

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3848319号

(P3848319)

(45) 発行日 平成18年11月22日(2006.11.22)

(24) 登録日 平成18年9月1日(2006.9.1)

(51) Int. Cl.		F I			
HO4N	5/91	(2006.01)	HO4N	5/91	R
G1OL	15/00	(2006.01)	G1OL	15/00	200G

請求項の数 25 (全 42 頁)

(21) 出願番号	特願2003-381637 (P2003-381637)	(73) 特許権者	000001007
(22) 出願日	平成15年11月11日(2003.11.11)		キヤノン株式会社
(65) 公開番号	特開2005-150841 (P2005-150841A)		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成17年6月9日(2005.6.9)	(74) 代理人	100076428
審査請求日	平成15年11月11日(2003.11.11)		弁理士 大塚 康徳
		(74) 代理人	100112508
			弁理士 高柳 司郎
		(74) 代理人	100115071
			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		(72) 発明者	深田 俊明
			東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

最終頁に続く

(54) 【発明の名称】 情報処理方法及び情報処理装置

(57) 【特許請求の範囲】

【請求項1】

画像データと音声データを対応付ける情報処理方法であって、
 文字を含む前記画像データから文字領域を検出する検出工程と、
 検出された前記文字領域に含まれる文字情報を認識して複数の文字情報を取得する第1の認識工程と、
 前記音声データ中の複数の音声区間のそれぞれに対応する音声認識結果を取得する第2の認識工程と、
 前記第1の認識工程により取得された前記複数の文字情報と前記第2の認識工程により取得された前記複数の音声認識結果とをそれぞれ照合することによって、該文字情報と該音声認識結果を対応付ける対応付け工程と
 を有することを特徴とする情報処理方法。

【請求項2】

前記対応付け工程は、前記複数の文字情報と前記複数の音声認識結果の対応付けに基づいて、該文字情報に対応する前記文字領域と該音声認識結果、該文字情報と該音声認識結果に対応する音声区間、又は該文字情報に対応する前記文字領域と該音声認識結果に対応する音声区間のいずれかを対応付けることを特徴とする請求項1記載の情報処理方法。

【請求項3】

前記文字情報を発音列に変換する発音列変換工程をさらに有し、
 前記対応付け工程は、前記文字情報に基づく発音列と前記音声認識結果の発音列との照

10

20

合結果に基づいて前記文字情報と前記音声認識結果を対応付けることを特徴とする請求項 1 に記載の情報処理方法。

【請求項 4】

前記音声認識結果を文字列に変換する文字列変換工程をさらに有し、前記対応付け工程は、前記文字情報の文字列と前記音声認識結果に基づく文字列との照合結果に基づいて前記文字情報と前記音声認識結果を対応付けることを特徴とする請求項 1 に記載の情報処理方法。

【請求項 5】

前記文字情報を発音列に変換する発音列変換工程と、前記音声認識結果を文字列に変換する文字列変換工程をさらに有し、前記対応付け工程は、前記文字情報に基づく発音列と前記音声認識結果の発音列との照合結果と、前記文字情報の文字列と前記音声認識結果に基づく文字列との照合結果とに基づいて、前記文字情報と前記音声認識結果を対応付けることを特徴とする請求項 1 に記載の情報処理方法。

10

【請求項 6】

前記第 1 の認識工程が、前記複数の文字情報のそれぞれについて、候補と該候補の度合いを取得し、

前記第 2 の認識工程が、前記複数の音声認識結果のそれぞれについて、候補と該候補の度合いを取得し、

前記文字情報の候補の度合いと前記音声認識結果の候補の度合いとに基づいて、それぞれの候補間の関連の度合いを算出する算出工程とをさらに有し、

20

前記対応付け工程が、前記関連の度合いの高さに応じて、前記文字情報の候補と前記音声認識結果の候補とを対応付ける

ことを特徴とする請求項 1 に記載の情報処理方法。

【請求項 7】

前記候補の度合いは、前記候補の認識確率又は認識尤度であることを特徴とする請求項 6 に記載の情報処理方法。

【請求項 8】

前記算出工程が、前記文字情報の候補又は前記音声認識結果の候補に重み付けを付与して前記候補間の関連の度合いを算出することを特徴とする請求項 6 に記載の情報処理方法。

30

【請求項 9】

前記第 2 の認識工程は、前記音声データ中の複数の音声区間それぞれに対応する音声認識結果を文字列に変換し、該音声認識結果の文字列中の前記第 1 の認識工程で取得した文字情報に含まれない文字列を除外したものを、前記音声認識結果として取得することを特徴とする請求項 1 に記載の情報処理方法。

【請求項 10】

前記第 2 の認識工程は、前記第 1 の認識工程で取得した文字情報を音声認識対象として音声認識を行い、前記音声認識結果を取得することを特徴とする請求項 1 に記載の情報処理方法。

40

【請求項 11】

前記第 2 の認識工程は、前記音声データ中の複数の音声区間それぞれに対応する音声認識結果から、前記第 1 の認識工程で取得した文字情報を発音列に変換したものに含まれない音声認識結果を除外したものを、前記音声認識結果として取得することを特徴とする請求項 1 に記載の情報処理方法。

【請求項 12】

前記第 2 の認識工程は、前記第 1 の認識工程で取得した文字情報を発音列に変換したものを音声認識対象として音声認識を行い、前記音声認識結果を取得することを特徴とする請求項 1 に記載の情報処理方法。

【請求項 13】

50

少なくとも自立語を含む重要語を抽出するためのデータに基づいて、前記文字情報に含まれる重要語を抽出する重要語抽出工程をさらに有し、

前記第2の認識工程は、前記重要語をキーワードスポッティングの対象とするか、前記重要語の音声認識用言語モデルの確率値を増加させるかの少なくともいずれかを行って音声認識を行うことで、前記音声認識結果を取得する

ことを特徴とする請求項1に記載の情報処理方法。

【請求項14】

前記第1の認識工程により認識された前記文字情報について、該文字情報のフォントサイズ、色、アンダーライン、太字、斜体、又はフォント種の少なくとも何れか1つを含むフォント情報を抽出するフォント情報抽出工程をさらに有し、

前記第2の認識工程が、前記フォント情報を利用して特定された文字列をキーワードスポッティングの対象するか、特定された文字列の統計的言語モデルの確率値を増加させるかの少なくともいずれかを行って音声認識を行い、前記音声認識結果を取得する

ことを特徴とする請求項1に記載の情報処理方法。

【請求項15】

前記第1の認識工程は、前記第2の認識工程により取得された音声認識結果を文字列に変換したものに含まれない文字列を、前記検出された文字領域に含まれる文字情報を認識して取得した複数の文字情報から除外したものを、前記複数の文字情報として取得することを特徴とする請求項1記載の情報処理方法。

【請求項16】

前記第1の認識工程は、前記第2の認識工程により取得された音声認識結果を文字列に変換したものを文字認識対象として文字認識を行い、前記文字情報を取得することを特徴とする請求項1記載の情報処理方法。

【請求項17】

前記画像データを複数の領域に分割して分割画像を取得する画像分割工程をさらに有し、

それぞれの分割画像に関して文字情報を認識する

ことを特徴とする請求項1に記載の情報処理方法。

【請求項18】

前記第1の認識工程により認識された前記文字情報を文字概念表現に変換する文字概念変換工程と、

前記第2の認識工程により認識された前記音声認識結果を音声概念表現に変換する音声概念変換工程と、

前記文字概念表現と前記音声概念表現とを照合する概念対応工程とをさらに有し、

前記対応付け工程が、前記概念対応工程によって得られる概念間の照合結果に基づいて、前記文字情報と前記音声認識結果とを対応付ける

ことを特徴とする請求項1に記載の情報処理方法。

【請求項19】

画像データと音声データを対応付ける情報処理方法であって、

前記画像データに含まれるオブジェクト領域を検出する第1の検出工程と、

検出された前記オブジェクト領域からオブジェクト情報を認識する第1の認識工程と、

前記音声データ中の複数の音声区間のそれぞれに対応する音声認識結果を取得する第2の認識工程と、

前記第1の認識工程により認識された前記オブジェクト情報の特徴情報に対応する文字情報と前記第2の認識工程により認識された前記音声認識結果とを照合することによって該オブジェクト情報と該音声認識結果とを対応付ける対応付け工程と

を有することを特徴とする情報処理方法。

【請求項20】

前記オブジェクト情報は図形情報であり、

前記オブジェクト情報の特徴情報は、前記図形情報の形状、色の少なくともいずれかで

10

20

30

40

50

あることを特徴とする請求項 19 記載の情報処理方法。

【請求項 21】

画像データと音声データを対応付ける情報処理方法であって、
前記画像データに含まれる人物領域を検出する第 1 の検出工程と、
検出された前記人物領域から人物又は少なくとも人物の性別、年代のいずれかを含む人物のクラスを認識する第 1 の認識工程と、
前記音声データ中の複数の音声区間それぞれに対応する話者又は少なくとも話者の性別、年代のいずれかを含む話者クラスを認識する第 2 の認識工程と、
前記第 1 の認識工程により認識された人物又は人物のクラスと、前記第 2 の認識工程により認識された話者又は話者クラスとを対応付ける対応付け工程と
を有することを特徴とする情報処理方法。

10

【請求項 22】

請求項 1 乃至 21 に記載の情報処理方法をコンピュータに実行させるための制御プログラム。

【請求項 23】

画像データと音声データを対応付ける情報処理装置であって、
文字を含む前記画像データから文字領域を検出する検出手段と、
検出された前記文字領域に含まれる文字情報を認識して複数の文字情報を取得する第 1 の認識手段と、
前記音声データ中の複数の音声区間それぞれに対応する音声認識結果を取得する第 2 の認識手段と、
前記第 1 の認識手段により取得された前記複数の文字情報と前記第 2 の認識手段により取得された前記複数の音声認識結果とをそれぞれ照合することによって、該文字情報と該音声認識結果を対応付ける対応付け手段と
を有することを特徴とする情報処理装置。

20

【請求項 24】

画像データと音声データを対応付ける情報処理方法であって、
前記画像データに含まれるオブジェクト領域を検出する第 1 の検出手段と、
検出された前記オブジェクト領域からオブジェクト情報を認識する第 1 の認識手段と、
前記音声データ中の複数の音声区間それぞれに対応する音声認識結果を取得する第 2 の認識手段と、
前記第 1 の認識手段により認識された前記オブジェクト情報の特徴情報に対応する文字情報と前記第 2 の認識手段により認識された前記音声認識結果とを照合することによって該オブジェクト情報と該音声認識結果とを対応付ける対応付け手段と
を有することを特徴とする情報処理装置。

30

【請求項 25】

画像データと音声データを対応付ける情報処理方法であって、
前記画像データに含まれる人物領域を検出する第 1 の検出手段と、
検出された前記人物領域から人物又は少なくとも人物の性別、年代のいずれかを含む人物のクラスを認識する第 1 の認識手段と、
前記音声データ中の複数の音声区間それぞれに対応する話者又は少なくとも話者の性別、年代のいずれかを含む話者クラスを認識する第 2 の認識手段と、
前記第 1 の認識手段により認識された人物又は人物のクラスと、前記第 2 の認識手段により認識された話者又は話者クラスとを対応付ける対応付け工程と
を有することを特徴とする情報処理装置。

40

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像データと音声データを対応付ける情報処理方法及び情報処理装置に関する。

50

【背景技術】

【0002】

近年、デジタルカメラで静止画像を撮影するとともに、撮影された当該静止画像に対して音声メモ機能を用いて当該静止画像についてのコメント等を録音するといった、画像データと音声データを関連付ける技術が開発されている。例えば、E x i f (EXchangeable Image File Format) と呼ばれるデジタルカメラ用画像ファイルの標準規格では、1枚の静止画像ファイルの中に付属情報として音声データを関連付けることができる。このようにして静止画像に関連付けられた音声データは、単に静止画像に音声データが付与されたものであるだけでなく、当該音声データを音声認識することによって付与された音声データを認識して文字情報に変換し、文字又は音声をキーとして、複数の静止画像の中から所望の静止画像を検索するといった目的に利用することができる。

10

【0003】

また、ボイスレコーダー機能を搭載したデジタルカメラや、デジタルカメラ機能を搭載したボイスレコーダーでは、最大で数時間程度の音声データを録音することが可能である。

【発明の開示】

【発明が解決しようとする課題】

【0004】

しかしながら、上述したような従来の技術は、1枚の静止画像全体に対して1つ又は複数の音声データを関連付けるに留まっており、1枚の静止画像中の特定の部分領域と、それに対する音声データ中の特定の音声区間とを関連付ける技術ではない。すなわち、デジタルカメラで撮影された静止画像中の部分領域とボイスレコーダーで録音された音声データ中の部分音声データとを関連付けるような技術については、現時点において出願人は発見していない。

20

【0005】

ここで、展示会等において1枚のパネルを用いて、発表者が口頭で製品のプレゼンテーションをしている場面を想定する。このとき、当該プレゼンテーションの聴講者は、ボイスレコーダーで発表者のプレゼンテーションに関する音声を録音する一方で、展示されているポスター（例えば、ポスター全体）をデジタルカメラで静止画像として撮影することがある。そして、その後、当該聴講者が、プレゼンテーション終了後、自宅等において、プレゼンテーション時に撮影した静止画像と録音した音声を再生して、撮影した静止画像中のある部分領域（例えば、展示されていたポスター中の一部に記載されている「製品の特徴」に関する部分）に関するプレゼンテーションを聞く場合を考える。

30

【0006】

この場合、当該聴講者は、録音した音声データから該当する部分領域についての録音音声を人手によって探す必要があるため、非常に時間のかかる作業となるという問題がある。特に、当日プレゼンテーション会場に居合わせておらず、自宅等で初めて当該プレゼンテーションを視聴する人にとっては、撮影されたポスターの上記部分領域に対するプレゼンテーションが、録音された音声データ全体のどのあたりに記録されているのかが全く分からないため、所望の部分音声を探すためには、録音音声を最初から聞いていく必要があり非常に手間がかかるという問題がある。

40

【0007】

本発明は、このような事情を考慮してなされたものであり、画像データ中の部分画像領域と音声データ中の部分音声データとを好適に対応付けることができる情報処理方法及び情報処理装置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】

【0009】

上記課題を解決するために、本発明は、画像データと音声データを対応付ける情報処理方法であって、

文字を含む前記画像データから文字領域を検出する検出工程と、

50

検出された前記文字領域に含まれる文字情報を認識して複数の文字情報を取得する第 1 の認識工程と、

前記音声データ中の複数の音声区間のそれぞれに対応する音声認識結果を取得する第 2 の認識工程と、

前記第 1 の認識工程により取得された前記複数の文字情報と前記第 2 の認識工程により取得された前記複数の音声認識結果とをそれぞれ照合することによって、該文字情報と該音声認識結果を対応付ける対応付け工程と

を有することを特徴とする。

【0011】

また、上記課題を解決するために、本発明は、画像データと音声データを対応付ける情報処理装置であって、

文字を含む前記画像データから文字領域を検出する検出手段と、

検出された前記文字領域に含まれる文字情報を認識して複数の文字情報を取得する第 1 の認識手段と、

前記音声データ中の複数の音声区間それぞれに対応する音声認識結果を取得する第 2 の認識手段と、

前記第 1 の認識手段により取得された前記複数の文字情報と前記第 2 の認識手段により取得された前記複数の音声認識結果とをそれぞれ照合することによって、該文字情報と該音声認識結果を対応付ける対応付け手段と

を有することを特徴とする。

【発明の効果】

【0012】

本発明によれば、画像データ中の部分画像領域と音声データ中の部分音声データとを好適に対応付けることができる。これにより、例えば、画像データと音声データとから画像データ中の部分画像領域に関連する音声データ中の部分音声区間を人手によって探す手間が省け、大幅に時間を節約することができる。

【発明を実施するための最良の形態】

【0013】

以下、図面を参照して、本発明の好適な実施例について詳細に説明する。

【実施例 1】

【0014】

図 1 は、本発明の第 1 の実施例に係る画像データと音声データの部分データ同士を対応付ける静止画像・音声処理装置の構成を示すブロック図である。図 1 において、CPU 101 は、ROM 102 に記憶された制御プログラム或いは外部記憶装置 104 から RAM 103 にロードされた制御プログラムに従って、本実施例の静止画像・音声処理装置の各種制御・処理を行う。ROM 102 は、各種パラメータや CPU 101 が実行する制御プログラム等を格納している。RAM 103 は、CPU 101 による各種制御の実行時に作業領域を提供するとともに、CPU 101 により実行される制御プログラムを記憶する。

【0015】

外部記憶装置 104 は、ハードディスク、フレキシブルディスク、CD-ROM、DVD-ROM、メモリカード等で実現される固定式記憶装置或いは着脱可能な可搬記憶装置であり、例えば、外部記憶装置 104 がハードディスクの場合には、CD-ROM やフレキシブルディスク等からインストールされた各種プログラムが記憶される。105 は、マイクロフォン等による音声入力装置であり、音声入力装置 105 から取り込まれた音声は、CPU 101 による音声認識又は音響信号分析によって静止画像に関連した音声認識又は分析される。106 は、デジタルカメラ等による画像入力装置であり、取り込まれた画像は静止画像に変換され、文字認識やオブジェクト認識が行われる。

【0016】

107 は、CRT、液晶ディスプレイ等の表示装置であり、処理内容の設定・入力に関する表示・出力を行う。108 は、ボタン、テンキー、キーボード、マウス、ペン等の補

助入出力装置である。109は、上記各部を互いに接続するバスである。尚、静止画像及び当該静止画像に対応付けられるための音声データは、それぞれ画像入力装置106及び音声入力装置105によって入力してもよいし、別の装置等によって獲得したものをROM102、RAM103、外部記憶装置104若しくはネットワークを介して接続された外部装置に記憶しておいてもよい。

【0017】

図2は、第1の実施例で互いに部分データ同士の対応付け処理の対象となる静止画像(a)と当該静止画像に関連する音声(b)の一例について示す図である。図2に示すように、この静止画像には、白地に「春」、「夏」、「秋」、「冬」という4つの文字が撮像されている(以降、静止画像の左下を原点として、水平方向をx軸、垂直方向をy軸とする座標軸を用いる。尚、座標単位にはピクセルを用いることができるが、特にこれに限定されることはない。)。また、この静止画像に関連した音声は、「フユ」、「ハル」、「アキ」、「ナツ」という4つの発声がこの順で録音されている(以降、音声の開始時間を0とした時間軸を用いる。尚、時間単位としてはサンプル数や秒を用いることができるが、特にこれらに限定されることはない。)。また、この音声は、各発声間に十分な無音区間を含んでいるものとする。

10

【0018】

尚、この音声は、発声場所、発声時間、発声者に制限はない。すなわち、静止画像を撮影した場所、時間、撮影者は、当該音声の発声場所、発声時間、発声者と同じであっても異なってもよい。また、音声データは、Exif等のように静止画像ファイルの一部として含まれていてもよいし、静止画像とは別のファイルであってもよい。さらに、静止画像データと音声データは、同じ装置又は同じ記憶媒体に記憶されていてもよいし、ネットワーク等を介して別の場所に格納されているものであってもよい。

20

【0019】

図3は、本発明の第1の実施例において静止画像と音声を入力して静止画像と音声との対応関係(画像音声対応情報)を求める際のモジュール構成を示すブロック図である。図2において、201は文字検出部であり、静止画像から文字部分を含む所定領域(文字領域)を検出する。図2の例では、「春」、「夏」、「秋」、「冬」の4つの文字領域が矩形の部分画像として、座標情報(図2のx、yの値)と共に検出される。尚、文字検出部201で検出される当該部分画像は、あくまで画像データであって、文字データではない。ここで、図7は、図2に示す静止画像と音声の例に対する文字認識結果情報と音声認識結果情報に対応させた結果を示す図である。図7(a)に示すように、各部分画像データの座標情報は、各文字領域(部分画像)の中心座標を表している。

30

【0020】

また、図3において、202は文字認識部であり、文字検出部201で検出された各文字領域に対して文字認識を行う。尚、文字認識処理自体については、既存の技術を用いることが可能である。図2の例では、4つの文字領域の部分画像データから、文字認識部202によって、「春」、「夏」、「秋」、「冬」の4文字の文字データが認識される。ここで、図6は、図7に示す文字認識結果情報と音声認識結果情報の例を示す図である。図6(a)に示すように、文字認識部202によって、各文字データと中心座標とが認識結果から対応付けられている。

40

【0021】

図3において、203は音声検出部であり、音声データから例えば人が発声した部分(音声区間)を検出する。図2の例では、「フユ」、「ハル」、「アキ」、「ナツ」の4つの音声区間が部分音声データとして、時間情報(図2のtの値)と共に検出される。図7(b)に示すように、各音声区間の時間情報は、各音声区間の開始及び終了時間を表している。

【0022】

図3において、204は音声認識部であり、音声検出部203で検出された各音声区間に対して音声認識を行う。尚、音声認識処理自体については、既存の技術を用いることが

50

できる。ここでは、簡単のため、「春（ハル）」、「夏（ナツ）」、「秋（アキ）」、「冬（フユ）」の4単語のみを認識対象語彙とする単語音声認識を行った場合について考える。この場合、図2の例では、4つの音声区間の音声データが、音声認識部204によって、「冬」、「春」、「秋」、「夏」の4単語の文字データに変換される。図6(b)に示すように、音声認識部204によって、各音声区間の音声データと時間情報とが認識結果から対応付けられている。

【0023】

図3において、205は静止画音声対応部であり、文字検出部201と文字認識部202の処理結果として得られる静止画像内の文字認識結果及びその座標情報（文字認識結果情報）と、音声検出部203と音声認識部204の処理結果として得られる音声内の音声認識結果及びその時間情報（音声認識結果情報）を用いて、静止画像と音声データの対応付けを行う。例えば、図2に示す静止画像と音声の例では、図6(a)に示される文字認識結果情報による文字列と、図6(b)に示される音声認識結果情報に基づく文字列とを比較・照合する。図8は、第1の実施例における静止画像と音声との対応付けの一例を示す図である。

10

【0024】

図9は、静止画像と音声との対応結果を用いたアプリケーションの例である。図9に示す例では、静止画像中の文字が位置する部分（例えば、図9では座標(x1, y1)付近)にマウスカーソル（図9の矢印マーク）を持っていくと、この文字に対応した音声データ（図9では、図7(b)に示す時刻s2からe2までの音声データ）が再生され、スピーカー等の音声出力装置から出力される。

20

【0025】

尚、図9に示す例とは逆に、音声を先頭から、或いはマウス、キーボード等で任意の時間を指定することによってその間の音声を再生し、再生されている音声区間に対応する静止画像の対応部分に枠を付与して表示することも可能である。図60は、図2に示す静止画像と音声との対応結果を用いた別のアプリケーションに基づく表示例を示す図である。図60に示す例では、利用者が「フユ」と音声認識された音声区間（図7のs1からe1）にマウスカーソル（図60の矢印マーク）を持っていくと、当該音声区間に対応した文字領域（すなわち、「冬」）に文字領域分の外枠が生成・表示される。この結果、本装置の操作者は、出力されている音声の静止画像のどの部分に対応しているかを容易に理解することができる。

30

【0026】

以下、図3に示す文字検出部201から静止画音声対応部205の各モジュールの動作についてさらに詳細に説明する。

【0027】

文字検出部201は、静止画像中から写真、絵、文字、図形、図表等の所定領域を切り出す技術（セグメンテーション）を用いる。セグメンテーションの方法としては、文書中に存在する文字の部分と他の図表や画像等の部分と区別するための技術である文書認識技術といった既存の技術を用いることができる。尚、上述した文字領域の検出に関する説明では、簡単のため、文字領域の座標情報として、図7(a)に示したように文字領域の中心座標としているが、矩形領域を表すことが可能な座標（2点の座標）とするものが一般的であり、融通性があるため好適である。

40

【0028】

文字認識部202は、文字検出部201で検出された文字領域からなる部分画像データを入力として、これに含まれる文字を認識する。文字認識の方法としては、既存の文字認識技術を用いればよいが、本実施例では静止画像を入力としているため、オンライン文字認識技術は適用することはできず、オフライン文字認識又はOCR (Optical Character Recognition) 技術を用いる必要がある。また、文字の種類が文字認識を行う前に分かっている場合、或いは文字認識時に利用者等によって与えることが可能な場合には、その文字の種類に応じた文字認識方法を適用することができる。

50

【 0 0 2 9 】

ここでいう文字の種類とは、例えば、手書き文字と印刷活字文字である。手書き文字は、さらに、制限付き手書き文字（点線上にそって文字が書かれる文字等）、常用手書き文字、自由手書き文字に分類することもできる。また、印刷活字文字は、さらに、フォント種が1つのシングルフォント、複数のフォント種が混在するマルチフォントに分類することもできる。また、文字の種類が予め分からない場合には、これらの全ての手法を適用して最も信頼度やスコアの高い結果を利用する方法や、各文字の種類を文字認識前に判定して、判定結果に基づいた文字認識方法を適用する方法等を用いればよい。

【 0 0 3 0 】

図4は、第1の実施例における文字認識部202の細部モジュール構成を示すブロック図である。図4において、301は前処理部であり、文字認識処理を行い易くするための各種処理を施し、正規化データとして出力する。具体的には、雑音成分の除去、文字の大きさの正規化等を行う。302は特徴抽出部であり、その文字が表わす特徴を抽出する。これは正規化データを、よりその文字の特徴を捉えた次元数の低いデータへ変換・圧縮する。例えば、2値画像の輪郭におけるchain-code等を特徴として抽出する。

10

【 0 0 3 1 】

また、303は識別部であり、特徴抽出部302で得られた入力の特徴量を文字認識用テンプレート305と比較・照合（マッチング）することによって、入力特徴量の文字の識別を行う。マッチング方法としては、DPマッチング法や2次元HMM（Hidden Markov Model）法等を用いればよい。ここで、文字間の言語的な関係を言語知識として確率的に利用することにより、文字認識性能が向上する場合がある。306は、この場合に用いる文字認識用の言語モデルであり、具体的には、2つ組み文字の出現確率（文字バイグラム）等である。しかし、文字認識用言語モデル306は必ずしも必要なものではない。304は文字認識結果情報出力部であり、識別部303で得られる文字認識の結果と、対応する文字領域の静止画像における座標情報を文字認識結果情報として出力する。

20

【 0 0 3 2 】

204は、音声検出部203で検出された音声区間からなる音声データを入力としてこれを音声認識する音声認識部である。音声認識部204における音声認識の方法としては、HMMに基づく方法等の既存の音声認識技術を用いればよい。音声認識の方法としては、単語音声認識、文法ベースの連続音声認識、N-gramベースの大語彙連続音声認識、単語単位を用いない音素認識もしくは音節認識を用いることがある。上述した音声認識の説明では、簡単のため、単語音声認識を用いたが、実際には、単語単位で発声される保障はなく、発声内容も事前に分からないため、大語彙連続音声認識又は音素認識（音節認識）による方法を利用することが望ましい。

30

【 0 0 3 3 】

図5は、第1の実施例における音声認識部204の細部モジュール構成を示すブロック図である。図5において、401は、音声分析部で、音声を変換し、特徴量を求める。音声分析の方法としては、MFCC分析（Mel-Frequency Cepstrum Coefficient）や線形予測分析等を用いればよい。405は、音声認識を行う際の辞書（表記と読み）及び言語制約（単語N-gramや音素N-gram等の確率値）が格納されている。402は、探索部で、401で得られた入力音声の特徴量を404の音声認識用音響モデルと405の音声認識用言語モデルを用いることによって音声認識結果を得る。403は、音声認識結果情報出力部で、403で得られる音声認識の結果と、対応する音声区間の音声における時間情報を音声認識結果情報として出力する。

40

【 0 0 3 4 】

205は、文字認識部202から得られる文字認識結果情報と、音声認識部204から得られる音声認識結果情報を入力として静止画像と音声を対応付け、静止画音声対応情報を出力する。対応付けは、文字認識の結果得られる文字もしくは文字列と、音声認識の結果得られる表記（単語）から得られる文字もしくは文字列のマッチングを取ることによって行う。或いは、文字認識の結果得られた文字列の発音列と音声認識の結果得られた発音

50

列との照合によって行う。尚、これらの詳細については、以降の実施例において詳細に説明する。図2の例では、説明を簡単にするため、静止画像中の文字と音声の発声が1対1に対応している例を示した。

【0035】

よって、文字列のマッチングは、完全に一致するものを探すことにより対応付けが行える。しかしながら、実際にプレゼンテーション等で録音される音声は、静止画像の文字をそのまま発声することはほとんどないと考えられる。このような場合には、文字認識の結果得られる文字列を音声認識の結果得られる文字列に対して部分マッチングさせて対応付けを行う。

【0036】

例えば、「ここにある春は、...」や、「つまり、これは夏になると...」という発声がない場合、文字認識結果の「春」は前者、「夏」は後者の音声認識結果の部分文字列と一致するため、これらに対応付ける。さらに一般的には、文字領域に対する音声区間がない、文字領域とは関係のない音声区間がある、文字認識結果に誤りがある、音声認識結果に誤りがあることが考えられるため、一致するか否かといった決定的なマッチングではなく、どの程度マッチングするかといった確率的な柔軟なマッチングを行う必要がある。

【0037】

以上の説明から明らかなように、本実施形態によれば、静止画像データから静止画像の部分画像領域と抽出し、音声データから音声の部分音声区間を抽出し、お互いに関連のあるものを好適に対応付けることができるようになり、その結果として、画像データ中の部分画像領域に関連した音声データ中の音声区間(部分音声データ)を従来のように人手によって探す手間が省け、大幅に時間を節約することが可能となる。

【実施例2】

【0038】

上述した第1の実施例における静止画音声対応部205では、文字認識の結果として得られる文字列と、音声認識の結果として得られる文字列とを直接比較して対応付けていた。しかし、音声認識方法が音素(音節)認識であったり、同音異表記が出力された場合には、文字列の直接比較を行うことができない。例えば、文字認識結果が「春」であり、音声認識結果が「haru」、「ハル」、「張る」等の場合である。そこで、一般に、音声認識では入力音声の読み情報(発音列)が分かっていることから、文字認識結果を読み情報(発音列)に変換した後に、発音列同士でマッチングを取ることによって、文字列同士の比較ができないような場合においても文字認識結果情報と音声認識結果情報の対応を取ることが可能となる。

【0039】

図10は、本発明の第2の実施例における発音列マッチングによる静止画音声対応部の細部モジュール構成を示すブロック図である。図10において、501は、文字認識部202から得られる文字認識結果情報の文字認識結果を発音列に変換する文字認識結果発音列変換部である。502は、文字列を発音に変換するために文字認識結果発音列変換部501で用いられる発音変換辞書である。ここで、文字と発音の対応は、一般に、1対1ではなく1対多となるため、1つの文字列に対する発音列は多くの場合1種類ではなく、発音列候補として1つ又は複数出力される。

【0040】

具体的には、図6(a)に示される文字認識結果情報の、「春」、「夏」、「秋」、「冬」という文字列から、それぞれ「ハル/シュン」、「ナツ/カ」、「アキ/シュウ」、「フユ/トウ」というような発音列候補を得る。図11は、第2の実施例における文字認識結果と音声認識結果に対する発音列の例を示す図である。すなわち、図6(a)に示される文字認識結果情報から図11(a)に示されるような発音列候補を得る。

【0041】

図10において、503は、音声認識部204から得られる音声認識結果情報から発音

10

20

30

40

50

列を抽出する音声認識結果発音列抽出部である。具体的には、図6(b)に示される音声認識結果情報から、図11(b)に示すように、「フユ」、「ハル」、「アキ」、「ナツ」という発音列を抽出する。

【0042】

また、図10において、504は発音列マッチング部であり、文字認識結果の文字列を発音列に変換したものと音声認識結果の発音列とのマッチングを取る。このマッチング処理によって、図11に示す例では、文字認識結果の複数の発音列候補から「ハル」、「ナツ」、「アキ」、「フユ」が選択され、音声認識結果の発音列と対応付けられる。

【0043】

さらに、図10において、505は静止画音声対応情報出力部であり、マッチング結果を図8に示すような静止画音声対応情報として出力する。尚、この例では、発音列としてカタカナ表記を用いているが、これに限らず音素表現等別の表記を用いてもよいことは言うまでもない。また、文字認識結果の発音列候補は「シュウ」や「トウ」と書き言葉の発音列を生成していたが、「シュー」や「トー」といった話し言葉の発音列に変換した結果や、これを書き言葉の発音列に加えた結果を用いてもよい。

10

【0044】

以上の説明から明らかなように、本実施例によれば、文字認識結果の文字列と音声認識結果の文字列が直接比較できない場合においても、静止画像と音声の対応付けを行うことが可能となる。

【実施例3】

20

【0045】

上述した第2の実施例では、文字認識の結果として得られる文字列を発音列に変換し、音声認識の結果として得られる発音列とマッチングしていたが、これとは逆に、音声認識の発音列を文字列に変換し、文字認識結果の文字列とマッチングすることも可能である。

【0046】

図12は、本発明の第3の実施例における文字列マッチングを行う静止画音声対応部205の細部モジュール構成を示すブロック図である。図12において、601は、文字認識部202から得られる文字認識結果情報の文字認識結果から文字列を抽出する文字認識結果文字列抽出部である。具体的には、図6(a)に示される文字認識結果情報から、図13(a)に示されるように「春」、「夏」、「秋」、「冬」という文字列を抽出する。すなわち、図13は、第3の実施例における文字認識結果と音声認識結果に対する文字列の例である。

30

【0047】

図12において、602は、音声認識部204から得られる音声認識結果情報の音声認識結果(発音列)を文字列に変換する音声認識結果文字列変換部である。また、603は、音声認識結果文字列変換部602で発音列を文字列に変換する際に用いられる文字変換辞書である。ここで、発音と文字の対応は、一般に、1対1ではなく1対多となるため、1つの発音列に対する文字列は1種類ではなく、文字列候補として複数出力する。

【0048】

具体的には、図6(b)に示される音声認識結果情報の、「フユ」、「ハル」、「アキ」、「ナツ」という発音列から、図13(b)に示すように、それぞれ「冬/不輸」、「春/張る/貼る」、「空/飽き/秋」、「夏/奈津/捺」という文字列候補を得る。

40

【0049】

604は、文字列マッチング部であり、文字認識結果の文字列と音声認識結果の発音列を文字列に変換したものととのマッチングを行う。このマッチング処理によって、図13に示す例では、音声認識結果の複数の文字列候補から「冬」、「春」、「秋」、「夏」が選択され、文字認識結果の文字列と対応付けられる。また、605は、静止画音声対応情報出力部であり、文字列マッチング部604によるマッチング結果を図8に示すように静止画音声対応情報として出力する。

【0050】

50

以上の説明から明らかなように、本実施例によれば、文字認識結果の文字列と音声認識結果の文字列が直接比較できないような場合においても、発音列でのマッチングを行うことによって、静止画像と音声の対応付けを行うことが可能となる。

【実施例 4】

【0051】

上述した実施例では、文字認識結果及び音声認識結果はいずれも1つの認識結果のみであり、また、静止画像と音声との対応付け処理では、認識結果の文字列又は発音列のみを用いて対応付けを行っていたが、認識結果に尤度や確率等のスコア情報を保持した複数候補を出力し、このスコア付きの複数候補を用いて文字認識結果と音声認識結果を対応付けることも可能である。

10

【0052】

ここで、 N 個の文字領域 I_1, \dots, I_N に対して、 M 個の音声区間 S_1, \dots, S_M の1つと対応付けを行った結果を C_1, \dots, C_N (但し、 $C_n = (I_n, S_m), 1 \leq n \leq N, 1 \leq m \leq M$) とするとき、 C_n は、

$$C_n = \arg \max (P I_{n i}, P S_{m j}, R I_{n i}, R S_{m j})$$

によって求めることができる。

【0053】

ここで、 $P I_{n i}$ は文字領域 I_n の i 番目の文字認識結果候補のスコア ($1 \leq i \leq K$, 但し、 K は文字認識結果の候補数。)、 $P S_{m j}$ は音声区間 S_m の j 番目の音声認識結果候補のスコア ($1 \leq j \leq L$, 但し、 L は音声認識結果の候補数) である。また、 I_n の第 i 位の文字認識結果の文字列 (又は、発音列) を $R I_{n i}$ 、 S_m の第 j 位の音声認識結果の文字列 (又は、発音列) を $R S_{m j}$ とするとき、 $R I_{n i}, R S_{m j}$ は、 $R I_{n i} = R S_{m j}$ の場合は $R I_{n i}, R S_{m j} = 1$ 、それ以外の場合は $R I_{n i}, R S_{m j} = 0$ という関数で与えられる。さらに、 $\arg \max$ は、 $P I_{n i}, P S_{m j}, R I_{n i}, R S_{m j}$ を最大にする i, m, j の組を求める演算を表し、これを求めることによって、 I_n に対する S_m 、すなわち C_n を決めることができる。

20

【0054】

以下、図14、15、16を用いて、対応付けの具体例について説明する。

【0055】

図14は、第4の実施例における文字認識結果 (a) と音声認識結果 (b) のスコア情報 (尤度や確率等で表された認識結果) を保持した複数候補の例を示す図である。図14に示す例では、 $N = 4, M = 4, K = 3, L = 3$ である。ここで、第1の実施例で説明したように、文字認識結果と音声認識結果の文字列を直接比較することによって、静止画と音声の対応付けを行うことにする。例えば、図14に示すように、 I_1 は「春」、 S_1 は「冬」、 $P I_{11} = 0.7, P S_{43} = 0.1, R I_{13}$ は「空」、 $R S_{32}$ は「足」等となる。

30

【0056】

このとき、 $n = 1$ 、すなわち「春」、「香」、「空」と文字認識された文字領域に対する音声区間は、 $i = 1, m = 2, j = 1$ の場合、 $P I_{11} = 0.7, P S_{21} = 0.7, R I_{11}$ は「春」、 $R S_{32}$ は「春」で $R I_{11}, R S_{21} = 1$ となり、上記 $\arg \max$ の中が最大 $0.49 (= 0.7 \times 0.7 \times 1)$ となる。尚、その他の場合は、いずれも $R I_{n i}, R S_{m j} = 0$ となるため、 $\arg \max$ の中は0となる。よって、 $C_1 = (I_1, S_2)$ と決定される。同様の計算を行うことによって、 $C_2 = (I_2, S_3), C_3 = (I_3, S_4), C_4 = (I_4, S_1)$ と対応付けがなされる。

40

【0057】

次に、第2の実施例で説明したように、文字認識結果を発音列に変換し、これと音声認識結果の発音列を比較することによって静止画と音声の対応付けを行う際に、スコア付きの複数候補を用いる例について説明する。

【0058】

図15は、第4の実施例における文字認識結果を発音列に変換した結果 (a) と音声認

50

識結果から得られる発音列 (b) のスコア情報を保持した複数候補の例を示す図である。この場合、文字認識結果のスコア情報をそのまま発音列のスコア情報とする。また、1つの文字認識結果から複数の発音列が得られる場合には、それぞれの発音列に対して同じスコア情報を用いる。

【 0 0 5 9 】

例えば、 $n = 1$ の場合は、 $i = 1$ で「ハル」と「シュン」の2通り、 $i = 2$ で「カ」と「コウ」の2通り、 $i = 3$ で「ソラ」と「アキ」と「クウ」の3通りの発音列に対して、図14に示した例の場合と同様の計算を行う。この結果、例えば、 $n = 1$, $i = 1$ の「ハル」と $m = 2$, $j = 1$ の「ハル」の $\arg \max$ の中は $0.49 (= 0.7 \times 0.7 \times 1)$ であり、 $n = 1$, $i = 3$ の「アキ」と $m = 3$, $j = 1$ の「アキ」の $\arg \max$ の中は $0.06 (= 0.1 \times 0.6 \times 1)$ であり、 $C1 = (I1, S2)$ と対応付けられる。また、 $n = 4$, $i = 2$ の「フユ」と $m = 1$, $j = 1$ の「フユ」は $0.15 (= 0.3 \times 0.5 \times 1)$ であり、 $n = 4$, $i = 3$ の「ツ」と $m = 4$, $j = 2$ の「ツ」は $0.02 (= 0.2 \times 0.1 \times 1)$ であり、 $C4 = (I4, S1)$ と対応付けられる。同様に、 $C2 = (I2, S3)$ 、 $C3 = (I3, S4)$ と対応付けがなされる。

10

【 0 0 6 0 】

次に、第3の実施例で説明したように、音声認識結果を文字列に変換し、これと文字認識結果の文字列を比較することによって静止画と音声の対応付けを行う際に、スコア付きの複数候補を用いる例について説明する。

【 0 0 6 1 】

図16は、第4の実施例における文字認識結果から得られる文字列 (a) と音声認識結果を文字列に変換した結果 (b) のスコア情報を保持した複数候補の例を示す図である。この場合も、図15で示した発音列の対応付けと同様であり、例えば、 $n = 1$, $i = 1$ の「春」と $m = 2$, $j = 1$ の「春」は $0.49 (= 0.7 \times 0.7 \times 1)$ であり、 $n = 1$, $i = 3$ の「空」と $m = 3$, $j = 1$ の「空」は $0.06 (= 0.1 \times 0.6 \times 1)$ であり、 $C1 = (I1, S2)$ と対応付けられる。

20

【 0 0 6 2 】

尚、上述したように本実施例では、 は完全に一致する場合に1、一致しない場合は0という2値の値のいずれかをとり関数を用いていたが、これに限らず、例えば一致の度合いに応じた値とする等、別の定義でもよい。また、文字認識結果のスコアと音声認識結果のスコアは同等に扱っているが、例えば、文字認識のスコアを音声認識のスコアよりも重視する等、これらのスコアに重みをつけてもよい。

30

【 0 0 6 3 】

以上の説明から明らかなように、本実施例によれば、文字認識結果と音声認識結果をスコア付きで複数候補出力することで、1位の候補に正解の認識結果が含まれない場合でも、より正確に静止画像と音声の対応付けを行うことが可能となる。

【 実施例 5 】

【 0 0 6 4 】

上述した第2～第4の実施例では、静止画音声対応部205において、発音列又は文字列のどちらかに変換された結果に基づいて静止画像と音声の対応付けを行っていたが、これらの両方を用いて対応付けを行うこともできる。すなわち、文字認識結果を読みに変換した文字認識結果発音列と音声認識結果として得られる音声認識結果発音列のマッチングと、文字認識結果として得られる文字認識結果文字列と音声認識結果を文字列に変換した音声認識結果文字列のマッチングの両方を用いる。これは、図10と図12のそれぞれで示されるモジュール構成を併用することによって実現することができる。

40

【 実施例 6 】

【 0 0 6 5 】

上述した実施例では、文字認識に関する処理と音声認識に関する処理は、それぞれ独立に行われていたが、文字認識の結果を音声認識で利用することも可能である。この際、以下に説明するように様々な利用の仕方が考えられる。

50

【0066】

まず、文字認識結果を音声認識結果情報出力部で利用する場合について説明する。図17は、本発明の第6の実施例における静止画像・音声認識装置のモジュール構成を示すブロック図である。図17において、文字認識部701は文字認識部202と、また、音声分析部702、探索部703、音声認識用音響モデル704及び音声認識用言語モデル705は、それぞれ音声分析部401、探索部402、音声認識用音響モデル404、音声認識部言語モデル405と、さらに静止画音声対応部707は静止画音声対応部205と同じであるため説明は省略する。

【0067】

706は音声認識結果情報出力部であり、探索部703の探索結果に加えて、文字認識部701の文字認識で得られる結果も利用する。例えば、図14に示す場合、図14(b)に示される音声認識結果に対して、図14(a)の結果に含まれない「古」、「露」、「樽」、「白」、「足」、「薪」、「松」、「津」の8種類の文字列は音声認識結果候補としない。この結果、これらの8種類の文字列に対しては、第4の実施例で説明した計算を行う必要がなくなり、処理の効率化が図れる。

10

【0068】

次に、文字認識結果を音声認識の探索部で利用する場合について説明する。図18は、本発明の第6の実施例における文字認識結果を音声認識に利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【0069】

図18において、文字認識部801は文字認識部202と、また音声分析部802及び音声認識用音響モデル804から音声認識結果情報出力部806はそれぞれ音声分析部401及び音声認識用音響モデル404、音声認識用言語モデル405、音声認識結果情報出力部403と、さらに静止画音声対応部807は静止画音声対応部205と同じであるため説明は省略する。

20

【0070】

探索部803は、音声認識用音響モデル804と音声認識用言語モデル805の2つのモデルを用いて音声認識を行う際に、文字認識部801で得られる結果を利用する。例えば、図14(a)に示された結果が文字認識の結果として得られた場合、探索部803は、これらの12種類の文字列(単語)のみを用いた探索処理を行う。すなわち、探索部803は、音声認識用言語モデル805に含まれる音声認識対象語としてこれらの12種類のみを用いて音声認識を行う。この結果、探索部803の計算が大幅に低減され、文字認識の結果候補に正解が含まれている場合、文字認識と音声認識を独立に行うものと比較して、音声認識の性能も一般に向上させることができる。

30

【0071】

次に、文字認識結果を発音列に変換し、これを音声認識結果情報出力部で利用する場合について説明する。図19は、第6の実施例における文字認識結果を発音列に変換して利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【0072】

図19において、文字認識部901は文字認識部202と、また文字認識結果発音列変換部902は文字認識結果発音列変換部501と、さらに音声分析部903から音声認識用言語モデル906はそれぞれ音声分析部401、探索部402、音声認識用音響モデル404、音声認識用言語モデル405と、さらにまた静止画音声対応部908は静止画音声対応部205と同じであるため説明は省略する。尚、図19では、文字認識結果発音列変換部902の処理を行う際に必要な発音変換辞書502は省略している。

40

【0073】

図19において、音声認識結果情報出力部907は、探索部903の結果に加えて、文字認識部901の文字認識結果を発音列に変換した結果も利用する。例えば、図15に示す例の場合、図15(b)に示される音声認識結果に対して、図15(a)の結果に含まれない「フル」、「ツユ」、「タル」、「ハク」、「アシ」、「マキ」、「マツ」の7種

50

類の発音列は音声認識結果候補としない。この結果、これらの7種類の文字列に対しては、第4の実施例で説明した計算を行う必要がなくなる。

【0074】

次に、文字認識結果から得られる発音列を音声認識の探索部1004で利用する場合について説明する。図20は、第6の実施例における文字認識結果を発音列に変換して探索部で利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【0075】

図20において、文字認識部1001は文字認識部202と、また文字列認識結果発音列変換部1002は文字列認識結果発音列変換部501と、さらに音声分析部1003及び音声認識用音響モデル1005から音声認識結果情報出力部1007はそれぞれ音声分析部401及び音声認識用音響モデル404、音声認識用言語モデル405、音声認識結果情報出力部403と、さらにまた静止画音声対応部1008は静止画音声対応部205と同じであるため説明は省略する。尚、図20では、文字認識結果発音列変換部1002の処理を行う際に必要な発音変換辞書502は省略している。

10

【0076】

図20において、探索部1004は、音声認識用音響モデル1005と音声認識用言語モデル1006の2つのモデルを用いて音声認識を行う際に、文字認識結果発音列変換部1002で文字認識結果を発音列に変換した結果も利用する。例えば、図15(a)に示された結果が文字認識の結果から得られる発音列であるとき、探索部1004は、これらの25種類の発音列のみを用いた探索処理を行う、すなわち、探索部1004は、音声認識用言語モデル1006に含まれる音声認識対象語として、これらの25種類のみを用いて音声認識を行う。

20

【0077】

この結果、探索部1004の計算が大幅に低減され、文字認識の結果から得られる発音列候補に正解が含まれている場合、文字認識と音声認識を独立に行うものと比較して、音声認識の性能も一般に向上させることができる。

【0078】

次に、文字認識結果から得られる文字列を音声認識結果から文字列に変換する際に利用する静止画音声対応処理について説明する。

【0079】

図21は、第6の実施例における文字認識結果の文字列を音声認識結果を文字列に変換する際に利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。図21において、文字認識結果文字列抽出部1101は文字認識結果文字列抽出部601と、文字変換辞書1103から静止画音声対応情報出力部1105はそれぞれ文字変換辞書603から静止画音声対応情報出力部605と同じであるため説明は省略する。

30

【0080】

図21において、1102は音声認識結果文字列変換部であり、音声認識結果を文字列に変換する際に、文字認識結果文字列抽出部1101による文字認識結果から抽出される文字列も利用する。例えば、図16(a)に示された結果が文字認識の結果から抽出される文字列であるとき、音声認識結果文字列変換部1102の音声認識結果を文字列に変換する際に、これらの16種類の文字列に変換しうる音声認識結果のみを文字列変換候補として選択する。

40

【0081】

以上の説明から明らかなように、本実施例によれば、文字認識で得られる結果を音声認識において利用することで、計算量の低減や、音声認識性能を向上させることが可能となる。

【実施例7】

【0082】

前述した実施例における、文字認識の結果を音声認識の探索部で利用する処理は、文字認識の結果の文字列をそのまま用いることによって行われていたが、一般に文字認識の結

50

果通りに音声が発声されるとは限らないため、文字認識の結果から音声として発声されると予想される重要語を抽出し、これを音声認識の探索部で利用することが好ましい。

【0083】

図22は、第7の実施例における文字認識結果から重要語を抽出して探索部で利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。図22において、文字認識部1201は文字認識部202と、また音声分析部1203及び音声認識用音響モデル1205から音声認識結果情報出力部1207はそれぞれ音声分析部401及び音声認識用音響モデル404、音声認識用言語モデル405、音声認識結果情報出力部403と、さらに静止画音声対応部1208は静止画音声対応部205と同じであるため説明は省略する。

10

【0084】

図22において、1202は重要語抽出部であり、文字認識の結果から重要語を抽出する。例えば、文字認識の結果が、「提案法は、統計的言語モデルを用いたアプローチです。」といった文字列であり、重要語の抽出方法が文字列を形態素解析し、この結果から自立語を抽出するものとすると、この結果から、「提案法」、「統計的」、「言語」、「モデル」、「アプローチ」の5単語が重要語として抽出される。

【0085】

また、1204は探索部であり、音声認識用音響モデル1205と音声認識用言語モデル1206の2つのモデルを用いて音声認識を行う際に、重要語抽出部1202で抽出された単語を利用する。具体的には、上述した5単語をキーワードとしたキーワードスポットティングによる音声認識を行う、または、大語彙連続音声認識の場合には、音声認識結果からこれらの5単語が含まれる部分を抽出する、若しくは、上記5単語に関する音声認識用言語モデルの確率値を増加させて音声認識を行う。尚、重要語の抽出規則は、ここでは自立語の抽出としたが、これに限らず他の規則や手法を適用してもよい。また、1209は重要語を抽出ための規則やデータ(単語辞書)である。

20

【0086】

以上の説明から明らかなように、本実施例によれば、文字認識の結果通りの音声でない場合においても、静止画像と音声の対応付けを好適行うことが可能となる。

【実施例8】

【0087】

一般に、静止画像に含まれる文字に関する情報は、単なる文字列のみではなく、フォントサイズ、文字種、色、斜体やアンダーライン等スタイルや文字飾りに関する情報も含まれているため、これらのフォント情報を抽出し、これを音声認識で利用することによって、より正確に静止画像と音声を対応付けることができる。

30

【0088】

そこで、例えば、図23に示されるような静止画像からフォント情報を抽出し、これを音声認識で利用する実施例を考える。図23は、第8の実施例における種々のフォント情報をもった静止画像の例である。また、図24は、第8の実施例における文字領域からフォント情報を抽出して文字認識結果情報として出力する静止画像・音声認識装置のモジュール構成を示すブロック図である。

40

【0089】

図24において、1301はフォント情報抽出部であり、文字領域に対して、フォントサイズ、文字種、色、斜体やアンダーラインの有無等のフォント情報を抽出する。また、他のモジュールは、図4に示す例と同じであるため省略する。

【0090】

図25は、図23に示す静止画像からの文字認識結果と各文字領域のフォント情報を示す図である。次に、図25に示されるフォント情報を音声認識で利用する。尚、このときのモジュール構成は、図18に示す装置と同様である。但し、図18の文字認識部801は、図24に示した構成となる点で異なる。

【0091】

50

ここで、フォント情報の音声認識での利用の仕方は様々であるが、例えば、フォントサイズが大きい文字列や斜体やアンダーラインが施されている文字列は、キーワードスポットティングの対象とする、又は統計的言語モデルの確率値を増加させて音声認識を行う。他にも、黒以外の色については、色の情報を音声認識の対象語彙に追加するといったことができる。

【0092】

以上の説明から明らかなように、本実施例によれば、静止画像に含まれる文字領域のフォント情報を音声認識で利用することによって、より正確に静止画像と音声を対応付けることが可能となる。

【実施例9】

【0093】

上述した第6の実施例では、文字認識の結果を音声認識で利用する場合について説明したが、これとは逆に、音声認識の結果を文字認識で利用することもできる。この際、以下に説明するように様々な利用の仕方が考えられる。

【0094】

まず、音声認識結果を文字認識結果情報出力部で利用する場合について説明する。図26は、第9の実施例における文字認識結果情報出力部の細部モジュール構成を示すブロック図である。図26において、音声認識部1401は音声認識部204と、また前処理部1402から文字認識用言語モデル1406はそれぞれ前処理部301、特徴抽出部302、識別部303、文字認識用テンプレート305、文字認識用言語モデル306と、さらに静止画音声対応部1408は静止画音声対応部205と同じであるため説明は省略する。

【0095】

1407は文字認識結果情報出力部であり、識別部1404の識別結果に加えて、音声認識部1401の音声認識で得られる結果も利用する。例えば、図14の場合、図14(a)に示される文字認識結果に対して、図14(b)の結果に含まれない「香」、「空」、「科」、「和」、「新」、「厚」、「各」、「尽」の8種類の文字列は文字認識結果候補としない。この結果、これらの8種類の文字列に対しては、第4の実施例で説明した計算を行う必要がなくなる。

【0096】

次に、音声認識結果を文字認識の識別部で利用する場合について説明する。図27は、第9の実施例における静止画像・音声認識装置のモジュール構成を示すブロック図である。図27において、音声認識部1501は音声認識部204と、また前処理部1502、特徴抽出部1503及び文字認識用テンプレート1505から文字認識結果情報出力部1507はそれぞれ前処理部301、特徴抽出部302及び文字認識用テンプレート305、文字認識用言語モデル306、文字認識結果情報出力部304と、さらに静止画音声対応部1508は静止画音声対応部205と同じであるため説明は省略する。

【0097】

識別部1504は、文字認識用テンプレート1505と文字認識用言語モデル1506の2つのモデルを用いて文字認識を行う際に、音声認識部1501の音声認識で得られる結果を利用する。例えば、図14(b)に示された結果が音声認識の結果として得られた場合、識別部1504は、これらの16種類の文字列のみを用いた識別処理を行う。すなわち、識別部1504に含まれる文字認識対象語としてこれらの16種類のみを用いて文字認識を行う。この結果、識別部の計算が大幅に低減され、音声認識の結果候補に正解が含まれている場合、文字認識と音声認識を独立に行うものと比較して、文字認識の性能も一般に向上させることができる。

【0098】

次に、音声認識結果を文字列に変換し、これを文字認識結果情報出力部で利用する場合について説明する。図28は、第9の実施例における音声認識結果を文字列に変換して利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。図28におい

10

20

30

40

50

て、音声認識部 1601 は音声認識部 204 と、また音声認識結果文字列変換部 1602 は音声認識結果文字列変換部 602 と、さらに前処理部 1603 から文字認識用言語モデル 1607 はそれぞれ図 4 に示す前処理部 301、特徴抽出部 302、識別部 303、文字認識用テンプレート 305、文字認識用言語モデル 306 と、さらにまた静止画音声対応部 1609 は図 3 に示す静止画音声対応部 205 と同じであるため説明は省略する。尚、図 28 では、音声認識結果文字列変換部 1602 の処理を行う際に必要な文字変換辞書 602 は省略している。

【0099】

図 28 において、1608 は文字認識結果情報出力部であり、識別部 1605 の識別結果に加えて、音声認識部 1602 の音声認識結果を文字列に変換した結果も利用する。例えば、図 16 に示す例の場合、図 16 (a) に示される文字認識結果に対して、図 16 (b) の結果に含まれない「香」、「科」、「和」、「真」、「厚」、「各」、「尽」の 7 種類の文字列は文字認識結果候補としない。この結果、これらの 7 種類の文字列に対しては、第 4 の実施例で説明した計算を行う必要がなくなる。

【0100】

次に、音声認識結果から得られる文字列を文字認識の識別部で利用する場合について説明する。図 29 は、第 9 の実施例における音声認識結果から得られる文字列を文字認識で利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。図 29 において、音声認識部 1701 は図 3 に示す音声認識部 204 と、また音声認識結果文字列変換部 1702 は図 12 に示す音声認識結果文字列変換部 602 と、さらに前処理部 1703、特徴抽出部 1704 及び文字認識用モデル 1706 から文字認識結果情報出力部 1708 はそれぞれ図 4 に示す前処理部 301、特徴抽出部 302 及び文字認識用テンプレート 305、文字認識用言語モデル 306、文字認識結果情報出力部 304 と、さらにまた静止画音声対応部 1709 は図 3 に示す静止画音声対応部 205 と同じであるため説明は省略する。尚、図 29 では、音声認識結果文字列変換部 1702 の処理を行う際に必要な図 12 に示す文字変換辞書 603 は省略している。

【0101】

識別部 1705 は、文字認識用モデル 1706 と文字認識用言語モデル 1707 の 2 つのモデルを用いて文字認識を行う際に、音声認識結果文字列変換部 1702 の音声認識結果を文字列に変換した結果も利用する。例えば、図 16 (b) に示された結果が音声認識の結果から得られる文字列であるとき、識別部 1705 は、これらの 32 種類の文字列のみを用いた識別処理を行う。すなわち、識別部 1705 は、文字認識用モデル 1706 や文字認識用言語モデル 1707 に含まれる文字認識対象語としてこれらの 32 種類のみを用いて文字認識を行う。

【0102】

この結果、識別部の計算が大幅に低減され、音声認識の結果から得られる文字列候補に正解が含まれている場合、文字認識と音声認識を独立に行うものと比較して、文字認識の性能も一般に向上させることができる。

【0103】

次に、音声認識結果から得られる発音列を文字認識結果の発音列の変換で利用する際の静止画音声対応手段について説明する。図 30 は、第 9 の実施例における音声認識結果から得られる発音列を文字認識結果の発音列の変換に利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。図 30 において、音声認識結果発音列抽出部 1801 は図 10 に示す音声認識結果発音列抽出部 503 と、発音変換辞書 1803 から静止画音声対応情報出力部 1805 はそれぞれ発音変換辞書 502、発音列マッチング部 504、静止画音声対応情報出力部 505 と同じであるため説明は省略する。

【0104】

図 30 において、1802 は文字認識結果発音列変換部であり、文字認識結果を発音列に変換する際に、音声認識結果発音列抽出部 1801 の音声認識結果から抽出される発音列も利用する。例えば、図 15 (b) に示された結果が音声認識の結果から抽出される発

10

20

30

40

50

音列であるとき、文字認識結果発音列変換部 1802 の文字認識結果を発音列に変換する際に、これらの 16 種類の発音列に変換しうる文字認識結果のみを発音列変換候補として選択する。

【0105】

以上の説明から明らかなように、本実施例によれば、音声認識で得られる結果を文字認識において利用することで、計算量の低減や、文字認識性能を向上させることが可能となる。

【実施例 10】

【0106】

前述した実施例で示した図 2 (a) や図 23 に示す静止画像の例は、非常に単純なものであったが、本発明をより複雑な静止画像に対して音声との対応付けを行うためには、静止画像から単純に文字領域を見つけるだけでなく、静止画像の構造を正確に捉える必要がある。すなわち、まず大局的に静止画像を複数の領域に分割し、次に個々の分割静止画像領域に対して文字領域を抽出し、文字認識を行う。

【0107】

図 31 は、より複雑な静止画像 (a) とそれに関連付けられた音声 (b) の一例を示す図である。また、図 32 は、図 31 に示すような複雑な静止画像を分割する機能を有する静止画像・音声認識装置のモジュール構成を示すブロック図である。また、図 62 は、第 10 の実施例に係る静止画像・音声認識装置の処理動作を説明するためのフローチャートである。

【0108】

図 32 に示すように、静止画像分割部 1901 は、1 つの静止画像を複数の静止画像領域に分割する (ステップ S11)。尚、この分割方法としては、既存の技術を利用すればよい。また、文字検出部 1902 から静止画音声対応部 1906 は、図 3 に示す文字検出部 201 から静止画音声対応部 205 と同様であるため説明は省略する。但し、図 3 の文字検出部 201 は静止画像全体が入力であったが、本実施例では静止画像分割部 1901 で分割された個々の静止画像が入力となる点で異なる。

【0109】

図 33 は、図 31 (a) の静止画像を静止画像分割部 1901 によって 5 つの領域に分割された結果を示す図である。また、図 34 は、図 33 に示す各分割領域に対して領域 ID を付与して図 32 (a) の座標系における座標情報を表した図である。尚、図 35 は、図 34 と図 32 (a) の対応関係について示す図である。

【0110】

次に、各分割画像に対して文字検出部 1902 によって文字領域の検出処理を施す (ステップ S12)。図 36 は、文字検出部 1902 による文字領域の検出結果例を示す図である。さらに、図 36 に示す文字領域に対して文字認識部 1903 による文字認識を行うとともに (ステップ S13)、図 31 (b) の音声に対して音声検出部 1904 による音声検出 (ステップ S14) 及び音声認識部 1905 による音声認識を行う (ステップ S15)。尚、ステップ S11 ~ S13 の文字認識、ステップ S14 ~ S15 の音声認識は、両認識処理を同時にしてもよく、どちらの処理を先にしてもよい。

【0111】

図 37 は、文字認識及び音声認識の結果得られる文字認識結果情報 (a) 及び音声認識結果情報を (b) 示す図である。尚、文字認識結果情報の座標情報は、図 38 に示すような矩形領域として 2 点の座標で表している。すなわち、図 38 は、図 36 に示す文字領域の検出結果に文字認識結果情報を対応させた結果を示す図である。そして、静止画音声対応部 1906 で、前述した実施例と同様の方法によって、図 37 (a) の文字認識結果と図 37 (b) の音声認識結果を対応付け、静止画音声対応部 1906 を得る (ステップ S16)。

【0113】

以上の説明から明らかなように、本実施例によれば、静止画像領域を大局的に分割する

10

20

30

40

50

ことにより、複雑な静止画像に対しても文字領域と音声との対応付けを行うことが可能となる。

【実施例 1 1】

【0 1 1 4】

前述した実施例で示した図 2 (b) や図 3 1 (b) の音声の例は、各発声区間の間には十分な無音区間が含まれ、また、発声される内容は静止画像の文字領域のいずれかと全く同じであるという非常に単純なものであった。しかしながら、実際の音声は、文字領域の内容と同じ発声をするとは限らず、さらに、ある文字領域の内容に関する発声は全くされていなかったり、どの文字領域とも関係のない発声が含まれていたりする場合がある。また、複数の文字領域に関する音声が十分な無音区間なしに連続的に発声される場合や、雑音や音楽といった音声以外のものも含まれている場合もある。よって、本発明をより一般的な音声であっても静止画像との対応付けが行えるようにするためには、音声区間の正確な抽出、音声認識結果と文字認識結果の柔軟なマッチングを行う必要がある。

10

【0 1 1 5】

そこで、音声区間の正確な抽出に関しては、まず、雑音や音楽といった音声以外のものが入力音声に含まれている場合の対処について説明する。このような音声が入力される場合には、最初に音声を複数のセグメントに大局的に分割し、次に個々の音声セグメントに対して、音声 / 非音声の判定や音声区間の検出を行うことが望ましい。

【0 1 1 6】

図 4 0 は、第 1 1 の実施例に係る静止画像・音声認識装置のモジュール構成を示すブロック図である。図 4 0 において、文字検出部 2 0 0 1、文字認識部 2 0 0 2 及び音声認識部 2 0 0 5、静止画音声対応部 2 0 0 6 は、それぞれ図 3 に示す文字検出部 2 0 1、文字認識部 2 0 2、音声認識部 2 0 4、静止画音声対応部 2 0 5 と同じであるため説明は省略する。

20

【0 1 1 7】

図 4 0 において、2 0 0 3 は音声分割部であり、音声を大局的に捉え、複数のセグメントに分割する。具体的には、音声信号をフレーム処理し、スペクトル情報を求め、複数フレーム間のスペクトルの類似性から着目しているフレームをセグメント境界とするか否かを判定する等の方法によって分割することができる。

【0 1 1 8】

次に、音声検出部 2 0 0 4 では、音声分割部 2 0 0 3 で分割された各セグメントに音声が含まれるか否かを判定し、音声が含まれる場合には音声区間を検出する。具体的には、音声及び非音声のそれぞれに対して事前に G M M (Gaussian Mixture Model) を作成し、入力音声をフレーム処理することによって得られるスペクトル情報とこれらの G M M を用いて、当該セグメントに音声が含まれるか否かを判定する。そして、音声が含まれていないと判定された場合は音声認識の対象とせず、音声が含まれていると判定された場合は、2 0 0 4 における次の処理として音声区間を検出し、検出された音声区間を 2 0 0 5 の音声認識部に入力する。

30

【0 1 1 9】

ここで、セグメント数は、セグメント間又はセグメント境界における音声スペクトルに関する尤度基準を用いて音声から決定する方法が考えられるが、これに限らず、静止画像分割、文字領域、文字認識結果によって得られる情報を用いて決定することもできる。具体的には、静止画像分割および文字領域の情報としては、分割数又は領域数に応じてセグメント数を変更する。文字認識結果の情報としては、文字認識結果全体の確からしさが高い場合にはセグメント数を増やすといった方法である。

40

【0 1 2 0】

次に、文字領域の内容と同じ発声されていない場合、一部の文字領域の内容に関する発声は全くされていない場合、どの文字領域とも関係のない発声が含まれていたりする場合、複数の文字領域に関する音声が十分な無音区間なしに連続的に発声される場合について説明する。

50

【 0 1 2 1 】

図 3 9 は、図 3 1 (a) に示す静止画像に関連する音声を説明するための図である。この例では、図 3 6 の文字領域の内容と同じ発声になされておらず、また、図 3 9 の 3 番目の音声区間である「これまでの研究では、...」の部分は静止画像のどの文字領域とも関係のない発声であるとする。さらに、図 3 9 に示されるように、2 番目から 4 番目の音声には十分な無音区間が存在しないものとする。

【 0 1 2 2 】

図 3 9 に示すような発声に対しては、音声分割部 2 0 0 3 又は音声検出部 2 0 0 4 が、正確に静止画像の文字領域に対応した音声分割又は音声区間を検出することは困難である。そこで、音声検出部 2 0 0 4 で検出された音声区間に対して音声認識部 2 0 0 5 で音声認識を行い、音声認識の結果から音声検出部 2 0 0 4 で決定された音声区間を必要に応じてさらに分割するようにする。

10

【 0 1 2 3 】

具体的には、無音区間が十分に存在しない音声に対しては、音声認識部 2 0 0 5 による音声認識として大語彙連続音声認識に基づく方法を用いれば、句点を推定することによって文の区切りが分かるため、図 4 1 に示すように、この情報を用いて音声区間を分割することができる。ここで、図 4 1 は、図 3 1 の例に対する文字認識結果情報と音声認識結果情報の一例を示す図である。また、文字領域の内容に関する発声がない場合、又はどの文字領域とも関係のない発声になされている音声に対しては、音声認識結果と文字認識の結果をそれぞれ部分マッチングすることによって対応付けを行うことが可能である。

20

【 0 1 2 4 】

また、第 7 の実施例で説明したように、文字認識の結果から重要語を検出すれば、この重要語をキーワードとしたワードスポッティングに基づく方法を音声認識部 2 0 0 5 による音声認識とすれば、より直接的に文字認識の結果と音声認識の結果を対応付けることが可能となる。図 4 2 は、重要語抽出によるワードスポッティングを用いた場合の音声認識結果情報の一例を示す図である。図 4 2 に示す例では、文字認識結果から重要語として抽出された「音声認識」、「文字認識」、「統計的言語モデル」、「目的」等の言葉を音声認識のワードスポッティングとしている。なお、図 4 2 における「*」は、これらのキーワード以外の音声区間を表し、また、「NO_RESULTS」は、この音声区間に対してはどのキーワードもマッチングしなかったことを表している。このワードスポッティング結果と文字認識結果から得られる重要語をマッチングさせることによって、文字領域と音声の対応付けを行うことができる。

30

【 0 1 2 5 】

以上の説明から明らかなように、本実施例によれば、音声に雑音や音楽といった音声以外のものが含まれている場合や、無音区間が十分に存在しない場合、文字領域の内容に関する発声がない場合、どの文字領域とも関係のない発声になされている場合の音声であっても文字領域と音声との対応付けを行うことが可能となる。

【 実施例 1 2 】

【 0 1 2 6 】

上記第 1 0 の実施例では、複雑な静止画像に対しても文字領域と音声との対応付けを行えるようにするために、静止画像領域を大局的に分割する方法について説明した。本実施例では、この静止画像分割処理を分割数の異なる分割静止画像を階層的な構造として得ることによって、より柔軟な対応付けを行うことができることを説明する。

40

【 0 1 2 7 】

図 4 3 は、図 3 3 で示した静止画像の分割をさらに行った場合の分割結果 (a) (一点破線)、(a) をさらに分割した場合の結果 (b) (二点破線) を示す図である。尚、分割数の増減は、分割するか否かの基準 (例えば、尤度基準に対する閾値) を変化させることによって制御することができる。ここで、図 4 3 (a) は図 3 3 の結果を元に分割されており、また、図 4 3 (b) は図 4 3 (a) の結果を元に分割されているため、分割は階層的に行われている。

50

【 0 1 2 8 】

図 4 4 は、階層的な静止画像の分割を木構造で表現した例を示す図である。図 4 4 において、黒丸はルートノードであって静止画像全体を表している。また、I 1 ~ I 5 の 5 個のノードは、図 3 3 の分割領域に対する静止画像であり、I 1 は、図 3 3 の分割領域の「音声認識・文字認識のための統計的言語モデルの利用」を含む画像領域、I 2 は、「目的」、「音声認識性能の向上」、「文字認識性能の向上」を含む画像領域、I 3 は、「提案法」、「統計的言語モデルの利用」、「単語間、文字間の...可能となる」を含む画像領域、I 4 は、「実験結果」、「認識率」、「音声認識」、「文字認識」を含む画像領域、I 5 は、「結論」、「統計的言語モデルは、...分かった。」を含む画像領域である。

【 0 1 2 9 】

また、次の階層の I 2 1 ~ I 5 2 の 1 1 個のノードは、図 4 3 (a) の分割領域に対する静止画像であり、I 2 1 は、「目的」を含む画像領域、I 2 2 は、「音声認識性能の向上」及び「文字認識性能の向上」を含む画像領域、I 3 1 は、「提案法」を含む画像領域、I 3 2 は、「統計的言語モデルの利用」を含む画像領域、I 3 3 は、下矢印記号を含む画像領域である。尚、図 4 3 (a) の分割時には I 1 の画像領域分割が施されていないため、I 1 のノード分割はない。

【 0 1 3 0 】

同様に、最下階層の I 2 2 1 ~ I 4 3 2 の 4 個のノードは、図 4 3 (b) の分割領域に対する静止画像であり、I 2 2 1 は、「音声認識性能の向上」を含む画像領域、I 2 2 2 は、「文字認識性能の向上」を含む画像領域、I 4 3 1 は、「音声認識」を含む画像領域、I 4 3 2 は、「文字認識」を含む画像領域である。

【 0 1 3 1 】

本実施例では、音声のセグメント分割又は音声区間検出は必ずしも階層的に行う必要はないが、ここでは階層的に行った場合の例を示す。図 4 5 は、階層的に音声分割を行った場合の例を示す図である。また、図 4 6 は、図 4 5 で階層的に分割された音声を木構造で表現した例である。

【 0 1 3 2 】

次に、前述した実施例で説明したいずれかの方法によって、図 4 4 に示す各ノードに対応する画像領域から文字領域を抽出し、文字認識を施すことによって、文字認識結果情報を得ることができる。同様に、前述した実施例で説明したいずれかの方法によって、図 4 6 に示す各ノードに対応する音声セグメントから音声区間を検出し、音声認識を施すことによって、音声認識結果情報を得ることができる。

【 0 1 3 3 】

そして、これらの文字認識結果情報に音声認識結果情報を対応付ける。対応付けの方法は、前述した実施例で説明したいずれの方法を用いればよい。また、木構造の特徴を生かした対応付けの方法として、静止画像の上位ノードから下位ノードの順に対応付けを行い、その際に、上位ノードの対応付けの結果を下位ノードの対応付けにおいて制約として利用することができる。例えば、下位ノードの音声に対応付ける際に、上位ノードで対応付けられた音声区間に含まれる音声を優先的にもしくは限定的に選択する。他にも、上位ノードほど時間的に長い音声区間を優先的に選択し、下位ノードほど時間的に短い音声区間を優先的に選択する等の方法を用いることができる。

【 0 1 3 4 】

図 4 7 は、静止画像の木構造ノードに複数候補の分割音声を対応付けた結果の一例を示す図である。図 4 7 において、「NULL」は音声区間の候補がなかった場合を示しており、特に I 3 3 に対しては、どの音声区間にも対応付けられなかったことを表している。図 4 8 は、図 3 1 の例に対する静止画像と音声の対応結果を用いたアプリケーションの一例を示す図である。図 4 8 に示す例では、静止画像の文字の場所にマウスカーソル(矢印マーク)を持っていくと、この文字に対応した音声データが再生され、スピーカー等の音声出力装置から出力される。

【 0 1 3 5 】

10

20

30

40

50

また、図 4 8 とは逆に、音声を先頭から、或いはマウス等で任意の時間を指定することによって音声を再生し、再生されている音声区間に対応する静止画像に枠を付与して表示することも可能である。図 6 1 は、図 4 3 に示す静止画像と音声との対応結果を用いた別のアプリケーションに基づく表示例を示す図である。この例では、利用者が「そこで、本研究では、統計的...」と音声認識された音声区間 (s 4 から e 4) にマウスカーソル (矢印マーク) を持っていくと、この音声区間に対応した文字領域の座標に文字領域分の外枠が生成・表示される。この結果、出力されている音声と静止画像のどの部分に対応しているかを理解することができる。

【 0 1 3 6 】

本実施例で説明した静止画像を木構造表現すること、また、複数候補の音声を対応付けることは、静止画像と音声の対応付けに誤りを含む場合に特に有効である。図 4 9 は、静止画像の木構造の結果及び複数候補音声を利用する際のユーザインタフェースの一例を示す図である。図 4 9 では、上位候補の音声出力に左矢印キー「 ← 」を、下位候補の音声出力に右矢印キー「 → 」を、静止画像の親ノードへ移動して 1 位候補の音声出力をするために上矢印キー「 ↑ 」を、静止画像の子ノードへ移動して 1 位候補の音声出力をするために下矢印キー「 ↓ 」をそれぞれ割り当てている。そして、利用者がマウス等によって所望の画像領域を選択 (クリック等) すると、選択領域に含まれる画像領域の最下位ノードに対応する文字領域を枠で囲み画面上に表示し、さらに 1 位候補の音声出力をする。この際、音声又は画像領域が所望のものでない場合には、これら 4 つのキーのみを用いて他を選択する簡単な操作によって、他の候補を効率よく探すことが可能となる。

【 実施例 1 3 】

【 0 1 3 7 】

前述した実施例では、文字認識の結果又はこれから抽出された重要語と音声認識の結果をマッチングしていたため、文字認識から得られる文字列と音声認識結果から得られる文字列が少なくとも部分的には同じである必要があった。すなわち、例えば、「題目」という文字認識結果に対して「タイトル」という発声になされたり、「夏」に対して「暑い」という発声になされた場合には対応付けを行うことはできない。そこで、本実施例は、このような場合においても静止画像と音声を対応付けることが可能となる方法を提供する。

【 0 1 3 8 】

図 5 0 は、第 1 3 の実施例における静止画像と音声の例を示す図である。図 5 0 より明らかのように、静止画像に含まれる「春」、「夏」、「秋」、「冬」という単語列は、音声の中に一切含まれていない。この場合、文字認識の結果と、音声認識の結果をそれぞれ抽象化、すなわち概念に変換し、それぞれの概念レベルでマッチングを行うことによって図 5 0 のような場合であっても静止画像と音声を対応付けることが可能となる。

【 0 1 3 9 】

図 5 1 は、第 1 3 の実施例における文字概念変換機能及び音声概念変換機能を有する静止画像・音声認識装置のモジュール構成を示すブロック図である。図 5 1 において、文字検出部 2 1 0 1、文字認識部 2 1 0 2、音声検出部 2 1 0 4、音声認識部 2 1 0 5 は、それぞれ図 3 に示す静止画像・音声認識装置のモジュールと同様であるため説明は省略する。図 5 1 において、2 1 0 3 は文字概念変換部であり、文字認識部 2 1 0 2 で得られる文字認識の結果を予め定められた概念に抽象化する。

【 0 1 4 0 】

また、2 0 1 6 は音声概念変換部であり、音声認識部 2 1 0 5 で得られる音声認識の結果を予め定められた概念に抽象化する。2 1 0 7 は概念対応部であり、文字概念変換部 2 1 0 3 と音声概念変換部 2 1 0 6 で得られる結果に対して概念レベルでマッチングを行う。静止画音声対応部 2 1 0 8 は、概念対応部 2 1 0 7 で対応付けられた概念に対して静止画像と音声を対応付ける。

【 0 1 4 1 】

例えば、\$SPRING、\$SUMMER、\$AUTUMN、\$WINTER という 4 つの概念が定義されており、各概念に含まれる文字列として、\$\$SPRING = { 春、spring、桜、入学式、... }、\$SUMMER =

10

20

30

40

50

{夏、summer、hot、暑、...}、\$AUTUMN={秋、autumn、fall、紅葉、...}、\$WINTER={冬、winter、cold、寒、...}が定義されているとする。図52は、文字概念変換結果と静止画像の座標情報、及び音声概念変換結果と音声の時間情報の一例を示す図である。そこで、図50における静止画像及び音声に対して、図52に示すような関係があるとす。尚、この例の場合は、音声認識として英語が認識できるものを用いているとする。

【0142】

そこで、この結果を概念対応部2107で対応付けることによって、\$SPRING同士、\$SUMMER同士等がそれぞれ対応付けられ、静止画音声対応部2108では、「春」の画像領域に対して「入学式の...」の音声が、「夏」の画像領域に対して「暑くなって...」の音声が、「秋」の画像領域に対して「紅葉狩りに...」の音声が、「冬」の画像領域に対して「Winter is a...」の音声がそれぞれ対応付けられる。

10

【0143】

以上の説明から明らかなように、本実施例によれば、文字列ではなく概念レベルでマッチングを行うことによって、文字認識から得られる文字列と音声認識結果から得られる文字列が全く一致しない場合であっても文字領域と音声との対応付けを好適に行うことが可能となる。

【実施例14】

【0144】

前述した実施例では、静止画像の文字領域の部分に対してのみ音声と対応付けることが可能であり、静止画像中の文字以外の、例えば円や三角形等の図形や、人、車等のオブジェクトに対しては音声を対応付けることはできなかった。そこで、本実施例では、このような場合においても静止画像と音声を対応付けることが可能な方法を提供する。

20

【0145】

図53は、第14の実施例において用いる静止画像とそれに対応付けられる音声の例を示す図である。図53より明らかなように、静止画像には文字列が一切含まれていない。この場合、前述した実施例における文字認識の代わりに、オブジェクト認識を行い、その認識結果と音声認識の結果をマッチングすることによって図53のような場合であっても静止画像と音声を対応付けることが可能となる。

【0146】

図54は、本発明の第14の実施例に係るオブジェクト認識処理機能を有する静止画像・音声処理装置のモジュール構成を示すブロック図である。図54において、音声検出部2203及び音声認識部2204は、図3に示すそれぞれのモジュールと同様であるため説明は省略する。図54において、2201はオブジェクト検出部であり、静止画像からオブジェクト領域を抽出する。また、2202はオブジェクト認識部であり、オブジェクト検出部2201で抽出されたオブジェクトを認識する。尚、オブジェクト検出処理及びオブジェクト認識処理については、既存の技術を用いることができる。

30

【0147】

本実施例では、例えば、円、三角形、長方形、正方形等の図形の形状、棒グラフ、折れ線グラフ、円グラフ等のグラフの形状、およびそれぞれの形状に対する代表的な色の抽出が可能でオブジェクト検出処理及びオブジェクト認識処理が実施できるとする。この場合、図54(a)の静止画像に対して、図55(a)に示されるようなオブジェクト認識結果情報が得られる。

40

【0148】

図55は、オブジェクト認識結果情報の例(a)とオブジェクト認識結果情報から得られる画像領域の例(b)を示す図である。図55に示すように、オブジェクト認識結果として得られる「長方形」、「黒」、「正方形」、「白」といったオブジェクトの形状や色を表す言葉を文字列とし、この文字列と音声認識結果を2205で比較することによって、静止画像と音声を対応付けることができる。この結果、図55(b)で示されるように、静止画像のオブジェクトと音声に対応付けられる。

【0149】

50

以上の説明から明らかなように、本実施例によれば、オブジェクトを検出・認識する機能を備えることによって、静止画像に文字列が含まれない場合であっても音声との対応付けを好適に行うことが可能となる。

【実施例 15】

【0150】

前述した実施例では、静止画像と音声を対応付ける場合に、音声は音声認識を行っていたが、静止画像に人物が含まれ、この人物もしくは人物のクラスが特定でき、さらに、音声は、静止画像の人物もしくは人物クラスに関連している場合には、音声認識を行う代わりに、話者もしくは話者クラスの識別を行うことによって、静止画像と音声を対応付けることが可能となる。

10

【0151】

図56は、第15の実施例において用いる静止画像とそれに対応付けられる音声の例を示す図である。図56より明らかなように、静止画像には文字列が一切含まれていない。また、音声は、高齢者・男性音声で「戦争の頃は...」、成人・男性音声で「僕は来年受験が...」、子供・女性音声で「今日の給食は...」、成人・女性音声で「今夜のドラマは...」という発声になされているものとする。

【0152】

図57は、本発明の第15の実施例に係る人物認識機能及び話者認識機能を有する静止画像・音声認識装置のモジュール構成を示すブロック図である。図57において、2301は人物検出部であり、静止画像から人物に関する画像領域を検出する。2302は人物認識部であり、人物検出部2301で検出された画像領域に対して、人物又は人物クラスの認識を行う。2303は音声検出部であり、音声区間の検出を行う。2304は話者認識部であり、音声検出部2303で検出された音声区間に対して、話者又は話者クラスの認識を行う。

20

【0153】

いま、人物認識部2302が、男性/女性の性別、及び子供/成人/高齢者の年代からなる人物クラスが認識できるとし、話者認識部2304も同様に男性/女性の性別、及び子供/成人/高齢者の年代からなる話者クラスが認識できるものとする。図58は、第15の実施例における人物認識結果情報及び話者認識結果情報の一例を示す図である。ここで、静止画音声対応部2305は、人物クラスと話者クラスのマッチングをとることによって、図59に示すように静止画像と音声の対応付けをすることができる。すなわち、図59は、人物認識結果情報から得られる画像領域を示す図である。

30

【0154】

以上の説明から明らかなように、本実施例によれば、人物又は人物クラスを検出・認識する機能と話者又は話者クラスを認識する機能を備えることによって、静止画像に文字列が含まれない場合に、音声認識を行うことなく音声との対応付けを行うことが可能となる。

【実施例 16】

【0155】

前述した実施例では、静止画像と音声それぞれ1つずつ存在する場合の対応方法について説明したが、本発明の適用はこれだけに限られることなく、例えば静止画像2つと音声3つを対応付ける等、任意の数の静止画像と音声を対応付けるようにしてもよい。

40

【0156】

尚、上述した第1～第15の実施例では静止画像を対象として説明したが、動画像が例えば複数のカテゴリ等に分割されており、各カテゴリの代表的なフレーム(静止画像)に対して本発明を適用することで、所望の動画像を検索することも可能である。

【実施例 17】

【0157】

以上、実施形態例を詳述したが、本発明は、例えば、システム、装置、方法、プログラムもしくは記憶媒体等としての実施態様をとることが可能であり、具体的には、複数の機

50

器から構成されるシステムに適用しても良いし、また、一つの機器からなる装置に適用しても良い。

【0158】

尚、本発明は、前述した実施形態の機能を実現するソフトウェアのプログラム（実施形態では図に示すフローチャートに対応したプログラム）を、システムあるいは装置に直接あるいは遠隔から供給し、そのシステムあるいは装置のコンピュータが該供給されたプログラムコードを読み出して実行することによっても達成される場合を含む。

【0159】

従って、本発明の機能処理をコンピュータで実現するために、該コンピュータにインストールされるプログラムコード自体も本発明を実現するものである。つまり、本発明は、本発明の機能処理を実現するためのコンピュータプログラム自体も含まれる。

10

【0160】

その場合、プログラムの機能を有していれば、オブジェクトコード、インタプリタにより実行されるプログラム、OSに供給するスクリプトデータ等の形態であっても良い。

【0161】

プログラムを供給するための記録媒体としては、例えば、フロッピー（登録商標）ディスク、ハードディスク、光ディスク、光磁気ディスク、MO、CD-ROM、CD-R、CD-RW、磁気テープ、不揮発性のメモリカード、ROM、DVD（DVD-ROM、DVD-R）などがある。

【0162】

その他、プログラムの供給方法としては、クライアントコンピュータのブラウザを用いてインターネットのホームページに接続し、該ホームページから本発明のコンピュータプログラムそのもの、もしくは圧縮され自動インストール機能を含むファイルをハードディスク等の記録媒体にダウンロードすることによっても供給できる。また、本発明のプログラムを構成するプログラムコードを複数のファイルに分割し、それぞれのファイルを異なるホームページからダウンロードすることによっても実現可能である。つまり、本発明の機能処理をコンピュータで実現するためのプログラムファイルを複数のユーザに対してダウンロードさせるWWWサーバも、本発明に含まれるものである。

20

【0163】

また、本発明のプログラムを暗号化してCD-ROM等の記憶媒体に格納してユーザに配布し、所定の条件をクリアしたユーザに対し、インターネットを介してホームページから暗号化を解く鍵情報をダウンロードさせ、その鍵情報を使用することにより暗号化されたプログラムを実行してコンピュータにインストールさせて実現することも可能である。

30

【0164】

また、コンピュータが、読み出したプログラムを実行することによって、前述した実施形態の機能が実現される他、そのプログラムの指示に基づき、コンピュータ上で稼動しているOSなどが、実際の処理の一部または全部を行ない、その処理によっても前述した実施形態の機能が実現され得る。

【0165】

さらに、記録媒体から読み出されたプログラムが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれた後、そのプログラムの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行ない、その処理によっても前述した実施形態の機能が実現される。

40

【図面の簡単な説明】

【0166】

【図1】本発明の第1の実施例に係る画像データと音声データの部分データ同士を対応付ける静止画像・音声処理装置の構成を示すブロック図である。

【図2】第1の実施例で互いに部分データの対応付け処理の対象となる静止画像（a）と当該静止画像に関連する音声（b）の一例について示す図である。

50

【図3】本発明の第1の実施例において静止画像と音声を入力して静止画像と音声との対応関係（画像音声対応情報）を求める際のモジュール構成を示すブロック図である。

【図4】第1の実施例における文字認識部202の細部モジュール構成を示すブロック図である。

【図5】第1の実施例における音声認識部204の細部モジュール構成を示すブロック図である。

【図6】図7に示す文字認識結果情報と音声認識結果情報の例を示す図である。

【図7】図2に示す静止画像と音声の例に対する文字認識結果情報と音声認識結果情報を対応させた結果を示す図である。

【図8】第1の実施例における静止画像と音声との対応付けの一例を示す図である。 10

【図9】静止画像と音声との対応結果を用いたアプリケーションの例である。

【図10】本発明の第2の実施例における発音列マッチングによる静止画音声対応部205の細部モジュール構成を示すブロック図である。

【図11】第2の実施例における文字認識結果と音声認識結果に対する発音列の例を示す図である。

【図12】本発明の第3の実施例における文字列マッチングを行う静止画音声対応部205の細部モジュール構成を示すブロック図である。

【図13】第3の実施例における文字認識結果と音声認識結果に対する文字列の例である。

【図14】第4の実施例における文字認識結果（a）と音声認識結果（b）のスコア情報（尤度や確率等で表された認識結果）を保持した複数候補の例を示す図である。 20

【図15】第4の実施例における文字認識結果を発音列に変換した結果（a）と音声認識結果から得られる発音列（b）のスコア情報を保持した複数候補の例を示す図である。

【図16】第4の実施例における文字認識結果から得られる文字列（a）と音声認識結果を文字列に変換した結果（b）のスコア情報を保持した複数候補の例を示す図である。

【図17】本発明の第6の実施例における静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図18】本発明の第6の実施例における文字認識結果を音声認識に利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図19】第6の実施例における文字認識結果を発音列に変換して利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。 30

【図20】第6の実施例における文字認識結果を発音列に変換して探索部で利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図21】第6の実施例における文字認識結果の文字列を音声認識結果を文字列に変換する際に利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図22】第7の実施例における文字認識結果から重要語を抽出して探索部で利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図23】第8の実施例における種々のフォント情報をもった静止画像の例である。

【図24】第8の実施例における文字領域からフォント情報を抽出して文字認識結果情報として出力する静止画像・音声認識装置のモジュール構成を示すブロック図である。 40

【図25】図23に示す静止画像からの文字認識結果と各文字領域のフォント情報を示す図である。

【図26】第9の実施例における文字認識結果情報出力部の細部モジュール構成を示すブロック図である。

【図27】第9の実施例における静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図28】第9の実施例における音声認識結果を文字列に変換して利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図29】第9の実施例における音声認識結果から得られる文字列を文字認識で利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。 50

【図30】第9の実施例における音声認識結果から得られる発音列を文字認識結果の発音列の変換に利用する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図31】より複雑な静止画像(a)とそれに関連付けられた音声(b)の一例を示す図である。

【図32】図31に示すような複雑な静止画像を分割する機能を有する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図33】図31(a)の静止画像を静止画像分割部1901によって5つの領域に分割された結果を示す図である。

【図34】図33に示す各分割領域に対して領域IDを付与して図32(a)の座標系における座標情報を表した図である。

10

【図35】図34と図32(a)の対応関係について示す図である。

【図36】文字検出部1902による文字領域の検出結果例を示す図である。

【図37】文字認識及び音声認識の結果得られる文字認識結果情報(a)及び音声認識結果情報を(b)示す図である。

【図38】図36に示す文字領域の検出結果に文字認識結果情報を対応させた結果を示す図である。

【図39】図31(a)に示す静止画像に関連する音声を説明するための図である。

【図40】第11の実施例に係る静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図41】図31の例に対する文字認識結果情報と音声認識結果情報の一例を示す図である。

20

【図42】重要語抽出によるワードスポットティングを用いた場合の音声認識結果情報の一例を示す図である。

【図43】図33で示した静止画像の分割をさらに行った場合の分割結果(a)(一点破線)、(a)をさらに分割した場合の結果(b)(二点破線)を示す図である。

【図44】階層的な静止画像の分割を木構造で表現した例を示す図である。

【図45】階層的に音声分割を行った場合の例を示す図である。

【図46】図45で階層的に分割された音声を木構造で表現した例である。

【図47】静止画像の木構造ノードに複数候補の分割音声を対応付けた結果の一例を示す図である。

30

【図48】図31の例に対する静止画像と音声の対応結果を用いたアプリケーションの一例を示す図である。

【図49】静止画像の木構造の結果及び複数候補音声を利用する際のユーザインタフェースの一例を示す図である。

【図50】第13の実施例における静止画像と音声の例を示す図である。

【図51】第13の実施例における文字概念変換機能及び音声概念変換機能を有する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図52】文字概念変換結果と静止画像の座標情報、及び音声概念変換結果と音声の時間情報の一例を示す図である。

【図53】第14の実施例において用いる静止画像とそれに対応付けられる音声の例を示す図である。

40

【図54】本発明の第14の実施例に係るオブジェクト認識処理機能を有する静止画像・音声処理装置のモジュール構成を示すブロック図である。

【図55】オブジェクト認識結果情報の例(a)とオブジェクト認識結果情報から得られる画像領域の例(b)を示す図である。

【図56】第15の実施例において用いる静止画像とそれに対応付けられる音声の例を示す図である。

【図57】本発明の第15の実施例に係る人物認識機能及び話者認識機能を有する静止画像・音声認識装置のモジュール構成を示すブロック図である。

【図58】第15の実施例における人物認識結果情報及び話者認識結果情報の一例を示す

50

図である。

【図59】人物認識結果情報から得られる画像領域を示す図である。

【図60】図2に示す静止画像と音声との対応結果を用いた別のアプリケーションに基づく表示例を示す図である。

【図61】図43に示す静止画像と音声との対応結果を用いた別のアプリケーションに基づく表示例を示す図である。

【図62】第10の実施例に係る静止画像・音声認識装置の処理動作を説明するためのフローチャートである。

【符号の説明】

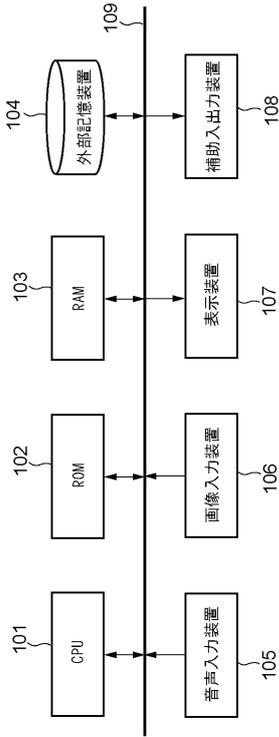
【0167】

- 201 文字検出部
- 202 文字認識部
- 203 音声検出部
- 204 音声認識部
- 205 静止画音声対応部
- 301 前処理部
- 302 特徴抽出部
- 303 識別部
- 304 文字認識結果情報出力部
- 305 文字認識用テンプレート
- 306 文字認識用言語モデル
- 401 音声分析部
- 402 探索部
- 403 音声認識結果情報出力部
- 404 音声認識用音響モデル
- 405 音声認識用言語モデル

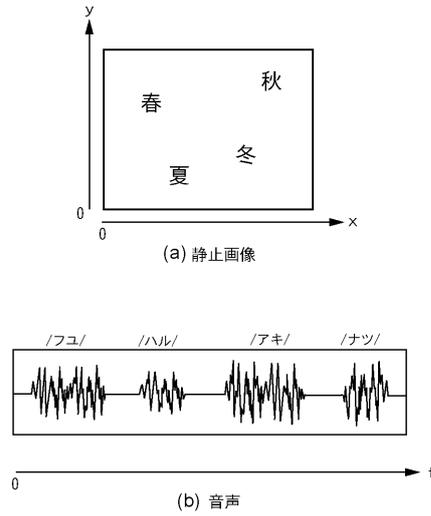
10

20

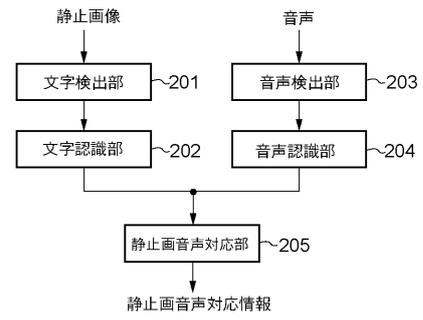
【図1】



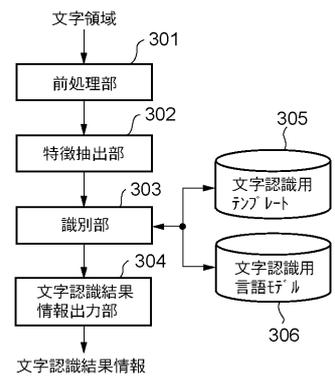
【図2】



【図3】



【図4】



【図6】

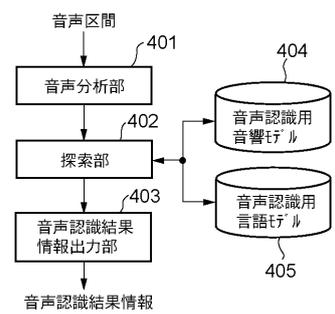
文字認識結果	座標情報
春	(x ₁ , y ₁)
夏	(x ₂ , y ₂)
秋	(x ₃ , y ₃)
冬	(x ₄ , y ₄)

(a) 文字認識結果情報の列

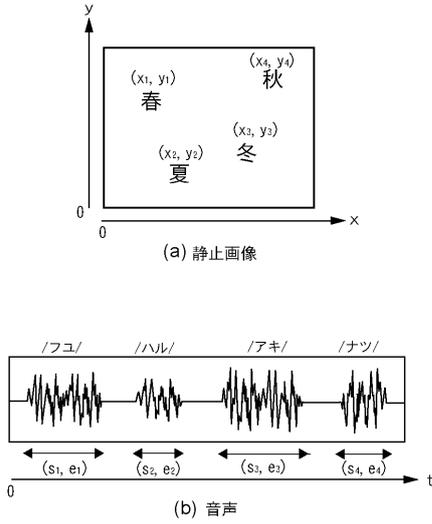
音声認識結果	時間情報
冬(フユ)	(s ₁ , e ₁)
春(ハル)	(s ₂ , e ₂)
秋(アキ)	(s ₃ , e ₃)
夏(ナツ)	(s ₄ , e ₄)

(b) 音声認識結果情報の列

【図5】



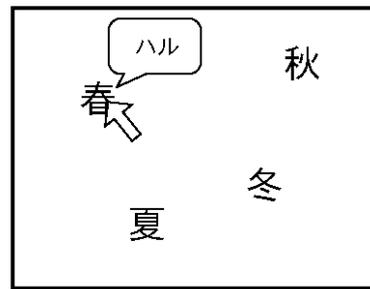
【 図 7 】



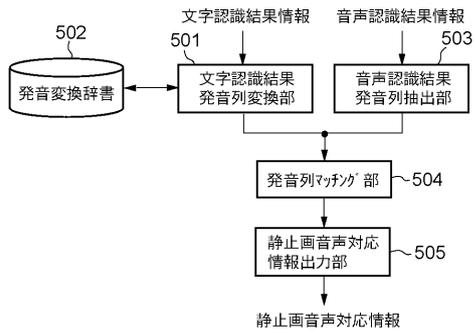
【 図 8 】

静止画像	音声
(x ₁ , y ₁)	(s ₂ , e ₂)
(x ₂ , y ₂)	(s ₄ , e ₄)
(x ₃ , y ₃)	(s ₃ , e ₃)
(x ₄ , y ₄)	(s ₁ , e ₁)

【 図 9 】



【 図 10 】



【 図 11 】

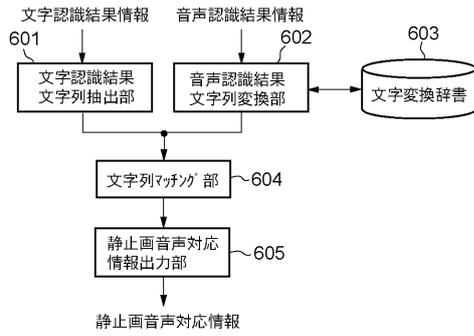
発音候補		座標情報
ハル	シュン	(x ₁ , y ₁)
ナツ	カ	(x ₂ , y ₂)
アキ	シュウ	(x ₃ , y ₃)
フユ	トウ	(x ₄ , y ₄)

(a) 文字認識結果に対する発音列候補

発音候補	時間情報
フユ	(s ₁ , e ₁)
ハル	(s ₂ , e ₂)
アキ	(s ₃ , e ₃)
ナツ	(s ₄ , e ₄)

(b) 音声認識結果に対する発音列

【 図 1 2 】



【 図 1 3 】

文字候補	座標情報
春	(x ₁ , y ₁)
夏	(x ₂ , y ₂)
秋	(x ₃ , y ₃)
冬	(x ₄ , y ₄)

(a) 文字認識結果に対する文字列

文字候補			時間情報
冬	不輸	—	(s ₁ , e ₁)
春	張る	貼る	(s ₂ , e ₂)
空	飽き	秋	(s ₃ , e ₃)
夏	奈津	捺	(s ₄ , e ₄)

(b) 音声認識結果に対する文字列候補

【 図 1 4 】

文字認識結果候補			
n/i	1	2	3
1	春(0.7)	香(0.2)	空(0.1)
2	科(0.5)	秋(0.3)	和(0.2)
3	夏(0.6)	真(0.3)	厚(0.1)
4	各(0.5)	冬(0.3)	尽(0.2)

(a) 文字認識結果候補の例

音声認識結果候補			
m/j	1	2	3
1	フユ(冬)(0.5)	フル(古)(0.4)	ツユ(露)(0.1)
2	ハル(春)(0.7)	タル(樽)(0.2)	ハク(白)(0.1)
3	アキ(秋)(0.6)	アシ(足)(0.2)	マキ(薪)(0.2)
4	マツ(松)(0.8)	ツ(津)(0.1)	ナツ(夏)(0.1)

(b) 音声認識結果候補の例

【 図 1 5 】

文字認識結果候補			
n/i	1	2	3
1	<u>ハル(0.7)</u> シュン(0.7)	カ(0.2) コウ(0.2)	ソラ(0.1) <u>アキ(0.1)</u> クウ(0.1)
2	カ(0.5) シナ(0.5)	アキ(0.3) シュウ(0.3)	ワ(0.2) ナゴ(0.2)
3	ナツ(0.6) カ(0.6)	シン(0.3) マ(0.3)	コウ(0.1) アツ(0.1)
4	カワ(0.5) オノ(0.5)	<u>フユ(0.2)</u> トウ(0.3)	ジン(0.2) <u>ツ(0.2)</u>

(a) 文字認識結果候補に対する発音列候補

音声認識結果候補			
m/j	1	2	3
1	<u>フユ(0.5)</u>	フル(0.4)	ツユ(0.1)
2	<u>ハル(0.7)</u>	タル(0.2)	ハク(0.1)
3	<u>アキ(0.6)</u>	アシ(0.2)	マキ(0.2)
4	マツ(0.8)	<u>ツ(0.1)</u>	ナツ(0.1)

(b) 音声認識結果候補に対する発音列

【 図 1 6 】

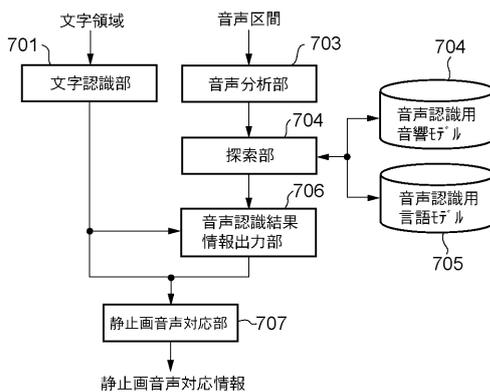
文字認識結果候補			
n/i	1	2	3
1	春(0.7)	香(0.2)	空(0.1)
2	科(0.5)	秋(0.3)	和(0.2)
3	夏(0.6)	真(0.3)	厚(0.1)
4	各(0.5)	冬(0.3)	尽(0.2)

(a) 文字認識結果候補に対する文字列

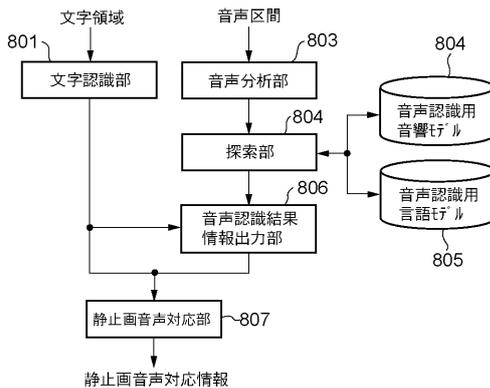
音声認識結果候補			
m/j	1	2	3
1	冬(0.5) 不輸(0.5)	古(0.4) 振る(0.4) 降る(0.4)	露(0.1) 梅雨(0.1) 汁(0.1)
2	春(0.7) 張る(0.7) 貼る(0.7)	樽(0.2) 足る(0.2)	白(0.1) 吐く(0.1) 履く(0.1)
3	空(0.6) 飽き(0.6) 秋(0.6)	足(0.2) 脚(0.2) 芦(0.2)	薪(0.2) 巻き(0.2) 真希(0.2)
4	松(0.8) 末(0.8) 待つ(0.8)	津(0.1)	夏(0.1) 奈津(0.1) 擦(0.1)

(b) 音声認識結果候補に対する文字列

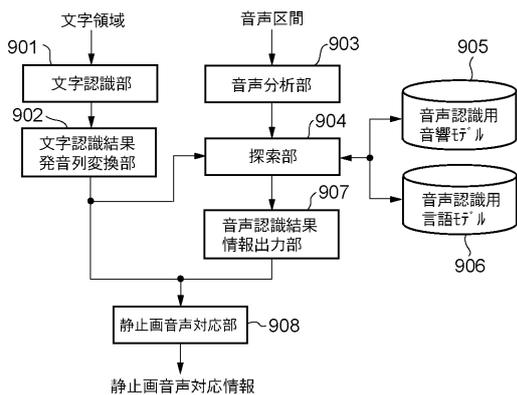
【 図 1 7 】



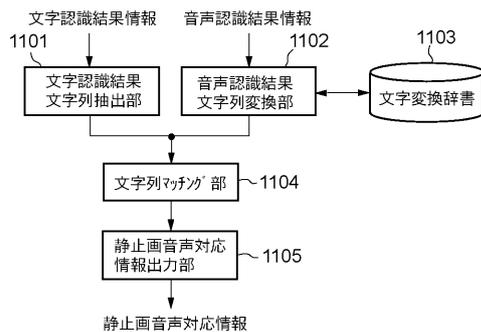
【 図 1 8 】



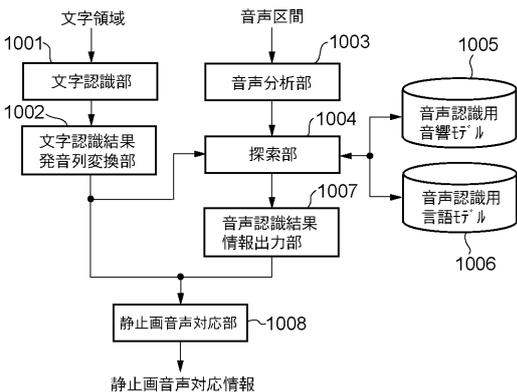
【 図 1 9 】



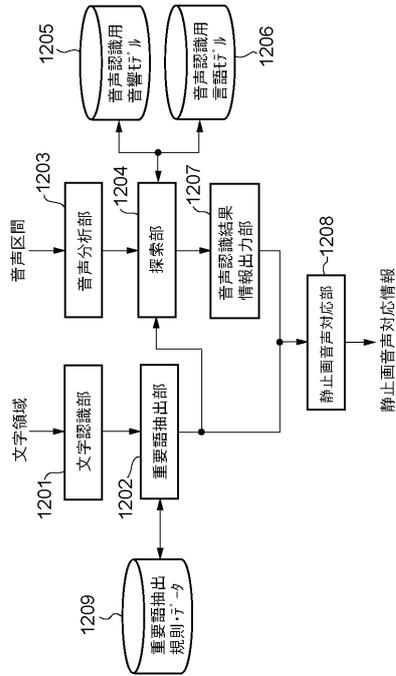
【 図 2 1 】



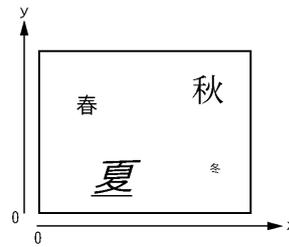
【 図 2 0 】



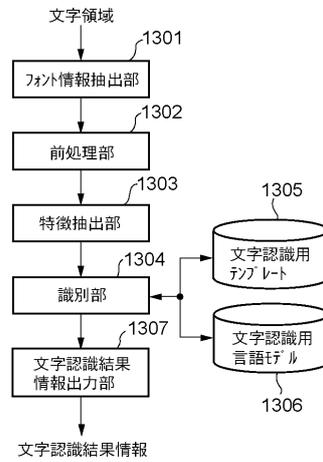
【図 2 2】



【図 2 3】



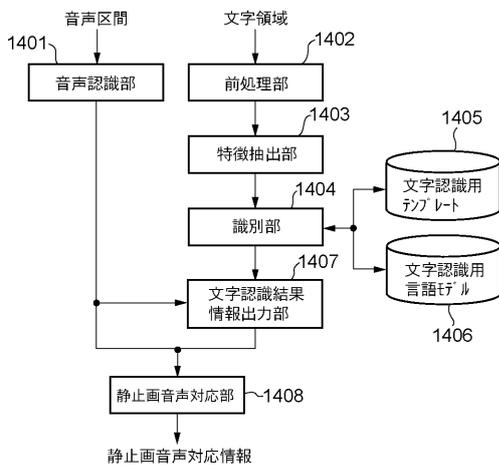
【図 2 4】



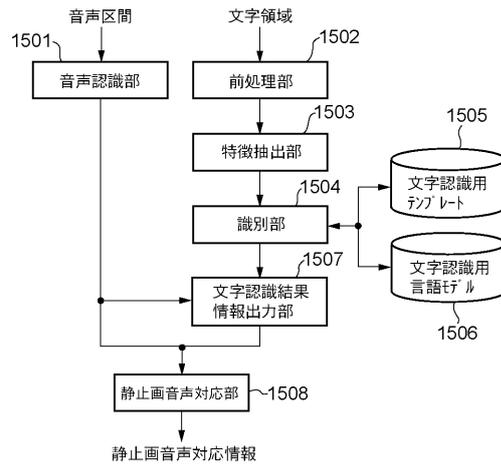
【図 2 5】

文字認識結果	フォント情報					座標情報
	サイズ	文字種	色	斜体	アンダーライン	
春	24	ゴシック	黒	N	N	(x ₁ , y ₁)
夏	44	ゴシック	黒	Y	Y	(x ₂ , y ₂)
秋	36	明朝	赤	N	N	(x ₃ , y ₃)
冬	12	ゴシック	黒	N	N	(x ₄ , y ₄)

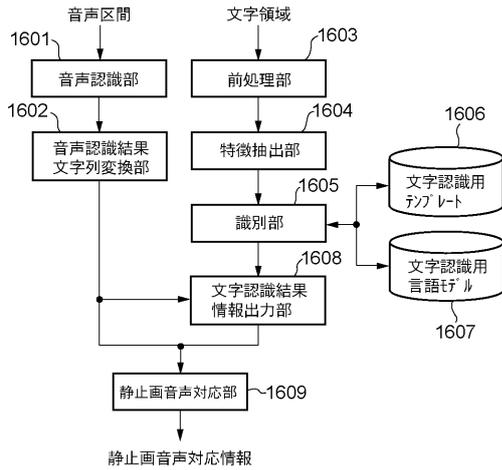
【図 2 6】



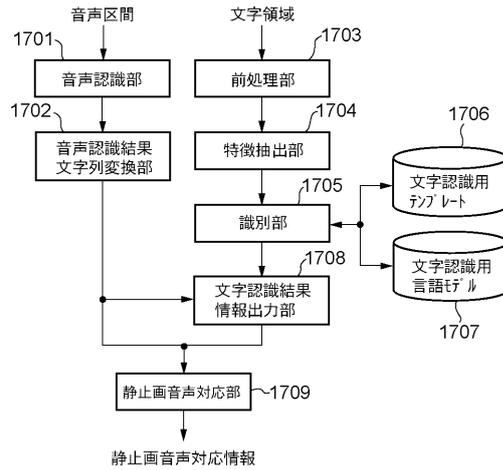
【図 2 7】



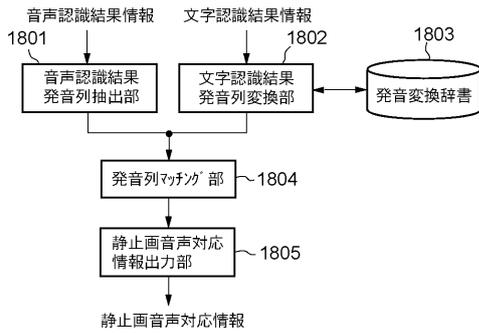
【図 28】



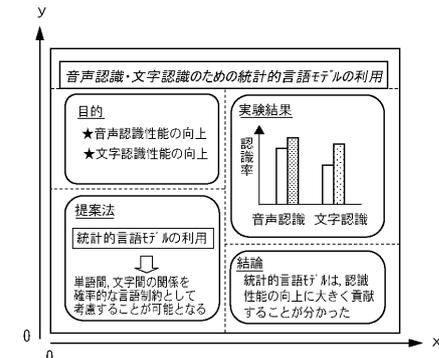
【図 29】



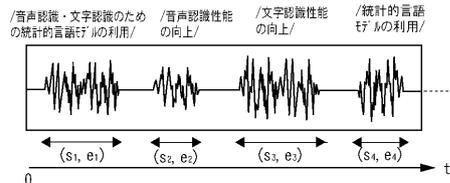
【図 30】



【図 31】

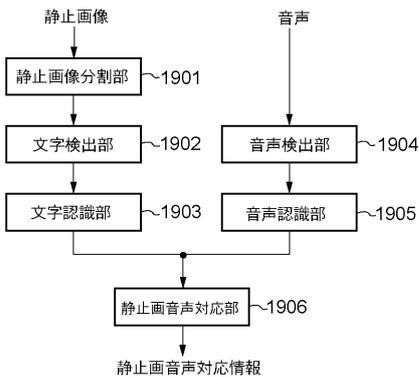


(a) 静止画像

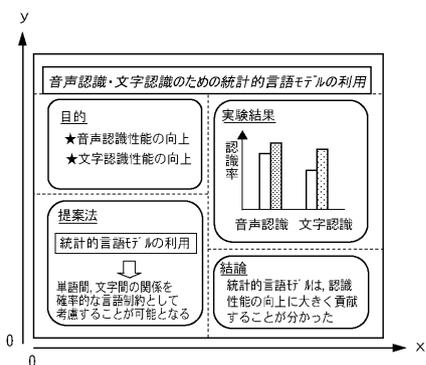


(b) 音声

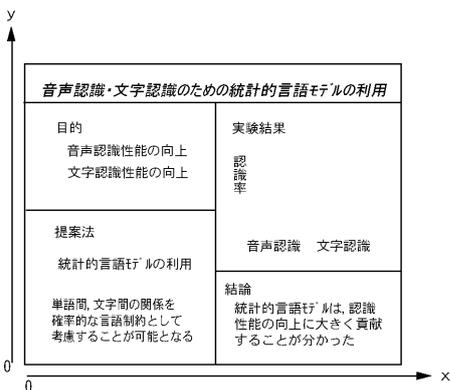
【図 3 2】



【図 3 3】



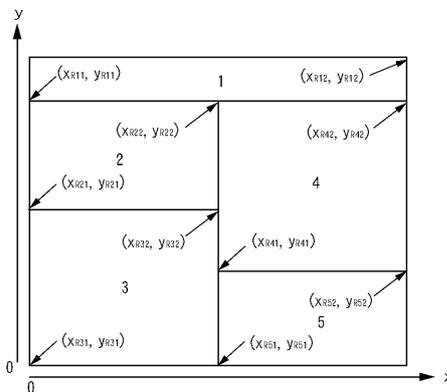
【図 3 6】



【図 3 4】

静止画分割結果(領域 D)	座標情報
1	(X _{R11} , Y _{R11} , X _{R12} , Y _{R12})
2	(X _{R21} , Y _{R21} , X _{R22} , Y _{R22})
3	(X _{R31} , Y _{R31} , X _{R32} , Y _{R32})
4	(X _{R41} , Y _{R41} , X _{R42} , Y _{R42})
5	(X _{R51} , Y _{R51} , X _{R52} , Y _{R52})

【図 3 5】



【図 3 7】

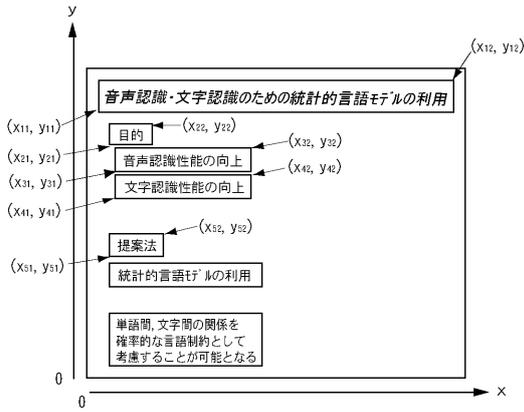
文字認識結果	座標情報
音声認識・文字認識のための統計的言語モデルの利用	(X ₁₁ , Y ₁₁ , X ₁₂ , Y ₁₂)
目的	(X ₂₁ , Y ₂₁ , X ₂₂ , Y ₂₂)
音声認識性能の向上	(X ₃₁ , Y ₃₁ , X ₃₂ , Y ₃₂)
文字認識性能の向上	(X ₄₁ , Y ₄₁ , X ₄₂ , Y ₄₂)
提案法	(X ₅₁ , Y ₅₁ , X ₅₂ , Y ₅₂)
...	...

(a) 文字認識結果情報の列

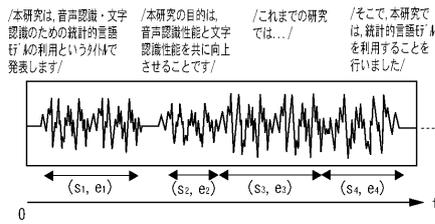
音声認識結果	時間情報
音声認識・文字認識のための統計的言語モデルの利用	(S ₁ , e ₁)
目的	(S ₂ , e ₂)
音声認識性能の向上	(S ₃ , e ₃)
文字認識性能の向上	(S ₄ , e ₄)
提案法	(S ₅ , e ₅)
...	...

(b) 音声認識結果情報の列

【 図 3 8 】



【 図 3 9 】



【 図 4 1 】

文字認識結果	座標情報
音声認識・文字認識のための統計的言語モデルの利用	(X11, Y11, X12, Y12)
目的	(X21, Y21, X22, Y22)
音声認識性能の向上	(X31, Y31, X32, Y32)
文字認識性能の向上	(X41, Y41, X42, Y42)
提案法	(X51, Y51, X52, Y52)
...	...

(a) 文字認識結果情報の列

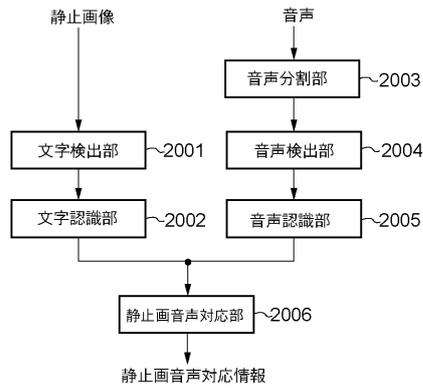
音声認識結果	時間情報
本研究、音声認識・文字認識の...	(s1, e1)
本研究の目的は、音声認識性能と...	(s2, e2)
これまでの研究では、...	(s3, e3)
そこで、本研究では、統計的言語...	(s4, e4)
...	...

(b) 音声認識結果情報の列

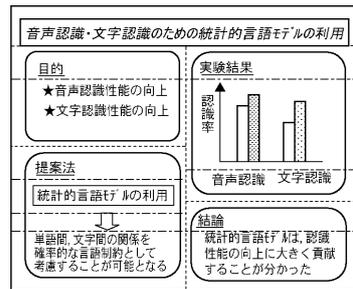
【 図 4 2 】

音声認識結果	時間情報
*音声認識*文字認識*統計的言語モデル*	(s1, e1)
*目的*音声認識*文字認識*	(s2, e2)
NO_RESULTS	(s3, e3)
統計的言語モデル	(s4, e4)
...	...

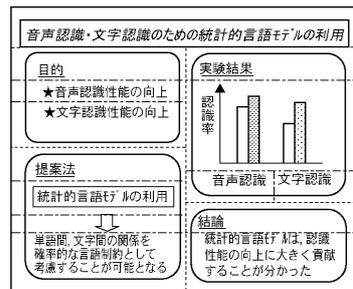
【 図 4 0 】



【 図 4 3 】

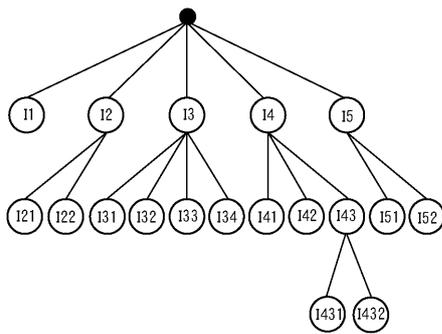


(a) 更に画像分割を行った結果

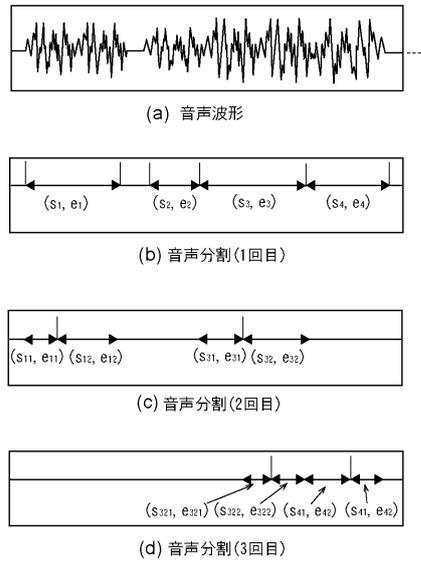


(b) (a)を更に画像分割した結果

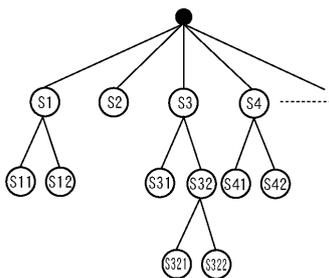
【 図 4 4 】



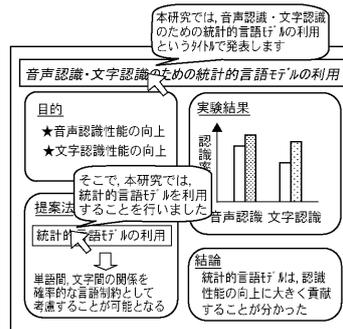
【 図 4 5 】



【 図 4 6 】



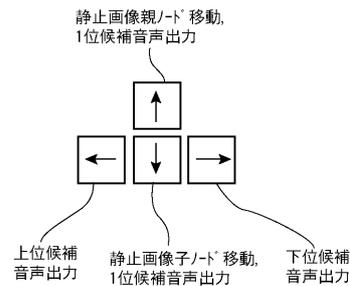
【 図 4 8 】



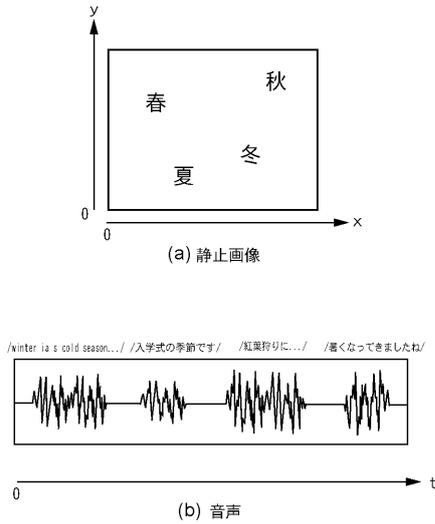
【 図 4 7 】

静止画像	音声		
	第1候補	第2候補	第3候補
11	S2	S31	NULL
12	S3	S31	S321
13	S4	S41	S42
14	S3	S321	S31
15	S1	S11	S12
121	S31	NULL	NULL
122	S3	S2	S42
131	S12	S1	S2
132	S322	S321	NULL
133	NULL	NULL	NULL
...

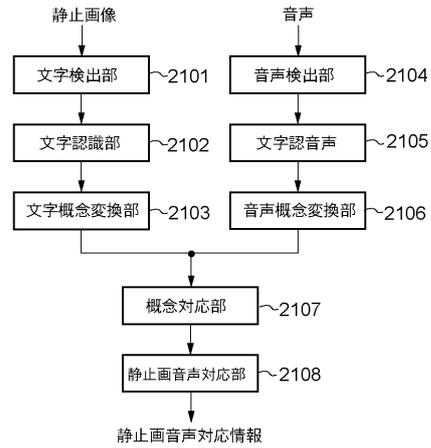
【 図 4 9 】



【 図 5 0 】



【 図 5 1 】



【 図 5 2 】

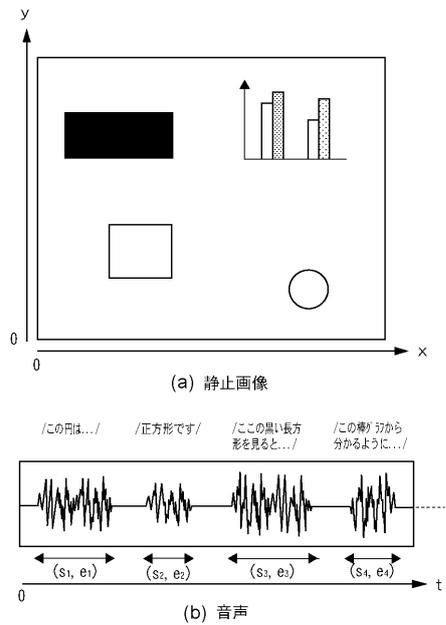
文字概念変換結果 (文字認識結果)	座標情報
\$ SPRING(春)	(x_1, y_1)
\$ SUMMER(夏)	(x_2, y_2)
\$ AUTUMN(秋)	(x_3, y_3)
\$ WINTER(冬)	(x_4, y_4)

(a) 文字概念変換結果と座標情報の列

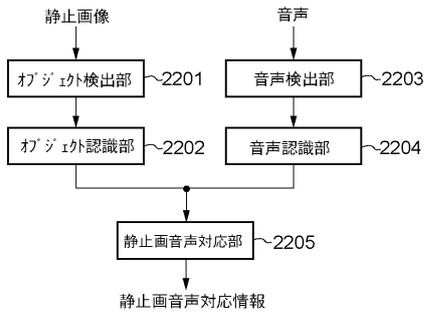
音声概念変換結果 (音声認識結果)	時間情報
\$ WINTER(Winter is a ...)	(s_1, e_1)
\$ SPRING(入学式の...)	(s_2, e_2)
\$ AUTUMN(紅葉狩りに...)	(s_3, e_3)
\$ SUMMER(暑くなつて...)	(s_4, e_4)

(b) 音声概念変換結果と時間情報の列

【 図 5 3 】



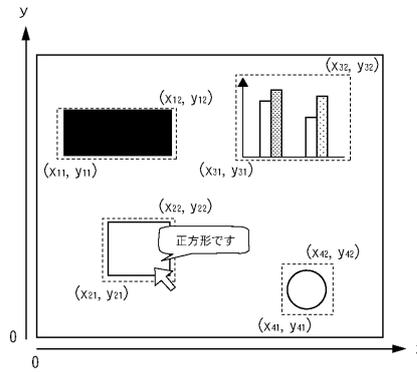
【図54】



【図55】

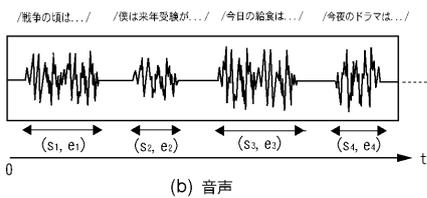
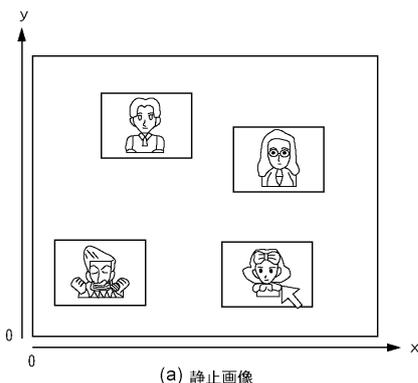
オブジェクト認識結果	座標情報
長方形, 黒	(X11, Y11, X12, Y12)
正方形, 白	(X21, Y21, X22, Y22)
棒グラフ, 灰	(X31, Y31, X32, Y32)
円, 白	(X41, Y41, X42, Y42)

(a) オブジェクト認識結果情報の例

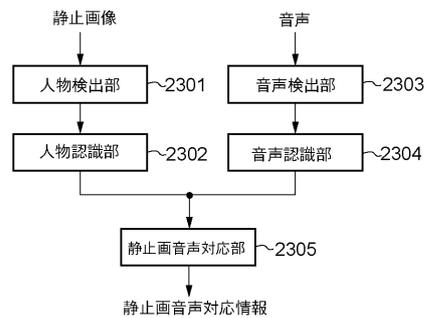


(b) オブジェクト認識結果情報から得られる画像領域

【図56】



【図57】



【図58】

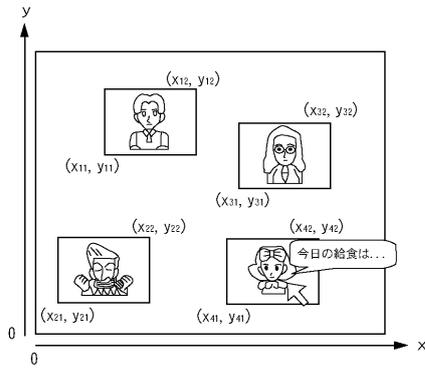
人物認識結果	座標情報
成人, 男性	(X11, Y11, X12, Y12)
高齢者, 男性	(X21, Y21, X22, Y22)
成人, 女性	(X31, Y31, X32, Y32)
子供, 女性	(X41, Y41, X42, Y42)

(a) 話者認識結果情報の列

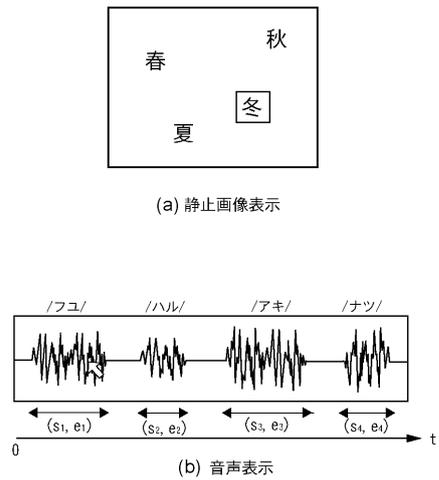
話者認識結果	時間情報
高齢者, 男性	(s1, e1)
成人, 男性	(s2, e2)
子供, 女性	(s3, e3)
成人, 女性	(s4, e4)

(b) 話者認識結果情報の列

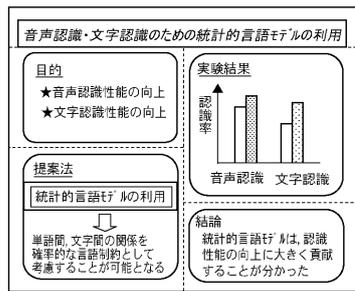
【 図 5 9 】



【 図 6 0 】

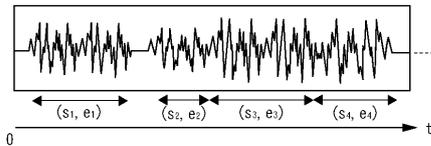


【 図 6 1 】



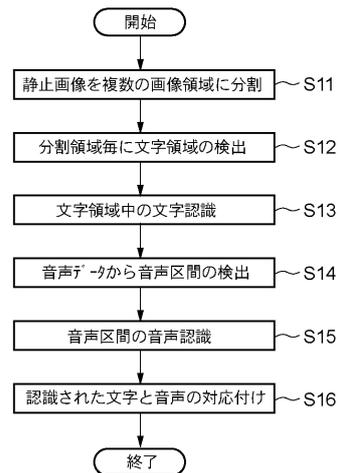
(a) 静止画像表示

/本研究は、音声認識・文字認識のための統計的言語モデルの利用という外は、発表します/
 /本研究の目的は、音声認識性能と文字認識性能を共に向上させることです/
 /これまでの研究では.../
 /そこで、本研究では、統計的言語モデルを利用することをを行いました/



(b) 音声表示

【 図 6 2 】



フロントページの続き

審査官 加藤 恵一

- (56)参考文献 特開平09 - 218955 (JP, A)
特開平09 - 233442 (JP, A)
特開2001 - 034151 (JP, A)
特開2002 - 197103 (JP, A)
特開2003 - 085572 (JP, A)
特開2003 - 242150 (JP, A)

- (58)調査した分野(Int.Cl., DB名)
H04N 5/76 - 5/956