

相互結合網 RSOT のルーティング方式

2 L-6

秋山知之 田中英彦
東京大学工学系研究科

1はじめに

相互結合網 RSOT(Recursive Orthogonal Torus)は次数8、65536ノード構成で平均距離7.97という特長を持つ。また、木構造のマルチキャストが可能、メッシュを内蔵するなどの特長も有する[1]。

本稿ではデッドロックを回避するためのルーティングや仮想チャネルの構成を示す。一般に、デッドロック回避の目的で仮想チャネルを設けると実効バッファサイズが小さくなるのが問題であるが、その問題を解決するために用いた適応フロー制御について説明する。

2 RSOT のトポロジ

RSOTは図1(a)に示すレベル0~4の5種類の二次元トーラス網から構成される。各ノードはレベル0トーラスとレベル1から4のうちの一つのトーラスに接続され、次数は8となる。ノード配列を図1(b)に示す。

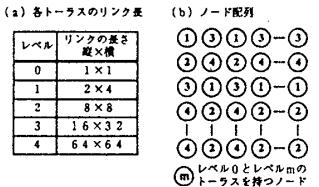


図1: RSOTのトポロジ

3 RSOT のルーティング

k-ary n-cubeは、使用する物理チャネルの方向(次元)に順序を設けることによってデッドロックフリーのルーティングが実現できる。RSOTでは、k-ary n-cubeの次元はレベルとトーラス上の方向(次元)に対応する。RSOTは、ディレクトリベースキヤッシュコンバーチョンプロトコルを支援することを主目的とする木構造のマルチキャストの機能を備えるが、そのとき使用するレベルの順序は上位レベルから下位レベルであるのに対して、マルチキャストされたコマンドパケットに対するAckパケットは逆方向のレベルの使用順序となる。これによるデッドロックを回避するためには2つの仮想チャネルを備える必要がある。1対1通信に関しては、基本的に最適ルーティング(平均距離の最も小さいルーティング)を用いるが、デッドロックフリーであるためには、レベルの使用順序についての制限があるため、すべての始点・終点間で最適ルーティングが可能ではない。しかし、最適ルーティングが不可能となる場合は全体の13%にまで抑えられることが分かっており(4.3節)、この場合平均距離の増加は0.13程度に抑えられる(16384ノード構成時)。

4 物理チャネルの分類と仮想チャネルの設定法

物理チャネル、仮想チャネルの分類を図2に示す。

4.1 物理チャネルの分類

縦方向物理チャネルをV(vertical)、横方向物理チャネルをH(horizontal)と分類する。レベル0の物理チャネルについてはさらに、その物理チャネルによって低レベルのノードから高レベルのノードに移動するときはR(rise)、その逆のときはF(fall)のように分類する。レベル0の物理チャネルについて、さらに細かく分類するときは、その物理チャネルの始点、終点のレベル L_s, L_e を用いた(L_s, L_e)により区別する。

4.2 仮想チャネルの設定

レベル1~4の物理チャネルにはun, uc, dn, dcの4つの仮想チャネルを設定する。レベル0の物理チャネルにはその4つとさらにtn, tcの2つの仮想チャネルを設定する。

The Routing Scheme of RSOT
Tomoyuki AKIYAMA, Hidehiko TANAKA
Faculty of Engineering, University of Tokyo

物理チャネル		仮想チャネル
レベル	方向	
1		
2	H	u d
3	V	n c
4		

物理チャネル		仮想チャネル		
レベル	方向	レベル移動の向き	始点・終点のレベル	仮想チャネル
		R	(13) (24)	
		F	(31) (42)	
0	H	R	(12) (23) (34) (14)	u d t
	V	F	(21) (32) (43) (41)	n c

図2: 物理チャネル、仮想チャネルの分類

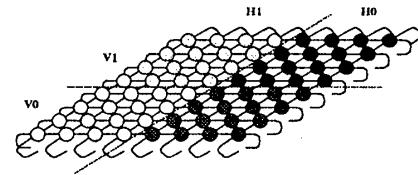


図3: 二次元トーラス網の領域分割

n(not crossed)、c(crossed)は二次元トーラス網中の横[縦]方向物理チャネルのみによって構成されるループによって生じるデッドロックを回避するためのもので、次のように使用される。あらかじめ、図3の様に各二次元トーラスに対して領域分割を行なう。パケットがあるレベルにおける二次元トーラスの横[縦]方向の移動を始めた時は、必ずnの仮想チャネルを用いるが、領域H1からH0[V1からV0]に移る時に、仮想チャネルをcに切り替える(図6(a))。こうすることにより、特定レベルの特定方向内でのデッドロックは回避できる。次に特定レベル内において、使用する物理チャネルの方向にH → V → Hのようなループが生じないことを保証する必要がある。そのためには、あるレベル内ではH → VあるいはV → Hの順序で使用するという条件を満足するルーティングを行なえば良い。最後に、使用するレベルに例えば0 → 1 → 2 → 3 → 1のようなループが生じないことを保証する必要がある。そのためには、レベルを下位から上位の順に用いるなど、ある特定の順序に従ったルーティングを行なえば良い。厳密には、RSOTではレベル1の次にレベル2を用いるということは出来ず、レベル1の次にレベル0によってレベル2を持つノードへ移動した後でレベル2を用いる必要がある。そのため、レベルを下位から上位の順に用いる場合を考えると、レベル間依存関係は0 → 1 → 0 → 2 → 0 → 3 → 0 → 4のようになる。このままで依存関係にループが生じるので、レベル0をレベル間移動に用いるときは専用の仮想チャネルt(transit)を用いる。0t(レベル0の物理チャネル中の仮想チャネルt、以下チャネルをこの様に記述する)の間の依存関係によってループが生じないことを保証する必要がある。前記の場合では、0tはそれぞれ0(12)t, 0(23)t, 0(34)t, 0(13)t, 0(24)t, 0(14)tであり、すべて異なるため、ループは生じない。RSOTは、3節に述べたように、少なくとも上位から下位、下位から上位の2つの順序でレベルを使用する必要がある。そのためには上位から下位と下位から上位の場合で仮想チャネルを使い分ける必要がある。レベルを上位から下位の順で用いる場合仮想チャネルd(down)を使い、逆の場合はu(up)を使うことにする。

4.3 仮想チャネル間の依存関係

4.2節で述べた仮想チャネル間の依存関係を図4に示す。(a)は基本形を示しており、この場合ルーティングはレベルの使用順序において、上位→下位、下位→上位の2つの順序に限定される。そのため、このルーティングに従った場合、平均距離は0.347増加する(16384ノード構成時)。

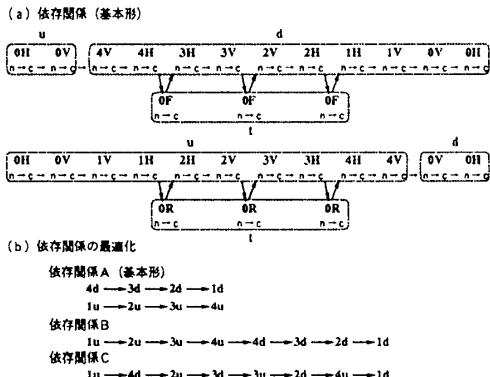


図4: 仮想チャネル間の依存関係
ド時)。(b)の依存関係Aは(a)を簡略化して記述したもので、レベル0や、n、cの区別、H、Vの区別は省略してある。(b)の依存関係B、Cのような依存関係を持たせ、この範囲でルーティングを行なうことにより、平均距離の増加をAより小さく抑えることができる。依存関係B、Cを用いた時、レベル間移動に用いる仮想チャネル0t間の依存関係にループが生じないことが保証される(依存関係Cに対する証明を図5に示す)。平均距離の増加は依存関係B、Cでそれぞれ0.141、0.131(16384ノード時)となる。そこで依存関係Cを用いることにする。このとき平均距離は16384ノード構成で6.95となる。

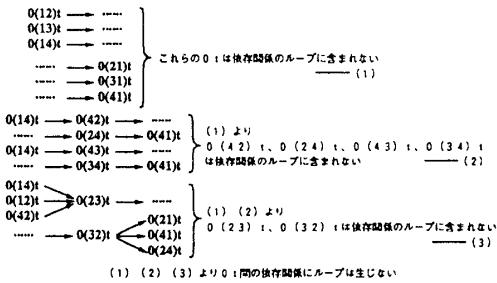


図5: 0t間の依存関係にループが生じないことの証明

5 適応フロー制御による実効バッファサイズ'低下の抑制
デッドロックを回避するために複数の仮想チャネルを設けると、一般的に、仮想チャネルの使用頻度には偏りがあるため、実効的なバッファサイズが小さくなるのが問題である。仮想チャネルのサイズを使用頻度に合わせて選択する方法は、アクセスパターンによって使用頻度が異なる場合が多く、最適なサイズを選択するのが難しい。RSOTでは次に示す適応的なフロー制御を行なうことにより、実効バッファサイズの減少を抑制する。ただし、Virtual Cut-Throughを用いたパケット単位のフロー制御を前提としており、ここではフロー制御は仮想チャネルの割り当てという意味で用いる。

5.1 トーラス網における適応フロー制御

(a) 基本方式



(b) 適応フロー制御



図6: トーラス網における適応フロー制御

図6は図3の横[縦]方向の断面を表している。通常(a)に示すようなフロー制御を行なうことによって、二次元トーラス網中の横[縦]方向物理チャネルのみによって構成されるループにより生じるデッ

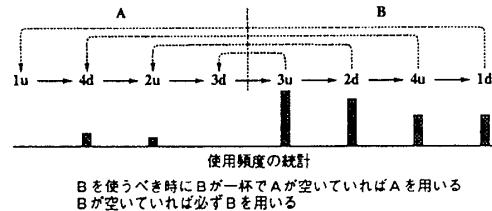


図7: レベル間依存関係における適応フロー制御
ドロックを回避できることは4.2節で述べたが、このような方式をとると、仮想チャネルcの使用頻度がnよりも小さいため、cが最大限利用できず、実効バッファサイズが小さくなるのが問題である。(b)に示す適応フロー制御を行なうことにより、c、nをともに最大限利用することが可能となる。この方式では、領域1から0に移らないパケットは、もしもnが空いていれば必ずnを使うが、nが一杯でcが空いている時はcを使う。一度cを使っても、次のホップの時にnが空いていれば必ずnを使う。nが一杯のためにcに移されたパケットは次のホップではnとcの双方に行く可能性があり、一見cからnに依存関係があるよう見えるが、そのパケットが進むための条件は、前方のnあるいはcのどちらかが空いていることであり、前方のcが空くことが保証されれば、cからnへの依存関係は無い。また、nが一杯の時cへ移すことのできるパケットは、領域1から0に移らないパケットに限定しているため、仮想チャネルcの間で依存関係のループが生じることもない。したがって、デッドロックフリーであることが分かる。

5.2 レベル間依存関係における適応フロー制御

図7は図4(b)の依存関係Cを満足する1対1通信のルーティングを行なった場合の、各仮想チャネルの使用頻度を示している。図を見ると明らかにグループA(1u, 4d, 2u, 3d)よりもグループB(3u, 2d, 4u, 1d)の方が使用頻度が高い。あるノード対に対して、最適ルーティング(ホップ数が最小のルーティング)は一般に複数存在するが、その中からBの使用頻度が高いものを選択したためである。このような選択をした理由は後述する。このとき、次のような適応フロー制御を行なうことにより、バッファを最大限利用することが可能となる。すなわち、パケットがBに属するあるレベルの仮想チャネルを使用すべき時にその仮想チャネルが一杯で、Aに属するそのレベルの仮想チャネルが空いている時は、Aの方を用いる。Bが空いている時は必ずBを用いるようにする。Bが一杯であるためにAに転送したパケットは、次のホップでBに空きがある場合は必ずBに戻す。このような操作を行なっても、BからAへの依存関係は生じないため、デッドロックは起こらない。

次に頻度分布をBに偏らせた理由を述べる。Bに使用頻度を偏らせるという操作を行なわなかった場合、Bに属する仮想チャネルよりもAに属する仮想チャネルの方が頻度が大きいという状況が生じ、またその逆の状況も生じる(レベルによって異なる)。前者の場合、Aを使うべき時にAが一杯でBが空いていればBを用いるというような適応フロー制御も可能だが、常にBを用いて良い訳ではなく、次のような条件を満たしているパケットについてのみ、Bを用いることができる。すなわち、例えば4d, 2uの順に進むべきとき、4dが一杯であるため4uを用いるということを行なうと、4uから2uに依存関係が生じるため、デッドロックが起こる可能性が生じる。したがって、4uを使うためには、その後に4uより左の仮想チャネルを用いることがあるかどうかを調べなければならない。このようなコストをかけないために、あらかじめBの使用頻度が高くなるようにルーティングテーブルを設定し、Aが一杯のときBが空ならBを用いるというような操作を必要としない状況を作り出す。

6 おわりに

本稿ではRSOTにおいてデッドロックを回避するためのルーティング、仮想チャネルの構成を示した。デッドロックを回避するために仮想チャネルを設けると、実効バッファサイズが小さくなるが問題であるが、その問題を解決する適応フロー制御について述べた。

参考文献

- [1] 秋山知之、小池帆平、田中英彦、分散共有メモリ型超並列計算機におけるディレクトリ方式を相互結合網について、信学技報CPSY94-49, 65-72.