

キューブ系ネットワークの特性

1P-5

村松 晃¹、田中輝雄²、林 剛久²、吉原都夫¹、前田栄一郎³¹(株)日立製作所システム開発研究所 ²同、中央研究所³日立ニュークリアエンジニアリング(株)

1. はじめに

疎結合並列計算機の相互結合ネットワークとしては、要素プロセッサ (PE) 間の通信トポロジーをできるだけ多種類内包しているネットワークが望ましい。この種のネットワークとして評価の高いのが、数値計算によく現れる格子、リング、バタフライ等のトポロジーを内包し、これらのパターンは中継無しで通信できるハイパーキューブである。しかし、中継を行うとデッドロックが起きる可能性がある。

また、一般に、ネットワークのハードウェア規模は、PE台数Nの1次以上のオーダーで増加する。増加の割合は高性能ネットワークほど大きく、通信性能とハードウェア規模はトレードオフ関係にある。

本稿では、キューブ系ネットワークのデッドロックを回避する通信手順と、性能とハードウェア規模のトレードオフ関係を評価した結果について述べる。

2. アルファネットワーク

ここでは、一般化されたハイパーキューブであるアルファネットワーク¹⁾ (α NW) を取り上げる。 α NWは、PE台数Nの因数分解： $N = m_1 \cdot m_2 \cdots m_n$ にもとずいて構成される $m_1 \times m_2 \times \cdots \times m_n$ に格子分割されたn次元超直方体を考え、その各格子点に各PEを対応させ、各次元方向に完全結合したものである(図1)。普通、完全結合ネットワークに対してPEは逐次的に入出力を行うから、出力側結合リンクを選択するデコーダと入力側結合リンクを選択するセレクタを持つ。これは完全結合をクロスバスイッチで実現することと等価である。このとき、1次元の α NWはフルクロスバスイッチに、また($N = 2^n$ のとき)n次元の α NWは2進のハイパーキューブに相当する。

2進のハイパーキューブの問題点として、

- ① 拡張性 (PE台数が2のべき乗に制限される)
- ② 実装性 (VLSIでのレイアウトがしにくい)
- ③ ハードウェア規模 (自動中継機能付きの場合)
- ④ 中継段数 (最大でPE台数の対数)
- ⑤ 問題の写像容易性 (3次元以下が望ましい)

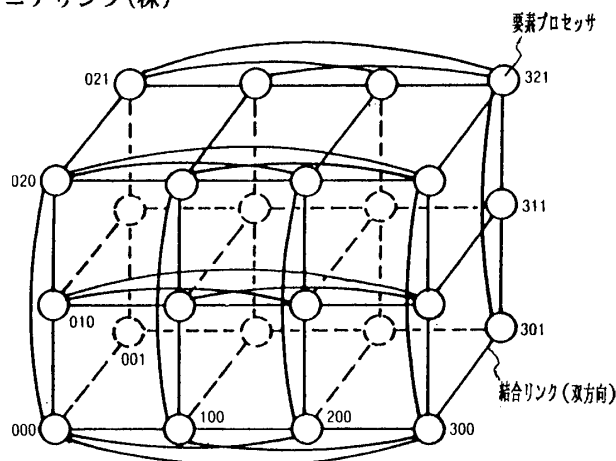


図1 4×3×2アルファネットワーク

⑥ デッドロックの可能性

等があるが、 α NWまで拡張して考えれば①～⑤は改善される。また、次に述べるような通信手順によれば⑥も回避できる。

3. デッドロックフリーの送信手順

α NWにおける送信手順として、送信側PEのアドレス座標 (i_1, i_2, \dots, i_n) と宛先アドレス座標 (j_1, j_2, \dots, j_n) に関し最も低位の次元の不一致座標対を選び、該座標対を結合しているクロスバスイッチを選択してパケットを出力するものを考える。これを繰り返せば宛先に送信することができる。この手順によればデッドロックを起こさないことが、以下のようにして示される。

Dally&Seitz²⁾によれば、PEを枝に、結合リンクを節点にして相互結合関係を示した依存グラフがサイクルを含まなければ、デッドロックフリーである。帰納法を用いれば、① 1次元の α NWにおいては、任意のPE間は高々1ステップで送信可能である。従って、依存グラフは節点のみで構成されるのでサイクルは含まれない。② n次元の α NWは、 $n-1$ 次元の α NWがデッドロックフリーであればデッドロックフリーである。なぜなら、上記手順によれば第n次元に属する結合リンクは転送経路の最後

に選択されるため、サイクルができることは決してないからである。図2に例を示す。(a)は2次元の α NWにおいて、各PEが斜め向いのPEにパケットを送ろうとしている場合である。結合リンクa, a', c, c'は第1次元の、b, b', d, d'は第2次元の結合リンクである。上記手順によらず、リンクa, b, c, dのみを用いると(b)のように依存グラフにサイクルが現れる。上記手順によれば(c)のようにサイクルは現れない。このとき、b, b', d, d'は最後に選択されていることに注意されたい。

4. 性能評価

(1) ハードウェア規模：パケットの送信・中継・受信を行うルータ(クロスバスイッチ)を各PE毎に持つ α NWのハードウェア量を、クロスポイント数で評価すると、

$$m_1 \times \dots \times m_n \times (n^2 + m_1 + \dots + m_n)$$

となる。図3に次元構成とハードウェア量の関係を示す。フルクロスバスイッチとハイパーキューブの間にハードウェア量最小の構成があることが分かる。

(2) 通信性能：GPSで記述したシミュレータにより、 α NWの動作時性能を評価した。シミュレーションの条件は以下のとおりである。①パケットは固定長。②PEは一定間隔でメッセージパケット発行命令を実行し、ルータがビジーの場合は待たされる。③PE、ルータ、次元方向のクロスバスイッチは1パケット長分のバッファを持ち、パイプライン的に送信する。④到着パケットは、一定時間後に消費される。⑤通信パターンは規則型(リング、格子、バタフライ型通信)、不規則型(各PEが毎回ランダムに宛先PEを決める)の2パターンとする。

結果は、総通信パケット数はPE台数に比例して増大するが、予期に反して、同一PE台数では高次元構成の方が多い(不規則型の場合)(図4)。これは、 α NWを通信用メッセージキューと見立てたとき、高次元構成の方がキュー容量が大きく、パケットの送信ピッチを短くできるからである。ただし、次元構成による差が目立つのはパケット発行頻度が非常に高い場合である。一方、トラベルタイム(パケットを発行してから宛先に届く迄の時間)は低次元構成の方が短い。

5. おわりに

疎結合並列計算機向きのキューブ系ネットワークとして α NWを取り上げ、その特性を評価した。その結果、低位次元結合リンクから送信していく手順を用いればデッドロックフリーであること、フルクロスバとハイパーキューブの間にハードウェア量最

小の構成が存在すること、トラベルタイムは低次元構成の方が短い総通信パケット数は高次元構成の方が多い傾向にあり、パケット発行頻度の高い場合にこの傾向が目立つこと等の知見を得た。実装性を考えると、数千台以下の規模の疎結合並列計算機では2~4次元の α NWが有効であると思われる。

6. 参考文献

- 1)Agrawal & Janakiram: Evaluating the Performance of Multicomputer Configurations, IEEE COMPUTER, May 1986
- 2)Dally & Seitz: Deadlock-free Message Routing, IEEE Trns.on Computers, C-36, 5, 1987

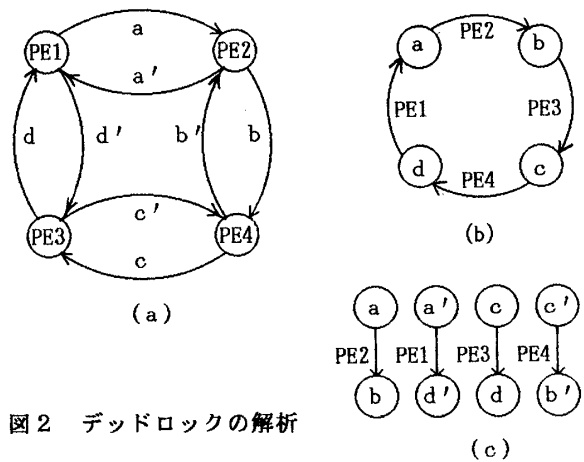


図2 デッドロックの解析

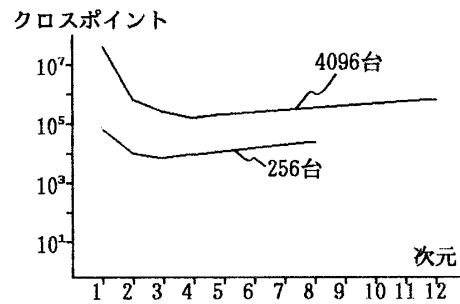


図3 ハードウェア規模

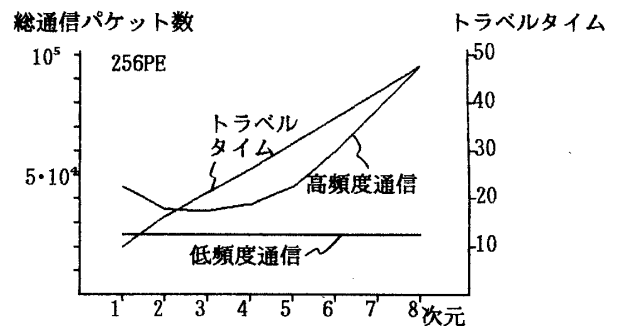


図4 次元構成と性能