

FPGA を用いた手書き漢字認識ニューラルネットのハードウェア支援

竹谷 史生* 吉永 努* 荒井 正之†
大津 金光* 馬場 敬信*

*宇都宮大学工学部、†帝京大学理工学部

本稿では、書き換え可能な FPGA を用い、様々なニューラルネットアルゴリズムに対応可能なシステムを提案する。また、一例として、本システムにバックプロパゲーションネットワーク構造の手書き漢字認識アルゴリズムを適用したアーキテクチャについて述べる。4 ニューロンを集積する PU の FPGA を XILINX 社の XC40125XV-1 で実現する場合の予備評価を行った。その結果、演算サイクル数はワークステーション上の C プログラムと比べて約 1/200 となり、学習演算処理能力は 129.2MCUPS、フォワード演算処理能力は 600.6MCUPS という結果が得られた。

FPGA Hardware Support Tool for Recognition of Handwritten Kanji Using Neural Network

Fumio TAKEYA*, Tsutomu YOSHINAGA*, Masayuki ARAI†,
Kanemitsu OOTSU* and Takanobu BABA*

*Utsunomiya University, †Teikyo University

This paper presents a system which is able to implement various neural network algorithms. It utilizes the reconfigurability of FPGAs. We apply the system to an algorithm for hand-written kanji recognition using a back propagation network. Preliminary evaluation is made by supposing to use XILINX XC40125XV-1 for a processing unit. As a result, we show its performance of 129.2MCUPS for learning calculations and 600.6MCUPS for forward calculations.

1 はじめに

ニューラルネットワークの学習計算は、膨大な計算量を必要とする。そのため、汎用のワークステーション等では、学習計算に非常に長い時間を必要とする。これに対して、ASIC を用いたニューロチップ [1][2][3] やニューロシステム [4] などが開発されてきている。しかし、それらのニューロチップを使って高速に学習計算をしようとする、そのチップ(ハードウェア)構成に合わせたアルゴリズムを採用しなくてはならない。さらに、新たなアルゴリズムに対しては、ハードウェアを再開発しなければならないことも多い。

そのような中、近年プログラム可能なゲートアレ

イ FPGA の技術的な発展はめざましく、高集積、高速な FPGA が開発されている。1990 年頃から現在までに、書き換え可能な FPGA を用いた様々なリコンフィギュラブル・コンピューティング・システム [5](RC システム) の開発が報告されている [6][7][8][9]。これらのシステムの中で最も有名なものは、SPLASH[6] である。汎用計算機が苦手とするパターンマッチング問題を高速化し、一躍注目を浴びた。

本稿では、書き換え可能な FPGA を用い、様々なニューラルネットアルゴリズムに対応可能なシステムを提案する。また、一例として、本システムにバックプロパゲーション(以下 BP) ネットワーク構造の手書き漢字認識アルゴリズム [10] を適用し、その予備評価を行った。

2 BP ネットワーク構造の手書き漢字認識アルゴリズム

2.1 ネットワーク構造

ネットワーク構造は、3層 BP ネットワークとする。3層 BP アルゴリズムは、外部からの入力データを入力層を通じてネットワークに入力し、中間層を経て出力層から出力する。そして、出力データとその入力データに対しての正しい出力(教師信号)に近づくように各層間の重みを修正し、その誤差が許容値以下になるまで繰り返し学習計算する方法である。

本研究では、手書き漢字の認識率を向上するため入力層の特徴ベクトルには外郭方向寄与度特徴(Peripheral Direction Contributivity 以下 PDC)[11]を用いる。また、クラスタリング手法(k-means 法[12])により、17文字の類似文字群(聞, 聞, 聞, 開, 聞, 間, 閑, 閑, 閑, 閃, 陶, 陶, 閑, 聞, 閑, 間, 閃)を選択する。そして、以下の理由より、ネットワーク構造は入力層 768、中間層 34、出力層 17 とする。

1. 入力層

PDC 特徴は文字を大局的にみて、文字線の複雑さ、方向、接続関係、相対位置関係を反映する特徴である。文字を 45 度おきに 8 方向から走査し、横切る各文字線の最初に横切る輪郭点での水平・垂直の各方向寄与度を投影する。この場合、横切る文字線数(外郭深度)、3 本目までの方向寄与度が成分別に投影される。投影軸を 8 区間に等分割し、各区間内の各方向寄与度成分の値を平均することによって区間の寄与度成分とする。さらに、± 45 度の方向寄与度を合計し、入力次元数は $8 \times 2 \times 3 \times 8 \times 2 = 768$ とする。

2. 中間層

出力層の 2 倍の 34 とする。

3. 出力層

識別文字数と同数の 17 とする。

2.2 3層 BP アルゴリズム

3層 BP アルゴリズムを 8 つのステージに分ける。ステージ 0: 入力層-中間層の積和演算 u_j を計算する。

$$u_j = \sum_{i=0}^{767} W_{im_{ij}} x_i \quad (1)$$

ここで、 $W_{im_{ij}}$ は入力層-中間層間荷重係数、 x_i は入力層データである。

ステージ 1: ステージ 0 で計算した u_j のシグモイド関数 m_j を計算する。

$$m_j = f(u_j) = \frac{1}{1 + \exp(-u_j)} \quad (2)$$

ステージ 2: 中間層-出力層の積和演算 v_k を計算する。

$$v_k = \sum_{j=0}^{33} W_{mo_{jk}} m_j \quad (3)$$

ここで、 $W_{mo_{jk}}$ は中間層-出力層間荷重係数である。

ステージ 3: ステージ 2 で計算した v_k のシグモイド関数 o_k を計算する。

$$o_k = f(v_k) = \frac{1}{1 + \exp(-v_k)} \quad (4)$$

ステージ 4: 出力層から中間層への荷重係数を修正する学習データ δo_k を計算する。

$$\delta o_k = (t_k - o_k) o_k (1 - o_k) \quad (5)$$

ステージ 5: 中間層から入力層への荷重係数を修正する学習データ δm_j を計算する。

$$\delta m_j = m_j (1 - m_j) \sum_{k=0}^{16} \delta o_k W_{mo_{jk}} \quad (6)$$

ステージ 6: 中間層-出力層間の荷重係数、結合係数 $\Delta W_{mo_{jk}}$ を修正する。

$$\begin{aligned} \Delta W_{mo_{jk}} &= \eta \delta o_k m_j + \alpha \Delta W_{mo_{jk}} \\ W_{mo_{jk}} &= W_{mo_{jk}} + \Delta W_{mo_{jk}} \end{aligned} \quad (7)$$

ここで、 η は学習定数、 α は安定化定数である。

ステージ 7: 入力層-出力層間の荷重係数、結合係数 $\Delta W_{im_{ij}}$ を修正する。

$$\begin{aligned} \Delta W_{im_{ij}} &= \eta \delta m_j x_i + \alpha \Delta W_{im_{ij}} \\ W_{im_{ij}} &= W_{im_{ij}} + \Delta W_{im_{ij}} \end{aligned} \quad (8)$$

学習計算はステージ 0~7 の計算(フィールドフォワード処理と学習処理)を繰り返す。また、識別時にはステージ 0~3 のフィールドフォワード処理のみを行い、出力値が最大となる出力層ニューロンに学習させた文字を識別結果とする。

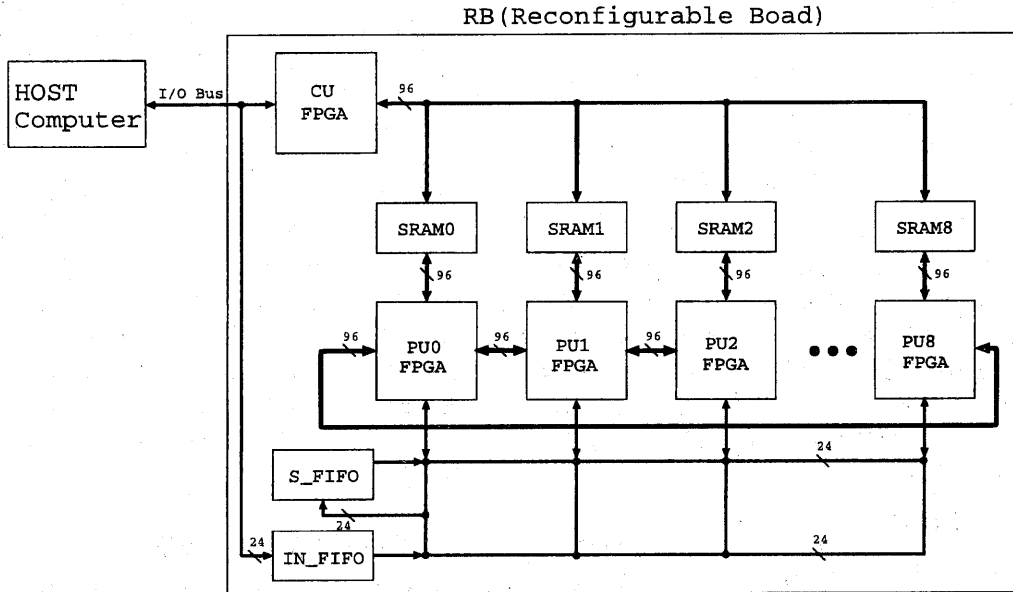


図 1: ニューロシステムの全体構成図

3 システムの全体構成

図1にシステム全体の構成図を示す。本システムは、ホストコンピュータとリコンフィギャラブルボード(以下RB)で構成する。RBは、RB上のバスやI/Oバスなどを制御するControl Unit(以下CU)、ニューラルネットの学習計算などの計算集約的な処理を行う9個のProcessing Unit(以下PU)、データ格納用の9個のSRAMと2個のFIFOで構成する。CUとPUはFPGAである。また、PU-PU間のデータ通信を柔軟にするため、すべてのPUにブロードキャストやマルチキャスト転送を可能にするデータバスと、隣接するPU間にリング状のデータバスを設けた。

各ブロックは、先に説明したBPネットワーク構造の手書き漢字認識アルゴリズムを実現するために、次のように用いる。

1. CU
RB全体の動作を規定するステートマシンを構成。
2. PU
ニューロンを構成。
3. SRAM
結合係数、荷重係数、教師信号、出力層データなどを格納。
4. FIFO
入力層データを格納。

次に、学習計算時のシステム全体の流れを簡単に説明する。

1. 各SRAMを初期設定する。
2. ホストコンピュータは、IN_FIFOに向けて、そのハーフ・フラグがアクティブになるまで一文字分の入力層データ(768ワード)を転送する。
3. CUは、IN_FIFOのエンプティ・フラグがネガティブになったらフィールドフォワード処理を制御する。
4. 各PUは、フィールドフォワード処理を行い、出力層データをSRAMに格納する。CUはホストコンピュータに出力層データを転送する。また、`study_flag`が立っていれば5へ、立ってなければ3へ移行する。ホストコンピュータは次の入力層データ用の教師信号を転送する。
※ `study_flag`はホストコンピュータが制御し、現在の計算が学習計算であるかフィールドフォワード処理のみかを示す。
5. CUが学習処理を制御する。また、ホストコンピュータが誤差($err = err + \sum_{k=0}^{16} (t_k - o_k)^2$)を求める。

これを17(文字)×100(セット)行う。そして、誤差の平均値を求め、前の誤差平均と比較し、許容値より大きければ引続き学習計算を行う。また、小さければ学習計算を終了する。

このように、本システムではハードウェアとホストコンピュータの両方で処理を行う。

4 ハードウェアアーキテクチャ

PUは、24ビット浮動小数点乗算器4個、加算器4個、除算器1個と24ビット幅のレジスタ26個を搭載する。また、各ステージの状態遷移レジスタを持ち、その状態遷移により各演算が切り替わり高速な計算ができる。また、入力層-中間層間の計算は1PU4ニューロン、中間層-出力層間の計算は1PU2ニューロンを構成する。

4.1 FIFO

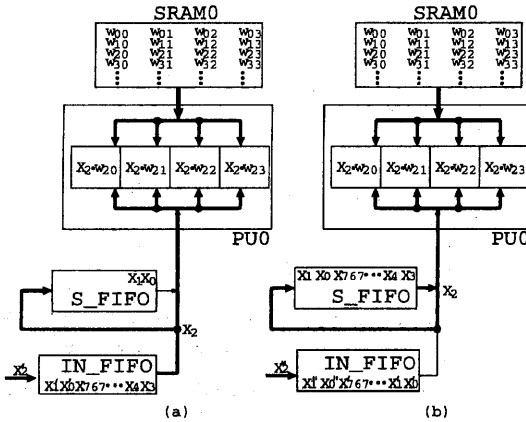


図2: 2つのFIFO

FIFOは入力層データを格納するのに用いる。入力層データを格納するのにFIFOを用いた理由は、以下の通りである。

1. データを格納するとエンプティ・フラグがネガティブになる。
これにより、ホストコンピュータから転送された入力層データがFIFOに格納してあるかどうか分る。
2. FIFOにデータを半分格納するとハーフ・フラグがアクティブになる。
これにより、ホストコンピュータは1文字分の入力層データ(768ワード)を転送できるかどうか分る。
3. アドレスバスが不必要。
RB側の読み出しアドレスが不要となる。

また、FIFOを2つ用いる。1つ(IN_FIFO)はホストコンピュータから転送してきた入力層データを格納するために、もう1つ(S_FIFO)は現在学習している入力層データを格納するために用いる。このようにすることにより、図2に示すようにステージ0、ステージ7で計算中でも、IN_FIFOに次の入力層データ(X_i')、その次の入力層データ(X_i'')とI/O Busが空いるときに格納できる。また、現在学習している入力層データ(X_i)は、図2の(a)のように最初に入力層データを使うときにS_FIFOに格納する。そして、次に使うときからS_FIFOにある入力層データを読み出し、S_FIFOに書き込む。このようにすることにより、ニューロン数が増え36ニューロンを越えた場合でも対応ができる。

4.2 ブロードキャスト転送

中間層データはステージ1で各PUの内部レジスタmregに格納する。そして、図3に示すようにステージ2、ステージ6において、CUに選択されたPUがその中間層データを他のPUにブロードキャストする。

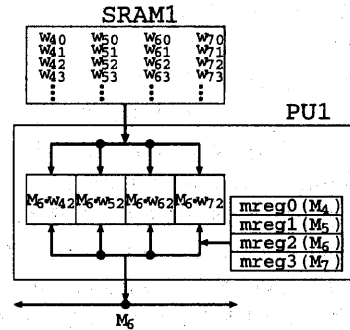


図3: ブロードキャスト転送

4.3 出力層-中間層間学習信号

ステージ5で行う計算(6)式で、各PUのローカルメモリSRAMはすべての結合係数データ Wm_{ojk} を持っていない。そこで、図4に示した後向き積和演算の計算方法[1](中間層-出力層の計算を1PU4ニューロンとしたときの例)のように、内部レジスタbregに格納してある中間層-出力層間学習データ δo_k とSRAMから読み込んだ結合係数データ Wm_{ojk} から $\delta o_k Wm_{ojk}$ を求め、隣接するPUに転送しながら計算していく。

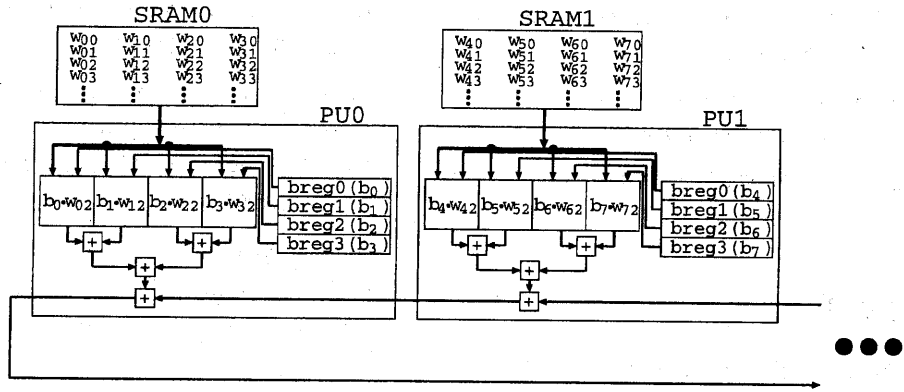


図 4: 後向き積和演算の計算方法

4.4 シグモイド計算

シグモイド関数を計算するために(9)式のようにテイラー展開を用いる。

$$f(x) = \frac{1}{1+e^{-x}} = \frac{1}{2 + \frac{-x}{1!} + \frac{(-x)^2}{2!} + \dots + \frac{(-x)^n}{n!}} \quad (9)$$

ステージ1とステージ3では、内部レジスタに格納してある重み付き和 u_j, v_k を入力とし、4個の乗算器と加算器、そして1個の除算器で2段のパイプラインを構成し計算する。また、シグモイド関数は点(0,0.5)に対称であることと、(9)式の場合入力データが負の数の方が収束しやすいことより、入力が負の数の場合そのまま出力し、正の数の場合は $1-f(-x)$ を出力する。

5 予備評価

表 1: ワークステーションとの比較

ステージ	WS(O2)	RB	比率
S0	261834	769	340.5
S1	---	47	---
S2	5563	37	150.3
S3	---	27	---
S4	179	3	59.7
S5	5202	46	113.1
S6	11757	88	133.6
S7	522589	3075	226.5
Sum-(S1+S3)	807124	4018	200.9

5.1 演算サイクル数

本システムの性能を評価するために、各演算ステージの演算サイクル数を計算し、ワークステーション(以下 WS)の演算サイクル数と比較した結果を表1に

示す。表1からわかるようにほとんどの計算時間は、ステージ0とステージ7である。本システムの演算サイクル数は、WSの約1/200であることが分る。なお、WSの演算サイクルは、C言語プログラムコンパイル結果から、一命令を一サイクルと仮定して求めた。コンパイル時の最適化にはO2を用いた。また、WSでのシグモイド計算は、指数関数がxの値によって演算サイクル数が違うのでS1とS3は比較をしていない。

5.2 論理合成

RBをVerilog-HDLを用いて記述した。そして、Verilog-XLシミュレータで動作確認を行い、PUを論理合成した。論理合成の結果を表2に示す。論理合成には、SynopsysのFPGA EXPRESSを用い、ターゲットFPGAはXILINXのXC40125XV-1(4,624clb)とした。

表 2: PUの仕様

ターゲットFPGA	XC40125XV-1
パッケージ	499ピンPG
CLB数	4,241
クロック	18MHz(50.5ns)

5.3 パフォーマンス評価

本稿でシステムに実装したネットワークは、入力層768、中間層34、出力層17である。これの総結合数は26,690である。また、表1より学習演算処理での演算サイクル数は4092、フォワード演算処理(S0からS3)は880、表2より一演算サイクルは50.5nsである。よって、学習演算処理能力CUPS(Connections Update Per Second)とフォーワ

ド演算処理能力 CPS(Connections Per Second) は、以下のように計算できる。

$$\begin{aligned}CUPS &= 26,690/(4092 * 50.5 * 10^{-9}) \\ &= 129.2MCUPS \\ CPS &= 26,690/(880 * 50.5 * 10^{-9}) \\ &= 600.6MCPS\end{aligned}$$

6 むすび

本稿では、書き換え可能な FPGA を用い、様々なニューラルネットワークアルゴリズムに対応可能なシステムを提案した。また、一例として、本システムに BP ネットワーク構造の手書き漢字認識アルゴリズムを適用した時のアーキテクチャについて述べた。そして、処理を行う PU の FPGA を XILINX 社の XC40125XV-1 として論理合成し、予備評価を行った。その結果、演算サイクル数はワークステーション上の C プログラムと比べて約 1/200 となり、学習演算処理能力が 129.2MCUPS、フォワード演算処理能力が 600.6MCPS という結果が得られた。

今後は、ニューラルネットワークに認識させる漢字数を増やし、本システムで大規模 BP ネットワークを構成する場合について評価する。また、他のニューラルネットワークアルゴリズムを適用し、本システムが様々なアルゴリズムに対応できることを示していく予定である。

謝辞日頃より御指導、御助力を頂いた宇都宮大学情報工学科馬場、吉永、大津研究室の諸氏に感謝する。

本研究の一部は東京大学 大規模集積システム設計教育研究センターより提供して頂いた CAD ツールを使用した。深く感謝する。

本研究は一部文部省科学研究費 基盤研究 (C) 課題番号 09680324、基盤研究 (B) 課題番号 10558039、奨励研究 (A) 課題番号 09780237 の援助による。

参考文献

- [1] 広瀬佳生, 安仏英明, 山下公一, 後藤源助, “バックプロパゲーション専用プロセッサのアーキテクチャ”, 信学技報 ICD92-17, pp.40-46, 1992.
- [2] Y.Kondo, et al., “A 1.2FLOPS Neural Network Chip Exhibiting Fast Convergence,” ISSCC '94 IEEE International Solid-State Circuits Conference, TP 13.1, pp.218-219, Feb.1994.
- [3] 齋藤修, 相原公久, 藤田修, 内村国治, “大規模ニューラルネットワークに対応可能な学習機能内蔵デジタルニューロチップ” 信学技報 CPSY98-4, pp.25-32, 1998.
- [4] 柴田克成, 安永守利, 大山光男, 益田昇, 柳生正義, 浅井光男, 山田稔, 坂口隆宏, 橋本雅, “高速学習型ニューロ WSI のシステム設計”, 信学技報 ICD90-127, pp.49-56, 1990.
- [5] 末吉敏則, “Reconfigurable Computing System の現状と課題-Computer Evolution へ向けて-”, 信学技報 CPSY96-91, pp.111-118, 1996-12.
- [6] Maya Gokhale, Bill Holmes, Andrew Kopser, Dick Kunze, Dan Lopresti, Sara Lucas, Ron Minnich, “SPLASH: A Reconfigurable Linear Logic Array”, January 26, 1994.
- [7] M.Wazlowski, L.Agarwal, T.Lee, A.Smith, E.Lam, P.Athanas, H.Silverman, S.Ghosh, “PRISM-II Compiler and Architecture”, Proc. IEEE Workshop on FPGAs for Custom Computing Machines, pp.147-157, 1994.
- [8] Jean-Luc Beuchat, Jacques-Oliver Haenni and Eduardo Sanchez, “Hardware Reconfigurable Neural Networks”, LNCS1388. Rolim Parallel and Distributed Processing, 1998.
- [9] Michael J. Wirthlin, Brad L. Hutchings, “DISC: The dynamic instruction set computer”, Proc.SPIE-Int.Soc.Opt.Eng.(USA), vol.2607, pp.92-103, 1995.
- [10] 荒井正之, 王晋申, 奥田健三, 宮道壽一, “Honeycomb ネットによる多字種の手書き漢字認識”, 信学論 (D-II), J77-D-II, 9, pp.1708-1715, 1994-9.
- [11] 萩田泰一, 内藤誠一郎, 増田功, “外郭方向寄与度特徴による手書き漢字の識別”, 信学論, J66-D, 10, pp.1185-1192, 1983-10.
- [12] ハイ, 森下, 蕪山, 伊崎, 山本, “手書き漢字認識におけるテンプレート複数化の検討”, 信学技報, PRL81-42, pp.49-56, 1981-9.