

全ハードウェア形高速ニューロプロセッサの設計

苦米地 宣 裕

Design of Full-Hardware High-Speed Neuro-Processors

Nobuhiro TOMABECHI

Abstract

This paper shows the design of the high speed neuro-processors in which 200~1,000 neurons are integrated and all of the arithmetic operations and the signal transmission are realized by hardware technology. The sizes of the processors become wafer scale. It is concluded that a neuro-processor consisted with 500 neurons may be implemented on a single wafer.

Keywords: full-hardware/high-speed/neuro-processor/design

1. ま え が き

ニューラルネットワークのリアルタイム制御など高速処理への応用が期待されている^[1]。通常、ニューラルネットワーク機能は汎用プロセッサを用いてソフトウェアによって実現されているが、高速処理を実現するためには、ハードウェアで構成されたニューロン(実ニューロン)を搭載した専用プロセッサが必要となる。これまで、数十ないし数百個の実ニューロンを搭載したプロセッサの開発が報告されているが^{[2]-[5]}、より一層多くのニューロンを集積したプロセッサが望まれている。また、ニューロプロセッサではニューロン間の相互配線が大きなチップ面積となるため、信号伝送を時分割で行うことが多く、その分プロセッサの演算速度が低下する。

本論文では、200~1,000個の実ニューロンを搭載し、しかも、演算機能や信号伝送機能は時分割することなしに、すべてハードウェアで構成するニューロプロセッサの設計を行っている^{[6]-[7]}。これらのプロセッサはWSIの規模となる。結論として、現在の技術水準で、ニューロンを500個搭載したプロセッサが1枚のウエーハ上に構成できる可能性のあることが分かった。

なお、WSIの製造においては、欠陥の増加による歩留り低下が最大の問題となる。その対策として、冗長構成を導入して欠陥救済を行う方法が考えられるが、本論文では、第1段階として、欠陥救済を行わない場合について論じている。

2. 仕様の設定

本論文では、リアルタイム制御などに用いる高速実行型のニューロプロセッサを対象とする。その仕様を表1に示している。高速性を達成するために、演算機能や信

号伝送機能は時分割することなしに、すべてハードウェアで実現することとしている。また、学習機能はなしとしている。このとき、学習は他の汎用プロセッサを用いて行い、そこで得られたシナプスの重み係数を本プロセッサ内部の重みメモリに入力するという使い方をする。表1に示す3層フィードフォワードニューロプロセッサの動作は次式で表される。但し、 y_j^k , $w_{i,j}^k$, h_j^k , $f(\cdot)$, および、 N は、それぞれ、第 k 層第 j 番目のニューロンの出力、第 k 層第 j 番目のニューロンの第 i 番目の入力に対する重み、第 k 層第 j 番目のニューロンのしきい値、シグモイド関数、および、一つの層のニューロンの数を表している。

$$y_j^k = f\left(\sum_{i=1}^N w_{i,j}^k y_i^{k-1} - h_j^k\right) \quad k=1, 2, \quad j=1, 2, \dots, N \quad (1)$$

なお、 $k=1$ は中間層、 $k=2$ は出力層に対応し、 y_i^0 は入力層 i 番目の出力を表している。

表1 ニューロプロセッサの仕様

ソフトウェア/ハードウェア	ハードウェア
アナログ/デジタル	デジタル
時分割/全ハードウェア	全ハードウェア
形式	3層フィードフォワード型
ニューロンの数 N	中間層・出力層とも 100~500
ニューロン間の配線	完全配線
1ニューロンのシナプスの数	各層のニューロン数に同じ
入出力データの精度	16ビット
重みメモリビット数	16ビット
素子	CMOS
学習機能	なし(シミュレーションで学習)

3. アーキテクチャ設計

式(1)に基づくニューロプロセッサのシステム構成を図1に、1個のニューロンの回路構成を図2に示している

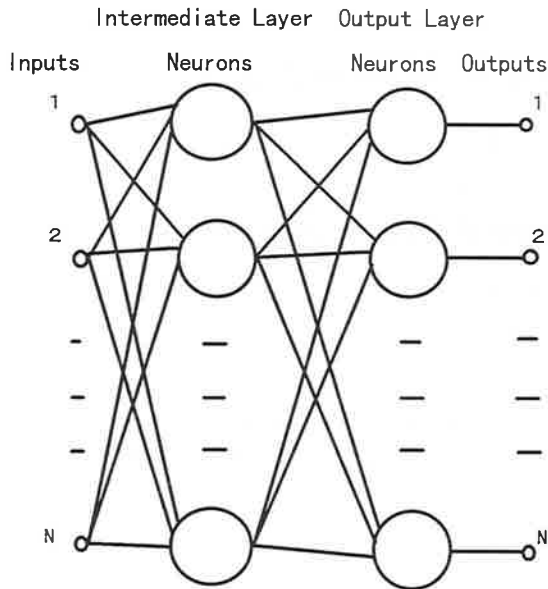


図1 ニューロプロセッサの全体構成

Multipliers

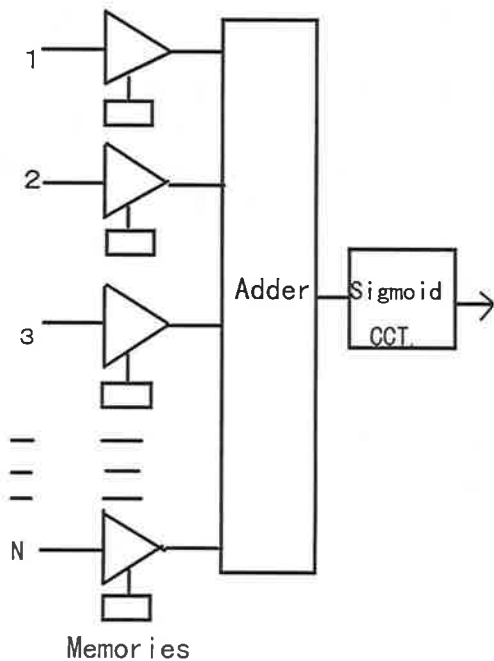


図2 ニューロンの回路構成

る。本論文では、ニューロン内の各回路を以下のように構成する。

① 乗算器

負数の演算が補数のままで実行できる Baugh-Wooley の乗算器を用いる。本乗算器は、入力信号を n ビット $\times n$ ビットとすると、 $n(n-1)+3 \approx n^2$ 個の全加算器 (但し、ANDゲート1個付) で構成される。図3に、 $(a_4 a_3 a_2 a_1 a_0) \times (b_4 b_3 b_2 b_1 b_0)$ の構成例を示している^[8]。

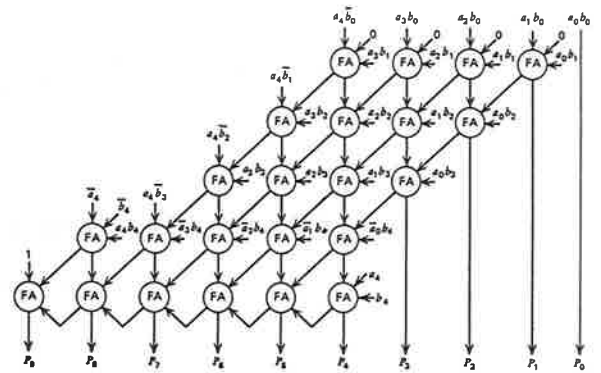


図3 乗算器の構成

② 重みメモリ

プロセッサの外部からデータを入力するのに好都合となるようシフトレジスタで構成し、すべてのレジスタを直列接続して用いる。以下、重みメモリは乗算器に含めて取り扱う。

③ 加算器

$N=100 \sim 500$ 個という多数のデータをできるだけ少ないハードウェアと少ない段数で加算する点が問題となる。本論文では、2入力加算器 (但し、各入力が16ビット) を用いて、2個ずつのデータを加算していく方法、いわゆる加算器の2進木構造を用いることとする。本構成に必要な2入力加算器の数 N_A は次式で与えられる。

$$N_A = N - 1 \quad (2)$$

2入力加算器が直列接続される段数 S は次式で与えられる。但し、 $\text{INT}[X]$ は、 X 以上の値となる最小の整数を表している。

$$S = \text{INT}[\log_2(N-1)] \quad (3)$$

2入力加算器の具体的回路は、リップルキャリアダーとする。リップキャリアダーはキャリーの伝播に時間がかかるという問題を有するが、多数の加算器を用いて2進木構造とした場合は最終段のアッダー1個分のキャリー伝播時間が問題となるだけである。

④ シグモイド関数回路

いくつかの直線、または、いくつかの2次曲線により近似する^[9]。ここでは5個の直線で近似することとする。このとき、直線を $y = a_i x + b_i (i=1, 2, \dots, 5)$ と表わすと、シグモイド関数回路は、 x の値の範囲を判定する5個の比較器、 a_i, b_i の選択ゲート、および、 $(a_i x + b_i)$ の計算回路で構成される。なお、しきい値のメモリレジスタとしきい値の減算器はシグモイド関数回路に含めて取り扱う。

4. レイアウト設計

(1) フロアプラン

① ニューロプロセッサ全体のフロアプラン

図4に示している。100~500個のニューロンが縦方向に並ぶので、ニューロン1個のレイアウトは、可能な限り横長に設計する必要がある。

② ニューロン1個のフロアプラン

図5に示している。レイアウトを横長にするため、すべての乗算器を横方向に1列に配置している。入力バスラインはバラバラにして、各々の乗算器に入力される配線だけをそれぞれの乗算器の直前に配置している。2入力加算器は、乗算器2個ごとに1個挿入する。また、2入力加算器2個に対して1個の2入力加算器を付加していく。2入力加算器の入出力配線は、16ビットをまとめで、1ビットごとに全加算器(以下、FAという)とFAの間を通過させる。

(2) レイアウト設計

レイアウト設計規則は、Mead-Conwayの設計規則^[10]をCMOSに拡張した規則^[11]を用いる。図6に、ANDゲート1個付きFAのレイアウトを示している。但し、 λ は基準寸法を表している。FAを元にした各回路の形状を表2に示している。図7は、ニューロン1個のレイアウトを示している。但し、回路の種類ごとにまとめた寸法を示している。結局、ニューロプロセッサの縦の長さ h 、横の長さ l 、および、面積 A_p は次のようになる。

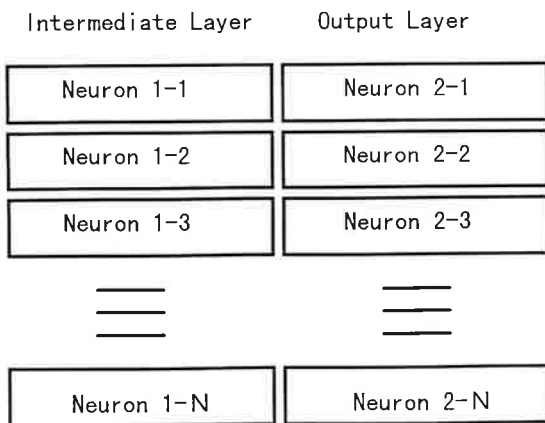
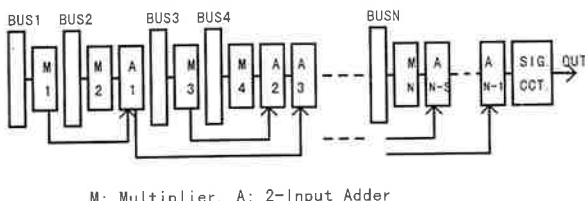


図4 ニューロプロセッサのフロアプラン



M: Multiplier, A: 2-Input Adder

図5 ニューロンのフロアプラン

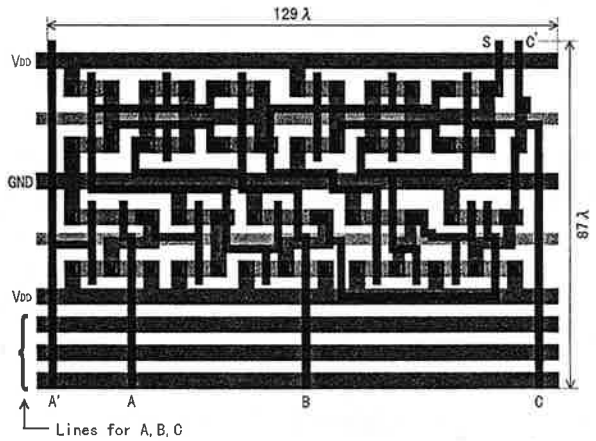


図6 FAのレイアウト

表2 各回路の形状とチップ面積

回路	形状 縦(λ)×横(λ)	チップ面積 (λ ²)
FA	129×87	1.12×10 ⁴
乗算器	2,064×1,392	2.88×10 ⁶
重みメモリ	41×66	2.71×10 ³
2入力加算器	2,064×87	1.80×10 ⁵
シグモイド回路	2,064×2,349	4.85×10 ⁶
入力バスライン	2,064×(112×N)	2.31 N×10 ⁵

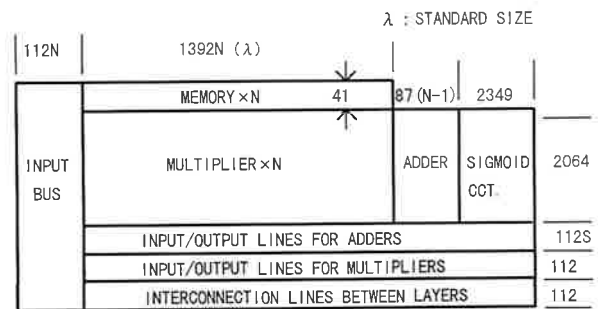


図7 ニューロンのレイアウト

$$h = (112S + 2329)N \quad (\lambda)$$

$$l = 2(1591N + 2262) \quad \lambda$$

$$A_p = 2N(112S + 2329)(1591N + 2262) \quad (\lambda^2) \quad (4)$$

例として、 $N=100$ とすると、縦横比は、縦 311,300 対横 322,724 \approx 1 対 1 となり、これは WSI に適した形状といえる。式(4)より、 h, l とともに N にほぼ比例するので、 h 対 $l \approx 1$ 対 1 の関係は、 N の値によらず成り立つ。従って、ここに示したレイアウトは、任意の N の値に対して有効であるということが出来る。

また、 $N=100$ の場合、ニューロプロセッサの面積 A_p は $1.00 \times 10^{11} (\lambda^2)$ となるが、これは、基準 VLSI^[12] のチップ面積の 33 倍となり、本プロセッサが WSI に近い規模であることを示している。

5. 考 察

式(4)より、ニューロプロセッサの面積 A_P は N^2 で増大する。図8に、全ニューロン数 $T(=2N)$ 対 A_P の関係を示している。今、利用可能なウエーハの最大寸法を10インチ、 $\lambda=0.2\mu$ とすると、WSIの最大面積 A_M は、 $(25.4/(0.2 \times 10^{-4}))^2/2=8.06 \times 10^{11}(\lambda^2)$ となる。図8より、 $A_P=A_M$ となる T の値を求めると約500であることが分かる。よって、本論文で得られた結果を次のようにまとめることができる。

[結論] ニューロンを500個搭載したプロセッサが1枚のウエーハ上に集積できる可能性がある。

WSIの製造では、製造欠陥の増加による歩留りの低下が最大の問題となる。その対策として、冗長構成を導入して欠陥救済を行うことが考えられる^[7]。この点については別報にて論ずる。

6. む す び

本論文では、ニューロンを200~1,000個搭載し、しかも、すべての演算機能をハードウェアで構成する高速ニューロプロセッサの設計を行った。その結果、現在の技術水準で、ニューロンを500個搭載したプロセッサが1枚のウエーハ上に集積できる可能性があることが分かった。

今後、学習回路を組み込んだ場合について設計を行う予定である。

また、WSIの製造では、製造欠陥の増加による歩留りの低下が最大の問題となるが、この問題に対しては、冗長構成を導入して欠陥救済を行う方法を検討しており、別途報告する予定である。

なお、本研究は平成8年度文部省科学研究費(基盤研究C)の補助を一部受けたことを付記する。

文 献

[1] 岩田, 雨宮, “ニューラルネットワーク LSI,” 電子情報通信学会, 1995.

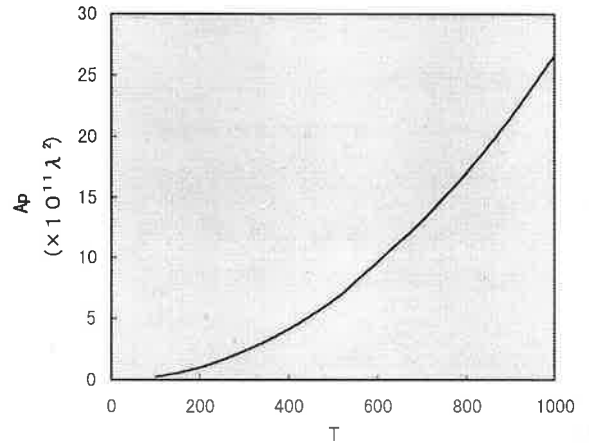


図8 全ニューロン数 T 対 A_P の関係

[2] H. Kato, H. Yoshizawa, et. al., “A parallel neurocomputer architecture towards billion connection update per second,” Proc. *IJCNN '90*, pp. 47-50, 1990.

[3] M. Yasunaga, N. Masuda, et. al., “Design fabrication and evaluation of a 5-inch wafer scale neural network LSI composed of 576 digital neurons,” Proc. *IJCNN '90*, pp. 527-535, 1990.

[4] M. Griffin, G. Tahara et. al., “An 11-million transistor neural network execution engine,” *ISSCC '91*, pp. 180-181, 1991.

[5] Y. Kondo, Y. Koshiba et. al., “A 1.2GFLOPS neural network chip exhibiting fast convergence,” *ISSCC '94*, pp. 218-219, 1994.

[6] 吉米地宣裕, “フィードフォワードニューロ WSI の階層冗長構成,” 電子情報通信学会 1997 年総合大会講演論文集, p. D-2-20, 1997.

[7] 吉米地宣裕, “冗長化による WSI 規模ニューロプロセッサの構成法,” 電子情報通信学会技術研究報告, vol. ICD97-30, pp. 57-64, May 1997.

[8] Kai Hwang, “Computer Arithmetic,” John Wisely & Sons, Inc., 1979.

[9] 林原香織他, “シグモイド関数の連続性/離散性とニューラルネットワークのマシ能力について,” 電子情報通信学会論文誌 D-II, vol. J73-D-II, no. 8, pp. 1220-1226, Aug. 1990.

[10] C. Mead and L. Conway, “Introduction to VLSI systems,” Addison-Wesley Pub., 1980.

[11] 松山泰男, 富沢 孝, “VLSI 設計入門,” 共立出版, 1983.

[12] 吉米地宣裕, 金沢正治, “WSI 規模高速 FFT プロセッサの冗長化,” 電子情報通信学会論文誌 D-I, vol. J80-D-I, no. 1, pp. 110-120, Jan. 1997.