

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7628112号
(P7628112)

(45)発行日 令和7年2月7日(2025.2.7)

(24)登録日 令和7年1月30日(2025.1.30)

(51)国際特許分類

F I

G 0 6 F	13/16	(2006.01)	G 0 6 F	13/16	5 2 0 B
G 0 6 F	12/00	(2006.01)	G 0 6 F	13/16	5 1 0 E
G 0 6 F	15/167	(2006.01)	G 0 6 F	12/00	5 6 0 F
G 0 6 N	3/063	(2023.01)	G 0 6 F	15/167	6 1 0 Z
			G 0 6 N	3/063	

請求項の数 14 (全29頁)

(21)出願番号 特願2022-517123(P2022-517123)
 (86)(22)出願日 令和2年9月14日(2020.9.14)
 (65)公表番号 特表2022-548641(P2022-548641 A)
 (43)公表日 令和4年11月21日(2022.11.21)
 (86)国際出願番号 PCT/US2020/050713
 (87)国際公開番号 WO2021/055280
 (87)国際公開日 令和3年3月25日(2021.3.25)
 審査請求日 令和4年5月13日(2022.5.13)
 審判番号 不服2024-2308(P2024-2308/J1)
 審判請求日 令和6年2月9日(2024.2.9)
 (31)優先権主張番号 16/573,805
 (32)優先日 令和1年9月17日(2019.9.17)
 (33)優先権主張国・地域又は機関 米国(US)

(73)特許権者 595168543
 マイクロン テクノロジー, インク .
 アメリカ合衆国, アイダホ州 8 3 7 1
 6 - 9 6 3 2 , ボイズ, サウス フェデ
 ラル ウェイ 8 0 0 0
 (74)代理人 100121083
 弁理士 青木 宏義
 (74)代理人 100138391
 弁理士 天田 昌行
 (74)代理人 100074099
 弁理士 大菅 義之
 (72)発明者 エイラート ショーン エス .
 アメリカ合衆国 カリフォルニア州 9 5
 6 6 3 ペンリン イングリッシュ コロ
 ニー ウェイ 1 7 2 1

最終頁に続く

(54)【発明の名称】 システムオンチップ及びアクセラレータチップを接続するメモリチップ

(57)【特許請求の範囲】

【請求項1】

アクセラレータチップと、
 メモリチップと、
 を含むシステムであって、
 前記メモリチップは、

配線を介して前記アクセラレータチップ及びシステムオンチップ(SoC)に接続するよ
 うに構成される単一セットのピンであって、前記アクセラレータチップと前記メモリチ
 ップとが、前記アクセラレータチップの1セットのピンと前記メモリチップの前記単一セ
 ットのピンの一部とを介して互いに直接接続され、前記SoCと前記メモリチップとが、前
 記SoCの1セットのピンと前記メモリチップの前記単一セットのピンの他の一部とを介
 して互いに直接接続される、前記単一セットのピンと、

前記単一セットのピンの前記他の一部を介して前記SoCから受信される計算入力デー
 タを格納して提供するように構成される複数の第一メモリセルであって、前記計算入力デー
 タは計算入力として前記アクセラレータチップによって使用される、前記複数の第一メ
 モリセルと、

を含み、

前記SoCは、前記メモリチップを介して間接的に前記アクセラレータチップと通信し、
 前記SoCは、バス又は配線を介して前記アクセラレータチップと直接通信することがな
 い、前記システム。

【請求項 2】

前記アクセラレータチップは、人工知能（AI）アクセラレータチップであり、前記複数の第一メモリセルは、前記単一セットのピンの前記他の一部を介して前記SOCから受信されるAI計算入力データを格納して提供するように構成され、前記AI計算入力データは、AI計算入力として前記AIアクセラレータチップによって使用される、請求項1に記載のシステム。

【請求項 3】

前記メモリチップは、

前記アクセラレータチップから前記単一セットのピンの前記一部を介して受信する第一計算出力データを格納して提供するように構成される複数の第二メモリセルであって、前記第一計算出力データは、前記SOCによって取得されるか、または計算入力として前記アクセラレータチップによって再使用される、前記複数の第二メモリセル、
を含む、請求項1に記載のシステム。

10

【請求項 4】

前記メモリチップは、前記単一セットのピンの前記他の一部を介して前記SOCから受信する第二計算出力データを格納するように構成される複数の第三メモリセルを含み、前記第二計算出力データは前記SOCによって取得される、請求項3に記載のシステム。

【請求項 5】

前記複数の第一メモリセル、前記複数の第二メモリセル、及び前記複数の第三メモリセルは、ダイナミックランダムアクセスメモリ（DRAM）セルを含む、請求項4に記載のシステム。

20

【請求項 6】

前記複数の第一メモリセル、前記複数の第二メモリセル、及び前記複数の第三メモリセルは、不揮発性ランダムアクセスメモリ（NVRAM）セルを含む、請求項4に記載のシステム。

【請求項 7】

前記NVRAMセルは、3D X Pointメモリセルを含む、請求項6に記載のシステム。

【請求項 8】

前記アクセラレータチップは、ベクトルプロセッサを含み、前記ベクトルプロセッサは、前記複数の第一メモリセル及び前記複数の第二メモリセルをメモリとして使用して、前記SOCについてのベクトル及び行列に対して数値計算を実行するように設定される、請求項3に記載のシステム。

30

【請求項 9】

前記アクセラレータチップは、特定用途向け集積回路（ASIC）を含み、前記ASICは、前記ベクトルプロセッサを含み、前記ベクトルプロセッサを介してAI計算を高速化するための専用ハードワイヤードである、請求項8に記載のシステム。

【請求項 10】

前記アクセラレータチップは、フィールドプログラマブルゲートアレイ（FPGA）を含み、前記FPGAは、前記ベクトルプロセッサを含み、前記ベクトルプロセッサを介してAI計算を高速化するための専用ハードワイヤードである、請求項8に記載のシステム。

40

【請求項 11】

アクセラレータチップと、
システムオンチップ（SOC）と、
メモリチップと、
を含むシステムであって、

前記メモリチップは、
配線を介して前記アクセラレータチップに接続するように構成される1セットのピンであって、前記アクセラレータチップと前記メモリチップとが、前記アクセラレータチップの1セットのピンと前記メモリチップの前記1セットのピンと前記配線とを介して互いに直

50

接続され、前記SoCと前記メモリチップとが、バスを介して互いに通信可能に接続される、前記1セットのピンと、

前記バスを介して前記SoCから受信される計算入力データを格納して提供するように構成される複数の第一メモリセルであって、前記計算入力データは計算入力として前記アクセラレータチップによって使用される、前記複数の第一メモリセルと、

を含み、

前記SoCは、前記メモリチップを介して間接的に前記アクセラレータチップと通信し、前記SoCは、バス又は配線を介して前記アクセラレータチップと直接通信することがない、前記システム。

【請求項12】

10

前記メモリチップは、

前記アクセラレータチップから前記メモリチップの前記1セットのピンを介して受信する第一計算出力データを格納して提供するように構成される複数の第二メモリセルであって、前記第一計算出力データは、前記SoCによって取得されるか、または第一計算入力として前記アクセラレータチップによって再使用される、前記複数の第二メモリセル、

を含む、請求項11に記載のシステム。

【請求項13】

前記SoCは、グラフィックスプロセッシングユニット(GPU)を含み、前記アクセラレータチップは、前記複数の第一メモリセル及び前記複数の第二メモリセルをメモリとして使用して前記GPUについての計算を実行して高速化するように設定される、請求項12に記載のシステム。

20

【請求項14】

前記アクセラレータチップは、ベクトルプロセッサを含み、前記ベクトルプロセッサは、前記複数の第一メモリセル及び前記複数の第二メモリセルをメモリとして使用して、前記GPUについてのベクトル及び行列に対して数値計算を実行するように設定される、請求項13に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願

30

本出願は、2019年9月17日に出願され、「MEMORY CHIP CONNECTING A SYSTEM ON A CHIP AND AN ACCELERATOR CHIP」と題された米国特許出願第16/573,805号に優先権を主張し、その開示全体は、参照により本明細書に援用される。

【0002】

本明細書に開示される少なくともいくつかの実施形態は、SoC及びアクセラレータチップ(例えば、AIアクセラレータチップ)を接続するメモリチップに関する。本明細書に開示される少なくともいくつかの実施形態は、メモリ階層及びメモリチップストラリングを使用してメモリを形成することに関する。

【背景技術】

40

【0003】

メインメモリなどのメモリは、コンピュータまたはコンピューティングデバイスでの即時使用のために情報を格納するコンピュータハードウェアである。一般に、メモリはコンピュータストレージよりも高速で動作する。コンピュータストレージによって、情報にアクセスする速度が遅くなるが、その容量が増え、データの信頼性が高くなることもできる。メモリの1つのタイプであるランダムアクセスメモリ(RAM)は、高い動作速度を有することができる。

【0004】

通常、メモリは、アドレス指定可能な半導体メモリユニットまたはセルで構成される。メモリIC及びそのメモリユニットは、シリコンベースの金属酸化物半導体電界効果トラ

50

ンジスタ (M O S F E T) によって少なくとも部分的に実装されることができる。

【 0 0 0 5 】

メモリには、揮発性及び不揮発性という2つの主なタイプがある。不揮発性メモリは、フラッシュメモリ (ストレージとして使用されることもできる)、ならびに R O M、P R O M、E P R O M、及び E E P R O M (ファームウェアを格納するために使用されることができる) を含むことができる。別のタイプの不揮発性メモリは、不揮発性ランダムアクセスメモリ (N V R A M) である。揮発性メモリは、ダイナミックランダムアクセスメモリ (D R A M) などのメインメモリテクノロジー、及び通常はスタティックランダムアクセスメモリ (S R A M) を使用して実装されるキャッシュメモリを含むことができる。

【 0 0 0 6 】

A I アクセラレータは、マイクロプロセッサまたはコンピュータシステムの1つのタイプであり、このタイプは、人工ニューラルネットワーク、マシンビジョン、及び機械学習などのA I アプリケーションを含む、A I アプリケーションについての計算を高速化するように設定される。A I アクセラレータは、データ集約型またはセンサ駆動型タスクについてのデータ処理を向上させるハードワイヤードであることができる。A I アクセラレータは、1つ以上のコアを含むことができ、低精度演算及びインメモリコンピューティング用に配線されることができる。A I アクセラレータは、スマートフォン、タブレット、及びあらゆるタイプのコンピュータ (特に、センサ、ならびにグラフィックス及び光学処理などのデータ集約型タスクを有するコンピュータ) などの多くのデバイスに見いだされることができる。また、A I アクセラレータは、A I アプリケーションで使用される数値シミュレーション及び他のタイプのタスクに関するパフォーマンスを向上させるために、ベクトルプロセッサまたはアレイプロセッサを含むことができる。

【 0 0 0 7 】

S o C は、コンピュータコンポーネントをシングルチップに集積する集積回路 (I C) である。S o C における一般的なコンピュータコンポーネントは、中央処理装置 (C P U)、メモリ、入出力ポート、及びセカンダリストレージを含む。S o C は、そのすべてのコンポーネントを単一の基板またはマイクロチップ上に含むことができ、一部のチップは25セント硬貨よりも小さくなることができる。S o C は、さまざまな信号処理機能を有することができる、グラフィックスプロセッシングユニット (G P U) など、専用のプロセッサまたはコプロセッサを含むことができる。緊密に集積されることにより、S o C は、同等の機能を有する従来のマルチチップシステムよりも電力の消費をはるかに少なくすることができる。これにより、S o C は、モバイルコンピューティングデバイス (スマートフォン及びタブレットなどの内の) の統合に有益になる。また、S o C は、組み込みシステム及びモノのインターネット (特にスマートデバイスが小さい場合) に有用であることができる。

【 0 0 0 8 】

メモリに戻り参照すると、コンピューティングシステムのメモリは、階層であることができる。コンピュータアーキテクチャではメモリ階層と称されることが多い、メモリ階層は、応答時間、複雑さ、容量、永続性及びメモリ帯域幅など、ある特定の要因に基づいて、コンピュータメモリを階層に分離することができる。それらのような要因は関連していることができ、多くの場合、メモリ階層の有用性をさらに強調するトレードオフであることができる。

【 0 0 0 9 】

一般に、メモリ階層はコンピュータシステムでのパフォーマンスに影響する。他の要因に優先してメモリ帯域幅と速度を優先順位付けするには、応答時間、複雑さ、容量、及び永続性などのメモリ階層の制限を考慮する必要がある場合がある。このような優先順位付けを管理するために、さまざまなタイプのメモリチップを組み合わせ、より高速なチップと、より信頼性の高い、または費用効果の高いチップなどとのバランスをとることができる。さまざまなチップのそれぞれをメモリ階層の一部と見なすことができる。そして、例えば、より高速なチップでのレイテンシを減らすために、メモリチップの組み合わせで

10

20

30

40

50

の他のチップは、バッファを充填してから、チップ間のデータ転送をアクティブにする信号を送ることによって応答することができる。

【0010】

メモリ階層は、さまざまなタイプのメモリユニットまたはセルを含むチップで構成されることができる。例えば、メモリセルはDRAMユニットであることができる。DRAMは、データの各ビットをメモリセルに格納するランダムアクセス半導体メモリの1つのタイプであり、メモリセルは、通常、コンデンサ及びMOSFETを含む。コンデンサは、充電されるか放電されるかいずれかが可能であり、これは、1ビットの中の2値、「0」及び「1」などで表される。DRAMでは、コンデンサの電荷が漏れ出すため、DRAMは、コンデンサごとに元の電荷を回復することによってコンデンサ内のデータを定期的にリライトする外部メモリリフレッシュ回路を必要とする。DRAMは、電源が切断されると、そのデータが急速に失われるため、揮発性メモリと見なされる。これは、データストレージがより永続的である、フラッシュメモリ、及びNVRAMなどの他のタイプの不揮発性メモリとは異なる。

10

【0011】

NVRAMの1つのタイプは3D XPointメモリである。3D XPointメモリでは、メモリユニットは、積層可能な交差格子状データアクセスアレイと組み合わせて、バルク抵抗の変化に基づいてビットを格納する。3D XPointメモリは、DRAMよりも費用効果が高いが、フラッシュメモリよりも費用効果が低い場合がある。また、3D XPointは、不揮発性メモリ及びランダムアクセスメモリである。

20

【0012】

フラッシュメモリは、別のタイプの不揮発性メモリである。フラッシュメモリの利点は、電氣的に消去されて再プログラムされることができることである。フラッシュメモリは、NAND型フラッシュメモリ及びNOR型フラッシュメモリという2つの主なタイプを有すると考えられており、これらは、フラッシュメモリのメモリユニットを実装することができるNAND及びNOR論理ゲートにちなんで名付けられている。フラッシュメモリユニットまたはセルは、対応するゲートのものと同様の内部特性を示す。NAND型フラッシュメモリはNANDゲートを含む。NOR型フラッシュメモリはNORゲートを含む。NAND型フラッシュメモリは、デバイス全体よりも小さくすることができるブロックに読み書きされてもよい。NOR型フラッシュは、シングルバイトを、消去した位置に書き込む、または独立して読み出すことを可能にする。NAND型フラッシュメモリの利点により、このようなメモリは、メモリカード、USBフラッシュドライブ、及びソリッドステートドライブによく利用されている。ただし、一般にフラッシュメモリを使用することの主なトレードオフは、DRAM及びNVRAMなどの他のタイプのメモリと比較して、特定のブロックに比較的少数の書き込みサイクルしかできないことである。

30

【0013】

本開示は、以下に示す詳細な説明及び本開示の様々な実施形態の添付図面から、より十分に理解される。

【図面の簡単な説明】

【0014】

40

【図1】SoC及びメモリチップを接続するアクセラレータチップ（例えば、AIアクセラレータチップ）を含む関連システムの一例を示す。

【図2】図1に示されるアクセラレータチップを含む関連システム、及び別個のメモリの例を示す。

【図3】図1に示されるアクセラレータチップを含む関連システム、及び別個のメモリの例を示す。

【図4】SoC及びアクセラレータチップ（例えば、AIアクセラレータチップ）を接続するメモリチップを含む、本開示のいくつかの実施形態によるシステムの一例を示す。

【図5】図4に示されるメモリチップを含むシステム、及び別個のメモリの例を示す。

【図6】図4に示されるメモリチップを含むシステム、及び別個のメモリの例を示す。

50

【図 7】図 4 に示されるメモリチップを含むシステム、及び別個のメモリの例を示す。

【図 8】本開示のいくつかの実施形態による、例示的なコンピューティングデバイスのパーツ配置の一例を示す。

【図 9】本開示のいくつかの実施形態による、例示的なコンピューティングデバイスのパーツ配置の別の例を示す。

【図 10】図 2 ~ 3 及び図 5 ~ 7 に示される別個のメモリに使用されることができるメモリチップストリングの例を示す。

【図 11】図 2 ~ 3 及び図 5 ~ 7 に示される別個のメモリに使用されることができるメモリチップストリングの例を示す。

【発明を実施するための形態】

【0015】

本明細書に開示される少なくともいくつかの実施形態は、SoC及びアクセラレータチップ（例えば、AIアクセラレータチップ）を接続するメモリチップ（例えば、DRAM）に関する。換言すれば、本明細書に開示される少なくともいくつかの実施形態は、メモリチップを介してアクセラレータチップ（例えば、AIアクセラレータチップ）をSoCに接続することに関する。アクセラレータチップは、メモリチップを介して間接的にSoCと通信する。メモリチップによってSoC及びアクセラレータチップを接続するメモリチップに置かれるデータは、アクセラレータチップへの要求であると解釈される。また、SoCは、SoC及びアクセラレータチップを接続するメモリチップを、アクセラレータチップに関与しないその動作のために任意選択で使用してもよい。したがって、SoC及びアクセラレータチップを接続するメモリチップは、SoCに使用される目的、及びアクセラレータチップに使用される目的という2つの一般的な目的を有することができる。それらのような実施形態のいくつかの例については、図4~7に示される、第一メモリチップ402、アクセラレータチップ404、及びSoC406を参照されたい。また、図8~9に示されるSoC806及び特定用途向けコンポーネント807を参照されたい。特定用途向けコンポーネント807は、デバイス800及び900のいくつかの実施形態では、第一メモリチップ402及びアクセラレータチップ404を含むことができる。

【0016】

図4~7に示されるように、SoC及びアクセラレータチップを接続するメモリチップは、論理的に（場合によっては物理的に）SoC及びアクセラレータチップの間にあることができる。また、SoC及びアクセラレータチップの間にあるアクセラレータ用のメモリチップは、2セットのピンを有する必要がない場合がある。いくつかの実施形態では、アクセラレータチップ及びメモリチップは、物理的に同じバス上にあることができる。ただし、中間にあるメモリチップを使用するいかなる状況でも、SoCがバスまたは配線を介してアクセラレータチップと直接通信することはない。したがって、SoC及びアクセラレータチップを接続するメモリチップは、少なくとも論理的にはアクセラレータチップとSoCとの間にある。また、メモリチップによって提供される、SoC及びアクセラレータチップの接続は、論理接続にすぎない場合がある。

【0017】

SoC及びアクセラレータチップを接続するメモリチップは、別個の2セットのピンを含むことができる。1セットは配線を介してアクセラレータチップに直接接続するためのものであり（例えば、図4、5、及び7に示される1セットのピン414、及び配線424を参照）、もう1セットは配線を介してSoCに直接接続するためのものである（例えば、図4~5に示される1セットのピン416、及び配線426を参照）。

【0018】

アクセラレータチップがメモリチップを介してSoCに接続されていると、SoCについての、一般に、またはより具体的には、いくつかの実施形態では、SoCに含まれるGPU（例えば、図4~7に示されるGPU408を参照）についての、特定用途向け計算（AI計算など）が高速化されることができる。いくつかの実施形態では、SoC内のGPUと、SoC及びアクセラレータチップを接続するメモリチップとを直接接続すること

10

20

30

40

50

ができる。いくつかの実施形態では、GPU及びアクセラレータチップを接続するメモリチップは、1セットのピンを含むことができ、この1セットのピン及び配線（例えば、1セットのピン414、及び配線424を参照）を介してアクセラレータチップに直接接続されることができる。アクセラレータチップは、対応する1セットのピン（例えば、1セットのピン415を参照）も含むことができる。そして、SoC及びアクセラレータチップを接続するメモリチップは、第二セットのピンを含むことができ、この第二セットのピン及び配線（例えば、1セットのピン416、及び配線426を参照）を介してGPUに直接接続されることができる。また、SoC内のGPUは、1セットのピンを含むことができ、この1セットのピン及び配線（例えば、1セットのピン417、及び配線426を参照）を介してメモリチップに直接接続されることができる。

10

【0019】

本開示の目的のために、本明細書に記載されるアクセラレータチップのいずれか1つが専用アクセラレータチップであること、またはそれを含むこと、またはその一部であることができることを理解されたい。専用アクセラレータチップの例は、低レイテンシまたは高帯域幅のメモリアクセスを提供することができる、人工知能(AI)アクセラレータチップ、仮想現実アクセラレータチップ、拡張現実アクセラレータチップ、グラフィックスアクセラレータチップ、機械学習アクセラレータチップ、またはいずれかの他のタイプのASICもしくはFPGAを含むことができる。例えば、本明細書に記載のアクセラレータチップのいずれか1つは、AIアクセラレータチップであること、またはそれを含むこと、またはその一部であることができる。

20

【0020】

アクセラレータチップは、人工ニューラルネットワーク、マシンビジョン、及び機械学習を含む、AIアプリケーションのハードウェア高速化のために設計されるマイクロプロセッサチップまたはSoC自体であることができる。いくつかの実施形態では、アクセラレータチップは、ベクトル及び行列に対して数値計算を実行するように設定される（例えば、ベクトル及び行列に対して数値計算を実行するように設定されることができる、図4に示されるベクトルプロセッサ412を参照）。アクセラレータチップは、ASICまたはFPGAであること、またはそれを含むことができる。アクセラレータチップのASIC実施形態では、アクセラレータチップは、特定用途向け計算(AI計算など)の高速化に専用のハードワイヤードであることができる。いくつかの他の実施形態では、アクセラレータチップは、変更されていないFPGAまたはGPUを超えた特定用途向け計算の高速化のために変更されている、変更されたFPGAまたはGPUであることができる。いくつかの他の実施形態では、アクセラレータチップは、変更されていないFPGAまたはGPUであることができる。

30

【0021】

アクセラレータチップに直接接続されているメモリチップ（例えば、第一メモリチップ402を参照）は、システム全体の複数のメモリチップを説明する際に明確にするために、本明細書では特定用途向けメモリチップとも称される。特定用途向けメモリチップは、必ずしも特定用途向け計算(AI計算など)専用のハードワイヤードであるとは限らない。特定用途向けメモリチップのそれぞれは、DRAMチップまたはNVRAMチップであることができる。そして、特定用途向けメモリチップのそれぞれは、アクセラレータチップに直接接続されることができ、SoCまたはアクセラレータチップによって特定用途向けメモリチップが構成された後、アクセラレータによる特定用途向け計算の高速化専用のメモリユニットを含むことができる。

40

【0022】

いくつかの実施形態では、SoCは、メインプロセッサ（例えば、CPU）を含むことができる。例えば、図4～7に示されるメインプロセッサ110を参照されたい。それらのような実施形態では、SoC内のGPUは、特定用途向けタスク及び計算（例えば、AIタスク及び計算）のための命令を実行することができ、メインプロセッサは、非特定用途向けタスク及び計算（例えば、非AIタスク及び計算）のための命令を実行することが

50

できる。そして、それらのような実施形態では、アクセラレータは、GPU専用の特定用途向けタスク及び計算の高速化を提供することができる。また、SoCは、SoCのコンポーネントを相互接続する（メインプロセッサ及びGPUを接続するなどの）ための独自のバスを含むことができる。そのうえ、SoCのバスは、SoCをSoCの外部のバスに接続するように構成されることのできるため、SoCのコンポーネントは、別個のメモリチップなどのSoCの外部のチップ及びデバイスと結合することができる。

【0023】

GPUの非特定用途向け計算及びタスク（例えば、非AI計算及びタスク）、またはアクセラレータチップを使用しないそれらのような計算及びタスクは、メインプロセッサによって実行される従来のタスクではない可能性があるが、別個のメモリチップなどの別個のメモリ（特定用途向けメモリであることのできる）を使用することができる。そして、メモリは、DRAM、NVRAM、フラッシュメモリ、またはそれらの任意の組み合わせで実装されることのできる。例えば、別個のメモリまたはメモリチップを、SoCの外部のバスを介してSoC及びメインプロセッサに接続することができる（例えば、図5に示されるメモリ204及びバス202を参照）。それらのような実施形態では、別個のメモリまたはメモリチップは、メインプロセッサ専用のメモリユニットを有することができる。また、別個のメモリまたはメモリチップを、SoCの外部のバスを介してSoC及びGPUに接続することができる（例えば、図5～7に示される第二メモリチップ204及びバス202を参照）。それらのような実施形態では、別個のメモリまたはメモリチップは、メインプロセッサまたはGPUにメモリユニットを含むことができる。

【0024】

本開示の目的のために、特定用途向けメモリチップ及び別個のメモリチップがメモリチップストリング（例えば、図10及び11に示されるメモリチップストリングを参照）などのメモリチップ群によって各置換されることのできることを理解されたい。例えば、別個のメモリチップは、少なくともNVRAMチップ及びそのNVRAMチップの下流にあるフラッシュメモリチップを含むメモリチップストリングで置換されることのできる。また、別個のメモリチップは、少なくとも2つのメモリチップで置換されることのでき、これらのチップのうちの1つはメインプロセッサ（例えば、CPU）用であり、もう1つのチップは非AI計算及び/またはタスクのためのメモリとして使用するためのGPU用である。

【0025】

さらに、本明細書に開示される少なくともいくつかの実施形態は、ベクトルプロセッサ（例えば、図4～7に示されるベクトルプロセッサ412を参照）を有するアクセラレータチップ（例えば、AIアクセラレータチップ）に関する。そして、本明細書に開示される少なくともいくつかの実施形態は、メモリ階層及びメモリチップストリングを使用してメモリを形成することに関する（例えば、図10及び11を参照）。

【0026】

本開示の目的のために、本明細書に記載されるアクセラレータチップのいずれか1つが専用アクセラレータチップであること、またはそれを含むこと、またはその一部であることのできることを理解されたい。専用アクセラレータチップの例は、低レイテンシまたは高帯域幅のメモリアクセスを提供することができる、AIアクセラレータチップ、仮想現実アクセラレータチップ、拡張現実アクセラレータチップ、グラフィックスアクセラレータチップ、機械学習アクセラレータチップ、またはいずれかの他のタイプのASICもしくはFPGAを含むことができる。

【0027】

図1は、SoC及びメモリチップを接続するアクセラレータチップ（例えば、AIアクセラレータチップ）を含む関連システムの一例を示す。

【0028】

図1は、例示的なシステム100を示し、このシステムは、ある程度、システム400に関連している。システム100は、第一メモリチップ104及びSoC106を接続す

10

20

30

40

50

るアクセラレータチップ102（例えば、AIアクセラレータチップ）を含む。示されるように、SoC106は、GPU108及びメインプロセッサ110を含む。メインプロセッサ110は、CPUである、またはそれを含むことができる。そして、アクセラレータチップ102は、ベクトルプロセッサ112を含む。

【0029】

システム100では、アクセラレータチップ102は、第一セットのピン114及び第二セットのピン116を含む。第一セットのピン114は、配線124を介して第一メモリチップ104に接続するように構成される。第二セットのピン116は、配線126を介してSoC106に接続するように構成される。示されるように、第一メモリチップ104は、対応する1セットのピン115を含み、このセットは、配線124を介してメモリチップをアクセラレータチップ102に接続する。SoC106のGPU108は、対応する1セットのピン117を含み、このセットは、配線126を介してSoCをアクセラレータチップ102に接続する。

10

【0030】

アクセラレータチップ102は、SoC106の特定用途向け計算（例えば、AI計算）を実行して高速化するように設定される。また、アクセラレータチップ102は、第一メモリチップ104を特定用途向け計算のためのメモリとして使用するように構成される。特定用途向け計算の高速化は、ベクトルプロセッサ112によって実行されることができる。アクセラレータチップ102内のベクトルプロセッサ112は、SoC106についてのベクトル及び行列に対して数値計算を実行するように設定されることができる。アクセラレータチップ102は、ASICを含むことができ、このASICは、ベクトルプロセッサ112を含み、ベクトルプロセッサ112を介して特定用途向け計算（例えば、AI計算）を高速化するための専用ハードワイヤードである。あるいは、アクセラレータチップ102は、FPGAを含むことができ、このFPGAは、ベクトルプロセッサ112を含み、ベクトルプロセッサ112を介して特定用途向け計算を高速化するための専用ハードワイヤードである。いくつかの実施形態では、アクセラレータチップ102は、GPUを含むことができ、このGPUは、ベクトルプロセッサ112を含み、ベクトルプロセッサ112を介して特定用途向け計算を高速化するための専用ハードワイヤードである。それらのような実施形態では、GPUは、ベクトルプロセッサ112を介して特定用途向け計算を高速化するための専用に変更されることができる。

20

30

【0031】

示されるように、SoC106はGPU108を含む。そして、アクセラレータチップ102は、GPU108についての特定用途向け計算（例えば、AI計算）を実行して高速化するように設定されることができる。例えば、ベクトルプロセッサ112は、GPU108についてのベクトル及び行列に対して数値計算を実行するように設定されることができる。また、GPU108は、特定用途向けタスク及び計算（例えば、AIタスク及び計算）を実行するように設定されることができる。

【0032】

また、示されるように、SoC106は、非AIタスク及び計算を実行するように設定されるメインプロセッサ110を含む。

40

【0033】

いくつかの実施形態では、メモリチップ104はDRAMチップである。それらのような例では、第一セットのピン114は、配線124を介してDRAMチップに接続するように構成されることができる。また、アクセラレータチップ102は、特定用途向け計算（例えば、AI計算）のためのメモリとして、DRAMチップ内のDRAMセルを使用するように構成されることができる。いくつかの他の実施形態では、メモリチップ104はNVRAMチップである。それらのような実施形態では、第一セットのピン114は、配線124を介してNVRAMチップに接続するように構成されることができる。また、アクセラレータチップ102は、特定用途向け計算のためのメモリとして、NVRAMチップ内のNVRAMセルを使用するように構成されることができる。さらに、NVRAMチ

50

ップは、3D XPointメモリチップである、またはそれを含むことができる。それらのような例では、第一セットのピン114は、配線124を介して3D XPointメモリチップに接続するように構成されることができ、アクセラレータチップ102は、3D XPointメモリチップ内の3D XPointメモリセルを、特定用途向け計算のためのメモリとして使用するように構成されることができ。

【0034】

いくつかの実施形態では、システム100は、配線を介して第一メモリチップ104に接続されるアクセラレータチップ102を含み、第一メモリチップ104は、特定用途向けメモリチップであることができる。また、システム100は、SoC106を含み、このSoCは、GPU108(AIタスクを実行するように設定されることができ)、及びメインプロセッサ110(非AIタスクを実行し、AIタスクをGPU108にデリゲートするように設定されることができ)を含む。それらのような実施形態では、GPU108は、配線126を介してアクセラレータチップ102に接続するように構成される1セットのピン117を含み、アクセラレータチップ102は、GPU108についてのAIタスクのAI計算を実行して高速化するように設定される。

10

【0035】

それらのような実施形態では、アクセラレータチップ102は、GPU108についてのベクトル及び行列に対して数値計算を実行するように設定されるベクトルプロセッサ112を含むことができる。そして、アクセラレータチップ102は、ASICを含み、このASICは、ベクトルプロセッサ112を含み、ベクトルプロセッサ112を介してAI計算を高速化するための専用ハードワイヤードである。または、アクセラレータチップ102は、FPGAを含み、このFPGAは、ベクトルプロセッサ112を含み、ベクトルプロセッサ112を介してAI計算を高速化するための専用ハードワイヤードである。または、アクセラレータチップ102は、GPUを含み、このGPUは、ベクトルプロセッサ112を含み、ベクトルプロセッサ112を介してAI計算を高速化するための専用ハードワイヤードである。

20

【0036】

また、システム100は、メモリチップ104を含み、アクセラレータチップ102は、配線124を介してメモリチップ104に接続されることができ、AIタスクのAI計算を実行して高速化するように設定されることができ。メモリチップ104は、DRAMセルを有するDRAMチップである、またはそれを含むことができ、DRAMセルは、アクセラレータチップ102によって、AI計算の高速化のためにデータを格納するように構成されることができ。または、メモリチップ104は、NVRAMセルを有するNVRAMチップである、もしくはそれを含むことができ、NVRAMセルは、アクセラレータチップ102によって、AI計算の高速化のためにデータを格納するように構成されることができ。NVRAMチップは、3D XPointメモリセルを含むことができ、3D XPointメモリセルは、アクセラレータチップ102によって、AI計算の高速化のためにデータを格納するように構成されることができ。

30

【0037】

図2~3は、それぞれ例示的なシステム200及び300、ならびに別個のメモリ(NVRAMなど)の例を示し、各システムは、図1に示されるアクセラレータチップ102を含む。

40

【0038】

図2では、バス202は、システム100(アクセラレータチップ102を含む)をメモリ204に接続する。メモリ204は、いくつかの実施形態ではNVRAMであることができ、システム100の第一メモリチップ104のメモリとは別のメモリである。そして、メモリ204は、いくつかの実施形態ではメインメモリであることができる。

【0039】

システム200では、システム100のSoC106は、バス202を介してメモリ204に接続される。そして、システム200の一部としてのシステム100は、アクセラ

50

レータチップ102、第一メモリチップ104、及びSoC106を含む。システム100のこれらのパーツは、バス202を介してメモリ204に接続される。また、図2に示されるように、SoC106に含まれるメモリコントローラ206は、システム100のSoC106によるメモリ204のデータアクセスを制御する。例えば、メモリコントローラ206は、GPU108及び/またはメインプロセッサ110によるメモリ204のデータアクセスを制御する。いくつかの実施形態では、メモリコントローラ206は、システム200内のすべてのメモリのデータアクセス（第一メモリチップ104及びメモリ204のデータアクセスなど）を制御することができる。そして、メモリコントローラ206は、第一メモリチップ104及び/またはメモリ204に通信可能に結合されることができる。

10

【0040】

メモリ204は、システム100の第一メモリチップ104によって提供されるメモリとは別のメモリであり、それは、メモリコントローラ206及びバス202を介して、SoC106のGPU108及びメインプロセッサ110にメモリとして使用されることができる。また、メモリ204は、GPU108及びメインプロセッサ110に、アクセラレータチップ102によって実行されない非特定用途向けタスクまたは特定用途向けタスク（非AIタスクまたはAIタスクなど）のためのメモリとして使用されることができる。それらのようなタスクについてのデータは、メモリコントローラ206及びバス202を介してメモリ204によってアクセスされ、そのメモリとの間で通信されることができる。

20

【0041】

いくつかの実施形態では、メモリ204は、システム200をホストするデバイスなどのデバイスのメインメモリである。例えば、システム200では、メモリ204は、図8に示されるメインメモリ808であることができる。

【0042】

図3では、バス202は、システム100（アクセラレータチップ102を含む）をメモリ204に接続する。また、システム300では、バス202は、アクセラレータチップ102をSoC106に接続し、アクセラレータチップ102をメモリ204に接続する。そのうえ示されるように、システム300では、バス202は、アクセラレータチップの第二セットのピン116、ならびにSoC106及びGPU108の配線126及び1セットのピン117の代わりにする。システム300内のアクセラレータチップ102は、システム200と同様に、システム100の第一メモリチップ104及びSoC106を接続する。ただし、この接続は、第一セットのピン114、及びバス202を介する。

30

【0043】

また、システム200と同様に、システム300では、メモリ204は、システム100の第一メモリチップ104のメモリとは別のメモリである。システム300では、システム100のSoC106は、バス202を介してメモリ204に接続される。そして、システム300では、システム300の一部としてのシステム100は、アクセラレータチップ102、第一メモリチップ104、及びSoC106を含む。システム100のこれらのパーツは、システム300では、バス202を介してメモリ204に接続される。また、図3に示されるものと同様に、SoC106に含まれるメモリコントローラ206は、システム100のSoC106によるメモリ204のデータアクセスを制御する。いくつかの実施形態では、メモリコントローラ206は、システム300内のすべてのメモリのデータアクセス（第一メモリチップ104及びメモリ204のデータアクセスなど）を制御することができる。そして、メモリコントローラは、第一メモリチップ104及び/またはメモリ204に接続されることができる。そして、メモリコントローラ206は、第一メモリチップ104及び/またはメモリ204に通信可能に結合されることができる。

40

【0044】

また、システム300では、メモリ204（いくつかの実施形態ではNVRAMである

50

ことができる)は、システム100の第一メモリチップ104によって提供されるメモリとは別のメモリであり、それは、メモリコントローラ206及びバス202を介して、SoC106のGPU108及びメインプロセッサ110にメモリとして使用されることができる。さらに、アクセラレータチップ102は、いくつかの実施形態及び状況では、バス202を介してメモリ204を使用することができる。そして、メモリ204は、GPU108及びメインプロセッサ110に、アクセラレータチップ102によって実行されない非特定用途向けタスクまたは特定用途向けタスク(非AIタスクまたはAIタスクなど)のためのメモリとして使用されることができる。それらのようなタスクについてのデータは、メモリコントローラ206及び/またはバス202を介してメモリ204によってアクセスされ、そのメモリとの間で通信されることができる。

10

【0045】

いくつかの実施形態では、メモリ204は、システム300をホストするデバイスなどのデバイスのメインメモリである。例えば、システム300では、メモリ204は、図9に示されるメインメモリ808であることができる。

【0046】

図4は、本開示のいくつかの実施形態による、アクセラレータチップ404(例えば、AIAアクセラレータチップ)及びSoC406を接続する第一メモリチップ402を含むシステム400の一例を示す。示されるように、SoC406は、GPU408及びメインプロセッサ110を含む。メインプロセッサ110は、システム400内のCPUである、またはこのCPUを含むことができる。そして、アクセラレータチップ404はベクトルプロセッサ412を含む。

20

【0047】

システム400では、メモリチップ402は、第一セットのピン414及び第二セットのピン416を含む。第一セットのピン414は、配線424を介してアクセラレータチップ404に接続するように構成される。第二セットのピン416は、配線426を介してSoC406に接続するように構成される。示されるように、アクセラレータチップ404は、対応する1セットのピン415を含み、このセットは、配線424を介して第一メモリチップ402をアクセラレータチップに接続する。SoC406のGPU408は、対応する1セットのピン417を含み、このセットは、配線426を介してSoCを第一メモリチップ402に接続する。

30

【0048】

第一メモリチップ402は、第二セットのピン416を介してSoC406から受信する計算入力データ(例えば、AI計算入力データ)を格納して提供するように構成される複数の第一メモリセルを含み、この計算入力データは、計算入力(例えば、AI計算入力)としてアクセラレータチップ404によって使用される。計算入力データは、複数の第一メモリセルからアクセスされ、第一メモリチップ402から、第一セットのピン414を介して送信され、アクセラレータチップ404によって受信されて使用される。複数の第一メモリセルは、DRAMセル及び/またはNVRAMセルを含むことができる。NVRAMセルを有する例では、NVRAMセルは、3D XPointメモリセルである、または3D XPointメモリセルを含むことができる。

40

【0049】

また、第一メモリチップ402は、第一セットのピン414を介してアクセラレータチップ404から受信する計算出力データ(例えば、AI計算出力データ)を格納して提供するように構成される複数の第二メモリセルを含み、この計算出力データは、SoC406によって取得される、または計算入力(例えば、AI計算入力)としてアクセラレータチップ404によって再使用される。計算出力データは、複数の第二メモリセルからアクセスされ、第一メモリチップ402から、第一セットのピン414を介して送信され、アクセラレータチップ404によって受信されて使用されることができる。また、計算出力データは、複数の第二メモリセルからアクセスされ、SoC406またはSoC内のGPU408から、第二セットのピン416を介して送信され、SoCまたはSoC内のGP

50

Uによって受信されて使用されることができる。複数の第二メモリセルは、DRAMセル及び/またはNVRAMセルを含むことができる。NVRAMセルを有する例では、NVRAMセルは、3D XPointメモリセルである、または3D XPointメモリセルを含むことができる。

【0050】

また、第一メモリチップ402は、1セットのピン416を介してSoC406から受信する非AIタスクに関連する非AIデータを格納するように構成される複数の第三メモリセルを含み、この非AIデータは非AIタスクのためのSoC406によって取得される。非AIデータは、複数の第三メモリセルからアクセスされ、第一メモリチップ402から第二セットのピン416を介して送信され、SoC406、SoC内のGPU408、またはSoC内のメインプロセッサ110によって受信されて使用されることができる。複数の第三メモリセルは、DRAMセル及び/またはNVRAMセルを含むことができる。NVRAMセルを有する例では、NVRAMセルは、3D XPointメモリセルである、または3D XPointメモリセルを含むことができる。

10

【0051】

アクセラレータチップ404は、SoC406についての特定用途向け計算（例えば、AI計算）を実行して高速化するように設定される。また、アクセラレータチップ404は、第一メモリチップ402を特定用途向け計算のためのメモリとして使用するように構成される。特定用途向け計算の高速化は、ベクトルプロセッサ412によって実行されることができる。アクセラレータチップ404内のベクトルプロセッサ412は、SoC406についてのベクトル及び行列に対して数値計算を実行するように設定されることができる。例えば、ベクトルプロセッサ412は、複数の第一メモリセル及び複数の第二メモリセルをメモリとして使用して、SoC406についてのベクトル及び行列に対して数値計算を実行するように設定されることができる。

20

【0052】

アクセラレータチップ404は、ASICを含むことができ、このASICは、ベクトルプロセッサ412を含み、ベクトルプロセッサ412を介して特定用途向け計算（例えば、AI計算）を高速化するための専用ハードワイヤードである。あるいは、アクセラレータチップ404は、FPGAを含むことができ、このFPGAは、ベクトルプロセッサ412を含み、ベクトルプロセッサ412を介して特定用途向け計算を高速化するための専用ハードワイヤードである。いくつかの実施形態では、アクセラレータチップ404は、GPUを含むことができ、このGPUは、ベクトルプロセッサ412を含み、ベクトルプロセッサ412を介して特定用途向け計算を高速化するための専用ハードワイヤードである。それらのような実施形態では、GPUは、ベクトルプロセッサ412を介して特定用途向け計算を高速化するため、専用に変更されることができる。

30

【0053】

示されるように、SoC406はGPU408を含む。そして、アクセラレータチップ402は、GPU408についての特定用途向け計算を実行して高速化するように設定されることができる。例えば、ベクトルプロセッサ412は、GPU408についてのベクトル及び行列に対して数値計算を実行するように設定されることができる。また、GPU408は、特定用途向けタスク及び計算を実行するように設定されることができる。また、示されるように、SoC406は、非AIタスク及び計算を実行するように設定されるメインプロセッサ110を含む。

40

【0054】

いくつかの実施形態では、システム400は、メモリチップ402、アクセラレータチップ404、及びSoC406を含み、メモリチップ402は、配線424を介してアクセラレータチップ404に接続するように構成される第一セットのピン414、及び配線426を介してSoC406に接続するように構成される第二セットのピン416を少なくとも含む。そして、メモリチップ402は、1セットのピン416を介してSoC406から受信するAI計算入力データを格納して提供するように構成される複数の第一メモ

50

リセルであって、このA I 計算入力データはA I 計算入力としてアクセラレータチップ404によって使用される、これら複数の第一メモリセルと、他のセットのピン414を介してアクセラレータチップ404から受信するA I 計算出力データを格納して提供するように構成される複数の第二メモリセルであって、このA I 計算出力データはS o C 406によって取得される、またはA I 計算入力としてアクセラレータチップ404によって再使用される、これら複数の第二メモリセルと、を含むことができる。そして、メモリチップ402は、非A I 計算のためのメモリに使用される複数の第三セルを含むことができる。

【0055】

また、S o C 406は、G P U 408を含み、アクセラレータチップ404は、複数の第一メモリセル及び複数の第二メモリセルをメモリとして使用して、G P U 408について10
のA I 計算を実行して高速化するように設定されることができる。そして、アクセラレータチップ404はベクトルプロセッサ412を含み、このベクトルプロセッサは、複数の第一メモリセル及び複数の第二メモリセルをメモリとして使用して、S o C 406についてのベクトル及び行列に対して数値計算を実行するように設定されることができる。

【0056】

また、システム400では、メモリチップ402内の複数の第一メモリセルは、1セットのピン416を介してS o C 406から受信するA I 計算入力データを格納して提供するように構成され、このA I 計算入力データは、A I 計算入力としてアクセラレータチップ404（例えば、A I アクセラレータチップ）によって使用されることができる。そして、メモリチップ402内の複数の第二メモリセルは、他のセットのピン414を介して20
アクセラレータチップ404から受信するA I 計算出力データを格納して提供するように構成され、このA I 計算出力データは、S o C 406によって取得される、またはA I 計算入力としてアクセラレータチップ404によって再使用されることができる。そして、メモリチップ402内の複数の第三メモリセルは、1セットのピン416を介してS o C 406から受信する非A I タスクに関連する非A I データを格納するように構成され、この非A I データは、非A I タスクのためのS o C 406によって取得されることができる。

【0057】

メモリチップ402内の複数の第一メモリセル、複数の第二メモリセル、及び複数の第三メモリセルは、それぞれ、D R A Mセル及び/またはN V R A Mセルを含むことができ、N V R A Mセルは3 D X P o i n tメモリセルを含むことができる。30

【0058】

図5～7は、それぞれシステム500、600、及び700、ならびに別個のメモリの例を示し、各システムは、図4に示されるメモリチップ402を含む。

【0059】

図5では、バス202は、システム400（メモリチップ402及びアクセラレータチップ404を含む）をメモリ204に接続する。メモリ204（例えば、N V R A M）は、システム400の第一メモリチップ402のメモリとは別のメモリである。そして、メモリ204はメインメモリであることができる。

【0060】

システム500では、システム400のS o C 406は、バス202を介してメモリ204に接続される。そして、システム500の一部としてのシステム400は、第一メモリチップ402、アクセラレータチップ404、及びS o C 406を含む。システム400のこれらのパーツは、バス202を介してメモリ204に接続される。また、図5に示されるように、S o C 406に含まれるメモリコントローラ206は、システム400のS o C 406によるメモリ204のデータアクセスを制御する。例えば、メモリコントローラ206は、G P U 408及び/またはメインプロセッサ110によるメモリ204のデータアクセスを制御する。いくつかの実施形態では、メモリコントローラ206は、システム500内のすべてのメモリのデータアクセス（第一メモリチップ402及びメモリ204のデータアクセスなど）を制御することができる。そして、メモリコントローラ206は、第一メモリチップ402及び/またはメモリ204に通信可能に結合されること40

10

20

30

40

50

ができる。

【0061】

メモリ204は、システム400の第一メモリチップ402によって提供されるメモリとは別のメモリであり、それは、メモリコントローラ206及びバス202を介して、SOC406のGPU408及びメインプロセッサ110にメモリとして使用されることができる。また、メモリ204は、GPU408及びメインプロセッサ110に、アクセラレータチップ404によって実行されない非特定用途向けタスクまたは特定用途向けタスク（非AIタスクまたはAIタスクなど）のためのメモリとして使用されることができる。それらのようなタスクについてのデータは、メモリコントローラ206及びバス202を介してメモリ204によってアクセスされ、そのメモリとの間で通信されることができる。

10

【0062】

いくつかの実施形態では、メモリ204は、システム500をホストするデバイスなどのデバイスのメインメモリである。例えば、システム500では、メモリ204は、図8に示されるメインメモリ808であることができる。

【0063】

図6では、図5と同様に、バス202は、システム400（メモリチップ402及びアクセラレータチップ404を含む）をメモリ204に接続する。システム500及び700に関してシステム600に一意である、第一メモリチップ402は、単一セットのピン602を含み、この単一セットのピンは、アクセラレータチップ404及びSOC406の両方にそれぞれ配線614及び616を介して第一メモリチップ402を直接接続する。また示されるように、システム600では、アクセラレータチップ404は、配線614を介してアクセラレータチップ404を第一メモリチップ402に直接接続する単一セットのピン604を含む。さらに、システム600では、SOCのGPUは、配線606を介してSOC406を第一メモリチップ402に直接接続する1セットのピン606を含む。

20

【0064】

システム600では、システム400のSOC406は、バス202を介してメモリ204に接続される。そして、システム600の一部としてのシステム400は、第一メモリチップ402、アクセラレータチップ404、及びSOC406を含む。システム400のこれらのパーツは、バス202を介してメモリ204に接続される（例えば、アクセラレータチップ404及び第一メモリチップ402はSOC406及びバス202を介したメモリ204への間接接続を有し、SOC406はバス202を介したメモリ204への直接接続を有する）。また、図6に示されるように、SOC406に含まれるメモリコントローラ206は、システム400のSOC406によるメモリ204のデータアクセスを制御する。例えば、メモリコントローラ206は、GPU408及び/またはメインプロセッサ110によるメモリ204のデータアクセスを制御する。いくつかの実施形態では、メモリコントローラ206は、システム600内のすべてのメモリのデータアクセス（第一メモリチップ402及びメモリ204のデータアクセスなど）を制御することができる。そして、メモリコントローラ206は、第一メモリチップ402及び/またはメモリ204に通信可能に結合されることができる。

30

40

【0065】

メモリ204は、システム400の第一メモリチップ402によって提供されるメモリとは別のメモリ（例えば、NVRAM）であり、それは、メモリコントローラ206及びバス202を介して、SOC406のGPU408及びメインプロセッサ110にメモリとして使用されることができる。また、メモリ204は、GPU408及びメインプロセッサ110に、アクセラレータチップ404によって実行されない非特定用途向けタスクまたは特定用途向けタスク（非AIタスクまたはAIタスクなど）のためのメモリとして使用されることができる。それらのようなタスクについてのデータは、メモリコントローラ206及びバス202を介してメモリ204によってアクセスされ、そのメモリとの間

50

で通信されることができる。

【0066】

いくつかの実施形態では、メモリ204は、システム600をホストするデバイスなどのデバイスのメインメモリである。例えば、システム600では、メモリ204は、図8に示されるメインメモリ808であることができる。

【0067】

図7では、バス202は、システム400（メモリチップ402及びアクセラレータチップ404を含む）をメモリ204に接続する。また、システム700では、バス202は、第一メモリチップ402をSoC406に接続するだけでなく、第一メモリチップ402をメモリ204にも接続する。そのうえ示されるように、システム700では、バス202は、第一メモリチップ402の第二セットのピン416、ならびにSoC406及びGPU408の配線426及び1セットのピン417の代わりをする。システム700内の第一メモリチップ402は、システム500及び600と同様に、システム400のアクセラレータチップ404及びSoC406を接続する。ただし、この接続は、第一セットのピン414、及びバス202を介する。

10

【0068】

また、システム500及び600と同様に、システム700では、メモリ204は、システム400の第一メモリチップ402のメモリとは別のメモリである。システム700では、システム400のSoC406は、バス202を介してメモリ204に接続される。そしてシステム700では、システム700の一部としてのシステム400は、第一メモリチップ402、アクセラレータチップ404、及びSoC406を含む。システム400のこれらのパーツは、システム700では、バス202を介してメモリ204に接続される。また、図7に示されるものと同様に、SoC406に含まれるメモリコントローラ206は、システム400のSoC406によるメモリ204のデータアクセスを制御する。いくつかの実施形態では、メモリコントローラ206は、システム700内のすべてのメモリのデータアクセス（第一メモリチップ402及びメモリ204のデータアクセスなど）を制御することができる。そして、メモリコントローラ206は、第一メモリチップ402及び/またはメモリ204に通信可能に結合されることができる。

20

【0069】

またシステム700では、メモリ204は、システム400の第一メモリチップ402によって提供されるメモリとは別のメモリ（例えば、NVRAM）であり、それは、メモリコントローラ206及びバス202を介して、SoC406のGPU408及びメインプロセッサ110にメモリとして使用されることができる。さらに、アクセラレータチップ404は、いくつかの実施形態及び状況では、第一メモリチップ402及びバス202を介してメモリ204を使用することができる。それらのような例では、第一メモリチップ402は、アクセラレータチップ404及びメモリ204についてのキャッシュを含むことができる。そして、メモリ204は、GPU408及びメインプロセッサ110に、アクセラレータチップ404によって実行されない非特定用途向けタスクまたは特定用途向けタスク（非AITaskまたはAITaskなど）のためのメモリとして使用されることができる。それらのようなタスクについてのデータは、メモリコントローラ206及び/またはバス202を介してメモリ204によってアクセスされ、そのメモリとの間で通信されることができる。

30

40

【0070】

いくつかの実施形態では、メモリ204は、システム700をホストするデバイスなどのデバイスのメインメモリである。例えば、システム700では、メモリ204は、図9に示されるメインメモリ808であることができる。

【0071】

本明細書に開示されるアクセラレータチップの実施形態（例えば、図1～3及び図4～7にそれぞれ示されるアクセラレータチップ102及びアクセラレータチップ404を参照）は、マイクロプロセッサチップまたはSoCなどであることができる。アクセラレー

50

タッチの実施形態は、人工ニューラルネットワーク、マシンビジョン、及び機械学習を含む、AIアプリケーションのハードウェア高速化のために設計されることができる。いくつかの実施形態では、アクセラレータチップ（例えば、AIアクセラレータチップ）は、ベクトル及び行列に対して数値計算を実行するように設定されることができる。そのような実施形態では、アクセラレータチップは、ベクトル及び行列に対して数値計算を実行するベクトルプロセッサを含むことができる（例えば、ベクトル及び行列に対して数値計算を実行するように設定されることができる、図1～3及び図4～7にそれぞれ示されるベクトルプロセッサ112及び412を参照）。

【0072】

本明細書に開示されるアクセラレータチップの実施形態は、ASICもしくはFPGAである、またはそれを含むことができる。アクセラレータチップのASIC実施形態では、アクセラレータチップは、特定用途向け計算（AI計算など）の高速化に専用のハードワイヤードである。いくつかの他の実施形態では、アクセラレータチップは、変更されていないFPGAまたはGPUを超えた特定用途向け計算（AI計算など）の高速化のために変更されている、変更されたFPGAまたはGPUであることができる。いくつかの他の実施形態では、アクセラレータチップは、変更されていないFPGAまたはGPUであることができる。

10

【0073】

本明細書に説明されているASICは、特定用途向け計算（AI計算など）の高速化など、特定の使用または用途にカスタマイズされているICを含むことができる。これは、CPU、または一般にグラフィックス処理のためのものであるGPUなどの別のタイプの汎用プロセッサによって通常実装される汎用用途とは異なる。

20

【0074】

本明細書に記載のFPGAは、IC及びFPGAの製造後に設計される、及び/または設定されるICに含まれることができる。したがって、IC及びFPGAはフィールドプログラマブルである。FPGA設定は、ハードウェア記述言語（HDL）を使用して指定されることができる。同様に、ASIC設定はHDLを使用して指定されることができる。

【0075】

本明細書で説明されるGPUは、ICを含み、このICは、メモリを迅速に操作して変更し、表示装置に出力されるフレームバッファ内の画像の生成及び更新を高速化するように設定されることができる。そして、本明細書で説明されるシステムは、GPUに接続される表示装置、ならびに表示装置及びGPUに接続されるフレームバッファを含むことができる。本明細書に説明されるGPUは、組み込みシステム、モバイルデバイス、パーソナルコンピュータ、ワークステーション、もしくはゲームコンソールの一部、または表示装置に接続され、この表示装置を使用する任意のデバイスであることができる。

30

【0076】

本明細書に記載のマイクロプロセッサチップの実施形態は、それぞれ、少なくとも中央処理装置の機能を組み込む1つ以上の集積回路である。各マイクロプロセッサチップは、多目的であり、少なくともクロック及びレジスタを含むことができ、これらのクロック及びレジスタは、入力としてバイナリデータを受け入れ、マイクロプロセッサチップに接続されたメモリに格納された命令に従ってレジスタ及びクロックを使用してデータを処理することによってチップを実装する。データを処理すると、マイクロプロセッサチップは入力及び命令の結果を出力として提供することができる。そして、この出力は、マイクロプロセッサチップに接続されたメモリに提供されることができる。

40

【0077】

本明細書で説明されるSoCの実施形態は、それぞれ、コンピュータまたは他の電子システムのコンポーネントを集積する1つ以上の集積回路である。いくつかの実施形態では、SoCは単一のICである。他の実施形態では、SoCは、分離され接続された集積回路を含むことができる。いくつかの実施形態では、SoCは、独自のCPU、メモリ、入出力ポート、セカンダリストレージ、またはそれらの任意の組み合わせを含むことができ

50

る。それらのような1つ以上のパーツは、本明細書で説明されるSoC内の単一の基板またはマイクロプロセッサチップ上にあることができる。いくつかの実施形態では、SoCは、25セント硬貨、5セント硬貨、または10セント硬貨よりも小さい。SoCのいくつかの実施形態は、モバイルデバイス(スマートフォンまたはタブレットコンピュータなど)、組み込みシステム、またはモノのインターネット内のデバイスの一部であることができる。一般に、SoCは、機能に基づいてコンポーネントを分離させてこれらのコンポーネントを、中央のインタフェース回路基板を介して接続するマザーボードベースのアーキテクチャを有するシステムとは異なる。

【0078】

アクセラレータチップ(例えば、AIアクセラレータチップ)に直接接続される、本明細書に記載のメモリチップの実施形態(例えば、図1~3に示される第一メモリチップ104、または図4~7に示される第一メモリチップ402を参照)は、システム全体の複数のメモリチップを説明するときに明確にするために、本明細書では特定用途向けメモリチップとも称される。本明細書に記載される特定用途向けメモリチップは、必ずしも特定用途向け計算(AI計算など)専用のハードワイヤードであるとは限らない。特定用途向けメモリチップのそれぞれは、DRAMチップもしくはNVRAMチップ、またはDRAMチップかNVRAMチップかいずれかと同様の機能を有するメモリデバイスであることができる。そして、特定用途向けメモリチップのそれぞれは、アクセラレータチップ(例えば、AIアクセラレータチップ、例えば、図1~3に示されるアクセラレータチップ102、及び図4~7に示されるアクセラレータチップ404を参照)に直接接続されることができ、特定用途向けメモリチップがアクセラレータチップまたは別個のSoCもしくはプロセッサ(例えば、図1~3及び図4~7にそれぞれ示されるSoC106及び406を参照)によって構成された後、アクセラレータチップによって特定用途向け計算(AI計算など)の高速化専用のメモリユニットまたはセルを含むことができる。

【0079】

本明細書で説明されるDRAMチップは、コンデンサ及びトランジスタ(MOSFETなど)を有するメモリセルまたはユニットにデータの各ビットを格納するランダムアクセスメモリを含むことができる。本明細書で説明されるDRAMチップは、ICチップの形態を取り、数十億個のDRAMメモリユニットまたはセルを含むことができる。各ユニットまたはセルでは、コンデンサは充電されるか、放電されるかいずれかであることができる。これにより、1ビットの中の2値を表すために使用される2つの状態を提供することができる。コンデンサでの電荷はコンデンサから緩徐に漏れる可能性があるため、コンデンサ及びメモリユニットの状態を維持するには、コンデンサ内のデータを定期的によりライトする外部メモリリフレッシュ回路が必要である。また、DRAMは、電源が切断されるとすぐにそのデータを失うという点で、揮発性メモリであり、フラッシュメモリまたはNVRAMなどの不揮発性メモリではない。DRAMチップの利点は、低コストで大容量のコンピュータメモリを必要とするデジタル電子機器でDRAMチップが使用されることができることである。DRAMは、GPU専用のメインメモリまたはメモリとして使用するのにも役立つ。

【0080】

本明細書に説明されるNVRAMチップは、DRAMとの主な差別化特徴である不揮発性のランダムアクセスメモリを含むことができる。本明細書で説明される実施形態に使用されることができるNVRAMユニットまたはセルの一例は、3DXPointユニットまたはセルを含むことができる。3DXPointユニットまたはセルでは、ビットストレージは、積層可能な交差格子状データアクセスアレイと組み合わせて、バルク抵抗の変化に基づく。

【0081】

本明細書で説明されるSoCの実施形態は、メインプロセッサ(CPUまたはCPUを含むメインプロセッサなど)を含むことができる。例えば、図1~3に示されるSoC106、及び図4~7に示されるSoC406だけでなく、図1~7に示されるメインプロ

10

20

30

40

50

セッサ110も参照されたい。それらのような実施形態では、SoC内のGPU（例えば、図1～3に示されるGPU108、及び図4～7に示されるGPU408を参照）は、特定用途向けタスク及び計算（AIタスク及び計算など）のための命令を実行することができ、メインプロセッサは、非特定用途向けタスク及び計算（非AIタスク及び計算など）のための命令を実行することができる。そして、それらのような実施形態では、SoCに接続されるアクセラレータチップ（例えば、図1～7に示されるアクセラレータチップのいずれか1つを参照）は、GPU専用の特定用途向けタスク及び計算（AIタスク及び計算など）の高速化を提供することができる。本明細書で説明されるSoCの実施形態のそれぞれは、SoCのコンポーネントを相互接続する（メインプロセッサ及びGPUを接続するなどの）ために独自のバスを含むことができる。また、SoCのバスは、SoCをSoCの外部のバスに接続するように構成されることができるため、SoCのコンポーネントは、別のメモリまたはメモリチップ（例えば、図2～3及び図5～7に示されるメモリ204、ならびに図8～9に示されるメインメモリ808を参照）などのSoCの外部のチップ及びデバイスと結合することができる。

10

【0082】

GPUの非特定用途向け計算及びタスク（例えば、非AI計算及びタスク）、またはアクセラレータチップを使用しない特定用途向け計算及びタスク（例えば、AI計算及びタスク）は、メインプロセッサによって実行される従来のタスクではない可能性があるが、別個のメモリチップなどの別個のメモリ（特定用途向けメモリであることができる）を使用することができ、このメモリは、DRAM、NVRAM、フラッシュメモリ、またはそれらの任意の組み合わせによって実装されることができる。例えば、図2～3及び図5～7に示されるメモリ204だけでなく、図8～9に示されるメインメモリ808も参照されたい。別個のメモリまたはメモリチップは、SoCの外部のバスを介してSoC及びメインプロセッサ（例えば、CPU）に接続されることができる（例えば、図2～3及び図5～7に示されるメモリ204だけでなく、図8～9に示されるメインメモリ808も参照、そして図2～3及び図5～7に示されるバス202だけでなく、図8～9に示されるバス804も参照）。それらのような実施形態では、別個のメモリまたはメモリチップは、メインプロセッサ専用のメモリユニットを有することができる。また、別個のメモリまたはメモリチップは、SoCの外部のバスを介してSoC及びGPUに接続されることができる。それらのような実施形態では、別個のメモリまたはメモリチップは、メインプロセッサまたはGPUにメモリユニットまたはセルを含むことができる。

20

30

【0083】

本開示の目的のために、本明細書に記載の特定用途向けメモリまたはメモリチップ（例えば、図1～3に示される第一メモリチップ104または図4～7に示される第一メモリチップ402を参照）、及び本明細書に記載の別個のメモリまたはメモリチップ（例えば、図2～3及び図5～7に示されるメモリ204だけでなく、図8～9に示されるメインメモリ808も参照）がそれぞれ、メモリチップストリング（例えば、図10及び11に示されるメモリチップストリングを参照）などのメモリチップ群で置換されることができることを理解されたい。例えば、別個のメモリまたはメモリチップは、少なくともNVRAMチップ及びそのNVRAMチップの下流にあるフラッシュメモリチップを含むメモリチップストリングで置換されることができる。また、別個のメモリチップは、少なくとも2つのメモリチップで置換されることができ、これらのチップのうちの1つはメインプロセッサ（例えば、CPU）用であり、もう1つのチップは非AI計算及び/またはタスクのためのメモリとして使用するためのGPU用である。

40

【0084】

本明細書に記載のメモリチップの実施形態は、メインメモリの一部であることができる、及び/またはコンピュータでの即時使用のために、または本明細書に記載のプロセッサのいずれか1つ（例えば、本明細書に記載の任意のSoCまたはアクセラレータチップ）による即時使用のために情報を格納するコンピュータハードウェアであることができる。本明細書に説明されるメモリチップは、コンピュータストレージよりも高速で動作するこ

50

とができる。コンピュータストレージによって、情報にアクセスする速度が遅くなるが、その容量が増え、データの信頼性が高くなることもできる。本明細書で説明されるメモリチップは、高い動作速度を有することができるメモリの1つのタイプであるRAMを含むことができる。メモリは、アドレス指定可能な半導体メモリユニットまたはセルで構成されることができ、そのユニットまたはセルは、MOSFETによって少なくとも部分的に実装されることができる。

【0085】

さらに、本明細書に開示される少なくともいくつかの実施形態は、ベクトルプロセッサ（例えば、図1～3及び図4～7にそれぞれ示されるベクトルプロセッサ112及び412を参照）を有するアクセラレータチップ（例えば、AIアクセラレータチップ）に関する。そして、本明細書に開示される少なくともいくつかの実施形態は、メモリ階層及びメモリチップストリングを使用してメモリを形成することに関する（例えば、図10及び11を参照）。

10

【0086】

本明細書で説明されるベクトルプロセッサの実施形態はそれぞれICであり、各ICは、ベクトルと称される一次元配列のデータ、または行列と称される多次元配列のデータ上で動作する命令を含む命令セットを実装することができる。ベクトルプロセッサは、命令がシングルデータ項目上で動作するスカラープロセッサとは異なる。いくつかの実施形態では、ベクトルプロセッサは、単に命令をパイプライン化するだけでなく、データ自体をパイプライン化することができる。パイプライン化は、命令、またはベクトルプロセッサの場合にはデータ自体が、複数のサブユニットを順に通過するプロセスを含むことができる。いくつかの実施形態では、ベクトルプロセッサは、数のベクトルまたは行列に対して同時に算術演算を指令する命令を供給される。連続的に命令を復号してから、それらの命令を完了するために必要なデータをフェッチしなければならない代わりに、ベクトルプロセッサは、メモリから単一の命令を読み出し、命令自体の定義では、命令が最後より1インクリメント大きいアドレスで別のデータ項目上で再度動作することが単に黙示される。これにより、復号時間を大幅に節約できる。

20

【0087】

図8は、本開示のいくつかの実施形態による、例示的なコンピューティングデバイス800のパーツ配置の一例を示す。コンピューティングデバイス800のパーツ配置の一例は、図1に示されるシステム100、図2に示されるシステム200、図4に示されるシステム400、図5に示されるシステム500、及び図6に示されるシステム600を含むことができる。コンピューティングデバイス800では、特定用途向けコンポーネント（例えば、図8の特定用途向けコンポーネント807を参照）は、AIコンポーネントであることができ、図1、2、4、5及び6にそれぞれ配置されて示される第一メモリチップ104または402及びアクセラレータチップ102または404だけでなく、図1、2、4、5及び6にそれぞれ構成されて示されるSoC106または406を含むことができる。コンピューティングデバイス800では、配線は、特定用途向けコンポーネントのコンポーネントを相互に直接接続する（例えば、図1～2及び図4～6にそれぞれ示される配線124及び424ならびに配線614を参照）。そして、コンピューティングデバイス800では、配線は、特定用途向けコンポーネントをSoCに直接接続する（例えば、特定用途向けコンポーネントをSoC806に直接接続する配線817を参照）。特定用途向けコンポーネントをSoCに直接接続する配線は、図1及び2に示されるような配線126、または図4及び5に示されるような配線426を含むことができる。また、特定用途向けコンポーネントをSoCに直接接続する配線は、図6に示されるような配線616を含むことができる。

30

40

【0088】

コンピューティングデバイス800は、図8に示されるようなコンピュータネットワーク802を介して他のコンピューティングデバイスに通信可能に結合されることができる。コンピューティングデバイス800は、少なくともバス804（メモリバスとペリフェ

50

ラルバスの組み合わせなど、1つ以上のバスであることができる)、S o C 8 0 6 (S o C 1 0 6 または 4 0 6 である、またはそれを含むことができる)、特定用途向けコンポーネント 8 0 7 (アクセラレータチップ 1 0 2 及び第一メモリチップ 1 0 4 または第一メモリチップ 4 0 2 及びアクセラレータチップ 4 0 4 であることができる)、及びメインメモリ 8 0 8 (メモリ 2 0 4 である、またはそれを含むことができる)だけでなく、ネットワークインタフェース 8 1 0 及びデータストレージシステム 8 1 2 も含む。バス 8 0 4 は、S o C 8 0 6、メインメモリ 8 0 8、ネットワークインタフェース 8 1 0、及びデータストレージシステム 8 1 2 を通信可能に結合する。そして、バス 8 0 4 は、バス 2 0 2、及び/または配線 1 2 6、4 2 6、または 6 1 6 などのポイントツーポイントメモリ接続を含むことができる。コンピューティングデバイス 8 0 0 は、コンピュータシステムを含み、このコンピュータシステムは、少なくとも、S o C 8 0 6 内の1つ以上のプロセッサ、メインメモリ 8 0 8 (例えば、読み出し専用メモリ (R O M)、フラッシュメモリ、同期 D R A M (S D R A M) または R a m b u s D R A M (R D R A M) などの D R A M、N V R A M、S R A M など)、及びデータストレージシステム 8 1 2 を含み、これらは、バス 8 0 4 (1つ以上のバス及び配線を含むことができる)を介して相互に通信する。

10

【 0 0 8 9 】

メインメモリ 8 0 8 (メモリ 2 0 4 である、それを含む、またはそれに含まれることができる)は、図 1 0 に示されるメモリストリング 1 0 0 0 を含むことができる。また、メインメモリ 8 0 8 は、図 1 1 に示されるメモリストリング 1 1 0 0 を含むことができる。いくつかの実施形態では、データストレージシステム 8 1 2 は、メモリストリング 1 0 0 0 またはメモリストリング 1 1 0 0 を含むことができる。

20

【 0 0 9 0 】

S o C 8 0 6 は、マイクロプロセッサ、C P U などのような1つ以上の汎用処理デバイスを含むことができる。また、S o C 8 0 6 は、G P U、A S I C、F P G A、デジタルシグナルプロセッサ (D S P)、ネットワークプロセッサ、プロセッサインメモリ (P I M) などのような1つ以上の専用処理デバイスを含むことができる。S o C 8 0 6 は、複合命令セットコンピューティング (C I S C) マイクロプロセッサ、縮小命令セットコンピューティング (R I S C) マイクロプロセッサ、超長命令語 (V L I W) マイクロプロセッサを有する1つ以上のプロセッサ、または他の命令セットを実施するプロセッサ、または命令セットの組み合わせを実施するプロセッサとすることができる。S o C 8 0 6 のプロセッサは、本明細書で論じられる動作及びステップを遂行するための命令を実行するように構成することができる。S o C 8 0 6 はさらに、1つ以上の通信ネットワーク(例えば、ネットワーク 8 0 2)を介して通信するために、ネットワークインタフェースデバイス、例えばネットワークインタフェース 8 1 0 を含むことができる。

30

【 0 0 9 1 】

データストレージシステム 8 1 2 は、本明細書で説明する方法または機能のうちのいずれか1つ以上を具現化する1つ以上の命令セットまたはソフトウェアが記憶されるマシン可読記憶媒体(コンピュータ可読媒体としても知られている)を含むことができる。また命令は、コンピュータシステムがそれを実行する間に、メインメモリ 8 0 8 内に、及び/または S o C 8 0 6 のプロセッサのうちの1つ以上の内に、完全に、または少なくとも部分的に存在することができる。またメインメモリ 8 0 8 及び S o C 8 0 6 の1つ以上のプロセッサ 5 0 6 はマシン可読記憶媒体を構成する。

40

【 0 0 9 2 】

メモリ、プロセッサ、及びデータ記憶装置部分を、例示的な実施形態においてそれぞれ単一部分であると示しているが、各部分は、命令を格納してそのそれぞれの動作を実行できる単一部分または複数部分を含むと解釈されるべきである。また用語「マシン可読記憶媒体」には、任意の媒体であって、マシンが実行するように命令のセットを記憶または符号化することができ、本開示の方法のいずれか1つ以上をマシンに行わせる媒体が含まれると解釈すべきである。したがって、用語「マシン可読記憶媒体」は、ソリッドステートメモリ、光媒体、及び磁気媒体を含むが、これらに限定されるものではないと解釈される

50

ものとする。

【0093】

図9は、本開示のいくつかの実施形態による、例示的なコンピューティングデバイス900のパーツ配置の別の例を示す。コンピューティングデバイス900のパーツ配置の例は、図3に示されるシステム300、及び図7に示されるシステム700を含むことができる。コンピューティングデバイス900では、特定用途向けコンポーネント（例えば、図9の特定用途向けコンポーネント807を参照）は、AIコンポーネントであることができ、図3及び7にそれぞれ配置されて示されるような第一メモリチップ104または402及びアクセラレータチップ102または404だけでなく、図3及び7にそれぞれ構成されて示されるようなSOC106または406も含むことができる。コンピューティングデバイス900では、配線は、特定用途向けコンポーネントのコンポーネントを相互に直接接続する（例えば、図3及び7にそれぞれ示される配線124及び424を参照）。ただし、コンピューティングデバイス900では、配線は、特定用途向けコンポーネントをSOCに直接接続しない。代替に、コンピューティングデバイス900では、1つ以上のバスは、特定用途向けコンポーネントをSOCに接続する（例えば、図9に構成されて示されるバス804、ならびに図3及び7に構成されて示されるバス202を参照）。

10

【0094】

図8及び9に示されるように、デバイス800及び900は、多くの同様のコンポーネントを含む。コンピューティングデバイス900は、図9に示されるようなコンピュータネットワーク802を介して他のコンピューティングデバイスに通信可能に結合されることができる。同様に、図9に示されるように、コンピューティングデバイス900は、少なくとも、バス804（メモリバス及びペリフェラルバスの組み合わせなど、1つ以上のバスであることができる）、SOC806（SOC106もしくは406である、またはそれを含むことができる）、特定用途向けコンポーネント807（アクセラレータチップ102及び第一メモリチップ104または第一メモリチップ402及びアクセラレータチップ404であることができる）、及びメインメモリ808（メモリ204である、またはそれを含むことができる）だけでなく、ネットワークインタフェース810及びデータストレージシステム812も含む。同様に、バス804は、SOC806、メインメモリ808、ネットワークインタフェース810、及びデータストレージシステム812を通信可能に結合する。そして、バス804は、バス202、及び/または配線126、426、または616などのポイントツーポイントメモリ接続を含むことができる。

20

30

【0095】

上述のように、本明細書に開示される少なくともいくつかの実施形態は、メモリ階層及びメモリチップストリングを使用してメモリを形成することに関する。

【0096】

図10及び11は、それぞれメモリチップストリング1000及び1100の例を示し、これらは、図2～3及び図5～7に示される別個のメモリ（すなわち、メモリ204）で使用されることができる。

【0097】

図10では、メモリチップストリング1000は、第一メモリチップ1002及び第二メモリチップ1004を含む。第一メモリチップ1002は、第二メモリチップ1004に直接配線され（例えば、配線1022を参照）、第二メモリチップと直接インタラクトするように構成される。メモリチップストリング1000内の各チップは、このストリング内の上流チップ及び/または下流チップに接続するために、1セット以上のピンを含むことができる（例えば、ピン1012及び1014のセットを参照）。いくつかの実施形態では、メモリチップストリング1000内の各チップは、ICパッケージ内に封入される単一のICを含むことができる。

40

【0098】

図10に示されるように、1セットのピン1012は第一メモリチップ1002の一部であり、配線1022、及び第二メモリチップ1004の一部である1セットのピン10

50

14を介して第一メモリチップ1002を第二メモリチップ1004に接続する。配線1022は、2セットのピン1012及び1014を接続する。

【0099】

いくつかの実施形態では、第二メモリチップ1004は、ストリング1000内のチップの中で最も低いメモリ帯域幅を有することができる。それらのような実施形態及び他の実施形態では、第一メモリチップ1002は、ストリング1000内のチップの中で最も高いメモリ帯域幅を有することができる。いくつかの実施形態では、第一メモリチップ1002は、DRAMチップである、またはそれを含む。いくつかの実施形態では、第一メモリチップ1002は、NVRAMチップである、またはそれを含む。いくつかの実施形態では、第二メモリチップ1004は、DRAMチップである、またはそれを含む。いくつかの実施形態では、第二メモリチップ1004は、NVRAMチップである、またはそれを含む。そして、いくつかの実施形態では、第二メモリチップ1004は、フラッシュメモリチップである、またはそれを含む。

10

【0100】

図11では、メモリチップストリング1100は、第一メモリチップ1102、第二メモリチップ1104、及び第三メモリチップ1106を含む。第一メモリチップ1102は、第二メモリチップ1104に直接配線され（例えば、配線1122を参照）、第二メモリチップと直接インタラクトするように構成される。第二メモリチップ1104は、第三メモリチップ1106に直接配線され（例えば、配線1124を参照）、第三メモリチップと直接インタラクトするように構成される。それらのような方法では、第一及び第三メモリチップ1102及び1106は、第二メモリチップ1104を介して間接的に相互にインタラクトする。

20

【0101】

メモリチップストリング1100内の各チップは、このストリング内の上流チップ及び/または下流チップに接続するために、1セット以上のピンを含むことができる（例えば、ピン1112、1114、1116、及び1118のセットを参照）。いくつかの実施形態では、メモリチップストリング1100内の各チップは、ICパッケージ内に封入される単一のICを含むことができる。

【0102】

図11に示されるように、1セットのピン1112は第一メモリチップ1102の一部であり、配線1122、及び第二メモリチップ1104の一部である1セットのピン1114を介して、第一メモリチップ1102を第二メモリチップ1104に接続する。配線1122は、2セットのピン1112及び1114を接続する。また、1セットのピン1116は、第二メモリチップ1104の一部であり、配線1124、及び第三メモリチップ1106の一部である1セットのピン1118を介して、第二メモリチップ1104を第三メモリチップ1106に接続する。配線1124は、2セットのピン1116及び1118を接続する。

30

【0103】

いくつかの実施形態では、第三メモリチップ1106は、ストリング1100内のチップの中で最も低いメモリ帯域幅を有することができる。それらのような実施形態及び他の実施形態では、第一メモリチップ1102は、ストリング1100内のチップの中で最も高いメモリ帯域幅を有することができる。また、それらのような実施形態及び他の実施形態では、第二メモリチップ1104は、ストリング1100内のチップの中でその次に最も高いメモリ帯域幅を有することができる。いくつかの実施形態では、第一メモリチップ1102は、DRAMチップである、またはそれを含む。いくつかの実施形態では、第一メモリチップ1102は、NVRAMチップである、またはそれを含む。いくつかの実施形態では、第二メモリチップ1104は、DRAMチップである、またはそれを含む。いくつかの実施形態では、第二メモリチップ1104は、NVRAMチップである、またはそれを含む。いくつかの実施形態では、第二メモリチップ1104は、フラッシュメモリチップである、またはそれを含む。いくつかの実施形態では、第三メモリチップ1106

40

50

また、例えば、メモリチップストリングの一実施形態は、DRAMからDRAMからNVRAM、またはDRAMからNVRAMからNVRAM、またはDRAMからフラッシュメモリからフラッシュメモリを含むことができる。ただし、DRAMからNVRAMからフラッシュメモリは、マルチティアメモリとして柔軟にプロビジョニングされるメモリチップストリングに、より効果的なソリューションを提供することができる。

【0110】

また、本開示の目的のために、DRAM、NVRAM、3D XPointメモリ、及びフラッシュメモリが個々のメモリユニットのための技法であること、そして本明細書に記載のメモリチップのいずれか1つのためのメモリチップがコマンド及びアドレスの復号のための論理回路、ならびにDRAM、NVRAM、3D XPointメモリ、またはフラッシュメモリのメモリユニットアレイを含むことができることを理解されたい。例えば、本明細書で説明されるDRAMチップは、コマンド及びアドレスの復号のための論理回路、ならびにDRAMのメモリユニットアレイを含む。例えば、本明細書で説明されるNVRAMチップは、コマンド及びアドレスの復号のための論理回路、ならびにNVRAMのメモリユニットアレイを含む。例えば、本明細書で説明されるフラッシュメモリチップは、コマンド及びアドレスの復号のための論理回路、ならびにフラッシュメモリのメモリユニットアレイを含む。

10

【0111】

また、本明細書で説明されるメモリチップのいずれか1つのためのメモリチップは、着信及び/または発信データのためのキャッシュまたはバッファメモリを含むことができる。いくつかの実施形態では、キャッシュまたはバッファメモリを実装するメモリユニットは、キャッシュまたはバッファメモリをホストするチップ上のユニットとは異なってもよい。例えば、キャッシュまたはバッファメモリを実装するメモリユニットは、SRAMのメモリユニットであることができる。

20

【0112】

前述の明細書では、本開示の実施形態は、その特定の例示的な実施形態を参照して説明されてきた。以下の請求項に述べる本開示の実施形態のより広い趣旨及び範囲から逸脱することなく、様々な変更を加えることができることが明らかである。したがって、明細書及び図面は限定的な意味ではなく例示的な意味で考慮されるべきである。

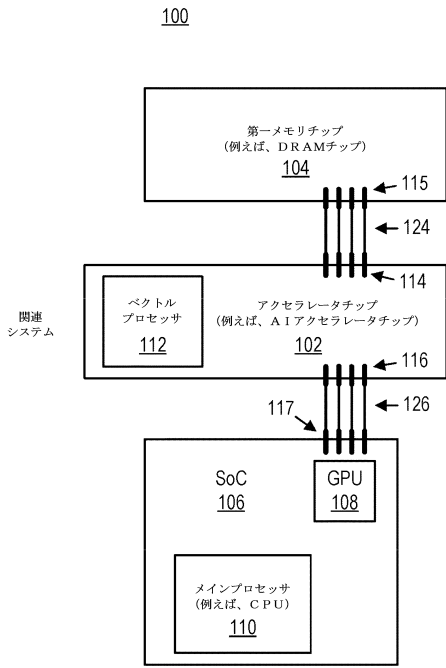
30

40

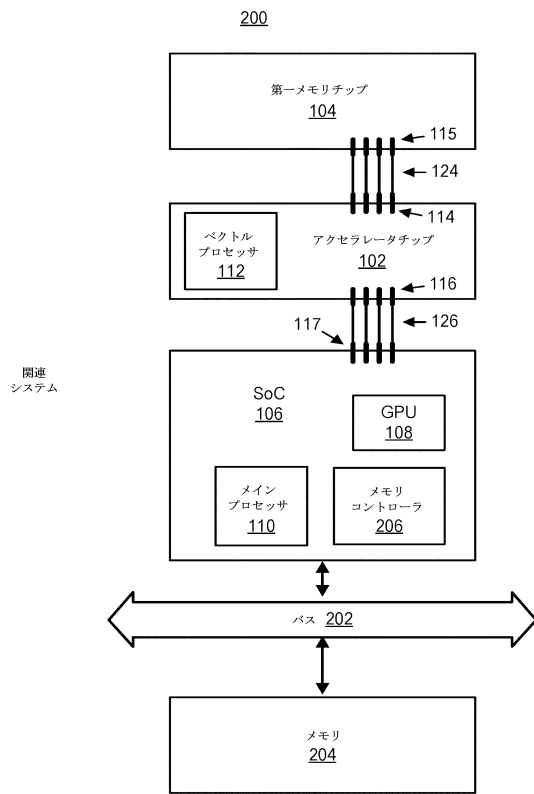
50

【図面】

【図 1】



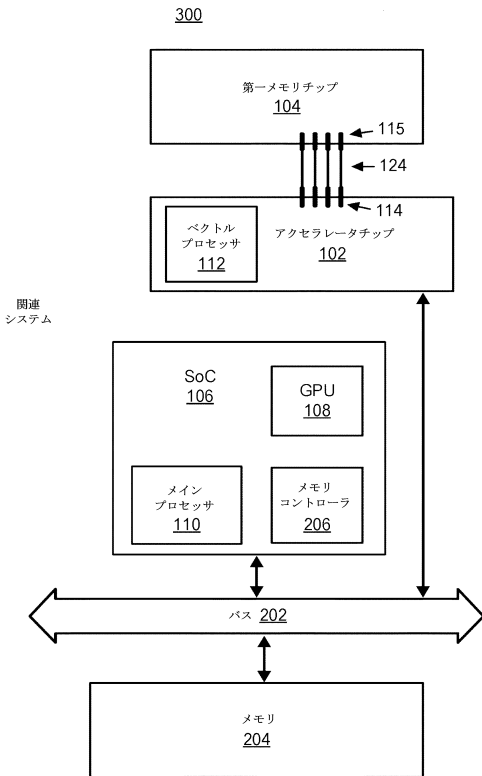
【図 2】



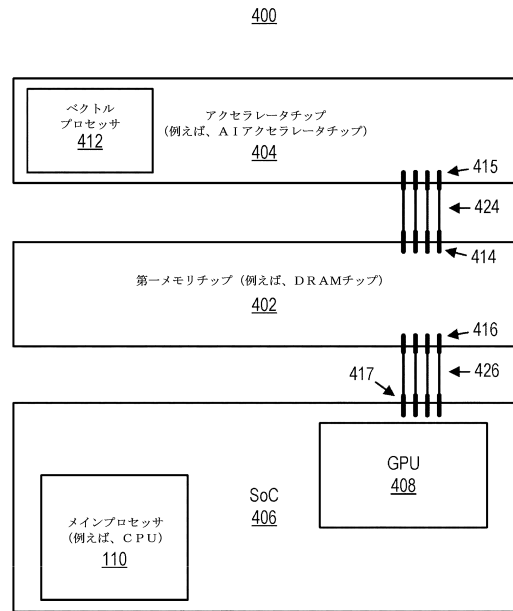
10

20

【図 3】



【図 4】

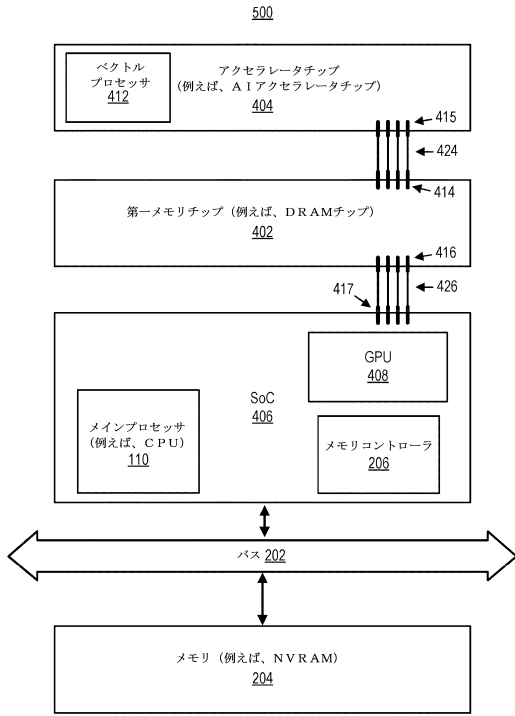


30

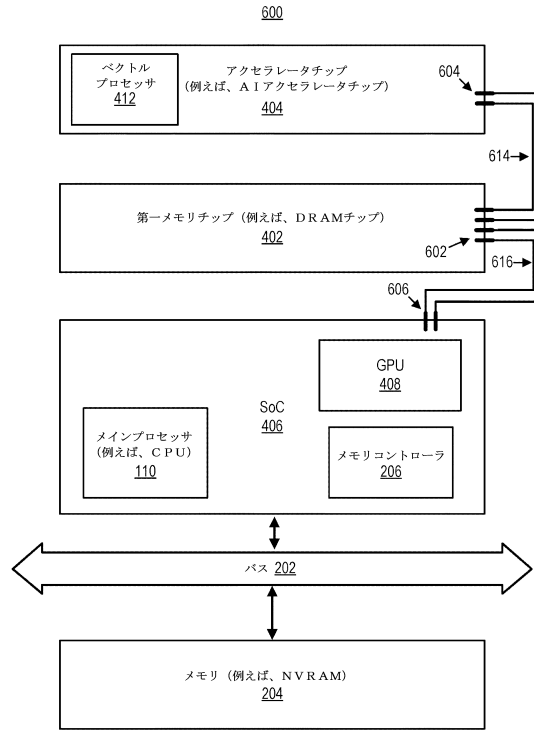
40

50

【図 5】



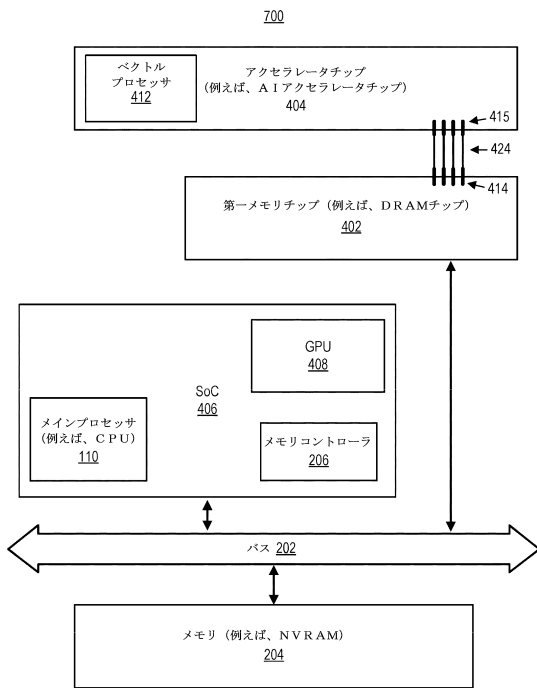
【図 6】



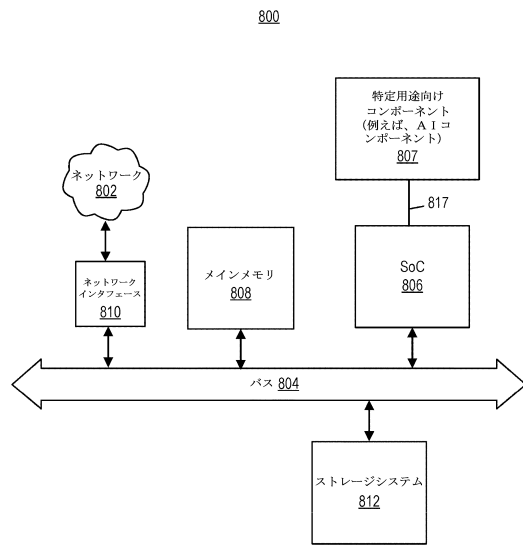
10

20

【図 7】



【図 8】

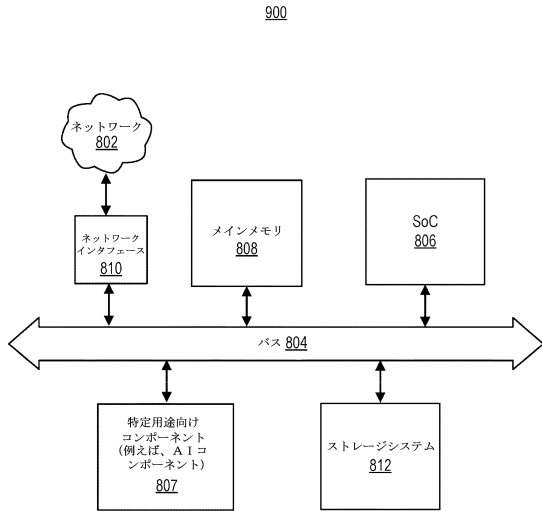


30

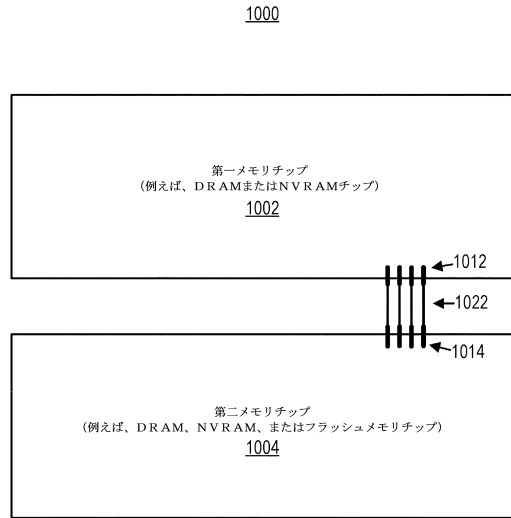
40

50

【 図 9 】

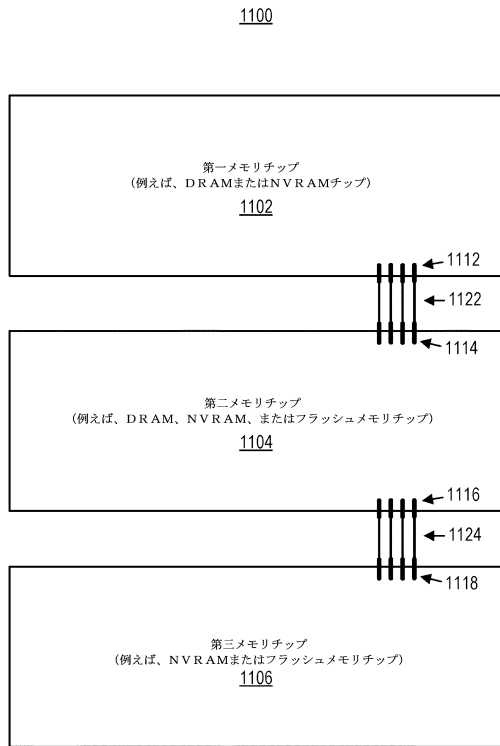


【 図 10 】



10

【 図 11 】



20

30

40

50

フロントページの続き

(72)発明者 クレウィッツ ケネス マリオン
アメリカ合衆国 カリフォルニア州 9 5 6 8 2 キャメロン パーク アッシュランド コート 4 0 8

(72)発明者 エノ ジャスティン エム .
アメリカ合衆国 カリフォルニア州 9 5 7 6 2 エル ドラド ヒルズ イエローストーン レーン
3 8 4 4

合議体

審判長 吉田 美彦

審判官 須田 勝巳

審判官 大塚 俊範

(56)参考文献 米国特許出願公開第 2 0 1 9 / 0 1 1 4 5 3 4 (U S , A 1)
特表 2 0 1 9 - 5 2 5 2 7 7 (J P , A)

(58)調査した分野 (Int.Cl. , D B 名)

G06F12/00

G06F13/16

G06F15/167

G06N 3/063