

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6309615号  
(P6309615)

(45) 発行日 平成30年4月11日 (2018. 4. 11)

(24) 登録日 平成30年3月23日 (2018. 3. 23)

(51) Int. Cl. F I  
G 1 O L 15/10 (2006.01) G 1 O L 15/10 2 O O W

請求項の数 15 (全 33 頁)

(21) 出願番号	特願2016-512922 (P2016-512922)	(73) 特許権者	595020643
(86) (22) 出願日	平成26年4月24日 (2014. 4. 24)		クアルコム・インコーポレイテッド
(65) 公表番号	特表2016-526178 (P2016-526178A)		QUALCOMM INCORPORATED
(43) 公表日	平成28年9月1日 (2016. 9. 1)		ED
(86) 国際出願番号	PCT/US2014/035247		アメリカ合衆国、カリフォルニア州 92
(87) 国際公開番号	W02014/182460		121-1714、サン・ディエゴ、モア
(87) 国際公開日	平成26年11月13日 (2014. 11. 13)		ハウス・ドライブ 5775
審査請求日	平成29年3月29日 (2017. 3. 29)	(74) 代理人	100108855
(31) 優先権主張番号	61/820, 498		弁理士 蔵田 昌俊
(32) 優先日	平成25年5月7日 (2013. 5. 7)	(74) 代理人	100109830
(33) 優先権主張国	米国 (US)		弁理士 福原 淑弘
(31) 優先権主張番号	61/859, 058	(74) 代理人	100158805
(32) 優先日	平成25年7月26日 (2013. 7. 26)		弁理士 井関 守三
(33) 優先権主張国	米国 (US)	(74) 代理人	100194814
			弁理士 奥村 元宏

最終頁に続く

(54) 【発明の名称】 ターゲットキーワードを検出するための方法および装置

(57) 【特許請求の範囲】

【請求項 1】

ターゲットキーワードを検出するための方法であって、前記ターゲットキーワードが、冒頭の部分と複数の後続の部分とを含み、前記方法が、

電子デバイスにおいて、前記ターゲットキーワードの前記複数の後続の部分のうちの1つから始まる入力音声に基づいて入力音声ストリームを生成することと、前記入力音声ストリームは、デューティサイクルに従って時間期間の間に生成される、

前記入力音声ストリームに基づいて前記入力音声に関連付けられた音声特徴を決定することと、

状態ネットワークを記述しているデータを取得することと、ここにおいて、前記状態ネットワークは、開始状態と、複数のエントリ状態と、前記開始状態から前記複数のエントリ状態の各々への遷移とを含む、

前記音声特徴に基づいておよび前記データにさらに基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することと

を備える、方法。

【請求項 2】

前記入力音声の前記ターゲットキーワードに対応すると決定することに応答して、

前記電子デバイスのボイスアシスタントを起動することと、

前記ボイスアシスタントを使用し、前記電子デバイスにおいてメッセージを生成することと、

10

20

前記電子デバイスの１つまたは複数の機能の起動を示すユーザ入力を受け取ることとをさらに備える、請求項１に記載の方法。

【請求項３】

前記複数のエントリ状態に対応するデータが、

前記ターゲットキーワードの前記冒頭の部分と前記複数の後続の部分とに対応する基準入力音声を受け取ることと、

前記基準入力音声に対する複数の基準状態シーケンスを決定することと、

前記複数の基準状態シーケンスにおける複数の状態に対する状態時間期間を決定することと、

前記複数のエントリ状態を決定することと

10

によって前記電子デバイスに記憶され、

前記複数の基準状態シーケンスにおける前記複数の状態に対する前記状態時間期間が、前記複数の基準状態シーケンスをバックトラックすることによって決定される、請求項１に記載の方法。

【請求項４】

前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、複数のキーワードスコアを決定することを備え、各キーワードスコアが、前記開始状態から前記複数のエントリ状態のうちの１つへの遷移を含むそれぞれの状態シーケンスに対応する、請求項１に記載の方法。

【請求項５】

20

前記状態ネットワークが、複数の状態シーケンスを含み、各状態シーケンスが、前記開始状態と、前記複数のエントリ状態のうちの１つと、１つまたは複数の後続の状態とを含む状態を備え、前記複数の状態シーケンスの各状態シーケンスが、隠れマルコフモデルと、前記状態シーケンスの前記複数の状態についての遷移情報とに関連付けられる、請求項４に記載の方法。

【請求項６】

前記複数のキーワードスコアを決定することが、

前記状態ネットワークに基づいて、前記音声特徴の各々についての前記状態の各々の観測スコアを決定することと、

前記状態ネットワークの遷移情報に基づいて、前記複数の状態シーケンスの各々における前記状態の各々から次の状態への遷移スコアを取得することと

30

を備え、

前記複数のキーワードスコアが、前記観測スコアと前記遷移スコアとに基づいて決定される、請求項５に記載の方法。

【請求項７】

前記複数のキーワードスコアの中の最も大きいキーワードスコアが、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定するために使用され、前記最も大きいキーワードスコアがしきい値スコアよりも大きい場合、前記入力音声は、前記ターゲットキーワードに対応すると決定される、請求項４に記載の方法。

【請求項８】

40

前記状態ネットワークが非キーワード状態シーケンスを含み、前記複数のキーワードスコアを決定することが、前記非キーワード状態シーケンスについての非キーワードスコアを決定することを備え、

前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、前記複数のキーワードスコアの中から最も大きいキーワードスコアを選択することと、前記最も大きいキーワードスコアと前記非キーワードスコアとの間の差に基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することとをさらに備える、

前記差に基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、

前記差に基づいて、信頼性値を決定することと、

50

前記信頼性値がしきい値信頼性値よりも大きい場合、前記入力音声が入力ターゲットキーワードに対応すると決定することと

を備える、請求項 4 に記載の方法。

【請求項 9】

冒頭の部分と複数の後続の部分とを含むターゲットキーワードを検出するための電子デバイスであって、

前記ターゲットキーワードの前記複数の後続の部分のうちの 1 つから始まる入力音声に基づいて入力音声ストリームを生成するように構成された音声センサと、前記音声センサは、デューティサイクルに従って時間期間中に前記入力音声ストリームを生成するようにさらに構成される、

10

前記入力音声ストリームに基づいて前記入力音声に関連付けられた音声特徴を決定し、状態ネットワークを記述しているデータを取得し、前記音声特徴と前記データとに基づいて、前記入力音声が入力ターゲットキーワードに対応するかどうかを決定するように構成されたボイスアクティベーションユニットと

を備え、

前記状態ネットワークは、開始状態と、複数のエントリ状態と、前記開始状態から前記複数のエントリ状態の各々への遷移とを含む、電子デバイス。

【請求項 10】

前記デューティサイクルに基づいて、前記電子デバイスの音声センサを起動するためにアクティベーション信号を与えることと、

20

前記デューティサイクルに基づいて、前記電子デバイスの前記音声センサを非起動するためにデアクティベーション信号を与えることと、

をさらに備える、請求項 1 に記載の方法。

【請求項 11】

前記入力音声が入力ターゲットキーワードに対応すると決定することに応答して、前記電子デバイスのボイスアシスタントを起動するためにアクティベーション信号を与えることをさらに備える、請求項 10 に記載の方法。

【請求項 12】

前記時間期間は、前記デューティサイクルに関連付けられたアクティベーション時間間隔に対応し、前記アクティベーション信号は、前記アクティベーション時間間隔の間に与えられ、前記デアクティベーション信号は、前記デューティサイクルに関連付けられたデアクティベーション時間間隔の間に与えられ、前記冒頭の部分は、前記デアクティベーション時間間隔の間に話され、前記複数の後続の部分のうちの前記 1 つは、前記アクティベーション時間間隔の間に話される、請求項 10 に記載の方法。

30

【請求項 13】

前記デューティサイクルは、前記電子デバイスのデューティサイクル機能に関連付けられ、前記方法は、前記冒頭の部分がスピーチを含まないと決定することに応答して、前記デューティサイクル機能を起動するためにアクティベーション信号を与えることをさらに備える、請求項 1 に記載の方法。

【請求項 14】

40

前記時間期間は、1 つまたは複数のアクティブ時間間隔を備え、前記ターゲットキーワードの前記複数の後続の部分は、前記 1 つまたは複数のアクティブ時間間隔の間にユーザによって話される、請求項 1 に記載の方法。

【請求項 15】

前記音声センサは、前記時間期間中に前記入力音声ストリームを生成するために前記入力音声を記録するように構成されたマイクロフォンを含み、前記音声センサは、前記マイクロフォンに結合され、前記入力音声ストリームの一部が音声強度しきい値を満たすかどうかを決定するように構成された音声検出器をさらに含み、前記電子デバイスは、前記音声検出器と前記ボイスアクティベーションユニットとに結合されたスピーチ検出器をさらに備える、請求項 9 に記載の電子デバイス。

50

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

## 関連出願の相互参照

[0001]本出願は、各々の内容全体が参照により本明細書に明確に組み込まれる、同一出願人が所有する、2013年5月7日に出願された米国仮特許出願第61/820,498号、2013年7月26日に出願された米国仮特許出願第61/859,058号、および2013年11月22日に出願された米国非仮特許出願第14/087,939号の優先権を主張する。

## 【0002】

[0002]本開示は一般にオーディオ処理に関し、より詳細には、電子デバイスにおけるオーディオ入力からターゲットキーワードを検出することに関する。

## 【背景技術】

## 【0003】

[0003]近年、スマートフォン、タブレットコンピュータ、およびスマートTVなどの、音声キャプチャ機能を有する電子デバイスの使用が着実に増加している。そのような電子デバイスは、1つまたは複数のアプリケーションまたは機能がボイスキーワードに応答して起動されることを可能にし得る。そのようなデバイスのユーザは通常、ランダムな時間フレームにおいてそのようなボイスアクティベーション機能を使用するので、これらのデバイスはしばしば、そのような入力音声キャプチャされ処理されるのを可能にするために、アクティブ状態で継続的に動作する。

## 【0004】

[0004]そのようなデバイスをアクティブ状態で継続的に動作させることは、一般に、入力音声からキーワードを検出する際にキーワード全体を受け取れることを可能にする。しかしながら、キーワードを検出するためにそのようなデバイスを使用して入力音声を継続的にキャプチャし処理することは通常、モバイルデバイスの場合、電力消費の増加またはバッテリー寿命の低下をもたらす。いくつかのデバイスは、入力音声ストリームが受け取られ処理される時間を低減するために、デューティサイクルを実施している。しかしながら、ユーザからのキーワードの先頭がデューティサイクルの非アクティブ状態にあたる場合、キーワードの検出の失敗を引き起し、ユーザがキーワードを繰り返すことを必要とする可能性がある。

## 【0005】

[0005]加えて、異なるユーザが電子デバイスにおける機能を起動するために同じキーワードを話すとき、ユーザは通常、異なる速度でキーワードを話し、その結果、異なる時間期間がキーワードの部分の各々に充てられ得る。ユーザによる差に対処するために、いくつかの電子デバイスは、キーワードを検出するためにすべての考えられる時間期間のすべての組合せを処理している。しかしながら、そのような音声処理方法は一般に、計算複雑性 (computational complexity) を増大させ、かなり大きいコンピューティングリソースを必要とする。

## 【発明の概要】

## 【0006】

[0006]本開示は、機能またはアプリケーションを起動するためにターゲットキーワードを検出するための方法および装置を提供する。

## 【0007】

[0007]本開示の一態様によれば、電子デバイスにおける機能を起動するために複数の部分を含むターゲットキーワードを検出する方法が開示される。方法は、ターゲットキーワードの複数の部分のうちの1つから始まる入力音声を受け取り、入力音声から複数の音声特徴を抽出する。さらに、方法は、状態ネットワークを記述しているデータを取得し、ここにおいて、状態ネットワークは、単一の開始状態と、複数のエントリ状態と、単一の開始状態から複数のエントリ状態の各々への遷移とを含む。抽出された音声特徴と状態ネッ

10

20

30

40

50

トワークとに基づいて、入力音声ターゲットキーワードとして検出され得る。本開示は、本方法に係る装置、デバイス、システム、手段の組合せ、およびコンピュータ可読媒体についても説明する。

【0008】

[0008]本開示の別の態様によれば、機能を起動するためにターゲットキーワードを検出するための電子デバイスは、音声センサと、ボイスアクティベーションユニットとを含む。ターゲットキーワードは、複数の部分を含む。音声センサは、ターゲットキーワードの複数の部分のうちの1つから始まる入力音声を受け取るように構成される。ボイスアクティベーションユニットは、入力音声から複数の音声特徴を抽出し、状態ネットワークを記述しているデータを取得し、抽出された音声特徴と状態ネットワークとに基づいて、入力音声ターゲットキーワードに対応するかどうかを決定するように構成される。状態ネットワークは、単一の開始状態と、複数のエントリ状態と、単一の開始状態から複数のエントリ状態の各々への遷移とを含む。

10

【0009】

[0009]本開示のさらに別の態様によれば、電子デバイスにおいてターゲットキーワードを検出するための方法が開示される。ターゲットキーワードは、複数の部分を含む。方法は、入力音声を受け取り、入力音声から複数の音声特徴を抽出する。さらに、方法は、ターゲットキーワードの複数の部分に関連付けられた状態情報を取得する。状態情報は、ターゲットキーワードの部分に関連付けられた複数の状態に対する複数の持続時間範囲を含む。抽出された音声特徴と状態情報とに基づいて、入力音声ターゲットキーワードとして検出される。本開示は、本方法に係る装置、デバイス、システム、手段の組合せ、およびコンピュータ可読媒体についても説明する。

20

【0010】

[0010]本開示のまた別の態様によれば、複数の部分を含むターゲットキーワードを検出するための電子デバイスが開示される。電子デバイスは、音声センサと、ボイスアクティベーションユニットとを含む。音声センサは、入力音声を受け取るように構成される。ボイスアクティベーションユニットは、入力音声から複数の音声特徴を抽出し、ターゲットキーワードの複数の部分に関連付けられた状態情報を取得するように構成される。状態情報は、ターゲットキーワードの部分に関連付けられた複数の状態に対する複数の持続時間範囲を含む。ボイスアクティベーションユニットは、抽出された音声特徴と状態情報とに基づいて、入力音声をターゲットキーワードとして検出するようにさらに構成される。

30

【0011】

[0011]本開示の発明的態様の実施形態は、以下の詳細な説明を参照しながら、添付の図面とともに読むことで理解されよう。

【図面の簡単な説明】

【0012】

【図1】[0012]本開示の一実施形態による、入力音声においてターゲットキーワードが検出されたときにボイスアシスタントアプリケーションを起動するモバイルデバイスを示す図。

【0013】

【図2】[0013]本開示の一実施形態による、電子デバイスにおける機能を起動するために入力音声ストリームにおいてターゲットキーワードを検出するように構成された電子デバイスのブロック図。

40

【0014】

【図3】[0014]本開示の一実施形態による、ボイスアクティベーションユニットを起動し、入力音声ストリームをボイスアクティベーションユニットに継続的に与えるように構成された音声センサおよびスピーチ検出器のブロック図。

【0015】

【図4】[0015]本開示の一実施形態による、最初はデューティサイクルに基づいて受け取られ、デューティサイクルのアクティブ状態の間に受け取られた一部分が十分な強度の対象となる音声を含むと決定されると、継続的に受け取られる入力音声ストリームの図。

【0016】

【図5】[0016]本開示の一実施形態による、入力音声を受け取り、入力音声においてターゲットキーワードを検出するように構成されたボイスアクティベーションユニットのより

50

詳細なブロック図。

【図 6】[0017]本開示の一実施形態による、ボイスアシスタントユニットを起動するために入力音声ストリームにおいてターゲットキーワードを検出するための方法のフローチャート。

【図 7】[0018]一実施形態による、ターゲットキーワードの複数の状態についての複数の音声特徴の各々について決定された観測スコアのセットの図。

【図 8】[0019]一実施形態による、ターゲットキーワードの検出に使用するための複数のエントリ状態を含む、マルコフ連鎖モデルの下でのターゲットキーワードに対する複数の状態シーケンスの図。

【図 9】[0020]一実施形態による、各音声特徴に関連付けられた複数の状態の各々において終わる、各状態シーケンスについての最も大きいキーワードスコアを示す図。

10

【図 10】[0021]本開示の一実施形態による、入力音声においてターゲットキーワードを検出するために状態シーケンスについてのキーワードスコアを決定するための方法のフローチャート。

【図 11】[0022]本開示の一実施形態による、ターゲットキーワードに対するエントリ状態の数を決定する際に使用するためのターゲットキーワードに対する基準状態シーケンスの図。

【図 12】[0023]本開示の一実施形態による、ターゲットキーワードに対する基準入力音声进行处理することによってターゲットキーワードに対するエントリ状態の数を決定するための方法のフローチャート。

20

【図 13】[0024]本開示の別の実施形態による、入力音声を受け取り、ターゲットキーワードに関連付けられた複数の状態の各々に対する持続時間の所定の範囲に基づいて、入力音声においてターゲットキーワードを検出するように構成されたボイスアクティベーションユニットのブロック図。

【図 14】[0025]一実施形態による、ターゲットキーワードに関連付けられた各状態に対する持続時間の範囲に基づいて生成された、マルコフ連鎖モデルの下でのターゲットキーワードに対する複数の状態シーケンスのブロック図。

【図 15】[0026]一実施形態による、ターゲットキーワードを検出するために各状態に対する持続時間の所定の範囲に基づいて最も大きいキーワードスコアを決定するための方法のフローチャート。

30

【図 16】[0027]本開示のいくつかの実施形態による、機能を起動するために入力音声からターゲットキーワードを検出するための方法および装置が実装され得る、ワイヤレス通信システムにおけるモバイルデバイスのブロック図。

【発明を実施するための形態】

【0013】

[0028]次に、添付の図面にその例が示されている、様々な実施形態を詳細に参照する。以下の詳細な説明では、本主題の完全な理解を与えるために多数の具体的な詳細が記載される。ただし、本主題はこれらの具体的な詳細なしに実施され得ることが当業者には明らかであろう。他の事例では、様々な実施形態の態様を不必要に不明瞭にしないように、よく知られている方法、手順、システム、および構成要素については詳細に説明していない。

40

【0014】

[0029]図 1 は、本開示の一実施形態による、入力音声においてターゲットキーワードが検出されたときにボイスアシスタントアプリケーション 130 を起動するモバイルデバイス 120 を示す。最初に、モバイルデバイス 120 は、モバイルデバイス 120 におけるボイスアシスタントアプリケーション 130 などのアプリケーションまたは機能を起動するためのターゲットキーワードを記憶する。図示の実施形態では、ユーザ 110 が「START ASSISTANT」などのターゲットキーワードをモバイルデバイス 120 に話すと、モバイルデバイス 120 は入力音声をキャプチャし、入力音声においてターゲットキーワードを検出する。

50

## 【 0 0 1 5 】

[0030]ターゲットキーワードを検出すると、モバイルデバイス 1 2 0 はボイスアシスタントアプリケーション 1 3 0 を起動し、ボイスアシスタントアプリケーション 1 3 0 はユーザ 1 1 0 からの他のコマンドに応答して追加の機能を実行することができる。いくつかの実施形態では、モバイルデバイス 1 2 0 は、ユーザの音声入力からターゲットキーワードを検出する際に使用するための複数のターゲットキーワードを記憶し得る。ターゲットキーワードの各々は、モバイルデバイス 1 2 0 におけるターゲットキーワードに関連付けられたアプリケーションまたは機能を起動するように構成され得る。

## 【 0 0 1 6 】

[0031]図 2 は、本開示の一実施形態による、電子デバイス 2 0 0 における機能を起動するために入力音声ストリームにおいてターゲットキーワードを検出するように構成された電子デバイス 2 0 0 のブロック図を示す。本明細書で使用する「音声ストリーム」という用語は、1 つまたは複数の音声信号または音声データのシーケンスを指す。さらに、「ターゲットキーワード」という用語は、電子デバイス 2 0 0 における機能またはアプリケーションを起動するために使用され得る、1 つまたは複数の言葉または音声の任意のデジタルまたはアナログ表現を指す。電子デバイス 2 0 0 は、音声センサ 2 1 0 と、I/O ユニット 2 2 0 と、ストレージユニット 2 3 0 と、通信ユニット 2 4 0 と、プロセッサ 2 5 0 とを含む。電子デバイス 2 0 0 は、セルラーフォン、スマートフォン（たとえば、モバイルデバイス 1 2 0 ）、パーソナルコンピュータ、ラップトップコンピュータ、タブレットパーソナルコンピュータ、スマートテレビジョン、ゲームデバイス、マルチメディアプレーヤなどの、音声キャプチャおよび処理機能を備えた任意の適切なデバイスであり得る。

## 【 0 0 1 7 】

[0032]プロセッサ 2 5 0 は、デジタル信号プロセッサ (DSP) 2 5 2 と、ボイスアシスタントユニット 2 6 0 とを含み、電子デバイス 2 0 0 を管理し動作させるためのアプリケーションプロセッサまたは中央処理装置 (CPU) であり得る。DSP 2 5 2 は、スピーチ検出器 2 5 4 と、ボイスアクティベーションユニット 2 5 6 とを含む。一実施形態では、DSP 2 5 2 は、音声ストリームを処理する際の電力消費を低減するための低電力プロセッサである。この構成では、DSP 2 5 2 中のボイスアクティベーションユニット 2 5 6 は、入力音声ストリームにおいてターゲットキーワードが検出されたときにボイスアシスタントユニット 2 6 0 を起動するように構成される。図示の実施形態ではボイスアクティベーションユニット 2 5 6 がボイスアシスタントユニット 2 6 0 を起動するように構成されるが、ボイスアクティベーションユニット 2 5 6 はターゲットキーワードに関連付けられ得る任意の機能またはアプリケーションも起動し得る。

## 【 0 0 1 8 】

[0033]音声センサ 2 1 0 は、入力音声ストリームを受け取り、入力音声ストリームを DSP 2 5 2 中のスピーチ検出器 2 5 4 に与えるように構成され得る。音声センサ 2 1 0 は、1 つまたは複数のマイクロフォンあるいは電子デバイス 2 0 0 への音声入力を受け取る、キャプチャする、感知する、および/または検出するために使用され得る任意の他のタイプの音声センサを含み得る。加えて、音声センサ 2 1 0 は、そのような機能を実行するための任意の適切なソフトウェアおよび/またはハードウェアを利用し得る。

## 【 0 0 1 9 】

[0034]一実施形態では、音声センサ 2 1 0 は、デューティサイクルに従って入力音声ストリームを周期的に受け取るように構成され得る。この場合、音声センサ 2 1 0 は、入力音声ストリームの受け取られた部分がしきい値音声強度を超えるかどうかを決定し得る。音声ストリームの受け取られた部分がしきい値強度を超えるとき、音声センサ 2 1 0 はスピーチ検出器 2 5 4 を起動し、受け取られた部分を DSP 2 5 2 中のスピーチ検出器 2 5 4 に与える。代替的に、受け取られた部分がしきい値音声強度を超えるかどうかを決定することなしに、音声センサ 2 1 0 は、入力音声ストリームの一部分を周期的に受け取り、受け取られた部分をスピーチ検出器 2 5 4 に与えるためにスピーチ検出器 2 5 4 を起動し得る。

## 【 0 0 2 0 】

[0035]ターゲットキーワードを検出するために、ストレージユニット230は、ターゲットキーワードと、ターゲットキーワードの複数の部分に関連付けられた複数の状態についての状態情報とを記憶する。一実施形態では、ターゲットキーワードは、単音(phones)、音素(phonemes)などの音声の複数の基本単位、またはそれらの副次的単位に分割され得、ターゲットキーワードを表す複数の部分は、音声の基本単位に基づいて生成され得る。次いで、ターゲットキーワードの各部分は、隠れマルコフモデル(hidden Markov model) (「HMM」)、半マルコフモデル(semi-Markov model) (「SMM」)、またはそれらの組合せなどのマルコフ連鎖モデルの下での状態に関連付けられる。状態情報は、ターゲットキーワードに関連付けられた状態の中からの所定の数のエントリ状態と、これらの状態の各々からそれ自体を含む次の状態への遷移情報とを含み得る。ストレージユニット230は、RAM(ランダムアクセスメモリ)、ROM(読取り専用メモリ)、EEPROM(登録商標)(電氣的消去可能プログラマブル読取り専用メモリ)、フラッシュメモリ、またはSSD(ソリッドステートドライブ)などの任意の適切なストレージまたはメモリデバイスを使用して実装され得る。

10

## 【 0 0 2 1 】

[0036]DSP252中のスピーチ検出器254は、起動されると、音声センサ210から入力音声ストリームの部分を受け取る。一実施形態では、スピーチ検出器254は、受け取られた部分から複数の音声特徴を抽出し、ガウス混合モデル(Gaussian mixture model) (GMM)ベースの分類器、ニューラルネットワーク、HMM、グラフィカルモデル、およびサポートベクターマシン(Support Vector Machine) (SVM)技法などの任意の適切な音声分類方法を使用することによって、抽出された音声特徴がスピーチなどの対象となる音声を示すかどうかを決定する。受け取られた部分が対象となる音声であると決定された場合、スピーチ検出器254はボイスアクティベーションユニット256を起動し、入力音声ストリームの受け取られた部分と残りの部分とはボイスアクティベーションユニット256に与えられる。いくつかの他の実施形態では、スピーチ検出器254はDSP252中で省略され得る。この場合、受け取られた部分がしきい値強度を超えると、音声センサ210はボイスアクティベーションユニット256を起動し、入力音声ストリームの受け取られた部分と残りの部分とを直接ボイスアクティベーションユニット256に与える。

20

30

## 【 0 0 2 2 】

[0037]ボイスアクティベーションユニット256は、起動されると、入力音声ストリームを継続的に受け取り、入力音声ストリームからターゲットキーワードを検出するように構成される。入力音声ストリームが受け取られると、ボイスアクティベーションユニット256は入力音声ストリームから複数の音声特徴を連続的に抽出し得る。加えて、ボイスアクティベーションユニット256は、複数の状態と、所定の数のエントリ状態と、ターゲットキーワードの遷移情報とを含む状態情報をストレージユニット230から取得し得る。各音声特徴について、GMM、ニューラルネットワーク、およびSVMなどの任意の適切な確率モデルを使用することによって、状態の各々に対して観測スコアが決定され得る。

40

## 【 0 0 2 3 】

[0038]遷移情報から、ボイスアクティベーションユニット256は、ターゲットキーワードについて考えられる複数の状態シーケンスにおける状態の各々から次の状態への遷移スコアを取得し得る。遷移情報は、ターゲットキーワードの音声の基本単位に関連付けられた複数のエントリ状態に基づいてボイスアクティベーションユニット256が入力音声ストリームにおいてターゲットキーワードを検出することを可能にするための、所定の数のエントリ状態についての遷移スコアも含み得る。そのような複数のエントリ状態をターゲットキーワードに与えることによって、ボイスアクティベーションユニット256は、ターゲットキーワードの先頭の後に始まる入力音声ストリームを処理することによって、ターゲットキーワードを検出し得る。

50



## 【 0 0 2 4 】

[0039]観測スコアを決定し、遷移スコアを取得した後、ボイスアクティベーションユニット256は、考えられる状態シーケンスについてのキーワードスコアを決定する。一実施形態では、決定されたキーワードスコアの中で最も大きいキーワードスコアが所定のしきい値スコアを超える場合、ボイスアクティベーションユニット256は入力音声ストリームをターゲットキーワードとして検出する。ターゲットキーワードを検出すると、ボイスアクティベーションユニット256は、ボイスアシスタントユニット260をオンにするためのアクティベーション信号を生成および送信し、ボイスアシスタントユニット260はターゲットキーワードに関連付けられる。

## 【 0 0 2 5 】

[0040]ボイスアシスタントユニット260は、ボイスアクティベーションユニット256からのアクティベーション信号に応答して起動される。起動されると、ボイスアシスタントユニット260は、タッチスクリーン上および/またはI/Oユニット220のスピーカーを通じて「MAY I HELP YOU?」などのメッセージを出力することによって、ボイスアシスタント機能を実行し得る。それに応答して、ユーザは電子デバイス200の様々な関連機能を起動するためにボイスコマンドを話してもよい。たとえば、インターネット検索のためのボイスコマンドが受け取られると、ボイスアシスタントユニット260はボイスコマンドを検索コマンドとして認識し、ネットワーク270を通じて通信ユニット240を介してウェブ検索を実行し得る。

## 【 0 0 2 6 】

[0041]図3は、本開示の一実施形態による、ボイスアクティベーションユニット256を起動し、入力音声ストリームをボイスアクティベーションユニット256に継続的に与えるように構成された音声センサ210およびスピーチ検出器254のブロック図を示す。音声センサ210は、マイクロフォン310と、音声検出器320とを含む。一実施形態では、マイクロフォン310および音声検出器320は、デューティサイクルに基づいて入力音声ストリームを周期的に受け取るように構成される。たとえば、マイクロフォン310および音声検出器320は、時間の10%(たとえば、200ms期間中の20ms)で入力音声ストリームを受け取るように、10%デューティサイクルで動作してもよい。図示したように、マイクロフォン310は、デューティサイクルの各アクティブ状態の間に入力音声ストリームを受け取り、入力音声ストリームの受け取られた部分を音声検出器320に与えるように構成され得る。このプロセスでは、マイクロフォン310はまた、受け取られた音声ストリームをデジタル形式に変換し得る。

## 【 0 0 2 7 】

[0042]音声検出器320は、入力音声ストリームの受け取られた部分の信号特性を解析し、受け取られた部分がしきい値音声強度を超えるかどうかを決定し得る。たとえば、音声検出器320は、受け取られた部分の音声強度を決定するために、受け取られた部分の音声エネルギー値または振幅値を解析し得る。受け取られた部分がしきい値音声強度を超える音声であると決定された場合、音声検出器320はアクティベーション信号と受け取られた部分とをスピーチ検出器254に与え得る。

## 【 0 0 2 8 】

[0043]スピーチ検出器254は、起動されると、音声検出器320から入力音声ストリームの部分を受け取る。次いで、スピーチ検出器254は、MFCC(メル周波数ケプストラム係数(Mel-frequency cepstral coefficients))法、LPC(線形予測符号(linear predictive coding))法、またはLSP(線スペクトル対(line spectral pair))法などの任意の適切な信号処理方式を使用することによって、受け取られた部分から1つまたは複数の音声特徴を抽出し得る。抽出された音声特徴を使用して、スピーチ検出器254は、受け取られた部分がスピーチを含むかどうかを決定するために、入力音声ストリームの受け取られた部分を分類する。

## 【 0 0 2 9 】

[0044]図示の実施形態では、入力音声ストリームの受け取られた部分がしきい値音声強

10

20

30

40

50

度を超え、スピーチを含むと決定された場合、デューティサイクル機能は、（たとえば、フルデューティサイクルまたは100%デューティサイクルを使用して）入力音声ストリームの残りの部分をボイスアクティベーションユニット256に継続的に与えるために無効にされ得る。一実施形態によれば、受け取られた部分がしきい値音声強度を超えると音声検出器320が決定した場合、マイクロフォン310および音声検出器320が、入力音声ストリームの残りの部分を受け取り、スピーチ検出器254に送信するために、アクティブ状態で動作し続け得るように、音声検出器320は音声センサ210のデューティサイクル機能を無効にする。デューティサイクルのアクティブ状態の間に受け取られた冒頭の部分（initial portion）がスピーチを含むとスピーチ検出器254が決定した場合、スピーチ検出器254は、入力音声ストリームの冒頭の部分と残りの部分とをボイスアクティベーションユニット256に与える。一方、冒頭の部分がスピーチを含まないとスピーチ検出器254が決定した場合、スピーチ検出器254はアクティベーション信号を生成し、アクティベーション信号は、マイクロフォン310および音声検出器320のデューティサイクル機能を起動するために音声センサ210に与えられる。

10

【0030】

[0045]別の実施形態では、受け取られた部分がしきい値音声強度を超えると音声検出器320が決定した場合、音声検出器320はスピーチ検出器254を起動し、受け取られた部分をスピーチ検出器254に与える。受け取られた部分がスピーチを含むとスピーチ検出器254が決定した場合、スピーチ検出器254はマイクロフォン310と音声検出器320とにデューティサイクル機能のデアクティベーション信号を与える。デアクティベーション信号を受け取ると、マイクロフォン310および音声検出器320は、入力音声ストリームの残りの部分を受け取り、スピーチ検出器254に送信するために、アクティブ状態で動作し続け得、その結果、スピーチ検出器254はその部分をボイスアクティベーションユニット256に与える。一方、デューティサイクルのアクティブ状態の間に受け取られた冒頭の部分がスピーチを含まないとスピーチ検出器254が決定した場合、スピーチ検出器254は、マイクロフォン310および音声検出器320がデューティサイクル機能に従って動作するように、アクティベーション信号をマイクロフォン310と音声検出器320とに与えない。

20

【0031】

[0046]図4は、本開示の一実施形態による、最初はデューティサイクルに基づいて受け取られ、デューティサイクルのアクティブ状態の間に受け取られた一部分が十分な強度の対象となる音声を含むと決定されると、継続的に受け取られる入力音声ストリームの図を示す。図示のように、入力音声ストリームの複数の部分410、420、および430が、デューティサイクルに基づいて周期的に受け取られる。部分410および420は、十分な強度のまたは対象となる音声（たとえば、スピーチ）を含まない。一方、部分430は、しきい値音声強度を超えるとともにスピーチを含む音声を含む。したがって、入力音声ストリームの部分430と残りの部分440とは、継続的に受け取られ、ターゲットキーワードを検出するためのボイスアクティベーションユニット256に与えられる。

30

【0032】

[0047]示される図では、ある部分（たとえば、部分410）の開始と次の部分（たとえば、部分420）の開始との間の期間は、デューティサイクル期間T1を表す。期間T2は、デューティサイクルがアクティブ状態である持続時間を表す。デューティサイクルの非アクティブ状態はT1 - T2によって示され、その時間期間の間、入力音声ストリームは受け取られない。デューティサイクルに従って入力音声ストリームを受け取ることにより、ターゲットキーワードを検出するためのボイスアクティベーションユニット256に与えられる入力音声ストリームの冒頭の部分430は、ターゲットキーワードの先頭部分に続くターゲットキーワードの一部分に対応し得る。

40

【0033】

[0048]図5は、本開示の一実施形態による、入力音声を受け取り、入力音声においてターゲットキーワードを検出するように構成されたボイスアクティベーションユニット25

50

6 のより詳細なブロック図を示す。ボイスアクティベーションユニット 2 5 6 は、セグメント化ユニット 5 1 0 と、特徴抽出器 5 2 0 と、スコア決定ユニット 5 3 0 と、キーワード検出ユニット 5 4 0 とを含む。スコア決定ユニット 5 3 0 は、観測スコア決定ユニット 5 5 0 と、最大キーワードスコア決定ユニット 5 6 0 とを含む。

【 0 0 3 4 】

[0049]セグメント化ユニット 5 1 0 は、スピーチ検出器 2 5 4 から入力音声を受け取り、受け取られた入力音声を等しい時間期間の複数の連続したフレームにセグメント化する。特徴抽出器 5 2 0 は、セグメント化ユニット 5 1 0 からフレームを連続的に受け取り、フレームの各々から音声特徴を抽出する。一実施形態では、特徴抽出器 5 2 0 は、M F C C 法などの任意の適切な特徴抽出方法を使用して、フレームから音声特徴を抽出し得る。たとえば、M F C C 法の場合、N 次元ベクトル中の成分はセグメント化されたフレームの各々から計算され、ベクトルは音声特徴として使用される。

10

【 0 0 3 5 】

[0050]スコア決定ユニット 5 3 0 において、観測スコア決定ユニット 5 5 0 は、音声特徴を連続的に受け取り、ストレージユニット 2 3 0 からターゲットキーワードの状態情報を受け取る。一実施形態によれば、ターゲットキーワードの状態情報は、ターゲットキーワードの複数の部分に関連付けられた複数の状態と、状態の各々に与えられる G M M などの確率モデル（たとえば、確率関数）とを含み得る。上記で説明したように、ターゲットキーワードは音声の複数の基本単位に分割され得、ターゲットキーワードを表す複数の部分は、状態に対応する音声の基本単位に基づいて生成され得る。いくつかの実施形態では、ターゲットキーワードの状態は、ターゲットキーワードを表す複数の部分のうちのいずれにも関連付けられず、開始状態として使用され得る、非キーワード状態（たとえば、「フィラー」状態）も含み得る。たとえば、各々が単音などの音声の基本単位に対応する所定の数の状態を含むターゲットキーワードの場合、非キーワード状態は、ターゲットキーワードに含まれる基本単位以外の音声の基本単位を表し得る。

20

【 0 0 3 6 】

[0051]各音声特徴が受け取られると、観測スコア決定ユニット 5 5 0 は、音声特徴とストレージユニット 2 3 0 から受け取られた状態情報とに基づいて、ターゲットキーワードに関連付けられた状態の各々についての観測スコアを決定する。一実施形態では、状態の各々についての観測スコアは、関連する状態の確率モデルに従って確率値を計算することによって、受け取られた音声特徴について決定される。このようにして計算された確率値の各々は、関連する状態についての観測スコアとして使用され得る。状態についての高い観測スコアは、その状態に対して、音声特徴が音声の基本単位に対応する確率が高いことを示す。観測スコア決定ユニット 5 5 0 は、ターゲットキーワードについて、考えられる複数の状態シーケンスについてのキーワードスコアを決定するために、受け取られた音声特徴の各々についての観測スコアを最大キーワードスコア決定ユニット 5 6 0 に与える。

30

【 0 0 3 7 】

[0052]最大キーワードスコア決定ユニット 5 6 0 は、音声特徴の各々についての観測スコアを受け取り、ストレージユニット 2 3 0 から状態情報を取得する。この構成では、状態情報は、非キーワード状態（単一の非キーワード開始状態を含む）を含む、ターゲットキーワードの複数の状態と、複数の状態の中からの所定の数のエントリ状態と、状態の各々からそれ自体を含む次の状態への遷移情報とを含み得る。エントリ状態は、ターゲットキーワードの考えられる状態シーケンスの各々において非キーワード状態（または開始状態）が遷移し得る先の最初の状態を表す。

40

【 0 0 3 8 】

[0053]状態情報における遷移情報は、ターゲットキーワードの考えられる状態シーケンスの各々における、状態の各々から次の状態への遷移スコアを含む。遷移スコアは、考えられる状態シーケンスの各々における、状態の各々が次の状態に遷移する確率値を表し得る。遷移スコアは、非キーワード状態から所定の数のエントリ状態への遷移スコアも含む。

50

## 【 0 0 3 9 】

[0054] 受け取られた観測スコアと遷移スコアとに基づいて、最大キーワードスコア決定ユニット 5 6 0 は、考えられる状態シーケンスの各々についてのキーワードスコアを計算する。この場合、非キーワード状態は入力音声を受け取られる前に割り当てられているので、状態シーケンスは非キーワード状態（すなわち、開始状態）から始まり得る。したがって、遷移スコアは、非キーワード状態からエントリ状態のいずれか 1 つへの遷移スコアを含み、状態シーケンスにおける非キーワード状態からそれ自体への遷移スコアも含む。各音声特徴についての観測スコアのセットが観測スコア決定ユニット 5 5 0 から受け取られると、最大キーワードスコア決定ユニット 5 6 0 は、上記で説明したような方法で、次の状態を各状態シーケンスに追加し、更新された状態シーケンスの各々についてのキーワードスコアを決定する。

10

## 【 0 0 4 0 】

[0055] 状態シーケンスについてのキーワードスコアのセットが計算されると、最大キーワードスコア決定ユニット 5 6 0 は、キーワードスコアの中から最も大きいキーワードスコアを選択する。キーワードスコアは、ビタビアルゴリズムなどの任意の適切な方法を使用することによって、最も大きいキーワードスコアを決定するように計算され得る。最も大きいキーワードスコアを決定した後、最大キーワードスコア決定ユニット 5 6 0 は最も大きいキーワードスコアをキーワード検出ユニット 5 4 0 に与える。一実施形態では、最も大きいキーワードスコアは、最も大きいキーワードスコアを有する状態シーケンスの最後の状態がターゲットキーワードの音声の最後の基本単位（たとえば、最後の単音）に対応するときのみ、キーワード検出ユニット 5 4 0 に与えられる。

20

## 【 0 0 4 1 】

[0056] 最大キーワードスコア決定ユニット 5 6 0 から最も大きいキーワードスコアを受け取ると、キーワード検出ユニット 5 4 0 は、最も大きいキーワードスコアに基づいて、入力音声においてターゲットキーワードを検出する。たとえば、キーワード検出ユニット 5 4 0 は、ターゲットキーワードを検出するためのしきい値スコアをストレージユニット 2 3 0 から受け取り、最も大きいキーワードスコアが受け取られたしきい値スコアよりも大きい場合、ターゲットキーワードを検出し得る。この場合、しきい値スコアは、所望の信頼性レベル内でターゲットキーワードを検出するために最小キーワードスコアに設定され得る。

30

## 【 0 0 4 2 】

[0057] いくつかの実施形態では、最大キーワードスコア決定ユニット 5 6 0 は、非キーワード状態シーケンスについての非キーワードスコアを決定する。非キーワードスコアは、非キーワード状態シーケンスを含む、考えられる状態シーケンスについてのキーワードスコアから取得され、キーワード検出ユニット 5 4 0 に与えられ得る。キーワード検出ユニット 5 4 0 は、最も大きいキーワードスコアと非キーワードスコアとの間の差に基づいて信頼性値を決定し、入力音声においてターゲットキーワードを検出する。この場合、キーワード検出ユニット 5 4 0 は、ストレージユニット 2 3 0 からしきい値信頼性値を受け取り、信頼性値がしきい値信頼性値よりも大きい場合、ターゲットキーワードを検出し得る。ターゲットキーワードを検出する際に最も大きいキーワードスコアと非キーワードスコアとの間の差を使用することは、入力音声がキーワードスコアに影響を及ぼす可能性がある雑音などの周囲音を含むときは特に、検出精度を改善し得る。ターゲットキーワードが検出されると、キーワード検出ユニット 5 4 0 は、ボイスアシスタントユニット 2 6 0 をオンにするためのアクティベーション信号を生成し、これを与え、ボイスアシスタントユニット 2 6 0 はターゲットキーワードに関連付けられる。

40

## 【 0 0 4 3 】

[0058] 図 6 は、本開示の一実施形態による、ボイスアシスタントユニット 2 6 0 を起動するために入力音声においてターゲットキーワードを検出するための、ボイスアクティベーションユニット 2 5 6 によって実行される方法 6 0 0 のフローチャートである。ターゲットキーワードは、冒頭の部分と、複数の後続（subsequent portions）の部分とを含み

50

得る。最初に、610において、ボイスアクティベーションユニット256は、ターゲットキーワードの後続の部分のうちの1つから始まる入力音声を受け取る。受け取られた入力音声は複数のフレームにセグメント化された後、620において、ボイスアクティベーションユニット256は、MFCC法などの任意の適切な信号処理方式を使用することによって、複数のフレームから複数の音声特徴を抽出する。

#### 【0044】

[0059]次いで、630において、ボイスアクティベーションユニット256は、ストレージユニット230からターゲットキーワードの冒頭の部分と後続の部分とに関連付けられた状態情報を取得する。図2および図5に関して上記で説明したように、状態情報は、所定の数のエントリ状態と、確率モデルと、遷移情報とを含み得る。640において、抽出された音声特徴と状態情報とに基づいて、ボイスアクティベーションユニット256は、入力音声をターゲットキーワードとして検出する。ターゲットキーワードが検出されると、650において、ボイスアクティベーションユニット256は、ターゲットキーワードに関連付けられたボイスアシスタントユニット260を起動する。

#### 【0045】

[0060]図7は、一実施形態による、ターゲットキーワード（たとえば、「START ASSISTANT」）の複数の状態についての複数の音声特徴F1～F5の各々について観測スコア決定ユニット550によって生成された観測スコアのセットの図700を示す。図700に示すターゲットキーワードの状態は「F」、「S」、「T」、「A」、「R」、「T」などを含み、ここで、状態「F」は非キーワード状態またはフィルタ状態を示す。観測スコア決定ユニット550は、各フレームから抽出された音声特徴を連続的に受け取る。各音声特徴について、観測スコア決定ユニット550は、たとえば、図700の各状態および音声特徴の数字によって示されるように、GMMなどの確率モデルを使用することによって、状態の各々についての観測スコアを決定する。

#### 【0046】

[0061]図示の実施形態では、観測スコア決定ユニット550が、所定の時間間隔に連続的に受け取られた音声特徴F1、F2、F3、F4、およびF5の各々を受け取ると、状態についての観測スコアのセットが決定される。たとえば、音声特徴F1の場合、観測スコアのセット710は状態（すなわち、状態「F」、「S」、「T」、「A」、「R」、「T」など）について決定される。音声特徴F1についての観測スコアのセット710を決定した後、観測スコア決定ユニット550は、音声特徴F2～F5についてそれぞれ、複数の観測スコアのセット720～750を連続的に決定し得る。観測スコア決定ユニット550は、それぞれ音声特徴F1～F5についての観測スコアのセット710～750を、ターゲットキーワードを検出するための最大キーワードスコア決定ユニット560に連続的に与え得る。

#### 【0047】

[0062]図8は、一実施形態による、ターゲットキーワードの検出に使用するための複数のエントリ状態「S」、「T」、「A」、および「R」を含む、マルコフ連鎖モデルの下でのターゲットキーワードについて、考えられる複数の状態シーケンスを含む、状態ネットワークの図800を示す。一実施形態では、エントリ状態の数はあらかじめ決定され得る。図8の図示の実施形態では、図800のエントリ状態のセット810によって示されるように、所定の数のエントリ状態は4である。さらに、図800は、ターゲットキーワードについて、考えられる複数の状態シーケンスにおける現在の音声特徴（たとえば、音声特徴F1）に関連付けられた現在の状態の各々から次の音声特徴（たとえば、音声特徴F2）に関連付けられた複数の次の状態への遷移ラインを示す。

#### 【0048】

[0063]最初に、図800の音声特徴F0は、入力音声を受け取られていないことを示す。入力音声を受け取られると、音声特徴F1～F5は受け取られた入力音声から連続的に抽出される。したがって、非キーワード状態「F」は音声特徴F0のみに割り当てられ、すべての状態シーケンスの単一の開始状態として働く。図800は、音声特徴F0の開始

状態「F」から、次の音声特徴F 1に関連付けられた、考えられる次の状態、すなわち、エントリ状態「S」、「T」、「A」、および「R」の各々への遷移ラインも示す。次いで、音声特徴F 1において、考えられる状態シーケンスの各々について、音声特徴F 1の状態の各々から次の音声特徴F 2（すなわち、次の状態）の状態の各々への遷移ラインが示されている。図800に示すように、そのような遷移ラインは、ターゲットキーワードについてあらかじめ決定され、受け取られた入力音声の残りの音声特徴F 2、F 3、F 4、およびF 5の各々に同様の方法で適用され得る。

【0049】

[0064]この構成では、各遷移ラインは遷移スコアに関連付けられる。状態のいくつかは、次の状態のいくつかへの遷移ラインを有しない場合がある。たとえば、ある音声特徴の状態「S」から次の音声特徴の状態「A」、「R」、および「T」への遷移ラインは与えられていない。一実施形態では、現在の状態から次の状態への遷移ラインがない場合、現在の状態から次の状態への状態シーケンスが生成されないことがある。別の実施形態では、ある状態から次の状態への遷移ラインが与えられない場合、遷移スコアは、そのような遷移スコアを含む状態シーケンスがターゲットキーワードを検出する際に使用するための最も大きいキーワードスコアを有さなくてもよいことを保証するために、大きい負の数（たとえば、-10）に設定され得る。

10

【0050】

[0065]図800に基づいて、遷移ラインに関連付けられた、および遷移ラインに関連付けられない遷移スコアの例示的な表が次のように与えられ得る。

20

【0051】

【表 1】

表1

次の状態 現在の状態	F	S	T	A	R	T	...
F	0.7	0.8	0.8	0.7	0.9	-10	...
S	-10	0.9	0.8	-10	-10	-10	...
T	-10	-10	0.7	0.8	-10	-10	...
A	-10	-10	-10	0.8	0.9	-10	...
R	-10	-10	-10	-10	0.7	0.6	...
T	-10	-10	-10	-10	-10	0.8	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	

## 【 0 0 5 2 】

表 1 に示すように、現在の状態（たとえば、状態「S」）から次の状態（たとえば、状態「A」）への遷移ラインがない場合、- 10 の関連する遷移スコアが割り当てられる。最初に、現在の状態「F」はエントリ状態「S」、「T」、「A」、および「R」への4つの遷移ラインを有するので、0.8、0.8、0.7、および0.9の遷移スコアが現在の状態「F」からエントリ状態「S」、「T」、「A」、および「R」への4つの遷移ラインに割り当てられる。複数のエントリ状態を与えることによって、ボイスアクティベーションユニット256がターゲットキーワードの先頭部分（すなわち、状態「S」）から始まらない入力音声を受け取った場合でも、ターゲットキーワードはそのような入力音声から検出され得る。

## 【 0 0 5 3 】

[0066] 図 9 は、一実施形態による、各音声特徴に関連付けられた複数の状態の各々において終わる、各シーケンスについての最も大きいキーワードスコアを示す図 900 である。この図 900 では、ターゲットキーワード「START ASSISTANT」は、音声の基本単位の各々に関連付けられた状態を含む。説明を容易にするために、図 900 は音声の基本単位として単音（「S」、「T」、「A」、「R」、「T」など）を使用して示されている。

## 【 0 0 5 4 】

[0067] この実施形態では、各音声特徴に関連付けられた状態の各々において終わる、考えられる状態シーケンスについて決定されるキーワードスコアの中で、最も大きいキーワ

ードスコアが決定される。次いで、そのような最も大きいキーワードスコアは、たとえば、候補キーワードスコアの中から最も大きいキーワードスコアを選択し、選択された最も大きいキーワードスコアをしきい値スコアと比較することによって、ターゲットキーワードを検出する際に候補キーワードスコアとして使用される。さらに、候補キーワードスコアは、各次の音声特徴に関連付けられた状態の各々において終わる、考えられる状態シーケンスについての最も大きいキーワードスコア（すなわち、次の候補キーワードスコア）を決定する際に使用され得る。候補キーワードスコアおよび選択された最も大きいキーワードスコアは、図6で与えられる観測スコアと表1で与えられる遷移スコアとに基づいてビタビアルゴリズムを使用することによって、決定され得る。

【0055】

10

[0068]図900では、音声特徴F1の場合、図示した状態「F」、「S」、「T」、「A」、「R」、および「T」の各々は、音声特徴F0の開始状態「F」からの1つの、考えられる状態シーケンスを有し、状態についての最も大きいスコア（すなわち、候補スコア）はそれぞれ、1.0、1.6、1.7、0.8、1.0、および-9.9と決定される。候補キーワードスコアは、非キーワード状態「F」から音声特徴F1の状態の各々への遷移スコアと、音声特徴F1の関連する状態の観測スコアとを合計することによって決定され得る。たとえば、このようにして、エントリ状態のセット810に含まれる状態「T」の候補スコア1.7は、非キーワード状態「F」から状態「T」への遷移スコアと、状態「T」の観測スコアとを合計することによって決定される。図示の例では、音声特徴F1の時点での状態「T」を含むいかなる連続する状態シーケンスも、ターゲットキーワードを検出する際に使用するための最も大きいキーワードスコアを有することができないように、大きい負の数（たとえば、-10）が遷移スコアとして、エントリ状態のセット810に含まれない状態「T」に割り当てられる。

20

【0056】

[0069]音声特徴F2の場合、最後の状態としての音声特徴F1の状態の各々を有する状態シーケンスについての候補キーワードスコアと、音声特徴F1の各状態から音声特徴F2の状態への遷移スコアと、音声特徴F2の状態の観測スコアとを合計することによって、候補キーワードスコアが各状態について決定される。音声特徴F2に関連付けられた状態の各々について、上記の合計の中で最も大きい合計が、最後の状態としての状態を有する、考えられる状態シーケンスについての候補キーワードスコアとして選択される。たとえば、最後の状態として音声フレームF2の状態「A」を有する、考えられる状態シーケンスについてのキーワードスコアは、表2において次のように与えられる。

30

【0057】



【表 2】

表2

前の状態	前の候補 キーワード スコア	現在の状態 「A」への 遷移スコア	音声特徴F2に 関連付けられた 状態「A」についての 観測スコア	合計
F	1.0	0.7	0.7	2.4
S	1.6	-10	0.7	-7.7
T	1.7	0.8	0.7	3.2
A	0.8	0.8	0.7	2.3
R	1.0	-10	0.7	-8.3
T	-9.9	-10	0.7	-19.2

## 【 0 0 5 8 】

上記の表 2 では、上記の合計の中で最も大きいキーワードスコアの 3.2 が、音声特徴 F 2 に関連付けられた状態「A」を有する、考えられる状態シーケンスについての候補キーワードスコアとして選択される。同様の方法で、その他の図示した状態、すなわち、「F」、「S」、「T」、「R」、「T」についての候補キーワードスコアは 1.9、2.9、2.7、2.0、および 1.8 と決定され得る。次いで、図 900 に示すように、状態「A」についての候補キーワードスコアの 3.2 が、音声特徴 F 2 の最も大きいキーワードスコアとして選択され得る。

## 【 0 0 5 9 】

[0070] 残りの音声特徴 F 3、F 4、および F 5 の各々について、最後の状態として残りの音声特徴 F 3、F 4、および F 5 の各々に関連付けられた状態の各々を有する、考えられる状態シーケンスについての候補キーワードスコアのセットは、音声特徴 F 2 と同様の方法で計算され得る。候補キーワードスコアの中で、音声特徴 F 3、F 4、および F 5 の各々についての最も大きいキーワードスコアも同様の方法で決定され得る。音声特徴 F 1 ~ F 5 の各々についての最も大きいキーワードスコアは、ターゲットキーワードを検出するために使用され得る。図 9 の図示の例では、矢印と太線の円とで示される、このようにして決定された最も大きいキーワードスコアの状態を含む状態シーケンスは、ターゲットキーワードの複数の連続した部分に対応し得る。

## 【 0 0 6 0 】

[0071] 図 10 は、本開示の一実施形態による、入力音声においてターゲットキーワードを検出するための、ボイスアクティベーションユニット 256 によって実行される詳細な方法 640 のフローチャートである。最初に、1010 において、観測スコア決定ユニッ

ト 5 5 0 は、入力音声の複数の音声特徴を受け取り、ターゲットキーワードに関連付けられた複数の状態の各々についての観測スコアを決定する。観測スコアは、最大キーワードスコア決定ユニット 5 6 0 に与えられる。1 0 2 0 において、最大キーワードスコア決定ユニット 5 6 0 は、所定の数のエントリ状態と、ターゲットキーワードについて、考えられる複数の状態シーケンスにおける状態の各々から次の状態への遷移スコアとを取得する。遷移スコアは、非キーワード状態からエントリ状態の各々への遷移スコアを含み得る。

#### 【 0 0 6 1 】

[0072] 1 0 3 0 において、最大キーワードスコア決定ユニット 5 6 0 は、観測スコアと遷移スコアとを使用することによって、考えられる状態シーケンスについてのキーワードスコアを決定する。このプロセスでは、図 5 および図 9 に関して上記で説明した方法で、決定されたキーワードスコアの中から最も大きいキーワードスコアが選択され、キーワード検出ユニット 5 4 0 に与えられ得る。次いで、1 0 4 0 において、受け取られた最も大きいキーワードスコアに基づいて、入力音声においてターゲットキーワードが検出される。

#### 【 0 0 6 2 】

[0073] 図 1 1 は、本開示の一実施形態による、ターゲットキーワードに対するエントリ状態の数を決定する際に使用するためのターゲットキーワード「START ACTION」に対する基準状態シーケンス 1 1 0 2 の図を示す。一実施形態では、エントリ状態の数の決定は、スコア決定ユニット 5 3 0 において実行され得る。図示のように、ターゲットキーワードに対する基準入力音声は、キーワードのすべての部分が受け取られるように、デューティサイクルを使用することなしに（すなわち、フルデューティサイクルに基づいて）時間期間 T 0 にわたってキャプチャされる。基準入力音声は、時間期間 T 0 にわたってユーザによって話された音声であり、等しい時間期間の複数の連続したフレームにセグメント化され得る。図示の実施形態では、複数のフレームは、ターゲットキーワードに対する基準状態シーケンス 1 1 0 2 における状態に対応する。各フレームについて、音声特徴が抽出され、ストレージユニット 2 3 0 に記憶され得る。

#### 【 0 0 6 3 】

[0074] 一実施形態では、スコア決定ユニット 5 3 0 は、図 5 ~ 図 9 に関して上記で説明した方法で、抽出された音声特徴をストレージユニット 2 3 0 から受け取り、すべての考えられる状態シーケンスについての最も大きいキーワードスコアを決定する。この場合、単一のエントリ状態 1 1 0 4 から始まる基準状態シーケンス 1 1 0 2 全体が、最も大きいキーワードスコアの決定に使用され得る。基準状態シーケンス 1 1 0 2 における状態および各状態の数は、基準状態シーケンス 1 1 0 2 をバックトラックすることによって決定され得る。この実施形態では、基準状態シーケンス 1 1 0 2 のエントリ状態の数は、T 2 / T 1 のデューティサイクルにおける非アクティブ状態 T 1 - T 2 の間にターゲットキーワードのいくつかの部分が受け取られないことがある時間期間に基づいて、決定され得る。たとえば、ユーザによって話された入力音声はデューティサイクルに従って受け取られるとき、非アクティブ状態に対応するターゲットキーワードの複数の部分に関連付けられた複数の状態 1 1 1 0（すなわち、「S」、「T」、「A」、「A」、および「R」）が受け取られないことがある。

#### 【 0 0 6 4 】

[0075] 図 1 1 に示すように、状態「A」は、ユーザのスピーチ特性（たとえば、スピーチの速度）により、状態 1 1 1 0 において 2 回現れることがある。デューティサイクルの冒頭の非アクティブ期間 T 1 - T 2 に対応する状態 1 1 1 0 は 4 つの冒頭の状態「S」、「T」、「A」、および「R」を含むので、入力音声はデューティサイクルに基づいて受け取られるとき、入力音声は最大で 4 つの冒頭の状態を含まないことがある。この例では、スピーチが非アクティブ状態内で始まり、デューティサイクルの次のアクティブ状態において受け取られるように始まるときでも、ユーザによって話されたスピーチにおいてターゲットキーワードが検出され得るように、エントリ状態の数は 5 以上であると決定され得る。この場合、受け取られたスピーチの先頭部分は、エントリ状態のうちのいずれか 1

10

20

30

40

50

つから始まり得る。

【 0 0 6 5 】

[0076]一実施形態では、スコア決定ユニット 5 3 0 はターゲットキーワードに対する複数の基準入力音声を受け取り得、これらの基準入力音声の各々はフルデューティサイクルに基づいてキャプチャされる。基準入力音声の各々について、スコア決定ユニット 5 3 0 は、基準状態シーケンス 1 1 0 2 に関して上記で説明したのと同様の方法で、基準状態シーケンスと、基準状態シーケンスにおける複数の状態に対する状態時間期間とを決定することによって、デューティサイクルにおける冒頭の非アクティブ期間に従ってエントリ状態を決定する。次いで、基準状態シーケンスにおけるエントリ状態の数は、入力音声からターゲットキーワードを検出する際に使用するためのエントリ状態を決定するために、平均され得る。

10

【 0 0 6 6 】

[0077]図 1 2 は、本開示の一実施形態による、ターゲットキーワードに対する複数の基準入力音声を処理することによってターゲットキーワードに対するエントリ状態の数を決定するための、ボイスアクティベーションユニット 2 5 6 によって実行される方法 1 2 0 0 のフローチャートである。最初に、1 2 1 0 において、ボイスアクティベーションユニット 2 5 6 は、フルデューティサイクルに基づいてターゲットキーワードに対する基準入力音声を受け取る。1 2 2 0 において、ボイスアクティベーションユニット 2 5 6 は、ターゲットキーワードの冒頭の部分と複数の後続の部分とに関連付けられた状態情報を取得する。

20

【 0 0 6 7 】

[0078]1 2 3 0 において、基準入力音声に対する複数の基準状態シーケンスが状態情報に基づいて決定される。1 2 4 0 において、ボイスアクティベーションユニット 2 5 6 は、図 1 1 に関して詳細に説明したように、基準状態シーケンスをバックトラックすることによって、基準状態シーケンスにおける複数の状態に対する複数の状態時間期間を決定する。次いで、1 2 5 0 において、ボイスアクティベーションユニット 2 5 6 は、状態時間期間とデューティサイクルの非アクティブ期間とに基づいて、複数のエントリ状態を決定する。

【 0 0 6 8 】

[0079]ボイスアクティベーションユニット 2 5 6 は、入力音声においてターゲットキーワードを検出するために、入力音声を可変持続時間の部分にセグメント化するように構成され得る。いくつかの実施形態では、ターゲットキーワードに関連付けられた状態の各々に対するそのような持続時間の考えられる範囲は、あらかじめ決定され得る。各状態に関連付けられた各持続時間について、スコア（たとえば、確率値（probability value））は最も大きいキーワードスコアを決定するために割り当てられ得る。

30

【 0 0 6 9 】

[0080]図 1 3 は、本開示の別の実施形態による、入力音声を受け取り、SMMの下で持続時間の所定の範囲に基づいて、入力音声においてターゲットキーワードを検出するように構成されたボイスアクティベーションユニット 2 5 6 のブロック図を示す。ボイスアクティベーションユニット 2 5 6 は、可変セグメント化ユニット 1 3 1 0 と、特徴抽出器 1 3 2 0 と、スコア決定ユニット 1 3 3 0 と、キーワード検出ユニット 1 3 4 0 とを含む。スコア決定ユニット 1 3 3 0 は、観測スコア決定ユニット 1 3 5 0 と、最大キーワードスコア決定ユニット 1 3 6 0 とを含む。

40

【 0 0 7 0 】

[0081]可変セグメント化ユニット 1 3 1 0 は、音声センサ（たとえば、音声センサ 2 1 0）から入力音声を受け取り、受け取られた入力音声を、複数の持続時間を有する複数のフレームにセグメント化する。たとえば、30ms 持続時間の入力音声を受け取られた場合、入力音声は、入力音声の 0ms から 10ms に位置する 10ms 持続時間の第 1 のフレームにセグメント化され得る。同様の方法で、入力音声は、10から 20ms の 10ms 持続時間の第 2 のフレームと、20ms から 30ms の 10ms 持続時間の第 3 のフレ

50

ームと、0msから20msの20ms持続時間の第4のフレームと、10msから30msの20ms持続時間の第5のフレームとにセグメント化され得る。

【0071】

[0082]可変セグメント化ユニット1310は、複数のセグメント化されたフレームを特徴抽出器1320に与え、入力音声フレームとしての入力音声全体（たとえば、上記の例の30ms持続時間）を特徴抽出器1320に与え得る。このプロセスでは、可変セグメント化ユニット1310はまた、特徴抽出器1320に、各フレームの持続時間とロケーションとについてのフレーム情報を与える。フレームとフレーム情報とを受け取ると、特徴抽出器1320は、任意の適切な特徴抽出方法を使用して、フレームの各々から音声特徴を抽出し、出力する。上記の例では、特徴抽出器1320は、合計で6つのフレームを受け取り、合計で6つの音声特徴をフレームから抽出する。

10

【0072】

[0083]スコア決定ユニット1330は、抽出された音声特徴に基づいて、キーワード検出ユニット1340に与えられるべき最も大きいキーワードスコアを生成するように構成される。スコア決定ユニット1330中の観測スコア決定ユニット1350は、特徴抽出器1320から音声特徴とフレーム情報とを受け取る。観測スコア決定ユニット1350はまた、ストレージユニット230からターゲットキーワードの状態情報を受け取る。状態情報は、図5に関して上記で説明したように、ターゲットキーワードに関連付けられた複数の状態と、状態の各々の確率モデルとを含む。

20

【0073】

[0084]この実施形態では、状態情報は、状態の各々についてあらかじめ決定された持続時間の範囲についての持続時間情報をさらに含む。各状態に対する持続時間の所定の範囲は、状態に対する、考えられる時間範囲に設定され得る。たとえば、あるユーザは比較的短い持続時間で状態「S」に対応する音声を話すことがあるが、別のユーザはその音声を話すのにより長くかかることがある。したがって、状態に対する持続時間の所定の範囲は、ユーザが通常、状態に関連付けられた音声を話すのにかかり得る持続時間の範囲を含むように設定され得る。

【0074】

[0085]指定された持続時間に関連付けられた各受け取られた音声特徴について、観測スコア決定ユニット1350は、受け取られた状態情報に基づいて状態の各々についての観測スコアを決定するかどうかについて決定し得る。一実施形態では、観測スコア決定ユニット1350は、持続時間を含むフレーム情報に基づいて、持続時間が各状態に対する持続時間の所定の範囲内にあるかどうかを決定する。持続時間が持続時間の所定の範囲内ないと決定された場合、観測スコア決定ユニット1350は観測スコアを決定しないと決定する。さもなければ、観測スコア決定ユニット1350は観測スコアを決定することに進む。

30

【0075】

[0086]観測スコアを決定する際、状態の各々の確率値は、関連する状態の確率モデルに従って計算され得る。計算された確率値は、関連する状態についての観測スコアとして使用され得る。次いで、観測スコア決定ユニット1350は、各音声特徴についての決定された観測スコアを、ターゲットキーワードについて、考えられる複数の状態シーケンスについてのキーワードスコアを決定するための最大キーワードスコア決定ユニット1360に与える。

40

【0076】

[0087]次いで、最大キーワードスコア決定ユニット1360は、音声特徴の各々についての観測スコアと、それぞれ音声特徴に関連付けられた、フレームの持続時間とロケーションとを含むフレーム情報とを受け取る。加えて、最大キーワードスコア決定ユニット1360は、（図5に関して上記で説明したように）遷移スコアを含む状態情報と、持続時間の所定の範囲を含む持続時間情報とを受け取る。一実施形態では、遷移スコアは、非キーワード状態から単一のエントリ状態への遷移スコアを含む。代替的に、遷移スコアは、

50

非キーワード状態から複数のエントリ状態の各々への遷移スコアを含み得る。

【0077】

[0088]持続時間情報は、状態の各々について、持続時間の各々について決定された持続時間スコアをさらに含む。各状態についての所定の持続時間スコアは、それぞれ、状態が持続時間に入る確率値に関連して設定され得る。フレーム情報と持続時間の所定の範囲とに基づいた、受け取られた観測スコアと、遷移スコアと、持続時間スコアとを使用した計算により、最大キーワードスコア決定ユニット1360は、考えられる状態シーケンスについてのキーワードスコアを決定する。決定されたキーワードスコアの中で、最大キーワードスコア決定ユニット1360は、ターゲットキーワードを検出するために使用するための最も大きいキーワードスコアを決定する。

10

【0078】

[0089]最大キーワードスコア決定ユニット1360は、持続時間が制限された半マルコフモデル方式に従ってキーワードスコアを決定し得る。たとえば、考えられる状態シーケンスのうちの1つがある状態を含み、その持続時間がその状態に対する持続時間の所定の範囲内でない場合、最大キーワードスコア決定ユニット1360は、その状態シーケンスについてのキーワードスコアを決定しないことがある一方で、そのような状態を含まない他の状態シーケンスについてのキーワードスコアを決定することがある。この例では、最大キーワードスコア決定ユニット1360は、ビタビアルゴリズムなどの任意の適切な方法を使用することによって、決定されたキーワードスコアの中から最も大きいキーワードスコアを次のように選択し得る。

20

【0079】

【数1】

$$V(t, s) = \max_{(dmin(s) \leq d \leq dmax(s), s')} \left( V(t-d, s') + T(s', s) + O(t, d, s) + D(d, s) \right) \quad (式1)$$

【0080】

ここで、 $t$ は現在の入力音声のフレームサイズを示し、 $s$ は現在の状態を表し、 $V(t, s)$ は時間フレーム1～ $t$ のうちの1つから始まる最後の状態としての状態 $s$ を有するすべての考えられる状態シーケンスについての最も大きいキーワードスコアを示し、 $d$ は状態 $s$ の持続時間を表し、 $dmin(s)$ は状態 $s$ に対する持続時間の所定の範囲の中で状態 $s$ の最も小さい持続時間を表し、 $dmax(s)$ は状態 $s$ に対する持続時間の所定の範囲の中で状態 $s$ の最も大きい持続時間を表し、 $s'$ は状態 $s$ の前の状態を表し、 $T(s', s)$ は前の状態 $s'$ から現在の状態 $s$ への遷移スコアを示し、 $O(t, d, s)$ は状態 $s$ についてのフレーム $t-d+1$ からフレーム $t$ の時間フレームについての観測スコアを表し、 $D(d, s)$ は状態 $s$ の持続時間が $d$ であるときの $s$ についての持続時間スコアを示す。

30

【0081】

[0090]SMM方式の下での式1に示すように、ターゲットキーワードについての最も大きいキーワードスコアは、各状態に対する持続時間の制限された範囲に基づいて計算される。そのような持続時間の制限された範囲を使用することによって、最大キーワードスコア決定ユニット1360は比較的低い計算複雑性で実装され得る。さらに、持続時間の制限された範囲は、ユーザの異なる発声速度(speaking speeds)を考慮してあらかじめ決定されるので、キーワード検出性能は大幅に劣化しない。

40

【0082】

[0091]最も大きいキーワードスコアを決定した後、最大キーワードスコア決定ユニット1360は最も大きいキーワードスコアをキーワード検出ユニット1340に与える。最大キーワードスコア決定ユニット1360から最も大きいキーワードスコアを受け取ると

50

、キーワード検出ユニット 1340 は、最も大きいキーワードスコアに基づいて、入力音声においてターゲットキーワードを検出する。たとえば、キーワード検出ユニット 1340 は、図 5 に関して上記で説明したように、ストレージユニット 230 からターゲットキーワードを検出するためのしきい値スコアを使用し得る。ターゲットキーワードを検出すると、キーワード検出ユニット 1340 は、機能またはアプリケーションをオンにするためのアクティベーション信号（「ON 信号」）を生成し、これを与え、機能またはアプリケーションはターゲットキーワードに関連付けられる。

#### 【0083】

【0092】一実施形態では、状態の各々に対する持続時間の範囲を決定するために、スコア決定ユニット 1330 は、ターゲットキーワードに対する複数の基準入力音声を受け取り得る。基準入力音声の各々について、図 11 および図 12 に関して上記で説明したのと同様の方法で、スコア決定ユニット 1330 は基準状態シーケンスを決定し、基準状態シーケンスをバックトラックすることによって、基準状態シーケンスにおける状態の各々に関連付けられた持続時間も決定する。したがって、スコア決定ユニット 1330 は、基準状態シーケンスからの状態に対する決定された持続時間に基づいて（たとえば、決定された持続時間を平均することによって）、各状態に対する持続時間の範囲を決定し得る。次いで、このようにして決定された持続時間の範囲は、ストレージユニット 230 に記憶され得る。

#### 【0084】

【0093】図 14 は、一実施形態による、ターゲットキーワードに関連付けられた複数の状態「S」、「T」、「A」、「R」、「T」などの各々に対する持続時間の範囲に基づいて生成された、マルコフ連鎖モデルの下でのターゲットキーワードに対する複数の状態シーケンスのブロック図 1400 を示す。状態の各々に対する持続時間の範囲は、図 13 に関して上記で説明した方法で制限されるようにあらかじめ決定され得る。説明を簡単にするために、図 1400 の時点（たとえば、「T1」）と次の時点（たとえば、「T2」）との間の期間は 10ms であり、セグメント化された音声特徴は  $10 \times N$  ms の持続時間を有し、ここで N は正整数であり、セグメント化された音声特徴の持続時間は入力音声の時間期間以下であると仮定され得る。

#### 【0085】

【0094】図示の実施形態では、前の状態から現在の状態への遷移ラインの横断時間（traverse times）は、状態の各々に対する所定の範囲における持続時間を示し得る。たとえば、時間 T4 および時間 T5 における状態「A」の場合、時間 T2 における前の状態「T」から時間 T4 および時間 T5 における現在の状態「A」への遷移ラインはそれぞれ、20ms および 30ms で横断される。この場合、状態「A」の持続時間の所定の範囲は、状態シーケンスにおいて 20ms から 30ms である。したがって、状態「A」の場合、20ms または 30ms のセグメント化された特徴についてのみ、観測スコアが決定され得、持続時間スコアがあらかじめ決定され得る。次いで、観測スコアおよびあらかじめ決定された持続時間スコアは、ターゲットキーワードを検出するためにキーワードスコアを決定し、次いでキーワードスコアの中から最も大きいキーワードスコアを決定するために使用され得る。したがって、ターゲットキーワードを検出する際に使用するためのキーワードスコアは、持続時間の制限された範囲に基づいて計算されるので、ターゲットキーワードを検出するための計算時間は、SMM 方式の下で大幅に低減され得る。

#### 【0086】

【0095】図 15 は、一実施形態による、SMM の下で持続時間の所定の範囲に基づいて入力音声においてターゲットキーワードを検出するための、スコア決定ユニット 1330 によって実行される方法 1500 のフローチャートである。最初に、1510 において、スコア決定ユニット 1330 は、可変時間フレームから抽出された音声特徴と、各フレームの持続時間とロケーションとについてのフレーム情報を受け取る。1520 において、スコア決定ユニット 1330 は、複数の状態と各状態の確率モデルとについての状態情報と、ターゲットキーワードについて、考えられる複数の状態シーケンスにおける状態の各々

10

20

30

40

50

から次の状態への遷移情報と、持続時間の所定の範囲と持続時間の各々について決定された持続時間スコアとについての持続時間情報とを受け取る。

【 0 0 8 7 】

[0096] 音声特徴の各々について、1530において、スコア決定ユニット1330は、各状態の確率モデルに基づいて、状態の各々の観測スコアを決定する。1540において、遷移情報に基づいて、スコア決定ユニット1330は、考えられる状態シーケンスにおける状態の各々から次の状態への遷移スコアを取得する。1550において、スコア決定ユニット1330は、持続時間の所定の範囲と、持続時間の各々について決定された持続時間スコアとを持続時間情報から取得する。1560において、フレーム情報と持続時間の所定の範囲とに基づいた、受け取られた観測スコアと、遷移スコアと、持続時間スコアとを使用した計算により、スコア決定ユニット1330は、考えられる状態シーケンスについてのキーワードスコアを決定し、決定されたキーワードスコアの中から最も大きいキーワードスコアを選択する。

10

【 0 0 8 8 】

[0097] 図16は、いくつかの実施形態による、機能を起動するために入力音声からターゲットキーワードを検出するための本開示の方法および装置が実装され得る、ワイヤレス通信システムにおけるモバイルデバイス1600のブロック図を示す。モバイルデバイス1600は、セルラーフォン、端末、ハンドセット、携帯情報端末(PDA)、ワイヤレスモデム、コードレスフォン、タブレットなどであり得る。ワイヤレス通信システムは、符号分割多元接続(CDMA)システム、モバイル通信用グローバルシステム(GSM(登録商標))システム、広帯域CDMA(W-CDMA(登録商標))システム、ロングタームエボリューション(LTE)システム、LTE Advancedシステムなどであり得る。

20

【 0 0 8 9 】

[0098] モバイルデバイス1600は、受信経路および送信経路を介して双方向通信を行うことが可能であり得る。受信経路上で、基地局によって送信された信号は、アンテナ1612によって受信され、受信機(RCVR)1614に与えられる。受信機1614は、受信信号を調整し、デジタル化し、さらなる処理のために調整およびデジタル化された信号をデジタルセクション1620に与える。送信経路上で、送信機(TMTR)は、デジタルセクション1620から送信されるべきデータを受信し、データを処理し、調整し、変調信号を生成し、変調信号はアンテナ1612を介して基地局に送信される。受信機1614および送信機1616は、CDMA、GSM、W-CDMA、LTE、LTE Advancedなどをサポートするトランシーバの一部である。

30

【 0 0 9 0 】

[0099] デジタルセクション1620は、たとえば、モデムプロセッサ1622、縮小命令セットコンピュータ/デジタル信号プロセッサ(RISC/DSP)1624、コントローラ/プロセッサ1626、内部メモリ1628、一般化オーディオエンコーダ1632、一般化オーディオデコーダ1634、グラフィックス/ディスプレイプロセッサ1636、および/または外部バスインターフェース(EBI)1638など、様々な処理、インターフェース、およびメモリユニットを含む。モデムプロセッサ1622は、データ送信および受信のための処理、たとえば、符号化、変調、復調、および復号を実行する。RISC/DSP1624は、モバイルデバイス1600のための一般的処理と特殊処理とを実行する。コントローラ/プロセッサ1626は、デジタルセクション1620内の様々な処理およびインターフェースユニットの動作を制御する。内部メモリ1628は、デジタルセクション1620内の様々なユニットのためのデータおよび/または命令を記憶する。

40

【 0 0 9 1 】

[00100] 一般化オーディオエンコーダ1632は、オーディオソース1642、マイクロフォン1643などからの入力信号に対して符号化を実行する。一般化オーディオデコーダ1634は、コード化オーディオデータに対して復号を実行し、出力信号をスピーカ

50

ーノヘッドセット１６４４に与える。一般化オーディオエンコーダ１６３２および一般化オーディオデコーダ１６３４は、必ずしも、オーディオソース、マイクロフォン１６４３およびスピーカーノヘッドセット１６４４とのインターフェースのために必要とされとは限らず、したがって、モバイルデバイス１６００に示されていないことに留意されたい。グラフィックスノディスプレイプロセッサ１６３６は、ディスプレイユニット１６４６に提示されるグラフィックス、ビデオ、画像、およびテキストのための処理を実行する。ＥＢＩ１６３８は、デジタルセクション１６２０とメインメモリ１６４８との間のデータの転送を可能にする。

【００９２】

[00101]デジタルセクション１６２０は、１つまたは複数のプロセッサ、ＤＳＰ、マイクロプロセッサ、ＲＩＳＣなどを用いて実装される。デジタルセクション１６２０はまた、１つまたは複数の特定用途向け集積回路（ＡＳＩＣ）および／または何らかの他のタイプの集積回路（ＩＣ）上に作製される。

【００９３】

[00102]一般に、本明細書で説明する任意のデバイスは、ワイヤレスフォン、セルラーフォン、ラップトップコンピュータ、ワイヤレスマルチメディアデバイス、ワイヤレス通信パーソナルコンピュータ（ＰＣ）カード、ＰＤＡ、外部または内部モデム、ワイヤレスチャネルを介して通信するデバイスなど、様々なタイプのデバイスを示す。デバイスは、アクセス端末（ＡＴ）、アクセスユニット、加入者ユニット、移動局、クライアントデバイス、モバイルユニット、モバイルフォン、モバイル、リモート局、リモート端末、リモートユニット、ユーザデバイス、ユーザ機器、ハンドヘルドデバイスなど、様々な名前を有し得る。本明細書で説明する任意のデバイスは、命令とデータとを記憶するためのメモリ、ならびにハードウェア、ソフトウェア、ファームウェア、またはそれらの組合せを有し得る。

以下に本願の出願当初の特許請求の範囲に記載された発明を付記する。

【Ｃ１】

ターゲットキーワードを検出するための方法であって、前記ターゲットキーワードが、冒頭の部分と複数の後続の部分とを含み、前記方法が、

電子デバイスにおいて、前記ターゲットキーワードの前記後続の部分のうちの１つから始まる入力音声を受け取ることと、

前記入力音声から音声特徴を抽出することと、

状態ネットワークを記述しているデータを取得することと、ここにおいて、前記状態ネットワークは、単一の開始状態と、複数のエントリ状態と、前記単一の開始状態から前記複数のエントリ状態の各々への遷移とを含む、

前記抽出された音声特徴と前記状態ネットワークとに基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することとを備える、方法。

【Ｃ２】

前記入力音声を受け取ることが、前記電子デバイスに関連付けられたデューティサイクルに基づいて入力音声ストリームを受け取ることを備える、Ｃ１に記載の方法。

【Ｃ３】

前記エントリ状態に対応するデータが、

前記電子デバイスのフルデューティサイクルに基づいて、前記ターゲットキーワードの前記冒頭の部分と前記複数の後続の部分とに対応する基準入力音声を受け取ることと、

前記基準入力音声に対する複数の基準状態シーケンスを決定することと、

前記基準状態シーケンスにおける複数の状態に対する状態時間期間を決定することと、

前記状態時間期間と前記フルデューティサイクルの非アクティブ期間とに基づいて、前記エントリ状態を決定することとによって前記電子デバイスに記憶される、Ｃ２に記載の方法。

【Ｃ４】

前記基準状態シーケンスにおける前記複数の状態に対する前記状態時間期間が、前記基

10

20

30

40

50



準状態シーケンスをバックトラックすることによって決定される、C 3 に記載の方法。

[ C 5 ]

前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、複数のキーワードスコアを決定することを備え、各キーワードスコアが、前記単一の開始状態から前記複数のエン트리状態のうちの1つへの遷移を含むそれぞれの状態シーケンスに対応する、C 1 に記載の方法。

[ C 6 ]

前記状態ネットワークが、複数の状態シーケンスを含み、複数の状態を含む各状態シーケンスが、前記単一の開始状態と、前記複数のエン트리状態のうちの1つと、1つまたは複数の後続の状態とを備える、C 5 に記載の方法。

[ C 7 ]

前記複数の状態シーケンスの各状態シーケンスが、隠れマルコフモデルと、前記状態シーケンスの前記状態についての遷移情報とに関連付けられる、C 6 に記載の方法。

[ C 8 ]

前記キーワードスコアを決定することが、  
前記状態ネットワークに基づいて、前記抽出された音声特徴の各々についての前記状態の各々の観測スコアを決定することと、

前記状態ネットワークの遷移情報に基づいて、前記状態シーケンスの各々における前記状態の各々から次の状態への遷移スコアを取得することとを備え、

前記キーワードスコアが、前記観測スコアと前記遷移スコアとに基づいて決定される、C 6 に記載の方法。

[ C 9 ]

前記複数のキーワードスコアの中の最も大きいキーワードスコアが、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定するために使用される、C 5 に記載の方法。

[ C 10 ]

前記最も大きいキーワードスコアがしきい値スコアよりも大きい場合、前記入力音声の前記ターゲットキーワードに対応すると決定される、C 9 に記載の方法。

[ C 11 ]

前記状態ネットワークが非キーワード状態シーケンスを含み、前記複数のキーワードスコアを決定することが、前記非キーワード状態シーケンスについての非キーワードスコアを決定することを備える、C 5 に記載の方法。

[ C 12 ]

前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、

前記複数のキーワードスコアの中から最も大きいキーワードスコアを選択することと、

前記最も大きいキーワードスコアと前記非キーワードスコアとの間の差に基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することとをさらに備える、C 11 に記載の方法。

[ C 13 ]

前記差に基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、

前記差に基づいて、信頼性値を決定することと、

前記信頼性値がしきい値信頼性値よりも大きい場合、前記入力音声の前記ターゲットキーワードに対応すると決定することとを備える、C 12 に記載の方法。

[ C 14 ]

ターゲットキーワードを検出するための方法であって、前記ターゲットキーワードが複数の部分を含み、前記方法が、

電子デバイスにおいて入力音声を受け取ることと、

前記入力音声から音声特徴を抽出することと、

前記ターゲットキーワードの前記複数の部分に関連付けられた状態情報を取得すること

10

20

30

40

50

と、前記状態情報は、前記ターゲットキーワードの前記部分に関連付けられた複数の状態の各状態に対する持続時間範囲を含む、

前記抽出された音声特徴と前記状態情報とに基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することとを備える、方法。

[ C 1 5 ]

前記音声特徴を抽出することが、

前記入力音声をフレームにセグメント化することと、各フレームは持続時間を有する、

前記フレームから前記音声特徴を抽出することとを備える、C 1 4 に記載の方法。

[ C 1 6 ]

前記複数の状態が半マルコフモデルに関連付けられる、C 1 4 に記載の方法。

10

[ C 1 7 ]

前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することが、

前記音声特徴の各々について、前記状態情報に基づいて、前記複数の状態の各状態の観測スコアを決定することと、

前記音声特徴の各々について、前記状態情報に基づいて、前記複数の状態の各状態の持続時間スコアを取得することと、

遷移情報に基づいて、複数の状態シーケンスの各状態シーケンスにおける特定の状態から次の状態への遷移スコアを取得することと、

前記観測スコアと、前記持続時間スコアと、前記遷移スコアとに基づいて、前記複数の状態シーケンスの各状態シーケンスについてのキーワードスコアを決定することとを備える、C 1 6 に記載の方法。

20

[ C 1 8 ]

特定の状態シーケンスについての前記キーワードスコアを決定することが、前記状態シーケンスにおける状態の持続時間が前記状態に対する前記持続時間範囲内にあるかどうかを決定することとを備える、C 1 7 に記載の方法。

[ C 1 9 ]

各状態に対する前記持続時間範囲が、

前記ターゲットキーワードに対する基準入力音声を受け取ることと、

前記基準入力音声に対する基準状態シーケンスを決定することと、

前記基準状態シーケンスにおける状態に対する状態時間期間を決定することと、

前記状態時間期間に基づいて前記持続時間範囲を決定することとによって前記電子デバイスに記憶される、C 1 4 に記載の方法。

30

[ C 2 0 ]

冒頭の部分と複数の後続の部分とを含むターゲットキーワードを検出するための電子デバイスであって、

前記ターゲットキーワードの前記後続の部分のうちの1つから始まる入力音声を受け取るように構成された音声センサと、

前記入力音声から音声特徴を抽出し、状態ネットワークを記述しているデータを取得し、前記抽出された音声特徴と前記状態ネットワークとに基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定するように構成されたボイスアクティベーションユニットとを備え、

40

前記状態ネットワークが、単一の開始状態と、複数のエントリ状態と、前記単一の開始状態から前記複数のエントリ状態の各々への遷移とを含む、電子デバイス。

[ C 2 1 ]

前記ボイスアクティベーションユニットが、複数のキーワードスコアを決定することによって、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定するように構成され、各キーワードスコアが、前記単一の開始状態から前記複数のエントリ状態のうちの1つへの遷移を含むそれぞれの状態シーケンスに対応する、C 2 0 に記載の電子デバイス。

[ C 2 2 ]

50

前記状態ネットワークが、複数の状態シーケンスを含み、複数の状態を含む各状態シーケンスが、前記単一の開始状態と、前記複数のエントリ状態のうちの1つと、1つまたは複数の後続の状態とを備える、C 2 1に記載の電子デバイス。

[ C 2 3 ]

前記ボイスアクティベーションユニットが、  
前記状態ネットワークに基づいて、前記抽出された音声特徴の各々についての前記状態の各々の観測スコアを決定することと、

前記状態ネットワークの遷移情報に基づいて、前記状態シーケンスの各々における前記状態の各々から次の状態への遷移スコアを取得することとによって、前記キーワードスコアを決定するように構成され、

前記キーワードスコアが、前記観測スコアと前記遷移スコアとに基づいて決定される、C 2 2に記載の電子デバイス。

[ C 2 4 ]

前記複数のキーワードスコアの中の最も大きいキーワードスコアが、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定するために使用され、前記最も大きいキーワードスコアがしきい値スコアよりも大きい場合、前記入力音声の前記ターゲットキーワードに対応すると決定される、C 2 1に記載の電子デバイス。

[ C 2 5 ]

前記状態ネットワークが非キーワード状態シーケンスを含み、前記複数のキーワードスコアを決定することが、前記非キーワード状態シーケンスについての非キーワードスコアを決定することを備える、C 2 1に記載の電子デバイス。

[ C 2 6 ]

前記ボイスアクティベーションユニットが、  
前記複数のキーワードスコアの中から最も大きいキーワードスコアを選択することと、  
前記最も大きいキーワードスコアと前記非キーワードスコアとの間の差に基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することとによって、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定するように構成される、C 2 5に記載の電子デバイス。

[ C 2 7 ]

複数の部分を含むターゲットキーワードを検出するための電子デバイスであって、  
入力音声を受け取るように構成された音声センサと、  
前記入力音声から音声特徴を抽出することと、前記ターゲットキーワードの前記複数の部分に関連付けられた状態情報を取得することと、前記状態情報は、前記ターゲットキーワードの前記部分に関連付けられた複数の状態の各状態に対する持続時間範囲を含む、前記音声特徴と前記状態情報とに基づいて、前記入力音声の前記ターゲットキーワードに対応するかどうかを決定することとを行うように構成されたボイスアクティベーションユニットとを備える、電子デバイス。

[ C 2 8 ]

前記ボイスアクティベーションユニットが、  
前記入力音声をフレームにセグメント化するように構成されたセグメント化ユニットと、各フレームは持続時間を有する、  
前記フレームから前記音声特徴を抽出するように構成された特徴抽出器とを備える、C 2 7に記載の電子デバイス。

[ C 2 9 ]

前記ボイスアクティベーションユニットが、  
前記音声特徴の各々について、前記状態情報に基づいて、前記複数の状態の各状態の観測スコアを決定することと、

前記音声特徴の各々について、前記状態情報に基づいて、前記複数の状態の各状態の持続時間スコアを取得することと、

遷移情報に基づいて、複数の状態シーケンスの各状態シーケンスにおける特定の状態か

10

20

30

40

50

ら次の状態への遷移スコアを取得することと、

前記観測スコアと、前記持続時間スコアと、前記遷移スコアとに基づいて、前記複数の状態シーケンスの各状態シーケンスについてのキーワードスコアを決定することによって、前記入力音声が入記ターゲットキーワードに対応するかどうかを決定するように構成される、C 2 7に記載の電子デバイス。

[ C 3 0 ]

前記ボイスアクティベーションユニットが、前記状態シーケンスにおける状態の持続時間が前記状態に対する前記持続時間範囲内にあるかどうかを決定することによって、特定の状態シーケンスについての前記キーワードスコアを決定するように構成される、C 2 9に記載の電子デバイス。

10

【図 1】

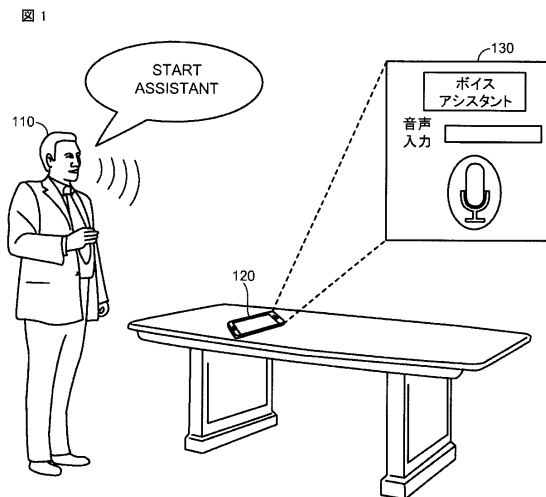


FIG. 1

【図 2】

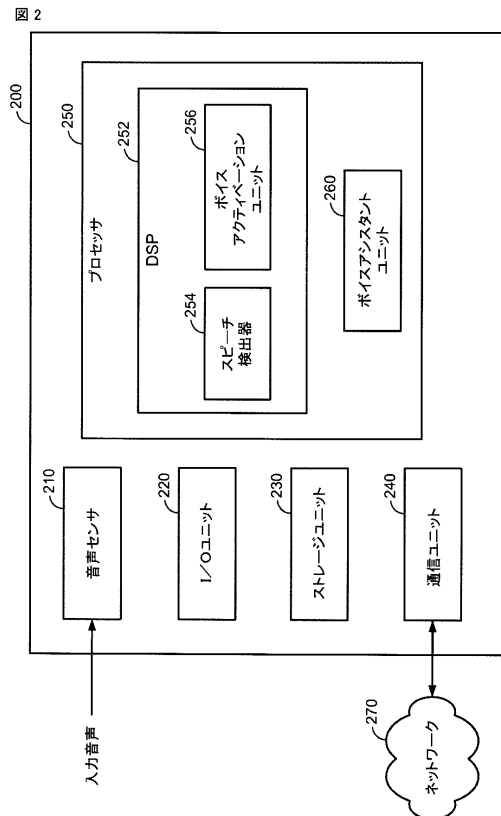


FIG. 2

【図 3】

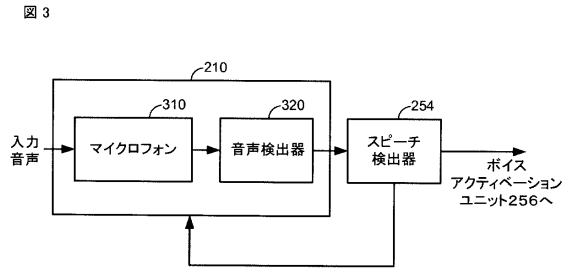


FIG. 3

【図 4】

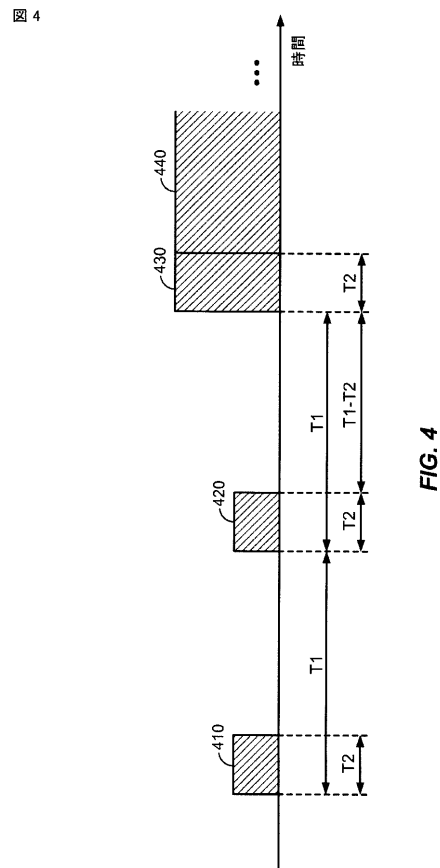


FIG. 4

【図 5】

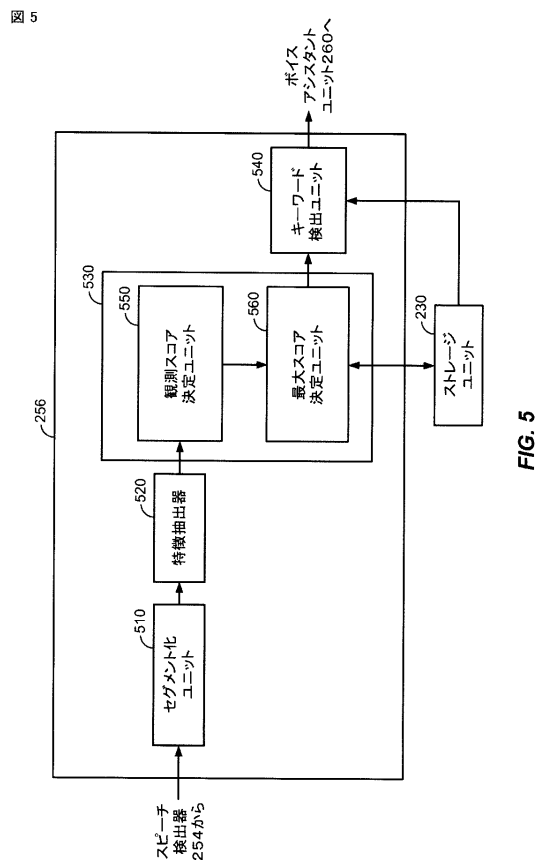


FIG. 5

【図 6】

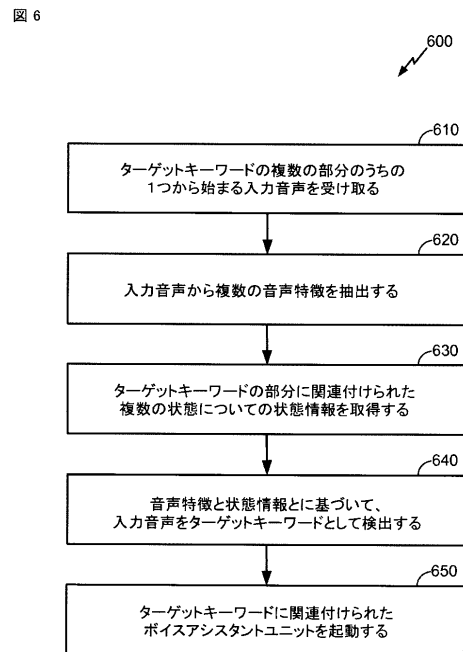
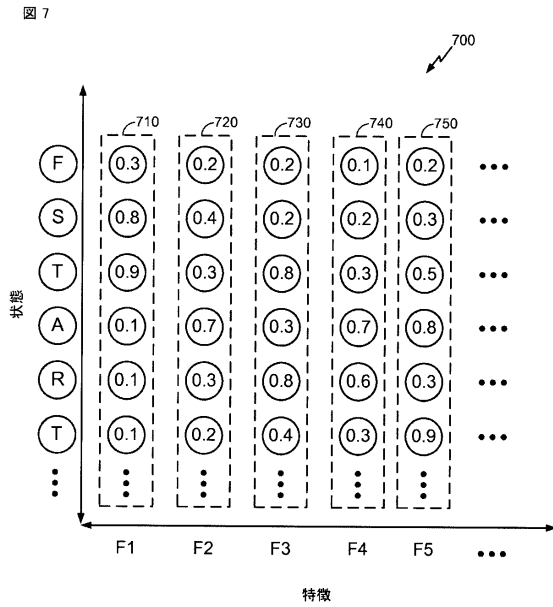
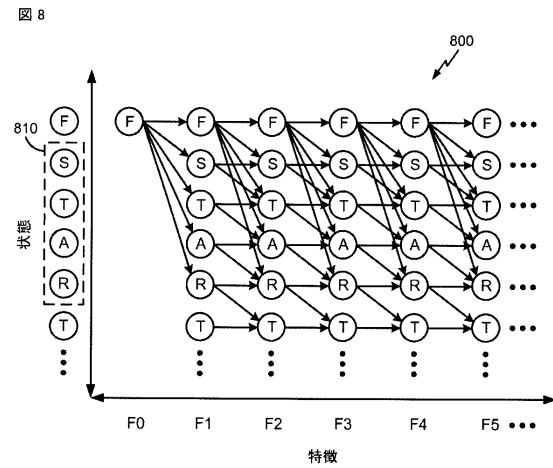


FIG. 6

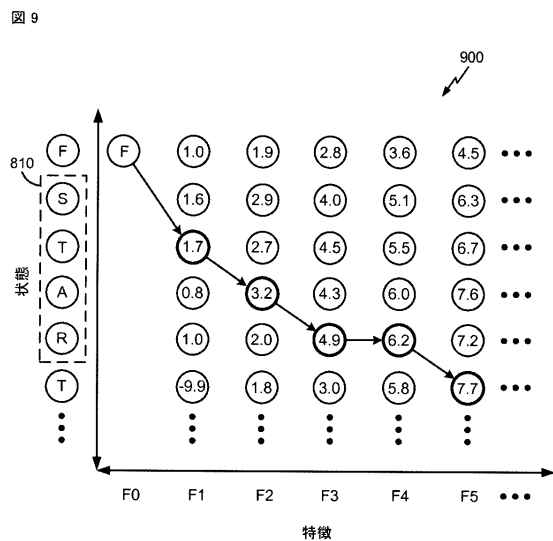
【図 7】



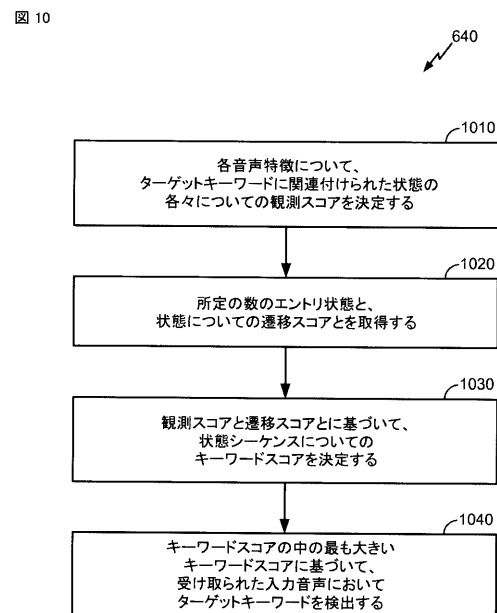
【図 8】



【図 9】



【図 10】



【図 1 1】

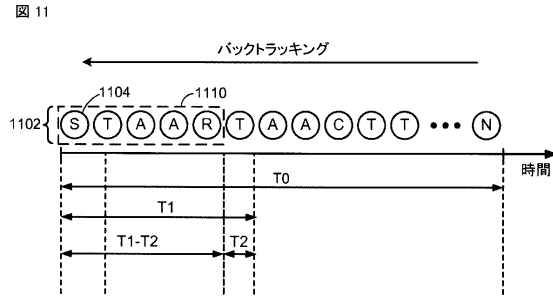


FIG. 11

【図 1 2】

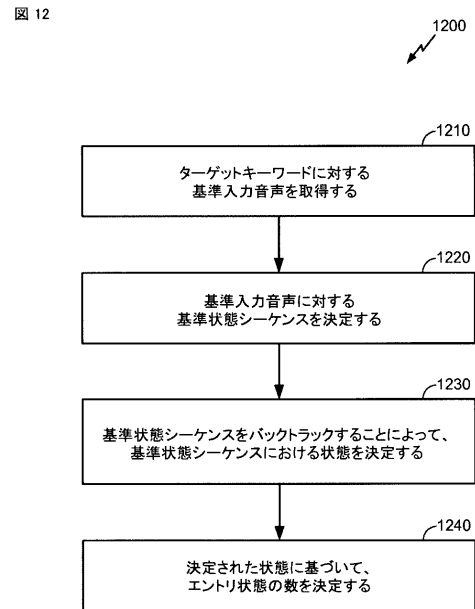


FIG. 12

【図 1 3】

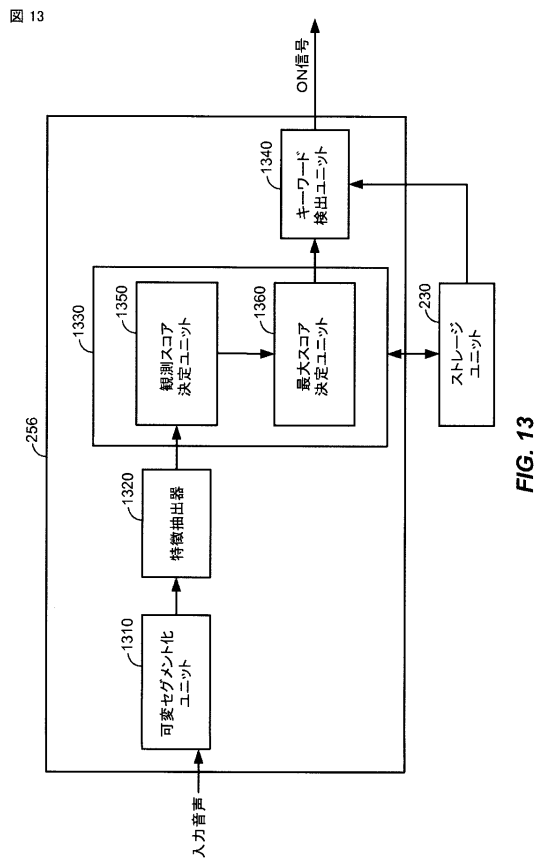


FIG. 13

【図 1 4】

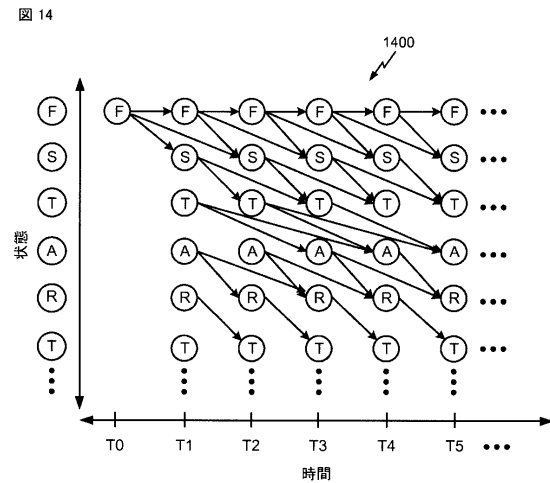


FIG. 14

【図 15】

図 15

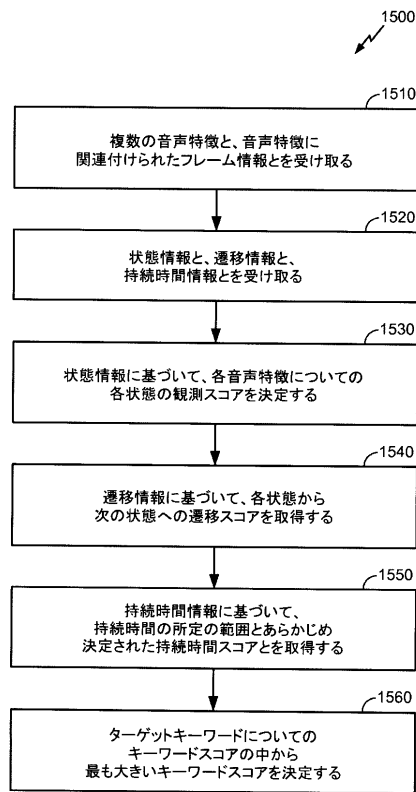


FIG. 15

【図 16】

図 16

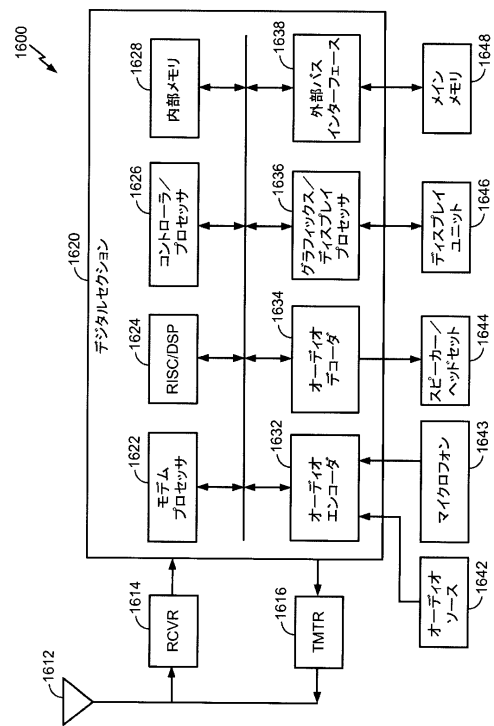


FIG. 16



## フロントページの続き

(31)優先権主張番号 14/087,939

(32)優先日 平成25年11月22日(2013.11.22)

(33)優先権主張国 米国(US)

## 早期審査対象出願

(72)発明者 キム、ソンウン

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 リ、ミンスブ

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 キム、テス

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ジン、ミンホ

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ホワン、キュ・ウォン

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

審査官 菊池 智紀

(56)参考文献 特開平10-312194(JP,A)

特開平5-11798(JP,A)

特開2000-89792(JP,A)

(58)調査した分野(Int.Cl.,DB名)

G10L 15/00-15/34