

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7548421号
(P7548421)

(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 F 9/50 (2006.01) G 0 6 F 9/50 1 5 0 C

請求項の数 8 (全17頁)

(21)出願番号	特願2023-515978(P2023-515978)	(73)特許権者	000004226 日本電信電話株式会社 東京都千代田区大手町一丁目5番1号
(86)(22)出願日	令和3年4月22日(2021.4.22)	(74)代理人	110002147 弁理士法人酒井国際特許事務所
(86)国際出願番号	PCT/JP2021/016313	(72)発明者	横野 智也 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
(87)国際公開番号	WO2022/224409	(72)発明者	山部 芳朗 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
(87)国際公開日	令和4年10月27日(2022.10.27)	(72)発明者	石崎 晃朗 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
審査請求日	令和5年8月18日(2023.8.18)	審査官	田中 幸雄

最終頁に続く

(54)【発明の名称】 アクセラレータ制御システム、アクセラレータ制御方法およびアクセラレータ制御プログラム

(57)【特許請求の範囲】

【請求項1】

アクセラレータ制御装置と、複数のアクセラレータとを有し、
前記アクセラレータ制御装置は、
演算処理の対象のデータの所在と、演算処理を行うアクセラレータの情報と、該アクセラレータの演算処理の内容を指定する情報とを含む制御用データを記憶する記憶部と、
前記アクセラレータにより演算処理済みの制御用データが前記記憶部に格納されている場合に、該アクセラレータによる演算処理の完了を判定する判定部と、を有し、
前記アクセラレータは、
前記記憶部から前記制御用データを取得する取得部と、
取得された前記制御用データに含まれる前記演算処理の対象のデータの所在とアクセラレータの演算処理の内容を指定する情報に従って、前記演算処理の対象のデータに対する演算処理を行う演算部と、
前記演算処理が完了した場合であって、前記制御用データに、転送先のアクセラレータの情報が含まれていない場合に、前記制御用データを前記記憶部に格納する格納部と、
前記制御用データに、転送先のアクセラレータの情報が含まれている場合に、該アクセラレータに該制御用データを転送する転送部と、
他のアクセラレータから前記制御用データを受信した場合に、該制御用データに含まれる前記演算処理の対象のデータの所在と、該演算処理の内容を指定する情報に従って、前記演算処理の対象のデータに対する演算処理を行う第2の演算部と、を有する

10

20

ことを特徴とするアクセラレータ制御システム。

【請求項 2】

他のアクセラレータからの前記制御用データの受信可否を確認する確認部を、さらに有することを特徴とする請求項 1 に記載のアクセラレータ制御システム。

【請求項 3】

前記転送部は、所定の時間内に転送先の前記アクセラレータから前記制御用データの受信可否が通知されない場合に、前記制御用データを再送することを特徴とする請求項 2 に記載のアクセラレータ制御システム。

【請求項 4】

前記格納部は、転送先の前記アクセラレータに前記制御用データが転送された場合に、前記演算処理対象のデータに対する前記演算処理の結果を前記記憶部に格納し、

10

転送先の前記アクセラレータにおいて、前記第 2 の演算部が、前記記憶部に格納された前記演算処理の結果を読み込んで前記演算処理の対象とする、

ことを特徴とする請求項 1 に記載のアクセラレータ制御システム。

【請求項 5】

前記記憶部に前記制御用データが有るか否かを監視する監視部をさらに有することを特徴とする請求項 1 に記載のアクセラレータ制御システム。

【請求項 6】

前記判定部は、前記記憶部に前記演算処理済みの制御用データが有るか否かを監視することを特徴とする請求項 1 に記載のアクセラレータ制御システム。

20

【請求項 7】

アクセラレータ制御装置と、複数のアクセラレータとを有するアクセラレータ制御システムが実行するアクセラレータ制御方法であって、

前記アクセラレータ制御システムは、演算処理の対象のデータの所在と、演算処理を行うアクセラレータの情報と、該アクセラレータの演算処理の内容を指定する情報とを含む制御用データを記憶する記憶部を有し、

前記記憶部から前記制御用データを取得する取得工程と、

取得された前記制御用データに含まれる前記演算処理の対象のデータの所在とアクセラレータの動作を指定する情報に従って、前記演算処理の対象のデータに対する演算処理を行う演算工程と、

30

前記演算処理が完了した場合であって、前記制御用データに、転送先のアクセラレータの情報が含まれていない場合に、前記制御用データを前記記憶部に格納する格納工程と、

演算処理済みの制御用データが前記記憶部に格納されている場合に、演算処理の完了を判定する判定工程と、

前記制御用データに、転送先のアクセラレータの情報が含まれている場合に、該アクセラレータに該制御用データを転送する転送工程と、

他のアクセラレータから前記制御用データを受信した場合に、該制御用データに含まれる前記演算処理の対象のデータの所在と、該演算処理の内容を指定する情報に従って、前記演算処理の対象のデータに対する演算処理を行う第 2 の演算工程と、

を含むことを特徴とするアクセラレータ制御方法。

40

【請求項 8】

コンピュータを請求項 1 ~ 6 のいずれか 1 項に記載のアクセラレータ制御システムとして機能させるためのアクセラレータ制御プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、アクセラレータ制御システム、アクセラレータ制御方法およびアクセラレータ制御プログラムに関する。

【背景技術】

【0002】

50

従来、演算処理システムに複数の演算加速器（アクセラレータ）を組み込んで処理を高速化する技術が知られている。半導体の微細化技術の限界により、プロセッサの動作周波数の向上や演算機の高集積化が困難となってきたことから、この技術により、特定の処理に特化してアクセラレータに処理をオフロードして、高性能な計算を実現する。さらに、複数台のアクセラレータを協調して動作させることにより、複数の演算で構成される処理の高速化を実現する（非特許文献1～3参照）。

【0003】

具体的には、あるアクセラレータの演算結果を別のアクセラレータに送信するアクセラレータ間通信を制御することが必要となる。そこで、転送元アクセラレータにおける演算処理の終了後に、CPUがアクセラレータ間のデータ通信の制御を行う。CPUは、転送先アクセラレータにデータ受付の準備を要求する制御信号を送信した後に、転送元アクセラレータに転送開始の制御信号を送信することで、アクセラレータ間のデータ通信を駆動する。アクセラレータ間のデータ通信の後、転送先アクセラレータがデータに対する演算処理を行う。

【先行技術文献】

【非特許文献】

【0004】

【文献】J. Yang, D. B. Minton, F. Hady, “When Poll is Better than Interrupt”, FAST, vol.12, pp.3-3, [online], 2012年, [2021年3月29日検索], インターネット URL: <https://www.usenix.org/system/files/conference/fast12/yang.pdf>

【文献】“CUDA Toolkit Documentation v11.2.2”, [online], NVIDIA, [2021年3月29日検索], インターネット URL: <https://docs.nvidia.com/cuda/>

【文献】K. Gulati, S. P. Khatri, “GPU Architecture and the CUDA Programming Model”, Hardware acceleration of EDA algorithms, pp.23-30, [online], Springer, [2021年3月29日検索], インターネット URL: https://link.springer.com/content/pdf/10.1007%2F978-1-4419-0944-2_3.pdf

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、従来技術では、複数のアクセラレータを用いてリアルタイム処理を高速化することが困難な場合がある。例えば、アクセラレータはCPUによって制御されるため、アクセラレータ間のデータ転送の処理が発生するたびに、転送元アクセラレータ・転送先アクセラレータの排他制御や処理完了の割り込み通知、同期制御が発生する。これらの処理がオーバーヘッドとなり、遅延時間が増大する。また、多数のアクセラレータを用いる場合には、上記のようなオーバーヘッドに加え、バスの仕様によりアクセラレータ間で直接データ転送が不可能な場合が存在するため、スケーラビリティが低下する。したがって、アクセラレータをリアルタイム処理や多様な処理の領域で用いることが困難である。

【0006】

本発明は、上記に鑑みてなされたものであって、複数のアクセラレータを用いてリアルタイム処理を高速化することを目的とする。

【課題を解決するための手段】

【0007】

上述した課題を解決し、目的を達成するために、本発明に係るアクセラレータ制御システムは、アクセラレータ制御装置と、複数のアクセラレータとを有し、前記アクセラレータ制御装置は、演算処理の対象のデータの所在と、演算処理を行うアクセラレータの情報と、該アクセラレータの演算処理の内容を指定する情報とを含む制御用データを記憶する記憶部と、前記アクセラレータにより演算処理済みの制御用データが前記記憶部に格納されている場合に、該アクセラレータによる演算処理の完了を判定する判定部と、を有し、前記アクセラレータは、前記記憶部から前記制御用データを取得する取得部と、取得され

10

20

30

40

50

た前記制御用データに含まれる前記演算処理の対象のデータの所在とアクセラレータの演算処理の内容を指定する情報に従って、前記演算処理の対象のデータに対する演算処理を行う演算部と、前記演算処理が完了した場合であって、前記制御用データに、転送先のアクセラレータの情報が含まれていない場合に、前記制御用データを前記記憶部に格納する格納部と、前記制御用データに、転送先のアクセラレータの情報が含まれている場合に、該アクセラレータに該制御用データを転送する転送部と、他のアクセラレータから前記制御用データを受信した場合に、該制御用データに含まれる前記演算処理の対象のデータの所在と、該演算処理の内容を指定する情報に従って、前記演算処理の対象のデータに対する演算処理を行う第2の演算部と、を有することを特徴とする。

【発明の効果】

10

【0008】

本発明によれば、複数のアクセラレータを用いてリアルタイム処理を高速化することが可能となる。

【図面の簡単な説明】

【0009】

【図1】図1は、本実施形態のアクセラレータ制御システムの概略構成を例示する模式図である。

【図2】図2は、本実施形態のアクセラレータ制御システムを説明するための図である。

【図3】図3は、メタデータのデータ構成を例示する図である。

【図4】図4は、メタデータバッファの構造を例示する図である。

20

【図5】図5は、データバッファを説明するための図である。

【図6】図6は、アクセラレータ間通信を説明するための図である。

【図7】図7は、アクセラレータにおけるアクセラレータ制御処理手順を示すフローチャートである。

【図8】図8は、アクセラレータ制御プログラムを実行するコンピュータの一例を示す図である。

【発明を実施するための形態】

【0010】

以下、図面を参照して、本発明の一実施形態を詳細に説明する。なお、この実施形態により本発明が限定されるものではない。また、図面の記載において、同一部分には同一の符号を付して示している。

30

【0011】

[アクセラレータ制御システムの構成]

図1は、本実施形態のアクセラレータ制御システムの概略構成を例示する模式図である。図1に示すように、アクセラレータ制御システム1は、アクセラレータ制御装置10と複数のアクセラレータ20とを有する。アクセラレータ制御装置10とアクセラレータ20とは、バス配線等の通信回線を介して相互にデータ通信が可能である。

【0012】

アクセラレータ制御装置10は、CPU (Central Processing Unit) 等の汎用プロセッサを用いて実現される情報処理装置である。アクセラレータ制御装置10において、ユーザ空間上で、ユーザが定義した処理であるアプリケーションが実行される。またカーネル空間上で、アクセラレータ20等の物理デバイスの操作や管理、ユーザ空間処理における物理デバイスの抽象化処理が行われる。

40

【0013】

ユーザメモリ空間、カーネルメモリ空間は、RAM (Random Access Memory)、フラッシュメモリ (Flash Memory) 等の半導体メモリ素子、または、ハードディスク、光ディスク等の記憶装置によって実現され、記憶部として機能する。記憶部には、アクセラレータ制御装置10を動作させる処理プログラムや、処理プログラムの実行中に使用されるデータなどが予め記憶され、あるいは処理の都度一時的に記憶される。

【0014】

50

汎用プロセッサを用いて実現される制御部は、メモリに記憶された処理プログラムを実行することにより、図 1 に例示するように、デバイスドライバのメモリ管理機能 1 2、メタデータ制御機能 1 3、後述する API 1 4 等の各機能部として機能する。なお、制御部は、不図示の他の機能部を備えてもよい。

【 0 0 1 5 】

アクセラレータ制御装置 1 0 では、演算処理の対象のデータ（以下、処理データとも記す）が、ユーザ空間上のユーザメモリ空間のデータバッファ 1 1 a に配置される。また、演算処理の対象のデータの所在と、アクセラレータの演算処理の内容を指定する情報とを含む制御用データ（以下、メタデータとも記す）が、カーネル空間上のカーネルメモリ空間のメタデータバッファ 1 1 b に配置される。

10

【 0 0 1 6 】

アクセラレータ 2 0 は、GPU（Graphics Processing Unit）、FPGA（Field Programmable Gate Array）、DPU（Data Processing Unit）等で実現される演算加速器であり、図 1 に例示するように、メタデータアクセス機構 2 1、機能制御ブロック 2 2、データ転送機能 2 3、演算機能 2 4 等の各機能部として機能する。

【 0 0 1 7 】

アクセラレータ 2 0 は、処理データとメタデータとを読み込んで、演算処理のオフロードを実現する。具体的には、アクセラレータ 2 0 は、メタデータバッファ 1 1 b を自律的な Polling により、メタデータが配置されているか否かを監視している。アクセラレータ 2 0 は、メタデータが配置されている場合に、配置されているメタデータを読み込んで、解析してアクセラレータ 2 0 を駆動する。また、アクセラレータ 2 0 は、処理データに対する演算処理の結果をデータバッファ 1 1 a に格納する。

20

【 0 0 1 8 】

なお、アクセラレータ制御装置 1 0 は、アクセラレータ 2 0 とは非同期に、次の処理に進めることが可能である。アクセラレータ 2 0 は、演算処理が完了した場合にメタデータをメタデータバッファ 1 1 b に配置する。これにより、アクセラレータ制御装置 1 0 は、演算処理が完了したとみなす。また、アクセラレータ制御装置 1 0 は、データバッファ 1 1 a から処理データに対する演算処理の結果を取得して、次の処理に進める。

【 0 0 1 9 】

具体的には、アクセラレータ制御装置 1 0 は、デバイスドライバに実装されたメタデータ制御機能 1 3 を介して、カーネルメモリ空間上のメタデータバッファ 1 1 b に処理データの所在や演算処理の内容を含むメタデータを配置することで、アクセラレータ 2 0 の動作を決定する。

30

【 0 0 2 0 】

アクセラレータ制御装置 1 0 は、アクセラレータ制御装置 1 0 とアクセラレータ 2 0 とが参照可能なデータバッファ 1 1 a を、メモリ管理機能 1 2 を介してユーザメモリ空間に作成し、処理データを配置する。

【 0 0 2 1 】

アクセラレータ 2 0 は、メタデータアクセス機構 2 1 を有する。メタデータアクセス機構 2 1 は、カーネルメモリ空間上のメタデータバッファ 1 1 b を Polling する機能、メタデータの読込機能、メタデータバッファへのメタデータの書込機能を備え、アクセラレータ制御装置 1 0 からの明示的な制御信号によらず自律的にメタデータを読み込むことが可能である。

40

【 0 0 2 2 】

メタデータバッファ 1 1 b のアドレス等の情報は、デバイスドライバのロードの際に、アクセラレータ初期化機能によってアクセラレータ 2 0 に通知される。

【 0 0 2 3 】

アクセラレータ 2 0 は、メタデータバッファ 1 1 b からメタデータを読み込み、メタデータ解析機能 2 2 a がデータを解析し、その記述に従って動作する。アクセラレータ 2 0 は、メタデータに含まれる処理データのアドレスや長さの情報を元に、データ転送機能 2

50

3におけるデータ読込機能23aを駆動し、処理データをデータバッファ11aから読みだして、演算機能24での演算に用いる。演算機能24では、メタデータに記述されている演算情報を元に、演算処理が実行される。メタデータ解析機能22aは、解析結果を演算機能24に渡し、演算機能24による演算処理内容を設定する。

【0024】

演算処理の終了後には、演算処理の結果が演算機能24からデータ書込機能23bに転送され、データ転送制御機能22bにより指定されたデータバッファに書き込まれる。その後、データ転送が終了した旨の通知がデータ転送機能23から完了通知機能22cに送信され、メタデータアクセス機構21の書込機能によりメタデータがメタデータバッファ11bに書き込まれ、アクセラレータ20の動作が完了する。

10

【0025】

このように、アクセラレータ制御システム1では、データ転送や演算処理を制御するためのメタデータをメインメモリに配置して、アクセラレータ20が自律的にPollingおよび読み込みを行う。これにより、割込み等のアクセラレータ制御のオーバーヘッドを低減し、リアルタイム処理が可能となる。

【0026】

また、アクセラレータ20が処理データをアクセラレータ制御装置10のユーザ空間に直接転送可能であり、アクセラレータ制御装置10とアクセラレータ20とが非同期に駆動可能であるため、遅延もCPU使用率も削減される。

【0027】

さらに、アクセラレータ20がアクセラレータ制御装置10のメインメモリにアクセスして動作可能であるため、複数のアクセラレータ20への拡張が可能であり、スケーラビリティが大きく向上する。

20

【0028】

また、複数のアクセラレータ20を用いて演算処理が行われる場合には、メタデータには、さらに、演算を行うアクセラレータの情報、アクセラレータでのオペレーション、読み込む処理データの所在等の情報が含まれる。これらの情報を元に、転送元アクセラレータ20aが転送先アクセラレータ20bの設定や制御情報(以下、メタデータ、またはアクセラレータ間メタデータと記す)を作成し、アクセラレータ制御装置10を介さずにアクセラレータ間で直接、メタデータのデータ通信を行う。

30

【0029】

転送元アクセラレータ20aは、演算処理の終了後、演算処理結果をカーネルメモリ空間上の中間データバッファ11cに出力し、転送先アクセラレータ20bにメタデータを送信する。

【0030】

転送先アクセラレータ20bは、現在のアクセラレータ20bの状態から、送信されたメタデータの受信が可能か否かを判定する。メタデータを受信可能な場合には、転送先アクセラレータ20bは、転送元アクセラレータ20aに受信完了を通知する。転送先アクセラレータ20bから受信完了を通知されるまで、転送元アクセラレータ20aは待機状態であり、所定の時間以内に受信完了の通知がない場合には、再度、メタデータを転送先アクセラレータ20bに送信する。

40

【0031】

転送先アクセラレータ20bは、受信したメタデータを解析することにより、中間データバッファ11cの中間データを読み込んで、演算処理を実行する。

【0032】

具体的には、複数のアクセラレータ20で連携する場合には、メタデータに複数のアクセラレータ20で動作するための情報が記述されている。アクセラレータ20は、メタデータ解析機能22aで読み込んだメタデータを解析することで、演算処理終了後にアクセラレータ間のデータ通信を行う。

【0033】

50

まず、転送元アクセラレータ 20 a では、演算処理終了後に、カーネルメモリ空間上の中間データバッファ 11 c に、転送先アクセラレータ 20 b に転送するための処理データの演算結果を出力する。

【0034】

その後、転送元アクセラレータ 20 a は、転送先アクセラレータ通知機能 22 d により、バスを介して転送先アクセラレータ 20 b に、演算処理結果の位置情報や演算内容、アクセラレータの連携に関する情報を含むメタデータを送信する。

【0035】

転送先アクセラレータ 20 b は、アクセラレータ間メタデータ通信機構 21 d の受信機能により、メタデータを受け取る。メタデータは、解析機能に送信され、アクセラレータの状態に基づいて、メタデータを受信するか否かを判断する。その際に、アクセラレータ間メタデータ通信機構 21 d の受信通知機能が、受信したか否かの情報を生成し、転送元アクセラレータ 20 a に返送する。

【0036】

待機状態の転送元アクセラレータ 20 a は、転送先アクセラレータ 20 b から受信完了の通知を受け取ったことにより、転送が完了したとみなして、次の処理に移行する。所定の時間以内に受信完了の通知がない場合に、転送元アクセラレータ 20 a は、再度、メタデータを送信する。

【0037】

転送先アクセラレータ 20 b は、転送されたメタデータを受信可能であれば、アクセラレータ間メタデータ通信機構 21 d の解析機能から読み込んだメタデータを基に、演算処理を実行する。

【0038】

このように、アクセラレータ制御システム 1 は、複数のアクセラレータ 20 を用いた処理において、アクセラレータ制御装置 10 を介さずに実行されるため、転送元アクセラレータ / 転送先アクセラレータの排他処理や処理完了の割込み通知、同期制御等のオーバヘッドを低減して、スケーラビリティを低下させない。また、通信処理においても、アクセラレータ制御装置 10 は、各アクセラレータ 20 と非同期に駆動可能であるため、それぞれの演算リソースを有効に活用することが可能である。また、処理データの受け渡しにメインメモリを使用することにより、アクセラレータ制御装置 10 の制御がなくても複数のアクセラレータ 20 からデータソースを参照可能である。そのためバスの仕様に依らずにデータ転送が可能であり、スケーラブルかつ幅広いワークロードに対応したシステムを構築することが可能である。

【0039】

[アクセラレータ制御装置]

図 2 は、本実施形態のアクセラレータ制御システムを説明するための図である。図 2 に例示するように、アクセラレータ制御装置 10 において、API 14 は、アプリケーションから呼び出し可能な機能群であり、アクセラレータ 20 の動作に必要な機能を提供する。例えば、API は、データバッファ確保機能、データバッファ解放機能、メタデータ作成機能、メタデータ送信機能、演算完了感知機能、通信用メタデータ作成機能として機能する。

【0040】

データバッファ確保機能は、アクセラレータ 20 とデータの送受信を行うための記憶領域であるデータバッファ 11 a を生成する。データバッファ確保機能は、メモリ管理機能 12 を呼び出して、メモリ管理機能 12 からデータバッファ 11 a の領域情報を取得する。

【0041】

データバッファ解放機能は、データバッファ 11 a を削除する。データバッファ解放機能は、メモリ管理機能 12 を呼び出して、解放するデータバッファ 11 a の領域情報を通知して、当該領域を解放する。

【0042】

10

20

30

40

50

メタデータ作成機能は、アクセラレータ 20 の制御情報が付加された制御用データであるメタデータを作成する。メタデータ送信機能は、メタデータ作成機能が作成したメタデータをメタデータ制御機能 13 に送信する。演算完了感知機能は、メタデータ制御機能 13 を呼び出して、メタデータ制御機能 13 の返り値により、任意のアクセラレータ 20 の処理が終了しているかを確認する。通信用メタデータ作成機能は、アクセラレータ間通信を行うための通信路情報を作成する。

【0043】

デバイスドライバは、カーネル空間上に存在し、特定のデバイスを扱うためのメモリ管理機能 12、メタデータ制御機能 13 等の機能群を提供する。

【0044】

メモリ管理機能 12 は、ユーザ空間上の API 14 から呼び出され、データバッファ 11 a の領域の確保、解放、管理を行う。

【0045】

メタデータバッファ 11 b は、アクセラレータ 20 とメタデータの送受信を行うための記憶領域であり、アクセラレータ初期化機能により、アクセラレータ 20 の数に応じて生成される。例えば、メタデータバッファ 11 b は、演算処理の対象のデータの所在と、演算処理を行うアクセラレータ 20 の情報と、該アクセラレータ 20 の演算処理の内容を指定する情報とを含む制御用データを記憶する。

【0046】

メタデータバッファ 11 b は、アクセラレータ 20 にメタデータを送信する RQ (Request Queue) バッファと、アクセラレータ 20 からメタデータを受信する CQ (Completion Queue) バッファからなる。

【0047】

メタデータ制御機能 13 は、ユーザ空間上の API 14 から呼び出され、メタデータバッファ 11 b の RQ バッファ / CQ バッファへの読み込み、書き込み、管理を行う。メタデータ制御機能 13 は、Polling 機能 13 a と書込機能 13 b とを有する。

【0048】

Polling 機能 13 a は、CQ バッファを Polling し、演算の完了を確認する。すなわち、Polling 機能 13 a は、判定部として機能して、アクセラレータ 20 により演算処理済みのメタデータがメタデータバッファ 11 b に格納されているか否かを監視する。そして、Polling 機能 13 a は、アクセラレータ 20 により演算処理済みのメタデータがメタデータバッファ 11 b に格納されている場合に、該アクセラレータ 20 による演算処理の完了を判定する。また、書込機能 13 b は、RQ バッファに書き込みを行う。

【0049】

[アクセラレータ]

アクセラレータ 20 において、メタデータアクセス機構 21 は、カーネル空間上のメタデータバッファ 11 b にアクセスする機能として、Polling 機能 21 a、読込機能 21 b、書込機能 21 c、アクセラレータ間メタデータ通信機構 21 d を含む。

【0050】

Polling 機能 21 a は、メタデータバッファ 11 b の RQ バッファを Polling して、RQ バッファにメタデータが有るか否かを確認する。すなわち、Polling 機能 21 a は、監視部として機能して、メタデータバッファ 11 b にメタデータが有るか否かを監視する。

【0051】

読込機能 21 b は、取得部として機能して、メタデータバッファ 11 b からメタデータを取得する。すなわち、読込機能 21 b は、Polling 機能 21 a により RQ バッファにメタデータが有ることが確認された場合に、メタデータを読み込み、機能制御ブロック 22 に転送する。

【0052】

10

20

30

40

50

書込機能 2 1 c は、メタデータバッファ 1 1 b の C Q バッファにメタデータを書き込む。すなわち、書込機能 2 1 c は、格納部として機能して、演算処理が完了した場合であって、メタデータに、転送先のアクセラレータ 2 0 b の情報が含まれていない場合に、メタデータをメタデータバッファ 1 1 b に格納する。

【 0 0 5 3 】

機能制御ブロック 2 2 は、アクセラレータ 2 0 の各機能部に対する必要なデータの共有、駆動・休止タイミングを制御する。機能制御ブロック 2 2 は、メタデータ解析機能 2 2 a、データ転送制御機能 2 2 b、完了通知機能 2 2 c、転送先アクセラレータ通知機能 2 2 d を含む。

【 0 0 5 4 】

メタデータ解析機能 2 2 a は、メタデータアクセス機構 2 1 から送信されたメタデータから、オペレーション、データバッファ 1 1 a のアドレス・長さの情報を抽出する。データ転送制御機能 2 2 b は、データ転送機能 2 3 におけるデータの読み込み、書き込みを制御する機能を提供する。完了通知機能 2 2 c は、メタデータアクセス機構 2 1 に C Q バッファに書き込むメタデータを転送する。

【 0 0 5 5 】

演算機能 2 4 は、演算部として機能して、取得されたメタデータに含まれる演算処理の対象の処理データの所在とアクセラレータ 2 0 の演算処理の内容を指定する情報に従って、演算処理の対象の処理データに対する演算処理を行う。

【 0 0 5 6 】

具体的には、演算機能 2 4 は、入力制御機能、演算回路、出力制御機能を含む。演算回路は複数の独立した演算回路を有し、それぞれに対して演算を定義することが可能である。入力制御機能は、機能制御ブロック 2 2 から送信されるメタデータの情報を基に、データ転送機能 2 3 のデータ読込機能 2 3 a から送信されたしデータを適切な演算回路に入力する。

【 0 0 5 7 】

出力制御機能は、演算回路から出力された処理データに対する演算処理の結果を、適切なタイミングでデータ転送機能 2 3 のデータ書込機能 2 3 b に転送する。

【 0 0 5 8 】

データ転送機能 2 3 のデータ書込機能 2 3 b は、格納部として機能して、処理データに対する演算処理の結果をデータバッファ 1 1 a に格納する。

【 0 0 5 9 】

次に、図 3 は、メタデータのデータ構成を例示する図である。図 3 に示すように、メタデータは、3 2 b i t のタスク I D、3 2 b i t のオペレーション、3 2 b i t の読み出しデータの長さ、6 4 b i t の読み出しデータのアドレス、3 2 b i t の書き込みデータの長さ、6 4 b i t の書き込みデータのアドレスの計 2 5 6 b i t のデータで構成される。

【 0 0 6 0 】

また、図 4 は、メタデータバッファの構造を例示する図である。図 4 に示すように、図 3 に示したメタデータが、R Q バッファ / C Q バッファに格納される。図 4 に示す R Q バッファの h e a d 情報は、メタデータ制御機能 1 3 により制御される。メタデータ制御機能 1 3 は、A P I 1 4 のメタデータ送信機能から送信されたメタデータを、図 3 に示した構成に成形し、R Q バッファに格納した後、h e a d の情報を次に位置にずらす。

【 0 0 6 1 】

同様に、メタデータ制御機能 1 3 は、C Q バッファにおいても、h e a d の情報を用いて、A P I 1 4 の演算完了感知機能に呼び出された場合に、C Q バッファからメタデータを読みだして、h e a d を次の位置にずらす。

【 0 0 6 2 】

また、図 5 は、データバッファを説明するための図である。メモリ管理機能 1 2 は、図 5 に例示するテーブルで、データバッファ 1 1 a のアドレスと状態とを管理する。図 5 に示すように、テーブルでは、データバッファ 1 1 a の物理アドレスとそのアドレスの状態

10

20

30

40

50

とが管理されている。アクセラレータ 20 に対するメタデータをセットする際に、key を参照することにより、所望のデータバッファ 11 a の物理アドレスを取得して、読み出しデータアドレスと書き出しデータアドレスとをセットする。これにより、カーネル空間およびユーザ空間における仮想アドレスの違いを吸収して管理できる。

【0063】

図 2 の説明に戻る。上記した書込機能 21 c は、転送部として機能して、メタデータに、転送先のアクセラレータ 20 b の情報が含まれている場合に、該アクセラレータ 20 b に該メタデータを転送する。

【0064】

アクセラレータ間メタデータ通信機構 21 d は、受信機能、解析機能、受信通知機能を含み、アクセラレータ間通信に必要な機能を提供する。受信機能は、確認部として機能して、他のアクセラレータ 20 a からのメタデータの受信可否を確認する。すなわち、受信機能は、他のアクセラレータ 20 (転送元アクセラレータ 20 a) から送られてきたメタデータの受信可否を通知するためのデータを生成し、転送元アクセラレータ 20 a に送信する。

10

【0065】

機能制御ブロック 22 の転送先アクセラレータ通知機能 22 d は、転送先アクセラレータ 20 b に、メタデータを通知する機能を提供する。メタデータ解析機能 22 a は、メタデータアクセス機構 21 から送られたメタデータから、オペレーション、データバッファのアドレス・長さの情報を抽出する。加えて、メタデータ解析機能 22 a は、メタデータから転送先アクセラレータ 20 b の情報を抽出する。また、メタデータ解析機能 22 a は、転送先アクセラレータ 20 b からの受信完了の通知がない場合に、転送先アクセラレータ通知機能 22 d および書込機能 21 c を介してメタデータを再送する。書込機能 21 c が、所定の時間内に転送先のアクセラレータ 20 b からメタデータの受信可否が通知されない場合に、メタデータを再送する。

20

【0066】

また、転送先のアクセラレータ 20 b において、演算機能 24 は、第 2 の演算部として機能して、後述するように、他のアクセラレータ 20 a からメタデータを受信した場合に、演算機能 24 が、メタデータに含まれる演算処理の対象の中間データの所在と、該演算処理の内容を指定する情報に従って、演算処理の対象の中間データに対する演算処理を行う。

30

【0067】

データ転送機能 23 のデータ書込機能 23 b は、転送先のアクセラレータ 20 b にメタデータが転送された場合に、演算処理対象の処理データに対する演算処理の結果の中間データを中間データバッファ 11 c に格納する。転送先のアクセラレータ 20 b では、演算機能 24 が、中間データバッファ 11 c に格納された演算処理の結果を読み込んで演算処理の対象とする。

【0068】

ここで、図 6 は、アクセラレータ間通信を説明するための図である。アクセラレータ制御システム 1 では、メタデータを用いて、アクセラレータ間通信を実現する。図 6 (a) に示すように、アクセラレータ間通信の通信路は、各アクセラレータ 20 を頂点 V_0 、 V_1 、 V_2 、... とするグラフで表現する。このグラフを隣接行列として表現することにより、データとして取り扱う。例えば、図 6 (a) に例示したグラフは、図 6 (b) に示すように、隣接行列として表現できる。

40

【0069】

これらの情報は、図 6 (c) に示すデータ構造にすることで、アクセラレータ 20 で扱うことができる。図 6 (c) に示すデータ構造では、先頭の 8 bit は行列の大きさを表し、アクセラレータ 20 が行列本体のデータの大きさを知ることができる。続く 256 bit は行列本体のデータであり、アクセラレータ間の接続情報が含まれる。続くフィールドには、転送先アクセラレータ 20 b を駆動するための演算内容や処理データの所在等の

50

情報が記述される。

【 0 0 7 0 】

[アクセラレータ制御処理]

次に、図 7 は、アクセラレータ 2 0 におけるアクセラレータ制御処理手順を示すフローチャートである。図 7 のフローチャートは、例えば、アプリケーションにより開始が指示されたタイミングで開始される。

【 0 0 7 1 】

図 7 のフローチャートは、例えば、所定の間隔で開始される。アクセラレータ 2 0 は、メタデータバッファ 1 1 b の R Q バッファを P o l l i n g により監視している (ステップ S 1)。また、アクセラレータ 2 0 は、他のアクセラレータ 2 0 からメタデータが転送されてきているか否かを確認する (ステップ S 2)。

10

【 0 0 7 2 】

他のアクセラレータ 2 0 からメタデータが転送されていない場合には (ステップ S 2、N o)、アクセラレータ 2 0 は、ステップ S 3 に処理を進める。一方、他のアクセラレータ 2 0 からメタデータが転送されている場合には (ステップ S 2、Y e s)、アクセラレータ 2 0 (転送先アクセラレータ 2 0 b) は、転送元アクセラレータ 2 0 a に、メタデータの受信完了の通知を送信するとともに、メタデータを解析し (ステップ S 4)、ステップ S 5 に処理を進める。

【 0 0 7 3 】

ステップ S 3 において、R Q バッファにメタデータがセットされていない場合には (ステップ S 3、N o)、アクセラレータ 2 0 は、ステップ S 1 に処理を戻す。一方、R Q バッファにメタデータがセットされている場合には (ステップ S 3、Y e s)、アクセラレータ 2 0 は、メタデータを読み込んで解析する (ステップ S 5)。

20

【 0 0 7 4 】

また、アクセラレータ 2 0 は、解析結果を基に、データバッファ 1 1 a の処理データを取得して (ステップ S 6)、演算処理を行う (ステップ S 7)。

【 0 0 7 5 】

また、転送先のアクセラレータ 2 0 が存在しない場合には (ステップ S 8、N o)、アクセラレータ 2 0 は、演算処理済みの処理データをデータバッファ 1 1 a に書き込み (ステップ S 9)、C Q バッファにメタデータを送信して (ステップ S 1 0)、処理を完了する。その後、アクセラレータ 2 0 は、ステップ S 1 に処理を戻す。

30

【 0 0 7 6 】

一方、転送先のアクセラレータ 2 0 が存在する場合には (ステップ S 8、Y e s)、アクセラレータ 2 0 は、転送元アクセラレータ 2 0 a として、処理データの演算処理結果を中間データバッファ 1 1 c に書き出し (ステップ S 1 1)、転送先アクセラレータ 2 0 b にメタデータを送信する (ステップ S 1 2)。そして、アクセラレータ 2 0 は、転送先アクセラレータ 2 0 b から受信完了の通知を受け取った場合に、メタデータの転送が完了したものとみなし、ステップ S 1 に処理を戻す。

【 0 0 7 7 】

[効果]

以上、説明したように、本実施形態のアクセラレータ制御システム 1 において、アクセラレータ制御装置 1 0 では、メタデータバッファ 1 1 b が、演算処理の対象のデータの所在と、演算処理を行うアクセラレータ 2 0 の情報と、該アクセラレータ 2 0 の演算処理の内容を指定する情報とを含む制御用データを記憶する。P o l l i n g 機能 1 3 a が、アクセラレータ 2 0 により演算処理済みのメタデータがメタデータバッファ 1 1 b に格納されている場合に、該アクセラレータ 2 0 による演算処理の完了を判定する。

40

【 0 0 7 8 】

また、アクセラレータ 2 0 では、読込機能 2 1 b が、メタデータバッファ 1 1 b からメタデータを取得する。演算機能 2 4 が、取得されたメタデータに含まれる演算処理の対象の処理データの所在とアクセラレータ 2 0 の演算処理の内容を指定する情報に従って、演

50

算処理の対象の処理データに対する演算処理を行う。演算処理が完了した場合であって、メタデータに、転送先のアクセラレータの情報が含まれていない場合に、書込機能 2 1 c が、メタデータをメタデータバッファ 1 1 b に格納する。メタデータに、転送先のアクセラレータ 2 0 b の情報が含まれている場合に、書込機能 2 1 c が、該アクセラレータ 2 0 b に該制御用データを転送する。他のアクセラレータ 2 0 a からメタデータを受信した場合に、演算機能 2 4 が、メタデータに含まれる演算処理の対象の中間データの所在と、該演算処理の内容を指定する情報に従って、演算処理の対象の中間データに対する演算処理を行う。

【 0 0 7 9 】

これにより、アクセラレータ制御システム 1 は、複数のアクセラレータ 2 0 を用いた処理において、アクセラレータ制御装置 1 0 を介さずに実行されるため、転送元アクセラレータ / 転送先アクセラレータの排他処理や処理完了の割込み通知、同期制御等のオーバーヘッドを低減して、スケーラビリティを低下させず、リアルタイム処理を高速に行うことが可能となる。

【 0 0 8 0 】

また、アクセラレータ間メタデータ通信機構 2 1 d が、他のアクセラレータ 2 0 a からのメタデータの受信可否を確認する。これにより、アクセラレータ間でのデータ転送が確実に実行される。

【 0 0 8 1 】

また、書込機能 2 1 c が、所定の時間内に転送先のアクセラレータ 2 0 b からメタデータの受信可否が通知されない場合に、メタデータを再送する。これにより、アクセラレータ間でのデータ転送が確実に実行される。

【 0 0 8 2 】

また、データ書込機能 2 3 b は、転送先のアクセラレータ 2 0 b にメタデータが転送された場合に、演算処理対象の処理データに対する演算処理の結果の中間データを中間データバッファ 1 1 c に格納する。転送先のアクセラレータ 2 0 b において、演算機能 2 4 が、中間データバッファ 1 1 c に格納された演算処理の結果を読み込んで演算処理の対象とする。これにより、アクセラレータ間でのデータ転送が効率よく完了する。

【 0 0 8 3 】

また、Polling 機能 2 1 a が、メタデータバッファ 1 1 b にメタデータが有るか否かを監視する。これにより、アクセラレータ 2 0 が自律的に処理を行うことが可能となる。

【 0 0 8 4 】

また、Polling 機能 1 3 a が、は、メタデータバッファ 1 1 b に演算処理済みのメタデータが有るか否かを監視する。これにより、アクセラレータ制御装置 1 0 は、アクセラレータ 2 0 と非同期に処理を行うことが可能となる。

【 0 0 8 5 】

[プログラム]

上記実施形態に係るアクセラレータ制御装置 1 0 が実行する処理をコンピュータが実行可能な言語で記述したプログラムを作成することもできる。一実施形態として、アクセラレータ制御装置 1 0 は、パッケージソフトウェアやオンラインソフトウェアとして上記のアクセラレータ制御処理を実行するアクセラレータ制御プログラムを所望のコンピュータにインストールさせることによって実装できる。例えば、上記のアクセラレータ制御プログラムを情報処理装置に実行させることにより、情報処理装置をアクセラレータ制御装置 1 0 として機能させることができる。ここで言う情報処理装置には、デスクトップ型またはノート型のパーソナルコンピュータが含まれる。また、その他にも、情報処理装置にはスマートフォン、携帯電話機や P H S (Personal Handyphone System) などの移動体通信端末、さらには、P D A (Personal Digital Assistant) などのスレート端末などがその範疇に含まれる。また、アクセラレータ制御装置 1 0 の機能を、クラウドサーバに実装してもよい。

10

20

30

40

50

【0086】

図8は、アクセラレータ制御プログラムを実行するコンピュータの一例を示す図である。コンピュータ1000は、例えば、メモリ1010と、CPU1020と、ハードディスクドライブインタフェース1030と、ディスクドライブインタフェース1040と、シリアルポートインタフェース1050と、ビデオアダプタ1060と、ネットワークインタフェース1070とを有する。これらの各部は、バス1080によって接続される。

【0087】

メモリ1010は、ROM(Read Only Memory)1011およびRAM1012を含む。ROM1011は、例えば、BIOS(Basic Input Output System)等のブートプログラムを記憶する。ハードディスクドライブインタフェース1030は、ハードディスクドライブ1031に接続される。ディスクドライブインタフェース1040は、ディスクドライブ1041に接続される。ディスクドライブ1041には、例えば、磁気ディスクや光ディスク等の着脱可能な記憶媒体が挿入される。シリアルポートインタフェース1050には、例えば、マウス1051およびキーボード1052が接続される。ビデオアダプタ1060には、例えば、ディスプレイ1061が接続される。

10

【0088】

ここで、ハードディスクドライブ1031は、例えば、OS1091、アプリケーションプログラム1092、プログラムモジュール1093およびプログラムデータ1094を記憶する。上記実施形態で説明した各情報は、例えばハードディスクドライブ1031やメモリ1010に記憶される。

20

【0089】

また、アクセラレータ制御プログラムは、例えば、コンピュータ1000によって実行される指令が記述されたプログラムモジュール1093として、ハードディスクドライブ1031に記憶される。具体的には、上記実施形態で説明したアクセラレータ制御装置10が実行する各処理が記述されたプログラムモジュール1093が、ハードディスクドライブ1031に記憶される。

【0090】

また、アクセラレータ制御プログラムによる情報処理に用いられるデータは、プログラムデータ1094として、例えば、ハードディスクドライブ1031に記憶される。そして、CPU1020が、ハードディスクドライブ1031に記憶されたプログラムモジュール1093やプログラムデータ1094を必要に応じてRAM1012に読み出して、上述した各手順を実行する。

30

【0091】

なお、アクセラレータ制御プログラムに係るプログラムモジュール1093やプログラムデータ1094は、ハードディスクドライブ1031に記憶される場合に限られず、例えば、着脱可能な記憶媒体に記憶されて、ディスクドライブ1041等を介してCPU1020によって読み出されてもよい。あるいは、アクセラレータ制御プログラムに係るプログラムモジュール1093やプログラムデータ1094は、LANやWAN(Wide Area Network)等のネットワークを介して接続された他のコンピュータに記憶され、ネットワークインタフェース1070を介してCPU1020によって読み出されてもよい。

40

【0092】

以上、本発明者によってなされた発明を適用した実施形態について説明したが、本実施形態による本発明の開示の一部をなす記述および図面により本発明は限定されることはない。すなわち、本実施形態に基づいて当業者等によりなされる他の実施形態、実施例および運用技術等は全て本発明の範疇に含まれる。

【符号の説明】

【0093】

- 10 アクセラレータ制御装置
- 11 a データバッファ(記憶部)
- 11 b メタデータバッファ(記憶部)

50

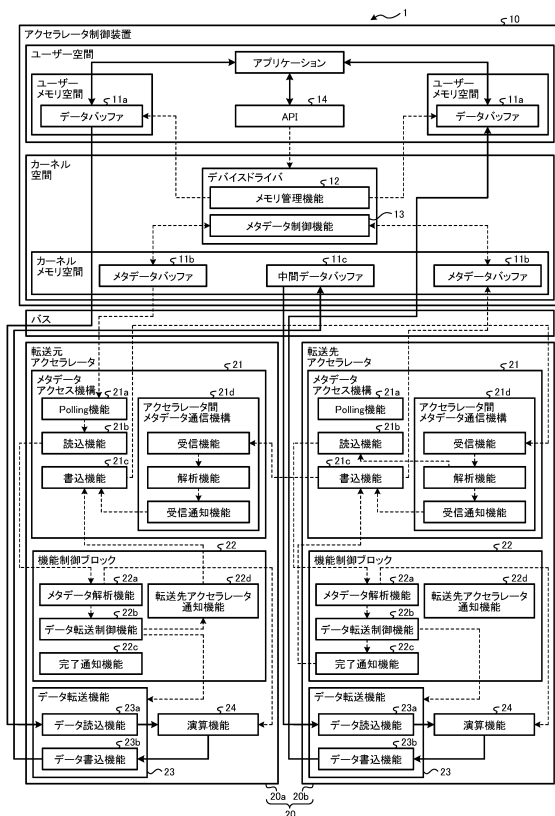
- 1 1 c 中間データバッファ (記憶部)
- 1 2 メモリ管理機能
- 1 3 メタデータ制御機能
- 1 3 a Polling機能 (判定部)
- 1 3 b 書込機能
- 1 4 API
- 2 0 アクセラレータ
- 2 0 a 転送元アクセラレータ
- 2 0 b 転送先アクセラレータ
- 2 1 メタデータアクセス機構
- 2 1 a Polling機能 (取得部)
- 2 1 b 読込機能
- 2 1 c 書込機能 (格納部、転送部)
- 2 1 d アクセラレータ間メタデータ通信機構 (確認部)
- 2 2 機能制御ブロック
- 2 2 a メタデータ解析機能
- 2 2 b データ転送制御機能
- 2 2 c 完了通知機能
- 2 2 d 転送先アクセラレータ通知機能
- 2 3 データ転送機能
- 2 3 a データ読込機能
- 2 3 b データ書込機能
- 2 4 演算機能 (演算部、第 2 の転送部)

10

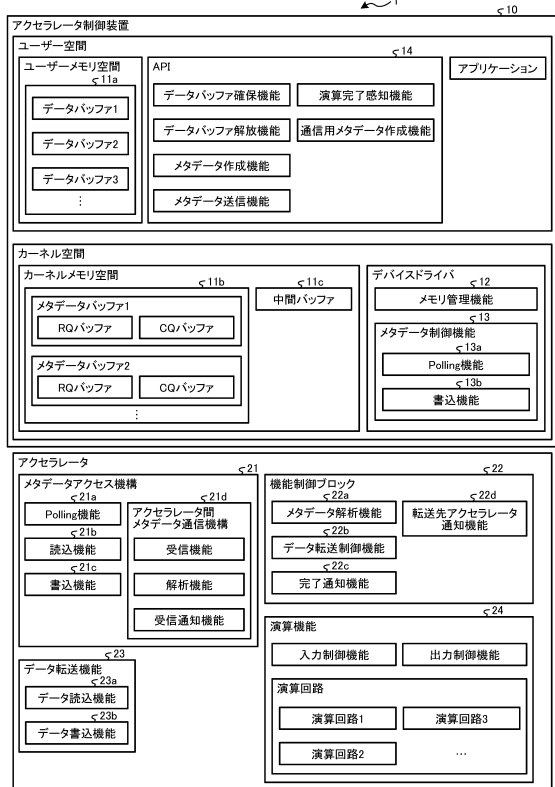
20

【図面】

【図 1】



【図 2】



30

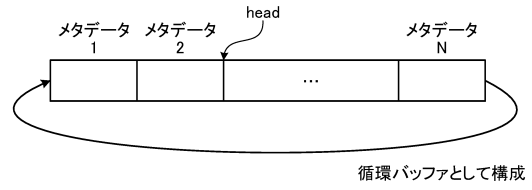
40

50

【 図 3 】

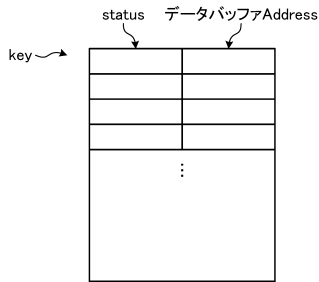
63	0
タスクID	オペレーション
読出データの長さ	読出データのアドレス
読出データのアドレス	書込データの長さ
書込データのアドレス	

【 図 4 】

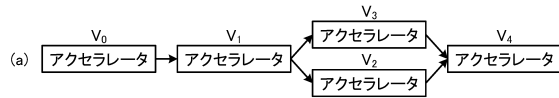


10

【 図 5 】



【 図 6 】



20

(b)

$$\begin{matrix}
 V_0 & V_1 & V_2 & V_3 & V_4 \\
 \begin{pmatrix}
 V_0 & 0 & 1 & 0 & 0 & 0 \\
 V_1 & 0 & 0 & 1 & 1 & 0 \\
 V_2 & 0 & 0 & 0 & 0 & 1 \\
 V_3 & 0 & 0 & 0 & 0 & 1 \\
 V_4 & 0 & 0 & 0 & 0 & 0
 \end{pmatrix}
 \end{matrix}$$

行列の大きさ 8bit	行列本体データ 256bit
----------------	-------------------

(c)

アクセラレータ駆動に必要な制御情報 N bit

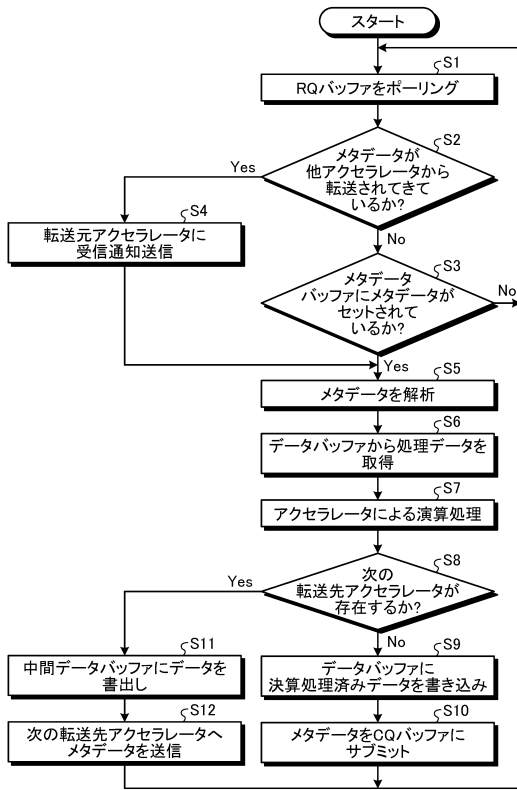
30

⋮

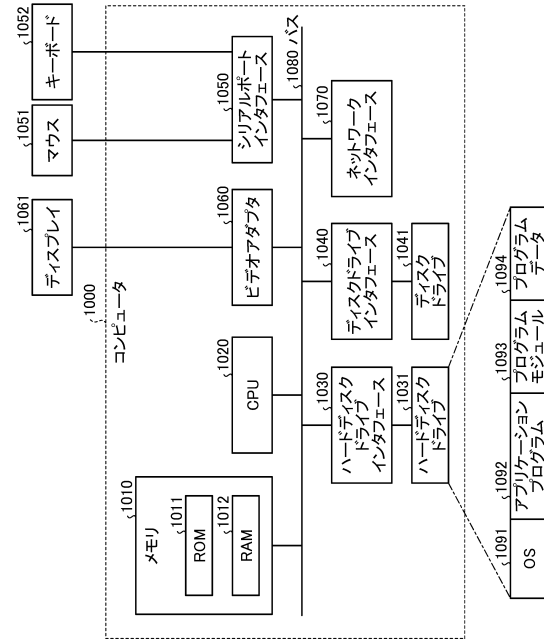
40

50

【 図 7 】



【 図 8 】



10

20

30

40

50

フロントページの続き

- (56)参考文献 特開 2 0 1 9 - 1 4 9 0 8 6 (J P , A)
特表 2 0 2 0 - 5 3 7 2 3 5 (J P , A)
- (58)調査した分野 (Int.Cl. , D B 名)
G 0 6 F 9 / 5 0