



- (51) 国際特許分類:
G06N 20/00 (2019.01)
- (21) 国際出願番号 : PCT/JP2019/044770
- (22) 国際出願日 : 2019年11月14日(14.11.2019)
- (25) 国際出願の言語 : 日本語
- (26) 国際公開の言語 : 日本語
- (71) 出願人: 富士通株式会社 (FUJITSU LIMITED)
[JP/JP]; 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 Kanagawa (JP).
- (72) 発明者: 山田 萌(YAMADA, Moyuru); 〒2118588
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP).
- (74) 代理人: 酒井 昭徳(SAKAI, Akinori); 〒1020094
東京都千代田区紀尾井町3番12号 紀尾井町ビル7階 酒井総合特許事務所 Tokyo (JP).
- (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,

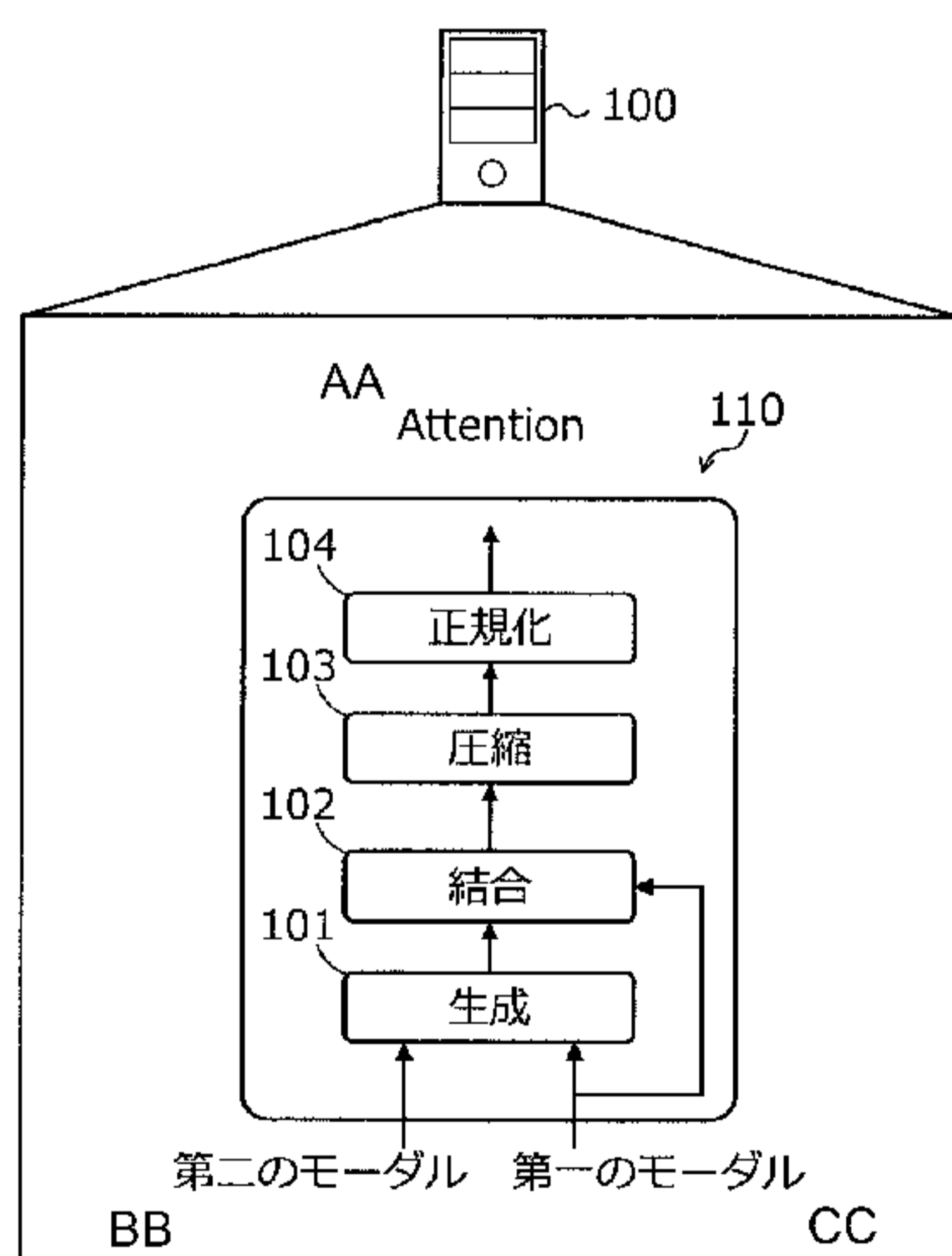
BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, 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).

(54) Title: OUTPUT METHOD, OUTPUT PROGRAM, AND OUTPUT DEVICE

(54) 発明の名称 : 出力方法、出力プログラム、および出力装置

[図1]



- 101 Generation
- 102 Coupling
- 103 Compression
- 104 Normalization
- AA Attention
- BB Second modal
- CC First modal

(57) Abstract: An output device (100) generates, by using a generation model (101), a correction vector for correcting a vector based on first modal information, the correction vector being generated on the basis of a correlation between a vector based on the first modal information and a vector based on second modal information. The output device (100) couples, by using a coupling model (102), the generated correction vector to the vector based on the first modal information. The output device (100) compresses, by using a compression model (103), the vector based on the first modal information after coupling, in accordance with a prescribed rule. The output device (100) carries out, by using a normalization model (104), a normalization process on the vector based on the first modal information after compression. The output device (100) outputs the vector obtained through the normalization process.

WO 2021/095212 A1

添付公開書類：

- 一 国際調査報告（条約第21条(3)）

(57) 要約：出力装置（100）は、生成モデル（101）を用いて、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成する。出力装置（100）は、結合モデル（102）を用いて、生成した補正ベクトルを、第一のモーダルの情報に基づくベクトルに結合する。出力装置（100）は、圧縮モデル（103）を用いて、所定のルールに従って、結合後の第一のモーダルの情報に基づくベクトルを圧縮する。出力装置（100）は、正規化モデル（104）を用いて、圧縮後の第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施する。出力装置（100）は、正規化処理により得たベクトルを出力する。

明 細 書

発明の名称：出力方法、出力プログラム、および出力装置

技術分野

[0001] 本発明は、出力方法、出力プログラム、および出力装置に関する。

背景技術

[0002] 従来、複数のモーダルの情報を用いて問題を解く技術がある。この技術は、例えば、文書翻訳や質疑応答、物体検出、状況判断などの問題を解く際に利用される。ここで、モーダルとは、情報の様式や種類を示す概念であり、具体例としては、画像、文書（テキスト）、音声などを挙げることができる。複数のモーダルを用いた機械学習はマルチモーダル学習と呼ばれる。

[0003] 先行技術としては、例えば、Attentionにより情報を変換するTransformerと呼ばれるものがある。Attentionは、具体的には、第一のモーダルの情報に基づくベクトルから得たクエリと、第二のモーダルの情報に基づくベクトルから得たキーとの相関に基づいて、第二のモーダルの情報に基づくベクトルから得たバリューの重み付け和を算出し、第一のモーダルの情報に基づくベクトルに加算する。

先行技術文献

非特許文献

[0004] 非特許文献1：Vaswani, Ashish, et al. “Attention is all you need.” Advances in neural information processing systems. 2017.

発明の概要

発明が解決しようとする課題

[0005] しかしながら、従来技術では、複数のモーダルの情報を用いて問題を解いた際の解の精度が悪い場合がある。例えば、画像と文書とを基に状況を判断する問題を解くにあたり、Attentionにより、画像に関するモーダ

ルの情報に基づくベクトルに、文書に関するモーダルの情報に基づくベクトルから得たバリューの重み付け和を、単純に加算すると、問題の解決に有用な情報が失われやすい。このため、問題を解いた際の解の精度が悪くなりやすい。

[0006] 1つの側面では、本発明は、複数のモーダルの情報を用いて問題を解いた際の解の精度の向上を図ることを目的とする。

課題を解決するための手段

[0007] 1つの実施態様によれば、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、前記第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成し、生成した前記補正ベクトルを、前記第一のモーダルの情報に基づくベクトルに結合し、所定のルールに従って、結合後の前記第一のモーダルの情報に基づくベクトルを圧縮し、圧縮後の前記第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施し、前記正規化処理により得たベクトルを出力する出力方法、出力プログラム、および出力装置が提案される。

発明の効果

[0008] 一態様によれば、複数のモーダルの情報を用いて問題を解いた際の解の精度の向上を図ることが可能になる。

図面の簡単な説明

[0009] [図1]図1は、実施の形態にかかる出力方法の一実施例を示す説明図である。

[図2]図2は、情報処理システム200の一例を示す説明図である。

[図3]図3は、出力装置100のハードウェア構成例を示すブロック図である。

[図4]図4は、出力装置100の機能的構成例を示すブロック図である。

[図5]図5は、Co-Attention Network 500の具体例を示す説明図である。

[図6]図6は、SA層600の具体例と、TA層610の具体例とを示す説明図である。

[図7]図7は、画像TA層501の具体例を示す説明図である。

[図8]図8は、画像TA層501の別の具体例を示す説明図である。

[図9]図9は、画像TA層501と文書TA層503との比較例を示す説明図である。

[図10]図10は、CAN500を用いた動作の一例を示す説明図である。

[図11]図11は、出力装置100の利用例1を示す説明図（その1）である。

[図12]図12は、出力装置100の利用例1を示す説明図（その2）である。

[図13]図13は、出力装置100の利用例2を示す説明図（その1）である。

[図14]図14は、出力装置100の利用例2を示す説明図（その2）である。

[図15]図15は、学習処理手順の一例を示すフローチャートである。

[図16]図16は、推定処理手順の一例を示すフローチャートである。

[図17]図17は、アテンション処理手順の一例を示すフローチャートである。

発明を実施するための形態

[0010] 以下に、図面を参照して、本発明にかかる出力方法、出力プログラム、および出力装置の実施の形態を詳細に説明する。

[0011] （実施の形態にかかる出力方法の一実施例）

図1は、実施の形態にかかる出力方法の一実施例を示す説明図である。出力装置100は、複数のモーダルの情報を用いて、問題の解決に有用な情報を得やすくすることにより、問題を解いた際の解の精度の向上を図るためのコンピュータである。

[0012] 従来、問題を解くための手法として、例えば、Attentionにより情報を変換するTransformerを利用した、BERT (Bidirectional Encoder Representations f

rom Transformers) と呼ばれるものがある。BERTは、具体的には、TransformerのEncoder部を積み重ねて形成される。BERTについては、例えば、下記非特許文献2を参照することができる。

[0013] 非特許文献2 : Devlin, Jacob et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.” NAACL-HLT (2019).

[0014] ここで、BERTは、文書に関するモーダルの情報を用いて問題を解くような状況に適用することが想定されており、複数のモーダルの情報を用いて問題を解くような状況に適用することができない。

[0015] これに対し、例えば、VideoBERTと呼ばれる手法がある。VideoBERTは、具体的には、BERTを、文書に関するモーダルの情報と、画像に関するモーダルの情報とを用いて問題を解くような状況に適用可能に拡張したものである。VideoBERTについては、例えば、下記非特許文献3を参照することができる。

[0016] 非特許文献3 : Sun, Chen, et al. “Videobert: A joint model for video and language representation learning.” arXiv preprint arXiv:1904.01766 (2019).

[0017] また、例えば、MCAN (Modular Co-Attention Network) と呼ばれる手法がある。MCANは、文書に関するモーダルの情報に基づくベクトルと、文書に関するモーダルの情報に基づくベクトルを基に補正した、画像に関するモーダルの情報に基づくベクトルとを参照し、問題を解くものである。MCANについては、例えば、下記非特許文献4を参照することができる。

[0018] 非特許文献4 : Yu, Zhou, et al. “Deep M

odular Co-Attention Networks for Visual Question Answering.” Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019.

[0019] また、例えば、ViLBERT (Vision-and-Language Bidirectional Encoder Representations from Transformers) と呼ばれる手法がある。ViLBERTは、画像に関するモーダルの情報に基づくベクトルを基に補正した、文書に関するモーダルの情報に基づくベクトルと、文書に関するモーダルの情報に基づくベクトルを基に補正した、画像に関するモーダルの情報に基づくベクトルとを参照し、問題を解く技術である。

[0020] 非特許文献5 : Lu, Jiasen, et al. “vilbert: Pretraining task-agnostic vision-linguistic representations for vision-and-language tasks.” arXiv preprint arXiv:1908.02265 (2019).

[0021] しかしながら、上述したVideoBERT、MCAN、およびViLBERTなどの手法でも、複数のモーダルの情報を用いて問題を解いた際の解の精度が悪い場合がある。具体的には、いずれの手法でも、Attentionにより、画像に関するモーダルの情報に基づくベクトルに、文書に関するモーダルの情報に基づくベクトルから得たバリューの重み付け和を、単純に加算するため、問題の解決に有用な情報が失われやすいという性質が存在する。このため、いずれの手法でも、問題を解いた際の解の精度が悪くなりやすい。また、VideoBERTでは、問題を解くにあたり、文書に関するモーダルの情報と、画像に関するモーダルの情報とを明示的に区別せずに扱うため、問題を解いた際の解の精度が悪い。

[0022] そこで、本実施の形態では、問題を解くにあたり有用なベクトルを生成可

能にすることにより、複数のモーダルの情報を用いて問題を解くような状況に適用可能でありつつ、問題を解いた際の解の精度を向上可能にすることができる出力方法について説明する。

[0023] 図1において、出力装置100は、例えば、Attentionを実現する変換モデル110を有する。変換モデルは、生成モデル101と、結合モデル102と、圧縮モデル103と、正規化モデル104とを含む。

[0024] 出力装置100は、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを取得する。モーダルは、情報の様式を意味する。第一のモーダルと、第二のモーダルとは、それぞれ異なるモーダルである。第一のモーダルは、例えば、画像に関するモーダルである。第二のモーダルは、例えば、文書に関するモーダルである。

[0025] 第一のモーダルの情報に基づくベクトルは、例えば、第一のモーダルに従って表現されたベクトルである。第一のモーダルの情報に基づくベクトルは、例えば、第一のモーダルの情報に基づいて生成される。第一のモーダルの情報は、例えば、画像である。第一のモーダルの情報に基づくベクトルは、例えば、画像に基づいて生成されたベクトルである。

[0026] 第二のモーダルの情報に基づくベクトルは、例えば、第二のモーダルに従って表現されたベクトルである。第二のモーダルの情報に基づくベクトルは、例えば、第二のモーダルの情報に基づいて生成される。第二のモーダルの情報は、例えば、文書である。第二のモーダルの情報に基づくベクトルは、例えば、文書に基づいて生成されたベクトルである。

[0027] (1-1) 出力装置100は、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成する。出力装置100は、例えば、生成モデル101を用いて、第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成する。

[0028] 相関は、例えば、第一のモーダルの情報に基づくベクトルから得たベクトルと、第二のモーダルの情報に基づくベクトルから得たベクトルとの類似度

によって表現される。第一のモーダルの情報に基づくベクトルから得たベクトルは、例えば、クエリである。第二のモーダルの情報に基づくベクトルから得たベクトルは、例えば、キーである。類似度は、例えば、内積によって表現される。類似度は、例えば、差分の二乗和などによって表現されてもよい。

[0029] (1-2) 出力装置100は、生成した補正ベクトルを、第一のモーダルの情報に基づくベクトルに結合する。出力装置100は、例えば、結合モデル102を用いて、生成した補正ベクトルを、第一のモーダルの情報に基づくベクトルに結合する。

[0030] (1-3) 出力装置100は、所定のルールに従って、結合後の第一のモーダルの情報に基づくベクトルを圧縮する。出力装置100は、例えば、圧縮モデル103を用いて、結合後の第一のモーダルの情報に基づくベクトルを圧縮する。圧縮は、次元数を低減しない変換を含む。

[0031] (1-4) 出力装置100は、圧縮後の第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施する。出力装置100は、例えば、正規化モデル104を用いて、正規化処理を実施する。正規化処理を実施する具体例については、例えば、図7を用いて後述する。

[0032] (1-5) 出力装置100は、正規化処理により得たベクトルを出力する。出力形式は、例えば、ディスプレイへの表示、プリンタへの印刷出力、他のコンピュータへの送信、または、記憶領域への記憶などである。これにより、出力装置100は、第一のモーダルの情報に基づくベクトルと第二のモーダルの情報に基づくベクトルとのうち問題の解決に有用な情報が反映される傾向があるベクトルを生成し、利用可能にすることができる。結果として、出力装置100は、以降の、問題を解いた際の解の精度の向上を図ることができる。

[0033] ここで、例えば、第一のモーダルが画像に関し、第二のモーダルが文書に関する場合、第二のモーダルは、第一のモーダルの上位の階層であるという特徴を有していると考えることができる。具体的には、「りんご(単語)」

は、複数の「りんご（画像）」を包含する概念である。

[0034] 出力装置100は、この特徴を利用し、画像に関する第一のモーダルの情報に基づくベクトルに、文書に関する第二のモーダルの情報に基づくベクトルに基づく補正ベクトルを結合した上で、圧縮することができる。このため、出力装置100は、圧縮後のベクトルにおいて、画像と文書とのうち問題の解決に有用な情報が失われ辛く、反映され易くすることができる。出力装置100は、例えば、実世界の画像や文書の特徴のうち、問題の解決に有用な特徴を、コンピュータ上で効果的に表現した圧縮後のベクトルを利用可能にすることができる。結果として、出力装置100は、複数のモーダルの情報を用いて問題を解くにあたり、有用なベクトルを得ることができ、問題を解いた際の解の精度を向上可能にすることができる。

[0035] ここでは、第一のモーダルと、第二のモーダルとが、それぞれ異なるモーダルである場合について説明したが、これに限らない。例えば、第一のモーダルと、第二のモーダルとが同一のモーダルである場合があってもよい。

[0036] (情報処理システム200の一例)

次に、図2を用いて、図1に示した出力装置100を適用した、情報処理システム200の一例について説明する。

[0037] 図2は、情報処理システム200の一例を示す説明図である。図2において、情報処理システム200は、出力装置100と、クライアント装置201と、端末装置202とを含む。

[0038] 情報処理システム200において、出力装置100とクライアント装置201とは、有線または無線のネットワーク210を介して接続される。ネットワーク210は、例えば、LAN (Local Area Network)、WAN (Wide Area Network)、インターネットなどである。また、情報処理システム200において、出力装置100と端末装置202とは、有線または無線のネットワーク210を介して接続される。

[0039] 出力装置100は、第一のモーダルの情報に基づくベクトルと、第二のモ

ーダルの情報に基づくベクトルとに基づいて、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを統合した統合ベクトルを生成するCo-Attention Networkを有する。第一のモーダルは、例えば、画像に関するモーダルである。第二のモーダルは、例えば、文書に関するモーダルである。Co-Attention Networkは、例えば、図1に示した変換モデル110を用いて形成される。

[0040] 出力装置100は、教師データに基づいて、Co-Attention Networkを更新する。教師データは、例えば、標本となる第一のモーダルの情報に基づくベクトルを生成する元となる第一のモーダルの情報と、標本となる第二のモーダルの情報に基づくベクトルを生成する元となる第二のモーダルの情報と、正解データとを対応付けた対応情報である。教師データは、例えば、出力装置100のユーザにより出力装置100に入力される。正解データは、例えば、問題を解いた場合についての正解を示す。例えば、第一のモーダルが、画像に関するモーダルであれば、第一のモーダルの情報は、画像である。例えば、第二のモーダルが、文書に関するモーダルであれば、第二のモーダルの情報は、文書である。

[0041] 出力装置100は、例えば、第一のモーダルの情報となる教師データの画像から、第一のモーダルの情報に基づくベクトルを生成することにより取得し、第二のモーダルの情報となる教師データの文書から、第二のモーダルの情報に基づくベクトルを生成することにより取得する。そして、出力装置100は、取得した第一のモーダルの情報に基づくベクトルと、取得した第二のモーダルの情報に基づくベクトルと、教師データの正解データとに基づいて、誤差逆伝搬などにより、Co-Attention Networkを更新する。出力装置100は、誤差逆伝搬以外の学習方法により、Co-Attention Networkを更新してもよい。

[0042] 出力装置100は、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを取得する。そして、出力装置100は、

Co-Attention Networkを用いて、取得した第一のモーダルの情報に基づくベクトルと、取得した第二のモーダルの情報に基づくベクトルとに基づいて、統合ベクトルを生成し、生成した統合ベクトルに基づいて、問題を解く。その後、出力装置100は、問題を解いた結果を、クライアント装置201に送信する。

[0043] 出力装置100は、例えば、出力装置100のユーザにより出力装置100に入力された第一のモーダルの情報に基づくベクトルを取得する。また、出力装置100は、第一のモーダルの情報に基づくベクトルを、クライアント装置201または端末装置202から受信することにより取得してもよい。また、出力装置100は、例えば、第一のモーダルの情報を、クライアント装置201または端末装置202から受信し、受信した第一のモーダルの情報から、第一のモーダルの情報に基づくベクトルを生成することにより取得してもよい。

[0044] 出力装置100は、例えば、出力装置100のユーザにより出力装置100に入力された第二のモーダルの情報に基づくベクトルを取得する。また、出力装置100は、第二のモーダルの情報に基づくベクトルを、クライアント装置201または端末装置202から受信することにより取得してもよい。また、出力装置100は、例えば、第二のモーダルの情報を、クライアント装置201または端末装置202から受信し、受信した第二のモーダルの情報から、第二のモーダルの情報に基づくベクトルを生成することにより取得してもよい。

[0045] そして、出力装置100は、Co-Attention Networkを用いて、取得した第一のモーダルの情報に基づくベクトルと、取得した第二のモーダルの情報に基づくベクトルとに基づいて、統合ベクトルを生成し、生成した統合ベクトルに基づいて、問題を解く。その後、出力装置100は、問題を解いた結果を、クライアント装置201に送信する。出力装置100は、例えば、サーバやPC (Personal Computer) などである。

[0046] クライアント装置201は、出力装置100と通信可能なコンピュータである。クライアント装置201は、例えば、第一のモーダルの情報に基づくベクトルを、出力装置100に送信してもよい。また、クライアント装置201は、例えば、第一のモーダルの情報を、出力装置100に送信してもよい。クライアント装置201は、例えば、第二のモーダルの情報に基づくベクトルを、出力装置100に送信してもよい。また、クライアント装置201は、例えば、第二のモーダルの情報を、出力装置100に送信してもよい。

[0047] クライアント装置201は、出力装置100が問題を解いた結果を受信して出力する。出力形式は、例えば、ディスプレイへの表示、プリンタへの印刷出力、他のコンピュータへの送信、または、記憶領域への記憶などである。クライアント装置201は、例えば、PC、タブレット端末、またはスマートフォンなどである。

[0048] 端末装置202は、出力装置100と通信可能なコンピュータである。端末装置202は、例えば、第一のモーダルの情報に基づくベクトルを、出力装置100に送信してもよい。また、端末装置202は、例えば、第一のモーダルの情報を、出力装置100に送信してもよい。端末装置202は、例えば、第二のモーダルの情報に基づくベクトルを、出力装置100に送信してもよい。また、端末装置202は、例えば、第二のモーダルの情報を、出力装置100に送信してもよい。端末装置202は、例えば、PC、タブレット端末、スマートフォン、電子機器、IoT機器、またはセンサ装置などである。端末装置202は、具体的には、監視カメラであってもよい。

[0049] ここでは、出力装置100が、Co-Attention Networkを更新し、かつ、Co-Attention Networkを用いて、問題を解く場合について説明したが、これに限らない。例えば、他のコンピュータが、Co-Attention Networkを更新し、出力装置100が、他のコンピュータから受信したCo-Attention Networkを用いて、問題を解く場合があってもよい。また、例えば、出力

装置100が、Co-Attention Networkを更新し、他のコンピュータに提供し、他のコンピュータで、Co-Attention Networkを用いて、問題を解く場合があってもよい。

[0050] ここでは、教師データが、第一のモーダルの情報に基づくベクトルを生成する元となる第一のモーダルの情報と、第二のモーダルの情報に基づくベクトルを生成する元となる第二のモーダルの情報と、正解データとを対応付けた対応情報である場合について説明したが、これに限らない。例えば、教師データが、標本となる第一のモーダルの情報に基づくベクトルと、標本となる第二のモーダルの情報に基づくベクトルと、正解データとを対応付けた対応情報である場合があってもよい。

[0051] ここでは、出力装置100が、クライアント装置201や端末装置202とは異なる装置である場合について説明したが、これに限らない。例えば、出力装置100が、クライアント装置201と一体である場合があってもよい。また、例えば、出力装置100が、端末装置202と一体である場合があってもよい。

[0052] ここでは、出力装置100が、ソフトウェア的に、Co-Attention Networkを実現する場合について説明したが、これに限らない。例えば、出力装置100が、Co-Attention Networkを、電子回路的に実現する場合があってもよい。

[0053] (情報処理システム200の適用例1)

適用例1において、出力装置100は、画像と、画像についての質問文となる文書とを記憶する。質問文は、例えば、「画像内で何を切っているか」である。そして、出力装置100は、画像と文書とに基づいて、質問文に対する回答文を推定する問題を解く。出力装置100は、例えば、画像と文書とに基づいて、画像内で何を切っているかの質問文に対する回答文を推定し、クライアント装置201に送信する。

[0054] (情報処理システム200の適用例2)

適用例2において、端末装置202は、監視カメラであり、対象を撮像し

た画像を、出力装置100に送信する。対象は、具体的には、試着室の外観である。また、出力装置100は、対象についての説明文となる文書を記憶している。説明文は、具体的には、人間が試着室を利用中は、試着室のカーテンが閉まっている傾向があることの説明文である。そして、出力装置100は、画像と文書とに基づいて、危険度を判断する問題を解く。危険度は、例えば、試着室に避難が未完了の人間が残っている可能性の高さを示す指標値である。出力装置100は、例えば、災害時に、試着室に避難が未完了の人間が残っている可能性の高さを示す危険度を判断する。

[0055] (情報処理システム200の適用例3)

適用例3において、出力装置100は、動画を形成する画像と、画像についての説明文となる文書を記憶している。動画は、例えば、料理の様子を写した動画である。説明文は、具体的には、料理の手順についての説明文である。そして、出力装置100は、画像と文書とに基づいて、危険度を判断する問題を解く。危険度は、例えば、料理中の危険性の高さを示す指標値である。出力装置100は、例えば、料理中の危険性の高さを示す危険度を判断する。

[0056] (出力装置100のハードウェア構成例)

次に、図3を用いて、出力装置100のハードウェア構成例について説明する。

[0057] 図3は、出力装置100のハードウェア構成例を示すブロック図である。

図3において、出力装置100は、CPU (Central Processing Unit) 301と、メモリ302と、ネットワークI/F (Interface) 303と、記録媒体I/F 304と、記録媒体305とを有する。また、各構成部は、バス300によってそれぞれ接続される。

[0058] ここで、CPU301は、出力装置100の全体の制御を司る。メモリ302は、例えば、ROM (Read Only Memory)、RAM (Random Access Memory) およびフラッシュROMなどを有する。具体的には、例えば、フラッシュROMやROMが各種プログラ

ムを記憶し、RAMがCPU301のワークエリアとして使用される。メモリ302に記憶されるプログラムは、CPU301にロードされることで、コーディングされている処理をCPU301に実行させる。

[0059] ネットワーク1/F303は、通信回線を通じてネットワーク210に接続され、ネットワーク210を介して他のコンピュータに接続される。そして、ネットワーク1/F303は、ネットワーク210と内部のインターフェースを司り、他のコンピュータからのデータの入出力を制御する。ネットワーク1/F303は、例えば、モデムやLANアダプタなどである。

[0060] 記録媒体1/F304は、CPU301の制御に従って記録媒体305に対するデータのリード/ライトを制御する。記録媒体1/F304は、例えば、ディスクドライブ、SSD (Solid State Drive)、USB (Universal Serial Bus) ポートなどである。記録媒体305は、記録媒体1/F304の制御で書き込まれたデータを記憶する不揮発メモリである。記録媒体305は、例えば、ディスク、半導体メモリ、USBメモリなどである。記録媒体305は、出力装置100から着脱可能であってもよい。

[0061] 出力装置100は、上述した構成部のほか、例えば、キーボード、マウス、ディスプレイ、プリンタ、スキャナ、マイク、スピーカーなどを有してもよい。また、出力装置100は、記録媒体1/F304や記録媒体305を複数有していてもよい。また、出力装置100は、記録媒体1/F304や記録媒体305を有していなくてもよい。

[0062] (クライアント装置201のハードウェア構成例)

クライアント装置201のハードウェア構成例は、具体的には、図3に示した出力装置100のハードウェア構成例と同様であるため、説明を省略する。

[0063] (端末装置202のハードウェア構成例)

端末装置202のハードウェア構成例は、具体的には、図3に示した出力装置100のハードウェア構成例と同様であるため、説明を省略する。

[0064] (出力装置100の機能的構成例)

次に、図4を用いて、出力装置100の機能的構成例について説明する。

[0065] 図4は、出力装置100の機能的構成例を示すブロック図である。出力装置100は、記憶部400と、取得部401と、生成部402と、結合部403と、変換部404と、正規化部405と、出力部406とを含む。

[0066] 記憶部400は、例えば、図3に示したメモリ302や記録媒体305などの記憶領域によって実現される。以下では、記憶部400が、出力装置100に含まれる場合について説明するが、これに限らない。例えば、記憶部400が、出力装置100とは異なる装置に含まれ、記憶部400の記憶内容が出力装置100から参照可能である場合があってもよい。

[0067] 取得部401～出力部406は、制御部の一例として機能する。取得部401～出力部406は、具体的には、例えば、図3に示したメモリ302や記録媒体305などの記憶領域に記憶されたプログラムをCPU301に実行させることにより、または、ネットワークI/F303により、その機能を実現する。各機能部の処理結果は、例えば、図3に示したメモリ302や記録媒体305などの記憶領域に記憶される。

[0068] 記憶部400は、各機能部の処理において参照され、または更新される各種情報を記憶する。記憶部400は、Attentionを実現し、第一のモーダルの情報に基づくベクトルを、第二のモーダルの情報に基づくベクトルに基づいて補正し、補正後の第一のモーダルの情報に基づくベクトルを出力する変換モデルを記憶する。

[0069] 例えば、第一のモーダルは、画像に関するモーダルであり、第二のモーダルは、文書に関するモーダルである。例えば、第一のモーダルは、画像に関するモーダルであり、第二のモーダルは、音声に関するモーダルである。例えば、第一のモーダルは、第一の言語の文書に関するモーダルであり、第二のモーダルは、第二の言語の文書に関するモーダルである。例えば、第一のモーダルは、第二のモーダルと同一であってもよい。

[0070] 取得部401は、各機能部の処理に用いられる各種情報を取得する。取得

部401は、取得した各種情報を、記憶部400に記憶し、または、各機能部に出力する。また、取得部401は、記憶部400に記憶しておいた各種情報を、各機能部に出力してもよい。取得部401は、例えば、ユーザの操作入力に基づき、各種情報を取得する。取得部401は、例えば、出力装置100とは異なる装置から、各種情報を受信してもよい。

[0071] 取得部401は、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを取得する。取得部401は、例えば、ユーザによる、第一のモーダルの情報に基づくベクトルを生成する元となる第一のモーダルの情報と、第二のモーダルの情報に基づくベクトルを生成する元となる第二のモーダルの情報との入力を受け付ける。そして、取得部401は、入力された各種情報に基づいて、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを生成する。

[0072] 取得部401は、具体的には、第一のモーダルの情報として、画像を取得し、第一のモーダルの情報に基づくベクトルとして、取得した画像に関する特徴量ベクトルを生成する。画像に関する特徴量ベクトルは、例えば、画像に写る物体ごとの特徴量ベクトルを並べたものである。また、取得部401は、具体的には、第二のモーダルの情報として、文書を取得し、第二のモーダルの情報に基づくベクトルとして、取得した文書に関する特徴量ベクトルを生成する。文書に関する特徴量ベクトルは、例えば、文書に含まれる単語ごとの特徴量ベクトルを並べたものである。

[0073] 取得部401は、例えば、第一のモーダルの情報に基づくベクトルを生成する元となる第一のモーダルの情報と、第二のモーダルの情報に基づくベクトルを生成する元となる第二のモーダルの情報とを、クライアント装置201または端末装置202から受信してもよい。そして、取得部401は、取得した各種情報に基づいて、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを生成する。

[0074] 取得部401は、具体的には、第一のモーダルの情報として、画像を取得し、第一のモーダルの情報に基づくベクトルとして、取得した画像に関する

特徴量ベクトルを生成する。画像に関する特徴量ベクトルは、例えば、画像に写る物体ごとの特徴量ベクトルを並べたものである。また、取得部401は、具体的には、第二のモーダルの情報として、文書を取得し、第二のモーダルの情報に基づくベクトルとして、取得した文書に関する特徴量ベクトルを生成する。文書に関する特徴量ベクトルは、例えば、文書に含まれる単語ごとの特徴量ベクトルを並べたものである。

[0075] 取得部401は、例えば、ユーザによる、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの入力を受け付けることにより、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを取得してもよい。取得部401は、例えば、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを、クライアント装置201または端末装置202から受信することにより取得してもよい。

[0076] 取得部401は、いずれかの機能部の処理を開始する開始トリガーを受け付けてもよい。開始トリガーは、例えば、ユーザによる所定の操作入力があったことである。開始トリガーは、例えば、他のコンピュータから、所定の情報を受信したことであってもよい。開始トリガーは、例えば、いずれかの機能部が所定の情報を出力したことであってもよい。取得部401は、例えば、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとを取得したことを、各機能部の処理を開始する開始トリガーとして受け付ける。

[0077] 生成部402は、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成する。相関は、例えば、第一のモーダルの情報に基づくベクトルから得たベクトルと、第二のモーダルの情報に基づくベクトルから得たベクトルとの類似度によって表現される。第一のモーダルの情報に基づくベクトルから得たベクトルは、例えば、クエリである。第二のモーダルの情報に基づくベクトルから得たベクトルは、例え

ば、キーである。類似度は、例えば、内積によって表現される。類似度は、例えば、差分の二乗和などによって表現されてもよい。

[0078] 生成部402は、例えば、第一のモーダルの情報に基づくベクトルから得たベクトルと、第二のモーダルの情報に基づくベクトルから得たベクトルとの内積に基づいて、補正ベクトルを生成する。生成部402は、具体的には、第一のモーダルの情報に基づくベクトルから得たクエリと、第二のモーダルの情報に基づくベクトルから得たキーとの内積に基づいて、第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成する。

[0079] 生成部402は、より具体的には、画像に関するモーダルの情報に基づくベクトルから得たクエリと、文書に関するモーダルの情報に基づくベクトルから得たキーとの内積に基づいて、画像に関するモーダルの情報に基づくベクトルを補正する補正ベクトルを生成する。ここで、補正ベクトルを生成する一例は、例えば、図7を用いて後述する動作例に示す。これにより、生成部402は、第二のモーダルの情報に基づくベクトルのうち、第一のモーダルの情報に基づくベクトルと相対的に関連深い成分ほど、第一のモーダルの情報に基づくベクトルに強く反映されるように、第一のモーダルの情報に基づくベクトルを補正可能な補正ベクトルを生成することができる。

[0080] 結合部403は、生成した補正ベクトルを、第一のモーダルの情報に基づくベクトルに結合する。結合部403は、例えば、補正ベクトルを、第一のモーダルの情報に基づくベクトルに加算せず、第一のモーダルの前後いずれかに結合する。これにより、結合部403は、第一のモーダルの情報に基づくベクトルと第二のモーダルの情報に基づくベクトルとのうち問題の解決に有用な情報が失われ辛く、反映され易いように、第一のモーダルの情報に基づくベクトルを加工することができる。

[0081] 変換部404は、所定のルールに従って、結合後の第一のモーダルの情報に基づくベクトルを圧縮する。所定のルールは、例えば、学習により自動で設定される。変換部404は、例えば、多層ニューラルネットワークを用いて、結合後の第一のモーダルの情報に基づくベクトルを圧縮する。これによ

り、変換部404は、結合後の第一のモーダルの情報に基づくベクトルの次元数を、扱いやすい次元数に変換することができる。

[0082] 正規化部405は、圧縮後の第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施する。正規化部405は、例えば、第一のモーダルの情報に基づくベクトルと、補正ベクトルとの和を正規化し、当該正規化により得たベクトルと、圧縮後の第一のモーダルの情報に基づくベクトルとの和を正規化する。これにより、正規化部405は、第一のモーダルの情報に基づくベクトルと第二のモーダルの情報に基づくベクトルとのうち問題の解決に有用な情報が効率よく反映された、問題の解決に有用なベクトルを得ることができる。

[0083] 正規化部405は、例えば、結合後の第一のモーダルの情報に基づくベクトルと、圧縮後の第一のモーダルの情報に基づくベクトルとの和を正規化する。これにより、正規化部405は、第一のモーダルの情報に基づくベクトルと第二のモーダルの情報に基づくベクトルとのうち問題の解決に有用な情報が効率よく反映された、問題の解決に有用なベクトルを得ることができる。

[0084] 出力部406は、いずれかの機能部の処理結果を出力する。出力形式は、例えば、ディスプレイへの表示、プリンタへの印刷出力、ネットワークI/F303による外部装置への送信、または、メモリ302や記録媒体305などの記憶領域への記憶である。これにより、出力部406は、各機能部の処理結果をユーザに通知可能にし、出力装置100の利便性の向上を図ることができる。

[0085] 出力部406は、正規化処理により得たベクトルを出力する。これにより、出力部406は、正規化処理により得たベクトルを利用し、Attentionを実現することができる。そして、出力部406は、Attentionにより、Co-Attention Networkを実現可能にすることができる。

[0086] 出力部406は、例えば、Attentionにより、問題の解決に有用

に、正規化処理により得られたベクトルを出力することができる。このため、出力部406は、Co-Attention Networkを、問題の解決に有用になるように学習可能にすることができる。また、出力部406は、問題を解いた際の解の精度を向上可能にすることができる。

[0087] (出力装置100の動作例)

次に、図5～図7を用いて、出力装置100の動作例について説明する。まず、図5を用いて、出力装置100によって用いられるCo-Attention Network 500の具体例について説明する。

[0088] 図5は、Co-Attention Network 500の具体例を示す説明図である。以下の説明では、Co-Attention Network 500を「CAN 500」と表記する場合がある。また、ターゲットアテンションを「TA」と表記する場合がある。また、セルフアテンションを「SA」と表記する場合がある。

[0089] 図5に示すように、CAN 500は、画像TA層501と、画像SA層502と、文書TA層503と、文書SA層504と、結合層505と、統合SA層506とを有する。

[0090] 図5において、CAN 500は、文書に関する特徴量ベクトルLと画像に関する特徴量ベクトルIとが入力されたことに応じて、ベクトルZ_Tを出力する。文書に関する特徴量ベクトルLは、例えば、文書に関するM個の特徴量ベクトルを並べたものである。M個の特徴量ベクトルは、例えば、文書に含まれるM個の単語を示す特徴量ベクトルである。画像に関する特徴量ベクトルIは、例えば、画像に関するN個の特徴量ベクトルを並べたものである。N個の特徴量ベクトルは、例えば、画像に写ったN個の物体を示す特徴量ベクトルである。

[0091] 具体的には、画像TA層501は、画像に関する特徴量ベクトルIと、文書に関する特徴量ベクトルLとの入力を受け付ける。画像TA層501は、画像に関する特徴量ベクトルIから得たクエリと、文書に関する特徴量ベクトルLから得たキーおよびバリューとに基づいて、画像に関する特徴量ベク

トル I を補正する。画像 T A 層 5 0 1 は、補正後の画像に関する特徴量ベクトル I を、画像 S A 層 5 0 2 に出力する。画像 T A 層 5 0 1 の具体例については、例えば、図 7 および図 8 を用いて後述する。

[0092] また、画像 S A 層 5 0 2 は、補正後の画像に関する特徴量ベクトル I の入力を受け付ける。画像 S A 層 5 0 2 は、補正後の画像に関する特徴量ベクトル I から得たクエリ、キーおよびバリューに基づいて、補正後の画像に関する特徴量ベクトル I をさらに補正し、新たな特徴量ベクトル Z_I を生成し、結合層 5 0 5 に出力する。画像 S A 層 5 0 2 を実現する S A 層の具体例については、例えば、図 6 を用いて後述する。

[0093] また、文書 T A 層 5 0 3 は、文書に関する特徴量ベクトル L と、画像に関する特徴量ベクトル I との入力を受け付ける。文書 T A 層 5 0 3 は、文書に関する特徴量ベクトル L から得たクエリと、画像に関する特徴量ベクトル I から得たキーおよびバリューとに基づいて、文書に関する特徴量ベクトル L を補正する。文書 T A 層 5 0 3 は、補正後の文書に関する特徴量ベクトル L を、文書 S A 層 5 0 4 に出力する。文書 T A 層 5 0 3 を実現する T A 層の具体例については、例えば、図 6 を用いて後述する。

[0094] また、文書 S A 層 5 0 4 は、補正後の文書に関する特徴量ベクトル L の入力を受け付ける。文書 S A 層 5 0 4 は、補正後の文書に関する特徴量ベクトル L から得たクエリ、キーおよびバリューに基づいて、補正後の文書に関する特徴量ベクトル L をさらに補正し、新たな特徴量ベクトル Z_L を生成して出力する。文書 S A 層 5 0 4 を実現する S A 層の具体例については、例えば、図 6 を用いて後述する。

[0095] また、結合層 5 0 5 は、集約用ベクトル H と、特徴量ベクトル Z_I と、特徴量ベクトル Z_L との入力を受け付ける。結合層 5 0 5 は、集約用ベクトル H と、特徴量ベクトル Z_I と、特徴量ベクトル Z_L とを結合し、結合ベクトル C を生成し、統合 S A 層 5 0 6 に出力する。

[0096] また、統合 S A 層 5 0 6 は、結合ベクトル C の入力を受け付ける。統合 S A 層 5 0 6 は、結合ベクトル C から得たクエリ、キーおよびバリューに基づ

いて、結合ベクトルCを補正し、特徴量ベクトル Z_T を生成して出力する。特徴量ベクトル Z_T は、集約ベクトル Z_H と、文書に関する統合特徴量ベクトル $Z_1 \sim Z_M$ と、画像に関する統合特徴量ベクトル $Z_{M+1} \sim Z_{M+N}$ とを含む。これにより、出力装置100は、問題を解いた際の解の精度を向上させる観点で有用な集約ベクトル Z_H を含む特徴量ベクトル Z_T を生成し、参照可能にすることができる。このため、出力装置100は、問題を解いた際の解の精度を向上可能にすることができる。

[0097] ここでは、説明の簡略化のため、画像TA層501と、画像SA層502と、文書TA層503と、文書SA層504とのグループ510が、1段である場合について説明したが、これに限らない。例えば、画像TA層501と、画像SA層502と、文書TA層503と、文書SA層504とのグループ510が、複数段存在する場合があってもよい。これによれば、出力装置100は、問題を解いた際の解の精度のさらなる向上を図ることができる。

[0098] ここでは、CAN500が、画像TA層501と、画像SA層502と、文書TA層503と、文書SA層504と、結合層505と、統合SA層506とを有する場合について説明したが、これに限らない。例えば、CAN500が、結合層505と、統合SA層506とを有していない場合があってもよい。この場合、出力装置100は、例えば、問題を解くにあたり、画像SA層502の出力と、文書SA層504の出力とを利用する。

[0099] 次に、図6の説明に移行し、CAN500を形成する画像SA層502と文書SA層504と統合SA層506となどを実現するSA層600の具体例と、CAN500を形成する文書TA層503などを実現するTA層610の具体例とについて説明する。CAN500を形成する画像TA層501の具体例については、図7を用いて後述する。

[0100] 図6は、SA層600の具体例と、TA層610の具体例とを示す説明図である。以下の説明では、Multi-Head Attentionを「MHA」と表記する場合がある。また、Add&Normを「A&N」と表

記する場合がある。また、Feed Forwardを「FF」と表記する場合がある。

[0101] 図6に示すように、SA層600は、MHA層601と、A&N層602と、FF層603と、A&N層604とを有する。MHA層601は、入力ベクトルXから得たクエリQとキーKとバリューVとに基づいて、入力ベクトルXを補正する補正ベクトルRを生成し、A&N層602に出力する。MHA層601は、具体的には、入力ベクトルXを、Head個のベクトルに分割して処理する。Headは、1以上の自然数である。

[0102] A&N層602は、入力ベクトルXと補正ベクトルRとを加算した上で正規化し、正規化後のベクトルを、FF層603とA&N層604とに出力する。FF層603は、正規化後のベクトルを圧縮し、圧縮後のベクトルを、A&N層604に出力する。A&N層604は、正規化後のベクトルと、圧縮後のベクトルとを加算した上で正規化し、出力ベクトルZを生成して出力する。

[0103] また、TA層610は、MHA層611と、A&N層612と、FF層613と、A&N層614とを有する。MHA層611は、入力ベクトルXから得たクエリQと、入力ベクトルYから得たキーKとバリューVとに基づいて、入力ベクトルXを補正する補正ベクトルRを生成し、A&N層612に出力する。A&N層612は、入力ベクトルXと補正ベクトルRとを加算した上で正規化し、正規化後のベクトルを、FF層613とA&N層614とに出力する。FF層613は、正規化後のベクトルを圧縮し、圧縮後のベクトルを、A&N層614に出力する。A&N層614は、正規化後のベクトルと、圧縮後のベクトルとを加算した上で正規化し、出力ベクトルZを生成して出力する。

[0104] 上述したMHA層601やMHA層611は、より具体的には、Headの個数分のAttention層620により形成される。Attention層620は、MatMul層621と、Scale層622と、Mask層623と、SoftMax層624と、MatMul層625とを有す

る。

[0105] MatMul層621は、クエリQとキーKとの内積を算出し、Scoreに設定する。Scale層622は、Score全体を定数aで除算し、更新する。Mask層623は、更新後のScoreをマスク処理してもよい。SoftMax層624は、更新後のScoreを、正規化し、Attに設定する。MatMul層625は、AttとバリューVとの内積を算出し、補正ベクトルRに設定する。次に、図7および図8を用いて、CAN500を形成する画像TA層501の具体例について説明する。

[0106] 図7は、画像TA層501の具体例を示す説明図である。図7において、画像TA層501は、MHA層701と、A&N層702と、Con層703と、FF層704と、A&N層705とを含む。MHA層701は、入力ベクトルXから得たクエリQと、入力ベクトルYから得たキーKとバリューVとに基づいて、入力ベクトルXを補正する補正ベクトルRを生成し、A&N層702およびCon層703に出力する。A&N層702は、入力ベクトルXと補正ベクトルRとを加算した上で正規化し、正規化後のベクトルを、A&N層705に出力する。

[0107] Con層703は、入力ベクトルXと補正ベクトルRとを結合し、結合ベクトルをFF層704に出力する。FF層704は、結合ベクトルを圧縮し、圧縮後のベクトルを、A&N層705に出力する。A&N層705は、正規化後のベクトルと、圧縮後のベクトルとを加算した上で正規化し、正規化で得た出力ベクトルを出力する。次に、図8を用いて、画像TA層501の別の具体例について説明する。

[0108] 図8は、画像TA層501の別の具体例を示す説明図である。図8において、画像TA層501は、MHA層801と、Con層802と、FF層803と、A&N層804とを含む。MHA層801は、入力ベクトルXから得たクエリQと、入力ベクトルYから得たキーKとバリューVとに基づいて、入力ベクトルXを補正する補正ベクトルRを生成し、Con層802に出力する。

- [0109] Con層802は、入力ベクトルXと補正ベクトルRとを結合し、結合ベクトルをFF層803およびA&N層804に出力する。FF層803は、結合ベクトルを圧縮し、圧縮後のベクトルを、A&N層804に出力する。A&N層804は、結合ベクトルと、圧縮後のベクトルとを加算した上で正規化し、正規化で得た出力ベクトルを出力する。次に、図9を用いて、画像TA層501と文書TA層503との比較例について説明する。
- [0110] 図9は、画像TA層501と文書TA層503との比較例を示す説明図である。図9に示すように、画像TA層501と、文書TA層503とは、文書に関する特徴量ベクトルLと、画像に関する特徴量ベクトルIとの入力を受け付ける。しかしながら、画像TA層501と、文書TA層503とは、それぞれ、異なる手法で、文書に関する特徴量ベクトルLと、画像に関する特徴量ベクトルIとを扱うことになる。
- [0111] 例えば、画像TA層501は、画像に関する特徴量ベクトルIに、ベクトル Z_{I1} を結合することにより、新たな特徴量ベクトル Z_{I2} を生成する。一方で、文書TA層503は、文書に関する特徴量ベクトルLに、ベクトル Z_{L1} を加算することにより、新たな特徴量ベクトル Z_{L2} を生成する。これにより、出力装置100は、それぞれ性質が異なる、文書に関する特徴量ベクトルLと、画像に関する特徴量ベクトルIとに対し、異なる扱い方をすることができる。
- [0112] そして、出力装置100は、画像TA層501において、文書に関する特徴量ベクトルLと、画像に関する特徴量ベクトルIとのうち、問題の解決に有用な情報が失われ辛くすることができる。結果として、出力装置100は、複数のモーダルの情報を用いて問題を解くにあたり有用なベクトルを得ることができ、問題を解いた際の解の精度を向上可能にすることができる。
- [0113] ここでは、画像TA層501を、図7および図8に示す具体例のように形成する場合について説明したが、これに限らない。例えば、画像SA層502と、文書TA層503と、文書SA層504と、統合SA層506との少なくともいずれかを、図7および図8に示す具体例と同様に形成する場合が

あってもよい。次に、図10を用いて、出力装置100による、CAN500を用いた動作の一例について説明する。

[0114] 図10は、CAN500を用いた動作の一例を示す説明図である。図10において、出力装置100は、文書1000を取得し、画像1010を取得する。出力装置100は、文書1000をトークン化し、トークン集合1001をベクトル化し、文書1000に関する特徴量ベクトル1002を生成し、CAN500に入力する。また、出力装置100は、画像1010から物体を検出し、物体ごとの部分画像の集合1011をベクトル化し、画像1010に関する特徴量ベクトル1012を生成し、CAN500に入力する。

[0115] 出力装置100は、CAN500から、特徴量ベクトル Z_T を取得し、特徴量ベクトル Z_T に含まれる集約ベクトル Z_H を、危険度推定器1030に入力する。出力装置100は、危険度推定器1030から推定結果 N_O を取得する。これにより、出力装置100は、画像と文書との特徴が反映された集約ベクトル Z_H を用いて、危険度推定器1030に危険であるか否かを推定させることができ、危険であるか否かを精度よく推定可能にすることができる。危険度推定器1030は、例えば、銃を持った人物が写っている画像1010があるが、ミュージアムの展示物であることを示す文書もあるため、推定結果 N_O = 危険ではないと推定することができる。

[0116] (出力装置100の利用例)

次に、図11～図14を用いて、出力装置100の利用例について説明する。

[0117] 図11および図12は、出力装置100の利用例1を示す説明図である。図11において、出力装置100は、学習フェーズを実施し、CAN500を学習する。出力装置100は、例えば、何らかのシーンを写した画像1100と、画像1100に対応する字幕となる文書1110とを取得する。画像1100は、例えば、りんごを切るシーンを写す。

[0118] 出力装置100は、画像1100を変換器1120により特徴量ベクトル

に変換し、CAN500に入力する。また、出力装置100は、文書1110の単語appleをマスクした上で、変換器1130により特徴量ベクトルに変換し、CAN500に入力する。

[0119] 出力装置100は、CAN500により生成された特徴量ベクトルを、識別器1140に入力し、マスクされた単語を予測した結果を取得し、マスクされた単語の正解「apple」との誤差を算出する。出力装置100は、算出した誤差に基づいて、誤差逆伝搬によりCAN500を学習する。さらに、出力装置100は、誤差逆伝搬により、変換器1120, 1130や識別器1140を学習してもよい。

[0120] これにより、出力装置100は、画像1100と字幕となる文書1110の文脈とを考慮して単語を推定する観点で有用なように、CAN500、および変換器1120, 1130や識別器1140を更新することができる。次に、図12の説明に移行する。

[0121] 図12において、出力装置100は、試験フェーズを実施し、学習した変換器1120, 1130と、学習したCAN500とを用いて、回答を生成して出力する。出力装置100は、例えば、何らかのシーンを写した画像1200と、画像1200に対応する質問文となる文書1210とを取得する。画像1200は、例えば、りんごを切るシーンを写す。

[0122] 出力装置100は、画像1200を変換器1120により特徴量ベクトルに変換し、CAN500に入力する。また、出力装置100は、文書1210を変換器1130により特徴量ベクトルに変換し、CAN500に入力する。出力装置100は、CAN500により生成された特徴量ベクトルを、回答生成器1220に入力し、回答となる単語を取得して出力する。これにより、出力装置100は、画像1200と質問文となる文書1210の文脈とを考慮して、精度よく回答となる単語を推定することができる。

[0123] 図13および図14は、出力装置100の利用例2を示す説明図である。図13において、出力装置100は、学習フェーズを実施し、CAN500を学習する。出力装置100は、例えば、何らかのシーンを写した画像13

00と、画像1300に対応する字幕となる文書1310とを取得する。画像1300は、例えば、りんごを切るシーンを写す。

[0124] 出力装置100は、画像1300を変換器1320により特徴量ベクトルに変換し、CAN500に入力する。また、出力装置100は、文書1310の単語appleをマスクした上で、変換器1330により特徴量ベクトルに変換し、CAN500に入力する。

[0125] 出力装置100は、CAN500により生成された特徴量ベクトルを、識別器1340に入力し、画像に写ったシーンの危険度を予測した結果を取得し、危険度の正解との誤差を算出する。出力装置100は、算出した誤差に基づいて、誤差逆伝搬によりCAN500を学習する。また、出力装置100は、誤差逆伝搬により、変換器1320, 1330や識別器1340を学習する。

[0126] これにより、出力装置100は、画像1300と字幕となる文書1310の文脈とを考慮して危険度を予測する観点で有用なように、CAN500、および変換器1120, 1130や識別器1140を更新することができる。次に、図14の説明に移行する。

[0127] 図14において、出力装置100は、試験フェーズを実施し、学習した変換器1320, 1330や識別器1340と、学習したCAN500とを用いて、危険度を予測して出力する。出力装置100は、例えば、何らかのシーンを写した画像1400と、画像に対応する説明文となる文書1410とを取得する。画像1400は、例えば、ももを切るシーンを写す。

[0128] 出力装置100は、画像1400を変換器1320により特徴量ベクトルに変換し、CAN500に入力する。また、出力装置100は、文書1410を変換器1330により特徴量ベクトルに変換し、CAN500に入力する。出力装置100は、CAN500により生成された特徴量ベクトルを、識別器1340に入力し、危険度を取得して出力する。これにより、出力装置100は、画像1400と説明文となる文書1410の文脈とを考慮して、精度よく危険度を予測することができる。

[0129] (学習処理手順)

次に、図15を用いて、出力装置100が実行する、学習処理手順の一例について説明する。学習処理は、例えば、図3に示したCPU301と、メモリ302や記録媒体305などの記憶領域と、ネットワークI/F303とによって実現される。

[0130] 図15は、学習処理手順の一例を示すフローチャートである。図15において、出力装置100は、画像の特徴量ベクトルと、文書の特徴量ベクトルとを取得する(ステップS1501)。

[0131] 次に、出力装置100は、取得した画像の特徴量ベクトルから生成したクエリと、取得した文書の特徴量ベクトルから生成したキーおよびバリューとに基づいて、画像TA層501を用いて、画像の特徴量ベクトルを補正する(ステップS1502)。ここで、出力装置100は、具体的には、図14に後述するアテンション処理を実行することにより、画像の特徴量ベクトルを補正する。

[0132] そして、出力装置100は、補正後の画像の特徴量ベクトルに基づいて、画像SA層502を用いて、補正後の画像の特徴量ベクトルをさらに補正し、新たに画像の特徴量ベクトルを生成する(ステップS1503)。

[0133] 次に、出力装置100は、取得した文書の特徴量ベクトルから生成したクエリと、取得した画像の特徴量ベクトルから生成したキーおよびバリューとに基づいて、文書TA層503を用いて、文書の特徴量ベクトルを補正する(ステップS1504)。

[0134] そして、出力装置100は、補正後の文書の特徴量ベクトルに基づいて、文書SA層504を用いて、補正後の文書の特徴量ベクトルをさらに補正し、新たに文書の特徴量ベクトルを生成する(ステップS1505)。

[0135] 次に、出力装置100は、集約用ベクトルを初期化する(ステップS1506)。そして、出力装置100は、集約用ベクトルと、生成した画像の特徴量ベクトルと、生成した文書の特徴量ベクトルとを結合し、結合ベクトルを生成する(ステップS1507)。

[0136] 次に、出力装置100は、結合ベクトルに基づいて、統合SA層506を用いて、結合ベクトルを補正し、集約ベクトルを生成する（ステップS1508）。そして、出力装置100は、集約ベクトルに基づいて、CAN500を学習する（ステップS1509）。

[0137] その後、出力装置100は、学習処理を終了する。これにより、出力装置100は、CAN500を用いて問題を解くにあたり、問題を解いた際の解の精度が向上するように、CAN500のパラメータを更新することができる。

[0138] ここで、出力装置100は、図15の一部ステップの処理の順序を入れ替えて実行してもよい。例えば、ステップS1502、S1503の処理と、ステップS1504、S1505の処理との順序は入れ替え可能である。また、出力装置100は、ステップS1502～S1505の処理を繰り返し実行してもよい。

[0139] （推定処理手順）

次に、図16を用いて、出力装置100が実行する、推定処理手順の一例について説明する。推定処理は、例えば、図3に示したCPU301と、メモリ302や記録媒体305などの記憶領域と、ネットワークI/F303とによって実現される。

[0140] 図16は、推定処理手順の一例を示すフローチャートである。図16において、出力装置100は、画像の特徴量ベクトルと、文書の特徴量ベクトルとを取得する（ステップS1601）。

[0141] 次に、出力装置100は、取得した画像の特徴量ベクトルから生成したクエリと、取得した文書の特徴量ベクトルから生成したキーおよびバリューとに基づいて、画像TA層501を用いて、画像の特徴量ベクトルを補正する（ステップS1602）。ここで、出力装置100は、具体的には、図14に後述するアテンション処理を実行することにより、画像の特徴量ベクトルを補正する。

[0142] そして、出力装置100は、補正後の画像の特徴量ベクトルに基づいて、

画像 S A 層 5 0 2 を用いて、補正後の画像の特徴量ベクトルをさらに補正し、新たに画像の特徴量ベクトルを生成する（ステップ S 1 6 0 3）。

[0143] 次に、出力装置 1 0 0 は、取得した文書の特徴量ベクトルから生成したクエリと、取得した画像の特徴量ベクトルから生成したキーおよびバリューとに基づいて、文書 T A 層 5 0 3 を用いて、文書の特徴量ベクトルを補正する（ステップ S 1 6 0 4）。

[0144] そして、出力装置 1 0 0 は、補正後の文書の特徴量ベクトルに基づいて、文書 S A 層 5 0 4 を用いて、補正後の文書の特徴量ベクトルをさらに補正し、新たに文書の特徴量ベクトルを生成する（ステップ S 1 6 0 5）。

[0145] 次に、出力装置 1 0 0 は、集約用ベクトルを初期化する（ステップ S 1 6 0 6）。そして、出力装置 1 0 0 は、集約用ベクトルと、生成した画像の特徴量ベクトルと、生成した文書の特徴量ベクトルとを結合し、結合ベクトルを生成する（ステップ S 1 6 0 7）。

[0146] 次に、出力装置 1 0 0 は、結合ベクトルに基づいて、統合 S A 層 5 0 6 を用いて、結合ベクトルを補正し、集約ベクトルを生成する（ステップ S 1 6 0 8）。そして、出力装置 1 0 0 は、集約ベクトルに基づいて、識別モデルを用いて、状況を推定する（ステップ S 1 6 0 9）。

[0147] 次に、出力装置 1 0 0 は、推定した状況を出力する（ステップ S 1 6 1 0）。そして、出力装置 1 0 0 は、推定処理を終了する。これにより、出力装置 1 0 0 は、C A N 5 0 0 を用いて、問題を解いた際の解の精度を向上させることができる。

[0148] ここで、出力装置 1 0 0 は、図 1 6 の一部ステップの処理の順序を入れ替えて実行してもよい。例えば、ステップ S 1 6 0 2, S 1 6 0 3 の処理と、ステップ S 1 6 0 4, S 1 6 0 5 の処理との順序は入れ替え可能である。また、出力装置 1 0 0 は、ステップ S 1 6 0 2 ~ S 1 6 0 5 の処理を繰り返し実行してもよい。

[0149] (アテンション処理手順)

次に、図 1 7 を用いて、画像 T A 層により、出力装置 1 0 0 が実行する、

アテンション処理手順の一例について説明する。アテンション処理は、例えば、図3に示したCPU301と、メモリ302や記録媒体305などの記憶領域と、ネットワークI/F303とによって実現される。

[0150] 図17は、アテンション処理手順の一例を示すフローチャートである。図17において、出力装置100は、ベクトルXとなる画像の特徴量ベクトルと、ベクトルYとなる文書の特徴量ベクトルとを取得する（ステップS1701）。

[0151] 次に、出力装置100は、取得した画像の特徴量ベクトルからベクトルQueryを生成する（ステップS1702）。そして、出力装置100は、取得した文書の特徴量ベクトルからベクトルkeyとベクトルValueを生成する（ステップS1703）。

[0152] 次に、出力装置100は、生成したベクトルQueryと、生成したベクトルkeyとの内積を算出する（ステップS1704）。そして、出力装置100は、内積のsoftmaxによりベクトルAttを生成する（ステップS1705）。

[0153] 次に、出力装置100は、ベクトルAttとベクトルValueとの内積によりベクトルRを生成する（ステップS1706）。そして、出力装置100は、ベクトルRとベクトルXとを結合したベクトルX'を生成する（ステップS1707）。

[0154] 次に、出力装置100は、多層ニューラルネットワークにより、ベクトルX'を、ベクトルXと同じ次元に圧縮し、ベクトルX''を生成する（ステップS1708）。そして、出力装置100は、ベクトルRとベクトルXとを用いて、ベクトルX''を正規化し、正規化後のベクトルを取得する（ステップS1709）。

[0155] 次に、出力装置100は、取得した正規化後のベクトルを出力する（ステップS1710）。そして、出力装置100は、アテンション処理を終了する。これにより、出力装置100は、画像と文書とのうち問題の解決に有用な情報が失われ辛いように、正規化後のベクトルを生成して取得することが

できる。

[0156] ここで、出力装置100は、図17の一部ステップの処理の順序を入れ替えて実行してもよい。例えば、ステップS1702の処理と、ステップS1703の処理との順序は入れ替え可能である。

[0157] 以上説明したように、出力装置100によれば、第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成することができる。出力装置100によれば、生成した補正ベクトルを、第一のモーダルの情報に基づくベクトルに結合することができる。出力装置100によれば、所定のルールに従って、結合後の第一のモーダルの情報に基づくベクトルを圧縮することができる。出力装置100によれば、圧縮後の第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施することができる。出力装置100によれば、正規化処理により得たベクトルを出力することができる。これにより、出力装置100は、第一のモーダルの情報に基づくベクトルと第二のモーダルの情報に基づくベクトルとのうち問題の解決に有用な情報を残して、問題を解くのに有用なベクトルを得ることができ、問題を解いた際の解の精度を向上可能にすることができる。

[0158] 出力装置100によれば、第一のモーダルの情報に基づくベクトルから得たベクトルと、第二のモーダルの情報に基づくベクトルから得たベクトルとの内積に基づいて、補正ベクトルを生成することができる。これにより、出力装置100は、アテンションを実現することができる。また、出力装置100は、問題を解くのに有用な補正ベクトルを得ることができる。

[0159] 出力装置100によれば、第一のモーダルの情報に基づくベクトルと、補正ベクトルとの和を正規化し、当該正規化により得たベクトルと、圧縮後の第一のモーダルの情報に基づくベクトルとの和を正規化することができる。これにより、出力装置100は、正規化処理を実現することができる。

[0160] 出力装置100によれば、結合後の第一のモーダルの情報に基づくベクトルと、圧縮後の第一のモーダルの情報に基づくベクトルとの和を正規化する

ことができる。これにより、出力装置100は、正規化処理を実現することができる。

[0161] 出力装置100によれば、第一のモーダルとして、画像に関するモーダルを採用することができる。出力装置100によれば、第二のモーダルとして、文書に関するモーダルを採用することができる。これにより、出力装置100は、ターゲットアテンション層を実現することができる。また、出力装置100は、画像と文書とに基づいて問題を解く場合に適用可能にすることができる。

[0162] 出力装置100によれば、第一のモーダルとして、画像に関するモーダルを採用することができる。出力装置100によれば、第二のモーダルとして、音声に関するモーダルを採用することができる。これにより、出力装置100は、ターゲットアテンション層を実現することができる。また、出力装置100は、画像と音声とに基づいて問題を解く場合に適用可能にすることができる。

[0163] 出力装置100によれば、第一のモーダルとして、第一の言語の文書に関するモーダルを採用することができる。出力装置100によれば、第二のモーダルとして、第二の言語の文書に関するモーダルを採用することができる。これにより、出力装置100は、ターゲットアテンション層を実現することができる。また、出力装置100は、異なる言語の2つの文書に基づいて問題を解く場合に適用可能にすることができる。

[0164] 出力装置100によれば、第一のモーダルと、第二のモーダルとに、同一のモーダルを採用することができる。これにより、出力装置100は、セルフアテンション層を実現することができる。また、出力装置100は、同一のモーダルの異なる情報に基づいて問題を解く場合に適用可能にすることができる。

[0165] なお、本実施の形態で説明した出力方法は、予め用意されたプログラムをPCやワークステーションなどのコンピュータで実行することにより実現することができる。本実施の形態で説明した出力プログラムは、コンピュータ

で読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行される。記録媒体は、ハードディスク、フレキシブルディスク、CD (Compact Disc) -ROM、MO、DVD (Digital Versatile Disc) などである。また、本実施の形態で説明した出力プログラムは、インターネットなどのネットワークを介して配布してもよい。

符号の説明

- [0166] 100 出力装置
- 101 生成モデル
- 102 結合モデル
- 103 圧縮モデル
- 104 正規化モデル
- 110 変換モデル
- 200 情報処理システム
- 201 クライアント装置
- 202 端末装置
- 210 ネットワーク
- 300 バス
- 301 CPU
- 302 メモリ
- 303 ネットワーク I/F
- 304 記録媒体 I/F
- 305 記録媒体
- 400 記憶部
- 401 取得部
- 402 生成部
- 403 結合部
- 404 変換部

405 正規化部
406 出力部
500 CAN
501 画像TA層
502 画像SA層
503 文書TA層
504 文書SA層
505 結合層
506 統合SA層
510 グループ
600 SA層
601, 611, 701, 801 MHA層
602, 604, 612, 614, 702, 705, 804 A&N層
603, 613, 704, 803 FF層
610 TA層
620 Attention層
621, 625 MatMul層
622 Scale層
623 Mask層
624 SoftMax層
703, 802 Con層
1000, 1110, 1210, 1310, 1410 文書
1001 トークン集合
1002, 1012 特徴量ベクトル
1010, 1100, 1200, 1300, 1400 画像
1011 集合
1030 危険度推定器
1120, 1130, 1320, 1330 変換器

1 1 4 0, 1 3 4 0 識別器

1 2 2 0 回答生成器

請求の範囲

- [請求項1] 第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、前記第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成し、
- 生成した前記補正ベクトルを、前記第一のモーダルの情報に基づくベクトルに結合し、
- 所定のルールに従って、結合後の前記第一のモーダルの情報に基づくベクトルを圧縮し、
- 圧縮後の前記第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施し、
- 前記正規化処理により得たベクトルを出力する、
- 処理をコンピュータが実行することを特徴とする出力方法。
- [請求項2] 前記生成する処理は、
- 前記第一のモーダルの情報に基づくベクトルから得たベクトルと、前記第二のモーダルの情報に基づくベクトルから得たベクトルとの内積に基づいて、前記補正ベクトルを生成する、ことを特徴とする請求項1に記載の出力方法。
- [請求項3] 前記正規化処理を実施する処理は、
- 前記第一のモーダルの情報に基づくベクトルと、前記補正ベクトルとの和を正規化し、
- 当該正規化により得たベクトルと、圧縮後の前記第一のモーダルの情報に基づくベクトルとの和を正規化する、ことを特徴とする請求項1または2に記載の出力方法。
- [請求項4] 前記正規化処理を実施する処理は、
- 結合後の前記第一のモーダルの情報に基づくベクトルと、圧縮後の前記第一のモーダルの情報に基づくベクトルとの和を正規化する、ことを特徴とする請求項1または2に記載の出力方法。
- [請求項5] 前記第一のモーダルと前記第二のモーダルとの組は、画像に関する

モーダルと文書に関するモーダルとの組、画像に関するモーダルと音声に関するモーダルとの組、第一の言語の文書に関するモーダルと第二の言語の文書に関するモーダルとの組のうちいずれかの組である、ことを特徴とする請求項1～4のいずれか一つに記載の出力方法。

[請求項6] 前記第一のモーダルは、前記第二のモーダルと同一である、ことを特徴とする請求項1～4のいずれか一つに記載の出力方法。

[請求項7] 第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、前記第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成し、

生成した前記補正ベクトルを、前記第一のモーダルの情報に基づくベクトルに結合し、

所定のルールに従って、結合後の前記第一のモーダルの情報に基づくベクトルを圧縮し、

圧縮後の前記第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施し、

前記正規化処理により得たベクトルを出力する、

処理をコンピュータに実行させることを特徴とする出力プログラム

。

[請求項8] 第一のモーダルの情報に基づくベクトルと、第二のモーダルの情報に基づくベクトルとの相関に基づいて、前記第一のモーダルの情報に基づくベクトルを補正する補正ベクトルを生成し、

生成した前記補正ベクトルを、前記第一のモーダルの情報に基づくベクトルに結合し、

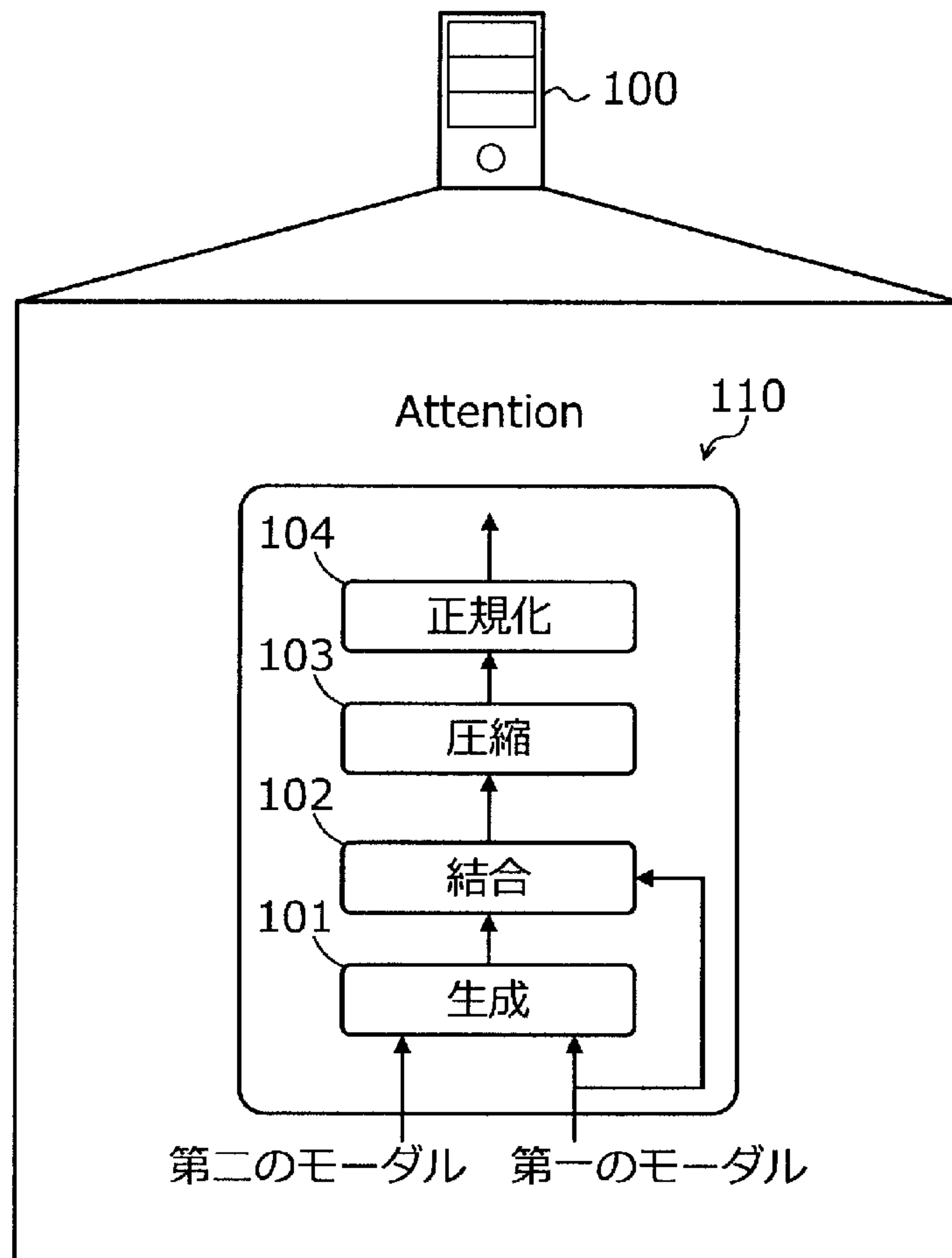
所定のルールに従って、結合後の前記第一のモーダルの情報に基づくベクトルを圧縮し、

圧縮後の前記第一のモーダルの情報に基づくベクトルに対して、正規化処理を実施し、

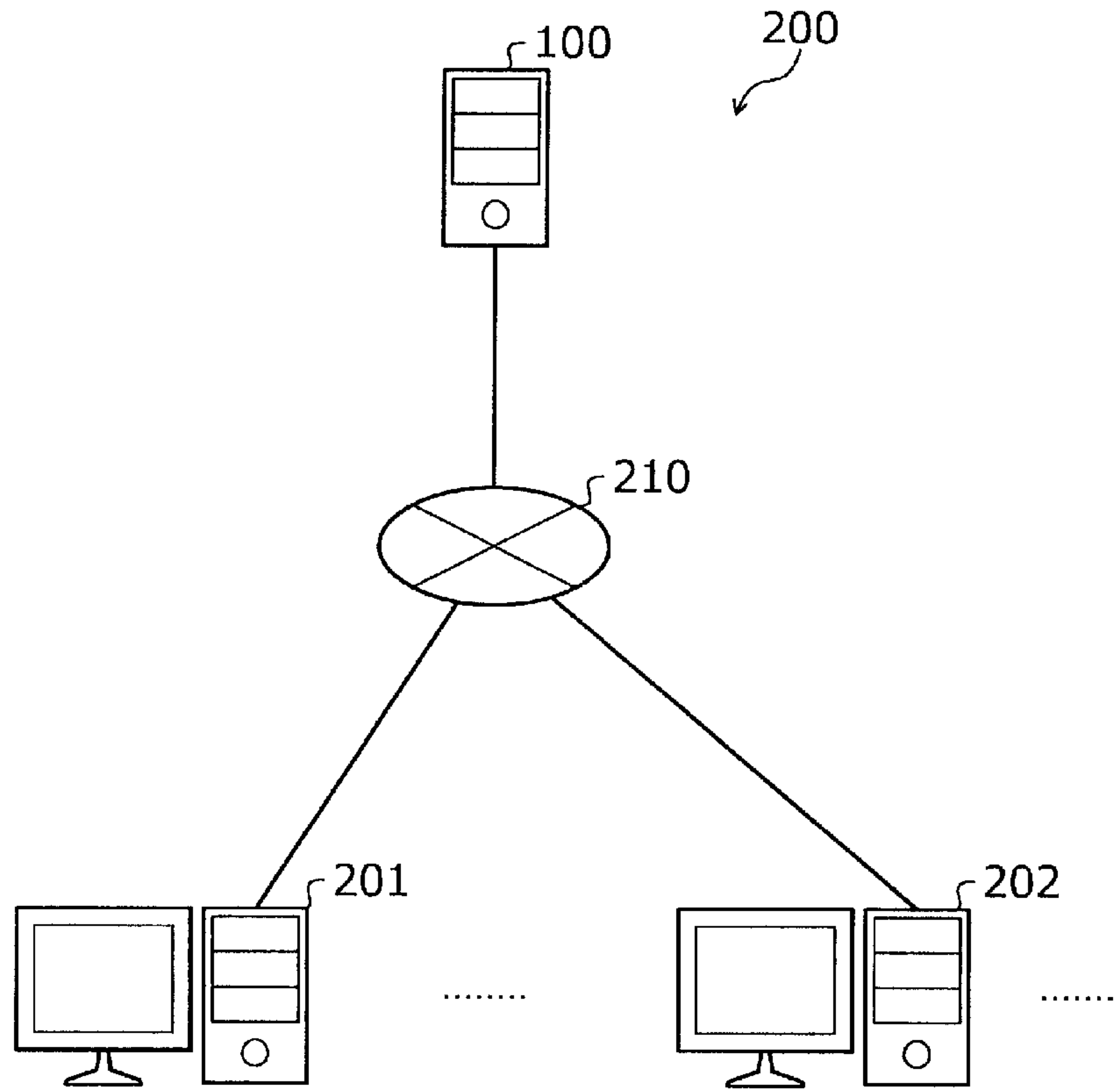
前記正規化処理により得たベクトルを出力する、

制御部を有することを特徴とする出力装置。

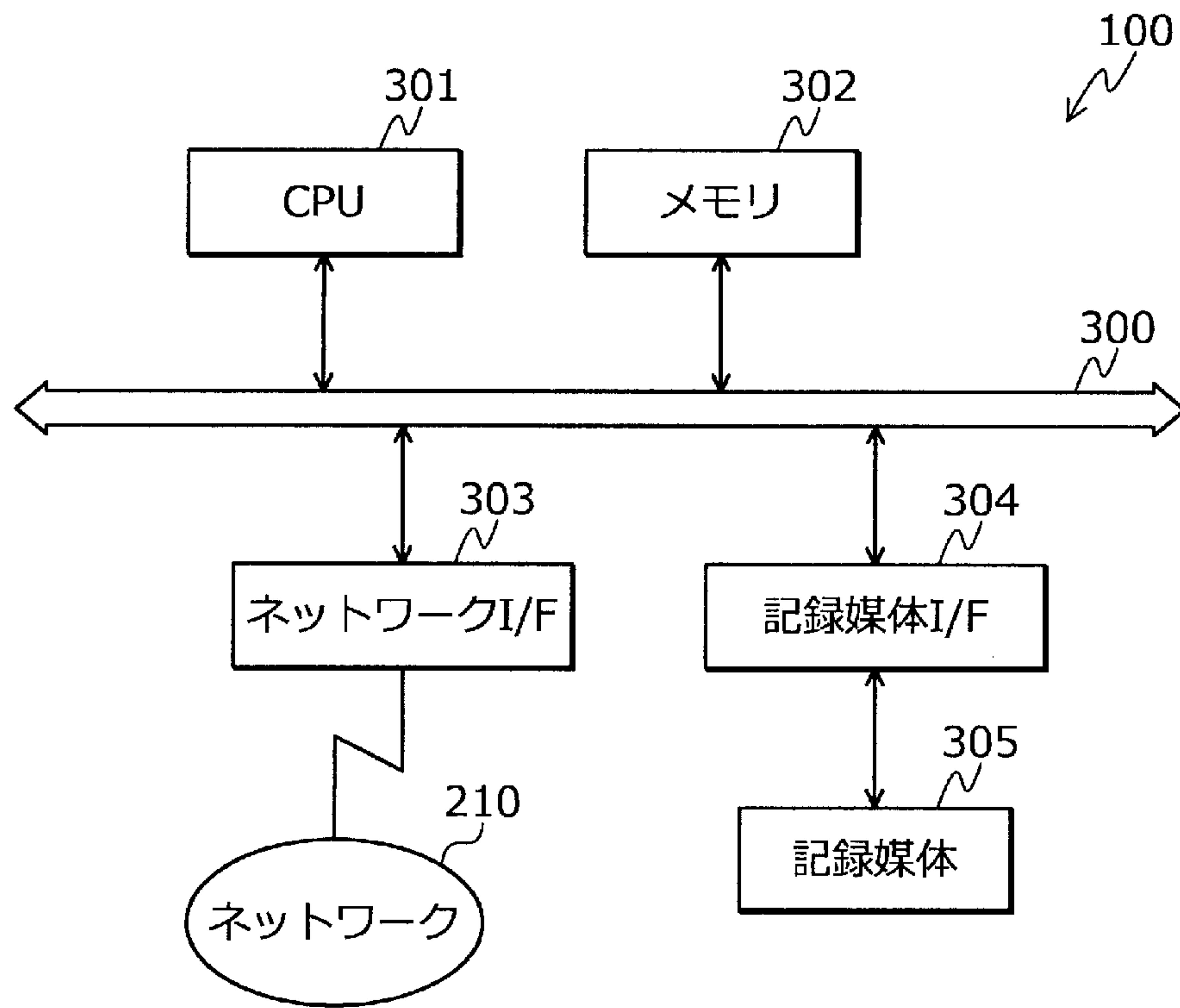
[図1]



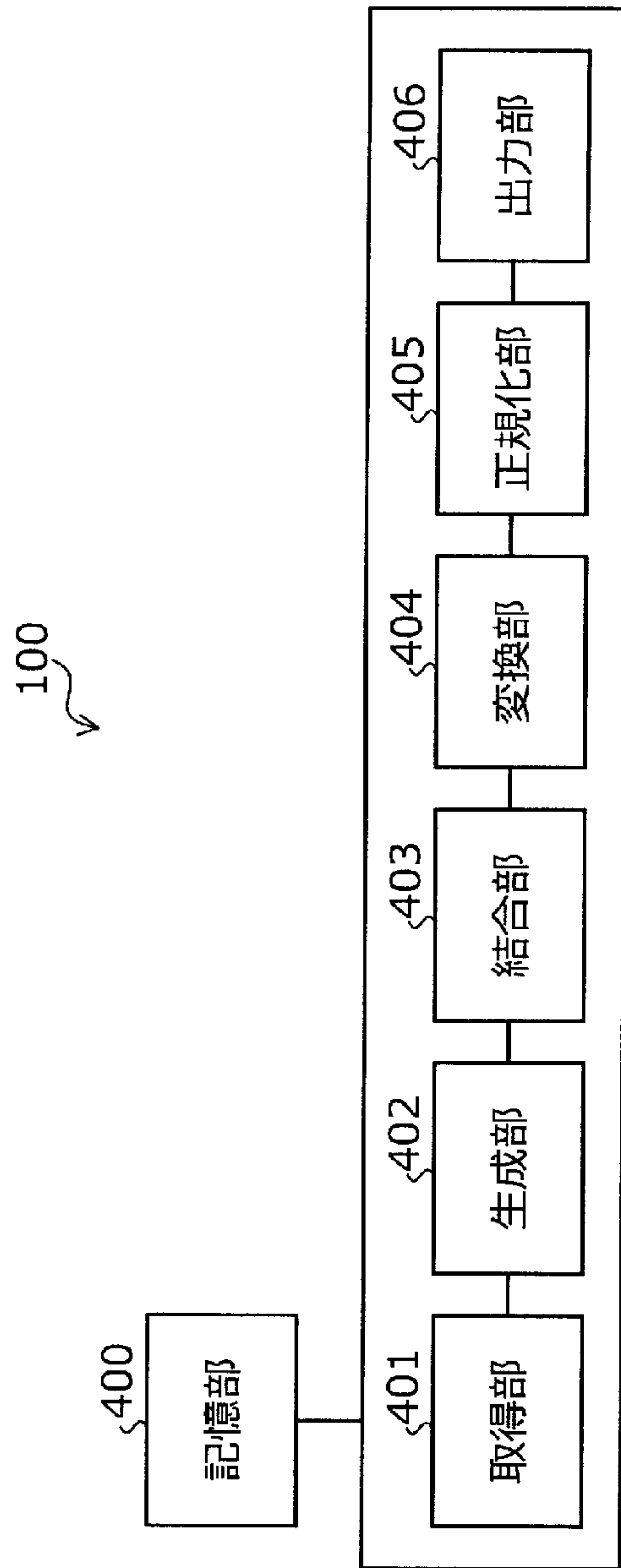
[図2]



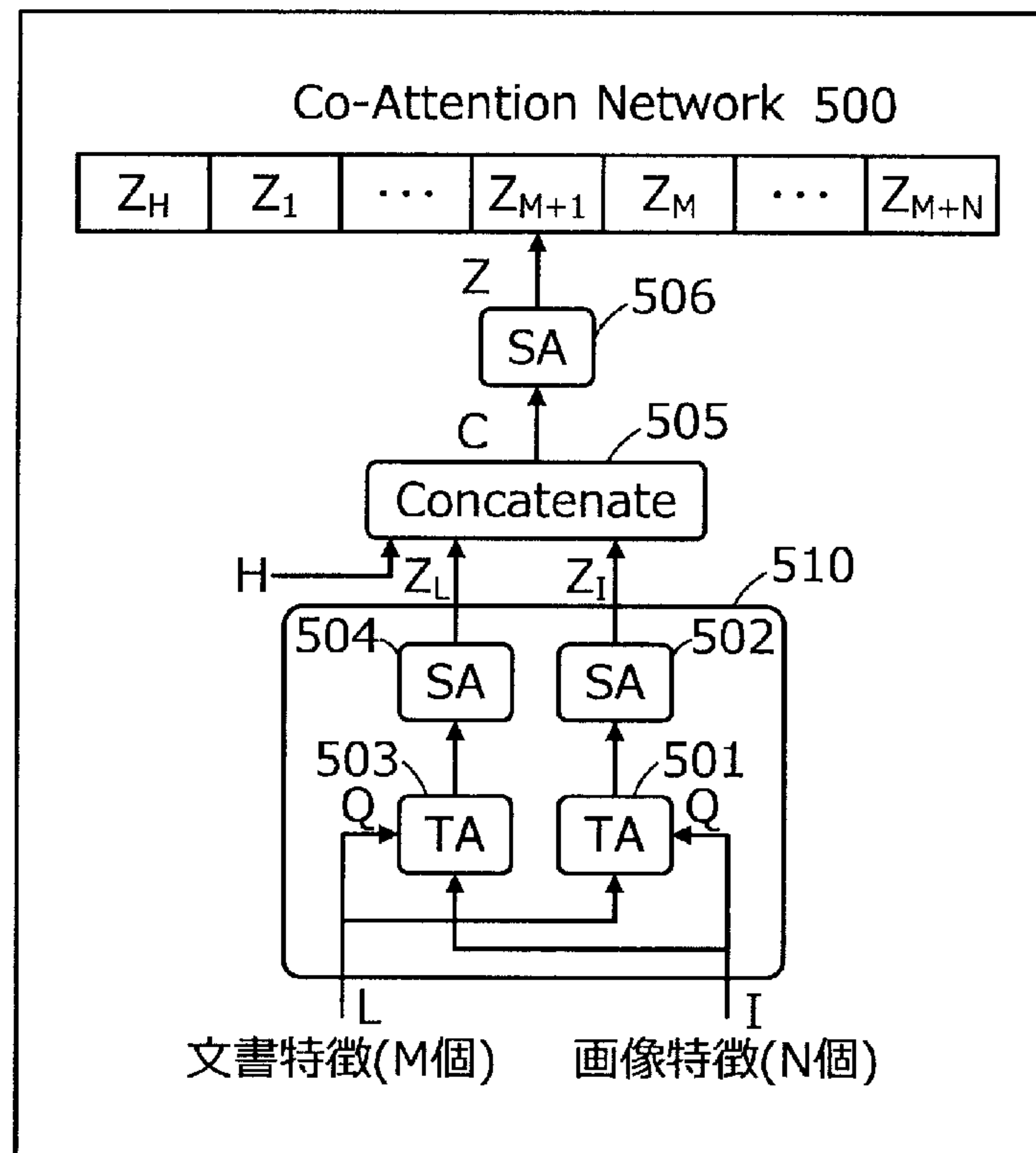
[図3]



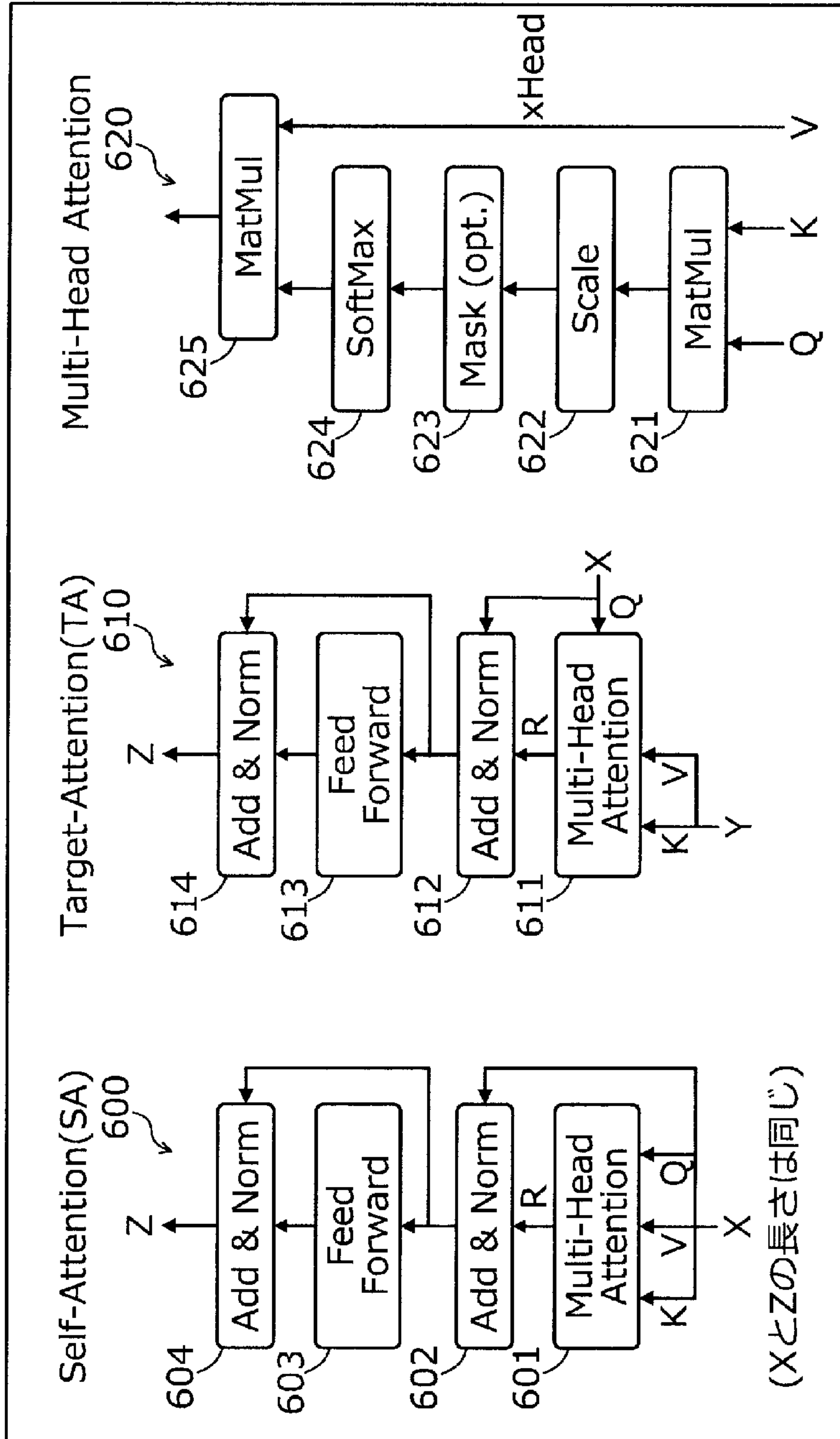
[図4]



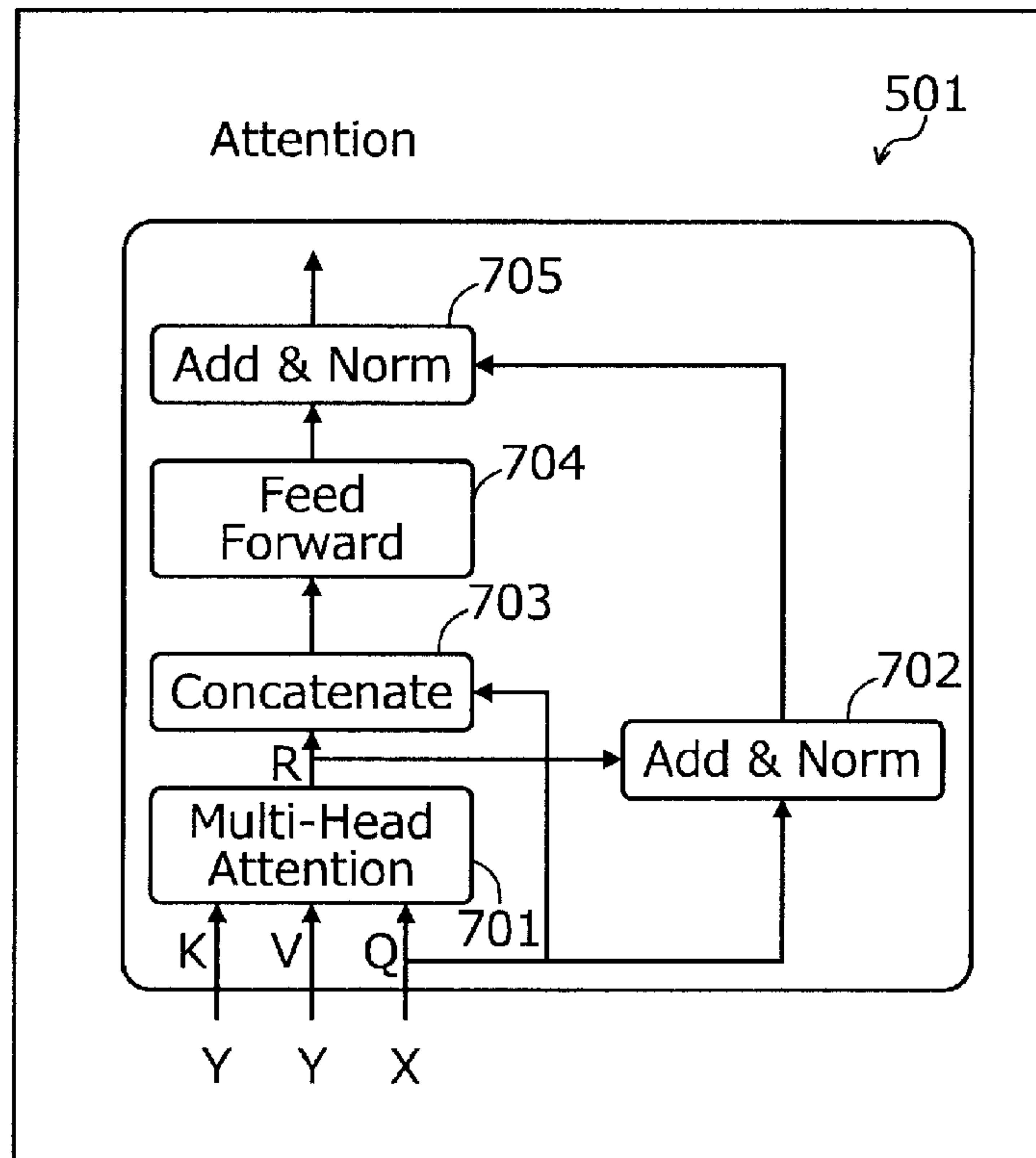
[図5]



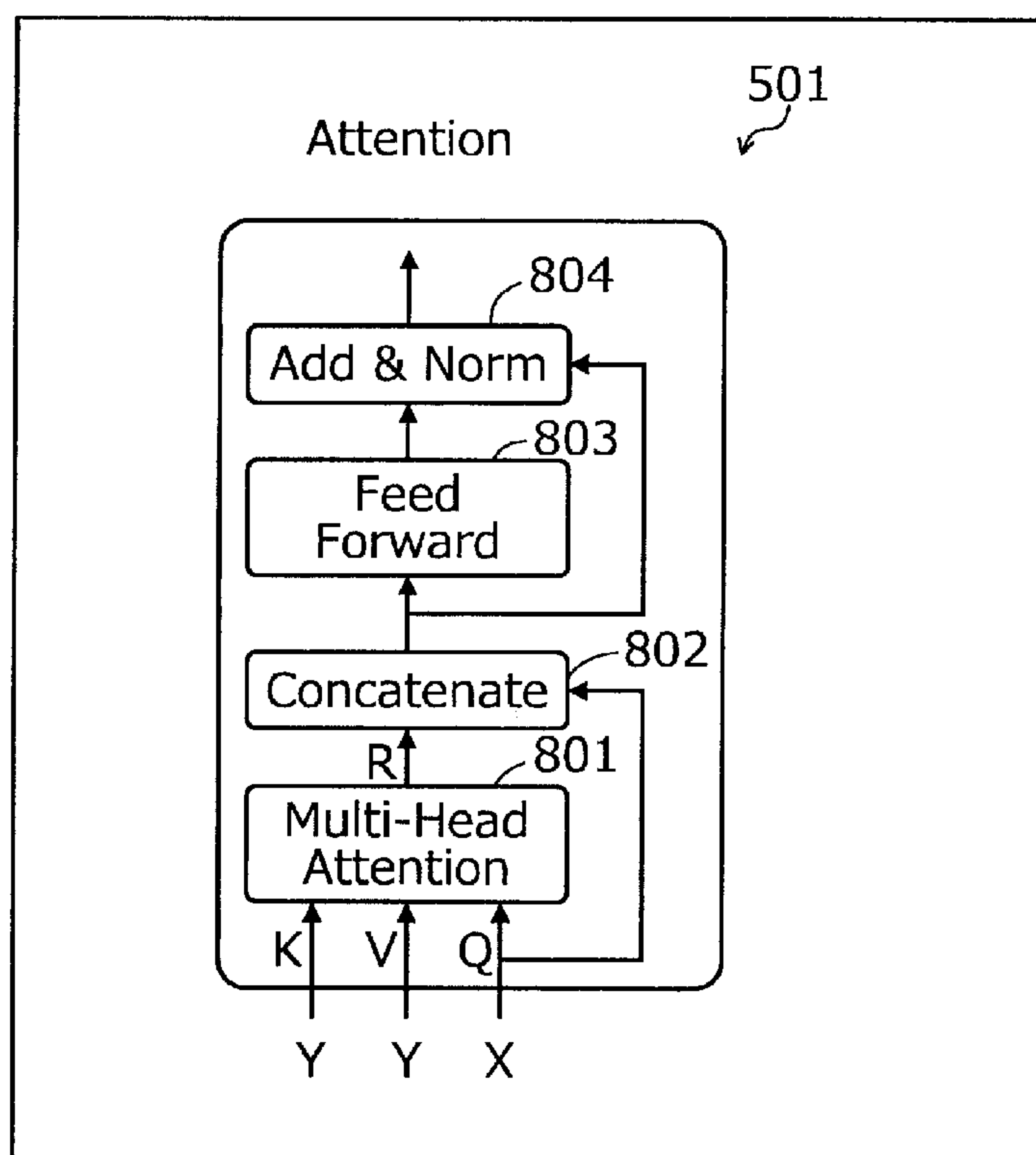
[図6]



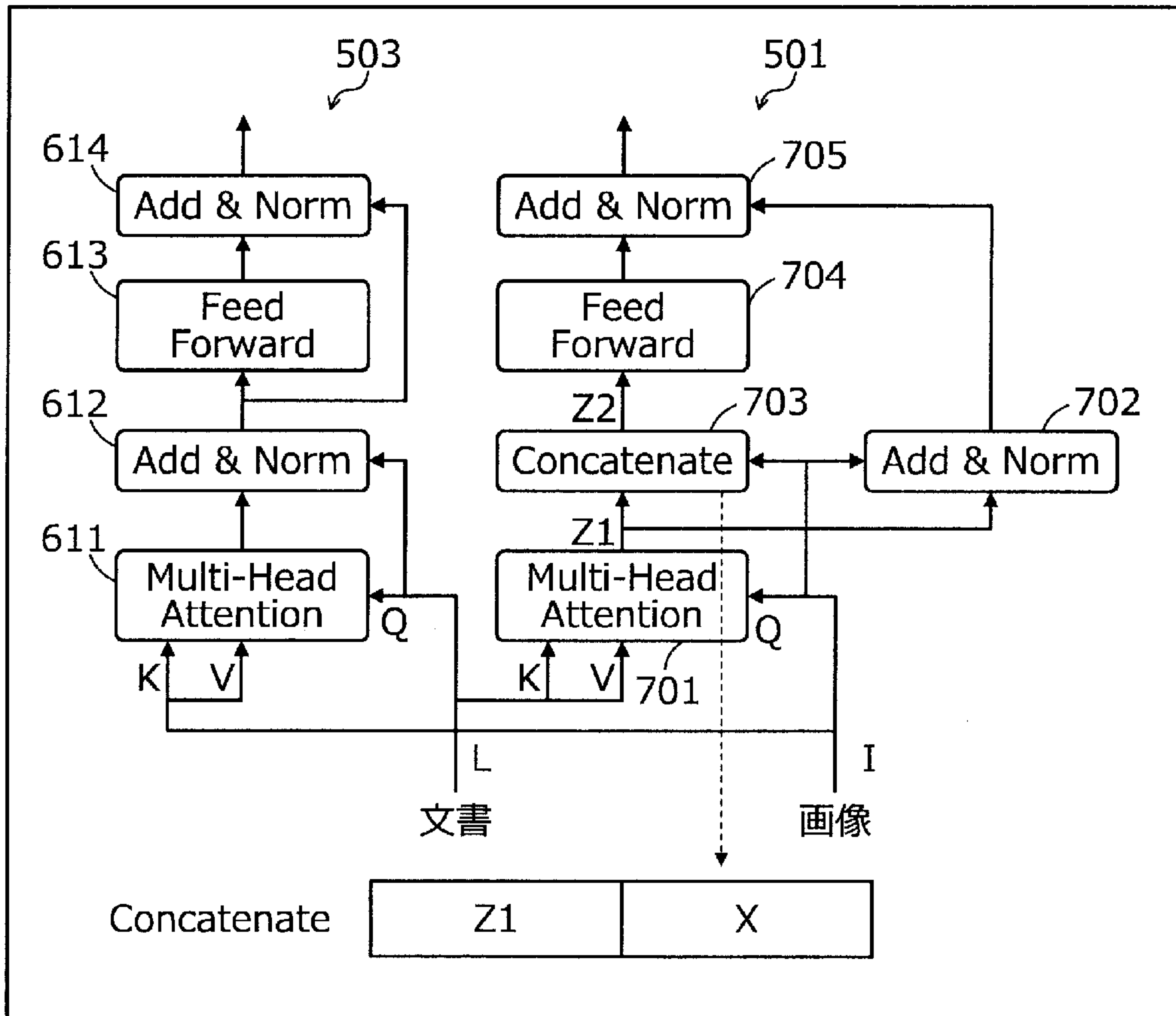
[図7]



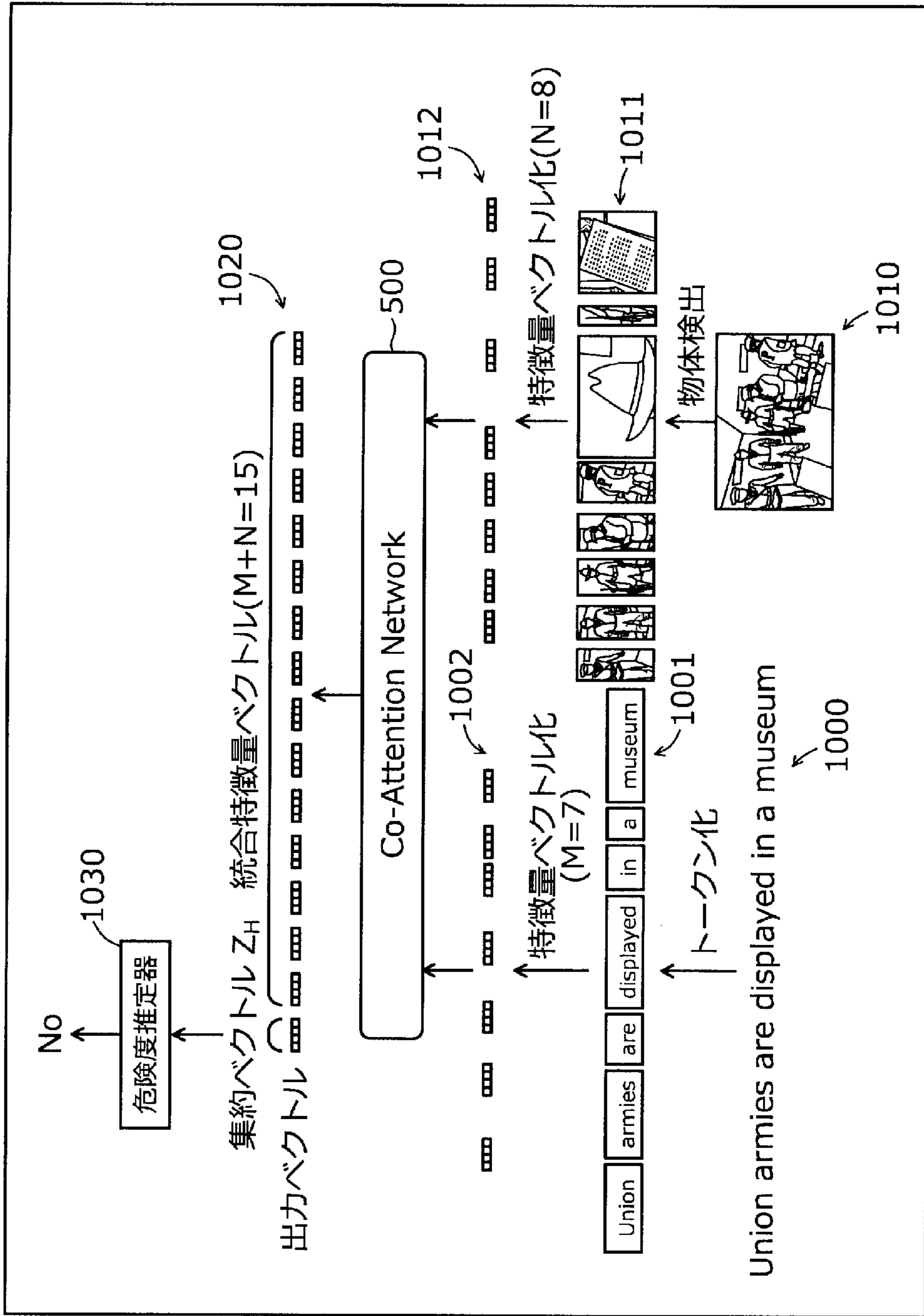
[図8]



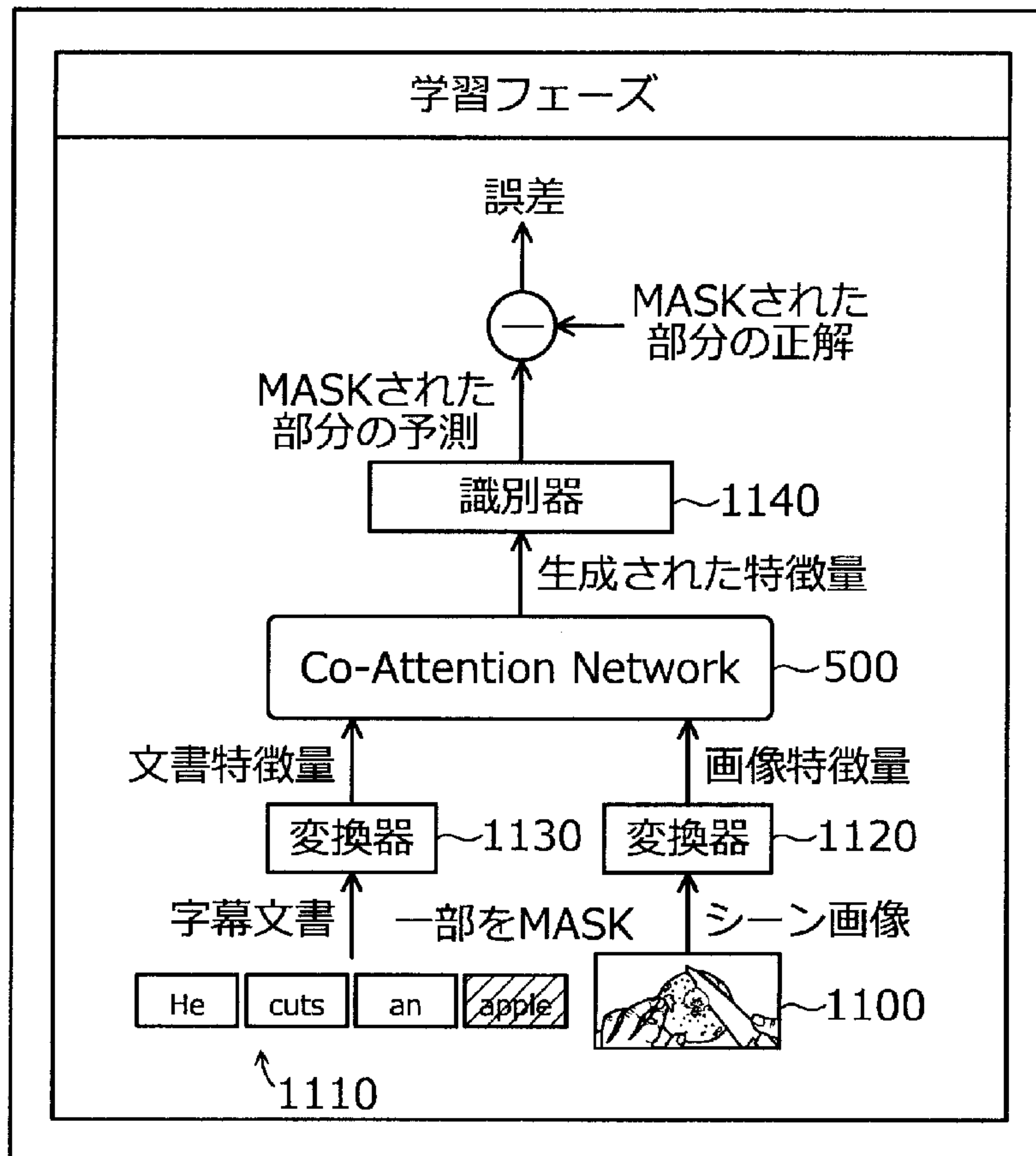
[図9]



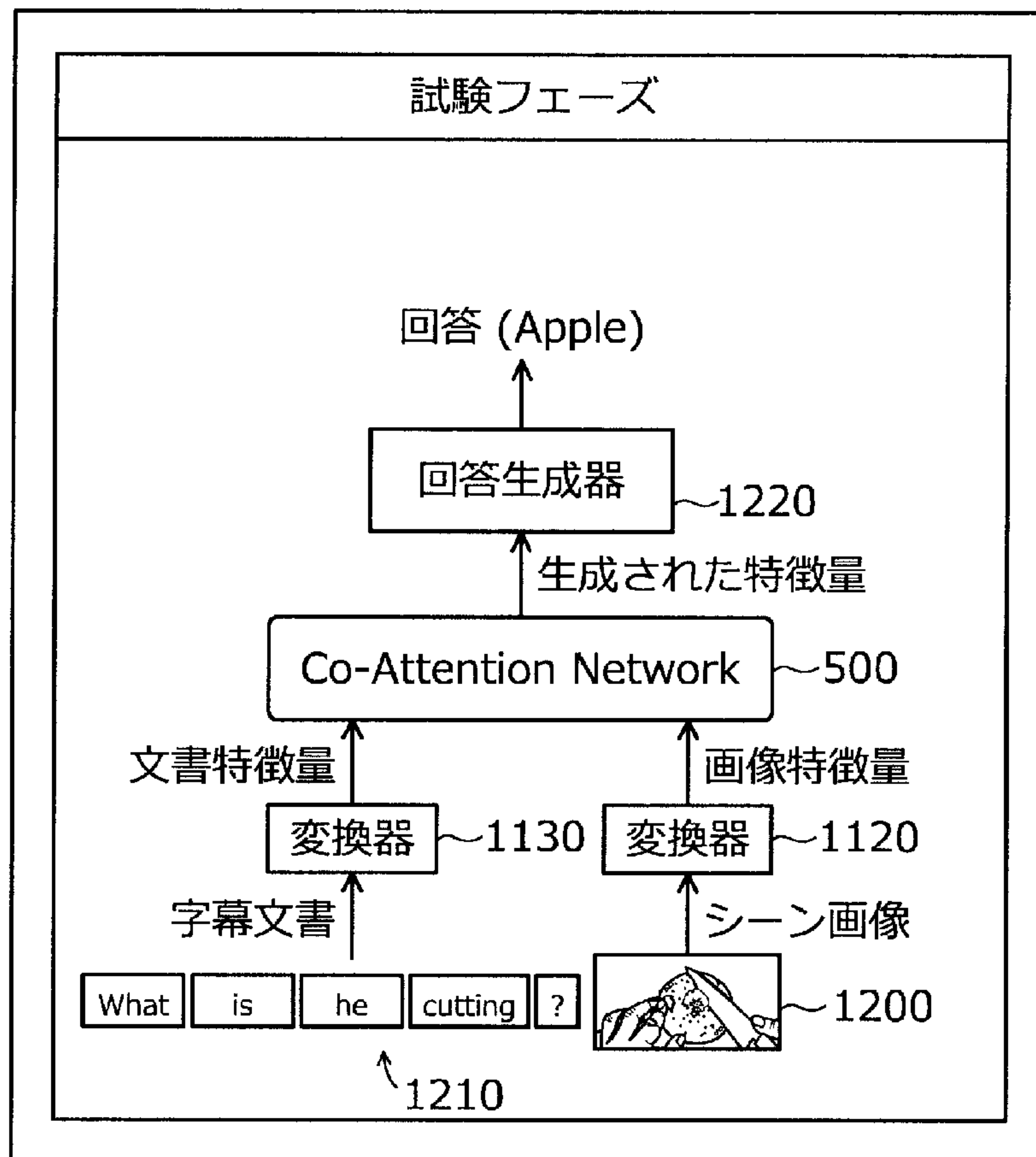
[図10]



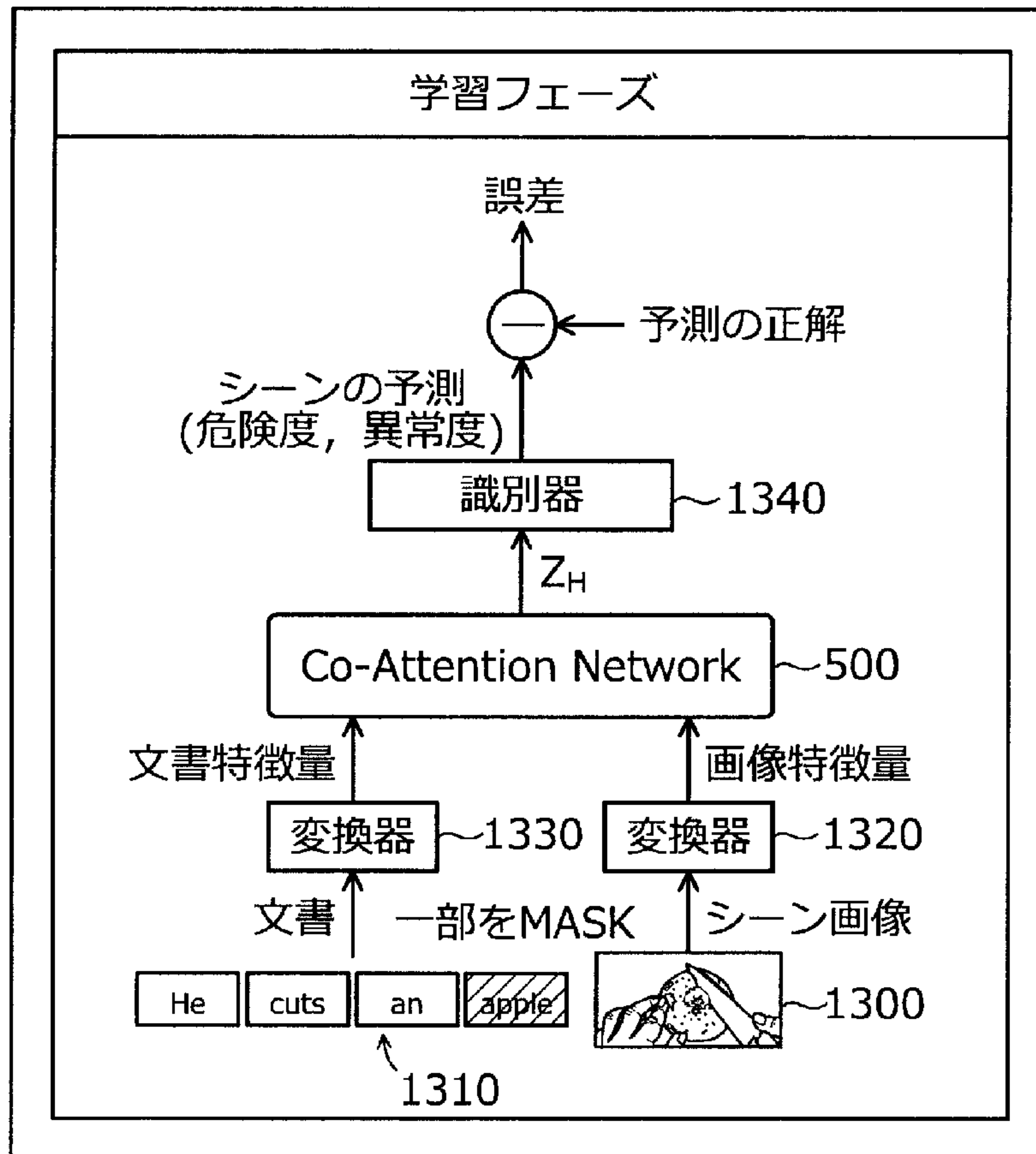
[図11]



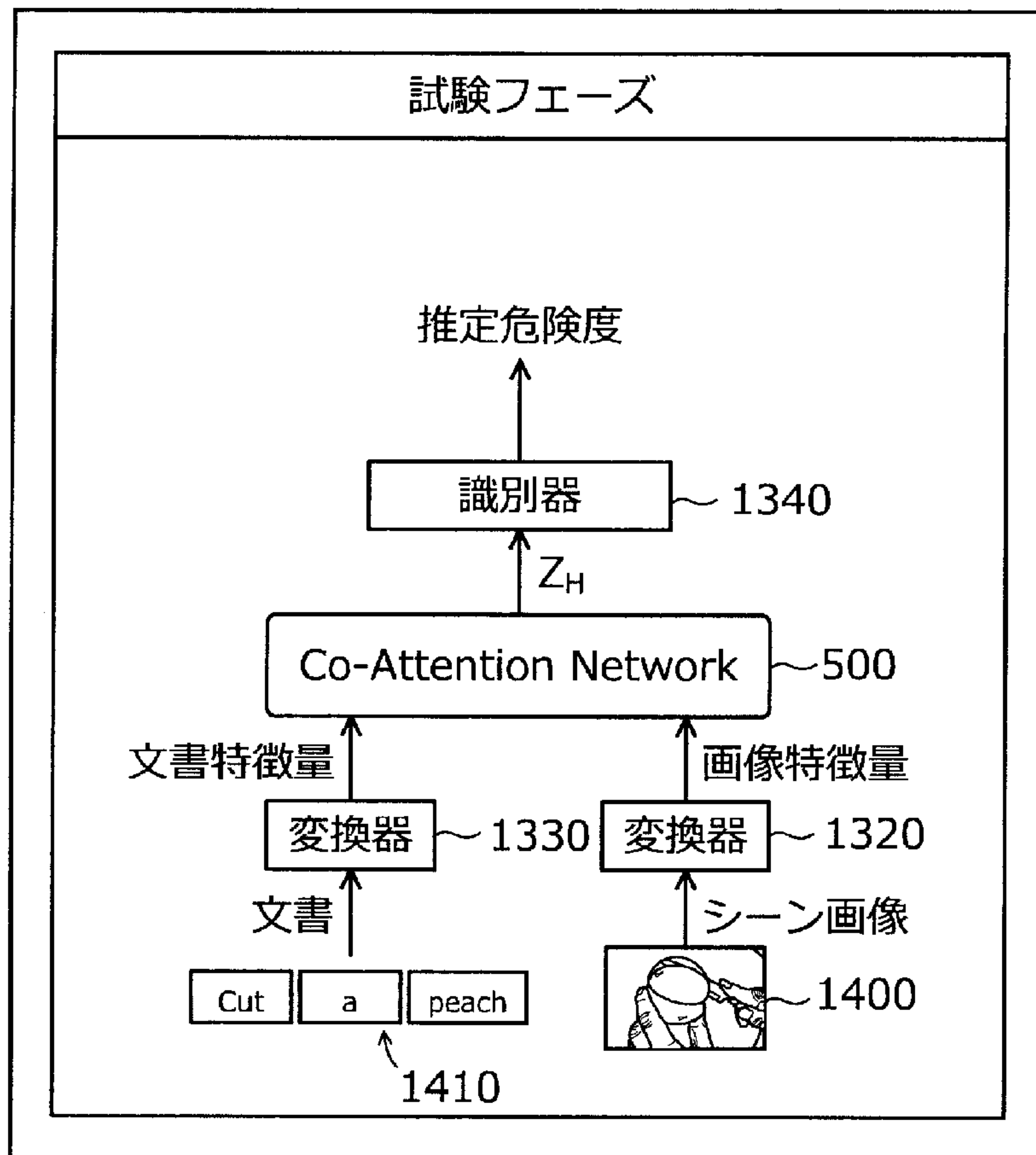
[図12]



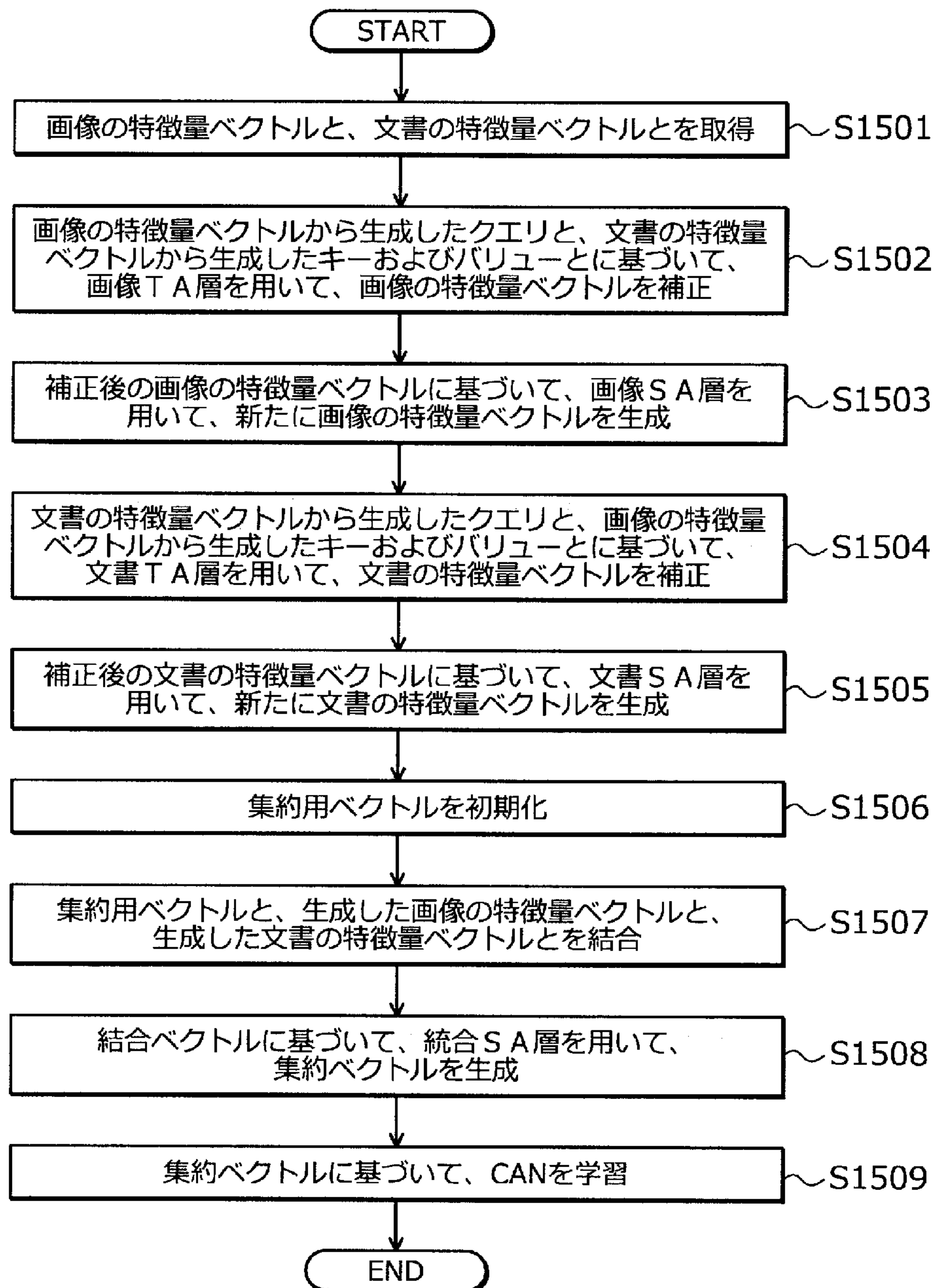
[図13]



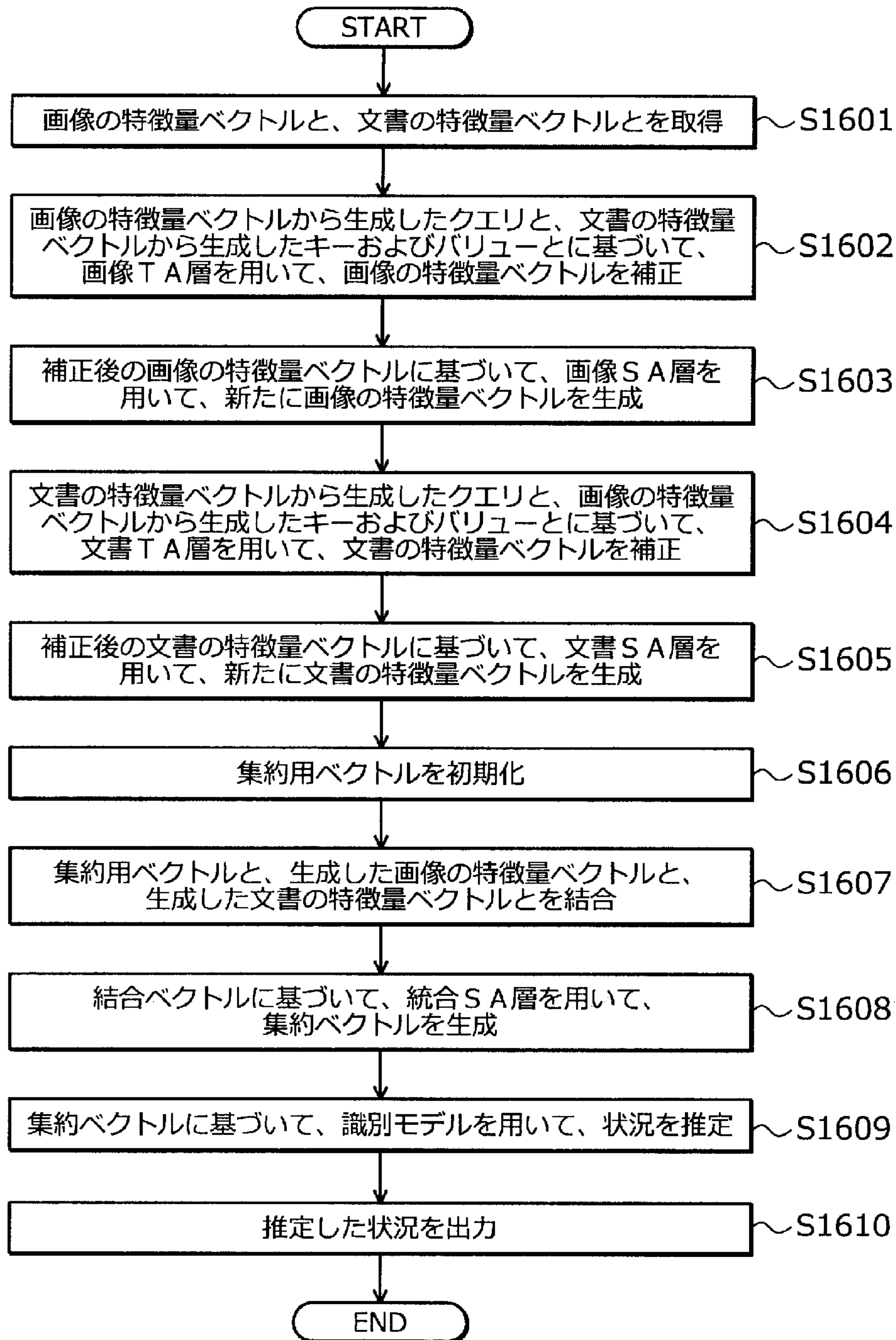
[図14]



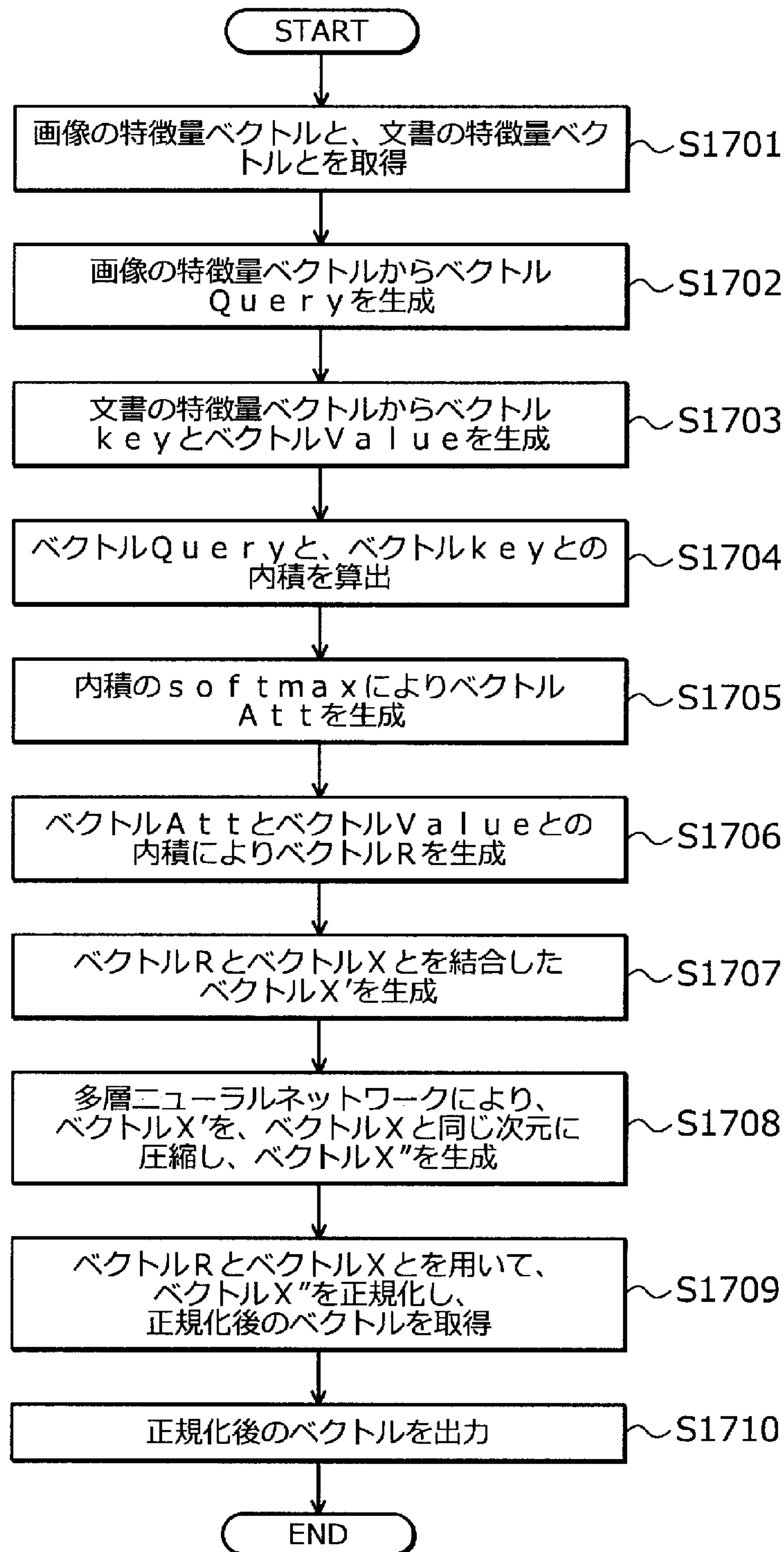
[図15]



[図16]



[図17]



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2019/044770

A. CLASSIFICATION OF SUBJECT MATTER

Int.Cl. G06N20/00 (2019.01) i

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)

Int.Cl. G06N3/00-99/00

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	1971-2019
Registered utility model specifications of Japan	1996-2019
Published registered utility model applications of Japan	1994-2019

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	LU, J. S. et al., ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks, arXiv.org [online], 06 August 2019, pp. 1-11, [retrieved on 13 December 2019], Internet: <URL:https://arxiv.org/pdf/1908.02265v1.pdf>, particularly, chapter 2	1-8

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents:	“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
“A” document defining the general state of the art which is not considered to be of particular relevance	“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
“E” earlier application or patent but published on or after the international filing date	“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
“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)	“&” document member of the same patent family
“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	

Date of the actual completion of the international search
18.12.2019

Date of mailing of the international search report
07.01.2020

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/JP2019/044770

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	NGUYEN, D. K. et al., Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering, [online], 2018, pp. 6087-6096, [retrieved on 13 December 2019], Internet: <URL:http://openaccess.thecvf.com/content_cvpr_2018/html/Nguyen_Improved_Fusion_of_CVPR_2018_paper.html>, particularly, chapters 1-3	1-8

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int.Cl. G06N20/00(2019.01) i

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))

Int.Cl. G06N3/00-99/00

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	1971-2019年
日本国実用新案登録公報	1996-2019年
日本国登録実用新案公報	1994-2019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	LU, Jiasen, et al., ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks, arXiv.org [online], 2019.08.06, pp.1-11, [検索日 2019.12.13], インターネット:<URL:https://arxiv.org/pdf/1908.02265v1.pdf>, 特に第2章	1-8

C欄の続きにも文献が列挙されている。

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	の日の後に公表された文献
「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	「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矛盾するものではなく、発明の原理又は理論の理解のために引用するもの
「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	「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
「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)	「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
「O」口頭による開示、使用、展示等に言及する文献	「&」同一パテントファミリー文献
「P」国際出願日前で、かつ優先権の主張の基礎となる出願	

国際調査を完了した日 18.12.2019	国際調査報告の発送日 07.01.2020
国際調査機関の名称及びあて先 日本国特許庁 (ISA/J P) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号	特許庁審査官 (権限のある職員) 多賀 実 電話番号 03-3581-1101 内線 3545

C (続き) . 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	NGUYEN, Duy-Kien, et al., Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering, [online], 2018, pp.6087-6096, [検索日 2019.12.13], インターネット : <URL : http://openaccess.thecvf.com/content_cvpr_2018/html/Nguyen_Improved_Fusion_of_CVPR_2018_paper.html >, 特に第1章-第3章	1-8