

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2022-62876

(P2022-62876A)

(43)公開日 令和4年4月21日(2022.4.21)

(51)国際特許分類 F I テーマコード(参考)
 H 0 4 R 3/00 (2006.01) H 0 4 R 3/00 3 2 0 5 D 2 2 0

審査請求 未請求 請求項の数 14 O L (全13頁)

(21)出願番号	特願2020-171052(P2020-171052)	(71)出願人	000004075 ヤマハ株式会社 静岡県浜松市中区中沢町10番1号
(22)出願日	令和2年10月9日(2020.10.9)	(74)代理人	110000970 特許業務法人 楓国際特許事務所
		(72)発明者	鶴飼 訓史 静岡県浜松市中区中沢町10番1号 ヤマハ株式会社内
		(72)発明者	田中 良 静岡県浜松市中区中沢町10番1号 ヤマハ株式会社内
		Fターム(参考)	5D220 BA06 BC02 BC08

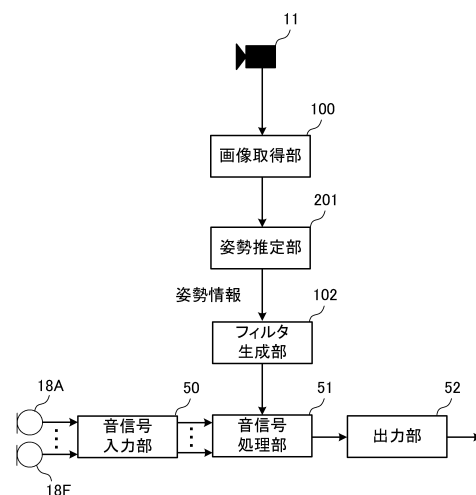
(54)【発明の名称】 音信号処理方法および音信号処理装置

(57)【要約】

【課題】話者の姿勢に応じて適切に話者の音声を取得できる音信号処理方法および音信号処理装置を提供する。

【解決手段】音信号処理方法は、話者の音声に係る音信号を入力し、話者画像を取得し、前記話者画像から前記話者の姿勢情報を推定し、推定した前記姿勢情報に応じた補正フィルタを生成し、前記補正フィルタに係るフィルタ処理を前記音信号に施し、前記フィルタ処理を施した後の音信号を出力する。

【選択図】図10



【特許請求の範囲】**【請求項 1】**

話者の音声に係る音信号を入力し、
話者画像を取得し、
前記話者画像から前記話者の姿勢情報を推定し、
推定した前記姿勢情報に応じた補正フィルタを生成し、
前記補正フィルタに係るフィルタ処理を前記音信号に施し、
前記フィルタ処理を施した後の音信号を出力する、
音信号処理方法。

【請求項 2】

前記姿勢情報は、前記話者の顔の向きを含み、
前記補正フィルタは、前記顔の向きに応じて減衰するレベルを補償する処理を含む、
請求項 1 に記載の音信号処理方法。

10

【請求項 3】

前記補正フィルタは、イコライザを含む、
請求項 1 または請求項 2 に記載の音信号処理方法。

【請求項 4】

前記姿勢情報は、顔の左右の向きを示す情報を含み、
前記顔の左右の向きに応じて前記補正フィルタを生成する、
請求項 1 乃至請求項 3 のいずれか 1 項に記載の音信号処理方法。

20

【請求項 5】

前記補正フィルタは、前記顔の左右の向きが大きいほど高域のレベルを高くする、または
低域のレベルを低くする処理を含む、
請求項 4 に記載の音信号処理方法。

【請求項 6】

前記姿勢情報は、後ろ向きの姿勢の情報を含む、
請求項 1 乃至請求項 5 のいずれか 1 項に記載の音信号処理方法。

【請求項 7】

前記話者画像から前記話者の位置情報を推定し、
前記位置情報に基づいて前記補正フィルタを生成し、
前記位置情報の推定速度は、前記姿勢情報の推定速度よりも速く、
前記補正フィルタは、前記位置情報を推定した時、および前記姿勢情報を推定した時、の
それぞれのタイミングで生成される、
請求項 1 乃至請求項 6 のいずれか 1 項に記載の音信号処理方法。

30

【請求項 8】

話者の音声に係る音信号を入力する音信号入力部と、
話者画像を取得する画像取得部と、
前記話者画像から前記話者の姿勢情報を推定する位置推定部と、
推定した前記姿勢情報に応じた補正フィルタを生成するフィルタ生成部と、
前記補正フィルタに係るフィルタ処理を前記音信号に施す音信号処理部と、
前記フィルタ処理を施した後の音信号を出力する出力部と、
備えた音信号処理装置。

40

【請求項 9】

前記姿勢情報は、前記話者の顔の向きを含み、
前記補正フィルタは、前記顔の向きに応じて減衰するレベルを補償する処理を含む、
請求項 8 に記載の音信号処理装置。

【請求項 10】

前記補正フィルタは、イコライザを含む、
請求項 8 または請求項 9 に記載の音信号処理装置。

【請求項 11】

50

前記姿勢情報は、顔の左右の向きを示す情報を含み、
前記フィルタ生成部は、前記顔の左右の向きに応じて前記補正フィルタを生成する、
請求項 8 乃至請求項 10 のいずれか 1 項に記載の音信号処理装置。

【請求項 12】

前記補正フィルタは、前記顔の左右の向きが大きいほど高域のレベルを高くする、または低域のレベルを低くする処理を含む、
請求項 11 に記載の音信号処理装置。

【請求項 13】

前記姿勢情報は、後ろ向きの姿勢の情報を含む、
請求項 8 乃至請求項 12 のいずれか 1 項に記載の音信号処理装置。

10

【請求項 14】

前記話者画像から前記話者の位置情報を推定する位置推定部を備え、
前記フィルタ生成部は、前記位置情報に基づいて前記補正フィルタを生成し、
前記位置情報の推定速度は、前記姿勢情報の推定速度よりも速く、
前記補正フィルタは、前記位置情報を推定した時、および前記姿勢情報を推定した時、のそれぞれのタイミングで生成される、
請求項 8 乃至請求項 13 のいずれか 1 項に記載の音信号処理装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の一実施形態は、音源の位置に基づいてマイクで取得した音信号を処理する音信号処理方法および音信号処理装置に関する。

20

【背景技術】

【0002】

特許文献 1 には、カメラで撮影した映像から話者の位置情報を検出し、検出した位置情報に基づいて、話者の音声が増強されるような処理を行なう音処理システムが開示されている。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】特開 2012 - 29209 号公報

30

【発明の概要】

【発明が解決しようとする課題】

【0004】

話者の音声は、話者の姿勢に応じて変化する。しかし、特許文献 1 の音処理システムは、話者の姿勢を考慮していない。

【0005】

そこで、本発明の一実施形態の目的は、話者の姿勢に応じて適切に話者の音声を取得できる音信号処理方法および音信号処理装置を提供することにある。

【課題を解決するための手段】

40

【0006】

音信号処理方法は、音信号処理方法は、話者の音声に係る音信号を入力し、話者画像を取得し、前記話者画像から前記話者の姿勢情報を推定し、推定した前記姿勢情報に応じた補正フィルタを生成し、前記補正フィルタに係るフィルタ処理を前記音信号に施し、前記フィルタ処理を施した後の音信号を出力する。

【発明の効果】

【0007】

本発明の一実施形態によれば、話者の姿勢に応じて適切に話者の音声を取得できる。

【図面の簡単な説明】

【0008】

50

【図 1】音信号処理装置の構成を示すブロック図である。

【図 2】音信号処理方法の動作を示すフローチャートである。

【図 3】音信号処理装置の機能的構成を示すブロック図である。

【図 4】カメラ 11 が撮影した画像の一例を示す図である。

【図 5】話者の位置情報の一例を示す図である。

【図 6】音信号処理部 51 の機能的構成を示すブロック図である。

【図 7】残響特性を取得する場合の音信号処理部 51 の機能的構成を示すブロック図である。

【図 8】机 T の認識結果に応じて補正フィルタを生成する場合の例を示す図である。

【図 9】姿勢情報に基づいて補正フィルタを生成する場合の、音信号処理方法の動作を示すフローチャートである。 10

【図 10】音信号処理装置の機能的構成を示すブロック図である。

【図 11】姿勢情報の一例を示す図である。

【図 12】音信号処理部 51 の機能的構成を示すブロック図である。

【図 13】残響特性を取得する場合の音信号処理部 51 の機能的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】

(第 1 実施形態)

図 1 は、音信号処理装置 1 の構成を示すブロック図である。図 2 は、音信号処理方法の動作を示すフローチャートである。 20

【0010】

音信号処理装置 1 は、カメラ 11、CPU 12、DSP 13、フラッシュメモリ 14、RAM 15、ユーザインタフェース (I/F) 16、スピーカ 17、6 個のマイク 18A ~ 18F、および通信部 19 を備えている。なお、本実施形態において、信号とはデジタル信号を意味する。

【0011】

カメラ 11、スピーカ 17、およびマイク 18A ~ 18F は、例えば表示器 (不図示) の上または下に配置される。カメラ 11 は、表示器 (不図示) の前に居る利用者の画像を取得する。マイク 18A ~ 18F は、表示器 (不図示) の前に居る利用者の音声を取得する。スピーカ 17 は、表示器 (不図示) の前に居る利用者に対して、音声を出力する。なお、マイクの数 は 6 個に限らない。マイクは、1 つのマイクであってもよい。本実施形態のマイクの数 は 6 個であり、アレイマイクを構成する。DSP 13 は、マイク 18A ~ 18F で取得した音信号にビームフォーミング処理を施す。 30

【0012】

CPU 12 は、フラッシュメモリ 14 から動作プログラムのプログラムを RAM 15 に読み出すことにより、音信号処理装置 1 の動作を統括的に制御する制御部として機能する。なお、プログラムは自装置のフラッシュメモリ 14 に記憶しておく必要はない。CPU 12 は、例えばサーバ等から都度ダウンロードして RAM 15 に読み出してもよい。

【0013】

DSP 13 は、CPU 12 の制御に従って、映像信号および音信号をそれぞれ処理する信号処理部である。DSP 13 は、例えば映像信号から話者の画像を切り出すフレーミング処理を行なう画像処理部として機能する。また、DSP 13 は、例えば話者の音声の減衰を補償するための補正フィルタ処理を行うフィルタ処理部としても機能する。 40

【0014】

通信部 19 は、DSP 13 により処理された後の映像信号および音信号を、他の装置に送信する。また、通信部 19 は、他の装置から映像信号および音信号を受信する。通信部 19 は、受信した映像信号を表示器 (不図示) に出力する。通信部 19 は、受信した音信号をスピーカ 17 に出力する。表示器は、他の装置のカメラで取得した映像を表示する。スピーカ 17 は、他の装置のマイクで取得した話者の音声出力する。他の装置は、例えば 50

遠隔地に設置された音信号処理装置である。これにより、音信号処理装置 1 は、遠隔地との音声会話を行うためのコミュニケーションシステムとして機能する。

【 0 0 1 5 】

図 3 は、音信号処理装置 1 の機能的ブロック図である。これら機能的構成は、CPU 1 2 および DSP 1 3 により実現される。図 3 に示す様に、音信号処理装置 1 は、機能的に、音信号入力部 5 0、音信号処理部 5 1、出力部 5 2、画像取得部 1 0 0、位置推定部 1 0 1、およびフィルタ生成部 1 0 2 を備えている。

【 0 0 1 6 】

音信号入力部 5 0 は、マイク 1 8 A ~ 1 8 F から音信号を入力する (S 1 1)。また、画像取得部 1 0 0 は、カメラ 1 1 から話者画像を含む画像を取得する (S 1 2)。位置推定部 1 0 1 は、取得した話者画像から話者の位置情報を推定する (S 1 3)。

10

【 0 0 1 7 】

位置情報の推定は、人物の顔認識処理を含む。人物の顔認識処理は、例えばニューラルネットワーク等の所定のアルゴリズムにより、カメラ 1 1 が撮影した画像から複数の人物の顔の位置を認識する処理である。以下、本実施形態において話者とは、会議に参加しかつ現在会話している人を意味し、利用者とは会議に参加している人を意味し、話者を含む。非利用者とは、会議に参加していない人を意味し、人物とは、カメラ 1 1 に映る全ての人を意味する。

【 0 0 1 8 】

図 4 は、カメラ 1 1 が撮影した画像の一例を示す図である。図 4 の例では、カメラ 1 1 は、机 T の長手方向 (奥行き方向) に沿って居る複数の人物の顔画像を撮影している。

20

【 0 0 1 9 】

机 T は、平面視して長形状である。カメラ 1 1 は、机 T を短手方向に挟んで左側および右側に居る 4 人の利用者、および机 T よりも遠い位置に居る非利用者を撮影している。

【 0 0 2 0 】

位置推定部 1 0 1 は、この様なカメラ 1 1 の撮影した画像から人物の顔を認識する。図 4 の例では、画像の左下に居る利用者 A 1 が発話している。位置推定部 1 0 1 は、複数フレームの画像に基づいて、発話中の利用者 A 1 の顔を、話者の顔として認識する。なお、他の人物 A 2 ~ A 5 は、顔認識されているが、話者ではない。したがって、位置推定部 1 0 1 は、利用者 A 1 の顔を、話者の顔として認識する。

30

【 0 0 2 1 】

位置推定部 1 0 1 は、認識した話者の顔の位置に図中の四角で示す様な境界ボックス (B o u n d i n g B o x) を設定する。位置推定部 1 0 1 は、境界ボックスの大きさに基づいて話者との距離を求める。フラッシュメモリ 1 4 には、予め境界ボックスの大きさと距離との関係を示したテーブルまたは関数等が記憶されている。位置推定部 1 0 1 は、設定した境界ボックスの大きさと、フラッシュメモリ 1 4 に記憶されているテーブルを比較し、話者との距離を求める。

【 0 0 2 2 】

位置推定部 1 0 1 は、設定した境界ボックスの 2 次元座標 (X , Y 座標) および話者との距離を、話者の位置情報として求める。図 5 は、話者の位置情報の一例を示す図である。話者の位置情報は、話者を示すラベル名、2 次元座標、および距離を含む。2 次元座標は、カメラ 1 1 の撮影した画像の所定位置 (例えば左下) を原点とした X , Y 座標 (直交座標) である。距離は、例えばメートル等で示す値である。位置推定部 1 0 1 は、フィルタ生成部 1 0 2 に、話者の位置情報を出力する。なお、位置推定部 1 0 1 は、複数の話者の顔を認識した場合、複数の話者の位置情報を出力する。

40

【 0 0 2 3 】

なお、位置推定部 1 0 1 は、カメラ 1 1 で撮影した画像だけでなく、さらにマイク 1 8 A ~ 1 8 F で取得した音信号に基づいて人物の位置情報を推定してもよい。この場合、位置推定部 1 0 1 は、マイク 1 8 A ~ 1 8 F で取得した音信号を音信号入力部 5 0 から入力する。例えば、位置推定部 1 0 1 は、複数のマイクで取得した音信号の相互相関を求めるこ

50

とにより、人物の音声マイクに到達したタイミングを求めることができる。位置推定部101は、各マイクの位置関係および音声の到達タイミングに基づいて、人物の音声の到来方向を求めることができる。この場合、位置推定部101は、カメラ11の撮影した画像から、顔認識を行なうだけでもよい。例えば図4の例では、位置推定部101は、机Tを短手方向に挟んで左側および右側に居る4人の利用者、および机Tよりも遠い位置に居る非利用者の顔画像を認識する。そして、位置推定部101は、これらの顔画像から、話者の音声の到来方向に一致する顔画像を話者の位置情報として推定する。

【0024】

また、位置推定部101は、カメラ11の撮影した画像から人物の身体を推定し、人物の位置情報を推定してもよい。位置推定部101は、ニューラルネットワーク等の所定のアルゴリズムにより、カメラ11の撮影した画像から人の骨格(ボーン)を求める。ボーンは、目、鼻、首、肩、および手足等を含む。フラッシュメモリ14には、予めボーンの大きさと距離との関係を示したテーブルまたは関数等が記憶されている。位置推定部101は、認識したボーンの大きさと、フラッシュメモリ14に記憶されているテーブルを比較し、人物との距離を求める。

10

【0025】

次に、フィルタ生成部102は、話者の位置情報に応じて、補正フィルタを生成する(S14)。補正フィルタは、音声の減衰を補償するためのフィルタを含む。補正フィルタは、例えばゲイン補正、イコライザ、およびビームフォーミングを含む。話者の音声は、遠い距離ほど減衰する。また、話者の音声の高域成分は、話者の音声の低域成分に比べて、遠い距離ほど減衰する。したがって、フィルタ生成部102は、話者の位置情報のうち距離の値が大きいほど音信号のレベルを高くする様なゲイン補正フィルタを生成する。また、フィルタ生成部102は、話者の位置情報のうち距離の値が大きいほど高域のレベルを高くする様なイコライザのフィルタを生成してもよい。また、フィルタ生成部102は、話者の座標に指向性を向けるビームフォーミング処理を行なう補正フィルタを生成してもよい。

20

【0026】

音信号処理部51は、フィルタ生成部102で生成された補正フィルタに係るフィルタ処理を音信号に施す(S15)。出力部52は、フィルタ処理後の音信号を通信部19に出力する(S16)。音信号処理部51は、例えばデジタルフィルタからなる。音信号処理部51は、音信号を周波数軸上の信号に変換して、各周波数の信号のレベルを変更することにより、各種のフィルタ処理を行なう。

30

【0027】

図6は、音信号処理部51の機能的構成を示すブロック図である。音信号処理部51は、ビームフォーミング処理部501、ゲイン補正部502、およびイコライザ503を構成する。ビームフォーミング処理部501は、マイク18A~18Fで取得した音信号に、それぞれフィルタ処理を施して合成することによりビームフォーミングを行う。ビームフォーミングに係る信号処理は、遅延和(Delay Sum)方式、Griffiths Jim型、Sidelobe Canceller型、あるいはFrost型Adaptive Beamformer等、どのような手法であってもよい。

40

【0028】

ゲイン補正部502は、ビームフォーミング処理後の音信号のゲインを補正する。イコライザ503は、ゲイン補正後の音信号の周波数特性を調整する。ビームフォーミング処理のフィルタ、ゲイン補正部502のフィルタ、およびイコライザ503のフィルタは、全て補正フィルタに対応する。フィルタ生成部102は、話者の位置情報に応じて、補正フィルタを生成する。

【0029】

フィルタ生成部102は、話者の位置に向けて指向性を形成する様なフィルタ係数を生成し、ビームフォーミング処理部501に設定する。これにより、音信号処理装置1は、話者の音声を高い精度で取得することができる。

50

【 0 0 3 0 】

また、フィルタ生成部 1 0 2 は、話者の位置情報に基づいて、ゲイン補正部 5 0 2 のゲインを設定する。上述した様に、話者の音声は、遠い距離ほど減衰する。したがって、フィルタ生成部 1 0 2 は、話者の位置情報のうち距離の値が大きいほど音信号のレベルを高くする様なゲイン補正フィルタを生成し、ゲイン補正部 5 0 2 に設定する。これにより、音信号処理装置 1 は、話者との距離に関わらず、安定したレベルで話者の音声を取得することができる。

【 0 0 3 1 】

また、フィルタ生成部 1 0 2 は、話者の位置情報に基づいて、イコライザ 5 0 3 の周波数特性を設定する。上述した様に、フィルタ生成部 1 0 2 は、話者の位置情報のうち距離の値が大きいほど高域のレベルを高くする様なイコライザのフィルタを生成する。これにより、音信号処理装置 1 は、話者との距離に関わらず、安定した音質で話者の音声を取得することができる。

【 0 0 3 2 】

また、フィルタ生成部 1 0 2 は、ビームフォーミング処理部 5 0 1 から音声の到来方向の情報を取得してもよい。上述の様に、音声の到来方向は、複数のマイクの音信号に基づいて求めることができる。フィルタ生成部 1 0 2 は、人物の位置情報と、音声の到来方向の情報と、を対比して、ゲイン補正部 5 0 2 のゲインを設定してもよい。例えば、フィルタ生成部 1 0 2 は、話者の位置情報の示す話者の位置と、音声の到来方向との差（離角）が大きくなるほどゲインの値を小さく設定する。つまり、フィルタ生成部 1 0 2 は、離角に反比例するようなゲインを設定する。あるいは、フィルタ生成部 1 0 2 は、離角に応じて指数的にゲインが小さくなるような設定を行なってもよい。あるいは、フィルタ生成部 1 0 2 は、離角が所定の閾値以上となった場合にゲインが 0 になるような設定を行なってもよい。これにより、音信号処理装置 1 は、話者の音声をさらに高い精度で取得することができる。

【 0 0 3 3 】

また、フィルタ生成部 1 0 2 は、室内の残響特性を取得し、取得した残響特性に応じて補正フィルタを生成してもよい。図 7 は、残響特性を取得する場合の音信号処理部 5 1 の機能的構成を示すブロック図である。図 7 に示す音信号処理部 5 1 は、さらに適応エコーキャンセラ（A E C）7 0 1 を備えている。

【 0 0 3 4 】

A E C 7 0 1 は、スピーカ 1 7 から出力された音のうちマイク 1 8 A ~ 1 8 F に帰還する成分（エコー成分）を推定し、推定したエコー成分をキャンセルする。エコー成分は、スピーカ 1 7 に出力する信号に適応フィルタ処理を施すことで生成する。適応フィルタは、所定の適応アルゴリズムにより、室内の残響特性を模擬した F I R フィルタを構成する。適応フィルタは、当該 F I R フィルタでスピーカ 1 7 に出力する信号をフィルタ処理することによりエコー成分を生成する。

【 0 0 3 5 】

フィルタ生成部 1 0 2 は、A E C 7 0 1 の適応フィルタで模擬された残響特性（残響情報）を取得する。フィルタ生成部 1 0 2 は、取得した残響情報に応じて補正フィルタを生成する。例えば、フィルタ生成部 1 0 2 は、残響特性のパワーを求める。フィルタ生成部 1 0 2 は、残響特性のパワーに応じてゲイン補正部 5 0 2 のゲインを設定する。上述した様に、フィルタ生成部 1 0 2 は、離角に応じて指数的にゲインが小さくなるような設定を行なってもよい。また、フィルタ生成部 1 0 2 は、残響特性のパワーがより大きくなるほど減衰指数をよりゆっくり減衰するように設定してもよい。これらの場合、フィルタ生成部 1 0 2 は、残響特性のパワーが大きくなるほど閾値を大きく設定する。当該閾値を大きくすると、ビームフォーミング処理部 5 0 1 で生成されるビームの指向性が鈍化する。すなわち、フィルタ生成部 1 0 2 は、残響成分が大きい場合には、指向性を鈍化させる。残響成分が大きい場合、実際の話者の方向以外からも音声到来するため、到来方向の推定精度が低下する。つまり、推定した到来方向以外に人物が存在する可能性があり、上記離角

10

20

30

40

50

の値が大きくなる場合がある。したがって、フィルタ生成部 102 は、残響成分が大きい場合には指向性を鈍化させて、話者音声を取得できないことを防止する。

【0036】

なお、フィルタ生成部 102 は、人物の位置情報に加えて、さらに、フレーミング処理の結果を補正フィルタに反映してもよい。利用者 A1 は、ユーザ I/F 16 を用いてカメラ 11 の撮影した画像の中から特定の領域を切り出す操作を行なう。DSP 13 は、指定された領域を切り出すフレーミング処理を行なう。フィルタ生成部 102 は、切り出した領域の境界角度と、音声の到来方向に応じてゲイン補正部 502 のゲインを設定する。フィルタ生成部 102 は、音声の到来方向が、切り出した領域の境界角度を超えて、切り出した領域から出た場合にゲインを 0 にする。あるいは、フィルタ生成部 102 は、音声の到来方向が、切り出した領域の境界角度を超えて、切り出した領域から出た場合に、切り出した領域の境界角度を大きく超えれば超えるほどより 0 に近づくようなゲインを与えてもよい。また、境界角度は、左右両方に設けてもよいし、左右上下 4 方向に設けてもよい。これにより、音信号処理装置 1 は、利用者の指定した領域の話者の音声を高い精度で取得することができる。

10

【0037】

また、フィルタ生成部 102 は、特定のオブジェクトの認識結果に応じて補正フィルタを生成してもよい。例えば、位置推定部 101 は、特定のオブジェクトとして机 T を認識してもよい。図 8 は、机 T の認識結果に応じて補正フィルタを生成する場合の例を示す図である。位置推定部 101 は、ニューラルネットワーク等の所定のアルゴリズムにより、机 T を特定のオブジェクトとして認識する。位置推定部 101 は、机 T の位置情報をフィルタ生成部 102 に出力する。

20

【0038】

フィルタ生成部 102 は、机 T の位置情報に応じて補正フィルタを生成する。例えば、図 8 に示す様に、机 T の位置より上で、かつ机 T を短手方向に挟んで左側および右側の領域 S1 および領域 S2 に向けて指向性を形成する様なフィルタ係数を生成し、ビームフォーミング処理部 501 に設定する。あるいは、フィルタ生成部 102 は、領域 S1 および領域 S2 の位置と、音声の到来方向との差（離角）に応じてゲイン補正部 502 のゲインを設定してもよい。フィルタ生成部 102 は、離角が大きくなるほどゲインの値を小さく設定する。あるいは、フィルタ生成部 102 は、離角に応じて指数的にゲインが小さくなるような設定を行なってもよい。あるいは、フィルタ生成部 102 は、離角が所定の閾値以上となった場合にゲインが 0 になるような設定を行なってもよい。あるいは、フィルタ生成部 102 は、人物の位置が領域 S1 および領域 S2 の内部に存在するか外部に存在するかを判定して、人物の位置が外部に存在する場合にゲインが 0 になるようにゲイン補正部 502 のゲインを設定してもよい。

30

【0039】

これにより、音信号処理装置 1 は、機の位置より上でかつ机 T を短手方向に挟んで左側および右側の領域 S1 および領域 S2 の音声を高い精度で取得することができる。例えば、図 8 の例であれば、音信号処理装置 1 は、利用者 A3 の音声を取得せず、利用者 A1, A2, A4, A5 の音声のみ取得することができる。

40

【0040】

また、フィルタ生成部 102 は、人物と机との距離が所定値以上である場合に、対応する人物の音声をカットする補正フィルタを生成してもよい。例えば、図 8 の例で、利用者 A3 が発話した場合、位置推定部 101 は、利用者 A3 の位置を話者の位置情報として推定する。しかし、フィルタ生成部 102 は、人物との距離が所定値以上であるとして、利用者 A3 の音声をカットする補正フィルタを生成する。

【0041】

なお、所定値は、特定のオブジェクトの認識結果に基づいて求めてもよい。例えば図 8 の例では、フィルタ生成部 102 は、机 T よりも遠い位置の音声をカットする補正フィルタを生成する。

50

【 0 0 4 2 】

(第2実施形態)

次に、図9は、姿勢情報に基づいて補正フィルタを生成する場合の、音信号処理方法の動作を示すフローチャートである。図10は、姿勢情報に基づいて補正フィルタを生成する場合の、音信号処理装置1の機能的構成を示すブロック図である。この例の音信号処理装置1は、位置推定部101に代えて、姿勢推定部201を備える。ハードウェア構成は、図1に示した構成と同一である。

【 0 0 4 3 】

図9の例では、位置推定部101の位置推定処理(S13)に代えて、姿勢推定部201は、取得した話者画像から話者の姿勢情報を推定する(S23)。その他の処理は、図2に示したフローチャートと同様である。 10

【 0 0 4 4 】

姿勢情報の推定は、話者の顔認識処理を含む。話者の顔認識処理は、位置情報の推定と同様であり、例えばニューラルネットワーク等の所定のアルゴリズムにより、カメラ11が撮影した画像から話者の顔の位置を認識する処理である。姿勢推定部201は、カメラ11の撮影した画像から話者の顔を認識する。また、姿勢推定部201は、認識した顔のうち目の位置、口の位置、および鼻の位置等から、話者の向いている方向を推定する。例えば、フラッシュメモリ14には、顔に対する目の位置、口の位置、および鼻の位置のずれ(オフセット)と、姿勢情報とを対応付けたテーブルまたは関数等を記憶している。姿勢推定部201は、顔に対する目の位置、口の位置、および鼻の位置のオフセットと、フラッシュメモリ14に記憶されているテーブルとを比較し、話者の姿勢を求める。なお、姿勢推定部201は、顔の位置を認識しても目、口、および鼻を認識できない場合、後ろ向きの姿勢であると推定する。 20

【 0 0 4 5 】

図11は、姿勢情報の一例を示す図である。話者の姿勢は、顔の左右の向き(角度)を示す情報である。例えば、姿勢推定部201は、利用者A1の姿勢を15度と認識している。この例では、姿勢推定部201は、向かって正面に向いている場合を0度とし、向かって右側に向いている場合を正の角度、向かって左側に向いている場合を負の角度、真後ろを向いている場合を180度(または-180度)と認識する。

【 0 0 4 6 】

なお、姿勢推定部201は、カメラ11の撮影した画像から話者の身体を推定し、姿勢情報を推定してもよい。姿勢推定部201は、ニューラルネットワーク等の所定のアルゴリズムにより、カメラ11の撮影した画像から鼻のボーンと、身体(首、肩、および手足等)のボーンを認識する。フラッシュメモリ14には、予め鼻のボーンと、身体のボーンのずれ(オフセット)と、姿勢情報とを対応付けたテーブルまたは関数等を記憶している。姿勢推定部201は、身体のボーンに対する鼻のボーンのオフセットと、フラッシュメモリ14に記憶されているテーブルとを比較し、話者の姿勢を求めてもよい。 30

【 0 0 4 7 】

フィルタ生成部102は、姿勢情報に応じて、補正フィルタを生成する。補正フィルタは、顔の向きに応じて減衰するレベルを補償するためのフィルタを含む。補正フィルタは、例えばゲイン補正、イコライザ、およびビームフォーミングを含む。 40

【 0 0 4 8 】

図12は、音信号処理部51の機能的構成を示すブロック図である。図12に示すブロック図は、フィルタ生成部102が姿勢情報を入力する点以外は、図6に示したブロック図と同一の構成である。

【 0 0 4 9 】

話者の音声は、真正面を向いている場合に最も高いレベルを示し、左右の向きが大きくなるほど減衰する。また、左右の向きが大きくなるほど高域が低域に比べてより減衰する。したがって、フィルタ生成部102は、左右の向き(角度)が大きいかほど音信号のレベルを高くする様なゲイン補正フィルタを生成し、ゲイン補正部502に設定する。また、フ 50

フィルタ生成部 102 は、左右の向き（角度）が大きいほど高域のレベルを高くする、あるいは低域のレベルを低くする様なイコライザのフィルタを生成し、イコライザ 503 に設定してもよい。

【0050】

これにより、音信号処理装置 1 は、話者の姿勢に関わらず、安定したレベル、安定した音質で話者の音声を取得することができる。

【0051】

また、フィルタ生成部 102 は、姿勢情報に基づいてビームフォーミング処理部 501 の指向性を制御してもよい。残響成分は、話者が真正面を向いている場合に最も低いレベルを示し、左右の向きが大きくなるほど大きくなる。したがって、フィルタ生成部 102 は、左右の向き（角度）が大きい場合には、残響成分が大きいと判断して、指向性を鈍化させてもよい。これにより、音信号処理装置 1 は、話者の音声を高い精度で取得することができる。

10

【0052】

また、図 13 に示す様に、フィルタ生成部 102 は、残響情報を取得してもよい。図 13 の構成は、図 7 の例と同様である。フィルタ生成部 102 は、AEC701 から残響情報を取得する。フィルタ生成部 102 は、取得した残響情報に応じて補正フィルタを生成する。例えば、フィルタ生成部 102 は、残響特性のパワーを求める。フィルタ生成部 102 は、残響特性のパワーに応じてゲイン補正部 502 のゲインを設定してもよい。

【0053】

第 1 実施形態の音信号処理装置 1 は、位置情報に基づいて補正フィルタを生成する例を示し、第 2 実施形態の音信号処理装置 1 は、姿勢情報に基づいて補正フィルタを生成した。無論、音信号処理装置 1 は、位置情報および姿勢情報の両方に基づいて補正フィルタを生成してもよい。ただし、位置情報の推定速度と、姿勢情報の推定速度は、異なる場合がある。第 1 実施形態の音信号処理装置 1 における位置情報の推定速度は、第 2 実施形態の姿勢情報の推定速度よりも速い。この場合、フィルタ生成部 102 は、位置推定部 101 が位置情報を推定した時、および姿勢推定部 201 が姿勢情報を推定した時、のそれぞれのタイミングで補正フィルタを生成すればよい。

20

【0054】

第 1 実施形態および第 2 実施形態の説明は、すべての点で例示であって、制限的なものではないと考えられるべきである。本発明の範囲は、上述の実施形態ではなく、特許請求の範囲によって示される。さらに、本発明の範囲は、特許請求の範囲と均等の範囲を含む。

30

【符号の説明】

【0055】

- 1 ... 音信号処理装置
- 11 ... カメラ
- 12 ... CPU
- 13 ... DSP
- 14 ... フラッシュメモリ
- 15 ... RAM
- 16 ... ユーザ I/F
- 17 ... スピーカ
- 18A ~ 18F ... マイク
- 19 ... 通信部
- 50 ... 音信号入力部
- 51 ... 音信号処理部
- 52 ... 出力部
- 100 ... 画像取得部
- 101 ... 位置推定部
- 102 ... フィルタ生成部

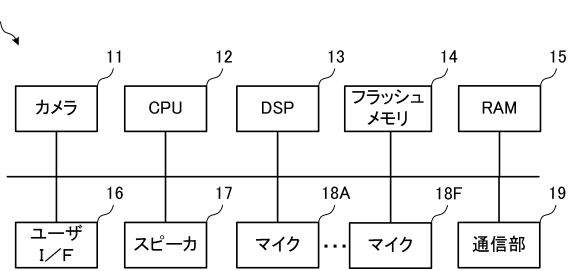
40

50

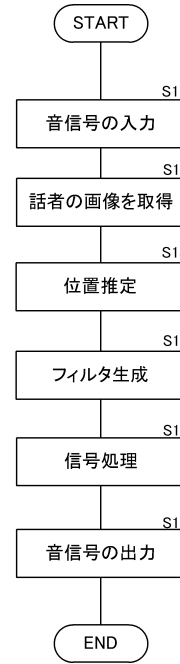
- 2 0 1 ... 姿勢推定部
- 5 0 1 ... ビームフォーミング処理部
- 5 0 2 ... ゲイン補正部
- 5 0 3 ... イコライザ
- 7 0 1 ... A E C

【図面】

【図 1】



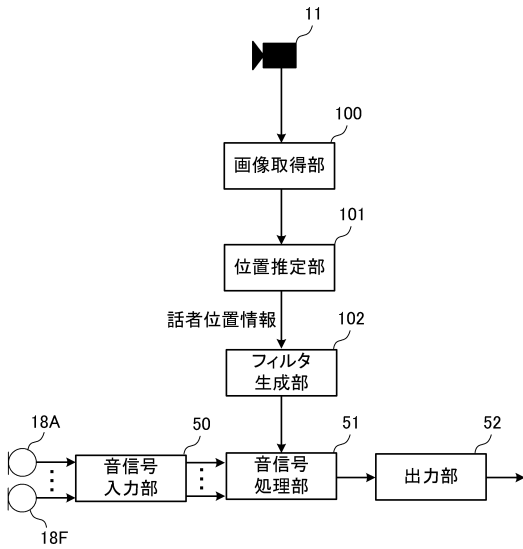
【図 2】



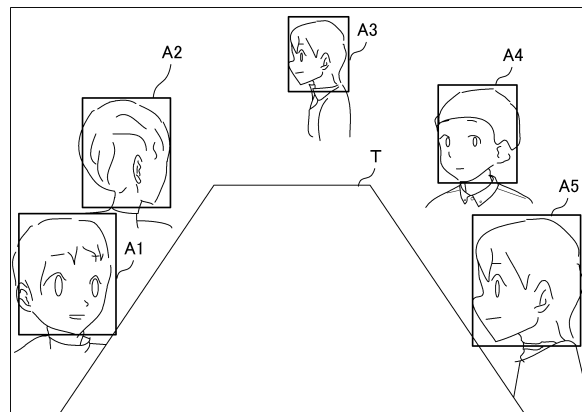
10

20

【図 3】



【図 4】



30

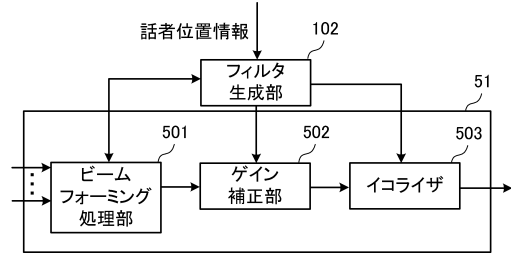
40

50

【図5】

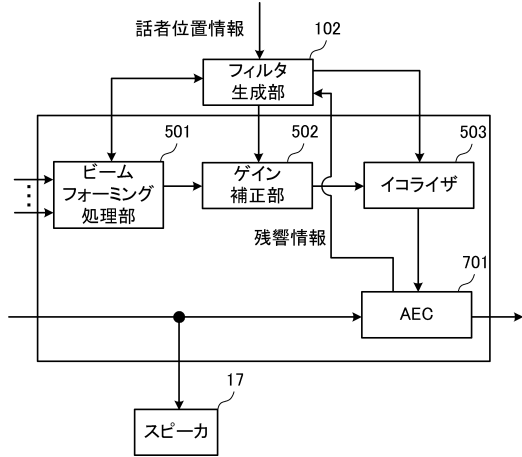
話者	X,Y座標	距離
A1	0.152,0.155	0.25

【図6】

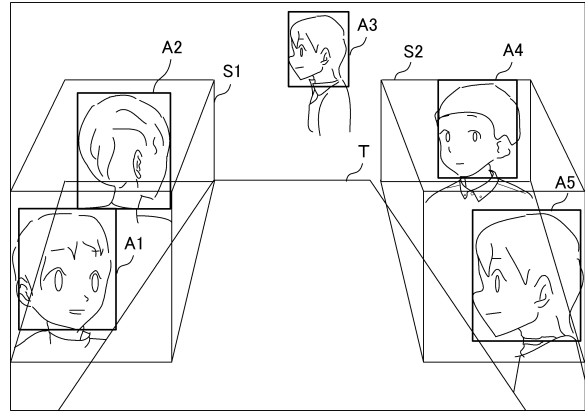


10

【図7】

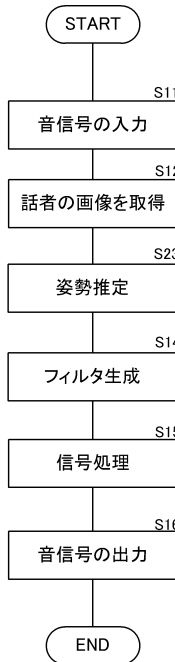


【図8】

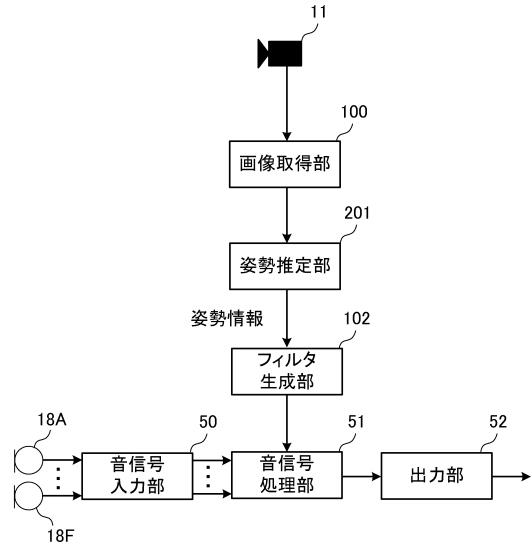


20

【図9】



【図10】



30

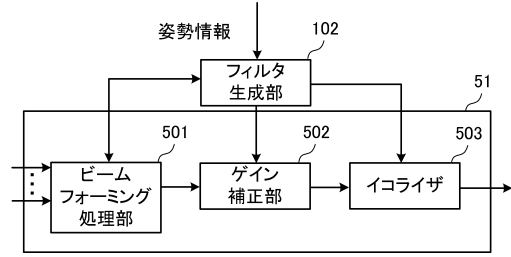
40

50

【図 1 1】

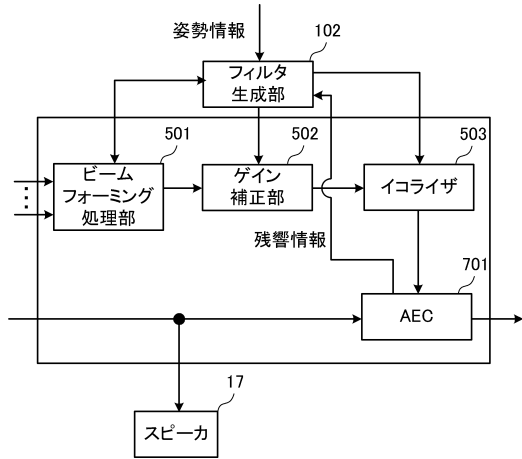
話者	顔の左右角度
A1	15

【図 1 2】



10

【図 1 3】



20

30

40

50