

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5027400号  
(P5027400)

(45) 発行日 平成24年9月19日 (2012.9.19)

(24) 登録日 平成24年6月29日 (2012.6.29)

(51) Int.Cl.

F I

H O 4 N 5/76 (2006.01)  
G O 6 T 7/00 (2006.01)H O 4 N 5/76 B  
G O 6 T 7/00 3 O O F

請求項の数 16 (全 16 頁)

(21) 出願番号 特願2005-305799 (P2005-305799)  
 (22) 出願日 平成17年10月20日 (2005.10.20)  
 (65) 公開番号 特開2006-129480 (P2006-129480A)  
 (43) 公開日 平成18年5月18日 (2006.5.18)  
 審査請求日 平成20年10月20日 (2008.10.20)  
 (31) 優先権主張番号 10/978, 172  
 (32) 優先日 平成16年10月30日 (2004.10.30)  
 (33) 優先権主張国 米国 (US)

(73) 特許権者 500046438  
 マイクロソフト コーポレーション  
 アメリカ合衆国 ワシントン州 9805  
 2-6399 レッドモンド ワン マイ  
 クロソフト ウェイ  
 (74) 代理人 100140109  
 弁理士 小野 新次郎  
 (74) 代理人 100089705  
 弁理士 社本 一夫  
 (74) 代理人 100075270  
 弁理士 小林 泰  
 (74) 代理人 100080137  
 弁理士 千葉 昭男  
 (74) 代理人 100096013  
 弁理士 富田 博行

最終頁に続く

(54) 【発明の名称】 録画された会議のタイムラインに使用するための自動顔領域抽出

(57) 【特許請求の範囲】

【請求項 1】

ビデオ・サンプル中の二人以上の顔画像を検出するステップと、  
 前記ビデオ・サンプルに対応するオーディオ・サンプル中の二人以上の話し手の位置を検出するステップと、  
 前記二人以上の話し手のうちの主要な話し手を検出するステップと、  
 前記オーディオ・サンプルを、時間的持続性及び話し手の位置によりクラスタ化して話し手タイムラインを作成し、タイムライン・データベースに格納するステップと、  
 話し手位置毎に少なくとも1つの顔画像を顔データベースに格納するステップと、  
 前記顔データベースに格納された前記少なくとも1つの顔画像の中から、話し手位置毎に1つの顔画像を選択するステップと、  
 二人以上の話し手をその中に含む前記オーディオ／ビデオ（A／V）サンプルを表示するステップと、  
 前記主要な話し手の顔画像、及び、前記話し手タイムラインに隣接して表示された前記選択された顔画像を別個に表示するステップと、  
 を備えることを特徴とする方法。

【請求項 2】

前記二人以上の顔画像を検出するステップは、フェイス・トラッキングを使用して前記二人以上の顔画像を検出するステップをさらに備えることを特徴とする請求項 1 に記載の方法。

10

20

## 【請求項 3】

前記二人以上の話し手の位置を検出するステップは、音源位置測定を使用して前記二人以上の話し手の位置を検出するステップをさらに備えることを特徴とする請求項 1 に記載の方法。

## 【請求項 4】

前記顔データベースに格納された前記少なくとも 1 つの顔画像の中から、話し手位置毎に 1 つの顔画像を選択するステップは、最良の顔画像を選択するステップを含むことを特徴とする請求項 1 に記載の方法。

## 【請求項 5】

前記最良の顔画像を選択するステップは、顔の最も正面からの撮影像を含む顔画像を前記最良の顔画像として選択するステップを含むことを特徴とする請求項 4 に記載の方法。

10

## 【請求項 6】

前記最良の顔画像を選択するステップは、最少の動きを示す顔画像を前記最良の顔画像として選択するステップを含むことを特徴とする請求項 4 に記載の方法。

## 【請求項 7】

前記最良の顔画像を選択するステップは、最大の対称性を示す顔画像を前記最良の顔画像として選択するステップを含むことを特徴とする請求項 4 に記載の方法。

## 【請求項 8】

前記話し手の位置は、ビデオ・サンプルの座標で識別される話し手バウンディング・ボックスによって表されることを特徴とする請求項 1 に記載の方法。

20

## 【請求項 9】

前記話し手の位置は、ビデオ・サンプルにおいて方位角および仰角によって識別される話し手の顔の角度によって示されることを特徴とする請求項 1 に記載の方法。

## 【請求項 10】

前記主要な話し手は、前記タイムラインに基づいて判断されることを特徴とする請求項 1 に記載の方法。

## 【請求項 11】

コンピュータに、請求項 1 ~ 10 の何れか 1 項に記載されたステップを実行させるためのプログラム。

## 【請求項 12】

30

コンピュータに、請求項 1 ~ 10 の何れか 1 項に記載されたステップを実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

## 【請求項 13】

ビデオ・サンプル中の二人以上の顔画像を検出する手段と、  
前記ビデオ・サンプルに対応するオーディオ・サンプル中の二人以上の話し手の位置を検出する手段と、

前記オーディオ・サンプルを、時間的持続及び話し手の位置によりクラスタ化して話し手タイムラインを作成し、タイムライン・データベースに格納する手段と、

話し手位置毎に少なくとも 1 つの顔画像を顔データベースに格納する手段と、

前記顔データベースに格納された前記少なくとも 1 つの顔画像の中から、話し手位置毎に 1 つの顔画像を選択する手段と、

40

前記タイムライン・データベースに格納されたタイムラインに基づいて、前記二人以上の話し手のうちの主要な話し手を検出する手段と、

二人以上の話し手をその中に含む前記オーディオ/ビデオ (A/V) サンプルを表示する手段と、

前記主要な話し手の顔画像、及び、前記話し手タイムラインに隣接して表示された前記選択された顔画像を別個に表示する手段と、

を備えることを特徴とするシステム。

## 【請求項 14】

前記オーディオ・サンプル中の二人以上の話し手の位置を検出する手段は、音源ローカ

50

ライザを備えることを特徴とする請求項 1 3 に記載のシステム。

【請求項 1 5】

前記ビデオ・サンプル中の二人以上の顔画像を検出する手段は、フェイス・トラッカを備えることを特徴とする請求項 1 3 に記載のシステム。

【請求項 1 6】

前記主要な話し手は、前記タイムラインに基づいて判断されることを特徴とする請求項 1 3 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

10

(関連出願の相互参照)

本出願は、2002年6月21日に出願し、本出願の譲受人であるMicrosoft Corp. に譲渡された本発明者による米国特許出願第10/177315号「A System and Method for Distributed Meetings」の一部継続出願である。本出願人は、参照によりその開示、教示する内容全てが本明細書に組み込まれる前記出願の出願日について優先権を主張するものである。

【0002】

以下の説明は、一般にはビデオ画像処理に関する。より詳細には、以下の説明はビデオ再生に用いるインデックス付タイムラインを提供することに関する。

20

【背景技術】

【0003】

2人以上の話し手(speaker)を含む状況を録画したビデオの再生、例えば、録画された会議の再生は、通常インデックス付のタイムライン(時間記録)と同時に表示される。タイムラインを使用することで、ユーザは1つまたは複数のタイムライン制御手段を操作して会議中の特定の時点に素早く移動することができる。ビデオが2人以上の話し手が含む場合、それぞれが特定の話し手に関連する複数のタイムラインを使うことができる。各タイムラインは、対応する話し手がいつ話しているかを示す。それによって、ユーザは、会議中特定の話し手が話している部分に進むことができる。

【0004】

このような複数のタイムラインは、例えば「話し手1」、「話し手2」等のように、一般的なやり方でラベルを付けて各話し手を識別することができる。タイムラインに特定の話し手の名前を自動的に付与する現行の技術は不正確であり、ユーザおよびユーザに関連する声紋(voiceprints)および顔紋(faceprints)のデータベースが必要となる可能性もあり、これらはセキュリティおよびプライバシーの問題を伴う可能性がある

30

【特許文献1】米国特許出願公開第2003/0234866号明細書「A System and Method for Distributed Meetings」

【発明の開示】

【発明が解決しようとする課題】

【0005】

本発明の目的は、ビデオ画像処理におけるビデオ再生に用いるインデックス付タイムラインを提供することにある。

40

【課題を解決するための手段】

【0006】

上記目的を達成するために、本発明によれば、(1)ビデオ・サンプル中の1つまたは複数の顔画像を検出すること、前記ビデオ・サンプルに対応するオーディオ・サンプル中の1人または複数の話し手を検出すること、話し手識別子によって話し手を識別する話し手タイムラインと、前記話し手タイムラインの各時点における話し手の位置とを格納すること、顔データベースに検出された各話し手毎に少なくとも1つの顔画像を格納すること、および検出された各話し手に話し手タイムラインと顔画像とを関連付けることを備えることを特徴とする方法、(2)1人または複数の話し手をその中に含むオーディオ/ビジ

50

ュアル（Ａ／Ｖ）サンプルを表示すること、各話し手に対応する話し手タイムラインであって、時間的な経過のどの時点で前記話し手タイムラインに対応する前記話し手が話しているかを示す、話し手タイムラインを表示すること、各話し手タイムラインに、前記話し手タイムラインに関連付けられた前記話し手に対応する前記話し手の顔画像を関連付けること、および前記顔画像を前記対応する話し手タイムラインと共に表示することを備えることを特徴とする方法、（３）実行可能な命令を含む１つまたは複数のコンピュータ可読媒体であって、前記命令は実行されたとき、Ａ／Ｖサンプル中の各話し手を話し手識別子で識別すること、前記Ａ／Ｖサンプル中の各話し手の位置を識別すること、前記Ａ／Ｖサンプル中で識別された各話し手毎に少なくとも１つの顔画像を抽出すること、前記Ａ／Ｖサンプル中で識別された各話し手毎に、それぞれが時間、話し手識別子および話し手の位置を示す話し手タイムラインを作成すること、および話し手の前記顔画像を、同じ話し手に対応する話し手タイムラインと関連付けることを備える方法を実施することを特徴とする１つまたは複数のコンピュータ可読媒体、（４）１つまたは複数のコンピュータ可読媒体であって、Ａ／Ｖサンプルの各話し手毎の話し手タイムラインであって、それぞれが時間的経過の複数の時点における話し手と話し手の位置とを識別する話し手タイムラインを含む話し手タイムライン・データベースと、話し手タイムラインで識別された各話し手毎に少なくとも１つの顔画像と、各顔画像を前記話し手タイムライン・データベース中の適切な話し手タイムラインに関連付ける話し手識別子とを含む顔データベースとを備えることを特徴とする１つまたは複数のコンピュータ可読媒体、ならびに（５）Ａ／Ｖサンプルと、前記Ａ／Ｖサンプルに登場する各話し手を識別する手段と、前記Ａ／Ｖサンプルで識別された各話し手毎に顔画像を識別する手段と、前記Ａ／Ｖサンプルで識別された各話し手毎に話し手タイムラインを作成する手段と、顔画像を適切な話し手タイムラインに関連付ける手段とを備えることを特徴とするシステム、が提供される。

#### 【０００７】

前述の諸態様および本発明のその他の利点の多くは、以下の詳細な説明を添付の図面と併せて参照することにより、より理解されよう。

#### 【発明を実施するための最良の形態】

#### 【０００８】

以下の説明は、複数の話し手が存在する環境において各話し手の顔を自動的に検出し、その話し手の１つまたは複数の顔画像を、その話し手に対応するタイムラインの部分に関連付けるための様々な実装および実施形態に関する。この種の特殊ラベル付けは、タイムラインのどの部分が複数の話し手のうち特定の１人に対応しているかを視聴者がより容易に判断できる点で、汎用的なラベル付けと比べて優れている。

#### 【０００９】

以下の議論では、２人以上の参加者および／または話し手が登場する会議の録画にパノラマ式カメラを使用する例について説明する。複数のカメラを含むパノラマ式カメラについて説明されているが、以下の説明は単一のカメラにも、および２台以上のカメラを有するマルチカメラ装置にも関連する。

#### 【００１０】

パノラマ画像が、会議内の顔を検出し追跡するフェイス・トラッカ（ＦＴ）に入力される。マイクロホン・アレイが、音に基づいて話し手の位置を検出する音源ローカライザ（ＳＳＬ）に入力される。フェイス・トラッカおよび音源ローカライザからの出力は、仮想シネマトグラファ（virtual cinematographer）に入力され、複数の話し手の位置を検出する。

#### 【００１１】

話し手は、話し手クラスタリング・モジュールで後処理され、これは、２つ以上の個別タイムラインを含む集約タイムラインをより明確に区別するために、話し手を時間的かつ空間的にクラスタ化する。（集約）タイムラインは、タイムライン・データベースに格納される。各話し手毎に１つまたは複数の画像を格納する顔データベースが作成されるが、それぞれの顔のうち少なくとも１つは話し手に関連付けられたタイムラインに使用される

。

## 【 0 0 1 2 】

本明細書で提示され特許請求される概念を、1つまたは複数の適切な動作環境に関連して以下に詳しく説明する。以下に説明する要素のうち幾つかは、「A System and Method for Distributed Meetings」の表題を有する、2002年6月21日に出願された、本出願の親出願特許文献1)にも記載されている。

## 【 0 0 1 3 】

## 例示的動作環境

図1は、汎用コンピュータ/カメラ・デバイスを示すブロック図である。コンピューティング・システム環境100は、適切なコンピューティング環境の一例に過ぎず、特許請求の主題の範囲または機能に関するいかなる限定を示唆するものでもない。また、コンピューティング・システム環境100も、例示的動作環境100に示される構成要素の1つまたはそれらの組合せに依存するか、あるいはそれらを必要としていると解釈されるべきではない。

## 【 0 0 1 4 】

説明される技術および対象物は、他の幾多の汎用または専用のコンピューティング・システム環境または構成で動作させることができる。使用に適した周知のコンピューティング・システム、環境および/または構成としては、パーソナル・コンピュータ、サーバ・コンピュータ、携帯型またはラップトップ型装置、マルチ・プロセッサ・システム、マイクロ・プロセッサに基づくシステム、セットトップボックス、プログラム可能な家庭用電化製品、ネットワークPC、ミニ・コンピュータ、メインフレーム・コンピュータ、上記のシステムまたは装置のうちいずれかを含む分散コンピューティング環境などがあるが、ただしこれに限定されない。

## 【 0 0 1 5 】

以下の説明は、コンピュータで実行される、例えばプログラム・モジュールなどのコンピュータ実行命令という一般的なコンテキストで表現することができる。一般的に、プログラム・モジュールとしては、特定のタスクを実行しまたは特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、構成要素、データ構造等がある。説明する実装形態は、通信ネットワークを介してリンクされたりリモート処理デバイスによってタスクが実行される分散コンピューティング環境で実施することもできる。分散コンピューティング環境では、プログラム・モジュールは、メモリ記憶装置を含むローカルまたはリモートのコンピュータ記憶媒体に置くことができる。

## 【 0 0 1 6 】

図1を参照すると、本発明を実施するための例示的システムは、コンピュータ110の形をとる汎用コンピューティング・デバイスを含む。コンピュータ110の構成要素は、処理装置120、システム・メモリ130、およびシステム・メモリを含めたシステムの様々な構成要素を処理装置120に連結するシステム・バス121を含むが、それだけに限定されない。システム・バス121は、メモリ・バスまたはメモリ・コントローラ、周辺バス、および様々なバス・アーキテクチャのいずれかを使用したローカル・バスを含めて、様々なバス構造を有するものでよい。限定ではなく例として挙げると、こうしたアーキテクチャとしては、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) ローカル・バス、およびメザニン・バス (Mezzanine bus) と呼ばれるPCI (Peripheral Component Interconnect) バスがある。

## 【 0 0 1 7 】

コンピュータ110は、一般には様々なコンピュータ可読媒体を含む。コンピュータ可読媒体は、コンピュータ110がアクセスできる利用可能などんな媒体でもよく、揮発性および不揮発性の媒体、着脱式または非着脱式の媒体のいずれも含まれる。限定ではなく例として挙げると、コンピュータ可読媒体は、コンピュータ記憶媒体および通信媒体を備えることができる。コンピュータ記憶媒体としては、コンピュータ可読命令、データ構造

10

20

30

40

50

、プログラム・モジュールまたは他のデータ等の情報の記憶に用いられる任意の方式または技術で実装される揮発性および不揮発性の媒体、着脱式または非着脱式の媒体が含まれる。コンピュータ記憶媒体としては、RAM、ROM、EEPROM、フラッシュ・メモリまたは他のメモリ技術、CD-ROM、DVD (digital versatile disks) または他の光学式ディスク記憶、磁気カセット、磁気テープ、磁気ディスク記憶または他の磁気記憶デバイス、あるいは所望の情報を格納するために使用され、コンピュータ110がアクセスできるような他のいかなる媒体なども含まれる、それだけには限定されない。通信媒体は、搬送波もしくは他の搬送機構等の変調データ信号の形で、一般に、コンピュータ可読命令、データ構造、プログラム・モジュール、または他のデータを具現化し、任意の情報送達媒体がこれに含まれる。「変調データ信号」という言葉は、その1つまたは複数の特性が、信号に含まれる情報を符号化するように設定または変更された信号を意味する。限定ではなく例として挙げると、通信媒体としては、有線ネットワークまたは直接配線接続等の有線媒体、および音響、RF、赤外線およびその他の無線媒体がある。上記の任意の組合せも、コンピュータ可読媒体の範囲に含まれるべきである。

#### 【0018】

システム・メモリ130は、例えば読取り専用メモリ (ROM) 131やランダム・アクセス・メモリ (RAM) 132など揮発性および/または不揮発性のメモリの形をとるコンピュータ記憶媒体を含む。ROM 131は、一般に、例えば起動の際にコンピュータ110内の要素間での情報の転送を助ける基本ルーチンを含む基本入出力システム (BIOS) 133を格納している。RAM 132は、一般に、ダイレクトに処理装置120からアクセス可能であり、かつ/または処理装置120が操作を加えているデータおよび/またはプログラム・モジュールを格納している。限定ではなく例として挙げると、図1はオペレーティング・システム134、アプリケーション・プログラム135、他のプログラム・モジュール136、およびプログラム・データ137を示す。

#### 【0019】

コンピュータ110は、他の着脱式/非着脱式、揮発性/不揮発性のコンピュータ記憶媒体を含むこともできる。ほんの一例として、図1は非着脱式で不揮発性の磁気媒体を読み書きするハードディスク・ドライブ141、着脱式で不揮発性の磁気ディスク152を読み書きする磁気ディスク・ドライブ151、およびCD-ROMや他の光学式媒体など着脱式で不揮発性の光ディスク156を読み書きする光ディスク・ドライブ155を示す。この例示的動作環境で使用できる他の着脱式/非着脱式、揮発性/不揮発性のコンピュータ記憶媒体としては、磁気テープ・カセット、フラッシュ・メモリ・カード、DVD、デジタル・ビデオ・テープ、ソリッド・ステートRAM (solid state RAM)、ソリッド・ステートROM (solid state ROM) などがあるが、それだけには限定されない。ハードディスク・ドライブ141は、一般にインターフェース140などの非着脱式メモリ・インターフェースを介してシステム・バス121に接続されており、磁気ディスク・ドライブ151および光ディスク・ドライブ155は、一般にインターフェース150などの着脱式メモリ・インターフェースを介してシステム・バス121に接続されている。

#### 【0020】

上に説明し、図1に示すドライブおよびそれらに関連するコンピュータ記憶媒体は、コンピュータ110用のコンピュータ可読命令、データ構造、プログラム・モジュールまたは他のデータを格納する。図1では、例えばハードディスク・ドライブ141は、オペレーティング・システム144、アプリケーション・プログラム145、他のプログラム・モジュール146、およびプログラム・データ147を記憶しているものとして示される。これらの構成要素は、オペレーティング・システム134、アプリケーション・プログラム135、他のプログラム・モジュール136、およびプログラム・データ137と同じものでも、別のものでもよいことに留意されたい。ここでは、オペレーティング・システム144、アプリケーション・プログラム145、他のプログラム・モジュール146、およびプログラム・データ147が少なくとも異なるコピーであることを示すために別の番号を付けた。ユーザは、キーボード162や、一般にマウス、トラックボールまたは

10

20

30

40

50

タッチパッドと称されるポインティング・デバイス 161 などの入力装置を介して、コンピュータ 110 に命令および情報を入力する。他の入力装置（図示せず）としては、マイクロホン、ジョイスティック、ゲーム・パッド、パラボラ・アンテナ、スキャナなどがある。これらおよび他の入力装置は、システム・バス 121 に結合されているユーザ入力インターフェース 160 を介して処理装置 120 に接続されることが多いが、パラレル・ポート、ゲーム・ポート、ユニバーサル・シリアル・バス（USB）などその他のインターフェースおよびバス構造によって接続することもできる。モニタ 191 または他の表示装置も、ビデオ・インターフェース 190 などのインターフェースを介してシステム・バス 121 に接続されている。モニタに加え、コンピュータは、出力周辺インターフェース 195 を介して接続することができるスピーカ 197 やプリンタ 196 など他の周辺出力装置を含むこともできる。本発明にとって特に重要なことに、パーソナル・コンピュータ 110 の入力装置として、一連の画像 164 を取り込むカメラ 163（例えばデジタル／電子のスチルまたはビデオ・カメラ、あるいはフィルム／写真式スキャナ）も含むことができる。さらに、カメラは 1 つのみ図示されているが、パーソナル・コンピュータ 110 の入力装置として複数のカメラを含むこともできる。1 つまたは複数のカメラの画像 164 は、適切なカメラ・インターフェース 165 を介してコンピュータ 110 に入力される。このインターフェース 165 はシステム・バス 121 に接続され、それによって、画像を RAM 132 またはコンピュータ 110 に関連する他のデータ記憶装置の 1 つに送り格納することが可能になる。しかし、画像データは、カメラ 163 の使用を必要とせず、前述のどんなコンピュータ可読媒体からもコンピュータ 110 に入力することができることに留意されたい。

#### 【0021】

コンピュータ 110 は、リモート・コンピュータ 180 など 1 つまたは複数のリモート・コンピュータへの論理接続を使用したネットワーク環境で動作することもできる。リモート・コンピュータ 180 は、パーソナル・コンピュータ、サーバ、ルータ、ネットワーク PC、ピア・デバイスまたは他の共通ネットワーク・ノードでよく、図 1 ではメモリ記憶装置 181 しか図示されていないが、一般にはコンピュータ 110 に関連して説明した上述の多くまたは全ての要素を含む。図 1 に描かれた論理接続は、ローカル・エリア・ネットワーク（LAN）171 およびワイド・エリア・ネットワーク（WAN）173 を含むが、他のネットワークを含むこともできる。これらのネットワーク環境は、オフィス、事業体規模のコンピュータ・ネットワーク、イントラネットおよびインターネットにおいて一般的となっている。

#### 【0022】

LAN ネットワーク環境で使用する場合、コンピュータ 110 はネットワーク・インターフェースまたはアダプタ 170 を介して LAN 171 に接続される。WAN ネットワーク環境で使用する場合、コンピュータ 110 は一般に、インターネット等の WAN 173 上での通信を確立するための、モデム 172 または他の手段を含む。内部、外部のどちらでもよいモデム 172 は、ユーザ入力インターフェース 160 または他の適切な機構を介してシステム・バス 121 に接続することができる。ネットワーク環境では、コンピュータ 110 に関連して描かれたプログラム・モジュールあるいはその一部は、リモート・メモリ記憶装置に格納することができる。限定ではなく例として挙げると、図 1 ではリモート・アプリケーション・プログラム 185 がメモリ装置 181 上にあるものとして示されている。図示されたネットワーク接続は例示的なものであり、コンピュータ間の通信リンクを確立する他の手段を使用することもできることが理解できよう。

#### 【0023】

例示的パノラマ式カメラおよびクライアント装置

図 2 は、例示的なパノラマ式カメラ機器 200 および例示的なクライアント装置 222 を示すブロック図である。特定の構成が示されているが、パノラマ式カメラ機器 200 は、パノラマ式カメラまたはその機能的等価品を含むどんな装置でもよいことに留意されたい。本明細書に記載される 1 つまたは複数の技術を組み込んだ実用的適用例では、この図

でパノラマ式カメラ機器 200 に含まれている構成要素よりも多くの、あるいは少ない構成要素を含むことができる。

【0024】

パノラマ式カメラ機器 200 は、プロセッサ 202 およびメモリ 204 を含む。パノラマ式カメラ機器 200 は、複数のカメラ 206 (206\_\_1 から 206\_\_n で示される) が撮影した画像の幾つかをつなぎ合わせることでパノラマ画像を作成する。パノラマ画像は、完全な 360 度のパノラマ画像でも、その一部のみでもよい。パノラマ式カメラ機器 200 が図示され説明されているが、ここで説明される技術は単一のカメラでも活用することに留意されたい。

【0025】

パノラマ式カメラ機器 200 はさらに、マイクロホン・アレイ 208 を含む。以下に詳細に説明するが、マイクロホン・アレイは、音の方向が位置測定されるように構成される。言い換えれば、マイクロホン・アレイに入力された音の分析によって、検出された音が発生した方向が得られる。スピーカフォン (speakerphone) を使用可能にし、あるいはユーザに通知信号などを発するため、スピーカ 210 をパノラマ式カメラ機器 200 に含めてもよい。

【0026】

メモリ 204 は、校正データ、露出設定、ステッチング表 (stitching table) 等の幾つかのカメラ設定 212 を格納する。メモリ 204 には、1 つまたは複数の他のカメラ・ソフトウェア・アプリケーション 216 と共に、カメラ機能を制御するオペレーティング・システム 214 も格納されている。

【0027】

パノラマ式カメラ機器 200 は、さらにパノラマ式カメラ機器 200 との間でのデータの送受信に用いる入出力 (I/O) モジュール 218、およびカメラの機能が必要とする可能性のある他の種々のハードウェア 220 要素も含む。

【0028】

パノラマ式カメラ機器 200 は、少なくとも 1 つのクライアント装置 222 と通信する。クライアント装置 222 は、プロセッサ 224、メモリ 226、大容量記憶装置 242 (ハードディスク・ドライブ等)、および図で下に示すクライアント装置 222 に帰せられる機能を実行するために必要となるかもしれない他のハードウェア 230 を含む。

【0029】

メモリ 226 は、フェイス・トラッカ (FT) モジュール 230 および音源位置測定 (SSL: sound source localization) モジュール 232 を格納している。フェイス・トラッカ・モジュール 230 および音源位置測定モジュール 232 は、仮想シネマトグラフィ 234 と共に使って、カメラ・シーン内の人物を検出し、その人物が話しているかどうか、話しているとすればいつなのかを決定する。音源位置測定に関する幾つかの従来方式を使用することができる。親出願 (特許文献 1) に記載のものを含めて、様々なフェイス・トラッカ方式 (または人物検出および追跡システム) が、本明細書の説明のように使用することができる。

【0030】

メモリ 226 は、また、話し手クラスタリング・モジュール 236 を格納し、これは、2 人以上の人物が話している際に主要な話し手を判断して、主要な話し手にタイムラインの特定部分を集中させるように、構成される。ほとんどの会議の状況では、2 人以上の人物が同時に話すことがある。通常は、主要な話し手が話している最中に別の人物が短時間話し手に横やりを入れたり、あるいは話し手に重なって話す。話し手クラスタリング・モジュール 236 は、話し手を時間的かつ空間的にクラスタ化してタイムラインを整理するように構成される。

【0031】

タイムライン 238 は、仮想シネマトグラフィ 234 によって作成される。タイムライン 238 は、大容量記憶装置 242 上のタイムライン・データベース 244 に格納される

10

20

30

40

50



。タイムライン・データベース 244 は、複数のフィールドを含む。これらのフィールドとしては、時間、話し手の番号、およびカメラ画像（x、y、幅、高さ）内の話し手バウンディング・ボックス（bounding box）等があるが、必ずしもそれだけには限定されない。タイムライン・データベース 244 はさらに、話し手の 1 つまたは複数の顔の角度（方位角と仰角）を含むことができる。

#### 【0032】

メモリ 226 には、フェイス・エクストラクタ（顔領域抽出）モジュール 240 も格納され、カメラ画像の（フェイス・トラッカ 230 が識別した）顔バウンディング・ボックスから話し手の顔の画像を抽出するように構成される。フェイス・エクストラクタ・モジュール 240 は、抽出した顔画像を大容量記憶装置 242 上の顔データベース 246 に格納する。

10

#### 【0033】

少なくとも 1 つの実装では、1 人または複数の話し手について複数の顔画像を格納することができる。どの顔画像がいつ使用されるかを決定するために、パラメータを指定することができる。あるいは、ユーザが、複数の顔画像から手動で特定の顔画像を選択することを可能とすることもできる。

#### 【0034】

少なくとも 1 つの代替実装では、それぞれの話し手に 1 つの顔画像のみが格納される。記憶される顔画像は、フェイス・エクストラクタ・モジュール 240 が抽出した単一の画像でもよいが、フェイス・エクストラクタ・モジュール 240 は、話し手の最良の画像を選択するように構成することもできる。

20

#### 【0035】

話し手の最良の画像の選択は、（正面の顔の画像が別の画像よりも優れた写真であると仮定して）正面の顔角度を識別することによって、最小の動きを示す顔画像を識別することによって、あるいは顔の対称性を最大化している顔画像を識別することによって、行うことができる。

#### 【0036】

録画された会議 248 も、後ほど呼び出して、再生できるように、大容量記憶装置 242 に記憶される。

#### 【0037】

30

図 2 に示され、図 2 と関連して説明した要素および機能は、後続の図面に関連して、以下に、より詳しく説明する。

#### 【0038】

例示的な再生画面

図 3 は、パノラマ画像 302 および顔画像タイムライン 304 を含む再生画面 300 の線画である。第 1 の会議参加者 303 および第 2 の会議参加者 305 を含むパノラマ画像 302 が、と示されている。図の再生画面 300 は、さらにタイトル・バー 306 および個人画像 308 も含んでいる。個人画像 308 は、オプション機能であり、そこでは、特定の個人、一般には主要な話し手、に焦点が当てられる。図 3 では、その個人画像 308 は、第 1 の会議参加者 303 の顔画像を表示している。

40

#### 【0039】

例示的な再生画面 300 は、さらに再生ボタン、早送りボタン、巻き戻しボタンなど、メディア・プレイヤーで通常見られる制御手段を備える制御セクション 310 を含む。再生画面 300 は、再生画面 300 の主題に関する情報を表示できる情報エリア 312 を含む。例えば、会議の題名、会議室の番号、会議出席者のリストなどを情報エリア 312 に表示してもよい。

#### 【0040】

顔画像タイムライン 304 は、第 1 の会議参加者 303 に対応する第 1 のサブタイムライン 314、および第 2 の会議参加者 305 に対応する第 2 のサブタイムライン 316 を含む。それぞれのサブタイムライン 314 および 316 は、ある時間的連続体（temporal

50

continuum)の中で対応する会議参加者が話している箇所(セクション)を示す。ユーザは、サブタイムライン314、316の任意の時点に直接アクセスして、特定の会議参加者が話している会議の部分に直ちにアクセスすることができる。

【0041】

第1の会議参加者303の第1の顔画像318が、第1のサブタイムライン314の隣に現れて、第1のサブタイムライン314が第1の会議参加者303に対応することを示す。第2の会議参加者305の顔画像320が第2のサブタイムライン316の隣に現れて、第2のサブタイムライン316が第2の会議参加者305に対応することを示す。

【0042】

図4は、図3に示し説明した例示的再生画面300に類似した要素を含む例示的再生画面400を示す。図3に関して示し説明した要素および参照番号を、図4の例示的再生画面400に関して使用する。

10

【0043】

例示的再生画面400は、パノラマ画像302および顔画像タイムライン304を含む。パノラマ画像302は、第1の会議参加者303および第2の会議参加者305を示す。タイトル・バー306が再生画面400の上部の両端間に延びており、個人画像308は第2の会議参加者305を示している。

【0044】

例示的再生画面400は、さらに、ホワイト・ボードの前に位置する会議参加者(この場合は、第2の会議参加者305)を表示するホワイト・ボード話し手画像を含む。ホワイト・ボード話し手画像は、図3の再生画面300に含まれていないが、ここではどのようにして、任意の特定再生画面300、400に他の画像をどのように含むことができるかを説明するために、使用している。

20

【0045】

制御セクション310は、マルチメディア制御を含み、情報エリア312は再生画面400で表示されている会議に関する情報を表示する。

【0046】

顔画像タイムライン304は、第1のサブタイムライン314、第2のサブタイムライン316、および第3のサブタイムライン402を含む。図3には2つのサブタイムラインしか示していないが、タイムラインは管理できる限り任意の数のサブタイムラインを含むことができることに留意されたい。例えば、図4は、3つのサブタイムラインが存在する。

30

【0047】

この例では会議参加者が2名しかいないにも関わらず、サブタイムラインが3つあることに留意されたい。これは、単一の話し手が2つ以上のサブタイムラインに関連付けられる可能性が存在するからである。この例では、第2のサブタイムライン316は、ホワイト・ボードの前にいるときの第2の会議参加者305に関連付けられており、第3のサブタイムライン404はホワイト・ボードの前以外の場所にいるときの第2の会議参加者305に関連付けられている。

【0048】

40

この状況は、会議参加者が会議中に複数の場所を占める場合に起こる。仮想シネマトグラフィ234はこの場合、3つの場所で話し手を検出した。仮想シネマトグラフィ234は、これらの場所には2人の話し手しかいないことを必ずしも知らない。この機能は、ユーザが、ある特定の位置にいるときの話し手に、主に関心がある場合に、ユーザの役に立つ。例えば、ユーザは、録画した会議の中で、ある話し手がホワイト・ボードの前にいる箇所のみを再生したいと思うかもしれない。

【0049】

例示的再生画面400はさらに、第1の会議参加者303の第1の顔画像318、および第2の会議参加者305の第2の顔画像320を含む。さらに、第3の顔画像404も含まれ、これは第3のサブタイムライン402に関連する。第3の顔画像404は、第2

50

の会議参加者 3 0 5 の第 2 の場所に対応する。

【 0 0 5 0 】

例示的再生画面 3 0 0、4 0 0 を示すのに使用した技術については、後続の図面に関して以下により詳しく説明する。

【 0 0 5 1 】

例示的な方法論的実装：顔画像タイムラインの作成

図 5 は、顔画像付きのタイムラインを作成するための方法論的実装の例示的なフロー図 5 0 0 である。例示的フロー図 5 0 0 についての以下の説明では、前の図に示される要素および参照番号を引き続き参照する。

【 0 0 5 2 】

ブロック 5 0 2 で、パノラマ式カメラ機器 2 0 0 が、1 つまたは複数のビデオ画像を抽出してパノラマ画像を作成する。パノラマ画像は、画像内の顔を検出して追跡するフェイス・トラッカ 2 3 0 に入力される（ブロック 5 0 4）。ほぼ同時に、ブロック 5 0 6 で、マイクロホン・アレイ 2 0 8 がパノラマ画像に対応する音を抽出し、その音をブロック 5 0 8 で、抽出した音に基づいて話し手の位置を検出する音源ローカライザ 2 3 2 に入力する。

【 0 0 5 3 】

仮想シネマトグラフィ 2 3 4 は、フェイス・トラッカ 2 3 0 および音源ローカライザ 2 3 2 からのデータを処理して、ブロック 5 1 0 でタイムライン 2 3 8 を作成する。ブロック 5 1 2 で、先に述べたように、話し手クラスタリング・モジュール 2 3 6 が、話し手を時間的かつ空間的にクラスタ化して、タイムライン 2 3 8 の諸部分を統合して明確にする。

【 0 0 5 4 】

タイムラインは、次に挙げるフィールド、すなわち、時間、話し手の番号、画像内の話し手バウンディング・ボックス（x、y、幅、高さ）、話し手の顔の角度（方位角と仰角）など、を付けて、タイムライン・データベース 2 4 4 に記憶される。

【 0 0 5 5 】

フェイス・トラッカ 2 3 0 が導出したパノラマ画像および顔認識座標（すなわち、顔バウンディング・ボックス）を使用して、ブロック 5 1 4 でフェイス・エクストラクタ 2 4 0 が話し手の顔画像を抽出する。抽出された顔画像は、顔データベース 2 4 6 に記憶され、話し手番号に関連付けられる。

【 0 0 5 6 】

先に述べたように、フェイス・エクストラクタ 2 4 0 は、各話し手毎に複数の画像を抽出して、フェイス・エクストラクタ 2 4 0 がタイムライン 2 3 8 の中で最良の画像であると判断した画像を使用するように、構成することができる。

【 0 0 5 7 】

「最良」の画像を選択し、顔データベース 2 4 6 を作成するための例示的な方法論的実装を、図 6 に関して以下に示し説明する。

【 0 0 5 8 】

例示的な方法論的実装：顔データベースの作成

図 6 は、顔データベースを作成するための方法論的実装を示す例示的なフロー図 6 0 0 である。図 6 についての以下の説明では、前の複数の図に示された要素および参照番号を引き続き参照する。

【 0 0 5 9 】

ブロック 6 0 2 で、フェイス・エクストラクタ 2 4 0 が先に述べたようにパノラマ画像から顔画像を抽出する。顔データベース 2 4 6 に話し手の顔画像がまだ記憶されていない場合（ブロック 6 0 4 の「No」ブランチ）、顔画像はブロック 6 1 0 で顔データベース 2 4 6 に記憶される。顔画像が記憶されているかどうかの判断は、顔画像に登場する人物の肖像の画像が既に記憶されている画像と類似するイメージを有するかどうかにも必ずしも依存せず、むしろ識別された話し手がその話し手に対応する既に記憶されている画像を有

10

20

30

40

50

するかどうかに依存することに留意されたい。したがって、第１の場所にいる、ある話し手が、既に格納された顔画像を有し、その後、その話し手が第２の場所で検出された場合、その話し手が既に記憶されている顔画像を有するかどうかを判断するのに、第２の場所にいる話し手の顔画像を第１の場所にいる話し手の記憶済み顔画像と比較することはない。

#### 【００６０】

話し手の顔画像が既に顔データベース２４６に記憶されている（以下、「記憶済み顔画像」と呼ぶ）場合（ブロック６０４の「Ｙｅｓ」ブランチ）、その顔画像はブロック６０６で記憶済み顔画像と比較される。フェイス・エクストラクタ２４０が、その顔画像が記憶済み顔画像よりも良いか、あるいはより好ましいと判断した場合（ブロック６０８の「

10

#### 【００６１】

顔画像が記憶済み顔画像よりも良くなかった場合（ブロック６０８の「Ｎｏ」ブランチ）、顔画像は破棄され、記憶済み顔画像が保持される。

#### 【００６２】

どの顔画像がより優れた顔画像かを判断する基準は数多く、様々である。例えば、話し手が最も正面を向いたところをとらえた顔画像を「最良」の顔画像と判断するようにフェイス・エクストラクタ２４６を構成することができる。あるいは、第１の顔画像が動きを示し、第２の顔画像が動きを示さない場合、第２の顔画像が「最良」の顔画像であるとフェイス・エクストラクタ２４６が判断してもよい。あるいは、話し手の複数の画像のうち、どれが最大の対称性を示しているかを決定し、その顔画像をタイムラインに使用するようにフェイス・エクストラクタ２４６を構成することもできる。タイムラインに使用する最適の顔画像の判断には、ここに列挙した以外の基準も使うことができる。

20

#### 【００６３】

他の話し手が存在する場合（ブロック６１２の「Ｙｅｓ」ブランチ）、処理はブロック６０２に戻り、一意の各話し手毎にこの処理が繰り返される。この場合も、このコンテキストで使用される「一意の話し手」は、異なる発言位置（speaking locations）に登場する人物が異なる話し手と解釈される可能性もあるので、必ずしも固有の人物を意味する必要はない。からである。識別すべき一意の話し手がそれ以上存在しない場合、処理は終了する（ブロック６１２の「Ｎｏ」ブランチ）。

30

#### 【００６４】

##### 結論

１つまたは複数の例示的な実装を図示し説明したが、本明細書に添付した特許請求の範囲の精神および範囲から逸脱することなしに、様々な変更が行えることが理解されよう。

#### 【図面の簡単な説明】

#### 【００６５】

【図１】例示的な汎用コンピューティング／カメラ装置を示すブロック図である。

【図２】例示的なパノラマ式カメラおよびクライアント装置を表すブロック図である。

【図３】パノラマ画像および顔画像のタイムラインを含む例示的な再生画面を表現した図である。

40

【図４】パノラマ画像および顔画像のタイムラインを含む例示的な再生画面を示す図である。

【図５】顔画像付きのタイムライン作成の方法論的な実装の例示的なフロー図である。

【図６】顔データベースの作成の方法論的な実装を示す例示的なフロー図である。

#### 【符号の説明】

#### 【００６６】

- １００ コンピューティング・システム環境
- １１０ コンピュータ
- １２０ 処理装置

50

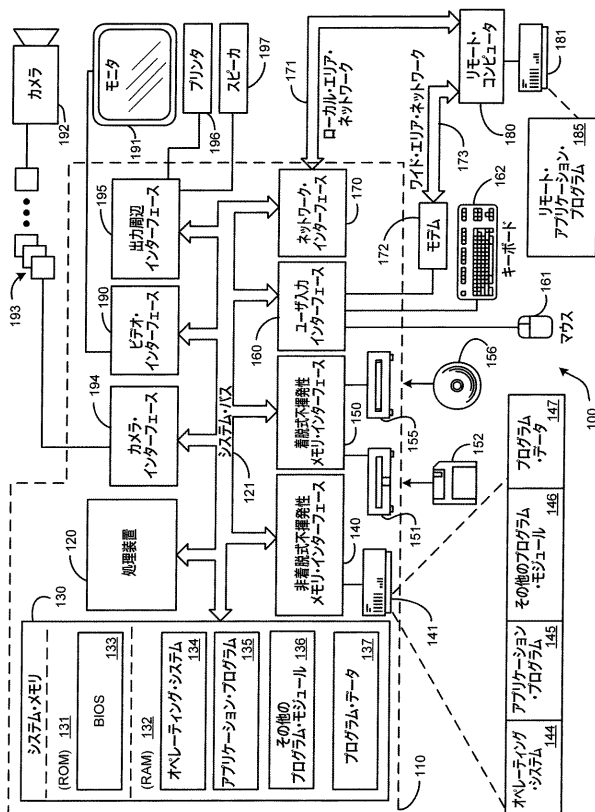
1 2 1	システム・バス	
1 3 0	システム・メモリ	
1 3 4	オペレーティング・システム	
1 3 5	アプリケーション・プログラム	
1 3 6	その他のプログラム・モジュール	
1 3 7	プログラム・データ	
1 4 0	非着脱式不揮発性メモリ・インターフェース	
1 4 1	ハードディスク・ドライブ	
1 4 4	オペレーティング・システム	
1 4 5	アプリケーション・プログラム	10
1 4 6	その他のプログラム・モジュール	
1 4 7	プログラム・データ	
1 5 0	着脱式不揮発性メモリ・インターフェース	
1 5 1	磁気ディスク・ドライブ	
1 5 2	磁気ディスク	
1 5 5	光ディスク・ドライブ	
1 5 6	光ディスク	
1 6 0	ユーザ入力インターフェース	
1 6 1	マウス	
1 6 2	キーボード	20
1 7 0	ネットワーク・インターフェース	
1 7 1	ローカル・エリア・ネットワーク	
1 7 2	モデム	
1 7 3	ワイド・エリア・ネットワーク	
1 8 0	リモート・コンピュータ	
1 8 1	メモリ装置	
1 8 5	リモート・アプリケーション・プログラム	
1 9 0	ビデオ・インターフェース	
1 9 1	モニタ	
1 9 2	カメラ	30
1 9 4	カメラ・インターフェース	
1 9 5	出力周辺インターフェース	
1 9 6	プリンタ	
1 9 7	スピーカ	
2 0 0	パノラマ式カメラ機器	
2 0 2	プロセッサ	
2 0 4	メモリ	
2 0 6	カメラ__1 . . . カメラ__N	
2 0 8	マイクロホン・アレイ	
2 1 0	スピーカ	40
2 1 2	設定	
2 1 4	OS	
2 1 6	アプリケーション	
2 1 8	I/O	
2 2 0	その他のハードウェア	
2 2 2	クライアント装置	
2 2 4	プロセッサ	
2 2 6	メモリ	
2 3 0	ハードウェア	
2 3 0	フェイス・トラッカ	50

- 2 3 2 音源ローカライザ  
 2 3 4 仮想シネマトグラフィ  
 2 3 6 話し手のクラスタリング  
 2 3 8 タイムライン  
 2 4 0 フェイス・エクストラクタ  
 2 4 2 大容量記憶装置  
 2 4 4 タイムライン・データベース  
 2 4 6 顔データベース  
 2 4 8 録画された会議  
 3 0 0 再生画面  
 3 0 2 パノラマ画像  
 3 0 3 第1の会議参加者  
 3 0 4 顔画像タイムライン  
 3 0 5 第2の会議参加者  
 3 0 6 タイトル・バー  
 3 0 8 個人画像  
 3 1 0 制御セクション(コントロール)  
 3 1 2 情報エリア  
 3 1 4 サブタイムライン  
 3 1 6 サブタイムライン  
 3 1 8 第1の会議参加者 3 0 3 の顔画像  
 3 2 0 第2の会議参加者 3 0 5 の顔画像  
 4 0 0 例示的再生画面

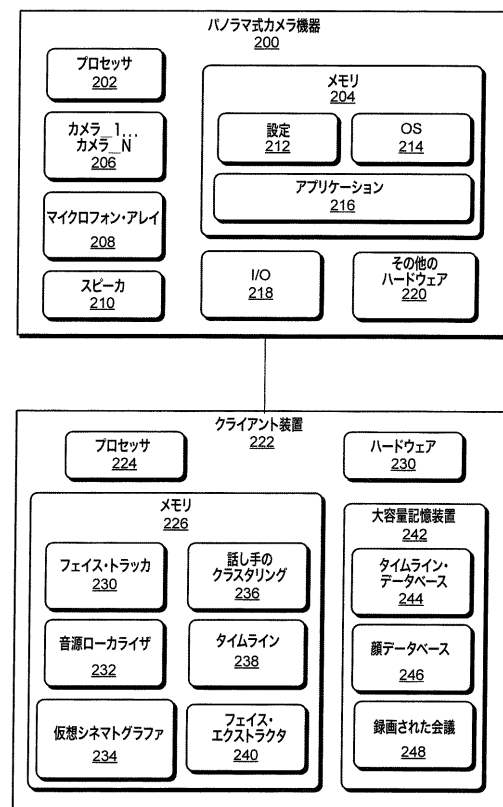
10

20

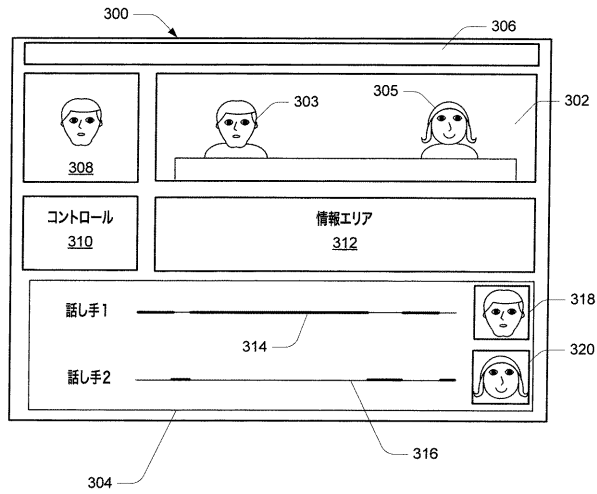
【図1】



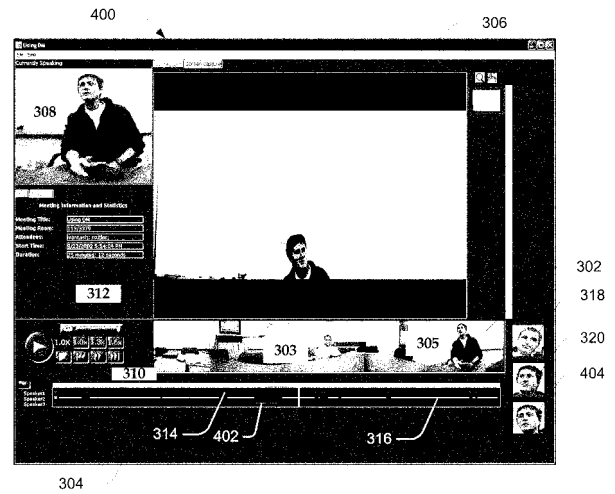
【図2】



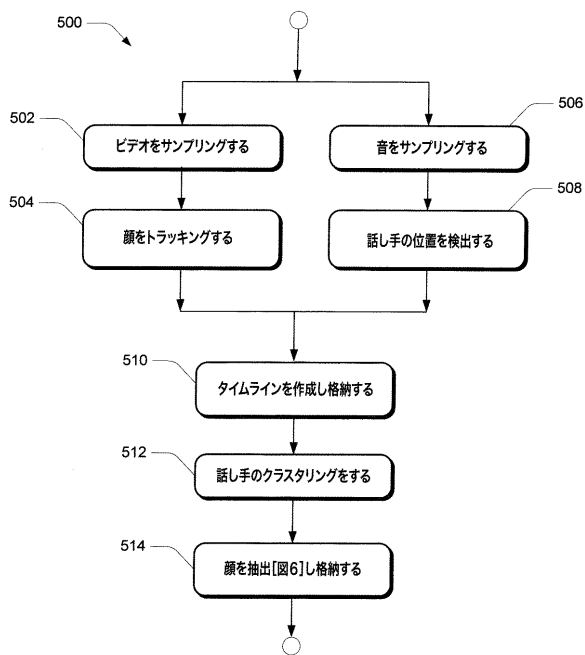
【図 3】



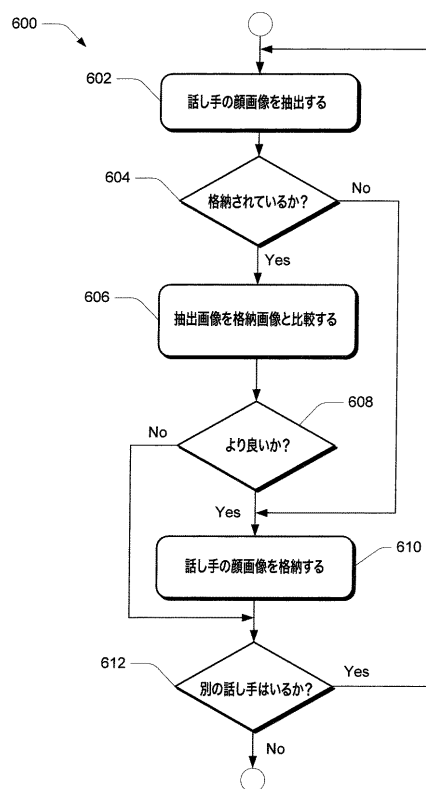
【図 4】



【図 5】



【図 6】



---

フロントページの続き

(74)代理人 100091063

弁理士 田中 英夫

(74)代理人 100153028

弁理士 上田 忠

(74)代理人 100120112

弁理士 中西 基晴

(74)代理人 100113974

弁理士 田中 拓人

(72)発明者 ロス ジー・カトラー

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション内

審査官 竹中 辰利

(56)参考文献 特開2000-125274(JP,A)

特開2003-230049(JP,A)

特開2000-036052(JP,A)

特開2002-251393(JP,A)

再公表特許第97/009683(JP,A1)

(58)調査した分野(Int.Cl., DB名)

H04N 5/76

G06T 7/00