

Process-Variation and Temperature Aware SoC Test Scheduling Using Particle Swarm Optimization

Nima Aghaee, Zebo Peng, and Petru Eles
Embedded Systems Laboratory (ESLAB), Linköping University, Sweden
{nima.ghaee, zebo.peng, petru.eles}@liu.se

Abstract—High working temperature and process variation are undesirable effects for modern systems-on-chip. It is well recognized that the high temperature should be taken care of during the test process. Since large process variations induce rapid and large temperature deviations, traditional static test schedules are suboptimal in terms of speed and/or thermal-safety. A solution to this problem is to use an adaptive test schedule which addresses the temperature deviations by reacting to them. We propose an adaptive method that consists of a computationally intense offline-phase and a very simple online-phase. In the offline-phase, a near optimal schedule tree is constructed and in the online-phase, based on the temperature sensor readings, an appropriate path in the schedule tree is traversed. In this paper, particle swarm optimization is introduced into the offline-phase and the implications are studied. Experimental results demonstrate the advantage of the proposed method.

Keywords: SoC test scheduling; particle swarm optimization; adaptive test; temperature aware; process variation

I. INTRODUCTION

Two main challenges for deep submicron integration are high power density and process variation [1]. The power density for a System-on-Chip (SoC) during test compared to its normal operation is high enough to put testing in trouble by considerably raising the overheating risk [2]. Efficient temperature-aware test scheduling techniques have been developed in order to minimize the test application time and avoid overheating [3, 4]. These methods neglect, however, the thermal consequences of the process variation and focus only on minimization of the test application time while maintaining the chips temperature under a given limit [3, 4].

The negative thermal consequence of process variation is unpredictability of the thermal behavior of the chip. It means that identical test vectors will result in a variety of different temperatures for different chips and cores. The difference between the expected temperature (estimated by simulation) and the actual temperature (measured by sensors) is called temperature error, which captures all errors generated due to different power/temperature-related effects. These effects include ambient temperature fluctuations, voltage variations, and process variation. For traditional technologies, temperature error is small enough to be negligible or to allow worst-case design with negligible performance penalty [3, 4].

The general trend of increase in power density and process variation will eventually lead to a situation where temperature

errors cannot be ignored any longer. Therefore, the thermal consequences of the process variation should be taken into account in order to develop efficient test. In [5], two process variation aware methods are proposed in order to maximize the test throughput by considering the thermal-safety as a part of the test cost. However, one of the proposed methods in [5] does not react to temperature deviations, and the other does not handle intra-chip process variation and rapid temperature error changes. In this paper an adaptive test scheduling method is introduced which navigates the tests according to the intra-chip process variation and temporal deviations in temperature errors. It makes use of multiple on-chip temperature sensors to provide on-line intra-chip temperature information.

A dynamic thermal-aware test scheduling technique using on-chip temperature sensors is proposed in [6] in order to cope with the power/temperature simulation inaccuracies in static scheduling. Thermal simulations are performed during the test in order to enable the earliest thermal-safe start of the next test [6]. This method does not handle the process variation and besides, it requires excessive Automated Test Equipment (ATE) resources to run the thermal simulation during test.

In [7], a dynamic test scheduling technique based on a tree structure is proposed in order to address thermal issues in presence of large process variation. The proposed method in [7] uses on-chip temperature sensors selectively and tries to maintain a low test cost while considering factors like sensor access overheads, temperature dependent leakage power, and ATE memory limit. The schedule with tree structure is very well suited for this problem, since it requires minimal ATE processing power.

This paper presents a technique to improve the quality of the test schedules by a new optimization strategy. The proposed method in this paper generates a near optimal schedule tree at design time (offline-phase). During test (online-phase), each chip traverses the schedule tree depending on the actual temperatures. The schedule indicates when a core is being tested and when it is put into the cooling state. The order of the test sequences is untouched and the schedule tree occupies a relatively small memory. Traversing the schedule tree requires a very small delay overhead for jumping from one point in the schedule tables to another point. This way, the complexity is moved into the offline-phase, leaving a simple online-phase. In this paper, *particle swarm optimization* is acquired for *sub-tree scheduling* and the implications are studied. *Sub-tree scheduling* is the computational bottleneck of

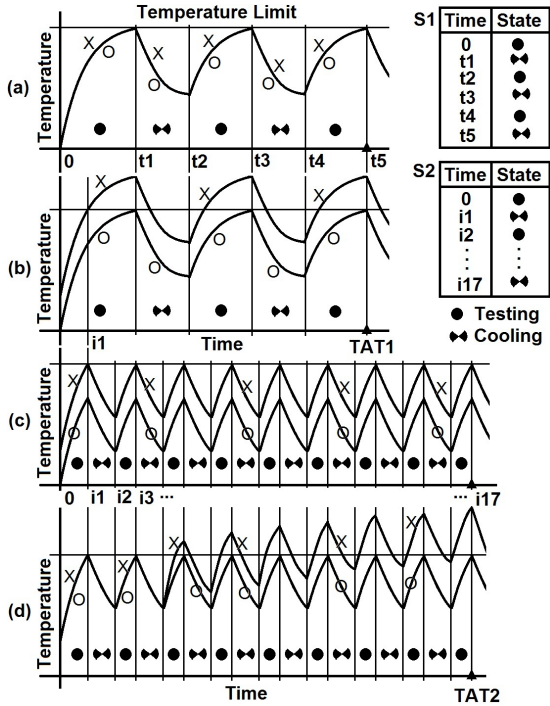


Figure 1. Test schedule examples (curves are only illustrative). (a) Temperature error is zero (no process variation). (b) Temperature error is constant and is not taken care of. (c) Temperature error is constant and is taken care of by a static schedule. (d) Temperature error is not constant and therefore can not be taken care of by a static schedule.

the algorithm since it requires thermal simulation for each candidate solution. Therefore, a better optimization heuristic which requires a lower number of iterations, to find a solution with a certain quality, will be very helpful. In the context of the optimization algorithms, this means that a better solution could be found with the same number of iterations and this is extremely appreciated.

The rest of the paper is organized as follows. A motivational example is given in section II. Section III provides the problem formulation and then introduces the cost function. Section IV describes the temperature error model. The linear schedule tables are discussed in section V. The tree scheduling is presented in section VI. The proposed heuristic, particle swarm optimization, is described in section VII. Experimental results are presented in section VIII. Section IX presents the conclusions.

II. MOTIVATIONAL EXAMPLE

Assume that there are two instances, o and x , from a set of chips manufactured for a given design. When the temperature error is negligible, the temperatures of o and x are equal and the same offline test schedule S1 is safely used for both of them, Fig. 1(a). Cooling periods for S1 are determined using thermal simulation. The simplest model of process variation only models the time invariant temperature errors. Assume that chip x is warmer than expected while chip o is normal; the result is overheating of chip x as shown in Fig. 1(b). In order to prevent overheating, a more conservative offline schedule S2 has to be designed considering x for both chips as illustrated in Fig. 1(c). S2 will lead to a longer test application time (TAT2 vs. TAT1). For chip o , S2 is unnecessarily long, since S1 was a safe

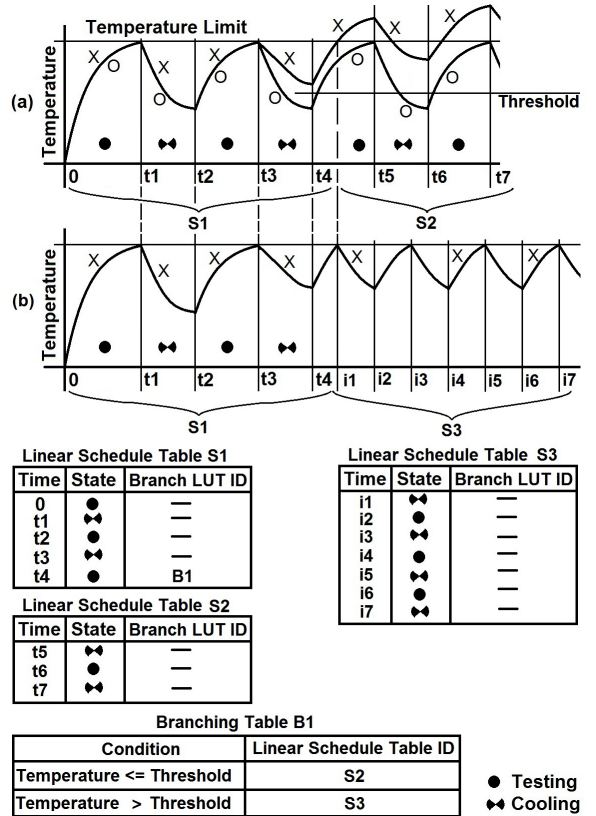


Figure 2. Schedule and branching tables (curves are only illustrative). (a) Chip x keeps warming up and will overheat at time $i1$. A good solution is to detect chip x at an earlier time, e.g. $t4$, and continue its test with schedule S3, as shown in (b).

schedule for o . In case of a set of manufactured chips with large temperature variations, a global thermal-safe offline schedule will be based on the hottest chip in the set. This test schedule will introduce unnecessary cooling time for most of the chips, leading to a very slow test.

We have proposed a technique, in [5], to address the above problem with the help of a chip classification scheme. This scheme consists of several test schedules for different temperature error ranges. After applying a short test sequence for warm up, the actual temperature is sensed and the proper test schedule is selected. Therefore, the hotter chips will test with a slower schedule, while the colder chips will test faster. The overheating issue is solved and the test application time is not unnecessarily long. This approach works fine assuming the simplest model of process variation (time invariant temperature error), as shown in Fig. 1(a-c).

However, in the real world with large process variation, the thermal behavior is time variant and the technique presented in [5] will not be able to achieve high-quality schedules. The variation of thermal response with time is illustrated in Fig. 1(d). In this case, the temperature of chip x gradually lifts up compared to chip o . A scheduling method capable of capturing temporal deviations is therefore required. The temperature behavior shown in Fig. 1(d) is captured in Fig. 2(a) with more details. The lift up of the temperatures of chip x becomes visible at $t3$, as shown in Fig. 2(a). Since x will only overheat after $t4$, both chips can be safely tested with schedule S1 up to $t4$. At $t4$, the actual temperature of the chip under test can be

obtained via sensors. The actual temperature can then be compared to a *Threshold* and two different situations can be identified:

$$\begin{cases} x; & \text{if } \text{Temperature}(t4) > \text{Threshold} \\ o; & \text{if } \text{Temperature}(t4) \leq \text{Threshold} \end{cases}$$

For the rest of the test, after $t4$, two dedicated schedules, S2 and S3, are generated in the offline-phase for o and x , respectively. Therefore, in the online-phase the test of o continues using schedule S2, as in Fig. 2(a), and the test of x continues using schedule S3, as in Fig. 2(b). In this illustrative example, at the end of S1, the schedule does a branching to either S2 or S3 based on the actual temperature. This information and the branching condition can be captured in a branching table, B1 in Fig. 2.

The segments of the schedule which are executed sequentially without branching are called linear schedules. An adaptive test schedule consists therefore of a number of branching tables in addition to multiple linear schedule tables. The cores will go through a very similar experience as mentioned above for chips. Therefore, a solution similar to what described above will work for the cores as well.

III. PROBLEM FORMULATION AND COST FUNCTION

Our goal is to generate an optimal adaptive test schedule, offline. The input consists of a SoC design with a set of cores and their corresponding test sequences. The floor plan, the thermal parameters, and the static/dynamic power parameters for the chip are given. The probability distributions that represent the temperature deviations are also given.

The desired adaptive schedule minimizes the test application time and overheating. These objectives are encapsulated into a cost function which is introduced later in this section. The desired schedule satisfies the two following constraints. The first constraint is the available test bus width; it limits the number of simultaneously active cores. The second constraint is the available ATE memory which limits the schedule size; indirectly, it also limits the total number of sensor accesses.

Similar to [7], a comprehensive cost function is introduced by combining the cost of the overheated chips and the cost of the test facility operation. These two contributors to cost go against each other. In order to prevent the overheating more cooling is required. The extra cooling cycles increase the test application time as shown in Fig. 1 (c) compared to Fig. 1 (a). Excessively long test application time means underutilization of the test facility. On the other hand, without enough cooling some chips will overheat. The Cost Function (CF) is defined as follows:

$$CF = TAT/ATS + BC \times TOP \quad (1)$$

The CF consists of two terms; the first one represents the test facility cost and is equal to Test Application Time (TAT) divided by the Applied Test Size (ATS). This term shows what volume of test could be applied by test facility per time unit. The second term represents the overheated chip cost and is equal to Test Overheating Probability (TOP) multiplied with a Balancing Coefficient (BC). The TOP is the number of overheated chips per number of chips entering the test facility.

The BC is used in order to balance the cost of the overheated chips against the cost of the test facility. Expensive chips will result in a larger BC and expensive test facility will result in a smaller BC. The exact reasoning behind the CF is given in [7].

IV. TEMPERATURE ERROR MODEL

The temperature error has various sources including process variation, ambient temperature fluctuations, voltage variations, simulator errors, and the temperature dependent errors, e.g. static power (leakage). Temperature error is broadly categorized into spatial error and temporal error. A temperature error model gives the probabilities of the temperature error values for each core (spatial) and for each test cycle (temporal). The spatial temperature error model gives the initial error distribution and the temporal temperature error model is used to recursively estimate the error distribution for the next test cycle.

The spatial temperature error is defined by a discrete statistical distribution, which assigns probabilities to temperature error ranges known as *error clusters*. The temporal temperature error is represented by a discrete-time model, i.e., the temperature error is fixed during a period and then it changes discretely from one period to the next. Therefore, the temporal temperature error model has two pieces of information, the period which is called *temporal error period* and a table of error change probabilities.

For a SoC with as many as C cores, the error clusters divide the C -dimensional error space into *error cluster cells* indexed using Cartesian system, i.e. $(l_0, l_1, \dots, l_{C-1})$. For example assume that in a 2-core SoC, each core has 2 *error clusters*. The 2-dimensional error space is divided into 4 *error cluster cells*, indexed with $(0, 0)$, $(0, 1)$, $(1, 0)$, and $(1, 1)$.

V. LINEAR SCHEDULE TABLES

A linear schedule table, as discussed in section II, captures a schedule without branching (offline). The linear schedule table entries (time) should be optimized in the offline-phase. In order to simplify the search space, the possible times are assumed to be multiples of a constant, denoted by *linear scheduling period*. The *states* in the linear schedule tables are generated using the heuristic proposed in [4].

The estimated temperature is updated periodically with *linear scheduling period* by correcting the cores' simulated temperatures with *representative temperature error* value for each core. The estimated temperature is then used to compute the static power and to determine the "state" of the cores for the next *linear scheduling period*. The *representative temperature error* is updated periodically with *temporal error period* while the estimated temperature, static power, and state of the cores are updated periodically with *linear scheduling period*. After updating the state of cores, the dynamic cycle-accurate power sequence for the next *linear scheduling period* is computed. Having dynamic and static power sequences, the next *linear scheduling period* is thermally simulated. A number of linear schedule tables (edges) which are connected using a number of branching tables (nodes) will form the schedule tree, similar to Fig. 3(a).

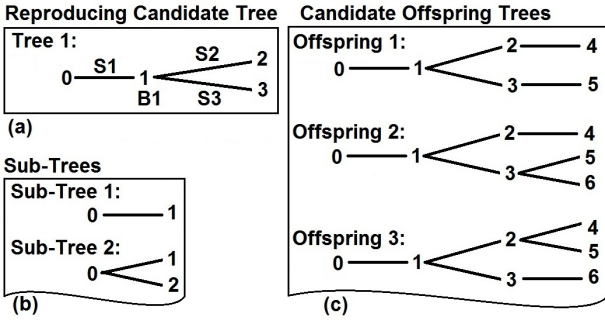


Figure 3. Reproducing tree, sub trees, and candidate offspring trees. For S1, S2, S3, and B1 in (a) refer to Fig. 2.

VI. ADAPTIVE TEST SCHEDULING

The adaptive method works as follows. During test, the actual temperatures (of selected cores) are read (at selected moments) and the gaps among sensor readouts are filled with thermal simulation. Chips are dynamically classified into one of the *chip clusters* and are tested using the corresponding schedule. At each adaptation moment the *chip clusters* change into a new scheme which is optimum for the new situation. The parameters that affect the efficiency of the adaptive method are the moments when branching/adaptation happens, the number of branches (linear schedule tables), and the branching condition (*chip clustering*). For example in Fig. 2, the adaptation is happening at t_4 , the number of branches is 2 (two linear schedule tables after the adaptation point), and the branching condition is a comparison with the *Threshold*. The problem is summarized into the two following sub-problems.

1. How many *chip clusters* (branches or linear schedule tables) at each possible adaptation point (node) are needed? One branch means no branching and no sensor reading.
2. What is the proper *chip clustering* into the given number of *chip clusters*? The number of *chip clusters* is known from question 1. Depending on the *chip clustering* some cores may do not need sensor readout.

When the answer to question 1 is one, question 2 is skipped. These two questions are then formulated into two different forms: the first question is described as a tree topology and the second question is to find the optimum *chip clustering* for the nodes of the specific tree topology.

A candidate schedule tree is generated by putting a possible tree topology together with a possible corresponding *clustering*. Since the number of candidate trees is the product of the tree topology alternatives and the *chip clustering* alternatives, the search space is unacceptably large. In order to reduce the search space, a constructive method is used. The schedule tree is constructed by adding sub-trees (small partial trees) to its leaves. A sub-tree consists of a small number of linear schedule and branching tables which makes it possible to be clustered and optimized (scheduled) at once.

For example, assume that there is a reproducing tree, *Tree 1*, as shown in Fig. 3(a). The linear schedule tables of Fig. 2 correspond to the edges of *Tree 1* while the branching table corresponds to node 1, as shown in Fig. 3(a). Two sub-trees with 1 and with 2 edges are shown in Fig. 3(b). *Tree 1* has two leaves, which combinations of sub-trees are added to them in

order to generate the offspring as shown in Fig. 3(c). *Offspring 2* is generated by attaching the *Sub-tree 1* to node 2 of *Tree 1* and attaching the *Sub-tree 2* to node 3 of *Tree 1*. The sub-tree scheduling is explained in section A. In section B, it is explained how the trees are constructed and selected.

A. Sub-tree Scheduling

A heuristic is used to find the near optimal sub-tree, by using the partial cost function of sub-tree *clustering* alternatives. When the schedule is a tree, the expected values of *test application time* (TAT), *applied test size* (ATS), and *test overheating probability* (TOP) which are denoted by ETAT, EATS, and ETOP should be used in the CF computation, Eq. (1), in order to utilize the temperature error statistics. The expected values are computed as each edge is being scheduled. The *chip clustering* at each node is done in a C-dimensional space and each *chip cluster* consists of a certain combination of *error cluster cells*. A candidate sub-tree *clustering* is a set of node *clustering* alternatives. For each candidate sub-tree topology there are a number of candidate *clustering* alternatives, which labels the nodes' *error cluster cells* with their corresponding *chip clusters*. Each *chip cluster* for a node corresponds to an edge branching out of that node (equivalent to a linear schedule table). Each node has its dedicated Error Cluster Cell Labeling (ECCL) as follows.

$$ECCL = \{ECCL(n, (l_0, l_1, \dots, l_{c-1})); n = 0, 1, \dots, N - 1\} \quad (2)$$

Having the *error cluster cell labeling* corresponding to an edge, the edge is scheduled (linear schedule table is determined). The candidate sub-tree *clustering* is evaluated based on the optimized linear schedule tables and optimized branching conditions. A heuristic explores the candidate *clustering* alternatives to find the optimal *clustering*. The *Error Cluster Cell Probabilities* of nodes and the nodes probabilities are computed using the temperature error model and based on the *clustering* of the ancestor nodes. This information is then used to compute ETAT, EATS, and ETOP. The procedure of computing *Error Cluster Cell Probabilities* could be found in [7]. Having the ETAT, EATS, and ETOP values, the partial cost function is computed. Particle swarm optimization is used to explore the *clustering* alternatives. The search space is the collection of different alternatives of Eq. (2). For example for a SoC with 2 cores and for a sub-tree similar to *offspring 3* in Fig. 3, an alternative solution is the following.

$$\{ECCL(0, (0,0)) = 0, ECCL(0, (0,1)) = 0, ECCL(0, (1,0)) = 1, \\ ECCL(0, (1,1)) = 1, ECCL(1, (0,0)) = 1, ECCL(1, (0,1)) = 1, \\ ECCL(1, (1,0)) = 1, ECCL(1, (1,1)) = 0\}$$

As it is shown above, a possible solution for sub-tree scheduling is coded by labeling the temperature *error cells* with a cluster label for each branching node in that sub-tree. An example of solution coding for a single node sub-tree similar to *sub-tree 2* in Fig. 3 (b) is illustrated in Fig. 4. The solution belongs to a 2-core design with 3 temperature *error clusters* per core. The number of branches, i.e. number of chip clusters, in the corresponding sub-tree is 2.

The 2-dimensional solution space is mapped into a vector for simplicity. As it is shown in Fig. 4, the cell order (mapping rule from 2-D space to a vector) is static and there is no need to include it in the solution vectors. (This is also true for the nodes and their order.) This solution space is then explored by

Core 1	2	Label = 1	Label = 1	Label = 1
	1	Label = 0	Label = 0	Label = 1
	0	Label = 0	Label = 0	Label = 1
		0	1	2

Solution Encoding: 0, 0, 0, 0, 1, 1, 1, 1, 1
Cell Order: (0,0), (0,1), (0,2), (1,0), (1,1), (1,2), (2,0), (2,1), (2,2)

Figure 4. Solution encoding for a single node in sub-tree scheduling.

particle swarm optimization in order to find a near optimal solution. Particle swarm optimization is presented in section VII.

B. Tree Construction

The construction starts with a root node and in each iteration the reproducing candidate tree extends and multiplies by adding possible combinations of sub-trees to its active leaf nodes, as shown in Fig. 3. Then, a small number of promising reproducing candidates (similar to Fig. 3(a)) are selected out of the candidate offspring trees (partially shown in Fig. 3(c)). The selection process guarantees the ATE memory constraint and provides the freedom to put more clusters in the more beneficial regions. Such a freedom is provided by the virtue of a Scaled Cost Function (SCF) which is used as the selection criterion. SCF is defined as:

$$SCF = CF \cdot (adjusting_offset + number_of_nodes) \quad (3)$$

The *cost function* (CF) is scaled by the tree's number of nodes plus *adjusting_offset*. Now, adding nodes to the tree is only beneficial if it gives a reasonable cost reduction otherwise a smaller tree may get a lower *scaled cost function* and be selected, while bigger trees are discarded. The effect of the number of nodes is adjusted by *adjusting_offset*. A larger *adjusting_offset* promotes having more branches, especially near the tree's root.

The number of the sub-tree topologies is controlled with the sub-tree length and the maximum allowed number of branches per node. Increase in the sub-tree length will improve the global optimality and increase in the allowed number of branches per node improves the *chip clustering* resolution, but both will increase the offline scheduling effort which is denoted by *CPU time* in this paper.

VII. PARTICLE SWARM OPTIMIZATION

The sub-tree scheduling, as explained in section VI-A, relies on particle swarm optimization to explore the solution space for a near optimal schedule. *Particle swarm optimization* is an optimization heuristic which mimics the social behavior of a *swarm* searching for food [8]. Each individual member of the *swarm* is called a *particle*. A *particle* is represented by two attributes, its *location* and its *velocity*. The *location* in fact is a *solution* which, usually, is represented by a Cartesian coordinate system. The dimensions in the coordinate system are analogous to the *genes* in a *chromosome* in a genetic algorithm. The *velocity* keeps the *particles* moving in the search space. Each *particle* remembers its *previous best* location, and in addition to this individual memory, the *swarm* remembers the best location any of its *particles* have visited before, the *global best*. The *previous bests* and the *global best* are then used to give a hint to the random *velocities*. A

canonical form of the *particle swarm optimization* is expressed with the following equations.

$$\begin{aligned} velocity = & 0.7298 * (velocity + \\ & + 2.05 * random1 * \\ & * (previousBestSolutions - solution) + \\ & + 2.05 * random2 * \\ & * (globalBestSolution - solution)) \\ solution = & Round(solution + velocity) \end{aligned} \quad (4)$$

This canonical form of the particle swarm optimization uses Eq. (4) to update the *velocity*. The coefficients in Eq. (4) (0.7298, 2.05, and 2.05) are given as a part of the chosen canonical form [8]. The *random1* and *random2* are two distinct random numbers between 0 and 1 which are renewed iteratively. The *solution* and *velocity* on the right hand side of Eq. (4) are the current values, and the left hand side *velocity* is the next value. Since the *solution*, in this paper, is a vector of natural number, the next *solution* is the rounded sum of the current *solution* and next *velocity*, as represented in Eq. (5). The *particle swarm optimization* is presented in the following as a pseudo code.

1. Generate the initial *swarm*.
2. Generate random initial *velocities*
 - a. Limit the range of the random number to the number of *chip clusters* for the corresponding node.
3. Evaluate the solutions. This is the most time consuming step.
4. Find the best solutions as follows.
 - a. Loop for all *particles*.
 - i. If the current *location* is better than the *previous best location* replace it and check if it is better than the *global best*, if so, replace the *global best* as well. (For the first iteration, copy the current solution as *previous best*, and find the *global best* among the *previous best* solutions.)
5. If the termination condition is met, exit with the *global best* as final solution. Otherwise continue with point 6.
6. Update the *Swarm* as follows.
 - a. Loop for *particles*.
 - i. Loop for *cells*.
 1. Update the *velocities* according to Eq. (4).
 2. Update the solution (*particle's location*) according to Eq. (5).
 3. Saturate (limit) the solution. It means that if the *location* is outside the valid search space, make it equal to the corresponding extreme and reset the corresponding component of the *velocity* to 0.
7. Check if all clusters exist, as follows.
 - a. Loop for all *particles*.
 - i. Loop for nodes.
 1. If a *chip cluster* is missing, move the *particle* to its *previous best* location.
8. GO TO point 3.

In order to see how *particle swarm optimization* works for the sub-tree scheduling, we focus on a single *error cluster cell*. Assume that the lower cluster labels correspond to more cooling. Therefore, a negative *velocity* guides the *error cluster cell* toward more cooling and if more cooling in this iteration helps to reduce the cost, it is reasonable to keep a negative *velocity* for the next iteration as well. If such a move makes the *particle* the best in the *swarm*, it will affect the *velocities* of other *particles* as well. The usefulness of this behavior has been demonstrated by the experimental results given in the next section.

VIII. EXPERIMENTAL RESULTS

In our experiments, the temperature simulation is done using HotSpot [9]. The static power is computed using the method given in [10]. Other elements of the experimental setup are similar to [7]. Experiments are performed using ITC'02 benchmarks. The proposed method which uses particle swarm optimization for sub-tree scheduling is compared with a similar method which uses genetic algorithm for the same purpose.

Scheduling for a SoC requires a large number of sub-trees to be optimized, however, only some of them are used in the construction of the finally selected tree. Each sub-tree is optimized using particle swarm optimization (or genetic algorithm) and as a result, a single SoC test scheduling includes a large number of executions of the particle swarm optimization algorithm (or genetic algorithm). This means that each row in the tables given in this section corresponds to multiple instances of the particle swarm optimization (or genetic algorithm).

In the following tables, **CF** is the cost of the schedule as given in Eq. (1) and **Mem** is the ATE memory volume which is required to store the schedule (not the test sequences). **CPU time** is given in "hour : minute" format.

A number of experiments with different problem sizes are reported in Table I. The problem size is the product of the number of the *error clusters* and *memory limit*. The results show that particle swarm optimization outperforms the genetic algorithm most of the time. It takes often less CPU time (row 1 to 4) to produce the same results or can obtain better solutions (row 6).

We have also studied the impact of the population size on the performance of the particle swarm optimization algorithm. The results are given in Table II for an identical experimental scenario. A smaller population size implies a shorter CPU time as shown in the table. Therefore, it is important to reduce the population size as long as it does not degrade the quality of the solution too much. Particle swarm optimization finds the best schedule with the best cost equal to 5.409 and size equal to 1880 with population size equal to 20 for this example.

Particle swarm optimization is able to find a solution with a cost of 5.497 (only 1.6% larger than the best) with a population size as small as 5 and with a CPU time as short as 5 hours. In contrast, the genetic algorithm requires, for this example, 11 hours with a population size of 30 to find such a solution. Another advantage of the particle swarm optimization is its smaller code size and therefore its shorter implementation (coding) time.

A comparison between the traditional and the proposed test scheduling methods is reported in Table III. The traditional methods include the Offline method (only one linear schedule is used) and the Hybrid method (similar to [5]). Two experimental scenarios are denoted by *es1* and *es2* in the table. The proposed adaptive method has reduced the cost compared to offline and hybrid methods. The problem size for *es1* is smaller than that of *es2*, and with similar memory limits, the cost reduction is larger for *es1*. The cost reduction offered by the proposed adaptive method is about two times the reduction that hybrid method offers and this demonstrates the advantage of the proposed method.

TABLE I. RESULTS FOR DIFFERENT PROBLEM SIZES

Problem size	Genetic Algorithm			Particle Swarm Optimization		
	CF	Mem	CPU Time	CF	Mem	CPU Time
1800	5.96	460	2 : 33	5.96	460	2 : 21
3000 (a)	5.58	920	3 : 58	5.58	920	3 : 26
3000 (b)	5.97	460	5 : 21	5.97	460	2 : 57
5000	5.82	980	5 : 55	5.82	980	5 : 03
6000	5.30	1540	4 : 15	5.30	1540	4 : 48
10000	5.50	1800	10 : 51	5.41	1880	11 : 20

TABLE II. RESULTS FOR DIFFERENT POPULATION SIZES FOR PARTICLE SWARM OPTIMIZATION

Population Sizes	CF	Mem	CPU Time
20	5.409	1880	11 : 20
10	5.497	1800	6 : 47
5	5.497	1800	4 : 54

TABLE III. COMPARISON OF TRADITIONAL AND PROPOSED METHOD

Methods	Results							
	CF		Mem		CF		Mem	
	<i>es1</i>	<i>es2</i>	<i>es1</i>	<i>es2</i>	<i>es1</i>	<i>es2</i>	<i>es1</i>	<i>es2</i>
Offline	5.96	5.97	460	460	<i>Relative to the Offline</i>			
Hybrid	5.58	5.82	920	980	94%	97%	200%	213%
Adaptive	5.30	5.50	1540	1800	89%	92%	335%	391%

es1 Experimental scenario number 1

es2 Experimental scenario number 2

IX. CONCLUSION

This paper presents an adaptive SoC test scheduling technique in order to deal with spatial and temporal temperature deviations, caused by large process variations. The proposed technique relies on particle swarm optimization in order to generate an efficient test schedule tree. During the test, on-chip temperature sensors are used to monitor the actual temperatures of the different cores and to guide the selection of proper test schedules accordingly. In this way, the overall test cost will be minimized.

Experiments are made in order to evaluate the particle swarm optimization in comparison with the well established genetic algorithm. The experiments demonstrate the superiority of the proposed approach over the traditional methods and show that the proposed particle swarm optimization algorithm outperforms the genetic algorithm most of the time.

REFERENCES

- [1] K.-T. Cheng, S. Dey, M. Rodgers, and K. Roy. "Test challenges for deep sub-micron technologies." pp. 142-149, DAC 2000.
- [2] D.R. Bild et Al. "Temperature-aware test scheduling for multiprocessor systems-on-chip." pp.59-66, ICCAD 2008.
- [3] C. Yao, K. K. Saluja, P. Ramanathan. "Partition based SoC test scheduling with thermal and power constraints under deep submicron technologies." pp. 281-286, ATS 2009.
- [4] Z. He, Z. Peng, P. Eles. "Simulation-driven thermal-safe test time minimization for System-on-Chip." pp. 283-288, ATS 2008.
- [5] N. Aghaee, Z. He, Z. Peng, and P. Eles. "Temperature-aware SoC test scheduling considering inter-chip process variation." ATS 2010.
- [6] C. Yao, K. K. Saluja, P. Ramanathan. "Thermal-aware test scheduling using on-chip temperature sensors." VLSI Design 2011.
- [7] N. Aghaee, Z. Peng, and P. Eles. "Adaptive Temperature-Aware SoC Test Scheduling Considering Process Variation." DSD 2011.
- [8] R. Poli, J. Kennedy, and T. Blackwell. "Particle swarm optimization, An overview." Swarm Intell., Vol. 1, No. 1, pp. 33-57, 2007.
- [9] W. Huang, M. R. Stan, K. Skadron, K. Sankaranarayanan, S. Ghosh, and S. Velusamy. "Compact thermal modeling for temperature-aware design." DAC 2004.
- [10] W. P. Liao, L. He, K. M. Lepak. "Temperature and supply voltage aware performance and power modeling at microarchitecture level." IEEE Trans. CAD, Vol. 24, No. 7, pp. 1042- 1053, 2005