

(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2020-34835
(P2020-34835A)

(43) 公開日 令和2年3月5日(2020.3.5)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 0 L 15/22 (2006.01)	G 1 0 L 15/22 3 0 0 U	5 E 5 5 5
G 1 0 L 13/00 (2006.01)	G 1 0 L 13/00 1 0 0 M	
G 1 0 L 15/10 (2006.01)	G 1 0 L 15/10 5 0 0 N	
G 0 6 F 3/16 (2006.01)	G 1 0 L 15/10 5 0 0 Z	
G 0 6 F 3/01 (2006.01)	G 0 6 F 3/16 6 5 0	

審査請求 未請求 請求項の数 11 O L (全 27 頁) 最終頁に続く

(21) 出願番号 特願2018-162774 (P2018-162774)
(22) 出願日 平成30年8月31日 (2018. 8. 31)

(71) 出願人 504132272
国立大学法人京都大学
京都府京都市左京区吉田本町 3 6 番地 1
(71) 出願人 000003207
トヨタ自動車株式会社
愛知県豊田市トヨタ町 1 番地
(74) 代理人 100103894
弁理士 家入 健
(72) 発明者 河原 達也
京都府京都市左京区吉田本町 3 6 番地 1
国立大学法人京都大学内
(72) 発明者 堀 達朗
愛知県豊田市トヨタ町 1 番地 トヨタ自動車株式会社内

最終頁に続く

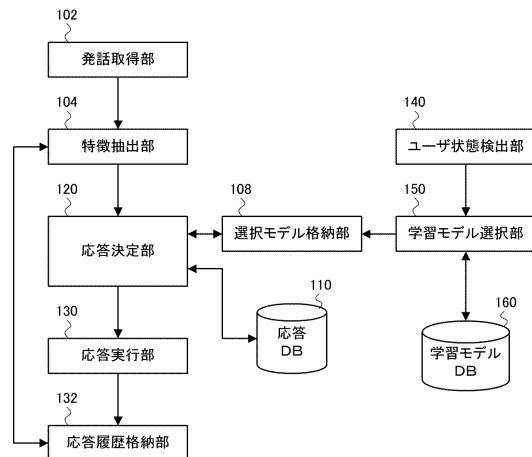
(54) 【発明の名称】 音声対話システム、音声対話方法、プログラム、学習モデル生成装置及び学習モデル生成方法

(57) 【要約】

【課題】 応答誤りが発生しないように適切に対処することが可能な音声対話システムを提供する。

【解決手段】 発話取得部 1 0 2 は、ユーザ発話を取得する。特徴抽出部 1 0 4 は、ユーザ発話の特徴を抽出する。応答決定部 1 2 0 は、複数の学習モデルのうちのいずれかを用いて、抽出された特徴ベクトルに応じた応答を決定する。応答実行部 1 3 0 は、決定された応答を実行する。ユーザ状態検出部 1 4 0 は、ユーザ状態を検出する。学習モデル選択部 1 5 0 は、検出されたユーザ状態に応じて、複数の学習モデルから学習モデルを選択する。応答決定部 1 2 0 は、選択された学習モデルを用いて、応答を決定する。

【選択図】 図 2



【特許請求の範囲】

【請求項 1】

ユーザと音声を用いた対話を行う音声対話システムであって、
 前記ユーザの発話であるユーザ発話を取得する発話取得部と、
 前記取得されたユーザ発話の特徴を少なくとも抽出する特徴抽出部と、
 予め機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、前記抽出された特徴に応じた応答を決定する応答決定部と、
 前記決定された応答を実行するための制御を行う応答実行部と、
 前記ユーザの状態であるユーザ状態を検出するユーザ状態検出部と、
 前記検出されたユーザ状態に応じて、前記複数の学習モデルから前記学習モデルを選択する学習モデル選択部と
 を有し、
 前記応答決定部は、前記学習モデル選択部によって選択された学習モデルを用いて、前記応答を決定する
 音声対話システム。

10

【請求項 2】

前記ユーザ状態検出部は、前記ユーザ状態として対話に対する前記ユーザの積極性の度合を検出し、
 前記学習モデル選択部は、前記ユーザの積極性の度合に対応する前記学習モデルを選択する
 請求項 1 に記載の音声対話システム。

20

【請求項 3】

前記ユーザ状態検出部は、予め定められた期間における前記ユーザの発話量、又は、前記期間において当該音声対話システムが応答として音声を出力した時間と前記ユーザが発話した時間との合計に対する前記ユーザが発話した時間の割合を検出し、
 前記学習モデル選択部は、前記ユーザの発話量又は前記ユーザが発話した時間の割合に対応する前記学習モデルを選択する
 請求項 2 に記載の音声対話システム。

【請求項 4】

前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの識別情報を検出し、
 前記学習モデル選択部は、前記ユーザの識別情報に対応する前記学習モデルを選択する
 請求項 1 に記載の音声対話システム。

30

【請求項 5】

前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの感情を検出し、
 前記学習モデル選択部は、前記ユーザの感情に対応する前記学習モデルを選択する
 請求項 1 に記載の音声対話システム。

【請求項 6】

前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの健康状態を検出し、
 前記学習モデル選択部は、前記ユーザの健康状態に対応する前記学習モデルを選択する
 請求項 1 に記載の音声対話システム。

40

【請求項 7】

前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの覚醒状態の度合を検出し、
 前記学習モデル選択部は、前記ユーザの覚醒状態の度合に対応する前記学習モデルを選択する
 請求項 1 に記載の音声対話システム。

【請求項 8】

ユーザと音声を用いた対話を行う音声対話システムを用いて行われる音声対話方法であって、
 前記ユーザの発話であるユーザ発話を取得し、

50

前記取得されたユーザ発話の特徴を少なくとも抽出し、
 予め機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、前記抽出された特徴に応じた応答を決定し、
 前記決定された応答を実行するための制御を行い、
 前記ユーザの状態であるユーザ状態を検出し、
 前記検出されたユーザ状態に応じて、前記複数の学習モデルから前記学習モデルを選択し、
 前記選択された学習モデルを用いて、前記応答を決定する
 音声対話方法。

【請求項 9】

ユーザと音声を用いた対話を行う音声対話システムを用いて行われる音声対話方法を実行するプログラムであって、
 前記ユーザの発話であるユーザ発話を取得するステップと、
 前記取得されたユーザ発話の特徴を少なくとも抽出するステップと、
 予め機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、前記抽出された特徴に応じた応答を決定するステップと、
 前記決定された応答を実行するための制御を行うステップと、
 前記ユーザの状態であるユーザ状態を検出するステップと、
 前記検出されたユーザ状態に応じて、前記複数の学習モデルから前記学習モデルを選択するステップと、
 前記選択された学習モデルを用いて、前記応答を決定するステップと
 をコンピュータに実行させるプログラム。

【請求項 10】

ユーザと音声を用いた対話を行う音声対話システムで用いられる学習モデルを生成する学習モデル生成装置であって、
 1以上の任意ユーザと対話を行うことによって前記任意ユーザの発話であるユーザ発話を取得する発話取得部と、
 前記取得されたユーザ発話の特徴を少なくとも示す特徴ベクトルを抽出する特徴抽出部と、
 前記ユーザ発話に対する応答を示す正解ラベルと前記特徴ベクトルとが対応付けられたサンプルデータを生成するサンプルデータ生成部と、
 前記ユーザ発話を発したときの前記任意ユーザの状態であるユーザ状態を取得して、前記取得されたユーザ状態を前記ユーザ発話に対応する前記サンプルデータに対応付けるユーザ状態取得部と、
 前記ユーザ状態ごとに前記サンプルデータを分類するサンプルデータ分類部と、
 前記分類された前記サンプルデータごとに、機械学習によって複数の学習モデルを生成する学習モデル生成部と
 を有する学習モデル生成装置。

【請求項 11】

ユーザと音声を用いた対話を行う音声対話システムで用いられる学習モデルを生成する学習モデル生成方法であって、
 1以上の任意ユーザと対話を行うことによって前記任意ユーザの発話であるユーザ発話を取得し、
 前記取得されたユーザ発話の特徴を少なくとも示す特徴ベクトルを抽出し、
 前記ユーザ発話に対する応答を示す正解ラベルと前記特徴ベクトルとが対応付けられたサンプルデータを生成し、
 前記ユーザ発話を発したときの前記任意ユーザの状態であるユーザ状態を取得して、前記取得されたユーザ状態を前記ユーザ発話に対応する前記サンプルデータに対応付け、
 前記ユーザ状態ごとに前記サンプルデータを分類し、
 前記分類された前記サンプルデータごとに、機械学習によって複数の学習モデルを生成

10

20

30

40

50

する

学習モデル生成方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声対話システム、音声対話方法、プログラム、学習モデル生成装置及び学習モデル生成方法に関し、特に、ユーザと音声を用いた対話を行うことが可能な音声対話システム、音声対話方法、プログラム、学習モデル生成装置及び学習モデル生成方法に関する。

【背景技術】

【0002】

ユーザが音声対話ロボット（音声対話システム）と日常会話を楽しむための技術が普及しつつある。この技術にかかる音声対話ロボットは、ユーザが発声した音声の音韻情報を解析して、解析結果に応じた応答を実行する。ここで、音声対話ロボットが学習モデルを用いて応答を決定することが、行われている。

【0003】

上記の技術に関連し、特許文献1は、ユーザの発する音声からそのユーザの感情状態を判断して適切な対応動作を実行する感情状態反応動作装置を開示する。特許文献1にかかる感情状態反応動作装置は、音声情報の音韻スペクトルに係る特徴量を抽出する音韻特徴量抽出手段と、音韻特徴量を入力して、予め備えた状態判断テーブルに基づき音声情報の感情状態を判断する状態判断手段と、感情状態を入力して、予め備えた対応動作選択テーブルに基づき対応動作処理を決定する対応動作選択手段とを有する。さらに、特許文献1にかかる感情状態反応動作装置は、感情状態学習用テーブルと感情状態学習手段を備え、感情状態学習手段は、感情状態学習テーブルに基づいて所定の機械学習モデルにより音韻特徴量と感情状態との関連を取得して状態判断テーブルに学習結果を保存し、状態判断手段は、状態判断テーブルに基づいて上記機械学習モデルによる感情状態の判断を行う。

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2005-352154号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

ユーザの状態（ユーザの違い、又は、ユーザの感情等）によっては、機械学習モデルが適切でないおそれがある。この場合、例えば、ユーザの発話と装置の発話とが衝突する発話衝突、又は、ユーザの発話と装置の発話との間の期間が長期に亘ってしまう長期沈黙といった、応答誤りが発生するおそれがある。一方、特許文献1にかかる技術では、1つの機械学習モデルを用いて対応動作処理を決定している。したがって、特許文献1にかかる技術では、学習モデルが適切でないために応答誤りが発生する場合であっても、応答誤りが発生しないように適切に対処することが困難である。

【0006】

本発明は、応答誤りが発生しないように適切に対処することが可能な音声対話システム、音声対話方法、プログラム、学習モデル生成装置及び学習モデル生成方法を提供するものである。

【課題を解決するための手段】

【0007】

本発明にかかる音声対話システムは、ユーザと音声を用いた対話を行う音声対話システムであって、前記ユーザの発話であるユーザ発話を取得する発話取得部と、前記取得されたユーザ発話の特徴を少なくとも抽出する特徴抽出部と、予め機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、前記抽出された特徴に応じた応答を決定

10

20

30

40

50

する応答決定部と、前記決定された応答を実行するための制御を行う応答実行部と、前記ユーザの状態であるユーザ状態を検出するユーザ状態検出部と、前記検出されたユーザ状態に応じて、前記複数の学習モデルから前記学習モデルを選択する学習モデル選択部とを有し、前記応答決定部は、前記学習モデル選択部によって選択された学習モデルを用いて、前記応答を決定する。

【0008】

また、本発明にかかる音声対話方法は、ユーザと音声を用いた対話を行う音声対話システムを用いて行われる音声対話方法であって、前記ユーザの発話であるユーザ発話を取得し、前記取得されたユーザ発話の特徴を少なくとも抽出し、予め機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、前記抽出された特徴に応じた応答を決定し、前記決定された応答を実行するための制御を行い、前記ユーザの状態であるユーザ状態を検出し、前記検出されたユーザ状態に応じて、前記複数の学習モデルから前記学習モデルを選択し、前記選択された学習モデルを用いて、前記応答を決定する。

10

【0009】

また、本発明にかかるプログラムは、ユーザと音声を用いた対話を行う音声対話システムを用いて行われる音声対話方法を実行するプログラムであって、前記ユーザの発話であるユーザ発話を取得するステップと、前記取得されたユーザ発話の特徴を少なくとも抽出するステップと、予め機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、前記抽出された特徴に応じた応答を決定するステップと、前記決定された応答を実行するための制御を行うステップと、前記ユーザの状態であるユーザ状態を検出するステップと、前記検出されたユーザ状態に応じて、前記複数の学習モデルから前記学習モデルを選択するステップと、前記選択された学習モデルを用いて、前記応答を決定するステップとをコンピュータに実行させる。

20

【0010】

応答誤りが発生する要因は、学習モデルが適切でないことが多い。本発明は、上記のように構成されているので、学習モデルが適切でない場合に、ユーザ状態に応じて適切な学習モデルに切り替えることができる。したがって、本発明は、応答誤りが発生しないように適切に対処することが可能となる。

【0011】

また、好ましくは、前記ユーザ状態検出部は、前記ユーザ状態として対話に対する前記ユーザの積極性の度合を検出し、前記学習モデル選択部は、前記ユーザの積極性の度合に対応する前記学習モデルを選択する。

30

本発明は、このように構成されていることによって、ユーザの対話に対する積極性の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの積極性に合わせて応答を実行することができる。

【0012】

また、好ましくは、前記ユーザ状態検出部は、予め定められた期間における前記ユーザの発話量、又は、前記期間において当該音声対話システムが応答として音声を出力した時間と前記ユーザが発話した時間との合計に対する前記ユーザが発話した時間の割合を検出し、前記学習モデル選択部は、前記ユーザの発話量又は前記ユーザが発話した時間の割合に対応する前記学習モデルを選択する。

40

本発明は、このように構成されていることによって、より正確に、ユーザの積極性の度合を判定することができる。

【0013】

また、好ましくは、前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの識別情報を検出し、前記学習モデル選択部は、前記ユーザの識別情報に対応する前記学習モデルを選択する。

本発明は、このように構成されていることによって、ユーザに適合した学習モデルを用いて対話を行うので、対話を行うユーザに合わせて応答を実行することができる。

【0014】

50

また、好ましくは、前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの感情を検出し、前記学習モデル選択部は、前記ユーザの感情に対応する前記学習モデルを選択する。

本発明は、このように構成されていることによって、ユーザの対話に対する感情の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの感情に合わせて応答を実行することができる。

【0015】

また、好ましくは、前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの健康状態を検出し、前記学習モデル選択部は、前記ユーザの健康状態に対応する前記学習モデルを選択する。

10

本発明は、このように構成されていることによって、ユーザの健康状態の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの健康状態に合わせて応答を実行することができる。

【0016】

また、好ましくは、前記ユーザ状態検出部は、前記ユーザ状態として前記ユーザの覚醒状態の度合を検出し、前記学習モデル選択部は、前記ユーザの覚醒状態の度合に対応する前記学習モデルを選択する。

本発明は、このように構成されていることによって、ユーザの覚醒状態の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの覚醒状態に合わせて応答を実行することができる。

20

【0017】

また、本発明にかかる学習モデル生成装置は、ユーザと音声を用いた対話を行う音声対話システムで用いられる学習モデルを生成する学習モデル生成装置であって、1以上の任意ユーザと対話を行うことによって前記任意ユーザの発話であるユーザ発話を取得する発話取得部と、前記取得されたユーザ発話の特徴を少なくとも示す特徴ベクトルを抽出する特徴抽出部と、前記ユーザ発話に対する応答を示す正解ラベルと前記特徴ベクトルとが対応付けられたサンプルデータを生成するサンプルデータ生成部と、前記ユーザ発話を発したときの前記任意ユーザの状態であるユーザ状態を取得して、前記取得されたユーザ状態を前記ユーザ発話に対応する前記サンプルデータに対応付けるユーザ状態取得部と、前記ユーザ状態ごとに前記サンプルデータを分類するサンプルデータ分類部と、前記分類された前記サンプルデータごとに、機械学習によって複数の学習モデルを生成する学習モデル生成部とを有する。

30

【0018】

また、本発明にかかる学習モデル生成方法は、ユーザと音声を用いた対話を行う音声対話システムで用いられる学習モデルを生成する学習モデル生成方法であって、1以上の任意ユーザと対話を行うことによって前記任意ユーザの発話であるユーザ発話を取得し、前記取得されたユーザ発話の特徴を少なくとも示す特徴ベクトルを抽出し、前記ユーザ発話に対する応答を示す正解ラベルと前記特徴ベクトルとが対応付けられたサンプルデータを生成し、前記ユーザ発話を発したときの前記任意ユーザの状態であるユーザ状態を取得して、前記取得されたユーザ状態を前記ユーザ発話に対応する前記サンプルデータに対応付け、前記ユーザ状態ごとに前記サンプルデータを分類し、前記分類された前記サンプルデータごとに、機械学習によって複数の学習モデルを生成する。

40

【0019】

本発明は、ユーザ状態ごとにサンプルデータを分類して機械学習によって複数の学習モデルを生成することによって、ユーザ状態に対応した複数の学習モデルを生成することができる。したがって、音声対話システムは、ユーザ状態に応じて学習モデルを選択することができる。

【発明の効果】

【0020】

本発明によれば、応答誤りが発生しないように適切に対処することが可能な音声対話シ

50

ステム、音声対話方法、プログラム、学習モデル生成装置及び学習モデル生成方法を提供できる。

【図面の簡単な説明】

【0021】

【図1】実施の形態1にかかる音声対話システムのハードウェア構成を示す図である。

【図2】実施の形態1にかかる音声対話システムの構成を示すブロック図である。

【図3】実施の形態1にかかる特徴抽出部によって生成される特徴ベクトルを例示する図である。

【図4】実施の形態1にかかる学習モデルの生成方法の概略を説明するための図である。

【図5】実施の形態1にかかる学習モデルの生成方法の概略を説明するための図である。

10

【図6】実施の形態1にかかる学習モデルの生成方法の概略を説明するための図である。

【図7】実施の形態1にかかる音声対話システムによってなされる音声対話方法を示すフローチャートである。

【図8】実施の形態1にかかる音声対話システムによってなされる音声対話方法を示すフローチャートである。

【図9】ユーザ状態がユーザの識別情報である場合における処理を示す図である。

【図10】ユーザ状態がユーザの対話に対する積極性の度合である場合における処理を示す図である。

【図11】積極性の度合を判定するためのテーブルを例示する図である。

【図12】ユーザ状態がユーザの感情である場合における処理を示す図である。

20

【図13】ユーザ状態がユーザの健康状態である場合における処理を示す図である。

【図14】ユーザ状態がユーザの覚醒状態の度合である場合における処理を示す図である。

【図15】実施の形態2にかかる音声対話システムの構成を示すブロック図である。

【図16】実施の形態2にかかる学習モデル生成装置の構成を示す図である。

【図17】実施の形態2にかかる学習モデル生成装置によって実行される学習モデル生成方法を示すフローチャートである。

【発明を実施するための形態】

【0022】

(実施の形態1)

30

以下、図面を参照して本発明の実施の形態について説明する。なお、各図面において、同一の要素には同一の符号が付されており、必要に応じて重複説明は省略されている。

【0023】

図1は、実施の形態1にかかる音声対話システム1のハードウェア構成を示す図である。音声対話システム1は、ユーザと音声を用いて対話を行う。具体的には、音声対話システム1は、ユーザからの発話（ユーザ発話）に応じて、ユーザに対して音声等の応答を実行することで、ユーザと対話を行う。音声対話システム1は、例えば、生活支援ロボット及び小型ロボット等のロボット、クラウドシステム及びスマートフォン等に搭載可能である。以下の説明では、音声対話システム1がロボットに搭載された例を示している。

【0024】

40

音声対話システム1は、周囲の音声を収集するマイク2と、音声を発するスピーカ4と、ユーザの状態を検出するために使用される検出装置6と、ロボットの首等を動作させるマニピュレータ8と、制御装置10とを有する。制御装置10は、例えばコンピュータとしての機能を有する。制御装置10は、マイク2、スピーカ4、検出装置6及びマニピュレータ8と、有線又は無線で接続されている。検出装置6は、例えば、カメラ及び生体センサの少なくとも1つを含む。生体センサは、例えば、血圧計、体温計、脈拍計等である。

【0025】

制御装置10は、主要なハードウェア構成として、CPU（Central Processing Unit）12と、ROM（Read Only Memory）14と、RAM（Random Access Memory）16と

50

、インタフェース部（ I F ; Interface ） 1 8 とを有する。 C P U 1 2、 R O M 1 4、 R A M 1 6 及びインタフェース部 1 8 は、データバスなどを介して相互に接続されている。

【 0 0 2 6 】

C P U 1 2 は、制御処理及び演算処理等を行う演算装置としての機能を有する。 R O M 1 4 は、 C P U 1 2 によって実行される制御プログラム及び演算プログラム等を記憶するための機能を有する。 R A M 1 6 は、処理データ等を一時的に記憶するための機能を有する。インタフェース部 1 8 は、有線又は無線を介して外部と信号の入出力を行う。また、インタフェース部 1 8 は、ユーザによるデータの入力を受け付け、ユーザに対して情報を表示する。

【 0 0 2 7 】

制御装置 1 0 は、マイク 2 によって集音されたユーザ発話を解析して、そのユーザ発話に応じて、ユーザに対する応答を決定して実行する。ここで、本実施の形態では、「応答」は、「沈黙」、「頷き」及び「発話」を含む。「沈黙」は、音声対話システム 1 が何もしないという動作である。「頷き」は、ロボットの首部を縦に振るという動作である。また、「発話」は、音声対話システム 1 が音声を出力するという動作である。決定された応答が「頷き」である場合、制御装置 1 0 は、マニピュレータ 8 を制御して、ロボットの首部を動作させる。また、決定された応答が「発話」である場合、制御装置 1 0 は、スピーカ 4 を介して、生成された応答に対応する音声（システム発話）を出力する。

【 0 0 2 8 】

図 2 は、実施の形態 1 にかかる音声対話システム 1 の構成を示すブロック図である。実施の形態 1 にかかる音声対話システム 1 は、発話取得部 1 0 2 と、特徴抽出部 1 0 4 と、選択モデル格納部 1 0 8 と、応答データベース 1 1 0（応答 D B ; Database）と、応答決定部 1 2 0 と、応答実行部 1 3 0 と、応答履歴格納部 1 3 2 とを有する。さらに、実施の形態 1 にかかる音声対話システム 1 は、ユーザ状態検出部 1 4 0 と、学習モデル選択部 1 5 0 と、学習モデルデータベース 1 6 0（学習モデル D B）とを有する。

【 0 0 2 9 】

図 2 に示す各構成要素は、マイク 2、スピーカ 4、マニピュレータ 8 及び制御装置 1 0 の少なくとも 1 つによって実現可能である。また、各構成要素の少なくとも 1 つは、例えば、 C P U 1 2 が R O M 1 4 に記憶されたプログラムを実行することによって実現可能である。また、必要なプログラムを任意の不揮発性記録媒体に記録しておき、必要に応じてインストールするようにしてもよい。なお、各構成要素は、上記のようにソフトウェアによって実現されることに限定されず、何らかの回路素子等のハードウェアによって実現されてもよい。さらに、図 2 に示す構成要素の全てが 1 つの装置に設けられている必要はなく、図 2 に示す構成要素の 1 つ以上は、他の構成要素とは物理的に別個の装置に設けられていてもよい。例えば、学習モデルデータベース 1 6 0 がサーバに設けられ、その他の構成要素が、サーバと通信可能な音声対話ロボット等に設けられるようにしてもよい。これらのことは、後述する他の実施の形態においても同様である。

【 0 0 3 0 】

発話取得部 1 0 2 は、マイク 2 を含み得る。発話取得部 1 0 2 は、ユーザ発話（及び音声対話システム 1 の発話）を取得する。具体的には、発話取得部 1 0 2 は、ユーザの発話（及び音声対話システム 1 の発話）を集音してデジタル信号に変換する。そして、発話取得部 1 0 2 は、ユーザ発話の音声データ（ユーザ音声データ）を、特徴抽出部 1 0 4 に対して出力する。

【 0 0 3 1 】

特徴抽出部 1 0 4 は、少なくともユーザ発話の特徴を抽出する。具体的には、特徴抽出部 1 0 4 は、ユーザ発話について、発話の具体的な意味内容を示す言語情報とは異なる非言語情報の解析を行う。また、特徴抽出部 1 0 4 は、非言語情報の解析結果である非言語情報解析結果として、後述する特徴ベクトルを生成する。そして、特徴抽出部 1 0 4 は、非言語情報解析結果（特徴ベクトル）を、応答決定部 1 2 0 に対して出力する。なお、特徴抽出部 1 0 4 は、ユーザ発話以外のユーザの特徴を抽出して特徴ベクトルを生成しても

10

20

30

40

50

よい。

【0032】

ここで、非言語情報とは、処理対象のユーザ発話の言語情報（文字列）とは異なる情報であり、ユーザ発話の韻律情報、及び、応答履歴情報の少なくとも一方を含む。韻律情報とは、ユーザ発話の音声波形の特徴を示す情報であり、例えば、基本周波数、音圧、周波数等の変化量、変動帯域、振幅の最大値及び平均値等である。また、応答履歴情報とは、応答決定部120によって決定（生成）され、応答実行部130によって実行された応答に関する過去の履歴を示す情報である。応答履歴格納部132は、応答実行部130によって応答が実行されると、この応答履歴情報を格納（更新）する。

【0033】

具体的には、特徴抽出部104は、発話取得部102によって取得されたユーザ音声データについて音声分析等を行って、音声波形から韻律情報を解析する。そして、特徴抽出部104は、韻律情報を示す特徴量を示す値を算出する。なお、特徴抽出部104は、ユーザ音声データについて、例えば32msごとに区切られたフレームごとに、基本周波数等を算出してもよい。また、特徴抽出部104は、応答履歴格納部132から応答履歴情報を抽出して、応答履歴の特徴を示す特徴量を算出する。

【0034】

なお、ユーザ発話の言語情報を用いた構文解析は、パターン認識等を用いるため、多大な時間を要することが多い。一方、非言語情報の解析（韻律情報の解析及び応答履歴情報の解析）については、解析に用いられるデータ量が構文解析と比較して少なく、演算手法が、構文解析と比較して単純である。したがって、非言語情報の解析に要する時間は、構文解析と比較してかなり短くなり得る。

【0035】

選択モデル格納部108は、後述する学習モデル選択部150によって選択された学習モデルを格納する。ここで、本実施の形態においては、学習モデル選択部150は、学習モデルデータベース160に格納された複数の学習モデルから、後述する方法によって、適切な学習モデルを選択する。なお、音声対話の開始前など、学習モデル選択部150によって学習モデルの選択処理がなされていない場合は、選択モデル格納部108は、任意の1つの学習モデルを格納してもよい。

【0036】

応答データベース110は、音声対話システム1が応答を行う際に必要なデータを格納する。例えば、応答データベース110は、応答が「発話」である場合のシステム発話を示す複数のシステム音声データを、予め記憶している。

【0037】

応答決定部120は、非言語情報解析結果（特徴ベクトル）に応じて、どの応答を実行するかを決定する。ここで、本実施の形態においては、応答決定部120は、予め、教師あり学習等の機械学習によって生成された複数の学習モデルのうちのいずれかを用いて、抽出された特徴（特徴ベクトル）に応じた応答を決定する。詳しくは後述する。

【0038】

本実施の形態においては、応答決定部120は、「沈黙」、「頷き」及び「発話」のうちの1つを、応答として決定する。応答決定部120は、決定された応答を示すデータ（応答データ）を、応答実行部130に対して出力する。なお、応答決定部120は、応答として「発話」を行うと決定した場合、応答データベース110に記憶された複数のシステム発話から、順番に、又はランダムに、システム発話（システム音声データ）を選択してもよい。応答決定部120は、選択されたシステム音声データを、応答実行部130に対して出力する。

【0039】

応答実行部130は、応答決定部120によって決定された応答を実行するための制御を行う。具体的には、応答決定部120から出力された応答データが「沈黙（沈黙応答）」を示す場合、応答実行部130は、スピーカ4及びマニピュレータ8を動作させないよ

10

20

30

40

50

うに制御する。また、応答決定部 120 から出力された応答データが「頷き（頷き応答）」を示す場合、応答実行部 130 は、マニピュレータ 8 を制御してロボットの首部を動作させる。また、応答決定部 120 から出力された応答データが「発話（発話応答）」を示す場合、応答実行部 130 は、スピーカ 4 を制御して、応答決定部 120 によって選択されたシステム音声データを示す音声を出力させる。

【0040】

応答履歴格納部 132 は、応答実行部 130 によって実行された応答を識別するデータを、応答履歴情報として格納する。さらに、応答履歴格納部 132 は、応答履歴情報として対話に関する時間を含む場合に、その時間を計測し、計測された時間を応答履歴情報として格納してもよい。

10

【0041】

図 3 は、実施の形態 1 にかかる特徴抽出部 104 によって生成される特徴ベクトルを例示する図である。なお、図 3 に例示する特徴ベクトルは、一例にすぎず、他の様々な特徴ベクトルが可能である。i 番目のユーザ発話についての特徴ベクトルを v_i とすると、特徴ベクトルの n 個の成分は、 $v_i = (v_{i1}, v_{i2}, \dots, v_{im-1}, v_{im}, v_{i(m+1)}, \dots, v_{in})$ と表される。ここで、i、n 及び m は整数（但し $n > m$ ）である。また、 $v_{i1} \sim v_{i(m-1)}$ が、i 番目のユーザ発話情報に関する韻律情報の解析結果に対応する。また、 $v_{im} \sim v_{in}$ が、応答履歴情報の解析結果に対応する。なお、 $v_{im} \sim v_{in}$ については、応答履歴格納部 132 に格納された情報そのものであってもよい。つまり、特徴抽出部 104 は、応答履歴情報については、応答履歴格納部 132 から応答履歴を抽出するのみでもよく、特別な解析を行わなくてもよい。

20

【0042】

図 3 に示す例では、例えば、 v_{i1} は、i 番目のユーザ発話の句末の $T1 \text{ msec}$ （ユーザ発話の終了時点から $T1 \text{ msec}$ （ T ミリ秒）遡った時間から終了時点までの期間）における基本周波数 f_0 （ f_{0T1} ）についてのパラメータを示す。また、 v_{i7} は、i 番目のユーザ発話の長さ（ユーザ発話長） $L1 [\text{sec}]$ を示す。なお、基本周波数 f_0 は、フレームごとに、SPTK（Speech Signal Processing Toolkit）の SWIPE（Saw-tooth Waveform Inspired Pitch Estimation）のロジックを用いて算出され得る。

【0043】

また、 v_{im} は、直前の応答タイプを示す。直前の応答タイプは、直前（i 番目のユーザ発話の直前）に、応答実行部 130 によって実行された応答のタイプ（「沈黙」、「頷き」、及び「発話」のいずれか）である。ここで、 v_{im} のような、数値ではない成分の成分値（特徴量）については、各タイプに数値が割り当てられている。例えば、 v_{im} において、成分値「1」は「沈黙」を示し、成分値「2」は「頷き」を示し、成分値「3」は「発話」を示す。

30

【0044】

ユーザ状態検出部 140（図 2）は、ユーザ発話を発したユーザの状態（ユーザ状態）を検出する。詳しくは後述する。ユーザ状態検出部 140 は、検出されたユーザ状態を、学習モデル選択部 150 に対して出力する。ここで、ユーザ状態とは、例えば、ユーザの識別情報、対話に対するユーザの積極性、ユーザの感情、ユーザの健康状態、又は、ユーザの覚醒状態であるが、これらに限定されない。なお、ユーザの感情とは、例えば、ユーザの喜怒哀楽、又は驚き等であるが、これらに限定されない。また、ユーザの健康状態とは、例えば、ユーザの脈拍、体温又は血圧等であるが、これらに限定されない。ユーザ状態検出部 140 は、カメラによって撮影されたユーザの画像、生体センサによって検出されたユーザの脈拍、体温若しくは血圧、又はマイク 2 によって集音されたユーザ音声を用いて、上述したようなユーザ状態を検出する。

40

【0045】

学習モデル選択部 150 は、ユーザ状態検出部 140 によって検出されたユーザ状態に応じて、学習モデルデータベース 160 に記憶された複数の学習モデルから学習モデルを選択する。詳しくは後述する。学習モデルデータベース 160 は、予め機械学習によって

50

生成された複数の学習モデルを記憶する。複数の学習モデルの生成方法の具体例については後述する。

【0046】

応答誤りが発生する要因は、学習モデルが適切でないことが多い。例えば、あるユーザにとって適切な学習モデルが、別のユーザにとっては適切でないことがある。また、同じユーザであっても、そのユーザの感情等の変化によって、適切であった学習モデルが適切でなくなることがある。ここで、学習モデルが適切でないとは、ユーザ発話に対する応答の精度が悪いことである。言い換えると、学習モデルが適切でないとは、対話のテンポ及びリズムが良好でないということである。ユーザによって、適した対話のテンポ及びリズムは異なり得るので、あるユーザにとって適切な学習モデルが、別のユーザにとっては適切でないことがある。また、同じユーザであっても、そのユーザの感情等の変化によって、適した対話のテンポ及びリズムは異なり得るので、適切であった学習モデルが適切でなくなることがある。学習モデルの応答の精度が悪いと、ロボットは、あるユーザ発話に対して「沈黙応答」を実行すべきときに「発話応答」を実行してしまい、又は、あるユーザ発話に対して「発話応答」を実行すべきときに「沈黙応答」を実行してしまう。

10

【0047】

これに対し、本実施の形態においては、ユーザ状態に応じて、学習モデルを適切なものに切り替えることができる。したがって、本実施の形態にかかる音声対話システム1は、応答誤りが発生しないように適切に対処することが可能となる。つまり、本実施の形態にかかる音声対話システム1は、応答の精度を良好にすることが可能となる。

20

【0048】

次に、学習モデルの生成方法の概略を説明する。

図4～図6は、実施の形態1にかかる学習モデルの生成方法の概略を説明するための図である。まず、学習モデルを生成するためのサンプルデータを取得する。図4で示すように、ユーザAの発話に対して、ロボット（音声対話システム1）が応答するといった、ユーザAとロボットとの対話によって、特徴ベクトルと正解ラベルとが対応付けられたサンプルデータ（教師データ）を収集する。このとき、オペレータは、ユーザAの発話に対して適切な応答をロボットが実行するように、ロボット（音声対話システム1）を操作する。

【0049】

図4に示す例では、正解ラベル「A」が、「沈黙応答」に対応する。正解ラベル「B」が、「頷き応答」に対応する。正解ラベル「C」が、「発話応答」に対応する。オペレータは、ユーザAの発話の途中では、沈黙応答が実行されるように、ロボットを操作する。このとき、オペレータは何もしなくてもよい。また、オペレータは、ユーザAの発話の読点レベルの切れ目で頷き応答が実行されるように、ロボットを操作する。このとき、ロボットは、オペレータの操作によって頷く動作を行う。また、オペレータは、ユーザAの発話の句点レベルの切れ目で発話応答が実行されるように、ロボットを操作する。このとき、ロボットは、オペレータの操作によって発話を行う。

30

【0050】

図4の例では、ユーザAの発話「結局」と「1人で聴いたよ。」との間には切れ目がないので、オペレータは、ユーザAの発話の途中であると判断し、沈黙応答が実行されるように、ロボットを操作する。また、ユーザAの発話「1人で聴いたよ。」が終了すると、句点レベルの切れ目があったと判断し、発話応答が実行されるようにロボットを操作する。このとき、ロボットは、発話「本当ですか」を出力する。

40

【0051】

さらに、ユーザAのユーザ状態が検出される。ユーザ状態は、例えばオペレータによって判断されてもよいし、上述したユーザ状態検出部140のような機能により自動的に検出されてもよい。これにより、ユーザ状態#1～#Nのいずれかが、サンプルデータに対応付けられる。ここで、Nは、2以上の整数であり、ユーザ状態の個数を示す。このNが、学習モデルの数に対応する。

50

【 0 0 5 2 】

図 5 は、図 4 の例によって取得された特徴ベクトルと正解ラベルとの組であるサンプルデータを例示する図である。ユーザ発話「結局」のユーザ発話長は 0.5 秒であったので、特徴ベクトルの成分（図 3 の v_{i7} ）に「0.5」が入力される。また、ユーザ発話「結局」に対する応答は「沈黙応答」であったので、ユーザ発話「結局」の特徴ベクトルには、正解ラベル「A」が対応付けられる。

【 0 0 5 3 】

また、ユーザ発話「1人で聞いたよ。」のユーザ発話長は 1.5 秒であったので、特徴ベクトルの成分（図 3 の v_{i7} ）に「1.5」が入力される。また、ユーザ発話「1人で聞いたよ。」に対する応答は「発話応答」であったので、ユーザ発話「1人で聞いたよ。」の特徴ベクトルには、正解ラベル「C」が対応付けられる。さらに、この一連のユーザ発話「結局1人で聞いたよ。」では、ユーザ状態（例えばユーザの識別情報）が「ユーザ状態 # 1」（例えば「ユーザ A」）であったので、ユーザ発話「結局1人で聞いたよ。」に対応するサンプルデータ群に、ユーザ状態 # 1 が対応付けられる。

10

【 0 0 5 4 】

図 6 は、分類されたサンプルデータ群から学習モデルが生成される態様を例示する図である。上記のようにして収集されたサンプルデータ群が、ユーザ状態 # 1 ~ # N ごとに、N 個のグループに分類される。ユーザ状態 # 1 のサンプルデータ群（例えば「ユーザ A」のユーザ発話に対応するサンプルデータ群）から、例えば教師あり学習等の機械学習によって、学習モデル # 1 が生成される。同様にして、ユーザ状態 # N のサンプルデータ群（例えば「ユーザ N」のユーザ発話に対応するサンプルデータ群）から、機械学習によって、学習モデル # N が生成される。学習モデル # 1 ~ # N の正解ラベル「A」、「B」、「C」の境界が互いに異なっているので、学習モデル # 1 ~ # N それぞれに同じ特徴ベクトルを入力した場合であっても、出力される応答は異なり得る。このように生成された複数の学習モデルが、学習モデルデータベース 160 に格納される。学習モデルを生成するために使用される機械学習の方法は、例えば、ランダムフォレストであってもよいし、サポートベクターマシン（SVM; Support Vector Machine）であってもよいし、ディープラーニングであってもよい。

20

【 0 0 5 5 】

なお、学習モデルは、ユーザ状態に応じて、正解ラベル「A」、「B」、「C」の境界が定められている。例えば、ユーザ状態が「積極性」である場合、積極性の度合いが大きなユーザ状態に対応する学習モデルであるほど、「発話」が選択される確率が低くなり得、「沈黙」が選択される確率が高くなり得る。これは、対話に対する積極性の度合いが大きいということは、ユーザが積極的に発話する傾向にあるということであるので、発話衝突を抑制するため、音声対話システム 1 はあまり発話しないようにするということである。逆に、対話に対する積極性の度合いが小さい場合には、ユーザがあまり積極的に発話しない傾向にあるということであるので、長期沈黙を抑制するため、音声対話システム 1 がより発話するようにする。

30

【 0 0 5 6 】

また、学習モデルは、上述したユーザ状態検出部 140 によって検出されるユーザ状態に対応している。例えば、ユーザ状態検出部 140 がユーザ状態として「積極性の度合い」を検出する場合、学習モデルは、積極性の度合いごとに、複数設けられている。また、ユーザ状態検出部 140 がユーザ状態として「ユーザの識別情報」を検出する場合、学習モデルは、ユーザの識別情報（ユーザ A、ユーザ B、・・・、ユーザ N 等）ごとに、複数設けられている。

40

【 0 0 5 7 】

図 7 及び図 8 は、実施の形態 1 にかかる音声対話システム 1 によってなされる音声対話方法を示すフローチャートである。まず、発話取得部 102 は、上述したようにユーザ発話を取得する（ステップ S102）。特徴抽出部 104 は、上述したように、取得されたユーザ発話について非言語情報（韻律情報及び応答履歴情報）の解析を行って、ユーザ発

50

話の特徴（特徴ベクトル）を抽出する（ステップ S 1 0 4）。

【 0 0 5 8 】

次に、応答決定部 1 2 0 は、現在の学習モデル（選択モデル格納部 1 0 8 に格納された学習モデル）を用いて、抽出された特徴ベクトルに応じた、ユーザ発話に対する応答を決定する（ステップ S 1 1 0）。応答実行部 1 3 0 は、上述したように、S 1 1 0 で決定された応答を実行する（ステップ S 1 2 0）。

【 0 0 5 9 】

図 8 は、S 1 1 0 の処理を示すフローチャートである。応答決定部 1 2 0 は、抽出された特徴ベクトルを、学習モデルに入力する（ステップ S 1 1 2）。応答決定部 1 2 0 は、学習モデルの出力を判定する（ステップ S 1 1 4）。

10

【 0 0 6 0 】

出力が「沈黙応答」である場合（S 1 1 4 の「沈黙」）、応答決定部 1 2 0 は、沈黙応答を実行すると決定する（ステップ S 1 1 6 A）。つまり、応答決定部 1 2 0 は、その特徴ベクトルに対応するユーザ発話に対して、何もしないと決定する。また、出力が「頷き応答」である場合（S 1 1 4 の「頷き」）、応答決定部 1 2 0 は、頷き応答を実行すると決定する（ステップ S 1 1 6 B）。つまり、応答決定部 1 2 0 は、その特徴ベクトルに対応するユーザ発話に対して、ロボットの首部を縦に振るようにマニピュレータ 8 を動作させると決定する。また、出力が「発話応答」である場合（S 1 1 4 の「発話」）、応答決定部 1 2 0 は、発話応答を実行すると決定する（ステップ S 1 1 6 C）。つまり、応答決定部 1 2 0 は、その特徴ベクトルに対応するユーザ発話に対して、システム発話を出力させるようにスピーカ 4 を動作させると決定する。

20

【 0 0 6 1 】

次に、ユーザ状態検出部 1 4 0 は、上述したように、ユーザ状態を検出する（ステップ S 1 3 0）。学習モデル選択部 1 5 0 は、S 1 3 0 の処理で検出されたユーザ状態に対応する学習モデルを選択する（ステップ S 1 4 0）。具体的には、現在の学習モデルが、検出されたユーザ状態に対応するものと異なる場合、学習モデル選択部 1 5 0 は、現在の学習モデルを、検出されたユーザ状態に対応する学習モデルに切り替える。一方、現在の学習モデルが、検出されたユーザ状態に対応するものである場合、学習モデル選択部 1 5 0 は、学習モデルを変更しない。このように、実施の形態 1 にかかる学習モデル選択部 1 5 0 は、ユーザ状態に応じた新たな学習モデルを選択するように構成されているので、応答の精度がより良くなる学習モデルを選択することが可能となる。

30

【 0 0 6 2 】

以下、ユーザ状態の具体例を説明する。第 1 の例は、ユーザ状態がユーザの識別情報である場合の例である。第 2 の例は、ユーザ状態が対話に対するユーザの積極性の度合である場合の例である。第 3 の例は、ユーザ状態がユーザの感情の度合である場合の例である。第 4 の例は、ユーザ状態がユーザの健康状態の度合である場合の例である。第 5 の例は、ユーザ状態がユーザの覚醒状態の度合である場合の例である。

【 0 0 6 3 】

（ユーザ状態の第 1 の例）

図 9 は、ユーザ状態がユーザの識別情報である場合における処理を示す図である。図 9 は、ユーザ状態がユーザの識別情報である場合における、S 1 3 0、S 1 4 0（図 7）の具体的な処理を示す。ユーザ状態検出部 1 4 0 は、カメラである検出装置 6 から、ユーザの画像を取得する（ステップ S 1 3 2 A）。なお、「画像」とは、情報処理の対象としての、画像を示す画像データをも意味し得る（以下の説明において同様）。

40

【 0 0 6 4 】

ユーザ状態検出部 1 4 0 は、画像に対して顔認識処理を行って、ユーザの識別情報を検出する（ステップ S 1 3 4 A）。具体的には、ユーザ状態検出部 1 4 0 は、例えば、画像の中からユーザの顔領域を決定し、顔特徴点の検出を行って、目、鼻、口端などの顔の特徴点位置を判定する。そして、ユーザ状態検出部 1 4 0 は、特徴点位置を用いて顔領域の位置及び大きさを正規化した後、予め登録された人物の画像との顔照合処理を行う。これ

50

により、ユーザ状態検出部 140 は、照合された人物の識別情報を取得する。

【0065】

次に、学習モデル選択部 150 は、検出された識別情報に対応する学習モデルを選択する（ステップ S142A）。なお、予め、ユーザの識別情報ごとに、複数の学習モデルが学習モデルデータベース 160 に格納されているとする。例えば、ユーザ状態検出部 140 によって「ユーザ A」の識別情報が検出された場合、学習モデル選択部 150 は、「ユーザ A」に対応する学習モデルを選択する。

【0066】

このようにして、第 1 の例にかかる音声対話システム 1 は、ユーザに適合した学習モデルを用いて対話を行うので、対話を行うユーザに合わせて応答を実行することができる。したがって、第 1 の例にかかる音声対話システム 1 は、応答誤りが発生しないように適切に対処することが可能となる。また、応答誤りが発生するということは、現在の対話のテンポ又はリズムが、そのユーザに適していないということである。第 1 の例にかかる音声対話システム 1 は、ユーザに対応する学習モデルを選択することによって、対話のテンポ又はリズムをそのユーザに適したものにすることが可能となる。

【0067】

また、第 1 の例においては、学習モデルを生成する際に、ユーザ状態としてユーザの識別情報に対応付けられる。言い換えると、ユーザの識別情報ごとに、複数の学習モデルが生成される。学習モデルを生成する際には、例えばオペレータが、ユーザの識別情報を入力することで、サンプルデータとユーザの識別情報とが対応付けられる。これにより、ユーザの識別情報ごとに、サンプルデータが分類され、分類されたサンプルデータを用いて、機械学習によって複数の学習データが生成される。したがって、例えば、ユーザ A に対応する学習モデル、ユーザ B に対応する学習モデル、及び、ユーザ C に対応する学習モデルが生成されることとなる。

【0068】

なお、上述した例では、画像を用いた顔認識処理によってユーザを識別するとしたが、ユーザを識別する方法は、この方法に限られない。ユーザ発話に対して話者認識処理を行うことによって、ユーザ発話を発したユーザを識別してもよい。さらに、ユーザの識別情報（ID）を入力することによって、ユーザを識別してもよい。

【0069】

（ユーザ状態の第 2 の例）

図 10 は、ユーザ状態がユーザの対話に対する積極性の度合である場合における処理を示す図である。図 10 は、ユーザ状態がユーザの積極性の度合である場合における、S130、S140（図 7）の具体的な処理を示す。ユーザ状態検出部 140 は、過去 T 分間におけるユーザ発話割合 R_s を取得する（ステップ S132B）。ここで、T は、予め定められた期間を示す。例えば、 $T=1$ [分] であるが、これに限定されない。「過去 T 分間」とは、現在から T 分間遡った時刻から現在までの期間である。ユーザ発話割合 R_s は、過去 T 分間における、音声対話システム 1 が応答として音声を出力した時間 t_r [分] とユーザ発話した時間 t_u [分] との合計 $t_u + t_r$ [分] に対するユーザが発話した時間 t_u の割合である。つまり、 $R_s [\%] = 100 * t_u / (t_u + t_r)$ である。

【0070】

ユーザ状態検出部 140 は、ユーザ発話割合 R_s に対応する積極性の度合を検出する（ステップ S134B）。具体的には、ユーザ状態検出部 140 は、図 11 に例示するテーブルを、予め記憶している。ユーザ状態検出部 140 は、このテーブルを用いて、ユーザ発話割合 R_s が積極性のどの段階に対応するのかを判定する。

【0071】

図 11 は、積極性の度合を判定するためのテーブルを例示する図である。図 11 に例示したテーブルでは、積極性の度合とユーザ発話割合 R_s とが対応付けられている。図 11 の例では、積極性の度合が、#1 ~ #4 の 4 つの段階で定められている。度合 #1 から度合 #4 にかけて、積極性の度合が大きくなる。ユーザ状態検出部 140 は、取得されたユ

10

20

30

40

50

ーザ発話割合 R_s が、度合 # 1 ~ # 4 のどの度合に対応するのかを判定する。例えば、 $R_s = 20$ [%] である場合、ユーザ状態検出部 140 は、積極性の度合を # 1 と判定する。また、 $R_s = 80$ [%] である場合、ユーザ状態検出部 140 は、積極性の度合を # 4 と判定する。

【0072】

次に、学習モデル選択部 150 は、検出された積極性の度合に対応する学習モデルを選択する（ステップ S142B）。なお、予め、ユーザの積極性の度合ごとに、複数の学習モデルが学習モデルデータベース 160 に格納されているとする。例えば、ユーザ状態検出部 140 によって「積極性の度合 # 1」が検出された場合、学習モデル選択部 150 は、「積極性の度合 # 1」に対応する学習モデルを選択する。また、ユーザ状態検出部 140 によって「積極性の度合 # 4」が検出された場合、学習モデル選択部 150 は、「積極性の度合 # 4」に対応する学習モデルを選択する。

10

【0073】

なお、上述した説明では、ユーザ発話割合に応じてユーザの積極性の度合を判定としたが、ユーザの発話量に応じてユーザの積極性の度合を判定してもよい。具体的には、ユーザ状態検出部 140 は、過去 T 分間におけるユーザ発話量 [分] を取得する（S132B）。ユーザ状態検出部 140 は、ユーザ発話量に対応する積極性の度合を検出する（S134B）。この場合、図 11 に例示したテーブルと同様に、ユーザ状態検出部 140 は、ユーザ発話量と積極性の度合（段階）とが対応付けられたテーブルを記憶していてもよい。ユーザ状態検出部 140 は、このテーブルを用いて、ユーザ発話量が積極性のどの段階に対応するのかを判定し得る。

20

【0074】

このようにして、第 2 の例にかかる音声対話システム 1 は、ユーザの対話に対する積極性の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの積極性に合わせて応答を実行することができる。したがって、第 2 の例にかかる音声対話システム 1 は、応答誤りが発生しないように適切に対処することが可能となる。また、応答誤りが発生するということは、現在の対話のテンポ又はリズムが、ユーザの積極性の度合に適していないということである。第 2 の例にかかる音声対話システム 1 は、ユーザの積極性の度合に対応する学習モデルを選択することによって、対話のテンポ又はリズムをユーザの積極性の度合に適したものすることが可能となる。また、ユーザの対話に対する積極性の度合は、対話の話題等によって変化し得る。第 2 の例にかかる音声対話システム 1 は、積極性の度合の変化に応じて学習モデルを変更することができる。

30

【0075】

また、第 2 の例においては、学習モデルを生成する際に、ユーザ状態としてユーザの積極性の度合が対応付けられる。言い換えると、積極性の度合ごとに、複数の学習モデルが生成される。学習モデルを生成する際には、例えばオペレータが、対話中のユーザの積極性の度合を入力することで、サンプルデータとユーザの積極性の度合とが対応付けられる。また、学習モデルの生成の際でも、図 10 に示したように、ユーザ発話割合又はユーザ発話量を用いて対話中のユーザの積極性の度合が判定されてもよい。この場合、オペレータが、期間 T を適宜設定してもよい。例えば、対話の話題が変更されたときに、ユーザの積極性の度合を計算するようにしてもよい。

40

【0076】

これにより、ユーザの積極性の度合ごとに、サンプルデータが分類され、分類されたサンプルデータを用いて、機械学習によって複数の学習データが生成される。したがって、例えば、積極性の度合 # 1 に対応する学習モデル、積極性の度合 # 2 に対応する学習モデル、積極性の度合 # 3 に対応する学習モデル、及び、積極性の度合 # 4 に対応する学習モデルが生成されることとなる。

【0077】

なお、上述したように、学習モデルは、ユーザ状態に応じて、正解ラベル「A（沈黙）」、「B（頷き）」、「C（発話）」の境界が定められている。ユーザ状態が「積極性」

50

である場合、積極性の度合いが大きなユーザ状態に対応する学習モデルであるほど、「発話」が選択される確率が低くなり得、「沈黙」が選択される確率が高くなり得る。つまり、学習モデル#1(度合#1)で「A(沈黙)」が選択される確率よりも、学習モデル#4(度合#4)で「A(沈黙)」が選択される確率の方が高くなるように、学習モデルが生成される。これにより、積極性の度合いが大きなユーザとの対話において、発話衝突を抑制するため、音声対話システム1はあまり発話しないようにすることが可能となる。また、積極性の度合いが小さいユーザとの対話において、長期沈黙を抑制するため、音声対話システム1はシステム発話を多くするようにすることが可能となる。

【0078】

なお、上述した例では、ユーザ発話割合又はユーザ発話量を用いてユーザの対話に対する積極性の度合いを検出するとしたが、ユーザの積極性の度合いを検出する方法は、この方法に限られない。例えば、ユーザ状態検出部140は、ユーザの画像を取得することで、積極性の度合いを検出してもよい。具体的には、ユーザ状態検出部140は、ユーザの顔画像に示されたユーザの表情及び視線を解析してユーザの積極性を判定し、積極性の度合いを数値化してもよい。また、例えば、ユーザ状態検出部140は、ユーザ発話を取得することで、積極性の度合いを検出してもよい。具体的には、ユーザ状態検出部140は、ユーザ発話の韻律を解析してユーザの積極性を判定し、積極性の度合いを数値化してもよい。しかしながら、上述したように、ユーザ発話割合又はユーザ発話量を用いて積極性の度合いを判定することにより、より正確に、ユーザの積極性の度合いを判定することができる。したがって、ユーザ発話割合又はユーザ発話量を用いることにより、第2の例にかかる音声対話システム1は、応答誤りが発生しないように、より適切に対処することが可能となる。

【0079】

(ユーザ状態の第3の例)

図12は、ユーザ状態がユーザの感情である場合における処理を示す図である。図12は、ユーザ状態がユーザの感情の度合いである場合における、S130、S140(図7)の具体的な処理を示す。「感情の度合い」とは、例えば「喜び」の度合いである。しかしながら、「感情の度合い」は、怒りの度合い、悲しみの度合い、又は驚きの度合いであってもよい。

【0080】

ユーザ状態検出部140は、カメラである検出装置6から、ユーザの顔画像を取得する(ステップS132C)。ユーザ状態検出部140は、顔画像を解析して、表情及び視線等から、ユーザの感情(喜び)の度合いを検出する(ステップS134C)。例えば、ユーザ状態検出部140は、「Affdex」又は「Emotion API」等の人工知能を用いて、ユーザの感情(喜び)を数値化してもよい。そして、ユーザ状態検出部140は、図11に例示したような、感情を示す数値と感情の度合いとを対応付けたテーブルを用いて、感情の度合いを検出してもよい。

【0081】

次に、学習モデル選択部150は、検出された感情(喜び)の度合いに対応する学習モデルを選択する(ステップS142C)。なお、予め、ユーザの感情の度合いごとに、複数の学習モデルが学習モデルデータベース160に格納されているとする。例えば、ユーザ状態検出部140によって「感情(喜び)の度合い#1」が検出された場合、学習モデル選択部150は、「感情(喜び)の度合い#1」に対応する学習モデルを選択する。また、ユーザ状態検出部140によって「感情(喜び)の度合い#4」が検出された場合、学習モデル選択部150は、「感情(喜び)の度合い#4」に対応する学習モデルを選択する。

【0082】

このようにして、第3の例にかかる音声対話システム1は、ユーザの対話に対する感情の度合いに適した学習モデルを用いて対話を行うので、対話を行うユーザの感情に合わせて応答を実行することができる。したがって、第3の例にかかる音声対話システム1は、応答誤りが発生しないように適切に対処することが可能となる。また、応答誤りが発生するということは、現在の対話のテンポ又はリズムが、ユーザの感情の度合いに適していないということである。第3の例にかかる音声対話システム1は、ユーザの感情の度合いに対応

10

20

30

40

50

する学習モデルを選択することによって、対話のテンポ又はリズムをユーザの感情の度合に適したものにすることが可能となる。また、ユーザの対話に対する感情の度合は、対話の話題等によって変化し得る。第3の例にかかる音声対話システム1は、感情の度合の変化に応じて学習モデルを変更することができる。

【0083】

また、第3の例においては、学習モデルを生成する際に、ユーザ状態としてユーザの感情の度合が対応付けられる。言い換えると、感情の度合ごとに、複数の学習モデルが生成される。学習モデルを生成する際には、例えばオペレータが、対話中のユーザの感情の度合を入力することで、サンプルデータとユーザの感情の度合とが対応付けられる。また、学習モデルの生成の際でも、ユーザの顔画像を用いて対話中のユーザの感情の度合が判定

10

【0084】

これにより、ユーザの感情の度合ごとに、サンプルデータが分類され、分類されたサンプルデータを用いて、機械学習によって複数の学習データが生成される。したがって、例えば、感情の度合#1に対応する学習モデル、感情の度合#2に対応する学習モデル、感情の度合#3に対応する学習モデル、及び、感情の度合#4に対応する学習モデルが生成されることとなる。

【0085】

なお、上述した例では、ユーザの顔画像を用いてユーザの感情の度合を検出するとしたが、ユーザの感情の度合を検出する方法は、この方法に限られない。例えば、ユーザ状態検出部140は、ユーザ発話を取得することで、感情の度合を検出してもよい。具体的には、ユーザ状態検出部140は、ユーザ発話の韻律を解析してユーザの感情を判定し、感情の度合を数値化してもよい。

20

【0086】

また、上述した例では、ユーザ状態検出部140が感情の度合を検出するとした。しかしながら、ユーザ状態検出部140は、感情の種類、つまり、喜び、悲しみ、怒り、驚き等を検出してもよい。具体的には、ユーザ状態検出部140は、喜び、悲しみ、怒り、驚きのそれぞれを示す数値を検出する。そして、ユーザ状態検出部140は、これらの数値のうち最も大きな値に対応する感情（例えば「怒り」）を、ユーザの感情として検出してもよい。この場合、学習モデルが、感情の種類ごとに複数設けられている。そして、学習

30

【0087】

また、ユーザ状態検出部140は、感情の種類ごとに度合を検出してもよい。この場合、学習モデルデータベース160は、例えば、怒りの度合がX1であり驚きの度合がY1である場合の学習モデル、怒りの度合がX1であり驚きの度合がY2である場合の学習モデル、怒りの度合がX2であり驚きの度合がY1である場合の学習モデル、怒りの度合がX2であり驚きの度合がY2である場合の学習モデルを格納してもよい。そして、学習

40

【0088】

(ユーザ状態の第4の例)

図13は、ユーザ状態がユーザの健康状態である場合における処理を示す図である。図13は、ユーザ状態がユーザの健康状態の度合である場合における、S130、S140(図7)の具体的な処理を示す。「健康状態の度合」とは、例えば、心拍数の度合である。しかしながら、「健康状態の度合」は、血圧の度合又は体温の度合であってもよい。

【0089】

ユーザ状態検出部140は、生体センサである検出装置6から、ユーザの生体系パラメータを取得する(ステップS132D)。生体系パラメータは、例えば心拍数である。ユーザ状態検出部140は、生体系パラメータから、ユーザの健康状態の度合を検出する(

50

ステップ S 1 3 4 D)。例えば、ユーザ状態検出部 1 4 0 は、図 1 1 に例示したような、健康状態を示す数値（心拍数）と健康状態の度合とを対応付けたテーブルを用いて、健康状態の度合を検出してもよい。

【 0 0 9 0 】

次に、学習モデル選択部 1 5 0 は、検出された健康状態（心拍数）の度合に対応する学習モデルを選択する（ステップ S 1 4 2 D）。なお、予め、ユーザの健康状態の度合ごとに、複数の学習モデルが学習モデルデータベース 1 6 0 に格納されているとする。例えば、ユーザ状態検出部 1 4 0 によって「健康状態（心拍数）の度合 # 1」が検出された場合、学習モデル選択部 1 5 0 は、「健康状態（心拍数）の度合 # 1」に対応する学習モデルを選択する。また、ユーザ状態検出部 1 4 0 によって「健康状態（心拍数）の度合 # 4」が検出された場合、学習モデル選択部 1 5 0 は、「健康状態（心拍数）の度合 # 4」に対応する学習モデルを選択する。

10

【 0 0 9 1 】

このようにして、第 4 の例にかかる音声対話システム 1 は、ユーザの健康状態の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの健康状態に合わせて応答を実行することができる。したがって、第 4 の例にかかる音声対話システム 1 は、応答誤りが発生しないように適切に対処することが可能となる。また、応答誤りが発生するということは、現在の対話のテンポ又はリズムが、ユーザの健康状態の度合に適していないということである。第 4 の例にかかる音声対話システム 1 は、ユーザの健康状態の度合に対応する学習モデルを選択することによって、対話のテンポ又はリズムをユーザの健康状態の度合に適したものすることが可能となる。また、ユーザの心拍数等の度合は、対話の話題等によって変化し得る。第 4 の例にかかる音声対話システム 1 は、心拍数等の健康状態の度合の変化に応じて学習モデルを変更することができる。

20

【 0 0 9 2 】

また、第 4 の例においては、学習モデルを生成する際に、ユーザ状態としてユーザの健康状態の度合が対応付けられる。言い換えると、健康状態の度合ごとに、複数の学習モデルが生成される。学習モデルを生成する際には、例えば、生体センサを用いて対話中のユーザの健康状態の度合を入力することで、サンプルデータとユーザの健康状態の度合とが対応付けられる。

【 0 0 9 3 】

これにより、ユーザの健康状態の度合ごとに、サンプルデータが分類され、分類されたサンプルデータを用いて、機械学習によって複数の学習データが生成される。したがって、例えば、健康状態の度合 # 1 に対応する学習モデル、健康状態の度合 # 2 に対応する学習モデル、健康状態の度合 # 3 に対応する学習モデル、及び、健康状態の度合 # 4 に対応する学習モデルが生成されることとなる。

30

【 0 0 9 4 】

なお、上述した例では、生体センサを用いてユーザの健康状態の度合を検出するとしたが、ユーザの健康状態の度合を検出する方法は、この方法に限られない。例えば、ユーザ状態検出部 1 4 0 は、カメラである検出装置 6 からユーザの顔画像を取得することで、ユーザの健康状態の度合を検出してもよい。この場合、ユーザ状態検出部 1 4 0 は、顔画像を解析してユーザの顔色（赤色、青色、白色、黄色、黒色）を検出してもよい。そして、ユーザ状態検出部 1 4 0 は、ユーザの顔色が赤色、青色、白色、黄色、黒色のどの色に近いかに応じて、健康状態を検出してもよい。この場合、ユーザの顔色ごとに、複数の学習モデルが格納されている。

40

【 0 0 9 5 】

また、ユーザ状態検出部 1 4 0 は、複数の生体系パラメータ（心拍数、血圧及び体温から、ユーザの健康状態が良好であるか劣悪であるか、又は、ユーザの疲労度を判定してもよい。また、ユーザ状態検出部 1 4 0 は、心拍数、血圧及び体温それぞれが予め定められた正常範囲にあるか否かを判定し、正常範囲を逸脱した生体系パラメータの数に応じて、健康状態の度合（健康状態が良好か劣悪かの度合）を判定してもよい。

50

【0096】

(ユーザ状態の第5の例)

図14は、ユーザ状態がユーザの覚醒状態の度合である場合における処理を示す図である。図14は、ユーザ状態がユーザの覚醒状態の度合である場合における、S130、S140(図7)の具体的な処理を示す。

【0097】

ユーザ状態検出部140は、カメラ又は生体センサである検出装置6から、ユーザの生体パラメータを取得する(ステップS132E)。生体パラメータは、例えば、瞬目、心拍及び脳波の少なくとも1つである。なお、瞬目は、カメラから取得されたユーザの顔画像を解析することによって取得され得る。心拍及び脳波は、それぞれ生体センサである心拍計及び脳波計を用いて取得され得る。

10

【0098】

ユーザ状態検出部140は、生体パラメータから、ユーザの覚醒状態の度合を検出する(ステップS134E)。例えば、ユーザ状態検出部140は、上述した生体パラメータから、覚醒度を算出する。例えば、ユーザ状態検出部140は、瞬目の間隔、瞬目の開眼時間、目の開度等により、覚醒度を算出し得る。そして、ユーザ状態検出部140は、図11に例示したような、覚醒度と覚醒状態の度合とを対応付けたテーブルを用いて、覚醒状態の度合を検出してよい。

【0099】

次に、学習モデル選択部150は、検出された覚醒状態の度合に対応する学習モデルを選択する(ステップS142E)。なお、予め、ユーザの覚醒状態の度合ごとに、複数の学習モデルが学習モデルデータベース160に格納されているとする。例えば、ユーザ状態検出部140によって「覚醒状態の度合#1」が検出された場合、学習モデル選択部150は、「覚醒状態の度合#1」に対応する学習モデルを選択する。また、ユーザ状態検出部140によって「覚醒状態の度合#4」が検出された場合、学習モデル選択部150は、「覚醒状態の度合#4」に対応する学習モデルを選択する。

20

【0100】

このようにして、第5の例にかかる音声対話システム1は、ユーザの覚醒状態の度合に適合した学習モデルを用いて対話を行うので、対話を行うユーザの覚醒状態に合わせて応答を実行することができる。したがって、第5の例にかかる音声対話システム1は、応答誤りが発生しないように適切に対処することが可能となる。また、応答誤りが発生するということは、現在の対話のテンポ又はリズムが、ユーザの覚醒状態の度合に適していないということである。第5の例にかかる音声対話システム1は、ユーザの覚醒状態の度合に対応する学習モデルを選択することによって、対話のテンポ又はリズムをユーザの覚醒状態の度合に適したものすることが可能となる。また、ユーザの覚醒度は、対話の話題等によって変化し得る。第5の例にかかる音声対話システム1は、覚醒度の変化に応じて学習モデルを変更することができる。

30

【0101】

また、第5の例においては、学習モデルを生成する際に、ユーザ状態としてユーザの覚醒状態の度合が対応付けられる。言い換えると、覚醒状態の度合ごとに、複数の学習モデルが生成される。学習モデルを生成する際には、例えば、カメラ又は生体センサを用いて対話中のユーザの覚醒状態の度合を入力することで、サンプルデータとユーザの覚醒状態の度合とが対応付けられる。

40

【0102】

これにより、ユーザの覚醒状態の度合ごとに、サンプルデータが分類され、分類されたサンプルデータを用いて、機械学習によって複数の学習データが生成される。したがって、例えば、覚醒状態の度合#1に対応する学習モデル、覚醒状態の度合#2に対応する学習モデル、覚醒状態の度合#3に対応する学習モデル、及び、覚醒状態の度合#4に対応する学習モデルが生成されることとなる。

【0103】

50

なお、上述した例では、カメラ又は生体センサを用いてユーザの覚醒状態の度合を検出するとしたが、ユーザの覚醒状態の度合を検出する方法は、この方法に限られない。ユーザ状態検出部140は、ユーザ発話を取得することで、覚醒状態の度合を検出してよい。具体的には、ユーザ状態検出部140は、ユーザ発話の韻律を解析してユーザの覚醒状態を判定し、覚醒状態の度合を数値化してもよい。

【0104】

(実施の形態2)

次に、実施の形態2について説明する。実施の形態2においては、音声対話システム1が複数の学習モデルを生成する点で、実施の形態1と異なる。なお、実施の形態2にかかる音声対話システム1のハードウェア構成については、図1に示した実施の形態1にかかる音声対話システム1のハードウェア構成と実質的に同様であるので、説明を省略する。

10

【0105】

図15は、実施の形態2にかかる音声対話システム1の構成を示すブロック図である。実施の形態2にかかる音声対話システム1は、発話取得部102と、特徴抽出部104と、選択モデル格納部108と、応答データベース110と、応答決定部120と、応答実行部130と、応答履歴格納部132とを有する。また、実施の形態1にかかる音声対話システム1は、ユーザ状態検出部140と、学習モデル選択部150と、学習モデルデータベース160とを有する。さらに、音声対話システム1は、学習モデル生成装置200を有する。学習モデル生成装置200以外の構成要素については、実施の形態1にかかるものと実質的に同様の機能を有するので、説明を省略する。

20

【0106】

学習モデル生成装置200は、後述する方法によって、複数の学習モデルを生成する。学習モデル生成装置200によって生成された複数の学習モデルは、学習モデルデータベース160に格納される。学習モデルは、学習モデル生成装置200によって自動的に格納されてもよいし、オペレータ等の作業者によって手動で格納されてもよい。

【0107】

なお、学習モデル生成装置200は、その他の構成要素と物理的に一体となっている必要はない。つまり、その他の構成要素が設けられた装置(ロボット等)と、学習モデル生成装置200が設けられた装置(コンピュータ等)とが、同一である必要はない。学習モデル生成装置200の具体的な構成について、以下に説明する。なお、学習モデル生成装置200の処理(後述する図17に示す処理)は、図4~図6に対応し、ユーザとの対話(図7の処理)の前段階で行われる。

30

【0108】

図16は、実施の形態2にかかる学習モデル生成装置200の構成を示す図である。また、図17は、実施の形態2にかかる学習モデル生成装置200によって実行される学習モデル生成方法を示すフローチャートである。学習モデル生成装置200は、発話取得部212、特徴抽出部214、サンプルデータ生成部216、ユーザ状態取得部218、サンプルデータ分類部220、及び、学習モデル生成部222を有する。なお、学習モデル生成装置200は、図1に示した音声対話システム1のハードウェア構成と実質的に同様のハードウェア構成を、独立して有しうる。

40

【0109】

発話取得部212は、1以上の任意ユーザと対話を行うことによって、図7のS102の処理と同様にして、任意ユーザの発話であるユーザ発話を取得する(ステップS202)。ここで、「任意ユーザ」とは、音声対話システム1が対話を行う相手のユーザに限られない、任意のユーザである。特徴抽出部214は、図7のS104の処理と同様にして、取得されたユーザ発話の特徴を少なくとも示す特徴ベクトルを抽出する(ステップS204)。

【0110】

次に、サンプルデータ生成部216は、ユーザ発話に対する応答を示す正解ラベルと特徴ベクトルとが対応付けられたサンプルデータを生成する(ステップS206)。具体的

50

には、サンプルデータ生成部 216 は、図 4 を用いて上述したように、オペレータによって判定された応答（正解ラベル）を、対応するユーザ発話の特徴ベクトルに対応付ける。これにより、サンプルデータ生成部 216 は、サンプルデータを生成する。なお、正解ラベルを自動的に判定することができれば、サンプルデータ生成部 216 は、ユーザ発話から正解ラベル（応答）を自動的に判定して、判定された正解ラベルをユーザ発話の特徴ベクトルに対応付けてもよい。次に、学習モデル生成装置 200（又は図 2 に示した応答実行部 130）は、図 7 の S120 の処理と同様にして、応答を実行する（ステップ S208）。

【0111】

ユーザ状態取得部 218 は、ユーザ発話を発したときの任意ユーザの状態であるユーザ状態を取得して、取得されたユーザ状態をユーザ発話に対応するサンプルデータに対応付ける（ステップ S210）。具体的には、ユーザ状態取得部 218 は、図 9～図 14 を用いて説明したように、ユーザの画像、ユーザ発話、又は生体系パラメータ等を用いて、任意ユーザのユーザ状態を取得し得る。ユーザ状態の取得方法は、ユーザ状態の種類（第 1 の例～第 5 の例）に応じて異なり得る。あるいは、ユーザ状態取得部 218 は、例えばオペレータによって判断された、任意ユーザのユーザ状態を取得してもよい。そして、ユーザ状態取得部 218 は、取得されたユーザ状態を、任意ユーザのユーザ発話に対応するサンプルデータに対応付ける。

【0112】

学習モデル生成装置 200 は、ユーザ発話の取得を終了するか否かを判定する（ステップ S212）。ユーザ発話の取得を継続する場合（S212 の NO）、学習モデル生成装置 200 は、S202～S210 の処理を繰り返す。一方、サンプルデータを十分に取得できたためにユーザ発話の取得を終了する場合（S212 の YES）、サンプルデータ分類部 220 は、図 6 を用いて説明したように、ユーザ状態ごとにサンプルデータを分類する（ステップ S220）。そして、学習モデル生成部 222 は、図 6 を用いて説明したように、分類されたサンプルデータごとに、例えばランダムフォレスト又はサポートベクターマシン等の機械学習によって、複数の学習モデルを生成する（ステップ S222）。

【0113】

このように、実施の形態 2 にかかる学習モデル生成装置 200 は、ユーザ状態ごとにサンプルデータを分類して機械学習によって複数の学習モデルを生成することによって、ユーザ状態に対応した複数の学習モデルを生成することができる。したがって、音声対話システム 1 は、上述したように、ユーザ状態に応じて学習モデルを選択することができる。

【0114】

（変形例）

なお、本発明は上記実施の形態に限られたものではなく、趣旨を逸脱しない範囲で適宜変更することが可能である。例えば、上述したフローチャートにおいて、複数の処理の順序は、適宜、変更可能である。また、上述したフローチャートにおいて、複数の処理のうちの一つは、省略されてもよい。例えば、図 7 の S130 の処理は、S102～S120 の間に行われてもよい。

【0115】

また、図 9～図 14 を用いて説明したユーザ状態の第 1 の例～第 5 の例は、相互に適用可能である。つまり、ユーザ状態検出部 140 は、複数の種類のユーザ状態を検出してもよい。そして、学習モデル選択部 150 は、検出された複数の種類のユーザ状態に対応する学習モデルを選択してもよい。例えば、ユーザ状態検出部 140 は、ユーザの識別情報及びユーザの積極性の度合を検出してもよい。この場合、学習モデルデータベース 160 は、例えば、ユーザ A の積極性の度合ごと、ユーザ B の積極性の度合ごとに、複数の学習モデルを格納し得る。そして、ユーザ状態検出部 140 が、「ユーザ A」の「積極性の度合 #1」を検出した場合に、学習モデル選択部 150 は、「ユーザ A」の「積極性の度合 #1」に対応する学習モデルを選択し得る。

【0116】

10

20

30

40

50

また、上述した実施の形態では、特徴ベクトル（図3）は、ユーザ発話の韻律情報等から生成されるとしたが、このような構成に限られない。特徴ベクトルの成分は、韻律にかかるものだけでなく、カメラである検出装置6から取得されたユーザの特徴も含み得る。例えば、特徴ベクトルの成分は、ユーザの視線及び対話ロボットに対するユーザの距離を含んでもよい。

【0117】

また、上述した実施の形態においては、音声対話システム1がロボットに搭載された例を示しているが、このような構成に限られない。音声対話システム1は、スマートフォン又はタブレット端末等の情報端末にも搭載可能である。この場合、「頷き応答」を行うときは、マニピュレータ8を動作させる代わりに、情報端末の表示画面に、人物、動物、又はロボット等が頷くような動画を表示させてもよい。

10

【0118】

また、上述の例において、プログラムは、様々なタイプの非一時的なコンピュータ可読媒体（non-transitory computer readable medium）を用いて格納され、コンピュータに供給することができる。非一時的なコンピュータ可読媒体は、様々なタイプの実体のある記録媒体（tangible storage medium）を含む。非一時的なコンピュータ可読媒体の例は、磁気記録媒体（例えばフレキシブルディスク、磁気テープ、ハードディスクドライブ）、光磁気記録媒体（例えば光磁気ディスク）、CD-ROM、CD-R、CD-R/W、半導体メモリ（例えば、マスクROM、PROM（Programmable ROM）、EPROM（Erasable PROM）、フラッシュROM、RAM）を含む。また、プログラムは、様々なタイプの一時的なコンピュータ可読媒体（transitory computer readable medium）によってコンピュータに供給されてもよい。一時的なコンピュータ可読媒体の例は、電気信号、光信号、及び電磁波を含む。一時的なコンピュータ可読媒体は、電線及び光ファイバ等の有線通信路、又は無線通信路を介して、プログラムをコンピュータに供給できる。

20

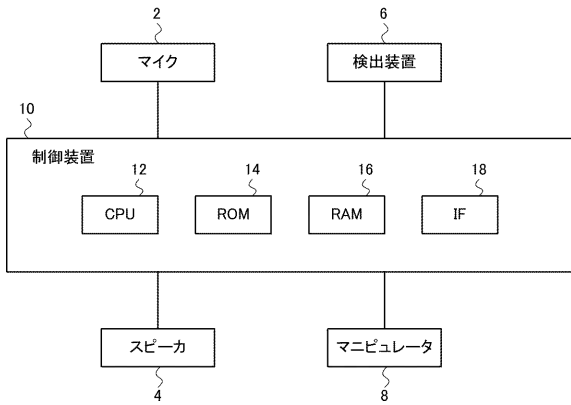
【符号の説明】

【0119】

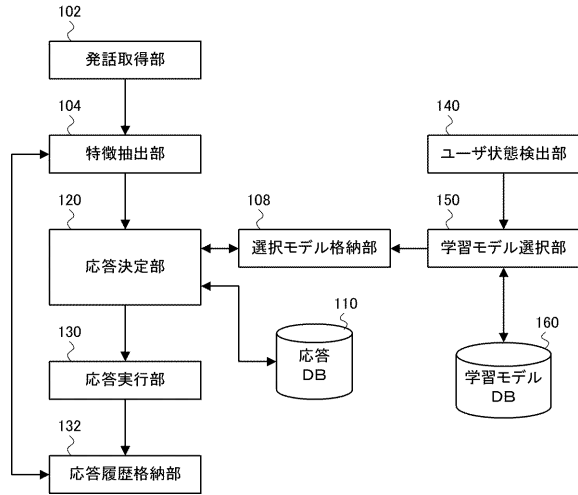
1・・・音声対話システム、2・・・マイク、4・・・スピーカ、6・・・検出装置、8・・・マニピュレータ、10・・・制御装置、102・・・発話取得部、104・・・特徴抽出部、108・・・選択モデル格納部、110・・・応答データベース、120・・・応答決定部、130・・・応答実行部、132・・・応答履歴格納部、140・・・ユーザ状態検出部、150・・・学習モデル選択部、160・・・学習モデルデータベース、200・・・学習モデル生成装置、212・・・発話取得部、214・・・特徴抽出部、216・・・サンプルデータ生成部、218・・・ユーザ状態取得部、220・・・サンプルデータ分類部、222・・・学習モデル生成部

30

【図1】



【図2】



【図3】

特徴ベクトル

$$V_i = (V_{i1}, V_{i2}, V_{i3}, V_{i4}, V_{i5}, V_{i6}, V_{i7}, \dots, V_{im}, \dots)$$

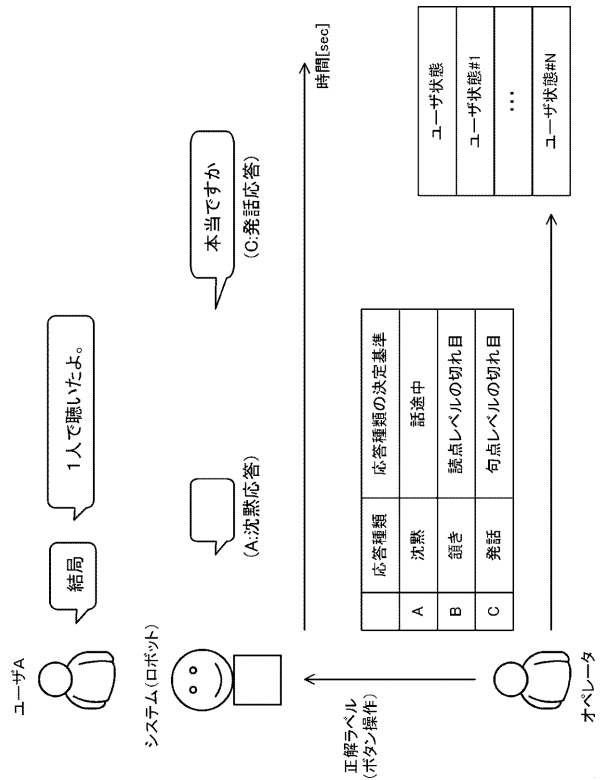
$$= (f0_{T1}, V_{T1}, f0_{T2}, V_{T2}, f0, V, L1, \dots, 1, \dots)$$

情報の種類	要素	成分値	成分
i番目の ユーザ発話 情報	句末T1msec	f0	f0 _{T1} v ₁
		ボリューム	V _{T1} v ₂
	句末T2msec	f0	f0 _{T2} v ₃
		ボリューム	V _{T2} v ₄
	発話区間 全体	f0	f0 v ₅
		ボリューム	V v ₆
ユーザ発話長	L1	v ₇	
...	
システム応答 の履歴	直前の応答タイプ	1	v _m

直前の応答タイプ:

- 1:沈黙、2:頷き、3:発話

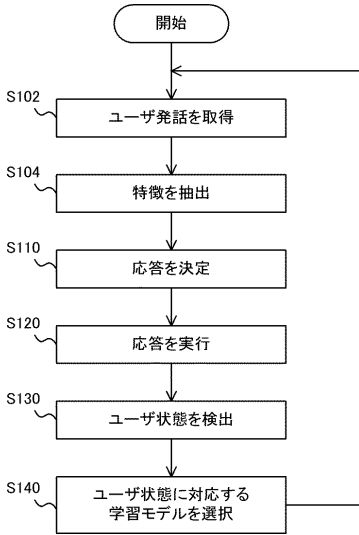
【図4】



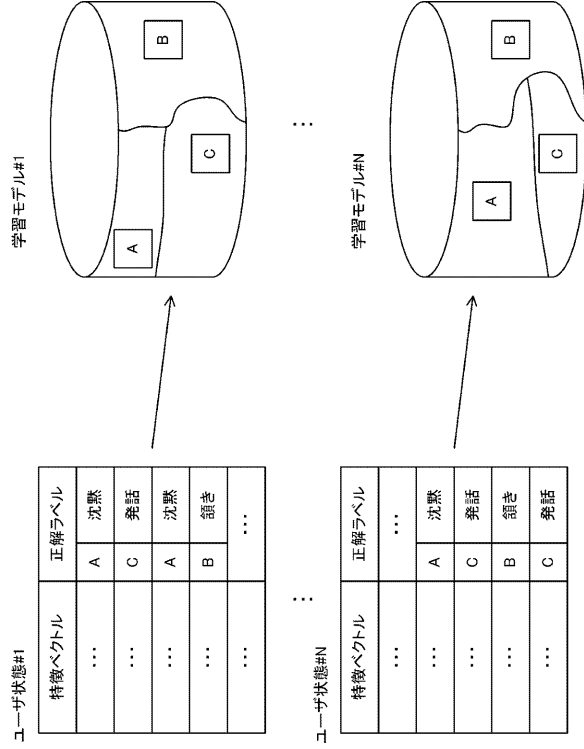
【図 5】

ユーザ状態	ユーザ発言	特徴ベクトル				正解ラベル		
		句末T1msec		句末T2msec		ユーザ発言長 (sec)		
		f0	ボリウム	f0	ボリウム	
ユーザ状態#1	「結局」	0.5	A	沈黙
	「1人で観いたよ。」	1.5	C	発言

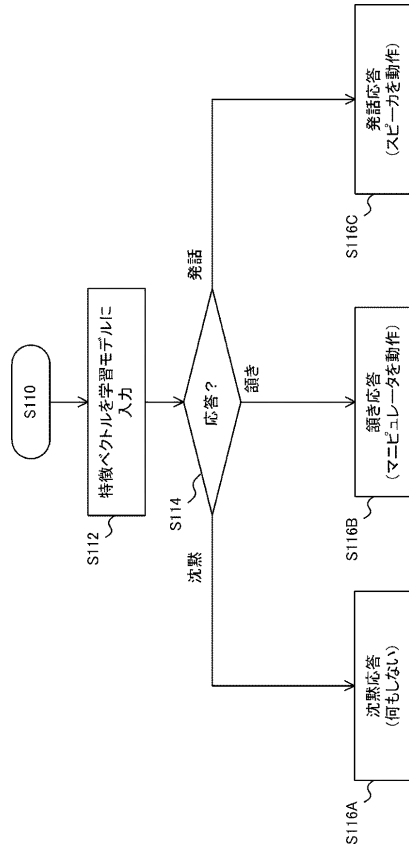
【図 7】



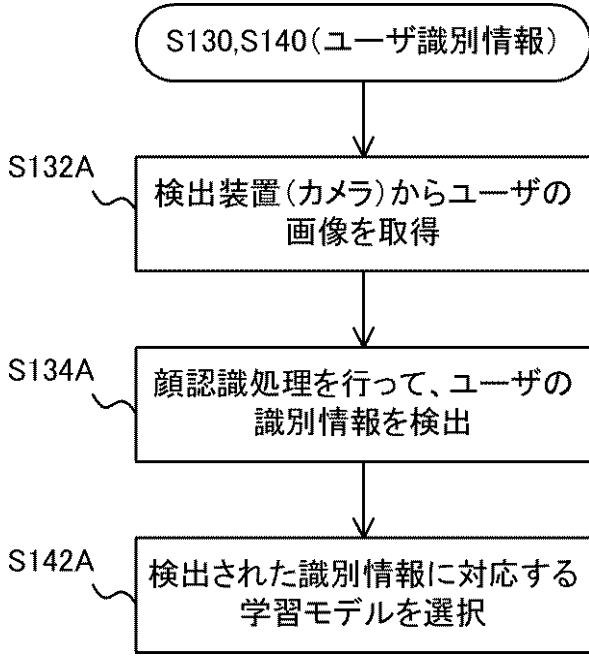
【図 6】



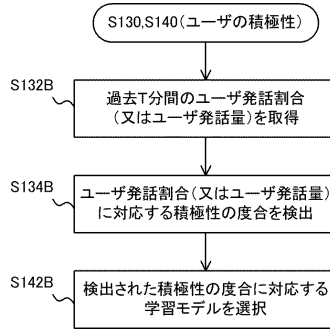
【図 8】



【図9】



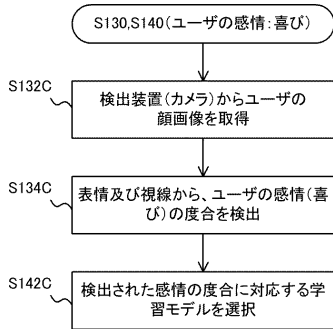
【図10】



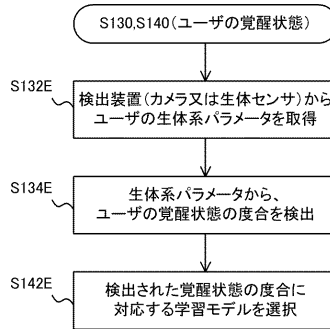
【図11】

積極性の度合	ユーザ発話割合
小 ↑ #1	$0 \leq R_s < 25$
#2	$25 \leq R_s < 50$
#3	$50 \leq R_s < 75$
大 ↓ #4	$75 \leq R_s \leq 100$

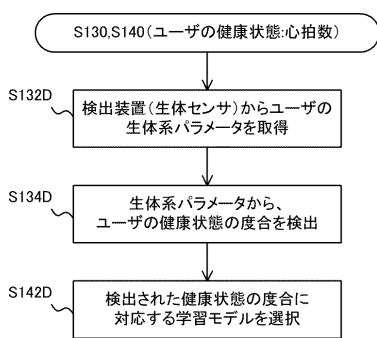
【図12】



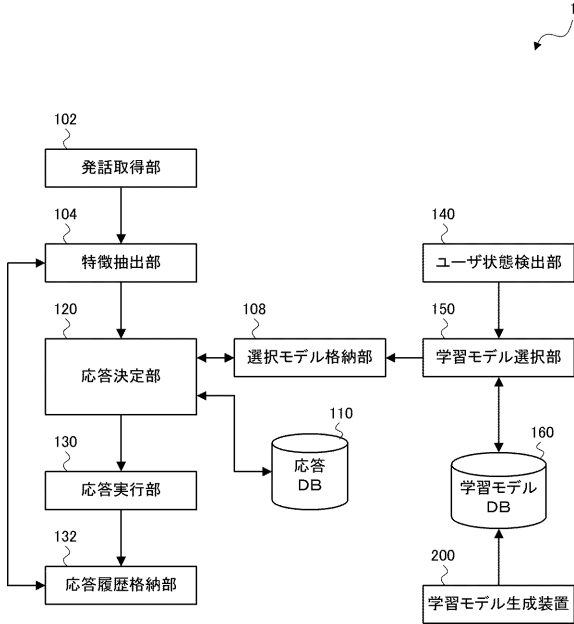
【図14】



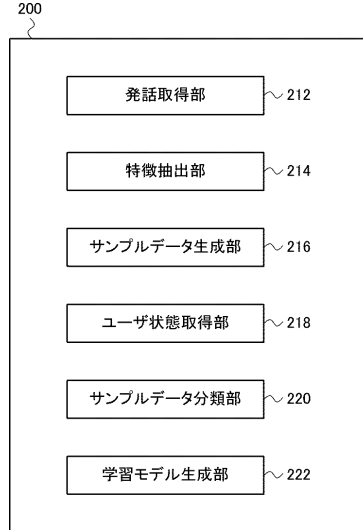
【図13】



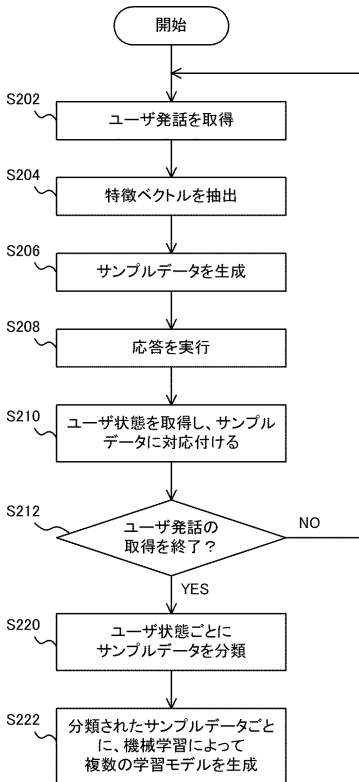
【図15】



【図16】



【図17】



フロントページの続き

(51)Int.Cl.	F I				テーマコード(参考)
	G 0 6 F	3/16		6 9 0	
	G 0 6 F	3/01		5 1 0	

(72)発明者 渡部 生聖

愛知県豊田市トヨタ町1番地 トヨタ自動車株式会社内

Fターム(参考) 5E555 AA11 AA48 AA76 BA01 BA06 BA88 BB01 BB06 BC01 BC08
BC17 CA41 CA42 CA47 CB64 CB66 CB69 CB76 DA23 EA02
EA05 EA19 EA20 EA22 EA23 EA28 FA00