

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2017-509982
(P2017-509982A)

(43) 公表日 平成29年4月6日(2017.4.6)

(51) Int.Cl.
G06N 3/04 (2006.01)

F I
G06N 3/04

テーマコード (参考)

審査請求 未請求 予備審査請求 有 (全 41 頁)

(21) 出願番号 特願2016-553381 (P2016-553381)
 (86) (22) 出願日 平成27年2月13日 (2015. 2. 13)
 (85) 翻訳文提出日 平成28年10月18日 (2016. 10. 18)
 (86) 国際出願番号 PCT/US2015/015917
 (87) 国際公開番号 W02015/178977
 (87) 国際公開日 平成27年11月26日 (2015. 11. 26)
 (31) 優先権主張番号 61/943, 155
 (32) 優先日 平成26年2月21日 (2014. 2. 21)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 14/273, 214
 (32) 優先日 平成26年5月8日 (2014. 5. 8)
 (33) 優先権主張国 米国 (US)

(71) 出願人 595020643
 クォアルコム・インコーポレイテッド
 QUALCOMM INCORPORATED
 アメリカ合衆国、カリフォルニア州 92
 121-1714、サン・ディエゴ、モア
 ハウス・ドライブ 5775
 (74) 代理人 100108855
 弁理士 蔵田 昌俊
 (74) 代理人 100109830
 弁理士 福原 淑弘
 (74) 代理人 100158805
 弁理士 井関 守三
 (74) 代理人 100112807
 弁理士 岡田 貴志

最終頁に続く

(54) 【発明の名称】 原位置ニューラルネットワークコプロセッシング

(57) 【要約】

ニューラルネットワークにおいてコプロセッシングを実行する方法は、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングすることを含む。本方法はまた、第1の処理ノードでニューラルネットワークの一部を実行することを含む。さらに、本方法は、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返すことを含む。さらに、本方法は、第2の処理ノードでニューラルネットワークの一部を実行することを含む。

【選択図】 図9

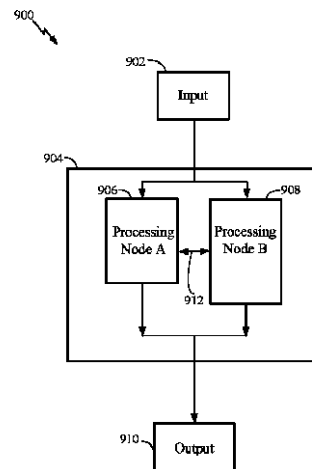


FIG. 9

【公報種別】特許法第17条の2の規定による補正の掲載
【部門区分】第6部門第3区分
【発行日】平成30年3月1日(2018.3.1)

【公表番号】特表2017-509982(P2017-509982A)
【公表日】平成29年4月6日(2017.4.6)
【年通号数】公開・登録公報2017-014
【出願番号】特願2016-553381(P2016-553381)
【国際特許分類】
 G 0 6 N 3/04 (2006.01)
【F I】
 G 0 6 N 3/04

【手続補正書】

【提出日】平成30年1月18日(2018.1.18)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

人工ニューラルネットワークにおいてコプロセッシングを実行する、コンピュータ実装方法であって、

一定時間期間にわたって、前記ニューラルネットワークの一部を第2の処理ノードから第1の処理ノードにスワッピングすること、ここにおいて、前記第1の処理ノードおよび前記第2の処理ノードは、相互に前記ニューラルネットワークの機能的特徴の処理を包含するように構成され、前記第1の処理ノードは、強化学習を実装するように構成される、学習処理コアを備える、と、

前記第1の処理ノードで前記ニューラルネットワークの前記一部を実行することと、
前記一定時間期間後に、前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すことと、

前記第2の処理ノードで前記ニューラルネットワークの前記一部を実行すること、ここにおいて、前記一部は、前記ニューラルネットワークに関する状態変数および接続性情報を含む、前記一部がそこから送られる前記ノードについての状態情報を備える、と
を備える、方法。

【請求項2】

前記第1の処理ノードは、第1のハードウェアコアに含まれ、前記第2の処理ノードは、第2のハードウェアコアに含まれ、前記第1のハードウェアコアは、前記第2のハードウェアコアとは別個である、

請求項1に記載の方法。

【請求項3】

前記学習処理コアは、前記第2の処理ノードよりも多くのリソースで構成される、
請求項1に記載の方法。

【請求項4】

前記第2の処理ノードは、前記ニューラルネットワークまたは前記一部を動作することに関連付けられる機能を実行するために構成された静的処理コアを備え、
スワッピングすることは、

前記静的処理コアの状態を前記学習処理コアにコピーすることと、

前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに

入力をルーティングすることと

を備え、

返すことは、

前記学習処理コアの状態を前記静的処理コアにコピーすることと、

前記静的処理コアに制御を返すことと

を備える、請求項 1 に記載の方法。

【請求項 5】

前記人工ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、

請求項 1 に記載の方法。

【請求項 6】

前記第 1 の処理ノードは、デバッグングコアを備える、

請求項 1 に記載の方法。

【請求項 7】

人工ニューラルネットワークにおいてコプロセッシングを実行するための装置であって、

メモリと、

前記メモリに結合された少なくとも 1 つのプロセッサと

を備え、前記少なくとも 1 つのプロセッサは、

一定時間期間にわたって、前記ニューラルネットワークの一部を第 2 の処理ノードから第 1 の処理ノードにスワッピングすること、
ここにおいて、前記第 1 の処理ノードおよび前記第 2 の処理ノードは、相互に前記ニューラルネットワークの機能的特徴の前記処理を包含するように構成され、前記第 1 の処理ノードは、強化学習を実装するように構成される、学習処理コアを備える、と、

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行することと、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を前記第 2 の処理ノードに返すことと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行すること、
ここにおいて、前記一部は、前記ニューラルネットワークに関する状態変数および接続性情報を含む、前記一部がそこから送られる前記ノードについての状態情報を備える、と

を行うように構成される、装置。

【請求項 8】

前記第 1 の処理ノードは、第 1 のハードウェアコアに含まれ、前記第 2 の処理ノードは、第 2 のハードウェアコアに含まれ、前記第 1 のハードウェアコアは、前記第 2 のハードウェアコアとは別個である、

請求項 7 に記載の装置。

【請求項 9】

前記学習処理コアは、前記第 2 の処理ノードよりも多くのリソースで構成される、

請求項 7 に記載の装置。

【請求項 10】

前記第 2 の処理ノードは、前記ニューラルネットワークまたは前記一部を動作することに関連付けられる機能を実行するために構成された静的処理コアを備え、前記少なくとも 1 つのプロセッサは、

前記静的処理コアの状態を前記学習処理コアにコピーすることと、

前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに入力をルーティングすることと、

前記学習処理コアの状態を前記静的処理コアにコピーすることと、

変更された静的処理コアに制御を返すことと

を行うようにさらに構成される、請求項 7 に記載の装置。

【請求項 11】

前記人工ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、

請求項7に記載の装置。

【請求項 1 2】

前記第 1 の処理ノードは、デバッグコアを備える、

請求項7に記載の装置。

【請求項 1 3】

人工ニューラルネットワークにおいてコプロセッシングを実行するためのプログラムコードを符号化した非一時的コンピュータ可読媒体であって、前記プログラムコードは、プロセッサによって実行され、請求項 1 ~ 請求項 6 に記載の方法のいずれかを実行するためのプログラムコードを備える、非一時的コンピュータ可読媒体。

【特許請求の範囲】**【請求項 1】**

ニューラルネットワークにおいてコプロセッシングを実行する方法であって、
一定時間期間にわたって、前記ニューラルネットワークの一部を第 1 の処理ノードにスワッピングすることと、

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行することと、
前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すことと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行することと
を備える、方法。

10

【請求項 2】

前記第 1 の処理ノードは、別個のハードウェアコアを備える、
請求項 1 に記載の方法。

【請求項 3】

前記第 1 の処理ノードは、学習処理コアを備える、
請求項 1 に記載の方法。

【請求項 4】

前記学習処理コアは、前記第 2 の処理ノードよりも高いレベルのリソースで構成される、

請求項 3 に記載の方法。

20

【請求項 5】

学習は、オフラインまたはオンラインで実装される、
請求項 3 に記載の方法。

【請求項 6】

前記学習処理コアの入力および出力は、学習がオフラインで実装される場合、前記ニューラルネットワークの他のレイヤを備える、
請求項 5 に記載の方法。

【請求項 7】

前記第 1 の処理ノードは、学習処理コアを備え、
前記第 2 の処理ノードは、静的処理コアを備え、
スワッピングすることは、

30

前記静的処理コアの状態を前記学習処理コアにコピーすることと、
前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに
入力をルーティングすることと

を備え、

返すことは、

前記学習処理コアの状態を前記静的処理コアにコピーすることと、
変更された静的処理コアに制御を返すことと

を備える、請求項 1 に記載の方法。

【請求項 8】

前記スワッピングすることは、前記第 1 の処理ノードから前記第 2 の処理ノードにリソ
ースを割り振ることを備える、

請求項 1 に記載の方法。

40

【請求項 9】

前記ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、
請求項 1 に記載の方法。

【請求項 10】

前記第 1 の処理ノードは、デバッグコアを備える、
請求項 1 に記載の方法。

【請求項 11】

50

前記スワッピングすることは、システム性能がしきい値を下回る場合に発生する、請求項 1 に記載の方法。

【請求項 1 2】

前記返すことは、システム性能がしきい値を上回る場合に発生する、請求項 1 に記載の方法。

【請求項 1 3】

前記スワッピングすること、または返すことは、電力がシステムに適用されると発生する、請求項 1 に記載の方法。

【請求項 1 4】

ニューラルネットワークにおいてコプロセッシングを実行するための装置であって、メモリと、前記メモリに結合された少なくとも 1 つのプロセッサとを備え、前記少なくとも 1 つのプロセッサは、一定時間期間にわたって、前記ニューラルネットワークの一部を第 1 の処理ノードにスワッピングすることと、前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行することと、前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すことと、前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行することとを行うように構成される、装置。

【請求項 1 5】

前記第 1 の処理ノードは、別個のハードウェアコアを備える、請求項 1 4 に記載の装置。

【請求項 1 6】

前記第 1 の処理ノードは、学習処理コアを備える、請求項 1 4 に記載の装置。

【請求項 1 7】

前記学習処理コアは、前記第 2 の処理ノードよりも高いレベルのリソースで構成される、請求項 1 6 に記載の装置。

【請求項 1 8】

学習は、オフラインまたはオンラインで実装される、請求項 1 6 に記載の装置。

【請求項 1 9】

前記学習処理コアの入力および出力は、学習がオフラインで実装される場合、前記ニューラルネットワークの他のレイヤを備える、請求項 1 8 に記載の装置。

【請求項 2 0】

前記第 1 の処理ノードは、学習処理コアを備え、前記第 2 の処理ノードは、静的処理コアを備え、前記少なくとも 1 つのプロセッサは、前記静的処理コアの状態を前記学習処理コアにコピーすることと、前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに入力をルーティングすることと、前記学習処理コアの状態を前記静的処理コアにコピーすることと、変更された静的処理コアに制御を返すこととを行うようにさらに構成される、請求項 1 4 に記載の装置。

【請求項 2 1】

前記少なくとも 1 つのプロセッサは、前記第 1 の処理ノードから前記第 2 の処理ノードにリソースを割り振ることを行うようにさらに構成される、

10

20

30

40

50

請求項 1 4 に記載の装置。

【請求項 2 2】

前記ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、請求項 1 4 に記載の装置。

【請求項 2 3】

前記第 1 の処理ノードは、デバッグコアを備える、請求項 1 4 に記載の装置。

【請求項 2 4】

前記少なくとも 1 つのプロセッサは、システム性能がしきい値を下回る場合に、前記ニューラルネットワークの前記一部を前記第 1 の処理ノードにスワッピングするようにさらに構成される、請求項 1 4 に記載の装置。

10

【請求項 2 5】

前記少なくとも 1 つのプロセッサは、システム性能がしきい値を上回る場合に、前記ニューラルネットワークの前記一部を前記第 2 の処理ノードに返すようにさらに構成される、請求項 1 4 に記載の装置。

【請求項 2 6】

前記少なくとも 1 つのプロセッサは、電力がシステムに適用されると、前記ニューラルネットワークの前記一部を前記第 1 の処理ノードにスワッピングする、または前記ニューラルネットワークの前記一部を前記第 2 の処理ノードに返すようにさらに構成される、請求項 1 4 に記載の装置。

20

【請求項 2 7】

ニューラルネットワークにおいてコプロセッシングを実行するための装置であって、一定時間期間にわたって、前記ニューラルネットワークの一部を第 1 の処理ノードにスワッピングするための手段と、前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行するための手段と、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すための手段と、

30

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行するための手段とを備える、装置。

【請求項 2 8】

ニューラルネットワークにおいてコプロセッシングを実行するためのコンピュータプログラム製品であって、

プログラムコードを符号化した非一時的コンピュータ可読媒体を備え、前記プログラムコードは、

一定時間期間にわたって、前記ニューラルネットワークの一部を第 1 の処理ノードにスワッピングするためのプログラムコードと、

40

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行するためのプログラムコードと、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すためのプログラムコードと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行するためのプログラムコードと

を備える、コンピュータプログラム製品。

【発明の詳細な説明】

【関連出願の相互参照】

【0 0 0 1】

50

[0001]本出願は、2014年2月21日に開示された「IN SITU NEURAL NETWORK CO-PROCESSING」と題する米国仮特許出願第61/943,155号の利益を主張し、その開示は、参照によりその全体が本明細書に明示的に組み込まれる。

【技術分野】

【0002】

[0002]本開示のいくつかの態様は、一般にニューラルシステムエンジニアリングに関し、より詳細には、原位置ニューラルネットワークコプロセッシング (in situ neural network co-processing) のためのシステムおよび方法に関する。

【背景技術】

【0003】

[0003]人工ニューロン(すなわち、ニューロンモデル)の相互結合されたグループを備え得る人工ニューラルネットワークは、計算デバイスであるか、または計算デバイスによって実行される方法を表す。人工ニューラルネットワークは、生物学的ニューラルネットワークにおける対応する構造および/または機能を有し得る。しかしながら、人工ニューラルネットワークは、従来の計算技法が厄介、実行不可能または不適切であるいくつかの適用例に革新的で有用な計算技法を提供することができる。人工ニューラルネットワークは観測から関数を推測することができるので、そのようなネットワークは、タスクまたはデータの複雑さが従来の技法による関数の設計を煩わしくする用途において、特に有用である。

【発明の概要】

【0004】

[0004]本開示のある態様では、ニューラルネットワークにおいてコプロセッシングを実行する方法が開示される。本方法は、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングすることを含む。本方法はまた、第1の処理ノードでニューラルネットワークの一部を実行することを含む。さらに、本方法は、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返すことを含む。本方法は、第2の処理ノードでニューラルネットワークの一部を実行することをさらに含む。

【0005】

[0005]本開示の別の態様では、ニューラルネットワーク内でコプロセッシングを実行するための装置が開示される。本装置は、メモリと、メモリに結合された少なくとも1つのプロセッサとを含む。本プロセッサは、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングするように構成される。本プロセッサはまた、第1の処理ノードでニューラルネットワークの一部を実行するように構成される。さらに、本プロセッサは、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返すように構成される。本プロセッサは、第2の処理ノードでニューラルネットワークの一部を実行するようにさらに構成される。

【0006】

[0006]本開示の別の態様では、ニューラルネットワーク内でコプロセッシングを実行するための装置が開示される。本装置は、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングするための手段を有する。本装置はまた、第1の処理ノードでニューラルネットワークの一部を実行するための手段を有する。さらに、本装置は、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返すための手段を有する。本装置は、第2の処理ノードでニューラルネットワークの一部を実行するための手段をさらに有する。

【0007】

[0007]本開示の別の態様では、ニューラルネットワーク内でコプロセッシングを実行するためのコンピュータプログラム製品が開示される。本コンピュータプログラム製品は、プログラムコードを符号化した非一時的コンピュータ可読媒体を含む。本プログラムコードは、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにス

10

20

30

40

50

ワッピングするためのプログラムコードを含む。本プログラムコードはまた、第1の処理ノードでニューラルネットワークの一部を実行するためのプログラムコードを含む。さらに、本プログラムコードは、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返すためのプログラムコードを含む。本プログラムコードは、第2の処理ノードでニューラルネットワークの一部を実行するためのプログラムコードをさらに含む。

【0008】

[0008]これは、以下の詳細な説明がより良く理解され得るために、本開示の特徴および技術的利点をかなり広く概説した。本開示の追加の特徴および利点は、以下で説明される。この開示は、本開示と同じ目的を実行するための他の構造を修正または設計するための基礎として容易に変更され得ることが、当業者によって理解されるべきである。また、添付の特許請求の範囲に記載されるように、そのような等価な構成は本開示の教示から逸脱しないことが、当業者によって理解されるべきである。本開示の特徴と考えられる新規な特徴は、その構成と動作の方法との両方に関して、さらなる目的および利点とともに、添付の図面と関連して考慮されるとき以下の説明からより良く理解されるであろう。しかしながら、図面の各々は単に例示および説明の目的のために提供されているにすぎず、本開示の制限の定義として意図されていないことが、明確には理解されるべきである。

10

【図面の簡単な説明】

【0009】

[0009]本開示の特徴、性質、および利点は、同様の参照文字が全体を通して相応して識別する図面を考慮した場合、以下に示される詳細な説明から、より明らかになるだろう。

20

【図1】本開示のいくつかの態様によるニューロンの例示的なネットワークを示す図。

【図2】本開示のいくつかの態様による、計算ネットワーク（ニューラルシステムまたはニューラルネットワーク）の処理ユニット（ニューロン）の一例を示す図。

【図3】本開示のいくつかの態様によるスパイクタイミング依存可塑性（STDP）曲線の一例を示す図。

【図4】本開示のいくつかの態様による、ニューロンモデルの挙動を定義するための正レジームおよび負レジームの一例を示す図。

【図5】本開示のある態様による、汎用プロセッサを使用してニューラルネットワークを設計することの例示的な実装形態を示す図。

【図6】本開示のいくつかの態様による、メモリが個々の分散処理ユニットとインターフェースされ得るニューラルネットワークを設計する例示的な実装形態を示す図。

30

【図7】本開示のいくつかの態様による、分散メモリおよび分散処理ユニットに基づいてニューラルネットワークを設計する例示的な実装形態を示す図。

【図8】本開示のいくつかの態様による、ニューラルネットワークの例示的な実装形態を示す図。

【図9】本開示の態様による、ニューラルネットワークの例示的なアーキテクチャを示すブロック図。

【図10A】本開示の態様による、ニューラルネットワーク内の原位置コプロセッシングを示す例示的なブロック図。

【図10B】本開示の態様による、ニューラルネットワーク内の原位置コプロセッシングを示す例示的なブロック図。

40

【図10C】本開示の態様による、ニューラルネットワーク内の原位置コプロセッシングを示す例示的なブロック図。

【図10D】本開示の態様による、ニューラルネットワーク内の原位置コプロセッシングを示す例示的なブロック図。

【図10E】本開示の態様による、ニューラルネットワーク内の原位置コプロセッシングを示す例示的なブロック図。

【図10F】本開示の態様による、ニューラルネットワーク内の原位置コプロセッシングを示す例示的なブロック図。

【図11】本開示の態様による、ニューラルネットワーク内でコプロセッシングを実行す

50

るための方法を示すブロック図。

【図 1 2】本開示の態様による、ニューラルネットワーク内でコプロセッシングを実行するための方法を示すブロック図。

【発明を実施するための形態】

【0010】

[0021]添付の図面に関連して以下に示される詳細な説明は、様々な構成の説明として意図されたものであり、本明細書において説明される概念が実現され得る唯一の構成を表すことを意図されるものではない。詳細な説明は、様々な概念の完全な理解を提供する目的で、具体的な詳細を含む。しかしながら、これらの概念がこれらの具体的な詳細なしで実施され得ることは、当業者にとっては明らかであろう。いくつかの事例では、よく知られている構造および構成要素が、そのような概念を曖昧にするのを避けるために、ブロック図形式で示される。

10

【0011】

[0022]本教示に基づいて、本開示の範囲は、本開示の任意の他の態様とは無関係に実装されるにせよ、本開示の任意の他の態様と組み合わせられるにせよ、本開示のいかなる態様をもカバーするものであることを、当業者なら諒解されたい。たとえば、記載される態様をいくつ使用しても、装置は実装され得、または方法は実施され得る。さらに、本開示の範囲は、記載される本開示の様々な態様に加えてまたはそれらの態様以外に、他の構造、機能、または構造および機能を使用して実施されるそのような装置または方法をカバーするものとする。開示する本開示のいずれの態様も、請求項の1つまたは複数の要素によって実施され得ることを理解されたい。

20

【0012】

[0023]「例示的」という単語は、本明細書では「例、事例、または例示の働きをすること」を意味するために使用される。「例示的」として本明細書で説明するいかなる態様も、必ずしも他の態様よりも好ましいまたは有利であると解釈されるべきであるとは限らない。

【0013】

[0024]本明細書では特定の態様について説明するが、これらの態様の多くの変形および置換は本開示の範囲内に入る。好ましい態様のいくつかの利益および利点が説明されるが、本開示の範囲は特定の利益、使用、または目的に限定されるものではない。むしろ、本開示の態様は、様々な技術、システム構成、ネットワーク、およびプロトコルに広く適用可能であるものとし、そのうちのいくつかを例として図および好ましい態様についての以下の説明で示す。発明を実施するための形態および図面は、本開示を限定するものではなく説明するものにすぎず、本開示の範囲は添付の特許請求の範囲およびその均等物によって定義される。

30

例示的なニューラルシステム、トレーニングおよび動作

【0014】

[0025]図 1 は、本開示のいくつかの態様による、複数のレベルのニューロンをもつ例示的な人工ニューラルシステム 100 を示す。ニューラルシステム 100 は、シナプス結合のネットワーク 104 (すなわち、フィードフォワード結合) を介してニューロンの別のレベル 106 に結合されたニューロンのあるレベル 102 を有し得る。簡単のために、図 1 には 2 つのレベルのニューロンのみが示されているが、ニューラルシステムには、より少ないまたはより多くのレベルのニューロンが存在し得る。ニューロンのいくつかは、ラテラル結合を介して同じレイヤの他のニューロンに結合し得ることに留意されたい。さらに、ニューロンのいくつかは、フィードバック結合を介して前のレイヤのニューロンに戻る形で結合し得る。

40

【0015】

[0026]図 1 に示すように、レベル 102 における各ニューロンは、前のレベル (図 1 に図示せず) のニューロンによって生成され得る入力信号 108 を受信し得る。信号 108 は、レベル 102 のニューロンの入力電流を表し得る。この電流は、膜電位を充電するた

50

めにニューロン膜上に蓄積され得る。膜電位がそのしきい値に達すると、ニューロンは、発火し、ニューロンの次のレベル（たとえば、レベル106）に転送されるべき出力スパイクを生成し得る。いくつかのモデリング手法では、ニューロンは、信号をニューロンの次のレベルに継続的に転送し得る。この信号は、典型的には膜電位の関数である。そのような挙動は、以下で説明するものなどのアナログおよびデジタル実装形態を含むハードウェアおよび/またはソフトウェアでエミュレートまたはシミュレートされ得る。

【0016】

[0027]生物学的ニューロンでは、ニューロンが発火するときに生成される出力スパイクは、活動電位と呼ばれる。電気信号は、約100mVの振幅と約1msの持続時間とを有する比較的急速で、一時的な神経インパルスである。一連の結合されたニューロンを有するニューラルシステムの特定の実施形態（たとえば、図1におけるあるレベルのニューロンから別のレベルのニューロンへのスパイクの転送）では、あらゆる活動電位が基本的に同じ振幅と持続時間とを有するので、信号における情報は、振幅によってではなく、スパイクの周波数および数、またはスパイクの時間によってのみ表され得る。活動電位によって搬送される情報は、スパイク、スパイクしたニューロン、および他の1つまたは複数のスパイクに対するスパイクの時間によって決定され得る。以下で説明するように、スパイクの重要性は、ニューロン間の接続に適用される重みによって決定され得る。

10

【0017】

[0028]図1に示されるように、ニューロンのあるレベルから別のレベルへのスパイクの転送は、シナプス結合（または、単純に「シナプス」）104のネットワークを介して達成され得る。シナプス104に関して、レベル102のニューロンはシナプス前ニューロンと考えられ得、レベル106のニューロンはシナプス後ニューロンと考えられ得る。シナプス104は、レベル102のニューロンから出力信号（すなわち、スパイク）を受信して、調整可能なシナプスの重み

20

【0018】

【数1】

$$w_1^{(i,i+1)}, \dots, w_P^{(i,i+1)}$$

【0019】

に応じてそれらの信号をスケールリングすることができ、上式で、Pはレベル102のニューロンとレベル106のニューロンとの間のシナプス結合の総数であり、iはニューロンレベルの指標である。図1の例では、iはニューロンレベル102を表し、i+1は、ニューロンレベル106を表す。さらに、スケールリングされた信号は、レベル106における各ニューロンの入力信号として合成され得る。レベル106におけるあらゆるニューロンは、対応する合成された入力信号に基づいて、出力スパイク110を生成し得る。出力スパイク110は、シナプス結合の別のネットワーク（図1には図示せず）を使用して、別のレベルのニューロンに転送され得る。

30

【0020】

[0029]生物学的シナプスは、シナプス後ニューロンにおける興奮性活動または抑制性（過分極化）活動のいずれかを調停することができ、ニューロン信号を増幅する役目を果たすことができる。興奮性信号は、膜電位を脱分極する（すなわち、静止電位に対して膜電位を増加させる）。しきい値を超えて膜電位を脱分極するために十分な興奮性信号が一定の時間期間内に受信された場合、シナプス後ニューロンに活動電位が生じる。対照的に、抑制性信号は一般に、膜電位を過分極する（すなわち、低下させる）。抑制性信号は、十分に強い場合、興奮性信号のすべてを相殺し、膜電位がしきい値に達するのを防止することができる。シナプス興奮を相殺することに加えて、シナプス抑制は、自然にアクティブなニューロンに対して強力な制御を行うことができる。自然にアクティブなニューロンは、たとえば、そのダイナミクスまたはフィードバックに起因するさらなる入力なしにスパイクするニューロンを指す。これらのニューロンにおける活動電位の自然な生成を抑圧

40

50

することによって、シナプス抑制は、一般にスカルプチャリングと呼ばれる、ニューロンの発火のパターンを形成することができる。様々なシナプス104は、望まれる挙動に応じて、興奮性シナプスまたは抑制性シナプスの任意の組合せとして働き得る。

【0021】

[0030]ニューラルシステム100は、汎用プロセッサ、デジタル信号プロセッサ(DSP)、特定用途向け集積回路(ASIC)、フィールドプログラマブルゲートアレイ(FPGA)もしくは他のプログラマブル論理デバイス(PLD)、個別ゲートもしくはトランジスタ論理、個別ハードウェア構成要素、プロセッサによって実行されるソフトウェアモジュール、またはそれらの任意の組合せによってエミュレートされ得る。ニューラルシステム100は、たとえば画像およびパターン認識、機械学習、モータ制御、および似ているなど、かなりの適用範囲において利用され得る。ニューラルシステム100における各ニューロンは、ニューロン回路として実装され得る。出力スパイクを開始するしきい値まで充電されるニューロン膜は、たとえば、そこを流れる電流を積分するキャパシタとして実装され得る。

10

【0022】

[0031]一態様では、キャパシタは、ニューロン回路の電流積分デバイスとして除去され得、その代わりにより小さいメモrista(memristor)要素が使用され得る。この手法は、ニューロン回路において、ならびにかさばるキャパシタが電流積分器として利用される様々な他の適用例において適用され得る。さらに、シナプス104の各々は、メモrista要素に基づいて実装され得、シナプス重みの変化は、メモrista抵抗の変化に関係し得る。ナノメートルの特徴サイズのメモristaを用いると、ニューロン回路およびシナプスの面積が大幅に低減され得、それによって、大規模なニューラルシステムハードウェア実装形態の実装がより実用的になり得る。

20

【0023】

[0032]ニューラルシステム100をエミュレートするニューラルプロセッサの機能は、ニューロン間の結合の強さを制御し得る、シナプス結合の重みに依存し得る。シナプス重みは、パワーダウン後にプロセッサの機能を維持するために、不揮発性メモリに記憶され得る。一態様では、シナプス重みメモリは、主たるニューラルプロセッサチップとは別個の外部チップ上に実装され得る。シナプス重みメモリは、交換可能メモ리카ードとしてニューラルプロセッサチップとは別個にパッケージ化され得る。これは、ニューラルプロセッサに多様な機能を提供することができ、特定の機能は、ニューラルプロセッサに現在取り付けられているメモ리카ードに記憶されたシナプス重みに基づき得る。

30

【0024】

[0033]図2は、本開示のいくつかの態様による、計算ネットワーク(たとえば、ニューラルシステムまたはニューラルネットワーク)の処理ユニット(たとえば、ニューロンまたはニューロン回路)202の例示的な図200を示す。たとえば、ニューロン202は、図1のレベル102のニューロンおよび106のニューロンのうちのいずれかに対応し得る。ニューロン202は、ニューラルシステムの外部にある信号、または同じニューラルシステムの他のニューロンによって生成された信号、またはその両方であり得る、複数の入力信号 $204_1 \sim 204_N$ を受信し得る。入力信号は、電流、コンダクタンス、電圧、実数値および/または複素数値であり得る。入力信号は、固定小数点表現または浮動小数点表現をもつ数値を備え得る。これらの入力信号は、調整可能なシナプス重み $206_1 \sim 206_N$ ($w_1 \sim w_N$)に従って信号をスケーリングするシナプス結合を通してニューロン202に伝えられ得、Nはニューロン202の入力接続の総数であり得る。

40

【0025】

[0034]ニューロン202は、スケーリングされた入力信号を合成し、合成された、スケーリングされた入力を使用して、出力信号208(すなわち、信号y)を生成し得る。出力信号208は、電流、コンダクタンス、電圧、実数値および/または複素数値であり得る。出力信号は、固定小数点表現または浮動小数点表現をもつ数値であり得る。出力信号208は、次いで、同じニューラルシステムの他のニューロンへの入力信号として、また

50

は同じニューロン 202 への入力信号として、またはニューラルシステムの出力として伝達され得る。

【0026】

[0035] 処理ユニット（ニューロン）202 は電気回路によってエミュレートされ得、その入力接続および出力接続は、シナプス回路をもつ電気接続によってエミュレートされ得る。処理ユニット 202 ならびにその入力接続および出力接続はまた、ソフトウェアコードによってエミュレートされ得る。処理ユニット 202 はまた、電気回路によってエミュレートされ得るが、その入力接続および出力接続はソフトウェアコードによってエミュレートされ得る。一態様では、計算ネットワーク中の処理ユニット 202 はアナログ電気回路であり得る。別の態様では、処理ユニット 202 はデジタル電気回路であり得る。さらに別の態様では、処理ユニット 202 は、アナログ構成要素とデジタル構成要素の両方をもつ混合信号電気回路であり得る。計算ネットワークは、上述の形態のいずれかにおける処理ユニットを含み得る。そのような処理ユニットを使用した計算ネットワーク（ニューラルシステムまたはニューラルネットワーク）は、たとえば画像およびパターン認識、機械学習、モータ制御など、かなりの適用範囲において利用され得る。

10

【0027】

[0036] ニューラルネットワークをトレーニングする過程で、シナプス重み（たとえば、図 1 の重み

【0028】

【数 2】

20

$$w_1^{(i,i+1)}, \dots, w_P^{(i,i+1)}$$

【0029】

および / または図 2 の重み $w_{1 \sim 206_N}$ がランダム値により初期化され得、学習ルールに従って増加または減少し得る。学習ルールの例は、これに限定されないが、スパイクタイミング依存可塑性（STDP）学習ルール、Hebb 則、Oja 則、Bienenstock-Copper-Munro（BCM）則等を含むことを当業者は理解するだろう。いくつかの態様では、重みは、2 つの値のうちの 1 つに安定または収束し得る（すなわち、重みの双峰分布）。この効果が利用されて、シナプス重みごとのビット数を低減し、シナプス重みを記憶するメモリとの間の読取りおよび書込みの速度を上げ、シナプスメモリの電力および / またはプロセッサ消費量を低減し得る。

30

シナプスタイプ

【0030】

[0037] ニューラルネットワークのハードウェアおよびソフトウェアモデルでは、シナプス関係機能の処理がシナプスタイプに基づき得る。シナプスタイプは、非塑性シナプス（non-plastic synapse）（重みおよび遅延の変化がない）、可塑性シナプス（重みが増減し得る）、構造遅延可塑性シナプス（重みおよび遅延が増減し得る）、完全可塑性シナプス（重み、遅延および結合性が増減し得る）、およびその変形（たとえば、遅延は増減し得るが、重みまたは結合性の変化はない）であり得る。複数のタイプの利点は、処理が再分割され得ることである。たとえば、非塑性シナプスは、可塑性機能の実行を含まない場合がある（またはそのような機能が完了するのを待つ）。同様に、遅延および重み可塑性は、一緒にまたは別々に、順にまたは並列に動作し得る動作に再分割され得る。異なるタイプのシナプスは、適用される異なる可塑性タイプの各々の異なるルックアップテーブルまたは式およびパラメータを有し得る。したがって、本方法は、シナプスのタイプについての関連する表、式、またはパラメータにアクセスする。

40

【0031】

[0038] スパイクタイミング依存構造可塑性がシナプス可塑性とは無関係に実行され得るという事実のさらなる含意がある。構造可塑性は、重みの大きさに変化がない場合（たとえば、重みが最小値または最大値に達したか、あるいはそれが何らかの他の理由により変

50

更されない場合) s 構造可塑性 (すなわち、遅延量の変化) は前後スパイク時間差 (pre-post spike time difference) の直接関数であり得ても実行され得る。代替的に、構造可塑性は、重み変化量に応じて、または重みもしくは重み変化の限界に関係する条件に基づいて設定され得る。たとえば、重み変化が生じたとき、または重みが最大値になるのではなく、重みがゼロに達した場合のみ、シナプス遅延が変化し得る。しかしながら、これらのプロセスが並列化され、メモリアクセスの数および重複を低減し得るように、独立した機能を有することが有利であり得る。

シナプス可塑性の決定

【0032】

[0039] 神経可塑性 (または単に「可塑性」) は、脳内のニューロンおよびニューラルネットワークがそれらのシナプス結合と挙動とを新しい情報、感覚上の刺激、発展、損傷または機能不全に回答して変える能力である。可塑性は、生物学における学習および記憶にとって、また計算論的神経科学およびニューラルネットワークにとって重要である。(たとえば、Hebb 則理論による) シナプス可塑性、スパイクタイミング依存可塑性 (STDP)、非シナプス可塑性、アクティビティ依存可塑性、構造可塑性および恒常的可塑性など、様々な形の可塑性が研究されている。

【0033】

[0040] STDP は、ニューロン間のシナプス結合の強さを調整する学習プロセスである。結合強度は、特定のニューロンの出力スパイクおよび受信入力スパイク (すなわち、活動電位) の相対的タイミングに基づいて調整される。STDP プロセスの下で、あるニューロンに対する入力スパイクが、平均して、そのニューロンの出力スパイクの直前に生じる傾向がある場合、長期増強 (LTP) が生じ得る。その場合、その特定の入力はいくらも強くなる。一方、入力スパイクが、平均して、出力スパイクの直後に生じる傾向がある場合、長期抑圧 (LTD) が生じ得る。その場合、その特定の入力はいくらも弱くなるので、「スパイクタイミング依存可塑性」と呼ばれる。したがって、シナプス後ニューロンの興奮の原因であり得る入力は、将来的に寄与する可能性がさらに高くなる一方、シナプス後スパイクの原因ではない入力は、将来的に寄与する可能性が低くなる。結合の初期セットのサブセットが残る一方で、その他の部分の影響がわずかなレベルまで低減されるまで、このプロセスは続く。

【0034】

[0041] ニューロンは一般に出力スパイクを、その入力の多くが短い期間内に生じる (すなわち、出力をもたらすのに十分な累積がある) ときに生成するので、通常残っている入力のサブセットは、時間的に相関する傾向のあった入力を含む。さらに、出力スパイクの前に生じる入力は強化されるので、最も早い十分に累積的な相関指示を提供する入力は結局、ニューロンへの最終入力となる。

【0035】

[0042] STDP 学習ルールは、シナプス前ニューロンのスパイク時間 t_{pre} とシナプス後ニューロンのスパイク時間 t_{post} との間の時間差 (すなわち、 $t = t_{post} - t_{pre}$) に応じて、シナプス前ニューロンをシナプス後ニューロンに結合するシナプスのシナプス重みを効果的に適合させ得る。STDP の通常の公式化は、時間差が正である (シナプス前ニューロンがシナプス後ニューロンの前に発火する) 場合にシナプス重みを増加させ (すなわち、シナプスを増強し)、時間差が負である (シナプス後ニューロンがシナプス前ニューロンの前に発火する) 場合にシナプス重みを減少させる (すなわち、シナプスを抑制する) ことである。

【0036】

[0043] STDP プロセスでは、経時的なシナプス重みの変化は通常、以下の式によって与えられるように、指数関数的減衰を使用して達成され得る。

【0037】

10

20

30

40

【数 3】

$$\Delta w(t) = \begin{cases} a_+ e^{-t/k_+} + \mu, t > 0 \\ a_- e^{t/k_-}, t < 0 \end{cases}, \quad (1)$$

【0038】

ここで、 k_+ および $k_- = \text{sign}(-t)$ はそれぞれ、正の時間差および負の時間差の時間定数であり、 a_+ および a_- は対応するスケーリングの大きさであり、 μ は正の時間差および/または負の時間差に適用され得るオフセットである。

【0039】

[0044] 図3は、STDPによる、シナプス前スパイクおよびシナプス後スパイクの相対的タイミングに応じたシナプス重み変化の例示的な図300を示す。シナプス前ニューロンがシナプス後ニューロンの前に発火する場合、グラフ300の部分302に示すように、対応するシナプス重みは増加し得る。この重み増加は、シナプスのLTPと呼ばれ得る。グラフ部分302から、シナプス前スパイク時間とシナプス後スパイク時間との間の時間差に応じて、LTPの量がほぼ指数関数的に減少し得ることが観測され得る。グラフ300の部分304に示すように、発火の逆の順序は、シナプス重みを減少させ、シナプスのLTDをもたらし得る。

【0040】

[0045] 図3のグラフ300に示すように、STDPグラフのLTP(原因)部分302に負のオフセット μ が適用され得る。 x 軸の交差306のポイント($y=0$)は、レイヤ $i-1$ からの原因入力的相关を考慮して、最大タイムラグと一致するように構成され得る。フレームベースの入力(すなわち、スパイクまたはパルスを用意する特定の持続時間のフレームの形態である入力)の場合、オフセット値 μ は、フレーム境界を反映するように計算され得る。直接的にシナプス後電位によってモデル化されるように、またはニューラル状態に対する影響の点で、フレームにおける第1の入力スパイク(パルス)が経時的に減衰することが考慮され得る。フレームにおける第2の入力スパイク(パルス)が特定の時間フレームと関連したまたは特定の時間フレームに関連したものと考えられる場合、フレームの前および後の関連する時間は、その時間フレーム境界で分離され、関連する時間の値が異なり得る(たとえば、1つのフレームよりも大きい場合は負、1つのフレームよりも小さい場合は正)ように、STDP曲線の1つまたは複数の部分をオフセットすることによって、可塑性の点で別様に扱われ得る。たとえば、曲線が、フレーム時間よりも大きい前後の時間で実際にゼロよりも下になり、結果的にLTPの代わりにLTDの一部であるようにLTPをオフセットするために負のオフセット μ が設定され得る。

ニューロンモデルおよび演算

【0041】

[0046] 有用なスパイクニューロンモデルを設計するための一般的原理がいくつかある。良いニューロンモデルは、2つの計算レジーム、すなわち、一致検出および関数計算の点で豊かな潜在的挙動を有し得る。その上、良いニューロンモデルは、時間コーディングを可能にするための2つの要素を有する必要がある: 入力の到着時間は出力時間に影響を与え、一致検出は狭い時間ウィンドウを有し得る。最終的に、計算上魅力的であるために、良いニューロンモデルは、連続時間に閉形式解と、ニアアトラクター(near attractor)と鞍点とを含む安定した挙動とを有し得る。言い換えれば、有用なニューロンモデルは、実用的なニューロンモデルであり、豊かで、現実的で、生物学的に一貫した挙動をモデル化でき、神経回路のエンジニアリングとリバースエンジニアリングの両方が可能なニューロンモデルである。

【0042】

[0047] ニューロンモデルは事象、たとえば入力の到着、出力スパイク、または内部的であるか外部的であるかを問わず他の事象に依存し得る。豊かな挙動レパートリーを実現するために、複雑な挙動を示すことができる状態機械が望まれ得る。入力寄与(ある場合)

10

20

30

40

50

とは別個の事象の発生自体が状態機械に影響を与え、事象の後のダイナミクスを制限し得る場合、システムの将来の状態は、単なる状態および入力関数ではなく、むしろ状態、事象および入力関数である。

【0043】

[0048]一態様では、ニューロン n は、下記のダイナミクスによって決定される膜電圧 $v_n(t)$ によるスパイクグリキー積分発火ニューロンとしてモデル化され得る。

【0044】

【数4】

$$\frac{dv_n(t)}{dt} = \alpha v_n(t) + \beta \sum_m w_{m,n} y_m(t - \Delta t_{m,n}), \quad (2) \quad 10$$

【0045】

ここで α および β は、シナプス前ニューロン m をシナプス後ニューロン n に結合するシナプスのパラメータ、 $w_{m,n}$ はシナプス重みであり、 $y_m(t)$ は、ニューロン n の細胞体に到着するまで $t_{m,n}$ に従って樹状遅延または軸索遅延によって遅延し得るニューロン m のスパイク出力である。

【0046】

[0049]シナプス後ニューロンへの十分な入力達成された時間からシナプス後ニューロンが実際に発火する時間までの遅延があることに留意されたい。イジケヴィッチの単純モデルなど、動的スパイクニューロンモデルでは、脱分極しきい値 v_r とピークスパイク電圧 v_{peak} との間に差がある場合、時間遅延が生じ得る。たとえば、単純モデルでは、電圧および復元のための1対の微分方程式、すなわち、

【0047】

【数5】

$$\frac{dv}{dt} = (k(v - v_r)(v - v_r) - u + I) / C, \quad (3) \quad 20$$

$$\frac{du}{dt} = a(b(v - v_r) - u). \quad (4) \quad 30$$

【0048】

によってニューロン細胞体ダイナミクス (neuron soma dynamics) が決定され得る。ここでは膜電位であり、 u は、膜復元変数であり、 k は、膜電位の時間スケールを記述するパラメータであり、 a は、復元変数 u の時間スケールを記述するパラメータであり、 b は、膜電位のしきい値下変動に対する復元変数 u の感度を記述するパラメータであり、 v_r は、膜静止電位であり、 I は、シナプス電流であり、 C は、膜のキャパシタンスである。このモデルによれば、ニューロンは $v > v_{peak}$ のときにスパイクすると定義される。 40

Hunzinger Coldモデル

【0049】

[0050]Hunzinger Coldニューロンモデルは、豊かな様々な神経挙動を再生し得る最小二重レジームスパイク線形動的モデルである。モデルの1次元または2次元の線形ダイナミクスは2つのレジームを有することができ、時間定数(および結合)はレジームに依存し得る。しきい値下レジームでは、時間定数は、慣例により負であり、一般に生物学的に一貫した線形方式で静止状態に細胞を戻す役目を果たすリーキーチャネルダイナミクスを表す。しきい値上レジームにおける時間定数は、慣例により正であり、一般にスパイク生成のレイテンシを生じさせる一方でスパイク状態に細胞を駆り立てる反 50

リーキーチャネルダイナミクスを反映する。

【0050】

[0051] 図4に示すように、モデル400のダイナミクスは2つの（またはそれよりも多くの）レジームに分割され得る。これらのレジームは、負のレジーム（negative regime）402（leaky-integrate-and-fire（LIF）ニューロンモデルと混同されないように、交換可能にLIFレジームとも呼ばれる）、および正のレジーム（positive regime）404（anti-leaky-integrate-and-fire（ALIF）ニューロンモデルと混同されないように、交換可能にALIFレジームとも呼ばれる）と呼ばれ得る。負レジーム402では、状態は将来の事象の時点における静止（ v_{rest} ）の傾向がある。この負レジームでは、モデルは一般に、時間的入力検出特性と他のしきい値下挙動とを示す。正レジーム404では、状態はスパイクング事象（ v_{spike} ）の傾向がある。この正レジームでは、モデルは、後続の入力事象に応じてスパイクにレイテンシを生じさせるなどの計算特性を示す。事象の点からのダイナミクスの公式化およびこれら2つのレジームへのダイナミクスの分離は、モデルの基本的特性である。

10

【0051】

[0052] 線形二重レジーム2次元ダイナミクス（状態 v および u の場合）は、慣例により次のように定義され得る。

【0052】

【数6】

$$\tau_p \frac{dv}{dt} = v + q_p \quad (5)$$

$$-\tau_u \frac{du}{dt} = u + r \quad (6)$$

20

【0053】

ここで q および r は、結合のための線形変換変数である。

30

【0054】

[0053] シンボル v_{rest} は、ダイナミクスレジームを示すためにここで使用され、特定のレジームの関係を論述または表現するときに、それぞれ負レジームおよび正レジームについて符号「-」または「+」にシンボル v_{rest} を置き換える慣例がある。

【0055】

[0054] モデル状態は、膜電位（電圧） v および復元電流 u によって定義される。基本形態では、レジームは基本的にモデル状態によって決定される。正確で一般的な定義の微妙だが重要な側面があるが、差し当たり、モデルが、電圧 v がしきい値（ v_{spike} ）を上回る場合に正レジーム404にあり、そうでない場合に負レジーム402にあると考える。

【0056】

[0055] レジーム依存時間定数は、負レジーム時間定数である τ_p と正レジーム時間定数である τ_u とを含む。復元電流時間定数 τ_u は通常、レジームから独立している。便宜上、 τ_u と同様に、指数および τ_p が一般に正となる正レジームの場合に、電圧発展（voltage evolution）に関する同じ表現が使用され得るように、減衰を反映するために負の量として負レジーム時間定数 τ_u が一般に指定される。

40

【0057】

[0056] 2つの状態要素のダイナミクスは、事象において、ヌルクラインから状態をオフセットする変換によって結合され得、ここで変換変数は、

【0058】

50

【数 7】

$$q_p = -\tau_p \beta u - v_p \quad (7)$$

$$r = \delta(v + \varepsilon) \quad (8)$$

【0059】

であり、 v_+ 、 v_- 、 ε および β 、 τ_p はパラメータである。 v_+ のための2つの値は、2つのレジームのための参照電圧のベースである。パラメータ β は、負レジームのためのベース電圧であり、膜電位は一般に、負レジームにおいて v_- に減衰する。パラメータ τ_p は、正レジームのためのベース電圧であり、膜電位は一般に、正レジームにおいて v_+ から離れる傾向となる。 10

【0060】

[0057] および u のためのヌルクラインは、それぞれ変換変数 q および r の負によって与えられる。パラメータ β は、 u ヌルクラインの傾きを制御するスケール係数である。パラメータ τ_p は通常、 τ_p に等しく設定される。パラメータ τ_p は、両方のレジームにおいてヌルクラインの傾きを制御する抵抗値である。時間定数パラメータは、指数関数的減衰だけでなく、各レジームにおいて別個にヌルクラインの傾きを制御する。 20

【0061】

[0058] モデルは、電圧 v が値 v_+ に達したときにスパイクするように定義され得る。続いて、状態は（スパイク事象と同じ1つのものであり得る）リセット事象でリセットされ得る。

【0062】

【数 8】

$$v = \hat{v}_- \quad (9)$$

$$u = u + \Delta u \quad (10)$$

【0063】

ここで、

【0064】

【数 9】

$$\hat{v}_-$$

【0065】

および u はパラメータである。リセット電圧

【0066】

【数 10】

$$\hat{v}_-$$

【0067】

は通常、 v_+ にセットされる。

【0068】

[0059] 瞬時結合の原理によって、状態について（また、単一の指数項による）だけでは 50

なく、特定の状態に到達するための時間についても、閉形式解が可能である。閉形式状態解は、次のとおりである。

【 0 0 6 9 】

【 数 1 1 】

$$v(t + \Delta t) = (v(t) + q_\rho) e^{\frac{\Delta t}{\tau_\rho}} - q_\rho \quad (11)$$

$$u(t + \Delta t) = (u(t) + r) e^{\frac{\Delta t}{\tau_u}} - r \quad (12)$$

10

【 0 0 7 0 】

[0060]したがって、モデル状態は、入力（シナプス前スパイク）または出力（シナプス後スパイク）などの事象に伴ってのみ更新され得る。また、演算が（入力があるか、出力があるかを問わず）任意の特定の時間に行われ得る。

【 0 0 7 1 】

[0061]その上、瞬時結合原理によって、反復的技法または数値解法（たとえば、オイラー数値解法）なしに、特定の状態に到達する時間が事前に決定され得るように、シナプス後スパイクの時間が予想され得る。前の電圧状態 v_0 を踏まえ、電圧状態 v_f に到達するまでの時間遅延は、次の式によって与えられる。

20

【 0 0 7 2 】

【 数 1 2 】

$$\Delta t = \tau_\rho \log \frac{v_f + q_\rho}{v_0 + q_\rho} \quad (13)$$

【 0 0 7 3 】

[0062]スパイクが、電圧状態 v_s が v_+ に到達する時間に生じると定義される場合、電圧が所与の状態 v にある時間から測定されたスパイクが生じるまでの時間量、または相対的遅延に関する閉形式解は、次のとおりである。

30

【 0 0 7 4 】

【 数 1 3 】

$$\Delta t_s = \begin{cases} \tau_+ \log \frac{v_s + q_+}{v + q_+} & \text{もし、 } v > \hat{v}_+ \text{ であれば} \\ \infty & \text{そうでなければ} \end{cases} \quad (14)$$

【 0 0 7 5 】

ここで、

【 0 0 7 6 】

【 数 1 4 】

$$\hat{v}_+$$

【 0 0 7 7 】

は通常、パラメータ v_+ にセットされるが、他の変形も可能であり得る。

【 0 0 7 8 】

[0063]モデルダイナミクスの上記の定義は、モデルが正レジームにあるか、それとも負

50

レジームにあるかに依存する。上述のように、結合およびレジームは、事象に伴って計算され得る。状態の伝搬のために、レジームおよび結合（変換）変数は、最後の（前の）事象の時間における状態に基づいて定義され得る。続いてスパイク出力時間を予想するために、レジームおよび結合変数は、次の（最新の）事象の時間における状態に基づいて定義され得る。

【0079】

[0064] Coldモデルの、適時にシミュレーション、エミュレーションまたはモデルを実行するいくつかの可能な実装形態がある。これは、たとえば、事象更新モード、ステップ事象更新モード、およびステップ更新モードを含む。事象更新は、（特定の瞬間における）事象または「事象更新」に基づいて状態が更新される更新である。ステップ更新は、

10

間隔（たとえば、1ms）においてモデルが更新される更新である。これは必ずしも、反復的技法または数値解法を含むとは限らない。また、事象がステップもしくはステップ間で生じる場合または「ステップ事象」更新によってモデルを更新するのみによって、ステップベースのシミュレータにおいて限られた時間分解能で事象ベースの実装形態が可能である。

【0080】

[0065] 本開示の態様は、ニューラルネットワークシミュレータを対象とし、より詳細には、原位置ニューラルコプロセッシングを対象とする。

20

【0081】

[0066] 一般的に、ニューラルネットワークシミュレータは、柔軟性と性能（たとえば、シミュレータの電力）との間でトレードオフを行う。たとえば、設計者は、しばしば、学習を可能にするチップを作成するか、より高速に実行するチップを作成するか、または消費電力がより少ないチップを作成するかを決定しなければならない場合がある。したがって、学習がオフラインで実装されている場合、学習をサポートしないシミュレータ上に実装されるトレーニングされたニューラルネットワークは、学習をサポートするシミュレータ上に実装されたネットワークと同じ入力を経験しない場合がある。これは、学習に関連付けられるネットワークへのリアルタイムの変化が、ニューラルネットワークの環境に影響を与える（ニューラルネットワークに関連付けられるエフェクタを介して）可能性があるためであり得、それは、今度は環境を表現してネットワークへの入力を提供するセンサ

30

【0082】

[0067] 本開示の態様によれば、複数のシミュレーションプラットフォームは、シミュレータの通常動作中にトレードオフが行われ得るように組み合わせられ得る。たとえば、学習を利用しないシミュレーションは、この機能を提供しないシミュレーションプラットフォーム上で実行され得る。これは、たとえば、第2のシミュレーションプラットフォームが消費する電力が、第1のシミュレーションプラットフォームが消費する電力よりも少ない場合に有益であり得る。

40

【0083】

[0068] 本開示のいくつかの態様では、相互にスワッピングし得るニューラルコプロセッサが提供され得る。いくつかの態様では、ニューラルコプロセッサは、異なる機能を備えたニューラル処理ユニットまたはノードであり得る。たとえば、あるニューラル処理ノードは学習動作を実行するように構成され得、他の処理コアは静的重みで構成される。

【0084】

[0069] 1つの例示的な態様では、より多くの機能を備えたコア（すなわち、コアより多くの機能（たとえば、メモリまたはプロセッサを有するコア）は、より少ない機能を備えたコア（すなわち、より少ない機能を有するコア）の機能を引き継ぐ、または包含するこ

50

とができる。機能の包含は、処理ノードの「ホットスワップ (hot swap)」の形で行われ得る。この「ホットスワップ」を行うことにより、柔軟性と性能が向上され得る。

【0085】

[0070] 図5は、本開示の特定の態様による、汎用プロセッサ502を使用して、ニューラルネットワークにおける上述の実行しているコプロセッシングの例示的な実装形態500を示す。変数(ニューラル信号)、シナプス重み、計算ネットワーク(ニューラルネットワーク)に関連付けられるシステムパラメータ、遅延、周波数ピン情報、性能メトリック、およびシステム状態情報は、メモリブロック504に記憶され得、汎用プロセッサ502で実行される命令はプログラムメモリ506からロードされ得る。本開示のある態様では、汎用プロセッサ502にロードされる命令は、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングして、第1の処理ノードでニューラルネットワークの一部を実行して、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返して、および/または、第2の処理ノードでニューラルネットワークの一部を実行するためのコードを備え得る。

10

【0086】

[0071] 図6は、本開示のいくつかの態様による、メモリ602が相互接続ネットワーク604を介して計算ネットワーク(ニューラルネットワーク)の個々の(分散型)処理ユニット(ニューラルプロセッサ)606とインターフェースされ得る、ニューラルネットワーク内の上述したコプロセッシングの実行の例示的な実装形態600を示している。計算ネットワーク(ニューラルネットワーク)、遅延、周波数ピン情報、性能メトリック、およびシステム状態情報に関連付けられる、変数(ニューラル信号)、シナプス重み、システムパラメータはメモリ602に記憶され得、相互接続ネットワーク604の接続を介してメモリ602から各処理ユニット(ニューラルプロセッサ)606にロードされ得る。本開示のある態様では、処理ユニット606は、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングして、第1の処理ノードでニューラルネットワークの一部を実行して、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返して、および/または、第2の処理ノードでニューラルネットワークの一部を実行するように構成され得る。

20

【0087】

[0072] 図7は、ニューラルネットワーク内の上述したコプロセッシングの実行の例示的な実装形態700を示している。図7に示されるように、1つのメモリバンク702は、計算ネットワーク(ニューラルネットワーク)の1つの処理ユニット704に直接インターフェースされ得る。各メモリバンク702は、対応する処理ユニット(ニューラルプロセッサ)704、遅延、周波数ピン情報、性能メトリック、およびシステム状態情報に関連付けられる変数(ニューラル信号)、シナプス重み、および/またはシナプスパラメータを記憶し得る。本開示のある態様では、処理ユニット704は、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングして、第1の処理ノードでニューラルネットワークの一部を実行して、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返して、および/または、第2の処理ノードでニューラルネットワークの一部を実行するように構成され得る。

30

40

【0088】

[0073] 図8は、本開示のいくつかの態様による、ニューラルネットワーク800の例示的な実装形態を示す。図8に示すように、ニューラルネットワーク800は、本明細書に記載される方法の様々な動作を実行し得る複数のローカル処理ユニット802を有することができる。各ローカル処理ユニット802は、ニューラルネットワークのパラメータを記憶する、ローカルステートメモリ804およびローカルパラメータメモリ806を備え得る。また、ローカル処理ユニット802は、ローカルモデルプログラムを記憶するためのローカル(ニューロン)モデルプログラム(LMP)メモリ808、ローカル学習プログラムを記憶するためのローカル学習プログラム(LLP)メモリ810、およびローカ

50

ル接続メモリ 812 を有し得る。さらに、図 8 に示されるように、各ローカル処理ユニット 802 は、ローカル処理ユニットのローカルメモリの構成を提供するための構成処理ユニット 814 と、またローカル処理ユニット 802 間のルーティングを提供するルーティング接続処理ユニット 816 とインターフェースされ得る。

【0089】

[0074]一構成では、ニューロンモデルは、一定時間期間にわたって、ニューラルネットワークの一部を第 1 の処理ノードにスワッピングして、第 1 の処理ノードでニューラルネットワークの一部を実行して、一定時間期間後に、ニューラルネットワークの一部を第 2 の処理ノードに返して、および/または、第 2 の処理ノードでニューラルネットワークの一部を実行するために構成される。ニューロンモデルは、スワッピング手段と、第 1 の処理ノードでニューラルネットワークの一部を実行するための手段と、返す手段と、第 2 の処理ノードでニューラルネットワークの一部を実行するための手段とを含む。一態様では、スワッピング手段、第 1 の処理ノードでニューラルネットワークの一部を実行するための手段、返す手段、および/または、第 2 の処理ノードでニューラルネットワークの一部を実行するための手段は、記載された機能を実行するように構成された汎用プロセッサ 502、プログラムメモリ 506、メモリブロック 504、メモリ 602、相互接続ネットワーク 604、処理ユニット 606、処理ユニット 704、ローカル処理ユニット 802、およびまたはルーティング接続処理ユニット 816 であり得る。別の構成では、上述の手段は、上記の手段によって記載された機能を実行するように構成された任意のモジュールまたは任意の装置であり得る。

10

20

【0090】

[0075]別の構成では、ニューロンモデルは、まず第 1 の処理コアでニューラルネットワークの一部を実行することによって、および/または、さらなる実行のためにニューラルネットワークの一部を第 2 の処理コアに移動させることによって、オフライン学習を共同設置するために構成される。ニューロンモデルは、共同設置手段と移動手段とを含む。一態様では、共同設置手段および/または移動手段は、記載された機能を実行するように構成された汎用プロセッサ 502、プログラムメモリ 506、メモリブロック 504、メモリ 602、相互接続ネットワーク 604、処理ユニット 606、処理ユニット 704、ローカル処理ユニット 802、およびまたはルーティング接続処理ユニット 816 であり得る。別の構成では、上述の手段は、上記の手段によって記載された機能を実行するように構成された任意のモジュールまたは任意の装置であり得る。

30

【0091】

[0076]本開示のいくつかの態様によれば、各ローカル処理ユニット 802 は、ニューラルネットワークの所望の 1 つまたは複数の機能的特徴に基づいて、ニューラルネットワークのパラメータを決定して、決定されたパラメータがさらに適応され、同調され、更新されるにつれて、所望の機能的特徴に向けて 1 つまたは複数の機能的特徴を開発するように構成され得る。

【0092】

[0077]図 9 は、本開示の態様による、ニューラルネットワークの例示的なアーキテクチャ 900 を示すブロック図である。アーキテクチャ 900 は、処理ノード A 906 と処理ノード B 908 とを含み得るコプロセッサ 904 を備え得る。いくつかの態様では、処理ノード A 906 と処理ノード B 908 とは、同じハードウェアコア内に含まれ得る。しかしながら、これは単なる例示であり、処理ノード A 906 と処理ノード B 908 とは、代替で、別個のハードウェアコアにおいて提供され得る。

40

【0093】

[0078]処理ノード A 906 と処理ノード B 908 とは、異なるように構成され得る。すなわち、いくつかの態様では、処理ノード A 906 と処理ノード B 908 とは、ニューラルネットワークの機能的特徴を効率的に実行することに適した異なる構成を有し得る。いくつかの構成では、処理ノード A 906 は、処理ノード B よりも大きいリソースで構成され得る。たとえば、処理ノード A 906 は、処理ノード B 908 よりも速いおよび/また

50

は高い処理能力（たとえば、複数のプロセッサ、またはより速い処理速度）で構成され得る。第2の例では、処理ノードB908は、より多数および/またはより速いメモリで構成され得る。

【0094】

[0079]処理ノードA906と処理ノードB908とは、入力ノード902を介して入力を受信するように構成され得る。処理ノードA906と処理ノードB908とはまた、出力ノード910に出力を供給するように構成され得る。入力902と出力910とは、センサ、アクチュエータ、および他の入力/出力デバイスを備え得る。

【0095】

[0080]さらに、処理ノードA906と処理ノードB908とは、処理ノード間でニューラルネットワークの機能的特徴を実行することのホットスワッピングを可能にするために、相互に通信可能に結合され得る。すなわち、実行時に、より多くの機能を持つ処理ノード（たとえば、906、908）が、より少ない特徴を有するコアの機能の処理を含ままたは引き継ぎ得る。

10

【0096】

[0081]いくつかの態様では、処理ノードA906の状態がコピーされて、通信経路912または任意の他の通信経路を介して処理ノードB908に供給され得る。処理ノードA906の状態は、たとえば、状態変数、接続性情報、および他の状態情報を含み得る。

【0097】

[0082]処理ノードB908のリソースは、処理ノードA906からニューラルネットワークの機能的特徴の処理を引き継ぐために割り振られ得る。さらに、入力ノード902を介して供給された入力は、処理ノードB908にルーティングされ得る。処理ノードA906からの状態情報および入力に基づいて、処理ノードB908は、以前は処理ノードA906によって処理されていたニューラルネットワークの機能的特徴を処理することを引き継ぎ得る。

20

【0098】

[0083]いくつかの態様では、処理ノードA906は、入力ノード902を介して、処理ノードB908に供給されたものと同じ入力を受信することを継続し得る。したがって、整合性チェックを提供するために、処理ノードA906の出力が処理ノードB908の出力と比較され得る。一例では、処理ノードB908は、処理ノードA906内の欠陥またはバグを識別して減少させるためのデバッグコアとして構成され得る。本開示の他の態様では、処理ノードA906は、ニューラルネットワークの他の機能的特徴を処理し得る。

30

【0099】

[0084]処理ノードB908は、あらかじめ定められた時間期間にわたって、または、いくつかの態様では、特定のタスクまたはタスクのセットの完了まで、処理ノードA906から包含されたニューラルネットワークの一部の処理を継続し得る。たとえば、処理ノードB908は、学習を実装するように構成され得、また、学習が達成されるまで処理ノードA906から包含されたニューラルネットワークの一部を処理することを継続し得る。別の例では、処理ノードB908は、スパイクタイミング依存可塑性を実装するように構成され得る。したがって、処理ノードBは、受信された状態情報を処理して、状態情報の更新（たとえば、重み更新）が決定されるまで可塑性ルールを適用し得る。

40

【0100】

[0085]いくつかの態様では、より多くの機能を備えた処理ノード（たとえば、906、908）は、システム性能メトリックに基づいて処理を引き継ぎ得る。たとえば、より多くの機能を備えた処理ノードは、より少ない機能を備えた処理ノードのシステム性能が低い値レベルを下回る場合、処理を包含し得る。他の態様では、電力がシステムに適用されるとスワッピングが実行され得る。もちろん、これらは単に例示的な基礎であり、他のシステムおよびネットワーク性能メトリックは、より少ない機能を備えた処理ノードから、より多くの機能を備えた処理ノードに、スワッピング処理のための基礎を提供し得る。

50

【0101】

[0086]タスクが完了すると、または時間期間が満了すると、処理ノードB908の状態がコピーされて、変更されたコアとして処理ノードA906に供給され得る。いくつかの態様では、ニューラルネットワークの一部を返すことは、システム性能メトリックに基づいて実行され得る。たとえば、システム性能がしきい値を上回る場合、処理ノードB908の状態がコピーされて、処理ノードA906に供給され得る。第2の例では、返すことは、電力がシステム（たとえば、プラグインシステム）に適用されると発生し得る。いくつかの態様では、入力ノード902を介して提供される入力は、処理ノードB908からの状態情報を含む変更されたコアを使用してニューラルネットワークの機能的特徴を処理することを再開するために処理ノードA906にルーティングされ得る。

10

【0102】

[0087]図10A～図10Fは、本開示の態様による、ニューラルネットワークにおける原位置コプロセッシングを示す例示的なブロック図1000である。例示的なブロック図の各々は、静的コア1008と学習コア1006とを含むコプロセッサ1004を示す。静的コア1008は、ニューラルネットワークまたはその一部を動作することに関連付けられる機能を実行するための静的重みで構成され得る。学習コア1006は、学習を実装して、学習動作を実行するように構成され得る。たとえば、いくつかの態様では、学習コア1006は、強化学習または他の学習モデルを実装するように構成され得る。

【0103】

[0088]いくつかの態様では、学習コア1006は、静的コア1008よりも大きいリソースで構成され得る。たとえば、学習コア1006は、静的コア1008よりも速いおよび/または多数の処理能力（たとえば、複数のプロセッサ、またはより速い処理速度）で構成され得る。別の例では、学習コア1006は、静的コア1008とは異なるメモリリソース（たとえば、より多数および/またはより速いメモリ）で構成され得る。異なるタイプのメモリリソースは、たとえば、パラメータ（たとえば、重み）に関してより高い（または、より低い）精度を可能にしてもよく、スパイク履歴をキャプチャするためのより多くのリソースを提供してもよく、学習ルールへのアクセス、ならびにスパイクタイミング依存可塑性および/またはビット割振りの実装を可能にしてもよい。もちろん、これらの処理および性能関連機能は単なる例示であり、他の処理および性能関連機能または強化は、学習コア1006および静的コア1008に異なるように含まれ得る。

20

30

【0104】

[0089]図10A～図10Fに含まれるブロック図の各々は1つだけの静的コア1008および学習コア1006を示すが、これは単なる例示であり、説明を容易にするためである。代わりに、たとえば設計効率の目的のために、任意の数の静的コア1008および学習コア1006が含まれ得る。さらに、静的コア1008および学習コア1006は、同一の処理コア内に含まれてもよく、代替で、別個の処理コアにおいて提供されてもよい。

【0105】

[0090]静的コア1008および学習コア1006は、入力ノード1002を介して入力を選択的に受信して、出力ノード1010に出力を供給し得る。いくつかの態様では、静的コア1008と学習コア1006との両方は、入力ノード1002を介して入力を受信し得る。同様に、静的コア1008と学習コア1006との両方は、整合性チェックまたは処理検証を可能にするために、出力を出力ノード1010に供給し得る。

40

【0106】

[0091]図10Aで、入力ノード1002からの入力が静的コア1008に提供されるが、学習コア1006には提供されない。この例示的な態様では、ニューラルネットワークの動作は、静的コア1008を介して実行のために合理化され得る。いくつかの態様では、学習は実装され得ない。

【0107】

[0092]図10Bで、静的コア1008の状態情報がコピーされて、通信経路1012を介して学習コア1006に提供され得る。状態情報は、たとえば、ニューロン状態変数、

50

シナプス状態情報、接続性情報（たとえば、図または表）、および重み情報を含み得る。

【0108】

[0093]図10Cで、入力ノード1002を介する入力が学習コア1006にルーティングされ得る。いくつかの態様では、入力は学習コア1006だけに提供され得る。もちろん、入力は、代替で学習コア1006と静的コア1008との両方に提供され得る。この構成では、たとえば、検証技法は、静的コア1008からの出力と学習コア1006からの出力とが一致している（たとえば、同一である）ことを保証するために実行され得る。

【0109】

[0094]図10Dで、学習コア1006は、静的コア1008によって実行されていたニューラルネットワーク（または、その一部）に関連付けられる処理機能を包含する、または引き継ぐ。学習コア1006は、あらかじめ定められた時間期間にわたって、あるいは特定のタスクまたは機能の実行の間に、処理を引き継ぎ得る。たとえば、いくつかの態様では、学習コア1006は、STDP、あるいはニューラルネットワークまたはその一部に関連する強化学習などの学習モデルを実装するために、より少ない機能を備えた静的コア1008から処理を引き継ぎ得る。

10

【0110】

[0095]別の例では、学習コア1006によって処理が包含されるニューラルネットワークの一部は、深層信念ネットワークのレイヤであり得る。深層信念ネットワークは、確率的潜在変数の複数のレイヤからなる確率的生成モデルである。深層信念ネットワークでは、学習は、たとえば、トップダウン様式で、レイヤごとに実装され得る。

20

【0111】

[0096]学習はオンラインで実装されてもよく、オフラインで実装されてもよい。オフライン学習が発生すると、学習コア1006の入力（たとえば、1002）および出力（たとえば、1010）は、ニューラルネットワークの他のレイヤを備え得る。さらに、学習コア1006の入力（たとえば、1002）および出力（たとえば、1010）はまた、センサ、アクチュエータ等を備え得る。

【0112】

[0097]いくつかの態様では、静的コア1008は、入力の受信を継続し得る。たとえば、静的コア1008は、教師付き学習を可能にするために監視コアとして動作され得る。したがって、静的コア1008の出力は、学習コア1006をトレーニングし得る。他の態様では、静的コア1008は入力の受信を継続し得、また、ニューラルネットワークまたはその一部の動作に関連付けられる他のタスクを実行するよう割り当てられ得る。他の態様では、静的コア1008は、入力の受信を停止し得る。

30

【0113】

[0098]図10Eで、あらかじめ定められた時間期間の満了後、あるいはタスクまたは実行された機能が完了する（たとえば、学習が達成される）と、学習コア1006は処理制御の静的コア1008への返却を開始し得る。学習コア1006の状態情報がコピーされて、通信経路1012を介して静的コア1008に供給され得る。いくつかの態様では、学習コア1006の状態情報は、静的コア1008の異なるインスタンスを備え得る。たとえば、異なるインスタンスは、達成された学習に基づいて拡張された、変更された静的コア1008であり得る。別の例では、変更された静的コア1008は、STDPルールの実装形態に基づいて静的重みの更新を含み得る。

40

【0114】

[0099]図10Fで、学習コア1006は、学習コア1006からの状態情報に基づいて、ニューラルネットワークまたはその一部の動作に関連付けられる機能の実行を再開するために、制御を静的コア1008に返す。

【0115】

[00100]図11は、ニューラルネットワークにおいてコプロセッシングを実行するための方法1100を示す。ブロック1102で、ニューロンモデルは、一定時間期間にわたって、ニューラルネットワークの一部を第1の処理ノードにスワッピングする。ブロック

50

1104で、ニューロンモデルは、第1の処理ノードでニューラルネットワークの一部を実行する。ブロック1106で、ニューロンモデルは、一定時間期間後に、ニューラルネットワークの一部を第2の処理ノードに返す。さらに、ブロック1108で、ニューロンモデルは、第2の処理ノードでニューラルネットワークの一部を実行する。

【0116】

[00101]図12は、ニューラルネットワークにおいてコプロセッシングを実行するための方法1200を示す。ブロック1202で、ニューロンモデルは、まず第1の処理コアでニューラルネットワークの一部を実行することによって、オフライン学習を共同設置する。ブロック1204で、ニューロンモデルは、さらなる実行のためにニューラルネットワークの一部を第2の処理コアに移動させる。

10

【0117】

[00102]上述した方法の様々な動作は、対応する機能を実行することが可能な任意の好適な手段によって実行され得る。それらの手段は、限定はしないが、回路、特定用途向け集積回路(AASIC)、またはプロセッサを含む、様々なハードウェアおよび/またはソフトウェア構成要素および/またはモジュールを含み得る。概して、図に示されている動作がある場合、それらの動作は、同様の番号をもつ対応するカウンターパートのミーンズプラスファンクション構成要素を有し得る。

【0118】

[00103]本明細書で使用する「決定」という用語は、多種多様なアクションを包含する。たとえば、「決定」は、計算すること、算出すること、処理すること、導出すること、調査すること、ルックアップすること(たとえば、テーブル、データベースまたは別のデータ構造においてルックアップすること)、確認することなどを含み得る。さらに、「決定」は、受信すること(たとえば、情報を受信すること)、アクセスすること(たとえば、メモリ中のデータにアクセスすること)などを含み得る。さらに、「決定」は、解決すること、選択すること、選定すること、確立することなどを含み得る。

20

【0119】

[00104]本明細書で使用する、項目のリスト「のうちの少なくとも1つ」を指す句は、単一のメンバーを含む、それらの項目の任意の組合せを指す。一例として、「a、b、またはcのうちの少なくとも1つ」は、a、b、c、a-b、a-c、b-c、およびa-b-cを包含するものとする。

30

【0120】

[00105]本開示および付録Aに関連して説明した様々な例示的な論理ブロック、モジュール、および回路は、汎用プロセッサ、デジタル信号プロセッサ(DSP)、特定用途向け集積回路(AASIC)、フィールドプログラマブルゲートアレイ信号(FPGA)または他のプログラマブル論理デバイス(PLD)、個別ゲートまたはトランジスタ論理、個別ハードウェア構成要素、あるいは本明細書で説明した機能を実行するように設計されたそれらの任意の組合せを用いて実装または実行され得る。汎用プロセッサはマイクロプロセッサであり得るが、代替として、プロセッサは、任意の市販のプロセッサ、コントローラ、マイクロコントローラまたは状態機械であり得る。プロセッサはまた、コンピューティングデバイスの組合せ、たとえば、DSPとマイクロプロセッサとの組合せ、複数のマイクロプロセッサ、DSPコアと連携する1つまたは複数のマイクロプロセッサ、あるいは任意の他のそのような構成として実装され得る。

40

【0121】

[00106]本開示および付録Aに関連して説明した方法またはアルゴリズムのステップは、ハードウェアで直接実施されるか、プロセッサによって実行されるソフトウェアモジュールで実施されるか、またはその2つの組合せで実施され得る。ソフトウェアモジュールは、当技術分野で知られている任意の形式の記憶媒体で存在し得る。使用され得る記憶媒体のいくつかの例は、ランダムアクセスメモリ(RAM)、読出し専用メモリ(ROM)、フラッシュメモリ、消去可能プログラマブル読出し専用メモリ(EPROM)、電気的消去可能プログラマブル読出し専用メモリ(EEPROM(登録商標))、レジスタ、ハ

50

ードディスク、リムーバブルディスク、CD-ROMなどを含む。ソフトウェアモジュールは、単一の命令、または多数の命令を備えることができ、いくつかの異なるコードセグメント上で、異なるプログラム間で、複数の記憶媒体にわたって分散され得る。記憶媒体は、プロセッサがその記憶媒体から情報を読み取ることができ、その記憶媒体に情報を書き込むことができるように、プロセッサに結合され得る。代替として、記憶媒体はプロセッサと一体化され得る。

【0122】

[00107]本明細書で開示する方法は、説明した方法を達成するための1つまたは複数のステップまたはアクションを備える。本方法のステップおよび/またはアクションは、特許請求の範囲から逸脱することなく互いに交換され得る。言い換えれば、ステップまたはアクションの特定の順序が指定されない限り、特定のステップおよび/またはアクションの順序および/または使用は、特許請求の範囲から逸脱することなく変更され得る。

10

【0123】

[00108]本明細書で説明した機能は、ハードウェア、ソフトウェア、ファームウェア、またはそれらの任意の組合せで実装され得る。ハードウェアで実装される場合、例示的なハードウェア構成はデバイス中に処理システムを備え得る。処理システムは、バスアーキテクチャを用いて実装され得る。バスは、処理システムの特定の適用例および全体的な設計制約に応じて、任意の数の相互接続バスとブリッジとを含み得る。バスは、プロセッサと、機械可読媒体と、バスインターフェースとを含む様々な回路を互いにリンクし得る。バスインターフェースは、ネットワークアダプタを、特に、バスを介して処理システムに接続し得る。ネットワークアダプタは、信号処理機能を実装し得る。いくつかの態様では、ユーザインターフェース(たとえば、キーボード、ディスプレイ、マウス、ジョイスティックなど)もバスに接続され得る。バスはまた、タイミングソース、周辺機器、電圧調整器、電力管理回路などの様々な他の回路にリンクし得るが、それらは当技術分野でよく知られており、したがってこれ以上は説明されない。

20

【0124】

[00109]プロセッサは、機械可読媒体に記憶されたソフトウェアの実行を含む、バスおよび一般的な処理を管理することを担当し得る。プロセッサは、1つまたは複数の汎用および/または専用プロセッサを用いて実装され得る。例としては、マイクロプロセッサ、マイクロコントローラ、DSPプロセッサ、およびソフトウェアを実行し得る他の回路を含む。ソフトウェアは、ソフトウェア、ファームウェア、ミドルウェア、マイクロコード、ハードウェア記述言語などの名称にかかわらず、命令、データ、またはそれらの任意の組合せを意味すると広く解釈されたい。機械可読媒体は、一例として、ランダムアクセスメモリ(RAM)、フラッシュメモリ、読出し専用メモリ(ROM)、プログラマブル読出し専用メモリ(PROM)、消去可能プログラマブル読出し専用メモリ(EPROM)、電氣的消去可能プログラム可能読出し専用メモリ(EEPROM)、レジスタ、磁気ディスク、光ディスク、ハードドライブ、または他の任意の適切な記憶媒体、あるいはそれらの任意の組合せを含み得る。機械可読媒体はコンピュータプログラム製品において実施され得る。コンピュータプログラム製品はパッケージング材料を備え得る。

30

【0125】

[00110]ハードウェア実装形態では、機械可読媒体は、プロセッサとは別個の処理システムの一部であり得る。しかしながら、当業者なら容易に理解するように、機械可読媒体またはその任意の部分は処理システムの外部にあり得る。例として、機械可読媒体は、すべてバスインターフェースを介してプロセッサによってアクセスされ得る、伝送線路、データによって変調された搬送波、および/またはデバイスとは別個のコンピュータ製品を含み得る。代替的に、または追加で、機械可読媒体またはその任意の部分は、キャッシュおよび/または汎用レジスタファイルがそうであり得るように、プロセッサに統合され得る。論じた様々な構成要素は、ローカル構成要素などの特定の位置を有するものとして説明され得るが、それらはまた、分散コンピューティングシステムの一部として構成されているいくつかの構成要素などの様々な方法で構成され得る。

40

50

【 0 1 2 6 】

[00111]処理システムは、すべて外部バスアーキテクチャを介して他のサポート回路と互いにリンクされる、プロセッサ機能を提供する1つまたは複数のマイクロプロセッサと、機械可読媒体の少なくとも一部分を提供する外部メモリとをもつ汎用処理システムとして構成され得る。あるいは、処理システムは、本明細書に記載のニューロンモデルとニューラルシステムのモデルとを実装するための1つまたは複数のニューロモルフィックプロセッサを備え得る。別の代替として、処理システムは、プロセッサを有する特定用途向け集積回路（ASIC）と、バスインターフェースと、ユーザインターフェースと、サポート回路と、単一のチップに統合された機械可読媒体の少なくとも一部とを用いて、あるいは1つまたは複数のフィールドプログラマブルゲートアレイ（FPGA）、プログラマブル論理デバイス（PLD）、コントローラ、状態機械、ゲート論理、個別ハードウェア構成要素、または他の任意の適切な回路、あるいは本開示全体を通じて説明した様々な機能を実行し得る回路の任意の組合せを用いて実装され得る。当業者なら、特定の適用例と、全体的なシステムに課される全体的な設計制約とに応じて、どのようにしたら処理システムについて説明した機能を最も良く実装し得るかを理解されよう。

10

【 0 1 2 7 】

[00112]機械可読媒体はいくつかのソフトウェアモジュールを備え得る。ソフトウェアモジュールは、プロセッサによって実行されたときに、処理システムに様々な機能を実行させる命令を含む。ソフトウェアモジュールは、送信モジュールと受信モジュールとを含み得る。各ソフトウェアモジュールは、単一の記憶デバイス中に常駐するか、または複数の記憶デバイスにわたって分散され得る。例として、トリガイベントが発生したとき、ソフトウェアモジュールがハードドライブからRAMにロードされ得る。ソフトウェアモジュールの実行中、プロセッサは、アクセス速度を高めるために、命令のいくつかをキャッシュにロードし得る。次いで、1つまたは複数のキャッシュラインが、プロセッサによる実行のために汎用レジスタファイルにロードされ得る。以下でソフトウェアモジュールの機能に言及する場合、そのような機能は、そのソフトウェアモジュールからの命令を実行したときにプロセッサによって実装されることが理解されよう。

20

【 0 1 2 8 】

[00113]ソフトウェアで実装される場合、機能は、1つまたは複数の命令またはコードとしてコンピュータ可読媒体上に記憶されるか、あるいはコンピュータ可読媒体を介して送信され得る。コンピュータ可読媒体は、ある場所から別の場所へのコンピュータプログラムの転送を可能にする任意の媒体を含む、コンピュータ記憶媒体と通信媒体の両方を含む。記憶媒体は、コンピュータによってアクセスされ得る任意の利用可能な媒体であり得る。限定ではなく例として、そのようなコンピュータ可読媒体は、RAM、ROM、EEPROM、CD-ROMまたは他の光ディスクストレージ、磁気ディスクストレージまたは他の磁気記憶デバイス、あるいは命令またはデータ構造の形態の所望のプログラムコードを搬送または記憶し得、コンピュータによってアクセスされ得る、任意の他の媒体を備えることができる。さらに、いかなる接続もコンピュータ可読媒体を適切に名づけられる。たとえば、ソフトウェアが、同軸ケーブル、光ファイバーケーブル、ツイストペア、デジタル加入者回線（DSL）、または赤外線（IR）、無線、およびマイクロ波などのワイヤレス技術を使用して、ウェブサイト、サーバ、または他のリモートソースから送信される場合、同軸ケーブル、光ファイバーケーブル、ツイストペア、DSL、または赤外線、無線、およびマイクロ波などのワイヤレス技術は、媒体の定義に含まれる。本明細書で使用するディスク（disk）およびディスク（disc）は、コンパクトディスク（disc）（CD）、レーザーディスク（登録商標）（disc）、光ディスク（disc）、デジタル多用途ディスク（disc）（DVD）、フロッピー（登録商標）ディスク（disk）、およびBlu-ray（登録商標）ディスク（disc）を含み、ディスク（disk）は、通常、データを磁気的に再生し、ディスク（disc）は、データをレーザーで光学的に再生する。したがって、いくつかの態様では、コンピュータ可読媒体は非一時的コンピュータ可読媒体（たとえば、有形媒体）を備え得る。さらに、他の態様では、コンピュータ可読媒体は一時的コンピ

30

40

50

ユーザ可読媒体（たとえば、信号）を備え得る。上記の組合せもコンピュータ可読媒体の範囲内に含まれるべきである。

【0129】

[00114]したがって、いくつかの態様は、本明細書で提示する動作を実行するためのコンピュータプログラム製品を備え得る。たとえば、そのようなコンピュータプログラム製品は、本明細書で説明する動作を実行するために1つまたは複数のプロセッサによって実行可能である命令を記憶した（および/または符号化した）コンピュータ可読媒体を備え得る。いくつかの態様では、コンピュータプログラム製品はパッケージング材料を含み得る。

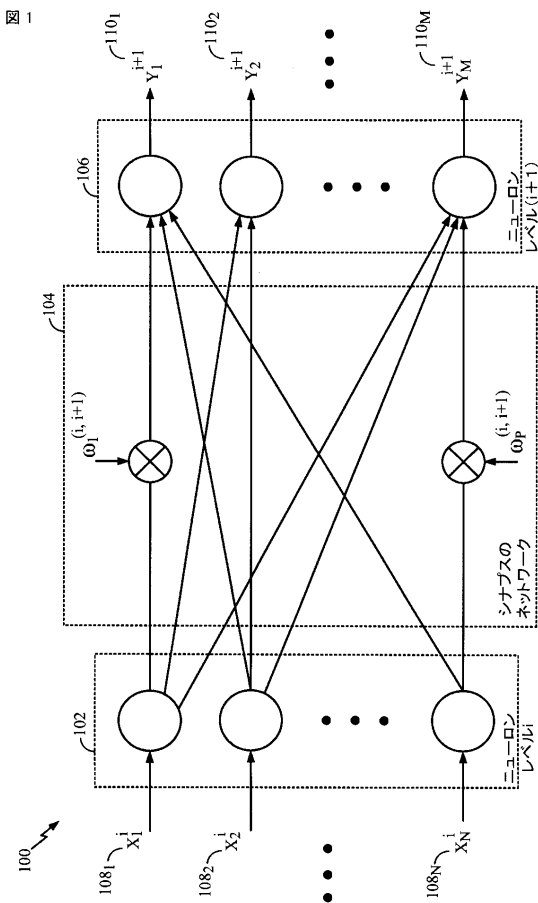
【0130】

[00115]さらに、本明細書で説明した方法および技法を実行するためのモジュールおよび/または他の適切な手段は、適用可能な場合にユーザ端末および/または基地局によってダウンロードされ、および/または他の方法で取得され得ることを諒解されたい。たとえば、そのようなデバイスは、本明細書で説明した方法を実施するための手段の転送を可能にするためにサーバに結合され得る。代替的に、本明細書で説明した様々な方法は、ユーザ端末および/または基地局が記憶手段をデバイスに結合または提供すると様々な方法を得ることができるよう、記憶手段（たとえば、RAM、ROM、コンパクトディスク（CD）またはフロッピーディスクなどの物理記憶媒体など）によって提供され得る。その上、本明細書で説明した方法および技法をデバイスに与えるための任意の他の好適な技法が利用され得る。

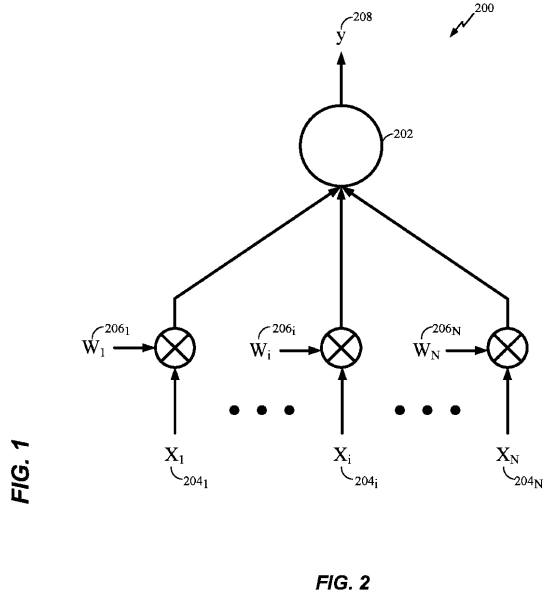
【0131】

[00116]特許請求の範囲は、上記で示した厳密な構成および構成要素に限定されないことを理解されたい。上記で説明した方法および装置の構成、動作および詳細において、特許請求の範囲から逸脱することなく、様々な改変、変更および変形が行われ得る。

【図1】



【図2】



10

20

【 図 3 】

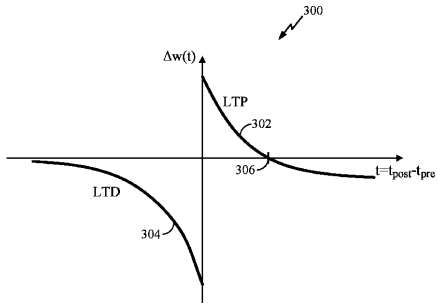


FIG. 3

【 図 4 】

図 4

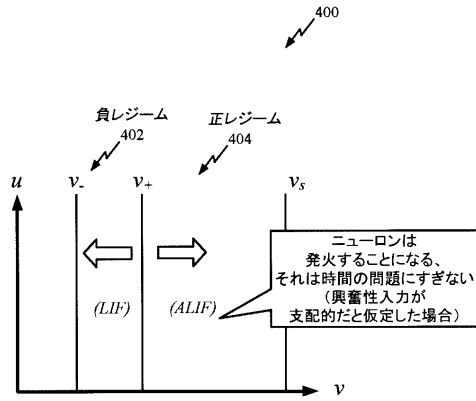


FIG. 4

【 図 5 】

図 5

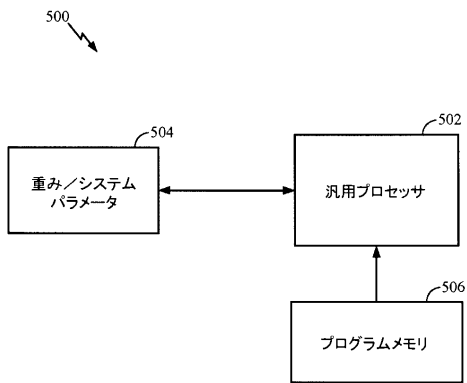


FIG. 5

【 図 6 】

図 6

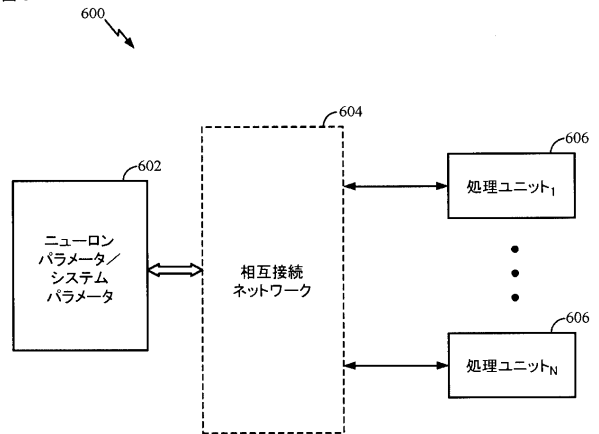


FIG. 6

【 図 7 】

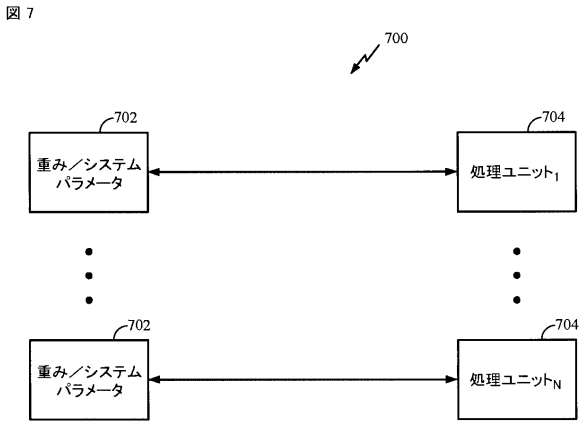


FIG. 7

【 図 8 】

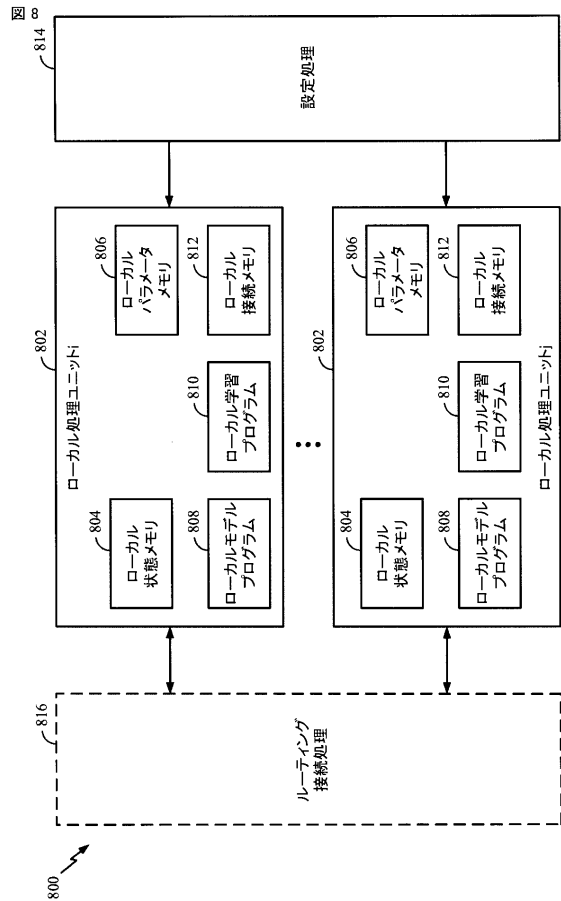


FIG. 8

【 図 9 】

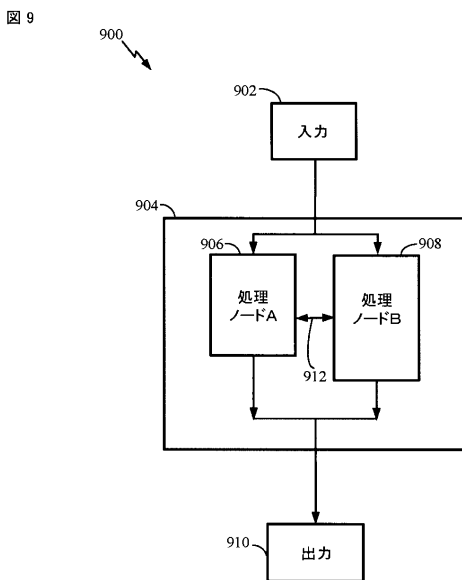


FIG. 9

【 図 10 A 】

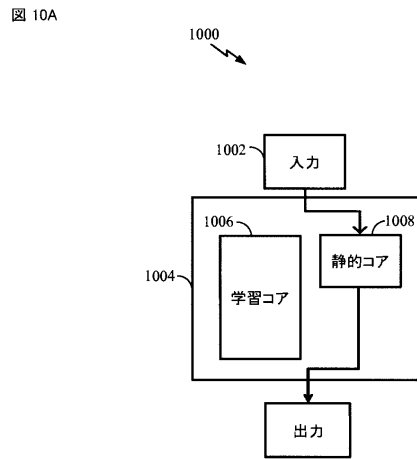


FIG. 10A

【図 10 B】

図 10B

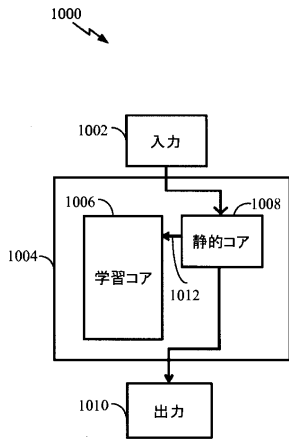


FIG. 10B

【図 10 C】

図 10C

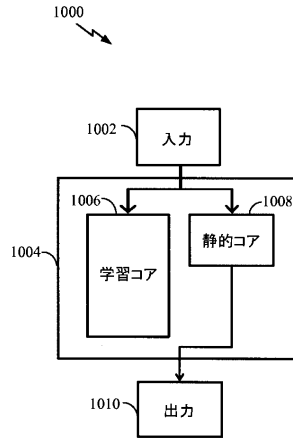


FIG. 10C

【図 10 D】

図 10D

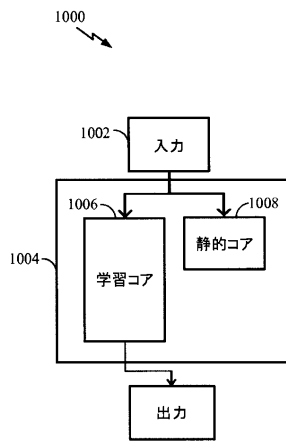


FIG. 10D

【図 10 E】

図 10E

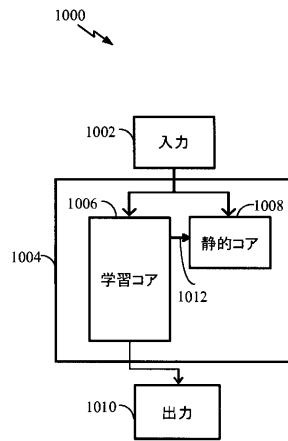


FIG. 10E

【 図 1 0 F 】

図 10F

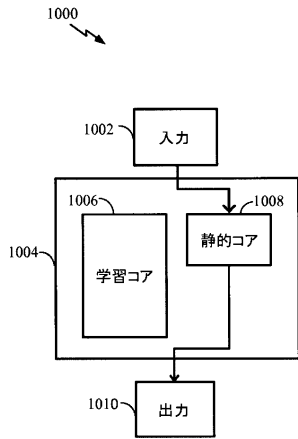


FIG. 10F

【 図 1 1 】

図 11

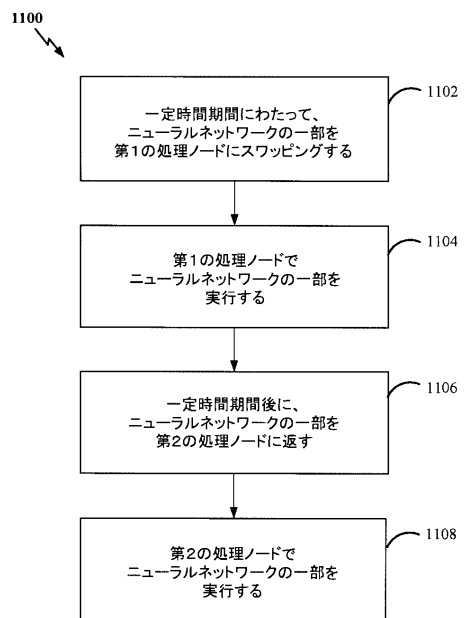


FIG. 11

【 図 1 2 】

図 12

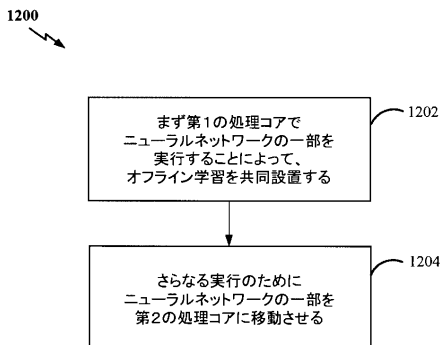


FIG. 12

【手続補正書】

【提出日】平成28年10月20日(2016.10.20)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

人工ニューラルネットワークにおいてコプロセッシングを実行する方法であって、一定時間期間にわたって、前記ニューラルネットワークの一部を第2の処理ノードから第1の処理ノードにスワッピングすることと、

前記第1の処理ノードで前記ニューラルネットワークの前記一部を実行することと、
前記一定時間期間後に、前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すことと、

前記第2の処理ノードで前記ニューラルネットワークの前記一部を実行することとを備える、方法。

【請求項2】

前記第1の処理ノードは、前記第2の処理ノードとは別個のハードウェアコアを備える

、
請求項1に記載の方法。

【請求項3】

前記第1の処理ノードは、学習処理コアを備える、
請求項1に記載の方法。

【請求項4】

前記学習処理コアは、前記第2の処理ノードよりも多くのリソースで構成される、
請求項3に記載の方法。

【請求項5】

学習は、オフラインで実装される、
請求項3に記載の方法。

【請求項6】

前記第1の処理ノードは、学習処理コアを備え、
前記第2の処理ノードは、静的処理コアを備え、
スワッピングすることは、

前記静的処理コアの状態を前記学習処理コアにコピーすることと、

前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに
入力をルーティングすることと

を備え、

返すことは、

前記学習処理コアの状態を前記静的処理コアにコピーすることと、

変更された静的処理コアに制御を返すことと

を備える、請求項1に記載の方法。

【請求項7】

前記スワッピングすることは、前記第2の処理ノードから前記人工ニューラルネットワークの前記一部を処理するための前記第1の処理ノードのリソースを割り振ることを備える、

請求項1に記載の方法。

【請求項8】

前記人工ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、

請求項 1 に記載の方法。

【請求項 9】

前記第 1 の処理ノードは、デバッグコアを備える、
請求項 1 に記載の方法。

【請求項 10】

前記スワッピングすることは、システム性能がしきい値を下回る場合にトリガされる、
請求項 1 に記載の方法。

【請求項 11】

前記返すことは、システム性能がしきい値を上回る場合にトリガされる、
請求項 1 に記載の方法。

【請求項 12】

前記スワッピングすること、または返すことは、電力がシステムに適用されるとトリガされる、

請求項 1 に記載の方法。

【請求項 13】

人工ニューラルネットワークにおいてコプロセッシングを実行するための装置であって

メモリと、

前記メモリに結合された少なくとも 1 つのプロセッサと

を備え、前記少なくとも 1 つのプロセッサは、

一定時間期間にわたって、前記ニューラルネットワークの一部を第 2 の処理ノードから
第 1 の処理ノードにスワッピングすることと、

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行することと、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を前記第 2 の処理ノード
に返すことと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行することと
を行うように構成される、装置。

【請求項 14】

前記第 1 の処理ノードは、前記第 2 の処理ノードとは別個のハードウェアコアを備える

、
請求項 13 に記載の装置。

【請求項 15】

前記第 1 の処理ノードは、学習処理コアを備える、

請求項 13 に記載の装置。

【請求項 16】

前記学習処理コアは、前記第 2 の処理ノードよりも多くのリソースで構成される、

請求項 15 に記載の装置。

【請求項 17】

学習は、オフラインで実装される、

請求項 15 に記載の装置。

【請求項 18】

前記第 1 の処理ノードは、学習処理コアを備え、前記第 2 の処理ノードは、静的処理コ
アを備え、前記少なくとも 1 つのプロセッサは、

前記静的処理コアの状態を前記学習処理コアにコピーすることと、

前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに
入力をルーティングすることと、

前記学習処理コアの状態を前記静的処理コアにコピーすることと、

変更された静的処理コアに制御を返すことと

を行うようにさらに構成される、請求項 13 に記載の装置。

【請求項 19】

前記少なくとも1つのプロセッサは、前記第2の処理ノードから前記人工ニューラルネットワークの前記一部を処理するための前記第1の処理ノードのリソースを割り振ることを行うようにさらに構成される、

請求項13に記載の装置。

【請求項20】

前記人工ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、

請求項13に記載の装置。

【請求項21】

前記第1の処理ノードは、デバッグコアを備える、

請求項13に記載の装置。

【請求項22】

前記少なくとも1つのプロセッサは、システム性能がしきい値を下回る場合に、前記ニューラルネットワークの前記一部を前記第1の処理ノードにスワッピングするようにさらに構成される、

請求項13に記載の装置。

【請求項23】

前記少なくとも1つのプロセッサは、システム性能がしきい値を上回る場合に、前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すようにさらに構成される、

請求項13に記載の装置。

【請求項24】

前記少なくとも1つのプロセッサは、電力がシステムに適用されると、前記ニューラルネットワークの前記一部を前記第1の処理ノードにスワッピングする、または前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すようにさらに構成される、

請求項13に記載の装置。

【請求項25】

人工ニューラルネットワークにおいてコプロセッシングを実行するための装置であって、

一定時間期間にわたって、前記ニューラルネットワークの一部を第2の処理ノードから第1の処理ノードにスワッピングするための手段と、

前記第1の処理ノードで前記ニューラルネットワークの前記一部を実行するための手段と、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すための手段と、

前記第2の処理ノードで前記ニューラルネットワークの前記一部を実行するための手段と

を備える、装置。

【請求項26】

人工ニューラルネットワークにおいてコプロセッシングを実行するためのプログラムコードを符号化した非一時的コンピュータ可読媒体であって、前記プログラムコードは、プロセッサによって実行され、

一定時間期間にわたって、前記ニューラルネットワークの一部を第1の処理ノードにスワッピングするためのプログラムコードと、

前記第1の処理ノードで前記ニューラルネットワークの前記一部を実行するためのプログラムコードと、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を第2の処理ノードに返すためのプログラムコードと、

前記第2の処理ノードで前記ニューラルネットワークの前記一部を実行するためのプログラムコードと

を備える、非一時的コンピュータプログラム可読媒体。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】0131

【補正方法】変更

【補正の内容】

【0131】

[00116]特許請求の範囲は、上記で示した厳密な構成および構成要素に限定されないことを理解されたい。上記で説明した方法および装置の構成、動作および詳細において、特許請求の範囲から逸脱することなく、様々な改変、変更および変形が行われ得る。

以下に、出願当初の特許請求の範囲に記載された発明を付記する。

[C 1]

ニューラルネットワークにおいてコプロセッシングを実行する方法であって、
一定時間期間にわたって、前記ニューラルネットワークの一部を第 1 の処理ノードにスワッピングすることと、

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行することと、
前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すことと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行することと
を備える、方法。

[C 2]

前記第 1 の処理ノードは、別個のハードウェアコアを備える、

C 1 に記載の方法。

[C 3]

前記第 1 の処理ノードは、学習処理コアを備える、

C 1 に記載の方法。

[C 4]

前記学習処理コアは、前記第 2 の処理ノードよりも高いレベルのリソースで構成される、

C 3 に記載の方法。

[C 5]

学習は、オフラインまたはオンラインで実装される、

C 3 に記載の方法。

[C 6]

前記学習処理コアの入力および出力は、学習がオフラインで実装される場合、前記ニューラルネットワークの他のレイヤを備える、

C 5 に記載の方法。

[C 7]

前記第 1 の処理ノードは、学習処理コアを備え、

前記第 2 の処理ノードは、静的処理コアを備え、

スワッピングすることは、

前記静的処理コアの状態を前記学習処理コアにコピーすることと、

前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに
入力をルーティングすることと

を備え、

返すことは、

前記学習処理コアの状態を前記静的処理コアにコピーすることと、

変更された静的処理コアに制御を返すことと

を備える、C 1 に記載の方法。

[C 8]

前記スワッピングすることは、前記第 1 の処理ノードから前記第 2 の処理ノードにリソースを割り振ることを備える、

C 1 に記載の方法。

[C 9]

前記ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、

C 1 に記載の方法。

[C 1 0]

前記第 1 の処理ノードは、デバッグコアを備える、

C 1 に記載の方法。

[C 1 1]

前記スワッピングすることは、システム性能がしきい値を下回る場合に発生する、

C 1 に記載の方法。

[C 1 2]

前記返すことは、システム性能がしきい値を上回る場合に発生する、

C 1 に記載の方法。

[C 1 3]

前記スワッピングすること、または返すことは、電力がシステムに適用されると発生する、

C 1 に記載の方法。

[C 1 4]

ニューラルネットワークにおいてコプロセッシングを実行するための装置であって、メモリと、

前記メモリに結合された少なくとも 1 つのプロセッサと

を備え、前記少なくとも 1 つのプロセッサは、

一定時間期間にわたって、前記ニューラルネットワークの一部を第 1 の処理ノードにスワッピングすることと、

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行することと、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すことと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行することと

を行うように構成される、装置。

[C 1 5]

前記第 1 の処理ノードは、別個のハードウェアコアを備える、

C 1 4 に記載の装置。

[C 1 6]

前記第 1 の処理ノードは、学習処理コアを備える、

C 1 4 に記載の装置。

[C 1 7]

前記学習処理コアは、前記第 2 の処理ノードよりも高いレベルのリソースで構成される、

C 1 6 に記載の装置。

[C 1 8]

学習は、オフラインまたはオンラインで実装される、

C 1 6 に記載の装置。

[C 1 9]

前記学習処理コアの入力および出力は、学習がオフラインで実装される場合、前記ニューラルネットワークの他のレイヤを備える、

C 1 8 に記載の装置。

[C 2 0]

前記第 1 の処理ノードは、学習処理コアを備え、前記第 2 の処理ノードは、静的処理コ

アを備え、前記少なくとも1つのプロセッサは、
前記静的処理コアの状態を前記学習処理コアにコピーすることと、
前記学習処理コアが前記静的処理コアの機能を包含するように、前記学習処理コアに
入力をルーティングすることと、
前記学習処理コアの状態を前記静的処理コアにコピーすることと、
変更された静的処理コアに制御を返すことと
を行うようにさらに構成される、C 1 4に記載の装置。

[C 2 1]

前記少なくとも1つのプロセッサは、前記第1の処理ノードから前記第2の処理ノード
にリソースを割り振ることを行うようにさらに構成される、
C 1 4に記載の装置。

[C 2 2]

前記ニューラルネットワークの前記一部は、深層信念ネットワークのレイヤを備える、
C 1 4に記載の装置。

[C 2 3]

前記第1の処理ノードは、デバッグコアを備える、
C 1 4に記載の装置。

[C 2 4]

前記少なくとも1つのプロセッサは、システム性能がしきい値を下回る場合に、前記ニューラルネットワークの前記一部を前記第1の処理ノードにスワッピングするようにさらに構成される、
C 1 4に記載の装置。

[C 2 5]

前記少なくとも1つのプロセッサは、システム性能がしきい値を上回る場合に、前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すようにさらに構成される、
C 1 4に記載の装置。

[C 2 6]

前記少なくとも1つのプロセッサは、電力がシステムに適用されると、前記ニューラルネットワークの前記一部を前記第1の処理ノードにスワッピングする、または前記ニューラルネットワークの前記一部を前記第2の処理ノードに返すようにさらに構成される、
C 1 4に記載の装置。

[C 2 7]

ニューラルネットワークにおいてコプロセッシングを実行するための装置であって、
一定時間期間にわたって、前記ニューラルネットワークの一部を第1の処理ノードにスワッピングするための手段と、
前記第1の処理ノードで前記ニューラルネットワークの前記一部を実行するための手段と、
前記一定時間期間後に、前記ニューラルネットワークの前記一部を第2の処理ノードに返すための手段と、
前記第2の処理ノードで前記ニューラルネットワークの前記一部を実行するための手段と
を備える、装置。

[C 2 8]

ニューラルネットワークにおいてコプロセッシングを実行するためのコンピュータプログラム製品であって、
プログラムコードを符号化した非一時的コンピュータ可読媒体を備え、前記プログラムコードは、
一定時間期間にわたって、前記ニューラルネットワークの一部を第1の処理ノードにスワッピングするためのプログラムコードと、

前記第 1 の処理ノードで前記ニューラルネットワークの前記一部を実行するためのプログラムコードと、

前記一定時間期間後に、前記ニューラルネットワークの前記一部を第 2 の処理ノードに返すためのプログラムコードと、

前記第 2 の処理ノードで前記ニューラルネットワークの前記一部を実行するためのプログラムコードと

を備える、コンピュータプログラム製品。

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No PCT/US2015/015917

A. CLASSIFICATION OF SUBJECT MATTER INV. G06N3/08 G06N3/10 ADD. G06N3/04		
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) G06N		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal, WPI Data, INSPEC		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 6 418 423 B1 (KAMBHATLA NANDAKISHORE [US] ET AL) 9 July 2002 (2002-07-09) columns 4-10	1-28
X	----- D. CALVERT ET AL: "Distributed artificial neural network architectures", PROCEEDINGS OF THE 19TH INTERNATIONAL SYMPOSIUM ON HIGH PERFORMANCE COMPUTING SYSTEMS AND APPLICATIONS (HPCS'05), 15 May 2005 (2005-05-15), pages 2-10, XP010800324, DOI: 10.1109/HPCS.2005.24 section II ----- -/--	1-28
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents : "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search		Date of mailing of the international search report
1 December 2015		14/12/2015
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer Douarche, Nicolas

3

INTERNATIONAL SEARCH REPORT

International application No PCT/US2015/015917

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	C. ROSSANT ET AL: "Playdoh: a lightweight Python library for distributed computing and optimisation", JOURNAL OF COMPUTATIONAL SCIENCE, vol. 4, no. 5, 6 July 2011 (2011-07-06), pages 352-359, XP055183262, DOI: 10.1016/j.jocs.2011.06.002 sections 2 and 4.1 -----	1-28
X	F. GALLUPPI ET AL: "A hierachical configuration system for a massively parallel neural hardware platform", PROCEEDINGS OF THE 9TH CONFERENCE ON COMPUTING FRONTIERS (CF'12), 15 May 2012 (2012-05-15), pages 183-192, XP055179253, DOI: 10.1145/2212908.2212934 section 4.1 -----	1-28
X	X. CHEN ET AL: "Pipelined back-propagation for context-dependent deep neural networks", PROCEEDINGS OF THE 13TH ANNUAL CONFERENCE OF THE INTERNATIONAL SPEECH COMMUNICATION ASSOCIATION (INTERSPEECH'12), 9 September 2012 (2012-09-09), pages 26-29, XP055113689, ISSN: 1990-9770 section 3 -----	1-28
X	E. SCHIKUTA, E. MANN: "N2Sky - neural networks as services in the clouds", PROCEEDINGS OF THE 2013 INTERNATIONAL JOINT CONFERENCE ON NEURAL NETWORKS (IJCNN'13), 4 August 2013 (2013-08-04), XP032542139, DOI: 10.1109/IJCNN.2013.6707113 the whole document -----	1-28
X	K. MINKOVICH ET AL: "HRLSim: a high performance spiking neural network simulator for GPGPU clusters", IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS, vol. 25, no. 2, 15 August 2013 (2013-08-15), pages 316-331, XP011536922, DOI: 10.1109/TNNLS.2013.2276056 section III -----	1-28

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2015/015917

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 6418423	B1	09-07-2002	NONE

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(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), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), 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, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, 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

(72)発明者 カンボス、マイケル

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 レウイス、アンソニー

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ラオ、ナビーン・ガンドハム

アメリカ合衆国、カリフォルニア州 9 2 1 3 1、サン・ディエゴ、カミニート・アレグラ 1 1 0 8 7