

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4738847号
(P4738847)

(45) 発行日 平成23年8月3日(2011.8.3)

(24) 登録日 平成23年5月13日(2011.5.13)

(51) Int.Cl.

F I

G 1 O L 15/00 (2006.01)

G 1 O L 15/00 2 O O T

G 1 O L 15/10 (2006.01)

G 1 O L 15/10 3 O O G

請求項の数 12 (全 18 頁)

(21) 出願番号 特願2005-63149 (P2005-63149)
 (22) 出願日 平成17年3月7日(2005.3.7)
 (65) 公開番号 特開2006-243673 (P2006-243673A)
 (43) 公開日 平成18年9月14日(2006.9.14)
 審査請求日 平成20年2月19日(2008.2.19)

(73) 特許権者 000001007
 キヤノン株式会社
 東京都大田区下丸子3丁目30番2号
 (74) 代理人 100076428
 弁理士 大塚 康德
 (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 山本 寛樹
 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】 データ検索装置および方法

(57) 【特許請求の範囲】

【請求項1】

検索対象のデータがそれぞれ音声データと関連付けられて記憶されたデータベースから、ユーザにより入力されたキーワードを基にデータを検索するデータ検索装置であって、前記データベース内の各データに関連付けられた音声データに対し音声認識を行い、サブワード表現形式で認識結果を出力する第1の音声認識手段と、

前記キーワードをサブワード表現形式に変換する変換手段と、

前記第1の音声認識手段により得られたサブワード表現形式の前記認識結果と、前記変換手段によりサブワード表現形式に変換された前記キーワードとに基づいて、前記キーワードと前記データベース内の各データに関連付けられた音声データとの類似度を計算する類似度計算手段と、

前記類似度計算手段により計算された前記類似度に基づき選択される1または2以上のデータの各々について、そのデータに関連付けられた音声データを入力とし、前記サブワード表現形式に変換された前記キーワードを認識対象語とする音声認識を行う第2の音声認識手段と、

前記第2の音声認識手段の認識スコアに基づいて検索スコアを計算する検索スコア計算手段と、

前記検索スコア計算手段により計算された前記検索スコアに基づいて選択される前記データベース内のデータを検索結果としてユーザに提示する検索結果提示手段と、

を有し、

10

20

前記検索スコア計算手段は、前記検索スコアとして、前記類似度計算手段により計算された類似度と前記第 2 の音声認識手段により得られた認識スコアとの重み付き和を計算する

ことを特徴とするデータ検索装置。

【請求項 2】

前記第 1 の音声認識手段および前記変換手段は、前記キーワードが入力される前にあらかじめ実行されるものであり、前記類似度計算手段、前記第 2 の音声認識手段、前記検索スコア計算手段、および前記検索結果提示手段は、前記キーワードが入力されたことに応じて動作することを特徴とする請求項 1 に記載のデータ検索装置。

【請求項 3】

前記類似度計算手段は、前記類似度として、前記変換手段によりサブワード表現形式に変換された前記キーワードを正解とする前記第 1 の音声認識手段により得られたサブワード表現形式の前記認識結果のサブワード正解率またはサブワード正解精度を計算することを特徴とする請求項 1 または 2 に記載のデータ検索装置。

【請求項 4】

前記サブワード正解精度は、正解サブワード数から挿入誤りサブワード数、置換誤りサブワード数、および削除誤りサブワード数をそれぞれ引いて得たサブワード数と、前記正解サブワード数との比でもって表されるものであって、前記挿入誤りサブワード数に所定の重み係数が乗じられることを特徴とする請求項 3 に記載のデータ検索装置。

【請求項 5】

前記サブワードは、音素または音節であることを特徴とする請求項 1 から 4 までのいずれか 1 項に記載のデータ検索装置。

【請求項 6】

前記第 2 の音声認識手段により実行される音声認識は、前記キーワードを認識対象語とするキーワードスポッティングであることを特徴とする請求項 1 から 5 までのいずれか 1 項に記載のデータ検索装置。

【請求項 7】

前記第 2 の音声認識手段は、前記類似度が大きい順に所定個数のデータを選択し、当該選択されたデータの各々について前記音声認識を行うことを特徴とする請求項 1 から 6 までのいずれか 1 項に記載のデータ検索装置。

【請求項 8】

前記第 2 の音声認識手段は、前記類似度が所定の値よりも大きい 1 または 2 以上のデータを選択し、当該選択されたデータの各々について前記音声認識を行うことを特徴とする請求項 1 から 7 までのいずれか 1 項に記載のデータ検索装置。

【請求項 9】

前記検索結果提示手段は、前記検索スコアが大きい順に所定個数のデータを検索結果として表示することを特徴とする請求項 1 から 8 までのいずれか 1 項に記載のデータ検索装置。

【請求項 10】

前記検索結果提示手段は、前記検索スコアが所定の値よりも大きいデータを検索結果として表示することを特徴とする請求項 1 から 9 までのいずれか 1 項に記載のデータ検索装置。

【請求項 11】

検索対象のデータがそれぞれ音声データと関連付けられて記憶されたデータベースから、ユーザにより入力されたキーワードを基にデータを検索するデータ検索装置によって実行されるデータ検索方法であって、

第 1 の音声認識手段が、前記データベース内の各データに関連付けられた音声データに対し音声認識を行い、サブワード表現形式で認識結果を出力する第 1 の音声認識ステップと、

変換手段が、前記キーワードをサブワード表現形式に変換する変換ステップと、

10

20

30

40

50

類似度計算手段が、前記第 1 の音声認識ステップにより得られたサブワード表現形式の前記認識結果と、前記変換ステップによりサブワード表現形式に変換された前記キーワードとに基づいて、前記キーワードと前記データベース内の各データに関連付けられた音声データとの類似度を計算する類似度計算ステップと、

第 2 の音声認識手段が、前記類似度計算ステップにより計算された前記類似度に基づき選択される 1 または 2 以上のデータの各々について、そのデータに関連付けられた音声データを入力とし、前記サブワード表現形式に変換された前記キーワードを認識対象語とする音声認識を行う第 2 の音声認識ステップと、

検索スコア計算手段が、前記第 2 の音声認識ステップでの認識スコアに基づいて検索スコアを計算する検索スコア計算ステップと、

検索結果提示手段が、前記検索スコア計算ステップにより計算された前記検索スコアに基づいて選択される前記データベース内のデータを検索結果としてユーザに提示する検索結果提示ステップと、

を有し、

前記検索スコア計算ステップでは、前記検索スコア計算手段が、前記検索スコアとして、前記類似度計算ステップで計算された類似度と前記第 2 の音声認識ステップで得られた認識スコアとの重み付き和を計算する

ことを特徴とするデータ検索方法。

【請求項 1 2】

請求項 1 に記載のデータ検索方法をコンピュータに実行させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ検索技術に関し、特に、音声認識を用いて、音声情報が付与されたマルチメディアデータを検索するデータ検索技術に関する。

【背景技術】

【0002】

音声はキー入力に不慣れな人やキーボードを設置するのが物理的に困難な小型機器での情報入力手段として有効である。現在市販されているデジタルカメラの多くは録音機能を搭載しており、撮影した画像に音声でメモをつけることができる。また、ICレコーダの利用事例として、ICレコーダを携帯し、日常のメモ代わりに音声メモ機能を利用する事例も増えている。これらの録音された音声は単に再生するにとどまらず、音声認識を利用して検索に利用することができる。

【0003】

音声認識を利用した音声データの検索のタイプは次の二通りが考えられる。

【0004】

(1) あらかじめ検索に用いるキーワード(あるいはフレーズ)が決められているもの。

(2) 任意のキーワード(あるいはフレーズ)を使用できるもの。

【0005】

(1) の場合、検索対象となる音声データに対して、あらかじめ決められたキーワードを含んだ認識辞書あるいは認識文法を用いて音声データを音声認識した認識結果を記憶しておき、ユーザが指定した検索用キーワードと一致あるいは部分一致する認識結果となる音声データを検索結果として提示する。音声認識の方法としてキーワードスポッティングを用いて音声データからキーワードを抽出してもよい。この方法の利点は、検索に先立って音声入力時に音声認識を実行できる点である。計算量の多い音声認識処理を事前に行っておけば、認識結果と検索用キーワードとの文字列同士の比較という比較的計算量の少ない処理だけで検索できる。しかしながら、認識辞書あるいは認識文法に含まれるあらかじめ決められたキーワードしか音声認識できないため、付与する音声の内容が制限され、音声メモの付与・検索という活用目的で利用する場合にはその活用範囲が制限されることになる。この方法を実現した従来技術の一例が、特開 2003-219327 公報(特許文

10

20

30

40

50

献 1) に開示されている。

【 0 0 0 6 】

(2) の場合、検索時にユーザが入力したキーワード (あるいはフレーズ) を認識辞書あるいは認識文法を用いて音声データを音声認識し、ユーザが指定した検索用キーワードと一致あるいは部分一致する認識結果となる音声データを検索結果として提示する。音声認識の方法としてキーワードスポッティングを用いて、ユーザが指定した任意のキーワードが音声データ内に含まれるかどうかを判定しても良い。この方法を実現した従来技術の一例が、特開平 1 0 - 1 7 3 7 6 9 号公報 (特許文献 2) に開示されている。この特許文献 2 に開示された音声メッセージ検索装置では、ユーザが入力したキーワード文字列を不特定話者の標準パターンに変換して、この標準パターンを用いてキーワードスポッティング音声認識を行い、音声メッセージ中でキーワードが検出された箇所が再生される。この方法の利点は、任意の内容を音声入力でき、任意のキーワードで検索できる点である。しかしながら、検索用のキーワードを入力してから音声認識を実行するので、検索時に計算量の多い音声認識を実行する必要がある、検索対象とする音声データが大量にある場合は検索に時間がかかるという問題がある。

10

【 0 0 0 7 】

特許文献 1 および 2 の欠点を解消する検索方法の一例が、特開 2 0 0 2 - 2 7 8 5 7 9 公報 (特許文献 3) に開示されている。特許文献 3 に開示されている音声データ検索装置では、検索対象となる音声データ (特許文献 3 では「音声波形データ」) を音声シンボル列に変換したものを記録しておき、キーワード (特許文献 3 中では「検索語」) を音声シンボル列に変換して、記憶されている音声データを変換した音声シンボル列と一致する部分を検索し、検索された部分に対してキーワードと一致しているか否かを判定する。この判定にキーワードスポッティング音声認識装置を用いることが特許文献 3 の請求項 4 に記載されている。この方法では、検索対象となる音声データの音声シンボル列への変換を検索に先立って事前に行うことにより、検索時にシンボルレベルでの高速な候補の絞り込みが実現され、さらに絞り込まれた候補に対してより正確な音声認識を用いた検索を実施することで、任意のキーワードに対して、正確で従来よりも高速な検索が実現できる。

20

【 0 0 0 8 】

【特許文献 1】特開 2 0 0 3 - 2 1 9 3 2 7 公報

【特許文献 2】特開平 1 0 - 1 7 3 7 6 9 号公報

【特許文献 3】特開 2 0 0 2 - 2 7 8 5 7 9 公報

30

【発明の開示】

【発明が解決しようとする課題】

【 0 0 0 9 】

特許文献 3 の方法では、シンボルレベルでのマッチングの際にキーワードと一致した部分のみを候補として絞り込むため、音声データを音声シンボル列に変換する際に誤った音声シンボル列に変換された候補は検索できないという問題がある。

【 0 0 1 0 】

特許文献 3 には、音声データから音声シンボル列への変換に音素認識結果を用いた例が開示されている。日本語の音素認識の場合、一般に 7 0 ~ 8 0 % の認識率で性能が良いとされている。つまり、いかに高性能な音素認識を用いても、音素認識結果には誤りが含まれることになる。特許文献 3 には、この対策として、音素認識結果から一般に誤りやすいとされる子音を無視して有音母音、長母音、発音、無音のみの音声シンボル列へ変換する方法や、コンフュージョンを起こしやすい子音をグループ化して 1 つの音節として扱う方法が開示されている。前者は例えば、「h a k o n e (ハコネ)」という音素認識結果に対して、「a a e (アアエ)」という音声シンボル列に変換する。後者は、例えば k と p と t を同じグループとして t に置き換え、「h a k o n e (ハコネ)」を「h a t o n e (ハトネ)」にマッピングする。入力されたキーワードにも同様のマッピングを行うので、子音の音素認識を誤っても検索することができる。

40

【 0 0 1 1 】

50

しかしながら、音素認識誤りは子音だけではなく母音にもしばしば起こる問題であるし、音素認識誤りは置換誤りだけではなく、脱落誤り、挿入誤りなども考えられる。前述の「h a k o n e (ハコネ)」の例では、例えば前述の「ハコネ」という音声データに対して、「h o k o n e (ホコネ)」(母音の認識誤り)、「h a k o n e e (ハコネー)」(挿入誤り)、「k o n e (コネ)」(脱落誤り)といった音素誤りが生じた場合、特許文献3に開示されている方法でキーワードを「ハコネ」として検索しても、検索できないという問題がある。これらの問題は、シンボルレベルで検索して候補を絞り込む際に、音声データの音声シンボル列からキーワードの音声シンボル列と一致する部分を検索するために生じている。

【0012】

10

本発明は、このような問題点に鑑みてなされたものであり、音声認識を用いたデータ検索における検索精度を高めることを目的とする。

【課題を解決するための手段】

【0013】

上記目的を達成するために、例えば本発明のデータ検索装置は以下の構成を備える。すなわち、本発明の一側面に係るデータ検索装置は、検索対象のデータがそれぞれ音声データと関連付けられて記憶されたデータベースから、ユーザにより入力されたキーワードを基にデータを検索するデータ検索装置であって、前記データベース内の各データに関連付けられた音声データに対し音声認識を行い、サブワード表現形式で認識結果を出力する第1の音声認識手段と、前記キーワードをサブワード表現形式に変換する変換手段と、前記第1の音声認識手段により得られたサブワード表現形式の前記認識結果と、前記変換手段によりサブワード表現形式に変換された前記キーワードとに基づいて、前記キーワードと前記データベース内の各データに関連付けられた音声データとの類似度を計算する類似度計算手段と、前記類似度計算手段により計算された前記類似度に基づき選択される1または2以上のデータの各々について、そのデータに関連付けられた音声データを入力とし、前記サブワード表現形式に変換された前記キーワードを認識対象語とする音声認識を行う第2の音声認識手段と、前記第2の音声認識手段の認識スコアに基づいて検索スコアを計算する検索スコア計算手段と、前記検索スコア計算手段により計算された前記検索スコアに基づいて選択される前記データベース内のデータを検索結果としてユーザに提示する検索結果提示手段とを有し、前記検索スコア計算手段は、前記検索スコアとして、前記類似度計算手段により計算された類似度と前記第2の音声認識手段により得られた認識スコアとの重み付き和を計算することを特徴とする。

20

30

【発明の効果】

【0014】

本発明によれば、音声認識を用いたデータ検索における検索精度を高めることができる。

【発明を実施するための最良の形態】

【0015】

以下、図面を参照して本発明の好適な実施形態について詳細に説明する。

【0016】

40

(実施形態1)

本実施形態では、データ検索装置の一例として、画像データに関連付けられた音声データを用いて画像データを検索する画像データ検索装置について説明する。なお、本発明に係る検索の対象は画像データに限定されるものではなく、文書、図形などその他の種類のデータにも適用が可能である。

【0017】

図1は、本実施形態における画像データ検索装置の構成を示すブロック図である。

【0018】

図1において、101は制御メモリ(ROM)、102は中央処理装置(CPU)、103はメモリ(RAM)、104はハードディスクなどの外部記憶装置、105はキー

50

やボタンなどの入力装置、１０６は液晶などの表示装置、１０７はバスである。画像データ検索処理を実現するための制御プログラム１０４aやその制御プログラムで用いるデータ（後述する言語処理用データ２１０、マルチメディアデータ２１１、音声認識用データ２１３）やサブワード認識結果２１２は、例えば外部記憶装置１０４に記憶される。このような構成であるから、本画像データ検索装置は汎用のコンピュータによっても実現することが可能である。

【００１９】

これらの制御プログラムやデータは、中央処理装置１０２の制御のもと、バス１０７を通じて適宜メモリ１０３に取り込まれ、中央処理装置１０２によって実行される。言うまでもないことであるが、制御プログラムやデータは制御メモリ１０１に記憶してもよい。

10

【００２０】

図２は、本実施形態における画像データ検索装置の機能構成を示すブロック図である。

【００２１】

制御部２０１は、他の機能モジュール間の連携やシステム全体の処理を制御する。キーワード取得部２０３は、画像検索時にユーザが入力装置１０５から入力するキーワードを取得する。キーワード・サブワード変換部２０３は、後述の言語処理用データ２１０を参照してキーワード取得部２０２で取得したキーワードを後述するサブワード音声認識部２０５で用いるサブワード表現形式に変換する。「サブワード」とは、単語を構成する単語よりも小さい音声の単位の総称であり、音声認識では音節や音素を用いる場合が多い。

【００２２】

20

言語処理用データ２１０は、言語処理用辞書に含まれ、単語からサブワードに変換する際に必要なデータが記録されている。例えば、検索キーワードとして漢字やかなが混在した文字列が入力された場合に、キーワード・サブワード変換部２０３では、漢字の読みを推測し、推測された読みからかな文字列に変換し、かな文字列からサブワードに変換するといった処理を行うが、このときに必要な漢字に対する読みやかな文字からサブワードへの変換ルールなどが言語処理用データ２１０として記録されている。また、英語をはじめとする日本語以外の言語の場合は、単語に対応するサブワードを記述した変換テーブルを直接記述したテーブルを言語処理用データ２１０として用いるのが一般的である。

【００２３】

なお、以後、説明の簡略化のため、ユーザが入力したキーワードをサブワード表現形式に変換したものを「クエリサブワード」と記述する。

30

【００２４】

以下では、サブワードとして音節を用いた場合について説明する。キーワード取得部２０３で取得したキーワードが「箱根山」であった場合、キーワード・サブワード変換部により、「箱根山」は「は こ ね や ま」というクエリサブワードに変換される。

【００２５】

音声データ取り込み部２０４は、検索対象であるマルチメディアデータ２１１に記憶された画像データに関連付けられた音声データを取り込む。

【００２６】

本実施形態における画像データ検索装置が検索対象とするマルチメディアデータ２１１は、音声に関連付けられた画像データであり、データベースとして外部記憶装置１０４に記憶されている。画像データと関連する音声データは一つのデータとして統合されていても良いし、それぞれ独立したデータとして記憶し、画像データと音声データの関連を管理する別のデータを記憶しておいても良い。また、関連する音声データと画像データの拡張子部分を除いたファイル名を同じにして画像データと音声データの関連を管理しても良い。

40

【００２７】

サブワード認識部２０５は、検索対象となる全ての画像データについて、音声取り込み部２０４で取り込んだ音声を入力とし、サブワード表現形式の認識結果を出力する音声認識を行い、得られた認識結果をサブワード認識結果２１２に記憶する。ここで用いる音

50

声認識は、認識結果としてサブワード表現形式で記述されたものが得られれば良いので、音素タイプライタや音節タイプライタなどサブワード表現形式の認識結果を出力するサブワードを単位とした音声認識を行っても良いし、一般にディクテーションとして知られる大語彙連続音声認識を行って得られた認識結果をサブワードに変換するようにしてもよい。

【0028】

サブワード音声認識は、各画像データごとに実行し、認識結果を記憶する際に画像データに関連付けて記憶する。画像データと認識結果を関連付ける方法は、前述した画像データと音声データの関連付けと同様の方法で実現できる。また、サブワード音声認識では、単一の認識結果だけでなく複数の認識結果候補を求めても良い。また、図1および2では、サブワード認識結果212と、画像および音声データを記憶するマルチメディアデータ211とを別々に図示しているが、サブワード認識結果を画像データや音声データと同様にマルチメディアデータ211に記憶しても良い。サブワード音声認識の際に必要な音響モデルをはじめとする各種データは音声認識用データ213に記憶されている。本実施形態の画像データ検索装置では、サブワード音声認識として音節認識を行い、認識スコアの高い上位3候補をサブワード認識結果212に記憶するものとする。

【0029】

図4に、「箱根山」という音声データに対して得られるサブワード認識結果の一例を示す。本実施形態では、サブワード認識結果として認識スコアの上位3候補を求め、401で示される順位、402で示される認識結果の音節列、403で示される認識スコアを、認識スコアの順に記憶する。

【0030】

サブワード類似度計算部206は、サブワード音声認識部205で得られたサブワード認識結果212とクエリサブワードを比較し、サブワード類似度を計算する。サブワード類似度は、マルチメディアデータ記憶部に記憶された各画像データごとに計算する。本実施形態の画像データ検索装置では、サブワード認識結果212に記憶された複数の音声認識候補について、クエリサブワードを正解とした場合の音節正解精度を求め、その最大値をサブワード類似度として計算する。図4に示したサブワード認識結果に対して、クエリサブワードが「は こ ね や ま」であった場合は、第1位(411)、第2位(412)、第3位(413)の音節認識結果に対して、それぞれクエリサブワードを正解とするサブワード正解精度としての音節正解精度を求める。音節正解精度は、クエリサブワードを正解とした場合に、正解に含まれる音節数をN、サブワード音声認識結果に含まれる挿入誤り音節数をI、置換誤り音節数をS、削除誤りをDとすると、 $((N - I - S - D) / N) \times 100$ で計算される。411に示した「ふぁ お ね や ま あ」を例にとると、正解音節数N=5、挿入誤り音節数I=1(「あ」)、置換誤り音節数S=2(「ふぁ」、「お」)、削除誤り音節数D=0であるので、 $((5 - 1 - 2 - 0) / 5) \times 100 = 40$ (%)となる。

【0031】

同様に、第2位(412)の、「ふぁ お ね や ま」(N=5、I=0、S=2、D=0)は60(%)、第3位(413)、「ふぁ こ ね や ま」(N=5、I=0、S=1、D=0)は80(%)となる。したがって、これらの最大値である80(%)を、クエリサブワードと図4に示したサブワード認識結果との類似度とする。

【0032】

以上のサブワード類似度計算部205におけるサブワード類似度の計算においては、サブワード音声認識で得られた複数のサブワード認識結果候補を用いたが、単に第1位のサブワード認識結果だけを用いても良い。また、複数のサブワード認識結果に対する音節正解精度の最大値を類似度として用いたが、各順位の音節正解精度に対して、順位(401)や認識スコア(403)に基づいて重み付けした値を求め、その最大値、あるいは和を類似度として用いても良い。また、音節正解精度ではなく、音節正解率や他の尺度を使って計算しても良い。また、言うまでもないことであるが、サブワード音声認識部205が

出力するサブワード認識結果の単位、音素その他のサブワードである場合には、サブワード類似度はその単位を基に計算される方法（音素正解精度など）を用いて計算することになる。

【0033】

キーワードスポッティング音声認識部207は、音声データ取り込み部204が取得した音声データを入力し、キーワード・サブワード変換部203でキーワードをサブワード表現形式に変換したクエリサブワードを認識対象語とするキーワードスポッティングを行い、認識スコアを求める。具体的には、音声認識用データ213に記憶されている音響モデルを用いて、クエリサブワードを表す音響モデルを構成し、この音響モデルを用いて、音声データ取り込み部204が取得した音声データを入力した際に計算される認識スコアを求める。キーワードスポッティングに必要な音響モデルをはじめとする各種データは音声認識用データ213に記憶されている。キーワードスポッティング音声認識はサブワード類似度計算部206で計算されたサブワード類似度に基づいて選択された画像データについてのみ行う。例えば、サブワード類似度が大きい所定個数の画像データを選択しても良いし、所定の閾値を越える画像データを選択しても良い。

10

【0034】

検索スコア計算部208は、キーワードスポッティング音声認識で求まる認識スコアを基に、検索スコアを求める。本実施形態における画像データ検索装置では、単にワードスポッティング音声認識で求まる認識スコアを検索スコアとして用いるが、これに加えて、サブワード類似度計算部206で計算されたサブワード類似度や、サブワード認識結果212に記憶されているサブワード音声認識で求めた認識スコアなどを組み合わせて検索スコアとして用いても良い。例えば、サブワード類似度をA、ワードスポッティング音声認識で求まる認識スコアをB、 $A + \alpha \times B$ のような重み付き和を検索スコアとしてもよい。

20

【0035】

検索結果表示部209は、検索スコア計算部208で求めた検索スコアに基づいて選択した画像データを表示装置106に表示する。この際、検索スコアの最も良い画像データのみを表示してもよいし、検索スコアが所定の値を超える画像データを表示してもよいし、検索スコアの良い所定個数の画像データを表示しても良い。また、表示の際に検索スコア順に並べ替えて表示してもよいし、検索スコアにかかわらず、ファイル名など画像データに付随したほかの情報を基準に並べえて表示しても良い。

30

【0036】

以上説明した機能モジュール構成で実現される画像データ検索装置による処理の流れを、図3のフローチャートを用いて説明する。このフローチャートに対応するプログラムは制御プログラム104aに含まれ、RAM103にロードされた後、CPU102によって実行される。

【0037】

まず、音声データ取り込み部204により、マルチメディアデータ211に記憶された画像データに関連付けられた音声データを取り込む（ステップS301）。

【0038】

次に、サブワード音声認識部206により、ステップS301で取り込んだ音声データを入力し、サブワード列を出力とする音声認識を行い（ステップS302）、その認識結果をサブワード認識結果212として記憶する（ステップS303）。前述した通り、ここで行う音声認識は、音素タイプライタや音節タイプライタなどのサブワードを単位とする音声認識を行ってもよいし、ディクテーションに代表される大語彙連続音声認識を行って得られた認識結果をサブワードに変換して出力するようにしてもよい。ステップS303では、音節認識を行い、認識スコアの高い上位3候補の認識結果の音節列や認識スコアをサブワード認識結果212として記憶する。

40

【0039】

ステップS301の音声の取り込み、ステップS302のサブワード音声認識、ステッ

50

ステップS303のサブワード認識結果の記憶の処理は、マルチメディアデータ211に記憶された全ての画像データについて、各画像データごとに実行する。

【0040】

次に、キーワード取得部202は、ユーザにより入力装置105を介して入力された画像検索用のキーワードを取得する(ステップS304)。取得したキーワードは、キーワード・サブワード変換部203でサブワード表現形式に変換する(ステップS305)。繰り返しになるが、サブワード表現形式に変換したキーワードをクエリサブワードと記述する。また、本実施形態の画像データ検索装置では、クエリサブワードの表現形式は音節とする。

【0041】

次に、サブワード類似度計算部205で、サブワード認識結果212として記憶されたサブワード認識結果とクエリサブワードとを比較し、サブワード類似度を計算する(ステップS306)。類似度の計算方法は、このサブワード類似度計算部205の機能に関して説明した通りである。このステップS306では、サブワード類似度を各画像データごとに計算する。

【0042】

次に、サブワード類似度が大きい所定個数の画像データを選択し(ステップS307)、キーワードスポッティング音声認識部207により、その選択された画像データに関連づけられた音声データについて、キーワード・サブワード変換部203で求めたクエリサブワードを認識対象語とするキーワードスポッティングを行い、認識スコアを求める(ステップS308)。

【0043】

次に、検索スコア計算部208により、ステップS308で求めた認識スコアを基に検索スコアを求め(ステップS309)、検索結果表示部209により、検索スコアのよい画像データを表示装置106に表示する(ステップS310)。

【0044】

本実施形態の画像データ検索装置ではキーワードスポッティング音声認識で求めた認識スコアを検索スコアとして用いる。キーワードスポッティング音声認識で求まる認識スコアは、入力されたキーワードをサブワード表現形式に変換したクエリサブワードを認識対象語として求めているため、ステップS302で行った言語制約が緩いサブワード音声認識で得られる認識結果に基づいてステップS306にてサブワード類似度計算部206で計算された類似度に比べ、音声データの“キーワードらしさ”をより正確に表していることが期待できる。すなわち、このキーワードスポッティング音声認識で求まる認識スコアに基づいて、画像データを選択することにより、検索精度が向上することが期待できる。

【0045】

以上説明した画像データ検索装置の処理では、説明の簡単のため、図3に示した全ての処理を一度に行う場合について説明したが、ステップS301～S303の処理については、ユーザが入力するキーワードに依存しない処理であり、また、ステップS304以降の処理とは独立に実行できるため、ステップS304の処理に先立って事前に行っておくことが好ましい。つまり、ステップS301～S303の処理は、画像データが更新された場合にのみ実行し、ステップS304以降のキーワードに依存した処理だけを、キーワードが入力される度に実行するように構成するのが望ましい。

【0046】

以上説明した実施形態によれば、シンボルレベルの検索で候補を絞り込む際に、サブワード表現形式に変換した音声データとキーワードの類似度を基準にすることにより、従来の技術では検索できなかったサブワード表現形式に変換したものに変換誤りが含まれる音声データを、検索することができるようになる。

【0047】

(実施形態2)

10

20

30

40

50

サブワード類似度計算部 206 における類似度の計算では、クエリサブワードがサブワード認識結果 212 に対して部分一致する場合に類似度が大きくなるような計算方法を用いても良い。その一例を以下で説明する。

【0048】

実施形態 1 で用いた音節正解精度の計算式 $((N - I - S - D) / N) \times 100$ のかわりに、挿入誤り音節数 I に対して重み α をかけた式 $((N - \alpha \cdot I - S - D) / N) \times 100$ を用い、 α の値を調節することでクエリサブワードが部分一致した場合のサブワード類似度の値を大きくすることができる。

【0049】

クエリサブワードが「は こ ね」、サブワード音声認識結果が「は こ ね や ま」であった場合、実施形態 1 で用いた音節正解精度をサブワード類似度とした場合は、正解音節数 $N = 5$ 、挿入誤り音節数 $I = 2$ (「や」「ま」)、置換誤り音節数 $S = 0$ 、削除誤り音節数 $D = 0$ となり、サブワード類似度は $((5 - 2 - 0 - 0) / 5) \times 100 = 60$ となる。挿入誤り音節数に重みをかけた式 $((5 - \alpha \cdot 2 - 0 - 0) / 5) \times 100$ を用いた場合、 $\alpha = 0.5$ にするとサブワード類似度は 80、 $\alpha = 0$ にするとサブワード類似度は 100 になり、音節正解精度を用いた場合に比べクエリサブワードが部分一致したときのサブワード類似度が大きくなる。このようにクエリサブワードが部分一致したときにサブワード類似度が大きくなるようなサブワード類似度の計算方法を用いることで、データに関連付けられた音声データに部分一致するような検索キーワードを使った検索が実現できる。なお、上式において、 $\alpha = 0$ の時に求まる値は音節正解率、 $\alpha = 1$ とおいた場合は実施形態 1 で用いた音節正解精度になる。

【0050】

また、音節正解精度や音節正解率と別の方法として、例えば音節認識結果から音節同士の間違いやすさを記述したコンフュージョンマトリクスを作成し、これを基にして音節間の距離を定義して、この距離を基にクエリサブワードとサブワード認識結果 212 の類似度を求めるようにしてもよい。言うまでもないが、これらの方法は音節に限らず、音素など他の認識単位を用いた場合でも適用可能である。

【0051】

(実施形態 3)

実施形態 1 では、サブワード認識結果 212 に記憶するサブワード音声認識結果として、402 (図 4) に示すような認識結果の音節列を記憶する場合について説明したが、本発明はこれに限るものではなく、ラティス構造やグラフ構造でサブワード認識結果を表現したものを記憶しても良い。その一例を図 5 に示す。図 5 は、ノードとリンクを用いたグラフ構造で表現した音節音声認識結果である。ノード 501、ノード 505 はそれぞれ認識結果の開始、終了を意味し、音節認識結果を構成する各音節は、501、505 の間のノードで表現されている。502 のノードを例に説明すると、音節名「お」が 503 に、さらに音節「お」の区間で計算された音声認識スコアが 504 に記述されている。開始ノード 501 から終了ノード 505 に至る経路を辿ることで、音節認識結果とその認識スコアを求めることができる。認識スコアは辿った経路上の各ノードに記録されている、その区間の音声認識スコアの和によって求まる。図 5 に示した例では、

「ふぁ お ね や ま あ」、

「ふぁ お ね や ま」、

「ふぁ こ ね や ま あ」、

「ふぁ こ ね や ま」、

の 4 個の音節認識結果が表現されており、それぞれの認識スコアは、

「ふぁ お ね や ま あ」: $50 + 41 + 40 + 50 + 30 + 22 = 233$ 、

「ふぁ お ね や ま」: $50 + 41 + 40 + 50 + 30 = 211$ 、

「ふぁ こ ね や ま あ」: $50 + 38 + 40 + 50 + 30 + 22 = 230$ 、

「ふぁ こ ね や ま」: $50 + 38 + 40 + 50 + 30 = 208$ 、

となる。

【 0 0 5 2 】

このようにサブワード認識結果 2 1 2 がグラフ構造で記憶されている場合、ステップ S 3 0 6 でサブワード類似度の計算を行う際に、グラフが表現する全てのサブワード列をサブワード類似度の計算対象にしても良い。例えば、図 5 で示したグラフ構造では、先に説明した 4 つの音節認識結果からサブワード類似度を求める。このかわりに、グラフが表現するサブワード列のうち認識スコアの高い所定数の候補のみを計算対象にしても良い。例えば、図 5 で示したグラフ構造において、認識スコアが高い上位 3 候補のみをサブワード類似度の計算対象にする場合には、「ふぁ お ね や ま あ」、「ふぁ お ね や ま」、「ふぁ こ ね や ま あ」を用いてサブワード類似度を計算する。あるいは、認識スコアに閾値を設け、閾値以上のサブワード列のみ計算対象にしてもよい。例えば、図 5 で示したグラフ構造において、認識スコアの閾値を 2 3 0 とする場合には、「ふぁ お ね や ま あ」、「ふぁ こ ね や ま あ」からサブワード類似度を計算する。サブワード認識結果がラティス構造で表現されている場合についても、グラフ構造同様にラティスから全てのサブワード列を抽出することにより、本実施形態と同様の処理が適用できることは言うまでもない。

10

【 0 0 5 3 】

(実施形態 4)

上述の実施形態 1 乃至実施形態 3 では、画像データ検索装置を例として、ユーザがキーボードなどで検索用キーワードを入力する態様を説明したが、同様の構成で、キーワードを音声によって入力する装置も実現が可能である。

20

【 0 0 5 4 】

実施形態 1 において、入力されたキーワードをサブワード表現形式のクエリサブワードに変換する処理を、ユーザが音声で入力したキーワードをサブワード表現形式のクエリサブワードに変換する処理に置き換えることで、キーワードの音声入力を実現できる。実施形態 1 で説明した画像データ検索装置は、音声データからサブワード認識結果を求めるサブワード音声認識部 2 0 5 を備えており、このサブワード音声認識部 2 0 5 を利用することにより、キーワードとして入力された音声をサブワード表現形式のクエリサブワードに変換することができる。

【 0 0 5 5 】

以下、この詳細を実施形態 1 と同様に画像データ検索装置を例に説明する。

30

【 0 0 5 6 】

図 6 は、本実施形態における画像データ検索装置のハードウェア構成を示すブロック図である。これは、図 1 に示した実施形態 1 の画像データ検索装置の構成に、マイクロフォンなどの音声を入力するための音声入力装置 1 0 8 が追加された構成である。また、本実施形態における画像データ検索装置の機能構成は、図 2 に示した実施形態 1 の画像データ検索装置と同様の機能構成を有するので、ここでは図 2 を援用する。ただし、実施形態 1 とは異なる処理を行う機能モジュールがあるので以下で説明する。

【 0 0 5 7 】

キーワード取得部 2 0 2 は、音声入力装置 1 0 8 を介して音声で入力されるキーワード（キーワード音声）を取得する。以後の説明では、画像データに関連付けられた音声データと、ユーザによりキーワードとして入力された音声データとを区別するため、前者をこれまでどおり単に「音声データ」と記述し、後者を「キーワード音声」と記述する。

40

【 0 0 5 8 】

キーワード・サブワード変換部 2 0 3 は、取得したキーワード音声をサブワード表現形式に変換したクエリサブワードに変換する。キーワード音声からクエリサブワードへの変換には、サブワード音声認識部 2 0 5 によるサブワード音声認識を用いる。キーワード音声データを入力としてサブワード音声認識を行い、得られたサブワード表現形式の認識結果をクエリサブワードとする。実施形態 1 で説明したように、ここで用いる音声認識は、認識結果としてサブワード表現形式で記述されたものが得られれば良いので、音素タイプライタや音節タイプライタなどサブワード表現形式の認識結果を出力するサブワードを単

50

位とした音声認識を行っても良いし、一般にディクテーションとして知られる大語彙連続音声認識を行って得られた認識結果をサブワードに変換するようにしてもよい。また、実施形態1における音声データのサブワード音声認識と同様に、キーワード音声のサブワード音声認識においても、単一の認識結果だけでなく複数の認識結果候補を求めてクエリサブワードとしてもよい。また、実施形態3で説明したようなラティス構造やグラフ構造で表現されたサブワード認識結果（以下「サブワードグラフ」と記述する。）を求めてクエリサブワードとしてもよい。すなわち、クエリサブワードの表現形式は、単一のサブワード列、複数のサブワード列、サブワードグラフのいずれの形式でも良い。

【0059】

サブワード音声認識部205では、キーワード・サブワード変換部203におけるキーワード音声のサブワード音声認識に加え、実施形態1と同様に画像データに関連付けられた音声データのサブワード音声認識も行うので、本実施形態におけるサブワード音声認識部205では、検索用に入力されたキーワード音声と画像データに関連付けられた音声データの両方に対してサブワード音声認識を行うことになる。

【0060】

サブワード類似度計算部206では、音声データをサブワード認識したサブワード認識結果212とキーワード・サブワード変換部203で求めたクエリサブワードとを比較し、サブワード類似度を計算する。前述したように、クエリサブワードの表現形式は、単一のサブワード列、複数のサブワード列、サブワードグラフのいずれの形式でも良い。同様に、実施形態1、実施形態3で説明したように、音声データのサブワード認識結果212の表現形式も、単一のサブワード列、複数のサブワード列、サブワードグラフのいずれの形式でも良い。

【0061】

クエリサブワードがQ個のサブワード列、サブワード認識結果がN個のサブワード列である場合のサブワード類似度の計算方法の一例を図7のフローチャートを用いて説明する。この方法では、クエリサブワードのq番目のサブワード列を正解とするときの、サブワード認識結果のn番目のサブワード列のサブワード正解精度 $acc(q, n)$ を、 $1 \leq q \leq Q$ 、 $1 \leq n \leq N$ について全て求め、その最大値をサブワード類似度として計算する。

【0062】

なお、図7および以下の説明において記述されている式は、C言語等の記法に従っていることに留意されたい。すなわち、単一の等号「=」は「右辺の値を左辺に代入する」ことを意味し、二重等号「==」は「左辺と右辺の値が等しい」ことを意味する。

【0063】

まず、ステップS701で、サブワード正解精度の最大値を示す変数max、クエリサブワードのサブワード列のインデックスqをそれぞれ0に初期化する。ステップS702で、クエリサブワードのサブワード列のインデックスqの値を1増分するとともに、サブワード認識結果のサブワード列のインデックスnを0に初期化する。続くステップS703では、サブワード認識結果のサブワード列のインデックスnの値を1増分する。

【0064】

次に、クエリサブワードのq番目のサブワードを正解として（ステップS704）、サブワード認識結果のn番目のサブワード列のサブワード正解精度 $acc(q, n)$ を計算する（ステップS705）。 $acc(q, n)$ がサブワード正解精度の最大値maxよりも大きい場合（ステップS706）はmaxの値を当該 $acc(q, n)$ で置き換える（ステップS707）。全てのサブワード認識結果についてサブワード正解精度の計算を終えた場合（ステップS708においてnがNに等しくなった場合）は、ステップS703～S708と同様の処理をインデックスqの値を1増分してクエリサブワードの次のサブワード列について行う。クエリサブワードの全てのサブワード列について、ステップS703～S708の処理を終了したとき（ステップS709においてqがQに等しくなったとき）に求まっているサブワード正解精度の最大値maxをサブワード類似度とする（ステップS710）。なお、実施形態1で示したサブワード類似度の計算方法は、図7のフ

10

20

30

40

50

ローチャートにQを1、Nを3、サブワードを音節として適用した場合に相当する。

【0065】

以上の説明では、サブワード正解精度 $acc(q, n)$ ($1 \leq q \leq Q, 1 \leq n \leq N$) の最大値をサブワード類似度として用いたが、 $acc(q, n)$ に認識スコアなどに基づいて重み付けした値を求めて、その最大値あるいは和をサブワード類似度として用いても良い。また、サブワード正解精度ではなく、実施形態2で説明したような音節正解率や他の尺度を使って計算しても良い。

【0066】

クエリサブワードあるいはサブワード認識結果がサブワードグラフで得られている場合は、実施形態3で説明したようにサブワードグラフから複数のサブワード列を抽出して、抽出した複数（あるいは単数）のサブワード列に対して図7で示した処理を適用することでサブワード類似度を計算することができる。

10

【0067】

キーワードスポッティング部207では、実施形態1と同様に音声データを入力、クエリサブワードを認識対象語とするキーワードスポッティングを行い、認識スコアを求める。クエリサブワードが複数のサブワード列として求まっている場合は、各サブワード列を認識対象語とし、クエリサブワードがサブワードグラフの形式で表される場合はサブワードグラフから抽出した複数のサブワード列を認識対象語とするキーワードスポッティングを行い、各認識対象語に対して求めた認識スコアの最大値を認識スコアとする。

【0068】

20

以上の機能モジュール構成で実現した画像データ検索装置の処理手順を図8に示すフローチャートに示す。図3に示した実施形態1の画像データ検索装置の処理手順との違いは、ステップS304およびS305のかわりにステップS804およびS805が実行される点であるので、この部分のみ説明する。

【0069】

実施形態1では、ステップS304においてキーワード取得部202が入力装置107から入力されたキーワードを取得したのに対し、本実施形態のステップS804では、キーワード取得部202は音声入力装置108から入力されたキーワード音声を取得する。

【0070】

続くステップS805では、ステップS804で取得したキーワード音声をキーワード・サブワード変換部203でサブワード表現形式に変換するという点でステップS305と同一の処理であるが、本実施形態では先に説明したように、キーワード・サブワード変換部203からサブワード音声認識部205を駆動し、キーワード音声を入力とするサブワード音声認識を行い、得られた認識結果からサブワード表現形式のクエリサブワードを得る。

30

【0071】

ステップS306以降の処理については、得られるクエリサブワードのサブワード列が複数になる場合の処理が追加されるが、実施形態1と同様の処理によって検索処理が実行される。

【0072】

40

以上説明したように、本実施形態によれば、実施形態1と同様の構成で、検索用キーワードを音声入力可能なデータ検索装置が実現される。

【0073】

(他の実施形態)

以上、本発明の実施形態を詳述したが、本発明は、複数の機器から構成されるシステムに適用してもよいし、また、一つの機器からなる装置に適用してもよい。

【0074】

なお、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その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

50

る。その場合、プログラムの機能を有していれば、その形態はプログラムである必要はない。

【 0 0 7 5 】

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、その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およびそのプログラムを格納した記憶媒体も本発明を構成することになる。つまり、本発明の特許請求の範囲には、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、およびそのプログラムを格納した記憶媒体も含まれる。

【 0 0 7 6 】

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、プログラムの形態を問わない。

10

【 0 0 7 7 】

プログラムを供給するための記憶媒体としては、例えば、フレキシブルディスク、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。

【 0 0 7 8 】

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、その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憶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、本発明のクレームに含まれるものである。

20

【 0 0 7 9 】

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

【 0 0 8 0 】

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。

【 0 0 8 1 】

さらに、記憶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】

【 0 0 8 2 】

【図1】本発明の実施形態における画像データ検索装置のハードウェア構成を示すブロック図である。

【図2】本発明の実施形態における画像データ検索装置の機能構成を示すブロック図である。

【図3】本発明の実施形態における画像データ検索装置によるデータ検索処理を示すフローチャートである。

50

【図 4】本発明の実施形態におけるサブワード認識結果の一例を示す図である。

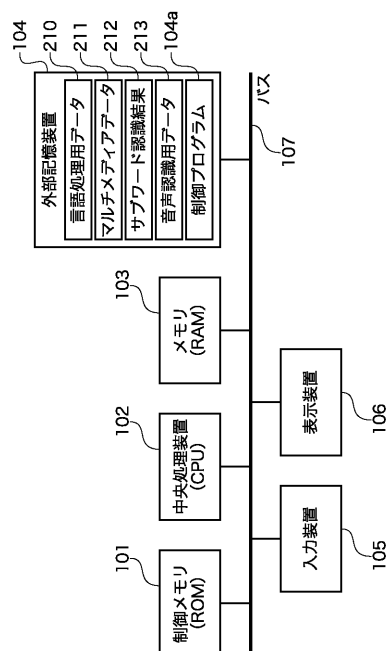
【図５】本発明の実施形態におけるグラフ構造で表現されたサブワード認識結果の一例を示す図である。

【図 6】本発明の実施形態における、キーワードを音声によって入力するタイプの画像データ検索装置のハードウェア構成を示すブロック図である。

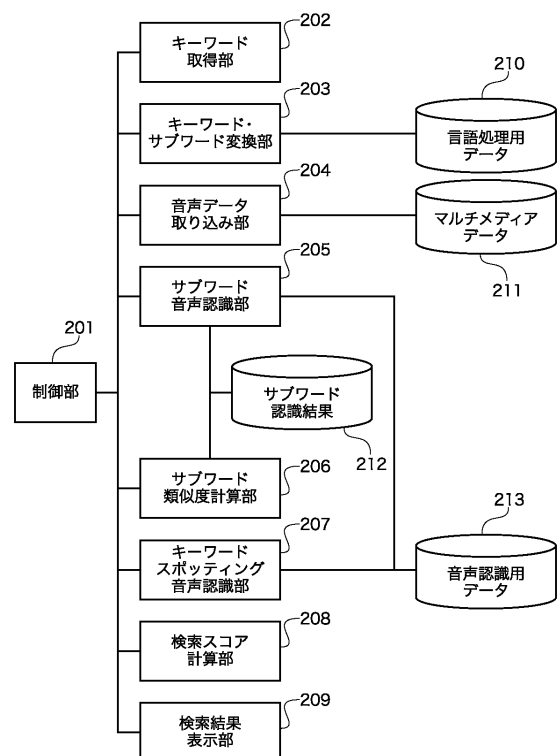
【図 7】本発明の実施形態におけるキーワードを音声によって入力するタイプの画像データ検索装置によるサブワード類似度の計算処理を示すフローチャートである。

【図 8】本発明の実施形態におけるキーワードを音声によって入力するタイプの画像データ検索装置による、データ検索処理を示すフローチャートである。

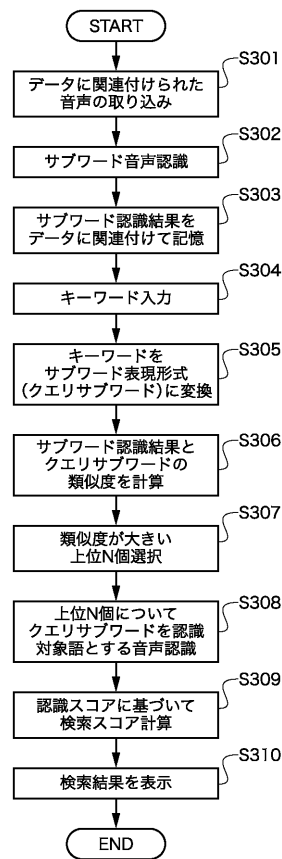
【圖 1】



【圖 2】



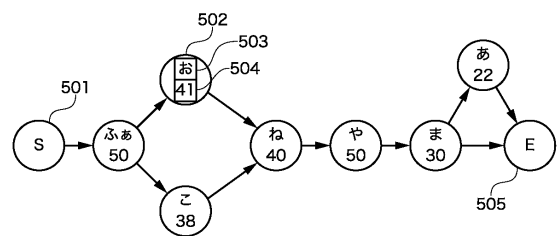
【図 3】



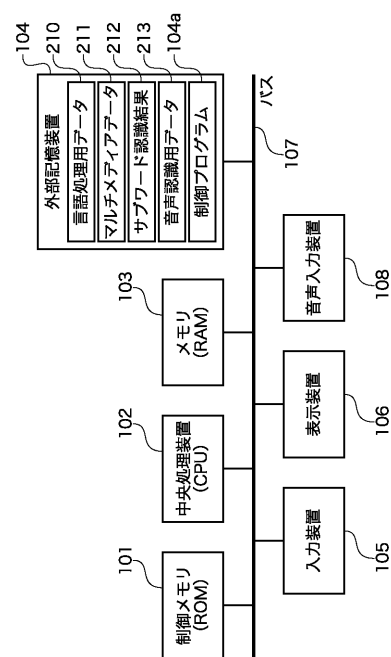
【図 4】

順位	認識結果の音節列	認識スコア
1	ふぁおねやまあ	233
2	ふぁおねやま	211
3	ふぁこねやま	208

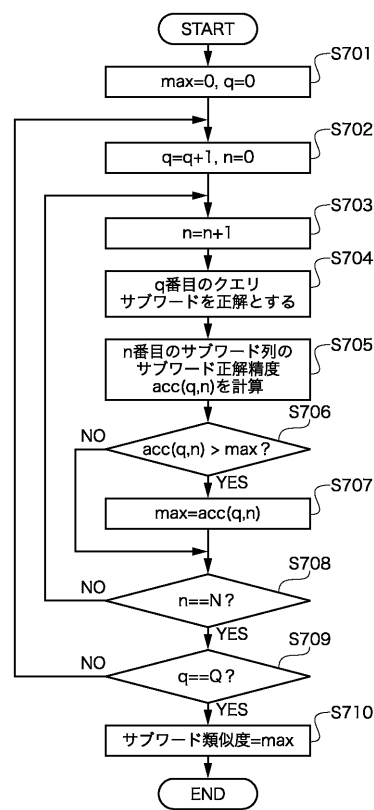
【図 5】



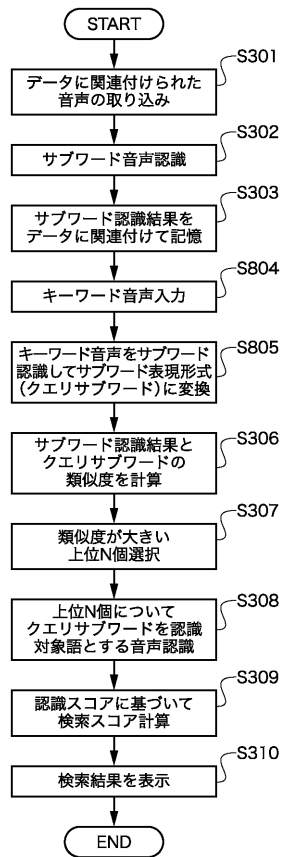
【図 6】



【図 7】



【図 8】



フロントページの続き

- (72)発明者 小森 康弘
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
- (72)発明者 山田 耕平
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

審査官 井上 健一

- (56)参考文献 特開2002-278579(JP,A)
特開2004-302175(JP,A)
特開2000-259645(JP,A)
特開平11-085187(JP,A)
特開昭63-239499(JP,A)
特開平10-173769(JP,A)
特開2003-219327(JP,A)
特開2006-040150(JP,A)
特開平08-211893(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G10L 15/00 - 15/28