

## VLAN を用いたマルチパス Ethernet における経路構築法

大塚智宏<sup>†</sup> 鯉渕道紘<sup>††</sup> 上楽明也<sup>†</sup>  
工藤知宏<sup>†††</sup> 天野英晴<sup>†</sup>

Ethernet を用いた PC クラスタにおいて、VLAN ルーティング法を用いることでホスト間に複数の経路を提供して経路分散を行うことができる。本稿では、VLAN ルーティング法を利用したトポロジにおいて、リンクレベルのフロー・コントロール、デッドロックフリールーティングの採用によりパケットの消失を抑えた効率の良いデータ転送ができる事を示す。そして、並列分散システムにおける代表的なトポロジである  $k$ -ary  $n$ -cube において経路分散されたデッドロックフリールーティングを実現する方法を提案し、それが  $n^{k-1}$  個の VLAN で実現できることを示す。

### Path Assignment Methods in VLAN-based Multi-path Ethernet

TOMOHIRO OTSUKA,<sup>†</sup> MICHIIRO KOIBUCHI,<sup>††</sup> AKIYA JOURAKU,<sup>†</sup>  
TOMOHIRO KUDOH<sup>†††</sup> and HIDEHARU AMANO<sup>†</sup>

In a PC cluster with Ethernet, well-distributed multiple paths among hosts can be obtained by using the VLAN-based routing method. In this paper, we show that efficient lossless data-transfer in a topology based on the VLAN-based routing method can be achieved using link-level flow control and deadlock-free routing. And we also show that deadlock-free well-distributed paths can be implemented with  $n^{k-1}$  VLANs on  $k$ -ary  $n$ -cube.

#### 1. はじめに

Ethernet はその高いコストパフォーマンスにより、最近では PC クラスタのインタコネクトとしても利用されている。初期のペオウルフ型クラスタと異なり、現在の Ethernet を用いたクラスタでは、TCP/IP を用いずにノード通信を行う専用ライブラリ<sup>1)</sup>や、システムエリアネットワーク (SAN) で用いられるゼロコピー通信またはワシコピー通信を利用できる。さらに、10Gigabit Ethernet (10GbE) の標準化など、CPU パワーの増加に伴ってリンクバンド幅も急速に大きくなっている。

一方で、Ethernet におけるスイッチ間のトポロジおよびルーティングに関する研究は、あまり行われていない。現在、Gigabit Ethernet (GbE) では 24 ポート、もしくはそれ以上のポート数を持つ大規模スイッチが主流になりつつあるが、以下の 3 つの理由によ

り、Ethernet のトポロジ、ルーティングは、今後も PC クラスタの性能を左右する要因の一つとなりうる。1) PC クラスタの大規模化がこのまま進んだ場合、やはり多数のスイッチによる接続が必要となる。2) 今後ホスト内の CPU のマルチコア化により、1 つのホスト内の各 CPU がそれぞれ異なるネットワークインターフェースを用いて通信するようになり、1 ホストあたり複数のスイッチポートが必要になる可能性がある。この場合、1 スイッチに接続可能なホスト数は限られる。3) 今後 10GbE のスイッチが市場に流通する場合、8 ポート程度の小規模なものから普及していくと考えられる。

Ethernet はループのない木構造トポロジを基本としており、トポロジ設定の自由度が低い。そこで、VLAN 技術を用いることで Ethernet においてもループ構造を含む様々なトポロジ、ルーティングを利用できるようにする VLAN ルーティング法<sup>2)</sup> が提案されている。

本稿では、まず、VLAN ルーティング法による経路設定をする場合に、パケットの循環依存によりパケットの消失が頻繁に発生し、バンド幅が大幅に低下することを示す。そして、この問題を解決するために、Ethernet における VLAN を用いた経路構築法としてデッドロックフリールーティングを使用することを提

<sup>†</sup> 慶應義塾大学理工学部

Faculty of Science and Technology, Keio University

<sup>††</sup> 国立情報学研究所

National Institute of Informatics

<sup>†††</sup> 産業技術総合研究所

National Institute of Advanced Industrial Science and Technology

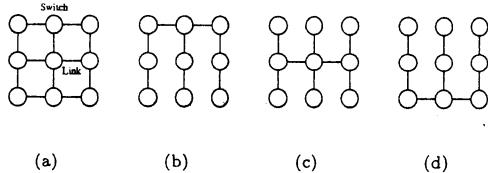


図 1  $3 \times 3$  2 次元メッシュ上の VLAN トポロジの例

案する。

次に、並列分散システムにおける代表的なトポロジである  $k$ -ary  $n$ -cubeにおいて、経路分散されたデッドロックフリールーティング<sup>3)</sup>が  $n^{k-1}$  個の VLAN で実現できることを示す。 $k$ -ary  $n$ -cube は、 $(k, n)$  の値を変化させることで対象とするスイッチの次数に合わせた構成をとることができ、多くの並列分散システムで採用されている基本的なトポロジであることから、VLAN ルーティング法を用いたクラスタにおいても有用なトポロジであると考えられる。

以降、第 2 節で VLAN ルーティング法を用いた場合のデッドロックの問題について議論し、実機による簡単な実験の結果を示す。次に、第 3 節では  $k$ -ary  $n$ -cubeにおいてデッドロックフリールーティングを実現する VLAN 集合の構築法について説明する。最後に、第 4 節にてまとめを述べる。

## 2. Ethernet におけるデッドロック問題

本節では、VLAN ルーティング法を用いた Ethernet 上の経路設定法を説明し、そこで発生する可能性のあるパケット間のデッドロックの問題を明らかにする。そして、その解決策としてデッドロックフリールーティングを使用することを提案し、その効果を示す。

### 2.1 VLAN ルーティング法

Ethernet では、基本的に木構造トポロジ上でパケットを転送するが、VLAN 技術を利用することにより、柔軟な経路設定が可能な VLAN ルーティング法<sup>2)</sup>が提案されている。VLAN ルーティング法では、異なる木構造トポロジの VLAN を組み合わせることにより、ホスト間に複数の物理経路を設定することができる。例えば、図 1 は、 $3 \times 3$  2 次元メッシュ(a)とその VLAN トポロジの例を 3 つ示したものである。図に示される通り、送信元ホストにおいて (b), (c), (d) から適切な VLAN を選択することにより、すべての宛先ホストに対して最短経路をとることができる。

### 2.2 デッドロックの発生と回避

前節で示した通り、Ethernet において VLAN を用

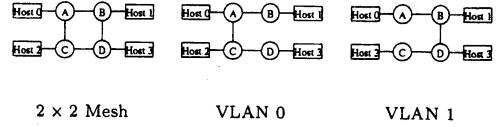


図 2  $2 \times 2$  2 次元メッシュ上の VLAN トポロジの例

いることにより、ループ構造を含む様々なトポロジを採用することができる。しかし、ループを含むトポロジでは、並列計算機の結合網や SAN と同様に Ethernetにおいてもパケット間でデッドロックが発生する可能性がある。

例えば、図 2において、ホスト 0 から 3、および 1 から 2 へのパケットが VLAN 1、2 から 1、3 から 0 へのパケットが VLAN 0 をそれぞれ用いるとする。一般的な Ethernet スイッチは、VLAN ID に対応した仮想チャネルなどは持たないため、この場合、4つの転送経路間にデッドロックを発生させる要因となる(物理)チャネル循環依存が存在することになる。

通常、Ethernet のフレーム転送においてスイッチやホスト上の NIC のバッファが一杯だった場合、フレームは単に破棄される。このとき、上位層のプロトコルによる再送制御などで実効バンド幅が低下する可能性はあるが、デッドロックは発生しない。しかし、IEEE 802.3x で規定されている PAUSE フレームを用いたリンクレベルのフローコントロールを用いた場合、フレーム破棄が抑制されるため、実際にデッドロックが発生してしまう可能性がある。

このデッドロックの問題は、チャネル依存グラフ (Channel Dependency Graph: CDG) においてチャネル間の循環依存が発生しないことを保証するデッドロックフリールーティングを用いることで解決することができる。Ethernet のように仮想チャネルを持たないネットワーク上すべての循環依存を除去するデッドロックフリールーティングは、並列計算機や SAN 向けに様々なものが提案されている<sup>4)</sup>。例えば、Dimension-order Routing<sup>3)</sup>では、図 2において、ホスト 0 から 3、2 から 1 へのパケットが VLAN 1、1 から 2、3 から 0 へのパケットが VLAN 0 を用いることによりデッドロックを回避する。

### 2.3 デッドロックフリールーティングの性能測定

本節では、VLAN ルーティング法を用いたループを含むトポロジにおいて、デッドロックフリールーティングの性能評価を行った結果を示す。

#### 2.3.1 評価環境

ホスト PC の構成は以下の通りである。

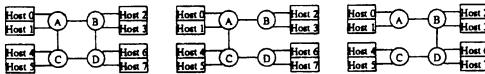


図 3 評価に用いたトポロジと VLAN の構成

- プロセッサ: Pentium4 2.4C (2.4GHz)
- マザーボード: Intel D865GLC
- メモリ: 512MB DDR400
- NIC: Intel 82547EI (オンボード, CSA 接続)
- OS: Red Hat Linux 9 (kernel 2.4.21)
- NIC ドライバ: Intel e1000 5.2.39

また, Ethernet スイッチには, DELL PowerConnect 5224 (Gigabit Ethernet × 24 ポート, ノンブロッキング) を用いた。

スイッチのトポロジには, 図 3 に示す 4 スイッチを用いたリングトポロジを用いた。スイッチ間, スイッチ-ホスト間のリンクはそれぞれ 1 本ずつとしている。

測定には, Iperf 2.0.2<sup>5)</sup> の UDP 転送機能を用いた。UDP を用いたのは, パケット消失の要因を限定するために, End-to-End フローコントロールを用いないようにするためである。ここで, UDP データグラムの転送レートは 1000Mbps, バッファサイズは 128KB とした。

まず, デッドロックフリーの通信パターンとして, 次の 4 通りのパケット転送を同時に行った。

- VLAN #3 を用いてホスト 1 から 6 へ
- VLAN #2 を用いてホスト 3 から 4 へ
- VLAN #3 を用いてホスト 5 から 2 へ
- VLAN #2 を用いてホスト 7 から 0 へ

次に, デッドロックが発生しうる通信パターンとして, 次の 4 通りのパケット転送を同時に行った。

- VLAN #3 を用いてホスト 1 から 6 へ
- VLAN #3 を用いてホスト 3 から 4 へ
- VLAN #2 を用いてホスト 5 から 2 へ
- VLAN #2 を用いてホスト 7 から 0 へ

これら 2 種類の通信パターンでは, 1 つの単方向チャネルを通過する経路数は同じであるため, 経路分散の度合は等しい。さらに, パケットのホップ数は全てのソース-デスティネーション対において同じである。

### 2.3.2 結 果

まず, 予備評価として, 単純な直線トポロジにおいて, 802.3x の PAUSE フレームを用いたリンクレベルのフローコントロールを用いる場合と用いない場

合それぞれのバンド幅およびパケット消失率の測定を行った。

表 1 直線トポロジにおけるバンド幅とパケット消失率

	FC なし		FC あり	
	BW [Mbps]	消失率 [%]	BW [Mbps]	消失率 [%]
1-SW-1	820	0.025	821	0.033
2-SW-2	1642	0.033	1640	0.092
4-SW-4	3283	0.063	3275	0.050
1-SW-SW-1	821	0.021	821	0.016
2-SW-SW-2	958	41.6	957	0
4-SW-SW-4	957	70.8	960	0
8-SW-SW-8	958	85.2	959	0

表 1 において, 「FC あり/なし」はリンクレベルのフローコントロールのあり/なしを, 「2-SW-SW-2」は 2 つのスイッチを介して 2 組のソース-デスティネーション対が接続されている場合を表している。スイッチ間は 1 本のリンクで接続されているため, 各ソース-デスティネーション対間のパケットは, スイッチ間のリンクにおいて経路が重なる。また, ここでの bandwidth (BW) は, 全通信対における bandwidth の総和である。

表 1 の結果より, フローコントロールなしの場合とありの場合とで bandwidth はほとんど変わらないが, フローコントロールなしの場合, 転送経路が重なったとき (2-SW-SW-2, 4-SW-SW-4, 8-SW-SW-8) に大量のパケット消失が発生している。一方, フローコントロールを用いた場合, これらのパケット消失を完全に防ぐことができている。

これらの結果から, リンクレベルのフローコントロールのみを用いても, 直線トポロジにおいてパケット消失を抑えることができる事が分かる。これにより, Ethernet において, TCP のような上位プロトコルによる End-to-End のフローコントロールを前提としないプロトコルを使える可能性が示されたと言える。ただし, 1-SW-1 から 1-SW-SW-1 までの 4 つの場合については, パケットの消失を完全に 0 にすることはできない。これは NIC からスイッチに対するフローコントロールがうまくいっていないことが原因と考えられるが, 詳細については現在調査中である。

次に, デッドロックフリールーティングの評価結果を表 2 に示す。

表 2 において, DF および DL は前述のデッドロックフリーの通信パターンとデッドロックが発生しうる通信パターンをそれぞれ表す。ここで, PAUSE フレームによるリンクレベルのフローコントロールを使用し

表 2 デッドロックフリーとデッドロックが発生しうる転送パターンにおけるバンド幅とパケット消失率

	BW[Mbps]	消失率[%]
DF/FC なし	816.5	0.100
DL/FC なし	412.4	49.5
DF/FC あり	814.8	0.249
DL/FC あり	-	-

た場合、デッドロックが発生しうるパターンでは測定開始後すぐにネットワークの状態が不安定になり、プログラムが終了しなかった。プログラムの強制終了後もネットワークは回復せず、スイッチをリセットしない限りノード間の通信が全くできない状態が続いた。これは、フローコントロールによりパケットの破棄が抑制され、実際に循環依存によるデッドロックが発生して回復不可能になってしまったためと考えられる。

表 2 の結果から、フローコントロールのあり／なしにかかわらず、デッドロックフリーの経路群が高い性能を示している。フローコントロールなしの場合、デッドロックが発生しうる経路群では約半分のパケットが消失しており、バンド幅が大幅に低下している。

並列計算機の結合網や SAN では、スイッチング技術としてカットスルーもしくはワームホールルーティングを用いている。一方、Ethernet は Store-and-Forward 方式を用いているため、並列計算機の結合網や SAN に比べてデッドロックは発生しにくいと考えられる。しかし、この結果からは、VLAN ルーティング法によりループを含む経路を導入し、リンクレベルのフロー コントロールを用いた場合には、Ethernet でもデッドロックが現実問題として十分発生しうることが分かる。また、フローコントロールを用いない場合でも、循環依存のある経路群ではパケットの消失により大幅に性能が低下しており、VLAN ルーティング法を用いた経路構築においてデッドロックフリールーティングは非常に効果的であるといえる。

### 3. $k$ -ary $n$ -cube におけるデッドロックフリールーティング

前節で述べた通り、Ethernet 上の VLAN ルーティング法によるトポロジにおいても、デッドロックフリールーティングが有効であることが分かった。そこで本節では、並列分散システムにおける代表的なトポロジである  $k$ -ary  $n$ -cube において、VLAN ルーティング法を用いたデッドロックフリールーティングを実現する方法を示し、その経路集合が必要とする VLAN 数を明らかにする。

IEEE 802.1Q の規定では、VLAN タグによって

4,094 ( $2^{12} - 2$ ) 個の VLAN を識別することができるが、商用のコストパフォーマンスの高い Ethernet スイッチはそれほど多くの VLAN 数をサポートしていない場合が多い。よって、必要となる VLAN 数は、大規模クラスタの構築を制限する要因となりうるため、正確に見積もる必要がある。

$k$ -ary  $n$ -cube は、 $(k, n)$  の値を変化させることで対象とするスイッチの次数に合わせた構成(ハイパー キューブから 1 次元リングトポロジまで)をとることができ、多くの並列分散システムで採用されている基本的なトポロジである。様々なルーティングアルゴリズム、マルチキャストアルゴリズムが  $k$ -ary  $n$ -cube 用に開発されており、RDT<sup>6)</sup> のように 2 次元格子トポロジを基に階層構造を作ることにより直径を小さくするような応用もできる。 $k$ -ary  $n$ -cube には、大きく分けて Wrap-around チャネルが存在するトラスと、存在しないメッシュの 2 種類があるが、本稿では、簡略化のため、Wrap-around チャネルのないメッシュに焦点をあてる。

$k$ -ary  $n$ -cube におけるデッドロックフリールーティングとしては、Dimension-order Routing (DOR)<sup>3)</sup> が挙げられる。Dimension-order Routing は、 $k$ -ary  $n$ -cube において最短経路を保証し、かつ全対全通信において最も経路が分散するルーティング手法の一つである。Dimension-order Routing では、パケットはもっとも低い次元から順に必要なホップ数だけ移動することでルーティングされる。

以下の各節では、 $M$  次元メッシュにおける Dimension-order Routing に従う最短経路の構築法を説明し、必要に応じて 2 次元メッシュの場合の例を示す。

#### 3.1 準 備

図 4 は、 $4 \times 4 \times 2$  次元メッシュとその VLAN トポロジの例を示したものである。図に示されるように、1 つのメッシュ上の VLAN トポロジ(木構造)には多くの種類がある。しかし、(d) のような規則性に乏しいトポロジを用いて最短経路集合を構築するのは難しい。そのため、提案手法では、(b) や (c) のような規則的なトポロジのみを用いる。

図 4において、各頂点および辺はそれぞれ Ethernet スイッチとリンクを表す。通常、各スイッチには何台かのホストが接続されるが、ここでは省略されている。

**定義 1** ( $M$  次元メッシュ) 各頂点(スイッチ)に対し、 $M$  次元座標  $(x_0, x_1, x_2, \dots, x_{M-1})$  を割り当てる。ここで、 $0 \leq x_0, x_1, x_2, \dots, x_{M-1} < N$  である。頂点  $(x_0, x_1, x_2, \dots, x_{M-1})$  を最大  $2M$  個の頂点  $\{(x_0, x_1, x_2, \dots, x_{i-1}, x_i \pm 1, x_{i+1}, \dots, x_{M-1}) \mid 0 \leq$

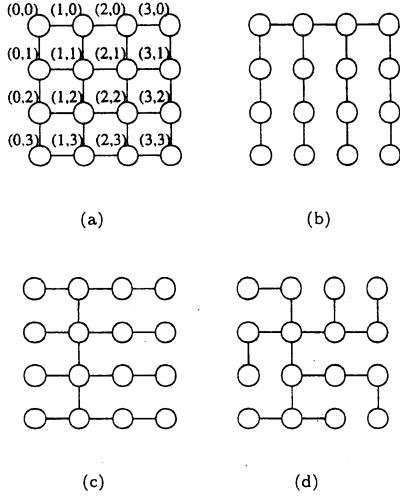


図 4  $4 \times 4$  2 次元メッシュとその VLAN トポロジの例

$i < M, x_{i-1} \geq 0, x_{i+1} < N\}$  とそれぞれ接続することにより,  $N^M$   $M$  次元メッシュが形成される.  $\square$

この  $M$  次元メッシュは, 一般には  $N$ -ary  $M$ -cube として定義されるものである. 例えば  $4 \times 4$  2 次元メッシュの場合, 図 4(a) のように座標が割り当てられる.

図 4 で, (b), (c), (d) の各 VLAN トポロジは物理ネットワーク (a) のスパンニングツリーになっており,  $N^2$  個のスイッチと  $N^2 - 1$  本のリンクから構成される. ここで, 各 VLAN トポロジを識別するため, 以下の記法を導入する.

**定義 2 (メッシュ上の線形接続)**  $N^M$   $M$  次元メッシュにおいて, 線形接続  $l(x_0, x_1, x_2, \dots, x_{i-1}, -, x_{i+1}, \dots, x_{M-1})$  は,  $N$  個の頂点  $\{(x_0, x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_{M-1}) \mid 0 \leq x_i < N\}$  が  $N-1$  本のリンクで接続された第  $i$  軸に平行な線形構造である. ここで, 頂点  $(x_0, x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_{M-1})$  は 1 個または 2 個の頂点  $\{(x_0, x_1, x_2, \dots, x_{i-1}, x_i \pm 1, x_{i+1}, \dots, x_{M-1}) \mid x_{i-1} \geq 0, x_{i+1} < N\}$  とそれぞれ接続される.  $\square$

以上の定義を用いて,  $M$  次元メッシュにおける VLAN トポロジは以下のように定義される.

**定義 3 (メッシュ上の VLAN トポロジ)**  $N^M$   $M$  次元メッシュにおける VLAN トポロジ  $VL(x_0, x_1, x_2, \dots, x_{i_0-1}, -, x_{i_0+1}, \dots, x_{M-1} \mid (i_0, i_1, i_2, \dots, i_{M-1}))$  (ただし,  $0 \leq i_0, i_1, i_2, \dots, i_{M-1} < M, j \neq$

$k$  のとき  $i_j \neq i_k$ ) は, 全頂点 (スイッチ) と, 以下の線形接続の集合からなる.

$$\begin{aligned} & \{l(x_0, x_1, x_2, \dots, x_{i_0-1}, -, x_{i_0+1}, \dots, x_{M-1})\} \\ & \cup \{l(x_0, x_1, x_2, \dots, x_{i_1-1}, -, x_{i_1+1}, \dots, x_{M-1}) \\ & \quad | 0 \leq x_{i_0} < N\} \\ & \cup \{l(x_0, x_1, x_2, \dots, x_{i_2-1}, -, x_{i_2+1}, \dots, x_{M-1}) \\ & \quad | 0 \leq x_{i_0}, x_{i_1} < N\} \\ & \vdots \\ & \cup \{l(x_0, x_1, x_2, \dots, x_{i_{M-1}-1}, -, x_{i_{M-1}+1}, \\ & \quad \dots, x_{M-1}) | 0 \leq x_{i_0}, x_{i_1}, x_{i_2}, \dots, x_{i_{M-2}} < N\} \end{aligned}$$

$\square$

上の定義において, 各 VLAN はそれぞれ第  $i_m$  軸に平行な線形接続を  $N^m$  本, 計  $(N^M - 1) / (N - 1)$  本の線形接続を持つ.

2 次元メッシュ上の例として, 図 4 のトポロジ (b) は  $VL(-, 0 \mid (0, 1))$  で表され, 第 0 軸 ( $x$  軸) に平行な線形接続  $l(-, 0)$  と, 第 1 軸 ( $y$  軸) に平行な線形接続  $l(0, -)$ ,  $l(1, -)$ ,  $l(2, -)$ ,  $l(3, -)$  を持つ. また, トポロジ (c) は  $VL(1, - \mid (1, 0))$  で表され, 第 1 軸に平行な線形接続  $l(1, -)$  と, 第 0 軸に平行な線形接続  $l(-, 0)$ ,  $l(-, 1)$ ,  $l(-, 2)$ ,  $l(-, 3)$  を持つ.

### 3.2 DOR の経路を構築するための VLAN 集合

本節では, メッシュ上で Dimension-order Routing (DOR) に従う最短経路を構築するための VLAN 集合を示し, 必要となる VLAN 数を示す.

**定義 4 (メッシュ上の DOR VLAN 集合)**  $N^M$   $M$  次元メッシュ上の DOR VLAN 集合は, 以下の  $N^{M-1}$  個の VLAN から構成される.

$$\begin{aligned} & \{VL(-, x_1, x_2, \dots, x_{M-1} \mid A) \\ & \quad | 0 \leq x_1, x_2, \dots, x_{M-1} < N\} \end{aligned}$$

ここで,  $A = (0, 1, 2, \dots, M-1)$  である.  $\square$

スイッチ  $(x_0, x_1, x_2, \dots, x_{M-1})$  からの経路は,  $VL(-, x_1, x_2, \dots, x_{M-1} \mid A)$  を用いることによってすべて DOR に従う最短経路となる. 図 5 に  $4 \times 4$  2 次元メッシュにおける DOR VLAN 集合の例を示す. 図に示された通り,  $VL(-, 0 \mid A)$ ,  $VL(-, 1 \mid A)$ ,  $VL(-, 2 \mid A)$ ,  $VL(-, 3 \mid A)$  の 4 個の VLAN が使用される. ここで,  $A = (0, 1)$  である.

**定理 1**  $N^M$   $M$  次元メッシュ上の DOR VLAN 集合は, Dimension-order Routing (DOR) に従う最短経路集合を提供する.

**証明** VLAN  $VL(-, x_1, x_2, \dots, x_{M-1} \mid A)$  は, 第  $i$  軸 ( $0 \leq i < M$ ) に平行なそれぞれ  $N^i$  本の線形

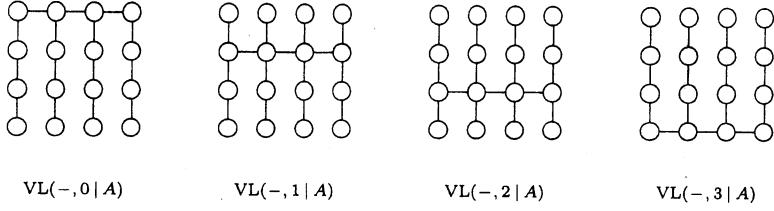


図 5  $4 \times 4$  2 次元メッシュ上の DOR VLAN 集合

接続  $\{l(x_0, x_1, x_2, \dots, x_{i-1}, -, x_{i+1}, \dots, x_{M-1}) \mid 0 \leq x_0, x_1, x_2, \dots, x_{i-1} < N\}$  で構成される。ゆえに、スイッチ  $(x_0, x_1, x_2, \dots, x_{M-1})$  からの経路は、 $VL(-, x_1, x_2, \dots, x_{M-1} \mid A)$  を使うことにより、第 0 軸、第 1 軸、第 2 軸、…、第  $M-1$  軸の順 (DOR と同じ) にルーティングされ、最短経路となる。DOR VLAN 集合は  $N^{M-1}$  個の VLAN  $\{VL(-, x_1, x_2, \dots, x_{M-1} \mid A) \mid 0 \leq x_1, x_2, \dots, x_{M-1} < N\}$  で構成されるため、すべてのスイッチからの経路が DOR に従う最短経路となる。  
□

通常の Ethernet スイッチは、並列計算機の相互結合網における仮想チャネルのように VLAN ID 毎のバッファは持たない。そのため、経路が異なる VLAN 間に分散しているかどうかはネットワーク全体の性能にはほとんど影響を与えない。

#### 4. ま と め

本稿では、まず、VLAN ルーティング法を利用した複数パスを持つ Ethernetにおいて、リンクレベルのフローコントロール、およびデッドロックフリールーティングの採用により、パケットの消失を抑えた効率の良いデータ転送ができるることを示した。性能測定の結果、リンクレベルのフローコントロールのみを用いることで、UDP での転送においても、パケットを消失せずに効率の良いデータ転送ができるることを確認した。また、4 スイッチのリングトポロジにおいて、循環依存のある経路群がパケット消失により大幅なバンド幅低下を招くのに対し、デッドロックフリーの経路群は高いバンド幅を維持できることが分かった。

次に、並列分散システムにおける代表的なトポロジである  $k$ -ary  $n$ -cube において、VLAN ルーティング法を利用したデッドロックフリールーティングを実現する方法を提案し、またそれが  $n^{k-1}$  個の VLAN で実現できることを示した。

今後の課題としては、フローコントロールやデッドロックに関するより詳しい解析、 $k$ -ary  $n$ -cube やその他の様々なトポロジにおけるデッドロックフリールーティングの実装および評価が挙げられる。

#### 謝 辞

クラスタおよびEthernetによるトポロジ構築を手伝って下さった産業技術総合研究所 清水敏行氏、岡崎史裕氏に感謝いたします。また、VLAN ルーティング法に関する相談に乗って下さった元産業技術総合研究所 手塚宏史氏に感謝いたします。

#### 参 考 文 献

- 1) T.Takahashi, S.Sumimoto, A.Hori, H.Harada and Y.Ishikawa: PM2: High Performance Communication Middleware for Heterogeneous Network Environment, *SC2000*, pp. 52–53 (2000).
- 2) 工藤知宏、松田元彦、手塚宏史、児玉祐悦、建部修見、関口智嗣: VLAN を用いた複数バスを持つクラスタ向き L2 Ethernet ネットワーク、情報処理学会論文誌コンピューティングシステム、Vol. 45, No. SIG 6(ACS 6), pp. 35–43 (2004).
- 3) Daily, W. J. and Seitz, C. L.: Deadlock-Free Message Routing in Multiprocessor Interconnection Networks, *IEEE Transactions on Computers*, Vol. 36, No. 5, pp. 547–553 (1987).
- 4) J.Duato, S.Yalamanchili and L.Ni: *Interconnection Networks: an engineering approach*, Morgan Kaufmann (2002).
- 5) Iperf - The TCP/UDP Bandwidth Measurement Tool, <http://dast.nlanr.net/Projects/Iperf/>.
- 6) Y. Yang, A. Funahashi, A. Jouraku, H. Nishi, H. Amano and T. Sueyoshi: Recursive Diagonal Torus: an interconnection network for massively parallel computers, *IEEE Transaction on Parallel and Distributed Systems*, Vol. 12, No. 7, pp. 701–715 (2001).