

Virtual Queue : 超並列計算機向きメッセージ通信機構

五島 正裕[†] 池田 泰敏 森 眞一郎

中島 浩 富田 眞治

京都大学 工学部

〒606-01 京都市 左京区 吉田本町

E-mail: {goshima, yiked, nakasima, tomita}@kuis.kyoto-u.ac.jp

内容梗概

Virtual Queue システムは、通信機構の仮想化とストリームのキャッシングを特長とするメッセージ通信のための専用ハードウェアである。仮想化によって、任意の数の仮想的通信機構を提供でき、通常のキャッシュと同様システム・コールなしで利用できる。このようにメッセージ送受信時のソフトウェア・オーバーヘッドを大幅に削減するので、細粒度の通信に耐える。また、メッセージそれ自体ではなく、メッセージのストリームへの参照の局所性を利用し、よくアクセスされるストリームを物理的通信機構にキャッシングすることによって高速化を図る。物理的通信機構はベクトル・データの転送にも十分なスループットを持つ。

[†]: 日本学術振興会特別研究員

Virtual Queue : A Message Communication Mechanism for Massively Parallel Computers

Masahiro Goshima[†] Yasutoshi Ikeda Shin-ichiro Mori

Hiroshi Nakashima Shinji Tomita

Faculty of Engineering, Kyoto University

Yoshida-hon-machi, Sakyo-ku, Kyoto 606-01 Japan

E-mail: {goshima, yiked, nakasima, tomita}@kuis.kyoto-u.ac.jp

Abstract

The **Virtual Queue** system is a specialized message communication hardware that is marked by the virtualization of the message communication mechanism and the caching of streams of messages. The virtualization provides arbitrary number of virtual communication mechanisms, and enables a user to access the hardware without system call like a usual cache memory. Since these reduce the message handling overhead of software drastically, it withstands the fine grain communication. Moreover the system makes use of the locality of reference, not to the message themselves but to the message stream, and speeds up by caching frequently accessed message streams to physical units. The units provide enough throughput for transmission of vector data.

[†]: Research Fellow of the Japan Society for the Promotion of Science

1 はじめに

ノード間の距離が大きくなる超並列規模の計算機では、メッセージ処理の高速化がますます重要になる。しかし、従来のメッセージ通信の方式はメッセージ処理のオーバーヘッドが大きく、細粒度の通信に適さない。また、通常のキャッシュはメッセージに対しては効果的ではなく、特にベクトル・データの転送に適さない。

本稿では、メッセージ送受信時のオーバーヘッドの低減と、キャッシュに代わるメッセージ向き的高速化方式の実現を目標とするメッセージ通信用ハードウェア、Virtual Queue システムを提案する。以降 VQ と略す。

2 従来のメッセージ通信の問題点

本章では、従来のメッセージ通信方式におけるメッセージ送受信のオーバーヘッドと、通常のキャッシュのメッセージに対する不適合について述べる。

2.1 メッセージ送受信のオーバーヘッド

共有メモリ/メッセージ、両方式において、メッセージ送受信のオーバーヘッドは数百サイクル程もある。このオーバーヘッドは、細粒度のメッセージを頻繁に交換する応用において、性能に対する影響が大きい。

2.1.1 共有メモリ方式のオーバーヘッド

一般に、共有メモリ・システム上でメッセージ通信を行う場合には、遠隔メモリ上のバッファに対するロック操作が必要となる。超並列規模の計算機では、遠隔ノードへのアクセス（ロック操作）のレイテンシは数百サイクルにもなる。

単なるメッセージの読み出しの場合とは異なり、更新型のキャッシュ一貫性制御やプリフェッチなどでは、ロック操作のレイテンシを隠蔽することはできない。

2.1.2 メッセージ方式のオーバーヘッド

メッセージ方式の計算機では、メッセージ処理を高速化するハードウェアを持つことが一般的となっているが、これらの計算機におけるオーバーヘッドはむしろメッセージ処理のソフトウェア部分で生じる。

一部のソフトウェア処理は最適化やハードウェアの工夫によって省略可能であるが、資源の獲得とアクセス保護のためのシステム・コールは省略することができない。また受信時には、必要なメッセージを検索したり、システム領域に置かれたメッセージをユーザ領域にコピーする必要のあるシステムも存在する。これらのソフトウェア処理のため、500~1000 サイクル程のオーバーヘッドが生じる [1]。

2.2 メッセージのキャッシング

共有メモリ/メッセージの両方式において、キャッシュ・メモリの果たす役割は大きい。しかし、通常のキャッシュはメッセージに対しては効果的でなく、その問題点は特に数値処理などの大規模ベクトル処理で顕著に現れる。

2.2.1 参照の時間的局所性とキャッシング

メッセージに対する参照はたかだか 1 回程度であることが多く、時間的局所性の観点からは、メッセージをキャッシングする意味は薄い。それにもかかわらず、メッセージをキャッシングすると、本来キャッシングしておくべき他のデータを追い出してしまい、かえって性能を低下させかねない。ベクトル処理などでは、長大で参照回数の少ないデータを交換するのでこの傾向が強くなる。

2.2.2 キャッシュを利用したバッファリング

一方、1つのメッセージに対する参照は時間的に集中して行われることが多く、バッファリングによるレイテンシの隠蔽が有効に働く可能性が高いと考えられる。

実際、キャッシュをバッファとして利用する方式がいくつか提案されており、これらの方式では以下のようにしてメッセージ・アクセスのレイテンシを隠蔽する：送信側ではキャッシュ上で生成したメッセージをキャッシュから直接送信し [1][2][3]、受信側ではプロセッサが参照する前にメッセージをキャッシュに移しておく [2][3][4]。しかしキャッシュをバッファとして利用する方式では、ベクトル処理のようにプロセッサの稼働率とメモリ・アクセスの頻度が高い場合、以下の問題が生じる：

1. メッセージが 2 度ずつキャッシュ・メモリのポートを通過するため、スループットが不足する。
2. プリフェッチの処理量の増加が性能に影響を及ぼす。
3. また、前節で述べた他のデータの追い出しの問題も解決されない。

3 VQ システムの設計方針

VQ システムは、前章で述べたの問題点を以下のような方針で解決する。

3.1 オーバーヘッドの低減

共有メモリ方式では、各プロセッサがメッセージのバッファを分散管理しているのでバッファに対するロック操作を避けることはできないが、一方、メッセージ方式では、受信ノードにおいてバッファを集中管理するので明示的なロック操作を必要としない。しかし、メッセージ処理の高速化のためにバッファ管理用のハードウェアを導入すると、その獲得とアクセス保護のためのシステム・コールが必要となる。

VQシステムは、バッファ管理用のハードウェアを導入し、それを仮想化することでOSの介入を避けるアプローチを採る。

さて、キャッシュ・メモリはシステムの共有資源であるにもかかわらず、OSをまったく介することなく利用できる。それは資源(キャッシュ・ブロック)の獲得/横取りと記憶保護をハードウェアで実現しているためである。キャッシュと同様に、VQシステムは以下のようにしてOSを介さないハードウェアの利用を可能にする。

3.1.1 仮想化

ハードウェア・ユニットの処理中の状態を定義するのに十分な情報をそのユニットのコンテキストと呼ぶ。例えば、サーキュラ・バッファの制御回路は、バッファの書き込み/読み出し位置やエントリ数などのバッファの管理情報を保持するレジスタを持つが、それらのレジスタの内容はその制御回路のコンテキストとなる。

VQシステムでは、通信機構のコンテキストを主記憶上の退避領域に置き、コンテキストを適宜切替えることによって、その通信機構を時分割多重利用する。コンテキストをキャッシュする専用のキャッシュを用意し、通信機構のコントローラがそのキャッシュを使用する形で、これを実現する。

コンテキストをキャッシングすることで、通常のキャッシュと同様に、仮想記憶方式を応用することによって記憶保護が実現でき、また、キャッシュ・ミスの処理をハードウェアで自動的に行うことによって資源(通信機構)の獲得/横取りをOSを介さずに実行することができる。

一方、通信機構の個々のコンテキストは主記憶上におかれるので、実際上任意の数の仮想的な通信機構をユーザに提供できる。ユーザは、メッセージのストリームごとに仮想通信機構を割り付けることによって、受信時にメッセージの検索やコピーを避けることができる。

このようにVQシステムでは、従来のメッセージ方式で必要であったソフトウェア処理のほとんどを省略することができる。

3.2 高速バッファ・メモリの小量化と高速化

高速なメモリにデータを移しておくことによってアクセスのレイテンシを短縮することは重要であるが、通常のキャッシュは主に参照の時間的局所性を利用して高速化を図るものであるから、メッセージのようにその性質の弱いデータに対しては効率的に働かない。メッセージに対するアクセスはただか一回程度であることが多く、一度アクセスされたメッセージをキャッシュに置いておくことは無駄であることが多い。

3.2.1 FIFOによる通信モデル

VQシステムでは、通信モデルをFIFOによるものだけに制限することによって、この問題を解決する。すなわち、VQシステムは通常のキャッシュと同様にメッセージを高速なバッファ・メモリに移しておくのだが、メッセージに対する参照の回数を1度に限定することによって、古いメッセージの占めていたバッファ・メモリ上の領域の再利用を図る。

FIFOによる通信のみを用いる言語も多く、この制限は決して厳し過ぎるものではないと考えられる。プロセス間チェイニングによるベクトル処理などにも対応可能である。

また、FIFOによる通信モデルを仮定することによって、プリフェッチを自動化できるという効果も得られる。

3.2.2 ストリームのキャッシング

一度アクセスされたメッセージは2度とアクセスされないことが多いが、一度アクセスされたメッセージのストリームは近い将来再びアクセスされる可能性が高いと考えられる。VQシステムは、メッセージそれ自体ではなく、メッセージのストリームに対する参照の局所性を仮定し、これを利用する。

前述したハードウェアのコンテキストのキャッシングは、よくアクセスされるストリームを高速な物理的通信路に置くことと捉えることができる。コンテキスト・キャッシュはプロセッサごとに4~16エントリ程度しか用意しないが、プロセッサで走行中のスレッドがそれ以下のストリームを使用するならばVQシステムへのアクセスはヒットし続ける。

3.2.3 バッファの構成

以上の2点によって、バッファリングのために必要なバッファの容量は大幅に削減される。したがって、デュアルポート・メモリを用いてこのバッファを構成し、メッセージのバッファリング専用を用いる。

デュアルポート・メモリを用いてバッファリングを行うので、バッファリングの実施によって系のスループットが低下することはない。また、メッセージのバッファリング専用を用いるので、通常のキャッシュ上の他のデータを追い出すことがない。

4 VQシステムの構成

本章では、以下の手順でVQシステムについて詳細に説明する。まず4.1節ではVQシステムがプログラムからどのように見えるかを示し、プログラムから見える動作を実現するハードウェアについて、4.2節以降で述べる。

4.1 プログラムに対するインターフェイス

4.1.1 VQ

図1に、VQシステムにおける通信の様子を示す。

VQは、送信スレッドと受信スレッドの間の共有仮想空間上に設定された仮想的な通信機構である。プログラムは任意の数のVQを設定できるので、通常メッセージのストリームごとにVQを設定する。

VQは、3種の仮想空間上の領域—ライト・ポート(WP)、サーキュラ・バッファ(CB)、および、リード・ポート(RP)を構成要素とする。1つのVQは、1つ以上のWP、および、1組のCBとRPからなる。

VQは、数十バイトのブロックを通信の単位とする。ブロックは不可分に扱われるので、転送経路各部の実効バンド幅を向上させ、管理情報を小量化する他、リモート・プロシージャ・コールなどの応用において有用である。

ブロックの送信はWPへのストアによって、ブロックの読み出しはRPへのロードによって行う。1つのWPから送信されたブロックは、送信された順序でRPから読み出されるが、複数のWPから送信されたブロック間の順序は非決定的となる。

4.1.2 VQのアドレス

VQがプログラムからどう見えるか、次項からより詳細に説明する。

ユーザは、VQの使用に先だってVQのデータ構造を設定する必要がある。データ構造は、VQシステムのハードウェア・ユニットがそのコンテキストの位置を特定するため、アクセスされたアドレスとコンテキスト回避領域のアドレスの間の簡単な写像を定義する。

(A) システムのアドレス空間の構成

1つのプロセス内のスレッドは1つの仮想アドレス空間を共有するものとする。ただし、遠隔ノードのメモリに対して直接アクセスが行える必要はない。ノード間にまたがる大域的仮想アドレス空間の構成方法は本稿の範囲を越えるので、ここでは触れない。

VQシステムの実アドレス空間中には、主記憶の物理アドレス空間とは別に、VQアクセスのためのアドレス空間—ポート・アドレス空間—を設ける(図2左)¹。ポート・アドレスは、主記憶の物理アドレスに対する一種の別名である。ポート・アドレス空間は、主記憶の物理アドレス空間と等しい大きさを持つ。また、ポート・アドレス空間上のページ—ポート・ページ—は、それぞれのアドレス空間内での変位が互いに等しい物理ページと対応する。

VQを使用するユーザ・プロセスは、UNIXではmmapに相当するシステム・コールにより、物理ペー

¹その他にメモリ・マップI/Oのアドレス空間があってもよい。

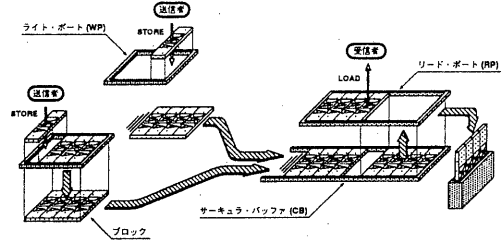


図1: プログラム・モデル

ジに対応するポート・ページを獲得し、仮想ページに写像する(図2左から中)。OSは、ユーザが対応する物理ページを既に獲得している場合に限り、ポート・ページへの写像を許可する(理由は後述)。

物理アドレス空間に写像された仮想ページを仮想メモリページ、ポート・アドレス空間に写像された仮想ページを仮想ポート・ページと呼ぶ。

これらのページ対の上に、VQのデータ構造を設定するのだが、一旦ページ対を割り付けられた後には仮想記憶によってアクセス保護が実現されるので、VQに関するすべてのアクセスはユーザ・モードで行うことができる。

(B) VQのデータ構造

プログラムは、図2右に示すように、VQのデータ構造を設定する。

VQの各要素は、アラインされた1ブロック・サイズの領域—ブロック枠—を単位とする。

プログラムはまず、仮想メモリ・ページ上に各要素の回避領域とCBを割り付ける。ただし、RPユニットとCBユニットのコンテキストは1つのブロック枠に詰め、CBはその次のブロック枠からはじまる1ページ内の連続領域に割り付ける(図2右下)。

WPとCB/RPは、送信先のCB/RPの仮想アドレスをWPのコンテキストとして指定することで結び付けられる(図2右下の矢印)。図2では、説明の都合上WPとCB/RPを同一ページに設定しているが、実際には相互の位置関係に制限はない。

仮想ポート・ページ上の、回避領域とページ内変位が互いに等しいブロック枠がそれぞれのポートとなる。ただし、RPは2ブロック枠を使用する(図2右上)。この領域に対するストア/ロードにより、VQによる送/受信を行う。

ポート・ページと物理ページは、各アドレス空間内での変位が互いに等しい位置にとった(図2左)ので、ポートと回避領域の実アドレスは各アドレス空間内での変位が等しい。したがって、VQシステムのハードウェアは、ポートへのアクセスがミスを起こした時、ポ

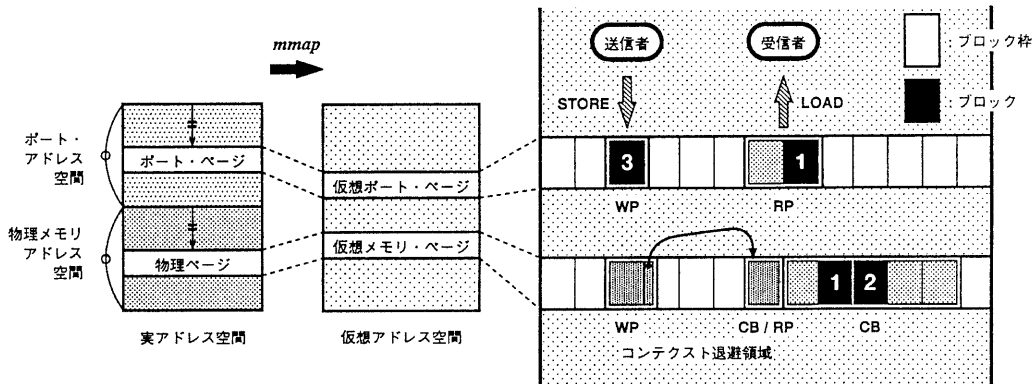


図 2: VQ のデータ構造

トの実アドレスのポート・アドレス空間内の変位から対応する退避領域の物理アドレスを知ることができる。

4.1.3 通信手順

上記のデータ構造を設定した後に、VQ による通信が可能となる。通信は以下を行う：

1. 送信 送信者は、WP に対するストアにより、WP 上にブロックを形成する。WP の最後のバイトが変更を受けた時、その時点での WP 上のブロックが CB の末尾に追加される。転送すべきデータが 1 ブロックに満たない場合には、最後のバイトへのダミー・ライトを行う。

最後のバイト以外の部分に対する書き込みには、順序、回数などの制限はない。

プログラムは送信処理の完了を確認する必要はない。

2. 受信 RP の 2 つのブロック枠を用いて 2 面バッファリングを行う：2 つのブロック枠のどちらか一方には CB の先頭のブロックが、他方には 2 番目のブロックが写像されている。先頭のブロックへの参照が終了し 2 番目のブロックへの参照を開始すると、参照の終わったブロックは自動的に CB からデキューされ、空いたブロック枠には新たに 2 番目となったブロックが写像される。

1 つのブロック枠内のデータに対する参照には、順序、回数などの制限はない。

4.2 ハードウェアの構成

VQ システムは、WP、CB、および RP それぞれに対応するユニットからなり、これらのユニットの協調動作によって上述の通信処理を実現する。

ポートへのアクセスに対して各ユニットは、そのコンテキストを主記憶上の退避領域からコンテキスト・キャッシュにキャッシュし、コンテキスト・キャッシュ

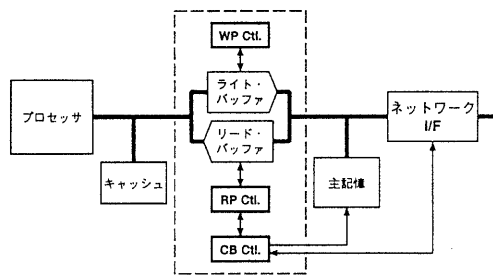


図 3: 処理ノードの構成

上のコンテキストを参照/更新しながら動作する。OS は、ユーザが対応する物理ページを既に獲得している場合に限り、ポート・ページへの写像を許可すると述べた。したがって、ポートへのアクセスが許可されれば、退避領域へのアクセスの正当性を改めて検査する必要はない。コンテキスト・キャッシュについては、次節で詳しく述べる。

図 3 に、対象計算機の処理ノードの構成 (主要なデータ・パス) を示す。図中には示していないが、通常のメモリ・アクセスのためのパイプス (・バッファ) がある。簡単のため、プロセッサが 1 次キャッシュやストア・バッファを内蔵しない場合について説明する。それらへの対応については 4.5 節でまとめる。

システムはデュアルポート・メモリで構成されたライト・バッファ/リード・バッファを備え、ブロックに対するアクセス・レイテンシを隠蔽する。

以下、各ユニットがどのように動作して前述の通信処理を実現するか説明する。同時に、記憶保護の実現について述べる。

(A) WP ユニット

WP に対してストアされるデータは、実際にはライト・バッファに書き込まれ、ライト・バッファ上でブ

ロックにまとめられる。

ライト・バッファの最後のバイトへのストアをトリガとして、保持するブロックを送出する。ブロックは、主記憶をバイパスして、ライト・バッファから直接ネットワーク I/F 部に送られる (図 3 を参照)。

ネットワーク I/F 部では、CB の仮想アドレスから CB のある処理ノードの物理アドレスへの変換が行われ、転送されるブロックには CB の仮想アドレスが付加される。大域的仮想アドレスの具体的な構成方法についてはここでは触れない。

(B) CB ユニットの構成

CB ユニットの構成は、ネットワークからブロックが届いた時、ブロックに付加された CB の仮想アドレスで CB ユニットのコンテキスト・キャッシュにアクセスし、主記憶上の CB の書き込み位置の物理アドレスを求める。したがって、OS が仮想アドレスと物理アドレスの写像を制御することによって、記憶保護が実現される。

書き込みに関する CB の管理情報の更新は CB ユニットのコントローラで行われる。

(C) RP ユニットの構成

RP ユニットのコントローラは、先頭と 2 番目のブロックがリード・バッファ上に存在するように制御する。RP に対するロードでは、実際にはリード・バッファから読み出しが行われる。

リード・バッファ上の 2 ブロック枠を用いて 2 面バッファリングを行い、主記憶上の CB に対するアクセス・レイテンシを隠蔽する。アクセスが先頭のブロックから次のブロックに移った時に、RP ユニットのコントローラは RP の実アドレスでコンテキスト・キャッシュにアクセスし、主記憶上の CB の読み出し位置を求める。CB は RP のコンテキスト退避領域と同一ページ内に割り付けるので、RP へのアクセスが許可されていれば、CB へのアクセスの正当性を改めて検証する必要はない。

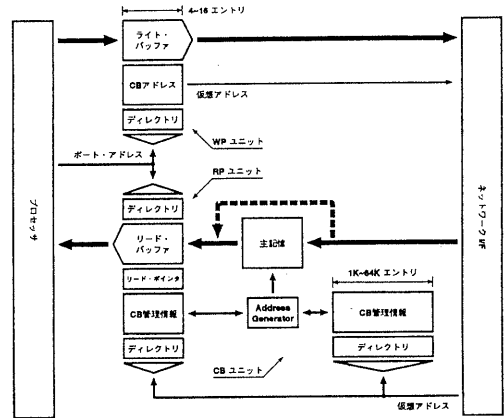
読み出しに関する CB の管理情報の更新は、RP ユニットのコントローラで行われる。

4.3 コンテキスト・キャッシュ

VQ システムの各ユニットは、それぞれのコンテキスト・キャッシュを持ち、コンテキストはそれぞれ独立にキャッシングされる。

前述のライト/リード・バッファは、実際には、コンテキスト・キャッシュの一部とみなすことができる。

図 4 にコンテキスト・キャッシュの構成と論理的なデータのながれを示す。また、各コンテキスト・キャッシュのエントリ数、連想度、アクセスするアドレスを、同図中の表に示す。



ユニット	エントリ数	連想度	アドレス
WP	4~16	完全	実
CB	1K~64K	1~4	仮想
RP	4~16	完全	実/仮想

図 4: コンテキスト・キャッシュ構成

以下、各ユニットの (A) コンテキスト と、(B) コンテキスト・キャッシュ について詳細に説明する。

4.3.1 WP ユニットの構成

(A) コンテキスト

ライト・バッファ上の生成途中のブロックは、WP のコンテキストとなる。

また、送信先の CB の仮想アドレスもコンテキスト・キャッシュにキャッシュし、送信処理を高速化する。

(B) コンテキスト・キャッシュ

WP ユニットのコンテキスト・キャッシュは、通常のプロセッサが持つストア・バッファなどとほぼ同じものと考えてよい。WP の実アドレスがディレクトリの保持するアドレスと比較され、ヒット/ミスが判断される。

4.3.2 CB ユニットの構成

(A) コンテキスト

以下の CB の管理情報がコンテキストとなる：

- サイズ (CB 末尾の、先頭アドレスからの変位)
- 読み出し/書き込み位置 (先頭アドレスからの変位)
- 現在のブロック数

(B) コンテキスト・キャッシュ

WP, RP へはプロセッサ上で走行中のスレッドのみがアクセスするのに対して、CB ユニットのアクセスは遠隔の多くのノードから非同期にブロックが届くので、1つ

の CB ユニットは性能上同時に多数のコンテキストをキャッシングする必要がある。したがって、CB ユニットのコンテキスト・キャッシュは、RAM チップを用いて構成する。

理由は後述するが、CB に対する書き込みはクリティカルではない。通常のメモリ・アクセスが VQ アクセスをバイパスできるとすれば、このキャッシュに対するアクセスのレイテンシを短縮する必要はない。スループットに関して、(このキャッシュのアクセス・サイクル数) = (主記憶のブロック・アクセス・サイクル数) ÷ (主記憶アクセスに占める VQ アクセスの割合) 程度であればよい。

複数サイクルをかけてアクセスすることで線幅を縮小できるので、このキャッシュは主記憶と同等のサイクル・タイムを持つ SRAM チップ 1 個程度を用いて構成できる。これは、ハードウェア量としては外部 2 次キャッシュのディレクトリなどと同程度である。

4.3.3 RP ユニット

(A) コンテキスト

WP ユニットとは異なり、リード・バッファ上のブロックは RP のコンテキストとする必要はない。なぜなら、リード・バッファ上のブロックは CB のコピーであり、CB のコンテキストが RP の状態を完全に指定するからである。

また、リード・ポートの 2 ブロック枠のうちどちらに先頭のブロックが写像されているかを示すビットが必要であるが、このビットは RP のコンテキストとなる。

(B) コンテキスト・キャッシュ

図 4 の、RP ユニットのコンテキスト・キャッシュのリードバッファから上側は、通常のキャッシュと同じものと考えてよい。

RP ユニットのコンテキスト・キャッシュは、CB ユニットとは独立に CB の管理情報をキャッシュし、新たに 2 番目となったブロックをリード・バッファにプリフェッチする時には、ここの情報を用いて高速に主記憶にアクセスする。2 面バッファリングを行っているとはいえ、リード・バッファへのプリフェッチは比較的高速に行う必要があるため、大容量で低速な CB ユニットのコンテキスト・キャッシュとの共用はしない。

このように、CB の管理情報は RP ユニットと CB ユニットの 2 箇所をキャッシュされるため、2 つのキャッシュの一貫性を保つ必要がある。RP ユニットのコンテキスト・キャッシュの下部分は仮想アドレスでアクセスできるようになっており、外部からブロックが到着した時には、CB ユニットと同時に RP ユニットのコンテキスト・キャッシュにもアクセスする。2 つのコンテキスト・キャッシュに同一のエントリが存在した時に

は、RP ユニット側を優先して用い、CB ユニット側のエントリは無効化する。

また、ある VQ に対してプロセッサがスピン・ウェイトしている最中にブロックが到着した場合には、RP ユニット側のコンテキスト・キャッシュがヒットするので、ブロックは主記憶をバイパスしてリード・バッファに送ることができる(図 3 参照)。前項で CB ユニットのコンテキスト・キャッシュに対するアクセスはクリティカルでないと述べたが、それはこの理由による。

CB の管理情報を 2 箇所をキャッシュすることには、プロセッサからの VQ アクセスとネットワークからの VQ アクセスを並列処理できるというメリットもある。

4.4 アンダーフロー/オーバーフロー

本節では、VQ のアンダーフロー/オーバーフローの検査の方法についてまとめる。以下で述べる方法の中には割り込みを用いるものがあるが、紙面の都合により割り込みからの回復の方法に関しては割愛する。

4.4.1 アンダーフローの検査

VQ のアンダーフローは、以下の 3 種の方法によって検査される：

1. **検査用コマンドを使用** コマンドにより、CB 中のブロックの数を検査する。
2. **使用しないパタンを応答** アンダーフローを起こした読み出しに対して、ヌル・ポインタ値 0、あるいは、NotANumber など、アプリケーションが使用しないパタンをブロックにつめて応答する。
例えば、次のようにするとよい：送信者はブロックの最終バイトへ 0 以外の値を書き込むことによって送信を行い、受信者はその位置が 0 以外であることを確認する。
3. **割り込みを発生する** アンダーフロー時に(しばらく待つ)割り込みを発生する。この方法を採れば、プログラムは VQ のアンダーフローを明示的に検査する必要がないので、VQ が空にならない可能性が高い場合に有効である。

2, 3 のうちのいずれの方法を使用するかは、RP のコンテキストとして指定する。1 はいつでも使用することができる。

4.4.2 オーバーフロー

プログラムは、原則として、以下のいずれかの方法によってオーバーフローの発生以前に送信側のスレッドと速度を調整するものとする：

- コマンドにより定期的にブロック数を検査する。
- ブロック数がユーザが設定した値になった時点で割り込みを発生する。

それでもオーバーフローを起こした時には、しばらく受信処理を中断した後に、プロセッサに割り込みをかけて知らせる。エンキューできるようになるまで受信処理を中断する方法も考えられるが、デッドロックを起こす可能性があるため採用しない。

4.5 市販のプロセッサへの対応

現在、市販の多くのプロセッサが1次キャッシュやストア・バッファを内蔵している。本節では、それらのプロセッサへの対応について述べる。

4.5.1 プロセッサの条件

VQシステムを実装するためには、要素プロセッサは以下の条件を揃えている必要がある：

1. プログラム・オーダでストアを観測する手段がある。
 2. 1次キャッシュ・ブロックが外部から無効化できる。
- 市販の多くの汎用プロセッサがこれらの条件を満たしている [5][6]。

4.5.2 対応の方法

プロセッサが上述の条件を満たしている場合、VQへのアクセスは内蔵1次キャッシュにまでキャッシュ可能であり、内蔵ストア・バッファによるストアのマージも可能である [7]。ただし、その場合、VQのブロックのサイズは、1次キャッシュのブロックのサイズと等しくする必要がある。

以下、WPとRPにおけるハードウェア、ソフトウェアの変更点について説明する。

(A) WP

WPでは、送信のトリガとなるブロックの最終バイトへのストアが、ブロックの他の部分へのストアを追い越さないようにする必要がある。ストア・バッファにおいてストアのオーダを変更するようなプロセッサでは、ストア・バリア命令の挿入が必要となる [6]。

(B) RP

RPでは、バッファリングを行う時の次のブロック枠への参照の移動が外部から観測できる必要がある。そのため、RPは3ブロック枠で構成し、3面バッファリングを行う。RPのコントローラは先頭のブロックと直前に先頭であったブロックのみが1次キャッシュ上に存在するように制御する。すると、2番目のブロックへのアクセスは必ず1次キャッシュ・ミスを起こすので、参照の移動を外部から観測できる。この1次キャッシュ・ミスに対して新たに先頭となったブロックを応答し、古いブロックを無効化する。

最適に動作した場合、VQに対する読み出しは、メッセージが常に2次キャッシュ上に存在する場合と同等の性能を示す。

5 おわりに

VQシステムは、仮想化により任意の数の仮想的通信機構を提供し、メッセージ送受信時のソフトウェア・オーバーヘッドを大幅に削減する。また、FIFOによる通信モデルとストリームのキャッシングにより、ハードウェアを小量化/高速化する。

しかし、これらの特長はプログラムに対して以下の制限を与えた結果として得られるものであり、これらの制限の妥当性を検証する必要がある：

1. **FIFOによる通信モデル** 送信者が一旦送信したブロック、受信者が一旦デキューしたブロックに対する参照は行えない。
2. **固定長メッセージ方式** 長いメッセージを不可分に転送すると、その間ハードウェアを横取りすることができず仮想化が困難になるので、可変長のメッセージは扱わない。

可変長メッセージを扱うためには、送信者が1人であるVQを別途生成するなどの処理が必要となる。

本機構の一部は、超並列計算機 JUMP-1 に実装される [7]。今後は、超並列計算機 JUMP-1 上で得られた実アプリケーション、および、シミュレーション [8] などで得られたデータを元に、上で示した制限について定量的評価を行う予定である。

謝辞

本研究の一部は、文部省科学研究費補助金（重点領域研究(1)「超並列ハードウェア・アーキテクチャの研究」、および、特別研究員奨励費「高速メッセージ通信機構に基づく並列計算機システム」）による。

参考文献

- [1] 清水俊幸, 堀江健志, 石畑宏明: 高速メッセージハンドリング機構—AP1000における実現—, in *JSP'92*, pp. 195-201 (1992).
- [2] Byrd, G. T. and Delagi, B. A.: StreamLine: Cache-Based Message Passing in Scalable Multiprocessor, in *ICPP '91*, pp. I-251-I-254 (1991).
- [3] 五島正裕, 森眞一郎, 富田眞治: プロセッサ間通信をサポートする On-Memory FIFO 機構, 情処研報 92-OS-56, pp. 111-118 (1992).
- [4] 平木敬他: 超並列計算機プロトタイプ JUMP-1 の構想, 情処研報 93-ARC-102 (1993).
- [5] Texus Instruments Inc.: *SuperSPARC User's Guide* (1992).
- [6] Digital Equipment Corp.: *DECchip™ 21064 Microprocessor Hardware Reference Manual* (1992).
- [7] 五島正裕他: 細粒度プロセッサ間通信をサポートする高機能キャッシュ・システム, 情処研報 93-ARC-101, pp. 121-128 (1993).
- [8] 五島正裕, 森眞一郎, 中島浩, 富田眞治: アーキテクチャ・シミュレータのための C++ クラス・ライブラリ, 情報処理学会 第 47 回 全国大会 論文集 (6) (1993).