

マルチプロセッサ・システムのフォールトトレラント結合方式について

田中 徳彦

黒川 恭一

古賀 義亮

防衛大学校 情報工学科

マルチプロセッサ・システムのプロセッサ間結合方式について、パイプライン並列処理を例として相互結合ネットワークとその結合回路（バンクメモリ結合回路）を提案し、フォールトトレラント・パイプライン処理システムを構成する。バンクメモリ結合回路によるマルチプロセッサ・システムは、結合しているプロセッサの機能障害を検出することができるため、フォールトトレラント化に有利な結合回路であるばかりでなく、パイプライン処理のように大量のデータを一方向に転送するアプリケーションに対して従来の結合回路の問題点を改善する。

A Method of Fault Tolerant Communication in a Multi-processor System

Norihiko TANAKA Takakazu KUROKAWA Yoshiaki KOGA

Department of Computer Science, The National Defense Academy

This paper proposes a new fault tolerant communication method and three new interconnection networks to construct a multi-processor system for a pipeline processing. The proposed communication scheme using bank memory switching technique has an advantage to make a fault tolerant pipeline system so that it can detect any failure caused in processing elements of the system. In addition, it can overcome the conventional problems caused in interconnection circuits to flow data with one way direction such as a pipeline processing.

## 1 まえがき

1つの処理を複数の処理ステージに分けて、流れ作業の原理で処理を行なうパイプライン処理は、 $n$ 台の処理要素（PE）を用いて処理速度を最大 $n$ 倍に向上することができる利点を有すると共に、現在までに多数提案されている逐次的アルゴリズムに比較的近いことなどから、低コストの並列処理機構として多くのシステムが提案され<sup>[1]</sup>汎用化システムの試みもなされている<sup>[2]</sup>。しかし、PEを直線状に結合した線形結合マルチプロセッサ・システム上で汎用パイプライン処理を実現する場合、PEや結合回路の故障によって容易にその処理能力が失われる。このためPEをネットワーク状に冗長結合して、故障が生じた場合には故障部位を切り離してシステムを再構成し、システム全体の機能を維持するフォールトトレラント化<sup>[3]</sup>が必要となる。本報告は、パイプライン処理を行なうマルチ・プロセッサシステムのフォールトトレラント化を考慮した結合方式として、2重リング結合、3重リング結合、木状リング結合の3種類の結合方式を提案し、これらの特徴を比較する。2.では、パイプライン処理を行なうマルチプロセッサ・システムのPE間結合について、従来から提案されて来た各方式の問題点を明確にした後、新しい結合方式を提案する。3.では、フォールトトレラント・パイプライン処理のネットワーク結合方式を提案し、3つのネットワークポロジについてそれらの特徴を比較する。

## 2 プロセッサの結合方式

### 2.1 従来の結合方式

パイプライン処理を行なうマルチプロセッサ・システムでは、処理速度の向上を図るために、各PEの仕事量（タスク粒度）を均一にしてPE間でのデータ待ち時間をできる限り少なくする必要がある。このため、非同期のパイプライン処理におけるデータの同期法として、結合回路にFIFOバッファや2ポートRAMのような記憶装置を用いる方法等が提案されている<sup>[4]</sup>。このような結合回路は、内部に含まれる記憶装置の形態によって特徴付けることができる。ここでは、この記憶装置として代表的なFIFOバッファと2ポートRAMを取り上げ、それら

の結合回路としての問題点とそのフォールトトレランス性について明確にする。

#### 2.1.1 FIFOバッファ

結合回路がFIFOバッファによって構成されている場合、送信PEがデータをFIFOバッファに書き込んだ後、今度は受信PEがFIFOバッファから読み出さなければならないため、PEとFIFOバッファ間で処理に無関係なデータ転送が2回必要となり、データ転送のオーバーヘッドが大きくなる。また、FIFOバッファの容量は比較的小さいので多くのデータを送る場合、ハンドシェイク待ち時間が生じてしまう。

フォールトトレランス性に関しては、FIFOバッファは基本的に単方向の通信機能しか有していないため、送信PEはデータを送り出すのみで受信のチェックをすることができない。このため、転送先PEの故障にも関わらずデータを送り続けるような事態が生じる。また、システムの再構成を行なう場合、双方向の通信機能が必要であることから、2つのFIFOバッファが必要となると共に、再構成を行なうためのメッセージ到着時間がバッファによって不確定になるため、確実な誤り回復が難しくなる。

#### 2.1.2 2ポートRAM

結合回路が2ポートRAMによって構成されている場合においてもFIFOバッファと同様に、PEと結合回路との間で2回の転送が必要である。しかし、2ポートRAMはFIFOバッファよりも大きなメモリ容量を実現することができるため、ハンドシェイクの待ち時間はFIFOバッファの場合よりも減少する。

フォールトトレランス性に関しては、デュアルポートRAMは双方向から自由に読み書きできるため、ソフトウェアによって同期を取ることが容易であり再構成に有利といえるが、逆の面ではハードウェアのサポートがないためにソフトウェアのバグ等で隣接PEに障害が波及する可能性がある。例えばデータの受信完了前にデータが新たに書き換えられる等の危険性を常に持っている。

## 2.2 バンクメモリ結合方式

従来からの結合方式の問題点を改善するために、本報告ではバンクメモリ結合方式を提案する。バンクメモリによる結合方式は、2つのPEを結合するために、PEに接続されているメモリブロック自体を交換することによってプロセッサ間でデータを交換する結合方式である。また、PE同士がスイッチによって分離されているため、PEの機能障害がシステム全体に波及しにくい。更に、結合回路の主要な部分は誤り検出機能を有しており、フォールトトレラントマルチプロセッサ・システムの構成に有利である。

### 2.2.1 バンクメモリ結合方式の概要

バンクメモリによるPEの結合方式を、2ポートPEの線形結合網を例として図1に示す。バンクメモリ結合回路は、2つのメモリブロックMB0、MB1とメモリブロック切り換え回路とから構成されている。プロセッサPE1、PE2は、メモリブロック切り換え回路によってメモリブロックMB0、MB1にそれぞれ接続されている。メモリブロック切り換え回路は2つの状態を持っており、図1の例ではPE1と左のMB0、PE2と2つのMB1、PE3と右のMB0とを接続している。従って、PE1が左のMB0に、PE2が左のMB1にデータを書き込んだ後、そのバンクメモリのスイッチを切り換えるとMB0とMB1が入れ換わり、PE1とPE2間でデータを交換する双方向リンクが実現できる。

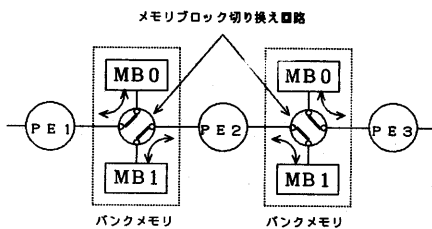


図1 バンクメモリによるPEの結合

### 2.2.2 PEの構成

PEは図2に示すように、CPU、ROM、RAM、I/Oポートとそれら結び付けるシステムバ

スとから構成されている。PEのシステムバスからバッファゲートによって分離されているバンクメモリ結合バス（結合バス）には、CPUの制御により結合バス上でDMA転送を行うDMAC (Direct Memory Access Controller) とバンクメモリとが接続されている。通常動作中バッファゲートは導通状態となっているが、結合バスでのDMA転送中は、ハイインピーダンス状態となって回路を遮断する。ネットワーク結合方式においては、非隣接プロセッサ間の通信のために、中間に位置するプロセッサの中継が必要であり、一般的には、この中継通信が並列処理の効率を低下させるオーバーヘッドとなっている。この問題点に対して本構成では、バッファゲートによってPE内バスと分離できる結合バス内でDMA転送を行うことによって、PEに中継通信のための負荷を負わせずに高速な通信を行なうことができる。

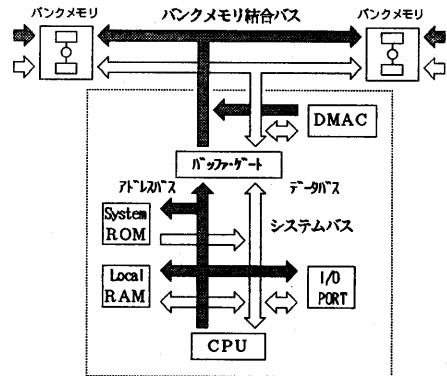


図2 ノードPEの構成

### 2.2.3 多ポートPEの構成

PEの多ポート化は、バンクメモリ結合バスにバンクメモリを増設して、各CPUのメモリ空間内に新たに接続したバンクメモリのための領域を確保することで可能となる。3ポート以上のPEの場合、結合バスにブロードキャストモードを付加し、複数のバンクメモリに同時にメッセージを書き込むと共に、DMA転送によって1つのバンクメモリから複数のバンクメモリに対して同時にメッセージを転送することが可能となる。3ポートのノードPEの構成法を図3に示す。

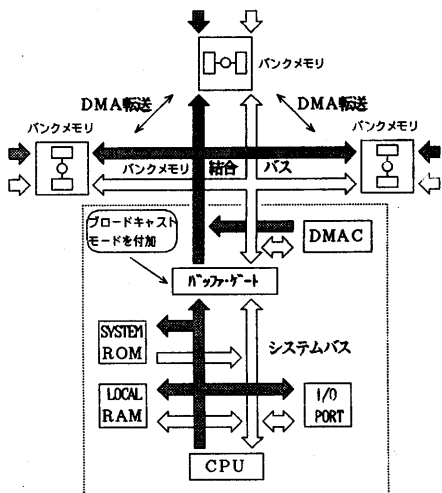


図3 3ポートノードPEの構成

### 2.2.4 メモリブロックの切り換え

メモリブロックの切り換えは、コレスポンデンスオペレータ (COP) とプリセットリミットタイム (PLT) とを用いて、結合されている両方のPEの合意のもとに行なわれる。COPは双安定出力を持っており、図4に示すように2つの入力の論理が一致したときのみ出力を変化させる。また、PLTはPE間のデータ転送時間の上限を規定して監視する。転送制限時間からの逸脱は、PEまたは結合回路の障害としてPLTが検出し故障PEを診断する。メモリブロック切り換えの同期機構は図5に示すように3線ハンドシェイクを用いており、双方の切り換え要求とCOPによる切り換え結果とをソフトウェア、ハードウェア双方によって相互に確認しながら確実に行う。

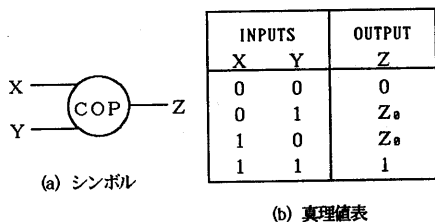


図4 COP (Correspondence OPerator)

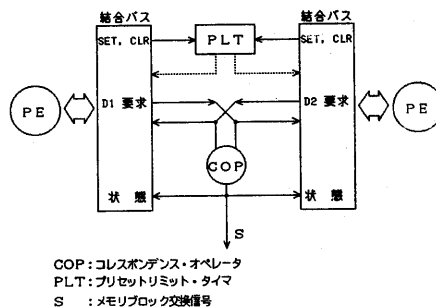


図5 メモリブロック切り換えの合意機構

## 2.3 従来の結合方式とバンクメモリ結合方式との比較

従来の結合方式では記憶装置が結合プロセッサで共有されているため、データを転送する場合、記憶装置への読み書きの動作が転送速度を低下させる原因となっていた。これに対して、バンクメモリは見かけ上隣接PEで共有されているが、実際はおのおの別個の2つのメモリが独立して入出力に使われているため、結合回路上のデータを直接PEで使用することができる。特にパイプライン処理のように、大量のデータを流しながら処理を行なう場合、結合回路の転送速度の差が問題となりバンクメモリによる結合方式が有利となる。

フォールトトレランス性に関しては、従来の方式では、転送されてきたデータの正当性が時間的にも量的にも厳密に保証されていないのに対して、バンクメモリは入出力のためのメモリが分離されているため、入出力の動作が互いに他に影響することなく安全であり、時間管理の同期機構をも備えていることから、転送データの正当性を保証する事が出来る。

## 3 フォールトトレラント・パイプライン処理方式

フォールトトレラントシステムを実現するための1つの方法として、システムにあらかじめ冗長な部位を付加しておき、故障が発生した場合はその故障部位を診断して置き換える方法がある。このためには、故障の検出、故障の診断、システムの再構成、誤りからの回復の各処置をオンライン中に自動的に

なうシステム構成やアルゴリズムが必要となる。

以下では、まず3.1で本報告における議論の対象を明確にするため、ネットワークの前提条件、故障モデル、ネットワークの再構成法の方針をそれぞれ示した後、3.2で3種類のネットワークを順に提案してその特徴を比較する。今回の提案では、3ポートのPEによるネットワークを考える。PEのポート数はネットワークの通信距離に関係しており、多いほどPE間の通信距離が短くなるが、あまり多くなるとPEに対して結合回路が指数関数的に多くなると共に制御ソフトウェアも複雑になり、結合回路の信頼性の低下によるシステム全体の信頼性の低下が問題となる。その中で3ポートは、多様なネットワークを構成することができる最小の次数であり、パイプライン処理のような全順序関係を持つ一次元的な処理には最低2ポートが必要であることから、最小の冗長構成であると言える。

### 3.1 準備事項

#### 3.1.1 ネットワークの前提条件

- 全てのPEは均質で、次数は3（3ポートのPE）であり、バンクメモリによってPE間を直接結合する。
- パイプライン処理はPE数  $n$  ( $n$  は偶数) のリングネットワーク上で行なう。リング上の0番のPEが入出力を担当しており、PE1から番号順にパイプラインの各ステージが実行されて  $n-1$  ステージで終了し、出力データをPE0に送る。このリングネットワークを基本リングと呼び、基本リングを構成するPEを基本リングPEと呼ぶ。
- 冗長プロセッサRPを基本リングPEの冗長なリンクに結合し冗長ネットワークを構成する。

#### 3.1.2 故障モデル

パイプライン処理を行なうマルチプロセッサ・システムの故障モデルとしては、PEとそのリンクを対象とし、PEが出力を出さず入力も受け付けなくなる沈黙の故障を考える。マルチプロセッサ・システムの一般的な故障モデルでは、故障PEはどのような動作も行ない得るとしているが、本報告で対象

とする故障モデルは、PEの入出力に対して厳重なプロトコルを設けることによって妥当性を持つと考えられる。また、このような沈黙の故障の仮定はPEの故障によってシステム全体に障害が波及することがなく、隣接PEから容易に故障検出・診断ができる故障モデルでもある。

#### 3.1.3 ネットワークの再構成法

故障が発生した後、システムの再構成を行ない誤りを回復するため、手段としては故障PEを検出・診断して切り離し、予備PEを新たに割り当てることになる。しかし、PEの故障が発生したときパイプライン中にはまだ汚染されていない未処理のデータ列があるため、迅速な誤り回復のためには正常なPEは初期の接続状態を維持していた方が、全PEの接続を変更して最適な再構成を行なうよりも効率が良いと考えられる。このため、故障PEのみを代替PEに最短経路で置き換えるものとする。

### 3.2 結合ネットワーク

ここではリング結合を基礎として、3.1.3で示したネットワークの再構成法に合致する結合ネットワークとして、2重リング結合、3重リング結合、木状リング結合の3種類を提案し、ネットワークの構成方法と再構成アルゴリズムを示す。以下の説明において、 $P_i$  は  $i$  によって番号付けられた  $n$  個の基本リングPEの1つであり、正常動作中におけるパイプライン処理の主体である。また、 $R_i$  は基本リングPEに結合する冗長プロセッサの1つであり、基本リングPEの代替プロセッサとする。

[定義1]

プロセッサ  $P_i, R_i$  の結合を  $(P_i, R_i)$  と表現する。

[定義2]

番号付けられたプロセッサ  $P_i (i = 0, 1, \dots, n-1)$  が番号順にリング結合したものを  $\langle P_i \mid i = 0, \dots, n-1 \rangle$  と表現する。

[定義3]

$m+1$  ビットの2進数  $n$  を  $(1, n_{m-1}, \dots, n_1, n_0)$  としたとき、置換  $\sigma$  を次のように定義する。

$$\sigma(1, n_{m-1}, \dots, n_1, n_0) = (0, n_{m-2}, \dots, n_0, n_{m-1})$$

また、逆置換  $\sigma^{-1}$  を次のように定義する。

$$\sigma^{-1}(0, n_{m-2}, \dots, n_0, n_{m-1}) = (1, n_{m-1}, \dots, n_1, n_0)$$

置換は、 $n$  の下位  $m$  ビットをローテイトシフトし、MSBをビット反転したものである。

### 3.2.1 2重リング結合

基本リングPEと冗長PEの数を同数とし、基本リングPEと冗長PEとで同じ処理を同時に実行しておけば、誤りが検出され次第冗長PEを基本リングPEと置き換えることができるので、誤りを速やかに回復することができる。この方式はホットスタンバイ方式と呼ばれており迅速な機能回復ができる。

#### [構成法]

基本リング  $\langle P_i \mid i = 0, \dots, n-1 \rangle$  と冗長プロセッサによる第2リング  $\langle R_i \mid i = 0, \dots, n-1 \rangle$  の間に  $\{(P_i, R_i) \mid i = 0, 1, \dots, n-1\}$  の結合を行なう。 $n = 8$  の時の例を図6に示す。図6の円はプロセッサを表わしており、2重の円は基本リングプロセッサ  $P_i (i = 0, 1, \dots, 7)$  を、また単円のプロセッサは冗長プロセッサ  $R_i (i = 0, 1, \dots, 7)$  をそれぞれ表わしている。円の中の数字は初期動作中そのプロセッサが実行するパイプライン処理のステージ番号である。第0ステージは入出力のタスクであり、パイプライン処理は1から7の7ステージで実行される。

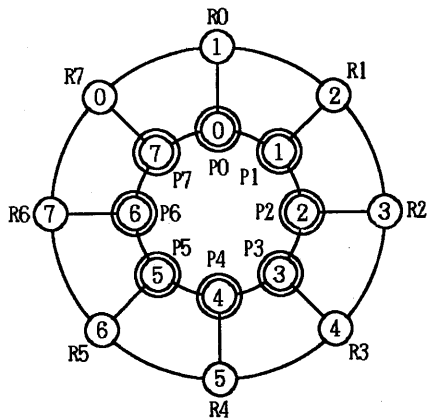


図6 2重リング結合 ( $n = 8$ )

#### [パイプライン処理]

パイプライン処理は、 $P_0$  のプロセッサが入力データをパイプに投入できる状態に準備して、 $P_1$  へデータ転送を始めることで開始され、 $P_{n-1}$  のプロセッサで最終ステージが終わり  $P_0$  プロセッサへ処理結果を送ることで終了する。処理の途中、プロセッサ  $P_i$  は

プロセッサ  $P_{i+1}$  へデータを送ると共に、プロセッサ  $R_i$  にも同じデータを送る。プロセッサ  $R_i$  は常に  $P_{i+1}$  プロセッサと同じ処理を行っているが、再構成の指示がない場合はその結果を捨て、新たなデータを  $P_i$  から受け付ける。

#### [再構成アルゴリズム]

基本リングPE中、任意の1つの基本リングプロセッサ  $P_i$  が故障した場合、故障プロセッサに接続されている前後の基本リングプロセッサ  $P_{i-1}, P_{i+1}$  と冗長プロセッサ  $R_i$  が故障を検出する。この位置関係を図7に示す。特に故障プロセッサ  $P_i$  の前段にある  $P_{i-1}$  は  $R_{i-1}$  に対して再構成のため、 $P_{i+1}$  との接続を要求し、 $R_{i-1}$  が新しい  $P_i$  として処理を補完する。 $R_{i-1}$  から  $P_{i+1}$  への接続経路は2通りあるが、近い経路(右回りの経路)では2個の中継プロセッサ  $R_i, R_{i+1}$  が必要であり、遠い経路(左回りの経路)では  $n-2$  個の中継PEが必要となる。故障の検出から再構成までの時間が短い点から、基本的に中継PE数の少ない右回り経路を使用する。

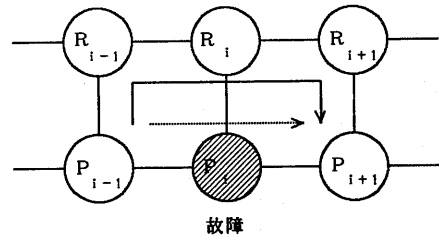


図7  $P_i$  と他のプロセッサの接続関係

#### [例]

図6においてプロセッサ  $P_2$  が故障した場合、 $P_1, P_3, R_2$  がそれを検出して  $P_2$  の代わりに  $R_1$  がステージ2の処理を実行し、処理データを  $P_3$  に転送する。

### 3.2.2 3重リング結合

本方式は、基本的には2重リング結合と同様であるが、基本リングPEの故障によって冗長PEを置き換える場合、再構成接続のための中継プロセッサ数を減少させるために、偶数番と奇数番の基本リングPEにそれぞれ結合する2つのリングを新たに構成するものである。

#### [構成法]

$n$  が偶数のとき、基本リング  $\langle P_i \mid i = 0, \dots, n-1 \rangle$  と冗長プロセッサによる第2リング  $\langle R_j \mid j =$

0, ..., n - 2, jは偶数 >, 第3リング<  $R_k \mid k = 1, \dots, n - 1, k$ は奇数 >を考え、基本リングと第2、第3リングの間に  $\{(P_i, R_i) \mid i = 0, 1, \dots, n - 1\}$  の結合を行なう。n = 8の時の例を図8に示す。

[パイプライン処理]

2重リング結合と同様に基本リング上で処理される。

[再構成アルゴリズム]

基本リングPEの1つ  $P_i$  が故障した場合、冗長プロセッサの  $R_{i-1}$  が処理を引き継ぎ、 $R_{i-1}$  から  $P_{i+1}$  への接続経路が第2リング、第3リングのどちらかのリングによって形成される。2重リング結合と異なり再構成接続のために2つのリングを使用することができるため、右回りの経路では1個の中継プロセッサが必要となり、左回りの経路では  $\frac{n}{2} - 1$  個の中継で経路が構成できる。

[例]

図8においてプロセッサ  $P_2$  が故障した場合、 $P_1, P_3, R_2$  がそれを検出して  $P_2$  の代わりに  $R_1$  がステージ2の処理を実行し、処理データを  $P_3$  に転送する。

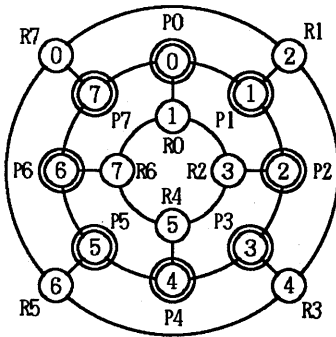


図8 3重リング結合 (n = 8)

### 3.2.3 木状リング結合

2重リング結合と3重リング結合を比較すると、3重リング結合の方が再構成接続のための中継PE数が  $\frac{1}{2}$  に減少している。木状リング結合は、再構成のための中継PEをなくすために直接代替PEを結合するものである。

[構成法]

レベル数  $m(m = 0, 1, \dots, l)$  の2進木の末端の葉ノードの番号  $2^m + i(i = 0, 1, \dots, 2^m - 1)$  を  $\sigma$  置換して、置換後の番号順に基本リングを構成する。n = 8 ( $m = 3$ ) の例を図9, 10に示す。 $R_1$  から  $R_{15}$  までの2進木結合網の葉ノード  $R_8$  から  $R_{15}$  の番号を  $\sigma$  置換して  $P_0$  から  $P_7$  に置き換え、リング結合させている。

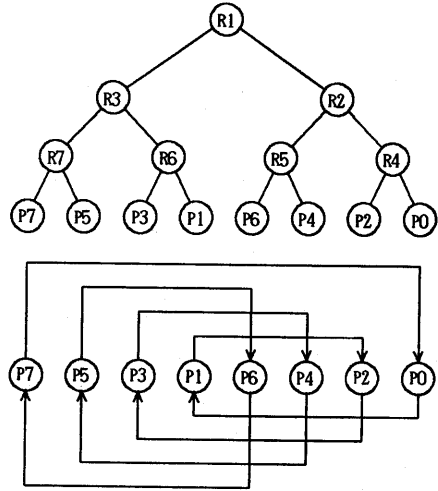


図9 木状リング結合の構成

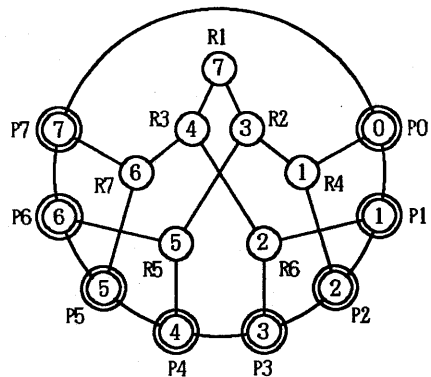


図10 木状リング結合 (n = 8)

[パイプライン処理]

2重リング結合と同様に基本リング上で処理される。

[再構成アルゴリズム]

基本リングPEの1つ  $P_i$  が故障した場合、 $i$  の  $\sigma^{-1}$

置換によって再構成経路が決定される。

【例】

図10においてプロセッサ  $P_2$  が故障した場合、 $P_1, P_3$  がそれを検出して  $P_2$  の代わりに  $R_6$  がステージ2の処理を実行し、処理データを  $P_3$  に転送する。

### 3.3 各ネットワークの比較

基本リングPE数を  $n$  としたときの冗長PE数、再構成時の中継PE数、また、中継PEの信頼度を  $r$  としたときの中継経路の信頼度によって各ネットワークを比較する。

#### 3.3.1 冗長PE数と中継PE数による比較

2重リング結合、3重リング結合は冗長PE数が  $n$  であり、中継PE数も前述のように2個と1個にそれぞれ定まっている。しかし、木状リング結合の場合、再構成のための中継PE数は故障するプロセッサの位置によって変化する。 $n = 2^m (m = 3, 4, \dots, l)$  のときの木状リング結合の冗長PE数と中継PE数は次のように示される。

(1) 冗長PE数  $n - 1$

(2) 中継PE数

$$2(m-1) - 2 \text{ であるものが } 4 \text{ 通り}$$

$$2(m-k) - 2 \text{ であるものが } 2^k \text{ 通り}$$

(ただし  $k = 2, 3, \dots, m-1$ )

また、これらの中継PE数を経路数で重み付けをして平均を求めると

(3) 中継PE数の平均  $2 - \frac{8}{n}$

#### 3.3.2 中継経路の信頼度による比較

中継PEの1つの信頼度を  $r$  としたとき、2重リング結合の中継経路の信頼度  $R_{2c}(r)$  と、3重リング結合の中継経路の信頼度  $R_{3c}(r)$  は中継PEの直列結合と見ることができるから次のように示される。

$$R_{2c}(r) = r^2$$

$$R_{3c}(r) = r$$

木状リング結合では、再構成の経路によって中継PE数が違うことからそれぞれの経路毎に信頼度を求めてその平均を取る。平均の信頼度  $R_{tc}(r)$  は

$$R_{tc}(r) = \frac{1}{n} (4r^{(2m-4)} + \sum_{k=3}^m 2^{k-1} r^{2(m-k)})$$

$n = 8$  の場合の平均の信頼度  $R_{tc}(r)_{(n=8)}$  は

$$R_{tc}(r)_{(n=8)} = \frac{1}{2} r^2 + \frac{1}{2}$$

となり、これは2重リング結合、3重リング結合の信頼度よりも高い。

表1 各ネットワークの比較

比較項目	2重リング結合	3重リング結合	木状リング結合
基本PE数	$n$	$n$	$n$
冗長PE数	$n$	$n$	$n-1$
中継PE数	2	1	$2-8/n$
信頼度( $n=8$ )	$r^2$	$r$	$r^2/2+1/2$

## 4 あとがき

汎用パイプライン処理のためのフォールトトレラント結合方式として2重リング結合、3重リング結合、木状リング結合の3種類を考え、それらの評価を行なった。バンクメモリ結合回路を用いたマルチプロセッサ・システムは、大量のデータをプロセッサ間で転送するパイプライン処理において、従来の結合回路よりも有利であると共に、故障プロセッサを検出して再構成を迅速に行なうことができる。また、各ネットワークは、ホットスタンバイ方式の冗長構成であるが再構成のアルゴリズムと再構成接続経路の信頼度等にそれぞれ特徴がある。その中で3重リング結合は、再構成アルゴリズムが簡単であると共に、再構成時の接続経路が短く有利であると考える。

## 参考文献

- [1] 高橋義造 編：並列処理機構，丸善(1989).
- [2] 遠藤，中村，重井：汎用パイプライン処理における機能割付の最適化，信学技報，EC82-38(1982).
- [3] 南谷 崇：並列処理におけるフォールトトレランス技術，情報処理，Vol.27, No.9, pp.1039-1048(Sep. 1986).