

## 並列ニューラルネットワークシミュレータ Parallel Neural Network Simulators

梶原 信樹, 田地野 真次†, 中田 登志之, 松下 智, 小池 誠彦  
Nobuki KAJIHARA, Shinji TAJINO†, Toshiyuki NAKATA,  
Satoshi MATSUSHITA, Nobuhiko KOIKE

日本電気(株) C&C システム研究所, †日本電気技術情報システム開発(株)  
C&C Systems Research Laboratories NEC Corp.,  
†NEC Scientific Information System Development Ltd.

あらまし 近年、ニューラルネットワークの研究が盛んでハードウェア化の研究も行われている。本文ではニューラルネットワークの持つ並列性を引き出し大規模なニューラルネットワークの高速なシミュレーションを実現するために現在開発中の2つの並列ニューラルネットワークシミュレータについて述べる。1つは汎用シミュレーションマシンとして開発した分散共有メモリ型のマルチプロセッサシステム上でのシミュレータでネットワーク分割によりニューラルネットワークの持つ並列性を引き出す。もう1つはニューラルネットワークシミュレーションに特化し、ネットワーク分割に加え専用プロセッサによるシミュレーションアルゴリズムのバイブライン化によりアルゴリズムの持つ並列性も引き出す。

Abstract This paper describes two parallel neural network simulators, CjNet and NeuMan, which were especially designed for high speed simulation of large-scaled neural networks. CjNet was developed on a general purpose multi processor system with distributed and shared memory. CjNet makes use of parallelism inherent in neural networks through sub-network parallelism. NeuMan on the other hand is a special-purpose processor dedicated to neural network simulation. Besides making use of sub-network parallelism, it makes use of parallelism in the simulation algorithm through pipelining.

### 1. はじめに

ニューラルネットワークを粒度の細かい並列計算のモデルと考え、情報連続、時間連続のニューラルネットワークモデルを提案し、ニューラルネットワーク記述言語、コンパイラ、シミュレータの試作、及びニューラルネットワークによる小規模な情報処理機能(簡単な制御システム、前向き推論システム、時系列パターンの認識、N\_Queen, 等)の実験を行なっている。ニューラルネットワークのプログラムはリンクの重みとして格納し、シミュレーションにより動作の確認を行なった。シミュレータは従来型の計算機(LISPマシン, VAX, PC9801)上にソフトウ

ェアでインプリメントした。LISPマシン上のソフトウェアによるシミュレーションでは90~400リンク/秒の処理速度である(但し、LISPマシン上のシミュレータは作成に当たって処理速度はあまり重視していないのでこの数値はもう1桁程度は向上すると思われる)。VAX8650上に高速性を重視して試作したシミュレータでは約140Kリンク/秒の速度を実現できた。しかし、より大規模な実験や実際のアプリケーションの開発には処理速度が不十分である。これは逐次型の計算機で並列な計算をシミュレートしているためで、ニューラルネットワークの並列性を引き出せるシミュレータの開発が必要である。

並列処理の研究の一貫として電気系のCADを対象とした汎用シミュレーションマシンを開発した。このシミュレーションマシンは分散共有メモリ型のマルチプロセッサシステムで、この上にニューラルネットワークシミュレータ(CjNet)を開発した。

さらに高速なシミュレーションを行なうためにニューラルネットワークのシミュレーションに特化した専用マシン(NeuMan)を開発中である。NeuManは複数の専用PE(プロセッサエレメント)とPE間通信のための多段接続ネットワークで構成される。ネットワーク分割に加えて、専用PEによるシミュレーションアルゴリズムのパイプライン化によりアルゴリズムの持つ並列性も引き出す。

以下ではシミュレータの対象とする時間連続のニューラルネットワークモデル、設計方針、ネットワーク分割による並列シミュレーション方式、汎用シミュレーションマシン上のニューラルネットワークシミュレータCjNet、及びニューラルネットワーク専用マシンNeuManについて述べる。

## 2. ニューラルネットワークモデル

図1にニューラルネットワークシステムのイメージ図を示す。「環境」はニューラルネットワークの動作する環境で、一般には時間とともに状態を変化させるダイナミックなシステムである。ニューラルネットワークは多数のノードと各ノードを接続する重み付のリンクから構成される。

ニューラルネットワークはダイナミックに変化する環境を認識しながら動作する。このため各ノードは活性度という内部状態を持ち、環境または他のノードの影響を受けながら時間とともに変化する。各ノードの活性度のパターンが現在の状況を表しているといえる。ニューラルネットワークは現在の状況(各ノードの活性度)と環境からの情報に従って行動を決定し、次の状態に変化するダイナミックなシステムである。

ノードにはセンサノード、コンテキストノード、モータノードの3種類のノードがある。センサノードは環境の状態に応じてその活性度を変化させ、モータノードはその活性度に応じて環境に操作を加える。コンテキストノードはセンサノードとモータノードの介在を行い、ここでシステムの短期

記憶と情報処理の機能が実現される。

ネットワークの構造には対称、層構造等の特定の構造は仮定しない。自分自身へのフィードバックも含めて任意の構造をとる。理論的な解析や学習には特定の構造を仮定した方が便利であるが、情報処理能力という観点から見ると構造を仮定しない方が有利である。

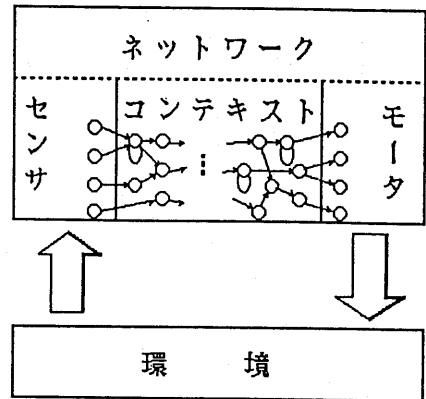


図1. ニューラルネットワークシステム

ダイナミックに変化するパターンを認識するためのニューラルネットワークモデルとして以下に説明する時間連続、情報連続のモデル<sup>[1]</sup>を提案する。図2は自分自身へのフィードバックを許すニューラルネットワークのモデルである。図2のノードxの活性度は他のノード $y_1, y_2, \dots, y_n$ 及び自分自身の活性度の影響を受けて変化する。 $W_{y_1x}$ は、ノード $y_1$ からxへのリンクの重み、 $W_{xx}$ はxからx自身への重みである。各ノードの活性度を $x, y_1, y_2, \dots, y_n$ とするとxは次の微分方程式に従って変化する。

$$\tau \frac{dx}{dt} = \sum W_{y_1x} \cdot \mu(y_1); \text{他のノードからの影響} \\ + W_{xx} \cdot \mu(x); \text{自分自身からの影響} \\ - x \quad ; \text{減衰項} \quad (1)$$

(1)式を動作方程式と呼ぶ。 $W_{y_1x} \cdot \mu(y_1)$ をノード $y_1$ からノードxへの投票と呼ぶ。 $\tau$ は時定数でこの値が小さいほど活性度の変化は大きくなる。出力関数 $\mu$ は(2)式、図3に示す様なリミッター特性を持つ広義の単調増加関数である。

$$\mu(x) = \begin{cases} 0 & ; x \leq 0 \\ 2x^2 & ; 0 < x \leq 0.5 \\ 1 - 2(1-x)^2 & ; 0.5 < x \leq 1 \\ 1 & ; 1 < x \end{cases} \quad (2)$$

動作方程式(1)はオイラー法によって逐次的に解くことができる。時間 $\Delta t$ 毎の積分は次の漸化式で実行できる。

$$\begin{aligned} accx &:= \sum W_{y_i x} \cdot \mu(y_i) + W_{xx} \cdot \mu(x) \\ \Delta x &:= \Delta t [accx - x] / \tau \\ x &:= x + \Delta x \end{aligned} \quad (3)$$

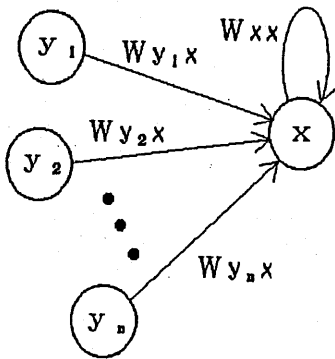


図2. ニューラルネットワークモデル

ニューラルネットワークの動作は各ノード毎に定義される動作方程式(1)の連立方程式で定義される。従ってニューラルネットワーク内の全てのノードに関し(3)式を実行するとニューラルネットワーク全体について $\Delta t$ のシミュレーションができる。これをシミュレーションの1ステップと呼ぶ。シミュレーションの1ステップにはニューラルネットワークの全ノードについて次の3つの処理を行なう必要がある。

【投票フェーズ】

ノードの活性度で出力関数を適用した結果にリンクの重みを乗じた値(投票値)を他のノードに伝える。

【累算フェーズ】

他のノードから送られた投票値を累算する。

【更新フェーズ】

累算結果からノードの活性度を更新する。

この3つの処理は本ニューラルネットワークモデルに限らず他のニューラルネットワークモデルのシミュレーションにも共通する処理である。2つの並列ニューラルネットワークシミュレータは本モデルの高速なシミュレーションを目標とするが、プログラム、マイクロプログラムの変更により他のモデルのシミュレーションも可能である。

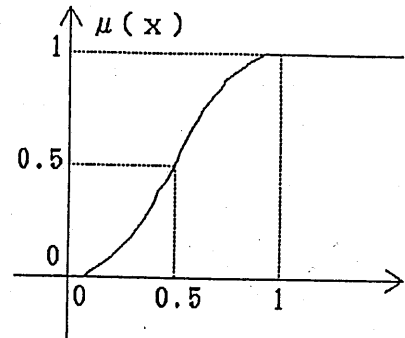


図3. 出力関数 $\mu$

### 3. 基本設計

ニューラルネットワークは行列によって表現することができる。スーパーコンピュータ、アレイプロセッサは行列演算を高速に行なうことができ、これらの高速プロセッサを用いたニューラルネットワークシミュレータが提案されている。ニューラルネットワークを行列で表現する方法ではニューラルネットワークのノードどうしが殆ど完全結合しているような場合には効率の良い処理ができる。しかしシミュレーションに必要な記憶容量、ハードウェア量はノード数の2乗に比例して増加するのでノード数が数千、数万のニューラルネットワークのシミュレーションは現実的でない。

一般にニューラルネットワークのノードはすべてのノードとリンクを持つわけではない。人の脳でもニューロンの数は数100億といわれているが1つのニューロンの持つシナプス結合は数千から数万で脳全体のニューロンの数に比べると非常に少ない。実用的な情報処理システムをニューラルネットワークで実現する場合は幾つかの機能モジュール(部分ネットワーク)で全体が構成される。このとき各機能モジュール内のノード間ではリンクの数は多いが、機能モジュール間ではそれ

ほど多くならない。このようなニューラルネットワークを行列で表現した場合、要素のほとんどは零のスパースな行列となりスーパーコンピュータやアレイプロセッサによるアプローチでは無駄な計算、メモリ、ハードウェアが多くなり効率的ではない。

提案する2つのシミュレータはニューラルネットワークの接続情報は1つの行列で表現するのではなくノード毎にそのノードをソースとするリンクの表 (FanOut表) を持つことで表現する。

ニューラルネットワークの持つ並列性をいかに引き出すかがシミュレータの性能を左右する。ここではニューラルネットワークの持つ並列性、及びシミュレーションアルゴリズムの持つ並列性に着目しシミュレーションの高速化を図った。

ニューラルネットワークの各ノードは並列に動作することができる。従って、ニューラルネットワークのノード1つに1つの物理的なプロセッサを割り当てると並列性は最も高くなる。しかしノードのシミュレーションに必要な処理をデジタル処理で実行するハードウェアをノード毎に割り当てるのは技術的にもコスト的にも現実的ではない。ニューラルネットワークを複数の部分ニューラルネットワークに分割しこれを複数のプロセッサエレメント (PE) に割り当てるアプローチを採用する。1つのPE内ではノードは逐次的に処理されるが各々のPEに割り当てられた部分ネットワークは並列に処理される。これはCADに於ける論理シミュレーションで回路分割として実用化されている手法である。

シミュレーションの1ステップを構成する3つの処理のうち投票フェーズと、累算フェーズは並列に実行できる。ニューラルネットワーク専用マシン NeuMan では、各PEは演算ユニットを2つ持ち2つのフェーズをパイプライン的に並列処理できるように専用化し、シミュレーションアルゴリズムの持つ並列性も引き出す。

#### 4. 並列シミュレーションアルゴリズム

2つのシミュレータ (CjNet, NeuMan) は、どちらもニューラルネットワークを複数の部分ニューラルネットワークに分割し複数のPEに割当て並列に処理することによりニューラルネットワークの持つ並列性を引き出す。シミュレーション用のデータ構造 (ニューラルネットワー

クの表現法)、シミュレーションアルゴリズムも基本的には同じである。

シミュレーション用のデータ構造を図4に示す。ノードの活性度を保持する活性度表 Activity[x]、他のノードからの投票値を累算する累算器表 Accumulator[x]、及びノードから出るリンクの情報を保持した FanOut表 (ノードxが影響を与えるノードyのアドレスとノードyへのリンクの重み Wxyの表) からなる。各PEは担当する部分ネットワークの情報を図4の形式で保持する。

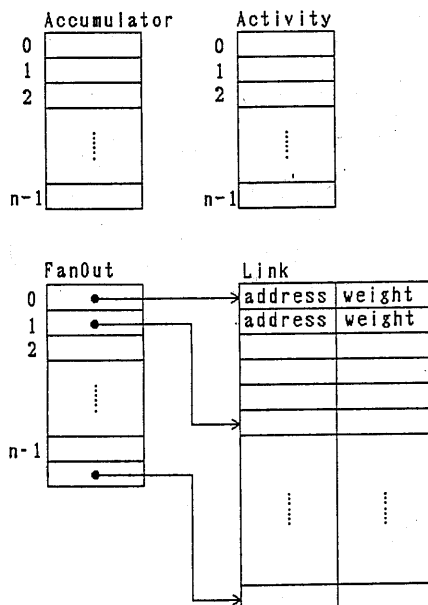


図4. シミュレーション用データ構造

各PEはこのデータ構造を参照して、他のPEと同期しながらシミュレーションを行なう。各フェーズでの処理は次のようになる。

#### 【投票フェーズ】

```
FOR x IN PEiが担当する全てのノード DO
  mx := μ(Activity[x]);
  FOR link IN ノードxの全てのFanOutリンク DO
    MSG.addr := link.address;
    MSG.vote := link.weight × mx;
  put(MSG);
sync();
```

【累算フェーズ】

```
FOR MSG in 送られてきた全てのメッセージ DO
    Accumulator[MSG.addr] :=
        Accumulator[MSG.addr] + MSG.vote;
```

【更新フェーズ】

```
FOR x IN PEiが担当する全てのノード DO
    Activity[x] := Activity[x] +
        Δt×{Accumulator[x] - Activity[x]}/τ;
    Accumulator[x] := 0.0;
```

【記号の説明】

- PEi: プロセッサエレメントi
- Activity[x]: ノードの活性度
- Accumulator[x]: ノードの累算器
- MSG.addr: メッセージのアドレス部
- MSG.vote: メッセージの投票部
- link.address: リンクのノードアドレス
- link.weight: リンクの重み
- sync(): 他のPEとの同期処理
- get(): メッセージの入力
- put(MSG): メッセージの出力
- Δt: 積分の刻み幅

投票フェーズでは活性度表とFanOut表を参照して投票処理を行なう。ノードの活性度 $\mu$ に出力関数 $f$ を適用しノードの出力を計算する。この値にリンクの重みを乗じた値（投票値）と投票先のノードのアドレスを投票メッセージにまとめ他のPEに送る。担当する全てのノードの投票処理が終了すると他のPEと同期のための処理を行なう。

累算フェーズでは、他のPEから届いた投票メッセージのアドレス部で指定されるノードの累算器にメッセージの投票部の値を加算する。

投票・累算フェーズでは、同一PEで処理されるノード間のリンクの処理には他のPEとの通信は必要ない。異なるPEに割り当てられた部分ニューラルネットワークに属するノード間にリンクがある場合はPE間の通信が必要になる。

投票フェーズと累算フェーズが終了すると更新フェーズを実行する。更新フェーズでは活性度表と累算器表を参照して、担当する全てのノードについて活性度の更新処理を行う。

5. 汎用シミュレーションマシン上のシミュレータ

電気系のCAD用に開発した汎用シミュレーションマシンのシステム構成を図5に示す。本シミュレーションマシンは複数のPEをVMEbus準拠のバスで結合した分散共有メモリ型のシステムである。各PEは汎用のマイクロプロセッサシステムで32ビットマイクロプロセッサ(68020)と浮動小数点演算コプロセッサ(68882)、浮動小数点アクセラレータ(WTL1167)を持つ。PEは4Mバイトのローカルメモリを持つ。このメモリは2ポートで他のPEからもアクセスできる。ローカルメモリには、ローカルなアドレスとグローバルなアドレスが割り当てられている。PEは自分自身のローカルメモリはローカルアドレス、グローバルアドレスのどちらでもアクセスできる。他のPEのローカルメモリはグローバルなアドレスによってアクセスできる。

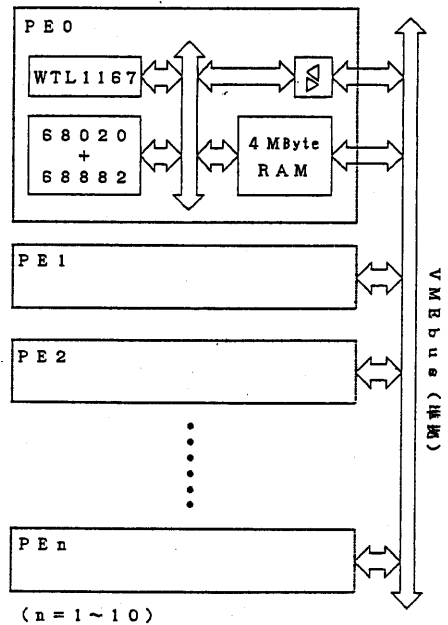


図5. 汎用シミュレーションマシンシステム構成

このシステムは、汎用性が高くこの上でいくつかのアプリケーションを開発している。アプリケーションの1つとしてニューラルネットワークの持つ並列性を引出せるニューラルネットワークシミュレータ(CjNet)を開発した。

CJNetでは、4節の並列アルゴリズムをプログラムによって実現した。

CJNetはPE間通信はバスを介して行なう。異なるPEに割り当てられた部分ニューラルネットワークに属するノード間のリンクの処理は次のようにバスのトラフィックを軽減するようにしている。

ある部分ニューラルネットワーク(N1とする、PE1が担当)の複数のノード $y_1, \dots, y_n$ から別の部分ニューラルネットワーク(N2, PE2が担当)のあるノード $x$ へリンクがある場合、各リンク毎にPE間通信を行なう必要はない。ノード $y_1, \dots, y_n$ からノード $x$ への投票値の累算はPE2ではなくPE1が行い、累算結果を投票値としてPE1に送ることにより、PE間通信を $n$ 回から1回に減らすことができる。

CJNetの性能評価はまだ行っていないが、PE単体の性能は約100Kリンク/秒で、VAX8650の0.7倍、次節で説明するニューラルネットワーク専用マシンNeuManのPEの0.1倍の処理速度である。

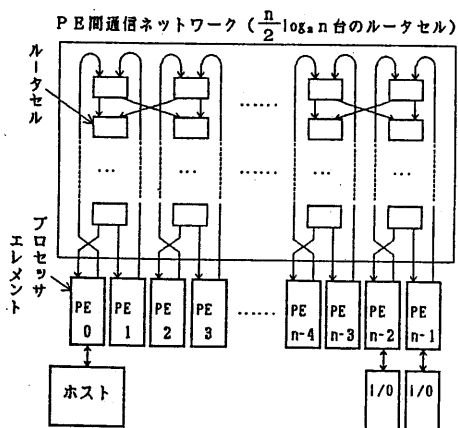


図6. NeuManシステム構成

## 6. ニューラルネットワーク専用マシン

ニューラルネットワーク専用マシンNeuManのシステム構成、PEのブロック図を図6、図7に示す。NeuManは複数の専用PEとPE間でメッセージを転送するためのPE間通信ネットワークで構成される。

NeuManはホストのバックエンドプロセッ

サとして動作し、シミュレーション用のニューラルネットワークのデータはPE間ネットワークを介してホストから各PEにロードされる。シミュレーションの制御、シミュレーション時の各ノードの活性度のモニタ、環境とのインタフェースもホストが行なう。また各PEはI/Oポートを持っておりホストを介さずに環境とインタフェースを行なうことも可能である。

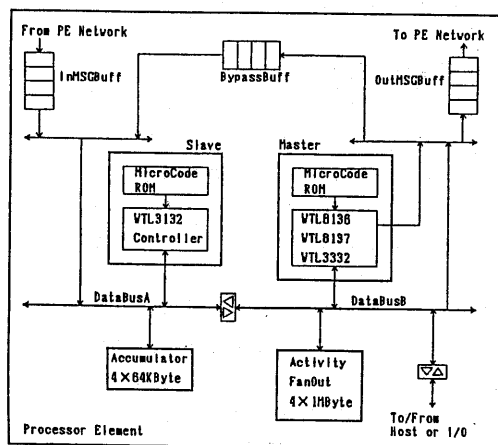


図7. プロセッサエレメント (PE)

PEは投票フェーズと更新フェーズをパイプライン的に並列に実行するためにマスタとスレーブの2つの処理ユニットを持つ。マスタ、スレーブはともにマイクロプログラムで制御されマイクロコードの変更によりいろいろなニューラルネットワークモデルや学習にも対応できる。マスタは32ビット浮動小数点プロセッサ、整数プロセッサ、シーケンサで構成されPE内のすべての資源にアクセス可能で投票フェーズと更新フェーズを実行する。スレーブは32ビット浮動小数点プロセッサと単純なコントローラで構成され、マスタが投票フェーズを実行中にこれと並列に累算フェーズを行なう。ニューラルネットワークのシミュレーションの大きな部分を占める投票フェーズと累算フェーズをマスタ、スレーブでパイプライン的に並列に処理することにより高速なシミュレーションを実現する。投票、累算フェーズ終了後の更新フェーズはマスタが実行する。

投票フェーズで参照する活性度表、FanOut表と累算フェーズで参照する累算器表は各々の

演算ユニットから独立にアクセスすることができる。更新フェーズでは累算器表と活性度表を参照する必要がある。マスクはスレーブを停止した状態で累算器表、活性度表をアクセスする。

PEは他のPEと投票メッセージの送受信用にPE間ネットワークの入出力ポートを持つ。また、2つのユニット間には投票メッセージ用のパイプがあり、同一PE内のノード間の投票メッセージはこのパイプを使いPE間ネットワークのトラフィックを軽減する。

活性度表・FanOut表の容量は4Mバイト、累算器表の容量は256Kバイトで、1つのPEでは最大64Kノード、440Kリンクのニューラルネットワークのシミュレーションが可能で、PEあたり約1Mリンク/秒の処理速度を実現する。

シミュレーション時のデータ転送、ホストによるシミュレーションの制御、PEの同期、PEへのニューラルネットワークデータのロードはPE間通信ネットワークを介してメッセージベースで行なう。

メッセージは可変語長（1語16ビット）でヘッダ（1語）と本体（0~127語）から成り、ヘッダにはメッセージのタイプ、送り先のPEのアドレス、メッセージ長の情報を含む。

ニューラルネットワークシミュレーション時の各ノード間の投票値の転送は投票メッセージで行なう。投票メッセージは3語構成で、投票の送り先のノードを識別するための情報と、投票値を含む。

PE間通信ネットワークは複数のルータセル（2×2の交換器）からなる多段接続ネットワークである。n台のPEを接続するためには $\lceil n/2 \rceil \times \log_2 n$ 台のルータセルが必要である。図8にルータセルのブロック図を示す。ルータセルは、2つの入力ポートA、Bと2つの出力ポートX、Yを持つ2×2の交換器で内部に4本のメッセージバッファを持つ。各ポートの幅は16ビットである。メッセージはそのアドレスに従って所定のバッファに入れられる。どちらのバッファに入れるかは前段のルータセルの出力ポートが制御する。

ルータセルはマルチプレクサへの入力となる2つのバッファが両方とも空の時以外は稼働し続けることができる。シミュレーション時に各PEが1Mリンク/秒の速度で動作するためにはルータ

セルは2つの入力ポートの各々から平均1Mメッセージ/秒の投票メッセージを受取り、2つの出力ポートのそれぞれへ同様のメッセージ密度で送り出さなければならない。本ルータセルは1つの入力ポートの平均メッセージ密度が1.5Mメッセージ/秒まではバッファが溢れることなく（溢れる確率は $1.5 \times 10^{-8}$ ）交換できる。

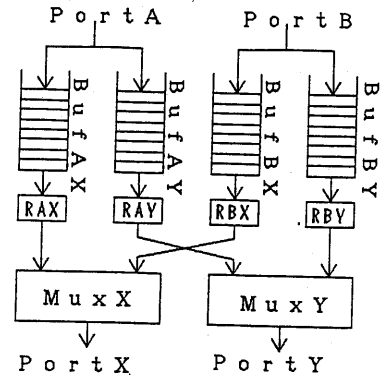


図8. ルータセル

## 7. おわりに

分散共有メモリ型のマルチマイクロプロセッサ上のニューラルネットワークシミュレータCjNet、及びニューラルネットワークシミュレーションに専用化したニューラルネットワーク専用マシンNeuManについて述べた。両シミュレータともにネットワーク分割によりニューラルネットワークの持つ並列性を引き出す。さらにNeuManではPEとPE間通信ネットワークを専用化することによりシミュレーションアルゴリズムの持つ並列性も引き出す。

今後は、CjNet上に大規模なニューラルネットワーク応用システム（画像処理システム）を構築する。NeuManは、PE4~8台構成のシステムを試作しアーキテクチャの妥当性を確認する予定である。また、シミュレータと密結合したニューラルネットワークプログラミング環境を開発する。ニューラルネットワークプログラミング環境にはニューラルネットワーク記述言語、コンパイラ、ネットワーク分割ローダ、モニタが必要で、特に記述性の高いニューラルネットワーク記述言語とネットワーク分割ローダは重要である。

ネットワーク分割によるニューラルネットワー

クシミュレータが効率良くシミュレーションを実行するためには部分ニューラルネットワーク間のリンクがなるべく少なくなる様にニューラルネットワークを部分ニューラルネットワークに分割しなければならない。ニューラルネットワークプログラムはニューラルネットワーク記述言語を使って機能毎にモジュール化して記述される。ニューラルネットワークの分割にはこのモジュール化の情報を利用することができる。

#### 8. 参考文献

[1] 梶原, 辻; "ニューラルネットワーク推論

システム"; 情報処理学会知識工学と人工知能研究会資料45-10(1986)

[2] 梶原, 中田, 小池; "専用並列マシンMAN-YOにおけるニューラルネットワークシミュレータ"; 情報処理学会第35回全国大会予稿集(1987)

[3] 梶原, 中田, 松下, 小池; "ニューラルネットワークシミュレーションマシン: NeuMan"; 情報処理学会「コンピュータアーキテクチャ」シンポジウム(1987-5)