

ニューラルネットワークアクセラレータ NEURO-TURBO

柳橋喜久治* 松田 聡* 吉田征夫* 岩田 彰** 鈴木宣夫**
*(株) マイテック **名古屋工業大学電気情報工学科

24ビット浮動小数点方式の汎用DSPとして開発されたMB86220を4個用いたニューラルネットワークアクセラレータNEURO-TURBOについて述べる。NEURO-TURBOは、DSP4個をデュアルポートメモリ(DPM)を介してリング状に結合したMIMD型並列処理プロセッサである。バックプロパゲーションの学習時の演算速度で2MCPS、また、前向きの演算速度で11MCPSを得た。また、NEURO-TURBOに、我々が先に提案した大規模4層ニューラルネットワーク(CombNET)をインプリメントした。その結果、印刷漢字認識の処理時間が、学習時に数時間、認識時には100文字/秒と高速化され、実用的な処理速度を得ることができた。

A NEURAL NETWORK ACCELERATOR NEURO-TURBO

Kikuji YANAGIHASHI*, Satoshi MATSUDA*, Masao YOSHIDA*
Akira IWATA** and Nobuo SUZUMURA**

* Mitec Corp., 3-32-1, Takada, Toshima-ku, Tokyo, 171, JAPAN

**Dept. of Electrical and Computer Eng., Nagoya Institute of Technology
Gokiso-cho, Showa-ku, Nagoya, 466, JAPAN

A Neural Network Accelerator named Neuro Turbo was developed using four general purpose 24 bits floating point Digital Signal Processor(DSP) MB86220. Neuro Turbo is a MIMD type parallel processor having ring coupled 4 DSPs and 4 Dual Port Memories(DPM). The performance was evaluated by constructing a neural network to recognize the 26 type fonts of the alphabet set. The processing speeds of 2 MCPS for learning procedure and 11 MCPS for forward pass were achieved. This paper also describes the implementation of CombNET on NEURO-TURBO. A network to recognize 2965 printed Chinese characters has been constructed. The recognizing rate of 99.5% has been achieved even for test data sets.

1 まえがき

近年、ニューラルネットを音声認識、画像処理、ロボット制御などさまざまな分野に適用した研究が多く報告されているが、これらの分野に実際にニューラルネットを適用する場合、実用的処理速度を得るにはかなり高速な処理能力を有するプロセッサが必要とされる。ニューラルネットにおける計算処理は元来並列処理であるため、計算速度を上げるためには並列処理アーキテクチャを採用することが有効である。究極的には、ニューロンごとに専用プロセッサを用意し、それらを適当なネットワークで連結してニューラルネットを実現する方式（ニューロコンピュータ）が見込まれるが、現状技術ではその方式はコストパフォーマンス面から考えると得策ではない。

ここで、積和演算を高速に行うことのできる信号処理プロセッサ（DSP）を複数個用いたニューラルネットアクセラレータがいくつか発表されている。それらのシステムでは、ニューラルネットを複数個のプロセッサに分割して配置し、プロセッサ間を結ぶバスによって互いのプロセッサが通信しながら並列的にニューラルネットの計算し実行する。近年、大きなアドレス空間を持ちまた浮動小数点演算のできる第3世代DSPが開発されるに至り、大規模なニューラルネットをDSPを用いたシステムで構築することができるようになった。汎用DSPを用いたニューラルネットアクセラレータとしては、TIのOdyssey [1]、IBMのNEP [2]などが発表されている。前者は共通バスに複数のDSPを配置した構成で2MCPSの速度を得ている。後者は複数のDSPをリングに結合したシステムで490KCPSの速度を得ている。また、Warp [3]では専用の浮動小数点プロセッサ10個をリニアアレイ状に配置したシステムで17MCPSを実現している。

我々は、汎用DSPとして最近開発されたMB86220を4個用いたニューラルネットアクセラレータ、NEURO-TURBOを構築した[4][5]。表1にMB86220の諸元を示す。MB86220は24ビット浮動小数点方式DSP（ただし、内部演算は30ビット浮動小数点方式）であるが、ニューラルネットの計算ではニューロン毎に特性関数によって0から1までの間に正規化されることと、実際応用の際、入力データとしては高々16ビット精度であることを考慮すると、24ビットの浮動小数点演算による演算精度で十分であると考えられる。MB86220は富士通のCMOS1.2 μ mルールのVLSIとして構成され安価に供給される。NEURO-TURBOは、このDSP4個を独自のリング結合したMIMD型並列処理プロセッサである。

2 DSPのリング結合アーキテクチャ

ニューラルネットは積和演算をベースとする単純な計算の集合で、並列的に構成される。このニューラルネットの並列性を生かした高速な演算処理装置が要求される。この目的のために従来から複数のコンピュータによる並列分散処理が考えられてきたがコンピュータ間のシステムバスのバッファ及びスイッチング回路等のハードウェアの増大とバスアービトレーション等の制御の複雑さが問題となっている。

ここでは最もハードウェアが単純化される並列分散処理トロボジを選択し、その上でニューラルネットの実行速度を落とさない最適な計算アルゴリズムを組み合わせて実行速度の向上を図った。図1にこの構成の概念的なブロックダイヤグラムを示す。図に示されるように4個のDSPはお互いにデュアルポートメモリ（DPM）を介してリング状に結合され、その間にバッファやスイッチング回路をなくし、ハードウェアの単純化を図っている。1つのDSPの外部データメモリとして、隣接する2つのDPMと1つのワーキングメモリ（WM）があり、それらはDSPのシステムバスで結合されている。

1つのDPMは隣接する2つのDSPからランダムにデータをアクセスすることができ、DSP間のデータ転送用として用いられる。WMは、それが接続されているDSPのワーキングメモリとして結合の重みやデータ及び中間結果をストアするのに用いられる。ニューラルネットは連続した単純な積和演算を繰り返し実行することが多く、計算アルゴリズムの工夫により、演算を分割しデータをパイプライン状にリングに沿って淀みなく転送し、並列処理することによって高速実行が可能となる。この時、一方のDPMのデータをDSPに入力し処理後、他方のDPMに出力することによって、データのメモリアクセスが実質的にデータ転送となり、データ転送時間を最小にして高速並行処理を実現している。

このリング結合アーキテクチャを用いた並列処理プロセッサNEURO-TURBOの構成を図2に示す。図3はNEURO-

表1. MB86220の諸元

DATA FORMAT	24 bits floating point 18 bits mantissa, 6 bits exponent
MULTIPLIER	18E6 x 18E6 → 24E6
ALU	18E6 + 18E6 → 24E6 24E6 + 24E6 → 24E6
MACHINE CYCLE	150nsec per floating point arithmetic instruction 75nsec per other instruction
PROGRAM MEMORY	2KW (internal) 64KW (external)
DATA MEMORY	256W x 2 banks (internal) 64KW x 4 banks (external)

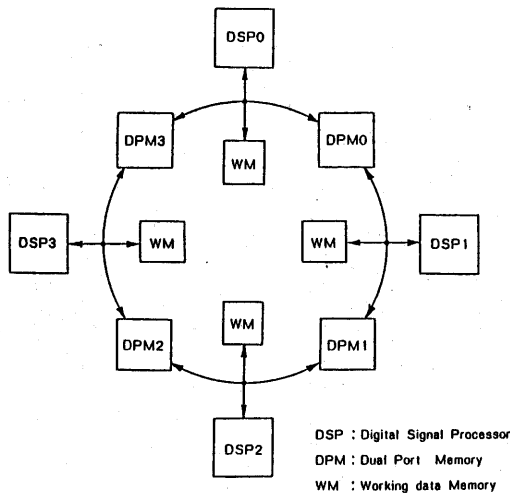


図1. デュアルポートメモリ (DPM) を介したリング結合アーキテクチャ

TURBO 外観である。ホストコンピュータとして、NEC のパーソナルコンピュータ PC-9801 を用い、その拡張スロットに本装置をプラグインする。ホストコンピュータと本装置のインターフェースとして別の DPM を用いて、データ転送を行う。ホストコンピュータからは、各 DSP のプログラムメモリ (PM) に、実行前に予めプログラムをロードすることができ、問題毎にプログラムを入れ換えて使用することが出来る。WM はアドレスを 2 ビット拡張し、メモリ空間を 64 kw ~ 256 kw としている。さらに、このシステムの DSP は内部システムクロック 12 MHz で動作し、外部データメモリとのアクセスは 1 ウェイトを必要とする。各 DSP は積和演算を 2 クロックサイクルで実行する。4 つの DSP を用いたこのシステムの最大スループットはほぼ 24 MOPS となる。

3 性能評価

NEURO-TURBO の性能評価は、バックプロパゲーションアルゴリズムを用いて、アルファベット 26 文字のタイプフォントの認識のための 3 層 (1 隠れ層) ニューラルネットワークの作成によって行った。入力パターンは、イメージスキャナによって読み込まれたタイプフォントイメージを 7x5 の小領域に分割し、小領域ごとの濃度値を 0 から 1 の値に正規化することによって得た。ネットワークの構成は、入力層ユニット 35 個、中間層ユニット 25 個、出力層ユニット 26 個、バイアスユニット 2 個であり、ユニット間の結合数は 1576 結合である。NEURO-TURBO を用いて、このネットワーク

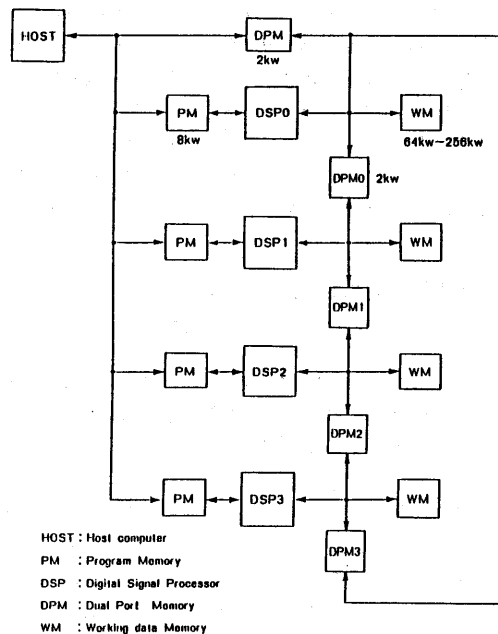


図2. NEURO-TURBO の構成図

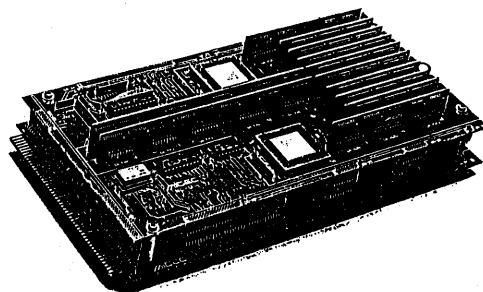


図3. NEURO-TURBO の外観

の一回の学習に要した時間は 780 μ sec であった。この速度は、ほぼ 2 MCPS に相当する。また、前向き演算速度は、このネットワークで 140 μ sec であり、これは、ほぼ 11 MCPS に相当する。ちなみに、同じネットワークのシミュレーションを SUN4/260 ワークステーションを用いて行ったところ、学習時 100 kCPS、前向き演算で 370 KCPS であった。すなわち、NEURO-TURBO は SUN4/260 の約 2 倍の高速演算能力を持っている。この速度は専用の浮動小数点プロセッサ 10 個をリニアアレイ状に配置した Warp には及ばないものの、本システムと同様に汎用 DSP を用いたニューラルネットアクセラレータ Odyssey や NEP より優っている。NEURO-TURBO で

用いたDSP、MB86220は6MOPSの積和演算速度を持っており、これを4個使用した本アクセラレータは最大24MOPSの積和演算速度を持っている。この点を考慮すると、3層ニューラルネットワークのバックプロパゲーションアルゴリズムの実行において上記の高速演算性能を得たことは、デュアル・ポート・メモリを介してリング状に結合された4個のDSPが効率よく動作していることを示している。

4 大規模4層ニューラルネット (CombNET)

これまでに行なわれたニューラルネットワークに関する研究は、分類カテゴリ数の少ない比較的小規模なニューラルネットを取り扱っており、実用的で大規模なニューラルネットにこれらの成果をそのまま拡張できるかは疑問である。例えばJIS第一水準の漢字を認識するような大規模なニューラルネットを単純なバックプロパゲーションアルゴリズムによって構築することは実際には多くの困難を伴う。そのような大規模なネットワークの学習では、ローカルミニマムに陥る可能性も高く、たとえ収束するにせよ膨大な計算量を費すことになる。

大規模なニューラルネットを実用的な計算時間で構築するには、ひとつにはニューラルネットのシミュレーションを高速に実行できるプロセッサを用いることであるが、ネットワークを学習を行い易い小規模なネットワークに分割し、それらの統合として大規模なネットワークを構築することも必要となる。そこで、我々は、従来のネットワークモデルでは困難であった多数のカテゴリを分類するための大規模ニューラルネットワーク (CombNET) を提案している [6][7]。

CombNETは、図4に示すように、第1層にベクトル量子化型ニューラルネットを配置し、後段の第2、3、4層には小グループ内のデータを分類する小規模な階層型ニューラルネットを多数並列に配置したネットワークである。第1層を櫛の幹、第2、3、4層を櫛の葉とみなすと、ちょうど櫛のような構造をもっているため、櫛型ネットワークという意味からCombNETと名付けた。そして、第1層をStem Network、第2、3、4層をBranch Networkと呼んでいる。

CombNETの学習は次の手順で行う。まず、Stem Networkの学習をKohonenの自己組織化アルゴリズム [8] により行う。Stem Networkの学習後、入力データの各カテゴリがStem Networkのどのニューロンと最適整合になるかを調べ、各ニューロンの分担すべきカテゴリを求め、全カテゴリをStem Networkのニューロン数と同数のグループに分割する。次に、その分割されたグループごとにその中にあるカテゴリを識別

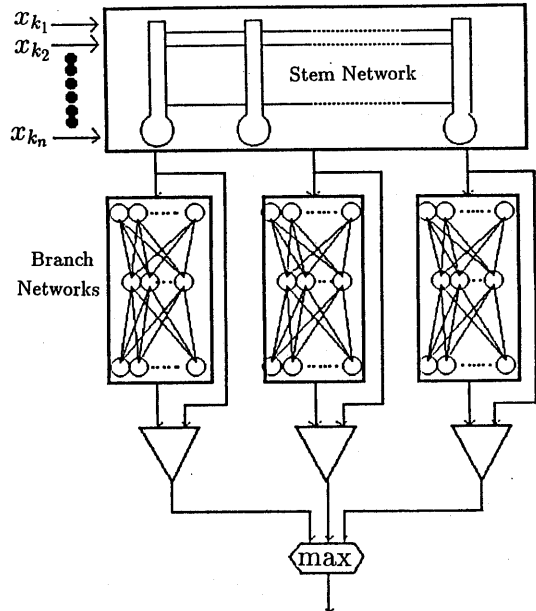


図4. CombNETの構成

するためのBranch Networkの学習をバックプロパゲーション法 [9] を用いて行う。

識別は次のようにして行う。まず、Stem Networkに入力データを通し、適合度の高いものから3ないし5番目までのニューロンを選ぶ。次に、選ばれたニューロンが担当するカテゴリグループを分割する後段のBranch Networkに入力データを入力し、出力層の中でもっとも高い発火レベルになったものの出力値を調べる。そして、

$$(\text{適合度})^\alpha \cdot (\text{出力値})^\beta$$

の値がもっとも高くなったものを識別結果としてを選ぶ。こうして、後段の複数の階層型ネットワークの出力と競合させることで識別能力を高めることを目指した。

このネットワークの利点は学習の容易さにある。バックプロパゲーション法による学習は小規模なネットワークについては収束も容易であるが、大規模なネットワークになるとローカルミニマムに陥る可能性も高く、たとえ収束するにせよ膨大な計算量を費すことになる。本方式は、Stem Networkによって大分類をしてからバックプロパゲーションによる学習を行なうため、Branch Networkの規模は小さくなり、したがって、学習は容易になる。また、Branch Networkは相互には結合のない独立したネットワークであるから、大規模なネットワークであっても全体の結合数を少なく抑えることができる。

5 CombNETによるJIS第1水準印刷漢字の識別

CombNETのパターン識別能力を検討するために、JIS第1水準の印刷漢字2965字種を識別するニューラルネットの構築を試みた。

学習データは次のように作成した。まず、JIS第1水準漢字を印刷したものを解像度が400dot/inchのイメージスキャナで読み込む。ここでは、原稿はA4縦置き、横書きとした。そして、イメージスキャナで読み込んだ画像から、

- (1) 横(行)方向に画素値を積分した周辺分布から、しきい値処理により行を切り出し、
- (2) 縦方向に画素値を積分した周辺分布から文字を切り出した。

1文字の大きさは4mm程度であり、これをイメージスキャナで読み取り切り出すと縦横がおよそ64ドット×64ドット程度の2値画像が得られる(図5(a))。これを16×16の小領域にまとめて、濃度正規化したパターンを作成した(図5(b))。今回、これを5セット作成し、そのうちの4セットの平均パターンを学習データにした。また、4セットのうちの1セット(学習パターン)と、残り1セット(未学習パターン)を用いてネットワークのパターン識別能力の評価を行なった。

Stem Networkのニューロン数は144個(12×12)とした。参照ベクトルの初期値は一樣乱数で与え、学習回数は1字種当たり100回の学習を行なった。その結果得られた各ニューロンの参照ベクトルを図6に示す。この結果は非常に興味深いものが得られた。図6には、参照ベクトルの値を16×16の2次元マトリックス状のパターンで示したので、各参照ベクトルがどのようなテンプレートを形成しているかが明らかとなっている。各参照ベクトルには、漢字の「へん」、「つくり」、「かまえ」が形成されている。自己組織化によって図6に示すような参照ベクトルが自律的に形成されたことは注目値する。

次に、後段のBranch Networkについて、各ニューロンに割り当てられた漢字のグループごとに、バックプロパゲーション法を用いて学習を行なった。各グループの字種数は最大でも91個であり、比較的小規模なニューラルネットとなるため、すべてのグループについて容易に学習が収束した。

このように学習を行なったネットワークの識別能力を検討した。さまざまな α 、 β の値を用いて識別を行なったところ、 α/β が5程度のところで未学習パターンについても99.5%の正当率を得られることがわかった。適合度や出力値はともに0~1の間の値であるから、 $\alpha/\beta=5$ というのはStem

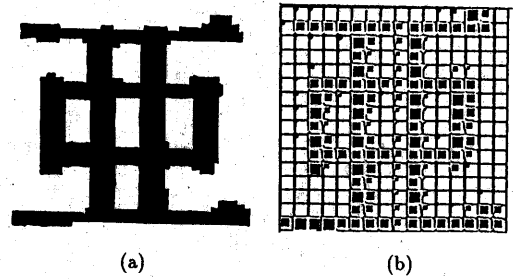


図5. 文字パターンの作成

(a) 切り出された文字イメージ

(b) 濃度正規化された文字パターン

Networkの適合度の結果よりも、Branch Networkの出力値をより強く考慮して判断していることを表している。これはBP学習による階層型ニューラルネットの2つの性質、1つは学習パターンについては特定の素子が強い出力を出し、未知パターンについてはどの素子も弱い出力を出す性質、もう1つは汎化能力に優れパターン変動に強いという性質をうまく利用したものになっている。

実際にこのことを基本の「基」という文字を例に上げて説明する。図7は、この文字をStem Networkに通した結果を示している。白抜き四角内黒い部分は入力パターンに対する各ニューロンの適合度を示しており、その下は各ニューロンの参照ベクトルを表わしている。そして、四角枠は、適合度の高いニューロンを3つ示している。

次に、この3つのニューロンに接続されたBranch Networkにこの文字を入力する。その様子を図8に示す。一番左の部分は、Stem Networkで最も整合のとれたニューロンとして選んだものにつながっているBranch Networkである。しかし、この中にはこの基本の「基」という文字が学習されていないので、Branch Networkの出力値はどれも小さくなっている。そのため、Stem NetworkとBranch Networkを結合した出力値(CombNETの棒グラフ)もそれほど大きくなっていない。一方、真中のNetworkにおいては、基本の「基」という文字が学習されているため、Branch Networkの出力値の中に非常に値の大きなものがある。このため、結合した出力値は他のものと比べて大きくなり正しく認識されたことになる。

これは第1層における誤りを階層型ネットワークで回復していることになり、全体として汎化能力に優れたネットワークとなっている。このように実用的な規模で実用的な識別能力を持ったネットワークを構築することができた。

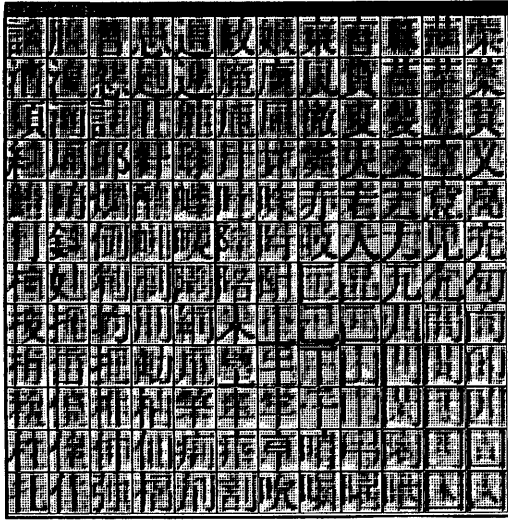


図6. 学習後の Stem Network の参照ベクトル



図7. 「基」を Stem Network に通した結果

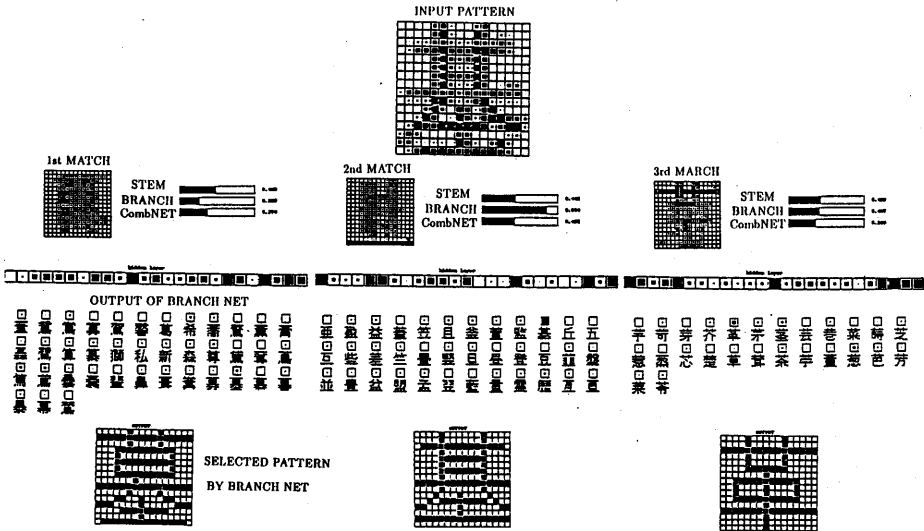


図8. 「基」を Stem Network で選んだ Branch Network に通した結果

6 NEURO TURBOへのCombNETのインプリメント

次にNEURO-TURBOにCombNETをインプリメントした。まず、Stem NetworkをNEURO-TURBO上にインプリメントする。ここでは、12×12の二次元格子状に配置した多数のニューロンを4つのグループに分割し、各DSPは担当するグループのニューロンのみを計算することで並列化を図っている。

次に、Branch Networkの前向き計算は、シストリックアレイの考え方と同様にしてリング結合された4つのDSPに分割して計算を行わせることができる。そこで、ニューロン間の結合の重みを分割して各DSPに割り当てる

今回の実験に用いたシステムは、図9に示すように、ホストコンピュータとしてPC-9801を用い、印刷漢字の画像データを読み込むためのイメージスキャナと、イメージスキャナから取り込まれた画像データなどを保存しておくためのRAMボードと、ネットワークの計算を行うためのNEURO-TURBOから構成されている。

まず、学習時において、NEURO-TURBO上にStem Networkを構築し、学習データを用いて学習を行う。次に、学習のされた参照ベクトルを用いてカテゴリの分割を行った後、各グループ毎にBranch NetworkをNEURO-TURBO上に構築し学習を行う。そしてこの学習によって得られた参照ベクトルの重み及びニューロン間の結合の重みを保存しておく。

また、認識時においては、ネットワークの学習時に生成されたStem Networkの参照ベクトル、及び、Branch Networkの結合重みをNEURO-TURBOのWMにまず送り込み、NEURO-TURBO上にStem Networkと144個のBranch Networkを構築する。そして、新たにイメージスキャナより取り込まれた画像から入力データを作成し、NEURO-TURBOに送り込みネットワークに通して一気に認識を行う。

本システムによる識別率は先に述べたワークステーション上でのシミュレーションと同様の結果が得られた。また、処理速度については、前段のStem Networkの学習時間は、2965文字で100回の学習を行うのに47分、後段のBranch Networkの学習時間は、1000回ずつの学習を行なうのにおおよそ3時間を要した。認識速度については、Stem Networkにおいて一文字当たり8.1ミリ秒であり、Branch Networkにおいては0.4ミリ秒であった。

この結果から、Stem Networkにより適合度の高いものを5番目まで選び、それをBranch Network通す場合を考えると、一文字当たり約10ミリ秒の時間を要することになる。すなわち、100文字毎秒の文字識別速度が得られたことになる。

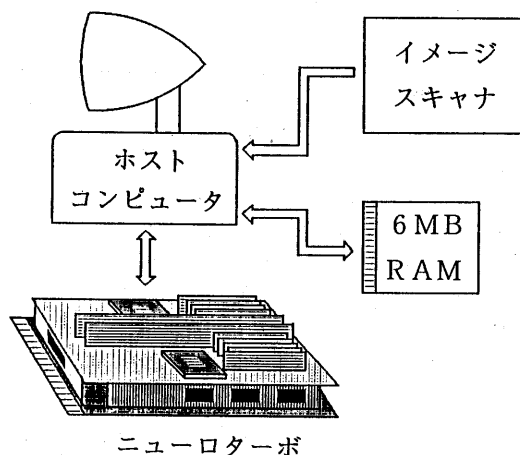


図9. 文字認識システム

表2. 処理時間

(単位:分)

	Stem Network	Branch Network
学習時間 (NEURO TURBO)	47	約 180
学習時間 (SUN4-260)	1,110	約 1,800
認識時間 (NEURO TURBO)	8.1 ms (123 文字/秒)	0.4 ms (2300 文字/秒)

表3. ネットワーク規模

	Stem Network	Branch Network
ニューロン数	144	43,429
結合数	36,864	228,146

このCombNETは表3のようにニューロン数は合計で約43,500個、ニューロン間の結合数は約26万コネクションである。このような大規模なネットワークにおいてもこのような高速処理能力を持つ専用のハードウェアを用いることにより実用的な速度を得られることが確認できた。

7 まとめ

本論文では、汎用DSPとして最近開発されたMB86220を4個用いたニューラルネットワークアクセラレータ、NEURO-TURBOについて述べた。NEURO-TURBOは、DSP4個をデュアルポートメモリを介してリング状に結合したMIMD型並列処理プロセッサである。NEURO-TURBOはバックプロパゲーションの学習時の演算速度で2MCPS、また、前向き演算速度で11MCPSを得た。この速度は、SUN4/260ワークステーションの約20倍であった。Neuro-Turboのリング結合アーキテクチャは、プロセッサの数を4以上に増加することができる。その場合、ニューラルネットワークの演算速度は、プロセッサにはほぼ比例して速くなる。

つぎに、NEURO-TURBOに、我々が先に提案した大規模4層ニューラルネットワーク(CombNET)をインプリメントした。その結果、これまでエンジニアワークステーション上で学習時に2~3日、認識時に数文字/秒かかっていた印刷漢字認識の処理時間が、学習時に数時間、認識時には、100文字/秒と高速化され、実用的な処理速度を得ることができた。

References

- [1] Penz, P.A. and Wiggins, R. : Digital Signal Processor Accelerators for Neural Network Simulations, Neural Networks for Computing, AIP Conf. Proc. 151, pp.345-355.
- [2] Cruz, C.A., Hanson, W.A. and Tam, J.Y.: Neural Network Emulation Hardware Design Considerations, ICNN87, pp.III427-434, 1997.
- [3] Pomerleau, D.A., Gusciora, G.L., Touretzky, D.S. and Kung, H.T. : Neural Network Simulation at Warp Speed : How We Got 17 Million ConnectionPer Second, ICNN88 pp. II-143-150, 1988.
- [4] 佐藤、岩田、鈴木、松田、吉田 : 「汎用浮動小数点DSPによるニューラルネットワークアクセラレータ」、電子情報通信学会技術研究報告、MEとバイオサイバネティクス MBE88-134, pp.83-88, 1989
- [5] Iwata,A., Yoshida,Y., Matsuda,S., Sato,Y. and Suzumura,N. : An Artificial Neural Network Accelerator Using General Purpose 24 Bits Floating Point Digital Signal Processors, Proc. of Int. Joint Conf. on Neural Networks, Washington D.C., Vol.2, pp171-175, 1989
- [6] 岩田、當麻、松尾、鈴木 : 「大規模4層ニューラルネットワークの構築手法」、電子情報通信学会技術研究報告、ニューロコンピューティング,NC89-7, pp.37-42, 1989
- [7] 當麻、岩田、堀田、松尾、鈴木 : 「大規模4層ニューラルネットワーク(CombNET)による印刷漢字認識」、電子情報通信学会技術研究報告、ニューロコンピューティング,NC89-39, pp.39-44, 1989
- [8] Kohonen,T. : Self-Organization and Associative Memory, Springer-Verlag, 1984 and 1988.
- [9] Rummelhart, D.E., Hinton, G.E. and Williams, R.J. : Learning Representations by Back-propagating Errors, Nature, 323-9, pp533-536, 1986.