

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2012-105199

(P2012-105199A)

(43) 公開日 平成24年5月31日(2012.5.31)

(51) Int.Cl.	F I	テーマコード (参考)
HO4R 3/00 (2006.01)	HO4R 3/00 320	5C164
HO4R 1/40 (2006.01)	HO4R 1/40 320A	5D018
HO4N 7/15 (2006.01)	HO4N 7/15 630Z	5D020

審査請求 有 請求項の数 10 O L (全 21 頁)

(21) 出願番号 特願2010-253947 (P2010-253947)  
 (22) 出願日 平成22年11月12日 (2010.11.12)

(71) 出願人 000003078  
 株式会社東芝  
 東京都港区芝浦一丁目1番1号  
 (74) 代理人 100076233  
 弁理士 伊藤 進  
 (72) 発明者 天田 皇  
 東京都港区芝浦一丁目1番1号 株式会社  
 東芝内  
 Fターム(参考) 5C164 FA10 UB92S VA04S VA06P VA07S  
 VA52P  
 5D018 BB22  
 5D020 BB04

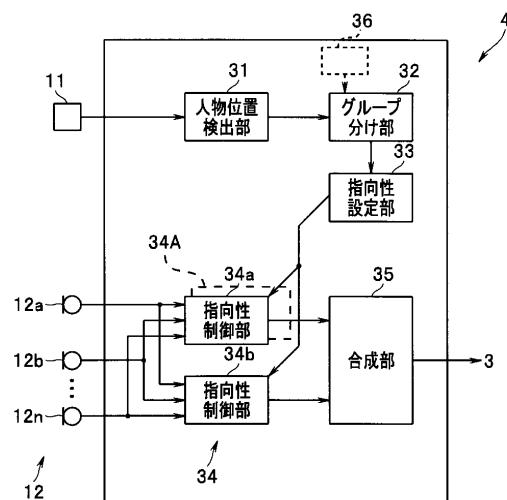
(54) 【発明の名称】 音響信号処理装置、テレビジョン装置及びプログラム

(57) 【要約】

【課題】リアルタイム性を確保できる音響信号処理の可能な音響信号処理装置を提供すること。

【解決手段】ユニット4は、所定空間内に存在する人物の位置を検出する人物位置検出部31と、人物位置検出部31により検出された1又は2以上の人物を所定数以下の数のグループに割り振るグループ分けを行うグループ分け部32と、設定されたそれぞれの指向性に基づいてマイクロホンアレー12の指向性を制御する複数の指向性制御部34と、グループ分け部32でグループ分けがされた各グループの指向性を、対応する指向性制御部34に設定する指向性設定部33とを有する。

【選択図】 図3



**【特許請求の範囲】****【請求項 1】**

所定空間内に存在する人物の位置を検出する人物位置検出部と、  
前記人物位置検出部により検出された 1 又は 2 以上の人物を所定数以下の数のグループに割り振るグループ分けを行うグループ分け部と、

設定されたそれぞれの指向性に基づいてマイクロホンアレーの指向性を制御する複数の指向性制御部と、

前記グループ分け部で前記グループ分けがされた各グループの指向性を、対応する指向性制御部に設定する指向性設定部と、

を有することを特徴とする音響信号処理装置。

10

**【請求項 2】**

前記グループ分け部は、前記マイクロホンアレーの雑音抑圧性能の和が最大になるように、前記グループ分けを行うことを特徴とする請求項 1 に記載の音響信号処理装置。

**【請求項 3】**

前記グループ分け部は、前記各グループに割り振られる人物の数が平均化するように、前記グループ分けを行うことを特徴とする請求項 1 に記載の音響信号処理装置。

**【請求項 4】**

前記グループ分けの結果のグループ数あるいは前記所定数を、変更するためのグループ数変更部を有することを特徴とする請求項 1 から 3 のいずれか 1 つに記載の音響信号処理装置。

20

**【請求項 5】**

発話者を検出する発話者検出部を有し、

前記グループ分け部は、前記 1 又は 2 以上の人物の中から、前記発話者だけを対象に、前記グループ分けを行うことを特徴とする請求項 1 から 4 のいずれか 1 つに記載の音響信号処理装置。

**【請求項 6】**

前記人物位置検出部により検出された人物の増減を検出する人物増減検出部を有し、

前記グループ分け部は、前記人物増減検出部により検出された前記人物の増減に応じて、前記グループ分けを行うことを特徴とする請求項 1 から 5 のいずれか 1 つに記載の音響信号処理装置。

30

**【請求項 7】**

前記指向性設定部によって指向性が設定された指向性制御部の出力監視あるいは前記人物位置検出部により検出された人物の発話監視を行う発話状態監視部を有し、

前記グループ分け部は、前記発話状態監視部により、前記指向性が設定された指向性制御部の出力あるいは前記検出された人物の発話が所定時間以上観測されなくなったときは、前記所定時間以上観測されなくなった指向性制御部あるいは人物を除いて、前記グループ分けを行うことを特徴とする請求項 1 から 6 のいずれか 1 つに記載の音響信号処理装置。

**【請求項 8】**

前記人物位置検出部は、前記所定空間をカメラにより撮像して得られた画像により、あるいは前記所定空間からの音響信号により、前記人物の位置を検出することを特徴とする請求項 1 から 7 のいずれか 1 つに記載の音響信号処理装置。

40

**【請求項 9】**

テレビジョン放送のコンテンツを表示する表示部と、

音響信号処理装置と、

通信回線を介して通信を行うための通信インターフェースと、

を有し、

前記音響信号処理装置は、

所定空間内に存在する人物の位置を検出する人物位置検出部と、

前記人物位置検出部により検出された 1 又は 2 以上の人物を所定数以下の数のグループ

50

に割り振るグループ分けを行うグループ分け部と、

設定されたそれぞれの指向性に基づいてマイクロホンアレーの指向性を制御する複数の指向性制御部と、

前記グループ分け部で前記グループ分けがされた各グループの指向性を、対応する指向性制御部に設定する指向性設定部と、

を有することを特徴とするテレビジョン装置。

【請求項 10】

音響信号を処理するプログラムであって、コンピュータに、

所定空間内に存在する人物の位置を検出する機能と、

検出された 1 又は 2 以上の人物を所定数以下の数のグループに割り振るグループ分けを行う機能と、

グループ分けがされた各グループの指向性を、マイクロホンアレーの指向性を制御する指向性制御部に設定する機能と、

を実現させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、音響信号処理装置、テレビジョン装置及びプログラムに関する。

【背景技術】

【0002】

従来より、マイクロホンアレーを用いて、目的とする方向からの音だけを強調する技術が知られている。

また、カメラを用いて人の位置を検出して、その検出した人の方向にマイクロホンアレーの指向性を向ける技術も提案されている。その提案に係る装置では、話者の位置を検出して、その検出された話者の位置関係に基づいて、話者毎に音声を抽出する方向と範囲が抽出され、音声抽出手段がその範囲内の音声を抽出する。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】特開 2005 - 274707 号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかし、そのような装置は、各話者に向かうようにマイクロホンアレーの指向性を制御して音声処理を行うようにするため、話者の数が増えると、抽出された各方向についての指向性制御部における音声信号の音声処理をするための演算量が増えるという問題がある。すなわち、話者数に応じて、指向性制御部の演算量が増える。

【0005】

よって、このような装置は、話者の増えた分だけ計算量が増えるため、特に、リアルタイムで応答が要求されるシステムに適用した場合、そのシステムのリアルタイム性が確保できなくなってしまうという問題がある。例えば、遠隔地の者同士がテレビ会議等を行うような場合、音声信号がリアルタイムで音声処理できなくなるため、音声の途切れ、雑音の混入等が発生してしまう。

また、所謂計算コストの面から、余裕を持った数の指向性制御部を予め想定してハードウェアの性能を高めることは現実的でない。

【0006】

そこで、本実施形態は、リアルタイム性を確保できる音響信号処理の可能な音響信号処理装置、テレビジョン装置及びプログラムを提供することを目的とする。

【課題を解決するための手段】

【0007】

10

20

30

40

50

実施形態の音響信号処理装置は、所定空間内に存在する人物の位置を検出する人物位置検出部と、前記人物位置検出部により検出された1又は2以上の人物を所定数以下の数のグループに割り振るグループ分けを行うグループ分け部と、設定されたそれぞれの指向性に基づいてマイクロホンアレーの指向性を制御する複数の指向性制御部と、前記グループ分け部で前記グループ分けがされた各グループの指向性を、対応する指向性制御部に設定する指向性設定部と、を有する。

【図面の簡単な説明】

【0008】

【図1】図1は、第1の実施の形態に係る音響信号処理装置が適用されるテレビ電話システムの例を説明するための図である。

【図2】第1の実施形態に係わる、各テレビ2と各ユニット4の構成を示すブロック図である。

【図3】第1の実施形態に係わる、ユニット4のソフトウェア構成を示すブロック図である。

【図4】第1の実施形態に係わる、人物位置検出部31によって、マイクロホンアレー12あるいはテレビ画面2aに対して、所定の位置を基準として、3人の人物が検出されたときに、2グループにグループ分けされた一の状態を説明するための図である。

【図5】第1の実施形態に係わる、マイクロホンアレー12の雑音抑圧性能の和が最大になるようにグループ分けを行う方法を説明するための図である。

【図6】第1の実施形態に係わる、ユニット4における指向性設定処理の流れの例を示すフローチャートである。

【図7】第1の実施形態に係わる、グループ数すなわちビーム数の変更を説明するための図である。

【図8】第2の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。

【図9】第2の実施の形態に係るユニット4の変形例に係るソフトウェア構成を示すブロック図である。

【図10】第2の実施の形態に係るユニット4における指向性設定処理の流れの例を示すフローチャートである。

【図11】第3の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。

【図12】第3の実施の形態に係るユニット4における指向性設定処理の流れの例を示すフローチャートである。

【図13】第4の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。

【図14】図13のユニット4の変形例に係るソフトウェア構成を示すブロック図である。

【図15】第4の実施の形態に係るユニット4における指向性設定処理の流れの例を示すフローチャートである。

【図16】第5の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。

【図17】マイクロホンアレー12の各マイクロホンからの音響信号のみから、人物の位置を検出するユニットのソフトウェア構成を示すブロック図である。

【発明を実施するための形態】

【0009】

以下、図面を参照して実施形態を説明する。

【0010】

(第1の実施形態)

(構成)

図1は、第1の実施の形態に係る音響信号処理装置が適用されるテレビ電話システムの

10

20

30

40

50

例を説明するための図である。

図 1 に示すように、テレビ電話システム 1 は、テレビジョン放送を受信するテレビジョン装置（以下、テレビという）を利用して、遠隔地の者同士が、テレビ画面に映し出された相手の画像を見ながら会話を行うことができるシステムである。テレビ電話システム 1 は、互いに離れた場所に設置された 2 台のテレビ 2A、2B を含み、テレビ 2A、2B は、互いに通信回線としてのインターネット 3 を介して通信可能に接続されている。

【 0 0 1 1 】

そして、各テレビ 2A、2B（以下、2 台のテレビの両方あるいは一方を指すとき、テレビ 2 ともいう）は、テレビジョン放送を受信可能であると共に、インターネット接続下で、インターネット 3 を介して互いに画像信号及び音声信号の通信が可能である。そのため、各テレビ 2A、2B には、それぞれカメラとマイクロホンアレーを内蔵するユニット 4A、4B（以下、2 つのユニットの両方あるいは一方を指すとき、ユニット 4 ともいう）が取り付けられている。

各テレビ 2 のテレビ画面 2 a には、テレビ 2 が放送受信モードのときは、テレビジョン放送のコンテンツが表示され、かつ後述するテレビ電話モードのときは、相手方の画像が表示される。

【 0 0 1 2 】

テレビ電話ユニットとしての各ユニット 4 は、対応するテレビ 2 の画面の前を撮像するカメラと、テレビ 2 の前の音声を取り込むマイクロホンアレーを有する（図 2）。ユーザは、リモコン 5A、5B（以下、2 つのリモコンの両方あるいは一方を指すとき、リモコン 5 ともいう）を操作することによって、テレビジョン放送をテレビ画面上に表示させたり、ユニット 4 のテレビ電話機能を動作させたりすることができる。ユニット 4 は、音響信号処理装置を構成する。

【 0 0 1 3 】

図 2 は、各テレビ 2 と各ユニット 4 の構成を示すブロック図である。

ユニット 4 は、カメラ 1 1、マイクロホンアレー 1 2、中央処理装置（CPU）1 3、ROM 1 4、RAM 1 5、及びインターフェース（以下、I/F と略す）1 6、1 7、1 8 を含む。CPU 1 3、ROM 1 4 及び RAM 1 5 は、バス 1 9 を介して互いに接続されている。テレビ 2 とユニット 4 は、インターネット 3 に接続するための通信ライン 2 0 に接続されている。

【 0 0 1 4 】

カメラ 1 1 は、CCD などの撮像素子を有して、デジタルの画像信号を出力する。マイクロホンアレー 1 2 は、複数の（ここでは、n 個の）マイクロホンを有し、各マイクロホンの音響信号を出力する。テレビ 2 に載置されたユニット 4 内において、カメラ 1 1 は、テレビ 2 のテレビ画面を視るユーザ側を撮像するように配置され、マイクロホンアレー 1 2 は、テレビ画面を視るユーザの音声を取り込むように配置されている。

【 0 0 1 5 】

I/F 1 6 は、カメラ 1 1 とバス 1 9 を接続するインターフェースである。I/F 1 7 は、マイクロホンアレー 1 2 とバス 1 9 を接続するインターフェースである。I/F 1 8 は、インターネット 3 に接続された通信ライン 2 0 とバス 1 9 とを接続するためのインターフェースである。通信ライン 2 0 には、テレビ 2 も接続されているので、テレビ 2 とユニット 4 は、互いに通信可能となっている。ユニット 4 と合わせてテレビ 2 は、テレビ電話可能なテレビジョン装置を構成する。

【 0 0 1 6 】

CPU 1 3 は、ROM 1 4 に格納された、後述する各種ソフトウェアプログラム（以下、単にプログラムという）を実行する処理部である。ROM 1 4 は、後述する各種プログラムを格納する不揮発性の記憶部である。RAM 1 5 は、CPU 1 3 が各種プログラムを実行するときに作業領域として利用する記憶部である。

【 0 0 1 7 】

ROM 1 4 には、テレビ電話システム 1 が機能するときに、ユニット 4 が音響信号処理装置としての機能を実行するための各種プログラムが含まれる。

10

20

30

40

50

なお、図1と図2に示すテレビ電話システム1では、テレビ2とユニット4が別体のものであり、ユニット4がテレビ2に載置されるようになっているが、ユニット4がテレビ2の本体内に内蔵されるような構成でもよい。

【0018】

図3は、ユニット4のソフトウェア構成を示すブロック図である。ユニット4内の各ブロックは、プログラムにより構成される。ユニット4は、人物位置検出部31、グループ分け部32、指向性設定部33、複数の(ここでは2つの)指向性制御部34a、34b、及び合成部35を有して構成される。

【0019】

カメラ11からのデジタルの画像信号が、人物位置検出部31に入力される。

10

人物位置検出部31は、入力された画像信号に基づいて、撮像された人物を判別し、各人物の位置を検出して、その位置情報をグループ分け部32に出力する。カメラ11は、テレビ2の前の所定空間を撮像するので、人物位置検出部31は、その所定空間内に存在する人物の位置を検出する処理部を構成する。人物位置検出部31は、画像の中から人の顔を認識する顔画像認識処理により、複数の人物のそれぞれの位置を検出する。検出された各顔の位置と、所定の基準位置との位置関係から、カメラ11が撮像した所定空間内における各人物の位置が算出される。すなわち、人物位置検出部31は、複数の顔が検出されれば、各顔の位置に対応する各人物の位置を算出する。

【0020】

グループ分け部32は、入力された位置情報に基づいて、検出された人物のグループ分けを行い、複数の人物が検出されれば、所定のグループ数のグループに纏められる。判別された1又は2以上の人物は、予め設定された上限のグループ数までグループ数にグループ分けされる。よって、グループ分け部32は、人物位置検出部31により検出された1又は2以上の人物を所定数以下の数のグループに割り振るグループ分けを行う処理部を構成する。

20

【0021】

例えば、所定数としての上限が2の場合、検出された人物が一人であれば、グループ数は1であり、検出された人物が二人であれば、グループ数は2となる。さらに、上限が2の場合、判別された人物が三人以上であっても、グループ分けされたグループ数は2となる。

30

【0022】

グループ数の上限は、CPU13の処理能力によって決定される。CPU13の処理能力に応じて、プログラムの処理時間は異なる。CPU13の処理能力が高ければ、CPU13における指向性制御部34における1つの指向性制御部の処理時間が短くなり、CPU13の処理能力が低ければ、1つの指向性制御部の処理時間が長くなる。

【0023】

特に、指向性制御部34において、指定された方向にマイクロホンアレー12のビームを形成して音声処理するための計算量が、テレビ電話システム1に要求されるリアルタイム性が維持できない計算量にならないように、上限が決定される。

【0024】

40

例えば、CPU13の音声処理のための処理能力が、100MIPS(百万命令毎秒)であるときに、1つの指向性制御部の処理時間が50MIPSであれば、2つの指向性制御部の処理がCPU13の処理能力の限界である。その場合、グループ数の上限は2となり、ユニット4は、指向性制御部を2つ有することができる。検出された人物についてのグループ分けの方法については、後述する。

【0025】

グループ分け部32は、グループ分けされた人物の位置情報に基づいて、グループ毎の指向性の情報(すなわちビームの情報)を算出して決定し、指向性設定部33に出力する。ビームの情報は、指向性の設定方向dsと設定範囲sの情報を含む。設定方向dsは、設定範囲sの中心方向である。例えば、設定範囲sは、設定方向dsを中心とする角度幅で

50

ある。

【0026】

なお、グループ分け部32は、所定のルールに従って検出された人物の位置情報に基づいてグループ分けを行うが、雑音抑圧能力すなわち性能が最も高くなるように、すなわち最適なグループ分けを行うようにしてもよい。最適なグループ分けの方法については後述する。

グループ分け部32は、グループ分けされたグループ毎のビームの情報を、指向性設定部33に供給する。

【0027】

指向性設定部33は、グループ毎のビームの情報に基づいて、各グループの話者の音声を強調するように、各ビームを形成するための設定情報をグループ毎に生成して、対応する指向性制御部34a、34bに各設定情報を供給する。例えば、検出された人物が一人の場合、指向性制御部34aのみに設定情報が供給されて設定され、指向性制御部34aのみがその設定情報に基づくビームを形成する。検出された人物が二人以上の場合、指向性制御部34a、34bのそれぞれに設定情報が供給されて設定され、指向性制御部34a、34bは、それぞれ設定情報に基づくビームを形成する。

すなわち、グループ分け部32は、人物のグルーピングを行い、指向性設定部33は、指向性をグループ単位で形成するように、各指向性制御部34を制御する。よって、指向性設定部33は、グループ分け部32でグループ分けがされた各グループの指向性を、対応する指向性制御部34に設定する処理部を構成する。

【0028】

複数の指向性制御部34a、34bは、設定されたそれぞれの指向性に基づいてマイクロホンアレーの指向性を制御する処理部を構成する。すなわち、各指向性制御部34は、設定された方向から到来する音声を強調する。指向性制御部34a、34bにおいて音声強調処理されたデジタルの音声信号は、合成部35で加算されて合成信号として、インターネット3へ送信される。指向性制御部34は、それぞれに設定された指向性を実現すべくn個の音声入力信号に対してアレー処理を行う。アレー処理の例は、例えば、特許第3795610号公報、特開2007-10897号公報に開示されている。

なお、指向性制御部34a、34bにおいて音声強調処理されて得られたデジタルの複数の音声信号は、合成部35で合成しないで、それぞれ個別に、インターネット3へ送信するようにしてもよい。

【0029】

上述したように、グループ分けにより生成された各グループの人物の音声を強調するのに適した指向性を形成するための設定が決定される。指向性制御部34がその設定に基づいて複数のマイクロホン12a~12nからの音声信号に対して所定のフィルタ演算を行い、その演算結果を加算する処理を行うことによって、マイクロホンアレー12の指向性すなわちビームは形成される。ここでは、上限数である2つの指向性制御部34a、34bは、それぞれ、互いに異なる特性のフィルターセットを動作させるようにしてビームを形成する。

【0030】

なお、本実施の形態及びこれに続く他の実施の形態においても、話者の位置は、水平方向における角度で特定する場合で説明するが、これに限られず、複数のマイクロホンを2次元配置等することによって、奥行き方向も含めた複数の話者のグループ化を行って、話者の存在する空間の奥行き方向の制御を行うようにしてもよい。

【0031】

(グループ分けの方法)

ここで、グループ分け部32におけるグループ分けの方法について説明する。図4は、人物位置検出部31によって、マイクロホンアレー12あるいはテレビ画面2aに対して、所定の位置を基準として、3人の人物が検出されたときに、2グループにグループ分けされた一の状態を説明するための図である。

10

20

30

40

50

## 【 0 0 3 2 】

図 4 は、3 人の人物 P1, P2, P3 が、それぞれ、マイクロホンアレー 1 2 の所定の中心位置 P0 から、方向 d1, d2, d3 に存在することが検出された場合であって、人物 P1 と P2 が第 1 のグループを構成し、人物 P3 が第 2 のグループを構成する例を示す。

## 【 0 0 3 3 】

例えば、検出された各人物の顔の中心の方向が、人物の方向として決定される。図 4 の場合、人物 P1 は、方向 d1 に存在する。同様に、人物 P2 と P3 は、それぞれ方向 d2 と d3 に存在する。

## 【 0 0 3 4 】

(ルール)

検出された人物に対するグループ化は、所定のルールで行われるが、所定のルールは、種々のルールが適用可能である。まず、検出された人物の数がグループ数の上限に達するまでは、検出された人毎にグループ分けが行われ、その人物の数が上限を超えると、所定のルールによるグループ化が行われる。所定のルールの中で簡単なルールは、例えば、マイクロホンアレー 1 2 の所定の基準方向 (0 度) から所定の方向 (180 度の方向) に向かって、検出された人物を、各グループが所定の人数になるようにグループ化する、というようなルールである。

## 【 0 0 3 5 】

図 4 の場合、検出された三人の人物 P1, P2, P3 が、0 度方向から 180 度方向に向かって、二人のグループと一人のグループにグループ分けされている。人物 P1 と P2 が第 1 のグループとなり、人物 P3 がもう一つの第 2 のグループとなっている。第 1 のグループでは、方向 d1 と d2 の真ん中の方向 D1 が、目的音源方向としての設定方向 ds となる。第 2 のグループでは、方向 d3 の方向 D2 が、目的音源方向としての設定方向 ds となる。さらに、第 1 のグループでは、目的音源方向 D1 の前後  $(\theta_1) / 2$  の範囲  $\theta_1$  が、設定範囲  $\theta_s$  となる。第 2 のグループでは、目的音源方向 D2 の前後  $(\theta_2) / 2$  の範囲  $\theta_2$  が、設定範囲  $\theta_s$  となる。

## 【 0 0 3 6 】

ここでは、図 4 の第 2 のグループのように、一つのグループに人物が一人だけいる場合、設定範囲  $\theta_s$  は、所定の範囲  $\theta_p$  を有するものとする。図 4 の第 1 のグループのように、一つのグループに人物が二人以上存在する場合、それぞれの所定の前後  $(\theta_p) / 2$  を加味した人物 P1, P2 間の角度  $\theta_d$  を含むように、設定範囲  $\theta_s$  は、角度  $\theta_1$  となる。

## 【 0 0 3 7 】

他にも、グループ分け部 3 2 は、各グループに割り振られる人物の数が平均化するというルールに基づいて、グループ分けを行うようにしてもよい。

以上のように、グループ分けは、所定のルールに従って行われる。その所定のルールは、グループ数の上限に達するまでは、検出された人毎にグループ分けが行われ、検出された人物の数がグループ数の上限を超えると、所定のルールで、グループ数が上限を超えないようにグループ分けするルールである。

## 【 0 0 3 8 】

(最適なグループ分け)

グループ分け部 3 2 は、マイクロホンアレー 1 2 の雑音抑圧性能の和が最大になるように最適なグループ分けを行うようにしてもよい。

図 5 は、マイクロホンアレー 1 2 の雑音抑圧性能の和が最大になるようにグループ分けを行う方法を説明するための図である。人物の配置は、図 4 と同じであるが、3 人の人物 P1, P2, P3 は、図 4 とは異なるグループにグループ分けされている。

## 【 0 0 3 9 】

図 5 の場合、人物 P1 が第 1 のグループであり、人物 P2, P3 が第 2 のグループとなっている。第 1 のグループ (人物 P1 のグループ) は、図 4 の第 2 のグループ (人物 P3 のグループ) と同じ設定範囲  $\theta_s$  を有する。第 2 のグループ (人物 P2, P3 のグループ) は、図 4 の第 1 のグループ (人物 P1, P2 のグループ) よりも、設定範囲  $\theta_s$  が狭くなっている。

## 【 0 0 4 0 】

10

20

30

40

50

図5において、第2のグループの設定範囲  $s$  は、角度  $3$  であり、図4の第1のグループの設定範囲  $s$  である  $1$  よりも狭い。

設定範囲  $s$  が狭い方が、目的とする方向からの音だけを強調する性能が高い。図4の場合、人物P2とP3の間には隙間がなく、仮にこの隙間から雑音が到来しても、抑圧することはできない。これに対して、図5の場合は、人物P1とP2の間には隙間があり、この隙間から雑音が到来しても抑圧することができる。

【0041】

よって、図5に示すグループ分けは、図4のグループ分けよりも、よりよいグループ化ということが言える。これは、互いに異なる設定範囲  $s$  を有する2つのグループ分けがあったときに、評価値EVとしての  $(1/s)$  の和が大きいグループ分けの方が、他のグループ分けよりも、システム全体としての強調性能が高いので、図5の方が図4よりも、よりよいグループ分けとすることができるからである。

10

【0042】

従って、最適なグループ分けの方法としては、検出された人物の数が、グループ数の上限を超えた場合、検出された全ての人物について取り得るグループ分けの組み合わせを仮定し、取り得る組み合わせの中で、各グループの評価値EV  $(= 1/s)$  の総和EVsが最も大きい組み合わせのグループ分けが、最適なグループ分けであるというルールを用いることができる。すなわち、グループ数の上限の範囲内で、全ての組み合わせについて所定の評価値の総和EVsを算出し、その総和EVsが最も大きな値のグループ分けを選択するというルールを、評価基準として、グループ分け部32が採用することができる。

20

【0043】

例えば、図4において、第1のグループの設定範囲  $s (= 1)$  が65度、第2のグループの設定範囲  $s (= p)$  が20度とすると、上記の評価値の総和EVs1は、 $((1/s) + (1/p)) = ((1/65) + (1/20))$  である。これに対して、図5において、第1のグループの設定範囲  $s (= p)$  が20度、第2のグループの設定範囲  $3$  が38度とすると、上記の評価値の総和EVs2は、 $((1/p) + (1/3)) = ((1/20) + (1/38))$  である。EVs2 > EVs1となるので、図5の組み合わせの方が、図4の組み合わせよりも、評価値EVが高い。

【0044】

よって、グループ分け部32は、このような評価値の総和EVsの比較を、全てのグループ分け可能な組み合わせ同士で行い、最も評価値の総和EVsが高い組み合わせのグループ分けを、最適なグループ分けとして決定する。

30

【0045】

(処理)

図6は、ユニット4における指向性設定処理の流れの例を示すフローチャートである。図6の処理は、テレビ電話システム1におけるテレビ電話機能がオンになると、CPU13によって実行される。ユーザは、リモコン5を操作して、テレビ電話機能をオンにすることができる。

【0046】

まず、CPU13は、人物位置検出部31により人物位置の検出を行い(S1)、続いて、グループ分け部32による、上述したようなグループ分けを行う(S2)。CPU13は、指向性設定部33による、グループ分けの結果に基づいて各指向性制御部34に対する指向性の設定を行う(S3)。

40

【0047】

指向性制御部34は、設定された指向性の情報に基づいて、ビームを制御して音声処理を行う。指向性制御部34で処理された音声信号は、合成部35で合成されて、通信回線であるインターネット3を介して、相手方のテレビ電話システムへ伝送される。

【0048】

図6の処理が実行された後、において、人物位置検出部31により検出された人物の位置に変化がなければ、図6の処理は実行されないが、人物の位置に変化があると、図6の

50

グループ分け及び指向性設定の処理が実行される。

【0049】

その結果、遠隔に離れた者同士が、テレビ2を利用して、リアルタイムで音声が届いたり等なく、テレビ電話を行うことができる。

【0050】

なお、上述した例では、上限数まではグループ数が増えていくが、ユーザによって、予め設定された上限を変更できるようにしてもよい。

図7は、グループ数すなわちビーム数の変更を説明するための図である。ユーザは、自己が見ているテレビ2のリモコン5に対して所定の操作を行うと、図7に示すようなビーム数上限の設定変更画面21を表示させることができる。設定変更画面21は、テレビ2の画面2a上に、サブウィンドウとして表示される。設定変更画面21を表示させるプログラムは、ROM14中に予め格納されている。

10

【0051】

設定変更画面21は、変更可能なビーム数を表示する表示部21Aを含み、ユーザは、リモコン5を操作して、カーソル(図7では斜線で示されている)を所望のビーム数の表示部21Aに移動させて選択することができる。例えば、図7では、上限数の「2」よりも少ない「1」の表示部が選択されている。選択の確定は、ユーザが、リモコン5の確定ボタンなどを操作することによって行うことができる。

なお、上限が「3」であれば、上限より少ない「2」と「1」が選択可能となるように、設定画面変更画面21には、上限の数と、上限よりも少ない数の表示部21Aが、選択可能に表示される。

20

【0052】

例えば、ユーザが上限「2」よりも少ないビーム数である「1」を選択すると、ビーム形成に必要な指向性制御部34の数が少なくなるので(1つの指向性制御部でよいことになるので)、上述したCPU13の演算量が少なくなる。その場合、CPU13の処理時間に余裕ができるので、より処理時間のかかる別の指向性制御部34A(図3では点線で示す)を利用して、例えば、より高音質な音声信号が得られる指向性制御を実行させるようにすることも可能である。

【0053】

例えば、一人しかいないような状況では、ユーザは、設定変更画面21を用いて、ビーム数の上限を1に設定して、より高性能な指向性制御部34Aによるビーム形成を行わせることができる。その結果、相手方には、より高音質で音声信号を送信することができる。

30

【0054】

ROM14に格納された図7の画面に関する処理を実行するプログラムは、上述したビーム数言い換えるグループ数を変更するためのグループ数変更部36(図3で点線で示す)を構成する。

なお、上述した例では、グループ数変更部36は、グループ数を上限数よりも少ない数に変更できる例であるが、予め設定された上限を、より少ない数に変更するようにしてもよい。よって、グループ数変更部36は、グループ分けの結果のグループ数あるいは所定の上限を変更するための処理部を構成する。

40

【0055】

以上のように、上述した本実施の形態によれば、CPUの処理能力などに応じたグループ数の上限内でグループ化が行われるので、リアルタイム性を確保できる音響信号処理が可能なテレビ電話システムを提供することができる。

【0056】

(第2の実施の形態)

第1の実施の形態では、顔検出された人物の位置を検出して、グループ分けを行っていたが、第2の実施の形態は、検出された人物の中で発話した人物のみを、グループ分けするようにした点が、第1の実施の形態と異なる。

50

以下、第2の実施の形態を説明するが、第1の実施の形態と同様の構成については、同じ符号を付し説明は省略し、異なる点を主として説明する。

【0057】

図8は、第2の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。図8は、発話者検出部41を含む点で、図3と異なる。発話者検出部41は、人物位置検出部31により検出された人物の中から、発話者を検出する処理部である。

【0058】

よって、本実施の形態は、発話者検出部41によって実際に発話したかどうかを検出し、発話が検出された人物の位置情報のみが、後段のグループ分け部32に出力される。グループ分け部32は、人物位置と検出された発話者との対応関係の情報を有している。

10

【0059】

なお、発話者検出部41にマイクロホンアレー12からの音声信号を入力させ、口の動きと共に、音声信号が入力された場合に、発話者であるという判定をするようにしてもよい。図9は、本実施の形態に係るユニット4の変形例に係るソフトウェア構成を示すブロック図である。図9では、口の動きと音声信号による判定のために、マイクロホンアレー12の各マイクロホンからの音声信号が、発話者検出部41に入力されている。

【0060】

図10は、本実施の形態に係るユニット4における指向性設定処理の流れの例を示すフローチャートである。図10の処理は、人物位置の検出の後に、発話者検出部41による発話者の検出処理を行い(S11)、発話者が検出されたか否かを判定し(S12)、発話者が検出された場合に(S12:YES)、グループ分けが行われるようになっている。そして、グループ分け部32は、人物位置検出部31により検出された1又は2以上の人物の中から、発話者だけを対象に、グループ分けを行う。

20

【0061】

よって、存在はしているが、発話しない人物にマイクロホンアレー12のビームを向けしておくような無駄をなくすことができ、さらに人物の誤検出による人でない対象に指向性が向けられるということがないようにすることができる。また、発話者が検出されると、グルーピングを再度行うので、常に最適なグループ分けの状態、テレビ電話システム1は動作することができる。

【0062】

30

また、発話者検出部41を追加することにより、人物位置検出部31の検出感度を、人物を検出し易いように高く設定できるので、人物位置検出部31において、誤検出はあっても、検出漏れのないようにして、検出漏れによる指向性の設定がされないという問題を回避することもできる。

【0063】

(第3の実施の形態)

第2の実施の形態では、顔検出された人物の中から発話した人物のみを、グループ分けするようにしているが、第3の実施の形態は、人物の増減を監視し、増減に応じてグループ分けを行うようにした点が、第2の実施の形態と異なる。

【0064】

40

以下、第3の実施の形態を説明するが、第1及び第2の実施の形態と同様の構成については、同じ符号を付し説明は省略し、異なる点を主として説明する。

【0065】

図11は、第3の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。図11は、人物増減検出部42を含む点で、図3と異なる。人物増減検出部42は、人物位置検出部31により検出された人物の増減を検出する処理部である。

【0066】

図12は、本実施の形態に係るユニット4における指向性設定処理の流れの例を示すフローチャートである。図12の処理は、人物位置の検出の後に、人物増減検出部42による人物の増減の検出処理を行い(S21)、検出された人物の数に増減があったか否かを判

50

定し (S22)、その増減があった場合に (S22:YES)、グループ分けが行われるようになっている。そして、グループ分け部 3 2 は、人物増減検出部 4 2 により検出された人物の増減に応じて、グループ分けを行う。

【0067】

よって、存在はしているが、途中から新たな人物がテレビ電話に参加したり、逆にそれまで居た人物が途中で退席してテレビ電話に参加しなくなったような場合にも、適切な指向性の制御が可能となる。また、人物の増減が検出されると、グルーピングを再度行うので、常に最適なグループ分けの状態で、テレビ電話システムは動作することができる。

【0068】

(第4の実施の形態)

第3の実施の形態では、人物の増減に応じてグループ分けするようにしているが、第4の実施の形態は、人物に増減はないが、途中から発話がなくなった人物がいる場合には、そのような人物の位置情報を一旦削除して、再度グループ分けを行うようにした点が、第1, 第2及び第3の実施の形態と異なる。

【0069】

以下、第4の実施の形態を説明するが、第1, 第2及び第3の実施の形態と同様の構成については、同じ符号を付し説明は省略し、異なる点を主として説明する。

【0070】

図13は、第4の実施の形態に係るユニット4のソフトウェア構成を示すブロック図である。ユニット4内の各ブロックは、プログラムにより構成される。図13は、不活性ビーム検出部43を含む点で、図3と異なる。不活性ビーム検出部43は各指向性制御部34の出力に基づいて、不活性ビームの検出を行う。

【0071】

不活性ビーム検出部43は、そのビーム方向からの発話が所定の時間以上に亘って観測されなくなったか否かを検出する。不活性ビーム検出部43は、指向性設定部33によって指向性が設定された指向性制御部34の出力監視を行い、発話状態を監視する発話状態監視部を構成する。不活性ビーム検出部43は、発話が観測されなくなったビームを検出すると、そのビームに対応する位置の人物の位置情報を削除した人物位置情報をグループ分け部32に出力する。グループ分け部32は、その人物位置情報に基づいて、再度グループ分けする。

【0072】

このような構成によれば、例えば途中で寝てしまった人物がいたような場合に、発話が所定の時間観測されなくなると、無駄なビームを削除することができる。

【0073】

なお、上述した図13の構成では、1つのグループの中に複数人いる場合は、全員が発話を止めた場合にのみ、再度グループ分けされるが、不活性ビーム検出部51の構成を変更することによって、グループの中の一人でも発話を止めた場合に、再グループ分けを行うようにしてもよい。

【0074】

図14は、図13のユニット4の変形例に係るソフトウェア構成を示すブロック図である。図14では、マイクロホンアレー12の各マイクロホンからの音声信号が、不活性ビーム検出部43aに入力されている。不活性ビーム検出部43aは、検出された人物の数kだけ、指向性制御部51-a~51-k(以下、複数の指向性制御部51-a等を指すとき、あるいは一つの指向性制御部を指すとき、指向性制御部51という)を有している。

【0075】

検出された人物毎に指向性制御部51が生成される。指向性制御部51は、人物位置検出部31からの人物位置の人物の発話が所定時間観測されなかったか否かを検出する。従って、不活性ビーム検出部43aは、人物位置検出部31により検出された人物の発話監視を行う発話状態監視部を構成する。

【0076】

10

20

30

40

50

なお、指向性制御部 5 1 のプログラムは、テレビ電話によるリアルタイムな通話目的のものではないので、より高性能で処理時間のかかるプログラムであってもよい。

【 0 0 7 7 】

不活性ビーム検出部 4 3 a は、人物位置検出部 3 1 で検出された人物の発話が検出されなくなると、その人物の位置情報を削除した人物位置情報をグループ分け部 3 2 に出力する。グループ分け部 3 2 は、その人物位置情報に基づいて、再度グループ分けする。

【 0 0 7 8 】

図 1 5 は、本実施の形態に係るユニット 4 における指向性設定処理の流れの例を示すフローチャートである。図 1 5 の処理は、人物位置の検出の後に、不活性ビーム検出部 4 3 による不活性ビームの検出処理（あるいは図 1 4 の不活性ビーム検出部 4 3 a による発話を止めた人物の検出処理）を行い（S31）、不活性ビーム（あるいは発話を止めた人物）が検出されたか否かを判定し（S32）、その不活性ビーム（あるいは図 1 4 の発話を止めた人物）があった場合に（S32:YES）、グループ分けが再度行われるようになっている。

10

【 0 0 7 9 】

よって、グループ分け部 3 2 は、発話状態監視部である不活性ビーム検出部 4 3 , 4 3 a により、指向性が設定された指向性制御部 3 4 の出力あるいは検出された人物の発話が所定時間以上観測されなくなったときは、所定時間以上観測されなくなった指向性制御部あるいは人物を除いて、グループ分けを行う。

【 0 0 8 0 】

なお、上述した不活性ビームの検出などにより、グループ分けの対象から外された人物も、上述した第 2 の実施の形態の処理を利用することによって、グループ分けの対象人物として、再度加入することができる。

20

【 0 0 8 1 】

さらになお、発話をしなくなったグループあるいは人物を画像データからのみに基づいて検出する場合は、不活性ビーム検出部 4 3 , 4 3 a への、図 1 3 及び図 1 4 に示すような、指向性制御部 3 4 及びマイクロホンアレー 1 2 からの入力は不要である。

【 0 0 8 2 】

よって、本実施の形態によれば、発話がなくなったビームあるいは人物を検出して、無駄なビームを削除することができる。また、発話の有無が検出されると、グルーピングを再度行うので、常に最適なグループ分けの状態、テレビ電話システムは動作することができる。

30

【 0 0 8 3 】

（第 5 の実施の形態）

第 5 の実施の形態に係るユニット 4 は、上述した第 2 から第 4 の実施の形態において説明した発話者検出部 4 1 , 4 1 a、人物増減検出部 4 2 , 及び不活性ビーム検出部 4 3 を含む検出統合部 4 4 を有する。

【 0 0 8 4 】

図 1 6 は、本実施の形態に係るユニット 4 のソフトウェア構成を示すブロック図である。図 1 6 に示すように、ユニット 4 は、発話者検出部 4 1 , 4 1 a、人物増減検出部 4 2 , 及び不活性ビーム検出部 4 3 を含む検出統合部 4 4 を有する。なお、これらの 3 つの検出部 4 1 , 4 2 , 4 3 の全てを含まなくてもよい。なお、図 1 6 は、m 個の指向性制御部 3 4 を有する例である。

40

【 0 0 8 5 】

その結果、ユニット 4 は、第 2 から第 4 の実施の形態で説明した利点を含むユニットとなるので、テレビ電話システム 1 は、ユーザにとって、より使いやすいものとなる。また、発話者の検出、人物の増減、あるいは発話の有無が検出されると、グルーピングを再度行うので、常に最適なグループ分けの状態、テレビ電話システムは動作することができる。

【 0 0 8 6 】

以上のように、上述した第 1 から第 5 の実施の形態によれば、リアルタイム性を確保で

50

きる音響信号処理が可能な音響信号処理装置、テレビジョン装置、テレビ電話システムを提供することができる。特に、指向性がグループ単位で形成されるため、人物が多い場合、計算リソースが限られている場合等においても、少ない計算量で話者全員をカバーし、かつ雑音抑圧性能の高い音響信号処理装置、及びその応用装置を実現することができる。

【0087】

なお、上述した各実施の形態では、カメラ11を用いて人物の位置を検出しているが、カメラを用いずに、人物を検出するようにしてもよい。

【0088】

図17は、マイクロホンアレー12の各マイクロホンからの音響信号のみから、人物の位置を検出するユニットのソフトウェア構成を示すブロック図である。

10

【0089】

人物位置検出部31Aは、複数のマイクロホンの音響信号から音の到来方向を推定する処理部である。推定方法としては、MUSIC法、ESPRIT法などを用いたDOA (Direction Of Arrival) 推定が利用可能である。

DOA推定については、例えば、菊間信良著、「アレーアンテナによる適応信号処理」(科学技術出版、2004年)の第10章などを参照されたし。

【0090】

図17の構成によれば、カメラが不要となるので、コストの低減を図ることができる。さらに、人物以外の音が少ない環境では、雑音方向に誤って人物を検出してしまう等の誤動作も発生しにくい。

20

【0091】

なお、以上説明した動作を実行するプログラムは、コンピュータプログラム製品として、フレキシブルディスク、CD-ROM等の可搬媒体や、ハードディスク等の記憶媒体に、その全体あるいは一部を記録し、あるいは記憶するようにしてもよい。そのプログラムは、コンピュータにより読み取られて、動作の全部あるいは一部が実行される。あるいは、そのプログラムの全体あるいは一部を通信ネットワークを介して流通または提供することができる。利用者は、通信ネットワークを介してそのプログラムをダウンロードしてコンピュータにインストールしたり、あるいは記録媒体からコンピュータにインストールすることで、容易に上述した実施形態の音響信号処理装置、テレビジョン装置、テレビ電話システムを実現することができる。

30

【0092】

本発明は、上述した実施の形態に限定されるものではなく、本発明の要旨を変えない範囲において、種々の変更、改変等が可能である。

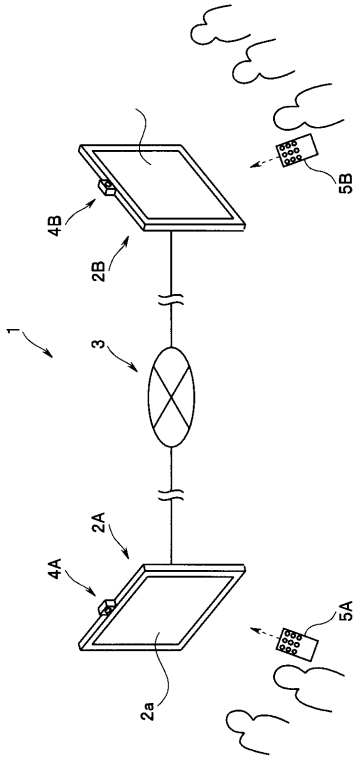
【符号の説明】

【0093】

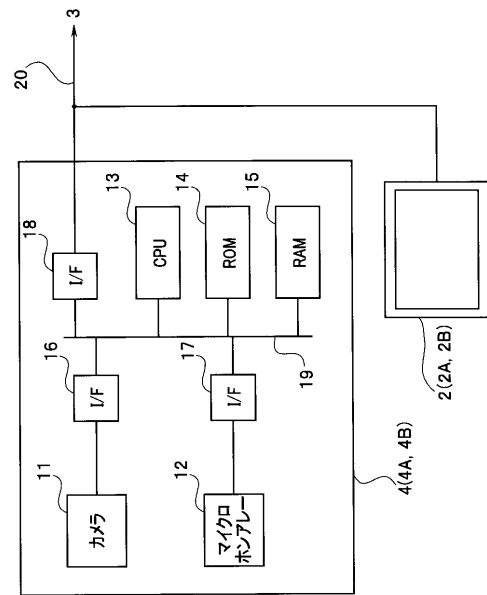
1 テレビ電話システム、2、2A、2B テレビ、2a テレビ画面、3 インターネット、4、4A、4B ユニット、5、5A、5B リモコン、11 カメラ、12 マイクロホンアレー、13 CPU、14 ROM、15 RAM、16、17、18 I/F、19 バス、20 通信ライン、31 人物位置検出部、32 グループ分け部、33 指向性設定部、34 指向性制御部、35 合成部、36 グループ数変更部、41、41a 発話者検出部、42 人物増減検出部、43、43a 不活性ビーム検出部

40

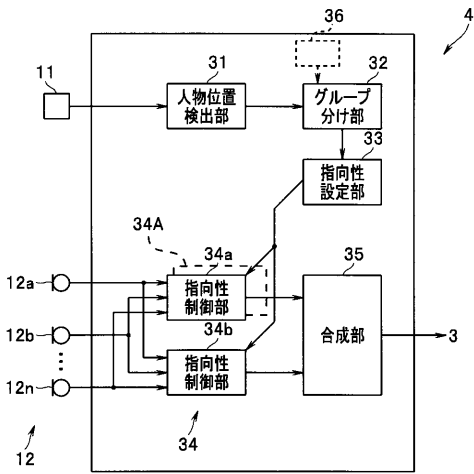
【図1】



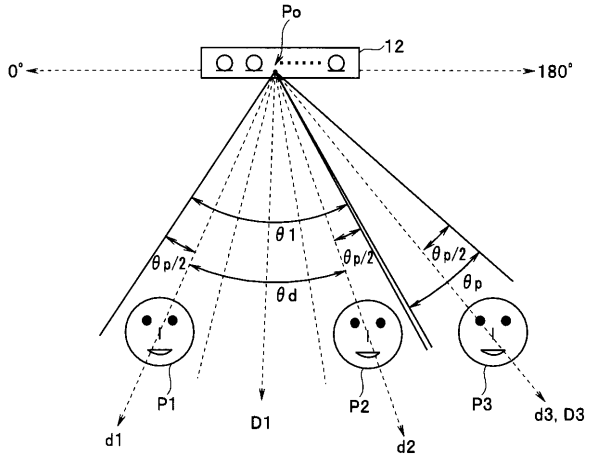
【図2】



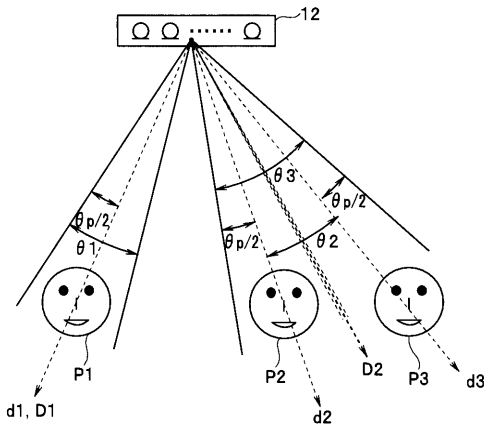
【図3】



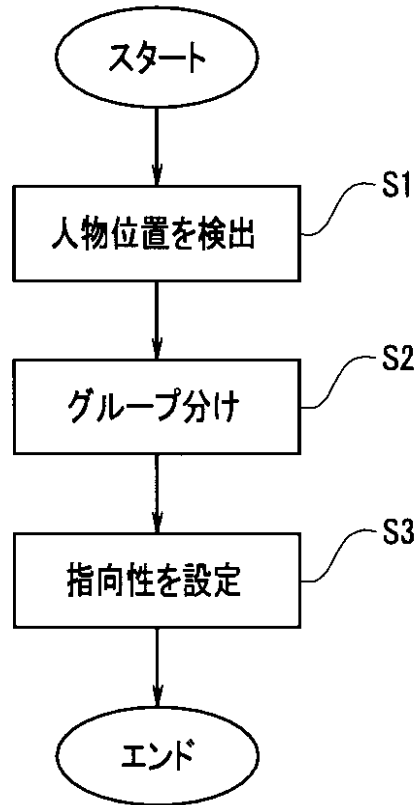
【図4】



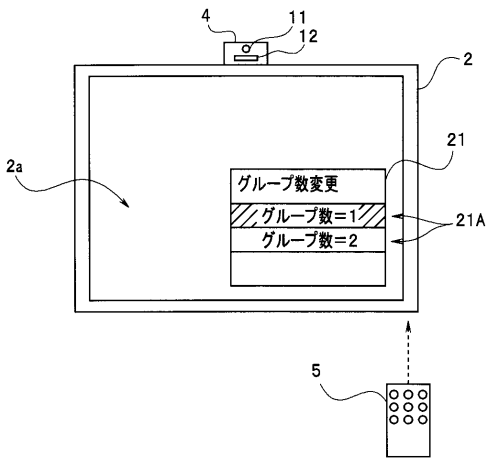
【図5】



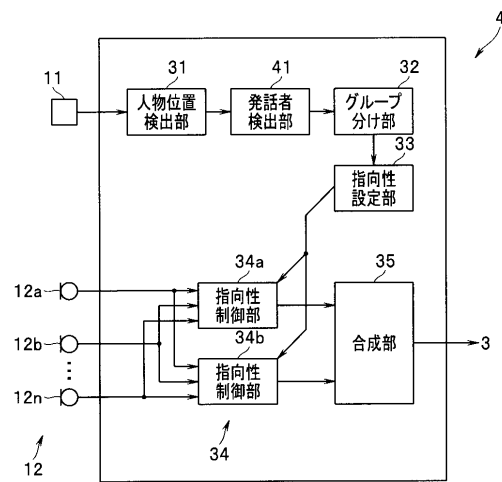
【図6】



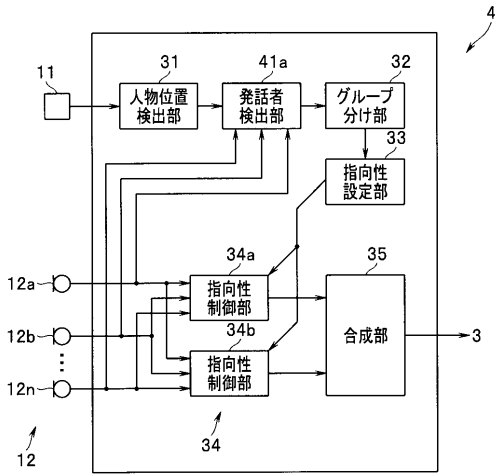
【図7】



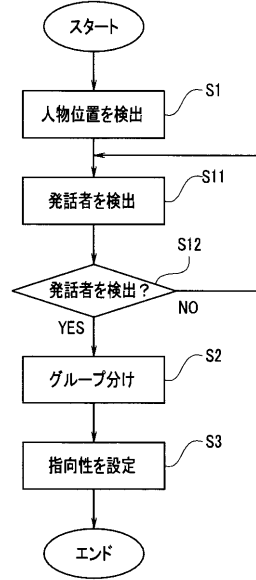
【図8】



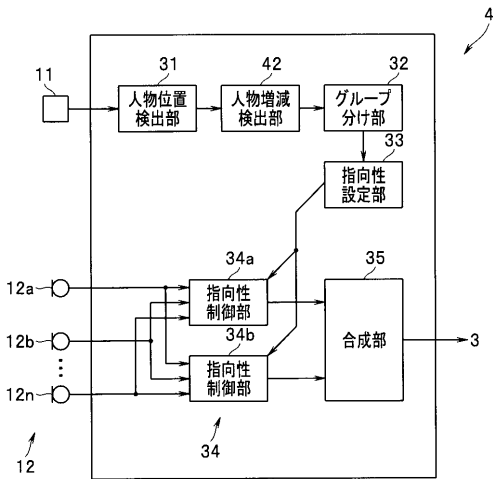
【図9】



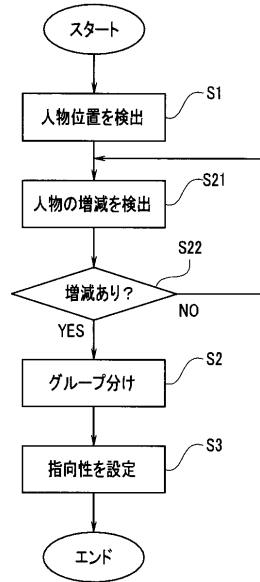
【図10】



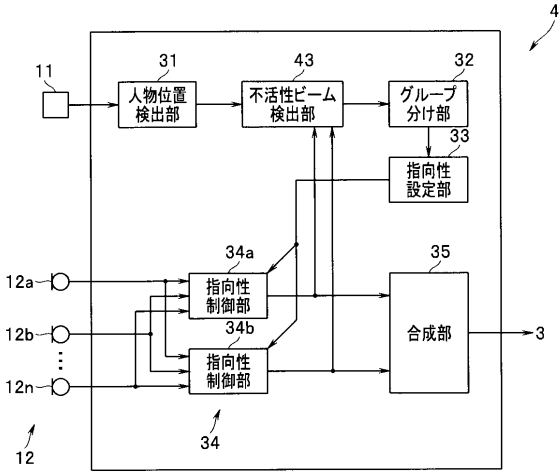
【図11】



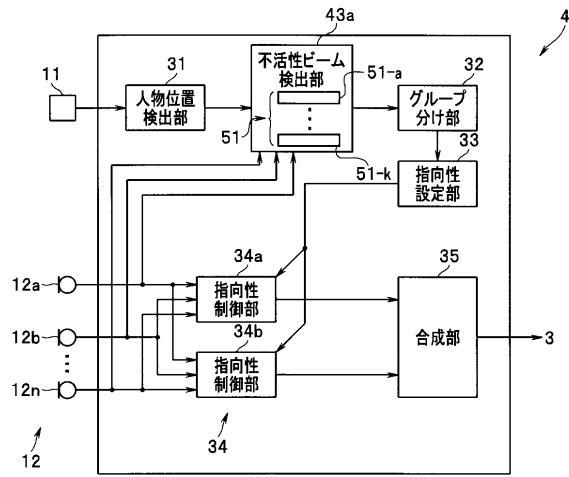
【図12】



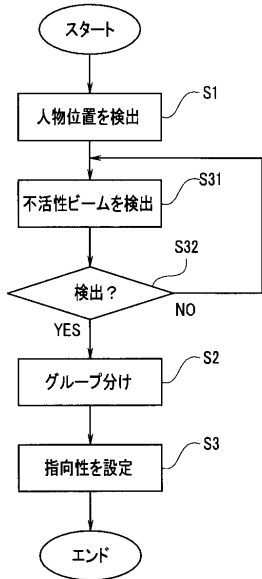
【図 1 3】



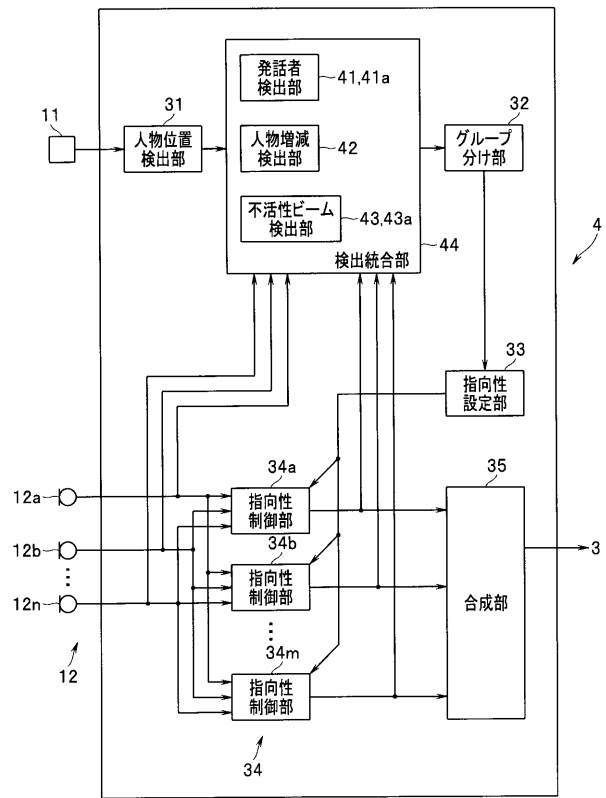
【図 1 4】



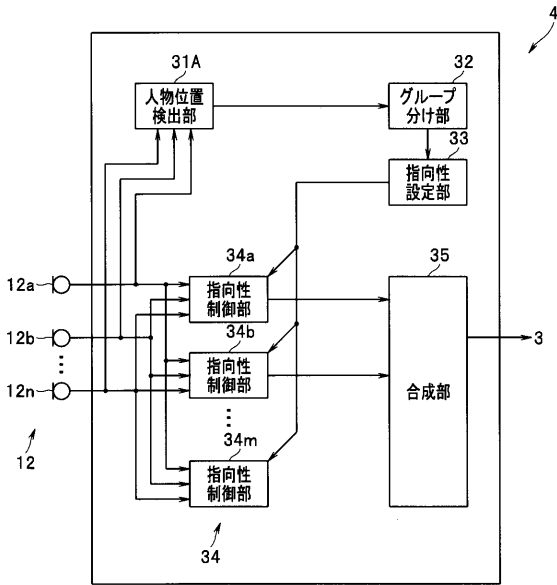
【図 1 5】



【図 1 6】



【図 17】



## 【手続補正書】

【提出日】平成23年11月22日(2011.11.22)

## 【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】0007

【補正方法】変更

【補正の内容】

【0007】

実施形態の音響信号処理装置は、所定空間内に存在する人物の位置を検出する人物位置検出部と、前記人物位置検出部により検出された2以上の人物の数が2以上の所定のグループ数を超えると、前記検出された2以上の人物を、前記所定のグループ数以下の数のグループに割り振るグループ分けを行うグループ分け部と、前記グループ分け部で前記グループ分けがされた各グループの指向性を、マイクロホンアレーの指向性を制御する複数の指向性制御部に設定する指向性設定部と、を有する。

## 【手続補正2】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

所定空間内に存在する人物の位置を検出する人物位置検出部と、

前記人物位置検出部により検出された2以上の人物の数が2以上の所定のグループ数を超えると、前記検出された2以上の人物を、前記所定のグループ数以下の数のグループに割り振るグループ分けを行うグループ分け部と、

前記グループ分け部で前記グループ分けがされた各グループの指向性を、マイクロホンアレーの指向性を制御する複数の指向性制御部に設定する指向性設定部と、  
を有することを特徴とする音響信号処理装置。

【請求項 2】

前記グループ分け部は、前記マイクロホンアレーの雑音抑圧性能の和が最大になるように、前記グループ分けを行うことを特徴とする請求項 1 に記載の音響信号処理装置。

【請求項 3】

前記グループ分け部は、前記各グループに割り振られる人物の数が平均化するように、前記グループ分けを行うことを特徴とする請求項 1 に記載の音響信号処理装置。

【請求項 4】

前記グループ分けの結果のグループ数あるいは前記所定数を、変更するためのグループ数変更部を有することを特徴とする請求項 1 から 3 のいずれか 1 つに記載の音響信号処理装置。

【請求項 5】

発話者を検出する発話者検出部を有し、

前記グループ分け部は、前記 2 以上の人物の中から、前記発話者だけを対象に、前記グループ分けを行うことを特徴とする請求項 1 から 4 のいずれか 1 つに記載の音響信号処理装置。

【請求項 6】

前記人物位置検出部により検出された人物の増減を検出する人物増減検出部を有し、

前記グループ分け部は、前記人物増減検出部により検出された前記人物の増減に応じて、前記グループ分けを行うことを特徴とする請求項 1 から 5 のいずれか 1 つに記載の音響信号処理装置。

【請求項 7】

前記指向性設定部によって指向性が設定された指向性制御部の出力監視あるいは前記人物位置検出部により検出された人物の発話監視を行う発話状態監視部を有し、

前記グループ分け部は、前記発話状態監視部により、前記指向性が設定された指向性制御部の出力あるいは前記検出された人物の発話が所定時間以上観測されなくなったときは、前記所定時間以上観測されなくなった指向性制御部あるいは人物を除いて、前記グループ分けを行うことを特徴とする請求項 1 から 6 のいずれか 1 つに記載の音響信号処理装置。

【請求項 8】

前記人物位置検出部は、前記所定空間をカメラにより撮像して得られた画像により、あるいは前記所定空間からの音響信号により、前記人物の位置を検出することを特徴とする請求項 1 から 7 のいずれか 1 つに記載の音響信号処理装置。

【請求項 9】

テレビジョン放送のコンテンツを表示する表示部と、

音響信号処理装置と、

通信回線を介して通信を行うための通信インターフェースと、  
を有し、

前記音響信号処理装置は、

所定空間内に存在する人物の位置を検出する人物位置検出部と、

前記人物位置検出部により検出された 2 以上の人物の数が 2 以上の所定のグループ数を超えると、前記検出された 2 以上の人物を、前記所定のグループ数以下の数のグループに割り振るグループ分けを行うグループ分け部と、

前記グループ分け部で前記グループ分けがされた各グループの指向性を、マイクロホンアレーの指向性を制御する複数の指向性制御部に設定する指向性設定部と、  
を有することを特徴とするテレビジョン装置。

【請求項 10】

音響信号を処理するプログラムであって、コンピュータに、

所定空間内に存在する人物の位置を検出する機能と、  
検出された2以上の人物の数が2以上の所定のグループ数を超えると、前記検出された2以上の人物を、前記所定のグループ数以下の数のグループに割り振るグループ分けを行う機能と、

グループ分けがされた各グループの指向性を、マイクロホンアレーの指向性を制御する複数の指向性制御部に設定する機能と、  
を実現させるためのプログラム。