

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6464449号  
(P6464449)

(45) 発行日 平成31年2月6日(2019.2.6)

(24) 登録日 平成31年1月18日(2019.1.18)

(51) Int.Cl.	F I
<b>G 1 O L 15/28 (2013.01)</b>	G 1 O L 15/28 4 0 0
<b>G 1 O L 15/20 (2006.01)</b>	G 1 O L 15/20 3 7 O E
<b>G 1 O L 21/0272 (2013.01)</b>	G 1 O L 21/0272 1 0 0 A
<b>G O 6 T 7/00 (2017.01)</b>	G O 6 T 7/00 Q
<b>G 1 O L 25/51 (2013.01)</b>	G 1 O L 25/51 4 0 0
請求項の数 10 (全 27 頁) 最終頁に続く	

(21) 出願番号	特願2014-176518 (P2014-176518)	(73) 特許権者	000005326
(22) 出願日	平成26年8月29日 (2014. 8. 29)		本田技研工業株式会社
(65) 公開番号	特開2016-51081 (P2016-51081A)		東京都港区南青山二丁目1番1号
(43) 公開日	平成28年4月11日 (2016. 4. 11)	(74) 代理人	100165179
審査請求日	平成28年11月29日 (2016. 11. 29)		弁理士 田▲崎▼ 聡
		(74) 代理人	100126664
			弁理士 鈴木 慎吾
		(74) 代理人	100154852
			弁理士 酒井 太一
		(74) 代理人	100194087
			弁理士 渡辺 伸一
		(74) 代理人	100064908
			弁理士 志賀 正武
		(74) 代理人	100146835
			弁理士 佐伯 義文
最終頁に続く			

(54) 【発明の名称】 音源分離装置、及び音源分離方法

(57) 【特許請求の範囲】

【請求項1】

音響信号を収録する収録部と、  
 画像を撮像する撮像部と、  
 前記収録された音響信号を評価する音響信号評価部と、  
 前記撮像された画像信号を評価する画像信号評価部と、  
 前記音響信号評価部と前記画像信号評価部とによって評価された結果に基づいて、音源方向推定部によって前記収録された音響信号に基づいて音源方向の推定を行うか、人位置推定部によって前記撮像された画像に対する音源方向の推定を行うか、を判定する判定部と、  
 前記判定部が判定した結果に基づいて、前記収録された音響信号に基づいて推定された音源方向を示す情報を用いて、前記画像から音源の方向を示す情報を推定する人位置推定部と、  
 前記判定部が判定した結果に基づいて、前記撮像された画像に基づいて推定された音源方向を示す情報を用いて、前記音響信号に対して音源の方向を推定する音源方向推定部と、  
 前記人位置推定部によって推定された前記音源の方向を示す情報、または、前記音源方向推定部によって推定された前記音源の方向を示す情報に基づいて、前記音源の方向に対応する音響信号を前記音響信号から抽出する音源分離部と、  
 を備える音源分離装置。

## 【請求項 2】

前記推定された音源方向に対応する領域の画像を、前記撮像された画像から抽出する画像抽出部と、

前記抽出された画像以外の領域の画像を変更し、前記変更した画像と前記抽出した画像とを合成する画像合成部と、

を備える請求項 1 に記載の音源分離装置。

## 【請求項 3】

前記画像合成部は、

前記抽出された画像以外の領域の画像の解像度を、前記抽出された画像の解像度より低くするように変更する請求項 2 に記載の音源分離装置。

10

## 【請求項 4】

前記画像信号評価部は、

前記撮像された画像のヒストグラムを算出し、算出した前記ヒストグラムにおいて、ピクセル数が所定の値以上の輝度の輝度範囲を算出し、算出した前記輝度範囲が所定の範囲以上の場合に画像の信頼性が高いと評価し、算出した前記輝度範囲が所定の範囲未満の場合に画像の信頼性が低いと評価する請求項 1 から請求項 3 のいずれか 1 項に記載の音源分離装置。

## 【請求項 5】

前記画像信号評価部は、

前記撮像された画像のヒストグラムを算出し、算出した前記ヒストグラムにおいて、ピクセル数が所定値以上ある輝度の数をカウントし、前記画像の総ピクセル数を前記カウントした値で除算して判定値を算出し、算出した判定値に基づいて画像の信頼性を評価する請求項 1 から請求項 3 のいずれか 1 項に記載の音源分離装置。

20

## 【請求項 6】

前記音響信号評価部は、

前記音響信号の雑音成分の大きさを、前記音響信号に対して雑音抑圧処理を行った結果に基づいて算出し、前記算出した雑音成分の大きさに基づいて前記音響信号の信頼性を評価する請求項 1 から請求項 5 のいずれか 1 項に記載の音源分離装置。

## 【請求項 7】

前記音響信号及び前記画像信号のうち、少なくとも一方の信号に基づいて、発話区間を検出する発話区間検出部と、

発話区間ごとに、前記推定された音源方向に対応する音響信号を、前記収録された音響信号から抽出する音源分離部と、

発話区間ごとに、前記画像信号から抽出された話者の顔を含む領域の画像と、抽出された音響信号とを関連付ける関連付け部と、

を備え、

前記音源分離部は、

発話区間ごとに、前記推定された音源方向に対応する音響信号を、前記収録された音響信号から抽出する請求項 1 から請求項 6 のいずれか 1 項に記載の音源分離装置。

30

## 【請求項 8】

発話区間ごとに、前記画像信号から抽出された話者の顔を含む領域の画像と、抽出された音響信号とが関連付けられた情報を送信する送信部、を備える請求項 7 に記載の音源分離装置。

40

## 【請求項 9】

前記音源分離部は、

前記音響信号を用いて音源の方向を推定し、または、前記人位置推定部によって推定された話者の方向を示す情報に基づいて音源の方向を推定し、前記音源の方向の推定結果を用いて、前記音響信号を音源毎に分離することで抽出し、

前記音源分離部によって分離した音源毎の音響信号の特徴量を算出する特徴量算出部と

50

前記音響信号の特徴量に基づいて発話内容を認識してテキスト情報に変換する音声認識部、を備える請求項 1 から請求項 8 のいずれか 1 項に記載の音源分離装置。

【請求項 10】

收音部が、音響信号を収録する收音手順と、

撮像部が、画像を撮像する撮像手順と、

音響信号評価部が、前記收音手順によって収録された音響信号を評価する音響信号評価手順と、

画像信号評価部が、前記撮像手順によって撮像された画像信号を評価する画像評価手順と、

選択部が、前記音響信号評価手順と前記画像評価手順とによって評価された結果に基づいて、音源方向推定手順によって前記収録された音響信号に基づいて音源方向の推定を行うか、人位置推定手順によって前記撮像された画像に対する音源方向の推定を行うかを判定する判定手順と、

人位置推定部が、前記判定手順が判定した結果に基づいて、前記收音手順によって収録された音響信号に基づいて推定された音源方向を示す情報を用いて、前記画像から音源の方向を示す情報を推定する人位置推定手順と、

音源方向推定部は、前記判定手順が判定した結果に基づいて、前記撮像手順によって撮像された画像に基づいて推定された音源方向を示す情報を用いて、前記音響信号に対して音源の方向を推定する音源方向推定手順と、

音源分離部が、前記人位置推定手順によって推定された前記音源の方向を示す情報、または、前記音源方向推定手順によって推定された前記音源の方向を示す情報に基づいて、前記音源の方向に対応する音響信号を前記音響信号から抽出する音源分離手順と、

を含む音源分離方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音源分離装置、及び音源分離方法に関する。

【背景技術】

【0002】

会議における出席者の音声を集音し、出席者の映像を撮像する装置が提案されている。例えば、テレビ会議では、会議を行う地点毎にテレビ会議装置を設置し、これらのテレビ会議装置をネットワークで接続し、各テレビ会議装置が収録した音声信号と撮像した映像信号とを通信している。

【0003】

例えば、特許文献 1 に記載のテレビ会議装置では、マイクアレイと、会議室全体を撮像する全体撮像手段と、各会議者をそれぞれ個別に撮像して各会議者に関連付けされた個別画像を生成する複数の特定会議者撮像手段と、マイクアレイの收音信号に基づいて話者方向を検出して話者方向データを生成する話者方向検出手段と、検出された話者方向に応じて話者音声信号を生成する話者音声信号生成手段と、全体画像、各個別画像、話者音声信号、および話者方向データを送信する送信手段と、を備えることが提案されている。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2007 - 274462 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、特許文献 1 に記載の技術では、音声信号を收音する環境の雑音が大きい場合、話者に対する音源定位の精度が低下する。また、特許文献 1 に記載の技術では、この音源定位させた話者方向データに対応する個別画像を選択していたので、音源定位の精

10

20

30

40

50

度が落ちた場合、正しい個別画像を選択できない場合があった。

【0006】

本発明は上記の点に鑑みてなされたものであり、雑音が多い環境下であっても音源の方向を推定する精度を向上することができる音源分離装置、及び音源分離方法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】

(1) 上記目的を達成するため、本発明の一態様に係る音源分離装置は、音響信号を収録する収録部と、画像を撮像する撮像部と、前記収録された音響信号を評価する音響信号評価部と、前記撮像された画像信号を評価する画像信号評価部と、前記音響信号評価部と前記画像信号評価部とによって評価された結果に基づいて、音源方向推定部によって前記収録された音響信号に基づいて音源方向の推定を行うか、人位置推定部によって前記撮像された画像に対する音源方向の推定を行うか、を判定する判定部と、前記判定部が判定した結果に基づいて、前記収録された音響信号に基づいて推定された音源方向を示す情報を用いて、前記画像から音源の方向を示す情報を推定する人位置推定部と、前記判定部が判定した結果に基づいて、前記撮像された画像に基づいて推定された音源方向を示す情報を用いて、前記音響信号に対して音源の方向を推定する音源方向推定部と、前記人位置推定部によって推定された前記音源の方向を示す情報、または、前記音源方向推定部によって推定された前記音源の方向を示す情報に基づいて、前記音源の方向に対応する音響信号を前記音響信号から抽出する音源分離部と、を備える。

10

20

【0008】

(2) また、本発明の一態様に係る音源分離装置であって、前記推定された音源方向に対応する領域の画像を、前記撮像された画像から抽出する画像抽出部と、前記抽出された画像以外の領域の画像を変更し、前記変更した画像と前記抽出した画像とを合成する画像合成部と、を備えるようにしてもよい。

【0009】

(3) また、本発明の一態様に係る音源分離装置であって、前記画像合成部は、前記抽出された画像以外の領域の画像の解像度を、前記抽出された画像の解像度より低くするように変更するようにしてもよい。

【0010】

(4) また、本発明の一態様に係る音源分離装置であって、前記画像信号評価部は、前記撮像された画像のヒストグラムを算出し、算出した前記ヒストグラムにおいて、ピクセル数が所定の値以上の輝度の輝度範囲を算出し、算出した前記輝度範囲が所定の範囲以上の場合に画像の信頼性が高いと評価し、算出した前記輝度範囲が所定の範囲未満の場合に画像の信頼性が低いと評価するようにしてもよい。

30

【0011】

(5) また、本発明の一態様に係る音源分離装置であって、前記画像信号評価部は、前記撮像された画像のヒストグラムを算出し、算出した前記ヒストグラムにおいて、ピクセル数が所定値以上ある輝度の数をカウントし、前記画像の総ピクセル数を前記カウントした値で除算して判定値を算出し、算出した判定値に基づいて画像の信頼性を評価するようにしてもよい。

40

【0012】

(6) また、本発明の一態様に係る音源分離装置であって、前記音響信号評価部は、前記音響信号の雑音成分の大きさを、前記音響信号に対して雑音抑圧処理を行った結果に基づいて算出し、前記算出した雑音成分の大きさに基づいて前記音響信号の信頼性を評価するようにしてもよい。

【0013】

(7) また、本発明の一態様に係る音源分離装置であって、前記音響信号及び前記画像信号のうち、少なくとも一方の信号に基づいて、発話区間を検出する発話区間検出部と、発話区間ごとに、前記推定された音源方向に対応する音響信号を、前記収録された音響信号

50

から抽出する音源分離部と、発話区間ごとに、前記画像信号から抽出された話者の顔を含む領域の画像と、抽出された音響信号とを関連付ける関連付け部と、を備え、前記音源分離部は、発話区間ごとに、前記推定された音源方向に対応する音響信号を、前記収録された音響信号から抽出するようにしてもよい。

【0014】

(8) また、本発明の一態様に係る音源分離装置であって、発話区間ごとに、前記画像信号から抽出された話者の顔を含む領域の画像と、抽出された音響信号とが関連付けられた情報を送信する送信部、を備えるようにしてもよい。

【0015】

(9) また、本発明の一態様に係る音源分離装置であって、前記音源分離部は、前記音響信号を用いて音源の方向を推定し、または、前記人位置推定部によって推定された話者の方向を示す情報に基づいて音源の方向を推定し、前記音源の方向の推定結果を用いて、前記音響信号を音源毎に分離することで抽出し、前記音源分離部によって分離した音源毎の音響信号の特徴量を算出する特徴量算出部と、前記音響信号の特徴量に基づいて発話内容を認識してテキスト情報に変換する音声認識部、を備えるようにしてもよい。

【0016】

(10) 上記目的を達成するため、本発明の一態様に係る音源分離方法は、收音部が、音響信号を収録する收音手順と、撮像部が、画像を撮像する撮像手順と、音響信号評価部が、前記收音手順によって収録された音響信号を評価する音響信号評価手順と、画像信号評価部が、前記撮像手順によって撮像された画像信号を評価する画像評価手順と、選択部が、前記音響信号評価手順と前記画像評価手順とによって評価された結果に基づいて、音源方向推定手順によって前記収録された音響信号に基づいて音源方向の推定を行うか、人位置推定手順によって前記撮像された画像に対する音源方向の推定を行うかを判定する判定手順と、人位置推定部が、前記判定手順が判定した結果に基づいて、前記收音手順によって収録された音響信号に基づいて推定された音源方向を示す情報を用いて、前記画像から音源の方向を示す情報を推定する人位置推定手順と、音源方向推定部は、前記判定手順が判定した結果に基づいて、前記撮像手順によって撮像された画像に基づいて推定された音源方向を示す情報を用いて、前記音響信号に対して音源の方向を推定する音源方向推定手順と、音源分離部が、前記人位置推定手順によって推定された前記音源の方向を示す情報、または、前記音源方向推定手順によって推定された前記音源の方向を示す情報に基づいて、前記音源の方向に対応する音響信号を前記音響信号から抽出する音源分離手順と、を含む。

【発明の効果】

【0017】

上述した(1)または(10)の構成によれば、撮像した画像と収録した音響信号の評価結果に応じて、少なくとも一方に基づいて音源の方向を推定することができる。このため、本構成によれば、画像信号の信頼性が低い場合には、音響信号を用いて音源の方向を推定ことができ、音響信号の信頼性が低い場合には、画像信号に基づいて音源の方向を推定することができる。そして、本構成によれば、推定した結果に基づいて、発話された音声信号を分離することができる。このように、本構成によれば、音響信号と画像信号とを相互補完して音源分離を行うことができる。

【0018】

上述した(2)の構成によれば、音源の方向の画像以外の領域の画像を変更し、音源の方向の画像と変更した画像とを合成するようにしたので、画像情報の容量を軽減することができる。

上述した(3)の構成によれば、音源の方向の画像以外の領域の画像の解像度を、音源の方向の画像の解像度より低く変更し、音源の方向の画像と変更した画像とを合成するようにしたので、画像情報の容量を軽減することができる。

【0019】

上述した(4)の構成によれば、画像信号をヒストグラムにおいて、ピクセル数が所定

10

20

30

40

50

の数以上ある輝度領域の幅に基づいて、画像信号の信頼性を評価するようにしたので、画像信号を定量的に且つ簡便に評価することができる。

上述した(5)の構成によれば、画像信号をヒストグラムにおいて、ピクセルがある輝度の数をカウントして、画像の総ピクセル数をカウントした値で除算した判定値に基づいて画像信号の信頼性を評価するようにしたので、画像信号を定量的に且つ簡便に評価することができる。

【0020】

上述した(6)の構成によれば、音響信号の雑音成分の大きさを、残響抑圧処理の結果に基づいて算出された値に基づいて音響信号の信頼性を評価するようにしたので、音響信号を定量的に且つ簡便に評価することができる。

10

【0021】

上述した(7)の構成によれば、画像信号または音響信号のうち、少なくとも一方の音源の方向の推定結果に基づいて、画像信号または音響信号から発話区間を抽出することができるので、画像信号または音響信号のいずれか一方の信頼性が低い場合であっても、精度良く発話区間を検出することができる。この結果、この構成によれば、発話区間毎に発話された音響信号を精度良く分離することができる。

【0022】

上述した(8)の構成によれば、発話区間毎に、画像信号または音響信号のうち、少なくとも一方の音源の方向の推定結果に基づいて分離された音響信号と抽出された話者の顔を含む画像とを関連付けて、他の装置へ送信することができる。この結果、本構成の音源分離装置をテレビ会議等に用いる場合、画像信号または音響信号のいずれか一方の信頼性が低い場合であっても、分離された音響信号と抽出された話者の顔を含む画像とを関連付けて、他の装置へ送信することができる。

20

【0023】

上述した(9)の構成によれば、画像信号または音響信号のうち、少なくとも一方の音源の方向の推定結果に基づいて分離された音響信号に対して音声認識を行いテキスト化することができる。この結果、本構成によれば、画像信号または音響信号のいずれか一方の信頼性が低い場合であっても、議事録におけるテキスト認識率の精度を向上させることができる。

【図面の簡単な説明】

30

【0024】

【図1】第1実施形態に係る音源分離装置の構成を示すブロック図である。

【図2】撮像部によって撮像された画像の一例を説明する図である。

【図3】図2に示した画像Ph1において顔認識を行った結果の一例を説明する図である。

【図4】図2及び図4の画像のヒストグラムを説明する図である。

【図5】第1実施形態に係る信頼性判定部による判定結果の一例を説明する図である。

【図6】関連付部によって関連付けられた議事録情報の一例を説明する図である。

【図7】暗すぎる画像Ph2の一例である。

【図8】明るすぎる画像Ph3の一例である。

40

【図9】図2、図7、及び図8の各画像のヒストグラムを説明する図である。

【図10】コントラストが低すぎる画像Ph21の一例である。

【図11】コントラストが高すぎる画像Ph22の一例である。

【図12】図10及び図11の画像のヒストグラムを説明する図である。

【図13】第1実施形態に係る画像信号評価部の処理手順のフローチャートである。

【図14】雑音が少ない音響信号の一例を示す図である。

【図15】雑音が多い音響信号の一例を示す図である。

【図16】第1実施形態に係る音響信号評価部の処理手順のフローチャートである。

【図17】第1実施形態に係る音源分離装置が行う処理手順の一例を説明するフローチャートである。

50

【図18】第1実施形態に係る人位置推定の推定結果を優先する処理手順のフローチャートである。

【図19】第1実施形態に係る音源定位の推定結果を優先する処理手順のフローチャートである。

【図20】第2実施形態に係る音源分離装置の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】

【0025】

まず、本発明の概要を説明する。

本実施形態の音源分離装置1は、例えば会議室に設置されている。音源分離装置1は、  
 10 收音部と撮像部とを備え、会議中の話者音響信号と画像信号とを取得する。音源分離装置1は、取得した音響信号を評価し、また取得した画像信号を評価する。ここで、音響信号の評価対象は、例えば雑音の大きさ、話者の発話に対応する音声信号に対する雑音信号の比(S/N比)等である。また、画像信号の評価対象は、画像における輝度、コントラスト等である。そして、音源分離装置1は、評価結果に応じて、音響信号に基づいて音源定位を行うか、画像信号に基づいて話者の位置を推定するかを決定する。さらに、音源分離装置1は、話者が発話している区間、発話している話者の顔を含む領域のみ解像度を落とさず、他の領域の解像度を落とす。また、音源分離装置1は、発話されている音声のみを抽出する。

以下、図面を参照しながら本発明の実施形態について説明する。

【0026】

<第1実施形態>

図1は、本実施形態に係る音源分離装置1の構成を示すブロック図である。図1に示すように、音源分離装置1は、撮像部10、画像信号処理部20、收音部30、音響信号処理部40、信頼性判定部50(選択部)を含んで構成される。また、音源分離装置1には、議事録作成部60が接続されている。なお、音源分離装置1は、議事録作成部60を含むようにしてもよい。

【0027】

撮像部10は、所定の間隔毎に画像を撮像し、撮像した画像信号を画像信号処理部20に送信する。撮像部10は、撮像された画像信号を無線で送信してもよいし、有線で送信してもよい。撮像部10が複数台の場合には、送信の際にチャンネル間で画像信号が同期していればよい。なお、画像は静止画であってもよく、動画であってもよい。また、撮像部10は、会議室全体を撮像できる位置に取り付けられていてもよい。または、撮像部10は、会議に参加している人が、例えば並列に配置されたテーブルに着席している場合、それぞれのテーブルに着席している参加者の少なくとも口元が撮像できる位置に複数台、取り付けられていてもよい。

図2は、撮像部10によって撮像された画像の一例を説明する図である。図2に示すように、画像Ph1は、画像全体が十分な輝度成分を有しているため、後述するように画像Ph1を画像認識させることで、話者数、各話者の位置、話者が発話しているか否かを推定することができる。なお、図2において、符号A1~A4に示す領域の画像は、話者Sp1~Sp4に対応する画像である。なお、以下において、話者Sp1~Sp4のうち、  
 40 特定しない場合は単に話者Spという。

【0028】

図1に戻って、音源分離装置1の構成の説明を続ける。

画像信号処理部20は、画像信号取得部21、画像信号評価部22、画像信号の事前情報生成部23、人位置推定部24、唇検出部25(人位置推定部)、発話区間検出部26(人位置推定部)、及び画像処理部27を含んで構成される。

【0029】

画像信号取得部21は、撮像部10が送信した画像信号を取得し、取得した画像信号をアナログ信号からデジタル信号に変換する。画像信号取得部21は、デジタル信号に変換した画像信号を、画像信号評価部22、画像信号の事前情報生成部23、人位置推定部2  
 50

4、及び画像処理部 27 に出力する。

【0030】

画像信号評価部 22 は、画像信号取得部 21 から入力された画像信号を評価する。例えば、画像信号評価部 22 は、入力された画像信号のヒストグラムを算出し、算出したヒストグラムに基づいて、画像信号の輝度が所定の値以上であるか否かを判別することで、映像信号を評価する。画像信号評価部 22 は、評価した評価結果を信頼性判定部 50 に出力する。なお、画像信号評価部 22 が行う評価方法については、後述する。

【0031】

画像信号の事前情報生成部 23 には、画像信号取得部 21 から画像信号が入力され、信頼性判定部 50 から判定結果が入力される。また、画像信号の事前情報生成部 23 には、音響信号処理部 40 の音源方向推定部 44 から音源方向の推定結果が入力され、発話区間検出部 48 から発話区間を示す情報が入力される。画像信号の事前情報生成部 23 は、判定結果に基づいて、画像信号の事前情報を生成するか否かを決定する。ここで、画像信号の事前情報とは、発話を行っている話者の位置を示す情報である。なお、音源方向推定部 44 から入力される情報は、世界座標系の座標に基づくものとする。このため、画像信号の事前情報生成部 23 は、画像信号に基づく画像の座標を、周知の座標変換技術を用いて、世界座標系の座標に変換する。画像信号の事前情報生成部 23 は、判定結果が人位置推定の推定結果を優先する場合、画像信号の事前情報の生成処理を行わず、事前情報を画像処理部 27 に出力しない。一方、画像信号の事前情報生成部 23 は、判定結果が音源定位の推定結果を優先する場合、音源方向の推定結果を用いて、発話区間毎に画像信号の事前情報の生成処理を行い、生成した画像信号の事前情報である話者の顔の領域を示す情報を画像処理部 27 に出力する。

【0032】

人位置推定部 24 には、画像信号取得部 21 から画像信号が入力され、信頼性判定部 50 から判定結果が入力される。人位置推定部 24 は、判定結果に基づいて、入力された画像信号から人の位置推定（以下、人位置推定ともいう）を行うか否かを決定する。人位置推定部 24 は、判定結果が人位置推定の推定結果を優先する場合、入力された画像信号を用いて、画像に写っている人毎の位置を人の顔の画像を周知の画像認識技術を用いて認識する。人の位置は、例えば顔を含む領域の位置である。なお、人位置推定部 24 は、人位置を示す情報の座標を、画像における座標系から世界座標系へ周知の技術を用いて変換する。人位置推定部 24 は、認識した領域を示す情報（人位置情報）を推定結果とし、推定結果、及び顔の領域を含む画像情報（以下、顔画像情報という）を唇検出部 25 に出力する。

なお、人位置推定部 24 には、顔認識を、例えば、画像から顔のパーツ（顔の外形、髪の毛、眉毛、目、鼻、口等）を検出し、人位置推定部 24 に予め記憶されている顔認識用のデータベースと各パーツの位置関係とを比較することで、人の顔として妥当であるか判別し、妥当であれば人の顔であると認識する。そして、人位置推定部 24 には、認識した顔を含む領域の画像を、全体の画像である画像信号から抽出することで、顔画像情報を抽出する。顔画像情報には、顔の領域を含む画像と、顔画像が画像全体のうちどの領域であるかを示す情報とが含まれる。

一方、人位置推定部 24 は、判定結果が音源定位の推定結果を優先する場合、人位置推定を行わず、推定結果を唇検出部 25 に出力しない。

また、人位置推定部 24 は、世界座標系で表される推定した人毎の顔を含む領域の位置を示す情報を、音響信号処理部 40 の音響処理の事前情報生成部 43 に出力する。

【0033】

図 3 は、図 2 に示した画像 Ph1 において顔認識を行った結果の一例を説明する図である。図 3 に示す例では、符号 A11 ~ A14 それぞれに示す領域の画像が、話者 Sp1 ~ Sp4 それぞれに対応する顔画像である。なお、顔画像の領域は、少なくとも顔の外形を含む範囲であればよく、例えば上半身であってもよく、さらには、話者 Sp に対応する全ての領域の画像であってもよい。

10

20

30

40

50

## 【 0 0 3 4 】

図 1 に戻って、音源分離装置 1 の構成の説明を続ける。

唇検出部 2 5 は、人位置推定部 2 4 から入力された顔画像情報に基づいて、話者の唇の形状を周知の技術（例えば、特開 2 0 1 1 - 1 9 1 4 2 3 号公報参照）を用いて検出することで、発話を行っている話者を推定する。唇検出部 2 5 は、検出した検出結果に応じて、発話している話者の顔画像情報を選択する。唇検出部 2 5 は、選択した発話している話者の顔画像情報を発話区間検出部 2 6 に出力する。また、唇検出部 2 5 は、選択した発話している話者の顔画像情報に含まれる話者の顔を含む領域の位置情報を音響処理の事前情報生成部 4 3 に出力する。

## 【 0 0 3 5 】

発話区間検出部 2 6 は、唇検出部 2 5 から入力された検出結果に基づいて、周知の技術（例えば、特開 2 0 1 1 - 1 9 1 4 2 3 号公報参照）を用いて発話区間を検出する。発話区間検出部 2 6 は、検出した発話区間を示す情報、及び話者の顔画像情報を画像処理部 2 7 に出力する。また、発話区間検出部 2 6 は、検出した発話区間を示す情報を音響処理の事前情報生成部 4 3 に出力する。

## 【 0 0 3 6 】

なお、本実施形態では、唇検出部 2 5、発話区間検出部 2 6 を備える例を説明したが、これらの機能部を人位置推定部 2 4 が備えていてもよい。この場合、人位置推定部 2 4 は、発話区間を示す情報、人毎の顔を含む領域の位置を示す情報、発話を行っている人の顔を含む領域の位置を示す情報、及び発話区間を示す情報を音響処理の事前情報生成部 4 3

## 【 0 0 3 7 】

画像処理部 2 7 には、画像信号取得部 2 1 から画像信号が入力され、発話区間検出部 2 6 から発話区間を示す情報、及び話者の顔画像情報が入力される。または、画像処理部 2 7 には、画像信号の事前情報生成部 2 3 から画像信号の事前情報である話者の顔の領域を示す情報が入力される。画像処理部 2 7 は、入力された情報を用いて、発話区間毎に、発話していない人及び他の領域の画像の解像度を、入力された画像の解像度より低くする。例えば、撮像された画像が 3 0 0 [ b p i ( ビット / インチ ) ] であった場合、画像処理部 2 7 は、抽出する画像の解像度を 3 0 0 [ b p i ] に維持し、発話していない人及び他の領域の画像の解像度を、例えば 1 / 1 0 の 3 0 [ b p i ] に落とす。そして、画像処理部 2 7 は、発話区間毎に、解像度を落とした発話していない人及び他の領域の画像と、解像度を変更していない話者の顔画像とを合成する。

なお、画像処理部 2 7 は、発話が行われていない無音区間のとき、画像全体の解像度を下げるとしてもよい。

## 【 0 0 3 8 】

図 4 は、本実施形態に係る話者 S p 4 が発話中において画像処理された後の画像 P h 1 1 の一例を説明する図である。図 4 において、符号 P h 1 2 が示す領域の画像は、発話していない人及び他の領域の画像である。また、符号 A 1 4 が示す領域の画像は、話者 S p 4 の顔画像として抽出された画像である。図 4 に示すように、発話していない人及び他の領域の画像の解像度を元の画像の解像度に対して落としても、会議に参加している人の輪郭が残っているので、画像を見ている人は、会議の参加者数、発話している話者 S p 4 を確認することができる。そして、このように発話を行っていない人や他の領域の画像の解像度を下げることで、画像データの大きさを削減することができる。このとき、撮像された画像全体の解像度を下げているのではなく、話者 S p 4 の顔画像の解像度が維持されているため、図 4 に示した画像情報を記録したり他の装置に送信したりする場合、話者が誰であり発話している様子を観察者は画像から確認することができる。なお、図 4 に示した画像は、動画の一部であってもよい。この場合も、画像処理部 2 7 は、話者の顔を含む領域の映像の解像度を維持し、発話していない人及び他の領域の映像の解像度を落とすことで、同様の効果を得ることができる。

## 【 0 0 3 9 】

図 1 に戻って、音源分離装置 1 の構成の説明を続ける。

收音部 30 は、M 個 (M は 1 よりも大きい整数、例えば 8 個) のチャネルの音響信号を収録し、収録した M チャネルの音響信号を音響信号処理部 40 に送信する。收音部 30 は、例えば周波数帯域 (例えば 200 Hz ~ 4 kHz) の成分を有する音波を受信する M 個のマイクロホン 31-1 ~ 31-M を備えている。以下、マイクロホン 31-1 ~ 31-M のうち、特定しない場合は、単にマイクロホン 31 という。M 個のマイクロホン 31 は、それぞれ異なる位置に配置されている。收音部 30 は、収録した M チャネルの音響信号を無線で送信してもよいし、有線で送信してもよい。M が 1 よりも大きい場合には、送信の際にチャネル間で音響信号が同期していればよい。

【0040】

音響信号処理部 40 は、音響信号取得部 41、音響信号評価部 42、音響信号の事前情報生成部 43、音源方向推定部 44、音源分離部 45、雑音抑圧部 46、音響特徴量抽出部 47、発話区間検出部 48、及び音響処理部 49 (音源分離部) を含んで構成される。

【0041】

音響信号取得部 41 は、收音部 30 から送信された音響信号をチャネル毎に受信する。音響信号取得部 41 は、取得した音響信号をアナログ信号からデジタル信号に変換し、変換した音響信号を音響信号評価部 42、音響信号の事前情報生成部 43、及び音源方向推定部 44 に出力する。

【0042】

音響信号評価部 42 は、音響信号取得部 41 から入力された音響信号を評価する。例えば、音響信号評価部 42 は、信号の振幅の確認、音響信号の周波数成分の解析等を行うことで、音響信号に含まれる雑音成分の大きさが、所定の大きさ以上であるか否かを評価する。音響信号評価部 42 は、評価した評価結果を信頼性判定部 50 に出力する。なお、音響信号評価部 42 が行う評価方法については、後述する。

【0043】

音響信号の事前情報生成部 43 には、音響信号取得部 41 から音響信号が入力され、信頼性判定部 50 から判定結果が入力される。また、音響信号の事前情報生成部 43 には、人位置推定部 24 から人毎の顔を含む領域の位置を示す情報が入力され、唇検出部 25 から発話を行っている人の顔を含む領域の位置を示す情報が入力され、発話区間検出部 26 から発話区間を示す情報が入力される。音響信号の事前情報生成部 43 は、判定結果に基づいて、音響信号の事前情報を生成するか否かを決定する。ここで、音響信号の事前情報とは、発話を行っている話者の方向 (音源方向) を示す情報である。

音響信号の事前情報生成部 43 は、判定結果が人位置推定の推定結果を優先する場合、事前情報の生成処理を行わず、音響信号の事前情報を音源分離部 45 に出力しない。一方、音響信号の事前情報生成部 43 は、判定結果が音源定位の推定結果を優先する場合、発話区間毎に、発話を行っている人の顔を含む領域の位置を示す情報を用いて、音響信号の事前情報を生成し、生成した事前情報である話者の方向を示す情報を音源分離部 45 に出力する。なお、作成される話者の方向は、世界座標系で表される方位角である。

【0044】

音源方向推定部 44 には、音響信号取得部 41 から音響信号が入力され、信頼性判定部 50 から判定結果が入力される。音源方向推定部 44 は、判定結果に基づいて、入力された音響信号から人の音源方向の推定 (以下、音源定位ともいう) を行うか否かを決定する。音源方向推定部 44 は、判定結果が音響信号の方が画像信号より信頼性が高い場合、入力された音響信号を用いて、例えば MUSIC (Multiple Signal Classification; 多重信号分類) 法、ビームフォーミング法等によって音源毎の方向を推定し、推定した推定結果及び音響信号を音源分離部 45 に出力する。

【0045】

音源分離部 45 は、音源方向推定部 44 から入力された推定結果及び音響信号を用いて、または、音響処理の事前情報生成部 43 から入力された音響処理の事前情報及び音響信号を用いて、周知の手法、例えばブラインド信号分離手法、独立成分分析に基づくブラ

10

20

30

40

50

ンド音源分離手法、信号のスパース性を用いたブラインド音源分離手法等によって、音源を分離する。なお、信号がスパースであるとは、信号がほとんどの時間周波数において0であることを指す。音源分離部45は、分離した分離結果及び音響信号を雑音抑圧部46に出力する。

#### 【0046】

雑音抑圧部46は、音源分離部45から入力された分離結果を用いて、音響信号に含まれる雑音成分を、周知の手法、例えばHRL E (Histogram-based Recursive Level Estimation) 法、室内インパルス応答の逆フィルタ処理による手法、音源パワースペクトラム推定による手法、MTF (変調伝達関数または振幅伝達関数; Modulation Transfer Function) 理論に基づく手法、GSS (Geometric Sound Separation; 幾何学的音源分離) による手法等により抑圧する。雑音抑圧部46は、音源毎に残響抑圧された音響信号である音声信号を音響特徴量抽出部47に入力する。

10

#### 【0047】

音響特徴量抽出部47は、雑音抑圧部46から入力された音源毎に残響抑圧された音声信号から音響特徴量である例えばMSLS (Mel Scale Logarithmic Spectrum; メルスケール対数スペクトル) を抽出する。なお、MSLSは、音響認識の特徴量としてスペクトル特徴量を用い、MFCC (メル周波数ケプストラム係数; Mel Frequency Cepstrum Coefficient) を逆離散コサイン変換することによって得られる。音響特徴量抽出部47は、入力された音声信号と、抽出した音響特徴量とを音源毎に発話区間検出部48に出力する。また、音声特徴量は、MFCCのみを用いることもある。

20

#### 【0048】

発話区間検出部48は、音響特徴量抽出部47から入力された音響特徴量に基づき発話と発話との間、すなわち無音区間である非発話の対数尤度を算出し、算出した非発話の対数尤度が予め定められている値以上のとき、無音区間であると判別する。非発話の対数尤度の算出は、既存の手法、例えばデータベース発話区間検出法を用いる。発話区間検出部48は、音響信号における発話区間検出 (Audio VAD (Voice Activity Detection); A-VAD) の途中結果である非発話の対数尤度を用いるようにしてもよい (例えば、特開2011-191423号公報参照)。発話区間検出部48は、この無音区間以外の区間を発話区間と判別し、判別した発話区間を示す情報、及び発話区間毎の音響特徴量を音響処理部49、画像処理の事前情報生成部23、及び議事録作成部60に出力する。

30

#### 【0049】

音響処理部49は、発話区間毎に、発話区間に発話された音響信号 (音源分離処理かつ雑音抑圧処理済み) と、発話に対応する音響特徴量とを抽出する。

#### 【0050】

信頼性判定部50には、画像信号評価部22から評価結果と、音響信号評価部42から評価結果とが入力される。信頼性判定部50は、画像信号評価部22から評価結果と音響信号評価部42から評価結果とを、図5のような予め定められている対応表を用いて判定することで、音源定位の推定結果を優先するか、人位置推定の推定結果を優先するかを、決定する。

40

#### 【0051】

図5は、本実施形態に係る信頼性判定部50による判定結果の一例を説明する図である。図5に示すように信頼性判定部50は、画像信号評価部22の評価結果が「信頼性が高い」場合かつ音響信号評価部42の評価結果が「信頼性が高い」場合、音源定位の推定結果を優先する。なお、この場合は、人位置推定の推定結果を優先するようにしてもよい。

信頼性判定部50は、画像信号評価部22の評価結果が「信頼性が低い」場合かつ音響信号評価部42の評価結果が「信頼性が高い」場合、音源定位の推定結果を優先する。

信頼性判定部50は、画像信号評価部22の評価結果が「信頼性が高い」場合かつ音響

50

信号評価部 4 2 の評価結果が「信頼性が低い」場合、または、画像信号評価部 2 2 の評価結果が「信頼性が低い」場合かつ音響信号評価部 4 2 の評価結果が「信頼性が低い」場合、人位置推定の推定結果を優先する。画像信号評価部 2 2 の評価結果が「信頼性が低い」場合かつ音響信号評価部 4 2 の評価結果が「信頼性が低い」場合に人位置推定の推定結果を優先する理由は、後述する図 7、7、9、及び 10 に示した例のように、画像が暗すぎたり明るすぎたりコントラストが高すぎたりコントラストが低すぎても人の輪郭を推定できる場合があるからである。

#### 【 0 0 5 2 】

なお、図 5 に示した例は一例であり、画像信号評価部 2 2 の評価結果と音響信号評価部 4 2 の評価結果に応じた判定結果は、どちらの推定結果を優先するかを予め実験により決定しておいてもよい。

10

また、信頼性判定部 5 0 に入力される判定結果は、後述するように画像信号評価部 2 2 及び音響信号評価部 4 2 それぞれで算出された評価値であってもよい。信頼性判定部 5 0 は、入力された評価値に基づいて、どちらの推定結果を優先するかを決定するようにしてもよい。この場合においても、評価値の値に応じてどちらの推定結果を優先するかを、図 5 に示したような判定可能な表形式で自部に予め記憶させておく。この場合、信頼性判定部 5 0 は、画像信号評価部 2 2 及び音響信号評価部 4 2 それぞれで算出された評価値を正規化しておき、両方の値を比較することで判定するようにしてもよい。

#### 【 0 0 5 3 】

図 1 に戻って、音源分離装置 1 の説明を続ける。

20

議事録作成部 6 0 は、発話認識部 6 1、関連付部 6 2、及び記憶部 6 3 を含んで構成される。

#### 【 0 0 5 4 】

発話認識部 6 1 には、音響信号処理部 4 0 から検出された発話区間情報と、発話区間中の音響信号に対応する音響特徴量とが入力される。発話認識部 6 1 は、入力された発話区間情報と音響特徴量の M S L S 情報とを用いて発話認識を行う。発話認識は、例えば、汎用大語彙連続音声認識エンジンであるストリーム重み付を指定可能なマルチバンド Julius ( Y . Nishimura , et al . , " Speech recognition for a humanoid with motor noise utilizing missing feature theory , " Humanoids 2006 , pp . 26 - 33 ) を用いて行う。なお、発話認識部 6 1 は、周知の構文解析、係り受け解析等を行うことで、発話認識を行うようにしてもよい。なお、認識結果はテキスト情報である。発話認識部 6 1 は、認識した認識結果を関連付部 6 2 に出力する。

30

#### 【 0 0 5 5 】

関連付部 6 2 には、発話認識部 6 1 から認識結果が入力され、画像処理部 2 7 で処理された画像情報が入力される。関連付部 6 2 は、図 6 に示すように、発話内容、及び発話中の画像を関連付けて議事録情報を生成する。関連付部 6 2 は、生成した議事録情報を記憶部 6 3 に記憶させる。ここで、発話中の画像とは、図 4 に示した例のように、話者の領域の解像度を変更しない画像と、発話していない人及び他の領域の画像の解像度を低くしたとを合成した画像である。なお、関連付部 6 2 は、上記の情報にさらに話者を示す情報、発話区間中の話者の音声信号も関連付けて記憶部 6 3 に記憶させるようにしてもよい。この場合、音響信号処理部 4 0 から発話区間中の話者を示す情報、及び話者の音声信号も入力される。

40

#### 【 0 0 5 6 】

図 6 は、関連付部 6 2 によって関連付けられた議事録情報の一例を説明する図である。図 6 に示す例では、発話内容が「それでは、本日の会議を始めます。Bさん、本日の議事は何か?」と、発話中の画像として話者 Sp 1 の領域の解像度を変更せず他の人や他の領域の解像度を低くした画像とが関連付けられている。議事録の閲覧者は、このような議事録情報により、会議に参加している人数、その中の誰が話者なのか画像によって知る

50

ことができる。

【 0 0 5 7 】

記憶部 6 3 には、図 6 に示したような発話内容、及び発話中の画像が関連付けられて記憶される。なお、関連付部 6 2 が議事録情報として、話者識別情報、及び発話区間の音声情報を関連付けた場合、これらの情報も関連付けて記憶するようにしてもよい。

【 0 0 5 8 】

以上のように、本実施形態では、画像信号の方が音響信号より信頼性が高い場合、人位置推定部 2 4 が、入力された画像信号を用いて、画像に写っている人毎の位置を周知の技術を用いて推定する。そして、音響信号処理部 4 0 は、この推定結果を用いて音響信号の事前情報を生成し、生成した音響信号の事前情報を用いて音源定位処理を行う。

10

一方、判定結果が音響信号の方が画像信号より信頼性が高い場合、音源方向推定部 4 4 が、入力された音響信号に対して、周知の技術を用いて音源定位の推定を行う。そして、画像信号処理部 2 0 は、この推定結果を用いて画像信号の事前情報を生成し、生成した画像信号の事前情報を用いて人位置推定処理を行う。

すなわち、本実施形態の音源分離装置 1 は、画像信号と音響信号とによる情報を用いて相互に補完し合うことで、話者の位置の検出、音源定位を行う。

【 0 0 5 9 】

< 画像信号の評価 >

次に、画像信号評価部 2 2 が行う処理について説明する。

図 7 は、暗すぎる画像 P h 2 の一例である。図 8 は、明るすぎる画像 P h 3 の一例である。

20

図 7 に示す例の画像全体が暗すぎるため、または図 8 に示す例の画像全体が明るすぎるため、画像 P h 2 及び画像 P h 3 を画像認識させた場合、話者数、各話者の位置、話者が発話しているか否かを精度良く推定することができない場合がある。

【 0 0 6 0 】

図 9 は、図 2、図 7、及び図 8 の各画像のヒストグラムを説明する図である。図 9 において、横軸は輝度、縦軸はピクセル数である。なお、輝度は、左側が最小値であり、右側が最大値である。符号 g 1 0 1 が示す画像は、図 2 の画像 P h 1 のヒストグラムの図である。符号 g 1 0 1 が示す画像のように、画像 P h 1 の画像信号の成分は、輝度の最小値から最大値の範囲に分布している。

30

【 0 0 6 1 】

符号 g 1 0 2 が示す画像は、図 7 の画像 P h 2 のヒストグラムの図である。符号 g 1 1 2 が示す領域の画像のように、画像 P h 2 の画像信号の成分は、輝度の最小値から中間値以下の範囲に分布している。

符号 g 1 0 3 が示す画像は、図 8 の画像 P h 3 のヒストグラムの図である。符号 g 1 1 3 が示す領域の画像のように、画像 P h 3 の画像信号の成分は、輝度の中間値以上から最大値の範囲に分布している。すなわち、明るすぎる画像及び暗すぎる画像は、ヒストグラムにおいて、輝度が最小値側または最大値側に偏っている。

このように、画像信号評価部 2 2 は、ヒストグラムを解析することで、入力された画像信号による画像の輝度成分が輝度の最小値側か最大値側に偏っている場合、暗すぎる画像または明るすぎる画像であると評価することができる。

40

【 0 0 6 2 】

図 1 0 は、コントラストが低すぎる画像 P h 2 1 の一例である。図 1 1 は、コントラストが高すぎる画像 P h 2 2 の一例である。

図 1 0 に示す例の画像 P h 2 1 はコントラストが低すぎるため、または図 1 1 に示す画像 P h 2 2 のコントラストが高すぎるため、画像 P h 2 1 及び画像 P h 2 2 を画像認識させた場合、話者数、各話者の位置、話者が発話しているか否かを精度良く推定することができない場合がある。

【 0 0 6 3 】

図 1 2 は、図 1 0 及び図 1 1 の画像のヒストグラムを説明する図である。図 1 2 におい

50

て、横軸は輝度、縦軸はピクセル数である。

符号 g 1 2 1 が示す画像は、図 1 0 の画像 P h 2 1 のヒストグラムの図である。符号 g 1 3 1 が示す領域の画像のように、画像 P h 2 1 の画像信号の成分は、輝度の中間値を中心に分布し、輝度が最小値及び最大値に近い領域には分布していない。

符号 g 1 2 2 が示す画像は、図 1 1 の画像 P h 2 2 のヒストグラムの図である。符号 g 1 4 1 及び g 1 4 2 が示す領域の画像のように、画像 P h 2 2 の画像信号の成分は、輝度の最小値付近と最大値付近のみに分布している。すなわち、コントラストが低すぎる画像及びコントラストが高すぎる画像は、ヒストグラムにおいて、輝度の中心付近のみに分布、または最小値付近と最大値付近のみに分布する。

このように、画像信号評価部 2 2 は、ヒストグラムを解析することで、入力された画像信号による画像の輝度成分に所定の輝度範囲より狭い範囲のみに画像信号の成分が分布している場合、コントラストが低すぎる画像またはコントラストが高すぎる画像であると評価することができる。そして、本実施形態によれば、上述したように評価を行うことで、画像信号を定量的に且つ簡便に評価することができる。

#### 【 0 0 6 4 】

図 1 3 は、本実施形態に係る画像信号評価部 2 2 の処理手順のフローチャートである。  
(ステップ S 1) 画像信号評価部 2 2 は、入力された画像信号の輝度毎のピクセル数であるヒストグラムを算出する。

(ステップ S 2) 画像信号評価部 2 2 は、ピクセル数が所定の値以上の連続する輝度の範囲を検出する。

#### 【 0 0 6 5 】

(ステップ S 3) 画像信号評価部 2 2 は、算出した範囲が、所定の範囲以上であるか否かを判別する。画像信号評価部 2 2 は、算出した範囲が所定の範囲以上であると判定した場合 (ステップ S 3 ; Y E S)、ステップ S 4 に進み、算出した範囲が所定の範囲以上ではないと判定した場合 (ステップ S 3 ; N O)、ステップ S 5 に進む。

#### 【 0 0 6 6 】

(ステップ S 4) 画像信号評価部 2 2 は、画像の信頼性が高いと判別し、判別した結果を、評価結果を示す情報として信頼性判定部 5 0 に出力し、処理を終了する。

(ステップ S 5) 画像信号評価部 2 2 は、画像の信頼性が低いと判別し、判別した結果を、評価結果を示す情報として信頼性判定部 5 0 に出力し、処理を終了する。

#### 【 0 0 6 7 】

なお、本実施形態では、図 7 ~ 図 1 2 において、画像の信頼性を輝度とコントラストに基づいて判別する例を説明したが、少なくともどちらか一方に基づいて判別するようにしてもよい。

また、図 1 3 に示した処理手順は一例であり、これに限られない。例えば、画像信号評価部 2 2 は、ヒストグラムを算出後、ピクセルが所定値以上ある輝度の数をカウントし、総ピクセル数をカウントした値 (以下、カウント値という) で除算して判定値を算出するようにしてもよい。例えば、最小輝度が 0、最大輝度が 2 5 5 の場合、カウント値は 0 ~ 2 5 5 の値になる。総ピクセル数は一定のため、カウント値が多いほど判定値が小さくなり、カウント値が少ないほど判定値が大きくなる。具体的には、図 2 の総ピクセル数が 2 6 万ピクセルであるとすると、図 2 に示した画像 P h 1 の輝度の範囲が 0 ~ 2 5 5 であり、図 7 に示した画像 P h 2 の輝度の範囲は 0 ~ 1 1 1 であり、図 1 1 に示した画像 P h 2 2 の輝度範囲は 0 ~ 1 5 と 2 4 0 ~ 2 5 5 である。このため、画像 P h 1 の判定値は約 1 0 2 0 であり、画像 P h 2 の判定値は約 2 3 4 2 であり、画像 P h 2 2 の判定値は約 8 6 6 7 である。この場合、画像信号評価部 2 2 は、判定値が所定の値以上の場合に画像の信頼性が低いと判別し、判定値が所定の値未満の場合に画像の信頼性が高いと判別するようにしてもよい。そして、画像信号評価部 2 2 は、この判別した結果を、評価結果を示す情報として信頼性判定部 5 0 に出力するようにしてもよい。本実施形態によれば、このように判定値を算出して評価を行うことで、画像信号を定量的に且つ簡便に評価することができる。

10

20

30

40

50

## 【 0 0 6 8 】

## &lt; 音響信号の評価 &gt;

次に、音響信号評価部 4 2 が行う処理について説明する。

図 1 4 は、雑音が少ない音響信号の一例を示す図である。図 1 5 は、雑音が多い音響信号の一例を示す図である。図 1 4 及び図 1 5 において、横軸は時刻 [ s ( 秒 ) ]、縦軸は信号レベル [ V ] である。また、符号 S g 1 ~ S g 3 に示す領域の波形は、話者 S p による発話による音響信号の波形を表している。符号 S g 4 及び S g 1 1 が示す領域の波形は、音響信号に含まれる雑音信号の波形を表している。なお、図 1 4 と図 1 5 とにおける発話による音響信号は、同じタイミングかつ同じ信号レベルである。

## 【 0 0 6 9 】

図 1 4 に示す例では、符号 S g 4 に示す領域の波形のように、雑音信号の信号レベルの振幅は、 $0.01 [V_{p-p}]$  以下である。時刻約  $20.7 [s] \sim 21.3 [s]$  の区間に符号 S g 1 に示される領域の波形が観測される。また、時刻約  $23.0 [s] \sim 23.8 [s]$  の区間に符号 S g 2 に示される領域の波形が観測され、時刻約  $25.5 [s] \sim 26.3 [s]$  の区間に符号 S g 3 に示される領域の波形が観測される。

一方、図 1 5 に示す例では、符号 S g 1 1 に示す領域の波形のように、雑音信号の信号レベルの振幅は、約  $0.1 [V_{p-p}]$  である。このため、発話による波形 ( S g 1 ~ S g 3 ) は、 $\pm 0.05 [V]$  を越える区間のみ観測される。

図 1 3 に示す音響信号を用いて音源定位や音声認識を行った場合と比較して、図 1 4 に示す音響信号を用いて音響信号を用いて音源定位や音声認識を行った場合の方が、音源方向の推定 ( 音源定位 ) の精度が落ち、さらに音声認識の精度が落ちる。

このため、本実施形態の音響信号評価部 4 2 は、例えば H R L E 法を用いて雑音パワーを算出し、算出した雑音パワーに基づいて、音響信号の信頼性を評価する。

## 【 0 0 7 0 】

ここで、H R L E 法の概要について説明する。

音響信号評価部 4 2 は、入力された音響信号を周波数領域の複素入力スペクトル  $Y ( k , l )$  に変換する。k は周波数を表すインデックスであり、l は各フレームを表すインデックスである。次に、音響信号評価部 4 2 は、複素入力スペクトル  $Y ( k , l )$  に基づいてパワースペクトル  $| Y ( k , l ) | ^ 2$  を算出する。 $| \dots |$  は、複素数  $\dots$  の絶対値を示す。次に、音響信号評価部 4 2 は、パワースペクトル  $| Y ( k , l ) | ^ 2$  に含まれる雑音成分のパワースペクトル  $( k , l )$  を、H R L E 法を用いて算出する。H R L E 法は、ある周波数について、パワー毎の頻度を計数してヒストグラムを生成し、生成したヒストグラムにおいて計数した頻度をパワーについて累積した累積頻度を算出し、予め定めた累積頻度を与えるパワーを雑音パワー  $( k , l )$  と定める方法である。従って、H R L E 法では、累積頻度が大きいほど、推定される雑音パワーが大きくなり、累積頻度が小さいほど、推定される雑音パワーが小さくなる ( 例えば特願 2 0 1 3 - 0 1 3 2 5 1 号公報参照 ) 。

## 【 0 0 7 1 】

図 1 6 は、本実施形態に係る音響信号評価部 4 2 の処理手順のフローチャートである。(ステップ S 1 1 ) 音響信号評価部 4 2 は、例えば H R L E 法を用いて雑音パワーを算出する。

(ステップ S 1 2 ) 音響信号評価部 4 2 は、雑音パワーが所定の値以上であるか否かを判別する。音響信号評価部 4 2 は、雑音パワーが所定の値以上であると判別した場合 (ステップ S 1 2 ; Y E S )、ステップ S 1 3 に進み、雑音パワーが所定の値以上ではないと判別した場合 (ステップ S 1 2 ; N O )、ステップ S 1 4 に進む。

## 【 0 0 7 2 】

(ステップ S 1 3 ) 音響信号評価部 4 2 は、音響信号の信頼性が低いと判別し、判別した結果を、評価結果を示す情報として信頼性判定部 5 0 に出力し、処理を終了する。

(ステップ S 1 4 ) 音響信号評価部 4 2 は、音響信号の信頼性が高いと判別し、判別した結果を、評価結果を示す情報として信頼性判定部 5 0 に出力し、処理を終了する。

10

20

30

40

50

## 【 0 0 7 3 】

以上のように、本実施形態では、音響信号の雑音成分の大きさを、残響抑圧処理の結果に基づいて算出された値に基づいて音響信号の信頼性を評価するようにしたので、音響信号を定量的に且つ簡便に評価することができる。

## 【 0 0 7 4 】

なお、上述した例では、音響信号評価部 4 2 が雑音パワーを算出する例を説明したが、雑音抑圧部 4 6 が算出し、算出した雑音パワーを示す値を音響信号評価部 4 2 に出力するようにしてもよい。

さらに、音響信号評価部 4 2 は、入力された音響信号と、雑音抑圧部 4 6 によって雑音成分が抑圧された後の音響信号との比を算出することで、雑音パワーが大きいか小さいかを評価するようにしてもよい。この場合、雑音抑圧部 4 6 は、雑音抑圧後の音響信号を音響信号評価部 4 2 に出力する。

10

## 【 0 0 7 5 】

< 音源分離装置の処理 >

次に、音源分離装置 1 が行う処理について説明する。

図 1 7 は、本実施形態に係る音源分離装置 1 が行う処理手順の一例を説明するフローチャートである。

## 【 0 0 7 6 】

(ステップ S 1 0 1) 画像信号取得部 2 1 は、所定の間隔毎に撮像部 1 0 によって撮像された画像信号を取得する。なお、画像信号は動画であっても静止画であってもよい。

20

(ステップ S 1 0 2) 音響信号取得部 4 1 は、收音部 3 0 によって収録された音響信号を取得する。なお、ステップ S 1 0 1 の処理とステップ S 1 0 2 との処理は、処理順番が逆であってもよく、同時に行われてもよい。

## 【 0 0 7 7 】

(ステップ S 1 0 3) 画像信号評価部 2 2 は、図 1 3 を用いて説明した画像信号を評価する処理を行う。

(ステップ S 1 0 4) 音響信号評価部 4 2 は、図 1 6 を用いて説明した音響信号を評価する処理を行う。

## 【 0 0 7 8 】

(ステップ S 1 0 5) 信頼性判定部 5 0 は、画像信号評価部 2 2 から入力された評価結果と、音響信号評価部 4 2 から入力された評価結果とに基づいて、音源定位の推定結果を優先するか、人位置推定の推定結果を優先するかを決定する。

30

(ステップ S 1 0 6) 信頼性判定部 5 0 は、人位置推定の推定結果を優先する場合(ステップ S 1 0 6 ; 人位置推定の推定結果を優先)、ステップ S 1 0 7 に進み、音源定位の推定結果を優先する場合(ステップ S 1 0 6 ; 音源定位の推定結果を優先)、ステップ S 1 0 8 に進む。

## 【 0 0 7 9 】

(ステップ S 1 0 7) 画像処理信号部 2 0 及び音響信号処理部 4 0 は、人位置推定の推定結果を優先する処理を行い、処理をステップ S 1 0 9 に進める。

(ステップ S 1 0 8) 画像処理信号部 2 0 及び音響信号処理部 4 0 は、音源方向の推定結果を優先する処理を行い、処理をステップ S 1 0 9 に進める。

40

(ステップ S 1 0 9) 議事録作成部 6 0 は、議事録作成の処理を行う。

以上で、音源分離装置 1 が行う処理を終了する。

## 【 0 0 8 0 】

< 人位置推定の推定結果を優先する処理 >

次に、人位置推定の推定結果を優先する処理を説明する。

図 1 8 は、本実施形態に係る人位置推定の推定結果を優先する処理手順のフローチャートである。

(ステップ S 2 0 1) 人位置推定部 2 4 は、信頼性判定部 5 0 から人位置推定の推定結果を優先することを示す判定結果が入力された場合、画像信号取得部 2 1 から入力された画

50

像信号を用いて、画像に写っている人毎の顔を含む領域の画像位置を周知の画像認識技術を用いて推定する。

【0081】

(ステップS202) 人位置推定部24は、推定した結果に基づいて、各人の顔を含む領域の画像(顔画像)を、画像信号取得部21から入力された画像信号(全体画像)から抽出する。

(ステップS203) 人位置推定部24は、世界座標系で表される推定した人毎の顔を含む領域の位置を示す情報を、音響処理の事前情報生成部43に出力する。

【0082】

(ステップS204) 唇検出部25は、人位置推定部24から入力された顔画像情報に基づいて、話者の唇の形状を周知の技術(例えば、特開2011-191423号公報参照)を用いて検出することで、発話を行っている話者を推定する。次に、唇検出部25は、検出した検出結果に応じて、発話している話者の顔画像情報を選択する。

10

(ステップS205) 発話区間検出部26は、唇検出部25によって検出された検出結果に基づいて、周知の技術(例えば、特開2011-191423号公報参照)を用いて発話区間を検出する。

【0083】

(ステップS206) 画像処理部27は、入力された情報を用いて、発話区間、発話していない人及び他の領域の画像の解像度を、入力された画像の解像度より低くする。

(ステップS207) 画像処理部27は、解像度を落とした発話していない人及び他の領域の画像と、解像度を変更していない話者の顔画像とを合成する。

20

【0084】

(ステップS208) 音響処理の事前情報生成部43は、信頼性判定部50から人位置推定の推定結果を優先することを示す判定結果が入力された場合、発話区間毎に、発話を行っている人の顔を含む領域の位置を示す情報を用いて、音響信号の事前情報を生成する。

【0085】

(ステップS209) 音源分離部45は、音響処理の事前情報及び音響信号を用いて、周知の手法によって、音源を分離する。

(ステップS210) 雑音抑圧部46は、音源分離部45から入力された分離結果を用いて、音響信号に含まれる雑音成分を、周知の手法により抑圧する。

30

【0086】

(ステップS211) 音響特徴量抽出部47は、雑音抑圧部46から入力された音源毎に残響抑圧された音響信号から音響特徴量である例えばMSLSを抽出する。

(ステップS212) 発話区間検出部48は、音響特徴量抽出部47から入力された特徴量に基づき発話と発話との間、すなわち無音区間を検出する。次に、発話区間検出部48は、検出した無音区間を用いて発話区間を検出する。

【0087】

(ステップS213) 音響処理部49は、発話区間毎に、発話区間に発話された音響信号(音源分離処理かつ雑音抑圧処理済み)と、発話に対応する音響特徴量とを抽出する。

(ステップS214) 発話認識部61は、入力された発話区間を示す情報と音響特徴量とを用いて発話認識を行う。

40

(ステップS215) 関連付部62は、発話内容、及び発話中の画像を関連付けて議事録情報を生成する。次に、関連付部62は、生成した議事録情報を記憶部63に記憶させる。

以上で、人位置推定の推定結果を優先する処理を終了する。

【0088】

<音源定位の推定結果を優先する処理>

次に、音源定位の推定結果を優先する処理を説明する。

図19は、本実施形態に係る音源定位の推定結果を優先する処理手順のフローチャートである。

50

(ステップS301) 音源方向推定部44は、信頼性判定部50から音源定位の推定結果を優先することを示す判定結果が入力された場合、入力された音響信号から音源毎の方向を、例えばMUSIC法、ビームフォーミング法等によって推定する。なお、話者の特定は、音源方向推定部44によって音源方向の推定結果に基づいて行われる。

【0089】

(ステップS302) 音源分離部45は、音源方向推定部44から入力された推定結果及び音響信号を用いて、周知の手法、例えばブラインド信号分離手法、独立成分分析に基づくブラインド音源分離手法、信号のスパース性を用いたブラインド音源分離手法等によって、音源を分離する。

【0090】

(ステップS303) 雑音抑圧部46は、音源分離部45から入力された分離結果を用いて、音響信号に含まれる雑音成分を、周知の手法、例えばHRL E法、室内インパルス応答の逆フィルタ処理による手法、音源パワースペクトラム推定による手法、MTF理論に基づく手法、GSSによる手法等により抑圧する。

【0091】

(ステップS304) 音響特徴量抽出部47は、雑音抑圧部46から入力された音源毎に残響抑圧された音声信号から音響特徴量である例えばMSLSを抽出する。

(ステップS305) 発話区間検出部48は、音響特徴量抽出部47から入力された音響特徴量に基づき無音区間と発話区間とを検出する。

【0092】

(ステップS306) 音響処理部49は、発話区間毎に、発話区間に発話された音響信号(音源分離処理かつ雑音抑圧処理済み)と、発話に対応する音響特徴量とを抽出する。

(ステップS307) 画像信号の事前情報生成部23は、音源方向の推定結果を用いて、発話区間毎に画像信号の事前情報の生成処理を行い、生成した画像信号の事前情報である話者の顔の領域を示す情報を画像処理部27に出力する。

【0093】

(ステップS308) 画像処理部27は、入力された情報を用いて、発話区間毎に、発話していない人及び他の領域の画像の解像度を、入力された画像の解像度より低くする。

(ステップS309) 画像処理部27は、解像度を落とした発話していない人及び他の領域の画像と、解像度を変更していない話者の顔画像とを合成する。

【0094】

(ステップS310) 発話認識部61は、入力された発話区間を示す情報と音響特徴量とを用いて発話認識を行う。

(ステップS311) 関連付部62は、発話内容、及び発話中の画像を関連付けて議事録情報を生成する。次に、関連付部62は、生成した議事録情報を記憶部63に記憶させる。

以上で、音源定位の推定結果を優先する処理を終了する。

【0095】

以上のように、本実施形態の音源分離装置(例えば音源分離装置1)は、音響信号を収録する収録部(例えば収録部30)と、画像を撮像する撮像部(例えば撮像部10)と、収録された音響信号を評価する音響信号評価部(例えば音響信号評価部42)と、撮像された画像信号を評価する画像信号評価部(例えば画像信号評価部22)と、音響信号評価部と画像信号評価部とによって評価された結果に基づいて、収録された音響信号に基づいて音源方向の推定を行うか、撮像された画像に対する音源方向の推定を行うかを選択する選択部(例えば信頼性判定部50)と、収録された音響信号に基づいて推定された音源方向を示す情報を用いて、画像から話者の方向を示す情報を推定する人位置推定部と、撮像された画像に基づいて推定された音源方向を示す情報を用いて、音響信号に対して音源の方向を推定する音源方向推定部(例えば音源方向推定部44)と、推定された音源の方向に基づいて、音源の方向に対応する音響信号を音響信号から抽出する音源分離部(例えば音源方向推定部44)と、を備える。

10

20

30

40

50

## 【 0 0 9 6 】

この構成によって、本実施形態の音源分離装置 1 は、撮像した画像と収録した音響信号の評価結果に応じて、少なくとも一方に基づいて音源の方向を推定することができる。このため、本実施形態によれば、画像信号の信頼性が低い場合には、音響信号を用いて音源の方向を推定することができ、音響信号の信頼性が低い場合には、画像信号に基づいて音源の方向を推定することができる。そして、本実施形態によれば、推定した結果に基づいて、発話された音声信号を分離することができる。このように、本実施形態によれば、音響信号と画像信号とを相互補完して音源分離を行うことができる。

## 【 0 0 9 7 】

また、本実施形態の音源分離装置（例えば音源分離装置 1）において、音源分離部（例えば音源分離部 4 5、音響処理部 4 9）は、音響信号を用いて音源の方向を推定し、または、人位置推定部（例えば人位置推定部 2 4）によって推定された話者の方向を示す情報に基づいて音源の方向を推定し、音源の方向の推定結果を用いて、音響信号を音源毎に分離することで抽出し、音源分離部によって分離した音源毎の音響信号の特徴量を算出する特徴量算出部（例えば音響特徴量抽出部 4 7）と、音響信号の特徴量に基づいて発話内容を認識してテキスト情報に変換する音声認識部（例えば発話認識部 6 1）、を備える。

10

## 【 0 0 9 8 】

この構成によって、本実施形態の音源分離装置 1 は、画像信号または音響信号のうち、少なくとも一方の音源の方向の推定結果に基づいて分離された音響信号に対して音声認識を行いテキスト化することができる。この結果、本構成によれば、画像信号または音響信号のいずれか一方の信頼性が低い場合であっても、議事録におけるテキスト認識率の精度を向上させることができる。

20

## 【 0 0 9 9 】

なお、本実施形態では、画像信号の評価を画像信号評価部 2 2 が行い、音響信号の評価を音響信号評価部 4 2 が行う例を説明したが、信頼性判定部 5 0 が評価するようにしてもよい。この場合、撮像部 1 0 または画像信号処理部 2 0 は、撮像された画像を信頼性判定部 5 0 に出力する。また、收音部 3 0 または音響信号処理部は、収録された音響信号を信頼性判定部 5 0 に出力する。

## 【 0 1 0 0 】

また、本実施形態では、図 4 に示したように、話者が 1 人の場合を例に説明したが、これに限られない。同時に発話する話者は複数であってもよい。例えば、図 4 において、話者 S p 2 と話者 S p 4 が同時に発話している場合、音源分離装置 1 は、話者 S p 2 と話者 S p 4 の顔を含む領域の画像の解像度を変更しないようにする。

30

## 【 0 1 0 1 】

また、本実施形態では、話者 S p の顔を含む領域の画像について、解像度を変更しない例を説明したが、これに限られない。撮像された画像の解像度が十分に高い（例えば 6 0 0 [ b p i ] 場合、話者を識別でき、発話していることがわかる程度の解像度に変更するようにしてもよい。

## 【 0 1 0 2 】

また、本実施形態では、撮像部 1 0 が 1 台の例を説明したが、撮像部 1 0 は複数のカメラを有していてもよい。例えば、左右のテーブルに、おのおの参加者が着席している場合、左右のテーブル毎にカメラを設置してもよい。この場合、画像信号処理部 2 0 の各部分は、複数の画像信号について、上述した処理を行う。そして、画像処理部 2 7 は、このようにして生成された複数の画像（例えば左側のカメラによる画像と右側カメラによる画像）をそのまま議事録作成部 6 0 に出力してもよく、または複数の画像をカメラの配置に応じて合成するようにしてもよい。

40

## 【 0 1 0 3 】

また、会議の出席者が多数であり、撮像部 1 0 で全ての参加者の撮像が困難な場合、例えば、少なくとも全体画像は話者を含む画像であればよい。この場合、音源分離装置 1 は、話者を限定し、限定した話者の方向に撮像部 1 0 を向けて撮像するようにしてもよい。

50

すなわち、音源方向推定結果または人位置推定結果に基づいて、音源分離装置 1 が撮像部 10 の撮像方向を制御するようにしてもよい。

さらに、参加者が多数の場合、音源分離装置 1 は、音源方向推定結果または人位置推定結果に基づいて、撮像部 10 の画角を話者及びその周辺の人を含む画角に制御するようにしてもよい。

#### 【0104】

< 第 2 実施形態 >

図 20 は、本実施形態に係る音源分離装置 1 A の構成を示すブロック図である。図 20 に示すように、音源分離装置 1 A は、撮像部 10、画像信号処理部 20、收音部 30、音響信号処理部 40、信頼性判定部 50（選択部）、議事録作成部 60、及び送信部 70 を含んで構成される。音源分離装置 1 A は、ネットワーク 90 を介して受信部 80 と接続される。ネットワーク 90 は、有線または無線であってもよい。また、受信部には、議事録作成部 60 が接続される。第 1 実施形態の音源分離装置 1（図 1）と同じ機能を有する機能部には、同じ符号を用いて説明を省略する。

10

#### 【0105】

送信部 70 には、画像信号処理部 20 から図 4 に示したような画像処理後の画像情報が入力される。また、送信部 70 には、音響信号処理部 40 から検出された発話区間情報と、発話区間中の音響信号に対応する音響特徴量の MSLS 情報が入力される。送信部 70 は、入力された情報を、ネットワーク 90 を介して受信部 80 に送信する。

#### 【0106】

受信部 80 は、送信部 70 からの情報を、ネットワーク 90 を介して受信し、受信した情報を議事録作成部 60 に出力する。

20

#### 【0107】

音源分離装置 1 A から議事録作成部 60 へ各種の情報を送信する場合、仮に撮像部 10 で撮像された画像をそのまま送信すると、送信データが大きくなる。これにより、ネットワーク 90 に負荷が増大し、また作成された議事録と記録される画像情報の容量も大きくなるため、議事録のファイル容量が増大する。ここで、話者の顔を含む領域のみの画像情報を音源分離装置 1 A から議事録作成部 60 へ送信した場合、ファイル容量は小さくなるが、話者以外の参加者が画像に写っていないため、会議の様子がわかりにくくなる。また、議事録の閲覧者は、話者が誰に話しかけているのかも画像から判断できなくなる。

30

一方、本実施形態のように、話者の顔を含む領域の画像の解像度を変更せず、話者以外の人と他の領域の画像の解像度を低くした画像を合成し、この合成した画像を発話データと関連づけて議事録作成部 60 に送信することで、ファイル容量を軽減でき、かつ会議の様子も把握可能な画像を提供することができる。

#### 【0108】

以上のように、本実施形態の音源分離装置（例えば音源分離装置 1 A）において、発話区間ごとに、画像信号から抽出された話者の顔を含む領域の画像と、抽出された音響信号とが関連付けられた情報を送信する送信部（例えば送信部 70）、を備える。

#### 【0109】

この構成により、本実施形態の音源分離装置 1 A では、発話区間毎に、画像信号または音響信号のうち、少なくとも一方の音源の方向の推定結果に基づいて分離された音響信号と抽出された話者の顔を含む画像とを関連付けて、他の装置へ送信することができる。この結果、本構成の音源分離装置をテレビ会議等に用いる場合、画像信号または音響信号のいずれか一方の信頼性が低い場合であっても、分離された音響信号と抽出された話者の顔を含む画像とを関連付けて、他の装置へ送信することができる。

40

#### 【0110】

第 2 実施形態に示した音源分離装置 1 A を、会議室毎に設置することで、ネットワークを介したテレビ会議を行うことができる。この場合、上述したように、送信される画像データの容量を小さくすることができる。

#### 【0111】

50

なお、本実施形態において、音源分離装置 1 A は、受信部 8 0 及び議事録作成部 6 0 を含んで構成されていてもよい。

【 0 1 1 2 】

なお、第 1 実施形態及び第 2 実施形態では、話者以外の人と他の領域の画像の解像度を低くする例を説明したが、これに限られない。画像は、例えば、話者 S p の画像のみカラー画像とし、話者以外の人と他の領域の画像をグレースケールまたは白黒の 2 値化した画像であってもよい。また、話者以外の人と他の領域の画像がぼやけるような画像フィルタ（例えばガウシアンフィルタ等）を用いた画像処理を行うようにしてもよい。この場合であっても、議事録の閲覧者は、会議中の様子と話者を画像から把握することができる。

【 0 1 1 3 】

なお、第 1 実施形態及び第 2 実施形態では、音源分離装置（ 1 または 1 A ）が会議室に取り付けられている例を説明したが、これに限られない。例えば、撮像部 1 0 と收音部 3 0 とが会議室に設置され、画像信号処理部 2 0 、音響信号処理部 4 0 、信頼性判定部 5 0 、及び議事録作成部 6 0 は、会議室とは別の場所に設置されていてもよい。

また、音源分離装置（ 1 または 1 A ）は、例えば議事録作成装置、人型のロボット、携帯端末（スマートフォン、タブレット、携帯ゲーム機等）、P C（パーソナルコンピュータ）等が備えていてもよい。

【 0 1 1 4 】

なお、第 1 実施形態及び第 2 実施形態では、議事録情報として、発話中の画像を関連付ける例を説明したが、これに限られない。議事録情報には、例えば、話者の顔の領域を含む画像のみを関連付けるようにしてもよい。

【 0 1 1 5 】

また、第 1 実施形態及び第 2 実施形態で説明した音源分離装置（ 1 または 1 A ）の構成は一例であり、音源分離装置（ 1 または 1 A ）は、これらの全ての機能部のうち、用途に応じて必要な機能部のみ備えるようにしてもよく、他の機能部を備えるようにしてもよい。

【 0 1 1 6 】

なお、本発明における音源分離装置（ 1 または 1 A ）の機能を実現するためのプログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録されたプログラムをコンピュータシステムに読み込ませ、実行することにより画像信号に対する処理や音響信号に対する処理等を行ってもよい。なお、ここでいう「コンピュータシステム」とは、OS や周辺機器等のハードウェアを含むものとする。また、「コンピュータシステム」は、ホームページ提供環境（あるいは表示環境）を備えた WWW システムも含むものとする。また、「コンピュータ読み取り可能な記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM、CD-ROM 等の可搬媒体、コンピュータシステムに内蔵されるハードディスク等の記憶装置のことをいう。さらに「コンピュータ読み取り可能な記録媒体」とは、インターネット等のネットワークや電話回線等の通信回線を介してプログラムが送信された場合のサーバやクライアントとなるコンピュータシステム内部の揮発性メモリ（RAM）のように、一定時間プログラムを保持しているものも含むものとする。

【 0 1 1 7 】

また、上記プログラムは、このプログラムを記憶装置等に格納したコンピュータシステムから、伝送媒体を介して、あるいは、伝送媒体中の伝送波により他のコンピュータシステムに伝送されてもよい。ここで、プログラムを伝送する「伝送媒体」は、インターネット等のネットワーク（通信網）や電話回線等の通信回線（通信線）のように情報を伝送する機能を有する媒体のことをいう。また、上記プログラムは、前述した機能の一部を実現するためのものであってもよい。さらに、前述した機能をコンピュータシステムにすでに記録されているプログラムとの組み合わせで実現できるもの、いわゆる差分ファイル（差分プログラム）であってもよい。

【 符号の説明 】

【 0 1 1 8 】

10

20

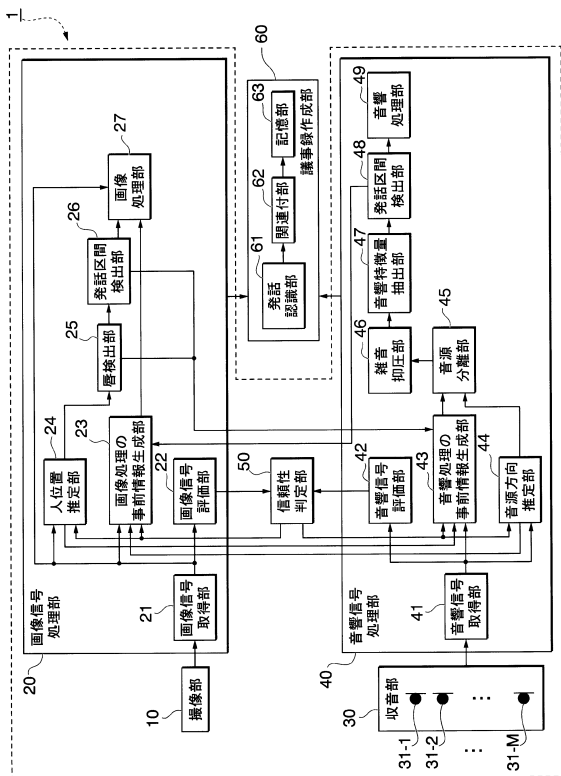
30

40

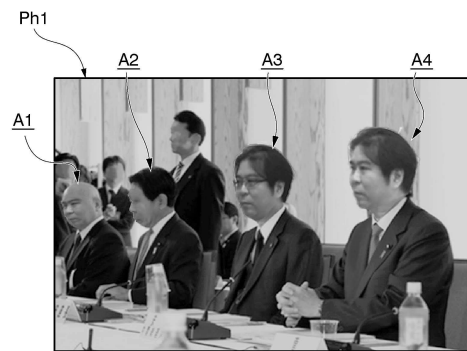
50

- 1、1 A ... 音源分離装置、10 ... 撮像部、20 ... 画像信号処理部、21 ... 画像信号取得部、22 ... 画像信号評価部、23 ... 画像信号の事前情報生成部、24 ... 人位置推定部、25 ... 唇検出部、26 ... 発話区間検出部、27 ... 画像処理部、30 ... 收音部、40 ... 音響信号処理部、41 ... 音響信号取得部、42 ... 音響信号評価部、43 ... 音響信号の事前情報生成部、44 ... 音源方向推定部、45 ... 音源分離部、46 ... 雑音抑圧部、47 ... 音響特徴量抽出部、48 ... 発話区間検出部、49 ... 音響処理部、50 ... 信頼性判定部、60 ... 議事録作成部、61 ... 発話認識部、62 ... 関連付部、63 ... 記憶部、70 ... 送信部、80 ... 受信部、90 ... ネットワーク

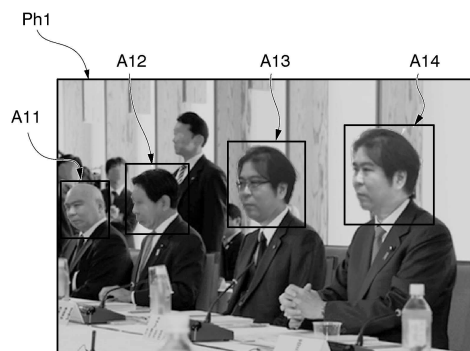
【図1】



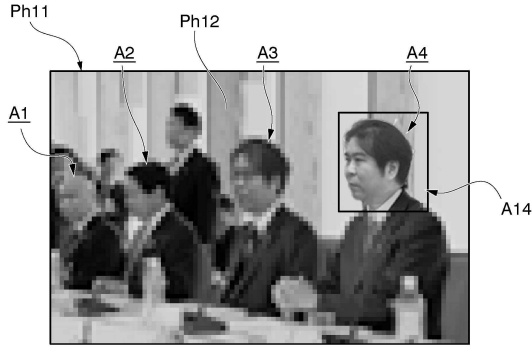
【図2】



【図3】



【 図 4 】



【 図 5 】

画像信号評価部の評価結果	音響信号評価部の評価結果	判定結果
信頼性が高い	信頼性が高い	音源定位の推定結果を優先
信頼性が高い	信頼性が低い	人位置推定の推定結果を優先
信頼性が低い	信頼性が高い	音源定位の推定結果を優先
信頼性が低い	信頼性が低い	人位置推定の推定結果を優先

【 図 6 】

発話内容	話者の顔面の鮮明度を 変更しない画像
「それでは、本日の会議を始めます。Bさん、本日の議題は何ですか？」	話者Sp1の画像
「本日の議題は、ヨーロッパの国々の消費税率についてです。」	話者Sp2の画像
「それではCさん、ヨーロッパの消費税率の現状について報告をお願いします。」	話者Sp1の画像
「日本では、消費税は食品でも玉石でも同じですが、多くの国では、一般の商品に対する消費税は、食品に対する消費税が算なっています。イギリスは、一般の商品に対する消費税率は17.5%ですが、食品に対する消費税率は0%です。フランスでは、一般の商品に対する消費税率は19.6%ですが、食品に対する消費税率は5.5%です。ドイツでは、一般の商品に対する消費税率は17%ですが、食品に対する消費税率は6%です。」	話者Sp3の画像
「食品と言ってもキャベツなど高級品に対しても、その税率が適用されるのか？」	話者Sp1の画像
「国によって異なるようです。例えば、フランスでは、キャベツに対する消費税率は19.6%ですが、ポアグラやトリュフに対する消費税率は5.6%です。」	話者Sp2の画像
...	...

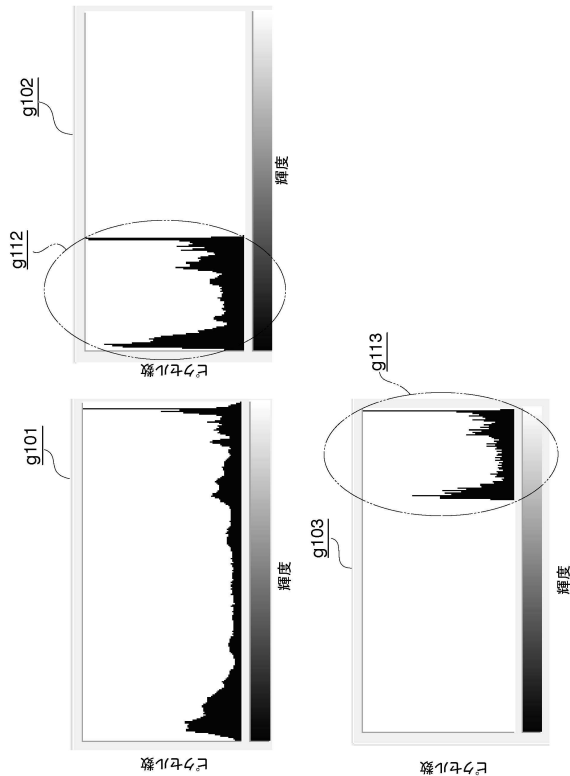
【 図 7 】



【 図 8 】



【 図 9 】



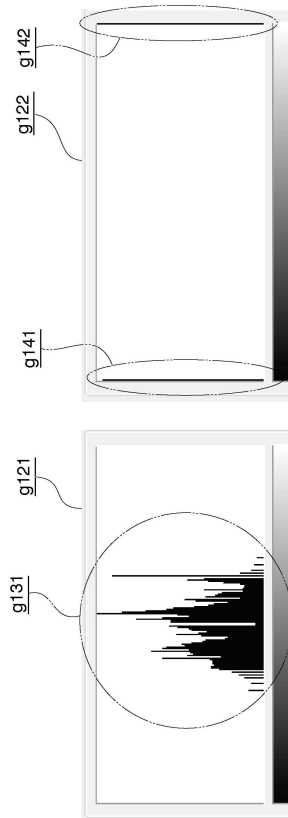
【図10】



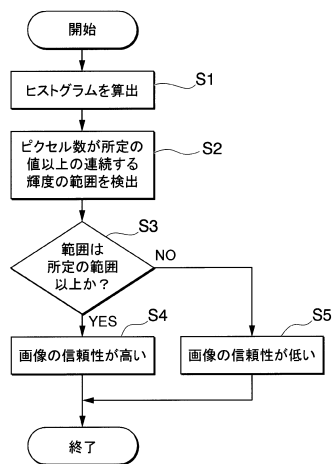
【図11】



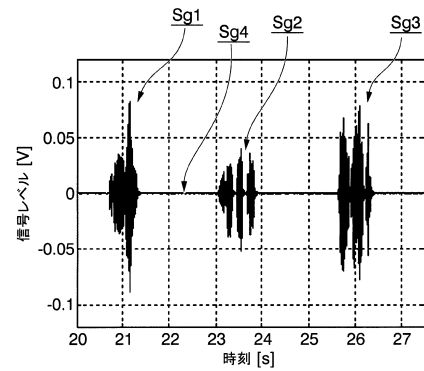
【図12】



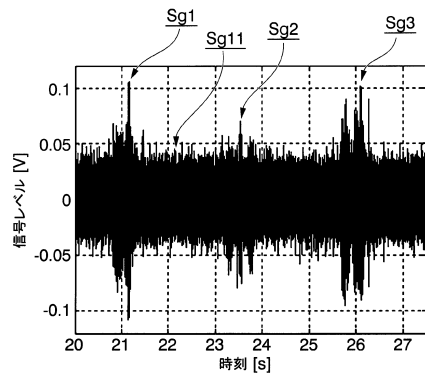
【図13】



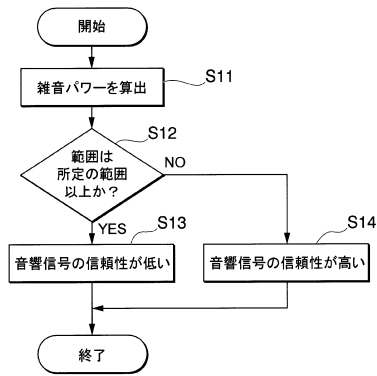
【図14】



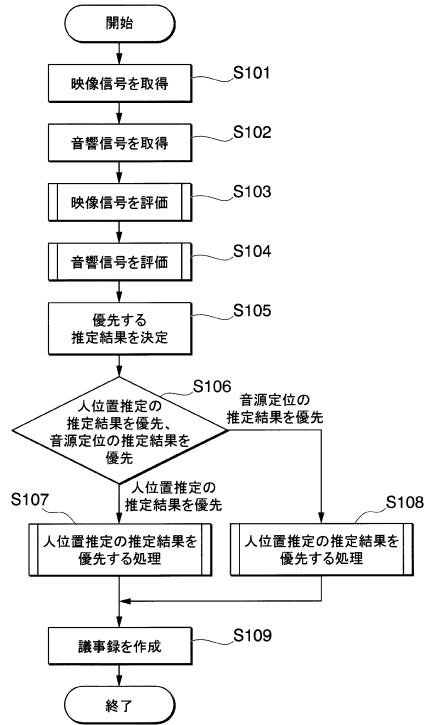
【図15】



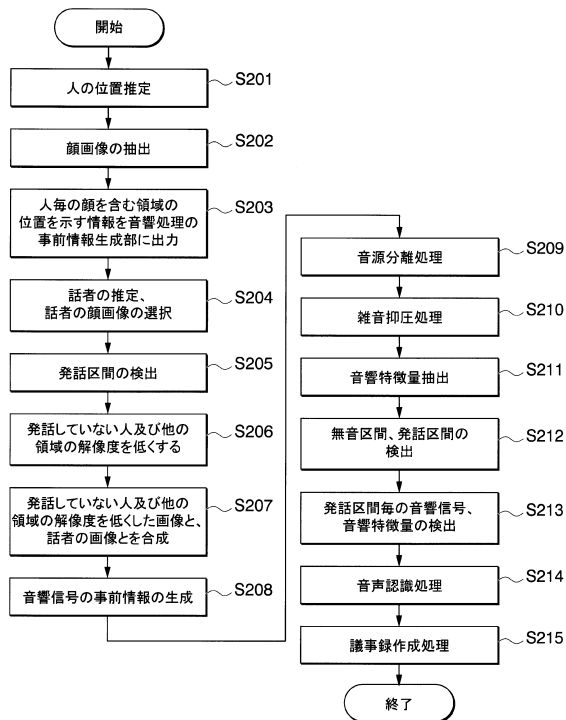
【図16】



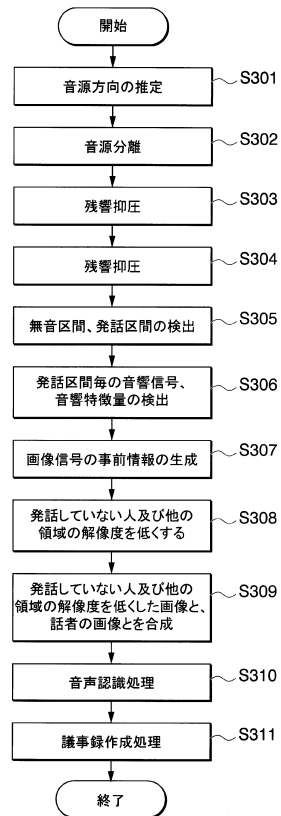
【図17】



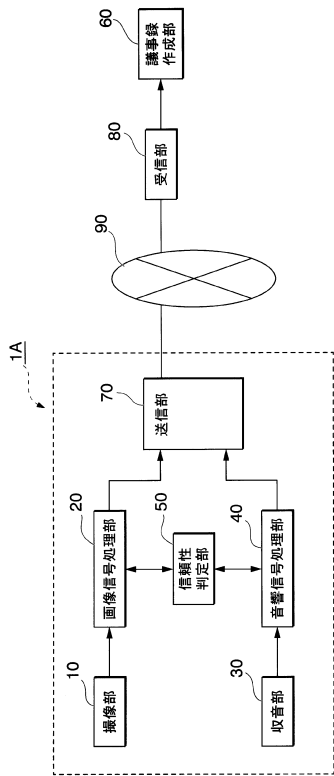
【図18】



【図19】



【 図 20 】



## フロントページの続き

(51)Int.Cl. F I  
H 0 4 R 3/00 (2006.01) H 0 4 R 3/00 3 2 0

(74)代理人 100175802

弁理士 寺本 光生

(74)代理人 100094400

弁理士 鈴木 三義

(72)発明者 水本 武志

埼玉県和光市本町 8 - 1 株式会社ホンダ・リサーチ・インスティテュート・ジャパン内

(72)発明者 中臺 一博

埼玉県和光市本町 8 - 1 株式会社ホンダ・リサーチ・インスティテュート・ジャパン内

審査官 上田 雄

(56)参考文献 特開 2 0 1 1 - 1 9 1 4 2 3 ( J P , A )

特開 2 0 0 2 - 3 1 2 7 9 6 ( J P , A )

特開 2 0 1 0 - 0 8 1 6 4 4 ( J P , A )

特開 2 0 0 8 - 1 4 1 6 9 9 ( J P , A )

特開 2 0 1 2 - 2 1 5 6 0 6 ( J P , A )

特開 2 0 0 7 - 0 3 3 9 2 0 ( J P , A )

特開 2 0 1 3 - 1 2 2 6 9 5 ( J P , A )

特開 2 0 0 6 - 1 2 3 1 6 1 ( J P , A )

特開 2 0 1 0 - 1 5 4 2 5 9 ( J P , A )

(58)調査した分野(Int.Cl. , DB名)

G 1 0 L 1 5 / 0 0 - 1 5 / 3 4

G 1 0 L 2 1 / 0 0 - 2 1 / 1 8

G 1 0 L 2 5 / 0 0 - 2 5 / 9 3

G 0 6 T 7 / 0 0