

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4986433号  
(P4986433)

(45) 発行日 平成24年7月25日(2012.7.25)

(24) 登録日 平成24年5月11日(2012.5.11)

(51) Int. Cl.	F I	
<b>G 1 0 L</b> 21/02 (2006.01)	G 1 0 L	21/02 2 0 1 B
<b>G 0 6 T</b> 7/00 (2006.01)	G 0 6 T	7/00 P
<b>H 0 4 R</b> 1/40 (2006.01)	G 1 0 L	21/02 2 0 2 A
<b>H 0 4 N</b> 5/232 (2006.01)	H 0 4 R	1/40 3 2 0 Z
<b>G 0 5 D</b> 1/02 (2006.01)	H 0 4 N	5/232 C
請求項の数 12 (全 38 頁) 最終頁に続く		

(21) 出願番号 特願2005-286754 (P2005-286754)  
 (22) 出願日 平成17年9月30日(2005.9.30)  
 (65) 公開番号 特開2006-123161 (P2006-123161A)  
 (43) 公開日 平成18年5月18日(2006.5.18)  
 審査請求日 平成18年12月27日(2006.12.27)  
 (31) 優先権主張番号 10-2004-0078019  
 (32) 優先日 平成16年9月30日(2004.9.30)  
 (33) 優先権主張国 韓国(KR)  
 (31) 優先権主張番号 10/998,984  
 (32) 優先日 平成16年11月30日(2004.11.30)  
 (33) 優先権主張国 米国(US)

(73) 特許権者 390019839  
 三星電子株式会社  
 Samsung Electronics  
 Co., Ltd.  
 大韓民国京畿道水原市靈通区三星路129  
 129, Samsung-ro, Yeon  
 gtong-gu, Suwon-si, G  
 yeonggi-do, Republic  
 of Korea  
 (74) 代理人 100070150  
 弁理士 伊東 忠彦  
 (72) 発明者 崔 昌 圭  
 大韓民国 ソウル特別市 鍾路區 舊基洞  
 62-1番地

最終頁に続く

(54) 【発明の名称】 物体を認識および追跡する装置及び方法

(57) 【特許請求の範囲】

【請求項1】

受信した音および映像を使用して物体を認識および追跡する装置において、  
 異なる方向から受信した複数のサウンドのそれぞれに基づいて、前記サウンドが観測され  
 された際に、ある方向に追跡する物体が存在することの尤もらしさを表す音声尤度を求める  
 音声尤度モジュールと、

映像内の異なる方向に配置された複数のイメージのそれぞれに基づいて、前記映像内の  
 イメージが観測された際に、ある方向に追跡する物体が存在することの尤もらしさを表す  
 映像尤度を求める映像尤度モジュールと、

前記音声尤度において尤もらしいとする方向と、前記映像尤度において尤もらしいとさ  
 れる方向とが一致するかどうかを判断し、一致すると判断されれば、前記音声尤度および  
 映像尤度の対を使用して、前記物体を認識して追跡し、一致しなければ、前記音源または  
 イメージ源には、追跡される物体が存在しないと判断する認識および追跡モジュールと、  
 を備え、

前記認識した物体それぞれに対して前記認識した物体それぞれに一意に該当するオーデ  
 ィオチャンネルを出力するために、前記サウンドから認識した物体の位置に該当する音声  
 を分離するビームフォーマを更に備え、

前記音声尤度モジュールは、それぞれの受信したサウンドに基づいて、その音声方向を  
 更に検出し、

前記映像尤度モジュールは、それぞれの観測されたイメージに基づいて、その映像方向

10

20

を更に検出し、

前記認識および追跡モジュールは、前記音声方向と映像方向とに基づいて、前記サウンドとイメージとの方向が一致するかどうかを更に判断し、

前記ビームフォーマによって出力されたオーディオチャンネルのそれぞれは、音声を検出される聴き取り期間とその聴き取り期間の間で音声を検出されない静寂期間を検出し、前記出力されたオーディオチャンネルに対して、それぞれの検出された聴き取り期間に対する開始および終了時間を検出する音声期間検出器を更に備え、

前記音声期間検出器は、隣接する前記聴き取り期間の間の近接性を検出し、前記近接性が所定の値より小さければ、前記隣接聴き取り期間を一つの連続した聴き取り期間として決定し、前記隣接聴き取り期間を連続的な聴き取り期間として形成するために連結し、あるいは、前記近接性が所定の値より大きければ、前記隣接聴き取り期間は静寂期間により分離されると決定し、前記隣接聴き取り期間を連結しないことを特徴とする装置。

10

【請求項2】

前記装置は、受信した第1個数の受信されたオーディオチャンネルを出力するマイクロフォンアレイを使用して前記サウンドを受信し、前記受信されたオーディオチャンネルは前記サウンドの要素を含み、前記ビームフォーマは、前記第1個数と異なる第2個数のオーディオチャンネルを出力し、前記第2個数は、認識した物体の個数に該当することを特徴とする請求項1に記載の装置。

【請求項3】

認識した物体のそれぞれに対し、前記ビームフォーマによって出力されたオーディオチャンネルを各物体と関連して分離されたオーディオトラックとして記録する記録装置を更に備えることを特徴とする請求項2に記載の装置。

20

【請求項4】

前記音声期間検出器は、前記聴き取り期間のそれぞれの長さを検出し、前記長さが所定の値より短ければ、前記聴き取り期間を静寂であると決定して、前記聴き取り期間を削除し、あるいは、前記長さが所定の値より長ければ、前記聴き取り期間が静寂期間ではないと決定し、前記聴き取り期間を削除しないことを特徴とする請求項1に記載の装置。

【請求項5】

前記音声期間検出器は、それぞれの聴き取り期間に対して前記検出した音声出力して、前記それぞれの静寂期間に対して、前記オーディオチャンネルから前記サウンドを削除することを特徴とする請求項1に記載の装置。

30

【請求項6】

前記ビームフォーマから受信した前記複数のオーディオチャンネルのそれぞれに対し、他のオーディオチャンネルの干渉により発生するクロスチャンネル干渉に該当する音声部分を検出し、かつ前記クロスチャンネル干渉を除去する後処理装置を更に備えることを特徴とする請求項1に記載の装置。

【請求項7】

音声および映像データを受信する少なくとも一つのコンピュータを使用して、物体を追跡および認識する方法において、

異なる方向から受信した複数のサウンドのそれぞれに基づいて、前記少なくとも一つのコンピュータで、前記サウンドが観測された際に、ある方向に追跡される物体が存在することの尤もらしさを表す音声尤度を求めるステップと、

40

前記映像内の異なる方向に配置された複数イメージのそれぞれに基づいて、前記少なくとも一つのコンピュータで、前記映像内のイメージが観測された際に、ある方向に追跡される物体が存在することの尤もらしさを表す映像尤度を求めるステップと、

前記音声尤度において尤もらしいとする方向と、前記映像尤度において尤もらしいとされる方向とが一致するかどうかを判断し、一致すると判断されれば、前記少なくとも一つのコンピュータで、前記音声尤度および映像尤度の対を使用して、前記物体のうち該当する一つを認識および追跡するステップと、

前記音声尤度において尤もらしいとする方向と、前記映像尤度において尤もらしいとさ

50

れる方向とが一致しなければ、前記少なくとも一つのコンピュータで、前記音源またはイメージ源は、追跡する物体ではないと認識するステップと、を含み、

前記認識した物体のそれぞれに対して、前記認識した物体のそれぞれに位置を決定することでビームフォーミングを行うステップと、前記認識した物体のそれぞれに対して一意に該当するオーディオチャンネルを出力するために、前記それぞれ認識された物体の位置に該当する音声を前記受信サウンドから分離するステップを更に含み、

前記音声尤度を求めるステップは、それぞれの受信したサウンドに基づいて、その音声方向を更に検出するステップを含み、

前記映像尤度を求めるステップは、それぞれの観測されたイメージに基づいて、その映像方向を更に検出するステップを含み、

10

前記認識および追跡するステップは、前記音声方向と前記映像方向とに基づいて前記サウンドと前記イメージとの方向が一致するかどうかを更に決定するステップを含み、

前記音声を受信サウンドから分離するステップは、音声が検出される聴き取り期間とその聴き取り期間の間で音声が検出されていない静寂期間を検出し、前記出力されたオーディオチャンネルに対して、それぞれの検出された聴き取り期間に対する開始および終了時間によりスピーチインターバルを検出するステップを含み、

前記スピーチインターバルを検出するステップは、隣接する前記聴き取り期間の間の近接性を検出するステップと、前記近接性が所定の値より小さければ、前記隣接する聴き取り期間を一つの連続的な聴き取り期間と決定し、前記隣接聴き取り期間を連結して連続的な聴き取り期間を形成するステップと、前記近接性が所定の値より大きければ、前記隣接する聴き取り期間を前記静寂期間により分離されると決定し、前記隣接する聴き取り期間を連結しないステップと、を含むことを特徴とする方法。

20

#### 【請求項 8】

前記少なくとも一つのコンピュータは、受信した第 1 個数の受信されたオーディオチャンネルを出力するマイクロフォンアレイを使用して前記サウンドを受信し、前記受信されたオーディオチャンネルは、前記サウンドの要素を含み、前記ビームフォーミングステップは、前記第 1 個数と異なる第 2 個数のオーディオチャンネルを出力するステップを含み、前記第 2 個数は、認識した物体の個数に該当することを特徴とする請求項 7 に記載の方法。

#### 【請求項 9】

認識した物体のそれぞれに対し、前記ビームフォーマによって出力されたオーディオチャンネルを各物体と関連して分離されたオーディオトラックとして保存するステップを更に含むことを特徴とする請求項 8 に記載の方法。

30

#### 【請求項 10】

前記スピーチインターバルを検出するステップは、前記各聴き取り期間の長さを検出するステップと、前記長さが所定の値より短ければ、前記聴き取り期間を静寂期間であると決定し、前記聴き取り期間を削除するステップと、

前記長さが所定の値より長ければ、前記聴き取り期間が静寂期間ではないと決定し、前記聴き取り期間を削除しないステップと、を含むことを特徴とする請求項 7 に記載の方法

40

#### 【請求項 11】

前記スピーチインターバルを検出するステップは、前記各聴き取り期間に対し、前記検出された音声を出力するステップと、前記各静寂期間に対し、前記オーディオチャンネルから前記サウンドを削除するステップと、

を含むことを特徴とする請求項 7 に記載の方法。

#### 【請求項 12】

複数のビームが形成されたオーディオチャンネルのそれぞれに対し、他のオーディオチャンネルの干渉により発生するクロスチャンネル干渉に該当する音声部分を除去し、前記

50

クロスチャンネル干渉を除去することで前記ビーム形成されたオーディオチャンネルを後処理するステップを含むことを特徴とする請求項7に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、対象物を検出する装置およびその方法に関し、より詳細には、同時に空間上に分離している複数の対象物をオーディオおよびビデオセンサーを用いて検出し、その対象物の位置を把握して追跡を可能にする方法および装置に関する。

【背景技術】

【0002】

一般に、対象物を検出する場合、従来の装置および方法は、映像信号または音声信号に依存している。音声追跡に対しては、TDE (time-delay estimates) が使用されている。しかしながら、周囲のノイズおよび反響に対処するために、最尤法 (maximum likelihood approach) および位相変換から重み関数を求めても、TDEに基づく技術は、明示的な方向性ノイズに対して脆弱である。

【0003】

一方、映像追跡に対しては、物体検出は非特許文献1に記載されているハウスドルフ距離 (Hausdorff Distance) を使用して、イメージを比較することで行われる。この方法は、スケールや変形などで簡単、かつ強固に処理できるが、多様なサイズについて任意の候補イメージを比較する場合に相当な時間を要する。

【0004】

また、他の対象物から出るスピーチ/サウンドが重なる場合に、対象物を分離および検出することに問題が生じる。重なった音声においては、非特許文献2で提示したように、音声を話者ごとに分割することが重要な位置を占める。マイクロフォンアレイを利用した重なった音声の分割に対する結果は、両耳暗黙信号分離、二重スピーカーHMM (Hidden Markov Model) およびTDE (Time Delay Estimates) をもって話者位置をモデリングするために、ガウス分布で具体化されたスピーチ/サイレント比率を使用することで得られる。このような結果の例は、非特許文献3、非特許文献4および非特許文献5に記載されている。5つのビデオストリームおよびマイクロフォンアレイから出るパノラマイメージを利用した話者追跡に関しては、非特許文献6および非特許文献7に記載されている。これら方法は、同時話者分割の両極端にあり、一方は、音声情報のみに依存する手法であり、他方は大部分を映像信号に依存している。

【0005】

しかし、前記の何れの接近法も、重なった音声を分離するために映像および音声入力を効果的に使用していない。さらに、Y. ChenおよびY. Ruiにより開示された方法は、受信したすべての音声データを記録するため、多くのメモリ容量を使用し、分離した音声がどの話者から出たかを識別するように、映像と音声入力を使用して同時に発生した複数の音声を個々に分割することは不可能である。

【非特許文献1】D. P. Huttenlocher, G. A. Klanderman and W. J. Rjcklidge, "Comparing Image using the Hasusdorff Distance under Translation," in Proc. IEEE Int. Conf. CVPR, 1992, pp. 654 to 656

【非特許文献2】E. Shirberg, A. Stolcke and D. baron, "Obeservations on Overlap: Findings and Implications for Automatic Processing of Multi-Party Coversation" Proc. Eruospech, 2001

【非特許文献3】C.Choi, "Read-tiem Binaural Blind Source Separation" Proc. Int. Symp. ICA and BSS, pp. 567 to 572

【非特許文献4】G. Lathoud and I. A. McCowna, "Location Based Speaker Segmentati on," Proc. ICASSP, 2003

【非特許文献5】G. Lathoud, I. A. MCCowan and D. C. Moore, "Segmenting Multiple Concurrent Speakers using Microphone Arrays," Proc.Eurospech,2003

10

20

30

40

50

【非特許文献6】R. Cutler et. al., "Distributed Meetings: A Meeting Capture and Broadcasting System", Proc. ACMInt. Conf. Multimedia, 2002

【非特許文献7】Y. Chen and Y. Rui, "Realtime Speaker Tracking using Particle filter Sensor Fusion", Proc. of the IEEE, vol.92 No. 3, pp. 485 to 494, 2004

【発明の開示】

【発明が解決しようとする課題】

【0006】

本発明が達成しようとする技術的課題は、オーディオおよびビデオセンサーを融合して、複数の音源として存在する特定の物体について正確な位置および方向を把握して追跡し、この音源を分離する装置および方法を提供することにある。

10

【課題を解決するための手段】

【0007】

前記技術的課題を解決するための本発明に係る装置は、受信した音および映像を使用して物体を認識および追跡する装置において、異なる方向から受信した複数のサウンドのそれぞれに対して、前記サウンドが追跡する物体のものである尤度を表す音声尤度を求める音声尤度モジュールと、映像内の異なる方向に配置された複数のイメージのそれぞれに対し、前記映像内の前記イメージが追跡される物体である尤度を表す映像尤度を決定する映像尤度モジュールと、前記音声尤度と前記映像尤度とが一致するかを判断し、前記音声尤度と前記映像尤度とが一致すると判断すれば、前記音声尤度および映像尤度の対を使用して、前記物体を認識して追跡し、前記音声尤度と前記音声尤度が一致しなければ、前記音源またはイメージ源には、追跡される物体が存在しないと判断する認識および追跡モジュールと、を備える。

20

【0008】

前記の技術的課題を解決するための本発明に係る方法は、音声および映像データを受信する少なくとも一つのコンピュータを使用して、物体を追跡および認識する方法において、異なる方向から受信した複数のサウンドのそれぞれに対し、前記少なくとも一つのコンピュータで、前記サウンドが追跡される物体のものである尤度を表す音声尤度を求めるステップと、前記映像内の異なる方向に配置された複数イメージのそれぞれに対し、前記少なくとも一つのコンピュータで、前記映像内の前記イメージが追跡される物体である尤度を表す映像尤度を求めるステップと、前記音声尤度と前記映像尤度とが一致するかを判断し、前記音声尤度と前記映像尤度とが一致すると判断すれば、前記少なくとも一つのコンピュータで、前記音声尤度および映像尤度の対を使用して、前記物体のうち該当する一つを認識および追跡するステップと、前記音声尤度と前記映像尤度とが一致しなければ、前記少なくとも一つのコンピュータで、前記音源またはイメージ源は、追跡する物体ではないと認識するステップと、を含む。

30

【発明の効果】

【0009】

本発明の一実施形態に係る方法は、他の方法に比べて幾つかの長所を有する。第一に、前記方法が精巧に測定された調整ベクトルを有するサブ空間方法が全体システムに具体化されるため、ノイズ耐性が強い。第二に、人の上半身に対する三つの形状モデルに基づくものであり、人をその上半身で認識することは、体全体で識別するよりも好都合である。なぜなら、散乱した環境においては、人の下半身は他の対象と度々誤認されるためである。しかしながら、環境によっては下半身を人の認識に使用してよいことは言うまでもない。第三に、前記方法は、プロファイルを人間形状モデルとして採用するため、姿勢予測が可能である。このような姿勢情報は、特に、特定フィルタリングに有効であるが、他の方法を使用してもよい。第四に、調整ベクトル不一致に強いが、実際調整ベクトルが現実で利用できなくても、対象物の音声を相殺する問題点は、対角線ローディング方法を有するターゲットフリー共分散マトリックスにより克服し、順に、これは、本発明の一実施形態によって提供される正確なセグメンテーションにより可能である。

40

50

## 【 0 0 1 0 】

また、本発明の一実施形態に係る装置の長所は、直観的でかつ簡単なセンサーの混合ストラテジーであり、ここで前記装置は、オーディオビデオのセンサー混合を利用して所望の対象をより正確に追跡するために、活動する話者から拡声器および人の絵を分離して維持できる。さらに、その性能は、適応クロスチャンネル干渉相殺により更に向上するため、その結果は、大容量語彙の連続音声認識システム、または自動会議記録を行う遠距離トークのために使用される書き取り機械に直接適用可能である。さらに、この装置は音声向上器だけでなく、終点検出器としても作動する。しかし、他の実施形態および長所は、前記内容から理解できるということは言うまでもない。

## 【 0 0 1 1 】

また、如何なる実施形態で要求されるものではないが、図 2 に示す方法またはその一部分は、少なくとも汎用または専用のコンピュータを使用して、コンピュータで読み取り可能な記録媒体にコード化された一つ以上のコンピュータプログラムを使用して実装される。また、カメラを利用した視覚追跡について記載したが、360°パイロセンサーのような検出手段を使用してもよい。

## 【 発明を実施するための最良の形態 】

## 【 0 0 1 2 】

本発明の実施形態の構成および動作を、添付図面を参照して詳細に説明する。図面の構成要素に参照番号を付与するにあたり、同一構成要素に対しては他の図面上にあって同一参照番号を付与した。

## 【 0 0 1 3 】

図 1 は、本発明の一実施形態によって、位置把握および追跡能力のあるオーディオおよびビデオセンサーを有するロボットのブロック図を示す図面である。

図 1 を参照すると、ロボットは、映像システム 100 と、音声システム 200 と、コンピュータ 400 とから構成される。本発明の如何なる態様において必要とするものではないが、図 1 に示すロボットは、映像システム 100 および音声システム 200 からの入力によりコンピュータ 400 によって制御されるロボット要素 300 を更に備える。ロボット要素 300 において、映像システム 100 および音声システム 200 は、コンピュータ 400 と共に統合される必要はなく、別個に配置されてもよいということは言うまでもない。

本発明の一実施形態に係る装置はロボットであり、未知の環境で動作または停止が可能である。ロボットは、周囲環境で観察される特徴を収集し、かつその制御を行うことができる。制御および観察順序に基づき、本発明の実施形態に係るロボットは、少なくとも 1 つの対象物を検出し、その位置を把握して追跡でき、更に複数の対象物を追跡し、かつ応答することができる。本発明の他の実施形態において、ロボットは複数の物体のうち、各対象物の様相、すなわち、目標話者のそれぞれの音声および顔に基づいて様相を分離できる。更に他の本発明の実施形態において、物体とロボットは  $x - y$  平面にあると見なされる。しかし、この方法は、本発明の一実施形態において、3次元空間まで容易に拡張できることは言うまでもない。

## 【 0 0 1 4 】

本実施形態に係る装置および方法は、視覚および音声上の特徴を有する複数の物体を追跡かつ隔離することによって、航空機、自動車、船などの衝突防止や、これらのナビゲーションを行うために用いることができ、固定された設備に適用する事も可能である。

## 【 0 0 1 5 】

映像システム 100 は、全方向カメラ 110 を備える。全方向カメラ 110 の出力は、USB 2.0 インターフェース 120 を介してコンピュータ 400 と連結される。図示するように、全方向カメラ 110 は、図 3 A に示すような 360°視野を提供する。しかし、カメラ 110 は、視野が 180°未満のテレビ会議用のビデオカメラのように、更に視野が制限されていてもよい。あるいは、360°パイロセンサーのような赤外線検出手段を備えたカメラでもよい。また、複数の視野が制限されたビューカメラおよび/または全

10

20

30

40

50

方向カメラを、図3Aに示すような1つの平面および追加された平面で視野を広げるために使用できる。さらに、本発明の一実施形態によって、インターフェースの他のタイプを、USB 2.0インターフェース120の代わりに用いたり、あるいは追加してたりして使用してもよい。また、コンピュータ400との接続は有線または無線接続で確立することができる。

**【0016】**

音声システム200は、8個のマイクロフォン210を有するマイクロフォンアレイを備える。8個のマイクロフォンは、カメラ110の中心点を含む装置の中心点に比例する角関数として均等に空間上に広がるように、カメラ110の中心を含む中心位置周辺に45°の間隔でそれぞれ配置される。しかし、他の構成も可能であり、中心位置にマイクロフォンが配置されなくてもよく、その代わりに、所定の位置の空間壁に位置してもよい。マイクロフォン210の数は実施形態に応じて異なってよく、マイクロフォン210は、実施形態に応じて異なる角度に配置されてよい。

10

**【0017】**

それぞれのマイクロフォン210は、それぞれの対応するチャンネルに出力をする。それにより、図1に示すマイクロフォンアレイは、8個のアナログオーディオデータチャンネルを出力する。アナログ-デジタル変換器220は、アナログオーディオデータを受信してデジタル化して、8個のデジタル化されたオーディオデータチャンネルを提供する。デジタル化された8個のチャンネルオーディオデータは、変換器220から出力され、USBインターフェース230を介してコンピュータ400で受信される。本発明の実施形態によって、インターフェースの他の形態がUSBインターフェース230に代わるか、または追加して使用されてもよく、有線および/または無線で接続される。また、1つ以上のマイクロフォン210は、該当するデジタルオーディオチャンネルを直接出力できる場合（すなわち、デジタルマイクロフォンを使用した場合）、アナログ-デジタル変換器220を使用する必要はない。

20

**【0018】**

コンピュータ400は、以下で説明する本発明の一実施形態において、図2に示す方法を実行する。本発明の実施形態に係るコンピュータ400は、ペンティアム（登録商標）IV 2.5GHzのシングルボードコンピュータである。しかし、本発明の実施形態において、異なるタイプの汎用または専用コンピュータ、あるいは複数のコンピュータおよびプロセッサを使用して実装可能なことは言うまでもない。

30

**【0019】**

図1に示す一実施形態で、前記装置は、検出された対象物に反応して動くロボットと共に使用される。コンピュータ400の出力は、RS232Cインターフェース330からモータコントローラ320を経てロボット要素300に提供される。モータコントローラ320は、コンピュータ400の指示によってロボットが動くように、2つのモータ310を制御する。このような方法で、コンピュータ400は、コンピュータ400により処理された音声および映像データによって、他の対象物と区別される認識音声によって特定の対象物を追跡するようにロボットを制御できる。しかし、異なる個数のモータを、ロボットの機能により使用してよいことは言うまでもない。このようなロボットの例は限定されるものではないが、家庭用ロボット、ロボット機能付き設備、産業用ロボットおよびおもちゃなどが挙げられる。

40

**【0020】**

モータ310は、必ずしもロボットに内蔵される必要はなく、その代わりに、テレビで放送するミーティングにおけるそれぞれの話者、収録されるコンサートでの歌手、テレビ会議での話者にそれぞれをフォーカスするために、外部カメラ（図示せず）を制御するか、または商店の周りをうろつく不信人物や侵入者を検出するために、家または事業場の保安システムから検出された物体に焦点を合わせて追跡するように外部カメラを制御してよい。

**【0021】**

50

図2は、本発明の1実施形態によって、コンピュータ400により行われる方法をフローチャートで示す。ビデオカメラ110の入力は、映像システム100から受信され、コンピュータ400は、後ほど詳述する数式27を利用して、多数の人を視覚的に検出する(ステップ500)。受信されたイメージからコンピュータ400は、それぞれの潜在的対象物600、610、620、630、640が、下式28を使用して、「人」らしさを表す尤度を計算する(ステップ510)。

#### 【0022】

例えば、図3Aに示すように、受信映像イメージとして、追跡される複数の潜在的対象物600、610、620、630、640が存在する。図示する例で対象物は、人であるかどうかあらかじめ選択される。第一対象物600はオーディオスピーカーであって、これはオーディオノイズを提供するが、人として識別される映像入力イメージは提供しない。対象物620、630、640は、何れも潜在的な人であり、それぞれを、コンピュータ400によって追跡する必要がある。対象物610は、人間の可能性がある形状をした視覚ノイズを提供する写真であるが、これは、コンピュータ400により音声ノイズを提供するものではないと解釈される。

#### 【0023】

図3Aのイメージは、図3Bおよび図3Cに示す二つのサブイメージに分割される。図3Bのエッジイメージは、図3Aの写真から検出される。図示する例で、エッジイメージは、人の上半身の所定形態だけでなく、所定個数の姿勢に基づく。これについては、以下で詳細に説明する。図3Bに示すように、人の上半身は、写真610および対象物620、630、640のエッジイメージとして示されるが、対象物600のエッジイメージについてははっきりと認識できない。これにより、コンピュータ400は、写真610および対象物620、630、640に対するエッジイメージが、図4Bに示す映像尤度グラフから分かるように、人であることをより容易に検出する。

#### 【0024】

さらに正確に人を検出するために、第二サブイメージが、本発明に係る実施形態によって使用される。特に、コンピュータ400は、人間と人間ではないものとを区別するために、色(すなわち、皮膚色)を検出する。図3Cに示すように、コンピュータ400は、対象物620、630、640を人間として識別できる可能性を高めるために、対象物620、630、640の皮膚色のよう、皮膚色に基づいて顔および手を把握する。

#### 【0025】

皮膚色は、写真に対するプロブ(blob)に基づくが、プロブはコンピュータ400により写真が人間として認識できる可能性を高める。しかし、オーディオスピーカー600は、図3Cで有効な皮膚色が足りなく、図3Bで準拠しないエッジイメージを有するため、人として登録されない。

#### 【0026】

また、如何なる実施形態で要求される訳ではないが、第二サブイメージで認識された特徴は、検出されたエッジイメージがあらかじめ選択されたエッジイメージと非常に密接に一致するように、第1サブイメージのエッジイメージを正規化するのに使用される。例として、図3Cに示すプロブの位置は、図3Bのエッジイメージで手と顔の位置と、あらかじめ選択されたエッジイメージのサイズが非常に近似して一致するように、コンピュータ400に保存された身体および姿勢のイメージを一致させるために使用されるか、または図3Bおよび図3Cに示す第一サブイメージおよび第二サブイメージを2つとも使用して、検出結果を向上させるのに使用される。

#### 【0027】

したがって、ステップ510でコンピュータ400は、図3Bに示すエッジイメージおよび図3Cに示すプロブイメージに基づいて映像尤度を計算し、後ほど詳述する相対度の関数として、図4Bに示す結合された映像尤度イメージを生成する。特に、図4Bに示すように、コンピュータが識別した対象物620、630、640および写真610は、追跡される可能性のある人として何れも識別されるが、オーディオスピーカー600は、追

10

20

30

40

50

跡される人（対象）として識別されない。

【0028】

図2の方法を使用して音声尤度を決定するために、音声システム200からコンピュータにより入力されるマイクロフォンアレイは、ノイズの位置を決定するために、ビームフォーミング技術を使用して受信角度の関数としてノイズを追跡するが、これについては後ほど詳述する。受信された音声データは、後記の数式20を利用して単一のサブ空間で計算され（ステップ520）、オーディオデータが人間である尤度は、後記の数式26を利用して決定される（ステップ530）。

【0029】

図4Aに示すように、例えば、コンピュータ400は、オーディオスピーカーだけでなく、対象物630、640もノイズを提供すると認識する。コンピュータ400は、オーディオスピーカー600および対象物630、640は追跡すべき潜在的な人間として認識する。

ステップ540で、コンピュータ400は、後記の数式31を利用して、ステップ530で検出された音声対象物、およびステップ510で検出された映像対象物が追跡されるべき人であるか否かを決定するために、映像および音声尤度を組み合わせる。さらに、映像および音声尤度は方向情報を含むため、それぞれの対象物は位置関数として認識される。

【0030】

図4Cに示す例のように、コンピュータ400は、ステップ530を行うことにより、現在話している人である対象物630、640を区別できる。それぞれの対象物630、640は、位置により識別され、前記位置は、角の位置に図示されているが、本発明の他の実施形態で異なって認識される可能性がある。オーディオスピーカー600は、ステップ500および510で検出される場合に、人であるという高い映像尤度を示さないため、図4Cに示されていない。対話していない対象物620と対話できない対象物610は、コンピュータ400により追跡される人である高い尤度を有すると決定されない。

ステップ540で、音声および映像データ尤度を組み合わせることで、コンピュータ400は、後記の数式31および数式37ないし数式39を使用して、ステップ550で各人を分離して追跡できる。このような方式で、各人は、位置により個別に認識され、音声データチャンネルは特定イメージにより識別される。さらに、対象物620が話していると、分離されたトラックは出力され、対象物620と関連付けされる。

【0031】

例として、話者1、2、3が何れも図5Dに示すように話している時、コンピュータ400は、角位置の関数として話者1、2、3のそれぞれの位置を認識できる。このような既存の角位置に基づいて、コンピュータ400は、図5Aに示すような話者1を検出した第一音声トラック、図5Bに示すような話者2を検出した第二音声トラック、図5Cに示すような話者3を検出した第三オーディオトラックのように、話者1ないし3のそれぞれの角位置で音声を分離する。このような方式で、本発明に係る実施形態の方式では、残りの音声データは記録したり、または伝送する必要がなく、それゆえ、帯域幅および保存空間を低減できる。また、各トラックは、視覚対象物として認識されるため、コンピュータ400は、話している人によって分離された音声を維持できる。

【0032】

また、コンピュータ400は、本発明の実施形態において、話者1ないし3が動いても分離されたトラックを維持できる。例として、個人に対するカラーコードのために、カラーヒストグラムを使用するように、音声尤度および映像尤度の様相を認識することにより、コンピュータ400は、同一にそれぞれ割当てられたチャンネルにそれぞれの音声トラックを維持して、互いに交差しても、話者1ないし3をそれぞれ追跡できる。本発明の一実施形態によって、コンピュータ400は、図4Cおよび図5Aないし図5Cを提供するために、数式1および数式31を使用した。しかし、他のアルゴリズムおよび数式を、本発明の他の実施形態に従い使用することができ、または有効に適用することができ、また

10

20

30

40

50

、対象物が停止するか、または正確な追跡が要求されなければ、数式を簡単にすることができる。数式 1 は下記に示され、数式 2 8 を参照することで理解することができる。

【 0 0 3 3 】

【 数 1 】

$$p(z_v^i(t)) = \alpha_i N(\theta_i, \sigma_i^2) \quad (1)$$

ここで、 $p(z_v^i(t))$  は、時間  $t$  において  $i$  番目の人に対する映像尤度関数であり、 $z_v^i(t)$  は、時間  $t$  において  $i$  番目の人に対する映像観察ベクトルであり、 $\alpha_i$  は、

$$\sum_i \alpha_i = 1$$

になる 0 から 1 までの値を有する加重値であり、 $N(\theta_i, \sigma_i^2)$  は平均が  $\theta_i$  であり、分散が  $\sigma_i^2$  のガウシアン関数である。

【 0 0 3 4 】

図 5 A ないし図 5 C は、コンピュータ 4 0 0 が、各対象物 6 2 0、6 3 0、6 4 0 から発せられた音声を三つの分離されたトラックに分離する例を示す。特に、音声尤度  $L_a(\text{audio} | )$  および音源の位置に基づく音声領域は、図 5 D に示されている。このオーディオ領域で、各話者は他の角位置に位置し、話者は互いに対話しており、図 1 に示す装置を使用してその対話を収録する。追加的な背景ノイズは、他の位置に存在する。音声および映像データ尤度  $L(\text{audio}, \text{video} | )$  を組み合わせることにより、コンピュータ 4 0 0 は、各検出された話者 1 ないし 3 の相対的な角位置に基づいて個別に音声を分離できる。さらに、コンピュータ 4 0 0 は、ビームフォーミング技術を使用して、図 5 A ないし図 5 C に示すように、分離されたトラックを出力できる。また、本発明の実施形態によれば、残りの追跡されていない話者、または背景ノイズを含まずに話者 1 ないし話者 3 のそれぞれの音声のみを記録できるため、メモリ空間および伝送帯域幅を低減させ、後処理過程は、各話者の記録された音声の選択性を向上させる。このような分離は、会議、上演、ミュージカル公演の記録のような多様な状況で、話者、歌手、および/または楽器について選択されたトラックを後で増幅するときに有効である。

【 0 0 3 5 】

発明の如何なる実施形態で要求されるものではないが、様々な人の音声を記録および伝送するために、音声自体を追跡する所で、ステップ 5 6 0 でオプション信号調節動作が行われる。図示の例で、コンピュータ 4 0 0 は、数式 4 1 ないし数式 4 9 と関連して、下記で説明するように、話者に対する音声パターンをスムーズにするために、ステップ 5 6 2 で各音声トラックに対する音声存在期間 (Speech Presence Intervals: SPI) を検出する。ステップ 5 6 4 で、適応的クロスキャンセル (adaptive cross cancellation) 技術を使用して向上した各対象物から出たそれぞれの対象となる音声は、数式 5 0 ないし数式 6 5 と関連して、下記で詳細に説明される。内容を簡略化するために、コンピュータ 4 0 0 により行われることを前提に説明したが、それぞれの対象の話者が識別されれば、他のコンピュータまたはプロセッサも信号調節のための処理過程を行うことに使用可能なことは言うまでもない。

【 0 0 3 6 】

ステップ 5 6 0 において、このような信号調節は、ミーティングの議事録作成、音楽または演劇を記録、および/または高音質が要求される TV ショー、またはミーティングの記録したり、伝送するために使用される。しかし、ステップ 5 6 2 および 5 6 4 は、互いに独立に行われるか、あるいは、音質を要求しないロボット、または対象物である人の音声パターンを向上させる必要がない状況では行われなくてよい。

ステップ 5 6 2 において、人の音声パターンは、音声停止として検出される所定の減少 (dip) を含む可能性があり、その場合には、記録または伝送された音に不愉快な不連続

10

20

30

40

50

を形成してしまう。咳などによる音声の突然のスパイクは、人の音声として好ましくない。例として、図6(c)で、話者は、時間80付近で音声を休止している。そのような休止は、図6(a)および図6(b)に示す話者のパターンに図示されていない。図7(c)に示すように、このような休止は、音質を向上させるために除去されるべき音声の不連続としてみなされる。しかし、これは、対話と関連していない背景ノイズを記録しないように、音声に対する開始および終了時間を記録することが好ましい。ステップ562で音声処理を行うことにより、音声包絡線を形成して、音声の実際の休止は保存されるか、または伝送されないようにする一方、特定の人音声の終了と比較するため、コンピュータ400は、図8(c)に示すように、時間80周辺の音声を保存できる。このようなタイプの信号調節のための処理は、数式41ないし数式49と関連して以下で説明する。しかし、音声で休止または突然のスパイクが重要ではない場合は、ステップ562で示された数式49は省略することができる。

10

#### 【0037】

如何なる実施形態で要求されるものではないが、コンピュータ400は、好ましくない対象物として指名された既知の音源の音を弱める一方、特定の好ましい対象物を分離して、更に向上した音声を検出するために話者の位置を使用できる。図9に示す例として、オーディオスピーカー600は、人ではないと識別されるため、コンピュータ400は、本発明の実施形態によって、特定方向に対する利得を減らすことにより、その根源から出るノイズを除去する。対象物630、640の音声は除去されるか、または無音になる所で、コンピュータ400は、対象物630、640方向での利得を減少させることにより、対象物630、640からのノイズが効果的に除去される。さらに、対象物620からの音声またはノイズを強調するために、利得を対象物620の方向で増加させる。これにより、様々対象の音声は、ユーザーの必要性によって変化する。

20

#### 【0038】

一実施形態において、コンピュータ400は、対象物620、630、640およびオーディオスピーカー600の位置を知っているため、それらについてそれぞれの利得を操作するのにビームフォーミング技術を使用する。ビームフォーミングについての詳細な説明は、以下で提供され、ビームフォーミング技術の例は、S. Shahbazpanahi, A. B. Gershman, Z. -Q. Luo and K. Wong, "Robust Adaptive Beam-forming using Worst-case SINR Optimization: A new Diagonal Loading-type Solution for General-rank Signal", Proc. ICASSP, 2003; H. L. V. Trees, Optimum Array Processing, Wiley, 2002に記載されている。しかし、音声位置測定の方法は、本発明の実施形態で限定されるものではない。

30

#### 【0039】

図10は、図1に示す装置に組み込まれるか、または接続された後処理装置を示し、この後処理装置は、音質向上のために出力音声データをスムーズにすることに使用される。具体的には、AVシステム700は、処理されるべきオーディオおよびビデオチャンネルを受信する。如何なる実施形態で要求されるものではないが、AVシステム700は、図1の映像システム100と、音声システム200と、コンピュータ400とを備える。

AVシステム700は、音声データの分離されたトラックを出力し、ここで、それぞれのトラックは、話者から出たそれぞれの音声に該当する。出力の実施形態は、図11(a)ないし図11(c)に示されている。後処理装置710は、各トラックに含まれたオーディオノイズを除去するために、本発明の実施形態に係る適応クロスチャンネル干渉除去を行う。後処理装置710は、これらの信号を処理して、他のチャンネルの干渉を除去したそれぞれのチャンネルに対する処理された信号を出力する。これは、数式50ないし数式65と関連して以下で説明するが、さらに、C. Choi, G.-J. Jang, Y. Lee and S. Kim, "Adaptive Cross-channel Interference Cancellation on Blind Signal Separation Outputs", Proc. Int. Symp. ICA and BSS, 2004で更に詳細に説明されている。

40

#### 【0040】

図11(a)ないし図11(c)に、AVシステム700により出力された3つのチャ

50

ンネルを示す。話者 1 の音声は図 1 1 ( a ) に、話者 2 の音声は図 1 1 ( b ) に、そして、話者 3 の音声は図 1 1 ( c ) に示す。図示するように、各トラックは隣接トラックからの干渉を含む。

処理後に、後処理装置 7 1 0 は、図 1 2 ( a ) に示す話者 1 に対して処理されたトラック、図 1 2 ( b ) に示す話者 2 に対して処理されたトラック、および図 1 2 ( c ) に示す話者 3 に対して処理されたトラックを出力する。図示するように、A V システム 7 0 0 に入力される信号対ノイズ比 (Signal-to-Noise ratio: S N R) は、0 d B より小さい。図 1 1 ( a ) ないし図 1 1 ( c ) に示すように、A V システム 7 0 0 からの出力は、1 1 . 4 7 d B 程度の S N R を有する。後処理装置 7 1 0 を通過後、図 1 2 ( a ) ないし図 1 2 ( c ) に示す出力は、1 6 . 7 5 d B の S N R を有する。これにより、本発明の実施形態によって、ステップ 5 6 4 で行われた分離されたチャンネルの後処理は、隣接トラックによる干渉を除去し、出力チャンネルの記録および伝送を向上させる。

一般的に、動いている対象は、励起力および摩擦力を受ける。は、直角座標系で x、y または z を表し、極座標系で r、または z を表し、球座標系で、または を表す。座標系で、単位質量を仮定した動きの離散式は、次の数式 2 ないし数式 4 で表せる。

【 0 0 4 1 】

【数 2】

$$\zeta(t) = \zeta(t-1) + \dot{\zeta}(t) \cdot \Delta T \quad (2)$$

【 0 0 4 2 】

【数 3】

$$\dot{\zeta}(t) = \dot{\zeta}(t-1) + u'_{\zeta}(t) \cdot \Delta T \quad (3)$$

【 0 0 4 3 】

【数 4】

$$\dot{\zeta}(t) = \dot{\zeta}(t-1) + \{u_{\zeta}(t) - f(\dot{\zeta}(t))\} \cdot \Delta T \quad (4)$$

数式 2 ないし数式 4 で、t は離散時間増加値であり、T は離散時間 t の間隔であり、 $u_{\zeta}(t)$  は、外部励起力であり、

$$f(\dot{\zeta}(t))$$

は、摩擦力である。

$$f(\dot{\zeta}(t))$$

が、線形であると仮定すれば、摩擦力は、

$$b \dot{\zeta}$$

に近似する。ここで、b は、摩擦定数である。したがって、数式 2 ないし数式 4 は、下記の数式 5 および数式 6 のように単純化される。

【 0 0 4 4 】

【数 5】

$$\zeta(t) = \zeta(t-1) + \dot{\zeta}(t) \cdot \Delta T \quad (5)$$

【 0 0 4 5 】

10

20

30

40

【数 6】

$$\dot{\zeta}(t) = \frac{\dot{\zeta}(t-1) + u_{\zeta}(t) \cdot \Delta T}{1 + b \cdot \Delta T} \quad (6)$$

対象の動きに突然の変化があれば、

$$\dot{\zeta}(t)$$

を計算するための数式 5 の逆方向近似は間違えである。エラーは、

$$\zeta(t)$$

10

を得るために、

$$\ddot{\zeta}(t)$$

を 2 回積分する時に更に大きくなる可能性がある。さらに、本発明によれば、( t + 1 ) と

$$\dot{\zeta}(t+1)$$

は、

20

$$\dot{\zeta}(t)$$

および

$$\ddot{\zeta}(t)$$

をそれぞれ近似するために、数式 7 および数式 8 で提示されたように組み込まれる。

【 0 0 4 6 】

【数 7】

$$\dot{\zeta}(t) = \frac{\zeta(t+1) - \zeta(t-1)}{2\Delta T} \quad (7) \quad 30$$

【 0 0 4 7 】

【数 8】

$$\ddot{\zeta}(t) = \frac{\dot{\zeta}(t+1) - \dot{\zeta}(t-1)}{2\Delta T} = u_{\zeta}(t) - b\dot{\zeta}(t) \quad (8)$$

前記数式に基づいて、図 1 に示す装置に対する動きの式は、数式 9 および数式 10 の通りである。

【 0 0 4 8 】

40

【数 9】

$$\zeta(t+1) = \zeta(t-1) + \dot{\zeta}(t) \cdot 2\Delta T \quad (9)$$

【 0 0 4 9 】

【数 10】

$$\dot{\zeta}(t+1) = -d \cdot 2\Delta T \cdot \dot{\zeta}(t) + \dot{\zeta}(t-1) + u_{\zeta}(t) \cdot 2\Delta T \quad (10)$$

行列形式にすれば、動きの式は、次の数式 11 ないし数式 14 になる。

【 0 0 5 0 】

50

【数 1 1】

$$\Xi(t+1) = \mathbf{F}(t)\Xi(t) + \mathbf{G}(t)u_{\xi}(t) \quad (11)$$

【0 0 5 1】

【数 1 2】

$$\Xi(t+1) = \left[ \xi(t) \quad \dot{\xi}(t) \quad \xi(t+1) \quad \dot{\xi}(t+1) \right]^T \quad (12)$$

【0 0 5 2】

10

【数 1 3】

$$\mathbf{F}(t) = \left[ \begin{array}{cc|cc} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ \hline 1 & 0 & 0 & 2\Delta T \\ 0 & 1 & 0 & -b \cdot 2\Delta T \end{array} \right] = \left[ \begin{array}{c|c} 0 & \mathbf{I} \\ \hline \mathbf{I} & \mathbf{F}_0 \end{array} \right] \quad (13)$$

【0 0 5 3】

【数 1 4】

$$\mathbf{G}(t) = \left[ \begin{array}{c} 0 \\ 0 \\ 0 \\ 2\Delta T \end{array} \right] \quad (14)$$

20

【0 0 5 4】

動く物体には、ロボット自体と、人を含む対象物体との2つがある。図1に示す装置を含むロボットに対し、外力は、制御命令  $u(t) = [u(t)]$  であり通常既知である。時間  $t$  でロボットの姿勢は、 $r(t)$  で表せる。例えば、平面環境でのロボットの動作に対し、この姿勢は平面で  $x - y$  位置およびその進路方向から構成される。姿勢が、数式 15 で特定された第一次マルコフ (Markov) 過程に従うと仮定する。ただし、図1に示す装置が動かない所では、 $r(t)$  は定数である。

30

【0 0 5 5】

【数 1 5】

$$p(r(t+1)|r(t), u(t)) \quad (15)$$

【0 0 5 6】

カルマンフィルタ等または後続 (successor) タイプフィルタ用いれば、姿勢を十分に測定できる。同時位置把握と地図作成 (Simultaneous Localization and Map Building: SLAM) アルゴリズムは、本発明の実施形態におけるノイズの観察と制御とを行う場合には、姿勢  $r(t)$  の予測だけでなく、地図を探すためにコンピュータを使用する。このようなアルゴリズムの例は、M.Montemerlo, "FastSLAM: A Factored Solution to the Simultaneous Localization and Mapping Problem with Unknown Data Association", Ph.D.dissertation, CMU, 2003 に詳細に記載されている。

40

時間  $t$  を使って対象の姿勢を表すと、 $s(t)$  になる。対象に対する挙動がわからないため、対象に加えられる外力  $v(t)$  は、数式 16 に提示されたガウス関数としてモデル化され、対象の姿勢は、数式 17 で提示された第一次マルコフ過程によってコンピュータにより推定される。

【0 0 5 7】

50

【数 1 6】

$$v(t) = N(v(t); 0, \Sigma) \quad (16)$$

【0 0 5 8】

【数 1 7】

$$p(s(t+1)|s(t), v(t)) \quad (17)$$

測定モデルについて、観察データ集合  $Z(t)$  は、 $m$  番目のマイクロフォン 210 により観察された要素  $z_m(t)$  ( $m = 1, \dots, m$ ) を含む複数のチャンネルオーディオストリーム  $z_a(t)$ 、およびカメラにより観察される極座標での全方向映像データ  $z_v(t) = I(r, \theta, t)$  を含む。したがって、観察データ集合  $Z(t)$  は、数式 18 により表わされる。

10

【0 0 5 9】

【数 1 8】

$$Z(t) = \{z_a(t), z_v(t)\} \quad (18)$$

【0 0 6 0】

観察データ集合  $Z(t)$  を求めるためのバックグラウンドとして、

J. Vermaak and A. Blake, "Nonlinear Filtering for Speaker Tracking in Noisy and Reverberant Environments", in Proc. ICASSP, 200,

20

C. Choi, "Real-time Binaural Blind Source Separation", in Proc. Int. Symp. ICA and BSS, 2003, pp. 567 to 572、

G. Lathoud and I. A. McCowan, "Location based Speaker Segmentation", in Proc. ICASSP, 2003、

G. Lathoud, I. A. McCowan and D. C. Moore, "Segmenting Multiple Concurrent Speakers using Microphone Arrays," in Proc. Eurospeech, 2003、

R. Cutler et al. "Distributed Meetings: A Meeting Capture and Broadcasting System", in Proc. ACM Int. Conf. Multimedia, 2002、

Y. Chen and Y. Rui, "Real-time Speaker Tracking using Particle Filter Sensor Fusion", Proc. of the IEEE, vol. 92, No.3, pp.485 to 494, 2004

30

に記載されているように、時間遅延測定 (Time-Delay Estimates: TDE) は、音声追跡のための構造を記述する。しかし、本発明の一実施形態によって使用可能であり、周辺のノイズおよび反響に対処するために、最大の尤度法および位相変換からの重み関数があっても、TDEに基づく技術は、M. Brandstein and D. Ward, EDS., Microphone Arrays: Signal Processing Techniques and Applications. Springer, 2001に記載されたように、明示的な方向のノイズに脆弱である。

【0 0 6 1】

逆に、信号部分空間法は、複数音源を前提としたシナリオを利用する長所がある。さらに、信号部分空間法は比較的簡単で、かつ明確であり、また、広帯域の信号に対して高解像度とノンバイアスの角の測定を提供する。このようなサブ空間の例は、G. Su and M. Morf, "The Signal Subspace Approach for Multiple Wide-band Emitter Location," IEEE Trans. ASSP, vol. 31, No. 6, pp. 1502 to 1522, 1983 and H. Wang and M. Kaveh, "Coherent Signal-subspace Processing for the Detection and Estimation of Angles of Arrival of Multiple Wide-band Sources," IEEE Trans. ASSP, vol. 33, No. 4, pp. 823 to 831, 1985に記載されている。さらに、本発明の一実施形態において、図 2 の方法およびコンピュータ 400 は、TDE の代わりにサブ空間アプローチを使用する。しかし、普遍性の喪失なしに、TDE に基づく方法は、信号部分空間法の代わりに、または共に使用でき、TDE に基づく方法は、さらに本発明の実施形態に係る帰納的ベイズ (recursive Bayesian) フィルタリングのフレームワークで行われるということは言

40

50

うまでもない。

【 0 0 6 2 】

観察データ集合  $Z(t)$  を求めるためのバックグラウンドとして、図 2 の方法およびコンピュータ 400 は、本発明の一実施形態に係るハウスドルフ距離を使用して、イメージを比較して物体検出を行う。ハウスドルフ距離の例は、D. P. Huttenlocher, G. A. Klamberman and W. J. Rucklidge, "Comparing Images Using the Hausdorff Distance under Translation," in Proc. IEEE Int. Conf. CVPR, 1992, pp.654 to 656に記載されている。この方法は、スケーリングおよび変換において簡単でかつ強固なため本発明の実施形態で使用するが、この方法は、多様なサイズのあらゆる候補イメージを比較するため相当な時間がかかる。

10

【 0 0 6 3 】

本発明の一実施形態において、より迅速な計算のために、単純さが特徴であるブーストカスケード構造を使用する。ブーストカスケード構造は、P. Viola and M. Honess, "Rapid Object Detection Using a Boosted Cascade of Simple Features," in Proc. CVPR, 2001に記載されている。さらに他の例は、歩行者検出システムのコンテキストに記載されており、P. Viola, M. Jones and D. Snow, "Detecting Pedestrians using Patterns of Motion and Appearance," in Proc. ICCV, 2003に記述されたように、動きおよび形状を一つのモデルに合わせる。たとえば、本発明の実施形態に使用できたとしても、ブーストカスケード構造は、速度および性能面で効果的であるが、一方、難しい学習と相当の練習サンプルを必要とする。

20

【 0 0 6 4 】

物体の認識を行うにあたって、色は本発明の実施形態において適切な識別因子である。人を検出する状況で、皮膚色は、人を探す魅力的な視覚の手がかりであることが分かる。このような認識の例は、M. Jones and J. M. Rehg, "Statistical Color Models with Application to Skin Detection," International Journal of Computer Vision, 2002に記述されている。したがって、ハウスドルフ距離およびブーストカスケード構造が、本発明の一実施形態によって使用できるとしても、コンピュータ 400 は、計算速度を上げるために、皮膚色の検出を使用し、難しい学習の負担を減らすように簡単な外観モデルを使用する。しかし、人または他の物体に対して他の色が、本発明の実施形態によって視覚手がかりとして使用可能なことは言うまでもない。

30

【 0 0 6 5 】

追跡は、Y. Bar-Shalom and X. R. Li, multitarget-multisensor Tracking: Principles and Techniques, Yaakov Bar-Shalom, 1995に記載されているように、長い間、航空宇宙工学で課題となっている。

最近、視覚で物体追跡を行うことに関連した分野で発展している。そのような方法の例として、mean shift方法、CAMSHIFT方法およびCONDENSATIONアルゴリズムが挙げられる。このような方法の例は、

D. Comaniciu, V. Ramesh, and P. Meer, "Real-time Tracking of Non-rigid Objects using Mean Shift," in Proc. CVPR, 2000と、

"Kernel-based Object Tracking," IEEE Trans.PAMI, 2003と、

40

G. R. Bradski, "Computer Vision Face Tracking for use in a Perceptual user Interface," Intel Technology Journal, 1998と、

M. Isard and A. Blake, "Contour Tracking by Stochastic Propagation of Conditional Density," in Proc.ECCV, 1996と、

"Condensation:Unifying Low-level and High-level Tracking in a stochastic Framework," in Proc. ECCV, 1998

に記載されている。

【 0 0 6 6 】

また、Y. Chen and Y. Ruiの "Real-time Speaker Tracking using particle Filter Sensor Fusion," Proc. of the IEEE 2004に記載されているような関心 (interest) 粒子

50

フィルタ追跡が増加している。それに対し、音出力源 ( sound emitter ) 追跡はあまり一般的ではないが、興味深い主題であり、J. Vermaak and A. Blake, "Nonlinear Filtering for Speaker Tracking in Noisy and Reverberant Environments," in Proc. ICASSP, 2001に記載されている。

【 0 0 6 7 】

位置把握および追跡について、本発明の一実施形態では、著名な回帰的ベイジアン ( celebrated recursive Bayesian ) フィルタリングを利用する。このフィルタリングは、根源的かつ基本であり、概略を言えば、他のアルゴリズムは、このフィルタリングの変更または近似バージョンである。

【 0 0 6 8 】

図 1 に示すように、マイクロフォンアレイのマイクロフォン 2 1 0 は、等方位的であるため音源ローカライザであり、あらゆる方向からの音源からの到達角度が分かる。本発明の一実施形態に係るコンピュータ 4 0 0 によるサブ空間アプローチは、それらの測定パラメータ ( すなわち、アレイマイクロフォンと各話者との角度が固定 ) を仮定すれば、期間上のアンサンプル平均を経て、観察された信号から出た空間共分散マトリックスに基づく。

測定された音声データは、次の数式 1 9 のように、周波数領域において  $m$  次元ベクトル (  $m$  センサー ) で得られる。図 1 に示すように、マイクロフォンの配列は、8 個のマイクロフォン 2 1 0 を含むため、図示する例で、 $m = 8$  である。しかし、マイクロフォンの数が異なることにより、異なる値が使用されることは言うまでもない。

【 0 0 6 9 】

【 数 1 9 】

$$\mathbf{z}_a(f, t) = \mathbf{A}(f, \theta) \mathbf{x}(f, t) + \mathbf{n}(f, t) \quad (19)$$

数式 1 9 で、 $\mathbf{z}_a(f, t)$  は、 $m \times 1$  の観察ベクトルであり、 $\mathbf{x}(f, t)$  は、 $d \times 1$  サイズの音源ベクトルで、 $\mathbf{n}(f, t)$  は、周波数  $f$  および離散時間  $t$  で、 $m \times 1$  の測定ノイズベクトルである。 $\mathbf{A}(f, \theta)$  は、調整ベクトル  $\mathbf{a}(f, \theta)$  を含む伝達関数マトリックスである。調整ベクトル  $\mathbf{a}(f, \theta)$  は、周波数  $f$  で方向  $\theta$  の配列への信号源の伝達を反映した減衰および遅延を表す。

本発明の一実施形態によれば、調整ベクトル  $\mathbf{A}(f, \theta)$  は、マイクロフォンアレイ構造に対して、実験的に  $5^\circ$  間隔で形成されたインパルス音に対する反応を測定することで決定される。しかし、そのベクトル  $\mathbf{A}(f, \theta)$  は、他の方法で導き出すことができるのは言うまでもない。

【 0 0 7 0 】

観察のための空間共分散マトリックスは、  

$$\mathbf{R}(f) = \mathbf{E} \{ \mathbf{z}_a(f, t) \cdot \mathbf{z}_a^H(f, t) \}$$
の式からあらゆる連続的なフレームに対して得ることができる。ここで、“ $H$ ” は、エルミート転置行列を表す。空間共分散マトリックス  $\mathbf{N}(f)$  は、明示的な方向性の音源がなければ、あらかじめ計算される。したがって、数式 2 0 に表されたような一般的な固有値の問題を解くことは、一般的な固有値マトリックス、およびその該当固有ベクトル  $\mathbf{E} = [ \mathbf{E}_S | \mathbf{E}_N ]$  に帰着する。

$\mathbf{E}_S = [ \mathbf{e}_S^1, \dots, \mathbf{e}_S^d ]$  および  $\mathbf{E}_N = [ \mathbf{e}_N^{d+1}, \dots, \mathbf{e}_N^m ]$  は、それぞれ信号サブ空間およびノイズサブ空間に及ぶ固有ベクトルマトリックスである。“ $d$ ” は、音源数の近似値であり、推定数字 ( 3 のように ) で表示される。要求されるものではないが、“ $d$ ” は、追跡される人数に基づいて入力される。しかし、一般化された固有値問題は、本発明の一実施形態に係る他の固有分析方法により代替される。そのような方

10

20

30

40

50

法の例は、固有値問題に限定されるものではなく、特異値分解および本発明の実施形態に係る一般化された特異値の分解を含む。

【 0 0 7 1 】

【 数 2 0 】

$$R(f) \cdot E = N(f) \cdot E \cdot \Lambda \quad (20)$$

10

音声システム 200 により受信された周波数  $f$  および各方向  $\theta$  に音源がある条件付尤度  $p(z_a(t) | f, \theta)$  は、本発明の一実施形態に係る MUSIC (Multiple Signal Classification) アルゴリズムを使用したコンピュータ 400 により獲得される。しかし、他の方法が使用可能なことは言うまでもない。数式 21 で、 $a(f, \theta)$  は、周波数  $f$  および方向  $\theta$  での調整ベクトルである。

【 0 0 7 2 】

【 数 2 1 】

$$p(z_a(t) | f, \theta) = p(z_a(f, t) | \theta) = \frac{a^H(f, \theta) a(f, \theta)}{a^H(f, \theta) E_N(f, \theta) E_N^H(f, \theta) a(f, \theta)} \quad (21)$$

20

特定の角度の方向  $\theta$  にある特定音源の尤度は、次の数式 22 ないし数式 24 に表されている。

【 0 0 7 3 】

【 数 2 2 】

$$p(z_a(t) | \theta) = \int_f p(z_a(t), f | \theta) df \quad (22)$$

【 0 0 7 4 】

【 数 2 3 】

$$p(z_a(t) | \theta) = \int_f p(z_a(t) | f, \theta) p(f | \theta) df \quad (23)$$

30

【 0 0 7 5 】

【 数 2 4 】

$$p(z_a(t) | \theta) = \int_f p(z_a(f, t) | \theta) p(f) df \quad (24)$$

数式 22 ないし数式 24 で表すように、 $p(f | \theta)$  は、 $p(f)$  に代替されるが、これは、周波数選択が根源信号の方向と無関係と仮定されたためである。装置が離散周波数領域にあり、周波数ピンの選択が何れも  $p(f_k) = 1 / N_f$  と同じであると仮定すれば、数式 24 の各信号源の方向  $\theta$  の尤度は、コンピュータ 400 が信号源に対する方向可能性を検出できるように、本発明の一実施形態によって、数式 25 および数式 26 で表される。

40

【 0 0 7 6 】

【 数 2 5 】

$$p(z_a(t) | \theta) = \frac{\sum_{f_k \in F} P(f_k, \theta)}{N_f} \quad (25)$$

【 0 0 7 7 】

50

【数 26】

$$p(z_a(t)|\theta) = \frac{\sum_{f_k \in F} \frac{a^H(f_k, \theta) a(f_k, \theta)}{a^H(f_k, \theta) E_N(f_k, \theta) E_N^H(f_k, \theta) a(f_k, \theta)}}{N_f} \quad (26)$$

数式 26 を使用して、コンピュータ 400 は、図 4 A に示す角関数として音声尤度を計算する。人の追跡について記述したが、他の物体（すなわち、車、在庫品、飛行機、船など）および動物も、本発明の一実施形態によって追跡可能である。

多数の人々を追跡するにあたって、図 1 に示す装置は、図 3 A に示すように、同時に全ての人が見えるように、360° の視野を有する全方向カラーカメラ 110 を使用する。多数の人を見つけるために、2 つの特徴である皮膚色およびイメージ形状が本発明の一実施形態によって使用される。皮膚領域は、ほぼ一定の色を有するため、顔および手の領域は、図 3 C に示すように、カラー分割を使用して容易に区分できる。多数の人種および皮膚色を追跡できるように、本発明の一実施形態によって多様な皮膚色が検出されるということは言うまでもない。皮膚色プロブが人であるか否かを区別するために、上半身の 3 つの形状が本発明の一実施形態によって、コンピュータ 400 により具体化されて使用される。

【0078】

具体的に、図 3 A ないし図 3 C に示すような入力カラーイメージは、図 3 C および図 3 B にそれぞれ示すように、コンピュータ 400 により色変換され、閾値設定されたイメージおよびエッジイメージの 2 つのイメージに変換される。第一イメージ（すなわち、図 3 C に示す例）は、色標準化および閾値による色変換により形成される。特に、

$$r = \frac{R}{R+G+B} ; \quad g = \frac{G}{R+G+B} ; \quad b = \frac{B}{R+G+B}$$

である。色変換は、2D ガウス関数  $N(m_r, \sigma_r; m_g, \sigma_g)$  で表現され、ここで  $(m_r, \sigma_r)$  および  $(m_g, \sigma_g)$  は、それぞれ赤色および緑色要素の平均および標準偏差である。標準化された色は、色認識過程に大きく影響する輝度効果を低減させ、色要素をそのままに残す。ピクセル値が皮膚と関連した色に更に接近する時、ピクセルは強い強度を有する。皮膚と関連した色による臨界設定は、第一イメージを形成する。他の色が選択されるか、または他の色調を獲得するために、前記図示する色に追加または代替して、他の選択された色が強い強度を有するような変換が適用されるということは言うまでもない。

【0079】

第二イメージ（すなわち、図 3 B に示す例）は、3 つのエッジイメージ（赤色、緑色および青色）の平均である。変換されて閾値設定されたイメージ（すなわち、図 3 C に示す例）の各皮膚色プロブの中心とサイズに基づいて、コンピュータ 400 は、エッジイメージで人の上半身に対するサイズが標準化された候補群を得る。しかし、他のテンプレートエッジイメージが、本発明の一実施形態によって、人の上半身に追加または代替して使用され、エッジイメージは、他の方法で標準化されるということは言うまでもない。例として、対象が動物または他の物体（すなわち、車、在庫品、飛行機、船など）を含めば、テンプレートは、そのような対象または動物を視覚的に認識するのに使用される形状または部分を反映する。

【0080】

本発明の一実施形態に係る人間形状のマッチングのために、コンピュータ 400 は、人の姿勢と一致する人の上半身の 3 つの形状モデルイメージ（すなわち、エッジイメージテンプレート）を使用する。使用される 3 つの形状モデルイメージは、前側、左側、および右側の形状を含む。形状モデルイメージと候補エッジイメージとの類似度を計算するために、コンピュータ 400 は、形状モデルイメージと候補エッジイメージとのハウズドルフ

10

20

30

40

50

距離を測定する。ハウスドルフ距離は、集合間の類似度の測定を定義する。ハウスドルフ距離の例は、D. P. Huttenlocher, G. A. Klanderman and W. J. Rucklidge, "Comparing Images Using the Hausdorff Distance under Translation" in Proc. IEEE Int. Conf. CVPR, 1992, pp.654 to 656に詳細に記載されている。

【0081】

ハウスドルフ距離は、二つの非対称距離を有する。形状モデルイメージの  $A = \{ a_1, \dots, a_p \}$  および候補エッジイメージである  $B = \{ b_1, \dots, b_q \}$  の二つの点集合が与えられれば、形状モデル  $A$  と候補エッジイメージ  $B$  との間のハウスドルフ距離  $H$  は、数式 27 で表すように決定される。

【0082】

【数27】

$$H(A, B) = \max(h(A, B), h(B, A)) \quad (27)$$

数式 27 で、

$$h(A, B) = \max_{a \in A} \min_{b \in B} \|a - b\|$$

である。関数  $h(A, B)$  は、 $A$  から  $B$  までの直接ハウスドルフ距離と呼ばれ、 $B$  のある点から最も遠くにある点を識別し、 $a$  から  $B$  の最も隣接する点までの距離を測定する。言い換えれば、 $A$  から  $B$  までの直接の距離は、 $A$  のあらゆる点が  $B$  の幾つかの点  $b$  と近い時に小さい。二つとも小さい場合には、コンピュータ 400 は、候補エッジイメージおよび形状モデルイメージが互いに似ていると判断する。如何なる実施形態で要求されるものではないが、ハウスドルフ距離の三角不等式は、複数の保存された形状モデルイメージを、カメラ 110 などから得たエッジイメージと比較する場合に特に有効である。この距離をもって、コンピュータ 400 は、保存された姿勢および人の身体のイメージを使用して、映像イメージから人の上半身および体の姿勢を検出できる。したがって、コンピュータ 400 により行われる方法は、図 3 A ないし図 3 C に示すように、複雑な背景および照明変化のある複雑な環境で多数の人々を検出することができる。

【0083】

本発明の一実施形態によれば、コンピュータ 400 は、それぞれの検出された人  $i$  の中心  $c_i$  を平均にし、検出された人が占める角度範囲の関数で定まる分散  $\sigma^2$  を有する 1D ガウス関数のガウス混合モデルを利用して、映像システム 100 を介して検出されたイメージの尤度関数を決定する。分散  $\sigma^2$  は、一般的に人のサイズ（すなわち、 $c_i$  で中心から人  $i$  が占める各サイズ）を反映する。分散  $\sigma^2$  は、検出された人の各範囲の増加関数である。したがって、人である映像イメージが検出された尤度は、数式 28 の通りである。

【0084】

【数28】

$$p(z_i(t) | \theta) = \sum_i \alpha_i N(\theta_i, \sigma_i^2) \quad (28)$$

数式 28 で、 $\alpha_i$  は、候補イメージに対する混合加重値であり、ハウスドルフ距離 (Hausdorff distance)  $B$  の減少関数である。例えば、 $\alpha_i$  は、ディスタンス ( $A, B$ ) に反比例する。ハウスドルフ距離の減少値は、一致に対する高い尤度を表し、候補イメージが形状モデルイメージのうち、何れか一つとの一致の度合いを表す。

【0085】

また、複数の対象物を検出し、位置把握および追跡するために、コンピュータ 400 は

10

20

30

40

50

、さらに数式 29 に提示された連続観察  $Z^t$  に対する対象ポーズ分布の回帰的測定を行う。コンピュータ 400 により行われる本発明に係る回帰は、数式 30 ないし数式 34 で提示される。

【0086】

【数29】

$$Z^t = \{Z(1), \dots, Z(t)\} \quad (29)$$

【0087】

【数30】

$$p(s(t)|Z^t) = p(s(t)|Z(t), Z^{t-1}) \propto p(Z(t)|s(t), Z^{t-1})p(s(t)|Z^{t-1}) \quad (30)$$

10

【0088】

【数31】

$$p(Z(t)|s(t), Z^{t-1}) = p(Z(t)|s(t)) = p(z_a(t)|s(t))p(z_v(t)|s(t)) \quad (31)$$

【0089】

【数32】

$$p(s(t)|Z^{t-1}) = \int p(s(t), s(t-1)|Z^{t-1}) ds(t-1) \quad (32)$$

20

【0090】

【数33】

$$p(s(t)|Z^{t-1}) = \int p(s(t)|s(t-1), Z^{t-1})p(s(t-1)|Z^{t-1}) ds(t-1) \quad (33)$$

【0091】

【数34】

$$p(s(t)|Z^{t-1}) = \int p(s(t)|s(t-1))p(s(t-1)|Z^{t-1}) ds(t-1) \quad (34)$$

30

また、本発明の一実施形態によれば、尤度  $p(s(t)|s(t-1))$  は、前記で提示した数式 5、数式 6、数式 9 および数式 10 の能動モデルに従うため、尤度  $p(s(t)|s(t-1))$  は、数式 35 で表すように、本発明の一実施形態に係るガウス分布により更に近似される。

【0092】

【数35】

$$p(s(t)|s(t-1)) = N(s(t); s(t-1), \Sigma) \quad (35)$$

したがって、数式 34 および数式 35 は、次の数式 36 のように、コンボリューション積分で結合され、コンピュータ 400 により行われるベイジアンフィルタリングは、数式 37 および数式 38 で表されるように要約される。

40

【0093】

【数36】

$$p(s(t)|Z^{t-1}) = \int N(s(t); s(t-1), \Sigma) p(s(t-1)|Z^{t-1}) ds(t-1) \quad (36)$$

【0094】

【数 3 7】

$$p(s(t)|Z^{t-1}) = N(s(t); s(t-1), \Sigma) * p(s(t-1)|Z^{t-1}) \quad (37)$$

【0095】

【数 3 8】

$$p(s(t)|Z^t) \propto p(z_a(t)|s(t))p(z_v(t)|s(t))p(s(t)|Z^{t-1}) \quad (38)$$

数式 37 で、演算子 \* は、本発明の一実施形態によって、コンピュータ 400 により使用されるコンボリユーション演算を表す。また、コンピュータ 400 により行われるベジアン回帰は、予測演算および補正演算を含む。特に、予測演算は、対象の移動に対する能動モデルに基づいて対象のポーズを予測するために数式 37 を使用する。補正演算は、予測された対象のポーズが現在の観察尤度によって調整される数式 38 を使用する。

10

【0096】

本発明の一実施形態によれば、コンピュータ 400 は、重なった音声を分離するためにビームフォーマを備える。これにより、コンピュータ 400 は、対話中にそれぞれの話者の音声を分離でき、トラックは、本発明の一実施形態に従って、各識別された話者に対して個別に出力される。しかし、もし分離された音声を出力する必要がなく、装置が各個人のみを認識することのみが必要ならば、ビームフォーミングを使用する必要がないか、または下記のような方法により、各個人の認識を行う。

20

【0097】

話者分割は、会話、会議および業務の対話において重要であり、なお、大語彙連続音声認識システム、対話システムおよびディクテーションシステムのような多くの音声処理設備でも有用である。重なった音声は、話者順に分割音声で中心位置を占める。これについては、E. Shirberg, A. Stolcke and D. Baron, "Observations on Overlap: Findings and Implications for Automatic Processing of Multi-party Conversation," in Proc. Eurospeech, 2001 に記載されている。マイクロフォンアレイによる重なった音声の分割結果は、両耳ブライント信号分離 (Binaural Blind Signal Separation: BSS)、デュアルスピーカが隠れたマルコフモデル (dual-speaker hidden Markov models)、および遅延時間予測を有するモデル化した話者位置にガウス分布を組み込んだ音声/サイレント比率を使用して報告される。この結果の例は、

30

C. Choi, "Real-time Binaural Blind Source Separation," in Proc. Int. Symp. ICA and BSS, pp. 567 to 572, 2003 と、

G. Lathoud and I. A. McCowan, "Location based Speaker Segmentation," in Proc. ICASSP, 2003 と、

G. Lathoud, I. A. McCowan and D. C. Moore, "Segmenting Multiple Concurrent Speakers using Microphone Arrays," in Proc. Eurospeech, 2003 と、

に開示されている。五つのビデオストリーム入力およびマイクロフォンアレイから出たパノラマ式イメージを使用した話者追跡は、R. Cutler et al., "Distributed Meetings: A Meeting Capture and Broadcasting System," in Proc. ACM Int. Conf. Multimedia, 2002 and Y. Chen and Y. Rui, "Real-time Speaker Tracking using Particle Filter Sensor Fusion," Proc. of the IEEE, vol. 92, No. 3, pp. 485 to 494, 2004 に報告されている。

40

【0098】

これらの方法は、同時的話者分離において両極端にある。一方は音声情報のみに依存するが、他方は映像情報に主に依存する。さらに、Chen および Y. Rui により開示された方法は、発話の音声部分のみを記録する能力を含まずに、その代わりに、対象の人が話すか否かに関係なく、任意のデータを記録し、更に特定の話者としてオーディオチャンネルを識別するために映像データを使用することはできない。本発明に係る一実施形態によ

50

れば、コンピュータ400は、複数の音声を話者毎に分割し、それぞれの音声を、対象の空間情報、干渉およびノイズの時間特性を使用して分離する。この方法によると、特定の対象が話している間に、開始時間と終了時間とを検出して記録する本発明の一実施形態では、特定の人と話しているか否かに基づいて、音声および/または映像を選択的に記録でき、それにより任意のデータを記録するシステムと比較して、メモリ空間および/または伝送帯域幅を低減できる。更に、特定の話者への選択性を向上させることにより、特定の注目する対象物に焦点を合わせることができる。

【0099】

本発明の一実施形態によれば、LCMBF (Linearly Constrained Minimum Variance Beam-Former) は、分割された複数の同時発生 of 音声から各対象物の音声を分離するために、コンピュータ400により使用される。ビームフォーマの使用は、実際の調整ベクトルと仮定推定ベクトル  $a(f, \theta)$  とが一致しないため、対象の音声を潜在的に相殺する深刻な問題点がある。一般的に、実際調整ベクトル  $a(f, \theta)$  およびターゲットフリー分散マトリックスは、何れも得難いものではない。さらに、相殺に対する強固さを達成するための一つの一般的な方法は、対角線ローディングであり、この例については、S. Shahbazpanahi, A. B. Gershman, Z. -Q. Luo and K. Wong, "Robust Adaptive Beamforming using Worst-case SINR Optimization: A new diagonal loading-type solution for general-rank signal," in Proc. ICASSP, 2003に記載されている。しかし、このような一般的な手法は、H. L. V. Trees, Optimum Array Processing. Wiley, 2002に記述されているように、音声干渉を効果的に解消させることができず、干渉対ノイズの比率が、低い場合に対象物の相殺に対して強固ではないという短所を有する。

【0100】

実際ベクトルと仮定ベクトル  $a(f, \theta)$  との不一致は、本発明の一実施形態による図1の装置では特に取り扱い難い。それゆえ、コンピュータ400は、ターゲットフリー分散マトリックスを精巧に得ることに焦点を合わせる。特に、音声映像混合システムと図1および図2の方法によって、ビームフォーマは、対象の音声が存在データスナップショットに存在している否かを非常に正確に知らせることができる。このような長所は、主に、大きいノイズに対するサブ空間位置把握アルゴリズムが持つ強固さより得られる。さらに、本発明の一実施形態に係るコンピュータ400により使用されるビームフォーマは、対象の音声が存在していない時のみに分散マトリックスをアップデートできるため、対象音声との相殺は避けることが可能である。ビームフォーマで使用される加重値は、数式39を使用して計算される。

【0101】

【数39】

$$W_k = \frac{(R_k + \lambda I)^{-1} a_k(\theta_o)}{a_k^H(\theta_o)(R_k + \lambda I)^{-1} a_k(\theta_o)} \quad (39)$$

数式39で、 $\theta_o$  は対象物の方向であり、 $\lambda$  は対角線ローディングファクターであり、 $R_k$  は、対象と関係ない期間に対する第k周波数ピンでの分散マトリックスであり、 $a_k(\theta_o)$  は、 $k^{\text{th}}$  周波数ピンでの対象の方向に対する調整ベクトルである。数式39で、対角線ローディングファクター  $\lambda$  は、実際および仮定調整ベクトルとの若干の不一致による対象信号の相殺を緩和する。

【0102】

例えば、図11(a)ないし図11(c)は、8個のマイクロフォン210を使用して検出した8個のチャンネル音声入力から、更に検出したビーム形成出力を示す図面である。図11(a)ないし図11(c)示すように、コンピュータ400のビームフォーマは、3人の話者が同時に話している状況で、マイクロフォン210から音声入力の8個のチャンネルを得るために話者1ないし3を分離し、図11(a)ないし図11(c)に示す3つのチャンネルに話者を特定した位置を出力する。

## 【 0 1 0 3 】

本発明の一実施形態によれば、映像尤度は、全方向カメラ 1 1 0 からの入力を使用して計算するが、視野が限定された他のカメラを使用して計算することも可能なことは言うまでもない。このような視野が限定されたカメラの例としては、TV、カムコーダ、ウェブベースカメラ（度々コンピュータに装着されるもの）および特定方向に向うレンズを利用して、限定された視野イメージのみを個別に撮像する他のカメラなどを含む。そのような視野が限定されたシステムに対し、尤度関数は、J. Vermaak and A. Blake, "Nonlinear Filtering for Speaker Tracking in Noisy and Reverberant Environments," in Proc. ICASSP, 2001の数式 6 および数式 7 を適用することができる。具体的には、結果の式は、下記に示す数式 4 0 となる。

10

## 【 0 1 0 4 】

## 【数 4 0】

$$L(\text{video} | \theta) \leftarrow L(\text{video} | \theta) * P(\text{detection}) + \text{constant}. \quad (40)$$

## 【 0 1 0 5 】

一般的に、方向検出を促進するために、本発明の一実施形態によって少なくとも二つのマイクロフォンを使用せねばならない。さらに、本発明の一実施形態は、二つのマイクロフォンの間の中間点に配置された視野が限定されたカメラ（ウェブカメラのような）を有するデスクトップ型パソコンを使用することで実装される。

20

## 【 0 1 0 6 】

さらに、音源が視野外に位置する所で尤度関数は調整されるため、音源には、その物体が追跡される（数式 4 0 の定数を使用するように）ことを保証するために、視野外に位置すれば、追跡される対象であるという可能性が高くなる。この情報を使用して、音源は追跡される。また、コンピュータ 4 0 0 は、カメラが回転して、あらかじめ視野外にあるノイズ源に焦点を合わせるように制御でき、ノイズ源が追跡されないものと決定すれば、本発明の一実施形態によってビームフォーミング過程は、音源を排除するように使用される。または、視野外の物体が無視されるべきものであれば、コンピュータ 4 0 0 は、音源位置の尤度を減少させるようにプログラミングする。

## 【 0 1 0 7 】

さらに他の実施形態で、数式 4 0 は、座標変換を利用して、制限された視野を有する複数のカメラを合成するのに使用される。特に、マイクロフォンが所定の位置に配置されれば、グローバル座標は配列の中央に配置される。これにより、各カメラは、グローバル座標と関連した座標に割当てられ、コンピュータ 4 0 0 は、全方向カメラを要求せずに、複数のカメラおよびマイクロフォン配列を使用して対象を追跡するために座標変換を使用する。

30

## 【 0 1 0 8 】

ステップ 5 6 2 に関する本発明の一実施形態によれば、音声パターン認識（Speech Pattern Identification: SPI）は、下記に提示された数式 4 1 ないし数式 4 9 を利用して、コンピュータ 4 0 0 により行われる。特に、それぞれの出力トラックに対し、コンピュータ 4 0 0 は、人の静寂時と話している時とを対比させた可能性を検出する。図 6 ( a ) ないし図 6 ( c ) のそれぞれのチャンネルに示すように、3 人の話者のそれぞれは、話す期間および沈黙する期間を有する。一般的に、音声を重ねることは対話で予想される。各人間が、話し始める、または話し終わる時を区別するために、内積  $Y(t)$  は、下記の数式 4 1 のように、特定話者が話している途中であるという尤度  $L(t)$ （図 5 A ないし図 5 C を参照）を使用して計算する。

40

## 【 0 1 0 9 】

## 【数 4 1】

$$Y(t) = L(t)^T L(t-1) \quad (41)$$

## 【 0 1 1 0 】

50

内積を使用する場合、音声は特定トラックに存在しているか、または存在していないという2つの状態が仮定される。特に、音声が存在しなければ、 $H_0$ は、 $Y = N$ である時に検出され、音声が存在すれば、 $H_1$ は、 $Y = S$ である時に検出される。音声が存在しないかについての密度モデルは、数式42であり、音声の存否についての密度モデルは、数式43である。両密度モデルは、音声は特定の時間に特定の話者(すなわち、トラック)に対して存在しているか否かについてのモデルである。

【0111】

【数42】

$$p(Y|H_0) = \frac{1}{\sqrt{2\pi\sigma_N^2}} \exp\left(-\frac{(Y-m_N)^2}{2\sigma_N^2}\right) \quad (42) \quad 10$$

【0112】

【数43】

$$p(Y|H_1) = \frac{1}{\sqrt{2\pi\sigma_S^2}} \exp\left(-\frac{(Y-m_S)^2}{2\sigma_S^2}\right) \quad (43)$$

【0113】

密度モデルを使用して、コンピュータ400は、特定の時間に特定のオーディオトラックに対して音声の存否を決定するために密度比率を求める。音声の存否は、その比率が数式44に示された所定定数  $\eta$  を超えるか否かに基づく。 20

【0114】

【数44】

$$\frac{p(Y|H_1)}{p(Y|H_0)} \geq \eta \quad (44)$$

ここで、 $Y$ はobservation  $Y$ に対するnull hypothesis  $H_0$ が真(true)であるにもかかわらず、偽(false)として判明される確率をユーザがどれまで許容するのかによって決定される。 30

【0115】

その比率を満足すれば、コンピュータ400は、音声が存在していると判断する。一方、その比率を満足しなければ、コンピュータは、音声が存在していないと判断し、特定のトラックに対する記録/伝送が中断される。さらに、特定話者の音声に対する開始および終了時間は、音声包絡線(すなわち、特定オーディオトラックで音声が存在している間の時間)を展開するために、コンピュータ400により検出されて記録される。本発明の如何なる実施形態で要求されるものではないが、背景ノイズが記録されることを防止し、保存空間および伝送帯域幅が背景ノイズのために浪費されないように、コンピュータ400は、隣接包絡線の間で検出されたこのようなノイズを除去できるため、包絡線の開始および終了時間の間に記録された音声のみがトラックに残る。 40

【0116】

数式44の結果に基づいて、本発明の一実施形態によってコンピュータ400が、数式42および数式43の $m$ および $\sigma^2$ をオンライン更新することが可能である。更新は、数式45および数式47を使用して行われる。数式45および数式46で、 $0 < \alpha < 1$ であるが、本発明の一実施形態によって、一般的に1に近い。さらに、数式44を満足する所で、 $m_S$ および $\sigma_S^2$ は更新される。それに対し、数式44を満足せずに、その比率がより小さければ、数式42の $m_N$ および $\sigma_N^2$ が更新される。このような方法で、コンピュータ400は、数式41の内積に基づいて、密度モデルの正確度を維持できる。

【0117】

【数 4 5】

$$m \leftarrow \lambda m + (1 - \lambda)Y \quad (45)$$

【0 1 1 8】

【数 4 6】

$$\sigma^2 \leftarrow \lambda \sigma^2 + (1 - \lambda)Y^2 \quad (46)$$

本発明の一実施形態による数式 4 1 ないし数式 4 6 を使用する時は、図 6 ( a ) ないし図 6 ( c ) に示す音声は、図 7 ( a ) ないし図 7 ( c ) で示す開始および終了時間を有するように決定される。さらに、音声が存在 (  $Y = S$  ) することを示す図示された包絡線内の音声データのみが記録または伝送される必要がある。

10

【0 1 1 9】

しかし、図 7 ( c ) を参照すれば、そのほかの連続的な音声での休止は、時間 8 0 付近で見られる。したがって、記録される時、図示された隣接包絡線間に瞬間的な不連続が存在するが、これは、トラックの再生時に認識される。このような不連続は、本発明の一実施形態によって受け入れられることもあるが、本発明の一実施形態は、コンピュータ 4 0 0 が図 7 ( c ) に示す包絡線を訂正するようにすることにより、話者の音声は呼吸のための休止または大きな衝撃により不均一にならないようにする。特に、コンピュータ 4 0 0 は、長さ  $L_1$  を有する小さな静寂により分離される音声セグメントを集める。例えば、小さな静寂は、4 フレームの長さ  $L_1$  を有する。しかし、他の長さ  $L_1$  が休止を定義するために使用されるということは言うまでもない。

20

コンピュータ 4 0 0 は、連続音声部分 ( すなわち、それぞれの  $L_1$  - フレーム内 ) として見なされ、十分に近くて時間に敏感な隣接音声包絡線を結合して、その包絡線を拡張するために  $L$  フレーム拡張演算子を利用して、それぞれの検出された SPI に 2 進拡張を行う。2 進シーケンス  $u$  のために、コンピュータ 4 0 0 により使用される  $L$  フレーム拡張演算子の例は数式 4 7 で表される。

【0 1 2 0】

【数 4 7】

$$u = \{u_n\} \rightarrow v = f_{dil}^L(u), \text{ where } \forall n \quad v_n = \max(u_{n-L}, \dots, u_{n+L}) \quad (47)$$

30

【0 1 2 1】

図 8 ( a ) ないし図 8 ( c ) に示すように、コンピュータ 4 0 0 は、拡張演算を行った時、図 8 ( c ) で時間 8 0 の近くに挿入された他の休止は除去され、結合された包絡線が形成されて、音声が時間 6 0 の以後から 8 0 の以後の間で、8 0 の以前に含まれる他の休止 ( すなわち、不連続記録 ) せずに第三話者に対して連続的に記録される。

また、本発明の如何なる実施形態で要求されるものではないが、コンピュータ 4 0 0 は、通常対話の一部ではないノイズから分離したスパイクを除去する。例として、このような分離したノイズスパイクは、一般的に記録されるに当たって好ましくない咳または他の突如のノイズにより発生される。これにより、コンピュータ 4 0 0 は、本発明の一実施形態によって 2 進浸食 ( erosion ) 演算子を使用して、このようなスパイクを認識して除去で

40

【0 1 2 2】

【数 4 8】

$$u = \{u_n\} \rightarrow v = f_{ero}^L(u), \text{ where } \forall n \quad v_n = \min(u_{n-L}, \dots, u_{n+L}) \quad (48)$$

【0 1 2 3】

本発明の如何なる実施形態で要求されるものではないが、浸食演算子を行う前に、2 進

50

拡張演算子を行うことが一般的により好ましい。そうでなければ、音声期間を分離する休止が小さな記録包絡線を発生させる可能性がある。このような小さな包絡線は、浸食演算子により連続的な音声の一部に対照的なスパイクとして誤認されて、好ましくなく削除される可能性がある。

【0124】

要すれば、本発明の一実施形態によってコンピュータ400は、図7(a)ないし図7(c)に示す検出された音声包絡線に基づいて、図8(a)ないし図8(c)に示す出力を提供するために結合された数式49を使用して数式47および数式48を用いた。図8(c)に示すように、時間80付近の休止により発生する音声包絡線の不連続が除去されて、第三話者の音声の全体は、不愉快な休止なしに記録される。

10

【0125】

【数49】

$$SPI_{-} = f_{dil}^{L_2} \left( f_{ero}^{L_1+L_2} \left( f_{dil}^{L_1} (SPI) \right) \right) \tag{49}$$

【0126】

図10に示す本発明の一実施形態によれば、後処理装置710は、コンピュータ400の出力またはコンピュータ400に含まれたAVプロセッサ700を向上させるために、ブラインド音源分離(blind source separation: BSS)に適応的クロスチャンネル干渉相殺を行う。具体的に、いくつかのセンサーで記録された信号の重畳から複数の信号の分離することは、通信、生物医学および音声処理のような多様な適用において表れる重要な問題である。複数の混合されたソースの外にソース信号情報を要求していない分離方法の部類は、度々、ブラインド音源分離と呼ばれる。複数のマイクロフォンを有する実際の記録状況で、各ソース信号は、あらゆる方向に広がっており、各マイクロフォンに“直接経路”および“反響経路”を通じて到達する。観察された信号は、次の数式50のように表現できる。

20

【0127】

【数50】

$$x_j(t) = \sum_{i=1}^N \sum_{\tau=0}^{\infty} h_{ji}(\tau) s_i(t-\tau) + n_j(t) = \sum_{i=1}^N h_{ji}(t) \times s_i(t) + n_j(t) \tag{50}$$

30

数式50で、 $s_i(t)$ は、第*i*ソースシグナルであり、*N*はソースの個数、 $x_j(t)$ は、観察された信号、 $h_{ji}(t)$ は、ソース*i*からセンサー*j*までの伝達関数である。ノイズ項 $n_j(t)$ は、記録装置の特性による非線形歪曲に関連する。ソースが決して移動しないという仮定は、音響物体の能動的な性質のために度々覆る。さらに、実際のシステムは、インパルス応答の長さに限界を設定せねばならず、限定された長さは、実際の状況で度々主要な性能の障害になる。これにより、実環境のための周波数領域の暗黙ソース分離アルゴリズムは、本来の時間領域フィルタリングアーキテクチャを、周波数領域の瞬間的なBSS問題に変換するために行われる。短期フーリエ変換を使用して、数式50は、数式51のように再作成される。

40

【0128】

【数51】

$$X(\omega, n) = H(\omega) S(\omega, n) + N(\omega, n) \tag{51}$$

【0129】

説明を単純化するために、2x2の場合を例にして以下の説明を行う。しかし、一般的に、NxNの場合まで容易に展開できることは言うまでもない。数式51で、 $X(\omega, n)$ は、周波数インデックスであり、 $H(\omega)$ は、2x2正方混合マトリクスであり、 $X(\omega, n) = [X_1(\omega, n) X_2(\omega, n)]^T$  および

$$X_j(\omega, n) = \sum_{\tau=0}^{T-1} e^{\frac{-2\pi\omega\tau}{T}} x_j(t_n + \tau)$$

は、時間

$$t_n = \left\lfloor \frac{T}{2} \right\rfloor (n-1) + 1$$

(ここで、

“ $\lfloor \cdot \rfloor$ ”

は、フローリング演算子)で始まるシフト長( $T/2$ )を有するサイズ $T$ のフレームに対するDFTを表し、該当表現が、 $S(\cdot, n)$ および $N(\cdot, n)$ に適用される。混合されない過程は、次の数式52を使用して周波数ビンの公式で表すことができる。

【0130】

【数52】

$$Y(\omega, n) = W(\omega) + N(\omega, n) \quad (52)$$

数式52で、ベクトル $Y(\omega, n)$ は、 $2 \times 1$ ベクトルであり、ノイズ $N(\cdot, n)$ の効果を見捨てた原ソース $S(\cdot, n)$ についての予測である。時間領域のコンポリューション演算子は、周波数領域の複素数の掛け算に該当する。即席ICAアルゴリズムは、数式53に与えられた直角解を保証する情報最大化である。

【0131】

【数53】

$$\Delta W \propto \{\varphi(Y)Y^H - \text{diag}(\varphi(Y)Y^H)\} \quad (53)$$

数式53で、“ $^H$ ”は、複素数共役転置行列に該当し、極性の非線形関数( $\varphi(Y)$ )は、

$$\varphi(Y) = [Y_1 / |Y_1| Y_2 / |Y_2|]^T$$

で定義される。この分解の短所は、それぞれ独立した周波数ビンでの置換問題が発生するところである。しかし、この問題は、時間領域スペクトルスムージングを使用して解決される。

【0132】

第 $i$ BS S出力の各フレームに対し、 $Y_i(n) = \{Y_i(\cdot, n) \mid \cdot = 1, \dots, T\}$ によるフレームに対するあらゆる周波数要素の集合および下記の数式54で表すような最初のソースの存否をそれぞれ示す二つの仮定 $H_{i,0}$ および $H_{i,1}$ が与えられる。

【0133】

【数54】

$$H_{i,0} : Y_i(n) = \bar{S}_j(n)$$

$$H_{i,1} : Y_i(n) = \bar{S}_i(n) + \bar{S}_j(n), \quad i \neq j \quad (54)$$

数式54で、

10

20

30

40

$$\bar{S}_i$$

は、 $S_i$  のフィルタ処理されたバージョンである。 $Y_i(n)$  に条件を設定すれば、ソースの存在 / 否存在の確率は、次の数式 55 の通りである。

【 0 1 3 4 】

【 数 5 5 】

$$p(H_{i,m} | Y_i(n)) = \frac{p(Y_i(n) | H_{i,m}) p(H_{i,m})}{p(Y_i(n) | H_{i,0}) p(H_{i,0}) + p(Y_i(n) | H_{i,1}) p(H_{i,1})} \quad (55)$$

10

数式 55 で、 $p(H_{i,0})$  は、ソース  $i$  の否存在に対する先行確率であり、 $p(H_{i,1}) = 1 - p(H_{i,0})$  は、ソース音源  $i$  の存在に対する先行確率である。周波数要素のうち確率的独立を仮定すれば、数式 55 は、数式 56 になり、音源否存在確率は数式 57 になる。

20

【 0 1 3 5 】

【 数 5 6 】

$$p(Y_i(n) | H_{i,m}) = \prod_{\omega} p(Y_i(\omega, n) | H_{i,m}) \quad (56)$$

【 0 1 3 6 】

【 数 5 7 】

$$p(H_{i,0} | Y_i(n)) = \left[ 1 + \frac{P(H_{i,1})}{P(H_{i,0})} \prod_{\omega} \frac{p(Y_i(\omega, n) | H_{i,1})}{p(Y_i(\omega, n) | H_{i,0})} \right]^{-1} \quad (57)$$

30

$H_{i,1}$  の事後確率は、簡単に  
 $p(H_{i,1} | Y_i(n)) = 1 - p(H_{i,0} | Y_i(n))$   
 であり、これは、第  $i$  BSS 出力でクロスチャンネル干渉の量を表す。下記で説明するように、後処理装置 710 は、相互チャンネル干渉の相殺および要素密度  $p(Y_i(\omega, n) | H_{i,m})$  に対する統計的モデルを行う。

【 0 1 3 7 】

40

ANC の仮定された混合モデルは、FIR フィルタアーキテクチャであるため、ANC の直接適用は、実際の条件で線形フィルタの不一致を作れない。具体的に、無限フィルタ長およびセンサーノイズによる非線形性は、モデルで問題点を発生させ得る。数式 58 および数式 59 で提示されたように、後処理装置 710 により使用されるモデルに含まれる非線形特性は、差スペクトルに含まれる。

【 0 1 3 8 】

【数58】

$$|U_i(\omega, n)| = f(|Y_i(\omega, n) - a_i b_{ij}(\omega) Y_j(\omega, n)|) \quad (58)$$

$$\angle U_i(\omega, n) = \angle Y_i(\omega, n). \quad i \neq j$$

【0139】

【数59】

$$f(a) = \begin{cases} a & \text{if } a \geq \varepsilon \\ \varepsilon & \text{if } a < \varepsilon \end{cases} \quad (59)$$

数式58および数式59で、 $a_i$ はover-subtractionファクターであり、 $Y_i(\omega, n)$ は、BSS出力 $Y(\omega, n)$ の第*i*要素であり、 $b_{ij}(\omega)$ は、チャンネル*j*から*i*まで周波数 $\omega$ に対するクロス-チャンネル干渉相殺である。さらに、非線形演算子 $f(a)$ は、BSSの残っているエラーを抑えるが、大部分のスペクトル減算技術で発生するものと類似した音楽ノイズが挿入される可能性がある。

【0140】

クロス相殺が、数式58を使用して問題無く行われれば、スペクトルサイズ $|U_i(\omega, n)|$ は、ある非活性フレームに対しては0である。 $|U_i(\omega, n)|$ の複素数ガウス分布による各仮定に与えられた $Y_i(\omega, n)$ の事後確率は、次の数式60の通りである。

【0141】

【数60】

$$p(Y_i(\omega, n) | H_{i,m}) \approx p(U_i(\omega, n) | H_{i,m}) \propto \exp\left[-\frac{|U_i(\omega, n)|^2}{\lambda_{i,m}(\omega)}\right] \quad (60)$$

数式60で、 $\lambda_{i,m}$ は、減算されたフレームの分散である。 $m=1$ である時、 $\lambda_{i,m}$ は、最初のソースの分散である。 $m=0$ である時、 $\lambda_{i,m}$ は、第二ソースの分散である。分散 $\lambda_{i,m}$ は、数式61の以下の確率平均によりフレーム毎に更新される。

【0142】

【数61】

$$\lambda_{i,m} \leftarrow \{1 - \eta \lambda p(H_{i,m} | Y_i(n))\} \lambda_{i,m} + \eta \lambda p(H_{i,m} | Y_i(n)) |U_i(\omega, n)|^2 \quad (61)$$

数式61で、正の定数 $\eta$ は、適応フレーム率を表す。最初のソース信号は、BSSにより少なくとも“強調”されるものと予想される。したがって、最初のソースのサイズは、他のBSS出力チャンネルでの最初のソースである干渉ソースのサイズより大きいと仮定される。モデルパラメータを更新する間、向上したソース $\lambda_{i,1}$ の分散は、 $\lambda_{i,0}$ より小さくなるのが可能である。このような場合は好ましくないため、二つのモデルは、数式62のように変化する。

【0143】

【数62】

$$\sum_{\omega} \lambda_{i,0}(\omega) > \sum_{\omega} \lambda_{i,1}(\omega) \quad (62)$$

次いで、後処理装置710は、干渉相殺ファクターを更新する。第一に、後処理装置7

10

20

30

40

50

10 は、次の数式 6 3 ないし数式 6 5 を使用して、周波数 およびフレーム  $n$  で  $Y_i$  および  $Y_j$  のスペクトルサイズの差を計算する。数式 6 4 は、フレーム  $n$  に差の  $v$  - ノームを乗算した費用関数  $J$  を定義し、数式 6 5 は、 $b_{ij}$  に対するグラディエント減少学習規則 (gradient-descent learning rules) を定義する。

【 0 1 4 4 】

【数 6 3】

$$\delta_i(\omega, n) = |Y_i(\omega, n)|^a - \sum_{j \neq i} b_{ij}(\omega) |Y_j(\omega, n)|^a \quad (6 3)$$

【 0 1 4 5 】

【数 6 4】

$$J(\omega, n) = p(H_{i,0} | Y_i(n)) |\delta_i(\omega, n)| \quad (6 4)$$

【 0 1 4 6 】

【数 6 5】

$$\Delta b_{ij}(\omega) \propto -\frac{\partial J(\omega, n)}{\partial b_{ij}(\omega)} = p(H_{i,0} | Y_i(n)) \text{sign}(\delta_i(\omega, n)) |Y_j(\omega, n)|^a \quad (6 5)$$

このような方法論を使用して、後処理装置 7 1 0 は、図 1 1 ( a ) ないし図 1 1 ( c ) に示す入力に基づいて、図 1 2 ( a ) ないし図 1 2 ( c ) に示す向上した出力を提供する。しかし、他のタイプのクロス相殺技術の後処理装置 7 1 0 に使用して、音質を向上させることが可能なことは理解できるであろう。

【 0 1 4 7 】

以上、本発明を図面に示した実施形態を用いて説明したが、これらは例示的なものに過ぎず、本技術分野の当業者ならば、本発明の範囲および趣旨から逸脱しない範囲で多様な変更および変形が可能なことは理解できるであろう。また、本実施形態に係る方法は、プログラムとして記録媒体の形態で実装され、コンピュータによって実行されてよい。したがって、本発明の範囲は、説明された実施形態によって定められず、特許請求の範囲に記載された技術的趣旨により定められねばならない。

【産業上の利用可能性】

【 0 1 4 8 】

本発明は、会議等で話者別に音声分離および記録する装置、およびコンサートで歌手の歌からノイズを排除して記録する装置に使用することができる。

【図面の簡単な説明】

【 0 1 4 9 】

【図 1】本発明の一実施形態に係る、物体を追跡するために映像および音声情報を合成する装置を示す図面である。

【図 2】本発明の一実施形態に係る、複数の物体を追跡するために映像および音声情報を合成する方法を示すフローチャートである。

【図 3 A】本発明の一実施形態に係る、図 1 の装置により受信されて追跡される潜在的目標のイメージを含む映像の例である。

【図 3 B】本発明の一実施形態に係る、図 3 A から抽出されて追跡されるエッジイメージを示すサブイメージである。

【図 3 C】本発明の一実施形態に係る、図 3 A から抽出されて追跡される所定のカラーを含むイメージの一部分を示すサブイメージである。

【図 4 A】本発明の一実施形態に係る、特定の期間に追跡される音源の位置および追跡される対象の音声尤度を示す図面である。

【図 4 B】本発明の一実施形態に係る、特定の期間に追跡される物体の位置および追跡さ

10

20

30

40

50

れる物体の映像尤度を示す図面である。

【図 4 C】本発明の一実施形態に係る、図 4 A および図 4 B の音声尤度および映像尤度が結合された尤度を示す図面である。

【図 5 A】本発明の一実施形態に係る、結合された音声尤度および映像尤度に基づいて識別された話者 1 の位置に基づいて、話者 1 の音声尤度を示す図面である。

【図 5 B】本発明の一実施形態に係る、結合された音声尤度および映像尤度に基づいて識別された話者 2 の位置に基づいて、話者 2 の音声尤度を示す図面である。

【図 5 C】本発明の一実施形態に係る、結合された音声尤度および映像尤度に基づいて識別された話者 3 の位置に基づいて、話者 3 の音声尤度を示す図面である。

【図 5 D】本発明の一実施形態に係る、音声尤度に基づいて位置および時間関数として音声領域を示す図面である。

10

【図 6】本発明の一実施形態に係る、それぞれの該当するチャンネルを形成するために、分離された話者の音声グラフであり、( a ) は話者 1 の音声グラフ、( b ) は話者 2 の音声グラフ、( c ) は話者 3 の音声グラフである。

【図 7】本発明の一実施形態に係る、音声期間に対する開始時間と終了時間を定義する音声包絡線を示す図面であり、( a ) は図 6 ( a ) の音声に基づいた音声包絡線を示す図面、( b ) は図 6 ( b ) の音声に基づいた音声包絡線を示す図面、( c ) は図 6 ( c ) の音声に基づいた音声包絡線を示す図面である。

【図 8】本発明の一実施形態に係る、休止および突然の発声除去して、音声期間の開始および終了時間を再定義するために生成された該当音声包絡線を示す図面であり、( a ) は図 7 ( a ) の音声包絡線に基づいた該当音声包絡線を示す図面、( b ) は図 7 ( b ) の音声包絡線に基づいた該当音声包絡線を示す図面、( c ) は図 7 ( c ) の音声包絡線に基づいた該当音声包絡線を示す図面である。

20

【図 9】本発明の一実施形態に係る、選択した対象の位置を把握して集中するために、選択していない対象から発生したノイズを除去するためのビームフォーミングの使用を示す図面である。

【図 10】本発明の一実施形態に係る、図 1 の装置の出力に対して適応的クロスチャンネル干渉除去を行う後処理装置を示すブロック図である。

【図 11】本発明の一実施形態によって、図 10 の AV システムのそれぞれの音声データ出力に該当し、隣接チャンネルから干渉を受けるチャンネルを示す図面である。

30

【図 12】本発明の一実施形態に係る、それぞれのチャンネルで干渉が除去された後処理された音声データを示す図面である。

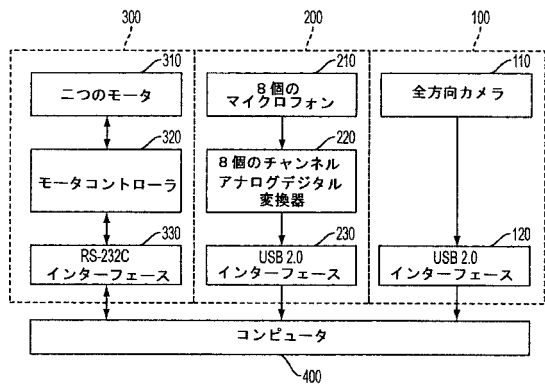
【符号の説明】

【 0 1 5 0 】

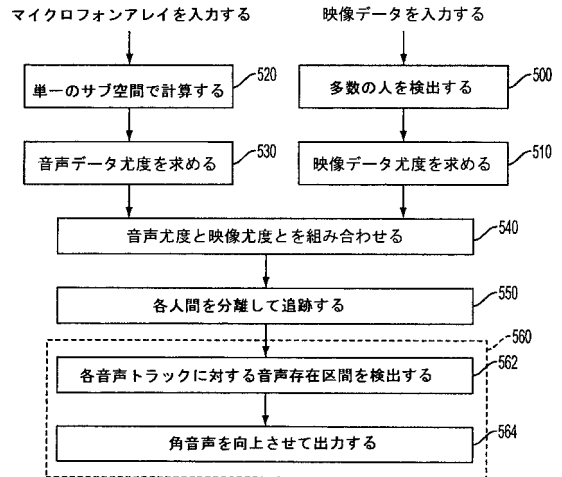
- 1 0 0 映像システム
- 1 1 0 全方向カメラ
- 1 2 0 USB 2.0 インターフェース
- 2 0 0 音声システム
- 2 1 0 8 個のマイクロフォン
- 2 2 0 アナログ - デジタル変換器
- 2 3 0 USB インターフェース
- 3 0 0 ロボット
- 3 1 0 二つのモータ
- 3 2 0 モータコントローラ
- 3 3 0 RS 2 3 2 C インターフェース
- 4 0 0 コンピュータ

40

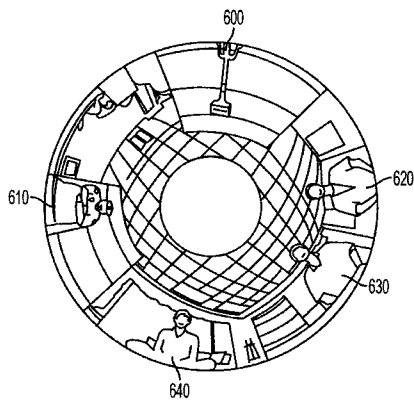
【図1】



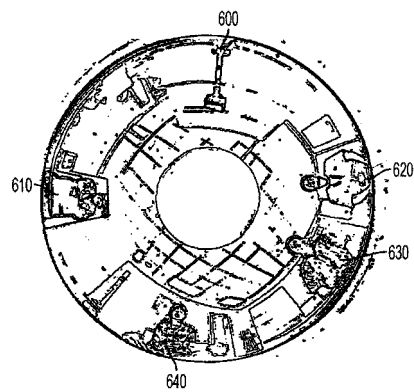
【図2】



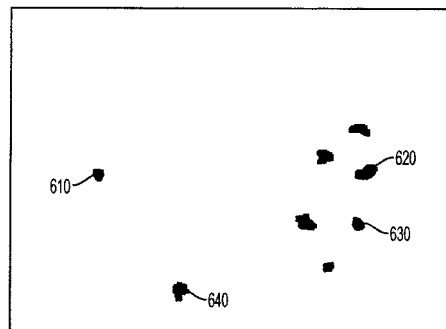
【図3A】



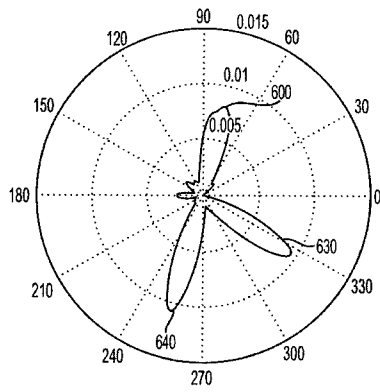
【図3B】



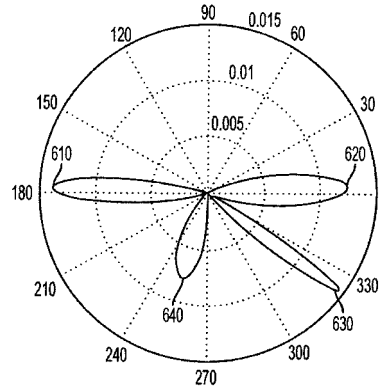
【図3C】



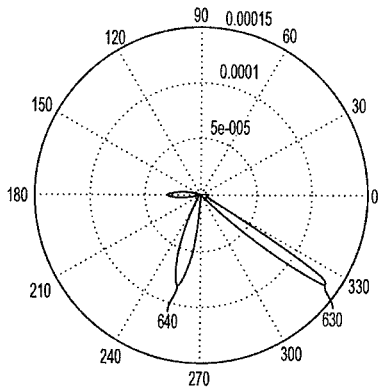
【 図 4 A 】



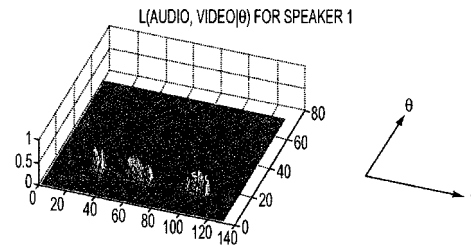
【 図 4 B 】



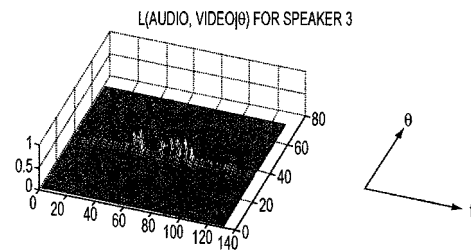
【 図 4 C 】



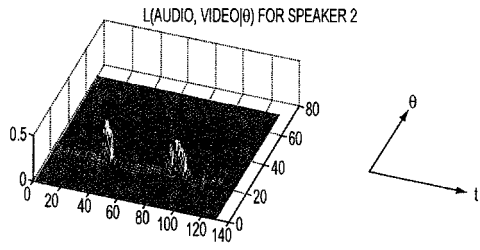
【 図 5 A 】



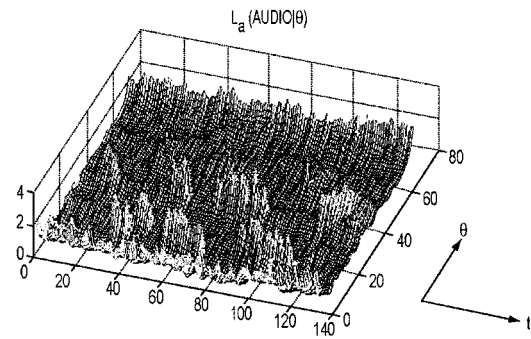
【 図 5 B 】



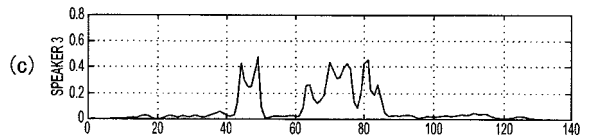
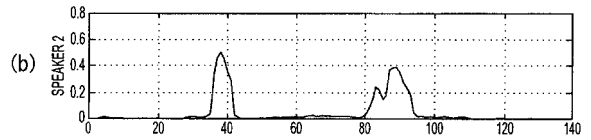
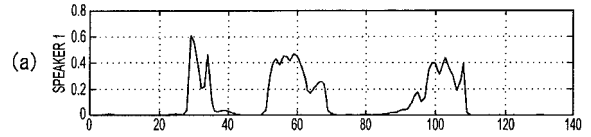
【 5 C 】



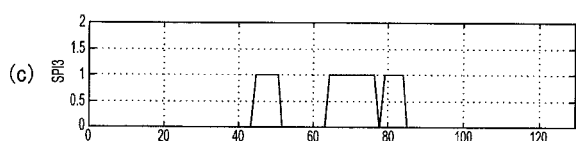
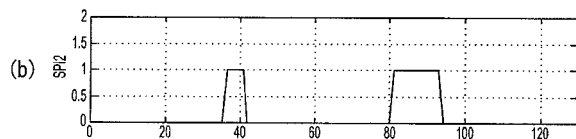
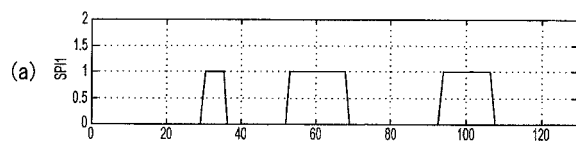
【 5 D 】



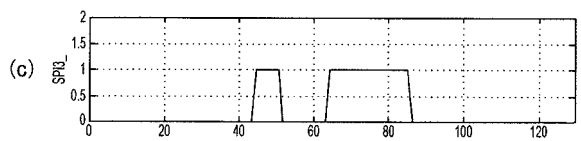
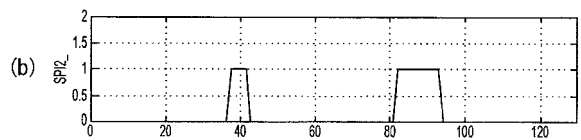
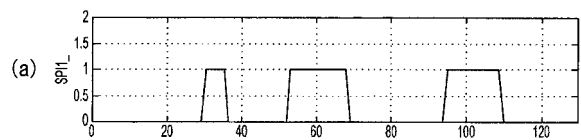
【 6 】



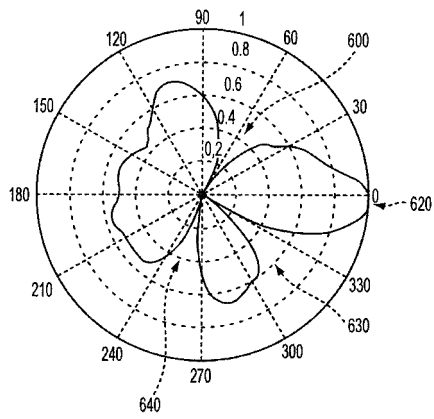
【 7 】



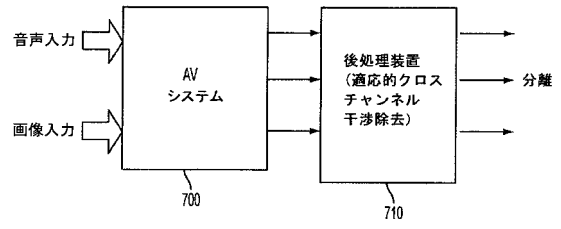
【 8 】



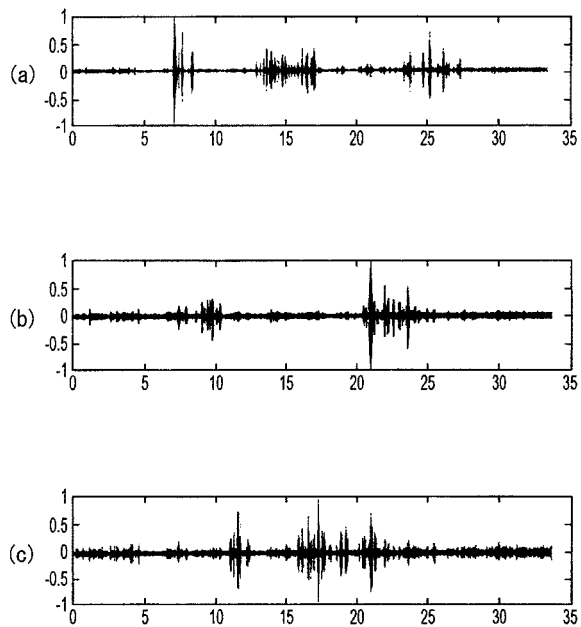
【図9】



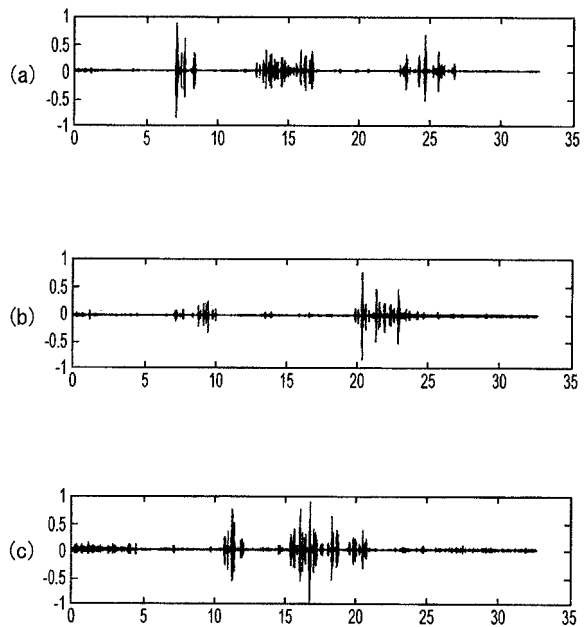
【図10】



【図11】



【図12】



## フロントページの続き

(51)Int.Cl. F I  
**G 0 5 D 1/12 (2006.01)** G 0 5 D 1/02 J  
 G 0 5 D 1/12 G

(72)発明者 李 炯 機  
 大韓民國 京畿道 水原市 靈通區 靈通洞 1 0 4 6 - 1 番地 清 明マウル4 團地アパ - ト  
 4 0 1 棟 1 9 0 2 號

(72)発明者 尹 相 民  
 大韓民國 京畿道 龍仁市 器興邑 農書里 山 1 4 - 1 番地 三星綜合技術院内

(72)発明者 孔 棟 建  
 大韓民國 京畿道 龍仁市 器興邑 農書里 山 1 4 - 1 番地 三星綜合技術院内

審査官 松浦 陽

(56)参考文献 国際公開第 0 3 / 0 3 5 3 3 4 ( W O , A 1 )  
 特開 2 0 0 4 - 2 4 9 3 8 9 ( J P , A )  
 特開 2 0 0 2 - 3 6 1 5 8 4 ( J P , A )  
 特開 2 0 0 4 - 0 3 4 9 2 4 ( J P , A )  
 Dmitry N. Zotkin, Ramani Duraiswami, Larry S. Davis , Joint Audio-Visual Tracking Using  
 Particle Filters , EURASIP Journal on Applied Signal Processing , Hindawi Publishing Co  
 rporation , 2 0 0 2 年 1 1 月 , pp. 1154-1164  
 天田皇、金澤博史、音声認識のためのマイクロホンアレー技術、東芝レビュー、日本、株式会社  
 東芝、2 0 0 4 年 9 月 1 日、Vol. 59, No. 9, pp. 42-44, U R L , [http://www.toshiba.co.jp/tech/review/2004/09/59\\_09pdf/a10.pdf](http://www.toshiba.co.jp/tech/review/2004/09/59_09pdf/a10.pdf)  
 S. Shahbazpanahi, A. B. Gershman, Z. -Q.Luo and K. Wong , Robust Adaptive Beam-forming  
 using Worst-case SINR Optimization: A new Diagonal Loading-type Solution for General-r  
 ank Signal , Proc. ICASSP , 米国 , IEEE , 2 0 0 3 年 , V-333 - V336  
 Gil-Jin Jang, Changkyu Choi, Yong-Beom Lee, Yung-Hwan Oh , Adaptive Cross-Channel Inter  
 ference Cancellation on Blind Signal Separation Outputs Using Source Absence/Presence  
 Detection and Spectral Subtraction , Proc. Interspeech 2004 - ICSLP , International Spe  
 ech Communication Association , 2 0 0 4 年 1 0 月 4 日  
 Changkyu Choi, Donggeon Kong, Hyoung-Ki Lee, Sang Min Yoon , Separation of Multiple Co  
 ncurrent Speeches using Audio-visual Speaker Localization and Minimum Variance Beam-fo  
 rming , Proc. Interspeech 2004 - ICSLP , International Speech Communication Association  
 , 2 0 0 4 年 1 0 月 4 日  
 Kiyong Park, Changkyu Choi, Jeongsu Kim , Voice Activity Detection using Global Soft  
 Decision with Mixture of Gaussian Model , Proc. Interspeech 2004 - ICSLP , International  
 Speech Communication Association , 2 0 0 4 年 1 0 月 4 日

(58)調査した分野(Int.Cl. , D B 名)

G 1 0 L 2 1 / 0 2  
 G 0 5 D 1 / 0 2  
 G 0 5 D 1 / 1 2  
 G 0 6 T 7 / 0 0  
 H 0 4 N 5 / 2 3 2  
 H 0 4 R 1 / 4 0  
 Science Direct  
 IEEE Xplore

C i N i i  
J S T P l u s ( J D r e a m I I )