# Topological Description and Interconnect Delay Modeling for Performace-Driven Partitioning and Placement of MCM Design

Christoforos Harvatis    Yusuf Tekmen    Sharad Mehrotra    Griff Bilbro i    Paul Franzon

Department of Electrical and Computer Engineering
North Carolina State University
Box 7911
Raleigh, NC 27695

## 1   Introduction

In the early design of a high-performance custom or semi-custom MCM, some of the required tasks that must be completed include the following: (1) assignment of functions to chips, (2) determine the floorplan of the chips, and (3) determine the floorplan of the chips within the MCM. We wish to complete this task so that, for the completed system, cost is minimized, and that there is a high probability that timing and noise constraints are met, thermal, DFT, DFM, etc. requirements are satisfied and that the design is routable with minimum effort.

In this paper, we discuss some issues related to partitioning and placement (or floorplanning) for timing-driven MCM design. In particular, we describe how a topological description is used to describe a combined partitioning and placement, and how linearized circuit responses can be used to drive the tool. We give some results that show that using these approaches, 80% or more of the timing and congestion constraints are met.

## 2   Topological description of the MCM architecture

In order to capture a topological description of the MCM architecture the approach in [?] is extended as follows: Grid lines are placed on the MCM board in both the horizontal and the vertical directions. The MCM board is divided into bins, whose number is equal to the product of horizontal grid lines times the number of vertical ones. For example, a 4X4 grid divides the board into 16 bins, while a 10X10 grid divides the board into 100 bins (Figure 1). All the bins have equal size, but their module capacity may vary according to the case. No matter how many modules are into a bin, the location of all of them and their pins is considered to be the center of the bin.

There is no specific number of cut lines. If the so far applied grid lines do not give a precise enough characterization of the MCM board and a good enough number of bins (number of bins is proportional to the griding lines), additional grid lines are placed automatically. The number of bins needed for a precise enough characterization of the MCM board is proportional to the design size and the accuracy we expect on our results. The more bins applied, the more efficient the wire length estimation becomes. However, the increment of the number of bins simultaneously increases the complexity of the problem and thus, the CPU time that the software tools need.

The periphery bins are reserved only for the MCM I/O pads or connectors (Figure 1). No other module is allowed to move into them. Depending on if the location of the MCM connectors is fixed or not, the connector/pads modules are allowed to move or swap positions into the periphery bins. The capacity and the footprint of the chips on the MCM board may vary. Our topological MCM description takes that under consideration, by providing flexibility during the chip assignment level. Each chip has a specific location on the MCM board according to its assigned bins, that also determine its capacity and shape. For example, in

Figure 1, a 4 chip assignment is shown, where all 4 chips are of the same size. In general, it is not required that all chips have the same size.

## 3  Wiring Delay Modeling for Early Design

In the early design phase we are dealing with the following issues:

- the assignment of macromodules, or cells, to ICs (partitioning);

- the assignment of the modules' approximate placement within the MCM and IC (floorplanning); and

- the assignment of wires to MCM or IC interconnect.

The objective of the early design is to determine these assignments so that performance/cost is maximized, and the resulting design is feasible in terms of having a high likelihood, when completed, of meeting timing and noise constraints, thermal constraints, and routability constraints.

If this problem is to be solved automatically, it is necessary to evaluate the objectives and constraints hundreds or thousands of times a second. To achieve this speed, the objectives and constraints must be linearized. For example, we must predict electrical delay and noise as a function of distance between I/Os. This prediction must be carried out without pre-supposing details such as topology constraints, exact stub length limits, routing details, etc.

Since the follow-on tools being developed for this project (pin assignment, global router, detailed router, and driver re-sizer) rely on an accurate pre-characterization strategy [?], this same strategy was investigated for the early design tools. The advantages of a pre-characterization approach are greater accuracy, and high speed evaluation at design time, no matter how complex the simulation model.

To determine, if a pre-characterization could be adequately linearized for application to the early design problem, a number of evaluations were conducted. Simulations were run on nets with one driver and various number of receivers. Different wiring delay models were constructed for on- and off- chip nets. During the off-chip wiring modeling, lossy transmission lines on a distributed RLC network were used, while, for the on-chip, a simple distributed RLC network was applied. A sufficiently large number of points within the design space was simulated, while the length of each individual branch was varying. The calculation of the total net length is made by adding the individual branches and making a scatter plot of this total branch length versus the settling delay to the far most receiver. A least squares fit on this plot gives a single linear equation for settling delay. A generic wiring delay model for all the nets - heavily loaded or not - can be obtained by relaxing the driver strength for different net classes. Driver tuning makes the wiring delay equation insensitive to the load of each net. It is up to the partitioning/placement tool to make the decision whether to use one generic wiring delay model for all the nets or different models according to the number of pins connected to each net. It should be noted that at this level we don't expect an exact routing, but an accurate enough wiring delay estimation for the partitioning/placement procedure.

One of the new features of the approach described on this paper is that different wiring delay models are used for on- and off-chip nets. During the wiring delay estimation phase, the partitioning/placement tool checks if the pins on the currently examined net are all on the same chip or not. Based on that checking, different wiring delay models are applied for each net and thus, a distinction between on- and off- chip wiring delay is made.

## 4  Combined Partitioning and Placement

In a timing-driven design, it is essential to combine partitioning and placement. If, instead, partitioning is conducted before placement, then it is quite likely that a non-feasible partitioning will be produced.

In order to do this, we divide the MCM into an N×N mesh of bins as shown in Figure 1. Each bin may, or may not, correspond to a portion of an IC. Partitioning and placement are determined simultaneously by assigning macromodules or cells to the bins such that the objective is maximized and the constraints satisfied. The partitioning and placement algorithm starts with an initial random partition and, after estimating its total wiring length and timing penalty, it performs multiple module movements from bin to bin until a partition with less timing penalty is reached and within the area constraints. After achieving a partition that meets all the timing constraints, it starts looking for a smaller total wiring length estimation.
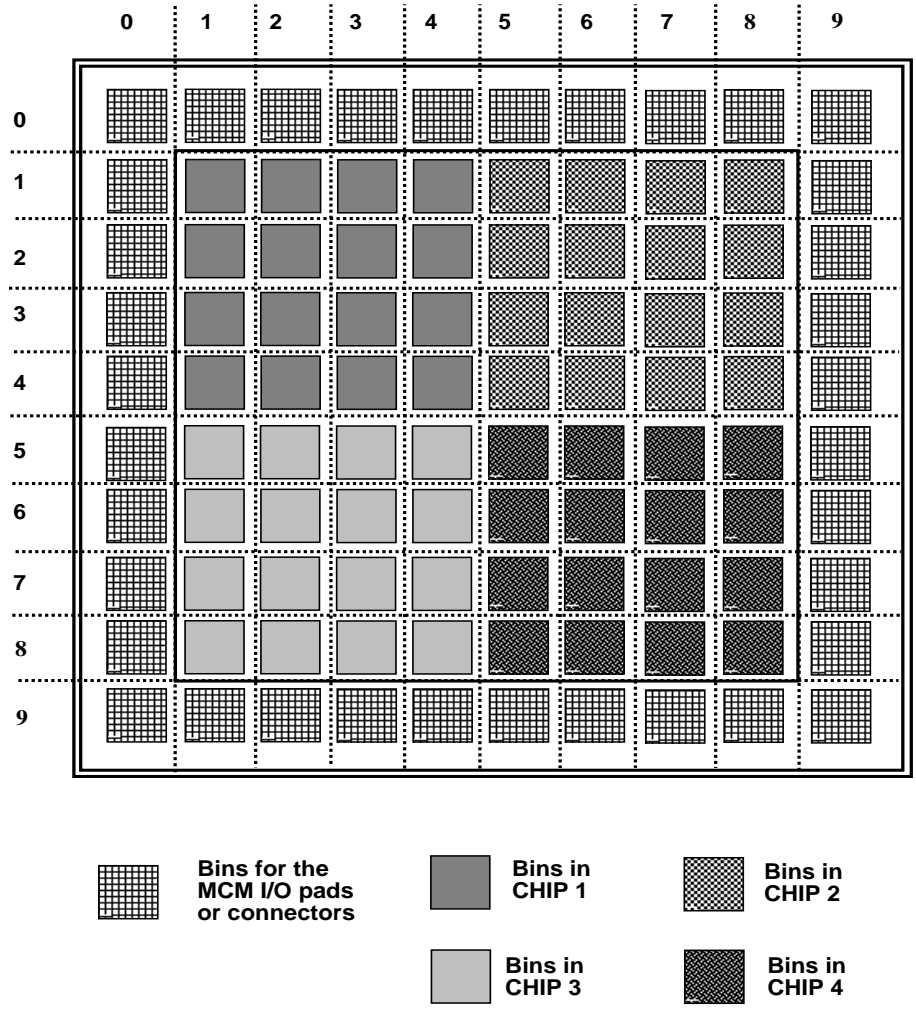
Figure 1: Topological Description of the MCM architecture (10X10 grid, 4 chips floorplan)

| Results based on the classical "2-stage" partitioning/placement approach | | | |
|---|---|---|---|
| Number of Chips | Average Total Wiring Length | % of Placements with 0 Timing Penalty | % of Globally Routed Designs |
| 4 | 6948cm | 57% | 40% |
| 13 | 6804cm | 63% | 51% |
| 16 | 6720cm | 60% | 50% |
| Results based on the combined partitioning/placement approach | | | |
| Number of Chips | Average Total Wiring Length | % of Placements with 0 Timing Penalty | % of Globally Routed Designs |
| 4 | 6890cm | 80% | 75% |
| 13 | 6671cm | 89% | 88% |
| 16 | 6532cm | 92% | 91% |

Table 1: Comparison between the two partitioning/placement approaches

Our purpose here is to demonstrate that this tool (which is based on the here proposed topological MCM description and interconnect delay modeling) produces a solution that is demonstrably routable and satisfies electrical constraints. This demonstration was done by testing results on the follow-on tool, the global router [?].

A benchmark MCC design (with timing constraints) was introduced to the partitioning/placement tool-which performs a timing-driven partitioning and placement at the same time on the so far described topological description. It's aim is to minimize total wire length and satisfy signal path timing constraints, while taking area constraints under consideration.

The MCC benchmark design models a next generation supercomputer on a 6 x 6 inch substrate with 37 Honeywell VHSIC gate arrays (treated as distinct macromodules) and 18 high density connectors. The chips are 1.5 x 1.5 cm with 35 mil TAB leads on a 4 mil pitch. The connectors are placed around the perimeter of the substrate. The net list contains 7118 signal nets and 14659 pins, and the multiplicity of each net varies from 1 to 188.

The MCC design was also introduced to a two stage partitioning and placement tool, that uses exactly the same partitioning and placement algorithms.

Using an 10X10 griding, many test cases were run under various alternatives as to number of chips, chip sizes and shape.The results are summarized in Table 1. The column "% of Placements with 0 Timing Penalty" describes the number of test cases that, after partitioning/placement, met detailed timing constraints. These results show that the combined partitioning/placement tool performed better than classical "2-stage" approach on the percentage of placed designs with zero timing delay penalty and their total wire length. These successfully placed designs (with zero timing delay penalty) were given as an input to the global router in order to check their feasibility. The column "% of Globally Routed Designs" indicates the percentage of test cases that, when routed, met detailed timing and congestion constraints. These results show that this tool is very likely to converge on a good result with only a few iterations.

## 5   Conclusions

A multi-layer topological description is required for the problem of custom/semi-custom MCM/IC partitioning and floorplanning (placement). Linearized delay measures are also required. In this paper we show that such a combination of measures lead a high probability of producing a partitioning/placement that meets timing and routing congestion constraints.

## Acknowledgements