

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2023-83207
(P2023-83207A)

(43)公開日 令和5年6月15日(2023.6.15)

(51)国際特許分類 F I
 G 0 6 N 3/04 (2023.01) G 0 6 N 3/04
 G 0 6 N 3/084(2023.01) G 0 6 N 3/08 1 4 0

審査請求 未請求 請求項の数 20 O L (全19頁)

(21)出願番号	特願2022-129509(P2022-129509)	(71)出願人	390019839 三星電子株式会社 Samsung Electronics Co., Ltd. 大韓民国京畿道水原市靈通区三星路129 129, Samsung-ro, Yeongtong-gu, Suwon-si, Gyeonggi-do, Republic of Korea
(22)出願日	令和4年8月16日(2022.8.16)	(74)代理人	100107766 弁理士 伊東 忠重
(31)優先権主張番号	10-2021-0171979	(74)代理人	100070150 弁理士 伊東 忠彦
(32)優先日	令和3年12月3日(2021.12.3)	(74)代理人	100135079
(33)優先権主張国・地域又は機関	韓国(KR)		

最終頁に続く

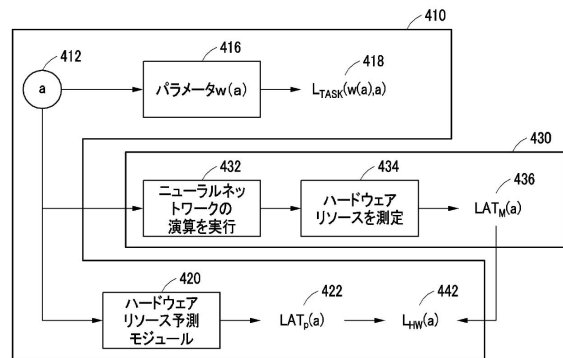
(54)【発明の名称】 ニューラルネットワークの最適なアーキテクチャーを探索する装置及び方法

(57)【要約】

【課題】ニューラルネットワークの最適なアーキテクチャーを探索する方法及び装置を提供すること。

【解決手段】ニューラルネットワークの最適なアーキテクチャーを探索する装置は、プロセッサを含み、プロセッサは、ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定し、候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定し、ハードウェアリソース予測モジュールを用いて、候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測し、測定されたハードウェアリソースと予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定し、ニューラルネットワーク損失とハードウェアリソース損失に基づいてニューラルネットワークのターゲットアーキテクチャーを決定することができる。

【選択図】 図4



【特許請求の範囲】**【請求項 1】**

ニューラルネットワークの最適なアーキテクチャーを探索する装置であって、
プロセッサを含み、
前記プロセッサは、
前記ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定し、
前記候補アーキテクチャーのニューラルネットワークが動作するとき要求されるハードウェアリソースを測定し、
ハードウェアリソース予測モジュールを用いて前記候補アーキテクチャーのニューラルネットワークが動作するとき要求されるハードウェアリソースを予測し、
前記測定されたハードウェアリソースと前記予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定し、
前記ニューラルネットワーク損失と前記ハードウェアリソース損失に基づいて前記ニューラルネットワークのターゲットアーキテクチャーを決定する、
装置。

10

【請求項 2】

前記ハードウェアリソース予測モジュールは、前記候補アーキテクチャーのパラメータを入力とし、前記入力されたパラメータに基づいて、前記候補アーキテクチャーのニューラルネットワークの前記ハードウェアリソースを予測してハードウェアリソース予測値を出力するニューラルネットワークである、請求項 1 に記載の装置。

20

【請求項 3】

前記プロセッサは、
前記測定されたハードウェアリソースと前記予測されたハードウェアリソースとの間の差に基づいて前記ハードウェアリソース損失を決定し、
前記ハードウェアリソース損失が最小になるように前記候補アーキテクチャーのパラメータをアップデートする、請求項 1 又は 2 に記載の装置。

【請求項 4】

前記プロセッサは、前記候補アーキテクチャー及び前記候補アーキテクチャーのパラメータによる前記ニューラルネットワーク損失と、前記候補アーキテクチャーのパラメータによる前記ハードウェアリソース損失の加重和を最適化損失として決定し、前記最適化損失を最小にする前記ターゲットアーキテクチャーを決定する、請求項 1 に記載の装置。

30

【請求項 5】

前記プロセッサは、前記ニューラルネットワーク損失と前記ハードウェアリソース損失を減少させる前記ターゲットアーキテクチャーとターゲットパラメータを決定する、請求項 1 に記載の装置。

【請求項 6】

前記ニューラルネットワークの各レイヤごとに、各レイヤが有する候補演算のいずれか 1 つの候補演算を選択することによって前記候補アーキテクチャーが決定される、請求項 1 乃至 5 のいずれか一項に記載の装置。

40

【請求項 7】

前記ハードウェアリソース予測モジュールには、前記ニューラルネットワークの各レイヤで実行可能な候補演算のうち、前記候補アーキテクチャーを構成する選択された候補演算に関する情報が入力される、請求項 1 乃至 6 のいずれか一項に記載の装置。

【請求項 8】

前記ハードウェアリソースは、前記候補アーキテクチャーのニューラルネットワークが動作するときの電力消費、メモリ要求量、演算の数、及び処理時間のうち少なくとも 1 つを含む、請求項 1 乃至 7 のいずれか一項に記載の装置。

【請求項 9】

前記プロセッサは、

50

前記ニューラルネットワーク損失と前記ハードウェアリソース損失を含む最適化損失を決定し、

前記候補アーキテクチャーのニューラルネットワークに含まれている各レイヤの候補演算のうち、前記最適化損失を最小にするターゲット演算を選択することによって前記ターゲットアーキテクチャーを決定する、請求項 1 に記載の装置。

【請求項 10】

前記プロセッサは、前記候補アーキテクチャーのニューラルネットワークが学習データを処理して導き出した結果データと検証データとの間の差に基づいて前記ニューラルネットワーク損失を決定する、請求項 1 乃至 9 のいずれか一項に記載の装置。

【請求項 11】

拡張現実提供装置であって、
ターゲットアーキテクチャーを有するニューラルネットワークを用いて処理動作を行うプロセッサを含み、

前記プロセッサは、
前記ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定し、

前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定し、

ハードウェアリソース予測モジュールを用いて、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測し、

前記測定されたハードウェアリソースと前記予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定し、

前記ニューラルネットワーク損失と前記ハードウェアリソース損失に基づいて前記ターゲットアーキテクチャーを決定する、

拡張現実提供装置。

【請求項 12】

前記拡張現実提供装置は、映像を撮影するカメラをさらに含み、

前記プロセッサは、前記ターゲットアーキテクチャーを有するニューラルネットワークを用いて前記映像を処理することによって拡張現実コンテンツを生成する、請求項 11 に記載の拡張現実提供装置。

【請求項 13】

前記ハードウェアリソース予測モジュールは、前記候補アーキテクチャーのパラメータを入力とし、前記入力されたパラメータに基づいて、前記候補アーキテクチャーのニューラルネットワークの前記ハードウェアリソースを予測してハードウェアリソース予測値を出力するニューラルネットワークである、請求項 11 又は 12 に記載の拡張現実提供装置。

【請求項 14】

前記プロセッサは、前記候補アーキテクチャー及び前記候補アーキテクチャーのパラメータによる前記ニューラルネットワーク損失と、前記候補アーキテクチャーのパラメータによる前記ハードウェアリソース損失の加重和を最小にする前記ターゲットアーキテクチャーを決定する、請求項 11 乃至 13 のいずれか一項に記載の拡張現実提供装置。

【請求項 15】

ニューラルネットワークの最適なアーキテクチャーを探索する方法であって、

前記ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定する動作と、

前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定する動作と、

ハードウェアリソース予測モジュールを用いて、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測する動作と、

前記測定されたハードウェアリソースと前記予測されたハードウェアリソースに基づい

10

20

30

40

50

てハードウェアリソース損失を決定する動作と、

前記ニューラルネットワーク損失と前記ハードウェアリソース損失に基づいて前記ニューラルネットワークのターゲットアーキテクチャーを決定する動作と、
を含む方法。

【請求項 16】

前記ハードウェアリソース予測モジュールは、前記候補アーキテクチャーのパラメータを入力とし、前記入力されたパラメータに基づいて、前記候補アーキテクチャーのニューラルネットワークの前記ハードウェアリソースを予測してハードウェアリソース予測値を出力するニューラルネットワークである、請求項 15 に記載の方法。

【請求項 17】

前記ターゲットアーキテクチャーを決定する動作は、前記候補アーキテクチャー及び前記候補アーキテクチャーのパラメータによる前記ニューラルネットワーク損失と、前記候補アーキテクチャーのパラメータによる前記ハードウェアリソース損失の加重和を最小にする前記ターゲットアーキテクチャーを決定する動作を含む、請求項 15 又は 16 に記載の方法。

【請求項 18】

前記ターゲットアーキテクチャーを決定する動作は、前記候補アーキテクチャーのニューラルネットワークに含まれる各レイヤの候補演算のうちターゲット演算を選択することによって前記ターゲットアーキテクチャーを決定する動作を含む、請求項 15 又は 16 に記載の方法。

【請求項 19】

前記ハードウェアリソースは、前記候補アーキテクチャーのニューラルネットワークが動作するときの電力消費、メモリ要求量、演算の数及び処理時間のうち少なくとも 1 つを含む、請求項 15 乃至 18 のいずれか一項に記載の方法。

【請求項 20】

請求項 15 乃至 19 のいずれか一項に記載の方法を行うための命令語を含む 1 つ以上のコンピュータプログラムを格納したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

以下の開示は、ニューラルネットワークの最適なアーキテクチャーを探索する技術に関する。

【背景技術】

【0002】

ニューラルアーキテクチャー探索 (neural architecture search、NAS) は、所定の目的のニューラルネットワークにおける最適なアーキテクチャーを自動で探すための方法論のうちの一つである。NAS は、特定の問題を解決するための最も適切なニューラルネットワークのアーキテクチャーの構造及び形態を、ディープラーニングを介して探索する方法である。NAS におけるニューラルネットワークは、探索空間と呼ぶ、予め定義された演算子及び関数で構成されたプリミティブ演算 (primitive operations) を選択及び組み合わせることで生成されてもよい。ここで、演算子の例示として、畳み込み (convolution)、プーリング (pooling)、併合 (concatenation)、スキップ接続 (skip connection) などが挙げられる。

【発明の概要】

【発明が解決しようとする課題】

【0003】

本発明の目的は、ニューラルネットワークの最適なアーキテクチャーを探索する技術を提供する。

【課題を解決するための手段】

10

20

30

40

50

【 0 0 0 4 】

一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する装置は、プロセッサを含み、前記プロセッサは、前記ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定し、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定し、ハードウェアリソース予測モジュールを用いて前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測し、前記測定されたハードウェアリソースと前記予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定し、前記ニューラルネットワーク損失と前記ハードウェアリソース損失に基づいて前記ニューラルネットワークのターゲットアーキテクチャーを決定する。

10

【 0 0 0 5 】

前記ハードウェアリソース予測モジュールは、前記候補アーキテクチャーのパラメータを入力とし、前記入力されたパラメータに基づいて、前記候補アーキテクチャーのニューラルネットワークの前記ハードウェアリソースを予測してハードウェアリソース予測値を出力するニューラルネットワークであってもよい。

【 0 0 0 6 】

前記プロセッサは、前記候補アーキテクチャー及び前記候補アーキテクチャーのパラメータによる前記ニューラルネットワーク損失と、前記候補アーキテクチャーのパラメータによる前記ハードウェアリソース損失の加重和を最適化損失として決定し、前記最適化損失を最小にする前記ターゲットアーキテクチャーを決定することができる。

20

【 0 0 0 7 】

前記ハードウェアリソースは、前記候補アーキテクチャーのニューラルネットワークが動作するときの電力消費、メモリ要求量、演算の数、及び処理時間のうち少なくとも1つを含むことができる。

【 0 0 0 8 】

前記プロセッサは、前記ニューラルネットワーク損失と前記ハードウェアリソース損失を含む最適化損失を決定し、前記候補アーキテクチャーのニューラルネットワークに含まれている各レイヤの候補演算のうち、前記最適化損失を最小にするターゲット演算を選択することによって前記ターゲットアーキテクチャーを決定することができる。

30

【 0 0 0 9 】

一実施形態に係る拡張現実提供装置は、ターゲットアーキテクチャーを有するニューラルネットワークを用いて処理動作を行うプロセッサを含み、前記プロセッサは、前記ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定し、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定し、ハードウェアリソース予測モジュールを用いて、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測し、前記測定されたハードウェアリソースと前記予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定し、前記ニューラルネットワーク損失と前記ハードウェアリソース損失に基づいて前記ターゲットアーキテクチャーを決定する。

40

【 0 0 1 0 】

一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する方法は、前記ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定する動作と、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定する動作と、ハードウェアリソース予測モジュールを用いて、前記候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測する動作と、前記測定されたハードウェアリソースと前記予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定する動作と、前記ニューラルネットワーク損失と前記ハードウェア

50

リソース損失に基づいて前記ニューラルネットワークのターゲットアーキテクチャーを決定する動作とを含む。

【発明の効果】

【0011】

一実施形態によれば、ハードウェアリソース制限事項を考慮してニューラルネットワークの最適なアーキテクチャーを探索することができ、短い時間で最適化を行うことができる。

【図面の簡単な説明】

【0012】

【図1】一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する探索フレームワークを説明するための図である。 10

【図2】一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する探索装置の構成を示すブロック図である。

【図3】一実施形態に係る各レイヤの候補演算のうち、最適なターゲット演算を選択する過程を説明するための図である。

【図4】一実施形態に係るニューラルネットワークのターゲットアーキテクチャーを決定するための探索過程を説明するための図である。

【図5】一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する方法の動作を説明するためのフローチャートである。

【図6】一実施形態に係る電子装置の構成を示す図である。 20

【発明を実施するための形態】

【0013】

実施形態に対する特定の構造的又は機能的な説明は、単なる例示のための目的として開示されたものであって、様々な形態に変更されることができる。したがって、実施形態は特定の開示形態に限定されるものではなく、本明細書の範囲は技術的な思想に含まれる変更、均等物ないし代替物を含む。

【0014】

第1又は第2などの用語を、複数の構成要素を説明するために用いることがあるが、このような用語は1つの構成要素を他の構成要素から区別する目的としてのみ解釈されなければならない。例えば、第1構成要素は第2構成要素と命名することができ、同様に、第2構成要素は第1構成要素にも命名することができる。 30

【0015】

いずれかの構成要素が他の構成要素に「連結」されているか「接続」されていると言及されたときには、その他の構成要素に直接的に連結されているか又は接続されているが、中間に他の構成要素が存在し得るものと理解されなければならない。

【0016】

単数の表現は、文脈上、明白に異なる意味をもたない限り複数の表現を含む。本明細書において、「含む」又は「有する」等の用語は、明細書上に記載した特徴、数字、ステップ、動作、構成要素、部品又はこれらを組み合わせたものが存在することを示すものであって、1つ又はそれ以上の他の特徴や数字、ステップ、動作、構成要素、部品、又はこれらを組み合わせたものなどの存在又は付加の可能性を予め排除しないものとして理解しなければならない。 40

【0017】

異なるように定義されない限り、技術的又は科学的な用語を含め、ここで用いる全ての用語は、本実施形態が属する技術分野で通常の知識を有する者によって一般的に理解されるものと同じ意味を有する。一般的に用いられる予め定義された用語は、関連技術の文脈上で有する意味と一致する意味を有するものと解釈されなければならないが、本明細書で明白に定義しない限り、理想的又は過度に形式的な意味として解釈されることはない。

【0018】

以下、添付する図面を参照しながら実施形態を詳細に説明する。図面を参照して説明す 50

る際に、図面符号に拘わらず同じ構成要素は同じ参照符号を付与し、これに対する重複する説明は省略する。

【0019】

図1は、一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する探索フレームワークを説明するための図である。

【0020】

図1を参照すると、探索フレームワーク100は、機械学習を介して基本ニューラルネットワーク120に対する最適なアーキテクチャー（又は、ニューラルネットワーク構造）を探索するフレームワークである。基本ニューラルネットワーク120は、学習される以前のニューラルネットワーク（又は、学習されていないニューラルネットワーク）として、各レイヤの演算及びパラメータ（例えば、連結加重値）などが確定していないニューラルネットワークである。基本ニューラルネットワーク120は、複数のニューラルネットワークレイヤ（又は、簡単に「レイヤ」）を含んでもよい。基本ニューラルネットワーク120は、深層ニューラルネットワーク（`deep neural network`、`DNN`）、畳み込みニューラルネットワーク（`convolutional neural network`、`CNN`）、再帰的ニューラルネットワーク（`recurrent neural network`、`RNN`）、`RBM`（`restricted boltzmann machine`）、`DBN`（`deep belief network`）、`BRDNN`（`bidirectional recurrent deep neural network`）、深層Q-ネットワーク（`deep Q-networks`）又は、このうち2以上の組み合わせの1つであってもよいが、前述した例に限定されない。基本ニューラルネットワーク120は、ハードウェア構造及び/又はソフトウェア構造を含んでもよい。

【0021】

探索フレームワーク100は、データベース110に格納された学習データを用いて基本ニューラルネットワーク120に対して機械学習を行う。機械学習は、教師あり学習（`supervised learning`）又は部分的教師あり学習（`partial supervised learning`）方式により実行され得る。

【0022】

一実施形態において、探索フレームワーク100は、教師あり学習を介して基本ニューラルネットワーク120を学習させることができる。探索フレームワーク100は、確率的勾配下降方式（`stochastic gradient descent`）のような調整アルゴリズム及び損失関数を用いて学習を行ってもよい。学習に使用される学習データは、ニューラルネットワークに入力される入力データと、該当入力データに対応する検証データを含んでもよい。基本ニューラルネットワーク120は、学習データに含まれる入力データを処理して結果データを出力することができる。探索フレームワーク100は、基本ニューラルネットワーク120から出力された結果データと検証データとの間の比較結果に基づいてニューラルネットワーク損失を決定し、ニューラルネットワーク損失を最小化する最適なアーキテクチャーを探索することができる。

【0023】

探索フレームワーク100は、多重目的（`multiple objective`）のニューラルアーキテクチャー探索（`NAS`）方式を行って、ターゲットニューラルネットワーク130のための最適なアーキテクチャーを探索することができる。探索フレームワーク100は、基本ニューラルネットワーク120のアーキテクチャーをサンプリングするのではなく、基本ニューラルネットワーク120の各レイヤごとに様々な候補演算を設定し、この候補演算のうち最も適切な候補演算であるターゲット演算をレイヤごとに選択する方式に基づいて最適なアーキテクチャーを探索することができる。このような探索方法によって探索フレームワーク100は、早期に最適化を行うことができる。

【0024】

探索フレームワーク100は、学習過程を通じた目的（例えば、オブジェクト分類、オ

10

20

30

40

50

プロジェクト認識、音声認識など)による最適なアーキテクチャーを有するターゲットニューラルネットワーク130を導き出すことができる。最適なアーキテクチャーを探索することは、ニューラルネットワークの各レイヤで実行される演算を決定し、ニューラルネットワークパラメータの最適値を決定することを含む。探索フレームワーク100は、本明細書で説明されるニューラルネットワークの最適なアーキテクチャーを探索する装置(例えば、図2に示す探索装置200)によって実行され得る。

【0025】

探索フレームワーク100は、ターゲットニューラルネットワーク130のための最適なアーキテクチャーを探索することにおいて、ハードウェアリソース制限事項(hardware resource constraint)を考慮する。探索フレームワーク100は、ニューラルネットワークが行うタスク(task)に対する検証損失だけでなく、ニューラルネットワークが実行されるとき用いられるハードウェアリソースを考慮して最適化することができる。探索フレームワーク100は、ニューラルネットワークが動作時に必要とするハードウェアリソースを考慮して、ターゲットニューラルネットワーク130を探索することができる。ハードウェアリソースは、例えば、電力消費、メモリ要求量、演算の数(number of operations)(例えば、multiply-accumulate(MAC)演算の数)、処理時間、及びGPUの占有率などであり得る。探索フレームワーク100は、1つ又は2以上のハードウェアリソースを考慮し、上記の例のハードウェアリソース以外にも数値で観測され得るハードウェアリソースであればいずれも制限されることなく考慮できる。

【0026】

探索フレームワーク100は、基本ニューラルネットワーク120に対する候補アーキテクチャーが決定されれば、候補アーキテクチャーに対するニューラルネットワーク損失とハードウェアリソースに対するハードウェアリソース損失を決定し、ニューラルネットワーク損失とハードウェアリソース損失を最小化するターゲットアーキテクチャーを探索することができる。ニューラルネットワーク損失とハードウェアリソース損失は、ターゲットアーキテクチャーを決定するための最適化損失を構成する。

【0027】

探索フレームワーク100は、ハードウェアリソース損失を決定するとき、候補アーキテクチャーのニューラルネットワークが要求するハードウェアリソースに対する実際の測定値と、ハードウェアリソース予測モジュール(例えば、図4に示すハードウェアリソース予測モジュール420)を用いて導き出されたハードウェアリソースに対する予測値に基づいて、ハードウェアリソース損失を決定することができる。下記でより詳しく説明されるが、ハードウェアリソース予測モジュールは、現在の学習の対象である候補アーキテクチャーのニューラルネットワークに対するハードウェアリソースを予測した予測値を提供するモジュールである。ハードウェアリソース予測モジュールは、入力された候補アーキテクチャーのパラメータに基づいて、該当候補アーキテクチャーのニューラルネットワークが必要とするハードウェアリソースに対する予測値を出力するよう学習された、ニューラルネットワークによって実現されることができる。ハードウェアリソース予測モジュールは、微分可能な特性を有し、ハードウェアリソース予測モジュールを介して探索過程における微分可能性が保持され得る、微分可能性が保持されることにより、エンドツーエンド学習(end-to-end learning)が可能になる。探索フレームワーク100は、ハードウェアリソース予測モジュールを介してニューラルネットワークのアーキテクチャーに対するハードウェアリソースを最適化損失に反映する。

【0028】

上記のように、探索フレームワーク100は、ハードウェアリソースの制限事項を考慮して、ニューラルネットワークの最適なアーキテクチャーを探索することができ、短い時間で最適化を行うことができる。また、探索フレームワーク100は、実際の測定したハードウェアリソース測定値を考慮して最適なアーキテクチャーを探索することができる。

【0029】

10

20

30

40

50

図2は、一実施形態に係るニューラルネットワークの最適なアーキテクチャーを探索する探索装置の構成を示すブロック図である。

【0030】

図2を参照すると、探索装置200は、ニューラルネットワークに対する最適なアーキテクチャーを探索する装置であって、図1を参照して説明した探索フレームワーク100を行うことができる。探索装置200は、アーキテクチャー探索に関連して本明細書で記載されるか、又は、図示した1つ以上の動作を行ってもよい。探索装置200は、プロセッサ210及びメモリ220を含む。格納装置230は、アーキテクチャー探索のためのデータ（例えば、学習データ）を格納し、学習に使用されるニューラルネットワークを格納することができる。

10

【0031】

メモリ220は、探索装置200の構成要素（例えば、プロセッサ210）によって用いられる様々なデータを格納することができる。データは、例えば、ソフトウェア、及びこれに関する命令に対する入力データ又は出力データを含んでもよい。メモリ220は、揮発性メモリ及び不揮発性メモリのうちの1つ以上を含んでもよい。

【0032】

プロセッサ210は、探索装置200の全体的な動作を制御し、探索装置200の動作を行うための命令を実行する。プロセッサ210は、例えば、ソフトウェアを実行してプロセッサ210に接続された探索装置200の少なくとも1つの他の構成要素（例えば、ハードウェア、又は、ソフトウェア構成要素）を制御し、様々なデータ処理又は演算を行うことができる。

20

【0033】

一実施形態によれば、データ処理又は演算の少なくとも一部として、プロセッサ210は、命令又はデータをメモリ220に格納し、メモリ220に格納された命令又はデータを処理し、結果データをメモリ220に格納することができる。プロセッサ210は、メインプロセッサ（例えば、中央処理装置又はアプリケーションプロセッサ）、又はこれとは独立的又は共に運営可能な補助プロセッサ（例えば、グラフィック処理装置、神経網処理装置（NPU：neural processing unit）、イメージ信号処理部、センサハブプロセッサ、又は、コミュニケーションプロセッサ）を含んでもよい。

【0034】

プロセッサ210は、学習データを用いてニューラルネットワーク（例えば、図1に示す基本ニューラルネットワーク120）の候補アーキテクチャーに対する学習を行い、ニューラルネットワーク損失を決定することができる。学習される以前のニューラルネットワークは、それぞれ1つ以上の人工ニューロンを含む複数のレイヤを含み、各レイヤは、実行可能な候補演算が予め定義されていることができる。候補演算は、例えば、 3×3 カーネル（kernel）基盤の畳み込み演算、 5×5 カーネル基盤の畳み込み演算、プーリング演算などを含み得るが、これに限定されることはない。ニューラルネットワークの各レイヤごとに、各レイヤが有する候補演算のいずれか1つの候補演算を選択することで、候補アーキテクチャーが決定されてもよい。プロセッサ210は、ニューラルネットワークの候補アーキテクチャーに対するパラメータに基づいてニューラルネットワーク損失を決定することができる。プロセッサ210は、候補アーキテクチャーのニューラルネットワークが学習データを処理して導出した結果データと検証データとの間の差に基づいて、ニューラルネットワーク損失を決定することができる。ニューラルネットワーク損失は、予め定義された損失関数によって決定され得る。

30

40

【0035】

プロセッサ210は、候補アーキテクチャーのニューラルネットワークが動作するとき要求（又は、使用）されるハードウェアリソースを測定できる。測定されるハードウェアリソースは、例えば、候補アーキテクチャーのニューラルネットワークが動作するときの電力消費、メモリ要求量、演算の数、及び処理時間のうちの1つ又は2以上を含むことができるが、これに限定されることはない。プロセッサ210は、ハードウェアリソース

50

を測定してハードウェアリソース測定値を決定することができる。

【0036】

プロセッサ210は、ハードウェアリソース予測モジュール（例えば、図4に示すハードウェアリソース予測モジュール420）を用いて、候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測できる。ハードウェアリソース予測モジュールは、候補アーキテクチャーのパラメータ（例えば、各レイヤの候補演算のうち選択された候補演算に関する情報）を入力とし、入力されたパラメータに基づいて候補アーキテクチャーのニューラルネットワークのハードウェアリソースを予測し、ハードウェアリソース予測値を出力するニューラルネットワークであり得る。

【0037】

プロセッサ210は、測定されたハードウェアリソースと予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定できる。プロセッサ210は、測定されたハードウェアリソースと予測されたハードウェアリソースとの間の差と、予め定義された損失関数に基づいてハードウェアリソース損失を決定することができる。例えば、プロセッサ210は、候補アーキテクチャーのニューラルネットワークが動作するとき、実際に測定された処理時間測定値とハードウェアリソース予測モデルから出力された処理時間予想値との間の差を、損失関数に適用してハードウェアリソース損失を決定することができる。

【0038】

プロセッサ210は、ニューラルネットワーク損失とハードウェアリソース損失に基づいて、ニューラルネットワークのターゲットアーキテクチャーを決定できる。プロセッサ210は、ニューラルネットワーク損失とハードウェアリソース損失を減少させるターゲットアーキテクチャー及びターゲットパラメータを決定する。プロセッサ210は、ハードウェアリソース損失が最小になるよう候補アーキテクチャーのパラメータをアップデートすることができる。プロセッサ210は、ニューラルネットワーク損失とハードウェアリソース損失を含む最適化損失を決定し、候補アーキテクチャーのニューラルネットワークに含まれている各レイヤの候補演算のうち、最適化損失を最小にするターゲット演算を選択することで、ターゲットアーキテクチャーを決定することができる。プロセッサ210は、ニューラルネットワーク損失とハードウェアリソース損失を最も最小化する各レイヤのターゲット演算を選択し、ニューラルネットワークのパラメータをアップデートすることができる。プロセッサ210は、候補アーキテクチャー及び候補アーキテクチャーのパラメータによるニューラルネットワーク損失と、候補アーキテクチャーのパラメータによるハードウェアリソース損失の加重和を最適化損失として決定し、該当最適化損失を最小にするターゲットアーキテクチャーを決定することができる。

【0039】

以上で説明した探索装置200により実行される動作は、ウェアラブル装置、スマートフォンのようなモバイル装置だけでなく、埋め込みシステム（embedded system）で作動可能なニューラルネットワーク基盤のアルゴリズムに多様に応用できる。

【0040】

図3は、一実施形態に係る各レイヤの候補演算のうち最適なターゲット演算を選択する過程を説明するための図である。

【0041】

図3を参照すると、ステップS310において、学習される以前のニューラルネットワーク（例えば、図1に示す基本ニューラルネットワーク120）は、複数のレイヤ312、314、316、318を含み、各レイヤ312、314、316、318は、実行可能な候補演算が定義されている。図示した実施形態において、ニューラルネットワーク310は、各レイヤ312、314、316、318に対して3個の候補演算を有するものと定義されている。レイヤ間で実行される候補演算は、互いに異なる演算方式であってもよい。例えば、レイヤ312とレイヤ314との間で実行される候補演算は、互いに異なる方式の演算であってもよい。

10

20

30

40

50

【 0 0 4 2 】

探索装置（例えば、図 2 に示す探索装置 2 0 0）は、学習過程において、この候補演算のうち最適な候補演算であるターゲット演算を選択することができる。ステップ S 3 2 0 において、探索装置は、各レイヤ 3 1 2, 3 1 4, 3 1 6, 3 1 8 ごとに、候補演算のいずれか 1 つの候補演算 3 2 2, 3 2 4, 3 2 6, 3 2 8, 3 2 9 を選択し、選択された候補演算 3 2 2, 3 2 4, 3 2 6, 3 2 8, 3 2 9 の組み合わせで構成された候補アーキテクチャに対する最適化損失を決定することができる。探索装置は、探索空間で各レイヤ 3 1 2, 3 1 4, 3 1 6, 3 1 8 の候補演算を数回組み合わせ、それぞれの組み合わせにおける最適化損失を算出し、最適化損失を最小にする候補演算（言い換えれば、ターゲット演算）の組み合わせを決定することができる。様々な組み合わせに対する学習過程が完了すれば、各レイヤ 3 1 2, 3 1 4, 3 1 6, 3 1 8 ごとに選択されたターゲット演算に基づいて、ターゲットアーキテクチャを決定することができる。ターゲットアーキテクチャは、各レイヤのターゲット演算を含む。

10

【 0 0 4 3 】

図 4 は、一実施形態に係るニューラルネットワークのターゲットアーキテクチャを決定するための探索過程を説明するための図である。

【 0 0 4 4 】

図 4 を参照すると、ニューラルネットワークのターゲットアーキテクチャの探索過程において、ニューラルネットワークの候補アーキテクチャ a 4 1 2 が与えられれば、候補アーキテクチャ a 4 1 2 のパラメータ $w(a)_{416}$ が決定され得る。候補アーキテクチャ a 4 1 2 は、各レイヤごとに選択された候補演算の集合を含むニューラルネットワーク構造を示す。

20

【 0 0 4 5 】

候補アーキテクチャ a 4 1 2 のパラメータは、ニューラルネットワークの各レイヤの候補演算のうちから選択された候補演算に対するパラメータと、選択された各候補演算に対する演算特性を示すパラメータを含んでもよい。例えば、ニューラルネットワークに含まれるいずれかのレイヤの候補演算が、 3×3 カーネル基盤の畳み込み演算及び 5×5 カーネル基盤の畳み込み演算があると仮定すれば、候補アーキテクチャ a 4 1 2 のパラメータは、2 つの畳み込み演算のいずれかの畳み込み演算を選択したかを示すパラメータと、選択された畳み込み演算のカーネルパラメータ (kernel parameter) などを含んでもよい。ここで、畳み込み演算は、畳み込みレイヤとして実現されてもよい。

30

【 0 0 4 6 】

探索装置（例えば、図 2 に示す探索装置 2 0 0）は、パラメータ $w(a)_{416}$ に基づいて、候補アーキテクチャ a 4 1 2 のニューラルネットワークの特定タスクに対するニューラルネットワーク損失 $L_{TASK}(w(a), a)_{418}$ を決定することができる。ニューラルネットワーク損失 $L_{TASK}(w(a), a)_{418}$ は、ニューラルネットワークが行うタスクの損失を最小化するための損失である。

【 0 0 4 7 】

探索装置は、候補アーキテクチャ a 4 1 2 を有するニューラルネットワークの演算を実行し (4 3 2)、該当演算の実行過程で要求される全体ニューラルネットワークのハードウェアリソースを実際に測定する (4 3 4)。測定されるハードウェアリソースは、例えば、電力消費、メモリ要求量、演算の数、及び処理時間などを含んでもよい。ハードウェアリソースの測定 4 3 4 を介して候補アーキテクチャ a 4 1 2 に対するハードウェアリソース測定値 $L_{ATM}(a)_{436}$ が決定される。実施形態により、マックス演算を介して決定されたアーキテクチャに対するニューラルネットワーク演算を行うことで、ハードウェアリソースを測定することができる。このような過程を含む過程 4 3 0 は、微分可能性が成立されないため、該当過程 4 3 0 に対しては、フォワード演算及びバックワード演算を定義することができない。これを解決するために、ハードウェアリソース予測モジュール 4 2 0 を使用してもよい。

40

50

【 0 0 4 8 】

探索装置は、ハードウェアリソース予測モジュール 4 2 0 を用いて、候補アーキテクチャ a 4 1 2 を有するニューラルネットワークがこの演算を行うときに要求するものと予想される、全体ニューラルネットワークのハードウェアリソースを予測することができる。ハードウェアリソース予測モジュール 4 2 0 を介して候補アーキテクチャ a 4 1 2 に対するハードウェアリソース予測値 $LAT_P(a)$ 4 2 2 が決定され得る。ハードウェアリソース予測モジュール 4 2 0 は、候補アーキテクチャ a 4 1 2 のパラメータを入力とし、入力されたパラメータに基づいて候補アーキテクチャ a 4 1 2 のニューラルネットワークのハードウェアリソースを予測し、ハードウェアリソース予測値 $LAT_P(a)$ 4 2 2 を出力することができる。ハードウェアリソース予測モジュール 4 2 0 における演算は、微分可能な演算からなる。

【 0 0 4 9 】

ハードウェアリソース予測モジュール 4 2 0 は、ニューラルネットワークのアーキテクチャのパラメータに基づいて該当アーキテクチャが要求するか、使用するものと予想されるハードウェアリソースの予測値を出力するよう、学習過程を介して学習されたニューラルネットワークであってもよい。但し、ハードウェアリソース予測モジュール 4 2 0 は、ニューラルネットワーク以外に、候補アーキテクチャ a 4 1 2 に基づいてニューラルネットワークのハードウェアリソースを予測できる他の手段により実現されてもよい。

【 0 0 5 0 】

探索装置は、ハードウェアリソース測定値 $LAT_M(a)$ 4 3 6 とハードウェアリソース予測値 $LAT_P(a)$ 4 2 2 に基づいて、候補アーキテクチャ a 4 1 2 に対するリソース損失 $L_{HW}(a)$ 4 4 2 を決定することができる。ハードウェアリソース測定値 $LAT_M(a)$ 4 3 6 とハードウェアリソース予測値 $LAT_P(a)$ 4 2 2 との間の差が最小化されるよう、リソース損失 $L_{HW}(a)$ 4 4 2 が定義される。

【 0 0 5 1 】

一実施形態において、リソース損失 $L_{HW}(a)$ 4 4 2 は、ハードウェアリソース測定値 $LAT_M(a)$ 4 3 6 とハードウェアリソース予測値 $LAT_P(a)$ 4 2 2 との間の差による損失を示す $L_{HW1}(a)$ と、ハードウェアリソースの最適化のための要素（例えば、レイテンシを最小化するための要素）の $L_{HW2}(a)$ に基づいて決定される。 $L_{HW1}(a)$ と $L_{HW2}(a)$ はそれぞれ次の数式 (1) 及び数式 (2) によって決定され得る。

【 数 1 】

$$L_{HW1}(a) = (LAT_M(a) - LAT_P(a))^2 \quad (1)$$

【 数 2 】

$$L_{HW2}(a) = (LAT_M(a))^2 \quad (2)$$

リソース損失 $L_{HW}(a)$ 4 4 2 は、例えば、次の数式 (3) のように $L_{HW1}(a)$ と $L_{HW2}(a)$ との間の加重和として決定されてもよい。

【 数 3 】

$$L_{HW}(a) = L_{HW1}(a) + w \times L_{HW2}(a) \quad (3)$$

ここで、 w は、 $L_{HW2}(a)$ に適用される加重値として、例えば、予め設定された定数である。実施形態により、 $L_{HW1}(a)$ のみに加重値が適用されてもよく、 $L_{HW1}(a)$ と $L_{HW2}(a)$ にそれぞれ互いに異なる加重値が適用されてもよい。

【 0 0 5 2 】

探索装置は、ハードウェアリソース測定値 $LAT_M(a)$ 4 3 6 とハードウェアリソー

ス予測値 $L_{ATP}(a)$ 422 との間の差を予め定義された損失関数に適用し、ハードウェアリソース損失 $L_{HW}(a)$ 442 を決定することができる。

【0053】

探索装置は、ニューラルネットワーク損失 $L_{TASK}(w(a), a)$ 418 とハードウェアリソース損失 $L_{HW}(a)$ 442 を含む最適化損失を決定し、候補アーキテクチャ a 412 のニューラルネットワークに含まれている各レイヤの候補演算のうち、最適化損失を最小にするターゲット演算を選択することで、ターゲットアーキテクチャを決定することができる。例えば、ターゲットアーキテクチャは、次の数式(4)のように最適化損失 $L_{TASK}(w(a), a) + \lambda \cdot L_{HW}(a)$ を最小にする候補アーキテクチャ a 412 と候補アーキテクチャ a のパラメータ $w(a)$ 416 を探索することによって決定されてもよい。数式(4)の例として、最適化損失は、ハードウェアリソース損失 $L_{HW}(a)$ 442 に加重値 λ が適用されたニューラルネットワーク損失 $L_{TASK}(w(a), a)$ 418 とハードウェアリソース損失 $L_{HW}(a)$ 442 との間の加重和として決定され得る。

【数4】

$$\min_a \min_w L_{TASK}(w(a), a) + \lambda \cdot L_{HW}(a) \quad (4)$$

探索装置は、少ない探索時間でニューラルネットワークのターゲットアーキテクチャを探索し、ターゲットアーキテクチャを選定することにおいて、ニューラルネットワークが必要とするハードウェアリソースを最適化制限事項として考慮する。図4において、過程410は、微分可能な特性を有し、過程430は、微分不可能な(non-differentiable)特性を有する。探索装置は、ニューラルネットワークの各レイヤの演算に基づいた全体ニューラルネットワークのハードウェアリソースをニューラルネットワークとして実現できるハードウェアリソース予測モジュール420を用いて予測することで、微分可能性を保持することができる。微分可能性を保持することによりエンドツーエンド学習が可能になる。また、上述した探索過程は、全体ニューラルネットワークのハードウェアリソースの制限事項を反映して、ニューラルネットワークのアーキテクチャを微分可能に最適化することができ、一回の学習過程を介してターゲットアーキテクチャを探し得るため、最適化時間が短いという長所を有する。そして、ターゲットアーキテクチャの探索過程で生じるフォワード演算とバックワード演算との間の一貫性も保持され得る。

【0054】

図5は、一実施形態に係るニューラルネットワークの最適なアーキテクチャを探索する方法の動作を説明するためのフローチャートである。該当方法の動作は、図2に示す探索装置200によって実行され得る。

【0055】

図5を参照すると、動作510において、探索装置は、ニューラルネットワーク(例えば、図1に示す基本ニューラルネットワーク120)の候補アーキテクチャを選択する。探索装置は、ニューラルネットワークの各レイヤに対して定義された候補演算のいずれか1つの候補演算を選択することで、候補アーキテクチャを選択することができる。

【0056】

動作520において、探索装置は、ニューラルネットワークの候補アーキテクチャに対するパラメータに基づいて、ニューラルネットワーク損失を決定する。探索装置は、学習データを用いてニューラルネットワークの候補アーキテクチャに対する学習を行い、ニューラルネットワーク損失を決定する。探索装置は、候補アーキテクチャのニューラルネットワークが学習データを処理して導出した結果データと検証データとの間の差に基づいて、ニューラルネットワーク損失を決定することができる。候補アーキテクチャのニューラルネットワークから導き出された結果データと目的とする検証データとの間の差が大きくなるほど、ニューラルネットワーク損失は大きくなる。

10

20

30

40

50

【 0 0 5 7 】

動作 5 3 0 において、探索装置は、候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを測定する。測定されるハードウェアリソースは、例えば、候補アーキテクチャーのニューラルネットワークが動作するときの電力消費、メモリ要求量、演算の数、及び処理時間のうちの 1 つ又は 2 以上を含むが、これに限定されることはない。

【 0 0 5 8 】

動作 5 4 0 において、探索装置は、ハードウェアリソース予測モジュール（例えば、図 4 に示すハードウェアリソース予測モジュール 4 2 0 ）を用いて候補アーキテクチャーのニューラルネットワークが動作するときに要求されるハードウェアリソースを予測する。ハードウェアリソース予測モデルには、ニューラルネットワークの各レイヤで実行可能な候補演算のうち、候補アーキテクチャーを構成する選択された候補演算に関する情報が入力され、ハードウェアリソース予測モデルは、入力された情報に基づいて該当ニューラルネットワークを必要とするハードウェアリソースの予測値を提供することができる。

10

【 0 0 5 9 】

動作 5 5 0 において、探索装置は、測定されたハードウェアリソースと予測されたハードウェアリソースに基づいて、ハードウェアリソース損失を決定する。探索装置は、測定されたハードウェアリソースと予測されたハードウェアリソースとの差異と、予め定義された損失関数に基づいてハードウェアリソース損失を決定することができる。

【 0 0 6 0 】

動作 5 6 0 において、探索装置は、ニューラルネットワーク損失とハードウェアリソース損失に基づいて、ニューラルネットワークのターゲットアーキテクチャーとターゲットパラメータを決定する。探索装置は、ニューラルネットワーク損失とハードウェアリソース損失を含む最適化損失を最小にするターゲットアーキテクチャーとターゲットパラメータを決定し得る。探索装置は、候補アーキテクチャーのニューラルネットワークに含まれている各レイヤの候補演算のうち、最適化損失を最小にするターゲット演算を選択することでターゲットアーキテクチャーを決定できる。探索装置は、候補アーキテクチャー及び候補アーキテクチャーのパラメータによるニューラルネットワーク損失と、候補アーキテクチャーのパラメータによるハードウェアリソース損失の加重和を最適化損失として決定し、該当最適化損失を最小にするターゲットアーキテクチャーを決定する。

20

30

【 0 0 6 1 】

図 6 は、一実施形態に係る電子装置の構成を示す図である。

【 0 0 6 2 】

図 6 を参照すると、電子装置 6 0 0 は、様々な形態の電子装置であり得る。例えば、電子装置 6 0 0 は、ウェアラブル装置（例えば、AR グラスのような拡張現実提供装置、HMD (head mounted display)) スマートフォン、タブレットコンピュータ、ネットブック、ラップトップ、製品検査装置、パーソナルコンピュータ、又はサーバであってもよいが、これに制限されることはない。

【 0 0 6 3 】

電子装置 6 0 0 は、プロセッサ 6 1 0、メモリ 6 2 0、カメラ 6 3 0、センサ 6 4 0、出力装置 6 5 0、及び通信装置 6 6 0 を含む。電子装置 6 0 0 の各構成要素のうち少なくとも一部は、周辺機器間の通信インターフェース 6 7 0（例えば、バス、GPIO (general purpose input and output) インターフェース、SPI (serial peripheral interface)、又は、MIPI (mobile industry processor interface)) を介して接続されて信号（例えば、命令又はデータ）を相互間に交換することができる。

40

【 0 0 6 4 】

プロセッサ 6 1 0 は、電子装置 6 0 0 の全体的な動作を制御し、電子装置 6 0 0 内で実行するための機能及び命令を実行する。プロセッサ 6 1 0 は、図 1 ~ 図 5 を参照して前述した探索装置（例えば、図 2 の探索装置 2 0 0 ）の動作を行う。

50

【 0 0 6 5 】

メモリ 6 2 0 は、プロセッサ 6 1 0 によって実行可能な命令と入力 / 出力されるデータを格納する。メモリ 6 2 0 は、RAM、DRAM、SRAM のような揮発性メモリ及び / 又は ROM、フラッシュメモリ のようなこの技術分野で知られた不揮発性メモリを含んでもよい。

【 0 0 6 6 】

カメラ 6 3 0 は、映像を撮影する。カメラ 6 3 0 は、例えば、カラー映像、白黒映像、グレイ映像、赤外線映像、又は深度映像などを取得してもよい。

【 0 0 6 7 】

センサ 6 4 0 は、電子装置 6 0 0 の作動状態（例えば、電力又は温度）、又は、外部の環境状態（例えば、ユーザ状態）を検出し、検出された状態に対応する電気信号又はデータ値を生成することができる。センサ 6 4 0 は、例えば、ジェスチャーセンサ、ジャイロセンサ、気圧センサ、マグネチックセンサ、加速度センサ、グリップセンサ、近接センサ、カラーセンサ、IR (infrared) センサ、生体センサ、温度センサ、湿度センサ、又は、照度センサを含んでもよい。

10

【 0 0 6 8 】

出力装置 6 5 0 は、視覚的、聴覚的、又は触覚的なチャンネルを通じてユーザに電子装置 6 0 0 の出力を提供する。出力装置 6 5 0 は、例えば、液晶ディスプレイや LED / OLED ディスプレイ、マイクロ LED (micro light emitting diode、microLED)、タッチスクリーン、スピーカー、振動発生装置又はユーザ

20

【 0 0 6 9 】

通信装置 6 6 0 は、電子装置 6 0 0 と外部装置との間の直接（例えば、有線）通信チャンネル又は無線通信チャンネルの確立、及び確立された通信チャンネルを通じた通信実行を支援する。一実施形態によれば、通信装置 6 6 0 は、無線通信モジュール（例えば、セルラー通信モジュール、近距離無線通信モジュール、又は GNSS (global navigation satellite system) 通信モジュール（例えば、LAN (local area network) 通信モジュール、又は電力線通信モジュール）を含むことができる。無線通信モジュールは、近距離通信ネットワーク（例えば、ブルートゥース（登録商標）、WiFi (wireless fidelity) direct 又は IrDA (infrared data association)）又は遠距離通信ネットワーク（例えば、レガシーセルラーネットワーク、5G ネットワーク、次世代通信ネットワーク、インターネット、又は、コンピュータネットワーク（例えば、LAN 又は WAN））を通じて外部の装置と通信できる。

30

【 0 0 7 0 】

一実施形態において、電子装置 6 0 0 は、ニューラルネットワーク基盤のアルゴリズムを使用する拡張現実提供装置（例えば、AR (augmented reality) グラス）であってもよい。拡張現実提供装置は、ユーザの顔面に着用され、ユーザに拡張現実サービス及び / 又は仮想現実サービスに関するコンテンツを提供する。プロセッサ 6 1 0 は、ターゲットアーキテクチャーを有するニューラルネットワークを用いて処理動作を行う。カメラ 6 3 0 は、拡張現実コンテンツの生成のための映像を撮影し、プロセッサ 6 1 0 は、ターゲットアーキテクチャーを有するニューラルネットワークを用いて映像を処理することで、拡張現実コンテンツを生成できる。例えば、プロセッサ 6 1 0 は、カメラ 6 3 0 を介して取得された映像から特定のオブジェクトを認識し、認識したオブジェクト領域又はオブジェクト周辺領域に仮想のコンテンツを重複して表現することで、拡張現実コンテンツを生成することができる。

40

【 0 0 7 1 】

プロセッサ 6 1 0 は、図 2 及び図 5 を参照して説明したような過程を介してニューラルネットワークのターゲットアーキテクチャーを決定することができる。例えば、プロセッサ 6 1 0 は、ニューラルネットワーク（例えば、図 1 に示す基本ニューラルネットワーク

50

120)の候補アーキテクチャに対するパラメータに基づいてニューラルネットワーク損失を決定し、候補アーキテクチャのニューラルネットワークが動作するとき要求されるハードウェアリソースを測定することができる。プロセッサ610は、ハードウェアリソース予測モジュール(例えば、図4に示すハードウェアリソース予測モジュール420)を用いて候補アーキテクチャのニューラルネットワークが動作するとき要求されるハードウェアリソースを予測し、測定されたハードウェアリソースと予測されたハードウェアリソースに基づいてハードウェアリソース損失を決定することができる。プロセッサ610は、ニューラルネットワーク損失とハードウェアリソース損失に基づいてターゲットアーキテクチャを決定する。プロセッサ610は、候補アーキテクチャ及び候補アーキテクチャのパラメータによるニューラルネットワーク損失と、候補アーキテクチャのパラメータによるハードウェアリソース損失の加重和に基づいて最適化損失を決定し、最適化損失を最小にするターゲットアーキテクチャと候補アーキテクチャのパラメータを決定することができる。

10

【0072】

以上で説明された実施形態は、ハードウェア構成要素、ソフトウェア構成要素、又はハードウェア構成要素及びソフトウェア構成要素の組み合わせで具現される。例えば、本実施形態で説明した装置及び構成要素は、例えば、プロセッサ、コントローラ、ALU(arithmetic logic unit)、デジタル信号プロセッサ(digital signal processor)、マイクロコンピュータ、FPA(field programmable array)、PLU(programmable logic unit)、マイクロプロセッサ、又は命令(instruction)を実行して応答する異なる装置のように、1つ以上の汎用コンピュータ又は特殊目的コンピュータを用いて具現される。処理装置は、オペレーティングシステム(OS)及びオペレーティングシステム上で実行される1つ以上のソフトウェアアプリケーションを実行する。また、処理装置は、ソフトウェアの実行に回答してデータをアクセス、格納、操作、処理、及び生成する。理解の便宜のために、処理装置は1つが使用されるものとして説明する場合もあるが、当技術分野で通常の知識を有する者は、処理装置が複数の処理要素(processing element)及び/又は複数種類の処理要素を含むことが把握する。例えば、処理装置は、複数のプロセッサ又は1つのプロセッサ及び1つのコントローラを含む。また、並列プロセッサ(parallel processor)のような、他の処理構成も可能である。

20

30

【0073】

ソフトウェアは、コンピュータプログラム、コード、命令、又はそのうちの1つ以上の組み合わせを含み、希望の通りに動作するよう処理装置を構成するか、独立的又は結合的に処理装置を命令することができる。ソフトウェア及び/又はデータは、処理装置によって解釈されるか、処理装置に命令又はデータを提供するために、いずれかの種類の機械、構成要素、物理的装置、仮想装置、コンピュータ記録媒体又は装置、又は送信される信号波に永久的又は一時的に具体化することができる。ソフトウェアはネットワークに連結されたコンピュータシステム上に分散され、分散した方法で格納されるか実行され得る。ソフトウェア及びデータは1つ以上のコンピュータ読み取り可能な記録媒体に格納され得る。

40

【0074】

本実施形態による方法は、様々なコンピュータ手段を介して実施されるプログラム命令の形態で具現され、コンピュータ読み取り可能な記録媒体に記録される。記録媒体は、プログラム命令、データファイル、データ構造などを単独又は組み合わせで含む。記録媒体及びプログラム命令は、本発明の目的のために特別に設計して構成されたものでもよく、コンピュータソフトウェア分野の技術を有する当業者にとって公知のものであり使用可能なものであってもよい。コンピュータ読み取り可能な記録媒体の例として、ハードディスク、フロッピー(登録商標)ディスク及び磁気テープのような磁気媒体、CD-ROM、DVDのような光記録媒体、プロプティカルディスクのような磁気-光媒体、及びROM

50

、RAM、フラッシュメモリなどのようなプログラム命令を保存して実行するように特別に構成されたハードウェア装置を含む。プログラム命令の例としては、コンパイラによって生成されるような機械語コードだけでなく、インタプリタなどを用いてコンピュータによって実行される高級言語コードを含む。

【0075】

上記で説明したハードウェア装置は、本発明に示す動作を実行するために1つ以上のソフトウェアモジュールとして作動するように構成してもよく、その逆も同様である。

【0076】

上述したように実施形態を例として限定された図面によって説明したが、当技術分野で通常の知識を有する者であれば、上記の説明に基づいて様々な技術的な修正及び変形を適用することができる。例えば、説明された技術が説明された方法と異なる順に実行され、及び/又は説明されたシステム、構造、装置、回路などの構成要素が説明された方法とは異なる形態に結合又は組み合わせられてもよく、他の構成要素又は均等物によって置き換え又は代替されたとしても適切な結果を達成することができる。

【0077】

したがって、他の具現、他の実施形態及び特許請求の範囲と均等なものも後述する特許請求の範囲に属する。

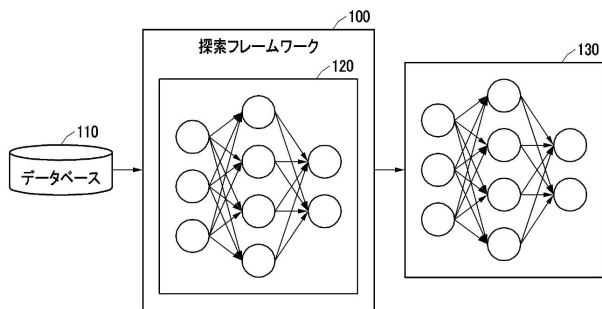
【符号の説明】

【0078】

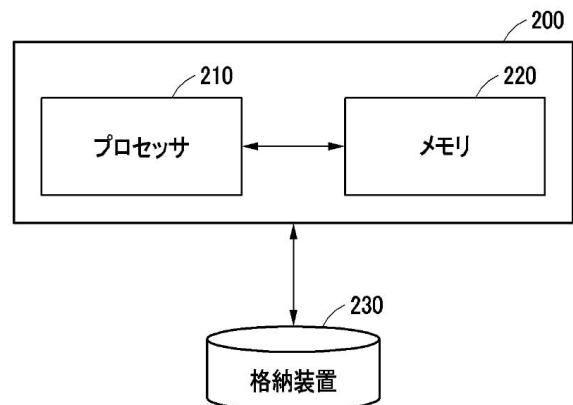
- 100：探索フレームワーク
- 110：データベース
- 120：基本ニューラルネットワーク
- 130：ターゲットニューラルネットワーク
- 200：探索装置
- 210：プロセッサ
- 220：メモリ
- 230：格納装置
- 600：電子装置
- 610：プロセッサ
- 620：メモリ
- 630：カメラ
- 640：センサ
- 650：出力装置
- 660：通信装置
- 670：通信インターフェース

【図面】

【図1】



【図2】



10

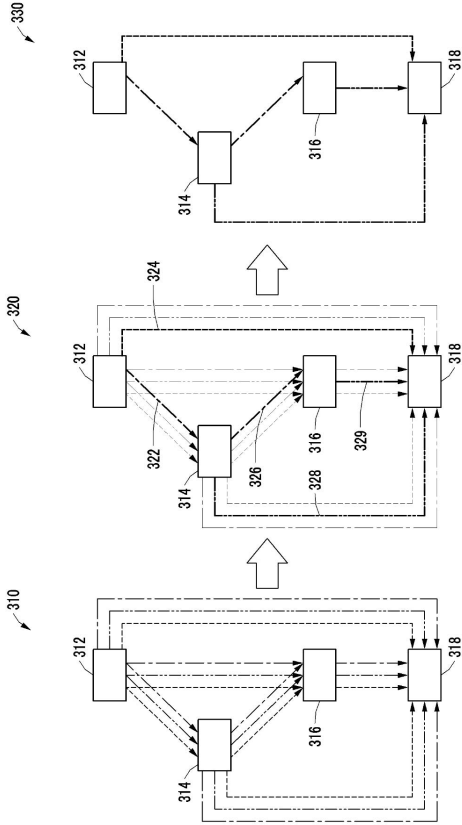
20

30

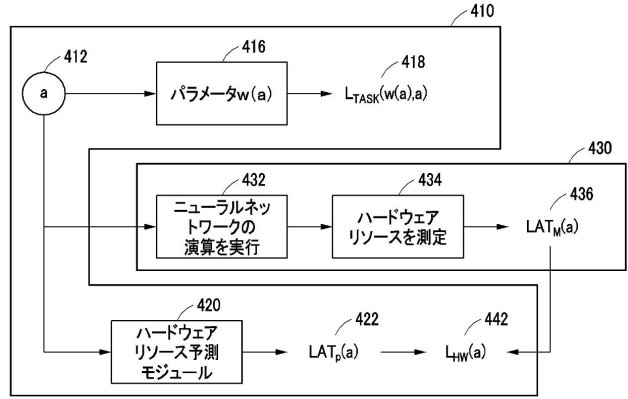
40

50

【 図 3 】



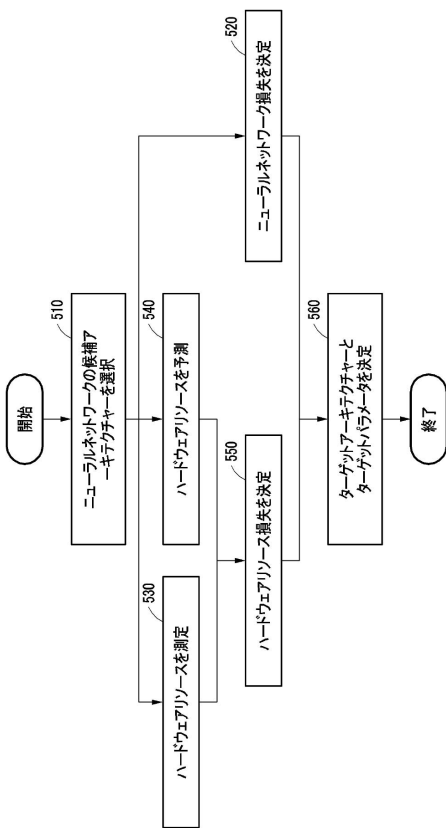
【 図 4 】



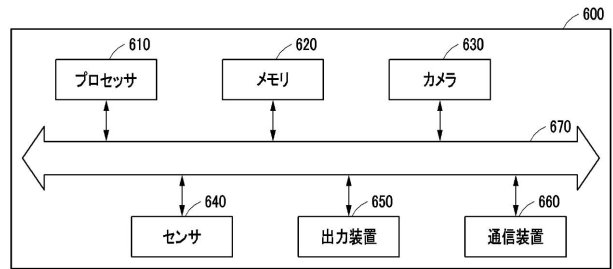
10

20

【 図 5 】



【 図 6 】



30

40

50

フロントページの続き

弁理士 宮崎 修

(72)発明者 李 元熙

大韓民国京畿道水原市靈通区三星路 1 3 0 三星綜合技術院内