

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4474013号
(P4474013)

(45) 発行日 平成22年6月2日 (2010.6.2)

(24) 登録日 平成22年3月12日 (2010.3.12)

(51) Int. Cl.

F I

G 1 O L 15/24 (2006.01)

G 1 O L 15/24 Q

G 1 O L 15/28 (2006.01)

G 1 O L 15/28 4 O O

G 1 O L 15/00 (2006.01)

G 1 O L 15/00 2 O O U

G O 6 F 3/16 (2006.01)

G O 6 F 3/16 3 2 O H

G O 6 F 17/30 (2006.01)

G O 6 F 17/30 1 7 O B

請求項の数 78 外国語出願 (全 54 頁)

(21) 出願番号 特願2000-86807 (P2000-86807)
 (22) 出願日 平成12年3月27日 (2000.3.27)
 (65) 公開番号 特開2000-352996 (P2000-352996A)
 (43) 公開日 平成12年12月19日 (2000.12.19)
 審査請求日 平成19年3月27日 (2007.3.27)
 (31) 優先権主張番号 9907103.7
 (32) 優先日 平成11年3月26日 (1999.3.26)
 (33) 優先権主張国 英国 (GB)
 (31) 優先権主張番号 9908546.6
 (32) 優先日 平成11年4月14日 (1999.4.14)
 (33) 優先権主張国 英国 (GB)

(73) 特許権者 000001007
 キヤノン株式会社
 東京都大田区下丸子3丁目30番2号
 (74) 代理人 100076428
 弁理士 大塚 康德
 (74) 代理人 100112508
 弁理士 高柳 司郎
 (74) 代理人 100116894
 弁理士 木村 秀二
 (74) 代理人 100115071
 弁理士 大塚 康弘

最終頁に続く

(54) 【発明の名称】 情報処理装置

(57) 【特許請求の範囲】

【請求項 1】

画像データ及び音声データを処理する装置であって、
 少なくとも1台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物を三次元で追跡する画像処理手段と、
 音声データを処理して、音声の到達方向を判定する音声処理手段と、
 画像処理手段により実行される処理の結果と、音声処理手段により実行される処理の結果とに基づいて、どの人物が話しているかを判定する話者識別手段と、
 受信した音声データを処理して、話者識別手段により実行される処理の結果に従って音声データからテキストデータを生成する音声認識処理手段とを備えることを特徴とする装置。

【請求項 2】

音声認識処理手段は、人物ごとの対応する音声認識パラメータを格納する記憶手段と、話者識別手段により話していると判定された人物に従って音声データを処理するために使用すべき音声認識パラメータを選択する手段とを含むことを特徴とする請求項1記載の装置。

【請求項 3】

画像処理手段は、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより各人物を追跡するように構成されていることを特徴とする請求項1又は2記載の装置。

【請求項 4】

画像処理手段は、各人物の頭部を追跡することにより各人物を追跡するように構成されていることを特徴とする請求項 1 乃至 3 のいずれか 1 項に記載の装置。

【請求項 5】

画像処理手段は、少なくとも話をしている各人物がどこを見ているかを判定するために画像データを処理するように構成されていることを特徴とする請求項 1 乃至 4 のいずれか 1 項に記載の装置。

【請求項 6】

話者識別手段は、受信した画像データの所定のフレームについて画像処理手段及び音声処理手段により実行される処理の結果を使用して話者を識別できない場合、少なくとも 1 つの他のフレームに対して画像処理手段及び音声処理手段により実行される処理の結果を使用して所定のフレームにおいて話をしている人物を識別するように構成されていることを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載の装置。

10

【請求項 7】

受信した画像データ、音声データ、音声認識処理手段により生成されるテキストデータ、及び少なくとも話をしている各人物がどこを見ているかを定義する視線データの少なくとも一部を格納するデータベースを更に具備し、前記データベースは、対応するテキストデータと視線データとが互いに関連すると共に、対応する画像データ及び音声データとも関連するようにデータを格納する構成であることを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の装置。

20

【請求項 8】

データベースに格納するために、画像データ及び音声データを圧縮する手段を更に具備することを特徴とする請求項 7 記載の装置。

【請求項 9】

画像データ及び音声データを圧縮する手段は、画像データ及び音声データを MPEG データとして符号化する手段であることを特徴とする請求項 8 記載の装置。

【請求項 10】

所定の期間にわたり、その所定の期間中に所定の人物がその他の人物の各々を見るのに費やした時間の割合を定義するデータを生成する手段を更に具備し、データベースは、そのデータが対応する画像データ、音声データ、テキストデータ及び視線データと関連するようにデータを格納する構成であることを特徴とする請求項 7 乃至 9 のいずれか 1 項に記載の装置。

30

【請求項 11】

所定の期間は所定の人物が話をしていた期間である請求項 10 記載の装置。

【請求項 12】

画像データ及び音声データを処理する装置において、

少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物を三次元で追跡する画像処理手段と、

音声データを処理して、音声の到達方向を判定する音声処理手段と、

画像処理手段により実行される処理の結果と、音声処理手段により実行される処理の結果とに基づいて、どの人物が話しているかを判定する話者識別手段とを具備することを特徴とする装置。

40

【請求項 13】

画像処理手段は、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより各人物を追跡するように構成されていることを特徴とする請求項 12 記載の装置。

【請求項 14】

画像処理手段は、各人物の頭部を追跡することにより各人物を追跡するように構成されていることを特徴とする請求項 12 又は 13 に記載の装置。

【請求項 15】

50

画像処理手段は、少なくとも話をしている各人物がどこを見ているかを判定するために画像データを処理するように構成されていることを特徴とする請求項 1 2 乃至 1 4 のいずれか 1 項に記載の装置。

【請求項 1 6】

話者識別手段は、受信した画像データの所定のフレームについて画像処理手段及び音声処理手段により実行される処理の結果を使用して話者を識別できない場合、少なくとも 1 つの他のフレームに対して画像処理手段及び音声処理手段により実行される処理の結果を使用して所定のフレームにおいて話をしている人物を識別するように構成されていることを特徴とする請求項 1 2 乃至 1 5 のいずれか 1 項に記載の装置。

【請求項 1 7】

画像データ及び音声データを処理する方法において、
少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物を三次元で追跡する画像処理ステップと、
音声データを処理して、音声の到達方向を判定する音声処理ステップと、
画像処理ステップで実行される処理の結果と、音声処理ステップで実行される処理の結果とに基づいて、どの人物が話しているかを判定する話者識別ステップと、
受信した音声データを処理して、話者識別ステップで実行される処理の結果に従って音声データからテキストデータを生成する音声認識処理ステップとを備えることを特徴とする方法。

【請求項 1 8】

音声認識処理ステップは、人物ごとに格納されている音声認識パラメータから、話者識別ステップで話していると判定された人物に従って音声データを処理するために使用すべき音声認識パラメータを選択することを含むことを特徴とする請求項 1 7 記載の方法。

【請求項 1 9】

画像処理ステップでは、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより各人物を追跡することを特徴とする請求項 1 7 又は 1 8 記載の方法。

【請求項 2 0】

画像処理ステップでは、各人物の頭部を追跡することにより各人物を追跡することを特徴とする請求項 1 7 乃至 1 9 のいずれか 1 項に記載の方法。

【請求項 2 1】

画像処理ステップでは、少なくとも話をしている各人物がどこを見ているかを判定するために画像データを処理することを特徴とする請求項 1 7 乃至 2 0 のいずれか 1 項に記載の方法。

【請求項 2 2】

話者識別ステップでは、受信した画像データの所定のフレームについて画像処理ステップ及び音声処理ステップにより実行される処理の結果を使用して話者を識別できない場合、少なくとも 1 つの他のフレームに対して画像処理ステップ及び音声処理ステップにより実行される処理の結果を使用して所定のフレームにおいて話をしている人物を識別することを特徴とする請求項 1 7 乃至 2 1 のいずれか 1 項に記載の方法。

【請求項 2 3】

音声認識処理ステップで生成されるデータを搬送する信号を発生するステップを更に含むことを特徴とする請求項 1 7 乃至 2 2 のいずれか 1 項に記載の方法。

【請求項 2 4】

受信した画像データ、音声データ、音声認識処理ステップにより生成されるテキストデータ、及び少なくとも話をしている各人物がどこを見ているかを定義する視線データの少なくとも一部をデータベースに格納するステップを更に含み、データは、対応するテキストデータと視線データとが互いに関連すると共に、対応する画像データ及び音声データとも関連するようにデータベースに格納されることを特徴とする請求項 1 7 乃至 2 3 のいずれか 1 項に記載の方法。

【請求項 2 5】

画像データ及び音声データは圧縮された形態でデータベースに格納されることを特徴とする請求項 2 4 記載の方法。

【請求項 2 6】

画像データ及び音声データはMPEGデータとして格納されることを特徴とする請求項 2 5 記載の方法。

【請求項 2 7】

所定の期間にわたり、その所定の期間中に所定の人物がその他の人物の各々を見るのに費やした時間の割合を定義するデータを生成するステップと、そのデータが対応する画像データ、音声データ、テキストデータ及び視線データと関連するようにデータをデータベースに格納するステップとを更に含むことを特徴とする請求項 2 4 乃至 2 6 のいずれか 1 項に記載の方法。

10

【請求項 2 8】

所定の期間は所定の人物が話をしていた期間であることを特徴とする請求項 2 7 記載の方法。

【請求項 2 9】

データベースを格納されているデータと共に搬送する信号を発生するステップを更に含むことを特徴とする請求項 2 4 乃至 2 8 のいずれか 1 項に記載の方法。

【請求項 3 0】

信号の記録を生成するために信号を直接に又は間接的に記録するステップを更に含むことを特徴とする請求項 2 9 記載の方法。

20

【請求項 3 1】

画像データ及び音声データを処理する方法において、
少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物を三次元で追跡する画像処理ステップと、
音声データを処理して、音声の到達方向を判定する音声処理ステップと、
画像処理ステップにより実行される処理の結果と、音声処理ステップにより実行される処理の結果とに基づいて、どの人物が話しているかを判定する話者識別ステップとを備えることを特徴とする方法。

30

【請求項 3 2】

画像処理ステップでは、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより各人物を追跡することを特徴とする請求項 3 1 記載の方法。

【請求項 3 3】

画像処理ステップでは、各人物の頭部を追跡することにより各人物を追跡することを特徴とする請求項 3 1 又は 3 2 記載の方法。

【請求項 3 4】

画像処理ステップでは、少なくとも話をしている各人物がどこを見ているかを判定するために画像データを処理することを特徴とする請求項 3 1 乃至 3 3 のいずれか 1 項に記載の方法。

40

【請求項 3 5】

話者識別ステップでは、受信した画像データの所定のフレームについて画像処理ステップ及び音声処理ステップにより実行される処理の結果を使用して話者を識別できない場合、少なくとも 1 つの他のフレームに対して画像処理ステップ及び音声処理ステップにより実行される処理の結果を使用して所定のフレームにおいて話をしている人物を識別することを特徴とする請求項 3 1 乃至 3 4 のいずれか 1 項に記載の方法。

【請求項 3 6】

話者識別ステップで識別された話者のアイデンティティを搬送する信号を発生するステップを更に含むことを特徴とする請求項 3 1 乃至 3 5 のいずれか 1 項に記載の方法。

【請求項 3 7】

50

プログラム可能処理装置を請求項 1 乃至 16 の少なくとも 1 項に記載の装置として構成させるための命令を格納することを特徴とする記憶装置。

【請求項 38】

プログラム可能処理装置を請求項 17 乃至 36 の少なくとも 1 項に記載の方法を実行するように動作可能にさせるための命令を格納することを特徴とする記憶装置。

【請求項 39】

画像データ及び音声データを処理する装置であって、

少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物がどこを見ているかを判定すると共に、複数の人物がどこを見ているかに基づいてどの人物が話しているかを判定する画像処理手段と、

人物が話した言葉を定義する音声データを処理して、画像処理手段により実行される処理の結果に従って音声データからテキストデータを生成する音声処理手段とを備えることを特徴とする装置。

【請求項 40】

音声処理手段は、人物ごとの対応する音声認識パラメータを格納する記憶手段と、画像処理手段により話をしていると判定された人物に従って音声データを処理するために使用する音声認識パラメータを選択する手段とを含むことを特徴とする請求項 39 記載の装置。

【請求項 41】

画像処理手段は、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより、各人物がどこを見ているかを判定するように構成されていることを特徴とする請求項 39 又は 40 記載の装置。

【請求項 42】

画像処理手段は、各人物の頭部の位置と向きを三次元で追跡するために画像データを処理することにより、各人物がどこを見ているかを判定するように構成されていることを特徴とする請求項 39 乃至 41 のいずれか 1 項に記載の装置。

【請求項 43】

画像処理手段は、各々の人物を見ている人物の数に基づいてどの人物が話をしているかを判定するように構成されていることを特徴とする請求項 39 乃至 42 のいずれか 1 項に記載の装置。

【請求項 44】

画像処理手段は、各人物が誰を見ているかを定義する値を人物ごとに生成し且つそれらの値を処理して、話をしている人物を判定するように構成されていることを特徴とする請求項 43 記載の装置。

【請求項 45】

画像処理手段は、話をしている人物が他の人物の大半が見ている人物であることを判定するように構成されていることを特徴とする請求項 39 乃至 44 のいずれか 1 項に記載の装置。

【請求項 46】

画像データ、音声データ、音声処理手段により生成されるテキストデータ、及び各人物がどこを見ているかを定義する視線データを格納するデータベースを更に備え、前記データベースは、対応するテキストデータと視線データとが互いに関連すると共に、対応する画像データ及び音声データとも関連するようにデータを格納する構成であることを特徴とする請求項 39 乃至 45 のいずれか 1 項に記載の装置。

【請求項 47】

データベースに格納するために、画像データ及び音声データを圧縮する手段を更に備えることを特徴とする請求項 46 記載の装置。

【請求項 48】

画像データ及び音声データを圧縮する手段は、画像データ及び音声データを MPEG データとして符号化する手段であることを特徴とする請求項 47 記載の装置。

10

20

30

40

50

【請求項 4 9】

所定の期間にわたり、その所定の期間中に所定の人物がその他の人物の各々を見るのに費やした時間の割合を定義するデータを生成する手段を更に備え、データベースは、そのデータが対応する画像データ、音声データ、テキストデータ及び視線データと関連するようにデータを格納する構成であることを特徴とする請求項 4 6 乃至 4 8 のいずれか 1 項に記載の装置。

【請求項 5 0】

所定の期間は所定の人物が話をしていた期間であることを特徴とする請求項 4 9 記載の装置。

【請求項 5 1】

画像データを処理する装置において、少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物がどこを見ているかを判定すると共に、複数の人物がどこを見ているかに基づいてどの人物が話しているかを判定する画像処理手段を備えることを特徴とする装置。

【請求項 5 2】

画像処理手段は、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより、各人物がどこを見ているかを判定するように構成されていることを特徴とする請求項 5 1 記載の装置。

【請求項 5 3】

画像処理手段は、各人物の頭部の位置と向きを三次元で追跡するために画像データを処理することにより、各人物がどこを見ているかを判定するように構成されていることを特徴とする請求項 5 1 又は 5 2 記載の装置。

【請求項 5 4】

画像処理手段は、各々の人物を見ている人物の数に基づいてどの人物が話をしているかを判定するように構成されていることを特徴とする請求項 5 1 乃至 5 3 のいずれか 1 項に記載の装置。

【請求項 5 5】

画像処理手段は、各人物が誰を見ているかを定義する値を人物ごとに生成し且つそれらの値を処理して、話をしている人物を判定するように構成されていることを特徴とする請求項 5 4 記載の装置。

【請求項 5 6】

画像処理手段は、話をしている人物が他の人物の大半が見ている人物であることを判定するように構成されていることを特徴とする請求項 5 1 乃至 5 5 のいずれか 1 項に記載の装置。

【請求項 5 7】

画像データ及び音声データを処理する方法において、

少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物がどこを見ているかを判定すると共に、複数の人物がどこを見ているかに基づいてどの人物が話しているかを判定する画像処理ステップと、

人物が話した言葉を定義する音声データを処理して、画像処理手段により実行される処理の結果に従って音声データからテキストデータを生成する音声処理ステップとを備えることを特徴とする方法。

【請求項 5 8】

音声処理ステップは、人物ごとに格納されているそれぞれの音声認識パラメータから、画像処理ステップで話をしていると判定された人物に従って音声データを処理するために使用すべき音声認識パラメータを選択するステップを含むことを特徴とする請求項 5 7 記載の方法。

【請求項 5 9】

画像処理ステップでは、処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより、各人物がどこを見て

10

20

30

40

50

いるかを判定することを特徴とする請求項 5 7 又は 5 8 記載の方法。

【請求項 6 0】

画像処理ステップでは、各人物の頭部の位置と向きを三次元で追跡するために画像データを処理することにより、各人物がどこを見ているかを判定することを特徴とする請求項 5 7 乃至 5 9 のいずれか 1 項に記載の方法。

【請求項 6 1】

画像処理ステップでは、各々の人物を見ている人物の数に基づいてどの人物が話をしているかを判定することを特徴とする請求項 5 7 乃至 6 0 のいずれか 1 項に記載の方法。

【請求項 6 2】

画像処理ステップでは、各人物が誰を見ているかを定義する値を人物ごとに生成し且つそれらの値を処理して、話をしている人物を判定することを特徴とする請求項 6 1 記載の方法。

【請求項 6 3】

画像処理ステップでは、話をしている人物が他の人物の大半が見ている人物であることを判定することを特徴とする請求項 5 7 乃至 6 2 のいずれか 1 項に記載の方法。

【請求項 6 4】

画像データ、音声データ、音声処理手段により生成されるテキストデータ、及び各人物がどこを見ているかを定義する視線データをデータベースに格納するステップを更に含み、前記データベースは、対応するテキストデータと視線データとが互いに関連すると共に、対応する画像データ及び音声データとも関連するようにデータを格納することを特徴とする請求項 5 7 乃至 6 3 のいずれか 1 項に記載の方法。

【請求項 6 5】

画像データ及び音声データは圧縮された形態で格納されることを特徴とする請求項 6 4 記載の方法。

【請求項 6 6】

画像データ及び音声データはMPEGデータとして格納されることを特徴とする請求項 6 5 記載の方法。

【請求項 6 7】

所定の期間にわたり、その所定の期間中に所定の人物がその他の人物の各々を見るのに費やした時間の割合を定義するデータを生成するステップと、そのデータが対応する画像データ、音声データ、テキストデータ及び視線データと関連するようにデータをデータベースに格納するステップとを更に含むことを特徴とする請求項 6 4 乃至 6 6 のいずれか 1 項に記載の方法。

【請求項 6 8】

所定の期間は所定の人物が話をしていた期間であることを特徴とする請求項 6 7 記載の方法。

【請求項 6 9】

データベースを格納されているデータと共に搬送する信号を発生するステップを更に含むことを特徴とする請求項 6 4 乃至 6 8 のいずれか 1 項に記載の方法。

【請求項 7 0】

信号の記録を生成するために信号を直接に又は間接的に記録するステップを更に含むことを特徴とする請求項 6 9 記載の方法。

【請求項 7 1】

画像データを処理する方法において、少なくとも 1 台のカメラにより記録され且つ複数の人物の動きを示す画像データを処理して、各人物がどこを見ているかを判定すると共に、複数の人物がどこを見ているかに基づいてどの人物が話しているかを判定するステップを備えることを特徴とする方法。

【請求項 7 2】

処理される画像データを提供する各カメラの位置と向きを定義するカメラ校正データを使用して画像データを処理することにより、各人物がどこを見ているかを判定することを

10

20

30

40

50

特徴とする請求項 7 1 記載の方法。

【請求項 7 3】

各人物の頭部の位置と向きを三次元で追跡するために画像データを処理することにより、各人物がどこを見ているかを判定することを特徴とする請求項 7 1 又は 7 2 記載の方法。

【請求項 7 4】

各々の人物を見ている人物の数に基づいてどの人物が話をしているかを判定することを特徴とする請求項 7 1 乃至 7 3 のいずれか 1 項に記載の方法。

【請求項 7 5】

各人物が誰を見ているかを定義する値を人物ごとに生成し且つそれらの値を処理して、話をしている人物を判定することを特徴とする請求項 7 4 記載の方法。

10

【請求項 7 6】

話をしている人物が他の人物の大半が見ている人物であることを判定することを特徴とする請求項 7 1 乃至 7 5 のいずれか 1 項に記載の方法。

【請求項 7 7】

プログラム可能処理装置を請求項 3 9 乃至 5 6 の少なくとも 1 項に記載の装置として構成させるための命令を格納することを特徴とする記憶装置。

【請求項 7 8】

プログラム可能処理装置を請求項 5 7 乃至 7 6 の少なくとも 1 項に記載の方法を実行するように動作可能にさせるための命令を格納することを特徴とする記憶装置。

20

【発明の詳細な説明】

【0001】

【発明が属する技術分野】

本発明は、画像データのアーカイピングを補助するためのデータを生成する画像データの処理に関する。

【0002】

更に、本発明は、画像データ及び音声データのアーカイピングを補助するためのデータを生成する画像データ及び音声データの処理に関する。

【0003】

【従来の技術】

データを格納するためのデータベースは数多く存在している。しかし、既存のデータベースには、データベースから情報を検索するためにデータベースを問い合わせる方法が限られているという問題がある。

30

【0004】

【発明が解決しようとする課題】

本発明は、上記の問題に留意してなされた。

【0005】

【課題を解決するための手段】

本発明によれば、複数の人物の動きと話し言葉を記録した画像データ及び音声データを、画像処理と音声処理の組み合わせを利用して処理し、画像データ中に示されているどの人物が話をしているかを識別すると共に、音声データを処理し、識別された話者である参加者に従って選択された処理パラメータを使用して、話された言葉に対応するテキストデータを生成する装置又は方法が提供される。

40

【0006】

データベースからの情報の検索を容易にするために、この後、テキストデータを画像データ及び / 又は音声データと共にデータベースに格納しても良い。

【0007】

また、本発明は、画像データを処理することにより複数の人物の三次元位置を判定し、それらの人物が話した言葉を伝達する音声データを処理して音源の方向を三次元で判定し、生成された位置情報を使用して、それらの言葉の話者を識別し、且つ識別された話者に関

50

して、音声／テキスト変換処理を実行するための音声認識パラメータを選択する装置又は方法を提供する。

【 0 0 0 8 】

このようにして、話者である参加者を容易に識別し、音声データを処理することができる。

【 0 0 0 9 】

各人物の位置は、各人物の少なくとも頭部を追跡するために画像データを処理することにより判定されるのが好ましい。

【 0 0 1 0 】

更に、本発明は、そのようなシステムにおいて画像データ及び音声データを処理し、話者である参加者を識別する装置又は方法を提供する。

10

【 0 0 1 1 】

更に、本発明は、信号の形態をとる命令及び記録された形態をとる命令を含めて、プログラム可能処理装置をそのようなシステムにおいて装置として構成させるため又は方法を実行可能にさせるように構成するための命令を提供する。

【 0 0 1 2 】

本発明によれば、画像データを処理し、どの人物が画像中でその他の人物の注目を集めているかを判定することにより、画像中のどの人物が話しているかを判定し、且つ音声データを処理し、画像データを処理することにより識別された話者である参加者に従って選択される処理パラメータを使用して、その人物により話された言葉に対応するテキストデータを生成する装置又は方法も提供される。

20

【 0 0 1 3 】

また、本発明は、画像データを処理して、画像中の人物が誰を見ているかを判定すると共に、それに基づいてどの人物が話をしているかを判定し、且つ音声データを処理して、話者である参加者について音声認識を実行する装置又は方法を提供する。

【 0 0 1 4 】

このようにして、話者である参加者を容易に識別し、音声データを処理することができる。

【 0 0 1 5 】

更に、本発明は、そのようなシステムにおいて画像データを処理する装置又は方法を提供する。

30

【 0 0 1 6 】

更に、本発明は、信号の形態をとる命令及び記録された形態をとる命令を含めて、プログラム可能処理装置をそのようなシステムにおいて装置として構成させるため又は方法を実行可能にさせるように構成するための命令を提供する。

【 0 0 1 7 】

【発明の実施の形態】

以下、添付の図面を参照して、単なる例示として、本発明の実施形態を説明する。

【 0 0 1 8 】

< 第 1 の実施形態 >

40

図 1 を参照して説明すると、複数のビデオカメラ（図 1 に示す例では 3 台であるが、これとは異なる台数であっても良い）2 - 1、2 - 2、2 - 3 と、マイクロホンアレイ 4 とを使用して、複数の人物 6、8、10、12 の間で行われる会議から画像データと音声データをそれぞれ記録する。

【 0 0 1 9 】

マイクロホンアレイ 4 は、例えば、英国特許第 2 1 4 0 5 5 8 号、米国特許第 4 3 3 3 1 7 0 号及び米国特許第 3 3 9 2 3 9 2 号に記載されているような、入って来る音の方向を判定できるように配列されたマイクロホンのアレイから構成されている。

【 0 0 2 0 】

ビデオカメラ 2 - 1、2 - 2、2 - 3 からの画像データと、マイクロホンアレイ 4 からの

50

音声データは、ケーブル（図示せず）を介してコンピュータ 20 に入力され、コンピュータ 20 は受信したデータを処理し、データベースにデータを格納して、会議のアーカイブ記録を作成する。後に、このデータベースから情報を検索することができる。

【0021】

コンピュータ 20 は、従来のように、表示装置 26 や、この実施形態ではキーボード 28 及びマウス 30 であるユーザ入力装置と共に、1つ又は複数のプロセッサ、メモリ、サウンドカードなどを含む処理装置 24 を有する従来通りのパーソナルコンピュータである。

【0022】

コンピュータ 20 の構成要素と、それらの構成要素に対し入出力されるデータの流れを図 2 に概略的に示す。

10

【0023】

図 2 を参照すると、処理装置 24 は、例えば、ディスク 32 などのデータ記憶媒体に格納されたデータとして及び / 又は例えば、インターネットなどの通信ネットワーク（図示せず）を介する送信又は無線送信により遠隔データベースから処理装置 24 に入力され且つ / 又はキーボード 28 などのユーザ入力装置又は他の入力装置を介してユーザにより処理装置 24 に入力される信号 34 として入力されるプログラミング命令に従って動作するようにプログラムされている。

【0024】

プログラミング命令によりプログラムされると、処理装置 24 は、処理動作を実行するための複数の機能ユニットに有効に構成される。そのような機能ユニットの例とそれらの配線を図 2 に示す。しかし、図 2 に示すユニットと配線は概念的なもので、単に理解を助けるために例示を目的として示されているにすぎない。従って、図 2 のユニットと配線とは、処理装置 24 のプロセッサ、メモリなどが構成される実際のユニットと接続とを必ずしも表してはいない。

20

【0025】

図 2 に示す機能ユニットについて説明すると、中央制御装置 36 はユーザ入力装置 28、30 からの入力を処理すると共に、ユーザによりディスク 38 などの記憶装置に格納されたデータとして、又は処理装置 24 へ送信される信号 40 として処理装置 24 に入力されるデータを受信する。また、中央制御装置 36 はその他の機能ユニットに対して制御と処理を実行する。メモリ 42 は、中央制御装置 36 及びその他の機能ユニットにより使用されるメモリである。

30

【0026】

頭部追跡装置 50 はビデオカメラ 2-1、2-2、2-3 から受信した画像データを処理して、会議のそれぞれの参加者 6、8、10、12 の頭部の位置と向きを三次元で追跡する。この実施形態では、この追跡を実行するために、頭部追跡装置 50 は、後述するように、各々の参加者の頭部の三次元コンピュータモデルを定義するデータと、その特徴を定義するデータとを使用する。これらのデータは頭部モデル記憶装置 52 に格納されている。

【0027】

方向プロセッサ 53 はマイクロホンアレイ 4 から音声データを処理して、マイクロホンにより記録された音が来た方向を判定する。そのような処理は、例えば、英国特許第 2140558 号、米国特許第 4333170 号及び米国特許第 3392392 号に記載されているような従来的一种方式で実行される。

40

【0028】

音声認識プロセッサ 54 はマイクロホンアレイ 4 から受信された音声データを処理して、そこからテキストデータを生成する。すなわち、音声認識プロセッサ 54 は、「Dragon Dictate」又は IBM の「ViaVoice」などの従来の音声認識プログラムに従って動作し、参加者 6、8、10、12 により話された言葉に対応するテキストデータを生成する。音声認識処理を実行するために、音声認識プロセッサ 54 は、音声認識パラメータ記憶装置 56 に格納されている、参加者 6、8、10、12 ごとの音声認識パラメータを定義するデー

50

タを使用する。すなわち、音声認識パラメータ記憶装置 56 に格納されるデータは、音声認識プロセッサを従来の方式で訓練することにより生成される各参加者の音声プロファイルを定義するデータである。例えば、このデータは、訓練後にDragon Dictateの「ユーザファイル」に格納されるデータである。

【0029】

アーカイブプロセッサ 58 は、頭部追跡装置 50、方向プロセッサ 53 及び音声認識プロセッサ 54 から受信したデータを使用して、会議アーカイブデータベース 60 に格納すべきデータを生成する。すなわち、後述するように、カメラ 2-1、2-2 及び 2-3 からの映像データと、マイクロホンアレイ 4 からの音声データとを、音声認識プロセッサ 54 からのテキストデータ及び所定の時点で会議の各参加者が誰を見ていたかを定義するデータと共に会議アーカイブデータベース 60 に格納するのである。

10

【0030】

テキストサーチャ 62 は、中央制御装置 36 と関連して、会議アーカイブデータベース 60 を探索し、後に更に詳細に説明するように、ユーザにより指定される探索基準に適合する会議の 1 つ又は複数の部分を見出し、その部分の音声データ及び映像データを再生するために使用される。

【0031】

表示プロセッサ 64 は、中央制御装置 36 の制御の下に、ユーザに対し表示装置 26 を介して情報を表示すると共に、会議アーカイブデータベース 60 に格納された音声データと映像データを再生する。

20

【0032】

出力プロセッサ 66 はアーカイブデータベース 60 のデータの一部又は全てを、例えば、ディスク 68 などの記憶装置に又は信号 70 として出力する。

【0033】

会議を始める前に、処理装置 24 が必要な処理動作を実行できるようにするために必要なデータを入力することにより、コンピュータ 20 を初期設定する必要がある。

【0034】

図 3 は、この初期設定中に処理装置 24 により実行される処理動作を示す。

【0035】

図 3 を参照して説明すると、ステップ S1 では、中央制御装置 36 は表示プロセッサ 64 に、ユーザが会議に参加するであろう各人物の名前を入力することを要求するメッセージを表示装置 26 に表示させる。

30

【0036】

ステップ S2 では、例えば、ユーザがキーボード 28 を使用して入力した、名前を定義するデータを受信して、中央制御装置 36 は各参加者に独自の識別番号を割り当て、識別番号と参加者の名前との関係を定義するデータ、例えば、図 4 に示すテーブル 80 を会議アーカイブデータベース 60 に格納する。

【0037】

ステップ S3 では、中央制御装置 36 は表示プロセッサ 64 に、会議中のかなり長い時間にわたり人物が見ると考えられ、会議アーカイブデータベース 60 にアーカイブデータを格納することが望まれる物体それぞれの名前をユーザが入力することを要求するメッセージを表示装置 26 に表示させる。そのような物体としては、例えば、図 1 に示すフリップチャート 14 などのフリップチャート、ホワイトボード又は黒板、又はテレビなどが挙げられる。

40

【0038】

ステップ S4 では、例えば、ユーザがキーボード 28 を使用して入力した、物体の名前を定義するデータを受信して、中央制御装置 36 は各物体に独自の識別番号を割り当て、識別番号と物体の名前との関係を定義するデータ、例えば、図 4 に示すテーブル 80 を会議アーカイブデータベース 60 に格納する。

【0039】

50

ステップS 6では、中央制御装置3 6は頭部モデル記憶装置5 2を探索して、会議の参加者ごとに頭部モデルを定義するデータが既に格納されているか否かを判定する。

【0040】

ステップS 6で、1人または複数の参加者について頭部モデルがまだ格納されていないと判定されたならば、ステップS 8で、中央制御装置3 6は表示プロセッサ6 4に、頭部モデルがまだ格納されていない各参加者の頭部モデルを定義するデータをユーザが入力することを要求するメッセージを表示装置2 6に表示させる。

【0041】

これに回答して、ユーザは、例えば、ディスク3 8などの記憶媒体にあるデータを入力するか、又は接続している処理装置から信号4 0としてデータをダウンロードすることにより、必要な頭部モデルを定義するデータを入力する。このような頭部モデルは、従来の方式により、例えば、Valente他の「An Analysis / Synthesis Cooperation for Head Tracking and Video Face Cloning」(Proceedings ECCV ' 9 8 Workshop on Perception of Human Actionに掲載、ドイツ、フライブルク大学、1 9 9 8年6月6日)に記載されている方法で生成されれば良い。

【0042】

ステップS 1 0では、中央制御装置3 6は、ユーザにより入力されたデータを頭部モデル記憶装置5 2に格納する。

【0043】

ステップS 1 2では、中央制御装置3 6及び表示プロセッサ6 4はユーザにより入力された各三次元コンピュータ頭部モデルをレンダリングして、ユーザが各モデルにおいて少なくとも7つの特徴を識別することを要求するメッセージと共に、ユーザに対し表示装置2 6を介してモデルを表示する。

【0044】

これに回答して、ユーザは、各々のモデルの中で、参加者の頭部の正面、側面及び(可能であれば)背面にある顕著な特徴、例えば、目尻、鼻孔、口、耳又は参加者が掛けている眼鏡の特徴などに対応する3 0個の点をマウスを使用して指定する。

【0045】

ステップS 1 4では、中央制御装置3 6は、ユーザにより識別された特徴を定義するデータを頭部モデル記憶装置5 2に格納する。

【0046】

これに対し、ステップS 6で、参加者ごとに頭部モデルが既に頭部モデル記憶装置5 2に格納されていると判定された場合には、ステップS 8からS 1 4を省略する。

【0047】

ステップS 1 6では、中央制御装置3 6は音声認識パラメータ記憶装置を探索して、参加者ごとに音声認識パラメータが既に格納されているか否かを判定する。

【0048】

ステップS 1 6で、全ての参加者については音声認識パラメータを利用できないと判定されたならば、ステップS 1 8で、中央制御装置3 6は表示プロセッサ6 4に、パラメータがまだ格納されていない各参加者について音声認識パラメータを入力することをユーザに要求するメッセージを表示装置2 6に表示させる。

【0049】

これに回答して、ユーザは、例えば、ディスク3 8などの記憶媒体のデータを入力するか、又は遠隔処理装置からの信号4 0として入力することにより、必要な音声認識パラメータを定義するデータを入力する。先に述べた通り、これらのパラメータはユーザの話す音声のプロファイルを定義し、従来の方式で音声認識プロセッサを訓練することにより生成される。従って、例えば、Dragon Dictateを組み込んだ音声認識プロセッサの場合、ユーザにより入力される音声認識パラメータは、Dragon Dictateの「ユーザファイル」に格納されるパラメータに相当する。

【0050】

10

20

30

40

50

ステップS 2 0では、中央制御装置3 6は、ユーザにより入力された音声認識パラメータを定義するデータを音声認識パラメータ記憶装置5 6に格納する。

【0 0 5 1】

これに対し、ステップS 1 6で、参加者ごとに音声認識パラメータを既に利用できる状態にあると判定された場合には、ステップS 1 8からS 2 0を省略する。

【0 0 5 2】

ステップS 2 2では、中央制御装置3 6は表示プロセッサ6 4に、カメラ2 - 1、2 - 2及び2 - 3の校正(キャリブレーション)を可能にするためのステップをユーザが実行することを要求するメッセージを表示装置2 6に表示させる。

【0 0 5 3】

これに回答して、ユーザは必要なステップを実行し、ステップS 2 4では、中央制御装置3 6はカメラ2 - 1、2 - 2及び2 - 3を校正するための処理を実行する。すなわち、この実施形態においては、ユーザにより実行されるステップ及び中央制御装置3 6により実行される処理は、Wiles及びDavisonの「Calibrating and 3D Modelling with a Multi - Camera System」(1 9 9 9 IEEE Workshop on Multi - View Modelling and Analysis of Visual Scenes, ISBN 0 7 6 9 5 0 1 1 0 9)に記載されているような方式で実行される。これは、会議室に対する各カメラ2 - 1, 2 - 2及び2 - 3の位置及び向きを定義する校正データ(キャリブレーションデータ)と、各カメラ固有のパラメータ(横縦比、焦点距離、主点及び一次半径方向ひずみ係数)とを生成する。カメラ校正データ(カメラキャリブレーションデータ)は、例えば、メモリ4 2に格納される。

【0 0 5 4】

ステップS 2 5では、中央制御装置3 6は表示プロセッサ6 4に、ステップS 4で識別データが格納された物体それぞれの位置と向きを判定できるようにするためのステップをユーザが実行することを要求するメッセージを表示装置2 6に表示させる。

【0 0 5 5】

これに回答して、ユーザは必要なステップを実行し、ステップS 2 6では、中央制御装置3 6は、各物体の位置と向きを判定するための処理を実行する。すなわち、この実施形態においては、ユーザは、会議の参加者が見られる物体の面の周囲、例えば、フリップチャート1 4の紙の平面にカラーマーカーを置く。次に、中央制御装置3 6は、カメラ2 - 1、2 - 2及び2 - 3の各々により記録された画像データをステップS 2 4で格納されたカメラ校正データを使用して処理し、従来の方式で、各々のカラーマーカーの三次元位置を判定する。この処理はカメラ2 - 1、2 - 2及び2 - 3ごとに実行されるので、各カラーマーカーの位置は別個に推定され、各カメラ2 - 1、2 - 2及び2 - 3からのデータを使用して計算された位置から、各マーカーの位置について平均位置が判定される。各マーカーの平均位置を使用して、中央制御装置3 6は、従来の方式により、物体面の中心と、物体面の向きを定義するための面垂線とを計算する。物体ごとに判定された位置と向きは、例えば、メモリ4 2に物体校正データとして格納される。

【0 0 5 6】

ステップS 2 7では、中央制御装置3 6は表示プロセッサ6 4に、会議の次の参加者(初めてステップS 2 7を実行する場合には、これは最初の参加者である)が着席することを要求するメッセージを表示装置2 6に表示させる。

【0 0 5 7】

ステップS 2 8では、要求された参加者に着席する時間を与えるために、処理装置2 4は所定の期間待機し、ステップS 3 0では、中央制御装置3 6は各カメラ2 - 1、2 - 2及び2 - 3からのそれぞれの画像データを処理して、カメラごとに、着席した参加者の頭部の推定位置を判定する。すなわち、この実施形態においては、中央制御装置3 6は従来の方式でカメラごとに別個に処理を実行し、参加者の肌の色に対応する色(この色は、頭部モデル記憶装置5 2に格納されている参加者の頭部モデルを定義するデータから判定される)を有する、カメラからの画像データの1つのフレームにおける位置をそれぞれ識別し、次に、(頭部は人体の中で最も高い位置にある肌色の部分であると想定されるので)会

10

20

30

40

50

議室内の最も高い位置に相当する部分を選択する。画像中の識別された部分の位置と、ステップS 2 4で判定されたカメラ校正パラメータとを使用して、中央制御装置3 6は従来の方式により頭部の三次元推定位置を判定する。この処理はカメラ2 - 1、2 - 2及び2 - 3ごとに実行され、カメラごとに別個の推定頭部位置が得られる。

【0058】

ステップS 3 2では、中央制御装置3 6は、カメラ2 - 1、2 - 2及び2 - 3ごとに、参加者の頭部の三次元推定向きを判定する。すなわち、この実施形態においては、中央制御装置3 6は、頭部モデル記憶装置5 2に格納されている参加者の頭部の三次元コンピュータモデルをそのモデルの複数の異なる向きについてレンダリングして、向きごとに対応するモデルの二次元画像を作成する。この実施形態では、参加者の頭部のコンピュータモデルを108の異なる向きでレンダリングするので、108の対応する二次元画像が得られる。これらの向きは、頭部モデルを0°（正面を向いている場合）、+45°（上を向いている場合）及び-45°（下を向いている場合）のそれぞれについて10°ずつ36回回転させた向きに相当する。次に、中央制御装置3 6は、モデルの各々の二次元画像を参加者の頭部を示す、カメラ2 - 1、2 - 2、2 - 3からの映像フレームの部分と比較し、モデルの画像が映像データと最も良く整合する向きを選択する。この比較と選択はカメラごとに実行されるので、カメラごとに推定頭部向きが得られる。頭部モデルをレンダリングすることにより生成される画像データをカメラからの映像データと比較するときには、例えば、Schodl、Haro及びEssaの「Head Tracking Using a Textured Polygonal Model」（Proceedings 1998 Workshop on Perceptual User Interfacesに掲載）に記載されているような従来の技法を使用する。

【0059】

ステップS 3 4では、ステップS 3 0で生成された参加者の頭部のそれぞれの推定位置と、ステップS 5 2で生成された参加者の頭部のそれぞれの推定向きとを頭部追跡装置5 0に入力し、各々のカメラ2 - 1、2 - 2及び2 - 3から受信した画像データのフレームを処理して、参加者の頭部を追跡する。すなわち、この実施形態においては、頭部追跡装置5 0は、例えば、Valente他の「An Analysis / Synthesis Cooperation for Head Tracking and Video Face Cloning」（Proceedings EECV '98 Workshop on Perception of Human Action、ドイツ、フライブルク大学、1998年6月6日）に記載されているような従来の方式で頭部を追跡するための処理を実行する。

【0060】

図5は、ステップS 3 4で頭部追跡装置5 0により実行される処理動作の概要を示す。

【0061】

図5を参照すると、ステップS 4 2 - 1からS 4 2 - n（この実施形態では、カメラは3台であるので、「n」は3である）の各々においては、頭部追跡装置5 0は会議を記録しているカメラのうち対応する1台からの画像データを処理して、そのカメラからの画像データにおける参加者の頭部の特徴（ステップS 1 4で格納された）の位置を判定すると共に、それに基づき、そのカメラからの画像データの現在フレームについて参加者の頭部の三次元位置と向きを判定する。

【0062】

図6は、ステップS 4 2 - 1からS 4 2 - nの所定の1つで実行される処理動作を示す。この処理動作は各ステップで同一であるが、異なるカメラからの画像データに対して実行されることになる。

【0063】

図6を参照すると、ステップS 5 0では、頭部追跡装置5 0は参加者の頭部の現在推定3D位置及び現在推定3D向きを読み取る。初めてステップS 5 0を実行する場合、これらは図3のステップS 3 0及びS 3 2で生成される推定位置と推定向きである。

【0064】

ステップS 5 2では、頭部追跡装置5 0はステップS 2 4で生成されたカメラ校正データを使用して、ステップS 5 0で読み取られた推定位置と推定向きに従って、頭部モデル記

10

20

30

40

50

憶装置 5 2 に格納されている参加者の頭部の三次元コンピュータモデルをレンダリングする。

【 0 0 6 5 】

ステップ S 5 4 では、頭部追跡装置 5 0 は、カメラから受信された映像データの現在フレームについて画像データを処理して、ユーザにより識別され、ステップ S 1 4 で格納された頭部の特徴の中の 1 つの特徴の期待位置を取り囲む各領域からの画像データを取り出す。この期待位置はステップ S 5 0 で読み取られた推定位置及び推定向きと、ステップ S 2 4 で生成されたカメラ校正データとから判定される。

【 0 0 6 6 】

ステップ S 5 6 では、頭部追跡装置 5 0 はステップ S 5 2 で生成された、レンダリングされた画像データと、ステップ S 5 4 で取り出されたカメラ画像データとを整合し、レンダリングされた頭部モデルに最も良く整合するカメラ画像データを見出す。

10

【 0 0 6 7 】

ステップ S 5 8 では、頭部追跡装置 5 0 は、ステップ S 5 6 で識別された、レンダリングされた頭部モデルに最も良く整合するカメラ画像データを、ステップ S 2 4 (図 3) で格納されたカメラ校正データと共に使用して、映像データの現在フレームについて参加者の頭部の 3 D 位置と 3 D 向きを判定する。

【 0 0 6 8 】

再び図 5 に戻ると、ステップ S 4 4 では、頭部追跡装置 5 0 は、ステップ S 4 2 - 1 から S 4 2 - n の各々で識別された、レンダリングされた頭部モデルに最も良く整合するカメラ画像データ (図 6 のステップ S 5 8 で識別される) を使用して、映像データの現在フレームについて参加者の頭部の平均 3 D 位置と平均 3 D 向きを判定する。

20

【 0 0 6 9 】

ステップ S 4 4 を実行するのと同時に、ステップ S 4 6 では、ステップ S 4 2 - 1 から S 4 2 - n の各々で判定されたカメラ画像データにおける頭部の特徴の位置 (図 6 のステップ S 5 8 で識別される) を従来のカルマンフィルタに入力して、映像データの次のフレームについて参加者の頭部の推定 3 D 位置及び推定 3 D 向きを生成する。ビデオカメラ 2 - 1、2 - 2 及び 2 - 3 から映像データのフレームが受信されている間、その参加者についてステップ S 4 2 から S 4 6 を繰り返し実行する。

【 0 0 7 0 】

30

再び図 3 に戻ると、ステップ S 3 6 では、中央制御装置 3 6 は、会議に他の参加者がいるか否かを判定し、参加者ごとに先に説明したように処理が実行され終わるまでステップ S 2 7 から S 3 6 を繰り返す。しかし、参加者ごとにこれらのステップが実行されている間、ステップ S 3 4 では、頭部追跡装置 5 0 は既に着席した各参加者の頭部を追跡し続けている。

【 0 0 7 1 】

ステップ S 3 6 で、会議に他の参加者はなく、従って、各参加者の頭部は頭部追跡装置 5 0 により追跡されていると判定されたならば、ステップ S 3 8 で、中央制御装置 3 6 は、参加者の間で会議を始めて良いことを指示するために、処理装置 2 4 から可聴信号を出力させる。

40

【 0 0 7 2 】

図 7 は、参加者間で会議が行われている間に処理装置 2 4 により実行される処理動作を示す。

【 0 0 7 3 】

図 7 を参照すると、ステップ S 7 0 では、頭部追跡装置 5 0 は会議中の各参加者の頭部を追跡し続ける。ステップ S 7 0 で頭部追跡装置 5 0 により実行される処理は先にステップ S 3 4 に関して説明した処理と同じであるので、ここでは繰り返し説明しない。

【 0 0 7 4 】

頭部追跡装置 5 0 がステップ S 7 0 で各参加者の頭部を追跡しているのと同時に、ステップ S 7 2 では、データを生成し、会議アーカイブデータベース 6 0 に格納するための処理

50

を実行する。

【 0 0 7 5 】

図 8 は、ステップ S 7 2 で実行される処理動作を示す。

【 0 0 7 6 】

図 8 を参照すると、ステップ S 8 0 では、アーカイブプロセッサ 5 8 は、参加者ごとに、その参加者がどの人物又はどの物体を見ているかを定義するいわゆる「視線パラメータ」を生成する。

【 0 0 7 7 】

図 9 は、ステップ S 8 0 で実行される処理動作を示す。

【 0 0 7 8 】

図 9 を参照すると、ステップ S 1 1 0 では、アーカイブプロセッサ 5 8 は頭部追跡装置 5 0 から各参加者の頭部の現在三次元位置を読み取る。これは、ステップ S 4 4 (図 5) で頭部追跡装置 5 0 により実行される処理で生成された平均位置である。

【 0 0 7 9 】

ステップ S 1 1 2 では、アーカイブプロセッサ 5 8 は頭部追跡装置 5 0 から次の参加者 (初めてステップ S 1 1 2 を実行する場合、これは最初の参加者である) の頭部の現在向きを読み取る。ステップ S 1 1 2 で読み取られる向きは、ステップ S 4 4 (図 5) で頭部追跡装置 5 0 により実行される処理で生成された平均向きである。

【 0 0 8 0 】

ステップ S 1 1 4 では、アーカイブプロセッサ 5 8 は、参加者がどこを見ているかを定義する線 (いわゆる「視線」) と、参加者の頭部を他の参加者の頭部の中心と結ぶ概念上の各々の線とが成す角度を判定する。

【 0 0 8 1 】

更に詳細に説明するため、図 1 0 及び図 1 1 を参照すると、1 人の参加者、すなわち、図 1 の参加者 6 についてステップ S 1 1 4 で実行される処理の一例が示されている。図 1 0 を参照すると、ステップ S 1 1 2 で読み取られる参加者の頭部の向きは、参加者の両目の中央の一点から、参加者の頭部に対し垂直に延びる視線 9 0 を定義する。同様に、図 1 1 を参照すると、ステップ S 1 1 0 で読み取られる全ての参加者の頭部の位置は、参加者 6 の両目の中央の点からその他のそれぞれの参加者 8 、 1 0 、 1 2 の頭部の中心に至る概念上の線 9 2 、 9 4 、 9 6 を定義する。ステップ S 1 1 4 で実行される処理においては、アーカイブプロセッサ 5 8 は視線 9 0 と、それぞれの概念上の線 9 2 、 9 4 、 9 6 とが成す角度 9 8 、 1 0 0 、 1 0 2 を判定する。

【 0 0 8 2 】

再び図 9 に戻ると、ステップ S 1 1 6 では、アーカイブプロセッサ 5 8 は、最小値を有する角度 9 8 、 1 0 0 又は 1 0 2 を選択する。すなわち、図 1 1 に示す例でいえば、角度 1 0 0 が選択されることになるであろう。

【 0 0 8 3 】

ステップ S 1 1 8 では、アーカイブプロセッサ 5 8 は、ステップ S 1 1 6 で選択した角度が 1 0 ° より小さい値を有するか否かを判定する。

【 0 0 8 4 】

ステップ S 1 1 8 で、角度が 1 0 ° より小さいと判定されれば、ステップ S 1 2 0 で、アーカイブプロセッサ 5 8 は参加者の視線パラメータを、視線と最小の角度を成す概念上の線により結ばれている参加者の識別番号 (図 3 のステップ S 2 で割り当てられている) に設定する。すなわち、図 1 1 に示す例でいえば、角度 1 0 0 が 1 0 ° より小さければ、角度 1 0 0 は視線 9 0 と、参加者 6 を参加者 1 0 と結ぶ概念上の線 9 4 とが成す角度であるので、視線パラメータは参加者 1 0 の識別番号に設定されるであろう。

【 0 0 8 5 】

これに対し、ステップ S 1 1 8 で、最小の角度が 1 0 ° 以上であると判定された場合には、ステップ S 1 2 2 で、アーカイブプロセッサ 5 8 はステップ S 2 6 (図 3) で先に格納された各物体の位置を読み取る。

10

20

30

40

50

【 0 0 8 6 】

ステップ S 1 2 4 では、アーカイブプロセッサ 5 8 は、参加者の視線 9 0 がいずれかの物体の平面と交わるか否かを判定する。

【 0 0 8 7 】

ステップ S 1 2 4 で、視線 9 0 が 1 つの物体の平面と交わると判定されたならば、ステップ S 1 2 6 で、アーカイブプロセッサ 5 8 は参加者の視線パラメータを視線と交わる物体の識別番号（図 3 のステップ S 4 で割り当てられている）に設定する。視線 9 0 と交わる物体が 2 つ以上ある場合には、これは、視線と交わる物体のうち、参加者に最も近い物体ということになる。

【 0 0 8 8 】

これに対し、ステップ S 1 2 4 で、視線 9 0 が物体の平面と交わらないと判定されたならば、ステップ S 1 2 8 で、アーカイブプロセッサ 5 8 は参加者の視線パラメータを「 0 」に設定する。これは、（視線 9 0 が概念上の線 9 2、9 4、9 6 のいずれにも十分近接していないために）参加者はその他の参加者の誰をも見ておらず、また、（視線 9 0 が物体と交わらないために）どの物体をも見ていないと判定されたことを示している。このような状況は、例えば、参加者が会議室内の、ステップ S 4 でデータが格納されず且つステップ S 2 6 で校正が実行されなかった何らかの物体（例えば、図 1 に示す例において参加者 1 2 が手に持っているメモ）を見ている場合などに起こりうるであろう。

【 0 0 8 9 】

ステップ S 1 3 0 では、アーカイブプロセッサ 5 8 は会議に他の参加者がいるか否かを判定し、参加者ごとに先に説明した処理が実行され終わるまでステップ S 1 1 2 から S 1 3 0 を繰り返す。

【 0 0 9 0 】

再び図 8 に戻ると、ステップ S 8 2 では、中央制御装置 3 6 及び音声認識プロセッサ 5 4 は、映像データの現在フレームに対応する音声データがマイクロホンアレイ 4 から受信されたか否かを判定する。

【 0 0 9 1 】

ステップ S 8 2 で、音声データが受信されたと判定されれば、ステップ S 8 4 で、会議中の参加者のうち誰が話をしているかを判定するための処理を実行する。

【 0 0 9 2 】

図 1 2 は、ステップ S 8 4 で実行される処理動作を示す。

【 0 0 9 3 】

図 1 2 を参照すると、ステップ S 1 4 0 では、方向プロセッサ 5 3 はマイクロホンアレイ 4 からの音声データを処理して、音声が入っている方向を判定する。この処理は、例えば、英国特許第 2 1 4 0 5 5 8 号、米国特許第 4 3 3 3 1 7 0 号及び米国特許第 3 3 9 2 3 9 2 号に記載されているような従来の方式で実行される。

【 0 0 9 4 】

ステップ S 1 4 2 では、アーカイブプロセッサ 5 8 は、画像データの現在フレームについてステップ S 4 4（図 5）で頭部追跡装置 5 0 により判定された各参加者の頭部の位置を読み取り、それに基づいて、どの参加者の頭部がステップ S 1 4 0 で判定された方向、すなわち、音声が入っている方向に対応する位置にあるかを判定する。

【 0 0 9 5 】

ステップ S 1 4 4 では、アーカイブプロセッサ 5 8 は、音声が入っている方向に 2 人以上の参加者がいるか否かを判定する。

【 0 0 9 6 】

ステップ S 1 4 4 で、音声が入っている方向には 1 人しか参加者がいないと判定されれば、ステップ S 1 4 6 で、アーカイブプロセッサ 5 8 は、音声が入っている方向にいる参加者を画像データの現在フレームの話者として選択する。

【 0 0 9 7 】

これに対し、ステップ S 1 4 4 で、音声が入っている方向に対応する位置に 2 人以上の参加

10

20

30

40

50

者の頭部があると判定された場合には、ステップS 1 4 8で、アーカイブプロセッサ5 8は、画像データの直前フレームでそれらの参加者のうち1人が話者として識別されていたか否かを判定する。

【0098】

ステップS 1 4 8で、音声 coming している方向にいる参加者の1人が画像データの直前フレームで話者として選択されていたと判定されれば、ステップS 1 5 0で、アーカイブプロセッサ5 8は画像データの直前フレームで識別されていた話者を画像データの現在フレームについても話者として選択する。これは、画像データの直前フレームの話者が現在フレームの話者と同1人物である確率が高いからである。

【0099】

これに対し、ステップS 1 4 8で、音声 coming している方向にいる参加者がいずれも直前フレームで話者として識別された参加者ではないと判定された場合、又は直前フレームでは話者が識別されなかった場合には、ステップS 1 5 2で、アーカイブプロセッサ5 8は、音声 coming している方向にいるそれぞれの参加者を話者に「なりうる」参加者として選択する。

【0100】

再び図8に戻ると、ステップS 8 6では、アーカイブプロセッサ5 8は話者である参加者ごとの視線パラメータ値、すなわち、ステップS 8 0で生成された、話者である各参加者が誰を又は何を見ているかを定義する視線パラメータ値を後の解析に備えて、例えば、メモリ4 2に格納する。

【0101】

ステップS 8 8では、アーカイブプロセッサ5 8は、ステップS 8 4で判定された話者である各参加者のアイデンティティを音声認識プロセッサ5 4に報知する。これに回答して、音声認識プロセッサ5 4は話者である参加者の音声認識パラメータを音声認識パラメータ記憶装置5 6から選択し、選択されたパラメータを使用して、受信した音声データに対して音声認識処理を実行し、話者である参加者が話した言葉に対応するテキストデータを生成する。

【0102】

他方、ステップS 8 2で、受信した音声データが話し言葉に含まないと判定されたならば、ステップS 8 4からステップS 8 8を省略する。

【0103】

ステップS 8 9では、アーカイブプロセッサ5 8は、会議アーカイブデータベース6 0にどの画像データを格納すべきか、すなわち、どのカメラ2 - 1、2 - 2及び2 - 3からの画像データを格納すべきかを判定する。

【0104】

図1 3は、ステップS 8 9でアーカイブプロセッサ5 8により実行される処理動作を示す。

【0105】

図1 3を参照すると、ステップS 1 6 0では、アーカイブプロセッサ5 8は、画像データの現在フレームについてステップS 8 2（図8）で何らかの話し言葉が検出されたか否かを判定する。

【0106】

ステップS 1 6 0で現在フレームについては話し言葉が存在しないと判定されれば、ステップS 1 6 2で、アーカイブプロセッサ5 8は、画像データを格納すべきカメラとしてデフォルトカメラを選択する。すなわち、この実施形態においては、アーカイブプロセッサ5 8は直前フレームで画像データが記録されたカメラを選択する。処理中の現在フレームが全く初めてのフレームである場合には、アーカイブプロセッサ5 8はカメラ2 - 1、2 - 2、2 - 3のうち1台を無作為に選択する。

【0107】

他方、ステップS 1 6 0で、処理中の現在フレームに話し言葉があると判定された場合には、ステップS 1 6 4で、アーカイブプロセッサ5 8は、次の話者である参加者（初めて

10

20

30

40

50

ステップ S 1 6 4 を実行するときには、これは最初の話者である参加者である) についてステップ S 8 6 で先に格納された視線パラメータを読み取り、その話者である参加者が見ている人物又は物体を判定する。

【 0 1 0 8 】

ステップ S 1 6 6 では、アーカイブプロセッサ 5 8 は、現在考慮されている話者である参加者の頭部の位置と向き (図 5 のステップ S 4 4 で判定された) を、話者である参加者の視線の先にいる参加者の頭部の位置と向き (図 5 のステップ S 4 4 で判定された) 又は話者である参加者の視線の先にある物体の位置と向き (図 3 のステップ S 2 6 で格納された) と共に読み取る。

【 0 1 0 9 】

ステップ S 1 6 8 では、アーカイブプロセッサ 5 8 はステップ S 1 6 6 で読み取られた位置と向きを処理して、カメラ 2 - 1、2 - 2、2 - 3 のうちどのカメラが話者である参加者と、話者である参加者が見ている参加者又は物体の双方を最も良く示しているかを判定し、且つこのカメラを現在フレームの画像データを会議アーカイブデータベース 6 0 に格納すべきカメラとして選択する。

【 0 1 1 0 】

図 1 4 は、ステップ S 1 6 8 でアーカイブプロセッサ 5 8 により実行される処理動作を示す。

【 0 1 1 1 】

図 1 4 を参照すると、ステップ S 1 7 6 では、アーカイブプロセッサ 5 8 は次のカメラ (初めてステップ S 1 7 6 を実行するときには、これは第 1 のカメラである) の三次元位置と視野方向を読み取る。この情報は先に図 3 のステップ S 2 4 で生成、格納されている。

【 0 1 1 2 】

ステップ S 1 7 8 では、アーカイブプロセッサ 5 8 は、ステップ S 1 7 6 で読み取られた情報を、話者である参加者の頭部の三次元位置と向き (図 5 のステップ S 4 4 で判定された) を定義する情報及び話者である参加者が見ている参加者の頭部の三次元位置と向き (図 5 のステップ S 4 4 で判定された) 又は話者である参加者が見ている物体の三次元位置と向き (図 3 のステップ S 2 6 で格納された) を定義する情報と共に使用して、話者である参加者と、話者である参加者が見ている参加者又は物体の双方が現在考慮されているカメラの視野の中に入るか否か (すなわち、現在考慮されているカメラが話者である参加者と、話者である参加者が見ている参加者又は物体の双方を捉えることができるか否か) を判定する。すなわち、この実施形態においては、アーカイブプロセッサ 5 8 は下記の式を評価し、全ての不等式が成立すれば、カメラは話者である参加者と、話者である参加者が見ている参加者又は物体の双方を捉えることができると判定する。

【 0 1 1 3 】

【 数 1 】

$$\left| \arccos \left[\frac{1}{\sqrt{(X_{p1} - X_c)^2 + (Y_{p1} - Y_c)^2}} \begin{pmatrix} X_{p1} - X_c \\ Y_{p1} - Y_c \end{pmatrix} \cdot \begin{pmatrix} dX_c \\ dY_c \end{pmatrix} \right] \right| < \theta_h \quad \dots (1)$$

【 0 1 1 4 】

【 数 2 】

$$\left| \arccos \left[\frac{1}{\sqrt{(X_{p1} - X_c)^2 + (Y_{p1} - Y_c)^2 + (Z_{p1} - Z_c)^2}} \begin{pmatrix} X_{p1} - X_c \\ Y_{p1} - Y_c \\ Z_{p1} - Z_c \end{pmatrix} \cdot \begin{pmatrix} dX_c \\ dY_c \\ dZ_c \end{pmatrix} \right] \right| < \theta_v \quad \dots (2)$$

10

20

30

40

50

【 0 1 1 5 】

【 数 3 】

$$\left| \arccos \left[\frac{1}{\sqrt{(X_{p2} - X_c)^2 + (Y_{p2} - Y_c)^2}} \begin{pmatrix} X_{p2} - X_c \\ Y_{p2} - Y_c \end{pmatrix} \cdot \begin{pmatrix} dX_c \\ dY_c \end{pmatrix} \right] \right| < \theta_h \quad \cdots (3)$$

【 0 1 1 6 】

【 数 4 】

$$\left| \arccos \left[\frac{1}{\sqrt{(X_{p2} - X_c)^2 + (Y_{p2} - Y_c)^2 + (Z_{p2} - Z_c)^2}} \begin{pmatrix} X_{p2} - X_c \\ Y_{p2} - Y_c \\ Z_{p2} - Z_c \end{pmatrix} \cdot \begin{pmatrix} dX_c \\ dY_c \\ dZ_c \end{pmatrix} \right] \right| < \theta_v \quad \cdots (4)$$

10

【 0 1 1 7 】

式中、 (X_c, Y_c, Z_c) は、それぞれ、カメラの主点のx座標、y座標及びz座標（図3のステップS24で先に判定、格納されている）であり、

(dX_c, dY_c, dZ_c) は、それぞれ、x方向、y方向及びz方向のカメラの視野方向（同様に、図3のステップS24で先に判定、格納されている）を表し、

20

θ_h 及び θ_v は、それぞれ、水平方向と垂直方向のカメラの角視野（同様に、図3のステップS24で判定、格納されている）であり、

(X_{p1}, Y_{p1}, Z_{p1}) は、それぞれ、話者である参加者の頭部の中心のx座標、y座標及びz座標（図5のステップS44で判定されている）であり、

$(dX_{p1}, dY_{p1}, dZ_{p1})$ は、それぞれ、話者である参加者の視線90の向き（同様に、図5のステップS44で判定されている）を表し、

(X_{p2}, Y_{p2}, Z_{p2}) は、それぞれ、話者である参加者が見ている参加者の頭部の中心のx座標、y座標及びz座標（図5のステップS44で判定されている）又は話者である参加者が見ている物体の面の中心のx座標、y座標及びz座標（図3のステップS26で判定されている）であり、

30

$(dX_{p2}, dY_{p2}, dZ_{p2})$ は、それぞれ、話者である参加者が見ている参加者の視線90のx方向、y方向及びz方向の方向（同様に、図5のステップS44で判定されている）又は話者である参加者が見ている物体面に対する垂線のx方向、y方向及びz方向の方向（図3のステップS26で判定されている）を表す。

【 0 1 1 8 】

ステップS178で、カメラが話者である参加者と、話者である参加者が見ている参加者又は物体の双方を捉えることができる（すなわち、上記の式（1）、（2）、（3）及び（4）の不等式が成立する）と判定されれば、ステップS180で、アーカイブプロセッサ58は、現在考慮されているカメラが話者である参加者を捉えている視野の画質を表す値を計算し、格納する。すなわち、この実施形態においては、アーカイブプロセッサ58

40

【 0 1 1 9 】

【 数 5 】

$$Q1 = \frac{1}{\sqrt{(X_c - X_{p1})^2 + (Y_c - Y_{p1})^2 + (Z_c - Z_{p1})^2}} \begin{pmatrix} X_c - X_{p1} \\ Y_c - Y_{p1} \\ Z_c - Z_{p1} \end{pmatrix} \cdot \begin{pmatrix} dX_{p1} \\ dY_{p1} \\ dZ_{p1} \end{pmatrix} \quad \cdots (5)$$

【 0 1 2 0 】

50

式中、各項の定義は先の式(1)及び(2)に関して挙げた定義と同じである。

【0121】

ステップS180で計算される画質値Q1は、-1から+1の値をとるスカラーであり、話者である参加者の頭部の背面がカメラに直接向いている場合、その値は-1であり、話者である参加者の顔面が直接カメラに向いている場合には+1である。話者である参加者の頭部がその他の向きである場合には、-1と+1の間の値をとる。

【0122】

ステップS182では、アーカイブプロセッサ58は、現在考慮されているカメラが話者である参加者が見ている参加者又は物体を捉えている視野の画質を表す値を計算し、格納する。すなわち、この実施形態においては、アーカイブプロセッサ58は下記の式を使用して、画質値Q2を計算する。

10

【0123】

【数6】

$$Q2 = \frac{1}{\sqrt{(X_c - X_{p2})^2 + (Y_c - Y_{p2})^2 + (Z_c - Z_{p2})^2}} \begin{pmatrix} X_c - X_{p2} \\ Y_c - Y_{p2} \\ Z_c - Z_{p2} \end{pmatrix} \cdot \begin{pmatrix} dX_{p2} \\ dY_{p2} \\ dZ_{p2} \end{pmatrix} \quad \dots (6)$$

【0124】

式中、パラメータの定義は先の式(3)及び(4)に関して挙げた定義と同じである。

20

【0125】

Q2も、同様に、参加者の頭部の背面又は物体の面の背面が直接カメラに向いている場合に-1、参加者の顔面又は物体の正面が直接カメラに向いている場合には+1の値をとるスカラーである。参加者の頭部又は物体の面がその他の向きである場合には、-1と+1の間の値をとる。

【0126】

ステップS184では、アーカイブプロセッサ58はステップS180で計算した画質値Q1と、ステップS182で計算した画質値Q2とを比較し、最小値を選択する。この最小値は、カメラが話者である参加者、あるいは話者である参加者が見ている参加者又は物体を捉えている「最悪の視野」を示す(Q1がQ2より小さい場合、最悪の視野は話者である参加者の視野であり、Q2がQ1より小さい場合は、最悪の視野は話者である参加者が見ている参加者又は物体の視野である)。

30

【0127】

他方、ステップS178で、式(1)、(2)、(3)及び(4)の不等式の1つ又は2つ以上が成立しない(すなわち、カメラが話者である参加者と、話者である参加者が見ている参加者又は物体の双方を捉えることができない)と判定された場合には、ステップS180からS184を省略する。

【0128】

ステップS186では、アーカイブプロセッサ58は、画像データを提供していたカメラが他にあるか否かを判定する。カメラごとに上記の処理が実行され終わるまで、ステップS176からS186を繰り返す。

40

【0129】

ステップS188では、アーカイブプロセッサ58は、ステップS184で処理を実行したときにカメラごとに格納された「最悪の視野」の値(すなわち、ステップS184でカメラごとに格納されたQ1又はQ2の値)を比較し、格納されているそれらの値の中で最大の値を選択する。この最大値は「最良の最悪の視野」を表し、そこで、ステップS188で、アーカイブプロセッサ58は、ステップS184でこの「最良の最悪の視野」値が格納されていたカメラを会議アーカイブデータベースに画像データを格納すべきカメラとして選択する。これは、このカメラが話者である参加者と、話者である参加者が見ている参加者又は物体の双方を最良の視野で捉えているからである。

50

【0130】

ステップS170では、アーカイブプロセッサ58は、話者に「なりうる」参加者を含めて、他に話者である参加者がいるか否かを判定する。話者である参加者ごとに、また、話者に「なりうる」参加者ごとに上記の処理が実行され終わるまで、ステップS164からS170を繰り返す。

【0131】

再び図8に戻ると、ステップS90では、アーカイブプロセッサ58はステップS89で選択したカメラから受信された映像データの現在フレームと、マイクロホンアレイ4から受信された音声データとを従来の方式によりMPEG2データとして符号化し、符号化されたデータを会議アーカイブデータベース60に格納する。

10

【0132】

図15は、会議アーカイブデータベース60へのデータの格納を概略的に示す。図15に示す格納構造は概念上のものであり、格納される情報間のリンクを示すことにより理解を助けることを目的としている。従って、これは、会議アーカイブデータベース60を構成するメモリにデータが厳密にどのように格納されるかを必ずしも表してはいない。

【0133】

図15を参照すると、会議アーカイブデータベース60は水平軸200により表される時間情報を格納している。水平軸200上の各単位は所定の量の時間、例えば、カメラから受信される映像データの1フレーム分の周期を表す。(会議アーカイブデータベース60が一般には図15に示す数より多くの数の時間単位を含むことは言うまでもなく了解されるであろう。)ステップS90で生成されたMPEG2データは、タイミング情報(このタイミング情報は図15では水平軸200に沿ったMPEG2データ202の位置により概略的に表されている)と共に、データ202として会議アーカイブデータベース60に格納されている。

20

【0134】

再び図8に戻ると、ステップS92では、アーカイブプロセッサ58は、現在フレームについてステップS88で音声認識プロセッサ54により生成されたテキストデータを会議アーカイブデータベース60に格納する(図15には204で示す)。すなわち、テキストデータは対応するMPEG2データへのリンクを伴って格納される。図15においては、このリンクは、テキストデータがMPEG2データと同じ縦列に格納されることによって表されている。話をしていない参加者からは格納すべきテキストデータが得られないことは理解されるであろう。図15に示す例では、参加者1については初めの10個のタイムスロットにテキストが格納され(206で示す)、参加者3については12番目から20番目のタイムスロットに格納され(208で示す)、参加者4については21番目のタイムスロットに格納されている(210で示す)。この例では、参加者2は図15に示すタイムスロットの間は話をしていないので、参加者2のテキストは格納されていない。

30

【0135】

ステップS94では、アーカイブプロセッサ58は、ステップS80で現在フレームについて参加者ごとに生成された視線パラメータ値を会議アーカイブデータベース60に格納する(図15には212で示す)。図15を参照すると、視線パラメータ値は、参加者ごとに、関連するMPEG2データ202及び関連するテキストデータ204へのリンクと共に格納されている(このリンクは、図15では、視線パラメータ値が関連するMPEG2データ202及び関連するテキストデータ204と同じ縦列にあることによって表されている)。従って、一例として、図15の第1のタイムスロットに関していえば、参加者1の視線パラメータ値は、参加者1が参加者3を見ていることを指示する3であり、参加者2の視線パラメータ値は、参加者2がフリップチャート14を見ていることを指示する5であり、参加者3の視線パラメータ値は、参加者3が参加者1を見ていることを指示する1であり、参加者4の視線パラメータ値は、参加者4が他の参加者の誰も見ていないことを指示する「0」である(図1に示す例では、12で示される参加者は他の参加者ではなく、自分のメモを見ている)。

40

50

【 0 1 3 6 】

ステップ S 9 6 では、中央制御装置 3 6 及びアーカイブプロセッサ 5 8 は、会議の参加者の 1 人が話を止めたか否かを判定する。この実施形態においては、この検査は、所定の参加者のテキストデータが直前のタイムスロットには存在したが、現在タイムスロットには存在しないことを判定するためにテキストデータ 2 0 4 を検査することにより実行される。いずれかの参加者についてこの条件が満たされれば（すなわち、参加者が話を止めていれば）、ステップ S 9 8 で、アーカイブプロセッサ 5 8 は、話を止めた参加者ごとに、先にステップ S 8 6 を実行したときに格納されていた視線パラメータ値を処理し（それらの視線パラメータ値は、この時点で止まった話をしていた期間中にその参加者が誰を又は何を見ていたかを定義する）、視線ヒストグラムを定義するデータを生成する。すなわち、参加者が話をしていた期間の視線パラメータ値を処理して、その期間中に話者である参加者がその他の参加者及び物体のそれぞれを見ていた時間の割合（％）を定義するデータを生成するのである。

10

【 0 1 3 7 】

図 1 6 A 及び図 1 6 B は、図 1 5 のテキストと 2 0 6 及び 2 0 8 の期間にそれぞれ相当する視線ヒストグラムを示す。

【 0 1 3 8 】

図 1 5 及び図 1 6 A を参照すると、参加者 1 が話していた期間 2 0 6 の間、参加者 1 は、図 1 6 に 3 0 0 で示すように、1 0 個のタイムスロットのうち 6 個のタイムスロット（すなわち、参加者が話をしていた期間全体の長さの 6 0 ％）にわたり参加者 3 を見ており、図 1 6 A に 3 1 0 で示すように、1 0 個のタイムスロットのうち 4 個のタイムスロット（すなわち、時間の 4 0 ％）にわたり参加者 4 を見ていた。

20

【 0 1 3 9 】

同様に、図 1 5 及び図 1 6 B を参照すると、期間 2 0 8 の間、参加者 3 は、図 1 6 B に 3 2 0 で示すように、時間の約 4 5 ％にわたり参加者 1 を見ており、図 1 6 B に 3 3 0 で示すように、時間の約 3 3 ％にわたり物体 5（すなわち、フリップチャート 1 4）を見ており、図 1 6 B に 3 4 0 で示すように、時間の約 2 2 ％にわたり参加者 2 を見ていた。

【 0 1 4 0 】

再び図 8 に戻ると、ステップ S 1 0 0 では、ステップ S 9 8 で生成した各視線ヒストグラムを、それを生成する元になったテキストの関連する期間とリンクさせて、会議アーカイブデータベース 6 0 に格納する。図 1 5 を参照すると、格納される視線ヒストグラムは 2 1 4 で示され、テキスト期間 2 0 6 のヒストグラムを定義するデータは 2 1 6 で示され、テキスト期間 2 0 8 のヒストグラムを定義するデータは 2 1 8 で示されている。図 1 5 においては、視線ヒストグラムと関連するテキストとのリンクは、視線ヒストグラムがテキストデータと同じ縦列に格納されることにより表されている。

30

【 0 1 4 1 】

他方、ステップ S 9 6 で、現在時限について参加者の 1 人が話を止めていないと判定された場合には、ステップ S 9 8 及び S 1 0 0 を省略する。

【 0 1 4 2 】

ステップ S 1 0 2 では、アーカイブプロセッサ 5 8 は、映像フレームの直前フレーム（すなわち、ステップ S 8 0 から S 1 0 0 でデータが生成、格納されたばかりのフレームの直前のフレーム）及び他の先行フレームについて、会議アーカイブデータベース 6 0 に格納されているデータを必要に応じて修正する。

40

【 0 1 4 3 】

図 1 7 は、ステップ S 1 0 2 でアーカイブプロセッサ 5 8 により実行される処理動作を示す。

【 0 1 4 4 】

図 1 7 を参照すると、ステップ S 1 9 0 では、アーカイブプロセッサ 5 8 は、次の先行フレーム（初めてステップ S 1 9 0 を実行する場合には、これは現在フレームの直前のフレームであり、すなわち、現在フレームが「i」番目のフレームであれば、「i - 1」番目

50

のフレームである)について、話者に「なりうる」参加者のデータを会議アーカイブデータベース60に格納するかどうかを判定する。

【0145】

ステップS190で、考慮されている先行フレームについて話者に「なりうる」参加者のデータが格納されていないと判定されれば、会議アーカイブデータベース60のデータを修正する必要はない。

【0146】

他方、ステップS190で、考慮されている先行フレームについて話者に「なりうる」参加者のデータが格納されていると判定された場合には、ステップS192で、アーカイブプロセッサ58は、先行フレームについてデータが格納された話者に「なりうる」参加者の1人が現在フレームについて識別された話者である参加者(話者に「なりうる」参加者ではない)、すなわち、図12のステップS146で識別された話者である参加者と同1人物であるか否かを判定する。

10

【0147】

ステップS192で、先行フレームの話者に「なりうる」参加者がいずれも現在フレームについてステップS146で識別された話者である参加者と同じではないと判定されれば、考慮されている先行フレームについて会議アーカイブデータベース60に格納されているデータの修正を実行しない。

【0148】

他方、ステップS192で、先行フレームの話者に「なりうる」参加者が現在フレームについてステップS146で識別された話者である参加者と同1人物であると判定された場合には、ステップS194で、アーカイブプロセッサ58は、現在フレームの話者である参加者と同じではない話者に「なりうる」参加者のそれぞれについて、考慮されている先行フレームのテキストデータ204を会議アーカイブデータベース60から削除する。

20

【0149】

以上説明したようにステップS190、S192及びS194の処理を実行することにより、現在フレームについて画像データ及び音声データを処理することによって話者が明確に識別された場合、現在フレームの話者は先行フレームの話者と同1人物であるという仮定を利用して、話者に「なりうる」参加者(すなわち、曖昧さなく話者を識別することが不可能であったため)について格納された直前フレームのデータを更新するのである。

30

【0150】

ステップS194を実行した後、次の先行フレームについてステップS190からS194を繰り返す。すなわち、現在フレームが「i」番目のフレームであれば、初めてステップS190からS194を実行するときに「i-1」番目のフレームを考慮し、2度目にステップS190からS194を実行するときには「i-2」番目のフレームを考慮する。これ以降も同様である。ステップS190で、考慮されている先行フレームについて話者に「なりうる」参加者のデータが格納されていないと判定されるか、またはステップS192で、考慮されている先行フレームにおける話者に「なりうる」参加者がいずれも現在フレームについて曖昧さなく識別された話者である参加者と同じではないと判定されるまで、ステップS190からS194を繰り返し実行し続ける。このようにして、いくつかの連続するフレームにわたり話者に「なりうる」参加者が識別された場合には、次のフレームで話者に「なりうる」参加者の中から実際の話者である参加者が識別されれば、会議アーカイブデータベースに格納されているデータを修正する。

40

【0151】

再び図8に戻ると、ステップS104では、中央制御装置36は、カメラ2-1、2-2、2-3から映像データの別のフレームが受信されたか否かを判定する。カメラ2-1、2-2、2-3から画像データが受信されている間は、ステップS80からS104を繰り返し実行する。

【0152】

会議アーカイブデータベース60にデータが格納されている場合、会議アーカイブデータ

50

ベース 60 を問い合わせ、会議に関連するデータを検索しても良い。

【 0 1 5 3 】

図 18 は、ユーザにより指定される探索基準を満たす会議の各部分に関連するデータを検索する目的で会議アーカイブデータベース 60 を探索するために実行される処理動作を示す。

【 0 1 5 4 】

図 18 を参照すると、ステップ S 200 では、中央制御装置 36 は表示プロセッサ 64 に、要求する会議アーカイブデータベース 60 の探索を定義する情報をユーザが入力することを求めるメッセージを表示装置 26 に表示させる。すなわち、この実施形態においては、中央制御装置 36 は図 19A に示す表示を表示装置 26 に表示させる。

10

【 0 1 5 5 】

図 19A を参照すると、ユーザは、会議アーカイブデータベース 60 の中で見出すことを臨む会議の部分の情報を入力することを求められる。すなわち、この実施形態においては、ユーザは、話をしていた参加者を定義する情報 400、情報 400 の中で識別される参加者が口に出した 1 つ又は複数のキーワードから成る情報 410、及び情報 400 の中で識別される参加者が話している間に見ていた参加者又は物体を定義する情報 420 を入力することを求められる。更に、ユーザは、探索を実行すべき会議の部分の定義する時間情報を入力することができる。すなわち、ユーザは、その時間を越えたら探索を中止すべきである会議中の時間（すなわち、指定される時間に至るまでの会議の期間を探索すべきである）を定義する情報 430 と、その時間から探索を実行すべきである会議中の時間を定義する情報 440 と、探索を実行すべき期間の開始時間と終了時間をそれぞれ定義する情報 450 及び 460 とを入力できる。この実施形態においては、情報 430、440、450 及び 460 は、例えば、分単位で絶対期限として時間を指定するか、又は会議時間全体に占める割合を指示する小数値を入力することにより相対期限で時間を指定することにより入力されれば良い。例えば、情報 430 として値 0.25 を入力した場合、探索は会議の初めの四分の一に限定されるであろう。

20

【 0 1 5 6 】

この実施形態では、ユーザは 1 回の探索で情報 400、410 及び 420 の全てを入力する必要はなく、この情報のうち 1 つ又は 2 つを省いても良い。ユーザが情報 400、410 及び 420 の全てを入力すれば、会議の中で、情報 400 の中で識別される参加者が情報 420 の中で識別される参加者又は物体に向かって話していた部分及び情報 400 の中で識別される参加者が情報 410 の中で定義されるキーワードを話した部分をそれぞれ識別するための探索が実行されることになる。これに対し、情報 410 を省いた場合には、会議の中で、情報 400 の中で識別される参加者が何を言ったかに関わらず、情報 420 の中で定義される参加者又は物体に向かって話していた部分をそれぞれ識別するための探索が実行される。情報 410 及び 420 を省いた場合には、会議の中で、情報 400 の中で識別される参加者が何を誰に向かって話したかに関わらず、話していた部分をそれぞれ識別するための探索が実行される。情報 400 を省いた場合には、会議の中で、いずれかの参加者が情報 420 の中で定義される参加者又は物体を見ながら情報 410 の中で定義されるキーワードを話した部分をそれぞれ識別するための探索が実行される。情報 400 及び 410 を省いた場合には、会議の中で、いずれかの参加者が情報 420 の中で定義される参加者又は物体に向かって話した部分をそれぞれ識別するための探索が実行される。情報 420 を省いた場合には、会議の中で、情報 400 の中で定義される参加者が情報 410 の中で定義されるキーワードを誰に向かって話したかに関わらず、キーワードを話した部分をそれぞれ識別するための探索が実行される。同様に、情報 400 及び 420 を省いた場合には、会議の中で、誰が誰に向かって言ったかに関わらず、情報 410 の中で識別されるキーワードが話された部分をそれぞれ識別するための探索が実行される。

30

40

【 0 1 5 7 】

更に、ユーザは時間情報 430、440、450 及び 460 の全てを入力しても良いし、あるいはそのうちいくつかを省いても良い。

50

【 0 1 5 8 】

また、探索者が言葉の組み合わせ又はある言葉に代わる言葉を探索できるようにするために、情報 4 1 0 の中で入力されるキーワードと組み合わせで周知のブール演算子及び探索アルゴリズムを使用しても良い。

【 0 1 5 9 】

探索を定義するためにユーザが必要な全ての情報を入力したならば、マウス 3 0 などのユーザ入力装置を使用して領域 4 7 0 をクリックすることにより探索を開始する。

【 0 1 6 0 】

再び図 1 8 に戻ると、ステップ S 2 0 2 では、ユーザが入力した探索情報を中央制御装置 3 6 により読み取り、命令された探索を実行する。すなわち、この実施形態においては、中央制御装置 3 6 は情報 4 0 0 又は 4 2 0 の中で入力された参加者又は物体の名前をテーブル 8 0 (図 4) を使用して識別番号に変換し、情報 4 0 0 で定義される参加者 (情報 4 0 0 が入力されなかった場合は全ての参加者) についてテキスト情報 2 0 4 を考慮する。ユーザにより情報 4 2 0 が入力されていれば、テキストの期間ごとに、中央制御装置 3 6 は対応する視線ヒストグラムを定義するデータを検査して、情報 4 2 0 の中で定義される参加者又は物体のヒストグラムにおける注目時間の割合がこの実施形態では 2 5 % である閾値以上であるか否かを判定する。このように、話者である参加者が話をしている時間の少なくとも 2 5 % にわたって情報 4 2 0 の中で定義される参加者又は物体を見ていれば、話者である参加者が話しの間に他の参加者又は物体を見たとしても、話し言葉 (テキスト) の各期間を考慮して、情報 4 0 0 の中で定義される参加者は情報 4 2 0 の中で定義される参加者又は物体に話しかけていたという基準を満たす。従って、例えば、情報 4 2 0 の中で 2 人以上の参加者が識別されていれば、視線ヒストグラムの値が 2 人以上の参加者について 2 5 % 以上であるような話の期間が識別されるであろう。ユーザが情報 4 1 0 を入力した場合、中央制御装置 3 6 及びテキストサーチ 6 2 は、先に情報 4 0 0 及び 4 2 0 に基づいて識別されたテキストの各部分 (情報 4 0 0 及び 4 2 0 が入力されていなければ、テキストの全ての部分) を探索して、情報 4 1 0 の中で定義されるキーワードを含む各部文を識別する。ユーザが時間情報を入力していれば、上記の探索はそれらの期限により定義される会議時間に限られる。

【 0 1 6 1 】

ステップ S 2 0 4 では、中央制御装置 3 6 は表示プロセッサ 6 4 に、探索中に識別された関連話題のリストを表示装置 2 6 を介してユーザに対し表示させる。すなわち、中央制御装置 3 6 は、図 1 9 B に示すような情報をユーザに対し表示させる。図 1 9 B を参照すると、探索パラメータを満足させる各々の話題のリストが作成され、その話の開始時間を定義する情報が絶対期限と、会議時間全体に占める割合の双方で表示されている。そこで、ユーザは、例えば、リスト中の必要な話題をマウス 3 0 を使用してクリックすることにより、話題の 1 つを選択して、再生させることができる。

【 0 1 6 2 】

ステップ S 2 0 6 では、中央制御装置 3 6 はステップ S 2 0 4 でユーザが行った選択を読み取り、格納されている会議の関連部分の MPEG 2 データ 2 0 2 を会議アーカイブデータベース 6 0 から再生させる。すなわち、中央制御装置 3 6 及び表示プロセッサ 6 4 は MPEG 2 データ 2 0 2 を復号し、画像データと音声を表示装置 2 6 を介して出力するのである。再生すべき話の一部又は全てについて 2 台以上のカメラからの画像データが格納されている場合には、そのことを表示装置 2 6 によりユーザに指示し、ユーザは、例えば、キーボード 2 8 を使用して中央制御装置 3 6 に命令を入力することにより、再生すべき画像データを選択することができる。

【 0 1 6 3 】

ステップ S 2 0 8 では、中央制御装置 3 6 は、ユーザが会議アーカイブデータベース 6 0 の問い合わせを中止することを望むか否かを判定し、望まないのであれば、ステップ S 2 0 0 から S 2 0 8 を繰り返す。

【 0 1 6 4 】

以上説明した本発明の実施形態に対しては、様々な変形や変更を実施することができる。

【0165】

上記の実施形態では、ステップS34（図3）及びステップS70（図7）においては、会議中の各参加者の頭部を追跡していた。しかし、これに加えて、ステップS4及びS26でデータを格納した物体が移動する場合（そのような物体としては、例えば、参加者により回覧されるようなメモ又は参加者間で手渡されるべき物体などが考えられる）それらの物体を追跡することも可能であろう。

【0166】

上記の実施形態では、複数台のビデオカメラ2-1、2-2、2-3からの画像データを処理していた。しかし、その代わりに、1台のビデオカメラからの画像データをも良い。この場合、例えば、ステップS42-1（図5）のみを実行し、ステップS42-2からS42-nを省略する。同様に、ステップS44を省略し、ステップS42-1で実行される処理の間、画像データの現在フレームに関わる参加者の頭部の3D位置及び向きをステップS58（図6）で判定される3D位置及び向きであるとみなす。ステップS46では、カルマンフィルタに入力される頭部の特徴の位置はその1台のカメラからの画像データにおける位置になるであろう。更に、会議アーカイブデータベース60に画像データを記録すべきカメラを選択するためのステップS89（図8）も省略されるであろう。

【0167】

上記の実施形態では、ステップS168（図13）において、話者である参加者と、話者である参加者が見ている参加者又は物体とを最も良く捉えるカメラを識別するための処理を実行していた、しかし、上記の実施形態において説明したようにカメラを識別する代わりに、処理装置24の初期設定中に、会議テーブルを囲む2つずつの着席位置を最も良く捉え且つ／又は各々の着席位置と所定の物体（フリップチャート14など）を最も良く捉えるのはカメラ2-1、2-2、2-3のうちどれであるかをユーザが定義することも可能である。このようにして、話者である参加者と、話者である参加者が見ている参加者があらかじめ定義された着席位置にいと判定されれば、それらのあらかじめ定義された着席位置を最も良く捉えるとユーザにより定義されたカメラを画像データを格納すべきカメラとして選択することができる。同様に、話者である参加者があらかじめ定義された位置にあり且つある物体を見ている場合、そのあらかじめ定義された着席位置と物体を最も良く捉えるとユーザにより定義されたカメラを画像データを格納すべきカメラとして選択することができる。

【0168】

上記の実施形態では、ステップS162（図13）において、直前フレームで画像データが格納されたカメラとしてデフォルトカメラを選択していた。しかし、その代わりに、例えば、処理装置24の初期設定中に、ユーザがデフォルトカメラを選択しても良い。

【0169】

上記の実施形態では、ステップS194（図17）において、その時点で実際には話者である参加者として識別されなかった話者に「なりうる」参加者について、テキストデータ204を会議アーカイブデータベース60から削除していた。しかし、これに加えて、関連する視線ヒストグラムデータ214も共に削除して良い。更に、カメラ2-1、2-2、2-3のうち2台以上からのMPEG2データ202が格納されていた場合、話者に「なりうる」参加者に関連するMPEG2データも削除して良い。

【0170】

上記の実施形態では、話者である参加者を一意性をもって識別することが不可能である場合、話者に「なりうる」参加者を定義し、話者になりうる参加者についてデータを処理して会議アーカイブデータベース60に格納し、その後、会議アーカイブデータベース60に格納されたデータを修正していた（図8のステップS102）。しかし、話者に「なりうる」参加者についてデータを処理し、格納するのではなく、カメラ2-1、2-2及び2-3から受信した映像データと、マイクロホンアレイ4から受信した音声データとを、

後続フレームに関連するデータから話者である参加者が識別されたときの後の処理及びアーカイピングに備えて格納しておいても良い。あるいは、ステップS 1 4 4 (図 1 2) で実行された処理の結果、音声が入っている方向に2人以上の参加者がいることが指示された場合には、カメラ2 - 1、2 - 2及び2 - 3からの画像データを処理して、参加者の唇の動きを検出すると共に、音声が入っている方向にいて、唇が動いている参加者を話者である参加者として選択しても良い。

【0 1 7 1】

上記の実施形態では、各人物の頭部の位置と、各人物の頭部の向きと、各人物が誰を又は何を見ているかを定義する人物ごとの視線パラメータとを判定するための処理を実行していた。その後、画像データのフレームごとに、各人物の視線パラメータ値を会議アーカイブデータベース60に格納する。しかし、全ての人物について視線パラメータを判定する必要はない。例えば、話者である参加者のみの視線パラメータを判定し、画像データのフレームごとにこの視線パラメータ値のみを会議アーカイブデータベース60に格納することは可能である。従って、この場合、話者である参加者の頭部の位置を判定するだけで良いであろう。このようにすれば、処理及び格納に課される負担を軽減することができる。

10

【0 1 7 2】

上記の実施形態では、ステップS 2 0 2 (図 1 8)において、テキストの特定の部分の視線ヒストグラムを考慮し、その視線ヒストグラムにおいて別の参加者又は物体に注目している時間の割合が所定の閾値以上である場合に、参加者は別の参加者と話していた又は別の物体を見ていたと判定していた。しかし、閾値を使用する代わりに、テキスト(話)の期間中、話者である参加者が見ていた参加者又は物体は、視線ヒストグラムの中で最も大きな割合の注目時間を有する参加者又は物体(例えば、図 1 6 Aの参加者3及び図 1 6 Bの参加者1)であると定義しても良い。

20

【0 1 7 3】

上記の実施形態では、カメラ2 - 1、2 - 2及び2 - 3と、マイクロホンアレイ4とからデータが受信されている間、MPEG2データ202、テキストデータ204、視線パラメータ212及び視線ヒストグラム214をリアルタイムで会議アーカイブデータベース60に格納していた。しかし、その代わりに、映像データと音声データを格納し、リアルタイムではなくデータ202、204、2120及び214を生成して、会議アーカイブデータベース60に格納しても良い。

30

【0 1 7 4】

上記の実施形態では、会議の定義された部分についてデータを検索するために会議アーカイブデータベース60を問い合わせる前に、MPEG2データ202、テキストデータ204、視線パラメータ212及び視線ヒストグラム214を生成し、データベースに格納していた。しかし、探索の要求に先立ってデータを生成、格納するのではなく、会議アーカイブデータベース60の探索がユーザにより要求されるのに応答して、既に会議アーカイブデータベース60に格納されているデータを処理することにより、視線ヒストグラムデータ214の一部又は全てを生成しても良い。例えば、上記の実施形態では、視線ヒストグラム214はステップS 9 8及びS 1 0 0 (図 8)でリアルタイムで計算、格納されていたが、ユーザにより入力される探索要求に応答してそれらのヒストグラムを計算しても良い。

40

【0 1 7 5】

上記の実施形態では、テキストデータ204は会議アーカイブデータベース60に格納されていた。テキストデータ204の代わりに、音声データを会議アーカイブデータベース60に格納しても良い。その場合、格納された音声データ自体を音声認識処理を利用してキーワードを求めて探索しても良いし、あるいは音声認識処理を利用して音声データをテキストに変換し、従来のテキストサーチャを使用してそのテキストを探索しても良い。

【0 1 7 6】

上記の実施形態では、処理装置24は、アーカイブすべきデータを受信し、生成するための機能構成要素(例えば、中央制御装置36、頭部追跡装置50、頭部モデル記憶装置5

50

2、方向プロセッサ53、音声認識プロセッサ54、音声認識パラメータ記憶装置56及びアーカイブプロセッサ58)と、アーカイブデータを格納するための機能構成要素(例えば、会議アーカイブデータベース60)と、データベースを探索し、そこから情報を検索するための機能構成要素(例えば、中央制御装置36及びテキストサーチ62)とを含む。しかし、これらの機能構成要素を別個の装置に設けても良い。例えば、アーカイブすべきデータを生成する1つ又は複数の装置と、データベースを探索する1つ又は複数の装置とをインターネットなどのネットワークを介して1つ又は複数のデータベースに接続しても良い。

【0177】

また、図20を参照して説明すると、1つ又は複数の場所での会議500、510、520から得られた映像データと音声データをデータ処理・データベース記憶装置530(アーカイブデータを生成し、格納するための機能構成要素を具備する)に入力し、データベースを問い合わせ、そこから情報を検索するために、1つ又は複数のデータベース問い合わせ装置540、550をデータ処理・データベース記憶装置530に接続しても良い。

10

【0178】

上記の実施形態では、プログラミング命令により定義される処理ルーチンを使用して、コンピュータにより処理を実行していた。しかし、処理の一部又は全てをハードウェアを使用して実行することも可能であろう。

【0179】

複数の参加者の間で行われる会議に関して上記の実施形態を説明したが、本発明はこの用途には限定されず、フィルムセットについて画像データ及び音声データを処理するなどの他の用途にも適用することができる。

20

【0180】

上記の変形の異なる組み合わせも言うまでもなく可能であり、本発明の趣旨から逸脱せずにその他の変更や変形を実施することができる。

【0181】

<第2の実施形態>

図21を参照すると、この実施形態では、1台のビデオカメラ602と、1つ又は複数のマイクロホン604とを使用して、複数の人物606、608、610、612の間で行われている会議から画像データと音声データをそれぞれ記録している。

30

【0182】

ビデオカメラ602からの画像データと、マイクロホン604からの音声データはケーブル(図示せず)を介してコンピュータ620に入力され、コンピュータ620は受信したデータを処理し、データをデータベースに格納して、会議のアーカイブ記録を作成する。後に、このデータベースから情報を検索することができる。

【0183】

コンピュータ620は、従来のように、表示装置626や、この実施形態においてはキーボード628及びマウス630であるユーザ入力装置と共に、1つ又は複数のプロセッサ、メモリ、サウンドカードなどを含む処理装置624を有する従来通りのパーソナルコンピュータである。

40

【0184】

コンピュータ620の構成要素と、それらの構成要素に対し入出力されるデータの流れを図22に概略的に示す。

【0185】

図22を参照すると、処理装置624は、例えば、ディスク632などのデータ記憶媒体に格納されたデータとして及び/又は例えば、インターネットなどの通信ネットワーク(図示せず)を介する送信又は大気中を通る送信により遠隔データベースから処理装置624に入力され且つ/又はキーボード628などのユーザ入力装置又は他の入力装置を介してユーザにより処理装置624に入力される信号634として入力されるプログラミング

50

命令に従って動作するようにプログラムされている。

【 0 1 8 6 】

プログラミング命令によりプログラムされると、処理装置 6 2 4 は処理動作を実行するための複数の機能ユニットとして有効に構成される。そのような機能ユニットの例と、それらの配線を図 2 2 に示す。しかし、図 2 2 に示すユニットと配線は概念上のものであり、単に理解を助けるために例示を目的として示されているにすぎない。従って、図 2 2 の機能ユニット及び配線は、処理装置 6 2 4 のプロセッサ、メモリなどが実際に構成される厳密なユニットや接続関係を必ずしも表してはいない。

【 0 1 8 7 】

図 2 2 に示す機能ユニットに関して説明すると、中央制御装置 6 3 6 はユーザ入力装置 6 2 8 , 6 3 0 からの入力を処理し、且つユーザによりディスク 6 3 8 などの記憶装置に格納されたデータとして又は処理装置 6 2 4 へ送信される信号 6 4 0 として処理装置 6 2 4 に入力されるデータを受信する。また、中央処理装置 6 3 6 はその他の複数の機能ユニットに対して制御及び処理を実行する。メモリ 6 4 2 は、中央制御装置 6 3 6 及びその他の機能ユニットにより使用されるべきメモリである。

10

【 0 1 8 8 】

頭部追跡装置 6 5 0 はビデオカメラ 6 0 2 から受信した画像データを処理して、会議中の各々の参加者 6 0 6、6 0 8、6 1 0、6 1 2 の頭部の位置と向きを三次元で追跡する。この実施形態では、この追跡を実行するために、頭部追跡装置 6 5 0 は各々の参加者の頭部の三次元コンピュータモデルを定義するデータと、頭部の特徴を定義するデータとを使用する。それらのデータは、後述するように、頭部モデル記憶装置 6 5 2 に格納される。

20

【 0 1 8 9 】

音声認識プロセッサ 6 5 4 はマイクロホン 6 0 4 から受信される音声データを処理する。音声認識プロセッサ 6 5 4 は、「Dragon Dictate」又はIBMの「ViaVoice」などの従来の音声認識プログラムに従って動作し、参加者 6 0 6、6 0 8、6 1 0、6 1 2 により話された言葉に対応するテキストデータを生成する。音声認識処理を実行するために、音声認識プロセッサ 6 5 4 は、参加者 6 0 6、6 0 8、6 1 0、6 1 2 ごとの音声認識パラメータを定義するデータを使用する。このデータは音声認識パラメータ記憶装置 6 5 6 に格納される。すなわち、音声認識パラメータ記憶装置 6 5 6 に格納されるデータは、音声認識プロセッサを従来の方式で訓練することにより生成される各参加者の音声プロファイルを定義するデータである。例えば、このデータは、訓練後にDragon Dictateの「ユーザファイル」に格納されるデータである。

30

【 0 1 9 0 】

アーカイブプロセッサ 6 5 8 は、頭部追跡装置 6 5 0 及び音声認識プロセッサ 6 5 4 から受信したデータを使用して、会議アーカイブデータベース 6 6 0 に格納すべきデータを生成する。すなわち、後述するように、カメラ 6 0 2 からの映像データとマイクロホン 6 0 4 からの音声データを、音声認識プロセッサ 6 5 4 からのテキストデータ及び会議中の各参加者が所定の時点で誰を見ていたかを定義するデータと共に会議アーカイブデータベース 6 6 0 に格納するのである。

【 0 1 9 1 】

40

テキストサーチャ 6 6 2 は、中央制御装置 6 3 6 と関連して、会議アーカイブデータベース 6 6 0 を探索して、後に更に詳細に説明するように、ユーザにより指定される探索基準に適合する会議の 1 つ又は複数の部分に対応する音声データと映像データを見出し、再生するために使用される。

【 0 1 9 2 】

表示プロセッサ 6 6 4 は、中央制御装置 6 3 6 の制御の下に、表示装置 6 2 6 を介してユーザに情報を表示すると共に、会議アーカイブデータベース 6 6 0 に格納されている音声データと映像データを再生する。

【 0 1 9 3 】

出力プロセッサ 6 6 6 はアーカイブデータベース 6 6 0 からのデータの一部又は全てを、

50

例えば、ディスク 6 6 8 などの記憶装置へ出力するか、又は信号 6 7 0 として出力する。

【 0 1 9 4 】

会議を始める前に、処理装置 6 2 4 が要求される処理動作を実行できるようにするために必要なデータを入力することによりコンピュータ 6 2 0 を初期設定しなければならない。

【 0 1 9 5 】

図 2 3 は、この初期設定中に処理装置 6 2 4 により実行される処理動作を示す。

【 0 1 9 6 】

図 2 3 を参照すると、ステップ S 3 0 2 では、中央制御装置 6 3 6 は表示プロセッサ 6 6 4 に、ユーザが会議に参加する各人物の名前を入力することを要求するメッセージを表示装置 6 2 6 に表示させる。

10

【 0 1 9 7 】

ステップ S 3 0 4 では、中央制御装置 6 3 6 は、例えば、ユーザがキーボード 6 2 8 を使用して入力した名前を定義するデータを受信すると、各参加者に独自の参加者番号を割り当て、参加者番号と参加者の名前との関係を定義するデータ、例えば、図 2 4 に示すテーブル 6 8 0 を会議アーカイブデータベース 6 6 0 に格納する。

【 0 1 9 8 】

ステップ S 3 0 6 では、中央制御装置 6 3 6 は頭部モデル記憶装置 6 5 2 を探索して、会議の参加者ごとに頭部モデルを定義するデータが既に格納されているか否かを判定する。

【 0 1 9 9 】

ステップ S 3 0 6 で、1 人又は 2 人以上の参加者について頭部モデルがまだ格納されていないと判定されれば、ステップ S 3 0 8 で、中央制御装置 6 3 6 は表示プロセッサ 6 6 4 に、モデルがまだ格納されていない各参加者の頭部モデルを定義するデータをユーザが入力することを要求するメッセージを表示装置 6 2 6 に表示させる。

20

【 0 2 0 0 】

これに回答して、ユーザは、例えば、ディスク 6 3 8 などの記憶媒体に格納されたデータとして要求された頭部モデルを定義するデータを入力するか、又は接続している処理装置から信号 6 4 0 としてこのデータをダウンロードすることによりデータを入力する。そのような頭部モデルは、例えば、Valente 他「An Analysis / Synthesis Cooperation for Head Tracking and Video Face Cloning」(Proceedings ECCV ' 9 8 Workshop on Perception of Human Action、ドイツ、フライブルク大学、1 9 9 8 年 6 月 6 日に掲載)に記載されているような従来の方式で生成されれば良い。

30

【 0 2 0 1 】

ステップ S 3 1 0 では、中央制御装置 6 3 6 は、ユーザにより入力されたデータを頭部モデル記憶装置 6 5 2 に格納する。

【 0 2 0 2 】

ステップ S 3 1 2 では、中央制御装置 6 3 6 及び表示プロセッサ 6 6 4 は、ユーザにより入力された各々の三次元コンピュータ頭部モデルをレンダリングして、ユーザが各モデルにおいて少なくとも 7 つの特徴を識別することを要求するメッセージと共に、モデルをユーザに対し表示装置 6 2 6 を介して表示する。

【 0 2 0 3 】

これに回答して、ユーザは、マウス 6 3 0 を使用して、参加者の頭部の正面、側面及び(可能であれば)背面にある顕著な特徴、例えば、目尻、鼻孔、口、耳又は参加者がかけている眼鏡の特徴などに対応する点を各モデルで指定する。

40

【 0 2 0 4 】

ステップ S 3 1 4 では、中央制御装置 6 3 6 はユーザにより定義された特徴を頭部モデル記憶装置 6 5 2 に格納する。

【 0 2 0 5 】

他方、ステップ S 3 0 6 で、参加者ごとに頭部モデルが頭部モデル記憶装置 6 5 2 に既に記憶されていると判定された場合には、ステップ S 3 0 8 から S 3 1 4 を省略する。

【 0 2 0 6 】

50

ステップS 3 1 6では、中央制御装置6 3 6は音声認識パラメータ記憶装置6 5 6を探索して、参加者ごとに音声認識パラメータが既に格納されているか否かを判定する。

【0 2 0 7】

ステップS 3 1 6で、一部の参加者について音声認識パラメータを利用できないと判定されれば、ステップS 3 1 8で、中央制御装置6 3 6は表示プロセッサ6 6 4に、パラメータがまだ格納されていない各参加者についてユーザが音声認識パラメータを入力することを要求するメッセージを表示装置6 2 6に表示させる。

【0 2 0 8】

これに応答して、ユーザは、例えば、ディスク6 3 8などの記憶媒体に格納されたデータとして又は遠隔処理装置からの信号6 4 0として、必要な音声認識パラメータを定義するデータを入力する。先に述べた通り、これらのパラメータはユーザの音声のプロファイルを定義し、音声認識プロセッサを従来の方式で訓練することにより生成される。従って、例えば、Dragon Dictateを組み込んだ音声認識プロセッサの場合、ユーザにより入力される音声認識パラメータはDragon Dictateの「ユーザファイル」に格納されたパラメータに相当する。

10

【0 2 0 9】

ステップS 3 2 0では、中央制御装置6 3 6は、ユーザにより入力されたデータを音声認識パラメータ記憶装置6 5 6に格納する。

【0 2 1 0】

他方、ステップS 3 1 6で、参加者ごとに音声認識パラメータが既に利用可能な状態になっていると判定された場合には、ステップS 3 1 8及びS 3 2 0を省略する。

20

【0 2 1 1】

ステップS 3 2 2では、中央制御装置6 3 6は表示プロセッサ6 6 4に、ユーザがカメラ6 0 2の校正を可能にするためのステップを実行することを要求するメッセージを表示装置6 2 6に表示させる。

【0 2 1 2】

これに応答して、ユーザは必要なステップを実行し、ステップS 3 2 4では、中央制御装置6 3 6はカメラ6 0 2を校正するための処理を実行する。すなわち、この実施形態においては、ユーザにより実行されるステップ及び中央制御装置6 3 6により実行される処理は、Wiles及びDavisonの「Calibrating and 3D Modelling with a Multi - Camera System」(1 9 9 9 IEEE Workshop on Multi - View Modelling and Analysis of Visual Scenes、ISBN 0 7 6 9 5 0 1 1 0 9に掲載)に記載されているような方式で実行される。これにより、会議室に対するカメラ6 0 2の位置と向きを定義する校正データと、カメラの固有パラメータ(横縦比、焦点距離、主点及び一次半径方向ひずみ係数)とが生成される。校正データはメモリ6 4 2に格納される。

30

【0 2 1 3】

ステップS 3 2 6では、中央制御装置6 3 6は表示プロセッサ6 6 4に、会議の次の参加者(初めてステップS 3 2 6を実行する場合には、これは最初の参加者である)が着席することを要求するメッセージを表示装置6 2 6に表示させる。

【0 2 1 4】

40

ステップS 3 2 8では、要求された参加者に着席する時間を与えるために、処理装置6 2 4は所定の時間待機し、その後、ステップS 3 3 0で、中央制御装置6 3 6はカメラ6 0 2からの画像データを処理して、着席した参加者の頭部の推定位置を判定する。すなわち、この実施形態においては、中央制御装置6 3 6は、参加者の肌の色に対応する色(この色は、頭部モデル記憶装置6 5 2に格納されている参加者の頭部モデルを定義するデータから判定される)を有する、カメラ6 0 2からの画像データの1フレーム中の各部分を識別するために従来通りの処理を実行し、次に、会議室内の最も高い位置に相当する部分を選択する(頭部は人体の中で最も高い位置にある肌色の部分であると想定されるため)。画像中の識別された部分の位置と、ステップS 3 2 4で判定されたカメラ校正パラメータとを使用して、中央制御装置6 3 6は従来の方式により頭部の三次元推定位置を判定する

50

。

【0215】

ステップS332では、中央制御装置636は参加者の頭部の三次元推定向きを判定する。すなわち、この実施形態においては、中央制御装置636は頭部モデル記憶装置652に格納されている参加者の頭部の三次元コンピュータモデルを複数の異なる向きについてレンダリングして、向きごとにそれぞれ対応するモデルの二次元画像を生成し、モデルの各二次元画像を参加者の頭部を示す、カメラ602からの映像フレームの部分と比較し、モデルの画像が映像データと最も良く整合する向きを選択する。この実施形態では、参加者の頭部のコンピュータモデルを108の異なる向きでレンダリングして、カメラ602からの映像データと比較すべき画像データを生成する。これらの向きは頭部モデルを0°（正面を向いている）、+45°（上を向いている）及び-45°（下を向いている）に相当する3つの頭部の傾きのそれぞれについて10°ずつのステップで36回回転させた向きに相当する。頭部モデルをレンダリングすることにより生成された画像データをカメラ602からの映像データと比較するときには、例えば、Schodl、Haro及びEssaの「Head Tracking Using a Textured Polygonal Model」（Proceedings 1998 Workshop on Perceptual User Interfacesに掲載）に記載されているような従来技法を使用する。

10

【0216】

ステップS334では、ステップS330で生成した参加者の頭部の推定位置と、ステップS332で生成した参加者の頭部の推定向きとを頭部追跡装置650に入力し、カメラ602から受信した画像データのフレームを処理して、参加者の頭部を追跡する。すなわち、この実施形態においては、頭部追跡装置650は、例えば、Valente他の「An Analysis/Synthesis Cooperation for Head Tracking and Video Face Cloning」（Proceedings EECV '98 Workshop on Perception of Human Action、ドイツ、フライブルク大学、1998年6月6日）に記載されているような従来方式で頭部を追跡するための処理を実行する。

20

【0217】

図25は、ステップS334で頭部追跡装置650により実行される処理動作の概要を示す。

【0218】

図25を参照すると、ステップS350では、頭部追跡装置650は参加者の頭部の現在推定3D位置及び現在推定3D向きを読み取る。ステップS350を初めて実行する場合には、これらは図23のステップS330及びS332で生成される推定位置及び推定向きである。

30

【0219】

ステップS352では、頭部追跡装置650はステップS324で生成されたカメラ校正データを使用して、頭部モデル記憶装置652に格納されている参加者の頭部の三次元コンピュータモデルをステップS350で読み取った推定位置及び推定向きに従ってレンダリングする。

【0220】

ステップS354では、頭部追跡装置650はカメラ602から受信された映像データの現在フレームについて画像データを処理し、ユーザにより識別され、ステップS314で格納された頭部の特徴のうち1つの特徴の期待位置を取り囲む各々の領域から画像データを取り出す。それらの期待位置は、ステップS350で読み取った推定位置及び推定向きと、ステップS324で生成されたカメラ校正データとから判定される。

40

【0221】

ステップS356では、頭部追跡装置650はステップS352でレンダリングし、生成した画像データと、ステップS354で取り出したカメラ画像データとを整合し、レンダリングされた頭部モデルに最も良く整合するカメラ画像データを見出す。

【0222】

ステップS358では、頭部追跡装置650は、ステップS356でレンダリングされた

50

頭部モデルに最も良く整合すると識別されたカメラ画像データを使用して、映像データの現在フレームについて参加者の頭部の3D位置及び3D向きを判定する。

【0223】

ステップS358を実行すると同時に、ステップS360では、ステップS356で判定されたカメラ画像データにおける頭部の特徴の位置を従来のカルマンフィルタに入力して、映像データの次のフレームについて参加者の頭部の推定3D位置及び推定3D向きを生成する。ビデオカメラ602から映像データのフレームが受信されている間、その参加者についてステップS350からS360を繰り返し実行する。

【0224】

再び図23に戻ると、ステップS336では、中央制御装置636は会議に他の参加者がいるか否かを判定し、参加者ごとに先に説明したように処理が実行され終わるまでステップS326からS336を繰り返す。しかし、参加者ごとにこれらのステップを実行している間、ステップS334では、頭部追跡装置650は既に着席した各参加者の頭部を追跡し続けている。

10

【0225】

ステップS336で、会議にそれ以上の参加者はなく、従って、各参加者の頭部が頭部追跡装置650により追跡されていることが判定されると、ステップS338で、中央制御装置636は、参加者間で会議を始めて良いことを指示するために、処理装置624から可聴信号を出力させる。

【0226】

20

図26は、参加者間で会議が行われている間に処理装置624により実行される処理動作を示す。

【0227】

図26を参照すると、ステップS370では、頭部追跡装置650は会議中の各参加者の頭部を追跡し続ける。ステップS370で頭部追跡装置650により実行される処理は、先にステップS334に関して説明した処理と同じであるので、ここでは説明を省略する。

【0228】

頭部追跡装置650がステップS370で各参加者の頭部を追跡しているのと同時に、ステップS372では、データを生成し、会議アーカイブデータベース660にデータを格納するための処理を実行する。

30

【0229】

図27は、ステップS372で実行される処理動作を示す。

【0230】

図27を参照すると、ステップS380では、アーカイブプロセッサ658は、参加者が誰を見ているかを定義するいわゆる「視線パラメータ」を参加者ごとに生成する。

【0231】

図28は、ステップS380で実行される処理動作を示す。

【0232】

図28を参照すると、ステップS410では、アーカイブプロセッサ658は各参加者の頭部の現在三次元位置を頭部追跡装置650から読み取る。これは、ステップS358(図25)で頭部追跡装置650により実行される処理において生成された位置である。

40

【0233】

ステップS412では、アーカイブプロセッサ658は次の参加者(初めてステップS412を実行する場合には、これは最初の参加者である)の頭部の現在向きを頭部追跡装置650から読み取る。ステップS412で読み取られる向きは、ステップS358(図25)で頭部追跡装置650により実行される処理において生成された向きである。

【0234】

ステップS414では、アーカイブプロセッサ658は、参加者がどこを見ているかを定義する線(いわゆる「視線」)と、参加者の頭部を別の参加者の頭部の中心と結ぶ概念上

50

の各々の線とが成す角度を判定する。

【0235】

図29及び図30を参照して更に詳細に説明する。図29及び図30には、1人の参加者、すなわち、図21の参加者610についてステップS414で実行される処理の一例が示されている。図29を参照すると、ステップS412で読み取られる参加者の頭部の向きは、その参加者の両目の中心の間の一点から出る、参加者の頭部に対し垂直な視線690を定義する。同様に、図30を参照すると、ステップS410で読み取られる全ての参加者の頭部の位置は、参加者610の両目の中心の間の一点から他の各々の参加者606、608、612の頭部の中心に至る概念上の線692、694、696を定義する。ステップS414では、アーカイブプロセッサ658は視線690と、概念上の線692、694、696とがそれぞれ成す角度698、700、702を判定する。

10

【0236】

再び図28に戻ると、ステップS416では、アーカイブプロセッサ658は最小値を有する角度698、700又は702を選択する。すなわち、図30に示す例で言えば、角度700が選択されることになるであろう。

【0237】

ステップS418では、アーカイブプロセッサ658は選択した角度が10°より小さいか否かを判定する。

【0238】

ステップS418で、角度が10°より小さいと判定されれば、アーカイブプロセッサ658は参加者の視線パラメータを、視線と最小の角度を成す概念上の線により結ばれている参加者の番号(図23のステップS304で割り当てられている)に設定する。すなわち、図30に示す例で言えば、角度700が10°より小さい場合には、角度700は視線690と、参加者610を参加者606と結ぶ概念上の線694とが成す角度であるので、視線パラメータは参加者606の参加者番号に設定されることになるであろう。

20

【0239】

他方、ステップS418で、最小角度が10°以上であると判定された場合には、ステップS422で、アーカイブプロセッサ658は参加者の視線パラメータを「0」に設定する。これは、視線690が概念上の線692、694、696のいずれにも十分に近接していないために、参加者はその他の参加者の誰も見ていないと判定されたことを示す。そのような状況は、例えば、参加者がメモ又は会議室内の他の何らかの物体を見ているときに起こりうるであろう。

30

【0240】

ステップS424では、アーカイブプロセッサ658は会議に他の参加者がいるか否かを判定し、参加者ごとに上記の処理がそれぞれ実行され終わるまでステップS412からS424を繰り返す。

【0241】

再び図27に戻ると、ステップS382では、中央制御装置636及び音声認識プロセッサ654は、映像データの現在フレームについてマイクロホン604から音声データが受信されたか否かを判定する。

40

【0242】

ステップS382で、音声データが受信されていると判定されれば、ステップS384で、アーカイブプロセッサ658はステップS380で生成された視線パラメータを処理して、会議中のどの参加者が話をしているかを判定する。

【0243】

図31は、ステップS384でアーカイブプロセッサ658により実行される処理動作を示す。

【0244】

図31を参照すると、ステップS440では、ステップS380で生成された各視線パラメータ値の出現回数を判定し、ステップS442では、出現回数が最も多い視線パラメータ

50

タ値を選択する。すなわち、図 27 のステップ S 3 8 0 で実行される処理は、会議中の参加者ごとに、映像データの現在フレームについて 1 つの視線パラメータ値を生成するのである（従って、図 21 に示す例では、4 つの値が生成されることになるであろう）。各視線パラメータは、その他の参加者のうち 1 人の参加者番号に相当する値又は「0」を有する。従って、ステップ S 4 4 0 及び S 4 4 2 では、アーカイブプロセッサ 6 5 8 は、ステップ S 3 8 0 で生成された視線パラメータ値の中で、映像データの現在フレームについて最も多くの回数で出現する値はどれであることを判定する。

【0245】

ステップ S 4 4 4 では、最も出現回数の多い視線パラメータが「0」の値を有するか否かを判定し、「0」の値であれば、ステップ S 4 4 6 で、次に出現回数の多い視線パラメータ値を選択する。これに対し、ステップ S 4 4 4 で、選択された値が「0」ではないと判定された場合には、ステップ S 4 4 6 を省略する。

10

【0246】

ステップ S 4 4 8 では、選択された視線パラメータ値（すなわち、ステップ S 4 4 2 で選択された値、又はその値が「0」であれば、ステップ S 4 4 6 で選択された値）を話をしている参加者として識別する。これは、会議中の参加者の大半は話者である参加者を見ているからである。

【0247】

再び図 27 に戻ると、ステップ S 3 8 6 では、アーカイブプロセッサ 6 5 8 は話者である参加者の視線パラメータ値、すなわち、ステップ S 3 8 0 で生成された、話者である参加者が誰を見ているかを定義する視線パラメータ値を後の解析に備えて、例えば、メモリ 6 4 2 に格納する。

20

【0248】

ステップ S 3 8 8 では、アーカイブプロセッサ 6 5 8 はステップ S 3 8 4 で判定された話者である参加者のアイデンティティを音声認識プロセッサ 6 5 4 に報知する。これに 응답して、音声認識プロセッサ 6 5 4 は話者である参加者の音声認識パラメータを音声認識パラメータ記憶装置 6 5 6 から選択し、選択したパラメータを使用して、受信された音声データに対して音声認識処理を実行し、話者である参加者が話した言葉に対応するテキストデータを生成する。

【0249】

30

他方、ステップ S 3 8 2 で、受信された音声データが話し言葉を含まないと判定された場合には、ステップ S 3 8 4 から S 3 8 8 を省略する。

【0250】

ステップ S 3 9 0 では、アーカイブプロセッサ 6 5 8 はカメラ 6 0 2 から受信された映像データの現在フレームと、マイクロホン 6 0 4 から受信された音声データとを従来の方式で MPEG 2 データを符号化し、符号化されたデータを会議アーカイブデータベース 6 6 0 に格納する。

【0251】

図 32 は、会議アーカイブデータベース 6 6 0 へのデータの格納状態を概略的に示す。図 32 に示す格納構造は概念的なものであり、単に理解を助けるために例示を目的として提示されているにすぎない。従って、図 32 に示す構造は、データが実際に会議アーカイブデータベース 6 6 0 に格納される厳密な状態を必ずしも表してはいない。

40

【0252】

図 32 を参照すると、会議アーカイブデータベース 6 6 0 は水平軸 8 0 0 により表される時間情報を格納している。水平軸 8 0 0 に沿った各単位は所定の量の時間、例えば、カメラ 6 0 2 から受信される映像データの 1 つのフレームを表している。ステップ S 3 9 0 で生成される MPEG 2 データは、時間情報と共に、データ 8 0 2 として会議アーカイブデータベース 6 6 0 に格納されている（この時間情報は、図 32 には、水平軸 8 0 0 に沿った MPEG 2 データ 8 0 2 の位置により概略的に表されている）。

【0253】

50

再び図 27 に戻ると、ステップ S 392 では、アーカイブプロセッサ 658 は、現在フレームについてステップ S 388 で音声認識プロセッサ 654 により生成されたテキストデータを会議アーカイブデータベース 660 に格納する（図 32 には 804 で示されている）。すなわち、テキストデータは対応する MPEG2 データへのリンクを伴って格納される。このリンクは、図 32 には、テキストデータが MPEG2 データと同じ縦列に格納されることによって表されている。話をしていない参加者からは格納すべきテキストデータが得られないことは理解されるであろう。図 32 に示す例では、初めの 10 個のタイムスロットにわたり参加者 1 についてテキストが格納され（806 で示す）、12 番目から 20 番目のタイムスロットには参加者 3 のテキストが格納され（808 で示す）、21 番目のタイムスロットには参加者 4 のテキストが格納されている（810 で示す）。この例では、図 32 に示すタイムスロットの間、参加者 2 は話をしなかったので、参加者 2 のテキストは格納されていない。

10

【0254】

ステップ S 394 では、アーカイブプロセッサ 658 は、ステップ S 380 で生成された参加者ごとの視線パラメータ値を会議アーカイブデータベース 660 に格納する（図 32 には 812 で示す）。図 32 を参照すると、視線パラメータ値は、参加者ごとに、関連する MPEG2 データ 802 及び関連するテキストデータ 804 へのリンクと共に格納されている（このリンクは、図 32 では、視線パラメータ値が関連する MPEG2 データ 802 及び関連するテキストデータ 804 と同じ縦列に格納されることにより示されている）。従って、一例として第 1 のタイムスロットに関して言えば、参加者 1 の視線パラメータ値は、参加者 1 が参加者 3 を見ていることを指示する「3」であり、参加者 2 の視線パラメータ値は、参加者 2 が参加者 1 を見ていることを指示する「1」であり、参加者 3 の視線パラメータ値も、参加者 3 が同様に参加者 1 を見ていることを指示する「1」であり、参加者 4 の視線パラメータ値は、参加者 4 が他のどの参加者も見していない（図 21 に示す例では、612 で示される参加者は他の参加者ではなく、自分のメモを見ている）ことを指示する「0」である。

20

【0255】

ステップ S 396 では、中央制御装置 636 及びアーカイブプロセッサ 658 は、会議中の参加者の 1 人が話を止めたか否かを判定する。この実施形態においては、この検査は、所定の参加者のテキストデータが直前のタイムスロットには存在したが、現在タイムスロットには存在しないことを判定するためにテキストデータ 804 を検査することにより実行される。ある参加者についてこの条件が満たされれば（すなわち、参加者が話を止めたならば）、ステップ S 398 で、アーカイブプロセッサ 658 は、話を止めた参加者について、先にステップ S 386 を実行したときに格納されていた視線パラメータ値を処理して（それらの視線パラメータ値は、その時点で止まった話をしていた期間中にその参加者が誰を見ていたかを定義する）、視線ヒストグラムを定義するデータを生成する。すなわち、参加者が話をしていた期間の視線パラメータ値を処理して、その期間中に話者である参加者がその他の参加者の各々を見ていた時間の割合（％）を定義するデータを生成するのである。

30

【0256】

図 33A 及び図 33B は図 32 のテキスト 806 及び 808 の期間にそれぞれ対応する視線ヒストグラムを示す。

40

【0257】

図 32 及び図 33A を参照して説明すると、参加者 1 が話していた期間 806 の間、図 33A に 900 で示すように、参加者 1 は 10 個のタイムスロットのうち 6 個（すなわち、参加者 1 が話していた期間全体の長さの 60％）にわたり参加者 3 を見ており、また、図 33A に 910 で示すように、4 個のタイムスロット（すなわち、時間の 40％）にわたり参加者 4 を見ていた。

【0258】

同様に、図 32 及び図 33B を参照すると、期間 808 の間、図 33B に 920 で示すよう

50

に、参加者 3 は時間の約 4 5 % にわたり参加者 1 を見ており、図 3 3 B に 9 3 0 で示すように、時間の約 3 3 % にわたり参加者 4 を見ており、図 3 3 B に 9 4 0 で示すように、時間の約 2 2 % にわたり参加者 2 を見ていた。

【 0 2 5 9 】

再び図 2 7 に戻ると、ステップ S 4 0 0 では、ステップ S 3 9 8 で生成された視線ヒストグラムをそれが生成された関連するテキストの期間にリンクさせて、会議アーカイブデータベース 6 6 0 に格納する。図 3 2 を参照すると、格納された視線ヒストグラムは 8 1 4 で示されており、8 1 6 で示されるテキスト期間 8 0 6 に対応するヒストグラムを定義するデータと、8 1 8 で示されるテキスト期間 8 0 8 に対応するヒストグラムを定義するデータとを伴う。図 3 2 において、視線ヒストグラムと関連するテキストとの間のリンクは、視線ヒストグラムがテキストデータと同じ縦列に格納されることにより示されている。

10

【 0 2 6 0 】

他方、ステップ S 3 9 6 で、現在時限について、参加者の 1 人が話を止めていないと判定された場合には、ステップ S 3 9 8 及び S 4 0 0 を省略する。

【 0 2 6 1 】

ステップ S 4 0 2 では、中央制御装置 6 3 6 は、カメラ 6 0 2 から映像データの別のフレームが受信されたか否かを判定する。カメラ 6 0 2 から画像データが受信されている間は、ステップ S 3 8 0 から S 4 0 2 を繰り返し実行する。

【 0 2 6 2 】

会議アーカイブデータベース 6 6 0 にデータが格納されている場合、会議に関連するデータを検索するために会議アーカイブデータベース 6 6 0 を問い合わせても良い。

20

【 0 2 6 3 】

図 3 4 は、ユーザにより指定される探索基準を満たす会議の各部分に関連するデータを検索する目的で会議アーカイブデータベース 6 6 0 を探索するために実行される処理動作を示す。

【 0 2 6 4 】

図 3 4 を参照すると、ステップ S 5 0 0 では、中央制御装置 6 3 6 は表示プロセッサ 6 6 4 に、要求される会議アーカイブデータベース 6 6 0 の探索を定義する情報をユーザが入力することを要求するメッセージを表示装置 6 2 6 に表示させる。すなわち、この実施形態においては、中央制御装置 6 3 6 は図 3 5 A に示す表示を表示装置 6 2 6 に出現させる。

30

【 0 2 6 5 】

図 3 5 A を参照すると、ユーザは、会議アーカイブデータベース 6 6 0 の中で見出したい会議の 1 つ又は複数の部分を定義する情報を入力することを求められる。すなわち、この実施形態においては、ユーザは話をしていた参加者を定義する情報 1 0 0 0 と、情報 1 0 0 0 の中で識別される参加者が話した 1 つ又は複数のキーワードから成る情報 1 0 1 0 と、情報 1 0 0 0 の中で識別される参加者が話しかけていた参加者を定義する情報 1 0 2 0 とを入力することを求められる。更に、ユーザは、探索を実行すべき会議の 1 つ又は複数の部分を定義する時間情報を入力することができる。すなわち、ユーザは、その時間を越えたときに探索を中断すべき会議中の時間（すなわち、指定される時間以前の会議の期間を探索すべきである）を定義する情報 1 0 3 0 と、その時間の後に探索を実行すべきである会議中の時間を定義する情報 1 0 4 0 と、探索を実行すべき期間の開始時間と終了時間をそれぞれ定義する情報 1 0 5 0 及び 1 0 6 0 とを入力することができる。この実施形態では、情報 1 0 3 0、1 0 4 0、1 0 5 0 及び 1 0 6 0 は、例えば、分単位などの絶対期限で時間を指定するか、又は会議時間全体に占める割合を指示する小数値を入力することにより相対期限で時間を指定するかのいずれかにより入力されれば良い。例えば、情報 1 0 3 0 として 0 . 2 5 の値を入力すると、探索は会議の初めの四分の一に限られるであろう。

40

【 0 2 6 6 】

この実施形態では、ユーザは 1 回の探索に際して全ての情報 1 0 0 0、1 0 1 0 及び 1 0

50

20を入力することを求められるわけではなく、そのうち1つ又は2つの情報を省いても良い。ユーザが情報1000、1010及び1020の全てを入力すれば、会議中に情報1000の中で識別される参加者が情報1020の中で識別される参加者に話しかけていた各部分及び情報1010の中で定義されるキーワードを話していた各部分を識別するための探索が実行される。これに対し、情報1010を省いた場合には、会議中に情報1000の中で識別される参加者が、何を言ったかに関わらず、情報1020の中で定義される参加者に話しかけていた各部分を識別するための探索が実行されることになる。情報1010及び1020を省いた場合には、会議中に情報1000の中で定義される参加者が何を誰に向かって話したかに関わらず、話をしていた各部分を識別するための探索が実行される。情報1000を省いた場合には、会議中にいずれかの参加者が情報1010の中で定義されるキーワードを情報1020の中で定義される参加者に向かって話した各部分を識別するための探索が実行される。情報1000及び1010を省いた場合には、会議中にいずれかの参加者が情報1020の中で定義される参加者に話しかけた各部分を識別するための探索が実行される。情報1020を省いた場合には、会議中に情報1000の中で定義される参加者が、誰に向かって話したかに関わらず、情報1010の中で定義されるキーワードを話した各部分を識別するための探索が実行される。同様に、情報1000及び1020を省いた場合には、会議中に、誰が誰に向かって話したかに関わらず、情報1010の中で識別されるキーワードが話された各部分を識別するための探索が実行される。

10

【0267】

20

更に、ユーザは時間情報1030、1040、1050及び1060の全てを入力しても良いし、そのうち1つ又は複数の情報を省いても良い。

【0268】

探索を定義するために必要な情報を全て入力したならば、ユーザは、マウス630などのユーザ入力装置を使用して領域1070をクリックすることにより探索を開始する。

【0269】

再び図34に戻ると、ステップS802では、中央制御装置636はユーザにより入力された探索情報を読み取り、命令された探索を実行する。すなわち、この実施形態においては、中央制御装置636は情報1000又は1020の中で入力された参加者の名前をテーブル680(図24)を使用して参加者番号に変換し、情報1000の中で定義される参加者(情報1000が入力されていない場合には全ての参加者)についてテキスト情報804を考慮する。ユーザにより情報1020が入力されていれば、テキストの期間ごとに、中央制御装置636は対応する視線ヒストグラムを定義するデータを検査して、情報1020の中で定義される参加者のヒストグラムにおける注目時間の割合がこの実施形態では25%である閾値以上であるか否かを判定する。このようにして、話し言葉(テキスト)の期間を考慮して、話者である参加者が話している時間の少なくとも25%にわたって情報1020の中で定義される参加者を見ていたならば、情報1000の中で定義される参加者が話をしている間に他の参加者を見たとしても、情報1020の中で定義される参加者に話しかけていたという基準を満たす。従って、情報1020の中で2人以上の参加者が指定されていれば、視線ヒストグラムの値が2人以上の参加者について25%以上であるような話の期間が識別されるであろう。ユーザにより情報1010が入力されていた場合、中央制御装置636及びテキストサーチ662は、先に情報1000及び1020に基づいて識別されたテキストの各部分(情報1000及び1020が入力されていなければテキストの全ての部分)を探索して、情報1010の中で識別されるキーワードを含む各部分を識別する。ユーザにより時間情報が入力されていた場合、上記の探索はそれらの期限により定義される会議の時間に限られる。

30

40

【0270】

ステップS504では、中央制御装置636は表示プロセッサ664に、探索中に識別された関連話題のリストを表示装置626を介してユーザに対し表示させる。すなわち、中央制御装置636は図35Bに示すような情報をユーザに対し表示させる。図35Bを参照

50

すると、探索パラメータを満足させるそれぞれの話題のリストが作成されており、その話題について開始時間を絶対期限と、会議時間全体に占める割合の双方で定義する情報が表示される。そこで、ユーザは、マウス 6 3 0 を使用してリスト中の必要な話題をクリックすることにより、話題の 1 つを選択し、再生することができる。

【 0 2 7 1 】

ステップ S 5 0 6 では、中央制御装置 6 3 6 はステップ S 5 0 4 でユーザにより実行された選択を読み取り、格納されている会議の関連部分の MPEG 2 データ 8 0 2 を会議アーカイブデータベース 6 6 0 から再生させる。すなわち、中央制御装置 6 3 6 及び表示プロセッサ 6 6 4 は MPEG 2 データ 8 0 2 を復号し、画像データと音声を表示装置 6 2 6 を介して出力する。

10

【 0 2 7 2 】

ステップ S 5 0 8 では、中央制御装置 6 3 6 は、ユーザが会議アーカイブデータベース 6 6 0 の問い合わせを中止することを望むか否かを判定し、望まないのであれば、ステップ S 5 0 0 から S 5 0 8 を繰り返す。

【 0 2 7 3 】

以上説明した本発明の実施形態に対し、様々な変形や変更を実施することができる。

【 0 2 7 4 】

例えば、上記の実施形態では、マイクロホン 6 0 4 は会議室のテーブルの上に設けられていた。しかし、その代わりに、ビデオカメラ 6 0 2 のマイクロホンを使用して音声データを記録しても良い。

20

【 0 2 7 5 】

上記の実施形態では、1 台のビデオカメラ 6 0 2 からの画像データを処理していた。しかし、各参加者の頭部の追跡精度を向上させるために、複数台のビデオカメラからの映像データを処理しても良い。例えば、複数台のカメラからの画像データを図 2 5 のステップ S 3 5 0 から S 3 5 6 におけるように処理し、全てのカメラから得られたデータをステップ S 3 6 0 で従来のようにカルマンフィルタに入力して、カメラごとに映像データの次のフレームにおける各参加者の頭部の位置と向きを更に正確に推定しても良い。複数台のカメラを使用する場合、会議アーカイブデータベース 6 6 0 に格納される MPEG 2 データ 8 0 2 は全てのカメラからの映像データということになり、図 3 4 のステップ S 5 0 4 及び S 5 0 6 では、ユーザが選択した 1 台のカメラからの画像データが再生される。

30

【 0 2 7 6 】

上記の実施形態では、所定の参加者の視線パラメータは、その参加者が他のどの参加者を見ているかを定義していた。しかし、参加者が表示板、映写機のスクリーンなどのどの物体を見ているかを定義するために視線パラメータを使用しても良い。この場合、会議アーカイブデータベース 6 6 0 を問い合わせるときに、図 3 5 A の情報 1 0 2 0 を利用して、参加者が話している間に誰を又は何を見ていたかを指定することが可能になるであろう。

【 0 2 7 7 】

上記の実施形態では、ステップ S 5 0 2 (図 3 4) で、テキストの特定の部分の視線ヒストグラムを考慮し、視線ヒストグラムにおける別の参加者への注目時間の割合が所定の閾値以上である場合に、参加者はその別の参加者に話しかけていたと判定していた。しかし、閾値を使用せずに、テキストの期間中に話者である参加者が見ていた参加者を視線ヒストグラムにおいて最も大きな割合の注目値を有する参加者 (例えば、図 3 3 A の参加者 3 及び図 3 3 B の参加者 1) であると定義しても良い。

40

【 0 2 7 8 】

上記の実施形態では、カメラ 6 0 2 及びマイクロホン 6 0 4 からデータが受信されている間、MPEG 2 データ 8 0 2、テキストデータ 8 0 4、視線パラメータ 8 1 2 及び視線ヒストグラム 8 1 4 をリアルタイムで会議アーカイブデータベース 6 6 0 に格納していた。しかし、映像データと音声データを格納しておき、データ 8 0 2、8 0 4、8 1 2 及び 8 1 4 をリアルタイムではなく生成し、会議アーカイブデータベース 6 6 0 に格納しても良い。

【 0 2 7 9 】

50

上記の実施形態では、会議の定義された部分のデータを検索するために会議アーカイブデータベース 660 を問い合わせる前に、MPEG2 データ 802、テキストデータ 804、視線パラメータ 812 及び視線ヒストグラム 814 を生成し、会議アーカイブデータベース 660 に格納していた。しかし、探索の要求に先立ってデータを生成、格納するのではなく、会議アーカイブデータベース 660 の探索がユーザにより要求されるのに応答して、格納されている MPEG2 データ 802 を処理することによりデータ 804、812 及び 814 の一部又は全てを生成しても良い。例えば、上記の実施形態では、ステップ S398 及び S400 (図 27) で視線ヒストグラム 814 をリアルタイムで計算、格納していたが、ユーザにより探索要求が入力されるのに応答してそれらのヒストグラムを計算することもできるであろう。

10

【0280】

上記の実施形態では、テキストデータ 804 を会議アーカイブデータベース 660 に格納していた。テキストデータ 804 の代わりに、音声データを会議アーカイブデータベース 660 に格納しても良い。その場合、格納されている音声データ自体を音声認識処理を使用してキーワードを求めて探索するか、又は音声認識処理を使用して音声データをテキストに変換し、従来のテキストサーチャを使用してテキストを探索すれば良い。

【0281】

上記の実施形態では、処理装置 624 はアーカイブすべきデータを受信し、生成するための機能構成要素 (例えば、中央制御装置 636、頭部追跡装置 650、頭部モデル記憶装置 652、音声認識プロセッサ 654、音声認識パラメータ記憶装置 656 及びアーカイブプロセッサ 658) と、アーカイブデータを格納するための機能構成要素 (例えば、会議アーカイブデータベース 660) と、データベースを探索し、そこから情報を検索するための機能構成要素 (例えば、中央制御装置 636 及びテキストサーチャ 662) とを含んでいた。しかし、これらの機能構成要素を別個の装置に設けても良い。例えば、アーカイブすべきデータを生成する 1 つ又は複数の装置と、データベース探索のための 1 つ又は複数の装置とを、インターネットなどのネットワークを介して 1 つ又は複数のデータベースに接続しても良い。

20

【0282】

また、図 36 を参照して説明すると、一箇所又は複数箇所での会議 1100、1110、1120 からの映像データと音声データをデータ処理・データベース記憶装置 1130 (アーカイブデータを生成し且つ格納するための機能構成要素を具備する) に入力し、データベースを問い合わせ、そこから情報を検索するために、1 つ又は複数のデータベース問い合わせ装置 1140、1150 をデータ処理・データベース記憶装置 1130 に接続しても良い。

30

【0283】

上記の実施形態では、プログラミング命令により定義される処理ルーチンを使用して、コンピュータにより処理を実行していた。しかし、処理の一部又は全てをハードウェアを使用して実行しても良い。

【0284】

以上、複数の参加者の間で行われる会議に関して実施形態を説明したが、本発明はこの用途には限定されず、フィルムセットなどについて画像データ及び音声データを処理するなど、他の用途にも適用することができる。

40

【0285】

上記の変形例の異なる組み合わせも言うまでもなく可能であり、本発明の趣旨から逸脱せずにその他の変更や変形を実施することができる。

【図面の簡単な説明】

【図 1】第 1 の実施形態における複数の参加者の間の会議からの音声データ及び映像データの記録を示す図である。

【図 2】第 1 の実施形態の処理装置内部の概念上の機能構成要素の一例を示すブロック線図である。

50

【図 3 A】図 1 に示す参加者間の会議が始まる以前に図 2 の処理装置 2 4 により実行される処理動作を示す図である。

【図 3 B】図 1 に示す参加者間の会議が始まる以前に図 2 の処理装置 2 4 により実行される処理動作を示す図である。

【図 3 C】図 1 に示す参加者間の会議が始まる以前に図 2 の処理装置 2 4 により実行される処理動作を示す図である。

【図 4】図 3 のステップ S 2 及びステップ S 4 で会議アーカイブデータベース 6 0 に格納されるデータを概略的に示す図である。

【図 5】図 3 のステップ S 3 4 及び図 7 のステップ S 7 0 で実行される処理動作を示す図である。

10

【図 6】図 5 のステップ S 4 2 - 1、S 4 2 - 2 及び S 4 2 - n のそれぞれで実行される処理動作を示す図である。

【図 7】参加者間で会議が行われている間に図 2 の処理装置 2 4 により実行される処理動作を示す図である。

【図 8 A】図 7 のステップ S 7 2 で実行される処理動作を示す図である。

【図 8 B】図 7 のステップ S 7 2 で実行される処理動作を示す図である。

【図 9 A】図 8 のステップ S 8 0 で実行される処理動作を示す図である。

【図 9 B】図 8 のステップ S 8 0 で実行される処理動作を示す図である。

【図 1 0】図 9 のステップ S 1 1 4 及びステップ S 1 2 4 で実行される処理で使用する参加者の視線を示す図である。

20

【図 1 1】図 9 のステップ S 1 1 4 で実行される処理において計算される角度を示す図である。

【図 1 2】図 8 のステップ S 8 4 で実行される処理動作を示す図である。

【図 1 3】図 8 のステップ S 8 9 で実行される処理動作を示す図である。

【図 1 4】図 1 3 のステップ S 1 6 8 で実行される処理動作を示す図である。

【図 1 5】会議アーカイブデータベース 6 0 への情報の格納を概略的に示す図である。

【図 1 6 A】会議アーカイブデータベース 6 0 に格納されたデータにより定義される視線ヒストグラムの例を示す図である。

【図 1 6 B】会議アーカイブデータベース 6 0 に格納されたデータにより定義される視線ヒストグラムの例を示す図である。

30

【図 1 7】図 8 のステップ S 1 0 2 で実行される処理動作を示す図である。

【図 1 8】会議アーカイブデータベース 6 0 から情報を検索するために処理装置 2 4 により実行される処理動作を示す図である。

【図 1 9 A】図 1 8 のステップ S 2 0 0 でユーザに対し表示される情報を示す図である。

【図 1 9 B】図 1 8 のステップ S 2 0 4 でユーザに対し表示される情報の一例を示す図である。

【図 2 0】1 つのデータベースが複数の会議からの情報を格納し、1 つ又は複数の遠隔装置からこのデータベースを問い合わせる第 1 の実施形態の変形例を概略的に示す図である。

【図 2 1】第 2 の実施形態における複数の参加者の間の会議からの音声データ及び映像データの記録を示す図である。

40

【図 2 2】第 2 の実施形態の処理装置内部の概念上の機能構成要素の一例を示すブロック線図である。

【図 2 3 A】図 2 1 に示す参加者間の会議が始まる以前に図 2 2 の処理装置 6 2 4 により実行される処理動作を示す図である。

【図 2 3 B】図 2 1 に示す参加者間の会議が始まる以前に図 2 2 の処理装置 6 2 4 により実行される処理動作を示す図である。

【図 2 3 C】図 2 1 に示す参加者間の会議が始まる以前に図 2 2 の処理装置 6 2 4 により実行される処理動作を示す図である。

【図 2 4】図 2 3 のステップ S 3 0 4 で会議アーカイブデータベース 6 6 0 に格納される

50

データを概略的に示す図である。

【図 2 5】図 2 3 のステップ S 3 3 4 で実行される処理動作を示す図である。

【図 2 6】参加者間で会議が行われている間に図 2 2 の処理装置 6 2 4 により実行される処理動作を示す図である。

【図 2 7 A】図 2 6 のステップ S 3 7 2 で実行される処理動作を示す図である。

【図 2 7 B】図 2 6 のステップ S 3 7 2 で実行される処理動作を示す図である。

【図 2 8】図 2 7 のステップ S 3 8 0 で実行される処理動作を示す図である。

【図 2 9】図 2 8 のステップ S 4 1 4 で実行される処理で使用する参加者の視線を示す図である。

【図 3 0】図 2 8 のステップ S 4 1 4 で実行される処理において計算される角度を示す図である。

10

【図 3 1】図 2 7 のステップ S 3 8 4 で実行される処理動作を示す図である。

【図 3 2】会議アーカイブデータベース 6 6 0 への情報の格納を概略的に示す図である。

【図 3 3 A】会議アーカイブデータベース 6 6 0 に格納されたデータにより定義される視線ヒストグラムの例を示す図である。

【図 3 3 B】会議アーカイブデータベース 6 6 0 に格納されたデータにより定義される視線ヒストグラムの例を示す図である。

【図 3 4】会議アーカイブデータベース 6 6 0 から情報を検索するために処理装置 6 2 4 により実行される処理動作を示す図である。

【図 3 5 A】図 3 4 のステップ S 5 0 0 でユーザに対し表示される情報を示す図である。

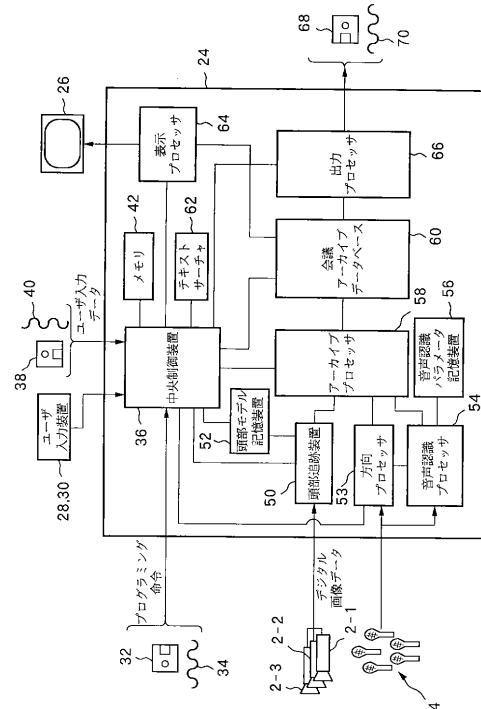
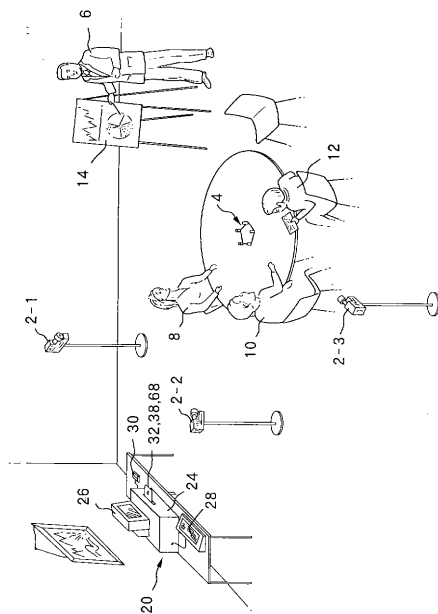
20

【図 3 5 B】図 3 4 のステップ S 5 0 4 でユーザに対し表示される情報の一例を示す図である。

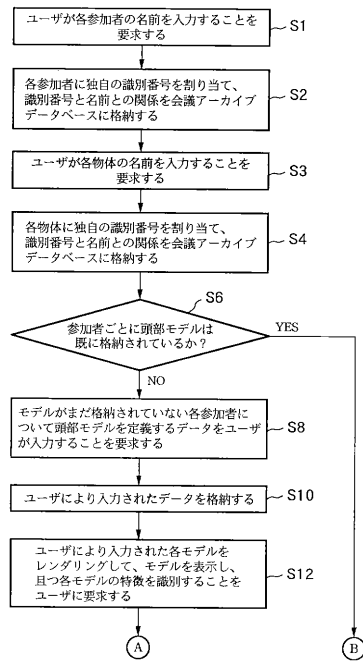
【図 3 6】1 つのデータベースが複数の会議からの情報を格納し、1 つ又は複数の遠隔装置からこのデータベースを問い合わせる第 2 の実施形態の変形例を概略的に示す図である。

【図 1】

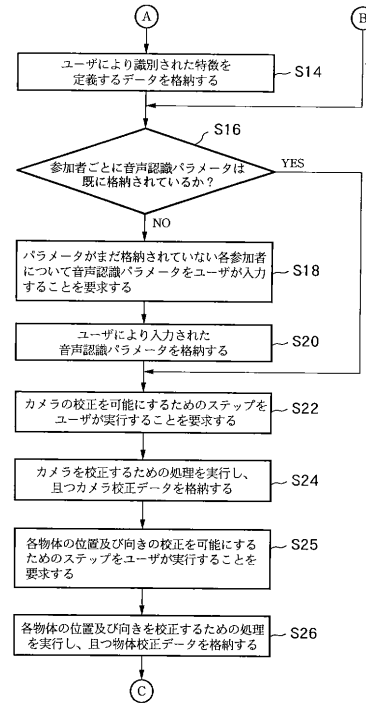
【図 2】



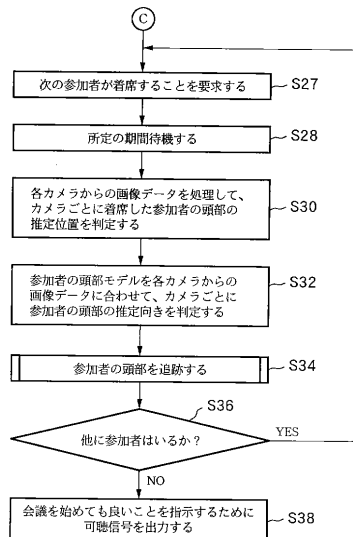
【図3A】



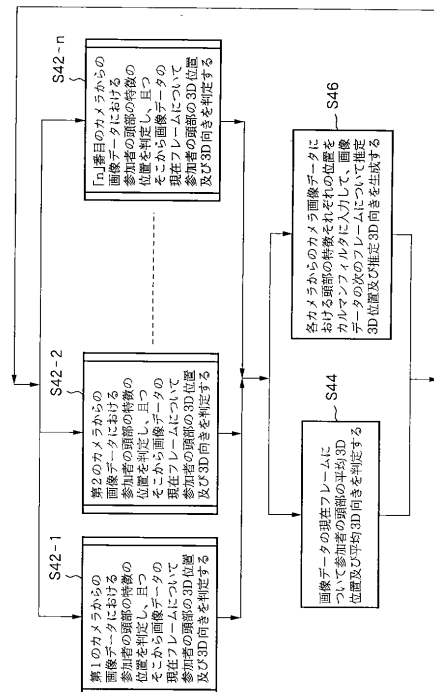
【図3B】



【図3C】



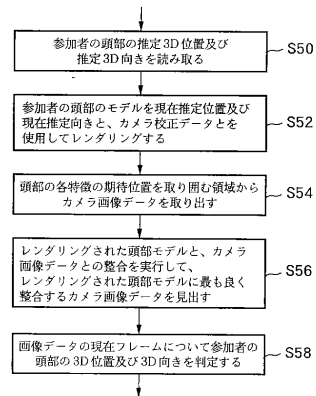
【図5】



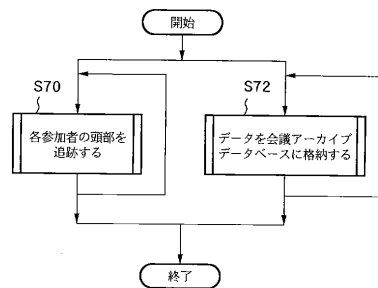
【図4】

番号	名前
1	A氏
2	Bさん
3	C氏
4	Dさん
5	フリップチャート

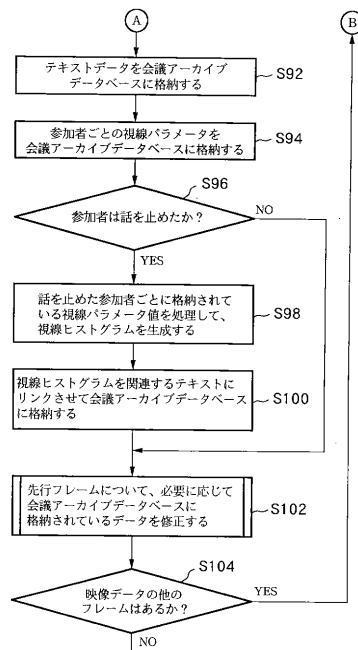
【図 6】



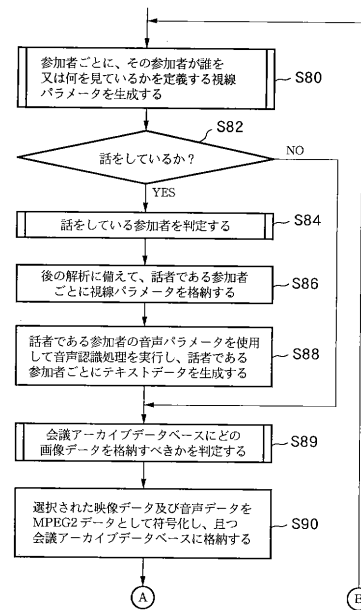
【図 7】



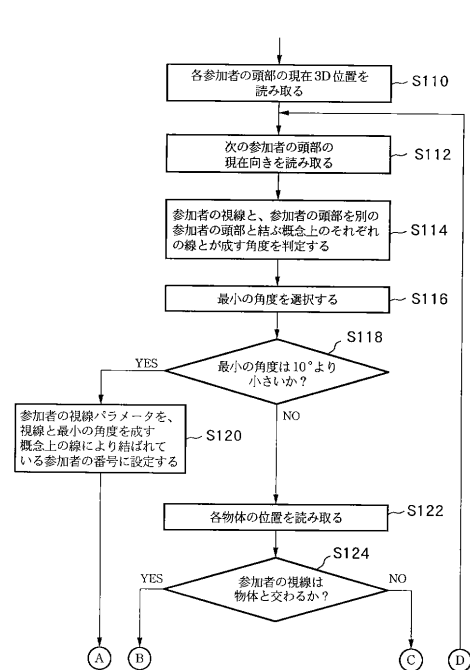
【図 8 B】



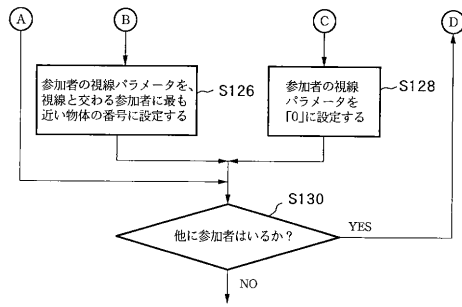
【図 8 A】



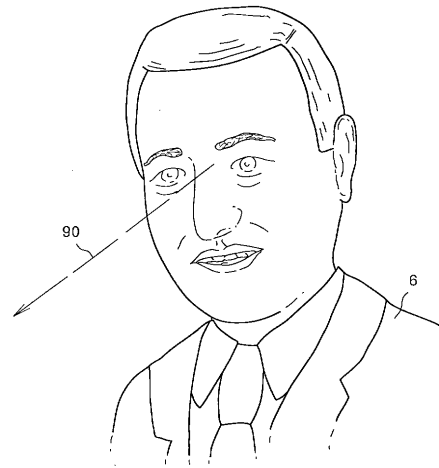
【図 9 A】



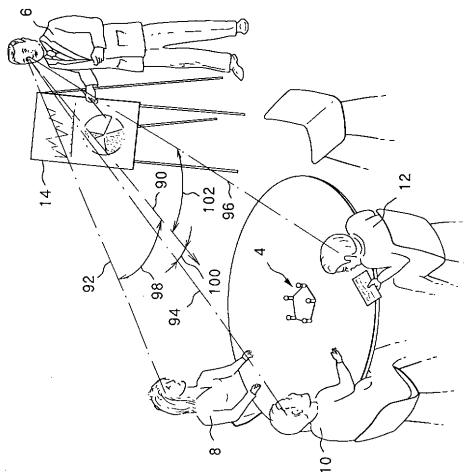
【図 9 B】



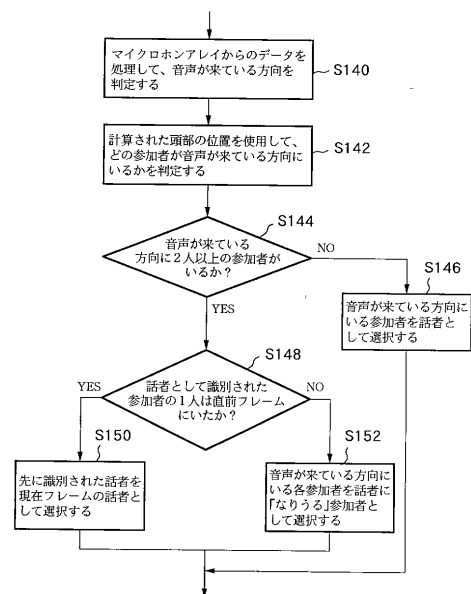
【図 10】



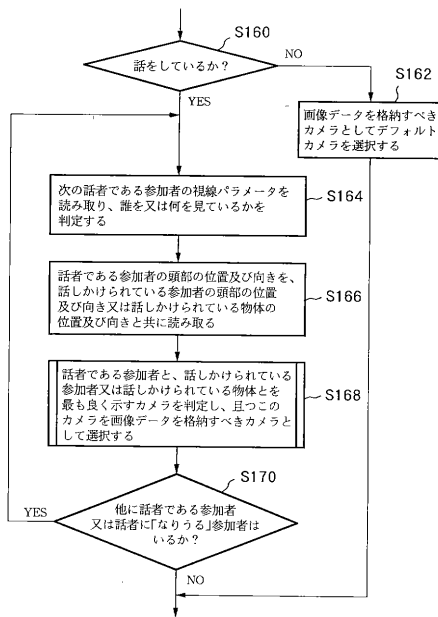
【図 11】



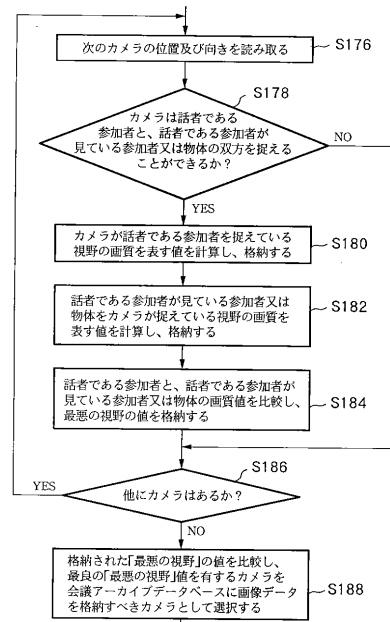
【図 12】



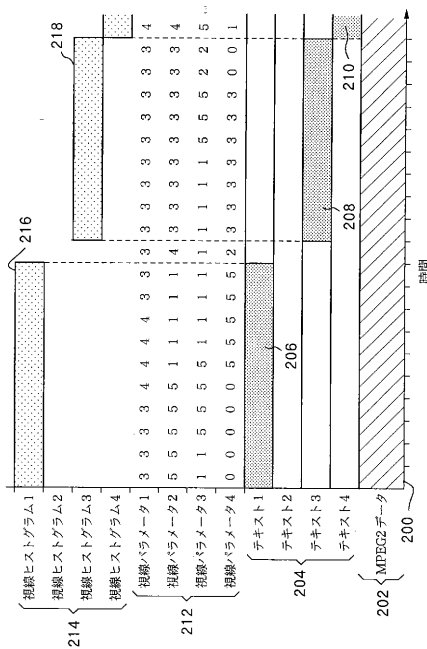
【図 13】



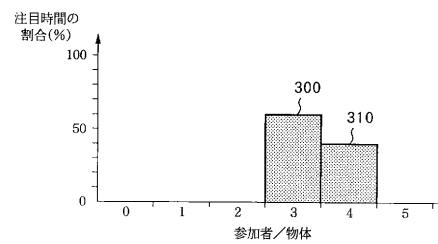
【図 14】



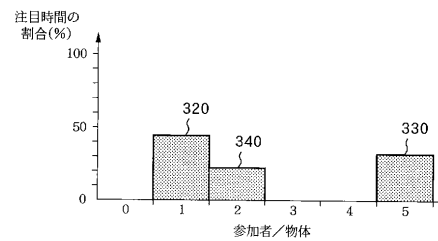
【図 15】



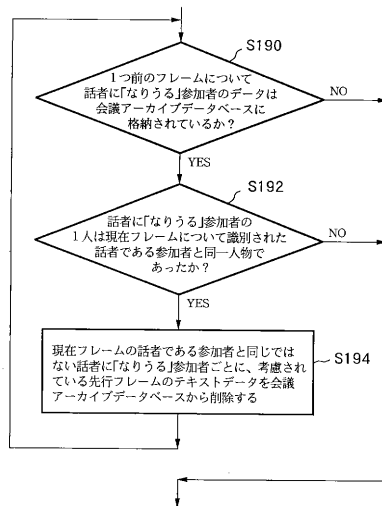
【図 16 A】



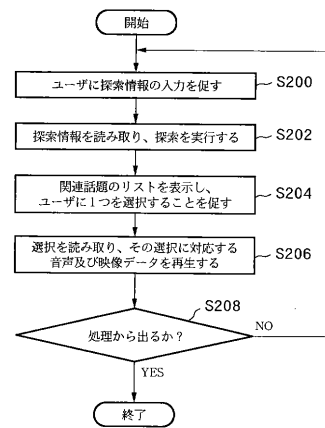
【図 16 B】



【図 17】



【図 18】



【図 19 A】

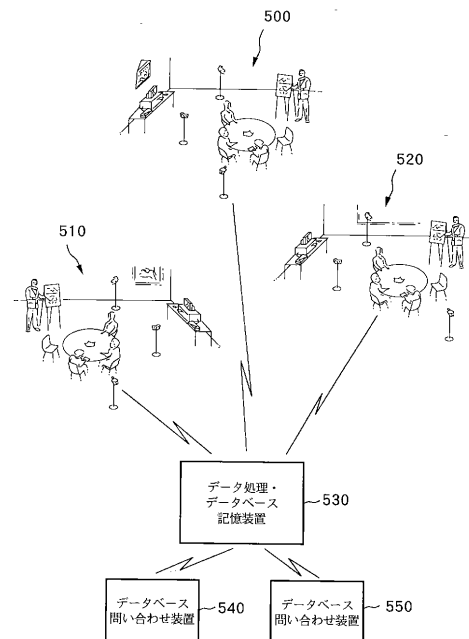
Please enter search parameters

400 talking about 410 to 420

Time limits : Before 430 After 440 Between 450 and 460

470 START

【図 20】



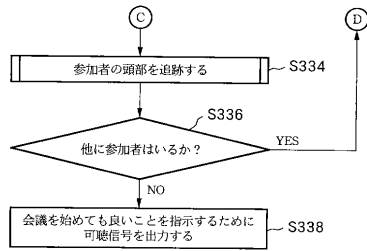
【図 19 B】

The following parts of the meeting are relevant. Please select one for playback :

1. Speech starting at 10 mins 0 secs (0.4 x full meeting time)

2. Speech starting at 12 mins 30 secs (0.5 x full meeting time)

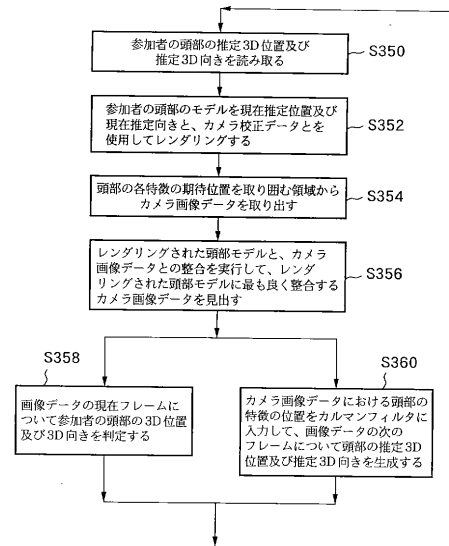
【図23C】



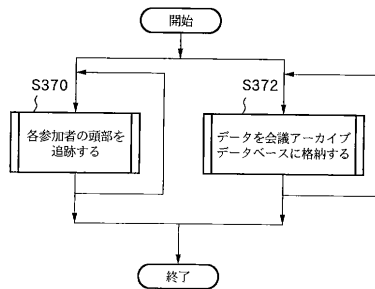
【図24】

番号	参加者	680
1	A氏	
2	B氏	
3	C氏	
4	Dさん	

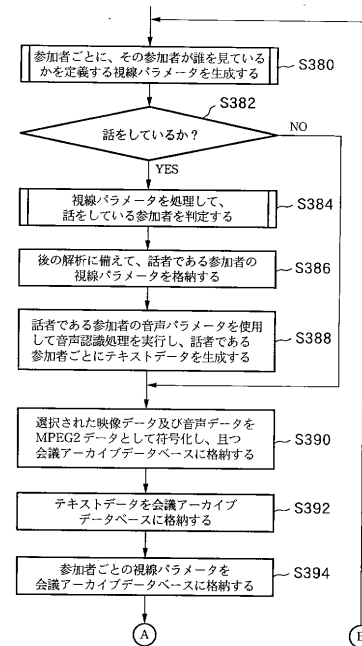
【図25】



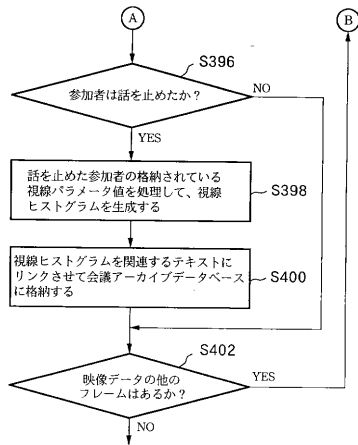
【図26】



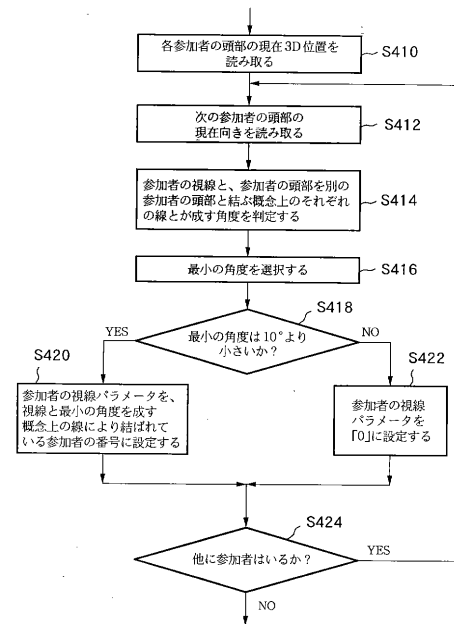
【図27A】



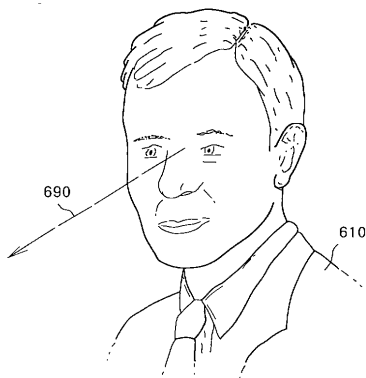
【図 27 B】



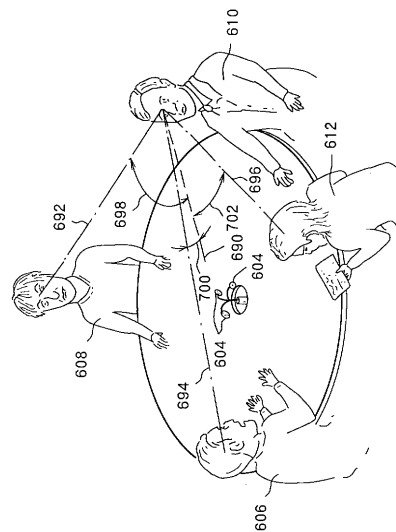
【図 28】



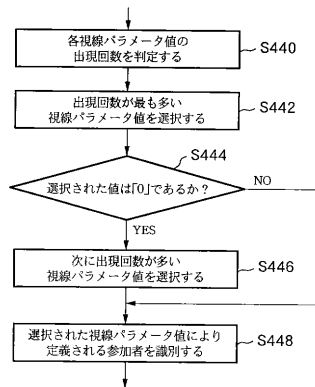
【図 29】



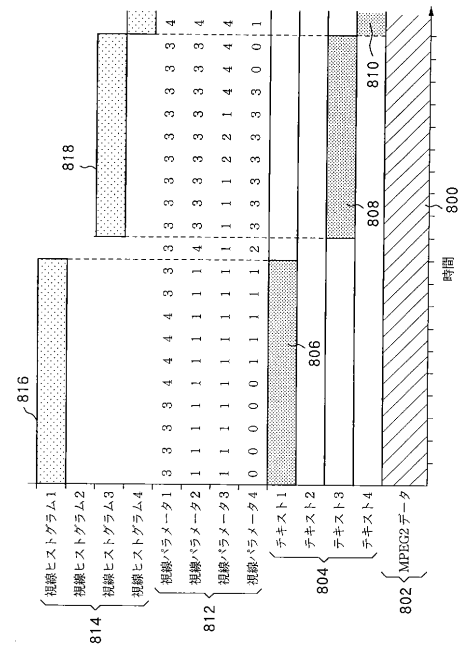
【図 30】



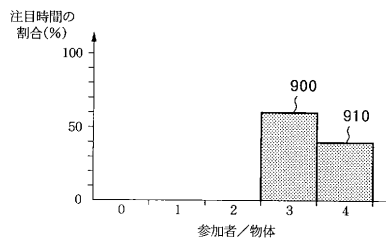
【図 3 1】



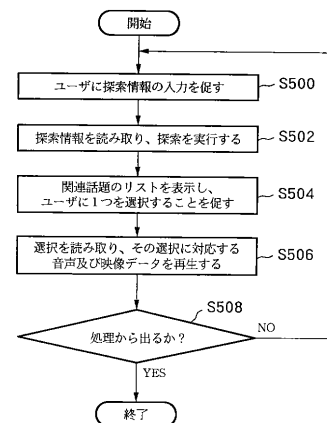
【図 3 2】



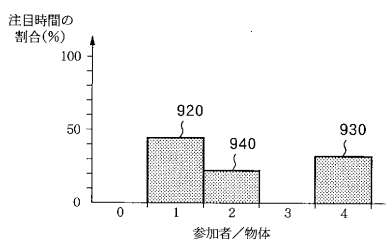
【図 3 3 A】



【図 3 4】



【図 3 3 B】



【図 35 A】

Please enter search parameters

1000 1010 1020

_____ talking about _____ to _____

Time limits :

Before _____ 1030

After _____ 1040

Between _____ and _____

1050 1060

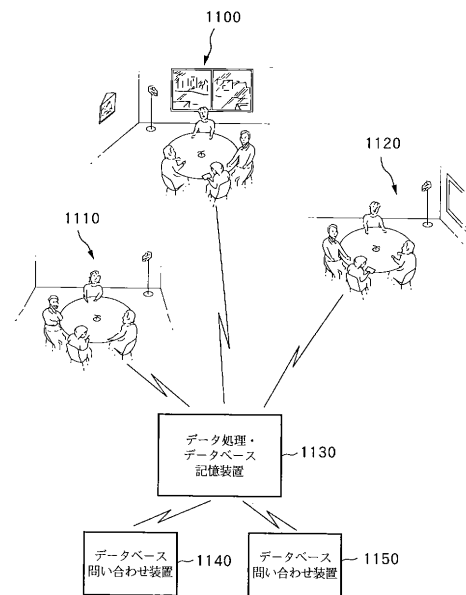
1070 **START**

【図 35 B】

The following parts of the meeting are relevant. Please select one for playback :

1. Speech starting at 10 mins 0 secs (0.4 × full meeting time)
2. Speech starting at 12 mins 30 secs (0.5 × full meeting time)

【図 36】



フロントページの続き

- (72)発明者 サイモン マイケル ロウ
イギリス国 ジーユー2 5ワイジェイ サリー, ギルドフォード, サリー リサーチ パーク, オッカム ロード, オッカム コート 1 キヤノン リサーチ センター ヨーロッパ リミテッド 内
- (72)発明者 マイケル ジェームス テイラー
イギリス国 ジーユー2 5ワイジェイ サリー, ギルドフォード, サリー リサーチ パーク, オッカム ロード, オッカム コート 1 キヤノン リサーチ センター ヨーロッパ リミテッド 内
- (72)発明者 ジェブ ジェイコブ ラジャン
イギリス国 ジーユー2 5ワイジェイ サリー, ギルドフォード, サリー リサーチ パーク, オッカム ロード, オッカム コート 1 キヤノン リサーチ センター ヨーロッパ リミテッド 内

審査官 井上 健一

- (56)参考文献 特開平10-145763(JP,A)
特開平02-206825(JP,A)
特開平05-035441(JP,A)
特開平04-082357(JP,A)
特開平03-029555(JP,A)
特開平04-181300(JP,A)
特開平11-259501(JP,A)

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