

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6392745号
(P6392745)

(45) 発行日 平成30年9月19日(2018.9.19)

(24) 登録日 平成30年8月31日(2018.8.31)

(51) Int.Cl.

F I

H O 4 L 12/931 (2013.01)
G O 6 F 15/173 (2006.01)H O 4 L 12/931
G O 6 F 15/173 6 8 5 S

請求項の数 24 (全 23 頁)

(21) 出願番号 特願2015-512707 (P2015-512707)
 (86) (22) 出願日 平成25年5月10日 (2013.5.10)
 (65) 公表番号 特表2015-522991 (P2015-522991A)
 (43) 公表日 平成27年8月6日 (2015.8.6)
 (86) 国際出願番号 PCT/US2013/040508
 (87) 国際公開番号 W02013/173181
 (87) 国際公開日 平成25年11月21日 (2013.11.21)
 審査請求日 平成28年4月14日 (2016.4.14)
 (31) 優先権主張番号 13/470,847
 (32) 優先日 平成24年5月14日 (2012.5.14)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 591016172
 アドバンスト・マイクロ・ディバイズ・
 インコーポレイテッド
 ADVANCED MICRO DEVI
 CES INCORPORATED
 アメリカ合衆国 95054 カリフォル
 ニア州、 サンタ クララ、 オーガスティ
 ン ドライブ 2485
 (74) 代理人 100108833
 弁理士 早川 裕司
 (74) 代理人 100111615
 弁理士 佐野 良太
 (74) 代理人 100162156
 弁理士 村雨 圭介

最終頁に続く

(54) 【発明の名称】 サーバノード相互接続デバイス及びサーバノード相互接続方法

(57) 【特許請求の範囲】

【請求項 1】

複数の電子デバイスと通信するスイッチファブリックと、
 前記スイッチファブリックと通信するランデブームメモリと、
 前記スイッチファブリックと1つ以上のネットワークインターフェースカード (NICs) との間に接続された入出力 (I/O) プロセッサと、を備え、
 前記複数の電子デバイスのうち第1の電子デバイスのメモリからデータが出力可能な状態にあるとの判定に応じて、且つ、前記ランデブームメモリに前記データ用の位置が割り当てられたことに応じて、前記データが前記第1の電子デバイスから前記ランデブームメモリに転送され、

前記 I/O プロセッサは、前記 NICs に接続されたネットワーク上のリモート電子デバイスへの前記ランデブームメモリからの前記 NICs を介したデータの転送に関連する動作を実行する、

集約デバイス。

【請求項 2】

前記第1の電子デバイスは、前記スイッチファブリックに接続したサーバノードである、請求項 1 に記載の集約デバイス。

【請求項 3】

前記第1の電子デバイスは、前記集約デバイスと通信するリモートコンピュータである、請求項 1 に記載の集約デバイス。

【請求項 4】

前記データは、前記ランデブームメモリから前記複数の電子デバイスのうち第 2 の電子デバイスに前記データが出力可能な状態にあるとの判定に応じて、且つ、前記第 2 の電子デバイスのメモリに前記データ用の位置が割り当てられたことに応じて、前記ランデブームメモリから前記第 2 の電子デバイスに転送される、請求項 1 に記載の集約デバイス。

【請求項 5】

前記データが第 2 の電子デバイスに提供可能になるまでに、前記ランデブームメモリにおいて前記データ用の位置を割り当てるランデブー管理モジュールをさらに含む、請求項 1 に記載の集約デバイス。

【請求項 6】

前記 I/O プロセッサは、出力可能な状態の前記データが受信される位置であって、前記ランデブームメモリに割り当てられた位置を識別する、請求項 1 に記載の集約デバイス。

【請求項 7】

前記複数の電子デバイスの中でデータを送信するスイッチングコントロールプレーンを管理するコントロールプレーンプロセッサをさらに含む、請求項 1 に記載の集約デバイス。

【請求項 8】

前記ランデブームメモリは TSV メモリを含む、請求項 1 に記載の集約デバイス。

【請求項 9】

前記ランデブームメモリはオンダイメモリを含む、請求項 1 に記載の集約デバイス。

【請求項 10】

複数のサーバノードと、

前記複数のサーバノードに接続した集約デバイスとを備え、

前記集約デバイスは、

前記複数のサーバノードと通信するスイッチファブリックと、

前記複数のサーバノードのうち第 1 のサーバノードのローカルメモリからデータが出力可能な状態にあるとの判定と、ランデブームメモリに前記データ用の位置が割り当てられたこととに
応じた、前記第 1 のサーバノードのローカルメモリを用いたデータの転送に関するランデブームメモリと、

前記スイッチファブリックと 1 つ以上のネットワークインターフェースカード (NICs) との間に接続された入出力 (I/O) プロセッサであって、前記 NICs に接続されたネットワーク上のリモート電子デバイスへの前記ランデブームメモリからの前記 NICs を介したデータの転送に関連する動作を実行する I/O プロセッサと、

を含む、サーバ集約システム。

【請求項 11】

前記データは、前記ランデブームメモリから宛先電子デバイスに前記データが出力可能な状態にあるとの判定と、前記データを受信するための位置が前記宛先電子デバイスに割り当てられたこととに
応じて、前記ランデブームメモリから前記宛先電子デバイスに転送される、請求項 10 に記載のサーバ集約システム。

【請求項 12】

前記宛先電子デバイスは、前記複数のサーバノードのうち第 2 のサーバノード、又は、前記集約デバイスと通信するリモートコンピュータである、請求項 11 に記載のサーバ集約システム。

【請求項 13】

前記データが宛先電子デバイスに提供可能になるまでに、前記ランデブームメモリにおいて前記データ用の位置を割り当てるランデブー管理モジュールをさらに含む、請求項 10 に記載のサーバ集約システム。

【請求項 14】

前記 I/O プロセッサは、出力可能な状態の前記データが受信される位置であって、前記ランデブームメモリに割り当てられた位置を識別する、請求項 10 に記載のサーバ集約シ

10

20

30

40

50

ステム。

【請求項 15】

前記第1のサーバノードの前記ローカルメモリと前記ランデブームメモリとの間で前記データの転送を確立するコントロールプレーンプロセッサをさらに含む、請求項10に記載のサーバ集約システム。

【請求項 16】

複数の電子デバイスと集約システムとの間で通信する方法であって、

第2の電子デバイスに送信するデータを、前記複数の電子デバイスのうち第1の電子デバイスのローカルメモリに提供することと、

前記第1の電子デバイスの前記ローカルメモリから前記データを受信するための前記集約システムのランデブームメモリを提供することと、

前記第1の電子デバイスの前記ローカルメモリに前記データが伝送用として存在するとの前記集約システムによる判定と、前記第1の電子デバイスの前記ローカルメモリから受信される前記データを記憶するための位置が前記ランデブームメモリに割り当てられているとの判定とに応じて、前記データを、前記第1の電子デバイスの前記ローカルメモリから前記ランデブームメモリに転送することと、

前記データを前記ランデブームメモリからリモートコンピュータである前記第2の電子デバイスのローカルメモリに転送することであって、前記データを、前記集約システムのI/Oプロセッサと1つ以上のネットワークインターフェースカード(NICs)を介して、前記NICs及び前記第2の電子デバイスに接続されたネットワークを経由して前記第2の電子デバイスに転送することを含む、ことと、

を含む方法。

【請求項 17】

前記第1の電子デバイスの前記ローカルメモリにおいて前記データが伝送用に使用可能であるという前記第1の電子デバイスからの通知を生成することと、

前記ランデブームメモリにて使用可能なバッファを識別することと、

前記第1の電子デバイスの前記ローカルメモリから前記データを受信するために前記使用可能なバッファを割り当てることと、

前記データを、前記第1の電子デバイスの前記ローカルメモリから前記ランデブームメモリの前記使用可能なバッファに出力することと、

をさらに含む請求項16に記載の方法。

【請求項 18】

前記第2の電子デバイスのローカルメモリが前記データを受信可能であるとの判定を行うことと、

前記判定に基づいて、前記データを、前記ランデブームメモリから前記第2の電子デバイスの前記ローカルメモリに転送することと、

をさらに含む請求項16に記載の方法。

【請求項 19】

前記ランデブームメモリにおいて前記データが伝送用に使用可能であるという前記ランデブームメモリからの通知を生成することと、

前記第2の電子デバイスの前記ローカルメモリにて使用可能なバッファを識別することと、

前記第2の電子デバイスの前記ローカルメモリにて前記使用可能なバッファを割り当てることと、

前記割り当てに基づいて、前記データを、前記ランデブームメモリから前記第2のデバイスの前記ローカルメモリの前記使用可能なバッファに転送することと、

をさらに含む請求項18に記載の方法。

【請求項 20】

前記第1の電子デバイスはサーバノードを含む、請求項16に記載の方法。

【請求項 21】

前記第2の電子デバイスはリモートコンピュータである、請求項16に記載の方法。

【請求項22】

複数の電子デバイスと集約システムとの間で通信する方法であって、

第2の電子デバイスに送信するデータを、前記複数の電子デバイスのうち第1の電子デバイスのローカルメモリに提供することと、

前記集約システムの仮想ネットワークインターフェースカード(vNIC)に対して、前記第1の電子デバイスの前記ローカルメモリにおいて前記データが伝送用として使用可能であると通知することと、

前記第1の電子デバイスの前記ローカルメモリに伝送用として前記データがあるとの判定と、前記集約システムのランデブームメモリに前記データ用の位置が割り当てられたこととに
10 応じて、前記データを、前記第1の電子デバイスの前記ローカルメモリから前記ランデブームメモリに転送することと、

前記第2の電子デバイスにおける少なくとも一つの使用可能なバッファの通知を生成することと、

前記データを受信するために前記第2の電子デバイスにおける少なくとも一つの使用可能なバッファを割り当てることと、

前記データを前記ランデブームメモリからリモートコンピュータである前記第2の電子デバイスにおける少なくとも一つの使用可能なバッファに転送することであって、前記データを、前記集約システムのI/Oプロセッサと1つ以上のネットワークインターフェースカード(NICs)を介して、前記NICs及び前記第2の電子デバイスに接続されたネット
20 ネットワークを経由して前記第2の電子デバイスに転送することを含む、ことと、

を含む方法。

【請求項23】

前記電子デバイスの少なくとも一つはサーバノードを含む、請求項22に記載の方法。

【請求項24】

コンピュータ可読プログラムコードを記憶するコンピュータ可読記憶媒体であって、

前記コンピュータ可読プログラムコードは、

集約システムを介して第2の電子デバイスにデータを送信するために、第1の電子デバイスのローカルメモリに前記データを提供するように構成されたコンピュータ可読プログラムコードと、
30

前記集約システムの仮想ネットワークインターフェースカード(vNIC)に対して、前記第1の電子デバイスの前記ローカルメモリにおいて前記データが伝送用として使用可能であると通知するように構成されたコンピュータ可読プログラムコードと、

前記第1の電子デバイスの前記ローカルメモリに前記データが伝送用として存在するとの判定と、前記ランデブームメモリに前記データ用の位置が割り当てられたこととに
40 応じて、前記データを、前記第1の電子デバイスの前記ローカルメモリから前記ランデブームメモリに転送するように構成されたコンピュータ可読プログラムコードと、

前記第2の電子デバイスにおける少なくとも一つの使用可能なバッファの通知を生成するように構成されたコンピュータ可読プログラムコードと、

前記データを受信するために前記第2の電子デバイスにおける前記少なくとも一つの使用可能なバッファを割り当てるように構成されたコンピュータ可読プログラムコードと、
40

前記vNICが、前記データを、前記ランデブームメモリからリモートコンピュータである前記第2の電子デバイスにおける前記少なくとも一つの使用可能なバッファに転送することであって、前記データを、前記集約システムのI/Oプロセッサと1つ以上のネットワークインターフェースカード(NICs)を介して、前記NICs及び前記第2の電子デバイスに接続されたネットワークを経由して前記第2の電子デバイスに転送することを含む、ことを行うように構成されたコンピュータ可読プログラムコードと、

を含むコンピュータ可読記憶媒体。

【発明の詳細な説明】

【技術分野】

10

20

30

40

50

【 0 0 0 1 】

本発明は、概して、ネットワークスイッチアーキテクチャに関し、より詳細には、スイッチファブリックを用いてサーバノードを相互接続するシステム及び方法に関する。

【背景技術】

【 0 0 0 2 】

データセンタは、概して、ビジネス及び組織をサポートするインターネット及びイントラネットサービスを提供する、一元管理型の設備である。典型的なデータセンタは、様々なタイプの電子機器（例えば、コンピュータ、サーバ（例えば、電子メールサーバ、プロキシサーバ及びDNSサーバ）、ネットワークスイッチ、ルーター、データ記憶デバイス、及び、その他の関連する構成要素）を格納している。所与のデータセンタは、スイッチ及びルーターを含むスイッチングアーキテクチャを介して相互に及び外部デバイスと通信する数百又は数千の相互接続されたサーバノードを有することが可能である。また、従来のデータセンタは、仮想化のための構成をとることにより、サーバノードに対して、ネットワークインタフェ이스カード（NIC）、ハードディスクドライブ又はその他のハードウェアの共有を可能にすることもできる。このようにする際には、通信データセンタアーキテクチャを介したノード間通信を容易にするために、複雑なスイッチファブリックが必要である。

10

【発明の概要】

【課題を解決するための手段】

【 0 0 0 3 】

一態様によれば、スイッチファブリックとランデブームメモリとを含む集約デバイスが提供される。スイッチファブリックは、複数の電子デバイスと通信する。ランデブームメモリは、スイッチファブリックと通信する。データは、複数の電子デバイスのうち第1の電子デバイスのメモリから前記データが出力可能な状態にあるという判定と、前記データ用のランデブームメモリに割り当てられた位置とに応じて、前記第1の電子デバイスから前記ランデブームメモリに転送される。

20

【 0 0 0 4 】

別の態様によれば、複数のサーバノードと、前記複数のサーバノードに接続した集約デバイスとを含むサーバ集約システムが提供される。集約デバイスは、前記複数のサーバノードのうち第1のサーバノードのローカルメモリを用いたデータ転送に関与するランデブームメモリを含む。このデータ転送は、前記第1のサーバノードのメモリからデータが出力可能であるという判定と、前記ランデブームメモリに割り当てられた前記データ用の位置とに応じて行われる。

30

【 0 0 0 5 】

別の態様によれば、複数の電子デバイスと集約システムとの間で通信する方法が提供される。この方法は、複数の電子デバイスのうち第2の電子デバイスに送信するデータを、複数の電子デバイスのうち第1の電子デバイスのローカルメモリに提供することと、前記第1の電子デバイスのローカルメモリからデータを受信するために、ランデブームメモリを集約システムに設けることと、伝送用のデータが前記第1の電子デバイスのローカルメモリに存在するとの集約システムによる判定と、前記第1の電子デバイスのローカルメモリから受信したデータを記憶するための位置がランデブームメモリに割り当てられたとの判定とに応じて、第1の電子デバイスのローカルメモリからランデブームメモリにデータを転送することと、を含む。

40

【 0 0 0 6 】

別の態様によれば、複数の電子デバイスと集約システムとの間で通信する方法が提供される。この方法は、複数の電子デバイスのうち第2の電子デバイスに送信するデータを、複数の電子デバイスのうち第1の電子デバイスのローカルメモリに提供することと、第1の電子デバイスのローカルメモリにおいてデータが伝送用に使用可能であることを、集約システムの仮想ネットワークインタフェースカード（vNIC）に通知することと、第2の電子デバイスにおける少なくとも一つの使用可能なバッファの通知を生成することと

50

、データを受信するために前記少なくとも一つの使用可能なバッファを割り当てることと、データを、第2の電子デバイスにおける前記少なくとも一つの使用可能なバッファにVNICによって出力することと、伝送用のデータが前記第1の電子デバイスのローカルメモリに存在するとの判定と、前記第2の電子デバイスにて前記少なくとも一つの使用可能なバッファを割り当てたとの判定とに応じて、第1の電子デバイスのローカルメモリからランデブーメモリにデータを転送することと、を含む。

【0007】

別の態様によれば、コンピュータ可読プログラムコードを有するコンピュータ可読記憶媒体を含むコンピュータプログラム製品であって、前記記憶媒体を用いてプログラムコードが具体化されるコンピュータプログラム製品が提供される。コンピュータ可読プログラムコードは、集約システムを介して第2の電子デバイスに送信するためのデータを、第1の電子デバイスのローカルメモリに提供するように構成されたコンピュータ可読プログラムコードと、第1の電子デバイスのローカルメモリにおいてデータが伝送用に使用可能であることを、集約システムの仮想ネットワークインターフェースカード(VNIC)に通知するように構成されたコンピュータ可読プログラムコードと、第2の電子デバイスにおける少なくとも一つの使用可能なバッファの通知を生成するように構成されたコンピュータ可読プログラムコードと、データを受信するために前記少なくとも一つの使用可能なバッファを割り当てるように構成されたコンピュータ可読プログラムコードと、第2の電子デバイスにおける少なくとも一つの使用可能なバッファに対して、VNICによってデータを出力するように構成されたコンピュータ可読プログラムコードと、伝送用のデータが前記第1の電子デバイスのローカルメモリに存在するとの判定と、前記第2の電子デバイスにて前記少なくとも一つの使用可能なバッファを割り当てたとの判定とに応じて、第1の電子デバイスのローカルメモリからランデブーメモリにデータを転送するように構成されたコンピュータ可読プログラムコードと、を含む。

【0008】

本発明の上記の利点及び更なる利点は、添付の図面に関連した後述の記載を参照することによって、より良く理解し得る。これらの図面において同様に表された数字は、同様の構成要素及び特徴を示している。図面は、必ずしも縮尺が揃っていないが、本発明の原理を例示することに重点が置かれている。

【図面の簡単な説明】

【0009】

【図1】従来のデータセンタネットワーク階層のブロック図である。

【図2】従来のインターコネクトファブリックスイッチに結合した複数のサーバノードを含むデータセンタのブロック図である。

【図3】一実施形態における、ランデブーメモリを有するサーバ集約システムと通信している複数のサーバノードを含む計算インフラストラクチャーのブロック図である。

【図4】一実施形態における、図3のランデブー管理装置のブロック図である。

【図5】一実施形態における、図3のネットワークプロセッサのブロック図である。

【図6】一実施形態における、サーバノードと、ランデブーメモリを有するサーバ集約システムとの間で電子通信を行う方法のフロー図である。

【図7】一実施形態における、ソースサーバノードと、ランデブーメモリを有するサーバ集約システムとの間で電子通信する方法のフロー図である。

【図8】一実施形態における、サーバ集約システムと、宛先サーバノードとの間で通信を確立する方法のフロー図である。

【図9】他の実施形態における、ランデブースイッチング用に構成されたサーバ集約システムに結合した第1及び第2のサーバノード間で通信を確立する方法のフロー図である。

【発明を実施するための形態】

【0010】

以下の記載では具体的な詳細が記載されるが、本発明にかかるシステム及び方法が少なくともいくつかの詳細無しで実行できることは、当業者にとって認識し得るはずである。

場合によっては、本発明が曖昧にならないように、既知の特徴又は工程は、詳細に記載されていない。

【0011】

図1は、従来のデータセンタネットワーク階層10のブロック図である。図1では、複数のサーバノード12-1~12-N(Nは、1より大きい整数)は、例えば、ラック18-1のイーサネット(登録商標)ローカルエリアネットワーク(LAN)又は関連するデータネットワークなどのLANを介して、ラックスイッチ14と通信することができる。ラック18-1は、データセンタネットワーク階層10における他の一つ以上のラック18-2, 18-N(Nは、1より大きい整数)とともに、クラスタの一部として構成されてもよい。各クラスタは、集約スイッチ22-1, 22-2, 22-N(Nは、1より大きい整数)を含むことができ、これらは、イーサネット(登録商標)又は他のネットワーク接続を介して、コアルータ24に接続されている。ユーザコンピュータ32(例えば、ラップトップ、スマートフォン又はその他の電子デバイス)は、ネットワーク26(例えば、IPネットワーク、インターネットなど)を介して、サーバノード12-1~12-N(概して、12)と通信することができる。データセンタスイッチに関連する規模、経費及び電力の削減を求める一方で、データセンタの拡張性の増大を求めるという、止むことのない要求が存在する。

10

【0012】

図2は、従来のインターコネクトファブリックスイッチ120に結合した複数のサーバノード112-1~112-N(Nは、1より大きい整数)を含むデータセンタ20のブロック図である。インターコネクトファブリックスイッチ120は、単一のシャシのもとでの集約スイッチ及びラックスイッチの一元管理機能によって、スイッチの数を削減することができる。従来のサーバラックにおいて使用される従来のイーサネット(登録商標)スイッチングの要求の多くを置き換えることができる。

20

【0013】

サーバノード112-1~112-N(概して、112)は、プロセッサクラスタとして、又は、他の周知の装置として構成することができる。サーバノード112は、単一のソケットサーバ、又は、共用インフラストラクチャーを共有する低消費電力プロセッサ102を含むことができる。サーバノードプロセッサ102は、一つ以上のマイクロプロセッサ、中央演算処理装置(CPU)、画像処理装置(GPU)、デジタルシグナルプロセッサ(DSP)、特定用途向け集積回路(ASIC)、メモリコントローラ、マルチコアプロセッサ、及び、その他のタイプのデータ処理デバイスだけでなく、これらのデバイスや他のデバイスの一部及び/又は組合せを含むことができる。

30

【0014】

また、サーバノード112は、ローカルメモリ104と、I/Oロジック106とを含む。ローカルメモリ104は、不揮発性又は揮発性メモリ(例えば、一つ以上のチャネルのダイナミックランダムアクセスメモリ(DRAM)又はスタティックRAM(SRAM))を含むことができる。I/Oロジック106は、コンピュータのI/O機能(例えば、サーバノード112とインターコネクトファブリックスイッチ120との間のデータ転送)を管理するために、サウスブリッジ等を含むI/Oコントローラを有するように構成されてもよい。I/Oロジック106は、インターコネクトファブリックスイッチ120と電子通信を確立するために、イーサネット(登録商標)、PCIe又はその他のネットワークコネクタ114を含んでいてもよい。

40

【0015】

インターコネクトファブリックスイッチ120は、複数の入出力ポートと、可変長又は固定長のフレーム、データパケット、セル等を入出力ポート間でルーティングし得るクロスバー124とを含み、サーバノード112、例えばNIC142又はハードドライブ144等の共有デバイス、及び/又は、例えばユーザコンピュータ152等の外部電子デバイス間での通信を容易にすることができるものである。ポートは、仮想ポート、物理ポート又はこれらの組合せを含んでいてもよい。ポートは、単方向又は双方向通信用の構成で

50

あってもよい。

【 0 0 1 6 】

クロスバー 1 2 4 は、行と列の配列に組織化され得る。データ伝送の間には、所与の行のいくつかの入力ポートは、ある列における出力ポートを奪い合う場合がある。入力ポートは、入力バッファ 1 2 6 と通信することができる。入力バッファは、クロスバー 1 2 4 が、使用可能な出力ポートにデータを送信することができるまで、サーバノード 1 1 2 から受信したデータを一時的に記憶するものである。出力ポートは、出力バッファ 1 2 8 を含む。出力バッファは、ネットワークバス 1 4 0 (例えば、イーサネット(登録商標)バス、P C I e バス等)を介してパケットを宛先に送信するために、所望の出力ポートが使用可能になるまで、一つ以上の入力ポートから受信したデータを一時的に記憶する。

10

【 0 0 1 7 】

しかしながら、インターコネクトファブリックスイッチ 1 2 0 は、輻輳に関連する問題(例えば、ヘッドライン(H O L)ブロッキング)が生じやすい。この問題は、スイッチ 1 2 0 の複数の入力ポートが同一の出力ポートを奪い合う場合に発生するものである。また、スケールングに関連する課題も発生するが、これは、スイッチ 1 2 0 での I / O ポートの追加によって、共有のリソースを奪い合うというリスクが増加するからである。関連する課題は、サーバノード 1 1 2 と、インターコネクトファブリックスイッチ 1 2 0 との間でバンド幅を十分に使用しないことである。例えば、H O L ブロッキングは、特に、入力ポートで待機しているパケット、セル等が、ファーストイン・ファーストアウト(F I F O)のキューで記憶されている場合に、クロスバー 1 2 4 でかなりの量のバンド幅を消費する可能性がある。また、入力及び出力バッファ 1 2 6 , 1 2 8 でのキューはすぐに満杯になり、その結果、望ましくない待ち時間、ジッター又はパケット喪失のみならず、オーバーラン及び/又はアンダーラン条件が生じ、これらは、バンド幅の不十分な使用及び性能の課題を生じる。

20

【 0 0 1 8 】

フロー制御技術は、容易に使用可能であり、データセンタ 2 0 でのネットワーク輻輳を緩和することができる。例えば、複雑なスケジューリング技術を適用して、クロスバー 1 2 4 を介したデータトラヒックを管理することができる。しかしながら、そのような技術は、概して、実現するには高価である。

【 0 0 1 9 】

他の周知である技術は、イーサネット(登録商標)をベースにした再送信を実現することである。しかしながら、イーサネット(登録商標)をベースにした再送信は、貴重なバンド幅を消費する。

30

【 0 0 2 0 】

他のフロー制御メカニズムを上流のデバイスに適用して、パケットの通過を停止させるように当該デバイスに要求することができる。そのような手法は、バッファオーバーフローを低減させるのに効果的であるものの、ネットワーク輻輳を完全に緩和しない。

【 0 0 2 1 】

他の手法は、一時記憶用の外部メモリデバイスと、パケットのキューイング(q u e u i n g)とに依るものである。しかしながら、この手法は、インターコネクトファブリックスイッチ 1 2 0 に対して、さらなるピン、及び、すぐにでも使用可能なメモリデバイスへの「配線」又はコネクタを必要とし、その結果として、拡張性の制限と、バンド幅に関連した問題とが、さらなる電力消費を求める関連要求とともに生じる。

40

【 0 0 2 2 】

本発明の概念は、従来のデータセンタ相互接続スイッチに付随する輻輳に関連した課題を、ランデブーデバイス(r e n d e z v o u s d e v i c e)を含むサーバ集約システムを導入することによって低減又は取り除くものである。ランデブーデバイスは、制御された方法で、高バンド幅スイッチファブリックを介して、複数のサーバノード又は関連する電子デバイスのローカルメモリと通信する。ランデブーデバイスは、出力可能なデータを有する伝送デバイスと、当該データの受信用に使用可能なメモリバッファを有する受

50

信デバイスとの間で通信し、伝送デバイス及び受信デバイス間で効率的なデータ伝送経路を提供する。例えば、データが、伝送デバイスのローカルメモリから明示的に転送され得るのは、データが出力可能であるとサーバ集約システムが判定した後であって、データ記憶用のランデブーデバイスのメモリ記憶位置に特定のメモリ記憶位置が割り当てられた後である。これらの条件が真であると判定された場合、すなわち、送信するデータを伝送デバイスが有して、使用可能なバッファ（空きバッファ）を受信デバイスが有している場合に、伝送デバイスから受信デバイスへの効率的で待ち時間の短いデータ移動が生じる。

【 0 0 2 3 】

図 2 に記載された従来のネットワークスイッチ環境では、データは、クロスバースイッチキュー 1 2 6 , 1 2 8 が使用可能な空間を有しているか否かに関係なく、キュー 1 2 6 , 1 2 8 に出力される。キュー 1 2 6 , 1 2 8 が使用可能な空間を有していない場合には、データは通常破棄され、複雑でバンド幅を消費する再送信過程が実行される。

【 0 0 2 4 】

一実施形態では、サーバ集約システムは、ソースサーバノード等からのデータをいつ受信するのかを判定して、「ランデブー (rendezvous)」する位置を確立することができる。この位置は、宛先サーバノード等のメモリ記憶位置と通信して、宛先サーバノードにデータをいつ送信するのかを提供することができる。事前に割り当てられた位置は、サーバ集約システムのあらゆる入力ポートからデータを着信するために、ランデブーメモリ及び / 又は宛先サーバノードメモリに提供される。これを行う場合、ソースサーバノードメモリからランデブーメモリに、又は、ランデブーメモリから宛先サーバノードメモリにデータを移動させるために、解放されているデータバッファが割り当てられる。ソースサーバノードから送信されることが意図されるデータ用のランデブーメモリにおいて使用可能な位置が存在しない場合には、ランデブーメモリに記憶位置の空きができるまで、データはランデブーメモリに転送されない。また、宛先サーバノードが、データを受信する宛先サーバノードのローカルメモリに使用可能な記憶位置があることをサーバ集約システムに通知するまで、受信データをランデブーメモリに記憶することができる。

【 0 0 2 5 】

他の実施形態では、リモートコンピュータは、外部ネットワーク上のサーバ集約システムへのデータを、外部ネットワークとサーバ集約システムとの間に結合された NIC 又は関連インターフェースに提供する。ここでは、NIC は、（例えば、NIC の受信リングでの記述子により識別される）ランデブーメモリの受信バッファにデータを転送する。使用可能なバッファが不十分である場合には、データは、脱落又はフィルタリングされる。そうでなければ、受信バッファは使用可能であり、データはランデブーメモリに提供され、ランデブーメモリと通信している処理複合体に通知される。次いで、メモリバッファは検査され、メモリバッファをどこにルーティングするべきか、そしてスイッチングファブリックへの処理複合体接続における記述子上のどこに置くべきかが判定される。この時点で、データは、あたかもそれがイングレス v NIC にあるのと同様な方法で移動される。ここでは、ランデブーメモリ管理装置は、宛先サーバノードによって、例えば v NIC を介して、宛先サーバノードメモリに空き記憶位置があるかどうか通知される。ランデブーメモリ管理装置は、受信バッファが宛先サーバノードメモリに割り当てられていない限り、宛先サーバノードメモリへの伝送用のデータをフェッチすることはしない。

【 0 0 2 6 】

従って、本発明の概念の特徴は、データ転送におけるランデブーメモリの入力及び / 又は出力領域、すなわち、ソースサーバノードと通信するランデブーメモリの一端と、宛先サーバノードと通信するランデブーメモリのもう一端とで、フロー制御を生じさせることができるということである。また、スイッチファブリックにより実行されるバッファリング (buffering) は、ほとんど又は全く存在しないことから、スイッチファブリックのイングレスポートからエグレスポートへの処理待ち時間が低い。例えば、データトラフィックをスケジューリングすることにより、投機的データ移動が原因となる輻輳

10

20

30

40

50

又はその他のトラヒックに起因するブロッキングのリスクが低減するから、待ち時間を改善することができる。

【 0 0 2 7 】

本発明の概念の他の特徴は、バンド幅を、サーバ集約システム内及びサーバノード間で効率的に分配可能なことである。これは、出力することになるデータをソースサーバノードのローカルメモリが有しているかどうか、ランデブーメモリがそのデータを受信できるかどうか、及び/又は、データのメモリ間交換を行っているサーバノードでのローカルメモリが、データを受信するのに十分な空間を有しているかどうかを、サーバ集約システムが実際のデータ転送に先立って判定することができるからである。

【 0 0 2 8 】

図 3 は、本実施形態における、ランデブーメモリ (R - V メモリ) 3 0 8 を有するサーバ集約システム 3 0 0 と通信している複数のサーバノード 3 1 2 - 1 ~ 3 1 2 - N を含む計算インフラストラクチャー 3 0 のブロック図である。計算インフラストラクチャー 3 0 は、大規模なデータセンタ、クラウド等を含むことができる。計算インフラストラクチャーは、図 2 を参照して記載されるデータセンタ 2 0 と同様の仮想化した構成をとることができる。従って、仮想化関連の構成要素に関する詳細は、簡潔にするために図 3 では省略される。

【 0 0 2 9 】

サーバ集約システム 3 0 0 は、スイッチファブリック 3 0 2 と、ランデブー管理モジュール 3 0 4 と、I / O プロセッサ 3 0 6 と、ランデブーメモリ 3 0 8 と、コントロールプレーンプロセッサ 3 1 0 とを含む。サーバ集約システム 3 0 0 は、仮想化した構成とすることができる。

【 0 0 3 0 】

サーバノード 3 1 2 - 1 ~ 3 1 2 - N (概して、3 1 2)、又は、マイクロサーバ及び/若しくは少なくとも一つのプロセッサを有するその他の電子デバイスは、コネクタ 3 1 4、好ましくは P C I e バス、又は、その他のネットワークコネクタを介して、サーバ集約システム 3 0 0 と通信することができる。各コネクタ 3 1 4 は、一つ以上のサーバノード 3 1 2 とサーバ集約システム 3 0 0 のと間のデータ経路を提供することができる。サーバ集約システム 3 0 0 及びサーバノード 3 1 2 は、同一のマルチプロセッシングユニット (例えばチップ、計算デバイス又はラック) に一緒に設置されてもよい。他の実施形態では、サーバノード 3 1 2 は、一つ以上のユニット上に形成され、サーバ集約システム 3 0 0 は、独立ユニット上 (例えば、チップ上) に形成される。

【 0 0 3 1 】

スイッチファブリック 3 0 2 は、複数の入力ポート及び出力ポートを含んでおり、データ又はその他の電子情報を、その入力ポートと通信するサーバノード 3 1 2 と、そのエグレスポートと通信するランデブーメモリ 3 0 8 との間で移動させる構成とすることができる。スイッチファブリック 3 0 2 は、イングレスポートとエグレスポートとの間でデータを移動させるスイッチング構成 (例えば、クロスバー) を含むことができ、このことは当業者には周知である。従来のスイッチファブリックとは異なり、スイッチファブリック 3 0 2 では、従来のキューイング技術が必要ないことから、バッファリングがほとんど又は全く必要ない。スイッチファブリック 3 0 2 及びランデブーメモリ 3 0 8 は、相互接続することができるので、それらの間での I / O バンド幅通信チャネルを有しており、従って、ランデブーメモリ 3 0 8 にデータ用として割り当てられた十分な空間がある限り、いくつかのサーバノード 3 1 2 からのデータも受信することができる。

【 0 0 3 2 】

ランデブー管理モジュール 3 0 4 は、コントロールプレーンプロセッサ 3 1 0 及び I / O プロセッサ 3 0 6 とともに、サーバノード 3 1 2 におけるメモリ 1 0 4 - 1 ~ 1 0 4 - N (概して、1 0 4) の一つ以上とランデブーメモリ 3 0 8 との間のデータパケットの明示的な転送を監視する。ランデブー管理モジュール 3 0 4 は、ランデブーメモリ 3 0 8 にバッファ空間を割り当てて、宛先サーバノード又は外部計算デバイスが、ソースサーバノ

10

20

30

40

50

ードから送信されたデータをランデブームメモリ308から読み出すことができるまで、それを「駐車 (parking) 」させておくことができる。ランデブー管理モジュール304は、記述子等を使用して、ランデブームメモリ308から宛先メモリにデータが伝送される方法を制御することができる。ランデブー管理モジュール304は、ランデブームメモリ308の受信バッファの空きを監視し、受信バッファがポスト (post) されるまで待機する。換言すれば、ランデブー管理モジュール304は、転送操作における一方の当事者 (すなわち、ソースサーバノード又はランデブームメモリ) が伝送用データを有すること、転送操作における他方の当事者 (すなわち、宛先サーバ又はランデブームメモリ) がデータを受信するのに十分な空間を有していることを保証することができる。

【0033】

ランデブームメモリ308は、複数のネットワークスイッチのリンクと通信する場合には、バンド幅が十分に高い (例えば、10Gb以上) 貫通ピアシリコン (TSV)、SRAM又はオンダイ (on - die) メモリを含むことができる。ランデブームメモリ308は、例えば、メモリのブロックに組織化された解放された複数のデータバッファを含むように構成され得るものであって、これらのデータバッファは、データをローカルメモリ104からランデブームメモリ308に移動させる場合に、一つ以上のvNIC334に割り当てられ得る。

【0034】

I/Oプロセッサ306は、ランデブームメモリ308と、一つ以上の電子デバイス (例えば、サーバノード312及び/又はリモート計算デバイス352) との間で転送されたデータを処理する。I/Oプロセッサ306は、マルチプレクサ、並びに、ランデブームメモリ308へのデータ転送及びランデブームメモリ308からのデータ転送を実行するその他のロジックを含むことができるが、転送の実行は、ランデブー管理モジュール304、コントロールプレーンプロセッサ310又はそれらの組合せを用いて形成された通信に従ってなされる。このように、I/Oプロセッサ306は、サーバ集約システム300へ及びサーバ集約システム300から (例えば、二つ以上のサーバノード312の間、又はサーバノード312とリモート計算デバイス352の間) データを移動させる集結地 (staging area) としての役割を果たすことができる。

【0035】

I/Oプロセッサ306は、使用可能なバッファを示す受信記述子をポストすることによって、ソースサーバノード312からの伝送可能状態にあると判定された特定のデータが一時的に (例えば、宛先サーバノード312がデータを読み出すまで) 記憶されることになるランデブームメモリ308におけるバッファを識別することができる。I/Oプロセッサ306は、ランデブームメモリ308からバッファ記述子を処理するvNIC334へのデータ移動のために、例えばランデブームメモリ308の一つ以上のチャネル用のバッファ記述子を保持することができる。I/Oプロセッサ306は、受信バッファのリソースを受信記述子に追加する。

【0036】

コントロールプレーンプロセッサ310は、ネットワークルーティングプロトコルを処理することによって、サーバノード312とサーバ集約システム300との間のデータ送信のスイッチングコントロールプレーン (図示省略) を管理し、サーバ集約システム300等で受信されるデータパケット、フレーム、セル等の転送に関与する。コントロールプレーンプロセッサ310の他の機能は、データレディ (data ready) 通知、受信バッファの空き通知、バッファ解放等の生成を含む。コントロールプレーンプロセッサ310は、サーバ集約システム300内のブロック、ポート及びノードの間の通信を提供することができ、ポート間のデータ移動のためにスイッチファブリック302と通信する。コントロールプレーンプロセッサ310は、ランデブームメモリ308及び/又は一つ以上のサーバノードメモリ104への書き込み及び/又は読み出しを行うことができる。コントロールプレーンは、使用可能な受信バッファの数を問い合わせるためのデータ可用メッセージ (例えば、送信可能メッセージ) 及びバッファ解放メッセージを送信するよう

10

20

30

40

50

に構成され得る。

【 0 0 3 7 】

コントロールプレーンプロセッサ 3 1 0 は、ゼロ又は一つ以上の v N I C 3 3 4 を、好ましくはダイレクトメモリアクセス (D M A) エージェント (a g e n t) (図示省略)、又は、スイッチファブリック 3 0 2 における関連するエージェントと組み合わせて実現及び管理することができる。様々なタイプの通信が、コントロールプレーンプロセッサ 3 1 0 の管理下で、コントロールプレーン上で行われてもよい。例えば、このような通信は、サーバノードインターフェースドライバから書き込まれ、バッファ記述子メモリ記憶位置の状態の変化を指示する伝送 / 受信記述子のドアベル (d o o r b e l l) 通知であってもよい。他の例では、コントロールプレーンは、伝送を終えたバッファを、さらなる伝送のために又は受信バッファとして再利用させるなどのバッファ解放メッセージを管理することができる。コントロールプレーンの他の機能は、データレディ通知の提供、バッファ空き通知の受信、及び、バッファ解放等を含むことができる。

10

【 0 0 3 8 】

図 4 は、一実施形態における、図 3 のランデブー管理モジュール 3 0 4 のブロック図である。ランデブー管理モジュール 3 0 4 は、メモリ割り当てモジュール 4 0 2、通知モジュール 4 0 4、データ配信モジュール 4 0 6、割り込み生成モジュール 4 0 8 及び / 又はタイマ 4 1 0 を含むことができる。ランデブー管理モジュール 3 0 4 は、同一デバイス (例えば、チップ、ラック等) の一部として示されている。代替として、ランデブー管理モジュール 3 0 4 のいくつかの構成要素は、図 3 に示された計算インフラストラクチャー 3 0 内の別の位置に物理的に配置することができる。

20

【 0 0 3 9 】

メモリ割り当てモジュール 4 0 2 は、ソースサーバノード 3 1 2 - 1 からデータをフェッチするときに v N I C 3 3 4 が使用するバッファ空間を、ランデブーメモリ 3 0 8 内に割り当てる。より詳細には、メモリ割り当てモジュール 4 0 2 は、ランデブーメモリ 3 0 8 のブロックを示すインデックスを提供することができる。このインデックスは、ランデブーメモリの空きデータバッファブロックを示すバッファ記述子を含む。

【 0 0 4 0 】

通知モジュール 4 0 4 は、ランデブーメモリ 3 0 8 にデータが入力された場合に、バッファ割り当てに関連した通知を生成する。例えば、通知モジュール 4 0 4 は、バッファを要求する v N I C 3 3 4 に応答して、受信バッファ割り当てを v N I C 3 3 4 に通知することができる。通知モジュール 4 0 4 は、サーバノード 3 1 2 からの伝送用にデータが使用可能であることを示す通知を、例えば v N I C 3 3 4 から受信することができる。通知モジュール 4 0 4 は、ランデブーメモリ 3 0 8 から宛先デバイス用の v N I C 3 3 4 への記述子及び / 又はデータをフェッチした後に新規データが使用可能であることを、 v N I C 3 3 4 に通知することができる。

30

【 0 0 4 1 】

データ配信モジュール 4 0 6 は、割り当て用のランデブーメモリ 3 0 8 から v N I C 3 3 4 への記述子、データ等をフェッチすることができる。この v N I C 3 3 4 は、次いで、データを宛先ローカルメモリに提供する。

40

【 0 0 4 2 】

割り込み生成モジュール 4 0 8 は、新規データがランデブーメモリ 3 0 8 にある場合に、割り込み信号を I / O プロセッサ 3 0 6 に出力する。これにより、I / O プロセッサ 3 0 6 は、データをランデブーメモリからプル (p u l l) し、当該データを提供する。また、割り込み生成モジュール 4 0 8 は、伝送記述子の解放後 (例えば、ランデブーメモリ 3 0 8 から宛先ローカルメモリ 1 0 4 - N にデータを移動させた後) に、割り込み信号をネットワークプロセッサに出力することができる。

【 0 0 4 3 】

タイマ 4 1 0 は、受信バッファが、ランデブーメモリ 3 0 8 にて v N I C 3 3 4 への割り当て用に使用可能でない (空きがない) 場合に起動し得る。タイマ 4 1 0 は、データを

50

ランデブーメモリ 308 に転送することが不可能であることをサーバ集約システム 300 の構成要素に示し、これにより、フロー制御のレベルを提供する。

【0044】

図5は、一実施形態における、図3のI/Oプロセッサ306のブロック図である。I/Oプロセッサ306は、プロセッサ502、記述子処理モジュール504及び/又はスケジューラ508を含む。I/Oプロセッサは、メモリ(図示省略)を含んでもよい。I/Oプロセッサ306は、同一のデバイス(例えばチップ、ラック等)の一部として示されているが、I/Oプロセッサ306の構成要素のいくつか又はすべては、図3に示された計算インフラストラクチャー30の他の構成要素にあってもよい。

【0045】

プロセッサ502は、データ出力に関連する構成要素(例えば、NIC342)とのインターフェースをとるのに使用されるプログラムのプログラムコードを実行することが可能である。プロセッサ502は、ルーティングの決定を行い、受信バッファを、受信記述子リングから宛先伝送記述子リングに移動させることができる。

【0046】

記述子処理モジュール504は、ソースサーバノード312-1からデータを受信するランデブーメモリ308の空き記憶位置を示す受信記述子を、生成及び/又はポストすることができる。受信記述子は、ランデブーメモリ308において、バッファが常駐する位置、バッファのサイズ、複数のセグメント等を示することができる。I/Oプロセッサ306は、現在解放されている(データを受信するのに空いている)データバッファの記述子を、ランデブーメモリ308からI/Oプロセッサ306へ受け渡すことによって、ランデブーメモリ308内の使用可能なデータバッファに関する情報を受信することができる。また、記述子処理モジュール504は、ランデブーメモリ308から宛先サーバノード312-Nのローカルメモリ104-Nにデータを転送する命令、分散収集リスト等を含む伝送記述子を生成することができる。

【0047】

スケジューラ508は、例えば、NIC342若しくはPCIeコネクタ314を介してサーバ集約システム300が新規データを受信するとの通知、又は、宛先デバイスへの伝送用としてランデブーメモリ308内で新規データが使用可能であるとの通知を受信した場合に、起動する。スケジューラ508は、ランデブーメモリと、一つ以上のサーバノード312及び/又はリモート計算デバイス352との間のデータの転送を、例えばラウンドロビンの順、ファーストインファーストアウトの順又は本発明の技術分野で周知の他の順で調整することができる。スケジューラは、所定の方針に従ってデータを転送するように調整することができる。例えば、方針は、フロー又はポートを優先度付きでタグ付けするメカニズムを含むことができる。I/Oプロセッサ306は、あるレベルの packets 検査及び分類を実行してフローを差別化することができる。厳格な優先度スケジューラ508を実現することができる。代替方法として、スケジューラ508は、最小バンド幅割り当てを有する異なるフローを提供することができる。スケジューラ508は、サーバ集約システム300のvNICを管理することにより、ある宛先用のデータを選択することができる。

【0048】

図6は、一実施形態における、サーバノードと、ランデブーメモリを有するサーバ集約システムとの間で電子通信を行う方法600のフロー図である。方法600は、図3の一つ以上のサーバノード312のメモリ104及び/又はサーバ集約システム300に記憶された命令に従うことができる。従って、図2~5が参照される。方法600のいくつか又はすべてを、ASIC、システムオンチップ(SOC)又は関連するデバイスにて、オンダイで実行することができる。

【0049】

ブロック602では、ソースサーバノード312-1は、セル、パケット、フレーム又は他のローカルメモリ104-1内のデータのユニットが、宛先(例えば、他のサーバノ

10

20

30

40

50

ード 3 1 2 - N 又はリモート計算デバイス 3 5 2) への N I C 3 4 2 を介した出力に使用可能であると告知する。ソースサーバノード 3 1 2 - 1 は、伝送用のデータがメモリ 1 0 4 - 1 内にあって伝送用に使用可能であることを、v N I C 3 3 4 に通知することができる。サーバノード 3 1 2 は、データを処理するためにメモリ 1 0 4 - 1 内で定義されているリングバッファ等を維持して、データ伝送又は受信処理を実行することができる。

【 0 0 5 0 】

ひし型の判断部 6 0 4 では、ランデブーメモリ 3 0 8 が、ソースサーバノード 3 1 2 - 1 のメモリ 1 0 4 - 1 内の伝送用データを受信することができるかどうかの判定が行われる。この判定は、ランデブー管理モジュール 3 0 4 によって行うことができる。ランデブー管理モジュール 3 0 4 は、ソースサーバノード 3 1 2 - 1 から直接的にデータを伝送することができるバッファ空間が、ランデブーメモリ 3 0 8 内に空いている（使用可能である）かどうかを確立することができる。ランデブーメモリ 3 0 8 がデータを受信することができない場合には、次にブロック 6 0 8 において、ランデブー管理モジュール 3 0 4 は、受信バッファがポストされるまで待機し得る。一方、ランデブーメモリ 3 0 8 がデータを受信することができる場合には、ブロック 6 0 6 にて、データ、記述子等を、ソースサーバノード 3 1 2 - 1 のメモリ 1 0 4 - 1 からランデブーメモリ 3 0 8 に転送することができる。

【 0 0 5 1 】

ひし形の判断部 6 1 0 では、I / O プロセッサ 3 0 6 は、ソースサーバノード 3 1 2 - 1 からランデブーメモリ 3 0 8 に転送されるデータが、宛先サーバノード（例えば、サーバノード 3 1 2 - N ）に転送されることになるかどうかを判定する。この判定は、宛先アドレスに基づいて行われてもよいし、伝送記述子又は宛先ポートを識別するパケットに関するメタ情報をプリペンド（p r e p e n d ）することによって行われてもよい。データを宛先サーバノードに転送しないという判定結果であれば、次にブロック 6 1 2 において、I / O プロセッサ 3 0 6 がデータを消費することができ、例えば、データはノプロセッサ 3 0 6 において終了する。代替方法として、I / O プロセッサ 3 0 6 は、データを、データ内容に応じて（例えば、宛先アドレスに基づいて）N I C 3 4 2 に転送する。例えば、データパケットは、スイッチファブリック 3 0 2 を介してフロー管理に関連付けられることが可能であり、I / O プロセッサ 3 0 6 の構成要素、データのルーティングを制御する処理複合体と相互作用する。そうでなければ、ひし型の判断部 6 1 4 では、宛先サーバノード 3 1 2 - N が、ランデブーメモリ 3 0 8 から伝送可能な状態のデータを受信するローカルメモリ 1 0 4 - N に空き空間を有しているかどうかという判定がなされる。この判定は、コントロールプレーンプロセッサ 3 1 0 が行うことができ、この場合には、仮想 N I C 3 3 4 が、宛先ローカルメモリ 1 0 4 - N においてバッファ空間が空いているかどうかを確定することができる。判断部 6 1 4 にて「いいえ」の場合には、ブロック 6 1 6 において、ランデブー管理モジュール 3 0 4 は、宛先サーバノード C P U 1 0 2 によって受信バッファがポストされるまで、宛先ローカルメモリ 1 0 4 - N へのデータ送信を待機することができる。判断部 6 1 4 にて「はい」の場合には、ブロック 6 1 8 において、ランデブーメモリ 3 0 8 にてデータが転送用として使用可能であって、且つ、ランデブーメモリ 3 0 8 用に受信バッファが割り当てられた、との判定がなされた場合に、ランデブーメモリ 3 0 8 から宛先サーバノードメモリ 1 0 4 - N にデータを転送することができる。

【 0 0 5 2 】

図 7 は、一実施形態における、サーバノードと、ランデブーメモリを有するサーバ集約システムとの間で電子通信を行う方法 7 0 0 のフロー図である。方法 7 0 0 は、図 3 の一つ以上のサーバノード 3 1 2 のメモリ 1 0 4 及びノ又はサーバ集約システム 3 0 0 に記憶された命令に従うことができる。従って、図 2 ~ 5 が参照される。方法 7 0 0 のいくつか又はすべては、A S I C、システムオンチップ（S O C ）又は関連する集積回路にて、オンダイで実行することができる。

【 0 0 5 3 】

まず、サーバノード 3 1 2 - 1 は、セル、パケット、フレーム又は他のデータのユニッ

10

20

30

40

50

トがローカルメモリ104-1にあって、宛先（例えば、他のサーバノード312-N）又はNIC342を介してリモート計算デバイス352に出力されるかどうかの判定を行う。ローカルメモリ104-1は、データを処理するためのキュー、リングバッファ、リンクされたりリスト等を含む。

【0054】

ブロック702では、サーバノード312-1のCPU102は、サーバノード312-1から出力されるローカルメモリ104-1内のデータに関する記述子又は関連情報を生成することができる。記述子は、ランデブーメモリ308に転送されるデータに関連したコマンド、分散収集リスト等を含むことができる。関連する記述子情報は、データが常駐するメモリ104-1における記憶位置、データが伝送される宛先アドレス、移動させるデータのバイト数、及び/又は、CPU102と他のサーバ集約システム300におけるvNIC334との間で通信を確立するための関連情報を識別することができる。

10

【0055】

ブロック704では、第1のサーバノード312-1のCPU102は、新規記述子がメモリ104-1にあって伝送に使用可能であるということを、vNIC334に通知する。CPU102は、メールボックス書き込み又は関連イベントの通知を、通信経路314を介してvNIC334に送信することによって、伝送用データの使用可能性をvNIC334に通知することができる。

【0056】

ブロック706では、I/Oプロセッサ306は、ランデブーメモリ308内の使用可能なデータバッファ（空きデータバッファ）を示す受信記述子をポストする。受信記述子は、バッファのアドレス、長さ又は関連情報を含むことができる。

20

【0057】

ブロック708では、I/Oプロセッサ306は、vNIC334用の一つ以上のランデブーメモリバッファの使用可能性（空き状況）を、ランデブー管理モジュール304に通知する。

【0058】

ブロック710では、コントロールプレーンプロセッサ310は、ランデブーメモリ308の一つ以上のバッファ記憶位置にvNIC334用の空きを要求するように、ランデブー管理モジュール304にメッセージを送信する。

30

【0059】

ブロック712では、ランデブー管理モジュール304は、ランデブーメモリ308に対して、vNIC334用の一つ以上の受信バッファを割り当てて、vNIC334が、サーバノードメモリ104-1からランデブーメモリ308にvNIC334がデータを転送、コピー又は移動させることができるようにする。ランデブー管理モジュール304は、バッファ割り当てに関する通知をコントロールプレーンプロセッサ310に送信する。現在空いている受信バッファがない場合には、ランデブー管理モジュール304は、受信バッファが空くまで待機することができる。記述子及び/又はデータは、この待ちの期間にフェッチされない。一実施形態では、入力記述子は、待ち時間を削減するためにプリフェッチ（prefetch）され得る。これは、入力記述子が、転送待ち時間を削減するために多くのメモリ資源を消費するということがないためである。

40

【0060】

ブロック714では、vNIC334は、データをフェッチするのに使用されるサーバノードメモリ104-1から伝送記述子をフェッチする。vNIC334は、一つ以上の記述子をフェッチする、又は、フェッチするための一連の記述子に従うことができる。vNIC334は、サーバノードメモリ104-1からの記述子情報（例えば、アドレス）に従ってデータをフェッチし、データを、スイッチファブリック302を介して、ランデブー管理モジュール304によって割り当てられた空き（使用可能な）ランデブーメモリ308に移動させる。例えば、コントロールプレーンプロセッサ310は、PCIeコネクタ314を介してメモリ104-1に読み出しを発令し、データをフェッチしてラン

50

デブーメモリ 308 に移動させることができる。

【0061】

ブロック 716 では、コントロールプレーンプロセッサ 310 は、ランデブー管理モジュール 304 に対して、フェッチされたデータがランデブーメモリ 308 において使用可能であるとの通知を送信することができる。ランデブー管理モジュール 304 は、I/O プロセッサ 306 に伝送される割り込みを生成することができる。これにより、I/O プロセッサ 306 は、データをランデブーメモリ 308 からプル (pull) し、当該データを提供することができる。

【0062】

ブロック 718 では、I/O プロセッサ 306 は、ランデブーメモリ 308 内のデータを処理して、例えば、データを消費するかどうか、NIC 342 に当該データを転送するかどうか、又は、計算インフラストラクチャー 30 の一部である他のサーバノード 312 の CPU 102 に当該データを転送するかどうかを判定することができる。

10

【0063】

示されていないが、コントロールプレーンプロセッサ 310 は、サーバノード 312 - 1 で生成された伝送記述子を解放することができる。コントロールプレーンプロセッサ 310 は、任意に、サーバノード CPU 102 に割り込むことができる。I/O プロセッサ 306 は、例えば、スケジューラ 508 に従って、ランデブーメモリ 308 内のデータの送信を制御することができる。このように、I/O プロセッサ 306 は、その宛先 (例えば、宛先サーバノード 312 - N) へのデータの送信を保証することができ、これは図 8 に記載されるとおりである。

20

【0064】

図 8 は、一実施形態における、ランデブースイッチングを行う構成のサーバ集約システムに結合した第 1 及び第 2 のサーバノード間で通信を確立する方法 800 のフロー図である。方法 800 は、サーバ集約システムのネットワークプロセッサによる図 7 のブロック 718 での判定に応じて実行することができる。ネットワークプロセッサは、ランデブーメモリ 308 内の使用可能なフェッチされたデータを、他のサーバノード、又は、CPU を有する関連する電子デバイス (例えば、図 3 に示されたサーバノード 312 - N) に転送する。方法 800 は、図 3 のサーバノード 312 のメモリ 104 及び / 又はサーバ集約システム 300 に記憶された命令により制御することができる。この場合には、図 2 ~ 5 及び図 7 が参照される。

30

【0065】

ブロック 802 では、I/O プロセッサ 306 の記述子処理モジュール 504 は、図 7 に記載される方法 700 に従ってランデブーメモリ 308 に移動されるデータに関する記述子又は関連する情報を提供することができる。記述子は、宛先サーバノード 312 - N のローカルメモリ 104 - N にデータを転送する命令、分散収集リスト等を含むことができる。関連する記述子情報は、データが常駐するメモリ 104 - 1 における記憶位置、データが伝送される宛先アドレス、移動させるデータのバイト数、及び / 又は、I/O プロセッサ 306 と、ローカルメモリ 104 - N へのデータの転送に関与する vNIC 334 との間で通信を確立する他の関連情報を識別することができる。

40

【0066】

ブロック 804 では、I/O プロセッサ 306 は、ランデブーメモリ 308 内のデータが伝送用に使用可能であることを、ランデブー管理モジュール 304 に通知する。

【0067】

ブロック 806 では、宛先サーバノード 312 - N の CPU 102 は、宛先サーバノード 312 - N のメモリ 104 - N 内の使用可能なデータバッファ (空きデータバッファ) を示す受信記述子をポストする。受信記述子は、バッファのアドレス、長さ又は関連情報を含むことができる。

【0068】

ブロック 808 では、宛先サーバノード 312 - N の CPU 102 は、コントロールプ

50

レーンプロセッサ310に対して、メモリ104 - Nにおけるランデブー管理モジュール304用の一つ以上のホストメモリバッファの使用可能性（空き状況）を通知する。

【0069】

ブロック810では、宛先サーバノード312 - N内のCPU102は、コントロールプレーンプロセッサ310に要求を送信して、宛先ノードメモリ104 - Nにおける一つ以上のバッファ記憶位置を空きにするようにランデブー管理モジュール304に要求する。

【0070】

ブロック812では、コントロールプレーンプロセッサ310は、ランデブー管理モジュール304用に宛先ノードメモリ104 - Nの受信バッファを割り当てて、ランデブー管理モジュール304が、ランデブーメモリ308から宛先ノードメモリ104 - Nへのデータの転送、コピー等を行えるようにする。コントロールプレーンプロセッサ310は、バッファ割り当てに関する通知を、ランデブー管理モジュール304に送信することができる。現在空いている受信バッファがない場合には、ランデブー管理モジュール304は、宛先サーバノードメモリ104 - Nの受信バッファが空くまで待機することができる。一実施形態では、記述子及び/又はデータは、待ち期間の間フェッチされない。

【0071】

ブロック814では、ランデブー管理モジュール304は、ランデブーメモリ308からのデータをフェッチするのに使用される伝送記述子をフェッチする。ランデブー管理モジュール304は、ランデブーメモリ308からの記述子情報（例えば、アドレス）に従ってデータをフェッチし、データを、スイッチファブリック302を介して、vnic634に移動させる。ランデブー管理モジュール304は、コントロールプレーンプロセッサ310に対して、宛先サーバノード312 - N用のランデブーメモリ308からvNIC334への記述子及び/又はデータをフェッチした後に新規データが使用可能であることを、通知することができる。

【0072】

ブロック816では、vNIC/CPは、受信バッファ用の記述子をフェッチ及び処理し、データを宛先サーバノードメモリ104 - Nに移動させる。コントロールプレーンプロセッサ310は、メモリ104 - Nに移動したデータを提供するために、宛先サーバノード312 - NのCPU102に出力される割り込みを生成することができる。ランデブー管理モジュール304は、伝送記述子を解放することができ、サーバノードI/Oプロセッサ306に任意に割り込みすることができる。

【0073】

本明細書に記載の方法の実施形態により、データが宛先に送信できる状態にあるとサーバ集約システム300が判定した場合には、要求に応じてデータを転送することができる。ランデブーメモリ308がデータを受信できる状態にない場合には、データは宛先に伝送されない。これは、仮に、データがソースサーバノード312 - 1において転送できる状態にあってもそうである。同様に、宛先サーバノード312 - Nがデータを受信できる状態にない場合には、データは、ランデブーメモリ308から伝送されない。このように、従来のキューイング技術を必要とすることなく、バッファ空間がないことに起因してデータが失われるということはない。むしろ、データの移動は記述子の使用可能性と結びついている。記述子が確立される場合には、サーバ集約システム300は、移動させることになるデータが存在するとの判定を行うことができる。

【0074】

図9は、一実施形態における、ランデブースイッチングを行う構成のサーバ集約システムに結合した第1及び第2のサーバノード間で電子通信する方法900のフロー図である。方法900は、一つ以上のサーバノード312のメモリ104及び/又は上記のサーバ集約システム300に記憶された命令により制御することができる。方法900の一つ以上の構成要素は、上述したものと同様であってもよい。図2～5が参照されるが、方法900は、ランデブーメモリ308無しに実行される。特に、ランデブー管理モジュール3

10

20

30

40

50

04及びコントロールプレーンプロセッサ310は、方法900に關与することができ、この方法では、データ転送が、二つのサーバノード212-1, 212-Nの間で、ランデブーメモリ308の介在無しに行われる。

【0075】

ブロック902では、記述子は、第1のサーバノード312-1で生成される。

【0076】

ブロック904では、vNIC334は、データが伝送用に使用可能であることを通知される。

【0077】

ブロック906では、受信記述子は、第2のサーバノード312-Nの使用可能なバッファ（空きバッファ）でポストされる。

10

【0078】

ブロック908では、宛先サーバノード312-Nは、コントロールプレーンプロセッサ310に対して、ランデブー管理モジュール304用のメモリ104-Nの一つ以上のホストメモリバッファの使用可能性（空き状況）を通知する。

【0079】

ブロック910では、宛先メモリ104-Nのバッファの要求がなされる。

【0080】

ブロック912では、受信バッファが割り当てられる。バッファは、vNIC334用に割り当てることができ、これにより、コントロールプレーンプロセッサ310のvNIC334が、サーバノードメモリ104-1から宛先ノードメモリ104-Nにデータを転送、コピー又は移動させることができる。

20

【0081】

ブロック914では、vNIC334は、受信バッファ用の記述子をフェッチ及び処理して、データを宛先サーバノードメモリ104-Nに移動させる。

【0082】

当業者が認識し得るように、本発明の態様は、システム、方法又はコンピュータプログラム製造物として具体化されてもよい。従って、本発明の態様は、完全にハードウェアの実施形態、完全にソフトウェアの実施形態（ファームウェア、常駐ソフトウェア、マイクロコード等を含む）、又は、ソフトウェア及びハードウェアの態様を組み合わせた実施形態であって、本明細書では「回路」、「モジュール」又は「システム」と称される形態であってもよい。さらには、本発明の態様は、コンピュータプログラム製品の形態であって、一つ以上のコンピュータ可読媒体で具体化されてもよい。コンピュータ可読媒体は、コンピュータ可読プログラムコードを有していてもよい。

30

【0083】

一つ以上のコンピュータ可読媒体のあらゆる組合せが利用できる。コンピュータ可読媒体は、コンピュータ可読信号媒体、又は、コンピュータ可読記憶媒体であってもよい。コンピュータ可読記憶媒体は、例えば、電氣的、磁氣的、光学的、電磁氣的、赤外線であってもよく、半導体システム、装置若しくはデバイスであってもよく、これらのあらゆる好適な組合せであってよいが、これらに限定されるものではない。コンピュータで可読記憶媒体のさらに特定の例（非包括的な列挙）は、以下のものを含む。すなわち、一つ以上の配線を有する電氣的接続、ポータブルコンピュータディスク、ハードディスク、ランダムアクセスメモリ（RAM）、リードオンリーメモリ（ROM）、イレーザブルプログラマブルリードオンリーメモリ（EPROM若しくはフラッシュメモリ）、光ファイバー、ポータブルコンパクトディスクリードオンリーメモリ（CD-ROM）、光記憶デバイス、磁気記憶デバイス、又は、これらのあらゆる好適な組合せである。本文献の文脈においては、コンピュータ可読記憶媒体は、あらゆる有形の媒体であって、命令実行システム、装置若しくはデバイスにより使用されるプログラム、又は、これらに接続して使用するプログラムを含む又は記憶するものであってもよい。

40

【0084】

50

コンピュータ可読信号媒体は、その内部で具体化されるコンピュータ可読プログラムコード（例えば、ベースバンド又は搬送波の一部として伝搬されるデータ信号）を含んでもよい。このような伝搬信号は、あらゆる多様な形態（例えば、電磁的、光学的又はこれらのあらゆる好適な組合せの形態）をとってもよいが、これらには限定されない。コンピュータ可読信号媒体は、コンピュータで読み取り可能なあらゆる媒体であって、コンピュータ可読記憶媒体ではなく、命令実行システム、装置若しくはデバイスによって使用されるプログラム、又は、これらと接続して使用するプログラムを、通信、伝搬又は移動することが可能な媒体であってもよい。コンピュータ可読媒体上で具体化されるプログラムコードは、あらゆる適切な媒体、例えば、限定はされないが、無線、有線、光ファイバケーブル、RF等又は前述のあらゆる好適な組合せを用いて伝送することができる。

10

【0085】

本発明の態様の動作を実行するコンピュータプログラムコードは、Java（登録商標）、Smalltalk、C++等のオブジェクト指向プログラミング言語及び「C」プログラミング言語等の従来の手続型言語、又は、同様のプログラミング言語を含む一つ以上のプログラミング言語のあらゆる組合せで記述されていてもよい。プログラムコードは、全部がユーザのコンピュータ上で実行されてもよいし、一部がユーザのコンピュータ上で実行されてもよい。また、プログラムコードは、スタンドアロンのソフトウェアパッケージとして、一部がユーザのコンピュータ上又はリモートコンピュータ上で実行されてもよいし、全部がリモートコンピュータ上又はサーバ上で実行されてもよい。後者のシナリオでは、リモートコンピュータは、ローカルエリアネットワーク（LAN）又はワイドエリアネットワーク（WAN）を含むあらゆるタイプのネットワークを介してユーザのコンピュータに接続されていてもよいし、（例えば、インターネットサービスプロバイダを用いたインターネットを介して）外部のコンピュータと接続されていてもよい。

20

【0086】

本発明の態様は、本発明の実施形態の方法、装置（システム）及びコンピュータプログラム製品のフローチャート並びに／又はブロック図を参照しつつ記載されている。フローチャート図及び／又はブロック図の各ブロックと、フローチャート図及び／又はブロック図におけるブロックの組合せとは、コンピュータプログラムの命令により実現可能であることは理解されるであろう。コンピュータプログラムのこれらの命令は、汎用コンピュータ、専用コンピュータ、又は、機械を製造する他のプログラマブルデータ処理装置のプロセッサに提供され、これらのコンピュータ又は他のプログラマブルデータ処理装置のプロセッサにおいて実行され、フローチャート及び／又はブロック図のブロックにおいて特定された機能／作用を実現する手段を生成するようにしてもよい。

30

【0087】

また、コンピュータプログラムのこれらの命令は、コンピュータ可読媒体に記憶されてもよく、この媒体が、コンピュータ、他のプログラマブルデータ処理装置、又は、特定の方法で機能する他のデバイスに命令して、コンピュータ可読媒体に記憶された命令によって製品が製造され、その製品に、フローチャート及び／又はブロック図のブロックにおいて特定された機能／作用を実現する命令が含まれるようにしてもよい。さらに、コンピュータプログラムの命令は、コンピュータ、他のプログラマブルデータ処理装置、又は、他のデバイスにロードされて、コンピュータ、他のプログラマブル装置、又は、他のデバイス上で一連の処理ステップを実行させ、コンピュータで実現されるプロセスを生成させ、コンピュータ又は他のプログラマブル装置上で実行されるこれらの命令が、フローチャート及び／又はブロック図のブロックにおいて特定される機能／作用を実現するプロセスを提供するようにしてもよい。

40

【0088】

図におけるフローチャート及びブロック図は、本発明の様々な実施形態によるシステム、方法及びコンピュータプログラム製品の実施可能なアーキテクチャ、機能性及び動作を例示する。この点で、フローチャート又はブロック図における各ブロックは、モジュール、セグメント又はコードの一部を表してもよく、これらは、特定された論理機能を実現す

50

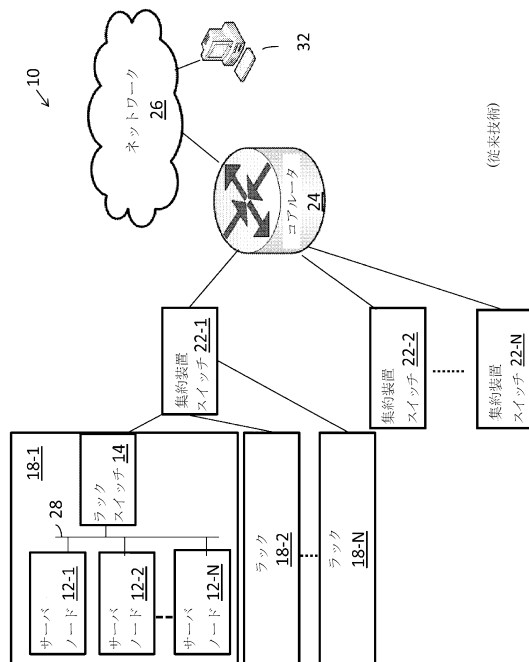
る一つ以上の実行可能な命令を含む。いくつかの代替の実施物では、ブロックで指摘される機能は、図において指摘される順序に従うことなく生じてよいことも留意されたい。例えば、連続して示された二つのブロックは、実際には、実質的に同時に実行されてもよいし、ブロックは、場合によっては逆順に実行されてもよく、これは関係する機能による。ブロック図及び／又はフローチャートの各ブロックと、ブロック図及び／又はフローチャートにおけるブロックの組合せとは、専用のハードウェア及びコンピュータ命令の特定された機能若しくは作用又は組合せを実行する、専用のハードウェアをベースにしたシステムによって実現されてもよいことも留意されるであろう。

【 0 0 8 9 】

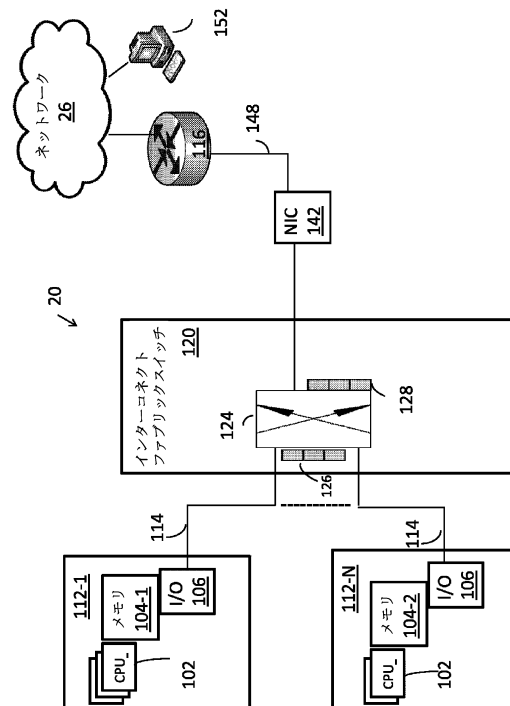
特定の実施形態を参照して、本発明を示し記載してきたが、形態及び詳細の様々な変更が、本発明の精神及び範囲を逸脱することなく、当業者により行われてもよいことは理解されるべきである。

10

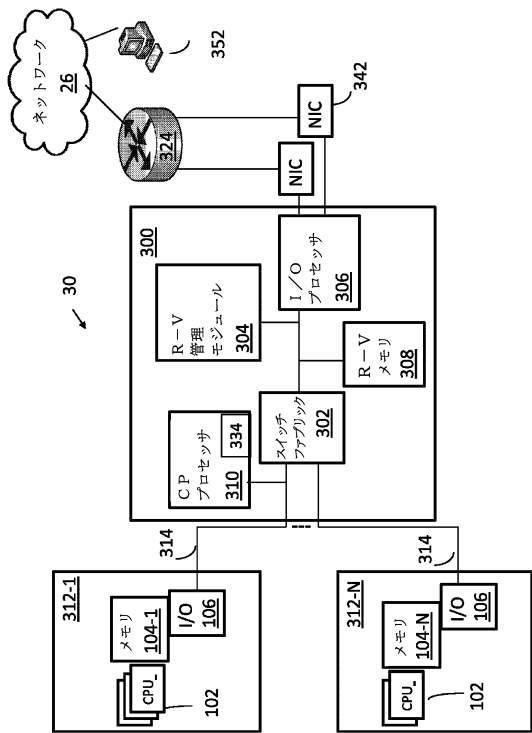
【 図 1 】



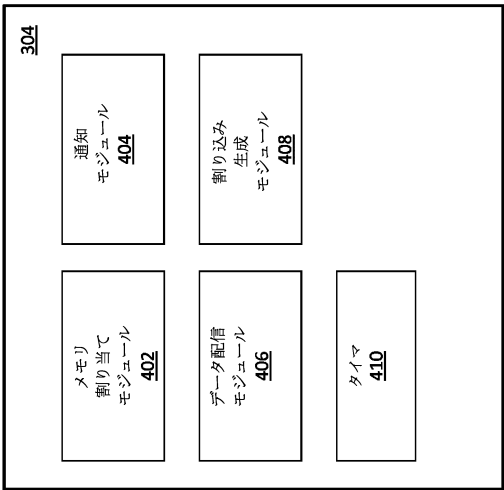
【 図 2 】



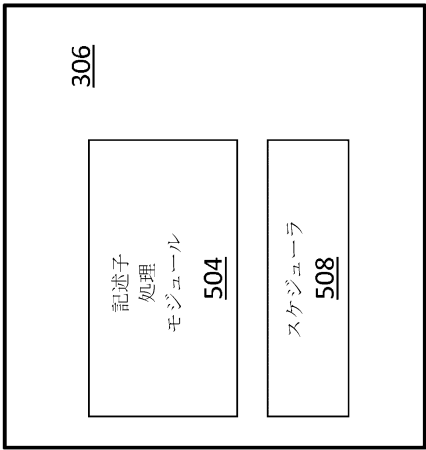
【図 3】



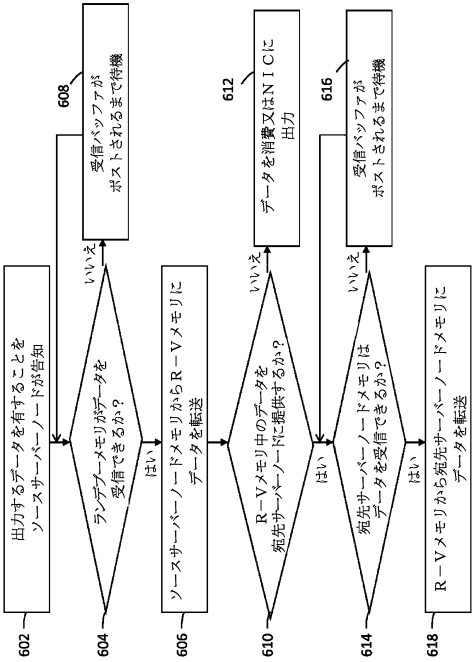
【図 4】



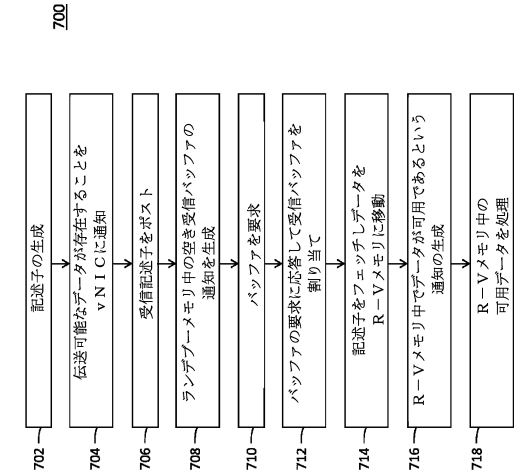
【図 5】



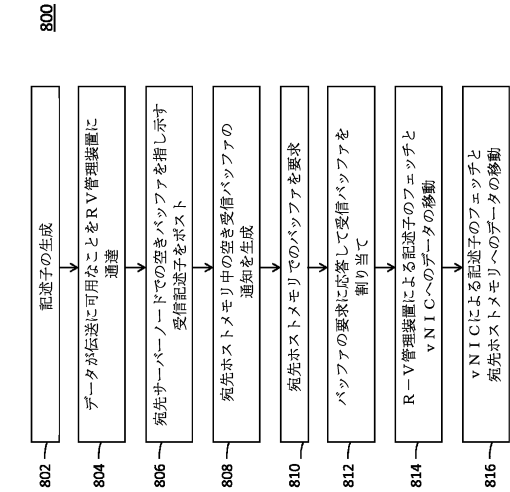
【図 6】



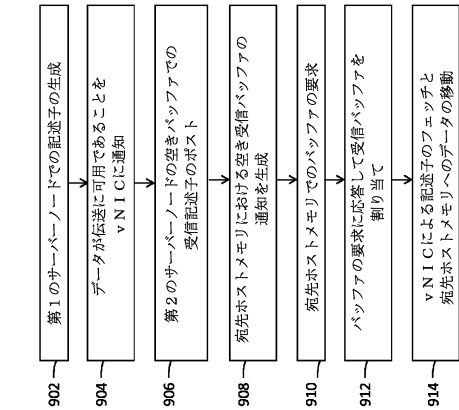
【図 7】



【図 8】



【図 9】



フロントページの続き

- (72)発明者 マーク ハメル
アメリカ合衆国 02038 マサチューセッツ州、フランクリン、スチュワート ストリート 68
- (72)発明者 デイビッド メイヒュー
アメリカ合衆国 01532 マサチューセッツ州、ノースボロー、プレザント ストリート 159
- (72)発明者 マイケル オズボーン
アメリカ合衆国 03049 ニューハンプシャー州、ホリス、ブラック オーク ドライブ 50

審査官 衣鳩 文彦

- (56)参考文献 特開平10-229429(JP,A)
特開2005-204089(JP,A)
米国特許出願公開第2005/0041510(US,A1)
米国特許出願公開第2011/0202701(US,A1)
米国特許出願公開第2003/0208631(US,A1)
特開2009-282917(JP,A)
特開平03-082244(JP,A)

- (58)調査した分野(Int.Cl., DB名)
H04L 12/931
G06F 15/173