

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7549153号
(P7549153)

(45)発行日 令和6年9月10日(2024.9.10)

(24)登録日 令和6年9月2日(2024.9.2)

(51)国際特許分類		F I	
G 0 6 N	3/10 (2006.01)	G 0 6 N	3/10
G 0 6 F	8/40 (2018.01)	G 0 6 F	8/40
G 0 6 F	8/30 (2018.01)	G 0 6 F	8/30

請求項の数 35 (全31頁)

(21)出願番号	特願2023-549884(P2023-549884)	(73)特許権者	514322098
(86)(22)出願日	令和4年11月2日(2022.11.2)		ベイジン バイドゥ ネットコム サイエ ンス テクノロジー カンパニー リミテ ッド
(65)公表番号	特表2024-527451(P2024-527451 A)		Beijing Baidu Netco m Science Technolog y Co., Ltd.
(43)公表日	令和6年7月25日(2024.7.25)		中華人民共和国 ベキン 100085, ハイディアン ディストリクト, シャン ディ テンス ストリート, 10番, バ イドゥ キャンパス 2階
(86)国際出願番号	PCT/CN2022/129228		2/F Baidu Campus, N o.10, Shangdi 10th Street, Haidian Dis trict, Beijing 1000
(87)国際公開番号	WO2023/221408		最終頁に続く
(87)国際公開日	令和5年11月23日(2023.11.23)		
審査請求日	令和5年8月23日(2023.8.23)		
(31)優先権主張番号	202210560831.9		
(32)優先日	令和4年5月19日(2022.5.19)		
(33)優先権主張国・地域又は機関	中国(CN)		

(54)【発明の名称】 深層学習フレームワークの演算子処理方法及び装置、電子機器、記憶媒体ならびにコンピュータプログラム

(57)【特許請求の範囲】

【請求項1】

深層学習フレームワークの演算子処理方法であって、
演算子カーネル関数、及び深層学習フレームワークに依存しないテンプレートパラメータを含む処理対象演算子を取得することと、
前記処理対象演算子に対する入力情報を受信したことに応答して、前記入力情報を用いて前記テンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得ることと、
複数の完全なテンプレートパラメータに基づいて、前記演算子カーネル関数を処理して、前記深層学習フレームワーク用の利用可能な演算子を得ることと、を含む
深層学習フレームワークの演算子処理方法。

【請求項2】

前記処理対象演算子に対する入力情報を受信したことに応答して、前記入力情報を用いて前記テンプレートパラメータを解析し、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得ることは、
マクロを用いて可変パラメータリストを作成することと、
前記可変パラメータリストに基づいて前記テンプレートパラメータを反復に解析することと、を含み、
ここで、前記入力情報を前記マクロのマクロパラメータとする
請求項1に記載の方法。

10

20

【請求項 3】

前記可変パラメータリストは、
 処理対象演算子の名称と、
 前記演算子カーネル関数の名称と、
 前記演算子カーネル関数のデータレイアウトと、
 前記深層学習フレームワークを実装するための機器タイプと、
 前記深層学習フレームワークを実装するためのデータタイプと、を含む
 請求項 2 に記載の方法。

【請求項 4】

前記可変パラメータリストに基づいて前記テンプレートパラメータを反復に解析することは、
 前記マクロに含まれるパラメータ数を特定することと、
 前記パラメータ数に応じて、反復終了位置を特定することと、
 可変パラメータリストを用いて、前記反復終了位置に応じて前記入力情報を反復に処理することと、を含む
 請求項 2 に記載の方法。

10

【請求項 5】

前記可変パラメータリストに基づいて前記テンプレートパラメータを反復に解析することは、
 前記可変パラメータリストを用いて前記テンプレートパラメータを反復に解析して、複数の
 内部テンプレートパラメータを得ることと、
 前記複数の内部テンプレートパラメータをスライシングして、前記複数の完全なテンプレ
 ートパラメータを得ることと、をさらに含む
 請求項 2 に記載の方法。

20

【請求項 6】

前記複数の内部テンプレートパラメータは複数の文字列を含み、前記複数の内部テンプレ
 ートパラメータをスライシングして、前記複数の完全なテンプレートパラメータを得る
 ことは、
 前記入力情報に基づいて、前記複数の文字列をスライシングして、前記処理対象演算子
 に対する複数の完全なテンプレートパラメータを得ることを含む
 請求項 5 に記載の方法。

30

【請求項 7】

複数の完全なテンプレートパラメータに基づいて、前記演算子カーネル関数を処理して、
 前記深層学習フレームワーク用の利用可能な演算子を得ることは、
 前記複数の完全なテンプレートパラメータを明示的にインスタンス化することで、前記演
 算子カーネル関数を処理して、前記深層学習フレームワーク用の利用可能な演算子を得る
 ことを含む
 請求項 5 に記載の方法。

【請求項 8】

前記複数の完全なテンプレートパラメータの各々は、テンプレート関数記述情報及びテン
 プレート関数情報を含み、前記複数の完全なテンプレートパラメータを明示的にインスタ
 ンス化することで、前記演算子カーネル関数を処理して、前記深層学習フレームワーク用
 の利用可能な演算子を得ることは、
 前記テンプレート関数記述情報をヘッダファイルとしてインスタンス化することと、
 前記テンプレート関数情報をソースコードファイルとしてインスタンス化することと、を
 含む
 請求項 7 に記載の方法。

40

【請求項 9】

前記複数の完全なテンプレートパラメータに基づいて、前記演算子カーネル関数を処理し
 て、前記深層学習フレームワーク用の利用可能な演算子を得ることは、

50

前記複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、該完全なテンプレートパラメータの関数パラメータタイプを特定することと、前記関数パラメータタイプに基づいて、前記完全なテンプレートパラメータの入力情報を対応するパラメータ情報に変換することと、前記対応するパラメータ情報を記録することと、を含む請求項 5 に記載の方法。

【請求項 10】

前記複数の完全なテンプレートパラメータに基づいて、前記深層学習フレームワークに基づく登録対象演算子を作成することと、前記登録対象演算子を前記深層学習フレームワーク内部のグローバル演算子テーブルに登録することと、をさらに含む請求項 5 に記載の方法。

10

【請求項 11】

前記登録対象演算子は登録対象演算子記述情報及び登録対象演算子カーネル関数を含み、前記登録対象演算子カーネル関数には正規化された入力情報及び正規化された関数ポインタが含まれ、前記複数の完全な内部テンプレート関数に基づいて、前記深層学習フレームワークに基づく登録対象演算子を作成することは、前記複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、前記完全なテンプレートパラメータに対応する、正規化された形式を有する静的関数を含む構造体を特定することと、前記静的関数の入力情報を前記正規化された入力情報とすることと、前記静的関数の関数ポインタを前記正規化された関数ポインタとすることと、を含む請求項 10 に記載の方法。

20

【請求項 12】

前記静的関数の入力情報は入力情報リストを構成し、前記入力情報リストは入力テンソルリストを含み、前記構造体は少なくとも 1 つの特化サブ構造体を含み、前記少なくとも 1 つの特化サブ構造体の各々是对応するデータタイプを有する請求項 11 に記載の方法。

【請求項 13】

前記グローバル演算子テーブルは、演算子カーネル工場クラスと、演算子カーネル名称クラス及びカーネルキー値クラスと、演算子カーネルクラスと、を含む請求項 10 に記載の方法。

30

【請求項 14】

前記演算子カーネルクラスは、演算子カーネル関数のポインタと、演算子カーネル関数入力パラメータの記述情報と、演算子カーネル関数出力パラメータの記述情報と、の少なくとも 1 つを含む請求項 13 に記載の方法。

40

【請求項 15】

前記深層学習フレームワークに依存しないテンプレートパラメータは、前記処理対象演算子を配布するために用いられ、前記深層学習フレームワークに依存しないテンプレートパラメータは、前記深層学習フレームワークの実装機器及び前記深層学習フレームワークのデータタイプのいずれにも依存しないテンプレートパラメータを含む請求項 1 に記載の方法。

【請求項 16】

前記利用可能な演算子に対する呼び出しコマンドにตอบสนองして、前記利用可能な演算子を呼び出すことと、呼び出された前記利用可能な演算子をコンパイルすることで、前記利用可能な演算子に対

50

応する関数機能を実行することと、をさらに含む
請求項 1 に記載の方法。

【請求項 17】

深層学習フレームワークの演算子処理装置であって、
演算子カーネル関数、及び深層学習フレームワークに依存しないテンプレートパラメータ
を含む処理対象演算子を取得するための取得モジュールと、
前記処理対象演算子に対する入力情報を受信したことに応答して、前記入力情報を用いて
前記テンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全
なテンプレートパラメータを得るための解析モジュールと、
複数の完全なテンプレートパラメータに基づいて、前記演算子カーネル関数を処理して、
前記深層学習フレームワーク用の利用可能な演算子を得るための処理モジュールと、を含む
深層学習フレームワークの演算子処理装置。

10

【請求項 18】

前記解析モジュールは、
マクロを用いて可変パラメータリストを作成するための作成サブモジュールと、
前記可変パラメータリストに基づいて前記テンプレートパラメータを反復に解析するた
めの解析サブモジュールと、を含み、
ここで、前記入力情報を前記マクロのマクロパラメータとする
請求項 17 に記載の装置。

【請求項 19】

前記可変パラメータリストは、
処理対象演算子の名称と、
前記演算子カーネル関数の名称と、
前記演算子カーネル関数のデータレイアウトと、
前記深層学習フレームワークを実装するための機器タイプと、
前記深層学習フレームワークを実装するためのデータタイプと、を含む
請求項 18 に記載の装置。

20

【請求項 20】

前記解析サブモジュールは、
前記マクロに含まれるパラメータ数を特定するための第 1 の特定ユニットと、
前記パラメータ数に応じて、反復終了位置を特定するための第 2 の特定ユニットと、
可変パラメータリストを用いて、前記反復終了位置に応じて前記入力情報を反復に処理す
るための処理ユニットと、を含む
請求項 18 に記載の装置。

30

【請求項 21】

前記解析サブモジュールは、
前記可変パラメータリストを用いて前記テンプレートパラメータを反復に解析して、複数
の内部テンプレートパラメータを得るための解析ユニットと、
前記複数の内部テンプレートパラメータをスライシングして、前記複数の完全なテンプレ
ートパラメータを得るためのスライシングユニットと、をさらに含む
請求項 18 ~ 20 のいずれか一項に記載の装置。

40

【請求項 22】

前記複数の内部テンプレートパラメータは複数の文字列を含み、前記スライシングユニ
ットは、
前記入力情報に基づいて、前記複数の文字列をスライシングして、前記処理対象演算子
に対する複数の完全なテンプレートパラメータを得るためのスライシングサブユニット
を含む
請求項 21 に記載の装置。

【請求項 23】

前記処理モジュールは、

50

前記複数の完全なテンプレートパラメータを明示的にインスタンス化することで、前記演算子カーネル関数を処理して、前記深層学習フレームワーク用の利用可能な演算子を得るための明示的なインスタンス化サブモジュールを含む

請求項 2 1 に記載の装置。

【請求項 2 4】

前記複数の完全なテンプレートパラメータの各々は、テンプレート関数記述情報及びテンプレート関数情報を含み、前記明示的なインスタンス化サブモジュールは、

前記テンプレート関数記述情報をヘッダファイルとしてインスタンス化するための第 1 のインスタンス化ユニットと、

前記テンプレート関数情報をソースコードファイルとしてインスタンス化するための第 2 のインスタンス化ユニットと、を含む

請求項 2 3 に記載の装置。

【請求項 2 5】

前記処理モジュールは、

前記複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、該完全なテンプレートパラメータの関数パラメータタイプを特定するための第 1 の特定サブモジュールと、

前記関数パラメータタイプに基づいて、前記完全なテンプレートパラメータの入力情報を対応するパラメータ情報に変換するための変換サブモジュールと、

前記対応するパラメータ情報を記録するための記録サブモジュールと、を含む

請求項 2 1 に記載の装置。

【請求項 2 6】

前記複数の完全なテンプレートパラメータに基づいて、前記深層学習フレームワークに基づく登録対象演算子を作成するための作成モジュールと、

前記登録対象演算子を前記深層学習フレームワーク内部のグローバル演算子テーブルに登録するための登録モジュールと、をさらに含む

請求項 2 1 に記載の装置。

【請求項 2 7】

前記登録対象演算子は登録対象演算子記述情報及び登録対象演算子カーネル関数を含み、前記登録対象演算子カーネル関数には正規化された入力情報及び正規化された関数ポインタが含まれ、前記作成モジュールは、

前記複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、前記完全なテンプレートパラメータに対応する、正規化された形式を有する静的関数を含む構造体を特定するための第 2 の特定サブモジュールと、

前記静的関数の入力情報を前記正規化された入力情報とするための第 1 の正規化サブモジュールと、

前記静的関数の関数ポインタを前記正規化された関数ポインタとするための第 2 の正規化モジュールと、含む

請求項 2 6 に記載の装置。

【請求項 2 8】

前記静的関数の入力情報は入力情報リストを構成し、前記入力情報リストは入力テンソルリストを含み、前記構造体は少なくとも 1 つの特化サブ構造体を含み、前記少なくとも 1 つの特化サブ構造体の各々は対応するデータタイプを有する

請求項 2 7 に記載の装置。

【請求項 2 9】

前記グローバル演算子テーブルは、

演算子カーネル工場クラスと、

演算子カーネル名称クラス及びカーネルキー値クラスと、

演算子カーネルクラスと、を含む

請求項 2 6 に記載の装置。

10

20

30

40

50

【請求項 30】

前記演算子カーネルクラスは、
演算子カーネル関数のポインタと、
演算子カーネル関数入力パラメータの記述情報と、
演算子カーネル関数出力パラメータの記述情報と、の少なくとも1つを含む
請求項 29 に記載の装置。

【請求項 31】

前記深層学習フレームワークに依存しないテンプレートパラメータは、前記処理対象演算子を配布するために用いられ、前記深層学習フレームワークに依存しないテンプレートパラメータは、前記深層学習フレームワークの実装機器及び前記深層学習フレームワークのデータタイプのいずれにも依存しないテンプレートパラメータを含む
請求項 17 に記載の装置。

10

【請求項 32】

前記利用可能な演算子に対する呼び出しコマンドにตอบสนองして、前記利用可能な演算子を呼び出すための呼び出しモジュールと、
呼び出された前記利用可能な演算子をコンパイルすることで、前記利用可能な演算子に対応する関数機能を実行するための実行モジュールと、をさらに含む
請求項 17 に記載の装置。

【請求項 33】

少なくとも1つのプロセッサと、
前記少なくとも1つのプロセッサと通信接続するメモリと、を含み、
前記メモリに、前記少なくとも1つのプロセッサによって実行され得るコマンドが記憶されており、前記コマンドが前記少なくとも1つのプロセッサによって実行されることで、前記少なくとも1つのプロセッサが請求項 1 ~ 16 のいずれか一項に記載の方法を実行することができる、
電子機器。

20

【請求項 34】

コンピュータに請求項 1 ~ 16 のいずれか一項に記載の方法を実行させるためのコンピュータコマンドを記憶している、
非一時的なコンピュータ読取可能な記憶媒体。

30

【請求項 35】

プロセッサにより実行される場合に、請求項 1 ~ 16 のいずれか一項に記載の方法を実現するコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本願は、2022年5月19日に提出した中国特許出願No. 202210560831.9の優先権を要求し、その内容をまとめてここで参考にする。

【0002】

<技術分野>

本開示はコンピュータ技術分野に関し、特に深層学習等の人工知能技術分野に関する。より具体的に、本開示は深層学習フレームワークの演算子処理方法、装置、電子機器、記憶媒体ならびにコンピュータプログラムを提供する。

40

【背景技術】

【0003】

人工知能技術の発展に伴い、深層学習フレームワークはますます注目されてきている。深層学習フレームワークの評価指標はフレームワークの拡張性を含む。

【発明の概要】

【課題を解決するための手段】

【0004】

50

本開示は、深層学習フレームワークの演算子処理方法及び装置、電子機器、記憶媒体ならびにコンピュータプログラムを提供する。

【0005】

本開示の一態様によれば、深層学習フレームワークの演算子処理方法を提供し、該方法は、演算子カーネル関数、及び深層学習フレームワークに依存しないテンプレートパラメータを含む処理対象演算子を取得することと、処理対象演算子に対する入力情報を受信したことに応答して、入力情報を用いてテンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得ることと、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得ることと、を含む。

10

【0006】

本開示の別の態様によれば、深層学習フレームワークの演算子処理装置を提供し、該装置は、演算子カーネル関数、及び深層学習フレームワークに依存しないテンプレートパラメータを含む処理対象演算子を取得するための取得モジュールと、処理対象演算子に対する入力情報を受信したことに応答して、入力情報を用いてテンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得るための解析モジュールと、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得るための処理モジュールと、を含む。

20

【0007】

本開示の別の態様によれば、少なくとも1つのプロセッサと、該少なくとも1つのプロセッサと通信接続するメモリと、を含み、該メモリに、該少なくとも1つのプロセッサによって実行され得るコマンドが記憶されており、該コマンドが該少なくとも1つのプロセッサによって実行されることで、該少なくとも1つのプロセッサが本開示が提供した方法を実行することができる、電子機器を提供した。

【0008】

本開示の別の態様によれば、コンピュータに本開示が提供した方法を実行させるためのコンピュータコマンドを記憶している、非一時的なコンピュータ読取可能な記憶媒体を提供した。

【0009】

本開示の別の態様によれば、プロセッサにより実行される場合に、本開示が提供した方法を実現するコンピュータプログラムを提供した。

30

【0010】

理解されるべきこととして、本部分に記載された内容は、本開示の実施例のキーポイント又は重要な特徴を示すことを意図するものではなく、本開示の範囲を限定するものでもない。本開示の他の特徴は、以下の説明により容易に理解される。

【0011】

ここで、図面は、本技術案をよりよく理解するために用いられ、本開示を限定するものではない。

【図面の簡単な説明】

40

【0012】

【図1】本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0013】

【図2】本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0014】

【図3】本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0015】

50

【図4】本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0016】

【図5】本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0017】

【図6】本開示の一実施例によるグローバル演算子テーブルの例示的模式図である。

【0018】

【図7】本開示の一実施例による深層学習フレームワークの演算子処理装置のブロック図である。

【0019】

【図8】本開示の一実施例による深層学習フレームワークの演算子処理方法を適用可能な電子機器のブロック図である。

【発明を実施するための形態】

【0020】

以下、図面を参照して本開示の例示的な実施例を説明する。ここで、より理解しやすいために本開示の実施例の様々な詳細は含まれ、それらが例示的なものであると考えられるべきである。したがって、当業者であれば、ここで記載される実施例に対して様々な変更・修正を行うことができ、本開示の範囲及び精神から逸脱することはないと分るべきである。同様に、明確かつ簡潔に説明するために、以下の記載において周知の機能や構成に対する説明を省略する。

【0021】

深層学習フレームワークは、人工知能分野のインフラストラクチャであり、より成熟した発展段階に入った。上層拡張および二次開発をより効率的かつ便利的にサポートすることができるかどうかは、深層学習フレームワーク製品が「インフラストラクチャ」として機能することができるかどうか、人工知能産業の着地を実現することができるかどうか、および分野拡張を行うことができるかどうかを評価する重要な指標である。

【0022】

深層学習フレームワークは、PaddlePaddle（パドル）フレームワーク、PyTorchフレームワーク、TensorFlowフレームワーク、MindSporeフレームワーク及び他の複数のニッチフレームワークを含む。これらのフレームワークは、類似の機能を提供することができる。しかしながら、人工知能時代の生態建設の「基石」として、深層学習フレームワークは、フレームワークがさらに加工されて改造されるように、より多くの上流開発者を吸引する必要がある。そのため、深層学習フレームワークは、上流開発者がフレームワークの生態を共同構築するコストを低減するために、設計が明確で、メンテナンスしやすく、理解しやすく、カスタマイズ化拡張開発しやすいなどの特徴を有する必要がある。演算子体系は、1つのフレームワークの半数以上のコードを占めることができ、深層学習フレームワークの主体とすることができる。低コストで演算子体系を拡張して追加することができるか否かは、そのうちのフレームワークの生態を共同構築するコストを低減する重要な一環である。

【0023】

演算子は、深層学習フレームワークにおけるテンソル計算ユニットである。演算子は、特定のテンソル組み合わせを入力とし、特定の計算ロジックを完成し、計算結果を返すことができる。深層学習フレームワークを用いてモデルを構築するプロセスは、異なる演算子を組み合わせる特定の計算ロジックを満たすプロセスであってもよい。

【0024】

演算子の多様性と豊富さは、深層学習フレームワークを評価する重要な指標とすることができる。深層学習フレームワークにおける演算子は、継続的に動的に追加されてもよい。異なる演算子の間には、関連もあるし、区別もある。複雑な演算子は、簡単な演算子の組み合わせによって実装されることが多い。深層学習フレームワークがこのような演算子の

10

20

30

40

50

組み合わせによる実装のモードをサポートすると、コスト削減及び効率アップに寄与し、外部開発者を、該フレームワークを二次開発するように吸引することにも寄与する。

【 0 0 2 5 】

関数は、プログラミング言語における基礎パラダイムである。関数同士は互いに呼び出されてもよい。深層学習フレームワークに便利な演算子組合せ開発能力を持たせるために、関数式演算子体系に基づいて深層学習フレームワークを構築することができる。

【 0 0 2 6 】

関数式演算子体系の応用が少ない。関数と演算子との間には矛盾がある。例えば、関数自体が簡単すぎて、情報と特徴を携帯する能力がない。また、関数が柔軟になりすぎ、形式が多様であり、パラメータが多様であり、正規化が困難である。演算子は、比較的複雑な概念であり、特徴記述を有し、様々な機器に適合する必要がある。また、演算子は、フレームワークの統合スケジューリングユニットとして、グローバルスケジューリングに適合するように、一致する形態を抽象化する必要がある。

【 0 0 2 7 】

現実的に考慮すると、深層学習フレームワーク製品は、トレードオフの設計を行うことが多く、関数式のパラダイムを放棄し、「構造体+内部計算方法」の形式を使用し、組み合わせ多重化の利便性を犠牲にして、演算子のパラダイムの統一を保持する。しかしながら、このようなトレードオフ設計は、フレームメンテナンスコスト及び拡張コストの増大を招く。

【 0 0 2 8 】

いくつかの深層学習フレームワーク製品は、演算子の多様性を遮蔽するために、トップレベル演算子関数をカプセル化することができる。これらの深層学習フレームワーク製品は、関数内部で配布選択を行う方式を利用して、演算子同士の多重化を実現することができる。演算子の間の多重化時に、これらの深層学習フレームワーク製品は、配布選択のスケジューリングオーバーヘッドが発生し、性能が損なわれる。これらの深層学習フレームワーク製品は、演算子の組み合わせ多重化の利便性を保有するが、一部の性能を犠牲にする。演算子の組み合わせを行わなければ、性能をより良くすることができ、演算子の多重化方式を放棄し、複数の演算子カーネルのコードを直接にまとめて、性能を向上させることができる。

【 0 0 2 9 】

したがって、関数式演算子同士の多重化性能の問題を解決すれば、関数式演算子の優位性を永続的に発揮させることができる。

【 0 0 3 0 】

「関数式」と「高性能多重化」を両立する演算子体系を形成するために、形式が簡単で、柔軟で多様な「関数」を形式が統一され且つ記述特徴を有する演算子として自動的にパッケージし、「関数式」の開発を実現する必要がある、それにより関数式演算子の多重化時に追加のスケジューリング配布オーバーヘッドを引き込まないことを確保し、「高性能多重化」を実現する。

【 0 0 3 1 】

関数式演算子体系を実装するために、関数式演算子のパラダイム設計と関数式演算子の解析及び登録が必要である。関数式演算子のパラダイム設計は、関数の書き方により演算子の計算カーネルを実装することである。関数において、一般的にC++プログラミング言語を用いて演算子のコア計算ロジックを作成する。関数式演算子は、自然に演算子同士の相互多重化をサポートする。関数式演算子の解析及び登録は、キーコンポーネントを設計し、形式が簡単で、柔軟で多様な関数を形式が統一され且つ記述特徴を有する演算子として自動的にパッケージし、フレームワークの統合スケジューリングに供する。

【 0 0 3 2 】

また、多重化可能な関数式演算子体系を実装するために、関数式演算子同士が機器非依存の多重化をサポートする必要もある。同一の演算子が異なる機器に適用される場合、異なるカーネル関数が必要となる。例えば、演算子がCPU (Central Process

10

20

30

40

50

sing Unit, 中央処理装置)に適用される場合、CPUの演算子カーネルが必要であり、GPU(Graphics Processing Unit, グラフィックプロセッサ)がGPUの演算子カーネルが必要であり、関数式演算子であっても、この2つのカーネルは異なる関数である必要があり、1つの演算子を新たに開発し、それが他の演算子を多重化して組み合わせる実装することができる場合、多重化される演算子は機器非依存である必要があり、そうでない場合、重複開発を招き、簡単な例を挙げる。

【0033】

2つの演算子が既に存在すると仮定する。

【0034】

- 演算子A :

- 演算子A CPUカーネル : A_CPU()

- 演算子A GPUカーネル : A_GPU()

【0035】

- 演算子B :

- 演算子B CPUカーネル : B_CPU()

- 演算子B GPUカーネル : B_GPU()

【0036】

例えば、演算子A及び演算子Bに基づいて1つの新たな演算子を実装する。

【0037】

1つの実装方式は、以下の通りである。

【0038】

- 演算子C (演算子AとBを多重化して実装する) :

- 演算子C CPUカーネル : C_CPU(){ A_CPU();B_CPU(); }

- 演算子C GPUカーネル : C_GPU(){ A_GPU();B_GPU(); }

【0039】

該実装方式は合理的な関数式演算子体系ではなく、冗長コードが多く、メンテナンスコストが高い。また、このような方式による多重化演算子は、機器分けて多重化することしかできず、各機器の演算子Cカーネルの実装コードがほぼ同じであることを容易に見出すことができる。関数式演算子体系は、演算子の多重化効率を向上させ、メンテナンスコストを低減し、拡張開発を容易にするためのものである。このような各機器カーネルは、1つのカーネルをほぼコピーして実装する必要があり、厳密な関数式演算子体系ではない。

【0040】

上述したように、異なる深層学習フレームワークは、例えばPyTorchフレームワーク、TensorFlowフレームワーク、及びMindSporeフレームワークを含んでもよい。

【0041】

PyTorchフレームワークでは、演算子は最終的に1つの関数式のインターフェースにカプセル化され、関数式演算子同士は多重化されてもよく、PyTorchの1つの演算子は複数の層に分けて実装され、上層が下層を呼び出す場合、次の層の演算子を検索して実装する必要がある。PyTorchフレームワークでは、演算子同士の多重化、及び演算子内の多重化の場合、演算子カーネルマッピングテーブルを用いて対応するカーネルを検索して選択する方式を採用する。このような方式のオーバーヘッドは比較的大きく、1つの演算子の実行過程において、複数回の演算子検索を行う必要がある可能性がある。

【0042】

TensorFlowフレームワークでは、演算子は依然として構造体演算子であり、関数式演算子ではない。構造体形式の演算子カーネルの実装は、直感的ではなく、追加の概念を増加し、且つ演算子同士の多重化には、より複雑な書き方及びより大きいスケジューリングオーバーヘッドを引き込む必要があり、設計上で、新たな演算子の組み合わせ開発に不便な形式に自然に属する。

【0043】

10

20

30

40

50

一般的に、構造体形式の演算子体系では、演算子カーネル関数を実装する際に、異なる演算子が当該演算子カーネル関数を多重化できるように、関数を別途定義する必要がある。このような方式は、ある程度の一般化を実現することができるが、メンテナンスコストが高く、規範的管理が困難であり、フレームワークの長期発展により多くの冗長なコードが現れることになり、メンテナンスコストが徐々に増大する。

【0044】

TensorFlowフレームワークの派生製品 `tensorflow/runtime` フレームワークにおいて、演算子は関数式演算子であってもよい。しかし、`tensorflow/runtime` フレームワークは1つの実験的製品に過ぎず、演算子の数が少なく、まだ体系になっていない。また、演算子パラメータリストが混乱し、規範性が悪く、高性能な多重化を実現できない。

10

【0045】

MindSporeフレームワークでの演算子は、TensorFlowフレームワークと類似し、説明を繰り返さない。

【0046】

図1は、本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0047】

図1に示すように、該方法100は、操作S110～操作S130を含んでもよい。

【0048】

20

操作S110において、処理対象演算子を取得する。

【0049】

例えば、処理対象演算子は、演算子カーネル関数、及び深層学習フレームワークに依存しないテンプレートパラメータを含む。

【0050】

例えば、取得された処理対象演算子は、例えば `Scale` (スケーリング) 演算子であってもよい。`Scale` 演算子は、演算子カーネル関数 `scale_kernel` (スケーリング演算子カーネル関数) 及びテンプレートパラメータ `T` を含んでもよい。`scale_kernel` は、`scale` 演算を行うことができる。

【0051】

30

例えば、深層学習フレームワークに依存しないテンプレートパラメータは、例えば深層学習フレームワークの実装機器に依存しないテンプレートパラメータを含んでもよく、深層学習フレームワークのデータタイプに依存しないテンプレートパラメータを含んでもよい。

【0052】

例えば、実装機器は、例えばCPU又はGPUであってもよい。また例えば、データタイプは、例えば `Float` (浮動小数点) タイプ、`Int` (整数) タイプであってもよい。

【0053】

例えば、`Scale` 演算子を処理対象演算子とする場合、処理対象演算子は、例えば `ScaleKernel T, Context ()` と表すことができる。

【0054】

40

操作S120において、処理対象演算子に対する入力情報を受信したことに応じて、入力情報を用いてテンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得る。

【0055】

例えば、1つの完全なテンプレートパラメータは、`Float, CPUContext` であってもよい。`Float` はデータタイプであり、`CPUContext` は機器タイプである。

【0056】

操作S130において、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得る。

50

【 0 0 5 7 】

例えば、処理対象演算子を取得した後、解析により具体的な深層学習フレームワークの利用可能な演算子を得て、該深層学習フレームワークに該利用可能な演算子を登録することで、深層学習フレームワークにおける他のルーチン又はタスクが、実際の深層学習アプリケーションに応じて該利用可能な演算子を呼び出して対応するアルゴリズム又は機能、例えばスケールリング、畳み込み等を実行することができる。

【 0 0 5 8 】

例えば、完全なテンプレートパラメータ `Float`、`CPUContext` に基づいて、演算子カーネル関数 `scale_kernel` を処理することで、利用可能な演算子を得ることができる。利用可能な演算子は、例えば `ScaleKernel Float`、`CPUContext ()` と表すことができる。

10

【 0 0 5 9 】

本開示の実施例によれば、演算子が多重化される時、スケジューリングオーバーヘッドがないことを実現し、高性能多重化を実現することができる。

【 0 0 6 0 】

幾つかの実施例において、方法 100 は、利用可能な演算子に対する呼び出しコマンドに回答して、利用可能な演算子を呼び出すことと、呼び出された利用可能な演算子をコンパイルすることで、利用可能な演算子に対応する関数機能を実行することとをさらに含んでもよい。

【 0 0 6 1 】

例えば、利用可能な演算子 `ScaleKernel float`、`CPUContext ()` に対する呼び出しコマンドに回答して、利用可能な演算子 `ScaleKernel float`、`CPUContext ()` を呼び出し、呼び出された利用可能な演算子 `ScaleKernel float`、`CPUContext ()` をコンパイルすることで、`scale` 関数の機能を実行する。

20

【 0 0 6 2 】

幾つかの実施例において、深層学習フレームワークの実装機器に依存しない演算子カーネル関数の宣言は、例えば以下の通りであってもよい。

【 0 0 6 3 】

【 数 1 】

```
template <typename T, typename Context>
void ScaleKernel(const Context& dev_ctx,
                 const DenseTensor& x,
                 const Scalar& scale,
                 float bias,
                 bool bias_after_scale,
                 DenseTensor* out);
```

30

40

【 0 0 6 4 】

`Template` はテンプレートである。 `typename T` は、深層学習フレームワークのデータタイプに依存しないテンプレートパラメータである。 `typename Context` は、深層学習フレームワークの実装機器に依存しないテンプレートパラメータである。 `typename T` 及び `typename Context` は、具体的な実装機器にバインドされなくてもよい。

【 0 0 6 5 】

50

また例えば、基本演算を多重化可能な演算子 `Calc` は、例えば以下の通りであってもよい。

【0066】

【数2】

```
template <typename T, typename Context>
void CalcKernel(const Context& dev_ctx, ...) {
    ScaleKernel<T, Context>(...)
}
```

10

【0067】

演算子 `Calc` が呼び出した演算子は、演算子 `ScaleKernel T, Context ()` である。演算子 `ScaleKernel T, Context ()` は、実装機器に依存しない演算子である。コンパイル時に、テンプレートパラメータ `T` 及び `Context` を取得し、コンパイラにより具体的なデータタイプ及び実装機器に変換する。一例において、演算子 `ScaleKernel T, Context ()` が `ScaleKernel float, CPUContext ()` に変換されてもよい。実行時に、演算子 `Calc` は、`scale` 演算の `float` タイプを呼び出し、`CPU` 機器のコマンドセットに基づいて、他のロジックを必要とせずに判断及び選択を行い、スケジューリングオーバーヘッド無しの多重化ができて、高性能多重化を実現することができる。

20

【0068】

幾つかの実施例において、処理対象演算子に対する入力情報を受信したことに応じて、入力情報を用いてテンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得ることは、マクロを用いて可変パラメータリストを作成することと、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析することを含み、ここで、入力情報をマクロのマクロパラメータとする。

【0069】

例えば、マクロは、`C` 言語のメカニズムの一種であり、重複のコードを簡略化することに用いられ、コンパイル前処理時に、該マクロを対応するコードセグメントに置換えて実行する。

30

【0070】

例えば、上述した処理対象演算子 `ScaleKernel T, Context ()` に対して、入力情報に基づいて、可変パラメータリストを特定してもよい。マクロにより可変パラメータリストに基づいて反復的に解析する。一例において、マクロ `PD_REGISTER_KERNEL` は、可変パラメータリストを解析することができる。

【0071】

【数3】

40

```

PD_REGISTER_KERNEL(scale,
                    CPU,
                    ALL_LAYOUT,
                    ScaleKernel,
                    float,
                    double,
                    bfloat16,
                    uint8_t,
                    int8_t,
                    int16_t,
                    int,
                    int64_t) {}

```

【0072】

幾つかの実施例において、可変パラメータリストは、処理対象演算子の名称と、演算子カーネル関数の名称と、演算子カーネル関数のデータレイアウトと、深層学習フレームワークを実装するための機器タイプと、深層学習フレームワークを実装するためのデータタイプとを含む。

【0073】

例えば、上述した `Scale` 演算子に対して、マクロを用いて作成されたパラメータリストにおいて、`scale` は演算子名称であり、`CPU` は深層学習フレームワークを実装するための機器タイプであり、`ALL_LAYOUT` は演算子カーネル関数 `ScaleKernel` がすべてのデータレイアウトに適合し得ることを表し、`ScaleKernel` は登録及び解析しようとする演算子カーネル関数名を表し、`float`、`double`、`bfloat16`、`uint8_t`、`int16_t`、`int`、`int64_t` は深層学習フレームワークがサポートするデータタイプを表す。

【0074】

幾つかの実施例において、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析することは、マクロに含まれるパラメータ数を特定することと、パラメータ数に応じて、反復終了位置を特定することと、可変パラメータリストを用いて、反復終了位置に応じて入力情報を反復的に処理することとを含む。

【0075】

例えば、上述した処理対象演算子 `ScaleKernel T, Context ()` に対して、マクロ `PD_REGISTER_KERNEL` を用いてパラメータリストを解析する時に、別のマクロ `PD_NARGS` を用いてマクロ `PD_REGISTER_KERNEL` に含まれるパラメータ数を特定してもよい。

【0076】

また例えば、パラメータ数を特定した後、パラメータ数をマクロの添え字とすることで、下層マクロ方法をスプリシングし、その後、該マクロ方法呼び出して解析することができる。

【0077】

例えば、マクロ `PD_NARGS` は、例えば以下の通りであってもよい。

10

20

30

40

50

【 0 0 7 8 】

【数 4】

```
#define PD_NARGS(...) _PD_NARGS((__VA_ARGS__, _PD_RESQ_N()))

#define _PD_NARGS(...) _PD_ARG_N(__VA_ARGS__)

#define _PD_ARG_N_EXPAND(
    \
    _1, _2, _3, _4, _5, _6, _7, _8, _9, _10, _11, _12, _13, _14, _15, N, ...) \
    N

#define _PD_ARG_N(args) _PD_ARG_N_EXPAND args

#define _PD_RESQ_N() 15, 14, 13, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1, 0
```

10

【 0 0 7 9 】

例えば、PD_NARGS(...)に対して、「...」は上述した可変パラメータリストであってもよい。本開示の実施例によれば、マクロに幾つのパラメータが存在するかを特定し、対応するパラメータ数の位置からマクロパラメータを反復的に解析することで、反復終了位置を特定することができる。

20

【 0 0 8 0 】

他の実施例において、マクロPD_NARGSが関連処理を行うために、マクロPD_NARGSを可変パラメータマクロVA_ARGSの外にネストしてもよい。

【 0 0 8 1 】

例えば、以下の方式によってマクロPD_NARGSを可変パラメータマクロVA_ARGSの外にネストすることができる。

【 0 0 8 2 】

【数 5】

```
_PD_KERNEL_REGISTRAR_INIT(PD_NARGS(__VA_ARGS__), \
    \
    reg_type, \
    \
    kernel_name, \
    \
    backend, \
    \
    context, \
    \
    layout, \
    \
    args_def_fn, \
```

30

40

【 0 0 8 3 】

反復終了位置が特定された後、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析することができる。

【 0 0 8 4 】

幾つかの実施例において、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析することは、可変パラメータリストを用いて、入力情報に基づいてテンプレートパラメータを反復的に解析して、複数の内部テンプレートパラメータを得ることと、複数の内部テンプレートパラメータをスプライシングして、複数の完全なテンプレートパラメータを得ることとをさらに含んでもよい。以下、図2を参照して詳細に説明する。

50

【 0 0 8 5 】

図 2 は、本開示の別の実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【 0 0 8 6 】

図 2 に示すように、方法 2 2 1 は、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析する。以下、操作 S 2 2 1 1 及び操作 S 2 2 1 2 を参照して詳細に説明する。

【 0 0 8 7 】

操作 S 2 2 1 1 において、可変パラメータリストを用いてテンプレートパラメータを反復的に解析して、複数の内部テンプレートパラメータを得る。

10

【 0 0 8 8 】

例えば、マクロ V A _ A R G S を用いてマクロパラメータを反復することで、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析する。

【 0 0 8 9 】

例えば、可変パラメータリストを用いて、以下の方式によって反復的に解析することで、テンプレートパラメータを得ることができる。

【 0 0 9 0 】

【 数 6 】

```
#define _PD_KERNEL_REGISTRAR_INIT_3(cpp_dtype, ...)
#define _PD_KERNEL_REGISTRAR_INIT_4(...)
    _PD_KERNEL_REGISTRAR_INIT_3(_VA_ARGS_)
```

20

【 0 0 9 1 】

マクロ P D _ K E R N E L _ R E G I S T R A R _ I N I T _ 4 に対して、そのマクロパラメータ「...」は、上述した可変パラメータリストである。一例において、マクロパラメータは、例えば f l o a t 、 d o u b l e 、 b f l o a t 1 6 、 u i n t 8 _ t 、 i n t 8 _ t 、 i n t 1 6 _ t 、 i n t 、 i n t 6 4 _ t を含む。

【 0 0 9 2 】

マクロ P D _ K E R N E L _ R E G I S T R A R _ I N I T _ 4 の内部は、マクロ P D _ K E R N E L _ R E G I S T R A R _ I N I T _ 3 を呼び出すことができる。

30

【 0 0 9 3 】

マクロ P D _ K E R N E L _ R E G I S T R A R _ I N I T _ 3 は、マクロ V A _ A R G S の外にネストして、可変パラメータリストを全体として反復するために用いられる。マクロ P D _ K E R N E L _ R E G I S T R A R _ I N I T _ 3 は、可変パラメータリストにおける最初のパラメータを解析し、例えば f l o a t を解析することができる。この場合、 c p p _ d t y p e は f l o a t であってもよく、マクロ P D _ K E R N E L _ R E G I S T R A R _ I N I T _ 4 のマクロパラメータ「...」には、 d o u b l e 、 b f l o a t 1 6 、 u i n t 8 _ t 、 i n t 8 _ t 、 i n t 1 6 _ t 、 i n t 、 i n t 6 4 _ t が含まれてもよい。次回の反復時に、 d o u b l e は可変パラメータリストにおける最初のパラメータになる。

40

【 0 0 9 4 】

操作 S 2 2 1 2 において、複数の内部テンプレートパラメータをスプライシングして、複数の完全なテンプレートパラメータを得る。

【 0 0 9 5 】

本開示の実施例において、複数の内部テンプレートパラメータは複数の文字列を含む。

【 0 0 9 6 】

本開示の実施例において、複数の内部テンプレートパラメータをスプライシングして、複数の完全なテンプレートパラメータを得ることは、入力情報に基づいて、複数の文字列をスプライシングして、処理対象演算子に対する複数の完全なテンプレートパラメータを得

50

ることを含む。

【0097】

例えば、解析された内部テンプレートパラメータは文字列である。複数の文字列は、例えば文字列 `Float` 及び文字列 `CPUContext` を含む。この2つの文字列をスライシングして、完全なテンプレートパラメータ `Float`、`CPUContext` を得ることができる。

【0098】

幾つかの実施例において、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得ることは、複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータについて、該完全なテンプレートパラメータの関数パラメータタイプを特定することと、関数パラメータタイプに基づいて、完全なテンプレートパラメータの入力情報を対応するパラメータ情報に変換することと、対応するパラメータ情報を記録することとを含む。以下、図3を参照して詳細に説明する。

10

【0099】

図3は、本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【0100】

図3に示すように、該方法330は、例えばI個の完全なテンプレートパラメータをトラバースし、第i個の完全なテンプレートパラメータに対して操作S331～操作S335を含む操作を実行することで、完全なテンプレートパラメータの入力情報を対応するパラメータ情報に変換する。Iは1以上の整数であり、iはI以下の整数である。

20

【0101】

操作S331において、第i個の完全なテンプレートパラメータの関数パラメータタイプを特定する。

【0102】

例えば、第i個の完全なテンプレートパラメータに基づいて、幾つかの関数パラメータのタイプを関数パラメータタイプとして特定することができる。

【0103】

操作S332において、関数パラメータタイプに基づいて、第i個の完全なテンプレートパラメータの入力情報を第i組のパラメータ情報に変換する。

30

【0104】

例えば、異なる関数パラメータタイプは、異なる変換方式に対応していてもよい。関数パラメータタイプと変換方式との間の対応関係は予め設定されてもよい。

【0105】

また例えば、上述した処理対象演算子に対する入力情報には、第i個の完全なテンプレートパラメータに対応する第i個の入力情報が含まれてもよい。第i個の入力情報には、複数の関数パラメータが含まれてもよい。対応する関数パラメータタイプに基づいて、入力情報を対応する第i組のパラメータ情報に変換することができる。

【0106】

操作S333において、第i組のパラメータ情報を記録する。

40

【0107】

例えば、第i組のパラメータ情報をメモリに記録してもよい。

【0108】

操作S334において、iがIに等しいか否かを判定する。

【0109】

i=Iの場合、操作S335を実行し、フローを終了する。i=Iの場合、既にI個の完全なテンプレートパラメータがトラバースされ、フローを終了することができる。

【0110】

iがIより小さい場合、操作S331に戻り、第i+1個の完全なテンプレートパラメータ

50

に対して操作 S 3 3 1 ~ 操作 S 3 3 4 を実行する。

【 0 1 1 1 】

一例において、データタイプが 8 種類であり、且つ機器種別が 2 種類である場合、反復トラバースの方式によって、16 組のパラメータ情報を得て、異なる機器タイプ及びデータタイプの組み合わせをサポートすることができる。

【 0 1 1 2 】

一例において、上述した第 i 個の関数の入力情報に対して、第 i 個の関数の関数パラメータタイプを特定することができ、そのタイプ及び関連する必要な情報を記録する。例えば、C++ テンプレートメタプログラミング技術を用いることで、コンパイル期間の計算、判断、変換、照会等の機能を実現する。

10

【 0 1 1 3 】

幾つかの実施例において、例えば上述した方法 3 0 0 は、複数の完全なテンプレートパラメータに基づいて、深層学習フレームワークに基づく登録対象演算子を作成することと、登録対象演算子を深層学習フレームワーク内部のグローバル演算子テーブルに登録することとをさらに含んでもよい。

【 0 1 1 4 】

深層学習フレームワーク内部において、統合スケジューリングを実現するために、深層学習フレームワーク内部の演算子に対応する関数書き方は正規化されもよく、且つ該正規化された書き方とユーザのカスタマイズ演算子の計算関数の作成方式とは大きく異なる。

【 0 1 1 5 】

例えば、異なるシーンにおいて、異なるユーザの計算関数の書き方は相違する可能性があり、例えば 1 つのテンソルを入力とする場合、関数の書き方は以下の通りである。

20

【 0 1 1 6 】

`Tensor` リスト演算子計算関数名(入力 `Tensor 1`) { ... } を返す。

【 0 1 1 7 】

一方、2 つの入力テンソルがあれば、関数の書き方は、以下の通りである。

【 0 1 1 8 】

`Tensor` リスト演算子計算関数名(入力 `Tensor 1` , 入力 `Tensor 2`) { ... } を返す。

【 0 1 1 9 】

より多くの入力がある場合、より多くの関数宣言もある。C++ プログラミング言語において、異なる関数宣言は、異なる関数ポインタタイプ、つまり、異なるデータタイプを意味し、フレームワークはユーザが作成し得る関数ポインタタイプを記憶する必要があり、フレームワークの最下層で呼び出すことができ、このような書き方の柔軟性が極めて高いパラダイムは、C++ プログラミング言語において、簡潔で汎用的な形式で記憶することができない。全体の簡潔を実現するために、ユーザが作成した様々な計算関数を正規化する必要もある。

30

【 0 1 2 0 】

例えば、1 つのカスタマイズ計算関数の関数形式は以下の通りである。

【 0 1 2 1 】

`Tensor` リスト演算子計算関数名(入力 `Tensor 1`) { ... } を返す。

【 0 1 2 2 】

別のカスタマイズ計算関数の関数形式は以下の通りである。

【 0 1 2 3 】

`Tensor` リスト演算子計算関数名(入力 `Tensor 1` , 入力 `Tensor 2`) { ... } を返す。

【 0 1 2 4 】

上記 2 つのカスタマイズ計算関数の関数形式が統一されておらず、本実施例において、深層学習フレームワーク内部で統合スケジューリングを容易にするために、異なるカスタマイズ計算関数を正規化する必要がある。

40

50

【 0 1 2 5 】

本開示の実施例において、登録対象演算子は、登録対象演算子記述情報及び登録対象演算子カーネル関数を含み、登録対象演算子カーネル関数は、正規化された入力情報及び正規化された関数ポインタを含み、複数の完全な内部テンプレート関数に基づいて、深層学習フレームワークに基づく登録対象演算子を作成することは、複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、完全なテンプレートパラメータに対応する、正規化された形式を有する静的関数を含む構造体を特定することと、静的関数の入力情報を正規化された入力情報することと、静的関数の関数ポインタを正規化された関数ポインタとすることとを含む。

【 0 1 2 6 】

例えば、深層学習フレームワーク内の各種の演算子は、`OpKernel`クラスを継承し、自身の計算関数(`Compute`)を充填して実装することができる。演算子の対応する入出力テンソルは、実行コンテキスト(`ExecutionContext`)に記憶されてもよく、実行コンテキストの`Input`及び`Output`機能によって取得することができる。一例において、登録対象演算子カーネル関数を構築する時、正規化された形式の関数に基づいて構築することができ、該正規化された関数は、完全なテンプレートパラメータに対応する構造体における静的関数であってもよい。

【 0 1 2 7 】

また例えば、登録対象演算子カーネル関数を構築する時、上記静的関数の入力情報を作成することと、上記静的関数の関数ポインタを作成することと、他の標準処理、例えば、`context`から入力情報を取得して`inputs`変数に入れる処理を作成することとを含む。

【 0 1 2 8 】

登録対象演算子カーネル関数を構築する時に統一された形式の静的関数を採用することで、深層学習フレームワークの統合スケジューリングのために、カスタマイズ計算関数を統一することができる。

【 0 1 2 9 】

さらに、幾つかの実施例において、静的関数の入力情報は入力情報リストを構成し、入力情報リストは入力テンソルリストを含み、構造体は少なくとも1つの特化サブ構造体を含み、少なくとも1つの特化サブ構造体の各々は対応するデータタイプを有し、データタイプはテンソル及び非テンソルを含む。

【 0 1 3 0 】

例えば、正規化された形式の静的関数の入力情報は入力情報リストであってもよく、さらに、入力情報リストは入力テンソルリスト及び属性リストを含んでもよく、属性リストは他のデータタイプの入力変数を記憶するために用いられる。一例において、入力変数のデータタイプはテンソルタイプ(`tensor`)、整数タイプ(`int`)、浮動小数点タイプ(`float`)を含み、`tensor1`、`tensor2`、`int`、`float`と表す。静的関数の入力情報リストは、テンソルリスト(`tensor1`、`tensor2`)及び属性リスト(`int`、`float`)を含んでもよい。

【 0 1 3 1 】

例えば、入力情報リストは、C++のテンプレートパラメータ導出メカニズムを用いて取得することができる。

【 0 1 3 2 】

以下、図4を参照して、本開示の一実施例による深層学習フレームワークの演算子処理方法を詳細に説明する。

【 0 1 3 3 】

図4は、本開示の一実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

【 0 1 3 4 】

図4に示すように、方法441は、演算子カーネル関数入力情報と関数ポインタとの正規

10

20

30

40

50

化を実現することができる。以下、操作 S 4 4 1 1 ~ 操作 S 4 4 1 5 を参照して、詳細に説明する。

【 0 1 3 5 】

操作 S 4 4 1 1 において、登録対象演算子を取得する。

【 0 1 3 6 】

例えば、登録対象演算子は、登録対象演算子記述情報及び登録対象演算子カーネル関数を含む。登録対象演算子記述情報は、登録対象演算子入力情報及び登録対象演算子出力情報を含む。登録対象演算子カーネル関数は、該演算子の主な計算ロジックを記述するために用いられる。

【 0 1 3 7 】

操作 S 4 4 1 2 において、登録対象演算子の完全なテンプレートパラメータに対応する構造体に基づいて、少なくとも 1 つのサブ構造体を特定する。

【 0 1 3 8 】

例えば、登録対象演算子は、複数の完全なテンプレートパラメータに基づいて作成されてもよい。各完全なテンプレートパラメータは、1 つの構造体に対応してもよい。構造体は、例えば予め設定されてもよい。該構造体は、関数戻り値、可変関数パラメータリスト及び 1 つの関数ポインタをテンプレートパラメータとしてもよく、ここで、関数ポインタの戻り値タイプとパラメータリスト及びテンプレートパラメータにおけるタイプとが一致している。該構造体は 1 つの計算用の静的関数を含む。該静的関数は正規化された形式を有する。

【 0 1 3 9 】

また例えば、該構造体は例えば少なくとも 1 つの特化サブ構造体を含んでもよい。特化サブ構造体はテンプレートパラメータ導出を行うために用いられる。

【 0 1 4 0 】

操作 S 4 4 1 3 において、1 つの特化サブ構造体と登録対象演算子の入力情報とをマッチングして、構造体の静的関数の入力サブ情報を得る。

【 0 1 4 1 】

例えば、特定のデータタイプの構造体にマッチングするように、各特化サブ構造体の特化実装を予め設定することができる。各サブ構造体の特化実装は、いずれも該構造体にマッチングするデータタイプを最初のテンプレートパラメータとし、その後、残りのテンプレートパラメータリストとする。

【 0 1 4 2 】

また例えば、少なくとも 1 つの特化サブ構造体は、終了マーカを有する 1 つの特化サブ構造体を含む。該終了マーカを有する特化サブ構造体は、終了マーカをテンプレートパラメータとすることができる。

【 0 1 4 3 】

また例えば、特化サブ構造体はサブ静的関数を含んでもよい。サブ静的関数はパラメータインデックスをテンプレートパラメータとし、現在何番目のパラメータにマッチングしたかをマークすることに用いられる。

【 0 1 4 4 】

操作 S 4 4 1 4 において、終了マーカにマッチングしたか否かを判定する。

【 0 1 4 5 】

終了マーカにマッチングした場合、操作 S 4 4 1 5 を実行し、フローを終了する。

【 0 1 4 6 】

終了マーカにマッチングしなかった場合、操作 S 4 4 1 3 に戻り、次のサブ構造体に基づいてマッチングする。

【 0 1 4 7 】

例えば、コンパイル過程において、コンパイラは、演算子カーネル関数の入力情報を一つずつ解析する。現在解析されている入力情報が前述定義されたサブ構造体のある特化実装の第 1 個のテンプレートパラメータとマッチングした場合、該特化サブ構造体の特化実装

10

20

30

40

50

を呼び出し、現在パラメータに対する解析を完了する。次に、引き続き残りの入力情報を他の特化サブ構造体とマッチングして、後の入力情報を解析し、終了マークをテンプレートパラメータとする特化サブ構造体とマッチングした場合、マッチングが完了する。

【0148】

さらに、幾つかの実施例において、正規化された形式を有する関数の関数形式は、以下の通りであってもよい。

【0149】

Tensor リスト演算子計算関数名(入力 Tensor リスト){ ...}を返す。

【0150】

図5は、本開示の別の実施例による深層学習フレームワークの演算子処理方法のフローチャートである。

10

【0151】

操作S510において、処理対象演算子を取得する。

【0152】

例えば、処理対象演算子は、演算子カーネル関数及びテンプレートパラメータを含む。該テンプレート関数は、深層学習フレームワークに依存しておらず、処理対象演算子を配布することに用いられる。

【0153】

次に、方法500は、処理対象演算子に対する入力情報を受信したことに応じて、入力情報を用いてテンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得る。以下、操作S521～操作S523を参照して詳細に説明する。

20

【0154】

操作S521において、マクロを用いて可変パラメータリストを作成する。

【0155】

操作S522において、可変パラメータリストを用いてテンプレートパラメータを反復的に解析して、複数の内部テンプレートパラメータを得る。

【0156】

操作S523において、複数の内部テンプレートパラメータをスプライシングして、複数の完全なテンプレートパラメータを得る。

30

【0157】

理解できることとして、操作S522及び操作S523の詳細な説明について、上述した操作S2211及び操作S2212を参照することができ、本開示はここで説明を繰り返さない。

【0158】

次に、方法500は、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得ることができる。以下、操作S531、S532及び操作S533を参照して詳細に説明する。

【0159】

操作S531において、複数の完全なテンプレートパラメータを明示的にインスタンス化する。

40

【0160】

本開示の実施例において、複数の完全なテンプレートパラメータの各々は、テンプレート関数記述情報及びテンプレート関数情報を含む。

【0161】

例えば、テンプレート関数記述情報をヘッダファイルとしてインスタンス化する。テンプレート関数記述情報は、テンプレート関数宣言であってもよい。ヘッダファイルは、例えば「.h」ファイルであってもよい。

【0162】

例えば、テンプレート関数情報をソースコードファイルとしてインスタンス化する。テン

50

プレート関数情報は、テンプレート関数実装であってもよい。ソースコードファイルは、例えば「.cc」ファイルであってもよい。

【0163】

C++プログラミング言語において、テンプレート関数を適用するカーネル形式には一定の制限がある。コンパイル時に、テンプレート関数は、実際のテンプレートパラメータタイプに基づいてインスタンス化し、テンプレート関数の実装も対応する呼び出し位置に展開する必要がある。テンプレート関数記述情報をヘッダファイルにインスタンス化し、ソースコードファイルをソースコードファイルにインスタンス化することで、コンパイル時に大量のテンプレート関数コードを導入して展開することによりコンパイル効率に影響することを回避することができる。

10

【0164】

操作S532において、完全なテンプレートパラメータの入力情報を対応するパラメータ情報に変換する。

【0165】

理解できることとして、操作S532の詳細な説明について、上述した方法330を参照することができる。本開示はここで説明を繰り返さない。

【0166】

操作S533において、登録対象演算子カーネル関数の入力情報及び関数ポインタの正規化を実現する。

【0167】

理解できることとして、操作S533の詳細な説明について、上述した方法441を参照することができる。本開示はここで説明を繰り返さない。

20

【0168】

操作S540において、演算子を深層学習フレームワーク内部のグローバル演算子テーブルに登録する。

【0169】

例えば、操作S532又は操作S533を実行した後、演算子をグローバル演算子テーブルに登録する情報を取得することができる。これらの情報は、例えば演算子カーネル関数の対応するパラメータ情報及び/又は正規化された関数ポインタを含んでもよい。

【0170】

操作S550において、グローバル演算子テーブルに基づいて、利用可能な演算子を呼び出す。

30

【0171】

幾つかの実施例において、上述したソースコードファイルにおいて、テンプレート関数を明示的にインスタンス化することができる。

【0172】

例えば、ヘッダファイルにおいて、テンプレート関数記述情報は、例えば以下の通りであってもよい。

【0173】

【数7】

40

```
template <typename T, typename Context>
```

```
void ScaleKernel(const Context& dev_ctx,
                 const DenseTensor& x,
                 const Scalar& scale,
                 float bias,
                 bool bias_after_scale,
                 DenseTensor* out)。
```

10

【 0 1 7 4 】

また例えば、ソースコードファイルにおいて、明示的にインスタンス化された関数 `ScaleKernel float, CPUContext` は、例えば以下の通りであってもよい。

【 0 1 7 5 】

【 数 8 】

```
template ScaleKernel<float, CPUContext>(const Context& CPUContext,
                                        const DenseTensor& x,
                                        const Scalar& scale,
                                        float bias,
                                        bool bias_after_scale,
                                        DenseTensor* out)。
```

20

【 0 1 7 6 】

データタイプが 8 種類であり、且つ機器種別が 2 種である場合、明示的にインスタンス化された関数に対応するコードは、16 組を含んでもよい。

【 0 1 7 7 】

他の実施例において、複数の完全なテンプレートパラメータに基づいて、以下のような複数のインスタンス化された宣言文を生成することができる。

【 0 1 7 8 】

【 数 9 】

```
Template decltype(
    カーネル関数名<データタイプ, 機器コンテキストタイプ>。)
```

40

【 0 1 7 9 】

ここで、`decltype` は、関数タイプを自動的に解析する方法である。このように、コンパイル効率を向上するとともに、書き方の間接性も向上し、この明示的なインスタンス化プロセスを自動的に完了する。

【 0 1 8 0 】

幾つかの実施例において、グローバル演算子テーブルは、演算子カーネル工場クラスと、演算子カーネル名称クラス及びカーネルキー値クラスと、演算子カーネルクラスとを含む。例えば、演算子カーネルクラスは、演算子カーネル関数のポインタ、演算子カーネル関

50

数入力パラメータの記述情報、及び演算子カーネル関数出力パラメータの記述情報の少なくとも1つを含む。

【0181】

以下、図6を参照してグローバル演算子テーブルを詳細に説明する。

【0182】

図6は、本開示の一実施例によるグローバル演算子テーブルの例示的模式図である。

【0183】

図6に示すように、グローバル演算子テーブル600のデータ構造は、第1級の演算子カーネル工場クラス610、第2級の演算子カーネル名称クラス620、カーネルキー値クラス630、及び演算子カーネルクラス640を含んでもよい。

10

【0184】

図6の例示において、演算子カーネルクラス640は、第3級の演算子カーネル関数のポインタ641及び演算子カーネル関数パラメータの記述情報642を含んでもよい。演算子カーネル関数パラメータの記述情報642は、例えば演算子カーネル関数入力パラメータの記述情報及び演算子カーネル関数出力関数の記述情報を含んでもよい。

【0185】

幾つかの実施例において、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得ることは、複数の完全なテンプレートパラメータを明示的にインスタンス化することで、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得ることを含む。

20

【0186】

幾つかの実施例において、複数の完全なテンプレートパラメータの各々は、テンプレート関数記述情報及びテンプレート関数情報を含み、複数の完全なテンプレートパラメータを明示的にインスタンス化することで、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得ることは、テンプレート関数記述情報をヘッダファイルとしてインスタンス化することと、テンプレート関数情報をソースコードファイルとしてインスタンス化することを含む。

【0187】

図7は、本開示の一実施例による深層学習フレームワークの演算子処理装置のブロック図である。

30

【0188】

図7に示すように、該装置700は、取得モジュール710、解析モジュール720、及び処理モジュール730を含んでもよい。

【0189】

取得モジュール710は、処理対象演算子を取得するために用いられる。例えば、処理対象演算子は、演算子カーネル関数、及び深層学習フレームワークに依存しないテンプレートパラメータを含む。

【0190】

解析モジュール720は、処理対象演算子に対する入力情報を受信したことに応じて、入力情報を用いてテンプレートパラメータを解析して、深層学習フレームワークに依存する複数の完全なテンプレートパラメータを得るために用いられる。

40

【0191】

処理モジュール730は、複数の完全なテンプレートパラメータに基づいて、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得るために用いられる。

【0192】

幾つかの実施例において、解析モジュールは、マクロを用いて可変パラメータリストを作成するための作成サブモジュールと、可変パラメータリストに基づいてテンプレートパラメータを反復的に解析するための解析サブモジュールとを含み、ここで、入力情報をマクロのマクロパラメータとする。

50

【0193】

幾つかの実施例において、可変パラメータリストは、処理対象演算子の名称と、演算子カーネル関数の名称と、演算子カーネル関数のデータレイアウトと、深層学習フレームワークを実装するための機器タイプと、深層学習フレームワークを実装するためのデータタイプとを含む。

【0194】

幾つかの実施例において、解析サブモジュールは、マクロに含まれるパラメータ数を特定するための第1の特定ユニットと、パラメータ数に応じて、反復終了位置を特定するための第2の特定ユニットと、可変パラメータリストを用いて、反復終了位置に応じて入力情報を反復に処理するための処理ユニットとを含む。

10

【0195】

幾つかの実施例において、解析サブモジュールは、可変パラメータリストを用いてテンプレートパラメータを反復に解析して、複数の内部テンプレートパラメータを得るための解析ユニットと、複数の内部テンプレートパラメータをスプライシングして、複数の完全なテンプレートパラメータを得るためのスプライシングユニットとをさらに含む。

【0196】

幾つかの実施例において、複数の内部テンプレートパラメータは複数の文字列を含み、スプライシングユニットは、入力情報に基づいて、複数の文字列をスプライシングして、処理対象演算子に対する複数の完全なテンプレートパラメータを得るためのスプライシングサブユニットを含む。

20

【0197】

幾つかの実施例において、処理モジュールは、複数の完全なテンプレートパラメータを明示的にインスタンス化することで、演算子カーネル関数を処理して、深層学習フレームワーク用の利用可能な演算子を得るための明示的なインスタンス化サブモジュールを含む。

【0198】

幾つかの実施例において、複数の完全なテンプレートパラメータの各々は、テンプレート関数記述情報及びテンプレート関数情報を含み、明示的なインスタンス化サブモジュールは、テンプレート関数記述情報をヘッダファイルとしてインスタンス化するための第1のインスタンス化ユニットと、テンプレート関数情報をソースコードファイルとしてインスタンス化するための第2のインスタンス化ユニットとを含む。

30

【0199】

幾つかの実施例において、処理モジュールは、複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、該完全なテンプレートパラメータの関数パラメータタイプを特定するための第1の特定サブモジュールと、関数パラメータタイプに基づいて、完全なテンプレートパラメータの入力情報を対応するパラメータ情報に変換するための変換サブモジュールと、対応するパラメータ情報を記録するための記録サブモジュールとを含む。

【0200】

幾つかの実施例において、装置700は、複数の完全なテンプレートパラメータに基づいて、深層学習フレームワークに基づく登録対象演算子を作成するための作成モジュールと、登録対象演算子を深層学習フレームワーク内部のグローバル演算子テーブルに登録するための登録モジュールとをさらに含む。

40

【0201】

幾つかの実施例において、登録対象演算子は、登録対象演算子記述情報及び登録対象演算子カーネル関数を含み、登録対象演算子カーネル関数に正規化された入力情報及び正規化された関数ポインタが含まれ、作成モジュールは、複数の完全なテンプレートパラメータのうちの各完全なテンプレートパラメータに対して、完全なテンプレートパラメータに対応する、正規化された形式を有する静的関数を含む構造体を特定するための第2の特定サブモジュールと、静的関数の入力情報を正規化された入力情報とするための第1の正規化サブモジュールと、静的関数の関数ポインタを正規化された関数ポインタとするための第

50

2の正規化モジュールとを含む。

【0202】

幾つかの実施例において、静的関数の入力情報は入力情報リストを構成し、入力情報リストは入力テンソルリストを含み、構造体は少なくとも1つの特化サブ構造体を含み、少なくとも1つの特化サブ構造体の各々は対応するデータタイプを含み、データタイプはテンソル及び非テンソルを含む。

【0203】

幾つかの実施例において、グローバル演算子テーブルは、演算子カーネル工場クラスと、演算子カーネル名称クラス及びカーネルキー値クラスと、演算子カーネルクラスとを含む。

【0204】

幾つかの実施例において、演算子カーネルクラスは、演算子カーネル関数のポインタと、演算子カーネル関数入力パラメータの記述情報と、演算子カーネル関数出力パラメータの記述情報との少なくとも1つを含む。

【0205】

幾つかの実施例において、深層学習フレームワークに依存しないテンプレートパラメータは、処理対象演算子を配布するために用いられ、深層学習フレームワークに依存しないテンプレートパラメータは、深層学習フレームワークの実装機器および深層学習フレームワークのデータタイプのいずれにも依存しないテンプレートパラメータを含む。

【0206】

幾つかの実施例において、装置700は、利用可能な演算子に対する呼び出しコマンドに
20 応答して、利用可能な演算子を呼び出すための呼び出しモジュールと、呼び出された利用可能な演算子をコンパイルすることで、利用可能な演算子に対応する関数機能を実行するための実行モジュールとをさらに含む。

【0207】

本開示の技術案において、係れたユーザ個人情報の収集、記憶、使用、加工、伝送、提供、開示及び応用等の処理は、関連法律や法規の規定に合致しており、必要なセキュリティ対策を取っており、公序良俗に反していない。本開示の技術案において、ユーザの個人情報を取得するか又は収集する前に、いずれもユーザの許可又は同意を得た。

【0208】

本開示の実施例によれば、本開示は、電子機器、読取可能な記憶媒体及びコンピュータプログラムをさらに提供する。
30

【0209】

図8は、本開示の実施例を実施するための例示的電子機器800の模式的ブロック図を示している。電子機器は、様々な形式のデジタルコンピュータを示すことを目的とし、例えば、ラップトップコンピュータ、デスクトップコンピュータ、ワークステーション、パーソナルデジタルアシスタント、サーバ、ブレードサーバ、大型コンピュータ及び他の適切なコンピュータである。電子機器は、さらに様々な形式の移動装置を示してもよく、例えば、パーソナルデジタルアシスタント、携帯電話、スマートフォン、ウェアラブル機器及び他の類似の演算装置である。本明細書に示された部材、それらの接続及び関係、並びにそれらの機能は、例示に過ぎず、本明細書に記載された及び/又は要求された本開示の実現を限定しない。
40

【0210】

図8に示すように、電子機器800は、計算手段801を含み、計算手段801は、リードオンリーメモリ(ROM)802に記憶されたコンピュータプログラム又は記憶手段808からランダムアクセスメモリ(RAM)803にロードされたコンピュータプログラムに基づいて、様々な適切な動作及び処理を実行してもよい。RAM803には、さらに電子機器800の操作に必要な様々なプログラム及びデータを記憶してもよい。計算手段801、ROM802、及びRAM803は、バス804を介して相互に接続される。入出力(I/O)インターフェース805も、バス804に接続される。

【0211】

10

20

30

40

50

電子機器 800 における複数の部品は、I/O インターフェース 805 に接続され、例えばキーボード、マウス等の入力手段 806 と、例えば様々な種類のディスプレイ、スピーカ等の出力手段 807 と、例えば磁気ディスク、光ディスク等の記憶手段 808 と、例えばネットワークカード、モデム、無線通信トランシーバ等の通信手段 809 とを含む。通信手段 809 は、電子機器 800 がインターネット等のコンピュータネットワーク及び/又は各種の電気ネットワークを介して他の機器と情報・データをやり取りすることを可能にする。

【0212】

計算手段 801 は、処理及び演算能力を有する各種の汎用及び/又は専用の処理モジュールであってもよい。計算手段 801 の幾つかの例として、中央処理ユニット (CPU)、GPU (Graphics Processing Unit)、各種専用の人工知能 (AI) 演算チップ、各種機械学習モデルアルゴリズムをランニングする演算ユニット、DSP (Digital Signal Processor)、並びに任意の適切なプロセッサ、コントローラ、マイクロコントローラ等が挙げられるが、これらに限定されない。計算手段 801 は、前文で記載された各方法及び処理、例えば深層学習フレームワークの演算子処理方法を実行する。例えば、幾つかの実施例において、深層学習フレームワークの演算子処理方法は、例えば記憶手段 808 のような機械可読媒体に有形的に含まれるコンピュータソフトウェアプログラムとして実現されてもよい。いくつかの実施例において、コンピュータプログラムの一部又は全部は、ROM 802 及び/又は通信手段 809 を介して電子機器 800 にロード及び/又はインストールされてもよい。コンピュータプログラムが RAM 803 にロードされて計算手段 801 により実行される場合、前文に記載の深層学習フレームワークの演算子処理方法の 1 つ又は複数のステップを実行してもよい。代替的に、他の実施例において、計算手段 801 は、他の任意の適切な方式 (例えば、ファームウェアを介する) により深層学習フレームワークの演算子処理方法を実行するように構成されてもよい。

【0213】

本明細書で以上に説明されたシステム及び技術の様々な実施形態は、デジタル電子回路システム、集積回路システム、フィールドプログラマブルゲートアレイ (FPGA)、特定用途向け集積回路 (ASIC)、特定用途向け標準製品 (ASSP)、システムオンチップ (SOC)、コンプレックスプログラマブルロジックデバイス (CPLD)、コンピュータハードウェア、ファームウェア、ソフトウェア、及び/又はそれらの組み合わせにおいて実現されてもよい。これらの様々な実施形態は、1 つ又は複数のコンピュータプログラムにおいて実施され、該 1 つ又は複数のコンピュータプログラムは、少なくとも 1 つのプログラマブルプロセッサを含むプログラマブルシステムで実行され及び/又は解釈されることが可能であり、該プログラマブルプロセッサは、専用又は汎用のプログラマブルプロセッサであってもよく、記憶システム、少なくとも 1 つの入力装置、及び少なくとも 1 つの出力装置からデータ及び命令を受信し、かつデータ及び命令を該記憶システム、該少なくとも 1 つの入力装置、及び該少なくとも 1 つの出力装置に伝送することができることを含んでもよい。

【0214】

本開示の方法を実施するためのプログラムコードは、1 つ又は複数のプログラミング言語の任意の組み合わせで作成されてもよい。これらのプログラムコードは、汎用コンピュータ、専用コンピュータ又は他のプログラマブルデータ処理装置のプロセッサ又はコントローラに提供されてもよく、それによって、プログラムコードがプロセッサ又はコントローラにより実行される時に、フローチャート及び/又はブロック図に規定された機能・操作が実施される。プログラムコードは、機器に完全に実行されてもよく、部分的に機器で実行されてもよく、独立したソフトウェアパッケージとして部分的に機器で実行され、かつ部分的に遠隔機器で実行されるか又は完全に遠隔機器又はサーバで実行されてもよい。

【0215】

本開示のコンテキストにおいて、機械可読媒体は、有形の媒体であってもよく、命令実行

10

20

30

40

50

システム、装置又は電子機器に使用され、又は命令実行システム、装置又は電子機器と組み合わせて使用されるプログラムを含んで又は記憶してもよい。機械可読媒体は、機械可読信号媒体又は機械可読記憶媒体であってもよい。機械可読媒体は、電子の、磁氣的、光学的、電磁的、赤外線、又は半導体システム、装置又は電子機器、又は上記内容の任意の適切な組み合わせを含んでもよいが、それらに限定されない。機械可読記憶媒体のより具体的な例としては、1つ以上の線による電氣的接続、携帯式コンピュータディスク、ハードディスク、ランダムアクセスメモリ（RAM）、読み出し専用メモリ（ROM）、消去可能なプログラマブルリードオンリーメモリ（EPROM又はフラッシュメモリ）、光ファイバ、コンパクトディスクリードオンリーメモリ（CD-ROM）、光学記憶装置、磁気記憶装置、又は上記内容の任意の適切な組み合わせを含む。

10

【0216】

ユーザとの対話を提供するために、コンピュータにここで説明されたシステム及び技術を実施させてもよく、該コンピュータは、ユーザに情報を表示するための表示装置（例えば、CRT（陰極線管）又はLCD（液晶ディスプレイ）モニター）と、キーボード及びポインティングデバイス（例えば、マウス又はトラックボール）とを備え、ユーザは、該キーボード及び該ポインティングデバイスを介して入力をコンピュータに提供することができる。他の種類の装置は、さらにユーザとの対話を提供してもよく、例えば、ユーザに提供されたフィードバックは、いかなる形式のセンシングフィードバック（例えば、視覚フィードバック、聴覚フィードバック、又は触覚フィードバック）であってもよく、かついかなる形式（音声入力、語音入力又は、触覚入力を含む）でユーザからの入力を受信してもよい。

20

【0217】

ここで説明されたシステム及び技術は、バックグラウンド部品を含むコンピューティングシステム（例えば、データサーバとする）、又はミドルウェア部品を含むコンピューティングシステム（例えば、アプリケーションサーバ）、又はフロントエンド部品を含むコンピューティングシステム（例えば、グラフィカルユーザインタフェース又はウェブブラウザを有するユーザコンピュータ、ユーザが該グラフィカルユーザインタフェース又は該ネットワークブラウザを介してここで説明されたシステム及び技術の実施形態と対話することができる）、又はこのようなバックグラウンド部品、ミドルウェア部品、又はフロントエンド部品のいずれかの組み合わせを含むコンピューティングシステムに実施されることが可能である。任意の形式又は媒体のデジタルデータ通信（例えば、通信ネットワーク）によりシステムの部品を互いに接続することができる。通信ネットワークの例としては、ローカルエリアネットワーク（LAN）、ワイドエリアネットワーク（WAN）及びインターネットを例示的に含む。

30

【0218】

コンピュータシステムは、クライアント及びサーバを含んでもよい。クライアントとサーバ同士は、一般的に離れており、通常、通信ネットワークを介して対話する。クライアントとサーバとの関係は、該当するコンピュータ上でランニングし、クライアント-サーバの関係性を有するコンピュータプログラムによって生成される。サーバは、クラウドサーバであってもよく、分散型システムのサーバであってもよく、又はブロックチェーンを組合せたサーバであってもよい。

40

【0219】

理解されるべきこととして、以上に示された様々な形式のフローを使用してもよく、操作を改めてソーティングしたり、追加したり又は削除してもよい。例えば、本開示に記載の各操作は、並列に実行されたり、順次に実行されたり、又は異なる順序で実行されてもよく、本開示に開示された技術案の所望の結果を実現することができれば、本明細書はここで限定されない。

【0220】

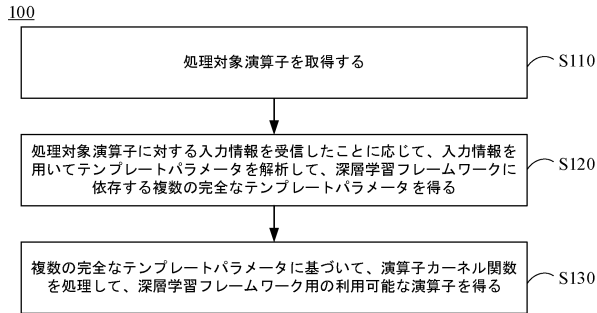
上記具体的な実施形態は、本開示の保護範囲を限定するものではない。当業者であれば、設計要件及び他の要因に応じて、様々な修正、組み合わせ、サブコンビネーション及び代

50

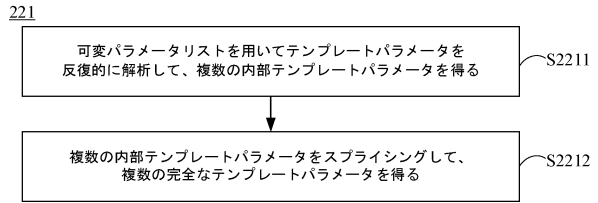
替を行うことが可能であると理解すべきである。本開示の精神と原則内で行われる任意の修正、均等置換及び改良などは、いずれも本開示の保護範囲内に含まれるべきである。

【 図面 】

【 図 1 】

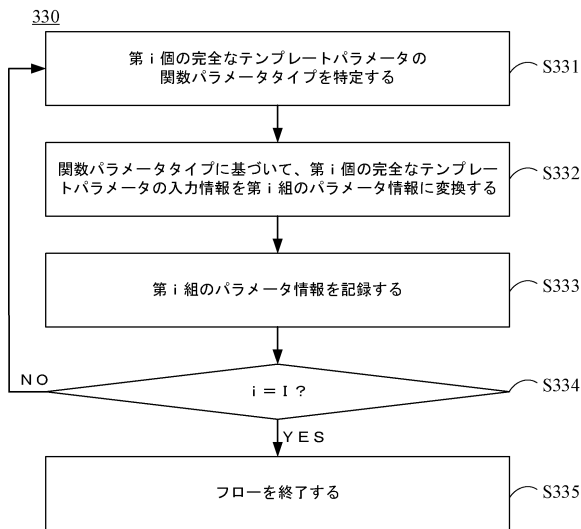


【 図 2 】

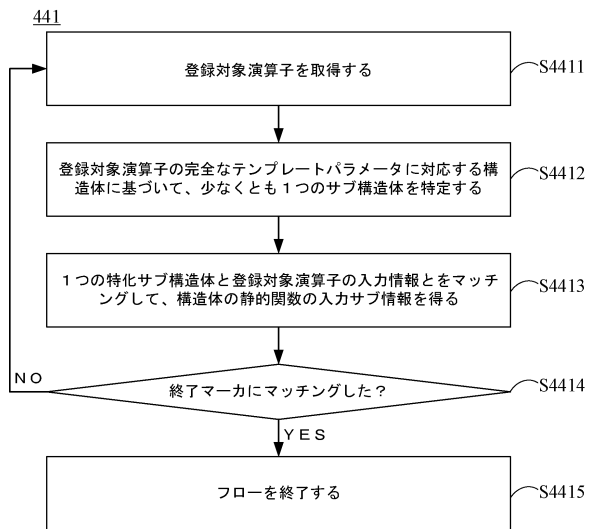


10

【 図 3 】



【 図 4 】



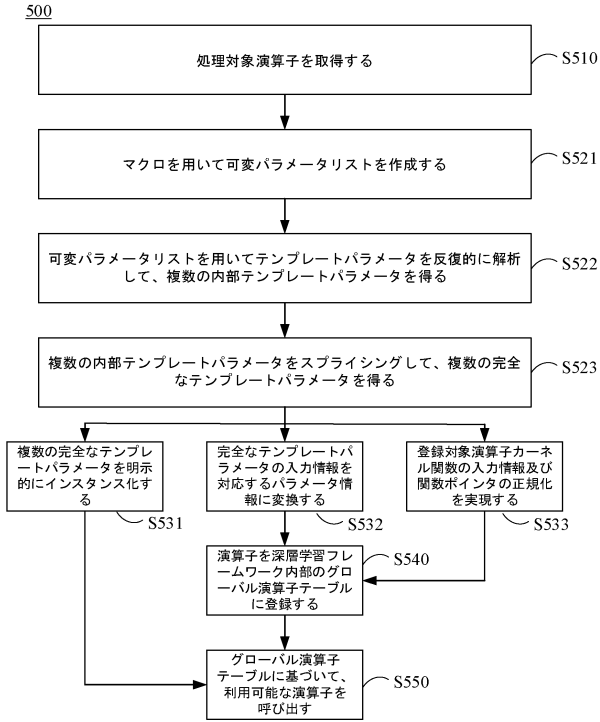
20

30

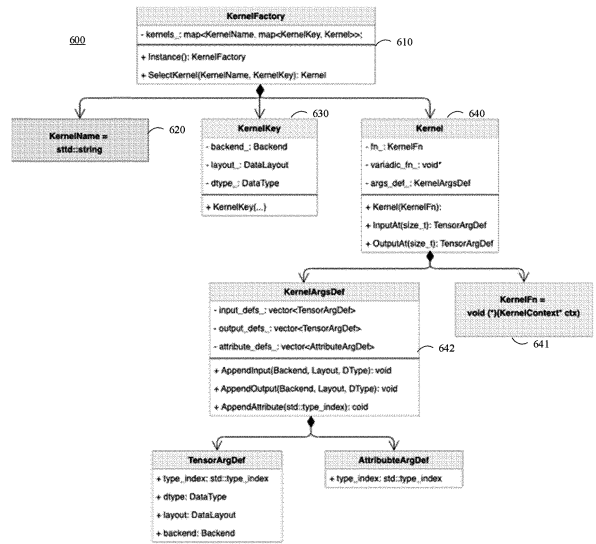
40

50

【図5】



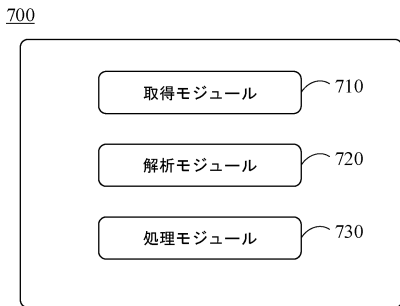
【図6】



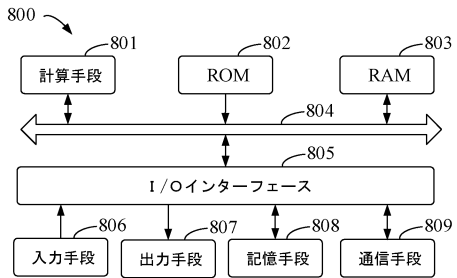
10

20

【図7】



【図8】



30

40

50

フロントページの続き

- 85, China
- (74)代理人 110000914
弁理士法人WisePlus
- (72)発明者 チェン, ウェイハン
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 ワン, ハイフォン
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 ジャン, ユンフェイ
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 ユェン, リーシオン
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 チェン, ティエンユー
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 リウ, ホンユー
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 フー, シャオグアン
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 ユー, ディエンハイ
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- (72)発明者 マー, イェンジュン
中華人民共和国 ペキン 100085, ハイディアン ディストリクト, シャンディ テンス ス
トリート, 10番, バイドウ キャンパス 2階
- 審査官 宮司 卓佳
- (56)参考文献 中国特許出願公開第113342346 (CN, A)
米国特許出願公開第2019/0392296 (US, A1)
国際公開第2020/263421 (WO, A1)
中国特許出願公開第113705798 (CN, A)
中国特許出願公開第111310934 (CN, A)
米国特許第10956132 (US, B1)
中国特許出願公開第113031966 (CN, A)
- (58)調査した分野 (Int.Cl., DB名)
G06N 3/10
G06F 8/40
G06F 8/30