

(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2009-42910

(P2009-42910A)

(43) 公開日 平成21年2月26日(2009.2.26)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G 0 6 F</b> 3/033 (2006.01)	G 0 6 F 3/033 3 1 0 Y	3 C 0 0 7
<b>G 1 0 L</b> 17/00 (2006.01)	G 1 0 L 17/00 2 0 0 C	5 B 0 8 7
<b>B 2 5 J</b> 13/08 (2006.01)	G 1 0 L 17/00 4 0 0	5 D 0 1 5
<b>G 0 6 F</b> 3/01 (2006.01)	B 2 5 J 13/08 Z	5 E 5 0 1
<b>B 2 5 J</b> 19/02 (2006.01)	G 0 6 F 3/01	
審査請求 未請求 請求項の数 25 O L (全 50 頁) 最終頁に続く		

(21) 出願番号 特願2007-205646 (P2007-205646)  
 (22) 出願日 平成19年8月7日(2007.8.7)

(71) 出願人 000002185  
 ソニー株式会社  
 東京都港区港南1丁目7番1号  
 (74) 代理人 100093241  
 弁理士 宮田 正昭  
 (74) 代理人 100101801  
 弁理士 山田 英治  
 (74) 代理人 100086531  
 弁理士 澤田 俊夫  
 (74) 代理人 100095496  
 弁理士 佐々木 榮二  
 (72) 発明者 澤田 務  
 東京都港区港南1丁目7番1号 ソニー株式会社内

最終頁に続く

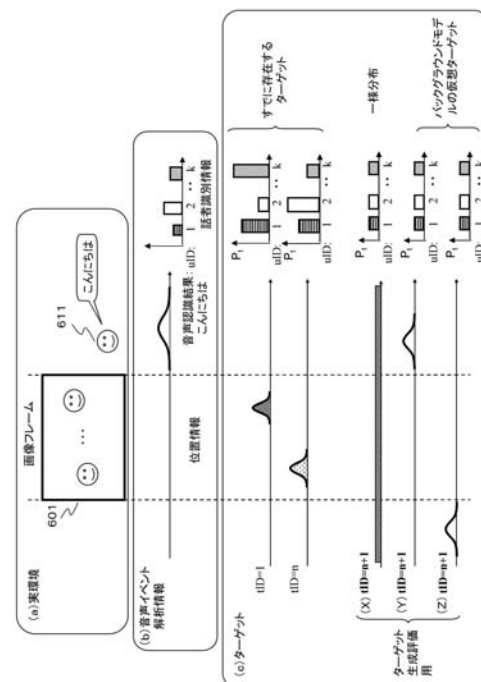
(54) 【発明の名称】 情報処理装置、および情報処理方法、並びにコンピュータ・プログラム

## (57) 【要約】

【課題】不確実で非同期な入力情報に基づく情報解析により、精度の高いユーザ位置およびユーザ識別情報を効率的に生成する構成を実現する

【解決手段】カメラやマイクによって取得される画像情報や音声情報に基づいてユーザの推定位置および推定識別データを含むイベント情報を入力して、複数ターゲットを設定した複数パーティクルを適用したパーティクルフィルタリング処理を行い仮説の更新取捨選択によりユーザ位置および識別情報を生成する。また、カメラの画像フレーム外に仮想ターゲットを設定した暫定ターゲットとイベント検出部の生成するイベント情報との尤度を検証し、検証結果に応じて暫定ターゲットを各パーティクルに追加する。本構成により、フレーム外ユーザの音声入力に対応した処理が可能となり、ユーザ位置や識別の正確な推定処理が実現される。

【選択図】図13



**【特許請求の範囲】****【請求項 1】**

実空間における画像情報または音声情報のいずれかを含む情報を入力する複数の情報入力部と、

前記情報入力部から入力する情報の解析により、前記実空間に存在するユーザの推定位置情報および推定識別情報を含むイベント情報を生成するイベント検出部と、

ユーザの位置および識別情報についての仮説 (Hypothesis) の確率分布データを設定し、前記イベント情報に基づく仮説の更新および取捨選択により、前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報の生成を実行する情報統合処理部を有し、

10

前記情報統合処理部は、

前記イベント検出部の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットと前記イベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント - ターゲット間尤度より大きい値である場合に、前記暫定ターゲットを各パーティクルに新規追加する処理を行うことを特徴とする情報処理装置。

**【請求項 2】**

20

前記情報統合処理部は、

前記画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとして、画像フレームの異なる方向のフレーム外部位置に仮想ターゲットを設定した複数の異なる暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出して、算出した暫定ターゲットのイベント - ターゲット間尤度の最大値が、既存ターゲットに対応するイベント - ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うことを特徴とする請求項 1 に記載の情報処理装置。

**【請求項 3】**

30

前記情報統合処理部は、

前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットの他、均一データによって構成されるユーザ位置情報、ユーザ識別情報を持つ暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出し、算出した暫定ターゲットのイベント - ターゲット間尤度の最大値が、既存ターゲットに対応するイベント - ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うことを特徴とする請求項 1 に記載の情報処理装置。

**【請求項 4】**

前記イベント検出部は、

ガウス分布からなるユーザの推定位置情報と、ユーザ対応の確率値を示すユーザ確信度情報を含むイベント情報を生成する構成であり、

40

前記情報統合処理部は、

仮想的なユーザに対応するガウス分布からなるユーザ位置情報と、ユーザ対応の確率値を示すユーザの確信度情報を有するターゲットを複数設定したパーティクルを保持し、各パーティクルに設定されたターゲットと、前記イベント情報との類似度の指標値であるイベント - ターゲット間尤度を算出して、イベント - ターゲット間尤度の高いターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行する構成であることを特徴とする請求項 1 ~ 3 いずれかに記載の情報処理装置。

**【請求項 5】**

前記情報統合処理部は、

50

前記イベント - ターゲット間尤度と、各パーティクルに設定したパーティクル重みとの総和データをターゲット重みとして算出し、ターゲット重みの大きいターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行する構成であることを特徴とする請求項 4 に記載の情報処理装置。

【請求項 6】

前記情報統合処理部は、

各パーティクルに設定したイベント発生源仮説ターゲットと、前記イベント検出部から入力するイベント情報との尤度を算出し、該尤度の大小に応じた値をパーティクル重みとして各パーティクルに設定する構成であることを特徴とする請求項 4 に記載の情報処理装置。

10

【請求項 7】

前記情報統合処理部は、

前記パーティクル重みの大きいパーティクルを優先的に再選択するリサンプリング処理を実行して、パーティクルの更新処理を行う構成であることを特徴とする請求項 6 に記載の情報処理装置。

【請求項 8】

前記情報統合処理部は、

各パーティクルに設定したターゲットについて、経過時間を考慮した更新処理を実行する構成であることを特徴とする請求項 1 ~ 3 いずれかに記載の情報処理装置。

【請求項 9】

前記情報統合処理部は、

各パーティクルに設定したイベント発生源仮説ターゲットについて、前記イベント検出部の生成するイベント情報を適用した更新処理を行う構成であることを特徴とする請求項 4 に記載の情報処理装置。

20

【請求項 10】

前記情報統合処理部は、

前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和を、各ターゲット対応のユーザ位置情報およびユーザ識別情報としたターゲット情報を生成する構成であることを特徴とする請求項 6 に記載の情報処理装置。

【請求項 11】

前記情報統合処理部は、

前記パーティクルの各々に設定したイベント発生源仮説ターゲットの数に応じて、イベント発生源の確率値としてのシクナル情報の生成を行う構成であることを特徴とする請求項 4 に記載の情報処理装置。

30

【請求項 12】

前記情報統合処理部は、

前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和に含まれるユーザ位置情報としてのガウス分布データのピーク値が予め設定した閾値未満である場合に、該ターゲットを削除する処理を実行する構成であることを特徴とする請求項 6 に記載の情報処理装置。

40

【請求項 13】

情報処理装置において情報解析処理を実行する情報処理方法であり、

複数の情報入力部が、実空間における画像情報または音声情報のいずれかを含む情報を入力する情報入力ステップと、

イベント検出部が、前記情報入力ステップにおいて入力する情報の解析により、前記実空間に存在するユーザの推定位置情報および推定識別情報を含むイベント情報を生成するイベント検出ステップと、

情報統合処理部が、ユーザの位置および識別情報についての仮説 (Hypothesis) の確率分布データを設定し、前記イベント情報に基づく仮説の更新および取捨選択により、前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情

50

報の生成を実行する情報統合処理ステップを有し、

前記情報統合処理ステップは、

前記イベント検出部の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットと前記イベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント - ターゲット間尤度より大きい値である場合に、前記暫定ターゲットを各パーティクルに新規追加する処理を行うステップであることを特徴とする情報処理方法。

10

【請求項 14】

前記情報統合処理ステップは、

前記画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとして、画像フレームの異なる方向のフレーム外部位置に仮想ターゲットを設定した複数の異なる暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出して、算出した暫定ターゲットのイベント - ターゲット間尤度の最大値が、既存ターゲットに対応するイベント - ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うステップであることを特徴とする請求項 13 に記載の情報処理方法。

20

【請求項 15】

前記情報統合処理ステップは、

前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットの他、均一データによって構成されるユーザ位置情報、ユーザ識別情報を持つ暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出し、算出した暫定ターゲットのイベント - ターゲット間尤度の最大値が、既存ターゲットに対応するイベント - ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うステップであることを特徴とする請求項 13 に記載の情報処理方法。

【請求項 16】

前記イベント検出ステップは、

ガウス分布からなるユーザの推定位置情報と、ユーザ対応の確率値を示すユーザ確信度情報を含むイベント情報を生成するステップであり、

30

前記情報統合処理部は、仮想的なユーザに対応するガウス分布からなるユーザ位置情報と、ユーザ対応の確率値を示すユーザの確信度情報を有するターゲットを複数設定したパーティクルを保持し、

前記情報統合処理ステップは、

各パーティクルに設定されたターゲットと、前記イベント情報との類似度の指標値であるイベント - ターゲット間尤度を算出して、イベント - ターゲット間尤度の高いターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行するステップであることを特徴とする請求項 13 ~ 15 いずれかに記載の情報処理方法。

40

【請求項 17】

前記情報統合処理ステップは、

前記イベント - ターゲット間尤度と、各パーティクルに設定したパーティクル重みとの総和データをターゲット重みとして算出し、ターゲット重みの大きいターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行するステップであることを特徴とする請求項 16 に記載の情報処理方法。

【請求項 18】

前記情報統合処理ステップは、

各パーティクルに設定したイベント発生源仮説ターゲットと、前記イベント検出部から入力するイベント情報との尤度を算出し、該尤度の大小に応じた値をパーティクル重みと

50

して各パーティクルに設定するステップであることを特徴とする請求項 16 に記載の情報処理方法。

【請求項 19】

前記情報統合処理ステップは、

前記パーティクル重みの大きいパーティクルを優先的に再選択するリサンプリング処理を実行して、パーティクルの更新処理を行うステップであることを特徴とする請求項 18 に記載の情報処理方法。

【請求項 20】

前記情報統合処理ステップは、

各パーティクルに設定したターゲットについて、経過時間を考慮した更新処理を実行するステップであることを特徴とする請求項 13 ~ 15 いずれかに記載の情報処理方法。

10

【請求項 21】

前記情報統合処理ステップは、

各パーティクルに設定したイベント発生源仮説ターゲットについて、前記イベント検出部の生成するイベント情報を適用した更新処理を行うステップであることを特徴とする請求項 16 に記載の情報処理方法。

【請求項 22】

前記情報統合処理ステップは、

前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和を、各ターゲット対応のユーザ位置情報およびユーザ識別情報としたターゲット情報を生成するステップであることを特徴とする請求項 18 に記載の情報処理方法。

20

【請求項 23】

前記情報統合処理ステップは、

前記パーティクルの各々に設定したイベント発生源仮説ターゲットの数に応じて、イベント発生源の確率値としてのシクナル情報の生成を行うステップであることを特徴とする請求項 16 に記載の情報処理方法。

【請求項 24】

前記情報統合処理ステップは、

前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和に含まれるユーザ位置情報としてのガウス分布データのピーク値が予め設定した閾値未満である場合に、該ターゲットを削除する処理を実行するステップを含むことを特徴とする請求項 18 に記載の情報処理方法。

30

【請求項 25】

情報処理装置において情報解析処理を実行させるコンピュータ・プログラムであり、

複数の情報入力部に、実空間における画像情報または音声情報のいずれかを含む情報を入力させる情報入力ステップと、

イベント検出部に、前記情報入力ステップにおいて入力する情報の解析により、前記実空間に存在するユーザの推定位置情報および推定識別情報を含むイベント情報を生成させるイベント検出ステップと、

情報統合処理部に、ユーザの位置および識別情報についての仮説 (Hypothesis) の確率分布データを設定し、前記イベント情報に基づく仮説の更新および取捨選択により、前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報の生成を実行させる情報統合処理ステップを有し、

40

前記情報統合処理ステップは、

前記イベント検出部の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットと前記イベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応

50

するイベント - ターゲット間尤度より大きい値である場合に、前記暫定ターゲットを各パ  
ーティクルに新規追加する処理を行わせるステップであることを特徴とするコンピュ  
ータ・プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置、および情報処理方法、並びにコンピュータ・プログラムに  
関する。さらに詳細には、外界からの入力情報、例えば画像、音声などの情報を入力し、入  
力情報に基づく外界環境の解析、具体的には言葉を発している人物の位置や誰であるか等  
の解析処理を実行する情報処理装置、および情報処理方法、並びにコンピュータ・プロ  
ラムに関する。

10

【背景技術】

【0002】

人とPCやロボットなどの情報処理装置との相互間の処理、例えばコミュニケーション  
やインタラクティブ処理を行うシステムはマン - マシン インタラクション システムと  
呼ばれる。このマン - マシン インタラクション システムにおいて、PCやロボット等  
の情報処理装置は、人のアクション例えば人の動作や言葉を認識するために画像情報や音  
声情報を入力して入力情報に基づく解析を行う。

【0003】

人が情報を伝達する場合、言葉のみならずしぐさ、視線、表情など様々なチャネルを情  
報伝達チャネルとして利用する。このようなすべてのチャネルの解析をマシンにおいて行  
うことができれば、人とマシンとのコミュニケーションも人と人とのコミュニケーション  
と同レベルに到達することができる。このような複数のチャネル（モダリティ、モータル  
とも呼ばれる）からの入力情報の解析を行うインタフェースは、マルチモダルインタフ  
ェースと呼ばれ、近年、開発、研究が盛んに行われている。

20

【0004】

例えばカメラによって撮影された画像情報、マイクによって取得された音声情報を入力  
して解析を行う場合、より詳細な解析を行うためには、様々なポイントに設置した複数の  
カメラおよび複数のマイクから多くの情報を入力することが有効である。

【0005】

30

具体的なシステムとしては、例えば以下のようなシステムが想定される。情報処理装置  
（テレビ）が、カメラおよびマイクを介して、テレビの前のユーザ（父、母、姉、弟）の  
画像および音声を入力し、それぞれのユーザの位置やどのユーザが発した言葉であるか等  
を解析し、テレビが解析情報に応じた処理、例えば会話をを行ったユーザに対するカメラの  
ズームアップや、会話をを行ったユーザに対する的確な応答を行うなどのシステムが実現可  
能となる。

【0006】

従来の一般的なマン - マシン インタラクション システムの多くは、複数チャネル（  
モータル）からの情報を決定論的に統合して、複数のユーザが、それぞれどこにいて、そ  
れらは誰で、誰がシグナルを発したのかを決定するという処理を行っていた。このような  
システムを開示した従来技術として、例えば特許文献1（特開2005 - 271137号  
公報）、特許文献2（特開2002 - 264051号公報）がある。

40

【0007】

しかし、従来システムにおいて行われるマイクやカメラから入力される不確実かつ非  
同期なデータを利用した決定論的な統合処理方法ではロバスト性にかけ、精度の低いデー  
タしか得られないという問題がある。実際のシステムにおいて、実環境で取得可能なセン  
サ情報、すなわちカメラからの入力画像やマイクから入力される音声情報には様々な余分  
な情報、例えばノイズや不要な情報が含まれる不確実なデータであり、画像解析や音声解  
析処理を行う場合には、このようなセンサ情報から有効な情報を効率的に統合する処理が  
重要となる。

50

【特許文献１】特開２００５－２７１１３７号公報

【特許文献２】特開２００２－２６４０５１号公報

【発明の開示】

【発明が解決しようとする課題】

【０００８】

本発明は、上述の問題点に鑑みてなされたものであり、複数のチャネル（モダリティ、モダル）からの入力情報の解析、具体的には、例えば周囲にいる人物の位置などの特定処理を行うシステムにおいて、画像、音声情報などの様々な入力情報に含まれる不確実な情報に対する確率的な処理を行ってより精度の高いと推定される情報に統合する処理を行うことによりロバスト性を向上させ、精度の高い解析を行う情報処理装置、および情報処理方法、並びにコンピュータ・プログラムを提供することを目的とする。

10

【課題を解決するための手段】

【０００９】

本発明の第１の側面は、

実空間における画像情報または音声情報のいずれかを含む情報を入力する複数の情報入力部と、

前記情報入力部から入力する情報の解析により、前記実空間に存在するユーザの推定位置情報および推定識別情報を含むイベント情報を生成するイベント検出部と、

ユーザの位置および識別情報についての仮説（Hypothesis）の確率分布データを設定し、前記イベント情報に基づく仮説の更新および取捨選択により、前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報の生成を実行する情報統合処理部を有し、

20

前記情報統合処理部は、

前記イベント検出部の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットと前記イベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント・ターゲット間尤度より大きい値である場合に、前記暫定ターゲットを各パーティクルに新規追加する処理を行うことを特徴とする情報処理装置にある。

30

【００１０】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとして、画像フレームの異なる方向のフレーム外部位置に仮想ターゲットを設定した複数の異なる暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出して、算出した暫定ターゲットのイベント・ターゲット間尤度の最大値が、既存ターゲットに対応するイベント・ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うことを特徴とする。

【００１１】

40

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットの他、均一データによって構成されるユーザ位置情報、ユーザ識別情報を持つ暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出し、算出した暫定ターゲットのイベント・ターゲット間尤度の最大値が、既存ターゲットに対応するイベント・ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うことを特徴とする。

【００１２】

さらに、本発明の情報処理装置の一実施態様において、前記イベント検出部は、ガウス

50

分布からなるユーザの推定位置情報と、ユーザ対応の確率値を示すユーザ確信度情報を含むイベント情報を生成する構成であり、前記情報統合処理部は、仮想的なユーザに対応するガウス分布からなるユーザ位置情報と、ユーザ対応の確率値を示すユーザの確信度情報を有するターゲットを複数設定したパーティクルを保持し、各パーティクルに設定されたターゲットと、前記イベント情報との類似度の指標値であるイベント・ターゲット間尤度を算出して、イベント・ターゲット間尤度の高いターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行する構成であることを特徴とする。

【0013】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記イベント・ターゲット間尤度と、各パーティクルに設定したパーティクル重みとの総和データをターゲット重みとして算出し、ターゲット重みの大きいターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行する構成であることを特徴とする。

10

【0014】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、各パーティクルに設定したイベント発生源仮説ターゲットと、前記イベント検出部から入力するイベント情報との尤度を算出し、該尤度の大小に応じた値をパーティクル重みとして各パーティクルに設定する構成であることを特徴とする。

【0015】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記パーティクル重みの大きいパーティクルを優先的に再選択するリサンプリング処理を実行して、パーティクルの更新処理を行う構成であることを特徴とする。

20

【0016】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、各パーティクルに設定したターゲットについて、経過時間を考慮した更新処理を実行する構成であることを特徴とする。

【0017】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、各パーティクルに設定したイベント発生源仮説ターゲットについて、前記イベント検出部の生成するイベント情報を適用した更新処理を行う構成であることを特徴とする。

30

【0018】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和を、各ターゲット対応のユーザ位置情報およびユーザ識別情報としたターゲット情報を生成する構成であることを特徴とする。

【0019】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記パーティクルの各々に設定したイベント発生源仮説ターゲットの数に応じて、イベント発生源の確率値としてのシクナル情報の生成を行う構成であることを特徴とする。

【0020】

さらに、本発明の情報処理装置の一実施態様において、前記情報統合処理部は、前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和に含まれるユーザ位置情報としてのガウス分布データのピーク値が予め設定した閾値未満である場合に、該ターゲットを削除する処理を実行する構成であることを特徴とする。

40

【0021】

さらに、本発明の第2の側面は、  
情報処理装置において情報解析処理を実行する情報処理方法であり、  
複数の情報入力部が、実空間における画像情報または音声情報のいずれかを含む情報を  
入力する情報入力ステップと、  
イベント検出部が、前記情報入力ステップにおいて入力する情報の解析により、前記実

50



空間に存在するユーザの推定位置情報および推定識別情報を含むイベント情報を生成するイベント検出ステップと、

情報統合処理部が、ユーザの位置および識別情報についての仮説 (Hypothesis) の確率分布データを設定し、前記イベント情報に基づく仮説の更新および取捨選択により、前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報の生成を実行する情報統合処理ステップを有し、

前記情報統合処理ステップは、

前記イベント検出部の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットと前記イベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント - ターゲット間尤度より大きい値である場合に、前記暫定ターゲットを各パーティクルに新規追加する処理を行うステップであることを特徴とする情報処理方法にある。

10

#### 【0022】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとして、画像フレームの異なる方向のフレーム外部位置に仮想ターゲットを設定した複数の異なる暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出して、算出した暫定ターゲットのイベント - ターゲット間尤度の最大値が、既存ターゲットに対応するイベント - ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うステップであることを特徴とする。

20

#### 【0023】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットの他、均一データによって構成されるユーザ位置情報、ユーザ識別情報を持つ暫定ターゲットを生成し、生成した複数の暫定ターゲットと前記イベント情報との尤度を個別に算出し、算出した暫定ターゲットのイベント - ターゲット間尤度の最大値が、既存ターゲットに対応するイベント - ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行うステップであることを特徴とする。

30

#### 【0024】

さらに、本発明の情報処理方法の一実施態様において、前記イベント検出ステップは、ガウス分布からなるユーザの推定位置情報と、ユーザ対応の確率値を示すユーザ確信度情報を含むイベント情報を生成するステップであり、前記情報統合処理部は、仮想的なユーザに対応するガウス分布からなるユーザ位置情報と、ユーザ対応の確率値を示すユーザの確信度情報を有するターゲットを複数設定したパーティクルを保持し、前記情報統合処理ステップは、各パーティクルに設定されたターゲットと、前記イベント情報との類似度の指標値であるイベント - ターゲット間尤度を算出して、イベント - ターゲット間尤度の高いターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行するステップであることを特徴とする。

40

#### 【0025】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記イベント - ターゲット間尤度と、各パーティクルに設定したパーティクル重みとの総和データをターゲット重みとして算出し、ターゲット重みの大きいターゲットを優先的にイベント発生源仮説ターゲットとしたパーティクル設定処理を実行するステップであることを特徴とする。

50

## 【 0 0 2 6 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、各パーティクルに設定したイベント発生源仮説ターゲットと、前記イベント検出部から入力するイベント情報との尤度を算出し、該尤度の大小に応じた値をパーティクル重みとして各パーティクルに設定するステップであることを特徴とする。

## 【 0 0 2 7 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記パーティクル重みの大きいパーティクルを優先的に再選択するリサンプリング処理を実行して、パーティクルの更新処理を行うステップであることを特徴とする。

## 【 0 0 2 8 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、各パーティクルに設定したターゲットについて、経過時間を考慮した更新処理を実行するステップであることを特徴とする。

## 【 0 0 2 9 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、各パーティクルに設定したイベント発生源仮説ターゲットについて、前記イベント検出部の生成するイベント情報を適用した更新処理を行うステップであることを特徴とする。

## 【 0 0 3 0 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和を、各ターゲット対応のユーザ位置情報およびユーザ識別情報としたターゲット情報を生成するステップであることを特徴とする。

## 【 0 0 3 1 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記パーティクルの各々に設定したイベント発生源仮説ターゲットの数に応じて、イベント発生源の確率値としてのシクナル情報の生成を行うステップであることを特徴とする。

## 【 0 0 3 2 】

さらに、本発明の情報処理方法の一実施態様において、前記情報統合処理ステップは、前記パーティクルの各々に設定したターゲットデータと前記パーティクル重みとの積算総和に含まれるユーザ位置情報としてのガウス分布データのピーク値が予め設定した閾値未満である場合に、該ターゲットを削除する処理を実行するステップを含むことを特徴とする。

## 【 0 0 3 3 】

さらに、本発明の第3の側面は、

情報処理装置において情報解析処理を実行させるコンピュータ・プログラムであり、

複数の情報入力部に、実空間における画像情報または音声情報のいずれかを含む情報を入力させる情報入力ステップと、

イベント検出部に、前記情報入力ステップにおいて入力する情報の解析により、前記実空間に存在するユーザの推定位置情報および推定識別情報を含むイベント情報を生成させるイベント検出ステップと、

情報統合処理部に、ユーザの位置および識別情報についての仮説 (Hypothesis) の確率分布データを設定し、前記イベント情報に基づく仮説の更新および取捨選択により、前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報の生成を実行させる情報統合処理ステップを有し、

前記情報統合処理ステップは、

前記イベント検出部の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して前記実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、前記情報入力部を構成するカメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットと前記イベント検出部の生成するイベ

10

20

30

40

50

ント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント・ターゲット間尤度より大きい値である場合に、前記暫定ターゲットを各パーティクルに新規追加する処理を行わせるステップであることを特徴とするコンピュータ・プログラムにある。

【0034】

なお、本発明のコンピュータ・プログラムは、例えば、様々なプログラム・コードを実行可能な汎用コンピュータ・システムに対して、コンピュータ可読な形式で提供する記憶媒体、通信媒体によって提供可能なコンピュータ・プログラムである。このようなプログラムをコンピュータ可読な形式で提供することにより、コンピュータ・システム上でプログラムに応じた処理が実現される。

【0035】

本発明のさらに他の目的、特徴や利点は、後述する本発明の実施例や添付する図面に基づくより詳細な説明によって明らかになるであろう。なお、本明細書においてシステムとは、複数の装置の論理的集合構成であり、各構成の装置が同一筐体内にあるものには限らない。

【発明の効果】

【0036】

本発明の一実施例の構成によれば、カメラやマイクによって取得される画像情報や音声情報に基づいてユーザの推定位置および推定識別データを含むイベント情報を入力して、複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を行い、フィルタリングによる仮説の更新および取捨選択に基づいてユーザの位置および識別情報を生成する。また、カメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとイベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント・ターゲット間尤度より大きい値である場合に、暫定ターゲットを各パーティクルに新規追加する処理を行う。本構成により、カメラの取得する画像フレームの外部にいるユーザからの音声入力イベントに対応した正しい推定処理が可能となり、ユーザ位置やユーザ識別情報を効率的に確実に生成することが可能となる。

【発明を実施するための最良の形態】

【0037】

以下、図面を参照しながら本発明の実施形態に係る情報処理装置、および情報処理方法、並びにコンピュータ・プログラムの詳細について説明する。

【0038】

まず、図1を参照して本発明に係る情報処理装置の実行する処理の概要について説明する。本発明の情報処理装置100は、環境情報を入力するセンサ、ここでは一例としてカメラ21と、複数のマイク31～34から画像情報、音声情報を入力し、これらの入力情報に基づいて環境の解析を行う。具体的には、複数のユーザ1, 11～4, 14の位置の解析、およびその位置にいるユーザの識別を行う。

【0039】

図に示す例において、例えばユーザ1, 11～ユーザ4, 14が家族である父、母、姉、弟であるとき、情報処理装置100は、カメラ21と、複数のマイク31～34から入力する画像情報、音声情報の解析を行い、4人のユーザ1～4の存在する位置、各位置にいるユーザが父、母、姉、弟のいずれであるかを識別する。識別処理結果は様々な処理に利用される。例えば、例えば会話をを行ったユーザに対するカメラのズームアップや、会話をを行ったユーザに対してテレビから応答を行うなどの処理に利用される。

【0040】

なお、本発明に係る情報処理装置100の主要な処理は、複数の情報入力部（カメラ21, マイク31～34）からの入力情報に基づいて、ユーザの位置識別およびユーザの特定処理としてのユーザ識別処理を行うことである。この識別結果の利用処理については特に限定するものではない。カメラ21と、複数のマイク31～34から入力する画像情報

10

20

30

40

50

、音声情報には様々な不確実な情報が含まれる。本発明の情報処理装置 100 では、これらの入力情報に含まれる不確実な情報に対する確率的な処理を行って、精度の高いと推定される情報に統合する処理を行う。この推定処理によりロバスト性を向上させ、精度の高い解析を行う。

#### 【0041】

図2に情報処理装置100の構成例を示す。情報処理装置100は、入力デバイスとして画像入力部(カメラ)111、複数の音声入力部(マイク)121a~dを有する。画像入力部(カメラ)111から画像情報を入力し、音声入力部(マイク)121から音声情報を入力し、これらの入力情報に基づいて解析を行う。複数の音声入力部(マイク)121a~dの各々は、図1に示すように様々な位置に配置されている。

10

#### 【0042】

複数のマイク121a~dから入力された音声情報は、音声イベント検出部122を介して音声・画像統合処理部131に入力される。音声イベント検出部122は、複数の異なるポジションに配置された複数の音声入力部(マイク)121a~dから入力する音声情報を解析し統合する。具体的には、音声入力部(マイク)121a~dから入力する音声情報に基づいて、発生した音の位置およびどのユーザの発生させた音であるかのユーザ識別情報を生成して音声・画像統合処理部131に入力する。

#### 【0043】

なお、情報処理装置100の実行する具体的な処理は、例えば図1に示すように複数のユーザが存在する環境で、ユーザA~Dがどの位置にいて、会話をを行ったユーザがどのユーザであるかを識別すること、すなわち、ユーザ位置およびユーザ識別を行うことであり、さらに声を発した人物などのイベント発生源を特定する処理である。

20

#### 【0044】

音声イベント検出部122は、複数の異なるポジションに配置された複数の音声入力部(マイク)121a~dから入力する音声情報を解析し、音声の発生源の位置情報を確率分布データとして生成する。具体的には、音源方向に関する期待値と分散データ $N(m_e, \sigma_e)$ を生成する。また、予め登録されたユーザの声の特徴情報との比較処理に基づいてユーザ識別情報を生成する。この識別情報も確率的な推定値として生成する。音声イベント検出部122には、予め検証すべき複数のユーザの声についての特徴情報が登録されており、入力音声と登録音声との比較処理を実行して、どのユーザの声である確率が高いかを判定する処理を行い、全登録ユーザに対する事後確率、あるいはスコアを算出する。

30

#### 【0045】

このように、音声イベント検出部122は、複数の異なるポジションに配置された複数の音声入力部(マイク)121a~dから入力する音声情報を解析し、音声の発生源の位置情報を確率分布データと、確率的な推定値からなるユーザ識別情報とによって構成される[統合音声イベント情報]を生成して音声・画像統合処理部131に入力する。

#### 【0046】

一方、画像入力部(カメラ)111から入力された画像情報は、画像イベント検出部112を介して音声・画像統合処理部131に入力される。画像イベント検出部112は、画像入力部(カメラ)111から入力する画像情報を解析し、画像に含まれる人物の顔を抽出し、顔の位置情報を確率分布データとして生成する。具体的には、顔の位置や方向に関する期待値と分散データ $N(m_e, \sigma_e)$ を生成する。また、予め登録されたユーザの顔の特徴情報との比較処理に基づいてユーザ識別情報を生成する。この識別情報も確率的な推定値として生成する。画像イベント検出部112には、予め検証すべき複数のユーザの顔についての特徴情報が登録されており、入力画像から抽出した顔領域の画像の特徴情報と登録された顔画像の特徴情報との比較処理を実行して、どのユーザの顔である確率が高いかを判定する処理を行い、全登録ユーザに対する事後確率、あるいはスコアを算出する。

40

#### 【0047】

なお、音声イベント検出部122や画像イベント検出部112において実行する音声識

50

別や、顔検出、顔識別処理は従来から知られる技術を適用する。例えば顔検出、顔識別処理としては以下の文献に開示された技術の適用が可能である。

佐部 浩太郎，日台 健一，"ピクセル差分特徴を用いた実時間任意姿勢顔検出器の学習"，第10回画像センシングシンポジウム講演論文集，pp. 547 - 552，2004

特開2004-302644 (P2004-302644A) [発明の名称：顔識別装置、顔識別方法、記録媒体、及びロボット装置]

#### 【0048】

音声・画像統合処理部131は、音声イベント検出部122や画像イベント検出部112からの入力情報に基づいて、複数のユーザが、それぞれどこにいて、それらは誰で、誰が音声等のシグナルを発したのかを確率的に推定する処理を実行する。この処理については後段で詳細に説明する。音声・画像統合処理部131は、音声・画像統合処理部131は、音声イベント検出部122や画像イベント検出部112からの入力情報に基づいて、

(a) 複数のユーザが、それぞれどこにいて、それらは誰であるかの推定情報としての[ターゲット情報]

(b) 例えば話しをしたユーザなどのイベント発生源を[シグナル情報]として、処理決定部132に出力する。

#### 【0049】

これらの識別処理結果を受領した処理決定部132は、識別処理結果を利用した処理を実行する、例えば、例えば会話をを行ったユーザに対するカメラのズームアップや、会話をを行ったユーザに対してテレビから応答を行うなどの処理を行う。

#### 【0050】

上述したように、音声イベント検出部122は、音声の発生源の位置情報を確率分布データ、具体的には、音源方向に関する期待値と分散データ $N(m_e, \sigma_e)$ を生成する。また、予め登録されたユーザの声の特徴情報との比較処理に基づいてユーザ識別情報を生成して音声・画像統合処理部131に入力する。また、画像イベント検出部112は、画像に含まれる人物の顔を抽出し、顔の位置情報を確率分布データとして生成する。具体的には、顔の位置や方向に関する期待値と分散データ $N(m_e, \sigma_e)$ を生成する。また、予め登録されたユーザの顔の特徴情報との比較処理に基づいてユーザ識別情報を生成して音声・画像統合処理部131に入力する。

#### 【0051】

図3を参照して、音声イベント検出部122および画像イベント検出部112が生成し音声・画像統合処理部131に入力する情報の例について説明する。図3(A)は図1を参照して説明したと同様のカメラやマイクが備えられた実環境の例を示し、複数のユーザ1~k，201~20kが存在する。この環境で、あるユーザが話しをしたとすると、マイクで音声が入力される。また、カメラは連続的に画像を撮影している。

#### 【0052】

音声イベント検出部122および画像イベント検出部112が生成し音声・画像統合処理部131に入力する情報は、基本的に同様の情報であり、図3(B)に示す2つの情報によって構成される。すなわち、

(a) ユーザ位置情報

(b) ユーザ識別情報(顔識別情報または話者識別情報)

これらの2つの情報である。これらの2つの情報は、イベントの発生毎に生成される。音声イベント検出部122は、音声入力部(マイク)121a~dから音声情報が入力された場合に、その音声情報に基づいて上記の(a)ユーザ位置情報、(b)ユーザ識別情報を生成して音声・画像統合処理部131に入力する。画像イベント検出部112は、例えば予め定めた一定のフレーム間隔で、画像入力部(カメラ)111から入力された画像情報に基づいて(a)ユーザ位置情報、(b)ユーザ識別情報を生成して音声・画像統合処理部131に入力する。なお、本例では、画像入力部(カメラ)111は1台のカメラを設定した例を示しており、1つのカメラに複数のユーザの画像が撮影される設定であり

、この場合、１つの画像に含まれる複数の顔の各々について（ａ）ユーザ位置情報、（ｂ）ユーザ識別情報を生成して音声・画像統合処理部１３１に入力する。

【００５３】

音声イベント検出部１２２が音声入力部（マイク）１２１ａ～ｄから入力する音声情報に基づいて、

（ａ）ユーザ位置情報

（ｂ）ユーザ識別情報（話者識別情報）

これらの情報を生成する処理について説明する。

【００５４】

音声イベント検出部１２２による（ａ）ユーザ位置情報の生成処理

10

音声イベント検出部１２２は、音声入力部（マイク）１２１ａ～ｄから入力された音声情報に基づいて解析された声を発したユーザ、すなわち〔話者〕の位置の推定情報を生成する。すなわち、話者が存在すると推定される位置を、期待値（平均）〔 $m_e$ 〕と分散情報〔 $\sigma_e$ 〕からなるガウス分布（正規分布）データ $N(m_e, \sigma_e)$ として生成する。

【００５５】

音声イベント検出部１２２による（ｂ）ユーザ識別情報（話者識別情報）の生成処理

音声イベント検出部１２２は、音声入力部（マイク）１２１ａ～ｄから入力された音声情報に基づいて話者が誰であるかを、入力音声と予め登録されたユーザ１～ｋの声の特徴情報との比較処理により推定する。具体的には話者が各ユーザ１～ｋである確率を算出する。この算出値を（ｂ）ユーザ識別情報（話者識別情報）とする。例えば入力音声の特徴と最も近い登録された音声特徴を有するユーザに最も高いスコアを配分し、最も異なる特徴を持つユーザに最低のスコア（例えば０）を配分する処理によって各ユーザである確率を設定したデータを生成して、これを（ｂ）ユーザ識別情報（話者識別情報）とする。

20

【００５６】

画像イベント検出部１１２が画像入力部（カメラ）１１１から入力する画像情報に基づいて、

（ａ）ユーザ位置情報

（ｂ）ユーザ識別情報（顔識別情報）

これらの情報を生成する処理について説明する。

【００５７】

30

画像イベント検出部１１２による（ａ）ユーザ位置情報の生成処理

画像イベント検出部１１２は、画像入力部（カメラ）１１１から入力された画像情報に含まれる顔の各々について顔の位置の推定情報を生成する。すなわち、画像から検出された顔が存在すると推定される位置を、期待値（平均）〔 $m_e$ 〕と分散情報〔 $\sigma_e$ 〕からなるガウス分布（正規分布）データ $N(m_e, \sigma_e)$ として生成する。

【００５８】

画像イベント検出部１１２による（ｂ）ユーザ識別情報（顔識別情報）の生成処理

画像イベント検出部１１２は、画像入力部（カメラ）１１１から入力された画像情報に基づいて、画像情報に含まれる顔を検出し、各顔が誰であるかを、入力画像情報と予め登録されたユーザ１～ｋの顔の特徴情報との比較処理により推定する。具体的には抽出された各顔が各ユーザ１～ｋである確率を算出する。この算出値を（ｂ）ユーザ識別情報（顔識別情報）とする。例えば入力画像に含まれる顔の特徴と最も近い登録された顔の特徴を有するユーザに最も高いスコアを配分し、最も異なる特徴を持つユーザに最低のスコア（例えば０）を配分する処理によって各ユーザである確率を設定したデータを生成して、これを（ｂ）ユーザ識別情報（顔識別情報）とする。

40

【００５９】

なお、カメラの撮影画像から複数の顔が検出された場合には、各検出顔に応じて、

（ａ）ユーザ位置情報

（ｂ）ユーザ識別情報（顔識別情報）

これらの情報を生成して、音声・画像統合処理部１３１に入力する。

50

また、本例では、画像入力部 1 1 1 として 1 台のカメラを利用した例を説明するが、複数のカメラの撮影画像を利用してもよく、その場合は、画像イベント検出部 1 1 2 は、各カメラの撮影画像の各々に含まれる各顔について、

( a ) ユーザ位置情報

( b ) ユーザ識別情報 ( 顔識別情報 )

これらの情報を生成して、音声・画像統合処理部 1 3 1 に入力する。

#### 【 0 0 6 0 】

次に、音声・画像統合処理部 1 3 1 の実行する処理について説明する。音声・画像統合処理部 1 3 1 は、上述したように、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から、図 3 ( B ) に示す 2 つの情報、すなわち、

( a ) ユーザ位置情報

( b ) ユーザ識別情報 ( 顔識別情報または話者識別情報 )

これらの情報を逐次入力する。なお、これらの各情報の入力タイミングは様々な設定が可能であるが、例えば、音声イベント検出部 1 2 2 は新たな音声が入力された場合に上記 ( a ) , ( b ) の各情報を音声イベント情報として生成して入力し、画像イベント検出部 1 1 2 は、一定のフレーム周期単位で、上記 ( a ) , ( b ) の各情報を音声イベント情報として生成して入力するといった設定が可能である。

#### 【 0 0 6 1 】

音声・画像統合処理部 1 3 1 の実行する処理について、図 4 以下を参照して説明する。音声・画像統合処理部 1 3 1 は、ユーザの位置および識別情報についての仮説 ( Hypothesis ) の確率分布データを設定し、その仮説を入力情報に基づいて更新することで、より確からしい仮説のみを残す処理を行う。この処理手法として、パーティクル・フィルタ ( Particle Filter ) を適用した処理を実行する。

#### 【 0 0 6 2 】

パーティクル・フィルタ ( Particle Filter ) を適用した処理は、様々な仮説、本例では、ユーザの位置と誰であるかの仮説に対応するパーティクルを多数設定し、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から、図 3 ( B ) に示す 2 つの情報、すなわち、

( a ) ユーザ位置情報

( b ) ユーザ識別情報 ( 顔識別情報または話者識別情報 )

これらの入力情報に基づいて、より確からしいパーティクルのウェイトを高めていくという処理を行う。

#### 【 0 0 6 3 】

パーティクル・フィルタ ( Particle Filter ) を適用した基本的な処理例について図 4 を参照して説明する。例えば、図 4 に示す例は、あるユーザに対応する存在位置をパーティクル・フィルタにより推定する処理例を示している。図 4 に示す例は、ある直線上の 1 次元領域におけるユーザ 3 0 1 の存在する位置を推定する処理である。

#### 【 0 0 6 4 】

初期的な仮説 ( H ) は、図 4 ( a ) に示すように均一なパーティクル分布データとなる。次に、画像データ 3 0 2 が取得され、取得画像に基づくユーザ 3 0 1 の存在確率分布データが図 4 ( b ) のデータとして取得される。この取得画像に基づく確率分布データに基づいて、図 4 ( a ) のパーティクル分布データが更新され、図 4 ( c ) の更新された仮説確率分布データが得られる。このような処理を、入力情報に基づいて繰り返し実行して、ユーザのより確からしい位置情報を得る。

#### 【 0 0 6 5 】

なお、パーティクル・フィルタを用いた処理の詳細については、例えば [ D . Sch ulz , D . Fox , and J . Hightower . People Tracking with Anonymous and ID - sensors Using Rao - Blackwellised Particle Filters . Proc . of the International Joint Conference

10

20

30

40

50

e on Artificial Intelligence (IJCAI-03)]  
に記載されている。

【0066】

図4に示す処理例は、ユーザの存在位置のみについて、入力情報を画像データのみとした処理例として説明しており、パーティクルの各々は、ユーザ301の存在位置のみの情報を有している。

【0067】

一方、本発明に従った処理は、音声イベント検出部122および画像イベント検出部112から、図3(B)に示す2つの情報、すなわち、

(a) ユーザ位置情報

(b) ユーザ識別情報(顔識別情報または話者識別情報)

これらの入力情報に基づいて、複数のユーザの位置と複数のユーザがそれぞれ誰であるかを判別する処理を行うことになる。従って、本発明におけるパーティクル・フィルタ(Particle Filter)を適用した処理では、音声・画像統合処理部131が、ユーザの位置と誰であるかの仮説に対応するパーティクルを多数設定して、音声イベント検出部122および画像イベント検出部112から、図3(B)に示す2つの情報に基づいて、パーティクル更新を行うことになる。

【0068】

図5を参照して、本処理例で設定するパーティクルの構成について説明する。音声・画像統合処理部131は、予め設定した数 $=m$ のパーティクルを有する。図5に示すパーティクル1~ $m$ である。各パーティクルには識別子としてのパーティクルID( $PID=1\sim m$ )が設定されている。

【0069】

各パーティクルに、位置および識別を行うオブジェクトに対応する仮想的なオブジェクトに対応する複数のターゲットを設定する。本例では、例えば実空間に存在すると推定される人数以上の仮想のユーザに対応する複数のターゲットを各パーティクルに設定する。 $m$ 個のパーティクルの各々はターゲット単位でデータをターゲット数分保持する。図5に示す例では、1つのパーティクルに $n$ 個のターゲットが含まれる。各パーティクルに含まれるターゲット各々が有するターゲットデータの構成を図6に示す。

【0070】

各パーティクルに含まれる各ターゲットデータについて図6を参照して説明する。図6は、図5に示すパーティクル1( $pid=1$ )に含まれる1つのターゲット(ターゲットID: $tID=n$ )311のターゲットデータの構成である。ターゲット311のターゲットデータは、図6に示すように、以下のデータ、すなわち、

(a) 各ターゲット各々に対応する存在位置の確率分布[ガウス分布: $N(m_{1n}, \sigma_{1n})$ ]、

(b) 各ターゲットが誰であることを示すユーザ確信度情報( $uid$ )

$uid_{1n1} = 0.0$

$uid_{1n2} = 0.1$

:

$uid_{1nk} = 0.5$

これらのデータによって構成される。

【0071】

なお、(a)に示すガウス分布: $N(m_{1n}, \sigma_{1n})$ における $[m_{1n}, \sigma_{1n}]$ の $(1n)$ は、パーティクルID: $pid=1$ におけるターゲットID: $tID=n$ に対応する存在確率分布としてのガウス分布であることを意味する。

また、(b)に示すユーザ確信度情報( $uid$ )における、 $[uid_{1n1}]$ に含まれる $(1n1)$ は、パーティクルID: $pid=1$ におけるターゲットID: $tID=n$ の、ユーザ=ユーザ1である確率を意味する。すなわちターゲットID= $n$ のデータは、ユーザ1である確率が0.0、

10

20

30

40

50



ユーザ 2 である確率が 0.1、  
:

ユーザ k である確率が 0.5、  
であることを意味している。

#### 【0072】

図 5 に戻り、音声・画像統合処理部 131 の設定するパーティクルについての説明を続ける。図 5 に示すように、音声・画像統合処理部 131 は、予め決定した数 = m のパーティクル (PID = 1 ~ m) を設定し、各パーティクルは、実空間に存在すると推定されるターゲット (tID = 1 ~ n) 各々について、

(a) 各ターゲット各々に対応する存在位置の確率分布 [ ガウス分布:  $N(m, \quad)$  ]

10

(b) 各ターゲットが誰であることを示すユーザ確信度情報 (uID)

これらのターゲットデータを有する。

#### 【0073】

音声・画像統合処理部 131 は、音声イベント検出部 122 および画像イベント検出部 112 から、図 3 (B) に示すイベント情報、すなわち、

(a) ユーザ位置情報

(b) ユーザ識別情報 (顔識別情報または話者識別情報)

これらのイベント情報を入力して m 個のパーティクル (PID = 1 ~ m) の更新処理を行う。

20

#### 【0074】

音声・画像統合処理部 131、これらの更新処理を実行して、

(a) 複数のユーザが、それぞれどこにいて、それらは誰であるかの推定情報としての [ ターゲット情報 ]、

(b) 例えば話をしたユーザなどのイベント発生源を示す [ シグナル情報 ]、

これらを生成して処理決定部 132 に出力する。

#### 【0075】

[ ターゲット情報 ] は、図 5 の右端のターゲット情報 305 に示すように、各パーティクル (PID = 1 ~ m) に含まれる各ターゲット (tID = 1 ~ n) 対応データの重み付き総和データとして生成される。各パーティクルの重みについては後述する。

30

#### 【0076】

ターゲット情報 305 は、音声・画像統合処理部 131 が予め設定した仮想的なユーザに対応するターゲット (tID = 1 ~ n) の

(a) 存在位置

(b) 誰であるか (uID1 ~ uIDk のいずれであるか)

これらを示す情報である。このターゲット情報は、パーティクルの更新に伴い、順次更新されることになり、例えばユーザ 1 ~ k が実環境内で移動しない場合、ユーザ 1 ~ k の各々が、n 個のターゲット (tID = 1 ~ n) から選択された k 個にそれぞれ対応するデータとして収束することになる。

#### 【0077】

40

例えば、図 5 に示すターゲット情報 305 中の最上段のターゲット 1 (tID = 1) のデータ中に含まれるユーザ確信度情報 (uID) は、ユーザ 2 ( $uID_{12} = 0.7$ ) について最も高い確率を有している。従って、このターゲット 1 (tID = 1) のデータは、ユーザ 2 に対応するものであると推定されることになる。なお、ユーザ確信度情報 (uID) を示すデータ [ $uID_{12} = 0.7$ ] 中の ( $uID_{12}$ ) 内の (12) は、ターゲット ID = 1 のユーザ = 2 のユーザ確信度情報 (uID) に対応する確率であることを示している。

#### 【0078】

このターゲット情報 305 中の最上段のターゲット 1 (tID = 1) のデータは、ユーザ 2 である確率が最も高く、このユーザ 2 は、その存在位置が、ターゲット情報 305 中

50

の最上段のターゲット 1 (  $tID = 1$  ) のデータに含まれる存在確率分布データに示す範囲にいと推定されることとなる。

【 0 0 7 9 】

このように、ターゲット情報 3 0 5 は、初期的に仮想的なオブジェクト ( 仮想ユーザ ) として設定した各ターゲット (  $tID = 1 \sim n$  ) の各々について、

( a ) 存在位置

( b ) 誰であるか (  $uID 1 \sim uID k$  のいずれであるか )

の各情報を示す。従って、各ターゲット (  $tID = 1 \sim n$  ) の  $k$  個のターゲット情報の各々は、ユーザが移動しない場合は、ユーザ 1 ~  $k$  に対応するように収束する。

【 0 0 8 0 】

ターゲット (  $tID = 1 \sim n$  ) の数がユーザ数  $k$  より大きい場合、どのユーザにも対応しないターゲットが発生する。例えば、ターゲット情報 3 0 5 中の最下段のターゲット (  $tID = n$  ) は、ユーザ確信度情報 (  $uID$  ) も最大で 0 . 5 であり、存在確率分布データも大きなピークを有していない。このようなデータは特定のユーザに対応するデータではないと判定される。なお、このようなターゲットについては、削除するような処理が行われる場合もある。ターゲットの削除処理については後述する。

【 0 0 8 1 】

先に説明したように、音声・画像統合処理部 1 3 1 は、入力情報に基づくパーティクルの更新処理を実行して、

( a ) 複数のユーザが、それぞれどこにいて、それらは誰であるかの推定情報としての [ ターゲット情報 ] 、

( b ) 例えば話をしたユーザなどのイベント発生源を示す [ シグナル情報 ] 、

これらを生成して処理決定部 1 3 2 に出力する。

【 0 0 8 2 】

ターゲット情報は、図 5 のターゲット情報 3 0 5 を参照して説明した情報である。音声・画像統合処理部 1 3 1 は、このターゲット情報の他に話をしたユーザなどのイベント発生源を示す [ シグナル情報 ] についても生成して出力する。イベント発生源を示す [ シグナル情報 ] は、音声イベントについては、誰が話をしたか、すなわち [ 話者 ] を示すデータであり、画像イベントについては、画像に含まれる顔が誰であるかを示すデータである。なお、画像イベントの場合のシグナル情報は、本例では結果としてターゲット情報のユーザ確信度情報 (  $uID$  ) から得られるものと一致することになる。

【 0 0 8 3 】

音声・画像統合処理部 1 3 1 が、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から、図 3 ( B ) に示すイベント情報、すなわち、ユーザ位置情報と、ユーザ識別情報 ( 顔識別情報または話者識別情報 ) 、これらのイベント情報を入力して、

( a ) 複数のユーザが、それぞれどこにいて、それらは誰であるかの推定情報としての [ ターゲット情報 ] 、

( b ) 例えば話をしたユーザなどのイベント発生源を示す [ シグナル情報 ] 、

これらの情報を生成して処理決定部 1 3 2 に出力する処理について、図 7 以下を参照して説明する。

【 0 0 8 4 】

図 7 は、音声・画像統合処理部 1 3 1 の実行する処理シーケンスを説明するフローチャートを示す図である。まず、ステップ S 1 0 1 において、音声・画像統合処理部 1 3 1 は、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から、図 3 ( B ) に示すイベント情報、すなわち、ユーザ位置情報と、ユーザ識別情報 ( 顔識別情報または話者識別情報 ) 、これらのイベント情報を入力する。

【 0 0 8 5 】

イベント情報の取得に成功した場合は、ステップ S 1 0 2 に進み、イベント情報の取得に失敗した場合は、ステップ S 1 2 1 に進む。ステップ S 1 2 1 の処理については後段で説明する。

10

20

30

40

50

## 【 0 0 8 6 】

イベント情報の取得に成功した場合は、音声・画像統合処理部 1 3 1 は、ステップ S 1 0 2 以下において、入力情報に基づくパーティクル更新処理を行うことになるが、パーティクル更新処理の前にステップ S 1 0 2 において、図 5 に示す  $m$  個のパーティクル ( $pID = 1 \sim m$ ) の各々にイベントの発生源の仮説を設定する。イベント発生源とは、例えば、音声イベントであれば、話をしたユーザがイベント発生源であり、画像イベントであれば、抽出した顔を持つユーザがイベント発生源である。

## 【 0 0 8 7 】

図 5 に示す例では、各パーティクルの最下段にイベント発生源の仮説データ ( $tID = x \times$ ) を示している。図 5 の例では、

パーティクル 1 ( $pID = 1$ ) は、 $tID = 2$ 、

パーティクル 2 ( $pID = 2$ ) は、 $tID = n$ 、

:

パーティクル  $m$  ( $pID = m$ ) は、 $tID = n$ 、

このように各パーティクルについて、イベント発生源がターゲット 1 ~  $n$  のいずれであるかの仮説を設定する。図 5 に示す例では、各パーティクルについて、仮説として設定したイベント発生源のターゲットデータを二重線で囲んで示している。

## 【 0 0 8 8 】

このイベント発生源の仮説設定は、入力イベントに基づくパーティクル更新処理を行う前に毎回実行する。すなわち、各パーティクル 1 ~  $m$  各々にイベントの発生源仮説を設定して、その仮説の下で、イベントとして音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から、図 3 (B) に示すイベント情報、すなわち、

(a) ユーザ位置情報

(b) ユーザ識別情報 (顔識別情報または話者識別情報)

これらのイベント情報を入力して  $m$  個のパーティクル ( $PID = 1 \sim m$ ) の更新処理を行う。

## 【 0 0 8 9 】

パーティクル更新処理が行われた場合は、各パーティクル 1 ~  $m$  各々に設定されていたイベントの発生源の仮説はリセットされて、各パーティクル 1 ~  $m$  各々に新たな仮説の設定が行われる。この仮説の設定態様としては、

(1) ランダムな設定、

(2) 音声・画像統合処理部 1 3 1 の有する内部モデルに従って設定、

上記 (1) , (2) のいずれかの手法で設定することが可能である。なお、パーティクルの数 :  $m$  は、ターゲットの数 :  $n$  より大きく設定されているので、複数のパーティクルが同一のターゲットをイベント発生源とした仮説に設定される。例えば、ターゲットの数 :  $n$  が 10 とした場合、パーティクル数 :  $m = 100 \sim 1000$  程度に設定した処理などが行われる。

## 【 0 0 9 0 】

上記の (2) 音声・画像統合処理部 1 3 1 の有する内部モデルに従って仮説を設定する処理の具体的処理例について説明する。

音声・画像統合処理部 1 3 1 は、まず、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から取得したイベント情報、すなわち、図 3 (B) に示す 2 つの情報、すなわち、

(a) ユーザ位置情報

(b) ユーザ識別情報 (顔識別情報または話者識別情報)

これらのイベント情報と、

音声・画像統合処理部 1 3 1 の保持するパーティクルのターゲットの持つデータとの比較によって、各ターゲットの重み [ $W_{tID}$ ] を算出し、算出した各ターゲットの重み [ $W_{tID}$ ] に基づいて、各パーティクル ( $pID = 1 \sim m$ ) に対するイベント発生源の仮説を設定する。以下、具体的な処理例について説明する。

10

20

30

40

50

## 【0091】

なお、初期状態では、各パーティクル ( $pID = 1 \sim m$ ) に設定されるイベント発生源の仮説は均等な設定とする。すなわち  $n$  個のターゲット ( $tID = 1 \sim n$ ) を持つ  $m$  個のパーティクル ( $pID = 1 \sim m$ ) が設定されている構成では、

ターゲット 1 ( $tID = 1$ ) をイベント発生源とするパーティクルを  $m/n$  個、

ターゲット 2 ( $tID = 2$ ) をイベント発生源とするパーティクルを  $m/n$  個、

:

ターゲット  $n$  ( $tID = n$ ) をイベント発生源とするパーティクルを  $m/n$  個、

というように、各パーティクル ( $pID = 1 \sim m$ ) に設定する初期的なイベント発生源の仮説ターゲット ( $tID = 1 \sim n$ ) を均等に割り振る設定とする。

10

## 【0092】

図 7 に示すフローのステップ S 101 において、音声・画像統合処理部 131 が音声イベント検出部 122 および画像イベント検出部 112 からイベント情報、すなわち、図 3 (B) に示す 2 つの情報、すなわち、

(a) ユーザ位置情報

(b) ユーザ識別情報 (顔識別情報または話者識別情報)

これらのイベント情報を取得して、イベント情報の取得に成功すると、ステップ S 102 において、音声・画像統合処理部 131 は、 $m$  個のパーティクル ( $PID = 1 \sim m$ ) の各々に対して、イベント発生源の仮説ターゲット ( $tID = 1 \sim n$ ) を設定する。

20

## 【0093】

ステップ S 102 におけるパーティクル対応の仮説ターゲットの設定の詳細について説明する。音声・画像統合処理部 131 は、まず、ステップ S 101 で入力したイベント情報と、音声・画像統合処理部 131 の保持するパーティクルのターゲットの持つデータとの比較を行い、比較結果を用いて、各ターゲットのターゲット重み [ $W_{tID}$ ] を算出する。

## 【0094】

ターゲット重み [ $W_{tID}$ ] の算出処理の詳細について図 8 を参照して説明する。ターゲット重みの算出は、図 8 の右端に示すように、各パーティクルに設定されるターゲット 1 ~  $n$  の各々に対応する  $n$  個のターゲット重みの算出処理として実行される。この  $n$  個のターゲット重みの算出に際しては、まず、図 8 (1) に示す入力イベント情報、すなわち、音声・画像統合処理部 131 が、音声イベント検出部 122 および画像イベント検出部 112 から入力したイベント情報と、各パーティクルの各ターゲットデータとの類似度の指標値としての尤度算出を行う。

30

## 【0095】

図 8 (2) に示す尤度算出処理例は、(1) 入力イベント情報と、パーティクル 1 の 1 つのターゲットデータ ( $tID = n$ ) との比較によるイベント - ターゲット間尤度の算出例を説明する図である。なお、図 8 には、1 つのターゲットデータとの比較例を示しているが、各パーティクルの各ターゲットデータについて、同様の尤度算出処理を実行する。

## 【0096】

図 8 の下段に示す (2) 尤度算出処理について説明する。図 8 (2) に示すように、尤度算出処理は、まず、

40

(a) ユーザ位置情報についてのイベントと、ターゲットデータとの類似度データとしてのガウス分布間尤度 [ $DL$ ]、

(b) ユーザ識別情報 (顔識別情報または話者識別情報) についてのイベントと、ターゲットデータとの類似度データとしてのユーザ確信度情報 ( $UID$ ) 間尤度 [ $UL$ ]

これらを個別に算出する。

## 【0097】

まず、(a) ユーザ位置情報についてのイベントと、ターゲットデータとの類似度データとしてのガウス分布間尤度 [ $DL$ ] の算出処理について説明する。

図 8 (1) に示す入力イベント情報中の、ユーザ位置情報に対応するガウス分布を  $N$  (

50

$m_e, \sigma_e$ )とし、

音声・画像統合処理部131の保持する内部モデルのあるパーティクルが持つあるターゲットのユーザ位置情報に対応するガウス分布を $N(m_t, \sigma_t)$ とする。図8に示す例では、パーティクル1( $pID=1$ )のターゲット $n(tID=n)$ のターゲットデータに含まれるガウス分布を $N(m_t, \sigma_t)$ とする。

【0098】

これら2つのデータのガウス分布の類似度を判定する指標としてのガウス分布間尤度 $[DL]$ は、以下の式によって算出する。

$$DL = N(m_t, \sigma_t + \sigma_e) \times |m_e|$$

上記式は、中心 $m_t$ で分散 $\sigma_t + \sigma_e$ のガウス分布において $x = m_e$ の位置の値を算出する式である。

10

【0099】

次に、(b)ユーザ識別情報(顔識別情報または話者識別情報)についてのイベントと、ターゲットデータとの類似度データとしてのユーザ確信度情報( $uID$ )間尤度 $[UL]$ の算出処理について説明する。

図8(1)に示す入力イベント情報中の、ユーザ確信度情報( $uID$ )の各ユーザ1~ $k$ の確信度の値(スコア)を $P_e[i]$ とする。なお、 $i$ はユーザ識別子1~ $k$ に対応する変数である。

音声・画像統合処理部131の保持する内部モデルのあるパーティクルが持つあるターゲットのユーザ確信度情報( $uID$ )の各ユーザ1~ $k$ の確信度の値(スコア)を $P_t[i]$ とする。図8に示す例では、パーティクル1( $pID=1$ )のターゲット $n(tID=n)$ のターゲットデータに含まれるユーザ確信度情報( $uID$ )の各ユーザ1~ $k$ の確信度の値(スコア)を $P_t[i]$ とする。

20

【0100】

これら2つのデータのユーザ確信度情報( $uID$ )の類似度を判定する指標としてのユーザ確信度情報( $uID$ )間尤度 $[UL]$ は、以下の式によって算出する。

$$UL = P_e[i] \times P_t[i]$$

上記式は、2つのデータのユーザ確信度情報( $uID$ )に含まれる各対応ユーザの確信度の値(スコア)の積の総和を求める式であり、この値をユーザ確信度情報( $uID$ )間尤度 $[UL]$ とする。

30

【0101】

もしくは、ユーザ確信度情報( $uID$ )間尤度 $[UL]$ として、各積の最大値、すなわち、

$$UL = \arg \max (P_e[i] \times P_t[i])$$

上記の値を算出し、この値をユーザ確信度情報( $uID$ )間尤度 $[UL]$ として利用する構成としてもよい。

【0102】

入力イベント情報とあるパーティクル( $pID$ )が持つ1つのターゲット( $tID$ )との類似度の指標としてのイベント-ターゲット間尤度 $[L_{pID, tID}]$ は、上記の2つの尤度、すなわち、

40

ガウス分布間尤度 $[DL]$ と、

ユーザ確信度情報( $uID$ )間尤度 $[UL]$

これら2つの尤度を利用して算出する。すなわち重み( $\alpha = 0 \sim 1$ )を用いて、イベント-ターゲット間尤度 $[L_{pID, tID}]$ は下式によって算出する。

$$[L_{pID, tID}] = UL^\alpha \times DL^{1-\alpha}$$

としてイベントとターゲットとの類似度の指標であるイベント-ターゲット間尤度 $[L_{pID, tID}]$ を算出する。

ただし、 $\alpha = 0 \sim 1$ とする。

【0103】

このイベント-ターゲット間尤度 $[L_{pID, tID}]$ は、各パーティクルの各ターゲ

50

ットについて各々算出し、このイベント - ターゲット間尤度  $[L_{pID}, tID]$  に基づいて各ターゲットのターゲット重み  $[W_{tID}]$  を算出する。

【0104】

なお、イベント - ターゲット間尤度  $[L_{pID}, tID]$  の算出に適用する重み  $[ ]$  は、予め固定された値としてもよいし、入力イベントに応じて値を変更する設定としてもよい。例えば入力イベントが画像である場合において、顔検出に成功し位置情報は取得できたが顔識別に失敗した場合などは、 $= 0$  の設定として、ユーザ確信度情報 ( $uID$ ) 間尤度:  $UL = 1$  としてガウス分布間尤度  $[DL]$  のみに依存してイベント - ターゲット間尤度  $[L_{pID}, tID]$  を算出して、ガウス分布間尤度  $[DL]$  のみに依存したターゲット重み  $[W_{tID}]$  を算出する構成としてもよい。

10

【0105】

また、入力イベントが音声である場合において、話者識別に成功し話者情報破取得できたが、位置情報の取得に失敗した場合などは、 $= 0$  の設定として、ガウス分布間尤度  $[DL] = 1$  として、ユーザ確信度情報 ( $uID$ ) 間尤度  $[UL]$  のみに依存してイベント - ターゲット間尤度  $[L_{pID}, tID]$  を算出して、ユーザ確信度情報 ( $uID$ ) 間尤度  $[UL]$  のみに依存したターゲット重み  $[W_{tID}]$  を算出する構成としてもよい。

【0106】

イベント - ターゲット間尤度  $[L_{pID}, tID]$  に基づく、ターゲット重み  $[W_{tID}]$  の算出式は、以下の通りである。

【数1】

20

$$W_{tID} = \sum_{pID}^m W_{pID} L_{pID, tID}$$

【0107】

とする。なお、上記式において、 $[W_{pID}]$  は、各パーティクル各々に設定されるパーティクル重みである。パーティクル重み  $[W_{pID}]$  の算出処理については後段で説明する。パーティクル重み  $[W_{pID}]$  は初期状態では、すべてのパーティクル ( $pID = 1 \sim m$ ) において均一な値が設定される。

30

【0108】

図7に示すフローにおけるステップS101の処理、すなわち、各パーティクル対応のイベント発生源仮説の生成は、上記のイベント - ターゲット間尤度  $[L_{pID}, tID]$  に基づいて算出したターゲット重み  $[W_{tID}]$  に基づいて実行する。ターゲット重み  $[W_{tID}]$  は、パーティクルに設定されるターゲット  $1 \sim n$  ( $tID = 1 \sim n$ ) に対応した  $n$  個のデータが算出される。

【0109】

$m$  個のパーティクル ( $pID = 1 \sim m$ ) 各々に対するイベント発生源仮説ターゲットは、ターゲット重み  $[W_{tID}]$  の比率に応じて割り振る設定とする。

40

例えば  $n = 4$  で、ターゲット  $1 \sim 4$  ( $tID = 1 \sim 4$ ) に対応して算出されたターゲット重み  $[W_{tID}]$  が、

ターゲット1: ターゲット重み = 3

ターゲット2: ターゲット重み = 2

ターゲット3: ターゲット重み = 1

ターゲット4: ターゲット重み = 5

である場合、 $m$  個のパーティクルのイベント発生源仮説ターゲットを

$m$  個のパーティクル中の30%をイベント発生源仮説ターゲット1、

$m$  個のパーティクル中の20%をイベント発生源仮説ターゲット2、

$m$  個のパーティクル中の10%をイベント発生源仮説ターゲット3、

50

m個のパーティクル中の50%をイベント発生源仮説ターゲット4、  
このような設定とする。

すなわちパーティクルに設定するイベント発生源仮説ターゲットをターゲットの重みに  
応じた配分比率とする。

#### 【0110】

この仮説設定の後、図7に示すフローのステップS103に進む。ステップS103では、各パーティクル対応の重み、すなわちパーティクル重み $[W_{pID}]$ の算出を行う。このパーティクル重み $[W_{pID}]$ は前述したように、初期的には各パーティクルに均一な値が設定されるが、イベント入力に応じて更新される。

#### 【0111】

図9、図10を参照して、パーティクル重み $[W_{pID}]$ の算出処理の詳細について説明する。パーティクル重み $[W_{pID}]$ は、イベント発生源の仮説ターゲットを生成した各パーティクルの仮説の正しさの指標に相当する。パーティクル重み $[W_{pID}]$ は、m個のパーティクル( $pID = 1 \sim m$ )の各々において設定されたイベント発生源の仮説ターゲットと、入力イベントとの類似度であるイベント-ターゲット間尤度として算出される。

#### 【0112】

図9には、音声・画像統合処理部131が、音声イベント検出部122および画像イベント検出部112から入力するイベント情報401と、音声・画像統合処理部131が、保持するパーティクル411~413を示している。核パーティクル411|413には、前述した処理、すなわち、図7に示すフローのステップS102におけるイベント発生源の仮説設定において設定された仮説ターゲットが1つずつ設定されている。図9中に示す例では、

パーティクル1( $pID = 1$ )411におけるターゲット2( $tID = 2$ )421、  
パーティクル2( $pID = 2$ )412におけるターゲットn( $tID = n$ )422、  
パーティクルm( $pID = m$ )413におけるターゲットn( $tID = n$ )423、  
これらの仮説ターゲットである。

#### 【0113】

図9の例において、各パーティクルのパーティクル重み $[W_{pID}]$ は、  
パーティクル1：イベント情報401とターゲット2( $tID = 2$ )421とのイベント-ターゲット間尤度、  
パーティクル2：イベント情報401とターゲットn( $tID = n$ )422とのイベント-ターゲット間尤度、  
パーティクルm：イベント情報401とターゲットn( $tID = n$ )423とのイベント-ターゲット間尤度、  
これらのイベント-ターゲット間尤度に対応することになる。

#### 【0114】

図10は、パーティクル1( $pID = 1$ )のパーティクル重み $[W_{pID}]$ 算出処理例を示している。図10(2)に示すパーティクル重み $[W_{pID}]$ 算出処理は、先に、図8(2)を参照して説明したと同様の尤度算出処理であり、本例では、(1)入力イベント情報と、パーティクルから選択された唯一の仮説ターゲットとの類似度指標としてのイベント-ターゲット間尤度の算出として実行される。

#### 【0115】

図10の下段に示す(2)尤度算出処理も、先に図8(2)を参照して説明したと同様、

(a) ユーザ位置情報についてのイベントと、ターゲットデータとの類似度データとしてのガウス分布間尤度 $[DL]$ 、

(b) ユーザ識別情報(顔識別情報または話者識別情報)についてのイベントと、ターゲットデータとの類似度データとしてのユーザ確信度情報( $uID$ )間尤度 $[UL]$

これらを個別に算出する。

10

20

30

40

50

## 【0116】

(a) ユーザ位置情報についてのイベントと、仮説ターゲットとの類似度データとしてのガウス分布間尤度  $[DL]$  の算出処理は以下の処理となる。

入力イベント情報中の、ユーザ位置情報に対応するガウス分布を  $N(m_e, \sigma_e)$ 、パーティクルから選択された仮説ターゲットのユーザ位置情報に対応するガウス分布を  $N(m_t, \sigma_t)$ 、

として、ガウス分布間尤度  $[DL]$  を、以下の式によって算出する。

$$DL = N(m_t, \sigma_t + \sigma_e) \times |m_e|$$

上記式は、中心  $m_t$  で分散  $\sigma_t + \sigma_e$  のガウス分布において  $x = m_e$  の位置の値を算出する式である。

10

## 【0117】

(b) ユーザ識別情報（顔識別情報または話者識別情報）についてのイベントと、仮説ターゲットとの類似度データとしてのユーザ確信度情報（ $UID$ ）間尤度  $[UL]$  の算出処理は以下の処理となる。

入力イベント情報中の、ユーザ確信度情報（ $UID$ ）の各ユーザ  $1 \sim k$  の確信度の値（スコア）を  $P_e[i]$  とする。なお、 $i$  はユーザ識別子  $1 \sim k$  に対応する変数である。

パーティクルから選択された仮説ターゲットのユーザ確信度情報（ $UID$ ）の各ユーザ  $1 \sim k$  の確信度の値（スコア）を  $P_t[i]$  として、ユーザ確信度情報（ $UID$ ）間尤度  $[UL]$  は、以下の式によって算出する。

$$UL = P_e[i] \times P_t[i]$$

20

上記式は、2つのデータのユーザ確信度情報（ $UID$ ）に含まれる各対応ユーザの確信度の値（スコア）の積の総和を求める式であり、この値をユーザ確信度情報（ $UID$ ）間尤度  $[UL]$  とする。

## 【0118】

パーティクル重み  $[W_{PID}]$  は、上記の2つの尤度、すなわち、

ガウス分布間尤度  $[DL]$  と、

ユーザ確信度情報（ $UID$ ）間尤度  $[UL]$

これら2つの尤度を利用し、重み（ $= 0 \sim 1$ ）を用いて下式によって算出する。

$$W_{PID} = UL \times DL^{-1}$$

上記式により、パーティクル重み  $[W_{PID}]$  を算出する。

30

ただし、 $= 0 \sim 1$  とする。

このパーティクル重み  $[W_{PID}]$  は、各パーティクルについて各々算出する。

## 【0119】

なお、パーティクル重み  $[W_{PID}]$  の算出に適用する重み  $[ ]$  は、前述したイベント - ターゲット間尤度  $[L_{PID, tID}]$  の算出処理と同様、予め固定された値としてもよいし、入力イベントに応じて値を変更する設定としてもよい。例えば入力イベントが画像である場合において、顔検出に成功し位置情報は取得できたが顔識別に失敗した場合などは、 $= 0$  の設定として、ユーザ確信度情報（ $UID$ ）間尤度： $UL = 1$  としてガウス分布間尤度  $[DL]$  のみに依存してパーティクル重み  $[W_{PID}]$  を算出する構成としてもよい。また、入力イベントが音声である場合において、話者識別に成功し話者情報破取得できたが、位置情報の取得に失敗した場合などは、 $= 0$  の設定として、ガウス分布間尤度  $[DL] = 1$  として、ユーザ確信度情報（ $UID$ ）間尤度  $[UL]$  のみに依存してパーティクル重み  $[W_{PID}]$  を算出する構成としてもよい。

40

## 【0120】

図7のフローにおけるステップS103の各パーティクル対応の重み  $[W_{PID}]$  の算出は、このように図9、図10を参照して説明した処理として実行される。次に、ステップS104において、ステップS103で設定した各パーティクルのパーティクル重み  $[W_{PID}]$  に基づくパーティクルのリサンプリング処理を実行する。

## 【0121】

このパーティクルリサンプリング処理は、 $m$  個のパーティクルから、パーティクル重み

50



[  $W_{pID}$  ] に応じてパーティクルを取捨選択する処理として実行される。具体的には、例えば、パーティクル数：  $m = 5$  のとき、

パーティクル 1：パーティクル重み [  $W_{pID}$  ] = 0.40

パーティクル 2：パーティクル重み [  $W_{pID}$  ] = 0.10

パーティクル 3：パーティクル重み [  $W_{pID}$  ] = 0.25

パーティクル 4：パーティクル重み [  $W_{pID}$  ] = 0.05

パーティクル 5：パーティクル重み [  $W_{pID}$  ] = 0.20

これらのパーティクル重みが各々設定されていた場合、

パーティクル 1 は、40% の確率でリサンプリングされ、パーティクル 2 は 10% の確率でリサンプリングされる。なお、実際には  $m = 100 \sim 1000$  といった多数であり、リサンプリングされた結果は、パーティクルの重みに応じた配分比率のパーティクルによって構成されることになる。

#### 【0122】

この処理によって、パーティクル重み [  $W_{pID}$  ] の大きなパーティクルがより多く残存することになる。なお、リサンプリング後もパーティクルの総数 [  $m$  ] は変更されない。また、リサンプリング後は、各パーティクルの重み [  $W_{pID}$  ] はリセットされ、新たなイベントの入力に応じてステップ S101 から処理が繰り返される。

#### 【0123】

ステップ S105 では、各パーティクルに含まれるターゲットデータ（ユーザ位置およびユーザ確信度）の更新処理を実行する。各ターゲットは、先に図 6 等を参照して説明したように、

(a) ユーザ位置：各ターゲット各々に対応する存在位置の確率分布 [ ガウス分布：  $N(m_t, \sigma_t)$  ]、

(b) ユーザ確信度：各ターゲットが誰であることを示すユーザ確信度情報 (  $uID$  ) として各ユーザ 1 ~  $k$  である確立値 ( スコア )：  $P_t[i]$  (  $i = 1 \sim k$  )、すなわち、

$$uID_{t1} = P_t[1]$$

$$uID_{t2} = P_t[2]$$

：

$$uID_{tk} = P_t[k]$$

これらのデータによって構成される。

#### 【0124】

ステップ S105 におけるターゲットデータの更新は、(a) ユーザ位置、(b) ユーザ確信度の各々について実行する。まず、(a) ユーザ位置の更新処理について説明する。

#### 【0125】

ユーザ位置の更新は、

(a1) 全パーティクルの全ターゲットを対象とする更新処理、

(a2) 各パーティクルに設定されたイベント発生源仮説ターゲットを対象とした更新処理、

これらの 2 段階の更新処理として実行する。

#### 【0126】

(a1) 全パーティクルの全ターゲットを対象とする更新処理は、イベント発生源仮説ターゲットとして選択されたターゲットおよびその他のターゲットのすべてを対象として実行する。この処理は、時間経過に伴うユーザ位置の分散が拡大するという仮定に基づいて実行され、前回の更新処理からの経過時間とイベントの位置情報によってカルマン・フィルタ ( Kalman Filter ) を用い更新される。

#### 【0127】

以下、位置情報が 1 次元の場合の更新処理例について説明する。まず、前回の更新処理時間からの経過時間 [  $dt$  ] とし、全ターゲットについての、 $dt$  後のユーザ位置の予測分布を計算する。すなわち、ユーザ位置の分布情報としてのガウス分布：  $N(m_t, \sigma_t)$

）の期待値（平均）： $[m_t]$ 、分散 $[P_t]$ について、以下の更新を行う。

$$m_t = m_t + x_c \times d_t$$

$$P_t = P_t + c^2 \times d_t$$

なお、

$m_t$ ：予測期待値（predicted state）

$P_t$ ：予測共分散（predicted estimate covariance）

$x_c$ ：移動情報（control model）

$c^2$ ：ノイズ（process noise）

である。

10

なお、ユーザが移動しない条件の下で処理する場合は、 $x_c = 0$ として更新処理を行うことができる。

上記の算出処理により、全ターゲットに含まれるユーザ位置情報としてのガウス分布： $N(m_t, P_t)$ を更新する。

#### 【0128】

さらに、各パーティクルに1つ設定されているイベント発生源の仮説となったターゲットに関しては、音声イベント検出部122や画像イベント検出部112から入力するイベント情報に含まれるユーザ位置を示すガウス分布： $N(m_e, P_e)$ を用いた更新処理を実行する。

$K$ ：カルマンゲイン（Kalman Gain）

20

$m_e$ ：入力イベント情報： $N(m_e, P_e)$ に含まれる観測値（Observed state）

$P_e$ ：入力イベント情報： $N(m_e, P_e)$ に含まれる観測値（Observed covariance）

として、以下の更新処理を行う。

$$K = P_t / (P_t + P_e)$$

$$m_t = m_t + K(x_c - m_t)$$

$$P_t = (1 - K) P_t$$

#### 【0129】

次に、ターゲットデータの更新処理として実行する（b）ユーザ確信度の更新処理について説明する。ターゲットデータには上記のユーザ位置情報の他に、各ターゲットが誰であるかを示すユーザ確信度情報（uID）として各ユーザ1～kである確立値（スコア）： $P_t[i]$ （ $i = 1 \sim k$ ）が含まれている。ステップS105では、このユーザ確信度情報（uID）についても更新処理を行う。

30

#### 【0130】

各パーティクルに含まれるターゲットのユーザ確信度情報（uID）： $P_t[i]$ （ $i = 1 \sim k$ ）についての更新は、登録ユーザ全員分の事後確率と、音声イベント検出部122や画像イベント検出部112から入力するイベント情報に含まれるユーザ確信度情報（uID）： $P_e[i]$ （ $i = 1 \sim k$ ）によって、予め設定した0～1の範囲の値を持つ更新率 $[w]$ を適用して更新する。

40

#### 【0131】

ターゲットのユーザ確信度情報（uID）： $P_t[i]$ （ $i = 1 \sim k$ ）についての更新は、以下の式によって実行する。

$$P_t[i] = (1 - w) \times P_t[i] + w \times P_e[i]$$

ただし、

$i = 1 \sim k$

$w : 0 \sim 1$

である。なお、更新率 $[w]$ は、0～1の範囲の値であり予め設定する。

#### 【0132】

ステップS105では、この更新されたターゲットデータに含まれる以下のデータ、す

50

なわち、

( a ) ユーザ位置：各ターゲット各々に対応する存在位置の確率分布 [ ガウス分布：  $N(m_t, \sigma_t)$  ]、

( b ) ユーザ確信度：各ターゲットが誰であることを示すユーザ確信度情報 ( u I D ) と  
して各ユーザ 1 ~ k である確立値 ( スコア ) :  $P_t[i]$  (  $i = 1 \sim k$  )、すなわち、

$$u I D_{t_1} = P_t[1]$$

$$u I D_{t_2} = P_t[2]$$

:

$$u I D_{t_k} = P_t[k]$$

これらのデータと、各パーティクル重み [  $W_{p I D}$  ] とに基づいて、ターゲット情報を生成して、処理決定部 1 3 2 に出力する。

10

【 0 1 3 3 】

なお、ターゲット情報の生成は、図 5 を参照して説明したように、各パーティクル (  $P I D = 1 \sim m$  ) に含まれる各ターゲット (  $t I D = 1 \sim n$  ) 対応データの重み付き総和データとして生成される。図 5 の右端のターゲット情報 3 0 5 に示すデータである。ターゲット情報は、各ターゲット (  $t I D = 1 \sim n$  ) 各々の

( a ) ユーザ位置情報、

( b ) ユーザ確信度情報、

これらの情報を含む情報として生成される。

20

【 0 1 3 4 】

例えば、ターゲット (  $t I D = 1$  ) に対応するターゲット情報中の、ユーザ位置情報は、

【 数 2 】

$$\sum_{i=1}^m W_i \cdot N(m_{i1}, \sigma_{i1})$$

30

【 0 1 3 5 】

上記式で表される。上記式において、 $W_i$  は、パーティクル重み [  $W_{p I D}$  ] を示している。

【 0 1 3 6 】

また、ターゲット (  $t I D = 1$  ) に対応するターゲット情報中の、ユーザ確信度情報は、

、

【数 3】

$$\begin{aligned} & \sum_{i=1}^m W_i \cdot uID_{i1} \\ & \sum_{i=1}^m W_i \cdot uID_{i2} \\ & \vdots \\ & \sum_{i=1}^m W_i \cdot uID_{ik} \end{aligned}$$

10

【0137】

上記式で表される。上記式において、 $W_i$  は、パーティクル重み  $[W_{PID}]$  を示している。

20

音声・画像統合処理部131は、これらのターゲット情報を  $n$  個の各ターゲット ( $tID = 1 \sim n$ ) 各々について算出し、算出したターゲット情報を処理決定部132に出力する。

【0138】

次に、図7に示すフローのステップS106の処理について説明する。音声・画像統合処理部131は、ステップS106において、 $n$  個のターゲット ( $tID = 1 \sim n$ ) の各々がイベントの発生源である確率を算出し、これをシグナル情報として処理決定部132に出力する。

【0139】

先に説明したように、イベント発生源を示す[シグナル情報]は、音声イベントについては、誰が話をしたか、すなわち[話者]を示すデータであり、画像イベントについては、画像に含まれる顔が誰であるかを示すデータである。

30

【0140】

音声・画像統合処理部131は、各パーティクルに設定されたイベント発生源の仮説ターゲットの数に基づいて、各ターゲットがイベント発生源である確率を算出する。すなわち、ターゲット ( $tID = 1 \sim n$ ) の各々がイベント発生源である確率を  $[P(tID = i)]$  とする。ただし  $i = 1 \sim n$  である。このとき、各ターゲットがイベント発生源である確率は、以下のように算出される。

$P(tID = 1)$  :  $tID = 1$  を割り当てた数 /  $m$

$P(tID = 2)$  :  $tID = 2$  を割り当てた数 /  $m$

:

$P(tID = n)$  :  $tID = n$  を割り当てた数 /  $m$

40

音声・画像統合処理部131は、この算出処理によって、生成した情報、すなわち、各ターゲットがイベント発生源である確率を[シグナル情報]として、処理決定部132に出力する。

【0141】

ステップS106の処理が終了したら、ステップS101に戻り、音声イベント検出部122および画像イベント検出部112からのイベント情報の入力の待機状態に移行する。

【0142】

50

以上が、図 7 に示すフローのステップ S 1 0 1 ~ S 1 0 6 の説明である。ステップ S 1 0 1 において、音声・画像統合処理部 1 3 1 が、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から、図 3 ( B ) に示すイベント情報を取得できなかった場合も、ステップ S 1 2 1 において、各パーティクルに含まれるターゲットの構成データの更新が実行される。この更新は、時間経過に伴うユーザ位置の変化を考慮した処理である。

#### 【 0 1 4 3 】

このターゲット更新処理は、先に、ステップ S 1 0 5 の説明において ( a 1 ) 全パーティクルの全ターゲットを対象とする更新処理と同様の処理であり、時間経過に伴うユーザ位置の分散が拡大するという仮定に基づいて実行され、前回の更新処理からの経過時間とイベントの位置情報によってカルマン・フィルタ ( K a l m a n F i l t e r ) を用い

10

#### 【 0 1 4 4 】

位置情報が 1 次元の場合の更新処理例について説明する。まず、前回の更新処理時間からの経過時間 [ d t ] とし、全ターゲットについての、d t 後のユーザ位置の予測分布を計算する。すなわち、ユーザ位置の分布情報としてのガウス分布：N ( m t , t ) の期待値 ( 平均 ) : [ m t ] 、分散 [ t ] について、以下の更新を行う。

$$m_t = m_t + x_c \times d_t$$

$$t^2 = t^2 + c^2 \times d_t$$

なお、

m t : 予測期待値 ( p r e d i c t e d s t a t e )

t 2 : 予測共分散 ( p r e d i c t e d e s t i m a t e c o v a r i a n c e )

20

x c : 移動情報 ( c o n t r o l m o d e l )

c 2 : ノイズ ( p r o c e s s n o i s e )

である。

なお、ユーザが移動しない条件の下で処理する場合は、x c = 0 として更新処理を行うことができる。

上記の算出処理により、全ターゲットに含まれるユーザ位置情報としてのガウス分布：N ( m t , t ) を更新する。

#### 【 0 1 4 5 】

30

なお、各パーティクルのターゲットに含まれるユーザ確信度情報 ( u I D ) については、イベントの登録ユーザ全員分の事後確率、もしくはイベント情報からスコア [ P e ] が取得できない限りは更新しない。

#### 【 0 1 4 6 】

ステップ S 1 2 1 の処理が終了したら、ステップ S 1 0 1 に戻り、音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 からのイベント情報の入力の待機状態に移行する。

#### 【 0 1 4 7 】

以上、図 7 を参照して音声・画像統合処理部 1 3 1 の実行する処理について説明した。音声・画像統合処理部 1 3 1 は、図 7 に示すフローに従った処理を音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 からのイベント情報の入力ごとに繰り返し実行する。この繰り返し処理により、より信頼度の高いターゲットを仮説ターゲットとして設定したパーティクルの重みが大きくなり、パーティクル重みに基づくリサンプリング処理により、より重みの大きいパーティクルが残存することになる。結果として音声イベント検出部 1 2 2 および画像イベント検出部 1 1 2 から入力するイベント情報に類似する信頼度の高いデータが残存することになり、最終的に信頼度の高い以下の各情報、すなわち、

40

( a ) 複数のユーザが、それぞれどこにいて、それらは誰であるかの推定情報としての [ ターゲット情報 ] 、

( b ) 例えば話をしたユーザなどのイベント発生源を示す [ シグナル情報 ] 、

これらが生成されて処理決定部 1 3 2 に出力される。

50

## 【 0 1 4 8 】

## [ ターゲットの生成および削除 ]

上述した実施例において、音声・画像統合処理部 1 3 1 では、予め  $m$  個のパーティクルにそれぞれ  $n$  個のターゲットを設定して処理を行う構成を説明したが、ターゲットの数は、適宜変更する設定としてよい、すなわち、必要に応じて、新たなターゲットの生成や、ターゲットの削除を行う構成としてもよい。

## 【 0 1 4 9 】

## ( ターゲットの生成 )

まず、音声・画像統合処理部 1 3 1 における新たなターゲットの生成処理について、図 1 1 を参照して説明する。新たなターゲットの生成は、例えば各パーティクルに対するイベント発生源仮説の設定時に行う。

10

## 【 0 1 5 0 】

イベントと既存の  $n$  個の各ターゲットとのイベント - ターゲット間尤度を計算する際、暫定的に  $n + 1$  番目のターゲットとして図 1 1 に示すような「位置情報」、「識別情報」に一樣分布（「分散が十分大きいガウス分布」と「全  $P_t[i]$  が等しい  $UserID$  分布」）に設定した新たな暫定新規ターゲット 5 0 1 を生成する。

## 【 0 1 5 1 】

この暫定的な新規ターゲット（ $tID = n + 1$ ）を設定した後、新たなイベントの入力に基づいて、図 7 を参照して説明したフローにおけるステップ S 1 0 2 のイベント発生源仮説の設定が行われ、この処理の際に、入力イベント情報と各ターゲット間の尤度算出が実行されて、各ターゲットのターゲット重み  $[W_{tID}]$  の算出が行われる。このとき、図 1 1 に示す暫定ターゲット（ $tID = n + 1$ ）についても、入力イベント情報との尤度算出を実行して、暫定的な  $n + 1$  番目のターゲットのターゲット重み（ $W_{n+1}$ ）を算出する。

20

## 【 0 1 5 2 】

この暫定的な  $n + 1$  番目のターゲットのターゲット重み（ $W_{n+1}$ ）が、既存の  $n$  個のターゲットのターゲット重み（ $W_1 \sim W_n$ ）より大きいと判断された場合は、その新規ターゲットを全パーティクルに対して設定する。

## 【 0 1 5 3 】

なお、例えばカメラの撮影する 1 つの画像中に複数の顔イベントがあり、1 つ 1 つの顔イベントに対して、図 7 に示すフローの処理を行う構成において、1 画像中の顔の数（＝イベント数）が、各パーティクルに設定されたターゲット数（ $n$ ）より少ない場合、 $tID = n + 1$  の暫定ターゲットの重み  $W_{n+1}$  が、他のターゲットの重み（ $W_1 \sim W_n$ ）より大きくなくても、そのまま新規ターゲットとして全パーティクルに対して生成する処理を行う構成としても良い。

30

## 【 0 1 5 4 】

なお、新規ターゲットが生成された場合、イベント発生源の仮説の生成は事前に計算したターゲット重み  $[W_{tID}]$  に基づいて確率的に行っても良いし、全てのパーティクルにおいてイベント発生源の仮説を新規ターゲットにしても良い。

## 【 0 1 5 5 】

## ( ターゲットの削除 )

次に、音声・画像統合処理部 1 3 1 におけるターゲットの削除処理について、図 1 2 を参照して説明する。ターゲットの削除は、例えば図 7 に示す処理フローにおけるステップ S 1 0 5 のターゲットデータの更新処理に際して実行する。

40

## 【 0 1 5 6 】

ステップ S 1 0 5 では、先に説明したように、ターゲットデータの更新を実行して更新されたターゲットデータと、各パーティクル重み  $[W_{pID}]$  とに基づいて、ターゲット情報を生成して、処理決定部 1 3 2 に出力する処理が行われる。例えば図 1 2 に示すターゲット情報 5 2 0 が生成される。ターゲット情報は、各ターゲット（ $tID = 1 \sim n$ ）各々の

50

- ( a ) ユーザ位置情報、
- ( b ) ユーザ確信度情報、

これらの情報を含む情報として生成される。

#### 【 0 1 5 7 】

音声・画像統合処理部 1 3 1 は、このように更新ターゲットに基づいてして生成したターゲット情報中のユーザ位置情報に着目する。ユーザ位置情報は、ガウス分布  $N(m, \sigma)$  として設定される。このガウス分布に一定のピークが検出されない場合は、特定のユーザの位置を示す有効な情報とはならない。音声・画像統合処理部 1 3 1 は、このようなピークを持たない分布データとなるターゲットを削除対象として選択する。

#### 【 0 1 5 8 】

例えば、図 1 2 に示すターゲット情報 5 2 0 には、ターゲット 1, 2, n の 3 つのターゲット情報 5 2 1, 5 2 2, 5 2 3 を示しているが、これらのターゲット情報中のユーザ位置を示すガウス分布データのピークと予め定めた閾値 5 3 1 との比較を実行し、閾値 5 3 1 以上のピークを持たないデータ、すなわち、図 1 2 の例では、ターゲット情報 5 2 3 を削除ターゲットとする。

#### 【 0 1 5 9 】

この例ではターゲット ( $tID = n$ ) が削除ターゲットとして選択され。すべてのパーティクルから削除される。このようにユーザ位置を示すガウス分布 (確率密度分布) の最大値が、削除の閾値よりも小さいときに、全パーティクルに対してそのターゲットを削除する。なお、適用する閾値は、固定値でも良いし、インタラクション対象ターゲットに関しては閾値を下げて削除されにくくするなど、ターゲット毎に変える構成としてもよい。

#### 【 0 1 6 0 】

[ 画像フレーム外に仮想ターゲットを生成する処理例 ]

上述した [ ターゲットの生成および削除 ] の説明では、新たなターゲットの生成および削除構成について説明したが、図 2 に示す画像イベント検出部 1 1 2 において、画像入力部 (カメラ) 1 1 1 から入力する画像情報、すなわちカメラの撮影している画像フレームの外にユーザが存在する場合は、そのユーザに対する画像イベントを取得できないため、画像イベントからそのターゲットを生成することはできないという問題がある。

#### 【 0 1 6 1 】

そのような状態でそのユーザが音声イベントを発生しても、そのユーザ対応のターゲットが生成されず、カメラフレーム内の他のターゲットから音声イベントが発生したと推定してしまい、この場合、誤った推定結果を生成することになる。

#### 【 0 1 6 2 】

すなわち、暫定的に  $n + 1$  番目のターゲットとして図 1 1 に示すような「位置情報」と、「識別情報」として一様分布 (「分散が十分大きいガウス分布」と「全  $P t [ i ]$  が等しい  $U s e r I D$  分布」) のデータを設定した新たな暫定新規ターゲット 5 0 1 を生成し、この暫定的な新規ターゲット ( $tID = n + 1$ ) を設定した後、新たなイベントの入力に基づいて、図 7 を参照して説明したフローにおけるステップ S 1 0 2 のイベント発生源仮説の設定が行われ、この処理の際に、入力イベント情報と各ターゲット間の尤度算出を実行して、各ターゲットのターゲット重み  $[ W_{tID} ]$  の算出を行う。

#### 【 0 1 6 3 】

このとき、図 1 1 に示す暫定ターゲット ( $tID = n + 1$ ) についても、入力イベント情報との尤度算出を実行して、暫定的な  $n + 1$  番目のターゲットのターゲット重み ( $W_{n+1}$ ) を算出する。この暫定的な  $n + 1$  番目のターゲットのターゲット重み ( $W_{n+1}$ ) が、既存の  $n$  個のターゲットのターゲット重み ( $W_1 \sim W_n$ ) より大きいと判断された場合は、その新規ターゲットを全パーティクルに対して設定する構成である。

#### 【 0 1 6 4 】

しかし、この方法を適用した場合、カメラフレーム外からの音声イベントのようにその位置情報の平均値と既に存在するターゲットの位置情報の平均値がある程度離れていても、位置情報の分散が大きい場合はガウス分布間尤度がそれほど小さくならない傾向がある

。

## 【0165】

その結果、システムがターゲットとして認識していないユーザ、すなわち、図2に示す画像入力部（カメラ）111から入力する画像フレームの外のユーザからの音声イベントであっても、「イベント」と「一様分布のターゲット（ $n+1$ ）」間の尤度が最大にならずターゲットを生成することができないため、既に存在するターゲットのみで音声イベント発生源である確率の計算を行ってしまうことがある。

## 【0166】

そこで、各パーティクルでのイベント発生源の仮説生成においてターゲットの生成を確認する際、画像フレーム外に仮想ターゲットを生成する。以下、この処理例について説明する。

10

## 【0167】

本処理例では、バックグラウンドモデル（Background Model）として画像入力部（カメラ）111から入力する画像フレーム外に仮想のターゲットを配置し、「イベント」と「既に存在するターゲット（ $1 \sim n$ ）」と一様分布のターゲット間の尤度計算に加え、画像フレーム外に仮想のターゲットを配置したバックグラウンドモデル（Background Model）の仮想ターゲットとも尤度計算を行う。なお、ユーザID間尤度の計算においては、一様分布のターゲットと同様、図11に示す「全Pt[i]が等しいUser ID分布」を持つ様のデータを用いる。

## 【0168】

20

新たなターゲットの生成は、例えば各パーティクルに対するイベント発生源仮説の設定時に行う。イベントと既存の $n$ 個の各ターゲットとのイベント-ターゲット間尤度を計算する際、暫定的に $n+1$ 番目のターゲットとして、画像フレーム外に仮想のターゲットを配置したバックグラウンドモデル（Background Model）の暫定的な仮想ターゲット（ $tID = n+1$ ）を生成する。

## 【0169】

この暫定的な新規ターゲット（ $tID = n+1$ ）を設定した後、新たなイベントの入力に基づいて、図7を参照して説明したフローにおけるステップS102のイベント発生源仮説の設定を行う。

## 【0170】

30

すなわち、入力イベント情報と各ターゲット間の尤度算出を実行して各ターゲットのターゲット重み $[W_{tID}]$ の算出を行う際に、バックグラウンドモデル（Background Model）の暫定的な仮想ターゲット（ $tID = n+1$ ）についても、入力イベント情報との尤度算出を実行して、暫定的な $n+1$ 番目のターゲットのターゲット重み（ $W_{n+1}$ ）を算出する。

## 【0171】

この暫定的な $n+1$ 番目のターゲットのターゲット重み（ $W_{n+1}$ ）が、既存の $n$ 個のターゲットのターゲット重み（ $W_1 \sim W_n$ ）より大きいと判断された場合は、その新規ターゲットを全パーティクルに対して設定する。

## 【0172】

40

図13に、画像フレーム外に仮想のターゲットを配置したバックグラウンドモデル（Background Model）を含めたイベント-ターゲット間尤度の計算例を示す。

## 【0173】

図13(a)はイベント検出を行う実環境を示している。画像入力部（カメラ）111から入力する画像情報、すなわちカメラの撮影している画像フレーム601の外に声を発したユーザ611が存在する。

## 【0174】

図13(b)は、図2に示す音声イベント検出部122において検出された音声イベント情報を示している。音声イベント検出部122は、複数の異なるポジションに配置された複数の音声入力部（マイク）121a~dから入力する音声情報を解析し、音声の発生

50



源の位置情報を確率分布データとして生成する。具体的には、音源方向に関する期待値と分散データ  $N(m_e, \sigma_e)$  を生成する。また、予め登録されたユーザの声の特徴情報との比較処理に基づいてユーザ識別情報を生成する。この識別情報も確率的な推定値として生成する。音声イベント検出部 122 には、予め検証すべき複数のユーザの声についての特徴情報が登録されており、入力音声と登録音声との比較処理を実行して、どのユーザの声である確率が高いかを判定する処理を行い、全登録ユーザに対する事後確率、あるいはスコアを算出する。

【0175】

図 13(c) は、音声画像統合処理部 131 が保持する既存の  $n$  個のターゲット ( $tID = 1 \sim n$ ) と、暫定的に  $n + 1$  番目のターゲットとして生成した ( $X$ ), ( $Y$ ), ( $Z$ ) の 3 つのターゲットを示している。

10

【0176】

ターゲット ( $X$ ) は、先に図 11 を参照して説明した暫定的な新規ターゲット ( $tID = n + 1$ ) であり、「位置情報」、「識別情報」に一樣分布（「分散が十分大きいガウス分布」と「全  $P_t[i]$  が等しい  $UserID$  分布」）に設定した新たな暫定新規ターゲットである。

【0177】

ターゲット ( $Y$ ), ( $Z$ ) は、上述したバックグラウンドモデル (Background Model) のターゲットであり、画像フレーム外に仮想のターゲットを配置した新規ターゲット ( $ID = n + 1$ ) である。ターゲット ( $Y$ ) は、「位置情報」が、画像フレームの外の左側の位置に高い存在確率を持つ情報であり、「識別情報」は、( $X$ ) の一樣分布ターゲットと同様、「全  $P_t[i]$  が等しい  $UserID$  分布」を持つ一様のデータである。

20

【0178】

ターゲット ( $Z$ ) は、「位置情報」が、画像フレームの外の右側の位置に高い存在確率を持つ情報であり、「識別情報」は、( $X$ ) の一樣分布ターゲットと同様、「全  $P_t[i]$  が等しい  $UserID$  分布」を持つ一様のデータである。

【0179】

これらの暫定的な新規ターゲット ( $tID = n + 1$ ) を設定した後、新たなイベントの入力に基づいて、図 7 を参照して説明したフローにおけるステップ S102 のイベント発生源仮説の設定が行われ、この処理の際に、入力イベント情報と各ターゲット間の尤度算出が実行されて、各ターゲットのターゲット重み [ $W_{tID}$ ] の算出が行われる。このとき、図 13 に示す 3 つの暫定ターゲット ( $X$ ), ( $Y$ ), ( $Z$ ) についても、入力イベント情報との尤度算出を実行して、暫定的な  $n + 1$  番目のターゲットとしてのターゲット重み ( $W_{n+1}$ ) を算出する。

30

【0180】

この暫定的な  $n + 1$  番目のターゲット ( $X$ ), ( $Y$ ), ( $Z$ ) のいずれかのターゲット重み ( $W_{n+1}$ ) が、既存の  $n$  個のターゲットのターゲット重み ( $W_1 \sim W_n$ ) より大きいと判断された場合は、その新規ターゲットを全パーティクルに対して設定する。

【0181】

ターゲット重みの算出例を図 14 に示す。ターゲット重みは、図 14 の右端に示すように、各パーティクルに設定されるターゲット  $1 \sim n$  の各々に対応する  $n$  個のターゲット重みの算出処理として実行される。この  $n$  個のターゲット重みの算出処理に際しては、先に図 8 を参照して説明したように、まず、入力イベント情報、すなわち、音声・画像統合処理部 131 が、音声イベント検出部 122 および画像イベント検出部 112 から入力したイベント情報と、各パーティクルの各ターゲットデータとの類似度の指標値としての尤度算出を行う。

40

【0182】

図 8 を参照して説明したように、尤度算出処理は、

(a) ユーザ位置情報についてのイベントと、ターゲットデータとの類似度データとし

50

てのガウス分布間尤度 [ D L ]、

( b ) ユーザ識別情報 ( 顔識別情報または話者識別情報 ) についてのイベントと、ターゲットデータとの類似度データとしてのユーザ確信度情報 ( u I D ) 間尤度 [ U L ]

これらを個別に算出する。

#### 【 0 1 8 3 】

次に、入力イベント情報とあるパーティクル ( p I D ) が持つ 1 つのターゲット ( t I D ) との類似度の指標としてのイベント - ターゲット間尤度 [ L<sub>p I D , t I D</sub> ] は、上記の 2 つの尤度、すなわち、

ガウス分布間尤度 [ D L ] と、

ユーザ確信度情報 ( u I D ) 間尤度 [ U L ]

10

これら 2 つの尤度を利用して算出する。すなわち重み ( = 0 ~ 1 ) を用いて、イベント - ターゲット間尤度 [ L<sub>p I D , t I D</sub> ] は下式によって算出する。

$$[ L_{p I D , t I D} ] = U L \times D L^{1 -}$$

としてイベントとターゲットとの類似度の指標であるイベント - ターゲット間尤度 [ L<sub>p I D , t I D</sub> ] を算出する。

ただし、 = 0 ~ 1 とする。

#### 【 0 1 8 4 】

このイベント - ターゲット間尤度 [ L<sub>p I D , t I D</sub> ] を、各パーティクルの各ターゲットについて各々算出し、このイベント - ターゲット間尤度 [ L<sub>p I D , t I D</sub> ] に基づいて各ターゲットのターゲット重み [ W<sub>t I D</sub> ] を算出する。

20

#### 【 0 1 8 5 】

イベント - ターゲット間尤度 [ L<sub>p I D , t I D</sub> ] に基づく、ターゲット重み [ W<sub>t I D</sub> ] の算出式は、先に説明した通り、以下の算出式である。

#### 【 数 4 】

$$W_{tID} = \sum_{pID}^m W_{pID} L_{pID,tID}$$

30

#### 【 0 1 8 6 】

図 1 4 に示すターゲット重みの算出例において、上段に記載の W<sub>1</sub> ~ W<sub>n</sub> は、すでに設定済みのターゲットについて算出したイベント - ターゲット間尤度である。下段の ( X ) , ( Y ) , ( Z ) として示す 3 つの W<sub>n + 1</sub> は、図 1 3 を参照して説明した暫定的な新規ターゲット ( t I D = n + 1 ) に対応するイベント - ターゲット間尤度である。

#### 【 0 1 8 7 】

すなわち、( X ) は、「位置情報」、「識別情報」に一様分布 ( 「分散が十分大きいガウス分布」と「全 P t [ i ] が等しい U s e r I D 分布」 ) に設定した新たな暫定新規ターゲット、( Y ) , ( Z ) は、上述したバックグラウンドモデル ( B a c k g r o u n d M o d e l ) のターゲットであり、画像フレーム外に仮想のターゲットを配置した新規ターゲット ( I D = n + 1 ) であり、これらに対応するイベント - ターゲット間尤度も算出する。

40

#### 【 0 1 8 8 】

この暫定的な n + 1 番目のターゲット ( X ) , ( Y ) , ( Z ) のいずれかのターゲット重み ( W<sub>n + 1</sub> ) が、既存の n 個のターゲットのターゲット重み ( W<sub>1</sub> ~ W<sub>n</sub> ) より大きいと判断された場合は、その新規ターゲットを全パーティクルに対して設定する。

#### 【 0 1 8 9 】

なお、ターゲット生成確認時に用いたカメラフレーム外に仮想のターゲットは、他の処理では用いない。この処理例に従えば、カメラによって撮影された画像フレーム外のユーザからの音声イベントに対して、各ターゲットがイベント発生源である確率推定の性能が

50

向上する。

【0190】

このように、本処理例では、図1に示す情報処理装置100の音声・画像統合処理部131がイベント検出部112、122の生成するイベント情報を入力し、仮想的なユーザに対応する複数のターゲットを設定した複数のパーティクルを適用したパーティクルフィルタリング処理を実行して実空間に存在するユーザのユーザ位置情報およびユーザ識別情報を含む解析情報を生成する構成を有し、カメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとイベント検出部112、122の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント・ターゲット間尤度より大きい値である場合に、暫定ターゲットを各パーティクルに新規追加する処理を行う。

10

【0191】

また、音声・画像統合処理部131は、暫定ターゲットとして、図13、図14を参照して説明したように、

(X)均一データによって構成されるユーザ位置情報、ユーザ識別情報を持つ暫定ターゲット

(Y)、(Z)画像フレームの異なる方向のフレーム外部位置に仮想ターゲットを設定した複数の異なる暫定ターゲット

これらの異なるタイプの暫定ターゲットを生成し、生成した複数の暫定ターゲットとイベント情報との尤度を個別に算出して、算出した暫定ターゲットのイベント・ターゲット間尤度の最大値が、既存ターゲットに対応するイベント・ターゲット間尤度より大きい値を有する場合に、その最大値に対応する暫定ターゲットを各パーティクルに新規追加する処理を行う。本構成により、カメラの取得する画像フレームの外部にいるユーザからの音声入力イベントに対応した正しい推定処理が可能となり、ユーザ位置やユーザ識別情報を効率的に確実に生成することが可能となる。

20

【0192】

以上、特定の実施例を参照しながら、本発明について詳解してきた。しかしながら、本発明の要旨を逸脱しない範囲で当業者が実施例の修正や代用を成し得ることは自明である。すなわち、例示という形態で本発明を開示してきたのであり、限定的に解釈されるべきではない。本発明の要旨を判断するためには、特許請求の範囲の欄を参酌すべきである。

30

【0193】

また、明細書中において説明した一連の処理はハードウェア、またはソフトウェア、あるいは両者の複合構成によって実行することが可能である。ソフトウェアによる処理を実行する場合は、処理シーケンスを記録したプログラムを、専用のハードウェアに組み込まれたコンピュータ内のメモリにインストールして実行させるか、あるいは、各種処理が実行可能な汎用コンピュータにプログラムをインストールして実行させることが可能である。例えば、プログラムは記録媒体に予め記録しておくことができる。記録媒体からコンピュータにインストールする他、LAN(Local Area Network)、インターネットといったネットワークを介してプログラムを受信し、内蔵するハードディスク等の記録媒体にインストールすることができる。

40

【0194】

なお、明細書に記載された各種の処理は、記載に従って時系列に実行されるのみならず、処理を実行する装置の処理能力あるいは必要に応じて並列的にあるいは個別に実行されてもよい。また、本明細書においてシステムとは、複数の装置の論理的集合構成であり、各構成の装置が同一筐体内にあるものには限らない。

【産業上の利用可能性】

【0195】

以上、説明したように、本発明の一実施例の構成によれば、カメラやマイクによって取得される画像情報や音声情報に基づいてユーザの推定位置および推定識別データを含むイベント情報を入力して、複数のターゲットを設定した複数のパーティクルを適用したパー

50

ティクルフィルタリング処理を行い、フィルタリングによる仮説の更新および取捨選択に基づいてユーザの位置および識別情報を生成する。また、カメラの取得する画像フレームの外部に仮想ターゲットを設定した暫定ターゲットとイベント検出部の生成するイベント情報との尤度が、画像フレームの内部にターゲットを設定した既存ターゲットに対応するイベント・ターゲット間尤度より大きい値である場合に、暫定ターゲットを各パーティクルに新規追加する処理を行う。本構成により、カメラの取得する画像フレームの外部にいるユーザからの音声入力イベントに対応した正しい推定処理が可能となり、ユーザ位置やユーザ識別情報を効率的に確実に生成することが可能となる。

【図面の簡単な説明】

【0196】

10

【図1】本発明に係る情報処理装置の実行する処理の概要について説明する図である。

【図2】本発明の一実施例の情報処理装置の構成および処理について説明する図である。

【図3】音声イベント検出部122および画像イベント検出部112が生成し音声・画像統合処理部131に入力する情報の例について説明する図である。

【図4】パーティクル・フィルタ(Particle Filter)を適用した基本的な処理例について説明する図である。

【図5】本処理例で設定するパーティクルの構成について説明する図である。

【図6】各パーティクルに含まれるターゲット各々が有するターゲットデータの構成について説明する図である。

【図7】音声・画像統合処理部131の実行する処理シーケンスを説明するフローチャートを示す図である。

20

【図8】ターゲット重み $[W_{tID}]$ の算出処理の詳細について説明する図である。

【図9】パーティクル重み $[W_{pID}]$ の算出処理の詳細について説明する図である。

【図10】パーティクル重み $[W_{pID}]$ の算出処理の詳細について説明する図である。

【図11】音声・画像統合処理部131における新たなターゲットの生成処理について説明する図である。

【図12】音声・画像統合処理部131におけるターゲットの削除処理について説明する図である。

【図13】画像フレーム外に仮想ターゲットを生成する処理例について説明する図である。

30

【図14】画像フレーム外に仮想ターゲットを生成する処理におけるイベント・ターゲット間尤度の算出処理例について説明する図である。

【符号の説明】

【0197】

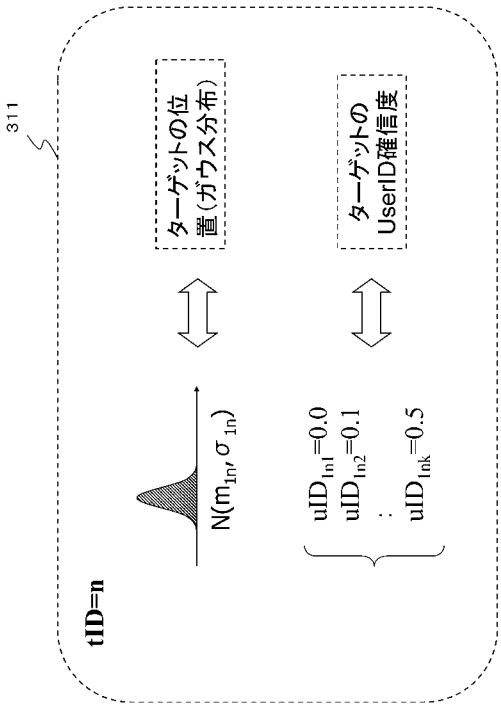
- 11 ~ 14 ユーザ
- 21 カメラ
- 31 ~ 34 マイク
- 100 情報処理装置
- 111 画像入力部
- 112 画像イベント検出部
- 121 音声入力部
- 122 音声イベント検出部
- 131 音声・画像統合処理部
- 132 処理決定部
- 201 ~ 20k ユーザ
- 301 ユーザ
- 302 画像データ
- 305 ターゲット情報
- 311 ターゲットデータ
- 401 イベント情報

40

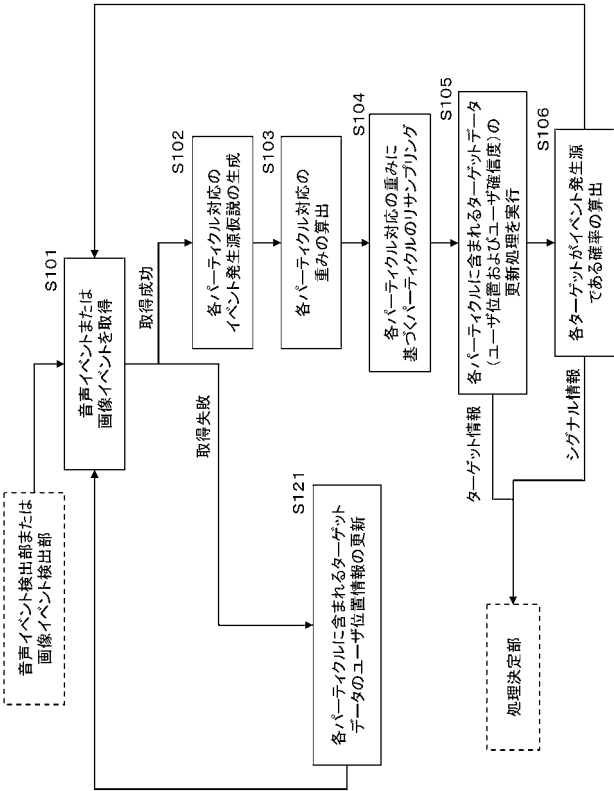
50

- 4 1 1 ~ 4 1 3    パーティクル
- 4 2 1 ~ 4 2 3    ターゲット
- 5 0 1    暫定新規ターゲット
- 5 2 0    ターゲット情報
- 5 2 1 ~ 5 2 3    ターゲット情報
- 5 3 1    閾値
- 6 0 1    画像フレーム
- 6 1 1    ユーザ

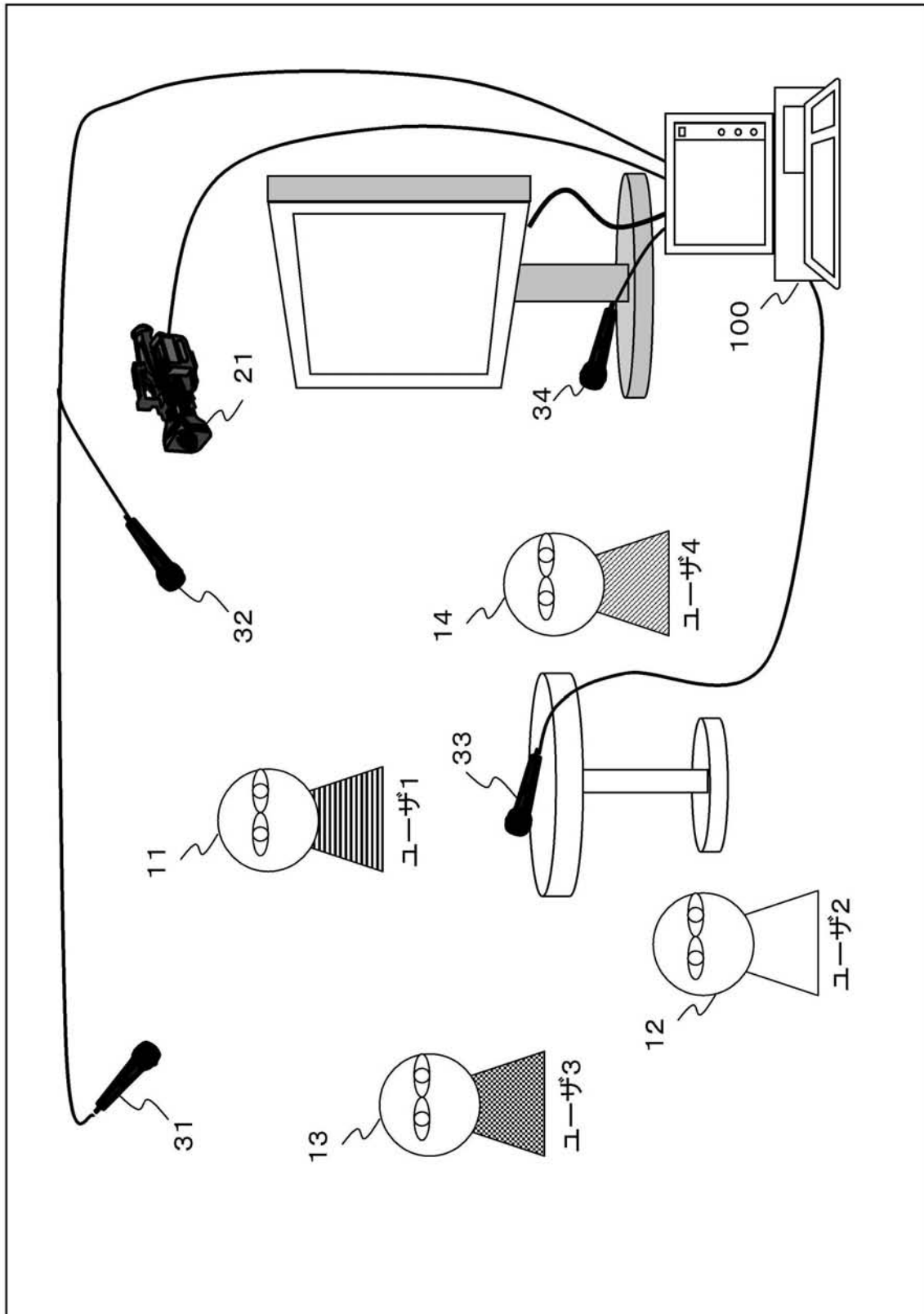
【 図 6 】



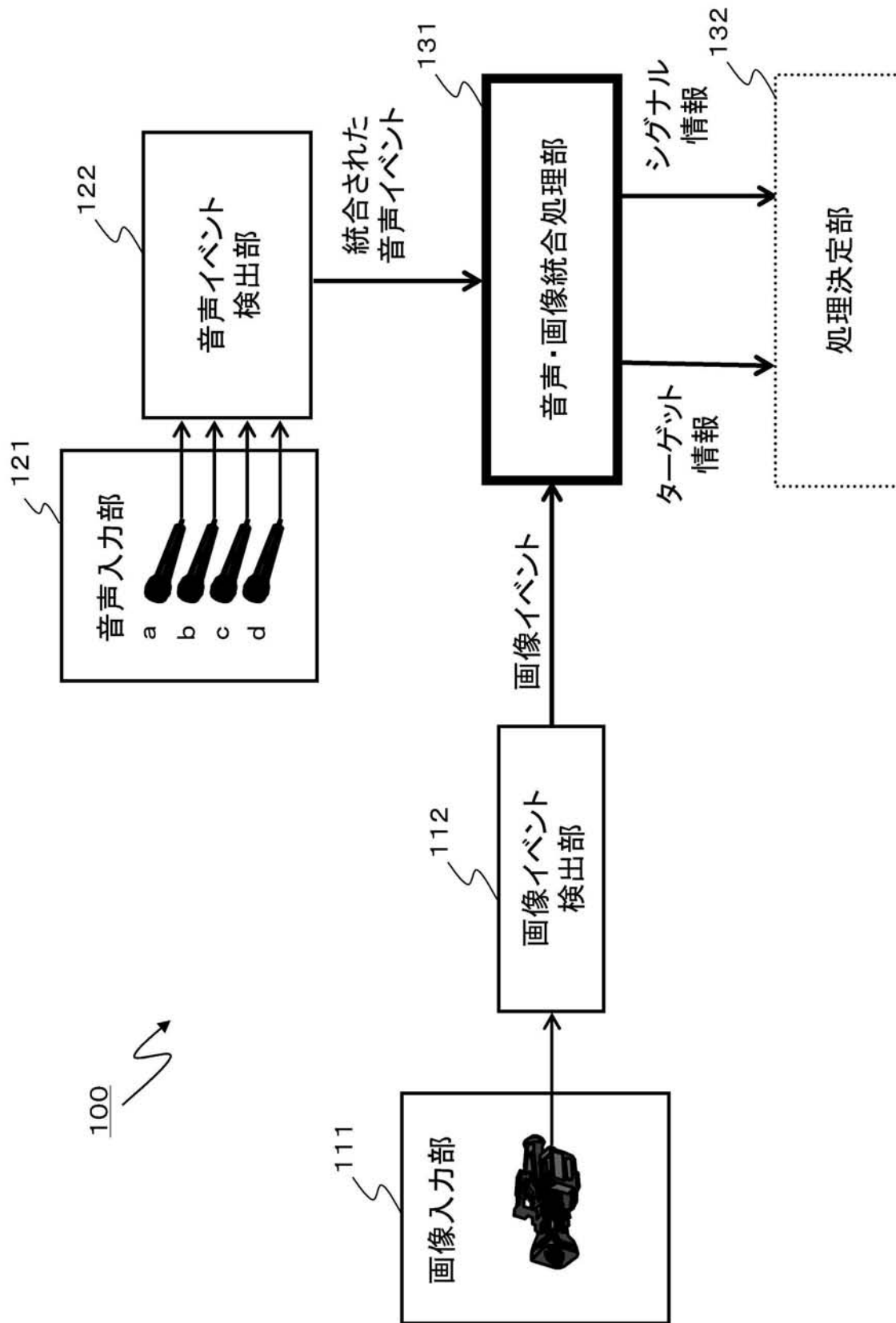
【 図 7 】



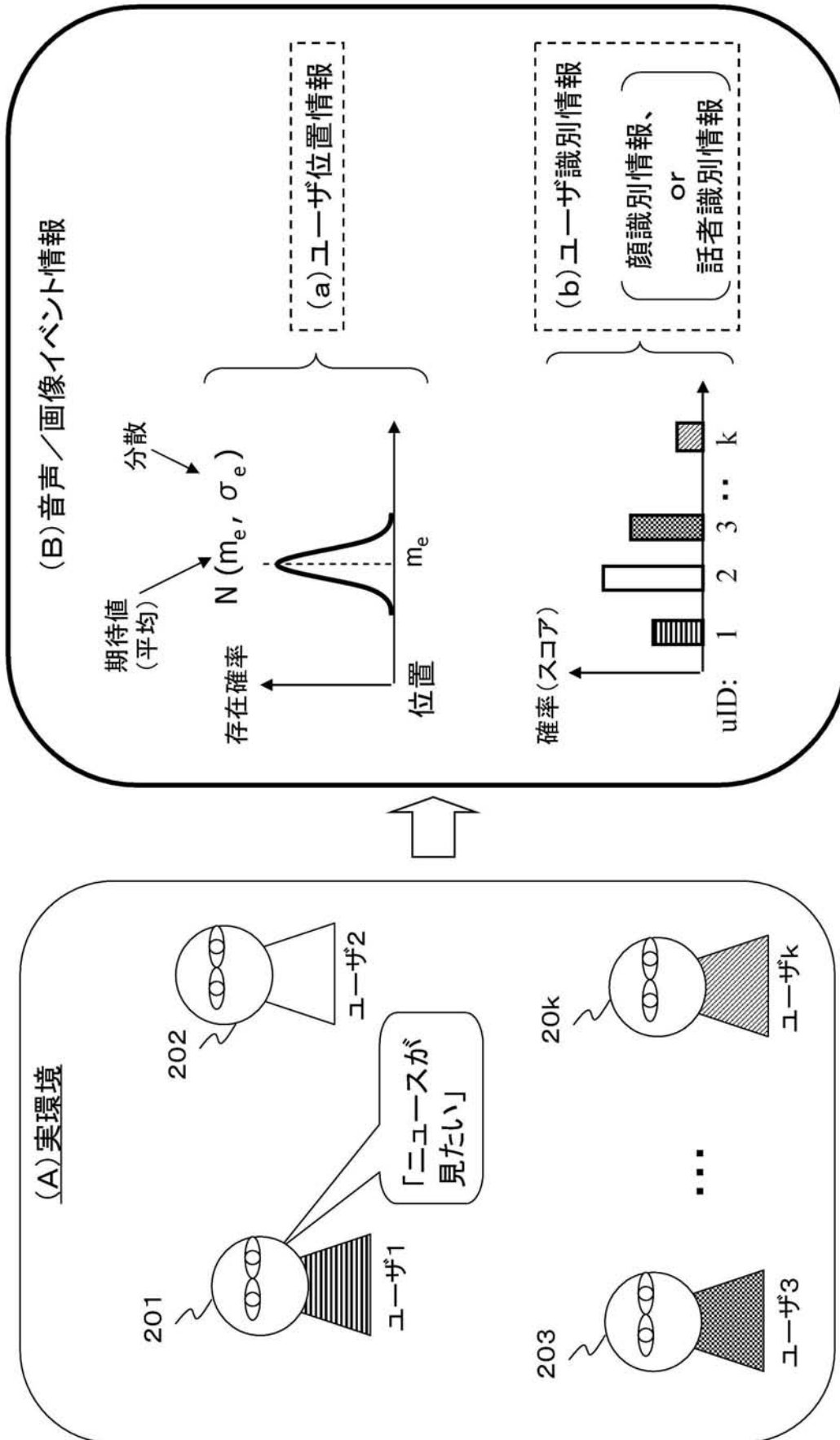
【図 1】



【図 2】

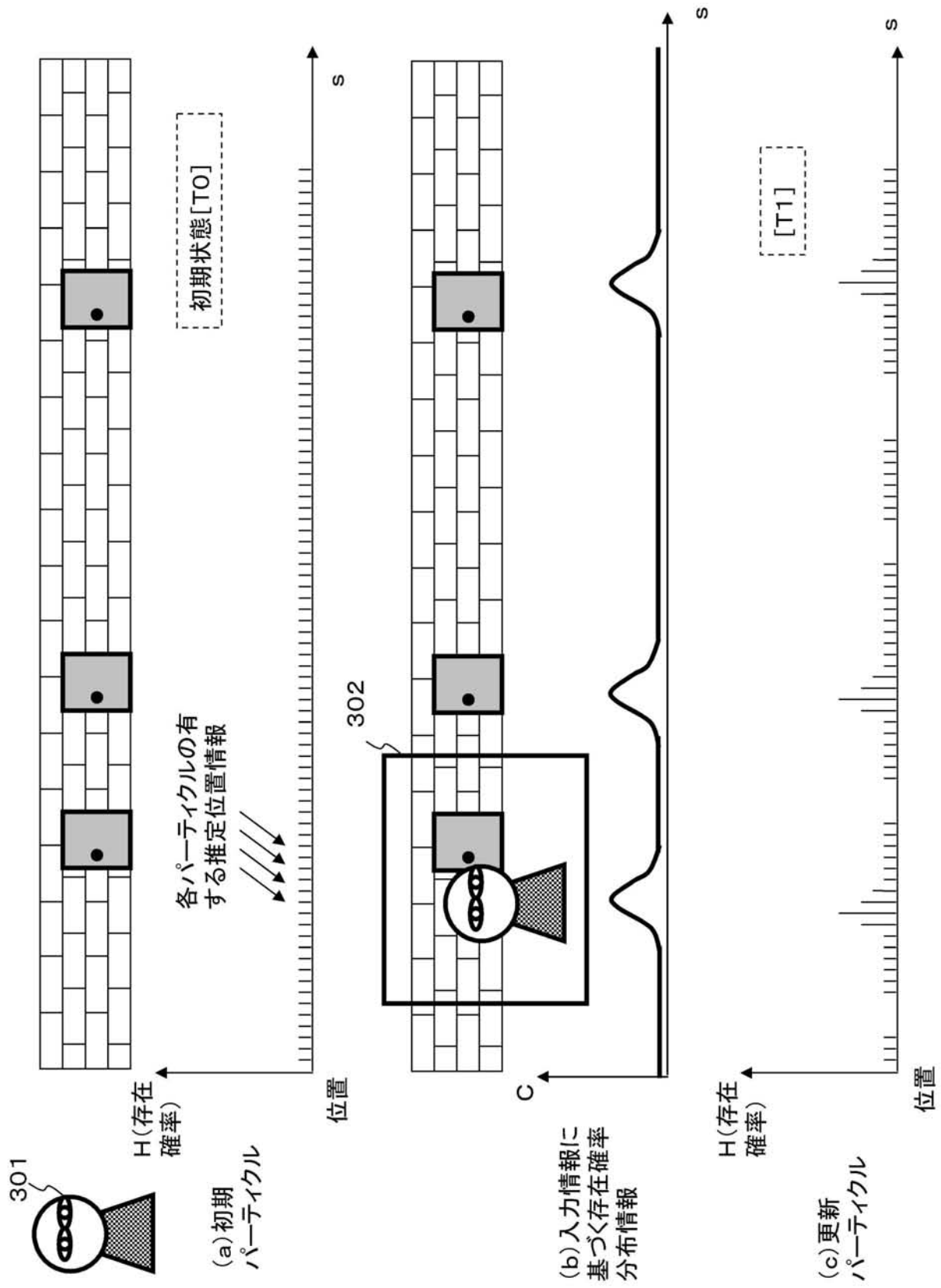


【図 3】



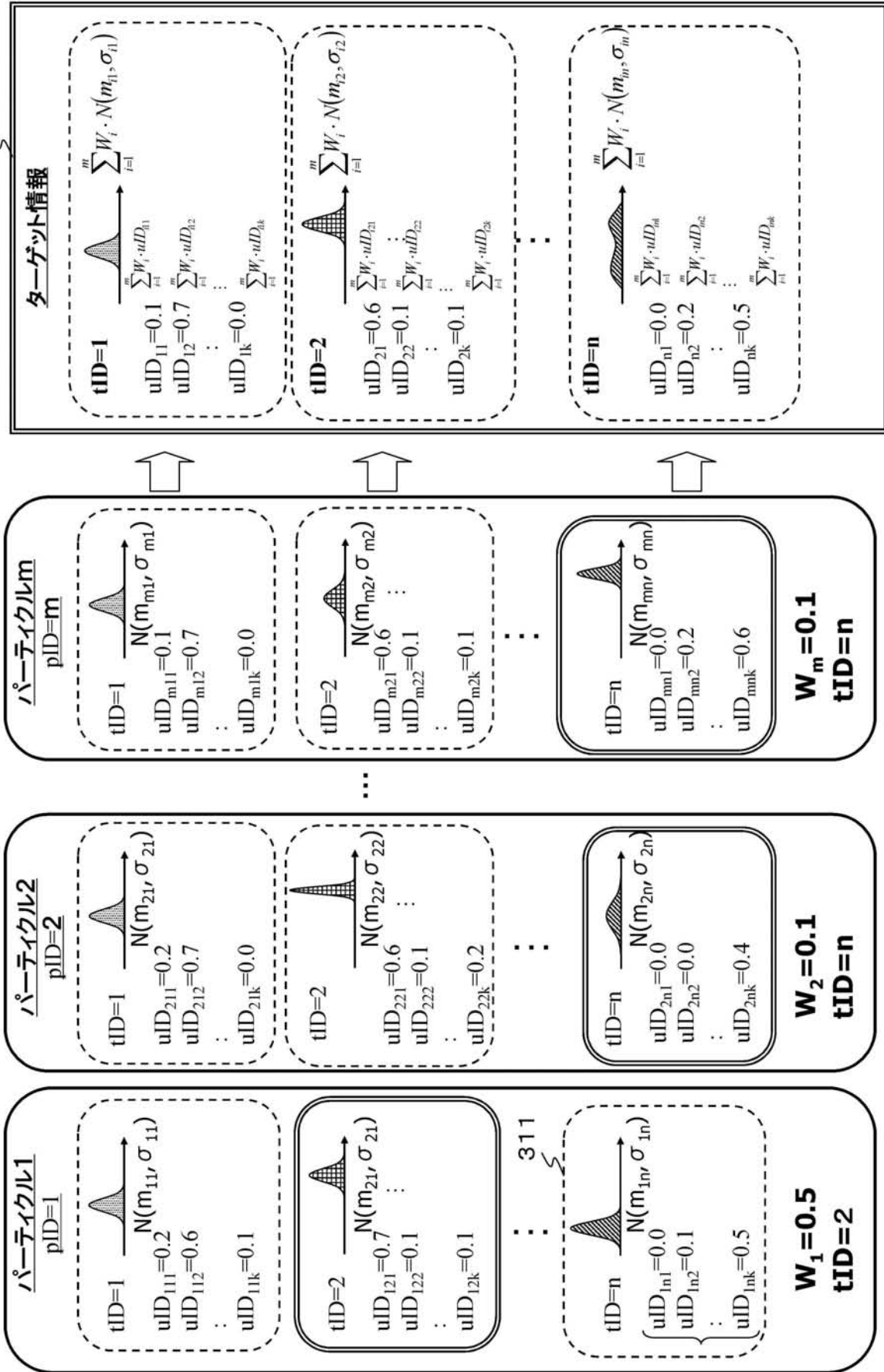


【 図 4 】

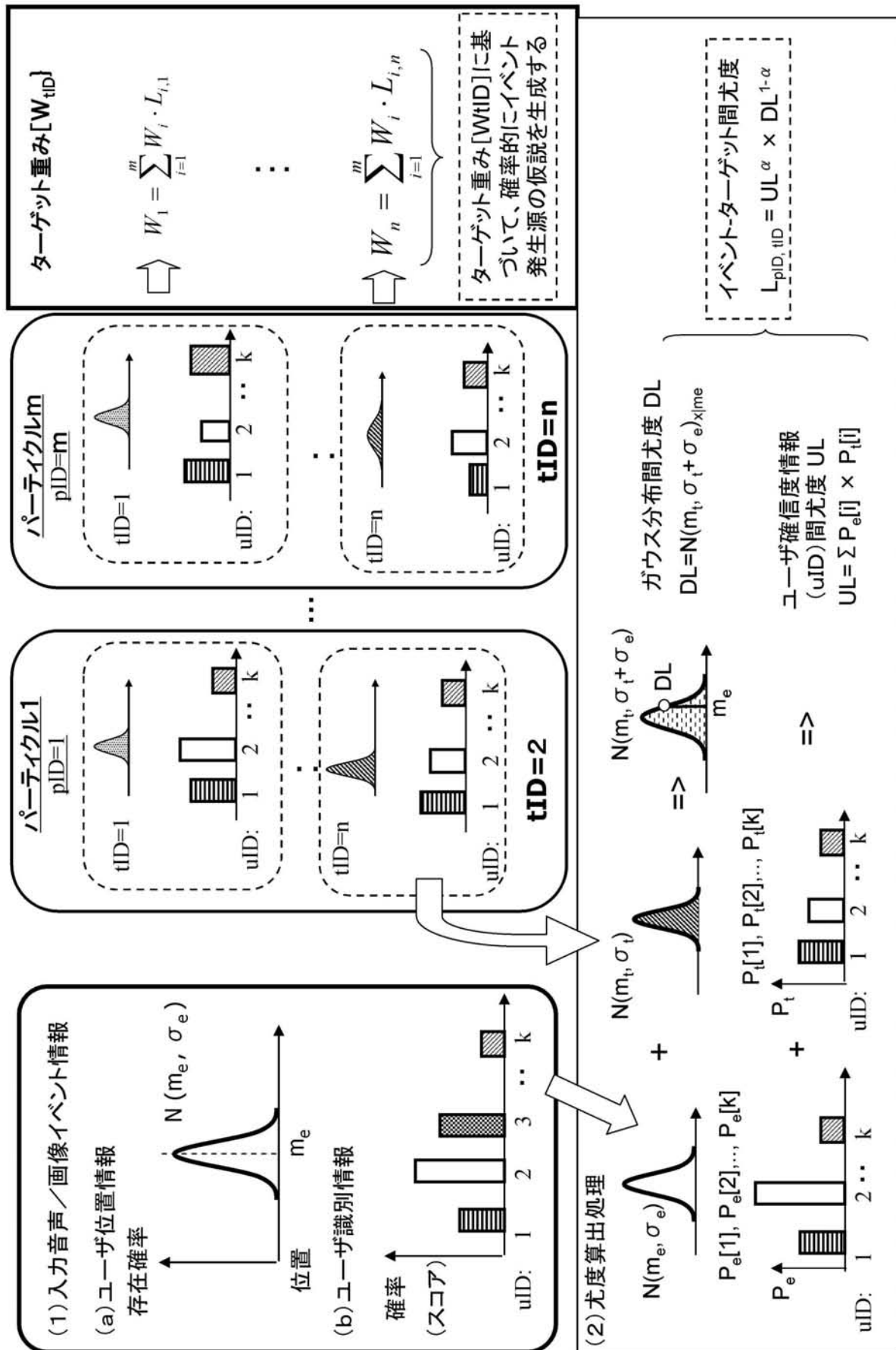


【図5】

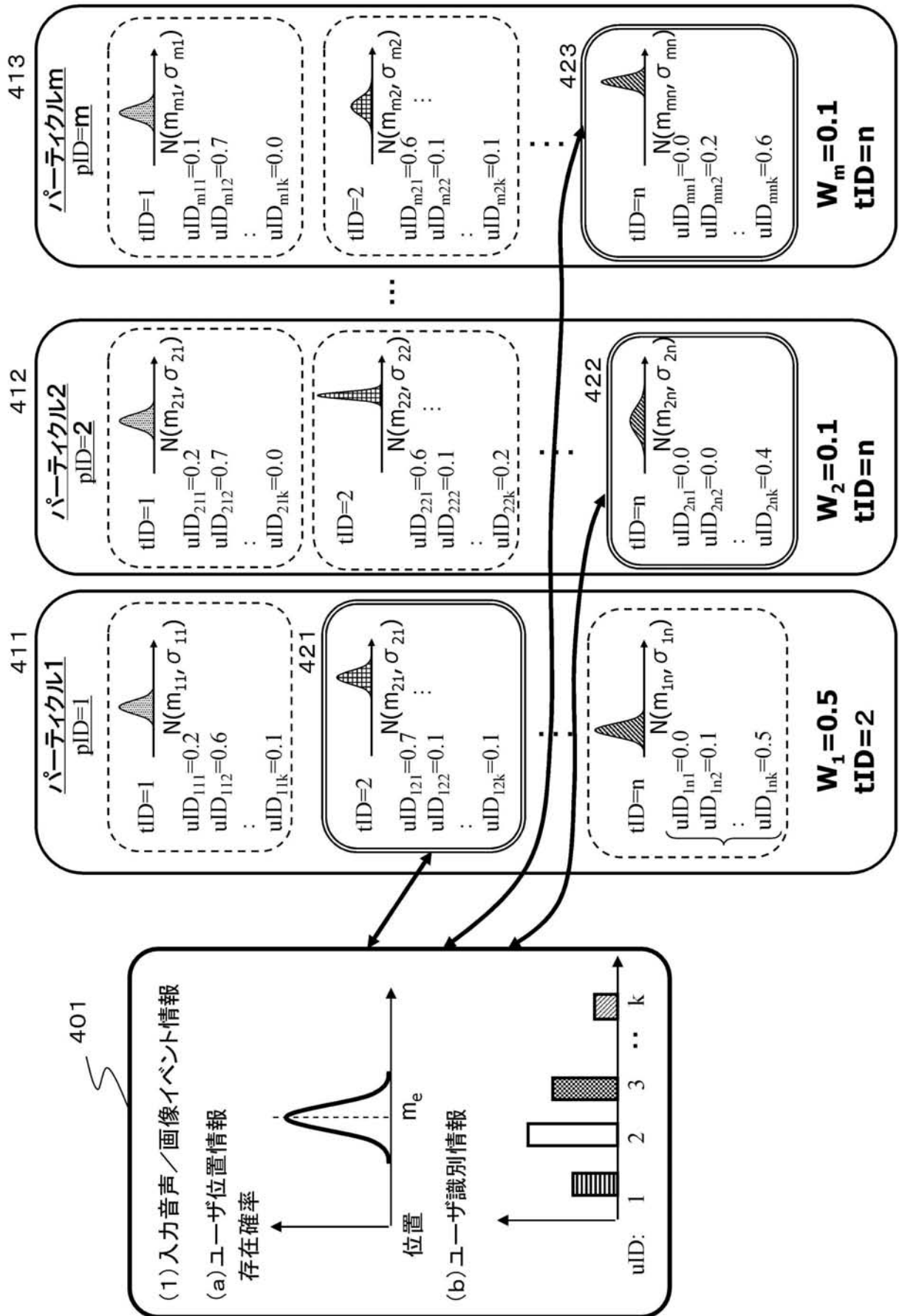
305



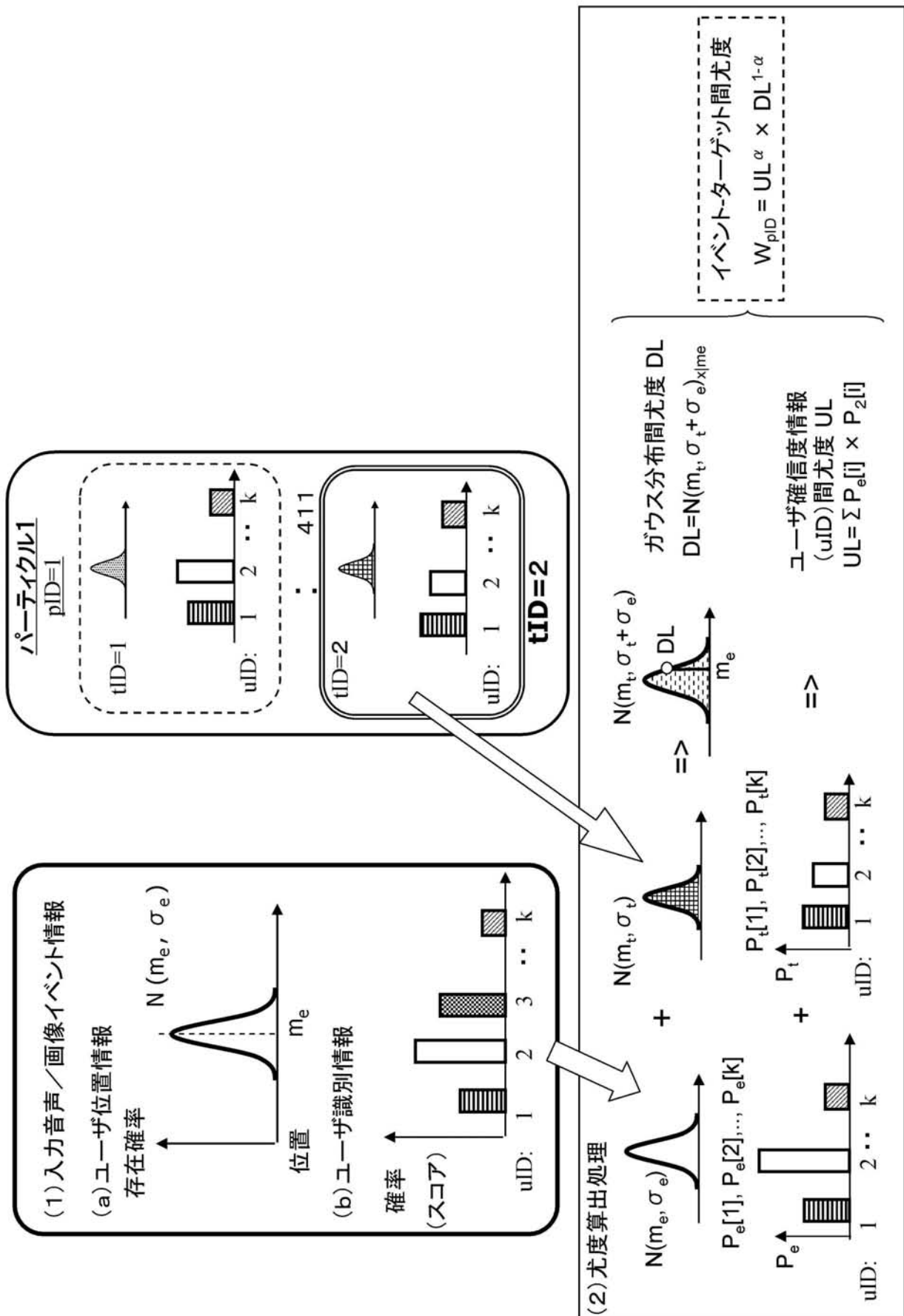
【図 8】



【図 9】



【図 10】



【図 11】

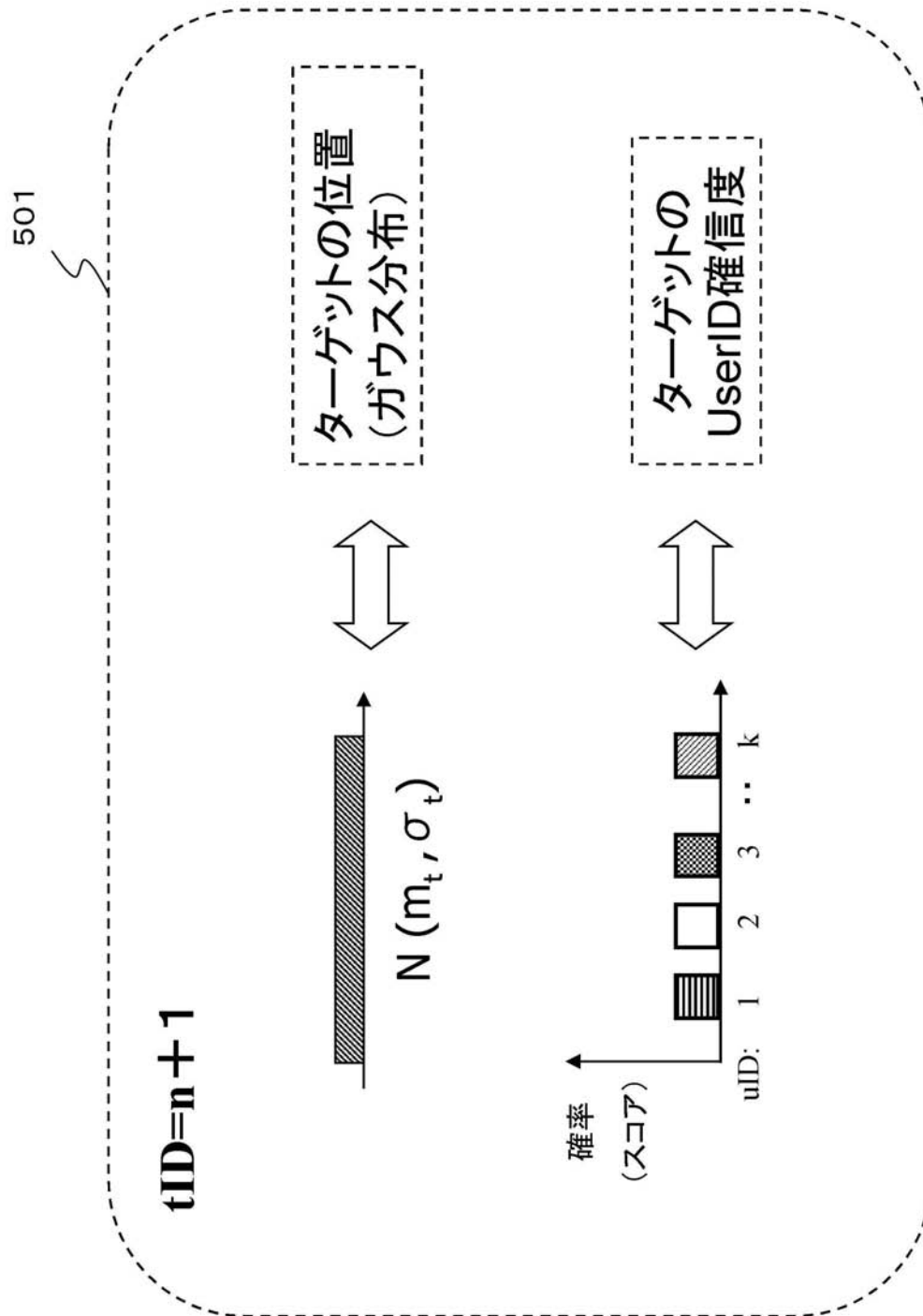








Figure 1 illustrates the generation of target evaluation values for a target generation system. The diagram is organized into three main columns, each representing a different particle type (1, 2, and m), and a central section for the target generation process.

**Particle Type 1 (パーティクル1):** The initial state (tID=1) is shown with a normal distribution  $N(m_{11}, \sigma_{11})$  and a histogram of target evaluation values  $P_t$  for  $ulD: 1, 2, \dots, k$ . The process evolves through intermediate states (tID=n+1) to the final state (tID=n+1), where the target evaluation values are updated.

**Particle Type 2 (パーティクル2):** The initial state (tID=1) is shown with a normal distribution  $N(m_{21}, \sigma_{21})$  and a histogram of target evaluation values  $P_t$  for  $ulD: 1, 2, \dots, k$ . The process evolves through intermediate states (tID=n+1) to the final state (tID=n+1), where the target evaluation values are updated.

**Particle Type m (パーティクルm):** The initial state (tID=1) is shown with a normal distribution  $N(m_{m1}, \sigma_{m1})$  and a histogram of target evaluation values  $P_t$  for  $ulD: 1, 2, \dots, k$ . The process evolves through intermediate states (tID=n+1) to the final state (tID=n+1), where the target evaluation values are updated.

**Target Generation Process (ターゲット生成評価用):** The process involves generating target evaluation values for a target generation system. The initial state (tID=1) is shown with a normal distribution  $N(m_{11}, \sigma_{11})$  and a histogram of target evaluation values  $P_t$  for  $ulD: 1, 2, \dots, k$ . The process evolves through intermediate states (tID=n+1) to the final state (tID=n+1), where the target evaluation values are updated.

**Target Evaluation Values (ターゲット評価値):** The target evaluation values are generated for a target generation system. The initial state (tID=1) is shown with a normal distribution  $N(m_{11}, \sigma_{11})$  and a histogram of target evaluation values  $P_t$  for  $ulD: 1, 2, \dots, k$ . The process evolves through intermediate states (tID=n+1) to the final state (tID=n+1), where the target evaluation values are updated.

**Target Generation System (ターゲット生成システム):** The target generation system is used to generate target evaluation values. The initial state (tID=1) is shown with a normal distribution  $N(m_{11}, \sigma_{11})$  and a histogram of target evaluation values  $P_t$  for  $ulD: 1, 2, \dots, k$ . The process evolves through intermediate states (tID=n+1) to the final state (tID=n+1), where the target evaluation values are updated.

---

フロントページの続き

(51)Int.Cl.

F I

テーマコード(参考)

B 2 5 J 19/02

F ターム(参考) 3C007 AS36 KS11 KS13 KS39 KT01 LW03 WB13 WB14 WB17  
5B087 AA09 CC33  
5D015 AA03  
5E501 AA01 AC37 BA05 CB14 CB15