

A Memory Grouping Method for Sharing Memory BIST Logic

Masahide Miyazaki, Tomokazu Yoneda, and Hideo Fujiwara

Graduate School of Information Science, Nara Institute of Science and Technology (NAIST),

8916-5 Takayama, Ikoma, Nara 630-0101, Japan

Email: {masah-mi, Yoneda, fujiwara}@is.naist.jp

Abstract - With the increasing demand for SoCs to include rich functionality, SoCs are being designed with hundreds of small memories with different sizes and frequencies. If memory BIST logics were individually added to these various memories, the area overhead would be very high. To reduce the overhead, memory BIST logic must therefore be shared. This paper proposes a memory-grouping method for memory BIST logic sharing. A memory-grouping problem is formulated and an algorithm to solve the problem is proposed. Experimental results showed that the proposed method reduced the area of the memory BIST wrapper by up to 40.55%. The results also showed that the ability to select from two types of connection methods produced a greater reduction in area than using a single connection method.

I. Introduction

With the increasing number of functions being included in SoCs, many memories with different sizes and frequencies are being used. The latest SoCs contain hundreds of memories. Testing all the memories in these SoCs sequentially would take a long time. Therefore, a memory BIST design that allows two or more memories to be tested simultaneously is needed. However, due to power-consumption constraints, not all memories can be activated at the same time. To solve this problem, a scheduling technique for minimizing the test application time under power-consumption constraints is needed. Adding individual circuits for memory BISTs to lots of small memories would result in huge area overheads. To reduce these overheads, memory BIST logic must be able to be shared.

A BIST architecture, based on a single micro-programmable BIST processor and a set of memory wrappers, was proposed to simplify the testing of systems containing many distributed SRAMs of different sizes [1]. To reduce the BIST area overhead, it was proposed to share a single wrapper between a cluster of SRAMs (same type, width, and addressing space). However, in some cases, memories that have different widths or addressing spaces can be connected and share BIST logic. There can also be two or more connection methods. To achieve a satisfactory solution, the memory-connection type should be considered

along with decisions on memory groups.

In this paper, we propose two types of memory-connection methods for BIST wrapper sharing. A memory-grouping problem for test circuit minimization under constraints of power consumption and test application time is also formulated together with an algorithm that solves the problem. In addition, the effectiveness of this technique is demonstrated experimentally. This paper is organized as follows. In section II, our method for memory BIST logic sharing is described. In section III, the memory-grouping problem and an algorithm to solve the problem are presented. The experimental results are shown in section IV.

II. Memory BIST Logic sharing

In this section, we describe our method of BIST logic sharing for single port and word access memory. Figure 1 shows an example of a memory BIST wrapper. The data generator generates input test sequences. The address generator generates read and write addresses and the response analyzer captures test output responses and detects faults. The by-pass FFs are not used to test memory, but are used to care the memory interface signal during a scan test. The area of the address generator, data generator, and response analyzer are almost proportional to the bit width of the address, input data, and output data, respectively. However, some of these logics can be shared by different memories wherever the number of words or the data bit width are the same; hence, the area of test circuits can be reduced. In this paper, we treat the following two memory connection methods for memory BIST logic sharing: parallel connection and serial connection. Parallel connection can be used to connect memories that have the same number of words. Figure 2 shows an example of parallel connection.

In this example, three data and address generators are reduced to one by distributing the same test data and address signals from a couple of data and address generators to (1) - (4), enabling four memories to be tested simultaneously.

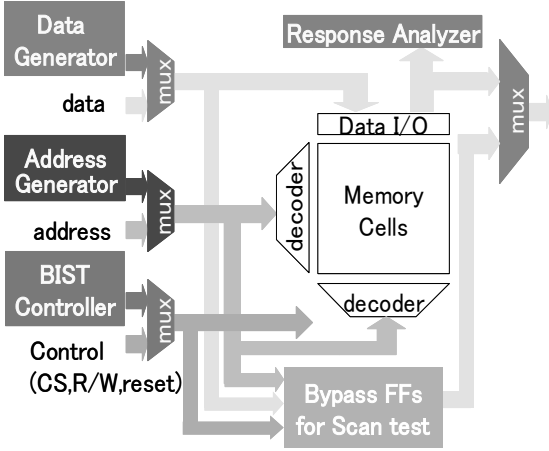


Fig.1 Memory BIST Wrapper

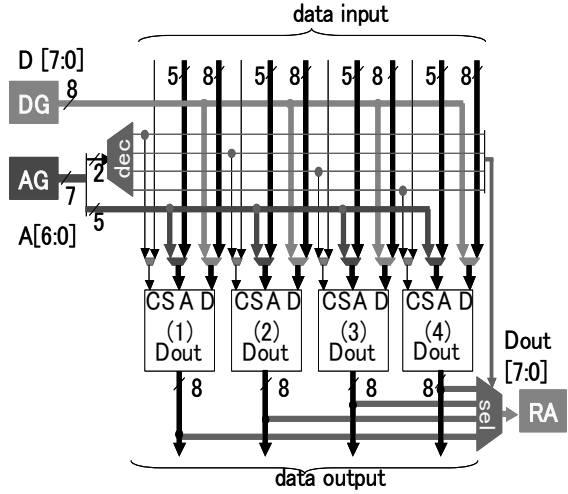


Fig.3 Serial connection of memories

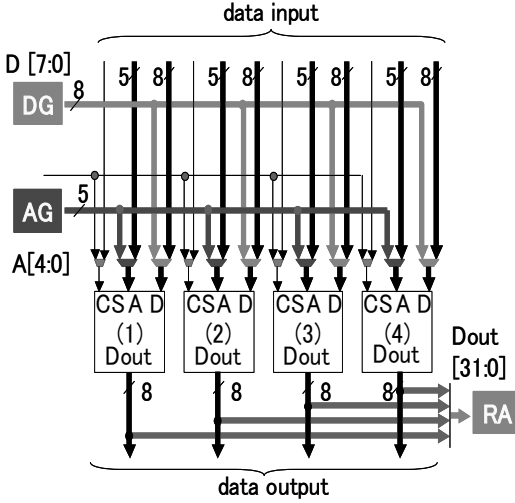


Fig.2 Parallel connection of memories

Serial connection allows memories with the same bit width to be connected. Figure 3 shows an example of four serially connected 8x32 word memories. In this example, the four memories are tested as an 8x128 word memory. The address generator generates an additional 2bit signal, and the signal is used to select the memories from (1) - (4), enabling the four memories to be tested serially. If all the memories have individual BIST logic, a 32-bit data generator and response analyzer are required, but in this example, all the memories can be tested using a shared 8bit generator and 8bit response analyzer.

Serial connection reduces the area more than parallel connection and also uses less power than parallel connection. However, the time required for serial connection testing is longer than that for parallel connection testing. To achieve the minimum area and a reasonable test application time under power consumption constraints, the type of memory connection should be considered during decisions on memory grouping. The layout design must also take into account distance constraints in relation to these connections.

III. Memory-Grouping Problem and Algorithm

A. Formulation of Memory-Grouping Problem

In this subsection, we present a memory-grouping problem. We assume that the following information for each memory m_i is given:

- b_i : data bit width of m_i
- w_i : word depth of m_i
- p_i : maximum power consumption of testing m_i
- f_i : operating frequency of m_i
- x_i : X coordinate of m_i , y_i : Y coordinate of m_i

We define two types of compatibility, namely p-compatibility and s-compatibility, as follows: Given a set of memories $V = \{m_1, m_2, \dots, m_n\}$, a pair of memories $m_i, m_j \in V$ is **p-compatible** if they satisfy the following conditions:

$$w_i = w_j \quad (1)$$

$$f_i = f_j \quad (2)$$

$$\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2} < D \quad (3)$$

D is a constraint value that the designer decides according to the design condition.

P-compatibility is represented by a graph $G_p = (V, E_p)$, where V is a set of a memory and the edge between a pair of vertices $(m_i, m_j) \in E_p$ exists if m_i and m_j are **p-compatible**. If a set of memories can be connected in parallel, the graph induced on G_p by the memories has to be a clique.

In the same way, a pair of memories $m_i, m_j \in V$ is **s-compatible** if they satisfy the following conditions:

$$b_i = b_j \quad (4)$$

$$f_i = f_j \quad (5)$$

$$\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2} < D \quad (6)$$

S-compatibility is represented by a graph $G_s = (V, E_s)$, where V is a set of memories and the edge between a pair of vertices $(m_i, m_j) \in E_s$ exists if m_i and m_j are *s-compatible*. If a set of memories can be connected serially, the graph induced on G_s by the memories has to be a clique.

To design memory BIST wrappers using these techniques for memory BIST logic sharing, we have to find a partition of V such that the memories that share the wrapper are included in the same block. Moreover, the partition $\pi = \{B_1, B_2, \dots, B_k\}$ has to satisfy the following conditions:

G_{ip} is the graph induced on G_p by block B_i .

G_{is} is the graph induced on G_s by block B_i .

G_{ip} or G_{is} is a clique.

When only the graph G_{ip} (G_{is}) is a clique, the memories included in B_i are connected in parallel (serially). If G_{ip} and G_{is} are both clique, we have to select the type of connection.

For a partition π , we can calculate the area of the BIST wrapper, test application time, and power consumption of each block. The area and test application time depend on the test-pattern algorithm. In this work, these were calculated according to a published design [4] using an 8N algorithm as follows.

If the connection type of block $B_i = \{m_1, m_2, \dots, m_k\}$ is a parallel connection,

$$\begin{aligned} \text{Area } S_{B_i} = & 0.75(\log_2(w_{B_i}))^2 + 2k \log_2(w_{B_i}) + 18 \sum_{l=1}^k b_l \\ & + 25 \log_2(w_{B_i}) + 3 \max(b_i) + 66 \end{aligned} \quad (7)$$

$$\text{Power consumption } P_{B_i} = \sum_{l=1}^k P_l \quad (8)$$

$$\text{Test application time } T_{B_i} = 8 \times w_{B_i} / f_{B_i} \quad (9)$$

$$(f_{B_i} = f_1 = f_2 = \dots = f_k)$$

If the connection type of block $B_j = \{m_1, m_2, \dots, m_k\}$ is a serial connection,

$$\begin{aligned} \text{Area } S_{B_j} = & 0.75 \left(\log_2 \left(\sum_{l=1}^k w_l \right) \right)^2 + 2k \log_2 \left(\sum_{l=1}^k w_l \right) \\ & + 25 \log_2 \left(\sum_{l=1}^k w_l \right) + k(\log_2 k) + 9b_{B_j}k + 14b_{B_j} + 8k + 61 \end{aligned} \quad (10)$$

$$(b_{B_j} = b_1 = b_2 = \dots = b_k)$$

$$\text{Power consumption } P_{B_j} = \max_l(p_l) \quad (11)$$

$$\text{Test application time } T_{B_j} = 8 \times \left(\sum_{l=1}^k w_l \right) / f_{B_j}$$

$$\times (\text{number of background patterns}) \quad (12)$$

The expressions for area calculation (7) and (10) do not consider the influence of timing conditions, but feedback is available from previous designs.

Parallel-connected memories are tested concurrently, and the power consumption is the sum of the power consumption of each memory. In contrast, serial-connected memories are activated one by one. Therefore, the power consumption is the maximum power consumption of the connected memories.

When a partition π is found, the area, power consumption and test application time of each block are calculated using the above expression.

The total area of the memory BIST wrappers S_{total} is calculated as the sum of S_{B_i} .

$$S_{total} = \sum_{i=1}^k S_{B_i} \quad (13)$$

To control each memory BIST wrapper, at least one BIST controller must be used. In this study, the number of memory BIST wrappers was reduced by using the proposed connections. There was therefore no increase in the number of controllers. In addition, our target design includes a lot of memories so that the area of the memory BIST wrappers is predominant. Therefore the area of the BIST controllers is disregarded.

To calculate the total test application time of a memory BIST under a power-consumption constraint, we used a rectangle packing algorithm that has been described elsewhere [5]. The algorithm optimizes the test schedule of each core so that the total test application time of an SoC is minimized under maximum power constraints. The inputs of the scheduling algorithm are the maximum allowed power consumption, the test application time, and the power consumption of each core. In this study, we considered a block to be a core. Therefore, we input $\{P_{B_j}\}$ $\{T_{B_j}\}$ as the information for each core. In addition, we assumed the bit width of the inter-connect between each wrapper and control logic remained unchanged. We therefore disregarded the maximum TAM width.

To reduce the total area of memory BIST wrappers by memory BIST logic sharing, we formulated the following memory-grouping problem.

Inputs:

a) A set of memories S and

Information for each memory:

$$M = M_i(b_i, w_i, p_i, f_i, x_i, y_i)$$

where, b_i , w_i , p_i , f_i , x_i and y_i are as follows:

b_i : data bit width of m_i

w_i : word depth of m_i

p_i : maximum power consumption of testing m_i

f_i : frequency of m_i

x_i : X coordinate of m_i

y_i : Y coordinate of m_i

Outputs:

- A partition π of a given set of memories S for which all the blocks satisfy the following conditions:
 - G_{ip} is the graph induced on G_p by block B_i .
 - G_{is} is the graph induced on G_s by block B_i .
 - G_{ip} or G_{is} is a clique.
- Type of connection of each block
- Test schedule of each memory

Constraints:

- Maximum distance of memory connection: D
- Maximum available peak power of the SoC: P
- Maximum test application time of memory: T

Objective:

To minimize S_{total} .

To solve this problem, an algorithm is proposed below.

B. Memory Grouping Algorithm

Fig.4 shows the pseudo code of the Memory Grouping Algorithm. Our proposed algorithm repeats division from 0-partition that only one block includes all memory to obtaining a target partition. As the algorithm divides the block, S_{total} increases. The min-cut method [2][3] is used to leave the possibility of the area reduction as much as possible. Moreover, it uses the following strategies to decide the compatibility of each block of the partition. Serial connection can reduce the area than parallel connection, and the power consumption is smaller than that of parallel connection. Therefore, it is possible that giving priority to serial connection reduces S_{total} . Based on this prospect, proposed algorithm searches for the partition that minimizes S_{total} only using s-compatibility in the first search.

First, the algorithm initializes variables. The minimum value of S_{total} is stored into S_{min} , and, in the first step, S_{min} is set to the total area of memory BIST wrapper without sharing. The partition of a set of memory S is stored into π , and the initial partition is set to 0-partition of S. (line 1-2).

Next, the algorithm creates two compatibility graphs (line 3), and select s-compatibility graph as the graph G that is used to find partition (line 4).

In order to check the compatibility of each block, the algorithm construct a set of graph C_{all} (line 6). Each graph G_i that is the member of C_{all} is induced on G by block B_i that is the member of π .

Then, for all B_i that include two or more memories, execute the following operations (line 7-21).

The minimum cut edge is calculated and delete them from G_i . By this operation, the vertex set B_i is divided into two blocks, leaving much possibility of the area reduction. If all the graph of new graph set C_{all} are clique, calculate S_{total} and test schedule of the new partition π_{imp} . If $S_{min} > S_{total}$ and the test scheduling succeeded, π_{imp} is stored into π_{best} as the best partition, and S_{total} is stored into S_{min} (line 8-17). If there is a graph G_i that is not a clique, or the test scheduling failed, π_{imp} is stored into π_{next} (line 18-20).

If there is no partition that should be tried, the first search is end (line22-24). Then the algorithm stores p-compatibility graph into G , and collects the blocks that have only one memory into one block (line25-29). Then, the algorithm searches for the partition that S_{total} is minimized using p-compatibility (line5-24). In the second search, it doesn't touch the blocks in which two or more memories are included after first search. Their connection type is fixed to serial connection. The connection type of the rest is determined to be parallel connection.

This algorithm performs $n(n-1)$ times division and scheduling in the worst case. The complexity of the scheduling algorithm and min-cut algorithm are $O(V \log V)$ and $O(V^2 \log V)$, respectively. Therefore the complexity of this algorithm is $O(V^3 \log V)$.

In this paper, we described our method for a single port and the word access memory. However, this method is applicable to other memories if the compatibility is defined about the memory type and the connecting method, and the area, power consumption, test application time can be shown by expression.

IV. Experimental Results

We carried out experiments to evaluate the proposed method. The proposed algorithm was implemented in C and the experiments were conducted on a 600-MHz Windows PC. Table 1 shows the information in each memory used in the experiment. The 2-4th columns denote the data bit width, word depth, and operating frequencies, respectively. The 5th column shows the power consumption. In this experiment, the power consumption of each memory was a relative value in which memory No. 1 was assumed to be 100 under the following assumption:

- The area is proportional to (number of words \times number of bits).
- The power consumption is proportional to the area.
- The power consumption is proportional to the frequency.

Table1. Information on Memories

No.	# data bit width	#Words	Frequency (MHz)	Power *1	Location	
					X	Y
1	16	128	133	100	10	10, 20, 30, 40, 50
2	16	128	133	100	20	
3	16	128	266	200	30	
4	16	128	266	200	40	
5	16	256	133	200	50	
6	16	256	133	200	60	
7	16	256	133	200	70	
8	16	256	133	200	80	
9	32	512	133	400	90	
10	32	512	133	400	100	

*1 Relative values in which memory No.1 is assumed to be 100

Procedure Memory Grouping (M, P, T, D)

```

1   $S_{min}$  = the total area of memory BIST wrapper without sharing;  $maxedgenum=0$ ;  $edgenum=0$ ;
2   $\pi = \{B_j\}$ ,  $B = \{m_1, m_2, \dots, m_n\}$ ;  $\pi_{tmp} = \phi$ ;  $\pi_{next} = \phi$ ;  $\pi_{best} = \phi$ ;  $\pi_{s-compatible} = \phi$ ;
3   $G_s = s\text{-compatibility\_graph}$  of  $B$ ;  $G_p = p\text{-compatibility\_graph}$  of  $B$ ;
4   $G = G_s$ ;
5  loop:
6  Construct a set of graph  $C_{all} = \{G_i \mid G_i \text{ is induced graph on } G \text{ by } B_i \in \pi\}$ 
7  for(  $\{B_i \in (\pi - \pi_{s-compatible}) \mid \text{which includes two or more memories}\}$  ) {
8      delete min-cut edge from  $G_i$ , make a set of graph  $C_{min} = \{G_{i1}, G_{i2}\}$ ;
9       $C_{all} = (C_{all} - G_i) \cup C_{min}$ ;
10     Set  $edgenum = \sum_j$  (the number of edges of  $G_j \in C_{all}$ );
11     Set a partition  $\pi_{tmp} = \{B_j \mid \text{vertex set of } G_j \in C_{all}\}$ ; if all  $G_j$  are clique, calculate  $S_{total}$  of  $\pi_{tmp}$ 
12     if(  $(\forall G_j \in C_{all}, G_j \text{ is clique}) \wedge (S_{min} > S_{total} \text{ of } \pi_{tmp})$  ) {
13         calculate  $T_{total} = \text{Schedule}(P, \{P_{B_j}\}, \{T_{B_j}\})$ ;
14         if(  $(\text{Schedule succeeded}) \wedge (T_{total} \leq T)$  ) {
15              $S_{min} = S_{total}$ ;  $\pi_{best} = \pi_{tmp}$ ;
16         }
17     }
18     if(  $edgenum > maxedgenum \wedge ((\text{Schedule failed, or } T_{total} \leq T) \vee (\exists G_j \in C_{all}, G_j \text{ is not a clique}))$  ) {
19          $\pi_{next} = \pi_{tmp}$ ;  $maxedgenum = edgenum$ ;
20     }
21 }
22 if(  $\pi_{next} \neq \phi$  ) {
23      $\pi = \pi_{next}$ ;  $\pi_{next} = \phi$ ; go to loop;
24 }
25 else if(  $G = G_s$  ) {
26      $G = G_p$ ;
27      $\pi_{s-compatible} = \{B_j \in \pi_{best} \mid \text{which includes two or more memories}\}$ ;
28      $B_s = \bigcup_k B_k$  ( $B_k \in (\pi_{best} - \pi_{s-compatible})$ );
29      $\pi = \{B_s\} \cup \pi_{s-compatible}$ ; go to loop;
30 } else { end; }

```

Fig.4 Memory Grouping Algorithm

The 6th and 7th columns show location. In this experiment, the number of memories was varied between 3 and 50, and the program was executed respectively. When the number of memories was $N < 11$, we used No. 1 to N, and for the rest, we extended the same set of No. 1-10, with the Y coordinate changing between 20 to 50.

In an actual test, several background patterns (e.g. marching, checker, checker-bar) are used, but in this experiment, the test application time was calculated by assuming the number of background patterns=1. In addition, the following constraint values were used:

Maximum distance of memory connection: $D=40$

Maximum available peak power of the SoC:

$P=5000$

Maximum test application time of memory:

$T=300 \mu s$

Experiments were carried out for the following five cases: (1) Not shared (all the memories had individual BIST wrappers); (2) parallel connection (memory BIST logic was shared using only parallel connection as described in the proposed technique); (3) serial connection (memory BIST logic was shared using only serial connection as described in the proposed technique); (4) parallel and serial connection (memory BIST logic was shared using both parallel and serial connection as described in the proposed technique); and (5) exhaustive search (memory BIST logic was shared using only parallel connection after an exhaustive search). Table 2 shows the experimental results. The first column

shows the number of memories and the second column shows the total area of memory BIST wrappers without sharing. Columns 3-5 shows the total area of memory BIST wrappers using the proposed techniques. The third column shows the results of using only parallel connection, while the fourth column shows the results of using only serial connection. The fifth column shows the results of using both parallel and serial connection and the sixth column shows the minimum solution obtained using an exhaustive search.

We were only able to complete an exhaustive search when the number of memories was less than 7. In these cases, the results of the exhaustive search showed that the memory BIST logic sharing technique reduced the area of the BIST wrappers by between 21.59 and 47.83% as minimum solutions. However, the technique achieved only 64.45% of the minimum solution in these cases, so there is room for improving the quality of the solution.

The average reduction ratio for parallel connection, serial connection, and parallel and serial connection were 21.08%, 37.25%, and 40.55%, respectively. In all cases, parallel and serial connection achieved the best solution. This result demonstrates that selection from two types of connection methods reduces the area more than using a single connection method.

Finally, Figure 5 shows the execution time of the implemented memory-grouping program. In all cases, the program was executed within 10 seconds using the proposed algorithm. The technique thus obtained good results within a very short CPU time so it is suitable for practical application.

Table2. Area of Memory-BIST Logic

#mem	not shared	Proposed algorithm			exhaustive
		P only	S only	S&P	
3	2289	1967	1660	1660	1660
4	2913	2591	2284	2284	2284
5	3537	2893	2279	2279	2279
6	4203	3559	3044	2722	2415
7	4869	3863	3690	3368	2540
8	5535	4529	3719	3397	N/A
9	6201	4793	3828	3506	
10	7242	5427	3122	3122	
11	8283	6021	6447	5678	
12	8907	7539	4703	4703	
13	9531	7394	5411	5089	
14	10155	7696	5406	5406	
15	10779	7998	5401	5401	
20	14484	10854	6769	6769	
30	21726	16281	10455	10455	
40	28968	22070	23784	19662	
50	36210	27497	22964	21551	

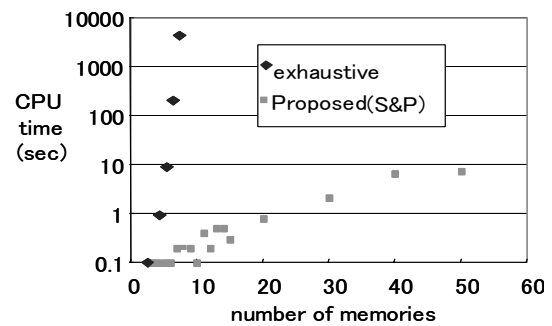


Fig.5 CPU Time of Memory Grouping program

V. Summary and Conclusions

A memory grouping problem was formulated and an algorithm to solve the problem was proposed. Experimental results showed that the proposed method reduced the area of memory BIST wrappers by up to 40.55%. It was also shown that the ability to select from two types of connection methods reduced the area more than using a single connection method.

In future work we will investigate improving the quality of the solution and minimizing the test application time.

Acknowledgements

Authors would like to thank Prof. Michiko Inoue and Prof. Satoshi Ohtake and members of Computer Design and Test Lab. (Nara Institute of Science and Technology) for their valuable discussions.

References

- [1] A. Benso, S. Di Carlo, G. Di Natale and P. Prinetto, "A Programmable BIST Architecture for Clusters of Multiple-Port SRAMs," in *Proc. International Test Conf.*, pp. 557-566, October 2000.
- [2] H. Nagamochi and T. Ibaraki, "A linear-time algorithm for finding a sparse k-connected spanning subgraph of a k-connected graph," *Algorithmica*, vol. 7, 1992, pp. 583--596.
- [3] H. Nagamochi and T. Ibaraki, "Computing the edge-connectivity of multigraphs and capacitated graphs," *SIAM J. Discrete Mathematics*, vol. 5, 1992, pp. 54--66.
- [4] Charles E. Stroud, *A Designer's Guide to Built-In Self-Test*, Kluwer Academic Publishers, The Netherlands, 2002.
- [5] V. Iyengar, K. Chakrabarty and E. J. Marinissen, "On using rectangle packaging for SOC wrapper/TAM co-optimization," in *Proc. VLSI Test Symposium*, pp. 253-258, May 2002.
- [6] Y. Huang, N. Mukherjee, S. Reddy, C. Tsai, W. Cheng, O. Samman, P. Reuter, and Y. Zaidan, "Optimal Core Wrapper Width Selection and SOC Test Scheduling Based On 3-Dimensional Bin Packing Algorithm," in *Proc. International Test Conf.*, pp. 74-82, October 2002.