

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局

(43) 国際公開日
2021年1月14日(14.01.2021)



(10) 国際公開番号
WO 2021/006117 A1

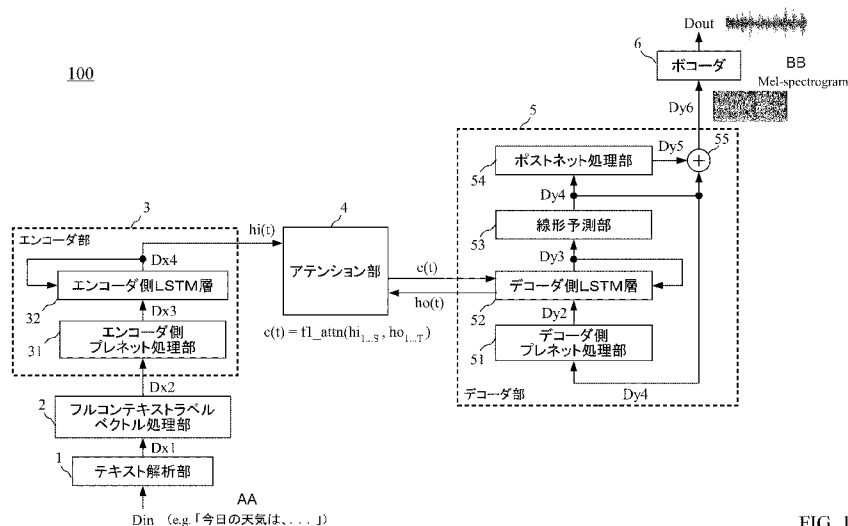
- (51) 国際特許分類:
G10L 13/08 (2013.01) G10L 25/30 (2013.01)
G10L 13/10 (2013.01)
- (21) 国際出願番号: PCT/JP2020/025682
- (22) 国際出願日: 2020年6月30日(30.06.2020)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:
特願 2019-125726 2019年7月5日(05.07.2019) JP
特願 2019-200440 2019年11月5日(05.11.2019) JP
- (71) 出願人: 国立研究開発法人情報通信研究機構 (NATIONAL INSTITUTE OF INFORMATION AND COMMUNICATIONS TECHNOLOGY) [JP/JP]; 〒1848795

東京都小金井市貫井北町4-2-1 Tokyo (JP).

- (72) 発明者: 岡本 拓磨 (OKAMOTO Takuma); 〒1848795 東京都小金井市貫井北町4-2-1 国立研究開発法人情報通信研究機構内 Tokyo (JP). 戸田 智基 (TODA Tomoki); 〒4801103 愛知県長久手市岩作三ヶ峯38-1 Aichi (JP). 志賀 芳則 (SHIGA Yoshinori); 〒1848795 東京都小金井市貫井北町4-2-1 国立研究開発法人情報通信研究機構内 Tokyo (JP). 河井 恒 (KAWAI Hisashi); 〒1848795 東京都小金井市貫井北町4-2-1 国立研究開発法人情報通信研究機構内 Tokyo (JP).
- (74) 代理人: 中西 健, 外 (NAKANISHI Ken et al.); 〒5300001 大阪府大阪市北区梅田二丁目5番6号 桜橋八千代ビル3F りのわ国際特許事務所 Osaka (JP).

(54) Title: VOICE SYNTHESIS PROCESSING DEVICE, VOICE SYNTHESIS PROCESSING METHOD, AND PROGRAM

(54) 発明の名称: 音声合成処理装置、音声合成処理方法、および、プログラム



- 1 Text analysis unit
- 2 Full-context label vector processing unit
- 3 Encoder unit
- 4 Attention unit
- 5 Decoder unit
- 6 Vocoder
- 31 Encoder-side pre-net processing unit
- 32 Encoder-side LSTM layer
- 51 Decoder-side pre-net processing unit
- 52 Decoder-side LSTM layer
- 53 Linear prediction unit
- 54 Post-net processing unit
- AA e.g. "The weather today is ..."
- BB Mel-spectrogram

FIG. 1

(57) Abstract: The present invention enables a voice synthesis processing device which performs learning and optimization and enables a high-quality voice synthesis process by means of a neural network model that can turn a processing target language into an arbitrarily defined language and that is for text-to-speech synthesis using the sequence-to-sequence method. In a voice synthesis processing device (100), a text analysis process is performed in accordance with a processing target language, optimized full-context label data that is suitable for processing with a neural network model is acquired from full-context label data acquired in the text analysis process, and the acquired optimized full-context label data is used to perform a process, whereby a highly accurate voice synthesis process for an arbitrarily defined processing target language can be performed.

WO 2021/006117 A1

- (81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, WS, ZA, ZM, ZW.
- (84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類 :

- 一 国際調査報告 (条約第21条(3))

(57) 要約 : 処理対象言語を任意の言語にできる、sequence-to-sequence方式を用いたテキスト音声合成用のニューラルネットワークのモデルにより、学習・最適化を行い、高品質な音声合成処理を実現する音声合成処理装置を実現する。音声合成処理装置(100)では、処理対象言語に応じたテキスト解析処理を行い、当該テキスト解析処理で取得されたフルコンテキストラベルデータから、ニューラルネットワークのモデルで処理するのに適した最適化フルコンテキストラベルデータを取得し、取得した最適化フルコンテキストラベルデータを用いて処理を行うことで、任意の処理対象言語について、高精度な音声合成処理を行うことができる。

明 細 書

発明の名称：

音声合成処理装置、音声合成処理方法、および、プログラム

技術分野

[0001] 本発明は、音声合成処理技術に関する。特に、テキストを音声に変換するテキスト音声合成（TTS：text-to-speech）技術に関する。

背景技術

[0002] テキストから自然な音声を合成するテキスト音声合成（TTS）技術において、近年、ニューラルネットワークの導入により高品質な音声合成が可能となっている。このようなテキスト音声合成技術を用いたシステムでは、英語音声を合成する場合、音素継続長と音響モデルとを同時に学習・最適化するsequence-to-sequence方式を用いたテキスト音声合成技術により、英語テキストからメルスペクトログラムを推定し、推定したメルスペクトログラムから、ニューラルボコーダにより音声波形を取得する。このように処理することで、上記テキスト音声合成技術を用いたシステムでは、処理対象言語が英語である場合、人間の音声と同等の品質の音声合成が可能となる（例えば、非特許文献1を参照）。

先行技術文献

非特許文献

[0003] 非特許文献1：Jonathan Shen, R Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerry-Ryan, R. A. Saurous, Y. Agiomyriannakis, and Y. Wu, "Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions," Proc. ICASSP, Apr. 2018, pp. 4779-4783.

発明の概要

発明が解決しようとする課題

[0004] しかしながら、上記のsequence-to-sequence方式を用いたテキスト音声合

成を日本語に適用するのは困難である。日本語は、漢字を使う言語であり、漢字の数が膨大であり、また、漢字の読みのバリエーションも多いので、日本語テキストを、sequence-to-sequence方式を用いたテキスト音声合成のモデルに、直接入力し、処理言語を英語としたときと同様に、当該モデルの学習・最適化を行うことは困難である。

[0005] そこで本発明は、上記課題に鑑み、日本語等の英語以外の言語を処理対象言語とする場合においても（処理対象言語を任意の言語にできる）、sequence-to-sequence方式を用いたテキスト音声合成用のニューラルネットワークのモデルにより、学習・最適化を行い、高品質な音声合成処理を実現する音声合成処理装置、音声合成処理方法、および、プログラムを実現することを目的とする。

課題を解決するための手段

[0006] 上記課題を解決するための第1の発明は、任意の言語を処理対象言語とし、エンコーダ・デコーダ方式のニューラルネットワークを用いて音声合成処理を実行する音声合成処理装置であって、テキスト解析部と、フルコンテキストラベルベクトル処理部と、エンコーダ部と、デコーダ部と、を備える。

[0007] テキスト解析部は、処理対象言語のテキストデータに対してテキスト解析処理を実行し、コンテキストラベルデータを取得する。

[0008] フルコンテキストラベルベクトル処理部は、テキスト解析部により取得されたコンテキストラベルデータから、コンテキストラベルデータを取得する処理において処理対象とされた音素である単独音素についてのコンテキストラベルを取得することで、ニューラルネットワークの学習処理に適した最適化フルコンテキストラベルデータを取得する。

[0009] エンコーダ部は、最適化フルコンテキストラベルデータに基づいて、ニューラルネットワークのエンコード処理を実行することで、隠れ状態データを取得する。

[0010] デコーダ部は、隠れ状態データに基づいて、ニューラルネットワークのデコード処理を実行することで、最適化フルコンテキストラベルデータに対応

する音響特徴量データを取得する。

[0011] ボコーダは、デコーダ部により取得された音響特徴量から音声波形データを取得する。

[0012] この音声合成処理装置では、ニューラルネットワークのモデルで処理するのに適した最適化フルコンテキストラベルデータを用いて、ニューラルネットワークによる処理（学習処理、予測処理）を実行するため、高精度な音声合成処理を実行することができる。つまり、この音声合成処理装置では、従来技術とは異なり、処理対象の音素に先行する、あるいは、後続する音素についてのデータを含まないコンテキストラベルデータを最適化フルコンテキストラベルデータとして取得し、取得した最適化フルコンテキストラベルデータにより、ニューラルネットワークのモデルの処理を行う。ニューラルネットワーク（特に、sequence-to-sequence方式のニューラルネットワーク）では、時系列のデータを用いた処理が実行されるので、従来の音声合成処理に用いるコンテキストラベルデータに含める必要があった、処理対象の音素に先行する、あるいは、後続するデータは、ニューラルネットワークのモデルの処理において冗長となり、処理効率を落とす原因となる。この音声合成処理装置100では、最適化フルコンテキストラベルデータ（単独音素についてのデータから構成されるコンテキストラベルデータ）を用いるので、ニューラルネットワークのモデルの処理が非常に効果的に実行できる。その結果、この音声合成処理装置では、高精度の音声合成処理を実行できる。

[0013] また、この音声合成処理装置では、処理対象言語に応じたテキスト解析処理を行い、当該テキスト解析処理で取得されたフルコンテキストラベルデータから、ニューラルネットワーク（例えば、sequence-to-sequence方式を用いたニューラルネットワーク）のモデルで処理するのに適した最適化フルコンテキストラベルデータを取得し、取得した最適化フルコンテキストラベルデータを用いて処理を行うことで、任意の処理対象言語について、高精度な音声合成処理を行うことができる。

[0014] したがって、この音声合成処理装置では、日本語等の英語以外の言語を処

理対象言語とする場合においても（処理対象言語を任意の言語にできる）、例えば、sequence-to-sequence方式を用いたテキスト音声合成用のニューラルネットワークのモデルにより、学習・最適化を行い、高品質な音声合成処理を実現することができる。

[0015] なお、「単独音素」とは、テキスト解析処理においてコンテキストラベルデータを取得するときに、処理対象とした音素のことをいう。

[0016] また、「最適化」とは、厳密な意味での最適化の他に、所定の誤差範囲を許容する範囲内に収めることを含む概念である。

[0017] 第2の発明は、第1の発明であって、音響特徴量は、メルスペクトログラムのデータである。

[0018] これにより、この音声合成処理装置では、入力されたテキストに対応するメルスペクトログラムのデータにより、音声合成処理を実行できる。

[0019] 第3の発明は、第1または第2の発明であって、ボコーダは、ニューラルネットワークのモデルを用いた処理を実行することで、音響特徴量から音声波形データを取得する。

[0020] これにより、この音声合成処理装置では、ニューラルネットワーク処理ができるボコーダを用いて、音声合成処理を実行できる。

[0021] 第4の発明は、第3の発明であって、ボコーダは、可逆変換ネットワークにより構成されたニューラルネットワークのモデルを用いた処理を実行することで、音響特徴量から音声波形データを取得する。

[0022] この音声合成処理装置では、ボコーダが、可逆変換ネットワークにより構成されたニューラルネットワークのモデルを用いた処理を行うので、ボコーダの構成をシンプルにできる。その結果、この音声合成処理装置では、ボコーダでの処理を高速化でき、音声合成処理をリアルタイムで実行できる。

[0023] 第5の発明は、第1から第4のいずれかの発明であって、音素単位のコンテキストラベルデータから音素継続長を推定する音素継続長推定部をさらに備える。

[0024] フルコンテキストラベルベクトル処理部は、音素継続長推定部により推定

された音素継続長である推定音素継続長に対応する期間において、当該推定音素継続長に対応する音素の最適化フルコンテキストラベルデータを継続してエンコーダ部へ出力する。

[0025] この音声合成処理装置では、エンコーダ部への入力データ（最適化フルコンテキストラベルデータ）を、音素継続長推定部により取得（推定）した音素ごとの音素継続長に基づいて、引き延ばす処理（音素 $p h_k$ の音素継続長 $d u r (p h_k)$ に相当する期間、音素 $p h_k$ の最適化フルコンテキストラベルデータを、繰り返しエンコーダ部 3 へ入力する処理）を実行する。つまり、この音声合成処理装置では、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理を実行して取得した音素継続長を用いて予測処理を実行するので、注意機構予測が失敗することに起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することはない。

[0026] すなわち、この音声合成処理装置では、（1）音素継続長については、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理（音素継続長推定部による処理）により取得し、（2）音響特徴量については、sequence-to-sequence方式を用いたニューラルネットワークのモデルで処理することにより取得する。

[0027] したがって、この音声合成処理装置では、注意機構予測が失敗することに起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することを適切に防止するとともに、高精度な音声合成処理を実行することができる。

[0028] 第6の発明は、任意の言語を処理対象言語とし、エンコーダ・デコーダ方式のニューラルネットワークを用いて音声合成処理を実行する音声合成処理方法であって、テキスト解析ステップと、フルコンテキストラベルベクトル処理ステップと、エンコード処理ステップと、デコード処理ステップと、ボコーダ処理ステップと、を備える。

[0029] テキスト解析ステップは、処理対象言語のテキストデータに対してテキス

- ト解析処理を実行し、コンテキストラベルデータを取得する。
- [0030] フルコンテキストラベルベクトル処理ステップは、テキスト解析ステップにより取得されたコンテキストラベルデータから、コンテキストラベルデータを取得する処理において処理対象とされた音素である単独音素についてのコンテキストラベルを取得することで、ニューラルネットワークの学習処理に適した最適化フルコンテキストラベルデータを取得する。
- [0031] エンコード処理ステップは、最適化フルコンテキストラベルデータに基づいて、ニューラルネットワークのエンコード処理を実行することで、隠れ状態データを取得する。
- [0032] デコード処理ステップは、隠れ状態データに基づいて、ニューラルネットワークのデコード処理を実行することで、最適化フルコンテキストラベルデータに対応する音響特徴量データを取得する。
- [0033] ボコーダ処理ステップは、デコード処理ステップにより取得された音響特徴量から音声波形データを取得する。
- [0034] これにより、第1の発明と同様の効果を奏する音声合成処理方法を実現することができる。
- [0035] 第7の発明は、第6の発明である音声合成処理方法をコンピュータに実行させるためのプログラムである。
- [0036] これにより、第1の発明と同様の効果を奏する音声合成処理方法をコンピュータに実行させるためのプログラムを実現することができる。
- [0037] 第8の発明は、任意の言語を処理対象言語とし、エンコーダ・デコーダ方式のニューラルネットワークを用いて音声合成処理を実行する音声合成処理装置であって、テキスト解析部と、フルコンテキストラベルベクトル処理部と、エンコーダ部と、音素継続長推定部と、強制アテンション部と、内分処理部と、コンテキスト算出部と、デコーダ部と、ボコーダと、を備える。
- [0038] テキスト解析部は、処理対象言語のテキストデータに対してテキスト解析処理を実行し、コンテキストラベルデータを取得する。
- [0039] フルコンテキストラベルベクトル処理部は、テキスト解析部により取得さ

れたコンテキストラベルデータから、コンテキストラベルデータを取得する処理において処理対象とされた音素である単独音素についてのコンテキストラベルを取得することで、ニューラルネットワークの学習処理に適した最適化フルコンテキストラベルデータを取得する。

- [0040] エンコーダ部は、最適化フルコンテキストラベルデータに基づいて、ニューラルネットワークのエンコード処理を実行することで、隠れ状態データを取得する。
- [0041] 音素継続長推定部は、音素単位のコンテキストラベルデータから音素継続長を推定する。
- [0042] 強制アテンション部は、音素継続長推定部により推定された音素継続長に基づいて、第1重み付け係数データを取得する。
- [0043] アテンション部は、エンコーダ部により取得された隠れ状態データに基づいて、第2重み付け係数データを取得する。
- [0044] 内分処理部は、第1重み付け係数データと第2重み付け係数データとに対して内分処理を行うことで、合成重み付け係数データを取得する。
- [0045] コンテキスト算出部は、合成重み付け係数データにより、エンコーダ部により取得された隠れ状態データに対して重み付け合成処理を実行することで、コンテキスト状態データを取得する。
- [0046] デコーダ部は、コンテキスト状態データに基づいて、ニューラルネットワークのデコード処理を実行することで、最適化フルコンテキストラベルデータに対応する音響特徴量データを取得する。
- [0047] ボコーダは、デコーダ部により取得された音響特徴量から音声波形データを取得する。
- [0048] この音声合成処理装置では、音素継続長については、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理（音素継続長推定部による処理）により取得した音素継続長を用いて処理することで、音素継続長の予測精度を保証する。つまり、この音声合成処理装置では、安定して音素継続長を適切に推定することができる、隠れマ

ルコフモデル等のモデルを用いた推定処理（音素継続長推定部による処理）により取得した音素継続長を用いて強制アテンション部により取得した重み付け係数データと、アテンション部により取得された重み付け係数データを適度に合成した重み付け係数データにより生成したコンテキスト状態データを用いて予測処理を実行する。したがって、この音声合成処理装置では、注意機構の予測が失敗する場合（アテンション部により適切な重み付け係数データが取得できない場合）であっても、強制アテンション部により取得した重み付け係数データによる重み分の重み付け係数データが取得できるため、注意機構の予測の失敗が音声合成処理に影響を及ぼさないようにできる。

[0049] さらに、この音声合成処理装置では、音響特徴量については、sequence-to-sequence方式を用いたニューラルネットワークのモデルで処理することにより取得できるので、高精度な音響特徴量の予測処理が実現できる。

[0050] したがって、この音声合成処理装置では、注意機構予測が失敗すること起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することを適切に防止するとともに、高精度な音声合成処理を実行することができる。

[0051] なお、この音声合成処理装置において、内分処理を実行するときの内分比は、固定値であってもよいし、動的に変化する（更新される）値であってもよい。

発明の効果

[0052] 本発明によれば、日本語等の英語以外の言語を処理対象言語とする場合においても（処理対象言語を任意の言語にできる）、sequence-to-sequence方式を用いたテキスト音声合成用のニューラルネットワークのモデルにより、学習・最適化を行い、高品質な音声合成処理を実現する音声合成処理装置、音声合成処理方法、および、プログラムを実現することができる。

図面の簡単な説明

[0053] [図1]第1実施形態に係る音声合成処理装置100の概略構成図。

[図2]処理対象言語を日本語とした場合のテキスト解析処理により取得される

フルコンテキストラベルデータに含まれる情報（パラメータ）（一例）を示す図。

[図3]最適化フルコンテキストラベルデータに含まれる情報（パラメータ）（一例）を示す図。

[図4]第1実施形態の第1変形例の音声合成処理装置のボコーダ6の概略構成を示す図。

[図5]第1実施形態の第1変形例の音声合成処理装置のボコーダ6の概略構成を示す図。

[図6]第1実施形態の第1変形例の音声合成処理装置によりTTS処理（処理対象言語：日本語）実行し、取得した音声波形データのメルスペクトログラム（予測データ）と、入力テキストの実際の音声波形データのメルスペクトログラム（オリジナルデータ）とを示す図。

[図7]第2実施形態に係る音声合成処理装置200の概略構成図

[図8]推定された音素継続長に基づいて、エンコーダ部3に入力するデータ $D \times 2$ を生成する処理を説明するための図。

[図9]第3実施形態に係る音声合成処理装置300の概略構成図。

[図10]アテンション部4Aにより取得された重み付け係数データ $w_a(t)$ と、強制アテンション部8により取得された重み付け係数データ $w_f(t)$ とから取得した合成重み付け係数データ $w(t)$ を用いてコンテキスト状態データ $c(t)$ を取得する処理について説明するための図。

[図11]アテンション部4Aにより取得された重み付け係数データ $w_a(t)$ と、強制アテンション部8により取得された重み付け係数データ $w_f(t)$ とから取得した合成重み付け係数データ $w(t)$ を用いてコンテキスト状態データ $c(t)$ を取得する処理について説明するための図（時刻 t_2 の処理）。

[図12]アテンション部4Aにより取得された重み付け係数データ $w_a(t)$ と、強制アテンション部8により取得された重み付け係数データ $w_f(t)$ とから取得した合成重み付け係数データ $w(t)$ を用いてコンテキスト

状態データ $c(t)$ を取得する処理について説明するための図（時刻 t_3 の処理）。

[図13]時刻 t_2 における処理で、注意機構の予測が失敗している場合を説明するための図。

[図14]本発明に係る音声合成処理装置を実現するコンピュータのハードウェア構成を示すブロック図。

発明を実施するための形態

[0054] [第1実施形態]

第1実施形態について、図面を参照しながら、以下説明する。

[0055] <1. 1：音声合成処理装置の構成>

図1は、第1実施形態に係る音声合成処理装置100の概略構成図である。

[0056] 音声合成処理装置100は、図1に示すように、テキスト解析部1と、フルコンテキストラベルベクトル処理部2と、エンコーダ部3と、アテンション部4と、デコーダ部5と、ボコーダ6とを備える。

[0057] テキスト解析部1は、処理対象言語のテキストデータ D_{in} を入力とし、入力されたテキストデータ D_{in} に対して、テキスト解析処理を実行し、様々な言語情報からなるコンテキストを含む音素ラベルであるコンテキストラベルの系列を取得する。なお、日本語のように、アクセントやピッチによって、同じ文字（例えば、漢字）であっても、発音されたときの音声波形が異なる言語では、当該音素（処理対象の音素）の前後の音素についての言語情報も、コンテキストラベルに含める必要がある。テキスト解析部1は、上記のように、テキストが発音されたときの音声波形を特定するためのコンテキストラベル（処理対象言語によって必要となる先行する音素、および／または、後続する音素のデータを含めたコンテキストラベル）をフルコンテキストラベルデータ D_x として、フルコンテキストラベルベクトル処理部2に出力する。

[0058] フルコンテキストラベルベクトル処理部2は、テキスト解析部1から出力

されるデータ $D \times 1$ （フルコンテキストラベルのデータ）を入力する。フルコンテキストラベルベクトル処理部 2 は、入力されたフルコンテキストラベルデータ $D \times 1$ から、sequence-to-sequence方式のニューラルネットワークのモデルの学習処理に適したフルコンテキストラベルデータを取得するためのフルコンテキストラベルベクトル処理を実行する。そして、フルコンテキストラベルベクトル処理部 2 は、フルコンテキストラベルベクトル処理により取得したデータをデータ $D \times 2$ （最適化フルコンテキストラベルデータ $D \times 2$ ）として、エンコーダ部 3 のエンコーダ側プレネット処理部 3 1 に出力する。

[0059] エンコーダ部 3 は、図 1 に示すように、エンコーダ側プレネット処理部 3 1 と、エンコーダ側 LSTM 層 3 2（LSTM: Long short-term memory）とを備える。

[0060] エンコーダ側プレネット処理部 3 1 は、フルコンテキストラベルベクトル処理部 2 から出力されるデータ $D \times 2$ を入力する。エンコーダ側プレネット処理部 3 1 は、入力したデータ $D \times 2$ に対して、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU 関数（ReLU: Rectified Linear Unit）による処理）を実行し、エンコーダ側 LSTM 層 3 2 に入力可能なデータを取得する。そして、エンコーダ側プレネット処理部 3 1 は、上記処理（プレネット処理）により取得したデータをデータ $D \times 3$ としてエンコーダ側 LSTM 層 3 2 に出力する。

[0061] エンコーダ側 LSTM 層 3 2 は、リカーレントニューラルネットワークの隠れ層（LSTM 層）に対応する層であり、エンコーダ側プレネット処理部 3 1 から、現時刻 t において出力されるデータ $D \times 3$ （これをデータ $D \times 3(t)$ と表記する）と、1 つ前の時間ステップにおいて、エンコーダ側 LSTM 層 3 2 から出力されたデータ $D \times 4$ （これをデータ $D \times 4(t-1)$ と表記する）とを入力する。エンコーダ側 LSTM 層 3 2 は、入力されたデータ $D \times 3(t)$ 、データ $D \times 4(t-1)$ に対して、LSTM 層による処理

を実行し、処理後のデータをデータ $D \times 4$ (データ $D \times 4 (t)$) としてアテンション部4に出力する。

[0062] アテンション部4は、エンコーダ部3から出力されるデータ $D \times 4$ と、デコーダ部5のデコーダ側LSTM層52から出力されるデータ h_o (出力側隠れ状態データ h_o) とを入力する。アテンション部4は、エンコーダ部3から出力されるデータ $D \times 4$ 、すなわち、入力側隠れ状態データ (これをデータ h_i という。また、時刻 t の入力側隠れ状態データをデータ $h_i (t)$ と表記する。) を所定の時間ステップ分記憶保持する。時間ステップ $t = 1$ から $t = S$ (S : 自然数) の期間において、エンコーダ部3により取得され、アテンション部4に出力されたデータ $D \times 4 (= h_i)$ の集合を、 $h_{i,1 \dots s}$ と表記する。つまり、アテンション部4は、下記に相当するデータ $h_{i,1 \dots s}$ を記憶保持する。

$$h_{i,1 \dots s} = \{D \times 4 (1), D \times 4 (2), \dots, D \times 4 (S)\}$$

また、アテンション部4は、デコーダ部5のデコーダ側LSTM層52から出力されるデータ $D_y 3$ 、すなわち、出力側隠れ状態データ (これをデータ h_o という) を所定の時間ステップ分記憶保持する。時間ステップ $t = 1$ から $t = T$ (T : 自然数) の期間において、デコーダ側LSTM層52により取得され、アテンション部4に出力されたデータ $D_y 3 (= h_o)$ の集合を、 $h_{o,1 \dots T}$ と表記する。つまり、アテンション部4は、下記に相当するデータ $h_{o,1 \dots T}$ を記憶保持する。

$$h_{o,1 \dots T} = \{D_y 3 (1), D_y 3 (2), \dots, D_y 3 (T)\}$$

そして、アテンション部4は、入力側隠れ状態データの集合データ $h_{i,1 \dots s}$ と、出力側隠れ状態データの集合データ $h_{o,1 \dots T}$ と、に基づいて、例えば、

$$c(t) = f1_attn(h_{i,1 \dots s}, h_{o,1 \dots T})$$

$f1_attn()$: コンテキスト状態データを取得する関数

に相当する処理を実行して、現時刻 t のコンテキスト状態データ $c(t)$ を取得する。そして、アテンション部4は、取得したコンテキスト状態データ

$c(t)$ をデコーダ側 LSTM 層 52 に出力する。

デコーダ部 5 は、図 1 に示すように、デコーダ側プレネット処理部 51 と、デコーダ側 LSTM 層 52 と、線形予測部 53 と、ポストネット処理部 54 と、加算器 55 と、を備える。

[0063] デコーダ側プレネット処理部 51 は、線形予測部 53 から出力される、1 時間ステップ前のデータ Dy_4 (これを $Dy_4(t-1)$ という) を入力する。デコーダ側プレネット処理部 51 は、例えば、複数層 (例えば、2 層) の全結合層を有しており、データの正規化処理 (例えば、線形予測部 53 から出力されるデータ (ベクトルデータ) の次元数が $2N$ であり、デコーダ側 LSTM 層に入力されるデータ (ベクトルデータ) の次元数が N である場合、データの次元数を N にするように、例えば、ドロップアウト処理を行うことを含む)、活性化関数による処理 (例えば、ReLU 関数 (ReLU: Rectified Linear Unit) による処理) を実行し、デコーダ側 LSTM 層 52 に入力可能なデータを取得する。そして、デコーダ側プレネット処理部 51 は、上記処理 (プレネット処理) により取得したデータをデータ Dy_2 としてデコーダ側 LSTM 層 52 に出力する。

[0064] デコーダ側 LSTM 層 52 は、リカーレントニューラルネットワークの隠れ層 (LSTM 層) に対応する層である。デコーダ側 LSTM 層 52 は、デコーダ側プレネット処理部 51 から、現時刻 t において出力されるデータ Dy_2 (これをデータ $Dy_2(t)$ と表記する) と、1 つ前の時間ステップにおいて、デコーダ側 LSTM 層 52 から出力されたデータ Dy_3 (これをデータ $Dy_3(t-1)$ と表記する) と、アテンション部 4 から出力される時刻 t のコンテキスト状態データ $c(t)$ とを入力する。

[0065] デコーダ側 LSTM 層 52 は、入力されたデータ $Dy_2(t)$ 、データ $Dy_3(t-1)$ 、および、コンテキスト状態データ $c(t)$ を用いて、LSTM 層による処理を実行し、処理後のデータをデータ Dy_3 (データ $Dy_3(t)$) として線形予測部 53 に出力する。また、デコーダ側 LSTM 層 52 は、データ $Dy_3(t)$ 、すなわち、時刻 t の出力側隠れ状態データ h_o

(t) をアテンション部4に出力する。

- [0066] 線形予測部53は、デコーダ側LSTM層52から出力されるデータDy3を入力する。線形予測部53は、所定の期間（例えば、メルスペクトログラムを取得するための1フレーム期間に相当する期間）内に、デコーダ側LSTM層52から出力されるデータDy3（複数のデータDy3）を記憶保持し、当該複数のデータDy3を用いて線形変換することで、所定期間におけるメルスペクトログラムの予測データDy4を取得する。そして、線形予測部53は、取得したデータDy4をポストネット処理部54、加算器55、および、デコーダ側プレネット処理部51に出力する。
- [0067] ポストネット処理部54は、例えば、複数層（例えば、5層）のコンボリユーション層を有しており、コンボリユーション処理（コンボリユーションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理やtanh関数による処理）を実行し、予測データ（予測メルスペクトログラム）の残差データ（residual）を取得し、取得した残差データをデータDy5として加算器55に出力する。
- [0068] 加算器55は、線形予測部53から出力される予測データDy4（予測メルスペクトログラムのデータ）と、ポストネット処理部54から出力される残差データDy5（予測メルスペクトログラムの残差データ）とを入力する。加算器55は、予測データDy4（予測メルスペクトログラムのデータ）と、残差データDy5（予測メルスペクトログラムの残差データ）とに対して加算処理を実行し、加算処理後のデータ（予測メルスペクトログラムのデータ）をデータDy6としてボコーダ6に出力する。
- [0069] ボコーダ6は、音響特徴量のデータを入力とし、入力された音響特徴量のデータから、当該音響特徴量に対応する音声信号波形を出力する。本実施形態において、ボコーダ6は、ニューラルネットワークによるモデルを用いたボコーダを採用する。ボコーダ6は、入力される音響特徴量を、メルスペクトログラムのデータとし、出力を当該メルスペクトログラムに対応する音声

信号波形とする。ボコーダ6は、学習時において、メルスペクトログラムと、当該メルスペクトログラムにより実現される音声信号波形（教師データ）として、ニューラルネットワークのモデルを学習させ、当該ニューラルネットワークのパラメータの最適化パラメータを取得することで、当該ニューラルネットワークのモデルを最適化する処理を行う。そして、ボコーダ6は、予測時において、最適化したニューラルネットワークのモデルを用いて、処理を行うことで、入力されるメルスペクトログラムのデータ（例えば、デコーダ部5から出力されるデータD_{y6}）から、当該メルスペクトログラムに対応する音声信号波形を予測し、予測した音声信号波形のデータをデータD_{out}として出力する。

[0070] <1. 2：音声合成処理装置の動作>

以上のように構成された音声合成処理装置100の動作について以下説明する。

[0071] 以下では、音声合成処理装置100の動作を、（1）学習処理（学習時の処理）と、（2）予測処理（予測時の処理）とに分けて説明する。

[0072] （1. 2. 1：学習処理）

まず、音声合成処理装置100による学習処理について、説明する。なお、説明便宜のため、処理対象言語を日本語として、以下、説明する。

[0073] 処理対象言語である日本語のテキストデータD_{in}をテキスト解析部1に入力する。また、当該テキストデータD_{in}に対応するメルスペクトログラム（音響特徴量）のデータを教師データとして用意する。

[0074] テキスト解析部1は、入力されたテキストデータD_{in}に対して、テキスト解析処理を実行し、様々な言語情報からなるコンテキストを含む音素ラベルであるコンテキストラベルの系列を取得する。

[0075] 日本語は、アクセントやピッチによって、同じ文字（例えば、漢字）であっても、発音されたときの音声波形が異なる言語であるので、当該音素（処理対象の音素）の前後の音素についての言語情報も、コンテキストラベルに含める必要がある。テキスト解析部1は、処理対象を日本語とする場合、テ

キストデータ D_{in} に対して、日本語用のテキスト解析処理を実行し、テキストが発音されたときの音声波形を特定するためのパラメータについて、必要に応じて、(1) 当該音素のみのデータ、(2) 先行する音素、および／または、後続する音素についてのデータを取得し、取得したデータをまとめてフルコンテキストラベルデータを取得する。

[0076] 図2は、処理対象言語を日本語とした場合のテキスト解析処理により取得されるフルコンテキストラベルデータに含まれる情報(パラメータ)(一例)を示す図である。

[0077] 図2に示す場合では、フルコンテキストラベルデータの各パラメータは、図2の「概要」に示した内容を特定するためのデータであり、図2の表に示した次元数、音素数分のデータである。

[0078] 図2に示すように、テキスト解析部1は、図2の表の全てのパラメータのデータをまとめて、フルコンテキストラベルデータ(ベクトルのデータ)として、取得する。図2の場合、フルコンテキストラベルデータは、478次元のベクトルデータとなる。

[0079] 上記のようにして取得されたフルコンテキストラベルデータ $D \times 1$ は、テキスト解析部1からフルコンテキストラベルベクトル処理部2に出力される。

[0080] フルコンテキストラベルベクトル処理部2は、入力されたフルコンテキストラベルデータ $D \times 1$ から、sequence-to-sequence方式のニューラルネットワークのモデルの学習処理に適したフルコンテキストラベルデータを取得するためのフルコンテキストラベルベクトル処理を実行する。具体的には、フルコンテキストラベルベクトル処理部2は、先行する音素についてのパラメータ(データ)、後続する音素についてのパラメータ(データ)を削除することで、最適化フルコンテキストラベルデータ $D \times 2$ を取得する。例えば、フルコンテキストラベルデータ $D \times 1$ が図2に示すパラメータを含むデータである場合、先行する音素についてのパラメータ(データ)、後続する音素についてのパラメータ(データ)を削除することで、最適化フルコンテキス

トラベルデータ $D \times 2$ を取得する。

[0081] 図3は、上記のようにして取得した最適化フルコンテキストトラベルデータに含まれる情報（パラメータ）（一例）を示す図である。

[0082] 図3の場合、最適化フルコンテキストトラベルデータは、130次元のベクトルデータとなり、478次元のベクトルデータであるフルコンテキストトラベルデータ $D \times 1$ と比べると、次元数が著しく低減されていることが分かる。

[0083] 音声合成処理装置100で用いられているニューラルネットワークのモデルが、sequence-to-sequence方式のニューラルネットワーク（リカーレントニューラルネットワーク）のモデルであり、エンコーダ側LSTM層32、デコーダ側LSTM層52を有しているので、入力されるデータ列について、時系列の関係を考慮した学習処理、予測処理ができるため、従来技術で必要とされていた先行する音素、後続する音素のデータは、冗長となり、学習処理の効率、予測処理の精度を悪化させる原因となる。そのため、音声合成処理装置100では、上記のように、当該音素についてのパラメータ（データ）のみを残して取得した最適化フルコンテキストトラベルデータ $D \times 2$ を取得し、取得した最適化フルコンテキストトラベルデータ $D \times 2$ を用いて、学習処理、予測処理を行うことで、高速かつ高精度に処理を実行することができる。

[0084] 上記により取得されたデータ $D \times 2$ （最適化フルコンテキストトラベルデータ $D \times 2$ ）は、フルコンテキストトラベルベクトル処理部2からのエンコーダ部3のエンコーダ側プレネット処理部31に出力される。

[0085] エンコーダ側プレネット処理部31は、フルコンテキストトラベルベクトル処理部2から入力したデータ $D \times 2$ に対して、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理）を実行し、エンコーダ側LSTM層32に入力可能なデータを取得する。そして、エンコーダ側プレネット処理部3

1は、上記処理（プレネット処理）により取得したデータをデータ $D \times 3$ としてエンコーダ側LSTM層32に出力する。

[0086] エンコーダ側LSTM層32は、エンコーダ側プレネット処理部31から、現時刻 t において出力されるデータ $D \times 3(t)$ と、1つ前の時間ステップにおいて、エンコーダ側LSTM層32から出力されたデータ $D \times 4(t-1)$ とを入力する。そして、エンコーダ側LSTM層32は、入力されたデータ $D \times 3(t)$ 、データ $D \times 4(t-1)$ に対して、LSTM層による処理を実行し、処理後のデータをデータ $D \times 4$ （データ $D \times 4(t)$ ）としてアテンション部4に出力する。

[0087] アテンション部4は、エンコーダ部3から出力されるデータ $D \times 4$ と、デコーダ部5のデコーダ側LSTM層52から出力されるデータ h_o （出力側隠れ状態データ h_o ）とを入力する。アテンション部4は、エンコーダ部3から出力されるデータ $D \times 4$ 、すなわち、入力側隠れ状態データ h_i を所定の時間ステップ分記憶保持する。例えば、アテンション部4は、時間ステップ $t=1$ から $t=S$ （ S ：自然数）の期間において、エンコーダ部3により取得され、アテンション部4に出力されたデータ $D \times 4(=h_i)$ の集合を、 $h_{i_{1..s}}(=\{D \times 4(1), D \times 4(2), \dots, D \times 4(S)\})$ として記憶保持する。

[0088] また、アテンション部4は、デコーダ部5のデコーダ側LSTM層52から出力されるデータ D_y3 、すなわち、出力側隠れ状態データ h_o を所定の時間ステップ分記憶保持する。例えば、アテンション部4は、時間ステップ $t=1$ から $t=T$ （ T ：自然数）の期間において、デコーダ側LSTM層52により取得され、アテンション部4に出力されたデータ $D_y3(=h_o)$ の集合を、 $h_{o_{1..T}}(=\{D_y3(1), D_y3(2), \dots, D_y3(T)\})$ として記憶保持する。

[0089] そして、アテンション部4は、入力側隠れ状態データの集合データ $h_{i_{1..s}}$ と、出力側隠れ状態データの集合データ $h_{o_{1..T}}$ と、に基づいて、例えば、

$$c(t) = f1_attn(h_{i_{1..s}}, h_{o_{1..T}})$$

$f1_attn()$: コンテキスト状態データを取得する関数

に相当する処理を実行して、現時刻 t のコンテキスト状態データ $c(t)$ を取得する。

[0090] そして、アテンション部4は、取得したコンテキスト状態データ $c(t)$ をデコーダ側LSTM層52に出力する。

デコーダ側プレネット処理部51は、線形予測部53から出力される、1時間ステップ前のデータ $Dy_4(t-1)$ を入力する。デコーダ側プレネット処理部51は、例えば、複数層（例えば、2層）の全結合層を有しており、データの正規化処理（例えば、線形予測部53から出力されるデータ（ベクトルデータ）の次元数が $2N$ であり、デコーダ側LSTM層に入力されるデータ（ベクトルデータ）の次元数が N である場合、データの次元数を N にするように、例えば、ドロップアウト処理を行うことを含む）、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理）を実行し、デコーダ側LSTM層52に入力可能なデータを取得する。そして、デコーダ側プレネット処理部51は、上記処理（プレネット処理）により取得したデータをデータ Dy_2 としてデコーダ側LSTM層52に出力する。

[0091] デコーダ側LSTM層52は、デコーダ側プレネット処理部51から、現時刻 t において出力されるデータ $Dy_2(t)$ と、1つ前の時間ステップにおいて、デコーダ側LSTM層52から出力されたデータ $Dy_3(t-1)$ と、アテンション部4から出力される時刻 t のコンテキスト状態データ $c(t)$ とを入力する。

[0092] デコーダ側LSTM層52は、入力されたデータ $Dy_2(t)$ 、データ $Dy_3(t-1)$ 、および、コンテキスト状態データ $c(t)$ を用いて、LSTM層による処理を実行し、処理後のデータをデータ $Dy_3(t)$ として線形予測部53に出力する。また、デコーダ側LSTM層52は、データ $Dy_3(t)$ 、すなわち、時刻 t の出力側隠れ状態データ $h_o(t)$ をアテンシ

ョン部4に出力する。

- [0093] 線形予測部53は、デコーダ側LSTM層52から出力されるデータDy3を入力する。線形予測部53は、所定の期間（例えば、メルスペクトログラムを取得するための1フレーム期間に相当する期間）内に、デコーダ側LSTM層52から出力されるデータDy3（複数のデータDy3）を記憶保持し、当該複数のデータDy3を用いて線形変換することで、所定期間におけるメルスペクトログラムの予測データDy4を取得する。そして、線形予測部53は、取得したデータDy4をポストネット処理部54、加算器55、および、デコーダ側プレネット処理部51に出力する。
- [0094] ポストネット処理部54は、例えば、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理やtanh関数による処理）を実行し、予測データ（予測メルスペクトログラム）の残差データ（residual）を取得し、取得した残差データをデータDy5として加算器55に出力する。
- [0095] 加算器55は、線形予測部53から出力される予測データDy4（予測メルスペクトログラムのデータ）と、ポストネット処理部54から出力される残差データDy5（予測メルスペクトログラムの残差データ）とを入力する。加算器55は、予測データDy4（予測メルスペクトログラムのデータ）と、残差データDy5（予測メルスペクトログラムの残差データ）とに対して加算処理を実行し、加算処理後のデータ（予測メルスペクトログラムのデータ）をデータDy6として出力する。
- [0096] そして、音声合成処理装置100では、上記のように取得されたデータDy6（予測メルスペクトログラムのデータ）と、テキストデータDinに対応するメルスペクトログラム（音響特徴量）の教師データ（正解のメルスペクトログラム）とを比較し、両者の差（比較結果）（例えば、差分ベクトルのノルムやユークリッド距離により表現する差）が小さくなるように、エンコーダ部3、デコーダ部5のニューラルネットワークのモデルのパラメータ

を更新する。音声合成処理装置100では、このパラメータ更新処理を繰り返し実行し、データDy6（予測メルスペクトログラムのデータ）と、テキストデータDinに対応するメルスペクトログラム（音響特徴量）の教師データ（正解のメルスペクトログラム）との差が十分小さくなる（所定の誤差範囲におさまる）、ニューラルネットワークのモデルのパラメータを最適化パラメータとして取得する。

[0097] 音声合成処理装置100では、上記のようにして取得した最適化パラメータに基づいて、エンコーダ部3、デコーダ部5のニューラルネットワークのモデルの各層に含まれるシナプス間の結合係数（重み係数）を設定することで、エンコーダ部3、デコーダ部5のニューラルネットワークのモデルを最適化モデル（学習済みモデル）とすることができる。

[0098] 以上により、音声合成処理装置100において、入力をテキストデータとし、出力をメルスペクトログラムとするニューラルネットワークの学習済みモデル（最適化モデル）を構築できる。

[0099] また、ボコーダ6として、ニューラルネットワークによるモデルを用いたボコーダを採用する場合、入力される音響特徴量を、メルスペクトログラムのデータとし、出力を当該メルスペクトログラムに対応する音声信号波形として学習処理を実行する。つまり、ボコーダ6において、メルスペクトログラムのデータを入力し、音声合成処理をニューラルネットワークによるモデルを用いた処理により実行し、音声波形データを出力させる。ボコーダ6から出力される当該音声波形データと、ボコーダに入力したメルスペクトログラムに対応する音声波形データ（正解の音声波形データ）とを比較し、両者の差（比較結果）（例えば、差分ベクトルのノルムやユークリッド距離により表現する差）が小さくなるように、ボコーダ6のニューラルネットワークのモデルのパラメータを更新する。ボコーダ6では、このパラメータ更新処理を繰り返し実行し、ボコーダの入力データ（メルスペクトログラムのデータ）と、ボコーダ6に入力されたメルスペクトログラムに対応する音声波形データ（正解の音声波形データ）との差が十分小さくなる（所定の誤差範囲

におさまる)、ニューラルネットワークのモデルのパラメータを最適化パラメータとして取得する。

[0100] ボコーダ6では、上記のようにして取得した最適化パラメータに基づいて、ボコーダ6のニューラルネットワークのモデルの各層に含まれるシナプス間の結合係数(重み係数)を設定することで、ボコーダ6のニューラルネットワークのモデルの最適化モデル(学習済みモデル)とすることができる。

[0101] 以上により、ボコーダ6において、入力をテキストデータとし、出力をメルスペクトログラムとするニューラルネットワークの学習済みモデル(最適化モデル)を構築できる。

[0102] なお、音声合成処理装置100において、(1)エンコーダ部3、デコーダ部5の学習処理と、(2)ボコーダ6の学習処理とを連携させて学習処理を実行してもよいし、上記のように、個別に学習処理を実行してもよい。音声合成処理装置100において、(1)エンコーダ部3、デコーダ部5の学習処理と、(2)ボコーダ6の学習処理とを連携させて学習処理を実行する場合、入力をテキストデータとし、当該テキストデータに対応する音声波形データ(正解の音声波形データ)とを用いて、(1)エンコーダ部3、デコーダ部5のニューラルネットワークのモデルと、(2)ボコーダ6のニューラルネットワークのモデルの最適化パラメータを取得することで学習処理を実行すればよい。

[0103] (1. 2. 2: 予測処理)

次に、音声合成処理装置100による予測処理について、説明する。なお、予測処理においても、説明便宜のため、処理対象言語を日本語として、以下、説明する。

[0104] 予測処理を実行する場合、音声合成処理装置100では、上記の学習処理により取得された学習済みモデル、すなわち、エンコーダ部3、デコーダ部5のニューラルネットワークの最適化モデル(最適化パラメータが設定されているモデル)、および、ボコーダ6のニューラルネットワークの最適化モデル(最適化パラメータが設定されているモデル)が構築されている。そし

て、音声合成処理装置100では、当該学習済みモデルを用いて予測処理が実行される。

[0105] 音声合成処理の対象とする日本語のテキストデータ D_{in} をテキスト解析部1に入力する。

[0106] テキスト解析部1は、入力されたテキストデータ D_{in} に対して、日本語用のテキスト解析処理を実行し、例えば、図2に示すパラメータを含む478次元のベクトルデータとして、フルコンテキストラベルデータ $D \times 1$ を取得する。

[0107] そして、取得されたフルコンテキストラベルデータ $D \times 1$ は、テキスト解析部1からフルコンテキストラベルベクトル処理部2に出力される。

[0108] フルコンテキストラベルベクトル処理部2は、入力されたフルコンテキストラベルデータ $D \times 1$ に対して、フルコンテキストラベルベクトル処理を実行し、最適化フルコンテキストラベル $D \times 2$ を取得する。なお、ここで取得される最適化フルコンテキストラベル $D \times 2$ は、エンコーダ部3、デコーダ部5のsequence-to-sequence方式のニューラルネットワークのモデルの学習処理を行うときに設定した最適化フルコンテキストラベルデータ $D \times 2$ と同じ次元数を有し、かつ、同じパラメータ（情報）を有するデータである。

[0109] 上記により取得されたデータ $D \times 2$ （最適化フルコンテキストラベルデータ $D \times 2$ ）は、フルコンテキストラベルベクトル処理部2からエンコーダ部3のエンコーダ側プレネット処理部31に出力される。

[0110] エンコーダ側プレネット処理部31は、フルコンテキストラベルベクトル処理部2から入力したデータ $D \times 2$ に対して、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理）を実行し、エンコーダ側LSTM層32に入力可能なデータを取得する。そして、エンコーダ側プレネット処理部31は、上記処理（プレネット処理）により取得したデータをデータ $D \times 3$ としてエンコーダ側LSTM層32に出力する。

- [0111] エンコーダ側LSTM層32は、エンコーダ側プレネット処理部31から、現時刻 t において出力されるデータ $D \times 3(t)$ と、1つ前の時間ステップにおいて、エンコーダ側LSTM層32から出力されたデータ $D \times 4(t-1)$ とを入力する。そして、エンコーダ側LSTM層32は、入力されたデータ $D \times 3(t)$ 、データ $D \times 4(t-1)$ に対して、LSTM層による処理（ニューラルネットワーク処理）を実行し、処理後のデータをデータ $D \times 4$ （データ $D \times 4(t)$ ）としてアテンション部4に出力する。
- [0112] アテンション部4は、エンコーダ部3から出力されるデータ $D \times 4$ と、デコーダ部5のデコーダ側LSTM層52から出力されるデータ h_o （出力側隠れ状態データ h_o ）とを入力する。アテンション部4は、エンコーダ部3から出力されるデータ $D \times 4$ 、すなわち、入力側隠れ状態データ h_i を所定の時間ステップ分記憶保持する。例えば、アテンション部4は、時間ステップ $t=1$ から $t=S$ （ S ：自然数）の期間において、エンコーダ部3により取得され、アテンション部4に出力されたデータ $D \times 4$ （ $=h_i$ ）の集合を、 $h_{i_{1..s}}$ （ $=\{D \times 4(1), D \times 4(2), \dots, D \times 4(S)\}$ ）として記憶保持する。
- [0113] また、アテンション部4は、デコーダ部5のデコーダ側LSTM層52から出力されるデータ D_y3 、すなわち、出力側隠れ状態データ h_o を所定の時間ステップ分記憶保持する。例えば、アテンション部4は、時間ステップ $t=1$ から $t=T$ （ T ：自然数）の期間において、デコーダ側LSTM層52により取得され、アテンション部4に出力されたデータ D_y3 （ $=h_o$ ）の集合を、 $h_{o_{1..T}}$ （ $=\{D_y3(1), D_y3(2), \dots, D_y3(T)\}$ ）として記憶保持する。
- [0114] そして、アテンション部4は、入力側隠れ状態データの集合データ $h_{i_{1..s}}$ と、出力側隠れ状態データの集合データ $h_{o_{1..T}}$ と、に基づいて、例えば、

$$c(t) = f1_attn(h_{i_{1..s}}, h_{o_{1..T}})$$

$f1_attn()$ ：コンテキスト状態データを取得する関数

に相当する処理を実行して、現時刻 t のコンテキスト状態データ $c(t)$ を取得する。

[0115] そして、アテンション部4は、取得したコンテキスト状態データ $c(t)$ をデコーダ側LSTM層52に出力する。

デコーダ側プレネット処理部51は、線形予測部53から出力される、1時間ステップ前のデータ $Dy_4(t-1)$ を入力する。デコーダ側プレネット処理部51は、例えば、複数層（例えば、2層）の全結合層を有しており、データの正規化処理（例えば、線形予測部53から出力されるデータ（ベクトルデータ）の次元数が $2N$ であり、デコーダ側LSTM層に入力されるデータ（ベクトルデータ）の次元数が N である場合、データの次元数を N にするように、例えば、ドロップアウト処理を行うことを含む）、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理）を実行し、デコーダ側LSTM層52に入力可能なデータを取得する。そして、デコーダ側プレネット処理部51は、上記処理（プレネット処理）により取得したデータをデータ Dy_2 としてデコーダ側LSTM層52に出力する。

[0116] デコーダ側LSTM層52は、デコーダ側プレネット処理部51から、現時刻 t において出力されるデータ $Dy_2(t)$ と、1つ前の時間ステップにおいて、デコーダ側LSTM層52から出力されたデータ $Dy_3(t-1)$ と、アテンション部4から出力される時刻 t のコンテキスト状態データ $c(t)$ とを入力する。

[0117] デコーダ側LSTM層52は、入力されたデータ $Dy_2(t)$ 、データ $Dy_3(t-1)$ 、および、コンテキスト状態データ $c(t)$ を用いて、LSTM層による処理を実行し、処理後のデータをデータ $Dy_3(t)$ として線形予測部53に出力する。また、デコーダ側LSTM層52は、データ $Dy_3(t)$ 、すなわち、時刻 t の出力側隠れ状態データ $h_o(t)$ をアテンション部4に出力する。

[0118] 線形予測部53は、デコーダ側LSTM層52から出力されるデータ Dy

3を入力する。線形予測部53は、所定の期間（例えば、メルスペクトログラムを取得するための1フレーム期間に相当する期間）内に、デコーダ側LSTM層52から出力されるデータDy3（複数のデータDy3）を記憶保持し、当該複数のデータDy3を用いて線形変換することで、所定期間におけるメルスペクトログラムの予測データDy4を取得する。そして、線形予測部53は、取得したデータDy4をポストネット処理部54、加算器55、および、デコーダ側プレネット処理部51に出力する。

[0119] ポストネット処理部54は、例えば、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理やtanh関数による処理）を実行し、予測データ（予測メルスペクトログラム）の残差データ（residual）を取得し、取得した残差データをデータDy5として加算器55に出力する。

[0120] 加算器55は、線形予測部53から出力される予測データDy4（予測メルスペクトログラムのデータ）と、ポストネット処理部54から出力される残差データDy5（予測メルスペクトログラムの残差データ）とを入力する。加算器55は、予測データDy4（予測メルスペクトログラムのデータ）と、残差データDy5（予測メルスペクトログラムの残差データ）とに対して加算処理を実行し、加算処理後のデータ（予測メルスペクトログラムのデータ）をデータDy6として、ボコーダ6に出力する。

[0121] ボコーダ6は、デコーダ部5の加算器55から出力されるデータDy6（予測メルスペクトログラムのデータ（音響特徴量のデータ））を入力とし、入力されたデータDy6に対して、学習済みモデルを用いたニューラルネットワーク処理による音声合成処理を実行し、データDy6（予測メルスペクトログラム）に対応する音声信号波形データを取得する。そして、ボコーダ6は、取得した音声信号波形データを、データDoutとして出力する。

[0122] このように、音声合成処理装置100では、入力されたテキストデータDinに対応する音声波形データDoutを取得することができる。

[0123] 以上のように、音声合成処理装置100では、処理対象言語（上記では日本語）のテキストを入力とし、当該処理対象言語に応じたテキスト解析処理により、フルコンテキストラベルデータを取得し、取得したフルコンテキストラベルデータからsequence-to-sequence方式を用いたニューラルネットワークのモデルで処理（学習処理、および／または、予測処理）を実行するのに適したデータである最適化フルコンテキストラベルデータを取得する。そして、音声合成処理装置100では、入力を最適化フルコンテキストラベルデータとし、出力をメルスペクトログラム（音響特徴量の一例）として、エンコーダ部3、アテンション部4、および、デコーダ部5において、ニューラルネットワークのモデルを用いた処理（学習処理、予測処理）を実行することで、高精度な処理を実現できる。さらに、音声合成処理装置100では、ボコーダ6により、上記により取得したメルスペクトログラム（音響特徴量の一例）から、当該メルスペクトログラムに対応する音声信号波形データを取得し、取得したデータを出力することで、音声波形データ（データOutput）を取得する。これにより、音声合成処理装置100では、入力されたテキストに相当する音声波形データを取得することができる。

[0124] つまり、音声合成処理装置100では、sequence-to-sequence方式を用いたニューラルネットワークのモデルで処理するのに適した最適化フルコンテキストラベルデータを用いて、ニューラルネットワークによる処理が実行されるため、高精度な音声合成処理を実行することができる。また、音声合成処理装置100では、処理対象言語に応じたテキスト解析処理を行い、当該テキスト解析処理で取得されたフルコンテキストラベルデータから、sequence-to-sequence方式を用いたニューラルネットワークのモデルで処理するのに適した最適化フルコンテキストラベルデータを取得し、取得した最適化フルコンテキストラベルデータを用いて処理を行うことで、任意の処理対象言語について、高精度な音声合成処理を行うことができる。

[0125] したがって、音声合成処理装置100では、日本語等の英語以外の言語を処理対象言語とする場合においても（処理対象言語を任意の言語にできる）

、sequence-to-sequence方式を用いたテキスト音声合成用のニューラルネットワークのモデルにより、学習・最適化を行い、高品質な音声合成処理を実現することができる。

[0126] 《第1変形例》

次に、第1実施形態の第1変形例について、説明する。なお、上記実施形態と同様の部分については、同一符号を付し、詳細な説明を省略する。

[0127] 本変形例の音声合成処理装置では、ボコーダ6が、例えば、下記先行技術文献に開示されているような、可逆変換が可能なニューラルネットワークのモデルを用いた処理を行う。この点が第1実施形態と相違し、それ以外については、本変形例の音声合成処理装置は、第1実施形態の音声合成処理装置100と同様である。

(先行技術文献A) :

R. Prenger, R. Valle, and B. Catanzaro, "WaveGlow: A flowbased generative network for speech synthesis," in Proc. ICASSP, May 2019.

図4は、第1実施形態の第1変形例の音声合成処理装置のボコーダ6の概略構成を示す図であり、学習処理時におけるデータの流れを明示した図である。

[0128] 図5は、第1実施形態の第1変形例の音声合成処理装置のボコーダ6の概略構成を示す図であり、予測処理時におけるデータの流れを明示した図である。

[0129] 本変形例のボコーダ6は、図4に示すように、ベクトル処理部61と、アップサンプリング処理部62と、 m 個 (m :自然数)の可逆処理部63a~63xとを備える。

[0130] まず、本変形例のボコーダ6の学習処理について、説明する。

[0131] 本変形例のボコーダ6は、学習処理において、音響特徴量としてメルスペクトログラム(これをデータ h とする)と、当該メルスペクトログラムに対応する音声信号波形データ(正解データ)(これをデータ x とする)とを入力し、ガウス白色ノイズ(これをデータ z とする)を出力する。

[0132] ベクトル処理部 6 1 は、学習処理時において、音声信号波形データ x を入力し、入力したデータ x に対して、例えば、コンボリューション処理を施して、可逆処理部 6 3 a（学習処理時において最初にデータ入力される可逆処理部）に入力可能な次元数のベクトルデータ $D \times 1$ に変換する。そして、ベクトル処理部 6 1 は、変換したベクトルデータ $D \times 1$ を可逆処理部 6 3 a に出力する。

[0133] アップサンプリング処理部 6 2 は、音響特徴量としてメルスペクトログラムのデータ h を入力し、入力されたメルスペクトログラムのデータ h に対して、アップサンプリング処理を実行し、処理後のデータ（アップサンプリングされたメルスペクトログラムのデータ）をデータ $h 1$ として、可逆処理部 6 3 a ~ 6 3 x のそれぞれの WN 変換部 6 3 2 に出力する。

[0134] 可逆処理部 6 3 a は、図 4 に示すように、可逆 1×1 畳み込み層と、アフィンカップリング層とを備える。

[0135] 可逆 1×1 畳み込み層は、ベクトル処理部 6 1 から出力されるデータ $D \times 1$ を入力とし、入力されたデータに対して、重み係数行列 W_k ($k = 1$)（シナプス間の結合係数（重み係数）を規定する行列）により、ニューラルネットワーク処理を実行する、つまり、

$$D \times A_1 = W_1 \times D \times 1$$

に相当する処理を実行して、データ $D \times A_1$ を取得する。

[0136] なお、重み係数行列 W_k は、直行行列となるように設定されており、したがって、逆変換が可能となる。

[0137] このようにして取得されたデータ $D \times A_1$ は、可逆 1×1 畳み込み層からアフィンカップリング層に出力される。

[0138] アフィンカップリング層では、データ分割部 6 3 1 により、

$$x = D \times A_1$$

$$x_a, x_b = \text{split}(x)$$

$\text{split}()$: データ分割をする関数

に相当する処理を実行し、入力データ x を 2 分割し、分割データ x_a と x_b を

取得する。例えば、 x が $n_1 \times 2$ (n_1 :自然数)のビット数のデータである場合、 x_a は、 x の上位 n_1 ビット分のデータであり、 x_b は、 x の下位 n_1 ビット分のデータである。

[0139] そして、データ x_a は、MN変換部632およびデータ合成部634に出力される。また、データ x_b は、アフィン変換部633に出力される。

[0140] MN変換部632は、データ分割部631から出力されるデータ x_a と、アップサンプリング処理部62から出力されるアップサンプリングされたメルスペクトログラムのデータ h_1 とを入力する。そして、MN変換部632は、データ x_a と、データ h_1 とに対して、任意の変換であるMN変換(例えば、WaveNetによる変換)を実行し、アフィン変換のパラメータとするデータ s_j 、 t_j (s_j :アフィン変換用の行列、 t_j :アフィン変換用のオフセット)を取得する。取得されたアフィン変換のパラメータとするデータ s_j 、 t_j は、WN変換部632からアフィン変換部633に出力される。

[0141] アフィン変換部633は、MN変換部632により取得されたデータ s_j 、 t_j を用いて、データ分割部631から入力されるデータ x_b に対して、アフィン変換を行う。つまり、アフィン変換部633は、

$$\begin{aligned} x_b' &= \text{Affin}(s_j, t_j, x_b) \\ &= s_j \times x_b + t_j \end{aligned}$$

に相当する処理を実行することで、データ x_b のアフィン変換後のデータ x_b' を取得し、取得したデータ x_b' をデータ合成部634に出力する。

[0142] データ合成部634では、データ分割部631から出力されるデータ x_a と、アフィン変換部633から出力されるデータ x_b' とを入力し、データ x_a と、データ x_b' とを合成する処理、すなわち、

$$Dx_2 = \text{concat}(x_a, x_b')$$

に相当する処理を実行し、データ Dx_2 を取得する。なお、データ合成部634でのデータ合成処理は、例えば、 x_a 、 x_b' が、それぞれ、 n_1 ビットのデータである場合、上位 n_1 ビットが x_a となり、下位 n_1 ビットが x_b' となる $n_1 \times 2$ ビットのデータを取得する処理である。

[0143] このようにして取得されたデータ $D \times_2$ は、可逆処理部 63 a から、可逆処理部 63 b（2番目の可逆処理部）に出力される。

[0144] 可逆処理部 63 b ~ 63 x では、可逆処理部 63 a と同様の処理が実行される。つまり、本変形例のボコーダ 6 では、図 4 に示すように、可逆処理部 63 a の処理が m 回繰り返し実行される。そして、最終段の可逆処理部 63 x からのデータ z が出力される。なお、本変形例のボコーダ 6 は、m 個の可逆処理部を備えるものとする。

[0145] そして、本変形例のボコーダ 6 では、出力データ z が、ガウス白色ノイズとなるように、ニューラルネットワークのモデルの学習を行う。つまり、x を入力としたときの z を $z(x)$ とすると、 $z(x)$ がガウス分布 $N(\mu, \sigma)$ (μ は平均値であり $\mu = 0$ 、 σ は標準偏差) に従うガウス確率変数となるように、本変形例のボコーダ 6 のニューラルネットワークのモデルのパラメータを設定する。なお、 σ は、例えば、入力される音響特徴量としてメルスペクトログラムのデータの情報量 I に相関のあるデータとする。

[0146] つまり、本変形例のボコーダ 6 では、x が入力されたときの尤度 (θ : ニューラルネットワークのパラメータ) $p_\theta(x)$ を、下記数式により規定することができ、当該尤度 $p_\theta(x)$ を最大にするパラメータ θ_{opt} を取得することで、学習処理を実行する。

[数1]

$$-\log p_\theta(\mathbf{x}) = \frac{\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x})}{2\sigma_{WG}^2} - \sum_{j=1}^{m1} \log s_j(\mathbf{x}, \mathbf{h}) - \sum_{k=1}^{m2} \log |\det(\mathbf{W}_k)|$$

$p_\theta(x)$: x が入力されたときの尤度 (θ : ニューラルネットワークのパラメータ)

$s_j(x, h)$: x、h が入力されたときの j 番目のアフィンカップリング層の出力係数ベクトル

W_k : k 番目の可逆 1×1 畳み込み層の係数行列 (重み付け係数の行列)

$z(x)$: x が入力されたときの出力値 (出力ベクトル)。

h : 音響特徴量 (ここでは、メルスペクトログラム)

σ_{WG}^2 : ガウス分布の予測分散値

なお、 $z(x)$ は、ガウス分布 $N(\mu, \sigma)$ (μ は平均値であり $\mu = 0$ 、 σ は標準偏差) に従うガウス確率変数に相当するものである。すなわち、 $z \sim N(\mu, \sigma) = N(0, \sigma)$ である。また、 m_1 は、アフィンカップリング層の処理の回数、 m_2 は、可逆 1×1 畳み込み層の処理の回数であり、本変形例のボコーダ 6 では、 $m_1 = m_2 = m$ である。

[0147] 本変形例のボコーダ 6 では、下記数式に相当する処理を実行することで、本変形例のボコーダ 6 のニューラルネットワークのモデルの最適化パラメータ θ_{opt} を取得する。

[数2]

$$\theta_{opt} = \underset{\theta}{\operatorname{argmax}} p_{\theta}(x)$$

本変形例のボコーダ 6 では、上記の学習処理により取得した最適化パラメータ θ_{opt} により、ニューラルネットワークのモデルのパラメータが設定され (各可逆処理部 63b ~ 63x のアフィンカップリング層、可逆 1×1 畳み込み層のパラメータが設定され)、学習済みモデルが構築される。

[0148] 次に、本変形例のボコーダ 6 の予測処理について、説明する。

[0149] 本変形例のボコーダ 6 は、予測処理において、音響特徴量としてメルスペクトログラム (これをデータ h とする) と、当該メルスペクトログラムの情報量 l に相関のあるデータを標準偏差 σ とし、平均値を「0」とするガウス白色ノイズ z とを入力とする。

[0150] 本変形例のボコーダ 6 では、予測処理時において、図 5 に示すように、学習処理時とは、逆の処理が実行される。

[0151] メルスペクトログラムのデータ (例えば、デコーダ部 5 から出力されるデータ Dy_6) がアップサンプリング処理部 62 に入力される。

- [0152] アップサンプリング処理部62は、音響特徴量としてメルスペクトログラムのデータ h を入力し、入力されたメルスペクトログラムのデータ h に対して、アップサンプリング処理を実行し、処理後のデータ（アップサンプリングされたメルスペクトログラムのデータ）をデータ h_1 として、可逆処理部63a~63xのそれぞれのWN変換部632に出力する。
- [0153] また、ガウス白色ノイズ z （データ z という）が可逆処理部63xに入力される。
- [0154] そして、可逆処理部63xにおいて、入力されたデータ z に対して、アフィンカップリング層の処理、可逆 1×1 畳み込み層の層の処理が実行される。この処理が、図5に示すように、 m 回繰り返し実行される。各処理は、同様であるので、可逆処理部63aでの処理について、説明する。
- [0155] データ合成部634では、可逆処理部63bから出力されるデータ Dx'_2 を入力し、学習処理時とは逆の処理、すなわち、データ分割処理を実行する。つまり、データ合成部634では、
- $$x = Dx'_2$$
- $$x_a, x_b' = \text{split}(x)$$
- $\text{split}()$: データ分割をする関数
- に相当する処理を実行し、入力データ x を2分割し、分割データ x_a と x_b' を取得する。
- [0156] そして、データ合成部634は、取得したデータ x_a をMN変換部632およびデータ分割部631に出力し、データ x_b' をアフィン変換部633に出力する。
- [0157] MN変換部632は、データ合成部634から出力されるデータ x_a と、アップサンプリング処理部62から出力されるアップサンプリングされたメルスペクトログラムのデータ h_1 とを入力する。そして、MN変換部632は、データ x_a と、データ h_1 とに対して、任意の変換であるMN変換（例えば、WaveNetによる変換）を実行し、アフィン変換のパラメータとするデータ s_j, t_j （ s_j : アフィン変換用の行列、 t_j : アフィン変換用のオフ

セット)を取得する。取得されたアフィン変換のパラメータとするデータ s_j , t_j は、WN変換部632からアフィン変換部633に出力される。

[0158] アフィン変換部633は、MN変換部632により取得されたデータ s_j , t_j を用いて、データ合成部634から入力されるデータ x'_b に対して、アフィン逆変換（学習処理時に行ったアフィン変換の逆変換）を行う。つまり、アフィン変換部633は、

$$x_b = \text{Affin}^{-1}(s_j, t_j, x'_b)$$

に相当する処理を実行することで、データ x'_b のアフィン逆変換後のデータ x_b を取得し、取得したデータ x_b をデータ分割部631に出力する。

[0159] データ分割部631は、データ合成部634から出力されるデータ x_a と、アフィン変換部633から出力されるデータ x_b とを入力し、データ x_a と、データ x_b とを合成する処理、すなわち、

$$Dx'_1 = \text{concat}(x_a, x_b)$$

に相当する処理を実行し、データ Dx'_1 を取得する。そして、データ分割部631は、取得したデータ Dx'_1 を出力する。

[0160] 上記のようにして可逆処理部63x~63aにより処理されることで取得されたデータ Dx'_1 が、ベクトル処理部61に入力される。

[0161] ベクトル処理部61は、学習処理時と逆の処理を実行することで、データ Dx'_1 から、予測音声信号波形データ x を取得し、出力する。

[0162] 以上のように処理することで、本変形例のボコーダ6では、入力 z （ガウス白色ノイズ z ）と、メルスペクトログラムのデータ h から、予測音声信号波形データ x を取得することができる。

[0163] 本変形例のボコーダ6では、ニューラルネットワークを可逆変換できる構成を採用している。このため、本変形例のボコーダ6では、（1）ガウス白色ノイズが入力されたときに出力される音声波形データの尤度と、（2）音声波形データが入力されたときに出力されるガウス白色ノイズの尤度とを等価にし、学習処理を行いやすい（計算が容易である）後者（音声波形データが入力されたときに出力されるガウス白色ノイズの尤度）により、学習処理

を行うことで、効率良く学習処理を行うことができる。

[0164] そして、本変形例のボコーダ6では、ニューラルネットワークを可逆変換できる構成を有しているため、上記学習処理により取得した学習済みモデルにより、予測処理を、学習処理時とは逆の処理（逆変換）により実現できる。

[0165] このように、本変形例のボコーダ6では、音響特徴量としてメルスペクトログラムのデータから音声波形データを直接予測（取得）できる構成をシンプルな構成で実現できる。そして、本変形例のボコーダ6では、このようなシンプルな構成を有しているため、処理精度を保ちながら、予測処理を高速に行うことができ、音声合成処理をリアルタイムで実行することが可能になる。

[0166] 図6は、本変形例の音声合成処理装置によりTTS処理（処理対象言語：日本語）実行し、取得した音声波形データのメルスペクトログラム（予測データ）と、入力テキストの実際の音声波形データのメルスペクトログラム（オリジナルデータ）とを示す図である。

[0167] 図6から分かるように、本変形例の音声合成処理装置によりTTS処理では、非常に高精度な音声波形データが予測（取得）できる。

[0168] [第2実施形態]

次に、第2実施形態について、説明する。なお、上記実施形態と同様の部分については、同一符号を付し、詳細な説明を省略する。

[0169] 第1実施形態では、エンコーダ・デコーダ方式（sequence-to-sequence方式）を用いた音声合成処理装置100について、説明した。第1実施形態の音声合成処理装置100は、注意機構（アテンション部4）を備えており、音素継続長と音響モデルとを注意機構を用いて同時に最適化するニューラル音声合成処理装置100では、自然音声クラスの高音質なテキスト音声合成を実現できる。しかしながら、第1実施形態の音声合成処理装置100では、推論時（予測処理時）に、まれに注意機構予測が失敗することがあり、これに

より合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題がある。

[0170] 第2実施形態では、上記問題を解決するための技術について、説明する。

[0171] <2. 1 : 音声合成処理装置の構成>

図7は、第2実施形態に係る音声合成処理装置200の概略構成図である。

[0172] 第2実施形態に係る音声合成処理装置200は、第1実施形態の音声合成処理装置100において、アテンション部4を削除し、音素継続長推定部7を追加した構成を有している。そして、第2実施形態に係る音声合成処理装置200は、第1実施形態の音声合成処理装置100において、テキスト解析部1をテキスト解析部1Aに置換し、フルコンテキストラベルベクトル処理部2をフルコンテキストラベルベクトル処理部2Aに置換し、デコーダ部5をデコーダ部5Aに置換した構成を有している。

[0173] テキスト解析部1Aは、第1実施形態のテキスト解析部1と同様の機能を有しており、さらに、音素のコンテキストラベルを取得する機能を有している。テキスト解析部1Aは、処理対象言語のテキストデータD_{in}から音素のコンテキストラベルを取得し、取得した音素のコンテキストラベルのデータをデータD_{x01}として、音素継続長推定部7に出力する。

[0174] 音素継続長推定部7は、テキスト解析部1Aから出力されるデータD_{x01}（音素のコンテキストラベルのデータ）を入力する。音素継続長推定部7は、データD_{x01}（音素のコンテキストラベルのデータ）から、データD_{x01}に対応する音素の音素継続長を推定（取得）する音素継続長推定処理を実行する。具体的には、音素継続長推定部7は、例えば、隠れマルコフモデル（HMM：Hidden Markov Model）、ニューラルネットワークモデル等を用いた、音素のコンテキストラベルから当該音素の音素継続長を推定（予測）するモデル（処理システム）により、音素継続長推定処理を実行する。

[0175] そして、音素継続長推定部7は、音素継続長推定処理により取得（推定）

した音素継続長のデータをデータD×02として、フルコンテキストラベルベクトル処理部2Aに出力する。

[0176] フルコンテキストラベルベクトル処理部2Aは、第1実施形態のフルコンテキストラベルベクトル処理部2と同様の機能を有しており、さらに、音素継続長推定部7により推定された音素継続長に相当する期間において、当該音素継続長に対応する音素についての最適化フルコンテキストラベルデータをエンコーダ部3に継続して出力する機能を有する。

[0177] フルコンテキストラベルベクトル処理部2Aは、テキスト解析部1から出力されるデータD×1（フルコンテキストラベルのデータ）と、音素継続長推定部7から出力されるデータD×02（音素の音素継続長のデータ）とを入力する。

[0178] フルコンテキストラベルベクトル処理部2Aは、テキスト解析部1Aから出力されるデータD×1（フルコンテキストラベルのデータ）を入力する。フルコンテキストラベルベクトル処理部2Aは、入力されたフルコンテキストラベルデータD×1から、sequence-to-sequence方式のニューラルネットワークのモデルの学習処理に適したフルコンテキストラベルデータを取得するためのフルコンテキストラベルベクトル処理を実行する。そして、フルコンテキストラベルベクトル処理部2Aは、フルコンテキストラベルベクトル処理により取得したデータをデータD×2（最適化フルコンテキストラベルデータD×2）として、エンコーダ部3のエンコーダ側プレネット処理部3-1に出力する。このとき、フルコンテキストラベルベクトル処理部2Aは、音素継続長推定部7により推定された音素継続長に相当する期間において、当該音素継続長に対応する音素についての最適化フルコンテキストラベルデータをエンコーダ部3に継続して出力する。

[0179] デコーダ部5Aは、第1実施形態のデコーダ部5において、デコーダ側LSTM層5-2をデコーダ側LSTM層5-2Aに置換した構成を有している。それ以外は、デコーダ部5Aは、第1実施形態のデコーダ部5と同様である。

[0180] デコーダ側LSTM層52Aは、デコーダ側LSTM層52と同様の機能を有している。デコーダ側LSTM層52Aは、デコーダ側プレネット処理部51から、現時刻 t において出力されるデータ $Dy2$ （これをデータ $Dy2(t)$ と表記する）と、1つ前の時間ステップにおいて、デコーダ側LSTM層52Aから出力されたデータ $Dy3$ （これをデータ $Dy3(t-1)$ と表記する）と、エンコーダ部3から出力される時刻 t の入力側隠れ状態データ $hi(t)$ を入力する。

[0181] デコーダ側LSTM層52Aは、入力されたデータ $Dy2(t)$ 、データ $Dy3(t-1)$ 、および、入力側隠れ状態データ $hi(t)$ を用いて、LSTM層による処理を実行し、処理後のデータをデータ $Dy3$ （データ $Dy3(t)$ ）として線形予測部53に出力する。

[0182] <2.2:音声合成処理装置の動作>

以上のように構成された音声合成処理装置200の動作について以下説明する。

[0183] 図8は、推定された音素継続長に基づいて、エンコーダ部3に入力するデータ $Dx2$ を生成する処理を説明するための図である。

[0184] 以下では、音声合成処理装置200の動作を、(1)学習処理（学習時の処理）と、(2)予測処理（予測時の処理）とに分けて説明する。

[0185] (2.2.1:学習処理)

まず、音声合成処理装置200による学習処理について、説明する。なお、説明便宜のため、処理対象言語を日本語として、以下、説明する。

[0186] 処理対象言語である日本語のテキストデータ Din をテキスト解析部1Aに入力する。また、当該テキストデータ Din に対応するメルスペクトログラム（音響特徴量）のデータを教師データとして用意する。

[0187] テキスト解析部1Aは、第1実施形態と同様に、入力されたテキストデータ Din に対して、テキスト解析処理を実行し、様々な言語情報からなるコンテキストを含む音素ラベルであるコンテキストラベルの系列を取得する。

[0188] テキスト解析部1Aは、第1実施形態と同様に、取得したフルコンテキス

トラベルデータをデータ $D \times 1$ としてフルコンテキストラベルベクトル処理部 2 に出力する。

[0189] また、テキスト解析部 1 A は、処理対象言語のテキストデータ D_{in} から音素のコンテキストラベルを取得し、取得した音素のコンテキストラベルのデータをデータ $D \times 0 1$ として、音素継続長推定部 7 に出力する。

[0190] 音素継続長推定部 7 は、テキスト解析部 1 A から出力されるデータ $D \times 0 1$ (音素のコンテキストラベルのデータ) から、データ $D \times 0 1$ に対応する音素の音素継続長を推定 (取得) する音素継続長推定処理を実行する。具体的には、音素継続長推定部 7 は、例えば、隠れマルコフモデル (HMM: Hidden Markov Model)、ニューラルネットワークモデル等を用いた、音素のコンテキストラベルから当該音素の音素継続長を推定 (予測) するモデル (処理システム) により、音素継続長推定処理を実行する。

[0191] そして、音素継続長推定部 7 は、音素継続長推定処理により取得 (推定) した音素継続長のデータをデータ $D \times 0 2$ として、フルコンテキストラベルベクトル処理部 2 A に出力する。

[0192] フルコンテキストラベルベクトル処理部 2 A は、テキスト解析部 1 A から出力されるデータ $D \times 1$ (フルコンテキストラベルのデータ) から、sequence-to-sequence方式のニューラルネットワークのモデルの学習処理に適したフルコンテキストラベルデータを取得するためのフルコンテキストラベルベクトル処理 (第 1 実施形態と同様のフルコンテキストラベルベクトル処理) を実行する。そして、フルコンテキストラベルベクトル処理部 2 A は、フルコンテキストラベルベクトル処理により取得したデータをデータ $D \times 2$ (最適化フルコンテキストラベルデータ $D \times 2$) として、エンコーダ部 3 のエンコーダ側プレネット処理部 3 1 に出力する。このとき、フルコンテキストラベルベクトル処理部 2 A は、音素継続長推定部 7 により推定された音素継続長に相当する期間において、当該音素継続長に対応する音素についての最適化フルコンテキストラベルデータをエンコーダ部 3 に継続して出力する。

- [0193] フルコンテキストラベルベクトル処理部 2 A により取得されたデータ $D \times 2$ （最適化フルコンテキストラベルデータ $D \times 2$ ）は、フルコンテキストラベルベクトル処理部 2 からエンコーダ部 3 のエンコーダ側プレネット処理部 3 1 に出力される。
- [0194] エンコーダ側プレネット処理部 3 1 は、フルコンテキストラベルベクトル処理部 2 から入力したデータ $D \times 2$ に対して、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU 関数（ReLU: Rectified Linear Unit）による処理）を実行し、エンコーダ側 LSTM 層 3 2 に入力可能なデータを取得する。そして、エンコーダ側プレネット処理部 3 1 は、上記処理（プレネット処理）により取得したデータをデータ $D \times 3$ としてエンコーダ側 LSTM 層 3 2 に出力する。
- [0195] エンコーダ側 LSTM 層 3 2 は、エンコーダ側プレネット処理部 3 1 から、現時刻 t において出力されるデータ $D \times 3 (t)$ と、1 つ前の時間ステップにおいて、エンコーダ側 LSTM 層 3 2 から出力されたデータ $D \times 4 (t-1)$ とを入力する。そして、エンコーダ側 LSTM 層 3 2 は、入力されたデータ $D \times 3 (t)$ 、データ $D \times 4 (t-1)$ に対して、LSTM 層による処理を実行し、処理後のデータをデータ $D \times 4$ （データ $D \times 4 (t)$ （=入力側隠れ状態データ $h_i (t)$ ））としてデコーダ部 5 A のデコーダ側 LSTM 層 5 2 A に出力する。
- [0196] デコーダ側プレネット処理部 5 1 は、線形予測部 5 3 から出力される、1 時間ステップ前のデータ $D y 4 (t-1)$ を入力する。デコーダ側プレネット処理部 5 1 は、例えば、複数層（例えば、2 層）の全結合層を有しており、データの正規化処理（例えば、線形予測部 5 3 から出力されるデータ（ベクトルデータ）の次元数が $2N$ であり、デコーダ側 LSTM 層に入力されるデータ（ベクトルデータ）の次元数が N である場合、データの次元数を N にするように、例えば、ドロップアウト処理を行うことを含む）、活性化関数による処理（例えば、ReLU 関数（ReLU: Rectified Li

near Unit) による処理) を実行し、デコーダ側LSTM層52に
入力可能なデータを取得する。そして、デコーダ側プレネット処理部51は
、上記処理(プレネット処理)により取得したデータをデータDy2として
デコーダ側LSTM層52に出力する。

[0197] デコーダ側LSTM層52Aは、デコーダ側プレネット処理部51から、
現時刻tにおいて出力されるデータDy2(t)と、1つ前の時間ステップ
において、デコーダ側LSTM層52から出力されたデータDy3(t-1)
)と、エンコーダ部3から出力される時刻tの入力側隠れ状態データhi(
t)(=Dx4(t))とを入力する。

[0198] デコーダ側LSTM層52Aは、入力されたデータDy2(t)、データ
Dy3(t-1)、および、入力側隠れ状態データhi(t)を用いて、L
STM層による処理を実行し、処理後のデータをデータDy3(t)として
線形予測部53に出力する。

線形予測部53、ポストネット処理部54、および、加算器55では、第
1実施形態と同様の処理が実行される。

[0199] そして、音声合成処理装置200では、上記のように取得されたデータD
y6(予測メルスペクトログラムのデータ)と、テキストデータDinに対
応するメルスペクトログラム(音響特徴量)の教師データ(正解のメルスペ
クトログラム)とを比較し、両者の差(比較結果)(例えば、差分ベクトル
のノルムやユークリッド距離により表現する差)が小さくなるように、エン
コーダ部3、デコーダ部5Aのニューラルネットワークのモデルのパラメー
タを更新する。音声合成処理装置100では、このパラメータ更新処理を繰
り返し実行し、データDy6(予測メルスペクトログラムのデータ)と、テ
キストデータDinに対応するメルスペクトログラム(音響特徴量)の教師
データ(正解のメルスペクトログラム)との差が十分小さくなる(所定の誤
差範囲におさまる)、ニューラルネットワークのモデルのパラメータを最適
化パラメータとして取得する。

[0200] 音声合成処理装置200では、上記のようにして取得した最適化パラメー

タに基づいて、エンコーダ部3、デコーダ部5Aのニューラルネットワークのモデルの各層に含まれるシナプス間の結合係数（重み係数）を設定することで、エンコーダ部3、デコーダ部5Aのニューラルネットワークのモデルを最適化モデル（学習済みモデル）とすることができる。

[0201] 以上により、音声合成処理装置200において、入力をテキストデータとし、出力をメルスペクトログラムとするニューラルネットワークの学習済みモデル（最適化モデル）を構築できる。

[0202] なお、音声合成処理装置200において、第1実施形態の音声合成処理装置100における学習処理により取得したニューラルネットワークの学習済みモデル（最適化モデル）を用いてもよい。つまり、音声合成処理装置200において、第1実施形態の音声合成処理装置100における学習処理により取得したニューラルネットワークの学習済みモデルのエンコーダ部3およびデコーダ部5の最適パラメータを用いて、音声合成処理装置200のエンコーダ部3およびデコーダ部5Aのパラメータを設定することで、音声合成処理装置200において、学習済みモデルを構築するようにしてもよい。

[0203] また、ボコーダ6として、ニューラルネットワークによるモデルを用いたボコーダを採用する場合、その学習処理は、第1実施形態と同様である。

[0204] これにより、第1実施形態と同様に、ボコーダ6において、入力をテキストデータとし、出力をメルスペクトログラムとするニューラルネットワークの学習済みモデル（最適化モデル）を構築できる。

[0205] なお、音声合成処理装置200において、（1）エンコーダ部3、デコーダ部5Aの学習処理と、（2）ボコーダ6の学習処理とを連携させて学習処理を実行してもよいし、上記のように、個別に学習処理を実行してもよい。音声合成処理装置200において、（1）エンコーダ部3、デコーダ部5Aの学習処理と、（2）ボコーダ6の学習処理とを連携させて学習処理を実行する場合、入力をテキストデータとし、当該テキストデータに対応する音声波形データ（正解の音声波形データ）とを用いて、（1）エンコーダ部3、デコーダ部5Aのニューラルネットワークのモデルと、（2）ボコーダ6の

ニューラルネットワークのモデルの最適化パラメータを取得することで学習処理を実行すればよい。

[0206] (2. 2. 2 : 予測処理)

次に、音声合成処理装置200による予測処理について、説明する。なお、予測処理においても、説明便宜のため、処理対象言語を日本語として、以下、説明する。

[0207] 予測処理を実行する場合、音声合成処理装置200では、上記の学習処理により取得された学習済みモデル、すなわち、エンコーダ部3、デコーダ部5Aのニューラルネットワークの最適化モデル（最適化パラメータが設定されているモデル）、および、ボコーダ6のニューラルネットワークの最適化モデル（最適化パラメータが設定されているモデル）が構築されている。そして、音声合成処理装置200では、当該学習済みモデルを用いて予測処理が実行される。

[0208] 音声合成処理の対象とする日本語のテキストデータ D_{in} をテキスト解析部1Aに入力する。

[0209] テキスト解析部1Aは、入力されたテキストデータ D_{in} に対して、日本語用のテキスト解析処理を実行し、例えば、図2に示すパラメータを含む478次元のベクトルデータとして、フルコンテキストラベルデータ D_{x1} を取得する。

[0210] そして、取得されたフルコンテキストラベルデータ D_{x1} は、テキスト解析部1Aからフルコンテキストラベルベクトル処理部2Aに出力される。

[0211] また、テキスト解析部1Aは、処理対象言語のテキストデータ D_{in} から音素のコンテキストラベルを取得し、取得した音素のコンテキストラベルのデータをデータ D_{x01} として、音素継続長推定部7に出力する。

[0212] 音素継続長推定部7は、テキスト解析部1Aから出力されるデータ D_{x01} （音素のコンテキストラベルのデータ）から、データ D_{x01} に対応する音素の音素継続長を推定（取得）する音素継続長推定処理を実行する。具体的には、音素継続長推定部7は、例えば、隠れマルコフモデル（HMM：H

idden Markov Model)、ニューラルネットワークモデル等を用いた、音素のコンテキストラベルから当該音素の音素継続長を推定(予測)するモデル(処理システム)により、音素継続長推定処理を実行する。

- [0213] 例えば、図8に示すように、入力データ D_{in} が「今日の天気は. . .」である場合、データ D_{x01} に含まれる各音素のデータを、
- (1) $ph_0 = 「k」$ 、(2) $ph_1 = 「y」$ 、(3) $ph_2 = 「ou」$ 、(4) $ph_3 = 「n」$ 、(5) $ph_{04} = 「o」$ 、(6) $ph_{sil} = \text{無音状態}$ 、(7) $ph_5 = 「t」$ 、(8) $ph_6 = 「e」$ 、(9) $ph_{07} = 「n」$ 、. . . とし、音素 ph_k (k : 整数)の推定された音素継続長を $dur(ph_k)$ とすると、音素継続長推定部7は、音素 ph_k (k : 整数)のコンテキストラベルを用いて、音素継続長推定処理を実行することで、音素 ph_k の推定された音素継続長 $dur(ph_k)$ を取得する。例えば、上記の各音素(音素 ph_k)について、音素継続長推定部7により取得(推定)された音素継続長 $dur(ph_k)$ が、図8に示す時間の長さ(継続長)を有するものとする。
- [0214] そして、音素継続長推定部7は、音素継続長推定処理により取得(推定)した音素継続長のデータ(図8の場合、 $dur(ph_k)$)をデータ D_{x02} として、フルコンテキストラベルベクトル処理部2Aに出力する。
- [0215] フルコンテキストラベルベクトル処理部2Aは、入力されたフルコンテキストラベルデータ D_{x1} に対して、フルコンテキストラベルベクトル処理を実行し、最適化フルコンテキストラベル D_{x2} を取得する。なお、ここで取得される最適化フルコンテキストラベル D_{x2} は、エンコーダ部3、デコーダ部5Aのsequence-to-sequence方式のニューラルネットワークのモデルの学習処理を行うときに設定した最適化フルコンテキストラベルデータ D_{x2} と同じ次元数を有し、かつ、同じパラメータ(情報)を有するデータである。
- [0216] 上記により取得されたデータ D_{x2} (最適化フルコンテキストラベルデータ D_{x2})は、フルコンテキストラベルベクトル処理部2からエンコーダ部

3のエンコーダ側プレネット処理部31に出力される。このとき、フルコンテキストラベルベクトル処理部2Aは、音素継続長推定部7により推定された音素継続長に相当する期間において、当該音素継続長に対応する音素についての最適化フルコンテキストラベルデータをエンコーダ部3に継続して出力する。例えば、図8に示すように、音素 ph_k についての最適化フルコンテキストラベルデータをデータ $D \times 2 (ph_k)$ とすると、フルコンテキストラベルベクトル処理部2Aは、音素 ph_k についての最適化フルコンテキストラベルデータ $D \times 2 (ph_k)$ を、当該音素 ph_k の推定された音素継続長 $dur(ph_k)$ に相当する期間において、継続してエンコーダ部3に出力する。

[0217] つまり、音素 ph_k についての最適化フルコンテキストラベルデータ $D \times 2 (ph_k)$ は、推定された音素継続長 $dur(ph_k)$ に相当する期間、繰り返しエンコーダ部3に出力される。すなわち、フルコンテキストラベルベクトル処理部2Aでは、推定された音素継続長 $dur(ph_k)$ に基づいて、エンコーダ部3へ入力するデータ（最適化フルコンテキストラベルデータ $D \times 2 (ph_k)$ ）の時間引き延ばし処理が実行される。

[0218] エンコーダ側プレネット処理部31は、フルコンテキストラベルベクトル処理部2Aから入力したデータ $D \times 2$ に対して、コンボリューション処理（コンボリューションフィルタによる処理）、データの正規化処理、活性化関数による処理（例えば、ReLU関数（ReLU: Rectified Linear Unit）による処理）を実行し、エンコーダ側LSTM層32に入力可能なデータを取得する。そして、エンコーダ側プレネット処理部31は、上記処理（プレネット処理）により取得したデータをデータ $D \times 3$ としてエンコーダ側LSTM層32に出力する。

[0219] エンコーダ側LSTM層32は、エンコーダ側プレネット処理部31から、現時刻 t において出力されるデータ $D \times 3 (t)$ と、1つ前の時間ステップにおいて、エンコーダ側LSTM層32から出力されたデータ $D \times 4 (t-1)$ とを入力する。そして、エンコーダ側LSTM層32は、入力されたデータ $D \times 3 (t)$ 、データ $D \times 4 (t-1)$ に対して、LSTM層による

処理（ニューラルネットワーク処理）を実行し、処理後のデータをデータ $D \times 4$ （データ $D \times 4(t)$ （=入力側隠れ状態データ $h_i(t)$ ））としてデコーダ部 5 A のデコーダ側 LSTM 層 5 2 A に出力する。

[0220] デコーダ側 LSTM 層 5 2 A は、入力されたデータ $D y 2(t)$ 、データ $D y 3(t-1)$ 、および、入力側隠れ状態データ $h_i(t)$ を用いて、LSTM 層による処理を実行し、処理後のデータをデータ $D y 3(t)$ として線形予測部 5 3 に出力する。

線形予測部 5 3、ポストネット処理部 5 4、および、加算器 5 5 では、第 1 実施形態と同様の処理が実行される。

[0221] ボコーダ 6 は、デコーダ部 5 A の加算器 5 5 から出力されるデータ $D y 6$ （予測メルスペクトログラムのデータ（音響特徴量のデータ））を入力とし、入力されたデータ $D y 6$ に対して、学習済みモデルを用いたニューラルネットワーク処理による音声合成処理を実行し、データ $D y 6$ （予測メルスペクトログラム）に対応する音声信号波形データを取得する。そして、ボコーダ 6 は、取得した音声信号波形データを、データ $D o u t$ として出力する。

[0222] このように、音声合成処理装置 2 0 0 では、入力されたテキストデータ $D i n$ に対応する音声波形データ $D o u t$ を取得することができる。

[0223] 以上のように、音声合成処理装置 2 0 0 では、処理対象言語（上記では日本語）のテキストを入力とし、当該処理対象言語に応じたテキスト解析処理により、フルコンテキストラベルデータを取得し、取得したフルコンテキストラベルデータから sequence-to-sequence 方式を用いたニューラルネットワークのモデルで処理（学習処理、および／または、予測処理）を実行するのに適したデータである最適化フルコンテキストラベルデータを取得する。そして、音声合成処理装置 2 0 0 では、入力を最適化フルコンテキストラベルデータとし、出力をメルスペクトログラム（音響特徴量の一例）として、エンコーダ部 3、および、デコーダ部 5 A において、ニューラルネットワークのモデルを用いた処理（学習処理、予測処理）を実行することで、高精度な処理を実現できる。さらに、音声合成処理装置 2 0 0 では、ボコーダ 6 によ

り、上記により取得したメルスペクトログラム（音響特徴量の一例）から、当該メルスペクトログラムに対応する音声信号波形データを取得し、取得したデータを出力することで、音声波形データ（データDout）を取得する。これにより、音声合成処理装置200では、入力されたテキストに相当する音声波形データを取得することができる。

[0224] さらに、音声合成処理装置200では、エンコーダ部3への入力データ（最適化フルコンテキストラベルデータ）を、音素継続長推定部7により取得（推定）した音素ごとの音素継続長に基づいて、引き延ばす処理（音素 ph_k の音素継続長 $dur(ph_k)$ に相当する期間、音素 ph_k の最適化フルコンテキストラベルデータを、繰り返しエンコーダ部3に入力する処理）を実行する。つまり、音声合成処理装置200では、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理を実行して取得した音素継続長を用いて予測処理を実行するので、注意機構予測が失敗することに起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することはない。

[0225] すなわち、音声合成処理装置200では、（1）音素継続長については、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理（音素継続長推定部7による処理）により取得し、（2）音響特徴量については、sequence-to-sequence方式を用いたニューラルネットワークのモデルで処理することにより取得する。

[0226] したがって、音声合成処理装置200では、注意機構予測が失敗することに起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することを適切に防止するとともに、高精度な音声合成処理を実行することができる。

[0227] [第3実施形態]

次に、第3実施形態について、説明する。なお、上記実施形態と同様の部分については、同一符号を付し、詳細な説明を省略する。

[0228] <3. 1：音声合成処理装置の構成>

図9は、第3実施形態に係る音声合成処理装置300の概略構成図である。

- [0229] 第3実施形態に係る音声合成処理装置300は、第1実施形態の音声合成処理装置100において、テキスト解析部1をテキスト解析部1Aに置換し、アテンション部4をアテンション部4Aに置換し、デコーダ部5をデコーダ部5Bに置換した構成を有している。そして、音声合成処理装置300は、音声合成処理装置100において、音素継続長推定部7と、強制アテンション部8と、内分処理部9と、コンテキスト算出部10とを追加した構成を有している。
- [0230] テキスト解析部1A、および、音素継続長推定部7は、第2実施形態のテキスト解析部1Aと同様の構成、機能を有している。
- [0231] なお、音素継続長推定部7は、音素継続長推定処理により取得（推定）した音素継続長のデータをデータ $D \times 02$ として、強制アテンション部8に出力する。
- [0232] アテンション部4Aは、エンコーダ部3から出力されるデータ $D \times 4$ と、デコーダ部5Bのデコーダ側LSTM層52Bから出力されるデータ h_o （出力側隠れ状態データ h_o ）とを入力する。アテンション部4Aは、エンコーダ部3から出力されるデータ $D \times 4$ 、すなわち、入力側隠れ状態データ h_i を所定の時間ステップ分記憶保持する。時間ステップ $t = 1$ から $t = S$ （ S ：自然数）の期間において、エンコーダ部3により取得され、アテンション部4Aに出力されたデータ $D \times 4$ （ $= h_i$ ）の集合を、 $h_{i,1 \dots s}$ と表記する。つまり、アテンション部4Aは、下記に相当するデータ $h_{i,1 \dots s}$ を記憶保持する。
- $$h_{i,1 \dots s} = \{ D \times 4 (1), D \times 4 (2), \dots, D \times 4 (S) \}$$
- また、アテンション部4Aは、デコーダ部5Bのデコーダ側LSTM層52Bから出力されるデータ $D_y 3$ 、すなわち、出力側隠れ状態データ h_o を所定の時間ステップ分記憶保持する。時間ステップ $t = 1$ から $t = T$ （ T ：自然数）の期間において、デコーダ側LSTM層52Bにより取得され、ア

テンション部4 Aに出力されたデータ $Dy3$ ($=ho$)の集合を、 $ho_{1...T}$ と表記する。つまり、アテンション部4 Aは、下記に相当するデータ $ho_{1...T}$ を記憶保持する。

$$ho_{1...T} = \{Dy3(1), Dy3(2), \dots, Dy3(T)\}$$

そして、アテンション部4 Aは、入力側隠れ状態データの集合データ $hi_{1...s}$ と、出力側隠れ状態データの集合データ $ho_{1...T}$ と、に基づいて、例えば、

$$w_{att}(t)_{1...s} = f2_attn(hi_{1...s}, ho_{1...T})$$

$f2_attn()$: 重み付け係数データを取得する関数

に相当する処理を実行して、現時刻 t の重み付け係数データ $w_{att}(t)_{1...s}$ を取得する。そして、アテンション部4 Aは、取得した重み付け係数データ $w_{att}(t)_{1...s}$ を内分処理部9に出力する。なお、入力側隠れ状態データの集合データ $hi_{1...s}$ の各要素データに対する重み付け係数データの集合データを重み付け係数データ $w_{att}(t)_{1...s}$ と表記する。

[0233] また、アテンション部4 Aは、データ $Dx4$ ($=hi$)の集合データ $hi_{1...s}$ をコンテキスト算出部10に出力する。

[0234] 強制アテンション部8は、音素継続長推定部7から出力される推定された音素継続長のデータ $Dx02$ を入力する。強制アテンション部8は、音素継続長データ $Dx02$ に対応する音素についてのエンコーダ部3により処理されたデータが出力されるとき、当該音素の推定された音素継続長(音素継続長データ $Dx02$)に相当する期間、重み付け係数を強制的に所定の値(例えば、「1」)にした重み付け係数データ $w_f(t)$ を生成する。なお、入力側隠れ状態データの集合データ $hi_{1...s}$ の各要素データに対する重み付け係数データと対応づけるために、時刻 t を中心として、 S 個にデータを拡張(同一データを複製して拡張)した重み付け係数データ $w_f(t)$ を重み付け係数データ $w_f(t)_{1...s}$ と表記する。

[0235] 強制アテンション部8は、上記により生成した重み付け係数データ $w_f(t)_{1...s}$ を内分処理部9に出力する。

[0236] 内分処理部9は、アテンション部4Aから出力される重み付け係数データ $w_{att}(t)_{1...s}$ と、強制アテンション部8から出力される重み付け係数データ $w_f(t)_{1...s}$ とを入力する。そして、内分処理部9は、重み付け係数データ $w_{att}(t)_{1...s}$ と、重み付け係数データ $w_f(t)_{1...s}$ とに対して、内分処理を実行することで、合成重み付け係数データ $w(t)$ を取得する。具体的には、内分処理部9は、

$$w(t)_{1...s} = (1 - \alpha) \times w_{att}(t)_{1...s} + \alpha \times w_f(t)_{1...s}$$

$$0 \leq \alpha \leq 1$$

に相当する処理を実行することで、合成重み付け係数データ $w(t)$ を取得する。なお、上記数式（内分処理）は、それぞれ対応する要素ごとに、内分処理を実行することを表している。つまり、 j 番目 ($1 \leq j \leq S$) のデータについては、

$$w(t)_j = (1 - \alpha) \times w_{att}(t)_j + \alpha \times w_f(t)_j$$

に相当する処理が実行されることで、 j 番目の合成重み付け係数データ $w(t)_j$ が取得される。

[0237] そして、内分処理部9は、取得した合成重み付け係数データ $w(t)_{1...s}$ をコンテキスト算出部10に出力する。

[0238] コンテキスト算出部10は、アテンション部4Aから出力されるデータ $D \times 4 (= h_i)$ の集合データ $h_i_{1...s}$ と、内分処理部9から出力される合成重み付け係数データ $w(t)_{1...s}$ とを入力する。そして、コンテキスト算出部10は、合成重み付け係数データ $w(t)_{1...s}$ に基づいて、データ $D \times 4 (= h_i)$ の集合データ $h_i_{1...s}$ に対して、重み付け加算処理を実行することで、コンテキスト状態データ $c(t)$ を取得する。そして、コンテキスト算出部10は、取得したコンテキスト状態データ $c(t)$ をデコーダ部5Bのデコーダ側LSTM層52Bに出力する。

[0239] デコーダ部5Bは、第1実施形態のデコーダ部5において、デコーダ側LSTM層52をデコーダ側LSTM層52Bに置換した構成を有している。

それ以外は、デコーダ部5Bは、第1実施形態のデコーダ部5と同様である。

[0240] デコーダ側LSTM層52Bは、デコーダ側LSTM層52と同様の機能を有している。デコーダ側LSTM層52Bは、デコーダ側プレネット処理部51から、現時刻 t において出力されるデータ $Dy2$ （これをデータ $Dy2(t)$ と表記する）と、1つ前の時間ステップにおいて、デコーダ側LSTM層52Bから出力されたデータ $Dy3$ （これをデータ $Dy3(t-1)$ と表記する）と、コンテキスト算出部10から出力される時刻 t のコンテキスト状態データ $c(t)$ とを入力する。

[0241] デコーダ側LSTM層52Bは、入力されたデータ $Dy2(t)$ 、データ $Dy3(t-1)$ 、および、コンテキスト状態データ $c(t)$ を用いて、LSTM層による処理を実行し、処理後のデータをデータ $Dy3$ （データ $Dy3(t)$ ）として線形予測部53に出力する。また、デコーダ側LSTM層52Bは、データ $Dy3(t)$ 、すなわち、時刻 t の出力側隠れ状態データ $h_o(t)$ をアテンション部4Aに出力する。

[0242] <3.2: 音声合成処理装置の動作>

以上のように構成された音声合成処理装置300の動作について以下説明する。

[0243] 図10～図12は、アテンション部4Aにより取得された重み付け係数データ $w_{att}(t)$ と、強制アテンション部8により取得された重み付け係数データ $w_f(t)$ とから取得した合成重み付け係数データ $w(t)$ を用いてコンテキスト状態データ $c(t)$ を取得する処理について説明するための図である。

[0244] (3.2.1: 学習処理)

まず、音声合成処理装置300による学習処理について、説明する。なお、説明便宜のため、処理対象言語を日本語として、以下、説明する。

[0245] 処理対象言語である日本語のテキストデータ Din をテキスト解析部1Aに入力する。また、当該テキストデータ Din に対応するメルスペクトログ

ラム（音響特徴量）のデータを教師データとして用意する。

- [0246] テキスト解析部 1 A は、第 1 実施形態と同様に、入力されたテキストデータ D_{in} に対して、テキスト解析処理を実行し、様々な言語情報からなるコンテキストを含む音素ラベルであるコンテキストラベルの系列を取得する。
- [0247] テキスト解析部 1 A は、第 1 実施形態と同様に、取得したフルコンテキストラベルデータをデータ $D \times 1$ としてフルコンテキストラベルベクトル処理部 2 へ出力する。
- [0248] また、テキスト解析部 1 A は、処理対象言語のテキストデータ D_{in} から音素のコンテキストラベルを取得し、取得した音素のコンテキストラベルのデータをデータ $D \times 0 1$ として、音素継続長推定部 7 へ出力する。
- [0249] 音素継続長推定部 7 は、テキスト解析部 1 A から出力されるデータ $D \times 0 1$ （音素のコンテキストラベルのデータ）から、データ $D \times 0 1$ に対応する音素の音素継続長を推定（取得）する音素継続長推定処理を実行する。具体的には、音素継続長推定部 7 は、例えば、隠れマルコフモデル（HMM：Hidden Markov Model）、ニューラルネットワークモデル等を用いた、音素のコンテキストラベルから当該音素の音素継続長を推定（予測）するモデル（処理システム）により、音素継続長推定処理を実行する。
- [0250] そして、音素継続長推定部 7 は、音素継続長推定処理により取得（推定）した音素継続長のデータをデータ $D \times 0 2$ として、強制アテンション部 8 へ出力する。
- [0251] フルコンテキストラベルベクトル処理部 2 A は、テキスト解析部 1 A から出力されるデータ $D \times 1$ （フルコンテキストラベルのデータ）から、sequence-to-sequence方式のニューラルネットワークのモデルの学習処理に適したフルコンテキストラベルデータを取得するためのフルコンテキストラベルベクトル処理（第 1 実施形態と同様のフルコンテキストラベルベクトル処理）を実行する。そして、フルコンテキストラベルベクトル処理部 2 A は、フルコンテキストラベルベクトル処理により取得したデータをデータ $D \times 2$ （最適

化フルコンテキストラベルデータ $D \times 2$) として、エンコーダ部 3 のエンコーダ側プレネット処理部 3 1 に出力する。

[0252] フルコンテキストラベルベクトル処理部 2 A により取得されたデータ $D \times 2$ (最適化フルコンテキストラベルデータ $D \times 2$) は、フルコンテキストラベルベクトル処理部 2 からのエンコーダ部 3 のエンコーダ側プレネット処理部 3 1 に出力される。

[0253] エンコーダ側プレネット処理部 3 1 は、フルコンテキストラベルベクトル処理部 2 から入力したデータ $D \times 2$ に対して、コンボリユーション処理 (コンボリユーションフィルタによる処理)、データの正規化処理、活性化関数による処理 (例えば、ReLU 関数 (ReLU: Rectified Linear Unit) による処理) を実行し、エンコーダ側 LSTM 層 3 2 に入力可能なデータを取得する。そして、エンコーダ側プレネット処理部 3 1 は、上記処理 (プレネット処理) により取得したデータをデータ $D \times 3$ としてエンコーダ側 LSTM 層 3 2 に出力する。

[0254] エンコーダ側 LSTM 層 3 2 は、エンコーダ側プレネット処理部 3 1 から、現時刻 t において出力されるデータ $D \times 3 (t)$ と、1 つ前の時間ステップにおいて、エンコーダ側 LSTM 層 3 2 から出力されたデータ $D \times 4 (t-1)$ とを入力する。そして、エンコーダ側 LSTM 層 3 2 は、入力されたデータ $D \times 3 (t)$ 、データ $D \times 4 (t-1)$ に対して、LSTM 層による処理を実行し、処理後のデータをデータ $D \times 4$ (データ $D \times 4 (t)$ (=入力側隠れ状態データ $h_i (t)$)) としてアテンション部 4 A に出力する。

[0255] アテンション部 4 A は、エンコーダ部 3 から出力されるデータ $D \times 4$ と、デコーダ部 5 B のデコーダ側 LSTM 層 5 2 B から出力されるデータ h_o (出力側隠れ状態データ h_o) とを入力する。アテンション部 4 A は、エンコーダ部 3 から出力されるデータ $D \times 4$ 、すなわち、入力側隠れ状態データ h_i を所定の時間ステップ分記憶保持する。例えば、アテンション部 4 A は、時間ステップ $t = 1$ から $t = S$ (S : 自然数) の期間において、エンコーダ部 3 により取得され、アテンション部 4 A に出力されたデータ $D \times 4 (= h_i)$

の集合を、 $h_{i_{1...s}} (= \{D \times 4 (1), D \times 4 (2), \dots, D \times 4 (S)\})$ として記憶保持する。

[0256] また、アテンション部4 Aは、デコーダ部5 Bのデコーダ側LSTM層5 2 Bから出力されるデータ D_{y3} 、すなわち、出力側隠れ状態データ h_o を所定の時間ステップ分記憶保持する。例えば、アテンション部4 Aは、時間ステップ $t = 1$ から $t = T$ (T :自然数)の期間において、デコーダ側LSTM層5 2 Bにより取得され、アテンション部4 Aに出力されたデータ D_{y3} ($= h_o$)の集合を、 $h_{o_{1...T}} (= \{D_{y3} (1), D_{y3} (2), \dots, D_{y3} (T)\})$ として記憶保持する。

[0257] そして、アテンション部4 Aは、入力側隠れ状態データの集合データ $h_{i_{1...s}}$ と、出力側隠れ状態データの集合データ $h_{o_{1...T}}$ と、に基づいて、例えば、

$$w_{att} (t)_{1...s} = f_{2_attn} (h_{i_{1...s}}, h_{o_{1...T}})$$

$f_{2_attn} ()$: 重み付け係数データを取得する関数

に相当する処理を実行して、現時刻 t の重み付け係数データ $w_{att} (t)_{1...s}$ を取得する。

[0258] そして、アテンション部4 Aは、取得した重み付け係数データ $w_{att} (t)_{1...s}$ を内分処理部9に出力する。また、アテンション部4 Aは、データ $D \times 4 (= h_i)$ の集合データ $h_{i_{1...s}}$ をコンテキスト算出部10に出力する。

[0259] 強制アテンション部8は、音素継続長データ $D \times 0 2$ に対応する音素についてのエンコーダ部3により処理されたデータが出力されるとき、当該音素の推定された音素継続長(音素継続長データ $D \times 0 2$)に相当する期間、重み付け係数を強制的に所定の値(例えば、「1」)にした重み付け係数データ $w_f (t)$ を生成する。そして、強制アテンション部8は、入力側隠れ状態データの集合データ $h_{i_{1...s}}$ の各要素データに対する重み付け係数データと対応づけるために(内分処理ができるようにするために)、時刻 t を中心として、 S 個にデータを拡張(同一データを複製して拡張)した重み付け係数データ $w_f (t)_{1...s}$ を生成する。

[0260] 強制アテンション部 8 は、上記により生成した重み付け係数データ $w_f(t)_{1..s}$ を内分処理部 9 に出力する。

[0261] 内分処理部 9 は、アテンション部 4 A から出力される重み付け係数データ $w_{att}(t)_{1..s}$ と、強制アテンション部 8 から出力される重み付け係数データ $w_f(t)_{1..s}$ とを入力する。そして、内分処理部 9 は、重み付け係数データ $w_{att}(t)_{1..s}$ と、重み付け係数データ $w_f(t)_{1..s}$ とに対して、内分処理を実行することで、合成重み付け係数データ $w(t)$ を取得する。具体的には、内分処理部 9 は、

$$w(t)_{1..s} = (1 - \alpha) \times w_{att}(t)_{1..s} + \alpha \times w_f(t)_{1..s}$$

$$0 \leq \alpha \leq 1$$

に相当する処理を実行することで、合成重み付け係数データ $w(t)_{1..s}$ を取得する。そして、内分処理部 9 は、取得した合成重み付け係数データ $w(t)_{1..s}$ をコンテキスト算出部 10 に出力する。

[0262] コンテキスト算出部 10 は、アテンション部 4 A から出力されるデータ $D \times 4 (= h_i)$ の集合データ $h_i_{1..s}$ と、内分処理部 9 から出力される合成重み付け係数データ $w(t)_{1..s}$ とを入力する。そして、コンテキスト算出部 10 は、合成重み付け係数データ $w(t)_{1..s}$ に基づいて、データ $D \times 4 (= h_i)$ の集合データ $h_i_{1..s}$ に対して、重み付け加算処理を実行することで、コンテキスト状態データ $c(t)$ を取得する。そして、コンテキスト算出部 10 は、取得したコンテキスト状態データ $c(t)$ をデコーダ部 5 B のデコーダ側 LSTM 層 52 B に出力する。

[0263] なお、学習処理時において、内分比 α を「0」に固定してもよい。この場合（内分比 α を「0」に固定した場合）、音声合成処理装置 300 では、第 1 実施形態と同様の構成により学習処理が実行されることになる。また、学習処理時において、内分比 α を所定の値（例えば、0.5）に固定して、音声合成処理装置 300 において、学習処理を実行してもよい。

[0264] ここで、学習処理時において、内分比 α を所定の値に固定する場合につい

て、図10～図12を用いて説明する。なお、説明便宜のため、内分比 α を「0.5」に固定する場合について、説明する。以下では、(1)音素に対応する音声が出力される期間内の処理(図11の場合)と、(2)無音状態である期間内の処理(図12の場合)とについて説明する。

[0265] まず、「(1)音素に対応する音声が出力される期間内の処理(図11の場合)」について、説明する。

[0266] 例えば、図10に示すように、入力データ D_{in} が「今日の天気は、...」である場合、データ D_{x01} に含まれる各音素のデータを、

(1) ph_0 = 「k」、(2) ph_1 = 「y」、(3) ph_2 = 「ou」、(4) ph_3 = 「n」、(5) ph_{04} = 「o」、(6) ph_{sil} = 無音状態、(7) ph_5 = 「t」、(8) ph_6 = 「e」、(9) ph_{07} = 「n」、... とし、音素 ph_k (k :整数)の推定された音素継続長を $dur(ph_k)$ とすると、音素継続長推定部7は、音素 ph_k (k :整数)のコンテキストラベルを用いて、音素継続長推定処理を実行することで、音素 ph_k の推定された音素継続長 $dur(ph_k)$ を取得する。例えば、上記の各音素(音素 ph_k)について、音素継続長推定部7により取得(推定)された音素継続長 $dur(ph_k)$ が、図10に示す時間の長さ(継続長)を有するものとする。

[0267] 強制アテンション部8は、音素継続長データ D_{x02} に対応する音素についてのエンコーダ部3により処理されたデータが出力されるとき、当該音素の推定された音素継続長(音素継続長データ D_{x02})に相当する期間、重み付け係数を強制的に所定の値(例えば、「1」)にした重み付け係数データ $w_f(t)$ を生成する。図10の場合、強制アテンション部8は、音素 ph_k についてのエンコーダ部3により処理されたデータが出力されるとき、音素 ph_k の音素継続長 $dur(ph_k)$ に相当する期間、重み付け係数を強制的に所定の値(例えば、「1」)にした重み付け係数データ $w_f(t)$ を内分処理部9に出力し続ける(図10において、 $w_f(t)[ph_k]$ と表記した部分に相当)。

[0268] また、図10において、処理対象の音素に対応付けて、アテンション部4

Aにより取得された重み付け係数データ $w_{att}(t)$ を示している。具体的には、図10において、音素 ph_k に対応する、アテンション部4Aにより取得された重み付け係数データ $w_{att}(t)$ が出力される期間を「 $w_{att}(t)[ph_k]$ 」として示している。なお、説明便宜のため、図10では、アテンション部4Aによる音素継続長の予測が正しくなされた場合を示している。

[0269] また、図10において、音素 ph_k に対応する合成重み付け係数データ $w(t)$ を「 $w(t)[ph_k]$ 」として示している。

[0270] 図11は、時刻 t_2 （時間ステップ t_2 ）における処理を説明するための図であり、図10において処理対象音素が「ou」であるときの期間の一部を時間軸方向に拡大して示した図である。なお、説明便宜のため、音声合成処理装置300において、データ $D \times 4 (=hi)$ の集合データ $hi_{1...s}$ は、9個のデータ（すなわち、 $S=9$ ）（図11において、期間 $T(t_2)$ において取得され、記憶保持されているデータ）であるものとする（以下、同様）。

[0271] ここで、時刻 t_2 における処理について、説明する。

[0272] 強制アテンション部8は、時刻 t_2 において、音素継続長 $D \times 0_2$ から、音素「ou」に相当する音声出力継続される期間であることを認識し、時刻 t_2 の重み付け係数データ $w_f(t)$ を「1」に設定する。さらに、強制アテンション部8は、入力側隠れ状態データの集合データ $hi_{1...s}$ の各要素データに対する重み付け係数データと対応づけるために（内分処理ができるようにするために）、時刻 t_2 を中心として、 $S (=9)$ 個にデータを拡張（同一データを複製して拡張）した重み付け係数データ $w_f(t)_{1...s}$ を生成する。なお、 $w_f(t)_{1...s}$ は、

$$w_f(t)_{1...s} = \{w_{01}, w_{02}, w_{03}, w_{04}, w_{05}, w_{06}, w_{07}, w_{08}, w_{09}\}$$

$$0 \leq w_{0j} \leq 1 \quad (1 \leq j \leq S)$$

$$t = t_2$$

であるものとし、 $w_f(t_2)_{1...s}$ において、 $w_{01} \sim w_{09}$ は、すべて「1」

に設定される（図11参照）。

[0273] 強制アテンション部8は、上記により生成した重み付け係数データ $w_f(t_2)_{1..s}$ を内分処理部9に出力する。

[0274] アテンション部4Aは、入力側隠れ状態データの集合データ $h_{i_{1..s}}$ と、出力側隠れ状態データの集合データ $h_{o_{1..T}}$ と、に基づいて、例えば、

$$w_{att}(t)_{1..s} = f2_attn(h_{i_{1..s}}, h_{o_{1..T}})$$

$f2_attn()$: 重み付け係数データを取得する関数

に相当する処理を実行して、時刻 t_2 の重み付け係数データ $w_{att}(t_2)_{1..s}$ を取得する。時刻 t_2 の重み付け係数データ $w_{att}(t)_{1..s}$ が図11に示すデータ（一例）であるものとする。なお、 $w_{att}(t)_{1..s}$ は、

$$w_{att}(t)_{1..s} = \{w_{11}, w_{12}, w_{13}, w_{14}, w_{15}, w_{16}, w_{17}, w_{18}, w_{19}\}$$

$$0 \leq w_{1j} \leq 1 \quad (1 \leq j \leq S)$$

$$t = t_2$$

であるものとし、 $w_{11} \sim w_{19}$ は、例えば、アテンション部4Aにより、以下の値として、取得されたものとする（図11参照）。

$$w_{11} = 0.0, w_{12} = 0.2, w_{13} = 0.4, w_{14} = 0.8, w_{15} = 1.0$$

$$w_{16} = 0.8, w_{17} = 0.4, w_{18} = 0.2, w_{19} = 0.0$$

アテンション部4Aは、上記により取得された重み付け係数データ $w_{att}(t_2)_{1..s}$ を内分処理部9に出力する。

[0275] 内分処理部9は、アテンション部4Aから出力される重み付け係数データ $w_{att}(t_2)_{1..s}$ と、強制アテンション部8から出力される重み付け係数データ $w_f(t_2)_{1..s}$ とを入力する。そして、内分処理部9は、重み付け係数データ $w_{att}(t_2)_{1..s}$ と、重み付け係数データ $w_f(t_2)_{1..s}$ とに対して、内分処理を実行することで、合成重み付け係数データ $w(t_2)_{1..s}$ を取得する。具体的には、内分処理部9は、

$$w(t_2)_{1..s} = (1 - \alpha) \times w_{att}(t_2)_{1..s} + \alpha \times w_f(t_2)_{1..s}$$

) $1 \dots s$

$$0 \leq \alpha \leq 1$$

に相当する処理を実行することで、合成重み付け係数データ $w(t_2)_{1 \dots s}$ を取得する。

[0276] ここでは、 $\alpha = 0.5$ であるので、 $w_{att}(t_2)_{1 \dots s}$ と、 $w_f(t_2)_{1 \dots s}$ との平均値が合成重み付け係数データ $w(t)_{1 \dots s}$ となる。なお、 $w(t)_{1 \dots s}$ は、

$$w(t)_{1 \dots s} = \{w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8, w_9\}$$

$$0 \leq w_{1j} \leq 1 \quad (1 \leq j \leq S)$$

$$t = t_2$$

であるものとする、 $w_1 \sim w_9$ は、内分処理部 9 により、以下の値として、取得される (図 11 参照)。

$$w_1 = 0.5 \times w_{01} + 0.5 \times w_{11} = 0.5 + 0 = 0.5$$

$$w_2 = 0.5 \times w_{02} + 0.5 \times w_{12} = 0.5 + 0.1 = 0.6$$

$$w_3 = 0.5 \times w_{03} + 0.5 \times w_{13} = 0.5 + 0.2 = 0.7$$

$$w_4 = 0.5 \times w_{04} + 0.5 \times w_{14} = 0.5 + 0.4 = 0.9$$

$$w_5 = 0.5 \times w_{05} + 0.5 \times w_{15} = 0.5 + 0.5 = 1.0$$

$$w_6 = 0.5 \times w_{06} + 0.5 \times w_{16} = 0.5 + 0.4 = 0.9$$

$$w_7 = 0.5 \times w_{07} + 0.5 \times w_{17} = 0.5 + 0.2 = 0.7$$

$$w_8 = 0.5 \times w_{08} + 0.5 \times w_{18} = 0.5 + 0.1 = 0.6$$

$$w_9 = 0.5 \times w_{09} + 0.5 \times w_{19} = 0.5 + 0 = 0.5$$

そして、内分処理部 9 は、取得した合成重み付け係数データ $w(t_2)_{1 \dots s}$ をコンテキスト算出部 10 に出力する。

[0277] コンテキスト算出部 10 は、アテンション部 4 A から出力されるデータ $D \times 4 (= h_i)$ の集合データ $h_i_{1 \dots s}$ と、内分処理部 9 から出力される合成重み付け係数データ $w(t_2)_{1 \dots s}$ とを入力する。そして、コンテキスト算出部 10 は、合成重み付け係数データ $w(t_2)_{1 \dots s}$ に基づいて、データ $D \times 4 (= h_i)$ の集合データ $h_i_{1 \dots s}$ に対して、重み付け加算処理を実行す

ることで、コンテキスト状態データ $c(t)$ を取得する。つまり、コンテキスト算出部 10 は、以下の数式に相当する処理を実行することで、コンテキスト状態データ $c(t)$ を取得する。

[数3]

$$c(t) = \sum_{j=1}^S (w_j \times hi_j)$$

$t = t_2$

w_j : 合成重み付け係数データ $w(t_2)_{1..s}$ の j 番目の要素データ ($1 \leq j \leq S$)

そして、コンテキスト算出部 10 は、取得したコンテキスト状態データ $c(t_2)$ をデコーダ部 5B のデコーダ側 LSTM 層 52B に出力する。

[0278] 次に、「(2) 無音状態である期間内の処理 (図 12 の場合)」について、説明する。

[0279] 図 12 は、時刻 t_3 (時間ステップ t_3) における処理を説明するための図であり、図 10 において無音状態の期間 (図 10 において、「silent (無音状態)」で示した期間) の一部を時間軸方向に拡大して示した図である。

[0280] ここで、時刻 t_3 における処理について、説明する。

[0281] 強制アテンション部 8 は、時刻 t_3 において、音素継続長 $D \times 0_2$ から、無音状態 (発声すべき音素がない状態) の期間であることを認識し、時刻 t_3 の重み付け係数データ $w_f(t)$ を「0」に設定する。さらに、強制アテンション部 8 は、入力側隠れ状態データの集合データ $hi_{1..s}$ の各要素データに対する重み付け係数データと対応づけるために (内分処理ができるようにするために)、時刻 t_2 を中心として、 $S (= 9)$ 個にデータを拡張 (同一データを複製して拡張) した重み付け係数データ $w_f(t)_{1..s}$ を生成する

。なお、 $w_f(t)_{1..s}$ は、

$$w_f(t)_{1..s} = \{w_{01}, w_{02}, w_{03}, w_{04}, w_{05}, w_{06}, w_{07}, w_{08}, w_{09}\}$$

$$0 \leq w_{0j} \leq 1 \quad (1 \leq j \leq S)$$

$$t = t_2$$

であるものとし、 $w_f(t_3)_{1..s}$ において、 $w_{01} \sim w_{09}$ は、すべて「0」に設定される（図12参照）。

[0282] 強制アテンション部8は、上記により生成した重み付け係数データ $w_f(t_3)_{1..s}$ を内分処理部9に出力する。

[0283] アテンション部4Aは、入力側隠れ状態データの集合データ $h_i_{1..s}$ と、出力側隠れ状態データの集合データ $h_o_{1..T}$ と、に基づいて、例えば、

$$w_{att}(t)_{1..s} = f2_attn(h_i_{1..s}, h_o_{1..T})$$

$f2_attn()$ ：重み付け係数データを取得する関数

に相当する処理を実行して、時刻 t_3 の重み付け係数データ $w_{att}(t_3)_{1..s}$ を取得する。時刻 t_3 の重み付け係数データ $w_{att}(t)_{1..s}$ が図12に示すデータ（一例）であるものとする。なお、 $w_{att}(t)_{1..s}$ は、

$$w_{att}(t)_{1..s} = \{w_{11}, w_{12}, w_{13}, w_{14}, w_{15}, w_{16}, w_{17}, w_{18}, w_{19}\}$$

$$0 \leq w_{1j} \leq 1 \quad (1 \leq j \leq S)$$

$$t = t_2$$

であるものとし、 $w_{11} \sim w_{19}$ は、例えば、アテンション部4Aにより、すべて値が「0」として、取得されたものとする（図12参照）。

[0284] アテンション部4Aは、上記により取得された重み付け係数データ $w_{att}(t_3)_{1..s}$ を内分処理部9に出力する。

[0285] 内分処理部9は、アテンション部4Aから出力される重み付け係数データ $w_{att}(t_3)_{1..s}$ と、強制アテンション部8から出力される重み付け係数データ $w_f(t_3)_{1..s}$ とを入力する。そして、内分処理部9は、重み付け係数データ $w_{att}(t_3)_{1..s}$ と、重み付け係数データ $w_f(t_3)_{1..s}$

とに対して、内分処理を実行することで、合成重み付け係数データ $w(t_3)_{1..s}$ を取得する。具体的には、内分処理部 9 は、

$$w(t_3)_{1..s} = (1 - \alpha) \times w_{att}(t_3)_{1..s} + \alpha \times w_f(t_3)_{1..s}$$

$$0 \leq \alpha \leq 1$$

に相当する処理を実行することで、合成重み付け係数データ $w(t_3)_{1..s}$ を取得する。

[0286] ここでは、 $\alpha = 0.5$ であるので、 $w_{att}(t_3)_{1..s}$ と、 $w_f(t_3)_{1..s}$ との平均値が合成重み付け係数データ $w(t)_{1..s}$ となる。なお、 $w(t)_{1..s}$ は、

$$w(t)_{1..s} = \{w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8, w_9\}$$

$$0 \leq w_{1j} \leq 1 \quad (1 \leq j \leq S)$$

$$t = t_2$$

であるものとする。と、 $w_1 \sim w_9$ は、内分処理部 9 により、すべて値が「0」として、取得される (図 12 参照)。

[0287] そして、内分処理部 9 は、取得した合成重み付け係数データ $w(t_2)_{1..s}$ をコンテキスト算出部 10 に出力する。

[0288] コンテキスト算出部 10 は、アテンション部 4A から出力されるデータ $D \times 4 (= h_i)$ の集合データ $h_i_{1..s}$ と、内分処理部 9 から出力される合成重み付け係数データ $w(t_3)_{1..s}$ とを入力する。そして、コンテキスト算出部 10 は、合成重み付け係数データ $w(t_2)_{1..s}$ に基づいて、データ $D \times 4 (= h_i)$ の集合データ $h_i_{1..s}$ に対して、重み付け加算処理を実行することで、コンテキスト状態データ $c(t)$ を取得する。つまり、コンテキスト算出部 10 は、以下の数式に相当する処理を実行することで、コンテキスト状態データ $c(t)$ を取得する。

[数4]

$$c(t) = \sum_{j=1}^S (w_j \times hi_j)$$

$t = t_2$

w_j : 合成重み付け係数データ $w(t_3)_{1 \dots s}$ の j 番目の要素データ ($1 \leq j \leq S$)

そして、コンテキスト算出部 10 は、取得したコンテキスト状態データ $c(t_3)$ をデコーダ部 5 B のデコーダ側 LSTM 層 5 2 B に出力する。

[0289] 図 12 の場合、無音状態であるので、アテンション部 4 A、および、強制アテンション部 8 により取得される重み付け係数データがすべて 0 であるので、コンテキスト状態データ $c(t_3)$ も「0」となる。つまり、上記により、無音状態であることを適切に示すコンテキスト状態データ $c(t_3)$ が取得される。

[0290] 上記のように取得されたコンテキスト状態データ $c(t)$ がデコーダ部 5 B のデコーダ側 LSTM 層 5 2 B に出力される。

[0291] そして、アテンション部 4 A は、取得した重み付け係数データ $w_{att}(t)_{1 \dots s}$ を内分処理部 9 に出力する。また、アテンション部 4 A は、データ $D \times 4 (= h_i)$ の集合データ $h_{i_{1 \dots s}}$ をコンテキスト算出部 10 に出力する。

[0292] デコーダ側プレネット処理部 5 1 での処理は、第 1 実施形態と同様である。

[0293] デコーダ側 LSTM 層 5 2 B は、デコーダ側プレネット処理部 5 1 から、現時刻 t において出力されるデータ $Dy_2(t)$ と、1 つ前の時間ステップにおいて、デコーダ側 LSTM 層 5 2 から出力されたデータ $Dy_3(t-1)$ と、コンテキスト算出部 10 から出力される時刻 t のコンテキスト状態データ $c(t)$ とを入力する。

[0294] デコーダ側LSTM層52Aは、入力されたデータ $Dy_2(t)$ 、データ $Dy_3(t-1)$ 、および、コンテキスト状態データ $c(t)$ を用いて、LSTM層による処理を実行し、処理後のデータをデータ $Dy_3(t)$ として線形予測部53に出力する。

線形予測部53、ポストネット処理部54、および、加算器55では、第1実施形態と同様の処理が実行される。

[0295] そして、音声合成処理装置200では、上記のように取得されたデータ Dy_6 （予測メルスペクトログラムのデータ）と、テキストデータ Din に対応するメルスペクトログラム（音響特徴量）の教師データ（正解のメルスペクトログラム）とを比較し、両者の差（比較結果）（例えば、差分ベクトルのノルムやユークリッド距離により表現する差）が小さくなるように、エンコーダ部3、デコーダ部5Bのニューラルネットワークのモデルのパラメータを更新する。音声合成処理装置100では、このパラメータ更新処理を繰り返し実行し、データ Dy_6 （予測メルスペクトログラムのデータ）と、テキストデータ Din に対応するメルスペクトログラム（音響特徴量）の教師データ（正解のメルスペクトログラム）との差が十分小さくなる（所定の誤差範囲におさまる）、ニューラルネットワークのモデルのパラメータを最適化パラメータとして取得する。

[0296] 音声合成処理装置300では、上記のようにして取得した最適化パラメータに基づいて、エンコーダ部3、デコーダ部5Bのニューラルネットワークのモデルの各層に含まれるシナプス間の結合係数（重み係数）を設定することで、エンコーダ部3、デコーダ部5Aのニューラルネットワークのモデルを最適化モデル（学習済みモデル）とすることができる。

[0297] 以上により、音声合成処理装置300において、入力をテキストデータとし、出力をメルスペクトログラムとするニューラルネットワークの学習済みモデル（最適化モデル）を構築できる。

[0298] なお、音声合成処理装置300において、第1実施形態の音声合成処理装置100における学習処理により取得したニューラルネットワークの学習済

みモデル（最適化モデル）を用いてもよい。つまり、音声合成処理装置200において、第1実施形態の音声合成処理装置100における学習処理により取得したニューラルネットワークの学習済みモデルのエンコーダ部3およびデコーダ部5の最適パラメータを用いて、音声合成処理装置200のエンコーダ部3およびデコーダ部5Bのパラメータを設定することで、音声合成処理装置300において、学習済みモデルを構築するようによい。

[0299] また、ボコーダ6として、ニューラルネットワークによるモデルを用いたボコーダを採用する場合、その学習処理は、第1実施形態と同様である。

[0300] これにより、第1実施形態と同様に、ボコーダ6において、入力をテキストデータとし、出力をメルスペクトログラムとするニューラルネットワークの学習済みモデル（最適化モデル）を構築できる。

[0301] なお、音声合成処理装置300において、（1）エンコーダ部3、デコーダ部5Bの学習処理と、（2）ボコーダ6の学習処理とを連携させて学習処理を実行してもよいし、上記のように、個別に学習処理を実行してもよい。音声合成処理装置300において、（1）エンコーダ部3、デコーダ部5Bの学習処理と、（2）ボコーダ6の学習処理とを連携させて学習処理を実行する場合、入力をテキストデータとし、当該テキストデータに対応する音声波形データ（正解の音声波形データ）とを用いて、（1）エンコーダ部3、デコーダ部5Bのニューラルネットワークのモデルと、（2）ボコーダ6のニューラルネットワークのモデルの最適化パラメータを取得することで学習処理を実行すればよい。

[0302] （3. 2. 2：予測処理）

次に、音声合成処理装置300による予測処理について、説明する。なお、予測処理においても、説明便宜のため、処理対象言語を日本語として、以下、説明する。

[0303] 予測処理を実行する場合、音声合成処理装置300では、上記の学習処理により取得された学習済みモデル、すなわち、エンコーダ部3、デコーダ部5Bのニューラルネットワークの最適化モデル（最適化パラメータが設定さ

れているモデル)、および、ボコーダ6のニューラルネットワークの最適化モデル(最適化パラメータが設定されているモデル)が構築されている。そして、音声合成処理装置300では、当該学習済みモデルを用いて予測処理が実行される。

[0304] 音声合成処理の対象とする日本語のテキストデータ D_{in} をテキスト解析部1Aに入力する。

[0305] テキスト解析部1Aは、入力されたテキストデータ D_{in} に対して、日本語用のテキスト解析処理を実行し、例えば、図2に示すパラメータを含む478次元のベクトルデータとして、フルコンテキストラベルデータ D_{x1} を取得する。

[0306] そして、取得されたフルコンテキストラベルデータ D_{x1} は、テキスト解析部1Aからフルコンテキストラベルベクトル処理部2に出力される。

[0307] また、テキスト解析部1Aは、処理対象言語のテキストデータ D_{in} から音素のコンテキストラベルを取得し、取得した音素のコンテキストラベルのデータをデータ D_{x01} として、音素継続長推定部7に出力する。

[0308] 音素継続長推定部7は、テキスト解析部1Aから出力されるデータ D_{x01} (音素のコンテキストラベルのデータ)から、データ D_{x01} に対応する音素の音素継続長を推定(取得)する音素継続長推定処理を実行する。具体的には、音素継続長推定部7は、例えば、隠れマルコフモデル(HMM: Hidden Markov Model)、ニューラルネットワークモデル等を用いた、音素のコンテキストラベルから当該音素の音素継続長を推定(予測)するモデル(処理システム)により、音素継続長推定処理を実行する。

[0309] 例えば、図8に示すように、入力データ D_{in} が「今日の天気は、...」である場合、データ D_{x01} に含まれる各音素のデータを、

(1) $ph_0 = 「k」$ 、(2) $ph_1 = 「y」$ 、(3) $ph_2 = 「ou」$ 、(4) $ph_3 = 「n」$ 、(5) $ph_{04} = 「o」$ 、(6) $ph_{s_{i1}} = \text{無音状態}$ 、(7) $ph_5 = 「t」$ 、(8) $ph_6 = 「e」$ 、(9) $ph_{07} = 「n」$ 、...

とし、音素 $p h_k$ (k : 整数) の推定された音素継続長を $d u r (p h_k)$ とすると、音素継続長推定部 7 は、音素 $p h_k$ (k : 整数) のコンテキストレベルを用いて、音素継続長推定処理を実行することで、音素 $p h_k$ の推定された音素継続長 $d u r (p h_k)$ を取得する。例えば、上記の各音素 (音素 $p h_k$) について、音素継続長推定部 7 により取得 (推定) された音素継続長 $d u r (p h_k)$ が、図 8 に示す時間の長さ (継続長) を有するものとする。

[0310] そして、音素継続長推定部 7 は、音素継続長推定処理により取得 (推定) した音素継続長のデータ (図 8 の場合、 $d u r (p h_k)$) をデータ $D \times 0 2$ として、強制アテンション部 8 に出力する。

[0311] 強制アテンション部 8 は、音素継続長データ $D \times 0 2$ に対応する音素についてのエンコーダ部 3 により処理されたデータが出力される時、当該音素の推定された音素継続長 (音素継続長データ $D \times 0 2$) に相当する期間、重み付け係数を強制的に所定の値 (例えば、「1」) にした重み付け係数データ $w_f (t)$ を生成する。そして、強制アテンション部 8 は、入力側隠れ状態データの集合データ $h_{i_1, \dots, s}$ の各要素データに対する重み付け係数データと対応づけるために (内分処理ができるようにするために)、時刻 t を中心として、 S 個にデータを拡張 (同一データを複製して拡張) した重み付け係数データ $w_f (t)_{1, \dots, s}$ を生成する。

[0312] 強制アテンション部 8 は、上記により生成した重み付け係数データ $w_f (t)_{1, \dots, s}$ を内分処理部 9 に出力する。

[0313] エンコーダ部 3 では、第 1 実施形態と同様の処理が実行される。

[0314] アテンション部 4 A は、エンコーダ部 3 から出力されるデータ $D \times 4$ と、デコーダ部 5 B のデコーダ側 LSTM 層 5 2 B から出力されるデータ h_o (出力側隠れ状態データ h_o) とを入力する。アテンション部 4 A は、エンコーダ部 3 から出力されるデータ $D \times 4$ 、すなわち、入力側隠れ状態データ h_i を所定の時間ステップ分記憶保持する。例えば、アテンション部 4 A は、時間ステップ $t = 1$ から $t = S$ (S : 自然数) の期間において、エンコーダ部 3 により取得され、アテンション部 4 A に出力されたデータ $D \times 4$ ($= h_i$)

の集合を、 $h_{i_{1...s}} (= \{D \times 4 (1), D \times 4 (2), \dots, D \times 4 (S)\})$ として記憶保持する。

[0315] また、アテンション部4 Aは、デコーダ部5 Bのデコーダ側LSTM層5 2 Bから出力されるデータ D_{y3} 、すなわち、出力側隠れ状態データ h_o を所定の時間ステップ分記憶保持する。例えば、アテンション部4 Aは、時間ステップ $t = 1$ から $t = T$ (T :自然数)の期間において、デコーダ側LSTM層5 2 Bにより取得され、アテンション部4 Aに出力されたデータ D_{y3} ($= h_o$)の集合を、 $h_{o_{1...T}} (= \{D_{y3} (1), D_{y3} (2), \dots, D_{y3} (T)\})$ として記憶保持する。

[0316] そして、アテンション部4 Aは、入力側隠れ状態データの集合データ $h_{i_{1...s}}$ と、出力側隠れ状態データの集合データ $h_{o_{1...T}}$ と、に基づいて、例えば、

$$w_{att} (t)_{1...s} = f_{2_attn} (h_{i_{1...s}}, h_{o_{1...T}})$$

$f_{2_attn} ()$: 重み付け係数データを取得する関数

に相当する処理を実行して、現時刻 t の重み付け係数データ $w_{att} (t)_{1...s}$ を取得する。

[0317] そして、アテンション部4 Aは、取得した重み付け係数データ $w_{att} (t)_{1...s}$ を内分処理部9に出力する。また、アテンション部4 Aは、データ $D \times 4 (= h_i)$ の集合データ $h_{i_{1...s}}$ をコンテキスト算出部10に出力する。

[0318] 内分処理部9は、アテンション部4 Aから出力される重み付け係数データ $w_{att} (t)_{1...s}$ と、強制アテンション部8から出力される重み付け係数データ $w_f (t)_{1...s}$ とを入力する。そして、内分処理部9は、重み付け係数データ $w_{att} (t)_{1...s}$ と、重み付け係数データ $w_f (t)_{1...s}$ とに対して、内分処理を実行することで、合成重み付け係数データ $w (t)$ を取得する。具体的には、内分処理部9は、

$$w (t)_{1...s} = (1 - \alpha) \times w_{att} (t)_{1...s} + \alpha \times w_f (t)_{1...s}$$

s

$$0 \leq \alpha \leq 1$$

に相当する処理を実行することで、合成重み付け係数データ $w(t)_{1..s}$ を取得する。そして、内分処理部 9 は、取得した合成重み付け係数データ $w(t)_{1..s}$ をコンテキスト算出部 10 に出力する。

[0319] コンテキスト算出部 10 は、アテンション部 4 A から出力されるデータ $D \times 4 (= h_i)$ の集合データ $h_i_{1..s}$ と、内分処理部 9 から出力される合成重み付け係数データ $w(t)_{1..s}$ とを入力する。そして、コンテキスト算出部 10 は、合成重み付け係数データ $w(t)_{1..s}$ に基づいて、データ $D \times 4 (= h_i)$ の集合データ $h_i_{1..s}$ に対して、重み付け加算処理を実行することで、コンテキスト状態データ $c(t)$ を取得する。そして、コンテキスト算出部 10 は、取得したコンテキスト状態データ $c(t)$ をデコーダ部 5 B のデコーダ側 LSTM 層 5 2 B に出力する。

[0320] デコーダ側 LSTM 層 5 2 B は、入力されたデータ $D y_2(t)$ 、データ $D y_3(t-1)$ 、および、コンテキスト状態データ $c(t)$ を用いて、LSTM 層による処理を実行し、処理後のデータをデータ $D y_3(t)$ として線形予測部 5 3 に出力する。

[0321] 線形予測部 5 3、ポストネット処理部 5 4、および、加算器 5 5 では、第 1 実施形態と同様の処理が実行される。

[0322] ボコーダ 6 は、デコーダ部 5 B の加算器 5 5 から出力されるデータ $D y_6$ (予測メルスペクトログラムのデータ (音響特徴量のデータ)) を入力とし、入力されたデータ $D y_6$ に対して、学習済みモデルを用いたニューラルネットワーク処理による音声合成処理を実行し、データ $D y_6$ (予測メルスペクトログラム) に対応する音声信号波形データを取得する。そして、ボコーダ 6 は、取得した音声信号波形データを、データ $D o u t$ として出力する。

[0323] このように、音声合成処理装置 300 では、入力されたテキストデータ $D i n$ に対応する音声波形データ $D o u t$ を取得することができる。

[0324] 音声合成処理装置 300 では、図 10~図 12 を用いて説明したのと同様に、予測処理時においても、アテンション部 4 A により取得された重み付け係数データ $w_{att}(t)$ と、強制アテンション部 8 により取得された重み付け

係数データ $w_f(t)$ とを内分処理により合成した重み付け係数データを用いて、コンテキスト状態データ $c(t)$ を生成する。そして、音声合成処理装置 300 では、上記のようにして生成されたコンテキスト状態データ $c(t)$ を用いて、デコーダ部 5 B、ボコーダ 6 による処理が実行されるため、注意機構予測が失敗することに起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することを適切に防止できる。

[0325] 例えば、図 13 に示すように、時刻 t_2 における処理で、注意機構の予測が失敗している場合、すなわち、図 13 に示すように、アテンション部 4 により取得された重み付け係数データが「0」（あるいは所定の値以下）である場合（ $w_{att}(t)_{1...s}$ のすべての要素データの値が「0」（あるいは所定の値以下）である場合）であっても、音声合成処理装置 300 では、強制アテンション部 8 により取得された重み付け係数データ $w_f(t)$ の重みにより、注意機構の予測の失敗が音声合成処理に影響を及ぼさないようにできる合成重み付け係数データ $w(t)_{1...s}$ を取得することができる（図 13 の場合。合成重み付け係数データ $w(t)_{1...s}$ の各要素データの値は、すべて「0.5」）。

[0326] このように、音声合成処理装置 300 では、音素継続長については、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理（音素継続長推定部 7 による処理）により取得した音素継続長を用いて処理することで、音素継続長の予測精度を保証する。つまり、音声合成処理装置 300 では、安定して音素継続長を適切に推定することができる、隠れマルコフモデル等のモデルを用いた推定処理（音素継続長推定部 7 による処理）により取得した音素継続長を用いて強制アテンション部 8 により取得した重み付け係数データと、アテンション部 4 A により取得された重み付け係数データとを適度に合成した重み付け係数データにより生成したコンテキスト状態データ $c(t)$ を用いて予測処理を実行する。したがって、音声合成処理装置 300 では、注意機構の予測が失敗する場合（ア

テンション部4により適切な重み付け係数データが取得できない場合)であっても、強制アテンション部8により取得した重み付け係数データによる重み分の重み付け係数データが取得できるため、注意機構の予測の失敗が音声合成処理に影響を及ぼさないようにできる。

[0327] さらに、音声合成処理装置300では、音響特徴量については、sequence-to-sequence方式を用いたニューラルネットワークのモデルで処理することにより取得できるので、高精度な音響特徴量の予測処理が実現できる。

[0328] したがって、音声合成処理装置300では、注意機構予測が失敗することに起因する、合成発話が途中で止まってしまう、同じフレーズを何回も繰り返してしまう、等の問題が発生することを適切に防止するとともに、高精度な音声合成処理を実行することができる。

[0329] なお、上記では、内分比 α を固定値(例えば、0.5)に設定した場合について、説明したが、これに限定されることはなく、内分比 α は動的に更新されるものであってもよい。例えば、内分処理部9において、アテンション部4Aから入力される重み付け係数データ $w_{att}(t)_{1..s}$ が所定の期間、継続して、所定の値よりも小さい、あるいは、略0であり、かつ、強制アテンション部8から入力される重み付け係数データ $w_f(t)_{1..s}$ が「1」である場合、アテンション部4による処理が失敗している(注意機構予測が失敗している)と判定し、 α の値をより大きな値(重み付け係数データ $w_f(t)_{1..s}$ の重みが大きくなる値)に調整(更新)するようにしてもよい。

[0330] また、音声合成処理装置300において、エンコーダ部3、デコーダ部5は、上記の構成に限定されるものではなく、他の構成のものであってもよい。例えば、下記文献Aに開示されているトランスフォーマーモデルのアーキテクチャによるエンコーダ、デコーダの構成を採用して、エンコーダ部3、デコーダ部5を構成するようにしてもよい。この場合、トランスフォーマーモデルのアーキテクチャによるエンコーダとデコーダの間に設置されるアテンション機構を、本実施形態で説明した機構、すなわち、アテンション部4、強制アテンション部8、内分処理部9、コンテキスト算出部10により、ア

テンション機構が取得した重み付け係数データと、強制アテンション部8が取得した重み付け係数データとを内分処理により合成し、合成した重み付け係数データによりコンテキスト状態データを取得する機構に置換する構成を採用すればよい。

(文献A) : A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, AN. Gomez, L. Kaiser, I. Polosukhin, "Attention is all you need" 31st Conference on Neural Information Processing System (NIPS 2017), Long Beach, CA, USA.

[他の実施形態]

上記実施形態(変形例を含む)の音声合成処理装置において、エンコーダ側LSTM層32、デコーダ側LSTM層52は、それぞれ、複数のLSTM層を備えるものであってもよい。また、エンコーダ側LSTM層32、デコーダ側LSTM層52は、それぞれ、双方向LSTM層(順伝搬、逆伝搬をLSTM層)で構成されるものであってもよい。

[0331] また、上記実施形態(変形例を含む)では、音声合成処理装置が、テキスト解析部1と、フルコンテキストラベルベクトル処理部2とを備え、テキスト解析部1で取得したフルコンテキストラベルデータから、フルコンテキストラベルベクトル処理部2により、最適化フルコンテキストラベルデータを取得する場合について説明したが、これに限定されることはなく、例えば、音声合成処理装置において、最適化フルコンテキストラベルデータを取得する、テキスト解析部を設け、フルコンテキストラベルベクトル処理部を省略する構成としてもよい。

[0332] また、上記実施形態(変形例を含む)を適宜組み合わせてもよい。

[0333] また上記実施形態(変形例を含む)で説明した音声合成処理装置において、各ブロックは、LSIなどの半導体装置により個別に1チップ化されても良いし、一部または全部を含むように1チップ化されても良い。

[0334] なおここではLSIとしたが、集積度の違いにより、IC、システムLSI、スーパーLSI、ウルトラLSIと呼称されることもある。

- [0335] また集積回路化の手法はLSIに限るものではなく、専用回路または汎用プロセサで実現してもよい。LSI製造後にプログラムすることが可能なFPGA (Field Programmable Gate Array) や、LSI内部の回路セルの接続や設定を再構成可能なリコンフィギュラブル・プロセッサを利用してよい。
- [0336] また上記各実施形態の各機能ブロックの処理の一部または全部は、プログラムにより実現されるものであってもよい。そして上記各実施形態の各機能ブロックの処理の一部または全部は、コンピュータにおいて、中央演算装置 (CPU) により行われる。また、それぞれの処理を行うためのプログラムは、ハードディスク、ROMなどの記憶装置に格納されており、ROMにおいて、あるいはRAMに読み出されて実行される。
- [0337] また上記実施形態の各処理をハードウェアにより実現してもよいし、ソフトウェア (OS (オペレーティングシステム)、ミドルウェア、あるいは所定のライブラリとともに実現される場合を含む。) により実現してもよい。さらにソフトウェアおよびハードウェアの混在処理により実現してもよい。
- [0338] 例えば上記実施形態の各機能部をソフトウェアにより実現する場合、図14に示したハードウェア構成 (例えばCPU、GPU、ROM、RAM、入力部、出力部、通信部、記憶部 (例えば、HDD、SSD等により実現される記憶部)、外部メディア用ドライブ等をバスBusにより接続したハードウェア構成) を用いて各機能部をソフトウェア処理により実現するようにしてもよい。
- [0339] また上記実施形態の各機能部をソフトウェアにより実現する場合、当該ソフトウェアは、図14に示したハードウェア構成を有する単独のコンピュータを用いて実現されるものであってもよいし、複数のコンピュータを用いて分散処理により実現されるものであってもよい。
- [0340] また上記実施形態における処理方法の実行順序は、必ずしも上記実施形態の記載に制限されるものではなく、発明の要旨を逸脱しない範囲で、実行順序を入れ替えることができるものである。

[0341] 前述した方法をコンピュータに実行させるコンピュータプログラム、及びそのプログラムを記録したコンピュータ読み取り可能な記録媒体は、本発明の範囲に含まれる。ここでコンピュータ読み取り可能な記録媒体としては、例えば、フレキシブルディスク、ハードディスク、CD-ROM、MO、DVD、DVD-ROM、DVD-RAM、大容量DVD、次世代DVD、半導体メモリを挙げることができる。

[0342] 上記コンピュータプログラムは、上記記録媒体に記録されたものに限らず、電気通信回線、無線または有線通信回線、インターネットを代表とするネットワーク等を経由して伝送されるものであってもよい。

[0343] なお本発明の具体的な構成は、前述の実施形態に限られるものではなく、発明の要旨を逸脱しない範囲で種々の変更および修正が可能である。

符号の説明

[0344] 100、200、300 音声合成処理装置

- 1 テキスト解析部
- 2、2A フルコンテキストラベルベクトル処理部
- 3 エンコーダ部
- 4、4A アテンション部
- 5 デコーダ部
- 6 ボコーダ
- 7 音素継続長推定部
- 8 強制アテンション部
- 9 内分処理部
- 10 コンテキスト算出部

請求の範囲

[請求項1] 任意の言語を処理対象言語とし、エンコーダ・デコーダ方式のニューラルネットワークを用いて音声合成処理を実行する音声合成処理装置であって、

前記処理対象言語のテキストデータに対してテキスト解析処理を実行し、コンテキストラベルデータを取得するテキスト解析部と、

前記テキスト解析部により取得された前記コンテキストラベルデータから、コンテキストラベルデータを取得する処理において処理対象とされた音素である単独音素についてのコンテキストラベルを取得することで、前記ニューラルネットワークの学習処理に適した最適化フルコンテキストラベルデータを取得するフルコンテキストラベルベクトル処理部と、

前記最適化フルコンテキストラベルデータに基づいて、ニューラルネットワークのエンコード処理を実行することで、隠れ状態データを取得するエンコーダ部と、

前記隠れ状態データに基づいて、ニューラルネットワークのデコード処理を実行することで、前記最適化フルコンテキストラベルデータに対応する音響特徴量データを取得するデコーダ部と、

前記デコーダ部により取得された音響特徴量から音声波形データを取得するボコーダと、
を備える音声合成処理装置。

[請求項2] 前記音響特徴量は、メルスペクトログラムのデータである、
請求項1に記載の音声合成処理装置。

[請求項3] 前記ボコーダは、
ニューラルネットワークのモデルを用いた処理を実行することで、
音響特徴量から音声波形データを取得する、
請求項1または2に記載の音声合成処理装置。

[請求項4] 前記ボコーダは、

可逆変換ネットワークにより構成されたニューラルネットワークのモデルを用いた処理を実行することで、音響特徴量から音声波形データを取得する、

請求項3に記載の音声合成処理装置。

[請求項5]

音素単位のコンテキストラベルデータから音素継続長を推定する音素継続長推定部をさらに備え、

前記フルコンテキストラベルベクトル処理部は、前記音素継続長推定部により推定された音素継続長である推定音素継続長に対応する期間において、当該推定音素継続長に対応する音素の前記最適化フルコンテキストラベルデータを継続して前記エンコーダ部へ出力する、

請求項1から4のいずれかに記載の音声合成処理装置。

[請求項6]

任意の言語を処理対象言語とし、エンコーダ・デコーダ方式のニューラルネットワークを用いて音声合成処理を実行する音声合成処理方法であって、

前記処理対象言語のテキストデータに対してテキスト解析処理を実行し、コンテキストラベルデータを取得するテキスト解析ステップと、

前記テキスト解析ステップにより取得された前記コンテキストラベルデータから、コンテキストラベルデータを取得する処理において処理対象とされた音素である単独音素についてのコンテキストラベルを取得することで、前記ニューラルネットワークの学習処理に適した最適化フルコンテキストラベルデータを取得するフルコンテキストラベルベクトル処理ステップと、

前記最適化フルコンテキストラベルデータに基づいて、ニューラルネットワークのエンコード処理を実行することで、隠れ状態データを取得するエンコード処理ステップと、

前記隠れ状態データに基づいて、ニューラルネットワークのデコード処理を実行することで、前記最適化フルコンテキストラベルデータ

に対応する音響特徴量データを取得するデコード処理ステップと、

前記デコード処理ステップにより取得された音響特徴量から音声波形データを取得するボコーダ処理ステップと、

を備える音声合成処理方法。

[請求項7] 請求項6に記載の音声合成処理方法をコンピュータに実行させるためのプログラム。

[請求項8] 任意の言語を処理対象言語とし、エンコーダ・デコーダ方式のニューラルネットワークを用いて音声合成処理を実行する音声合成処理装置であって、

前記処理対象言語のテキストデータに対してテキスト解析処理を実行し、コンテキストラベルデータを取得するテキスト解析部と、

前記テキスト解析部により取得された前記コンテキストラベルデータから、コンテキストラベルデータを取得する処理において処理対象とされた音素である単独音素についてのコンテキストラベルを取得することで、前記ニューラルネットワークの学習処理に適した最適化フルコンテキストラベルデータを取得するフルコンテキストラベルベクトル処理部と、

前記最適化フルコンテキストラベルデータに基づいて、ニューラルネットワークのエンコード処理を実行することで、隠れ状態データを取得するエンコーダ部と、

音素単位のコンテキストラベルデータから音素継続長を推定する音素継続長推定部と、

前記音素継続長推定部により推定された音素継続長に基づいて、第1重み付け係数データを取得する強制アテンション部と、

前記エンコーダ部により取得された隠れ状態データに基づいて、第2重み付け係数データを取得するアテンション部と、

前記第1重み付け係数データと第2重み付け係数データとに対して内分処理を行うことで、合成重み付け係数データを取得する内分処理

部と、

前記合成重み付け係数データにより、前記エンコーダ部により取得された前記隠れ状態データに対して重み付け合成処理を実行することで、コンテキスト状態データを取得するコンテキスト算出部と、

前記コンテキスト状態データに基づいて、ニューラルネットワークのデコード処理を実行することで、前記最適化フルコンテキストラベルデータに対応する音響特徴量データを取得するデコーダ部と、

前記デコーダ部により取得された音響特徴量から音声波形データを取得するボコーダと、

を備える音声合成処理装置。

【図1】

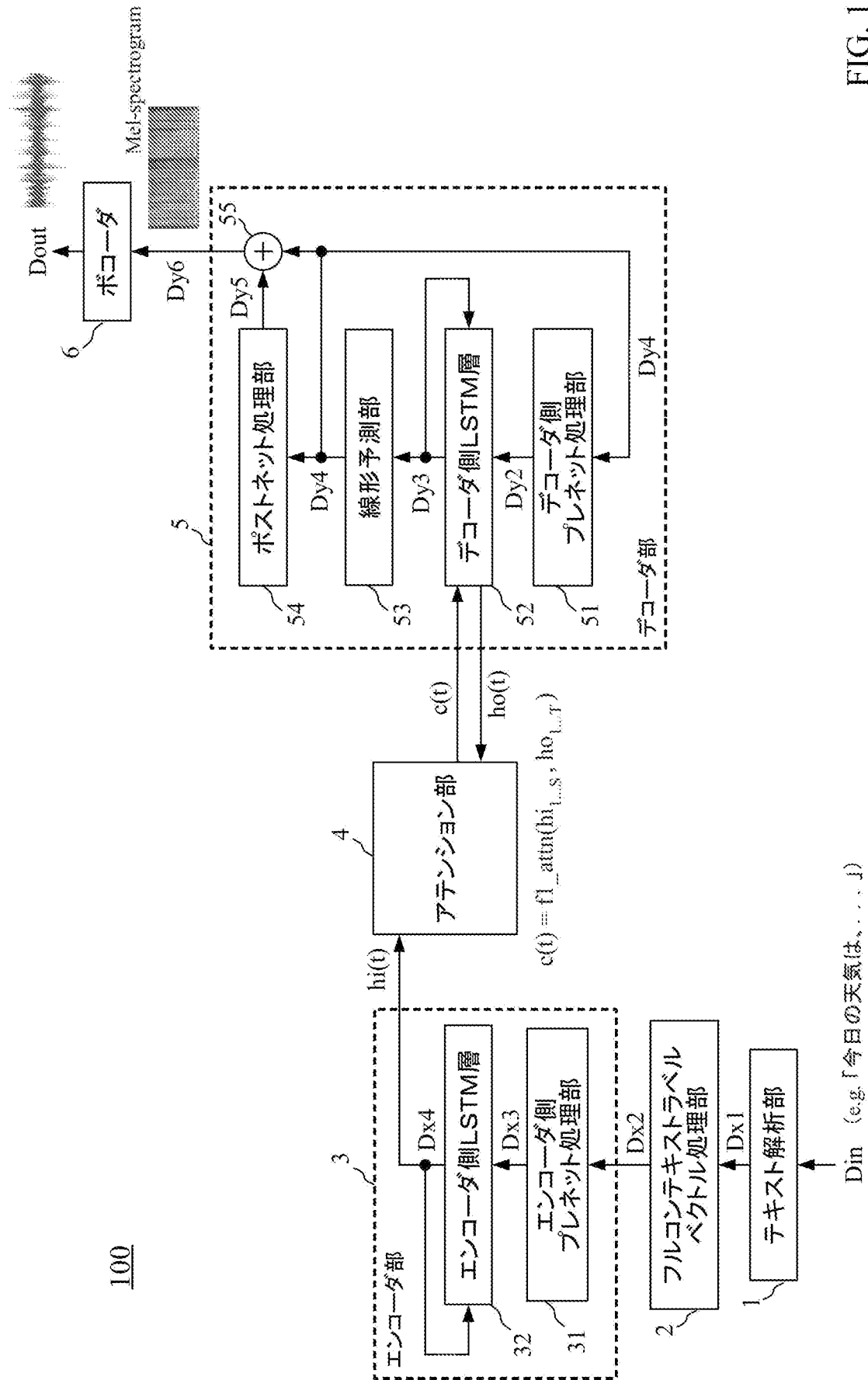


FIG. 1

100

Dim (e.g. 「今日の天気は、...」)

[図2]

フルコンテキストラベルに含まれる情報(一例)

次元	音素数	総計	概要
43	5	215	[2つ先行/1つ先行/当該/1つ後続/2つ後続]の音素インデックス
1	1	1	前から数えた当該アクセント句内のアクセント位置
2	1	2	[前/後ろ]から数えた当該アクセント句内のモーラ位置
14	1	14	[1つ先行]のPOS(品詞)
15	1	15	[1つ先行]の名詞のカテゴリ
7	1	7	[1つ先行]のPOS(品詞)活用例
14	1	14	[当該]のPOS
15	1	15	[当該]の名詞のカテゴリ
7	1	7	[当該]のPOS活用例
14	1	14	[1つ後続]のPOS
15	1	15	[1つ後続]の名詞のカテゴリ
7	1	7	[1つ後続]のPOS活用例
1	1	1	[1つ先行]のアクセント句のモーラ数
1	1	1	[1つ先行]のアクセント句のアクセント位置
1	1	1	[1つ先行]のアクセント句の語尾上げ判定
1	1	1	[1つ先行]のアクセント句のポーズ有無
1	1	1	[当該]のアクセント句のモーラ数
1	1	1	[当該]のアクセント句のアクセント位置
1	1	1	[当該]のアクセント句の語尾上げ判定
2	1	2	[前/後ろ]から数えた当該呼気段落内のアクセント句位置 (アクセント句単位)
2	1	2	[前/後ろ]から数えた当該呼気段落内のアクセント句位置 (モーラ単位)
1	1	1	[1つ後続]のアクセント句のモーラ数
1	1	1	[1つ後続]のアクセント句のアクセント位置
1	1	1	[1つ後続]のアクセント句の語尾上げ判定
1	1	1	[1つ後続]のアクセント句のポーズ有無
2	1	2	先行呼気段落の[アクセント句数/モーラ数]
1	1	1	文頭判定
2	1	2	当該呼気段落の[アクセント句数/モーラ数]
2	1	2	[前/後ろ]から数えた当該呼気段落の文中位置(呼気段落単位)
2	1	2	[前/後ろ]から数えた当該呼気段落の文中位置(アクセント句単位)
2	1	2	[前/後ろ]から数えた当該呼気段落の文中位置(モーラ単位)
2	1	2	後続呼気段落の[アクセント句数/モーラ数]
1	1	1	文末判定
1	1	1	文章内の呼気段落数
1	1	1	文章内のアクセント句数
1	1	1	文章内のモーラ数
22	5	110	[2つ先行/1つ先行/当該/1つ後続/2つ後続]の音素カテゴリ
1	5	5	[2つ先行/1つ先行/当該/1つ後続/2つ後続]音素の調音位置 によるカテゴライズ
1	5	5	[2つ先行/1つ先行/当該/1つ後続/2つ後続]音素の調音方法 によるカテゴライズ

合計: 478 次元

FIG. 2

[図3]

最適化フルコンテキストラベルに含まれる情報(一例)

次元	音素数	総計	概要
43	1	43	[当該]の音素インデックス
1	1	1	前から数えた当該アクセント句内のアクセント位置
2	1	2	[前/後ろ]から数えた当該アクセント句内のモーラ位置
14	1	14	[当該]のPOS
15	1	15	[当該]の名詞のカテゴリ
7	1	7	[当該]のPOS活用型
1	1	1	[当該]のアクセント句のモーラ数
1	1	1	[当該]のアクセント句のアクセント位置
1	1	1	[当該]のアクセント句の語尾上げ判定
2	1	2	[前/後ろ]から数えた当該呼気段落内のアクセント句位置 (アクセント句単位)
2	1	2	[前/後ろ]から数えた当該呼気段落内のアクセント句位置 (モーラ単位)
2	1	2	先行呼気段落の[アクセント句数/モーラ数]
1	1	1	文頭判定
2	1	2	当該呼気段落の[アクセント句数/モーラ数]
2	1	2	[前/後ろ]から数えた当該呼気段落の文中位置(呼気段落単位)
2	1	2	[前/後ろ]から数えた当該呼気段落の文中位置(アクセント句単位)
2	1	2	[前/後ろ]から数えた当該呼気段落の文中位置(モーラ単位)
2	1	2	後続呼気段落の[アクセント句数/モーラ数]
1	1	1	文末判定
1	1	1	文章内の呼気段落数
1	1	1	文章内のアクセント句数
1	1	1	文章内のモーラ数
22	1	22	[当該]の音素カテゴリ
1	1	1	[当該]音素の調音位置によるカテゴリ
1	1	1	[当該]音素の調音方法によるカテゴリ

合計: 130 次元

FIG. 3

[図4]

学習処理時

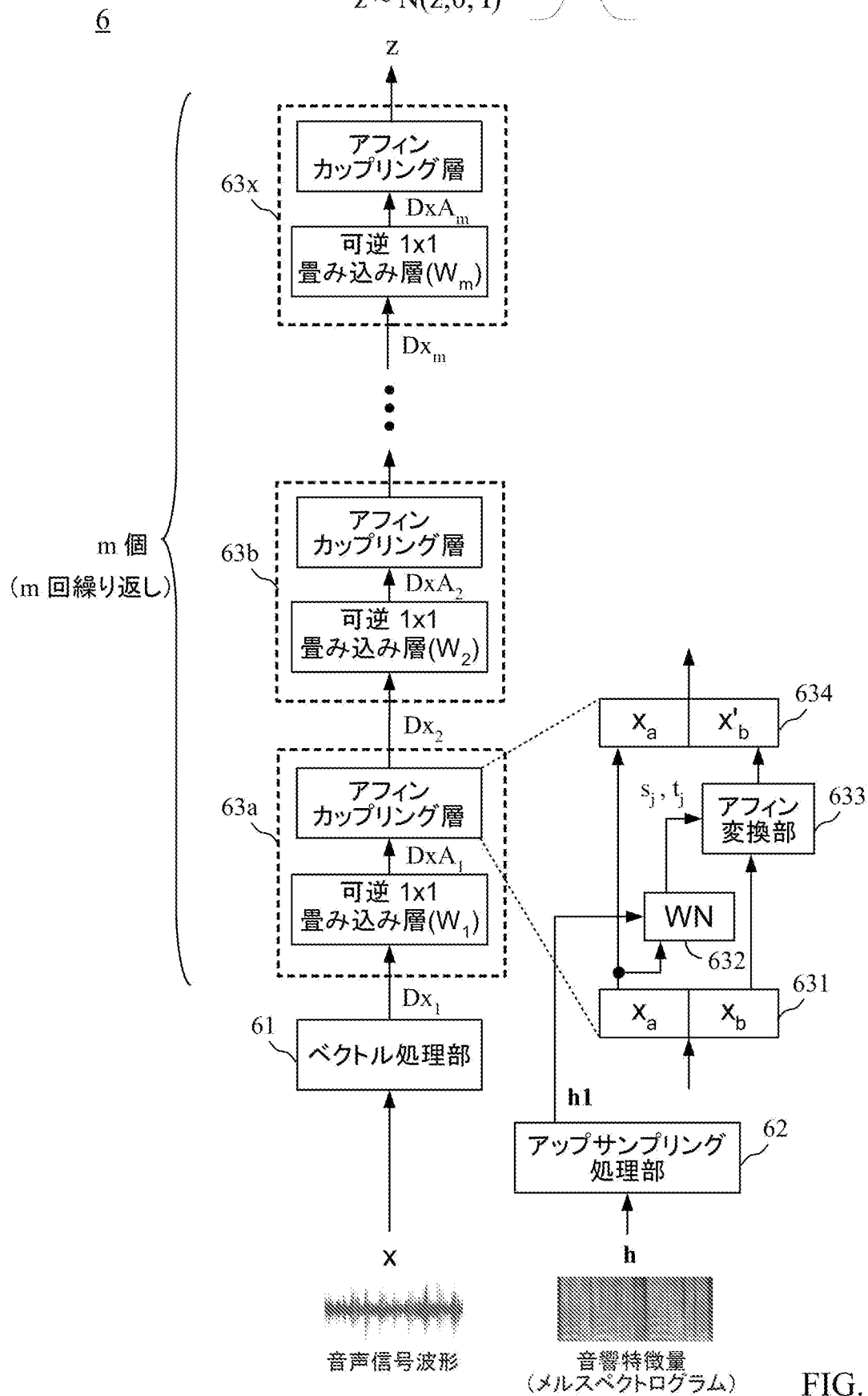


FIG. 4

[図5]

予測処理時

6
m 個
(m 回繰り返し)

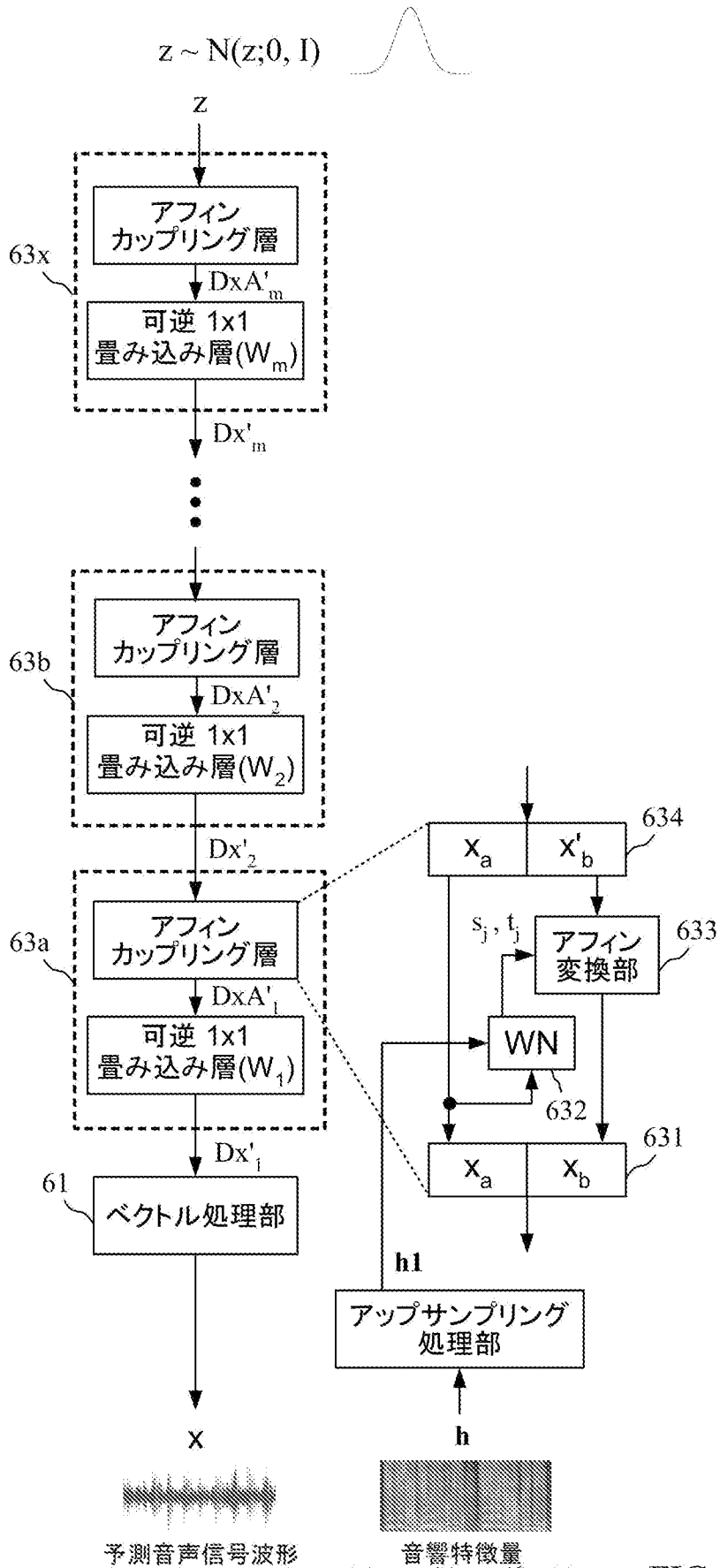


FIG. 5

[図6]

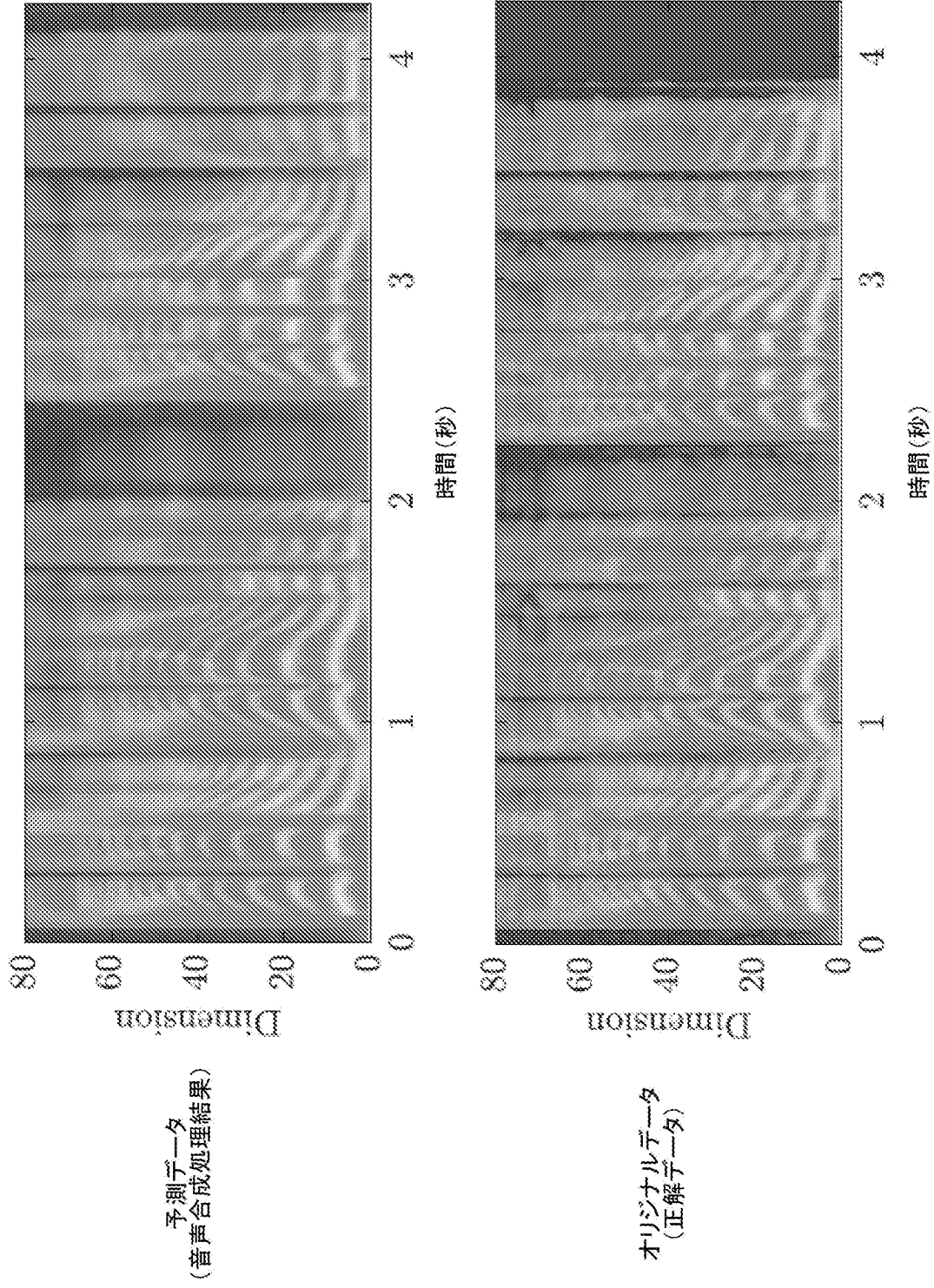


FIG. 6

[図7]

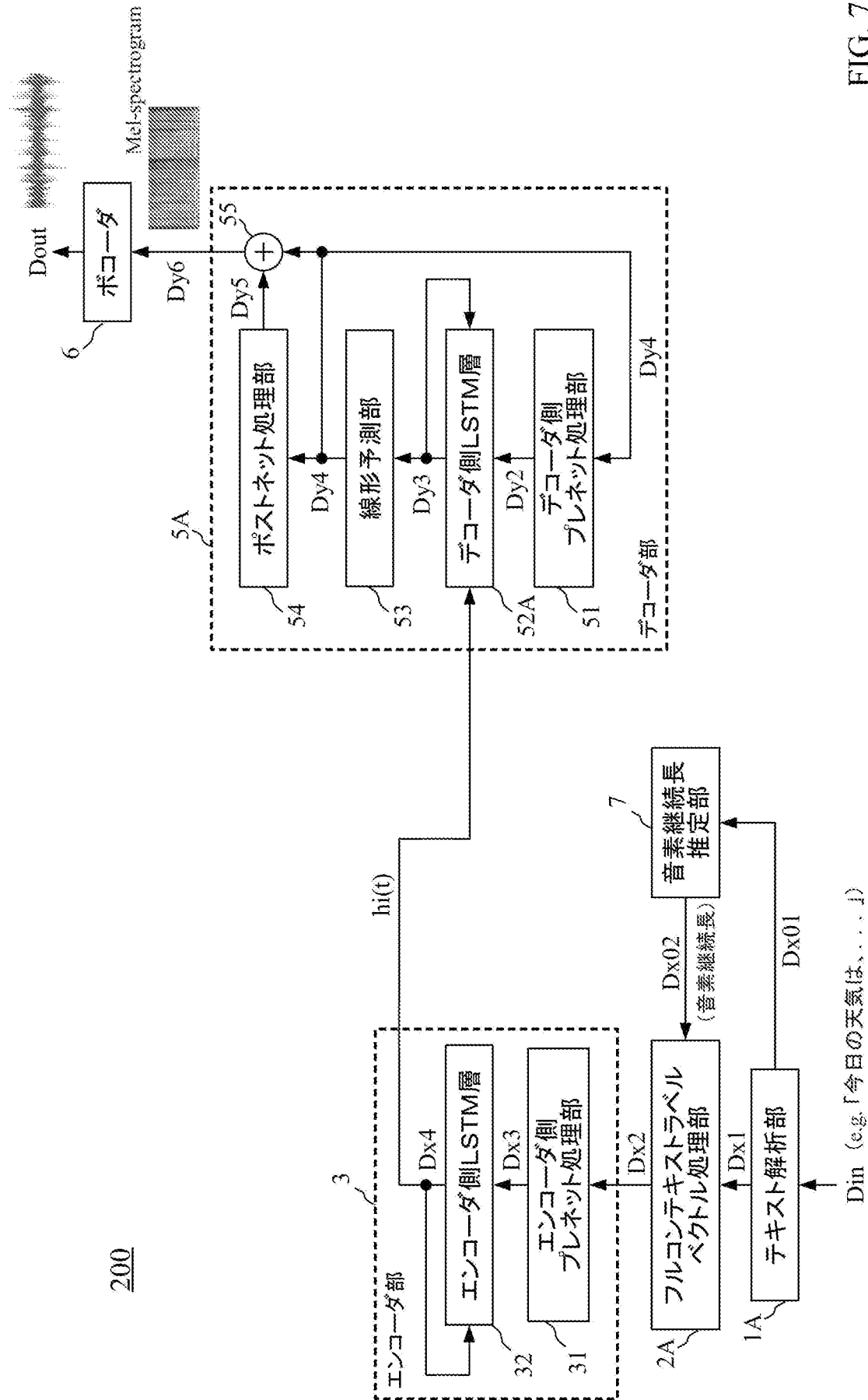


FIG. 7

[図8]

Din : 「今日の天気は、...」

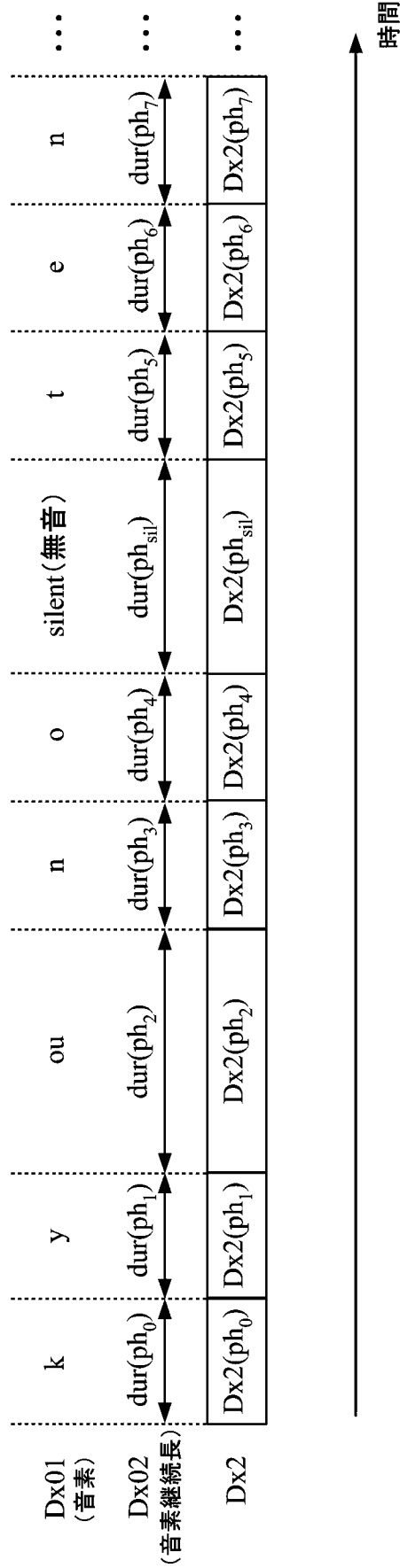


FIG. 8

図11

$\alpha = 0.5$ の場合

Din : 「今日の天気は、...」

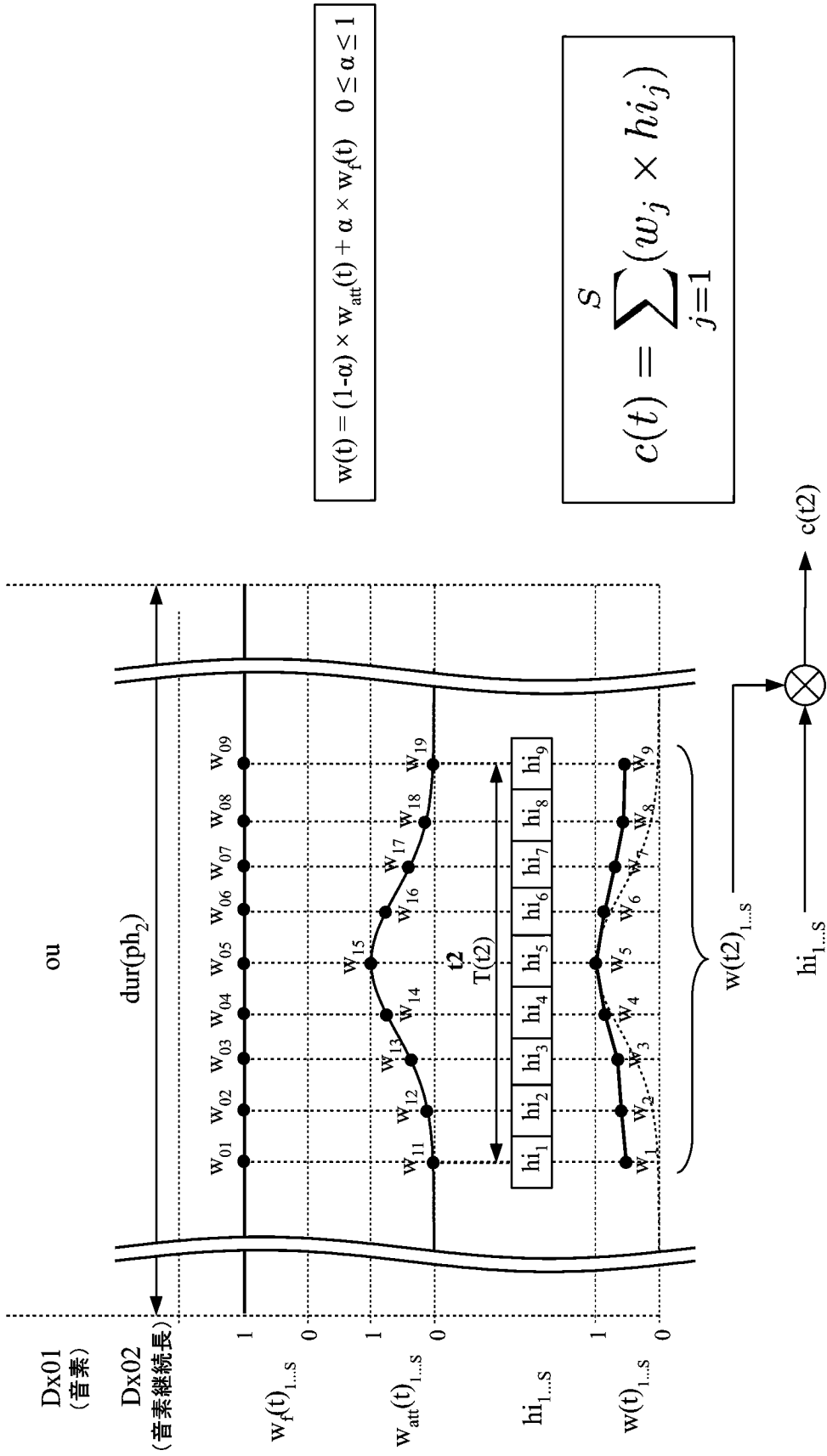


FIG. 11

図12

$\alpha = 0.5$ の場合

Din : 「今日の天気は、...」

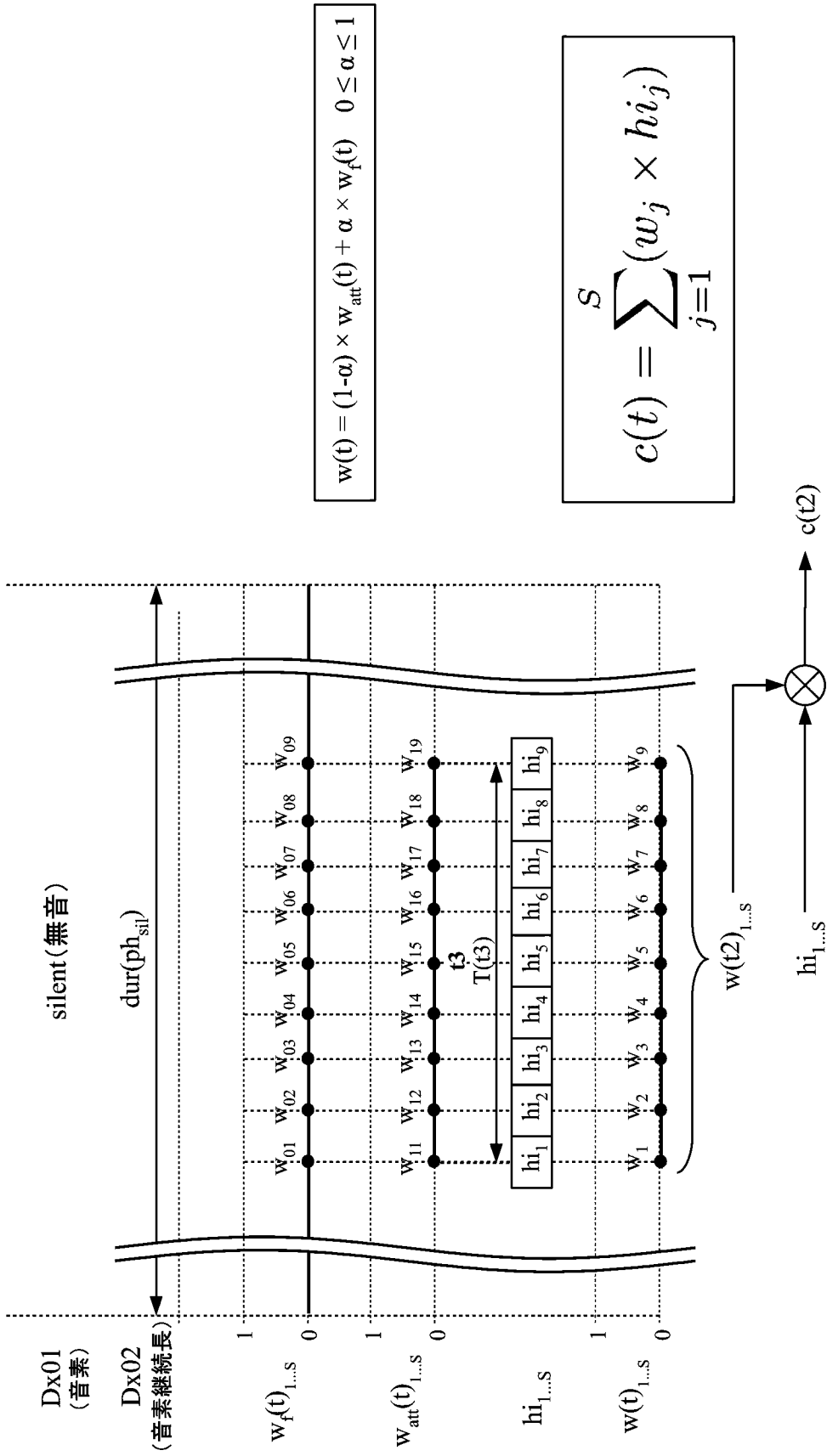


FIG. 12

[図13]

$\alpha = 0.5$ の場合

Din : 「今日の天気は、...」

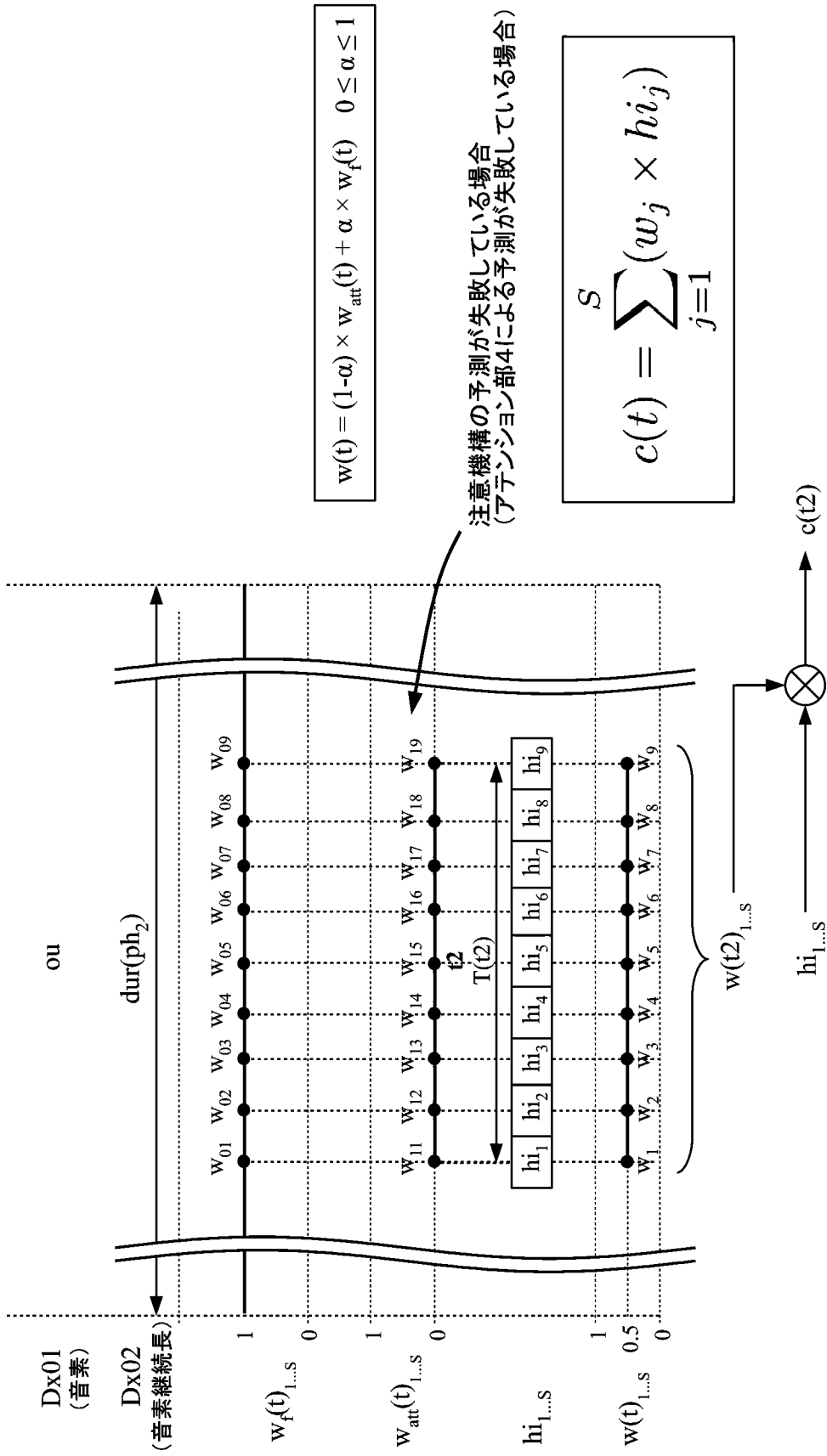


FIG. 13

[図14]

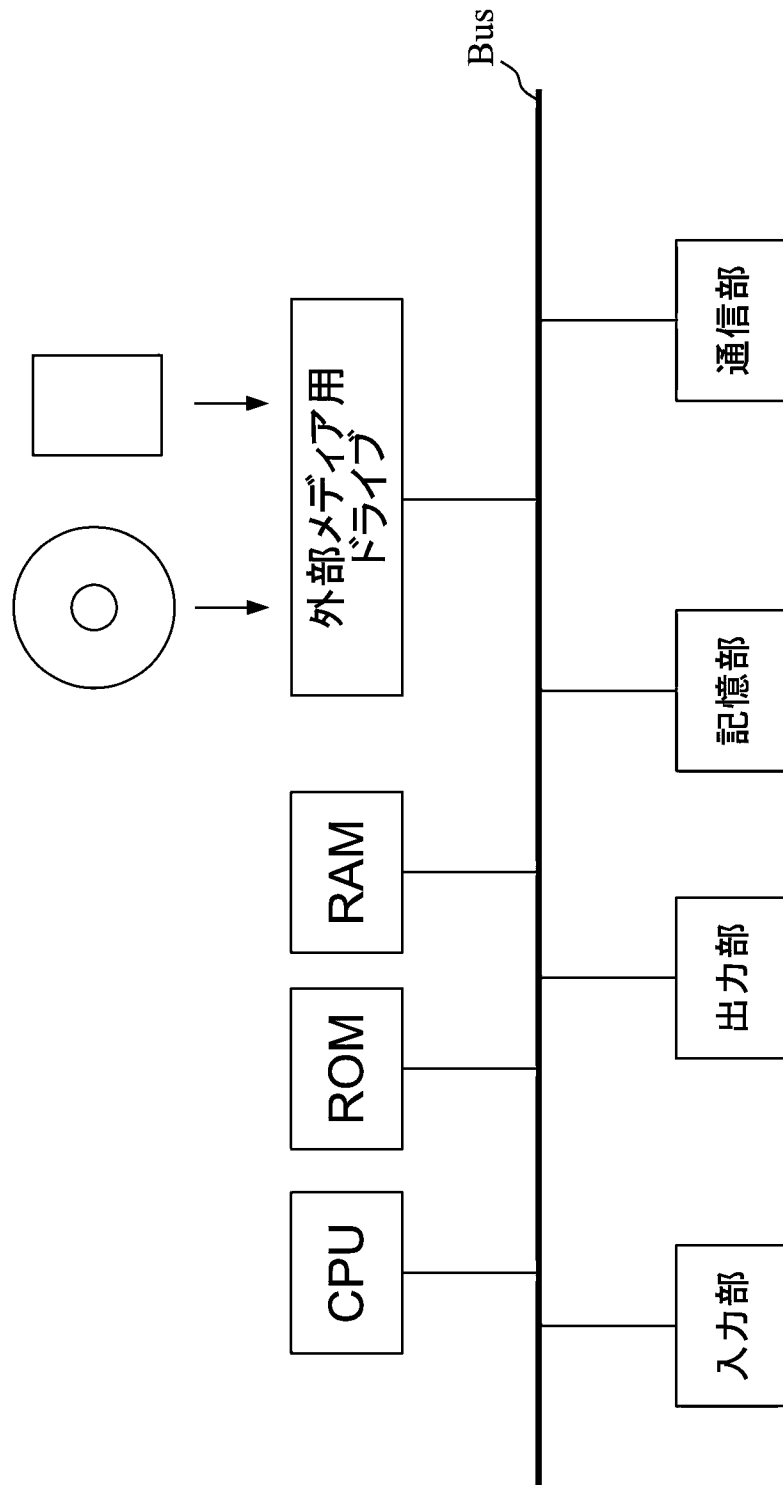


FIG. 14

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2020/025682

A. CLASSIFICATION OF SUBJECT MATTER

G10L 13/08(2013.01)i; G10L 13/10(2013.01)i; G10L 25/30(2013.01)i
 FI: G10L13/08 110Z; G10L25/30; G10L13/08 150Z; G10L13/10 111F

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G10L13/08; G10L13/10; G10L25/30

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Published examined utility model applications of Japan	1922-1996
Published unexamined utility model applications of Japan	1971-2020
Registered utility model specifications of Japan	1996-2020
Published registered utility model applications of Japan	1994-2020

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	藤本崇人、外 4 名, 日本語 End-to-End 音声合成における入力言語特徴量の影響, 日本音響学会 2019 年春季研究発表会講演論文集 CD-ROM, 19 February 2019, pp. 1061-1062, pp. 1061-1062, non-official translation (FUJIMOTO, Takahito et al., "Effect of input language feature values on Japanese End-to-End speech synthesis", Lecture proceedings of 2019 spring research conference of the Acoustical Society of Japan CD-ROM)	1-8
A	安田裕介、外 3 名, 日本語エンドツーエンド音声合成へむけて - 日本語 Tacotron の初期的検討, 日本音響学会 2018 年秋季研究発表会講演論文集, 29 August 2018, pp. 1167-1168, pp. 1167-1168, non-official translation (YASUDA, Yusuke et al., "Towards end-to-end Japanese speech synthesis -An initial consideration of Japanese Tacotron", Lecture Proceedings of 2018 Autumn Research Conference of the Acoustical Society of Japan)	1-8

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search
26 August 2020 (26.08.2020)

Date of mailing of the international search report
08 September 2020 (08.09.2020)

Name and mailing address of the ISA/
Japan Patent Office
3-4-3, Kasumigaseki, Chiyoda-ku,
Tokyo 100-8915, Japan

Authorized officer

Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2020/025682

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>栗原清、外 3 名, 読み仮名と韻律記号を入力とする日本語 end-to-end 音声合成の音質評価, 電子情報通信学会技術研究報告, 03 December 2018, vol. 118, no. 354, pp. 89-94, pp. 89-94, (KURIHARA, Kiyoshi et al., "Evaluation of Japanese end-to-end speech synthesis method inputting kana and prosodic symbols", IEICE technical report)</p>	1-8
A	<p>SHEN, Jonathan et al., "NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS", Proc. ICASSP 2018, April 2018, pp. 4779-4783, pp. 4779-4783</p>	1-8
A	<p>WO 2018/183650 A2 (GOOGLE LLC) 04.10.2018 (2018-10-04) entire text, all drawings</p>	1-8
A	<p>PRENGER, Ryan et al., "WAVEGLOW: A FLOW-BASED GENERATIVE NETWORK FOR SPEECH SYNTHESIS", Proc. ICASSP 2019, May 2019, pp. 3617-3621, pp. 3617-3621</p>	4

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/JP2020/025682

Patent Documents referred in the Report	Publication Date	Patent Family	Publication Date
WO 2018/183650 A2	04 Oct. 2018	US 2019/0311708 A1 EP 3533594 A2 AU 2018244917 A CA 3058433 A1 CN 110476206 A KR 10-2019-0130634 A	

A. 発明の属する分野の分類（国際特許分類（IPC）） G10L 13/08(2013.01)i; G10L 13/10(2013.01)i; G10L 25/30(2013.01)i FI: G10L13/08 110Z; G10L25/30; G10L13/08 150Z; G10L13/10 111F		
B. 調査を行った分野 調査を行った最小限資料（国際特許分類（IPC）） G10L13/08; G10L13/10; G10L25/30 最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2020年 日本国実用新案登録公報 1996-2020年 日本国登録実用新案公報 1994-2020年		
国際調査で使用した電子データベース（データベースの名称、調査に使用した用語）		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	藤本 崇人、外4名、日本語End-to-End 音声合成における入力言語特徴量の影響、 日本音響学会 2019年 春季研究発表会講演論文集CD-ROM, 2019.02.19, pp. 1061-1062 pp. 1061-1062	1-8
A	安田 裕介、外3名、日本語エンドツーエンド音声合成へむけてー日本語Tacotronの 初期的検討、日本音響学会 2018年 秋季研究発表会講演論文集, 2018.08.29, pp. 1167-1168 pp. 1167-1168	1-8
A	栗原 清、外3名、読み仮名と韻律記号を入力とする日本語end-to-end音 声合成の音質評価、電子情報通信学会技術研究報告, 2018.12.03, Vol. 118, No. 354, pp. 89-94 pp. 89-94	1-8
A	SHEN, Jonathan, et al., NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS, Proc. ICASSP 2018, 2018.04, pp. 4779-4783 pp. 4779-4783	1-8
A	WO 2018/183650 A2 (GOOGLE LLC) 04.10.2018 (2018-10-04) 全文, 全図	1-8
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input checked="" type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー “A” 特に関連のある文献ではなく、一般的技術水準を示すもの “E” 国際出願日前の出願または特許であるが、国際出願日以後に 公表されたもの “L” 優先権主張に疑義を提起する文献又は他の文献の発行日若し くは他の特別な理由を確立するために引用する文献（理由を 付す） “O” 口頭による開示、使用、展示等に言及する文献 “P” 国際出願日前で、かつ優先権の主張の基礎となる出願の日の 後に公表された文献	“T” 国際出願日又は優先日後に公表された文献であって出願と抵 触するものではなく、発明の原理又は理論の理解のために引 用するもの “X” 特に関連のある文献であって、当該文献のみで発明の新規性 又は進歩性がないと考えられるもの “Y” 特に関連のある文献であって、当該文献と他の1以上の文献 との、当業者にとって自明である組合せによって進歩性がな いと考えられるもの “&” 同一パテントファミリー文献	
国際調査を完了した日 26.08.2020	国際調査報告の発送日 08.09.2020	
名称及びあて先 日本国特許庁(ISA/JP) 〒100-8915 日本国 東京都千代田区霞が関三丁目4番3号	権限のある職員（特許庁審査官） 上田 雄 5Z 5095 電話番号 03-3581-1101 内線 3591	

C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	PRENGER, Ryan, et al., WAVEGLOW: A FLOW-BASED GENERATIVE NETWORK FOR SPEECH SYNTHESIS, Proc. ICASSP 2019, 2019.05, pp. 3617-3621 pp. 3617-3621	4

国際調査報告
パテントファミリーに関する情報

国際出願番号

PCT/JP2020/025682

引用文献			公表日	パテントファミリー文献			公表日
WO	2018/183650	A2	04.10.2018	US	2019/0311708	A1	
				EP	3583594	A2	
				AU	2018244917	A	
				CA	3058433	A1	
				CN	110476206	A	
				KR	10-2019-0130634	A	
<hr/>							