

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6900536号  
(P6900536)

(45) 発行日 令和3年7月7日(2021.7.7)

(24) 登録日 令和3年6月18日(2021.6.18)

(51) Int.Cl. F I  
G 1 O L 13/047 (2013.01) G 1 O L 13/047 Z

請求項の数 14 (全 21 頁)

<p>(21) 出願番号 特願2020-56285 (P2020-56285)                  (22) 出願日 令和2年3月26日(2020.3.26)                  (65) 公開番号 特開2021-56489 (P2021-56489A)                  (43) 公開日 令和3年4月8日(2021.4.8)                  審査請求日 令和2年3月26日(2020.3.26)                  (31) 優先権主張番号 201910927040.3                  (32) 優先日 令和1年9月27日(2019.9.27)                  (33) 優先権主張国・地域又は機関                  中国 (CN)</p>	<p>(73) 特許権者 513224353                  バイドゥ オンライン ネットワーク テ                  クノロジー (ベイジン) カンパニー                  リミテッド                  中華人民共和国、100085 ベイジン                  ハイディエン ディストリクト、シャン                  ディ 10ティーエイチ ストリート、バ                  イドゥ キャンパス、ナンバー 10、3                  /フロア                  (74) 代理人 100118913                  弁理士 上田 邦生                  (74) 代理人 100142789                  弁理士 柳 順一郎                  (74) 代理人 100163050                  弁理士 小栗 真由美</p>
---	---

最終頁に続く

(54) 【発明の名称】 音声合成モデルのトレーニング方法、装置、電子機器及び記憶媒体

(57) 【特許請求の範囲】

【請求項 1】

音声合成モデルのトレーニング方法であって、  
前記方法は、

現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップと、

符号化標記された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを融合して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの一つの加重組み合わせを取得するステップと、

前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得するステップと、

前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を前記トレーニング対象モデルのデコーダの入力とし、前記デコーダの出力端で前記現在のサンプルの音声メルスペクトル出力を取得するステップと、を含むことを特徴とする、音声合成モデルのトレーニング方法。

【請求項 2】

10

20

前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップは、

前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを一つの共有のエンコーダに入力するステップと、

前記共有のエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップと、を含むことを特徴とする、請求項 1 に記載の音声合成モデルのトレーニング方法。

【請求項 3】

前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップは、

前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立した畳み込み層変換モジュールにそれぞれ入力し、各独立した畳み込み層変換モジュールの出力端で畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをそれぞれ取得するステップと、

畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをシーケンス変換ニューラルネットワークモジュールの入力とし、前記シーケンス変換ニューラルネットワークモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップと、を含むことを特徴とする、請求項 1 に記載の音声合成モデルのトレーニング方法。

【請求項 4】

前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップは、

前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力するステップと、

各独立したエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップと、を含むことを特徴とする、請求項 1 に記載の音声合成モデルのトレーニング方法。

【請求項 5】

前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とする前に、前記方法は、

前記現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換するステップと、

前記音節のベクトル標記と前記漢字のベクトル標記とを前記音素のベクトル標記と同じ長さに変換して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを取得し、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを前記トレーニング対象モデルのエンコーダの入力とする操作を実行するステップと、をさらに含むことを特徴とする、請求項 1 に記載の音声合成モデルのトレーニング方法。

【請求項 6】

前記音素入力シーケンスは、声調入力シーケンスと、児化音入力シーケンスと、句読点入力シーケンスと、35個の独立した韻母の入力シーケンスとを含み、

前記音素入力シーケンスは、106個の音素単位を含み、各音素単位は、106ビットを含み、前記106ビットにおいて、有効ビットの値は1であり、非有効ビットの値は0であり、

10

20

30

40

50

前記漢字入力シーケンスは、3000個の漢字の入力シーケンスを含み、  
前記音節入力シーケンスは、508個の音節の入力シーケンスを含むことを特徴とする、請求項1に記載の音声合成モデルのトレーニング方法。

【請求項7】

音声合成モデルのトレーニング装置であって、  
前記装置は、入力モジュールと、融合モジュールと、出力モジュールと、を含み、  
前記入力モジュールは、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得し、

10

前記融合モジュールは、符号化標記された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを融合して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの一つの加重組み合わせを取得し、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得し、

前記出力モジュールは、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を前記トレーニング対象モデルのデコーダの入力とし、前記デコーダの出力端で前記現在のサンプルの音声メルスペクトル出力を取得することを特徴とする、音声合成モデルのトレーニング装置。

20

【請求項8】

前記入力モジュールは、具体的には、  
前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを一つの共有のエンコーダに入力し、  
前記共有のエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得することを特徴とする、請求項7に記載の音声合成モデルのトレーニング装置。

【請求項9】

前記入力モジュールは、具体的には、  
前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立した畳み込み層変換モジュールにそれぞれ入力し、各独立した畳み込み層変換モジュールの出力端で畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをそれぞれ取得し、

30

畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをシーケンス変換ニューラルネットワークモジュールの入力とし、前記シーケンス変換ニューラルネットワークモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得することを特徴とする、請求項7に記載の音声合成モデルのトレーニング装置。

【請求項10】

40

前記入力モジュールは、具体的には、  
前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力し、  
各独立したエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得することを特徴とする、請求項7に記載の音声合成モデルのトレーニング装置。

【請求項11】

前記装置は、変換モジュールをさらに含み、  
前記変換モジュールは、前記現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換し、前記音節のベクトル標記と前記漢字のベクトル

50

ル標記とを前記音素のベクトル標記と同じ長さに変換して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを取得し、

前記入力モジュールは、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを前記トレーニング対象モデルのエンコーダの入力とする操作を実行することを特徴とする、請求項7に記載の音声合成モデルのトレーニング装置。

#### 【請求項12】

前記音素入力シーケンスは、声調入力シーケンスと、児化音入力シーケンスと、句読点入力シーケンスと、35個の独立した韻母の入力シーケンスとを含み、

前記音素入力シーケンスは、106個の音素単位を含み、

各音素単位は、106ビットを含み、前記106ビットにおいて、有効ビットの値は1であり、非有効ビットの値は0であり、

前記漢字入力シーケンスは、3000個の漢字の入力シーケンスを含み、

前記音節入力シーケンスは、508個の音節の入力シーケンスを含むことを特徴とする、請求項7に記載の音声合成モデルのトレーニング装置。

#### 【請求項13】

電子機器であって、

少なくとも一つのプロセッサと、

前記少なくとも一つのプロセッサと通信可能に接続されるメモリと、を含み、

前記メモリには前記少なくとも一つのプロセッサによって実行可能な命令が記憶され、

前記命令が前記少なくとも一つのプロセッサによって実行される場合に、前記少なくとも一つのプロセッサが請求項1から6のいずれかに記載の方法を実行することを特徴とする、電子機器。

#### 【請求項14】

コンピュータ命令が記憶されている非一時的なコンピュータ読み取り可能な記憶媒体であって、

前記コンピュータ命令は、前記コンピュータに請求項1から6のいずれかに記載の方法を実行させることを特徴とする、非一時的なコンピュータ読み取り可能な記憶媒体。

#### 【発明の詳細な説明】

##### 【技術分野】

##### 【0001】

本出願は、人工知能技術の分野に関し、さらに、コンピュータ知能音声の分野に関し、特に、音声合成モデルのトレーニング方法、装置、電子機器及び記憶媒体に関する。

##### 【背景技術】

##### 【0002】

音声合成の分野では、WaveNetやWaveRNNなどのニューラルネットワークに基づく方法は、合成音声の音質及び自然度を大きく改善する。このような方法は、通常、フロントエンドシステムテキストに基づいて言語特徴を抽出し、基本周波数及び時間などの情報を予測する必要がある。Google(Google)(登録商標)が提供するエンドツーエンドモデリングのTacotronモデルは、大量の専門知識を必要とする複雑なフロントエンドシステムを脱却し、シーケンス変換モデルを介してサウンドライブラリにおける音声のリズム及び感情などの情報を自動的に学習するため、合成された音声は、表現力の面で特に優れている。しかし、Tacotronモデルの中国語への応用は、多くの挑戦がある。主に中国語では漢字の数が多く、一般的に使用される漢字は数千個あり、同じ音の漢字の現象が非常に一般的であり、同音異字の発音方式には違いがあり、同じ漢字であっても異なる単語又は文脈中では発音方式も異なるからである。

##### 【0003】

従来技術において、Tacotronモデルを中国語での応用を実現した技術案では、多くが三つに分けられる。(1)Tacotronの英語での応用と類似し、直接漢字を

10

20

30

40

50

入力要素とし、(2)漢字を音節に縮めて入力要素とし、(3)音節を音素に分割して入力要素とする。上記の技術案(1)を採用すると、漢字の数が多いため、通常、音声合成のトレーニングに使用されるサウンドライブラリは数時間から数十時間程度の規模であるので、直接漢字をモデルの入力要素とすると、データが少ないため、多くの低頻度の漢字の発音が十分に学習することができない。上記の技術案(2)及び(3)を採用して、音素又は音節を入力要素とすると、漢字の数が少ないという問題を解決することができ、漢字中の同じ音の漢字は、共有ユニットによってより十分なトレーニングを取得することができる。しかし、異なる漢字は、同じ発音であっても、発音方式が明確に違い、よく見られる虚語では、通常、発音が弱く、実語は発音が明確である。このため、Tacotronモデルでは、虚語の発音方式で学習する傾向があるため、合成効果が良くないという問題がある。また、音素を入力要素とすると、ある韻母は、単独に一つの完全な音節とすることができる場合があり、この二つのケースでは、韻母の発音が実際に一定の違いがあるので、独立した音節の韻母としてはより完全な発音プロセスが必要となるが、音素に基づくモデルは、この二つの場合を区別できないため、韻母が独立した場合の発音が不十分になるという別の問題も存在する。

10

【発明の概要】

【発明が解決しようとする課題】

【0004】

これを考慮して、本出願の実施例は、発音効果を効果的に改善し、音声製品に高い表現力と高い自然度の中国語合成音声を提供することができる音声合成モデルのトレーニング方法、装置、電子機器及び記憶媒体を提供する。

20

【課題を解決するための手段】

【0005】

第1の態様では、本出願の実施例は、音声合成モデルのトレーニング方法を提供し、前記方法は、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの符号化標記を取得するステップと、符号化標記された前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスを融合して、前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの一つの加重組み合わせを取得するステップと、前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの加重組み合わせの各時点における加重平均を取得するステップと、前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの加重組み合わせの各時点における加重平均を前記トレーニング対象モデルのデコーダの入力とし、前記デコーダの出力端で前記現在のサンプルの音声メル(Mel)スペクトル出力を取得するステップと、を含む。

30

【0006】

上記の実施例は、トレーニング対象モデルの入力端の入力テキスト及びトレーニング対象モデルの出力端の出力音声によって、トレーニング対象モデルのエンコーダ及びデコーダの共同トレーニングを実現する。本出願では、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合する技術的手段を採用するため、従来技術では、音節入力シーケンス又は音素入力シーケンス又は漢字入力シーケンスのみしか採用しないために音声合成効果が良くないという技術問題を克服し、発音効果を効果的に改善し、音声製品に高い表現力と高い自然度の中国語合成音声の技術的效果を提供するという利点又は有益な効果を奏する。

40

【0007】

上記の実施例では、前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エン

50

コーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップは、前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスを一つの共有のエンコーダに入力するステップと、前記共有のエンコーダの出力端で前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの符号化標記を取得するステップと、を含む。

【0008】

上記の実施例は、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを一つの共有のエンコーダに入力することによって、共有のエンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得することができるという利点又は有益な効果を奏する。

10

【0009】

上記の実施例では、前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの符号化標記を取得するステップは、前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスを入力到三つの独立した畳み込み層変換モジュールにそれぞれ入力し、各独立した畳み込み層変換モジュールの出力端で畳み込み変換された前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスをそれぞれ取得するステップと、畳み込み変換された前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスをシーケンス変換ニューラルネットワークモジュールの入力とし、前記シーケンス変換ニューラルネットワークモジュールの出力端で前記音節入力シーケンス、前記音素入力シーケンス、及び前記漢字入力シーケンスの符号化標記を取得するステップと、を含む。

20

【0010】

上記の実施例は、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを、三つの独立した畳み込み層変換モジュールにそれぞれ入力し、畳み込み変換された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをシーケンス変換ニューラルネットワークモジュールの入力とすることにより、シーケンス変換ニューラルネットワークモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得することができるという利点又は有益な効果を奏する。

30

【0011】

上記の実施例では、前記現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップは、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力するステップと、各独立したエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップと、を含む。

40

【0012】

上記の実施例は、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力し、各独立したエンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得するという利点又は有益な効果を有する。実験により、音素の独立した韻母及び音素、音節、漢字を融合する三つの技術案は、エンドツーエンドの中国語音声合成の問題をある程度で解決することができ、そのうち、独立エンコーダの効果が最適である。測定の結果は、発音問題の割合が2%から0.4%に減少することを示す。

【0013】

上記の実施例では、前記現在のサンプルの音節入力シーケンスと、音素入力シーケンス

50

と、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とする前に、前記方法は、前記現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換するステップと、前記音節のベクトル標記と前記漢字のベクトル標記とを前記音素のベクトル標記と同じ長さに変換して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを取得し、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスと、を前記トレーニング対象モデルのエンコーダの入力とする操作を実行するステップと、をさらに含む。

【0014】

上記の実施例は、音節のベクトル標記と漢字のベクトル標記を音素のベクトル標記と同じ長さに変換して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを取得できるため、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とする操作を実行することができるという利点又は有益な効果を有する。

10

【0015】

上記の実施例では、前記音素入力シーケンスは、声調入力シーケンスと、児化音入力シーケンスと、句読点入力シーケンスと、35個の独立した韻母の入力シーケンスとを含み、前記音素入力シーケンスは、106個の音素単位を含み、各音素単位は、106ビットを含む。前記106ビットにおいて、有効ビットの値は1であり、非有効ビットの値は0であり、前記漢字入力シーケンスは、3000個の漢字の入力シーケンスを含み、前記音節入力シーケンスは、508個の音節の入力シーケンスを含む。

20

【0016】

上記の実施例は、それぞれ音節及び漢字から有効な情報を抽出して発音効果を改善し、特に、同音異字の場合は、発音問題を著しく減少するという利点又は有益な効果を有する。本出願は、製品に高い表現力高い自然度の中国語合成音声を提供することができ、ユーザのヒューマン-コンピュータ-インタラクション体験を効果的に向上させることができ、ユーザの粘着性を向上させ、バイドゥAPP、スマートスピーカ、及びマップナビゲーションシステムのプロモーションに有利である。

【0017】

第2の態様では、本出願は、音声合成モデルのトレーニング装置をさらに提供する。前記装置は、入力モジュールと、融合モジュールと、出力モジュールと、を含む。

30

前記入力モジュールは、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

前記融合モジュールは、符号化標記された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを融合して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの一つの加重組み合わせを取得し、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得する。

40

前記出力モジュールは、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を前記トレーニング対象モデルのデコーダの入力とし、前記デコーダの出力端で前記現在のサンプルの音声メルスペクトル出力を取得する。

【0018】

上記の実施例では、前記入力モジュールは、具体的には、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを一つの共有のエンコーダに入力するステップと、前記共有のエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得するステップと、を含

50

む。

【 0 0 1 9 】

上記の実施例では、前記入力モジュールは、具体的には、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを、三つの独立した畳み込み層変換モジュールにそれぞれ入力し、各独立した畳み込み層変換モジュールの出力端で畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをそれぞれ取得し、畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをシーケンス変換ニューラルネットワークモジュールの入力とし、前記シーケンス変換ニューラルネットワークモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

10

【 0 0 2 0 】

上記の実施例では、前記入力モジュールは、具体的には、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力し、各独立したエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

【 0 0 2 1 】

上記の実施例では、前記装置は、前記現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換し、前記音節のベクトル標記と前記漢字のベクトル標記とを前記音素のベクトル標記と同じ長さに変換して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを取得するための変換モジュールをさらに含み、前記入力モジュールは、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを前記トレーニング対象モデルのエンコーダの入力とする操作を実行する。

20

【 0 0 2 2 】

上記の実施例では、前記音素入力シーケンスは、声調入力シーケンスと、児化音入力シーケンスと、句読点入力シーケンスと、35個の独立した韻母の入力シーケンスとを含み、前記音素入力シーケンスは、106個の音素単位を含む。各音素単位は、106ビットを含み、前記106ビットにおいて、有効ビットの値は1であり、非有効ビットの値は0であり、前記漢字入力シーケンスは、3000個の漢字の入力シーケンスを含み、前記音節入力シーケンスは、508個の音節の入力シーケンスを含む。

30

【 0 0 2 3 】

第3の態様では、本出願の実施例は、電子機器を提供する。電子機器は、一つ又は複数のプロセッサと、一つ又は複数のプログラムを記憶するためのメモリと、を含み、前記一つ又は複数のプログラムが前記一つ又は複数のプロセッサによって実行される場合、前記一つ又は複数のプロセッサが、本出願の任意の実施例に記載の音声合成モデルのトレーニング方法を実現する。

【 0 0 2 4 】

第4の態様では、本出願の実施例は、コンピュータ命令が記憶されている記憶媒体を提供する。記憶媒体は、当該プログラムがプロセッサによって実行される場合に、本出願の任意の実施例に記載の音声合成モデルのトレーニング方法が実現される。

40

【 0 0 2 5 】

上記の出願中の一つの実施例は、以下のような利点又は有益な効果を有する。本出願で提供される音声合成モデルのトレーニング方法、装置、電子機器及び記憶媒体は、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得し、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得し、再び音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケ

50

スとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得し、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出力を取得する。つまり、本出願は、トレーニング対象モデルの入力端の入力テキスト及びトレーニング対象モデルの出力端の出力音声により、トレーニング対象モデルのエンコーダ及びデコーダの共同トレーニングを実現する。本出願では、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合する技術的手段を採用するため、従来技術では音節入力シーケンス又は音素入力シーケンス又は漢字入力シーケンスのみを採用するために音声合成効果が良くないという技術問題を克服し、発音効果を効果的に改善し、音声製品に高い表現力と高い自然度の中国語合成音声の技術的效果を提供することができる。また、本出願の実施例の技術案は、簡単で便利に実現され、普及しやすいので、適用範囲がより広くなる。

10

【0026】

上記の選択可能な方式が有する他の効果は、以下で具体的な実施例を組み合わせで説明される。

【図面の簡単な説明】

【0027】

図面は、本技術案をよりよく理解するために使用され、本出願の構成を限定するものではない。

20

【図1】本出願の実施例1により提供される音声合成モデルのトレーニング方法の概略フローチャートである。

【図2】本出願の実施例2により提供される音声合成モデルのトレーニング方法の概略フローチャートである。

【図3】本出願の実施例2により提供されるTacotronモデルの概略構成図である。

【図4】本出願の実施例3により提供される音声合成モデルのトレーニング装置の概略構成図である。

【図5】本出願の実施例の音声合成モデルのトレーニング方法を実現するための電子機器のブロック図である。

30

【発明を実施するための形態】

【0028】

以下、図面を組み合わせで本出願の例示的な実施例を説明し、理解を容易にするために本出願の実施例の様々な詳細を含んでいるが、それらは単なる例示であると見なすべきである。したがって、当業者は、本出願の範囲及び精神から逸脱することなく、ここで説明される実施例に対して様々な変更と修正を行うことができることを認識されたい。同様に、明確及び簡潔にするために、以下の説明では、周知の機能及び構造の説明を省略する。

【0029】

実施例1

40

図1は、本出願の実施例1により提供される音声合成モデルのトレーニング方法の概略フローチャートである。当該方法は、音声合成モデルのトレーニング装置又は電子機器により実行することができ、当該装置又は電子機器は、ソフトウェア及び/又はハードウェアの方式によって実現することができ、当該装置又は電子機器は、任意のネットワーク通信機能を有するスマートデバイスに集積することができる。図1に示すように、音声合成モデルのトレーニング方法は、以下のステップを含むことができる。

S101：現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得する。

50



## 【 0 0 3 6 】

本出願の具体的な実施例では、電子機器は、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出力を取得することができる。

## 【 0 0 3 7 】

本出願の具体的な実施例では、Tacotronモデルは、古典的なエンコーダ・デコーダ構造に基づく。エンコーダは、入力された要素シーケンス（英語では、通常アルファベット及び句読点など）に基づき、畳み込み及びシーケンス変換ニューラルネットワークの変換を経て、各入力要素の符号化標記を取得する。デコーダは、前のフレームの出力メルスペクトルを入力とし、アテンションメカニズムを利用してエンコーダ出力の一つの加重組み合わせ表現を取得し、次に、LSTMなどの変換を経て、二つの出力を生成する。一つは、現在のフレーム出力のメルスペクトルであり、もう一つは、終了するか否かを判断する停止確率である。停止確率が50%より大きい場合には合成が終了し、そうでなければ現在の出力を次のフレームの入力とし、この自己回帰プロセスを継続する。このモデルでは、エンコーダは、各入力要素の符号化を担当し、デコーダは、符号化に基づいて現在合成された音を決定するとともに、LSTMの記憶機能を利用して秩序的に生成する。当該モデルは、典型的な一対多のマッピングモデルであり、同じ内容は、異なるリズム、異なる感情の音声に対応する。トレーニングセットにおける異なる音声（出力）が同じ文字（入力）に対応する場合、モデルが最終的に学習した発音は、統計的な平均効果を反映する。本出願は、このような一対多のマッピング関係を減少するために、モデルが異なる文脈で適切な発音方式で合成することを学習できるようにする。音素シーケンスが最高のカバー能力を有することを考慮すると、セット以外の発音要素が発生する問題はなく、106個の音素単位をモデル入力の基礎要素として選択することができ、各要素は、十分なデータを取得して十分にトレーニングすることができる。入力はone-hotの形式であり、embedding層を経て固定次元の稠密なベクトル標記に変換され、声調、兒化音、及び句読点などの特徴は、同様にembeddingを経て同じ次元のベクトルに変換されて、音素ベクトルと加算されてニューラルネットワークに送られる。韻母が独立した場合の発音特性をよりよく学習するために、本出願は、音素における35個の独立した韻母を単独でモデリングし、声母後に出現された韻母とは、二つの異なる要素と見なされる。実験により、このようなモデリング戦略は、独立した韻母発音が不明確であるという問題をうまく解決する。さらに、同音異字の発音特性を区分するために、本出願は、音節及び漢字を補助情報としてネットワークに入力し、補助モデルが異なる字の発音特性を区分できるようにサポートする。そのうち、無調音節の数は508であり、漢字要素は、トレーニングセットにおける高頻度3000漢字及び508音節の計3508個の要素を選択し、ある漢字が高頻度3000字に属していない場合には、対応する音節要素に縮め、より高いカバー率を確保する。

## 【 0 0 3 8 】

本出願の実施例により提供される音声合成モデルのトレーニング方法は、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得し、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得し、再び音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得し、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出

10

20

30

40

50

力を取得する。つまり、本出願は、トレーニング対象モデルの入力端の入力テキスト及びトレーニング対象モデルの出力端の出力音声によって、トレーニング対象モデルのエンコーダ及びデコーダの共同トレーニングを実現する。本出願では、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合する技術的手段を採用するため、従来技術では、音節入力シーケンス又は音素入力シーケンス又は漢字入力シーケンスのみが採用されるために音声合成効果が良くないという技術問題を克服し、発音効果を効果的に改善し、音声製品に高い表現力と高い自然度の中国語合成音声の技術的效果を提供する。また、本出願の実施例の技術案は、簡単で便利に実現されるので、普及しやすく、適用範囲がより広くなる。

**【 0 0 3 9 】**

## 実施例 2

図 2 は、本出願の実施例 2 により提供される音声合成モデルのトレーニング方法の概略フローチャートである。図 2 に示すように、音声合成モデルのトレーニング方法は、以下のようなステップを含む。

S 2 0 1 : 現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換する。

**【 0 0 4 0 】**

本出願の具体的な実施例では、電子機器は、現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換することができる。具体的には、電子機器は、現在のサンプルにおける音素を第 1 の長さのベクトル標記に変換することができ、現在のサンプルにおける音節及び漢字を第 2 の長さのベクトル標記に変換することができる。第 1 の長さが第 2 の長さより大きい。

**【 0 0 4 1 】**

S 2 0 2 : 音節のベクトル標記と漢字のベクトル標記を音素のベクトル標記と同じ長さに変換して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを取得する。

**【 0 0 4 2 】**

本出願の具体的な実施例では、電子機器は、音節のベクトル標記と漢字のベクトル標記を音素のベクトル標記と同じ長さに変換して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを取得することができる。具体的には、電子機器は、第 1 の長さの音素のベクトル標記を音素入力シーケンスとし、音節のベクトル標記と漢字のベクトル標記を第 2 の長さから第 1 の長さに変換し、変換された音節のベクトル標記及び漢字のベクトル標記を音節入力シーケンス及び漢字入力シーケンスとすることができる。

**【 0 0 4 3 】**

S 2 0 3 : 現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得する。

**【 0 0 4 4 】**

本出願の具体的な実施例では、電子機器は、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得することができる。具体的には、電子機器は、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを一つの共有のエンコーダに入力し、共有のエンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得することができる。好ましくは、電子機器は、さらに、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを、三つの独立した畳み込み層変換モジュールにそれぞれ入力し、各独立した畳み込み層変換モジュールの出力端で畳み込み変換された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをそれぞれ取得し、畳み込み変換された音節入力シーケンスと、音

10

20

30

40

50

素入力シーケンスと、漢字入力シーケンスとをシーケンス変換ニューラルネットワークモジュールの入力とし、シーケンス変換ニューラルネットワークモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得することができる。ここのシーケンス変換ニューラルネットワークは、RNN、LSTM、GRU、Transformerを含むが、これらに限定されない。好ましくは、電子機器は、さらに、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力するステップと、各独立したエンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得するステップとを含むことができる。

【0045】

S204：符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得する。

【0046】

本出願の具体的な実施例では、電子機器は、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得することができる。例えば、電子機器は、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを線形的に重ね合わせて、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得することができる。

【0047】

S205：音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得する。

【0048】

本出願の具体的な実施例では、電子機器は、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得することができる。

【0049】

S206：音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出力を取得する。

【0050】

図3は、本出願の実施例2により提供されるTacotronモデルの概略構成図である。図3に示すように、Tacotronモデルは、古典的なエンコーダ・デコーダ構造に基づいている。現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得し、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得し、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得し、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出力を取得する。

【0051】

10

20

30

40

50

本出願の具体的な実施例では、音素入力シーケンスは、声調入力シーケンス、児化音入力シーケンスと、句読点入力シーケンスと、35個の独立した韻母の入力シーケンスとを含むことができる。音素入力シーケンスは、106個の音素単位を含み、各音素単位は、106ビットをむ。106ビットにおいて、有効ビットの値は1であり、非有効ビットの値は0である。漢字入力シーケンスは、3000個の漢字の入力シーケンスを含み、音節入力シーケンスは、508個の音節の入力シーケンスを含む。

#### 【0052】

実験により、音素の独立した韻母と音素、音節、漢字を融合する三つの技術案は、エンドツーエンドの中国語音声合成の問題をある程度で解決することができる。そのうち、独立エンコーダの効果が最も優れている。測定の結果は、発音問題の割合が2%から0.4%に減少することを示す。詳細の分析結果は、音素が発音のタイプを基本的に決定するが、場合によって、音節を変更すると、発音に一定の影響があり、漢字を変更すると発音方式にのみ影響する。これらの結果から、モデルは、それぞれ音節及び漢字から有効な情報を抽出して発音効果を改善し、特に、同音異字の場合は、発音問題を著しく減少することを証明する。本出願は、製品に高い表現力高い自然度の中国語合成音声を提供することができ、ユーザのヒューマン-コンピュータ・インタラクション体験を効果的に向上させることができるので、ユーザの粘着性を向上させ、バイドゥAPP、スマートスピーカ、及びマップナビゲーションシステムのプロモーションに有利である。

#### 【0053】

本出願の実施例により提供される音声合成モデルのトレーニング方法は、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得し、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得し、再び音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得し、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出力を取得する。つまり、本出願は、トレーニング対象モデルの入力端の入力テキスト及びトレーニング対象モデルの出力端の出力音声によって、トレーニング対象モデルのエンコーダ及びデコーダの共同トレーニングを実現する。本出願では、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合する技術的手段を採用したため、従来技術では、音節入力シーケンス又は音素入力シーケンス又は漢字入力シーケンスのみを採用したことで音声合成効果が良くないという技術問題を克服し、発音効果を効果的に改善し、音声製品に高い表現力と高い自然度の中国語合成音声の技術的效果を提供することができる。また、本出願の実施例の技術案は、簡単で便利に実現され、普及しやすく、適用範囲がより広がる。

#### 【0054】

##### 実施例3

図4は、本出願の実施例3により提供される音声合成モデルのトレーニング装置の概略構成図である。図4に示すように、前記装置400は、入力モジュール401と、融合モジュール402と、出力モジュール403と、を含む。

前記入力モジュール401は、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、前記エンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

前記融合モジュール402とは、符号化標記された前記音節入力シーケンスと、前記音

10

20

30

40

50

素入力シーケンスと、前記漢字入力シーケンスとを融合して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの一つの加重組み合わせを取得し、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得する。

前記出力モジュール403は、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの加重組み合わせの各時点における加重平均を前記トレーニング対象モデルのデコーダの入力とし、前記デコーダの出力端で前記現在のサンプルの音声メルスペクトル出力を取得する。

10

【0055】

さらに、前記入力モジュール401は、具体的には、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを一つの共有のエンコーダに入力し、前記共有のエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

【0056】

さらに、前記入力モジュール401は、具体的には、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを、三つの独立した畳み込み層変換モジュールにそれぞれ入力し、各独立した畳み込み層変換モジュールの出力端で畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをそれぞれ取得し、畳み込み変換された前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとをシーケンス変換ニューラルネットワークモジュールの入力とし、前記シーケンス変換ニューラルネットワークモジュールの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

20

【0057】

さらに、前記入力モジュール401は、具体的には、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを三つの独立したエンコーダにそれぞれ入力し、各独立したエンコーダの出力端で前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとの符号化標記を取得する。

30

【0058】

さらに、前記装置は、前記現在のサンプルにおける音素と、音節と、漢字とを個々の固定次元のベクトル標記にそれぞれ変換し、前記音節のベクトル標記と前記漢字のベクトル標記とを前記音素のベクトル標記と同じ長さに変換して、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを取得するための変換モジュール404（図示せず）をさらに含む。

前記入力モジュール401は、前記音節入力シーケンスと、前記音素入力シーケンスと、前記漢字入力シーケンスとを前記トレーニング対象モデルのエンコーダの入力とする操作を実行する。

【0059】

さらに、前記音素入力シーケンスは、声調入力シーケンスと、児化音入力シーケンスと、句読点入力シーケンスと、35個の独立した韻母の入力シーケンスとを含む。前記音素入力シーケンスは、106個の音素単位を含み、各音素単位は、106ビットを含む。前記106ビットにおいて、有効ビットの値は1であり、非有効ビットの値は0である。前記漢字入力シーケンスは、3000個の漢字の入力シーケンスを含み、前記音節入力シーケンスは、508個の音節の入力シーケンスを含む。

40

【0060】

上記の音声合成モデルのトレーニング装置は、本発明の任意の実施例により提供される方法を実行することができ、実行方法に対応する機能モジュール及び有益な効果を有する。本実施例で詳細に説明されていない技術的詳細は、本発明の任意の実施例により提供さ

50

れる音声合成モデルのトレーニング方法を参照することができる。

【0061】

実施例4

本出願の実施例によれば、本出願は、電子機器及び読み取り可能な記憶媒体をさらに提供する。

【0062】

図5には、本出願の実施例に係る音声合成モデルのトレーニング方法の電子機器のブロック図が示されている。電子機器は、ラップトップコンピュータ、デスクトップコンピュータ、ワークステーション、パーソナルデジタルアシスタント、サーバ、ブレードサーバ、大型コンピュータ、及び他の適切なコンピュータなどの様々な形式のデジタルコンピュータを表すことを目的とする。電子機器は、パーソナルデジタル処理、携帯電話、スマートフォン、ウェアラブルデバイス、他の同様のコンピューティングデバイスなどの様々な形式のモバイルデバイスを表すこともできる。本明細書で示されるコンポーネント、それらの接続と関係、及びそれらの機能は単なる例であり、本明細書の説明及び/又は要求される本出願の実現を制限することを意図しない。

10

【0063】

図5に示すように、当該電子機器は、一つ又は複数のプロセッサ501と、メモリ502と、高速インターフェースと低速インターフェースを含む各コンポーネントを接続するためのインターフェースと、を含む。各コンポーネントは、異なるバスで相互に接続され、共通のマザーボードに取り付けられるか、又は必要に応じて他の方式で取り付けることができる。プロセッサは、外部入力/出力装置(インターフェースに結合されたディスプレイデバイスなど)にGUIの図形情報をディスプレイするためにメモリに記憶されている命令を含む、電子機器内に実行される命令を処理することができる。他の実施方式では、必要であれば、複数のプロセッサ及び/又は複数のバスを、複数のメモリと複数のメモリとともに使用することができる。同様に、複数の電子機器を接続することができ、各機器は、部分的な必要な操作(例えば、サーバレイ、ブレードサーバ、又はマルチプロセッサシステムとする)を提供することができる。図5では、一つのプロセッサ501を例とする。

20

【0064】

メモリ502は、本出願により提供される非一時的なコンピュータ読み取り可能な記憶媒体である。その中、前記メモリには、少なくとも一つのプロセッサによって実行される命令を記憶して、前記少なくとも一つのプロセッサが本出願により提供される音声合成モデルのトレーニング方法を実行することができるようにする。本出願の非一時的なコンピュータ読み取り可能な記憶媒体は、コンピュータが本出願により提供される音声合成モデルのトレーニング方法を実行するためのコンピュータ命令を記憶する。

30

【0065】

メモリ502は、非一時的なコンピュータ読み取り可能な記憶媒体として、本出願の実施例における音声合成モデルのトレーニング方法に対応するプログラム命令/モジュール(例えば、図4に示す入力モジュール401、融合モジュール402、出力モジュール403)ように、非一時的なソフトウェアプログラム、非一時的なコンピュータ実行可能なプログラム及びモジュールを記憶するために用いられる。プロセッサ501は、メモリ502に記憶されている非一時的なソフトウェアプログラム、命令及びモジュールを実行することによって、サーバの様々な機能アプリケーション及びデータ処理を実行する。すなわち上記の方法の実施例における音声合成モデルのトレーニング方法を実現する。

40

【0066】

メモリ502は、ストレージプログラム領域とストレージデータ領域とを含むことができる。ストレージプログラム領域は、オペレーティングシステムや、少なくとも一つの機能に必要なアプリケーションプログラムを記憶することができ、ストレージデータ領域は、音声合成モデルのトレーニング方法に基づく電子機器の使用によって作成されたデータなどを記憶することができる。また、メモリ402は、高速ランダム存取メモリを含むこ

50

とができ、非一時的なメモリをさらに含むことができる。例えば、少なくとも一つのディスクストレージデバイス、フラッシュメモリデバイス、又は他の非一時的なソリッドステートストレージデバイスである。いくつかの実施例では、メモリ502は、プロセッサ501に対して遠隔に設置されたメモリを含むことができ、これらの遠隔メモリは、ネットワークを介して音声合成モデルのトレーニング方法の電子機器に接続されることができる。上記のネットワークの例は、インターネット、イントラネット、ローカルエリアネットワーク、モバイル通信ネットワーク、及びその組み合わせを含むが、これらに限定しない。

#### 【0067】

音声合成モデルのトレーニング方法の電子機器は、入力装置503と出力装置504とをさらに含むことができる。プロセッサ501と、メモリ502と、入力装置503と、出力装置504とは、バス又は他の方式を介して接続することができ、図5では、バスを介して接続することを例示している。

10

#### 【0068】

入力装置503は、入力された数字又は文字情報を受信することができ、及び音声合成モデルのトレーニング方法の電子機器のユーザ設定及び機能制御に関するキー信号入力を生成することができる。例えば、タッチスクリーン、キーパッド、マウス、トラックパッド、タッチパッド、指示杆、一つ又は複数のマウスボタン、トラックボール、ジョイスティックなどの入力装置である。出力装置504は、ディスプレイデバイス、補助照明デバイス（例えば、LED）、及び触覚フィードバックデバイス（例えば、振動モータ）などを含むことができる。当該ディスプレイデバイスは、液晶ディスプレイ（LCD）、発光ダイオード（LED）ディスプレイ、及びプラズマディスプレイを含むことができるが、これらに限定しない。いくつかの実施方式では、ディスプレイデバイスは、タッチスクリーンであってもよい。

20

#### 【0069】

本明細書で説明されるシステムと技術の様々な実施方式は、デジタル電子回路システム、集積回路システム、特定用途向けASIC（特定用途向け集積回路）、コンピュータハードウェア、ファームウェア、ソフトウェア、及び/又はそれらの組み合わせで実現することができる。これらの様々な実施方式は、一つ又は複数のコンピュータプログラムで実施されることを含むことができ、当該一つ又は複数のコンピュータプログラムは、少なくとも一つのプログラマブルプロセッサを含むプログラム可能なシステムで実行及び/又は解釈されることができ、当該プログラマブルプロセッサは、特定用途向け又は汎用プログラマブルプロセッサであってもよく、ストレージシステム、少なくとも一つの入力装置、及び少なくとも一つの出力装置からデータ及び命令を受信し、データ及び命令を当該ストレージシステム、当該少なくとも一つの入力装置、及び当該少なくとも一つの出力装置に伝送することができる。

30

#### 【0070】

これらのコンピューティングプログラム（プログラム、ソフトウェア、ソフトウェアアプリケーション、又はコードとも呼ばれる）は、プログラマブルプロセッサの機械命令、高レベルのプロセス及び/又はオブジェクト指向プログラミング言語、及び/又はアセンブリ/機械言語でこれらのコンピューティングプログラムを実施することを含む。本明細書に使用されるように、用語「機械読み取り可能な媒体」及び「コンピュータ読み取り可能な媒体」は、機械命令及び/又はデータをプログラマブルプロセッサに提供するために使用される任意のコンピュータプログラム製品、機器、及び/又は装置（例えば、磁気ディスク、光ディスク、メモリ、プログラマブルロジックデバイス（PLD））を指し、機械読み取り可能な信号である機械命令を受信する機械読み取り可能な媒体を含む。用語「機械読み取り可能な信号」は、機械命令及び/又はデータをプログラマブルプロセッサに提供するための任意の信号を指す。

40

#### 【0071】

ユーザとのインタラクションを提供するために、コンピュータ上で、ここで説明されて

50

いるシステム及び技術を実施することができる。当該コンピュータは、ユーザに情報を表示するためのディスプレイ装置（例えば、CRT（陰極線管）又はLCD（液晶ディスプレイ）モニタ）と、キーボード及びポインティングデバイス（例えば、マウス又はトラックボール）とを有し、ユーザは、当該キーボード及び当該ポインティングデバイスによって入力をコンピュータに提供することができる。他の種類の装置は、ユーザとのインタラクションを提供するために用いられることもでき、例えば、ユーザに提供されるフィードバックは、任意の形式のセンシングフィードバック（例えば、視覚フィードバック、聴覚フィードバック、又は触覚フィードバック）であってもよく、任意の形式（音響入力と、音声入力と、触覚入力とを含む）でユーザからの入力を受信することができる。

**【0072】**

ここで説明されるシステム及び技術は、バックエンドコンポーネントを含むコンピューティングシステム（例えば、データサーバとする）、又はミドルウェアコンポーネントを含むコンピューティングシステム（例えば、アプリケーションサーバ）、又はフロントエンドコンポーネントを含むコンピューティングシステム（例えば、グラフィカルユーザインタフェース又はウェブブラウザを有するユーザコンピュータ、ユーザは、当該グラフィカルユーザインタフェース又は当該ウェブブラウザによってここで説明されるシステム及び技術の実施方式とインタラクションする）、又はこのようなバックエンドコンポーネントと、ミドルウェアコンポーネントと、フロントエンドコンポーネントの任意の組み合わせを含むコンピューティングシステムで実施することができる。任意の形式又は媒体のデジタルデータ通信（例えば、通信ネットワーク）によってシステムのコンポーネントを相互に接続されることができる。通信ネットワークの例は、ローカルエリアネットワーク（LAN）と、ワイドエリアネットワーク（WAN）と、インターネットとを含む。

**【0073】**

コンピュータシステムは、クライアントとサーバとを含むことができる。クライアントとサーバは、一般に、互いに離れており、通常に通信ネットワークを介してインタラクションする。対応するコンピュータ上で実行され、互いにクライアント-サーバ関係を有するコンピュータプログラムによってクライアントとサーバとの関係が生成される。

**【0074】**

本出願の実施例の技術案によれば、現在のサンプルの音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとをトレーニング対象モデルのエンコーダの入力とし、エンコーダの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの符号化標記を取得し、符号化標記された音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合して、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの一つの加重組み合わせを取得し、再び音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせをアテンションモジュールの入力とし、アテンションモジュールの出力端で音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均を取得し、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとの加重組み合わせの各時点における加重平均をトレーニング対象モデルのデコーダの入力とし、デコーダの出力端で現在のサンプルの音声メルスペクトル出力を取得する。つまり、本出願は、トレーニング対象モデルの入力端の入力テキスト及びトレーニング対象モデルの出力端の出力音声によって、トレーニング対象モデルのエンコーダ及びデコーダの共同トレーニングを実現する。本出願では、音節入力シーケンスと、音素入力シーケンスと、漢字入力シーケンスとを融合する技術的手段を採用するため、従来技術では音節入力シーケンス又は音素入力シーケンス又は漢字入力シーケンスのみを採用したために音声合成効果が良くないという技術問題を克服し、発音効果を効果的に改善し、音声製品に高い表現力と高い自然度の中国語合成音声の技術的效果を提供する。また、本出願の実施例の技術案は、簡単に実現され、普及しやすく、適用範囲がより広くなる。

**【0075】**

上記に示される様々な形式のフローを使用して、ステップを並べ替え、追加、又は削除

10

20

30

40

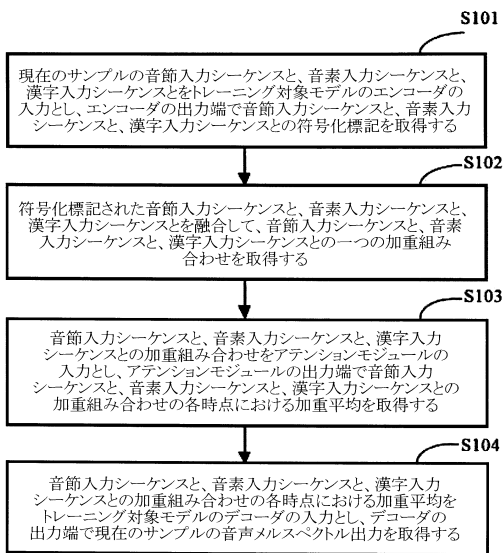
50

することができることを理解されたい。例えば、本出願に記載されている各ステップは、並列に実行されてもよいし、順次的に実行されてもよいし、異なる順序で実行されてもよく、本出願で開示されている技術案が所望の結果を実現することができれば、本明細書では限定されない。

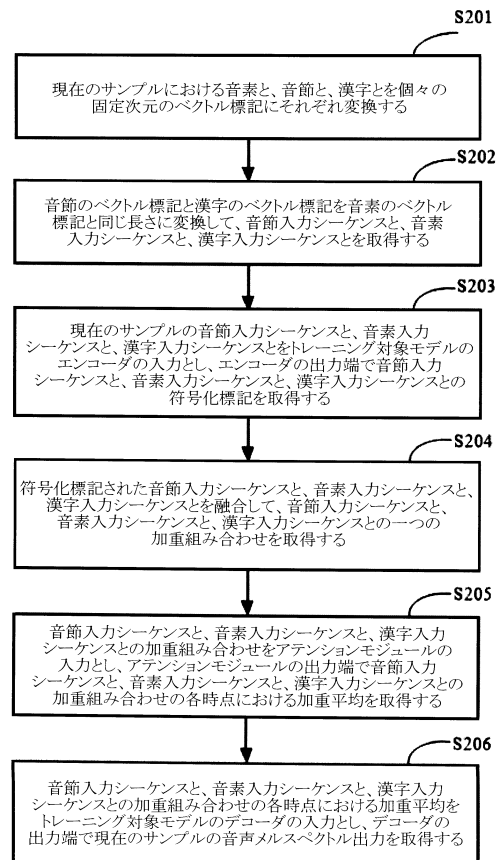
【 0 0 7 6 】

上記の具体的な実施方式は、本出願に対する保護範囲の制限を構成するものではない。当業者は、設計要求と他の要因に応じて、様々な修正、組み合わせ、サブコンビネーション、及び代替を行うことができる。任意の本出願の精神と原則内で行われる修正、同等の置換、及び改善などは、いずれも本出願の保護範囲内に含まれなければならない。

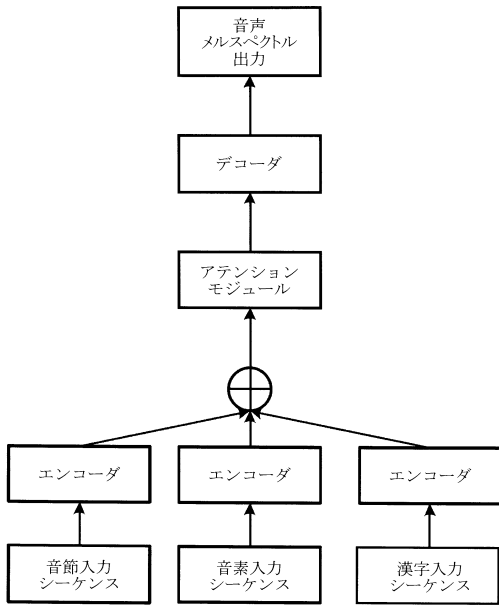
【 図 1 】



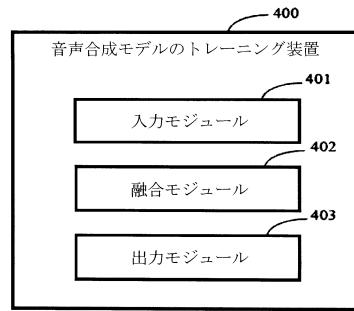
【 図 2 】



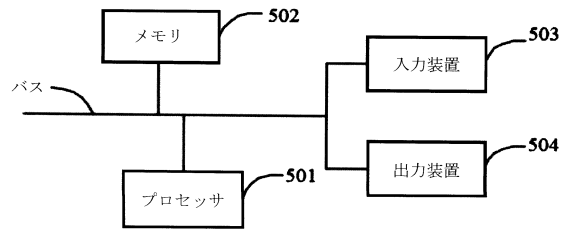
【図3】



【図4】



【図5】



## フロントページの続き

(74)代理人 100201466

弁理士 竹内 邦彦

(72)発明者 チェン, ジーペン

中華人民共和国 100085 ベイジン ハイディエン ディストリクト シャンディ 10  
イーエイチ ストリート パイドゥ キャンパス ナンバー 10 3 / フロア

(72)発明者 バイ, ジンフェン

中華人民共和国 100085 ベイジン ハイディエン ディストリクト シャンディ 10  
イーエイチ ストリート パイドゥ キャンパス ナンバー 10 3 / フロア

(72)発明者 ジア, レイ

中華人民共和国 100085 ベイジン ハイディエン ディストリクト シャンディ 10  
イーエイチ ストリート パイドゥ キャンパス ナンバー 10 3 / フロア

審査官 渡部 幸和

(56)参考文献 LU Chunhui et al., Self-attention Based Prosodic Boundary Prediction for Chinese Speech Synthesis, ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2019年04月17日

DING Chuang et al., Automatic prosody prediction for Chinese speech synthesis using BLASTM-RNN and embedding features, 2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), IEEE, 2016年02月11日

YASUDA Yusuke et al., Investigation of Enhanced Tacotron Text-to-Speech Synthesis Systems with Self-attention for Pitch Accent Language, ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2019年04月17日

(58)調査した分野(Int.Cl., D B名)

G10L 13/00 - 13/10