

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4274962号
(P4274962)

(45) 発行日 平成21年6月10日(2009.6.10)

(24) 登録日 平成21年3月13日(2009.3.13)

(51) Int.Cl.		F I			
G 1 0 L	15/06	(2006.01)	G 1 0 L	15/06	3 1 0 T
G 1 0 L	15/14	(2006.01)	G 1 0 L	15/06	4 0 0 V
G 1 0 L	15/28	(2006.01)	G 1 0 L	15/14	2 0 0 C
			G 1 0 L	15/28	3 6 0 C
			G 1 0 L	15/28	3 7 0 Z

請求項の数 13 (全 26 頁)

(21) 出願番号	特願2004-28542 (P2004-28542)	(73) 特許権者	393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2
(22) 出願日	平成16年2月4日(2004.2.4)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2005-221678 (P2005-221678A)	(72) 発明者	松田 繁樹 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
(43) 公開日	平成17年8月18日(2005.8.18)	(72) 発明者	實廣 貴敏 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
審査請求日	平成17年6月13日(2005.6.13)	(72) 発明者	コンスタンティン・マルコフ 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
(出願人による申告)平成15年度通信・放送機構、研究テーマ「大規模コーパス音声対話翻訳技術の研究開発」に関する委託研究、産業活力再生特別措置法第30条の適用を受ける特許出願			
前置審査		最終頁に続く	

(54) 【発明の名称】 音声認識システム

(57) 【特許請求の範囲】

【請求項1】

それぞれ所定の音響特徴量をパラメータとする複数の音響モデル群を記憶するための記憶手段を含む音声認識システムであって、前記複数の音響モデル群の各々は、それぞれ異なる発話環境での発話音声のデコードに最適化された、同種の複数の音響モデルを含み、

前記音声認識システムはさらに、

入力される音声から前記所定の音響特徴量を算出するための特徴量算出手段と、

前記入力される音声のうち、発話がない部分の前記音響特徴量に基づいて、前記複数の音響モデル群の各々に対して、それぞれ前記複数の音響モデル群の混合重み適応化により、前記入力される音声の発話環境に適応化された複数の適応化音響モデルを作成するためのモデル適応化手段と、

前記複数の音響モデル群の各々に対して設けられ、前記入力される音声の発話部分に回答し、前記複数の適応化音響モデルをそれぞれ用いて、前記入力される音声の前記発話部分の前記所定の音響特徴量をデコードし、複数の仮説を出力するための複数のデコード手段と、

前記複数のデコード手段が出力する前記複数の仮説を統合することにより音声認識結果を出力するための統合手段とを含み、

前記統合手段は、

前記複数の音響モデル群の各々に対し、前記複数のデコード手段により得られる前記複数の仮説から、各単語の音響言語尤度の和が最大となる仮説を選択するための仮説選択手

段と、

前記仮説選択手段によって前記複数の音響モデル群の各々に対して選択された仮説から、単語ラティスを作成するためのラティス作成手段と、

前記ラティス内の単語列の中で、単語の音響尤度とN - グラム単語列（Nは1以上の整数）の言語尤度とから算出される値が所定条件を満足するものを認識結果として選択するための単語列選択手段とを含む、音声認識システム。

【請求項2】

前記複数の音響モデル群の各々は、それぞれ異なる雑音が重畳された発話音声のデコードに最適化された複数の音響モデルを含む、請求項1に記載の音声認識システム。

【請求項3】

前記モデル適応化手段は、

前記入力される音声の前記音響特徴量に基づいて、前記複数の音響モデル群の各々について、当該音響モデル群に含まれる複数の音響モデルのうち、前記入力される音声の前記音響特徴量に関連する予め定める条件を充足する、所定個数の音響モデルを選択するための手段と、

前記選択するための手段により前記複数の音響モデル群の各々について選択された前記所定個数の音響モデルから、雑音ガウス混合分布の混合重み適応化手法により前記適応化音響モデルを作成するための手段とを含む、請求項2に記載の音声認識システム。

【請求項4】

前記複数の音響モデル群は、

互いに異なる複数種類の雑音が第1のSNR（信号対雑音比）で重畳された発話音声のデコードに最適化された複数の音響モデルを含む第1の音響モデル群と、

前記複数種類の雑音が、前記第1のSNRと異なる第2のSNRで重畳された発話音声のデコードに最適化された複数の音響モデルを含む第2の音響モデル群とを含む、請求項1に記載の音声認識システム。

【請求項5】

前記複数の音響モデル群は、

それぞれ異なる発話環境での発話音声のデコードに最適化された、第1の音響モデル構造に基づく第1の音響モデル群と、

それぞれ異なる発話環境での発話音声のデコードに最適化された、前記第2の音響モデル構造とは異なる第2の音響モデル構造に基づく第2の音響モデル群とを含む、請求項1に記載の音声認識システム。

【請求項6】

前記第1の音響モデル構造は、通常発話に対して想定される音響モデル構造である、請求項5に記載の音声認識システム。

【請求項7】

前記第2の音響モデル構造は、言直し発話に対して想定される音響モデル構造である、請求項5又は請求項6に記載の音声認識システム

【請求項8】

前記複数の音響モデル群は、

それぞれ異なる発話環境での発話音声のデコードに最適化された、第1の種類の音響特徴量をパラメータとする第1の音響モデル群と、

それぞれ異なる発話環境での発話音声のデコードに最適化された、前記第1の種類の音響特徴量と異なる第2の音響特徴量をパラメータとする第2の音響モデル群とを含む、請求項1に記載の音声認識システム。

【請求項9】

前記第1の種類の音響特徴量はM F C C（メル周波数ケプストラム係数）である、請求項8に記載の音声認識システム。

【請求項10】

前記第2の種類の音響特徴量はD M F C C（差分メル周波数ケプストラム係数）である、

10

20

30

40

50

請求項 8 又は請求項 9 に記載の音声認識システム。

【請求項 1 1】

前記モデル適応化手段は、

前記入力される音声の前記音響特徴量に基づいて、前記第 1 の音響モデル群に含まれる音響モデルのうち、前記入力される音声の前記音響特徴量に関連する予め定める条件を充足する、所定個数の音響モデルを選択するための第 1 の音響モデル選択手段と、

前記第 1 の音響モデル選択手段により選択された音響モデルから、雑音 GMM の混合重み適応化手法により第 1 の適応化音響モデルを作成するための手段と、

前記入力される音声の前記音響特徴量に基づいて、前記第 2 の音響モデル群に含まれる音響モデルのうち、前記入力される音声の前記音響特徴量に関連する予め定める条件を充足する、所定個数の音響モデルを選択するための第 2 の音響モデル選択手段と、

前記第 2 の音響モデル選択手段により選択された音響モデルから、雑音 GMM の混合重み適応化手法により第 2 の適応化音響モデルを作成するための手段とを含む、請求項 4 ~ 請求項 10 のいずれかに記載の音声認識システム。

【請求項 1 2】

前記単語列選択手段は、前記ラティス内の単語列の中で、前記算出される値が最大となるものを認識結果として選択するための手段を含む、請求項 1 ~ 請求項 11 のいずれかに記載の音声認識システム。

【請求項 1 3】

前記単語列選択手段は、

前記ラティス内の単語の音響尤度と、N - グラム単語列の言語尤度とを、それぞれ所定の正規化方式により正規化するための正規化手段と、

前記ラティス内の単語列ごとに、前記正規化手段により正規化された当該単語列内の単語の音響尤度と前記 N - グラム単語列の言語尤度とにそれぞれ所定の重みを加算して得られる値が前記所定条件を満足するものを認識結果として選択するための手段とを含む、請求項 1 ~ 請求項 11 のいずれかに記載の音声認識システム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は大語彙の連続音声認識装置及び方法に関し、特に、雑音に強く、発話スタイルの変動に対しても頑健に音声を認識することが可能な連続音声認識システムに関する。

【背景技術】

【0002】

近年、雑音又は発話スタイルに対して頑健な音声認識の研究が盛んに行なわれている。実環境において音声認識を使用するためには、通行する自動車などの乗り物から発せられるエンジン雑音や風切り音、駅、オフィス内などの人の声、コンピュータからのファンの音など、多種多様な雑音環境において高精度な音声認識が実現されなければならない。

【0003】

さらに雑音だけでなく、使用者の年齢や性別、また感情や体調によってその発話スタイルは刻一刻と変化する。音声認識装置は、そのような発話スタイルの変動に対しても雑音と同様に頑健でなければならない。

【0004】

雑音又は発話スタイルなど個別の変動に対する頑健化手法が従来から数多く提案されてきた。これについては後掲の非特許文献 1 を参照されたい。本明細書では以下、音声の音響的言語的特徴に影響する要因のことを総じて「発話環境」と呼ぶこととする。

【0005】

雑音に対して頑健な音響特徴量の分析手法として、「SS (Spectrum Subtraction) 法 (後掲の非特許文献 2 を参照されたい)」を音声認識の前処理として用いる手法が提案されている。これ以外にも、RASTA (Relative Spectra)、DMFCC (Differential Mel Frequency C

10

20

30

40

50

epstrum Coefficient) など、いくつかの音響分析手法が提案されている。

【0006】

SS法では、雑音重畳音声のスペクトルに対して雑音スペクトルを減算することにより、SNR(信号対雑音比)を改善している。RASTA法では、個々の周波数バンドの値の変化に対して、音声情報が多く含まれている1から12Hzの変調スペクトラム成分を抽出することにより雑音の影響を軽減している。またDMFCCはFFT(高速フーリエ変換)によって得られるフーリエ係数に対して、隣り合う係数間で差分をとり、音声などのピッチを持つスペクトルを強調することによって耐雑音性を改善している。

【0007】

雑音に頑健な音響モデルの研究としては、PMC(Parallel Model Combination)法(後掲の非特許文献5を参照されたい。)、ヤコビ適応法(後掲の非特許文献6を参照されたい。)、MLLR(Maximum Likelihood Linear Regression)(後掲の非特許文献7を参照されたい。)による雑音適応などが提案されている。

【0008】

これらのうち、PMC法は、HMM(隠れマルコフモデル)の出力確率分布を線形スペクトル領域に変換し雑音スペクトルを重畳することにより、環境雑音への適応を行なう手法である。このPMC法につき簡単に説明する。

【0009】

PMC法の概念を図28を参照して説明する。図28を参照して、PMC法の対象となるもとの音響モデルが、音響の特徴量からなる音響空間600において領域610の付近に存在する音響をモデル化したものであるものとする。このとき、音声認識対象の雑音を含んだ音声データ領域612は、雑音のためにもとの領域610からずれたものとなる。そこで、領域612と領域610との差分を考え、この差分に相当する量を音響モデル610に加えることにより音響モデルの音響空間600内における位置を領域612まで移動するよう音響モデルを変換する。

【0010】

このようにして変換した後の音響モデルを用いれば、領域612の付近に存在する雑音を含んだ音声については、もとの音響モデルを用いたものより高い精度で認識できる。

【0011】

ヤコビ適応法は、雑音の変化に伴う出力確率分布の非線形変換を線形近似することにより、雑音環境へ高速に適応する手法である。

【0012】

MLLRを用いた雑音適応は、無雑音音声と雑音重畳音声との間の分布移動を回帰行列を用いて表現し、音響モデル全体を雑音モデルに適応化する手法である。

【0013】

さらに、雑音の分布の時間変動を逐次的に推定することにより、非定常雑音に対する認識精度を改善する手法(後掲の非特許文献9を参照されたい。)が提案されている。

【0014】

発話スタイルに対する頑健性の改善手法としては、発話スタイル依存の音響モデルを用いる手法の他、ロンバード効果によるスペクトルの変形を考慮した手法(非特許文献8を参照されたい。)及び個々の母音HMMの最後に無音状態を追加することにより音声強調発声や言直し発話に頑健な音響モデルを構築する手法(非特許文献10を参照されたい。)などが提案されている。そのほかにも、講演音声などの音素継続時間の短い発声を含む音声に対して、分析フレーム周期又はウィンドウ幅を自動選択することにより認識精度を改善する手法(非特許文献11、12参照)が提案されている。

【0015】

これらの頑健化手法は主として、雑音や発話スタイルなどの個別の変動に対する頑健化である。音声認識を実環境で用いるためには、複数の発話環境が刻一刻と変化する状況で

10

20

30

40

50

あっても頑健に音声を認識することができなければならない。このような種々の外乱に対して頑健な音声認識を実現するための方法は大きく2つに分類することができると考えられる。発話環境の変動に頑健な音響モデル及び言語モデルを用いて単数のデコーダで認識を行なうシングルタイプの方法と、お互いに異なる環境に適応化された複数の音響モデル及び言語モデルを使用して得られた複数の仮説を統合するパラレルタイプの手法とである。

【0016】

シングルタイプの音声認識システムを構築するためには、広い発話環境の音声を頑健に認識する音響モデル及び言語モデルが必要である。そのために、男性及び女性双方の学習データから性別独立な音響モデルを推定するなど、複数の発話環境のデータを用いてHMMのモデルパラメータ推定を行なうことにより頑健性を改善する手法がある。しかし、男性女性などのお互いの音響的特徴が大きく異なる場合ではなく、種々のSNRのデータを用いて学習する場合、個々の音素モデルの分布が過度に広がることにより音素分類精度の低下が懸念される。従って、このようなモデル化法には頑健化の限界があると考えられる。

10

【0017】

セグメントモデル（非特許文献13参照）では、時間的に離れた音響特徴ベクトル間の相関を計算することで音声の非定常な振舞いのモデル化を試みている。時間的に離れた特徴ベクトル間の相関として発話環境の変動をモデル化することができるならば、セグメントモデルにおいて広い発話環境の音声を頑健に認識できる可能性がある。しかし、効率的な相関の計算方法やモデルパラメータの増大などの問題により十分な精度は得られていない。

20

【0018】

一方、パラレルタイプによる音声認識は、個々の音響モデルや言語モデルの利用可能な発話環境が限られていたとしても、それらを複数個使用しパラレルにデコーディングすることにより、個々の音素間の分類精度を低下させることなく広い発話環境の音声を頑健に認識できる可能性がある。

【0019】

このような音声認識の例としては、SNRに依存した音響モデルを用いて得られた複数の仮説を最大尤度基準で選択する手法、複数のお互いに異なる音響特徴量を用いて音声認識を行ない、得られた複数の仮説を単語単位で統合する仮説統合法（非特許文献15参照）が提案されている。

30

【0020】

【非特許文献1】中村、『実音響環境に頑健な音声認識を目指して』、信学技報、EA2002-12、pp.31-36、2002。

【非特許文献2】S.F.ボル、『スペクトル減算を用いた音声中の音響雑音の抑制』、IEEE音響音声信号処理論文集、第ASSP-27巻、第113-120頁、1979年。（S.F. Boll, 『Suppression of Acoustic Noise in Speech Using Spectral Subtraction,』IEEE Trans. Acoust., Speech, Signal Processing, vol. ASSP-27, pp. 113-120, 1979.）

40

【非特許文献3】H.ヘルマンスキ及びN.モーガン、「音声のRASTA処理」、IEEE音声及び音響処理トランザクション、第2巻、第4号、第587-589頁（H. Hermansky and N. Morgan, 『RASTA Processing of Speech,』IEEE Trans. Speech and Audio Processing, vol. 2, no. 4, pp. 587-589, 1994.）

【非特許文献4】J.チェン、K.K.パリワル、S.ナカムラ、『頑健な音声認識のための差分パワースペクトル由来のケプストラム』、音声コミュニケーション、第41巻第2-3号、第469-484頁、2003年。（J. Chen, K.K. Paliwal, S. Nakamura, 『Cepstrum Derived from Diffe

50

rentiated Power Spectrum for Robust Speech Recognition,」Speech Communication, vol. 41, no. 2-3, pp. 469-484, 2003.)

【非特許文献5】M. ゲールズ及びS. ヤング、『パラレルモデルの組合せを用いた頑健な連続音声認識』、IEEE音声及び音響処理論文集、第4巻、第5号、第352-359頁、1996年。(M. Gales and S. Young, 『Robust Continuous Speech Recognition Using Parallel Model Combination,』IEEE Trans. on Speech and Audio Processing, vol. 4, No. 5, pp. 352-359, 1996.)

10

【非特許文献6】Y. ヤマグチ、S. タカハシ及びS. サガヤマ、『ヤコビアン適応アルゴリズムを用いた環境雑音への音響モデルの高速適応』、ユーロスピーチ予稿集、97、第2051-2054頁、1997年。(Y. Yamaguchi, S. Takahashi and S. Sagayama, 『Fast Adaptation of Acoustic Models to Environmental Noise Using Jacobian Adaptation Algorithm,』Proc. Eurospeech, 97, pp. 2051-2054, 1997.)

【非特許文献7】C. J. レゲッタ及びP. C. ウッドランド、『連続密度隠れマルコフモデルの話者適応のための最大尤度線形回帰』、コンピュータ音声及び言語、第9巻、第171-185頁、1995年。(C. J. Leggetter and P. C. Woodland, 『Maximum Likelihood Linear Regression for Speaker Adaptation of Continuous Density Hidden Markov Models,』Computer Speech and Language, vol. 9, pp. 171-185, 1995.)

20

【非特許文献8】J. C. ジャンカ、『ロンバード効果とその聴者及び自動音声認識装置に対する役割』、アメリカ音響学会誌、第93巻、第510-524頁、1993年。(J. C. Junqua, 『The Lombard Reflex and its Role on Human Listeners and Automatic Speech Recognizer,』J. Acoustic Soc. Amer., vol. 93, pp. 510-524, 1993.)

30

【非特許文献9】K. ヤオ、B. E. シー、S. ナカムラ及びZ. カオ、『非定常雑音における頑健な音声認識のための連続EMアルゴリズムによる残存雑音の補償』、ICSLP2000予稿集、第1巻、第770-773頁、2000年。(K. Yao, B. E. Shi, S. Nakamura and Z. Cao, 『Residual Noise Compensation by a Sequential EM Algorithm for Robust Speech Recognition in Nonstationary Noise,』Proc. ICSLP2000, vol. 1, pp. 770-773, 2000.)

【非特許文献10】奥田、松井、中村、『誤認識時の言い直し発話における発話スタイルの変動に頑健な音響モデル構築法』信学論、vol. J86-DII, no. 1, pp. 42-51, 2003.

40

【非特許文献11】奥田、河原、中村、『ゆう度基準による分析周期・窓長の自動選択手法を用いた発話速度の補正と音響モデルの構築』信学論、vol. J86-DII, no. 2, pp. 204-211, 2003.

【非特許文献12】南條、河原、『発話速度に依存したデコーディングと音響モデルの適応』信学技報、SP2001-103, 2001.

【非特許文献13】M. オステンドルフ、V. ディガラキス及びO. キンバル、『HMMからセグメントモデルへ：音声認識のためのストカスティックモデリングの統一見解』、IEEE音声及び音響処理論文集、第4巻、第5号、第360-378頁、1996年。

50

(M. Ostendorf, V. Digalakis and O. Kimball, 『From HMMs to Segment Models: A Unified View of Stochastic Modeling for Speech Recognition,』IEEE Trans. Speech and Audio Proc., vol. 4, no. 5, pp. 360 - 378, 1996.)

【非特許文献14】伊田、中村、『雑音GMMの適応化とSN比別マルチパスモデルを用いたHMM合成による高速な雑音環境適応化』信学論、vol. J86-D-II, no. 2, pp. 195 - 203, 2003.

【非特許文献15】K. マルコフ、T. マツイ、R. グルーン、J. ツァン、S. ナカムラ、『DARPA SPINE2用の雑音及びチャネル歪に頑健なASRシステム』、IEICE情報及びシステム論文集、第E86-D巻、第3号、2003年。(K. Markov, T. Matsui, R. Gruhn, J. Zhang, S. Nakamura, 『Noise and Channel Distortion Robust ASR System for DARPA SPINE2 Task,』IEICE Trans. Inf. & Syst., vol. E86-D, no. 3, 2003.)

【非特許文献16】M. オステンドルフ及びH. シンガー、『最大尤度連続状態分割を用いたHMMトポロジー設計』、コンピュータ音声及び言語、第11巻、第1号、第17 - 41頁1997年。(M. Ostendorf and H. Singer, 『HMM Topology Design Using Maximum Likelihood Successive State Splitting,』Computer Speech and Language, vol. 11, no. 1, pp. 17 - 41, 1997.)

【発明の開示】

【発明が解決しようとする課題】

【0021】

しかしながら上述したいずれの方法においても、例えばモデルの変換に時間を要すること、雑音又は発話スタイルなど、個別の要素の変動に的確に対応することが難しいこと、などから、実環境における雑音を含んだ音声や、発話スタイルが変動する音声に対して音声認識を精度よく行なうことは未だ可能でない。

【0022】

それゆえにこの発明の目的は、雑音などの個別の変動に実時間で追従して精度高く認識することができる音声認識システムを提供することである。

【0023】

この発明の他の目的は、雑音などの個別の変動だけでなく、発話スタイルの変動に対しても頑健に音声認識することができる音声認識システムを提供することである。

【課題を解決するための手段】

【0024】

本発明に係る音声認識システムは、それぞれ異なる発話環境での発話音声のデコードに最適化された、それぞれ所定の音響特徴量をパラメータとする複数の音響モデル群を記憶するための記憶手段と、入力される音声から所定の音響特徴量を算出するための特徴量算出手段と、入力される音声の音響特徴量に基づいて、それぞれ複数の音響モデル群の混合重み適応化により、入力される音声の発話環境に適応化された複数の適応化音響モデルを作成するためのモデル適応化手段と、複数の適応化音響モデルを用いて、入力される音声の所定の音響特徴量をデコードし複数の仮説を出力するためのデコード手段と、デコード手段が出力する複数の仮説を統合することにより音声認識結果を出力するための統合手段とを含む。

【0025】

デコード手段により出力される複数の仮説が互いに相補的である場合、統合手段により仮説を統合することにより、より精度の高い音声認識結果を得られる可能性が高い。

【0026】

好ましくは、複数の音響モデル群の各々は、それぞれ異なる雑音が重畳された発話音声のデコードに最適化された複数の音響モデルを含む。

【0027】

複数の音響モデル群の各々が含む音響モデルは、互いに異なる雑音が重畳された発話音声のデコードに適している。種々の雑音について適応化音響モデルが得られるので、雑音が異なる発話環境での音声認識の精度が向上することが期待できる。

【0028】

より好ましくは、モデル適応化手段は、入力される音声の音響特徴量に基づいて、複数の音響モデル群の各々について、当該音響モデル群に含まれる複数の音響モデルのうち、入力される音声の音響特徴量に関連する予め定める条件を充足する、所定個数の音響モデルを選択するための手段と、選択するための手段により複数の音響モデル群の各々について選択された所定個数の音響モデルから、雑音ガウス混合分布の混合重み適応化手法により適応化音響モデルを作成するための手段とを含む。

【0029】

モデル適応化にあたって、入力音声の発話環境と所定のある発話環境に対応する音響モデルを複数個選び、それらに対し混合重み適応化を行なって適応化環境モデルを作成する。適応化モデル作成時の計算量が少なく済み、またモデルの精度は十分に高くなる。

【0030】

複数の音響モデル群は、互いに異なる複数種類の雑音が第1のSNRで重畳された発話音声のデコードに最適化された複数の音響モデルを含む第1の音響モデル群と、複数種類の雑音が、第1のSNRと異なる第2のSNRで重畳された発話音声のデコードに最適化された複数の音響モデルを含む第2の音響モデル群とを含んでもよい。

【0031】

種々のSNRに最適化された音響モデルに基づいて、それぞれ適応化モデルが作成される。より広い発話環境に対し適応化モデルが作成されるので、入力音声の発話環境に近いものが得られる可能性が高くなる。その結果、音声認識精度の向上が期待できる。

【0032】

好ましくは、複数の音響モデル群は、それぞれ異なる発話環境での発話音声のデコードに最適化された、第1の音響モデル構造に基づく第1の音響モデル群と、それぞれ異なる発話環境での発話音声のデコードに最適化された、第2の音響モデル構造とは異なる第2の音響モデル構造に基づく第2の音響モデル群とを含む。

【0033】

第1及び第2の音響モデル構造に基づく音響モデル群を用いることにより、雑音以外の要因の変化に対しても頑健に音声認識を行なうことができる。

【0034】

好ましくは、第1の音響モデル構造は、通常発声に対して想定される音響モデル構造であり、さらに好ましくは第2の音響モデル構造は、言直し発話に対して想定される音響モデル構造である。

通常発声と言直し発話という二種類の発話に対応する音響モデル構造を用いることにより、話者の発話スタイルの変化に対しても頑健に音声認識を行なうことができる。

【0035】

複数の音響モデル群は、それぞれ異なる発話環境での発話音声のデコードに最適化された、第1の種類の音響特徴量をパラメータとする第1の音響モデル群と、それぞれ異なる発話環境での発話音声のデコードに最適化された、第1の種類の音響特徴量と異なる第2の音響特徴量をパラメータとする第2の音響モデル群とを含んでもよい。

【0036】

好ましくは、第1の種類の音響特徴量はMFCCであり、さらに好ましくは、第2の種類の音響特徴量はDMFCCである。

【0037】

10

20

30

40

50

第1及び第2の音響特徴量を用いる音響モデル群を用いることにより、種々発話環境の変化に対しても頑健に音声認識を行なうことができる。

【0038】

デコード手段は、第1及び第2の適応化音響モデルの各々に対し、入力される音声の所定の音響特徴量をデコードし複数の仮説を出力するための手段を含んでもよい。

【0039】

好ましくは、デコード手段はさらに、第1及び第2の適応化音響モデルの各々について、出力するための手段により出力された複数の仮説のうち、各単語の音響言語尤度の和が最大となる仮説を選択し、統合手段に与えるための手段を含む。

【0040】

統合手段に与える前に、適応化音響モデルごとに音響言語尤度の和が最大となる仮説を選択することにより、統合処理の際の探索空間が小さくなり統合処理が容易かつ高速になる。

【0041】

さらに好ましくは、適応化音響モデルを作成するための手段は、入力される音声の音響特徴量に基づいて、第1の音響モデル群に含まれる音響モデルのうち、入力される音声の音響特徴量に関連する予め定める条件を充足する、所定個数の音響モデルを選択するための第1の音響モデル選択手段と、第1の音響モデル選択手段により選択された音響モデルから、雑音GMM (Gaussian Mixture Model) の混合重み適応化手法により第1の適応化音響モデルを作成するための手段と、入力される音声の音響特徴量に基づいて、第2の音響モデル群に含まれる音響モデルのうち、入力される音声の音響特徴量に関連する予め定める条件を充足する、所定個数の音響モデルを選択するための第2の音響モデル選択手段と、第2の音響モデル選択手段により選択された音響モデルから、雑音GMMの混合重み適応化手法により第2の適応化音響モデルを作成するための手段とを含む。

【0042】

複数の音響モデル群からそれぞれ適応化音響モデルが作成され、それらを用いてデコードが行われ、かつそれらのデコード結果が統合されるので、それぞれの音響モデルの特徴群の特長を組み合わせた頑健な音声認識を行なうことができる。

【0043】

より好ましくは、仮説統合手段は、デコード手段が出力する複数の仮説から、単語ラティスを作成するためのラティス作成手段と、ラティス内の単語列の中で、単語の音響尤度とN-グラム単語列 (Nは1以上の整数) の言語尤度とから算出される尤度が所定条件を満足するもの、例えば最大となるもの、を認識結果として選択するための選択手段とを含む。

【0044】

さらに好ましくは、選択手段は、ラティス内の単語の音響尤度と、N-グラム単語列の言語尤度とを、それぞれ所定の正規化方式により正規化するための正規化手段と、単語列ごとに、正規化手段により正規化された当該単語列内の単語の音響尤度とN-グラム単語列の言語尤度とにそれぞれ所定の重みを加算して得られる尤度が所定条件を満足するものを認識結果として選択するための手段とを含む。

【0045】

音響尤度と言語尤度とはそれぞれ別の概念に基づく尤度であるから、両者の尤度を単純に加算するのは意味がない。両者を予め正規化し、正規化後の尤度から算出される尤度によって単語列を選択することにより、尤度が正しく算出されることになり、頑健な音声認識を行なうことができる。

【0046】

この発明の他の局面に係る音声認識システムは、それぞれ異なる発話環境での発話音声のデコードに最適化された、DMFCCをパラメータとする音響モデル群を記憶するための記憶手段と、入力される音声からDMFCCを算出するためのDMFCC算出手段と、

10

20

30

40

50

入力される音声から算出されたDMFCCに基づいて、音響モデル群の混合重み適応化により、入力される音声の発話環境に適応化された適応化音響モデルを作成するためのモデル適応化手段と、適応化音響モデルを用いて、入力される音声から算出されるDMFCCをデコードするためのデコード手段とを含む。

【0047】

混合重み適応化を用いるため、実際の適応化時には音響モデル群を混合するための重みを算出するだけでよく、適応を高速に行なえる。その結果、DMFCCを音響特徴量として、発話環境の変化に対して頑健な音声認識を行なうことができる。

【発明を実施するための最良の形態】

【0048】

雑音環境が頻繁に変動する状況では、音響モデルを高速に雑音環境に適応させることが可能でなければならない。以下に述べる本発明の一実施の形態では、高速な雑音環境適応として、非特許文献14において提案されている雑音GMMの混合音適応化によるHMM合成法を用いる。

【0049】

図1～図3を参照して、この手法の概略について説明する。図1を参照して、あらかじめ準備した種々の雑音からなる雑音DB100から、個々の雑音を混合成分とする雑音GMM102と、個々の雑音に対して別々に適応化された雑音重畳音声用HMM104, 106, ...とを推定する。次に図2に示すように、短時間の未知雑音110を用いて雑音GMM102の混合ウェイト W_{N1} , W_{N2} , ...のみを推定する。そして、図3に示すように、この混合ウェイト W_{N1} , W_{N2} , ...を用いて、雑音重畳音声用HMM104, 106, ...を状態レベルで複数混合化する。例えばHMM104の状態120と、HMM106の状態122とに対して、それぞれのガウス混合分布に対し図2に示すステップにより計算された混合ウェイトを乗算して足し合わせて状態出力確率分布124を算出し、雑音適応されたHMMの状態126の状態出力確率分布とする。

【0050】

図1～図3において N_i は第*i*番目の雑音、 w_i は第*i*番目の雑音に対する雑音重畳音声用HMMを表す。 P_{Ni} と w_{Ni} は雑音GMMにおける第*i*番目の雑音の分布とその分布に対する混合ウェイトとをそれぞれ示す。さらに w_{ij} と p_{ij} は第*i*番目の雑音用の雑音重畳音声用HMMにおける第*j*番目の混合分布*N*の分岐確率と混合成分とを表す。

【0051】

この手法の利点として、適応の計算時間がGMMの混合ウェイトの推定時間のみであり大変高速である点と、雑音適応されたHMMが複数の雑音環境の分布を含んでおり、単一の雑音から推定された音響モデルよりも雑音の短時間の変動に対する頑健性が高い点とを挙げることができる。

【0052】

上記した混合重み適応化によるHMM合成法を用いる場合、音響特徴量としてはMFCCを用いることが考えられる。しかし、MFCCのみでは認識精度を高めることが難しいことが実験的に判明した。そこで本実施の形態では、MFCCとは異なる音響特徴量を用いた音声認識を行ない、その結果とMFCCによる音声認識の結果とを統合することを考える。本実施の形態では、雑音の変動に対して頑健な特徴量として非特許文献4において提案されたDMFCC特徴量を用いることとする。以下、DMFCC特徴について述べる。なお、以下の処理では、音声データは所定サンプリング周波数及び所定窓長でサンプリングしたフレームとして準備されているものとする。

【0053】

DMFCC特徴量は、式(1)に示すDPS(differential power spectrum)を基礎とする特徴量である。式(1)中の $Y(i, k)$ は、第*i*番目のフレームにおける第*k*番目のパワースペクトラム係数を表す。同様に $D(i, k)$ は第*i*番目のフレームにおける第*k*番目のDPS係数を表す。DMFCC特徴量は、このDPS係数に対してDCT(discrete cosine transform)を行

10

20

30

40

50

なうことにより抽出される。

$$D(i, k) = |Y(i, k) - Y(i, k + 1)| \quad (1)$$

有声母音などのピッチを含む音声から抽出されたパワースペクトラムは、基本周波数の高調波の影響によって櫛型の形状を持つ。このようなパワースペクトラムからDPS係数を計算した場合、隣り合うパワースペクトラム係数間の差が大きいため、DPS係数の値も同様に大きなパワーとして計算される。一方、雑音などの特徴を持たない波形のパワースペクトラムから計算されるDPS係数は、隣り合うパワースペクトラム係数間の差が小さいため、DPS係数の値も小さくなると考えられる。雑音重畳音声のパワースペクトラムを無雑音音声のパワーと雑音のパワーの和であると仮定した場合、DPS係数を計算することによって、音声と比較してなだらかに変化する雑音のパワー成分を減衰させることができると考えられる。

【0054】

本実施の形態では、上述のようにMFCC特徴量とDMFCC特徴量とを用いて、平行にデコーディングを行ない、得られた仮説の統合による音声認識精度の改善を試みている。

【0055】

図4に、本実施の形態に係る音声認識システム130の概略ブロック図を示す。図4を参照して、このシステム130は、初期HMM150と、雑音データベース(DB)152と、雑音が重畳された学習データ153とから、平行に音声をデコードするためのMFCC・HMM群156及びDMFCC・HMM群158を作成するためのHMM作成部154と、HMM作成部154により作成されたMFCC・HMM群156及びDMFCC・HMM群158を用いて、入力音声144に対する音声認識を行ない、音声認識結果146を出力するための認識処理部142とを含む。

【0056】

図5はHMM作成部154のブロック図である。図5を参照して、HMM作成部154は、初期HMM150と雑音DB152とから、前述したPMC法を用いて雑音重畳音声用MFCC・HMM群156を作成するためのMFCC雑音重畳音声用HMM推定部170と、雑音重畳済みの学習データ153を用いて初期HMM150に対する学習を行なうことにより、雑音重畳音声用DMFCC・HMM群158を作成するためのDMFCC雑音重畳音声用HMM推定部172とを含む。

【0057】

本実施の形態では、雑音DB152としては12種類の異なる雑音を用いる。学習データ153についても、無雑音学習データに上記したものと同種の雑音を重畳したのものを用いる。なお、雑音の重畳に際しては、10dB、20dB及び30dBの三種のSNRを用いている。初期HMM150としては、無雑音音響モデルとして学習済みのものを準備する。

【0058】

MFCC雑音重畳音声用HMM推定部170は、従来技術の項で説明した通りのPMC法を用いて各雑音に対応する雑音重畳音声用HMMを推定する機能を持つ。同様にDMFCC雑音重畳音声用HMM推定部172は、学習データ153を用いて最尤推定を行なうことにより雑音重畳音声用DMFCC・HMM群158の学習を行なう。DMFCC特徴量に対しては、MFCC特徴量と異なりPMC法が適用できないためである。

【0059】

図6に、MFCC雑音重畳音声用HMM推定部170による雑音重畳音声用MFCC・HMM群156の概念について示す。図6を参照して、MFCC用の初期HMM180は、無雑音通常発声用MFCC・HMM190と、無雑音言直し発話用MFCC・HMM192とを含む。本実施の形態では、発話スタイルの変動への対応としてシステムへの言直し時に頻繁に観測される音節強調発話に対する頑健性の改善を試みている。言直し発話用

10

20

30

40

50

のHMMはこのためのものである。

【0060】

音声認識ソフトウェアが認識誤りを起こした場合、そのソフトウェアの使用者はもう一度同じ発声を繰返さなければならない。このような言直し発話では、母音の後に短時間のポーズが挿入されるなど、通常発声とは異なる音響的特徴を持つことが報告されている。この言直し発話を頑健に認識するため、図17に示すような構造を持つ音響モデル440が提案されている。図17を参照して、この母音モデルは、母音の後に短時間ポーズを挿入するため、例えばt - a + s i lの状態パス(図17において、「t - a + k」などの表記は、先行音素が / t /、後続音素が / k /、当該音素が / a / の環境依存音素を表す。「s i l」は無音状態を表わす。)及び、その母音モデルの後にポーズ状態を追加した状態パスの合計3つの成分を有するマルチパス音響モデルの構造を持つ。さらに、このモデルでは、子音モデルの前に短時間ポーズの挿入を許すため、通常の子音モデルに加えてs i l - k + iの状態パスへの遷移が追加されている。このような音響モデルを用いることにより、通常発声の音声以外にも言直しや音節強調発声などの音声を頑健に認識することが可能となる。

10

【0061】

再び図6を参照して、雑音DB152は、本実施の形態では12種類の雑音データ200, 202, ..., 206を含む。MFCC雑音重畳音声用HMM推定部170はこれら12種類の雑音の各々について、3種類のSNR(10dB、20dB、及び30dB)ごとにPMCを用いて初期HMM180を適応化することにより、雑音重畳音声用MFCC

20

【0062】

生成される雑音重畳音声用MFCC・HMM群156は、男声通常発声用MFCC・HMM群210と、男声言直し発話用MFCC・HMM群212と、女声通常発声用MFCC・HMM群214と、女声言直し発話用MFCC・HMM群216と、通常発声用無雑音MFCC・HMM215と、言直し発話用無雑音MFCC・HMM217とを含む。すなわち本実施の形態では、雑音重畳音声用MFCC・HMM群156は、男声女声、12種類の雑音、3種類のSNR、及び通常発声、言直し発話用の、 $2 \times 12 \times 3 \times 2 = 144$ 種類と通常発声用及び言直し発話用の無雑音音声用モデルの計146種類のHMMを含む。

30

【0063】

図7に、MFCC雑音重畳音声用HMM推定部170により作成される音響モデルが、音響空間270中に占める領域を模式的に示す。図7に示すのは、12個の音響モデルに対応する領域280~302のみである。しかし、上述したように作成される音響モデルは146種類であるので、音響空間270にはこれら領域280~302と同様のものが合計で146個作成されることになる。

【0064】

図8に、DMFCC雑音重畳音声用HMM推定部172による雑音重畳音声用DMFCC・HMM群158の作成を概念的に示す。図8を参照して、初期DMFCC・HMM182は、無雑音通常発声用DMFCC・HMM230及び無雑音言直し発話用DMFCC

40

【0065】

また雑音重畳学習データ153は、前述した12種類の雑音を、前述した3種類のSNRで学習データに重畳したものであり、 $3 \times 12 = 42$ 種類の雑音重畳学習データ240~246を含む。DMFCC雑音重畳音声用HMM推定部172は、無雑音通常発声用DMFCC・HMM230及び無雑音言直し発話用DMFCC・HMM232に対し、上記した雑音重畳学習データ153を用いて学習を行なうことにより、男声通常発声用DMFCC・HMM群250、男声言直し発話用DMFCC・HMM群252、女声通常発声用DMFCC・HMM群254、女声言直し発話用DMFCC・HMM群256と、通常発声用無雑音DMFCC・HMM255と、言直し発話用無雑音DMFCC・HMM257

50

とを生成する。

【 0 0 6 6 】

例えば男声通常発声用 DMFCC・HMM群 250 は、各種類及び各 SNR の雑音重畳学習データに対して学習した結果得られた、複数個の男声雑音重畳通常発声用 DMFCC・HMM 260, 262, ..., 266 を含む。他の DMFCC・HMM群 252、254、256 も、男声か女声か、通常発声用モデルか言直し発話用モデルかを除き同様の構成である。

【 0 0 6 7 】

本実施の形態では、雑音重畳音声用 DMFCC・HMM群 158 は雑音重畳音声用 MFCC・HMM群 156 と同様の構成となっている。しかし、当業者であれば容易に理解できるように、MFCCを用いる音声認識と、DMFCCを用いる音声認識とで同様の構成をとる必要は全くない。それぞれ別々のデータに基づき HMM を作成してもよい。最終的に作成される HMM の数が等しくなる必要もない。

10

【 0 0 6 8 】

図 9 は、図 4 に示す認識処理部 142 の詳細な構造を示すブロック図である。図 9 を参照して、認識処理部 142 は、入力音声 144 に対し MFCC・HMM群を用いて音声認識を行なう MFCC 処理部 310 と、入力音声 144 に対し DMFCC・HMM群を用いた音声認識を行ない認識結果を出力するための DMFCC 処理部 312 と、MFCC 処理部 310 及び DMFCC 処理部 312 の出力を統合し、統合された認識結果を出力するための認識結果統合部 314 とを含む。

20

【 0 0 6 9 】

図 10 は MFCC 処理部 310 のより詳細なブロック図である。図 10 を参照して MFCC 処理部 310 は、入力音声 144 から MFCC パラメータを音響特徴量として算出するための MFCC 算出部 320 と、MFCC 算出部 320 から出力される MFCC パラメータに対し、MFCC・HMM群を用いて認識処理を行ない、HMMごとに認識結果を出力するための MFCC 通常発声認識処理部 322 と、MFCC 算出部 320 から与えられる MFCC パラメータに対し、言直し発話用 HMM を用いて認識処理を行ない、HMMごとに認識結果を出力するための MFCC 言直し発話認識処理部 324 と、MFCC 通常発声認識処理部 322 及び MFCC 言直し発話認識処理部 324 の出力のうち、尤度が最も高いものを選択して出力するための最尤選択部 326 とを含む。

30

【 0 0 7 0 】

図 11 は、DMFCC 処理部 312 のより詳細なブロック図である。図 11 を参照して DMFCC 処理部 312 は、入力音声 144 から音響特徴量として DMFCC パラメータを算出するための DMFCC 算出部 330 と、DMFCC 算出部 330 から与えられる DMFCC パラメータに対し DMFCC 通常発声用 HMM群を用いて認識処理を行ない、認識結果を HMMごとに出力するための DMFCC 通常発声認識処理部 332 と、DMFCC 算出部 330 から DMFCC パラメータを受取り、言直し発話用 DMFCC・HMM群を用いて認識処理を行ない、HMMごとに認識結果を出力するための DMFCC 言直し発話認識処理部 334 と、DMFCC 通常発声認識処理部 332 及び DMFCC 言直し発話認識処理部 334 から出力される認識結果のうち、尤度が最も高いものを選択して出力するための最尤選択部 336 とを含む。

40

【 0 0 7 1 】

図 10 及び図 11 を参照してわかるように、MFCC 処理部 310 及び DMFCC 処理部 312 の構造は互いに平行である。使用する音響特徴量が MFCC か DMFCC かによる差異があるにすぎない。従って以下では、MFCC 処理部 310 の構造の詳細についてのみ説明する。

【 0 0 7 2 】

図 12 は図 10 に示す MFCC 通常発声認識処理部 322 のより詳細なブロック図である。図 12 を参照して、MFCC 通常発声認識処理部 322 は、MFCC 算出部 320 から与えられる MFCC パラメータに基づき、男声通常発声用 MFCC・HMM群 210 及

50

び女声通常発声用MFCC・HMM群214に対する雑音GMMの混合重み適応化によるHMM合成を重畳された雑音のSNRごとに行ない、男声通常発声用適応化MFCC・HMM群354及び女声通常発声用適応化MFCC・HMM群352を生成するための雑音適応化処理部350と、男声通常発声用適応化MFCC・HMM群354を用いて、入力されるMFCCパラメータに対するデコードを行なうことにより、適応化されたHMMごとにデコード結果を出力するためのMFCC男声通常発声デコーダ部358と、入力されるMFCCパラメータに対し女声通常発声用適応化MFCC・HMM群を用いてデコードし、HMMごとにデコード結果を出力するためのMFCC女声通常発声デコーダ部356とを含む。

【0073】

男声通常発声用適応化MFCC・HMM群354及び女声通常発声用適応化MFCC・HMM群352はそれぞれ、3種類のSNRごとに一つ、合計三個のHMMを含む。デコードには無雑音HMMも使用するので、デコーダ部356及び358はそれぞれデコード結果を4つずつ出力する。その結果、MFCC通常発声認識処理部322全体としては8つのデコード結果を出力する。

【0074】

ここで、図12に示す雑音適応化処理部350の処理について図15及び図16を参照して説明する。図15を参照して、雑音適応化処理部350は、入力されるMFCCパラメータに基づき、音響空間270中における入力音声に対応する領域420を推定する。そしてこの領域420と、予め求められている各種の雑音が占める領域280～302との距離を算出する。そして、距離が最も近いものを所定個数（本実施の形態では4つ）だけ選択する。図15の例で示せば領域290、292、296及び298により示される雑音が、入力される音声の音響空間中の領域420に最も近い。従って、この4つの雑音に対応する音響モデルが採用される。

【0075】

続いて図16を参照して、これら4つの領域290、292、296及び298に対応するHMMのガウス混合分布の重みを計算し、加算することにより、入力される音声の音響空間270中における領域420をカバーするような音響モデルをHMMの形で算出する。この音響モデルを用いて入力音声に対するデコードを行なう。このように各雑音に対する音響モデル自体は変化させず今後のための重みのみを計算して音声認識用のHMMの適応化を行なえばよい。そのため適用の計算時間が短く、大変高速に適応化を行なうことができる。さらに、適応化されたHMMが複数の雑音環境の分布を含んでいる。従って単数の雑音から推定された音響モデルを用いた場合よりも、雑音の短時間の変動に対する頑健性がより高くなるという利点がある。

【0076】

図13はMFCC言直し発話認識処理部324の構成を示す。MFCC言直し発話認識処理部324は、入力されるMFCCパラメータを用いて、男声言直し発話用MFCC・HMM群212及び女声言直し発話用MFCC・HMM群216に対し雑音GMMの混合重み適応化によるHMM合成法を重畳された雑音のSNRごとに行ない、男声、女声及びSNRごとに適応化されたHMMを出力することにより、男声言直し発話用適応化MFCC・HMM群374及び女声言直し発話用適応化MFCC・HMM群372を出力するための雑音適応化処理部370と、与えられるMFCCパラメータを、女声言直し発話用適応化MFCC・HMM群372を用いてデコードし、HMMごとに出力するためのMFCC女声言直し発話デコーダ部376と、入力されるMFCCパラメータを男声言直し発話用適応化MFCC・HMM群374を用いてデコードし、HMMごとにデコード結果を出力するためのMFCC男声通常発声デコーダ部378とを含む。

【0077】

女声言直し発話用適応化MFCC・HMM群216は、SNRごとに合成される3つのHMMを含む。男声言直し発話用適応化MFCC・HMM群も同様に、SNRごとの3つのHMMを含む。また、デコードには無雑音HMMも使用される。従って、デコーダ部3

10

20

30

40

50

76及び378はそれぞれ4つずつのデコード結果を出力する。その結果MFCC言直し発話認識処理部324の出力は8つとなる。

【0078】

図12及び図13を参照して明らかなように、MFCC通常発声認識処理部322とMFCC言直し発話認識処理部324との構成はパラレルである。従って以下ではMFCC通常発声認識処理部322の詳細な構造のみを説明する。また図12及びこれ以前の説明から明らかなように、MFCC女声通常発声デコーダ部356及びMFCC男声通常発声デコーダ部358の構成も互いにパラレルである。従って以下では女声についてのみMFCC通常発声認識処理部322の詳細な構成を説明する。

【0079】

図14は、MFCC女声通常発声デコーダ部356及び女声通常発声用適応化MFCC・HMM群352の詳細な構成を示す。図14を参照して、女声MFCC・HMM群352は、無雑音HMM402、及びそれぞれ10dB、20dB、及び30dBのSNRで雑音が重畳された雑音重畳HMMから合成された10dB雑音HMM404、20dB雑音HMM406、及び30dB雑音HMM408とを含む。

【0080】

MFCC女声通常発声デコーダ部356は、入力されるMFCCパラメータを、無雑音HMM402、10dB雑音HMM404、20dB雑音HMM406、及び30dB雑音HMM408をそれぞれ用いてデコードし、デコード結果を出力するためのデコーダ390、392、394、及び396を含む。

【0081】

図18に、図9に示す認識結果統合部314のより詳細な構成を示す。図9に示すMFCC処理部310及びDMFCC処理部312からは複数の仮説が認識結果統合部314に与えられる。認識結果統合部314は、これら複数の仮説を単語単位で統合する。その原理について図19～図21を参照して説明する。

【0082】

複数の音声認識デコーダから得られた仮説が互いに相補的である場合、それぞれの仮説の正しい部分を抽出して組み合わせることにより、より正しい単語列が得られる可能性がある。ここで「相補的」とは、あるデコーダの認識結果の前半は正しいが後半は間違いであったとしても、別のデコーダの認識結果の後半部分が正しいならば、それぞれの正しい部分をつなぎあわせることによりその認識誤りを補償することができるという意味である。

【0083】

図19を参照して、2つの仮説470及び472が得られたものとする。仮説470の前半部分は誤っているが後半部分は正しい認識結果である。一方、仮説472については、前半の認識結果は正しいが後半は誤りである。従って仮説472の前半部分と仮説470の後半部分をつなぎ合わせることにより、正しい結果が得られるはずである。

【0084】

図20を参照して、上記した結果を得るために、まず図20に示されるような単語ラティスを、与えられた2つの仮説から再構成する。この再構成では、個々の単語の開始及び終了時間情報を用いる。

【0085】

続いて図21に示されるように、この単語ラティス480に含まれる単語列経路のうち、音響尤度と言語尤度とから算出される尤度が最も大きくなるような単語列482を再探索する。通常、仮説のうちでも正しい部分の尤度は高く、誤っている可能性が高い部分の尤度は低くなっている。従って、このような再探索を行なうことにより2つの仮説を統合して正しい結果を得ることができる可能性が高くなる。

【0086】

なお本実施の形態では、MFCCとDMFCC特徴量から得られた仮説に対する仮説統合を認識結果統合部314で行なっている。この場合、MFCCの音響モデルから計算さ

10

20

30

40

50

れる音響尤度と、DMFCCの音響モデルから計算される尤度とを直接比較することはできない。そのため、音響モデルの尤度を比較するためには尤度の正規化が必要である。本実施の形態では、そのために、認識文全体の音響尤度で個々の単語の音響尤度を割ることにより、各単語の尤度を正規化する。さらに、仮説統合の際には、言語モデルを用いた尤度計算も行なう。この場合、音響モデルの尤度計算と言語モデルによる尤度計算との間での重み付けを考慮しなければならない。本実施の形態では、仮説統合時における言語モデルウェイトを0.06とした。

【0087】

図18を参照して、認識結果統合部314は、上記したような機能を実現するために以下の各処理部を含む。すなわち認識結果統合部314は、MFCC及びDMFCCのそれぞれの仮説の単語の音響尤度を正規化するための尤度正規化部450と、2つの仮説から個々の単語の開始及び終了時間情報を用いて単語ラティス480(図20参照)を作成するための単語ラティス作成部452と、統合の際に参照される言語モデルを記憶するための言語モデル記憶部456と、統合の際の言語モデルの尤度の、音響モデルの尤度に対するウェイトを記憶するためのウェイト記憶部454と、単語ごとの音響尤度及び言語モデルに基づく単語列の尤度に基づいて単語ラティス480中の、音響尤度と言語尤度の和が最大となるような単語列を再探索することにより認識結果を統合するための最尤経路探索部458とを含む。

【0088】

上記した音声認識システム130は以下のように動作する。図22に、このシステムの動作の概略の流れについて示す。大きく分けて、このシステムは2つの動作局面を持つ。第一の局面は、雑音重畳音声用のHMMを準備するステップ500である。第二の局面は、このようにして準備された雑音重畳音声用のHMMと無雑音用のHMMとを用いて、入力される音声の認識を行なうステップ(502~508)である。

【0089】

ステップ500では、図4に示すような初期HMM150と、雑音DB152とを用いて、MFCC・HMM群156が作成され、また雑音重畳学習データ153を用いてDMFCC・HMM群158が作成される。

【0090】

このようにして、雑音重畳音声用のHMM群が作成された後は、いつでもこのMFCC・HMM群156及びDMFCC・HMM群158を用いた音声認識を行なうことができる。図4に示す入力音声144が与えられると、その入力音声からMFCCパラメータ及びDMFCCパラメータが算出される(ステップ502)。それらを用いて、予め準備されたMFCC・HMM群156及びDMFCC・HMM群158のうち入力音声144の発話環境に最も類似した発話環境に対応する所定個数(本実施の形態では4個)のHMMがMFCC及びDMFCCのそれぞれについて選択される。これらHMMからMFCC及びDMFCCの各々について、雑音GMMの混合重み適応化によるHMMが合成される。合成されるHMMは、男声・女声、通常発声・言直し発話、及び4種類のSNR(10dB、20dB、30dB、無雑音)の組み合わせの各々に対してであるから、全部で $2 \times 2 \times 4 = 16$ 通りである。

【0091】

続いてステップ504で発話入力があったか否かが判定される。発話入力があればステップ506に進むが、発話入力がない場合は、再び重み推定502を行なう。本実施の形態では、発話入力があった場合には、その直前の1秒間の期間における雑音を用いて重み推定を行なっている。

【0092】

ステップ506では、合成されたHMMを用いた認識と、それら認識結果の統合とが行なわれる。その認識結果がステップ508で出力される。その後再度重み推定502の処理から繰り返される。

【0093】

10

20

30

40

50

図23を参照して、発話522に対しては、発話522の直前の雑音524を用いて合成されたHMMによる音声認識が行なわれる。同様に次の発話526に対しては、発話526の直前の雑音528により推定されたHMMを用いて音声認識が行なわれる。

【0094】

なお、上記した男声女声、MFCC及びDMFCC、通常発声及び言直し発話などの組合せは任意に選ぶことができる。MFCC又はDMFCCのいずれか一方のみを用いるシステムも可能である。

[実験1]

上記した実施の形態に係る雑音適応化手法の評価を行なうため、日本語大語彙連続音声認識実験を行なった。実験においては、予め出願人において作成した言語モデルを準備した。言語モデルの作成に使用された自然発話音声・言語データベースに含まれていた単語は670万語程度である。実験に使用した音声波形は、サンプリング周波数16kHz、分析窓長20ms、分析周期10msで分析を行ない、MFCC及びDMFCC特徴量を抽出した。MFCCの音響特徴パラメータは、12次元MFCC、 c_0 、12次元DMFCCの計25次元である。DMFCCの音響特徴パラメータは、12次元DMFCC、 pow 、12次元DMFCCの計25次元である。使用した音素は、日本語分析でよく用いられる26種類の音素である。

【0095】

音響モデルの状態共有構造は、ML-SSS（非特許文献16を参照されたい）より生成した2100状態のHMnetを使用した。各状態の混合数は5である。

【0096】

学習データとして、出願人において準備した旅行会話データベースTRAを用いた。このデータベースTRAは、407名が発声した対話及び音素バランス503文の計30時間である。

【0097】

雑音適応元の音響モデルは、様々な場所で採取した12種類の雑音を用いて生成した。MFCCの音響モデルは、雑音とSNR毎にPMC法を用いて無雑音音声HMMを適応化することにより生成した。DMFCCの音響モデルは、雑音を重畳した学習データを用いて生成した。雑音重畳音声のSNRは、10dB、20dB、30dBである。

【0098】

MFCCとDMFCCの音響モデルはそれぞれ、男声女声、12種類の雑音、及び3種類のSNRとの組合せからなる、 $2 \times 12 \times 3 = 72$ 種類と無雑音音声モデルとの計73種類である。

【0099】

評価用音声データは、出願人において準備したATR旅行会話基本表現集BTEC testset-01（510文、男性4名、女性6名、それぞれ51文の発声データ）を使用し、10dB、20dB、30dBのSNRで雑音を重畳した。評価用に重畳した雑音はHMMの合成に用いた雑音とは異なる複数の場所で採取した雑音である。雑音GMMの混合ウェイト推定には1秒間の雑音を使用して個々の混合ウェイトの上位4つの雑音を用いて雑音重畳音声用音響モデルを生成した。

【0100】

図24に、3種類の評価用雑音重畳音声データに対する平均単語正解精度を示す。図中のMAXは個々の音響モデル（10dB、20dB、30dB、無雑音）を用いて得られた仮説を最大尤度基準で選択した場合の単語正解精度である。図24に示すように、最大尤度基準による選択を行なうことで、実験に用いたSNR全てにおいて平均90%以上の単語正解精度が得られた。DMFCCの音響モデルを用いるとMFCCの音響モデルを用いた場合よりも単語正解精度が低下している。しかしDMFCCの無雑音音声音響モデルを用いた場合、雑音重畳音声の単語正解精度がMFCCの無雑音音声音響モデルよりも高い。従って、雑音の種類や雑音SNRに対する正解精度への影響がMFCCよりも小さいことがわかる。

10

20

30

40

50

[実験 2]

さらに、言直し発話に対し頑健な音響モデルに対して雑音と発話スタイルの変動に対する単語正解精度への影響を調べるため、日本語大語彙連続音声認識実験を行なった。評価用音声として、実験 1 で用いた通常発声の音声と、意図的に音節ごとに区切って発声した音節強調発声の音声とを用いた。音節強調発声データは、旅行会話文、男性 2 名女性 2 名、各話者 10 文の計 40 文である。評価用音声には 30 dB、20 dB、10 dB の SNR で、実験 2 で用いた 3 種類の雑音が重畳されている。

【 0 1 0 1 】

言直し発話に頑健な音響モデルは、環境依存音素モデル数が通常発声モデルよりも多い。そのため探索空間が大きく広がり、通常発声音声に対して単語正解精度の低下が懸念される。そこで、上記実施の形態で説明した通り、言直し発話用音響モデルと通常発声用音響モデルとを用いて別々にデコーディングし、最大尤度基準による仮説の選択を行なった。

10

【 0 1 0 2 】

図 25 に、通常発声用音響モデルの場合、言直し発話用音響モデル単独の場合、2 つの音響モデルをパラレルデコーディングした場合それぞれに対する単語正解精度を示す。図 25 に示すように、言直し発話用音響モデルを単独で使用した場合、その単語正解精度は若干低下する。それに対しパラレルデコーディングを行なうことにより、通常発声の音声に対してもほぼ同等の正解精度が得られた。

【 0 1 0 3 】

次に、音節強調発声の音声に対する単語正解精度を図 26 に示す。図 26 に示すように、言直し発話用音響モデルは、通常発声用音響モデルよりも高い単語正解精度が得られた。雑音重畳音声に対しても、実験 1 で得られた結果同様、10 dB の音声に対しても無雑音音声や 30 dB の音声と同程度の単語正解精度が得られた。

20

【 0 1 0 4 】

[実験 3]

最後に、MFCC 特徴量と DMFCC 特徴量のデコーダから得られた仮説を統合することによる性能の改善を調べるための評価実験を行なった。予備実験から、上記実施の形態で述べたように仮説統合時における言語モデルウェイトを 0.06 とした。図 27 に、仮説統合を行なった場合の単語正解精度を示す。図 27 に示すように、通常発声に対しては MFCC 特徴量の正解精度と同等の結果が得られた。さらに、音節強調発声に対しては、MFCC と DMFCC の各々の正解精度以上の性能が得られた。これは、仮説統合により、MFCC による仮説と DMFCC による仮説とが互いに相補的であったため、仮説統合によって精度が高くなったためと考えられる。

30

【 0 1 0 5 】

以上のように本実施の形態の音声認識システム 130 では、雑音と発話スタイルの変動に頑健な音声認識を実現することを目指した。本システムでは、雑音の変動に頑健な音響特徴量としての DMFCC、予め種々の雑音環境に適応化した HMM を用いて雑音 GMM の混合ウェイトから雑音適応 HMM を高速に生成する雑音適応手法、言直し発話に頑健な音響モデル、及び複数の仮説を統合する手法を用いた。その結果、10 dB から 30 dB の SNR で雑音を重畳した通常発声の評価データに対して、平均 90% 以上の単語正解精度が得られた。また、言直し発話などの発話スタイルの変動に対しても、通常発声用音響モデルのみを用いた場合よりも高い単語正解精度が得られた。

40

【 0 1 0 6 】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味および範囲内のすべての変更を含む。

【 図面の簡単な説明 】

【 0 1 0 7 】

50

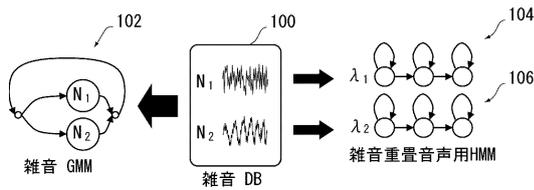
- 【図 1】雑音 G M M 及び雑音重畳音声 H M M の作成を説明するための図である。
- 【図 2】混合重みの推定を説明するための図である。
- 【図 3】適応化 H M M の生成を説明するための図である。
- 【図 4】本発明の一実施の形態に係る音声認識システムのブロック図である。
- 【図 5】H M M 作成部のより詳細なブロック図である。
- 【図 6】本発明の一実施の形態における雑音重畳音声用 M F C C ・ H M M 群の作成を説明するための図である。
- 【図 7】雑音 G M M の混合重み適応化において、P M C 法により準備される雑音 H M M を説明するための図である。
- 【図 8】本発明の一実施の形態において雑音重畳音声用 D M F C C ・ H M M を作成する方法を説明するための図である。 10
- 【図 9】認識処理部のより詳細な構成を示すブロック図である。
- 【図 10】M F C C 処理部 3 1 0 の詳細な構成を示すブロック図である。
- 【図 11】D M F C C 処理部 3 1 2 の詳細な構成を示すブロック図である。
- 【図 12】M F C C 通常発声認識処理部の詳細な構成を示すブロック図である。
- 【図 13】M F C C 言直し発話認識処理部の詳細な構成を示すブロック図である。
- 【図 14】M F C C 女声通常発声デコーダ部 3 5 6 及び女声通常発声用適応化 M F C C ・ H M M 群 3 5 2 の詳細な構成を示すブロック図である。
- 【図 15】本実施の形態における入力音声の発話環境から、予め準備された雑音 H M M の発話環境までの距離を概念的に説明するための図である。 20
- 【図 16】入力音声の発話環境に類似した雑音を含む雑音 H M M から適応化 H M M を合成する概念を示す図である。
- 【図 17】言直し発話に頑健な音響モデルの構成を示す図である。
- 【図 18】認識結果統合部の詳細な構成を示すブロック図である。
- 【図 19】仮説統合の経過を説明するための、2 つの仮説を示す図である。
- 【図 20】仮説統合の過程で生成される単語ラティスを示す図である。
- 【図 21】仮説統合の際に行なわれる最尤単語列の探索を説明するための図である。
- 【図 22】本発明の一実施の形態に係る音声認識システムの動作を説明するための図である。
- 【図 23】発話ごとの音声認識に用いられる雑音の位置を説明するための図である。 30
- 【図 24】本発明の一実施の形態を用いて行なわれた、雑音適応による頑健化の評価実験の結果を示すグラフである。
- 【図 25】本発明の一実施の形態の音声認識システムを用いて行なわれた、言直し発話用音響モデルによる頑健化の評価実験の結果を示すグラフである。
- 【図 26】音節強調発声の音声に対する単語正解精度を示すグラフである。
- 【図 27】本発明の一実施の形態に係る音声認識システムにおいて行なわれる仮説統語による単語正解精度を調べる実験結果を示すグラフである。
- 【図 28】P M C 法の概念を説明するための図である。

【符号の説明】

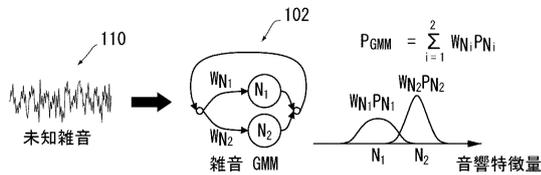
- 【 0 1 0 8 】 40
- 1 3 0 音声認識システム、1 5 0 初期 H M M、1 5 2 雑音 D B、1 5 3 雑音重畳学習データ、1 5 4 H M M 作成部、1 5 6 雑音重畳音声用 M F C C ・ H M M 群、1 5 8 雑音重畳音声用 D M F C C ・ H M M 群、1 4 4 入力音声、1 4 2 認識処理部、1 4 6 音声認識結果、1 9 0 無雑音通常発声用 M F C C ・ H M M、1 9 2 無雑音言直し発話用 M F C C ・ H M M、2 1 0 男声通常発声用 M F C C ・ H M M 群、2 1 2 男声言直し発話用 M F C C ・ H M M 群、2 1 4 女声通常発声用 M F C C ・ H M M 群、2 1 6 女声言直し発話用 M F C C ・ H M M 群、2 3 0 無雑音通常発声用 D M F C C ・ H M M、2 3 2 無雑音言直し発話用 D M F C C ・ H M M、2 5 0 男声通常発声用 D M F C C ・ H M M 群、2 5 2 男声言直し発話用 D M F C C ・ H M M 群、2 5 4 女声通常発声用 D M F C C ・ H M M 群、2 5 6 女声言直し発話用 D M F C C ・ H M M 群、3 1 0 M 50

FCC 処理部、312 DMFCC 処理部、314 認識結果統合部、320 MFCC 算出部、322 MFCC 通常発声認識処理部、324 MFCC 言直し発話認識処理部、326 最尤選択部、330 DMFCC 算出部、332 DMFCC 通常発声認識処理部、334 DMFCC 言直し発話認識処理部、336 最尤選択部、350 雑音適応化処理部、356 MFCC 女声通常発声デコーダ部、358 MFCC 男声通常発声デコーダ部、370 雑音適応化処理部、376 MFCC 女声言直し発話デコーダ部、378 MFCC 男声言直し発話デコーダ部、450 尤度正規化部、452 単語ラティス作成部、458 最尤経路探索部、480 単語ラティス

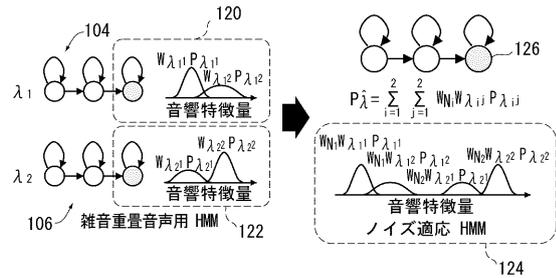
【図1】



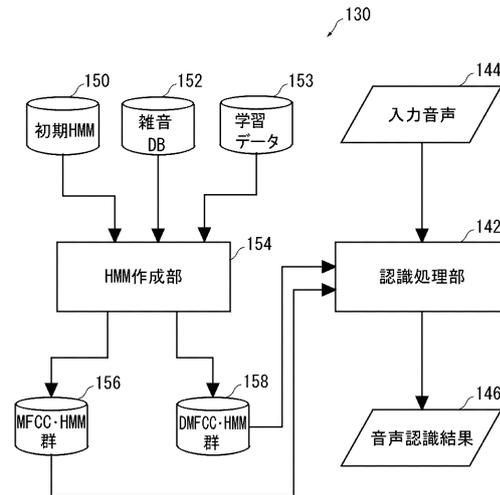
【図2】



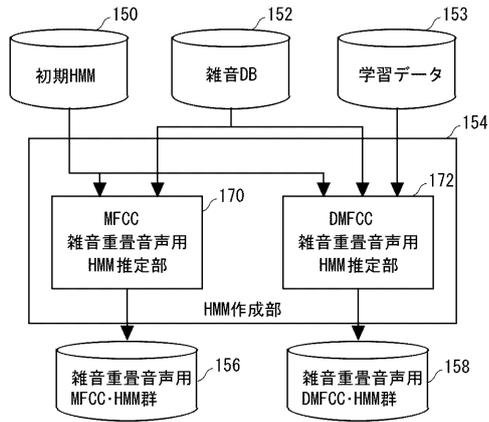
【図3】



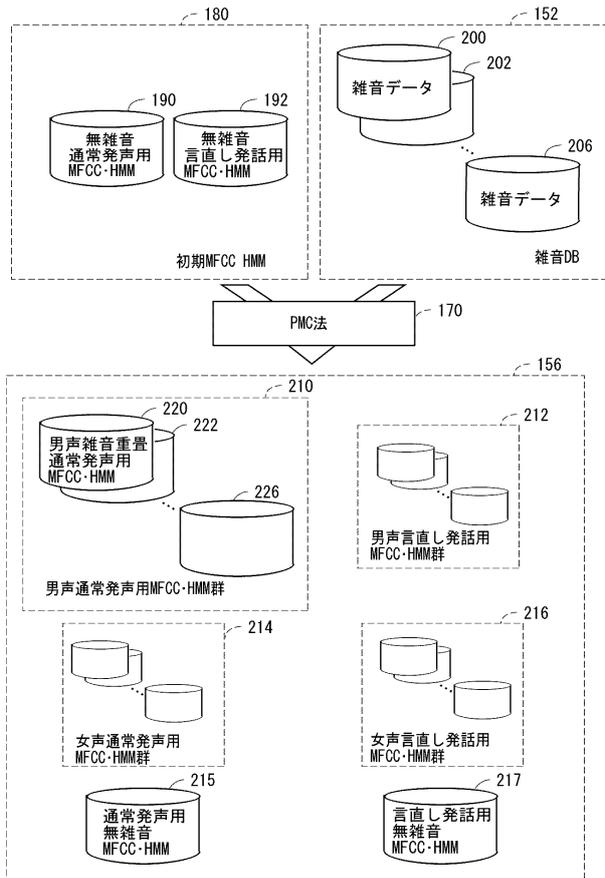
【図4】



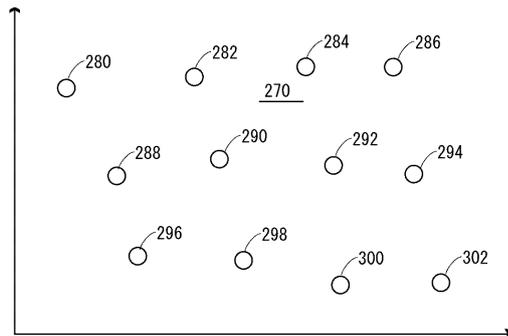
【図5】



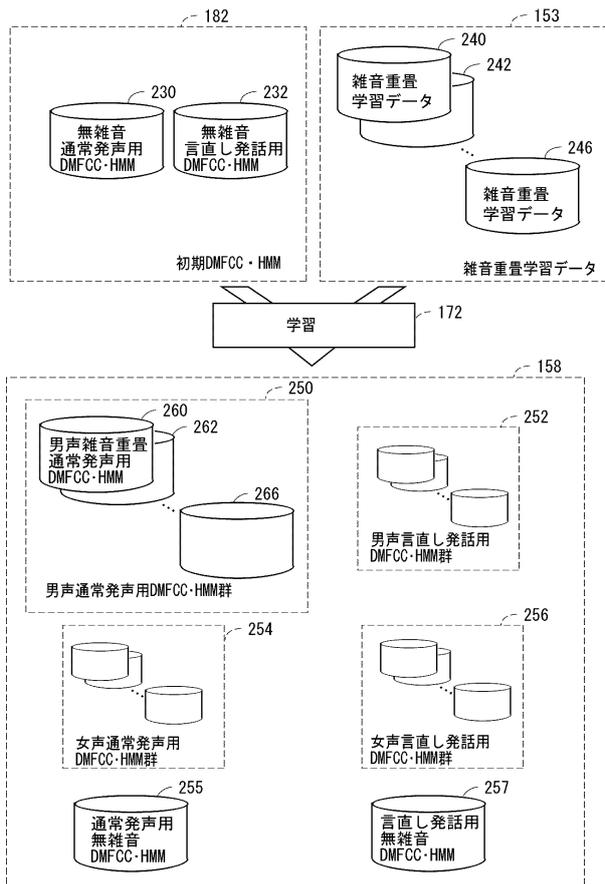
【図6】



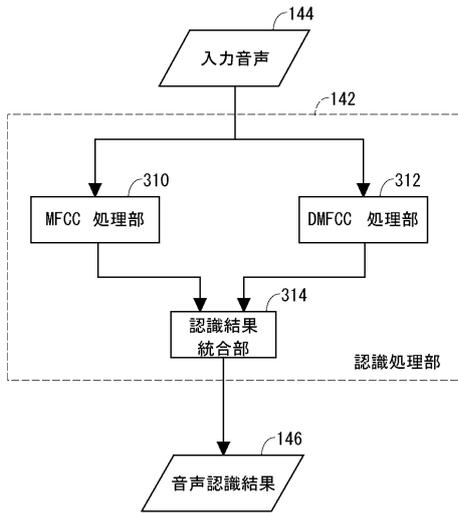
【図7】



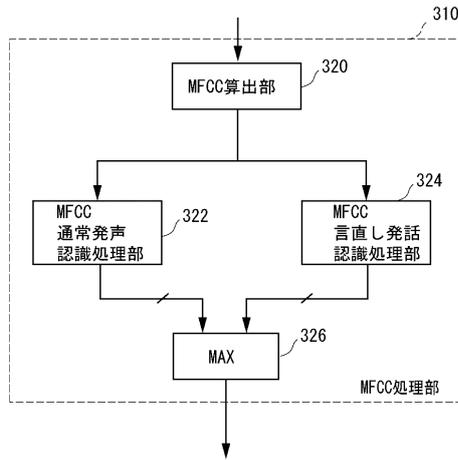
【図8】



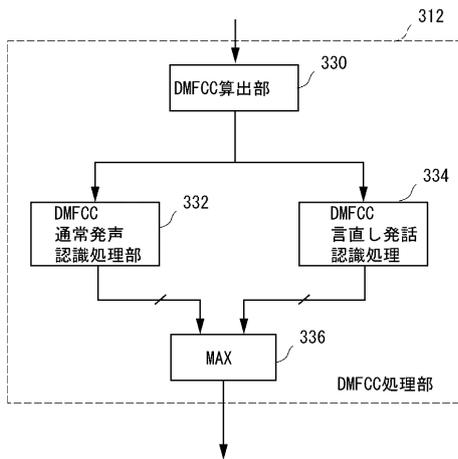
【図9】



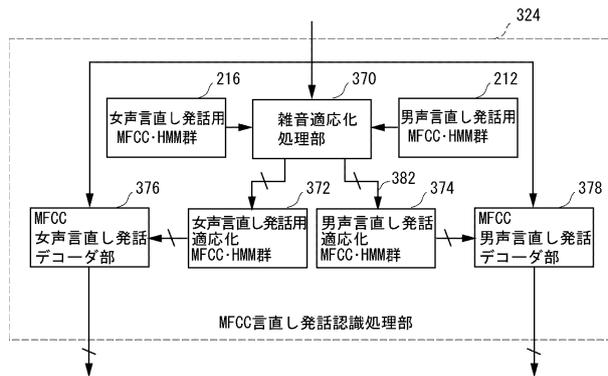
【図10】



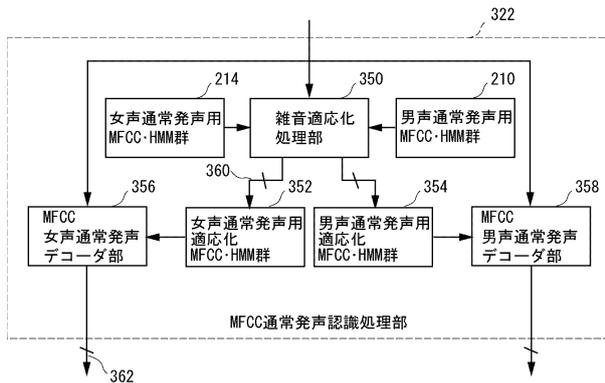
【図11】



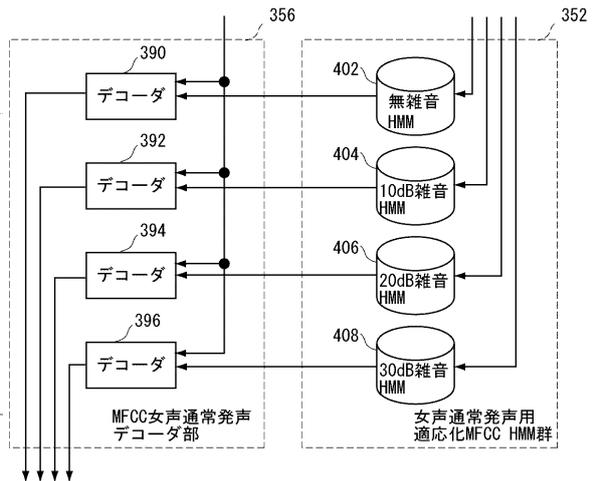
【図13】



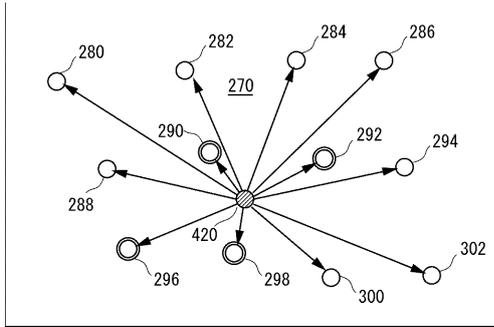
【図12】



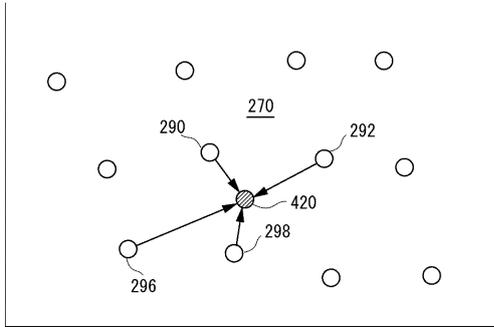
【図14】



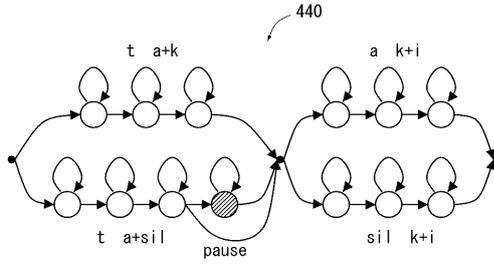
【図15】



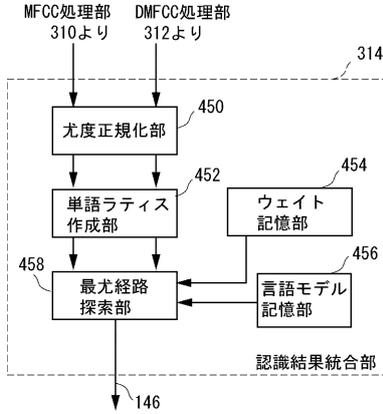
【図16】



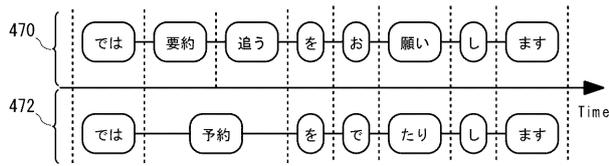
【図17】



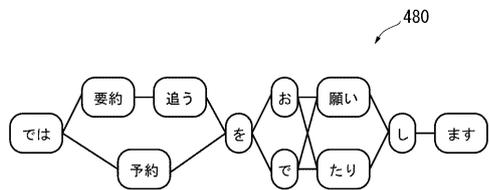
【図18】



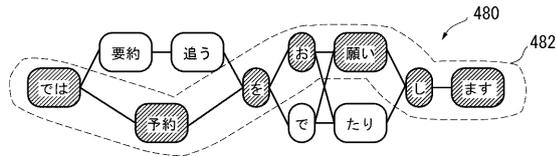
【図19】



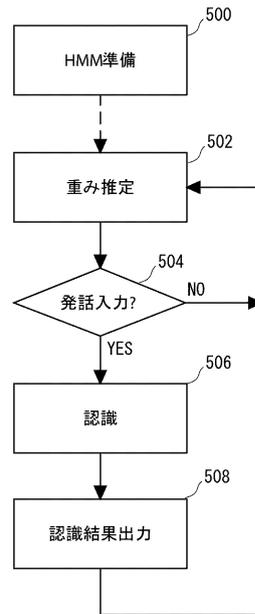
【図20】



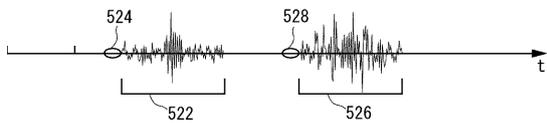
【図21】



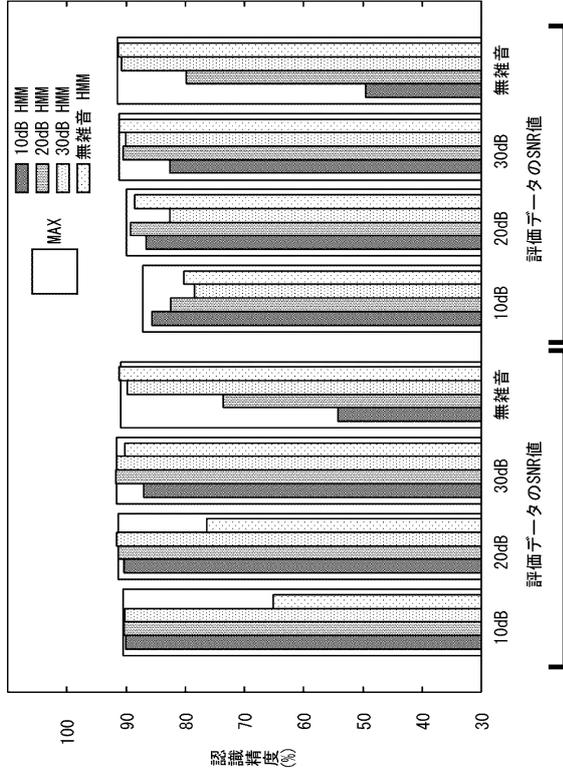
【図22】



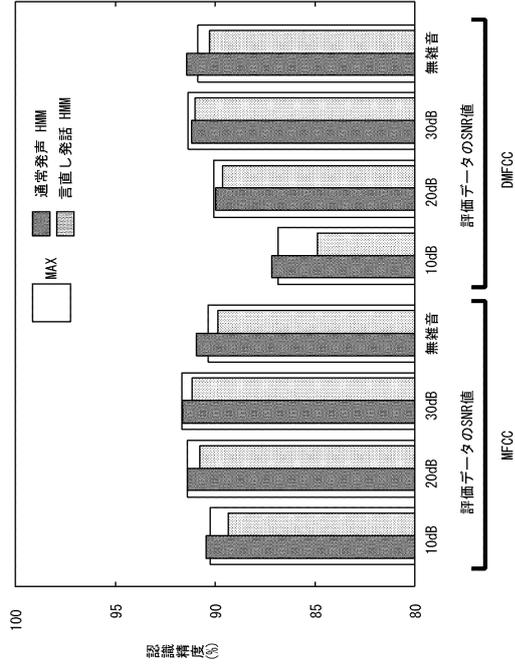
【図23】



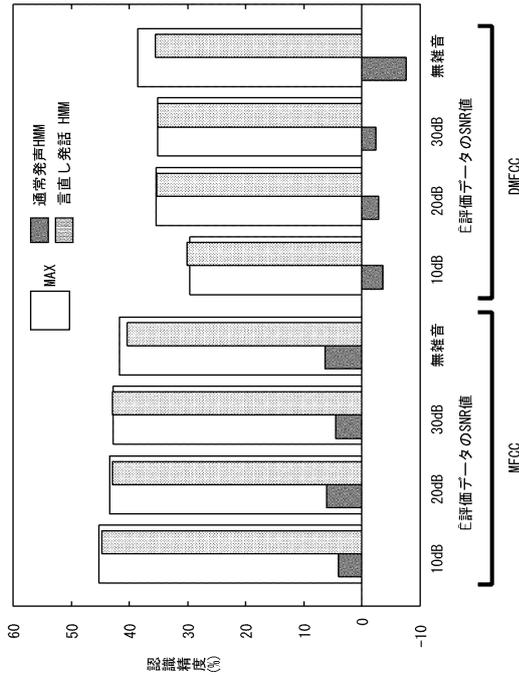
【図 24】



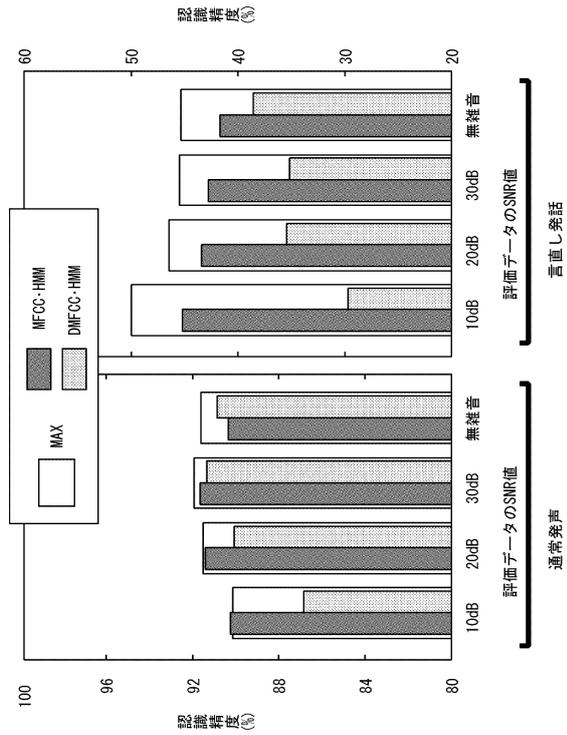
【図 25】



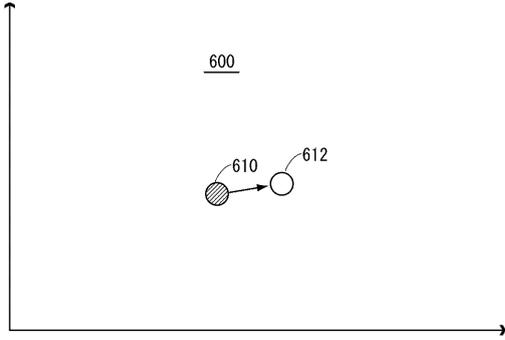
【図 26】



【図 27】



【 28】



フロントページの続き

(72)発明者 中村 哲

京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

審査官 菊池 智紀

(56)参考文献 特開2003-177781(JP,A)

特開2002-189494(JP,A)

特開2004-004509(JP,A)

Konstantin MARKOV et al., "Noise and Channel Distortion Robust ASR System for DARPA SPINE2 Task", IEICE Trans. Inf. & Syst., 2003年 3月 1日, Vol.E86-D, No.3, p.497-504

伊田政樹 他, "雑音GMMの適応化とSN比別マルチパスモデルを用いたHMM合成による高速な雑音環境適応化", 電子情報通信学会論文誌D-II, 2003年 2月 1日, Vol.J86-D-II, No.5, p.195-203

(58)調査した分野(Int.Cl., DB名)

G10L 15/00 - 17/00

JSTPlus(JDreamII)