

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第6996497号

(P6996497)

(45)発行日 令和4年1月17日(2022.1.17)

(24)登録日 令和3年12月20日(2021.12.20)

(51)国際特許分類

F I

G 0 6 N 3/02 (2006.01)

G 0 6 N 3/02

G 0 6 N 3/10 (2006.01)

G 0 6 N 3/10

G 0 6 N 3/08 (2006.01)

G 0 6 N 3/08

1 2 0

請求項の数 8 (全27頁)

(21)出願番号 特願2018-514172(P2018-514172)
(86)(22)出願日 平成29年3月7日(2017.3.7)
(86)国際出願番号 PCT/JP2017/008988
(87)国際公開番号 WO2017/187798
(87)国際公開日 平成29年11月2日(2017.11.2)
審査請求日 令和2年2月17日(2020.2.17)
(31)優先権主張番号 特願2016-91418(P2016-91418)
(32)優先日 平成28年4月28日(2016.4.28)
(33)優先権主張国・地域又は機関
日本国(JP)

(73)特許権者 000002185
ソニーグループ株式会社
東京都港区港南1丁目7番1号
(74)代理人 110002147
特許業務法人酒井国際特許事務所
(72)発明者 大淵 愉希夫
東京都港区港南1丁目7番1号 ソニー
株式会社内
(72)発明者 高松 慎吾
東京都港区港南1丁目7番1号 ソニー
株式会社内
(72)発明者 福井 啓
東京都港区港南1丁目7番1号 ソニー
株式会社内
(72)発明者 井手 直紀

最終頁に続く

(54)【発明の名称】 情報処理装置、及び情報処理方法

(57)【特許請求の範囲】

【請求項1】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行う判定部と、

前記ニューラルネットワークの設計を制御し、前記判定部により前記制約を満たさないと判定された場合に、前記制約を満たすように、前記ニューラルネットワークにおけるレイヤーの再配置を行う設計制御部と、
を備える情報処理装置。

【請求項2】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行う判定部と、

前記判定部により前記制約を満たさないと判定された場合に、前記ニューラルネットワークにおいて前記制約が満たされない部分を提示する警告画面を表示させる表示制御部と、
を備える情報処理装置。

【請求項3】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行う判定部と、

前記制約に基づいて、前記ニューラルネットワークにおけるレイヤーごとに、前記レイヤーに対応付けられたハードウェアに応じた学習を行う学習部と、
を備える情報処理装置。

【請求項 4】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行う判定部と、
前記判定部により前記制約を満たすと判定されたニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成する生成部と、
を備え、
前記判定部は、前記生成部により生成されるニューラルネットワークが前記制約を満たすか否か判定を行い、
前記生成部は、前記判定部により前記制約を満たすと判定されるニューラルネットワークが生成されるまで、前記別のニューラルネットワークの生成を繰り返す、
情報処理装置。

10

【請求項 5】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行うことと、
前記ニューラルネットワークの設計を制御することと、
前記制約を満たさないと判定した場合に、前記制約を満たすように、前記ニューラルネットワークにおけるレイヤーの再配置を行うことと、
を含む情報処理方法。

【請求項 6】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行うことと、
前記制約を満たさないと判定した場合に、前記ニューラルネットワークにおいて前記制約が満たされない部分を提示する警告画面を表示させることと、
を含む情報処理方法。

20

【請求項 7】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行うことと、
前記制約に基づいて、前記ニューラルネットワークにおけるレイヤーごとに、前記レイヤーに対応付けられたハードウェアに応じた学習を行うことと、
を含む情報処理方法。

30

【請求項 8】

ニューラルネットワークが、複数のハードウェアに係る制約を満たすか否か判定を行うことと、
前記制約を満たすと判定したニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成することと、
生成したニューラルネットワークが前記制約を満たすか否か判定することと、
前記制約を満たすと判定されるニューラルネットワークが生成されるまで、前記別のニューラルネットワークの生成を繰り返すことと、
を含む情報処理方法。

【発明の詳細な説明】

40

【技術分野】**【0001】**

本開示は、情報処理装置、及び情報処理方法に関する。

【背景技術】**【0002】**

近年、脳神経系の仕組みを模したニューラルネットワークが注目されている。また、ニューラルネットワークを開発するための種々の手法が提案されている。例えば、非特許文献 1 には、ニューラルネットワークによる学習過程をモニタリングするライブラリが開示されている。

【先行技術文献】

50

【非特許文献】

【 0 0 0 3 】

【文献】M. Abadi、外 3 9 名、「TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems」、2 0 1 5 年 1 1 月 9 日、[Online]、[平成 2 8 年 4 月 2 2 日検索]、インターネット <http://download.tensorflow.org/paper/whitepaper2015.pdf>

【発明の概要】

【発明が解決しようとする課題】

【 0 0 0 4 】

しかし、非特許文献 1 に記載のライブラリでは、ニューラルネットワークが単一のハードウェア上で実行されることを前提としており、複数のハードウェアによる処理に適したニューラルネットワークを設計することは困難であった。

【 0 0 0 5 】

そこで、本開示では、複数のハードウェアによる処理に適したニューラルネットワークをより効率的に設計することが可能な情報処理装置、及び情報処理方法を提案する。

【課題を解決するための手段】

【 0 0 0 6 】

本開示によれば、複数のハードウェアに係る制約を取得する取得部と、ニューラルネットワークが、前記制約を満たすか否か判定を行う判定部と、を備える情報処理装置が提供される。

【 0 0 0 7 】

また、本開示によれば、ニューラルネットワークが複数のハードウェアに係る制約を満たすか否かの判定結果を受信する受信部と、前記判定結果に基づいて処理を行う処理部と、を備える情報処理装置。

【 0 0 0 8 】

また、本開示によれば、複数のハードウェアに係る制約を取得することと、ニューラルネットワークが、前記制約を満たすか否か判定を行うことと、を含む情報処理方法が提供される。

【発明の効果】

【 0 0 0 9 】

以上説明したように本開示によれば、複数のハードウェアによる処理に適したニューラルネットワークをより効率的に設計することが可能である。

【 0 0 1 0 】

なお、上記の効果は必ずしも限定的なものではなく、上記の効果とともに、または上記の効果に代えて、本明細書に示されたいずれかの効果、または本明細書から把握され得る他の効果が奏されてもよい。

【図面の簡単な説明】

【 0 0 1 1 】

【図 1】ニューラルネットワークによる認識処理を複数のハードウェアで実行する例を説明するための説明図である。

【図 2】本開示の第一の実施形態に係る情報処理システムの構成例を説明するための説明図である。

【図 3】複数のハードウェアに係る制約を入力するための制約入力画面の一例を示す説明図である。

【図 4】同実施形態に係るニューラルネットワークの設計を行うための設計画面の一例を示す説明図である。

【図 5】同実施形態に係るサーバ 2 の構成例を説明するための説明図である。

【図 6】同実施形態による情報処理システム 1 0 0 0 の処理フロー例を示すフローチャート図である。

【図 7】同実施形態に係る情報処理システム 1 0 0 0 の処理フロー例を示すフローチャー

10

20

30

40

50

ト図である。

【図 8】同実施形態に係るクライアント端末 1 に表示される画面例を示す説明図である。

【図 9】同実施形態に係るクライアント端末 1 に表示される画面例を示す説明図である。

【図 10】同実施形態に係るクライアント端末 1 に表示される画面例を示す説明図である。

【図 11】本開示の第二の実施形態に係るサーバ 2 - 2 の構成例を説明するための説明図である。

【図 12】同実施形態の動作例を説明するためのフローチャート図である。

【図 13】同実施形態に係るニューラルネットワーク生成に係る処理フローを説明するためのフローチャート図である。

【図 14】同実施形態に係る生成部 217 による突然変異を用いたネットワーク生成を説明するためのフローチャートである。

10

【図 15】同実施形態に係る生成部 217 による交叉を用いたネットワーク生成を説明するためのフローチャート図である。

【図 16】ハードウェア構成例を示す説明図である。

【発明を実施するための形態】

【0012】

以下に添付図面を参照しながら、本開示の好適な実施の形態について詳細に説明する。なお、本明細書及び図面において、実質的に同一の機能構成を有する構成要素については、同一の符号を付することにより重複説明を省略する。

【0013】

20

<< 1. 第一の実施形態 >>

< 1 - 1. 背景 >

< 1 - 2. 構成例 >

< 1 - 3. 動作例 >

< 1 - 4. 効果 >

<< 2. 第二の実施形態 >>

< 2 - 1. 構成例 >

< 2 - 2. 動作例 >

< 2 - 3. 効果 >

<< 3. 変形例 >>

30

< 3 - 1. 変形例 1 >

< 3 - 2. 変形例 2 >

<< 4. ハードウェア構成例 >>

<< 5. むすび >>

【0014】

<< 1. 第一の実施形態 >>

< 1 - 1. 背景 >

本開示の第一の実施形態に係る情報処理装置について説明する前に、まず、本実施形態の創作に至った背景を説明する。

【0015】

40

近年、人間の脳神経回路を模したモデルであり、人間が持つ学習能力をコンピュータ上で実現しようとするニューラルネットワークが注目されている。上述したとおり、ニューラルネットワークは学習能力を有することを特徴の一つとする。ニューラルネットワークでは、シナプスの結合によりネットワークを形成した人工ニューロン（ノード）が、学習によりシナプスの結合強度を変化させることで、問題に対する解決能力を獲得することが可能である。すなわち、ニューラルネットワークは、学習を重ねることで、問題に対する解決ルールを自動的に推論することができる。

【0016】

ニューラルネットワークによる学習の例としては、画像認識や音声認識が挙げられる。ニューラルネットワークでは、例えば、手書きの数字パターンを繰り返し学習することで、

50

入力される画像情報を、0～9の数字のいずれかに分類することが可能となる。ニューラルネットワークの有する上記のような学習能力は、人工知能（Artificial Intelligence）の発展を押し進める鍵としても注目されている。また、ニューラルネットワークが有するパターン認識力は、種々の産業分野における応用が期待される。

【0017】

一方、ニューラルネットワークを設計する際、ニューラルネットワークによる認識処理におけるパフォーマンス（例えば実行時間や消費電力等）は重要な指標となる。また、パフォーマンスを向上させるため、異なる種類の複数のハードウェアを用いてニューラルネットワークによる認識処理を行うことが考えられる。

【0018】

例えば、ニューラルネットワークによる認識処理が、処理毎に（例えばニューラルネットワークに含まれるレイヤーごと）に異なるハードウェアにより実行されることで、全体としてのパフォーマンスが向上する場合がある。

【0019】

図1は、ニューラルネットワークによる認識処理を複数のハードウェアで実行する例を説明するための説明図である。図1に示すニューラルネットワークにおいて、入力レイヤー（Input layer）と出力レイヤー（Output layer）の間には処理P1～P3が存在する。上記の処理P1～P3が、それぞれ異なるハードウェアで実行されてもよく、例えば、それぞれニューロチップ、CPU、及びGPUで実行されてもよい。例えば、CPUは消費電力が大きく、GPUは消費電力が小さい場合、処理とハードウェアの組み合わせに応じて、パフォーマンスが異なり得る。

【0020】

しかし、単一のハードウェアで認識処理が実行されることを前提としたニューラルネットワークの設計ツールを用いて、複数のハードウェアで実行されるニューラルネットワークの設計を行うことは困難であった。また、複数のハードウェアで実行されるニューラルネットを設計できたとしても、それぞれのハードウェアには、制約が存在し得るため、例えば設計されたニューラルネットワークが、当該複数のハードウェアの実行に適合しない恐れがあった。例えば、ハードウェアによっては、処理可能なノード数に制約が存在する場合があります。当該ノード数以上のノード数を有するレイヤーが当該ハードウェアに割り当てられてしまうと、認識処理を実行することができない恐れがあった。

【0021】

そこで、上記事情を一着眼点にして本実施形態を創作するに至った。本実施形態は、複数のハードウェアに係る制約に基づき、ニューラルネットワークが当該制約を満たすか否かを判定し、例えばユーザに判定結果に基づく警告を提供する。本実施形態によれば、複数のハードウェアによる処理に適合したニューラルネットワークをより効率的に設計することが可能である。以下、上記のような効果を実現するための本開示の第一の実施形態の構成例について説明する。

【0022】

< 1 - 2 . 構成例 >

図2を参照しながら、本開示の第一の実施形態の構成例を説明する。図2は、本開示の第一の実施形態に係る情報処理システムの構成例を説明するための説明図である。本実施形態に係る情報処理システム1000は、ユーザによるニューラルネットワークの設計のための情報処理システムであり、例えばビジュアルプログラミングにより、ニューラルネットワークを設計することが可能なツールを提供してもよい。

【0023】

本開示において、ビジュアルプログラミングとは、ソフトウェア開発において、プログラムコードをテキストで記述することなく、視覚的なオブジェクトを用いて作成する手法を指す。ビジュアルプログラミングでは、例えば、GUI（Graphical User Interface）上で、オブジェクトを操作することで、プログラムを作成することができる。

10

20

30

40

50

【 0 0 2 4 】

図 2 に示すように、本実施形態に係る情報処理システム 1 0 0 0 は、クライアント端末 1、サーバ 2、及び通信網 5 を含み、クライアント端末 1 とサーバ 2 とは、通信網 5 を介して互いに通信を行えるように接続される。

【 0 0 2 5 】

クライアント端末 1 は、ユーザがニューラルネットワークの設計を行うための情報処理装置である。例えば、クライアント端末 1 は、ビジュアルプログラミングによりニューラルネットワークの設計を行うための設計画面をユーザに提示（表示）してもよい。また、クライアント端末 1 は、サーバ 2 からニューラルネットワークが複数のハードウェアに係る制約を満たすか否かの判定結果を受信し、当該判定結果に基づいて、表示制御処理を行ってもよい。

10

【 0 0 2 6 】

サーバ 2 は、クライアント端末 1 にニューラルネットワークの設計を行うための設計画面を提供し、クライアント端末 1 を介したユーザの入力に基づいて、ニューラルネットワークに係るプログラムを作成する情報処理装置である。また、サーバ 2 は、複数のハードウェアに係る制約に基づいて、ニューラルネットワークが当該制約を満たすか否かを判定を行い、判定結果をクライアント端末 1 に提供する。また、ニューラルネットワークが当該制約を満たさない場合、サーバ 2 は当該ニューラルネットワークに係るプログラムを作成しなくてもよい。

【 0 0 2 7 】

通信網 5 は、通信網 5 に接続されている装置、またはシステムから送信される情報の有線、または無線の伝送路である。例えば、通信網 5 は、インターネット、電話回線網、衛星通信網等の公衆回線網や、Ethernet（登録商標）を含む各種の LAN（Local Area Network）、WAN（Wide Area Network）等を含んでもよい。また、通信網 5 は、IP-VPN（Internet Protocol-Virtual Private Network）等の専用回線網を含んでもよい。

20

【 0 0 2 8 】

本実施形態に係る情報処理システム 1 0 0 0 によれば、ユーザにより設計されるニューラルネットワークが、複数のハードウェアに係る制約を満たすか否かの判定が行われ、判定結果がユーザに提供される。また、本実施形態に係る情報処理システム 1 0 0 0 によれば、複数のハードウェアに係る制約を満たすニューラルネットワークに係るプログラムが作成されるようなニューラルネットワークの設計ツールが提供される。

30

【 0 0 2 9 】

（クライアント端末）

続いて、本実施形態に係るクライアント端末 1 について詳細に説明する。図 2 に示すように、本実施形態に係るクライアント端末 1 は、制御部 1 0、通信部 1 2、表示部 1 4、及び操作部 1 6 を備える情報処理装置である。なお、クライアント端末 1 は、例えば PC（Personal Computer）、タブレット PC 等であってもよい。

【 0 0 3 0 】

制御部 1 0 は、クライアント端末 1 の各構成を制御する。例えば、制御部 1 0 は、通信部 1 2 による通信を制御する通信制御部としての機能を有する。係る構成により、通信部 1 2 は、例えば各種画面や、ニューラルネットワークが複数のハードウェアに係る制約を満たすか否かの判定結果をサーバ 2 から受信することが出来る。

40

【 0 0 3 1 】

また、制御部 1 0 は、表示部 1 4 による表示の制御処理を行う表示制御部としての機能を有する。例えば、制御部 1 0 は、各種画面を表示部 1 4 に表示させてもよい。図 3、図 4 を参照して、制御部 1 0 が表示部 1 4 に表示させる画面の例を説明する。

【 0 0 3 2 】

図 3 は、複数のハードウェアに係る制約を入力するための制約入力画面の一例を示す説明図である。図 3 に示すように、制約入力画面は、例えば、認識処理に用いられるハードウ

50

エアを示すハードウェア入力フォーム G 1 1 ~ G 1 4 と、ハードウェア間の通信速度を示す通信速度入力フォーム G 1 5 ~ G 1 7 を含む。ユーザは、ハードウェア入力フォーム G 1 1 ~ G 1 4 を用いて、ハードウェア (H W : H a r d w a r e) の種類 (T y p e) の選択や、ハードウェアの情報 (例えば演算性能等) の入力を行うことが可能である。例えば、ハードウェア入力フォーム G 1 1 ~ G 1 4 によりそれぞれハードウェア H W 1 ~ H W 4 が選択されることで、予め用意された、ハードウェアに係る制約が設定され得る。また、ユーザがハードウェア入力フォーム G 1 1 ~ G 1 4 を用いて各ハードウェアの情報をさらに入力することで、制約をカスタマイズすることも可能である。

【 0 0 3 3 】

また、ユーザは、通信速度入力フォーム G 1 5 ~ G 1 7 を用いて、ハードウェア間の通信速度の入力 (カスタマイズ) を行うことが可能である。なお、ハードウェア入力フォーム G 1 1 ~ G 1 4 の間の接続関係は、ユーザ操作により変更可能であってもよい。

10

【 0 0 3 4 】

図 4 は、ニューラルネットワークの設計を行うための設計画面の一例を示す説明図である。図 4 に示す設計画面は、通信部 1 2 を介して、サーバ 2 から受信されてもよい。図 4 に示すように、設計画面は、ニューラルネットワークにおける複数のレイヤー G 2 0 ~ G 3 0 が配置される。配置されたレイヤー G 2 0 ~ G 3 0 は、例えばデータの取り込み、データの処理、データの出力をそれぞれ意味していてもよい。ユーザは、図 4 に示すような設計画面を用いて、レイヤーの追加、削除、配置変更等を行い、ニューラルネットワークを設計することが出来る。

20

【 0 0 3 5 】

また、入力レイヤー G 2 0 と出力レイヤー G 3 0 の間のレイヤー G 2 1 ~ G 2 9 は、配置された順に逐次、または並列に処理されてもよい。例えば、レイヤー G 2 4 ~ G 2 6 の処理と、レイヤー G 2 7 ~ G 2 8 の処理とは、並列に処理されてもよい。

【 0 0 3 6 】

また、制御部 1 0 (処理部) は、表示部 1 4 を制御して、サーバ 2 から受信される判定結果に基づく表示制御処理を行ってもよい。例えば、サーバ 2 により、ニューラルネットワークが複数のハードウェアに係る制約を満たさないと判定された場合に、制御部 1 0 (表示制御部) は、制約が満たされないことを示す警告画面を表示させてもよい。また、制御部 1 0 が表示させる警告画面は、ニューラルネットワークにおいて当該制約を満たさない部分を提示する画面であってもよい。係る構成により、ユーザは複数のハードウェアに係る制約を満たすようにニューラルネットワークをより効率的に設計することが可能となる。

30

【 0 0 3 7 】

通信部 1 2 (受信部) は、制御部 1 0 により制御されて、他の装置との間の通信を仲介する通信インタフェースである。通信部 1 2 は、任意の無線通信プロトコルまたは有線通信プロトコルをサポートし、例えば図 2 に示す通信網 5 を介して他の装置との間の通信接続を確立する。例えば、通信部 1 2 はニューラルネットワークの設計を行うための設計画面や、ニューラルネットワークが複数のハードウェアに係る制約を満たすか否かの判定結果をサーバ 2 から受信する。また、通信部 1 2 は、表示部 1 4 に表示される各種画面に対するユーザの入力に係る情報をサーバ 2 に送信させる。

40

【 0 0 3 8 】

表示部 1 4 は制御部 1 0 に制御されて、各種画面を表示するディスプレイである。例えば、表示部 1 4 は、上述した制約入力画面、設計画面、警告画面等を表示してもよい。なお、表示部 1 4 は、例えば、C R T (C a t h o d e R a y T u b e) ディスプレイ装置、液晶ディスプレイ (L C D : L i q u i d C r y s t a l D i s p l a y) 装置、O L E D (O r g a n i c L i g h t E m i t t i n g D i o d e) 装置により実現されてもよい。

【 0 0 3 9 】

操作部 1 6 は、ユーザの入力を受け付け、制御部 1 0 に提供する。例えば、ユーザは、操作部 1 6 を操作して、複数のハードウェアに係る制約のカスタマイズや、ニューラルネッ

50

トワークの設計を行うための入力を行ってもよい。なお、操作部 16 は、例えばマウス、キーボード、タッチパネル、ボタン、スイッチ、視線入力装置、ジェスチャ入力装置、音声入力装置などにより実現されてもよい。

【0040】

(サーバ)

以上、本実施形態に係るクライアント端末 1 の構成例を説明した。続いて、図 5 を参照して、本実施形態に係るサーバ 2 の構成例を説明する。図 5 は、本実施形態に係るサーバ 2 の構成例を説明するための説明図である。図 5 に示すように、サーバ 2 は、制御部 20、通信部 22、及び記憶部 24 を備える情報処理装置である。

【0041】

制御部 20 は、サーバ 2 の各構成を制御する。また、制御部 20 は、図 5 に示すように、通信制御部 201、取得部 202、判定部 203、設計制御部 204、学習部 205、及び認識部 206 としても機能する。

【0042】

通信制御部 201 は、通信部 22 による通信を制御する。例えば、通信制御部 201 は、通信部 22 を制御して、ニューラルネットワークの設計画面、判定部 203 による判定結果等を、クライアント端末 1 に送信させてもよい。また、通信制御部 201 は、通信部 22 を制御して、複数のハードウェアに係る制約のカスタマイズや、ニューラルネットワークの設計を行うためのユーザの入力に係る情報を受信してもよい。

【0043】

取得部 202 は、複数のハードウェアに係る制約を取得する。例えば、取得部 202 は、記憶部 24 から、複数のハードウェアに係る制約を取得してもよい。また、取得部 202 は、通信部 22 を介して受信される複数のハードウェアに係る制約のカスタマイズを行うためのユーザの入力に基づいて、複数のハードウェアに係る制約を取得してもよい。また、取得部 202 は、ハードウェアの選択に係るユーザの入力と、記憶部 24 に予め記憶されたハードウェアに係る制約に基づいて、複数のハードウェアに係る制約を取得してもよい。

【0044】

以下では、図 1 を参照して説明した、ニューロチップ、CPU、及びGPUに係る制約の例について説明を行うが、ハードウェアに係る制約は以下の例に限定されない。

【0045】

取得部 202 により取得される、複数のハードウェアに係る制約は、例えば後述する判定部 203 による判定に用いられる制約であってもよい。判定に用いられる制約の例としては、ハードウェア間の接続に係る制約や、ハードウェア間の通信速度に係る制約、ハードウェアの処理能力に係る制約等であってもよい。以下に判定部 203 による判定に用いられる制約の例を示す。

【0046】

- ・ニューロチップはセンサ(入力レイヤー)に接続される
- ・ニューロチップが処理可能なノード数は10以下である
- ・ニューロチップが処理可能なレイヤーはコンボリューションレイヤーのみである
- ・ニューロチップとCPUとの間の通信速度は所定の速度である
- ・CPUが利用可能なRAMは所定値以下である
- ・CPUとGPUとの間の通信速度は所定の速度である
- ・GPUはニューロチップと直接接続できない

【0047】

また、取得部 202 により取得される、複数のハードウェアに係る制約は、例えば後述する学習部 205 がハードウェアに応じた学習を行うために用いられる制約であってもよい。学習を行うために用いられる制約の例としては、ハードウェアの特性に係る制約や、演算の種類に係る制約等であってもよい。以下に、学習部 205 による学習に用いられる制約の例を示す。

10

20

30

40

50

【 0 0 4 8 】

- ・ニューロチップのニューロンの特性はスパイクングである
- ・CPUは整数演算のみ可能（浮動小数点演算不可能）である

【 0 0 4 9 】

判定部 2 0 3 は、ニューラルネットワークが、取得部 2 0 2 により取得される複数のハードウェアに係る制約を満たすか否か判定を行う。判定部 2 0 3 は、例えば、ユーザの入力に基づいて設計制御部 2 0 4 により設計されたニューラルネットワークが、上記制約を満たすか否か判定を行う。

【 0 0 5 0 】

なお、判定部 2 0 3 は、予め設定される、またはユーザにより入力される、所定の処理時間にさらに基づき、ニューラルネットワークが、取得部 2 0 2 により取得される複数のハードウェアに係る制約を満たすか否か判定を行ってもよい。係る場合、判定部 2 0 3 は、ニューラルネットワークが取得部 2 0 2 により取得される複数のハードウェアに係る制約を満たし、かつ当該所定の処理時間内に処理が完了すると判定された場合に、制約を満たすと判定してもよい。

10

【 0 0 5 1 】

また、判定部 2 0 3 は、ユーザの入力に基づいて、設計制御部 2 0 4 によりニューラルネットワークが変更された（例えばレイヤーの追加、削除、配置変更等が発生した）場合に、変更後のニューラルネットワークが上記制約を満たすか否か判定を行ってもよい。なお、判定部 2 0 3 は判定結果を通信制御部 2 0 1 と設計制御部 2 0 4 に提供してもよい。

20

【 0 0 5 2 】

設計制御部 2 0 4 は、通信部 2 2 を介して取得されるユーザの入力に基づいて、ニューラルネットワークの設計を制御する。例えば、設計制御部 2 0 4 は、図 4 を参照して説明した設計画面を生成し、通信部 2 2 を介してクライアント端末 1 に提供してもよい。また、設計制御部 2 0 4 は、ユーザの入力に基づいて、ニューラルネットワークにおけるレイヤーの配置を制御してもよい。

【 0 0 5 3 】

また、設計制御部 2 0 4 は、通信部 2 2 を介して取得されるユーザの入力に基づいて、ニューラルネットワークにおけるレイヤーと、ハードウェアの対応付けを行う。対応付けに係るユーザの入力は、ニューラルネットワークの設計画面において範囲選択により行われてもよいし、レイヤーの処理順における順番の指定により行われてもよい。

30

【 0 0 5 4 】

また、設計制御部 2 0 4 は、設計制御部 2 0 4 により設計されたニューラルネットワークが、判定部 2 0 3 により、制約を満たすと判定された場合に、当該ニューラルネットワークを構築するプログラムを作成してもよい。係る構成によれば、複数のハードウェアに係る制約を満たし、当該複数のハードウェアにより認識を実行可能なニューラルネットワークを構築するプログラムが生成される。

【 0 0 5 5 】

また、設計制御部 2 0 4 は、設計制御部 2 0 4 により設計されたニューラルネットワークが、判定部 2 0 3 により、制約を満たさないと判定された場合に、当該制約を満たすように、ニューラルネットワークの変更に係るレイヤーの再配置を行ってもよい。例えば、設計制御部 2 0 4 は、追加されたレイヤーまたは配置変更されたレイヤーを、当該レイヤーが現在対応付けられているハードウェアから通信可能なハードウェアに再配置してもよい。また、設計制御部 2 0 4 は、再配置後のニューラルネットワークが制約を満たすか否かの判定結果を判定部 2 0 3 から取得して、制約が満たされるまで、再配置を繰り返してもよい。係る構成によれば、より効率的に、複数のハードウェアに係る制約を満たしたニューラルネットワークを設計することが可能となる。

40

【 0 0 5 6 】

学習部 2 0 5 は、設計制御部 2 0 4 により設計されたニューラルネットワークの学習を行う。学習部 2 0 5 は、例えば取得部 2 0 2 により取得される制約に基づいて、ニューラル

50

根とワークにおけるレイヤーごとに、レイヤーに対応付けられたハードウェアに応じた学習を行ってもよい。係る構成によれば、後述する認識部 206 により行われる認識の実行パフォーマンスが向上し得る。

【0057】

例えば、ニューロチップと対応付けられたレイヤーに係る学習は、当該ニューロチップの特性に応じた学習手法により行われてもよい。ニューロチップの特性に応じた学習手法は限定されないが、例えばニューロチップの特性がスパイクングである場合には、下記非特許文献 2 に記載される学習手法を用いることも可能である。

【0058】

非特許文献 2 : O. Peter、外 4 名、「Real-time classification and sensor fusion with a spiking deep belief network」、2013 年、Neuromorphic Engineering 7: 178.

10

【0059】

また、整数演算のみ可能なハードウェアと対応付けられたレイヤーに係る学習は、整数演算のみで処理可能となるような学習手法により行われてもよい。係る学習手法は限定されないが、例えば下記非特許文献 3 に記載される学習手法を用いることも可能である。

【0060】

非特許文献 3 : M. Courbariaux、外 2 名、「BinaryConnect: Training Deep Neural Networks with binary weights during propagations TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems」、2015 年 11 月 12 日、[Online]、[平成 28 年 4 月 21 日検索]、インターネット <http://arxiv.org/pdf/1511.00363.pdf>

20

【0061】

また、浮動小数点演算可能なハードウェアと対応付けられたレイヤーに係る学習は、浮動小数点演算を利用可能な多様な学習手法により行われ得る。

【0062】

なお、学習部 205 は、記憶部 24 に記憶される学習データに基づいて学習を行ってもよいし、通信部 22 を介して外部から取得される学習データに基づいて学習を行ってもよい。

【0063】

認識部 206 は、学習部 205 による学習に基づき、認識を実行する。認識部 206 はニューラルネットワークにおけるレイヤーごとに対応付けられたハードウェアで、フィードフォワード計算を行うことで認識を実行してもよい。なお、認識部 206 は、記憶部 24 に記憶されるデータに対して認識を行ってもよいし、通信部 22 を介して外部から取得されるデータに対して認識を行ってもよい。

30

【0064】

通信部 22 は、他の装置との間の通信を仲介する通信インタフェースである。通信部 22 は、任意の無線通信プロトコルまたは有線通信プロトコルをサポートし、例えば図 2 に示した通信網 5 を介して他の装置との間の通信接続を確立する。それにより、例えば、サーバ 2 が通信網 5 に接続されたクライアント端末 1 からユーザの入力等を受信し、及びクライアント端末 1 に設計画面、及び判定部 203 による判定結果等を送信することが可能となる。

40

【0065】

記憶部 24 は、サーバ 2 の各構成が機能するためのプログラムやパラメータを記憶する。また、記憶部 24 は、ハードウェアに係る制約、学習データ、認識用のデータ等を記憶してもよい。

【0066】

なお、上記ではサーバ 2 が学習部 205、認識部 206 を備えて、学習と認識を行う例を説明したが、本実施形態は上記に限定されない。例えば、学習、及び認識は、通信網 5 に接続される他の装置で行われてもよく、学習、及び認識のそれぞれが異なる装置により行われてもよい。係る場合、取得部 202 は、通信部 22 を介して、当該認識を行う装置が

50

ら、当該認識を行う装置が有する複数のハードウェアに係る制約を取得してもよい。

【 0 0 6 7 】

< 1 - 3 . 動作例 >

以上、本実施形態による情報処理システム 1 0 0 0 の構成例について説明した。続いて、本実施形態による情報処理システム 1 0 0 0 の動作例について、図 6 ~ 1 0 を参照して説明する。以下では、まず情報処理システム 1 0 0 0 の処理フローについて図 6、図 7 を参照して説明した後、本実施形態においてクライアント端末 1 に表示される画面遷移例について図 8 ~ 1 0 を参照して説明する。

【 0 0 6 8 】

(処理フロー例 1)

図 6 は、本実施形態による情報処理システム 1 0 0 0 の処理フロー例を示すフローチャート図である。図 6 に示すフローチャートは、本実施形態に係る動作のうち、特にニューラルネットワークの設計に係る処理フローを示す。

【 0 0 6 9 】

まず、ユーザが図 3 を参照して説明した制約入力画面を用いて、ハードウェアに係る制約を入力し、取得部 2 0 2 が制約を取得する (S 1 1 0 1)。

【 0 0 7 0 】

続いて、ユーザにより、ハードウェアとレイヤーの対応に関するハードウェアごとの範囲が設定される (S 1 1 0 2)。なお、ハードウェアごとの範囲の設定は、ニューラルネットワークの設計画面において、表示上の範囲を選択することで行われてもよいし、レイヤーの処理順における順番の指定により行われてもよい。

【 0 0 7 1 】

続いて、ユーザにより、ニューラルネットワークの設計画面を用いたニューラルネットワークの設計が行われる (S 1 1 0 3)。

【 0 0 7 2 】

続いて、判定部 2 0 3 により、設計されたニューラルネットワークが、複数のハードウェアに係る制約を満たすか否かの判定が行われる (S 1 1 0 4)。判定部 2 0 3 により、ニューラルネットワークが制約を満たさないと判定された場合 (S 1 1 0 4 : N O)、制約を満たさない箇所 (部分) を提示する警告画面が表示され (S 1 1 0 5)、処理はステップ S 1 1 0 3 に戻る。一方、判定部 2 0 3 により、ニューラルネットワークが制約を満たすと判定された場合 (S 1 1 0 4 : Y E S)、設計処理は終了する。

【 0 0 7 3 】

(処理フロー例 2)

図 7 は、本実施形態による情報処理システム 1 0 0 0 の処理フロー例を示すフローチャート図である。図 7 に示すフローチャートは、本実施形態に係る動作のうち、特にニューラルネットワークの設計の変更に係る処理フローを示す。例えば以下に説明する処理フローは、図 6 に示したフローチャートの処理により設計されたニューラルネットワークの設計を変更する際の処理フローであってもよい。

【 0 0 7 4 】

図 7 に示すステップ S 1 2 0 1、S 1 2 0 2 の処理は、図 6 を参照して説明したステップ S 1 1 0 1、S 1 1 0 2 の処理と同様であるため、説明を省略する。

【 0 0 7 5 】

続いて、ユーザにより、ニューラルネットワークの設計の変更 (例えば、レイヤーの追加、削除、配置変更) が行われる (S 1 2 0 3)。続いて、判定部 2 0 3 により、変更後のニューラルネットワークが、複数のハードウェアに係る制約を満たすか否かの判定が行われる (S 1 2 0 4)。

【 0 0 7 6 】

判定部 2 0 3 により、ニューラルネットワークが制約を満たさないと判定された場合 (S 1 2 0 4 : N O)、設計制御部 2 0 4 による自動的な再配置が可能であるか否かが判定される (S 1 2 0 5)。自動的な再配置が可能ではない場合 (S 1 2 0 5 : N O)、処理は

10

20

30

40

50

ステップ S 1 2 0 3 に戻る。

【 0 0 7 7 】

一方、自動的な再配置が可能である場合、(S 1 2 0 5 : N O)、設計制御部 2 0 4 は、自動的に再配置を行う(S 1 2 0 6)。また、自動的な再配置が行われたニューラルネットワークが設計画面に表示される(S 1 2 0 7)

【 0 0 7 8 】

また、ステップ S 1 2 0 4 において制約を満たすと判定された場合、設計が変更されたニューラルネットワークが設計画面に表示される(S 1 2 0 7)。

【 0 0 7 9 】

設計が変更されたニューラルネットワーク、または自動的な再配置が行われたニューラルネットワークを設計画面において確認したユーザにより、設計終了の操作入力が行われると(S 1 2 0 8 : Y E S)、処理は終了する。一方、ユーザが続けて設計の変更を行う場合(S 1 2 0 8 : N O)、処理はステップ S 1 2 0 3 に戻る。

【 0 0 8 0 】

(画面遷移例)

以上、本実施形態による情報処理システム 1 0 0 0 の処理フローを説明した。続いて、図 7 を参照して説明した処理フローにおいて、クライアント端末 1 に表示される画面の遷移例を図 3、4 及び図 8 ~ 1 0 を参照して説明する。図 8 ~ 1 0 は、本実施形態に係るクライアント端末 1 に表示される画面例を示す説明図である。なお、図 8 ~ 1 0 は図 4 を参照して説明した設計画面に含まれるニューラルネットワークの設計変更に係る画面遷移の一例である。また、以下では、図 7 に示した処理ステップを適宜参照しながら説明を行う。

【 0 0 8 1 】

まず、図 7 のステップ S 1 2 0 1 において、図 3 を参照して説明した制約入力画面が表示される。なお、ハードウェアに係る制約をユーザが変更する必要のない場合には、ユーザの入力なしに、取得部 2 0 2 が例えば記憶部 2 4 から制約を取得してもよい。

【 0 0 8 2 】

続いて、図 7 のステップ S 1 2 0 2 において、図 8 に示す設計画面のように、ハードウェアごとの範囲選択が行われる。ここで、図 8 に示す範囲 G 3 1 ~ G 3 4 に含まれるレイヤーは、設計制御部 2 0 4 により、各範囲に対応するハードウェアと対応付けられる。なお、図 8 に示す範囲 G 3 1、G 3 2、G 3 3、G 3 4 は、例えば、それぞれ図 3 におけるハードウェア H W 1、H W 2、H W 4、H W 3 に対応している。

【 0 0 8 3 】

続いて、図 7 のステップ S 1 2 0 3 において、図 9 に示す画面のように、ニューラルネットワークの変更が行われる。図 9 に示す例では、レイヤー G 4 2 がレイヤー G 4 1 とレイヤー G 4 3 の間に追加されている。なお、図 8 を参照すれば、図 9 において新たに追加されるレイヤー G 4 2 は、図 8 に示す範囲 G 3 3 に含まれるため、H W 4 に対応付けられる。

【 0 0 8 4 】

続いて、図 7 のステップ S 1 2 0 4 において、制約を満たすと判定された場合、ステップ S 1 2 0 7 において、図 1 0 に示すように、変更によりレイヤー G 5 2 が新たに追加されたニューラルネットワークを含む設計画面が表示される。一方、図 7 のステップ S 1 2 0 4 において、制約を満たさないと判定された場合、図 4 に示した画面のように、変更前のニューラルネットワークを含む設計画面が表示される(表示が図 4 の設計画面に戻る)。

【 0 0 8 5 】

< 1 - 4 . 効果 >

以上説明したように、本開示の第一の実施形態によれば、複数のハードウェアに係る制約に基づき、ニューラルネットワークが当該制約を満たすか否かを判定し、例えばユーザに判定結果に基づく警告画面を提供する。係る構成により、複数のハードウェアによる処理に適合したニューラルネットワークをより効率的に設計することが可能である。また、本実施形態によれば、変更されたニューラルネットワークが当該制約を満たさない場合には、制約を満たすように自動的に再配置が行われることで、制約を満たしたニューラルネッ

10

20

30

40

50

トワークの設計を支援することが可能である。

【 0 0 8 6 】

< 2 . 第二の実施形態 >

以上、本開示の第一の実施形態を説明した。続いて、本開示の第二の実施形態を説明する。本開示の第二の実施形態は、設計されたニューラルネットワークの評価結果に基づいて、ネットワーク構造の異なる別のニューラルネットワークを生成する。また、本開示の第二の実施形態は、生成されたニューラルネットワークの評価結果に基づいて、評価済のニューラルネットワークに係るパレート最適解を更新する。さらに、本開示の第二の実施形態は、ネットワークの生成とパレート最適解の更新を繰り返すことで、効率の良いネットワーク構造を探索することが可能である。

10

【 0 0 8 7 】

以下、上記のような効果を実現するための本開示の第二の実施形態の構成例について説明する。なお、以下の説明においては、第一の実施形態との差異について説明し、第一の実施形態と共通するクライアント端末 1、及び通信網 5 については説明を省略する。

【 0 0 8 8 】

< 2 - 1 . 構成例 >

図 1 1 は、本開示の第二の実施形態に係るサーバ 2 - 2 の構成例を説明するための説明図である。図 1 1 に示すように、サーバ 2 - 2 は、制御部 2 1、通信部 2 2、及び記憶部 2 4 を備える情報処理装置である。図 1 1 に示すように、本実施形態に係るサーバ 2 - 2 は、制御部 2 1 の機能構成が図 5 の制御部 2 0 の機能構成と一部異なる点で、図 5 のサーバ 2 と異なる。なお、図 1 1 に示す各構成のうち、図 5 に示した各構成と実質的に同様の構成については同一の符号を付してあるため、説明を省略する。以下では、本実施形態に係る制御部 2 1 が有する判定部 2 1 3、生成部 2 1 7、及び評価部 2 1 8、としての機能について説明する。

20

【 0 0 8 9 】

本実施形態に係る判定部 2 1 3 は、第一の実施形態で説明した判定部 2 0 3 と同様に、ユーザの入力に基づいて設計制御部 2 0 4 により設計された、または変更されたニューラルネットワークが、取得部 2 0 2 により取得される制約を満たすか否か判定を行う。また、本実施形態に係る判定部 2 1 3 は、後述する生成部 2 1 7 により生成されるニューラルネットワークが、取得部 2 0 2 により取得される制約を満たすか否か判定を行う。

30

【 0 0 9 0 】

生成部 2 1 7 は、元となるネットワークからネットワーク構造の異なる別のニューラルネットワークを生成する機能を有する。例えば、生成部 2 1 7 は、ユーザによる入力に基づいて設計され、判定部 2 1 3 により制約を満たすと判定されたニューラルネットワーク（以降、シードネットワークとも呼ぶ）から、ネットワーク構造の異なる別のニューラルネットワークを生成してもよい。また、生成部 2 1 7 は、パレート最適解に係るニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成してもよい。

【 0 0 9 1 】

本実施形態に係る生成部 2 1 7 によるニューラルネットワークの生成は、例えば、突然変異や交叉（または、交差、とも呼ぶ）などを含む遺伝的操作により実現されてもよい。ここで、上記の突然変異とは、生物に見られる遺伝子の突然変異をモデル化したものであってよい。すなわち、本実施形態では、ネットワークを構成する各レイヤーを遺伝子と見立て、レイヤーを突然変異させることで、ネットワーク構造の異なる別のニューラルネットワークを生成することができる。例えば、本実施形態に係る突然変異は、レイヤーの挿入、レイヤーの削除、レイヤー種類の変更、パラメータの変更、グラフ分岐、またはグラフ分岐の削除のうち少なくともいずれか一つを含んでもよい。

40

【 0 0 9 2 】

また、上記の交叉とは、生物の交配における染色体の部分的交換をモデル化したものであってよい。すなわち、本開示に係る情報処理方法では、2つのネットワークのレイヤー構

50

成を部分的に交換することで、上記の別のニューラルネットワークを生成することができる。

【 0 0 9 3 】

なお、上記では、遺伝的操作により別のニューラルネットワークを生成する場合を例に説明したが、本実施形態に係るニューラルネットワークの生成方法は、係る例に限定されない。本実施形態に係る別のニューラルネットワークの生成は、例えば、入力されたネットワークのネットワーク構造を変化させるニューラルネットワークを用いて実現されてもよい。ニューラルネットワークの生成には、上記の例を含む種々の方法が適用され得る。

【 0 0 9 4 】

また、生成部 2 1 7 は、判定部により制約を満たすと判定されたニューラルネットワークが生成されるまで、別のニューラルネットワークの生成を繰り返してもよい。係る構成によれば、生成部 2 1 7 は、複数のハードウェアに係る制約を満たすようなニューラルネットワークを生成することが可能となる。

【 0 0 9 5 】

評価部 2 1 8 は、生成されたニューラルネットワークの評価結果を取得する機能を有する。評価部 2 1 8 は、例えば、生成されたニューラルネットワークを認識部 2 0 6 に実行させ、上記の評価結果を取得してもよい。なお、評価部 2 1 8 による評価結果の取得は上記に限定されず、通信網 5 を介して接続される各種のデバイスに生成されたニューラルネットワークを実行させ、評価結果を取得してもよい。

【 0 0 9 6 】

また、評価部 2 1 8 が取得する評価結果には、生成されたニューラルネットワークに係る演算量、及び学習誤差またはヴァリデーション誤差（以下、まとめて誤差と表現することがある）のうち少なくとも一方が含まれてよい。評価部 2 1 8 は、生成されたニューラルネットワークのネットワーク構造に基づいて、上記の演算量を取得することができる。なお、本実施形態に係る評価結果は、上記に限定されず、例えば、ハードウェアに係る使用メモリ量、発熱量、消費電力量、演算量から算出されるハードウェアのトータルコストや、サーバ費用などを含むトータルサービスコスト等を含んでもよい。評価部 2 1 8 は、予め記憶されたハードウェアやサービスに係る情報を基に、上記の値を算出することができる。

【 0 0 9 7 】

また、評価部 2 1 8 は、生成されたニューラルネットワークの評価結果に基づいて、評価済のニューラルネットワークに係るパレート最適解を更新する機能を有する。すなわち、評価部 2 1 8 は、生成部 2 1 7 が生成したニューラルネットワークの評価結果を取得し、当該評価結果に基づいてパレート最適解の更新を繰り返し実行する。

【 0 0 9 8 】

< 2 - 2 . 動作例 >

以上、本実施形態に係るサーバ 2 - 2 の構成例について説明した。続いて、本実施形態の動作例について、図 1 2 ~ 1 5 を参照して説明する。

【 0 0 9 9 】

図 1 2 は、本実施形態の動作例を説明するためのフローチャート図である。図 7 に示すステップ S 2 1 0 0、S 2 2 0 0 の処理は、図 6 を参照して説明したステップ S 1 1 0 1、S 1 1 0 2 の処理と同様であるため、説明を省略する。

【 0 1 0 0 】

続いて、ユーザによるニューラルネットワークの設計と、判定部 2 1 3 による判定が行われる（S 2 3 0 0）。なお、ステップ S 2 3 0 0 の処理は、例えば、図 6 を参照して説明したステップ S 1 1 0 3 ~ S 1 1 0 5 の処理、または図 7 を参照して説明したステップ S 1 2 0 3 ~ S 1 2 0 8 の処理と同様の処理を含んでもよい。

【 0 1 0 1 】

続いて、生成部 2 1 7 は、ステップ S 2 3 0 0 において、判定部 2 1 3 により制約を満たすと判定されたニューラルネットワーク（シードネットワーク）から、ネットワーク構造

10

20

30

40

50

の異なる別のニューラルネットワークを生成する（Ｓ２４００）。なお、ステップＳ２４００におけるニューラルネットワークの生成に係る詳細な処理フローについては図１３～１５を参照して後述する。

【０１０２】

続いて、評価部２１８が、生成されたニューラルネットワークの評価結果を取得する（Ｓ２５００）。評価部２１８のネットワーク構造の探索が終了していない場合（Ｓ２６００：ＮＯ）、処理はステップＳ２７００に進む。

【０１０３】

ステップＳ２７００において、評価部２１８は、生成されたニューラルネットワークの評価結果に基づいて、評価済のニューラルネットワークに係るパレート最適解を更新する。続いて、処理はステップＳ２４００に戻り、生成部２１７は、当該パレート最適解に係るニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成する。

10

【０１０４】

上記のステップＳ２４００～Ｓ２７００の処理により、ニューラルネットワークの生成と、パレート最適解の更新が繰り返し実行されて、ネットワーク構造の探索が終了すると（Ｓ２６００：ＹＥＳ）、処理は終了する。

【０１０５】

なお、ネットワーク構造の探索が終了した際、例えば誤差が最小（最高性能）であるニューラルネットワーク、または演算量が最小であるニューラルネットワーク、または中間解であるニューラルネットワークが得られてもよい。中間解の定義は、条件に応じて適宜設計されてよい。なお、上記のようなニューラルネットワークがユーザに提示されて、ユーザがいずれかのニューラルネットワークを選択してもよい。

20

【０１０６】

以上、本実施形態の動作例の全体的な処理フローについて説明した。続いて、図１２に示すステップＳ２４００のニューラルネットワーク生成に係る処理フローについて、図１３を参照して説明する。図１３はニューラルネットワーク生成に係る処理フローを説明するためのフローチャート図である。

【０１０７】

図１３を参照すると、まず、生成部２１７は、元となるニューラルネットワークに適用する別のニューラルネットワークの生成方法をランダムで決定する（Ｓ２４１０）。この際、元となるニューラルネットワークは、ユーザによる入力に基づいて設計され、判定部２１３により制約を満たすと判定されたシードネットワークであってもよい。また、元となるニューラルネットワークは、評価部２１８が更新したパレート最適解に係るニューラルネットワークから生成部２１７がランダムに選択したネットワークであってもよい。

30

【０１０８】

続いて、生成部２１７は、ステップＳ２４１０で選択した生成方法に基づいて、元となるニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成する。図１３に示す一例を参照すると、本実施形態に係る生成部２１７は、元となるニューラルネットワークを突然変異させることで、上記の別のニューラルネットワークを生成してもよい（Ｓ２４２０）。

40

【０１０９】

また、生成部２１７は、元となるニューラルネットワークを交叉させることで、上記の別のニューラルネットワークを生成してもよい（Ｓ２４３０）。ステップＳ２４２０及びステップＳ２４３０における突然変異と交叉の詳細な流れについては、それぞれ図１４、図１５を参照して後述する。

【０１１０】

続いて、生成部２１７は、ステップＳ２４２０またはステップＳ２４３０で生成したニューラルネットワークの整合性を判定する（Ｓ２４４０）。この際、生成部２１７は、生成したニューラルネットワークのレイヤー構成にエラーが生じているか否かを判定してもよ

50

い。生成部 217 は、例えば、Max - Pooling 処理に際し、入力されるデータが小さすぎる場合などに、ネットワークの整合性がない、と判定してよい。このように、生成したニューラルネットワークの整合性がないと判定した場合 (S2450 : No)、生成部 217 は、生成したニューラルネットワークを破棄し、ステップ S2410 に復帰する。

【0111】

一方、生成されたニューラルネットワークに整合性が認められる場合 (S2450 : Yes)、判定部 213 が、生成されたニューラルネットワークが、取得部 202 により取得された制約を満たすか否かを判定する。判定部 213 が、生成されたニューラルネットワークが制約を満たさないと判定した場合 (S2450 : NO)、生成部 217 は、生成したニューラルネットワークを破棄し、ステップ S2410 に復帰する。

10

【0112】

一方、生成されたニューラルネットワークが制約を満たすと判定した場合、生成部 217 は、生成したニューラルネットワークと、元となるニューラルネットワークと、の入出力が同一であるか否かを判定する (S2460)。ここで、両者の入出力が異なる場合 (S2460 : No)、想定する認識問題を処理することが困難となるため、生成部 217 は、生成したニューラルネットワークを破棄し、ステップ S2410 に復帰する。一方、生成されたニューラルネットワークと、元となるニューラルネットワークとの入出力が同一である場合 (S2460 : Yes)、生成部 217 は、ネットワーク生成に係る処理を正常に終了する。

20

【0113】

以上、本実施形態に係るニューラルネットワークの生成について説明した。上述したとおり、本実施形態に係る生成部 217 は、シードネットワークやパレート最適解に係るネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成することが可能である。なお、図 13 では、生成部 217 が突然変異または交叉を用いた遺伝的操作により別のニューラルネットワークを生成する場合を例に説明したが、本実施形態に係るネットワークの生成は係る例に限定されない。本実施形態に係る生成部 217 は、入力されたニューラルネットワークのネットワーク構造を変化させるニューラルネットワークを用いて、上記の別のニューラルネットワークを生成してもよい。生成部 217 によるニューラルネットワークの生成には、種々の手法が適用されてよい。

30

【0114】

続いて、本実施形態に係る突然変異によるネットワーク生成の流れについて説明する。図 14 は、生成部 217 による突然変異を用いたネットワーク生成を説明するためのフローチャートである。すなわち、図 14 に示すフローチャートは、図 13 に示したステップ S2420 における生成部 217 の詳細な制御を示している。図 14 を参照すると、本実施形態に係る突然変異は、レイヤーの挿入、レイヤーの削除、レイヤー種類の変更、パラメータの変更、グラフ分岐、グラフ分岐の削除を含んでよい。

【0115】

図 14 を参照すると、まず、生成部 217 は、元となるニューラルネットワークに適用する突然変異の手法をランダムで決定する (S2421)。続いて、生成部 217 は、ステップ S2421 で選択した手法に基づいて、元となるニューラルネットワークのネットワーク構造を変化させる。

40

【0116】

生成部 217 は、新規レイヤーを挿入する処理を行ってもよい (S2422)。生成部 217 は、例えば、元となるニューラルネットワークに、ReLU などの活性化関数を新たに挿入することで、ネットワーク構造の異なる別のニューラルネットワークを生成することができる。

【0117】

また、生成部 217 は、既存レイヤーを削除する処理を行ってもよい (S2423)。生成部 217 は、例えば、元となるニューラルネットワークから、Max - Pooling

50

に係るレイヤーを削除することで、ネットワーク構造の異なる別のニューラルネットワークを生成することができる。

【0118】

また、生成部217は、既存レイヤーのレイヤー種類を変更する処理を行ってもよい（S2424）。生成部217は、例えば、元となるニューラルネットワークに存在する活性化関数を別の活性化関数に置換することで、ネットワーク構造の異なる別のニューラルネットワークを生成することができる。

【0119】

また、生成部217は、既存レイヤーに係るパラメータを変更する処理を行ってもよい（S2425）。生成部217は、例えば、既存するConvolutionレイヤーのカーネルシェイプを変更することで、ネットワーク構造の異なる別のニューラルネットワークを生成することができる。

10

【0120】

また、生成部217は、新たなグラフ分岐を作成する処理を行ってもよい（S2426）。生成部217は、例えば、既存レイヤーの一部をコピーすることでグラフ分岐を作成し、当該グラフ分岐の結合部としてConcatenateレイヤーを挿入することで、別のニューラルネットワークを生成することができる。

【0121】

また、生成部217は、既存のグラフ分岐を削除する処理を行ってもよい（S2427）。生成部217は、例えば、既存するグラフ分岐の1ルートを削除し、当該削除により分岐が消失した場合にはConcatenateレイヤーも削除することで、別のニューラルネットワークを生成することができる。

20

【0122】

以上、本実施形態に係る生成部217による突然変異を用いたネットワーク生成について説明した。なお、上記では、生成部217がランダムで選択したステップS2422～S2427の処理を実行する場合を例に説明したが、本実施形態に係る突然変異の制御は、係る例に限定されない。生成部217は、ステップS2422～S2427に係る処理を同時に2つ以上行ってもよいし、ステップS2422～S2427の実行判断をそれぞれ実施してもよい。また、生成部217は、図14の例に示した以外の処理を実行してもよい。生成部217による突然変異の制御は、柔軟に変更され得る。

30

【0123】

続いて、本実施形態に係る交叉によるネットワーク生成の流れについて説明する。図15は、生成部217による交叉を用いたネットワーク生成を説明するためのフローチャート図である。すなわち、図15に示すフローチャートは、図13に示したステップS2430における生成部217の詳細な制御を示している。

【0124】

図15を参照すると、まず、生成部217は、交叉を実行するために、元となる2つのネットワークを選択する（S2431）。ここで、生成部217は、ユーザによる入力に基づいて設計され、判定部213により制約を満たすと判定された2つのシードネットワークを選択してもよい。また、生成部217は、ユーザによる入力に基づいて設計され、判定部213により制約を満たすと判定された1つのシードネットワークと、予め登録された交叉用のネットワークと、を選択することもできる。さらには、生成部217は、ユーザによる入力に基づいて設計され、判定部213により制約を満たすと判定されたシードネットワークから突然変異により生成した別のニューラルネットワークを選択してもよい。

40

【0125】

続いて、生成部217は、ステップS2431で選択した2つのネットワークを交叉させ、ネットワーク構造の異なる別のニューラルネットワークを生成する（S2432）。この際、生成部217は、種々の手法により交叉を実行してよい。生成部217は、例えば、一点交叉、二点交叉、多点交叉、一様交叉などにより、上記の別のニューラルネットワークを生成することができる。

50

【 0 1 2 6 】

以上、本実施形態に係るニューラルネットワークの生成について説明した。上述したとおり、本実施形態に係る生成部 2 1 7 は、突然変異及び交叉を含む遺伝的操作などにより、元となるニューラルネットワークからネットワーク構造の異なる別のニューラルネットワークを生成することができる。すなわち、本実施形態に係る情報処理方法では、生成部 2 1 7 が生成したニューラルネットワークの評価結果に基づいてパレート最適解の更新を繰り返すことで、より効率のよいネットワーク構造を探索することが可能となる。

【 0 1 2 7 】

< 2 - 3 . 効果 >

以上説明したように、本開示の第二の実施形態によれば、設計されたニューラルネットワークの評価結果に基づいて、ネットワーク構造の異なる別のニューラルネットワークを生成することができる。また、本実施形態は、ネットワークの生成と、パレート最適解の更新を繰り返すことで、効率の良いネットワーク構造を探索することが可能である。また、本実施形態において、生成されるニューラルネットワークは、判定部 2 1 3 により、複数のハードウェアに係る制約を満たすと判定されたニューラルネットワークである。したがって、本実施形態によれば、複数のハードウェアによる処理に適合し、かつ効率の良いネットワーク構造を探索することが可能である。

【 0 1 2 8 】

< < 3 . 変形例 > >

以上、本開示の実施形態を説明した。以下では、本開示に係る幾つかの変形例を説明する。なお、以下に説明する各変形例は、単独で各実施形態に適用されてもよいし、組み合わせで各実施形態に適用されてもよい。また、各変形例は、上記実施形態で説明した構成に代えて適用されてもよいし、上記実施形態で説明した構成に対して追加的に適用されてもよい。

【 0 1 2 9 】

< 3 - 1 . 変形例 1 >

上記実施形態ではビジュアルプログラミングによりニューラルネットワークの設計が行われる例を説明したが、本技術は係る例に限定されない。例えば、本技術に係るニューラルネットワークの設計は、テキストによるプログラミングや、CUI (Command User Interface) 上の操作により行われてもよい。また、ハードウェアに係る制約や、ハードウェアとレイヤーの対応付けも、上記で説明された例に限定されず、テキスト、またはCUIにより入力されてもよい。

【 0 1 3 0 】

< 3 - 2 . 変形例 2 >

また、上記実施形態では、図 2、図 5、及び図 1 1 を参照してクライアント端末 1、サーバ 2、及びサーバ 2 - 2 が有する機能を説明したが、本技術は係る例に限定されない。上記実施形態で説明したクライアント端末 1 の機能をサーバ 2、またはサーバ 2 - 2 が有してもよいし、上記実施形態で説明したサーバ 2、またはサーバ 2 - 2 の機能をクライアント端末 1 が有してもよい。

【 0 1 3 1 】

例えば、クライアント端末 1 が、図 5 を参照して説明した取得部 2 0 2、判定部 2 0 3、及び設計制御部 2 0 4 の機能を有し、クライアント端末 1 を用いて設計されたニューラルネットワークの情報が、サーバ 2 に提供されてもよい。

【 0 1 3 2 】

< < 4 . ハードウェア構成例 > >

以上、本開示の各実施形態を説明した。上述した表示制御処理、通信制御処理、取得処理、判定処理、設計制御処理、学習処理、認識処理、ネットワーク生成処理、評価結果取得処理等の情報処理は、ソフトウェアと、クライアント端末 1、サーバ 2、2 - 2 との協働により実現される。以下では、本実施形態に係るサーバ 2 のハードウェア構成例について説明する。

10

20

30

40

50

【0133】

図16は、サーバ2のハードウェア構成を示した説明図である。図16に示したように、サーバ2は、CPU(Central Processing Unit)2001と、DSP(Digital Signal Processor)2002と、GPU(Graphics Processing Unit)2003と、ニューロチップ2004と、ROM(Read Only Memory)2005と、RAM(Random Access Memory)2006と、入力装置2007と、出力装置2008と、ストレージ装置2009と、ドライブ2010と、通信装置2011とを備える。

【0134】

CPU2001は、演算処理装置および制御装置として機能し、各種プログラムに従ってサーバ2内の動作全般を制御する。また、CPU2001は、マイクロプロセッサであってもよい。また、DSP2002、GPU2003、ニューロチップ2004は、演算処理装置として機能する。例えば、CPU2001、DSP2002、GPU2003、及びニューロチップ2004は、本開示において、ニューラルネットワークによる認識処理を実行するハードウェアであってもよい。ROM2005は、CPU2001が使用するプログラムや演算パラメータなどを記憶する。RAM2006は、CPU2001、DSP2002、GPU2003、ニューロチップ2004の実行において使用するプログラムや、その実行において適宜変化するパラメータなどを一時記憶する。これらはCPUバスなどから構成されるホストバスにより相互に接続されている。主に、CPU2001、ROM2005及びRAM2006とソフトウェアとの協働により、通信制御部201、取得部202、判定部203、設計制御部204、学習部205の機能が実現される。

【0135】

入力装置2007は、マウス、キーボード、タッチパネル、ボタン、マイクロフォン、スイッチおよびレバーなどユーザが情報を入力するための入力手段と、ユーザによる入力に基づいて入力信号を生成し、CPU2001に出力する入力制御回路などから構成されている。サーバ2のユーザは、該入力装置2007を操作することにより、サーバ2に対して各種のデータを入力したり処理動作を指示したりすることができる。

【0136】

出力装置2008は、例えば、液晶ディスプレイ(LCD)装置、OLED(Organic Light Emitting Diode)装置およびランプなどの表示装置を含む。さらに、出力装置2008は、スピーカおよびヘッドホンなどの音声出力装置を含む。例えば、表示装置は、撮像された画像や生成された画像などを表示する。一方、音声出力装置は、音声データなどを音声に変換して出力する。

【0137】

ストレージ装置2009は、本実施形態にかかるサーバ2の記憶部24の一例として構成されたデータ格納用の装置である。ストレージ装置2009は、記憶媒体、記憶媒体にデータを記録する記録装置、記憶媒体からデータを読み出す読出し装置および記憶媒体に記録されたデータを削除する削除装置などを含んでもよい。このストレージ装置2009は、CPU2001が実行するプログラムや各種データを格納する。

【0138】

ドライブ2010は、記憶媒体用リーダライタであり、サーバ2に内蔵、あるいは外付けされる。ドライブ2010は、装着されている磁気ディスク、光ディスク、光磁気ディスク、または半導体メモリなどのリムーバブル記憶媒体に記録されている情報を読み出して、RAM2005に出力する。また、ドライブ2010は、リムーバブル記憶媒体に情報を書き込むこともできる。

【0139】

通信装置2011は、通信デバイスなどで構成された通信インタフェースである。また、通信装置2011は、無線LAN(Local Area Network)対応通信装置であっても、LTE(Long Term Evolution)対応通信装置であっても、有線による通信を行うワイヤー通信装置であってもよい。通信装置2011は、サーバ

10

20

30

40

50

２の通信部２２に対応する。

【０１４０】

なお、上記ではサーバ２のハードウェア構成を説明したが、クライアント端末１及び第二の実施形態によるサーバ２－２も、サーバ２と同様に、ＣＰＵ２００１、ＲＯＭ２０５およびＲＡＭ２０６などに相当するハードウェアを有する。そして、クライアント端末１のハードウェアとソフトウェアとの協働により例えば制御部１０の機能が実現される。また、例えば第二の実施形態によるサーバ２－２、のハードウェアとソフトウェアとの協働により判定部２１３、生成部２１７、及び評価部２１８に相当する機能が実現される。

【０１４１】

<<５．むすび>>

以上、説明したように、本開示の実施形態によれば、複数のハードウェアによる処理に適したニューラルネットワークをより効率的に設計することが可能である。

【０１４２】

以上、添付図面を参照しながら本開示の好適な実施形態について詳細に説明したが、本開示の技術的範囲はかかる例に限定されない。本開示の技術分野における通常の知識を有する者であれば、特許請求の範囲に記載された技術的思想の範疇内において、各種の変更例または修正例に想到し得ることは明らかであり、これらについても、当然に本開示の技術的範囲に属するものと了解される。

【０１４３】

例えば、上記実施形態における各ステップは、必ずしもフローチャート図として記載された順序に沿って時系列に処理する必要はない。例えば、上記実施形態の処理における各ステップは、フローチャート図として記載した順序と異なる順序で処理されても、並列的に処理されてもよい。例えば、図６において、ハードウェアに係る制約を入力または取得するステップＳ１１０１の処理、及びハードウェアごとに範囲を設定するステップＳ１１０２の処理は、ニューラルネットワークの設計に係るステップＳ１１０３の処理の後に行われてもよい。

【０１４４】

また、上記実施形態によれば、ＣＰＵ２００１、ＲＯＭ２０５、及びＲＡＭ２０６等のハードウェアを、上述したクライアント端末１、サーバ２、及びサーバ２－２の各構成と同様の機能を発揮させるためのコンピュータプログラムも提供可能である。また、該コンピュータプログラムが記録された記録媒体も提供される。

【０１４５】

また、本明細書に記載された効果は、あくまで説明的または例示的なものであって限定的ではない。つまり、本開示に係る技術は、上記の効果とともに、または上記の効果に代えて、本明細書の記載から当業者には明らかな他の効果を奏しうる。

【０１４６】

なお、以下のような構成も本開示の技術的範囲に属する。

(１)

複数のハードウェアに係る制約を取得する取得部と、
ニューラルネットワークが、前記制約を満たすか否かを判定を行う判定部と、
を備える情報処理装置。

(２)

ユーザの入力に基づく前記ニューラルネットワークの設計を制御する設計制御部をさらに備える、前記(１)に記載の情報処理装置。

(３)

前記判定部は、前記設計制御部により前記ニューラルネットワークが変更された場合に、前記判定を行う、前記(２)に記載の情報処理装置。

(４)

前記設計制御部は、前記判定部により、前記制約を満たさないと判定された場合に、前記制約を満たすように、前記変更に係るレイヤーの再配置を行う、前記(３)に記載の情報

10

20

30

40

50

処理装置。

(5)

前記設計制御部は、前記判定部により、前記制約を満たすと判定された場合に、前記ニューラルネットワークを構築するプログラムを作成する、前記(2) ~ (4)のいずれか一項に記載の情報処理装置。

(6)

前記判定部により、前記制約を満たさないと判定された場合に、前記制約が満たされないことを示す警告画面を表示させる、表示制御部をさらに備える、前記(1) ~ (5)のいずれか一項に記載の情報処理装置。

(7)

前記警告画面は、前記ニューラルネットワークにおいて、前記制約を満たさない部分を提示する、前記(6)に記載の情報処理装置。

(8)

前記制約に基づいて、前記ニューラルネットワークにおけるレイヤーごとに、前記レイヤーに対応付けられたハードウェアに応じた学習を行う学習部をさらに備える、前記(1) ~ (7)のいずれか一項に記載の情報処理装置。

(9)

前記判定部により前記制約を満たすと判定されたニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成する生成部と、
生成されたニューラルネットワークの評価結果を取得する評価部と、をさらに備え、
前記評価部は、生成されたニューラルネットワークの評価結果に基づいて、評価済のニューラルネットワークに係るパレート最適解を更新し、
前記生成部は、前記パレート最適解に係るニューラルネットワークから、ネットワーク構造の異なる別のニューラルネットワークを生成する、
前記(1) ~ (8)のいずれか一項に記載の情報処理装置。

(10)

前記判定部は、前記生成部により生成されるニューラルネットワークが、前記制約を満たすか否か判定を行い、
前記生成部は、前記判定部により前記制約を満たすと判定されたニューラルネットワークが生成されるまで、前記別のニューラルネットワークの生成を繰り返す、前記(9)に記載の情報処理装置。

(11)

前記生成部は、遺伝的操作により、前記別のニューラルネットワークを生成する、前記(9)または(10)に記載の情報処理装置。

(12)

前記遺伝的操作は、突然変異または交叉のうち少なくとも一方を含む、前記(11)に記載の情報処理装置。

(13)

前記突然変異は、レイヤーの挿入、レイヤーの削除、レイヤー種類の変更、パラメータの変更、グラフ分岐、またはグラフ分岐の削除のうち少なくともいずれか一つを含む、前記(12)に記載の情報処理装置。

(14)

前記情報処理装置は、前記判定部による判定結果を送信させる通信制御部をさらに備える、前記(1) ~ (13)のいずれか一項に記載の情報処理装置。

(15)

ニューラルネットワークが複数のハードウェアに係る制約を満たすか否かの判定結果を受信する受信部と、
前記判定結果に基づいて処理を行う処理部と、
を備える情報処理装置。

(16)

10

20

30

40

50

複数のハードウェアに係る制約を取得することと、
ニューラルネットワークが、前記制約を満たすか否か判定を行うことと、
を含む情報処理方法。

【符号の説明】

【 0 1 4 7 】

1 クライアント端末

2、2 - 2 サーバ

5 通信網

1 0 制御部

1 2 通信部

1 4 表示部

1 6 操作部

2 0、2 1 制御部

2 2 通信部

2 4 記憶部

2 0 1 通信制御部

2 0 2 取得部

2 0 3、2 1 3 判定部

2 0 4 設計制御部

2 0 5 学習部

2 0 6 認識部

2 1 7 生成部

2 1 8 評価部

1 0 0 0 情報処理システム

10

20

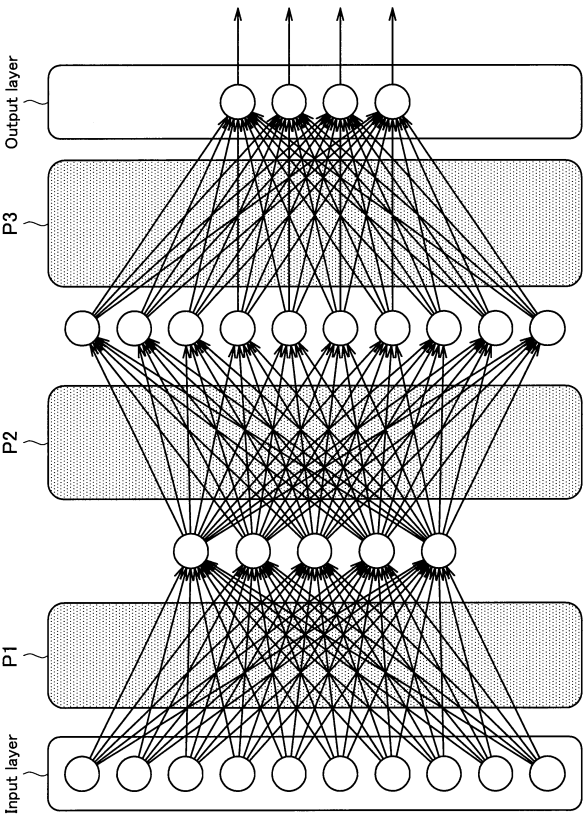
30

40

50

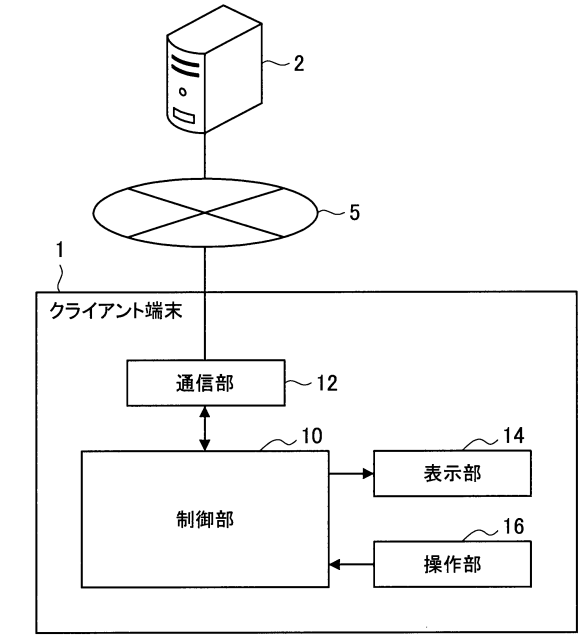
【図面】

【図 1】



【図 2】

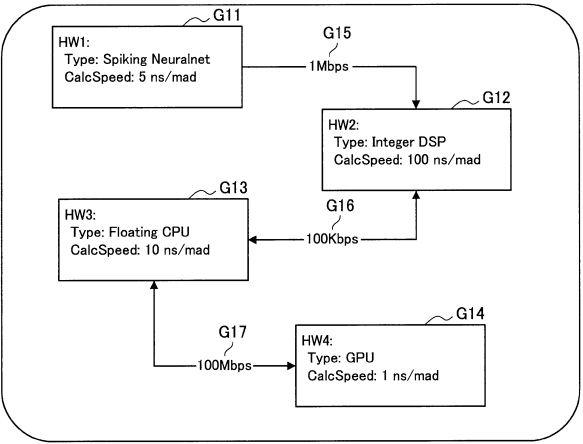
1000



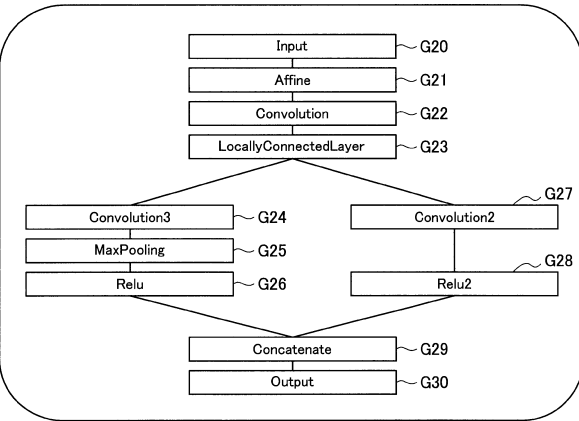
10

20

【図 3】



【図 4】

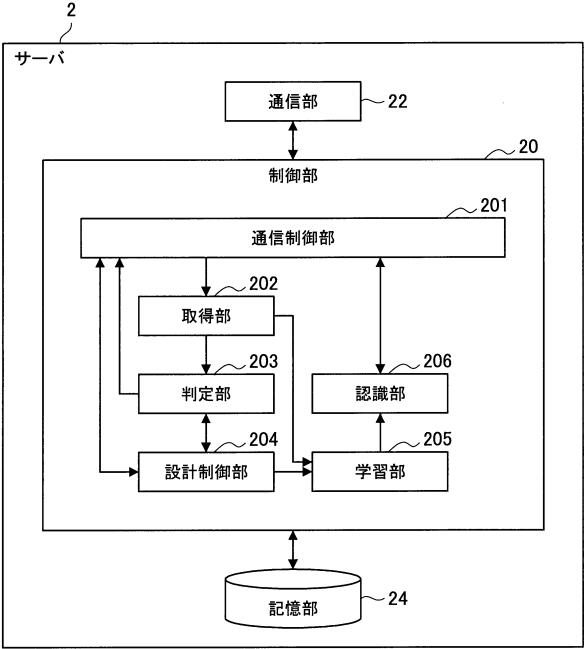


30

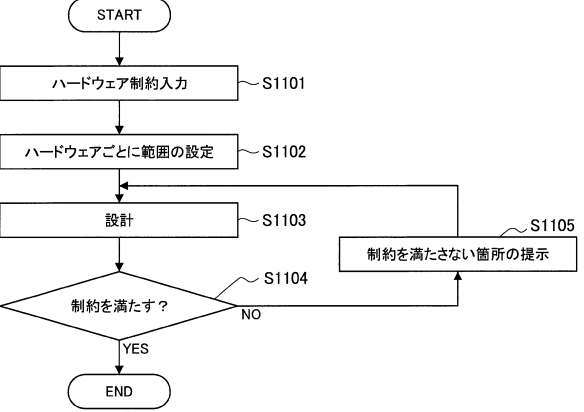
40

50

【図 5】



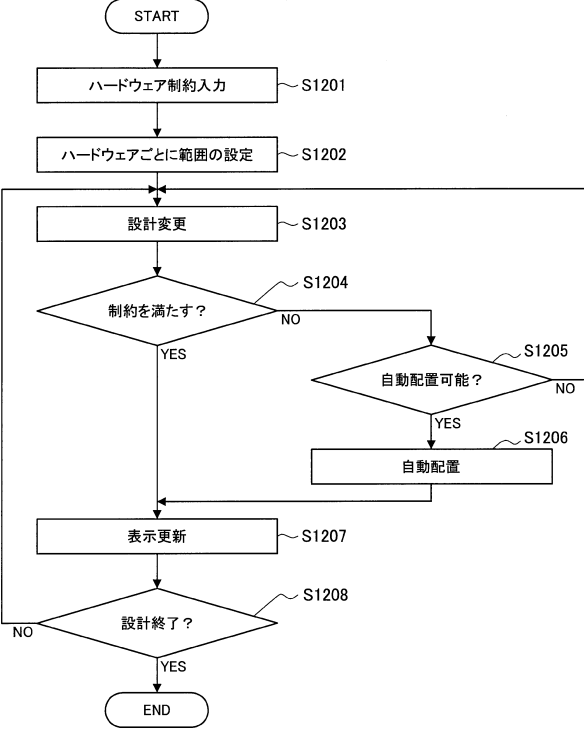
【図 6】



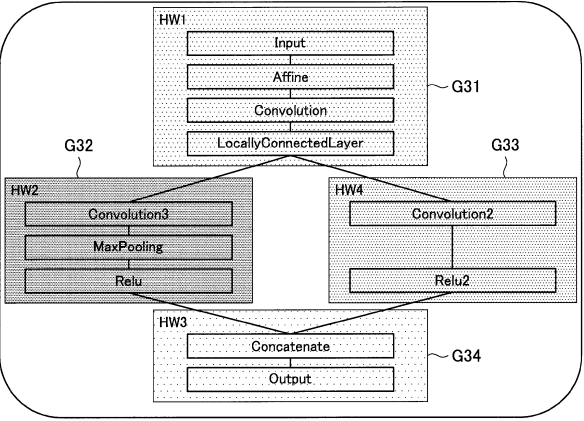
10

20

【図 7】



【図 8】

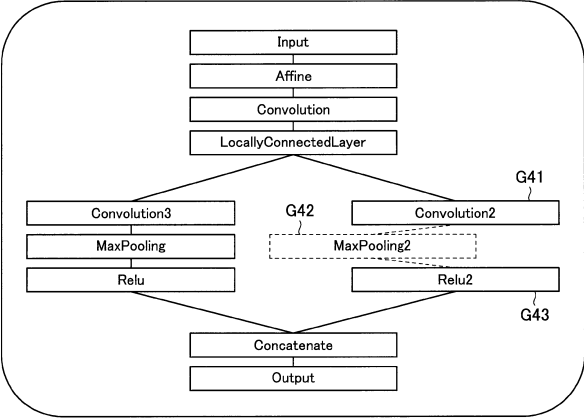


30

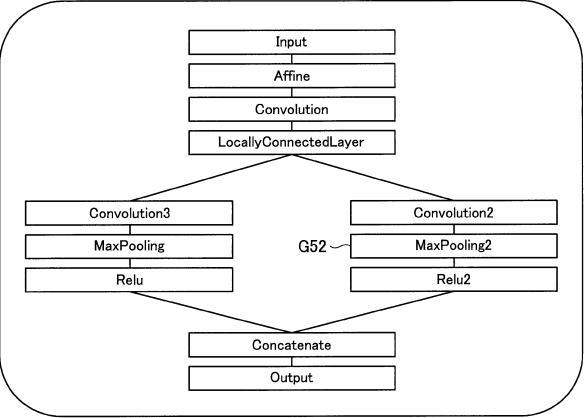
40

50

【図 9】

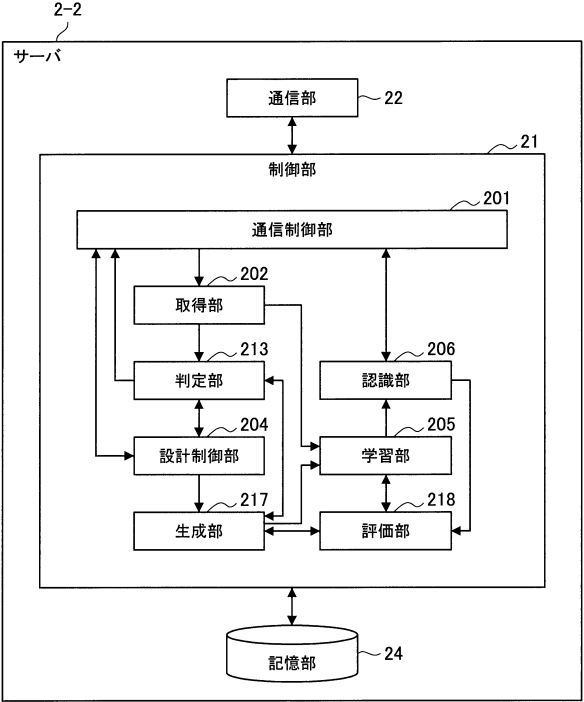


【図 10】

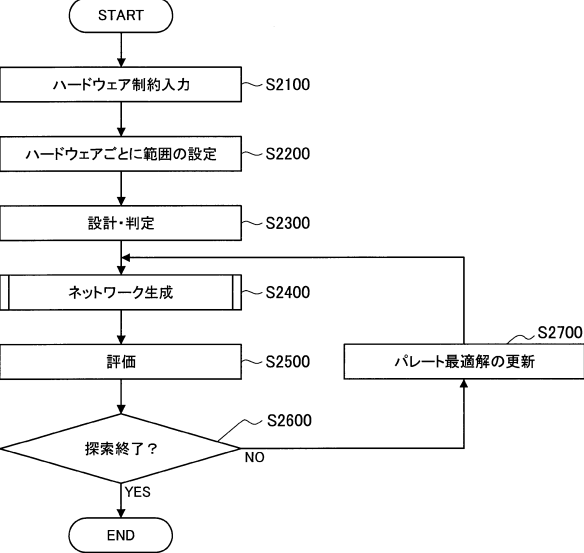


10

【図 11】



【図 12】



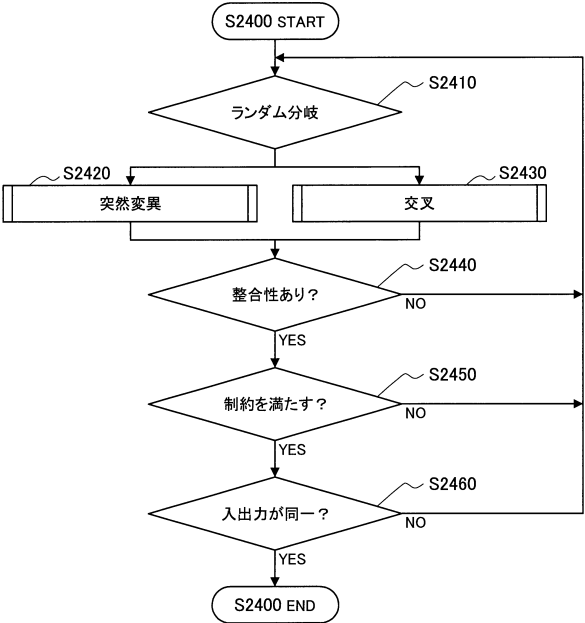
20

30

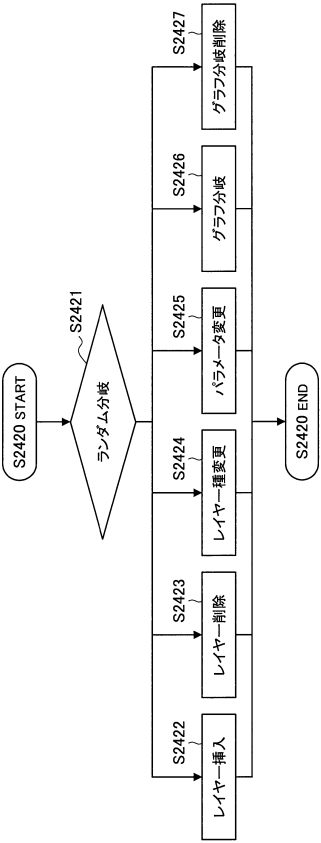
40

50

【図 1 3】



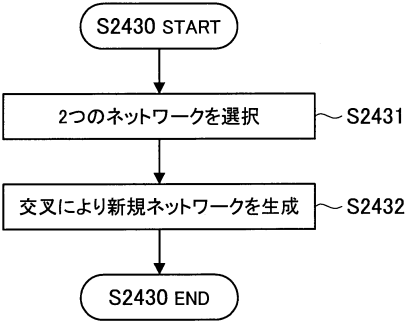
【図 1 4】



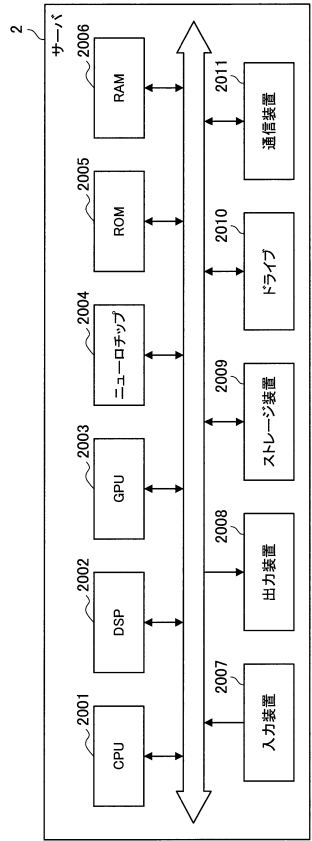
10

20

【図 1 5】



【図 1 6】



30

40

50

フロントページの続き

- 東京都港区港南 1 丁目 7 番 1 号 ソニー株式会社内
- (72)発明者 小林 由幸
東京都港区港南 1 丁目 7 番 1 号 ソニー株式会社内
- 審査官 多賀 実
- (56)参考文献 特開平 0 3 - 2 3 7 5 5 7 (J P , A)
特開平 1 0 - 0 9 1 6 7 6 (J P , A)
米国特許出願公開第 2 0 0 8 / 0 3 1 9 9 3 3 (U S , A 1)
片山 立 ほか, 「複数評価基準にもとづくファジィモデルとニューロモデルの総合評価」,
日本ファジィ学会誌, 日本ファジィ学会, 1992年10月31日, 第4巻, 第5号, pp.942-957
JIN, Yaochu ほか, "Neural network regularization and ensembling using multi-objective evolutionary algorithms", Proceedings of the 2004 Congress on Evolutionary Computation ,
米国, IEEE, 2004年06月23日, pp.1-8
- (58)調査した分野 (Int.Cl., D B 名)
G 0 6 F 1 1 / 3 4 - 1 1 / 3 6
G 0 6 F 3 0 / 0 0 - 3 0 / 3 9 8
G 0 6 N 3 / 0 2 - 3 / 1 2