

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5866728号
(P5866728)

(45) 発行日 平成28年2月17日(2016.2.17)

(24) 登録日 平成28年1月15日(2016.1.15)

(51) Int.Cl.

F I

G 0 6 F 13/00 (2006.01)

G 0 6 F 13/00 6 5 0 B

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 1 1 0 C

G 0 6 Q 50/00 (2012.01)

G 0 6 F 17/30 3 1 0 Z

G 1 0 L 15/28 (2013.01)

G 0 6 F 17/30 2 3 0 Z

G 1 0 L 15/00 (2013.01)

G 0 6 F 17/60

請求項の数 15 (全 73 頁) 最終頁に続く

(21) 出願番号 特願2011-226792 (P2011-226792)
 (22) 出願日 平成23年10月14日(2011.10.14)
 (65) 公開番号 特開2013-88906 (P2013-88906A)
 (43) 公開日 平成25年5月13日(2013.5.13)
 審査請求日 平成26年9月10日(2014.9.10)

(73) 特許権者 510003368
 サイバーアイ・エンタテインメント株式会
 社
 東京都世田谷区瀬田 1-5-1
 (74) 代理人 100092093
 弁理士 辻居 幸一
 (74) 代理人 100082005
 弁理士 熊倉 禎男
 (74) 代理人 100067013
 弁理士 大塚 文昭
 (74) 代理人 100086771
 弁理士 西島 孝喜
 (74) 代理人 100109070
 弁理士 須田 洋之

最終頁に続く

(54) 【発明の名称】 画像認識システムを備えた知識情報処理サーバシステム

(57) 【特許請求の範囲】

【請求項 1】

通信システムであって、
 サーバ装置と、

第 1 の画像と、当該第 1 の画像に関連付けされている第 1 のメッセージと、当該第 1 の画像に関連付けされている場所の情報を少なくとも含む第 1 の情報とをネットワークを介して前記サーバ装置に送信する第 1 の装置と、

前記ネットワークを介して前記サーバ装置に接続されている第 2 の装置と、
 を備え、

前記サーバ装置は、前記第 1 の画像に含まれている 1 以上の物体を特定し、前記第 1 のメッセージから 1 以上の単語を抽出し、当該 1 以上の単語を分析することによって前記第 1 の装置の第 1 のユーザが着目している物体を前記 1 以上の物体から特定し、前記第 1 のメッセージを当該着目物体に関連付けするよう構成され、

前記サーバ装置は、前記第 1 の画像と、前記第 1 のメッセージと、当該第 1 のメッセージが当該第 1 の画像における前記着目物体に関連付けされていることを示す情報とを前記第 2 の装置に前記ネットワークを介して送信するよう構成されている、通信システム。

【請求項 2】

前記第 2 の装置は、当該第 2 の装置に関連付けされている場所の情報を少なくとも含む第 2 の情報を前記サーバ装置に送信するよう構成され、

前記サーバ装置は、前記第 1 の情報及び前記第 2 の情報に基づいて、前記第 1 の画像と

10

20

、前記第 1 のメッセージと、当該第 1 のメッセージが当該第 1 の画像における前記着目物体に関連付けされていることを示す情報とを前記第 2 の装置に前記ネットワークを介して送信することを決定するよう構成されている、請求項 1 に記載の通信システム。

【請求項 3】

前記第 1 の装置は、前記第 1 の画像を送信した後に、前記第 1 のメッセージを前記サーバ装置に送信するよう構成されている、請求項 1 又は 2 に記載の通信システム。

【請求項 4】

前記第 2 の装置は、第 2 のメッセージを前記ネットワークを介して前記サーバ装置に送信するよう構成されている、請求項 1 ないし 3 のいずれか一つに記載の通信システム。

【請求項 5】

前記サーバ装置は、前記第 1 のメッセージ及び第 2 のメッセージを分析し、ユーザ間のインタレストグラフを取得するよう構成されている、請求項 4 に記載の通信システム。

【請求項 6】

前記第 1 の情報は、第 1 の時間情報を更に含み、前記サーバ装置は、前記着目物体を当該第 1 の時間情報に関連付けるよう構成されている、請求項 1 ないし 5 のいずれか一つに記載の通信システム。

【請求項 7】

前記サーバ装置は、前記第 1 の時間情報及び前記第 1 の画像を少なくとも用いて、アルバムを生成するよう構成されている、請求項 6 に記載の通信システム。

【請求項 8】

前記第 1 の装置及び / 又は前記第 2 の装置は、文字情報の書込み及び / 又はユーザの声による語りかけによって、前記第 1 のメッセージ及び / 又は第 2 のメッセージを入力するよう構成されている、請求項 1 ないし 7 のいずれか一つに記載の通信システム。

【請求項 9】

前記第 1 の装置及び / 又は前記第 2 の装置は、カメラ付き携帯電話を含む、請求項 1 ないし 8 のいずれか一つに記載の通信システム。

【請求項 10】

前記第 1 の装置及び / 又は前記第 2 の装置は、一以上のマイクロフォン、一以上のイヤフォン、及び一以上の画像撮像素子（カメラ）を少なくとも有するヘッドセットと、当該ヘッドセットに接続されているネットワーク端末とを含み、当該ネットワーク端末は、前記ネットワークを介して前記サーバ装置に接続されている、請求項 1 ないし 9 のいずれか一つに記載の通信システム。

【請求項 11】

前記ヘッドセットは、2 台以上の撮像視差を有するカメラ及び / 又は対象物体までの深度（距離）を測定可能な三次元カメラを含む、請求項 10 に記載の通信システム。

【請求項 12】

前記第 1 の装置及び / 又は前記第 2 の装置は、さらに、生体認証（バイオメトリクス）センサを含み、これにより、ユーザ固有の生体識別情報を生体認証装置に問い合わせるよう構成されている、請求項 1 ないし 11 のいずれか一つに記載の通信システム。

【請求項 13】

前記第 1 の装置、前記第 2 の装置及び / 又は前記サーバ装置は、前記ヘッドセットの着脱を監視するよう構成されている、請求項 12 に記載の通信システム。

【請求項 14】

前記第 1 の装置及び / 又は前記第 2 の装置は、さらに、生体情報（バイタルサイン）センサを含み、これにより、当該生体情報を前記サーバ装置に送信するよう構成されている、請求項 1 ないし 13 のいずれか一つに記載の通信システム。

【請求項 15】

サーバ装置であって、

第 1 の装置から、第 1 の画像と、当該第 1 の画像に関連付けされている第 1 のメッセージと、当該第 1 の画像に関連付けされている場所の情報を少なくとも含む第 1 の情報とを

10

20

30

40

50

ネットワークを介して受信し、

前記サーバ装置は、前記第１の画像に含まれている１以上の物体を特定し、前記第１のメッセージから１以上の単語を抽出し、当該１以上の単語を分析することによって前記第１の装置の第１のユーザが着目している物体を前記１以上の物体から特定し、前記第１のメッセージを当該着目物体に関連付けし、

前記第１の画像と、前記第１のメッセージと、当該第１のメッセージが当該第１の画像における前記着目物体に関連付けされていることを示す情報とを第２の装置に前記ネットワークを介して送信するよう構成されている、サーバ装置。

【発明の詳細な説明】

【技術分野】

10

【０００１】

本発明は、ユーザの頭部に装着可能なヘッドセットシステムに組み込まれたカメラから得られる当該ユーザの主観的な視野を反映した画像信号を、当該ユーザのネットワーク端末経由でネットワークを介して画像認識システムを備えた知識情報処理サーバシステム側に適宜アップロードする事で、当該ユーザが関心を持って着目した特定物体、一般物体、人、写真、或いはシーン等の１以上の対象（以降「対象」と呼称）が、上記カメラ映像中のいずれに当るのかを、前記サーバシステムと当該ユーザ間の音声による双方向のコミュニケーションにより抽出可能にした上で、それら対象の抽出過程及び画像認識結果を、上記サーバシステム側が当該ユーザのネットワーク端末経由で、上記ヘッドセットシステムに組み込まれたイヤフォンを通し、当該ユーザに対し音声情報により通知する事を特徴とする。

20

【０００２】

その上で当該ユーザが着目する様々な対象に対し、当該ユーザの音声によるメッセージやつぶやき、或いは質問等の音声タグを残す事を可能にする事で、異なる時空間内において自らを含む様々なユーザが当該対象に偶然遭遇する、或いはそれら対象を偶然目にした時に、前記サーバシステム側に蓄積された当該対象に係る様々なメッセージやつぶやき群を、当該対象への着目に同期して音声で受取る事を可能にし、それら個々のメッセージやつぶやきに対し、ユーザがさらなる音声応答を返す事を可能にする事で、様々なユーザの共通の着目対象に係る広範なソーシャル・コミュニケーションを喚起する事を特徴とする。

30

【０００３】

その上で、当該喚起された多数のユーザの視覚的関心に端を発する広範なソーシャル・コミュニケーションを、前記サーバシステム側で継続的に収集・解析・蓄積する事で、広範なユーザ、様々なキーワード、及び様々な対象を構成ノード群とする、動的なインタレストグラフとして獲得可能にし、それらを基に高度にカスタマイズされたサービスの提供、精度の高いリコメンデーションの提示、或いは動的な広告や告知等への効果的な情報提供サービスに繋げる事を可能にする、前記画像認識システムを備えた知識情報処理サーバシステムに関する。

【背景技術】

【０００４】

40

近年のインターネットの世界的な普及により、ネットワーク上の情報量が急激に増大しつつある事から、それら膨大な量の情報の海の中から目的とする情報を効果的且つ高速に探し出す手段としての検索技術が急速に進歩して来た。現在では、強力な検索エンジンを備えたポータルサイトがいくつも運営されている。また、閲覧者の検索キーワードやアクセス履歴等を解析し、閲覧者の嗜好にあったWebページや広告等を各々の検索結果に関連して配信する技術も開発され、閲覧者が多用するキーワードに基づく効果的なマーケティング活動等への応用も始まっている。

【０００５】

例えば、ユーザにとって有用な情報を精度良く且つ容易に提供する事が出来る情報提供装置がある（特許文献１）。この情報提供装置は、ユーザによる各コンテンツに対するア

50

クセスの頻度を表すアクセス頻度情報を、当該ユーザを識別するユーザ識別情報に対応付けて格納するアクセス履歴格納手段と、各ユーザ間における各コンテンツへのアクセス傾向の類似性を表すユーザ間類似度を、前記アクセス履歴格納手段に格納された前記アクセス頻度情報に基づいて算出するユーザ間類似度計算手段と、ユーザと各ユーザとの間の前記ユーザ間類似度により重み付けした、当該各ユーザの前記アクセス頻度情報から、当該ユーザにとってのコンテンツの有用度を表す情報であるコンテンツ・スコアを算出するコンテンツ・スコア計算手段と、前記コンテンツ・スコア計算手段によって算出された各コンテンツの前記コンテンツ・スコアを、前記ユーザ識別情報に対応付けて記憶するインデックス格納手段と、通信端末装置から送信されたユーザ識別情報を含むクエリの入力を受け付けるクエリ入力手段と、前記クエリ入力手段により受け付けられた前記クエリに適合するコンテンツのコンテンツ識別情報を取得し、当該クエリに含まれるユーザ識別情報に対応付けられて前記インデックス格納手段に記憶された前記コンテンツ・スコアを参照して、取得した前記コンテンツ識別情報から提供情報を生成する提供情報生成手段と、前記提供情報生成手段により生成された前記提供情報を、前記通信端末装置に出力する提供情報出力手段とを備える事を特徴とする、情報提供装置である。

10

【 0 0 0 6 】

これらのキーワード等の文字情報を検索クエリとする検索手段をさらに拡大する目的で、画像認識技術を備えた検索エンジンの開発が近年進み、文字に代わり画像そのものを入力クエリとする画像検索サービスが、広くインターネット上で提供されるようになって来ている。画像認識技術の研究の始まりは、一般に40年以上前に遡る事が出来る。以来、コンピュータの高速化と機械学習技術の進歩と共に、線画解釈(1970年代)、人手によるルールや幾何形状モデルによって構築された知識データベースに基づく認知モデル、3次元モデル表現(1980年代)といった研究が漸次行われる様になった。1990年代に入ると、特に顔画像の認識や学習による認識に関する研究が盛んになった。2000年代になると、コンピュータの処理能力の一層の向上により、統計処理や機械学習の為に必要となる膨大な計算処理が比較的安価に実行可能になった為、一般物体認識に関する研究が進んだ。一般物体認識とは、実世界のシーンを撮影した画像に対して、コンピュータがその画像中に含まれる物体を一般的な名称で認識する技術である。1980年代には、全て人手によってルールやモデルの構築を試みていたが、大量のデータを手軽に扱える様になったこの時期には、コンピュータを活用した統計的機械学習によるアプローチが注目され、近年の一般物体認識ブームのきっかけとなった。一般物体認識技術によって、画像に対するキーワードを対象画像に自動的に付与する事が可能になり、画像をその意味内容によって分類及び検索する事も可能になる。近い将来には、コンピュータによって全ての人間の画像認識機能を実現する事が目標とされている(非特許文献1)。一般物体認識技術は、画像データベースからのアプローチと統計的確率手法の導入によって急速に進歩した。その中でも先駆的な研究として、画像に人手でキーワードを付与したデータから個々の画像との対応付けを学習し物体認識を行なう手法(非特許文献2)や、局所特徴量に基づく手法(非特許文献3)等がある。また、局所特徴量による特定物体認識に関する研究にSIFT法(非特許文献4)、及びVideo Google(非特許文献5)等がある。その後、2004年に入り、「Bag-of-Keypoints」あるいは「Bag-of-Features」と呼ばれる手法が発表された。この手法は、対象となる画像をビジュアル・ワード(visual word)と呼ばれる代表的な局所パターン画像の集合として扱い、その出現頻度を多次元のヒストグラムで表現する。具体的には、SIFT法に基づいた特徴点抽出を行い、予め求められた複数のビジュアル・ワードに基づいてSIFT特徴ベクトルをベクトル量子化し、画像毎にヒストグラムを生成するものである。この様に生成されたヒストグラムの次元数は、通常、数百から数千次元のスパース(sparse)なベクトルになる。そして、これらのベクトルは、コンピュータ上の多次元ベクトルの分類問題として高速に処理される事により、一連の画像認識処理が行われる(非特許文献6)。

20

30

40

【 0 0 0 7 】

50

これらコンピュータによる画像認識技術の進展に伴い、カメラ付きネットワーク端末で撮影した画像を、ネットワーク経由でサーバ側に構築された画像認識システム側に問い合わせ、当該サーバ側に蓄積された膨大な画像データベースを基に、当該画像認識システム側がそれらの画像と、予め学習済みの物体毎の特徴を記述した画像特徴データベース群とを比較照合する事で、アップロードされた画像に含まれる主要な物体を画像認識し、その認識結果を前記ネットワーク端末側に速やかに提示するサービスが既に始まっている。画像認識技術の中でも特定の人間の顔の検出技術は、個々人を特定する手法の一つとして急速に応用開発が進んでいる。多数の顔画像の中から特定の人物の顔を精度良く抽出する為には、膨大な顔画像の事前学習が必要となる。その為に準備しなくてはならない知識データベースの量も極めて大きくなる事から、或る程度大規模な画像認識システムの導入が必要になる。一方、電子カメラにおけるオートフォーカスに用いられる様な一般的な「平均顔」の検出、或いは限られた人物の顔の特定であれば、電子カメラ等の小型の筐体内に十分収まる規模のシステムで今や容易に実現が可能である。また、近年供用が始まったインターネットを利用した地図提供サービスの中で、地図上の要所々々における路上写真（Street View）を居ながらにして俯瞰する事が出来る様になった。この様なアプリケーションでは、プライバシー保護の観点から偶然写り込んだ自動車のナンバープレートや歩行者の顔、或いは道路越しに垣間見えてしまう個人宅の様子等を、一定以上判別出来ない程度にフィルタ処理して再表示する必要性も出て来ている（非特許文献7）。

【0008】

近年、現実空間を拡張して、コンピュータによる情報空間としてのサイバー空間とを相互に融合しようとする拡張現実感（Augmented Reality：略称AR）というコンセプトが提案され、既に一部のサービスが始まっている。一例として、GPSや無線基地局等から取得可能な位置情報を利用した三次元位置測位システム、カメラ、及び表示装置等を一体として備えたネットワーク携帯端末を用い、上記三次元位置測位システムから割り出した自身の位置情報を基に、カメラで撮影した現実世界の映像と、サーバ上にデジタル情報として蓄積されている注釈（アノテーション：Annotation）とを重ね合わせ、サイバー空間に浮かぶエアタグ（Air tag）として現実世界の映像に貼り付ける事が可能になっている（非特許文献8）。

【0009】

1990年代後半になると、通信ネットワーク・インフラの整備拡張に伴い、インターネット上に構築されたユーザ相互の社会的関係を促進する目的で、ソーシャルネットワークに係るサイトが数多く開設され、数々のソーシャル・ネットワーキング・サービス（SNS）が生まれた。SNSにおいては、ユーザ検索機能、メッセージ送受信機能、掲示板等のコミュニティ機能によって、ユーザ間のコミュニケーションが有機的に促進される。例えばSNSのユーザは、趣味・嗜好を同じくするユーザが集う掲示板に積極的に参加して、文書や画像、音声等のパーソナル情報を交換し、また自分の友人を他の知人に紹介する事等により、人と人との相互の繋がりをさらに深め、ネットワーク上でコミュニケーションを有機的かつより広範に広げていく事が出来る。

【0010】

SNSにおけるサービスの一形態として、ネットワーク上にアップロードされた動画を複数のユーザが選択共有し、当該動画シーン上の任意の位置にユーザが自由に当該動画内容に関連するコメントをアップロードする事を可能にし、それらコメント群を当該動画画面上にスクロール表示する事で、複数のユーザ間で当該動画を媒介とした共有コミュニケーションを図る事が可能なコメント付き動画配信システムがある（特許文献2）。当該システムは、コメント情報をコメント配信サーバから受信し当該共有動画の再生を開始すると共に、当該コメント情報から再生する動画の、特定の動画再生時間に対応するコメントをコメント配信サーバから読み出し、読み出したコメント群に対応付けられた動画再生時間に、当該動画と共にそれらコメント群を表示可能にする。併せて、それらコメント情報をリストとしても個別に表示可能にし、表示されたコメント情報から特定のコメントデータが選択されると、選択されたコメントデータのコメント付与時間に対応する動画再生時間

10

20

30

40

50

から当該動画を再生し、読み出したコメントデータを表示部に再表示させる。また、ユーザによるコメントの入力操作を受け付けて、コメントが入力された時点の動画再生時間をコメント付与時間として、コメント内容と共に前記コメント配信サーバに送信する。

【 0 0 1 1 】

S N Sの中でも、ネットワーク上で交換可能な情報パケットサイズを大幅に限定する事で、コミュニケーションのリアルタイム性をより重視しようという動きもある。これらマイクロブログとも呼ばれるユーザの短いつぶやきや、それらに関連するURL等のアドレス情報を埋め込んだ140文字以内の文字データを、当該ユーザがインターネット上にリアルタイム且つ広範に発信する事で、当該ユーザのその時々体験を当該ユーザの文字によるつぶやきのみならず、画像や音声データを加えた一体的な情報として広範なユーザ間で共有可能にし、さらにユーザがそれらつぶやきの中から特定の発信者や特定の話題を選択してフォローする機能も提供する事で、地球規模でのリアルタイム・コミュニケーションを喚起するサービスが既に始まっている（非特許文献9）。

【 0 0 1 2 】

ネットワークを介した情報サービスとは異なるものの、特定の対象に対峙した時に当該対象に関する詳細な音声説明を受取る事が出来るサービスとして、博物館や美術館の「音声ガイド」システムがある。これらは、対象となる絵画等の近傍に設置された音声信号送出部から送出される赤外線変調された音声信号を、それら対象物に近接したユーザの端末装置に組込まれた赤外線受信部で復調し、当該ユーザのイヤフォンに当該絵画等に係る詳細な説明を音声として提供するもので、この方式以外にも極めて指向性の高い音声トランスミッターを用いて、ユーザの耳元に直接当該音声情報を送り込める様な音声ガイドシステムも実用化されている。

【 0 0 1 3 】

コンピュータ・システムに対する音声による情報入力やコマンド入力方法として、ユーザの発話音声を音声言語として認識し、テキストデータや各種のコンピュータコマンドに変換して入力処理する技術がある。当該入力処理には高速の音声認識処理が必要となるが、これらを可能にする音声認識技術群として、音響処理技術、音響モデル作成・適応化技術、適合・尤度演算技術、言語モデル技術、対話処理技術等があり、これらの要素技術をコンピュータ上で組み合わせる事で、近年では十分実用に耐える高速の音声認識システムが構築可能となっている。近年では、大規模語彙連続音声認識エンジンの開発によって、ユーザにより発話される音声言語認識処理を、ネットワーク端末上でほぼ実時間で処理する事も可能となっている。

【 0 0 1 4 】

音声認識技術の研究の歴史は、1952年に米国のベル研究所でのゼロ交差回数を用いた数字認識の研究に始まり、1970年代に入ると発声時間の長さの変動を、動的計画法を用いて非線形に正規化する手法（Dynamic Time Warping）が日本及びロシアの研究者によって提案され、米国においても統計確率的手法であるHMM（Hidden Markov Model：隠れマルコフモデル）を用いた音声認識の基礎的な研究が進んだ。現在では、利用者の音声の特徴を適応的に学習させる事より、明瞭な発声で読み上げられた文章をほぼ完全に口述筆記する事が可能なレベルにまで到達している。この様な高度の音声認識技術を応用した従来技術として、会議による発言音声を入力とする話し言葉から、文語としての議事録を自動作成する技術も開発されている（特許文献3）。

【 0 0 1 5 】

すわなち、特許文献3に開示された技術は、音声を入力して文書情報を作成し出力する音声文書変換装置であり、文書情報出力を受信して画面に表示する表示装置を備え、この音声文書変換装置が、入力する音声を認識する音声認識部と、入力音声を漢字仮名混じりの文語に変換する変換テーブルと、前記音声認識部から認識した音声を受信して整列させ前記変換テーブルを検索して文語に変換し所定の書式で文書に編集する文書形成部と、この編集済み文書を記憶保存する文書メモリと、この保存された文書情報を送信すると共に

10

20

30

40

50

他の情報・信号を前記表示装置との間で授受する送受信部とを有し、かつ前記表示装置が前記音声文書変換装置の送受信部との間で情報・信号を送受信する送受信部と、受信した文書情報を表示情報として記憶する表示情報メモリと、この記憶する表示情報を画面表示する表示盤とを有する事を特徴としている。

【0016】

また、コンピュータ上の文字情報からなる文章を、指定された言語で流暢に読み上げる音声合成システムは、近年最も進化の進んでいる領域の一つである。音声合成システムは、スピーチ・シンセサイザー (Speech Synthesizer) と呼ばれ、テキストを音声に変換するテキスト読み上げシステムや、発音記号を音声に変換するシステム等を含む。歴史的には、1960年代末以降、コンピュータによる音声合成システムの開発が進んだものの、初期のスピーチ・シンセサイザーによる発声はいかにもコンピュータによる音声だと感じさせる人間味のない無機質なものが多かった。以降研究が進むにつれ、後述する様に、場面、状況、前後の文脈関係により声の抑揚や調子を自在に変化させる事が出来る様になり、人間の肉声と比べてほとんど遜色がない高品質の音声合成が可能になっている。特に、サーバ側に構築された音声合成システムは、膨大な辞書を活用可能なばかりではなく、その発声アルゴリズム自体も人間に近い複雑な発音が可能な様に多数のデジタルフィルタ類を組み込む事も可能になり、ネットワーク端末機器の急速な普及に伴い、近年その応用可能な範囲が一段と拡大している。

10

【0017】

音声合成技術には、大きく分けてフォルマント合成と連結的合成とがある。フォルマント合成では、人間の音声を使用する事なく周波数や音色等のパラメータをコンピュータ上で調整して人工的な合成波形を生成する。これらは一般的に人工的な音声として聞こえる場合が多い。一方で連結的合成では、基本的に人間の音声を収録して、その音素断片等を滑らかに連結して肉声に近い音声を合成する方法である。具体的には、一定時間収録された音声を「音」「音節」「形態素」「単語」「成句」「文節」等に分割してインデックス化し、検索可能な音声ライブラリ群を作成する。こうした音声ライブラリは、テキスト読み上げシステム等により音声を合成する際に、適宜最適な音素や音節等が抽出され、適切なアクセントと共に最終的に人間の発話に近い流暢な一連の音声に変換される。

20

【0018】

係る従来技術に加え、声調機能を備えたテキスト読み上げシステム等の開発により、バリエーションに富んだ音声を合成する技術も続々実用化されている。例えば、高度な音声編成システムによって、アクセント調整や音の高低・長さの調整を行う事によって、「うれしさを伴った声」「悲しみを伴った声」「怒りを伴った声」「冷たさを伴った声」等の感情の抑揚を調整する事が出来る他、音声編成システムが備えるデータベースに登録された特定の人のクセを反映した音声を、これらシステム上で自在に合成する事も出来る様になっている。

30

【0019】

また、上述した音声合成についての先行技術に、合成音声区間と部分的に一致する肉声区間を検出して、その肉声区間の韻律 (抑揚・リズム) 情報を合成音声に付与し、肉声と合成音声を自然に結合させる技術も提案されている (特許文献4)。

40

【0020】

即ち、特許文献4に開示された技術は、録音音声格納手段、入力テキスト解析手段、録音音声選択手段、接続境界算出手段、規則合成手段、接続合成手段に加えて、合成音声区間のうちで録音済みの肉声と部分的に一致する区間を決定する肉声韻律区間決定手段と、その一致部分の肉声韻律を抽出する肉声韻律抽出手段と、抽出された肉声韻律を使って合成音声区間全体の韻律情報を生成する、ハイブリッド韻律生成手段を備える事を特徴としている。

【先行技術文献】

【特許文献】

【0021】

50

【特許文献1】特開2009-265754号公報

【特許文献2】特開2009-077443号公報

【特許文献3】特開1993-012246号公報

【特許文献4】特開2009-020264号公報

【非特許文献】

【0022】

【非特許文献1】柳井啓司, "一般物体認識の現状と今後", 情報処理学会論文誌, Vol.48, No.SIG16(CVIM19), pp.1-24, 2007

【非特許文献2】Pinar Duygulu, Kobus Barnard, Nando de Freitas, David Forsyth, "Object Recognition as Machine Translation: Learning a lexicon for a fixed image vocabulary," European Conference on Computer Vision (ECCV), pp.97-112, 2002.

【非特許文献3】R. Fergus, P. Perona, and A. Zisserman, "Object Class Recognition by Unsupervised Scale-invariant Learning," IEEE Conf. on Computer Vision and Pattern Recognition, pp.264-271, 2003.

【非特許文献4】David G. Lowe, "Object Recognition from Local Scale-Invariant Features," Proc. IEEE International Conference on Computer Vision, pp.1150-1157, 1999.

【非特許文献5】J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos", Proc. ICCV2003, Vol. 2, pp.1470-1477, 2003.

【非特許文献6】G. Csurka, C. Bray, C. Dance, and L. Fan, "Visual categorization with bags of keypoints," Proc. ECCV Workshop on Statistical Learning in Computer Vision, pp.1-22, 2004.

【非特許文献7】Ming Zhao, Jay Yagnik, Hartwig Adam, David Bau; Google Inc. "Large scale learning and recognition of faces in web videos" FG '08: 8th IEEE International Conference on Automatic Face & Gesture Recognition, 2008.

【非特許文献8】<http://jp.techcrunch.com/archives/20091221sekai-camera/>

【非特許文献9】Akshay Java, Xiaodan Song, Tim Finin, and Belle Tseng, "Why We Twitter: Understanding Microblogging Usage and Communities" Joint 9th WEBKDD and 1st SNA-KDD Workshop '07.

【発明の概要】

【発明が解決しようとする課題】

【0023】

しかしながら、従来の検索エンジンにおいては、検索対象に係るいくつかのキーワードを考え文字で入力する必要があった。それらの検索結果は、複数、時に夥しい数の候補群に係る文書タイトルと共に概略記述文章として提示される事から、目的とする検索結果に辿り着く為には、各候補群が示す情報の格納先をさらに個々に開いて読み進んでいく必要があった。近年は画像を直接入力クエリとする検索も可能になり、その検索出力として関連度の高い画像そのものを一覽的に閲覧可能な画像検索サービスも提供され始めている。しかし、ユーザが関心を持った対象や事象に対し、その好奇心をさらに喚起する様な関連情報を、快適且つ的確にユーザに提供出来る迄には至っていない。また従来の検索プロセスでは、PCやネットワーク端末等に向かって一時的ではあるにせよ集中的な入力操作を行う必要がある事から、ユーザがハンズフリーで何か別の事をしながら日常生活の中でふと誰かに語りかけ、身近な誰かが答えてくれる様な、普段我々が何気なく行っている自然なコミュニケーションが、従来のITシステム上ではまだ実現出来ていない。

【0024】

一例として、ユーザがふと調べたいと思った対象や事象に遭遇した場合、その名称等が判る場合には文字入力によるネットワーク検索を行うか、カメラ付き携帯電話やスマートフォン等を手に当該対象に近付き、当該ネットワーク端末に具備されているカメラで撮影した後、当該撮影画像を基に画像検索をかけるケースが多い。それでも思う様な検索結果が得られない場合は、ネットワーク上の他のユーザへ当該対象を問い合わせる事も可能で

10

20

30

40

50

はある。しかし、これら一連のプロセスはユーザにとって少々煩雑であるだけでなく、対象に直接携帯電話等をかざす等の行為が必要な事から時に対象から身構えられる、場合によっては失礼だと感じさせる、さらには携帯電話をかざす行為自体に対し周りから不審な目で見られる、といった嫌いがあった。また対象が動物や人物等の場合、対象と自分との間にカメラ付き携帯端末等が入る事により一種の視覚的な壁の様なものが出来てしまう点と、検索結果を先ずは当該携帯端末で確認しようとする事から、一時的にせよ当該対象や周囲の人々とのコミュニケーションが中断しがちであった。また、これら一連の検索プロセスには相応の時間がかかる事から、ユーザが外出中にふと目にした物体や人、動物、或いはシーン等に関心を持ったとしても、その場で上記一連の操作が完結出来ない場合も多く、一旦撮影した写真を自宅等に持ち帰って改めてＰＣ等で検索し直す必要もあった。

10

【 0 0 2 5 】

近年、実用化が始まった拡張現実と呼ばれるサービスにおいて、我々が存在する現実の空間と、コンピュータネットワーク網の中に構成されるサイバー空間とを紐付ける手法の一つとして、ＧＰＳ等から得られる測位情報に加え、カメラが向いている方位情報を併せて利用する手法がある。しかしこれら位置情報のみの利用では、対象物体自体の移動や、そもそも対象が観測時点で存在していない等、刻々と変化する現実の世界の状況に際し対応が困難なケースが多い。基本的に位置情報と固定的に紐付いている様々な建造物や都市のランドマーク等とは異なり、車などの移動・可搬可能な物体や、動き回る人や動物、或いは「夕焼け」等の概念的なシーンに対しては、当該システム内に画像認識機能を有していない場合には、本質的な意味での相互の対応付けが困難となる。

20

【 0 0 2 6 】

ＳＮＳにおけるサービスの一形態として、近年ユーザの間で人気のあるコメント付き動画共有サービスにおいては、共有視聴される動画が録画済みの動画である場合には、現実の世界で進行中の事象に対してリアルタイムの共有体験が得られないという問題がある。これに対し、ライブストリーム映像配信に対応したコメント付与サービスが既に始まっている。対象となるストリーム映像としては、記者会見、発表会、国会中継、イベント、スポーツ等に加えて、一般ユーザの投稿によるライブ映像配信がある。これらの動画共有サービスにおいては、ネットワークを介してリアルタイムで進行中の事象に係る「場」の共有が可能となる。しかし、延々と続くライブストリーム映像配信をフォローするには時間及び忍耐が必要である。そこからユーザ固有の或いは参加しているユーザ群に共通の関心の在り所等を効果的・効率的に抽出し、それらをインタレストグラフとして広範に体系付ける素材群として見ると、その収集可能な対象及び情報量には一定の限界があった。これは利用者数が急増しているネットワーク共有動画視聴サービスでも同じで、ユーザが様々な動画ファイルを連続視聴する為に消費する時間、及び配信サーバやネットワーク回線に係るコストに対し、ユーザが能動的に何か有用な情報をサーバ側に提供出来るチャンスはそれ程多くない。

30

【 0 0 2 7 】

これに対し、１４０文字以内という一定の制限は課されるものの、そのネットワーク上を流れるリアルタイムのトピックスの多彩さと参加者の急増も手伝って、これらマイクロブログと呼ばれるリアルタイム・メッセージ交換サービスから抽出可能なユーザ固有の、或いは特定のユーザ間で共通の、或いは広範なユーザ間において共通の、リアルタイムに収集可能なインタレストグラフの有用性に注目が集まっている。しかしながら、従来のマイクロブログにおいては、ユーザがその時点で自らが関心を持った対象や状況に係るつぶやきが中心で、当該ユーザの近傍或いは視野内に存在する他のユーザの関心の対象に対しては、有効な気付きを十分与える事が出来ていないとは言えない。これらマイクロブログにおけるつぶやきの内容は極めて多岐に亘る為、特定のユーザ、特定の話題、或いは特定の場所等を指定して、テーマやトピックスを絞り込む方向の機能は提供されているものの、逆にその関心の対象をさらに広げて行く方向として、個々のユーザ特有の潜在的な関心の反映や、当該ユーザの身近に存在する他のユーザによる顕在的な関心の在り処の通知等、さらに広範なＳＮＳを誘発する可能性については、まだ十分生かし切れていないと言えな

40

50

い。

【課題を解決するための手段】

【0028】

上記課題を解決するために、本発明に係るネットワーク・コミュニケーションシステムは、一形態として、インターネットに接続可能なネットワーク端末に対し、有線或いは無線で接続可能な多機能入出力デバイスであって、少なくとも一以上のマイクロフォン、一以上のイヤフォン、一以上の画像撮像素子（カメラ）を一体として有する、ユーザの頭部に装着可能なヘッドセットシステムから得られる当該ユーザの主観的な視野、及び視点を反映した画像、及び音声信号を、前記ネットワーク端末経由でインターネット上の前記画像認識システムを備えた知識情報処理サーバシステム側にアップロード可能にし、当該画像に内包されている当該ユーザが着目した特定物体、一般物体、人、写真、或いはシーンに対し、音声認識システムとの協調動作により、当該ユーザ自身の音声による当該着目対象の指定、選択、及び抽出操作を、前記サーバシステム上で可能にした上で、当該ユーザによる上記一連の画像認識プロセス及び画像認識結果を、音声合成システムとの協調動作により、前記サーバシステム側がインターネットを介し、当該ユーザのネットワーク端末経由で、当該画像認識結果及びその認識プロセスを当該ユーザのヘッドセットシステムに組み込まれたイヤフォンに対し音声情報として、及び／又は、当該ユーザのネットワーク端末に音声及び画像情報として通知する事を可能にし、当該画像認識可能になった対象に対し、当該ユーザが自らの声で語りかけたメッセージやつぶやきを、前記音声認識システムとの協調動作により、前記サーバシステム側がその内容を分析・分類・蓄積し、それらメッセージやつぶやきをネットワーク経由で、同様の対象を目にした自らを含む広範なユーザ間で共有可能にする事で、多数のユーザの視覚的な好奇心に端を発する広範なネットワーク・コミュニケーションを誘発させると共に、それら広範なユーザ間のコミュニケーションを、前記サーバシステム側で統計的に観察・蓄積・解析する事で、当該ユーザ特有の、或いは特定のユーザ群に特有の、或いはユーザ全体に共通の動的な関心や好奇心の在り所とその推移を、上記広範な「ユーザ」群、抽出可能な「キーワード」群、及び様々な着目「対象」に係るノード群との間を繋ぐ動的なインタレストグラフとして獲得可能にする事を特徴とする。

【0029】

また、前記ネットワーク・コミュニケーションシステムにおいて、ユーザが関心を持った着目対象がどのような特徴を有しているか、及び／又は、どのような位置関係にあるか、及び／又は、どのような運動状態にあるかを、前記画像認識システムを備えた知識情報処理サーバシステム側にユーザが明示的に指し示す手段として、当該ユーザの音声による対象の選択指定（ポインティング）操作を可能にし、これら一連の選択指定の過程で当該ユーザが発声する当該対象に係る様々な特徴群を基に、前記音声認識システムとの協調動作により前記サーバシステム側が当該対象を正確に抽出・認識し、その画像認識結果に係る前記サーバシステム側から当該ユーザに向けての再確認内容として、当該ユーザが前記サーバシステム側に対し明示的に音声で指し示した特徴群以外に、当該ユーザの主観的視野を反映したカメラ映像を基に、前記サーバシステム側が当該対象に共起する新たな物体や事象群を抽出し、当該対象をさらに正確に言い表す事が可能な共起事象として加え、それらを一連の文章に構成し、前記音声合成システムとの協調動作により、当該ユーザに対し音声により再確認を求める事を可能にする事を特徴とする。

【発明の効果】

【0030】

本発明は、ユーザの頭部に装着可能なヘッドセットシステムに組み込まれたカメラから得られるユーザの主観的な視野を反映した画像信号を、当該ユーザのネットワーク端末経由でネットワークを介し前記画像認識システムを備えた知識情報処理サーバシステム側に適宜アップロードする事で、当該ユーザが関心を持って着目した特定物体、一般物体、人、写真、或いはシーン等の1以上の対象（以降「対象」と呼称）が、前記カメラ映像中のいずれに当るのかを、前記サーバシステムと当該ユーザ間の音声による双方向のコミュニ

ケーションにより抽出可能にする事で、従来の画像認識システムが不得意として来たユーザの「主観」を反映した対象の抽出及び認識処理を可能にし、画像認識率そのものを向上させる効果を与えると同時に、そこにユーザの音声による対象指定（ポインティング）操作と、それに対するサーバ側からの音声による再確認という双方向のプロセスを組み入れる事で、当該画像認識システムに対し継続的な機械学習が可能となる。

【0031】

また、ユーザによる前記音声指示を前記サーバシステム側で適宜解析する事で、当該対象に係る有用なキーワード群の抽出、及び当該ユーザによる当該対象に対する関心の抽出を可能にし、そこから広範なユーザ、様々なキーワード、及び様々な対象を構成ノード群とする、動的なインタレストグラフが獲得可能になる。

10

【0032】

その上で、当該インタレストグラフの対象となるノード群をネットワーク上でさらに広範なユーザ、様々な対象、及び様々なキーワードに対し拡大取得する事により、当該インタレストグラフの対象領域のさらなる拡大に加え、その収集頻度をさらに高める事が出来る。これにより、コンピュータ・システムによる継続的な学習プロセスに、人類の「知」をより効果的に組み入れて行く事が可能となる。

【0033】

また本発明は、前記画像認識システムを備えた知識情報処理システムにより認識可能になったユーザの着目対象に対し、当該ユーザが残した音声によるメッセージやつぶやきをネットワーク経由で前記サーバシステム内にアップロードし分類・蓄積しておく事で、異なる時空間において同様或いは類似の対象に近付いた、或いは着目した他のユーザ、或いはユーザ群に対し、前記サーバシステム側がネットワークを介し、当該ユーザのネットワーク端末経由で、前記メッセージやつぶやきを、当該ユーザとの音声コミュニケーションにより、インタラクティブに送り込む事を可能にする。これにより、多数のユーザに及ぶ様々な視覚的好奇心に端を発する広範なユーザコミュニケーションを、ネットワーク上で継続的に喚起する事が可能になる。

20

【0034】

また、ユーザが様々な対象に対して残した前記メッセージやつぶやきに係る内容の解析及び分類を前記サーバシステム側でリアルタイムに実行する事で、当該サーバシステム内に保持されている前記インタレストグラフの記述を基に、当該メッセージやつぶやきに含まれる主たる話題を抽出し、当該抽出された話題を中心ノードとするさらに関連性の高い他の話題群を抽出し、それらを抽出された話題に関心の高い他のユーザ及びユーザ群と、ネットワークを介して相互に共有可能にする事で、広範なユーザが目にする様々な対象や事象に端を発したネットワーク・コミュニケーションを継続的に誘発する事が可能となる。

30

【0035】

また本発明においては、当該ユーザ側から発した前記メッセージやつぶやきのみならず、当該サーバシステム自身側から発する様々な関心、好奇心、或いは疑問を当該ユーザ、或いはユーザ群に対し提起する事が出来る。例えば前記インタレストグラフ内に記載の対象ノード間の関連性から想定可能な範囲を超えて、特定のユーザが特定の対象に対して一定以上の関心を示す場合や、或いは逆に一定以下の関心しか示さない場合や、当該サーバシステム側だけでは認識が困難な対象や事象が存在した場合、或いはそれらに遭遇した場合等に、当該サーバシステム側から関連する質問やコメントを、当該ユーザ、或いは特定のユーザ群、或いは広範なユーザ群に対し積極的に提起する事を可能にする。これにより、前記サーバシステム側が様々な事象を通じて人類の「知」を継続的に吸収し、学習の上で自らの知識データベース内に体系立てて取り込んで行くプロセスが構成可能となる。

40

【0036】

近年では超高速光ファイバー網によるネットワークのさらなる高速化と相俟って、巨大なデータセンタの敷設が進み、超並列演算可能なスーパーコンピュータの開発も一段と加速している事から、コンピュータ・システム自身の自動学習プロセスにおいて、そこに人

50

類の「知」が効果的、有機的、かつ継続的に加わって行く事で、ネットワークを介してこれらの高性能コンピュータ・システム群による様々な事象の自動認識、及び機械学習が急速に発展して行く可能性がある。その為には、人類の「知」をいかにコンピュータ側が効果的に取得し、ネットワークを介して広範に共有可能な「知」の体系として再利用可能な状態に整理して行けるかが重要となる。言い換えると、いかにコンピュータの「好奇心」を刺激し、人とのコミュニケーションの中で継続的にコンピュータ・システムが進化して行ける効果的な方法を見つけられるかが重要となる。本発明においては、これらサーバ側に構築されたコンピュータ・システム自身による学習を、広範な対象に対する人々の視覚的関心と直接結び付ける具体的な方法を与える。

【図面の簡単な説明】

10

【0037】

【図1】本発明の一実施形態におけるネットワーク・コミュニケーションシステムの構成に關しての説明図である。

【図2】本発明の一実施形態におけるヘッドセットシステム及びネットワーク端末の構成に關しての説明図である。

【図3A】本発明の一実施形態における音声による対象画像抽出処理に關しての説明図である。

【図3B】本発明の一実施形態における音声による対象画像抽出処理に關しての説明図である。

【図4A】本発明の一実施形態における音声によるポインティングに關しての説明図である。

20

【図4B】本発明の一実施形態における学習によるグラフ構造の成長に關しての説明図である。

【図4C】本発明の一実施形態における複数対象候補の選択優先度処理に關しての説明図である。

【図5】本発明の一実施形態における知識情報処理サーバシステムの構成に關しての説明図である。

【図6A】本発明の一実施形態における画像認識システムの構成に關しての説明図である。

【図6B】本発明の一実施形態における一般物体認識部の構成及び処理フローに關しての説明図である。

30

【図6C】本発明の一実施形態における一般物体認識システムの構成及び処理フローに關しての説明図である。

【図6D】本発明の一実施形態におけるシーン認識システムの構成及び処理フローに關しての説明図である。

【図6E】本発明の一実施形態における特定物体認識システムの構成及び処理フローに關しての説明図である。

【図7】本発明の一実施形態における生体認証手順に關する説明図である。

【図8A】本発明の一実施形態におけるインタレストグラフ部の構成及び処理フローに關する説明図である。

40

【図8B】本発明の一実施形態におけるグラフデータベースの基本要素及び構成に關する説明図である。

【図9】本発明の一実施形態における状況認識部の構成及びグラフ構造例に關する説明図である。

【図10】本発明の一実施形態におけるメッセージ保管部の構成及び処理フローに關する説明図である。

【図11】本発明の一実施形態における再生処理部の構成及び処理フローに關する説明図である。

【図12】本発明の一実施形態におけるACL（アクセス制御リスト）に關する説明図である。

50

【図 1 3 A】本発明の一実施形態におけるユースケース・シナリオに関する説明図である。

【図 1 3 B】本発明の一実施形態における共通の対象への視覚的な好奇心に誘起されるネットワーク・コミュニケーションに関する説明図である。

【図 1 4】本発明の一実施形態におけるインタレストグラフに関するグラフ構造の説明図である。

【図 1 5】本発明の一実施形態における画像認識プロセスからのグラフ抽出手順に関する説明図である。

【図 1 6】本発明の一実施形態におけるインタレストグラフの獲得に関する説明図である。

10

【図 1 7】本発明の一実施形態における獲得されたインタレストグラフのスナップショットの一部に関する説明図である。

【図 1 8 A】本発明の一実施形態における時空間及び対象を指定可能なメッセージやつぶやきの記録と再生手順に関する説明図である。

【図 1 8 B】本発明の一実施形態における時間／時間帯の指定手順に関する説明図である。

【図 1 8 C】本発明の一実施形態における場所／地域の指定手順に関する説明図である。

【図 1 9】本発明の一実施形態におけるユーザが指定した時空間でのメッセージやつぶやきの再生手順に関しての説明図である。

【図 2 0】本発明の一実施形態におけるユーザの手指による対象指示手順に関する説明図である。

20

【図 2 1】本発明の一実施形態における視野の固定による対象指示の手順に関する説明図である。

【図 2 2】本発明の一実施形態における写真の検出手法に関する説明図である。

【図 2 3 A】本発明の一実施形態における対象との対話手順に関する説明図である。

【図 2 3 B】本発明の一実施形態における会話エンジンの構成と処理フローに関する説明図である。

【図 2 4】本発明の一実施形態における複数のヘッドセットからの共有ネットワーク端末の利用に関する説明図である。

【図 2 5】本発明の一実施形態における音声による W i k i 利用に関する処理手順の説明図である。

30

【図 2 6】本発明の一実施形態における位置情報を利用した誤差補正に関する説明図である。

【図 2 7】本発明の一実施形態における視点マーカーのキャリブレーションに関する説明図である。

【図 2 8】本発明の一実施形態におけるサーバとのネットワーク接続が一時的に切断されている状況におけるネットワーク端末単体での処理に関する説明図である。

【図 2 9】本発明の一実施形態における同一の時空間内に撮影された画像から抽出された特定物体、及び一般物体の事例である。

【図 3 0】本発明の一実施形態におけるアップロードされた画像に含まれる特定の時空間情報の抽出及び特定の時間軸の選択指定表示に関する説明図である。

40

【図 3 1】本発明の一実施形態における特定の時空間への視点移動時に特定の対象に係る会話を促す仕組みに関する説明図である。

【発明を実施するための形態】

【0038】

以下、本発明の一実施形態を図 1 から図 3 1 を用いながら説明する。

【0039】

図 1 を用いて、本発明の一実施形態におけるネットワーク・コミュニケーションシステム 100 の構成に関し説明する。前記ネットワーク・コミュニケーションシステムは、ヘッドセットシステム 200、ネットワーク端末 220、知識情報処理サーバシステム 30

50

0、生体認証システム310、音声認識システム320、音声合成システム330から構成される。前記ヘッドセットシステムは1以上存在し、1以上の前記ヘッドセットシステムが1個の前記ネットワーク端末にネットワーク251で接続される。前記ネットワーク端末は1以上存在し、インターネット250に接続される。前記知識情報処理サーバシステムは、生体認証システム310、音声認識システム320、及び音声合成システム330と、各々ネットワーク252、253、及び254で接続される。前記生体情報処理システムは、インターネット250と接続されていても良い。本実施例におけるネットワークは専用回線であっても良いし、インターネットを含む公衆回線であっても良いし、公衆回線上にVPN技術を用いて仮想的な専用回線を構築したものであっても良い。以下、特に断らない限りネットワークを前記の通り定義する。

10

【0040】

図2(A)に、本発明の一実施形態におけるヘッドセットシステム200の構成例を示す。前記ヘッドセットシステムは、図2(B)に示す様な、ユーザが装着する事で当該ネットワーク・コミュニケーションシステム100を利用可能なインターフェース装置である。図1において、ヘッドセットシステム200aから200cは、接続251aから251cでネットワーク端末220aに対し接続され、ヘッドセットシステム200dから200eは、接続251dから251eでネットワーク端末220bに対し接続され、ヘッドセットシステム200fは、接続251fでネットワーク端末220cに接続されている。つまり、ヘッドセット200aから200fは、ネットワーク端末220aから220cを介して、インターネット経由で知識情報処理サーバシステム300に繋がっている様子を表わしている。以下、ヘッドセットシステム200と記載した場合にはヘッドセットシステム200aから200fのいずれか一台を指す。ヘッドセットシステム200aから200fは、全て同一機種である必要はない。同等の機能、或いは実施可能な最低限の機能を備えた同様の装置であれば良い。

20

【0041】

ヘッドセットシステム200は以下の要素群で構成されるが、これらに限らず、そのいくつかを選択して搭載しても良い。マイクロフォン201は1以上存在し、当該ヘッドセットシステムを装着したユーザの音声や、当該ユーザの周辺の音を収集する。イヤフォン202は1以上存在し、モノラル或いはステレオで、他のユーザのメッセージやつぶやき、サーバシステムからの音声による応答等を含む様々な音声情報を、当該ユーザに通知する。カメラ(画像撮像素子)203は1以上存在し、当該ユーザの主観的な視野を反映した映像以外に、ユーザの背後や側面、或いは上部等の死角となっているエリアからの映像も含んでも良い。また、静止画であるか動画であるかを問わない。生体認証センサ204は1個以上存在し、一実施例としてユーザの有用な生体識別情報の一つである静脈情報(鼓膜や外耳部から)を取得し、前記生体認証システム310と連携して、当該ユーザ、当該ヘッドセットシステム、及び前記知識情報処理サーバシステム300間を、認証し紐付ける。生体情報センサ205は1以上存在し、ユーザの体温、心拍、血圧、脳波、呼吸、眼球移動、発声、体の動き等の検出可能な各種生体情報(バイタルサイン)を取得する。深度センサ206は、前記ヘッドセットシステムを装着したユーザに近づく、人間を含む或る程度以上の大きさの生体の移動を検知する。画像出力装置207は、前記知識情報処理サーバシステム300からの各種通知情報を表示する。位置情報センサ208は、前記ヘッドセットシステムを装着したユーザの位置(緯経度、高度、向き)を検知する。一例として、当該位置情報センサに6軸モーションセンサ等を装備する事で、移動方向、向き、回転等を前記に追加して検出する様に構成しても良い。環境センサ209は、前記ヘッドセットシステム周辺の明るさ、色温度、騒音、音圧レベル、温湿度等を検知する。視線検出センサ210は、一実施例として前記ヘッドセットシステムの一部からユーザの瞳、又は網膜に向けて安全な光線を照射し、その反射光を計測する事で、ユーザの視線方向を直接検知する。無線通信装置211は、ネットワーク端末220との通信、及び前記知識情報処理サーバシステム300との通信を行う。電源部212は、前記ヘッドセットシステム全体に電力を供給する為の電池等を指すが、有線で前記ネットワーク端末に接続可能

30

40

50

な場合は、外部からの電力供給によっても良い。

【0042】

図2(C)に、本発明の一実施形態におけるネットワーク端末220の構成例を示す。図1において、ネットワーク端末220aから220fは広くユーザが利用するクライアント端末装置であり、PC、携帯情報端末(PDA)、タブレット、インターネット接続可能な携帯電話、スマートフォン等が含まれ、これらがインターネットに接続されている様子を表している。以下、ネットワーク端末220と記載した場合には、インターネットに接続されたネットワーク端末220aから220fのいずれか一台を指す。ネットワーク端末220aから220fは同一機種である必要はない。同等の機能、或いは実施可能な最低限の機能を備えた端末装置であれば良い。

10

【0043】

ネットワーク端末220は以下の要素群で構成されるが、これらに限らずそのいくつかを選択して搭載しても良い。操作部221は、表示部222と共にネットワーク端末220のユーザインターフェース部である。ネットワーク通信部223は、インターネットとの通信、及び1以上のヘッドセットシステムとの通信を担当する。前記ネットワーク通信部は、IMT-2000、IEEE802.11、Bluetooth、IEEE802.3、或いは独自の有線/無線規格、及びルータを経由したその混合形態であっても良い。認識エンジン224は、知識情報処理サーバシステム300の主要な構成要素である画像認識システム301が有する画像認識処理機能から、限定された対象に関する画像認識処理に特化した前記ネットワーク端末に最適化した画像認識プログラムを前記知識情報処理サーバシステム側からダウンロードし実行する。これにより、前記ネットワーク端末側にも一定の範囲内で画像検出・認識機能の一部を持たせる事で、前記サーバ側の画像認識システム側に対する処理負担の軽減、及びネットワーク回線の負荷の軽減を図る事が出来ると共に、その後のサーバ側での認識プロセスに際し、後述の図3Aにおけるステップ30-20から30-37に対応する予備的な前処理を実行する事が可能となる。同期管理部225は、ネットワークの不具合等により回線の一時的な切断が発生し、再び回線が復帰した際にサーバ側との同期処理を行う。CPU226は中央処理装置であり、記憶部227は主メモリ装置であり、又フラッシュメモリ等を含む一次、及び二次記憶装置である。電源部228は、当該ネットワーク端末全体に電力を供給する為の電池等の電源である。これらネットワーク端末は、ネットワーク網に対し緩衝的な役割を果たす。例えば、ユーザにとって重要ではない情報をネットワーク側にアップロードしても、それは知識処理サーバシステム300にとっては当該ユーザとの紐付けという意味ではノイズであり、ネットワーク回線に対しても不要なオーバーヘッドとなる。従って、可能な範囲で或る程度のスクリーニング処理をネットワーク端末側で行う事で、ユーザに対する有効なネットワークバンド幅の確保や、ローカルティが高い処理に関し応答速度の向上を図る事が可能になる。

20

30

【0044】

図3Aを用いて、本発明の一実施例としてユーザが関心を持った対象に着目する際のユーザの音声による対象画像抽出処理30-01のフローを説明する。前記で定義した様に本実施例では特定物体、一般物体、人、写真、或いはシーンを「対象」と総称する事にする。前記対象画像抽出処理は、ステップ30-02のユーザによる音声入力トリガで始まる。前記音声入力トリガには、特定の言葉や一連の自然言語を用いても良いし、音圧レベルの変化を検出する事によりユーザの発声を検出しても良いし、またネットワーク端末220上のGUI操作によっても良い。前記ユーザの音声入力トリガによりユーザのヘッドセットシステムに具備されているカメラの撮影が開始され、そこから取得可能になる動画像、連続した静止画、或いは静止画を、前記知識情報処理サーバシステム300に対しアップロードを開始し(30-03)、その後ユーザからの音声コマンド入力待ち状態(30-04)に入る。

40

【0045】

一連の対象画像抽出、及び画像認識処理フローは、音声認識処理、画像特徴抽出処理、

50

着目対象抽出処理、そして画像認識処理の順番で実行される。具体的には、音声入力コマンド待ち(30-04)からユーザの発話を認識し、当該音声認識処理によりユーザの発声した一連の言葉から単語列を抽出し、当該単語列に基づいて画像の特徴抽出処理を行い、抽出可能になった画像特徴群を基に画像認識処理を実行し、対象が複数に亘る場合や、対象自体からの特徴抽出が困難である場合等に、ユーザに対しさらなる画像特徴群の入力を求める事で、ユーザが着目した対象をサーバ側がより確実に認識するプロセスを構成する。上記ユーザの発話による「再確認」のプロセスを加える事で、画像認識システムの全ての処理プロセスをコンピュータ・システム側のみで対処しなくてはならないという従来の発想を転換して、従来画像認識システムが不得意として来た対象画像の正確な抽出、或いは従来の音声認識システムが不得意として来た同音異義語への対応問題等への効果的な対処が可能になる。実際の導入に当たっては、これらの一連の画像認識プロセスを、いかにユーザにとり煩わしい作業と思わせずに楽しいコミュニケーションと思わせられるかが重要となる。前記一連の画像特徴抽出処理では、図3Aに示す事例よりもさらに多様な画像特徴群に対応する画像特徴抽出処理部群を多数並列に配置して一気に並列処理する事が可能で、それにより画像認識精度の一層の向上と併せて処理の大幅な高速化を図る事が可能となる。

10

【0046】

ユーザの音声による対象のポインティング方法としては、当該ステップ30-06から30-15で例示した様な、各画像特徴群に対しユーザがそれらを各々単独に選択しながらポインティングして行く事例より、複数の画像特徴群を含んだ一連の言葉として一括してポインティングする事例の方が多いものと想定される。この場合は、複数の画像特徴群による対象の抽出処理が同時並列に行われ、そこから当該対象を表現する複数の画像特徴要素群が得られる可能性が高い。そこからより多くの特徴が抽出可能になれば、当該着目対象のポインティングの確度は一段と高まる。それら抽出可能になった画像特徴群を手掛かりに、前記画像認識システムによる画像認識処理30-16が開始される。画像認識は、一般物体認識システム106、特定物体認識システム110、及びシーン認識システム108により実行される。図3Aでは、これらを連続したフローで表現しているが、当該画像認識処理は各々並列、或いは各一般物体認識、特定物体認識、及びシーン認識処理の中でさらに並列化する事が可能で、当該画像認識処理の認識速度に係る処理時間を大幅に短縮する事が出来る。上記の結果として、当該画像認識された対象に係る様々な認識結果を、音声で当該対象に係る画像認識結果として、ユーザに通知する事が可能になる。

20

30

【0047】

この場合であっても、上記画像認識結果に加えて当該ユーザが指し示した特徴要素群のみを引用してユーザに再確認を求めたとしても、果たしてそれで本当にユーザが着目した対象をシステム側が正しく抽出したのか疑問が残る場合もある。例えば、ユーザの視野を反映したカメラ画像の中には、類似の物体が複数存在している可能性もある。本特許では、当該不確実性に対応する為、前記画像認識システムを備えた知識情報処理サーバシステム側が、当該対象の近傍状況を、当該カメラ映像を基に精査する事で当該対象と「共起」している新たな物体や事象を抽出(30-38)し、当該ユーザが明示的に指し示していないそれら新たな特徴要素群を上記再確認の要素に加え(30-39)、当該ユーザに対し音声による再確認(30-40)を求める事で、ユーザの着目対象と上記サーバシステム側が抽出した対象が同一であることを再確認する事を可能に構成することが出来る。

40

【0048】

上記一連の処理は、基本的に同一の対象に関する処理であり、ユーザはその行動において常に他の対象に興味を移行し得るので、図3Aにおける前記ステップ群を包含するさらに大きな外側の処理ループも存在する。なお、前記画像認識処理ループは、前記ヘッドセットシステムをユーザが装着した時点で開始しても良いし、ステップ30-02同様の音声トリガによっても開始しても良いし、前記ネットワーク端末を操作する事によって開始しても良いが、必ずしもそれらには限らない。前記処理ループの停止は、前記処理ループの開始における手段と同様に、前記ヘッドセットをユーザが外した時としても良いし、音

50

声トリガによっても良いし、前記ネットワーク端末を操作する事によって停止しても良いが、必ずしもそれらには限らない。さらに、ユーザの着目の結果認識された対象は、当該時空間情報を付して後述のグラフデータベース365に記録する事で、後日の問い合わせに回答出来る様に構成しても良い。前記図3Aに記載の対象画像抽出処理は本発明における重要なプロセスであり、以下その各ステップを説明する。

【0049】

最初に、ユーザによる音声入力トリガ(30-02)が発生し、カメラ画像のアップロード(30-03)開始後、音声認識処理30-05によりユーザの対象検出コマンドから単語列が抽出され、前記単語列が条件群30-07から30-15のいずれかの特徴に適合した場合には、係る画像特徴抽出処理に引き渡される。前記単語列が「対象の名称」である場合(30-06)、例えば、ユーザが当該対象に係る固有名詞を発話した場合、当該アノテーションはユーザの一定の認識判断を反映したものととして、係る特定物体認識の実行(110)処理を行う。その照合結果と、当該アノテーションに齟齬がある場合、或いは疑問がある場合は、当該ユーザによる誤認識の可能性もあるとして、当該ユーザに喚起を促す。或いはユーザが、当該対象に係る一般名詞を発話した場合、当該一般名詞に係る一般物体認識の実行(106)処理を行い、その画像特徴から対象を抽出する。或いはユーザが当該対象に係るシーンを発話した場合、当該シーンに係るシーン認識の実行(108)処理を行い、その画像特徴から対象領域を抽出する。またそれらの特徴を一つだけ指し示すのではなくて、複数の特徴を含む情景として指定しても良い。例えば、道路(一般物体)の左側(位置)を走る(状態)黄色い(色)タクシー(一般物体)、ナンバーは「1234(特定物体)」という様な指定の方法である。これらの対象指定を一連の言葉としても良いし、各々個別に指定を行っても良い。対象が複数個発見される場合には、前記画像認識システムによる再確認プロセスを経て、さらに新たな画像特徴を追加して対象を絞り込んで行く事が出来る。当該画像抽出結果は、一例としてユーザに対し音声による質問、例えば「それは~ですか?」を発行して再確認処理される(30-40)。当該再確認内容に対し、着目対象の抽出がユーザの意図通りである場合は、ユーザはその旨を示す言葉或いは単語を発話して、ステップ30-50「カメラ画像アップロード終了」を実行し、当該対象画像抽出処理を終了する(30-51)。一方、ユーザの意図とは違う場合には、再びステップ30-04「音声コマンド入力待ち」に戻り、さらなる画像特徴群を入力する。また、何度入力しても対象の特定に至らない場合や、そもそも対象自体が視野外に移動してしまった場合等には、処理を中断(QUIT)して当該対象画像抽出処理を終了する。

【0050】

例えば音声認識処理30-05の結果が図3Aで示す条件30-07に適合した場合、即ちユーザが対象の「色」に関する特徴を発話した場合には、色抽出処理30-20が行われる。当該色抽出処理には、RGB3原色において色毎に範囲を設定して抽出する手法を用いても良いし、それらをYUV色空間上で抽出しても良い。またこれら特定の色空間表現には限定されない。当該色抽出処理後に対象を分離抽出し(30-29)、セグメンテーション(切り出し領域)情報を得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理(30-16)を行う。その後は当該画像認識処理結果を利用して共起物体や共起事象を抽出(30-38)し、抽出可能になった全特徴群に関する記述を生成(30-39)し、当該記述をもってユーザに再確認を求める(30-40)。その結果がYESであれば、カメラ画像のアップロードを終了(30-50)し、音声による対象画像の抽出処理を終了(30-51)する。

【0051】

例えば音声認識処理30-05の結果が図3Aで示す条件30-08に適合した場合、即ちユーザが対象の「形状」に関する特徴を発話した場合には、形状特徴抽出30-21が行われる。当該形状特徴抽出処理では、対象に係るエッジ追跡を行いながら輪郭や主要な形状特徴を抽出後、形状のテンプレート・適合処理を行うが、それ以外の手法を用いても良い。当該形状抽出処理後に対象を分離し(30-30)、セグメンテーション情報を

得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理（３０－１６）を行う。その後は当該画像認識処理結果を利用して共起物体や共起事象を抽出（３０－３８）し、抽出可能になった全特徴群に関する記述を生成（３０－３９）し、当該記述をもってユーザに再確認を求める（３０－４０）。その結果がＹＥＳであれば、カメラ画像のアップロードを終了（３０－５０）し、音声による対象画像の抽出処理を終了（３０－５１）する。

【００５２】

例えば音声認識処理３０－０５の結果が図３Ａで示す条件３０－０９に適合した場合、即ちユーザが対象の「大きさ」に関する特徴を発話した場合には、物体サイズ検出処理３０－２２が行われる。その一例として、当該物体サイズ検出処理ではサイズ以外の他の特徴抽出処理等により切り分けされた当該対象物体に対し、周囲にある他の物体との相対的なサイズ比較がユーザとのインタラクティブな音声コミュニケーションにより実行される。例えば「左隣の～よりも大きな～」という様な指示である。その理由としては、対象が単独で存在する場合、その大きさの比較になる様な具体的な指標がないと、単に画角から見た大きさのみでそのサイズを一意に判断出来ない事によるが、それ以外の手法を用いても良い。当該サイズ検出後に対象を分離し（３０－３１）、セグメンテーション情報を得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理（３０－１６）を行う。その後は当該画像認識処理結果を利用して共起物体や共起事象を抽出（３０－３８）し、抽出可能になった全特徴群に関する記述を生成（３０－３９）し、当該記述をもってユーザに再確認を求める（３０－４０）。その結果がＹＥＳであれば、カメラ画像のアップロードを終了（３０－５０）し、音声による対象画像の抽出処理を終了（３０－５１）する。

【００５３】

例えば音声認識処理３０－０５の結果が図３Ａで示す条件３０－１０に適合した場合、即ちユーザが対象の「明るさ」に関する特徴を発話した場合には、輝度検出処理３０－２３が行われる。当該輝度検出処理では、ＲＧＢ３原色から、或いはＹＵＶ色空間から特定領域の輝度を求めるが、それら以外の手法を用いても良い。当該対象の輝度検出処理では、対象の周囲と比較した相対輝度の抽出が、ユーザとのインタラクティブな音声コミュニケーションにより実行される。例えば「周りより明るく輝いている～」という様な指示である。その理由としては、対象が単独で存在する場合、その明るさの比較になる様な具体的な指標がないと、単に画素が有する輝度値のみでユーザが感じた輝度を一意に判断出来ない理由によるが、それ以外の手法を用いても良い。当該輝度検出後に対象を分離し（３０－３２）、セグメンテーション情報を得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理（３０－１６）を行う。その後は当該画像認識処理結果を利用して共起物体や共起事象を抽出（３０－３８）し、抽出可能になった全特徴群に関する記述を生成（３０－３９）し、当該記述をもってユーザに再確認を求める（３０－４０）。その結果がＹＥＳであれば、カメラ画像のアップロードを終了（３０－５０）し、音声による対象画像の抽出処理を終了（３０－５１）する。

【００５４】

例えば音声認識処理３０－０５の結果が図３Ａで示す条件３０－１１に適合した場合、即ちユーザが「対象との距離」に関する特徴を発話した場合には、奥行き検出処理３０－２４が行われる。当該奥行き検出処理では、ユーザのヘッドセットシステム２００に具備された深度センサ２０６を用いて奥行きを直接測定しても良いし、２台以上のカメラ映像から得られる視差情報から計算により算出しても良い。また、これら以外の手法を用いても良い。当該距離検出後に対象を分離し（３０－３３）、セグメンテーション情報を得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理（３０－１６）を行う。その後は当該画像認識処理結果を利用して共起物体や共起事象を抽出（３０－３８）し、抽出可能になった全特徴群に関する記述を生成（３０－３９）し、当該記述をもってユーザに再確認を求める（３０－４０）。その結果がＹＥＳであれば、カメラ画像のアップロードを終了（３０－５０）し、音声による対象画像の抽出処理を終了（３０－５１）

する。

【 0 0 5 5 】

例えば音声認識処理 30 - 0 5 の結果が図 3 A で示す条件 30 - 1 2 に適合した場合、即ちユーザが「対象の存在する位置 / 領域」に関して発話した場合には、対象の領域検出 30 - 2 5 が行われる。当該領域検出処理では、一例としてユーザの主たる視野を反映したカメラ画像全体を予め等間隔のメッシュ状に領域分割し、ユーザからのインタラクティブな指示として「右上の～」という様な領域指定から対象を絞り込んでも良いし、「机の上の～」という様な、対象が存在する場所の指定で行っても良い。また、他の位置 / 領域に係る指定であっても良い。当該対象の存在する位置 / 領域検出後に対象を分離し (30 - 3 4)、セグメンテーション情報を得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理 (30 - 1 6) を行う。その後は当該画像認識処理結果を利用して他の共起物体や共起事象を抽出 (30 - 3 8) し、抽出可能になった当該共起特徴群を含む記述を生成 (30 - 3 9) し、当該記述をもってユーザに再確認を求める (30 - 4 0)。その結果が Y E S であれば、カメラ画像のアップロードを終了 (30 - 5 0) し、音声による対象画像の抽出処理を終了 (30 - 5 1) する。

10

【 0 0 5 6 】

例えば音声認識処理 30 - 0 5 の結果が図 3 A で示す条件 30 - 1 3 に適合した場合、即ちユーザが「対象と他物体との位置関係」に関して発話した場合には、当該対象に係る共起関係検出 30 - 2 6 が行われる。当該共起関係検出処理では、図 3 A に記載の各処理 (1 0 6、1 0 8、1 1 0、30 - 2 0 から 30 - 2 8) により抽出された対応特徴に係るセグメンテーション情報を用いて、それらのセグメンテーション情報に対応する各特徴との共起関係を精査する事で、対象の抽出を行う。一例として「～と一緒に写っている～」という様な指示であるが、これ以外の手法を用いても良い。これにより、当該対象と他物体との位置関係を基に対象を分離し (30 - 3 5)、当該対象に係るセグメンテーション情報を得る。次に当該セグメンテーション情報を手掛かりに対象の画像認識処理 (30 - 1 6) を行う。その後は当該認識結果を利用して他の共起物体や共起事象を抽出 (30 - 3 8) し、抽出可能になった当該共起特徴群を含む記述を生成 (30 - 3 9) し、当該記述をもってユーザに再確認を求める (30 - 4 0)。その結果が Y E S であれば、カメラ画像のアップロードを終了 (30 - 5 0) し、音声による対象画像の抽出処理を終了 (30 - 5 1) する。

20

30

【 0 0 5 7 】

例えば音声認識処理 30 - 0 5 の結果が図 3 A で示す条件 30 - 1 4 に適合した場合、即ちユーザが「対象の動き」に関して発話した場合には、動き検出処理 30 - 2 7 が行われる。当該動き検出処理では、時間軸上に連続的に展開された複数枚の画像を参照し、各画像を複数のメッシュ領域に分割し、当該領域を相互に比較する事によって、カメラ自体の移動による全体画像の平行移動以外に、相対的に個別移動している領域を見つけ出し、その領域の差分抽出 (30 - 3 6) 処理を行い、周囲に比べて相対的に移動している領域に係るセグメンテーション情報を得る。また、これら以外の手法を用いても良い。次に当該セグメンテーション情報を手掛かりに、対象の画像認識処理 (30 - 1 6) を行う。その後は当該画像認識処理結果を利用して他の共起物体や共起事象を抽出 (30 - 3 8) し、抽出可能になった当該共起特徴群を含む記述を生成 (30 - 3 9) し、当該記述をもってユーザに再確認を求める (30 - 4 0)。その結果が Y E S であれば、カメラ画像のアップロードを終了 (30 - 5 0) し、音声による対象画像の抽出処理を終了 (30 - 5 1) する。

40

【 0 0 5 8 】

例えば音声認識処理 30 - 0 5 の結果が図 3 A で示す条件 30 - 1 5 に適合した場合、即ちユーザが「対象の様子」に関して発話した場合には、状態検出処理 30 - 2 8 が行われる。当該状態検出処理では、物体の状態、例えば、運動状態 (静止、移動、振動、浮遊、上昇、下降、飛翔、回転、泳動、接近、離遠等)、動作状態 (走っている、跳んでいる、しゃがんでいる、座っている、寝ている、横たわっている、眠っている、食べている、

50

飲んでいる、観察可能な喜怒哀楽等を含む)を、当該状態に係る特徴を記述した知識データベース(未図示)を参照しながら、連続する複数の画像群から推定・抽出(30-37)し、セグメンテーション情報を得る。次に当該セグメンテーション情報を手掛かりに、対象の画像認識処理(30-16)を行う。その後は当該画像認識処理結果を利用して、他の共起物体や共起事象を抽出(30-38)し、抽出可能になった当該共起特徴群を含む記述を生成(30-39)し、当該記述をもってユーザに再確認を求める(30-40)。その結果がYESであれば、カメラ画像のアップロードを終了(30-50)し、音声による対象画像の抽出処理を終了(30-51)する。

【0059】

ユーザは前記ステップに係る音声による図3Aで示す再確認(30-40)のステップにおいて、前記対象画像抽出処理をユーザの発話により中止する事が出来る。音声認識処理30-05において、前記中止コマンドが認識された場合には、ステップ30-50に移行しカメラ画像アップロードを終了し、音声による対象画像抽出処理を終了する(30-51)。前記記載の各々の対象の検出、抽出、或いは認識処理において、処理時間が一定以上長引く場合には、ユーザに対して興味を引き続ける目的で処理の経過を示す状況や、関連する情報を音声で伝える事が出来る。例えば、「今着目している~の認識処理を、引き続きサーバに問い合わせ中です。現在~人の方が同様の対象に注目しています。もう少しお待ち下さい」「~までの処理が終わりました。途中経過は~です」の様な経過メッセージを当該ユーザに対し、音声で返す事が出来る。

【0060】

ここで、図3Bを用いて、図3Aをデータの流れから説明する。入力画像35-01と発話35-02である。認識・抽出処理制御35-03では、発話35-02の入力による図3Aにおけるステップ30-06から30-15を1以上実行し、画像35-01に対して図3Aにおけるステップ35-16を実行する際には、一般物体認識処理システム110による一般部隊認識処理、特定物体認識システム110による特定物体認識処理、及びシーン認識システム108によるシーン認識処理のいずれか1以上を実行する。画像認識システム106、108、110の各々の機能ブロックは、実行ユニット毎にさらなる並列化が可能であり、画像認識処理ディスパッチ35-04により1以上の処理に振り分けられて並列に実行される。また、発話35-02の入力に対し、図3Aにおけるステップ30-07から30-15を実行する場合には、特徴抽出処理30-20から30-28、及び分離抽出処理30-29から30-37を実行する。上記特徴抽出処理及び分離抽出処理は各々1以上存在し、特徴抽出ディスパッチ35-05により1以上の処理に分けられて並列に実行される。前記認識・抽出処理制御35-03では、ユーザの発話に処理順序に影響を与える単語が含まれている場合(例えば、「~の上」という場合には「~」を画像認識する必要がある、その次に「上」を処理する)には、順序制御を行う。

【0061】

入力画像35-01に関して、認識・抽出処理制御35-03は、後述のグラフデータベース365にアクセスして、代表ノード35-06を抽出(当該データベースに当該ノードが存在しなければ新しい代表ノードを生成)する。前記一連の処理より、画像35-01が発話35-02に従って処理され、前記同時実行される各認識・抽出処理群に係る結果のグラフ構造35-07がグラフデータベース365に蓄積される。この様にして、入力画像35-01に対する認識・抽出処理制御35-03による一連のデータの流れは、発話35-02が当該入力画像に関して有効にある限り続く。

【0062】

次に図4Aを用いて、本発明の一実施例におけるユーザの音声による対象のポインティング操作を説明する。これは、図3Aに記載の手順に対する応用例である。図4A(A)の場所は、ニューヨーク州マンハッタン島タイムズ・スクエア境界である。この場所にいるユーザ、或いはこの写真を見たユーザが仮に発話41「A yellow taxi on the road on the left side」をつぶやいたとする。ここから音声認識システム320は、当該発話41から複数の文字列或いは単語列を抽出す

る。当該発話から抽出可能な単語としては「一台」の「黄色」の「タクシー」が「左側」の「道路上」に見える、の5個である。ここから、前記図3Aで示した対象画像抽出フローにおける「対象の名称」「対象の色情報」「対象の位置」「対象の存在する領域」及び着目している対象が複数ではなく、単一の対象である事が判る。これらの手掛かりから、当該画像特徴群を有する対象の検出・抽出処理が開始され、それが点線円(50)のタクシーである可能性を前記画像認識システム側がユーザに音声により返答する事が可能となった場合、前述した様にその再確認内容として、上記ユーザが明示的に示した特徴要素群のみで再確認するだけでは、今一步確実性に欠ける場合がある。これらの不確実性に対処する為に、ユーザがまだ指し示していない当該対象に係る他の共起特徴要素群を検出し、それらを再確認内容に加える必要がある。例えば「それは手前の横断歩道に差し掛かっているタクシーで、前に人が見えますね?」という様に、前記画像認識システムを備えた知識情報処理サーバシステム側が検出した当該対象に係る新たな共起事象を加えユーザに再確認を求める事が出来れば、よりユーザの意に沿った対象の検出・抽出・絞り込み処理が可能となる。本事例では、点線円(50)を含む領域の拡大画像図4A(B)から、「横断歩道」(55)「人」(56)が検出可能となっている様子を示している。

10

【0063】

同様に、大きな看板があるビルを見上げているユーザが、発話45「I'm standing on the Times Square in NY now」とつぶやけば、カメラ画像を用いた適合処理により、そこが「ニューヨーク」州「タイムズ・スクウェア」で、ユーザが有名なランドマークとなっている建物を着目していると推測可能になる。

20

【0064】

同様に、発話42「A red bus on the road in front」という表現から、「1台(対象の数)」の「赤(対象の色特徴)」い「バス(対象の名称)」が「正面(対象の存在する位置)」の「道路(一般物体)」上(対象の位置関係)」が抽出可能になり、ユーザが点線円51内のバスを着目していると推定可能になる。

【0065】

同様に、発話44「The sky is fair in NY today」という表現から、「今日」の「NY」の天気は「晴れ」が抽出可能になり、ユーザが点線円(52)の領域「空」を見上げていると推定可能になる。

30

【0066】

少し複雑なつぶやき43「A big ad-board of "the Phantom of the Opera", top on the building on the right side」からは、「右端」に見える「ビル」の「屋上」にある、点線円(53)で示した「オペラ座の怪人」の「広告ボード」をユーザが着目していると推定可能になる。

【0067】

これら検出可能な単語列は、各々「固有の名称」「一般名詞」「シーン」「色」「位置」「領域」「場所」等を示しており、それらに対応した画像検出・画像抽出処理が実行される。その結果が当該時空間情報、及び画像情報と共に、前記知識情報処理サーバシステム300上に引き渡される。なお、図4Aに記載のイメージは本発明の一実施例を説明したもので、それに限定されない。

40

【0068】

ここで、図4Bを用いて、本発明の一実施例における図3Aに記載の手順を実行する過程の学習機能に関して、図4Aのシーンを例に説明する。図4B(A)は図4Aに記載のユーザの主たる視野を反映した画像に関して獲得されたグラフ構造(後述)の一部のスナップショットである。まず画像認識プロセスとグラフ構造との関係を説明する。

【0069】

ノード(60)は図4Aを代表するノードであり、図4Aの画像データを記録しているノード(61)とリンクしている。以下、ノードとノードのリンクを用いて情報を表現す

50

る。ノード(60)はまた、場所を表わすノード(62)と、時間を表わすノード(63)に対してもリンクしている事で、撮影場所と時間の情報を保持している。さらにノード(60)は、ノード(64)とノード(65)とリンクしている。ノード(64)は、図4A中の点線円(50)の対象を代表するノードであり、前記発話41により、特徴量T1(65)、特徴量T2(66)、色属性(67)、切り抜き画像(68)、及び画像内の位置座標(69)の各情報を保持している。前記特徴量は、図3Aの手順の過程における後述の一般物体認識システム106の処理結果として得られる。ノード(65)は、図4Aの点線円(51)の対象を代表するノードであり、前記ノード(64)と同様の情報を保持している。なお、ノード(60)即ち図4Aは、ユーザ1の主観視画像としてノード(77)とリンクしている。

10

【0070】

次に、ユーザ2を表すノード(80)の主観視を代表するノード(81)の保持する情報を、図4B(B)に示す。図では簡略化のため、図4B(A)に記載のノードのうちいくつかは省略している。ノード(82)は、ユーザ2の主観視における図4Aの点線円(51)に相当する対象の代表ノードである。同様に、特徴量C1(84)とC2(85)を情報として保持している。

【0071】

前記ノード(65)にリンクする特徴量であるB1(70)及びB2(71)と、前記ノード(82)にリンクする特徴量であるC1(84)及びC2(85)は、一般物体認識システム106において比較され、同一対象であると判断された場合(即ち同じカテゴリに属した場合)、或いは統計的に新たな重心となり得る場合には、代表特徴量D(91)が算出され学習に付される。本実施例では、当該学習結果をVisual Word辞書110-10に記録する。さらに、対象を代表するノード(90)、及びそのサブノード群(91から93と75から76)をリンクした部分グラフが生成され、ノード(60)は、ノード(65)とのリンクをノード(90)とのリンクに置き換える。同様にノード81は、ノード82とのリンクをノード90とのリンクに置き換える。

20

【0072】

次に、他のユーザが異なる時空間において、図4Aで点線円(50)に相当する対象に着目した場合には、前記同様のグラフ構造を構築するが、当該対象に対して一般物体認識システム106は、前記学習により当該対象の特徴量がノード(90)に記録された特徴量と同じクラスにも属すると判断出来るので、ノード(90)とリンクする様にグラフ構造を構築する事が出来る。

30

【0073】

図3Aに記載の、ステップ30-20から30-28に対応する特徴抽出処理において抽出された特徴群は、ユーザの発話と、セグメンテーション情報と、当該特徴とをノードに持つグラフ構造として表現出来る。例えば、図4Aの点線円(50)のセグメンテーション領域の場合で、特徴抽出処理がステップ30-20の場合には、色に関する特徴ノードを保持するグラフ構造となる。当該グラフ構造は、既に対象に関する代表ノードが存在する時には、その部分グラフと比較される。図4Bの例では、ノード(67)の色特徴“yellow”と近いと判断出来るので、当該グラフ構造は代表ノード(64)の部分グラフになる。この様なグラフ構造の統合を記録しておいても良い。それにより、当該例では、ユーザの発話と色特徴との関係を記録する事が出来るので、“yellow”に対応する色特徴の確からしさを高める事になる。

40

【0074】

上記記載の手順により、後述の画像認識に係るデータベース群(107、109、111、110-10)と、後述のグラフデータベース365は成長(新しいデータを獲得)する。上記記載では一般物体の場合を説明したが、特定物体、人、写真、或いはシーンであっても、同様に当該データベース群に対象に関する情報が蓄積される。

【0075】

次に図4Cを用いて、本発明の一実施形態におけるグラフデータベース365から複数

50

の対象候補ノードが抽出された場合に、ユーザがどれに着目しているかを算出する手段に関して説明する。当該手順は、例えば、図3Aにおける手順のステップ30-38及びステップ30-39において抽出可能になった複数の対象候補から、ユーザの着目対象を選び出す際に利用出来る。

【0076】

ステップ(S10)は、前記ステップ30-38の結果の共起物体・事象に対応する代表ノードをグラフデータベース365から抽出する(S11)。当該ステップは、図3Aに記載のステップ30-16、及びステップ30-20から30-28において、前記グラフデータベースをアクセスする事で、例えば色特徴抽出30-20では図4Aに係する色ノードから、対象ノード(64)と(65)を、図4Aノード60と、2つの色ノード(67)と(72)のリンクから抽出する事が出来る。

10

【0077】

前記ステップ(S11)では、1以上の代表ノードが抽出され得る。その全ての代表ノードに対して、次のステップを繰り返す(S12)。ステップ(S13)では、一つの代表ノードを変数*i*に格納する。そして、当該変数*i*の代表ノードを参照しているノード数を、変数*n_ref[i]*に格納する(S14)。例えば、図4B(C)ではノード(90)を参照しているノードからのリンクは点線円(94)のリンクであり、「3」となる。次に*n_all[i]*にノード*i*の部分グラフの全ノード数を代入(S15)する。図4B(C)のノード(90)では「5」を代入する。次に、*n_ref[i]*が規定値以上か?が判断される。YESの場合には*n_fea[i]*に1を代入(S17)し、NOの場合には0を代入(S18)する。ステップ(S19)では*n_fea[i]*に、前記ノード*i*の部分グラフ中で図3Aに記載の手順で、ユーザの発話した特徴に対応するノードの数を*n_all[i]*で除した数値を加算する。例えば、図4B(C)の例で、ノード(90)に関して、ユーザが“red”のみを発話した場合には1/5を加算し、ユーザは“red”と“on”と“road”を含む発話をした場合には3/5を加算する。その結果、{*n_all[i]*, *n_fea[i]*}の二項組を、ノード*i*に対する選択優先度とする。

20

【0078】

上記の構成により、前記画像認識プロセスによる学習結果を反映したグラフ構造を算出基準とする事になり、当該学習結果を選択優先度に反映する事が出来る。例えば、図3Aの記載の、ステップ30-20から30-28を含む特徴とユーザの発話が一致する場合には、代表ノードに当該特徴に関するノードが追加されるので、前記ステップにより算出された選択優先度が変化する。なお、選択優先度の算出は当該手法には限らない。例えばリンクの重みを考慮しても良い。また、図4B(C)ではノード(74)とノード(75)を他のノードを同じ重みとしてノード数をカウントしたが、当該ノード(74)とノード(75)は強関係にあるとして、1つのノードとしてカウントしても良い。この様にノード間の関係を考慮しても良い。

30

【0079】

ステップ30-39の抽出可能になった全特徴群の記述の生成では、前記選択優先度の第1項の値が大きな順に並べたノード群の中で、第2項が値「1」以上のノードを選び、後述の会話エンジン430を利用して、音声による再確認をユーザに対して行う事が出来る。当該第2項は、ステップ(S16)にて規定値との関係から算出している。即ち、前記代表ノードの非参照数から算出している。例えばステップ(S16)の規定値を「2」にした場合には、2以上の複数のユーザがリンクしている(即ち一度はユーザの着目対象になっている)代表ノードを選び出す。即ちユーザに対して再確認をする候補に加える事を意味している。以上記載の手順により、ステップ30-38の共起物体・事象の抽出による当該対象候補群の中から、よりユーザの意になかった対象を選び出す事が可能になる。

40

【0080】

なお、前記選択優先度に係る二項組の値は、前記組み合わせの利用手段以外を用いても良い。例えば、前記二項組で表現された選択優先度を2次元ベクトルとして正規化して比

50

較しても良い。また、例えば、代表ノードに係る部分グラフにある特徴量ノード、図4B(C)の例ではノード(91)の対応クラス内での代表特徴量(例えば、Visual Word辞書110-10における特徴量)との距離を考慮して、前記選択優先度を算出しても良い。

【0081】

さらに、前記再確認において、ユーザが規定時間無言の場合には、ユーザの意になった対象を認識した可能性と見做して、カメラ画像のアップロードを終了(30-50)しても良い。

【0082】

図5を用いて、本発明の一実施形態に係る知識情報処理サーバシステム300における機能ブロックを説明する。本発明では画像認識システム301、生体認証部302、インタレストグラフ部303、音声処理部304、状況認識部305、メッセージ保管部306、再生処理部307、ユーザ管理部308から構成しているが、これらに限定されず、そのいくつかを選択して構成しても良い。

【0083】

上記音声処理部304部は、ユーザが装着したヘッドセットシステム200が拾うユーザの発声を、音声認識システム320を利用して発話単語列に変換する。また、後述の再生処理部306からの出力を、音声合成システム330を利用して当該ユーザに前記ヘッドセットシステムを通して音声として通知する。

【0084】

次に図6Aから図6Eを用いて、本発明の一実施形態における画像認識システム301の機能ブロックを説明する。前記画像認識システムでは、ヘッドセットシステム200からの画像に対して、一般物体認識、特定物体認識、シーン認識等の画像認識処理を行う。

【0085】

最初に図6Aを用いて、本発明の一実施形態における画像認識システム301の構成例を説明する。画像認識システム301は、一般物体認識システム106、シーン認識システム108、特定物体認識システム110、画像カテゴリデータベース107、シーン構成要素データベース109、及びマザーデータベース(以下MDBと略す)111で構成される。一般物体認識システム106は、一般物体認識部106-01、カテゴリ検出部106-02、カテゴリ学習部106-03、及び新規カテゴリ登録部106-04とで構成され、シーン認識システム108は、領域抽出部108-01、特徴抽出部108-02、重み学習部108-03、及びシーン認識部108-04とで構成され、特定物体認識システム110は、特定物体認識部110-01、MDB検索部110-02、MDB学習部110-03、及び新規MDB登録部110-04とで構成され、画像カテゴリデータベース107は、カテゴリ分類データベース107-01、及び不特定カテゴリデータ107-02で構成され、シーン構成要素データベース109は、シーン要素データベース109-01、及びメタデータ辞書109-02とで構成され、MDB111は、詳細設計データ111-01、付帯情報データ111-02、特徴量データ111-03、及び不特定物体データ111-04とで構成される。画像認識システム301の機能ブロックは必ずしもこれらに限定されるものではないが、これら代表的な機能について簡単に説明する。

【0086】

一般物体認識システム106は、画像中に含まれる物体を一般的な名称、或いはカテゴリで認識する。ここでいうカテゴリは階層的であり、同じ一般物体として認識されているものでも、さらに細分化されたカテゴリ(同じ椅子でも4本足の「椅子」もあれば、全く足の無い「座椅子」の様なものまで含まれる)や、さらに大域的なカテゴリ(椅子も机もタンスも含めて、これらは全て「家具」のカテゴリとして大分類される)としても分類及び認識が可能である。カテゴリ認識は、この分類を意味する「Classification」、即ち既知のクラスに物体を分類するという命題であり、カテゴリはまたクラスとも呼ばれる。

10

20

30

40

50

【 0 0 8 7 】

一般物体認識プロセスにおいて、入力画像中の物体と参照物体画像との比較照合を行った結果、それらが同一形状であるか類似形状である場合、あるいは極めて類似した特徴を併せ持ち、他のカテゴリが有する主要な特徴において明らかに類似度が低いと認められる場合に、認識された物体に対し対応する既知のカテゴリ（クラス）を意味する一般名称を付与する。それらの各カテゴリを特徴付ける必須要素を詳細に記述したデータベースがカテゴリ分類データベース 1 0 7 - 0 1 であり、それらのいずれにも分類する事が出来ない物体は、不特定カテゴリデータ 1 0 7 - 0 2 として一時的に分類し、将来の新たなカテゴリ登録、あるいは既存カテゴリの定義範囲の拡大に備える。

【 0 0 8 8 】

一般物体認識部 1 0 6 - 0 1 では、入力された画像中の物体の特徴点から局所特徴量を抽出し、それらの局所特徴量が予め学習によって得られた所定の特徴量の記述と似ているか似ていないかを相互に比較して、前記物体が既知の一般物体であるかどうかを判別するプロセスを実行する。

【 0 0 8 9 】

カテゴリ検出部 1 0 6 - 0 2 では、一般物体認識可能となった物体がどのカテゴリ（クラス）に属するかを、カテゴリ分類データベース 1 0 7 - 0 1 との照合において特定あるいは推定し、その結果、特定カテゴリにおいてデータベースに追加あるいは修正を加える様な追加の特徴量が見出された場合には、カテゴリ学習部 1 0 6 - 0 3 において再学習した上で、カテゴリ分類データベース 1 0 7 - 0 1 の前記一般物体に関する記述をアップデートする。また一旦、不特定カテゴリデータ 1 0 7 - 0 2 とされた物体とその特徴量が別に検出された他の不特定物体の特徴量と極めて類似であると判定された場合には、それらは新たに発見された同一の未知のカテゴリ物体である可能性が高いとして、新規カテゴリ登録部 1 0 6 - 0 4 において、カテゴリ分類データベース 1 0 7 - 0 1 にそれらの特徴量が新規登録され、当該物体に対し新たな一般名称が付与される。

【 0 0 9 0 】

シーン認識システム 1 0 8 では、入力画像全体あるいは一部を支配している特徴的な画像構成要素を、性質の異なる複数の特徴抽出システムを用いて検出し、それらをシーン構成要素データベース 1 0 9 に記載されているシーン要素データベース 1 0 9 - 0 1 と多次元空間上で相互に参照する事で、各々の入力要素群が当該特定シーン内に検出されるパターンを統計処理により求め、画像全体あるいは一部を支配している領域が当該特定のシーンであるかどうかを認識する。併せて、入力画像に付帯しているメタデータ群と、シーン構成要素データベース 1 0 9 に予め登録済みのメタデータ辞書 1 0 9 - 0 2 に記載されている画像構成要素とを照合し、シーン検出の精度を一段と向上させる事が可能となる。領域抽出部 1 0 8 - 0 1 では、画像全体を必要に応じて複数の領域に分割して、領域毎にシーン判別を可能にする。例えば、都市空間内のビルの壁面や屋上に設置した監視カメラからは、交差点や数多くの店舗のエントランス等の複数のシーンを見渡す事が出来る。特徴抽出部 1 0 8 - 0 2 は、指定した画像領域内における検出された複数の特徴点の局所特徴量、色情報や物体の形状等、利用可能な様々な画像特徴量から得られる認識結果を後段の重み学習部 1 0 8 - 0 3 に入力し、各々の要素が特定のシーンにおいて共起する確率を求め、シーン認識部 1 0 8 - 0 4 に入力して最終的な入力画像に対するシーン判別を行う。

【 0 0 9 1 】

特定物体認識システム 1 1 0 は、入力された画像から検出された物体の特徴を、予め M D B 1 1 1 内に収納されている特定物体群の特徴と逐次照合し、最終的に物体を同定処理（ I d e n t i f i c a t i o n ）する。地球上に存在する特定物体の総数は膨大で、それら全ての特定物体との照合を行う事はおよそ現実的ではない。従って、後述する様に、特定物体認識システムの前段において、予め一定の範囲内に物体のカテゴリや探索範囲を絞り込んでおく必要がある。特定物体認識部 1 1 0 - 0 1 では、検出された画像特徴点における局所特徴量と、学習によって得られた M D B 1 1 1 内の特徴パラメータ群とを相互に比較し、前記物体がどの特定物体に当て嵌まるかの判別を統計処理により判別する。 M

10

20

30

40

50

ＤＢ１１１には、その時点で入手可能な当該特定物体に関する詳細なデータが保持されている。一例として、それら物体が工業製品であるならば、詳細設計データ１１１－０１として設計図やＣＡＤデータ等から抽出された物体の構造、形状、寸法、配置図、可動部、可動範囲、重量、剛性、仕上げ等、物体を再構成し製造する為に必要な基本情報等がＭＤＢ１１１内に保持される。付帯情報データ１１１－０２には、物体の名称、製造者、部品番号、日時、素材、組成、加工情報等、物体に関する様々な情報が保持される。特徴量データ１１１－０３には、設計情報に基づいて生成される個々の物体の特徴点や特徴量に係る情報が保持される。不特定物体データ１１１－０４は、その時点ではどの特定物体にも属していない不明な物体等のデータとして、将来の解析に備えＭＤＢ１１１内に暫定的に収納される。ＭＤＢ検索部１１０－０２は、当該特定物体に対応する詳細データを検索する機能を提供し、ＭＤＢ学習部１１０－０３は、適応的かつ動的な学習プロセスを通して、ＭＤＢ１１１内の当該物体に係る記載内容に対し追加・修正を行う。また一旦、不特定物体として不特定物体データ１１１－０４とされた物体も、その後に類似の特徴を有する物体が頻繁に検出された場合、新規ＭＤＢ登録部１１０－０４により、新たな特定物体として新規登録処理される。

10

【００９２】

図６Ｂに、本発明の一実施形態における一般物体認識部１０６－０１のシステム構成、及び機能ブロックの実施例を示す。一般物体認識部１０６－０１の機能ブロックは必ずしもこれらに限定されるものではないが、代表的な特徴抽出手法としてＢａｇ－ｏｆ－Ｆｅａｔｕｒｅｓ（以下、ＢｏＦと略す）を適用した場合の一般物体認識手法について、以下に簡単に説明する。一般物体認識部１０６－０１は、学習部１０６－１０、比較部１０６－１１、ベクトル量子化ヒストグラム部（学習）１１０－１１、ベクトル量子化ヒストグラム部（比較）１１０－１４、及びベクトル量子化ヒストグラム識別部１１０－１５で構成され、学習部１１０－１６は、局所特徴量抽出部（学習）１１０－０７、ベクトル量子化部（学習）１１０－０８、Ｖｉｓｕａｌ　Ｗｏｒｄ作成部１１０－０９、及びＶｉｓｕａｌ　Ｗｏｒｄ辞書（ＣｏｄｅＢｏｏｋ）１１０－１０とで構成される。

20

【００９３】

ＢｏＦは、画像中に現れる画像特徴点を抽出し、その相対位置関係を用いずに物体全体を複数の局所特徴量（Ｖｉｓｕａｌ　Ｗｏｒｄ）の集合体として表現し、それらを学習によって得られたＶｉｓｕａｌ　Ｗｏｒｄ辞書（ＣｏｄｅＢｏｏｋ）１１０－１０と比較照合して、それら局所特徴量の構成がどの物体に最も近いかを判別する。

30

【００９４】

図６Ｂを用いて、本発明の一実施形態における一般物体認識部１０６－０１における処理を説明する。学習部１０６－１０を構成する局所特徴量抽出部（学習）１１０－０７により得られた多次元の特徴ベクトルは、後段のベクトル量子化部（学習）１１０－０８によって一定次元数の特徴ベクトル群にクラスタ分割され、Ｖｉｓｕａｌ　Ｗｏｒｄ作成部１１０－０９で各々の重心ベクトルを元に、特徴ベクトル毎にＶｉｓｕａｌ　Ｗｏｒｄが生成される。クラスタリングの手法として、ｋ－ｍｅａｎｓ法やｍｅａｎ－ｓｈｉｆｔ法が知られている。生成されたＶｉｓｕａｌ　Ｗｏｒｄは、Ｖｉｓｕａｌ　Ｗｏｒｄ辞書（ＣｏｄｅＢｏｏｋ）１１０－１０に収納され、それを基に入力画像から抽出された局所特徴量を相互に照合し、ベクトル量子化部（比較）１１０－１３においてＶｉｓｕａｌ　Ｗｏｒｄ毎にベクトル量子化を行う。その後、ベクトル量子化ヒストグラム部（比較）１１０－１４において、全てのＶｉｓｕａｌ　Ｗｏｒｄに対するヒストグラムを生成する。

40

【００９５】

当該ヒストグラムの各ピンの総数（次元数）は通常数千から数万と多く、入力画像によっては特徴の一致が全くないヒストグラムのピンも数多く存在する一方、特徴の一致が顕著なピンもあり、それらを一括してヒストグラムの全ピンの値の総和が１になる様な正規化処理を行う。得られたベクトル量子化ヒストグラムは、後段のベクトル量子化ヒストグラム識別部１１０－１５へと入力され、一例として代表的な識別器であるＳｕｐｐｏｒｔ　Ｖｅｃｔｏｒ　Ｍａｃｈｉｎｅ（以下ＳＶＭと呼称）において、物体の属するクラス、

50

即ち当該対象が如何なる一般物体であるかを認識処理する。ここでの認識結果は、前記 Visual Word 辞書に対する学習プロセスとしても利用可能である。また、他の手法（メタデータや集合知の利用）から得られた情報も、同様に前記 Visual Word 辞書に対する学習フィードバックとして利用が可能で、同一クラスの特徴を最も適切に記述し、且つ他のクラスとの分離度を良好に保つ様に、適応的な学習を継続する事が重要となる。

【0096】

図6Cに、本発明の一実施形態における前記一般物体認識部106-01を含む一般物体認識システム106全体の概略構成ブロック図を示す。一般物体（クラス）は様々なカテゴリに属していて、それらは多重的な階層構造を成している。一例を挙げると、人間は「哺乳類」という上位カテゴリに属し、哺乳類は「動物」というさらに上位のカテゴリに属している。人間はまた、髪の色や目の色、大人か子供か？といった別のカテゴリでも認識が可能である。これらの認識判断を行うには、カテゴリ分類データベース107-01の存在が欠かせない。これは人類の「知」の集積庫であり、将来の学習や発見によって、そこにさらに新たな「知」が加わり継続的な進化が図られるものでもある。一般物体認識部106-01で同定されたクラス（およそ人類がこれまでに識別している全ての名詞の総数に及ぶ）は、様々な多次元的且つ階層的な構造体として、当該カテゴリ分類データベース107-01内に記述されている。継続的な学習において認識された一般物体は、カテゴリ分類データベース107-01と照合され、カテゴリ検出部106-02で所属カテゴリが認識される。その後、カテゴリ学習部106-03に当該認識結果が引き渡され、カテゴリ分類データベース107-01内の記述との整合性がチェックされる。一般物体認識された物体は、時に複数の認識結果を内包する場合が多い。例えば「昆虫」であると認識した場合に、目の構造や手足の数、触角の有無、全体の骨格構造や羽の大きさ、胴体の色彩や表面のテクスチャ等でも新たな認識・分類が可能で、前記カテゴリ分類データベース107-01内の詳細記述を基に照合される。カテゴリ学習部106-03では、これらの照合結果を基に、カテゴリ分類データベース107-01への追加・修正が必要に応じて適応的に行われる。その結果、既存カテゴリのいずれにも分類出来ない場合、「新種の昆虫」である可能性も高いとして、新規カテゴリ登録部106-04がこれらの新規物体情報をカテゴリ分類データベース107-01内に登録する。一方、その時点で不明な物体は、不特定カテゴリデータ107-02として、将来の解析や照合に備え一時的にカテゴリ分類データベース107-01内に収納される。

【0097】

図6Dに、本発明の一実施形態における入力画像に含まれるシーンを認識判別する、シーン認識システム108の本発明における代表的な実施例をブロック図で示す。学習画像及び入力画像からは、一般に複数の物体が認識可能となるケースが多い。例えば、「空」「太陽」「地面」等を表す領域と同時に、「木」や「草」そして「動物」等の物体が同時に認識可能となる場合、それらが「動物園」なのか「アフリカの草原」なのかは、全体の景色やそれ以外に発見される物体との共起関係等から類推する事になる。例えば、檻や案内板等が同時に発見され多くの見物客で賑わっていれば、そこが「動物園」である可能性が高まるが、全体のスケールが大きく、遠くに「キリマンジャロ」の様な雄大な景色を臨み、様々な動物が混在して草原上にいる様な場合には、そこが「アフリカの草原」である可能性が一気に高まる。この様な場合、さらに認識可能な物体や状況、共起事象等を知識データベースであるシーン構成要素データベース109に照合し、より総合的な判断を下す必要も出てくる。例えば、全画面の9割が「アフリカの草原」を指し示していると推定されても、後述の図22に記載の事例における手順と共に、それらが矩形の枠で切り取られ全体が平面状であれば、ポスターや写真である確率が極めて高くなる。

【0098】

シーン認識システム108は、領域抽出部108-01、特徴抽出部108-02、強識別器（重み学習部）108-03、シーン認識部108-04、及びシーン構成要素データベース109から構成され、特徴抽出部108-02は、局所特徴量抽出部108-

05、色情報抽出部108-06、物体形状抽出部108-07、コンテキスト抽出部108-08、及び弱識別器108-09から108-12とで構成され、シーン認識部108-04は、シーン分類部108-13、シーン学習部108-14、及び新規シーン登録部108-15で構成され、シーン構成要素データベース109は、シーン要素データベース109-01、及びメタデータ辞書109-02で構成される。

【0099】

領域抽出部108-01は、背景や他の物体の影響を受けずに目的とする物体の特徴を効果的に抽出する為に、対象画像に係る領域抽出を行う。領域抽出手法の例として、グラフベースの領域分割法(Efficient Graph-Based Image Segmentation)等が知られている。抽出された物体画像は、局所特徴量抽出部108-05、色情報抽出部108-06、物体形状抽出部108-07、及びコンテキスト抽出部108-08に各々入力され、それらの各抽出部から得られた特徴量が弱識別器108-09から108-12において識別処理され、多次元の特徴量群として統合的にモデリングされる。それらモデリング化された特徴量群を、重み付け学習機能を有する強識別器108-03に入力し、最終的な物体画像に対する認識判定結果を得る。前記の弱識別器の例としてSVM、強識別器の例としてはAdaBoost等が代表的である。

【0100】

一般に入力画像には複数の物体や、それらの上位概念である複数のカテゴリが含まれている場合が多く、人間はそこから一目で特定のシーンや状況(コンテキスト)を思い浮かべる事が出来る。一方、単独の物体や単一のカテゴリのみを提示された場合、それだけで入力画像がどういうシーンを表わしているのかを判断するのは困難である。通常は、それらの物体が存在している周囲の状況や相互の位置関係、また各々の物体やカテゴリの共起関係(同時に出現する確率)が、当該シーンの判別に対して重要な意味を持ってくる。前項で画像認識可能となった物体群やカテゴリ群は、シーン要素データベース109-01内に記述されているシーン毎の構成要素群の出現確率を基に照合処理され、後段のシーン認識部108-04において、係る入力画像がいかなるシーンを表現しているのかを統計的手法を用いて決定する。

【0101】

これとは別の判断材料として、画像に付帯しているメタデータも有用な情報源となり得る。しかし、時には人間が付したメタデータ自体が、思い込みや明らかな誤り、或いは比喩として画像を間接的に捉えている場合等もあり、必ずしも当該画像中に存在する物体やカテゴリを正しく表わしているとは限らない場合がある。このような場合にも、前記画像認識システムを備えた知識情報処理サーバシステムから抽出可能な当該対象に係る共起事象等を参考に総合的に判断し、最終的な物体やカテゴリの認識処理が行われる事が望ましい。また、一つの画像からは複数のシーンが得られる場合も多い。例えば、「夏の海」であると同時に「海水浴場」であったりもする。その場合は、複数のシーン名が当該画像に付される。さらに画像に付すべきシーン名として、例えば「夏の海」或いは「海水浴場」のいずれがより適当であるかは、当該画像のみからでは判断が難しく、前後の状況や全体との関係、各々の要素群の共起関係等を参考に、それらの要素間の関連性を記述した知識データベース(未図示)を基に最終的に判断が必要な場合もある。

【0102】

図6Eに、本発明の一実施形態における特定物体認識システム110のシステム全体の構成例、及び機能ブロックを示す。特定物体認識システム110は、一般物体認識システム106、シーン認識システム108、MDB111、特定物体認識部110-01、MDB検索部110-02、MDB学習部110-03、及び新規MDB登録部110-04とで構成され、特定物体認識部110-01は、二次元写像部110-05、個別画像切り出し部110-06、局所特徴量抽出部(学習)110-07、ベクトル量子化部(学習)110-08、Visual Word作成部110-09、Visual Word辞書(Code Book)110-10、ベクトル量子化ヒストグラム部(学習)110-11、局所特徴量抽出部(比較)110-12、ベクトル量子化部(比較)110

- 13、ベクトル量子化ヒストグラム部（比較）110-14、ベクトル量子化ヒストグラム識別部110-15、形状特徴量抽出部110-16、形状比較部110-17、色情報抽出部110-18、及び色彩比較部110-19とで構成される。

【0103】

一般物体認識システム106により、対象物体の属するクラス（カテゴリ）が認識可能になった時点で、物体がさらに特定物体としても認識可能か？という絞り込みのプロセスに移る事が出来る。クラスが或る程度特定されないと、無数の特定物体群からの検索を余儀なくされ、時間的にもコスト的にも実用的とは言えない。これらの絞り込みプロセスには、一般物体認識システム106によるクラスの絞り込み以外にも、シーン認識システム108の認識結果から当該対象の絞り込みを行う事も有用となる。また特定物体認識システム110から得られる特徴量を用いて、さらなる絞り込みが可能になるだけでなく、物体の一部にユニークな識別情報（商品名とか、特定の商標やロゴ等）が認識可能な場合、或いは有用なメタデータ等が予め付されているケースでは、さらなるピンポイントの絞り込みも可能となる。

【0104】

それら絞り込まれたいくつかの可能性の中から、複数の物体候補群に係る詳細データや設計データをMDB検索部110-02がMDB111内から順次引き出し、それらを基に入力画像との適合プロセスが実行される。物体が工業製品でない場合や、詳細な設計データ自体が存在していない場合においても、写真等があれば各々検出可能な画像特徴及び画像特徴量を詳細に突き合わせる事で、或る程度の特定物体認識が可能となる。しかし、入力画像と比較画像の見え方が全く同じというケースは稀で、例え同じであっても各々を違う物体として認識してしまう事例もある。反面、物体が工業製品であり、CAD等の詳細なデータベースが利用可能な場合には、一例として二次元写像部110-05が入力画像の見え方に応じMDB111内の三次元データを二次元画像に可視化（レンダリング）する事により、精度の高い特徴量の適合処理を行う事が可能になる。この場合、二次元写像部110-05における二次元画像へのレンダリング処理を全視点方向からくまなく写像して実行する事は、計算時間と計算コストの不要な増大を招く事から、入力画像の見え方に応じた絞り込み処理が必要となる。一方、MDB111を用いた高精度のデータから得られる各種特徴量は、学習プロセスにおいて予め求めておく事が可能である。

【0105】

特定物体認識部110-01では、物体の局所特徴量を局所特徴量抽出部110-07で検出し、ベクトル量子化部（学習）110-08で各々の局所特徴量を複数の類似特徴群に分離した後、Visual Word作成部110-09で多次元の特徴量セットに変換し、それらをVisual Word辞書110-10に登録する。これらは多数の学習画像に対し、十分な認識精度が得られるまで継続して行われる。学習画像が例えば写真等である場合は、画像の解像度不足やノイズの影響、オクルージョンの影響、対象以外の物体から受ける影響等が避けられないが、MDB111を基にしている場合は、ノイズのない高精度のデータを基に理想的な状態で対象画像の特徴抽出を行う事が可能な事から、従来の手法に比べて大幅に抽出・分離精度を高めた認識システムを構成する事が可能となる。入力画像は、個別画像切り出し部110-06で目的とする特定物体に係る領域が切り出された後に、局所特徴量抽出部（比較）110-12において局所特徴点及び特徴量が算出され、予め学習により用意されたVisual Word辞書110-10を用い個々の特徴量毎にベクトル量子化部（比較）110-13にてベクトル量子化された後に、ベクトル量子化ヒストグラム部（比較）110-14にて多次元の特徴量に展開され、ベクトル量子化ヒストグラム識別部110-15にて、物体が当該学習済み物体と同一か、似ているか、それとも否かが識別判断される。識別器の例として、SVM（Support Vector Machine）が広く知られているが、他にも識別判断の重み付けを学習の上で可能にするAdaBoost等も有効な識別器として広く活用されている。これらの識別結果は、MDB学習部110-03を通じてMDB自体への追加修正、或いは新たな項目の追加というフィードバックループにも利用可能となる。対象が依然と

して未確認となる場合には、新規MDB登録部110-04に保留され、次なる解析再開に備える。

【0106】

また、局所特徴量のみならず、検出精度をさらに向上させる目的で、物体の形状特徴を利用する事も有用となる。入力画像から切り出された物体は、形状特徴量抽出部110-16を経由して形状比較部110-17に入力され、物体の各部の形状的な特徴を用いた識別が行われる。その識別結果はMDB検索部110-02にフィードバックされ、それによりMDB111に対する絞り込み処理が可能となる。形状特徴量抽出手段の例として、HoG(Histograms of Oriented Gradients)等が知られている。形状特徴は、またMDB111を用いた二次元写像を得る為の多視点方向からのレンダリング処理を大幅に減らす目的でも有用となる。

10

【0107】

また、物体の色彩的な特徴やテクスチャ(表面処理)も、画像認識精度を上げる目的で有用である。切り出された入力画像は、色情報抽出部110-18に入力され、色彩比較部110-19で物体の色情報、あるいは当該テクスチャ等の抽出が行われ、その結果をMDB検索部110-02にフィードバックする事で、MDB111においてさらなる絞り込み処理を行う事が可能となる。これら、一連のプロセスを通じて、特定物体認識処理をより効果的に行う事が可能となる。

【0108】

次に、図7を用いて、本発明の一実施形態における生体認証部302の処理手順340を説明する。ユーザが前記ヘッドセットシステム200を装着する事で(341)、以下の生体認証処理が始まる。ユーザと前記知識情報処理サーバシステムとの間の通信において、個々のユーザに対応する生体認証情報や、個々のユーザのプロファイル等の個人情報をやり取りする場合には、通信途中でのデータの抜き取りや改竄等の不正な行為からの強力な保護が必須になる。そこで、まず上記生体認証システムとの間で、強力な暗号化通信路を確立する(342)。ここではSSL(Secure Sockets Layer)や、TLS(Transport Layer Security)等の技術(例えば、<http://www.openssl.org/>)を用いる事が可能になるが、他の同様の暗号化手法を導入しても良い。次に、前記ヘッドセットシステムに具備された生体認証センサ204から、生体認証情報345を取得する。生体認証情報には、前記ヘッドセットシステムを装着するユーザの外耳部や鼓膜における静脈パターン情報等を用いる事が出来るが、これらを選択して組み合わせても良いし、これらには限らない。前記生体認証情報はテンプレートとして、前記生体認証システムに送付される。図7のステップ355は、前記生体認証システム側での処理を説明している。ステップ356にて、当該テンプレートを知識情報処理サーバシステム300にユーザ登録する。ステップ357にて、当該テンプレートから署名+暗号化関数 $f(x, y)$ を生成し、ステップ358にて前記関数を当該ヘッドセットシステムに返す。ここで、 $f(x, y)$ における“ x ”は署名暗号化されるデータであり、“ y ”は署名暗号化の際に用いる生体認証情報である。判断345では、前記関数を入手出来たかどうかを確認し、YESの場合には当該ヘッドセットシステムと前記知識情報処理サーバシステム間の通信に前記関数を利用する(346)。判断345がNOの場合には、規定回数、前記判断345がNOであるかを判断(349)し、YESの場合には認証エラーをユーザに通知する(350)。当該判断349がNOの場合には、ステップ344から処理を繰り返す。その後、ステップ(347)で規定時間待ってから、ループ(343)を繰り返す。ユーザが当該ヘッドセットシステムを取り外した場合、或いは前記認証エラーの場合には、前記生体認証システムとの間の暗号化通信路を切断する(348)。

20

30

40

【0109】

図8Aに、本発明の一実施形態におけるインタレストグラフ部303の構成例を示す。本実施例においては、グラフデータベース365へのアクセスを、グラフデータベース365、及びユーザデータベース366への直接アクセスとして記述しているが、具体的な

50

実装においては、システムを利用中のユーザに係るインタレストグラフ適用処理の高速化を図る目的で、グラフ記憶部 360 はグラフデータベース 365 内に収納されているグラフ構造データの中から必要な部分のみ、及びユーザデータベース 366 内に記載の当該ユーザに係る必要な部分情報を自らの高速メモリ上に選択的に読み出し、内部にキャッシュする事が可能である。

【0110】

グラフ演算部 361 は、前記グラフ記憶部 360 から部分グラフの抽出、又は前記ユーザに係るインタレストグラフの演算を行う。関連性演算部 362 は、ノード間の関連性に関して、 $n(>1)$ 次繋がりノードの抽出、フィルタリング処理、及びノード間のリンクの生成・破壊等を行う。統計情報処理部 363 は、前記グラフデータベース内のノードとリンクデータを統計情報として処理し、新規の関連性を発見する。例えば、或る部分グラフが別の部分グラフと情報距離が近く、同じ様な部分グラフが同クラスタ内に分類出来る時は、新しい部分グラフは前記クラスタに含まれる確率が高いと判断可能になる。

10

【0111】

ユーザデータベース 366 は、当該ユーザに関する情報を保持しているデータベースであり、前記生体認証部 302 にて利用される。本発明では、前記ユーザデータベース内部の当該ユーザに対応したノードを中心としたグラフ構造を、当該ユーザのインタレストグラフとして扱う。

【0112】

図 8B を用いて、本発明の一実施形態におけるグラフデータベース (365) に関して説明する。図 8B (A) に、前記グラフデータベース (365) に対する基本アクセス手法を示す。value (371) は、key (370) から locate 演算 (372) により得られる。前記 key (370) は、value (373) をハッシュ (hash) 関数で計算して導出する。例えば、ハッシュ関数に SHA-1 アルゴリズムを用いた場合には、key (370) は 160 ビット長になる。Locate 演算 (372) には、分散ハッシュテーブル (Distributed Hash Table) 法を利用出来る。図 8B (B) に示す様に、本発明では、前記 key と value の関係を (key, {value}) で表現し、前記グラフデータベースへの格納単位とする。

20

【0113】

例えば、図 8B (C) の様に、2つのノードがリンクされている場合、ノード n1 (375) は、(n1, {ノード n1}) で、ノード n2 (376) は、(n2, {ノード n2}) で表現する。n1 や n2 は、各々ノード n1 (375)、ノード n2 (376) の key であり、ノード実体 n1 (375)、ノード実体 n2 (376) を各々 hash 演算し、各々の key を得る。また、リンク l1 (377) は、ノードと同様に (l1, {n1, n2}) で表現し、{n1, n2} を hash 演算する事で、その key (l1) 377 を得る。

30

【0114】

図 8B (D) は、前記グラフデータベースの構成要素の一例である。ノード管理部 380 は前記ノードを、リンク管理部 381 は前記リンクを管理し、各々をノード・リンク格納部 385 に記録する。データ管理部 382 は、ノードに関連したデータをデータ格納部 386 に記録すべく管理する。

40

【0115】

図 9 を用いて、本発明の一実施例における状況認識部 305 の構成例を説明する。図 9 (A) における履歴管理部 410 は、ユーザ毎にネットワーク・コミュニケーションシステム 100 内での利用履歴を管理する。例えば、対象に対する着目を足跡 (フットプリント) として残す事を可能にする。或いは、同じメッセージやつぶやきを繰り返して再生しない様に、前回どこまで再生したか? を記録する。或いは、メッセージやつぶやきの再生を途中で中止した時には、以降の継続再生の為に当該再生を中止した箇所を記録する。例えば、図 9 (B) では、その一実施例として、グラフデータベース 365 に記録されたグラフ構造の一部を示す。ユーザ (417) ノード、対象 (415) ノード、及びメッセー

50

ジやつぶやき(416)ノードは、各々リンクで繋がっている。ノード(416)に再生位置を記録したノード(418)をリンクする事で、ユーザ(417)の着目した対象(415)に関するメッセージやつぶやきの再生を、ノード(418)として記録した再生位置から再開する。なお、本実施例における前記利用履歴はこれらの手法には限定されず、同様の効果が期待出来る他の手法を用いても良い。

【0116】

メッセージ選択部411はユーザ毎に管理され、ユーザが着目した対象に複数のメッセージやつぶやきが記録されていた場合に、適切なメッセージやつぶやきを選択する。例えば、記録された時刻順で再生しても良い。当該ユーザに係るインタレストグラフから、当該ユーザの関心の高い話題を選択的に選び出し再生しても良い。また、当該ユーザを明示的に指定したメッセージやつぶやきを優先的に再生しても良い。なお、本実施例におけるメッセージやつぶやきの選択手順は、これらに限定されない。

10

【0117】

カレント・インタレスト412は、インタレストグラフ部303中の当該ユーザに係る現在の関心を表すノード群として、ユーザ毎に管理され収納されている。前記メッセージ選択部では、前記カレント・インタレストにおける当該ユーザの現在の関心に対応したノード群から上記グラフ構造を探索する事で、当該ユーザが当該時点において関心度の高いノード群を選び出し、後述の会話エンジン430の入力要素とし、それらを一連の文章に変換し再生する。

【0118】

20

当該ユーザの関心の対象や度合いは、例えば後述の図17におけるグラフ構造から求める。図17において、ユーザ(1001)ノードは、ノード(1005)とノード(1002)へのリンクを有している。即ち、このリンクから、「ワイン」と「車」に関心があるとする。前記ユーザが「ワイン」と「車」のどちらに関心が高いかは、ノード「ワイン」から繋がるグラフ構造と、ノード「車」から繋がるグラフ構造とを比較し、ノード数が多い方をより関心が高いとしても良いし、ノードに関連した着目履歴から、着目回数の多い方により関心が高いとしても良いし、前記ユーザが自らの関心の強さを指定しても良いし、これらには限定されない。

【0119】

図10を用いて、本発明の一実施形態におけるメッセージ保管部306に関して説明する。ユーザが発話したメッセージやつぶやき391、及び/又は、ヘッドセットシステム200で撮影した画像421は、当該メッセージ保管部によりメッセージデータベース420に記録される。メッセージノード生成部422は、インタレストグラフ部303から前記メッセージやつぶやきの対象となる情報を取得し、メッセージノードを生成する。メッセージ管理部423は、当該メッセージノードに前記メッセージやつぶやきを関連付けて、前記メッセージやつぶやきを前記グラフデータベース365に記録する。なお、前記ヘッドセットシステムで撮影した画像421を、同様に前記グラフデータベース365に記録しても良い。なお、前記メッセージやつぶやきの記録には、ネットワークを経由してネットワーク上の同様のサービスを利用しても良い。

30

【0120】

40

図11を用いて、本発明の一実施形態における再生処理部307に関して説明する。ユーザのメッセージやつぶやき391を含むユーザの発話は、音声認識システム320で認識処理され、単数の或いは複数の単語列に変換される。前記単語列は、状況認識部304において「ユーザが現在何かの対象に着目している?」「時空間情報を指示している?」「或いは何かの対象に向かい話しかけている?」という状況識別子を付与され、再生処理部306の構成要素である会話エンジン430に送付される。なお、前記状況認識部304の出力としての識別子は、前記の各々の状況には限定されないし、当該識別子を用いない手法で構成しても良い。

【0121】

前記再生処理部307は、前記会話エンジン430、着目処理部431、コマンド処理

50

部 4 3 2、ユーザメッセージ再生部 4 3 3 から構成されるが、これらを選択して構成しても良いし、新たな機能を追加して構成しても良く、当該構成には限定されない。前記着目処理部は、前記状況認識部から対象を着目中であるとの識別子が付された場合に実行され、図 3 A に記載の一連の処理を担う。前記ユーザメッセージ再生部は、対象に残されたメッセージやつぶやき、及び／又は、関連付けられた画像の再生を行う。

【 0 1 2 2 】

図 1 2 を用いて、本発明の一実施形態に係るユーザ管理部 3 0 8 に関し説明する。前記ユーザ管理部は、許可されたユーザの A C L (アクセス制御リスト) をグラフ構造で管理する。例えば、図 1 2 (A) は、一人のユーザ (4 5 1) ノードが、許可 (4 5 0) ノードとリンクを有している状態を示す。これにより、当該ユーザに対し、当該許可ノードとリンクしたノードに対する許可が与えられる。当該ノードがメッセージやつぶやきであれば、それらを再生する事が出来る。

10

【 0 1 2 3 】

図 1 2 (B) は、特定のユーザ群に許可を与えている例である。許可 (4 5 2) ノードは、ユーザグループ (4 5 3) ノードにリンクする、ユーザ 1 (4 5 4) ノード、ユーザ 2 (4 5 5) ノード、及びユーザ 3 (4 5 6) ノードに対し、一括して許可を与えている様子を示している。また、図 1 2 (C) は、全員 (4 5 8) ノードに対し、一括して許可 (4 5 7) ノードが与えられている例である。

【 0 1 2 4 】

さらに、図 1 2 (D) は、特定のユーザ (4 6 0) ノードに対し、特定の時間或いは時間帯 (4 6 1) ノード、特定の場所／地域 (4 6 2) ノードに限り許可 (4 5 9) ノードを与えている様子を示している。

20

【 0 1 2 5 】

なお、本実施例における A C L は、図 1 2 以外の構成をとっても良い。例えば、不許可ノードを導入して、許可を与えないユーザを明示する様に構成しても良い。また、前記許可ノードをさらに詳細化して、再生許可ノードと記録許可ノードを導入する事で、メッセージやつぶやきを再生する場合と記録する場合で、許可の形態を変える様に構成しても良い。

【 0 1 2 6 】

図 1 3 A を用いて、本発明の一実施形態に係るネットワーク・コミュニケーションシステム 1 0 0 を利用するユーザを中心とした、ユースケース・シナリオの一事例を説明する。

30

【 0 1 2 7 】

本発明では、ユーザが装着しているヘッドセットシステム 2 0 0 に具備されたカメラの撮影可能範囲を視野 5 0 3 と呼び、ユーザが主に見ている方向を当該ユーザの主観的な視野：主観視 5 0 2 と呼ぶ。ユーザは、ネットワーク端末 2 2 0 を装着しており、ユーザの発話 (5 0 6 又は 5 0 7) を前記ヘッドセットシステムに組み込まれたマイクロフォン 2 0 1 で拾い、ユーザの主観視を反映した前記ヘッドセットシステムに組み込まれたカメラ 2 0 3 が撮影する映像と共に、前記知識情報処理サーバシステム 3 0 0 側にアップロードされている。前記知識情報処理サーバシステム側からは、前記ヘッドセットシステムに組み込まれたイヤフォン 2 0 2、或いはネットワーク端末 2 2 0 に対し、音声情報、及び映像／文字情報等を返す事が可能になっている。

40

【 0 1 2 8 】

図 1 3 A において、ユーザ 5 0 0 は物体群 5 0 5 を見ているとし、ユーザ 5 0 1 はシーン 5 0 4 を見ているとする。例えば、ユーザ 5 0 0 に関して、図 3 A に記載の手順に従って当該ユーザのカメラの視野 5 0 3 には、物体群 5 0 5 が撮影され、その画像が前記知識情報処理サーバシステム 3 0 0 側にアップロードされる。前記画像認識システム 3 0 1 は、そこから認識可能な特定物体、及び／又は一般物体を抽出する。この時点で当該画像認識システムとしては、ユーザ 5 0 0 がどの対象に着目しているかまでは判らないので、ユーザ 5 0 0 は音声によって、例えば「右上」とか「ワイン」といった様な当該ユーザの音

50

声による着目対象のポインティング操作を行い、前記画像認識システムに当該ユーザが現在物体508に着目している事を通知する。この際、前記知識情報処理サーバシステム側は「アイスペールに入っているワインですね？」という様な当該ユーザが明示的に示していない共起事象を加えた再確認の問い合わせを、当該ユーザ500のヘッドセットシステム200に対し音声で通知する事を可能とする。その再確認通知内容がユーザの意とは違っていた場合には、一例として「違う」と発話して、ユーザの追加的な対象選択指示を前記サーバシステム側に音声で発行し、改めて着目対象の再検出を求めるプロセスを可能にしても良い。或いは、当該ユーザは、前記ネットワーク端末上のGUIにて着目中の対象を直接指定、又は修正しても良い。

【0129】

一例として、ユーザ501はシーン504を見ているが、ユーザの主観的視野503を反映したカメラ画像を、前記画像認識エンジンを備えた知識情報処理サーバシステム側にアップロードする事で、前記サーバシステム側に組み込まれた前記画像認識システムは、対象シーン504はおそらく「山の風景」であろうと推測する。ユーザ501は、前記シーンに対して自らのメッセージやつぶやき、例えば「懐かしい里山だ」を音声で発話する事で、当該ユーザのヘッドセットシステム200経由で、当該メッセージやつぶやきが前記サーバシステム側に当該カメラ映像と共に記録される。その後、他のユーザが異なる時空間内において同様、或いは類似のシーンに遭遇した場合に、当該ユーザに対して、前記ユーザ501のつぶやき「懐かしい里山だ」を前記サーバシステム側からネットワークを介して、当該ユーザに対し音声情報で送り込む事が可能となる。この事例の様に、実際目にした景色自体やその場所等は異なっても、誰でも思い浮かべる共通の印象的なシーン、例えば「夕焼け」等に対して、共有体験に係るユーザコミュニケーションを喚起する事が可能になる。

【0130】

また、ユーザの音声による指示、或いはネットワーク端末220上での直接操作により、上記ユーザが予め設定した条件に従い、上記ユーザ500やユーザ501が特定の対象に対して残したメッセージやつぶやきを、特定のユーザのみ、或いは特定のユーザグループのみ、或いはユーザ全員に対し、選択的に残す事を可能にする。

【0131】

また、ユーザの音声による指示、或いはネットワーク端末220上での直接操作により、当該ユーザが予め設定した条件に従い、当該ユーザ500やユーザ501が特定の対象に対して残したメッセージやつぶやきを、特定の時間、或いは時間帯、及び/又は、特定の場所、特定の地域、及び/又は、特定のユーザ、特定のユーザグループ、或いはユーザ全員に対し、選択的に残す事を可能にする。

【0132】

図13Bを用いて、前記ユースケース・シナリオから導出される、共通の対象への視覚的な好奇心により誘起されるネットワーク・コミュニケーションの事例を説明する。当該視覚的な好奇心により誘起されるネットワーク・コミュニケーションとして、異なる時空間内において、複数のユーザが各々に異なる状況で「桜」を眺めている様子で説明する。偶然桜の花(560)を目にしたユーザ1(550)が、「綺麗な桜だ」とつぶやき、別の時空間でユーザ2(551)が、「桜が満開だ」(561)とつぶやいている。一方で、離れた場所で水面を流れる花びらを見たユーザ4(553)が、「桜の花びらかな？」とつぶやくシーンである。この時、ユーザ3(552)が川面に桜の花びらが舞い落ちる様子を見て(562)、「花筏(はないかだ)だ」とつぶやいたとすると、このつぶやきは、同じ「花筏」を眺めている前記ユーザ4に、前記ユーザ3のつぶやきとして届ける事が可能になる。そして、偶然別の場所で桜の花を眺めているユーザ5(554)に対し、同じ時期に別の場所で「桜」を鑑賞している前記ユーザ1からユーザ4のつぶやきとして送り込む事が可能となり、その結果前記ユーザ5は「そうか、今週はちょうど桜の見頃を迎えているのだな」と、眼前の桜を前に各所の春の到来を感じる事が可能になる。この事例で示す様に、同様の対象やシーンに対し、それらを偶然目にする可能性のある異なる時

10

20

30

40

50

空間内に存在する複数のユーザ間で、共通する視覚的な関心に端を発した、広範な共有ネットワーク・コミュニケーションを誘起する事が可能となる。

【 0 1 3 3 】

図 1 4 で、リンク構造を用いて、本発明の一実施形態におけるユーザ、対象、キーワード、時間、時間帯、場所、地域、メッセージやつぶやき、及び／又は着目した対象が含まれる映像、及び特定のユーザ、特定のユーザ群、或いはユーザ全体をノードとした各要素間の許可の関係を説明する。本実施例では、これらの関係を全てグラフ構造で表現し、グラフデータベース 3 6 5 に記録する。全の関係をノード群とそれら相互のリンクからなるグラフ構造で表現する事で、例えば、リレーショナル・データベース（表構造）等を採用した場合に、事前に全てのノードの存在やノード間の関係や関連性を組み込んでおかなければならない、という実現不可能な要件から本質的に逃れる事が出来る。これらのノード群の中には、時間の経過と共に刻々と変化、及び成長する構造である性質を持っているノード群もある為、事前に全ても構造を予想し、設計しておく事は凡そ困難である。

10

【 0 1 3 4 】

図 1 4 に示す基本形では、対象 6 0 1 は、ユーザ（ 6 0 0 ）ノード、キーワード（ 6 0 2 ）ノード、対象画像特徴（ 6 0 3 ）ノード、時間／時間帯（ 6 0 4 ）ノード、場所／地域（ 6 0 5 ）ノード、メッセージやつぶやき 6 0 7 の各々のノードとリンクしている。対象 6 0 1 には、A C L（ 6 0 6 ）がリンクしている。メッセージやつぶやき（ 6 0 7 ）ノードには、A C L（ 6 0 8 ）ノード、時間／時間帯（ 6 0 9 ）ノード、場所／地域（ 6 1 0 ）ノードがリンクしている。即ち、図 1 4 は、ユーザの着目した対象と、その時間／時間帯、場所／地域、図 3 A に記載の手順 3 0 - 0 1 の過程で抽出された、及び／又は統計情報処理部 3 6 3 にて抽出された、及び／又は後述の会話エンジン 4 3 0 で抽出された、関連するキーワード及び着目対象に残されたユーザのメッセージやつぶやきが、A C L にて許可されている様子を表しているデータ構造である。なお、図 1 4 に記載のグラフ構造は、ノードを追加、或いは削除する事で、前記記載の時間／時間帯、場所／地域、A C L には限定されない情報を記録する事が出来る様に構成しても良い。

20

【 0 1 3 5 】

図 1 5 を用いて、本発明の一実施例における一般物体認識システム 1 0 6、特定物体認識システム 1 1 0、及びシーン認識システム 1 0 8 に係るグラフ構造の抽出プロセスを説明する。まず一般物体認識システム 1 0 6 において当該対象が属するカテゴリを検出（ 9 0 1 ）する。次に、グラフデータベース 3 6 5 からカテゴリノードを検索し（ 9 0 2 ）、当該カテゴリがグラフデータベース 3 6 5 上に存在しているかの確認を行う（ 9 0 3 ）。存在していなければ新規カテゴリノードが追加されグラフデータベースに記録される（ 9 0 4 ）。次に特定物体認識システム 1 1 0 にて特定物体の検出を行い（ 9 0 5 ）、前記グラフデータベース上に既に存在しているかの確認を行う（ 9 0 7 ）。存在していなければ新規当該特定物体ノードを追加し（ 9 0 8 ）、それらをグラフデータベース上に記録する（ 9 0 9 ）。もう一方のパスにおいては、シーン認識システム 1 0 8 においてシーンの検出（ 9 1 0 ）を行い、グラフデータベース 3 6 5 からシーンノードを検索して（ 9 1 1 ）、当該シーンがグラフデータベースに存在しているかの確認を行う（ 9 1 2 ）。存在していなければ当該シーンに係るノードを生成し、前記グラフデータベースに追加する（ 9 1 3 ）。これら一連の処理が終了した時点で、当該カテゴリノード、特定物体ノード、或いはシーンノードに、上記処理を行ったタイムスタンプ情報をグラフデータベース上に追加記録し（ 9 1 4 ）、当該処理を終了する。

30

40

【 0 1 3 6 】

前記図 1 5 に記載のグラフデータベース 3 6 5 への登録の為の新規ノード群生成は、図 3 A に記載のユーザによる再確認処理の際に行っても良い。前記再確認処理では、前記音声認識システムにより抽出された単語列と、前記画像認識システムを備えた知識情報処理サーバシステム側で抽出された各種特徴とを対応付ける事が可能である。一例として、図 4 A に記載のタクシー 5 0 に関し、前記サーバシステム側が、対象 5 1 に対する画像認識結果として「それは赤いバスですか？」とユーザに音声による確認を求めてきた場合、ユ

50

ーザが「いいえ、黄色いタクシーです」と答えたとすると、前記サーバシステム側が再追加的な画像特徴抽出処理を行う事で最終的にタクシー 50 を認識し、当該ユーザに対して「左側の黄色いタクシーを検出しました」と音声による再確認を発行し、それに対し当該ユーザは「そうです」と答えたとする。その結果、前記タクシー 50 に係る検出された全ての特徴群を当該ビュー（シーン）に係る関連ノード群として、当該ユーザが確認した単語「タクシー」「黄色」に係るノード群と共に、前記グラフデータベース 365 内に登録可能になる。

【0137】

また、前記図 15 に記載のカテゴリノード、特定物体ノード、或いはシーンノードにリンクされた上記タイムスタンプと、当該ユーザとの関係付けを行う事が出来る。この場合、当該ユーザの上記着目履歴を、上記獲得したインタレストグラフの部分グラフとして構成する事が出来る。これにより、当該対象に着目した特定の時空間における当該ユーザの着目対象、及びそれらに関連付けられた他のノード群に係る状況を、当該ユーザの音声或いはネットワーク端末 220 上の GUI 経由で、前記画像認識システムを備えた知識情報処理サーバシステム 300 側に問い合わせる事が可能になる。その結果として、前記サーバシステム側から、上記獲得したインタレストグラフの部分グラフにより導く事が可能な特定の時空間における当該着目対象に係る様々な状態を、当該ユーザに音声、或いは文字、写真、図形情報等で通知する事が可能となる。

【0138】

さらに、前記着目履歴は、画像認識システム 301 との協調動作により認識可能になった、特定物体、一般物体、人、写真、或いはシーンの名称に加え、当該操作を行った時空間情報、ユーザ情報、及び対象となる画像情報と共に、グラフデータベース 365 内にグラフ構造として蓄積される。従って前記着目履歴を、前記グラフ構造を直接参照・解析する事が可能な様に構成する事も可能となる。

【0139】

図 16 を用いて、本発明の一実施例における画像認識システムを備えた知識情報処理サーバシステム 300 において実行されるインタレストグラフの獲得に関して説明する。グラフ構造 (1000) は、或る時点でのユーザ (1001) ノードのインタレストグラフである。当該ユーザは特定物体としての車種 A (1003) ノードと車種 B (1004) ノードに興味があり、それらはカテゴリ「車」(1002) ノードに属している。当該ユーザは、また、3つの対象 (特定物体 1006 から 1008) ノードに興味があり、それらはワイン (1005) ノードに属している。次に、ユーザが対象車種 X (1011) ノードに着目したとする。前記対象車種 X (1011) ノードには、画像 (1012) ノードと、他のユーザのメッセージやつぶやき (1013) ノードがリンクしているとする。前記サーバシステムは、前記対象車種 X (1011) ノードを含むグラフ構造 (1010) を車 (1002) ノードに繋ぐリンク (1040) を生成する。一方、前記統計情報処理部 363 により、例えば共起確率を計算する事で、ワイン (1005) ノードに図中の 3本のワイン (1006 から 1008) ノードがリンクされている時には、囲み 1020 にある 2本のワイン (1021 から 1022) ノードも同様にリンクされている可能性が高まる。これにより前記サーバシステムは、当該ユーザに囲み (1020) を提案する事が出来る。その結果、当該ユーザが当該囲み (1020) に興味を示した場合には、それら囲み 1020 にある 2本のワイン (1021 から 1022) ノードをワイン (1005) ノードに直接繋ぐリンク (1041) を生成する事により、当該ユーザ (1001) に係るインタレストグラフを継続的に成長させる事が可能になる。

【0140】

前記図 16 に記載のインタレストグラフの成長がさらに進んだ状態における、ユーザ (1001) ノードを中心とするグラフ構造のスナップショット例を図 17 に示す。図は次の状態を表現している。ユーザ (1001) ノードは、車 (1002) ノードとワイン (1005) ノード以外に、特定のシーン (1030) ノードに関心がある。車 (1002) ノードでは、特に特定物体として車種 A (1003)、車種 B (1004)、及び車種

X (1 0 1 1) の各ノードに関心があり、ワイン (1 0 0 5) ノードでは5種のワイン (1 0 0 6、1 0 0 7、1 0 0 8、1 0 2 1、及び1 0 2 2) ノードに関心がある。特定のシーン (1 0 3 0) ノードは、画像 (1 0 3 1) ノードで代表されるシーンであり、特定の時間 (1 0 3 3) ノードにおいて、特定の場所 (1 0 3 4) ノードで撮影され、ACL (1 0 3 2) ノードにリストされたユーザに対してのみ再生が許されている。車種X (1 0 1 1) ノードは画像 (1 0 1 2) ノードで表現されており、そこに様々なユーザのメッセージやつぶやき (1 0 1 3) ノードが残されていて、ACL (1 0 3 6) ノードにリストされたユーザ群に対してのみ、それらの再生が許可されている。車種Aには、エンジンの仕様と色がノードとして記載されている。以下、5種のワイン (1 0 0 6、1 0 0 7、1 0 0 8、1 0 2 1、及び1 0 2 2) ノードに関しても同様の属性が記載されている。なお、これらのノードの一部は、他のユーザ2 (1 0 3 6) から直接リンクされても良い。

10

【0141】

図18Aを用いて、本発明の一実施形態におけるユーザのメッセージやつぶやきを音声として記録する手段、或いは再生する手段を説明する。まず、ユーザは図3Aに記載の手順で対象を特定 (1 1 0 1) して変数Oにバインドする。次に当該メッセージやつぶやきを記録した時間、或いは再生を可能にする時間/時間帯 (1 1 0 2) を指定して変数Tにバインドし、当該メッセージやつぶやきを記録した場所、或いは再生を可能にする場所/地域 (1 1 0 3) を指定して変数Pにバインドする。次に、それらメッセージやつぶやきを受取る事が可能な受領者を指定 (ACL) して変数Aにバインドする。そして、記録するか再生するかを選択 (1 1 0 5) し、記録処理の場合には当該メッセージやつぶやきの記録手順を実行する (1 1 0 6)。その後、前記4つの変数 (O、T、P、A) から必要なノード群を生成し、グラフデータベース365に記録する (1 1 0 7)。前記選択 (1 1 0 5) が再生処理の場合には、前記4つの変数 (O、T、P、A) から該当するノード群をグラフデータベース365から抽出 (1 1 0 8) し、前記ノードに残されたメッセージやつぶやきを再生する (1 1 0 9) 手順を実行して、一連の処理を終了する。

20

【0142】

図18Bに、図18Aにおける再生時のステップ1102を詳細化して説明する。ユーザは音声によって時間/時間帯を指定するか、或いはネットワーク端末220上のGUIによって直接時間/時間帯を指定するかを選択する (1 1 1 1)。発話による場合には、ユーザは時間/時間帯を発話 (1 1 1 2) し、前記音声認識システム320で認識処理 (1 1 1 3) される。その結果が時間/時間帯であるか確認 (1 1 1 4) し、その結果が正しい場合は、指定時間/時間帯データを変数Tに格納する (1 1 1 6)。違う場合は、時間/時間帯を発話 (1 1 1 2) に戻る。処理を中断 (QUIT) する場合は発話により終了する。一方、前記ネットワーク端末のGUIにより時間/時間帯を指定する場合 (1 1 1 5) には、入力された時間/時間帯を直接前記変数Tに格納 (1 1 1 6) して、一連の終了処理をする。

30

【0143】

図18Cに、図18Aにおける再生時のステップ1103を詳細化して説明する。ステップ1121で、ユーザは音声によって場所/地域を指定するか、ネットワーク端末220上のGUIによって直接場所/地域を指定するかを選択する。発話による場合には、ユーザは場所/地域を発話 (1 1 2 2) し、前記音声認識システム320で音声認識処理 (1 1 2 3) される。その結果が発話された場所/地域であるか確認 (1 1 2 4) し、その結果が正しい場合は、緯度・経度データに変換 (1 1 2 7) してから変数Pに格納する (1 1 2 8)。違う場合は、場所/地域を発話 (1 1 2 2) するに戻る。処理を中断 (QUIT) する場合は発話により終了する。一方、前記ネットワーク端末のGUIにて地図を表示 (1 1 2 5) し、当該ネットワーク端末の画面上で直接場所/地域を指定する場合 (1 1 2 6) し、当該緯度・経度データを変数Pに格納して、一連の処理を終了する (1 1 2 8)。

40

【0144】

図19を用いて、本発明の一実施例として、特定の対象に残された複数のメッセージや

50

つぶやきの中から、受領対象者がそれらメッセージやつぶやきが残された時間或いは時間帯、及び／又は、残された場所或いは地域、及び／又は、残したユーザ名を指定可能にする事で、絞り込み再生する手順を説明する。説明の為の前提条件として、上記受領対象となるユーザは、図3Aに記載した手順に従って当該対象に着目し、予め対応する対象となる各ノード群が選択されているとする(1140)。

【0145】

まず、当該対象に関して再生したい時間／時間帯、及び場所／地域を、図18B、及び図18Cに記載の手順で指定する(1201)。次に、誰の残したメッセージやつぶやきを再生するかを指定する(1202)。次にACLを確認し(1203)、当該指定条件に合致したメッセージやつぶやきに対応するノード、及び／又は、当該映像に対応したノードからデータを取り出す(1204)。この段階では、複数のノードが取り出される可能性があるので、その場合には、当該全ノードに関して次の処理を繰り返し適用する(1205)。

【0146】

次に当該メッセージやつぶやきを残したユーザに係る情報を、受領対象であるユーザに通知するか否かを選択する(1206)。通知する場合は、前記ノードに関連した当該メッセージやつぶやきを残したユーザ情報をグラフデータベース365から入手し、図11に記載の再生処理部306を利用して上記受領対象ユーザが装着しているヘッドセットシステム200、或いは／又は、上記受領対象ユーザに紐付けられているネットワーク端末220に音声、及び／又は、文字で通知する(1208)。通知内容が音声の場合には、ヘッドセットシステムに組み込まれたイヤホンから再生され、文字、写真、及び／又は図形の場合には、前記ネットワーク端末上にそれら音声以外の情報が当該メッセージやつぶやきに同期して表示される(1209)。上記ユーザ情報を通知しない場合には、当該音声ノードから上記メッセージやつぶやき、及び／又は、当該映像ノードから対応する画像データを取り出し、前記再生処理部306を利用して、上記受領対象ユーザが装着しているヘッドセットシステム200、及び／又は、上記受領対象ユーザに紐付けられているネットワーク端末220に、当該メッセージやつぶやきを残したユーザ情報を含まない音声、及び／又は、画像情報として送出し(1207)、それらの一連の処理を、前記取り出された全ノードに関して繰り返し終了する。

【0147】

前記実施例では、ループ(1205)で取り出された全ノードに関して繰り返し処理しているが、他の手段を用いても良い。例えば、状況認識部305を利用して受領対象ユーザに適切なメッセージやつぶやきを選び出し、上記メッセージやつぶやきのみ、及び／又は、付帯している映像情報と共に再生しても良い。前記、時間／時間帯と場所／地域の指定(1201)に係る説明では、過去に記録されたメッセージやつぶやき、及びそれらの基になる画像情報に関して時空間を過去に遡って受領する目的で、特定の時間／時間帯、及び場所／地域を指定する事例を示したが、逆に未来の時間／時間帯及び場所／地域を指定しても良い。その場合には、当該指定された未来の時空間に、当該メッセージやつぶやき、及びそれらの基となる映像情報を“タイムカプセル”に乗せて届ける事が可能になる。

【0148】

また、当該メッセージやつぶやきの再生に同期して、当該着目対象に関する詳細情報を前記ネットワーク端末上に表示しても良い。さらに、受領対象ユーザの主観的視野外となっている対象に向け、前記画像認識システムを備えた知識情報処理サーバシステム側が音声情報により、当該受領対象ユーザに対し、当該メッセージやつぶやきが残された対象に向け頭を動かす、或いは当該対象の存在する方向に向かって移動する等の指示を与え、その結果、受領対象ユーザが当該対象をその主観的視野内に捉えた時に、当該対象に残されたメッセージやつぶやきを再生する様に構成しても良い。また、類似の効果が得られる別的手段を用いても良い。

【0149】

上記、メッセージやつぶやきの再生においては、前記状況認識部の一構成要素である履歴管理部410によって、その時々再生位置が該当するノード内に記録されるので、受領対象ユーザが同一対象に再び着目した場合、以前と同一のメッセージやつぶやきを再び繰り返す事なく、前回の続きから、或いはそれ以降に更新されたメッセージやつぶやきを加え、受領する事を可能とする。

【0150】

次に、図20を用いて、ユーザが眼前のとある対象に着目している事を、前記画像認識システムを活用して前記知識情報処理サーバシステム側に明示的に指し示す一つの方法として、当該ユーザの音声による指示によらず、当該着目対象に向けユーザが直接手指でポインティングする、或いは当該対象に手指で直接触れる事により、当該ユーザのヘッドセ

10

【0151】

図20(A)は、ユーザの主観視(1300)事例である。ここでは、ワイン(1301)、アイスペール(1304)、及びそれ以外の2つの物体(1302、1303)が検出されている。ここでユーザは左側のワイン(1301)に着目している事を前記サーバシステム側に明示的に通知する為に、当該ユーザの手指(1310)でワインを直接指し示している状態を表している。ユーザはまた着目対象であるワイン(1301)に直接触れる事も出来る。また、指で指し示す代わりに、身近にある棒状の道具を使って指し示し、或いはレーザーポインター等の光線を対象に直接照射しても良い。

20

【0152】

図20(B)に、手指(1310)による対象のポインティング手順を説明する。前提条件として、図20(A)の画面はユーザの主観的な視野を反映したカメラからの映像であるとする。まず、画面中から、手指(1310)を含むユーザの手(1311)を検出する。当該カメラ映像を前記画像認識システムにより画像解析し、そこから検出された手指(1310)及び手(1311)の形状特徴から主要なオリエンテーション(1312)を求め、手指(1310)が指し示す方向を抽出する。上記オリエンテーション(1312)の検出は、ネットワーク端末220側に組み込まれた画像認識エンジン224によりローカルに実行しても良い。

【0153】

前記オリエンテーションが検出されれば(1322)、そのベクトル線上にユーザが指し示す対象が存在する可能性が高い。次に、図20(A)の画像から、前記画像認識システム301との協調動作により当該ベクトル線上に存在する物体を検出し(1323)、当該対象物体の画像認識処理を実行する(1324)。当該画像検出及び認識処理は、ユーザのネットワーク端末220側の一構成要素である認識エンジン224上で行う事も可能で、ネットワーク側の負荷を大幅に軽減する事が出来る。また、ユーザによる素早いポインティング操作に対しても、レイテンシ(時間遅れ)の少ない高速なトラッキング処理が可能になる。最終的な画像認識結果は、ネットワークを介して前記画像認識システムを備えた知識情報処理サーバシステム300側に問い合わせする事で確定され、ユーザに当該認識対象の名称等が通知される(1325)。当該ポインティング対象の画像認識結果がユーザの意にかなえば当該ポインティング処理を終了し(1325)、結果がユーザの意と異なる場合は、追加の指示要求を発行(1327)してステップ(1322)に戻り、引き続きポインティング操作を続ける。同様に、当該ユーザが着目対象のポインティングを明示的に確認しなかった場合に、当該検出結果がユーザの意図通りではなかったと推定して上記の処理を繰り返す、或いは無言の同意と見做して当該検出処理を終了するかを予め設定しておく、或いは前後の流れから、或いは個々のユーザの癖を学習する事により、適応的に当該判断内容を振り分ける事が出来る様に構成しておく事が可能である。これらのユーザによる確認にはユーザの音声による指示を用いるが、それに代わる同様の効果をもたらす手段を用いても良い。

30

40

【0154】

50

また、当該ユーザにおける前記一連のポインティング操作の過程で、前記画像認識システムを備えた知識情報処理サーバシステム300と当該ユーザの間で、インタラクティブなコミュニケーションを行う事が可能である。例えば図20(A)の画像において、前記オリエンテーション1312が指し示す方向が前記1302上に向かった時に、「対象は1302ですか？」と前記知サーバシステムが当該ユーザに対し確認する事で、当該ユーザが「そう。けれども、これは一体何かな？」と改めて質問し直す事も可能となる。

【0155】

次に、本発明の一実施例において、前記ヘッドセットシステム200に具備された位置情報センサ208を用い、当該ヘッドセットシステムの移動状態を都度検出する事で、当該ヘッドセットシステムを装着したユーザが、或る対象へ着目し始めた可能性を検出する手順を説明する。

【0156】

図21は、当該ヘッドセットシステム200の動作に関しての状態遷移を表している。動作開始(1400)状態は、当該ヘッドセットシステムが一定の静止状態から動き出す状態である。当該ヘッドセットシステムの動きには、当該ヘッドセットシステム自体の並行移動(上下、左右、前後)に加えて、当該ヘッドセットシステム自体の位置はそのまま、ユーザの首振り動作によりその向きを変える(左右を見る、上下を見る)動きを含む。停止(1403)は、当該ヘッドセットシステムが静止している状態である。短時間静止(1404)状態は、一時的に当該ヘッドセットシステムが静止している状態である。長時間静止(1405)状態は、当該ヘッドセットシステムがしばらくの間静止している状態である。当該ヘッドセットシステムが一定の動作状態から静止した場合、停止(1403)状態に遷移(1410)する。停止(1403)状態が一定時間以上続いた場合、短時間静止(1404)状態に遷移(1411)する。短時間制状態(1404)がその後一定時間以上継続し、さらに長時間静止している場合には、長時間静止状態(1405)に遷移(1413)する。短時間静止状態(1404)、或いは長時間静止状態(1405)から当該ヘッドセットシステムが再び動き出すと、再び動作開始(1400)状態に遷移(1412、或いは1414)する。

【0157】

これにより、例えば短時間静止(1404)状態に前記ヘッドセットがある時には、ユーザが何か眼前の対象を着目し始めている可能性があるとして判断して、前記画像認識システムを備えた知識情報処理サーバシステム300側に対し着目開始を予告すると同時に、前記ヘッドセットシステムに組込まれたカメラを自動的に撮影開始状態に投入し、引き続く一連の処理に備えるきっかけとする事が出来る。また、前記ヘッドセットシステムを装着したユーザの言外の反応、例えば首を傾げる(疑問)、首を左右に振る(否定)、首を上下に振る(同意)等の動作を、当該ヘッドセットシステムに具備された位置情報センサ208から検出可能なデータから検出する事も可能になる。これらのユーザが多用する首振りのジェスチャーは、地域の風習やユーザ毎の癖によって異なる可能性がある。従って、前記サーバシステム側で、それらユーザ個々の、或いは地域特有のジェスチャーを学習の上で取得して、当該属性を保持し反映する必要がある。

【0158】

図22に、本発明の一実施例における写真抽出の事例を示す。写真画像は、視点位置に従いアフィン変換された矩形領域に囲われている閉領域と想定し、当該領域内から検出される物体のサイズがその領域外にある物体のサイズと大幅に異なるスケールで存在している場合、或いは特定の領域に含まれる本来立体であるべき一般物体、或いは特定物体から抽出される各特徴点が、ユーザの視点移動に伴う相対位置変移を起こさず、当該特定の閉領域内で平行移動する場合、或いは画像の奥行き情報を直接検出可能なカメラから獲得可能な対象との距離情報、或いは複数のカメラ画像による両眼視差から獲得可能な物体の奥行き情報等が取得可能な場合において、本来立体であるべき物体やシーンに係る特徴点が同一平面上に存在する場合に、当該閉領域が平面的な印刷物や写真である可能性が高いと推定する事が可能となる。似た様な状況として、窓外の景色も同様の条件を満たし得るが

、それが窓であるか平面画像であるかは周囲の状況から或る程度推定可能になる場合もある。また、それらが写真である可能性が高いと推定された場合、それらの写真自体を一つの特定期間と見なして、前記画像認識システムを備えた知識情報処理サーバシステム300側に問い合わせる事で、類似写真の検索が可能になる。その結果、同様或いは類似の写真画像が発見されれば、以降異なる時空間内において同様或いは類似の写真画像を眺めている、或いは眺めた、或いは眺める可能性のある、他のユーザ群を繋ぐことが可能になる。

【0159】

図23A及び図23Bを用いて、本発明の一実施形態における着目対象との会話に関して説明する。前提としてユーザの着目画像をカメラが捉えているとする(1600)。ユーザの主観的視野を反映したカメラ画像から、ネットワーク上の画像認識システム301との協調動作により、図3Aに記載の着目対象の抽出プロセスにより、当該対象となる画像を認識する(1602)。次に、グラフデータベース365から着目対象に関するグラフ構造を抽出し、当該着目対象に残されたメッセージやつぶやきに係るノード群を抽出する(1603)。次に、それらメッセージやつぶやきの受領対象者を指定したACLを確認し(1604)、その結果として上記対象ノード群に関連付けられたメッセージやつぶやきを、当該ユーザのヘッドセットシステム200、或いはネットワーク端末220に、音声、画像、図形、イラスト、或いは文字情報で通知する(1605)事が出来る。

【0160】

本発明では、上記メッセージやつぶやきに対して、当該ユーザが発話(1606)によってさらに着目対象に向かい会話的に話しかける仕組みを提供する。前記発話内容は、前記音声認識システム320との協調動作により認識され(1607)、発話文字列に変換される。当該文字列は会話エンジン430に送られ、当該ユーザに係るインタレストグラフを基に、前記知識情報処理サーバシステム300側の前記会話エンジン430によって、時々最適な話題が選択され(1608)、前記音声合成システム330経由で当該ユーザのヘッドセットシステム201に、音声情報として届ける事が可能になる。これにより当該ユーザは、継続的な音声コミュニケーションを前記サーバシステムとの間で続ける事が可能になる。

【0161】

前記会話内容が、ユーザによる当該着目対象そのものに係る質問等の場合は、前記知識情報処理サーバシステム300が、当該質問に対する応答を前記MDB111内に記載の詳細情報、或いは当該着目対象に係る関連ノード群から引き出し、当該ユーザに音声情報により通知する。

【0162】

逆に、前記サーバシステム側から当該ユーザに対し、当該ユーザのインタレストグラフにその時々話題に係る関連ノード群を辿って継続的な話題を抽出し、タイムリーに提供する事が出来る。その場合には、同じ話題が不必要に繰り返し提供されない様に、当該会話の流れの中で以前触れた事のある話題に係るノード群それぞれに対し、上記会話の履歴情報を記録しておく事で回避が可能になる。また、当該ユーザにとり関心がない話題に不必要に向かう事により、当該ユーザの好奇心が殺がれない様にすることも大事となる事から、当該ユーザに係るインタレストグラフを基に、抽出される話題を選択する事が出来る。上記継続的な会話は、当該ユーザによる発話が続く限り、ステップ1606に戻り繰り返され、当該ユーザの発話がなくなるまで続き(1609)、その後終了する。

【0163】

上記における広範なユーザと前記知識情報処理サーバシステム300間の双方向の会話は、前記インタレストグラフ部303自体の学習パスとしても重要な役割を果たす事が出来る。特に、ユーザが特定の対象、或いは話題に対して頻繁に会話を促す場合には、当該ユーザが当該対象、或いは話題に対し極めて強い関心があるとして、それら関心に係るノードと当該ユーザに係るノードの直接或いは間接のリンクに対し、重み付けを加える事が可能となる。逆に、ユーザが特定の対象、或いは話題に対して継続的な会話を拒む場合に

は、当該ユーザが当該対象、或いは話題に対し興味を失った可能性があるとして、それら対象や話題に係るノードと当該ユーザに係るノードの直接或いは間接のリンクに対し、重み付けを減じる事も可能となる。

【0164】

前記実施例では、ユーザが着目対象をその視野内に捉えてからのステップを、順を追って説明したが、他の実施形態をとっても良い。例えば、図3Aに記載の手順において、途中のステップから当該ユーザと前記知識情報処理サーバシステム300間の双方向の会話を始める様に本実施形態を構成しても良い。

【0165】

図23Bに、本発明の一実施形態における会話エンジン430の一構成例を示す。前記会話エンジンへの入力、対象ノードを中心とするグラフ構造1640と、音声認識システム320からの発話文字列1641である。前者は関連ノード抽出1651により前記対象に関連する情報を取り出し、キーワード抽出1650に送る。ここでは、前記発話文字列と前記情報を基に、オントロジー辞書1652を参照して複数のキーワード群を抽出する。次に、話題抽出1653にて前記複数のキーワード群から1つを選択する。ここでは、同じ会話を繰り返さない為の話題の履歴管理を行う。また、上記キーワード抽出に当たっては、新しい、他のユーザにより参照頻度の高い、或いは当該ユーザの関心の高いキーワード群を優先して抽出するように構成する事も出来る。適切な話題が抽出された後は、反応文生成1654にて会話パターン辞書1655を参照しながら、自然な口語体に変換された反応文が作成1642され、後段の音声合成システム330に引き渡される。

【0166】

本実施例における前記会話パターン辞書1655は、前記キーワード群から想起される文章のルールを記述している。例えば、「Hello!」とのユーザ発話に対しては「I'm fine thank you. And you?」と返答するとか、「I」とのユーザ発話に際しては「you」と返答するか、「I like it.」とのユーザ発話に対しては「Would you like to talk about it?」と返答するといった代表的な会話のルールを記述している。返答のルールには変数を含めて良い。その場合、当該変数はユーザの発話から充当される。

【0167】

前記構成により、前記知識情報処理サーバシステム300側が、当該サーバシステム内に収納された前記インタレストグラフ部303内に記載の内容から、当該ユーザの関心に沿ったキーワード群を選び出し、前記インタレストグラフを基に適切な反応文を生成する事で当該ユーザにとって引き続き会話を続ける強い動機になると同時に、対象と会話しているような感覚を抱くように構成する事も可能になる。

【0168】

また、グラフデータベース365には、自らを含む特定のユーザ、或いは特定のユーザ群、或いはユーザ全体に対応するノード群が記録され、それらは、特定物体、一般物体、人、写真、或いはシーンに関するノード群、及びそれらに対して残されたメッセージやつぶやきを記録したノード群が相互にリンクされ、グラフ構造を構成している。統計情報処理部363により、前記メッセージやつぶやきに関連するキーワード群を抽出し、状況認識部305により選択的に当該ユーザのヘッドセットシステム200、或いはネットワーク端末220に、関連する音声や画像、図形、イラスト、或いは文字情報で通知する様に本実施例を構成しても良い。

【0169】

図24を用いて、本発明の一実施例として2以上のヘッドシステム200が一台のネットワーク端末220に接続された際の、前記ヘッドセットシステム間の協調動作に関して説明する。図24では、4人のユーザが各々ヘッドセットシステム200を装着しており、各々のユーザが見ている方向が図示されている。この際に、共有する前記ネットワーク端末上に位置のキャリブレーションを行うマーカー等を表示し(1701から1704)、それを各ユーザのヘッドセットシステムに組込まれたカメラで常時モニタリングする事

で、各々のユーザの相互の位置関係、及びその動きを把握する事が出来る。或いは、時間軸変調された画像パターンを当該共有ネットワーク端末の表示デバイス上に表示して、それらを各ユーザのヘッドセットシステムに具備されたカメラ映像で捉えた後に復調して、同様の位置関係を求めても良い。これらにより、各々のカメラの視野と視線のキャリブレーション、及び各ユーザのヘッドセットシステムと当該共有ネットワーク端末とのキャリブレーション、及びトラッキング処理を自動的に行う事で、前記ネットワーク端末は各々のユーザの位置を常に知る事が出来る。それにより、当該共有ネットワーク端末上のGUI操作に関して、どのユーザからの入力操作であるかを当該ネットワーク端末側が認識する事が可能になる。それにより、当該共有ネットワーク端末の共有表示デバイス上で、各々のユーザの位置を考慮した、各ユーザに向けたアライメントを有するサブ画面群の表示が可能になる。

10

【0170】

図25を用いて、本発明の一実施例として、前記画像認識システムを備えた知識情報処理サーバシステム300では認識出来なかった不明な着目対象に対し、当該ユーザが当該対象に係る質問をネットワーク上に残す事を可能にし、他のユーザがネットワーク経由でそれらの不明な対象に対する新たな情報や回答を寄せる事で、当該不明となった着目対象を、前記サーバシステム側が、それらユーザ間のやりとりの中から必要な情報を選択抽出し学習する手順を説明する。

【0171】

前記手順1800は、ユーザによる音声入力トリガ1801から始まる。前記音声入力トリガは、ユーザによる特定の単語の発話、マイクが拾う音圧レベルの急変、或いは前記ネットワーク端末部220のGUIによっても良い。また、それらの方法に制限されない。それによりカメラ画像のアップロードが開始され(1802)、音声コマンド待ち(1803)となる。次に、ユーザが着目対象抽出の為のコマンド群を音声により発話する事で、それらが音声認識処理され(1804)、例えば図3Aに記載の手段を使って音声による着目対象のポインティング処理が正しく完了したかが判断される(1805)。上記ポインティング処理が困難で認識対象をうまく指定出来ない場合には(1806)、新たな特徴追加による再試行が可能か判断される(1807)。再試行が可能な場合にはユーザからの音声コマンド入力待ち(1803)に戻り、再試行する。一方、特徴の追加が困難な場合には、ネットワーク上のWikiへの問い合わせを開始する(1808)。

20

30

【0172】

前記問い合わせ処理では、当該問い合わせ対象に係るカメラ画像、及びユーザの音声による質問やコメントをセットにして、ネットワーク上に発行する(1809)。それに対しWikiから新たな情報提供や回答があれば回収し(1810)、その内容を当該ユーザ、或いは多数のユーザ群、及び/又は、前記知識情報処理サーバシステム300側が検証する(1811)。当該検証処理では、寄せられた回答の正当性を判断する。検証に合格すれば、対象を新規登録する(1812)。当該新規登録に当たっては、前記質問、コメント、情報、回答に対応する各ノード群を生成し、当該対象に係るノード群として関連付け、グラフデータベース365に記録する。前記検証に不合格の場合には、保留処理1822を行う。当該保留処理では、ステップ1808或いはステップ1818におけるWikiへの問い合わせ処理が未完了である旨を記録し、前記検証に合格する回答が収集されるまでステップ1810のWikiからの情報・回答収集処理をバックグラウンドで続行する。

40

【0173】

前項ステップ1805にて、対象の音声によるポインティング処理が可能だった場合、当該対象の画像認識プロセスに移行する(1813)。当該画像認識処理は、本実施例では特定物体認識システム110にて特定物体認識を行い、認識出来なかった場合には一般物体認識システム106にて一般物体認識を行い、さらに認識出来なかった場合にはシーン認識システム108にてシーン認識を行う様を図示しているが、これらの各画像認識処理自体は、本事例のように必ずしも直列的に実行せず、各々を個別に並列、或いは各々の

50

認識ユニットの中をさらに並列化して実行しても良い。或いは、その各々を最適化した上で組み合わせても良い。

【0174】

前記画像認識処理が成功し、対象が認識可能となった場合、ユーザに対する音声による再確認のメッセージが発行され(1820)、それをユーザが正しく確認出来た場合には、カメラ画像のアップロードを終了(1821)して前記一連の対象画像認識処理を終了する(1823)。一方、ユーザが正しく確認出来なかった場合には、当該対象は未確認のままであるとして(1817)、ネットワーク上のWikiへの問い合わせが開始される(1818)。Wikiへの問い合わせに際しては、当該問い合わせ対象画像も一緒に発行する(1819)必要がある。ステップ1810では、Wikiから寄せられた新たな情報や回答群に対し、その内容及び正当性を検証する(1811)。検証に合格すれば、対象を登録する(1812)。当該登録に当たっては、前記質問・コメント及び情報・回答に対応するノード群を生成し、当該対象に係るノード群に関連付けてグラフデータベース365に記録する。

10

【0175】

図26を用い、前記ヘッドセットシステム200に具備された位置情報センサ208を利用する一実施例を説明する。前記位置情報センサには、GPS(Global Positioning System:全地球測位システム)を利用しても良いが、それには限定されない。前記位置情報センサで検出された位置情報及び絶対時間を、前記ヘッドセットシステムに具備されたカメラ203が撮影した画像に付加し、前記画像認識システムを備えた知識情報処理サーバシステム300側にアップロードする事で、グラフデータベース365が記録している情報を較正する事が出来る。図26(A)は、当該アップロード前の、前記グラフデータベースの画像504(図13A)に関係するグラフ構造の一実施例である。「太陽」が「真上」であるので、時間帯は昼頃であると推定可能になる。図26(B)は、前記画像アップロード後の、グラフ構造の一例である。「絶対時間」ノードの追加により、当該画像に対応した時刻が正確に確定可能になる。また、上記位置情報センサ208により検出された位置情報自体に内在する誤差を、カメラの撮像画像から前記サーバシステムによる認識結果により較正する事が可能になる。

20

【0176】

さらに、前記画像504が、前記グラフデータベース365内に存在しなかった場合、前記図25における一実施例と同様の手順を用いて、前記画像504に関係する情報をグラフ構造として前記グラフデータベース365に記録する。その際に、前記位置情報と絶対時間を利用して、近傍にいる他のユーザ群に対して、前記画像504に関する質問を発行する事で、ユーザ間の新たなネットワーク・コミュニケーションを誘発する事が可能になり、そこから得られる有用な情報群を、前記画像504に係るグラフ構造に追加する様に前記サーバシステムを構成する事が可能になる。

30

【0177】

さらに、前記画像認識システムを備えた知識情報処理サーバシステム300においてアップロードされた画像中の物体が不審物体として判断された場合には、当該不審物体を画像解析して入手可能になった情報を前記グラフデータベース365に、かかる不審物体に係る情報群として記録する事が出来る。当該不審物体の存在或いは発見を、事前に設定可能な特定のユーザ、或いは機関に速やかに自動通知しても良い。前記不審物体か否かの判断には、予め登録済みの不審物体、或いは平常状態における物体との照合を前記グラフデータベース365との協調動作により行う事が出来る。その他、不審な状況、或いは不審なシーンが検出された場合にも、係る不審な状況、或いはシーンが検出可能になる様に本システムを構成しても良い。

40

【0178】

また、ユーザが予め指定可能な発見対象とした特定物体、一般物体、人、写真、或いはシーンを、ユーザのヘッドセットシステム200に装着したカメラが偶然捉えた場合、当該ヘッドセットシステムに有線或いは無線で接続されるユーザのネットワーク端末220

50

上に、前記画像認識システムを備えた知識情報処理サーバシステム300側からネットワーク経由で予めダウンロードされ常駐可能となっている特定画像検出フィルタ群が、当該特定物体、一般物体、人、写真、或いはシーンの初期的な抽出及び対象の暫定的な認識を行い、その結果としてさらに詳細な画像認識処理が必要となった場合には、ネットワーク経由で前記サーバシステム側にそれらを詳細に問い合わせる事で、探し物や忘れ物等、或いは発見したい対象をユーザが前記サーバシステム側に登録しておく事で、効果的に見つけ出す事が可能になる。

【0179】

なお、当該発見対象の指定には、ユーザのネットワーク端末220上でのGUIを用いても良い。或いは、前記画像認識システムを備えた知識情報処理サーバシステム300側が、特定の発見対象画像に係るデータ、及び必要な検出フィルタ群を前記ユーザのネットワーク端末上にプッシュして、当該サーバシステム側が指定した発見対象を、広範なユーザ間で共同して探索する事が可能になる様に構成しても良い。

10

【0180】

前記特定画像検出フィルタ群を、前記画像認識システムを備えた知識情報処理サーバシステム300側から抽出する一実施事例として、前記サーバシステム内の前記グラフデータベース365内から前記指定された発見対象に係るノード群を部分グラフとして取り出し、当該指定された発見対象に係る画像特徴群を、それら部分グラフを基に抽出する事で、当該対象を検出する為に最適化された前記特定画像検出フィルタ群を獲得する事が可能になる様に構成しても良い。

20

【0181】

また、本発明に係る一実施例として、ユーザが装着しているヘッドセットシステム200とネットワーク端末220を一体として構成しても良い。また、前記ヘッドセットシステムにネットワークに直接接続可能な無線通信システム、及びユーザの視野の一部を覆う形で半透明の表示ディスプレイを組み込み、前記ヘッドセットシステム自体に前記ネットワーク端末の一部、或いは全体の機能を組み込んで一体として構成しても良い。これらの構成により前記ネットワーク端末を利用しなくとも、前記画像認識システムを備えた知識情報処理サーバシステム300側と直接通信する事が可能になる。その際には、前記ネットワーク端末に組み込まれたいくつかの構成要素は、一部統合・修正する必要がある。例えば、電源部227は当該ヘッドセットの電源部213と統合可能になる。また、表示部222も画像出力装置207に統合する事が可能になる。当該ヘッドセットシステムにおける無線通信装置211は、前記ネットワーク端末間の通信を担っていたが、それらもネットワーク通信部223に統合可能になる。その他の画像特徴検出部224、CPU225、及び記憶部226は、当該ヘッドセットに組み込む事が可能になる。

30

【0182】

図28に、サーバとのネットワーク接続が一時的に切断されている状況下における、ネットワーク端末220単体での処理の一実施例を示す。ネットワーク接続の一時的な中断は、トンネル内やコンクリートで覆われた建物内への移動、航空機での移動中等で頻繁に発生する可能性がある。また、様々な理由で電波状況が悪化する場合や、無線基地局毎に設定されているセル最大接続数を超過してしまった場合等に、ネットワーク接続速度が大幅に低下する傾向がある。この様な状況下でも、前記画像認識を行う対象の種類と数を必要最小限度に絞り込み、音声コミュニケーション機能を特定の会話内容に限定する事で、予めネットワーク接続が確立している時に、前記ネットワーク端末側の一時記憶メモリ容量内、或いはフラッシュメモリ等の二次記憶メモリ容量内にユーザが指定可能な限定された数の特定物体、一般物体、人、写真、或いはシーンの検出、判別、及び認識に必要な学習済みの特徴データ群、及び当該限定された数の対象群の検出・認識する為に最適な画像検出・認識プログラムのサブセットを、上記各特徴データ群と共に一体として前記サーバシステム側から前記ネットワーク端末側に予めダウンロードしておく事で、ネットワーク接続が一時的に中断した場合でも一定の基本動作が可能になる様に構成する事が出来る。

40

【0183】

50

上記の機能を実現する為の一実施例を以下に示す。図 28 (A) 及び (F) にユーザが装着するヘッドセットシステム 200、及びユーザのネットワーク端末 220 の主要機能ブロック構成を示す。一般的なネットワーク端末は、内蔵する CPU 226 により様々なアプリケーションがネットワーク・ダウンロード可能なソフトウェアの形で常駐可能となっている。それらの実行可能なプログラム規模や参照可能な情報量或いはデータ量自体は、サーバ上における構成に比べて大幅な制約は課されるものの、前記画像認識システムを備えた知識情報処理サーバシステム 300 側に構築される各種プログラムやデータの実行サブセットを一時的にユーザの前記ネットワーク端末に常駐させる事で、前記の様に最小限度の実行環境の構築が可能となる。

【0184】

図 28 (D) に、サーバ側に構築された画像認識システム 301 の主要機能ユニット構成を示す。この中で、特定物体認識システム 110、一般物体認識システム 106、シーン認識システム 108 においては、本来その要求される画像認識対象として、過去も含め現在に至るまで存在する、或いは存在していた全ての固有名詞 / 一般名詞を付す事が可能な、物体、人、写真、或いはシーン全体に及ぶ。これら無限とも言える種類及び対象に本来は備えなくてはならない事と、今後の継続的な物体や事象の発見や認識対象アイテムの増加に伴う追加学習も必要となり、その全体の実行環境自体は極めて限られた情報処理能力やメモリ容量しか持ち合わせないネットワーク端末の手に到底及ぶものではなく、それらの包括的な機能はネットワークを介しサーバ側の強力なコンピュータ・リソース、及び巨大なデータベースシステム上に置かれる事になる。その上で、その時々で都度必要な機能部分について、非力なクライアント機器でも実行可能な画像認識機能のサブセットや、予め学習済みの知識データ等の必要な部分を、ネットワーク経由で当該ネットワーク端末上に選択的にダウンロードする事で、ネットワーク接続の切断に或る程度備える事が出来る。これには、不測のネットワーク切断に備えると言う目的以外に、サーバ・リソースへのアクセス集中による負荷軽減や、ネットワーク回線の不要なトラフィックを抑制するという実用的な側面もある。

【0185】

これらを実現する一実施形態として、図 28 (D) に示す特定物体認識システム 110、一般物体認識システム 106、シーン認識システム 108 から選択した画像認識プログラムの必要なプログラム群を、ネットワークを介し図 28 (A) に示すネットワーク端末 220 上で実行可能な画像認識プログラム 229 として、認識エンジン 224 上にサーバ側からダウンロードの上で常駐させ、併せて各認識対象に即し必要な学習済みの特徴データ群を画像カテゴリデータベース 107、シーン構成要素データベース 109、及び M D B 111 から抽出し、同様にユーザのネットワーク端末 220 上の記憶部 227 上に選択的に常駐させる。これら対象となる認識対象候補群と、他のユーザによる当該対象候補群に対するメッセージやつぶやきを関連付ける為に、サーバ側の前記画像認識システムを備えた知識情報処理サーバシステム 300 側から、必要な当該対象との関連性を前記グラフデータベース 365 から抽出すると共に、前記メッセージデータベース 420 から必要な会話候補群を抽出し、ネットワークを介し予めユーザのネットワーク端末 220 上のメッセージ管理プログラム 232 上にダウンロードしておく。これらユーザのメッセージやつぶやきの候補群は、限られた容量のメモリを効果的に使用する目的で、圧縮して当該ネットワーク端末 220 上の記憶部 227 内に格納する事が出来る。

【0186】

一方、前記画像認識システムを備えた知識情報処理サーバシステム 300 側との双方向の音声による会話機能については、ネットワーク端末 220 上の音声認識プログラム 230、及び音声合成プログラム 231 により一定の制限下で実行可能になる。その為には前記一実施例において、前記サーバシステム側とのネットワーク接続が確立しているタイミングで、前記サーバシステムを構成する会話エンジン 430 内の音声認識システム 320、音声合成システム 330、及びそれらに対応する知識データベースである音声認識辞書データベース 321、会話パターン辞書 1655 から、必要最小限の実行プログラム群、

10

20

30

40

50

及びデータセットをユーザのネットワーク端末 220 上の記憶部 227 内に予めダウンロードしておく必要がある。

【0187】

上記において、ユーザのネットワーク端末 220 の処理能力、或いは記憶部 227 の記憶容量に十分な余裕がない場合には、予め会話の候補群をネットワーク上の音声合成システム 330 で音声化した後に、圧縮音声データとしてユーザのネットワーク端末 220 上の記憶部 227 上にダウンロードしておいても良い。これにより、ネットワーク接続に一時的に障害が生じて、主要な音声コミュニケーション機能は限定的ではあるが保持する事が可能になる。

【0188】

次に、ネットワークへの再接続時のプロセスについて説明する。ユーザが着目した様々な対象に係るカメラ画像、及び当該対象に対してユーザが残したメッセージやつぶやき等が、関連する様々な情報と共にユーザのネットワーク端末 220 上の記憶部 227 内に一時的に保持されているとする。そこで再びネットワーク接続が復帰した時点で、ネットワーク上の生体認証システム 310 内の生体認証処理サーバシステム 311、及び個々のユーザ毎の詳細な生体認証情報を保持している生体認証情報データベース 312 に対し、当該ユーザのヘッドセットシステム 200 に紐付けられたユーザのネットワーク端末 220 から得られる生体認証データを問い合わせる。その結果、紐付けされた当該ユーザのネットワーク端末 220 と、サーバ側の前記画像認識システムを備えた知識情報処理サーバシステム 300 内にそれまで蓄積されている情報及びデータとの同期処理を行う事で、関連するデータベース群を最新の状態に更新すると共に、ネットワークのオフライン時に先に進んだ会話ポイント等の更新も併せて行う事で、オフラインからオンライン、或いはオンラインからオフラインの状態への移行がシームレスに可能になる。

【0189】

また本発明により、PC やカメラ付きスマートフォン等に代表されるネットワーク端末、或いは前記ヘッドセットシステムから、インターネット経由で前記画像認識システムを備えた知識情報処理サーバシステム 300 側に様々な画像（カメラ画像、写真、動画等）をアップロードする事により、前記サーバシステム側が当該画像、或いは当該画像に内包されている、特定物体、一般物体、人、或いはシーン中から、認識可能になった様々な画像構成要素群に対応するノード群、及び／又は当該画像に付帯するメタデータ、及び／又は当該画像に係るユーザのメッセージやつぶやき、及び／又は当該画像に係るユーザ間のコミュニケーションから抽出可能なキーワード群を、ノード群として抽出する事が可能となる。

【0190】

これら抽出された各ノードを中心とする部分グラフから、前記グラフデータベース 365 に記載の関連ノード群を参照する事で、ユーザが指定可能な特定の対象やシーン、或いは特定の場所や地域に係る画像の選択・抽出を可能にし、それらを基に同様或いは類似の対象やシーンを集めたアルバムの作成、或いは一定の場所や地域に係る画像群の抽出処理を行う事が出来る。その上で、前記サーバシステム側が当該抽出された画像群に係る画像特徴群、或いはメタデータ群を基に、それらが特定の物体を撮影したものである場合には複数の視点方向からの映像、或いは異なる環境下で撮影した映像として集約、或いはそれらが特定の場所や地域に係る画像群であるなら、連続的、及び／又は離散的なパノラマ画像に繋ぎ合わせる事で、様々な視点の移動が可能とする。

【0191】

前記場所や地域を特定可能なパノラマ画像の構成要素群となっている、インターネット経由でアップロードされるそれぞれの画像に付帯しているメタデータ、或いは前記画像認識システムを備えた知識情報処理サーバシステム 300 により認識可能になった当該画像中の特定物体に関し、当該物体が存在していた時点或いは期間をインターネット上の各種知識データベース、或いはインターネットを介して広範なユーザに問い合わせる事で推定或いは獲得し、それら時間軸情報を基に当該画像群を時間軸に沿って振り分け、それら振

10

20

30

40

50

り分けられた画像群を基に、ユーザが指定可能な任意の時点或いは期間における前記パノラマ画像を再構成する事が可能となる。これにより、ユーザは任意の場所や地域を含む、任意の「時空間」を指定して、当該「時空間」上に存在していた現実世界の映像を、前記パノラマ画像として視点移動可能な状態で楽しむ事が出来る様になる。

【 0 1 9 2 】

その上で、特定の対象、或いは特定の場所や地域毎に編成された前記画像群を基に、当該対象に関心が高い、或いは特定の場所や地域に関わりの深いユーザ群を、前記グラフデータベース 3 6 5 を基に抽出し、それら多数のユーザ群による当該対象、或いは特定の場所や地域毎に編成されたネットワーク・コミュニケーションを誘発し、そこから特定の対象、或いは特定の場所や地域に係る様々なコメント、メッセージやつぶやきの共有、或いは参加ユーザによる新規情報の提供、或いは特定の不明・不足・欠落情報の探索要求等を可能にするネットワーク・コミュニケーションシステムが構築可能になる。

10

【 0 1 9 3 】

図 2 9 を用いて、本発明に係る一実施例における前記サーバシステム上にアップロードされた画像群の中から、特定の「時空間」を指定する事によって抽出した 3 枚の写真、写真 (A)、写真 (B)、写真 (C) を事例として示す。ここでは、1 9 0 0 年前半における東京日本橋界隈の様子を示す。

【 0 1 9 4 】

写真 (A) では、手前の「日本橋」に加えて、画面左側中央のランドマーク的な建物として知られている「野村証券」本社ビルが特定物体認識可能になり、また画面左側奥には「倉庫」らしき建物、橋の上には「路面電車」2 両が一般物体認識可能になっている様子を示す。

20

【 0 1 9 5 】

写真 (B) では、別の方向から俯瞰した「日本橋」であり、画面右側に同じく「野村証券」本社ビル、画面左手には「帝国製麻ビル」、また「日本橋」の橋上の装飾的な「外灯」が新たに特定物体認識可能になっている様子を示す。

【 0 1 9 6 】

写真 (C) では、画面左側に、同じ「帝国製麻ビル」と思われる建物がある事から、「野村証券」本社ビル屋上と思われる場所から「日本橋」方面を撮影したシーンである事が判り、画面上部に文字で『日本橋上ヨリ三越呉服店及び神田方面盛観』と読み取れる事からも、「日本橋」「三越呉服店」「神田」の 3 つのキーワード群が抽出可能となり、そこから画面奥の白い大きな建物は「三越呉服店」の可能性が高いと推定可能になっている様子を示す。

30

【 0 1 9 7 】

また、「日本橋」橋梁上に「路面電車」の形状がはっきり写っている事で前記画像認識システムによる精査が可能となり、この「路面電車」が写真 (D) と同じ「1 0 0 0 型」車両であると特定物体認識可能になっている様子を示す。

【 0 1 9 8 】

上記一連の画像認識処理は、前記画像認識システム 3 0 1 内に備わった特定物体認識システム 1 1 0、一般物体認識 1 0 6、シーン認識システム 1 0 8 との協調動作により実行される。

40

【 0 1 9 9 】

図 3 0 を用いて、アップロードされた画像群の中から、ユーザが任意の時空間情報を指定する事によって当該時空間内に撮影された画像群のみを抽出し、それらを基に当該時空間を連続的、或いは離散的なパノラマ画像に再構築して、ユーザが自由に当該空間内で視点の移動を行う、或いは自由に当該空間内で時間の移動が可能な、時空間移動表示システムについて、概略的な実施事例を用いて説明する。

【 0 2 0 0 】

最初に、インターネットを介し前記画像認識システムを備えた知識情報処理サーバシステム 3 0 0 側に、ユーザのネットワーク端末 2 2 0 経由で画像のアップロード (2 2 0 0

50

）が開始される。アップロードされた画像は前記画像認識システム301にて画像認識処理が開始される(2201)。当該画像ファイルに予めメタデータが付与されている場合は、メタデータ抽出処理(2204)が実行される。また、当該画像中に文字情報が発見された場合には、OCR(Optical Character Recognition)等を用いて、文字情報抽出処理(2203)が行われ、そこからメタデータ抽出処理(2204)を経て、有用なメタデータ群を得る。

【0201】

一方、アップロードされた一枚の画像の中から、ユーザのネットワーク端末220上のGUI、或いは図3Aに記載の前記音声による着目対象のポインティング処理により、当該画像中の個々の物体に係る画像の切り抜き(2202)処理を行い、当該対象に対して一般物体認識システム106、及びシーン認識システム108にて画像認識したクラス情報に従いMDB検索部110-02で物体の絞り込み処理を行い、当該画像に関する詳細情報を記述したMDB111を参照して、特定物体認識システム110により当該物体との比較照合処理を行い、最終的に同定された特定物体に関し、前記メタデータ群を参照して、当該画像に時間軸情報が存在するか否かを判別(2205)する。

10

【0202】

当該画像に時間軸情報が存在する場合、画像中の物体群が存在した時間情報をMDB111内の記述から抽出し、参照の上で物体が当該時間内に存在するか否かを判別(2206)する。前記存在が確認された場合は、当該物体以外に画像認識可能になった他の物体について、同様に当該時間内に存在し得ない物体がないかどうか(2207)前記同様にMDB111内の記述から判別し、当該全ての整合性が確認された時点で、当該画像に関する撮影時間の推定(2208)処理が行われる。それ以外の場合は、時間情報が不明(2209)として、当該ノード情報が更新される。

20

【0203】

次に、当該画像に場所に係る情報が存在する場合(2210)、画像中の物体群が存在した場所に係る情報をMDB111内の記述から抽出し、参照の上で物体が当該場所において存在するか否かを判別(2210)する。前記存在が確認された場合は、当該物体以外に画像認識可能になった他の物体について、同様に当該場所において存在し得ない物体がないかどうか(2211)前記同様にMDB111内の記述から判別し、当該全ての整合性が確認された時点で、当該画像に関する撮影された場所の推定(2212)処理が行われる。それ以外の場合は、場所情報が不明(2213)として、当該ノード情報が更新される。

30

【0204】

前記一連の処理に加えて、前記獲得可能になった当該画像自体から抽出可能な、或いは当該画像自体に付帯するメタデータ群と、前記推定可能になった時空間情報とを再度照合し、その整合性が確認された時点で、当該画像全体に係る時空間情報の獲得(2214)が完了し、当該時空間情報を当該画像に係るノードにリンク(2215)する。また上記整合性に齟齬のある場合には、メタデータ自体の誤り、画像認識システムの認識誤り、或いはMDB111内に記載の内容に誤りや不備があるとして、以降の再検証処理に備える。

40

【0205】

これらの時空間情報の付与が行われた画像群に対し、ユーザは任意の時空間を指定して当該条件に合致した画像群を抽出する事が可能になる(2216)。まず、多数の画像群の中から任意の場所(2217)、任意の時間(2218)に撮影された画像群を、当該指定した時空間に係るノードを辿って抽出する(2219)。これら抽出された複数の画像群を基に、画像中の共通の特定特徴点を探索する事で、検出された特定特徴点同士を連続的に繋いでパノラマ画像を再構成(2220)する事が可能になる。この場合、パノラマ画像中に欠落或いは欠損画像がある場合は、MDB111記載の地図、図面、或いは設計図等の利用可能な情報から広範に推定処理する事で、離散的なパノラマ画像として再構成が可能になる。

50

【 0 2 0 6 】

前記一連の時空間情報獲得の為の学習プロセスを、アップロードされる多数の写真（動画を含む）画像に対して、前記画像認識システムを備えた知識情報処理サーバシステム 300 が継続的に行う事により、時空間情報を有する連続的なパノラマ画像が取得可能になる。これにより、ユーザは任意の時間 / 空間を指定して、任意の視点移動、或いは同一空間における任意の時間に係る画像体験（2221）を楽しむ事が可能になる。

【 0 2 0 7 】

図 31 を用いて、本発明に係る一実施例における、ユーザが前記画像認識システムを備えた知識情報処理サーバシステムに対してアップロードした画像に対して、当該ユーザのネットワーク端末上の GUI 操作、或いは前記音声処理によるポインティング操作による当該ユーザが着目した特定物体、一般物体、人、或いはシーンに係る選択抽出処理により、前記サーバシステムが認識した結果を、当該入力画像と共に当該ユーザを含むあらかじめ指定可能な広範なユーザ間で共有可能にすることによるネットワーク・コミュニケーションシステムの構成を説明する。

【 0 2 0 8 】

当該時空間を指定したユーザの視点の移動により発見可能になった特定物体、一般物体、人、或いはシーンに対しても、これまで述べて来た様な特定の着目対象に係る一連のメッセージやつぶやきの記録、及び再生体験が可能になる。

【 0 2 0 9 】

当該ユーザによるアップロードされた画像 2101 は、前記サーバシステムにおいて選択・抽出処理 2103 が行われる。この際に、ユーザは図 3A に記載の手順での選択・抽出処理を実行しても良いし、図 30 に示した選択・抽出コマンドを、GUI 2104 を操作することによって選択・抽出処理を実行しても良い。当該選択・抽出処理により切り出された画像は、画像認識システム 301 において認識処理される。その結果は、インタレストグラフ部 303 において分析・分類・蓄積され、キーワード群や時空間情報と共にグラフデータベース 365 に記録される。当該ユーザは、画像のアップロードに際して、メッセージやつぶやき 2106、或いは文字情報 2105 による書き込みを行っても良い。これら当該ユーザの発したメッセージやつぶやき、或いは文字情報もインタレストグラフ部にて分析・分類・蓄積される。当該ユーザ、或いは当該ユーザを含むユーザ群、或いはユーザ全体は、前記対象に係るキーワード群、及び / 又は時空間情報（2106）を基に、インタレストグラフ部から記録された画像を選択する事が可能であり、当該画像に係る広範なネットワーク・コミュニケーションを誘発させることが出来る。さらに、前記広範なユーザ間のコミュニケーションを、前記サーバシステム側で観察・蓄積し、インタレストグラフ部 303 の 1 構成要素である統計情報処理部 363 において分析することで、当該ユーザ特有の、或いは特定のユーザ群に特有の、或いはユーザ全体に共通の動的な関心や好奇心の在り所とその推移を、上記広範なユーザ群、抽出可能なキーワード群、及び様々な着目対象に係るノード間を繋ぐ動的なインタレストグラフとして獲得する事が可能となる。

【 0 2 1 0 】

[周辺技術]

本発明に係るシステムは、既存の様々な技術と組み合わせる事によって、さらに利便性の高いシステムとして構成する事が可能となる。以下に、例示する。

【 0 2 1 1 】

本発明に係る一実施例として、ユーザの発話をヘッドセットシステム 200 に組み込まれたマイクロフォンが拾い、前記音声認識システム 320 により発話中に含まれる単語列及び構文を抽出した後、ネットワーク上の自動翻訳システムを活用する事で異なる言語に翻訳し、当該翻訳された単語列を前記音声合成システム 330 により音声変換した上で、他のユーザに当該ユーザのメッセージやつぶやきとして伝える事が可能になる。或いは、前記画像認識システムを備えた知識情報処理サーバシステム 300 側からの音声情報を、当該ユーザが指定可能な言語で受け取る事が出来る様に構成する事が出来る。

【 0 2 1 2 】

本発明に係る一実施例として、ユーザのヘッドセットシステムに組込まれたカメラがその視野内に捉えた映像の中から、規定の認識マーカと共に特定の画像変調パターンを抽出した場合、当該信号源の存在をユーザに喚起し、当該信号源が表示装置或いはその近傍にある場合、当該変調された画像パターンを前記認識エンジン 2 2 4 との協調動作により復調する事によって、そこから得られる URL 等のアドレス情報をインターネット経由で参照し、当該表示装置上に表示されている画像に係る音声情報を当該ユーザのヘッドセットシステム経由で送り込む事を可能にする。これにより、ユーザが偶然目にした様々な表示装置から、当該表示画像に係る音声情報を当該ユーザに効果的に送り込む事が可能になる。これにより、電子広告媒体としてのデジタル・サイネージの有効性を一段と高める事が出来る。反面、ユーザが目にする事が出来る全てのデジタル・サイネージから音声情報が一斉に送り届けられると、場合によってはそれらを不要なノイズと感じてしまう可能性もある事から、それぞれのユーザに係る前記インタレストグラフを基に、ユーザ毎に異なる嗜好を反映した広告等のみを選択して、個々のユーザ毎に異なる音声情報として送り届ける事が出来る様に構成しても良い。

10

【 0 2 1 3 】

本発明に係る一実施例として、様々な生体情報（バイタルサイン）をセンシング可能な複数の生体センサ群をユーザのヘッドセットシステムに組み込む事で、当該ユーザが関心を持って着目した対象と、当該生体情報との相関を、前記画像認識システムを備えた知識情報処理サーバシステム 3 0 0 側で統計処理した上で当該ユーザに係る特殊なインタレストグラフとして登録しておく事によって、当該ユーザが当該特定の対象或いは事象に遭遇した場合、或いは遭遇の可能性が高まった場合に、当該ユーザの生体情報値が急変する事態に備える事が出来る様に、前記サーバシステム側を構成する事が可能である。取得可能になる生体情報としては、ユーザの体温、心拍、血圧、発汗、皮膚表面の状態、筋電位、脳波、眼球運動、発声、頭の動き、体の動き等が含まれる。

20

【 0 2 1 4 】

この為の学習パスとして、カメラが捉えたユーザの主観視内に特定の特定物体、一般物体、人、写真、或いはシーンが現れた時に、測定可能な前記生体情報値が一定以上変化する場合、当該ユーザに関わる特異的な反応として係る事態を、前記画像認識システムを備えた知識情報処理サーバシステム 3 0 0 側に通知する事で、当該サーバシステム側は関連する生体情報の蓄積・分析を開始すると同時に、当該カメラ映像の解析を開始し、そこから抽出可能な画像構成要素群に係る事態に関連する可能性のある原因要素群として前記グラフデータベース 3 6 5、及びユーザデータベース 3 6 6 内に登録する事を可能にする。

30

【 0 2 1 5 】

以降、様々な事例で前記学習を繰り返す事で、前記各種生体情報値の変化に係る要因の分析・推定を統計処理から求める事が可能になる。

【 0 2 1 6 】

上記の一連の学習プロセスから、個々のユーザ毎に異なる当該生体情報値の異常な変化の要因となっていると予測可能な特定物体、一般物体、人、写真、或いはシーンに、当該ユーザが再び遭遇する、或いは遭遇する可能性が高いと予測可能な場合、前記サーバシステム側から当該ユーザに対し、ネットワークを介して音声、及び／又は、文字、画像、バイブレーション等で、係る可能性を速やかに通知する様に当該サーバシステムを構成する事が可能となる。

40

【 0 2 1 7 】

さらに、観測可能な前記生体情報値が急変し、ユーザの容体に一定以上の危機の可能性があると推定可能な場合、速やかに当該ユーザに係る事態の確認を求め、当該ユーザから一定の反応が得られない場合、当該ユーザに一定以上の緊急事態が発生した可能性が高いと判断し、予め設定可能な緊急連絡網、或いは特定の機関等に通知する事が可能な様に、前記画像認識システムを備えた知識情報処理サーバシステム 3 0 0 側を構成する事が出来る。

50

【 0 2 1 8 】

本発明に係る生体認証システムにおいて、ユーザが頭部に装着可能な前記ヘッドセットシステムから、ユーザ固有の声紋、静脈パターン、或いは網膜パターン等を取得して生体認証が可能な場合、ユーザと前記画像認識システムを備えた知識情報処理サーバシステム 300 側とを一意にバインドする様に本システムを構成する事が出来る。当該生体認証デバイスはユーザの前記ヘッドセットシステムに組み込み可能な事から、当該ヘッドセットシステムの着脱に合わせて自動的にログイン、ログアウト可能にする様に構成する事も可能になる。これら生体情報を活用した紐付けを常時上記サーバシステム側で監視する事により、異なるユーザによる不正なログイン、不正な利用が排除可能になる。当該ユーザ認証が正常に行われた場合、以下の情報群が当該ユーザにバインドされる。

10

(1) ユーザが設定可能なユーザプロフィール

(2) ユーザの音声

(3) カメラ画像

(4) 時空間情報

(5) 生体情報

(6) その他のセンサ情報

【 0 2 1 9 】

本発明に係る一実施例として、複数のユーザ間で共有される画像に関し、プライバシー保護の観点から、ユーザが予め指定可能なルールに従い、当該ユーザ毎の顔部分、及び/又は、当該ユーザを特定可能な画像の特定部分を、前記画像認識システムを備えた知識情報処理サーバシステム 300 側に組込まれた画像認識システム 301 により抽出及び検出し、それらの特定画像領域に対し、判別不能なレベルにまで自動的にフィルタ処理を施す様に構成する事が出来る。これにより、プライバシー保護を含む一定の閲覧制限を設ける事が可能となる。

20

【 0 2 2 0 】

本発明に係る一実施例として、ユーザが頭部に装着可能なヘッドセットシステムに複数のカメラを設置する事が出来る。この場合、一実施例として複数のカメラに撮像視差を設ける事が出来る。或いは、性質の異なる複数の撮像素子を使って、対象物体までの深度(距離)を直接測定可能な三次元カメラを組み込む様に構成する事も出来る。

その上で、前記画像認識システムを備えた知識情報処理サーバシステム 300 側からの音声による指示により、当該サーバシステムにより指定された特定のユーザに対し、当該サーバシステムが指定した特定の対象、或いは周囲の様子等を、当該サーバシステムが当該ユーザに対して様々な視点から撮影する様に依頼する事で、前記サーバシステム側が当該対象の立体的な把握、或いは周囲の状況等の立体的な把握が容易になると共に、当該画像認識結果により前記サーバシステム内の M D B 1 1 1 を含む関連データベース群の更新が可能となる様に、当該サーバシステムを構成する事が出来る。

30

【 0 2 2 1 】

本発明に係る一実施例として、ユーザが頭部に装着可能なヘッドセットシステムに、指向性を有する深度センサを組み込む事が出来る。これにより、当該ヘッドセットシステムを装着したユーザに近づく人間を含む生体や物体の動きを検知し、前記ユーザに音声で係る事態を通知する事が可能となる。同時に、当該ユーザのヘッドセットシステムに組み込まれたカメラ及び画像認識エンジンを自動的に起動し、不測の物体の急接近に即時に対応可能な様にリアルタイム処理が要求される部分をユーザのネットワーク端末側で、高度の情報処理を必要とする部分に関して前記画像認識システムを備えた知識情報処理サーバシステム 300 側で分担して実行可能にする様にシステムを構成する事で、ユーザに近づく特定の物体、特定の人間、特定の動物等を高速に識別/解析し、その結果を音声情報、或いはパイプレーション等により当該ユーザに速やかに喚起する事が可能となる。

40

【 0 2 2 2 】

本発明に係る一実施例として、ユーザが頭部に装着可能なヘッドセットシステムに、当該ユーザを中心とした周囲、或いはその上部や下部も含めた全方位を撮影する事が可能な

50

撮像システムを組み込む事が出来る。或いは、ユーザの主観的視野外となる後方や側面からの視野を撮影する事が可能な複数のカメラを、当該ユーザのヘッドセットシステムに追加する事が可能となる。この様な構成を採る事により、当該ユーザの主観視野外にあるものの、当該ユーザが特に関心や注意を払わなければならない対象が近傍に存在する場合に、当該ユーザに対して速やかに音声、或いはそれに代わる手段を用いて係る状況の喚起を促す事が可能になるように、当該画像認識システムを備えた知識情報処理サーバシステム300を構成する事が出来る。

【0223】

本発明に係る一実施例として、ユーザが頭部に装着可能なヘッドセットシステムに、以下の様な環境値を測定可能な環境センサ群を任意に組み込む事が可能である。

10

(1) 周囲の明るさ(光度)

(2) 照明や外光の色温度

(3) 周囲の環境騒音

(4) 周囲の音圧レベル

これにより周囲の環境雑音の低減、最適なカメラ露光状態への対応が可能になり、前記画像認識システムの認識精度、及び前記音声認識システムの認識精度を向上させる事が可能になる。

【0224】

本発明に係る一実施例として、ユーザが頭部に装着可能なヘッドセットシステムに、当該ユーザの視野の一部を覆う形で半透明のディスプレイ装置を組み込む事が出来る。或いは、当該ヘッドセットシステムをヘッドマウントディスプレイ(HMD)、或いはスカウター(Scouter)として表示ディスプレイと一体的に構成する事も出来る。この様な表示システムを可能とする装置には、ユーザの網膜に直接画像情報を走査投影するレチナール・センシングと呼ばれる画像投影システム、或いは眼前に配置した半透明の反射板に画像を投影するデバイス等が知られている。上記の様な表示システムを採用する事により、ユーザのネットワーク端末の表示画面に表示される画像の一部、或いは全部を、当該表示デバイス上に映し出す事が可能になり、前記ネットワーク端末をユーザの眼前に取り出す事なく、インターネット経由で直接前記画像認識システムを備えた知識情報処理サーバシステム300側とのコミュニケーションが可能となる。

20

【0225】

本発明の一実施形態としてユーザが頭部に装着可能な前記HMD、前記スカウター、或いはそれらに併設する形態で視線検出センサを具備しても良い。当該視線検出センサには光センサアレイを用いても良く、そこから照射される光線の反射光を計測する事で、当該ユーザの瞳の位置を検出し、当該ユーザの視線位置を高速に抽出する事が出来る。例えば、図27において、点線枠2001はユーザの装着する前記スカウター2002の視野画像であるとする。この時、当該ユーザの視線方向にある対象に対して、視点マーカー2003を重ねて表示しても良い。その場合、前記視点マーカーの位置が当該対象と同位置に表示される様に、ユーザの音声による指示でキャリブレーション可能にする事が出来る。

30

【符号の説明】

【0226】

- 100 ネットワーク・コミュニケーションシステム
- 106 一般物体認識システム
- 107 画像カテゴリデータベース
- 108 シーン認識システム
- 109 シーン構成要素データベース
- 110 特定物体認識システム
- 111 マザーデータベース
- 200 ヘッドセットシステム
- 220 ネットワーク端末
- 300 知識情報処理サーバシステム

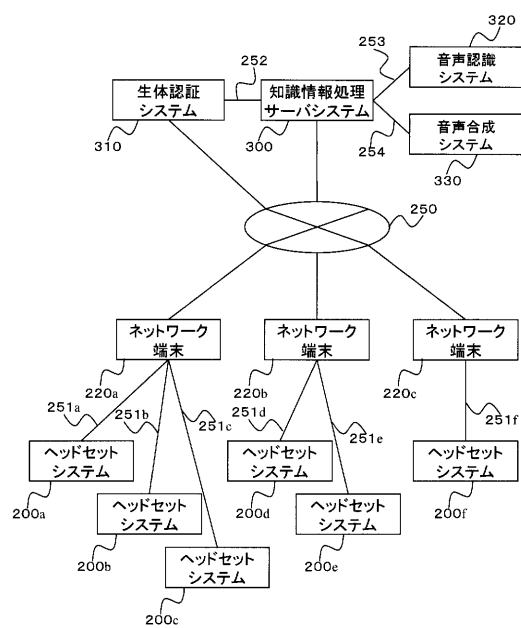
40

50

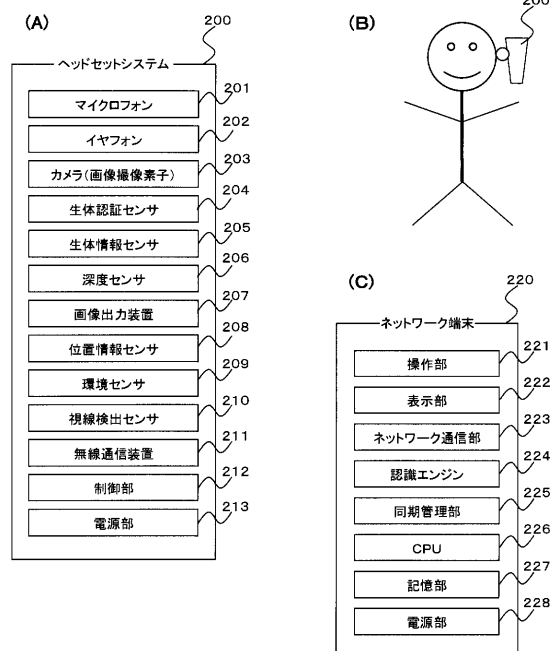
- 3 0 1 画像認識システム
- 3 0 3 インタレストグラフ部
- 3 0 4 状況認識部
- 3 0 6 再生処理部
- 3 1 0 生体認証システム
- 3 2 0 音声認識システム
- 3 3 0 音声合成システム
- 3 6 5 グラフデータベース
- 4 3 0 会話エンジン

【図 1】

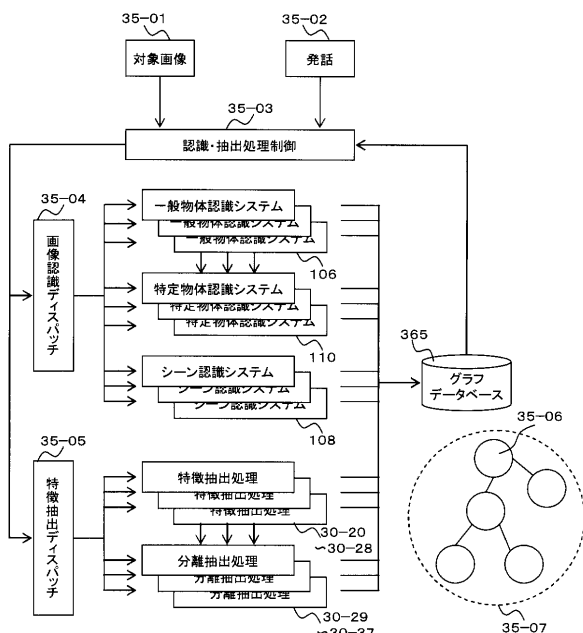
100



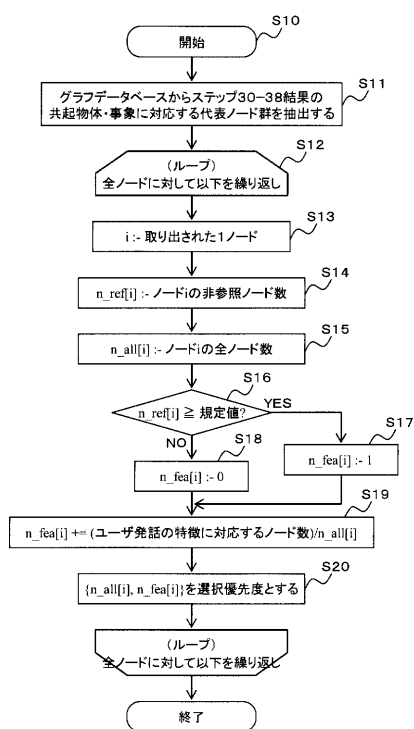
【図 2】



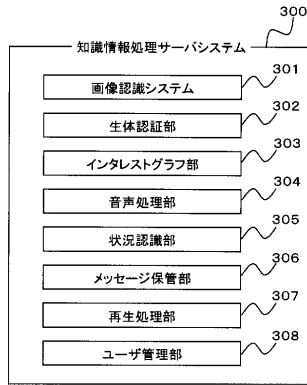
【 ㄨ 3 B 】



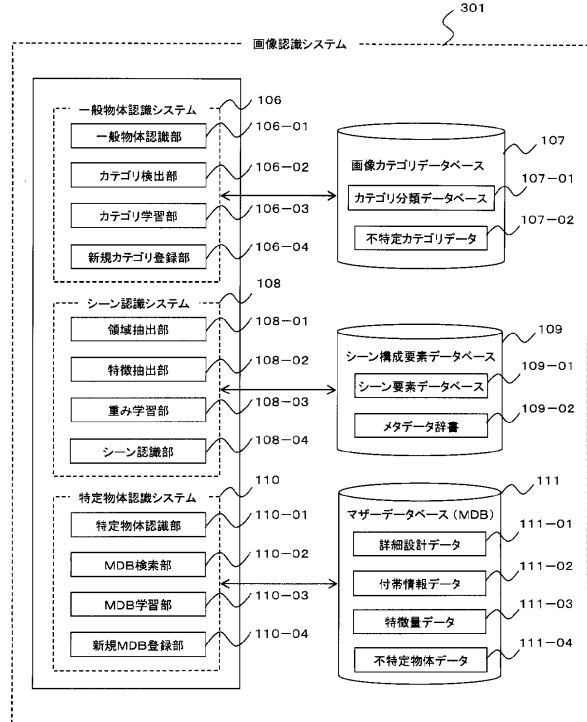
【 ㄨ 4 C 】



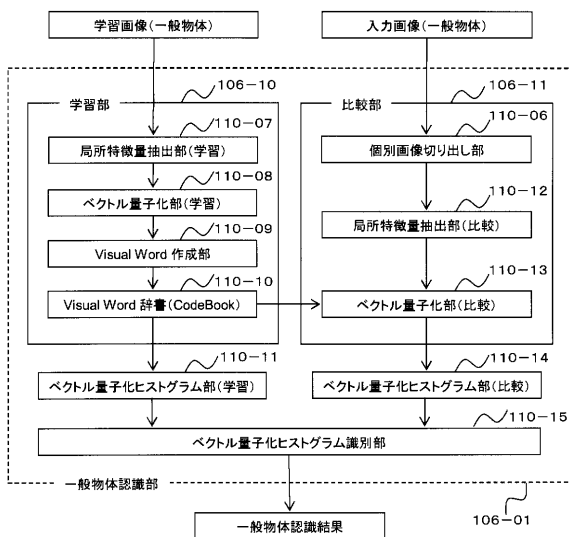
【図 5】



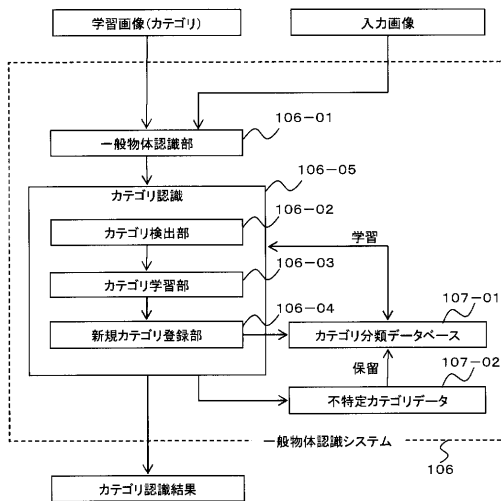
【図 6 A】



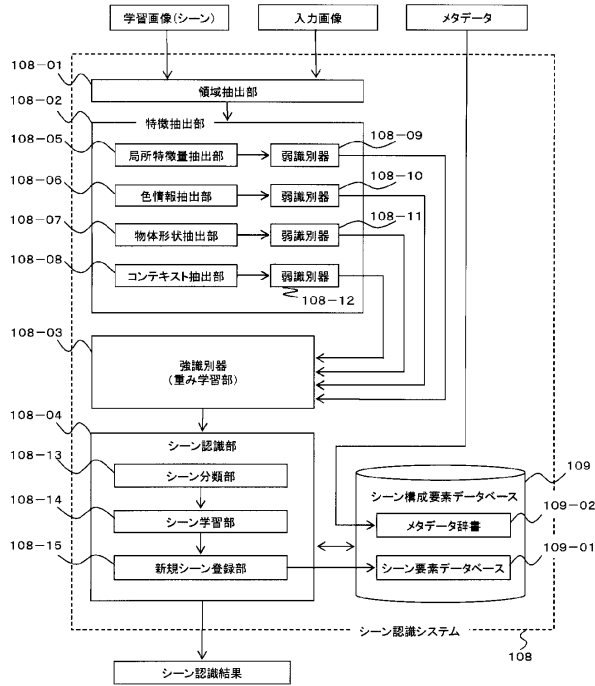
【図 6 B】



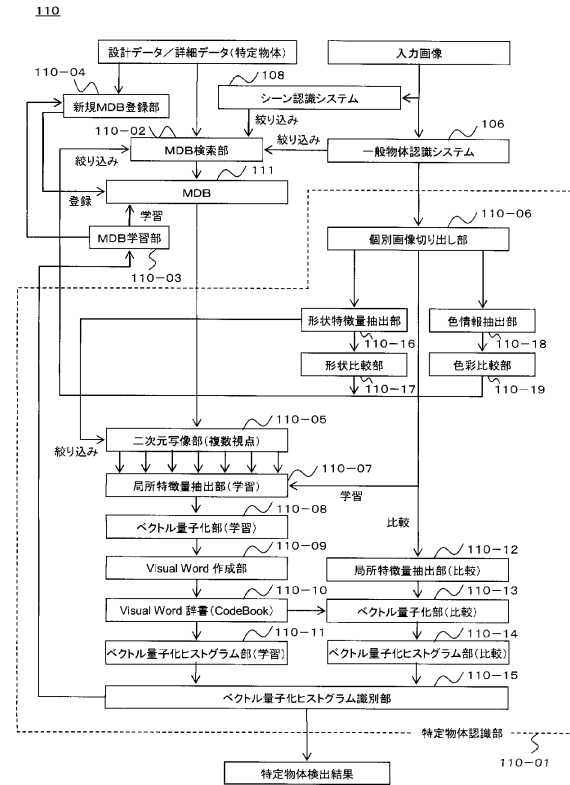
【図 6 C】



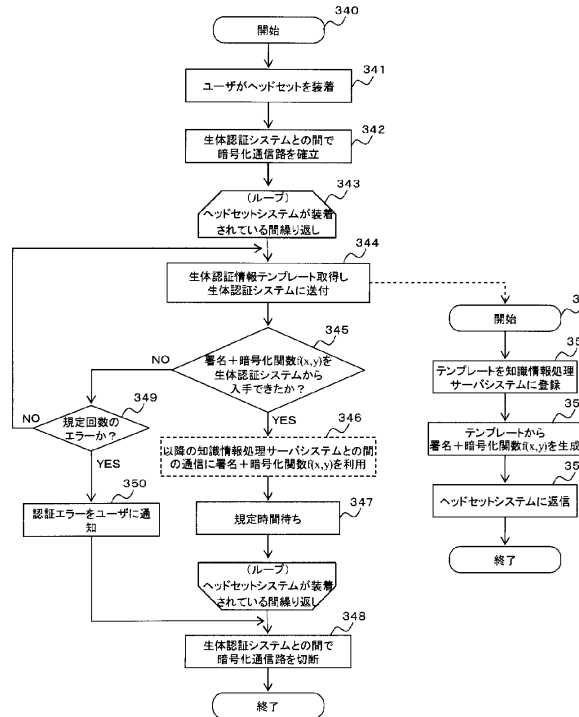
【図 6 D】



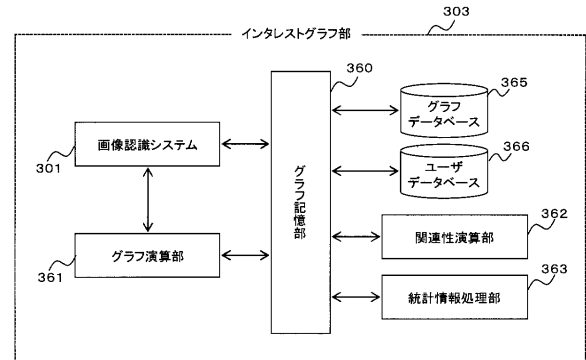
【図 6 E】



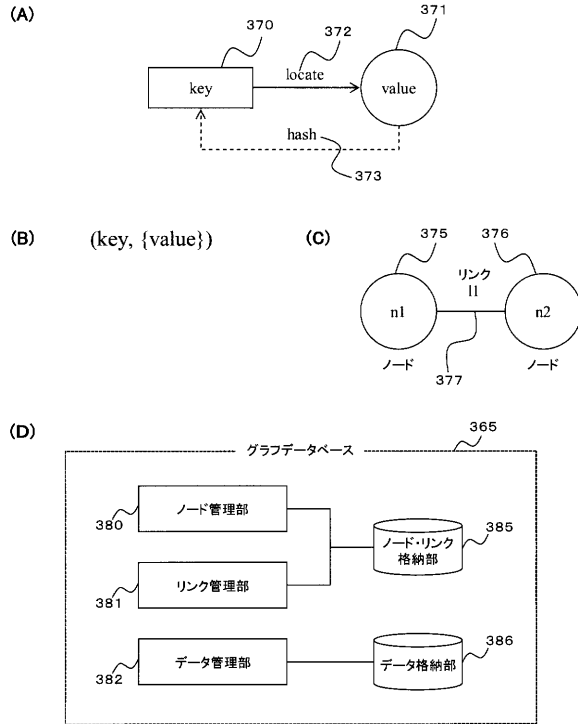
【図 7】



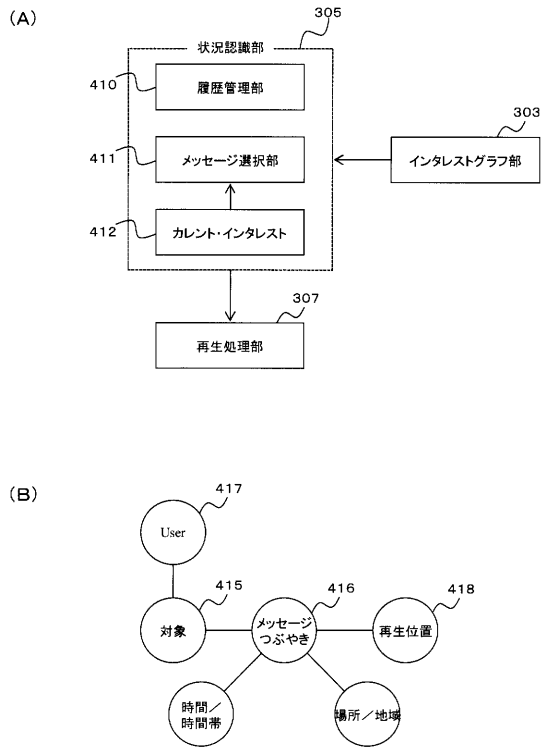
【図 8 A】



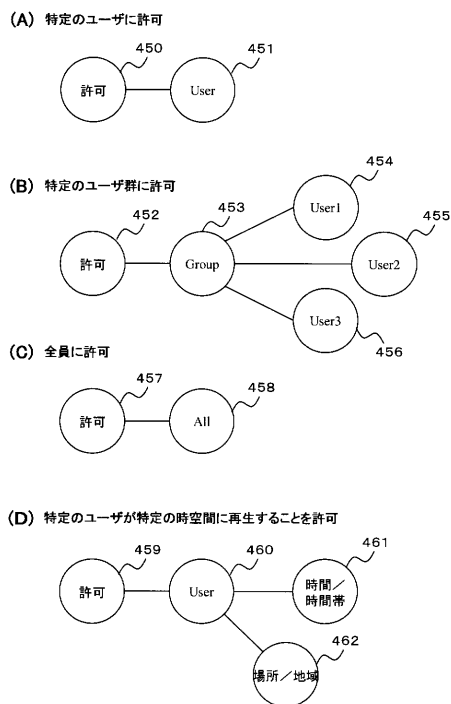
【図 8 B】



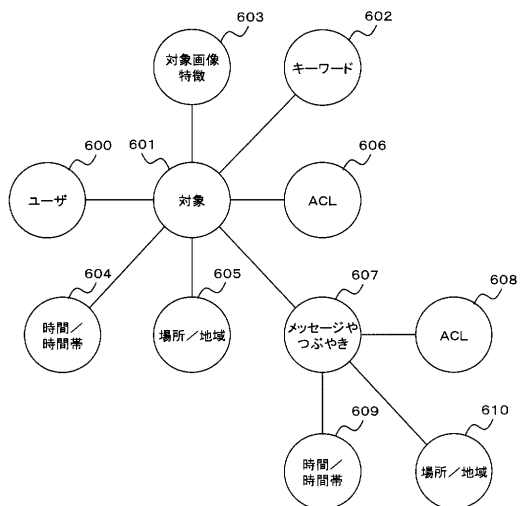
【図 9】



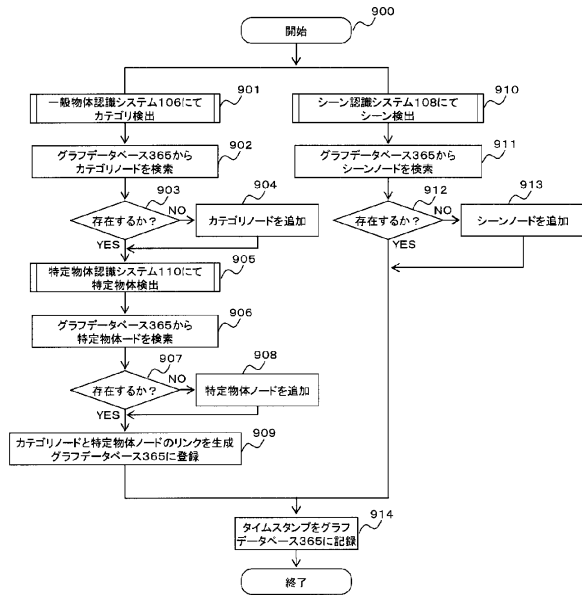
【図 12】



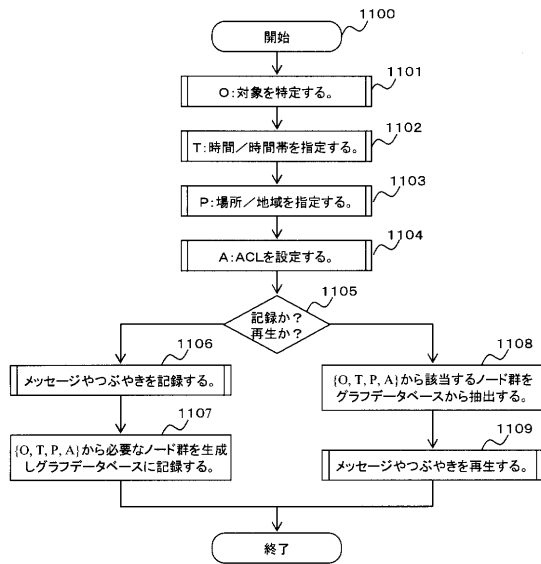
【図 14】



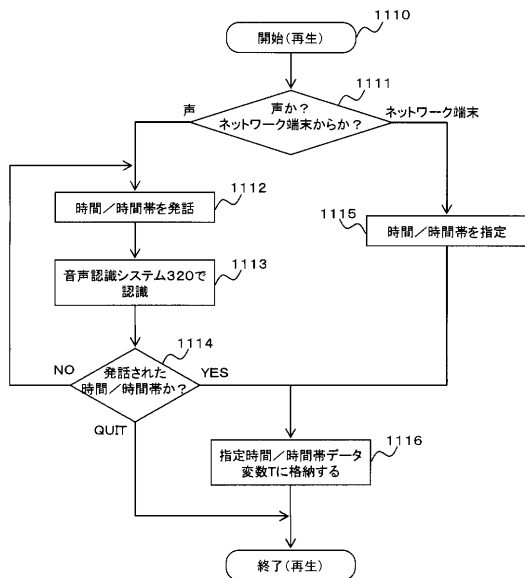
【図15】



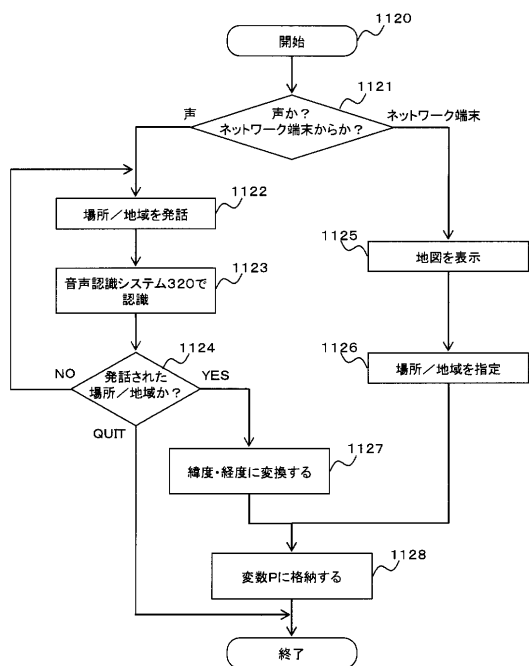
【図18A】



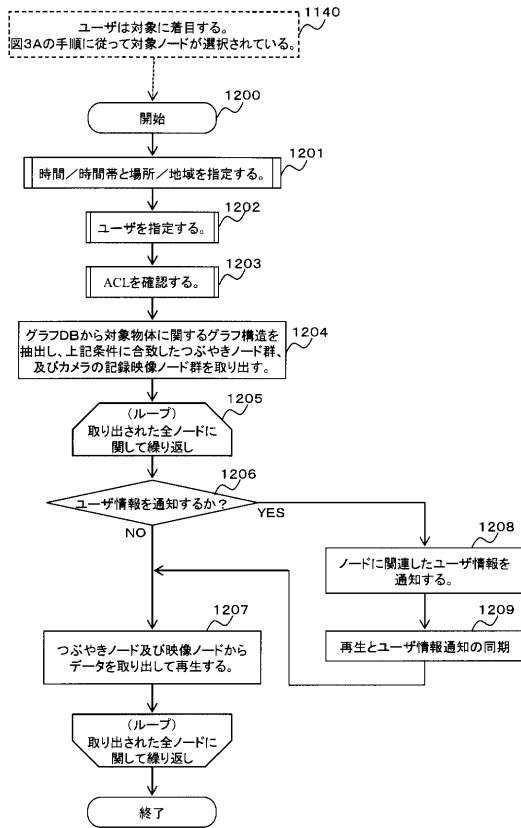
【図18B】



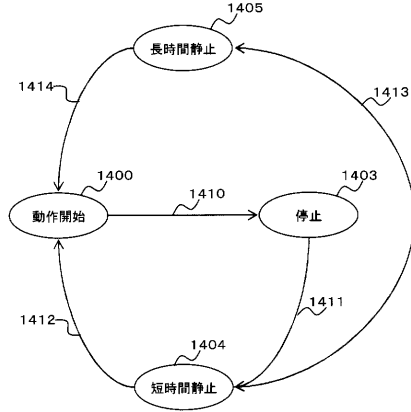
【図18C】



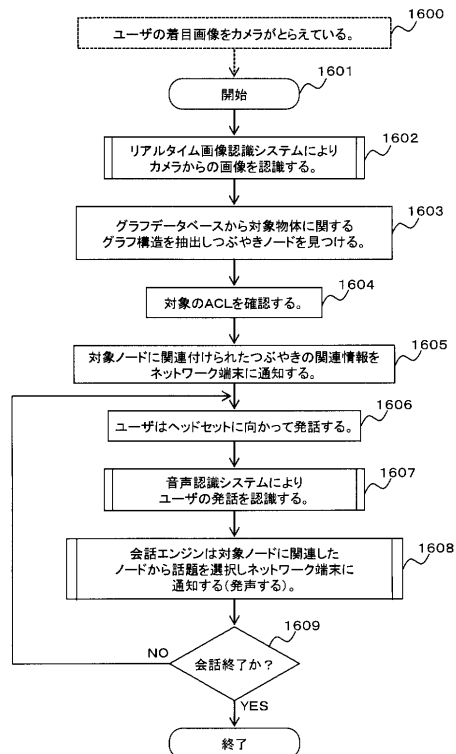
【 図 1 9 】



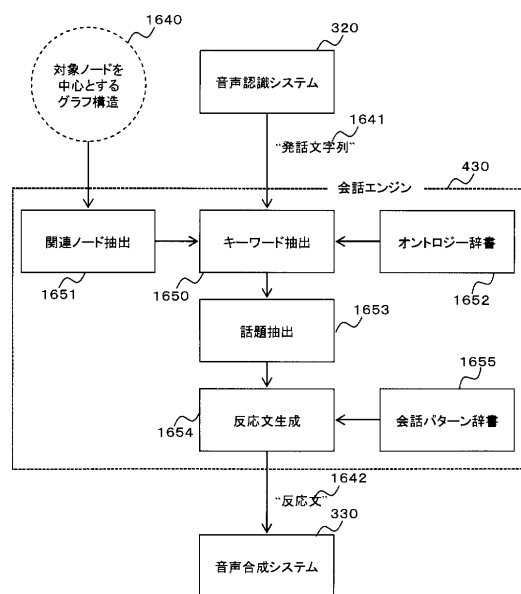
【 図 2 1 】



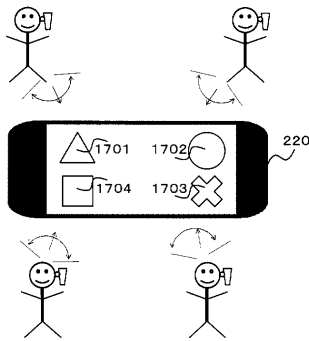
【 図 2 3 A 】



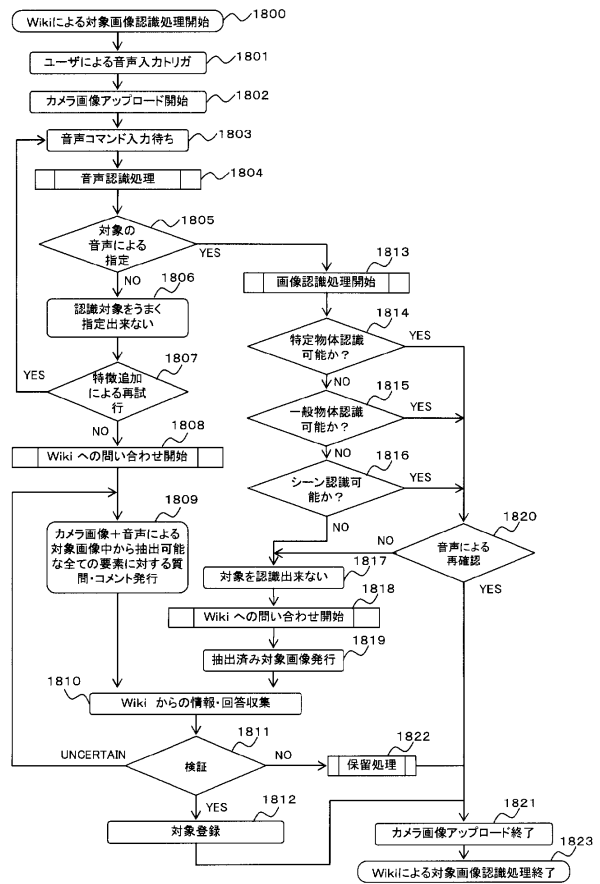
【 図 2 3 B 】



【図 24】

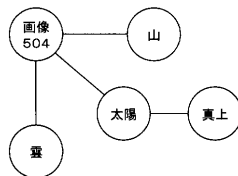


【図 25】

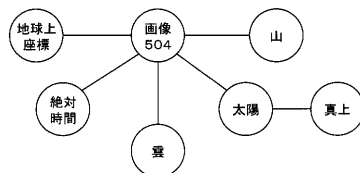


【図 26】

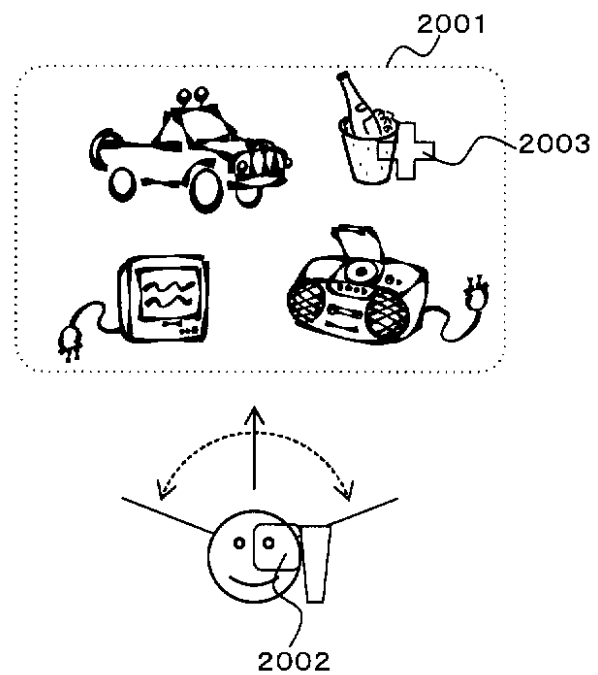
(A)



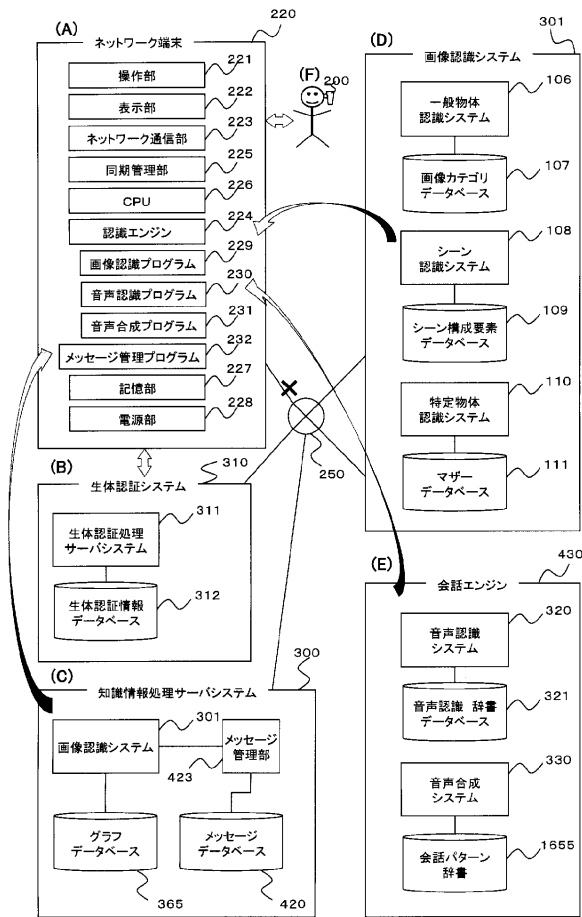
(B)



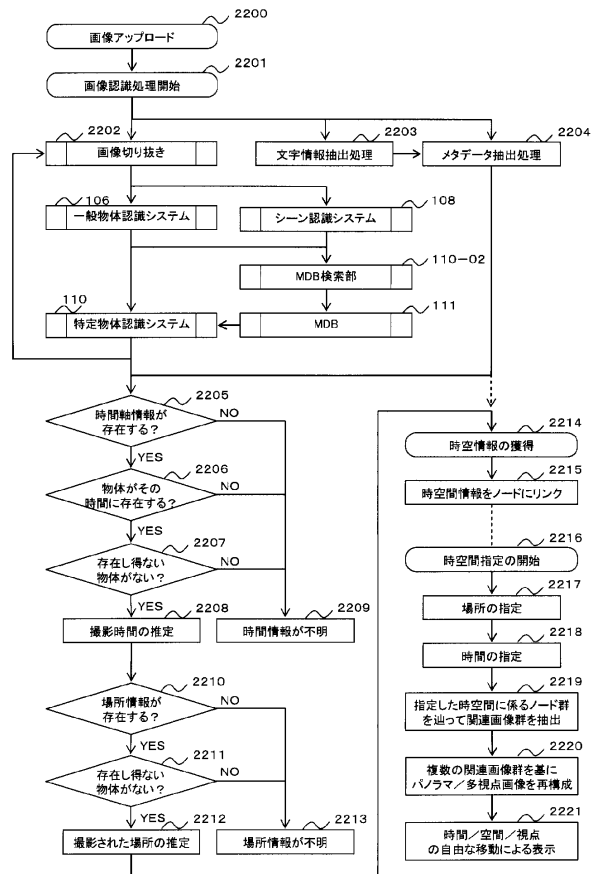
【図 27】



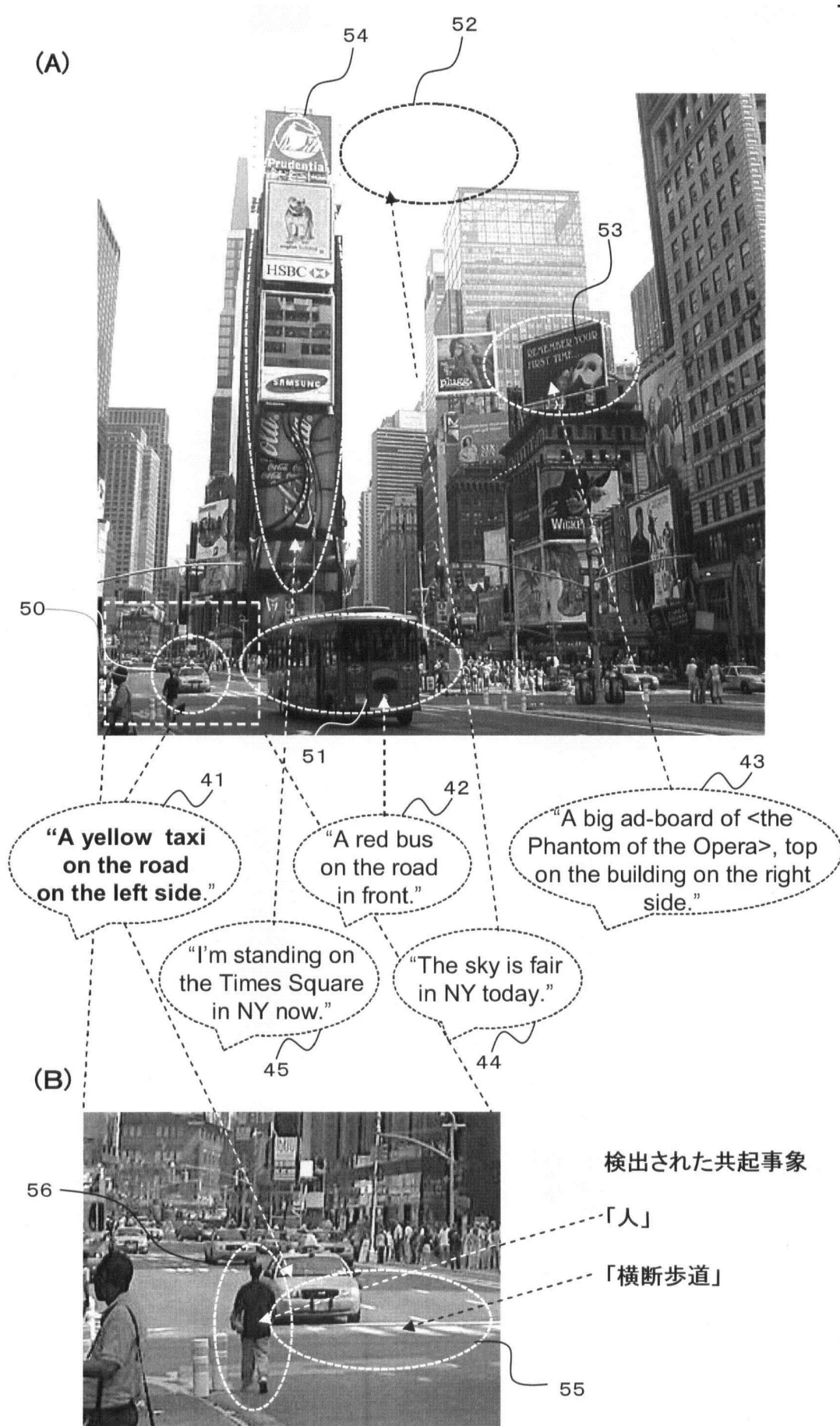
【図 28】



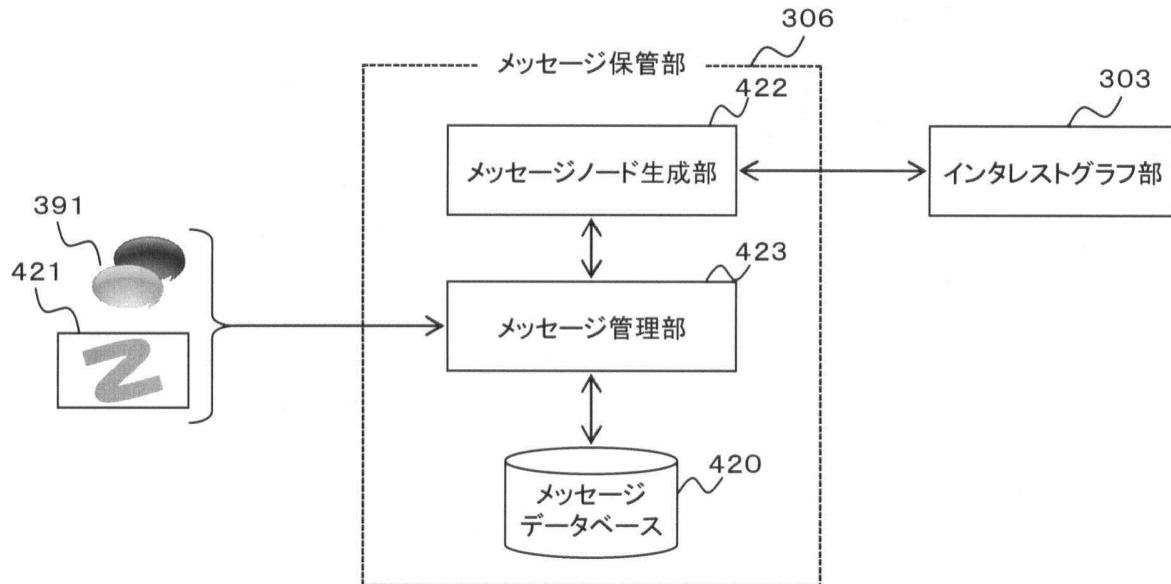
【図 30】



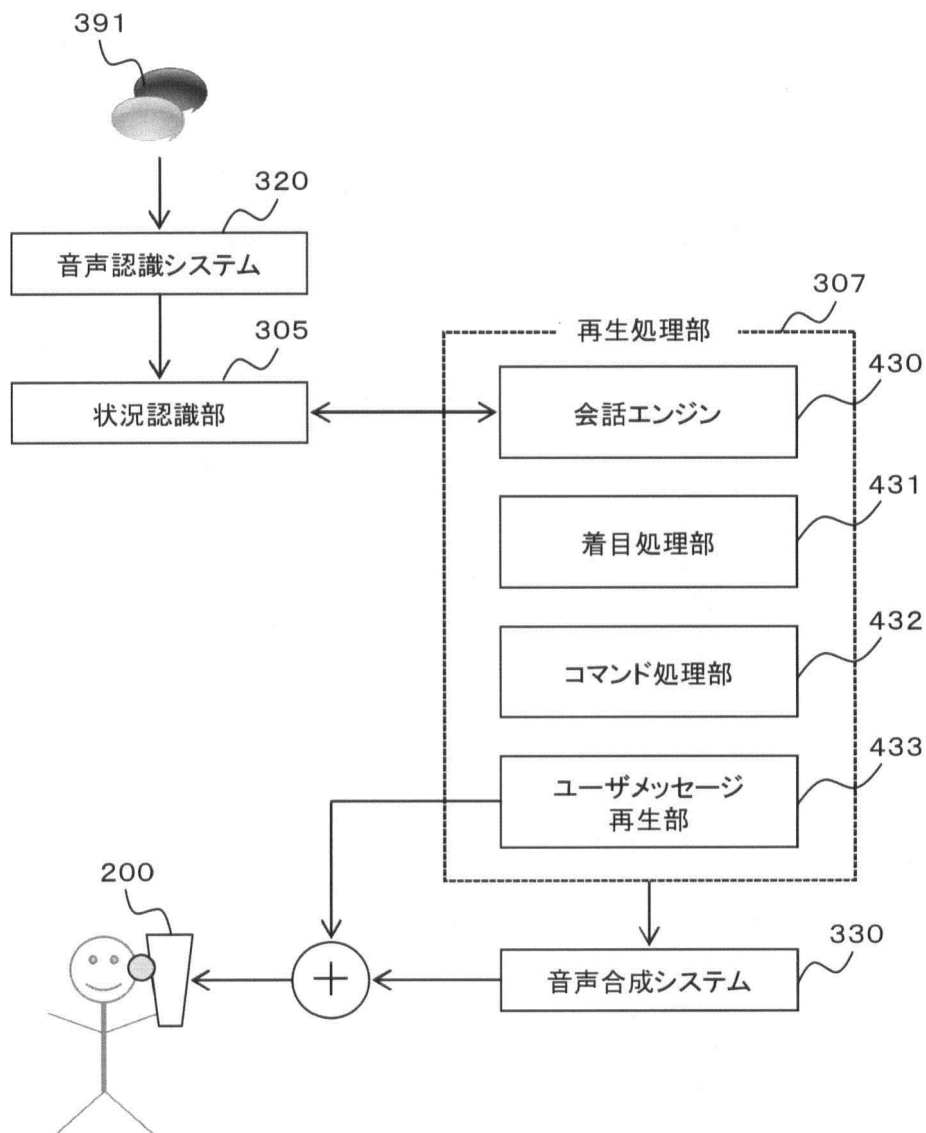
【図4A】



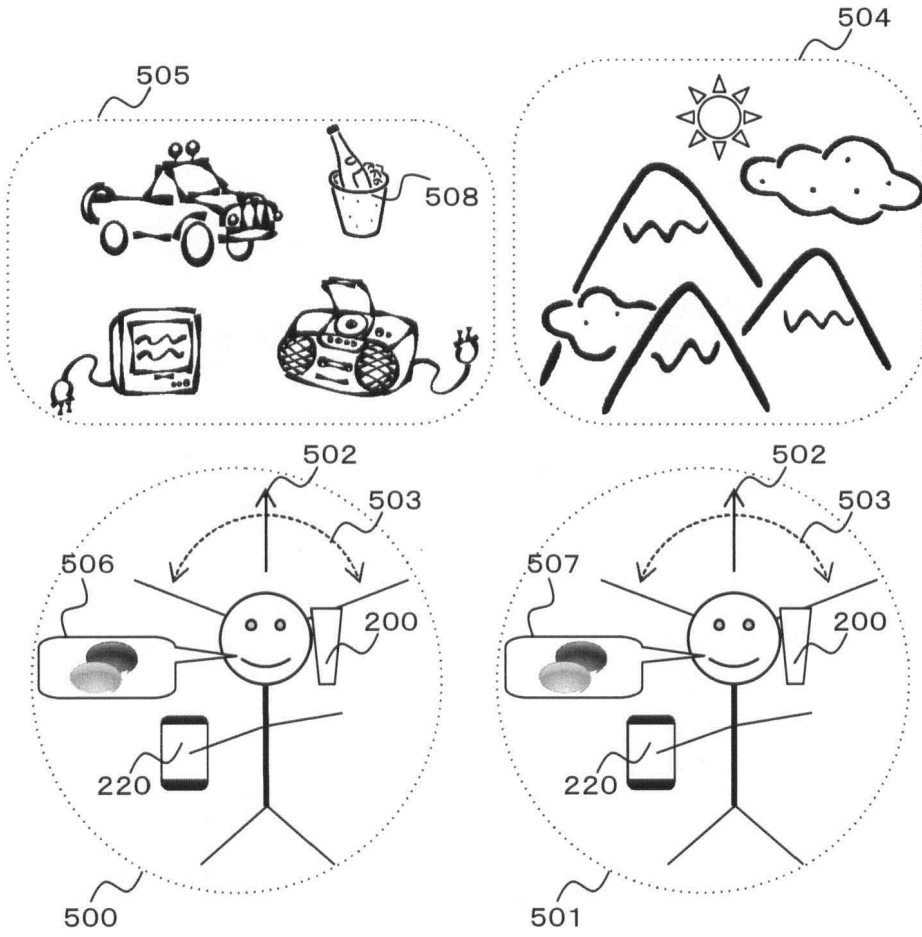
【図10】



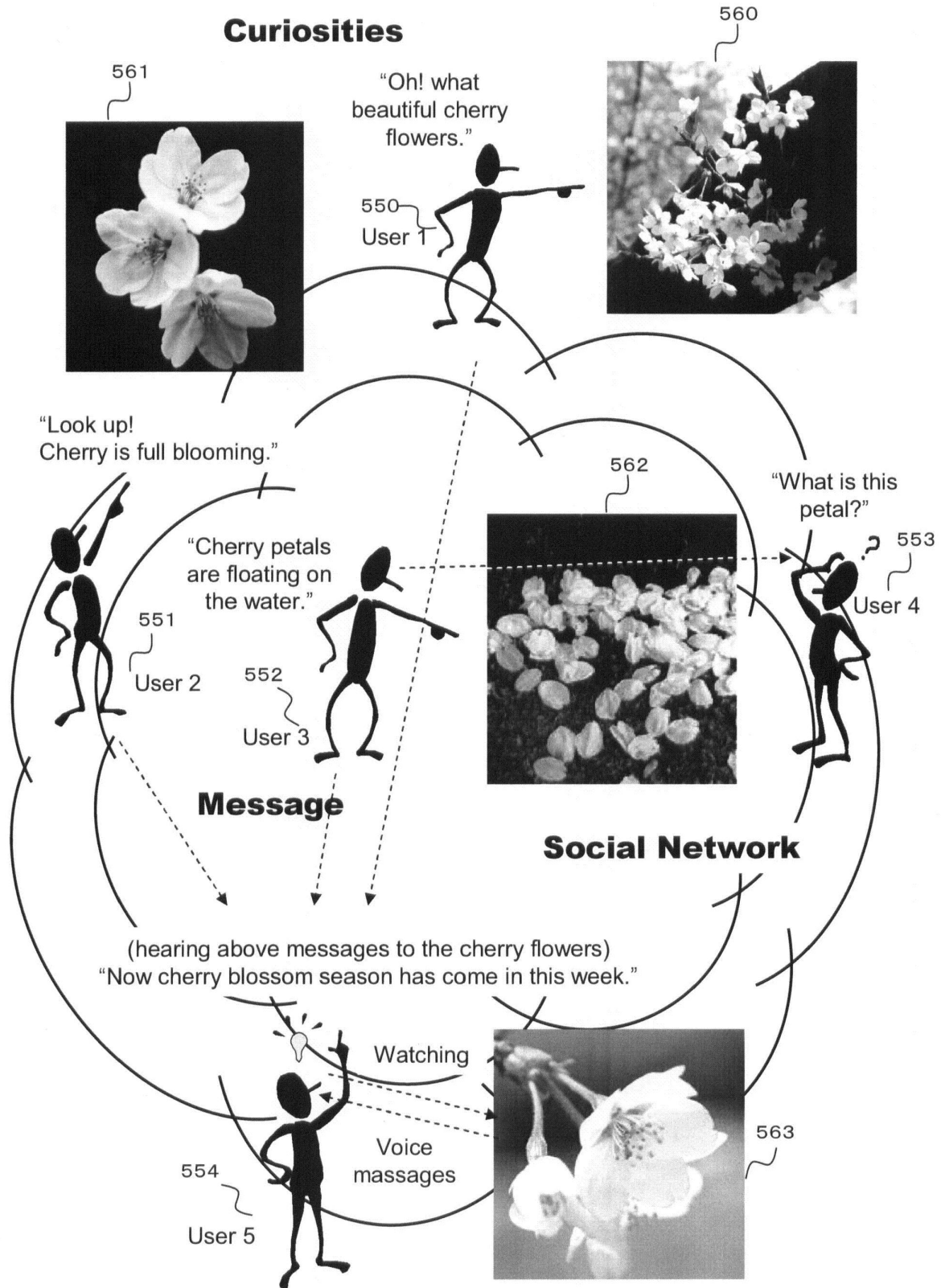
【図11】



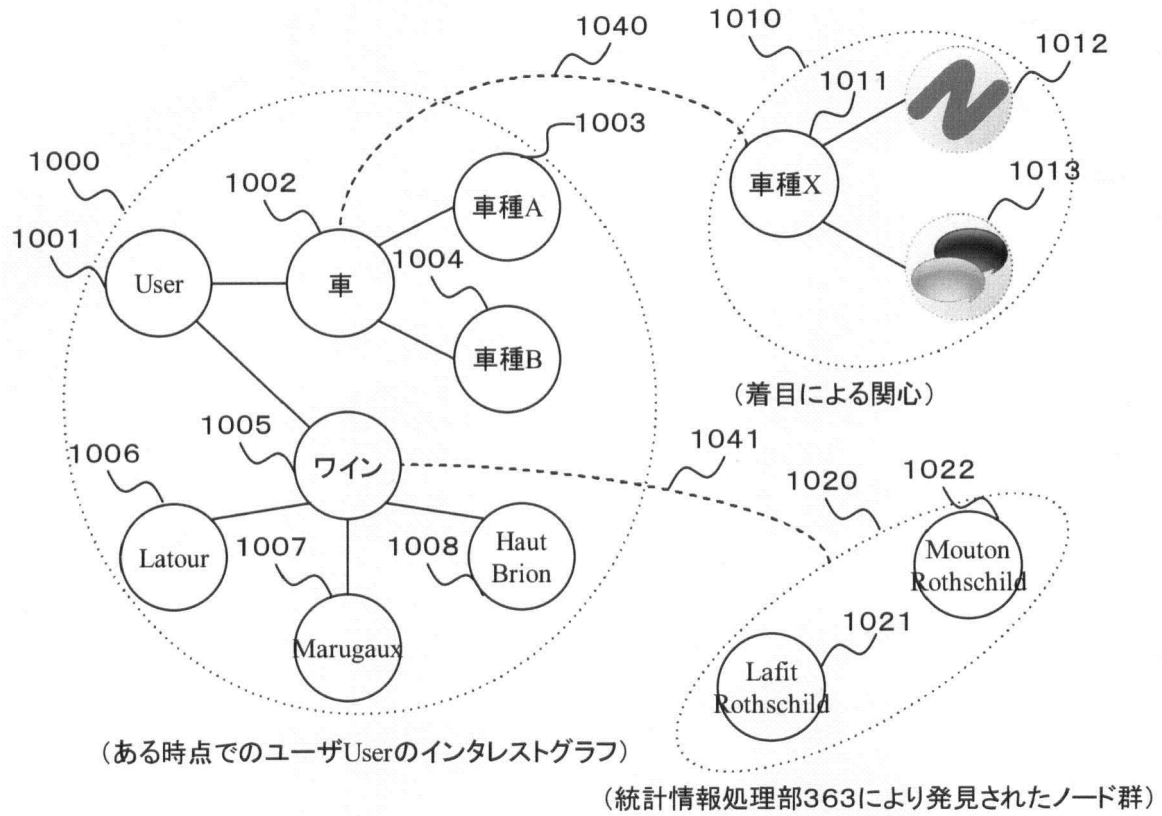
【図 13 A】



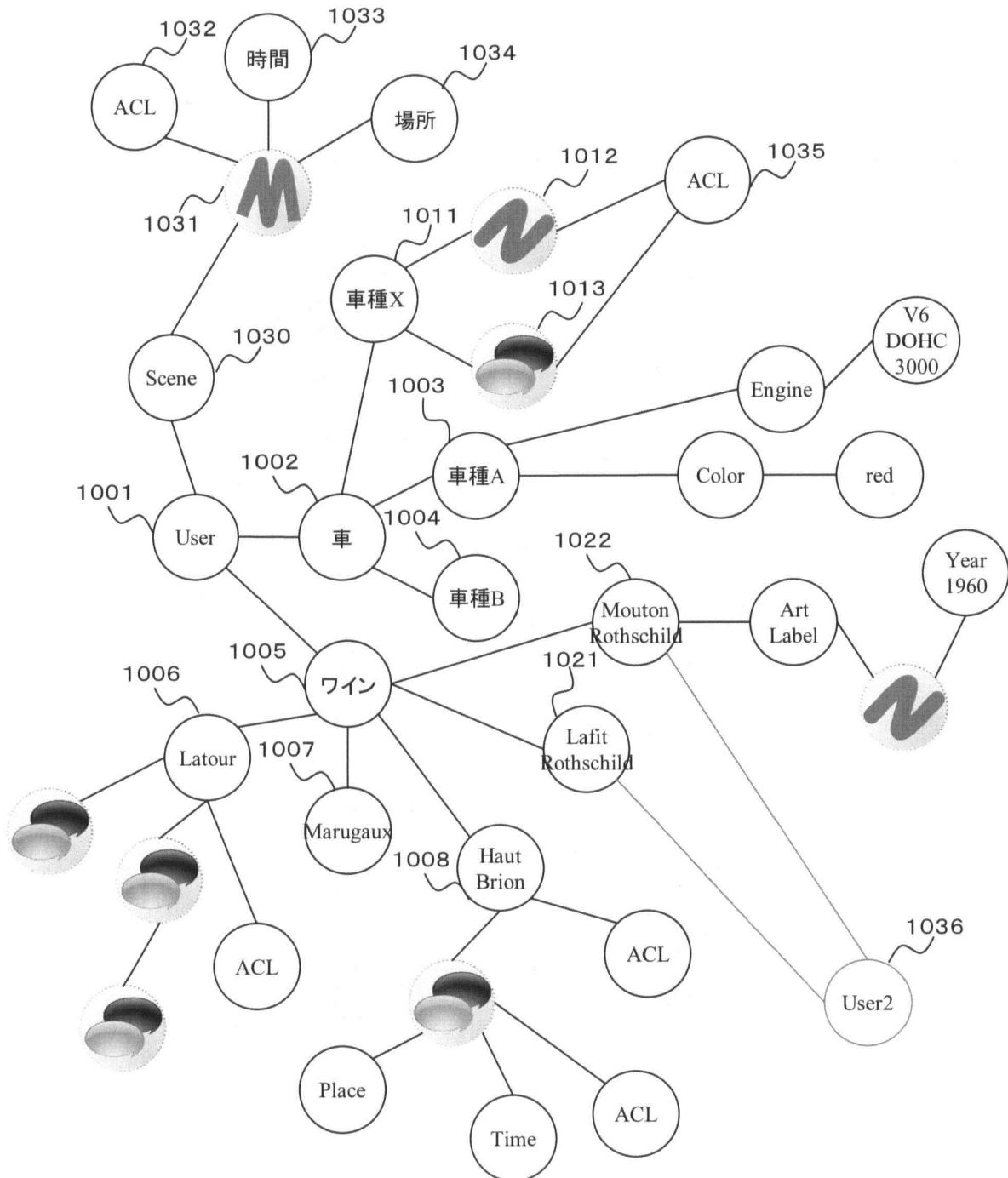
【図 13 B】



【図16】

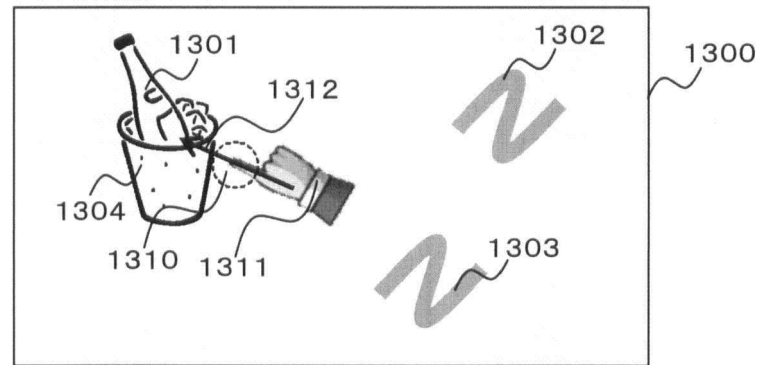


【図 17】

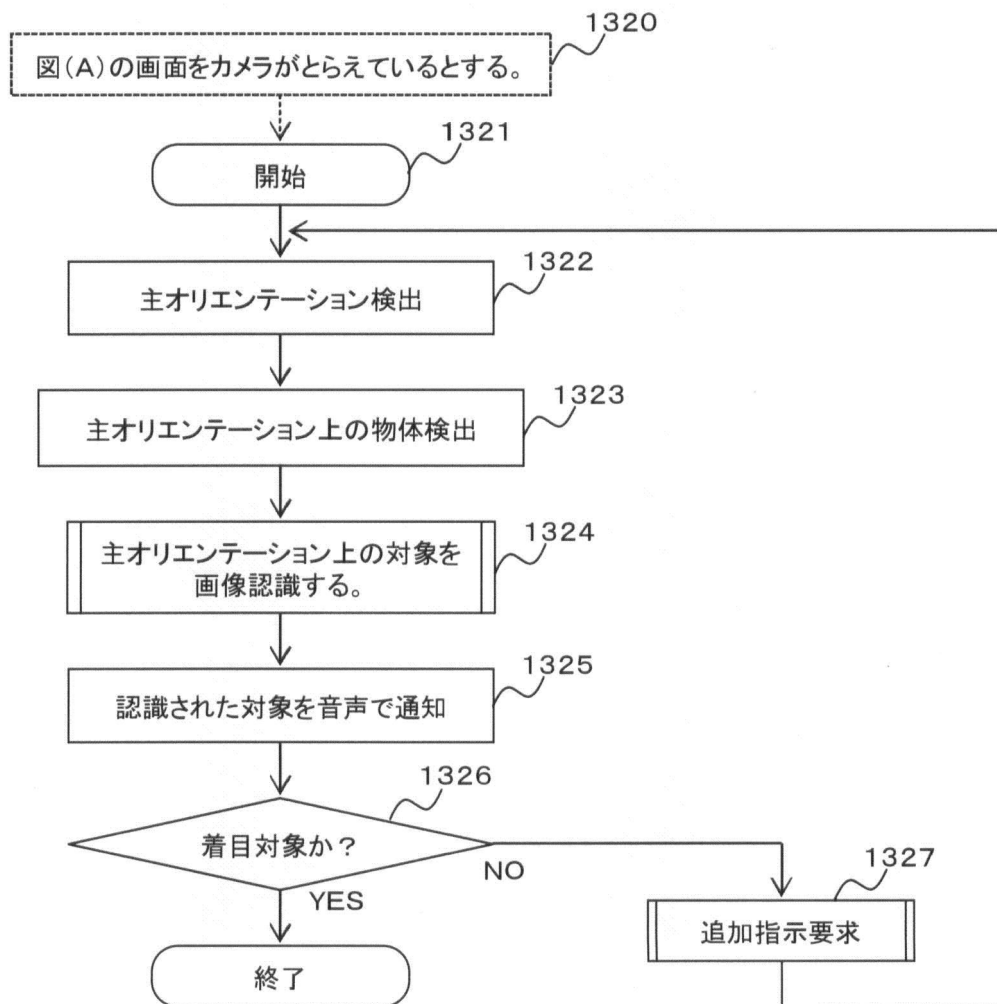


【図20】

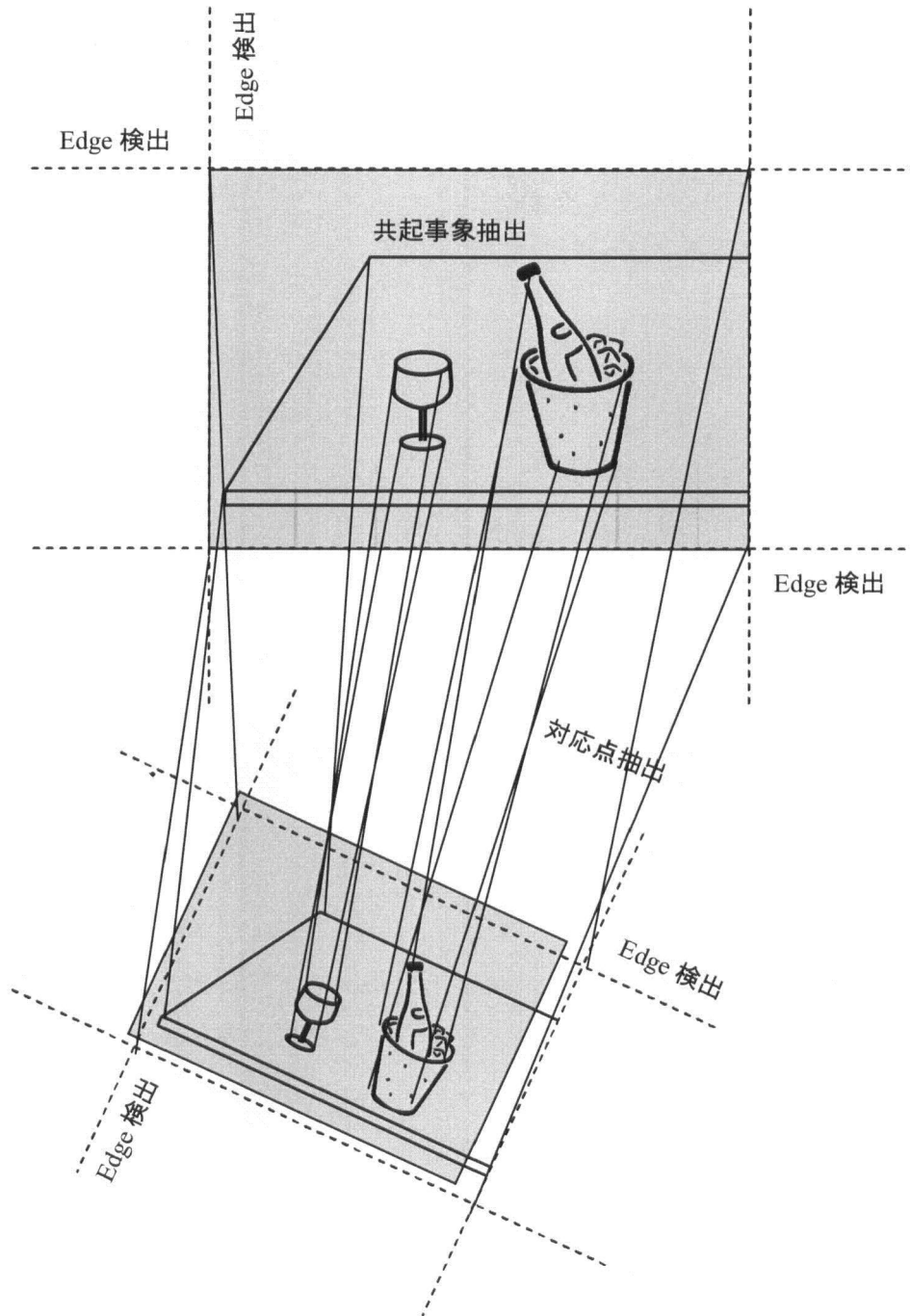
(A)



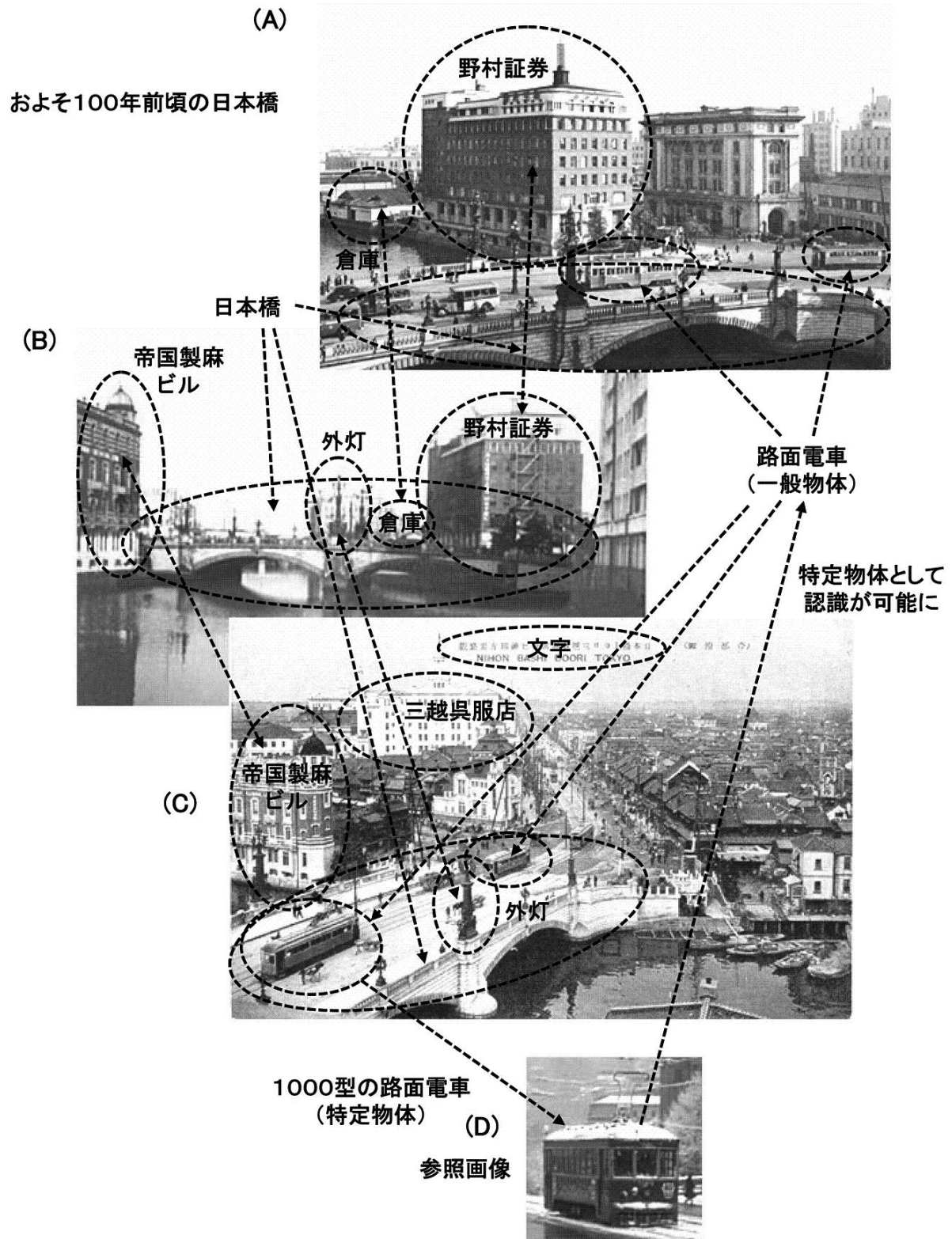
(B)



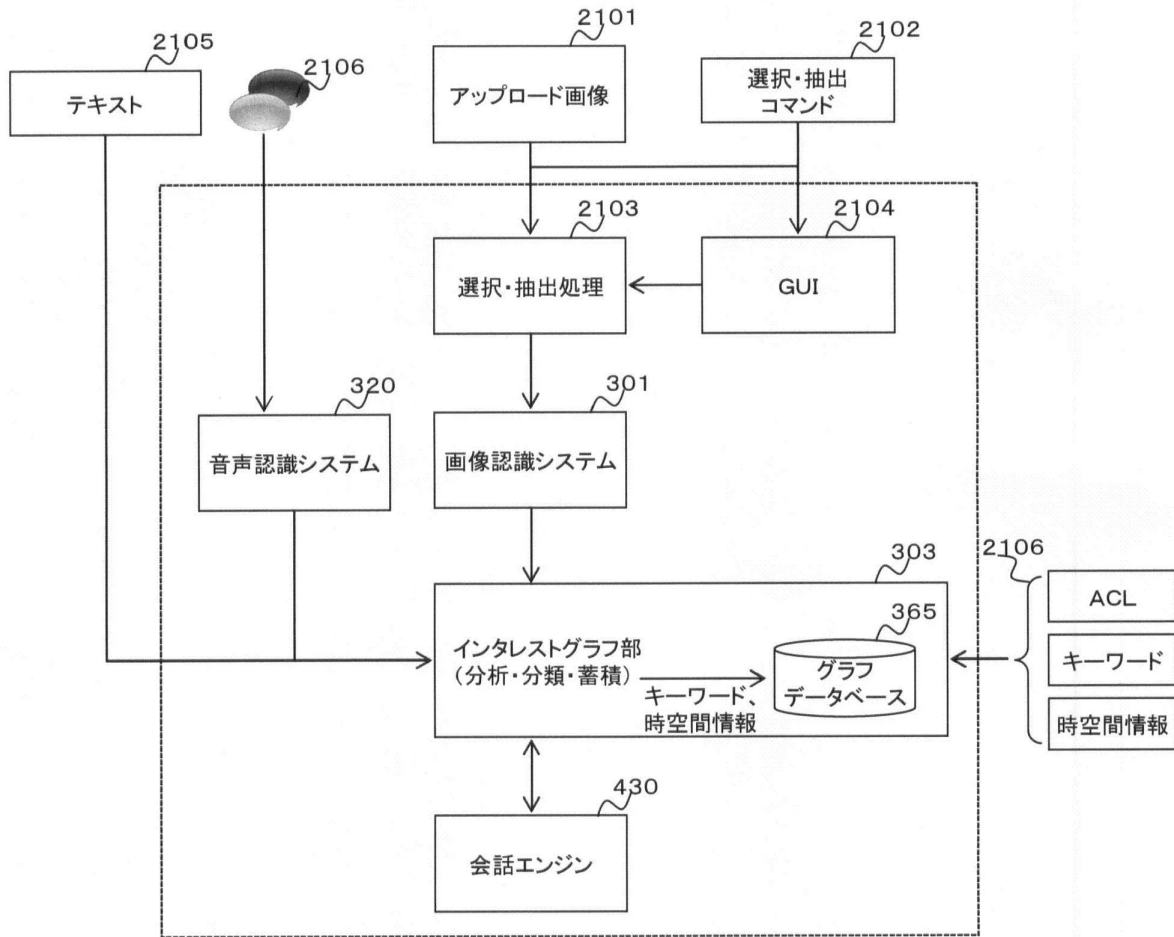
【図 22】



【図 29】



【図 3 1】



フロントページの続き

(51)Int.Cl.		F I
G 1 0 L 13/00 (2006.01)	G 1 0 L 15/28	5 0 0
	G 1 0 L 15/00	2 0 0 A
	G 1 0 L 13/00	

(74)代理人 100109335

弁理士 上杉 浩

(74)代理人 100136744

弁理士 中村 佳正

(72)発明者 久尋良木 健

東京都世田谷区瀬田 1 - 5 - 1 サイバーアイ・エンタテインメント株式会社内

(72)発明者 薄 隆

東京都世田谷区瀬田 1 - 5 - 1 サイバーアイ・エンタテインメント株式会社内

(72)発明者 横手 靖彦

東京都世田谷区瀬田 1 - 5 - 1 サイバーアイ・エンタテインメント株式会社内

審査官 木村 雅也

(56)参考文献 国際公開第 2 0 1 1 / 0 8 1 1 9 4 (W O , A 1)

特開 2 0 1 1 - 1 3 7 6 3 8 (J P , A)

特開 2 0 0 5 - 1 9 6 4 8 1 (J P , A)

特開 2 0 0 9 - 0 7 7 4 4 3 (J P , A)

特開 2 0 0 8 - 2 7 8 0 8 8 (J P , A)

国際公開第 2 0 1 1 / 0 0 4 6 0 8 (W O , A 1)

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 3 / 0 0

G 0 6 F 1 7 / 3 0

G 0 6 Q 5 0 / 0 0

G 1 0 L 1 3 / 0 0

G 1 0 L 1 5 / 0 0

G 1 0 L 1 5 / 2 8