

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2017-102247
(P2017-102247A)

(43) 公開日 平成29年6月8日(2017.6.8)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 0 L 15/10 (2006.01)	G 1 0 L 15/10 5 0 0 T	
G 1 0 L 15/22 (2006.01)	G 1 0 L 15/22 3 0 0 U	
G 1 0 L 13/00 (2006.01)	G 1 0 L 13/00 1 0 0 M	
G 1 0 L 15/02 (2006.01)	G 1 0 L 15/02 3 0 0 K	
G 1 0 L 15/197 (2013.01)	G 1 0 L 15/197	

審査請求 未請求 請求項の数 7 O L (全 11 頁)

(21) 出願番号 特願2015-234835 (P2015-234835)
(22) 出願日 平成27年12月1日 (2015.12.1)

(出願人による申告) 平成27年度国立研究開発法人科学技術振興機構研究成果展開事業 戦略的イノベーション創出推進プログラム「高齢者の記憶と認知機能低下に対する生活支援ロボットシステムの開発」委託研究、産業技術力強化法第19条の適用を受ける特許出願

(71) 出願人 301021533
国立研究開発法人産業技術総合研究所
東京都千代田区霞が関1-3-1
(72) 発明者 佐土原 健
茨城県つくば市東1-1-1 国立研究開発法人産業技術総合研究所つくばセンター内

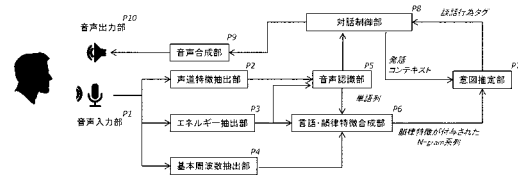
(54) 【発明の名称】 音声対話システム、音声対話制御法およびプログラム

(57) 【要約】 (修正有)

【課題】 談話行為識別精度の向上を図る。

【解決手段】 音声入力部 P 1 と、声道特徴系列を計算する声道特徴抽出部 P 2 と、エネルギー系列を計算するエネルギー抽出部 P 3 と、基本周波数系列を計算する基本周波数抽出部 P 4 と、声道特徴系列とエネルギー系列に基づいてタイムスタンプ付きの単語列を出力する音声認識部 P 5 と、タイムスタンプ付きの単語列とタイムスタンプ付きの前記エネルギー系列と前記基本周波数系列から韻律特徴が付加された拡張単語 N - g r a m を生成する言語・韻律特徴合成部 P 6 と、拡張単語 N - g r a m と該発話のコンテキスト情報から談話行為タグを推定する意図推定部 P 7 と、談話行為タグが表す発話意図と対話文脈を考慮してシステム発話を生成する対話制御部 P 8 と、該生成されたシステム発話を音声信号に変換する音声合成部 P 9 と、該変換された音声信号を音声出力装置で再生する音声出力部 P 1 0 とからなる。

【選択図】 図 1



【特許請求の範囲】

【請求項 1】

人との音声対話インタフェースを含む音声対話システムであって、
 音声入力装置からの音声入力を処理して音響信号に変換する音声入力部(P 1)と、
 その音響信号を処理して声道特徴系列を計算する声道特徴抽出部(P 2)と、
 その音響信号を処理してエネルギー系列を計算するエネルギー抽出部(P 3)と、
 その音響信号を処理して基本周波数系列を計算する基本周波数抽出部(P 4)と、
 該計算された声道特徴系列とエネルギー系列に基づいてタイムスタンプ付きの単語列を
 出力する音声認識部(P 5)と、

該出力されたタイムスタンプ付きの単語列とタイムスタンプ付きの前記エネルギー系列
 と前記基本周波数系列から韻律特徴が付加された拡張単語 N - g r a m を生成する言語・
 韻律特徴合成部(P 6)と、

該生成された拡張単語 N - g r a m と該発話のコンテキスト情報から該発話の談話行為
 タグを推定する意図推定部(P 7)と、

前記該発話のコンテキスト情報を提供し該推定された談話行為タグが表す発話意図と対
 話文脈を考慮してシステム発話を生成する対話制御部(P 8)と、

該生成されたシステム発話を音声信号に変換する音声合成部(P 9)と、

該変換された音声信号を音声出力装置で再生する音声出力部(P 10)とからなることを
 特徴とする音声対話システム。

【請求項 2】

言語・韻律特徴合成部(P 6)において、前記拡張単語 N - g r a m に付加された韻律特
 徴は離散化された韻律特徴であることを特徴とする請求項 1 に記載する音声対話システム
 。

【請求項 3】

意図推定部(P 7)において、前記該発話の談話行為タグを推定は、次の談話行為タグを
 表す確率変数 I の事後確率の最大化により行うことを特徴とする請求項 2 に記載の音声対
 話システム。

ただし、A は音響信号、 A_s はその声道成分、W は単語列、 W^f は前記離散化された韻律
 特徴が付与された拡張単語列とする。

【数 1 2】

$$P(A|I) \approx \sum_W P(A_s|W)P(W^f|I) \quad (12)$$

【請求項 4】

言語・韻律特徴合成部(P 6)において、前記離散化を、W に対応するフレーム列におい
 て、当該基本周波数の変化量が平均 + 標準偏差よりも大きい場合は + を、変化量が平均 -
 標準偏差よりも小さい場合は - を、変化量が平均 ± 標準偏差の範囲内であれば 0 を付与す
 る 3 値の離散化であることを特徴とする請求項 3 に記載の音声対話システム。

【請求項 5】

意図推定部(P 7)において、前記尤度 $P(W^f | I)$ は、次式により計算されることを特
 徴とする請求項 4 に記載の音声対話システム。

ただし W_i^f は W に含まれる N - g r a m、 W_1, \dots, W_m に前記離散韻律特徴を付与した
 、拡張 N - g r a m 列 W_1^f, \dots, W_m^f の 1 つとする。

【数 1 3】

$$P(W^f|I) = \prod_{i=1}^m P(W_i^f|I) \quad (13)$$

【請求項 6】

意図推定部 (P 7) において、前記尤度 $P(W_i^f | I)$ は、次式により計算されることを特徴とする請求項 5 に記載の音声対話システム。

ただし、 j は自然数、 $j > 0$ の場合は、基本周波数の勾配、強勢等、 f で用いられている離散韻律特徴が付与された $N - gram$ を表し、特に、 $j = 0$ の場合は、韻律特徴を用いない $N - gram$ を表す。

【数 14】

$$P(W_i^f | I) = \sum_{j=0}^J \lambda_j P(W_i^{fj} | I), \quad \sum_{j=0}^J \lambda_j = 1 \quad (14)$$

10

【請求項 7】

人との音声対話インタフェースを含む音声対話プログラムであって、請求項 1 乃至請求項 5 のいずれか 1 項に記載される音声対話システムの各処理を実行することを特徴とする音声対話プログラム、および当該プログラムを記憶したプログラム媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声対話を利用してユーザと情報のやりとりを行う、音声対話システム、音声対話制御手法、プログラムに関する。

20

【背景技術】

【0002】

従来、機器やシステムのインタフェースとして、ユーザが日常的に用いているコミュニケーション手段である音声対話を用いたインタフェースが利用されてきた。

音声対話インタフェースにおいては、ユーザの発話からユーザの意図を抽出し適切な応答を生成する対話制御技術が必要になる。

【0003】

ユーザの意図の中でも、対話文脈における発話の機能、例えば、相手の発言に対して、肯定しているのか、否定しているのか、あるいは疑問を発しているか等を判定する問題は、談話行為識別と呼ばれている。

30

例えば、非特許文献 1 では、発話意図を表す談話行為タグを予測するための確率モデルを、発話に談話行為タグが付与された音声対話データから、統計的に機械学習させ、学習されたモデルを用いて、音声対話中の発話に談話行為タグを付与する技術が開示されている。

【0004】

ところで、これまでの発話意図推定技術においては、音声に含まれる言語特徴が意図を推定するための重要な手がかりとして用いられてきた。

また、表情や韻律 (音声のリズム、抑揚、速度、強勢など) 等の非言語特徴も、発話者の意図を反映することが良く知られている。

40

【0005】

例えば、特許文献 1 では、表情と韻律を用いたユーザの感情の推定技術が開示されている他、特許文献 2 では、システムの間違いをユーザが指摘した箇所を同定する技術が開示されている。

また、非特許文献 1 および非特許文献 2 では、韻律に関する発話の統計量を言語情報に付加して、談話行為識別精度を向上させる技術が開示されている。

【0006】

さらに、特許文献 3 および非特許文献 3 においては、ピッチの変化等、韻律特徴の軌跡を離散符号化した系列を言語特徴とともに利用して談話行為識別精度を向上させる技術が開示されている。

50

いずれも、言語特徴に付加する形で、言語特徴とは独立な、発話文の音響的な特徴として韻律特徴が用いられている。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】特開2006-313287号公報

【特許文献2】特開2013-205842号公報

【特許文献3】米国特許出願公開第2006/0122834号明細書

【非特許文献】

【0008】

【非特許文献1】A. Stolcke et al. : "Dialogue act modeling for automatic tagging and recognition of conversational speech", Computational Linguistics, Vol.26, No.3, pp.339-373, 2000.

【非特許文献2】E. Shriberg et al. : "Can prosody aid the automatic classification of dialog acts in conversational speech?", Language and speech, Vol.41, No.3-4, pp.443-492,1998.

【非特許文献3】V.K.R. Sridhar et al. : "Combining lexical, syntactic and prosodic cues for improved online dialog act tagging", Computer Speech and Language, Vol.23, No.4, pp.407-422, 2009.

【非特許文献4】A. Black et al. : "Predicting the intonation of discourse segments from examples in dialogue speech.", Computing Prosody, pp.117-128, 1997.

【非特許文献5】K.Sadohara et al. : "Sub-lexical dialogue act classification in a spoken dialogue system support for the elderly with cognitive disabilities". In Proceedings of the Workshop on Speech and Language Processing for Assistive Technologies. pp.93-98, 2013.

【非特許文献6】鹿野清宏他：“音声認識システム”，オーム社，2001。

【非特許文献7】河原達也，荒木雅弘：“音声対話システム”，オーム社，2006。

【非特許文献8】P.Taylor: "The tilt intonation model", In Proceedings of the International Conference on Spoken Language Processing, Vol.4, pp.1383-1386,1998.

【非特許文献9】S.Ananthakrishnan: "Categorical prosody models for spoken language applications", PhD Thesis, University of Southern California, 2008.

【発明の概要】

【発明が解決しようとする課題】

【0009】

しかしながら、談話行為識別精度の向上のために用いられる、韻律特徴の従来の利用法は以下の理由で十分ではない。

まず、韻律特徴を発話文の特徴として抽出するためには、ユーザの一連の発話の中から、各発話文を正しく抜き出す必要がある。

しかし、文法や話し方のスタイルに制約のない自由発話においては、そもそも文の境界は不明瞭になりがちであり、しかも音声認識の間違いを含んだ発話を文の単位に正しく分節することは一般に難しい。

【0010】

例えば、「もう一回言ってくれない」は文の最後にかけてピッチを上げながら話せば依頼を意味するが、十分な長さの無音区間を置かずに直後に「もう一回」と念を押した場合、この発話のピッチは文末にかけて上昇するとは限らない。

もしも、両者を別個の文として分離できない場合は、依頼の発話であっても文末のピッチ上昇は観測できなくなってしまう。

【0011】

また、談話行為識別において、ある種の談話行為タグは、特定のフレーズに特定の韻律の変化を伴うことがしばしば観察される。

10

20

30

40

50

例えば、システムの発話をもう一度聞きたいという意図に対応する「言い直し要求」では、「言ってくれる?」、「言ってください?」のような特定のフレーズと同時にピッチの上昇が観察される場合が多い。

その場合、特徴の共起関係をより確実な識別の手がかりとすることで、より頑健かつ高精度な識別が期待できる。

【0012】

ところが、韻律特徴を言語特徴とは独立に発話の特徴としてモデル化する従来の技術では、言語特徴と韻律特徴の共起関係を捉えることができない。

このような問題点に鑑み、本発明は、言語特徴と韻律特徴の相関を直接モデル化し、もって談話行為識別精度の向上を図る技術を提案する。

【課題を解決するための手段】

【0013】

まず、入力されたユーザの音声から、音声認識に用いられる声道特徴量系列を抽出すると同時に、エネルギーや基本周波数など韻律に関する特徴量の時系列を抽出する。

【0014】

次に、得られた音声認識結果と韻律特徴時系列から、談話行為識別に用いる特徴系列を合成する。

【0015】

この合成特徴は、音声認識によって得られた単語 N - g r a m に対して、その時区間に
20 対応する離散的な韻律特徴を付与して得られる拡張単語 N - g r a m となっている。

【0016】

こうして得られた拡張単語 N - g r a m を入力として、言語特徴を利用した談話行為識別技術を適用し、識別モデルの学習やタグの予測を行う。

【発明の効果】

【0017】

このように、言語特徴と韻律特徴を合成した特徴を用いることで、両者の相関を考慮した談話行為識別が可能になるだけでなく、韻律特徴が、発話文ではなく、単語 N - g r a m に付与されていることで、発話文の正確な分節が得られない場合でも、韻律特徴を効果的に用いた談話行為識別が可能になる。

【図面の簡単な説明】

【0018】

【図1】音声対話システム装置構成をあらわす図である。

【図2】実験結果(談話行為タグ識別精度)をあらわす図である。

【図3】実験結果(情報要求の適合率)をあらわす図である。

【図4】実験結果(言い直し要求の適合率)をあらわす図である。

【発明を実施するための形態】

【0019】

次に、図1を参照して、本発明の音声対話装置の全体構成例を説明する。

【実施例1】

【0020】

音声入力部(P1)において、発話はマイクロホン等を用いてアナログ信号として取得された後、ただちにデジタル信号に変換される。

【0021】

声道特徴抽出部(P2)において、音声認識において用いられる、Mel Frequency Cepstral Coefficient(MFCC)等の声道特徴量が計算される。

【0022】

また、エネルギー抽出部(P3)において、フレーム毎のエネルギーが計算され、声道特徴系列と併せてエネルギー特徴系列が音声認識部(P5)に送られ、音声認識が行われ単語列に変換される。

この時、各単語の発話における時区間を表すタイムスタンプを同時に計算しておく。

10

20

30

40

50

また、ここで音声認識部の出力を単語列としているが、正確には当該音声認識システム(P5)の使用辞書に登録されている認識ユニットの列を意味しており、言語学的単語の列に限定されるものではない。

日本語のように単語に分かち書きされない言語の場合には、形態素解析のエラー等により、認識結果が正しく単語に分かち書きされない場合もあり、そのような場合には、音素やモーラのようなサブワードを認識ユニットとして用いて、サブワードユニットのN-gramを使って談話行為識別を行った方が良い場合もある。

また、識別に特徴的なフレーズ(単語列)がある場合は、フレーズを認識ユニットとして用いた方が良い場合もある。

本明細書では、典型的な認識ユニットである単語を用いて説明を行うが、認識ユニットは言語学的単語に限定されるわけではなく、単語N-gramは、認識ユニットのN-gramと読み替えることができる。

【0023】

声道特徴やエネルギーと並行して、基本周波数抽出部(P4)では、フレーム毎に基本周波数が計算される。

【0024】

次に、音声認識部(P5)で計算されたタイムスタンプ付きの単語列、およびタイムスタンプが付与されたエネルギー系列と基本周波数系列が言語・韻律特徴合成部(P6)に送られ、韻律特徴が付加された拡張単語N-gramが生成される。

【0025】

この拡張単語N-gramと、対話制御部(P8)が提供する発話のコンテキスト情報から、意図推定部(P7)において、発話の談話行為タグが推定される。

【0026】

引き続き、対話制御部(P8)では、談話行為タグが表す発話意図と対話文脈を考慮して、システム発話が生成され、引き続き音声合成部(P9)で音声信号に変換された後、スピーカー等の音声出力部(P10)を通して音声再生される。

【0027】

ここで、音声入力部(P1)、声道特徴抽出部(P2)、エネルギー抽出部(P3)、基本周波数抽出部(P4)、音声認識部(P5)、対話制御部(P8)、音声合成部(P9)、音声出力部(P10)には公知の技術を用いることができる(非特許文献6、非特許文献7)。

【0028】

以下では、言語・韻律特徴合成部(P6)および意図推定部(P7)についてのみ詳細に説明する。

【0029】

本発明の1つの実施形態で用いられる、談話行為識別のための基本的な原理は、談話行為タグを表す確率変数Iの事後確率の最大化である。

【0030】

【数1】

$$I^* = \underset{I}{\operatorname{argmax}} P(I|A) = \underset{I}{\operatorname{argmax}} P(A|I)P(I) \quad (1)$$

ここで、Aは音響信号を表す。

【0031】

音響信号Aが、声道成分A_sと韻律成分A_pに分離できると仮定すると、Aの尤度は以下のように書ける。

【0032】

10

20

30

40

【数 2】

$$P(A|I) = \sum_W P(AW|I) = \sum_W P(A_s A_p | WI) P(W|I) \quad (2)$$

ここでWは単語列を表す。

【0033】

声道成分と韻律成分が条件付き独立であると仮定すると、

【0034】

【数 3】

$$P(A|I) \approx \sum_W P(A_s | WI) P(A_p | WI) P(W|I) \quad (3)$$

10

【0035】

さらに、声道成分は、単語列のみに依存すると仮定すると、

【0036】

【数 4】

$$P(A|I) \approx \sum_W P(A_s | W) P(A_p | WI) P(W|I) \quad (4)$$

20

と書ける。

【0037】

非特許文献1では、ここからさらに、韻律成分は単語列に依存しないと仮定し、

【0038】

【数 5】

$$P(A|I) \approx \sum_W P(A_s | W) P(A_p | I) P(W|I) \quad (5)$$

30

というモデル化を行う。

【0039】

非特許文献2および非特許文献3においても、モデル化手法は異なるが、基本的に、韻律成分が単語列に依存せず、発話意図のみに依存するとしてモデル化を行っている。

本発明では、このような仮定(非特許文献1乃至非特許文献3)を置かず、数式(4)を、

【0040】

【数 6】

$$P(A|I) \approx \sum_W P(A_s | W) P(A_p | WI) P(W|I) = \sum_W P(A_s | W) P(A_p W | I) \quad (6)$$

40

と、単語列Wと韻律特徴A_pを同時にモデル化する。

【0041】

そのために、離散化された韻律特徴が付与された拡張単語列W^fを導入し、

【0042】

【数 7】

$$P(A|I) \approx \sum_W P(A_s|W)P(W^f|I) \quad (7)$$

とモデル化する。

【0043】

離散化された韻律特徴 f としては、例えば、 W に対応するフレーム列において、基本周波数の変化量が平均 + 標準偏差よりも大きい場合は + を、変化量が平均 - 標準偏差よりも小さい場合は - を、あるいは変化量が平均 ± 標準偏差の範囲内であれば 0 を付与する 3 値の離散化を用いることができる。 10

【0044】

あるいは、強勢であれば、エネルギーが平均 + 標準偏差よりも大きい場合は $s +$ を、平均 - 標準偏差よりも小さい場合は $s -$ を、あるいは平均 ± 標準偏差の範囲内であれば $s 0$ を付与する 3 値に離散化を用いることができる。

【0045】

例えば、「もう 1 回言って」という単語列に、強勢とピッチの上昇が観察されれば、「もう 1 回言って⁺ s^+ 」と単語列が拡張されることになる。

もちろん、離散化のやり方は、これ以外の方法を考えることもでき、離散化の粒度も 3 値に限るものではない。 20

【0046】

このような韻律特徴の離散化は、非特許文献 8 記載の、*t i l t* 特徴と基本的な考え方は同じであり、このような特徴を非特許文献 4 では音声合成に、非特許文献 9 では音声認識に用いているが、本発明では談話行為識別に用いる。

【0047】

また、非特許文献 3 では、韻律特徴量の軌跡を離散符号化し、符号の *N - g r a m* を談話行為識別に利用しているが、単語列との相関は考慮されていない。

【0048】

ところで、このような韻律特徴の離散化の際には、話者間の変動や発話環境の変動に対処するため、話者毎また発話環境毎に特徴量の正規化を行うことが望ましい。 30

例えば、同一話者の直近複数の発話を用いて、平均値や標準偏差を計算することができる。

【0049】

以上述べたような離散化を適用することで、言語・韻律特徴合成部 (P 6) において、音声認識により得られた単語列 $W = w_1, \dots, w_n$ は離散韻律特徴が付与された *N - g r a m* 列に拡張される。

【0050】

その過程をより詳細に述べる。まず、 W から *N - g r a m* ($N - 1$) を抽出する。

例えば、 $N = 2$ であれば、

【0051】

【数 8】 40

$$\langle s \rangle w_1, w_1 w_2, w_2 w_3, \dots, w_{n-1} w_n, w_n \langle /s \rangle \quad (8)$$

が抽出される。

ここで、 $\langle s \rangle$, $\langle /s \rangle$ はそれぞれ文頭、文末を表す記号である。

このとき、単語の一種として、短い無音区間を表す $w_i = \langle s p \rangle$ を含めれば、別種の韻律特徴を拡張 *N - g r a m* の中に取り込むことができる。

【0052】

次に、各単語に付与されたタイムスタンプに基づいて、各 N - g r a m 毎に、対応する時区間におけるエネルギー系列と基本周波数系列の部分区間を抽出し韻律特徴を計算する。

その際、欠損値があれば線形補完等で補い、当該時区間の平均や変化量等の統計量を計算し、前述の正規化を施した後に離散化し各 N - g r a m に付与される。

【 0 0 5 3 】

次に、意図推定部 (P 7) を説明する。

音声信号 A が所与のとき、数式 1 を最大化する談話行為タグ I を計算する。

【 0 0 5 4 】

事前確率 P (I) は、訓練データから予め計算した値を利用することができる。

10

このとき、発話のコンテキストを用いることでタグの識別精度が向上することが広く知られている。

例えば、非特許文献 5 では、直前のシステム発話の談話行為タグ C が所与であることを利用して、P (I) の代わりに P (I | C) を用いることでタグの識別精度を向上させている。

【 0 0 5 5 】

数式 1 の尤度 P (A | I) の計算には数式 7 を用いる。

数式 7 において P (A_s | W) は、音声認識で用いている音響モデルから計算される音響尤度を用いることができる。

つまり、音声認識部から出力される尤度上位 n 個の単語列を正しい認識の候補と考える場合、各候補 W_i (1 ≤ i ≤ n) の音響尤度 P (A_s | W_i) を重みとする P (W_i^f | I) の重みづけ和として P (A | I) が計算される。

20

【 0 0 5 6 】

この重みづけ和を計算する際、一般に音響尤度は非常に小さな値になるので、桁落ちを防ぐために、非特許文献 5 では、音響尤度を尤度の最大値 M で正規化して用いており、本発明でも有効な計算方法である。

【 0 0 5 7 】

【 数 9 】

$$P(A_S|W)/M, M = \max_i P(A_S|W_i) \quad (9)$$

30

【 0 0 5 8 】

尤度 P (W^f | I) は、N - g r a m でモデル化される。つまり、W に含まれる N - g r a m、W₁, …, W_m に対して離散韻律特徴を付与した、拡張 N - g r a m 列 W^f₁, …, W^f_m が条件付き独立であると仮定し、

【 0 0 5 9 】

【 数 1 0 】

$$P(W^f|I) = \prod_{i=1}^m P(W_i^f|I) \quad (10)$$

40

のように計算する。

ここで用いる、拡張 N - g r a m の尤度 P (W^f_i | I) は、予め訓練データから推定しておく。

【 0 0 6 0 】

ただし、離散韻律特徴で拡張されているので、訓練データが十分に多くない場合は、必ずしも拡張 N - g r a m W_i^f が訓練データに含まれない場合が想定される。

そのような場合は、以下のような平滑化を行うことが望ましい。

【 0 0 6 1 】

50

【数 1 1】

$$P(W_i^f | I) = \sum_{j=0}^J \lambda_j P(W_i^{f_j} | I), \quad \sum_{j=0}^J \lambda_j = 1 \quad (11)$$

ここで、 j は自然数であり、 $f_j (j > 0)$ は、基本周波数の勾配、強勢等、 f で用いられている離散韻律特徴である。

【0062】

特に、 $j = 0$ は、離散韻律特徴が付与されていない N -gram の生起確率であり、これ自身、ゼロ頻度問題に対処するために、Good-Turing 法など公知の方法(非特許文献 6)を用いて平滑化されているものとする。

10

【0063】

図 2 は、基本周波数 (F_0) の変化量を韻律特徴として用いた場合の、本発明による談話行為タグ識別率の向上を示している。

【0064】

図 2 で「 F_0 なし」として示されているのは、言語特徴のみを用いて識別した場合の識別精度を示している。

また、「 F_0 変化(発話単位)」として示されているものは、一つの発話における基本周波数の勾配を 3 値に離散化した韻律特徴を、言語特徴とは独立に用いた場合の識別精度を示している。「 F_0 変化(2 gram)」として示されているのは、単語 2-gram 毎に、3 値の離散韻律特徴を付与した拡張 2-gram 特徴を用いた場合の識別精度を示している。

20

【0065】

この結果から分かるように、韻律特徴を言語特徴と独立に付与しても、必ずしも識別精度は向上しない一方で、言語特徴と韻律特徴の相関を考慮する本発明によれば、識別率がおよそ 1% 向上していることが分かる。

【0066】

図 3 と図 4 は、韻律特徴が寄与すると予想される「情報要求」と「言い直し要求」の 2 つの談話行為タグの適合率、すなわち、それぞれのタグを予測した発話の中で、正しい予測の割合を示している。

30

タグ個別にみると、「情報要求」でおよそ 5%、「言い直し要求」で 2% 適合率が向上していることが分かる。

【0067】

本発明のシステムは、マイクロホンとスピーカーとパーソナルコンピュータを用い、図 1 に示した各処理部を実行するプログラムを C および Perl 言語で作成し、実行して確認した。

作成したプログラムは、上で述べたように、汎用計算機を用いた汎用的なプログラムであってもよいし、各種音声対話システム・装置・機器にのみ適合する固有のプログラムであってもよい。

また、プログラムは、内蔵式、埋め込み式(Imbedded)、読み込み式、ダウンロード方式、分散型、あるいはクラウドコンピューティングであってもよい。

40

【0068】

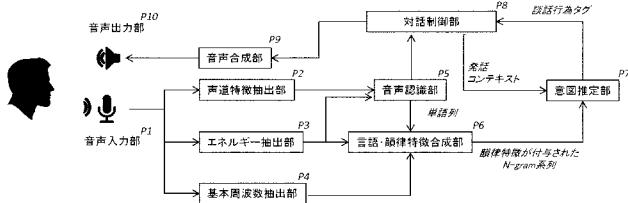
音声入力に使用したマイクロホンは、機器の一部として備わるマイクロホンであってもよいし、その設置場所は近接地・遠隔地を問わず、音声入力装置であれば足りる。

音声出力に使用したスピーカーは、機器の一部として備わるスピーカーや、イヤホンであってもよいし、その設置場所は近接地・遠隔地を問わず、音声出力装置であれば足りる。

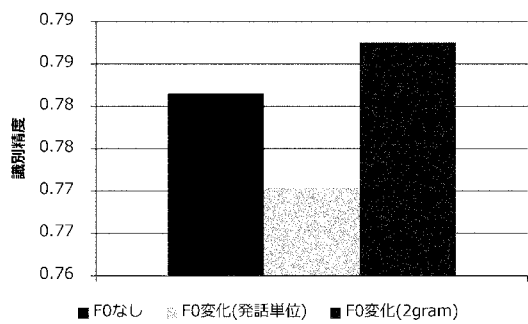
音声入力信号はアナログ音響信号だけでなく、本発明の内部処理に適してデジタル化された音響信号のいずれであってもよい。

50

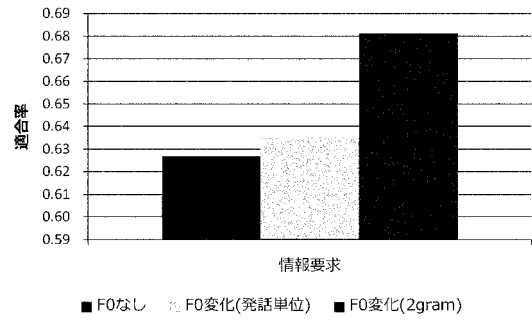
【 図 1 】



【 図 2 】



【 図 3 】



【 図 4 】

