

A New Built-In Self-Repair Approach to VLSI Memory Yield Enhancement by Using Neural-Type Circuits

Pinaki Mazumder, *Member, IEEE*, and Yih-Shyr Jih, *Member, IEEE*

Abstract—As VLSI chip size is rapidly increasing, more and more circuit components are becoming inaccessible for external testing, diagnosis, and repair. Memory arrays are widely used in VLSI chips, and restructuring of partially faulty arrays by the available spare rows and columns is a computationally intractable problem. Conventional software memory-repair algorithms cannot be readily implemented within a VLSI chip to diagnose and repair these faulty memory arrays. Intelligent hardware based on a neural-network model provides an effective solution for such built-in self-repair (BISR) applications. This paper clearly demonstrates how to represent the objective function of the memory repair problem as a neural-network energy function, and how to exploit the neural network's convergence property for deriving optimal repair solutions. Two algorithms have been developed using a neural network, and their performances are compared with the repair most (RM) algorithm that is commonly used by memory chip manufacturers. For randomly generated defect patterns, the proposed algorithm with a hill-climbing (HC) capability has been found to be successful in repairing memory arrays in 98% cases, as opposed to RM's 20% cases. The paper also demonstrates how, by using very small silicon overhead, one can implement this algorithm in hardware within a VLSI chip for BISR of memory arrays. The proposed auto-repair approach is shown to improve the VLSI chip yield by a significant factor, and it can also improve the life span of the chip by automatically restructuring its memory arrays in the event of sporadic cell failures during the field use.

I. INTRODUCTION

AS THE VLSI device technological feature width is rapidly decreasing to the range of hundreds of nanometers, the chip yield tends to reduce progressively due to increased chip area, complex fabrication processes, and shrinking device geometries. In order to salvage partially faulty chips, redundant circuit elements are incorporated in the chips, and an appropriate reconfiguration scheme is employed to bypass and replace the faulty elements. In high-density dynamic random-access mem-

ory (DRAM) chips, redundant rows and columns are added to reconfigure the memory subarrays, where the rows or columns in which defective cells appear, are eliminated by using techniques such as electrically programmable latches or laser personalization. The problem of optimal reconfiguration and spare allocation has been widely studied by many researchers [1], [9], [10], [14], [15], [29]. A review of these algorithms for memory diagnosis and repair can be seen in [4]. But all of these algorithms cannot be readily applied to embedded arrays, where they are neither controllable by external testers, nor are their responses readily observable. Built-in self-test (BIST) circuits are commonly used to comprehensively test such embedded arrays and bad chips are discarded when they fail to pass the relevant test procedures. In order to salvage the partially faulty chips, new built-in self-repair (BISR) circuits must be developed for optimal repair of the faulty memory arrays.

In this paper, a novel self-repair scheme is proposed that uses a built-in neural network to determine how to automatically repair a faulty memory array by utilizing the spare rows and columns. Two algorithms have been proposed to illustrate how a neural network can solve random defects, and their performances are compared with the conventional software repair algorithms, such as *Repair Most* [28]. The Repair Most (RM) algorithm is a greedy algorithm that iteratively assigns a spare row (column) to replace the row (column) which currently has the maximum uncovered cells until all the cells are covered or bypassed (providing a successful solution), or all the spare rows and columns are exhausted (failing to give a repairable solution, if all faulty cells are not bypassed), whichever occurs first. This algorithm is selected for comparison because it is relatively simple and can be implemented by digital hardware, unlike other software algorithms, such as simulated annealing or approximation algorithms. A hardware version of simulated annealing will require a Gaussian white noise generator to generate the random moves and a logarithmic circuit to schedule the temperature. The other heuristic-based algorithms (like FLCA, LECA, and BECA [15], and branch-and-bound [14]) will require microcode-driven complex digital circuits for solving the problem by hardware. The performance of neural algorithms has been compared with

Manuscript received May 1, 1990; revised January 16, 1992. This work was supported in part by the National Science Foundation under Grant MIPS 9013092, by the Office of Naval Research under Grant 00014-85-K-0122, and by the U.S. Army URI Program under Grant DAAL 03-87-K-0007. This paper was recommended by Associate Editor F. Brglez.

P. Mazumder is with the Department of Electrical Engineering and Computer Science, The University of Michigan, Ann Arbor, MI 48109.

Y.-S. Jih was with the University of Michigan, Ann Arbor, MI. He is now with IBM T.J. Watson Research Center, Yorktown Heights, NY.

IEEE Log Number 9200909.

RM by running these algorithms with repairable fault patterns and examining what percentages of these fault patterns can be repaired by them. The simulation studies made in this paper demonstrate that a neural network using a hill-climbing (HC) capability provides a fast and good solution for a repairable array. The technique is implemented within the chip by an electronic neural network, and it promises to be a potential application area for neural networks.

Historically, the memory repair problem dates back to the evolution of 64 Kbit DRAM's in the late 1970's [25], when, to improve the chip yield, two or three extra rows and extra columns were added to each quadrant consisting of 64 Kbit subarray of memory cells. Simple greedy [2], [28] and exhaustive search [1] algorithms were developed to repair the faulty memory subarrays. Since the search space was very small for such problem sizes, all these algorithms provided high throughput in memory diagnosis and repair. But as the memory size has increased to several megabits over the last few years, the defect patterns have become sufficiently complex and the search space has grown extensively. The problem of memory repair has been shown to be NP-complete [9], [14] by demonstrating that the repair problem is transformable in polynomial time to the Constrained Clique problem in a bipartite graph. A number of heuristic algorithms, such as branch-and-bound [14], approximation [14], best-first search [10], and others [15], [29], recently have been proposed to solve the memory array repair problem. The two key limitations of these algorithms are a) that their worst-case complexities are nearly exponential and b) they are not readily implementable within the chip for BISR. These algorithms are generally written in a high-level programming language and are executable on general-purpose computers. This paper addresses two fundamental aspects of the memory repair problem: how to devise efficient algorithms so that the overall production throughput improves along with the chip yield, and how to generate such repair algorithms in hardware so that they can be applied to repairing memories, embedded within the VLSI chips. The contribution of this paper is to demonstrate how to solve the array repair problem by using neural networks, and how to implement a BISR scheme in hardware. A neural network can produce solutions by the collective computing of its neuron processors with far faster speed than the abovementioned sequential algorithms running on conventional computers. Since these neural processors are simple threshold devices, the basic computing step of neural nets is comparable to the on-off switching of a transistor [7], [8], [22], [24]. Another potential advantage of using neural networks is that they are robust and fault-tolerant in the sense that they may compute correctly even if some components fail and/or the exciting signals are partially corrupted by noise. Thus the reliability of the self-repair circuit using electronic neural networks is very high. This has been experimentally verified in Section VI.

This paper has been organized as follows: Section II provides a brief overview of the neural network model

and its dynamic behavior. Section III provides a formal framework for the memory repair problem using the concepts of neural net computation. Two algorithms are developed by programming the synaptic weight matrix of the neural net. The first algorithm is greedy, and starting from any random configuration (i.e., arbitrary neuron states), it can solve the problem by monotonically reducing the energy function of the neural net. The second algorithm uses the HC technique to yield a near-optimum solution. Section IV gives the simulation results of the neural net algorithms. From the simulation experiments, it is seen that the neural net algorithms are superior to the RM algorithm. An electronic implementation of neural networks is demonstrated in Section V. The intrinsic fault-tolerant ability of a neural net is studied in Section VI. The yield improvement and hardware overhead caused by the neural-net self-repair circuitry are examined in Section VII.

II. NEURAL NET COMPUTATION MODEL

The earliest mathematical formalization of the human nervous system can be found in the work by McCulloch and Pitts [19], where it was modeled as a set of *neurons* (computation elements) interconnected by *synapses* (weighted links). Each neuron in the model is capable of receiving impulses as input from and firing impulses as output to potentially all neurons in the network. The output function of a neuron depends on whether the total amount of input excitation received exceeds its predetermined threshold value θ , or not. The state of neuron i is represented by an all-or-none binary value s_i , with $s_i = 0$ being a nonfiring condition and $s_i = 1$ denoting an impulse-firing condition. Note that one may represent a non-zero threshold by applying an external impulse bias b_i to the neuron i with $b_i = \theta_i$.

Interactions between two different neurons occur through a synapse serving as a duplex communication channel. The signals passing through the synapse in both directions are considered independent. Moreover, the signal traveling from neuron i to neuron j is amplified by a weight factor w_{ji} , i.e., if the impulses fired by a neuron correspond to a unit influence, then a firing neuron i produces w_{ji} amount of influence to neuron j .

Now, let s_i denote neuron i 's next state value. From the above description one can represent the neural net *state transition function* as follows:

$$s'_i = \begin{cases} 0, & \text{if } \sum_j w_{ij}s_j < \theta_i \\ 1, & \text{if } \sum_j w_{ij}s_j > \theta_i \\ s_i, & \text{otherwise.} \end{cases}$$

Actually, in the human nervous system or any other biological implementations, neural networks make smooth continuous transitions when the neurons change their states. Therefore, we assume that neurons which receive a larger influence in magnitude will have a faster state

transition. Neuron processors are expected to operate in an asynchronous and random fashion.

Following the description of the behavior of an individual neuron, let us look at the overall network state transition behavior. First, a network state is denoted by the vector of neuron state variables. Given a network state vector $X = (x_1, x_2, \dots, x_N)$, X is called a *fixed point*, if and only if $x'_i = x_i$ for all i 's. That is, of all the possible 2^N states of a neural net, only the fixed points are considered stable. If the current network state is not one of the fixed points, the neural net cannot maintain the present state. One immediate question about the network behavior could concern the state convergence: starting with an arbitrary state, will the network eventually reach a fixed point? To answer the question, an energy function has been formulated. First, let us consider the case of an excitatory or positive input to a neuron. According to the neuron's functional definition, this input will encourage the neuron to fire impulses. If the neuron is already in the firing state, this input can be looked upon as negative energy, since by convention, a system is more stable when the energy level is lower. Likewise, an inhibitory or negative input to a neuron in the nonfiring state should also be considered as negative energy. On the other hand, positive energy is created if a firing neuron receives an inhibitory input or a nonfiring neuron receives an excitatory input, since the network is potentially destabilized by the input.

If $w_{ij} = w_{ji}$, for all i 's and j 's, the total energy E^{NA} , commonly known as the Lyapunov function [11], is given by

$$E^{NA} = -\frac{1}{2} \sum_i \sum_j s_i w_{ij} s_j - \sum_i b_i s_i.$$

Note that the change of state by neuron i will result in an absolute decrease of $\sum_j w_{ij} s_j + f(s_i) b_i$ in energy, where $f(s_i) = 1$ if $s_i = 0$, and -1 otherwise. Together with the fact that the total energy is bounded, the network is guaranteed to reach a local minimal energy fixed point [12].

III. NEURAL NETWORK SOLUTIONS

The idea of using a neural network to tackle combinatorial optimization problems was first proposed by Hopfield [12], who designed and fabricated a neural network for solving the classic *Traveling Salesman Problem* [5]. The objective (cost) function of a combinatorial optimization problem can be represented by Lyapunov's energy function of the neural network, and the network's property of convergence from a random initial state to a local minimal energy state can be utilized to reduce the cost function of the combinatorial problem. However, the neural computing approach given in [12] leads only to the design of a greedy gradient descent (GD) algorithm which stops searching at the first local minimal solution encountered; consequently, the neural network solution is generally of low quality.

In this section we study the convergence behavior of neural networks. We demonstrate that more powerful

searching strategies can be developed by modifying the simple GD algorithm. By programming the neural network in an appropriate way, we find that the neural nets are capable of searching the solution by HC or tunneling. Thus, the probability of obtaining a globally optimal solution by using neural networks is very high.

A. Array Repair Problem Representation

Assume a memory array of size $N \times N$ with exactly p spare rows and q spare columns that can be utilized to replace the defective rows and columns. Assume an arbitrary fault pattern in which faulty cells randomly occur on $m (< N)$ distinct rows and $n (< N)$ distinct columns of the memory array, such that the compacted subarray of size $m \times n$ contains all relevant rows and columns which have at least one faulty cell. Let the matrix $D = \{d_{ij}\}$ of size $m \times n$ characterize the status of the memory cells, such that d_{ij} , which corresponds to the cell at row i and column j , is 1 if the cell is faulty, and $d_{ij} = 0$, otherwise. Normally, very few cells in the array are faulty, and the characteristic matrix, D , is highly sparse. A row repair scheme can be represented as a vector U of m bits, such that $u_i = 0$ if row i is to be replaced, otherwise $u_i = 1$, $0 \leq i \leq m - 1$. A column repair scheme is defined in the same fashion, and is denoted by a vector V of n bits. The numbers of 0's in U and V must be less than or equal to p and q , respectively. The memory is said to be repairable if the above constraints are satisfied, and essentially the repair problem is how to determine a pair of U and V , such that $U^T D V = 0$.

In addition to the characteristic matrix D , let the overall status of row i and column j of the array be characterized by r_i and k_j , respectively,

$$r_i = \begin{cases} 1, & \text{if row } i \text{ contains defective elements} \\ 0, & \text{otherwise} \end{cases}$$

$$k_j = \begin{cases} 1, & \text{if column } j \text{ contains defective elements} \\ 0, & \text{otherwise.} \end{cases}$$

To design a neural net for the memory repair problem, a neural net of size $M = m + n$ is used, where the states of the first m neurons are denoted by s_{1i} 's, $1 \leq i \leq m$, and the states of the remaining n neurons by s_{2j} 's, $1 \leq j \leq n$. For ease of interpretation, the first m neurons are addressed as *row neurons*, and the remaining n neurons as *column neurons*. The physical meaning of these neuron states are defined as follows:

$$s_{1i} = \begin{cases} 1, & \text{if row } i \text{ is suggested for replacement} \\ 0, & \text{otherwise} \end{cases}$$

$$s_{2j} = \begin{cases} 1, & \text{if column } j \text{ is suggested for replacement} \\ 0, & \text{otherwise.} \end{cases}$$

The cost function C^{MR} of the memory repair problem should represent both the nonfeasible repairing and incomplete coverage schemes by higher costs. These two

aspects can be modeled by the following two expressions:

$$C_1 = A/2 \left[\left(\sum_{i=1}^m s_{1i} \right) - p \right]^2 + A/2 \left[\left(\sum_{j=1}^n s_{2j} \right) - q \right]^2$$

$$C_2 = B/2 \left[\sum_{i=1}^m \sum_{j=1}^n d_{ij} (1 - s_{1i})(1 - s_{2j}) \right. \\ \left. + \sum_{j=1}^n \sum_{i=1}^m d_{ji} (1 - s_{1j})(1 - s_{2i}) \right]$$

so that the overall cost function is represented as an algebraic sum of C_1 and C_2 . In the above expressions, the values of A and B are empirically decided and as shown in Section IV, their relative values are critical for obtaining the appropriate percent of successful repairs.

The expression for C_1 encourages exactly p spare rows and q spare columns to be used in the repair scheme. For those nonfeasible schemes that require more than the available spare rows and/or columns, the cost will increase quadratically. If necessary, the optimum usage of spares can be experimentally attempted by repeating the search operation with successively reduced amounts of available spares (i.e., progressively decrementing the value of p and/or q in the expression for C_1), until no successful repair scheme can be found.

The expression for C_2 monotonically decreases as more and more defective cells are covered by the spare rows and columns. If all defective elements are covered by at least one spare row or column, the above expression is then minimized to zero. Unsuccessful repairing schemes which fail to cover all defective elements will yield positive costs. The two double summation terms in C_2 are equivalent. As it will be seen later in the expression for the synaptic weight matrix, this duplication facilitates the later transformation to neural net energy functions.

Another cost function can be defined by modifying the C_1 expression, i.e.,

$$C'_1 = A/2 \left(\sum_{i=1}^m s_{1i} \right)^2 + A/2 \left(\sum_{j=1}^n s_{2j} \right)^2.$$

Note that C'_1 is actually a simplified C_1 with p and q set to zero. Without considering whether all defects will be covered or not, the expression for C'_1 prefers repair schemes with fewer spares used. Intuitively, this way more efficient repair schemes may be encouraged. But the new $C'_1 + C_2$ combination has a potential danger of having a large amount of unsuccessful repair schemes with a small number of uncovered defects mapped to local minima just because they request less spares. The differences between the neural networks defined by these two cost functions will be further illustrated later in the simulation studies.

B. The Gradient Descent (GD) Approach

It may be recalled that the neural network energy function is $E^{NN} = -1/2 \sum_i \sum_j w_{ij} s_i s_j - \sum_i s_i b_i$. By rewriting C^{MR} in the form of E^{NN} , the weights of the interconnect-

ing synapses and the intensities of the external biases to neurons can be determined through simple algebraic manipulations.

But in order to keep the main diagonal of the neural net's synaptic weight matrix zero, all $A/2 \cdot \sum s_i^2$ terms in the array repair cost expression are rewritten as $A/2 \cdot \sum s_i$. Since $s_i \in \{0, 1\}$, these two terms are mathematically equivalent. However, $A/2 \cdot \sum s_i^2$ corresponds to self-feedback through a synapse of strength $-A$ for all neurons, and $A/2 \cdot \sum s_i$ means a bias of $-A/2$ for every neuron. The resulting neural network is shown as follows:

$$w_{1i,1j} = -A(1 - \delta_{ij}) \quad w_{2i,2j} = -A(1 - \delta_{ij})$$

$$w_{1i,2j} = -B \cdot d_{ij} \quad w_{2i,1j} = -B \cdot d_{ji}$$

$$b_{1i} = (p - 1/2) \cdot A + B \sum_j d_{ij}$$

$$b_{2j} = (q - 1/2) \cdot A + B \sum_i d_{ji}$$

where $\delta_{ij} = 0$ if $i \neq j$, otherwise 1.

In order to illustrate how to construct the synaptic weight matrix from these values, let us consider a pathological defect pattern that consists of 16 faulty cells as shown in Fig. 1. If there are only four spare rows and four spare columns, a greedy algorithm such as RM cannot produce a successful repair scheme [1]. The only successful allocation of spares for replacement is indicated in the figure by horizontal and vertical stripes. But due to the presence of four defects in both row 1 and column 9, the RM algorithm will be tempted to make the first two repairs on them. The algorithm will thereby fail to repair the memory because, in addition to the available six spare rows and columns, the algorithm will require two more columns or rows to cover all the defects along the diagonal direction. The neural network algorithm described in this section can repair successfully such a defect pattern by finding the unique solution shown in Fig. 1. For such a memory fault pattern, the neural net synaptic weights will be estimated by adjusting $A/B = 1$ (the value $A = 2$ has been chosen empirically and is not critical for the performance of the algorithm) in the above expressions. The resulting neural net synaptic weight matrix and the biases are given in Fig. 2.

C. The Hill-Climbing (HC) Approach

We propose here a neural net with simplified HC behavior that provides high-quality solutions similar to the powerful combinatorial technique, called *Simulated Annealing* [13]. First, let us note that the solutions to the memory array repair problem are classified as *success* and *failure*. In fact, it is commonly known that all optimization problems can be transformed to a series of yes/no questions [5]. Due to this specific criterion, we do not consider any random move that may increase the energy of the system until the move is absolutely necessary, i.e., the search has reached an unsuccessful local energy minimum. If the current local minimum does not yield a successful repair, moves that increase the system energy are

net will be stuck at the present state, making the repair scheme unsuccessful. On the other hand, if we have $A > B$, all the current row spares will cover only one faulty element exclusively will now receive a net negative influence equal to $|B - A|$, thus causing the network to switch to a new state and give up on the current scheme.

IV. SIMULATION STUDIES

In this section, simulation results are provided to demonstrate the superiority of the proposed neural computing approach for the problem of automatic array repair. Six hundred *reduced array fault patterns* of size 10×10 (Case 1) and 20×20 (Case 2) randomly generated arrays with about 10, 15, and 20% faulty elements were used in the experiments. As was explained in Section III, these small size arrays represent the actual defective rows and columns in large-size memory arrays, some as large as a few million bits.

The performance of two GD neural nets GD and GD' , defined by cost functions $C_1 + C_2$ and $C'_1 + C_2$, respectively, are compared. For both 10×10 and 20×20 arrays, the probabilities of finding a successful repair scheme versus the ratio B/A are depicted in Fig. 3. By controlling the value of B/A , the importance of the fault coverage over the spare usage can be manipulated.

According to the figure, the effectiveness of GD' is largely affected by how A and B are selected. When A is set equal to B , no successful repair scheme can be found by GD' for any one of the repairable fault patterns, just as was expected in Section III. For 10×10 arrays, the percentage of successful repairs improves to about 50% as B/A increases to 5, and thereafter converges to about 42%. GD' behaves similarly for 20×20 arrays, except that the peak performance happens at about $B/A = 12$. Now, let k be the maximum of p and q , where p and q are the maximum numbers of spare rows and columns allowed to form a successful repair scheme, respectively. The B/A ratio for a peak value to occur is found to be equal to k . Note that each firing row (column) neuron will send a $-A$ amount of influence to disable all other row (column) neurons, and a row or column neuron is encouraged to fire by as low as a $-B$ amount of influence to cover faults in a row or column. In order to allow the neural net to explore solutions using up to k spare rows and k spare columns, we must have $B \leq kA$ to keep up to k row or column neurons firing at the same time.

On the other hand, GD shows the advantage of its low sensitivity to how A and B are selected over the range of $B/A > 1$. In a hardware implementation of a neural network chip, the ratio of B/A is likely to vary over a wide range because of processing parameter variations and the wide tolerance of resistive elements. Thus the performance of the neural network will remain uniform from chip to chip if GD is implemented as opposed to GD' , whose behavior changes dramatically with different values of B/A . The percentage of successful repairs obtained by GD appears to have a peak value of about 50% at $A =$

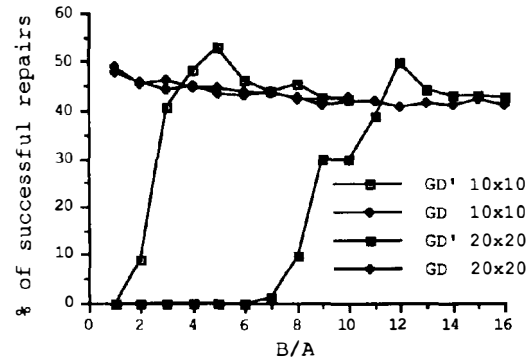


Fig. 3. Performance comparisons between GD and GD' .

B and declines asymptotically to about 42% as B/A increases. It is not surprising that the percentages of successful repairs obtained by GD' and GD are converging to the same value, since when B is much larger than A , the consideration for the spare usage is of little or no effect, compared with the consideration for fault coverage.

For the second experiment, the effectiveness of the RM algorithm is compared with GD . In order to provide an equal starting situation, the programmed neural net is started with all neuron processors in nonfiring states instead of other random combinations. We will use GD -zero to denote this special case of the GD convergence. The performances of these two algorithms are compared in Table I. Random defect patterns are generated for both cases, and the algorithms are applied to repair the faulty arrays with three different sets of available spares. From the results it was seen that on average, GD -zero is twice as successful as RM when the defect pattern is not very large (represented by Case 1) and few spare rows and columns are used. But GD -zero is about three to five times more successful when the defect pattern is large (represented by Case 2) and a relatively large number of spare rows and columns are used. As for the number of steps taken to find a repair scheme, it is about the same for both algorithms.

Second, tradeoffs between the use of two neural computing methods are examined, and the results are shown in Table II. For each defect pattern, a number of random trials are performed by each method. Average performance reaches consistency within one hundred random trials. As is expected, the simulation indicates that the HC approach is almost perfect in locating a successful repair scheme for repairable arrays due to its ability to perform global searches. The probability for the GD to reach a successful repair in a random trial is about 0.5. The average number of steps needed by HC is about two to three times that needed by GD . The runtime for the GD algorithm varies very little over a large number of experiments, but the HC algorithm is found to have a wide variance in the number of steps to execute the search successfully.

The chances of getting a successful scheme achieved by the four methods are shown in Fig. 4 as percentages.

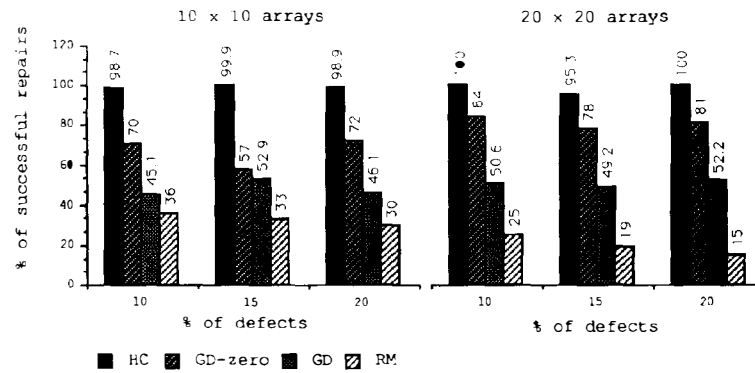


Fig. 4. Performance summary of the four repair methods.

TABLE I
AVERAGE % OF SUCCESSFUL REPAIRS DONE BY RM AND GD-ZERO

% of Defects	Case 1			Case 2		
	Spares	RM	GD-zero	Spares	RM	GD-zero
10	(3, 3)	36	70	(9, 9)	25	84
15	(3, 4)	33	57	(9, 12)	19	78
20	(4, 5)	30	72	(12, 12)	15	81

TABLE II
PERFORMANCE COMPARISONS BETWEEN GD AND HC NEURAL NETS

	% of Defects	Spares	% of Success		Avg # of Steps		σ	
			GD	HC	GD	HC	GD	HC
Case 1	10	(3, 3)	45.1	98.7	7.0	14.4	2.0	11.5
	15	(3, 4)	52.9	99.9	6.6	14.0	2.0	10.4
	20	(4, 5)	46.1	98.9	6.6	22.6	1.9	21.8
Case 2	10	(9, 9)	50.6	100.0	13.0	22.6	2.6	13.1
	15	(9, 12)	49.2	95.3	12.6	40.2	2.7	25.3
	20	(12, 12)	52.2	100.0	12.7	20.4	2.6	16.8

It can be seen that for 10×10 arrays HC achieves almost 100% success in all cases. The *GD-zero* performs better than *GD* where the neurons are randomly turned on at the initiation of the algorithm. Greedy techniques such as RM fail to cover more than 67% cases on the average. For 20×20 arrays, RM's performance deteriorates, and on average in more than 80% of the cases, it fails to repair the memory.

To summarize the simulation results, it may be pointed out that *GD* is fast in reaching a spare allocation, and has small variance in the number of search steps. On the other hand, the number of search steps used by the HC can vary widely, depending on the initial random starting states and how difficult the problem instances are. But on average, the number of search steps needed by HC does not escalate exponentially. For four out of six types of settings, the average number of search steps for HC is limited to about twice the number for *GD*. For the other two types of settings, the ratios of the average number of search steps are no worse than four. In fact, if we compare the equivalent average number of search steps expected for

TABLE III
GD AND HC IN AVERAGE NUMBER OF SEARCH STEPS

% of Defects	Case 1		Case 2	
	GD	HC	GD	HC
10	49.0	14.4	91.0	22.6
15	39.6	14.0	63.0	40.2
20	46.2	22.6	88.9	20.4

the *GD* approach to attain the same level of success as the HC approach does within 1%, the HC approach is found to be more efficient. The actual numbers can be seen in Table III.

V. ELECTRONIC IMPLEMENTATION

In this section, we demonstrate how to implement an electronic neural network that can repair a faulty memory by using the HC searching strategy. Electronic neural nets can be built by both analog and digital techniques. Murray and Smith have discussed their tradeoffs and devel-

oped a digital neural net representing the interactions by pulse stream of different rates [22]. Even though the digital network has better noise immunity, it requires large silicon area for its pulse generators and digital multipliers. Analog neural nets using current summing principles were developed by Graf *et al.* [7] and Sivilotti *et al.* [24]. Carver Mead's book on analog VLSI presents the design strategies for several types of analog neural systems [20]. These analog neural nets require small area and their intrinsic robustness and ability to compute correctly even in the presence of component failures, are particularly useful features for large-scale VLSI implementation.

In this design, a neuron is realized as a difference amplifier which produces binary output voltages. The firing neuron state is represented by the high voltage output, and the nonfiring state by the low voltage output. The synaptic influences between neurons, as well as the biases to neurons, are represented by electrical currents. Thus, an excitatory influence propagated to a neuron can be realized by injecting current to the corresponding amplifier input. Similarly, an inhibitory influence propagated to a neuron can be realized by sinking current from the corresponding amplifier input. As for the synaptic amplification, it is simulated by the branch resistance regulating the amount of current that passes through.

Note that for the array repair spare allocation problem discussed here, all the synaptic weights of the interconnection matrix are negative. Moreover, between two connected neurons, the synaptic weight of the connecting link can only be one of two different values, $-A$ or $-B$. If we assume that there are equal numbers of row and column neurons, an example electronic neural net with eight neurons can be implemented, as shown in Fig. 5.

According to Fig. 5, a neuron fires a negative influence to another neuron by turning on a transistor which sinks current from the target neuron's amplifier input. The amount of current involved is controlled by the size of the transistor. Note that owing to the assumption that equal numbers of neurons are reserved for row and column repairs, we are able to divide the interconnection matrix into four equal parts. It is not hard to find that only the first and third quadrants, as indicated in Fig. 5, need to be programmed based on the variable memory array defect pattern. For each programmable link, an additional transistor controlled by a memory cell is added to the pull-down transistor. A programmable link can be disconnected (connected) by storing a 0 (1) in the memory cell.

Fig. 6 shows the essential steps in programming the neural net's interconnection network. Let the memory array defect pattern be given in Fig. 6(a), with faulty cells represented by black dots. Next, by deleting fault-free rows and columns, a compressed defect pattern is obtained in Fig. 6(b), with a faulty cell denoted by a 1. Then the compressed defect pattern is used to program the first quadrant of the neural net interconnection matrix, and the third quadrant is programmed by the transposed compressed defect pattern, as is shown in Fig. 6(c).

Finally, the possibility of building the neural network

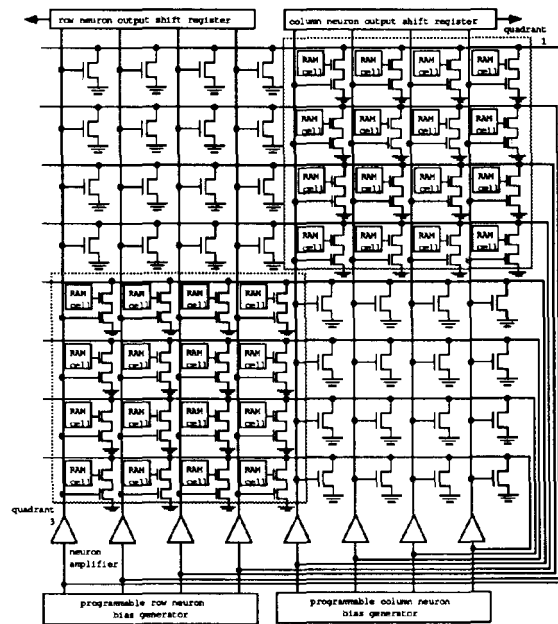


Fig. 5. An example electronic neural net for memory repair

alongside an embedded memory to achieve this self-repair purpose is demonstrated by the schematic diagram shown in Fig. 7.

The design to be discussed assumes that a built-in tester is available to provide fault diagnosis information. First of all, the status of each memory row and column is determined after the testing is done, and this information is stored in the faulty-row-indicator shift-register (FRISR) and faulty-column-indicator shift-register (FCISR), with a faulty row or column indicated by a 1. Then, to compress the memory array defect pattern, the detailed defect pattern of each faulty row is provided to the row-defect-pattern shift-register (RDPSR) one row at a time. As mentioned in Section III, the characteristic matrix D is highly sparse, and the fault pattern can be stored in the form of a compressed matrix each row and column of which will contain at least one faulty element. This is obtained by shifting those bits of RDPSR that correspond to the nonzero bits of FCISR, into the compressed-row-defect-pattern shift-register (CRDPSR). The content of CRDPSR is then used to program a row (column) in the first (third) quadrant of the neural net's interconnection network. The row (column) address of the quadrant is obtained by counting the nonzero bits in the FRISR. The stimulating biases to the row and column neuron inputs are generated by counting the total number of faults in each row (column) in the row (column) fault-count shift-registers. After a repair scheme is obtained by reading the outputs of the neuron amplifiers, the information is expanded in reverse order to the compression of defect patterns, and passed on to control logic of actual reconfiguration circuits. Normally, laser zapping [23], [3], focused ion-beam [21], or electron beam [6] techniques are used to restructure a

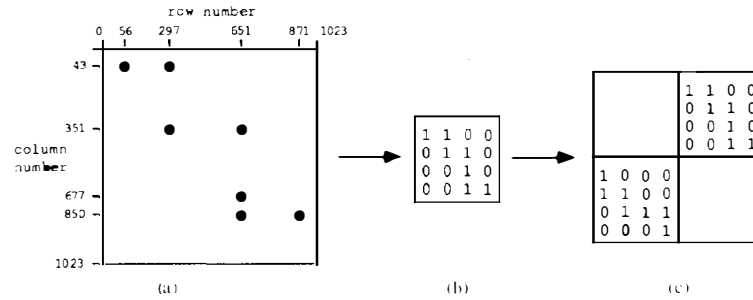


Fig. 6. Example programming of the neural net for memory repair. (a) Memory array defect pattern. (b) Compressed defect pattern. (c) Neural network RAM cell bitmap.

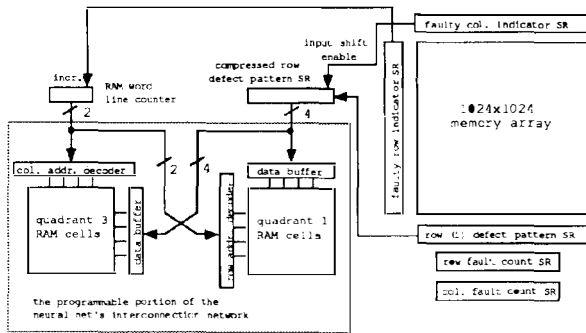


Fig. 7. Schematic of a neural net for embedded memory.

faulty memory array where the laser or charged beams are used for blowing out the programmable fuses to disconnect the faulty rows and columns. These schemes cannot be employed in automatic restructuring. Two viable techniques for self-restructuring are i) electronically programmable amorphous silicon fuses, which can be programmed by applying a 20-V pulse [27]; and ii) programmable electronic reconfiguration switches, which usually impose a small amount of penalties on circuit speed and area.

The circuit behavior was verified through SPICE simulations. Simulation output for a compressed 4×4 defect pattern with eight defective memory cells is shown in Fig. 8 as an example. The defects are represented by shaded squares. The initial state of the neural net was zero, i.e., there was no allocation of spares for any faulty rows or columns. Since every neuron represents either a faulty row or a faulty column in the compressed defect pattern, to cover the defects all neurons initially began transition toward the firing state. After 3 ns, the neurons representing row 2 and column 4 were successful in competition and remained in the firing state. This can be explained by the fact that row 2 and column 4 each have three defects to give the corresponding neuron the largest positive bias toward the firing state. All other neurons started to turn off, due to the mutual discouragement propagated through the negative synapses. When the remaining neurons were becoming low in activity, the mutual discouragement factor was also weakening, and the neurons started to move to-

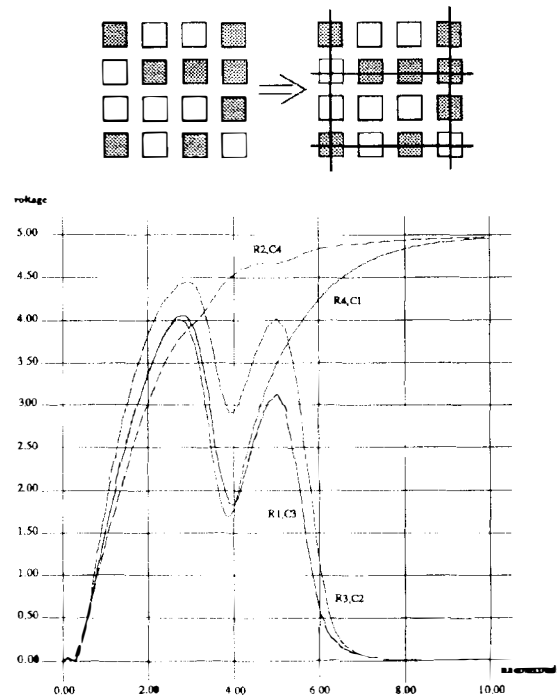


Fig. 8. Example memory repair SPICE simulation.

ward the firing state again. This time the neurons representing row 4 and column 1 were successful. The other four neurons then continued to the off state, since there was no remaining defect to cover and the spares all used up.

VI. NEURAL NET BEHAVIOR IN FAULT CONDITIONS

One unavoidable problem in adding extra hardware for BISR is that the extra hardware itself may be subject to component failures. In this section, we demonstrate the intrinsic fault-tolerant capability of neural networks as a reconfiguration control unit in the memory repair circuit.

Three types of component failures have been identified in neural networks, namely synapse-stuck faults, bias fluctuations, and neuron-stuck faults to serve as the fault model. For each faulty synapse, either of the synaptic

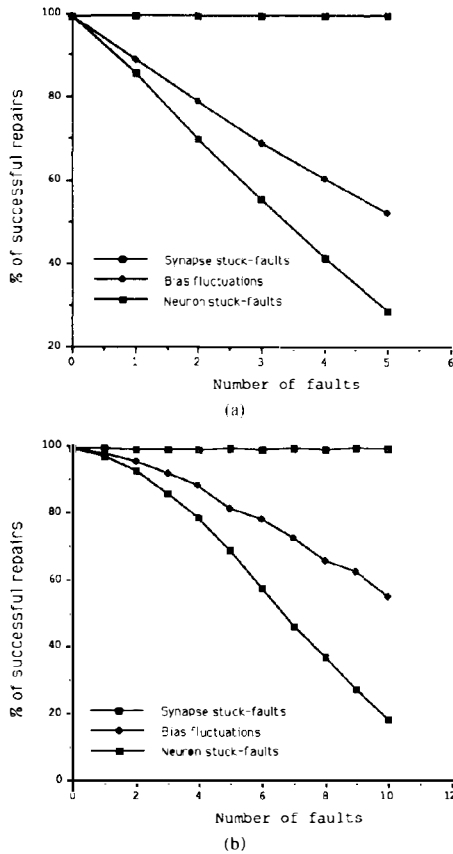


Fig. 9. Neural net performances in fault conditions

weights w_{ij} or w_{ji} can be assumed to be stuck-at- x , where x is a positive number in the range of synaptic strength values, due to transistor-stuck faults or defective memory cells that control the programmable synapses. Faulty bias generators are modeled to fluctuate within one unit of the predetermined biases, and faulty neurons will have stuck-at-firing or stuck-at-nonfiring states.

Rather than rendering the entire neural network useless, these faults will only reshape the energy (cost) distribution over the set of spare allocation schemes. Allocation schemes, which correspond to complete covers of all defective memory cells, may no longer be mapped to acceptable local energy minima, thus prolonging the search for an acceptable solution. On the other hand, incomplete allocation schemes may be mapped to false acceptable local energy minima, causing the neural net to stop searching without success.

Identical memory defect patterns prepared in Section IV are used here. Random faults are incrementally injected into a HC type of neural net to examine the achievable percentages of repairs for repairable defect patterns. To distinguish the degrees of seriousness among different types of faults, we inject only faults of the same type. The simulation results are shown in Fig. 9.

The results indicate that a small number of synapse-

stuck faults, up to five in the 20-neuron network and ten in the 40-neuron network, have almost no effect on the average performance. But the bias fluctuation faults and the neuron-stuck faults cause the average percentage of repairs out of repairable defect patterns to decrease steadily to zero when over one-fourth of total neurons are affected. The difference between a bias fluctuation fault and a neuron-stuck fault is in the amount of influence given to a particular faulty row or column for repair. While a bias fluctuation fault will encourage or discourage the corresponding row or column substitution slightly by one unit of influence, a neuron-stuck fault will actually insist the substitution be made.

VII. YIELD ANALYSIS

In this section, a quantitative analysis of yield enhancement due to neural network's self-repair capability is done. Faults are injected into memory arrays, spares, and neural networks to compute the resulting yield. The overhead of the self-repair logic is also estimated.

A well-known yield formula due to Stapper [25], [26] is used here to calculate the original yield, Y_0 , without the neural-net-controlled self-repair,

$$Y_0 = (1 + A\delta/\alpha)^{-\alpha}$$

where δ is the defect density, A is the memory array area, and α is some clustering factor of the defects. Let P_r be the probability function for a defect pattern to be repairable with respect to the remaining fault-free spares and the fault condition of the self-repair circuitry, and B be the area of overhead. Then, the yield, Y_N , with neural-net-controlled self-repair can be calculated as follows:

$$Y_N = Y'_0 + (1 - Y'_0) P_r,$$

where

$$Y'_0 = [1 + (A + B)\delta/\alpha]^{-\alpha}.$$

By minor algebraic manipulation, it can be seen that

$$Y'_0 = \left(\frac{Y_0}{1 + B\delta Y_0/\alpha} \right)^\alpha.$$

If $\alpha = 1$ then, the new yield formula can be simplified to

$$Y_N = P_r + \frac{(1 - P_r)Y_0}{1 + B\delta Y_0}.$$

We used 256×256 (64-Kbit) memory arrays for simulation with $\alpha = 1$ and $A\delta = 0.25, 0.5, 1, 2, 3, 4$, so that $Y_0 = 80\%, 66.7\%, 50\%, 33.3\%, 25\%$, and 20% , respectively. Six hundred memory arrays, each with at least one defect, are generated accordingly. For each set of memory arrays, up to four spare rows and four spare columns were provided to examine the neural net repairability with respect to the number of available spares. The spares and the neural net components are all subject to the same degree of defect density as the memory array. The resulting function values of P_r are shown in Fig. 10. Finally, the corresponding calculated yields are given in Fig. 11.

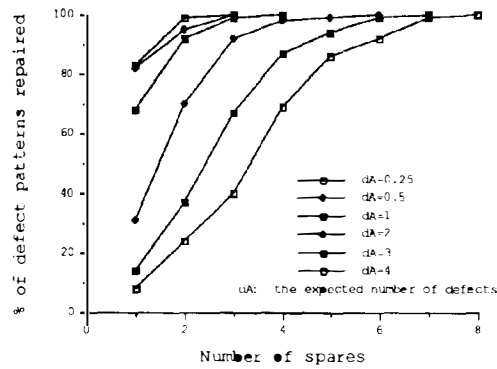
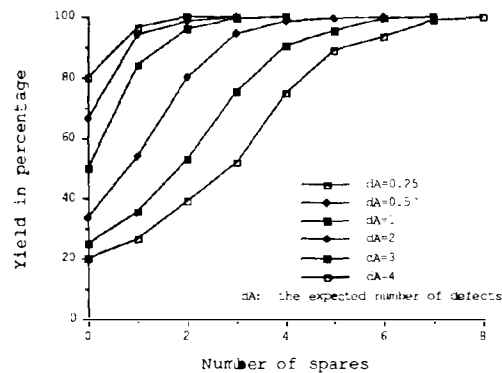
Fig. 10. Repairability of a defective 256×256 (64K) memory arrays.

Fig. 11. Chip yield w.r.t. fault density and available spares.

According to the schematic diagram shown in Fig. 7, except for the FC/RISR and the RDPSR, which are proportional to the dimension of the memory array, the complexity of the remaining self-repair logic basically depends on the dimensions of the compressed defect patterns. For instance, an $m \times n$ compressed defect pattern will require a neural net of $m + n$ neurons. Since the mean and variance of the number of faulty columns (rows) can be calculated according to the fault distribution function, an estimate on the required neural net size can be easily made to accommodate nearly all possible compressed defect patterns. For the case of 256×256 (64K) memory arrays with $\Delta\delta$ as high as 4, it is a near-certainty that the compressed array will not be larger than 10×10 . Given an $N \times N$ DRAM, let the size of the maximum compressed memory defect pattern be n^2 . An itemized account of the dynamic memory array and self-repair hardware based on transistor count is given in Table IV, and the percentages of overhead are listed in Table V for various memory and neural net sizes. As indicated by the results, the overhead is insignificant, compared with the yield improvement, and the overhead can be even smaller if static RAM's, with $6N^2$ transistors in the memory array, are considered. Here the overhead of the BIST circuit is not included. Typically this overhead is very low (about 28 flip-flops and ten gates for a 256-K RAM), as shown in Mazumder's earlier papers [16], [17].

TABLE IV
AN ITEMIZED ACCOUNT OF AN $N \times N$ RAM WITH BISR CIRCUITS

Component	Number of transistors
Memory cells	$4N^2$
Row decoder	$2N$
Column decoder	$2N$
Sense amp's	$6N$
Faulty column indicator SR	$8N$
Faulty row indicator SR	$8N$
Row defect pattern SR	$8N$
Row fault count SR	$8n \log n$
Column fault count SR	$8n \log n$
RAM word line counter	$8 \log n$
Compressed row defect pattern SR	$8n$
Programmable synapses SR's	$16n^2$
Non-programmable synapses	$2n^2$
Programmable synapses	$4n^2$
Neurons	$8n$
Bias generators	$2n(4 \log n + n)$

TABLE V
SELF-REPAIR HARDWARE OVERHEAD IN % OF RAM SIZE

	$n = 8$	$n = 12$	$n = 16$	$n = 20$
$N = 256$ (64 Kbit)	5.14	5.66	6.35	7.20
$N = 512$ (256 Kbit)	2.46	2.59	2.76	2.97
$N = 1024$ (1 Mbit)	1.20	1.23	1.28	1.33

VIII. FINAL REMARKS

The biological nervous system's ability to solve perceptual optimization problems is imitated here to tackle the VLSI array repair problem. In contrast to the current sequential repair algorithms, which run slowly on the conventional digital computers, the neural network's collective computational property provides very fast solutions. Methods to transform the problem instance into neural network computation model are demonstrated in detail. Of the two types of neural nets studied in this paper, the GD neural network has been found to be two to four times better than the RM algorithm in obtaining successful repair schemes. The GD minimizes the energy function of the network only in the locality of the starting energy value. The performance of the neural network is further improved by introducing an HC technique that allows the search to escape the traps of local minima. By generating random defect patterns and experimenting with a large number of arrays, it is seen that the HC algorithm finds a solution in a repairable memory array with near certainty (with a probability of 0.98 or more). For the same fault patterns, simple commercial algorithms, like RM, can yield feasible solutions for only 20% of the patterns. On the average, about twice as many search steps are used by the HC as opposed to GD.

Both the HC and GD neural networks can be implemented in hardware using very small overhead, typically less than 3% if the memory size is 256 Kbit or more. The payoff of this BISR approach is very high: the VLSI chip yield can increase from 10% without the BISR circuit to about 100% by using the proposed neural nets. The paper also proves that the neural networks are much more robust and fault-tolerant than the conventional logic circuits, and

thereby are ideally suited for self-repair circuits. The proposed designs of neural networks can operate correctly even in the presence of multiple faulty synaptic circuit elements, and as more numbers of neurons become permanently stuck to a firing or a nonfiring state, the networks gracefully degrade in their abilities to repair faulty arrays.

The paper shows how to solve a vertex-cover problem of graph theory, generally known to be an intractable problem, by using neural networks. A large number of problems in other domains, which can be modeled in similar graph-theoretic terms, can also be solved by the neural-network designs discussed in this paper. In the memory repair problem, an entire row or column is replaced to eliminate a defective cell. Such an approach is easy to implement and is cost-effective in redundant memory designs because the cells are very small in size. But a large class of array networks in systolic and arithmetic logic circuits employ a different strategy where a defective cell is exactly replaced by a spare or redundant cell. An appropriate graph model for such an array repair problem will be to construct a maximum matching of pairs between the vertices of a bipartite graph representing the set of faulty cells and the set of available spare cells. The technique described in this paper can be extended to solve the maximum matching algorithm by neural networks [18].

The overall goal of the proposed BISR circuits is to improve the chip yield by reconfiguring the faulty components at the time of chip fabrication, and also to extend the life span of the chip by automatically testing and repairing it whenever a failure is detected during the chip's normal operation. In space, avionics, and oceanic applications where manual maintenance, namely field testing and replacement, is not feasible, such an auto-repair technique will be very useful in improving the reliability and survivability of computer and communication equipment.

REFERENCES

- [1] J. R. Day, "A fault-driven comprehensive redundancy algorithm for repair of dynamic RAMs," *IEEE Design and Test*, vol. 2, pp. 35-44, June 1985.
- [2] R. C. Evans, "Testing repairable RAMs and mostly good memories," in *Proc. Int. Test Conf.*, Oct. 1981, pp. 49-55.
- [3] B. F. Fitzgerald and E. P. Thoma, "Circuit implementation of fusible redundant addresses of RAMs for productivity enhancement," *IBM J. Res. Develop.*, vol. 24, pp. 291-298, 1980.
- [4] W. K. Fuchs and M.-F. Chang, "Diagnosis and repair of large memories: a critical review and recent results," *Defect and Fault Tolerance in VLSI Systems*, I. Koren, Ed. New York: Plenum, 1989, pp. 213-225.
- [5] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. New York: Freeman, 1979.
- [6] P. Girard, F. M. Roche, and B. Pistoulet, "Electron beam effects on VLSI MOS: conditions for testing and reconfiguration," in *Wafer-Scale Integration*, G. Saucier and J. Trihle, Eds. New York: Elsevier, 1986, pp. 301-310.
- [7] H. P. Graf and P. deVegvar, "A CMOS implementation of a neural network model," in *Advanced Research in VLSI*, pp. 351-367, 1987.
- [8] H. P. Graf *et al.*, "VLSI implementation of a neural network memory with several hundreds of neurons," in *Proc. Conf. Neural Networks for Computing*, Amer. Inst. of Phys., 1986, pp. 182-187.
- [9] R. W. Haddad and A. T. Dahbura, "Increased throughput for the testing and repair of RAMs with redundancy," in *Proc. Int. Conf. Computer-Aided Design*, 1987, pp. 230-233.
- [10] N. Hasan and C.-L. Liu, "Minimum fault coverage in reconfigurable arrays," in *Proc. Int. Symp. Fault-Tolerant Comput.*, June 1988, pp. 348-353.
- [11] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities," in *Proc. Nat. Acad. Sci. USA*, vol. 79, April 1982, pp. 2554-2558.
- [12] J. J. Hopfield and D. W. Tank, "Neural computation of decisions in optimization problems," *Biol. Cybern.*, vol. 52, pp. 141-152, 1985.
- [13] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, "Optimization by simulated annealing," *Science*, vol. 220, pp. 671-680, 1983.
- [14] S. Y. Kuo and W. K. Fuchs, "Efficient spare allocation in reconfigurable arrays," *IEEE Design and Test*, vol. 4, pp. 24-31, 1987.
- [15] W.-K. Huang *et al.*, "Approaches for the repair of VLSI/WSI RAMs by row/column deletion," in *Proc. Int. Symp. Fault-Tolerant Comput.*, June 1988, pp. 342-347.
- [16] P. Mazumder, J. H. Patel, and W. K. Fuchs, "Methodologies for testing embedded content addressable memories," *IEEE Trans. Computer-Aided Design*, vol. 7, pp. 11-20, Jan. 1988.
- [17] P. Mazumder and J. H. Patel, "An efficient built-in self-testing algorithm for random-access memory," *IEEE Trans. Ind. Electron.*, vol. 36, May 1989, pp. 394-407.
- [18] P. Mazumder and Y. S. Jih, "Processor array self-reconfiguration by neural networks," in *Proc. IEEE Wafer-Scale Integration*, Jan. 1992.
- [19] W. S. McCulloch and W. H. Pitts, "A logical calculus of ideas immanent in nervous activity," *Bull. Math. Biol.*, vol. 5, pp. 115-133, 1943.
- [20] C. Mead, *Analog VLSI and Neural Systems*. Reading, MA: Addison-Wesley, 1989.
- [21] J. Melngailis, "Focused ion beam technology and applications," *J. Vac. Sci. Technol. B*, vol. 5, pp. 469-495, 1987.
- [22] A. F. Murry and A. V. W. Smith, "Asynchronous VLSI neural networks using pulse-stream arithmetic," *IEEE J. Solid-State Circuits*, vol. 23, pp. 688-697, 1988.
- [23] G. Nicholas, "Technical and economical aspect of laser repair of WSI memory," in *Wafer-Scale Integration*, G. Saucier and J. Trihle, Eds. New York: Elsevier, 1986, pp. 271-280.
- [24] M. A. Sivilotti, M. R. Emerling, and C. A. Mead, "VLSI architectures for implementation of neural networks," *Proc. Conf. Neural Networks for Computing*, Amer. Inst. of Phys., pp. 408-413, 1986.
- [25] C. H. Stapper, A. N. McLaren, and M. Dreckman, "Yield model for productivity optimization of VLSI memory chips with redundancy and partially good product," *IBM J. Res. Develop.*, vol. 24, no. 3, pp. 398-409, May 1980.
- [26] C. H. Stapper, "On yield, fault distributions, and clustering of particles," *IBM J. Res. Develop.*, vol. 30, no. 3, pp. 326-338, May 1986.
- [27] H. Stopper, "A wafer with electrically programmable interconnections," *Dig. IFFE Int. Solid-State Circuits Conf.*, 1985, pp. 268-269.
- [28] M. Tarr, D. Boudreau, and R. Murphy, "Defect analysis system speeds test and repair of redundant memories," *Electronics*, pp. 175-179, Jan. 1984.
- [29] C.-L. Wey *et al.*, "On the repair of redundant RAM's," *IEEE Trans. Computer-Aided Design*, vol. CAD-6, pp. 222-231, 1987.



Pinaki Mazumder (S'84-M'87) received the B.S.E.E. degree from the Indian Institute of Science in 1976, the M.Sc. degree in computer science from the University of Alberta, Canada in 1985, and the Ph.D. degree in electrical and computer engineering from the University of Illinois, Urbana-Champaign, in 1987. Presently he is an Assistant Professor at the Department of Electrical Engineering and Computer Science at the University of Michigan, Ann Arbor.

Prior to this, he was a research assistant at the Coordinated Science Laboratory, University of Illinois, and for over six years was with the Bharat Electronics Ltd., India, (a collaborator of RCA)

where he developed several types of analog and digital integrated circuits for consumer electronics products. During the summers of 1985 and 1986, he was Member of the Technical Staff in the Naperville branch of AT&T Bell Laboratories. His research interest includes VLSI testing, physical design automation, ultrafast digital circuit design, and neural networks. He is a recipient of Digital's Incentives for Excellence Award, National Science Foundation Research Initiation Award, and Bell Northern Research Laboratory Faculty Award. He has published over 50 papers in IEEE and ACM archival journals and reviewed international conference proceedings. He is a Guest Editor of the IEEE DESIGN AND TEST MAGAZINE's special issue on memory testing, to be published in 1993.

Dr. Mazumder is a member of Sigma Xi, Phi Kappa Phi and ACM SIGDA.



Yih-Shyr Jih (S'85-M'91) received the B.S. degree in computer science and information engineering from Nation Taiwan University, Taipei, Taiwan, Republic of China, in 1981, the M.S. degree in computer science from Michigan State University, East Lansing, in 1985, and the Ph.D. degree in computer science and engineering from the University of Michigan, Ann Arbor, in 1991.

He is currently a Research Staff Member in the IBM T. J. Watson Research Center, Yorktown Heights, NY. His research interests include VLSI computer-aided design, fault-tolerant computing, and high-bandwidth computer communication.