

## ハイパクロスバ・ネットワークにおける転送性能向上のための手法とその評価

朴 泰 祐<sup>†</sup> 曾 根 猛<sup>†</sup> 三 島 健<sup>†</sup>  
板 倉 憲 一<sup>†</sup> 中 澤 喜 三 郎<sup>†</sup> 中 村 宏<sup>†</sup>

ハイパクロスバ・ネットワークは、数千台以上のプロセッサを有する超並列計算機におけるプロセッサ間結合ネットワークとして、多くの優れた性質を持っている。これまで同ネットワークは、ネットワーク中でのデッドロックを回避するため、固定ルーティングを前提とした、store & forward あるいは wormhole 方式のルーティングのみを対象として研究されてきた。本論文では、同ネットワークにおけるデッドロック・フリーな適応ルーティングの手法を提案し、それによる性能向上について計算機シミュレーションにより評価する。また、単純な wormhole 方式の転送に代わり virtual-cut-through 方式のメッセージ転送を導入することにより、さらにその転送性能が向上することも評価する。これらの改良により、同ネットワークの転送性能が大幅に向上することが確認された。

### Advanced Techniques for Performance Improvement of Hyper-Crossbar Network

TAISUKE BOKU,<sup>†</sup> TAKESHI SONE,<sup>†</sup> TAKESHI MISHIMA,<sup>†</sup> KEN'ICHI ITAKURA,<sup>†</sup>  
KISABURO NAKAZAWA<sup>†</sup> and HIROSHI NAKAMURA<sup>†</sup>

Hyper-Crossbar Network has several excellent features as an inter-PU network of massively parallel processing systems. In researches on this network so far, the routing algorithm is based on store-and-forward or wormhole method with fixed routing to avoid deadlock on the network. In this paper, we propose a deadlock-free adaptive routing technique for Hyper-Crossbar Network, and evaluate the performance improvement with it. We also apply a virtual-cut-through technique to Hyper-Crossbar Network instead of simple wormhole routing, and evaluate its effect. We confirmed the network throughput is much improved introducing these techniques on Hyper-Crossbar Network.

#### 1. はじめに

数千台以上の高性能 RISC プロセッサを PU (Processing Unit) として、それらをネットワークによって相互結合し、非常に高い性能を持つ超並列計算機を実現しようという試みが現在盛んに行われている。特に、各 PU がローカルなメモリを持ち、メッセージ・パッシングによって交信する、いわゆる分散メモリ型並列計算機では、超並列システムの実装は現実のものとなってきている。

このようなアーキテクチャにおいて、ネットワーク・トポロジとルーティング・アルゴリズムの選択は最も重要な問題の一つである。本論文では、hypercube

系ネットワークの性質をうまく利用し、ランダム転送に対して高い柔軟性を持つ、ハイパクロスバ・ネットワーク (以下 HXB と略)<sup>1)</sup> を取り上げる。HXB は一種の閉鎖型多段クロスバ網で、PU 間距離が非常に小さく、また通信チャンネル数も多いため、大規模なシステムにおける複雑な転送パターンにおいても比較的高い交信性能を保つことができる。特に、転送パターンがランダムな場合、他の典型的な超並列向きネットワークに比べて高いスループットを実現することが可能である<sup>2)</sup>。

HXB に対する過去の研究では、wormhole または store&forward 方式ルーティングにおける静的解析と、ランダム転送時の動的解析が行われてきた。これらの研究では静的な経路選択および単純な flit 制御に基づく wormhole ルーティングが対象となってきたが、これらのルーティング手法を改良することにより、さ

<sup>†</sup> 筑波大学 電子・情報工学系  
Institute of Information Sciences and Electronics, University of Tsukuba

らにいろいろなパターンの転送においてその性能が向上すると考えられる。本研究では、このような性能向上の手法として以下の2通りを考える。

- virtual-channel を用いた適応型ルーティングの導入
- virtual-cut-through の導入

これらの手法は基本的に独立なものであり、本論文ではこれらをそれぞれ独立に評価する。適応型ルーティングに関しては、ネットワーク中でのメッセージ間のデッドロック・フリーを保证するルーティング方式を提案し、それを HXB に導入した場合のスループットの向上について評価する。また、virtual-cut-through の導入に関しては、HXB が得意とする比較的長いメッセージ・パケットの転送を行う場合に、従来の virtual-cut-through 方式をさらに改良したバッファリング方式を提案し、それを HXB に応用した場合について評価する。いずれの場合もランダムな転送を行った場合を対象とし、計算機シミュレーションによって評価する。

## 2. ハイパクロスバ・ネットワーク

HXB は  $n$  次元の binary-hypercube の各次元方向の PU 数を 2 から  $m$  ( $>2$ ) に拡張したものである。さらに、この  $m$  は、各次元ごとに任意のサイズをとることができる (このサイズが全次元方向で等しいものは base- $m$   $n$ -cube として知られている<sup>2)</sup>)。第  $i$  次元方向の PU 数を  $m_i$  とすると、その次元方向の  $m_i$  台の PU は 1 つのクロスバ・スイッチ ( $m_i \times m_i$ ) によって結合される。各 PU はこのような各次元方向のクロスバ・スイッチに対するリンクを次元数  $n$  だけ持つ。このようにすると、総 PU 数は  $\prod_{i=1}^n m_i$  となる。また、クロスバ・スイッチで直接結合されていない PU 間では最大  $n$  ホップの転送を繰り返すことによってメッセージ転送を行うことができる。

図1に  $4 \times 4 \times 4$  構成の3次元 HXB を示す。図中、PU とクロスバ・スイッチを結合している四角は、wormhole ルーティングを行うためのルータで、エクステンジャ (EX) と呼ばれる。これは小規模なクロスバ・スイッチで、X、Y または Z 方向のクロスバ・スイッチ (XB) と PU の接続、あるいは各 XB 間の接続を行う。XB 間のメッセージ転送を自動的に行うことにより、中継 PU を介さない、wormhole ルーティングが実現される。以下、X、Y、Z 方向の XB をそれぞれ X-XB、Y-XB、Z-XB と略す。

過去の研究の結果、HXB は数千台規模の並列システムにおけるネットワークとして、各種転送パターン

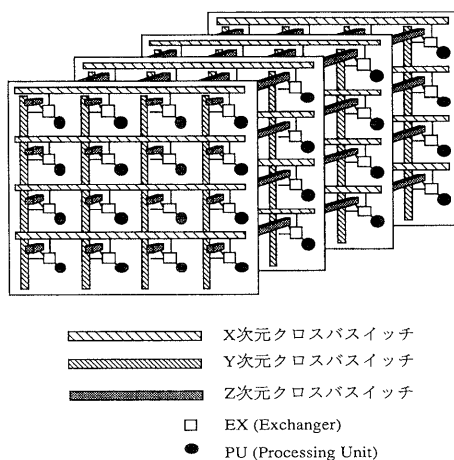


図1 3次元 HXB ( $4 \times 4 \times 4$ )  
Fig. 1 3 dimensional HXB ( $4 \times 4 \times 4$ ).

およびランダム転送における高い転送スループットを持つことが知られている<sup>3,4)</sup>。しかし、これらの評価においてはメッセージ・ルーティングが単純であり、以下のような問題点があった。

- 経路決定アルゴリズムが単純な固定ルーティング方式 (転送すべき次元の順番が常に一定) である。このため、メッセージが衝突した場合、他の空いている経路を利用せず、その場でブロックしてしまい、結果的にネットワークのスループットが低下する。
- メッセージ転送方式として単純な wormhole または store&forward 方式が用いられている。前者ではメッセージがブロックされた時、その後続部分がネットワークの他のチャンネルをブロックしてしまうため、余計な衝突が生じる。また後者では転送のホップ数が大きくなった場合、全体の転送レイテンシが大きくなってしまう。

本論文ではこれらの問題点を改善する2種類の手法を導入し、その評価を行う。

## 3. 性能向上の手法

前章で示した問題点の改善手法として以下の2種類の方法を導入する。

### 3.1 virtual-channel による適応ルーティング

単純な固定ルーティングに代わる経路決定法として適応ルーティング (adaptive routing) を実現する方法を提案する。従来の HXB において固定ルーティングが用いられていた最大の理由は、HXB において単純に適応ルーティングを用いて経路を動的に決定すると、複数のメッセージ間のチャンネルの連鎖に閉ループ

ができてしまい、そのまま wormhole ルーティングを行うと deadlock を生じるためである。そこで、ここでは virtual-channel の導入によって deadlock-free を保証する適応ルーティングを提案する。なお、ここでは迂回ルーティングは行わないものとする。すなわち、転送においては無駄な XB を通過せず、かつ各次元の XB を通過する回数はたかだか 1 回とする。

まず、 $n$  次元 HXB において、各 EX および XB の出力部に  $n$  枚のバッファを用意し、一つの物理チャンネル(リンク)を  $n$  本の virtual-channel として利用する。virtual-channel の割り当ては、そのメッセージがこれまでいくつの XB を経由したか(いくつの次元を通過したか)によって決定する。 $n$  次元 HXB ではメッセージは最大  $n$  個の XB を経由して転送されるため、 $n$  本の virtual-channel を用いれば、複数のメッセージのチャンネル連鎖中に閉ループができることはなくなる。例えば、3 次元 HXB において  $X \rightarrow Y \rightarrow Z$  という経路で転送されるメッセージ  $a$  と、 $Y \rightarrow Z \rightarrow X$  という経路で転送されるメッセージ  $b$  は、中継点の EX の座標が一致した場合、単純な物理チャンネル上では閉ループを構成する可能性がある。しかし、上述の virtual-channel を用いるとメッセージ  $a$  は  $X_1 \rightarrow Y_2 \rightarrow Z_3$  という経路で転送され、メッセージ  $b$  は  $Y_1 \rightarrow Z_2 \rightarrow X_3$  という経路で転送されるため、閉ループはなくなる( $X, Y, Z$  の添え字は virtual-channel 番号を表す)。よって deadlock-free な適応ルーティングが可能となる<sup>9)</sup>。以下に証明を示す。

ネットワーク中の全 virtual channel を  $C_{v,d,i,s}$  で表すとす。ここで、 $v$  は virtual 番号、 $d$  は次元番号(1, 2, ...,  $n$ )、 $i$  は次元  $d$  におけるディスティネーション EX の番号、 $s$  はそのチャンネルのソース EX の番号である( $i$  および  $s$  は  $n$  次元空間のアドレスであるが、各次元を順序づけし、線形化されているものとする)。メッセージは各 XB を通過する時、任意の次元(ただし過去通過していないもの)を選択でき、その際、 $v$  = 現在のステップ番号であるようなチャンネルを使用するものとする。このため  $v$  は単調増加し、結果として上述のチャンネル番号  $C_{v,d,i,s}$  は単調増加となるので、channel dependency graph にサイクルがなくなる。よってこのルーティングは deadlock-free である。(証明終わり)

以上の手法により、この適応ルーティングが deadlock を生じないことは保証されるが、HXB は間接網であるため、経路が busy であるかどうかをチェックする方法は直接網の場合より複雑になる。HXB における経路選択は EX 上で行われることになるが、たと

えここである次元方向へのチャンネルが free だったとしても、その先の XB における出力チャンネルが busy であれば、メッセージは結局ブロックされてしまう。

そこで、EX は XB の全出力チャンネルの状態を監視し、EX からそのチャンネルまでのリンクがすべて free の場合のみメッセージを転送するように制御する必要がある。さらに、実際の wormhole 転送ではメッセージの先頭が XB の出力チャンネルに到達するまでに遅延が生じるため、それまでの間はそのチャンネルを予約(reserve)するような制御も必要となる。HXB は PU 台数の増加に応じて次元数と各次元のサイズを調整できるため、各次元の XB のサイズはそれほど大きくない(例えば  $16 \times 16$  程度)。例えば、16 bit のデータ・パスに対し、1 bit の busy/free 状態を示す制御線を追加することは、それほど大きなコスト増にはつながらないと考えられる。

### 3.2 virtual-cut-through によるブロック率の低減

単純な wormhole ルーティングによる HXB の 2 つ目の問題点を解決するために virtual-cut-through 方式による転送を考える。これは wormhole と store&forward の両者の利点を併せ持つ転送方式で、基本的には wormhole 転送を行い、もしメッセージの先頭で衝突が起こりブロックされた場合は、そのノード(EX または XB)中にメッセージを溜め込み、少なくともメッセージの後続部分が他のメッセージと衝突するのを防ぎ、全体的なブロック率を低減させる手法である<sup>5)</sup>。

virtual-cut-through によって転送効率を上げるためには、メッセージは適当なサイズのパケットに分割される必要がある。これは、ノードへの store が速やかに行われるようにするためと、ノードに用意すべきバッファがあまり大きくならないようにするためである。メッセージは次のノードに進む際、そのノードの出力チャンネルが free であるか、またはそのノードのバッファが free であればそのノードに進むことができる。ここで、バッファが free であるというのは、そのノードのバッファに現在 store 中の他のメッセージが存在せず(一般に、複数の入力チャンネルから 1 つのバッファへの store が行われることは困難である)、かつそのメッセージを完全に収容できるだけの空きがある場合を指す。もし次のノードの出力チャンネルが free であれば、メッセージはそこに進み、さらに次の段階でその先のノードを窺うことができる。もし経路上のすべての出力チャンネルが free であれば、メッセージは結果的に wormhole 方式で転送される。

以上が基本的な virtual-cut-through 方式の転送であるが、ここに1つの問題点がある。それは、もしメッセージが先に進めない場合、メッセージはそのノードに store されるが、store している最中に出力チャンネルが free になった場合にどうするかである。そのような状況でも、そのメッセージの store が完了してから先に進むようにすると、制御方式は全体的に単純化されるが、store の回数が多くなり性能が低下する可能性がある。これに対し、たとえメッセージを store している最中でも、出力チャンネルが空けば直ちに先頭からの出力を開始するようにする方法が考えられる。この場合、バッファは先頭から次々に出力されると同時に、後続部分が前段のノードから次々送られてくるため制御が複雑になる。しかし、1つのバッファにデータが保存される時間が最低限で済むため、全体的な性能は向上すると考えられる。そこで、今回はこの両者の方法を比較し、評価する。前者の方法を static buffering、後者の方法を dynamic buffering と呼ぶことにする。

dynamic buffering 方式を採用することにより、全体的なスループットはかなり向上すると思われるが、バッファリングを工夫することにより、さらに大きな性能向上が期待できる。static buffering および dynamic buffering の両方式では、あるノードに存在しているメッセージは、経路上の次のノードのバッファにそのメッセージ自体を完全に格納し得るだけのバッファの空きがある場合のみ、メッセージを転送する。しかし、この制約を排除して、たとえ 1 flit 分であっても、とにかくバッファに空きがありさえすればメッセージの転送を開始するようにすると、メッセージがバッファ中に留まる時間がさらに短縮される。この方式では1つのメッセージのいろいろな部分が異なるノードのバッファに格納される可能性があり、制御が一層複雑になるが、メッセージは少しでも先に進める場合はどんどん進むため、dynamic buffering 方式よりさらにスループットが向上する。そこで、このバッファリング方式も前の2種類と同様に評価する。この方式を flexible buffering 方式と呼ぶ。

#### 4. 評価

適応ルーティングと virtual-cut-through の導入は基本的には独立な概念であり、両者を組み合わせることによってさらなる性能向上が期待される。しかし、ここでは基本的な性能評価を行うため、それぞれを独立に評価する。

評価は計算機シミュレーションを用いて行う。一般的なメッセージ衝突状態を想定するため、ここではラ

ンダム転送を基本とした。モデルとしては、システム中の全 PU がランダムに選んだ相手 PU に対して一定長のメッセージを転送し、転送が終了したら再び相手をランダムに選んで転送を繰り返す、という状況を想定した。PU に到着したメッセージは割り込み処理によって適宜受信されるものとする。通常、メッセージの送受信は DMA コントローラによって行われるので、ここでは送信処理と受信処理が同時にできるものとした。また、1 PU から同時に送信されるメッセージの数に関しては、適応ルーティングを用いた場合と、virtual-cut-through を用いた場合では若干異なる。これは、それぞれのルーティング方式の特性の違いに依存するもので、詳細は各方式の評価の部分でそれぞれ述べる。

システムの規模としては 4096 PU または 1024 PU の 3次元 HXB ( $16 \times 16 \times 16$  または  $8 \times 8 \times 16$ ) を対象とした。XB および EX の各チャンネルのバンド幅は 1 flit/clock に正規化し、メッセージ長は flit 単位で表すものとする。3次元であるため、PU 間距離の最大値は 3 であるが、実際には途中で最大 4 つの EX と最大 3 つの XB が経由されるため、PU → PU の最大距離は 8 となる。

評価は 1 PU 当りのメッセージ転送のスループットを、理想的な無衝突転送の場合を 1 として正規化して行った。ここで、理想的な無衝突転送とは、転送距離に依存するレイテンシと衝突による待ち時間がすべて 0 である場合を指す。この場合、 $L$  clock のシミュレーションにおいて 1 PU から転送されるメッセージの総量は  $L$  flit となる。実際には 10000 clock のシミュレーションを行ったので、スループットはこの間に PU が転送した総メッセージ量を 10000 で割ったものである。評価結果の各グラフは横軸にメッセージ長を取り、各メッセージ長におけるスループットを縦軸に取っている。

##### 4.1 適応ルーティングの評価

まず、適応ルーティングの転送性能の評価方法について述べる。ここではシステムの規模を 4096 PU とし、メッセージの転送方式は wormhole 方式とした(一般的に、HXB においては store&forward 方式より wormhole 方式の方が高いスループットを示す)。3.1 節に述べたように、全通信チャンネル (PU ↔ EX, EX ↔ XB 間) は各々 3 枚のバッファを用いて 3 重に仮想化され、この上で迂回なしの適応ルーティングをシミュレートする。また、ある EX から複数方向の XB へのチャンネルが free な場合は、X, Y, Z の優先順位で方向を決定することとする。

送信 PU の振舞いについては以下のような仮定を置いた。適応ルーティングでは、PU から送出されたメッセージは X, Y, Z のいずれの次元方向にも移動することができるため、ある EX において X 次元方向に転送されているメッセージがブロックされたときに、その EX を通過する後続のメッセージを virtual-channel により他の次元方向に転送することが可能となる。このことを考慮して、ここでは PU からのメッセージ送出に関して以下のように仮定した。各 PU はランダムに選んだ相手に一定長のメッセージを送る。1 メッセージの転送要求を DMA 回路に対して発行した後、一定時間の内部処理を行う(今回はこの長さを 50 クロックとする)。内部処理が終了して次のメッセージを送出する際、もし PU からネットワークへの出力チャネルが busy のときは、それが free になるまで内部処理を継続し、free になった時点で次のメッセージの転送を開始する。全 PU は以上の動作を繰り返す。

適応ルーティングにおける評価結果を図 2 に示す。図 2 の点線は固定ルーティング (FIX), 実線は適応ルーティング (ADP) の場合をそれぞれ示している。また FIX 1 と ADP 1 が、通常の HXB における、従来の固定ルーティングと今回提案した適応ルーティングの場合をそれぞれ表す。なお、固定ルーティングにおいては、経路決定は X, Y, Z の順に行われるものとした。

まず、固定ルーティング、適応ルーティングのどちらの場合もメッセージ長が長くなるにつれてスループットが増加していることがわかる。これは以下のような理由による。PU が 1 メッセージを送出するのに要する時間は、メッセージの先頭が相手 PU に到着するまでにかかるレイテンシと、それに続くメッセージ本体が転送され終わるまでの時間(今はバンド幅が正規化されているため、これはメッセージ長そのものに等しい)の和となる。また、レイテンシは PU 間の距離

とネットワーク中で他のメッセージと衝突した際のブロック時間の和となる。HXB における wormhole 転送では、ブロック時間はメッセージ長に対してほぼ線形に増加するため<sup>3)</sup>、メッセージ長が非常に短いときはブロック時間も小さくなり、1 メッセージの送出に要する時間が PU の内部処理時間 (50 クロック) より小さくなる。この場合、メッセージの送出間隔がある程度広がり、ネットワーク中の総メッセージ転送量は少なくなる。それに対して、メッセージ長が長くなるとブロック時間も大きくなり、1 メッセージの送出に要する時間が PU の内部処理時間より大きくなる。この結果、以前のメッセージを完全に送出し終わらないうちに次のメッセージの送出要求が生じ、ネットワーク中にメッセージが満たされ、総メッセージ転送量が増加し、結果的にスループットが向上する。

次に、固定ルーティングと適応ルーティングを比較する (FIX 1 と ADP 1 を参照)。メッセージ長が短いときはネットワークが空いているので、固定ルーティングと適応ルーティングで性能差は見られない。しかし、メッセージ長が長くなりネットワークが混雑してくると、適応ルーティングの場合は固定ルーティングに比べて最大で約 20% 程度性能が向上している。しかし、3 次元の場合、固定ルーティングでは X → Y, Y → Z, X → Z の 3 通りのチャンネルしか利用していないが、適応ルーティングでは 6 通りすべてのチャンネルを利用できる。このことを考えると、この程度の性能向上は十分とは言えない。

適応ルーティングにおいては、PU から送出されたメッセージは X, Y, Z のいずれの次元方向にも移動できることから、PU は最大で 3 つのメッセージを同時に送出することが可能となる。さらに、PU に到着するメッセージは X, Y, Z の各次元方向から同時に到着する可能性が高くなる。これらのことから、PU と EX 間のバンド幅がボトルネックになっていると予想される。そこで、PU と EX 間の物理的なバンド幅をネットワーク中のバンド幅の 2 倍、あるいは 3 倍に増加させることを考える。図 2 における凡例の数値 (1, 2, 3) は、ネットワーク中のバンド幅に対する、PU ↔ EX 間のバンド幅の比率を表す。

PU ↔ EX 間のバンド幅を増加させると、適応ルーティングにおいてはかなりの性能向上が見られるが、固定ルーティングではほとんど性能が向上しない。固定ルーティングにおいては、X → Y → Z の順でルーティングが行われるため、PU から送出されたメッセージはいつでも最初に X 方向に移動する。これにより、ほとんどのメッセージは EX における X-XB 方向へ

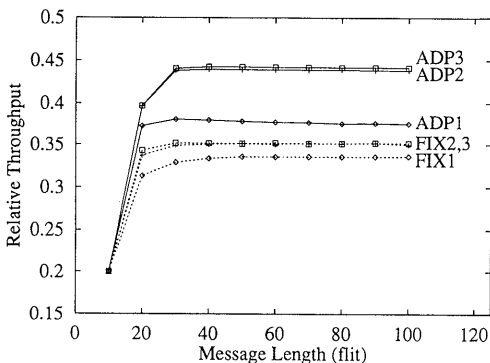


図 2 適応ルーティングによる性能向上

Fig. 2 Performance improvement with adaptive routing.

の出力において衝突し、そこでブロックされてしまう。逆に、受信 PU から見ると到着するメッセージはほとんどの場合 EX の Z-XB 方向から送られてくる。以上の理由により、固定ルーティングではネットワーク中の特定の次元方向の XB がボトルネックとなるため、PU → EX 間のバンド幅の増強はスループットに影響しないものと考えられる。

適応ルーティングにおいて、PU → EX 間のバンド幅をネットワーク中の 2 倍にしたとき、さらに性能は向上する (固定ルーティングに比べて約 30% の性能向上)。これは、先述したように、PU からネットワークに送出された複数のメッセージが異なる次元方向に移動することができるため、衝突する可能性が小さくなっているためである。また、受信 PU にとっても、メッセージはどの次元方向からも到着するため、PU と EX 間の 2 倍のバンド幅が有効に利用されている。それに対して、PU → EX 間のバンド幅をネットワーク中のその 3 倍に増強した場合には、それほど大きな性能向上は見られない。これはランダム転送で 1 PU に同時に 3 方向からメッセージが到着する可能性が比較的低いことと、ネットワークが非常に混雑しネットワーク中に溜っているメッセージが多くなっていることの 2 つによるものと考えられる。

以上、適応ルーティングにおいて PU → EX 間のバンド幅をネットワーク中のそれよりも増強することは大幅な転送性能の向上に寄与することがわかった。実装上は、EX における全入出力のうち PU に向けたリンクだけを増強すればよく、さらに、PU と EX は一対になっているため、実装上の物理的な距離が非常に近いことが予想されるため、その間のみ物理線の幅あるいは転送周波数を 2 倍に増強することは比較的容易であると考えられる。

#### 4.2 virtual-cut-through 転送方式の評価

virtual-cut-through 方式の転送性能の評価を、ランダム転送におけるスループットによって示す。ここでは、メッセージの経路決定は固定ルーティング方式で行われている。PU からのメッセージ送出のタイミングは、4.2 節に述べた適応ルーティングにおける評価と若干異なっている。

適応ルーティングの場合は、ソース PU の EX からのメッセージ送出の方向が最大 3 通り存在するが、固定ルーティングの場合は、X 方向への転送が必要な場合、常に最初に X 方向への送出が行われる。3 次元空間でのランダム転送では、ほとんどのメッセージは X 方向への転送が必要となる (1024 PU の場合で全体の約 88%)。このことから、たとえ PU から複数のメッセ

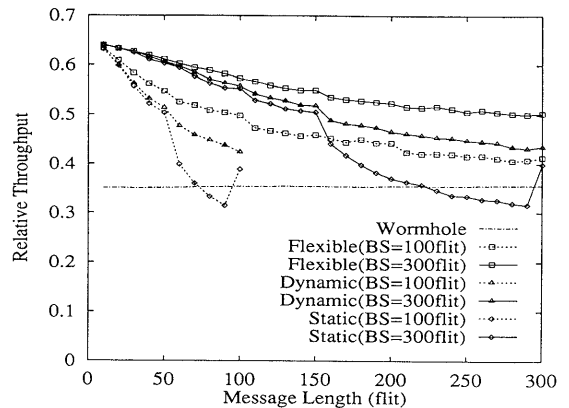


図3 virtual-cut-through による性能向上  
Fig. 3 Performance improvement with Virtual-Cut-Through.

ージを EX に送出できたとしても、結局 EX → X-XB へのパスがボトルネックになると考えられるため、適応ルーティングの評価で行ったような、1 PU からの複数のメッセージの同時送出や、PU → EX のバンド幅を増やすような工夫はスループットの向上に貢献しない。以上の理由により、virtual-cut-through 方式の評価においては、PU は 1 つのメッセージを送出した後、それが EX に完全に転送され終わるまでアイドル状態となり、PU → EX のパスが free になった段階で次のメッセージを送出するものとする。

シミュレーションによるスループットの評価を図 3 に示す。ここではシミュレータの実行時間の都合上、1024 PU の場合を評価する。各ノードにおけるバッファ容量として、300 flit の場合と 100 flit の場合の 2 通りを評価した。

図中、一点鎖線は通常の wormhole 転送におけるスループットを示している。実線および点線が、それぞれバッファ容量 (BS) が 300 flit の場合と 100 flit の場合の virtual-cut-through 方式を示している。また、プロットされた点の形が“□”のものが flexible、“△”のものが dynamic、“◇”のものが static の各 buffering 方式をそれぞれ示している。

全体的にメッセージ長が長くなるにつれスループットが低下することがわかる。予想されたように、基本的に flexible 方式が最も高いスループットを示し、次いで dynamic 方式、static 方式の順になっている。また、ノードにおけるバッファ容量が大きい程スループットが向上していることがわかる。

wormhole 方式のスループットがメッセージ長に依存せずほぼ一定である点が、適応ルーティングの評価において示した固定ルーティングの結果と異なってい

るが、これは前述したように、PUからのメッセージ送  
出のタイミングが異なっているためである。

以下に個々の buffering 方式の特徴を述べる。まず  
static buffering の場合であるが、メッセージ長がバッ  
ファサイズに比べ十分短い場合はある程度性能が向  
上するが、メッセージ長がバッファサイズの半分以上を越  
えたあたりから急激に下降する。これは、メッセージ  
が次のノードのバッファに入り切らなければそのノ  
ードに進まないという条件によるもので、バッファサイ  
ズの半分以上を越えた段階では、格納すべきバッファがな  
い状態でメッセージ長が大きくなっているため、ブロ  
ック率が高くなり性能が低下する。注目すべきは、メ  
ッセージ長がバッファサイズと等しい点では若干持ち  
直すという点である (V 字特性)。この現象については  
以下のように説明できる。

今、長さ  $M_a$  のメッセージ a が次のノードを窺って  
いる。次のノードのバッファには長さ  $M_b$  のメッセ  
ージ b が溜っているとす。バッファサイズを  $B$  とする  
( $B > M_a, M_b$ ) と、メッセージ b がさらに先のノ  
ードに転送される時、そのバッファ中の残りの容量が  
 $B - M_a$  となった時点で、メッセージ a はそのノ  
ードに進むことができる。しかし、メッセージ a の先頭が  
バッファに入った段階では、まだメッセージ b の残り  
がそのノードに残っているため、結局メッセージ a は  
そのノードで必ず待たされることになる。メッセ  
ージ長がバッファ容量の半分以上を越えると、この状況が発生  
し、それ以上メッセージ長が長くなるとメッセージの  
ブロック時間が増えるだけなので、全体的なスループ  
ットは低下する。

しかし、もし  $M_a = M_b = B$  だとすると、メッセ  
ージ a はメッセージ b が完全にノードから出ていくまでは  
そのノードに進めないが、そのノードに進んだ時はそ  
の前にメッセージ b は存在しないため、その先でメ  
ッセージ b がブロックしていなければ、メッセージ a は  
そのまま wormhole 的に進むことができる可能性が  
ある。このため、メッセージ長がバッファ容量に等し  
い場合は、上述のような無駄な store が減る可能性が  
高く、スループットが向上するものと思われる。

dynamic buffering 方式ではこの問題が改善されて  
いる。メッセージは、バッファに store 中であってもそ  
の先のノードが free になれば、store の完了を待つこ  
となしに転送が開始されるため、上述の状況になっ  
てもスループットは低下しない。static buffering の場合  
と同様に、メッセージ長がバッファサイズの半分以上を越  
えると性能は低下するが、上の理由により、低下はそ  
れほど小さくなく、また V 字型の特性も生じない。結

果として wormhole 方式の性能を常に上回り、メッセ  
ージ長がバッファ容量に等しい場合でも約 15% 程度  
の性能向上が達成できる。

また、flexible buffering 方式になると、スループ  
ットはさらに向上し、dynamic buffering 方式に比べ約  
10% から 20% の性能向上が見られる。これは、前方  
のバッファに少しでも空きがあればどんどんメッセ  
ージを送り込むという制御方式により、メッセージが先  
に進める可能性が高くなることでスループットが向上し  
ているためである。さらに注目すべき点は、flexible  
buffering 方式ではバッファ容量を越えるような長さ  
のメッセージの転送が可能になっている点である  
(BS=100 flit の場合を参照)。これは、次のノードにそ  
のメッセージが入り切るだけの余裕がなくてもメッセ  
ージを forward しているため、メッセージ長とバッ  
ファ容量の間に何の制約もなくなっているためである。  
このことから、flexible buffering 方式は、実装上の制  
約からバッファ容量がそれほど大きく取れない場合  
でも、HXB が本来得意とする比較的長いメッセージの  
転送において、その性能向上に十分貢献すると言える。

## 5. おわりに

本論文では、従来単純な固定ルーティングと  
store&forget あるいは wormhole 方式の転送のみを  
対象として評価されてきたハイパクロスバ・ネット  
ワークの転送性能に関し、virtual-channel を応用した動  
的な適応ルーティングと、virtual-cut-through による  
転送性能向上の手法を適用し、計算機シミュレー  
ションを通じてその性能を評価した。

まずメッセージ転送における経路選択方式に関して  
は、 $n$  次元のハイパクロスバ・ネットワークにおける最  
大転送ホップ数が  $n$  であることを利用し、各通信チャ  
ネルを virtual-channel を用いて  $n$  重に仮想多重化す  
ることにより、deadlock-free な適応ルーティングを  
実現する方法を提案した。このルーティング手法を用  
いることにより、基本的なランダム転送において、ネ  
ットワークの転送スループットが固定ルーティングの  
場合に比べて約 20% 向上することが計算機シミュ  
レーションにより確認された。さらに、送信 PU から同  
時に多数のメッセージが送出される場合、PU と EX  
間のバンド幅がボトルネックになることがわかり、そ  
の点を改善した結果、固定ルーティング方式に比べ、  
最大で約 30% の性能向上が確認された。

一方メッセージの転送方式に関しては、従来評価さ  
れてきた wormhole 方式に対し、古典的な virtual-  
cut-through 方式を適用した場合、特にメッセージ長

が短いランダム転送において大幅なスループットの向上が見られた。また、中継点のXBあるいはEXにおけるバッファへのflitの入出力を動的に処理するよう buffering 方式を改良することにより、スループットは wormhole 方式に比べ大幅に向上し、メッセージ長がバッファ容量に等しい場合で wormhole 方式に比べ約 40% の性能向上が達成できることがわかった。また、この場合メッセージ長がバッファ容量による制約を受けないため、任意の長さのメッセージを 1 パケットとして転送可能で、比較的長いメッセージの転送に向いているハイパクロスバ・ネットワークにおいて実用的であろう。

今回の評価ではこれら 2 つの方式を独立に扱ったが、これらの方式は互いに相反するものではない。今後はこれらの方式を融合した場合の性能評価と、ハードウェアコストを加味した場合のトレードオフについて評価していく予定である。また、適応ルーティングを実現する際、EX からその先の XB のチャンネルの状況を調べ、さらに必要に応じてそれを reserve するための具体的な機構についても考察していく予定である。

**謝辞** 本研究を進めるに当り貴重な御意見を頂いた、筑波大学西川博昭助教授に深く感謝します。なお、本研究の一部は文部省科学研究費(奨励(A)05780225)および創成的基礎研究費(05 NP 0601)の補助による。

### 参 考 文 献

- 1) 田中輝雄, 面田耕一郎: 高並列計算機による空気力学シミュレーションの構想, 第 8 回航空機計算空気力学シンポジウム論文集, pp. 99-108 (June 1990).
- 2) 鈴岡 節, 田邊 昇, 中村定雄, 藤田純一, 小柳 滋: 並列 AI マシン Prodigy の相互結合網の評価, 信学論 D, Vol. J 71-D, No. 8, pp. 1496-1501 (1988).
- 3) 朴 泰祐, 齊藤哲也, 板倉憲一, 中澤喜三郎, 中村宏: ハイパクロスバ・ネットワークの性能評価, 信学技報 CPSY 93-40 (1993).
- 4) 保田淑子, 藤井啓明, 田中輝雄, 稲神泰弘: ハイパクロスバネットワークの通信性能評価, 信学技報 CPSY 93-25 (1993).
- 5) Kermani, P. and Kleinrock, L.: Virtual Cut-Through: A New Computer Communication Switching Technique, *Computer Networks*, Vol. 3, No. 4, pp. 267-286 (1979).
- 6) Dally, W.J. and Seitz, C.L.: Deadlock-Free Message Routing in Multiprocessor Interconnection Networks, *IEEE Trans. Computer*, Vol. C-36, No. 5, pp. 547-553 (1987).



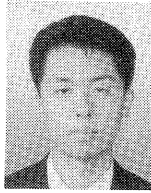
朴 泰祐 (正会員)

昭和 59 年慶應義塾大学工学部電気工学科卒業。平成 2 年同大学院理工学研究科電気工学専攻後期博士課程修了。工学博士。昭和 63 年慶應義塾大学理工学部物理学科助手。平成 4 年筑波大学電子・情報工学系講師。現在に至る。超並列計算機アーキテクチャおよび並列処理言語・アルゴリズムの研究に従事。電子情報通信学会会員。



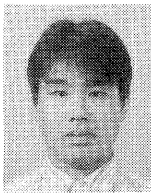
曾根 猛 (学生会員)

1971 年生。1994 年筑波大学第三学群情報学類卒業。現在同大学院工学研究科電子・情報工学専攻に在籍中。超並列計算機アーキテクチャ、並列処理に関する研究に従事。



三島 健 (学生会員)

平成 6 年筑波大学第三学群情報学類卒業。現在同大学院工学研究科電子・情報工学専攻修士課程に在籍中。並列計算機のネットワークの研究に従事。



板倉 憲一 (学生会員)

1969 年生。1993 年筑波大学第三学群情報学類卒業。1995 年同大学院工学研究科前期博士課程を修了。現在同研究科後期博士課程に在籍中。並列計算機アーキテクチャの研究に従事し、性能評価方法などに興味を持つ。



中澤喜三郎 (正会員)

昭和 30 年東京大学工学部応用物理卒業。昭和 35 年同大学院数物系博士課程応用物理修了。同年日立製作所入社。TAC, HITAC 5020, E/F, 8800/8700, M-200 H/280 H, 680 H, S-810 等, 超大型コンピュータ・スーパーコンピュータの開発に従事。平成元年より筑波大学電子・情報工学系教授。計算物理学センター向きの超並列処理システム CP-PACS の研究開発に従事。工学博士。IEEE, ACM 各会員。平成 5 年度本学会論文賞受賞。



**中村 宏 (正会員)**

昭和 38 年生。昭和 60 年東京大学工学部電子工学科卒業。平成 2 年同大学院工学系研究科電気工学専攻博士課程修了。工学博士。同年筑波大学電子・情報工学系助手。平成 3 年同講師，平成 7 年同助教授，現在に至る。計算機アーキテクチャ，並列処理，計算機の上位レベル設計支援に関する研究に従事。本学会平成 5 年度論文賞，本学会平成 6 年度山下記念研究賞各受賞。電子情報通信学会，IEEE 各会員。

---