

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4516395号
(P4516395)

(45) 発行日 平成22年8月4日 (2010.8.4)

(24) 登録日 平成22年5月21日 (2010.5.21)

(51) Int. Cl.	F I
G 0 6 F 13/38 (2006.01)	G O 6 F 13/38 3 1 O D
G 0 6 F 12/02 (2006.01)	G O 6 F 12/02 5 4 O
G 0 6 F 12/06 (2006.01)	G O 6 F 12/06 5 2 2 C
H 0 4 L 13/08 (2006.01)	G O 6 F 12/06 5 5 O B
	H O 4 L 13/08

請求項の数 10 外国語出願 (全 37 頁)

(21) 出願番号	特願2004-286229 (P2004-286229)	(73) 特許権者	596092698
(22) 出願日	平成16年9月30日 (2004.9.30)		アルカテルルーセント ユーエスエー
(65) 公開番号	特開2005-174286 (P2005-174286A)		インコーポレーテッド
(43) 公開日	平成17年6月30日 (2005.6.30)		アメリカ合衆国 07974 ニュージャ
審査請求日	平成19年9月28日 (2007.9.28)		ーシー, マレイ ヒル, マウンテン アヴ
(31) 優先権主張番号	10/699315	(74) 代理人	100064447
(32) 優先日	平成15年10月31日 (2003.10.31)		弁理士 岡部 正夫
(33) 優先権主張国	米国 (US)	(74) 代理人	100085176
			弁理士 加藤 伸晃
		(74) 代理人	100106703
			弁理士 産形 和央
		(74) 代理人	100096943
			弁理士 臼井 伸一

最終頁に続く

(54) 【発明の名称】 リンク・リスト・プロセッサを持つメモリ管理システム

(57) 【特許請求の範囲】

【請求項 1】

リンク・リスト・データ・ファイルを処理するように適合されたメモリ管理システムを操作する方法であって、前記システムは複数の高速小容量メモリおよび低速大容量メモリを備え、前記高速メモリは第1のデータ転送速度を持ち、前記大容量メモリは前記第1のデータ転送速度よりも遅い第2のデータ転送速度を持ち、前記システムはさらに、前記メモリによりリンク・リストの読み書きの要求を生成するためのアクセス・フロー・レギュレータを備え、該方法が、

前記アクセス・フロー・レギュレータから前記高速メモリに書き込み要求を送信することにより、前記高速メモリへのリンク・リストの書き込みを開始するステップ、

前記リンク・リストの先頭バッファおよび末尾バッファおよび少なくとも1つの中間バッファを前記高速メモリに書き込むステップ、及び

前記高速メモリから前記大容量メモリに前記少なくとも1つの中間バッファを転送する一方で、先頭バッファおよび末尾バッファを前記高速メモリ内に残すステップからなる方法。

【請求項 2】

リンク・リスト・データ・ファイルを処理するように適合されたメモリ管理システムを操作する方法であって、前記システムは複数の高速小容量メモリおよび低速大容量メモリを備え、前記高速メモリは第1のデータ転送速度を持ち、前記大容量メモリは前記第1のデータ転送速度よりも遅い第2のデータ転送速度を持ち、前記システムはさらに、前記メ

10

20

モリによりリンク・リストの読み書きの要求を生成するためのアクセス・フロー・レギュレータを備え、

指定されたリンク・リストの読み出し要求を前記アクセス・フロー・レギュレータから前記指定されたリンク・リストのバッファを含む前記高速メモリに送信するステップ、

前記指定されたリンク・リストの先頭バッファを読み込むステップ、

前記リンク・リストの前記少なくとも1つの中間バッファを前記大容量メモリから前記高速メモリのうちの1つに転送するステップ、

前記1つの高速メモリに転送された前記中間バッファを前記指定されたリンク・リストの交換先頭バッファとして指定するステップ、

前記1つの高速メモリから前記指定されたリンク・リストの前記中間バッファを読み出すステップ、及び

前記指定されたリンク・リストの前記読み出されたバッファを前記アクセス・フロー・レギュレータに送信するステップ

からなる方法。

【請求項3】

請求項1又は2の方法であって、

前記システムを前記アクセス・フロー・レギュレータからの複数の要求に対するリンク・リストを同時処理するように動作させるステップ、

前記システムを前記複数の高速メモリのうちの異なるメモリに格納されているリンク・リストのバッファを処理するように動作させるステップ、

前記システムを、末尾バッファを新しいリンク・リストに書き込まれた第1のバッファとして書き込むように動作させるステップ、及び

リンク・リストの先頭バッファを最初に読み込むステップ
をさらに含む方法。

【請求項4】

請求項1又は2の方法において、前記システムはさらに、複数の状態コントローラ(1804)を備え、それぞれの状態コントローラは前記複数の高速メモリのうちの対応するメモリに対する個別のコントローラであり、前記システムはさらに、前記アクセス・フロー・レギュレータを前記状態コントローラと接続する要求バス(1802)を備え、読み出し要求を送信する前記ステップが、

前記アクセス・フロー・レギュレータを動作させ、前記読み出し要求を受信すべきアイドル状態の高速メモリを選択するステップ、

前記読み出し要求を前記アクセス・フロー・レギュレータから前記要求バスを介して前記選択された高速メモリに対して個別の状態コントローラに送信するステップ、

前記状態コントローラを、前記選択された高速メモリの現在の占有レベルを判別するように動作させるステップ、

前記現在の占有レベルが所定のレベル以下の場合に前記要求を前記高速メモリに送信するステップ、及び

前記選択された高速メモリの前記現在の占有レベルが前記所定のレベルを超えた場合に前記大容量メモリへの接続を要求するステップ
を含み、

前記システムがさらに、バックグラウンド・アクセス・マルチプレクサ、および前記状態コントローラを前記マルチプレクサと接続するアクセス・バスを備え、前記システムがさらに、前記マルチプレクサを前記大容量メモリと接続するバスを備え、前記方法が、前記マルチプレクサを、

前記大容量メモリへの接続の要求を前記状態コントローラから受信し、

複数の要求状態コントローラのどれかに対し前記大容量メモリへのアクセスを許可するかを決定するように動作させるステップ、

前記1つの要求状態コントローラを前記大容量メモリに接続するステップ、

前記1つの高速メモリから前記大容量メモリにデータを転送する際の前記大容量メモリ

10

20

30

40

50

のオペレーションを制御するステップ、及び

前記アクセス・パスを介して前記状態コントローラから前記マルチプレクサに前記リンク・リストのバッファを転送するステップを含む方法。

【請求項 5】

請求項 1 又は 2 の方法において、前記大容量メモリから前記バッファを転送する前記ステップが、

リンク・リストの中間バッファを前記大容量メモリから前記高速メモリに、前記高速メモリのデータ転送速度に実質的に等しいデータ転送速度を持つバースト・モードで転送するステップ、

前記読み出されたバッファを前記高速メモリに格納するステップ、

その後、前記高速メモリから前記リンク・リストのバッファを読み出し、前記アクセス・フロー・レギュレータに転送するステップ、

バッファを既存のリンク・リストに書き込むのに、前記高速メモリから前記大容量メモリに前記既存のリンク・リスト既存の末尾を転送することにより書き込むステップ、及び

新しいバッファを前記既存のリンク・リストの新しい末尾バッファとして前記高速メモリに書き込むステップ

を含む方法。

【請求項 6】

請求項 4 記載の方法において、前記状態コントローラを動作させる前記ステップが、さらに、

複数のリンク・リストのバッファを同時に受け取るステップ、

前記アクセス・フロー・レギュレータに送られる複数のリンク・リストのバッファを分離するステップ、

前記アクセス・フロー・レギュレータによって受け取られた複数のアクセスを前記高速メモリに送るステップ、

前記アクセス・フロー・レギュレータからそれぞれの受信された要求に応答し、前記状態コントローラに個別の高速メモリの現在の占有レベルを判別するステップ、

前記現在の占有レベル以下の場合、前記関連する高速メモリに前記アクセスを送るステップ、

前記現在の占有レベルを超えた場合、前記要求をバッファリングするよう前記アクセス・フロー・レギュレータに信号を送るステップ、

前記高速メモリからバッファをバースト・モードで前記大容量メモリに転送するのを制御するステップ、

前記大容量メモリからバッファを前記高速メモリに転送するのを制御するステップ、

転送が要求されたときに前記大容量メモリがアイドル状態かどうか判別するステップ、

アイドル状態の場合前記バッファを前記大容量メモリに送るステップ、及び

前記大容量メモリが使用中の場合、前記転送をバッファリングするステップ

を含む方法。

【請求項 7】

請求項 4 の方法において、前記マルチプレクサを動作させる前記ステップが、さらに、複数のビッディング高速メモリのどれに対し前記大容量メモリへのアクセスを許可するかを決定するステップ、

他のバッファの要求を前記 1 つの高速メモリにバッファリングするステップ、

前記大容量メモリからのバッファの転送先の高速メモリに対する ID を決定するステップ、及び

前記大容量メモリから前記識別された高速メモリにバースト・モードで前記バッファを転送するのを制御するステップ

を含む方法。

【請求項 8】

請求項 1 又は 2 の方法であって、さらに、

それぞれの前記高速メモリの使用中／アイドル状態を示す、それぞれの高速メモリに固有の信号を発生するステップ、

それぞれの発生した信号を前記アクセス・フロー・レギュレータに送るステップ、

前記アクセス・フロー・レギュレータを、前記高速メモリによるリンク・リストの読み書きの要求を受け取るように動作させるステップ、

要求を受け取ったことに対して応答して前記アクセス・フロー・レギュレータを動作させることにより、前記高速メモリにより生成される前記使用中／アイドル状態信号を読み出すステップ、

前記読み出しに응答して前記アクセス・フロー・レギュレータを、前記いくつかの高速メモリのうちアイドル状態のメモリを識別するように動作させるステップ、

前記アクセス・フロー・レギュレータを、データ・ファイルの読み書きの要求を前記アイドル状態の1つの高速メモリに送るように動作させるステップからなる方法。

【請求項 9】

リンク・リストのデータ・ファイルを処理するように適合されているメモリ管理システム(1800)であって、前記システムが、

複数の高速小容量メモリ(1803-1、1803-2、1803-3、1803-4、1803-5、1803-6、1803-7、1803-8)および低速大容量メモリ(806-1、806-2、806-3)であって、前記高速メモリは第1のデータ転送速度を持ち、前記大容量メモリは前記第1のデータ転送速度よりも遅い第2のデータ転送速度を持つメモリ、

前記メモリ(1803、1806)によるリンク・リストの読み書きの要求を生成するためのアクセス・フロー・レギュレータ(1801)、

書き込み読み出し要求を前記複数の高速メモリ(1806-1、1806-2、1806-3、1806-4、1806-5、1806-6、1806-7、1806-8)のうちのアイドル状態のメモリに送ることにより前記メモリ内のリンク・リストの書き込みを開始するための装置(1819)、

前記リンク・リストの先頭バッファおよび末尾バッファおよび少なくとも1つの中間バッファを前記高速メモリに書き込むための装置(1919)、

前記高速メモリから前記大容量メモリに前記リンク・リストの前記少なくとも1つの中間バッファを転送し、その一方で前記リンク・リストの先頭バッファおよび末尾バッファを前記高速メモリに残すための装置(1804、1808、1810)、

その後前記アクセス・フロー・レギュレータからの前記リンク・リストの読み込みの要求を前記高速メモリに送るための装置(1802)、

前記リンク・リストの先頭バッファを前記高速メモリの1つに読み込むための装置(1804)、

前記リンク・リストの前記少なくとも1つの中間バッファを前記大容量メモリから前記高速メモリに転送するための装置(1808、1810、1804)、

前記高速メモリ内の転送されたバッファを新しい先頭バッファとして指定するための装置(1810)、

その後前記先頭バッファおよび前記末尾バッファだけでなく前記中間バッファを前記高速メモリから読み出すための装置(1804)、及び

前記リンク・リストの前記読み出されたバッファを前記アクセス・フロー・レギュレータに送信するための装置(1802)を備えるメモリ管理システム。

【請求項 10】

請求項 9 のメモリ管理システムであって、さらに、

それぞれの高速メモリの現在の使用中／アイドル状態を示す、それぞれの高速メモリに固有の信号を発生するための前記高速メモリ(1803)を備える装置(1804)、

前記信号を前記アクセス・フロー・レギュレータに送るための装置(1802)、

前記アクセス・フロー・レギュレータ(1801)により、前記高速メモリによるリンク・リストの読み書きの要求を受け取る装置(1814)、

前記アクセス・フロー・レギュレータ(1801)により、前記要求を受け取ったことに応答して、前記使用中/アイドル状態信号を読み込む装置(1821)、

前記アクセス・フロー・レギュレータ(1801)により、前記読み出しに対する応答して前記高速メモリのそれぞれの現在の使用中/アイドル状態を判別するための装置(1821)、及び

前記アクセス・フロー・レギュレータ(1801)により、前記メモリの1つが現在アイドル状態かどうかを判別したことに応答して、前記高速メモリによるリンク・リストの読み書きの要求を許可する装置(1821、1804)

10

を備えるメモリ管理システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、通信ネットワーク用のメモリ管理ファシリティに関するものであり、特に、ネットワーク・ノードでのブロッキングを減らすことによりネットワークのトラフィック処理能力を最適化するファシリティに関するものである。本発明は、さらに、ネットワーク・メモリ要素へのアクセス権の争奪時間を短縮することによりトラフィックの流れを改善するメモリ管理ファシリティに関するものである。本発明は、さらに、情報の格納を高速小容量メモリと低速大容量メモリとに割り当てることにより処理および争奪時間を短縮するための配置に関するものである。

20

【背景技術】

【0002】

ネットワーク・トラフィック処理能力を高めるために、マルチノード通信ネットワークをアクティブに管理する方法が知られている。予想されるトラフィックを適切に処理するため各ノードに十分なファシリティを置いたネットワークが設計されている。これは、通常のトラフィック量をこなすために必要なファシリティだけでなく、あまり頻繁には発生しないピーク・トラフィックを経済的に実現可能な範囲で処理するための追加ファシリティを設置することを含む。通信ネットワークは、通常、理論上はありえると思われるが、あったとしても発生がまれな、トラフィック・ピークを処理するために必要な数のファシリティを備えるような設計ではない。

30

【0003】

マルチノード通信ネットワークでは、ネットワークが全体として十分なレベルのトラフィックを処理するように設計されているとしてもトラフィック・ブロッキングが発生する場合がある。このブロッキングは、トラフィックの分布が不均等であることから来るもので、ネットワーク・ノードの一部または全部が法外なレベルのトラフィック発生によりオーバーロード状態になるのである。ネットワーク・ノードは、それがネットワーク接続要求の要求先である送信先ノードである場合にオーバーロード状態になる可能性がある。ネットワーク・ノードは、さらに、リンクを介して要求された送信先ノードに接続され、上流のノードから送信先ノードに送られる要求を受信した場合にオーバーロード状態になる可能性がある。不均等なトラフィック分布によるノードのオーバーロードを最小限に抑えるためネットワークにトラフィック・シェーピング・ファシリティを追加する方法が知られている。これらのトラフィック・シェーピング・ノードでは、各ノードでのトラフィックだけでなく、各ノードにより発生する接続要求も監視する。遠隔地のノードの輻輳は、送信側ノードがすでにオーバーロード状態にある遠隔地のノードへのアクセスに対して発生する可能性のある要求の数を絞ることにより防止される。

40

【0004】

マルチノード・ネットワークおよびそのトラフィック・シェーピング・ファシリティを利用すると、ネットワークは十分に低いレベルのブロッキングで通常のトラフィック・レベルを処理することができる。しかし、ネットワーク・トラフィックを管理し、制御する

50

ために必要なファシリティは、複雑で、高価であり、必要な処理オペレーションが複雑であることからネットワークのトラフィック・スループットが低下する。これらのファシリティでは、プロセッサにより制御されるリンク・リスト・エンジンをノードの入力と出力に装備し、着信および送信トラフィックをバッファリングする。リンク・リスト・エンジンのオペレーションでは、リンク・リスト・エンジン内の争奪に関する問題を最小限に抑えるために必要な複雑なデータ処理オペレーションを伴う。これらの争奪問題の複雑さゆえに、ネットワーク全体のトラフィック処理能力が低下する。

【 0 0 0 5 】

オーバーロードは、重いトラフィックが何回も発生したときに、各ノードでリンク・リスト・バッファが枯渇することで発生することがある。これにより、パケットが破棄され、システムのパフォーマンスがひどく低下することがある。バッファのオーバーロードは、バッファ・サイズが不適切であったり、着信トラフィックを処理するのに不十分な速度のバッファを使用したりして、生じる。システム設計者は、低速な大容量のバッファを使用するのか、それとも高速だが小容量のバッファを使用するのかの選択に直面している。低速な大容量のバッファだと、パケットが喪失してネットワーク・トラフィックの流れが阻害される。高速小容量バッファもまた、バッファのオーバーロードを引き起こし、重いバースト処理時に利用できるバッファが不足してパケットが破棄される。

【 0 0 0 6 】

これら両方のタイプのバッファに関連する根本の問題は、同じファシリティを使用して複数のアクセスが発生した場合に争奪問題が生じることにある。これは、例えば、特定のメモリ・バンクの読み書きアクセスに関して複数回のアクセスを受け取ったときに発生する。このような状況の下で、一方のアクセスは成功するが、他方のアクセスは、要求されたメモリ・バンクが利用可能になるのを待つ。この成功したアクセスに関連する呼び出しは、適切に処理されるが、遅延したアクセスに関連する呼び出しは、破棄されるか、または不適切な処理になる。

【 0 0 0 7 】

R A Mメモリ・バンクに対するアクセスの争奪は、不適切な数のR A Mメモリ・バンクを使用すること、および/またはメモリ・バンクを処理するために用意されている争奪ファシリティのせいで発生する。いくつかの争奪ファシリティは、アクセスが処理される速度を制限するアルゴリズムおよびプロセスに依存する。このような従来技術の配置の1つでは、アクセスからアクセスまでの間約250ナノ秒の最小遅延時間を要するアルゴリズムを使用している。これは、第1のメモリ・バンクへのアクセスの割り当ての後アクセスを受け取るために第2のR A Mメモリ・バンクが利用可能かどうかを判別するための機能を用意していないため著しい制限である。したがって、アクセスの処理に要する250ナノ秒の時間間隔では、システムのスループットは、R A Mメモリ・バンクが使用可能かどうかに関係なく、最大毎秒4,000,000回のアクセスの処理に制限される。既存の争奪の配置に関連する他の問題として、それらの多くは、複雑で、高価で、高トラフィック・レベルを処理するのには不適切なロジック・ファシリティを使用することが挙げられる。

【 発明の開示 】

【 発明が解決しようとする課題 】

【 0 0 0 8 】

本発明は、R A Mメモリ・バンクの数を増やす第1の可能な実施例により、これらの争奪問題を解消する。これはそれ自体で、争奪の可能性を減じる。本発明の持つ第2の特徴は、それぞれのR A Mメモリ・バンクに、状態コントローラと呼ばれる関連する制御要素を備えることである。この状態コントローラは、R A Mメモリ・バンクとアクセス要求の受信に使用されるシステム・バスとの間のインターフェースである。すべてのアクセス要求は、ノードによって生成されたすべてのアクセス要求を受け取るアクセス・フロー・レギュレータによりシステム・バスに適用され、R A Mメモリ・バンクがアクセス要求の処理に利用できるかどうかを判別し、指定されたR A Mメモリ・バンクが現在使用中であれ

ばアクセス要求をバッファに入れ、アイドル状態のときにアクセス要求をRAMメモリ・バンクに関連付けられた状態コントローラに適用する。RAMメモリ・バンクが使用中の場合、状態コントローラはアクセス・フロー・レギュレータに、RAMメモリ・バンクが現在他のアクセスの処理で使用中であり、当分の間さらにアクセス要求を処理するには使用できないことを示す信号を送る。

【0009】

書き込みアクセス要求の場合、アクセス・フロー・レギュレータは、アクセス要求をRAMメモリ・バンクに送ろうとしたときにすべての状態コントローラをスキャンする。その際に、その関連するRAMメモリ・バンクまたは格納のための空き領域のないRAMメモリ・バンクについて現在使用中信号を発生している状態コントローラを即座にバイパスする。アクセス・フロー・レギュレータは、使用中のまたは完全に使い果たされたRAMメモリ・バンクおよびその状態コントローラをハイパスし、格納領域用に利用可能なバッファがあるアイドルRAMメモリ・バンクにアクセス要求を送る。

10

【0010】

RAMメモリ・バンクのメモリは、高速だが比較的記憶容量の小さいタイプである。それぞれのRAMメモリ・バンクは、そのバンク宛のそれぞれのアクセス要求を高速に処理することができる。アクセス・サイクルが完了すると、その状態コントローラは、RAMの状態を示すアクセス・フロー・レギュレータへの使用中信号を除去する。使用中信号が除去されるとすぐに、アクセス・フロー・レギュレータは、RAMメモリ・バンクが新たなアクセス要求の処理に利用可能になったことを知る。

20

【0011】

本発明の他の可能な実施例は、RAMメモリ・バンクが使用中である短い時間間隔でのみ持続するRAMメモリ・バンク使用中信号の使用オペレーションである。この信号を送ることは、RAMメモリ・バンクを具現化するRAMデバイスの速度によってのみ制限される速度でアクセス要求を処理するのに都合がよく処理することができる争奪配置に含まれる。この争奪配置は、アクセスを処理可能な速度が不可欠な時間間隔を課すことで制限されるか、または実現される争奪ロジックの複雑さにより制限される。

【0012】

本発明を具現化する争奪ファシリティを使用すると、ダイナミックRAMメモリ・バンク・ファシリティが動作可能な最大速度は、争奪配置に内在する任意の制限ではなく使用されるRAMデバイスの速度によってのみ制限される。本発明を具現化する高速ダイナミックRAMメモリ・バンク・ファシリティは、パイプライン方式で動作させることで、光ファイバ伝送ファシリティのバス速度で到着するパケットを処理することができる。本発明により実現される争奪配置では、プロセッサにより制御されるリンク・リスト・エンジンが高トラフィック・レベルの処理時に最低の輻輳の着信および送信トラフィックを処理することができる速度が高まる。

30

【課題を解決するための手段】

【0013】

本発明の一態様は、リンク・リスト・データ・ファイルを処理するように適合されたメモリ管理システムを操作する方法であり、前記システムは複数の高速小容量メモリおよび低速大容量メモリを備え、前記高速メモリは第1のデータ転送速度を持ち、前記大容量メモリは前記第1のデータ転送速度よりも遅い第2のデータ転送速度を持ち、前記システムはさらに、前記メモリによりリンク・リストの読み書きの要求を生成するためのアクセス・フロー・レギュレータを備え、前記方法は、

40

【0014】

前記アクセス・フロー・レギュレータから前記高速メモリに書き込み要求を送信することにより、前記高速メモリへのリンク・リストの書き込みを開始する工程と、

前記リンク・リストの先頭バッファおよび末尾バッファおよび少なくとも1つの中間バッファを前記高速メモリに書き込む工程と、

前記高速メモリから前記大容量メモリに前記少なくとも1つの中間バッファを転送する

50

一方で、先頭バッファおよび末尾バッファを前記高速メモリ内に残す工程とを含む。

【0015】

本発明の他の態様は、リンク・リスト・データ・ファイルを処理するように適合されたメモリ管理システムを操作する方法であり、前記システムは複数の高速小容量メモリおよび低速大容量メモリを備え、前記高速メモリは第1のデータ転送速度を持ち、前記大容量メモリは前記第1のデータ転送速度よりも遅い第2のデータ転送速度を持ち、前記システムはさらに、前記メモリによりリンク・リストの読み書きの要求を生成するためのアクセス・フロー・レギュレータを備え、前記方法は、

【0016】

指定されたリンク・リストの読み取り要求を前記アクセス・フロー・レギュレータから前記指定されたリンク・リストのバッファを含む前記高速メモリに送信する工程と、

前記指定されたリンク・リストの先頭バッファを読み込む工程と、

前記リンク・リストの前記少なくとも1つの中間バッファを前記大容量メモリから前記高速メモリのうちの1つに転送する工程と、

前記1つの高速メモリに転送された前記中間バッファを前記指定されたリンク・リストの交換先頭バッファとして指定する工程と、

前記1つの高速メモリから前記指定されたリンク・リストの前記中間バッファを読み出す工程と、

前記指定されたリンク・リストの前記読み出されたバッファを前記アクセス・フロー・レギュレータに送信する工程を含む。

好ましくは、前記方法はさらに、

前記システムを前記アクセス・フロー・レギュレータからの複数の要求に対するリンク・リストを同時処理するように動作させる工程と、

前記システムを前記複数の高速メモリのうちの異なるメモリに格納されているリンク・リストのバッファを処理するように動作させる工程と、

前記システムを、末尾バッファを新しいリンク・リストに書き込まれた第1のバッファとして書き込むように動作させる工程と、

リンク・リストの先頭バッファをまず最初に読み込む工程を含む。

【0017】

好ましくは、前記システムはさらに、複数の状態コントローラを備え、それぞれの状態コントローラは前記複数の高速メモリのうちの対応するメモリに対する個別のコントローラであり、前記システムはさらに、前記アクセス・フロー・レギュレータを前記状態コントローラと接続する要求バスを備え、読み取り要求を送信する前記工程は、

前記アクセス・フロー・レギュレータを、前記読み取り要求を受信すべきアイドル状態の高速メモリを選択するように動作させる工程と、

前記読み取り要求を前記アクセス・フロー・レギュレータから前記要求バスを介して前記選択された高速メモリに対して個別の状態コントローラに送信する工程と、

前記状態コントローラを、前記選択された高速メモリの現在の占有レベルを判別するように動作させる工程と、

前記現在の占有レベルが所定のレベル以下の場合に前記要求を前記高速メモリに送信する工程と、

前記選択された高速メモリの前記現在の占有レベルが前記所定のレベルを超えた場合に前記大容量メモリへの接続を要求する工程を含み、

前記システムはさらに、バックグラウンド・アクセス・マルチプレクサ、および前記状態コントローラを前記マルチプレクサと接続するアクセス・バスを備え、前記システムはさらに、前記マルチプレクサを前記大容量メモリと接続するバスを備え、前記方法は、前記マルチプレクサを、

前記大容量メモリへの接続の要求を前記状態コントローラから受信し、

複数の要求状態コントローラのどれに対し前記大容量メモリへのアクセスを許可するかを決定するように動作させる工程と、

10

20

30

40

50

前記１つの要求状態コントローラを前記大容量メモリに接続する工程と、
前記１つの高速メモリから前記大容量メモリにデータを転送する際の前記大容量メモリのオペレーションを制御する工程と、
前記アクセス・バスを介して前記状態コントローラから前記マルチプレクサに前記リンク・リストのバッファを転送する工程とを含む。

【００１８】

好ましくは、前記大容量メモリから前記バッファを転送する前記工程は、
リンク・リストの中間バッファを前記大容量メモリから前記高速メモリに、前記高速メモリのデータ転送速度に実質的に等しいデータ転送速度を持つバースト・モードで転送する工程と、
前記読み出されたバッファを前記高速メモリに格納する工程と、
その後、前記アクセス・フロー・レギュレータに転送するため前記高速メモリから前記リンク・リストのバッファを読み出す工程と、
バッファを既存のリンク・リストに書き込むのに、前記高速メモリから前記大容量メモリに前記既存のリンク・リスト既存の末尾を転送することにより書き込む工程と、
新しいバッファを前記既存のリンク・リストの新しい末尾バッファとして前記高速メモリに書き込む工程とを含む。

【００１９】

好ましくは、前記状態コントローラを動作させる前記工程は、
複数のリンク・リストのバッファを同時に受け取る工程と、
前記アクセス・フロー・レギュレータに送られる複数のリンク・リストのバッファを分離する工程と、
前記アクセス・フロー・レギュレータによって受け取られた複数のアクセスを前記高速メモリに送る工程と、
前記アクセス・フロー・レギュレータからそれぞれの受信された要求に応答し、前記状態コントローラに個別の高速メモリの現在の占有レベルを判別する工程と、
前記現在の占有レベル以下の場合、前記関連する高速メモリに前記アクセスを送る工程と、
前記現在の占有レベルを超えた場合、前記要求をバッファリングするよう前記アクセス・フロー・レギュレータに信号を送る工程と、
前記高速メモリからバッファをバースト・モードで前記大容量メモリに転送するのを制御する工程と、
前記大容量メモリからバッファを前記高速メモリに転送するのを制御する工程と、
転送が要求されたときに前記大容量メモリがアイドル状態かどうか判別する工程と、
アイドル状態の場合前記バッファを前記大容量メモリに送る工程と、
前記大容量メモリが使用中の場合、前記転送をバッファリングする工程とを含む。
好ましくは、前記マルチプレクサを動作させる前記工程は、
複数のビディング高速メモリのどれに対し前記大容量メモリへのアクセスを許可するかを決定する工程と、
他のバッファの要求を前記１つの高速メモリにバッファリングする工程と、
前記大容量メモリからのバッファの転送先の高速メモリに対するＩＤを決定する工程と、

前記大容量メモリから前記識別された高速メモリにバースト・モードで前記バッファを転送するのを制御する工程とを含む。

【００２０】

好ましくは、前記方法は、さらに、
それぞれの前記高速メモリの使用中／アイドル状態を示す、それぞれの高速メモリに固有の信号を発生する工程と、
それぞれの発生した信号を前記アクセス・フロー・レギュレータに送る工程と、
前記アクセス・フロー・レギュレータを、前記高速メモリによるリンク・リストの読み

10

20

30

40

50

書きの要求を受け取るように動作させる工程と、

要求を受け取ったことに対する応答として前記アクセス・フロー・レギュレータを、前記高速メモリにより生成される前記使用中／アイドル状態信号を読み取るように動作させる工程と、

前記読み取りに対する応答として前記アクセス・フロー・レギュレータを、前記高速メモリのうちアイドル状態のメモリを識別するように動作させる工程と、

前記アクセス・フロー・レギュレータを、データ・ファイルの読み書きの要求を前記アイドル状態の１つの高速メモリに送るように動作させる工程とを含む。

【 0 0 2 1 】

本発明の他の態様は、リンク・リスト・データ・ファイルを処理するように適合されたメモリ管理システムを備え、前記システムは、

複数の高速小容量メモリおよび低速な大容量メモリであって、前記高速メモリは第１のデータ転送速度を持ち、前記大容量メモリは前記第１のデータ転送速度よりも遅い第２のデータ伝送速度を持つメモリと、

前記メモリによるリンク・リストの読み書きの要求を生成するためのアクセス・フロー・レギュレータと、

書き込み読みだし要求を前記複数の高速メモリのうちのアイドル状態のメモリに送ることにより前記メモリ内のリンク・リストの書き込みを開始するための装置と、

前記リンク・リストの先頭バッファおよび末尾バッファおよび少なくとも１つの中間バッファを前記高速メモリに書き込むための装置と、

前記高速メモリから前記大容量メモリに前記リンク・リストの前記少なくとも１つの中間バッファを転送し、その一方で前記リンク・リストの先頭バッファおよび末尾バッファを前記高速メモリに残すための装置と、

その後前記アクセス・フロー・レギュレータからの前記リンク・リストの読み込みの要求を前記高速メモリに送るための装置と、

前記リンク・リストの先頭バッファを前記高速メモリの１つに読み込むための装置と、

前記リンク・リストの前記少なくとも１つの中間バッファを前記大容量メモリから前記高速メモリに転送するための装置と、

前記高速メモリ内の転送されたバッファを新しい先頭バッファとして指定するための装置と、

その後前記先頭バッファおよび前記末尾バッファだけでなく前記中間バッファを前記高速メモリから読み出すための装置と、

前記リンク・リストの前記読み出されたバッファを前記アクセス・フロー・レギュレータに送信するための装置を備える。

【 0 0 2 2 】

好ましくは、前記メモリ管理システムはさらに、

それぞれの高速メモリの現在の使用中／アイドル状態を示す、それぞれの高速メモリに固有の信号を発生するための前記高速メモリを備える装置と、

前記信号を前記アクセス・フロー・レギュレータに送るための装置と、

前記高速メモリによるリンク・リストの読み書きの要求を受け取るための前記アクセス・フロー・レギュレータを備える装置と、

前記要求を受け取ったことに対する応答として、前記使用中／アイドル状態信号を読み込むため前記アクセス・フロー・レギュレータを備える装置と、

前記読み取りに対する応答として前記高速メモリのそれぞれの現在の使用中／アイドル状態を判別するための前記アクセス・フロー・レギュレータを備える装置と、

前記メモリの１つが現在アイドル状態かどうかを判別したことに対する応答として、前記高速メモリによるリンク・リストの読み書きの要求を許可する前記アクセス・フロー・レギュレータを備える装置とを備える。

本発明のこれらの態様および他の態様は、図面とともに詳細な説明を読むとよく理解できるであろう。

10

20

30

40

50

【発明を実施するための最良の形態】

【0023】

図1の説明

本発明は、マルチノード通信ネットワークのトラフィック・スループットを高めるための機能強化されたメモリ・インターフェースを備える。この機能強化されたメモリ・インターフェースは、図1に示されているような通信ネットワークへのパケット化された情報の放出を制御するトラフィック・シェーピング要素を具現化する。図1のネットワークは、互いに通信するノードと呼ばれるスイッチング要素を相互接続している。ノードは、A、B、C、D、E、F、およびGと指定されており、リンク1からリンク8で指定されている別々のリンクにより接続される。図1のノードでは、内向きポートから外向きポートにトラフィックを分散するネットワークを定義している。

10

【0024】

図2および3の説明

図2は、図1のノードを具現化する機器を開示している。図1は、ノードAとノードBを相互接続する単一経路（リンク1）を開示している。図1の各ノードは、着信リンクと送信リンクにより他のノードと接続されている。ノードAは、ノードAの着信リンクによりノードBから受信し、ノードAの送信リンクを使ってノードBに送信する。

【0025】

図2は、その着信リンクおよび送信リンクを含むそれぞれのノードの詳細を開示している。経路218は、ノードの送信リンクであり、経路223は、ノードの着信リンクである。図2は、ノードを定義する機器を示しており、図2の左半分では、ノードは着信リンク223、スプリッタ222、および各スプリッタと複数のポート201から205のそれぞれとを相互接続する経路221を含む。ポート201はその右側で外向きリンク218に接続され、他のノードに及ぶ。図2の左側に示されている機器は、外向きポート201から205を含んでおり、これらにより、ノードはリンク218に接続され、5つの外向きポート201から205を介して5つの異なるノードに及ぶようにできる。

20

【0026】

図2の右側は、外向きポート202を具現化する機器を詳しく例示している。ポート202は、複数のリンク・リスト・キュー215、各キューに個別の制御ロジック1800を備える。制御ロジックは、その後図18に詳しく例示されており、リンク・リスト情報を処理し、マルチプレクサ213に届けるのに必要な機器を備える。5つのキュー215は、経路231によりマルチプレクサ213に接続されている複数の制御ロジック要素1800のうちのそれぞれの要素により処理される。マルチプレクサ213は、経路231を、図1のリンク1から8のいずれかに対応する送信リンク218に接続する。図2のノードは、最大のパフォーマンスを発揮できるように出力キュー機能を備える。いかなる時点でも、外向きポート201～205のそれぞれについて、キュー215のうちの1つがマルチプレクサ213により選択され、キュー215からのパケットが送信リンク218に送られる。所定のポートに対するキュー1800の選択は、ネットワークで使用されるトラフィック・シェーピング・アルゴリズムによって異なる。そこで図1の複数のノードでは、ノードAなどの複数のノードのうちの1つにより処理される小さなコミュニティとの通信を望んでいると仮定する。オーバーロードは、ノードAへのリンクを介してノードAへの資源のトラフィック争奪として発生しうる。パケットはバッファリングされ、かわっているリンクが処理されると伝送を完了できる。

30

40

【0027】

この高トラフィック・シナリオは図3に詳しく示されており、ノードAは要求のあるノードAであり、暗色のリンク1、2、3、4、5、および6は要求のあるノードAへの可能な回線経路にサービスを提供するリンクを表す。

【0028】

ノードBは、自分自身、およびノードBおよびノードCをトラバースしノードAを取得するトラフィックを提供する3つのノード（C、E、およびF）のあわせて4つ

50

のノードのトラフィックを持つ。ノードBは、ノードAへの1つのリンク(リンク1)のみ持つ。ノードGおよびDは、完全リンク容量のトラフィックを供給することができる。したがって、リンク1はオーバーロードされる。ノードBは、高争奪時にノードAに向かうトラフィックをバッファリングしなければならない。したがって、トラフィックは、争奪の少ない状況に徐々に向かうにつれ送信されるようになる。

【0029】

高争奪間隔が十分長く持続すると、ノードBのバッファはオーバーフローする可能性がある。これを防止するために、トラフィック・シェーピング・アルゴリズムにより、リンクへのトラフィックの放出を調整することで、ネットワーク全体にわたってバッファが、要求のあるノードに向かう大きなトラフィックを排出し、それをいつでも吸収するようにできる。例えば、ノードEおよびFは、それぞれ、リンク5および6を満たすだけの十分なパケットがあるとしても、リンク5および6にトラフィックを放出しない。その後、もはや、それらのパケットを送信することはないが、そうする際に、ノードBはオーバーロードを回避し、ネットワークはより少ないパケットを破棄することになる。トラフィック・シェーピングは、前向きのフロー制御とみなすことができる。

【0030】

トラフィック・シェーピングは、正常動作させるためにはハイパフォーマンスのバッファリングを必要とする。このバッファリングは、ハードウェアのリンク・リスト・プロセッサに左右される。ハードウェアのリンク・リスト・プロセッサは、メモリを効率よく利用し、バッファを着信パケットに動的に割り当て、保持されているデータが正常に外向きリンクに転送された後バッファを回収する。

【0031】

図4～6の説明

リンク・リスト・バッファは、バッファリング情報に使用される本発明に基づく。初期化されたシステムにおいて、すべてのメモリは複数の汎用バッファに分けられる。各バッファには、そのバッファによって格納される内容401、および次のバッファへのポインタ402用の領域が確保される。これは、図4に示されている。

【0032】

初期化時に、すべてのバッファは、前のバッファのポインタ・フィールドを次のバッファのアドレスに設定することにより1つにつながれる。これは、自由リストと呼ばれ、図5に示されている。

【0033】

通信システムでは、これらの汎用バッファに1つずつ情報を書き込み、書き込まれたバッファを、何らかの特定の機能について情報を格納するキューにリンクする。システムの初期化が終わった後、すべてのキューは空である。それらのキューの長さは0であり、先頭と末尾はNULLを指す。図6の特定のキューについて情報が到着すると、汎用バッファが自由リストから取り出され、情報が埋め込まれ、キューのリストに追加される。末尾ポインタは、キューの追加された要素のアドレスに変更され、キューの長さカウンタは1つ増やされる。キューから情報が読み出されると、そのキューのリストの先頭にあるバッファの内容が読み出され、リストの先頭はリスト内の次のバッファに移される。キューの長さカウンタも1つ減らされる。

【0034】

図6のキューAでは、先頭バッファに内容Qがあり、末尾バッファに内容Zがあり、キューBでは、先頭バッファに内容Nがあり、末尾バッファに内容Gがあるが、キューCのバッファでは、先頭バッファも末尾バッファも内容HHが入っている。このシステムは、11個のバッファおよび3つのキューを持つ。リンク・リスト・バッファリング・システムの重要な特徴の1つに、バッファ割り当てが完全に動的であることが挙げられる。自由リストが空になっていない限り、バッファのキューへのどのような分配も許される。例えば、キューAは11個のバッファすべてを持ち、その後しばらくしてから、キューAは4つのバッファ、キューBは4つのバッファ、そしてキューCは3つのバッファを持つよう

10

20

30

40

50

にできる。その後しばらくして、すべてのキューは空になり、すべてのバッファが再び、自由リストに戻される。

【 0 0 3 5 】

リンク・リスト・バッファリングは、特に、通信アプリケーションに適しているのは、全能力で動作している所定のファシリティに対し、到着する情報が一定の速度でバッファを消費するからである。しかし、このファシリティは、多くの場合、複数の流れまたはストリームに関連する情報を多重化しなければならない。リンク・リストは、処理およびその後の外向き伝送媒体への多重化のためは内向き情報を逆多重化し、編成する効率のよい手段となっている。処理されたバッファは、送られた後、自由リストを通じてリサイクルすることができる。

10

【 0 0 3 6 】

図2のノードBへの内向きの合計トラフィックが毎秒10個のバッファを消費するが、ノードBはさらに毎秒10個のバッファを送信すると仮定する。このノードは、バッファが埋められるやいなや空にされるので、バランスがとれている。バッファリングの観点からは、どの流れが入って来るのか、出て行くのかは問題でない。例えば、ノードCからノードBへの内向きトラフィックの向かい先がノードGまたはノードAであるかは問題でない。ある時点において、リンク1をノードAに送るキューは、書き込み済みであり、リンクをノードGに送るキューは空でありえるが、その後、しばらくしてから、第1のグループは空になり、第2のグループは満杯になることがありえる。流れ全体がバッファ容量と一致する限り、システム完全性は維持される。しかしながら、トラフィック全体は固定された最大速度で媒体上を出入りするが、多くの流れはサポートされている通信媒体で多重化することができる。バッファの管理は柔軟性が高いため、一時的な機器状態が緩和される。

20

【 0 0 3 7 】

リンク・リスト・バッファは、トラフィック・シェーピングのシナリオでは有用である。ノードはそれぞれ、ある程度の全体的スループットに対応できる。トラフィック・シェーピングの実装により、バッファリングがノードに入り込むが、それは、シェーピング・プロファイルの条件を満たすため、いくつかの外向きリンクに向かう内向きトラフィックを一定時間バッファリングしてから送信しなければならないからである。コネクションレス型パケット指向ネットワークでは、内向きパケットの宛先は、そのパケットが到着してからでないといけない。特定のストリームまたは流れは全体的な最大速度を持ちえるが、所定の流れにより情報が爆発的にノードに送られる可能性がある。そのバースト送信時に、流れのバーストはリンク容量全体を消費する可能性があるため、他のノードを宛先とするリンクによりサポートされるパケットの流れは定義からアイドル状態である。

30

【 0 0 3 8 】

例として、図3のネットワークを再度考察する。ノードBは、3つのリンク1、4、および7に接続されている。リンク4上のトラフィックは、ノードAまたはノードGのいずれかを宛先とすることができる。ノードAを宛先とする一連のパケットは、ノードAからノードBへのストリームまたは流れをなす。ノードGを宛先とする一連のパケットは、ノードGからノードBへのストリームまたは流れをなす。

40

【 0 0 3 9 】

図7の説明

図7は、リンク4上の考えられる内向きパケット・シーケンスを説明するタイミング図を示している。

ノードAを宛先とする4つのパケットのバーストで、自由リストから切り離されたバッファはすべてリンク1をサポートするキューに追加される。このバーストが延長した場合、自由リストから取り出されたバッファはこのキューに流れ込む。しかし、いかなるときも、次の内向きパケットは、ノードAまたはノードGのいずれかを宛先とすることができる。次の内向きパケットの宛先は、確認されるまでわからない。したがって、効率のよいパケット交換機器のオペレーションには、バッファの自由度の高い動的な割り当てが必要

50

である。4つのパケットのバースト中に、ノードGへのリンク7をサポートするキューは、追加バッファを一切受け取らないが、バーストはノードAを宛先とするトラフィック専用であるため、それらのキューは、それ以上バッファを必要としないので内向きトラフィックを持たない。したがって、効率のよいパケット交換機器のオペレーションには、バッファの自由度の高い動的な割り当てが行えれば十分である。

【0040】

リンク・リストのオペレーション

リンク・リスト・エンジンは、半導体メモリをRAM記憶装置として採用している。バッファの追加または削除のためリンク・リストを操作する方法では、バッファ・メモリおよび関連するリンク・リスト・テーブル・メモリの両方に対し一連の読み書きを行う必要がある。リンク・リスト・テーブル・メモリは、リンク・リスト・プロセッサによりサポートされる各リンク・リストへの先頭および末尾のルックアップ・テーブルを含む静的構造である。例えば、流れに読み出したまたは書き込みトラフィックが含まれる場合、リンク・リスト・プロセッサは、最初に、流れの番号を使用して、注目しているリストの先頭および末尾のアドレスを検索する。注目しているリンク・リストのそれらのアドレスが判明すると、プロセッサはそのリストに対し所定のオペレーションを実行することができる。バッファがリンク・リストに追加される場合、空バッファが図6の自由リストの先頭から取り出され、自由リストの先頭がそのリスト内の次の空バッファとなるように書き直される。自由リストの新しい先頭のアドレスは、書き込みのためちょうど取り出されたばかりのバッファのリンクに含まれる。そのバッファの内容が書き込まれ、リンク・リストの末尾のバッファ内のリンク・フィールドは、書き込まれたばかりのバッファのアドレスを書き込まれる。その後、テーブル・メモリは、新しい末尾アドレスを書き込まれる。新しいバッファをリンク・リストに書き込むプロセスでは、テーブル・メモリは読み出しと書き込みを維持し、バッファ記憶装置では2つの書き込みを維持する。バッファをリンク・リストから読み出すプロセスでは、テーブル・メモリは読み出しと書き込みを維持し、バッファ記憶装置では読み出しおよび書き込みを維持する。バッファ記憶装置への書き込みは、空にされたバッファが自由リストに再度連結されなければならないと発生する。

【0041】

リンク・リスト・プロセッサはバッファ・メモリにランダムにアクセスする

このプロセスの重要な側面として、メモリへのアクセスがランダムであるという性質が挙げられる。ランダム化にはいくつかの要因がかかわっている。リンク・リストが通信ファシリティからのトラフィックをバッファリングする場合、一連のアクセスはそのファシリティによって伝送されるトラフィックに完全に左右される。Ethernet（登録商標）などのコネクションレス型ネットワークでは、ファシリティに到着するパケットは、その旅先で多数のキューを宛先としている可能性がある。将来のパケットに対する宛先キューは、一般的には予想できない。到着がランダムであるという性質から、所定のキュー内のアドレスはスクランブルされる。出て行くパケットの伝送は制御下にあるが、ここでもまた、ランダム化にネットワーク状態が絡んでいる。例えば、送信ファシリティで多数のキューを多重化するものとする。ときどき、それらのキューすべてが、一部にかかわることもあれば、1つにかかわることもあれば、何にもかかわらない場合もある。輻輳は、外向きファシリティの遠端フロー制御により、または外向きファシリティを宛先とするトラフィックが発生する多数の内向きファシリティにより、引き起こされることがある。

【0042】

ランダム化にかかわる第2の著しい要因は自由リストである。自由リストにかかわる要因は、送信ファシリティへのバッファ伝送の順序に完全に依存する。しかし、外向き伝送は予測不可能な条件に左右される。したがって、トラフィックの状態により、自由リスト上の空バッファのアドレス・シーケンスがランダム化される。

【0043】

リンク・リスト・プロセッサをバッファ管理の目的に使用する通常のシステムが1、2秒の間著しい負荷の下で動作すると、バッファ・メモリへのアクセスは、相関関係を完全

10

20

30

40

50

に欠いている。

【 0 0 4 4 】

バッファ・メモリのアクセス・パラメータおよびメカニクス

リンク・リスト処理ではメモリに対し、事前に作成されたスクリプトを使って一連のアクセスが実行されるため、リンク・リスト・プロセッサのパフォーマンスは、それらのアクセスの管理の仕方に大きく依存する。しかし、メモリ・コンポーネントの可用性およびアクセス配置などの最終的なバッファリング要件により、リンク・リスト・プロセッサに対するメモリ・システム設計が制約される。

【 0 0 4 5 】

リンク・リスト・プロセッサが中継交換アプリケーションで使用される場合、アタッチされたファシリティの長さおよびその容量がメモリ・システムの設計に課せられる。例えば、SONETネットワークではケーブル切断に関する標準報告間隔は60ミリ秒である。しかし、この標準は、一般にあるネットワーク・ケーブル配線路に対しては制約的であり、報告間隔は、さらに標準的には、200～300ミリ秒である。これは、数千キロメートルのケーブル配線長に対応する。最小サイズの packets をバースト伝送する場合、単一のOC-768ファイバは、2300万個を超える packets を300ミリ秒で送信する。光ファイバ・ケーブルのトラフィック伝達ファイバは、100本またはそれ以上の異なる個別の素線を持ちうる。したがって、そのような1つのケーブルを終端するシステムは、継ぎ目なくケーブル切断から回復できるように数十億個の packets の順序に基づいてバッファリングする必要がある。

【 0 0 4 6 】

知られているハードウェア・リンク・リスト・エンジンの基本的な問題

ハードウェア・リンク・リスト処理をサポートするメモリ・サブシステムは、大きく、高速であり、多数のキューをサポートし、クロックの任意のエッジに基づき任意のキューでメモリを操作し現在利用可能なトラフィック内での瞬間移動に対応できなければならない。ランダム・アクセス・メモリは、リンク・リスト・プロセッサのメモリを適当に操作できるが、十分高速にサイクル動作できないか、または十分に大きくないため最大容量ファシリティのバッファリングを行えない。

【 0 0 4 7 】

市販されている同期ダイナミック・メモリ(SDRAM)パッケージである2種類のメモリは、最大1メガビットの記憶域を持つが、そのランダムな読み書きサイクル時間は約60ナノ秒である。OC-768中継器にはメガビットのオーダーの記憶装置が必要であるため、SDRAMのサイズが好適であるが、最小サイズの packets をバースト送出する全負荷状態のOC-768中継器では、 packets は40ナノ秒おきに到着する。したがって、市販のSDRAMは遅すぎてOC-768ファシリティには使えない。

【 0 0 4 8 】

市販の同期スタティック・メモリ(SSRAM)は、約2ナノ秒でランダムな読み書きアクセス・サイクルを実行し、待ち時間は4ナノ秒である。これは、少数のOC-768ファシリティと制御オーバーヘッドを処理できる十分な速さである。しかし、SSRAMは16メガビットを超える容量のものは市販されていない。OC-768ファシリティのバッファリングを適切に行うためには約90個のSSRAMデバイスが必要であろう。そのような多数のSSRAMデバイスによって発生する熱量は問題になる。

【 0 0 4 9 】

結論として、ハードウェア・リンク・リスト・プロセッサにより形成される利用可能なメモリの基本的な問題は、大容量で低速(SDRAM)か、高速だが非常に小容量(SSRAM)のいずれかであるという点である。大容量でかつ高速なハードウェア・リンク・リスト・プロセッサが得られるような妥協点はない。

【 0 0 5 0 】

リンク・リスト・プロセッサの設計改善

本発明の改善されたリンク・リスト・プロセッサは、大容量高密度バッファだけでなく

10

20

30

40

50

、高速高性能バッファをどのようにしたら得られるかという問題に対する解を具現化する。これは、新規性のある方法で、メモリへのアクセスのストリームのランダム性に依存することにより実現される。

【 0 0 5 1 】

これまで、S D R A Mを使用するメモリ・サブシステムの争奪の問題により、システムは、利用可能なS D R A Mのパイプライン化された特徴を使用することによりアクセスをオーバーラップするのではなく、各連続するメモリ・サイクルが完了するのを待っていた。これにより、システムはかなり高速なバス・サイクル・レートで動作することが可能である。本発明では、R A Mメモリを、固有のランダム読み書きアクセス速度、つまり数メガヘルツではなく、そのポート速度、つまり、数百メガヘルツで動作させることにより、争奪問題を解消する。本発明の他の特徴は、バンク化R A Mメモリへのランダム・アクセス・ストリームが存在する場合に、バンクの個数が増えると、次のアクセスがすでに使用中のバンクに向かう確立が減じるので、より多くのメモリ・バンクが使用されるというものである。

【 0 0 5 2 】

図 8 および 9 の説明

ダイナミック R A Mは、複数のバンクとともに単一パッケージにパッケージングされている。バンクとは、別々にアドレス指定可能な 1 つの記憶ユニットのことである。バンクは物理パッケージ上のピンなどの入出力資源を共有するので、複数のバンクが同時に複数のアクセス要求を処理することができる。処理できる未解決要求の最大数は、クロック速度および同期 S D R A Mデバイスの構造によって決まる。近似的に、S D R A M内の 1 バンクへのアクセスが完了するのに 4 クロック・サイクルを要し、パッケージ内に 4 つまたはそれ以上のバンクがある場合、処理に 4 回のアクセスが同時に含まれる可能性がある。第 1 のアクセスを完了するのににかかわる 4 クロックのそれぞれの立ち上がりエッジで新しいアクセスが有効にされる。

【 0 0 5 3 】

図 8 は、円 8 0 3 の中に 4 つのバンク 8 1 0、8 1 1、8 1 2、および 8 1 3 を示している。これは、1 つの可能な S D R A Mの中の 4 つのバンクを表している。漏斗 8 0 2 および蛇口 8 0 4 は、S D R A Mの中にあるものにアクセスするための共有制御およびデータ・バス資源を表す。4 つのアクセス要求 A、B、C、D (8 0 1) が、S R A M 8 0 3 の漏斗 8 0 2 に入る状況を示している。

【 0 0 5 4 】

アクセス要求 A、B、C、または D のうち一度に 1 つだけが漏斗 8 0 2 に入ることができる。各バンク 8 1 0 ~ 8 1 3 は、アクセスに他と同じだけ時間を要する。バンクが要求 A を開始し、その後、他のバンクが要求 B を開始すると、要求 A に対する結果は、要求 B からの結果の前に蛇口 8 0 4 から現れる。しかし、しばらくの間、1 つのバンクが要求 A について動作し、他のバンクは同時に要求 B について動作する。

【 0 0 5 5 】

アクセスは、バンクの任意の組合せにより処理することができる。これらは、すべて、同じバンクにより処理するか、または異なるバンクにより処理することができる。いくつかのアクセスでは、そのアクセス自体のバンクとなるが、他のアクセスではバンクを共有する。アクセスのいくつかのグループでは、バンクがまったく使用されていない可能性がある。例えば、図 9 に示されているように、アクセス A および B は、同じバンク 8 1 0 により処理されるが、アクセス C は、バンク 8 1 2 により処理され、アクセス D は、バンク 8 1 3 により処理される場合がある。アクセスのグループを利用可能なバンクに分配するプロセスは、パーティション分割と呼ばれる。バンクの観点から、アクセスの回数は重要な情報であるが、それは、すべてのアクセスが一樣な特質を持つからである。したがって、パーティションは、バンクへのアクセスの会計といえる。例えば、{ 4 , 0 , 0 , 0 } は、4 つのアクセスが単一バンクを占有し、他の 3 つのバンクは占有されていないことを意味する。図 9 では、パーティションは { 2 , 1 , 1 , 0 } である。

【 0 0 5 6 】

図 1 0 ~ 1 3 の説明

図 1 0 の同期ダイナミック R A M (S D R A M) は、4 独立バンク・メモリ・アーキテクチャである。バンク毎に、オペレーションにはアクセス待ち時間、書き込みオペレーションの場合にピンからメモリ・アレイへの情報の送付、または読み出しオペレーションの場合にメモリ・アレイからピンへの情報の送付が伴う。また、メモリ内部のセンス・アンプが次の読み出しまたは書き込みサイクルを準備できるようにプリチャージ間隔も必要である。4 つのバンクのそれぞれが利用可能である。4 つのバンクはそれぞれ、専用のセンス・アンプを備えており、S D R A M の制御ポートおよびデータ・ポートについてのみアクセスの争奪が発生する。

10

【 0 0 5 7 】

図 1 0 および 1 1 は、制御バス 8 0 1 が S D R A M のバンク 1 および 2 に向けられた活動を受け入れていることを示している。アクセス・コマンド「A」と図 1 1 の関連する読み出しコマンド「R」との間に待ち時間がある。また、読み出しコマンドとデータが利用できるまでの間にも待ち時間がある。さらに、アクセス・コマンド A 1 と A 2 との間に待ち時間がある。図 1 1 では、1 0 ナノ秒のクロック・サイクルは、便宜上使用されている。全体的なサイクル時間は 8 0 ナノ秒である。S D R A M の 4 つすべての利用可能なバンクは、このときにアクセスできるが、どのバンクもこの時間間隔内では 1 回しかアクセスできない。このような待ち時間があるため、S D R A M は、パイプライン深さが 4 段であると言われる。衝突インスタンスの条件が満たされるためには、さらに 8 0 ナノ秒のメモリ・サイクルが必要である。衝突インスタンスは、使用中 S D R A M メモリ・バンクへのアクセス・フロー・レギュレータの到着として定義される。例えば、図 1 1 のクロック・サイクル 2 でバンク 1 に対するアクセス要求が到着した場合、クロック・サイクル 9 までメモリには適用できないが、それは、バンク 1 に対する要求 A がクロック・サイクル 1 で到着しているからである。

20

【 0 0 5 8 】

図 1 2 は、S D R A M 1 2 0 3 をバッファ記憶域として使用するハードウェア・リンク・リスト・プロセッサ 1 2 0 1 およびサポートされている流れ毎に先頭、末尾、およびカウンタを保持するテーブル記憶域用に小さな同期スタティック R A M 1 2 0 7 を示している。S D R A M 1 2 0 3 は、4 つの S D R A M 1 2 1 0、1 2 1 1、1 2 1 3、および 1 2 1 4 を備える。リンク・リスト・プロセッサ 1 2 0 1 では、クロック 1、3、5、1 1、および 1 5 について図 1 3 に示されているように 5 つのアクセス要求のストリームが発生すると仮定する。第 1 の 3 つのアクセス A 1、A 2、および A 3 は、異なるバンクに向かうのでパイプライン化することができるが、クロック 1 1 での第 4 のアクセス A 2 は、クロック 3 のアクセス A 2 を処理する使用中バンク 2 に進むので遅延させなければならない。

30

【 0 0 5 9 】

一般に、ランダムなアクセス・ストリームが存在する場合に発生する争奪の平均量は、4 バンク S D R A M に対する基準としてみなすことができる。バンクは入出力ファシリティのみを共有するので、4 バンク S D R A M への 1 ブロック分のアクセスに対し重み付きの平均アクセス時間を計算するほうが簡単である。例えば、パーティション { 3, 1, 0, 0 } を考察する。アクセスに (A, B, C, D) のラベルを付けた場合、可能なグループ分けとして、{ (A, B, C), (D) }, { (A, B, D), (C) }, { (A, C, D), (B) }, および { (B, C, D), (A) } がある。選択された 2 つのバンクは、{ 1, 2 }, { 1, 3 }, { 1, 4 }, { 2, 3 }, { 2, 4 }, または { 3, 4 } とすることができる。

40

【 0 0 6 0 】

3 つからなる 1 つのパーティションと 1 つからなる 1 つのパーティションを 2 つのバンクにマッピングする 2 つの方法がある。例えば、パーティション { (A, B, C), (D) } を { 1, 2 } にマッピングしようとした場合、バンク 1 を持つ (A, B, C) および

50

バンク 2 を持つ (D)、またはバンク 1 を持つ (D) およびバンク 1 を持つ (A, B, C) のいずれかがありえる。 $4 \times 6 \times 2 = 48$ であり、4 つのアクセスを配置する方法は 256 通りあるため、パーティション {3, 1, 0, 0} が生じる確率は、0.1875 である。漂遊アクセスは、2 つの衝突するアクセスの内側でパイプライン化することができ、したがって、完了するクロックの数は $8 \times 3 = 24$ となる。

【表 1】

パーティション	確率	完了するクロック	完了までの重み付き平均時間
4,0,0,0	0.0156	32	0.5
3,1,0,0	0.1875	24	4.5
2,2,0,0	0.1406	19	2.6714
2,1,1,0	0.5625	16	9.0
1,1,1,1	0.0938	14	1.3125

表 1: 4 つのトランザクションの 4 メモリへの配置

【0061】

表 1 は、それぞれの可能な分配について示されている関連するデータとともにアクセスを 4 メモリ・バンクに分配する様々な可能な方法に関係するデータを例示している。表 1 の第 1 欄では、様々なバンクへのアクセスが可能な様々な可能なパーティションの一覧を示している。表 1 の最上行は、第 1 のバンクが 4 つのアクセスを受け取り、残りバンクは何も受け取らないことを示している。第 1 欄 2 行目の分配は、3, 1, 0, 0 である。第 1 欄の下の方の分配は、4 つのアクセスすべてが平等に各バンクに分配されていることを示す。第 2 欄は、その行の各分配の確率を示す。第 3 欄は、その機能を完了するのに要するクロック・サイクル数を示している。第 4 欄は、完了までの重み付き平均時間を示す。第 1 行については、第 1 のバンクへの 4 つすべてのアクセスの分配の確率は 0.0156 である。これは、完了までの重み付き平均時間が 0.5 であれば完了までに 32 クロック・サイクルを要する。下の行は、1, 1, 1, 1 の分配、0.0938 の確率を持ち、重み付き平均時間 1.3125 で完了までに 14 クロック・サイクルを要する。この評価から、最も確率の高い分配は 2, 1, 1, 0 であり、確率は 0.5625 であることがわかる。最もよい可能な分配は、下の行に示されており、確率は 0.0938 である。

【0062】

表 1 の 4 つのバンクについて、4 つのアクセスを完了するまでの合計重み付き平均時間は 17.98 クロックである。この数字は、4 バンク・メモリへのアクセスの可能なすべての組合せについて完了までの重み付き時間の合計である。2 つの 10 ナノ秒システム・クロックのパイプライン・クロックおよび 16 ビット幅バスを仮定すると、4 バンク SDRAM に対する平均持続可能スループットは、1.24 Gビット/秒であるが、それは各メモリのトランザクションが、それぞれにつき 2 クロックに対し 16 ビットを要するからである。争奪オーバーヘッドは 28.5% である。この値は、アクセス完了までの平均時間とアクセス完了までの可能な最短時間との差をアクセス完了までの可能な最短時間で割るという計算で求められる。例えば、アクセス完了までの平均時間が 50 ナノ秒とし、アクセス完了までの可能な最短時間が 40 ナノ秒であったと仮定する。すると、オーバーヘッドは 25% になる。

図 13 および 14 は、争奪状態でないアクセスと争奪状態であるアクセスとの差を例示している。

【0063】

図 1 4 の説明

頻繁な争奪は設計の決定に影響を及ぼす。図 1 4 のタイミング図を考察する。ここで、メモリ・バンクへのアクセス・バスはスロット付きである。それぞれのアクセスには、読み出し間隔と書き込み間隔が設定されている。これらのアクセスは、メモリの制御バスのアイドル時間を最小限に抑えるためできる限り緊密にパックする必要がある。

【 0 0 6 4 】

このようにしてアクセスを構造化することは、争奪でメモリの平均読み書きサイクル時間が実質的に増大するため意味のあることである。平均読み書きサイクル時間は最大読み書きサイクル時間程度なので、メモリの読み書きサイクル時間が最大となるようにシステムを設計するとより効率的である。少ない読み書きサイクルを使用して効率および動作速度を高めることは、争奪を管理するために必要なハードウェアのコストに見合わない。

10

【 0 0 6 5 】

メモリは、バス速度で定義された速度でアクセス要求を受け付けることができるが、それは内部パイプライン設計によりそうするようにできるからである。しかし、争奪率 2 8 . 5 % では、争奪の少ないシステムとは反対に、争奪率がキューの深さに指数関数的に関係しているため、このシステムにおけるキュー操作が劇的に増える。様々なアクセス時間をサポートするこのキュー操作および状態制御装置は、一定の最大アクセス時間を予想するものと比べてかなり念入りである。例えば、一定の最大アクセス時間を持つ設計では、キューを必要とせず、全体にわたる状態制御はかなり単純である。アクセス時間が頻繁に変わる設計では、メモリを使用しているマシン内に複数のキュー、さらに争奪が発生したときにアプリケーションからこのマシンを起動、停止できるより精巧な状態ロジックを必要とする。

20

【 0 0 6 6 】

図 1 5 および 1 6 の説明

図 1 5 に示されているように、8 バンク分の S D R A M があると仮定する。この場合、メモリは、図 1 3 に示されている S D R A M のサイズの半分のブロックにパーティション分割される。例えば、図 1 3 のバンク 1 を宛先とするアクセスは、今度は、図 1 6 のバンク 1 またはバンク 2 のいずれかに入る。図 1 3 のバンク 2 を宛先とするアクセスは、図 1 6 のバンク 3 とバンク 4 との間に分割される、などとなる。図 1 3 では、バンク 2 への第 2 のアクセスは、バンク 2 への第 1 のアクセスとの争奪状態を引き起こし、遅延が発生する。図 1 6 では、これらのアクセスのうちの第 1 のものがバンク 3 に向かうが、それらのアクセスのうち第 2 のものはバンク 4 に向かう。

30

【 0 0 6 7 】

図 1 3 および 1 6 を比較すると、争奪の解消により、図 1 3 のクロック 1 3 (A 4) とクロック 1 6 (R 4) との間のアイドル・ギャップが閉じられることがわかる。図 1 6 では、データ・バスは、クロック 7 からクロック 1 6 へ連続的に占有される。争奪は、アクセスのランダムな内向きストリームが与えられると、異なる組合せで発生することがある。争奪は、8 つのバンクでは解消されないが、バンクの数が多くなるほど争奪の発生は減少する。

【表 2】

バンク	完了までの重み付き 平均時間	争奪オーバーヘッド	最大持続可能 スループット(Gビット/秒)
4	17.98	28.5%	1.24
5	17.06	21.9%	1.31
6	16.47	17.6%	1.36
7	16.06	14.7%	1.39
8	15.77	12.6%	1.42
9	15.55	11.1%	1.44
10	15.38	9.9%	1.45
11	15.24	8.9%	1.47
12	15.12	8.0%	1.48
13	15.03	7.4%	1.49
14	14.95	6.8%	1.5
15	14.88	6.3%	1.5
16	14.82	5.9%	1.51

表2: バンクの個数に対する4回アクセス完了までの重み付き平均時間

【0068】

表2は、バンクの数を増やした場合に完了までの重み付き平均時間がどのように減少するかを示している。最大持続可能スループットでは再び、10ナノ秒のシステム・クロックと16ビット幅のデータ・バスを仮定している。

【0069】

パイプライン深さも、パフォーマンスに影響を与える。例えば、8バンク分のメモリがあるが、パイプラインの中には2段しかない場合、2回アクセス完了までの重み付き平均時間は6.25クロックである。この場合のオーバーヘッドは4.2%であり、最大持続可能スループットは1.51Gビット/秒である。争奪率は5%に低下するものとする。その後、平均アクセス時間は3.5ナノ秒である。これは、かなりバス・サイクル時間に近い。この結果は、図14および図16に示されているアクセス制御配置を比較してみるとわかる。図14では制御バス活動の密度が高い、つまり、図14の場合よりも1クロック当たりのランダム読み書きアクセスが多いということがわかる。

【0070】

図 17 および 18 の説明

本発明では、独立連携状態コントローラ 1804 は RAM パッケージのそれぞれのバンク 1803 に割り当てられている。それぞれのバンクは、独立にサイクル動作し、その状態コントローラ 1804 を通じて、他の RAM バンク 1803 にあわせよく調整された形で結果を出すことができる。さらに、状態コントローラ 1804 の流れにより、ときおり発生する争奪のためにアクセス要求を保持するアクセス・フロー・レギュレータ 1801 のキューを制御する。これにより、アクセス要求の破棄が防止される。状態コントローラ 1804 は、さらに、バンク 1803 のリフレッシュなどのハウスキーピング機能を独立に管理する。状態コントローラ 1804 は、独立している。状態コントローラ 1804 を使用すると、RAM バンク 1803 との間のバックグラウンド・バースト転送を、他の RAM バンク 1803 のフォアグラウンドのアクセス活動と同時に行うようにできる。このため、RAM バンク 1803 からリンク・リストの中間部をリモート RAM 1806 に格納することができる。これにより、リンク・リストの先頭および末尾のみが RAM バンク 1803 に残る。例えば、図 6 のキュー 506 を見ると、バッファ Q および Z は、RAM バンク 1803 内のどこかに置かれるが、バッファ D および R は、リモート RAM 1806 内に格納される。リンク・リストの中間をリモートで格納できれば、開示されているシステムでは、市販のパッケージ RAM を使用して任意のサイズのリストをサポートすることができる。リンク・リストの大部分をリモート RAM 1806 にリモートで格納することができる場合、FPGA に埋め込まれている RAM 1803 を先頭と末尾に使用することができる。状態コントローラ 1804 は、先頭と末尾を保持する RAM 1803 と組み合わせる。この設計は、RAM 1803 と異なるパッケージに常駐する状態コントローラよりも効率がよい。RAM 1803 と状態コントローラとが同じ場所にある場合は、リストの先頭と末尾を格納するため技術を選択することになる。この選択結果は、少数のリンク・リスト・キュー用のオンボード・レジスタ、適度の量のリンク・リスト・キューのスタティック RAM 1803、または大量のリンク・リスト・キューのダイナミック RAM 1806 である。

【0071】

状態コントローラ 1804 のブロック図が図 17 に示されている。

状態コントローラ 1804 は、フォアグラウンド・ポートから流れ込む情報をゲートするアービトレーションおよびシーケンシング・ロジックが適用され、フォアグラウンドまたはバックグラウンド転送で使用中の場合に新しい着信活動を行わないように RAM メモリ・バンク 1803 をガードする。さらに、状態コントローラは、RAM メモリ・バンク 1803 内の状態を監視し、リモート RAM 1806 との相互作用が発生する場合にそれを判別する。状態コントローラ 1804 は、図 18 に示されているように、バンクグラウンド・アクセス・マルチプレクサ 1808、リモート RAM 1806、およびアクセス・フロー・レギュレータとともにシステム内に組み込まれる。

【0072】

図 17 の状態コントローラ 1804 は、関連する RAM メモリ・バンク 1803、アクセス・フロー・レギュレータ 1801、およびバックグラウンド・アクセス・マルチプレクサ 1808 の間のインターフェースとして機能する。これらの要素への接続は、図 17 に詳しく示されている。状態コントローラ 1804 は、マルチプレクサ 1702 ならびにアービトレーションおよびシーケンシング・ロジック要素 703 を備える。低位側では、マルチプレクサ 1702 は、経路 1710 および 1711 に接続され、図 18 のバス 1809 - 1 ~ 1809 - 9 の一部となる。この経路の上で、マルチプレクサは読み出しおよび書き込みオペレーションで経路 1710 を介して関連する RAM メモリ・バンクとデータを交換する。経路 1711 は、マルチプレクサ 1702 を介して状態コントローラ 1804 が関連する RAM メモリ・バンク 1803 のオペレーションを制御することを可能にする双方向制御経路である。RAM データ経路 1710 は、マルチプレクサを介してデータ経路 1704 またはバックグラウンド・アクセス・マルチプレクサ 1808 に届くデータ経路、またはデータ経路 1705 およびバス 1802 を介してアクセス・フロー・レギ

ュレータ 1801 に接続することができる。これらのデータ経路は、読み出しオペレーションと書き込みオペレーションの両方で使用することができる。

【0073】

マルチプレクサ 1702 の底部の RAM 制御経路 1711 は、経路 1712 ならびにアービトレーションおよび制御ロジック要素 1703 を介して経路 1707 および 1706 に接続されている。マルチプレクサの経路 1711 は、一度に経路 1707 および 1706 の一方のみに接続可能である。経路 1706 に接続された場合、さらに経路 1810 および、読み出しと書き込み両方のオペレーションでバックグラウンド・アクセス・マルチプレクサ 1808 および関連するリモート RAM 1806 のオペレーションを制御する。経路 1711 が要素 1703 を介して経路 1707 に接続された場合、さらに、バス 1802 を介してアクセス・フロー・レギュレータ 1801 に及ぶ。アービトレーションおよびシーケンシング・ロジック要素 1703 は、読み出しおよび書き込みの両方のオペレーションで状態コントローラ 1804 とのデータ交換後、アクセス・フロー・レギュレータ 1801 を制御するために必要な情報処理機能およびロジック機能を備える。アービトレーションおよびシーケンシング・ロジック 1703 はさらに、バス 1706 および 1810 を介して、バックグラウンド・アクセス・マルチプレクサ 1808 と通信し、リモート RAM 1806 が RAM メモリ・バンク 1803 からデータを受け取ったときのオペレーションだけでなくリモート RAM 1806 がデータを状態コントローラ 1804 に送信し、状態コントローラに関連付けられた RAM メモリ・バンクに入力するオペレーションを制御する。

【0074】

状態コントローラ 1804 は、関連する RAM メモリ・バンク 1803 と、およびアクセス・フロー・レギュレータ 1801 と、およびバックグラウンド・アクセス・マルチプレクサ 1808 を介してリモート RAM 1806 と制御およびデータを交換する 4 つの高水準機能を備える。これら 4 つの高水準機能について次に説明する。

【0075】

図 17 の状態コントローラ 1804 によって実行される第 1 の機能は、読み出しまたは書き込み要求が発生したときにアクセス・フロー・レギュレータ 1801 からの情報の転送に関連する内向きアクセス要求シーケンスを開始し、制御し、そうする際に、関連する RAM メモリ・バンク 1803 を制御し、書き込み要求ではデータを RAM メモリ・バンク 1803 に書き込み、アクセス・フロー・レギュレータ 1801 からの読み出し要求では RAM メモリ・バンク 1803 からデータを読み出すことである。

【0076】

図 17 の状態コントローラにより実行される第 2 の機能は、関連する RAM メモリ・バンク 1803 内のバッファ満杯レベルを検出するトリガ信号に応答することである。このトリガは、関連する RAM メモリ・バンク内のバッファが十分に消費されたか枯渇したことを示す。関連する RAM メモリ・バンク内のバッファが十分に消費されている場合、リモート RAM 1806 への書き込みがトリガされる。関連する RAM メモリ・バンク内のバッファが十分に枯渇している場合、リモート RAM 1806 からの読み出しがトリガされる。

【0077】

状態コントローラ 1804 により実行される第 3 の機能は、関連する RAM メモリ・バンク 1803 からリモート RAM 1806 への転送を開始し、管理し、さらに、リモート RAM 1806 から RAM メモリ・バンク 1803 に戻る逆の方向のデータ転送も管理することである。

【0078】

状態コントローラ 1804 により実行される第 4 の機能は、マルチプレクサ 1702 からの信号を待ち、マルチプレクサ 1702 からその信号を受信した後リモート RAM 1806 との間の転送を開始することである。

【0079】

10

20

30

40

50

マルチプレクサ 1702 により実行される他の機能は、複数の RAM メモリ・バンク 1803 が同時にリモート RAM 1806 へのアクセス権を要求している場合にリモート RAM 1806 へのアクセス権をどのビッディング RAM メモリ・バンク 1803 に与えるかを選択することである。マルチプレクサ 1702 により実行される他の機能は、リモート RAM 1806 とビッディング RAM メモリ・バンク 1803 との間のオペレーションに関連する転送およびスケジューリング機能を、前記転送の間に依存関係がある場合に、開始することである。このような依存関係は、メモリ・システムとの間のストリーミング化されたアクセスから生じることがある。

【0080】

マルチプレクサ 1702 により実行されるさらに他の機能は、RAM メモリ・バンク 1803 を制御し、リモート RAM 1806 からの書き込み入力の方角を決めることである。マルチプレクサ 1702 は、リモート RAM 1806 へのアクセス権を与え、リモート RAM 1806 と RAM メモリ・バンク 1803 との間の情報の経路を設定する。

【0081】

表 2 は、RAM バンク 1803 争奪があると、従来の SDRAM に基づくシステムのパフォーマンスが制限される場合があることを示している。図 14 の説明から、この制限は厳しいため、システムは RAM メモリ・バンク争奪を回避して設計されていることがわかる。図 14 の設計では、RAM バンクの利用可能性を予想することに基づく争奪とは反対に RAM メモリの全サイクル時間の時間間隔についてアクセスを構造化することによりすべてのサイクルでの争奪を回避している。つまり、RAM メモリ・バンクは設計上、最適な速度よりも遅い速度で動作するということである。図 12 および 13 は、バンクの利用可能性を予想するシステムを説明しているが、争奪が発生した場合のために実質的によいロジックを必要とする。この実装の問題は、争奪によりパフォーマンスが奪われ、必要なハードウェアを追加したとしても十分なパフォーマンス向上がないという点である。

【0082】

図 15 および 16 は、RAM メモリ・バンクの数を増やして争奪回数を減らすことによるオペレーションの向上を示している。これには、ハードウェアのロジックを増やす必要がある。パフォーマンスの向上とハードウェアの増大との関係は、ハードウェアを受け入れられるだけ追加するとパフォーマンスの違いが目立ってわかるというような関係である。多数の RAM メモリ・バンクへのアクセスを調整するには、図 17 に示されているような専用の状態ロジックが必要であり、ハードウェア資源が消費される。これらの資源は、市販の FPG A の中から見つかる。パフォーマンスを最大にするためには、マスカブル・データ・バスが好ましい。バースト・データを滑らかに吸収するために広くなければならず、また最小クオンタムのデータを格納しやすくするためにマスク可能でなければならない。このようにして構成可能なメモリは、市販の FPG A で利用可能である。しかし、このメモリは、前述のように、300 ミリ秒分の OC - 768 オプション・ファイバのバッファリングなど、実質的なバッファリング・ジョブには不十分な小さな量でしか利用できない。さらに、集積回路内では利用可能な領域の量は限られている。RAM メモリに使用される領域は、状態コントローラ・ロジックには使用できず、その逆もそうである。しかし、RAM バンクを増やすと、パフォーマンスが向上し、各バンクは専用の状態コントローラを持たなければならない。したがって、バッファリング要件を満たすこととシステムのパフォーマンスとの間には食い違いがある。

【0083】

このコンフリクトに対する解決策は、FPG A RAM バンク 1803 に載せるメモリの量を制限することである。リンク・リストの先頭および末尾のみがアクセスされ、リンク・リストの中間の要素は、リストの先頭に移動するまで、常にアイドル状態である。したがって、リンク・リストの中間要素をボードから移動して離し、FPG A RAM バンク 1803 からリモート RAM 1806 に移すと、パフォーマンスを向上させることができる。

【0084】

10

20

30

40

50

この妥協策を実装するシステムは、図 17 および 18 に示されている。

この解決策は、経済的な理由からも賢明な方法である。RAM メモリの構成可能なバンクは、1 ビット当たりの価格が市販の SDRAM に比べてかなり高い。リンク・リストの中間要素を保持するリモート RAM 1806 がアクセスされていなければ、そのメモリに関連するサイクル時間ペナルティはない。最終的には、リモート RAM 1806 にアクセスして、リストの中間への要素の格納および取得を行わなければならない。しかし、FPGA 上のメモリの多数のバンクには上述の構成可能な性質があるため、SDRAM RAM バンク 1803 のバースト・モードと互換性のある設計が可能である。メモリの多数のバンクからなるシステムのスループットを SDRAM に一致させ、バランスのとれた設計とすることも可能である。これにより、SDRAM 1803 をバースト・モードで動作させ、データ・バス上でクロックと同期がとられているデータとともにサイクル時間をパイプライン化して、SDRAM サイクル時間のコストを最小限に抑えることができる。これで、同様に、メモリはメモリが当初設計された動作モードに戻る。

【0085】

リンク・リスト・バッファは、バックグラウンド・アクセス・マルチプレクサ 1808 を介してリモート RAM 1806 に置いたり、取り出したりされる。このため、フォアグラウンド・アクセスを、この追加トラフィックに関係なく、進められる。これは、表 2 に例示されている確率論的モデルが図 18 に示されているように利用可能なフォアグラウンド・バスに依存しているため重要である。フォアグラウンド・バス 1810 で進行中のリモート RAM 1806 へのバックグラウンド転送をブロックすると、表 2 を得るために使用されるモデルはかなり込み入ったものとなる。したがって、パフォーマンスが低下する。バックグラウンド・アクセス・バス 1810 が示されている。

【0086】

図 18 は、本発明を具現化するリンク・リスト・エンジン 1800 を開示している。リンク・リスト・エンジン 1800 は、ポート 1817 および 1818 に接続されている着信経路 1812 および送信経路 1813 を持つ通信システム 1811 に接続されていることが示されている。このシステムは、プロセッサ 1814 を備え、さらに、経路 1819 を介してアクセス・フロー・レギュレータ 1801 にまで及ぶ経路 1815 および 1816 も含む。動作すると、システム 1811 はリンク・リスト・エンジン 1800 のメモリを使用して読み出しおよび書き込みオペレーションを実行し、動作中にポート 1817 および 1818 で必要とするデータを格納する。

【0087】

アクセス・フロー・レギュレータ 1801 およびバス 1802 は、複数の状態コントローラ 1804 に接続されるが、それぞれ複数の RAM メモリ・バンク 1803 のうちの 1 つと関連付けられている。アクセス・フロー・レギュレータ 1801 は、RAM メモリ・バンク 1803 の 1 つに情報を格納するよう要求する書き込み要求をシステム 1811 から受け取る。アクセス・フロー・レギュレータ 1801 は、それらのアクセス要求を受け取って格納し、データを関連する RAM メモリ・バンク 1803 に入力するため様々な状態コントローラ 1804 に選択的に分配する。このプロセスは、アクセス・フロー・レギュレータ 1801 がシステム 1811 から読み出し要求を受け取ったときにメモリ読み出しオペレーションで逆の方向に動作し、このプロセスにより、状態コントローラ 1804 を介して、要求されたデータを含む RAM メモリ・バンク 1803 が読み出され、状態コントローラ 1804 およびバス 1802 を介して、システム 1811 にデータを送信するアクセス・フロー・レギュレータに適用される。

【0088】

RAM メモリ・バンク・デバイス 1803 は、RAM メモリ・バンク 1803 から情報を受け取り、その情報が RAM メモリ・バンクですぐに必要ない場合に格納しておくことができるリモート RAM 1806 により増強された比較的小さなメモリ記憶容量を持つ高速な要素である。リモート RAM 1806 は、動作中、バックグラウンド・アクセス・マルチプレクサ 1808 およびバス経路 1810 により補助され、マルチプレクサおよび

バス経路はそれぞれ一意的な状態コントローラ 1804 および関連する RAM メモリ・バンク 1803 にまで及んでいる。この手段により、満杯または空になりつつある RAM メモリ・バンク 1803 は、このことを関連する状態コントローラ 1804 に通知し、このコントローラはさらに、バス経路 810 を介してバックグラウンド・アクセス・マルチプレクサ 1808 と通信することができる。マルチプレクサ 1808 は、状態コントローラ 1804 が RAM メモリ・バンク 1803 から情報を読み出して、リモート RAM 1806 に転送し、その情報が出所の RAM メモリ・バンク 1803 によって再び必要になるまで一時的記憶域に置く動作を補助する。このときに、RAM メモリ・バンク 1803 は、状態コントローラ 1804 に、リモート RAM 1806 が RAM メモリ・バンクによって必要になりそうな情報を含んでいることを通知する。バックグラウンド・アクセス・マルチプレクサ 1808 および状態コントローラ 1804 を連携動作し、それにより、リモート RAM 1806 の適切な部分は情報を読み出して出所の RAM バンク 1803 に送り返す。リモート RAM 1806 は、比較的低速の大容量メモリ要素であり、RAM メモリ・バンク 1803 からオーバーフローする情報を効率よく格納するか、または情報をアンダーフローになった RAM メモリ・バンク 1803 に供給することができる。

【0089】

本発明の一態様は、図 18 の方法および装置を含み、そこでは、RAM メモリ・バンク 1803 への書き込みオペレーションは状態コントローラ 1804 を介してリンク・リスト情報を RAM メモリ・バンク 1803 に書き込む工程と、満杯状態に近づくまで RAM メモリ・バンク 1803 への書き込みオペレーションを継続する工程と、状態コントローラ 1804 およびバックグラウンド・アクセス・マルチプレクサ 1808 を介して RAM メモリ・バンク 1803 から新規に受信された情報の一部を同時にリモート RAM 1806 に読み出しながらアクセス・フロー・レギュレータ 1801 から RAM メモリ・バンクに追加情報を書き込むことを継続する工程とにより実行される。この情報は、その後出所の RAM メモリ・バンク 1803 により必要になるまでリモート RAM 1806 内に残される。

【0090】

図 18 のシステムは、RAM メモリ・バンク 1803 からのデータの読み出しを実行するが、そのために、要求されたデータを含む RAM メモリ・バンク 1803 と関連付けられた状態コントローラ 1804 に信号を送る工程と、状態コントローラ 1804 およびバス 1802 を介して選択されたデータの読み出し結果をアクセス・フロー・レギュレータ 1801 に戻すことを開始する工程と、選択された RAM メモリ・バンク 1803 の読み出しを継続し、それと同時に、それらの読み出しで注目するデータの選択された RAM メモリ・バンク 1803 が尽きてしまったこと、およびその後の読み出しオペレーションで必要となる注目するデータの一部が現在、リモート RAM 1806 に格納されていることを判別する工程とを使用する。リモート RAM 1806 からプリフェッチ読み出しを開始し、情報を要求される前に出所の RAM メモリ・バンク 1803 に転送して戻す工程と、選択された RAM メモリ・バンク 1803 の読み出しを継続し、必要ならば、リモート RAM 1806 から読み出し中の RAM メモリ・バンク 1803 へのデータの読出しを継続する工程とを使用する。このオペレーションは、読み出しオペレーションについてアクセス・フロー・レギュレータ 1801 により要求された情報全部が履行されるまで継続される。

【0091】

RAM メモリ・バンク 1803 の高速メモリ・デバイスおよびリモート RAM 1806 の低速大容量メモリは連携動作し、高速 RAM メモリ・バンク 1803 の容量を超えるデータを格納する。リモート RAM 1806 は、書き込み中の RAM メモリ・バンクからの書き込みオペレーションでこのデータを受け取るが、その一方で、RAM メモリ・バンク 1803 はアクセス・フロー・レギュレータ 1801 から高速でさらにデータを受信し続けることができる。このプロセスは、RAM メモリ・バンク 1803 が最初に高速で読み出されると逆方向に読み出しオペレーションで動作し、RAM メモリ・バンクで必要

とされリモートRAM 1806に格納されるデータは、RAMメモリ・バンク1803に注目する情報が枯渇すると、プリフェッチ方式でRAMメモリ・バンク1803に転送して戻される。読出しオペレーションは継続し、リモートRAM1806内の関連情報すべてを高速RAMメモリ・バンク1803に転送して戻すことができ、アクセス・フロー・レギュレータ1801により要求された情報全部が状態コントローラ1804およびバス1802を介して高速なRAMメモリ・バンク1803によりアクセス・フロー・レギュレータ1801に送り返されるまで高速データ転送速度で読み出される続ける。または、状態コントローラ1804およびRAMメモリ・バンク1803は、バックグラウンドでの取得をトリガした前回の読み出しオペレーションと別の書き込みオペレーションを続行することができる。プリフェッチはトリガにより自動実行されるため、特定の状態コントローラ1804およびRAMメモリ・バンク1803がバックグラウンド転送で占有されている間、他の状態コントローラ1804およびRAMメモリ・バンク1803は自由に、このバックグラウンド・オペレーションとは無関係にオペレーションを実行できる。

【0092】

本発明の他の態様では、状態コントローラ1804を使用して経路1820を介して電位をアクセス・フロー・レギュレータ1801に印加し、それぞれの状態コントローラ1804に関連付けられているRAMメモリ・バンク1803が既存の読み出しまたは書き込みオペレーションで現在使用中かどうか、または新しい読み出しまたは書き込みオペレーションの受信された要求に利用可能かどうかを示す。RAMメモリ・バンク・デバイス1803のメモリ要素は光バスの高速データ転送速度で動作するので、RAMメモリ・バンク1803は、光ファイバ・バスに適した速度で読み出しまたは書き込みオペレーションを実行することができる。その結果、経路1820を経由して状態コントローラにより使用中信号がアクセス・フロー・レギュレータ1801に印加され、関連するRAMメモリ・バンクが利用可能である、または利用可能でないことを示す。この信号は、読み出しまたは書き込みオペレーションを実行するためにRAMメモリ・バンク1803デバイスに必要な数ナノ秒間しか持続しない。したがって、状態コントローラ1804によりこの使用中/アイドル電位が経路1820に印加される場合、争奪ファシリティによりアクセス・フロー・レギュレータ1801およびその要素1821はRAMメモリ・バンク1803の使用/アイドル状態を監視することができる。このため、読み書きでの任意の所定の最小時間間隔をRAMメモリ・バンク1803に課す従来技術の複雑なロジック回路要素または争奪配置による争奪の遅延は解消される。この手段により、アクセス・フロー・レギュレータ1801、経路1820、および状態コントローラ1804は、光ファイバ・バスのナノ秒転送速度で動作する効率的で高速な争奪ファシリティを実現する。この改善された高速な争奪配置では、ポート1817および1818に関連付けられた着信および送信キューはRAMメモリ・バンク1803の高速要素により処理されるためより高速にデータを交換し処理することができるので、システム1811により処理されるデータのスループットを高めることができる。したがって、図18のリンク・リスト・エンジン1800は、光ファイバ・リンクと同じ速度でポート1817および1818により必要とされるデータ・キューの読み書きを実行することができる。

【0093】

図19の説明

図19は、5つのバッファ1から5を含む標準的なリンク・リストを開示している。これら5つのバッファのリンク・リストは、図18の同じRAMバンク1803上には格納されない。その代わりに、これら5つのバッファは、5つの別々のバンク1803にランダムに格納される。それぞれのバッファは、システムによって格納され処理される物理的情報またはデータを含む第1の部分を持つ。各バッファの下位部分は、リンク・リストの次のバッファを格納するRAMバンクへのリンク・フィールド・アドレスを含む。図19のバッファ1は、上位部分に物理的情報を格納し、下位部分に0100のリンク・フィールド・アドレスを格納する。バッファ1を格納するRAMバンクのアドレスは、バッファ1の右に示されているように000/00である。

【 0 0 9 4 】

リンク・リストのバッファ 2 は、バッファ 1 のリンク・フィールドで指定されているように R A M アドレス 0 1 / 0 0 に格納される。バッファ 2 の最上位部分は、物理的情報（データ）を含む。下位部分は、R A M バンク I D およびリンク・リストのバッファ 3 を格納する場所を指定する 0 1 0 0 1 のリンク・フィールド・アドレスを含む。

【 0 0 9 5 】

バッファ 3 は、バッファ 2 のリンク・フィールドによって指定されているように R A M バンク・アドレス 0 1 0 / 0 1 に格納される。バッファ 3 のリンク・フィールドは、リンク・リストの第 4 のバッファの場所を指定する 1 1 0 1 0 のアドレスを含む。

【 0 0 9 6 】

第 4 のバッファは、バッファ 3 のリンク・フィールドによって指定されているように R A M アドレス 1 1 0 / 1 0 に格納される。同様に、リンク・リストの第 5 のバッファおよび最後のバッファは、バッファ 4 のリンク・フィールドによって指定されているように R A M アドレス 1 0 1 / 1 1 に格納される。バッファ 5 の最上位部分は、使用できる自由リストの先頭バッファであることを示す。

【 0 0 9 7 】

リンク・リスト上のバッファは、図 1 8 の別の R A M バンク上にランダムに格納される。このランダム性は、システムのデータ処理および制御オペレーションを効率よく行う上で望ましく、また必要である。このランダム性は、本発明の争奪ファシリティの所望のオペレーションを実現するために特に必要であり有用である。

【 0 0 9 8 】

図 2 0 および 2 1 の説明

図 2 0 は、R A M メモリ・バンク 1 8 0 3 内のリンク・リスト・キューの先頭がどのように枯渇し、大容量 R A M メモリ 1 8 0 6 のアクセスを必要とするかを示している。図に示されているリンク・リスト・キューは、それぞれの順序で、内容 Q、F、R、および Z を含むバッファを持つ。図 2 1 では、クロック期間 A 1 の第 1 のアクセスで、R A M 1 8 0 3 から図 2 1 のクロック期間 7 および 8 の内容 Q を保持するアドレスを読み出す。しかし、アクセス A 1 により、使用するキューの先頭が高速メモリ 1 8 0 3 から取り出され、リモート R A M 1 8 0 6 に次の要素 F が格納される。したがって、アクセス A 1 により、アクセス A 1 にかかわる R A M バンク 1 8 0 3 の状態コントローラ 1 8 0 4 を通じてクロック期間 2 でバックグラウンド要求 R q 1 がトリガされる。このリモート・メモリ 1 8 0 6 アクセス要求は、バックグラウンド・アクセス・マルチプレクサ 1 8 0 8 により処理され、図 2 1 のクロック期間 3 でリモート・メモリ・アクセス A 1 B によって返される、リモート R A M 1 8 0 6 は R A M バンク 1 8 0 3 へのアクセスの流れと同時に発生することに注意されたい。データ D 1 B は、クロック期間 9 および 1 0 でリモート R A M 1 8 0 6 から消費される。このデータはリンクされたリスト要素 F である。要素 F は、アドレスが R A M バンク 5 内にある自由リストから取り出された空要素に書き込まれ、アクセス A 5 はクロック・サイクル 1 1 で始まる。リンクは全体を通して保存される。したがって、キューの先頭は、要素 F を含み、再び、ハイパフォーマンス R A M バンク 1 8 0 3 ~ 1 8 0 5 に置かれる。要素 Q を保持している R A M バンク 1 8 0 3 ~ 1 8 0 5 および内容 F を以前保持していたリモート・メモリ要素は、それぞれの自由リストに戻される。

ハイパフォーマンス R A M 1 8 0 3 へのアクセスが無相関であるという性質は温存され、表 2 に至るモデルはそのままとなるだけでなく、アクセス制御のギャップを作ることがないことにより図 2 0 に示されているように、より効率的に利用できるのではないかという期待が持続する。ギャップは争奪によって生じるという前の議論を思い出そう。相関のあるトラフィック、つまり同じバンクへの順次アクセスが争奪を引き起こすということである。

【 0 0 9 9 】

図 2 2 の説明

図 2 2 は、バッファが図 1 8 の様々な R A M バンク 1 8 0 3 - 1 から 1 8 0 3 - 8 (こ

10

20

30

40

50

れ以降RAMバンク1803)にランダムに格納されるリンク・リストの読み出しを要求するアクセス・フロー・レギュレータ1801によって開始される読み出し要求を実行する本発明のプロセスを開示している。アクセス・フロー・レギュレータ1801が読み出しオペレーションを要求するプロセッサ1814から命令を受け取ったときに工程2201でプロセスが開始する。リンク・リストの個別バッファが一度に1つずつ順次読み出され、複数のRAMバンク1803のうち様々なバンクにランダムに格納される。各バッファの読み出しには、アクセス・フロー・レギュレータ1801による別々の読み出し要求が必要である。

【0100】

第1の読み出し要求は、工程2202で受け取られ、工程2203に進み、そこで、RAMバンク1803から読み出される要素の個数に関してしきい値を超えたかどうかを判別する。すでに述べたように、リンク・リストの読み出しでは、先頭バッファが格納されているRAMバンク1803からリストの初期バッファ(先頭)を読み出す必要がある。リンク・リストの実行の残り部分で、リンク・リストの中間バッファを読み出す必要があり、それらをリモートRAM1806に格納し、また取り出してRAMバンク1803に格納しなければならない。リモートRAM1806からRAMバンク1803にバッファを戻す効率よい転送には、そのような複数の要求をバックグラウンド・アクセス・マルチプレクサ1803に適用し、次に、その工程の実行時に効率のためリモートRAM1806に適用する必要がある。このような理由から、しきい値検出要素2203を用意し、個別に一度に1つずつではなく同時にそのような複数の要求をリモートRAM1806に送るようにする。

【0101】

最初に、要素2203がしきい値を超えていないことを判別すると仮定しよう。この場合、リモートRAM1806は即座にアクセスされず、プロセスは工程2204に進み、リンク・リストの第1のバッファ(先頭)の読み出し要求により識別されたRAMバンクを読み出す。このバッファの位置は、読み出され、一時的に格納され、プロセスは工程2205に進み、そこで、読み出された情報をアクセス・フロー・レギュレータ1801に返す。次に、このプロセスは、工程2206に進み、RAMバンク1803がアクセス・フロー・レギュレータ1801から次のアクセス要求を受け取る準備ができていることを示す。末尾バッファも、リンク・リストがただ1つのバッファで構成されている場合と同じように読み出される。

【0102】

次に、要素2203は、工程2202の到着読み出し要求により、読み出しオペレーションにRAMバンク1803内で利用可能なバッファの数に関してしきい値を超えたかどうかを判別すると仮定しよう。この場合、工程2211では、先頭バッファを読み出し、工程2220にさらに進み、工程2211の読み出し情報をアクセス・フロー・レギュレータ1801に送信する。次に、工程2221では、読み出すRAM1806内のリストの先頭付近にある中間バッファの読み出しを要求する。これは、リストの先頭の要求を含む。そうする際に、工程2221は、リストの新しい先頭バッファの読み出し要求を、指定されたRAMバンク1806を処理するバックグラウンド・アクセス・バス1810に出す。次に、工程2222では、リモートRAM1806からリストの複数の中間バッファを取り出す。次にこのプロセスは、工程2223に進み、リモートRAM1806が別のアクセスの実行準備ができていることを示す。このプロセスはさらに、工程2212に進み、そこで、RAMバンク1803に、指定されたRAMバンク1803に書き込まれた工程2222中に読み出されたリモートRAM1806の情報を書き込む。この情報は、指定されたRAMバンク1803内のリストの新しい先頭バッファの形成を含む。このプロセスは、次に、工程2205に進み、情報をアクセス・フロー・レギュレータ1801に送る。このプロセスは、次に、工程2206に進み、読み取り要求の完了を表し、またRAMバンク1803において次のアクセス要求の準備ができていることを示す。

【 0 1 0 3 】

図 2 3 の説明

図 2 3 は、アクセス・フロー・レギュレータ 1 8 0 1 から受け取った書き込み要求を実行するために必要な工程を開示している。このプロセスは、工程 2 3 0 1 から始まり、そこで、アクセス・フロー・レギュレータ 1 8 0 1 は書き込み要求をバス 1 8 0 2 に送る。工程 2 3 0 2 では、書き込み要求を図 1 8 の状態コントローラ 1 8 0 4 に入力する。工程 2 2 0 3 では、リストの末尾付近の書き込まれた要素に関してしきい値を超えたかどうかを判別し、そのしきい値を超えていない場合（新しいリストの場合はそうである）、リストの最後のバッファ（末尾）が R A M バンク 1 8 0 3 に書き込まれる。その後、このプロセスは、工程 2 3 0 5 に進み、そこで、リモート R A M 1 8 0 6 へのアクセス要求が必要でないことをアクセス・フロー・レギュレータ 1 8 0 1 に知らせる。次に、このプロセスは、工程 2 3 0 6 に進み、R A M バンク 1 8 0 3 において次のアクセスの実行準備ができていることを示す。

10

【 0 1 0 4 】

工程 2 3 0 3 が書き込み要求のしきい値を超えていることを判別すると仮定しよう。この場合、プロセスは工程 2 3 1 1 に進み、R A M バンク 1 8 0 3 に保持されているリストの末尾を取り出す。そこで、プロセスは工程 2 3 2 1 に進み、工程 2 3 1 1 で取り出された末尾バッファをバックグラウンド・アクセス・バス 1 8 1 0 に出し、リモート R A M 1 8 0 6 に入れる。次に、工程 2 3 2 2 で、リンク・フィールドをリモート R A M 1 8 0 6 内の最後から 2 番目の最終バッファから工程 2 3 2 1 で書き込まれたバッファの位置に更新するが、それは、最後から 2 番目のバッファのリンク・フィールド内に示されているバッファが、工程 2 3 1 1 および 2 3 2 1 により、R A M バンク 1 8 0 3 からリモート R A M 1 8 0 6 に位置が変更されているためである。工程 2 3 2 3 では、書き込まれたばかりのリモート R A M 1 8 0 3 はアクセスの準備が整っていることを示す。

20

【 0 1 0 5 】

工程 2 3 2 1 と同時に、工程 2 3 1 2 では、空のバッファがリンク・リストの末尾に書き込まれ、R A M バンク 1 8 0 3 に入る。工程 2 3 2 1 で書き込まれたバッファへのポインタは、工程 2 3 1 2 で書き込まれた空のバッファのリンク・フィールドに書き込まれる。このプロセスは、次に、工程 2 3 0 6 に進み、書き込み要求の完了を表し、また R A M バンク 1 8 0 3 において次のアクセス要求の準備ができていることを示す。

30

【 0 1 0 6 】

図 2 4 の説明

リンク・リスト・ファイル进行处理するためのメモリ管理ファシリティについて説明した。本発明の他の可能な実施形態では、開示されているメモリ管理ファシリティはさらに、高速 R A M バンク・メモリ 1 8 0 3 およびリモート・メモリ 1 8 0 6 を使用して、リンク・リスト型でないデータ・ファイル进行处理することもできる。これは、図 2 4 および 2 5 の処理工程について説明される。

【 0 1 0 7 】

以下では、図 1 8 の読み出しオペレーションが、高速小容量 R A M バンク 1 8 0 3 にアクセス・フロー・レギュレータ 1 8 0 1 から受け取った情報を格納するような方式で動作するように構成されている R A M バンク 1 8 0 3 およびリモート R A M 1 8 0 6 とともに説明されている図 2 4 の処理工程を説明する。リモート R A M 1 8 0 6 は、大きなファイル用の情報を格納する際のオーバーフローとして使用される。このプロセスは、工程 2 4 0 1 から始まり、工程 2 4 0 2 に進み、そこで、アクセス・フロー・レギュレータ 1 8 0 1 は読み出し要求をバス 1 8 0 2 に送り、R A M バンク 1 8 0 3 に格納されているその情報を要求する。

40

【 0 1 0 8 】

工程 2 4 0 3 で、取り出されるファイルのサイズが R A M 1 8 0 3 の現在の記憶容量を超えるかどうかを判別する。しきい値を超えていなかった場合、プロセスは工程 2 4 0 4 に進み、そこで、R A M バンク 1 8 0 3 は要求されたファイルを読み出し、アクセス・

50

フロー・レギュレータ 1801 に送り返す。次に、このプロセスは、工程 2405 に進み、状態コントローラ 1804 およびバス 1802 を介して RAM バンク 1803 から読み出された要求情報をアクセス・フロー・レギュレータ 1801 に返す。アクセス・フロー・レギュレータ 1801 は、情報を受け取り、それをプロセッサ 1814 に受け渡し、その際にコントローラでは要求された情報に関連する機能を使用する。次に、このプロセスは、工程 2406 に進み、RAM バンク 1803 が次のアクセス要求を受け取る準備ができていることをアクセス・フロー・レギュレータ 1801 にアドバイスする。

【0109】

要素 2403 で、読み出されるファイルのサイズがしきい値を超えたと判定された場合、プロセスは工程 2410 に進み、そこで、RAM バンク 1803 にあると思われるファイルの一部を読み出す。このプロセスは、さらに工程 2411 に進み、そこで、読み出された情報がアクセス・フロー・レギュレータ 1801 に返される。プロセスは工程 2412 に進み、リモート RAM 1806 からの読み出し要求をバックグラウンド・アクセス・バス 1810 に出し、このバスから、バックグラウンド・アクセス・マルチプレクサ 1808 を通じて、要求されたリモート RAM 1806 に要求が送られる。

【0110】

その後、このプロセスは、工程 2413 に進み、リモート RAM 1806 から要求された情報を取り出す。次に、工程 2415 で、工程 2413 で取り出された情報をリモート RAM 1806 から RAM バンク 1803 に送り、そこに格納する。次に、このプロセスは、工程 2414 および 2406 に進むが、これらの工程は両方とも、RAM バンク 1803 が他のアクセス要求を受け取る準備ができていることをアクセス・フロー・レギュレータ 1801 に知らせる。

【0111】

図 25 の説明

図 25 では、RAM バンク 1803 にすでに存在しているファイルに追加する書き込み要求を説明している。RAM バンク 1803 およびリモート RAM 1806 は連携動作し、図 24 の読み出しオペレーションについて説明されている方法と類似の方法でアクセス制御レギュレータ 1801 から受け取った大容量データを格納する。RAM バンク 1803 は、ファイルに対し既存のが選択された量のデータを格納し、大容量ファイルの残り部分はオーバーフローされ、リモート RAM 1806 に書き込まれる。

【0112】

このプロセスは工程 2501 から始まり、要素 2502 に進み、そこで、アクセス制御レギュレータ 1801 から受け取った書き込み要求を解析し、書き込む要求されたファイルのサイズが RAM バンク 1803 に格納することが可能な容量を超えているかどうかを判別する。しきい値を超えていない場合、要素 2502 により工程 2504 では、選択された RAM バンク 1803 にすでに存在している関連ファイルに追加データを書き込む。その後、このプロセスは要素 2505 に進み、そこで、要求されたファイルが RAM バンク 1803 に書き込まれたこと、およびリモート RAM 1803 への書き込みが不要であることを通知する肯定応答をアクセス制御レギュレータ 1801 に送り返す。

【0113】

しきい値を超えたと要素 2502 が判定した場合、プロセスは工程 2511 に進み、RAM バンク 1803 に格納されているファイルの格納済み部分を読み出す。さらに、工程 2522 に進み、工程 2511 で取り出されたファイルのこの部分をリモート RAM バンク 1806 に書き込むために必要なオペレーションを開始する。工程 2523 では、リモート RAM 1806 において再びアクセスの準備ができていることを示す。工程 2512 では、書き込みポインタを RAM バンク 1803 に書き込み、他の部分が RAM 1803 に書き込まれるファイルの残り部分を含むリモート RAM 1806 のアドレスを関連付けられるようにする。その後、このプロセスは、工程 2506 に進み、RAM バンク 1803 において他のアクセス要求を受け取る準備ができていることを示す。

【0114】

結末

ネットワーク・トラフィック・シェーピングには、柔軟なバッファリング機能が必要である。優先度の高い流れは、他の優先度の低い流れに優先する。例えば、リアルタイム・トラフィックの需要は、ファイル転送などの通常のデータ・トラフィックの需要よりも高い。しかし、トラフィックの一瞬一瞬の特性はランダムである。次に入って来るパケットは、どの流れにも送ることができる。例えば、再び図3を考察しよう。リンク1上でノードAに向かう内向きパケットのストリームは、リンク2またはリンク3のいずれかへ外向きとなる可能性がある。リンク1上で内向きのトラフィックは、リンク2に向かう数千個もの連続するパケットを伝送し、リンク3に向かう単一のパケットを伝送し、その後、数千個の連続パケットが再び、リンク2に向けられる。パケットを破棄しないこと、またケーブル切断などの大惨事からの復旧をサポートするという要件に加えて、遅延時間を変更できる、スケジュールされた外向きトラフィック・シェーピング・アルゴリズムをサポートしていなければならない。したがって、トラフィックの着信と送信の両方の性質は複雑であり、柔軟なバッファリングを必要とする。

【0115】

内向きストリームと外向きストリームに対するバッファ処理のバランスをとるために最も効率のよい配置は、ハードウェアによるリンク・リスト・エンジンである。しかし、ハードウェアによるリンク・リスト・エンジンの現在の実装は、大容量で低速か、または小容量で高速のいずれかである。本発明の改良されたリンク・リスト・エンジンは、安価にギガバイト・クラスのバッファ記憶容量が得られ、市販の半導体メモリの最大スループットで動作するため、従来のリスト・エンジンよりも優れている。したがって、本発明の改良されたリンク・リスト・エンジンは、最新技術のものと同じ大容量であり、実質的に高速化されている。速度がアップすることは、改良されたリンク・リスト・エンジンでは大容量（つまり、OC-768）ファイバ回線をサポートしており、同じ規模のハードウェア資源、つまりFPGAおよびSDRAMを利用できるため、トラフィック・シェーピング・アプリケーションには魅力的である。改良されたリンク・リスト・エンジンは、バッファがギガバイト程度の深さをとれるため、大容量の回線をサポートするようにできる。

【0116】

本発明では、市販のダイナミック・メモリのバースト・モードをサポートしている。このサポートは、本発明の2つの面で明らかである。第1に、市販のダイナミックRAMへのアクセスが連続アクセスであり、続いている場合、本発明では、市販のダイナミック・メモリのデータ・ピンを占有し、取り出された情報を読み出すか、または格納すべき情報を書き込み、そうしながら、次の連続アクセスを開始する。第2に、本発明では、市販のダイナミックRAMに書き込まれるバッファおよび市販のダイナミックRAMから読み出したバッファを状況に応じてキャッシュする。

【0117】

本発明では、現在のアクセス要求に対するデータ・ピンの処理と次のアクセス要求の開始を同時に行うことにより、市販のダイナミックRAMのデータ・バスのデータ・バス・ラッチ機能を利用する。この機能を使用すると、データ・バスの速度で、通常、166メガヘルツ以上で、ダイナミック・メモリのロー全体を読み出すか、または書き込むことができる。このローは、長さは128ビットとすることができるが、ダイナミック・メモリ上で利用できるデータ・バス・ピンは、8本だけである。したがって、メモリは内部的に、読み出されるメモリのローをラッチし、この情報を一度に8ビットずつピンに送出するか、または一度に8ビットずつ内向きデータをラッチしてから、ロー全体を内部メモリ・アレイに書き込む。本発明では、このようにしてアクセスをオーバーラップさせることにより、データ・バス上にデータを出しっぱなしにすることで、市販のダイナミック・メモリからの連続読み出しが図18のRAMバンク1803などの本発明のRAMバンクと互換性のある速度で実行されるようにする。これは、多数の長いバッファのリストに連続してアクセスする場合に有用である。このような場合、市販のダイナミック・メモリからストリーミングで情報を出力しなければならない。市販のダイナミック・メモリへのこのよ

うな拡張された一連のアクセス時にパフォーマンスを維持するために、市販のダイナミック・メモリとのインターフェースの定常状態パフォーマンスは他のメモリと整合性がとれていなければならない。

【0118】

書き込まれるバッファと読み出されたバッファをキャッシュすることにより、市販のダイナミック・メモリを効率よく使用することができる。市販のダイナミック・メモリにはロー毎にアクセスするのが最も効率的である。しかし、ローは、128ビット以上のオーダーと、通常大きい。格納されるまたは取り出される個々のバッファは、16ビットのみで構成できる。したがって、十分な情報が集まり完全なローを構成するまでバッファをキャッシュした場合、市販のダイナミック・メモリへの書き込みアクセスにより、完全なローが転送され、これが最も効率的である。同様に、読み出し後、完全なローがキャッシュされた場合、市販のメモリに対しロー毎に1回アクセスするだけでよく、最も効率的である。

10

【0119】

したがって、書き込みおよび読み出しをキャッシュすると、実質的なパフォーマンスの向上が得られるが、これは、市販のダイナミック・メモリに対するアクセス回数が少なく済み、アクセス回数を減らすことで、市販のダイナミック・メモリに対する争奪が少なくなるためである。

【0120】

本発明では、多数の同時リンク・リストをサポートしている。多数の同時リンク・リストのうちの任意のリンク・リストとのそれぞれの相互作用は、多数の同時リンク・リストの残りとは完全に独立している。

20

【0121】

添付の請求項では、RAMバンク1803を高速メモリと、またリモート・メモリ1806を大容量メモリと特徴付けている。

【0122】

上記の説明では、本発明の可能な実施例を開示している。当業者であれば、請求項に規定されているとおり文字通りに、または均一者の原則に従って、本発明を侵害する他の実施形態を設計することが可能であり、また設計するであろうことは予想される。

【図面の簡単な説明】

30

【0123】

【図1】マルチノード・ネットワークを開示する図である。

【図2】1つのノードを含むハードウェア要素を開示する図である。

【図3】ブロッキングのある図1のマルチノード・ネットワークを開示する図である。

【図4】リンク・リスト・バッファの配置を開示する図である。

【図5】リンク・リスト・バッファの配置を開示する図である。

【図6】リンク・リスト・バッファの配置を開示する図である。

【図7】図3のノード上の仮説的なトラフィック状態を開示する図である。

【図8】複数のRAMメモリ・バンクで処理されるネットワークによるアクセス要求の処理を開示する図である。

40

【図9】複数のRAMメモリ・バンクで処理されるネットワークによるアクセス要求の処理を開示する図である。

【図10】制御バスおよびデータ・バスに接続されている4つのRAMメモリ・バンクを開示する図である。

【図11】図8および9に示されているように図10のRAMメモリ・バンクがアクセス要求を処理するプロセスを説明するタイミング図である。

【図12】4つのRAMメモリ・バンクを持つリンク・リスト・プロセッサにより制御されるメモリ・システムを開示する図である。

【図13】図12のシステムのオペレーションを例示するタイミング図である。

【図14】図12のシステムの他のオペレーションを例示する他のタイミング図である。

50

【図 15】 8つのRAMメモリ・バンクを持つプロセッサにより制御されるリンク・リスト処理システムを開示する図である。

【図 16】 図 15 のシステムのオペレーションを例示するタイミング図である。

【図 17】 図 18 の状態コントローラ 1804 の要素を開示する図である。

【図 18】 本発明を具現化するプロセッサにより制御されるRAMメモリ・システムを開示する図である。

【図 19】 本発明によるリンク・リスト・バッファの配置を開示する図である。

【図 20】 バッファをリンク・リストの新しい先頭として作成することを例示する読み出しオペレーションを開示する図である。

【図 21】 図 18 のシステムのオペレーションを例示するタイミング図である。

10

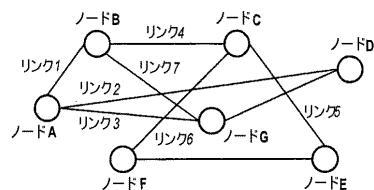
【図 22】 本発明のオペレーションを例示する流れ図である。

【図 23】 本発明のオペレーションを例示する流れ図である。

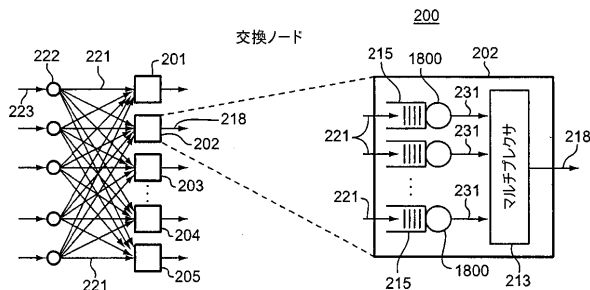
【図 24】 本発明のオペレーションを例示する流れ図である。

【図 25】 本発明のオペレーションを例示する流れ図である。

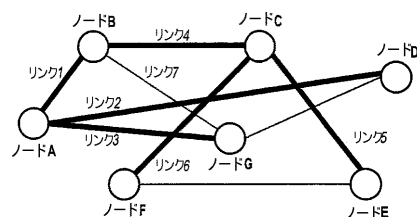
【図 1】



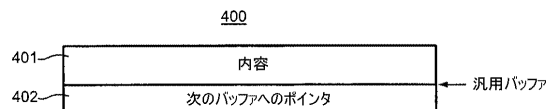
【図 2】



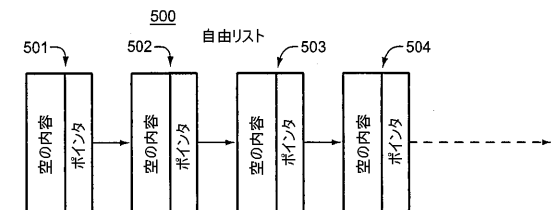
【図 3】



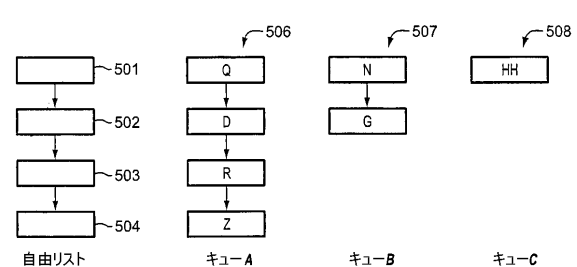
【図 4】



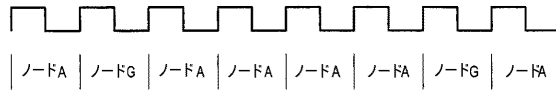
【図 5】



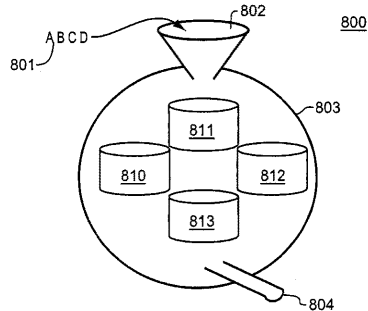
【図 6】



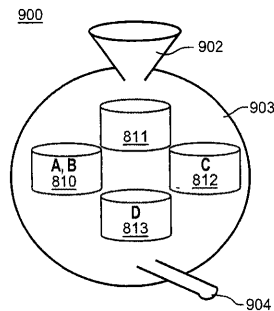
【図 7】



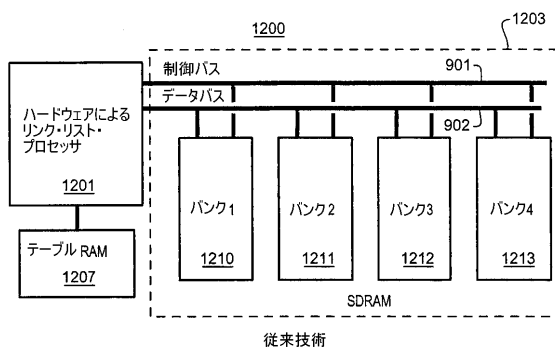
【図 8】



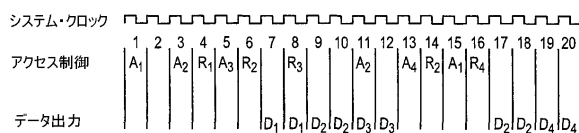
【図 9】



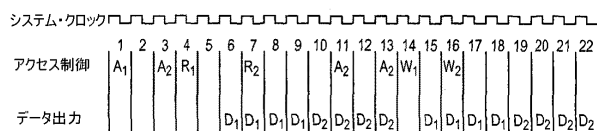
【図 12】



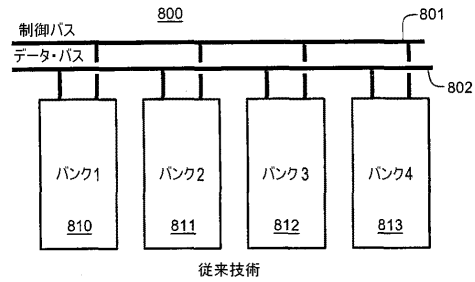
【図 13】



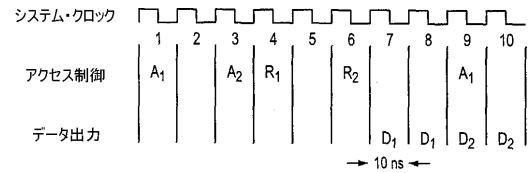
【図 14】



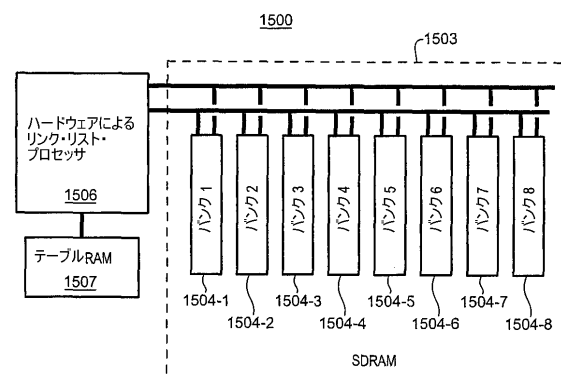
【図 10】



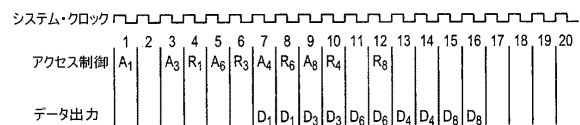
【図 11】



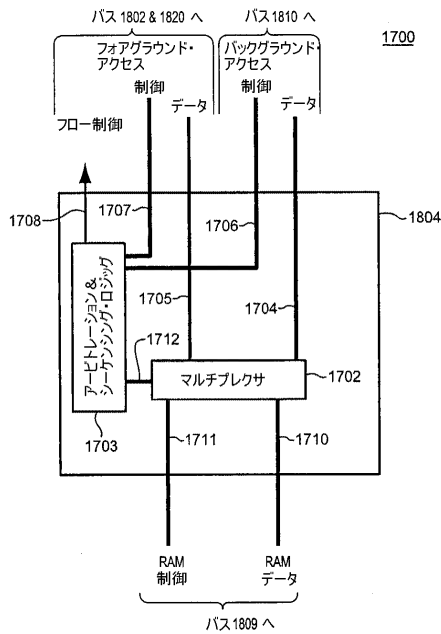
【図 15】



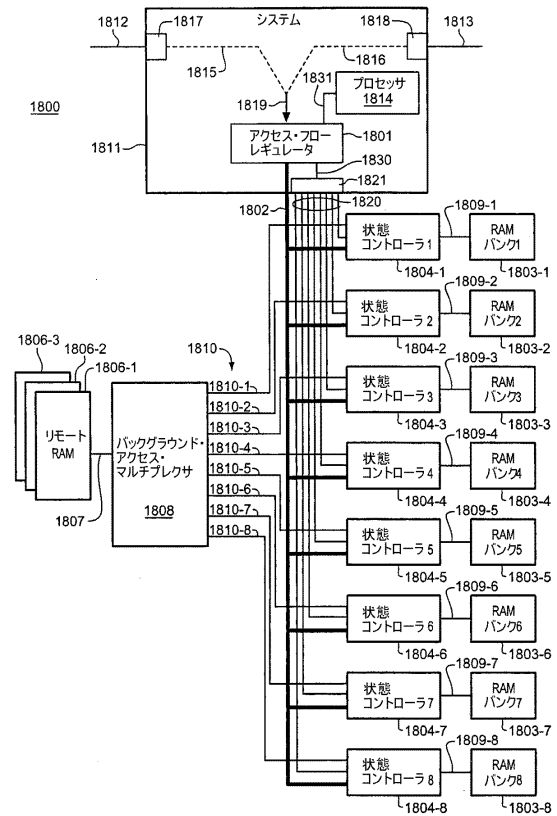
【図 16】



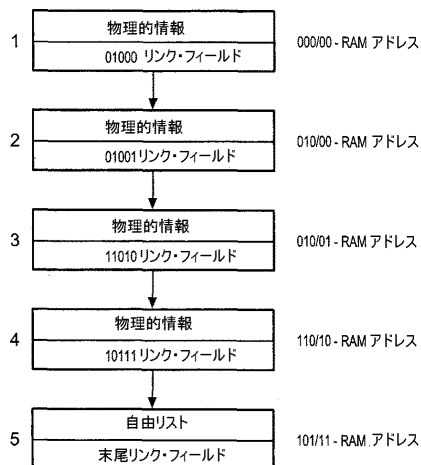
【図 17】



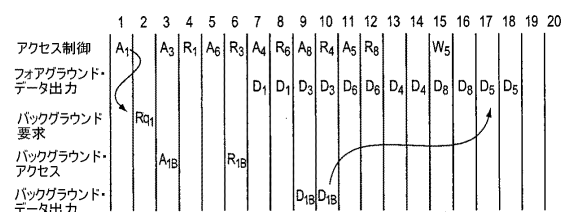
【図 18】



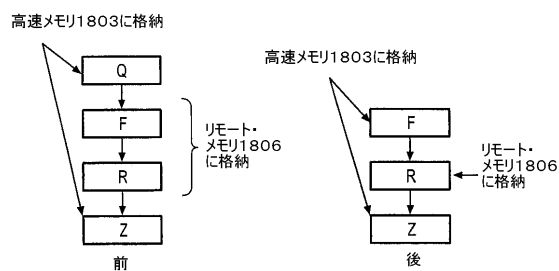
【図 19】



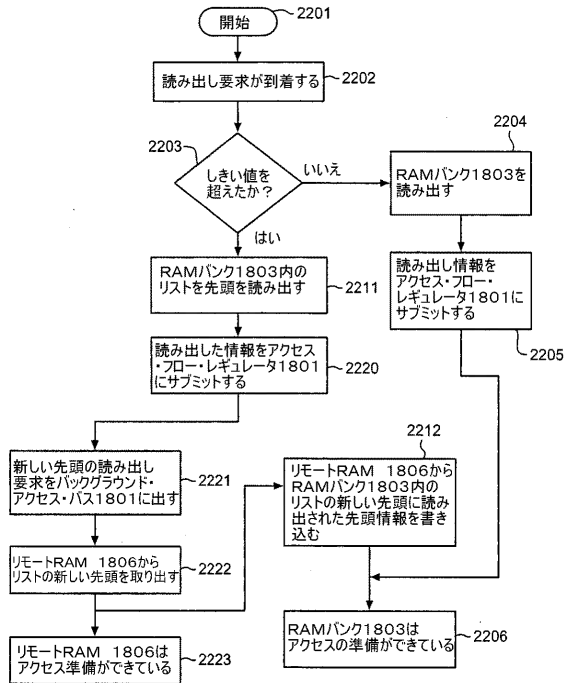
【図 21】



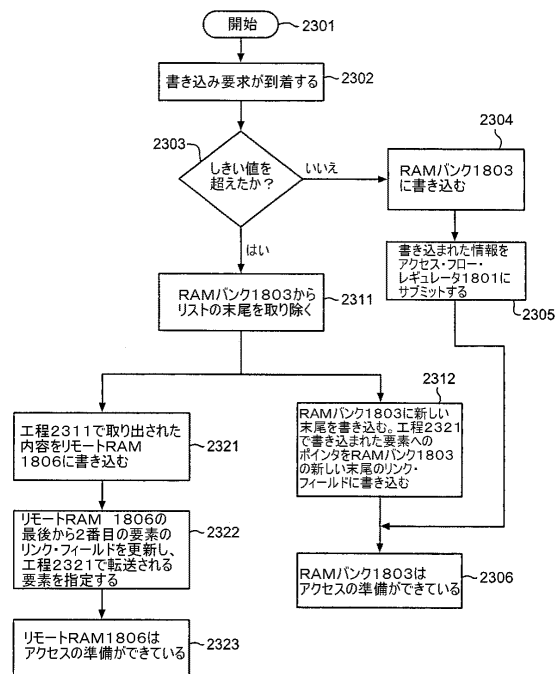
【図 20】



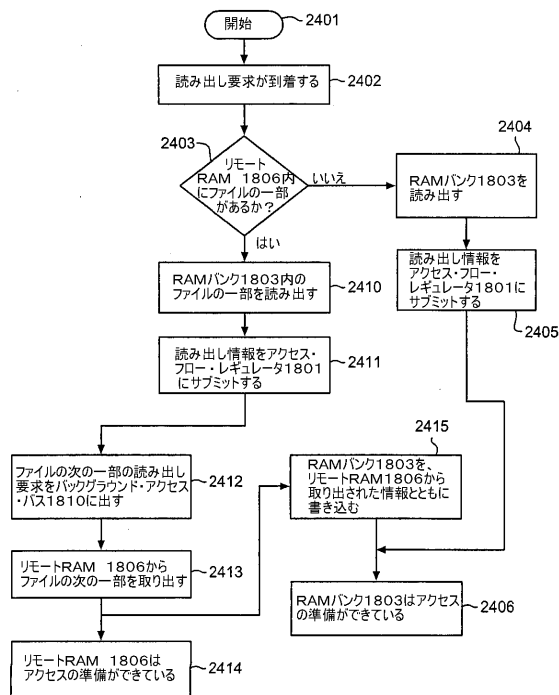
【図 22】



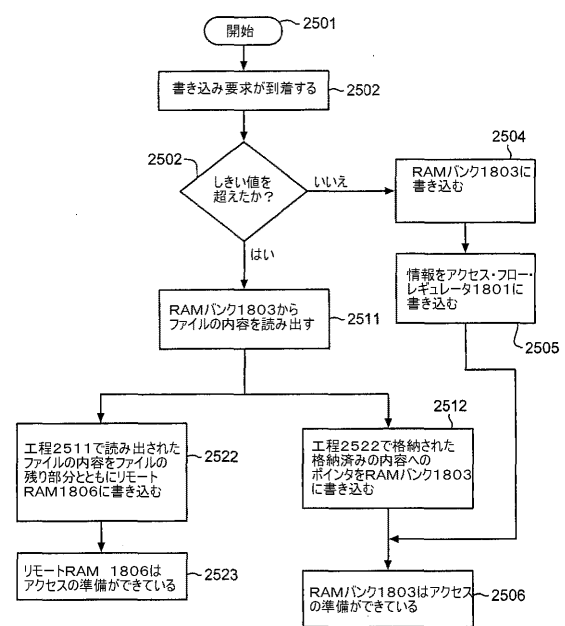
【図 23】



【図 24】



【図 25】



フロントページの続き

(74)代理人 100101498

弁理士 越智 隆夫

(74)代理人 100096688

弁理士 本宮 照久

(74)代理人 100104352

弁理士 朝日 伸光

(74)代理人 100128657

弁理士 三山 勝巳

(72)発明者 ピーター ジェー・ジーヴァース

アメリカ合衆国 60540 イリノイス, ネイパーヴィル, ブリドルスパー ドライヴ 498

審査官 坂東 博司

(56)参考文献 米国特許第6798777(US, B1)

米国特許出願公開第2001/0000815(US, A1)

特開平11-3596(JP, A)

特開2002-287757(JP, A)

特表2004-523017(JP, A)

特開平10-276220(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/38

G06F 12/02

G06F 12/06

H04L 13/08