

# 汎用スイッチング・ユニットを用いた 高並列計算機の相互結合網

## Interconnection Network For Highly Parallel Computers, Consisting Of Wide-Usable Switching Units

坂井 修<sup>†</sup> 計 宇生<sup>†</sup> 田中 英彦<sup>†</sup> 元岡 達<sup>†</sup>  
 Shuichi Sakai Usei Kei Hidehiko Tanaka Tohru Moto-oka  
<sup>†</sup> 東京大学 工学部  
 Faculty of Engineering University of Tokyo

### 1. まえがき

昨今のVLSI技術の飛躍により、高並列計算機への期待が高まっており、数十台から数千台、場合によっては数百万台という超多重システムの実現が検討されている。このようなシステムにおける相互結合網の役割は重要である。われわれはすでに、並列処理関係代数マシンGRACEの結合網の検討<sup>(1)</sup>、蓄積交換方式の汎用スイッチング・ユニットの試作と評価<sup>(2)</sup>、高並列推論エンジンPIEの結合網の検討<sup>(3)</sup>などを行ってきた。

一般に相互結合網に要求される性能として、高速性(低遅延)・大通信量(高スループット)・経済性(低いハードウェア・コスト)・信頼性・拡張性・VLSI化の容易さ・問題適応性などがある。並列計算機の設計時には、これらがそれぞれどのくらい必要とされるかを検討し、結合網の仕様が決定されねばならない。相互結合網は、動作モード・制御方式・交換形態・幾何学的形状などから分類される。<sup>(4)</sup>

本稿では、先に述べたスイッチング・ユニット(以下SUと略す<sup>(2)</sup>)の試作経験に基づいたSUの改良設計と、これを用いた結合網の転送性能のシミュレーション評価に関して報告する。当SUは、その入力ポートにFIFOバッファをもつ蓄積交換方式のスイッチであり、可変長パケットを低遅延・高スループットで転送する。また、柔軟性が高く、多段結合網・格子型結合網・CCC網・木状網など、さまざまな形状の

網に適用でき、さらに拡張性も高い。今回の改良設計は主に迂回制御に関するものであり、その他にバッファの動作を高速化した。

次章以下の内容は以下のとおりである。すなわち、2章では当SUの構成と動作を示し、3章でこれを用いた相互結合網の特徴を述べる。4章では当SUを用いて構成される結合網のシミュレーション評価を、2種類の多段結合網に関して報告する。5章では、より具体的な転送性能の検討と、ハードウェア量の評価・LSI化の検討を行う。

### 2. 汎用SUの構成

#### 2.1 SUの基本構成

対象とする相互結合網は、 $m$ 入力 $n$ 出力の比較的小規模なSUによって構成されるとする。本稿で報告するSUは、図1のような全体構成をとるが、以下にその設計の基本方針を列挙する。

- (1)  $m$ 個の入力ポート・ $n$ 個の出力ポートをもち、それらの間は各入力ポートに対応する内部バスで接続される。
  - (2) SU全体が同一の同期クロックで動作する。
  - (3) 制御は各入・出力ポートに分散し、入力ポート・コントローラルーティング制御を行い、出力ポートでスイッチングを行う。
  - (4) 各入力ポート内にはFIFO型のバッファ・メモリがあり、パケット単位にまとめられたデータの蓄積交換を行う。
- また、当SUの動作の特徴は次のようなもの

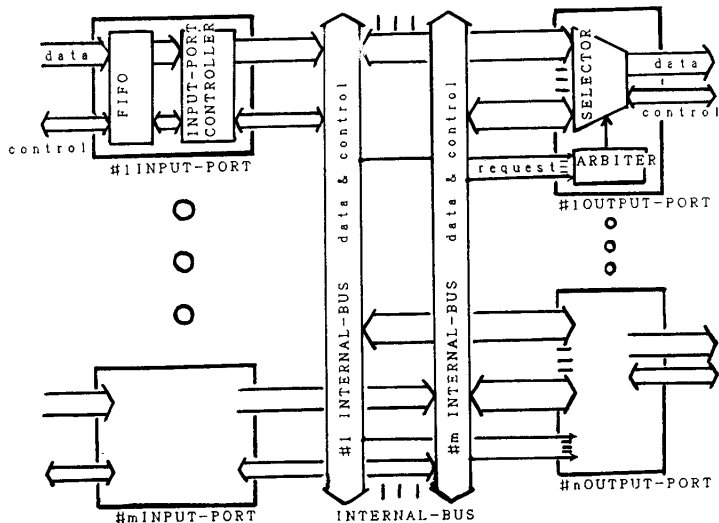


図1 SUの全体構成

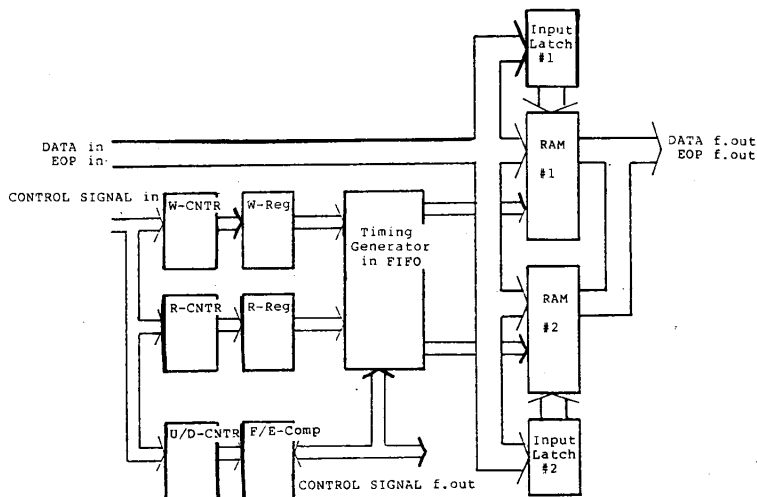


図2 入力ポートFIFOの構成

である。

- (1) データ転送は制御信号 (STB, ACK) のハンドシェイクによって行い、語単位の並列転送を行う。
- (2) ルーティングは、各行先アドレスにつき、あらかじめ複数の出力ポートを指定しておき、閉塞時は次々と候補を替えていく可変ルーティングである。
- (3) 転送はパケット単位に行う。1パケット内で転送をパイプライン化する。すなわち、パケットの末尾の語の到着を待たずに先頭の語が

ら次々に転送する方式をとる。

- (4) パケットは可変長で、先頭の1語 (または数語) が行先アドレスであり、パケットの終りはEOP (End of Packet) 信号によって示す。1語は9 bit であり、うち8 bit がデータ、1 bit がEOPである。
- (5) SUはパケット整形・デッドロックの検出などをせず (3. 2参照)、基本的なスイッチング操作のみを行う。

## 2. 2 入力ポートFIFOメモリの構成

入力ポートFIFOメモリはデータのバッファリングを行うことを前述した。バッファとして市販のFIFO-ICを用いることを検討したが、容量の点・拡張性の点で問題があったため、RAMを採用することにした。FIFOメモリ周辺の構成を、図2に示す。

RAMは2面設け、クロック・サイクルごとに読出し/書込みを交互に行い、パケットをインタリーブして格納する。これによって、1クロックで1語を転送することが可能になる。したがって、書込みカウンタ(次にwriteするRAMアドレス)・読出しカウンタ(次にreadするRAMアドレス)の指すRAMのバンクは、クロックごとに切り替えられる。両者の差はアップダウン・カウンタによって計算され、さらにF/EコンパレータによってメモリのFull/Empyが知らされる。

### 2.3 ルーティング制御

行先アドレスから出力ポートを選択する操作、いわゆるルーティングは、相互結合網の幾何学的形状や、集中制御・分散制御の別などにより方式が異なる。また、同じ形状の結合網でも、複数の経路選択法が考えられる場合がある。

今回設計したSUのルーティング制御は、入力ポート・コントローラ内で、以下の手順で行われる。

- (1) アドレス・レジスタが、パケットの先頭にある行先アドレスの情報を取り込む。
- (2) アドレス・コンバータが、アドレス・レジスタの内容を見て、出力ポートのアドレスを決定する。
- (3) 上記の出力ポート・アドレスをデコードし、実際のポートにREQ(Request)信号を出す。
- (4) 内部バスの駆動を行う。

アドレス・コンバータとして、(i) 組合せ論理回路を用いる方式、(ii) ルーティング・テーブルを記憶したメモリを用いる方式、(iii) その中間の方式を検討した。

(i) は、パケットのアドレス情報から組合せ論理回路で出力ポート・アドレスを得るもので、ハードウェア量も小さく、制御時間も短く(ゲート数段分の遅延)が、柔軟性に乏しく網の形状やSUの配置される位置によって回路を変え

ねばならないのが欠点である。

(ii) は、RAMあるいはPROMに、行先アドレスと出力ポート・アドレスの対応表を記憶し、これを用いてルーティングを行うものであり、柔軟性に富むが、ハードウェア量が大きくなる点、メモリアクセス時間が大きくなる危険がある点が問題である。また、メモリの書込み線が必要となる難点もある。

(iii) は、RAMに書かれた情報(自SUの位置、結合網の形状など)と行先アドレスの情報から、論理回路によって出力ポート・アドレスを得るもので、ハードウェア量も小さく、しかも柔軟性のある方式であり、メモリの書込み線も少なくてよい。

設計システムは、柔軟性と設計の容易さを考慮して(ii)のルーティング・テーブル方式を採用することにした。実際には(iii)の中間方式が最も有利であると考えられ、現在この方式に関して検討を始めている。

### 2.4 可変ルーティング機能の実現

行先アドレスに対して、たゞ一つの出力ポート・アドレスを生成する方式、すなわち固定ルーティング方式では、制御が簡単である反面、閉塞によって結合網全体の性能の著しい低下をまねくおそれがある。また、結合網の信頼性も低い。そこで、SUのレベルで可変ルーティング機能を持たせることにし、次のように実現した。

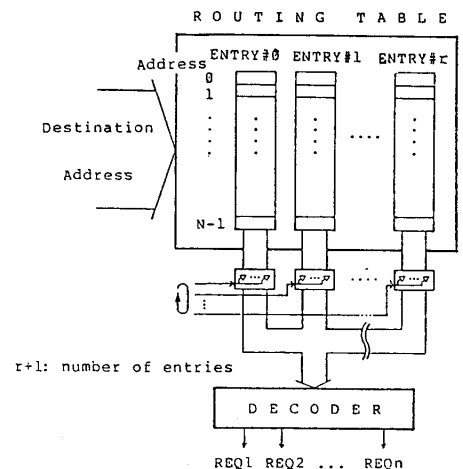


図3 可変ルーティングの実現

図3に示すように、アドレスコンバータ・メモリ(ルーティング・テーブルを格納する。)の容量を増し、転送先の出力ポート・アドレスを複数個( $\leq n$ )指定できるようにする。

アドレスコンバータ・メモリの各番地には行先アドレスを対応させ、各番地に複数( $r+1$ 個)のエントリを設ける。このエントリ内に、行先出力ポート・アドレスを格納する。行先出力ポートは、ルーティング・テーブルの第0エントリの指定したものから、 $0 \rightarrow 1 \rightarrow \dots \rightarrow r \rightarrow 0 \rightarrow 1 \rightarrow \dots$ と、接続が可能となるまで候補を順番に試みていく。

本方式は制御が簡単に柔軟性が大きい反面、アドレスコンバータ・メモリの容量が大きくなり必要になる欠点をもつ。前節で述べた(ii)方式を拡張すれば、比較的少い付加ハードウェアで可変ルーティングが行えるようになる可能性があり、検討中である。また、アドレスコンバータ・メモリを入力ポートごとに置くのではなく、SUに一つ設ける方式も考えられる。これは、網上の交通量が比較的小さい場合や、パケット長が大きい場合など、ルーティング・テーブルへのアクセス競合があまり起こらないときに有効である。

本方式では、ルーティング・テーブルに登録する順番が、すなわち経路の優先順位であり、これを決定することは重要である。(3.2参照)なお、行先の交通量によってルーティング・テーブルを書き換える、いわゆる適応型のルーティングは、今回の検討の対象としなかった。

## 2.5 出力ポートの構成

出力ポートは、各入力ポート・コントローラから発せられたREQ信号の一つを選択するアービタと、アービタからのセレクト信号により選択された入力ポートと当該出力ポートの信号線を接続するセレクトから成る。(図1)

アービタは、複数のREQ信号のうちの一つのみをランダムに選択するものとした。アービタの回路としては、

- (1) R-Sラッチを用いた方式
- (2) リングカウンタを用いた方式

の二つを検討した。

前者は、 $m$ 入力のNANDゲートをR-Sフリップフロップの形に接続するもので、きわめて簡単な回路で実現される。この回路を試験的

に製作し、動作を調べたところ、配線状態などによって選抜に優先のある傾向がみられたが、論理動作には問題がなかった。

後者は、自己補正型の $m$ ビットリングカウンタを用いてREQ信号をスキッピングする方式で、(1)より安定した動作を示すが、遅延時間が大きくなる危険がある。この問題には、カウンタ用のクロックを高速化して対処する。

今回は、主として設計の容易さの点から、(1)の方式を採用した。

## 2.6 SUの動作速度

当SUの動作速度をクロック数で示す。

- (1) 行先アドレスを知って出力ポートを決め、アービタレージョンを行って経路を設定し、実際の転送を始めるまでの制御時間 : 3クロック
- (2) 1回のルート変更に必要な時間 : 2クロック
- (3) 実際にFIFOメモリ間で1語の転送を行うのに必要な時間 : 1クロック

さらに、アービタの動作安定のため、1パケットの転送が終了したのちセレクト信号を出すのを1クロック遅らせている。

クロック周波数を決定する要因として、

- (i) アービタの遅延
- (ii) アドレスコンバータ・メモリのアクセス時間
- (iii) FIFOメモリのアクセス時間

があげられる。中でも(ii)が決定的であり、SUのクロックの周期は、FIFOメモリのアクセス時間の2倍以上になるように設定されねばならない。

## 3. 当SUを用いた相互結合網

### 3.1 当SUを用いた結合網の特徴

当SUは、各機能の独立性・並列性が高く、相互結合網を構成した場合、全体としても高速動作が可能である。パイプライン転送を行う蓋種交換方式であるため、データは小さなパケットを単位として転送される。ルーティングが表引きによって行われるため、適用可能な結合網の形状もさまざまであり、デルタ網・ガンマ網などの多段結合網、格子型網・超立方体網・C

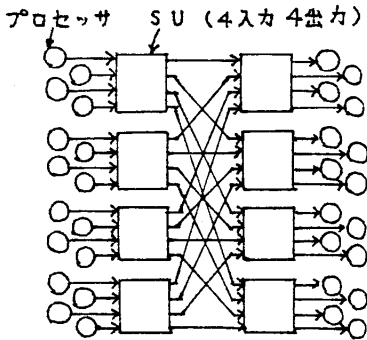


図4. 本SUを用いた多段結合網

CC網・木状網・有弦環網などに有効である。典型的な例として、図4に4×4のSUを用いて多段結合網を構成したものを、図5に5×5のSUを用いて格子型結合網を構成したものを示した。なお、本方式は網の拡張、変更、障害回復にも有利である。

### 3.2 ルーティング・テーブルの構成法

ルーティング・テーブルに登録する出力ポート・アドレスの優先順位を、与えられた網の形状と使用環境に応じて決定することは、転送効率の点から重要である(2.4参照)。その際考慮すべき諸点を列挙する。

- (1)行先のトラヒック : トラヒックを相互結合網全体にわたって均等に分布させるようにする。
- (2)SU内における閉塞 : 当該SUの出力ポートにおける競合をなるべく小さくする。
- (3)後の中継回数 : 今後の中継回数が少なくて済むものを優先する。
- (4)後の経路数 : 後の経路になるべく多くの選択余地を残すほうが有利な場合が多い。
- (5)デッドロック防止 : ルーティングを制限して間接ストアアンドフォワード・デッドロックを回避する方式が考えられる。(2)

## 4. 結合網のシミュレーション評価

以上で述べたSUの動作確認、当SUを用いた相互結合網の性能評価、ルーティング・テー

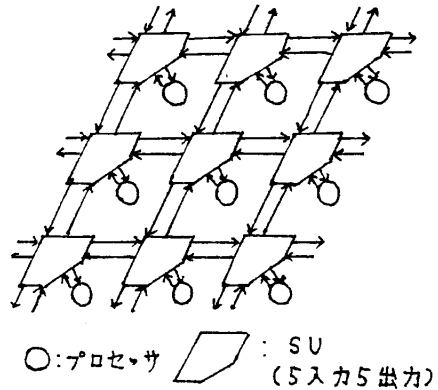


図5. 本SUを用いた格子型結合網

ブル構成法の具体的な検討などを目的とし、機能レベルシミュレータを作製した。シミュレータは、主として動作速度の点からFORTRAN 77で書かれており、M280Hの上で動作する。

最初に1SUのシミュレーションを行い、正しい論理動作をすることを確認した。続いて、次節に述べるモデルにもとづき、2種類の多段結合網に当SUを適用した場合の遅延とスループットのシミュレーション評価を行った。結合網は、オメガ網<sup>(4)</sup>とこれを改良した網、ガンマ網<sup>(5)</sup>を用いた。なお今回は、格子型網・超立方体網などの静的な網<sup>(6)</sup>は考慮の対象としなかった。

### 4.1 シミュレーション・モデル

シミュレーション・モデルは、N台の入力モジュールとN台の出力モジュールを相互結合網で結合した構成である。モデルの詳細を記す目的で、次にシミュレーションのパラメータを列挙し説明する。

- (1)網の形状
- (2)網の大きさN : 網の入力側・出力側のそれぞれに接続するモジュール台数
- (3)SUのポート数s
- (4)データ生成率 $\tau$  : 入力モジュールにおいて1クロックあたりデータ1語が生成される確率。本シミュレーションでは、パケットの到着はポアソン分布とした。
- (5)パケット長p : 1パケットに含まれる語

数。パケット長は最小2語とし、与えられた最大値までの範囲の値を一様乱数で与える。 $P$ はこの平均値である。

- (6) キュー長  $Q$  :  $SU$ の入力ポートFIFOメモリの語数。なお、入力モジュールには無限長のFIFOメモリがあり、出力モジュールは必ずデータを受理すると仮定した。
- (7) エントリ数  $e$  : 図3に示すルーティングテーブルのコラム数。すなわち、可能な行先出力ポート数の最大値である。 $e = 1$ のときが固定ルーティングに相当する。
- (8) 行先出力モジュールの決定方式 : 特にことわらないかぎり、0から $N-1$ までの整数値を一様乱数で与えることにする。
- (9) ルーティング・テーブルの構成法  
求めるデータは、主として次の二つである。
  - (1) 正規化スループット :  $1$ クロックの間に網の各出力ポートから出る語数の平均値をポートあたりで計算したもの。システム全体が定常状態になれば、正規化スループットは $\rho$ に等しくなる。なお、網が定常状態を保つための最大値を限界スループットと呼ぶ。
  - (2) 遅延 : パケットの結合網におけるクロック単位の転送遅延。すなわち、データが入力モジュールにおいて生成された時点から計算し、出力モジュールに到着するまでの時間である。

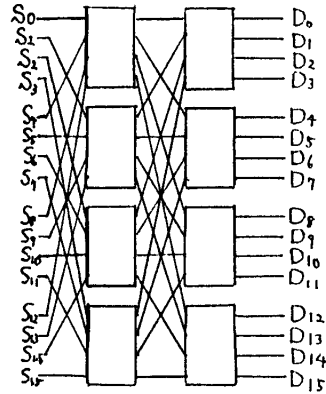


図6 オメガ網 ( $N=16, S=4$ )

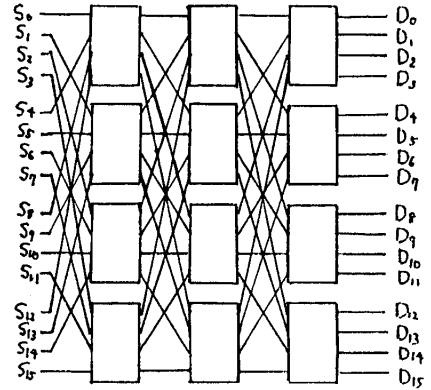


図7 オメガ網の冗長構成 ( $N=16, s=4$ )

4.2 オメガ網とその拡張について  
当 $SU$ をオメガ網<sup>(6)</sup>に適用した場合のシミュレーション評価に関して述べる。オメガ網の構成は図6に示すものであり、入出力4ポートずつの $SU$ を用いることにする( $s=4$ )。

オメガ網は1ルートの網であるが、これに1段分( $N/4$ 個)の $SU$ を追加した構造の網(図7)に關しても同様のシミュレーションを行い、比較検討した。後者は、一組の入出力モジュールの結合に關して、4本の経路を設けることが可能である。

シミュレーションでは、 $p=9$ (2語から16語までの一様乱数)、 $l=45$ (平均して5パケット分)の条件下で、他のパラメータとスループット・遅延の關係を求めた。ルーティングテーブルの構成法は、冗長構成をとった場合の初段の $SU$ でのみ問題になるが、今回はシミュ

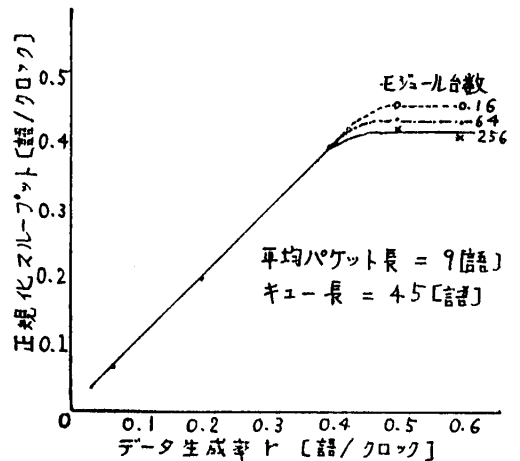


図8 オメガ網の正規化スループット

レータの容易さから、その内容を一樣乱数で与えた。

結果として、図8に網のスループットに関するものを、図9に網の遅延に関するものを示す。図は両者ともオメガ網に関するものであり、後に述べる理由から、その冗長構成に関するものは記さなかった。

最初に網の限界スループットに関して述べる。網の大きさが大きくなっていったとき、各段での閉塞により限界スループットは減少する。ただし、網の段数をさらに増していくときのスループットの減少幅は小さくなる。

次に遅延に関しても、網の大きさの増大につれて大幅に増大することが見てとれる。これは、遅延が、全段のFIFOメモリでの待ちと、各段ごとの閉塞による待ちの両方に依存するからである。

一方、オメガ網を1段分拡張した場合であるが、この場合、初段のみで多ルート化がなされている(2段目以後はオメガ網と同じ)。しかし、以後の段のトラヒックが均等になることから、オメガ網に比較して転送効率は上がらず、スループットに関するグラフは図8にほぼ重なり、遅延は増加した1段分だけ増すことになる。

ただし、この結果は先行モジュールが一樣乱数で与えられる場合のみに関してである。本方式の網の拡張は、モジュール間の結合に局所性があるとき有利であり、別のシミュレーションの結果、場合によっては20%近くのスループットの向上があることが示された。

多ルート化の効果は、スループットなどの性能面のみならず、信頼性の面において顕著である。すなわち、網の一部が故障しても、通信性能が低下するだけで、故障による経路の断絶が避けられる。

次にFIFOメモリのキュー長を変化させた時の、キュー長と限界スループットの関係を別のシミュレーションで調べたところ、キュー長の増大とともに限界スループットが増大するが、この傾向はしだいに頭打ちになることがわかった。たとえば、キュー長45語(平均5パケット分)のときの限界スループットは、キュー長2語のときより約60%大きい。スループットが飽和するのは、ほぼキュー長90語(平均10パケット分)のところである。

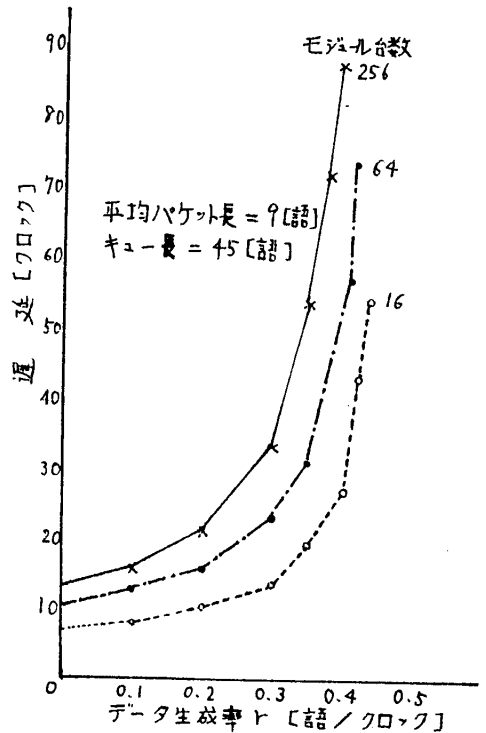


図9 オメガ網における遅延

さらに、パケット長を固定したときは、キュー長によらず、可変長のときにくらべて約0.02から0.03語/クロックだけ限界スループットが高くなることがわかった。

#### 4.3 ガンマ網について

前節のオメガ網の拡張構成では、ルートの選択が第1段においてのみ可能であった。本節では、可変ルーティング機能をほとんどすべての段で用いることの可能な多段結合網であるガンマ網(図10)を対象に、シミュレーション評価を行った結果を述べる。

ガンマ網は次のような性質を持っている。<sup>(4)</sup>

- (1)  $(\log_2 N + 1)$  段網である。
- (2) 第  $m$  ( $0 \leq m \leq \log_2 N - 1$ ) 段の  $k$  ( $0 \leq k \leq N - 1$ ) 番目のSUの  $i$  ( $i = 0, \pm 1$ ) 番目の出力ポートは、第  $m + 1$  段の  $k'$  番目のSUの  $j$  番目の入力ポートに接続される。ただし、 $j, k'$  は以下の式で表わされるとする。

$$j = -i$$

$$k' = (k + i \times 2^m) \bmod N$$

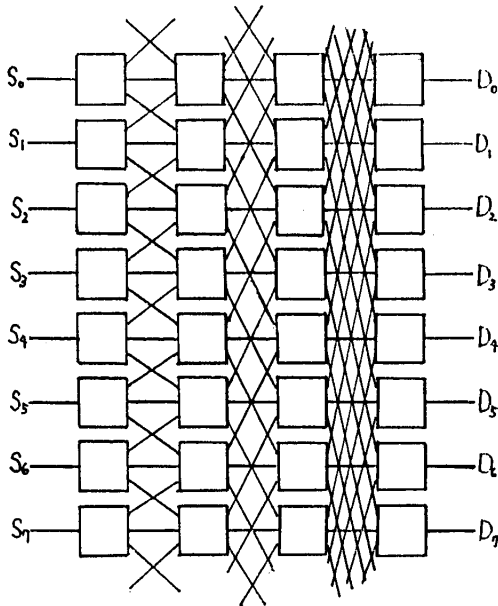


図10 ガンマ網 (N = 8, s = 3)

(3)ルーティングは以下のように行う。すなわち、行先モジュールのアドレスとソースモジュールのアドレスの差を、 $\{1, 0, \bar{1}\}$ を用いた冗長2進表現( $T = -1$ )によって表し、これをタグとして、網の各段において1桁ずつ参照していく。値が0となる段でのルートは唯一( $i = 0$ )で、その他のときは1と $\bar{1}$ の両方についての選択が可能であるが、以後の経路はこの選択によって変わる。

シミュレーションでは、 $p = 9$ の条件下で、他のパラメータとスループット・遅延の関係を調べた。ルーティング方式は以下の3種を比較した。

(方式1)タグは0と1しか含まない普通の2進数のみとし、 $i = -1$ の出力ポートを使わない固定式ルーティング。

(方式2)1および $\bar{1}$ の出力ポートを均等に使用する固定式ルーティング。

(方式3)タグの該当する桁が1または $\bar{1}$ の場合、行先の経路数が多く残されるような出力ポートを優先的に選択し、これが利用不可能のときにもう一方のルートを選択する可変ルーティング。

結果として図11にスループットに関するものを、図12に遅延に関するものを示す。

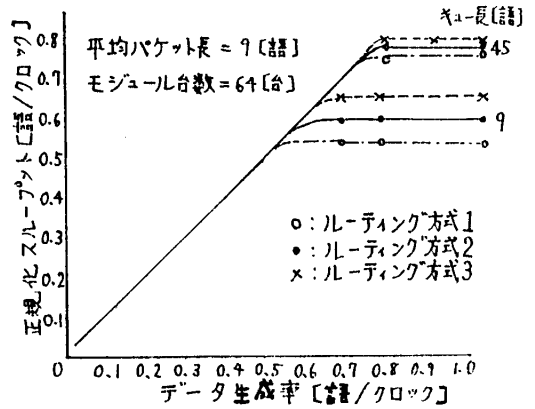


図11 ガンマ網の正規化スループット

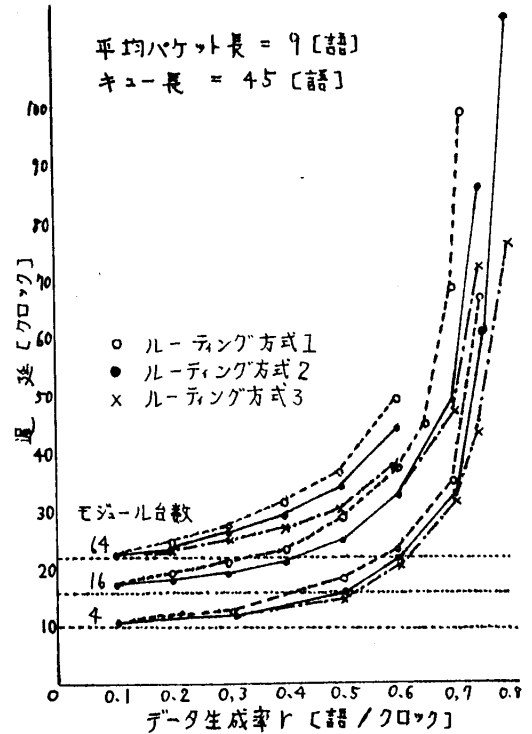


図12 ガンマ網における遅延

図より、ガンマ網ではルーティング方式によって性能に顕著な差が生じることが見てとれる。特に、可変ルーティングの効果が大きい。

限界スループットは、FIFOメモリのキュー長の増大により大幅に増大し、キュー長45語のときはキュー長9語のときの40%程度の増加がみられることがある。また、ルーティン



グ方式による限界スループットの差は、キュー長が短いほど大きい。(図11)

ルーティング方式による遅延の差は、データ生成率の増加とともに顕著になる。(図12)

前節で述べたオメガ網と比較して、ガンマ網は限界スループットが高く、転送遅延もきわめて小さい。

## 5. 検討・考察

### 5.1 1SUの転送性能

2.6で述べた当SUの動作速度より、次の三つの時間をクロック数で求める。

(1)データ転送遅延  $D_d$  : 長さ  $p$  のパケットの先頭がSUの入力ポートに到着してから、 $a$ 回のルート変更ののち、そのSUの出力ポートから出てくるまでの最短時間。

(2)パケット通過時間  $P_p$  : 同じ時刻から、パケットが出終わるまでの最短時間。

(3)パケット転送時間  $P_t$  :  $s$ 段のSUを経由し、総計  $A$ 回のルート変更を行ってパケットが転送されるのに要する最短時間

(1)(2)(3)は以下の式で表わされる。

$$D_d = 2a + 3$$

$$P_p = p + 2a + 4$$

$$P_t = 3s + p + 2A + 1$$

### 5.2 当SUを用いた相互結合網の性能

4章のシミュレーション結果から、各網の具体的な転送性能を論じる。

#### (1)スループット

限界スループット(正規化したもの)に関して検討する。本SUを用いた場合、理想的な限界スループットは、1クロックあたり(1- $\alpha$ )語である( $\alpha$ は制御オーバーヘッド)。実際には結合網内の閉塞が限界スループットを著しく低

下させる。

モジュール数256のオメガ網の場合、図8の条件下では限界スループットが約0.42語/クロックであり、これは理想的な場合の約60%である。FIFOメモリとして、アクセス速度が50nsec程度のRAMを用いた場合、1ポートあたり約4MB/sec程度の限界スループットである。網の大きさをさらに大きくしていけば、限界スループットはさらに減少するが、減少幅はしだいに小さくなる。たとえばモジュール数1000の結合網でも、0.4語/クロック程度の限界スループットが予想される。これは、先のFIFOメモリを用いた場合、網全体のスループットが約4GB/secとなることを示す。

モジュール数64のガンマ網(キュー長45)の場合、図11の条件下では、可変ルーティングを行ったとき、約0.8語/クロックの限界スループットを示す。これは理想的な場合の90%近い性能であり、ガンマ網では閉塞による性能低下がきわめて小さいことがわかる。また、ガンマ網の場合、網の大きさによる限界スループットの変化がほとんどないことから、先のRAMを用いた場合、モジュール数1000の網でも、8MB/sec程度のスループットを示すことが予想される。

#### (2)転送時間

モジュール数256のオメガ網の場合、図9より、データ生成率が0.2のときの遅延は約21クロック(理想的な場合の約1.6倍)となることを見てとれる。いま、10MHzのクロックを用いたとすれば、このとき9語の転送に約3 $\mu$ secの時間を要することになる。さらに、データ生成率が0.3になると、遅延が急激に増加する。

モジュール数64のガンマ網(可変ルーティング)の場合、図12より、データ生成率0.2

表1 当SUのハードウェア量

	1 Port		5 Ports	
Number of Gates in IP	840	(670)	4200	(3350)
FIFO - RAM in IP	2.30kbits	(2.30kbits)	11.5kbits	(11.5kbits)
Routing Table in IP	768(r+1)bits	(768bits)	3.84(r+1)kbits	(3.84kbits)
Number of Gates in OP	110	(120)	550	(600)

(x) : x is the value of [1]  
r+1 : number of entries

のときの遅延は約23クロックであり、先と同じ仮定で、約3.2 $\mu$ secでパケット転送ができる。データ生成率が0.7をこえたところで遅延が急激に増加する。

### (3) 網の信頼性

高並列計算機の相互結合網では、システムが大きくなるほど故障の発生する確率も高くなる。したがって、結合網もある程度故障に対処できるように設計されるべきである。故障はSU間のリンクとSU内部に起こる可能性がある。4.2で述べたオメガ網の拡張構成や、4.3で述べたガンマ網のように、網に冗長なルートを設け、可変ルーティングを行うことによって障害部分の迂回が可能になり、相互結合網の信頼性も向上する。

### 5.3 LSI化の検討

表1に、TTLを用いて設計した本SUのハードウェア量を示す。この他にも、内部バスの総配線数(5入力5出力で80本)、内部バス-OP間の総接続線数(同300本)を考慮に入れなければならない。

いま、ポート数5のSUのLSI化を考えると、総ゲート数は5000に不足し、必要なピン数は114本である。その場合、内部配線数とルーティング・テーブル用のRAMがやや大きい。十分に1チップ化が可能と思われる。

## 6. おわりに

可変ルーティングを行う蓄積交換方式の汎用スイッチング・ユニットの設計と、これを用いた多段結合網の転送シミュレーションに関して述べた。残された課題として、ルーティング方式の再検討(ルーティング情報の縮退化など)、静的な結合網の転送シミュレーションと間接ストアアンドフォワード・デッドロック対策の検討、障害対策の詳細化、非同期方式の検討、放送機能・パケット整形機能の付加、適応型ルーティングの検討などがあげられ、また本SUの試作、より大規模な網のシミュレーション評価なども今後行う予定である。

- (1) 坂井, 喜連川, 田中, 元岡: "データベースマシンGRACEに於けるモジュール間結合網", 信学技法, EC 83-14 (1983-06).
- (2) 前, 服部, 坂井, 田中, 元岡: "プロセッサ間結合網に於けるスイッチング・ユニットの試作と評価", 第27回情報学大会, 5N-1, (1983-10).
- (3) 坂井, 田中, 元岡: "PIEのゴールフレーム分配網とコマンド通信網", 第28回情報学大会, 6F-7, (1984-03).
- (4) 坂井, 計, 田中, 元岡: "可変ルーティングを行う分散制御スイッチング・ユニット", 第28回情報学大会, 1F-2, (1984-03).
- (5) Feng, T.: "A Survey of Interconnection Networks", IEEE COMP UTER, 14, 12, pp.12-27 (1981-12).
- (6) Lawrie, D. H.: "Access and Alignment of Data in an Array Processor", IEEE Trans. Comput. C-24, 12, pp.1145-1155 (Dec. 1975).
- (7) Patel, J. H.: "Performance of Processor-Memory Interconnections for Multiprocessors", IEEE Trans. Comput. C-30, 10, pp.771-780 (Oct. 1981).
- (8) Dias, D. M., Jump, J. R.: "Analysis and Simulations of Buffered Delta Networks", IEEE Trans. Comput. C-30, 4, pp.273-282 (Apr. 1981).
- (9) Adams, G. B., Siegel, H. J.: "The Extra Stage Cube: A Fault-Tolerant Interconnection Network for Supersystems", IEEE Trans. Comput. C-31, 5, pp.431-454 (May 1982).
- (10) Gajski, D. et al.: "Cedar - A Large Scale Multiprocessor", Proc. of the 1983 Int'l Conf. on Parallel Processing, pp.524-529 (Aug. 1983).
- (11) Parker, D. S., Raghavendra, C. S.: "The Gamma Network: A Multiprocessor Interconnection Network with Redundant Paths", The 9th Ann. Symp. on Comput. Arch., pp.73-80 (Apr. 1982).