

Enabling Efficient Post-Silicon Debug by Clustering of Hardware-Assertions

M. H. Neishaburi
McGill University
Montreal, QC, Canada
Mh.neishabouri@mail.mcgill.ca

Zeljko Zilic
McGill University
Montreal, QC, Canada
zeljko.zilic@mcgill.ca

Abstract—Bug-free first silicon is not guaranteed by the existing pre-silicon verification techniques. To have impeccable products, it is now required to identify any bug as soon as the first silicon becomes available. We consider the Assertion Based Verification techniques for the post-silicon debugging based on the insertion of hardware checkers in the debug infrastructure for complex systems on chip. This paper proposes a method to cluster hardware-assertion checkers using the graph partitioning approach. It turns out that having the clusters of hardware-assertions and controlling each cluster selectively during the debug mode and normal operation of the circuit makes integration of assertions inside the circuits easier, and causes lower energy consumption and efficient debug scheduling.

I. INTRODUCTION

With the increased complexity of systems, and driven by unquenchable demand for a large set of micro-architectural and other features in a system, design errors and bug become frequent and difficult to track. The existing pre-silicon verification techniques which rely on the formal or functional (dynamic verification) cannot guarantee that the first-silicon works perfectly without any bug. Almost two thirds of newly manufactured SoCs product suffers from bugs at the first silicon [1]. Logic errors are still among the main causes of failures, including the fact that pre-silicon verification is applied to the model of the IC, and not the actual product [6].

The increase in the time-to-market of new products may cause the significant loss of market share, or even complete loss of revenue [1]. Therefore, as soon as the first silicon becomes available, it is required to identify any bug either in the design (corner cases) or the bugs related to manufacturing such as process variation and packaging.

Assertion-Based Verification (ABV) is acknowledged as being the instrumental and efficient pre-silicon verification technique. Armed with temporal logic and extended regular expressions, the PSL (Property Specification Language, IEEE 1850 standard) and the SVA (System Verilog Assertions), perform as languages to describe the expected behavior of a design.

To expand the functionality of assertions beyond design and verification at the pre-silicon debugging, a *checker generator* tool can be exploited to get assertions to operate with the CUD described at the RTL level. Individual assertion, once converted into a circuit form, is also referred to as *Hardware-Assertion* or a checker. Here, we have used the MBAC checker

generator to generate hardware-checker or hardware-assertion from either PSL or SVA assertions[7].

Post-silicon validation involves three major activities: detecting problem through specific embedded modules for debugging or reutilizing already available scan architecture for the manufacturing test by applying proper stimulus; localizing and identifying the root cause of the problem; and, correcting or bypassing the problem. Post-silicon bug localization involves identifying the location-time pair for the bug and is the most time-consuming step in the process.

For activating these hardware-assertions and capturing their violation signals, debugger module inside a SoC must be provided by a suitable debug infrastructure [2][11].

As systems become larger and more distributed (and possibly wireless [3]), having clusters of hardware-assertions and controlling each cluster selectively during the debug mode and normal operation of the circuit makes integration of assertions easier, and causes lower energy consumption and efficient debug scheduling. In this work, we investigate a method to cluster hardware-assertions based on the graph partitioning approach by using the input cone graph of hardware-assertions.

II. RELATED WORK

Existing post-silicon debug techniques can be divided into two approaches: trace-based and scan-based. Various implementations for either of them have been realized by previous studies [4][6]. The primary goal in the scan based debug techniques is reusing the internal scan chains which were used during the manufacturing test. At the first step, whenever a specific trigger or hardware checker fires all the internal state elements using these available scan chains will be captured; thereafter these captured data can be offloaded using these scan-out operation of these scan chains [9]. Due to the consecutive stops and resumption during scan dump, this method will not provide the real-time debug information [6].

Trace buffer serves as a space to keep a snapshot of system under debug including signals and states of the system whenever a certain events take place [11]. They have been widely used in legacy debug and logic analysis systems.

The integration of hardware assertions in scan-based run-stop debug infrastructure and in a debug trace infrastructure has

been investigated in [2]. However, this study provides no solution for clustering of related hardware assertions. We have observed that there are still two challenging issues that need to be addressed. First how to cluster Hardware-Assertions in post-silicon chip; second how to use of information from firing signal of each cluster of assertion to spatially isolate the candidate error sites and speed up the debug process.

III. PROPOSED METHOD

Our method tries to cluster the hardware-checkers inside the CUD. To apply our proposed clustering approach to the CUD, at the first step, a directed graph from the circuit gate level net-list must be extracted. The idea of extracting graph from the CUD was used before for post-silicon debug in [10]. Each vertex in this graph represents the Flip-Flop in the design and the directed line shows the combinational circuits or wire which connects these Flip-Flops together. Fig. 1 shows the sample circuit and its corresponding directed graph. The combinational parts of the circuit between storage units were represented by edges and for each flip-flop there exist a vertex inside the directed graph.

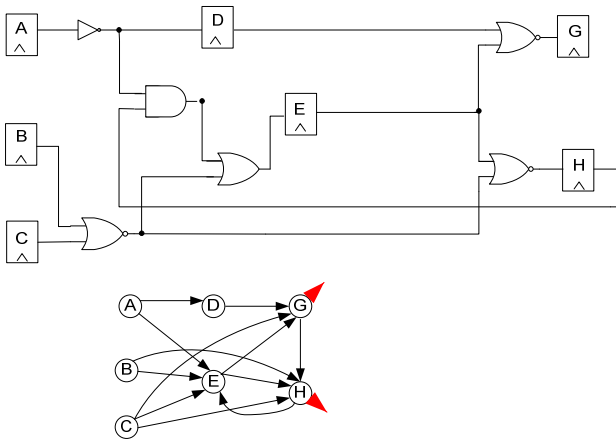


Figure 1. Directed graph Extraction from the given circuit

The second step is extraction of the Fan-in cone graph for each primary output. Fig. 2 shows the extracted fan-in cone graphs for primary outputs of the circuit in Fig. 1.

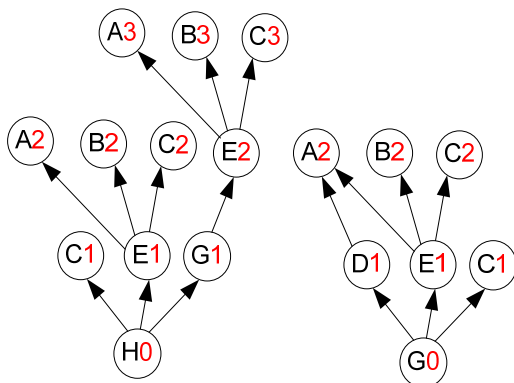


Figure 2. Fan-in cone graph of primary outputs

In this graph, each node represents a storage element inside the CUD. The first primary output is feed by the value of storage elements H0 the node at the first level of this graph. The second primary output is connected to node G0. The directed edge between nodes G1 and H0 in this figure indicates that if there is an error in G1 this error will be propagated after one cycle to the H0. Furthermore, the directed edge between node E2 and G1 implies that an error in E2 will be propagated to G1 and H0 after one and two cycle respectively.

The next step in the process of hardware-assertion clustering is finding out the fan-in cone of each assertion inside the CUD based on the Fan-in cone graph of each primary output. Related nodes to the hardware-assertions input signals should be specified at first in this step.

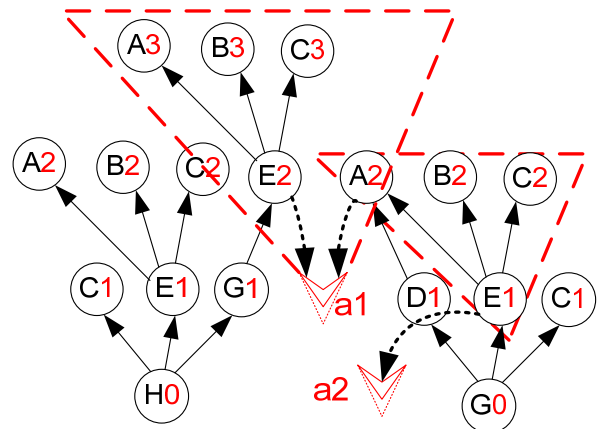


Figure 3. Fan-in cone graph of Hardware-Assertions

Fig. 2 illustrates fan-in cone graphs of each primary output of our simple example of CUD and assertions inside the CUD. Every triangle in this figure represents an assertion inside the CUD. The dotted lines specify the input signals of assertion which are connected to the specific nodes in CM. These nodes are corresponding to the storage elements inside the CUD.

For example in Fig. 3 the input signal of a1 (Assertion1) has been connected to E2 in fan-in cone graph of primary output H and A2 in fan-in cone graph of primary output G; also, a2 (Assertion2) is connected to E1 in the Fan-in cone graph of primary output G.

As soon as hardware-assertions get connected to their related nodes, the input cone of all of them inside the CUD can be specified. For instance, in Fig. 3, fan-in cone of A1 is {E2, A3, B3, C3, A2} and fan-in cone of A2 is {E1, A2, B2, C2}.

A. Clustering Hardware Assertions Based on their Fan-in cone graph

In our proposed approach, the placement of Hardware-Assertion and their fan-in cones inside the gate level net-list of CUD is modeled by a new weighted graph called CM (Checker Map) = (V, E).

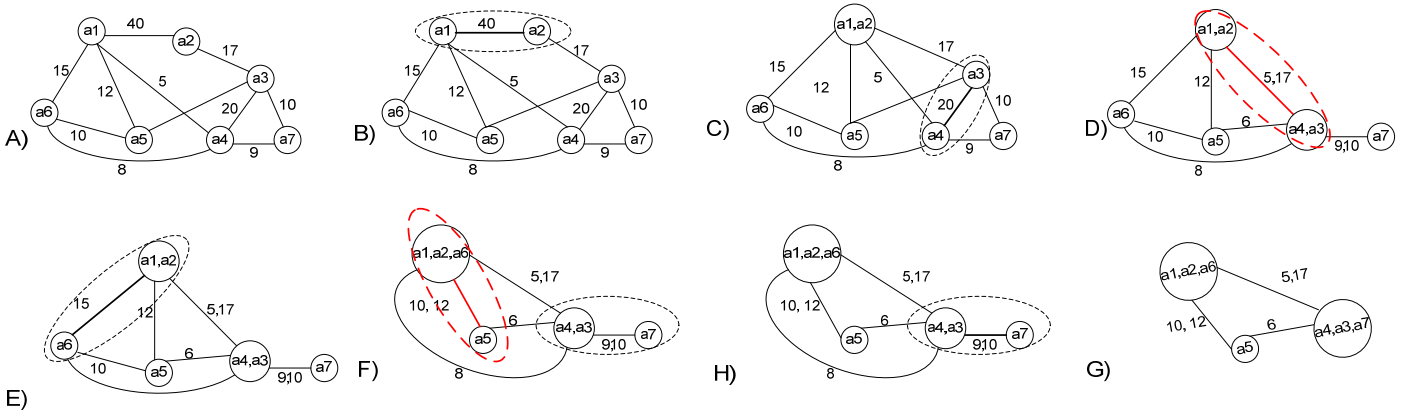


Figure 4. The process of graph clustering on the sample CM(Checker Map) graph

A set of vertices $V = \{v_1, v_2, v_3, \dots, v_n\}$ in CM represents the Hardware-Assertions inside the CUD. The common elements inside the fan-in cone of every coupled hardware-assertion will be indicated by an edge between corresponding two nodes. Furthermore, the weight of that edge represents the number of common elements in the fan-in cone graph of two related hardware-assertions. For instance, the CM graph of our sample circuit at Fig. 1 has two nodes $\{a1, a2\}$ and these two nodes are connected using an edge with the weight of one. To explain the partitioning algorithm on CM graph, we have considered a new CM graph shown in Fig. 4.

This algorithm takes as input the CM graph, the required number of partitions that we want to have, and the maximum number of elements that are allowed to be placed inside each cluster.

In fact, the partitioning algorithm should manage to create the required number of clusters which is given as the first parameter, while it should be also aware that during the iterative process of partitioning the number of elements inside each cluster does not exceed the maximum number of elements that are allowed to be placed inside each cluster.

As it was illustrated in Fig. 4, at each step of the partitioning, the edge with the largest weight is selected. Due to the fact that edge weight represents the common number of storage elements between fan-in cone of two assertions connected by that edge in CM, the larger the weight of the edge the higher the probability that two related assertions fire together and more increases in the chance of extracting required information from firing signal in a cluster of assertion to spatially isolate the candidate error sites.

After finding the edge with the largest weight, two connected nodes by this edge are chosen as a candidate to merge. Thereafter, the partitioning procedure will check whether by merging related nodes the numbers of nodes exceed from the maximum number of allowable nodes.

In Fig. 4 (B) since the weight of the edge between a1 and a2 is larger than that of the others, these two nodes will be merged together after making sure that the number of elements

in the new cluster, $\{a1, a2\}$ has less than the maximum number of allowable elements in each cluster.

After merging two nodes, the algorithm should update the CM graph. To update the CM graph any edge that went to $\{a1\}$ and $\{a2\}$ before, now should go to new composite node or a cluster $\{a1, a2\}$.

The iterative partition algorithm performs once again merge operation on the new updated CM graph and the node $\{a3\}$ and $\{a4\}$ is merged together in Fig. 4 (C). In the next iteration, by updating the related CM graph once again the edge with the largest weight is selected during the merge operation in Fig. 4 (D). However, since after merging two concerning clusters $\{a1, a2\}$, $\{a3, a4\}$ the number of elements in the new cluster exceeds the maximum number of allowable elements which is 3, merge Procedure refused to merge these two clusters. Therefore, the next largest edge is selected by the merge operation Fig. 4 (E).

The iterative partitioning algorithm based on the merge and update procedure continues until we create the required number of clusters. The final result of applying the iterative partitioning algorithm in our sample CM graph is shown in Fig. 4 (G). As this figure illustrates, eventually, we have three clusters of assertions and the hardware-assertion inside two of them share the maximum number storage elements. After clustering hardware-assertions inside the CUD into the specified number of clusters and taking into account the limitation on the number of elements in each cluster, we have to equip our post silicon debugging infrastructure by insertion of clusters of hardware-assertions in the debug infrastructure.

IV. INSERTION OF HARDWARE-ASSERTIONS CLUSTERS IN THE DEBUG INFRASTRUCTURE

Having partitioned hardware-assertions based on their input fan-in cone, we have to insert each cluster inside the debugging infrastructure. In fact, the related infrastructure should be equipped to control each cluster of hardware-assertion plus supplied by a mechanism to capture the violated signals of each cluster. To control a specific cluster of hardware-assertions, there must be a way to selectively enable

