

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-139300
(P2019-139300A)

(43) 公開日 令和1年8月22日(2019.8.22)

(51) Int.Cl.
G06N 3/063 (2006.01)

F I
G06N 3/063

テーマコード (参考)

審査請求 未請求 請求項の数 9 O L (全 47 頁)

(21) 出願番号 特願2018-19251 (P2018-19251)
(22) 出願日 平成30年2月6日 (2018.2.6)

(出願人による申告) 平成27年度 国立研究開発法人 科学技術振興機構 戦略的創造研究推進事業 (ACCCEL)、 「TCI 積層情報処理アクセラレータの研究開発」 委託研究、 産業技術力強化法第19条の適用を受ける特許出願

(71) 出願人 504173471
国立大学法人北海道大学
北海道札幌市北区北8条西5丁目

(74) 代理人 110000958
特許業務法人 インテクト国際特許事務所

(74) 代理人 100120189
弁理士 奥 和幸

(74) 代理人 100140763
弁理士 平野 隆之

(72) 発明者 高前田 伸也
北海道札幌市北区北8条西5丁目 国立大学法人北海道大学内

(72) 発明者 植吉 晃大
北海道札幌市北区北8条西5丁目 国立大学法人北海道大学内

最終頁に続く

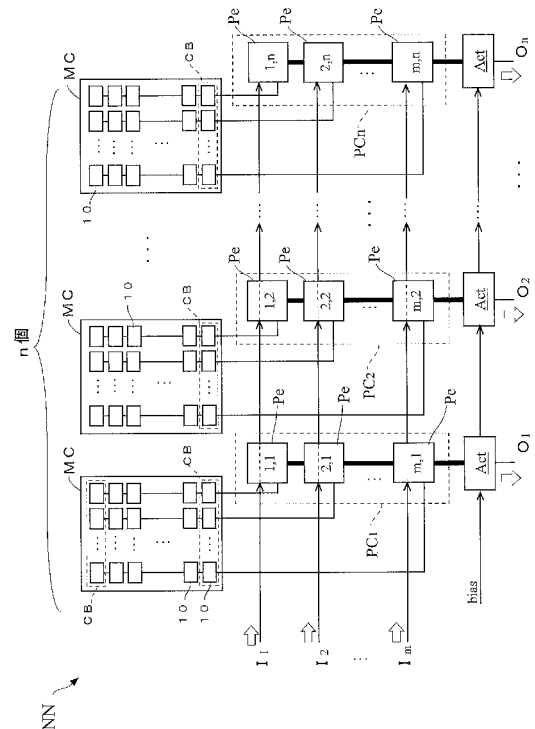
(54) 【発明の名称】 ニューラル電子回路

(57) 【要約】

【課題】 電子回路規模を縮小しつつ、 様々なタイプのニューラルネットワークを実現できるニューラル電子回路を提供する。

【解決手段】 並列で入力される1ビットの入力データであって、 並列の各入力データ (I_1 、 ...、 I_m) に応じて、「1」または「0」の1ビットの重み付け係数 (W_1 、 ...、 W_m) を記憶し、 当該重み付け係数を出力する記憶部 (MC) と、 並列の各入力データに設定された電子回路部であって、 記憶部 (MC) から出力された重み付け係数と入力データとを乗算する乗算機能を実現する第1電子回路部 (Pe) と、 並列の各入力データの第1電子回路部 (Pe) からの各乗算結果を加算し、 かつ、 当該加算結果に活性化関数を適用して1ビットの出力データを出力する加算・適用機能を実現する第2電子回路部 (Act) と、 を備える。

【選択図】 図3



【特許請求の範囲】**【請求項 1】**

並列で入力される 1 ビットの入力データであって、並列の各前記入力データに応じて、「1」または「0」の 1 ビットの重み付け係数を記憶し、当該重み付け係数を出力する記憶部と、

前記並列の各入力データに設定された電子回路部であって、前記記憶部から出力された重み付け係数と前記入力データとを乗算する乗算機能を実現する第 1 電子回路部と、

前記並列の各入力データの前記第 1 電子回路部からの各乗算結果を加算し、かつ、当該加算結果に活性化関数を適用して 1 ビットの出力データを出力する加算・適用機能を実現する第 2 電子回路部と、

を備え、

前記第 1 電子回路部が、前記入力データの値と前記記憶部から出力された値とが一致する場合に前記入力データの入力に対応して「1」を出力し、前記入力データの値と前記記憶部から出力された値とが異なる場合に前記入力データの入力に対応して「0」を出力することを特徴とするニューラル電子回路。

【請求項 2】

請求項 1 に記載のニューラル電子回路において、

前記記憶部および前記第 2 電子回路部が、並列で出力される各前記出力データに応じて設定されたことを特徴とするニューラル電子回路。

【請求項 3】

請求項 1 または請求項 2 に記載のニューラル電子回路において、

各前記第 1 電子回路部からの前記乗算結果を一時記憶する一時記憶部を前記第 1 電子回路部毎に更に備え、

前記各一時記憶部が、直列に設定され、前記乗算結果を前記第 2 電子回路部へ順次転送することを特徴とするニューラル電子回路。

【請求項 4】

請求項 1 から請求項 3 のいずれか 1 項に記載のニューラル電子回路において、

前記第 2 回路部は、前記並列の各入力データに設定された複数の前記第 1 電子回路部において、前記入力データが入力されるサイクル単位で、前記第 1 電子回路部が「1」を算出した回数から、前記第 1 電子回路部が「0」を算出した回数を減じた値が予め定められた閾値以上の場合に「1」を前記出力データとして出力し、前記減じた値が前記閾値未満の場合に「0」を前記出力データとして出力することを特徴とするニューラル電子回路。

【請求項 5】

請求項 1 から請求項 4 のいずれか 1 項に記載のニューラル電子回路において、

前記記憶部は、「1」または「0」の重み付け係数、および、ニューロン間の接続の有無を示す所定値の組を複数記憶し、当該重み付け係数および当該所定値を出力し、

前記第 1 電子回路部が、前記入力データの値と前記記憶部から出力された値とが一致する場合に前記入力データの入力に対応して「1」を出力し、前記入力データの値と前記記憶部から出力された値とが異なる場合に前記入力データの入力に対応して「0」を出力し、前記所定値が前記記憶部から出力された場合に前記入力データの入力に対応して当該所定値を出力することを特徴とするニューラル電子回路。

【請求項 6】

請求項 1 から請求項 5 のいずれか 1 項に記載のニューラル電子回路において、

前記記憶部が、前記第 1 電子回路部に順次入力される前記入力データに対応した重み付け係数を、前記第 1 電子回路部に順次出力することを特徴とするニューラル電子回路。

【請求項 7】

請求項 6 に記載のニューラル電子回路において、

前記第 1 電子回路部が、並列で入力される前記入力データの入力並列数分、前記乗算結果を加算した部分加算結果を出力し、

前記第 2 電子回路部が、前記部分加算結果から前記加算結果を算出することを特徴とす

10

20

30

40

50

るニューラル電子回路。

【請求項 8】

請求項 1 から請求項 5 のいずれか 1 項に記載のニューラル電子回路において、前記記憶部が、並列で入力される並列の各前記入力データに対応した重み付け係数を、各前記第 1 電子回路部へ出力することを特徴とするニューラル電子回路。

【請求項 9】

請求項 8 に記載のニューラル電子回路において、前記入力データを並列で一度に入力可能な入力可能並列数より、前記入力データの入力並列数が大きい場合、

前記第 1 電子回路部は、前記入力可能並列数の並列で前記入力データを受け入れた後、前記入力可能並列数の並列で受け入れできなかった残りの前記入力データを受け入れ、

前記記憶部は、前記残りの入力データに対応する前記重み付け係数を出力することを特徴とするニューラル電子回路。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ニューラルネットワークを電子回路で実現するニューラル電子回路の技術分野に属する。

【背景技術】

【0002】

近年、人の脳機能をモデル化した、いわゆるニューラルネットワーク回路についての研究開発が行われている。このとき、従来のニューラルネットワーク回路としては、例えば浮動小数点又は固定小数点を使った積和演算を用いて実現される場合が多く、この場合には、例えば演算コストが大きく、処理負荷が高いという問題点があった。

20

【0003】

そこで近年、上記入力データ及び上記重み付け係数をそれぞれ 1 ビットとする、いわゆる「バイナリニューラルネットワーク回路」のアルゴリズムが提案されている。ここで、上記バイナリニューラルネットワーク回路のアルゴリズムを示す先行技術文献としては、例えば下記非特許文献 1 及び非特許文献 2 が挙げられる。

【先行技術文献】

30

【非特許文献】

【0004】

【非特許文献 1】「XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks」論文、Mohammad Rastegari 他、arXiv:1603.05279v2 [cs.CV、2016年 4 月 19 日 (URL : <http://arxiv.org/abs/1603.05279>)

【非特許文献 2】「Binarized Neural Networks: Training Neural Networks with Weights and Activations Constrained to +1 or -1」論文、Matthieu Courbariaux 他、arXiv:1602.02830v3 [cs.LG]、2016年 3 月 17 日 (URL : <http://arxiv.org/abs/1602.02830>)

【発明の概要】

【発明が解決しようとする課題】

40

【0005】

しかしながら、上述したいずれの非特許文献においても、当該論文記載の理論を具体的にどのように実現するかについては、全く記載されていない。また、各論文記載の理論により単位演算コストが大幅に下がることを利用して並列演算を可能にしたいが、そのためのハードウェア構成も未知である。更に、複雑多岐なニューラルネットを構築するためにも、フル結合のニューラルネットワークの他にも、コンボリューション演算等も効率的に計算できる電子回路を実現する必要があった。

【0006】

そこで本発明は、上記の各問題点及び要請等に鑑みて為されたもので、その課題の一例は、上記バイナリニューラルネットワーク回路のアルゴリズムを用いて、電子回路規模を

50

縮小しつつ、様々なタイプのニューラルネットワークを実現できるニューラル電子回路を提供することにある。

【課題を解決するための手段】

【0007】

上記の課題を解決するために、請求項1に記載の発明は、並列で入力される1ビットの入力データであって、並列の各前記入力データに応じて、「1」または「0」の1ビットの重み付け係数を記憶し、当該重み付け係数を出力する記憶部と、前記並列の各入力データに設定された電子回路部であって、前記記憶部から出力された重み付け係数と前記入力データとを乗算する乗算機能を実現する第1電子回路部と、前記並列の各入力データの前記第1電子回路部からの各乗算結果を加算し、かつ、当該加算結果に活性化関数を適用して1ビットの出力データを出力する加算・適用機能を実現する第2電子回路部と、を備え、前記第1電子回路部が、前記入力データの値と前記記憶部から出力された値とが一致する場合に前記入力データの入力に対応して「1」を出力し、前記入力データの値と前記記憶部から出力された値とが異なる場合に前記入力データの入力に対応して「0」を出力することを特徴とする。

10

【0008】

請求項2に記載の発明は、請求項1に記載のニューラル電子回路において、前記記憶部および前記第2電子回路部が、並列で出力される各前記出力データに応じて設定されたことを特徴とする。

【0009】

20

請求項3に記載の発明は、請求項1または請求項2に記載のニューラル電子回路において、各前記第1電子回路部からの前記乗算結果を一時記憶する一時記憶部を前記第1電子回路部毎に更に備え、前記各一時記憶部が、直列に設定され、前記乗算結果を前記第2電子回路部へ順次転送することを特徴とする。

【0010】

請求項4に記載の発明は、請求項1から請求項3のいずれか1項に記載のニューラル電子回路において、前記第2回路部は、前記並列の各入力データに設定された複数の前記第1電子回路部において、前記入力データが入力されるサイクル単位で、前記第1電子回路部が「1」を算出した回数から、前記第1電子回路部が「0」を算出した回数を減じた値が予め定められた閾値以上の場合に「1」を前記出力データとして出力し、前記減じた値が前記閾値未満の場合に「0」を前記出力データとして出力することを特徴とする。

30

【0011】

請求項5に記載の発明は、請求項1から請求項4のいずれか1項に記載のニューラル電子回路において、前記記憶部は、「1」または「0」の重み付け係数、および、ニューロン間の接続の有無を示す所定値の組を複数記憶し、当該重み付け係数および当該所定値を出力し、前記第1電子回路部が、前記入力データの値と前記記憶部から出力された値とが一致する場合に前記入力データの入力に対応して「1」を出力し、前記入力データの値と前記記憶部から出力された値とが異なる場合に前記入力データの入力に対応して「0」を出力し、前記所定値が前記記憶部から出力された場合に前記入力データの入力に対応して当該所定値を出力することを特徴とする。

40

【0012】

請求項6に記載の発明は、請求項1から請求項5のいずれか1項に記載のニューラル電子回路において、前記記憶部が、前記第1電子回路部に順次入力される前記入力データに対応した重み付け係数を、前記第1電子回路部に順次出力することを特徴とする。

【0013】

請求項7に記載の発明は、請求項6に記載のニューラル電子回路において、前記第1電子回路部が、並列で入力される前記入力データの入力並列数分、前記乗算結果を加算した部分加算結果を出力し、前記第2電子回路部が、前記部分加算結果から前記加算結果を算出することを特徴とする。

【0014】

50

請求項 8 に記載の発明は、請求項 1 から請求項 5 のいずれか 1 項に記載のニューラル電子回路において、前記記憶部が、並列で入力される並列の各前記入力データに対応した重み付け係数を、各前記第 1 電子回路部に出力することを特徴とする。

【 0 0 1 5 】

請求項 9 に記載の発明は、請求項 8 に記載のニューラル電子回路において、前記入力データを並列で一度に入力可能な入力可能並列数より、前記入力データの入力並列数が大きい場合、前記第 1 電子回路部は、前記入力可能並列数の並列で前記入力データを受け入れた後、前記入力可能並列数の並列で受け入れできなかった残りの前記入力データを受け入れ、前記記憶部は、前記残りの入力データに対応する前記重み付け係数を出力することを特徴とする。

10

【発明の効果】

【 0 0 1 6 】

本発明によれば、1 ビットの入力データの値と記憶部から出力された値とが一致する場合に入力データの入力に対応して「1」を出力し、入力データの値と記憶部から出力された値とが異なる場合に入力データの入力に対応して「0」を出力することにより、乗算機能を実現するため、回路規模を縮小でき、並列で入力される並列の各入力データに応じて配置された第 1 電子回路部により、様々なタイプのニューラルネットワークを実現できる。

【図面の簡単な説明】

【 0 0 1 7 】

20

【図 1】実施形態に係るニューラルネットワークを説明する図であり、(a) は一つのニューロンをモデル化したユニットを示す図であり、(b) は複数のユニットの結合からなるニューラルネットワークの状態を示す図である。

【図 2】実施形態に係るニューラルネットワーク・システムの概要構成例を示すブロック図である。

【図 3】図 2 のニューラル電子回路の一例を示すブロック図である。

【図 4】図 3 のプロセスエレメント・コラムの一例を示すブロック図である。

【図 5】図 4 のプロセスエレメントの動作のタイミングの一例を示す模式図である。

【図 6 A】図 2 のプロセスエレメントのデジタル回路の一例を示す回路図である。

【図 6 B】図 2 の加算活性化部のデジタル回路の一例を示す回路図である。

30

【図 7】コンボリユーション演算におけるデータの関係の一例を示す模式図である。

【図 8】図 7 のコンボリユーション演算を実現するニューラル電子回路の一例を示すブロック図である。

【図 9】フル結合のニューラルネットワークの一例を示す模式図である。

【図 10】図 9 のフル結合のニューラルネットワークを実現するニューラル電子回路の一例を示すブロック図である。

【図 11】ニューラルネットワークの層内拡張の一例を示す模式図である。

【図 12】図 11 の層内拡張を実現するコア電子回路の接続の一例を示すブロック図である。

【図 13】ニューラルネットワークの層数拡張の一例を示す模式図である。

40

【図 14】図 13 の層数拡張を実現するコア電子回路の接続の一例を示すブロック図である。

【図 15】実施形態に係るニューラルネットワーク回路を示す図であり、(a) は当該ニューラルネットワーク回路に相当するニューラルネットワークを示す図であり、(b) は当該ニューラルネットワーク回路の構成を示すブロック図であり、(c) は当該ニューラルネットワーク回路に対応する真理値表である。

【図 16】実施形態に係るニューラルネットワーク回路の細部構成を示す図であり、(a) は当該細部構成に係るメモリセルの回路の一例を示す図であり、(b) は当該細部構成の回路の一例を示す図である。

【図 17】実施形態に係るニューラルネットワーク集積回路の第 1 例を示す図であり、(

50

a) は当該第 1 例に相当するニューラルネットワークを示す図であり、(b) は当該第 1 例の構成を示すブロック図である。

【図 18】実施形態に係るニューラルネットワーク集積回路の第 2 例を示す図であり、(a) は当該第 2 例に相当するニューラルネットワークを示す図であり、(b) は当該第 2 例の構成を示すブロック図である。

【図 19】実施形態に係るニューラルネットワーク集積回路の第 3 例を示す図であり、(a) は当該第 3 例に相当するニューラルネットワークを示す図であり、(b) は当該第 3 例の構成を示すブロック図である。

【図 20】実施形態に係るニューラルネットワーク集積回路の第 4 例を示す図であり、(a) は当該第 4 例に相当するニューラルネットワークを示す図であり、(b) は当該第 4 例の構成を示すブロック図であり、(c) は当該第 4 例に係るスイッチボックスの構成の一例を示すブロック図である。

10

【図 21】関連形態に係るニューラルネットワーク集積回路の第 1 例の一部を示す図であり、(a) は当該一部に相当するニューラルネットワークを示す図であり、(b) は当該一部の構成を示すブロック図であり、(c) の当該一部に対応する真理値表である。

【図 22】関連形態に係るニューラルネットワーク集積回路の第 1 例を示す図であり、(a) は当該第 1 例に相当するニューラルネットワークを示す図であり、(b) は当該第 1 例の構成を示すブロック図である。

【図 23】関連形態に係るニューラルネットワーク回路の第 1 例を示す図であり、(a) は当該第 1 例に相当するニューラルネットワークを示す図であり、(b) は当該第 1 例の構成を示すブロック図である。

20

【図 24】関連形態に係るニューラルネットワーク集積回路の第 2 例を示す図であり、(a) は当該第 2 例に相当するニューラルネットワークを示す図であり、(b) は当該第 2 例の構成を示すブロック図である。

【図 25】関連形態に係るニューラルネットワーク集積回路の第 3 例を示す図であり、(a) は当該第 3 例に相当するニューラルネットワークを示す図であり、(b) は当該第 3 例の構成を示すブロック図である。

【図 26】関連形態に係るニューラルネットワーク集積回路の第 4 例を示す図であり、(a) は当該第 4 例の構成を示すブロック図であり、(b) は当該第 4 例に相当する回路例を示す図である。

30

【図 27】関連形態に係るニューラルネットワーク集積回路の第 4 例の細部構成を示す図であり、(a) は当該第 4 例に係るパイプラインレジスタ等の回路の一例を示す図であり、(b) は当該第 4 例に係る多数判定入力回路及び直列多数判定回路それぞれの一例を示す図であり、(c) は当該第 4 例に係る並列多数判定回路の一例を示す図であり、(d) は当該第 4 例における動作を示すタイミングチャートである。

【発明を実施するための形態】

【0018】

次に、本発明に係る実施形態及び関連形態について、図面に基づいてそれぞれ説明する。なお以下に説明する実施形態等は、人の脳機能をモデル化したニューラルネットワークを電子的な回路で実現したニューラルネットワーク回路について本発明を適用した場合の実施形態等である。

40

【0019】

[1. ニューラルネットワークについて]

【0020】

まず、上記脳機能をモデル化したニューラルネットワークについて、一般的に図 1 を用いて説明する。

【0021】

一般に人の脳の中には、多数のニューロン（神経細胞）が存在しているとされている。脳の中で各ニューロンは、多数の他のニューロンからの電気信号を受信し、また更に他の多数のニューロンへ電気信号を送信している。そして脳は、各ニューロン間のこれら電気

50

信号の伝達によって、様々な情報処理を行っていると考えられている。このとき、各ニューロン間における電気信号の送受信は、シナプスと呼ばれる細胞を介して行われる。そして、脳における上記ニューロン間の電気信号の送受信をモデル化してコンピュータ内に脳機能を実現しようとしたものが、ニューラルネットワークである。

【0022】

より具体的にニューラルネットワークでは図1(a)に例示するように、外部から入力される複数の入力データ I_1 、入力データ I_2 、...、入力データ I_n (n は自然数。以下同様。)のそれぞれに対する乗算処理、各乗算結果とバイアス(閾値)の加算処理及び活性化関数の適用処理を一つのニューロン NR で実行し、その結果を出力データ O とすることで、脳機能における一つのニューロンに対する上記電気信号の送受信をモデル化する。なお以下の説明において、上記活性化関数の適用処理を、単に「活性化処理」と称する。このとき一つのニューロン NR では、複数の入力データ I_1 、入力データ I_2 、...、入力データ I_n それぞれに対応して予め設定された(つまり既定の)重み付け係数 W_1 、重み付け係数 W_2 、...、重み付け係数 W_n を当該入力データ I_1 、入力データ I_2 、...、入力データ I_n に対してそれぞれ乗算することで上記乗算処理が実行される。

10

【0023】

その後当該ニューロン NR は、各入力データ I_1 、入力データ I_2 、...、入力データ I_n に対する上記乗算処理の結果のそれぞれを加算してバイアス $bias$ の値を加算する上記加算処理を実行する。そして当該ニューロン NR は次に、上記加算処理の結果に既定の活性化関数 F を適用する上記活性化処理を実行し、その結果を上記出力データとして他の一又は複数のニューロン NR に出力する。上述した一連の乗算処理、加算処理及び活性化処理を数式で表すと、図1(a)に示す式(1)となる。このとき、重み付け係数 W_1 、重み付け係数 W_2 、...、重み付け係数 W_n を入力データ I_1 、入力データ I_2 、...、入力データ I_n にそれぞれ乗算する乗算処理が、ニューロン NR 間の上記電気信号のやり取りにおける上記シナプスの作用に相当すると共に、本発明に係る「乗算機能」の一例に相当する。また、上記加算処理及び活性化処理が本発明に係る「加算/適用機能」の一例に相当する。そして図1(b)に例示されるように、図1(a)に例示する一つのニューロン NR が多数集まってシナプスにより互いに接続されることにより、脳全体がニューラルネットワーク SS としてモデル化される。

20

【0024】

[2. ニューラルネットワーク・システムの構成および機能概要]

(2.1 ニューラルネットワーク・システムの構成および機能)

次に、本発明の一実施形態に係るニューラルネットワーク・システムの構成および概要機能について、図2を用いて説明する。

30

【0025】

図2は、本実施形態に係るニューラルネットワーク・システム NNS の概要構成例を示す模式図である。

【0026】

図2に示すように、ニューラルネットワーク・システム NNS は、様々なタイプのニューラルネットワークを電子回路で実現可能な複数のコア電子回路 $Core$ と、コア電子回路 $Core$ 同士を接続するシステムバス bus と、を備えている。

40

【0027】

コア電子回路 $Core$ は、様々なタイプのニューラルネットワークを電子回路で実現可能なニューラル電子回路 NN と、ニューラル電子回路 NN の重み付け係数等を設定するメモリアクセス制御部 $MCnt$ と、ニューラル電子回路 NN およびメモリアクセス制御部 $MCnt$ を制御する制御部 Cnt と、を有する。ここで、様々なタイプのニューラルネットワークの一例として、ニューロン層間のニューロン同士がフル結合したフル結合のタイプのニューラルネットワーク、コンボリューション演算をするニューラルネットワーク、ニューロン層における層内拡張をしたニューラルネットワーク、層数を拡張するニューラルネットワーク等が挙げられる。

50

【0028】

ニューラル電子回路NNは、入力データ I_1 、 \dots 、 I_m (m は自然数。以下同様)を並列で順次供給する入力メモリアレイ部MAiと、重み付け係数のデータを並列で順次供給するメモリセルレイ部MC (記憶部の一例)と、供給された入力データ I_1 、 \dots 、 I_m と重み付け係数とを乗算する乗算機能を実現して乗算結果を出力する複数のプロセスエレメント部Pe (第1電子回路部の一例)と、並列の各入力データのプロセスエレメント部Peからの各乗算結果を加算して加算結果に活性化関数を適用する加算活性化部Act (第2電子回路部の一例)と、各加算活性化部Actからの出力データ O_1 、 \dots 、 O_n (n は自然数。以下同様)をそれぞれ順次記憶する出力メモリアレイ部MAoと、各加算活性化部Actにバイアス用の値を順次提供するバイアス用メモリアレイMAbとを有する。

10

【0029】

メモリアクセス制御部MCntは、例えば、Direct Memory Access Controllerである。メモリアクセス制御部MCntは、制御部Cntの制御に従い、各プロセスエレメント部Peに逐次供給する入力データを、入力メモリアレイ部MAiに設定する。また、メモリアクセス制御部MCntは、制御部Cntの制御に従い、重み付け係数およびニューロン間の接続の有無を示す所定値を、各メモリセルレイ部MCに予め設定する。また、メモリアクセス制御部MCntは、制御部Cntの制御に従い、加算活性化部Actから出力された出力データを、出力メモリアレイ部MAoから取り出す。

20

【0030】

制御部Cntは、CPU (Central Processing Unit)等を有する。制御部Cntは、ニューラル電子回路NNの各素子の同期等のタイミングを計ったり、計算やデータの転送の同期を取ったりする。また、制御部Cntは、ニューラル電子回路NN内の後述のセレクト素子の切り替え制御を行う。

【0031】

制御部Cntは、メモリアクセス制御部MCntを制御して、他のコア電子回路Coreから出力されたデータを入力メモリアレイ部MAi用に整えて、入力データとして入力メモリアレイ部MAiに供給するように制御する。制御部Cntは、メモリアクセス制御部MCntが、出力メモリアレイ部MAoから取得した出力データを、他のコア電子回路Coreに転送するように制御する。

30

【0032】

なお、上位コントローラ (図示せず)が、ニューラルネットワーク・システムNNSや、各コア電子回路Coreの制御部Cntを制御してもよい。また、上位コントローラが、制御部Cntの代わりに、ニューラル電子回路NNおよびメモリアクセス制御部MCntを制御してもよい。上位コントローラは、外付けのコンピュータでもよい。

【0033】

バイアス用メモリアレイ部MAbは、各加算活性化部Actに提供するバイアス用のデータを予め記憶している。

【0034】

(2.2 ニューラル電子回路の構成および機能)

40

次に、ニューラル電子回路NNについて、図3を用いて説明する。

図3は、図2のニューラル電子回路の一例を示すブロック図である。

【0035】

図3に示すように、ニューラル電子回路NNは、例えば、入力 m 個 \times 出力 n 個の2層のニューラルネットワークを実現する。

【0036】

記憶部の一例であるメモリセルレイ部MCは、重み付け係数を記憶するメモリセル10を有する。メモリセル10は、構築するニューラルネットワークによる実現する脳機能に基づいて予め設定された「1」または「0」の1ビットの重み付け係数を記憶している。

50

【0037】

なお、メモリセルアレイ部MCは、上記脳機能に基づいて予め設定されたニューロン間の接続有無情報を記憶している別の接続有無情報用のメモリセル（図示せず）も有していてもよい。ここで、接続無情報は、例えば、NC（Not Connected（接続なし））を意味する1ビット所定値であり、所定値として「1」または「0」等が割り当てられる。

【0038】

メモリセル10が並んで、メモリセルの列が形成される。同時に各プロセスエレメント部Peに出力されるメモリセル10をまとめてメモリセルブロックCBが形成される。メモリセルブロックCBのメモリセル10は、並列で入力される各入力データに対応している。

10

【0039】

メモリセルアレイ部MCは、入力メモリアレイ部MAiから並列で入力される入力データ I_1 、 \dots 、 I_m の入力並列数 m 個以上のメモリセルブロックCBを有することが好ましい。メモリセルブロックCBにおいて、メモリセル10の数は、入力メモリアレイ部MAiから1ビット順次入力されるシリアルの入力データのサイクル数以上が好ましい。

【0040】

メモリセルアレイ部MCは、メモリセルブロックCB毎に、1ビットの重み付け係数を、1ビット順次入力されるシリアルの入力データに対応してプロセスエレメント部Peに、順次出力する。各プロセスエレメント部Peに、メモリセルブロックCBからの重み付け係数と、入力メモリアレイ部MAiからの入力データとが入力される。

20

【0041】

メモリセルブロックCBは、1ビットの重み付け係数と1ビットの接続有無情報とを、交互に、プロセスエレメント部Peに順次、出力してもよい。メモリセル10は、プロセスエレメント部Peに対して独立した結線を有し、別々に順次、プロセスエレメント部Peに出力してもよい。

【0042】

メモリセルアレイ部MCは、図2および図3に示すように、出力メモリアレイ部MAoに並列に出力される出力データの出力並列数 n 個、並列に出力される出力データ O_1 、 \dots 、 O_n に対応して、ニューラル電子回路NNに配置される。

【0043】

このように、メモリセルアレイ部MCは、並列で入力される1ビットの並列の各入力データに応じて、「1」または「0」の1ビットの重み付け係数を記憶し、当該重み付け係数を出力する記憶部の一例として機能する。メモリセルアレイ部MCは、並列で出力される前記並列の各出力データに応じて配置された記憶部の一例として機能する。メモリセルアレイ部MCは、「1」または「0」の重み付け係数、および、ニューロン間の接続の有無を示す所定値の組を複数記憶し、当該重み付け係数および当該所定値を出力する記憶部の一例として機能する。

30

【0044】

なお、メモリセル10、および、メモリセルブロックCBの構成および機能の詳細は、図15以降、特に、図15および図16のメモリセル1に関する記述部分、図21から図27のメモリセル10、および、メモリセルブロック15に関する記述部分等で後述する。また、メモリセルアレイ部MCは、後述するメモリセルアレイMC1、メモリセルアレイMC2に対応する。

40

【0045】

図3に示すように、並列で入力される並列の各入力データに配置された入力並列数 m 個のプロセスエレメント部Peは、ニューラル電子回路NNにおいて、プロセスエレメント・コラム（例えば、プロセスエレメント・コラムPC₁）を形成する。出力並列数 n 個のプロセスエレメント・コラムPC₁からPC_nは、並列に出力される出力データに対応して、ニューラル電子回路NNにおいて、 n 列に配置される。プロセスエレメント部Peは、図3に示すように、ニューラル電子回路NNにおいて、 m 行 \times n 列に2次元の演算器ア

50

レイとして設定される。

【0046】

行列(1, 1)、(1, 2)、・・・、(1, n)のプロセスエレメント部 P e には、入力データ I 1 が共通で入力される結線になっている。行列(2, 1)、(2, 2)、・・・、(2, n)のプロセスエレメント部 P e には、入力データ I 2 が共通で入力される結線になっている。行列(m, 1)、(m, 2)、・・・、(m, n)のプロセスエレメント部 P e には、入力データ I m が共通で入力される結線になっている。

【0047】

プロセスエレメント部 P e は、対応するメモリセル 10 から出力される 1 ビットの重み付け係数と、1 ビットの入力データとの排他的論理和 (X N O R) を乗算結果として算出して出力する。

10

【0048】

なお、接続有無情報用のメモリセルからの接続無情報(例えば、「N C」を意味する所定値)が出力された場合、加算活性化部 A c t において、乗算結果が加算されない。例えば、乗算結果と接続有無情報とが交互にペアで出力されてもよい。また、接続有無情報に関して、プロセスエレメント部 P e から加算活性化部 A c t へ、乗算結果とは独立の結線を有し、乗算結果と接続有無情報とは別々に出力されてもよい。

【0049】

なお、プロセスエレメント部 P e で、乗算結果の部分和を計算するとき、接続有無情報用のメモリセルからの接続無情報(例えば、「N C」を意味する所定値)が出力された場合、乗算結果の部分和に加算されない。

20

【0050】

このように、プロセスエレメント部 P e は、前記並列の各入力データに設定された電子回路部であって、前記記憶部から出力された重み付け係数と前記入力データとを乗算する乗算機能を実現する第 1 電子回路部の一例として機能する。プロセスエレメント部 P e は、前記入力データの値と前記記憶部から出力された値とが一致する場合に前記入力データの入力に対応して「1」を出力し、前記入力データの値と前記記憶部から出力された値とが異なる場合に前記入力データの入力に対応して「0」を出力する第 1 電子回路部の一例として機能する。プロセスエレメント部 P e は、前記入力データの値と前記記憶部から出力された値とが一致する場合に前記入力データの入力に対応して「1」を出力し、前記入力データの値と前記記憶部から出力された値とが異なる場合に前記入力データの入力に対応して「0」を出力し、前記所定値が前記記憶部から出力された場合に前記入力データの入力に対応して当該所定値を出力する第 1 電子回路部の一例として機能する。

30

【0051】

プロセスエレメント・コラム P C ₁、・・・、P C _n は、各プロセスエレメント部 P e からの乗算結果または一部の乗算結果を加算した部分和結果等を加算活性化部 A c t に出力する。

【0052】

なお、プロセスエレメント部 P e に対応する多数判定入力回路の構成および機能の詳細は、図 15 以降、特に、図 21 および図 22 等の多数判定入力回路 12 に関する記述部分等で後述する。ちなみに、多数判定入力回路 12 は、対応する接続有無情報用のメモリセルから出力される接続有無情報をそのまま出力したり、対応するメモリセル 10 から出力される重み付け係数と入力データとの排他的論理和を算出し、出力データとして出力したりする。

40

【0053】

図 2 および図 3 に示すように、加算活性化部 A c t は、並列で出力される各出力データ O ₁、・・・、O _n に応じて配置されている。

【0054】

加算活性化部 A c t は、プロセスエレメント・コラムから逐次出力される乗算結果を、接続有無情報に基づき加算して、加算結果に活性化関数を適用して 1 ビットの出力データ

50

を出力メモリアレイ部 M A o に出力する。プロセスエレメント部 P e が乗算結果の部分積を出力する場合、加算活性化部 A c t は、プロセスエレメント・コラムから逐次出力される乗算結果を加算して、加算結果に活性化関数を適用して 1 ビットの出力データを出力メモリアレイ部 M A o に出力する。

【 0 0 5 5 】

加算活性化部 A c t は、プロセスエレメント・コラムにおいて、入力データの 1 サイクル単位で、乗算結果として「 1 」を算出した回数から、乗算結果として「 0 」を算出した回数を減じた値が予め定められた閾値以上の場合に「 1 」を出力データとして出力し、減じた値が閾値未満の場合に「 0 」を前記出力データとして出力する。

【 0 0 5 6 】

このように、加算活性化部 A c t は、前記並列の各入力データの前記第 1 電子回路部からの各乗算結果を加算し、かつ、当該加算結果に活性化関数を適用して 1 ビットの出力データを出力する加算・適用機能を実現する第 2 電子回路部の一例として機能する。加算活性化部 A c t は、並列で出力される前記並列の各出力データに応じて配置された第 2 電子回路部の一例として機能する。加算活性化部 A c t は、前記並列の各入力データに設定された複数の第 1 電子回路部において、前記入力データが入力されるサイクル単位で、前記第 1 電子回路部が「 1 」を算出した回数から、前記第 1 電子回路部が「 0 」を算出した回数を減じた値が予め定められた閾値以上の場合に「 1 」を前記出力データとして出力し、前記減じた値が前記閾値未満の場合に「 0 」を前記出力データとして出力する第 2 電子回路部の一例として機能する。

【 0 0 5 7 】

なお、加算活性化部 A c t の構成および機能の詳細は、図 1 5 以降、特に図 2 1 から図 2 7 で後述する。ここで、加算活性化部 A c t は、例えば、後に詳細に説明する直列多数判定回路 1 3 に対応する。ちなみに、直列多数判定回路 1 3 は、時間的に直列に順次入力される「 1 」または「 0 」の乗算結果に対して、「 1 」または「 0 」の数に関する判定を行う。

【 0 0 5 8 】

図 3 に示すように、ビットシリアル入力の並列化が行われ、プロセスエレメント部 P e の行が入力データに対して共有し、プロセスエレメント部 P e の列である各プロセスエレメント・コラムが独立して出力データを出力する。

【 0 0 5 9 】

(1 . 3 プロセスエレメント・コラムの構成および機能)

次に、プロセスエレメント・コラムの構成および機能について、図 4 および図 5 を用いて説明する。

【 0 0 6 0 】

図 4 は、図 3 のプロセスエレメント・コラムの一例を示すブロック図である。図 5 は、図 4 のプロセスエレメントの動作のタイミングの一例を示す模式図である。

【 0 0 6 1 】

図 4 に示すように、プロセスエレメント・コラム P C ₁ のようなプロセスエレメント・コラムは、フェーズ 1 として計算を行う複数のプロセスエレメント部 P e と、フェーズ 1 での計算結果の転送を行うフェーズ 2 の複数のフリップフロップ F p (一時記憶部の一例) およびセクタ S e と、を有する。

【 0 0 6 2 】

フリップフロップ F p は、各プロセスエレメント部 P e の出力側に接続され、プロセスエレメント部 P e の乗算結果または部分積結果を一時的に記憶する。フリップフロップ F p は、 1 行目のプロセスエレメント部 P e から n 行目のプロセスエレメント部 P e に対応して、セクタ S e を介して直列に接続されている。 n 行目のフリップフロップ F p は、加算活性化部 A c t に接続されている。なお、これらの接続が、図 2 において、プロセスエレメント部 P e 間を太線で示した部分の機能の一例である。

【 0 0 6 3 】

10

20

30

40

50

セレクタ S_e は、プロセスエレメント部 P_e 間に配置され、上流のフリップフロップ F_p からのデータと、プロセスエレメント部 P_e からのデータとを切り替える。

【0064】

図5に示すように、フェーズ1において、各フリップフロップ F_p で乗算等の計算が行われ、セレクタ S_e は、入力側のフリップフロップ F_p のデータを選択して、フリップフロップ F_p に計算結果が出力される。次に、フェーズ2において、セレクタ S_e は、上流のフリップフロップ F_p のデータを選択する。これにより、入力並列数 m 個のサイクル単位のタイミングで、計算結果が加算活性化部 $A_c t$ に逐次転送され、計算結果の転送が行われる。

【0065】

このように、フリップフロップ F_p は、前記第1電子回路部毎に、各前記第1電子回路部からの前記乗算結果を一時記憶する一時記憶部の一例として機能する。フリップフロップ F_p は、直列に設定され、前記乗算結果を前記第2電子回路部へ順次転送する一時記憶部の一例として機能する。

【0066】

なお、各プロセスエレメント部 P_e からの乗算結果が直接、加算活性化部 $A_c t$ に供給されてもよい。この場合、加算活性化部 $A_c t$ は、各入力データが入力される並列多数判定回路である。並列多数判定回路20の詳細については、後述する。ちなみに、並列多数判定回路20は、並列で入力される「1」または「0」の乗算結果に対して、「1」または「0」の数に関する判定を行う。

【0067】

(2.4 プロセスエレメントおよび加算活性化部の回路構成)

次に、プロセスエレメントおよび加算活性化部の回路構成について、図6Aおよび図6Bを用いて説明する。

【0068】

図6Aに示すように、プロセスエレメント部 P_e は、XNOR素子 $p_e 1$ と、セレクタ $p_e 2$ と、加算器 $p_e 3$ と、フリップフロップ $p_e 4$ と、を有する。

【0069】

XNOR素子 $p_e 1$ は、1ビットの重み付け係数(例えば、 w_1)と、1ビットの入力データ(例えば、 i_1)との排他的論理和を乗算結果として算出して出力する。

【0070】

セレクタ $p_e 2$ は、乗算結果に基づき、値(0、1)を、正負の値(-1、1)または(1、-1)に変換する。例えば、セレクタ $p_e 2$ は、乗算結果が「0」の場合、「-1」を選択し、乗算結果が「1」の場合、「1」を選択する。

【0071】

加算器 $p_e 3$ は、フリップフロップ $p_e 4$ に一時記憶された前の値と、セレクタ $p_e 2$ の出力値(1または-1)とを、例えば、10ビット程度の精度で加算する。加算器 $p_e 3$ の加算結果は、フリップフロップ $p_e 4$ に記憶される。例えば、加算器 $p_e 3$ およびフリップフロップ $p_e 4$ は、入力データの並列数回ループして、乗算結果の部分積を、例えば、10ビットで出力する。

【0072】

次に、加算活性化部 $A_c t$ の回路構成について、図6Bを用いて説明する。

【0073】

図6Bに示すように、加算活性化部 $A_c t$ は、セレクタ $a_c 1$ と、加算器 $a_c 2$ と、フリップフロップ $a_c 3$ と、符号判定器 $a_c 4$ と、セレクタ $a_c 5$ と、フリップフロップ $a_c 6$ と、を有する。

【0074】

セレクタ $a_c 1$ は、加算器 $a_c 2$ への入力を制御する。セレクタ $a_c 1$ は、各プロセスエレメント・コラム $P C_1$ 、 \dots 、 $P C_n$ からの乗算結果の部分積のデータ(例えば、20ビット)を入力とし、最後にフリップフロップ $a_c 6$ からのバイアス $b i a s$ の値を

10

20

30

40

50

加算器 a c 2 で加算するように制御する。

【 0 0 7 5 】

加算器 a c 2 およびフリップフロップ a c 3 により加算を行い、加算結果のデータ（例えば、32ビット）を符号判定器 a c 4 に出力する。

【 0 0 7 6 】

符号判定器 a c 4 は、加算器 a c 2 およびフリップフロップ a c 3 で乗算結果の部分積の総和と、フリップフロップ a c 6 からのバイアス b i a s の値とを加算された値の符号を判定し、1ビットの出力データを出力する。なお、ゼロの値の場合、符号は正としてもよい。

【 0 0 7 7 】

なお、セクタ a c 5 およびフリップフロップ a c 6 は、バイアス用メモリアレイ M A b から出力されたバイアス b i a s の値（例えば、32ビット値）を、次の加算活性化部 A c t に転送するか保持するかを制御する。

【 0 0 7 8 】

[3 . ニューラル電子回路の適用例]

次に、ニューラル電子回路 N N により、様々なタイプのニューラルネットワークを実現する実施例について説明する。

【 0 0 7 9 】

(3 . 1 コンボリューション演算を実現するニューラル電子回路)

次に、コンボリューション演算を実現するニューラル電子回路について、図 7 および図 8 を用いて説明する。

【 0 0 8 0 】

図 7 は、コンボリューション演算におけるデータの関係の一例を示す模式図である。図 8 は、図 7 のコンボリューション演算を実現するニューラル電子回路の一例を示すブロック図である。

【 0 0 8 1 】

図 7 に示すように、チャンネル数 C I の入力画像 I i m の入力データと、フィルタ画像 P a 、 P b 、 . . . 、 P c o の種類数 C O のフィルタデータとに対して、コンボリューション演算が行われる。ここで、R G B のようなカラーが画像の場合、チャンネル数は、3 である。白黒画像の場合、チャンネル数は 1 である。C M Y K カラーモデルの画像の場合、チャンネル数は 4 である。

【 0 0 8 2 】

図 7 に示すように、1ビットの値の $k \times k$ ピクセルの入力画像 I i m に対して、個々の入力データ i_1 、 i_2 、. . .、 i_k 、. . .、 i_{k^2} が形成される。なお、原画像に対して、 $k \times k$ ピクセルの領域を順次切り出して、コンボリューション演算が原画像に対して行われる。ここで、コンボリューション演算は、第 1 関数を平行移動しながら第 2 関数に重ね足し合わせる二項演算である。例えば、入力画像 I i m が第 1 関数に対応して、フィルタ画像が第 2 関数に対応する。

【 0 0 8 3 】

ここで、図 7 において、入力データ I_1 は、チャンネル 1 の入力データの総称を示し、入力データ i_1 、 i_2 、. . .、 i_k 、. . .、 i_{k^2} は、チャンネル 1 に順次入力される 1 ビットの個々の入力データである。入力データ I_2 は、チャンネル 2 の入力データの総称を示し、灰色のボックスで示す入力データ i_1 、 i_2 、. . .、 i_k 、. . .、 i_{k^2} は、チャンネル 2 に順次入力される 1 ビットの個々の入力データである。

【 0 0 8 4 】

フィルタデータは、入力画像 I i m に対応して、1ビットの値の $k \times k$ ピクセルのフィルタ画像 P a 、 P b 、 . . . 、 P c o である。カラーの場合、例えば、3チャンネル分の要素画像が一组で、C O 個の種類フィルタ画像が用意される。

【 0 0 8 5 】

1 枚の $k \times k$ ピクセルの入力画像 I i m と、1 枚の $k \times k$ ピクセルのフィルタ画像（例

10

20

30

40

50

えば、1のフィルタ画像 P_a)から、コンボリューション演算により、1ビットの出力データが出力される。 C_I 個のチャンネル毎の1ビットの出力データに対して、フィルタ画像の C_O 種類分である $C_I \times C_O$ 個の出力データが生成される。

【0086】

図8に示すように、ニューラル電子回路 NN は、チャンネル数のプロセスエレメント部 P_e が並んだプロセスエレメント・コラム PC_1 、 PC_2 、 \dots 、 PC_{C_O} と、 C_O 個の種類の各フィルタ画像に対応したメモリセルアレイ部 MC とを有する。制御部 Cnt は、ニューラル電子回路 NN のうち、プロセスエレメント・コラム PC_1 、 PC_2 、 \dots 、 PC_{C_O} と、 C_O 個のメモリセルアレイ部 MC とを、使用するように制御する。

【0087】

メモリアクセス制御部 $MCnt$ は、メモリセルアレイ部 MC のメモリセル 10 の列(k^2 個の行のメモリセル 10)に、重み付け係数として、フィルタ画像の $k \times k$ ピクセルに対応した1ビットの値を、設定する。各チャンネルに対応した C_I 個のメモリセルの列に、データが設定される。

【0088】

メモリアクセス制御部 $MCnt$ は、入力メモリアレイ部 MAi に、チャンネル毎に k^2 個の入力データ i_1 、 i_2 、 \dots 、 i_k 、 \dots 、 i_{k^2} を設定する。

【0089】

ニューラル電子回路 NN は、先ず、チャンネル数 C_I 個分の入力データを順次処理する。

【0090】

具体的には、チャンネル1の入力データ I_1 のうち入力データ i_1 、 i_2 、 \dots 、 i_{C_I} が、行列 $(1, 1)$ 、 $(1, 2)$ 、 \dots 、 $(1, C_O)$ の各プロセスエレメント部 P_e に順次入力される。

【0091】

入力データ i_1 、 i_2 、 \dots 、 i_{C_I} に同期して、行列 $(1, 1)$ のプロセスエレメント部 P_e に、メモリセルアレイ部 MC から出力された重み付け係数 w_1 、 w_2 、 \dots 、 w_{C_I} も順次入力される。

【0092】

このように、メモリセルアレイ部 MC が、前記第1電子回路部に順次入力される前記入力データに対応した重み付け係数を、前記第1電子回路部に順次出力する記憶部の一例として機能する。

【0093】

図5に示すように、フェーズ1で、プロセスエレメント・コラム PC_1 に、おいては、行列 $(1, 1)$ のプロセスエレメント部 P_e が、乗算結果 $i_1 \times w_1$ 、 $i_2 \times w_2$ 、 \dots 、 $i_{C_I} \times w_{C_I}$ を算出し、チャンネル数 C_I 個分の和であるチャンネル1の部分 $i_1 \times w_1 + i_2 \times w_2 + \dots + i_{C_I} \times w_{C_I}$ を算出する。部分 $i_1 \times w_1 + i_2 \times w_2 + \dots + i_{C_I} \times w_{C_I}$ は、並列で入力される入力データの入力並列数分、乗算結果を加算した部分加算結果の一例である。

【0094】

チャンネル2の入力データ I_2 のうち、図8中、グレイの四角形で示した入力データ i_1 、 i_2 、 \dots 、 i_{C_I} が、行列 $(2, 1)$ 、 $(2, 2)$ 、 \dots 、 $(2, C_O)$ の各プロセスエレメント部 P_e に順次入力される。行列 $(2, 1)$ のプロセスエレメント部 P_e が、チャンネル2の入力データ I_2 に対して、乗算結果を算出し、チャンネル2の部分 $i_1 \times w_1 + i_2 \times w_2 + \dots + i_{C_I} \times w_{C_I}$ を算出する。

【0095】

チャンネル C_I に関しても、行列 $(C_I, 1)$ のプロセスエレメント部 P_e が、チャンネル C_I の入力データ I_{C_I} に対して、乗算結果を算出し、部分 $i_1 \times w_1 + i_2 \times w_2 + \dots + i_{C_I} \times w_{C_I}$ を算出する。

【0096】

このように、プロセスエレメント部 P_e が、並列で入力される前記入力データの入力並

10

20

30

40

50

列数分、前記乗算結果を加算した部分加算結果を出力する第1電子回路部の一例として機能する。

【0097】

次に、フェーズ2で、プロセスエレメント・コラム PC_1 は、行列 $(1, 1)$ 、 $(2, 1)$ 、 \dots 、 $(CI, 1)$ の各プロセスエレメント部 Pe から出力されたチャンネル毎の部分加算結果を順に加算活性化部 Act に転送する。

【0098】

フェーズ1の次の計算では、入力データ i_{CI+1} 、 i_{CI+2} 、 \dots 、 i_{2CI} に対して、行列 $(1, 1)$ のプロセスエレメント部 Pe が、乗算結果 $i_{CI+1} \times W_{CI+1}$ 、 $i_{CI+2} \times W_{CI+2}$ 、 \dots 、 $i_{2CI} \times W_{2CI}$ を算出し、部分加算結果 $i_{CI+1} \times W_{CI+1} + i_{CI+2} \times W_{CI+2} + \dots + i_{2CI} \times W_{2CI}$ を算出する。

10

【0099】

k^2 目の入力データまで、チャンネル数 CI 個分の入力データで、乗算結果および部分加算結果を算出して転送する。 $k \times k$ ピクセルの入力画像に対して k^2 サイクル単位のシリアル入力が形成され、各フィルタ画像に対して、1画素分、出力される。

【0100】

プロセスエレメント・コラム PC_2 等も、同様に部分加算結果を算出して加算活性化部 Act に転送する。

【0101】

加算活性化部 Act は、チャンネル毎に部分加算結果の和を計算して、入力データ k^2 個に対する総和をコンボリューション演算の結果として算出する。加算活性化部 Act は、入力データの重み付け総和に閾値であるバイアス $bias$ の値を加え、活性化関数を適用として、例えば、符号が正の値の場合（すなわち、所定の閾値を超えた場合）、「1」を、符号が負の値の場合（すなわち、所定の閾値を超えない場合）は「0」を、出力データ O_o 等として、出力メモリアレイ部 MA_o に出力する。入力画像 I_{im} とフィルタ画像 P_a とに対する出力結果が、出力データ o_a である。入力画像 I_{im} とフィルタ画像 P_b とに対する出力結果が、出力データ o_b である。各チャンネル毎に、出力データが計算される。なお、出力データ o_a 等が、コンボリューション演算の結果として設定されてもよい。

20

【0102】

このように、加算活性化部 Act は、前記部分加算結果から前記加算結果を算出する第2電子回路部の一例として機能する。

30

【0103】

出力メモリアレイ部 MA_o は、チャンネル数 CI 個分の出力データ o_a 、 \dots を1ワードとして記憶する。出力メモリアレイ部 MA_o は、フィルタ画像 P_a 、 P_b 、 \dots 、 P_c 毎に、1ワードの出力データ o_a 、 \dots 、1ワードの出力データ o_b 、 \dots 、 \dots を記憶する。

【0104】

(3.2 フル結合のニューラルネットワークを実現するニューラル電子回路)

次に、ニューロン層間のニューロン同士がフル結合したフル結合のニューラルネットワークを実現するニューラル電子回路について、図9および図10を用いて説明する。

40

【0105】

図9は、フル結合のニューラルネットワークの一例を示す模式図である。図10は、図9のフル結合のニューラルネットワークを実現するニューラル電子回路の一例を示すブロック図である。

【0106】

$M \times N$ のニューラル電子回路 NN で、入力並列数 M 個以上の入力で、出力並列数 N 個以上の出力のフル結合のニューラルネットワークを実現する場合について説明する。例えば、図9は、 $AM \times BN$ のフル結合のニューラルネットワークで $M = 3$ 、 $N = 2$ で、 $A = 2$ 、 $B = 3$ の一例を示している。

【0107】

50

図10に示すように、ニューラル電子回路NNは、入力並列数M個のプロセスエレメント部Peが並んだプロセスエレメント・コラムPC₁、PC₂、・・・、PC_Nと、出力並列数N個のプロセスエレメント・コラムPC₁、PC₂、・・・、PC_Nと、出力並列数N個のメモリセルアレイ部MCとを有する。制御部Cntは、ニューラル電子回路NNのうち、プロセスエレメント・コラムPC₁、PC₂、・・・、PC_Nと、N個のメモリセルアレイ部MCと、を使用するように制御する。ここで、入力並列数M個が、入力可能並列数の一例で、出力並列数N個が、出力可能並列数の一例である。

【0108】

メモリアクセス制御部MCntは、入力メモリアレイ部MAiに、並列に入力データi₁、i₂、・・・、i_Mを設定して、次のi_{M+1}、i_{M+2}、・・・、i_{2M}を設定して、次々に、入力データi_{A M}まで、を設定する。メモリアクセス制御部MCntは、入力データi₁から入力データi_{A M}までを、B個回繰り返して、入力メモリアレイ部MAiにデータを設定する。

10

【0109】

メモリアクセス制御部MCntは、メモリセルアレイ部MCのメモリセル10の列(A×B個の行のメモリセル10)に、A×B個の重み付け係数の1ビットの値を、予め設定する。例えば、メモリアクセス制御部MCntは、メモリセルアレイ部MCのメモリセルの列のメモリセル10に、入力データi₁、i_{M+1}、i_{2M+1}、・・・、i_{(A-1)M+1}に対応して重み付け係数w₁、w_{M+1}、w_{2M+1}、・・・、w_{(A-1)M+1}をB個回繰り返した重み付け係数を設定する。入力並列数M個のメモリセルの列に、データが予め設定される。

20

【0110】

ニューラル電子回路NNは、先ず、入力並列数M個分の入力データを並列処理する。

【0111】

具体的には、入力データi₁が、行列(1,1)、(1,2)、・・・、(1,N)の各プロセスエレメント部Peに入力される。入力データi₂が、行列(2,1)、(2,2)、・・・、(2,N)の各プロセスエレメント部Peに入力される。入力データi_Mが、行列(M,1)、(M,2)、・・・、(M,N)の各プロセスエレメント部Peに入力される。

【0112】

入力データi₁に同期して、行列(1,1)のプロセスエレメント部Peに、メモリセルアレイ部MCから出力された重み付け係数w₁も入力される。入力データi₂に同期して、行列(2,1)のプロセスエレメント部Peに、メモリセルアレイ部MCから出力された重み付け係数w₂も入力される。入力データi_Mに同期して、行列(M,1)のプロセスエレメント部Peに、メモリセルアレイ部MCから出力された重み付け係数w_Mも入力される。

30

【0113】

このように、メモリセルアレイ部MCが、並列で入力される並列の各前記入力データに対応した重み付け係数を、各前記第1電子回路部に出力する記憶部の一例として機能する。

40

【0114】

図5に示すように、フェーズ1で、プロセスエレメント・コラムPC₁においては、行列(1,1)のプロセスエレメント部Peが、乗算結果i₁×w₁を算出し、行列(2,1)のプロセスエレメント部Peが、乗算結果i₂×w₂を算出し、行列(M,1)のプロセスエレメント部Peが、乗算結果i_M×w_Mを算出する。

【0115】

次に、フェーズ2で、プロセスエレメント・コラムPC₁は、行列(1,1)、(2,1)、・・・、(M,1)の各プロセスエレメント部Peから出力された乗算結果i₁×w₁、乗算結果i₂×w₂、・・・、乗算結果i_M×w_Mを、乗算結果i_M×w_Mから順に加算活性化部Actに転送する。

50

【0116】

次に、加算活性化部 A c t は、出力データ O_1 に関して、 $A \times M$ 個の総和に対して、 M 個分の和である部分 $i_1 \times w_1 + i_2 \times w_2 + \dots + i_M \times w_M$ を生成する。

【0117】

プロセスエレメント・コラム PC_2 においては、フェーズ 1 で、行列 $(1, 2)$ のプロセスエレメント部 P e が、入力データ i_1 に関する乗算結果を算出し、行列 $(2, 2)$ のプロセスエレメント部 P e が、入力データ i_2 に関する乗算結果を算出し、行列 $(M, 2)$ のプロセスエレメント部 P e が、入力データ i_M に関する乗算結果を算出する。

【0118】

次に、フェーズ 2 で、プロセスエレメント・コラム PC_2 は、行列 $(1, 2)$ 、 $(2, 2)$ 、 \dots 、 $(M, 2)$ の各プロセスエレメント部 P e から出力された各乗算結果を、入力データ i_M に関する乗算結果から順に加算活性化部 A c t に転送する。

10

【0119】

次に、加算活性化部 A c t は、出力データ O_2 に関して、 M 個の分の和である部分 i_M を生成する。

【0120】

プロセスエレメント・コラム PC_N においても同様に、乗算結果が算出される。

【0121】

次の入力データを入力するタイミングで、フェーズ 1 で、プロセスエレメント・コラム PC_1 においては、行列 $(1, 1)$ のプロセスエレメント部 P e が、乗算結果 $i_{M+1} \times w_{M+1}$ を算出し、行列 $(2, 1)$ のプロセスエレメント部 P e が、乗算結果 $i_{M+2} \times w_{M+2}$ を算出し、行列 $(M, 1)$ のプロセスエレメント部 P e が、乗算結果 $i_{2M} \times w_{2M}$ を算出する。

20

【0122】

次に、フェーズ 2 で、プロセスエレメント・コラム PC_1 は、行列 $(1, 1)$ 、 $(2, 1)$ 、 \dots 、 $(M, 1)$ の各プロセスエレメント部 P e から出力された乗算結果 $i_{M+1} \times w_{M+1}$ 、乗算結果 $i_{M+2} \times w_{M+2}$ 、 \dots 、乗算結果 $i_{2M} \times w_{2M}$ を、乗算結果 $i_{2M} \times w_{2M}$ から順に加算活性化部 A c t に転送する。

【0123】

次に、加算活性化部 A c t は、出力データ O_1 に関して、部分 $i_{M+1} \times w_1 + i_{M+2} \times w_{M+2} + \dots + i_{2M} \times w_{2M}$ を生成する。

30

【0124】

以上、 $A \times M$ 番目の入力データ i_{AM} まで、繰り返し、各加算活性化部 A c t は、部分 i_{AM} の総和を計算して、閾値であるバイアス $bias$ の値を加え、活性化関数を適用として、例えば、符号が正の値の場合（すなわち、所定の閾値を超えた場合）、「1」を、符号が負の値の場合（すなわち、所定の閾値を超えない場合）は「0」とした出力データ o_1 、 \dots 、 o_N を算出して、出力メモリアレイ部 M A o に出力する。

【0125】

出力データ o_{N+1} 、 o_{N+2} 、 \dots 、 o_{N+1} に関しても、上記のように、入力データ i_1 から入力データ i_{AM} まで処理をして、各加算活性化部 A c t は、 A 個の部分 i_{AM} の総和を計算して、活性化関数を適用して、出力データ o_{N+1} 、 o_{N+2} 、 \dots 、 o_{N+1} を算出して、出力メモリアレイ部 M A o に出力する。

40

【0126】

同様な計算で、ニューラル電子回路 NN は、出力データ o_{BN} まで計算する。 $A \times M$ 個の入力データに対して、 M 個の並列入力で、 $A \times B$ サイクル単位のシリアル入力が形成され、 $B \times N$ 個の出力データに対して、 N 個の並列出力で、 B 個出力される。

【0127】

このように、メモリアレイ部 M C が、前記入力データを並列で一度に入力可能な入力可能並列数より、前記入力データの入力並列数が大きい場合、前記残りの入力データに対応する前記重み付け係数を出力する記憶部の一例として機能する。プロセスエレメント

50

部 P e が、前記入力可能並列数の並列で前記入力データを受け入れた後、前記入力可能並列数の並列で受け入れできなかった残りの前記入力データを入れ受け入れる前記第 1 電子回路部の一例として機能する。

【 0 1 2 8 】

(3 . 3 コア電子回路同士の連結)

次に、コア電子回路 C o r e 同士の連結により、ニューロン層における層内拡張をしたニューラルネットワーク、層数を拡張するニューラルネットワークを実現する実施例について図を用いて説明する。

【 0 1 2 9 】

図 1 1 は、ニューラルネットワークの層内拡張の一例を示す模式図である。図 1 2 は、図 1 1 の層内拡張を実現するコア電子回路の接続の一例を示すブロック図である。図 1 3 は、ニューラルネットワークの層数拡張の一例を示す模式図である。図 1 4 は、図 1 3 の層数拡張を実現するコア電子回路の接続の一例を示すブロック図である。

10

【 0 1 3 0 】

図 1 1 に示すように、出力側の層内拡張するためには、図 1 2 に示すように、入力データに対して、コア電子回路 C o r e を並列に接続すればよい。

【 0 1 3 1 】

図 1 1 に示すように、入力 3 個 × 出力 2 個の 2 層のニューラルネットワークと、入力 3 個 × 出力 4 個の 2 層のニューラルネットワークと、並列に接続してもよいし、入力 3 個 × 出力 3 個の 2 層のニューラルネットワークと、入力 3 個 × 出力 3 個の 2 層のニューラルネットワークと、並列に接続してもよい。

20

【 0 1 3 2 】

図 1 3 に示すように、層数拡張するためには、図 1 4 に示すように、コア電子回路 C o r e を直列に接続すればよい。

【 0 1 3 3 】

図 1 3 に示すように、入力 3 個 × 出力 2 個の 2 層のニューラルネットワークと、入力 2 個 × 出力 4 個の 2 層のニューラルネットワークと、を直列に接続する。

【 0 1 3 4 】

実際の接続は、メモリアクセス制御部 M C n t を介して、システムバス bus によって接続される。更に、メモリアクセス制御部 M C n t が、入力メモリアレイ部 M A i およびメモリセルアレイ部 M C を設定して、コア電子回路 C o r e の並列接続または直列接続を実現する。

30

【 0 1 3 5 】

以上説明したように、本実施形態によれば、1 ビットの入力データの値とメモリセルアレイ部 M C から出力された値とが一致する場合に入力データの入力に対応して「 1 」を出力し、入力データの値とメモリセルアレイ部 M C から出力された値とが異なる場合に入力データの入力に対応して「 0 」を出力することにより、乗算機能を実現するため、回路規模を縮小でき、並列で入力される並列の各入力データに応じて設定されたプロセスエレメント部 P e により、様々なタイプのニューラルネットワークを実現できる。

40

【 0 1 3 6 】

また、シナプスを独立に、ニューロンの入力を行方向に共有したアレイ構造のプロセスエレメント部 P e により、コンボリューション演算とフル結合演算の両計算を効率化させることができる。

【 0 1 3 7 】

また、メモリセルアレイ部 M C および加算活性化部 A c t が、並列で出力される並列の各出力データに応じて設定された場合、コンボリューション型や、フル結合等の多様性を有するニューラル電子回路を容易に実現することができる。

【 0 1 3 8 】

各プロセスエレメント部 P e からの乗算結果を一時記憶するフリップフロップ F p を備え、各フリップフロップ F p が、直列に設定され、乗算結果を加算活性化部 A c t へ順次

50

転送する場合、更に配線が単純になり、回路面積が小さくなるので、回路規模を縮小できる。また、配線が単純になるので、製造コストが減少する。

【0139】

また、部分加算等の途中計算結果を、プロセスエレメント・コラムにおいて、列方向に伝播させる制御で、演算器であるフリップフロップ F p の使用率を最大化させることができる。

【0140】

加算活性化部 A c t が、並列の各入力データに設定された複数のプロセスエレメント部 P e において、入力データが入力されるサイクル単位で、加算活性化部 A c t が「1」を算出した回数から、加算活性化部 A c t が「0」を算出した回数を減じた値が予め定められた閾値以上の場合に「1」を出力データとして出力し、減じた値が閾値未満の場合に「0」を出力データとして出力する場合、回路規模を縮小した活性化関数を実現できる。

10

【0141】

メモリセルアレイ部 M C が、「1」または「0」の重み付け係数、および、ニューロン間の接続の有無を示す所定値の組を複数記憶し、当該重み付け係数および当該所定値を出力し、プロセスエレメント部 P e が、入力データの値とメモリセルアレイ部 M C から出力された値とが一致する場合に入力データの入力に対応して「1」を出力し、入力データの値とメモリセルアレイ部 M C から出力された値とが異なる場合に入力データの入力に対応して「0」を出力し、所定値がメモリセルアレイ部 M C から出力された場合に入力データの入力に対応して当該所定値を出力する場合、「1」をプラス値、「0」をマイナス値に対応させ、シナプスの接続の有無を表現できるニューラル電子回路が実現できる。

20

【0142】

メモリセルアレイ部 M C が、プロセスエレメント部 P e に順次入力される入力データに対応した重み付け係数を、プロセスエレメント部 P e に順次出力する場合、フィルタ画像に対応したメモリセルアレイ部 M C の重み付け係数に、コンボリューション演算の一方の関数の値を設定し、入力画像に対応した入力データに、コンボリューション演算の他方の関数の値を設定することにより、コンボリューションニューラル電子回路が実現できる。

【0143】

プロセスエレメント部 P e が、並列で入力される入力データの入力並列数分、乗算結果を加算した部分加算結果を出力し、加算活性化部 A c t が、部分加算結果から加算結果を算出する場合、多チャンネルのコンボリューション演算を実現でき、カラー画像等の入力データに対応できる。

30

【0144】

メモリセルアレイ部 M C が、並列で入力される並列の各入力データに対応した重み付け係数を、各プロセスエレメント部 P e に出力する場合、フル結合のニューラル電子回路が実現できる。

【0145】

入力データを並列で一度に入力可能な入力可能並列数より、入力データの入力並列数が大きいとき、プロセスエレメント部 P e は、入力可能並列数の並列で入力データを受け入れた後、入力可能並列数の並列で受け入れできなかった残りの入力データを受け、メモリセルアレイ部 M C は、残りの入力データに対応する重み付け係数を出力する場合、少数の電子回路で、より多くの並列入力数を有するニューラル電子回路が実現できる。

40

【0146】

プロセスエレメント部 P e が格納されたコア電子回路 C o r e を、制御部 C n t またはメモリアクセス制御部 M C n t でコントロールし、任意の大きさのネットワークを計算することができる。また、コア電子回路 C o r e の入出力を制御部 C n t またはメモリアクセス制御部 M C n t でコントロールし、複数のコア電子回路 C o r e に展開することができる。

【0147】

[4 . メモリセル、メモリセルブロック等の詳細な構成および機能]

50

次に、メモリセル 10、接続有無情報用のメモリセル、メモリセルブロック CB に関するメモリセルブロック 15、メモリセルアレイ MC、プロセスエレメント部 Pe、および、加算活性化部 Act 等に関連する詳細な構成および機能について、図を用いて説明する。

【0148】

なお、プロセスエレメント部 Pe は、多数判定入力回路 12 に関係し、加算活性化部 Act は、下記に示される直列多数判定回路 13 に関係する。下記に示されるニューラルネットワーク回路およびニューラルネットワーク集積回路は、ニューラル電子回路 NN に関係する。

【0149】

(I) メモリセル等の実施形態

メモリセル等の実施形態について、図 15 乃至図 20 を用いて説明する。なお、図 15 は実施形態に係るニューラルネットワーク回路を示す図であり、図 16 は当該ニューラルネットワーク回路の細部構成を示す図である。また、図 17 は実施形態に係るニューラルネットワーク集積回路の第 1 例を示す図であり、図 18 は当該ニューラルネットワーク集積回路の第 2 例を示す図であり、図 19 は当該ニューラルネットワーク集積回路の第 3 例を示す図であり、図 20 は当該ニューラルネットワーク集積回路の第 4 例を示す図である。

【0150】

そして、以下に説明する実施形態等に係るニューラルネットワーク回路又はニューラルネットワーク集積回路は、図 1 を用いて説明した一般的なニューラルネットワークを、上記非特許文献 1 又は非特許文献 2 に記載されている手法によりバイナリ化したニューラルネットワーク回路又はニューラルネットワーク集積回路によりモデル化するものである。

【0151】

(A) 実施形態に係るニューラルネットワーク回路について

次に、実施形態に係るニューラルネットワーク回路について、図 15 及び図 16 を用いて例示しつつ説明する。ここで、入力データ I_1 乃至入力データ I_n 又は入力データ I_m に共通の事項を説明する場合、単に「入力データ I 」と称する。出力データ O_1 乃至出力データ O_n 又は出力データ O_m に共通の事項を説明する場合、単に「出力データ O 」と称する。重み付け係数 W_1 乃至重み付け係数 W_n 又は重み付け係数 W_m に共通の事項を説明する場合、単に「重み付け係数 W 」と称する。

【0152】

図 15 (a) に示すように、当該ニューラルネットワーク回路に相当するニューラルネットワーク S では、一つのニューロン NR に対して例えば四つの他のニューロン NR から 1 ビットの入力データ I がそれぞれ入力され、それに対応する出力データ O が当該ニューロン NR から出力される。このとき入力データ I は、その出力元のニューロン NR から見れば 1 ビットの出力データ O となる。また 1 ビットの出力データ O は、その出力先のニューロン NR から見ると 1 ビットの入力データ I となる。上記の通り入力データ I 及び出力データ O はそれぞれ 1 ビットであるため、入力データ I の値及び出力データ O の値は、いずれも「0」又は「1」のいずれかである。そして、図 15 (a) において四つの入力データ I が入力されているニューロン NR (図 15 (a) においてハッチングで示される) において実行される上記乗算処理等に相当する上記式 (1) は、上記式 (1) において $n = 4$ とした場合の式である。即ち上記ニューラルネットワーク S は、並列多入力 - 一出力型の一段ニューラルネットワークである。

【0153】

次に、図 15 (a) に示すニューラルネットワーク S においてハッチングで示されるニューロン NR に相当する実施形態に係るニューラルネットワーク回路の構成を、図 15 (b) にニューラルネットワーク回路 CS として示す。当該ニューラルネットワーク回路 CS は、各々が 1 ビットの入力データ I_1 乃至入力データ I_4 にそれぞれ対応する四つのメモリセル 1 と、多数判定回路 2 と、により構成される。このとき、各メモリセル 1 が本発明

10

20

30

40

50

に係る「第1回路部」の一例、「記憶部」の一例及び「出力部」の一例に、それぞれ相当する。また上記多数判定回路2が、本発明に係る「第2回路部」の一例に相当する。この構成において各メモリセル1は、「1」、又は「0」、或いは「NC」を意味する所定値の三つのいずれか一つを記憶値として記憶すると共に、比較機能を有する三値のメモリセルである。そして各メモリセル1は、それぞれへの入力データIの値と、それぞれの記憶値と、に応じた値を有する出力データ E_1 乃至出力データ E_n を、多数判定回路2にそれぞれ出力する。

【0154】

ここで、メモリセル1の記憶値の一つである上記所定値が意味する上記「NC」は、実施形態に係るニューラルネットワークSにおける二つのニューロンNR間に接続がない状態である。即ち、そのメモリセル1が対応している二つのニューロンNR（即ち入力ニューロンと出力ニューロン）が接続されていない場合、そのメモリセル1の記憶値は上記所定値に設定される。一方、メモリセル1の他の記憶値（「1」又は「0」）のいずれをそのメモリセル1に記憶させておくかは、そのメモリセル1が対応している接続により接続される二つのニューロンNR間の当該接続における重み係数Wに基づいて設定される。ここで各メモリセル1にどのような記憶値を記憶させておくかは、ニューラルネットワークSとしてどのような脳機能をモデル化するか（より具体的には、例えばニューラルネットワークSを構成するニューロンNR間の接続状態等）等に基づいて予め設定されている。なお以下の説明において、出力データ E_1 乃至出力データ E_n に共通の事項を説明する場合、単に「出力データE」と称する。

10

20

【0155】

そして、各メモリセル1における上記記憶値及びそれぞれに入力される入力データIの値と、各メモリセル1から出力される出力データEの値と、の間の関係は、図15(c)に真理値表を示す関係とされる。即ち各メモリセル1は、当該各メモリセル1の記憶値と入力データIの値との排他的否定論理和を当該各メモリセル1から出力データEとして出力する。また各メモリセル1の記憶値が上記所定値である場合、そのメモリセル1からは、入力データIがいずれの値であっても当該所定値が出力データEとして多数判定回路2に出力される。なお、各メモリセル1の細部構成については、後ほど図16(a)を用いて説明する。

30

【0156】

次に多数判定回路2は、各メモリセル1からの出力データEの値に基づき、値「1」の出力データEの数が値「0」の出力データEの数より大きい場合にのみ値「1」の出力データOを出力し、それ以外の場合に値「0」の出力データOを出力する。このとき、値「1」の出力データEの数が値「0」の出力データEの数より大きい場合以外の場合とは、具体的には、いずれかのメモリセル1から値「NC」が出力されている場合、又は各メモリセル1からの値「1」の出力データEの数が値「0」の出力データEの数以下の場合、のいずれかである。なお多数判定回路2及び各メモリセル1を含むニューラルネットワーク回路CSの細部構成については、後ほど図16(b)を用いて説明する。

【0157】

ここでニューラルネットワーク回路CSは上述したように、図15(a)においてハッチングで示されるニューロンNRにおける上記乗算処理、加算処理及び活性化処理をモデル化した回路である。そして、各メモリセル1からの上記排他的否定論理和としての出力データEの出力が上記重み付け係数Wを用いた上記乗算処理に相当する。また多数判定回路2は、値「1」の出力データEの数と値「0」の出力データEの数とを比較する比較処理の前提として、値「1」の出力データEの数を加算してその合計値を算出すると共に、値「0」の出力データEの数を加算してその合計値を算出する。これらの加算が上記加算処理に相当する。そして多数判定回路2において、値「1」の出力データE及び値「0」の出力データEそれぞれの数の上記合計値を比較し、前者の数から後者の数を減じた値が予め設定された多数判定閾値以上の場合にのみ、値「1」の出力データOを多数判定回路2から出力する。一方それ以外の場合、即ち値「1」の出力データEの数の合計値から値

40

50

「0」の出力データEの数の合計値を減じた値が多数判定閾値未満の場合に値「0」の出力データOを多数判定回路2から出力する。このとき、出力データEが上記所定値の場合に多数判定回路2は、当該出力データEを、値「1」の出力データEの数及び値「0」の出力データEの数のいずれにも加算しない。

【0158】

ここで、多数判定回路2における上記多数判定閾値を用いた処理について、より具体的に説明する。なお図15に例示するニューラルネットワーク回路CSでは、値「1」の出力データEの数と値「0」の出力データEの数との総数は「4」であるが、説明の明確化のため、当該総数が「10」である場合の上記処理について説明する。

【0159】

即ち、例えば多数判定閾値が「0」であり、値「1」の出力データEの数と値「0」の出力データEの数が共に「5」であるとする、値「1」の出力データEの数から値「0」の出力データEの数を減じた値は「0」であり、これは当該多数判定閾値と等しい。よってこの場合に多数判定回路2は、値「1」の出力データOを出力する。これに対して、多数判定閾値が「0」であり、値「1」の出力データEの数が「4」であり、値「0」の出力データEの数が「6」であるとする、値「1」の出力データEの数から値「0」の出力データEの数を減じた値は「-2」であり、これは当該多数判定閾値より小さい。よってこの場合に多数判定回路2は、値「0」の出力データOを出力する。

【0160】

他方、例えば多数判定閾値が「-2」であり、値「1」の出力データEの数と値「0」の出力データEの数が共に「5」であるとする、値「1」の出力データEの数から値「0」の出力データEの数を減じた値「0」は当該多数判定閾値より大きい。よってこの場合に多数判定回路2は、値「1」の出力データOを出力する。これに対して、多数判定閾値が「-2」であり、値「1」の出力データEの数が「4」であり、値「0」の出力データEの数が「6」であるとする、値「1」の出力データEの数から値「0」の出力データEの数を減じた値「-2」は当該多数判定閾値と等しい。よってこの場合も多数判定回路2は、値「1」の出力データOを出力する。

【0161】

以上具体的に説明した多数判定回路2における処理が上記活性化処理に相当する。以上の通り、図15(b)に示すニューラルネットワーク回路CSにより、図15(a)においてハッチングで示されるニューロンNRとしての各処理がモデル化される。

【0162】

次に、各メモリセル1の細部構成について、図16(a)を用いて説明する。図16(a)に示すように、各メモリセル1のそれぞれは、トランジスタ T_1 乃至トランジスタ T_4 と、インバータ IV_1 乃至インバータ IV_4 と、により構成されている。なお図16に示す各トランジスタ T_1 等のそれぞれは、例えばMOSFET(Metal Oxide semiconductor Field Effect Transistor)等により構成されている。そしてこれらの素子が、入力データ I_n に相当する接続線 LI_n 及び接続線/ LI_n 、ワード(Word)信号に相当する接続線 $W1$ 及び接続線 $W2$ 、並びにマッチ(match)信号に相当するマッチ線 M 及び反転マッチ線/ M により図16(a)に示す態様で接続されて、一つのメモリセル1が構成されている。このとき、トランジスタ T_1 及びトランジスタ T_2 、並びにインバータ IV_1 及びインバータ IV_2 により例えばSRAM(static random access memory)としての一のメモリ CL_1 が構成され、トランジスタ T_3 及びトランジスタ T_4 、並びにインバータ IV_3 及びインバータ IV_4 により例えばSRAMとしての一のメモリ CL_2 が構成される。また、トランジスタ T_5 乃至トランジスタ T_9 によりXNORゲート G_1 が構成され、トランジスタ T_{10} 乃至トランジスタ T_{14} によりXORゲート G_2 が構成される。

【0163】

次に、多数判定回路2及び各メモリセル1を含むニューラルネットワーク回路CSの細部構成について、図16(b)を用いて説明する。なお図16(b)は、図15(a)に対応して入力データIが四つである(即ち、メモリセル1を四つ備える)ニューラルネッ

10

20

30

40

50

トワーク回路CSの細部構成について示している。なお図16(b)に例示するニューラルネットワーク回路CSでは、上記多数判定閾値が「0」である場合について説明する。

【0164】

図16(b)に示すようにニューラルネットワーク回路CSは、四つのメモリセル1と、多数判定回路2を構成するトランジスタ T_{20} 乃至トランジスタ T_{30} (図16(b)破線参照)と、により構成されている。このとき、図16(b)において一点鎖線で示されるように、トランジスタ T_{25} 乃至トランジスタ T_{28} により、フリップフロップ型のセンスアンプSAが構成されている。そしてこれらの素子が、四つのメモリセル1に共通の上記接続線W1及び接続線W2並びに上記接続線M及び接続線/M、及び出力データOに相当する接続線LO及び接続線/LOにより図16(b)に示す態様で接続されて、一つのニューラルネットワーク回路CSが構成されている。また図16(b)に示すニューラルネットワーク回路CSには、当該ニューラルネットワーク回路CSとしての処理を規定するための予め設定されたタイミング信号 ϕ_1 、タイミング信号 ϕ_2 及びタイミング信号/ ϕ_2 並びにタイミング信号 ϕ_3 が外部から入力されている。このとき、タイミング信号 ϕ_1 はトランジスタ T_{20} 乃至トランジスタ T_{22} のゲート端子にそれぞれ入力され、タイミング信号 ϕ_2 及びタイミング信号/ ϕ_2 はトランジスタ T_{29} 及びトランジスタ T_{30} のゲート端子にそれぞれ入力され、タイミング信号 ϕ_3 はトランジスタ T_{23} 及びトランジスタ T_{24} のゲート端子にそれぞれ入力されている。以上の構成において、タイミング信号 ϕ_1 に基づいてそれぞれプリチャージされた各メモリセル1のマッチ線Mと反転マッチ線/Mとでは、入力データIの値並びにメモリ CL_1 及びメモリ CL_2 の記憶値に応じて、当該プリチャージされた電荷が引き抜かれるタイミングが異なる。そしてセンスアンプSAは、これらマッチ線M又は反転マッチ線/Mのどちらがより早く当該プリチャージされた電荷を引き抜かれるかを検出し、更に当該マッチ線Mと反転マッチ線/Mとの間の電圧差を増幅することにより、当該検知結果を接続線LO及び接続線/LOに出力する。ここで、接続線LOにおける値が「1」であることが、ニューラルネットワーク回路CSとしての出力データOの値「1」であることを意味することになる。以上の構成及び動作によりニューラルネットワーク回路CSは、上記タイミング信号 ϕ_1 等に基づいて、図15(a)においてハッチングで示されるニューロンNRとしての各処理をモデル化した処理を実行し、上記出力データOを出力する。

【0165】

(B) 実施形態に係るニューラルネットワーク集積回路の第1例について

次に、実施形態に係るニューラルネットワーク集積回路の第1例について、図17を用いて説明する。なお図17において、図15及び図16を用いて説明した実施形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

【0166】

以下の図17乃至図20を用いてそれぞれ説明する実施形態に係るニューラルネットワーク集積回路は、図15及び図16を用いて説明した実施形態に係るニューラルネットワーク回路を複数集積した集積回路である。そしてこれらのニューラルネットワーク集積回路は、より多くのニューロンNRからなる複雑なニューラルネットワークをモデル化するためのものである。

【0167】

まず、図17(a)に例示するニューラルネットワークS1をモデル化するための実施形態に係るニューラルネットワーク集積回路の第1例について説明する。当該ニューラルネットワークS1は、図17(a)においてハッチングで示されるm個のニューロンNRのそれぞれに対してn個のニューロンNRから1ビットの出力データOがそれぞれ出力されることにより、当該ハッチングで示されるニューロンNRから1ビットの出力データOがそれぞれ出力されるニューラルネットワークである。即ち上記ニューラルネットワークS1は、並列多入力・並列多出力型の一段ニューラルネットワークである。ここで図17(a)においては、各ニューロンNRの全てが入力信号I又は出力信号Oにより接続され

10

20

30

40

50

ている場合を示しているが、モデル化しようとする脳機能に応じて、各ニューロンNRのいずれかの間が接続されていなくてもよい。そしてこのことが、当該接続されていないニューロンNR間の接続に対応するメモリセル1の記憶値として上記所定値が記憶されていることにより表現される。なおこの点は、以降に図18(a)、図19(a)又は図20(a)を用いてそれぞれ説明するニューラルネットワークの場合においても同様である。

【0168】

上記ニューラルネットワークS1をモデル化する場合、図15及び図16を用いて説明した実施形態に係るニューラルネットワーク回路CSにおいて、1ビットの入力データIをn個とする。このとき、当該n個の入力データIが入力されるニューラルネットワーク回路CSのそれぞれが、図17(a)においてハッチングで示されるニューロンNRの機能をモデル化したものであり、上記乗算処理、加算処理及び活性化処理をそれぞれ実行する。なお以下の図17乃至図20を用いた説明においては、上記n個の入力データIが入力されるニューラルネットワーク回路CSを、「ニューラルネットワーク回路CS1」、「ニューラルネットワーク回路CS2」、...、と称する。そして実施形態に係るニューラルネットワーク集積回路の第1例としては、当該n個の入力データIが入力されるニューラルネットワーク回路CS1等をm個集積する。

10

【0169】

即ち図17(b)に示すように、実施形態に係るニューラルネットワーク集積回路の第1例であるニューラルネットワーク集積回路C1は、各1ビットのn個の入力データ I_1 乃至入力データ I_n がそれぞれ共通的に入力されるm個のニューラルネットワーク回路CS1乃至ニューラルネットワーク回路CSmが集積されて構成されている。そして、ニューラルネットワーク回路CS1乃至ニューラルネットワーク回路CSmのそれぞれには、上記タイミング信号 ϕ_1 等がタイミング生成回路TGから共通的に入力される。このときタイミング生成回路TGは、予め設定された基準クロック信号CLKに基づいて上記タイミング信号 ϕ_1 等を生成してニューラルネットワーク回路CS1乃至ニューラルネットワーク回路CSmに出力する。そしてニューラルネットワーク回路CS1乃至ニューラルネットワーク回路CSmの各々は、上記入力データ I_1 乃至入力データ I_n と、タイミング信号 ϕ_1 等と、に基づいて、各1ビットの出力データ O_1 、出力データ O_2 、...、出力データ O_m をそれぞれ出力する。

20

【0170】

以上説明した構成を備えるニューラルネットワーク集積回路C1において、m個のニューロンNRに対してn個のニューロンNRからそれぞれ出力データOが出力されることにより、m個のニューロンNRから出力データOが合計m個出力される図17(a)のニューラルネットワークS1がモデル化される。

30

【0171】

(C) 実施形態に係るニューラルネットワーク集積回路の第2例について

次に、実施形態に係るニューラルネットワーク集積回路の第2例について、図18を用いて説明する。なお図18において、図15及び図16を用いて説明した実施形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

40

【0172】

実施形態に係るニューラルネットワーク集積回路の第2例は、図18(a)に例示するニューラルネットワークSS1をモデル化するためのニューラルネットワーク集積回路である。当該ニューラルネットワークSS1は、図17(a)を用いて説明したニューラルネットワークS1において $n = m$ とした場合に相当する。即ちニューラルネットワークSS1は、図18(a)においてハッチングで示される $3 \times n$ 個のニューロンNRのそれぞれに対して相隣接する列の(n個の)ニューロンNRから出力データOがそれぞれ出力されることにより、図18(a)右端列のn個のニューロンNRから出力データOがそれぞれ出力されるニューラルネットワークである。上記ニューラルネットワークSS1は、並列多入力 - 並列多出力型の多段ニューラルネットワークである。

50

【 0 1 7 3 】

上記ニューラルネットワーク S S 1 をモデル化する場合も、図 1 7 を用いて説明したニューラルネットワーク S 1 と同様に、図 1 5 及び図 1 6 を用いて説明した実施形態に係るニューラルネットワーク回路 C S において、1 ビットの入力データ I を n 個とする。このとき、当該 n 個の入力データ I が入力されるニューラルネットワーク回路 C S のそれぞれが、図 1 8 (a) においてハッチングで示されるニューロン N R の機能をモデル化したものであり、上記乗算処理、加算処理及び活性化処理をそれぞれ実行する。そして実施形態に係るニューラルネットワーク集積回路の第 2 例としては、当該 n 個の入力データ I が入力されるニューラルネットワーク回路 C S 1 1 等を直列に接続して計 $3 \times n$ 個集積する。

【 0 1 7 4 】

即ち図 1 8 (b) に示すように、実施形態に係るニューラルネットワーク集積回路の第 2 例であるニューラルネットワーク集積回路 C C 1 は、各 1 ビットの n 個の入力データ I_1 乃至入力データ I_n がそれぞれ共通的に入力される n 個のニューラルネットワーク回路 C S 1 1 乃至ニューラルネットワーク回路 C S 1 n が集積されて一のニューラルネットワーク集積回路 C 1 が構成される (図 1 7 (b) 参照) 。そして、ニューラルネットワーク集積回路 C 1 を構成するニューラルネットワーク回路 C S 1 1 乃至ニューラルネットワーク回路 C S 1 n のそれぞれは、各 1 ビットの入力データ O_{11} 乃至入力データ O_{1n} を出力し、それらが次段の n 個のニューラルネットワーク回路 C S 2 1 乃至ニューラルネットワーク回路 C S 2 n に共通的に入力される。これらニューラルネットワーク回路 C S 2 1 乃至ニューラルネットワーク回路 C S 2 n により、他の一のニューラルネットワーク集積回路 C 2 が構成される。そして、ニューラルネットワーク集積回路 C 2 を構成するニューラルネットワーク回路 C S 2 1 乃至ニューラルネットワーク回路 C S 2 n のそれぞれは、各 1 ビットの入力データ O_{21} 乃至入力データ O_{2n} を出力し、それらが次段の n 個のニューラルネットワーク回路 C S 3 1 乃至ニューラルネットワーク回路 C S 3 n に共通的に入力される。これらニューラルネットワーク回路 C S 3 1 乃至ニューラルネットワーク回路 C S 3 n により、更に一のニューラルネットワーク集積回路 C 3 が構成される。ここで、各ニューラルネットワーク回路 C S 1 1 等に対しては、図 1 7 (a) に示す場合と同様に上記タイミング信号 ϕ_1 等が共通的に入力されているが、説明の簡略化のために図 1 8 (b) では図示を省略している。そしてニューラルネットワーク集積回路 C 1 は上記入力データ I_1 乃至入力データ I_n と、タイミング信号 ϕ_1 等と、に基づいて、出力データ O_{11} 、出力データ O_{12} 、...、出力データ O_{1n} をそれぞれ生成し、これらを次段のニューラルネットワーク集積回路 C 2 に共通的に出力する。次にニューラルネットワーク集積回路 C 2 は上記出力データ O_{12} 乃至出力データ O_{1n} と、タイミング信号 ϕ_1 等と、に基づいて、出力データ O_{21} 、出力データ O_{22} 、...、出力データ O_{2n} をそれぞれ生成し、これらを次段のニューラルネットワーク集積回路 C 3 に共通的に出力する。最後にニューラルネットワーク集積回路 C 3 は上記出力データ O_{21} 乃至出力データ O_{2n} と、タイミング信号 ϕ_1 等と、に基づいて、最終的な出力データ O_{31} 、出力データ O_{32} 、...、出力データ O_{3n} をそれぞれ生成して出力する。

【 0 1 7 5 】

以上説明した構成を備えるニューラルネットワーク集積回路 C C 1 において、n 個のニューロン N R から次段の n 個のニューロン N R に対してそれぞれが 1 ビットの出力データ O がそれぞれ出力されることが段階的に繰り返されることにより、最終的に出力データ O が合計 n 個出力される図 1 8 (a) のニューラルネットワーク S S 1 がモデル化される。

【 0 1 7 6 】

(D) 実施形態に係るニューラルネットワーク集積回路の第 3 例について

次に、実施形態に係るニューラルネットワーク集積回路の第 3 例について、図 1 9 を用いて説明する。なお図 1 9 において、図 1 5 及び図 1 6 を用いて説明した実施形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

【 0 1 7 7 】

10

20

30

40

50

実施形態に係るニューラルネットワーク集積回路の第3例は、図19(a)に例示するニューラルネットワークSS2をモデル化するためのニューラルネットワーク集積回路の例である。当該ニューラルネットワークSS2は、それぞれが図19(a)においてハッチングで示されるm個のニューロンNRからなる複数の組により構成され、これらのニューロンNRのそれぞれに対して共通のn個のニューロンNR(図19(a)において破線で示される)から1ビットの出力データOがそれぞれ出力されることにより、図19(a)においてハッチングで示される各ニューロンNRから各1ビットで合計m×組数個の出力データOが出力されるニューラルネットワークである。このニューラルネットワークSS2の場合、図19(a)においてハッチングで示される各ニューロンNRは、各1ビットで同じ数(n個)の出力データOをそれぞれ受信することになる。即ち上記ニューラルネットワークSS2は、並列多入力・並列多出力型の一段ニューラルネットワークである。

10

【0178】

上記ニューラルネットワークSS2をモデル化する場合も、図17を用いて説明したニューラルネットワークS1と同様に、図15及び図16を用いて説明した実施形態に係るニューラルネットワーク回路CSにおいて、1ビットの入力データIをn個とする。このとき、当該n個の入力データIが入力されるニューラルネットワーク回路CSのそれぞれが、図19(a)においてハッチングで示されるニューロンNRの機能をモデル化したものであり、上記乗算処理、加算処理及び活性化処理をそれぞれ実行する。そして実施形態に係るニューラルネットワーク集積回路の第3例としては、当該n個の入力データIが入力されるニューラルネットワーク回路CS11等を並列に接続して上記組数分集積する。

20

【0179】

即ち図19(b)に示すように、実施形態に係るニューラルネットワーク集積回路の第3例であるニューラルネットワーク集積回路CC2は、各1ビットのn個の入力データ I_1 乃至入力データ I_n がそれぞれ共通的に入力されるm個のニューラルネットワーク回路CS11乃至ニューラルネットワーク回路CS1mが集積されて、一のニューラルネットワーク集積回路C1が構成される(図17(b)参照)。また同じn個の入力データ I_1 乃至入力データ I_n がそれぞれ並列的且つ共通的に入力されるm個のニューラルネットワーク回路CS21乃至ニューラルネットワーク回路CS2mが集積されて、他の一のニューラルネットワーク集積回路C2が構成される(図17(b)参照)。これ以降、同様にn個の入力データ I_1 乃至入力データ I_n がそれぞれ並列的且つ共通的に入力されるm個のニューラルネットワーク回路が集積されて、図19(b)において図示を省略する他のニューラルネットワーク集積回路がそれぞれ構成される。ここで、各ニューラルネットワーク回路CS11等に対しては、図18を用いて説明した場合と同様に、図17(a)に示す場合と同様の上記タイミング信号 ϕ_1 等が共通的に入力されているが、説明の簡略化のために図19(b)では図示を省略している。そしてニューラルネットワーク集積回路C1は上記入力データ I_1 乃至入力データ I_n と、タイミング信号 ϕ_1 等と、に基づいて、各1ビットの出力データ O_{11} 、出力データ O_{12} 、...、出力データ O_{1m} をそれぞれ生成して出力する。一方ニューラルネットワーク集積回路C2は、同じ入力データ I_1 乃至入力データ I_n と、タイミング信号 ϕ_1 等と、に基づいて、各1ビットの出力データ O_{21} 、出力データ O_{22} 、...、出力データ O_{2m} をそれぞれ生成して出力する。これ以降、図示を省略する他のニューラルネットワーク集積回路も、それぞれにm個の出力データを出力する。

30

40

【0180】

以上説明した構成を備えるニューラルネットワーク集積回路CC2において、m×組数分のニューロンNRから並行的にそれぞれ出力データOが出力されることにより、最終的に出力データOが合計m×組数分出力される図19(a)のニューラルネットワークSS2がモデル化される。

【0181】

(E) 実施形態に係るニューラルネットワーク集積回路の第4例について

最後に、実施形態に係るニューラルネットワーク集積回路の第4例について、図20を

50

用いて説明する。なお図20において、図15及び図16を用いて説明した実施形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

【0182】

実施形態に係るニューラルネットワーク集積回路の第4例は、図20(a)に例示するニューラルネットワークSS3をモデル化するためのニューラルネットワーク集積回路の例である。当該ニューラルネットワークSS3は、これまで説明してきた実施形態に係るニューラルネットワークS1等に対して、ニューロンNRの数及びニューロンNR間の接続の態様についての自由度を更に向上させたニューラルネットワークである。なお図20(a)には、段階的に各1ビットの出力データO(入力データI)が授受されるニューロン群(図20(a)破線参照)に属するニューロンNRの数が相互に異なるニューラルネットワークSS3について例示している。

10

【0183】

上記ニューラルネットワークSS3をモデル化する場合、図2乃至図16を用いて説明した実施形態に係るニューラルネットワーク回路CSにおいて、1ビットの入力データIを例えばn個とする。このとき、当該n個の入力データIが入力されるニューラルネットワーク回路CSのそれぞれが、図20(a)に示す各ニューロンNRの機能をモデル化したものであり、上記乗算処理、加算処理及び活性化処理をそれぞれ実行する。そして実施形態に係るニューラルネットワーク集積回路の第4例としては、当該n個の入力データIが入力されるニューラルネットワーク回路CS11等をそれぞれに複数備えたニューラルネットワーク集積回路を複数備え、当該各ニューラルネットワーク集積回路を、後述する複数のスイッチ及びそれらを切り換えるスイッチボックスにより接続して集積する。

20

【0184】

即ち図20(b)に示すように、実施形態に係るニューラルネットワーク集積回路の第4例であるニューラルネットワーク集積回路CC3は、各1ビットのn個の入力データ I_1 乃至入力データ I_n がそれぞれ共通的に入力されるn個のニューラルネットワーク回路CS11乃至ニューラルネットワーク回路CS1nが集積されて一のニューラルネットワーク集積回路C1が構成される(図17(b)参照)。そして同様に、例えばm個のニューラルネットワーク回路CS21乃至ニューラルネットワーク回路CS2mが集積されて一のニューラルネットワーク集積回路C2が、ニューラルネットワーク回路CS31乃至ニューラルネットワーク回路CS3p(pは2以上の自然数。以下、同様。)が集積されて一のニューラルネットワーク集積回路C3が、ニューラルネットワーク回路CS41乃至ニューラルネットワーク回路CS4q(qは2以上の自然数。以下、同様。)が集積されて一のニューラルネットワーク集積回路C4が、それぞれ構成される。また、ニューラルネットワーク集積回路C1乃至ニューラルネットワーク集積回路C4のそれぞれは、図20(b)に例示するようにスイッチSW1乃至スイッチSW4を介して、相互に各1ビットの入力データI及び出力データOの授受が可能とされている。そして、ニューラルネットワーク集積回路C1乃至ニューラルネットワーク集積回路C4間における入力データI及び出力データOの授受の態様(即ちニューラルネットワーク集積回路C1乃至ニューラルネットワーク集積回路C4間の接続態様)が、スイッチSW1乃至スイッチSW4を介してスイッチボックスSB1乃至スイッチボックスSB4により切り換えられる。このとき、上記スイッチSW1乃至スイッチSW4及びスイッチボックスSB1乃至スイッチボックスSB4が、本発明に係る「スイッチ部」の一例に相当する。

30

40

【0185】

次に、上記スイッチボックスSB1乃至スイッチボックスSB4の細部構成について、図20(c)を用いて説明する。なお、スイッチボックスSB1乃至スイッチボックスSB4はいずれも同様の構成を備えるので、図20(c)では、これらを纏めてスイッチボックスSBとして説明する。

【0186】

図20(c)に示すように、ニューラルネットワーク集積回路CC3における各1ビット

50

トの入力データI又は出力データOの接続態様、及び結果的に有効なニューロンNRの数を制御するスイッチボックスSBは、セクタ M_1 乃至セクタ M_5 が図20(c)に示す態様により接続されて構成されている。この図20(c)に示すスイッチボックスSBの構成においては、上述してきた入力データIに対応する信号は図20(c)中左から入力される信号であり、上記出力データOに対応する信号は図20(c)中の上方及び下方からそれぞれ入力される信号である。そしてニューラルネットワーク集積回路C1乃至ニューラルネットワーク集積回路C4に対する入力データI等の切り換えは、当該切り換えを制御する切換制御信号 S_{c1} 乃至切換制御信号 S_{c5} が外部からそれぞれ入力されるセクタ M_1 乃至セクタ M_5 により、実行される。

【0187】

以上説明したように、入力データIに対応する出力データOを生成して出力する図20(a)のニューラルネットワークSS3が、図20(c)に示す構成を備えるスイッチボックスSB1乃至スイッチボックスSB4により各スイッチSW1乃至スイッチSW4が切り換えられる、図20(b)に示す構成のニューラルネットワーク集積回路CC3によりモデル化される。

【0188】

以上それぞれ説明したように、実施形態に係るニューラルネットワーク回路CS及びニューラルネットワーク集積回路C1等の構成及び動作によれば、図15及び図16に例示するように、対応すべき脳機能に基づいてその数が既定されているメモリセル1のそれぞれが、記憶値として「NC」を意味する所定値或いは「1」又は「0」のいずれかを記憶し、1ビットの入力データIの値と記憶値とが等しい場合に入力データIの入力に対応して「1」を出力し、入力データIの値と記憶値とが等しくない場合に入力データIの入力に対応して「0」を出力し、上記所定値が記憶されている場合にいずれの値の入力データIが入力されても当該所定値を出力する。そして多数判定回路2が、値「1」を出力するメモリセル1の総数が値「0」を出力するメモリセル1の総数より大きい場合に値「1」を出力データOとして出力し、値「1」を出力するメモリセル1の総数が値「0」を出力するメモリセル1の総数以下の場合に「0」を出力データOとして出力する。よって、メモリセル1においてニューラルネットワーク回路としての乗算処理を行い、一の多数判定回路2によりニューラルネットワーク回路としての加算処理及び活性化処理を行うので、回路規模及びそれに対応するコストを大幅に縮小しつつ、ニューラルネットワーク回路を効率的に実現することができる。

【0189】

また図17(b)に例示するように、それぞれが1ビットのn個の入力データIにそれぞれ対応してメモリセル1の数がnであるニューラルネットワーク回路CSをm個備え、各ニューラルネットワーク回路CSに対してn個の入力データIが並列且つ共通的に入力され、各ニューラルネットワーク回路CSから出力データOをそれぞれ出力させる場合は、図17(a)に例示するニューラルネットワークS1をモデル化し且つ入力がn個であり出力がm個である $n \times m$ のニューラルネットワーク集積回路C1を、回路規模及びそれに対応するコストを大幅に縮小しつつ効率的に実現することができる。更にこの場合に、図17(a)においてハッチングで示されるm個のニューロンNRと、当該m個に対して出力データOをそれぞれ出力するn個のニューロンNRと、の間に多様な接続のパターンがあっても、ニューラルネットワーク集積回路C1においてニューロンNR間に接続がない場合に対応するメモリセル1の記憶値として上記所定値を用いることで、ニューラルネットワーク集積回路C1をより効率的に実現することができる。更に図17に例示する場合には、n個の入力データIを各ニューラルネットワーク回路CSに対して並列且つ共通的に入力し、それに基づくm個の出力データOを並列に出力させることができるため、入力データI及び出力データOを逐次的に入出力させなければならない場合と比べて処理の大幅な高速化が可能となる。

【0190】

更にまた図18に例示するように、上記「n」と上記「m」とが等しいニューラルネッ

10

20

30

40

50

トワーク集積回路 C 1 等を直列に接続し、一のニューラルネットワーク集積回路 C 1 (又はニューラルネットワーク集積回路 C 2) からの出力データ O を、当該ニューラルネットワーク集積回路 C 1 (又はニューラルネットワーク集積回路 C 2) の直後に接続され他のニューラルネットワーク集積回路 C 2 (又はニューラルネットワーク集積回路 C 3) における入力データ I とする場合は、入力及び出力が共に並列型であるニューラルネットワーク集積回路 C C 1 を、回路規模及びそれに対応するコストを大幅に縮小しつつ効率的に実現することができる。

【0191】

また図 19 に例示するように、各ニューラルネットワーク集積回路 C S に対して n 個の入力データ I が並列且つ共通に入力され、各ニューラルネットワーク集積回路 C S から m 個の出力データ O を並列にそれぞれ出力する場合は、入力及び出力が共に並列型であり且つ出力データ O の数が入力データ I の数より多いニューラルネットワーク集積回路 C C 2 を、回路規模及びそれに対応するコストを大幅に縮小しつつ効率的に実現することができる。

10

【0192】

また図 20 に例示するように、ニューラルネットワーク集積回路 C 1 等を複数備え、各ニューラルネットワーク集積回路 C 1 等をアレイ状に且つ相互に接続するスイッチ S W 1 等により、各ニューラルネットワーク集積回路 C 1 等に対する入力データ I 及び出力データ O が切り換える場合には、スイッチ S W 1 等における切換動作が、対応すべき脳機能に基づいて既定されていれば、大規模なニューラルネットワーク集積回路 C C 3 を、対応するコストを大幅に縮小しつつ効率的に実現することができる。

20

【0193】

(II) 関連形態

次に、本発明に関連する関連形態について、図 21 乃至図 27 を用いて説明する。なお、図 21 及び図 22 は関連形態に係るニューラルネットワーク集積回路の第 1 例をそれぞれ示す図であり、図 23 は関連形態に係るニューラルネットワーク回路の第 1 例を示す図であり、図 24 は関連形態に係るニューラルネットワーク集積回路の第 2 例を示す図である。また、図 25 は当該ニューラルネットワーク集積回路の第 3 例を示す図であり、図 26 は当該ニューラルネットワーク集積回路の第 4 例を示す図であり、図 27 は当該第 4 例の細部構成を示す図である。

30

【0194】

以下に説明する関連形態は、図 1 および図 15 乃至図 20 を用いて上述してきたニューラルネットワーク S 等のモデル化の構成又は手法とは異なる構成又は手法により、当該ニューラルネットワーク S 等をモデル化しようとするものである。

【0195】

(A) 関連形態に係るニューラルネットワーク集積回路の第 1 例について

初めに、関連形態に係るニューラルネットワーク集積回路の第 1 例について、図 21 及び図 22 を用いて説明する。なお、図 21 は当該第 1 例としての上記乗算処理を行う当該第 1 例の一部を示す図であり、図 22 は当該第 1 例全体を示す図である。

【0196】

図 21 (a) に例示するように、当該第 1 例の一部によりモデル化されるネットワーク S' では、一つのニューロン N R から 1 ビットの出力データ O (換言すれば入力データ I) が入力される。そして、入力データ I の出力先となる図示しない複数の他のニューロンにそれぞれ対応する異なった重み付け係数 W_1 乃至重み付け係数 W_4 のうちの 하나가当該入力データ I に乗算され、上記図示しない他のニューロンに対して、出力データ E_1 乃至出力データ E_4 としてそれぞれ出力される。そしてこのときの出力データ E は、入力データ I と同様に 1 ビットの信号である。よって、図 21 に示す入力データ I の値、各重み付け係数 W の値及び出力データ E の値は、いずれも「0」又は「1」のいずれかである。

40

【0197】

次に、関連形態に係るニューラルネットワーク集積回路の第 1 例における図 21 (a)

50

に示すネットワーク S' に相当する部分の構成を、図 2 1 (b) にネットワーク回路 CS' として示す。当該ネットワーク回路 CS' は、各々が図 2 1 (a) に示す出力データ E_1 乃至出力データ E_4 にそれぞれ対応する四組のメモリセル 1 0 及びメモリセル 1 1 (接続有無情報用のメモリセル) と、各々が出力データ E (換言すれば、上記図示しない他のニューロンの入力データ I) にそれぞれ対応する四つの多数判定入力回路 1 2 と、により構成される。このとき、一つのメモリセル 1 0 と一つのメモリセル 1 1 からなるメモリセル対の数及びそれらに対応する多数判定入力回路 1 2 の数 (図 2 1 に例示する場合は共に四つ) は、関連形態に係るニューラルネットワーク集積回路の第 1 例として所望される出力データ O の数に等しい。なお以下の図 2 1 の説明において、出力データ O の数分の上記メモリセル対を、纏めて「メモリセルブロック 1 5」(図 2 1 (b) 破線参照) と称する。

10

【 0 1 9 8 】

以上の構成において各メモリセルブロック 1 5 内のメモリセル 1 0 は、ネットワーク回路 CS' を含む関連形態に係るニューラルネットワーク集積回路の第 1 例が対応すべき脳機能に基づいて予め設定された 1 ビットの重み付け係数 W をそれぞれに記憶している。これに対して各メモリセルブロック 1 5 内のメモリセル 1 1 は、上記脳機能に基づいて予め設定された 1 ビットの接続有無情報をそれぞれに記憶している。ここで当該接続有無情報は、上記実施形態におけるメモリセル 1 の記憶値「 NC 」に相当するものであり、関連形態に係るニューラルネットワークにおける二つのニューロン NR 間に接続がある状態であるか、又は当該接続がない状態であるか、のいずれかを表すための記憶値である。なお、各メモリセル 1 0 及びメモリセル 1 1 にどのような記憶値を記憶させておくかは、例えば

20

【 0 1 9 9 】

そして各メモリセル 1 0 は、それぞれの記憶値を、重み付け係数 W_1 、重み付け係数 W_2 、重み付け係数 W_3 及び重み付け係数 W_4 として多数判定入力回路 1 2 にそれぞれ出力する。このとき各メモリセル 1 0 は、それぞれの記憶値を、重み付け係数 W_1 乃至重み付け係数 W_4 として同時に多数判定入力回路 1 2 に出力する。なおこの同時出力の構成は、以下の図 2 2 乃至図 2 7 を用いてそれぞれ説明するニューラルネットワーク回路及びニューラルネットワーク集積回路における各メモリセル 1 0 においても同様である。一方各メモリセル 1 1 も、それぞれの記憶値を、接続有無情報 C_1 、接続有無情報 C_2 、接続有無情報 C_3 及び接続有無情報 C_4 として多数判定入力回路 1 2 にそれぞれ出力する。このとき各メモリセル 1 1 は、それぞれの記憶値を、接続有無情報 C_1 乃至接続有無情報 C_4 として同時に多数判定入力回路 1 2 に出力する。また各メモリセル 1 1 は、上記各メモリセル 1 0 からの記憶値の出力に対して例えば一サイクル前又は後にずらして、それぞれの記憶値を同時に多数判定入力回路 1 2 に出力する。なおこの同時出力の構成及び各メモリセル 1 0 からの記憶値の出力のタイミングと関係は、以下の図 2 2 乃至図 2 7 を用いてそれぞれ説明するニューラルネットワーク回路及びニューラルネットワーク集積回路における各メモリセル 1 1 においても同様である。更に以下の説明において、接続有無情報 C_1 、接続有無情報 C_2 、接続有無情報 C_3 、...、に共通の事項を説明する場合、単に「接続有無情報 C 」と称する。

30

40

【 0 2 0 0 】

他方各多数判定入力回路 1 2 には、図 2 1 (b) において図示しない他のノード NR (図 2 1 (a) 参照) からの 1 ビットの入力データ I が共通的に入力されている。そして各多数判定入力回路 1 2 は、対応するメモリセル 1 1 から出力される上記接続有無情報を、そのまま接続有無情報 C_1 乃至接続有無情報 C_4 としてそれぞれ出力する。

【 0 2 0 1 】

これらに加えて各多数判定入力回路 1 2 は、対応するメモリセル 1 0 から出力される重み付け係数 W_1 、重み付け係数 W_2 、重み付け係数 W_3 及び重み付け係数 W_4 と上記入力データ I との排他的論理和 ($XNOR$) を算出し、上記出力データ E_1 、出力データ E_2 、出力データ E_3 及び出力データ E_4 としてそれぞれ出力する。このとき、対応するメモリセル 1

50

1の記憶値(重み付け係数 W)と、入力データ I の値と、多数判定入力回路12から出力される出力データ E の値と、の間の関係は、図21(c)に真理値表を示す関係となる。なお図21(c)には、上記排他的否定論理和(XNOR)を算出する前提としての排他的論理和(XOR)についても記載している。

【0202】

ここで、図15を用いて説明した実施形態に係るニューラルネットワーク回路CSに対応する真理値表(図15(c)参照)と、上記図21(c)に示す真理値表とを比較する。このとき、メモリセル10における記憶値及び入力データ I の値をそれぞれ図15(c)に示した真理値表のものと同じとした場合、図21(b)に示す出力データ E の値は図2(b)に示す出力データ E の値と同一となる。これらにより、図21(b)に示すネットワーク回路CS'は、図15(b)に示したニューラルネットワーク回路CSにおける上記乗算処理と同様の論理により、図21(a)に示すネットワークS'における上記乗算処理をモデル化した回路となる。即ち、各メモリセル10から出力されてくる各々の記憶値(重み付け係数 W)と入力データ I の値との間の排他的論理和を多数判定入力回路12において算出することが上記乗算処理に相当する。以上の通り、図21(b)に示すネットワーク回路CS'により、図21(a)に示すネットワークS'における乗算処理がモデル化される。

10

【0203】

(B) 関連形態に係るニューラルネットワーク集積回路の第1例について

次に、関連形態に係るニューラルネットワーク集積回路の第1例について、図21及び図22を用いて説明する。なお図22において、図21を用いて説明した関連形態に係るネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

20

【0204】

図22を用いて説明する関連形態に係るニューラルネットワーク集積回路の第1例は、図21を用いて説明した関連形態に係るネットワーク回路CS'を複数集積した集積回路である。関連形態に係るニューラルネットワーク集積回路の第1例では、当該ネットワーク回路CS'に相当する上記乗算処理に加えて、上記加算処理及び上記活性化処理が実行される。

【0205】

先ず、関連形態に係るニューラルネットワーク集積回路の第1例によりモデル化されるニューラルネットワークの全体について、図22(a)を用いて説明する。図22(a)に示す当該ニューラルネットワークS1'は、図21を用いて説明したネットワークS'を m 個のニューロンNR分含んでいる。当該ニューラルネットワークS1'では、図22(a)においてハッチングで示される n 個のニューロンNRのそれぞれに対して、それぞれが上記ネットワークS'を構成する m 個のニューロンNRからそれぞれ1ビットの出力データ O (換言すれば入力データ I)が出力される。そしてこれにより、各出力データ O が出力データ E となって上記ハッチングで示される n 個のニューロンNRのそれぞれに入力され、当該ハッチングで示されるニューロンNRから出力データ O が一つずつ合計 n 個並列的に出力される。即ち上記ニューラルネットワークS1'は、直列(m)入力-並列(n)出力型の一段ニューラルネットワークである。

30

40

【0206】

上記ニューラルネットワークS1'をモデル化した関連形態に係るニューラルネットワーク集積回路の第1例は、図22(b)に示すニューラルネットワーク集積回路C1'となる。当該ニューラルネットワーク集積回路C1'は、上記メモリセル対及び多数判定入力回路12をそれぞれに n 個ずつ含む関連形態に係るニューラルネットワーク回路CS'(図21参照)を m 個備えると共に、各多数判定入力回路12及びメモリセル対に対応させて n 個の直列多数判定回路13を備えている。そして図22(b)に示すように、 $n \times m$ 個の上記メモリセル対(換言すれば、 m 個のメモリセルブロック15)により、メモリセルアレイMC1が構成されている。またニューラルネットワーク集積回路C1'では、

50

図 2 2 (b) に示すメモリセルアレイ M C 1 における横一行 (m 個) のメモリセル対で一つの多数判定入力回路 1 2 を共用する。なお、メモリセルアレイ M C 1 、各多数判定入力回路 1 2 及び各直列多数判定回路 1 3 には、それぞれ上記タイミング信号 T_1 等が共通的に入力されているが、説明の簡略化のために図 2 2 (b) では図示を省略している。

【 0 2 0 7 】

以上の構成において、それぞれのニューラルネットワーク回路 C S ' を構成するメモリセルブロック 1 5 のメモリセル 1 0 からは、上記重み付け係数 W が、一のメモリセルブロック 1 5 に含まれる各メモリセル 1 0 について同時に、且つ m 個のメモリセルブロック 1 5 について順次に (即ちシリアル形式で) 、それぞれ出力される。そして、これらの重み付け係数 W と、対応するタイミングによりシリアル形式で入力される m 個の入力データ I (各 1 ビットの入力データ I) と、の上記排他論理和を、上記共有する多数判定入力回路 1 2 において時分割的に演算し、それを出力データ E として、対応する直列多数判定回路 1 3 にシリアル形式で出力する。一方、それぞれのニューラルネットワーク回路 C S ' を構成するメモリセルブロック 1 5 のメモリセル 1 1 からは、上記接続有無情報 C が、一のメモリセルブロック 1 5 に含まれる各メモリセル 1 1 について同時に、且つ m 個のメモリセルブロック 1 5 について順次に (即ちシリアル形式で) 、それぞれ出力される。そしてこれらの接続有無情報 C が、上記共有する多数判定入力回路 1 2 を介し、各入力データ I の入力タイミングに対応したシリアル形式で、対応する直列多数判定回路 1 3 に出力される。なお、各メモリセルブロック 1 5 からの上記各重み付け係数 W の出力タイミングの態様、及び各メモリセルブロック 1 5 からの上記各接続有無情報 C の出力タイミングの態様のそれぞれは、以下の図 2 3 乃至図 2 7 を用いてそれぞれ説明するニューラルネットワーク集積回路における各メモリセル 1 1 においても同様である。

10

20

【 0 2 0 8 】

次に、各多数判定入力回路 1 2 から出力データ E 及び接続有無情報 C がそれぞれ入力される n 個の直列多数判定回路 1 3 はそれぞれ、同じタイミングで入力された接続有無情報 C が「接続あり」を示している最大 m 個の出力データ E について、値「 1 」の当該出力データ E の数を加算してその合計値を算出すると共に、値「 0 」の出力データ E の数を加算してその合計値を算出する。これらの加算が上記加算処理に相当する。そして各直列多数判定回路 1 3 はそれぞれ、値「 1 」の出力データ E 及び値「 0 」の出力データ E それぞれの数の上記合計値を比較し、前者の数から後者の数を減じた値が、実施形態に係る上記多数判定閾値と同様にして予め設定された多数判定閾値以上の場合にのみ、値「 1 」の出力データ O を出力する。一方それ以外の場合、即ち値「 1 」の出力データ E の数の合計値から値「 0 」の出力データ E の数の合計値を減じた値が上記多数判定閾値未満の場合に各直列多数判定回路 1 3 はそれぞれ、値「 0 」の出力データ O を出力する。これらの各直列多数判定回路 1 3 における処理が上記活性化処理に相当すると共に、各出力データ O は 1 ビットとなる。ここで、同じタイミングで出力される上記接続有無情報 C が「接続なし」を示している場合、直列多数判定回路 1 3 は、出力データ E を、値「 1 」の出力データ E の数及び値「 0 」の出力データ E の数のいずれにも加算しない。そして各直列多数判定回路 1 3 は、上述した各処理により 1 ビットの出力データ O を出力することを、入力データ I が入力されるタイミングに合わせて繰り返す。このときの出力データ O は、結果的に各直列多数判定回路 1 3 から並列的に出力される。この場合、出力データ O の総数は n 個となる。以上の通り、図 2 2 (a) においてハッチングで示される一つのニューロン N R に対応する上記乗算処理、加算処理及び活性化処理のそれぞれは、図 2 2 (b) に示すメモリセルアレイ M C 1 における横一行分のメモリセル対と、それらに対応する多数判定入力回路 1 2 及び直列多数判定回路 1 3 と、により実行されることになる。

30

40

【 0 2 0 9 】

以上説明したように、図 2 2 (a) においてハッチングで示される n 個のニューロン N R に対して m 個のニューロン N R から 1 ビットの出力データ O がそれぞれ出力され、これらにより当該 n 個のニューロン N R から出力データ O が合計 n 個出力されるニューラルネットワーク S 1 ' が、図 2 2 (b) に示す構成を備えるニューラルネットワーク集積回路

50

C 1' によりモデル化される。

【0210】

(C) 関連形態に係るニューラルネットワーク回路の第1例について

次に、関連形態に係るニューラルネットワーク回路の第1例について、図23を用いて説明する。

【0211】

図23(a)に例示するように、当該第1例に相当するニューラルネットワークSは、基本的には図2(a)に例示した実施形態に係るニューラルネットワークSと同一の構成である。但し図23(a)に示す例では、図23(a)においてハッチングで示される一つのニューロンNRに対して他の三つのニューロンNRから1ビットの入力データI(当該他のニューロンNRから見ると出力データO)が並列的に入力され、それに対応する出力データOが当該ニューロンNRからシリアル形式で一つ出力される構成となっている。このときの出力データOも、各入力データIと同様に1ビットの信号である。よって、図23に示す入力データIの値及び出力データOの値は、いずれも「0」又は「1」のいずれかである。そして、図23(a)に示すニューロンNRにおいて実行される上記乗算処理等に相当する上記式(1)は、上記式(1)において $n=3$ とした場合の式である。即ち上記ニューラルネットワークSは、並列入力-直列出力型の一段ニューラルネットワークである。

10

【0212】

次に、図23(a)においてハッチングで示したニューロンNRに相当する関連形態に係るニューラルネットワーク回路の第1例の構成を、図23(b)にニューラルネットワーク回路CCS'として示す。当該ニューロンNRに相当する関連形態に係るニューラルネットワーク回路CCS'は、各々が図23(a)に例示する入力データIにそれぞれ対応する三組のメモリセル10及びメモリセル11と、各入力データIが入力される並列多数判定回路20と、により構成される。このとき、一つのメモリセル10と一つのメモリセル11からなるメモリセル対の数及びそれらに対応する多数判定入力回路12の数(図23に例示する場合は共に三つ)は、図23(a)に示すニューラルネットワークSとして所望される入力データIの数に等しい。なお以下の図23の説明において、入力データIの数分の上記メモリセル対はそれぞれメモリセルブロック15と示されている(図23(b)破線参照)。

20

30

【0213】

以上の構成において各メモリセルブロック15内のメモリセル10は、ニューラルネットワーク回路CCS'が対応すべき脳機能に基づいて予め設定された1ビットの重み付け係数Wをそれぞれに記憶している。これに対して各メモリセルブロック15内のメモリセル11は、上記脳機能に基づいて予め設定された1ビットの接続有無情報をそれぞれに記憶している。ここで当該接続有無情報は、図21及び図22を用いて説明した関連形態に係るニューラルネットワーク回路の第1例における接続有無情報 C_n と同一の情報であるので、細部の説明は省略する。また、各メモリセル10及びメモリセル11にどのような記憶値を記憶させておくかは、例えば、図23(a)に示すニューラルネットワークSとしてどのような脳機能をモデル化するか等に基づいて予め設定されていけばよい。

40

【0214】

そして各メモリセル10は、それぞれの記憶値を、図21(b)に示す各メモリセル10と同様のタイミングで、重み付け係数 W_1 、重み付け係数 W_2 及び重み付け係数 W_3 として並列多数判定回路20にそれぞれ出力する。一方各メモリセル11も、それぞれの記憶値である接続有無情報Cを、図21(b)に示す各メモリセル11と同様のタイミングで並列多数判定回路20にそれぞれ出力する。

【0215】

他方並列多数判定回路20には、上述したように入力データ I_1 、入力データ I_2 及び入力データ I_3 (各1ビット)が並列的に入力されている。そして並列多数判定回路20は、図22を用いて説明した一組の多数判定入力回路12及び直列多数判定回路13と同様

50

の動作を含む動作（即ち、上記乗算処理、上記加算処理及び上記活性化処理）を行う。具体的に並列多数判定回路20は先ず、対応する接続有無情報Cが「接続あり」を示している場合に、それぞれが1ビットの各入力データIと、それらに対応する重み付け係数Wと、の上記排他論理和を当該各入力データIについて演算する。次に並列多数判定回路20は、各上記演算結果について、値「1」の当該演算結果の数を加算してその合計値を算出すると共に、値「0」の当該演算結果の数を加算してその合計値を算出する。そして並列多数判定回路20は、値「1」の当該演算結果及び値「0」の当該演算結果それぞれの数の上記合計値を比較し、前者の数から後者の数を減じた値が、実施形態に係る上記多数判定閾値と同様にして予め設定された多数判定閾値以上の場合にのみ、値「1」の出力データOをシリアル形式で出力する。一方それ以外の場合、即ち値「1」の出力データEの数の合計値から値「0」の出力データEの数の合計値を減じた値が上記多数判定閾値未満の場合に並列多数判定回路20は値「0」の出力データOをシリアル形式で出力する。この場合出力データOは1ビットとなる。ここで、対応する接続有無情報Cが「接続なし」を示している場合、並列多数判定回路20は上記排他的論理和を演算しない。なお、各入力データIと、対応する重み付け係数Wと、の上記排他論理和を全ての入力データIについて一旦演算し、対応する接続有無情報Cが「接続なし」を示している場合に、その演算結果を値「1」の演算結果の数及び値「0」の演算結果の数のいずれにも加算しないように構成してもよい。そして並列多数判定回路20は、上述した各処理により1ビットの出力データOをシリアル形式で出力することを、並列的に入力されている各入力データIの数ごとに繰り返す。以上の各処理により、図23(b)に示すニューラルネットワーク回路CCS'は、図23(a)においてハッチングで示されるニューロンNRにおける上記乗算処理、加算処理及び活性化処理をモデル化した回路となる。

【0216】

(D) 関連形態に係るニューラルネットワーク集積回路の第2例について

次に、関連形態に係るニューラルネットワーク集積回路の第2例について、図24を用いて説明する。なお図24において、図23を用いて説明した関連形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

【0217】

図24を用いて説明する関連形態に係るニューラルネットワーク集積回路の第2例は、図23を用いて説明した関連形態に係るニューラルネットワーク回路CCS'を複数集積した集積回路であり、より多くのニューロンNRからなる複雑なニューラルネットワークをモデル化するためのものである。

【0218】

先ず、関連形態に係るニューラルネットワーク集積回路の第2例によりモデル化されるニューラルネットワークについて、図24(a)を用いて説明する。図24(a)に示す当該ニューラルネットワークS2'は、図24(a)においてハッチングで示されるm個のニューロンNRのそれぞれに対してn個のニューロンNRから1ビットの出力データO(m個のニューロンNRから見た場合は入力データI)が並列的にそれぞれ入力され、それらに対応する出力データOが当該ニューロンNRからシリアル形式で出力される構成となっている。このときの出力データOも、各入力データIと同様に1ビットの信号である。よって、図24に示す入力データIの値及び出力データOの値は、いずれも「0」又は「1」のいずれかである。即ち上記ニューラルネットワークS2'は、並列入力-直列出力型の一段ニューラルネットワークである。

【0219】

上記ニューラルネットワークS2'をモデル化した関連形態に係るニューラルネットワーク集積回路の第2例は、図24(b)に示すニューラルネットワーク集積回路C2'となる。当該ニューラルネットワーク集積回路C2'は、上記メモリセル対をそれぞれにn個ずつ含む関連形態に係るニューラルネットワーク回路CCS'(図23参照)をm個備えると共に、上記並列多数判定回路20を備えている。そして図24(b)に示すように

、 $n \times m$ 個の上記メモリセル対（換言すれば、 m 個のメモリセルブロック15）により、メモリセルアレイMC2が構成されている。またニューラルネットワーク集積回路C2'では、一つの並列多数判定回路20を図24（b）に示すメモリセルアレイMC2における横一行（ m 個）のメモリセル対で共用する。なお、メモリセルアレイMC2及び並列多数判定回路20には、それぞれ上記タイミング信号 ϕ_1 等が共通的に入力されているが、説明の簡略化のために図24（b）では図示を省略している。

【0220】

以上の構成において、それぞれのニューラルネットワーク回路CCS'を構成するメモリセルブロック15のメモリセル10からは、上記重み付け係数 W が、図22（b）に示す各メモリセル10及び各メモリセルブロック15と同様のタイミングで並列多数判定回路20に出力される。一方、それぞれのニューラルネットワーク回路CCS'を構成するメモリセルブロック15のメモリセル11からは、上記接続有無情報 C が、図22（b）に示す各メモリセル11及び各メモリセルブロック15と同様のタイミングで並列多数判定回路20に出力される。

10

【0221】

そして並列多数判定回路20は、メモリセルアレイMC2から出力されてくる重み付け係数 W 及び接続有無情報 C と、それらに対応する入力データ I と、に基づき、接続有無情報 C が「接続あり」を示している重み付け係数 W 及び入力データ I を用いた上記排他的論理和の演算処理、その演算結果に基づく値「1」の演算結果及び値「0」の演算結果それぞれの数の加算処理、その加算結果に基づく上記合計数の比較処理（図23（b）参照）、及びその比較結果に基づく出力データ O の生成処理を、メモリセルアレイMC2における横一行（ m 個）についてそれぞれ行う。また並列多数判定回路20は、上記横一行についての演算処理、加算処理、比較処理及び生成処理を、各入力データ I について、メモリセルブロック15ごとにシリアル形式で実行し、それぞれの実行結果としての出力データ O をシリアル形式で出力する。ここで、対応する接続有無情報 C が「接続なし」を示している場合、並列多数判定回路20は上記演算処理、加算処理、比較処理及び生成処理を行わない。

20

【0222】

以上説明したように、図24（a）においてハッチングで示される m 個のニューロンNRに対して n 個のニューロンNRからそれぞれ出力データ O が出力され、これらにより当該 m 個のニューロンNRから1ビットの出力データ O がシリアル形式で出力されるニューラルネットワークS2'が、図24（b）に示す構成を備えるニューラルネットワーク集積回路C2'によりモデル化される。

30

【0223】

（E）関連形態に係るニューラルネットワーク集積回路の第3例について

次に、関連形態に係るニューラルネットワーク集積回路の第3例について、図25を用いて説明する。なお図25において、図21及び図23を用いてそれぞれ説明した関連形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

40

【0224】

図25を用いて説明する関連形態に係るニューラルネットワーク集積回路の第3例は、図22を用いて説明した関連形態に係るニューラルネットワーク集積回路C1'と、図24を用いて説明した関連形態に係るニューラルネットワーク集積回路C2'と、を組み合わせた集積回路である。ここで、上記ニューラルネットワーク集積回路C1'は上述した通り直列入力・並列出力型の一段ニューラルネットワークS1'をモデル化したニューラルネットワーク回路である。一方上記ニューラルネットワーク集積回路C2'は、上述した通り並列入力・直列出力型の一段ニューラルネットワークS2'をモデル化したニューラルネットワーク回路である。そしてこれらを組み合わせた関連形態に係るニューラルネットワーク集積回路の第3例は、全体として直列入力・並列処理・直列出力型の多段ニューラルネットワークをモデル化したニューラルネットワーク集積回路であり、更に多くの

50

ニューロンNRからなる複雑なニューラルネットワークをモデル化するためのものである。

【0225】

先ず、関連形態に係るニューラルネットワーク集積回路の第3例によりモデル化されるニューラルネットワークについて、図25(a)を用いて説明する。図25(a)に示す当該ニューラルネットワークS1-2は、図25(a)において45度のハッチングで示されるn個のニューロンNRのそれぞれに対してm個のニューロンNRからそれぞれ1ビットの出力データOがシリアル形式で出力され、当該45度のハッチングで示されるニューロンNRと図24(a)において135度のハッチングで示されるm個のニューロンNRとの間で出力データO及び入力データIの授受が行われ、結果的に135度のハッチングで示されるm個のニューロンNRからそれぞれ出力データOがシリアル形式で出力されるニューラルネットワークである。なお上記ニューラルネットワークS1-2は、全体として、図17を用いて説明したニューラルネットワークS1を複数並べたニューラルネットワークに相当する。

10

【0226】

上記ニューラルネットワークS1-2をモデル化した関連形態に係るニューラルネットワーク集積回路の第3例は、図25(b)に示すニューラルネットワーク集積回路C1-2となる。当該ニューラルネットワーク集積回路C1-2は、図22を用いて説明したニューラルネットワーク集積回路C1'の各出力データO(並列的に出力される各出力データO)を、図24を用いて説明したニューラルネットワーク集積回路C2'における並列多数判定回路20への入力データ(即ち図24(b)に示す入力データI)とし、これにより、当該並列多数判定回路20から上記出力データOをシリアル形式で出力する構成を備える。このように、上記ニューラルネットワーク集積回路C1'とニューラルネットワーク集積回路C2'とを組み合わせることにより、結果的に、図22(a)に例示するニューラルネットワークS1'と、図24(a)に例示するニューラルネットワークS2'と、が組み合わせられた上記ニューラルネットワークS1-2がモデル化される。なお、ニューラルネットワークS1-2に含まれる上記ニューラルネットワーク集積回路C1'及びニューラルネットワーク集積回路C2'それぞれの動作は、図22及び図24を用いてそれぞれ説明した動作と同様となる。なお図25(b)に示すニューラルネットワーク集積回路C1-2では、並列多数判定回路20に対応する直列多数判定回路16が、破線で示される一組の多数判定入力回路12及び直列多数判定回路13により、それぞれ構成されていることになる。

20

30

【0227】

以上説明したように、図25(a)に示すニューラルネットワークS1-2が、図25(b)に示す直列入力-並列処理-直列出力型の構成を備えるニューラルネットワーク集積回路C1-2によりモデル化される。

【0228】

(F) 関連形態に係るニューラルネットワーク集積回路の第4例について

次に、関連形態に係るニューラルネットワーク集積回路の第4例について、図26及び図27を用いて説明する。なお図26及び図27において、図22及び図24並びに図25を用いてそれぞれ説明した関連形態に係るニューラルネットワーク回路と同様の構成部材については、同様の部材番号を付して細部の説明を省略する。

40

【0229】

図26を用いて説明する関連形態に係るニューラルネットワーク集積回路の第4例は、図26(a)に示すように、図25を用いて説明した関連形態に係るニューラルネットワーク集積回路C1-2において、それを構成する上記ニューラルネットワーク集積回路C1'と上記ニューラルネットワーク集積回路C2'との間にパイプラインレジスタ21を介在させた構成を備えたニューラルネットワーク集積回路C1-3である。このときパイプラインレジスタ21は、メモリセルアレイMC1のビット幅に相当する数のデータを一時的に記憶すると共に、外部からのイネーブル信号ENにより、その出力動作が制御され

50

る。このイネーブル信号ENは、予め設定された基準クロック信号のうちの偶数基準クロックに相当するタイミング信号である。そしてニューラルネットワーク集積回路C1-3は、図26(b)に示すように全体として、ニューラルネットワーク集積回路C1'におけるメモリセルアレイMC1と、ニューラルネットワーク集積回路C2'におけるメモリセルアレイMC2と、の間に、1ビットの入力データIが例えばm個シリアル形式で入力されると共に上記イネーブル信号ENが入力され、これらに対応した1ビットの出力データOが例えばm個シリアル形式で出力される並列演算器PPを介在させた構成を備える。このとき、メモリセルアレイMC1及びメモリセルアレイMC2はそれぞれ、例えば256ビット幅で512ワード(Word)分の規模を備えており、アドレス指定用の例えば8ビットのアドレスデータADがそれぞれに入力される。そしてこの場合の並列演算器PPは、256ビット分の多数判定入力回路12及び直列多数判定回路13と、上記パイプラインレジスタ21と、256ビットに対応する並列多数判定回路20と、により構成される。

10

20

30

40

50

【0230】

以上の構成において、ニューラルネットワークS1-3に含まれる上記ニューラルネットワーク集積回路C1'及びニューラルネットワーク集積回路C2'それぞれの動作は、図22及び図24を用いて説明した動作と同様となる。一方パイプラインレジスタ21は、例えば、ニューラルネットワーク集積回路C2'のメモリセルアレイMC2から読み出した重み付け係数W及び接続有無情報Cに基づいて並列多数判定回路20において出力データOの生成/出力処理を行っているタイミングでニューラルネットワーク集積回路C1'のメモリセルアレイMC1から読み出した出力データOを一時的に記憶する。そして、上記重み付け係数W及び接続有無情報Cに基づく並列多数判定回路20の処理が完了したタイミングで、メモリセルアレイMC1から読み出して記憶していた出力データOを並列多数判定回路20に出力してそれに基づいた出力データOの生成/出力処理を行わせる。この処理により、見かけ上はメモリセルアレイMC1からの出力データOの読み出しと、メモリセルアレイMC2からの重み付け係数W及び接続有無情報Cの読み出しと、を同時に行わせることができ、結果的に、図25を用いて説明したニューラルネットワークS1-2に対して略二倍の処理速度を実現させることができる。

【0231】

次に、図26に示すニューラルネットワーク回路C1-3における特に並列演算器PPの細部構成について、図27を用いて説明する。

【0232】

先ず図27(a)に示すように並列演算器PPは、メモリセルアレイMC1のビット幅に相当する数の上記多数判定入力回路12及び上記直列多数判定回路13からなる直列多数判定回路16と、メモリセルアレイMC1のビット幅に相当する上記パイプラインレジスタ21と、出力フリップフロップ回路22を介して出力データOを出力する上記並列多数判定回路20と、により構成されている。この構成においてパイプラインレジスタ21は図27(a)に示すように、メモリセルアレイMC1のビット幅に相当する出力レジスタ21U及び入力レジスタ21Lにより構成されており、上記イネーブル信号ENが入力レジスタ21Lに入力される。そして入力レジスタ21Lは、イネーブル信号ENが入力されるタイミングでそれに記憶(ラッチ)されているデータを並列多数判定回路20に出力すると共に、当該タイミングで出力レジスタ21Uに記憶されているデータを引き出して(即ちシフトさせて)記憶(ラッチ)する。またこれにより出力レジスタ21Uは、入力レジスタ21Lによりそのデータが引き出されたタイミングで、次の出力データOを記憶(ラッチ)する。以上の入力レジスタ21L及び出力レジスタ21Uの動作が繰り返されることにより、上述したパイプラインレジスタ21としての動作が実現される。

【0233】

次に、上記多数判定入力回路12及び直列多数判定回路13の細部構成について、図27(b)を用いて説明する。図27(b)に示すように一の直列多数判定回路16内の多数判定入力回路12は、排他的論理和回路12Aと、マスクフリップフロップ回路12B

と、により構成されている。この構成において排他的論理和回路 12A には、メモリセルアレイ MC1 からの重み付け係数 W と、1 ビットの入力データ I と、が入力され、これらの排他的論理和の結果を上記出力データ E として直列多数判定回路 13 に出力する。またマスクフリップフロップ回路 12B は、メモリセルアレイ MC1 からの接続有無情報 C と、上記イネーブル信号 EN と、が入力され、イネーブル信号 EN が入力されたタイミングで上記接続有無情報 C を直列多数判定回路 13 に出力する。そして直列多数判定回路 13 は、上記出力データ E 及び上記接続有無情報 C に基づいた上述した動作により出力データ O を生成して、パイプラインレジスタ 21 の出力レジスタ 21U に出力する。このとき、上記既定の多数判定閾値を直列多数判定回路 13 内の図示しないレジスタ内に保持してそれを参照することで、上述した直列多数判定回路 13 としての動作を実現できる。

10

【0234】

次に、上記並列多数判定回路 20 の細部構成について、図 27(c) を用いて説明する。図 27(c) に示すように並列多数判定回路 20 は、排他的論理和回路 20A と、マスクフリップフロップ回路 20B と、並列多数決回路 20C と、により構成されている。この構成において排他的論理和回路 20A には、メモリセルアレイ MC2 からの 1 ビットの重み付け係数 W と、パイプラインレジスタ 21 の入力レジスタ 21L からの 1 ビットの出力データ O と、が入力され、これらの排他的論理和の結果を並列多数決回路 20C に出力する。またマスクフリップフロップ回路 20B は、メモリセルアレイ MC2 からの接続有無情報 C と、上記イネーブル信号 EN と、が入力され、イネーブル信号 EN が入力されたタイミングで上記接続有無情報 C を並列多数決回路 20C に出力する。そして並列多数決回路 20C は、メモリセルアレイ MC2 からの一組の重み付け係数 W 及び接続有無情報 C に対応する排他的論理和回路 12A 及びマスクフリップフロップ回路 20B それぞれからの出力に基づいた上述した動作をメモリセルアレイ MC1 からの出力データ O の数 (図 26 及び図 27 に例示する場合は 256) だけ繰り返し、出力フリップフロップ回路 22 を介してシリアル形式で出力データ O として出力する。このとき、上記既定の多数判定閾値を並列多数判定回路 20 内の図示しないレジスタ内に保持してそれを参照することで、上述した並列多数判定回路 20 としての動作を実現できる。

20

【0235】

このとき、上述したパイプラインレジスタ 21 の動作により並列演算器 PP では、例えば図 27(d) に示すように、メモリセルアレイ MC1 からの 256 ビット分の出力データ O に対する処理 (図 27(d) において「メモリセルブロック 15_{U1}」と示す) が終了すると、次にメモリセルアレイ MC1 からの次の 256 ビット分の出力データ O に対する処理 (図 27(d) において「メモリセルブロック 15_{U2}」と示す) と、メモリセルアレイ MC2 からの 256 ビット分の重み付け係数 W 及び接続有無情報 C に対する処理 (図 27(d) において「メモリセルブロック 15_{L1}」と示す) と、が、見かけ上同時並列的に実行される。そして、メモリセルブロック 15_{U2} に対応する出力データ O 並びにメモリセルブロック 15_{L1} に対応する重み付け係数 W 及び接続有無情報 C に対する処理が終了すると、次にメモリセルアレイ MC1 からの更に次の 256 ビット分の出力データ O に対する処理 (図 27(d) において「メモリセルブロック 15_{U3}」と示す) と、メモリセルアレイ MC2 からの次の 256 ビット分の重み付け係数 W 及び接続有無情報 C に対する処理 (図 27(d) において「メモリセルブロック 15_{L2}」と示す) と、が、見かけ上同時並列的に実行される。以降は、メモリセルアレイ MC1 及びメモリセルアレイ MC2 それぞれからの 256 ビット分の出力データ O 並びに重み付け係数 W 及び接続有無情報 C に対して、逐次的且つ同時並列的な処理が実行される。

30

40

【0236】

なお、図 27(b) に示す多数判定入力回路 12 及び直列多数判定回路 13 の細部構成、並びに図 27(c) に示す並列多数判定回路 20 の細部構成は、図 21 以降に示す各メモリセル 11 からの上記接続有無情報 C の出力タイミングが、図 21 以降に示す各メモリセル 10 からの上記重み付け係数 W の出力タイミングよりも例えば一サイクル早いことを前提とした構成である。この出力タイミングのずれを吸収するのが、図 27(b) 及び図

50

27(c)にそれぞれ示すマスクフリップフロップ回路12B及びマスクフリップフロップ回路20Bの機能である。一方、上記重み付け係数Wの出力タイミングと上記接続有無情報Cの出力タイミングとを同時並行的とすることも可能である。そしてこの場合、図27(b)及び図27(c)にそれぞれ示すマスクフリップフロップ回路12B及びマスクフリップフロップ回路20Bは、多数判定入力回路12及び並列多数判定回路20としては不要となる。

【0237】

以上説明したように、図26及び図27に示すニューラルネットワーク回路C1-3によれば、図25(a)に示すニューラルネットワークS1-2を、約二倍の処理速度をもってモデル化することができる。なお図27を用いて説明した直列多数判定回路16の細部構成は、図25を用いて説明したニューラルネットワーク集積回路C1-2に含まれる直列多数判定回路16の細部構成として適用することもできる。

10

【産業上の利用可能性】

【0238】

以上それぞれ説明したように、本発明はニューラルネットワークをモデル化したニューラルネットワーク回路等の分野に利用することが可能であり、特に、製造コストの低減や効率的なニューラルネットワーク回路等を開発する場合に適用すれば、特に顕著な効果が得られる。

【符号の説明】

【0239】

20

10、11：メモリセル

NN：ニューラル電子回路

NNS：ニューラルネットワーク・システム

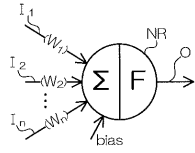
MC：メモリセルアレイ部（記憶部）

Pe：プロセスエレメント部（第1電子回路部）

PC₁・・・PC_n：プロセスエレメント・コラム

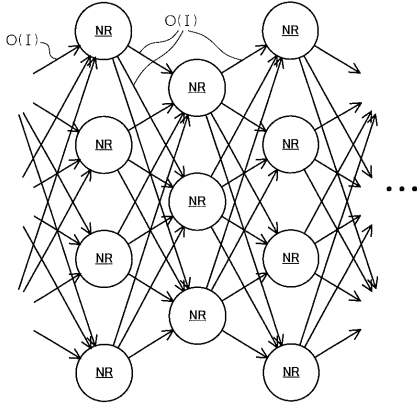
Act：加算活性化部（第2電子回路部）

【図 1】



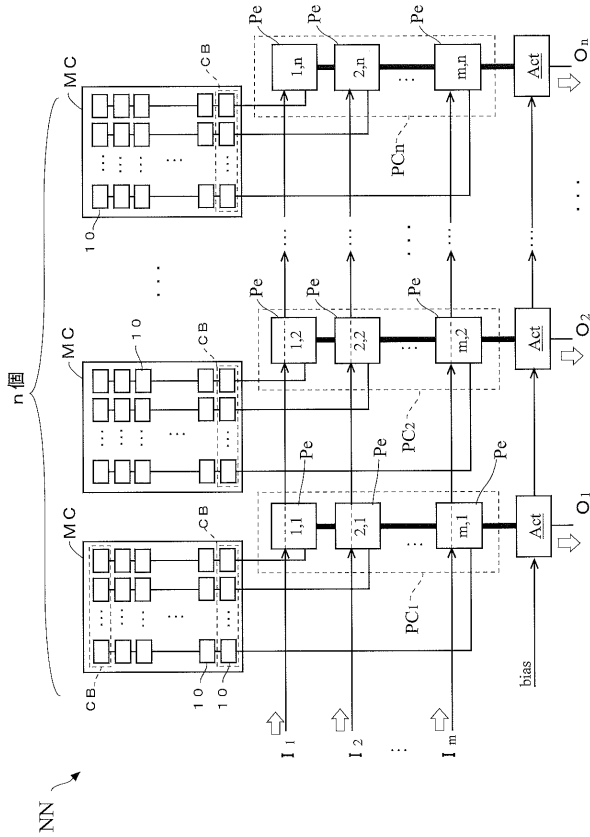
$$O = F(W_1 \times I_1 + W_2 \times I_2 + \dots + W_n \times I_n + \text{bias}) \dots (1)$$

(a)

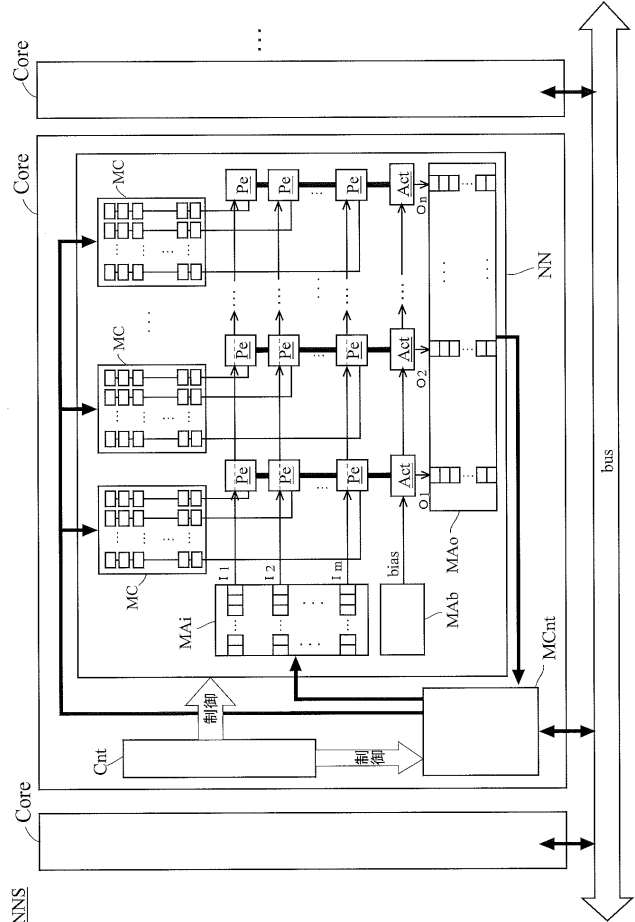


(b)

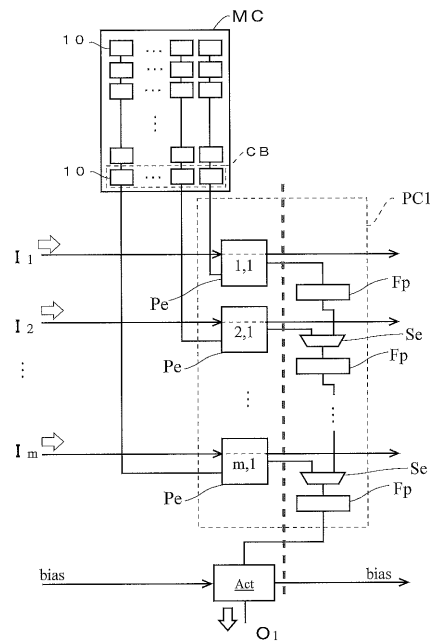
【図 3】



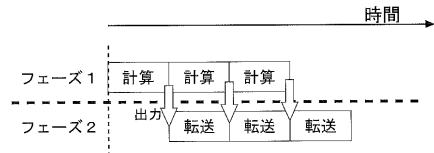
【図 2】



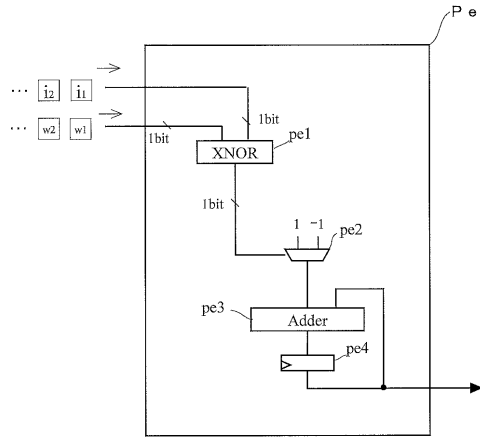
【図 4】



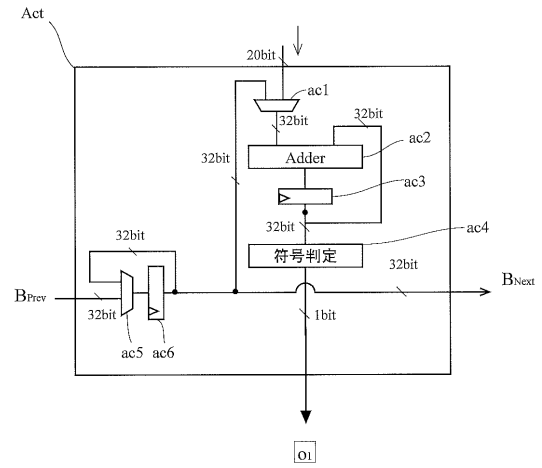
【図 5】



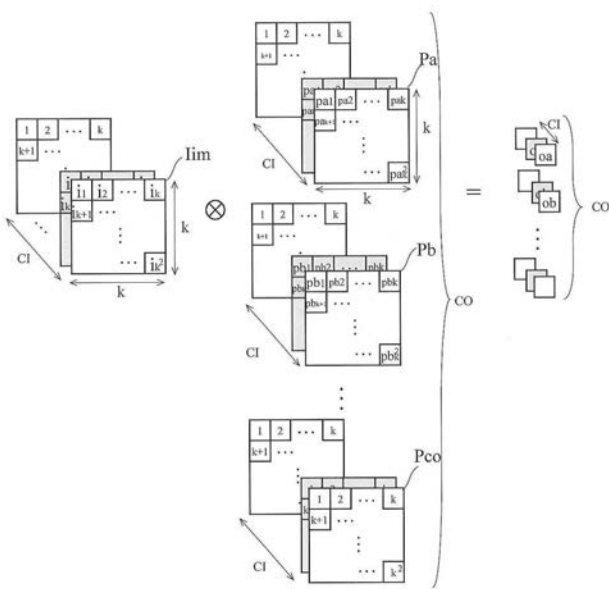
【 図 6 A 】



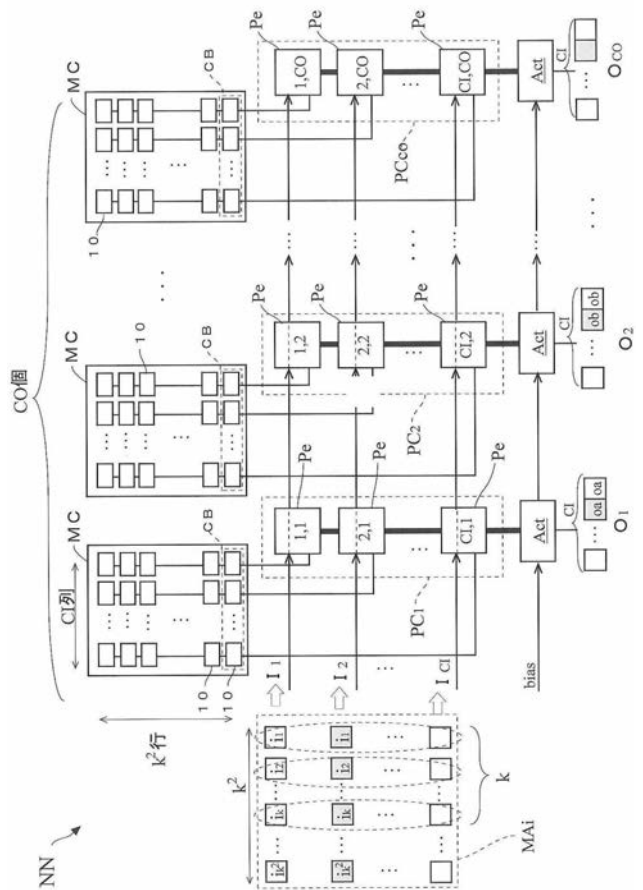
【 図 6 B 】



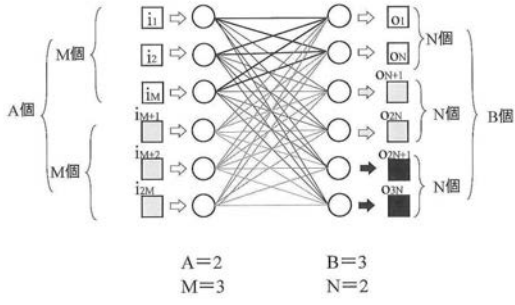
【 図 7 】



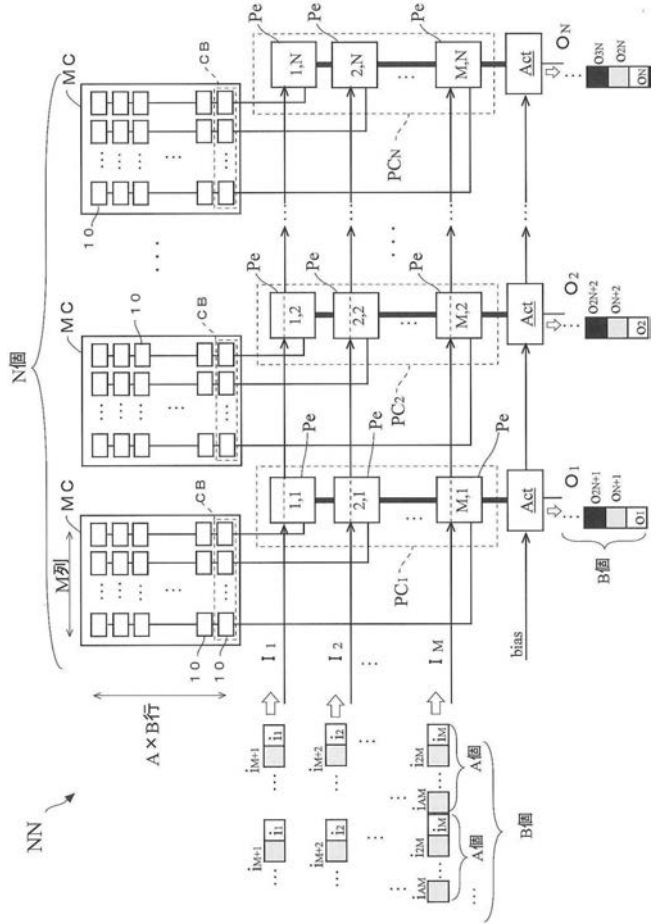
【 図 8 】



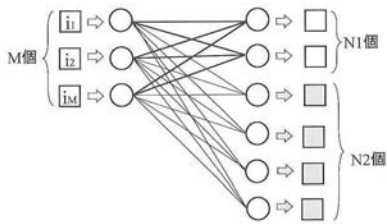
【 図 9 】



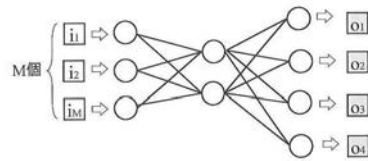
【 図 1 0 】



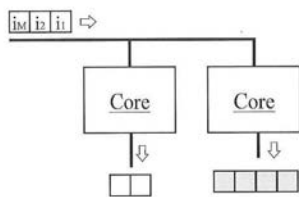
【 図 1 1 】



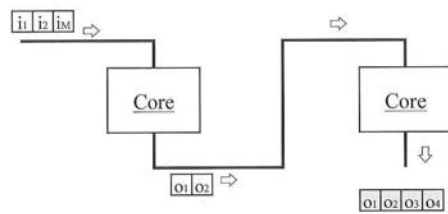
【 図 1 3 】



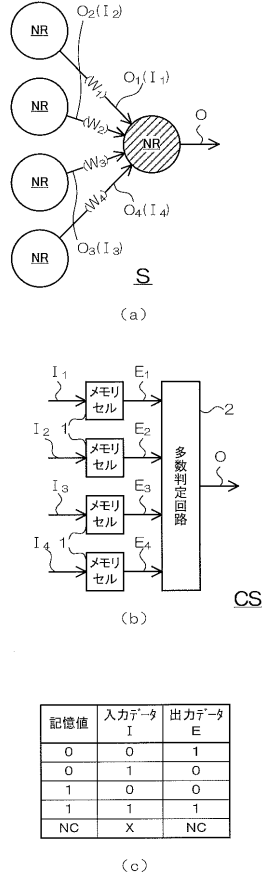
【 図 1 2 】



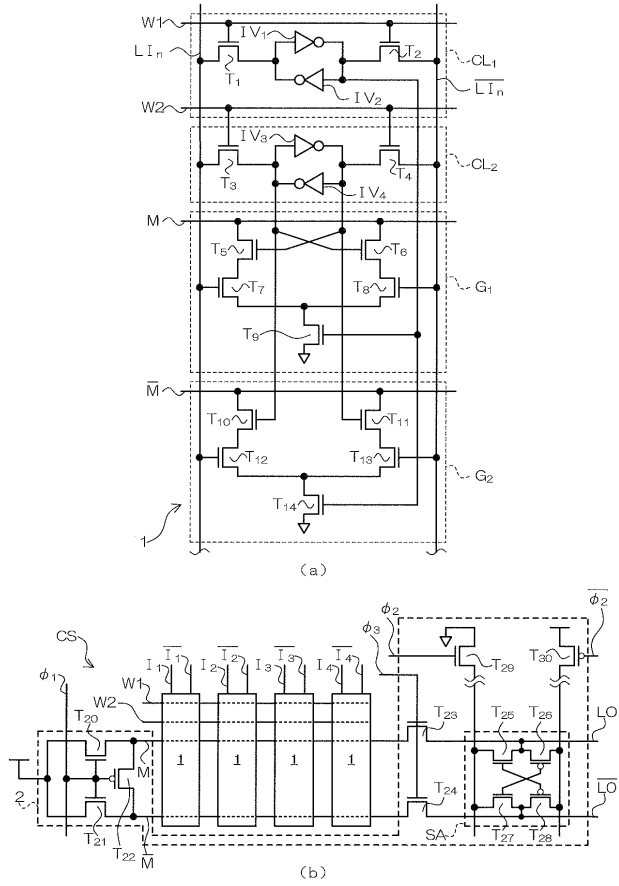
【 図 1 4 】



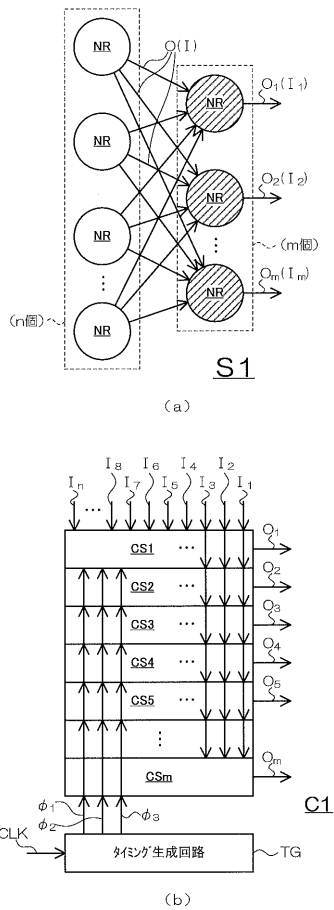
【図15】



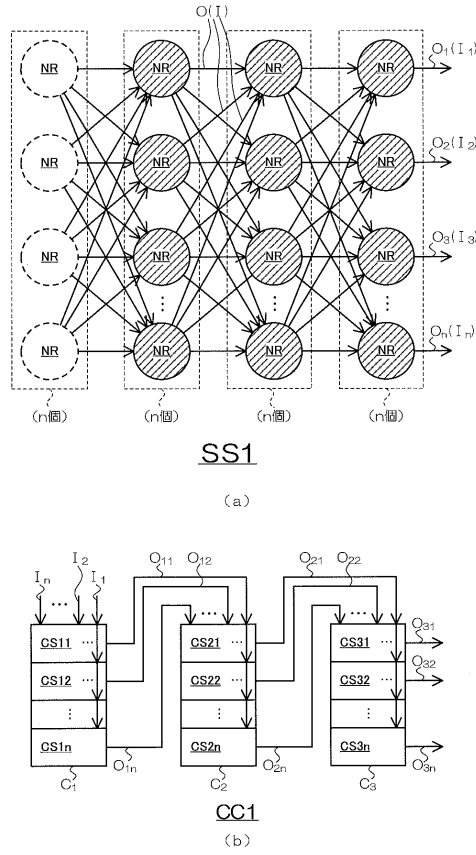
【図16】



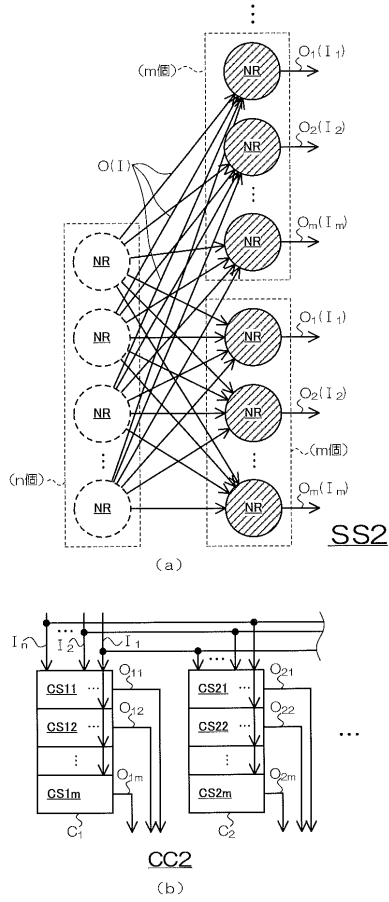
【図17】



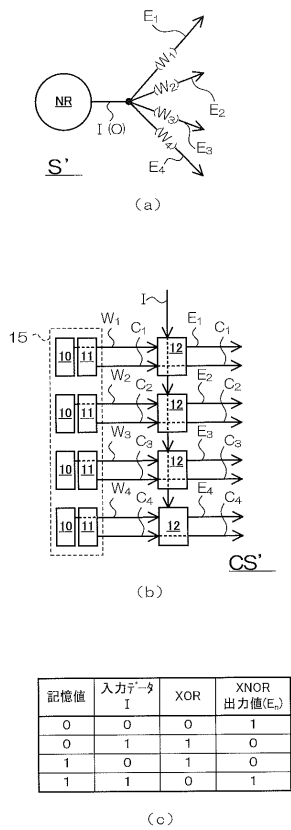
【図18】



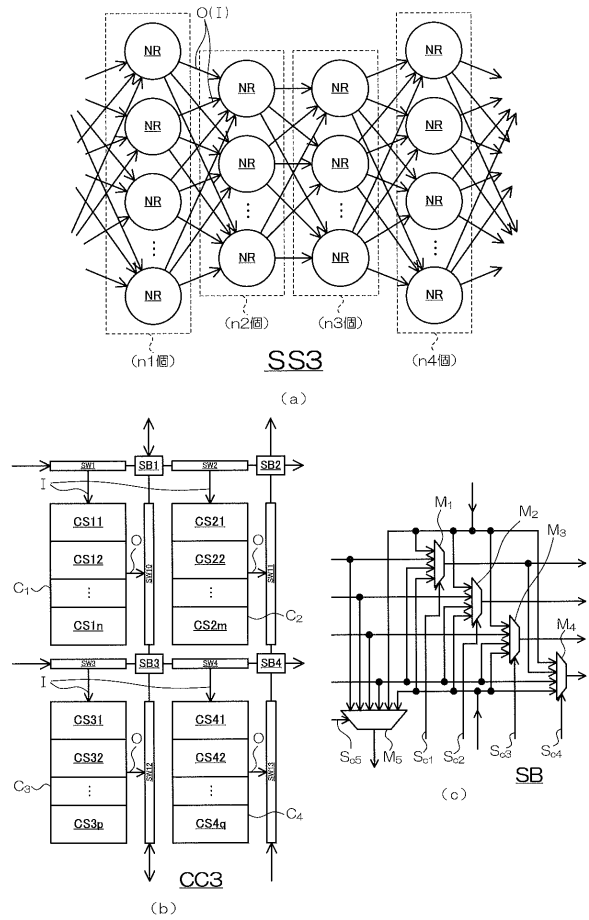
【図 19】



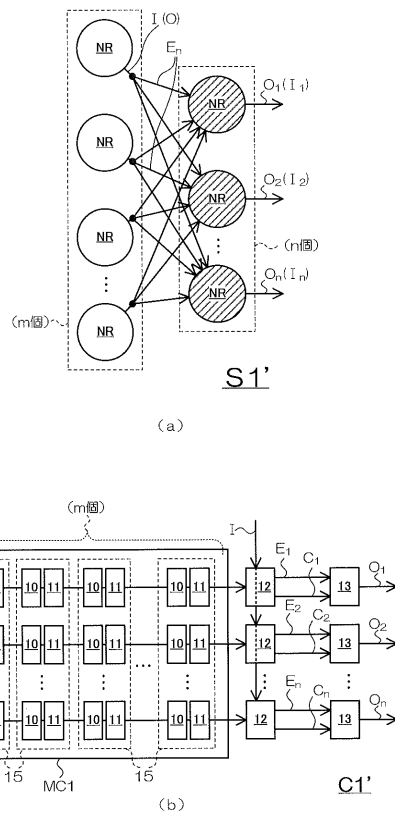
【図 21】



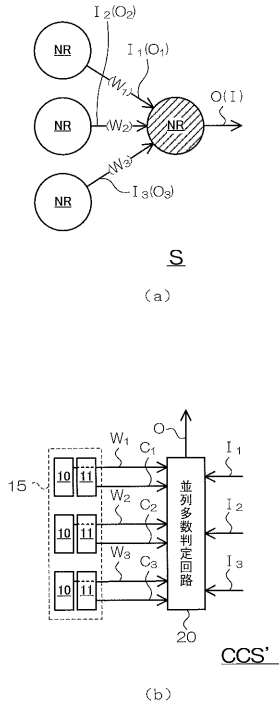
【図 20】



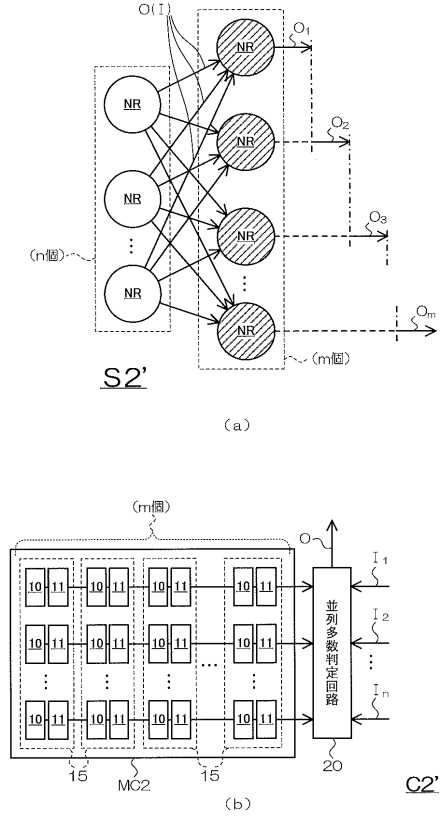
【図 22】



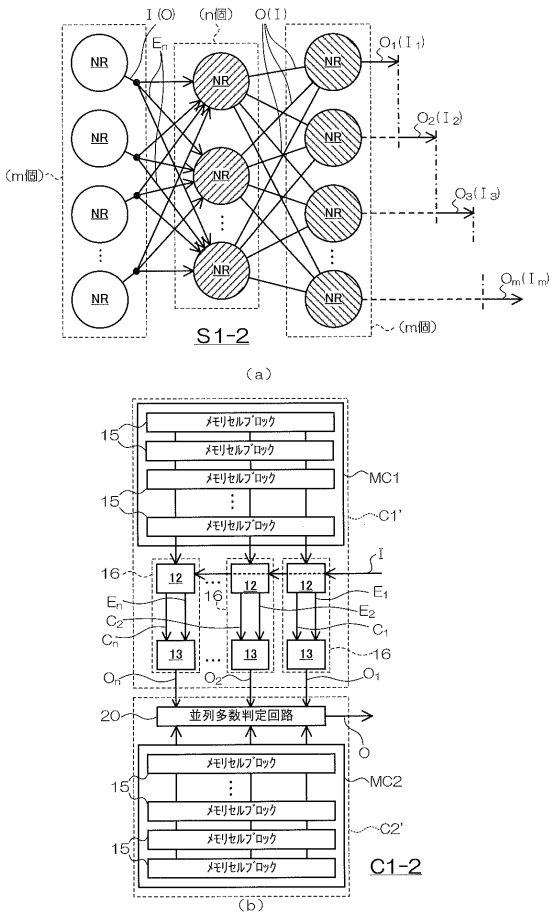
【図23】



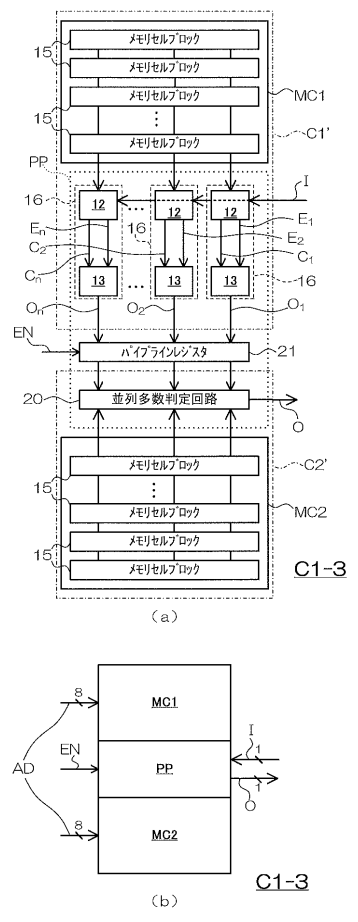
【図24】



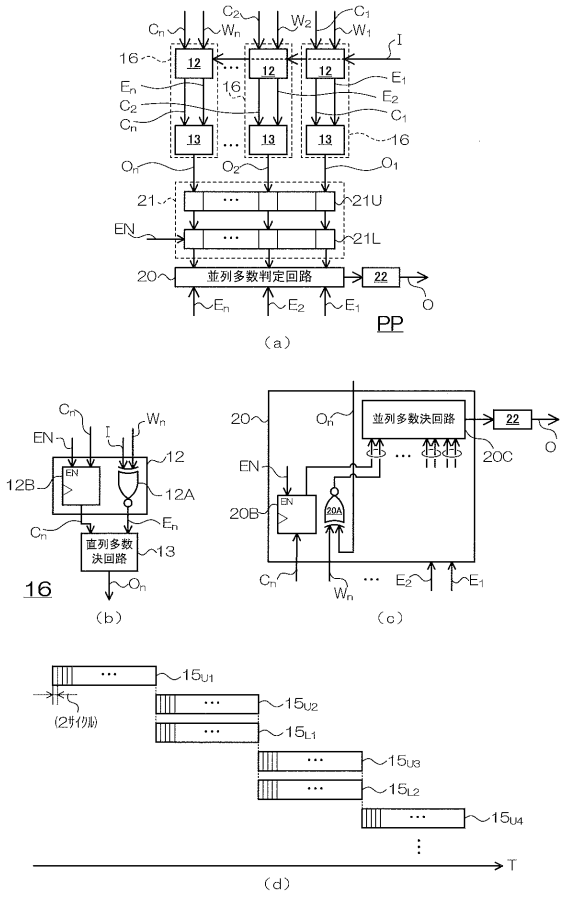
【図25】



【図26】



【 図 2 7 】



フロントページの続き

(72)発明者 本村 真人

北海道札幌市北区北8条西5丁目 国立大学法人北海道大学内