

(19)日本国特許庁(JP)

## (12)特許公報(B2)

(11)特許番号  
特許第7401513号  
(P7401513)

(45)発行日 令和5年12月19日(2023.12.19)

(24)登録日 令和5年12月11日(2023.12.11)

(51)国際特許分類 F I  
G 0 6 F 17/16 (2006.01) G 0 6 F 17/16 P

請求項の数 23 外国語出願 (全31頁)

(21)出願番号	特願2021-207147(P2021-207147)	(73)特許権者	502208397 グーグル エルエルシー Google LLC アメリカ合衆国 カリフォルニア州 94043 マウンテン ビュー アンフィシアター パークウェイ 1600 1600 Amphitheatre Parkway 94043 Mountain View, CA U.S.A.
(22)出願日	令和3年12月21日(2021.12.21)	(74)代理人	110001195 弁理士法人深見特許事務所
(65)公開番号	特開2022-181161(P2022-181161A)	(72)発明者	レイナー・アルウィン・ポーブ アメリカ合衆国、94043 カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、1600 最終頁に続く
(43)公開日	令和4年12月7日(2022.12.7)		
審査請求日	令和4年4月6日(2022.4.6)		
(31)優先権主張番号	17/329,259		
(32)優先日	令和3年5月25日(2021.5.25)		
(33)優先権主張国・地域又は機関	米国(US)		

(54)【発明の名称】 ハードウェアにおけるスパース行列乗算

## (57)【特許請求の範囲】

## 【請求項1】

システム入力行列とシステム入力ベクトルとの間の行列乗算を実行するように構成されたシステムであって、複数の乗算回路を含む回路の集合体を備えるスパースシャードを備え、前記スパースシャードは、

前記システム入力行列の複数の部分行列のいずれかであって予め定められた最大非ゼロ閾値以下の数の非ゼロ値を含むシャード入力行列を受けると構成され、前記シャード入力行列は、前記スパースシャードの最大行列入力サイズを指定する予め定められた次元閾値以下の次元を有し、前記スパースシャードはさらに、

前記システム入力ベクトルの複数の部分ベクトルのいずれかであって複数のベクトル値を含むシャード入力ベクトルを受け、

前記乗算回路の各々について、前記シャード入力行列のそれぞれの非ゼロ値を受け、前記複数の乗算回路によって、ベクトル値に前記シャード入力行列の前記それぞれの非ゼロ値を乗算した1つ以上の積を生成し、

前記スパースシャードへの出力として、前記1つ以上の積を用いて、前記シャード入力ベクトルを前記シャード入力行列に適用した積であるシャード出力ベクトルを生成するように構成されている、システム。

## 【請求項2】

前記シャード出力ベクトルの長さは1よりも大きい、請求項1に記載のシステム。

## 【請求項3】

10

20

前記スパスシャードは複数のスパスシャードのうちの1つであり、前記複数のスパスシャードは、

前記システム入力行列の部分行列である複数のシャード入力行列を受け、

前記システム入力ベクトルの部分ベクトルである複数のシャード入力ベクトルを受け、

前記複数のスパスシャードによって、前記システム入力ベクトルを前記システム入力行列に適用した積を表すシステム出力ベクトルを生成するように構成されている、請求項1に記載のシステム。

【請求項4】

前記複数のスパスシャードはシストリックアレイとして配列され、前記シストリックアレイは、前記シストリックアレイの列次元に沿ったスパスシャードの1つ以上のグループを含み、

10

前記システム出力ベクトルを生成するために、前記複数のスパスシャードはさらに、

前記シストリックアレイの前記列次元に沿ったグループごとに、前記グループ内の各スパスシャードのそれぞれのシャード出力ベクトルを合算してそれぞれの列出力ベクトルを生成し、

各グループの前記それぞれの列出力ベクトルを連結して前記システム出力ベクトルを生成するように構成されている、請求項3に記載のシステム。

【請求項5】

各乗算回路は、前記スパスシャードについての前記シャード入力行列のそれぞれからの前記それぞれの非ゼロ値を含むそれぞれのレジスタに結合される、請求項1に記載のシステム。

20

【請求項6】

前記複数の乗算回路における乗算回路の数は、前記予め定められた最大非ゼロ閾値と等しい、請求項5に記載のシステム。

【請求項7】

前記スパスシャードはクロスバー回路をさらに含み、前記スパスシャードはさらに、前記クロスバー回路によって、前記シャード入力ベクトルの前記複数のベクトル値を受け、

前記複数の乗算回路の各々への入力として、前記クロスバー回路によって、前記複数のベクトル値のうちの1つのベクトル値を送るように構成されている、請求項5に記載のシステム。

30

【請求項8】

前記スパスシャードはさらに、

前記シャード入力行列の同じ列の非ゼロ値を、前記複数の乗算回路のうちの隣接する乗算回路のレジスタにロードするように構成されている、請求項7に記載のシステム。

【請求項9】

前記スパスシャードはさらに、前記シャード入力行列の各列に沿った非ゼロ値の位置を少なくとも指定する1つ以上の制御値を受けするように構成されており、

前記スパスシャードの前記クロスバー回路はさらに、

前記1つ以上の制御値を受け、

40

前記1つ以上の制御値に従って、前記シャード入力行列の同じ列に沿った非ゼロ値が乗算されるベクトル値を、隣接する乗算回路に送るように構成されている、請求項8に記載のシステム。

【請求項10】

前記スパスシャードは複数の加算回路をさらに含み、

前記スパスシャードは1つ以上のセグメントマーカをさらに含み、各セグメントマーカは、前記セグメントマーカにロードされたそれぞれの制御値の値に基づいて、前記複数の加算回路のうちのそれぞれの加算回路への入力をゲート開閉するように構成されており、

前記スパスシャードはさらに、

前記1つ以上の制御値の少なくとも一部を前記1つ以上のセグメントマーカにロード

50

するように構成されており、前記シャード入力行列の 1 列目の非ゼロ値についての加算回路は、前記 1 列目とは異なる前記シャード入力行列の 2 列目の非ゼロ値を含む隣接する加算回路の入力を受けることからゲート開閉されており、前記スパスシャードはさらに、

前記複数の加算回路によって、前記 1 つ以上の積の 1 つ以上の合計を生成するように構成されており、前記 1 つ以上の合計の各々は、前記シャード入力行列の列の 1 つ以上の非ゼロ値に前記シャード入力ベクトルの 1 つ以上のそれぞれの値を乗算したそれぞれのセグメント化合計である、請求項 9 に記載のシステム。

【請求項 11】

前記複数の加算回路は並列セグメント化合計回路を形成し、1 つ以上の前記セグメント化合計の各々は、セグメントマーカによってゲート開閉されていない隣接する加算回路への出力の合計である、請求項 10 に記載のシステム。

10

【請求項 12】

前記クロスバー回路は第 1 のクロスバー回路であり、

前記スパスシャードは第 2 のクロスバー回路をさらに含み、前記第 2 のクロスバー回路は、

1 つ以上の前記セグメント化合計を受け、

前記 1 つ以上の制御値に従って 1 つ以上の前記セグメント化合計を配列して、前記スパスシャードについて前記シャード出力ベクトルのそれぞれを生成するように構成されている、請求項 10 に記載のシステム。

【請求項 13】

20

前記第 2 のクロスバー回路はベネシュネットワークを形成し、前記シャード入力行列は正方行列である、請求項 12 に記載のシステム。

【請求項 14】

前記複数の部分行列は、第 1 の部分行列および第 2 の部分行列を含み、

前記第 1 の部分行列の要素数は、前記第 2 の部分行列の要素数とは異なる、請求項 1 ~ 1.3 のいずれか 1 項に記載のシステム。

【請求項 15】

システム入力行列とシステム入力ベクトルとの間の行列乗算を実行する方法であって、

複数の乗算回路を含む回路の集合体を備えるスパスシャードが、前記システム入力行列の複数の部分行列のいずれかであって予め定められた最大非ゼロ閾値以下の数の非ゼロ値を含むシャード入力行列と、前記システム入力ベクトルの複数の部分ベクトルのいずれかであって複数のベクトル値を含むシャード入力ベクトルとを受けことを備え、前記シャード入力行列は、前記スパスシャードの最大行列入力サイズを指定する予め定められた次元閾値以下の次元を有し、前記方法はさらに、

30

前記乗算回路の各々について、前記シャード入力行列の前記非ゼロ値のそれぞれを受けると、

前記スパスシャードの前記複数の乗算回路が、それぞれのベクトル値に前記シャード入力行列の前記非ゼロ値のそれぞれを乗算した 1 つ以上の積を生成することと、

前記スパスシャードが、前記スパスシャードへの出力として、前記 1 つ以上の積を用いて、前記シャード入力ベクトルを前記シャード入力行列に適用した積であるシャード出力ベクトルを生成することとを備える、方法。

40

【請求項 16】

前記スパスシャードが前記シャード出力ベクトルを生成することは、1 よりも大きい長さの前記シャード出力ベクトルを生成することを含む、請求項 15 に記載の方法。

【請求項 17】

前記スパスシャードは複数のスパスシャードのうちの 1 つであり、前記複数のスパスシャードは、

前記システム入力行列の部分行列である複数のシャード入力行列を受け、

前記システム入力ベクトルの部分ベクトルである複数のシャード入力ベクトルを受け、

前記複数のスパスシャードによって、前記システム入力ベクトルを前記システム入力

50

行列に適用した積を表すシステム出力ベクトルを生成するように構成されている、請求項 1.5 に記載の方法。

【請求項 18】

前記複数のスパースシャードはシストリックアレイとして配列され、前記シストリックアレイは、前記シストリックアレイの列次元に沿ったスパースシャードの 1 つ以上のグループを含み、

前記システム出力ベクトルを生成することは、

前記シストリックアレイの前記列次元に沿ったグループごとに、前記グループ内の各スパースシャードのそれぞれのシャード出力ベクトルを合算してそれぞれの列出力ベクトルを生成することと、

各グループの前記それぞれの列出力ベクトルを連結して前記システム出力ベクトルを生成することを含む、請求項 1.7 に記載の方法。

【請求項 19】

各乗算回路は、前記スパースシャードについての前記シャード入力行列のそれぞれからのそれぞれの非ゼロ値を含むそれぞれのレジスタに結合される、請求項 1.5 に記載の方法。

【請求項 20】

前記複数の乗算回路における乗算回路の数は、前記予め定められた最大非ゼロ閾値と等しい、請求項 1.9 に記載の方法。

【請求項 21】

前記複数の部分行列は、第 1 の部分行列および第 2 の部分行列を含み、

前記第 1 の部分行列の要素数は、前記第 2 の部分行列の要素数とは異なる、請求項 1.5 ~ 2.0 のいずれか 1 項に記載の方法。

【請求項 22】

システム入力行列とシステム入力ベクトルとの間の行列乗算を実行するように構成されたシステムによって実行されると前記システムに動作を実行させる命令を格納した 1 つ以上のコンピュータプログラムであって、前記動作は、

複数の乗算回路を含む回路の集合体を備えるスパースシャードによって、前記システム入力行列の複数の部分行列のいずれかであって予め定められた最大非ゼロ閾値以下の数の非ゼロ値を含むシャード入力行列と、前記システム入力ベクトルの複数の部分ベクトルのいずれかであって複数のベクトル値を含むシャード入力ベクトルとを受けるとを備え、前記シャード入力行列は、前記スパースシャードの最大行列入力サイズを指定する予め定められた次元閾値以下の次元を有し、前記動作はさらに、

前記乗算回路の各々について、前記シャード入力行列のそれぞれの非ゼロ値を受けるとと、

前記スパースシャードの前記複数の乗算回路によって、それぞれのベクトル値に前記それぞれの非ゼロ値を乗算した 1 つ以上の積を生成することと、

前記スパースシャードへの出力として、前記 1 つ以上の積を用いて、前記シャード入力ベクトルを前記シャード入力行列に適用した積であるシャード出力ベクトルを生成することとを備える、1 つ以上のコンピュータプログラム。

【請求項 23】

前記複数の部分行列は、第 1 の部分行列および第 2 の部分行列を含み、

前記第 1 の部分行列の要素数は、前記第 2 の部分行列の要素数とは異なる、請求項 2.2 に記載の 1 つ以上のコンピュータプログラム。

【発明の詳細な説明】

【背景技術】

【0001】

背景

スパース行列は、行列の要素として非ゼロ値よりもゼロ値の方が割合が高い行列である。異なるスパース行列は、ゼロ値と非ゼロ値との割合に基づいてさまざまな程度のスパース性を有し得る。非ゼロ値よりもゼロ値の方が割合が高い行列は、非ゼロ値よりもゼロ値

10

20

30

40

50

の方が割合が低い行列よりもスパース性が高いと言われる。

【 0 0 0 2 】

ニューラルネットワークは、受けた入力に対する出力を予測するための1つ以上の非線形演算層を含む機械学習モデルである。入力層および出力層に加えて、1つ以上の隠れ層を含むニューラルネットワークもある。各隠れ層の出力は、ニューラルネットワークの別の隠れ層または出力層に入力され得る。ニューラルネットワークの各層は、当該層の1つ以上のモデルパラメータの値に従って、受けた入力からそれぞれの出力を生成することができる。モデルパラメータは、ニューラルネットワークが正確な出力を生成するように訓練アルゴリズムを通して求められる重みまたはバイアスであってもよい。ニューラルネットワークの層のモデルパラメータ値は、行列またはテンソルの要素として表すことができる。

10

【 発明の概要 】

【 課題を解決するための手段 】

【 0 0 0 3 】

簡単な概要

本開示の局面は、ハードウェアにおけるスパース行列密ベクトル乗算に関する。

【 0 0 0 4 】

本開示の一局面は、複数の乗算回路を含むスパースシャードを含むシステムを提供し、上記スパースシャードは、予め定められた最大非ゼロ閾値以下の数の非ゼロ値を含むシャード入力行列を受け、複数のベクトル値を含むシャード入力ベクトルを受け、上記乗算回路の各々について、上記シャード入力行列のそれぞれの非ゼロ値を受け、上記複数の乗算回路によって、ベクトル値に上記シャード入力行列の上記それぞれの非ゼロ値を乗算した1つ以上の積を生成し、上記スパースシャードへの出力として、上記1つ以上の積を用いて、上記シャード入力ベクトルを上記シャード入力行列に適用した積であるシャード出力ベクトルを生成するように構成されている。

20

【 0 0 0 5 】

本開示の別の局面は、複数のスパースシャードを含むシステムによって実行されると上記システムに動作を実行させる命令を格納した1つ以上の非一時的なコンピュータ読取可能記憶媒体を提供し、上記動作は、複数の乗算回路を含むスパースシャードによって、予め定められた最大非ゼロ閾値以下の数の非ゼロ値を含むシャード入力行列と、複数のベクトル値を含むシャード入力ベクトルとを受けると、上記乗算回路の各々について、上記シャード入力行列のそれぞれの非ゼロ値を受けると、上記スパースシャードの上記複数の乗算回路によって、それぞれのベクトル値に上記それぞれの非ゼロ値を乗算した1つ以上の積を生成することと、上記スパースシャードへの出力として、上記1つ以上の積を用いて、上記シャード入力ベクトルを上記シャード入力行列に適用した積であるシャード出力ベクトルを生成することを含む。

30

【 0 0 0 6 】

本開示の別の局面は、方法を提供し、上記方法は、複数の乗算回路を含むスパースシャードが、予め定められた最大非ゼロ閾値以下の数の非ゼロ値を含むシャード入力行列と、複数のベクトル値を含むシャード入力ベクトルとを受けると、上記乗算回路の各々について、上記シャード入力行列の上記非ゼロ値のそれぞれを受けると、上記スパースシャードの上記複数の乗算回路が、それぞれのベクトル値に上記シャード入力行列の上記非ゼロ値のそれぞれを乗算した1つ以上の積を生成することと、上記スパースシャードが、上記スパースシャードへの出力として、上記1つ以上の積を用いて、上記シャード入力ベクトルを上記シャード入力行列に適用した積であるシャード出力ベクトルを生成することを含む。

40

【 0 0 0 7 】

上述のおよび他の局面の各々は、任意に以下の特徴のうちの1つ以上を単独でまたは組み合わせで含み得る。一実現例は、以下の特徴のすべてを組み合わせで含み得る。

【 0 0 0 8 】

50

上記シャード出力ベクトルの長さは1よりも大きい。

上記スパースシャードは複数のスパースシャードのうちの1つであり、上記複数のスパースシャードは、システム入力行列の部分行列である複数のシャード入力行列を受け、システム入力ベクトルの部分ベクトルである複数のシャード入力ベクトルを受け、上記複数のスパースシャードによって、上記システム入力ベクトルを上記システム入力行列に適用した積を表すシステム出力ベクトルを生成するように構成されている。

【0009】

上記複数のスパースシャードはシストリックアレイとして配列され、上記シストリックアレイは、上記シストリックアレイの列次元に沿ったスパースシャードの1つ以上のグループを含み、上記システム出力ベクトルを生成するために、上記1つ以上のプロセッサはさらに、上記シストリックアレイの上記列次元に沿ったグループごとに、上記グループ内の各スパースシャードのそれぞれのシャード出力ベクトルを合算してそれぞれの列出力ベクトルを生成し、各グループの上記それぞれの列出力ベクトルを連結して上記システム出力ベクトルを生成するように構成されている。

10

【0010】

各乗算回路は、上記スパースシャードについての上記シャード入力行列のそれぞれからの上記それぞれの非ゼロ値を含むそれぞれのレジスタに結合される。

【0011】

上記複数の乗算回路における乗算回路の数は、上記予め定められた最大非ゼロ閾値と等しい。

20

【0012】

上記スパースシャードはクロスバー回路をさらに含み、上記スパースシャードはさらに、上記クロスバー回路によって、上記シャード入力ベクトルの上記複数のベクトル値を受け、上記複数の乗算回路の各々への入力として、上記クロスバー回路によって、上記複数のベクトル値のうちの1つのベクトル値を送るように構成されている。

【0013】

上記スパースシャードはさらに、上記シャード入力行列の同じ列の非ゼロ値を、上記複数の乗算回路のうちの隣接する乗算回路のレジスタにロードするように構成されている。

【0014】

上記スパースシャードはさらに、上記シャード入力行列の各列に沿った非ゼロ値の位置を少なくとも指定する1つ以上の制御値を受けるとともに構成されており、上記スパースシャードの上記クロスバー回路はさらに、上記1つ以上の制御値を受け、上記1つ以上の制御値に従って、上記シャード入力行列の同じ列に沿った非ゼロ値が乗算されるベクトル値を、隣接する乗算回路に送るように構成されている。

30

【0015】

上記スパースシャードは複数の加算回路をさらに含み、上記スパースシャードは1つ以上のセグメントマーカをさらに含み、各セグメントマーカは、上記セグメントマーカにロードされたそれぞれの制御値の値に基づいて、上記複数の加算回路のうちのそれぞれの加算回路への入力をゲート開閉するように構成されており、上記スパースシャードはさらに、上記1つ以上の制御値の少なくとも一部を上記1つ以上のセグメントマーカにロードするように構成されており、上記シャード入力行列の1列目の非ゼロ値についての加算回路は、上記1列目とは異なる上記シャード入力行列の2列目の非ゼロ値を含む隣接する加算回路の入力を受けないようにゲート開閉されており、上記スパースシャードはさらに、上記複数の加算回路によって、上記1つ以上の積の1つ以上の合計を生成するように構成されており、上記1つ以上の合計の各々は、上記シャード入力行列の列の1つ以上の非ゼロ値に上記シャード入力ベクトルの1つ以上のそれぞれの値を乗算したそれぞれのセグメント化合計である。

40

【0016】

上記複数の加算回路は並列セグメント化合計回路を形成し、1つ以上の上記セグメント化合計の各々は、セグメントマーカによってゲート開閉されていない隣接する加算回路へ

50

の出力の合計である。

【0017】

上記クロスバー回路は第1のクロスバー回路であり、上記スパスシャードは第2のクロスバー回路をさらに含み、上記第2のクロスバー回路は、1つ以上の上記セグメント化合計を受け、上記1つ以上の制御値に従って1つ以上の上記セグメント化合計を配列して、上記スパスシャードについて上記シャード出力ベクトルのそれぞれを生成するように構成されている。

【0018】

上記第2のクロスバー回路はベネシュネットワークを形成してもよく、上記シャード入力行列は正方行列である。

【0019】

本開示の別の局面は、1つ以上のプロセッサと、上記1つ以上のプロセッサによって実行されると上記1つ以上のプロセッサに動作を実行させる命令を格納した1つ以上のメモリデバイスを含むシステムを提供し、上記動作は、ゼロ値および非ゼロ値を含む入力行列を受けると、上記入力行列を複数の部分行列に区分することを含み、各部分行列の非ゼロ値の数は予め定められた最大非ゼロ閾値以下であり、各部分行列の次元は予め定められた次元閾値以下である。

【0020】

上記動作はさらに、部分行列ごとに、上記部分行列の各列に沿った非ゼロ値の位置を指定する1つ以上のそれぞれの制御値を生成することを含み得る。

【0021】

上記動作はさらに、各部分行列および各部分行列についての上記1つ以上のそれぞれの制御値を、各部分行列および上記部分行列についての上記それぞれの制御値を処理するように構成された複数のスパスシャードに送ることを含み得る。上記複数のスパスシャードの各々は、部分行列と、上記部分行列についての上記1つ以上のそれぞれの制御値と、入力ベクトルの少なくとも一部を受け、上記部分行列と上記入力ベクトルの上記一部との積を表すそれぞれの出力シャードベクトルを生成するように構成されている。上記入力行列を区分することは、上記行列を、上記複数のスパスシャードにおけるスパスシャードの数と等しい数の部分行列に区分することを含む。

【0022】

本開示の他の局面は、対応するシステム、装置、および1つ以上の非一時的なコンピュータ読取可能記憶媒体に格納されたコンピュータプログラムを含む。

【図面の簡単な説明】

【0023】

【図1】本開示の局面に係る、スパスシャードのアレイを含むシステムの一例のブロック図である。

【図2】本開示の局面に係る、スパスシャードの一例のブロック図である。

【図3A】スパスシャードのシャード入力行列の一例を示す図である。

【図3B】シャード入力行列の一例の非ゼロ値のベクトルを示す図である。

【図3C】シャード入力行列の一例についての制御値のベクトルを示す図である。

【図4】シャード入力行列と、シャード入力ベクトルと、シャード入力行列についての制御値とを受け、スパスシャードによる行列乗算の一例を示す図である。

【図5】本開示の局面に係る、スパスシャード上のスパス行列の部分行列にシステム入力ベクトルを乗算するためのプロセスの一例のフロー図である。

【図6】複数のスパスシャードからのシステム入力ベクトルとシステム入力行列との積を表すシステム出力ベクトルを生成するプロセスの一例のフロー図である。

【図7A】システム入力行列の一例およびシステム入力ベクトルを示す図である。

【図7B】システム入力行列およびシステム入力ベクトルの区分を示す図である。

【図7C】区分されたシステム入力行列、およびシステム入力行列とシステム入力ベクトルとを乗算した積を表すシステム出力ベクトルを示す図である。

10

20

30

40

50

【図 8 A】本開示の局面に係る、シャード入力行列の 1 つ以上の制御値を用いてスパースシャードを構成するためのプロセスの一例のフロー図である。

【図 8 B】図 8 A のプロセスに従って構成されたスパースシャードを用いて行列 - ベクトル乗算を実行するためのプロセスの一例のフロー図である。

【図 9】本開示の局面に係る、システム入力行列から部分行列を生成するためのプロセスの一例のフロー図である。

【図 10】本開示の局面に係る、スパース行列乗算システムを実現するコンピューティング環境の一例のブロック図である。

【発明を実施するための形態】

【0024】

詳細な説明

#### 概要

本開示の局面は、スパース行列密ベクトル乗算のために構成された 1 つ以上の集積回路を含むシステムに関する。複数のスパースシャード (sparse shard) のシステムは、スパースシャードごとに、入力システム行列およびベクトルの部分行列および部分ベクトルを受け取ることができる。各スパースシャードは、集積回路の少なくとも一部であってもよく、乗算回路および加算回路などの複数の演算ユニットを実現することができる。各スパースシャードは、最大非ゼロ閾値以下の数の非ゼロ値を有する部分行列を受け取るように構成されており、この最大非ゼロ閾値は、システムが、たとえば本明細書に記載されているようにスパースシャードと任意に 1 つ以上の他のコンポーネントとを含むチップとして実現される場合に、予め定めることができる。

【0025】

各スパースシャードはさらに、本明細書に記載されている制御値などのメタデータを受け取ることができ、スパースシャードはこのメタデータを用いて、異なる入力を各乗算回路または加算回路に導くことができ、あるユニットからの出力が別のユニットに渡されるときにゲート開閉することができる。受けたメタデータに従ってスパースシャードを構成することによって、スパースシャードは、予め定められた次元閾値までの任意のサイズの入力部分行列を効率的に処理し、対応する長さのベクトルとして積を入力に出力することができる。

【0026】

部分行列を受け取る一部として、各スパースシャードは、部分行列の各列における非ゼロ値の位置を表す 1 つ以上の制御値を受け取ることができる。この 1 つ以上の制御値を用いて、スパースシャードは、スパースシャードが部分行列と部分ベクトルとの積を表すシャード出力ベクトルを生成するように、部分行列および部分ベクトルの個々の値をどのように乗算、加算、および配列すべきかを調整するように構成され得る。システムはさらに、たとえばシステム入力ベクトルにシステム入力行列を乗算することによって、各シャードのシャード出力ベクトルから、システム入力ベクトルをシステム入力行列に適用した積を表すシステム出力ベクトルを生成するように構成され得る。

【0027】

また、本開示の局面は、スパースシャードのアレイによる処理のためにスパース行列を前処理するためのシステムに関する。1 つ以上のプロセッサのシステムは、入力行列を複数の部分行列に区分するように構成され得、各部分行列の次元は、スパースシャードの最大行列入力サイズを指定する予め定められた次元閾値以下である。区分の一部として、システムは、部分行列のうちのいずれかが、予め定められた最大非ゼロ閾値よりも大きい数の非ゼロ値を含むか否かを識別し、それに応じて、識別した部分行列と同じ行または列に沿って部分行列を再区分することができる。システムは、部分行列の次元が次元閾値内にあり、部分行列の非ゼロ値の数が予め定められた非ゼロ閾値以下であるように、スパースシャードごとに部分行列を生成するまでこのプロセスを繰り返すことができる。

【0028】

本開示の局面に従って実現されるシステムは、スパース行列とベクトルとの反復乗算を

10

20

30

40

50

伴う作業負荷をより効率的に実行することができる。たとえば、本開示の局面に係るシステムオンチップ（SoC：system-on-a-chip）を実現するデバイスは、少なくとも、最終的な積に寄与しない行列のゼロ要素の冗長な「ゼロ乗算」計算が省略されるので、従来のアプローチよりも少ない処理サイクルで、スパース行列にベクトルを乗算した積を生成することができる。他のアプローチでは、スパース比が高い大きな行列では性能損失が増大する可能性があるが、本明細書に記載されているように実現されるシステムは、少なくとも、入力行列のスパース比が高くなるにつれて省略される冗長計算の割合が高くなるので、これらの行列を用いてさらに効率的に積を計算することができる。

**【0029】**

ニューラルネットワークの実行または訓練などの一部の作業負荷は、行列乗算の実行に大きく依存している。ハードウェアアクセラレータまたはその他のデバイスは、行列乗算などの特定の演算を効率的に実行することができるが、処理条件が制限されていることが多い。たとえば、デバイスは、行列のゼロ値と非ゼロ値との予め定義されたスパース比を要する場合があったり、処理する入力サイズが厳しく制限されている場合がある。行列乗算は、多くのニューラルネットワーク作業負荷にとってユビキタスなタイプの計算であるが、処理条件が制限されているアクセラレータは、さまざまなスパース比の行列の行列乗算を伴う作業負荷など、アクセラレータがサポートできる作業負荷の種類数が制限されている。

10

**【0030】**

本開示の局面は、異なるサイズおよびスパース比の行列を柔軟に処理することができるスパース行列乗算のためのシステムを提供する。スパース行列乗算のために構成されたスパースシャードのシステムのスパースシャードは、複数の列にわたって所与の入力部分行列の非ゼロ値の位置を追跡し続けながら、入力行列のゼロ値を破棄することができる。スパースシャードは、予め定められた最大非ゼロ閾値と同数のレジスタおよび乗算回路のみで構成され、入力部分行列の非ゼロ値のみをメモリに格納することができる。その結果、スパースシャードは、フルサイズの行列を格納するのと比較して、より効率的に、より少ないリソースで、入力データを格納して処理することができる。

20

**【0031】**

さらに、スパースシャードは、非ゼロ値と、それらに入力ベクトルの値を乗算した積とを、隣接する乗算回路および加算回路に沿ってそれぞれ配列するように構成され得る。スパースシャードは、これらの隣接する回路の出力を組み合わせることによって入力行列の列ごとにセグメント化合計を効率的に生成し、このセグメント化合計を再配列してシャード出力ベクトルを生成するように構成され得る。スパースシャードは、非ゼロ値の順序を保存して、正確に非ゼロ値を加算してこれに入力ベクトルの対応する値を乗算し、入力部分行列と入力ベクトルとの積としてシャード出力ベクトルを正確に生成することができる。スパースシャードは、入力行列の形状を変更するための前処理演算を必要とせずに、入力行列に対して乗算を実行することができる。

30

**【0032】**

また、このシステムは、複数のスパースシャードによって入力ベクトルと乗算するためにスパース行列を前処理するように構成された1つ以上のプロセッサを含み得る。このシステムは、利用可能なスパースシャードの数に応じて異なる数の部分行列を生成すること、および異なる最大非ゼロ閾値についての部分行列を生成することなど、スパースシャードの異なる構成をサポートすることができる。

40

**【0033】**

本明細書に記載されている本開示の局面は、1つ以上のプロセッサと複数のスパースシャードとを含むシステムによって実現することができる。このシステムは、たとえば、コンピューティングプラットフォームのデータセンター内のサーバコンピューティングデバイスなどのコンピューティングデバイスにインストールされたチップとして実現することができる。

**【0034】**

50

### システムの例

図1は、スパースシャード101A~Pのアレイ101を含むシステム100の一例のブロック図である。アレイ101は、行列入力に対して行列乗算を実行するように構成されたスパース行列乗算システム100などのスパース行列乗算システムの少なくとも一部であってもよい。

#### 【0035】

図2を参照して本明細書により詳細に記載されているように、スパースシャードは、算術演算を実行するように、算術演算を実行する回路間に入力および出力を配列するように、ならびに/または他の回路間の入力および出力をゲート開閉するように構成された回路の集合体である。たとえば、スパースシャードは、スパース行列の矩形シャードまたは部分と、ベクトルの線形シャードまたは片との間の行列乗算を実行するように構成され得る。本明細書により詳細に記載されているように、スパースシャードは、予め定められた最大閾値までの任意の数の非ゼロ値を有する、スパース行列のさまざまな異なるシャードまたは部を受けることができる。

10

#### 【0036】

回路の集合体は、各ゲートの状態(すなわち、開もしくは閉)、計算されるオペランド、ならびに/または各回路への入力および各回路からの出力の配列を制御するように、本明細書に記載されているように構成され得る。この構成は、スパースシャードが入力として受けられるように構成されているシャード入力行列の次元に、ならびに/またはシャード入力行列内の非ゼロ値の配列および数に、少なくとも部分的に依存し得る。各スパースシャードは、スパース行列にベクトルを乗算するように構成されているシステムに対する個々の回路コンポーネントであってもよい。

20

#### 【0037】

各スパースシャード101A~Pは、アレイ101内の2つ以上の他のスパースシャードと通信するように構成されている。スパースシャードは、そのすぐ隣のスパースシャード、たとえば、スパースシャード101A~Pの矩形配列によって定義される次元に沿った当該スパースシャードの直前または直後のスパースシャードと通信することができる。スパースシャード間の接続は、たとえば、スパースシャードをその近隣のスパースシャードに物理的に接続するバスまたは1つ以上の回路相互接続を介して実現することができる。各スパースシャード101A~Pは、行列の少なくとも一部およびベクトルの少なくとも一部を受け、入力行列とベクトルとを乗算した積を表す出力ベクトルを生成するように構成された1つ以上の回路として実現することができる。いくつかの例では、スパースシャード101A~Pはシストリックアレイとして編成されるが、さまざまな実現例において、アレイ101は一般にスパースシャード101A~Pの矩形配列に従って構成または配列される。

30

#### 【0038】

アレイ101は、回路基板または他の材料の上に複数のコンポーネントおよび集積回路を実現するシステムオンチップの少なくとも一部であってもよい。アレイ101は、コンピューティングデバイスの一部としてインストールされ、メモリ、プロセッサ、ネットワークコンポーネント、および/または周辺機器など、デバイスの他のコンポーネントとやり取りするように構成され得る。たとえば、アレイ101は、システム100の一部として実現された1つ以上のメモリデバイスから、システム入力ベクトル105およびシステム入力行列110を受けることができる。システム100は、出力として、システム入力ベクトル105とシステム入力行列110とを乗算した積を表すシステム出力ベクトル115を生成することができる。

40

#### 【0039】

いくつかの例では、システム出力ベクトル115は、システム100を実現する他のデバイスまたはデバイスのコンポーネントへの入力として供給され得る。たとえば、システム出力ベクトル115が、モデルパラメータ値にニューラルネットワークについてのあるベクトル入力を乗算した積である場合は、システム出力ベクトル115は、出力ベクトル

50

115の活性化関数を計算するように構成された1つ以上のプロセッサへの入力として供給され得る。

【0040】

システム100は、前処理エンジン150からシステム入力ベクトル105およびシステム入力行列110を受けることができる。前処理エンジン150は、システム100を実現するデバイスであってもなくてもよい1つ以上のコンピューティングデバイス上に実現することができる。システム入力行列110は、一例として、ニューラルネットワークについてのモデルパラメータ値の少なくとも一部を表す値を含み得る。システム入力行列110は、多次元配列またはテンソルなどのより複雑なデータ構造の一部であってもよい。システム100は、3次元テンソルまたは行列などのより大きなデータ構造の少なくとも一部に対応する各行列を受け、この行列にシステム入力ベクトル105を乗算した対応する出力を生成するように構成され得る。システム入力ベクトル105は、たとえば、そのモデルパラメータ値がシステム入力行列110の値によって少なくとも部分的に表される訓練済みニューラルネットワークへの入力であってもよい。

10

【0041】

前処理エンジン150は、システム入力行列110を処理して1つ以上の制御値111および行列パーティションデータ112を生成するように構成され得る。前処理エンジン150は、システム入力行列110を、スパース行列を格納するためのさまざまな異なるフォーマットで、たとえばその完全な行・列形式(すべての非ゼロ値およびゼロ値を有する)で受けることができる。他の例として、前処理エンジン150は、圧縮されたスパース列フォーマット、座標リストフォーマット、またはスパース行列を格納するためのさまざまな他のフォーマットのうちのいずれかに従って、システム入力行列110を受けることができる。

20

【0042】

前処理エンジン150がシステム入力行列110をその完全な形で受けるいくつかの例では、前処理エンジン150は、システム入力行列110を、メモリに格納するのに適していると予め定められたフォーマットに変換するように構成されている。たとえば、前処理エンジン150は、行列からゼロ値を除去し、行列内の非ゼロ値の元の位置に対する非ゼロ値の位置を追跡するための制御値を生成してから、この予め定められたフォーマットに変換することができる。

30

【0043】

制御値は、図2を参照して本明細書により詳細に記載されているように、システム入力ベクトル105の部分ベクトルを受けてこれにシステム入力行列110の部分行列を乗算するための各スパースシャードを構成するために使用される。パーティションデータ112は、どのようにシステム入力行列110を部分行列に区分すべきかを指定するデータである。各部分行列はそれぞれのスパースシャードにおいて入力として受けられ、パーティションデータ112は、アレイ101内のスパースシャードと同数の部分行列についての区分を指定する。

【0044】

前処理エンジン150は、ベクトルパーティションデータ106を生成するように構成され得る。図5A~図5Cを参照してより詳細に記載されているように、同じ行または列に沿ったスパースシャードは、同じ部分ベクトルを、それぞれの部分行列と乗算するための入力として受けすることができる。

40

【0045】

システム入力行列110およびシステム入力ベクトル105はシステム100の左側および右側に沿って供給されるものとして示されているが、システム100に入力データを供給するバスの正確な位置決めおよび方向付けは実現例によって異なり得る。たとえば、システム100と同じチップ上の他のコンポーネントの位置に基づいて、システム100に入力を供給してシステム100から出力を受けるためのバスまたは回路相互接続は、それらの他のコンポーネントの位置を考慮してさまざまに方向付けられ得るかまたは位置決

50

めされ得る。

【0046】

図2は、本開示の局面に係る、スパースシャード200の一例のブロック図である。たとえば、システム100のスパースシャード101A~Pの各々は、スパースシャード200を参照して本明細書に記載されているように実現することができる。

【0047】

スパースシャード200は、シャード入力ベクトル205およびシャード入力行列210を受け取るように構成されている。シャード入力ベクトル205は、1つ以上のベクトル値を含み、 $1 \times R$ の最大次元を有し得る。シャード入力行列は、1つ以上のゼロ値および1つ以上の非ゼロ値を含み、 $R \times C$ の最大次元を有する。R(行)およびC(列)は、予め定められた、スパースシャードが入力として受け取ることができるベクトル/行列の最大入力サイズに対応する、次元閾値である。異なる実現例では、スパースシャードは異なる次元RおよびCについて構成され得る。RおよびCは互いに等しくても異なってもよく、たとえばスパースシャードが処理するように構成されている異なる作業負荷のデータの性質に応じて、異なる次元に異なるスパースシャードを実現することができる。スパースシャードのレイのスパースシャードは、同じ最大次元閾値内の入力を受け取るように構成され得る。

10

【0048】

スパースシャード200は、スパースシャード200に関連付けられたメモリ内の予め定められたアドレス範囲内のシャード入力ベクトル205およびシャード入力行列210を受け取るように構成され得る。たとえば、スパースシャード200は、結合されたメモリ内の第1のアドレス範囲からシャード入力ベクトル205を自動的に取り出し、メモリ内の同じまたは異なるアドレス範囲からシャード入力行列210を取り出すように構成されている。スパースシャード200を有するシステム100を実現するデバイスまたはデバイスのコンポーネントは、システムによって実現される1つ以上のスパースシャードの各々に対応するメモリ内の場所にシャード入力行列およびベクトルを送るよう構成され得る。たとえば、前処理エンジン150は、システム入力行列および/またはシステム入力ベクトルを処理生成した後、個々のシャード入力行列およびシャード入力ベクトルを、スパースシャード200を含む各スパースシャードに対応するアドレス範囲に格納するよう構成され得る。

20

【0049】

最大次元閾値に加えて、スパースシャード200は、予め定められた最大非ゼロ閾値以下の非ゼロカウントを有するシャード入力行列を受け取るように構成されている。次元閾値と同様に、最大非ゼロ閾値は、スパースシャード200およびその対応するレイの実現例によって異なる値に設定され得る。たとえば、システムが処理するように構成される作業負荷のデータが一般に、スパース比が高い処理行列を含む場合、スパースシャードのシステムは、比較的高い最大非ゼロ閾値で構成され得る。本明細書に記載されているように、スパースシャードの乗算回路および加算回路の数は、その最大非ゼロ閾値に対応し、したがって、最大非ゼロ閾値が相対的に低いスパースシャードは、最大非ゼロ閾値が相対的に高いスパースシャードと比べて、少ない回路で構築することができる。

30

40

【0050】

スパースシャード200は、クロスバー215および乗算回路220を含み得る。クロスバー215は、最大でR個のベクトル値を受け、その値をN個の乗算回路220に分散させるように構成され、Nはスパースシャード200についての最大非ゼロ閾値と等しい。乗算回路220A~C、Nが示されているが、異なる実現例ではスパースシャード200はより多いまたはより少ない乗算回路を含み得ると理解される。

【0051】

クロスバー215は、入力を受けてその入力を1つ以上の宛先に渡すように構成された任意の1つまたは複数の回路として実現することができ、この宛先は乗算回路220などの他の回路であってもよい。乗算回路220は、2つのオペランド間のハードウェア乗算

50

を実行するためのさまざまな異なる技術のうちのいずれかに従って実現することができる。乗算回路の第1のオペランドは、クロスバー215が受けたベクトル値であってもよい。乗算回路の第2のオペランドは、シャード入力行列210からの非ゼロ値であってもよい。各非ゼロ値は、それぞれの乗算回路220A～C、Nのそれぞれのレジスタ221A～C、Nにロードされる。最大非ゼロ閾値と同数の乗算回路を有することにより、非ゼロ閾値内のシャード入力行列ごとに利用可能な乗算回路が提供される。各乗算回路は、そのそれぞれのレジスタに格納された非ゼロ値に、クロスバー215が受けたベクトル値を乗算する。

#### 【0052】

また、スパースシャード200は加算回路225を含み得る。各加算回路225A～C、Nは、対応する乗算回路から入力を受けるとともに構成されている。各加算回路の間にセグメントマーカがある。加算回路225A～C、Nおよびセグメントマーカ226A～C、N-1が図2に示されているが、乗算回路220と同様に、加算回路およびセグメントマーカの数を実現例によって異なり得ると理解される。

10

#### 【0053】

加算回路は、2つのオペランドのハードウェア追加のための任意の技術を用いて実現することができる。加算回路からの第1のオペランドは、乗算回路から受けた積であってもよい。たとえば、乗算回路220Aは、ベクトル値と非ゼロ値とを乗算した積を加算回路225Aに渡す。セグメントマーカは、ゲート入力値に応じて、隣接する加算回路間の入力をゲート開閉するように構成された回路または他のハードウェアコンポーネントである。

20

#### 【0054】

制御値230は、シャード入力行列210およびシャード入力ベクトル205とともに入力として受けられ、クロスバー215、セグメントマーカ226、および/またはクロスバー235のうち1つ以上を構成するために使用され得る。制御値230は、シャード入力行列210のそれぞれの非ゼロ値にそれぞれ対応する値のシーケンスであってもよい。制御値230は、1などの第1のタイプの値を含み得、これらの値は、シャード入力行列のそれぞれの列の第1の非ゼロ値である非ゼロ値に対応し得る。制御値230は、0などの第2のタイプの値を含み得、これらの値は、シャード入力行列の同じ列の1つ以上の他の非ゼロ値の後ろに位置する非ゼロ値に対応し得る。また、制御値230は、図3A～図3Cを参照して本明細書に記載されているように、シャード入力行列210およびシャード入力ベクトル205を処理するようにスパースシャードを構成するために本明細書に記載されているように使用され得る値の1つ以上のベクトルを含み得る。図4およびその対応する本明細書中の記載は、スパースシャードを用いる行列乗算の一例を示す。

30

#### 【0055】

いくつかの実現例では、クロスバー215、235は、スパースシャード200の異なる最大次元閾値を利用するように実現することができる。たとえば、次元RおよびCが等しいまたはほぼ等しい場合は、クロスバー235は、たとえばベネシュネットワークとして、正方形またはほぼ正方形の入力に対するクロスバー再配列のための任意の技術に従って実現することができる。

#### 【0056】

図3A～図3Cは、シャード入力行列300A、シャード入力行列内の非ゼロ値のベクトル300B、およびシャード入力行列300Aに対応する制御値のベクトル300Cの一例を示す。また、図3Cは追加の制御ベクトル305Cおよび310Cを示している。

40

#### 【0057】

図3Aは、スパースシャードのシャード入力行列300Aの一例を示す図である。分かりやすくするために、非ゼロ値は影付きセルとして示されている。分かりやすくするために、シャード入力行列300Aの列および行に沿って、かつベクトル300BおよびCに沿って、インデックスが与えられている。たとえば、行列300Aにおいて、行2、列4(2,4)の値は1である。

#### 【0058】

50

図3Bは、シャード入力行列300Aの一例の非ゼロ値のベクトル300Bを示す図である。ベクトル300Bの非ゼロ値は、シャード入力行列300Aの非ゼロ値を左から右に読み取ったときの出現順序に対応するが、正確な読み取り順は実現例ごとに異なり、たとえば右から左であってもよい。

【0059】

図3Cは、シャード入力行列300Aの一例についての制御値のベクトル300C、305C、および310Cを示す図である。いくつかの例では、ベクトル300C、305C、および310Cは、予め定められた順序に従って同じベクトルの一部であってもよく、スパースシャードは、この予め定められた順序を用いて、本明細書に記載されているように、制御値のベクトルを受け、ベクトル300C、305C、および310Cに従って

10

【0060】

ベクトル300Cは、スパースシャードのセグメントマーカを構成するために使用される制御値に対応する。ベクトル300Cにおける1の値の制御値（この例ではビット）は、行列300Aの新しい列の開始に対応する。インデックス0の値は、制御値のベクトル300Cの開始として、自動的に1に設定され得る。いくつかの実現例では、開始制御値は、処理スパースシャードによって省略されて定数と仮定され得る。ハードウェア実現例は、この値が1であることが既知であるという事実を利用することによって、その回路を簡素化することができる。ベクトル300Cにおけるインデックス1の値も、行列300Aの次の列の最初の非ゼロ値であるベクトル300Bにおけるインデックス1の値に対応

20

【0061】

別の例として、ベクトル300Cにおけるインデックス3のビットは、対応する非ゼロ値（ベクトル300Bのインデックス3の値1）が行列300Aの次の列（具体的には列2）の最初の非ゼロ値であるので、1に設定される。ベクトル300Cにおけるシーケンスは、すべての列のすべての非ゼロ値が表されるまで、この説明したパターンに従う。

【0062】

ベクトル305Cは、スパースシャードの入力クロスバーを構成するための制御値に対応する。ベクトル305Cは、シャード入力行列300Aの非ゼロ値ごとに、部分行列300A内の非ゼロ値の「y」座標を指定する。この例では、「y」次元はスパース行列300Aを垂直方向に上下するが、他の例では「y」次元はたとえば水平に定義されるなど、異なって定義され得る。たとえば、ベクトル305Cにおける要素ゼロの値「3」は、行列300Aの1列目の非ゼロ値「1」の「y」座標に対応する。別の例として、ベクトル305Cにおける要素6の値「4」は、行列300Aの4列目の一番下の値「1」に対応する。

30

【0063】

ベクトル310Cは、スパースシャードの加算回路によって生成された合計が、シャード出力ベクトルを生成するために出力クロスバーによってどのように配列されるかを構成するための制御値に対応する。

40

【0064】

行列乗算の数学的定義により、スパース行列の位置(x, y)における各値に、位置yにおける入力ベクトルの値が乗算される。この乗算の結果は位置xにおける出力に加算される。入力クロスバーは、ベクトル305Cを用いて、同じ列の非ゼロ値をスパースシャード内の隣接する乗算回路に配列する。出力クロスバーは、ベクトル310Cを用いて、スパースシャードによってシャード入力行列とシャード入力ベクトルとを乗算した積を表すシャード出力ベクトルにおいて、計算した合計を正しい順序で配列する。

【0065】

図2に戻って、クロスバー215は、制御値230を受け、シャード入力ベクトルにシ

50

ヤード入力行列が乗算されると、対応するシャード入力行列の列に一致する 1 つ以上の乗算回路への入力として受けられるようにシャード入力ベクトル 2 0 5 についての各値を配列するように構成され得る。

【 0 0 6 6 】

クロスバー 2 3 5 は、加算回路から 1 つ以上の合計を受け、受けた合計を再配列して、入力シャード行列と入力シャードベクトルとの乗算に対応する正しい出力シャードベクトルを得るように構成され得る。クロスバー 2 1 5 と同様に、クロスバー 2 3 5 を 1 つ以上の回路として実現するためのさまざまな異なる技術のうちのいずれかを適用することができる。

【 0 0 6 7 】

本明細書に記載されているように、セグメントマーカは、セグメントマーカのゲート入力値に応じて、隣接する加算回路間の入力をゲート開閉するように構成されている。たとえば、セグメントマーカが 1 の値の制御値を受けた場合は、セグメントマーカは、セグメントマーカに隣接する第 1 の加算回路からの出力が、セグメントマーカに隣接する第 2 の加算回路への入力として渡されることを防止することができる。セグメントマーカが 0 の値の制御値を受けた場合は、セグメントマーカは、第 1 の加算回路からの出力を第 2 の加算回路に（または、実現例によってはその逆に）渡す。このセグメントマーカの構成は、異なる列の非ゼロ値の合計とは別に、同じ列の非ゼロ値に対応する合計のみを加算することに相当する。図 2 を参照して本明細書に記載されているように、この 1 つ以上の合計は、クロスバー 2 3 5 に渡され、シャード出力ベクトル 2 4 0 を生成するように再配列され得る。

【 0 0 6 8 】

クロスバー 2 3 5 は、入力を破棄するか受けるかを判断し、その入力をシャード出力ベクトル 2 4 0 内のその正しい位置に一致するように再配列するように構成され得る。図 2 に示されるように、各加算回路はクロスバー 2 3 5 に出力を渡すことができる（クロスバー 2 3 5 への矢印で示す）。ある加算回路の後ろの次のセグメントマーカがゲート開閉されていない場合は、隣接する加算回路間のランニングサムが終了していないので、クロスバー 2 3 5 はその加算回路への出力を破棄することができる。次のセグメントマーカがゲート開閉されている場合（または、最後の加算回路 2 2 5 N の場合はセグメントマーカがない場合）は、その列についてのランニングサムは終了しており、クロスバー 2 3 5 は、出力シャードベクトル 2 4 0 の一部になるべき入力としてランニングサムを受け、どの合計を無視し、どの合計をシャード出力ベクトル 2 4 0 の一部として含めるかを把握することによって、クロスバー 2 3 5 は、最大次元 C まで、列和を正確に追跡してさまざまな長さの出力ベクトルを生成することができる。

【 0 0 6 9 】

図 4 は、シャード入力行列 4 1 0 と、シャード入力ベクトル 4 0 5 と、シャード入力行列 4 1 0 についての制御値 4 3 0 とを受けけるスパスシャード 4 0 0 の計算の一例を示す。

【 0 0 7 0 】

シャード入力ベクトル 4 0 5 およびシャード入力行列 4 1 0 の値の一例について考える。

【 0 0 7 1 】

$$[ 1 \quad 3 \quad 2 ] \text{ (ベクトル 4 0 5)}$$

$$[ 0 \quad 3 \quad 0 \quad 2 \quad 0 \quad 4 \quad 1 \quad 0 \quad 0 ] \text{ (行列 4 1 0)}$$

この例における行列 4 1 0 についての対応する制御ベクトル 4 3 0 ~ 4 3 2 は、

$$[ 1 \quad 0 \quad 1 \quad 1 ] \text{ (制御ベクトル 4 3 0)}$$

$$[ 1 \quad 2 \quad 0 \quad 1 ] \text{ (制御ベクトル 4 3 1)}$$

$$[ 0 \quad 0 \quad 1 \quad 2 ] \text{ (制御ベクトル 4 3 2)}$$

である。

【 0 0 7 2 】

説明し易くするために、乗算回路 4 0 A ~ D を乗算器 A ~ D と略し、セグメントマーカ 4 3 A ~ C をセグメントマーカ A ~ C と略し、加算回路 4 2 A ~ D を加算器 A ~ D と略す

10

20

30

40

50

ることとする。

【 0 0 7 3 】

行列 4 1 0 の非ゼロ値に基づいて、乗算器 A には値 2 がロードされ、乗算器 B には値 1 がロードされ、乗算器 C には値 3 がロードされ、乗算器 D には値 4 がロードされる。なお、この例では、スパスシャードは 4 つの乗算器 A ~ D および 4 つの加算器 A ~ D のみを含む。

【 0 0 7 4 】

クロスバー 4 1 5 は、ベクトル値 1、3、および 2 を有するベクトル 4 0 5 を受ける。シャード入力ベクトル 4 0 5 からシャード出力ベクトル 4 4 0 へのデータの経路を示すために、破線 4 5 A、実線 4 5 B、および点線 4 5 C が示されている。クロスバー 4 1 5 は制御ベクトル 4 3 1 を受け、制御ベクトル 4 3 1 の各値はシャード入力行列 4 1 0 のそれぞれの非ゼロ値の「y」座標に対応する。制御ベクトル 4 3 1 の値は、1、2、0、および 1 を含む。なお、スパス入力行列 4 1 0 は次元  $3 \times 3$  を有するので、「y」座標の値は 0 から 2 の範囲である。制御ベクトル 4 3 1 は、乗算器 A ~ D の各々について、シャード入力ベクトル 4 0 5 の値のうちのどの値をどの乗算器に送るべきかを指定する。

10

【 0 0 7 5 】

たとえば、制御ベクトル 4 3 1 の最初の値は 1 であり、シャード入力行列 4 1 0 の最初の非ゼロ値の「y」座標に対応する。「y」座標 1 は（ゼロの後の）2 番目の座標であるので、クロスバー 4 1 5 は、シャード入力ベクトル 4 0 5 の 2 番目の値を第 1 の乗算器（ここでは乗算器 A）に導く。制御ベクトル 4 3 1 の 2 番目の値は 2 であり、シャード入力行列 4 1 0 の 2 番目の非ゼロ値の「y」座標に対応する。クロスバー 4 0 5 は、次に、シャード入力ベクトル 4 0 5 の 3 番目の値を乗算器 B に導くように構成され得る。別の例として、制御ベクトル 4 3 1 の 3 番目の値は 0 であり、0 の「y」座標を有する次の非ゼロ値に対応する。クロスバー 4 1 5 は、シャード入力ベクトル 4 0 5 の最初の値を乗算器 C に導く。

20

【 0 0 7 6 】

行列 4 1 0 の 1 列目について、クロスバー 4 1 5 は値 3 を乗算器 A に導き、値 2 を乗算器 B に導く。行列 4 1 0 の 2 列目について、クロスバー 4 1 5 は値 1 を乗算器 C に導く。行列 4 1 0 の最後の 3 列目について、クロスバー 4 1 5 は値 3 を乗算器 D に導く。乗算器 A ~ D の積は、乗算器 A は  $6 (3 \times 2)$ 、乗算器 B は  $2 (2 \times 1)$ 、乗算器 C は  $3 (1 \times 3)$ 、乗算器 D は  $12 (3 \times 4)$  である。

30

【 0 0 7 7 】

次に、乗算器 A ~ D が計算した積を加算器 A ~ D が受ける。加算器 A は、乗算器 A からの積、つまり 6 を受ける。最初の制御値（1）は破棄される。加算器 A の前には加算器がないので、加算器 A はセグメントマーカ A に合計を渡す。制御値 4 3 0 の 2 番目の値はゼロであるので、セグメントマーカ A はゲート開閉されていない。加算器 B は、加算器 A から現在の合計（6）を受け、それを乗算器 B の積（2）に加算する。セグメントマーカ B はゲート開閉されているので、加算器 B の出力（8）はクロスバー 4 2 0 に渡される（破線 4 5 A で示す）。セグメントマーカ B は加算器 B からの出力をゲート開閉するので、加算器 C は乗算器 C の積（3）にゼロを加算する。セグメントマーカ C はゲート開閉されているので、加算器 C の出力（3）はクロスバー 4 2 0 に渡される（点線 4 5 C で示す）。最後に、加算器 D は、乗算器 D の積（12）を受け、これはスパスシャード内の最後の加算器であるので、その出力（12）をクロスバー 4 2 0 に自動的に渡す（実線 4 5 B で示す）。

40

【 0 0 7 8 】

クロスバー 4 2 0 は、受けた合計 8、3、および 12 を、出力シャードベクトル 4 4 0 を出力するための正しい順序に従って再配列する。クロスバー 4 2 0 は、値 0、0、2、および 1 を有する制御ベクトル 4 3 2 を受ける。図 3 A ~ 図 3 C を参照して本明細書に記載されているように、出力クロスバー 4 2 0 についての制御ベクトルの値は、非ゼロ値の「x」座標位置に対応する。「y」座標と同様に、この値はこの例では 0 から 2 の範囲で

50

ある。制御ベクトル 4 3 2 の最初の 2 つの値は 0 である。したがって、クロスバー 4 2 0 は、受けた第 1 の合計を出力シャードベクトル 4 4 0 の第 1 の要素に導く。クロスバー 4 3 2 の 0 の次の値は 2 である。クロスバー 4 3 2 は、加算器 A ~ D から受けた第 2 の合計を出力ベクトル 4 4 0 の第 2 の要素に導き（線 4 5 C で示す）、第 3 の合計を要素に導く（線 4 5 B で示す）ように構成されている。いくつかの例では、クロスバー 4 2 0 は、たとえば最初の 2 つがゼロであるベクトル 4 3 2 に示されるような、ベクトル 4 3 2 の連続する重複制御値をスキップするように構成されている。いくつかの例では、連続する重複制御値をスキップするのではなく、クロスバー 4 2 0 は、受けた入力合計に対して包含的 OR 演算を実行し、この包含的 OR 演算の結果を、連続する重複制御値に対応する出力ベクトル 4 4 0 の位置に出力するように構成されている。

10

#### 【 0 0 7 9 】

たとえば、ライン 4 6 は、出力クロスバー 4 2 0 への加算器 A の潜在的な入力ソースを示す。セグメントマーカ 4 3 A はゼロの値を有するので、クロスバー 4 2 0 への加算器 A の出力は抑制され、たとえばマスクされるかまたはゼロに設定される。代わりに、加算器 A の出力はセグメントマーカ 4 3 A を通って加算器 B に渡される。いくつかの例では、出力クロスバー 4 2 0 が制御値 4 3 2 を受けると、出力クロスバー 4 2 0 は、（ライン 4 6 を通して）受けた第 1 の合計（ゼロの値を有する）と、（ライン 4 5 A を通して）受けた第 2 の合計（加算器 B からの 8 の値を有する）に対して包含的 OR 演算を実行する。クロスバー 4 2 0 は、非ゼロオペランドを出力するために包含的 OR 演算を実行するように構成され得る。クロスバー 4 2 0 は、受けた合計に対して包含的 OR 演算（たとえば、0 OR 8）を実行した後に 8 を出力し、その結果を出力ベクトル 4 4 0 内の最初の位置に渡す。いくつかの例では、出力クロスバー 4 2 0 は加算器 A ~ D のうちの少なくともいくつかから個々の出力を受けてそれらを合計することができる。

20

#### 【 0 0 8 0 】

加算回路は、セグメントを定義する連続した数値範囲を追加するためのさまざまな異なる回路構成のうちのいずれかに従って実現することができ、各セグメントは、対応するスパースシャードによって処理されるシャード入力行列のそれぞれの列内の値に対応する。たとえば、スパースシャード 2 0 0 または 4 0 0 の加算回路は、各列の非ゼロ値にシャード入力ベクトルのそれぞれの値を乗算することによって積の順次セグメント化合計を実行するための、（たとえば、スパースシャード 2 0 0 または 4 0 0 によって示されるような）1 つ以上の順次セグメント化合計回路として実現することができる。いくつかの実現例では、加算回路は、並列セグメント化合計回路として合計ツリーを実行するように構成され得る。個々の加算回路は、本明細書に記載されているように、対応する入力を並列に加算し、その合計を、任意の介在するセグメントマーカのゲート値に従って、出力クロスバーおよび/または隣接する加算回路に渡すように構成され得る。

30

#### 【 0 0 8 1 】

並列セグメント化合計回路は、特に合計している項の数が多い場合に、順次セグメント化合計回路と比べて回路の待ち時間を短縮することができる。セグメントマーカをゲート開閉するための制御値によって、並列セグメント化合計が可能になる。なぜなら、少なくとも、合計する値の範囲を、異なるセグメント内の値に対応する加算回路間の入力をゲート開閉するとともに同じセグメント内の追加される値に対応する加算回路間の入力を可能にする、セグメントマーカのゲート値に従って追跡することができるからである。

40

#### 【 0 0 8 2 】

いくつかの実現例では、出力クロスバー 4 2 0 は、各加算器によって計算された個々の合計を受け、その合計を加算し、制御ベクトル 4 3 2 を用いて、その合計を出力シャードベクトル 4 4 0 の対応する要素に導くように構成されている。2 つ以上の合計が同じ要素に導かれる場合、たとえば、制御ベクトル 4 3 2 がベクトル 4 3 2 の最初のゼロのように同じ値の重複を含む場合は、出力クロスバーは、図 4 の線 4 5 A で示すように 1 つの合計を受けるとはならず、出力シャードベクトル 4 4 0 内の同じ要素に導かれるように、受けた各合計を加算するように構成されている。

50

## 【 0 0 8 3 】

## 方法の例

図 5 は、本開示の局面に係る、スパースシャード上のスパース行列の部分行列にシステム入力ベクトルを乗算するためのプロセス 5 0 0 の一例のフロー図である。説明し易くするために、スパース行列の部分行列をシャード入力行列と呼ぶ。図 9 は、本明細書において、入力スパース行列を複数の部分行列に区分するためのプロセスの一例を示す。たとえば、図 2 のスパースシャード 2 0 0 などのスパースシャードがプロセス 5 0 0 を実行する。

## 【 0 0 8 4 】

ブロック 5 1 0 に従って、スパースシャードはシャード入力行列を受け、シャード入力行列は、予め定められた次元閾値内にあり、予め定められた最大非ゼロ閾値以下の非ゼロ値カウントを有する。

10

## 【 0 0 8 5 】

ブロック 5 2 0 に従って、スパースシャードは、複数のベクトル値を含むシャード入力ベクトルを受け、シャード入力ベクトルは、システム入力ベクトルの部分ベクトルであり、これは、図 1 および図 9 を参照して本明細書に記載されているように、前処理エンジンによってシステムに対する前処理入力の一部として生成することができる。

## 【 0 0 8 6 】

ブロック 5 3 0 に従って、スパースシャードは、それぞれのベクトル値にそれぞれの非ゼロ値を乗算した 1 つ以上の積を生成する。図 2 および 4 を参照して本明細書に記載されているように、スパースシャードは、シャード入力行列内の非ゼロ値の位置に対応する制御値で構成され得る。この制御値に基づいて、スパースシャードは、部分ベクトルの受信ベクトル値を、入力シャード行列の同じ列に沿った非ゼロ値が格納されている、対応する隣接する乗算回路に導くように構成され得る。本明細書に記載の図 8 A は、シャード入力行列についての制御値を用いてスパースシャードを構成するためのプロセスの一例を示す。

20

## 【 0 0 8 7 】

ブロック 5 4 0 に従って、スパースシャードは、この 1 つ以上の積の 1 つ以上の合計を生成する。図 2 を参照して本明細書に記載されているように、スパースシャードは、対応する乗算回路から入力を受けるとして構成された複数の加算回路を含む。加算回路はさらに、セグメントマーカに隣接する回路をゲート開閉するように設定された当該マーカに達するまで、隣接する加算回路に沿ってこれらの入力を加算するように構成されている。最大でセグメントマーカまでの加算回路入力の合計は、たとえば図 2 のクロスバー 2 3 5 などのクロスバーに渡され得る。

30

## 【 0 0 8 8 】

スパースシャードは、この 1 つ以上の合計から、シャード入力ベクトルにシャード入力行列を乗算した積であるシャード出力ベクトルを生成する。スパースシャードは、乗算回路によって生成された積のオペランドが第 1 のクロスバーによってどのように順序付けられたかに応じて、たとえば第 2 のクロスバーを通して、受けた合計を再配列することができる。

## 【 0 0 8 9 】

図 6 は、複数のスパースシャードからのシステム入力ベクトルとシステム入力行列との積を表すシステム出力ベクトルを生成するプロセス 6 0 0 の一例のフロー図である。図 1 のスパース行列乗算システム 1 0 0 などのスパースシャードのシステムが、プロセス 6 0 0 を実行することができる。

40

## 【 0 0 9 0 】

ブロック 6 1 0 に従って、スパースシャードのアレイの列次元に沿ったスパースシャードのグループごとに、グループ内の各スパースシャードのシャード出力ベクトルを加算してそのグループの列出力ベクトルを生成する。後述する図 7 C は、それぞれのシャード入力行列に従うスパースシャードのグループ化の一例を示す。スパースシャードのグループに沿って形成される次元は、実現例によって異なり得る。たとえば、システム入力行列およびシステム入力ベクトルがスパースシャードのアレイに供給される方向に応じて、グル

50

ープは、列とは対照的に、アレイの行に沿うことができる。

【 0 0 9 1 】

ブロック 6 2 0 に従って、システムは、各列出力ベクトルを連結してシステム出力ベクトルを生成する。システム出力ベクトルは、システム入力行列にシステム入力ベクトルを乗算した積である。いくつかの実現例では、図 9 を参照して本明細書により詳細に記載されているように、システムは、行列の列が置換されるシステム入力行列を受けて、たとえば、非ゼロ値の出現をスパースシャードに割り当てられた部分行列により均等に分散させることができる。それらの実現例では、システムは、連結されたシステム出力ベクトルの要素を再配列して元の順列を逆にするように構成され得る。システムは、システム入力行列、制御値、およびシステム入力行列の区分を定義するデータを受けの一部として、並び替えを定義するデータを受けすることができる。

10

【 0 0 9 2 】

図 7 A ~ 図 7 C および対応する記載は、スパース行列 7 0 0 の一例とベクトル 7 5 0 との間の乗算の一例を示す。説明のために、この乗算は、 $4 \times 4$  アレイの 1 6 個のスパースシャードを有するシステムに対して実行されるものとして記載されている。

【 0 0 9 3 】

図 7 A は、システム入力行列 7 0 0 の一例およびシステム入力ベクトル 7 5 0 を示す図である。この図では、システム入力行列 7 0 0 は整数値で示されているが、要素はたとえば浮動小数点値などの他の値であってもよいと理解される。また、図 7 A ~ 図 7 C では、システム入力行列 7 0 0 を含むさまざまな行列の非ゼロ値の要素は影付きセルとして示されている。

20

【 0 0 9 4 】

図 7 B は、システム入力行列 7 0 0 およびシステム入力ベクトル 7 5 0 の区分を示す図である。この例では、システム入力行列 7 0 0 は、 $4 \times 4$  アレイのスパースシャードごとに 1 つの部分行列があるように、1 6 個の部分行列 7 0 0 A ~ P に区分されている。システム入力ベクトル 7 5 0 は、 $4 \times 4$  アレイのスパースシャードの列ごとに 1 つの部分ベクトルがあるように、4 つの部分ベクトル 7 5 0 A ~ D に区分されている。

【 0 0 9 5 】

部分ベクトルおよび部分行列を 1 6 個のスパースシャード (スパースシャード A ~ P と呼ぶ) にマッピングする一例は、以下の表 1 に示す通りである。

30

【 0 0 9 6 】

40

50

【表 1】

スパーシャード	入力部分ベクトル	入力部分行列
A	750A	700A
B	750A	700B
C	750A	700C
D	750A	700D
E	750B	700E
F	750B	700F
G	750B	700G
H	750B	700H
I	750C	700I
J	750C	700J
K	750C	700K
L	750C	700L
M	750D	700M
N	750D	700N
O	750D	700O
P	750D	700P

表 1

## 【0097】

図 7 C は、区分されたシステム入力行列 700、およびシステム入力行列とシステム入力ベクトル 750 とを乗算した積を表すシステム出力ベクトル 770 を示す図である。図 7 C は、列 705 A ~ D に沿ってグループ化されている、区分された行列を示している。図 6 を参照して本明細書に記載されているように、システムは、スパーシャードのレイの列ごとにシャード出力ベクトルを加算し、列出力ベクトルを生成することができる。図 7 C において、列出力ベクトル 710 A ~ D は列 705 A ~ D に対応する。システムは、列出力ベクトル 710 A ~ D を連結して、システム入力行列にシステム入力ベクトルを乗算した積としてシステム出力ベクトルを生成することができる。

## 【0098】

図 8 A ~ 図 8 B は、本開示の局面に係る、入力行列の 1 つ以上の制御値を用いてスパーシャードを構成し、行列乗算を実行するためのプロセス 800 A ~ B の一例のフロー図である。

10

20

30

40

50

## 【 0 0 9 9 】

図 8 A は、本開示の局面に係る、シャード入力行列の 1 つ以上の制御値を用いてスパー  
スシャードを構成するためのプロセス 8 0 0 A の一例のフロー図である。たとえば図 2 の  
スパーシャード 2 0 0 などのスパーシャードが、プロセス 8 0 0 A を実行することが  
できる。

## 【 0 1 0 0 】

ブロック 8 1 0 に従って、スパーシャードは、シャード入力行列の各列に沿った非ゼ  
ロ値の位置を指定する 1 つ以上の制御値を受ける。

## 【 0 1 0 1 】

ブロック 8 2 0 に従って、スパーシャードは、シャード入力行列の非ゼロ値を乗算回  
路のレジスタにロードする。図 2 を参照して本明細書に記載されているように、スパー  
シャードは、システムについて予め定められた最大非ゼロ閾値と等しい数の乗算回路を  
実現することができる。スパーシャードは、左から右などの予め決められた読み取り方  
向に沿ってシャード入力行列を読み取ったときの出現順序で非ゼロ値をロードするこ  
とができる。

10

## 【 0 1 0 2 】

ブロック 8 3 0 に従って、スパーシャードは、この 1 つ以上の制御値をスパーシャ  
ードのクロスバーにロードする。図 2 のクロスバー 8 1 5 などの第 1 のクロスバーは、制  
御値を受け、行列 - ベクトル乗算の一部としてベクトル値が乗算される、対応するシャ  
ード入力行列の列に一致する 1 つ以上の乗算回路への入力として受けられるようにシャ  
ード入力ベクトルの各値を配列するように構成され得る。図 2 のクロスバー 2 3 5 などの第 2  
のクロスバーは、加算回路から 1 つ以上の合計を受け、受けた合計を再配列して、入力  
シャード行列と入力シャードベクトルとの乗算に対応する正しい出力シャードベクトル  
を得るように構成され得る。

20

## 【 0 1 0 3 】

ブロック 8 4 0 に従って、スパーシャードは、この 1 つ以上の制御値を、この制御値  
の値に基づいて加算回路への入力をゲート開閉するように構成された 1 つ以上のセグメン  
トマーカにロードする。

## 【 0 1 0 4 】

図 8 B は、図 8 A のプロセス 8 0 0 A に従って構成されたスパーシャードを用いて行  
列 - ベクトル乗算を実行するためのプロセスの一例のフロー図である。

30

## 【 0 1 0 5 】

ブロック 8 5 0 に従って、スパーシャードは、入力シャード行列の非ゼロ値を受けて  
ロードする。

## 【 0 1 0 6 】

ブロック 8 6 0 に従って、スパーシャードは、シャード入力ベクトルのベクトル値を  
受け、これらのベクトル値をシャード入力行列の同じ列に沿った非ゼロ値とともに乗算回  
路に送る。

## 【 0 1 0 7 】

ブロック 8 7 0 に従って、スパーシャードは、セグメントマーカによってゲート開閉  
されていない隣接する加算回路から 1 つ以上のセグメント化合計を生成する。図 2 および  
図 8 A を参照して本明細書に記載されているように、スパーシャードは、セグメントマ  
ーカによってゲート開閉されていない加算回路間の合計を集約し、シャード入力行列の異  
なる列からの計算を表す加算回路をゲート開閉するように制御値を用いてセグメントマ  
ーカを構成することができる。スパーシャードがアクティブゲートビットを有するセグメ  
ントマーカに達すると、スパーシャードは、セグメント化合計を再配列してシャード出  
力ベクトルを生成するように構成された第 2 のクロスバーにセグメント化合計を渡す。

40

## 【 0 1 0 8 】

ブロック 8 8 0 に従って、スパーシャードは、この 1 つ以上のセグメント化合計から  
シャード出力ベクトルを生成する。図 2 および図 8 A を参照して本明細書に記載されてい

50

るように、クロスバーは、スパースシャードの加算回路から1つ以上のセグメント化合計を受けると、制御値を用いて構成され得る。クロスバーはさらに、セグメント化合計を再配列して、シャード入力行列にシャード入力ベクトルを乗算した積を表す正しいシャード出力ベクトルを生成するように構成され得る。

#### 【0109】

図9は、本開示の局面に係る、システム入力行列から部分行列を生成するためのプロセス900の一例のフロー図である。1つ以上の場所にある1つ以上のプロセッサが、プロセス900を実行することができる。たとえば、図1の前処理エンジン150などの前処理エンジンが、プロセス900を実行することができる。

#### 【0110】

ブロック910に従って、前処理エンジンはシステム入力行列を受けると、システム入力行列は、たとえば、図7Aのシステム入力行列700で示されるような行列であってもよい。

#### 【0111】

ブロック920に従って、前処理エンジンは、システム入力行列を複数の候補部分行列に区分する。区分の一部として、前処理エンジンは、予め定められた次元閾値を指定するパラメータと、部分行列を受けるとシステムによって実現されるスパースシャードの数を示す1つ以上のパラメータ値とを受けるとすることができる。たとえば、前処理エンジンは、8行×8列の次元閾値内で(4×4アレイのスパースシャードについて)16個の候補部分行列を生成するように構成され得る。いくつかの例では、前処理エンジンは、たとえば、異なる次元閾値および/またはスパースシャードの構成を有する異なるシステムにわたる前処理入力のために、更新されたパラメータ値を受けるとすることができる。

#### 【0112】

いくつかの実現例では、前処理エンジンがシステム入力行列を区分する前に、前処理エンジンはシステム入力行列の列を置換して、非ゼロ値を候補部分行列に均等に分散させる。たとえば、非ゼロ値が、予め定められた許容差を超えて入力行列の一方側よりも他方側に頻繁に出現する場合、前処理エンジンは、非ゼロ値の出現がより広がることにより、区分後にスパースシャードにより均等に分散するように、入力行列の列の順序を変更するように構成され得る。

#### 【0113】

前処理エンジンがこの順序付けを実行する場合、前処理エンジンは、順序付けを表すデータを、追加入力として、たとえばスパースシャードのアレイを有するシステムに送られるパーティションデータの一部として、システムに渡す。システムは、出力ベクトルを、システム入力行列の列が置換される前のシステム入力行列とシステム入力ベクトルとを乗算した出力と一致させるために、順序付けに従ってシステム出力ベクトルの要素を並べ替えるように構成され得る。

#### 【0114】

システム入力行列の列を置換すると、システム入力行列がスパースシャードのシステムによって処理される全体速度を向上させることができる。たとえば、列を置換することによって、特に、いくつかの例では、一部のスパースシャードがゼロ値のみを有するシャード入力行列を受け得、他のシャードが非ゼロ値のみを有するシャード入力行列、または非ゼロ閾値までの数の非ゼロ値を有するシャード入力行列を受け得る場合に、各スパースシャードをより効率的に使用することができる。

#### 【0115】

ブロック930に従って、前処理エンジンは、予め定められた非ゼロ閾値よりも大きい非ゼロ値カウントを有する候補部分行列があるか否かを判定する。予め定められた非ゼロ閾値よりも大きい非ゼロ値カウントを有する候補部分行列があると前処理エンジンが判定した場合は、ブロック940に従って、前処理エンジンは、候補部分行列と同じ行または列に沿って部分行列を再区分する。

#### 【0116】

10

20

30

40

50

図 7 B に示されるように、部分行列は、入力行列におけるそれらの値の位置に基づいて、列および行に沿って編成され得る。たとえば、部分行列 7 0 0 J が非ゼロ閾値よりも高い非ゼロ値カウントを含むと前処理エンジンが判定した場合は、前処理エンジンは、部分行列 7 0 0 J の行（部分行列 7 0 0 I、7 0 0 K、および 7 0 0 L を含む）に沿って、ならびに / または部分行列 7 0 0 J の列（部分行列 7 0 0 B、7 0 0 F、および 7 0 0 N を含む）に沿って、部分行列を再区分することができる。再区分を実行する際、前処理エンジンは、予め定められた次元閾値を使用し、候補部分行列の数が変わらないように再区分を実行する。たとえば、前処理は、候補部分行列を分割し、決定された候補部分行列の行 / 列に沿って部分行列の行 / 列を再分散させることができる。

【 0 1 1 7 】

10

ブロック 9 3 0 および 9 4 0 に従って、部分行列を再区分した後、前処理エンジンは、非ゼロ閾値よりも大きい非ゼロ値カウントを有する候補部分行列があるか否かを再び判定する。前処理エンジンは、最大非ゼロ閾値を超える非ゼロ値カウントを有する候補部分行列がないと判定するまで、ブロック 9 3 0 および 9 4 0 に従って判定および再区分を繰り返す、ブロック 9 5 0 に進むことができる。

【 0 1 1 8 】

ブロック 9 5 0 に示されるように、前処理エンジンはシステム入力ベクトルを区分する。システム入力ベクトルの各部分ベクトルは、スパースシャードのアレイを供給するバスがどのように配列されるかに応じて、スパースシャードの行または列の各々に入力される。前処理エンジンは、ベクトル次元に、受けたスパースシャードの行列を乗算できるように、たとえば数学的に有効な行列乗算のための有効な次元を乗算できるように、ベクトルを区分する。

20

【 0 1 1 9 】

たとえば図 7 B ~ 図 7 C、および表 1 に示されるように、ベクトル 7 5 0 は部分ベクトル 7 5 0 A ~ D に区分され、これらの各々が入力として 1 つ以上のスパースシャードに渡される。また、図 7 B では、各部分ベクトル 7 5 0 A ~ D の各次元は、対応する部分行列 7 0 0 A ~ P と乗算するための正しい次元を有している。たとえば、部分ベクトル 7 5 0 B は  $1 \times 2$  (行  $\times$  列) であり、部分行列 7 0 0 E ~ H の各々は 2 行を有しているため、部分ベクトル 7 5 0 B と部分行列 7 0 0 E ~ H との間の有効な行列乗算が可能である。

【 0 1 2 0 】

30

ブロック 9 6 0 に示されるように、前処理エンジンは、候補部分行列ごとに制御値を生成する。図 3 A ~ 図 3 C を参照して記載されているように、前処理エンジンは、行列の列ごとに開始非ゼロ値を示す制御値のベクトルを生成することができる。前処理エンジンは、候補部分行列ごとにこの生成を繰り返して、部分行列ごとに対応する制御値を生成する。

【 0 1 2 1 】

ブロック 9 7 0 に従って、前処理エンジンは制御値および候補部分行列を出力する。前処理エンジンは、制御値と、システム入力行列の区分を指定するデータとを、たとえば図 1 に示されるようなシステム 1 0 0 に出力することができる。

【 0 1 2 2 】

#### コンピューティング環境の例

40

図 1 0 は、本開示の局面に係る、スパース行列乗算システム 1 0 0 および前処理エンジン 1 5 0 を実現するコンピューティング環境の一例のブロック図である。前処理エンジン 1 5 0 は、サーバコンピューティングデバイス 1 0 1 5 などにおいて、1 つ以上の場所に 1 つ以上のプロセッサを有する 1 つ以上のデバイス上に実現することができる。ユーザコンピューティングデバイス 1 0 1 2 およびサーバコンピューティングデバイス 1 0 1 5 は、ネットワーク 1 0 6 0 を介して 1 つ以上のストレージデバイス 1 0 3 0 に通信可能に結合され得る。ストレージデバイス (複数可) 1 0 3 0 は、揮発性メモリと不揮発性メモリとの組み合わせであってもよく、コンピューティングデバイス 1 0 1 2、1 0 1 5 と同じまたは異なる物理的位置にあってもよい。たとえば、ストレージデバイス (複数可) 1 0 3 0 は、ハードドライブ、ソリッドステートドライブ、テープドライブ、光ストレージ、

50

メモリカード、ROM、RAM、DVD、CD-ROM、書き込み可能メモリ、および読み取り専用メモリなどの、情報を格納することができる任意の種類非一時的なコンピュータ読取可能媒体を含み得る。

【0123】

サーバコンピューティングデバイス1015は、1つ以上のプロセッサ1013およびメモリ1014を含み得る。メモリ1014は、プロセッサ(複数可)1013によって実行され得る命令1021を含む、プロセッサ(複数可)1013がアクセス可能な情報を格納し得る。また、メモリ1014は、プロセッサ(複数可)1013が取り出す、操作する、または格納することができるデータ1023を含み得る。メモリ1014は、揮発性メモリおよび非揮発性メモリなどの、プロセッサ(複数可)1013がアクセス可能な情報を格納することができる非一時的なコンピュータ読取可能媒体の一種であってもよい。プロセッサ(複数可)1013は、1つ以上の中央処理装置(CPU)、グラフィックス処理ユニット(GPU)、フィールドプログラマブルゲートアレイ(FPGA)、および/またはテンソル処理ユニット(TPU)などの特定用途向け集積回路(ASIC)を含み得る。

10

【0124】

サーバコンピューティングデバイス1015は、スパース行列乗算システム100を、たとえばシステムオンチップとして、ハードウェア内に実現することができる。システム100は、サーバコンピューティングデバイス1015にスロットインされたまたはインストールされた物理的チップの一部として実現することができる。システム100は、サーバコンピューティングデバイス1015の他のコンポーネントと通信するように構成されている。

20

【0125】

命令1021は1つ以上の命令を含み得、この命令は、プロセッサ(複数可)1013によって実行されると、この1つ以上のプロセッサに、この命令によって定義された動作を実行させる。命令1021は、プロセッサ(複数可)1013によって直接処理するためのオブジェクトコード形式で、または、解釈可能なスクリプトもしくはオンデマンドで解釈されるか事前にコンパイルされる独立したソースコードモジュールの集合体を含む他の形式で、格納され得る。命令1021は、本開示の局面と一致するスパースシャード400を実現するための命令を含み得る。プリプロセッサエンジン105は、プロセッサ(複数可)1013を用いて、および/またはサーバコンピューティングデバイス1015から離れて位置する他のプロセッサを用いて実行され得る。

30

【0126】

データ1023は、命令1021に従ってプロセッサ(複数可)1013によって取り出され、格納され、または修正され得る。データ1023は、コンピュータレジスタ内に、複数の異なるフィールドおよびレコードを有するテーブルとしてリレーショナルもしくは非リレーショナルデータベース内に、またはJSON、YAML、proto、もしくはXML文書として、格納され得る。また、データ1023は、2進値、ASCIIまたはユニコードなどであるがこれらに限定されないコンピュータ読取可能形式にフォーマットされ得る。さらに、データ1023は、数字、説明文、プロプライエタリコード、ポイント、他のネットワーク場所を含む他のメモリに格納されたデータの参照などの、関連情報を識別するのに十分な情報、または関連データを計算する機能によって使用される情報を含み得る。

40

【0127】

ユーザコンピューティングデバイス1012も、1つ以上のプロセッサ1016と、メモリ1017と、命令1018と、データ1019とを有して、サーバコンピューティングデバイス1015と同様に構成され得る。また、ユーザコンピューティングデバイス1012は、ユーザ出力1026およびユーザ入力1024を含み得る。ユーザ入力1024は、キーボード、マウス、機械的アクチュエータ、ソフトアクチュエータ、タッチスクリーン、マイクロフォン、およびセンサなどの、ユーザから入力を受けるための任意の適

50

切なメカニズムまたは技術を含み得る。

【0128】

サーバコンピューティングデバイス1015は、ユーザコンピューティングデバイス1012にデータを送信するように構成され得、ユーザコンピューティングデバイス1012は、ユーザ出力1026の一部として実現されたディスプレイに受信データの少なくとも一部を表示するように構成され得る。また、ユーザ出力1026を用いて、ユーザコンピューティングデバイス1012とサーバコンピューティングデバイス1015との間のインターフェースを表示することができる。ユーザ出力1026は、これに代えてまたはこれに加えて、1つ以上のスピーカ、トランスデューサまたは他の音声出力、ユーザコンピューティングデバイス1012のプラットフォームユーザに非視覚情報および非可聴情報を提供するハプティックインターフェースまたは他の触覚フィードバックを含み得る。

10

【0129】

図10は、プロセッサ1013、1016およびメモリ1014、1017がコンピューティングデバイス1015、1012の内部にあると示しているが、プロセッサ1013、1016およびメモリ1014、1017を含む本明細書に記載されているコンポーネントは、異なる物理的位置で動作することができ、かつ同じコンピューティングデバイスの内部にない、複数のプロセッサおよびメモリを含み得る。たとえば、命令1021、1018およびデータ1023、1019の一部は、取り外し可能なSDカードに格納され、他の部分は読み取り専用のコンピュータチップ内に格納されてもよい。命令およびデータの一部またはすべてが、プロセッサ1013、1016から物理的に離れているがプロセッサ1013、1016がアクセス可能な場所に格納されてもよい。同様に、プロセッサ1013、1016は、同時および/または順次動作を実行することができるプロセッサの集合体を含み得る。コンピューティングデバイス1015、1012の各々は、タイミング情報を提供する1つ以上の内部クロックを含み得て、このタイミング情報は、コンピューティングデバイス1015、1012によって実行される動作およびプログラムの時間測定に使用され得る。

20

【0130】

サーバコンピューティングデバイス1015は、データの処理要求をユーザコンピューティングデバイス1012から受けるように構成され得る。たとえば、環境1000は、プラットフォームサービスを公開するさまざまなユーザインターフェースおよび/またはAPIを介して、さまざまなサービスをユーザに提供するように構成されたコンピューティングプラットフォームの一部であってもよい。サービスを実行する一部として、サーバコンピューティングデバイス1015は、システム100を用いて受信データを処理し得る。たとえば、サービスが機械学習モデルを訓練している場合は、サーバコンピューティングデバイス1015は、システム100を用いて、機械学習モデルを訓練する一部として乗算演算を実行するように構成され得る。

30

【0131】

デバイス1012、1015は、ネットワーク1060を介した直接および間接通信が可能であってもよい。デバイス1015、1012は、情報を送受信するための開始接続を受け入れ得るリスニングソケットをセットアップしてもよい。ネットワーク1060自体が、インターネット、ワールドワイドウェブ、イントラネット、仮想プライベートネットワーク、ワイドエリアネットワーク、ローカルネットワーク、および1つ以上の企業に独自の通信プロトコルを用いるプライベートネットワークなど、さまざまな構成およびプロトコルを含み得る。ネットワーク1060はさまざまな近距離および長距離接続をサポートし得る。この近距離および長距離接続は、2.402GHz~2.480GHz（一般にブルートゥース（登録商標）規格に関連する）、2.4GHzおよび5GHz（一般にWi-Fi（登録商標）通信プロトコルに関連する）などの異なる帯域幅にわたって行われてもよく、または無線ブロードバンド通信のLTE（登録商標）規格などのさまざまな通信規格で行われてもよい。また、ネットワーク1060は、これに加えてまたはこれに代えて、さまざまなタイプのイーサネット（登録商標）接続を介するなどして、デバ

40

50

イス1012、1015間の有線接続をサポートし得る。

【0132】

図10には1つのサーバコンピューティングデバイス1015およびユーザコンピューティングデバイス1012が示されているが、本開示の局面は、順次もしくは並列処理のためのパラダイムにおいて、または複数の装置の分散ネットワークを介してなど、さまざまな異なる構成および量のコンピューティングデバイスに従って実現することができる。いくつかの実現例では、本開示の局面は、1つのデバイス、およびその任意の組み合わせに対して実行され得る。さらに、前処理エンジンおよびスパス行列乗算システム100は同じサーバコンピューティングデバイス1015上に実現されるものとして示されているが、いくつかの実現例では、前処理エンジン150は、サーバコンピューティングデバイス1015とは別の、1つ以上のサーバコンピューティングデバイスおよび/またはユーザコンピューティングデバイス1012上に実現される。

10

【0133】

本開示の局面は、デジタル回路、コンピュータ読取可能記憶媒体において、1つ以上のコンピュータプログラム、または上述のうちの1つ以上の組み合わせとして実現することができる。コンピュータ読取可能記憶媒体は、たとえば、コンピューティングデバイスによって実行可能であり有形ストレージデバイスに格納される1つ以上の命令として、非一時的であってもよい。

【0134】

本明細書において、「～ように構成される」という表現は、コンピュータシステム、ハードウェア、またはコンピュータプログラム、エンジン、もしくはモジュールの一部に関連する異なる文脈で使用されている。システムが1つ以上の動作を実行するように構成される場合、これは、動作時に、システムにこの1つ以上の動作を実行させる適切なソフトウェア、ファームウェア、および/またはハードウェアが、システムにインストールされていることを意味する。あるハードウェアが1つ以上の動作を実行するように構成される場合、これは、動作時に、入力を受け、この入力に従ってこの1つ以上の動作に対応する出力を生成する1つ以上の回路を、ハードウェアが含むことを意味する。コンピュータプログラム、エンジン、またはモジュールが1つ以上の動作を実行するように構成される場合、これは、1つ以上のコンピュータによって実行されるとこの1つ以上のコンピュータにこの1つ以上の動作を実行させる1つ以上のプログラム命令を、コンピュータプログラムが含むことを意味する。

20

30

【0135】

図面に示されて請求項に記載されている動作は特定の順序で示されているが、これらの動作は示されている順序とは異なる順序で実行可能であること、また、一部の動作は省略可能であり、複数回実行可能であり、ならびに/または他の動作と並列におよび/もしくは同時に実行可能であることが理解される。さらに、異なる動作を実行するように構成された異なるシステムコンポーネントの分離は、コンポーネントの分離を要するものと理解されるべきではない。記載されているコンポーネント、モジュール、プログラム、およびエンジンは、1つのシステムとして統合されてもよく、または複数のシステムの一部であってもよい。

40

【0136】

特に明記しない限り、上述の代替例は互いに排他的ではなく、特有の利点を達成するためにさまざまな組み合わせで実現され得る。上記特徴のこれらのおよび他の変形例および組み合わせを、請求項によって定義される主題から逸脱することなく利用することができるので、実施例の上述の記載は、請求項によって定義される主題を限定するものとして解釈されるべきではなく、例示するものとして解釈されるべきである。加えて、本明細書に記載されている実施例の提供、ならびに、「などの」および「含む」などと表現されている節は、請求項の主題を特定の実施例に限定するものとして解釈されるべきではなく、むしろ、実施例は、多数の可能な実現例のうちの1つのみを示すことを意図している。さらに、さまざまな図における同一の参照番号は同一または同様の要素を識別し得る。

50

【図面】  
【図 1】

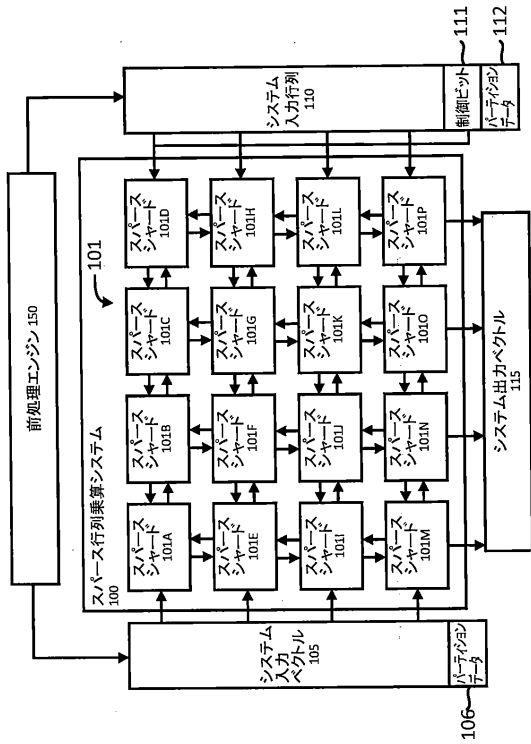


FIG. 1

【図 2】

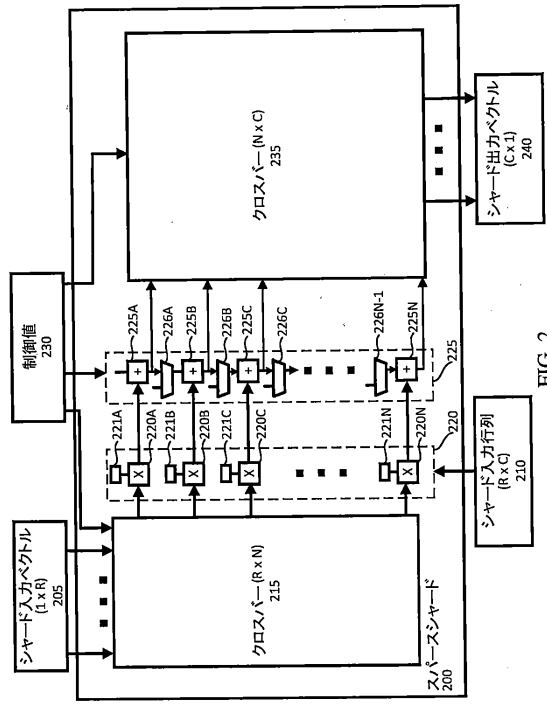


FIG. 2

【図 3 A】

	0	1	2	3	4
0 -	0	0	0	2	0
1 -	0	2	0	0	0
2 -	0	0	1	0	1
3 -	1	0	0	3	0
4 -	0	2	0	1	0

FIG. 3A

【図 3 B】

	0	1	2	3	4	5	6	7
1 -	1	2	2	1	2	3	1	1

FIG. 3B

10

20

30

40

50

【 図 3 C 】

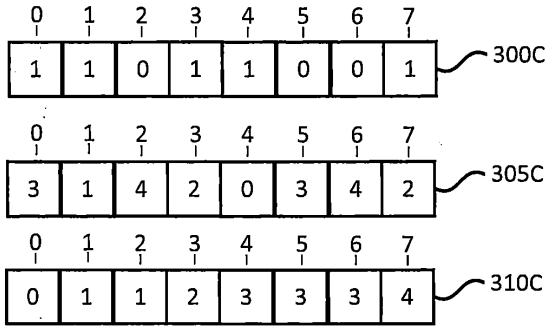


FIG. 3C

【 図 4 】

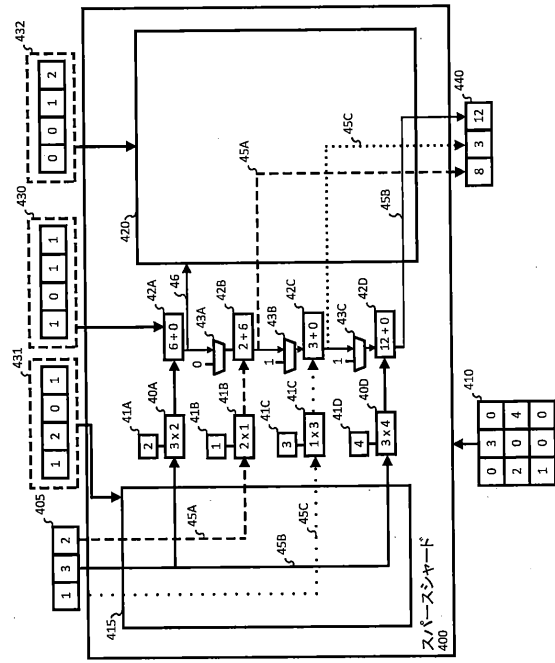


FIG. 4

【 図 5 】

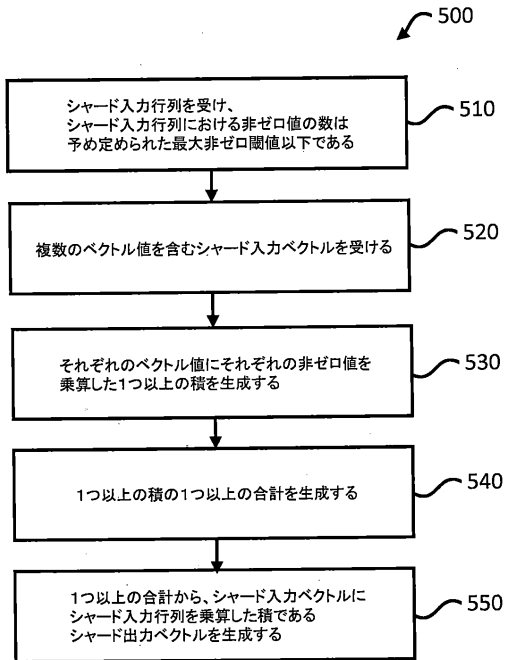


FIG. 5

【 図 6 】

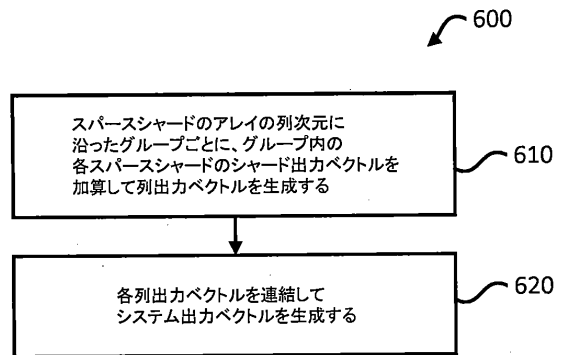


FIG. 6

10

20

30

40

50



【 図 8 B 】

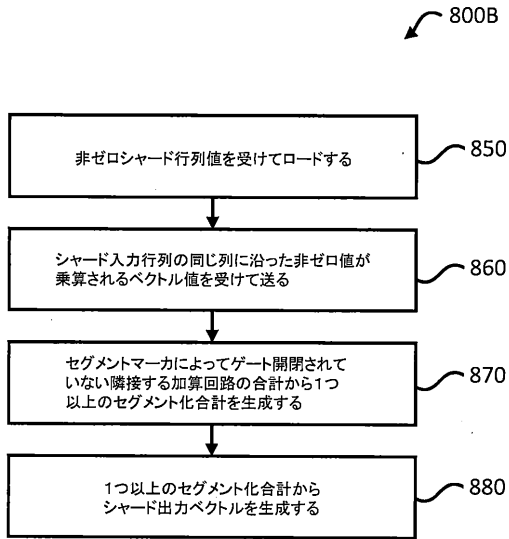


FIG. 8B

【 図 9 】

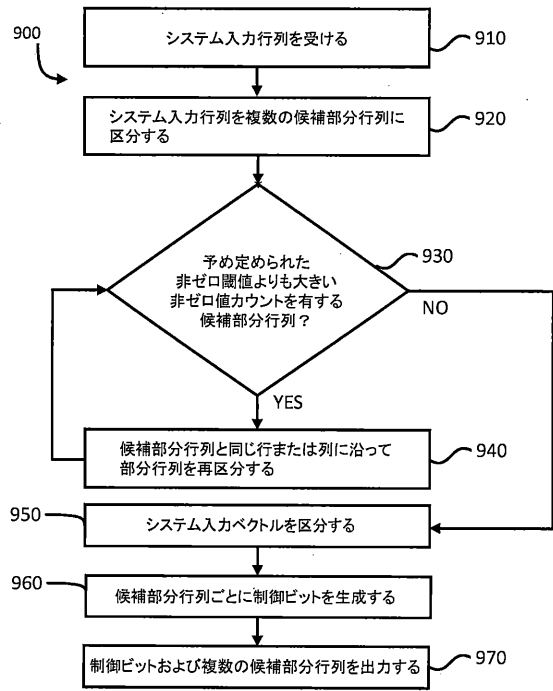


FIG. 9

【 図 1 0 】

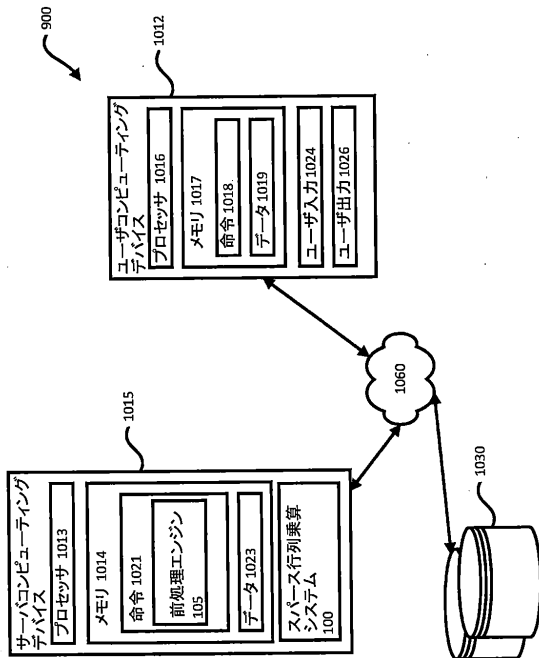


FIG. 10

## フロントページの続き

審査官 坂東 博司

(56)参考文献

YAN MINGYU ET AL , HyGCN: A GCN Accelerator with Hybrid Architecture , 2020 IEEE INTERNATIONAL SYMPOSIUM ON HIGH PERFORMANCE COMPUTER ARCHITECTURE (HPCA ) , 米国 , IEEE , 2020年02月22日 , pages 15-29 , [ 平成5年1月5日検索 ] , インターネット <URL : <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9065592> >

HE XIN ET AL , Sparse-TPU adapting systolic arrays for sparse matrices , PROCEEDINGS OF THE 34TH ACM INTERNATIONAL CONFERENCE ON SUPERCOMPUTING, ACM-PUB27 , 米国 , 2020年06月29日 , pages 1-12 , [ 平成5年1月5日検索 ] , インターネット <URL : [https://tnm.engin.umich.edu/wp-content/uploads/sites/353/2020/08/2020.6.sparse-tpu\\_ics2020.pdf](https://tnm.engin.umich.edu/wp-content/uploads/sites/353/2020/08/2020.6.sparse-tpu_ics2020.pdf) >

QIN ERIC ET AL , SIGMA: A Sparse and Irregular GEMM Accelerator with Flexible Interconnects for DNN Training , 2020 IEEE INTERNATIONAL SYMPOSIUM ON HIGH PERFORMANCE COMPUTER ARCHITECTURE (HPCA) , 米国 , IEEE , 2020年02月22日 , pages 58-70 , [ online ] , [ 平成5年1月5日検索 ] , インターネット <URL : <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9065523> >

(58)調査した分野

(Int.Cl. , D B 名)

G 0 6 F 1 7 / 1 6