

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7408518号  
(P7408518)

(45)発行日 令和6年1月5日(2024.1.5)

(24)登録日 令和5年12月22日(2023.12.22)

(51)国際特許分類		F I	
G 0 6 N	20/00 (2019.01)	G 0 6 N	20/00 1 3 0
G 0 6 N	3/08 (2023.01)	G 0 6 N	3/08
G 1 0 L	15/10 (2006.01)	G 1 0 L	15/10 5 0 0 Z
G 1 0 L	15/16 (2006.01)	G 1 0 L	15/16

請求項の数 6 (全24頁)

(21)出願番号	特願2020-155830(P2020-155830)	(73)特許権者	500257300 L I N E ヤフー株式会社 東京都千代田区紀尾井町 1 番 3 号
(22)出願日	令和2年9月16日(2020.9.16)	(74)代理人	110002147 弁理士法人酒井国際特許事務所
(65)公開番号	特開2022-49570(P2022-49570A)	(72)発明者	藤田 悠哉 東京都千代田区紀尾井町 1 番 3 号 ヤフー株式会社内
(43)公開日	令和4年3月29日(2022.3.29)	審査官	金田 孝之
審査請求日	令和3年8月19日(2021.8.19)		

最終頁に続く

(54)【発明の名称】 情報処理装置、情報処理方法、情報処理プログラム、端末装置、推論方法、及び推論プログラム

(57)【特許請求の範囲】

【請求項 1】

機械学習のモデルの学習に用いる音声データである入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記入力用データに含まれるノイズの多寡に基づく分類結果を示す分類ラベルとを含む学習用データを取得する取得部と、

前記学習用データを用いて、データの入力に応じて、前記正解データに対応する第 1 出力と前記分類ラベルに対応し、当該データに含まれるノイズの多寡に基づく第 2 出力とを出力する前記モデルを学習する学習部と、

を備え、

前記取得部は、

前記入力用データが検知された場所の分類結果を示す前記分類ラベルを含む前記学習用データを取得し、

前記学習部は、

前記第 1 出力と、入力されたデータが検知された場所の分類結果を示す前記第 2 出力とを出力する前記モデルを学習する

ことを特徴とする情報処理装置。

【請求項 2】

コンピュータが実行する情報処理方法であって、

機械学習のモデルの学習に用いる音声データである入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記入力用データに含まれるノイズの多寡に基

づく分類結果を示す分類ラベルとを含む学習用データを取得する取得工程と、

前記学習用データを用いて、データの入力に応じて、前記正解データに対応する第1出力と前記分類ラベルに対応し、当該データに含まれるノイズの多寡に基づく第2出力とを出力する前記モデルを学習する学習工程と、

を含み、

前記取得工程は、

前記入力用データが検知された場所の分類結果を示す前記分類ラベルを含む前記学習用データを取得し、

前記学習工程は、

前記第1出力と、入力されたデータが検知された場所の分類結果を示す前記第2出力とを出力する前記モデルを学習する

ことを特徴とする情報処理方法。

#### 【請求項3】

機械学習のモデルの学習に用いる音声データである入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記入力用データに含まれるノイズの多寡に基づく分類結果を示す分類ラベルとを含む学習用データを取得する取得手順と、

前記学習用データを用いて、データの入力に応じて、前記正解データに対応する第1出力と前記分類ラベルに対応し、当該データに含まれるノイズの多寡に基づく第2出力とを出力する前記モデルを学習する学習手順と、

をコンピュータに実行させ、

前記取得手順は、

前記入力用データが検知された場所の分類結果を示す前記分類ラベルを含む前記学習用データを取得し、

前記学習手順は、

前記第1出力と、入力されたデータが検知された場所の分類結果を示す前記第2出力とを出力する前記モデルを学習する

ことを特徴とする情報処理プログラム。

#### 【請求項4】

機械学習のモデルの学習に用いる音声データである入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記入力用データに含まれるノイズの多寡に基づく分類結果を示す分類ラベルとを含む学習用データであって、前記入力用データが検知された場所の分類結果を示す前記分類ラベルを含む前記学習用データを用いて生成されたモデルであって、データの入力に応じて、前記正解データに対応する第1出力と前記分類ラベルに対応し、当該データに含まれるノイズの多寡に基づく第2出力とを出力するモデルを受信する受信部と、

前記受信部により受信された前記モデルにデータを入力することにより、当該データに対応する前記第1出力と、入力されたデータが検知された場所の分類結果を示す前記第2出力とを生成する推論処理を行う推論部と、

を備えたことを特徴とする端末装置。

#### 【請求項5】

機械学習のモデルの学習に用いる音声データである入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記入力用データに含まれるノイズの多寡に基づく分類結果を示す分類ラベルとを含む学習用データであって、前記入力用データが検知された場所の分類結果を示す前記分類ラベルを含む前記学習用データを用いて生成されたモデルであって、データの入力に応じて、前記正解データに対応する第1出力と前記分類ラベルに対応し、当該データに含まれるノイズの多寡に基づく第2出力とを出力するモデルを受信する受信工程と、

前記受信工程により受信された前記モデルにデータを入力することにより、当該データに対応する前記第1出力と、入力されたデータが検知された場所の分類結果を示す前記第2出力とを生成する推論処理を行う推論工程と、

10

20

30

40

50

を含んだことを特徴とする推論方法。

【請求項 6】

機械学習のモデルの学習に用いる音声データである入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記入力用データに含まれるノイズの多寡に基づく分類結果を示す分類ラベルとを含む学習用データであって、前記入力用データが検知された場所の分類結果を示す前記分類ラベルを含む前記学習用データを用いて生成されたモデルであって、データの入力に応じて、前記正解データに対応する第 1 出力と前記分類ラベルに対応し、当該データに含まれるノイズの多寡に基づく第 2 出力とを出力するモデルを受信する受信手順と、

前記受信手順により受信された前記モデルにデータを入力することにより、当該データに対応する前記第 1 出力と、入力されたデータが検知された場所の分類結果を示す前記第 2 出力とを生成する推論処理を行う推論手順と、

を端末装置に実行させることを特徴とする推論プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置、情報処理方法、情報処理プログラム、端末装置、推論方法、及び推論プログラムに関する。

【背景技術】

【0002】

機械学習の技術により、様々な用途に用いられる学習モデル（以下「モデル」ともいう）を学習する技術が提供されている。例えば、モデルを音声認識に用いる音声認識システムが提供されている。

【先行技術文献】

【特許文献】

【0003】

【文献】特開 2019 - 159058 号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、上記の従来技術では、多様な出力を行うモデルを利用可能にすることが難しい。例えば、従来技術では音声認識の結果としてその文字データ（テキスト）の 1 つの種別の出力のみを行っているに過ぎない。そのため、複数種別の出力を行うモデルを利用可能にすることができるとは限らない。

【0005】

本願は、上記に鑑みてなされたものであって、複数種別の出力を行うモデルを利用可能にする情報処理装置、情報処理方法、情報処理プログラム、端末装置、推論方法、及び推論プログラムを提供することを目的とする。

【課題を解決するための手段】

【0006】

本願に係る情報処理装置は、機械学習のモデルの学習に用いる入力用データと、当該入力用データに含まれる認識対象を示す正解データと、前記認識対象に関連する分類結果を示す分類ラベルとを含む学習用データを取得する取得部と、前記学習用データを用いて、データの入力に応じて、前記正解データに対応する第 1 出力と前記分類ラベルに対応する第 2 出力とを出力する前記モデルを学習する学習部と、を備えたことを特徴とする。

【発明の効果】

【0007】

実施形態の一態様によれば、複数種別の出力を行うモデルを利用可能にすることができるという効果を奏する。

【図面の簡単な説明】

10

20

30

40

50

## 【 0 0 0 8 】

【図 1】図 1 は、実施形態に係る情報処理システムによる処理の一例を示す図である。

【図 2】図 2 は、実施形態に係る情報処理装置の構成例を示す図である。

【図 3】図 3 は、実施形態に係る学習用データ記憶部の一例を示す図である。

【図 4】図 4 は、実施形態に係るモデル情報記憶部の一例を示す図である。

【図 5】図 5 は、実施形態に係る端末装置の構成例を示す図である。

【図 6】図 6 は、実施形態に係る情報処理装置による処理の一例を示すフローチャートである。

【図 7】図 7 は、実施形態に係る端末装置による処理の一例を示すフローチャートである。

【図 8】図 8 は、ハードウェア構成の一例を示す図である。

10

【発明を実施するための形態】

## 【 0 0 0 9 】

以下に、本願に係る情報処理装置、情報処理方法、情報処理プログラム、端末装置、推論方法、及び推論プログラムを実施するための形態（以下、「実施形態」と呼ぶ）について図面を参照しつつ詳細に説明する。なお、この実施形態により本願に係る情報処理装置、情報処理方法、情報処理プログラム、端末装置、推論方法、及び推論プログラムが限定されるものではない。また、以下の各実施形態において同一の部位には同一の符号を付し、重複する説明は省略される。

## 【 0 0 1 0 】

〔 1 . はじめに 〕

20

近年、音声認識等に利用されるモデル（「音声認識モデル」ともいう）に、End-to-Endモデル（「E2Eモデル」ともいう）が用いられている。E2Eモデルは、例えば1つのニューラルネットワークで構成されるモデルである。E2Eモデルは、ユーザが利用するデバイス（端末装置10等）で完結する音声認識の処理に適している。音声データを入力してその音声データに対応する文字データを出力させる音声文字変換等の音声認識モデルにおいて、入力された音声データに対応する認識結果が出力される。このような音声認識モデルでは、従来は入力データに対応する文字データ（テキスト）等1つの種類の出力を行う。

## 【 0 0 1 1 】

一方で、以下に示す情報処理装置100が学習するモデルは、入力用データに含まれる認識対象の認識結果を示す出力（「第1出力」ともいう）と、認識対象に関連する分類結果を示す出力（「第2出力」ともいう）との複数の種類の出力を行う。これにより、情報処理装置100が学習するモデルは、認識対象に関連する分類結果に関連する特徴を加味して学習される。そのため、情報処理装置100が学習するモデルは、認識対象に関連する分類を加味しつつ、入力用データの特徴を抽出し、第1出力を出力するため、第1出力に関する認識精度を向上させることができる。

30

## 【 0 0 1 2 】

また、入力に音声データとその音声データが検知された場所を示すラベル等の複数種類の情報の入力を用いる場合、推論時にも複数種類の情報を入力する必要となる。そのため、モデルの利用する場面において利便性が低く、その情報を入力として用意できない場合、モデルを利用できなかつたり、推論の精度が低下したりする。一方で、情報処理装置100が学習するモデルは、入力用データとして、以下に示すように例えば音声データ等の1つのデータのみでよい場合、推論時にも複数種類の情報を入力する必要がない。そのため、情報処理装置100は、認識精度を向上させつつ、利便性の高いモデルを学習することができる。

40

## 【 0 0 1 3 】

（実施形態）

〔 2 . 情報処理 〕

ここから、図1を用いて、実施形態に係る情報処理の一例について説明する。図1は、実施形態に係る情報処理システムによる処理の一例を示す図である。まず、情報処理シス

50

テム 1 の構成について説明する。

【 0 0 1 4 】

図 1 に示すように、情報処理システム 1 は、端末装置 1 0 と、情報処理装置 1 0 0 とが含まれる。端末装置 1 0 と、情報処理装置 1 0 0 とは図示しない所定の通信網を介して、有線または無線により通信可能に接続される。なお、図 1 に示した情報処理システム 1 には、複数台の端末装置 1 0 や、複数台の情報処理装置 1 0 0 が含まれてもよい。

【 0 0 1 5 】

情報処理装置 1 0 0 は、機械学習のモデルの学習に用いる入力用データと、入力用データに含まれる認識対象を示す正解データと、認識対象に関連する分類結果を示す分類ラベルとを含む学習用データを用いて、正解データに対応する第 1 出力と分類ラベルに対応する第 2 出力とを出力するモデルを学習する情報処理装置である。情報処理装置 1 0 0 は、音声文字変換結果である第 1 出力と、分類結果を示す第 2 出力との 2 つの種別の出力を行うモデル M 1 を学習し、端末装置 1 0 に提供する。なお、モデル M 1 のネットワーク構成は、第 1 出力及び第 2 出力の出力が可能であればどのようなネットワーク構成であってもよく、E 2 E モデルであってもよい。

10

【 0 0 1 6 】

端末装置 1 0 は、ユーザによって利用されるデバイス（コンピュータ）である。端末装置 1 0 は、ユーザによる音声入力を受け付ける。端末装置 1 0 は、ユーザによる操作を受け付ける。端末装置 1 0 は、情報処理装置 1 0 0 から提供されたモデルを用いて推論を行う。

20

【 0 0 1 7 】

また、以下では、端末装置 1 0 をユーザと表記する場合がある。すなわち、以下では、ユーザを端末装置 1 0 と読み替えることもできる。なお、端末装置 1 0 は、例えば、スマートフォンや、タブレット型端末や、ノート型 P C ( Personal Computer ) や、デスクトップ P C や、携帯電話機や、P D A ( Personal Digital Assistant ) 等により実現される。図 1 の例では、端末装置 1 0 がタッチパネル機能を有するスマートフォンである場合を示す。

【 0 0 1 8 】

以下、図 1 を用いて、情報処理の一例を説明する。図 1 では、ユーザがユーザ I D 「 U 1 」により識別されるユーザ（以下、「ユーザ U 1 」とする場合がある）である場合を示す。また、図 1 では、情報処理装置 1 0 0 が音声データの入力に対して、その音声データが変換された文字データである第 1 出力、及びその音声データが検知（収集）された場所を分類する分類ラベルである第 2 出力を出力するモデル M 1 を学習する場合を一例として説明する。なお、分類ラベルは場所の分類に限らず、様々な対象の分類を示すものであってもよいが、この点については後述する。

30

【 0 0 1 9 】

まず、情報処理装置 1 0 0 は、機械学習に用いる学習用データ群 D S 1 を取得する（ステップ S 1 1）。学習用データ群 D S 1 には、モデルの入力として用いられる入力用データと、その入力用データに対応する正解データ及び分類ラベルとのセット（組合せ）が複数含まれる。例えば、入力用データである音声データ D T 1 は、正解データ R T 1 と分類ラベル C L 1 とが対応付けられている。この場合、正解データ R T 1 は、音声データ D T 1 が文字変換された文字データ（文字列）であり、分類ラベル C L 1 は、音声データ D T 1 が検知（収集）された場所の分類が「カフェ」であることを示す値（例えば 1 等）であるものとする。

40

【 0 0 2 0 】

そして、情報処理装置 1 0 0 は、学習用データ群 D S 1 を用いて、文字データである第 1 出力と、音声データの検知場所を示す分類ラベルである第 2 出力との 2 つの種別の出力を行うモデル M 1 を学習する（ステップ S 1 2）。図 1 では、音声データである入力用データ I N の入力層への入力に応じて、出力層から文字データである第 1 出力 O T 1 と、音声データの検知場所を示す分類ラベルである第 2 出力 O T 2 を出力するモデル M 1 を概念

50

的に示す。

【0021】

情報処理装置100は、音声データDT1が入力された場合に、第1出力OT1として正解データRT1が出力され、第2出力OT2として分類ラベルCL1が出力されるようにモデルM1の重み等のパラメータを学習する。また、音声データDT2が入力された場合に、第1出力OT1として正解データRT2が出力され、第2出力OT2として分類ラベルCL2が出力されるようにモデルM1の重み等のパラメータを学習する。これにより、情報処理装置100は、第1出力と第2出力との2つの種別の出力を行うモデルM1を生成する。モデルM1の学習処理には、任意の手法が採用可能である。

【0022】

例えば、情報処理装置100は、バックプロパゲーション（誤差逆伝播法）等の手法により学習処理を行う。例えば、情報処理装置100は、学習処理により、ノード間で値が伝達する際に考慮される重み（すなわち、接続係数）の値を調整する。このように、情報処理装置100は、モデルM1における出力（第1出力及び第2出力）と、入力に対応する正解（正解データ及び分類ラベル）との誤差が少なくなるようにパラメータ（接続係数）を補正するバックプロパゲーション等の処理によりモデルM1を学習する。例えば、情報処理装置100は、所定の損失（ロス）関数を最小化するようにバックプロパゲーション等の処理を行うことによりモデルM1を生成する。これにより、情報処理装置100は、モデルM1のパラメータを学習する学習処理を行うことができる。

【0023】

そして、情報処理装置100は、学習したモデルM1をユーザU1が利用する端末装置10に提供する（ステップS13）。モデルM1を受信した端末装置10は、モデルM1を利用してユーザU1の発話も文字に書き起こす音声文字変換処理（推論）を実行する。この点について以下説明する。

【0024】

まず、ユーザU1が「XXXX」と発話する。なお、「XXXX」は具体的な内容を含む発話であるものとする。端末装置10は、ユーザU1の発話PAを検知し、ユーザU1の発話PAである「XXXX」の音声データを入力として受け付ける（ステップS14）。

【0025】

そして、端末装置10は、入力として受け付けた「XXXX」の音声データと、モデルM1とを利用して推論処理を行う（ステップS15）。端末装置10は、「XXXX」の音声データをモデルM1に入力し、モデルM1に文字データ及び分類ラベルを出力させることにより、音声を文字に変換するとともに、ユーザU1が発話PAを行った場所の分類を推定する処理（推論処理）を行う。図1では、「XXXX」の音声データが入力されたモデルM1は、「XXXX」の文字データ（第1出力）と、発話PAを行った場所の分類結果がカフェであることを示す分類ラベル（第2出力）を出力する。

【0026】

そして、端末装置10は、推論結果である「XXXX」の文字データを表示してもよい（ステップS16）。例えば、端末装置10は、文字列「XXXX」を画面に表示する。なお、端末装置10は、分類ラベルを表示してもよい。

【0027】

また、端末装置10は、音声データ、その音声データに対応する文字データである正解データ、及びその音声データに対応する分類ラベルを学習用データとして情報処理装置100に送信してもよい（ステップS17）。この場合、情報処理装置100は、端末装置10から受信した学習用データを用いて、モデルM1のパラメータを更新してもよい。

【0028】

上述したように、情報処理装置100は、音声データを入力として、その音声データが変換された文字データと、音声データの検知場所を示す分類ラベルとの2つの種別の出力を行うモデルM1を適切に学習することができる。図1の例では、情報処理装置100は、分類ラベルを出力するE2Eの音声認識モデルであるモデルM1を適切に学習すること

10

20

30

40

50

ができる。したがって、情報処理装置 100 は、複数種別の出力を行うモデルを利用可能にすることができる。また、情報処理装置 100 は、分類ラベルを出力することで音声文字変換の精度が改善することができる。例えば、情報処理装置 100 は、認識対象に関連する分類を加味しつつ、入力用データの特徴を抽出し、音声文字変換の結果を出力するモデル M1 を学習するため、分類ラベルの出力により音声文字変換の精度を向上させたモデル M1 を学習することができる。また、端末装置 10 は、音声データが変換された文字データと、音声データの検知場所を示す分類ラベルとの 2 つの種別の出力するモデルを用いて推論処理を行うことができるため、適切なモデルを利用した処理を行うことができる。したがって、端末装置 10 は、複数種別の出力を行うモデルを利用した処理を行うことができる。

10

## 【0029】

## 〔2-1. 分類ラベル〕

なお、上記の例では、分類ラベルが音声データが検知（収集）された場所の分類を示す場合を示したが、分類ラベルは、認識対象に関連する分類結果を示すものであればどのような対象の分類であってもよい。分類ラベルは、入力用データに含まれる情報のうち、認識対象以外の情報の種別を示すものである。つまり、分類ラベルは、入力用データのうち認識対象となるデータ以外のデータの分類結果を示すものであってもよい。また、分類ラベルは、入力用データのうち、認識対象となるデータから認識される認識結果以外の各種情報であってもよい。また、分類ラベルは、入力用データが取得された際の各種コンテキストを示すものであってもよい。このように、分類ラベルは、認識対象となるデータから認識結果を認識する認識処理において影響を及ぼしうる任意の要素であって、入力用データから取得もしくは推定可能な要素、もしくは入力用データに付随する各種の要素の分類結果が採用可能である。この点について以下例示を列挙する。なお、図 1 と同様の点については適宜説明を省略する。また、以下に示す各モデルのネットワーク構成はモデル M1 と同様であってもよい。

20

## 【0030】

## 〔2-1-1. ユーザ関連〕

例えば、モデルが出力する分類ラベル（第 2 出力）は、入力用データが検知された場所に限らず、入力用データに関連する様々なコンテキストの分類結果であってもよい。例えば、モデルが出力する分類ラベル（第 2 出力）は、入力用データに含まれる発話を行ったユーザに関連するコンテキストの分類結果であってもよい。

30

## 【0031】

## 〔2-1-2. 周囲〕

モデルが出力する分類ラベル（第 2 出力）は、発話を行ったユーザの発話以外の分類結果を示してもよい。例えば、モデルが出力する分類ラベル（第 2 出力）は、入力用データに含まれる発話を行ったユーザの周囲の状況の分類結果であってもよい。この場合、情報処理装置 100 は、ユーザの周囲の状況の分類結果を示す分類ラベルを含む学習用データを用いて、ユーザの周囲の状況の分類結果を示す第 2 出力と第 1 出力とを出力するモデル（「モデル M2」とする）を学習する。

## 【0032】

情報処理装置 100 は、ユーザの周囲が騒がしいか否かを示す分類ラベル、すなわちユーザの発話以外のノイズが多いか否かを示す分類ラベルを第 2 出力として出力するモデル M2 を学習する。この場合、分類ラベルは、ユーザの周囲が騒がしい（ノイズが多い）程、大きい値となってもよい。

40

## 【0033】

情報処理装置 100 は、モデル M2 の入力として用いられる入力用データと、その入力用データに対応する正解データ及び分類ラベルとのセット（組合せ）を複数含む学習用データ（「学習用データ DS2」とする）を用いて、モデル M2 を学習する。例えば、学習用データ DS2 の正解データは、入力用データである音声データが変換された文字データであり、学習用データ DS2 の分類ラベルは、入力用データである音声データにおいて発

50

話したユーザの周囲の騒がしさの度合いを示す値である。

【 0 0 3 4 】

情報処理装置 1 0 0 は、学習用データ D S 2 を用いて、各入力用データが入力された場合に、その入力用データに対応する正解データ及び分類ラベルが出力されるようにモデル M 2 のパラメータを学習する。例えば、情報処理装置 1 0 0 は、音声データが入力された場合に、音声データが変換された文字データが第 1 出力として出力され、音声データに含まれるノイズの多寡を示す分類ラベルを第 2 出力として出力されるようにモデル M 2 を学習する。

【 0 0 3 5 】

〔 2 - 1 - 3 . ユーザ属性 〕

モデルが出力する分類ラベル（第 2 出力）は、入力用データに含まれる発話を行ったユーザの属性の分類結果であってもよい。この場合、情報処理装置 1 0 0 は、ユーザの属性の分類結果を示す分類ラベルを含む学習用データを用いて、ユーザの属性の分類結果を示す第 2 出力と第 1 出力とを出力するモデル（「モデル M 3」とする）を学習する。以下では、ユーザの年齢をユーザの属性の一例として説明するが、ユーザの属性は、年齢に限らず、性別、身長、出身地等の様々な属性（要素）であってもよい。

【 0 0 3 6 】

情報処理装置 1 0 0 は、ユーザの年齢を示す分類ラベル（例えば 1 0 代、2 0 代等の段階的な分類）を第 2 出力として出力するモデル M 3 を学習する。情報処理装置 1 0 0 は、モデル M 3 の入力として用いられる入力用データと、その入力用データに対応する正解データ及び分類ラベルとのセット（組合せ）を複数含む学習用データ（「学習用データ D S 3」とする）を用いて、モデル M 3 を学習する。例えば、学習用データ D S 3 の正解データは、入力用データである音声データが変換された文字データであり、学習用データ D S 3 の分類ラベルは、入力用データである音声データにおいて発話したユーザの年齢を示す値である。

【 0 0 3 7 】

情報処理装置 1 0 0 は、学習用データ D S 3 を用いて、各入力用データが入力された場合に、その入力用データに対応する正解データ及び分類ラベルが出力されるようにモデル M 3 のパラメータを学習する。例えば、情報処理装置 1 0 0 は、音声データが入力された場合に、音声データが変換された文字データが第 1 出力として出力され、音声データに含まれる発話を行ったユーザの年齢を示す分類ラベルを第 2 出力として出力されるようにモデル M 3 を学習する。

【 0 0 3 8 】

〔 2 - 1 - 4 . 端末装置 〕

モデルが出力する分類ラベル（第 2 出力）は、ユーザが利用する端末装置 1 0 の分類結果であってもよい。例えば、モデルが出力する分類ラベル（第 2 出力）は、ユーザの発話（音声データ）を検知（収集）した端末装置 1 0 の機種のカテゴリの分類結果であってもよい。この場合、情報処理装置 1 0 0 は、端末装置 1 0 の機種のカテゴリの分類結果を示す分類ラベルを含む学習用データを用いて、端末装置 1 0 の機種のカテゴリの分類結果を示す第 2 出力と第 1 出力とを出力するモデル（「モデル M 4」とする）を学習する。

【 0 0 3 9 】

情報処理装置 1 0 0 は、端末装置 1 0 の機種を示す分類ラベル（例えば製品 A、製品 B 等の製品の分類）を第 2 出力として出力するモデル M 4 を学習する。情報処理装置 1 0 0 は、モデル M 4 の入力として用いられる入力用データと、その入力用データに対応する正解データ及び分類ラベルとのセット（組合せ）を複数含む学習用データ（「学習用データ D S 4」とする）を用いて、モデル M 4 を学習する。例えば、学習用データ D S 4 の正解データは、入力用データである音声データが変換された文字データであり、学習用データ D S 4 の分類ラベルは、入力用データである音声データを検知した端末装置 1 0 の機種を示す値である。

【 0 0 4 0 】

10

20

30

40

50

情報処理装置 100 は、学習用データ DS4 を用いて、各入力用データが入力された場合に、その入力用データに対応する正解データ及び分類ラベルが出力されるようにモデル M4 のパラメータを学習する。例えば、情報処理装置 100 は、音声データが入力された場合に、音声データが変換された文字データが第 1 出力として出力され、音声データを検知した端末装置 10 の機種を示す分類ラベルを第 2 出力として出力されるようにモデル M4 を学習する。

#### 【0041】

なお、上述は一例に過ぎず、分類できるものであれば、どのような対象の分類ラベルを用いてもよい。

#### 【0042】

また、情報処理装置 100 は、複数の第 2 出力を出力するようにモデルを学習してもよい。すなわち、情報処理装置 100 は、1 つの第 1 出力と、2 つ以上の第 2 出力を出力するモデルを生成してもよい。この場合、情報処理装置 100 は、上述した各種の分類ラベルから選択された 2 つ以上の分類ラベルを出力とするモデルを学習する。例えば、情報処理装置 100 は、音声データが入力された場合に、音声データが変換された文字データと、音声データが検知（収集）された場所を示す第 1 分類ラベルと、音声データに含まれるユーザの属性を示す第 2 分類ラベルと出力するモデルを学習する。

#### 【0043】

##### 〔2-2. 推論対象〕

なお、学習するモデルの用途は、音声文字変換に限らず、他の音声認識に関する様々な用途であってもよい。また、モデルの入力は、音声データに限らず、画像データ等様々な種別のデータが対象であってもよい。例えば、モデルの入力が画像データである場合、学習されるモデルの用途は、一般物体認識等の各種の画像認識に関する用途であってもよい。この場合、分類ラベルは、画像に含まれる物体以外の分類を示すものであってもよい。例えば、分類ラベルは、画像に含まれる人（ユーザ）の年齢等のユーザの属性であってもよく、画像が示すシーンの状況（昼、夜、室内、屋外等）などのコンテキストであってもよい。

#### 【0044】

##### 〔3. 情報処理装置の構成〕

次に、図 2 を用いて、実施形態に係る情報処理装置 100 の構成について説明する。図 2 は、実施形態に係る情報処理装置 100 の構成例を示す図である。図 2 に示すように、情報処理装置 100 は、通信部 110 と、記憶部 120 と、制御部 130 とを有する。なお、情報処理装置 100 は、情報処理装置 100 の管理者等から各種操作を受け付ける入力部（例えば、キーボードやマウス等）や、各種情報を表示するための表示部（例えば、液晶ディスプレイ等）を有してもよい。

#### 【0045】

##### （通信部 110）

通信部 110 は、例えば、NIC（Network Interface Card）等によって実現される。そして、通信部 110 は、所定の通信網（ネットワーク）と有線または無線で接続され、端末装置 10 との間で情報の送受信を行う。

#### 【0046】

##### （記憶部 120）

記憶部 120 は、例えば、RAM（Random Access Memory）、フラッシュメモリ（Flash Memory）等の半導体メモリ素子、または、ハードディスク、光ディスク等の記憶装置によって実現される。実施形態に係る記憶部 120 は、図 2 に示すように、学習用データ記憶部 121 と、モデル情報記憶部 122 とを有する。

#### 【0047】

##### （学習用データ記憶部 121）

実施形態に係る学習用データ記憶部 121 は、学習に用いるデータに関する各種情報を記憶する。学習用データ記憶部 121 は、学習に用いる学習データ（データセット）を記

10

20

30

40

50

憶する。図3は、本開示の実施形態に係る学習用データ記憶部の一例を示す図である。例えば、学習用データ記憶部121は、学習に用いる学習データや精度評価（測定）に用いる評価用データ等の種々のデータに関する各種情報を記憶する。図3に、実施形態に係る学習用データ記憶部121の一例を示す。図3の例では、学習用データ記憶部121は、「データセットID」、「データID」、「データ」、「正解データ」、「分類ラベル」といった項目が含まれる。

【0048】

「データセットID」は、データセットを識別するための識別情報を示す。「データID」は、各学習用データを識別するための識別情報を示す。また、「データ」は、データIDにより識別されるデータを示す。「データ」は、モデルの入力として用いられるデータ（入力用データ）を示す。図3の例では、入力用データは、種別が「音声」である音声データの場合を示す。

10

【0049】

「正解データ」は、対応するデータ（入力用データ）に対応する正解を示す。図3の例では、「正解データ」は、入力用データである音声データが変換された文字データ（文字列）を示す。「正解データ」は、対応するデータ（入力用データ）がモデルに入力された場合に、モデルが出力することが期待される第1出力（文字データ）を示す。

【0050】

「分類ラベル」は、対応するデータ（入力用データ）に対応する分類結果を示す。図3の例では、「分類ラベル」は、入力用データである音声データが検知された際のコンテキストを推定するための分類を示す。「分類ラベル」は、対応するデータ（入力用データ）がモデルに入力された場合に、モデルが出力することが期待される第2出力（分類ラベル）を示す。

20

【0051】

例えば、「分類ラベル」は、入力用データである音声データが検知された場所を推定するための分類を示す。例えば、分類ラベルが「1」の場合は「カフェ」であることを示し、分類ラベルが「2」の場合は「自宅」であることを示してもよい。記憶部120は、分類ラベルと各コンテキストの対応付けを示す情報を記憶してもよい。なお、「分類ラベル」は、1つに限らず、モデルが出力する分類ラベルの数に応じた数であってもよい。例えば、場所と発話ユーザの属性を推定する場合、場所を示すラベルを登録する「分類ラベル#1」と、ユーザの属性を示すラベルを登録する「分類ラベル#2」の複数の項目が含まれてもよい。すなわち、「分類ラベル」は、2つ以上のコンテキストの各々に対応する分類ラベルが記憶されてもよい。なおユーザの属性は、年齢や性別等のデモグラフィック属性やサイコグラフィック属性の様々なユーザの属性を示す情報であってもよい。

30

【0052】

図3の例では、データセットID「DS1」により識別されるデータセット（データセットDS1）には、データID「DID1」、「DID2」、「DID3」等により識別される複数のデータが含まれることを示す。

【0053】

データID「DID1」により識別されるデータDT1は、正解データが「RT1」であることを示す。図3の例では「RT1」のように抽象的に図示するが、「正解データ」には、音声データ（入力用データ）に含まれる認識対象（ユーザの発話）が変換された文字データ（文字列）であるものとする。

40

【0054】

データDT1は、分類ラベルが「CL1」であることを示す。図3の例では「CL1」のように抽象的に図示するが、「分類ラベル」には、認識対象に関連する分類結果を示す分類ラベル（値）であるものとする。

【0055】

なお、学習用データ記憶部121は、上記に限らず、目的に応じて種々の情報を記憶してもよい。例えば、学習用データ記憶部121は、音声や画像等の教師データの種別を示

50

す情報を各データに対応付けて記憶する。例えば、学習用データ記憶部 1 2 1 は、データの種別を示す情報を各データに対応付けて記憶する。

【 0 0 5 6 】

例えば、学習用データ記憶部 1 2 1 は、各データが学習データであるか、評価用データであるか等を特定可能に記憶してもよい。例えば、学習用データ記憶部 1 2 1 は、学習データと評価用データとを区別可能に記憶する。学習用データ記憶部 1 2 1 は、各データが学習データや評価用データであることを識別する情報を記憶してもよい。情報処理装置 1 0 0 は、学習データとして用いられる各データと正解データと分類ラベルとに基づいて、モデルを学習する。情報処理装置 1 0 0 は、評価用データとして用いられる各データと正解データと分類ラベルとに基づいて、モデルの精度を測定する。情報処理装置 1 0 0 は、評価用データを入力した場合にモデルが出力する出力結果（第 1 出力、第 2 出力）と、正解データ及び分類ラベルとを比較した結果を収集することにより、モデルの精度を測定する。

10

【 0 0 5 7 】

（モデル情報記憶部 1 2 2 ）

実施形態に係るモデル情報記憶部 1 2 2 は、モデルに関する情報を記憶する。例えば、モデル情報記憶部 1 2 2 は、学習処理により学習（生成）された学習済みモデル（モデル）の情報（モデルデータ）を記憶する。図 4 は、本開示の第 1 の実施形態に係るモデル情報記憶部の一例を示す図である。図 4 に、第 1 の実施形態に係るモデル情報記憶部 1 2 2 の一例を示す。図 4 に示した例では、モデル情報記憶部 1 2 2 は、「モデル ID」、「用途」、「モデルデータ」といった項目が含まれる。

20

【 0 0 5 8 】

「モデル ID」は、モデルを識別するための識別情報を示す。「用途」は、対応するモデルの用途を示す。「モデルデータ」は、モデルのデータを示す。図 4 等では「モデルデータ」に「M D T 1」といった概念的な情報が格納される例を示したが、実際には、モデルの構成（ネットワーク構成）の情報やパラメータに関する情報等、そのモデルを構成する種々の情報が含まれる。例えば、「モデルデータ」には、ネットワークの各層におけるノードと、各ノードが採用する関数と、ノードの接続関係と、ノード間の接続に対して設定される接続係数とを含む情報が含まれる。

【 0 0 5 9 】

図 4 に示す例では、モデル ID「M 1」により識別されるモデル（モデル M 1）は、用途が「音声文字変換」、「コンテキスト推定」であることを示す。すなわち、モデル M 1 は、入力用データを文字起こしした文字データと、その入力用データに関連するコンテキストを推定する情報とを出力するモデルであることを示す。また、モデル M 1 のモデルデータは、モデルデータ M D T 1 であることを示す。

30

【 0 0 6 0 】

なお、モデル情報記憶部 1 2 2 は、上記に限らず、目的に応じて種々の情報を記憶してもよい。

【 0 0 6 1 】

（制御部 1 3 0 ）

図 2 の説明に戻って、制御部 1 3 0 は、コントローラ（controller）であり、例えば、CPU（Central Processing Unit）や MPU（Micro Processing Unit）等によって、情報処理装置 1 0 0 内部の記憶装置に記憶されている各種プログラム（情報処理プログラムの一例に相当）が RAM を作業領域として実行されることにより実現される。また、制御部 1 3 0 は、コントローラであり、例えば、ASIC（Application Specific Integrated Circuit）や FPGA（Field Programmable Gate Array）等の集積回路により実現される。

40

【 0 0 6 2 】

図 2 に示すように、制御部 1 3 0 は、取得部 1 3 1 と、決定部 1 3 2 と、学習部 1 3 3 と、提供部 1 3 4 とを有し、以下に説明する情報処理の機能や作用を実現または実行する。なお、制御部 1 3 0 の内部構成は、図 2 に示した構成に限られず、後述する情報処理を

50

行う構成であれば他の構成であってもよい。また、制御部 130 が有する各処理部の接続関係は、図 2 に示した接続関係に限られず、他の接続関係であってもよい。

【0063】

(取得部 131)

取得部 131 は、記憶部 120 から各種の情報を取得する。取得部 131 は、学習用データ記憶部 121 から学習に用いるデータを取得する。取得部 131 は、モデル情報記憶部 122 からモデルの情報を取得する。

【0064】

取得部 131 は、機械学習のモデルの学習に用いる入力用データと、当該入力用データに含まれる認識対象を示す正解データと、認識対象に関連する分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、音声データである入力用データと、当該入力用データに含まれる音声の音声認識の結果を示す正解データを含む学習用データを取得する。取得部 131 は、入力用データに対応する文字データである正解データを含む学習用データを取得する。取得部 131 は、入力用データに関連するコンテキストの分類結果を示す分類ラベルを含む学習用データを取得する。

10

【0065】

取得部 131 は、入力用データが検知された場所の分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、入力用データに含まれる発話を行ったユーザに関連するコンテキストの分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、ユーザの周囲の状況の分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、ユーザの発話以外の分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、ユーザの属性の分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、ユーザが利用する端末装置の分類結果を示す分類ラベルを含む学習用データを取得する。取得部 131 は、入力用データに含まれる情報のうち、認識対象以外の情報の種別を示す分類ラベルを取得する。

20

【0066】

取得部 131 は、通信部 110 を介して、端末装置 10 から情報を受信する。取得部 131 は、端末装置 10 から学習用データを取得する。取得部 131 は、端末装置 10 において、音声文字変換処理の対象となった音声データと、その音声データに対応する第 1 出力及び第 2 出力、またユーザが修正した修正結果とのセット(組合せ)を学習用データとして端末装置 10 から収集する。

30

【0067】

(決定部 132)

決定部 132 は、種々の情報を決定する。例えば、決定部 132 は、分類ラベルを決定する。決定部 132 は、認識対象に関連する分類結果を決定することにより、モデルに学習させる分類ラベルを決定する。決定部 132 は、入力用データに関連するコンテキストをモデルに分類させる対象に決定する。決定部 132 は、音声データが検知された場所をモデルに分類させる対象に決定する。

【0068】

(学習部 133)

学習部 133 は、モデルを学習する。学習部 133 は、外部の情報処理装置からの情報や記憶部 120 に記憶された情報に基づいて、各種情報を学習する。学習部 133 は、学習用データ記憶部 121 に記憶された情報に基づいて、各種情報を学習する。学習部 133 は、学習により生成したモデルをモデル情報記憶部 122 に格納する。

40

【0069】

学習部 133 は、学習用データを用いて、データの入力に応じて、正解データに対応する第 1 出力と分類ラベルに対応する第 2 出力とを出力するモデルを学習する。学習部 133 は、入力された音声データに対する音声認識の結果を示す第 1 出力と第 2 出力とを出力するモデルを学習する。学習部 133 は、入力された音声データが変換された文字データである第 1 出力と第 2 出力とを出力するモデルを学習する。

50

## 【 0 0 7 0 】

学習部 1 3 3 は、第 1 出力と、入力されたデータに関連するコンテキストの分類結果を示す第 2 出力とを出力するモデルを学習する。学習部 1 3 3 は、第 1 出力と、入力されたデータが検知された場所の分類結果を示す第 2 出力とを出力するモデルを学習する。学習部 1 3 3 は、第 1 出力と、入力されたデータに含まれる発話を行ったユーザに関連するコンテキストの分類結果を示す第 2 出力とを出力するモデルを学習する。学習部 1 3 3 は、第 1 出力と、ユーザの周囲の状況の分類結果を示す第 2 出力とを出力するモデルを学習する。学習部 1 3 3 は、第 1 出力と、ユーザの発話以外の分類結果を示す第 2 出力とを出力するモデルを学習する。学習部 1 3 3 は、第 1 出力と、ユーザの属性の分類結果を示す第 2 出力とを出力するモデルを学習する。学習部 1 3 3 は、第 1 出力と、ユーザが利用する端末装置の分類結果を示す第 2 出力とを出力するモデルを学習する。

10

## 【 0 0 7 1 】

学習部 1 3 3 は、モデル（ネットワーク）のパラメータを学習する。学習部 1 3 3 は、接続されたノード間の接続係数（重み）等のパラメータを学習する。学習部 1 3 3 は、種々の機械学習に関する技術を用いて、モデルを学習する。学習部 1 3 3 は、モデルに入力するデータと、そのデータが入力された場合の出力を示す正解データ及び分類ラベルとを用いて行う学習処理、すなわち教師有り学習の手法によりモデルのパラメータを学習する。なお、上記は一例であり、学習部 1 3 3 は、モデルのパラメータを学習可能であれば、どのような学習処理により、モデルのパラメータを学習してもよい。

20

## 【 0 0 7 2 】

## （提供部 1 3 4）

提供部 1 3 4 は、通信部 1 1 0 を介して、端末装置 1 0 へ情報を送信する。提供部 1 3 4 は、端末装置 1 0 へモデルを提供する。例えば、提供部 1 3 4 は、端末装置 1 0 へ音声文字変換に用いるモデル M 1 を送信する。

## 【 0 0 7 3 】

## 〔 4 . 端末装置の構成 〕

次に、図 5 を用いて、実施形態に係る端末装置 1 0 の構成について説明する。図 5 は、実施形態に係る端末装置 1 0 の構成例を示す図である。図 5 に示すように、端末装置 1 0 は、通信部 1 1 と、記憶部 1 2 と、入力部 1 3 と、表示部 1 4 と、制御部 1 5 とを有する。なお、端末装置 1 0 は、各種情報を音声出力するための音声出力部（例えばスピーカ等）を有してもよい。

30

## 【 0 0 7 4 】

## （通信部 1 1）

通信部 1 1 は、例えば、通信回路等によって実現される。そして、通信部 1 1 は、図示しない所定の通信網と有線または無線で接続され、情報処理装置 1 0 0 との間で情報の送受信を行う。

## 【 0 0 7 5 】

## （記憶部 1 2）

記憶部 1 2 は、例えば、RAM、フラッシュメモリ等の半導体メモリ素子、または、ハードディスク、光ディスク等の記憶装置によって実現される。記憶部 1 2 は、例えば、端末装置 1 0 にインストールされているアプリケーション（例えば音声文字変換アプリ等）に関する情報、例えばプログラム等を記憶する。また、記憶部 1 2 は、情報処理装置 1 0 0 から提供されたモデルを記憶する。例えば、記憶部 1 2 は、モデル M 1 を記憶する。

40

## 【 0 0 7 6 】

## （入力部 1 3）

入力部 1 3 は、ユーザからの各種操作を受け付ける。入力部 1 3 は、音声を検知する機能を有し、ユーザの発話による音声入力を受け付ける。入力部 1 3 は、音声を検知するマイクにより検知されたユーザによる発話を入力として受け付ける。

## 【 0 0 7 7 】

また、入力部 1 3 は、タッチパネル機能により表示面を介してユーザからの各種操作を

50

受け付けてもよい。また、入力部 1 3 は、端末装置 1 0 に設けられたボタンや、端末装置 1 0 に接続されたキーボードやマウスからの各種操作を受け付けてもよい。

【 0 0 7 8 】

例えば、入力部 1 3 は、端末装置 1 0 の表示部 1 4 を介してユーザの指定操作等の操作を受け付ける。例えば、入力部 1 3 は、タッチパネルの機能によりユーザの操作を受け付ける受付部として機能する。この場合、入力部 1 3 と受付部 1 5 2 とは一体であってもよい。なお、入力部 1 3 によるユーザの操作の検知方式には、タブレット端末では主に静電容量方式が採用されるが、他の検知方式である抵抗膜方式、表面弾性波方式、赤外線方式、電磁誘導方式など、ユーザの操作を検知できタッチパネルの機能が実現できればどのような方式を採用してもよい。

10

【 0 0 7 9 】

( 表示部 1 4 )

表示部 1 4 は、例えば液晶ディスプレイや有機 E L ( Electro-Luminescence ) ディスプレイ等によって実現されるタブレット端末等の表示画面であり、各種情報を表示するための表示装置である。

【 0 0 8 0 】

( 制御部 1 5 )

制御部 1 5 は、コントローラであり、例えば、CPU や MPU 等によって、端末装置 1 0 内部の記憶部 1 2 などの記憶装置に記憶されている各種プログラムが RAM を作業領域として実行されることにより実現される。例えば、この各種プログラムは、インストールされているアプリケーション ( 例えばメッセージアプリ等 ) のプログラムが含まれる。また、制御部 1 5 は、コントローラであり、例えば、ASIC や FPGA 等の集積回路により実現される。

20

【 0 0 8 1 】

図 5 に示すように、制御部 1 5 は、受信部 1 5 1 と、受付部 1 5 2 と、推論部 1 5 3 と、処理部 1 5 4 と、送信部 1 5 5 とを有し、以下に説明する情報処理の機能や作用を実現または実行する。なお、制御部 1 5 の内部構成は、図 5 に示した構成に限られず、後述する情報処理を行う構成であれば他の構成であってもよい。

【 0 0 8 2 】

( 受信部 1 5 1 )

受信部 1 5 1 は、通信部 1 1 を介して、情報処理装置 1 0 0 から情報を受信する。受信部 1 5 1 は、情報処理装置 1 0 0 から提供されたモデルを受信する。

30

【 0 0 8 3 】

受信部 1 5 1 は、機械学習のモデルの学習に用いる入力用データと、当該入力用データに含まれる認識対象を示す正解データと、認識対象に関連する分類結果を示す分類ラベルとを含む学習用データを用いて生成されたモデルであって、データの入力に応じて、正解データに対応する第 1 出力と分類ラベルに対応する第 2 出力とを出力するモデルを受信する。受信部 1 5 1 は、音声認識に関するモデルを受信する。受信部 1 5 1 は、音声データの入力に応じて、当該音声データに対応する文字データを第 1 出力として出力するモデルを受信する。

40

【 0 0 8 4 】

( 受付部 1 5 2 )

受付部 1 5 2 は、各種情報を受け付ける。例えば、受付部 1 5 2 は、入力部 1 3 を介してユーザによる入力を受け付ける。受付部 1 5 2 は、ユーザによる操作を受け付ける。受付部 1 5 2 は、表示部 1 4 により表示された情報に対するユーザの操作を受け付ける。受付部 1 5 2 は、ユーザによる発話を入力として受け付ける。例えば、受付部 1 5 2 は、ユーザ U 1 による「XXXX」という発話を入力として受け付ける。

【 0 0 8 5 】

( 推論部 1 5 3 )

推論部 1 5 3 は、推論処理を行う。推論部 1 5 3 は、記憶部 1 2 に記憶されたモデルを

50

用いて、推論処理を行う。推論部 153 は、受信部 151 により受信されたモデルを用いて推論を行う。推論部 153 は、受信部により受信されたモデルにデータを入力することにより、当該データに対応する第 1 出力と第 2 出力とを生成する推論処理を行う。推論部 153 は、モデルに音声データを入力することにより、当該音声データに対応する推論処理を行う。推論部 153 は、モデルに音声データを入力することにより、当該音声データに対応する文字データである第 1 出力と第 2 出力とを生成する推論処理を行う。

【0086】

(処理部 154)

処理部 154 は、推論部 153 の推論結果を用いて各種の処理を実行する。処理部 154 は、推論部 153 の推論結果を表示部 14 に表示する。また、処理部 154 は、推論において入力に用いた音声データと、その音声データを書き起こした文字データと分類ラベルとのセットを学習用データとして、情報処理装置 100 に提供する。処理部 154 は、推論において入力に用いた音声データと、その音声データに対応する出力結果をユーザが修正したデータとのセットを学習用データとして、情報処理装置 100 に提供する。処理部 154 は、学習用データを送信部 155 に送信することを要求する。

10

【0087】

(送信部 155)

送信部 155 は、通信部 11 を介して、情報処理装置 100 へ情報を送信する。送信部 155 は、処理部 154 からの要求に応じて、通信部 11 を介して、学習用データを情報処理装置 100 に送信する。送信部 155 は、推論において入力に用いた音声データと、その音声データを書き起こした文字データと分類ラベルとのセットを学習用データとして、情報処理装置 100 に送信する。送信部 155 は、推論において入力に用いた音声データと、その音声データに対応する出力結果をユーザが修正したデータとのセットを学習用データとして、情報処理装置 100 に送信する。

20

【0088】

なお、上述した制御部 15 による各処理は、例えば、JavaScript (登録商標) などにより実現されてもよい。また、上述した表示処理が所定のアプリケーション (例えば音声文字変換アプリ等) により行われる場合や推論処理等の処理が専用アプリにより行われる場合、制御部 15 は、例えば、所定のアプリや専用アプリを制御するアプリ制御部を有してもよい。

30

【0089】

{ 5 . 処理フロー }

次に、図 6 を用いて、実施形態に係る情報処理システム 1 による情報処理の手順について説明する。図 6 は、実施形態に係る情報処理装置による処理の一例を示すフローチャートである。

【0090】

図 6 に示すように、情報処理装置 100 は、機械学習のモデルの学習に用いる入力用データと、正解データと、分類ラベルとを含む学習用データを取得する (ステップ S101)。すなわち、情報処理装置 100 は、機械学習のモデルの学習に用いる入力用データと、当該入力用データに含まれる認識対象を示す正解データと、認識対象に関連する分類結果を示す分類ラベルとを含む学習用データを取得する。

40

【0091】

情報処理装置 100 は、データの入力に応じて、正解データに対応する第 1 出力と分類ラベルに対応する第 2 出力とを出力するモデルを学習する (ステップ S102)。情報処理装置 100 は、学習したモデルを端末装置 10 へ提供する (ステップ S103)。

【0092】

次に、図 7 を用いて端末装置 10 におけるモデルを用いた推論等の処理の流れを示す。図 7 は、実施形態に係る端末装置による処理の一例を示すフローチャートである。

【0093】

図 7 に示すように、端末装置 10 は、モデルを受信していない場合 (ステップ S201

50

: No)、モデルを受信するまで待機する。端末装置10は、モデルを受信した後(ステップS201: Yes)、音声入力を受け付けていない場合(ステップS202: No)、音声入力を受け付けるまで待機する。

【0094】

端末装置10は、音声入力を受け付けた場合(ステップS202: Yes)、モデルに音声入力に対応する音声データを入力することにより、音声データに対応する第1出力と第2出力とを生成する推論処理を実行する(ステップS203)。

【0095】

そして、端末装置10は、推論結果を表示する(ステップS204)。また、端末装置10は、音声データと正解データと分類ラベルとのセットを、学習用データとして情報処理装置100へ送信してもよい。

【0096】

〔6.効果〕

上述してきたように、実施形態に係る情報処理装置100は、取得部131と、学習部133とを有する。取得部131は、機械学習のモデルの学習に用いる入力用データと、当該入力用データに含まれる認識対象を示す正解データと、認識対象に関連する分類結果を示す分類ラベルとを含む学習用データを取得する。学習部133は、学習用データを用いて、データの入力に応じて、正解データに対応する第1出力と分類ラベルに対応する第2出力とを出力するモデルを学習する。

【0097】

これにより、実施形態に係る情報処理装置100は、入力用データに含まれる認識対象を示す第1出力と、認識対象に関連する分類結果を示す第2出力との複数の種別の出力を行うモデルを学習することができる。したがって、情報処理装置100は、複数種別の出力を行うモデルを利用可能にすることができる。

【0098】

また、実施形態に係る情報処理装置100において、取得部131は、音声データである入力用データと、当該入力用データに含まれる音声の音声認識の結果を示す正解データとを含む学習用データを取得する。学習部133は、入力された音声データに対する音声認識の結果を示す第1出力と第2出力とを出力するモデルを学習する。

【0099】

これにより、実施形態に係る情報処理装置100は、音声認識に関して複数の種別の出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0100】

また、実施形態に係る情報処理装置100において、取得部131は、入力用データに対応する文字データである正解データを含む学習用データを取得する。学習部133は、入力された音声データが変換された文字データである第1出力と第2出力とを出力するモデルを学習する。

【0101】

これにより、実施形態に係る情報処理装置100は、音声文字変換に関して複数の種別の出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0102】

また、実施形態に係る情報処理装置100において、取得部131は、入力用データに関連するコンテキストの分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、入力されたデータに関連するコンテキストの分類結果を示す第2出力とを出力するモデルを学習する。

【0103】

これにより、実施形態に係る情報処理装置100は、認識結果を示す第1出力と、入力用データに関連するコンテキストの分類結果を示す第2出力を行うモデルを柔軟に学習す

10

20

30

40

50

ることができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0104】

また、実施形態に係る情報処理装置100において、取得部131は、入力用データが検知された場所の分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、入力されたデータが検知された場所の分類結果を示す第2出力とを出力するモデルを学習する。

【0105】

これにより、実施形態に係る情報処理装置100は、認識結果を示す第1出力と、入力用データが検知された場所の分類結果を示す第2出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

10

【0106】

また、実施形態に係る情報処理装置100において、取得部131は、入力用データに含まれる発話を行ったユーザに関連するコンテキストの分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、入力されたデータに含まれる発話を行ったユーザに関連するコンテキストの分類結果を示す第2出力とを出力するモデルを学習する。

【0107】

これにより、実施形態に係る情報処理装置100は、認識結果を示す第1出力と、発話を行ったユーザに関連するコンテキストの分類結果を示す第2出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

20

【0108】

また、実施形態に係る情報処理装置100において、取得部131は、ユーザの発話以外の分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、ユーザの発話以外の分類結果を示す第2出力とを出力するモデルを学習する。

【0109】

これにより、実施形態に係る情報処理装置100は、認識結果を示す第1出力と、ユーザの発話以外の分類結果を示す第2出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0110】

また、実施形態に係る情報処理装置100において、取得部131は、ユーザの周囲の状況の分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、ユーザの周囲の状況の分類結果を示す第2出力とを出力するモデルを学習する。

30

【0111】

これにより、実施形態に係る情報処理装置100は、認識結果を示す第1出力と、発話を行ったユーザの周囲の状況の分類結果を示す第2出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0112】

また、実施形態に係る情報処理装置100において、取得部131は、ユーザの属性の分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、ユーザの属性の分類結果を示す第2出力とを出力するモデルを学習する。

40

【0113】

これにより、実施形態に係る情報処理装置100は、認識結果を示す第1出力と、発話を行ったユーザの属性の分類結果を示す第2出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0114】

また、実施形態に係る情報処理装置100において、取得部131は、ユーザが利用する端末装置の分類結果を示す分類ラベルを含む学習用データを取得する。学習部133は、第1出力と、ユーザが利用する端末装置の分類結果を示す第2出力とを出力するモデルを学習する。

【0115】

50

これにより、実施形態に係る情報処理装置 100 は、認識結果を示す第 1 出力と、発話を行ったユーザが利用する端末装置の分類結果を示す第 2 出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

【0116】

また、実施形態に係る情報処理装置 100 において、取得部 131 は、入力用データに含まれる情報のうち、認識対象以外の情報の種別を示す分類ラベルを取得する。

【0117】

これにより、実施形態に係る情報処理装置 100 は、入力用データに含まれる情報のうち、認識対象以外の情報の種別を示す第 2 出力を行うモデルを柔軟に学習することができ、複数種別の出力を行うモデルを利用可能にすることができる。

10

【0118】

また、実施形態に係る端末装置 10 は、受信部 151 と、推論部 153 とを有する。受信部 151 は、機械学習のモデルの学習に用いる入力用データと、当該入力用データに含まれる認識対象を示す正解データと、認識対象に関連する分類結果を示す分類ラベルとを含む学習用データを用いて生成されたモデルであって、データの入力に応じて、正解データに対応する第 1 出力と分類ラベルに対応する第 2 出力とを出力するモデルを受信する。推論部 153 は、受信部により受信されたモデルにデータを入力することにより、当該データに対応する第 1 出力と第 2 出力とを生成する推論処理を行う。

【0119】

これにより、実施形態に係る端末装置 10 は、入力したデータに含まれる認識対象を示す第 1 出力と、認識対象に関連する分類結果を示す第 2 出力との複数の種別の出力を行うモデルを用いて推論処理を行うことができる。したがって、端末装置 10 は、複数種別の出力を行うモデルを利用した処理を行うことができる。

20

【0120】

また、実施形態に係る端末装置 10 において、受信部 151 は、音声認識に関するモデルを受信する。推論部 153 は、モデルに音声データを入力することにより、当該音声データに対応する推論処理を行う。

【0121】

これにより、実施形態に係る情報処理装置 100 は、音声認識に関する第 1 出力及び第 2 出力を出力するモデルを用いて推論処理を行うことができるため、複数種別の出力を行うモデルを利用した処理を行うことができる。

30

【0122】

また、実施形態に係る端末装置 10 において、受信部 151 は、音声データの入力に応じて、当該音声データに対応する文字データを第 1 出力として出力するモデルを受信する。推論部 153 は、モデルに音声データを入力することにより、当該音声データに対応する文字データである第 1 出力と第 2 出力とを生成する推論処理を行う。

【0123】

これにより、実施形態に係る情報処理装置 100 は、音声文字変換の結果である第 1 出力と第 2 出力との 2 つの異なる種別の情報を出力するモデルを用いて推論処理を行うことができるため、複数種別の出力を行うモデルを利用した処理を行うことができる。

40

【0124】

〔7. ハードウェア構成〕

また、上述した実施形態に係る端末装置 10 や情報処理装置 100 は、例えば図 8 に示すような構成のコンピュータ 1000 によって実現される。以下、情報処理装置 100 を例に挙げて説明する。図 8 は、ハードウェア構成の一例を示す図である。コンピュータ 1000 は、出力装置 1010、入力装置 1020 と接続され、演算装置 1030、一次記憶装置 1040、二次記憶装置 1050、出力 I/F (Interface) 1060、入力 I/F 1070、ネットワーク I/F 1080 がバス 1090 により接続された形態を有する。

【0125】

演算装置 1030 は、一次記憶装置 1040 や二次記憶装置 1050 に格納されたプロ

50

グラムや入力装置 1 0 2 0 から読み出したプログラム等に基づいて動作し、各種の処理を実行する。演算装置 1 0 3 0 は、例えば CPU (Central Processing Unit)、MPU (Micro Processing Unit)、ASIC (Application Specific Integrated Circuit) や FPGA (Field Programmable Gate Array) 等により実現される。

**【 0 1 2 6 】**

一次記憶装置 1 0 4 0 は、RAM (Random Access Memory) 等、演算装置 1 0 3 0 が各種の演算に用いるデータを一次的に記憶するメモリ装置である。また、二次記憶装置 1 0 5 0 は、演算装置 1 0 3 0 が各種の演算に用いるデータや、各種のデータベースが登録される記憶装置であり、ROM (Read Only Memory)、HDD (Hard Disk Drive)、SSD (Solid State Drive)、フラッシュメモリ等により実現される。二次記憶装置 1 0 5 0 は、内蔵ストレージであってもよいし、外付けストレージであってもよい。また、二次記憶装置 1 0 5 0 は、USBメモリやSD (Secure Digital) メモリカード等の取り外し可能な記憶媒体であってもよい。また、二次記憶装置 1 0 5 0 は、クラウドストレージ (オンラインストレージ) やNAS (Network Attached Storage)、ファイルサーバ等であってもよい。

10

**【 0 1 2 7 】**

出力 I / F 1 0 6 0 は、ディスプレイ、プロジェクタ、及びプリンタ等といった各種の情報を出力する出力装置 1 0 1 0 に対し、出力対象となる情報を送信するためのインターフェイスであり、例えば、USB (Universal Serial Bus) やDVI (Digital Visual Interface)、HDMI (登録商標) (High Definition Multimedia Interface) といった規格のコネクタにより実現される。また、入力 I / F 1 0 7 0 は、マウス、キーボード、キーパッド、ボタン、及びスキャナ等といった各種の入力装置 1 0 2 0 から情報を受信するためのインターフェイスであり、例えば、USB 等により実現される。

20

**【 0 1 2 8 】**

また、出力 I / F 1 0 6 0 及び入力 I / F 1 0 7 0 はそれぞれ出力装置 1 0 1 0 及び入力装置 1 0 2 0 と無線で接続してもよい。すなわち、出力装置 1 0 1 0 及び入力装置 1 0 2 0 は、ワイヤレス機器であってもよい。

**【 0 1 2 9 】**

また、出力装置 1 0 1 0 及び入力装置 1 0 2 0 は、タッチパネルのように一体化していてもよい。この場合、出力 I / F 1 0 6 0 及び入力 I / F 1 0 7 0 も、入出力 I / F として一体化していてもよい。

30

**【 0 1 3 0 】**

なお、入力装置 1 0 2 0 は、例えば、CD (Compact Disc)、DVD (Digital Versatile Disc)、PD (Phase change rewritable Disk) 等の光学記録媒体、MO (Magneto-Optical disk) 等の光磁気記録媒体、テープ媒体、磁気記録媒体、又は半導体メモリ等から情報を読み出す装置であってもよい。

**【 0 1 3 1 】**

ネットワーク I / F 1 0 8 0 は、ネットワーク N を介して他の機器からデータを受信して演算装置 1 0 3 0 へ送り、また、ネットワーク N を介して演算装置 1 0 3 0 が生成したデータを他の機器へ送信する。

40

**【 0 1 3 2 】**

演算装置 1 0 3 0 は、出力 I / F 1 0 6 0 や入力 I / F 1 0 7 0 を介して、出力装置 1 0 1 0 や入力装置 1 0 2 0 の制御を行う。例えば、演算装置 1 0 3 0 は、入力装置 1 0 2 0 や二次記憶装置 1 0 5 0 からプログラムを一次記憶装置 1 0 4 0 上にロードし、ロードしたプログラムを実行する。

**【 0 1 3 3 】**

例えば、コンピュータ 1 0 0 0 が情報処理装置 1 0 0 として機能する場合、コンピュータ 1 0 0 0 の演算装置 1 0 3 0 は、一次記憶装置 1 0 4 0 上にロードされたプログラムを実行することにより、制御部 1 3 0 の機能を実現する。また、コンピュータ 1 0 0 0 の演算装置 1 0 3 0 は、ネットワーク I / F 1 0 8 0 を介して他の機器から取得したプログラ

50

ムを一次記憶装置 1040 上にロードし、ロードしたプログラムを実行してもよい。また、コンピュータ 1000 の演算装置 1030 は、ネットワーク I/F 1080 を介して他の機器と連携し、プログラムの機能やデータ等を他の機器の他のプログラムから呼び出して利用してもよい。

【0134】

〔8.その他〕

以上、本願の実施形態を説明したが、これら実施形態の内容により本発明が限定されるものではない。また、前述した構成要素には、当業者が容易に想定できるもの、実質的に同一のもの、いわゆる均等の範囲のものが含まれる。さらに、前述した構成要素は適宜組み合わせることが可能である。さらに、前述した実施形態の要旨を逸脱しない範囲で構成要素の種々の省略、置換又は変更を行うことができる。

10

【0135】

また、上記実施形態において説明した各処理のうち、自動的に行われるものとして説明した処理の全部又は一部を手動的に行うこともでき、あるいは、手動的に行われるものとして説明した処理の全部又は一部を公知の方法で自動的に行うこともできる。この他、上記文書中や図面中で示した処理手順、具体的名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。例えば、各図に示した各種情報は、図示した情報に限られない。

【0136】

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部又は一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的又は物理的に分散・統合して構成することができる。

20

【0137】

例えば、上述した情報処理装置 100 は、複数のサーバコンピュータで実現してもよく、また、機能によっては外部のプラットフォーム等を API (Application Programming Interface) やネットワークコンピューティング等で呼び出して実現するなど、構成は柔軟に変更できる。また、情報処理装置 100 と端末装置 10 とは一体であってもよい。この場合、例えばユーザが利用する端末装置 10 が情報処理装置 100 としての機能を有してもよい。

30

【0138】

また、上述してきた実施形態及び変形例は、処理内容を矛盾させない範囲で適宜組み合わせることが可能である。

【0139】

また、上述してきた「部 (section、module、unit)」は、「手段」や「回路」などに読み替えることができる。例えば、取得部は、取得手段や取得回路に読み替えることができる。

【符号の説明】

【0140】

- 1 情報処理システム
- 100 情報処理装置
- 120 記憶部
- 121 学習用データ記憶部
- 122 モデル情報記憶部
- 130 制御部
- 131 取得部
- 132 決定部
- 133 学習部
- 134 提供部
- 10 端末装置

40

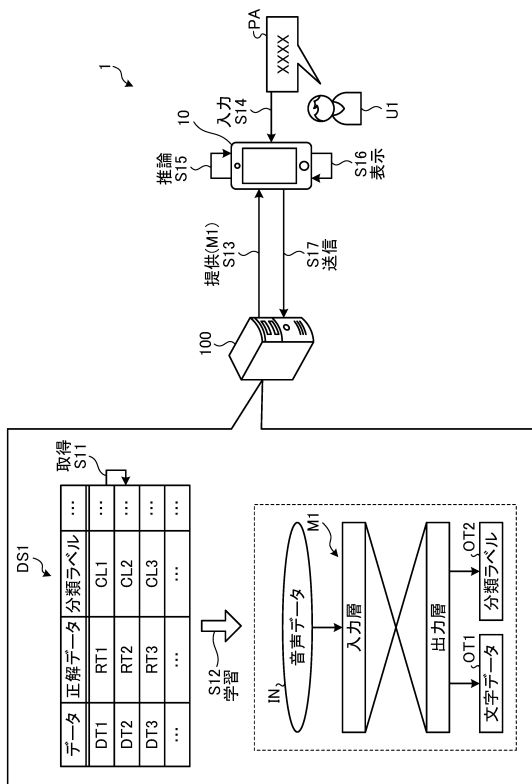
50

- 1 1 通信部
- 1 2 記憶部
- 1 3 入力部
- 1 4 表示部
- 1 5 制御部
- 1 5 1 受信部
- 1 5 2 受付部
- 1 5 3 推論部
- 1 5 4 処理部
- 1 5 5 送信部

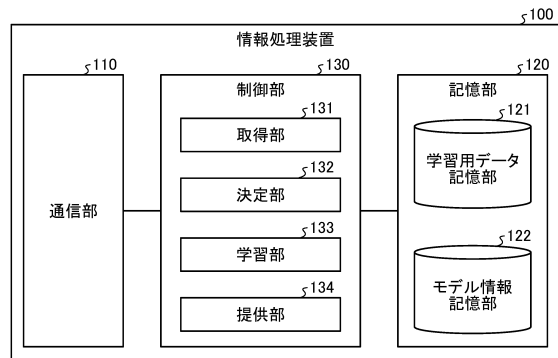
10

【図面】

【図 1】



【図 2】



20

30

40

50

【図3】

121  
↓

データセットID	データID	データ	正解データ	分類ラベル	...
DS1	DID1	DT1	RT1	CL1	...
	DID2	DT2	RT2	CL2	...
	DID3	DT3	RT3	CL3	...
	DID4	DT4	RT4	CL4	...
	DID5	DT5	RT5	CL5	...
	DID6	DT6	RT6	CL6	...
	DID7	DT7	RT7	CL7	...
	DID8	DT8	RT8	CL8	...
...	...	...	...	...	...
...	...	...	...	...	...

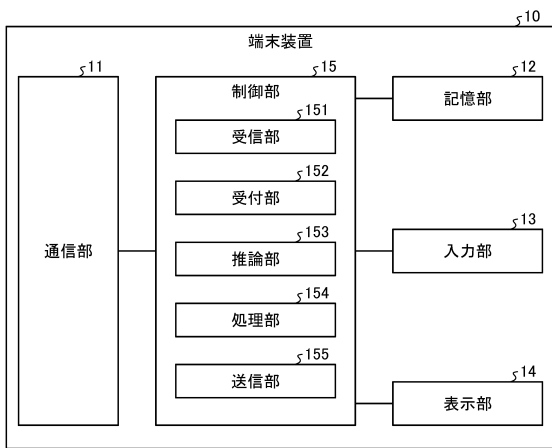
【図4】

122  
↓

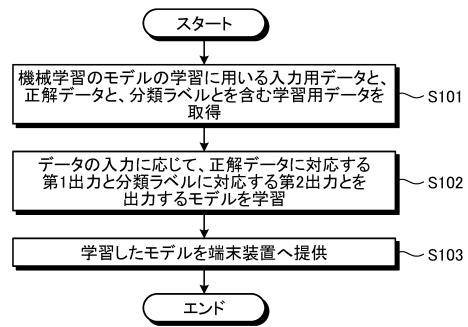
モデルID	用途	モデルデータ	...
M1	音声文字変換 コンテキスト推定	MDT1	...
...	...	...	...

10

【図5】



【図6】



20

30

40

50



## フロントページの続き

## (56)参考文献

特開2019-087229(JP,A)

特開2020-064253(JP,A)

特開2020-140673(JP,A)

大町 基, 単語の表記と素性を同時出力するend-to-end音声認識, 日本音響学会  
2020年 秋季研究発表会講演論文集CD-ROM [CD-ROM], 一般社団法人日本  
音響学会, 2020年08月26日, pp.815-818

早川 友瑛, End-to-End複数言語音声認識モデルにおける様々なマルチタスク学  
習の検討, 日本音響学会 2020年 秋季研究発表会講演論文集CD-ROM [CD-R  
OM], 一般社団法人日本音響学会, 2020年08月26日, pp.833-834

松原 拓未, CNN Autoencoderから抽出したボトルネック特徴量を用いた環  
境音分類, マルチメディア, 分散, 協調とモバイル(DICOMO2019)シンポジウ  
ム論文集[CD-ROM], 一般社団法人情報処理学会, 2019年06月26日, Vol.2019, No  
.1, p.339-346, ISSN: 1882-0840

## (58)調査した分野 (Int.Cl., DB名)

G06N 20/00

G06N 3/08

G10L 15/10

G10L 15/16