

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6619072号
(P6619072)

(45) 発行日 令和1年12月11日(2019. 12. 11)

(24) 登録日 令和1年11月22日(2019. 11. 22)

(51) Int. Cl.	F I	
G 1 O L 13/08 (2013.01)	G 1 O L 13/08	1 2 4
G 1 O L 13/00 (2006.01)	G 1 O L 13/00	1 0 0 V
G 1 O L 13/02 (2013.01)	G 1 O L 13/02	1 1 0 Z
H O 4 N 21/431 (2011.01)	H O 4 N 21/431	
H O 4 N 21/435 (2011.01)	H O 4 N 21/435	

請求項の数 7 (全 15 頁)

(21) 出願番号	特願2018-191616 (P2018-191616)	(73) 特許権者	000004226
(22) 出願日	平成30年10月10日 (2018.10.10)		日本電信電話株式会社
(62) 分割の表示	特願2015-80593 (P2015-80593) の分割		東京都千代田区大手町一丁目5番1号
原出願日	平成27年4月10日 (2015.4.10)	(74) 代理人	100121706
(65) 公開番号	特開2019-23747 (P2019-23747A)		弁理士 中尾 直樹
(43) 公開日	平成31年2月14日 (2019.2.14)	(74) 代理人	100128705
審査請求日	平成30年10月10日 (2018.10.10)		弁理士 中村 幸雄
		(74) 代理人	100147773
			弁理士 義村 宗洋
		(72) 発明者	鎌本 優
			東京都千代田区大手町一丁目5番1号 日 本電信電話株式会社内
		(72) 発明者	白木 善史
			東京都千代田区大手町一丁目5番1号 日 本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 音合成装置、音合成方法、及びそのプログラム

(57) 【特許請求の範囲】

【請求項1】

時系列の視覚情報に対し、当該視覚情報が表す意味に対応し、かつ、発話以外の音である音刺激を、当該時系列の視覚情報のタイミング及び数の少なくとも何れかに基づき、合成する音刺激合成部を含み、

上記時系列の視覚情報は上記音刺激とは異なる時系列の音信号と対応付けられており、上記音刺激合成部は上記音刺激を上記時系列の音信号に重畳して出力し、

上記音刺激合成部は、複数の上記視覚情報の中から、同じ意味を表す複数の視覚情報それぞれの入力時刻を抽出し、所定の時間区間毎に、抽出した入力時刻の平均値、最頻値、最小値及び最大値の少なくともいずれかである代表値に基づいて音刺激を重畳するタイミングを求める、

音合成装置。

【請求項2】

請求項1の音合成装置であって、

上記代表値は平均値であり、

上記音刺激合成部は、上記視覚情報の中から、同じ意味を表す複数の視覚情報それぞれの入力時刻を抽出し、抽出した入力時刻の、所定の時間区間毎の平均値と分散とを持つガウス分布に従う乱数に基づいて上記音刺激を重畳するタイミングとする、

音合成装置。

【請求項3】

請求項 2 の音合成装置であって、

音刺激が記憶される音刺激データベースを含み、

上記音刺激合成部は、上記視覚情報が映像表示部に表示されてから消えるまでの時間繰返して上記音刺激データベースに記憶された上記音刺激を重畳し、

フレームのインデックスを i とし、合成後の音刺激を Y_i とし、フレーム i に対して音刺激データベースから取り出した音刺激のテンプレートを T_i とし、 τ を音刺激のテンプレートの長さに対応する値とし、前記乱数を σ_i とし、音刺激を重畳するタイミングを表すインパルスを $(i \cdot \tau + \sigma_i)$ とし、 $*$ を畳み込み演算を表す演算子とし、上記音刺激合成部は、

【数 5】

$$Y_i = \delta(i \cdot \tau + \sigma_i) * T_i \quad (i = 0, 1, 2, \dots)$$

10

により音刺激を合成する、もしくは、

前記乱数を σ_m とし、映像表示部に表示される同じ意味を表す視覚情報の個数を M とし、上記音刺激合成部は、

上記音刺激が拍手の場合、音刺激を重畳するタイミングを表すインパルスを

【数 9】

$$\delta\left(\frac{1}{M} m \cdot \tau + \sigma_m\right)$$

とし、

【数 6】

$$Y_i = \sum_{m=1}^M \delta\left(\frac{1}{M} m \cdot \tau + \sigma_m\right) * T_i \quad (i = 0, 1, 2, \dots)$$

により、音刺激を合成し、

上記音刺激が手拍子の場合、音刺激を重畳するタイミングを表すインパルスを $(\tau + \sigma_m)$ とし、

【数 7】

$$Y_i = \sum_{m=1}^M \delta(\tau + \sigma_m) * T_i \quad (i = 0, 1, 2, \dots)$$

30

により、音刺激を合成し、

上記音刺激が笑い声の場合、音刺激を重畳するタイミングを表すインパルスを

【数 10】

$$\delta\left(\frac{1}{M} m \cdot \tau + \sigma_m\right)$$

とし、

【数 8】

$$Y_i = \sum_{m=1}^M \delta\left(\frac{1}{M} m \cdot \tau + \sigma_m\right) * T_i \quad (i = 0, 1, 2, \dots)$$

により、音刺激を合成する、

音合成装置。

【請求項 4】

請求項 1 から請求項 3 の何れかの音合成装置であって、

音刺激が記憶される音刺激データベースを含み、

50

上記音刺激合成部は、前記代表値を、一番初めに音刺激を音信号に重畳するタイミングとして利用し、上記視覚情報が映像表示部に表示されてから消えるまでの時間繰り返して上記音刺激を重畳する、
音合成装置。

【請求項 5】

請求項 1 から請求項 4 の何れかの音合成装置であって、
視覚情報と、個数と、視覚情報が表す意味に対応する、個数に応じた音量の音刺激とが対応付けて記憶される音刺激データベースを含み、
上記音刺激合成部は、上記時系列の視覚情報の中から抽出される同じ意味を表す複数の視覚情報とその視覚情報が映像表示部に表示される個数とに対応する音刺激を上記音刺激データベースから選択し、重畳する、
音合成装置。

【請求項 6】

音合成装置による音合成方法であって、
音刺激合成部が、時系列の視覚情報に対し、当該視覚情報が表す意味に対応し、かつ、発話以外の音である音刺激を、当該時系列の視覚情報のタイミング及び数の少なくとも何れかに基づき、合成する音刺激合成ステップを含み、
上記時系列の視覚情報は上記音刺激とは異なる時系列の音信号と対応付けられており、上記音刺激合成ステップは上記音刺激を上記時系列の音信号に重畳して出力し、
上記音刺激合成ステップは、複数の上記視覚情報の中から、同じ意味を表す複数の視覚情報それぞれの入力時刻を抽出し、所定の時間区間毎に、抽出した入力時刻の平均値、最頻値、最小値及び最大値の少なくともいずれかである代表値に基づいて音刺激を重畳するタイミングを求める、
音合成方法。

【請求項 7】

請求項 1 から請求項 5 の何れかの音合成装置として、コンピュータを機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、映像を見るものによって入力されるテキスト情報を、その映像に重畳して表示する技術に関する。

【背景技術】

【0002】

映像を見るものによって入力されるテキスト情報を、その映像に重畳して表示する従来技術として非特許文献 1 が知られている。非特許文献 1 では、視聴者は、動画を視聴しながら、コメントを投稿することができる。

【0003】

また、伝送元において収録された拍手や手拍子音、声援・掛け声などの環境音を効率よく伝送し、伝送先で伝送元の場の雰囲気再現する従来技術として特許文献 1 が知られている。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2014 - 63145 号公報

【非特許文献】

【0005】

【非特許文献 1】「動画の視聴 コメントの投稿」、[online]、NIWANGO.INC、[平成 27 年 2 月 2 日検索]、インターネット<URL : <http://info.nicovideo.jp/help/player/howto/>>

【発明の概要】

10

20

30

40

50

【発明が解決しようとする課題】

【0006】

しかしながら、従来技術では、テキスト情報が重畳された映像の雰囲気再現することは難しい。

【0007】

本発明は、テキスト情報が重畳された映像の雰囲気再現する音合成装置、音合成方法及びそのプログラムを提供することを目的とする。

【課題を解決するための手段】

【0008】

上記の課題を解決するために、本発明の一態様によれば、音合成装置は、時系列の視覚情報に対し、当該視覚情報が表す意味に対応し、かつ、発話以外の音である音刺激を、当該時系列の視覚情報のタイミング及び数の少なくとも何れかに基づき、合成する音刺激合成部を含み、時系列の視覚情報は音刺激とは異なる時系列の音信号と対応付けられており、音刺激合成部は音刺激を時系列の音信号に重畳して出力し、音刺激合成部は、複数の視覚情報の中から、同じ意味を表す複数の視覚情報それぞれの入力時刻を抽出し、所定の時間区間毎に、抽出した入力時刻の平均値、最頻値、最小値及び最大値の少なくともいずれかである代表値に基づいて音刺激を重畳するタイミングを求める。

【0009】

上記の課題を解決するために、本発明の他の態様によれば、音合成方法は、音刺激合成部が、時系列の視覚情報に対し、当該視覚情報が表す意味に対応し、かつ、発話以外の音である音刺激を、当該時系列の視覚情報のタイミング及び数の少なくとも何れかに基づき、合成する音刺激合成ステップを含み、時系列の視覚情報は音刺激とは異なる時系列の音信号と対応付けられており、音刺激合成ステップは音刺激を時系列の音信号に重畳して出力し、音刺激合成ステップは、複数の視覚情報の中から、同じ意味を表す複数の視覚情報それぞれの入力時刻を抽出し、所定の時間区間毎に、抽出した入力時刻の平均値、最頻値、最小値及び最大値の少なくともいずれかである代表値に基づいて音刺激を重畳するタイミングを求める。

【発明の効果】

【0010】

本発明によれば、テキスト情報が重畳された映像の雰囲気再現することができるという効果を奏する。

【図面の簡単な説明】

【0011】

【図1】第一実施形態に係る音合成装置の機能ブロック図。

【図2】第一実施形態に係る音合成装置の処理フローの例を示す図。

【図3】テキスト情報付き映像信号の例を示す図。

【図4】合成した音刺激を重畳した音信号の例を示す図。

【図5】音刺激データベースのデータ例を示す図。

【図6】テキスト情報が表示されるタイミングで、音刺激を音信号に重畳する例を示す図。

。

【図7】テキスト情報の個数が多いタイミングに合わせて、音刺激を音信号に重畳する例を示す図。

【図8】音刺激合成部の音刺激合成手順を例示する図。

【図9】テキスト情報の個数に応じて、音刺激の音量を変更し、変更後の音刺激を音信号に重畳する例を示す図。

【図10】音刺激データベースのデータ例を示す図。

【発明を実施するための形態】

【0012】

以下、本発明の実施形態について、説明する。なお、以下の説明に用いる図面では、同じ機能を持つ構成部や同じ処理を行うステップには同一の符号を記し、重複説明を省略す

る。

【 0 0 1 3 】

< 第一実施形態 >

図 1 は第一実施形態に係る音合成装置 1 0 0 の機能ブロック図を、図 2 はその処理フローを示す。

【 0 0 1 4 】

音合成装置 1 0 0 は、例えば、動画の視聴者によって操作される視聴者端末内に組み込まれる。なお、本実施形態において、動画とは、時間軸に同期させた音信号と共に提供される映像信号を意味する。視聴者端末は、入力部（キーボード、マウス、タッチパネル等）と、映像表示部（ディスプレイ、タッチパネル等）と音再生部（スピーカ等）を含み、例えば、パーソナルコンピュータ、スマートホン、タブレット等からなる。視聴者端末及び動画配信サーバ 9 2 は、通信回線を介して通信可能とされている。視聴者は、視聴者端末の入力部を介して、動画配信サーバ 9 2 に対して動画の再生を要求する（S 1）ことができる。

【 0 0 1 5 】

< 動画配信サーバ 9 2 >

動画配信サーバ 9 2 は、動画データベース及びビデオカメラから動画を受け取り、視聴者端末の要求に応じて、動画データベース内に格納されている動画、または、ビデオカメラで収録した動画をリアルタイムで配信する（S 2）。また、ビデオカメラで収録された動画に限らず、リアルタイムで合成・編集された CG やモーションキャプチャ等から合成された CG をリアルタイム配信することもある。動画データベースには、動画と共に動画に付加されたテキスト情報が記憶され、動画と共に配信される。さらに、テキスト情報にはメタデータが付加されている。メタデータとしては、テキスト情報の入力時刻、テキスト情報の大きさ、その色、その出現方法、その移動速度や、移動位置等がある。例えば、大きさ、色、出現方法、移動速度、移動位置等は、テキスト情報の入力者が選択できるものとしてもよく、視聴者端末がメタデータとしてテキスト情報と一緒に送信し、動画データベースに動画と共に記憶される。なお、テキスト情報付き動画に含まれるテキスト情報、音信号及び映像信号は、時間軸において同期しており、それぞれ時系列において対応付けられている。

【 0 0 1 6 】

< 音合成装置 1 0 0 >

音合成装置 1 0 0 には、テキスト情報付きの動画、より詳しく言うと、時間軸において同期しているテキスト情報、音信号及び映像信号が入力され、テキスト情報付き映像信号（図 3 参照）に合わせて、テキスト情報に対応する音刺激を合成し（S 1 2 0）、合成した音刺激を重畳した音信号（図 4 参照）を出力する。時間軸において同期しているテキスト情報、映像信号、音刺激及び音信号を併せて音刺激及びテキスト情報付き動画ともいう。

【 0 0 1 7 】

視聴者端末の映像表示部及び音再生部は、音刺激及びコメント情報付きの動画を再生し（S 4）、視聴者は、視聴者端末の映像表示部及び音再生部を介して、音刺激及びコメント情報付きの動画を視聴することができる。

【 0 0 1 8 】

例えば、非特許文献 1 のニコニコ動画（登録商標）では視聴者側から拍手や手拍子を表現するために「 8 」という文字を入力して手を打ったことを表現するテキスト情報が使われている。また、笑いを表すために「 w 」という文字が使われている。仮に、これらのテキスト情報に対して、従来の音声合成技術を適用した場合、これらのテキスト情報に対してそれぞれ、「ハチ」「ダブリュ」という音声合成される。一方、本実施形態では、「 8 」は拍手音及び手拍子音の何れか、「 w 」は笑い声というように変換し合成音を出力する。

【 0 0 1 9 】

10

20

30

40

50

音合成装置 100 は、音刺激データベース 110 と、音刺激合成部 120 とを含む。

【0020】

< 音刺激データベース 110 >

音刺激データベース 110 には、テキスト情報と、そのテキスト情報が表す意味に対応する音刺激のテンプレートとが対応付けられて記憶されている（図 5 参照）。なお、本実施形態ではテキスト情報は所定の行為を意味するものとする。また、音刺激は、発話以外の音である。ここで、「発話以外の音」とは、「言語を音声として発し、その結果として発せられた音声」以外の音を意味し、例えば、笑い声、拍手音である。例えば、笑いを意味するテキスト情報「w」に対して、「ダブリュ」という音声波形のテンプレートではなく、「笑い声」の音の波形のテンプレート（笑い声の場合、例えば数秒分のテンプレート）が記憶されている。また、拍手及び手拍子の何れかを意味するテキスト情報「8」に対して、「ハチ」という音声波形のテンプレートではなく、「拍手音及び手拍子音の何れか」の音の波形のテンプレート（拍手音及び手拍子音の場合、例えば、数百ミリ秒分のテンプレート）が記憶されている。なお、図 5 の例では、テキスト情報と音刺激のテンプレートとが 1 対 1 で対応しているが、1 対多、多対 1、多対多で対応してもよい。つまり、(1) 所定の意味（例えば笑い）を表す 1 つのテキスト情報（例えば「w」）とその意味に対応する複数の音刺激のテンプレート（複数の笑い声のバリエーションを用意する）とが 1 対多で対応してもよいし、(2) 所定の意味（例えば笑い）を表す複数のテキスト情報（例えば「w」「（笑）」「(^o^）」とその意味に対応する 1 つの音刺激のテンプレート（1 つの笑い声を用意する）とが多対 1 で対応してもよいし、(3) 所定の意味（例えば笑い）を表す複数のテキスト情報（例えば「w」「（笑）」「(^o^）」とその意味に対応する複数の音刺激のテンプレート（複数の笑い声のバリエーションを用意する）とが多対多で対応してもよい。なお、拍手音、手拍子音及び笑い声等は、その時々により、異なるほうがより自然に聞こえるため、テキスト情報が重畳された映像の雰囲気をもより自然に再現しようとするならば、複数の音刺激のテンプレートを用意するほうがよい。

【0021】

< 音刺激合成部 120 >

音刺激合成部 120 は、テキスト情報付き動画（テキスト情報+音信号+映像信号）を受け取り、時系列のテキスト情報に対し、テキスト情報が表す意味に対応する音刺激のテンプレートを音刺激データベース 110 から取り出し、時系列のテキスト情報のタイミングと数の少なくとも何れかに基づき、音刺激を合成し（S 120）、時系列の音信号に対応付けて出力する。なお、あるテキスト情報が表す意味に対応する音刺激のテンプレートが複数存在する場合には、その中から 1 つをランダムに選択すればよい。音刺激合成部 120 は、選択した音刺激のテンプレートを、必要に応じて前のフレームと補間をして、所定の時間長のフレーム単位（例えば映像の 1 フレームに対応する時間長）で、1 フレームごとに励起される音刺激を合成する。音刺激合成部 120 は、合成した音刺激を受け取った音信号に重畳して（時系列の音信号に対応付けて）、出力する。

【0022】

例えば、拍手及び手拍子の何れかを意味するテキスト情報「8」に対し、音刺激データベース 110 から対応する拍手音及び手拍子音の何れかの音の波形のテンプレート（例えば数百ミリ秒分のテンプレート）を取り出し、必要に応じて前のフレームと補間をして、1 フレームごとに励起される拍手音及び手拍子音の何れかの音刺激を合成する。そして、所定の時間分の拍手音及び手拍子音の何れかの音刺激を合成し、音信号に重畳する。同様の方法により、笑いを意味するテキスト情報「w」に対し、所定の時間分の笑い声の音刺激を合成し、音信号に重畳してもよい。

【0023】

なお、発話を意味するテキスト情報に合わせて従来の音声合成装置を用いて音声を合成してもよい。この場合、テキスト情報に対して、まず、本実施形態の音刺激合成を行い、音刺激合成の対象とならないテキスト情報に対して従来の音声合成装置を用いて音声を合成すればよい。例えば「素晴らしい 88」というテキスト情報が入力された場合、「すば

らしいハチハチ」という音声を合成するのではなく、「すばらしい(音声)+拍手音(音刺激)」という音を合成し、音信号に重畳する。このような構成とすることで、音信号に対して、従来の音声合成により合成された音声(話し声)と共に、音刺激(拍手音、手拍子音及び笑い声等)が重畳され、テキスト情報が重畳された映像の雰囲気をもより自然に再現することができる。

【0024】

なお、上述の通り、テキスト情報「8」は拍手及び手拍子の何れかを意味する。何れも手を叩く行為であるが、「手拍子」は一定のテンポに合わせて手を叩く行為であり、「拍手」は一定のテンポを持たずに手を叩く行為である。ここで、「手拍子」と「拍手」とは、手を叩く時間的間隔や音量的差異が異なるため(参考文献1)、例えば、音信号等に基づいて、テキスト情報が何れの行為を意味するのかを判別することができる。

(参考文献1) 鎌本優, 河原一彦, 尾本章, 守谷健弘, 「音楽鑑賞時に励起される拍手音・手拍子音の低遅延伝送に向けた基礎的検討」、日本音響学会 2014年秋季研究発表会, 1 Q 17, 2014年.

【0025】

例えば、音信号が曲を表し、一定のテンポがある場合には、テキスト情報「8」は手拍子を意味する可能性が高い。また、曲が終了後のテキスト情報「8」は拍手を意味する可能性が高い。また、テキスト情報「8」が一度の入力で連続している場合、つまりテキスト情報「88」が入力された場合には、「パチパチ」を意味し、拍手を意味する可能性が高い。また、テキスト情報「8」が周期的に入力される場合には、手拍子を意味する可能性が高い。

【0026】

(音刺激を重畳するタイミング)

(1)テキスト情報が表示されるタイミング(例えばテキスト情報の入力時刻)で、音刺激を音信号に重畳する(図6参照)。

【0027】

(2)テキスト情報の個数が多いタイミングに合わせて、音刺激を音信号に重畳する(図7参照)。

【0028】

複数のテキスト情報の中から、同じ意味を表すテキスト情報を抽出する。例えば、音刺激データベース110を参照して、音刺激合成部120は、同じ意味を表すテキスト情報毎に分類し、テキスト情報の入力時刻を抽出する。

【0029】

音刺激合成部120は、抽出した入力時刻の統計量に基づいて音刺激を音信号に重畳するタイミングを求める。例えば、抽出した入力時刻を用いて、所定の時間区間毎に、時間区間毎の代表値(平均値、最頻値、最小値及び最大値等の複数の入力時刻を代表する何らかの値)を求め、重畳するタイミングとして検出する。例えば、抽出した入力時刻を用いて、ヒストグラムを作成し、多数決により重畳するタイミングを求める。つまり、最頻値を重畳するタイミングとする。

【0030】

例えば、テキスト情報が手拍子を意味し、動画がミュージックビデオであり、曲のテンポが148BPM(Beats per Minutes)の場合、一拍の間隔は405ms程度なので、所定の時間区間を405msとする。また、例えば、テキスト情報が「拍手」または「笑い」を意味する場合、所定の時間区間を一連の行為「拍手」または「笑い」が、継続しうる最大の時間に設定する。例えば、何らかの事象に対して、「拍手」を送るのは、長くとも30秒程度であろうと想定される場合、最初に「拍手」を意味するテキスト情報が表示されてから1分以内に表示される「拍手」を意味する他のテキスト情報から代表値を求め、重畳するタイミングとして検出する。

【0031】

なお、この方法を用いる場合、所定の時間区間分のテキスト情報付き動画(テキスト情

報+音信号+映像信号)をバッファリングしておき、音刺激を重畳して、音刺激及びテキスト情報付き動画を出力すればよい。

【0032】

(繰り返し重畳する場合)

なお、一人の人間による一拍分の音刺激(拍手音、手拍子音、笑い声等)を音刺激データベース110に記憶しておき、音刺激を音信号に繰り返し重畳する構成としてもよい。その場合、(音刺激を重畳するタイミング)の(1)及び(2)で求めたタイミングを、一番初めに音刺激を重畳するタイミングとして利用する。所定の時間区間分(例えば、テキスト情報が映像表示部に表示されてから消えるまで)繰り返し重畳すればよい。

【0033】

例えば、同じ意味を表すテキスト情報が映像表示部に表示される個数が1個の場合は、図8Aのように、所定の間隔毎(例えば、拍手の場合約300msごと)に、音刺激を重畳する。なお、音刺激を重畳するタイミングに揺らぎを持たせてもよい。例えば、音刺激が拍手の場合、所定の間隔は約300msでよいが、より好ましくは300msを中心として時間間隔に揺らぎを持たせる。時間間隔に揺らぎを持たせることによってさらに自然な拍手音を合成することができる(参考文献1参照)。たとえば300msを中心としてガウス分布にしたがう乱数により、±数10msの揺らぎを持たせればよい。例えば、フレームのインデックスを*i*とし、合成後の音刺激(拍手音)を Y_i とし、フレーム*i*に対して音刺激データベース110から取り出した音刺激のテンプレートを T_i とし、テンプレート T_i の長さ(テンプレート T_i に含まれる全フレームに含まれる、音刺激のデータのサンプル数)を*P*とし、音刺激合成部120は、テンプレート $T_i = (t_i[1] \ t_i[2] \ \dots \ t_i[P])$ と拍手タイミングを表すインパルス($\delta(i \cdot \tau + \sigma_i)$)の畳み込み演算で Y_i を出力とする。テンプレートの長さは所定の間隔(拍手の場合では約300ms程度)よりも短いほうが、音が重ならないため好ましい。

【0034】

【数1】

$$Y_i = \delta(i \cdot \tau + \sigma_i) * T_i \quad (i = 0, 1, 2, \dots)$$

【0035】

ここで*は畳み込み演算を表す。ここで、 $\tau = 300 \text{ ms}$ であり、 σ_i は -10 ms から $+10 \text{ ms}$ の範囲で生成した乱数である。音刺激を重畳するタイミングはフレーム間隔で特定し、1フレームごとに励起される音刺激を合成し、その結果として、音の波形のテンプレート分(例えば、笑い声の場合、数秒分、拍手音及び手拍子音の場合、数百ミリ秒分)の音刺激を合成し重畳する。1人分の手拍子を合成する場合、音刺激を重畳する間隔は、手拍子の対象に応じて変化し、曲のテンポが148BPM(Beats Per Minute)の場合、405ms前後とする。さらに、時間間隔の揺らぎは、拍手の場合よりも手拍子の場合のほうが小さく設定したほうがよく、例えば、手拍子の場合の σ_i の範囲が拍手の場合の σ_i の範囲よりも小さくなるように設定する。

【0036】

同じ意味を表すテキスト情報が映像表示部に表示される個数に応じて、音刺激を重畳する時間間隔を変更してもよい。例えば、拍手を表す*M*個のテキスト情報が映像表示部に表示されている場合、図8Bのように、時間間隔を約 $300 / M \text{ (ms)}$ ごとに音刺激を重畳する。個数*M*の逆数を使って、時間間隔を約 $300 / M \text{ (ms)}$ と設定することで、拍手を表すテキスト情報の個数*M*が増えるに従って時間間隔が小さくなるように設定することができる。この場合もガウス分布やラプラス分布に従う乱数によって、揺らぎを持たせることができる。例えば音刺激合成部120は、

【0037】

【数 2】

$$Y_i = \sum_{m=1}^M \delta\left(\frac{1}{M}m \cdot \tau + \sigma_m\right) * T_i \quad (i = 0, 1, 2, \dots)$$

【0038】

によりテンプレートを変換した音刺激 Y_i ($i = 0, 1, 2, \dots$) を重畳する。テキスト情報の個数 M の手拍子を合成する場合、 M 個の手拍子のタイミングはほぼ同じなので、例えば、

【0039】

【数 3】

$$Y_i = \sum_{m=1}^M \delta(\tau + \sigma_m) * T_i \quad (i = 0, 1, 2, \dots)$$

【0040】

によりテンプレートを変換した音刺激 Y_i ($i = 0, 1, 2, \dots$) を重畳する。なお、この場合も時間間隔の揺らぎは、拍手の場合よりも手拍子の場合のほうが小さく設定したほうがよい(参考文献 1 参照)。

【0041】

音刺激の例として拍手音及び手拍子音の何れかを対象として説明したが、これに限らず拍手音及び手拍子音以外の音刺激(たとえば、一人の人間による笑い声)を対象としても良い。

【0042】

なお、音刺激合成部 120 において、音刺激のテンプレート $T_i = (t_i[1] \quad t_i[2] \quad \dots \quad t_i[P])$ と笑い声を表すインパルス $(m \cdot \tau + \sigma_m)$ の畳み込み演算で Y_i を出力としても良い。

【0043】

【数 4】

$$Y_i = \sum_{m=1}^M \delta\left(\frac{1}{M}m \cdot \tau + \sigma_m\right) * T_i \quad (i = 0, 1, 2, \dots)$$

【0044】

この場合、 τ は笑い声のテンプレートの長さ(数秒)に対応する値とする。

【0045】

(揺らぎのバリエーション)

揺らぎを持たせる際のバリエーションについて説明する。

【0046】

例えば、抽出した入力時刻を用いて、所定の時間区間毎に、入力時刻の平均値と分散とを求め、その平均値と分散とを持つガウス分布に従う乱数を重畳するタイミングとしてもよい。この方法により、音刺激を重畳するタイミングにゆらぎを与えることができ、より自然なタイミングで音刺激を再生することができる。

【0047】

テキスト情報が手拍子を意味する場合、予め手拍子を行う際に一般的に生じる分散の値を求めておき、その分散に基づき、重畳するタイミングを求めてもよい。例えば、上述の方法で重畳するタイミングを求め、その重畳するタイミングを中心として、手拍子を行う際に一般的に生じる分散を持つガウス分布に従う乱数を新たな(最終的に用いる)重畳するタイミングとする。

【0048】

10

20

40

50

(音刺激の音量を調整する方法)

音刺激の音量を調整する方法を説明する。

【0049】

(1)テキスト情報の個数に応じて、音刺激の音量を変更し、変更後の音刺激を音信号に重畳する(図9参照)。例えば、映像表示部に表示されるテキスト情報の個数が多くなるほど音量が大きくなるように変更する。

【0050】

(2)テキスト情報の大きさに応じて、音刺激の音量を変更し、変更後の音刺激を音信号に重畳する。例えば、映像表示部に表示されるテキスト情報の大きさが大きくなるほど音量が大きくなるように変更する。

【0051】

(3)音刺激データベース110に個数に応じて異なる音量の音刺激を收音し記憶しておく、音刺激合成部120は、テキスト情報とその個数に応じて音刺激を選択してもよい。この場合、音刺激データベース110には、テキスト情報と、個数と、そのテキスト情報が表す意味に対応する音刺激とが対応付けられて記憶されている。例えば、テキスト情報「w」と、個数1と、1人分の笑い声からなる音刺激とが対応付けられて記憶されており、テキスト情報「w」と、個数2と、2人分の笑い声からなる音刺激とが対応付けられて記憶されている。なお、テキスト情報と個数と音刺激とが1対1対1で対応してもよいし、1対1対多、多対1対1、多対1対多で対応してもよい。

【0052】

<効果>

以上の構成により、テキスト情報が重畳された映像の雰囲気再現することができる。

【0053】

<変形例>

本実施形態では、音合成装置100は、音刺激及びテキスト情報付き動画を出力しているが、本実施形態のポイントは、音刺激を合成することであり、少なくとも音刺激を出力すればよい。例えば、本実施例のように視聴者端末(パーソナルコンピュータ、スマートフォン、タブレット等)内に本実施形態の音合成装置100が組み込まれてもよいし、動画配信サーバ内に音合成装置100が組み込まれてもよい。また、音合成装置100を独立した装置として構成してもよい。少なくとも時間軸において音信号または映像信号に同期しているテキスト情報を受け取り、音刺激を合成し、出力することができればよい。音刺激を音信号に同期させる処理等は別装置において行ってもよい。

【0054】

本実施形態では、視聴者によって入力され、映像信号に重畳して表示される情報としてテキスト情報の例を示したが、他の視覚情報であってもよい。ここで、「視覚情報」とは、映像表示部を介して視覚的に認識可能な情報であって、例えば、文字、図形若しくは記号若しくはこれらの結合又はこれらと色彩との結合である。また、静止画に限らず、動く画像であってもよい。例えば、(1)本実施形態のように、「笑い」や「拍手」等の所定の行為を意味するテキスト情報(例えば「w」や「8」等)、(2)テキスト情報以外の「笑い」や「拍手」等の所定の行為を意味し、識別するためのコンピュータ上のビット情報、(3)顔文字、絵文字等、通常のテキスト情報で無いもの。例えば、キャリアの異なる携帯電話間で共通絵文字(参考文献2参照)、(4)アスキーアート等、全体としてはテキスト情報とテキスト情報の配置情報を用いた絵のようになっているもの(参考文献3参照)

(参考文献2)「docomo/au共通絵文字」、株式会社NTTドコモ、[online]、[平成27年2月9日検索]、インターネット<URL: https://www.nttdocomo.co.jp/service/developer/smart_phone/make_contents/pictograph/>

(参考文献3)「アスキーアート」、[online]、2015年2月2日、ウィキペディア、[平成27年2月9日検索]、インターネット<URL: <http://ja.wikipedia.org/wiki/%E3%82%A2%E3%82%B9%E3%82%AD%E3%83%BC%E3%82%A2%E3%83%BC%E3%83%88>>

このテキスト情報以外の視覚情報を含む場合の音刺激データベース110に格納されるデ

10

20

30

40

50

一タの例を図10に示す。

【0055】

<その他の変形例>

本発明は上記の実施形態及び変形例に限定されるものではない。例えば、上述の各種の処理は、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。その他、本発明の趣旨を逸脱しない範囲で適宜変更が可能である。

【0056】

<プログラム及び記録媒体>

また、上記の実施形態及び変形例で説明した各装置における各種の処理機能をコンピュータによって実現してもよい。その場合、各装置が有すべき機能の処理内容はプログラムによって記述される。そして、このプログラムをコンピュータで実行することにより、上記各装置における各種の処理機能がコンピュータ上で実現される。

【0057】

この処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。コンピュータで読み取り可能な記録媒体としては、例えば、磁気記録装置、光ディスク、光磁気記録媒体、半導体メモリ等のようなものでもよい。

【0058】

また、このプログラムの流通は、例えば、そのプログラムを記録したDVD、CD-ROM等の可搬型記録媒体を販売、譲渡、貸与等することによって行う。さらに、このプログラムをサーバコンピュータの記憶装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピュータにそのプログラムを転送することにより、このプログラムを流通させてもよい。

【0059】

このようなプログラムを実行するコンピュータは、例えば、まず、可搬型記録媒体に記録されたプログラムもしくはサーバコンピュータから転送されたプログラムを、一旦、自己の記憶部に格納する。そして、処理の実行時、このコンピュータは、自己の記憶部に格納されたプログラムを読み取り、読み取ったプログラムに従った処理を実行する。また、このプログラムの別の実施形態として、コンピュータが可搬型記録媒体から直接プログラムを読み取り、そのプログラムに従った処理を実行することとしてもよい。さらに、このコンピュータにサーバコンピュータからプログラムが転送されるたびに、逐次、受け取ったプログラムに従った処理を実行することとしてもよい。また、サーバコンピュータから、このコンピュータへのプログラムの転送は行わず、その実行指示と結果取得のみによって処理機能を実現する、いわゆるASP(Application Service Provider)型のサービスによって、上述の処理を実行する構成としてもよい。なお、プログラムには、電子計算機による処理の用に供する情報であってプログラムに準ずるもの(コンピュータに対する直接の指令ではないがコンピュータの処理を規定する性質を有するデータ等)を含むものとする。

【0060】

また、コンピュータ上で所定のプログラムを実行させることにより、各装置を構成することとしたが、これらの処理内容の少なくとも一部をハードウェア的に実現することとしてもよい。

【 図 1 】

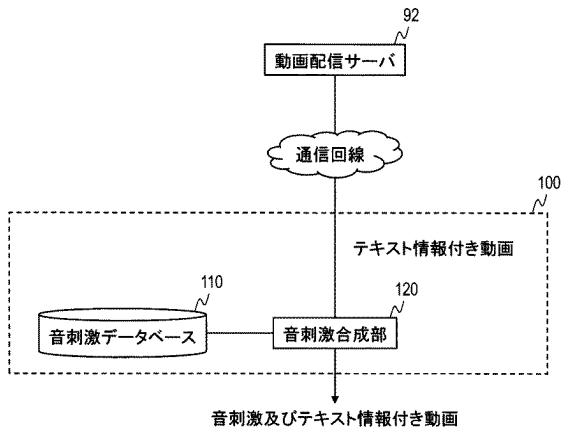


図1

【 図 2 】

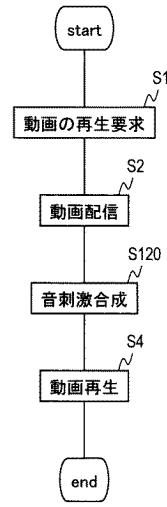


図2

【 図 3 】

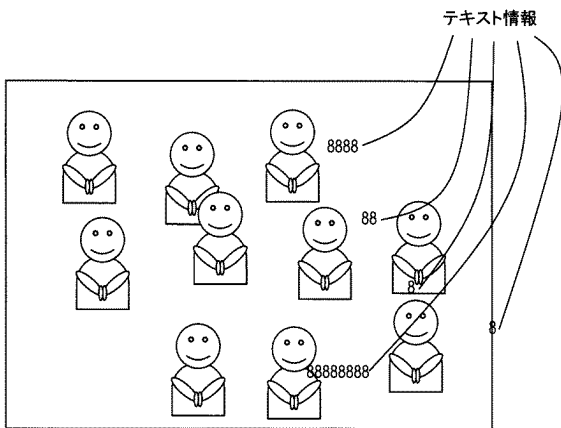


図3

【 図 4 】

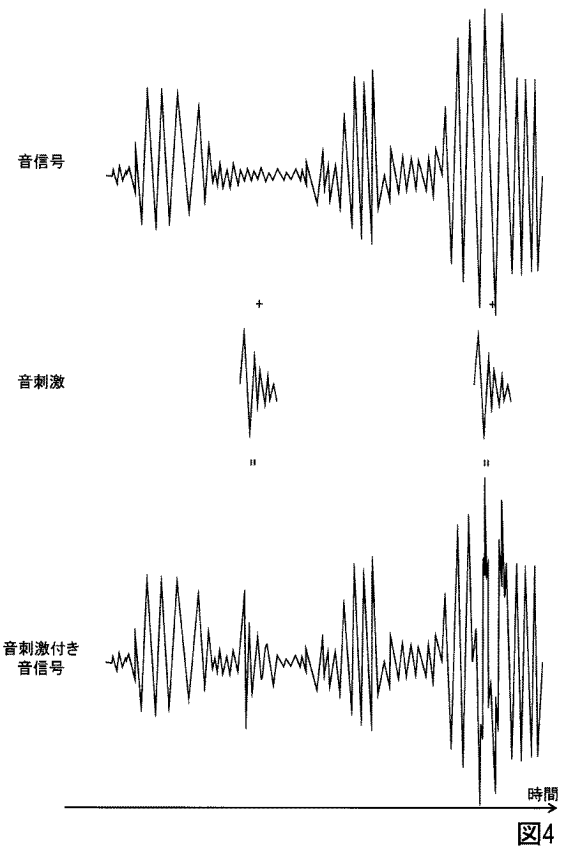


図4

【 図 5 】

	テキスト情報	音刺激
1	w	
2	8	
...

図5

【 図 6 】

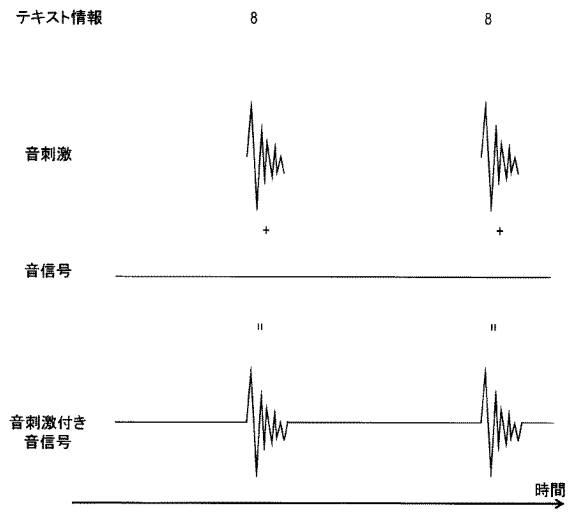


図6

【 図 7 】

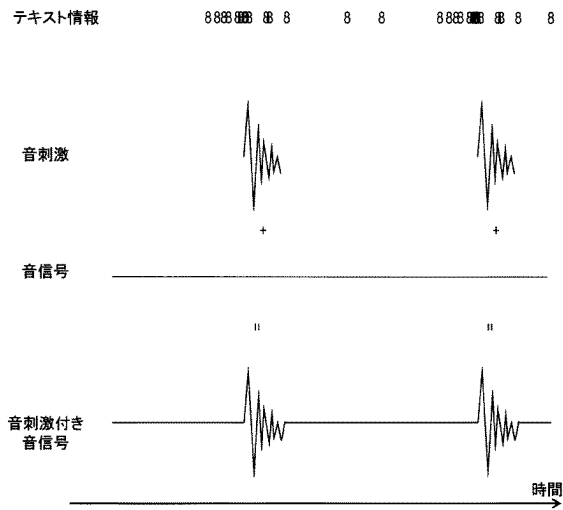


図7

【 図 8 】

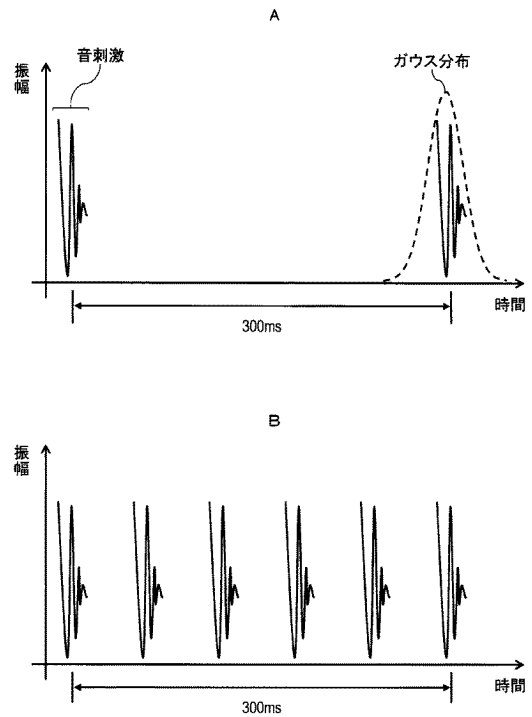


図8

【 図 9 】

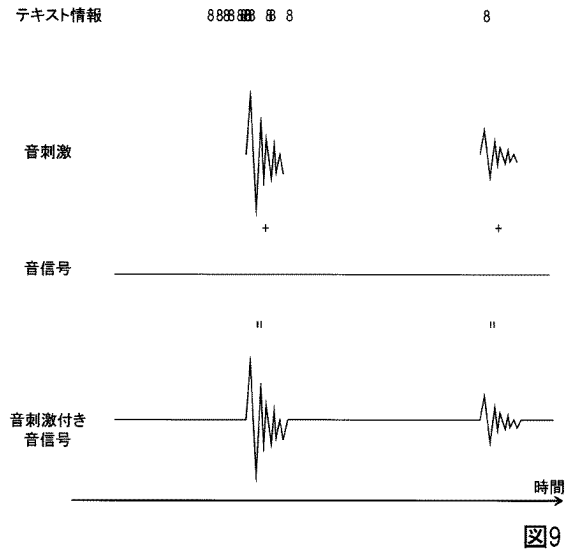


図9

【 図 10 】




	視覚情報	(意味)	音刺激
1	「w」	(笑い)	
2	「(笑)」		
3	☺		
4	「8」	(拍手)	
5			
..	...		

図10

フロントページの続き

- (72)発明者 佐藤 尚
東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
- (72)発明者 ガブリエル バプロ ナバ
東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
- (72)発明者 守谷 健弘
東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内

審査官 岩田 淳

- (56)参考文献 特開2014-212490(JP,A)
特開2014-063145(JP,A)
特開2012-133662(JP,A)

(58)調査した分野(Int.Cl., DB名)

G10L 13/00 - 13/10
H04N 21/00 - 21/858