

(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2005-275410

(P2005-275410A)

(43) 公開日 平成17年10月6日(2005.10.6)

(51) Int. Cl. ⁷	F I	テーマコード (参考)
G 1 0 L 11/00	G 1 0 L 11/00 2 0 1 A	5 D 0 1 5
G 1 0 L 15/20	G 1 0 L 21/02 1 0 1 B	
G 1 0 L 21/02	G 1 0 L 21/02 1 0 3 Z	
	G 1 0 L 3/02 3 0 1 Z	
	G 1 0 L 3/02 3 0 1 D	
審査請求 未請求 請求項の数 26 O L (全 21 頁)		

(21) 出願番号 特願2005-85040 (P2005-85040)
 (22) 出願日 平成17年3月23日 (2005. 3. 23)
 (31) 優先権主張番号 60/555, 582
 (32) 優先日 平成16年3月23日 (2004. 3. 23)
 (33) 優先権主張国 米国 (US)

(71) 出願人 504064319
 ハーマン ベッカー オートモーティブ
 システムズ-ウェイクメーカーズ, イン
 コーポレイテッド
 カナダ国 プイ6ビー 2ケー4 プリテ
 イッシュ コロンビア, バンクーバー,
 アボット ストリート ナンバー302
 -134
 (74) 代理人 100078282
 弁理士 山本 秀策
 (74) 代理人 100062409
 弁理士 安村 高明
 (74) 代理人 100113413
 弁理士 森下 夏樹

最終頁に続く

(54) 【発明の名称】 ニューラルネットワークを利用してスピーチ信号を分離する。

(57) 【要約】

【課題】 背景ノイズの存在において、スピーチ信号を分離し、再構築する分離スピーチ信号システムを提供する。

【解決手段】 スピーチ信号の周波数コンポーネントが、背景ノイズによってマスクされる環境において送信されるスピーチ信号を分離し、再構築するように構成されているスピーチ信号分離システム。スピーチ信号分離システム(10)は、オーディオソースからノイジーなスピーチ信号を取得する。ノイジーなスピーチ信号は、それから、背景ノイズからクリーンなスピーチ信号を分離し、再構築するように訓練されたニューラルネットワーク(20)を介して供給される。ノイジーなスピーチ信号が、ニューラルネットワーク(20)を介して供給されると、スピーチ信号分離システム(10)は、大幅に減少したノイズを有する推定されたスピーチ信号を生成する。

【選択図】 図14

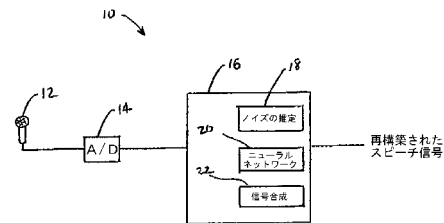


Figure 14

【特許請求の範囲】

【請求項 1】

オーディオ信号における背景ノイズからスピーチ信号を抽出するスピーチ信号分離システムであって、

複数の周波数に渡りオーディオ信号の背景ノイズの強度を推定するように適合された背景ノイズ推定コンポーネントと、

該背景ノイズからスピーチ推定信号を抽出するように適合されたニューラルネットワークコンポーネントと、

該背景ノイズの強度推定に基づいて該オーディオ信号および該抽出されたスピーチから再構築されたスピーチ信号を生成する合成コンポーネントと

を備えた、システム。

10

【請求項 2】

時系列の信号から周波数領域の信号に前記オーディオ信号を変換する周波数変換コンポーネントをさらに備えた、請求項 1 に記載のシステム。

【請求項 3】

周波数サブバンドの減少した数を有する圧縮されたオーディオ信号を生成する圧縮コンポーネントをさらに備えた、請求項 2 に記載のシステム。

【請求項 4】

前記ニューラルネットワークは、前記圧縮されたオーディオ信号における周波数サブバンドの数と等しい第 1 のセットの入力ノードであって、該圧縮されたオーディオ信号を受信する第 1 のセットの入力ノードを有する、請求項 3 に記載のシステム。

20

【請求項 5】

前記ニューラルネットワークは、周波数サブバンドの数と等しい第 2 のセットの入力ノードであって、前記背景ノイズの推定を受信する第 2 のセットの入力ノードを有する、請求項 4 に記載のシステム。

【請求項 6】

前記ニューラルネットワークは、前記圧縮されたオーディオ信号における周波数サブバンドの数と等しい第 2 のセットの入力ノードであって、以前の時間ステップから該圧縮されたオーディオ信号を受信する第 2 のセットの入力ノードを有する、請求項 4 に記載のシステム。

30

【請求項 7】

前記ニューラルネットワークは、前記圧縮されたオーディオ信号における周波数サブバンドの数と等しい第 2 のセットの入力ノードであって、以前の時間ステップから該ニューラルネットワークの出力を受信する第 2 のセットの入力ノードを有する、請求項 4 に記載のシステム。

【請求項 8】

前記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから中間結果を受信する第 2 のセットの入力ノードを有する、請求項 4 に記載のシステム。

【請求項 9】

合成コンポーネントは、前記背景ノイズの推定より大きい強度を有するオーディオ信号の一部分を該背景ノイズの推定より小さい強度を有する該オーディオ信号の一部分に対応する前記抽出されたスピーチの一部分と組み合わせるように適合された、請求項 1 に記載のシステム。

40

【請求項 10】

スピーチコンポーネントおよび背景ノイズを有するオーディオ信号からスピーチ信号を分離する方法であって、

時系列のオーディオ信号を周波数領域に変換することと、

複数の周波数帯域に渡り、該オーディオ信号における該背景を推定することと、

該オーディオ信号からスピーチ信号の推定を抽出することと、

50

該背景ノイズの推定に基づいてスピーチ信号の推定の一部分を該オーディオ信号の一部分と合成することにより、減少した背景ノイズを有する再構築されたスピーチ信号を提供することと

を包含した、方法。

【請求項 1 1】

前記オーディオ信号からスピーチ信号の推定を抽出することは、該オーディオ信号をニューラルネットワークへの入力として割り当てることを包含する、請求項 1 0 に記載の方法。

【請求項 1 2】

前記スピーチ信号の推定を前記オーディオ信号と合成することは、前記背景ノイズの推定より大きい、強度の上限しきい値を確立し、かつ、該強度の上限しきい値より大きい強度値を有する該オーディオ信号の一部分を該スピーチ信号の推定の一部分と組み合わせることを包含する、請求項 1 0 に記載の方法。

10

【請求項 1 3】

前記スピーチ信号の推定を前記オーディオ信号と合成することは、前記背景ノイズの推定であるか、もしくは付近の強度の下限しきい値を確立し、かつ、該強度の下限しきい値より小さい、強度値を有する該オーディオ信号の一部分に対応する該スピーチ信号の推定の一部分と組み合わせることを包含する、請求項 1 0 に記載の方法。

【請求項 1 4】

前記スピーチ信号の推定を前記オーディオ信号と合成することは、強度の上限および下限しきい値を確立し、かつ、該オーディオ信号の一部分を前記上限の強度のしきい値と前記下限のしきい値との間の強度値を有する該オーディオ信号の一部分に対応する該スピーチ信号の推定の一部分と組み合わせることを包含する、請求項 1 0 に記載の方法。

20

【請求項 1 5】

前記オーディオ信号の前記一部分を前記スピーチ信号の推定の一部分と組み合わせることは、該スピーチ信号の推定が、前記強度の下限しきい値に近い強度値を有する該オーディオ信号の一部分に対する該オーディオ信号より重みを置かれ、かつ、該オーディオ信号が、前記強度の上限しきい値に近い強度値を有する該オーディオ信号の一部分に対する該スピーチ信号の推定より重みを置かれるように、該オーディオ信号および該スピーチ信号に重みを置くことを包含する、請求項 1 4 に記載の方法。

30

【請求項 1 6】

前記背景ノイズの推定を前記ニューラルネットワークに供給することをさらに包含する、請求項 1 1 に記載の方法。

【請求項 1 7】

以前の時間ステップからの前記スピーチ信号の推定を前記ニューラルネットワークに供給することをさらに包含する、請求項 1 1 に記載の方法。

【請求項 1 8】

以前の時間ステップからの前記スピーチ信号の推定の中間結果を前記ニューラルネットワークに供給することをさらに包含する、請求項 1 1 に記載の方法。

【請求項 1 9】

以前の時間ステップからの前記オーディオ信号を前記ニューラルネットワークに供給することをさらに包含する、請求項 1 1 に記載の方法。

40

【請求項 2 0】

スピーチ信号をエンハンスするシステムであって、
スピーチコンテンツおよび背景ノイズの両方を有する時系列のオーディオ信号を提供するオーディオ信号出力ソースと、

時系列領域から周波数領域に該オーディオ信号を変換する周波数変換機能を提供する信号プロセッサと、

背景ノイズの推定器と、

ニューラルネットワークと、

50

信号コンバイナと

を備え、

該背景の推定器は、該オーディオ信号における該背景ノイズの推定を形成し、該ニューラルネットワークは、該オーディオ信号から、該スピーチ信号の推定を抽出し、該信号コンバイナは、該背景ノイズの推定に基づいて該スピーチ信号の推定を該オーディオ信号と組み合わせることにより、大幅に減少した背景ノイズを有する再構築されたスピーチ信号を生成する、システム。

【請求項 2 1】

前記ニューラルネットワークは、第 1 のセットの入力ノードであって、前記オーディオ信号を受信する第 1 のセットの入力ノードを包含した、請求項 2 0 に記載の方法。

10

【請求項 2 2】

前記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから前記オーディオ信号を受信する第 2 のセットの入力ノードを包含した、請求項 2 1 に記載の方法。

【請求項 2 3】

前記ニューラルネットワークは、第 2 のセットの入力ノードであって、前記背景ノイズの推定を受信する第 2 のセットの入力ノードを包含した、請求項 2 1 に記載の方法。

【請求項 2 4】

前記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから前記スピーチ信号の推定を受信する第 2 のセットの入力ノードを包含した、請求項 2 1 に記載の方法。

20

【請求項 2 5】

前記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから中間結果を受信する第 2 のセットの入力ノードを包含した、請求項 2 1 に記載の方法。

【請求項 2 6】

背景ノイズからスピーチ信号を分離する方法であって、

オーディオ信号を受信することと、

信号の正確さが、高い確実性を有すると知られている該オーディオ信号の一部分を識別することと、

30

ニューラルネットワークを訓練することにより、該オーディオ信号の正確さが不確かである該オーディオ信号の一部分に対して、著しく減少した背景ノイズを有する再構築された信号を推定することと

を包含する、方法。

【発明の詳細な説明】

【技術分野】

【0001】

(関連出願)

本出願は、2004年3月23日付けで出願された米国仮特許出願第60/555,582号の利益をクレームする。

40

【0002】

本発明は、概してスピーチ処理システム分野に関し、詳細には、ノイジーなサウンド環境におけるスピーチ信号の検出および分離に関する。

【背景技術】

【0003】

音は、固体、液体もしくは気体の任意の弾性材料を介して、送信される振動である。1つのタイプの共通の音は、人間のスピーチである。ノイジーな環境において、スピーチ信号を送信するとき、信号は、しばしば背景ノイズによってマスクされる。音は、周波数によって特徴付けられる。周波数は、時間単位上で起こる周期的な処理の完全なサイクルの数として定義される。信号は、時間を表すX軸および振幅を表すY軸に対してプロットさ

50

れる。典型的な信号は、その発生源から正のピークに上昇し、それから、負のピークへ下降する。信号は、それから、その初期の振幅へ戻り、それによって、第1の周期を完成させる。正弦波信号の周期は、信号が繰り返される間隔である。

【0004】

周波数は、一般的にヘルツ(Hz)で測定される。典型的な人間の耳は、20Hz~20,000Hzの周波数範囲の音を検出できる。音は、多くの周波数から成り得る。多重周波数サウンドの振幅は、各時間サンプルでの構成周波数の振幅の合計である。2つ以上の周波数が、調波関係によって互いに関連し得る。第1の周波数は、その第1の周波数が、第2の周波数の整数倍であるとき、第2の周波数の調波である。

【0005】

多重周波数サウンドは、その多重周波数サウンドを含む周波数パターンに従って特徴付けられる。一般的に、ノイズは、ある角度で周波数プロットにおいて低下する。この周波数パターンは、「ピンクノイズ」と名付けられる。ピンクノイズは、高強度の低周波数信号から成る。周波数が増加するにつれて、音の強度は減少する。「ブラウンノイズ」は、「ピンクノイズ」と同様であるが、より早い低下を示す。ブラウンノイズは、車両の音(例えば、ボディパネルから出る傾向のある低周波数ランブル)において見つけられ得る。すべての周波数で、同等のエネルギーを示す音は、「ホワイトノイズ」と呼ばれる。

10

【0006】

音は、また、通常、デシベル(dB)で測定される、その強度によって特徴付けられ得る。デシベルは、音の強度の対数単位であり、つまり音の強度のいくつかのリファレンス強度に対する比率の対数の10倍である。人間の聴力に対して、デシベルの大きさは、平均的な最小の知覚できる音に対するゼロ(dB)から、平均的な痛みのレベルのおよそ130(dB)で定義される。

20

【0007】

人間の音声は、声門で生成される。声門は、喉頭の上部での声帯間の開口部である。人間の声の音は、振動する声帯を介して、呼気によって作成される。声門の振動の周波数が、これらの音の特徴付ける。大半音声は、70Hz~400Hzの範囲に入る。典型的な男性は、およそ80Hz~150Hzの周波数範囲で話す。典型的な女性は、通常、125Hz~400Hzの周波数範囲で話す。

【0008】

人間のスピーチは、子音および母音から成る。「TH」および「F」といった子音は、ホワイトノイズによって特徴付けられる。これらの音の周波数スペクトラムは、卓上の扇風機と同様である。子音「S」は、通常、およそ3000Hzから始まり、およそ10,000Hzにまで及ぶ広帯域ノイズによって特徴付けられる。子音「T」、「B」および「P」は、「破裂音」と呼ばれ、また広帯域ノイズによって特徴付けられる。破裂音は、時間における急上昇によって「S」とは異なる。母音は、また一意の周波数スペクトラムを生成する。母音のスペクトラムは、フォルマント周波数によって特徴付けられる。フォルマントは、一意である母音のいくつかの共鳴帯域を含み得る。

30

【0009】

スピーチ検出および記録における大きな問題は、背景ノイズからのスピーチ信号の分離である。背景ノイズは、スピーチ信号に干渉し、低下させ得る。ノイジーな環境において、スピーチ信号の多くの周波数コンポーネントは、部分的にもしくは全体的にでさえ、背景ノイズの周波数によってマスクされ得る。

40

【発明の開示】

【発明が解決しようとする課題】

【0010】

従って、背景ノイズの存在において、スピーチ信号を分離し、再構築する分離スピーチ信号システムを提供する。

【課題を解決するための手段】

【0011】

50

本発明は、スピーチ信号の周波数コンポーネントが、背景ノイズによってマスクされる環境において、送信されるスピーチ信号を分離し、かつ、再構築することが可能であるスピーチ信号分離システムを開示する。本発明の1つの例において、ノイズなスピーチ信号が、ニューラルネットワークによって分析される。ニューラルネットワークは、クリーンなスピーチ信号を作成するように動作可能である。ニューラルネットワークは、背景ノイズから、スピーチ信号を分離するように訓練される。

【0012】

本発明の他のシステム、方法、特徴および利点が、以下の図面および詳細な記載の検討により当業者に明らかになる。すべてのこのような追加的なシステム、方法、特徴および利点が記載内および本発明の範囲内に含まれ、また請求項によって保護されることが意図される。

10

(項目1)

オーディオ信号における背景ノイズからスピーチ信号を抽出するスピーチ信号分離システムであって、

複数の周波数に渡りオーディオ信号の背景ノイズの強度を推定するように適合された背景ノイズ推定コンポーネントと、

上記背景ノイズからスピーチ推定信号を抽出するように適合されたニューラルネットワークコンポーネントと、

上記背景ノイズの強度推定に基づいて上記オーディオ信号および上記抽出されたスピーチから再構築されたスピーチ信号を生成する合成コンポーネントと

20

を備えた、システム。

(項目2)

時系列の信号から周波数領域の信号に上記オーディオ信号を変換する周波数変換コンポーネントをさらに備えた、項目1に記載のシステム。

(項目3)

周波数サブバンドの減少した数を有する圧縮されたオーディオ信号を生成する圧縮コンポーネントをさらに備えた、項目2に記載のシステム。

(項目4)

上記ニューラルネットワークは、上記圧縮されたオーディオ信号における周波数サブバンドの数と等しい第1のセットの入力ノードであって、上記圧縮されたオーディオ信号を受信する第1のセットの入力ノードを有する、項目3に記載のシステム。

30

(項目5)

上記ニューラルネットワークは、周波数サブバンドの数と等しい第2のセットの入力ノードであって、上記背景ノイズの推定を受信する第2のセットの入力ノードを有する、項目4に記載のシステム。

(項目6)

上記ニューラルネットワークは、上記圧縮されたオーディオ信号における周波数サブバンドの数と等しい第2のセットの入力ノードであって、以前の時間ステップから上記圧縮されたオーディオ信号を受信する第2のセットの入力ノードを有する、項目4に記載のシステム。

40

(項目7)

上記ニューラルネットワークは、上記圧縮されたオーディオ信号における周波数サブバンドの数と等しい第2のセットの入力ノードであって、以前の時間ステップから上記ニューラルネットワークの出力を受信する第2のセットの入力ノードを有する、項目4に記載のシステム。

(項目8)

上記ニューラルネットワークは、第2のセットの入力ノードであって、以前の時間ステップから中間結果を受信する第2のセットの入力ノードを有する、項目4に記載のシステム。

(項目9)

50

合成コンポーネントは、上記背景ノイズの推定より大きい強度を有するオーディオ信号の一部分を上記背景ノイズの推定より小さい強度を有する上記オーディオ信号の一部分に対応する上記抽出されたスピーチの一部分と組み合わせるように適合された、項目 1 に記載のシステム。

(項目 10)

スピーチコンポーネントおよび背景ノイズを有するオーディオ信号からスピーチ信号を分離する方法であって、

時系列のオーディオ信号を周波数領域に変換することと、

複数の周波数帯域に渡り、上記オーディオ信号における上記背景を推定することと、

上記オーディオ信号からスピーチ信号の推定を抽出することと、

10

上記背景ノイズの推定に基づいてスピーチ信号の推定の一部分を上記オーディオ信号の一部分と合成することにより、減少した背景ノイズを有する再構築されたスピーチ信号を提供することと

を包含した、方法。

(項目 11)

上記オーディオ信号からスピーチ信号の推定を抽出することは、上記オーディオ信号をニューラルネットワークへの入力として割り当てることを包含する、項目 10 に記載の方法。

(項目 12)

上記スピーチ信号の推定を上記オーディオ信号と合成することは、上記背景ノイズの推定より大きい、強度の上限しきい値を確立し、かつ、上記強度の上限しきい値より大きい強度値を有する上記オーディオ信号の一部分を上記スピーチ信号の推定の一部分と組み合わせることを包含する、項目 10 に記載の方法。

20

(項目 13)

上記スピーチ信号の推定を上記オーディオ信号と合成することは、上記背景ノイズの推定であるか、もしくは付近の強度の下限しきい値を確立し、かつ、上記強度の下限しきい値より小さい、強度値を有する上記オーディオ信号の一部分に対応する上記スピーチ信号の推定の一部分と組み合わせることを包含する、項目 10 に記載の方法。

(項目 14)

上記スピーチ信号の推定を上記オーディオ信号と合成することは、強度の上限および下限しきい値を確立し、かつ、上記オーディオ信号の一部分を上記上限の強度のしきい値と上記下限のしきい値との間の強度値を有する上記オーディオ信号の一部分に対応する上記スピーチ信号の推定の一部分と組み合わせることを包含する、項目 10 に記載の方法。

30

(項目 15)

上記オーディオ信号の上記一部分を上記スピーチ信号の推定の一部分と組み合わせることは、上記スピーチ信号の推定が、上記強度の下限しきい値に近い強度値を有する上記オーディオ信号の一部分に対する上記オーディオ信号より重みを置かれ、かつ、上記オーディオ信号が、上記強度の上限しきい値に近い強度値を有する上記オーディオ信号の一部分に対する上記スピーチ信号の推定より重みを置かれるように、上記オーディオ信号および上記スピーチ信号に重みを置くことを包含する、項目 14 に記載の方法。

40

(項目 16)

上記背景ノイズの推定を上記ニューラルネットワークに供給することをさらに包含する、項目 11 に記載の方法。

(項目 17)

以前の時間ステップからの上記スピーチ信号の推定を上記ニューラルネットワークに供給することをさらに包含する、項目 11 に記載の方法。

(項目 18)

以前の時間ステップからの上記スピーチ信号の推定の中間結果を上記ニューラルネットワークに供給することをさらに包含する、項目 11 に記載の方法。

(項目 19)

50

以前の時間ステップからの上記オーディオ信号を上記ニューラルネットワークに供給することをさらに包含する、項目 11 に記載の方法。

(項目 20)

スピーチ信号をエンハンスするシステムであって、
スピーチコンテンツおよび背景ノイズの両方を有する時系列のオーディオ信号を提供するオーディオ信号出力ソースと、
時系列領域から周波数領域に上記オーディオ信号を変換する周波数変換機能を提供する信号プロセッサと、
背景ノイズの推定器と、
ニューラルネットワークと、
信号コンバイナと
を備え、

上記背景の推定器は、上記オーディオ信号における上記背景ノイズの推定を形成し、上記ニューラルネットワークは、上記オーディオ信号から、上記スピーチ信号の推定を抽出し、上記信号コンバイナは、上記背景ノイズの推定に基づいて上記スピーチ信号の推定を上記オーディオ信号と組み合わせることにより、大幅に減少した背景ノイズを有する再構築されたスピーチ信号を生成する、システム。

(項目 21)

上記ニューラルネットワークは、第 1 のセットの入力ノードであって、上記オーディオ信号を受信する第 1 のセットの入力ノードを包含した、項目 20 に記載の方法。

(項目 22)

上記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから上記オーディオ信号を受信する第 2 のセットの入力ノードを包含した、項目 21 に記載の方法。

(項目 23)

上記ニューラルネットワークは、第 2 のセットの入力ノードであって、上記背景ノイズの推定を受信する第 2 のセットの入力ノードを包含した、項目 21 に記載の方法。

(項目 24)

上記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから上記スピーチ信号の推定を受信する第 2 のセットの入力ノードを包含した、項目 21 に記載の方法。

(項目 25)

上記ニューラルネットワークは、第 2 のセットの入力ノードであって、以前の時間ステップから中間結果を受信する第 2 のセットの入力ノードを包含した、項目 21 に記載の方法。

(項目 26)

背景ノイズからスピーチ信号を分離する方法であって、
オーディオ信号を受信することと、
信号の正確さが、高い確実性を有すると知られている上記オーディオ信号の一部を識別することと、
ニューラルネットワークを訓練することにより、上記オーディオ信号の正確さが不確かである上記オーディオ信号の一部に対して、著しく減少した背景ノイズを有する再構築された信号を推定することと
を包含する、方法。

(摘要)

スピーチ信号の周波数コンポーネントが、背景ノイズによってマスクされる環境において送信されるスピーチ信号を分離し、再構築するように構成されているスピーチ信号分離システム。スピーチ信号分離システムは、オーディオソースからノイジーなスピーチ信号を取得する。ノイジーなスピーチ信号は、それから、背景ノイズからクリーンなスピーチ信号を分離し、再構築するように訓練されたニューラルネットワークを介して供給される。

10

20

30

40

50

ノイジーなスピーチ信号が、ニューラルネットワークを介して供給されると、スピーチ信号分離システムは、大幅に減少したノイズを有する推定されたスピーチ信号を生成する。

【発明を実施するための最良の形態】

【0013】

本発明は、以下の図面および記載を参照して、より理解される。図中のコンポーネントは、縮尺に強調が置かれているのではなく、むしろ本発明の原理に強調が置かれている。さらに、図面において、同様の参照番号は、異なる見方の図面にわたって、対応するパーツを指し示す。

【0014】

本発明は、信号を背景ノイズから分離するためのシステムと方法に関するものである。そのシステムと方法は、特に、ノイズ環境の中で発せられたオーディオ信号からスピーチ信号を回復するのに効果的に適用される。しかしながら、この発明は、スピーチ信号のみに限られるものではなく、ノイズによって不明瞭となった任意の信号にも用いられ得る。

10

【0015】

図1は、スピーチ信号を背景ノイズから分離する方法100を説明している。方法100では、周波数成分が背景ノイズにマスクされているという環境において伝えられたスピーチ信号を再構築し分離することができる。以下の記述は、多くの具体的な詳細を説明することにより、スピーチ信号分離法100と、その方法を取り入れるための関連システム10について、より完全な説明を与えるものである。しかしながら、当業者にとっては、発明がこれらの具体的な詳細なしには実現されないということは明らかである。他の事例においては、本発明を不明瞭としないために、よく知られて特徴は詳述されない。背景ノイズからスピーチ信号を分離する方法10では、まずノイジーなスピーチ信号を受けとる（ステップ102）。第2のステップ104では、スピーチ信号を、ノイズを抑えたスピーチをノイズ入力信号から抽出するために採り入れられたニューラルネットワークを通して入力する。最後のステップ106は、スピーチ信号を推定することである。

20

【0016】

スピーチ信号分離システム10を図14に示す。スピーチ信号分離システムはマイクロフォン12のような、オーディオ信号装置やオーディオ信号を供給するために構成された任意の他のオーディオソースを含むこともある。A/Dコンバーター14は、マイクロフォン12から発せられたアナログのスピーチ信号をデジタル信号に変換し、そのデジタルスピーチ信号を信号処理ユニット16への入力として供給するためにある。オーディオ信号装置がデジタルオーディオ信号を供給する場合は、A/Dコンバーターは除外され得る。デジタル処理ユニット16は、デジタル処理ユニットや、コンピューター、あるいはオーディオ信号を供給することのできる他のタイプの回路やシステムであり得る。信号処理ユニットは、ニューラルネットワークコンポーネント18と、背景ノイズ評価コンポーネント20、信号ブレンド成分22を含んでいる。ノイズ評価コンポーネントは多数の周波サブバンドを通じて受け取られた信号のノイズレベルを測定するものである。ニューラルネットワークコンポーネント18は、オーディオ信号を受け取り、そのオーディオ信号のスピーチ成分を、オーディオ信号の背景ノイズコンポーネントから分離するために、構成されている。信号ブレンドコンポーネント22は、完全にノイズを取り除いたオーディオ信号を、分離されたスピーチコンポーネントとオーディオ信号のひとつの機能として再構築する。このように、オーディオ信号分離システム10はオーディオ信号を背景ノイズから分離し、背景ノイズをかなり抑制、あるいは除去した後、その背景ノイズが元の信号に存在していない場合、真のオーディオ信号がどのように見え、どのように響いたかの推定を与えることによって、完全なオーディオ信号を再構築するのである。

30

40

【0017】

図2は典型的な母音の周波スペクトラムを表したグラフであり、オーディオ信号がどのように特徴づけられるかの一例である。母音が特に興味深いのは、それらが概してオーディオ信号の最強度で構成されており、同様にオーディオ信号を妨害するノイズを超えるもっとも高い可能性を持つ。図2では母音について示しているが、オーディオ信号分離シ

50

テム10と方法100は入力された任意のタイプのオーディオ信号も処理する。

【0018】

母音、つまりオーディオ信号200はその構成周波数とそれぞれの周波数帯域の強さの両方によって特徴づけられる。オーディオ信号200が、周波(Hz)軸と強さ(dB)軸に座標で描かれている。周波数座標は一般に任意の数の不連続のbinあるいは帯域から成る。周波数バンク206は、256個の周波数バンク(256bins)がオーディオ信号200から取られたことを示している。信号帯域の数の選択は、当業者には方法論としてよく知られており、256周波数帯域の帯域長は図解のためだけに使われている、もちろん他の帯域長も同様であるけれども。おおむね水平な線208は、オーディオ信号200が獲得された環境における背景ノイズの強さを表している。オーディオ信号200はノイズ208を超える強度範囲において容易に見つけられる。しかしながら、スピーチ信号200はそのノイズレベル以下の強度レベルで背景ノイズから取り出されなければならない。さらに、ノイズレベル208の強度あるいはそれに近いノイズレベルでは、スピーチをノイズ208と区別することが難しくなる可能性がある。

10

【0019】

再度、図1と図14を見ると、ステップ102で、スピーチ信号は、スピーチ信号分離装置によってマイクロフォンなどといった外部装置から獲得され得る。通常の場合、スピーチ信号200は、背景ノイズ、たとえばコンサートでの群集のノイズ、あるいは自動車のノイズ、また他のノイズ源からのノイズを含み得る。図2の線208が示すように、背景ノイズがスピーチ信号200の一部にかぶっている。スピーチ信号200は線208上で1回から数回ピークに達するが、何回か分離線208以下に落ちるときは、背景ノイズのために、分析がより困難あるいは不可能になる。ブロック104においては、スピーチ信号200が、ノイズ環境におけるスピーチ信号の分離と再構築を教育されたニューラルネットワークを介したスピーチ信号分離システム10を通じて入力され得る。ステップ106においては、ニューラルネットワークによって背景ノイズから分離されたスピーチ信号200が、かなり抑制された、あるいは除外された背景ノイズで、推測されるスピーチ信号を発するために使われている。

20

【0020】

スピーチ検出の主な問題は、背景ノイズからスピーチ信号200を分離することである。ノイズ環境においては、スピーチ信号200の周波数成分の多くが、一部あるいは全体に、ノイズ周波数にマスクされ得る。この現象は明らかに図3に現れている。ノイズ302がスピーチ信号300を妨害しているので、スピーチ信号300は、304部分でノイズ302にマスクされていて、容易に検出可能であるのはノイズ302を超える306部分だけである。306領域が信号300の一部のみを含んでいるので、ノイズのせいでスピーチ信号300のいくらかが失われるか、ノイズにマスクされている。

30

【0021】

ここに参照されているように、ニューラルネットワークというのは、人間の脳の相互に連結するニューロン組織をモデルにしたコンピューター構造である。ニューラルネットワークはパターンを識別する脳の能力を模している。使用においては、ニューラルネットワークはネットワークに入力されたデータの基礎となる関連を抽出するのである。ニューラルネットワークは、子供や動物が仕事を教えられるように、これらの関連を認識するよう訓練される。ニューラルネットワークは、試行錯誤の方法論を通じて学ぶ。各レッスンの繰り返しにより、ニューラルネットワークの性能は進歩する。

40

【0022】

図4に、スピーチ信号分離システム10によって使われ得る典型的なニューラルネットワーク400を示す。ニューラルネットワーク400は3つの計算層から成る。入力層402は入力ニューロン404から成る。隠れ層406は、隠れニューロン408から成る。出力層410は、出力ニューロン412から成る。図のように、402、406、410それぞれの層にある404、408、412のニューロンそれぞれが、続いている層402、406、410にあるニューロン404、408、412のそれぞれと、完全に相

50

互関連しあっている。このように、入力ニューロン404の各々が、接続414によって隠れニューロン408の各々と接続される。さらに、隠れニューロン408のそれぞれが接続416によって出力ニューロン412のそれぞれと接続されている。414と416それぞれの接続が重量要因と関連している。

【0023】

それぞれのニューロンは、数値データの範囲内で活性化する。この範囲はたとえば0から1である。入力ニューロン404への入力、アプリケーションあるいは、ネットワーク環境設定によって決定される。隠れニューロン408への入力、接続414の負荷要因に入力ニューロン404を乗じたか、あるいはそれによって調整された状態である。出力ニューロン412への入力、入力ニューロン408に接続416の負荷要因を乗じるか、それによって調整された状態である。隠れ、あるいは出力ニューロン412のそれぞれの活性は、そのノードへの入力の合計に対し、スカッシング関数あるいはシグモイド関数を応用した結果であり得る。スカッシング関数は、入力合計を範囲内の値に限定する非線形の関数である。再度、その範囲は0から1である。

10

【0024】

ニューラルネットワークは、例（結果がわかっている）が示されているときに「学習する」。負荷要因は、出力を正しい結果に近づけるよう繰り返すことで調整されている。訓練の後、実際に、入力ニューロン404のそれぞれの状態は、アプリケーションあるいはネットワーク環境設定によって割り当てられている。入力ニューロン404の入力は負荷のかかった接続414を通じて、隠れニューロン408のそれぞれに広がる。隠れニューロン408の結果として生じる状態が、入力層402に呈せられるパターンへのネットワークのソリューションである。

20

【0025】

図5は、スピーチ信号分離システム10によって行われたスピーチ信号処理をさらに詳しく説明するブロック図である。ステップ500では、スピーチ信号は、マイクロフォンといった、外部のスピーチ信号装置から獲得される。そのスピーチ信号はおよそ46ミリ秒の時系列を例にとったものであるが、他の時系列でも同様に使うことができる。当業者は、スピーチ信号がいくつかの異なるタイプのソースから得られたものであるとの認識を持ち得る。たとえば、そのスピーチ信号は、だれかが背景ノイズを取り除くことによってきれいにしたいと思うオーディオ録音から獲得され得るし、うるさい自動車内で1つか

30

【0026】

ステップ502では、時間領域から周波数領域への変換が行われている。この変換は、高速フーリエ変換（FFT）であり得、またDFT、DCT、フィルターバンク、あるいは全周波数でのスピーチ信号の出力を推定する方法であり得る。FFTは加重したサイン、コサインの総計として波形を表現するテクニックである。FFTは一組の不連続データ値のフーリエ変換を計算するためのアルゴリズムである。任意の有限のデータポイント、たとえばスピーチ信号の定期的なサンプリングデータがある場合、FFTはそのデータを成分周波数によって表す。以下に述べるとおり、それはまた、時間領域信号を周波数データから再構築するという基本的に同一の逆の問題を解決する。

40

【0027】

さらに説明されているように、ステップ504ではスピーチ信号に含まれる背景ノイズが推定されている。背景ノイズは、任意の既知の手段によっても評価され得る。たとえば、沈黙の期間から、あるいはスピーチが検出されないところからも平均が計算される。その平均値は、ノイズを測定するためにそれぞれの周波数における信号の割合によって継続的に調整される。そこでは、ノイズに対する信号の割合が低い周波数において平均値が、より早く最新値にアップデートされる。あるいはニューラルネットワークそのものがノイズを測定するために使用され得る。

【0028】

ステップ502で発せられたスピーチ信号と504で行われたノイズ測定は、506の

50

ステップで圧縮される。1つの例として、「M e l 周波数尺度」アルゴリズムはスピーチ信号を圧縮するために使われ得る。スピーチは、高い周波数よりも低い周波数においてより大きな構造を持つ傾向がある。それで非線系圧縮は一樣に圧縮帯域全体に周波数情報を公平に配布する傾向にある。

【0029】

スピーチにおける情報は対数の形で減衰する。より高い周波数においては、「S」あるいは「T」のみが見出される。そのため、実に少ない情報で足りる。M e l 周波数尺度は、音声情報を保護するための圧縮を最適化する。より低周波数において直線的、より高周波数において対数的である。M e l 周波数尺度は次の方程式によって実際の周波数に関連し得る。

【0030】

$$m e l (f) = 2595 \log (1 + f / 700)$$

f はヘルツ (H z) で計測される。信号圧縮の結果として生じる値は、「M e l 周波数バンク」に蓄積される。M e l 周波数バンクは、中心周波数を等間隔におかれた M e l 値にセットすることによって作成される、フィルターバンクである。この圧縮の結果は、圧縮されたノイズ信号だけでなく音声信号の情報内容をも際立たせるスムーズな信号となる。

【0031】

M e l 尺度はピッチの心理音響的な比率尺度を表す。ログベース (l o g b a s e) 2 周波数尺度、あるいは B a r k 尺度や E R B (E q u i v a l e n t R e c t a n g l a r B a n d w i d t h) 尺度といった、他の圧縮尺度もまた使用され得る。後者の2つは、臨界帯域の心理音響的現象に基づく経験的尺度である。

【0032】

圧縮に先立ち、502からのスピーチ信号もまた、スムーズにされ得る。このスムージングは、圧縮信号のスムーズネス上での高いピッチの調波から生じる可変性の衝撃を抑制し得る。スムージングは L P C あるいはスペクトラム平均、あるいは補間を使うことによって実行される。

【0033】

ステップ508では、スピーチ信号は圧縮された信号を、信号処理ユニット16のニューラルネットワーク成分18への入力として割り当てることにより、背景ノイズから抽出される。抽出された信号は、背景ノイズのない状態での元のスピーチ信号の評価を表す。ステップ510では、ステップ508によって作成された抽出信号が、ステップ506で作成された圧縮信号と混合される。混合処理は、必要な時のみ抽出スピーチ評価に依存するものの、できるだけ元の圧縮スピーチ信号(ステップ506から)の多くを保持している。図3に戻ると、306のような元のスピーチ信号のいくつかの部分が明らかに背景ノイズ302のレベルを超えているものは容易に検出される。そのため、スピーチ信号のこういった部分は、できるだけ多くの元の信号の特性を保持するために混合信号において保持され得る。元の信号が完全に背景ノイズにマスクされている部分においては、もし抽出信号が背景ノイズ、あるいは元の信号の強さを超えない場合、ステップ508でニューラルネットワークによって抽出されたスピーチ信号評価に頼らざるを得ない。信号の強度が、背景ノイズと同じレベルかあるいはそれに近い領域では、できるだけ元の信号の評価に近づけるために、圧縮された元の信号とステップ508で抽出された信号が組み合わせられ得る。混合処理は、できるだけ元の自然のままのスピーチ信号の特性を多く残しつつ、背景ノイズをかなり取り除いた、圧縮再構築されたスピーチ信号となる。

【0034】

残りのブロックは、圧縮され、再構築されたスピーチ信号に実行され得るステップの概要を述べる。時間で再構築されたスピーチ信号に実行されるステップは、スピーチ信号が用いられる用途に依存して、変更し得る。例えば、再構築されたスピーチ信号は、自動スピーチ認識システムと互換性のある形状に直接的に変換され得る。ステップ520は、メル周波数ケプストラル係数 (M e l F r e q u e n c y C e p s t r a l C o e f f i c i e n t (M F C C)) 変換を示す。ステップ520の出力は、スピーチ認識シス

10

20

30

40

50

テムに直接的に入力され得る。もしくは、ステップ510において、生成された圧縮され、再構築されたスピーチ信号は、ステップ516で、圧縮され、再構築された信号に逆周波数領域 時系列変換を実行することによって、時系列すなわち可聴なスピーチ信号に直接的に変換され得る。このことは、著しく減少したもしくは完全に除かれた背景ノイズを有する時系列のスピーチ信号の結果になる。他の代替において、圧縮され、再構築されたスピーチ信号は、ステップ512で、解凍され得る。調波が、ステップ514で、信号に加えられ得、信号が、また合成され得る。この時、元の圧縮されていないスピーチ信号および合成信号が時系列のスピーチ信号に変換され得る。もしくは、信号は、追加的な合成なしで、調波が加えられた直後に、時系列の信号に変換され得る。

【0035】

第1の合成ステップ510からの出力、第2の合成ステップ522からの出力、もしくは、ステップ514で、追加的な調波が加えられた直後の出力であるスピーチ信号は、ステップ502で用いられる時間 領域変換の逆を用いて、ステップ516で、時間領域に変換され得る。

【0036】

図6は、図5において、ステップ506で表されるスピーチ信号圧縮処理の第1の段階を示す。スピーチ信号600は、構成周波数および各周波数帯域の強度の両方によって特徴付けられる。スピーチ信号600は、周波数(Hz)軸602および強度(dB)軸604に対してプロットされる。周波数プロットは、通常、任意的な数の離散帯域を含む。周波数バンク606は、256個の周波数帯域は、スピーチ信号600を含むことを示す。信号帯域の数の選択は、当業者によく知られる方法であり、256個の帯域長は、例示目的のためだけに用いられる。分離線608は、背景ノイズの強度を表す。

【0037】

スピーチ信号600は、多くの周波数スパイク610を含む。これらの周波数スパイク610は、スピーチ信号600内における調波によって引き起こされ得る。これら周波数スパイク610の存在が、リアルなスピーチ信号をマスクし、スピーチ分離処理を複雑にする。これらの周波数スパイク610は、平坦化処理によって除かれ得る。平坦化処理は、信号を、スピーチ信号おける調波間に補間することから成る。調波情報がわずかであるスピーチ信号600の領域において、補間アルゴリズムは、残りの信号上で、補間値を平均化する。補間信号612は、この平坦化処理の結果である。

【0038】

図7は、圧縮されたノイジーなスピーチ信号700を示す図である。圧縮されたスピーチ信号700は、Mel帯域軸702および強度(dB)軸704に対してプロットされる。圧縮されたノイズの推定706が、また示されている。信号圧縮の結果は、より少ない数の帯域によって表せられる信号である。この例において、帯域数は、20~36個の帯域であり得る。より低い周波数を表す帯域は、通常、圧縮されていない信号の4~5個の帯域を表す。中央値の周波数における帯域は、およそ20個の圧縮前の帯域を表す。より高い周波数でのそれらは、通常、およそ100個の圧縮前の帯域を表す。

【0039】

図7は、またステップ508の予想される結果を示す。圧縮されたノイジーなスピーチ信号700(実線)は、信号処理ユニット15のニューラルネットワークコンポーネント18に入力される(図14)。ニューラルネットワークからの出力は、圧縮されたスピーチ信号(点線)708である。信号708は、スピーチ信号上のノイズのすべての影響が、打ち消されるか、もしくは無効にされる、理想的なケースを表す。圧縮されたスピーチ信号708は、再構築されたスピーチ信号と言われる。

【0040】

図7は、またステップ510の合成処理に利用される強度のしきい値を示す。強度の上限しきい値710は、背景ノイズの強度より、大幅に大きい強度レベルを定義する。このしきい値より、大きい元のスピーチ信号のコンポーネントが、背景ノイズの除去なしに直ちに検出され得る。従って、強度の上限しきい値710より大きい強度レベルを有する元

10

20

30

40

50

のスピーチ信号の一部分に対して、合成処理は、元の信号だけ用いる。強度の下限しきい値 712 は、背景ノイズの平均強度よりほんのわずかに小さい強度レベルを定義する。強度の下限しきい値 712 より小さい強度レベルを有する元の信号のコンポーネントは、識別できない。背景ノイズと識別不能である。従って、強度の下限しきい値 712 より小さい強度レベルを有する元のスピーチ信号の一部分に対して、合成処理は、抽出された信号が、背景ノイズもしくは元の信号の強度を超えないという条件で、ステップ 508 から生成される再構築された信号だけを用いる。強度の下限しきい値 712 と強度の上限しきい値 710 との間の範囲である強度レベルを有する元のスピーチ信号の一部分に対して、元のスピーチ信号は、そのスピーチ信号の明瞭度および品質に寄与する情報を提供する点において依然貴重であるコンテンツを含む。しかし、元のスピーチ信号は、信頼性に欠ける。なぜなら、背景ノイズの平均値に近く、実際、ノイズのコンポーネントを含み得るからである。従って、強度の下限しきい値 712 と強度の上限しきい値 710 との間の範囲である強度レベルを有する元のスピーチ信号の一部分に対して、ステップ 510 での合成処理は、ステップ 508 から、圧縮された元のスピーチ信号と、圧縮され、再構築された信号両方のコンポーネントを用いる。強度の下限しきい値と強度の上限しきい値との間の範囲である強度レベルを有する再構築された信号の一部分に対して、ステップ 510 において、合成処理は、スライド制アプローチを用いる。強度の上限しきい値により近い元の信号から情報は、ノイズのしきい値からさらに遠くなり、強度の下限しきい値により近い元の信号から情報より信頼性がある。このことを説明するために、合成処理は、信号強度が、強度の下限しきい値 712 に近いとき、元のスピーチ信号により重みを置く。相互的な方法において、合成処理は、信号強度が、強度の下限しきい値 712 に近い強度レベルを有する強度レベルの一部分に対して、ステップ 508 からの、圧縮され、再構築されたスピーチ信号により重みを置き、かつ、強度の上限しきい値 710 に近づく強度レベルを有する元の信号一部分に対して、圧縮され、再構築されたスピーチ信号より少ない価値を置く。

【0041】

図 8 は、他の例示的スピーチ分離システムのニューラルネットワークを表す図である。ニューラルネットワーク 800 は、3つの処理層から成る。入力層 802、隠れ層 804 および出力層 806 である。入力層 802 は、入力ニューロン 808 を含み得る。隠れ層 804 は、隠れニューロン 810 を含み得る。出力層 806 は、出力ニューロン 812 を含み得る。入力層 802 における各入力ニューロン 808 は、1つ以上の接続 814 を介して、隠れ層 804 における各隠れニューロン 810 に完全に相互接続されている。隠れ層 804 における各隠れニューロン 810 は、1つ以上の接続 816 を介して、出力層 806 に各出力ニューロン 812 に完全に相互接続されている。

【0042】

詳細には示されていないが、入力層 802 における入力ニューロン 808 の数は、周波数バンク 702 における帯域の数に対応し得る。出力ニューロン 812 の数は、またに周波数バンク 702 における帯域の数と同等であり得る。隠れ層 804 における隠れニューロン 810 の数は、10個から 80個の間の数であり得る。入力ニューロン 808 の状態は、周波数バンク 702 における強度値によって決定される。実際には、ニューラルネットワーク 800 は、ノイジーなスピーチ信号 700 を、入力信号として取り、クリーンなスピーチ信号 708 を、出力として生成する。

【0043】

図 9 は、他の例示的なスピーチ分離システムもニューラルネットワーク 900 を表す図である。ニューラルネットワーク 900 は、3つの処理層を含む。入力層 902、隠れ層 904 および出力層 906 である。入力層 902 は、2つのセットの入力ニューロン、スピーチ信号の入力層 908 およびマスク入力層 910 を含み得る。スピーチ信号入力層 908 は、入力ニューロン 912 を含み得る。マスク入力層 910 は、入力ニューロン 914 を含み得る。隠れ層 904 は、隠れニューロン 916 を含み得る。出力層 906 は、出力ニューロン 918 を含み得る。スピーチ信号入力層 908 における各入力ニューロン 912

およびノイズ信号の入力層 910 における各入力ニューロン 914 は、1つ以上の接続 920 を介して、隠れ層 904 における各隠れニューロン 916 に完全に相互接続されている。隠れ層 904 における各隠れニューロン 916 は、1つ以上の接続 922 を介して、出力層 906 に各出力ニューロン 918 に完全に相互接続されている。

【0044】

スピーチ信号入力層 908 におけるニューロン 912 の数は、周波数バンク 702 における帯域の数に対応し得る。同様に、マスク信号の入力層 910 におけるニューロン 914 の数は、周波数バンク 702 における帯域の数に対応し得る。出力ニューロン 918 の数は、また周波数バンド 702 における帯域の数と同等であり得る。隠れ層 904 における隠れニューロン 916 の数は、10個から80個の間の数であり得る。入力ニューロン 912 および入力ニューロン 914 の状態は、周波数バンク 702 における強度値によって決定される。

10

【0045】

実際には、ニューラルネットワーク 900 は、入力としてノイジーなスピーチ信号 700 を取り、出力としてノイズが減少したスピーチ信号 708 を生成する。マスク入力層 910 は、506 からのスピーチ信号の品質についての情報を直接的に、もしくは間接的に、または 700 によって表される情報として、提供する。つまり、1つの例において、マスク入力層 910 は、入力して、圧縮されたノイズの推定 706 を取る。

【0046】

本発明の他の1つ例において、2進法のマスクが、ノイズの推定 706 と圧縮されたノイジーな信号 700 との比較から計算され得る。702 の各圧縮された周波数バンドで、マスクは、ノイジーな信号 700 とノイズの推定 706 との間の強度差異が、3dB といったしきい値を超えると、1にセットされ得、他のとき、0にセットされる。マスクは、スピーチを示す周波数帯域が信頼的もしくは有用的な情報を搬送するかどうかの指示を表す。506 の関数は、マスクによって0であると示される(つまり、ノイズの推定 706 によってマスクされる)ノイジーな信号 700 の一部分だけを再構築し得る。

20

【0047】

本発明の他の例において、マスクは、2進法ではなく、ノイジーな信号 700 とノイズの推定 706 との間の差異である。従って、この「ファジー」なマスクは、ニューラルネットワークに信頼性の自信度を示す。ノイジーな信号 700 がノイズの推定 706 に出会う領域は、2進法のマスクにおいてと同様に、0にセットされる。ノイジーな信号 700 がノイズの推定 706 に大変近い領域は、低い信頼性もしくは自信度を示すいくつかの小さい値を有し、またノイジーな信号 700 がノイズの推定 706 を大きく超える領域は、優れたスピーチ信号の品質を示す。

30

【0048】

ニューラルネットワークは、周波数に渡る関連性と同様に時間における関連性を学び得る。このことは、スピーチに対して重要であり得る。なぜなら、口、喉頭および声道の物理的なメカニズムは、どれだけ早く1つの音が他の音に続いて作成されるかに関して、制限を課すからである。従って、1つの時間枠から隣の時間枠への音は、相関している傾向があり、これらの相関を学び得るニューラルネットワークは、相関を学び得ないニューラルネットワークより、性能が優れている。

40

【0049】

図10は、他の例示的なスピーチ分離のニューラルネットワーク 1000 を表す図である。個々のニューロンは、簡略化のためにここに示されていない。ニューラルネットワーク 1000 は、3つの処理層を含む。入力層(1002~1008)、隠れ層 1010 および出力層 1012 である。ネットワーク 1000 は、入力層(1002~1006)におけるニューロンの起動値が、以前の時間ステップで、圧縮されたスピーチ信号から値を割り当てられ得ることを除いて、ニューラルネットワーク 900 と同一である。例えば、時間 t において、入力層 1002 は、 $t-2$ で、圧縮されたノイジーな信号 700 を割り当てられ、1004 は、 $t-4$ で、ノイジーな信号 700 に割り当てられ、時間 t で、入

50

力層 1006 は、ノイジーな信号 700 に割り当てられ、1008 は、上述のように、マスクを割り当てられ得る。従って、隠れ層 1010 は、圧縮されたスピーチ信号間の時間的な関連性を学び得る。

【0050】

図 11 は、他の例示的なスピーチ分離のニューラルネットワーク 1100 を表す図である。ニューラルネットワーク 1100 は、3つの処理層を含む。入力層 (1102 ~ 1106)、隠れ層 1108 および出力層 1110 である。ネットワーク 1100 は、入力層 1106 におけるニューロンの起動値が、以前の時間ステップで、出力層 1110 から抽出されたスピーチ信号から値を割り当てられ得ることを除いて、ニューラルネットワーク 900 と同一である。例えば、時間 t において、入力層 1102 は、 $t-1$ で、圧縮されたノイジーな信号 700 を割り当てられ、入力層 1104 は、マスクに割り当てられ、入力層 1106 は、時間 $t-1$ で、出力層 1110 の状態に割り当てられる。このネットワークは、ジョーダン (Jordan) ネットワークとして、学問においてよく知られ、かつ、現在の入力および依然の出力に依存して、その出力を変更することを学び得る。

10

【0051】

図 12 は、他の例示的なスピーチ分離のニューラルネットワーク 1200 を表す図である。ニューラルネットワーク 1200 は、3つの処理層を含む。入力層 (1202 ~ 1206)、隠れ層 1208 および出力層 1210 である。ニューラルネットワーク 1200 は、入力層 1206 におけるニューロンの起動値が、以前の時間ステップで、隠れ層 1208 から抽出されたスピーチ信号から値を割り当てられ得ることを除いて、ニューラルネットワーク 1100 と同一である。例えば、時間 t において、入力層 1202 は、 $t-1$ で、圧縮されたノイジーな信号 700 を割り当てられ、入力層 1204 は、マスクに割り当てられ、入力層 1206 は、時間 $t-1$ で、入力層 1206 の状態に割り当てられる。このネットワークは、エルマン (Elman) ネットワークとして、学問においてよく知られ、かつ、現在の入力および依然の内部的もしくは隠れ活動に依存して、その出力を変更することを学び得る。

20

【0052】

図 13 は、他の例示的なスピーチ分離のニューラルネットワーク 1300 を表す図である。ニューラルネットワーク 1300 は、そのニューラルネットワーク 1300 は、他の隠れユニット層 1310 を含むことを除いて、ニューラルネットワーク 1200 と同一である。この付加的な層は、スピーチをより良く抽出する、より高いオーダーの関連性の学習を可能にし得る。

30

【0053】

隠れもしくは出力ユニットの強度値は、その隠れもしくは出力ユニットが接続されている各入力ニューロンの強度とニューロン間の接続の重みの積の合計によって決定され得る。非線形関数は、隠れもしくは出力ニューロンの起動の範囲を減少させるために用いられる。この非線形関数は、S 字形関数、ロジスティック関数もしくは双曲線関数、または、絶対限度を有する線形のいずれかであり得る。これらの関数は、当業者にとってよく知られている。

【0054】

ニューロンネットワークは、リアルもしくはシュミレートされたノイズが加えられる複数参加型のクリーンなスピーチ信号に向けて訓練され得る。

40

【0055】

本発明のさまざまな実施形態が記載されてきたが、より多くの実施形態およびインプリメンテーションが本発明の範囲内で可能であることは当業者にとって明らかである。したがって、本発明は添付の請求項および均等物を含む。

【図面の簡単な説明】

【0056】

【図 1】スピーチ信号分離システムを示すブロック図である。

【図 2】典型的な母音の周波数スペクトラムを示す図である。

50

【図 3】ノイズによって部分的にマスクされる典型的な母音の周波数スペクトラムを示す図である。

【図 4】ニューラルネットワークの図である。

【図 5】スピーチ信号分離システムのスピーチ信号の処理方法を示すブロック図である。

【図 6】ノイズおよびその平坦化されたエンベロープによって部分的にマスクされる典型的な母音の例示である。

【図 7】圧縮されスピーチ信号を示す図である。

【図 8】スピーチ信号分離システムによって用いられる例示的なニューラルネットワークアーキテクチャの図である。

【図 9】本発明に従った他の例示的なニューラルネットワークアーキテクチャの図である 10

【図 10】他の例示的なニューラルネットワークアーキテクチャの図である。

【図 11】フィードバックを含む他の例示的なニューラルネットワークアーキテクチャの図である。

【図 12】フィードバックを含む他の例示的なニューラルネットワークアーキテクチャの図である。

【図 13】フィードバックおよび追加的な隠れ層を含む他の例示的なニューラルネットワークアーキテクチャの図である

【図 14】スピーチ信号分離システムのブロック図である。

【符号の説明】 20

【0057】

400、800、900、1000、1100、1200、1300 ニューラルネットワーク

404、808、912、914 入力ニューロン

406、804、904、1010、1108、1208 隠れ層

408、810、916 隠れニューロン

410、806、906、1012、1110、1210 出力層

412、812、918 出力ニューロン

802、902、1002、1004、1006、1008、1102、1104、1106、1202、1204、1206 入力層 30

814、816、920、922 接続

908 スピーチ信号入力層

910 マスク入力層

1310 隠れユニット層

【 図 1 】

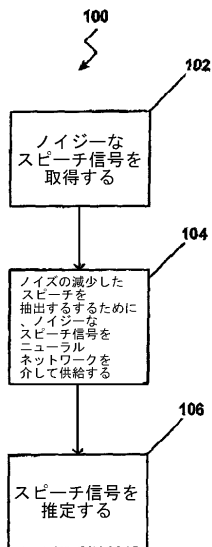


Figure 1

【 図 2 】

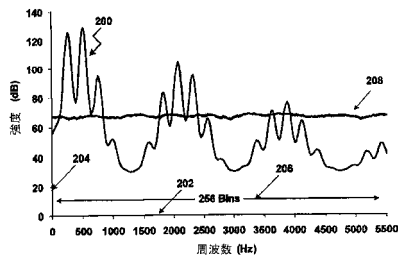


Figure 2

【 図 3 】

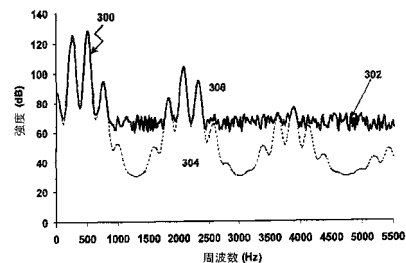


Figure 3

【 図 4 】

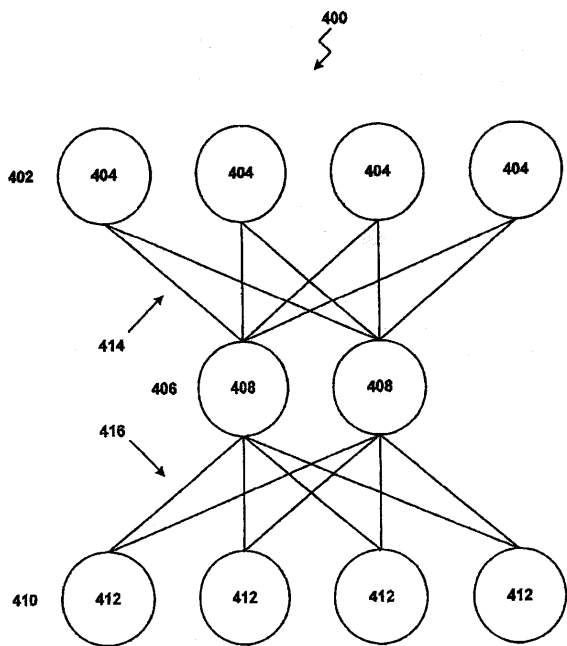


Figure 4

【 図 5 】

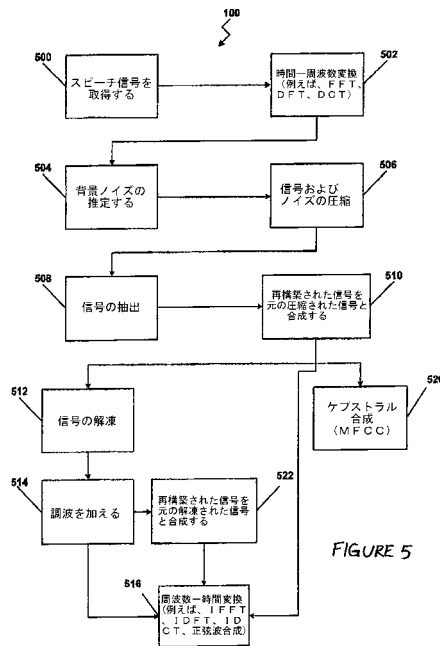


FIGURE 5

【 図 1 4 】

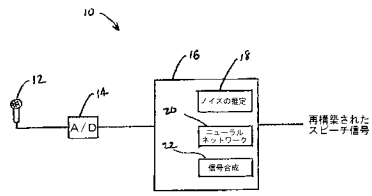


Figure 14

フロントページの続き

- (72)発明者 フィリップ ヘザーリントン
カナダ国 ブイ3エイチ 5エイチ7, ブリティッシュ コロンビア, ポート ムーディー,
ファーンウェイ ドライブ 23
- (72)発明者 ピアー ザカラウスカス
カナダ国 ブイ6エイチ 3アール4, ブリティッシュ コロンビア, バンクーバー, アイ
ロンワーク パッセージ 1015
- (72)発明者 シャーラ パービーン
カナダ国 ブイ5ダブリュー 3イー5, ブリティッシュ コロンビア, バンクーバー, ブ
リンズ アルバート ストリート 6163

Fターム(参考) 5D015 EE05