

(19)日本国特許庁(JP)

(12)特許公報(B1)

(11)特許番号

特許第7164793号

(P7164793)

(45)発行日 令和4年11月2日(2022.11.2)

(24)登録日 令和4年10月25日(2022.10.25)

(51)国際特許分類

F I

G 1 0 L 13/08 (2013.01)

G 1 0 L 13/08 1 2 4

G 1 0 L 13/033 (2013.01)

G 1 0 L 13/033 1 0 2 A

G 1 0 L 13/02 (2013.01)

G 1 0 L 13/02 1 3 0 Z

G 1 0 L 15/10 (2006.01)

G 1 0 L 15/10 5 0 0 N

請求項の数 12 (全25頁)

(21)出願番号 特願2021-190678(P2021-190678)

(22)出願日 令和3年11月25日(2021.11.25)

審査請求日 令和3年11月25日(2021.11.25)

早期審査対象出願

(73)特許権者 501440684

ソフトバンク株式会社

東京都港区海岸一丁目7番1号

(74)代理人 100079108

弁理士 稲葉 良幸

(74)代理人 100109346

弁理士 大貫 敏史

(74)代理人 100117189

弁理士 江口 昭彦

(74)代理人 100134120

弁理士 内藤 和彦

(72)発明者 中谷 敏之

東京都港区海岸一丁目7番1号

(72)発明者 末永 君慧

東京都港区海岸一丁目7番1号

最終頁に続く

(54)【発明の名称】 音声処理システム、音声処理装置及び音声処理方法

(57)【特許請求の範囲】

【請求項1】

第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、

前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、

前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する音声合成部と、

第2のユーザに対して前記合成音声を出力する音声出力部と、

前記発話音声信号に対応する第1のユーザの感情情報を生成する感情認識部と、

前記感情情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、を備え、

前記制御部は、前記感情情報が第2のユーザにストレスを与えうるものである場合に前記合成音声を出力する、

音声処理システム。

【請求項2】

第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、

前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、

前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する音声合成部と、

10

20

第 2 のユーザに対して前記合成音声を出力する音声出力部と、
前記第 2 のユーザのストレス状況に関するストレス情報を生成するストレス認識部と、
前記ストレス情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のど
ちらを出力するかを切り替える制御部と、を備え、
前記制御部は、前記ストレス情報がストレスの高い状態を示している場合に前記合成音
声を出力する、
音声処理システム。

【請求項 3】

第 1 のユーザの発話音声の信号である発話音声信号を取得する取得部と、
前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の
単語からなる単語列を含むテキストデータを生成する音声認識部と、
前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音
声の信号である合成音声信号を生成する音声合成部と、
第 2 のユーザに対して前記合成音声を出力する音声出力部と、
前記第 2 のユーザによって入力される切り替え情報に基づいて、前記音声出力部から前
記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、
前記第 2 のユーザによって入力された切り替え情報を、前記切り替え情報が入力された
際の発話音声信号と時間軸上で関連付けた情報を生成し、当該情報に基づいて、発話音声
信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデ
ータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合
わせを入力とし、前記合成音声と前記発話音声とを切り替えるタイミングを出力とする感
情抑制切替モデルを機械学習する学習部と、を備え、
前記制御部は、前記感情抑制切替モデルに、前記取得部が取得した発話音声信号、当該
発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該
テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力
することにより、前記合成音声と前記発話音声とを切り替えるタイミングを生成する、
音声処理システム。

【請求項 4】

前記感情認識部は、発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力とし、当該発話音声信号の発話者の感情情報を出力するよう機械学習された感情認識モデルに、前記取得部が取得した発話音声信号、当該発話音声信号から抽出した音声特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータに対応するテキスト特徴量、又はこれらの少なくとも二つの組み合わせを入力することにより、前記取得部が取得した発話音声信号に対応する第 1 のユーザの感情情報を生成する、請求項 1 に記載の音声処理システム。

【請求項 5】

前記発話音声信号に対応する第 1 のユーザの感情情報を生成する感情認識部と、
前記第 2 のユーザに対して前記感情情報を表示する表示部と、を備え、
前記表示部は、前記音声出力部による前記合成音声の出力タイミングに合わせて、前記合成音声に対応する前記感情情報を表示する、
請求項 1 又は 4 に記載の音声処理システム。

【請求項 6】

前記音声合成部は、前記感情認識部が生成した感情情報に基づいて、前記感情情報が示
す感情が前記合成音声に反映されるように、前記合成音声信号を生成する、
請求項 1 又は 4 に記載の音声処理システム。

【請求項 7】

第 1 のユーザの発話音声の信号である発話音声信号を取得する取得部と、
前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の
単語からなる単語列を含むテキストデータを生成する音声認識部と、

前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、第2のユーザに対して出力される合成音声の信号である合成音声信号を生成する音声合成部と、前記発話音声信号に対応する第1のユーザの感情情報を生成する感情認識部と、

前記感情情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、を備え、

前記制御部は、前記感情情報が第2のユーザにストレスを与えうるものである場合に前記合成音声を出力する、

音声処理装置。

【請求項8】

第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、

前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、

前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、第2のユーザに対して出力される合成音声の信号である合成音声信号を生成する音声合成部と、前記第2のユーザのストレス状況に関するストレス情報を生成するストレス認識部と、

前記ストレス情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、を備え、

前記制御部は、前記ストレス情報がストレスの高い状態を示している場合に前記合成音声

を出力する、

音声処理装置。

【請求項9】

第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、

前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、

前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、第2のユーザに対して出力される合成音声の信号である合成音声信号を生成する音声合成部と、前記第2のユーザによって入力される切り替え情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、

前記第2のユーザによって入力された切り替え情報を、前記切り替え情報が入力された際の発話音声信号と時間軸上で関連付けた情報を生成し、当該情報に基づいて、発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力とし、前記合成音声と前記発話音声とを切り替えるタイミングを出力とする感情抑制切替モデルを機械学習する学習部と、

を備え、

前記制御部は、前記感情抑制切替モデルに、前記取得部が取得した発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力することにより、前記合成音声と前記発話音声とを切り替えるタイミングを生成する、

音声処理装置。

【請求項10】

第1のユーザの発話音声の信号である発話音声信号を取得する工程と、

前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する工程と、

前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する工程と、

第2のユーザに対して前記合成音声を出力する工程と、

前記発話音声信号に対応する第1のユーザの感情情報を生成する工程と、

前記感情情報に基づいて、前記合成音声又は前記発話音声のどちらを出力するかを切り替える工程と、を含み、

10

20

30

40

50

前記切り替える工程は、前記感情情報が第２のユーザにストレスを与えうるものである場合に前記合成音声を出力する、
音声処理方法。

【請求項１１】

第１のユーザの発話音声の信号である発話音声信号を取得する工程と、
前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する工程と、
前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する工程と、
第２のユーザに対して前記合成音声を出力する工程と、
前記第２のユーザのストレス状況に関するストレス情報を生成する工程と、
前記ストレス情報に基づいて、前記合成音声又は前記発話音声のどちらを出力するかを切り替える工程と、を含み、
前記切り替える工程は、前記ストレス情報がストレスの高い状態を示している場合に前記合成音声を出力する、
音声処理方法。

10

【請求項１２】

第１のユーザの発話音声の信号である発話音声信号を取得する工程と、
前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する工程と、
前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する工程と、
第２のユーザに対して前記合成音声を出力する工程と、
前記第２のユーザによって入力される切り替え情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える工程と、
前記第２のユーザによって入力された切り替え情報を、前記切り替え情報が入力された際の発話音声信号と時間軸上で関連付けた情報を生成し、当該情報に基づいて、発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力とし、前記合成音声と前記発話音声とを切り替えるタイミングを出力とする感情抑制切替モデルを機械学習する工程と、を備え、
前記切り替える工程は、前記感情抑制切替モデルに、前記発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力することにより、前記合成音声と前記発話音声とを切り替えるタイミングを生成する、
音声処理方法。

20

30

【発明の詳細な説明】

【技術分野】

【０００１】

本発明は、音声処理システム、音声処理装置及び音声処理方法に関する。

40

【背景技術】

【０００２】

従来、顧客満足度（Customer Satisfaction：ＣＳ）向上のために、顧客の苦情等に対してオペレータが電話で対応する各種のコールセンターが運用されている。このような顧客対応業務では、顧客がオペレータに対して威圧的な言動や理不尽な要求を行う「カスタマーハラスメント」により、オペレータの精神不調を招いたり、オペレータの離職率が高くなったりすることが問題視されている。

【０００３】

近年、このようなカスタマーハラスメントから、企業側が従業員であるオペレータを守るための音声変換システムも検討されている。例えば、特許文献１では、入力音声信号が

50

ら音量及びピッチ変動量を算出し、音量及びピッチ変動量が所定値を超える場合に、音量及びピッチ変動量が所定内に収まるように音量及びピッチを変換して出力するように制御することが記載されている。

【先行技術文献】

【特許文献】

【0004】

【文献】特開2004-252085号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、例えば、特許文献1に記載の方法で話し手の発話音声を変換するだけでは、話し手（第1のユーザ）の感情が十分に抑制されず、聞き手（第2のユーザ）のストレスを十分に軽減できない恐れがある。一方、聞き手のストレスを軽減するために、聞き手に出力される話し手の発話音声を変換すると、聞き手が話し手の感情を十分に認識できず、聞き手が適切な対応を行うことができない恐れもある。

【0006】

そこで、本発明は、聞き手のストレスの十分な軽減、及び／又は、聞き手の適切な対応を可能とする音声処理システム、音声処理装置及び音声処理方法を提供する。

【課題を解決するための手段】

【0007】

本発明の一つの態様に係る音声処理システムは、第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する音声合成部と、第2のユーザに対して前記合成音声出力する音声出力部と、前記発話音声信号に対応する第1のユーザの感情情報を生成する感情認識部と、前記感情情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、を備え、前記制御部は、前記感情情報が第2のユーザにストレスを与えうるものである場合に前記合成音声出力する、を備える。

【0008】

本発明の一つの態様に係る音声処理システムは、第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する音声合成部と、第2のユーザに対して前記合成音声出力する音声出力部と、前記第2のユーザのストレス状況に関するストレス情報を生成するストレス認識部と、前記ストレス情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、を備え、前記制御部は、前記ストレス情報がストレスの高い状態を示している場合に前記合成音声出力する。

【0009】

本発明の一つの態様に係る音声処理システムは、第1のユーザの発話音声の信号である発話音声信号を取得する取得部と、前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する音声合成部と、第2のユーザに対して前記合成音声出力する音声出力部と、前記第2のユーザによって入力される切り替え情報に基づいて、前記音声出力部から前記合成音声又は前記発話音声のどちらを出力するかを切り替える制御部と、前記第2のユーザによって入力された切り替え情報を、前記切り替

10

20

30

40

50

え情報が入力された際の発話音声信号と時間軸上で関連付けた情報を生成し、当該情報に基づいて、発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力とし、前記合成音声と前記発話音声とを切り替えるタイミングを出力とする感情抑制切替モデルを機械学習する学習部と、を備え、前記制御部は、前記感情抑制切替モデルに、前記取得部が取得した発話音声信号、当該発話音声信号から抽出した特徴量、当該発話音声信号から生成したテキストデータ、当該テキストデータから抽出された特徴量、又はこれらの少なくとも二つの組み合わせを入力することにより、前記合成音声と前記発話音声とを切り替えるタイミングを生成する。

【図面の簡単な説明】

【 0 0 1 0 】

【図 1】本実施形態に係る音声処理システム 1 の概略の一例を示す図である。

【図 2】本実施形態に係る音声処理システム 1 を構成する各装置の物理構成の一例を示す図である。

【図 3】本実施形態に係る音声処理装置 10 の機能構成の一例を示す図である。

【図 4】本実施形態に係る合成音声信号の生成の一例を示す図である。

【図 5 A】本実施形態に係る顧客の感情情報の生成の一例を示す図である。

【図 5 B】本実施形態に係る顧客の感情情報の生成の一例を示す図である。

【図 6】本実施形態に係るオペレータ端末 20 の機能構成の一例を示す図である。

【図 7】本実施形態に係る画面 D 1 の一例を示す図である。

【図 8】本実施形態に係る画面 D 2 の一例を示す図である。

【図 9】本実施形態に係る感情抑制動作の一例を示すフローチャートである。

【図 10】本実施形態に係る感情抑制機能の自動切り替え動作を示すフローチャートである。

【図 11】本実施形態の変更例に係る合成音声信号の生成の一例を示す図である。

【図 12】本実施形態に係る画面 D 3 の一例を示す図である。

【発明を実施するための形態】

【 0 0 1 1 】

添付図面を参照して、本発明の実施形態について説明する。なお、各図において、同一の符号を付したものは、同一又は同様の構成を有する。

【 0 0 1 2 】

以下、本実施形態に係る音声処理システムをコールセンター等の顧客対応業務において使用することを想定して説明を行うが、本発明の適用形態はこれに限られない。本実施形態は、第 1 のユーザの発話音声の信号（以下、「発話音声信号」という）に所定の処理を施して生成される音声を第 2 のユーザに対して出力するどのような場面にも適用可能である。以下では、第 1 のユーザが顧客であり、第 2 のユーザがオペレータであるものとするが、これに限られない。

【 0 0 1 3 】

（音声処理システムの構成）

<全体構成>

図 1 は、本実施形態に係る音声処理システム 1 の概略の一例を示す図である。図 1 に示すように、音声処理システム 1 は、音声処理装置 10 と、第 2 のユーザ（以下、「オペレータ」という）によって使用される端末（以下、「オペレータ端末」という）20 と、第 1 のユーザ（以下、「顧客」という）によって使用される端末（以下、「顧客端末」という）30 と、を備える。

【 0 0 1 4 】

音声処理装置 10 は、顧客端末 30 で取得される発話音声信号を、ネットワーク 40 を介して受信する。ネットワーク 40 は、インターネット等の外部ネットワークであってもよいし、外部ネットワーク、及び、Local Access Network（LAN）等の内部ネットワークを含んでもよい。音声処理装置 10 は、顧客の発話音声信号に対して所定の処理を施

10

20

30

40

50

した音声をオペレータ端末 20 に送信する。なお、音声処理装置 10 は、一つ又は複数のサーバで構成されてもよい。

【0015】

オペレータ端末 20 は、例えば、電話、スマートフォン、パーソナルコンピュータ、タブレット等である。オペレータ端末 20 は、音声処理装置 10 で所定の処理で生成される音声信号又は顧客端末 30 からの発話音声信号に基づいて、音声をオペレータに出力する。

【0016】

顧客端末 30 は、例えば、電話、スマートフォン、パーソナルコンピュータ、タブレット等である。顧客端末 30 は、顧客の発話音声をマイクにより収音して、当該発話音声の信号である発話音声信号を音声処理装置 10 に送信する。

【0017】

<物理構成>

図 2 は、本実施形態に係る音声処理システム 1 を構成する各装置の物理構成の一例を示す図である。各装置（例えば、音声処理装置 10、オペレータ端末 20 及び顧客端末 30）は、演算部に相当するプロセッサ 10a と、記憶部に相当する RAM（Random Access Memory）10b と、記憶部に相当する ROM（Read Only Memory）10c と、通信部 10d と、入力部 10e と、表示部 10f と、カメラ 10g、音声入力部 10h と、音声出力部 10i と、を有する。これらの各構成は、バスを介して相互にデータ送受信可能に接続される。なお、図 2 で示す構成は一例であり、各装置はこれら以外の構成を有してもよいし、これらの構成のうち一部を有さなくてもよい。

【0018】

プロセッサ 10a は、例えば、CPU（Central Processing Unit）である。プロセッサ 10a は、RAM 10b 又は ROM 10c に記憶されているプログラムを実行することにより、各装置における各種処理を制御する制御部である。プロセッサ 10a は、各装置が備える他の構成と、プログラムとの協働により、各装置の機能を実現し、処理の実行を制御する。プロセッサ 10a は、入力部 10e や通信部 10d から種々のデータを受け取り、データの演算結果を表示部 10f に表示したり、RAM 10b に格納したりする。

【0019】

RAM 10b 及び ROM 10c は、各種処理に必要なデータ及び処理結果のデータを記憶する記憶部である。各装置は、RAM 10b 及び ROM 10c 以外に、ハードディスクドライブ等の大容量の記憶部を備えてもよい。RAM 10b 及び ROM 10c は、例えば、半導体記憶素子で構成されてもよい。

【0020】

通信部 10d は、各装置を他の機器に接続するインターフェースである。通信部 10d は、他の機器と通信する。入力部 10e は、ユーザからデータの入力を受け付けるためのデバイスや、各装置の外部からデータを入力するためのデバイスである。入力部 10e は、例えば、キーボード、マウス及びタッチパネル等を含んでよい。表示部 10f は、プロセッサ 10a による制御に従って、情報を表示するデバイスである。表示部 10f は、例えば、LCD（Liquid Crystal Display）により構成されてよい。

【0021】

カメラ 10g は、静止画像又は動画像を撮像する撮像素子を含み、所定の領域の撮像により撮像画像（例えば、静止画像又は動画像）を生成する。音声入力部 10h は、音声を収音するデバイスであり、例えば、マイクである。音声出力部 10i は、音声を出力するデバイスであり、例えば、スピーカーである。

【0022】

各装置を実行させるためのプログラムは、RAM 10b や ROM 10c 等のコンピュータによって読み取り可能な記憶媒体に記憶されて提供されてもよいし、通信部 10d により接続されるネットワーク 40 を介して提供されてもよい。各装置では、プロセッサ 10a が当該プログラムを実行することにより、各装置を制御するための様々な動作が実現される。なお、これらの物理的な構成は例示であって、必ずしも独立した構成でなくともよ

10

20

30

40

50

い。例えば、各装置は、プロセッサ 10 a と RAM 10 b や ROM 10 c が一体化した LSI (Large-Scale Integration) を備えていてもよい。

【0023】

<機能的構成>

音声処理装置

図3は、本実施形態に係る音声処理装置10の機能構成の一例を示す図である。音声処理装置10は、記憶部101、送受信部102、音声認識部103、除去部104、音声合成部105、感情認識部106、ストレス認識部107、制御部108、学習部109を含む。

【0024】

記憶部101は、各種情報、プログラム、アルゴリズム、モデル、操作ログ等を記憶する。具体的には、記憶部101は、後述する音声認識モデル101a、音声合成モデル101b、感情認識モデル101c、ストレス認識モデル101d、感情抑制切替モデル101e等を記憶する。

【0025】

送受信部102は、オペレータ端末20及び/又は顧客端末30との間で、種々の情報及び/又は信号を送信及び/又は受信する。例えば、送受信部102(取得部)は、顧客端末30で収音された顧客の発話音声の信号である発話音声信号を取得する。送受信部102は、オペレータ端末20に対して、合成音声信号及び/又は発話音声信号を送信する。また、送受信部102は、オペレータ端末20からオペレータによる操作ログを取得してもよい。操作ログにはオペレータによる顧客の感情の主観的評価に関する情報(以下、「主観的評価情報」という)、後述する「ストレスの度合い」、後述する「手動切替履歴データ」が含まれてよい。また、送受信部102は、オペレータ端末20に対して、顧客の感情に関する情報(以下、「感情情報」という)等を送信してもよい。

【0026】

音声認識部103は、送受信部102で取得された発話音声信号に基づいて抽出される特徴量(以下、「音声特徴量」という)を音声認識モデル101aに入力して、一以上の単語からなる単語列を含むテキストデータを生成する。具体的には、音声認識部103は、音声認識モデル101aの音響モデルを用いて上記音声特徴量から単語列を生成し、言語モデルを用いた単語列の分析結果に従って上記テキストデータを生成してもよい。音声認識部103は、発話音声信号に対して前処理(例えば、アナログ信号のデジタル化、ノイズの除去、フーリエ変換等)を実施して、音声特徴量を抽出してもよい。

【0027】

音声認識モデル101aは、音声信号に基づいて音声の内容を推定するアルゴリズムである。音声認識モデル101aは、ある単語がどのような音となって現れやすいかということモデル化した音響モデル、及び/又は、特定の言語においてある単語列がどのくらいの確率で現れるかをモデル化した言語モデルを含んでもよい。音響モデルとしては、例えば、隠れマルコフモデル(Hidden Markov Model: HMM)及び/又はディープニューラルネットワーク(Deep Neural Network: DNN)が用いられてもよい。言語モデルとしては、例えば、nグラム言語モデル等の確率的言語モデルが用いられてもよい。

【0028】

除去部104は、音声認識部103で生成されたテキストデータに含まれる特定の単語列を検出し、当該特定の単語列を除去又は前記特定の単語列を他の単語列に置換したテキストデータを生成し、音声合成部105に出力する。除去部104は、音声認識部103で生成されたテキストデータ内で特定の単語列を検出されない場合、当該テキストデータを音声合成部105に出力してもよい。

【0029】

当該特定の単語列は、例えば、聞き手を侮辱したり、聞き手の人格を否定したりする、聞き手を不快にする等、聞き手に心理的悪影響を与える一以上の単語であってもよい。ここで、各単語は、名詞、動詞、副詞、助詞、形容詞、助動詞等の少なくとも一つの品詞、

10

20

30

40

50

当該品詞が音変化したもの等を含んでもよい。例えば、特定の単語列は、「お前、ぶっ殺すぞ」というような「文」であってもよいし、「困るつつってんの」の「つつってん」等、乱暴な言葉遣いであることを示す「文の一部」であってもよい。除去部 104 は、テキストデータ内で検出された特定の単語列のみを他の単語列に置き換えたテキストデータを音声合成部 105 に出力してもよいし、又は、当該特定の単語列を含む文全体を他の単語列に置き換えたテキストデータを音声合成部 105 に出力してもよい。当該他の単語列は、空白等であってもよい。

【0030】

除去部 104 は、記憶部 101 に予め記憶された特定の単語列に基づいて、テキストデータ内の特定の単語列の検出及び／又は他の単語列への置き換えを実施してもよい。

10

【0031】

或いは、除去部 104 は、機械学習により学習されたモデルに基づいて、テキストデータ内の特定の単語列の検出、及び／又は、意味的感情を緩和した他の単語列への置き換えを実施してもよい。例えば、テキストデータ内の特定の単語列「お前」は、「あなた」に置換されてもよい。機械学習に基づくモデルに基づいて、テキストデータ内の特定の単語列の検出及び／又は他の単語列への置き換えを実施してもよい。

【0032】

なお、除去部 104 は、テキストデータ内で特定の単語列が検出される場合、当該特定の単語列の検出に関する情報（以下、「検出情報」という）を生成してもよい。当該検出情報は、例えば、当該特定の単語列が検出されたことを示す情報（例えば、「NGワード」又は「NGワード検出」という文字列）、当該特定の単語列を示す情報、及び、顧客に対する警告に関する情報（以下、「警告情報」という）の少なくとも一つを含んでもよい。当該警告情報は、例えば、オペレータに対する顧客の発話内容が侮辱罪、名誉棄損罪等の刑事告訴対象となり得ることを通知するための情報であってもよい。検出情報は、送受信部 102 によってオペレータ端末 20 に送信されてもよい。検出情報が生成された場合、音声処理装置 10 は、顧客端末 30 に対して警告情報（例えば、「当社オペレータに対して侮辱罪等の恐れがあります。当社の不手際もあるとは思いますが、当社オペレータに過度な負担になる場合がありますのでご協力を頂きますと幸いです。」）を出力させてもよい。このような警告情報は、カスタマーハラスメントに対する事前告知として利用することができる。

20

30

【0033】

音声合成部 105 は、除去部 104 から入力されるテキストデータに基づいて抽出される特徴量（以下、「テキスト特徴量」という）を音声合成モデル 101b に入力して、合成音声の信号（以下、「合成音声信号」という）を生成する。具体的には、除去部 104 は、テキスト特徴量に基づいて音声合成パラメータを予測し、予測された音声合成パラメータを用いて合成音声信号を生成してもよい。音声合成部 105 は、合成音声信号を送受信部 102 に出力する。合成音声信号は、テキストデータの内容を読み上げた音声の信号ともいえる。

【0034】

音声合成モデル 101b は、テキストデータを入力として当該テキストデータの内容に対応する合成音声信号を出力するアルゴリズムである。音声合成モデル 101b としては、例えば、上記 HMM 及び／又は DNN が用いられてもよい。

40

【0035】

当該音声合成モデル 101b は、複数の音声種別に対応してもよい。音声合成部 105 は、複数の音声種別の中から合成音声信号に用いる音声種別を選択し、選択した音声種別とテキストデータとを音声合成モデル 101b に入力して、選択した音声種別の合成音声信号を合成してもよい。当該複数の音声種別は、例えば、抑揚が少ない音声、機械音、キャラクターの音声、芸能人の音声及び声優の音声の少なくとも一つ等であってもよい。音声合成部 105 は、オペレータからオペレータ端末 20 を介して音声種別の選択を受け付けてもよい。

50

【 0 0 3 6 】

図 4 は、本実施形態に係る合成音声信号の生成の一例を示す図である。図 4 では、送受信部 1 0 2 で取得された発話音声信号 S 1 ~ S 3 に基づいて、音声認識部 1 0 3 においてテキストデータ T 1 ~ T 3 が生成されるものとする。例えば、図 4 では、除去部 1 0 4 は、テキストデータ T 1 内で特定の単語列を検出しないので、テキストデータ T 1 をそのまま音声合成部 1 0 5 に出力する。一方、除去部 1 0 4 は、テキストデータ T 2 及び T 3 内で特定の単語列（T 2 では「お前、ぶっ殺すぞ」、T 3 では「つつってん」）を検出するので、当該特定の単語列を除去又は置換したテキストデータ T 2 ' 及び T 3 ' を音声合成部 1 0 5 に出力する。例えば、テキストデータ T 2 ' では、テキストデータ T 2 内の特定の単語列が空白（ ）に置換される。また、テキストデータ T 3 ' では、テキストデータ T 3 内の特定の単語列「つつってん」が「という」に置換される。音声合成部 1 0 5 は、テキストデータ T 1、T 2 及び T 3 からそれぞれ合成音声信号 S 1、S 2 ' 及び S 3 ' を生成する。

10

【 0 0 3 7 】

感情認識部 1 0 6 は、送受信部 1 0 2 で取得された発話音声信号、音声認識部 1 0 3 で生成されたテキストデータ、及び、送受信部 1 0 2 で受信される主観的評価情報の少なくとも一つに基づいて、顧客の感情情報を生成する。感情認識部 1 0 6 は、発話音声信号に基づいて抽出された音声特徴量（例えば抑揚や音量など）に基づいて顧客の感情情報を生成してよい。感情認識部 1 0 6 は、発話音声信号に基づいて生成されたテキストデータに特定の単語列が検出されたこと、又は、特定の単語列が所定時間以上検出されなかったことに基づいて顧客の感情情報を生成してよい。感情認識部 1 0 6 は、カメラ 1 0 g で取得される顧客の撮像画像に基づいて、顧客の感情情報を生成してもよい。感情認識部 1 0 6 は感情認識モデル 1 0 1 c を用いて顧客の感情情報を生成してもよい。

20

【 0 0 3 8 】

感情認識モデル 1 0 1 c は、発話音声信号、当該発話音声信号から抽出した音声特徴量、当該発話音声信号から生成したテキストデータ、テキスト特徴量又はこれらの少なくとも二つの組み合わせを入力とし、当該発話音声信号に対応する顧客の感情である感情情報を出力するモデルである。

【 0 0 3 9 】

図 5 A は感情認識モデル 1 0 1 c の学習処理の説明図である。例えば、感情認識モデル 1 0 1 c の学習には、発話音声信号から抽出される音声特徴量、テキストデータから抽出されるテキスト特徴量、及び、オペレータによる「主観的評価情報」（又は主観的評価情報から抽出される特徴量）の少なくとも一つをそれぞれ含む複数のデータのセット（以下、「データセット」という）を用いてよい。主観的評価情報は、オペレータが顧客の発話音声信号を聞いて顧客の感情を主観で評価した情報である。例えば、怒りレベル 1 ~ 1 0 のように、オペレータが複数のレベルで顧客の怒りを評価するものであってもよい。感情認識モデル 1 0 1 c を学習するためのデータセットは例えば以下のように生成されてもよい。オペレータは、顧客の生の発話音声信号を聞いて、当該発話音声信号から推定される顧客の感情をアノテーションする（すなわち発話音声信号に対して「主観的評価情報」を付与する）。これにより、発話音声信号と当該発話音声信号から推定される顧客の感情とが時間軸上で関連付けされた情報が得られる。複数のオペレータが複数の発話音声信号に対して主観的評価情報の付与を行うことにより、このような情報の束であるデータセットが得られる。感情認識モデル 1 0 1 c は、このようなデータセットを用いて教師有り機械学習されてもよい。なお、感情認識モデル 1 0 1 c の学習に用いられるデータセットは、音声特徴量に加えて又は代えて発話音声信号を含んでもよいし、テキスト特徴量に加えて又は代えてテキストデータを含んでもよい。

30

40

【 0 0 4 0 】

図 5 B は感情認識モデル 1 0 1 c を用いた推定処理の説明図である。例えば、図 5 B に示すように、発話音声信号 S 1 から抽出した音声特徴量、及び / 又は、当該発話音声信号 S 1 から生成したテキストデータ T 1 から抽出したテキスト特徴量を感情認識モデル 1 0 1 c に入力することにより、入力に対応する出力、すなわち発話音声信号に対応する感情

50

情報が得られる。なお、感情認識モデル 101c には、音声特徴量に加えて又は代えて発話音声信号 S1 が入力されてもよいし、テキスト特徴量に加えて又は代えてテキストデータ T1 が入力されてもよい。

【0041】

主観的評価情報は、一以上の感情（例えば、「幸福」、「驚き」、「恐怖」、「怒り」、「嫌悪」及び「悲しみ」の少なくとも一つ等）の度合を数値で示すものであってもよい。又は、感情情報は、顧客が感じている可能性が高い特定の感情（例えば、「怒り」）を示すものであってもよい。

【0042】

ストレス認識部 107 は、オペレータのストレス状況に関する情報（以下、「ストレス情報」という）を生成する。例えば、ストレス認識部 107 は、オペレータの心拍数、発汗量、呼吸量などのバイタルデータあるいは、カメラを用いて収集したオペレータの視線、表情などの画像情報に基づいて、従来周知の方法によってオペレータのストレス状況を推定してよい。例えば、ストレス認識部 107 は、オペレータによる発話音声に基づいてオペレータのストレス状況を推定してよい。具体的には、ストレス認識部 107 は、オペレータの発話のトーンやスピードの変化、謝罪に関する単語の出現、顧客の発言に被せて発言すること等に基づいて、オペレータのストレス状況を推定してよい。例えば、ストレス認識部 107 は、オペレータ端末 20 の操作ログに基づいてオペレータのストレス状況を推定してよい。具体的には、ストレス認識部 107 は、マウス等の動きや、操作すべき場面で操作入力が無いことなどに応じて、オペレータのストレス状況を推定してよい。ストレス認識部 107 は、ストレス認識モデル 101d に基づいてストレス情報を生成してよい。ストレス認識モデル 101d は、発話音声信号、当該発話音声信号から抽出した音声特徴量、当該発話音声信号から生成したテキストデータ、テキスト特徴量又はこれらの少なくとも二つの組み合わせを入力とし、当該発話音声を聞いているオペレータが感じるストレスの推定値を出力するモデルである。ストレス認識モデル 101d の学習には、顧客の発話音声を聞いてオペレータが実際に感じたストレスの実測値を用いてよい。ストレス認識モデル 101d を学習するためのデータセットは例えば以下のように生成されてもよい。オペレータは、顧客の発話音声を聞いて感じたストレスの度合い（例えば 1 ~ 10 のようなレベル）をアノテーションする（すなわち発話音声信号に対して自身が感じた「ストレスの度合い」を付与する）。これにより、発話音声信号と当該発話音声信号を聞いた際のオペレータのストレスとが時間軸上で関連付けされた情報が得られる。複数のオペレータが複数の発話音声信号に対してストレスの度合いの付与を行うことにより、このような情報の束であるデータセットが得られる。ストレス認識モデル 101d は、このようなデータセットを用いて教師有り機械学習されてもよい。

【0043】

制御部 108 は、音声処理装置 10 に関する種々の制御を行う。具体的には、制御部 108 は、ストレス認識部 107 において生成されるストレス情報に基づいて、オペレータ端末 20 において音声合成部 105 で生成された合成音声又は顧客の発話音声のどちらを出力するかを切り替えてもよい。制御部 108 は、発話音声信号に基づいて合成音声信号を生成するか否かをストレス情報に基づいて切り替えてもよい。例えば、制御部 108 は、ストレス情報が示すストレス度数が所定の閾値以上又はより大きい場合、顧客の発話音声ではなく合成音声をオペレータに出力するように制御してもよい。一方、制御部 108 は、ストレス情報が示すストレス度数が所定の閾値より小さい又は以下である場合、発話音声をオペレータに出力するように制御してもよい。制御部 108 は、オペレータから感情抑制機能の自動切り替えについての指示情報が入力された場合、ストレス情報に基づいて上記切り替えを行ってもよい。感情抑制機能とは、顧客の発話音声に代えて合成音声をオペレータに出力する機能である。

【0044】

制御部 108 は、感情情報に基づいて上記切り替えを行ってもよい。制御部 108 は、当該切り替えを感情抑制切替モデル 101e の出力に基づいて行ってもよい。感情抑制切

10

20

30

40

50

替モデル 101e は、発話音声信号、音声特徴量、テキストデータ、テキスト特徴量又はこれらの少なくとも二つの組み合わせを入力として、感情抑制機能のオン・オフを切り替えるタイミングを出力とするモデルである。感情抑制切替モデル 101e は更にストレス情報又は感情情報を入力としてもよい。感情抑制切替モデル 101e の詳細については後述する。

【0045】

また、制御部 108 は、オペレータによって入力される切り替え情報に基づいて上記切り替えを行ってもよい。ここで、切り替え情報は、顧客の感情抑制機能の適用（オン）又は非適用（オフ）の切り替えに関する情報である。例えば、制御部 108 は、切り替え情報が顧客の感情抑制機能の適用を示す場合、合成音声をオペレータに出力するように制御してもよい。一方、制御部 108 は、切り替え情報が顧客の感情抑制機能の非適用を示す場合、発話音声をオペレータに出力するように制御してもよい。制御部 108 は、オペレータから感情抑制機能の手動切り替えについての指示情報が入力された場合、上記切り替え情報に基づいて上記切り替えを行ってもよい。

【0046】

学習部 109 は、感情認識モデル 101c、ストレス認識モデル 101d 及び感情抑制切替モデル 101e の学習処理を行ってよい。

【0047】

音声処理装置 10 は、以下 1) 乃至 7) に示すいずれかの情報、又は、少なくとも二つの情報の組み合わせを時間軸上で関連付け、送受信部 102 を介して、オペレータ端末 20 に対して送信してよい。1) 顧客の発話音声信号、2) 発話音声信号から生成されたテキストデータ、3) 除去部 104 の処理を経たあとのテキストデータ、4) 検出情報、5) 合成音声信号、6) 顧客の発話音声信号から推定される顧客の感情情報、7) 感情抑制機能のオン・オフを切り替えるタイミング。感情抑制機能がオンである場合、音声処理装置 10 は顧客の発話音声信号をオペレータ端末 20 に送らなくてもよい。感情抑制機能がオフである場合、音声処理装置 10 は合成音声信号をオペレータ端末 20 に送らなくてもよい。感情抑制機能のオン・オフに関わらず、音声処理装置 10 は顧客の発話音声信号と合成音声信号との両方をオペレータ端末 20 に送ってもよい。

【0048】

オペレータ端末

【0049】

図 6 は、本実施形態に係るオペレータ端末の機能構成の一例を示す図である。オペレータ端末 20 は、送受信部 201、入力受付部 202、制御部 203 を備える。なお、図 6 に示す機能構成は一例にすぎず、図示しない他の構成を備えてもよい。

【0050】

送受信部 201 は、音声処理装置 10 及び / 又は顧客端末 30 との間で、種々の情報及び / 又は信号を送信及び / 又は受信する。例えば、送受信部 201 は、顧客端末 30 で収音された顧客の発話音声の信号である発話音声信号を受信してもよい。送受信部 102 は、音声処理装置 10 から、合成音声信号を受信してもよい。また、送受信部 201 は、音声処理装置 10 に対して、主観的評価情報を送信してもよい。また、送受信部 201 は、音声処理装置 10 から、顧客の感情情報を受信してもよい。

【0051】

入力受付部 202 は、オペレータによる入力部 10e の操作に基づいて、種々の情報の入力を受け付ける。例えば、入力受付部 202 は、感情認識モデル 101c やストレス認識モデル 101d を学習するためのデータセットを生成するための作業の一環として、顧客の生の発話音声信号に対して主観的評価情報やストレスの度合いの入力を受け付けてもよい。以降、オペレータが、オペレータ端末 20 において主観的評価情報やストレスの度合いを入力する作業を「アノテーション作業」と呼ぶ。アノテーション作業は、通常のコールセンター業務とは別の業務として位置付けられていてもよい。また、入力受付部 202 は、顧客の感情抑制機能の切り替え情報の入力を受け付けてもよい。また、入力受付部

10

20

30

40

50

202は、感情抑制機能の手動切り替え又は自動切り替えのどちらかを指示する指示情報の入力を受け付けてもよい。

【0052】

制御部203は、オペレータ端末20に関する種々の制御を行う。例えば、制御部203は、表示部10fにおける情報及び/又は画像の表示を制御する。また、制御部203は、音声出力部10iにおける音声の出力を制御する。制御部203は、音声処理装置10から送信される情報に基づいて音声の出力を制御してもよいし、入力受付部202が受け付けた情報に基づいて音声の出力を制御してもよい。

【0053】

制御部203は、音声処理装置10から受信した合成音声信号に基づいて合成音声を音声出力部10iから出力させる。制御部203は、顧客端末30からの発話音声信号に基づいて発話音声を音声出力部10iから出力させてもよい。

10

【0054】

また、制御部203は、音声処理装置10から受信した感情情報に基づいて、合成音声信号に対応する感情情報を表示部10fに表示させてもよい。また、制御部203は、音声処理装置10から受信した合成音声信号に対応するテキストデータを表示部10fに表示させてもよい。例えば、制御部203は、感情情報、テキストデータ及び検出情報の少なくとも一つを含む画面D1を表示部10fに表示させてもよい。また、制御部203は、ストレス情報を表示部10fに表示させてもよい。例えば、制御部203は、ストレス情報を含む画面D2を表示部10fに表示させてもよい。

20

【0055】

図7は、本実施形態に係る画面D1の一例を示す図である。図7に示すように、画面D1において、制御部203は、音声出力部10iからの合成音声の出力タイミングTに合わせて、感情情報I1を表示部10fに表示させてもよい。合成音声の出力タイミングT毎に感情情報I1を表示させることにより、オペレータは、感情抑制機能により顧客の感情が抑制された合成音声を聞く場合でも、顧客の感情をリアルタイムで認識することができる。

【0056】

また、画面D1において、制御部203は、当該合成音声の出力タイミングTに合わせて、当該合成音声に対応するテキストデータI2の内容を表示部10fに表示させてもよい。テキストデータI2の内容を表示させることにより、オペレータは、合成音声だけでなく、視覚的にも顧客の発話内容を把握可能となる。

30

【0057】

また、画面D1では、制御部203は、音声処理装置10から受信した検出情報に基づいて、特定の単語列そのものの表示に代えて、特定の単語列の検出を示す情報I3（例えば、「NGワード検出」）を表示部10fに表示させてもよい。この機能を「NGワード非表示機能」と呼ぶ。これにより、心理的悪影響を与える顧客の発話の内容をそのままオペレータに認識させるのを回避できるのでオペレータのストレスを抑制できる。また、当該発話があったことはオペレータに通知できるので、オペレータが顧客に対する応対を適切に行うことができる。

40

【0058】

また、画面D1において、制御部203は、音声処理装置10からの感情情報に基づいて、合成音声の出力タイミングT毎に、顧客の特定の感情のレベルI4を時系列に表示部10fに表示させてもよい。例えば、図7では、合成音声の出力タイミングT毎の顧客の「怒り」のレベルI4が折れ線グラフで示される。これにより、オペレータが顧客の特定の感情（例えば、「怒り」）の遷移を容易に把握できるので、顧客に対するオペレータの応対の満足度を向上できる。

【0059】

画面D1において、制御部203は選択ボタンI5を表示部10fに表示させてもよい。選択ボタンI5は、感情抑制機能の適用（オン）又は非適用（オフ）を自動又は手動のど

50

ちらで切り替えるかをオペレータが選択可能とするインターフェースである。オペレータは選択ボタン I 5 に対してクリック、タップ又はスライド等の操作を行うことにより「自動切替モード」と「手動切替モード」を切り替えることができる。自動切替モードにおいては、例えば感情情報、ストレス情報、又は感情抑制切替モデル 1 0 1 e からの出力等に基づいて感情抑制機能のオン・オフが自動で切り替わる。

【 0 0 6 0 】

「手動切替モード」が選択された場合、制御部 2 0 3 は、感情抑制機能の適用又は非適用をオペレータが選択可能とするインターフェースである切替ボタン I 6 を表示部 1 0 f に表示させてよい。オペレータが感情抑制機能のオンとオフを切り替えたタイミングは、顧客の発話音声（及び/又は発話音声に基づいて抽出される各種特徴量）と時間軸上で関連付けられて「手動切替履歴データ」として不図示の記憶部に蓄積される。「手動切替履歴データ」には更にオペレータの識別情報が関連付けられてもよい。

10

【 0 0 6 1 】

切り替えボタン I 7 は、「NGワード非表示機能」のオン・オフを切り替えるためのボタンである。「NGワード非表示機能」がオフの場合には、テキストデータ I 2 の内に特定の単語列が検出された場合でも、除去部 1 0 4 による処理が行われる前のテキストデータ I 2 がそのまま表示部 1 0 f に表示される。感情抑制機能をオンにしつつNGワード非表示機能をオフにした場合、オペレータは顧客による特定の単語列を直接聞くことは無いのでストレスが軽減される一方で、顧客の発話内容を正確に把握することにより顧客の感情をより正確に把握することができる。

20

【 0 0 6 2 】

感情抑制切替モデル 1 0 1 e を学習するためのデータセットは、ストレス情報、感情情報、発話音声信号 S 1、音声特徴量、テキストデータ、テキスト特徴量又はこれらの少なくとも二つの組み合わせと、オペレータが感情抑制機能のオン・オフを切り替えたタイミングとが、時間軸上で関連付けされたデータの束であってよい。感情抑制切替モデル 1 0 1 e を学習する方法は、例えば下記 1) から 3) に述べるような様々な方法がある。1) 感情抑制切替モデル 1 0 1 e はオペレータ毎に学習されてもよい。すなわち、或るオペレータに対して適用される感情抑制切替モデル 1 0 1 e は、そのオペレータによる感情抑制機能の「手動切替履歴データ」のみに基づいて学習されてもよい。この方法によれば、感情抑制切替モデル 1 0 1 e はそのオペレータの好みに合わせたタイミングで感情抑制機能を切り替えることができるようになる。あるいは、2) 或るオペレータに対して適用される感情抑制切替モデル 1 0 1 e は、不特定多数のオペレータによる「手動切替履歴データ」に基づいて学習されてもよい。この方法によれば、学習に用いることができるデータが多くなるため、感情抑制切替モデル 1 0 1 e を早く学習することができるようになる。あるいは、3) 或るオペレータに対して適用される感情抑制切替モデル 1 0 1 e は、そのオペレータと年齢・性別・その他の特性が類似したオペレータによる「手動切替履歴データ」に基づいて学習されてもよい。この方法によれば、1) の方法と比較して学習に用いることができるデータが多くなるため感情抑制切替モデル 1 0 1 e を早く学習することができ、2) の方法と比較して自分の好みに合った切替タイミングを学習することができるようになる。

30

40

【 0 0 6 3 】

図 8 は、本実施形態に係る画面 D 2 の一例を示す図である。画面 D 2 において、制御部 2 0 3 は、音声処理装置 1 0 からのストレス情報を表示させてもよい。例えば、図 8 では、ストレス情報として、オペレータが感じるストレスの推定値を示す情報（例えば、「56 %」）と、当該オペレータの平常時の状態からの相対的な評価値を示す情報（例えば、「平常時より 8 . 1 % 減」）とが表示される。

【 0 0 6 4 】

図 1 2 は、本実施形態に係る画面 D 3 の一例を示す図である。画面 D 3 において、制御部 2 0 3 は、オペレータがアノテーション作業を行うためのインターフェース I 8 を表示させてもよい。オペレータは、例えば、顧客の生の音声（サンプル音声）を聞きながら、

50

サンプル音声から感じられる顧客の感情をインターフェース I 8 から都度選択する。図 12 において、顧客感情 I 1 はオペレータによる顧客感情の主観的評価情報である。例えば、オペレータが、サンプル音声「今日の夕方までにどうにかして届けてよ」に対して「怒り」という感情をアノテーションしたならば、図 12 に示すように、「今日の夕方までにどうにかして届けてよ」というサンプル音声と「怒り」という情報が時間軸上で関連付けられる。アノテーションは文単位で行われてもよいし所定の時間間隔ごとに行われてもよい。

【0065】

(音声処理システムの動作)

図 9 は、本実施形態に係る感情抑制動作の一例を示すフローチャートである。なお、図 9 は、例示にすぎず、少なくとも一部のステップ(例えば、ステップ S 106)の順番は入れ替えられてもよいし、不図示のステップが実施されてもよいし、一部のステップが省略されてもよい。

10

【0066】

音声処理装置 10 は、顧客端末 30 の音声入力部 10h で収音される顧客の発話音声の信号である発話音声信号を取得する(S 101)。

【0067】

音声処理装置 10 は、S 101 で取得された発話音声信号に基づいて抽出される特徴量を音声認識モデル 101a に入力して、一以上の単語からなる単語列を含むテキストデータを生成する(S 102)。

20

【0068】

音声処理装置 10 は、S 102 で生成されたテキストデータ内に特定の単語列が含まれるか否かを判定する(S 103)。当該テキストデータ内に特定の単語列が含まれる場合、音声処理装置 10 は、当該特定の単語列を除去又は前記特定の単語列を他の単語列に変換したテキストデータを生成する(S 104)。

【0069】

音声処理装置 10 は、テキストデータに基づいて抽出される特徴量を音声合成モデル 101b に入力して、合成音声の信号である合成音声信号を生成する(S 105)。

【0070】

音声処理装置 10 は、S 101 で取得された発話音声信号、S 102 で生成されたテキストデータ、及び、オペレータによって入力される顧客の感情の主観的評価情報の少なくとも一つに基づいて抽出される特徴量を感情認識モデル 101c に入力して、顧客の感情情報を生成する(S 106)。

30

【0071】

オペレータ端末 20 は、S 105 で生成された合成音声信号に基づいて合成音声を音声出力部 10i から出力させるとともに、当該合成音声の出力タイミング T に合わせて当該合成音声に対応する感情情報を表示部 10f に表示させる(S 107、例えば、図 7)。

【0072】

音声処理装置 10 は、処理を終了するか否かを判定する(S 108)。処理を終了しない場合(S 108: NO)、音声処理装置 10 は、処理 S 101 ~ S 107 を再び実行する。一方、音声変換処理を終了する場合(S 108: YES)、音声処理装置 10 は、処理を終了する。

40

【0073】

図 10 は、本実施形態に係る感情抑制機能の自動切り替え動作を示すフローチャートである。なお、図 10 は、例示にすぎず、少なくとも一部のステップの順番は入れ替えられてもよいし、不図示のステップが実施されてもよいし、一部のステップが省略されてもよい。

【0074】

音声処理装置 10 は、オペレータのストレス情報を生成する(S 201)。

【0075】

50

音声処理装置 10 は、ストレス情報が所定の条件を満たすか否かを判定する (S 202)。例えば、所定の条件は、ストレス情報が示すストレス度数が所定の閾値以上又はより大きいことであってもよい。

【0076】

音声処理装置 10 は、ストレス情報が所定の条件を満たす場合 (S 202: YES)、感情抑制機能を適用 (すなわち、オペレータ端末 20 から合成音声出力) してもよい (S 203)。一方、音声処理装置 10 は、ストレス情報が所定の条件を満たさない場合 (S 202: NO)、感情抑制機能を非適用 (すなわち、オペレータ端末 20 から顧客の発話音声出力) してもよい (S 204)。

【0077】

音声処理装置 10 は、処理を終了するか否かを判定する (S 205)。処理を終了しない場合 (S 205: NO)、音声処理装置 10 は、処理 S 201 ~ S 204 を再び実行する。一方、音声変換処理を終了する場合 (S 205: YES)、音声処理装置 10 は、処理を終了する。なお、S 201 及び S 202 において、音声処理装置 10 は、感情情報や感情抑制切替モデル 101e の出力に基づいて、感情抑制機能を適用するか否を決定してもよい。

【0078】

以上のように、本実施形態に係る音声処理システム 1 によれば、顧客の発話音声信号に基づいてテキストデータを生成し、当該テキストデータに基づいて生成される合成音声をオペレータに出力する。このため、顧客の発話音声に含まれる顧客の感情を十分に抑制した合成音声をオペレータに聞かせることができ、顧客の感情的発話に起因するオペレータのストレスを軽減できる。本発明の発明者は、約 50 名の被験者に対して、1) 顧客の発話音声そのもの、2) 顧客の発話音声の音量を調整した音声、3) 顧客の発話音声の声質を変換した音声、4) 顧客の発話音声をテキスト化してから生成した合成音声、の 4 種類の音声を聞き比べてもらい、音声から感じられる怒りの度合いを 7 段階の尺度で評価してもらう実験を行った。その結果、2) や 3) と比較して 4) が、被験者に伝わった怒りの軽減度合いが顕著であった。

【0079】

また、本実施形態に係る音声処理システム 1 によれば、オペレータに対して、合成音声出力するだけでなく顧客の感情情報を合成音声出力のタイミングに合わせて通知することができるので、合成音声を聞いたオペレータが顧客の感情をリアルタイムに認識でき、顧客に対して適切な対応を行うことができる。

【0080】

また、本実施形態に係る音声処理システム 1 によれば、オペレータのストレス情報又は顧客の感情情報等に基づいて、感情抑制機能を適用するか否か (すなわち、オペレータに対して合成音声又は発話音声のどちらを出力するか) が切り替えられるので、オペレータのストレスと顧客の満足度とのバランスを適切に図ることができる。

【0081】

(変更例)

上記音声処理システム 1 では、音声認識部 103 は、発話音声信号から、一つ又は複数の文として確定された単語列を含むテキストデータを生成したが、これに限られない。音声認識部 103 は、発話音声信号から認識された単語列が一つ又は複数の文として確定される前に、一つ又は複数の単語 (品詞又は形態素) からなる単語列を含むテキストデータを生成してもよい。除去部 104 は、当該文として確定されていないテキストデータ内の特定の単語列を除去し、音声合成部 105 は、当該文として確定されていないテキストデータから合成音声信号を生成してもよい。

【0082】

図 11 は、本実施形態の変更例に係る合成音声信号の生成の一例を示す図である。図 11 では、送受信部 102 で取得された発話音声信号 S4 に基づいて、音声認識部 103 においてテキストデータ T41 ~ T43 が生成されるものとする。図 11 に示すように、テ

10

20

30

40

50

キストデータ T 4 1 ~ T 4 3 は、「はやく送ってください」という一文の確定前に、意味を持つ形態素単位（「はやく」、「送って」、「ください」）でテキストデータが生成される点で、図 4 と異なる。除去部 1 0 4 は、テキストデータ T 4 1 ~ T 4 3 それぞれに対して特定の単語列が含まれるか否かを判定して、当該特定の単語列を除去して音声合成部 1 0 5 に出力する。音声合成部 1 0 5 は、テキストデータ T 4 1 ~ T 4 3 からそれぞれ合成音声信号 S 4 1 ~ S 4 3 を生成する。

【 0 0 8 3 】

図 1 1 に示すように、文の確定前に一つ又は複数の形態素単位でテキストデータを生成して合成音声を出力することにより、テキストデータの生成によりオペレータの応答遅延を軽減できる。なお、形態素単位での複数のテキストデータ（又は合成音声）が意味的に不自然でないかを判定するモデルなどが用いられてもよい。

10

【 0 0 8 4 】

また、応答遅延を軽減するために、図 4 に示す合成音声信号 S 1 ~ S 3、図 1 1 に示す合成音声信号 S 4 1 ~ S 4 3 それぞれの前及び / 又は後に、例えば、「あ～」、「え～」、「まあ」等のフィラー音が追加されてもよい。これにより、オペレータも応答遅延による顧客の満足度の低下を防止できる。

【 0 0 8 5 】

また、音声合成部 1 0 5 は、感情認識部 1 0 6 が推定した顧客の感情に基づいて、複数の音声合成モデル 1 0 1 b のうちから顧客の感情に合った音声合成モデル 1 0 1 b を選択してもよい。例えば、感情認識部 1 0 6 が推定した顧客の感情が「激昂」である場合、音声合成部 1 0 5 は、ピッチが速く抑揚が激しい音声合成モデル 1 0 1 b を用いてよい。例えば、感情認識部 1 0 6 が推定した顧客の感情が「号泣」である場合、音声合成部 1 0 5 は、泣き声のような音声を出力する音声合成モデル 1 0 1 b を用いてよい。或いは、音声合成部 1 0 5 は、感情認識部 1 0 6 が推定した顧客の感情に基づいて音声合成モデル 1 0 1 b のパラメータを変更し、顧客の感情に合った音声が出力されるように調整してよい。顧客が激昂している際の生の音声を直接聞いたオペレータは極めて強いストレスを感じてしまう。他方、オペレータは顧客対応業務を適切に遂行するために、顧客の感情をリアルタイムで正確に把握する必要がある。オペレータに発話音声を直接聞かせないことによりオペレータは過剰なストレスを感じることもなく、合成音声に顧客の感情を乗せることにより、オペレータは聴覚を通じて顧客の感情をリアルタイムに把握することができる。

20

30

【 0 0 8 6 】

（その他の実施形態）

上記実施形態では、顧客の発話音声信号をテキスト化して、合成音声信号をオペレータに出力するものとしたがこれに限られない。音声処理装置 1 0 は、顧客の発話音声信号に基づいて抽出される音声特徴量を音声変換モデルに入力して、変換音声の信号を生成し、オペレータ端末 2 0 から変換音声を出力してもよい。

【 0 0 8 7 】

特許請求の範囲に記載の「音声変換モデル」は、発話音声信号を一旦テキスト化して合成音声として出力するモデルと、発話音声信号をテキスト化せずに声質を変換させて出力するモデルとの両方を包含する概念である。顧客の発話音声に代えて合成音声または変換音声をオペレータに対して出力することにより、効果の程度の差こそあれ、オペレータが感じるストレスを軽減できる。他方で、顧客対応業務の遂行のためには、オペレータが顧客の感情をリアルタイムに把握することも欠かせない。

40

【 0 0 8 8 】

本変形例における音声処理装置 1 0 は、顧客の発話音声信号に基づいて抽出される音声特徴量を音声変換モデルに入力して、変換音声信号を生成する。音声処理装置 1 0 は、1) 変換音声信号と、2) 顧客の発話音声から推定される顧客の感情情報とを時間軸上で関連付けた情報を生成し、オペレータ端末 2 0 に対して送信する。音声処理装置 1 0 が送信する情報には、発話音声信号、発話音声信号から生成されたテキストデータ、除去部 1 0 4 の処理を経たあとのテキストデータ、検出情報、感情抑制機能のオン・オフを切り替え

50

るタイミングが関連付けられていてもよい。

【0089】

オペレータ端末20は、音声処理装置10から受信した変換音声の信号を音声出力部10iから出力し、且つ、音声出力部10iからの変換音声の出力タイミングTに合わせて、感情情報を示す情報を表示部10fに表示してよい。オペレータ端末20は更に、音声出力部10iからの変換音声の出力タイミングTに合わせて、テキストデータを表示部10fに表示してよい。かかる表示の態様は図7に図示するようであってよい。

【0090】

本変形例における音声処理装置10は、感情情報に基づいて、感情情報が示す感情が変換音声に反映されるように、変換音声の信号を生成してもよい。例えば感情情報が示す感情が「激昂」である場合、ピッチが速く抑揚が激しい音声変換モデルを用いてよい。例えば感情情報が示す感情が「号泣」である場合、泣き声のような音声を入力する音声変換モデルを用いてよい。音声処理装置10は、感情情報が示す感情が変換音声に反映されるように、変換音声の信号を生成してよい。オペレータに発話音声の直接聞かせないことによりオペレータは過剰なストレスを感じるがことなく、変換音声に顧客の感情を乗せることにより、オペレータは聴覚を通じて顧客の感情をリアルタイムに把握することができる。

10

【0091】

本変形例における音声処理システム1においては、オペレータによるアノテーション作業は、オペレータによる通常のコールセンター業務中において、変換音声に対して行われてもよい。オペレータが変換音声に対して「怒りの感情」をアノテーションした場合、当該アノテーションの結果に基づいて、音声変換モデルがより柔らかい音声を入力するようにリアルタイムに調整されてもよい。

20

【0092】

以上説明した実施形態では、第1のユーザが顧客であり、第2のユーザがオペレータであるコールセンターを想定したが、本実施形態の適用場面はコールセンターに限られない。例えば、Webミーティング等、第1のユーザの感情を抑制した音声を入力する第2のユーザに出力するどのような場面にも適用可能である。すなわち、本実施形態は、カスタマーハラスメント対策だけでなく、社内のパワーハラスメント等、様々なハラスメントに対する企業側の対策として利用可能である。

【0093】

以上説明した実施形態における、感情情報と合成音声とを「時間軸上で関連付け」する処理は、図7に示すように、合成音声または変換音声の出力タイミングに合わせて、それらの元となった発話音声から推定される感情情報を表示することが実現可能な態様であれば、その具体的な態様を問わない。以上説明した実施形態における「時間軸上で関連付け」する処理は、何時何分何秒といった時刻情報に基づいて関連付けする処理でもよいし、発話音声情報の開始から何分何秒経過時といった情報に基づいて関連付けする処理でもよいし、文単位、単語単位又は形態素単位で関連付けする処理でもよい。

30

【0094】

以上説明した実施形態における音声処理システム1において、顧客からは、自身の音声感情抑制されてオペレータに届いていることが分からないようにしてもよい。すなわち、感情抑制機能がオンになっているかオフになっているかは、顧客からは把握できないようにしてもよい。

40

【0095】

アノテーション作業は、オペレータがオペレータ端末20上で行ってもよいし、別途、アノテーション作業用の専用のアプリケーションや端末が用意されていてもよい。

【0096】

また、以上説明した実施形態は、本発明の理解を容易にするためのものであり、本発明を限定して解釈するためのものではない。実施形態が備える各要素並びにその配置、材料、条件、形状及びサイズ等は、例示したものに限定されるわけではなく適宜変更することができる。また、異なる実施形態で示した構成同士を部分的に置換し又は組み合わせるこ

50

とが可能である。また、音声処理装置 10 の機能として記載した機能をオペレータ端末 20 が備えていてもよい。また、オペレータ端末 20 の機能として記載した機能を音声処理装置 10 が備えていてもよい。

【符号の説明】

【0097】

1 ... 音声処理システム、10 ... 音声処理装置、20 ... オペレータ端末、30 ... 顧客端末、10a ... プロセッサ、10b ... RAM、10c ... ROM、10d ... 通信部、10e ... 入力部、10f ... 表示部、10g ... カメラ、10h ... 音声入力部、10i ... 音声出力部、101 ... 記憶部、102 ... 送受信部、103 ... 音声認識部、104 ... 除去部、105 ... 音声合成部、106 ... 感情認識部、107 ... ストレス認識部、108 ... 制御部、109 ... 学習部、201 ... 送受信部、202 ... 入力受付部、203 ... 制御部

10

20

30

40

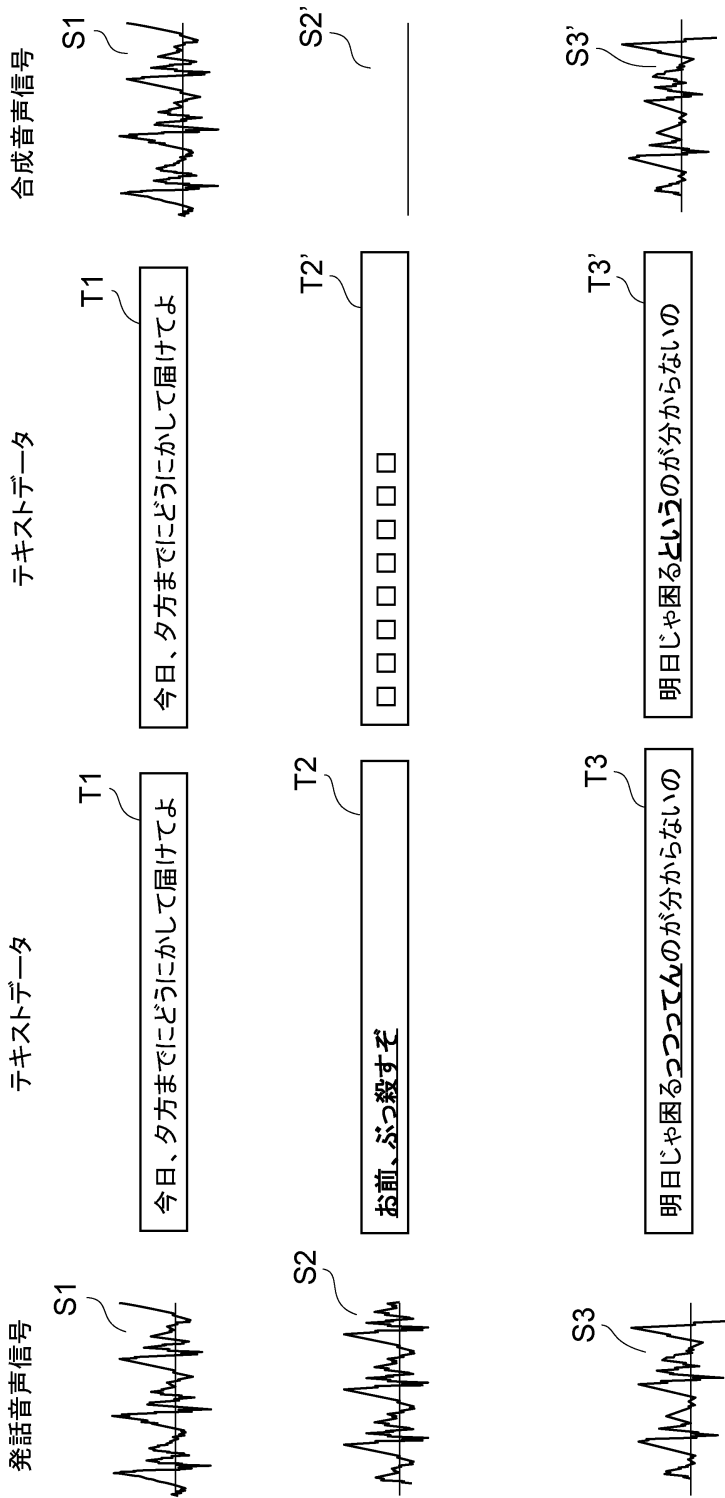
50

【要約】

【課題】聞き手のストレスの軽減を可能とすること。

【解決手段】音声処理システム 1 は、第 1 のユーザの発話音声の信号である発話音声信号を取得する取得部と、前記発話音声信号に基づいて抽出される特徴量を音声認識モデルに入力して、一以上の単語からなる単語列を含むテキストデータを生成する音声認識部と、前記テキストデータに基づいて抽出される特徴量を音声合成モデルに入力して、合成音声の信号である合成音声信号を生成する音声合成部と、第 2 のユーザに対して前記合成音声を出力する音声出力部と、を備える。

【選択図】図 4



10

20

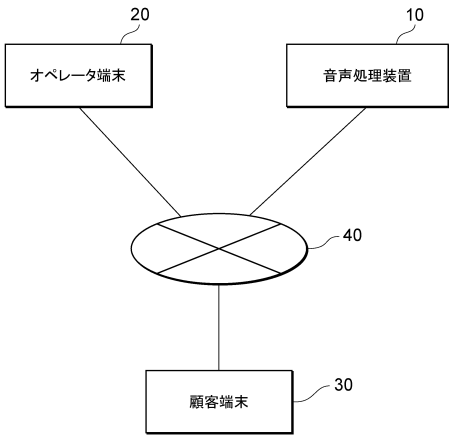
30

40

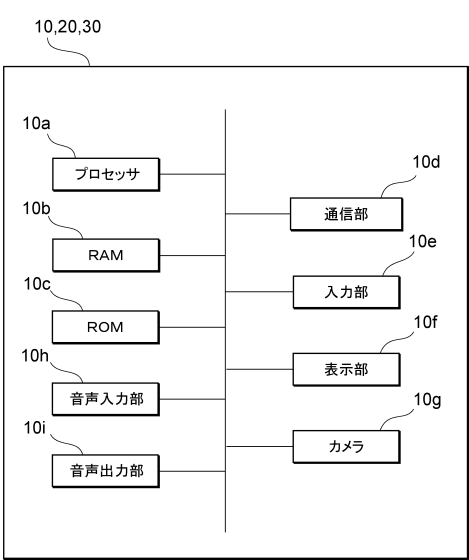
50

【図面】

【図 1】



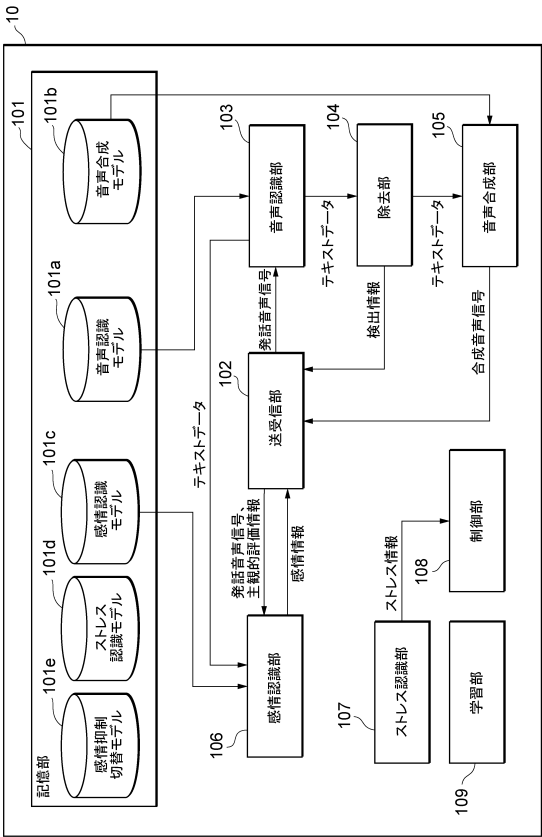
【図 2】



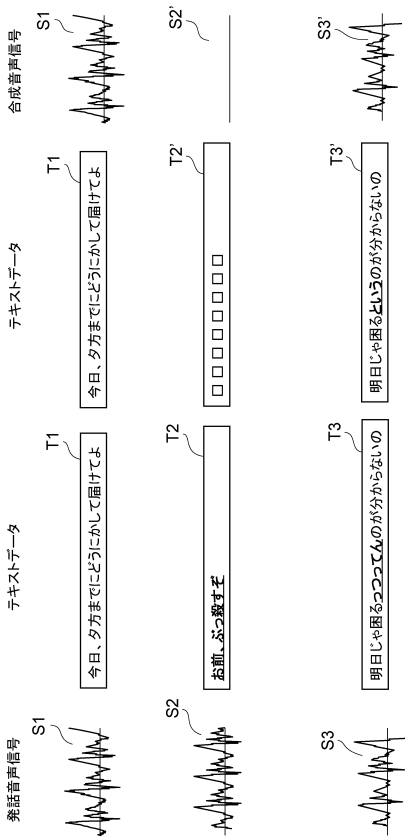
10

20

【図 3】



【図 4】



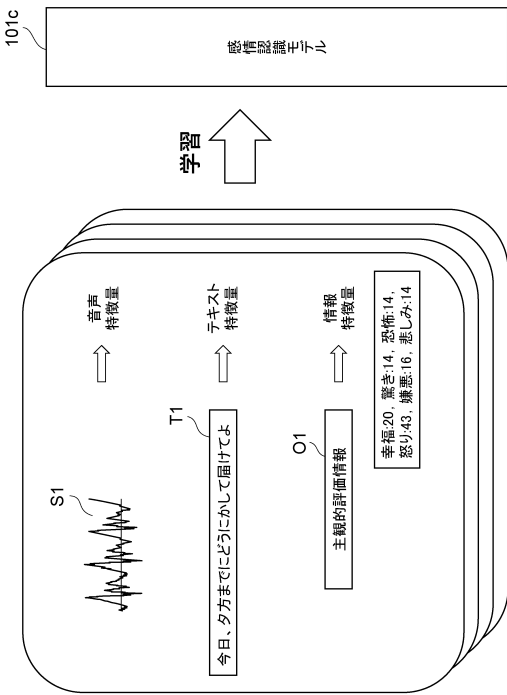
30

40

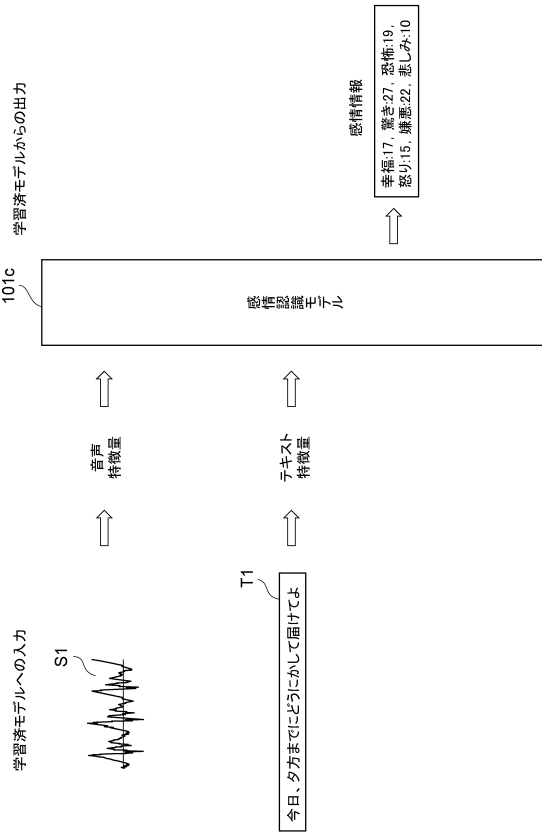
50

【図 5 A】

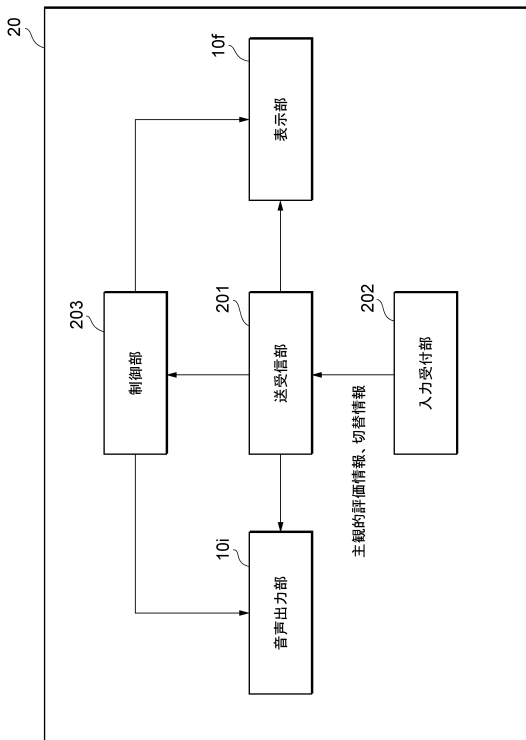
感情認識モデルを学習するためのデータセット



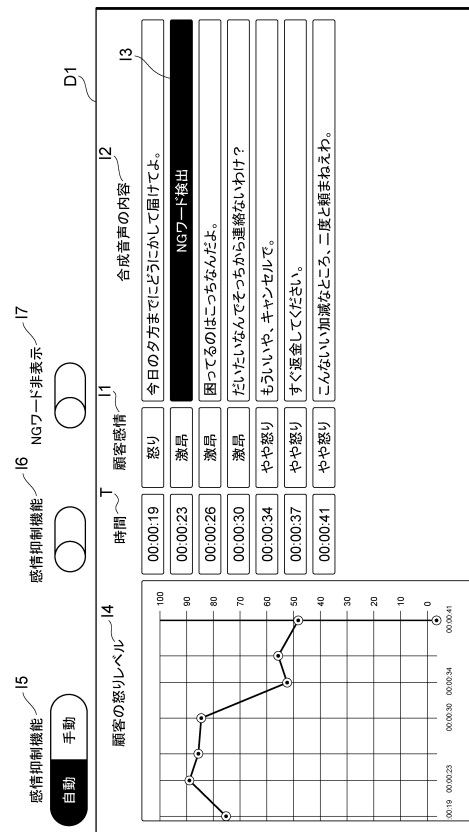
【図 5 B】



【図 6】



【図 7】



10

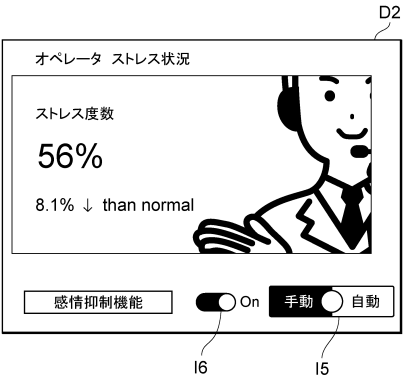
20

30

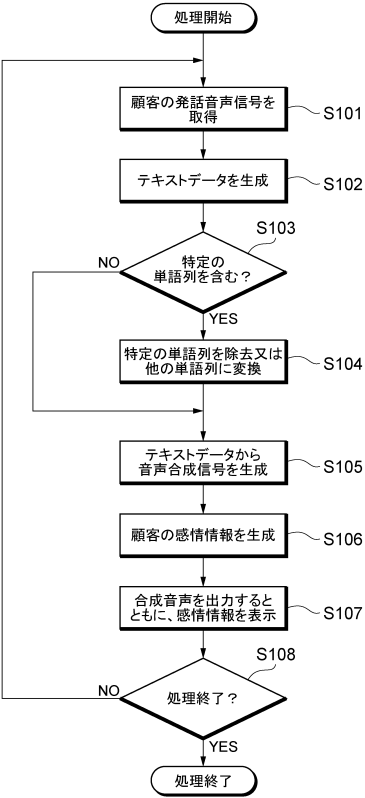
40

50

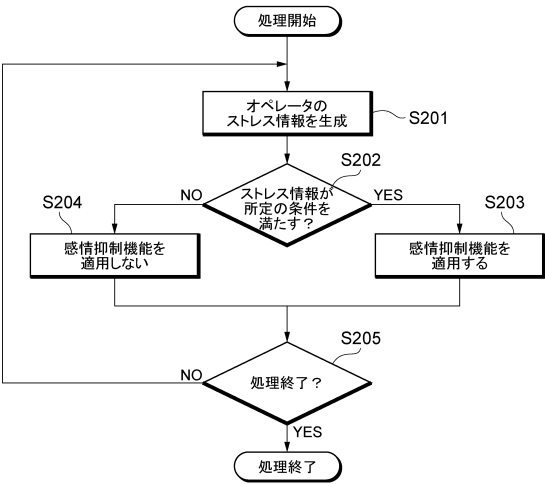
【図 8】



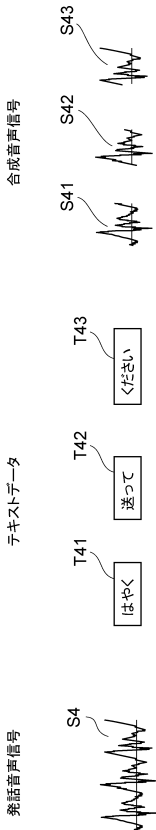
【図 9】



【図 10】



【図 11】



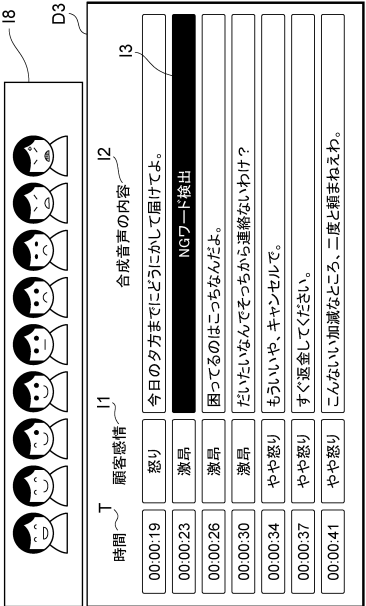
10

20

30

40

50



フロントページの続き

(72)発明者 今村 俊雄
東京都港区海岸一丁目7番1号

(72)発明者 阪下 啓祐
東京都港区海岸一丁目7番1号

(72)発明者 高 原 周平
東京都港区海岸一丁目7番1号

審査官 堀 洋介

(56)参考文献 特開2001-117752(JP,A)
特開2020-021025(JP,A)
特開2010-166324(JP,A)
国際公開第2019/111346(WO,A1)
特開2019-110451(JP,A)

(58)調査した分野 (Int.Cl., DB名)
G10L 13/00 - 15/34
G06F 3/16