

トーラス・リング・バス階層アーキテクチャ

伊藤 拓 小畑 正貴 宮垣 嘉也
岡山理科大学工学部

数万台規模の並列計算機の実装においては、そのネットワーク量とコネクタなどのピンヘッドの削減がコストと共に実装可否の重要な条件となってくる。これらの問題に対して、接続ネットワークを削減、局所化するクラスタ構造の導入が有効であるが、この手法は一般的にクラスタ内外での通信性能に不均衡を生じる事から通信性能の低下の問題がある。本稿では二次元トーラスとリングを用いたクラスタ構造にさらにバスによる通信径路の補強を行なった階層型ネットワークを提案しその基本的な性能について述べ、またこのネットワークが数万台程度のPE数による近接データを用いた数値シミュレーションに対しても有効性がある事を示す。

THE DESIGN AND EVALUATION OF TORUS-RING-BUS NETWORK

Taku Ito Masaki Kohata Yoshiya Miyagaki
Okayama University of Science

The increase of the number of interconnection wires is one of the serious problem in the implementation of large-scale multi-processor systems. In regard to this problem we propose a new cluster structure network called Torus-Ring-Bus Network (TRB). In a cluster, processors are connected by both of ring and bus, and the clusters are connected by torus network. This paper shows the network construction and evaluates the ability.

1. はじめに

近年では、数値シミュレーションなどに代表される大規模な科学技術計算を高速に実行するコンピュータとして数千から数万台の規模のPE数を擁する並列計算機に対する期待も高まりつつある。しかしながら数万台規模の並列計算機ではネットワーク量や基板間の接続コネクタのピンヘッドが、そのコストと合せて実装可否の重要な条件となってくると言える。

このような問題に対してネットワークとハードウェアコストを抑えつつ、多数のPEを実装するアプローチとして複数のプロセッサをローカルに接続して1つのネットワークノードとするクラスタ構造の導入が有効であると考えられるが、この手法の問題点としてはクラスタの内外で通信性能に不均衡を生じる点や複数のクラスタ間にまたがる通信を必用とするプロセスの配置に対して柔軟さに乏しい点などが上げられる。

本稿ではリングとバスで接続されるPE群を1クラスタとし、そのクラスタを二次元トーラスで接続した階層型のネットワーク (TRB) を提案し、その構成と基本的な性能について考察する。

2. ネットワークの構成と実装

2.1 ネットワークの構成

図1にTRBによって結合された並列プロセッサの構成例を、図2にそのクラスタ内部の構成例を示す。

図1ではN個のPEを双方向のリング、および2本のバスによって接続したものを1クラスタとし、各々のクラスタではバスに接続されているTC (Torus Controller) を介してN×Nの2次元のトーラスネットワークに結合される。TCは主に2次元トーラスからくる4本のチャンネルとクラスタ内部のバス間での通信のルーティング及び、クラスタ単位でPE間の同期を行なう。

各PEの内部は図3に示すような構成となりRC (Ring Controller) を介してクラスタ内のリングに接続される。通常クラスタ内部の通信はRCを介して隣接PE間でリングによって行ない異なるクラスタへの通信は、バスを介してTCによってルーティングされる。

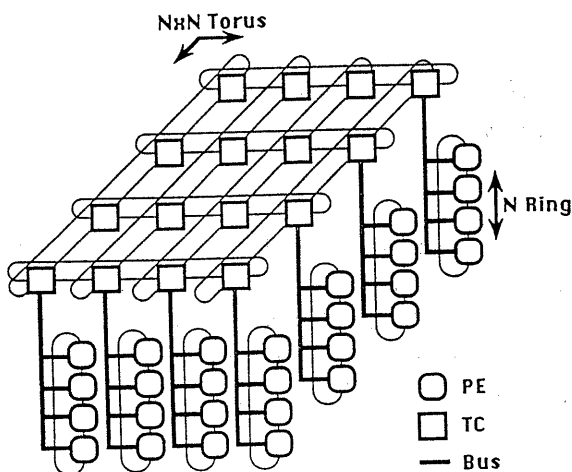


図1. TRBの構成

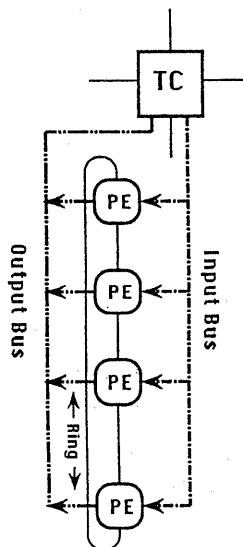


図2. クラスタ構成

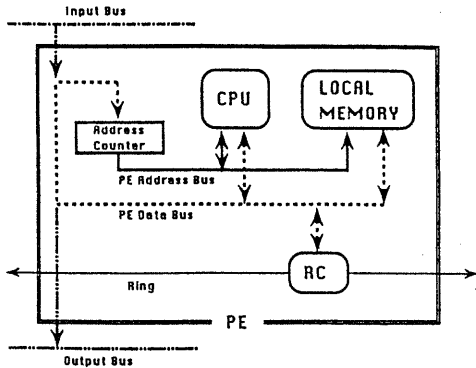


図3. PEの内部構成

PE台数 ()TRB一辺	Torus2D	Torus3D	TRB	
	Link	Link	Link	Bus
64 (4)	128	192	96	32
729 (9)	1458	2187	891	162
4096 (16)	8192	12288	4608	512
15625 (25)	31250	46875	16875	1250
46656 (36)	93312	139968	49248	2592

表1. 通信ネットワーク量の比較

図1、図2は共にクラスタ内PE数と、トーラスの1辺を4とおいた64プロセッサの構成例であるが、これ以外にも様々な組合せがありえる。また、トーラスの通信路が4本である事からクラスタ内部のバス数をそれに合わせる構成も考えられる。これについては、実装上でTCのハードウェア量と配線コストの増大を招く為、実装アプリケーションの性能向上とコストの両面で検討する必要がある。

2. 2ハードウェア量と実装

TRBでは総PE数をNとした場合の各PEを接続するリンク数は $2N^{2/3}+N$ となり、この時のバスは $2N^{2/3}$ 本となる。総PE数Nを等しくおいた場合の通信ネットワーク量をTRBと2次元トーラスの $2N$ 及び、3次元トーラスの $3N$ とで比較し表1に示す。

実際の実装に関してはTRBではバスによるクラスタ内結合の限界と、ランダムトラフィックで密な通信を行なった場合、通信路の共有による性能低下が顕著になると考えられる事から無制限にPE数を増やす事は無理である。しかしクラスタ化によって配線が局所化され、リンク数も少ない事からコネクタピンヘッドや配線面では有利であり、クラスタ数とトーラスの1辺が等しくnの場合の実装PE数 n^3 より $n=16$ で4096台、 $n=32$ で32767台程度のPE数での実装が可能であると考えられる。

3. ネットワークの通信特性

3.1 ランダムトラフィックにおける通信性能

TRBでは基本的なポイントツーポイントでの平均的な通信距離はPEの総数をNとして、クラスタ内を D_R とクラスタ外 D_T をとおくと、 $D_R \approx N^{1/3}/2$ 、 $D_T \approx N^{1/3}/2 + 2$ 。また、システム全体ではクラスタ内外の通信比率、 $N^{1/3} : N - N^{1/3}$ によって平均の通信距離Dは次の様になる。

$$D \approx \frac{N^{1/3}}{2} + 2 - \frac{2}{N^{2/3}} \approx \frac{N^{1/3}}{2} + 2$$

ランダムトラフィックな通信を行なう場合、複数の通信の競合が通信性能に及ぼす影響が深刻である為、これを調べる目的でC言語による通信シミュレータを作成し実験を行なった。

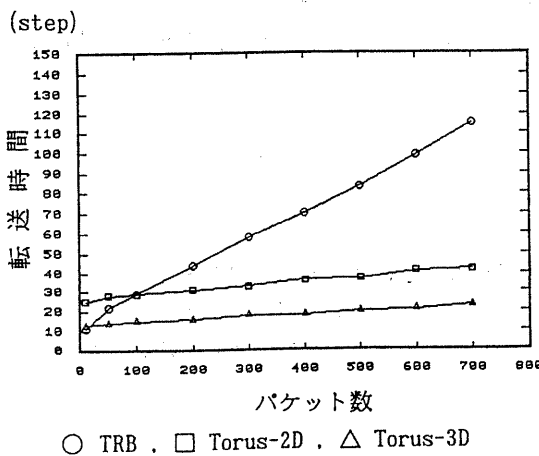
通信シミュレータでは以下の条件でのランダムなポイントツーポイントの通信を行なう。

- ネットワーク上の各ノード（TCを含む）は1単位時間当り1パケットの送信と受信が可能である。
- 通信の形態は不規則で局所性を持たない通信を基本とする。
- 送信PEと受信PEが同一のクラスタである場合のルーティングはリングを介してのみ行なわれ、また目的PEに対して常に近づく方向にルーティングされる。
- TCによるルーティングでは目的クラスタに常に近づく方向でルーティングされ、また近づく径路が複数存在し一方が使用中であれば他方に対してルーティングされる。
- 送信径路が全て使用中であるパケットはTC、RC共にブロックされる。

以上の条件のもとにある時点で複数のPEから同時に発生したパケットが全て目的のPEまで到達するまでの通信時間をパケット取拾時間とし、異なるネットワークでの通信性能の指標として比較を行なった。729台のPEでのパケット取拾時間を図4に示す。またクラスタ内外での通信比率が変化した場合のシミュレーション結果を図5に示す。なお、これらのグラフはPE台数が変化してもパケット量とPE台数の比によって同様な形状の特性を示した。

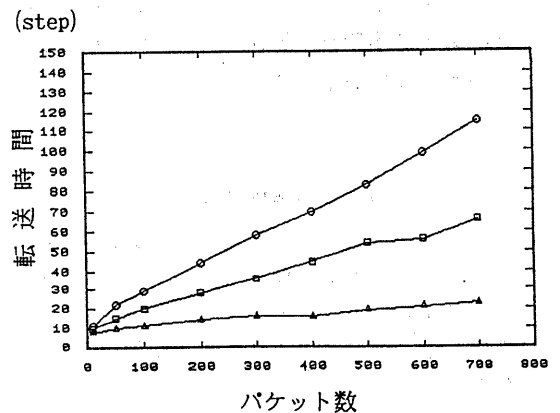
図4から、TRBが100パケット以下の領域で2次元トラスに対して優位である事がわかる。これはネットワーク上でのパケットの競合が少ない場合にPE台数を同じにとるとTRBが2次元トラスより平均的な通信距離が短い事によるものである。しかし、ランダムトラフィックにおいてはパケットが密になってくるとTRBはトラスの場合と異なり通信性能の低下が著しい。

図5ではクラスタ内外における通信比率をそれぞれ、内：外で1：9，5：5，9：1とした場合の通信性能の変化を調べたものである。



○ TRB , □ Torus-2D , △ Torus-3D

図4. 同時発生パケットと取拾時間
(729 Processor)



In Out In Out In Out
○ 1 : 9 , □ 5 : 5 , △ 9 : 1

図5. クラスタ内外での通信比率による
取拾時間の短縮 (729 Processor)

図より、クラスタ内の通信割合が大きくなるほどパケット取捨時間が減少しており、通信性能が向上する事がわかる。これはクラスタ内の通信割合を増加し、クラスタ外への通信が減少させる事によって、クラスタ外ネットワークの競合が減少し、クラスタ単位で通信多重度が上昇する事からも明らかである。この特性は一般にクラスタ構造によって複数プロセッサをローカルに接続し、グローバルなネットワーク量を削減したネットワークでは共通するものである。

3. 2 近接通信における通信性能

TRBは基本的に1辺を n とする3次元トーラスの2軸方向のネットワークを $1/n$ に削減しバスで括った構成となる。この構成での近接通信におけるTRBの性能を2次元、3次元の近接通信をそれぞれ例にとって考える。

3次元の近接通信において、3次元トーラスの一辺とTRBのクラスタ内のPE数を共に n とした時、3次元トーラスにおける東、西、南、北、上、下の6リンクからのデータ通信ステップを6ステップとすると、同等な通信条件でTRBでは2本のバスを用いて $4n+2$ ステップとなり、3次元トーラスに対して $3/(2n+1)$ 倍の性能になる。また2次元の近接通信においては、2次元トーラスの4リンクに対して東、西、南、北に同様な通信を仮定した場合、2次元トーラスの4ステップに対して、TRBでは2本のバスとリングの併用によって図6のように1軸方向のリングを複数クラスタ間に跨って構成した場合で $n+2$ ステップとなる事から2次元トーラスに対するTRB通信性能は、 $2/(n+1)$ 倍となる。この場合のTRBは図7に示す様な $n^2 \times n$ の長方形型の2次元トーラスを仮想している。

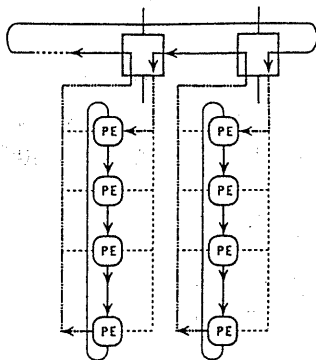


図6. 複数クラスタ間に跨るリングの構成

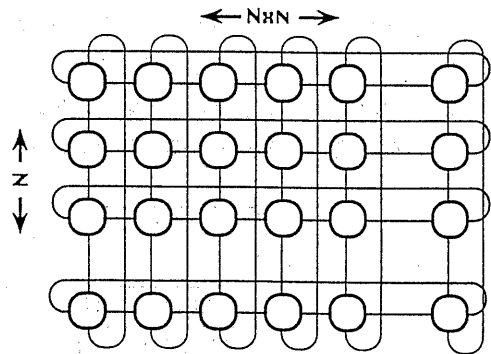


図7. TRBで構成される仮想的な $n^2 \times n$ の二次元トーラス

4. 応用プログラムでの性能評価

4. 1 システムの応用分野

一般に並列プロセッサの応用範囲としては、数値シミュレーション、グラフィクス、人工知能、データベースと言ったように様々な分野が期待されている。なかでも数値シミュレーションは静的な負荷分散が比較的容易で、データ並列性が高いと言った利点から並列処理の応用分野としての期待が高まっている分野である。

ここでは数値シミュレーションの中でも最も多用され、また応用範囲も広いと考えられる偏微分方程式の数値的求解を基本とする近接作用問題を例にとって、問題の処理粒度、負荷分散、TRB

上での通信性能等を考慮して実装時の台数効果を見積る。

近接作用問題とは連続空間に分布する温度や圧力といった物理量が、場を支配する偏微分方程式の形で表わす事が可能である場合に、空間を格子点状に分割し偏微分方程式を差分近似等で全ての格子点間の連立1次方程式の形に変形した後、この連立方程式を数値的に解く事で新たな分布を求めると言う問題である。この問題は各格子点の物理量とその微係数と近接する格子点の物理量によって計算できると言う特徴を持っている。

このような問題の計算処理はその主体が一般的に対角優位な疎行列となる連立1次方程式の求解である事からSOR法などに代表される緩和法によって収束解を求める手法に適している。しかし並列処理化に関しては緩和段階のループ処理にデータ依存関係がある事から、依存関係の無い奇数点の点緩和と偶数点の点緩和をそれぞれブロックにまとめて並列化するRed & Black SOR法やSLOR法などが用いられる場合が多い。

4. 2 近接作用問題での性能見積り

前章ではTRBとトラスの通信性能を近接通信で比較したが、実際に応用問題を実装して処理する場合に、その実行性能におよぼす影響は、各PEにおける処理の粒度とばらつき、通信量、通信速度、問題の規模、PE台数と言った種々の条件の総合によって決ってくるものである。したがってこれらを定量化する事が可能であれば台数効果を見積る事ができる。以下にそれぞれの条件の見積り方法とその値を示す。

- ①処理粒度：境界条件を考慮にいれた緩和法の差分式1ループに要するCPUの処理時間をC言語でコーディングし実測して用いた。対象CPUは、RISCチップとしてはSPARC 25MHz、DSP TMS320C30 40MHz、汎用CPUでは80386+387 20MHzの三つで、測定結果を表2に示す。ただしSPARCに関しては割り込みを止めていない為、やや厳密性に問題があるが、最近のRISCチップの処理速度の指標とした。

CPU		処理時間(μs)	
		2次元	3次元
SPARC	40MHz	7.0	9.0
TMS320C30	40MHz	13.8	15.9
80386+387	20MHz	76.0	89.2

表2. 緩和法1ループの処理時間

- ②処理のばらつき：境界条件の定義のしかたとコーディング方法によってやや変化するかも知れないがほとんど無視できると考える。
- ③通信量：一般的には2次元問題で自分自身を含む近接の5点、3次元問題では7点の格子点値である為、1ループ、1格子点毎にそれぞれ4ワード、6ワードの通信量を必用とする。
- ④通信速度：現在64プロセッサによるTRBのプロトタイプ的设计を行っており、1リンクあたりに通信速度は5~10M(ワード/sec.)程度となる予定である。また1通信あたりの同期処理とDMAの駆動の為のオーバーヘッドとして2(μs)程度を想定している。

⑤負荷分散：格子点をブロック単位でPEに分割し隣接する格子点をできるだけPE内部に集める。この分割によって通信が必用となるのはブロック分割の境界面に位置する格子点のみとなる。

⑥問題の規模：流体計算などの分野で慣性項が大きい計算では少なくとも3次元で800万格子点、2次元で8万格子点程度の規模が必用と言われている事から⁴⁾この値を問題の規模とする。しかし2次元問題の場合は数万台で実行される場合、各PEの負荷が小さくなり過ぎる為、⑤の通信量の削減が殆ど行なえなくなる事から通信時間の割合が増大してくる。したがって、この規模ではTRBによる数万台での実行には不向きである。

以上の条件から3次元問題に対して台数効果の見積りを行なった。この時の格子点数は800万であり、また並列化を行なっても点緩和の収束回数は殆ど変動しないものと仮定している。CPUがSPARC、DSP、80386のそれぞれの場合で通信速度が5M(w/sec.)、10M(w/sec.)の二通りについて求めてみた結果をそれぞれ図8、図9、図10に示す。SPARCでは3万台付近で台数効果はやや低下し5M(w/sec.)の場合で約63(%)、10M(w/sec.)で約74(%)程度の台数効果である。しかし、いずれの場合でもこの問題では3万台程度までのPE台数において比較的高い台数効果が得られている事がわかる。これより、3次元の近接作用問題に対しては、TRBのようにコストを抑えたネットワークでも数万台規模までであれば十分な有効性が得られる見通しが得られた。

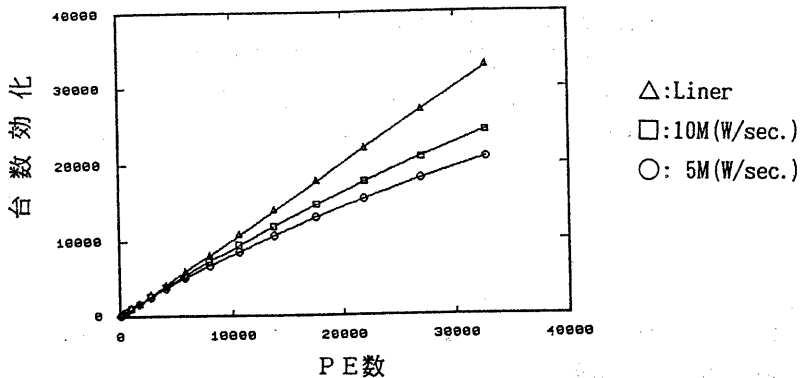


図8. SPARCでの台数効果

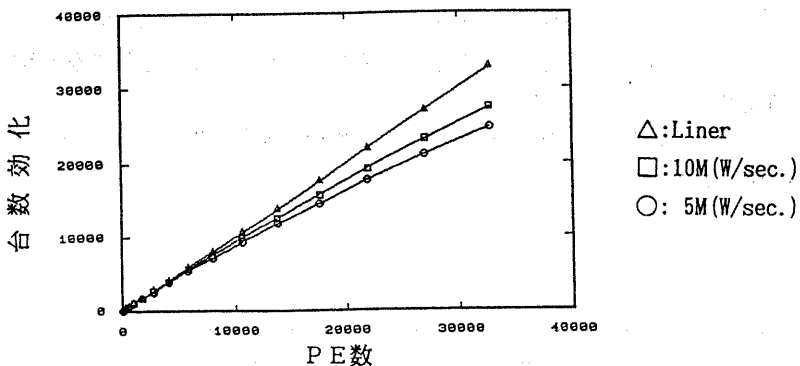


図9. TMS320C30での台数効果

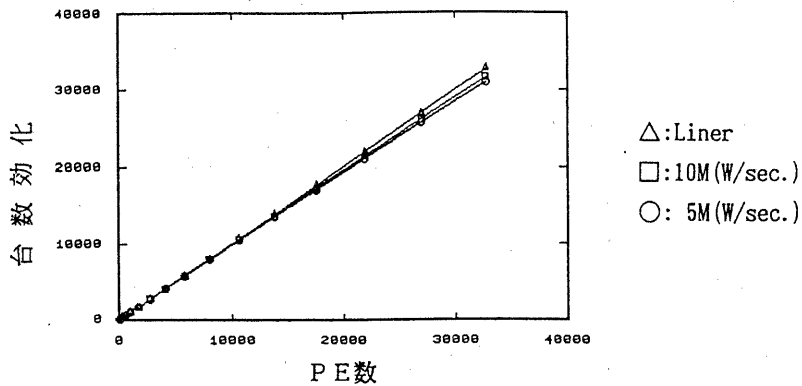


図10. 80386+387での台数効果

5. おわりに

クラスタ化による階層構造の採用によってネットワークのハードウェア量を削減し大規模なPE数での実装を志向するアーキテクチャとして二次元トーラス、リング、及びバスを用いた階層型のネットワークを提案し、その基本的な構成と通信性能について述べた。また比較的小さな処理粒度で頻りに通信を行なう必要がある大規模な3次元の近接作用問題に対して、このネットワークを用いて数万台規模までは十分に性能向上を見込める事を示した。

TRBによる並列システムの開発は現在CPUに32ビット浮動小数点DSP (TMS320C30)*を用い64プロセッサにからなるプロトタイプの実装を目指して、PE基板のパターン設計とルータなどの詳細な設計を平行に行なっている段階である。これらについてはまた別の機会に報告する予定である。

<参考文献>

- 1) 富田、末吉：並列処理マシン，オーム社
- 2) 戸川：数値計算，岩波書店
- 3) 大宮司、他：流れの数値シミュレーション，コロナ社
- 4) 桑原：スーパーコンピュータによる流れのシミュレーションとその可視化，JSPP'92招待講演
- 5) 小畑：256台DSPシステムによる並列計算機の実装と評価，電子情報通信学会，コンピュータシステム研究会，CPSY90-29 1990
- 6) 小畑：シグナルプロセッサアレイ-SPA-，電子情報通信学会，コンピュータシステム研究会，CPSY79-13 1979

*DSPを提供して頂いている日本テキサス・インスツルメンツ(株)に感謝します。
本研究は平成4年度科学研究費補助重点領域研究(04235210)の補助を受けています。