ターゲットドメインのGMMから変換したソースドメインのGMMを求める際にVector Taylor Series(VTS)近似を用いる、請求項5に記載の学習データ生成方法。

【請求項7】

20

学習された前記ターゲットドメインの混合ガウスモデル(GMM)は話者の性別ごと用意されており、入力として受け取ったソースドメインの音声データの1発話ごとに男性らしさ及び女性らしさを求めながら、前記ターゲットドメインの音声データにマッピングする、請求項1に記載の学習データ生成方法。

【請求項8】

前記マッピングするステップは、ソースドメインの各音声データについて、前記ターゲットドメインのGMMのガウス分布コンポーネントのうち、該音声データとの音響空間における距離が近いガウス分布コンポーネントの平均との差を求め、該差を、ガウス分布コンポーネントごとの尤度で重みづけした上で、時間方向の平均として求めて前記ソースドメインの音声データに足し合わせるステップを含む、請求項1に記載の学習データ生成方法。

30

【請求項9】

コンピュータに、請求項1乃至8のいずれか一項に記載の学習データ生成方法の各ステップを実行させるための学習データ生成プログラム。

【請求項10】

請求項1乃至8のいずれか一項に記載の学習データ生成方法の各ステップを実行するように適合された手段を備える、学習データ生成システム。

【請求項11】

ソースドメインの音声データを利用してターゲットドメインの学習用の音声データを生成するシステムであって、

40

ターゲットドメインの少量のクリーンな音声データを用いて学習された前記ターゲットドメインの混合ガウスモデル(Gaussian mixture model :GMM)を格納する混合ガウスモデル格納部と、

前記ターゲットドメインのGMMを参照して、入力として受け取ったソースドメインの音声データの補正量を求めることにより、前記ソースドメインの音声データを、前記ターゲットドメインの音声データに、該ターゲットドメインの音声データのチャンネル特性に基づきマッピングするチャンネルマッピング部と、

マッピングした前記ソースドメインの音声データに、前記ターゲットドメインのノイズを加えて、擬似的なターゲットドメインの音声データを出力するノイズ付加部と、
を含み、

50

前記チャンネルマッピング部は、前記ソースドメインの音声データを前記ターゲットドメインの音声データにチャンネルマッピング・パラメータによってマッピングする生成モデル式を決定し、該生成モデル式を参照して前記ターゲットドメインのGMMから変換したソースドメインのGMMに含まれるチャンネルマッピング・パラメータを、Expectation Maximization(EM)アルゴリズムを用いて推定する、

学習データ生成システム。

【請求項 1 2】

前記チャンネルマッピング・パラメータは、前記擬似的なターゲットドメインの音声データを求めるために前記ソースドメインの音声データから差し引くチャンネルバイアスと、前記ソースドメインの音声データに乘算するチャンネル振幅を含む、請求項 1 1 に記載の学習データ生成システム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ソースドメインの音声データを用いてターゲットドメインの音声データを生成する技術に関し、より詳細には、ソースドメインの音声データをターゲットドメインの音声データのチャンネル特性に基づいてマッピングする技術に関する。

【背景技術】

【0002】

音声認識の性能はターゲットドメインの音響環境に大きく依存する。即ち、音響モデルを学習する環境と、音声の評価する環境の間に音響的ミスマッチがあると、音声認識システムの性能は多くの場合低下する。環境間のミスマッチは、背景雑音、録音機器の音響特性、チャンネル歪みなどの様々な原因によって大きくなる。そこで従来は、ターゲットドメインの音響モデルを構築することで環境間のミスマッチを避けるため、特定の環境における音声データを録音することに非常に多くの時間と労力が費やされた。

20

【0003】

これとは対照的に、近年のスマートフォン等のハンドヘルド・デバイスを用いたインターネット・サービス（例えば、音声検索やボイスメールなど）によって、大量の生きた音声データを低コストで取得することが可能となった。そこでこのような種々の音響環境における豊富な音声データを再利用することが望まれる。

30

【0004】

音声認識におけるクロスドメインの問題は、伝統的に以下の 4 つのアプローチに大別される。

1. 再利用法
2. モデル適応法
3. 特徴量変換法
4. 正規化法

1. の再利用法は、ターゲットドメインの音響モデルを構築するために、ソースドメインの音声データを用いてターゲットドメインの音声データをシミュレートする手法である（例えば、非特許文献 1 ~ 2 を参照）。

40

2. のモデル適応法は、ソースドメインの音響モデルのパラメータを変更してテスト音声に合わせる手法であり、最大事後確率推定法（Maximum A Posteriori Estimation : MAP）や、最尤線形回帰法（Maximum Likelihood Linear Regression : MLLR）がこれに該当する（例えば、特許文献 1、非特許文献 3 ~ 非特許文献 5 を参照）。なお、該手法とは異なるが、同様にモデルを適応させる技術として特許文献 2 ~ 3、非特許文献 6 が存在する。

3. の特徴量変換法は、デコード時にテスト音声の特徴量をソースドメインの音響モデルに合わせるように変換する手法であり、特徴量空間最尤線形回帰法（Feature space Maximum Likelihood Linear Regression : fMLLR）や、特徴量空間相互情報量最小化法（Feature space Minimum Mutual Information : fMMI）がこれに該当する（例えば、非特許文献 3 ~ 5、非特許文献 7 を参照）。

50

4. の正規化法は、テスト音声の特徴量の分布を正規化してソースドメインの音響モデルに合わせる手法であり、ケプストラム平均正規化法 (Cepstral Mean Normalization: CMN)、分散平均正規化法 (Mean and Variance Normalization: MVN) がこの手法に該当する (例えば、非特許文献 8 を参照)。

【 0 0 0 5 】

これら 1. ~ 4. の各手法は組み合わせて用いることも可能である。また、2. ~ 4. の各手法は既に確立した技術である。一方で、1. の手法は全プロセスの出発点として重要な技術であるが、該手法に属する既存技術は、上述したインターネットを介して収集される音声データに対して適用することはできない。

【 0 0 0 6 】

非特許文献 1 は、ソースドメインのクリーン音声を入力として、まずターゲットドメインにおけるインパルス応答を畳み込み、その後雑音を加えて、ターゲットドメインの音声をシミュレートする手法を開示する (図 2 A を参照)。該手法は、チャンネルと雑音の特性を補償する最も直接的な手法であるが、インターネット上の音声データをソースデータとする場合には本手法は適さない。なぜならば、ソースデータはクリーン音声とはいえず、また、入力データのチャンネル特性が単一のインパルス応答に対しては多様すぎるためである。

【 0 0 0 7 】

非特許文献 2 は、ステレオデータを用いたマッピング手法を開示する。即ち、非特許文献 2 の技術は、ソースドメインの音声データとターゲットドメインの音声データを同時に記録することを必要とする。ソースとなる音声データをインターネット上のライブデータとする場合、ステレオデータを用意するのは困難であるため、本手法を利用することはできない。

【 0 0 0 8 】

また、1. の手法とは異なるが、目的タスクに合致した音声データを既存の音声コーパスから選択することで目的タスクの音声コーパスを構築する技術を開示するものもある (非特許文献 9 を参照)。

【 0 0 0 9 】

なお、以下の非特許文献 10 ~ 11 は、事前に用意したクリーン音声の混合ガウスモデル (Gaussian mixture model : GMM) と、クリーン音声と観測音声の関係式から、Vector Taylor Series (VTS) 近似より観測音声の GMM を算出する技術を示す背景技術として列挙するものである。

【 先行技術文献 】

【 特許文献 】

【 0 0 1 0 】

【 特許文献 1 】 特開 2 0 1 2 - 4 2 9 5 7 号公報

【 特許文献 2 】 特表 2 0 0 2 - 5 2 9 8 0 0 号公報

【 特許文献 3 】 特開平 1 0 - 1 4 9 1 9 1 号公報

【 非特許文献 】

【 0 0 1 1 】

【 非特許文献 1 】 V. Stahl, A. Fischer, R. Bippus, "Acoustic Synthesis of Training Data for Speech Recognition in Living Room Environments," Proc. of ICASSP, Vol. 1, pp. 285-288, 2001.

【 非特許文献 2 】 J. Droppo, L. Deng, A. Acero, "Evaluation of the SPLICE Algorithm on the Aurora2 Database", Proc. of Eurospeech, pp. 217-220, 2001.

【 非特許文献 3 】 P. J. Moreno, B. Raj, R. M. Stern, "A vector Taylor series approach for environment-independent speech recognition", Proc. of ICASSP, Vol. 2, 1996.

【 非特許文献 4 】 M. L. Seltzer, A. Acero, "Factored Adaptation for Separable Compensation of Speaker and Environmental Variability", Proc. of ASRU, pp. 146-151,

10

20

30

40

50

2011.

【非特許文献5】M. J. F. Gales, "Maximum likelihood linear transformations for HMM based speech recognition," *Computer Speech and Language*, Vol. 12, pp.75-98, 1998.

【非特許文献6】M. J. F. Gales, S. J. Young, "Robust continuous speech recognition using parallel model combination", *IEEE Trans. on Sp. and Audio Proc.*, Vol. 4, pp. 352-359, 1996.

【非特許文献7】B. Kingsbury, L. Mangu, G. Saon, H. Soltau, G. Zweig, "fMPE: Discriminatively Trained Features for Speech Recognition", *Proc. ICASSP*, Vol.1, pp. 961-964, 2005

【非特許文献8】小川 厚徳、高橋 敏、「ケプストラム正規化の実行単位に関する実験的検証」、*電子情報通信学会論文誌D*, Vol. J90-D No.9, pp.2648-2651

【非特許文献9】T. Cincarek, T. Toda, H. Saruwatari, K. Shikano, "Utterance-based Selective Training for the Automatic Creation of Task-Dependent Acoustic Models," *ICETRANSACTIONS on Information and Systems* Vol. E89-D No.3 pp.962-969

【非特許文献10】B. Raj, E. Gouvea, R. M. Stern, "Cepstral compensation using statistical linearization", *Proc. of the ETRW*, 1997.

【非特許文献11】D. Y. Kim, C. K. Un, N. S. Kim, "Speech recognition in noisy environments using first-order vector Taylor series", *Speech Communication*, Vol. 24, pp. 39-49, 1998.

【発明の概要】

【発明が解決しようとする課題】

【0012】

本発明は、上記従来技術における問題点に鑑みてなされたものであり、インターネット上の音声データのように豊富に存在する音響環境の異なる音声データを再利用してターゲットドメインの音声データをシミュレートする方法、装置、およびプログラム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】

【0013】

本願発明は、上記従来技術の課題を解決するために以下の特徴を有する音声データを生成する方法を提供する。本願発明の音声データを生成する方法は、ターゲットドメインのクリーンな音声データを用いて学習された前記ターゲットドメインの混合ガウスモデル(Gaussian mixture model :GMM)を読み出すステップと、前記ターゲットドメインのGMMを参照して、入力として受け取ったソースドメインの音声データを、前記ターゲットドメインの音声データに、該ターゲットドメインの音声データのチャンネル特性に基づきマッピングするステップと、マッピングした前記ソースドメインの音声データに、前記ターゲットドメインのノイズを加えて、擬似的なターゲットドメインの音声データを出力するステップとを含む。ここで、ターゲットドメインのGMMの学習に使用されるターゲットドメインのクリーンな音声データは少量であってよい。

【0014】

好ましくは、前記マッピングするステップは、チャンネルマッピング・パラメータによって前記ソースドメインの音声データを前記ターゲットドメインの音声データにマッピングする生成モデル式を決定するステップと、前記生成モデル式を参照して前記ターゲットドメインのGMMから変換したソースドメインのGMMに含まれるチャンネルマッピング・パラメータを、Expectation Maximization(EM)アルゴリズムを用いて推定するステップとを含む。

【0015】

好ましくは、前記チャンネルマッピング・パラメータは、前記擬似的なターゲットドメインの音声データを求めるために前記ソースドメインの音声データから差し引くチャンネルバイアスを含む。

【0016】

10

20

30

40

50

より好ましくは、前記チャンネルマッピング・パラメータは、前記擬似的なターゲットドメインの音声データを求めるために前記ソースドメインの音声データに乗算するチャンネル振幅を更に含む。

【0017】

好ましくは、前記EMアルゴリズムを用いて推定するステップは、変換した前記ソースドメインのGMMにソースドメインの観測値を入力して得られる音響尤度を求めるステップと、求めた音響尤度に基づいた目的関数を最小化する前記チャンネルマッピング・パラメータを求めるステップとを交互に繰り返すステップを含む。

【0018】

より好ましくは、前記音響尤度を求めるステップにおいて、現在の推定されている前記チャンネルマッピング・パラメータに基づき算出した雑音成分を参照する。

10

【0019】

更に好ましくは、前記生成モデル式を参照して前記ターゲットドメインのGMMから変換したソースドメインのGMMを求める際にVector Taylor Series (VTS)近似を用いる。

【0020】

好ましくは、学習された前記ターゲットドメインの混合ガウスモデル(GMM)は話者の性別ごと用意されており、入力として受け取ったソースドメインの音声データの1発話ごとに男性らしさ及び女性らしさを求めながら、前記ターゲットドメインの音声データにマッピングする。

【0021】

20

また好ましくは、前記マッピングするステップは、ソースドメインの各音声データについて、前記ターゲットドメインのGMMのガウス分布コンポーネントのうち、該音声データとの音響空間における距離が近いガウス分布コンポーネントの平均との差を求め、該差を、ガウス分布コンポーネントごとの尤度で重みづけした上で、時間方向の平均として求めて前記ソースドメインの音声データに足し合わせるステップを含む。

【0022】

以上、ターゲットドメインの学習用の音声データを生成する方法として本発明を説明した。しかし本発明は、そのような学習用音声データ生成方法の各ステップをコンピュータに実行させる学習用音声データ生成プログラム、及び該学習用音声データ生成プログラムを1以上のコンピュータにインストールして実現される学習用音声データ生成システムとして把握することもできる。

30

【発明の効果】

【0023】

本願発明によれば、クリーンな音声データを用いて学習されたターゲットドメインのGMMを参照して、入力として受け取ったソースドメインの音声データを、ターゲットドメインの音声データに、該ターゲットドメインの音声データのチャンネル特性に基づいてマッピングして、ターゲットドメインの音声データをシミュレートする。このため、両方の音声について書き起こしデータがある必要も音声認識する必要もない。そして、インターネット上の音声データのように豊富に存在する音響環境の異なる音声データを再利用してターゲットドメインの音声データをシミュレートすることが可能となる。本発明のその他の効果については、各実施の形態の記載から理解される。

40

【図面の簡単な説明】

【0024】

【図1】本願発明の実施形態による学習用音声データ生成システムを実現するのに好適な情報処理装置のハードウェア構成の一例を示す。

【図2A】従来の学習用音声データをシミュレートするシステムの機能ブロック図である。

。

【図2B】本発明の実施形態による学習用音声データ生成システムの機能ブロック図である。

【図3A】正規化された自動車のノイズを示すグラフである。

50

【図3B】正規化された自動車のノイズを示すグラフである。

【図4】本発明の実施形態による学習用音声データ生成処理の全体の流れの一例を示すフローチャートある。

【図5】本発明の実施形態によるマッピング処理の流れの一例を示すフローチャートある。

【図6】本発明を適用して生成した学習用音声データから構築された音響モデルを用いて行った音声認識の実験結果を示す図である。

【発明を実施するための形態】

【0025】

以下、本願発明を実施するための形態を図面に基づいて詳細に説明するが、以下の実施形態は特許請求の範囲にかかる発明を限定するものではなく、また実施形態の中で説明されている特徴の組み合わせの全てが発明の解決手段に必須であるとは限らない。なお、実施の形態の説明の全体を通じて同じ要素には同じ番号を付している。

【0026】

図1は、本発明を実施するためのコンピュータ100の例示的なハードウェア構成を示す。外部記憶装置114やROM106は、オペレーティング・システムと協働してCPU102に命令を与え、本発明を実施するための複数のコンピュータ・プログラムのコードや各種データを記録することができる。そして外部記憶装置114やROM106に格納された複数のコンピュータ・プログラムは各々RAM104にロードされることによってCPU102により実行される。なお、外部記憶装置114は、SCSIコントローラなどのコントローラ(図示しない)を経由してバス108へ接続されている。また、複数のコンピュータ・プログラムのコードには、本発明の実施形態に係る学習用音声データ生成プログラムが含まれる。また、各種データには、本発明においてソースドメインの音声データとするインターネット上の様々な音響環境の音声データや、ターゲットドメインの少量のクリーンな音声データを用いて学習されたターゲットドメインのGMMが含まれる。

【0027】

コンピュータ・プログラムは圧縮し、また複数に分割して複数の媒体に記録することもできる。なお、CPU102が、外部記憶装置114から渡されるデジタル信号に対して、学習用音声データ生成プログラムにより行う処理の詳細は後述する。

【0028】

コンピュータ100はまた、視覚データをユーザに提示するための表示装置116を含む。表示装置116は、グラフィックスコントローラ(図示しない)を経由してバス108へ接続されている。コンピュータ100は、通信インタフェース118を介してネットワークに接続し、他のコンピュータ等と通信を行うことが可能である。

【0029】

以上の説明により、コンピュータ100は、通常のパーソナルコンピュータ、ワークステーション、メインフレームなどの情報処理装置、又は、これらの組み合わせによって実現されることが容易に理解されるであろう。なお、上記説明した構成要素は例示であり、そのすべての構成要素が本発明の必須構成要素となるわけではない。同様に本発明を実施するためのコンピュータ100は、キーボードやマウスのような入力デバイス、スピーカ等他の構成要素を含むことも可能であることは言うまでもない。

【0030】

図2Bは、本発明の実施形態による学習用音声データ生成システム220の機能ブロック図を示す。学習用音声データ生成システム220は、ソースドメインの音声データをターゲットドメインの音声データにマッピングしてターゲットドメインの音声データを擬似的に生成するべく、ソースドメイン・データ格納部222と、混合ガウスモデル格納部224と、チャンネルマッピング部226と、ノイズデータ格納部228と、ノイズ付加部230と、ターゲットドメイン・データ格納部232とを備える。

【0031】

ソースドメイン・データ格納部222は、インターネット上で公開されている、あるいは

10

20

30

40

50

はインターネットを使った音声検索サービスなどのサービスにより提供される、種々の音響環境下における豊富な音声データをソースドメインの音声データとして格納する。したがって、ソースドメインの音声データには若干のノイズが重畳されている可能性が高い。このようなインターネットを使った音声データの収集は、例えば、携帯端末に音声認識を行うアプリケーションプログラムを実装し、その使用者の発声がインターネットを通じて送信され、サーバーに蓄積されることで行われる。

【0032】

混合ガウスモデル格納部224は、ターゲットドメインのクリーンな音声データを用いて学習されたターゲットドメインのGMMを格納する。即ち、ターゲットドメインのクリーンな音声データは、 k 混合のGMMとしてモデル化される。ターゲットドメインのGMM学習は、EMアルゴリズムを用いて行われる。なお、学習に用いられるターゲットドメインのクリーンな音声は少量であってよい（例えば、数名から数十名の男女話者によるトータル数時間の音声、精度に応じてより少ない時間の音声でもよい）。

10

【0033】

チャンネルマッピング部226は、ソースドメイン・データ格納部222から入力として読み出したソースドメインの音声データの1発話ごと、その補正量を混合ガウスモデル格納部224から読み出したターゲットドメインのGMMを参照して求め、ターゲットドメインの音声データにマッピングする。より具体的には、チャンネルマッピング部226は、ソースドメインの各音声データ（1発話）について、ターゲットドメインのGMMのガウス分布コンポーネントのうち、該音声データとの音響空間における距離が近いガウス分布コンポーネントの平均との差を求め、該差を、ガウス分布コンポーネントごとの尤度で重みづけした上で、時間方向の平均として求める。続いてチャンネルマッピング部226は、時間方向の平均として求めた量をチャンネル補正量としてソースドメインの音声データに足し合わせる。チャンネルマッピングの処理の更なる詳細は後述する。

20

【0034】

ノイズデータ格納部228は、ターゲットドメインにおけるノイズであって、マイクロフォンに音声以外の音源から混入し、音声データの波形に加法的な変形を加えるノイズを格納する。なお、ターゲットドメインにおけるノイズは、混合ガウスモデル格納部224に格納されるターゲットドメインのGMMの学習時と同じ条件（同一の音響環境）で記録する必要がある。

30

【0035】

ノイズ付加部230は、チャンネルマッピング部226から出力されるチャンネルマッピング後の音声データに、ノイズデータ格納部228からランダムに読み出したノイズを付加してターゲットドメインの音声データをシミュレートする。シミュレートされた擬似的なターゲットドメインの音声データはその後ターゲットドメイン・データ格納部232に格納される。なお、付加するノイズの割合は、最終的にターゲットドメインにおけるSN比の統計的分布に一致するよう1発話ごとに調整される。

【0036】

このように、音声データをターゲットドメインの音声データにマッピングするにあたり、チャンネル特性及びノイズ特性は主要な要素である。そして、上述した学習用音声データ生成システム220において、ノイズ付加の前にチャンネル補償を行うことが重要である。これを例証するために、2種類の音声データに同一の自動車のノイズを加えた。図3A、図3Bは、それぞれ、1発話ごとにCMN処理してチャンネル歪みによる影響を軽減した後のノイズの分布を示す。なお、縦軸は c_2 （2番目のケプストラム）を示し、横軸は c_1 （1番目のケプストラム）を示す（ただし、ケプストラムの最初の項は0番である）。

40

【0037】

図3Aは、停止した車中で録音された音声の特徴量ベクトルの平均値を0に正規化した後のノイズの分布を示す。図3Bは、携帯型機器に録音された音声の特徴ベクトルの平均値を0に正規化した後のノイズの分布を示す。同一のノイズを付加したにもかかわらず、ランタイム時のチャンネル正規化（CMN）処理により結果の信号は全く異なるものとなって

50

いる。この結果からも、デコーダにおいて類似した信号を取得するためには、最初にチャンネル補償を行う必要があることが分かる。

【 0 0 3 8 】

次に、チャンネルマッピング部 2 2 6 による処理の詳細を説明する。上述したように本発明ではターゲットドメインのチャンネル特性の事前知識として、ターゲットドメインの少量のクリーン音声データを用いて学習されたターゲットドメインのGMMを用いる。そしてターゲットドメインのGMMに似せるためのチャンネル補正量（以下、チャンネルマッピング・パラメータという）を大量の入力音声における1発話ごとに求めて、データのマッピングを行う。

【 0 0 3 9 】

ターゲットドメインのGMMは多少の話者バリエーションを含んでおり、これによって入力に話者依存要素が含まれることが許容される。入力は一切のノイズが重畳されていないクリーン音声为好適とされるが、実際の音声データには若干のノイズが含まれる。そこで、本発明では、VTS近似とEMアルゴリズムを用いてチャンネルマッピング・パラメータを推定する。以下、チャンネルマッピング・パラメータとして、特徴量のシフト量（以下、チャンネルバイアスという）のみを考慮する場合と、特徴量のシフト量に加えて、特徴量の大きさを変換する特徴量の係数（以下、チャンネル振幅という）をも考慮する場合の2つについて順に説明する。

【 0 0 4 0 】

<チャンネルバイアスのみの場合>

ここではマッピング・パラメータによってソースドメインの音声データをターゲットドメインの音声データにマッピングする生成モデル式において、チャンネルバイアスのみを考慮する場合について説明する。時間領域におけるソースドメインの観測音声の特徴量ベクトル Y_s は、チャンネルの特徴量ベクトル H_s 、クリーン音声の特徴量ベクトル X 、ノイズの特徴量ベクトル N を用いて下記式（1）のように表すことができる（添え字 s はソースドメインであることを示す）。なお N は、マイクロフォンに音声以外の音源から混入し波形に加法的な変形を与える雑音であり、 H は伝送系により加えられる乗法性の歪みである。

【数 1】

$$Y_s = H_s \cdot X + N \quad (1)$$

同様に、時間領域におけるターゲットドメインの観測音声の特徴量ベクトル Y_t は、チャンネルの特徴量ベクトル H_t 、クリーン音声の特徴量ベクトル X を用いて下記式（2）のように表すことができる（添え字 t はターゲットドメインであることを示す）。

【数 2】

$$Y_t = H_t \cdot X \quad (2)$$

【 0 0 4 1 】

上記式（1）をケプストラム領域に書き直すと下記式（3）、（4）になる。

【数 3】

$$y_s = h_s + x + G(x + h_s, n) \quad (3)$$

$$\text{但し、} G(x, n) = C \log(1 + \exp(C^{-1}(n - x))) \quad (4)$$

10

20

30

40

50

同様に上記式(2)をケプストラム領域に書き直すと下記式(5)になり、これを更に上記式(3)を用いて変形すると、最終的に下記式(8)、(9)が得られる。ここで式(9)より定義される c が求めるべきチャンネルバイアスである。ここで、行列 C は、Discrete Cosine Transform(DCT)行列を、 C^{-1} はその逆行列を表す。

【数4】

$$y_t = h_t + x - (5)$$

$$= h_t + (y_s - h_s - G(x + h_s, n)) - (6)$$

$$= (h_t - h_s) + y_s - G(x + h_s, n) - (7)$$

$$= y_s - c - G(y_t + c, n) - (8)$$

$$\text{但し、} c = h_s - h_t - (9)$$

【0042】

上記式(8)において、ターゲットドメインの観測音声の特徴量ベクトル y_t を

$$\hat{y}$$

、及びソースドメインの観測音声の特徴量ベクトル y_s を y と書き換えると、最終的に下記式(10)で表される生成モデル式が得られる。

【数5】

$$y = \hat{y} + c + G(\hat{y} + c, n) - (10)$$

上述したように、本発明では、チャンネルバイアス c をVTS近似とEMアルゴリズムを用いて推定する。具体的には、上記式(10)で表される生成モデル式を参照し、混合ガウスモデル格納部224に事前に用意したターゲットドメインのGMMから変換したソースドメインのGMMに含まれるチャンネルバイアス c を、EMアルゴリズムを用いて推定する。即ち、ターゲットドメインのGMMから変換したソースドメインのGMMにソースドメインの観測音声の特徴ベクトル y を入力して得られる音響尤度を求めるステップと、求めた音響尤度に基づいた目的関数を最小化するチャンネルバイアス c を求めるステップとを交互に繰り返す。以下、数式を用いてこの繰り返しステップを説明する。

【0043】

まず、ソースドメインのGMMにソースドメインの観測音声の特徴量ベクトル y を入力して得られる音響尤度 $p(y)$ の対数をとったものにマイナスを掛けたものを考え、下記式(11)を得る。ここで $\mu_{y,k}$ 、 $\sigma_{y,k}$ は順に、ソースドメインのGMMの k 番目の正規分布の事前確率、平均ベクトル、共分散行列を示す。式(11)から式(12)への変形は、イェンゼンの不等式を用いたものであり、更に変形を行って最終的に式(14)を得る。なお下記数式において d は、ソースドメインの観測音声の特徴量ベクトル y の d 番目の成分を示し、 D はその次元数を示す。

10

20

30

40

【数 6】

$$-\log(p(y)) = -\log\left(\sum_k \gamma_k \cdot \mathbf{N}(y; \mu_{y,k}, \Sigma_{y,k})\right) \quad (11)$$

$$\leq -\sum_k \gamma_k \cdot \log(\mathbf{N}(y; \mu_{y,k}, \Sigma_{y,k})) \quad (12)$$

$$= -\sum_k \gamma_k \cdot \log\left(\sqrt{2\pi}^{-D} \sqrt{|\Sigma_{y,k}|}^{-1} \exp\left(-\sum_d (y_d - \mu_{y,k,d})^2 / \Sigma_{y,k,d}\right)\right) \quad (13)$$

$$= \text{const} + \sum_k \gamma_k \left\{ \sum_d (y_d - \mu_{y,k,d})^2 / \Sigma_{y,k,d} + \log|\Sigma_{y,k}|^{\frac{1}{2}} \right\} \quad (14)$$

10

【0044】

チャネルバイアス c は音響尤度 $p(y)$ が最大になるように推定されるが、これは上記式(14)の右辺第2項が最小になるようにチャネルバイアス c を推定することに等しい。そこで、下記式(15)で表される目的関数を新たに導入する。

【数 7】

$$\Phi = E \left[\sum_k \rho_k(y) \cdot \left\{ \sum_d (y_d - \mu_{y,k,d})^2 / \Sigma_{y,k,d} + \log|\Sigma_{y,k}|^{\frac{1}{2}} \right\} \right] \quad (15)$$

20

ここで、事後確率は下記式(16)より定義される。

【数 8】

$$\rho_k(y) = \gamma_k \cdot \mathbf{N}(y; \mu_{y,k}, \Sigma_{y,k}) / \sum_{k'} \gamma_{k'} \cdot \mathbf{N}(y; \mu_{y,k'}, \Sigma_{y,k'}) \quad (16)$$

30

【0045】

チャネルバイアス c は、上記式(15)で表される目的関数を最小にするように推定される。そこで、事前に用意していたターゲットドメインのGMMと、上記式(10)の生成モデル式とから、共分散対角行列を仮定してVTS近似よりソースドメインのGMMを算出し、下記式(17)、(18)を得る。 μ_n 、 Σ_n は、それぞれノイズの平均ベクトル、共分散行列を示し、 $\delta_{d,l}$ はクロネッカーのデルタを示す。

【数 9】

$$\mu_{y,k} \cong \mu_{\hat{y},k} + c + G(\mu_{\hat{y},k} + c, \mu_n) \quad (17)$$

40

【数 10】

$$\Sigma_{y,k,d} \cong \sum_l (\delta_{d,l} - F(\mu_{\hat{y},k} + c, \mu_n)_{d,l})^2 \cdot \Sigma_{\hat{y},k,l} + \sum_l F(\mu_{\hat{y},k} + c, \mu_n)_{d,l}^2 \cdot \Sigma_{n,l} \quad (18)$$

なお、式(18)に表される共分散対角行列は、さらなる近似を用いて簡略化してもかまわない。例えば、ターゲットドメインの共分散行列と同一にしても、精度の劣化は少量に

50

とどまる。逆に、対角近似の条件を外して精密な共分散行列としても良い。上記式(17)に現れるベクトルG、及び、上記式(18)に現れるヤコビ行列Fは、それぞれ下記式(19)、(20)により定式化される雑音成分である。ここで、行列Cは、Discrete Cosine Transform(DCT)行列を、 C^{-1} はその逆行列を表す。

【数11】

$$G(x, n) = C \log(1 + \exp(C^{-1}(n - x))) - (19)$$

【数12】

10

$$F(x, n)_{i,j} = \sum_k C_{i,k} \cdot (C_{k,j}^{-1}) \cdot \exp\left(\sum_l C_{k,l}^{-1}(n_l - x_l)\right) / \left\{1 + \exp\left(\sum_l C_{k,l}^{-1}(n_l - x_l)\right)\right\} - (20)$$

【0046】

上記式(19)及び式(20)によりそれぞれ表される雑音成分のベクトルGとヤコビ行列Fは、ケプストラム領域において実装することが好ましいが、対数スペクトル領域において表すと、下記式(21)、(22)のようにそれぞれ表される。

20

【数13】

$$F(x, n) = \{1 + \exp(x - n)\}^{-1} - (21)$$

【数14】

$$G(x, n) = \log(1 + \exp(n - x)) - (22)$$

【0047】

30

次に上記式(17)～(20)を参照して、上記式(15)、(16)により表される目的関数 J をチャンネルバイアス c_d に関して微分したものを0に設定することで、チャンネルバイアス c の現在の推定値を得る。また、チャンネルバイアス c の現在の推定値を用いて雑音成分のベクトルG及びヤコビ行列Fを更新し、目的関数 J を求める。そして、求めた目的関数 J を最小にするように再びチャンネルバイアス c を推定する。この2ステップからなる処理をチャンネルバイアス c が収束するまで繰り返すことで、チャンネルバイアス c の最終的な推定値が得られる。

【0048】

収束したチャンネルバイアス c の推定値を生成モデル式の式(10)に代入し、更に式(10)の右辺第3項をMMSE推定により近似することで、最終的に下記式(23)が得られる。式(23)によりソースドメインの音声データからマップされたターゲットドメインのクリーン音声を得ることができる。

40

【数15】

$$\hat{y} = y - c - \sum_k^K \rho_k(y) \cdot G(\mu_{\hat{y},k} + c, \mu_n) - (23)$$

なお、式(23)では、ソースドメインのノイズ除去とチャンネル特性の補正を同時に行っている。ソースドメインのノイズを無視できる場合には、下記式(24)のように、ノイ

50

ズ除去を省略しても良い。

【数 1 6】

$$\hat{y} = y - c - (24)$$

【0049】

<チャンネル振幅を考慮する場合>

ここではマッピング・パラメータによってソースドメインの音声データをターゲットドメインの音声データにマッピングする生成モデル式において、チャンネルバイアス c に加えて、チャンネル振幅 a を新たに導入する場合について説明する。この場合、上記式 (10) の生成モデル式は、下記式 (25) のように拡張される。式中記号 $*$ は、ベクトル要素ごとの内積を表す。

【数 1 7】

$$y = a * \hat{y} + c + G(a * \hat{y} + c, n) - (25)$$

上述したように、本発明では、チャンネルバイアス c とチャンネル振幅 a をVTS近似とEMアルゴリズムを用いて推定する。具体的には、上記式 (25) で表される生成モデル式を参照し、混合ガウスモデル格納部 224 に事前に用意したターゲットドメインのGMMから変換したソースドメインのGMMに含まれるチャンネルバイアス c とチャンネル振幅 a を、EMアルゴリズムを用いて推定する。即ち、ターゲットドメインのGMMから変換したソースドメインのGMMにソースドメインの観測音声の特徴ベクトル y を入力して得られる音響尤度を求めるステップと、求めた音響尤度に基づいた目的関数を最小化するチャンネルバイアス c とチャンネル振幅 a を求めるステップとを交互に繰り返す。以下、数式を用いて説明する。

【0050】

上記式 (25) の生成モデル式を参照して、事前に用意していたターゲットドメインのGMMから、共分散対角行列を仮定してVTS近似より得られるソースドメインのGMMの平均ベクトル $\mu_{y,k}$ と共分散行列 $\Sigma_{y,k,d}$ は、それぞれ下記式 (26)、(27) のようになる。 μ_n 、 Σ_n は、それぞれノイズの平均ベクトル、共分散行列を示し、 δ はクロネッカーのデルタを示す。また式 (26)、(27) 中にそれぞれ現れる G と F は、上記式 (19) 及び式 (20) によりそれぞれ表される雑音成分のベクトル G とヤコビ行列 F である。

【数 1 8】

$$\mu_{y,k} \cong a * \mu_{\hat{y},k} + c + G(a * \mu_{\hat{y},k} + c, \mu_n) - (26)$$

【数 1 9】

$$\Sigma_{y,k,d} \cong \sum_l a_l^2 \cdot \left\{ \delta_{dl} - F(a * \mu_{\hat{y},k} + c, \mu_n)_{dl} \right\}^2 \cdot \Sigma_{\hat{y},k,l} + \sum_l F(a * \mu_{\hat{y},k} + c, \mu_n)_{dl}^2 \cdot \Sigma_{nl} - (27)$$

【0051】

なお、目的関数 J はチャンネルバイアス c のみを考慮する場合と同じであり上記式 (15)、(16) により表される。そこで上記式 (19) ~ (20)、(26) ~ (27) を参照して、上記式 (15)、(16) により表される目的関数 J をチャンネルバイアス c_d に関して微分したものを 0 に設定することで、チャンネルバイアス c の現在の推定値を得る。同様に上記式 (19) ~ (20)、(26) ~ (27) を参照して、上記式 (15)、(16) により表される目的関数 J をチャンネル振幅 a_d に関して微分したものを 0 に設定することで、チャンネル振幅 a の現在の推定値を得る。続いて、チャンネルバイアス c とチャ

10

20

30

40

50

ネル振幅 a のそれぞれの現在の推定値を用いて雑音成分のベクトル G 及びヤコビ行列 F を更新し、目的関数 を求める。そして、求めた目的関数 を最小にするように再びチャンネルバイアス c とチャンネル振幅 a とを推定する。この 2 ステップからなる処理をチャンネルバイアス c とチャンネル振幅 a とが収束するまで繰り返すことで、チャンネルバイアス c とチャンネル振幅 a の最終的な推定値が得られる。なお、チャンネルバイアス c とチャンネル振幅 a の初期値はそれぞれ順に値 0 と値 1 としてよい。また、チャンネルバイアス c については、通常その全ての成分を更新するが、チャンネル振幅 a については低次の成分、例えば 0 番目と 2 番目の成分のみを更新してもよい。これはケプストラムの低次の成分がチャンネル特性について支配的であるためである。

【 0 0 5 2 】

収束したチャンネルバイアス c 及びチャンネル振幅 a のそれぞれの推定値を生成モデル式の式 (2 5) に代入し、更に式 (2 5) の右辺第 3 項を M M S E 推定により近似して、最終的に下記式 (2 8) を得る。式 (2 8) によりソースドメインの音声データからマップされたターゲットドメインのクリーン音声を得ることができる。

【数 2 0】

$$\hat{y} = a^{-1} * \left(y - c - \sum_k^K \rho_k(y) \cdot G(a * \mu_{\hat{y},k} + c, \mu_n) \right) \quad (28)$$

10

20

式 (2 8) では、ソースドメインのノイズ除去とチャンネル特性の補正を同時に行っている。ソースドメインのノイズを無視できる場合には、下記式 (2 9) のように、ノイズ除去を省略しても良い。

【数 2 1】

$$\hat{y} = a^{-1} * (y - c) \quad (29)$$

【 0 0 5 3 】

< 話者の性別ごと用意された GMM を用いる場合 >

上記説明した 2 つのケースでは、混合ガウスモデル格納部 2 2 4 に事前に用意したターゲットドメインの GMM は 1 つであり、性別の区別なく多少の話者パリエーションを含んでいた。ここでは、ターゲットドメインの GMM は、話者の性別ごとに用意するものとする。そして、入力として受け取るソースドメインの音声データの 1 発話ごとに男性らしさ及び女性らしさを求めながら、ターゲットドメインの音声データにマッピングするものとする。すると目的関数 は下記式 (3 0) のように表される。

【数 2 2】

$$\Phi = E \left[\sum_d^D \sum_{g=f,m} \lambda_g \sum_k^K \rho_{g,k}(y) \cdot (y - \mu_{y,g,k})_d^2 / (\Sigma_{y,g,k})_d \right] \quad (30)$$

30

40

ここで g は gender インデックスであり、女性 (female) 又は男性 (male) のいずれかを示す。また、事後 gender 確率 λ_g は、ガウス分布の尤度を事後確率とみだてて合計を 1 になるように正規化したものである。また、 λ_g は上記式 (1 6) により定義される事後確率である。

【 0 0 5 4 】

ここで上記式 (1 0) の生成モデル式と上記式 (3 0) の目的関数 が与えられたとする。すると、上記式 (1 7) ~ (2 0) を参照して、上記式 (3 0) により表される目的関数 をチャンネルバイアス c_d に関して微分したものを 0 に設定することで、チャンネルバ

50

イアス c の現在の推定値を得る。また、チャンネルバイアス c の現在の推定値を用いて雑音成分のベクトル G 及びヤコビ行列 F を更新し、目的関数 を求める。そして、求めた目的関数 を最小にするように再びチャンネルバイアス c を推定する。この 2 ステップからなる処理をチャンネルバイアス c が収束するまで繰り返すことで、チャンネルバイアス c の最終的な推定値が得られる。

【 0 0 5 5 】

事後gender確率 λ_g もまた、チャンネルバイアス c を繰り返し更新する間、ターゲットドメインの性別ごとのGMMの事後確率に基づいて、下記式 (3 1)、(3 2) に示すように更新する。

【数 2 3】

$$\lambda'_g = E \left[\sum_k \gamma_{g,k} \cdot N(y; \mu_{y,g,k}, \Sigma_{y,g,k}) \right] - (31)$$

4/5

10

【数 2 4】

$$\lambda''_g = \lambda'_g / \sum_{g'} \lambda'_{g'} - (32)$$

20

オプションとして、下記式 (3 3) に示す事後gender確率 λ_g を用いてもよい。

【数 2 5】

$$\lambda_g = \exp(\beta \cdot \lambda''_g) / \sum_{g'} \exp(\beta \cdot \lambda''_{g'}) - (33)$$

但し β の値は定数である。なお、ここまで、ケプストラム領域を前提に定式化をしてきたが、よく行われるように対数メルスペクトル領域や対数スペクトル領域でも同様に定式化できることは当業者に明らかである (ケプストラム領域とは、対数メルスペクトル領域を離散コサイン変換したものである。)。

30

【 0 0 5 6 】

次に図 4 及び図 5 を参照して、本発明の実施形態による学習用音声データ生成処理の流れを説明する。図 4 は、学習用音声データ生成処理の全体の流れの一例を示すフローチャートある。図 5 は、図 4 に示すフローチャートのステップ 4 0 4 のマッピング処理の流れの一例を示すフローチャートある。

【 0 0 5 7 】

図 4 に示すフローチャートは、ステップ 4 0 0 で開始し、チャンネルマッピング部 2 2 6 は、ソースドメイン・データ格納部 2 2 2 からソースドメインの音声データを入力として取得する。続いてチャンネルマッピング部 2 2 6 は、混合ガウスモデル格納部 2 2 4 からターゲットドメインのGMMを読み出す (ステップ 4 0 2)。続いてチャンネルマッピング部 2 2 6 は、ターゲットドメインのGMMを参照して、ソースドメインの音声データを、ターゲットドメインの音声データに、ターゲットドメインの音声データのチャンネル特性に基づいてマッピングする (ステップ 4 0 4)。マッピング処理の詳細は図 5 を参照して後述する。

40

【 0 0 5 8 】

続いて、ノイズ付加部 2 3 0 は、ノイズデータ格納部 2 2 8 からランダムに読み出したノイズをマッピング後の音声データに付加してターゲットドメインの擬似音声データを生成し (ステップ 4 0 6)、これをターゲットドメイン・データ格納部に 2 3 2 へ出力する

50

(ステップ408)。そして処理は終了する。

【0059】

図5に示すフローチャートは、ステップ500で開始し、チャンネルマッピング部226はソースドメインの観測音声の特徴量ベクトルとターゲットドメインの観測音声の特徴量ベクトルとの関係を示す生成モデル式を決定する。上述したように利用可能な生成モデル式には、チャンネルマッピング・パラメータとしてチャンネルバイアス c のみを考慮する生成モデル式と、チャンネルマッピング・パラメータとしてチャンネルバイアス c とチャンネル振幅とを考慮する生成モデル式の2つがある。

【0060】

続いてチャンネルマッピング部226は、ターゲットドメインのGMMから変換したソースドメインのGMMに含まれるチャンネルマッピング・パラメータを、EMアルゴリズムを用いて推定する(ステップ502)。ここでソースドメインのGMMは、上述したように、事前に用意したターゲットドメインのGMMとステップ500で決定した生成モデル式とから、VTS近似により算出される。

10

【0061】

続いてチャンネルマッピング部226は、ステップ502で求めたチャンネルマッピング・パラメータを用いて、ステップ500で決定した生成モデル式に従いソースドメインの音声データをターゲットドメインの音声データにマッピングする(ステップ504)。そして処理は終了する。

【0062】

20

以上、実施形態を用いて本発明の説明をしたが、本発明の技術範囲は上記実施形態に記載の範囲には限定されない。上記の実施形態に、種々の変更又は改良を加えることが可能であることが当業者に明らかである。従って、そのような変更又は改良を加えた形態も当然に本発明の技術的範囲に含まれる。

【0063】

上述した実施形態の各機能は、C、C++、C#、Java(登録商標)などのオブジェクト指向プログラミング言語などで記述された装置実行可能なプログラムにより実現でき、本実施形態のプログラムは、ハードディスク装置、CD-ROM、MO、DVD、フレキシブルディスク、EEPROM、EPROMなどの装置可読な記録媒体に格納して頒布することができ、また他装置が可能な形式でネットワークを介して伝送することができる。

30

【実施例】

【0064】

以下、本発明について、実施例を用いてより具体的に説明を行なうが、本発明は、後述する実施例に限定されるものではない。

【0065】

先に開示した本発明の方法をコンピュータに実行させるためのコンピュータ・プログラムを作成し、一般社団法人情報処理学会(IPSJ)が提供する自動車内音声認識の評価用フレームワークを使用して、各コンピュータ・プログラムの性能を評価した。

40

【0066】

<実験条件>

本実験は、自動車内音声認識の評価用データベースCENSREC-3を用いて行った。評価条件は、学習データと評価データの双方に遠隔マイクロフォンを使用するミスマッチのない場合(A)と、学習データに接話マイクロフォンを使用し、評価データに遠隔マイクロフォンを使用するミスマッチのある場合(B~E)とした。更に、ミスマッチのある場合(B~E)は、本発明の適用の無い場合(B)と、本発明の適用のある場合(C~E)とし、本発明の適用のある場合は以下の3つの条件について検討した。学習データに対しチャンネルバイアスのみを考慮したチャンネルマッピングを適用した場合(C)。学習データに対し、チャンネルバイアスのみを考慮し、話者の性別を区別したチャンネルマッピングを適用した場合(D)。学習データに対し、チャンネルバイアスとチャンネル振幅を考慮し、話者の性

50

別を区別したチャンネルマッピングを適用した場合（E）。

【0067】

本発明を適用する場合、学習データとして、駐車した車内で記録された男性202人、女性91人の計293人のドライバーによる3608の発話を用いた。また、事前に用意するターゲットドメインのクリーン音声のGMMは、混合数は256とし、遠隔マイクロフォンで収録しランダムに選択した500の発話データで学習した。

【0068】

一方、評価データとしては、駐車した車内で記録された男性8人、女性10人の計18人のドライバーによる898の発話を用いた。

【0069】

また、実験に必要な様々な特徴量を出力するフロントエンドを用意し、学習データと評価データの双方に適用した。特徴量は、MFCC 12次元 + MFCC 12次元 + MFCC12次元 + 対数パワー の39次元で、発話単位のCMNを適用した場合と適用しない場合の両方の値を得た。音響モデルの作り方などバックエンドの構成は無変更とした(Category0)。

【0070】

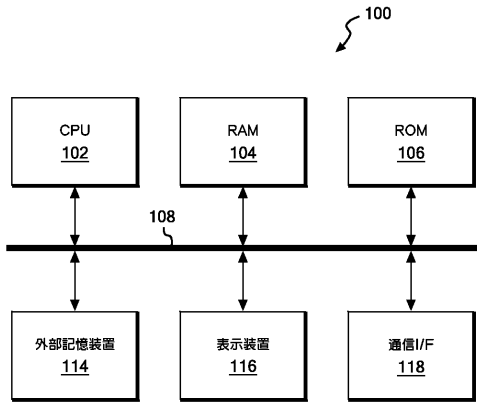
<実験結果>

図6は、上述したA～Eの6つの条件下それぞれでのCMNオンとCMNオフでの単語正解精度（％）を示す。条件Aはミスマッチの無い場合の単語正解精度であるため、その値は上限を示す。条件Bはミスマッチのある場合であって、本発明によるチャンネルマッピングがなされていないため、その値はベースラインとなる。条件C～条件Eは、ミスマッチのある場合であって、本発明によるチャンネルマッピングが適用されている場合である。いずれの場合も単語正解精度はベースラインよりも高くなっている。また、条件Dでの単語正解精度の値は条件Cのそれよりも改善されているため、本発明によるチャンネルマッピング手法は、チャンネル特性と話者特性の両方に有効であるといえる。また、条件C～条件EのすべてにおいてCMNオンの場合の単語正解精度の値はベースラインのそれよりも高くなっているため、本発明によるチャンネルマッピング手法はCMNと相性がよいといえる。

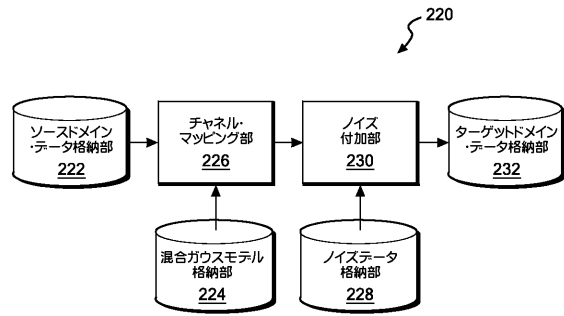
10

20

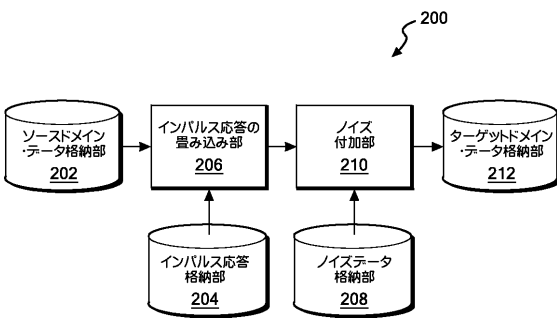
【図 1】



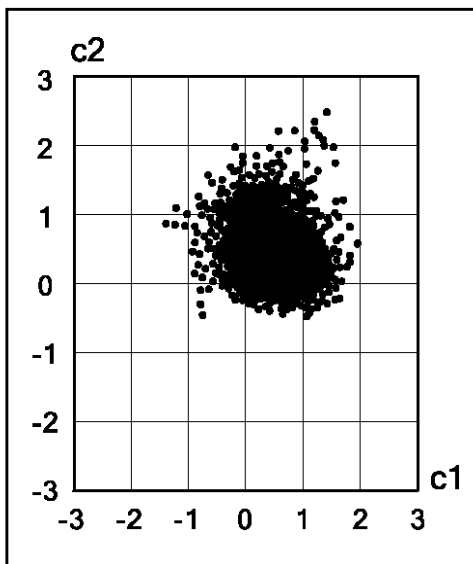
【図 2 B】



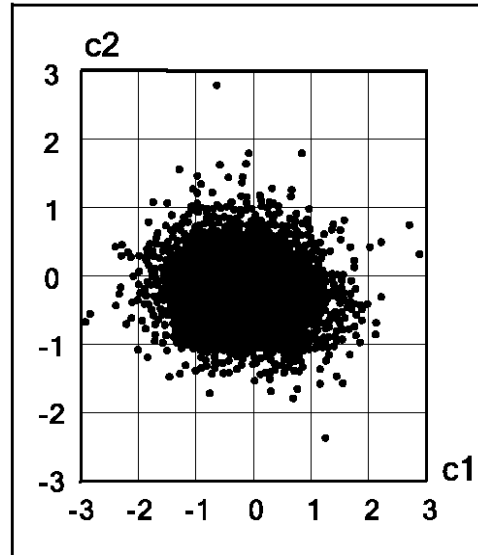
【図 2 A】



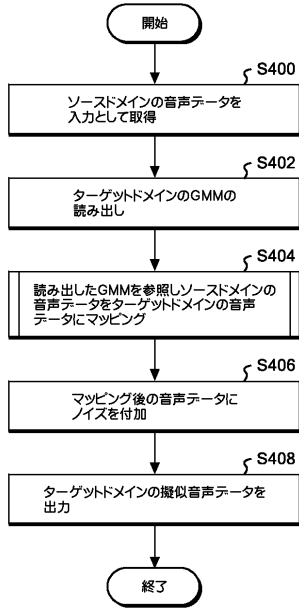
【図 3 A】



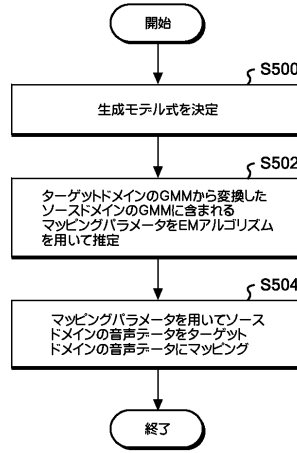
【図 3 B】



【 図 4 】



【 図 5 】



【 図 6 】

	学習データ	評価データ	CMN off	CMN on
A	遠隔マイク、アイドリング	遠隔マイク、アイドリング (ファン・オフ、窓閉)	99.7	99.7
B	接話マイク、アイドリング	遠隔マイク、アイドリング (ファン・オフ、窓閉)	92.3	95.2
C	接話マイク、アイドリング、 チャンネルマッピング (バイアスのみ)あり	遠隔マイク、アイドリング (ファン・オフ、窓閉)	93.1	96.0
D	接話マイク、アイドリング、 チャンネルマッピング (性別依存、バイアスのみ)あり	遠隔マイク、アイドリング (ファン・オフ、窓閉)	94.5	96.8
E	接話マイク、アイドリング、 チャンネルマッピング (性別依存、バイアス及び振幅)あり	遠隔マイク、アイドリング (ファン・オフ、窓閉)	96.4	98.0

フロントページの続き

(74)代理人 100112690

弁理士 太佐 種一

(72)発明者 市川 治

東京都江東区豊洲五丁目6番52号 NBF豊洲キャナルフロント 日本アイ・ビー・エム株式会社 東京基礎研究所内

(72)発明者 スティーブン・ジェイ・レニー

アメリカ合衆国10598ニューヨーク州ヨークタウン・ハイツルート134 キチャワン・ロード1101番地

審査官 鈴木 圭一郎

(56)参考文献 特開2008-026489(JP,A)

特開平05-073088(JP,A)

特開2005-196020(JP,A)

(58)調査した分野(Int.Cl., DB名)

G10L 15/00