

Mapping Applications onto the Chip

Midia Reshadi, Ph.D
CE Department of
Science and Research Branch
Islamic Azad University
(SRBIAU)
reshadi@iaau.ac.ir
<http://www.reshadi.net>

Background Picture
© DIGITAL STOCK

Agenda

- First step to the application specific Ics
- Partially Branch & Bound mapping
- Nmap

© Midia Reshadi

1st References

Energy-Aware Mapping for Tile-based NoC Architectures Under Performance Constraints *

Jingcao Hu Radu Marculescu
Department of Electrical and Computer Engineering
Carnegie Mellon University
Pittsburgh, PA 15213-3890, USA
e-mail: {jingcao.hu, radu}@ece.cmu.edu

Abstract — In this paper, we present an algorithm which automatically maps the *Devices* onto a generic regular *Network-on-Chip* (NoC) architecture such that the total communication energy is minimized. At the same time, the performance of the mapped system is guaranteed to satisfy the specified constraint through bandwidth reservation. As the main contribution, we first formulate the problem of energy-aware mapping, in a topological sense, and then propose an efficient branch-and-bound algorithm to solve it. Experimental results show that the proposed algorithm is very fast and robust, and significant energy savings can be achieved. For instance, for a complete adder-like test design, an average 83.4% energy savings have been observed compared to an ad-hoc implementation.

1. INTRODUCTION

With the advance of the semiconductor technology, the enormous number of transistors available on a single chip allows designers to integrate dozens of IP blocks together with large amounts of embedded memory. These IPs can be CPU or DSP core, video stream processor, high-bandwidth I/O, etc[1]. The richness of the computational resources places tremendous demands on the communication resources as well. Additionally, the shrinking feature size in the deep-sub-micron (DSM) era is continuously pushing interconnection delay and power consumption as the dominant factors in the optimization of silicon systems. Another consequence of the DSM effects is the difficulty in optimizing the interconnection because of the myriad manufacturing effects such as crosstalk, etc.

To mitigate these problems, Dally and Towles [2] have recently proposed a regular tile-based architecture where communications can be efficiently realized using an on-chip network (Fig. 1). As shown in the left part of Fig. 1, the chip is divided into regular tiles where each tile can be a general-purpose processor, a DSP, a memory subsystem, etc. A source is embedded within each tile with the objective of connecting it to its neighboring tiles. This, instead of routing design-specific global wires, can save the communication time given in [2], where a loose topology is adopted.

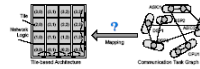


Fig. 1. Tile-based architecture and its mapping problem.

be achieved by routing packets via these embedded routers.

These last concepts come together to make this tile-based architecture very promising: structured network wiring, modularity and standard interfaces. More precisely, since the network wires are structured and wired beforehand, their electrical parameters can be very well controlled and optimized. In turn, these controlled electrical parameters make possible to use aggressively signaling circuits which reduce power dissipation and propagation delay significantly. Modularity and standard network interfaces facilitate re-usability and interoperability of the modules. Moreover, since the network platform can be designed in advance and later used for many applications, it makes sense to highly optimize this platform as its development cost can be amortized across many applications.

To exploit this regular tile-based architecture, the design flow needs the following three steps: First, the application needs to be divided into a graph of concurrent tasks. Second, using a set of available IPs, the application tasks are assigned and scheduled. Finally, the designer needs to decide to which tile each selected IP should be mapped such that the metrics of interest are optimized. More precisely, given the assigned/scheduled task graph which has been generated from previous two steps, this last phase determines the topological placement of these IPs onto different tiles. For instance, referring to Fig. 1, this step determines onto which tile (e.g. (0,0), (2,1), (3,2) etc.) each IP (e.g. ASIC, DSP, CPU, etc.) should be placed.

The first two steps described above are not new to the CAD community, as they have been addressed in the area of hardware-software co-design and IP-reuse [3]. However, the mapping phase (that is, the topological placement of the IPs onto the on-chip tiles) represents a new problem, especially in the context of the regular tile-based architecture, as it significantly

© Midia Reshadi

1st References

• Energy-Aware Mapping for Tile-based NoC Architectures Under Performance Constraints

– Authors : Jingcao Hu, Radu Marculescu

- Won Best paper award
- And
- Most Influential Papers of 10 Years in DATE Conference



© Midia Reshadi

2st reference

- **Bandwidth-Constrained Mapping of Cores onto NoC Architectures**
 – **Authors:** Srinivasan Murali, Giovanni De Micheli

© Midia Reshadi

Citation

- 293 Citation !!!!!



[Bandwidth-constrained mapping of cores onto NoC architectures](#)

S Murali, G De Micheli - 2004 - computer.org

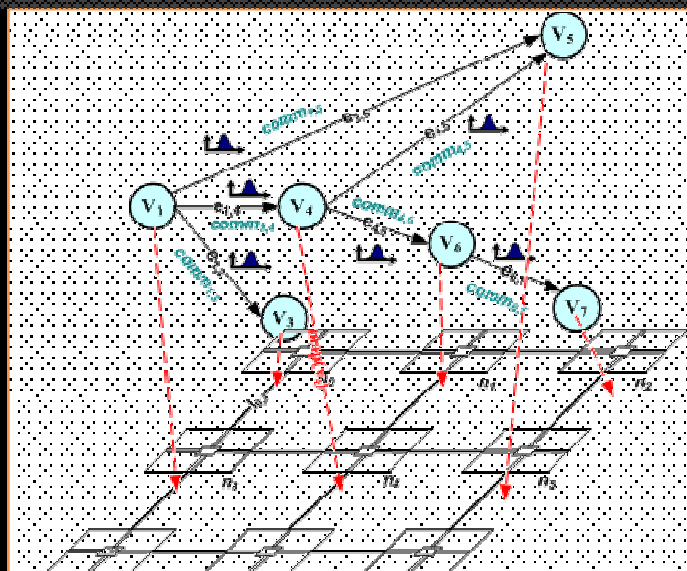
... ments on the links of a NoC are known for a particular ap- plication, thus the bandwidth constraints in the NoC archi- tecture need to be satisfied by the mapping. In [8], a branch and bound algorithm is proposed that maps cores onto a tile-based NoC architecture satisfying the ...

[Cited by 293](#) - [Related articles](#) - [All 19 versions](#)

[Cited by 293](#)

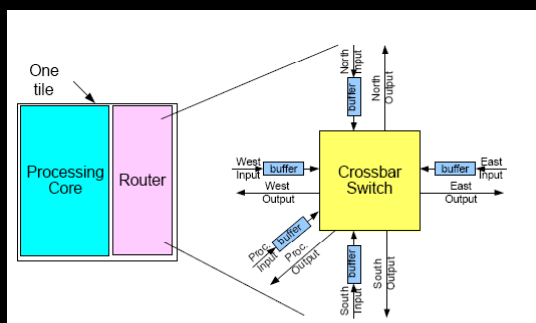
© Midia Reshadi

Mapping Basics



Midia Reshadi

Basic Definition:: tile



© Midia Reshadi

The Energy Model

- bit energy (E_{bit}) metric: for a router
 - T. T. Ye, L. Benini, G. De Micheli, "Analysis of power consumption on switch fabrics in network routers," *Proc. DAC, June 2002*.
- E_{bit} :

$$E_{\text{bit}} = E_{\text{Sbit}} + E_{\text{Bbit}} + E_{\text{Wbit}}$$
 - $E_{\text{Sbit}} \rightarrow$ The energy consumed by switch
 - $E_{\text{Bbit}} \rightarrow$ buffering
 - $E_{\text{Wbit}} \rightarrow$ interconnection wires
 - It is assumed that the buffers are implemented in SRAM or DRAM

© Midia Reshadi

The Energy Model

- Extending the router energy model to the tile-based network
- Note:
 - E_{Bbit} becomes dominant when congestion happens since accessing and refreshing the memory are very expensive in terms of power consumption.
- E_{Wbit} : The energy consumed on the wires inside the switch fabric.

© Midia Reshadi

The Energy Model

- E_{Lbit} : Energy consumed on the links
- **The average energy consumed in sending one bit of data from a tile to its neighboring tile**
 - $E_{bit} = E_{Sbit} + E_{Bbit} + E_{Wbit} + E_{Lbit}$
- Link length is typically in the order of *mm*,
 - $E_{Bbit} + E_{Wbit} \ll E_{Lbit}$
- **so:**
 - $E_{bit} = E_{Sbit} + E_{Lbit}$

© Midia Reshadi

The Energy Model

- Average energy consumption of sending one bit of data from tile t_i to tile t_j :

$$E_{bit}^{t_i, t_j} = n_{hops} \times E_{Sbit} + (n_{hops} - 1) \times E_{Lbit}$$

- n_{hops} : number of hops

© Midia Reshadi

Problem Formulation

- **Objective**
 - Mapping IPs onto different tiles such that
 - The total communication energy consumption is minimized
 - Guaranteeing the performance of the system

© Midia Reshadi

Problem Formulation

- **Definition 1:**
- **An Application Characterization Graph (APCG)**
- $G = G(C, A)$: A directed graph
 - c_i represents one selected IP/core
 - $a_{i,j}$ represents the communication from c_i to c_j .
- **Associated with each $a_{i,j}$ as arc properties:**
 - $v(a_{i,j})$: The communication volume (*bits*) from c_i to c_j
 - $b(a_{i,j})$: The bandwidth requirement from vertex c_i to c_j ,
 - The minimum bandwidth (*bits/sec.*) that should be allocated by the communication network.

© Midia Reshadi

Problem Formulation

- **Definition 2**
- **An Architecture Characterization Graph (ARCG)**
- $G' = G(T, P)$: A directed graph
 - Each vertex t_i represents one tile in the architecture
 - Each directed arc $p_{i,j}$ represents the routing path from t_i to t_j
- **Associated with each $p_{i,j}$ as arc properties:**
 - $e(p_{i,j})$: The average energy consumption (*joule*) of sending one bit of data from tile t_i to $t_j \rightarrow E_{i,j}^{t_i,t_j}$
 - $L(p_{i,j})$: the set of links that make up the path $p_{i,j}$

© Midia Reshadi

Problem Formulation

- Minimizing the communication energy consumption under performance constraints:
- **Given** an APCG and an ARCG that satisfy:

$$size(APCG) \leq size(ARCG)$$
- **Find** a mapping function $map()$ from APCG to ARCG which minimizes:

$$\min\{Energy = \sum_{\forall a_{i,j}} v(a_{i,j}) \times e(p_{map(c_i),map(c_j)})\}$$

© Midia Reshadi

Problem Formulation

- such that:

$$\forall c_i \in C, \quad \text{map}(c_i) \in T$$

$$\forall c_i \neq c_j \in C, \quad \text{map}(c_i) \neq \text{map}(c_j)$$

$$\forall \text{link } l_k, B(l_k) \geq \sum_{\forall a_{i,j}} b(a_{i,j}) \times f(l_k, p_{\text{map}(c_i), \text{map}(c_j)})$$

- Where $B(l_k)$ is the bandwidth of link l_k :

$$f(l_k, p_{m,n}) = \begin{cases} 0 & : l_k \notin L(p_{m,n}) \\ 1 & : l_k \in L(p_{m,n}) \end{cases}$$

© Midia Reshadi

Significance of the Problem

- Scenario:

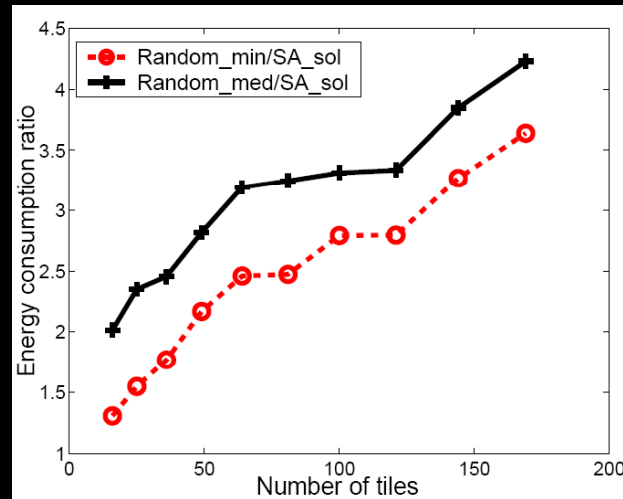
- The number of IPs: 3x3 to 13x13
- 3000 random generated mapping for each benchmark
- Optimizing method: Simulated Annealing (SA)

- Annotations:

- Random_min/SA_sol:
 - Minimum energy for random mapping/Simulated annealing
- Random_med/SA_sol:
 - Average energy for random mapping/Simulated annealing

© Midia Reshadi

Significance of the Problem



© Midia Reshadi

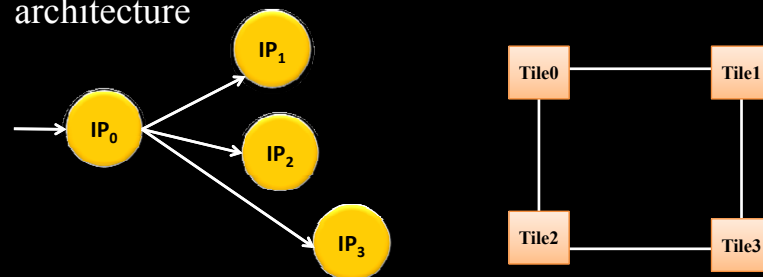
Significance of the Problem

- Unfortunately, the mapping problem is NP-hard problem
 - 4 x 4: tiles: 16! mappings
 - 10x10 tiles that are anticipated in five years or less
- This paper proposed the *branch-and-bound* algorithm

© Midia Reshadi

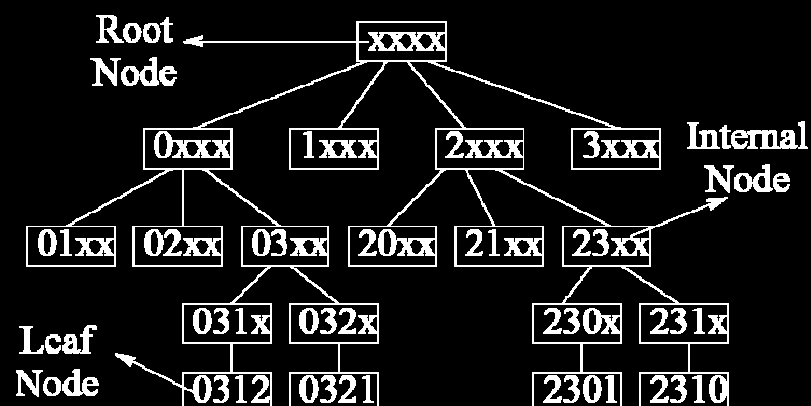
The Data Structure

- Walking through the search tree represents the whole searching space
- Example :
 - Mapping an application with 4 IPs onto a 2x2 tile architecture



© Midia Reshadi

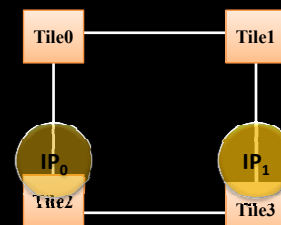
The Data Structure



© Midia Reshadi

The Data Structure

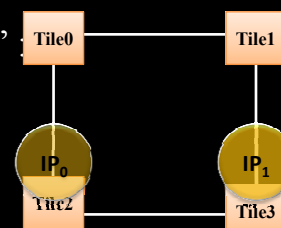
- **Example:**
 - “23xx” → IP0 and IP1 are mapped to Tile2 and Tile3 respectively, while IP2 and IP3 are still unmapped



© Midia Reshadi

The Data Structure

- **Definition 3**
- The **cost** of a node is the energy consumed by the communication among those IPs that have already been mapped.
- For instance
 - The cost of the node labeled “23xx”
 - $v(a_{0,1}) \cdot e(p_{2,3}) + v(a_{1,0}) \cdot e(p_{3,2})$



© Midia Reshadi

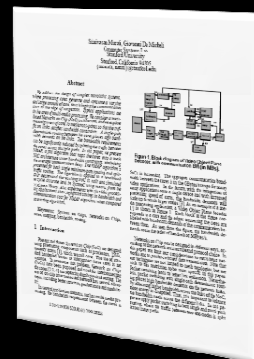
The Branch-and-Bound Algorithm

- **Branch:**
 - An unexpanded node is selected from the tree
 - The next unmapped IP is enumeratively assigned to the set of remaining unoccupied tiles
 - The corresponding new child nodes are generated.
- **Bound:**
 - Each of the newly generated child nodes are inspected to see if it is possible to generate the best leaf nodes later.

© Midia Reshadi

NMAP

- **Bandwidth-Constrained Mapping of Cores onto NoC Architectures**
 - **Authors:** Srinivasan Murali, Giovanni De Micheli



© Midia Reshadi

Mathematical Formulation of the Mapping Problem

- **Definition 1**
- **The core graph is a directed graph $\rightarrow G(V,E)$**
 - Each vertex $v_i \in V$
 - representing a core
 - Each edge $(v_i, v_j), e_{i,j} \in E$,
 - representing the communication between the cores v_i and v_j .
 - The weight of the edge $e_{i,j}$, denoted by $comm_{i,j}$,
 - represents the bandwidth of the communication from v_i to v_j .

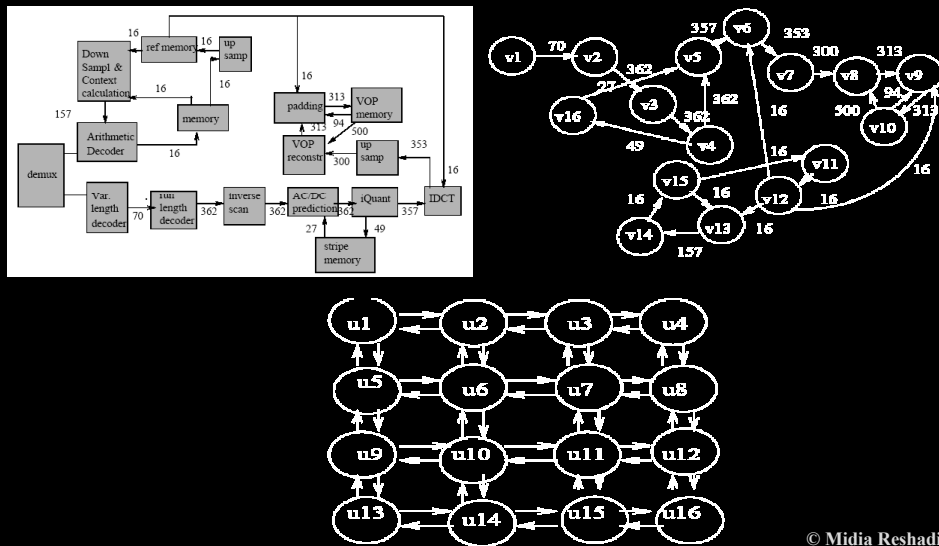
© Midia Reshadi

Mathematical Formulation of the Mapping Problem

- **Definition 2**
- **The NoC topology graph \rightarrow a directed graph**
 - $P(U, F)$ with each vertex $u_i \in U$
 - representing a node in the topology
 - Each edge (u_i, u_j) , denoted as $f_{i,j} \in F$
 - representing a direct communication between the vertices u_i and u_j .
 - The weight of the edge $f_{i,j}$, denoted by $bw_{i,j}$
 - represents the bandwidth available across the edge $f_{i,j}$.

© Midia Reshadi

Mathematical Formulation of the Mapping Problem



Mapping function

- Mapping function map:

$map : V \rightarrow U$, s.t. $map(v_i) = u_j, \forall v_i \in V, \exists u_j \in U$

$$|V| \leq |U|$$

- flow of single commodity (d^k), $k = 1, 2, \dots, |E|$.
 - The communication between each pair of cores (i.e. each edge $e_{i,j} \in E$)

Mapping function

- The value of d^k $vl(d^k)$
 - Represents the bandwidth of communication across the edge

$$D = \left\{ d^k : vl(d^k) = comm_{i,j}, k = 1, 2, \dots, |E|, \forall e_{i,j} \in E, \right. \\ \left. \text{with } source(d^k) = map(v_i), dest(d^k) = map(v_j) \right\}$$

- The bandwidth constraints

$$\sum_{k=1}^{|E|} x_{i,j}^k \leq bw_{i,j}, \forall i, j \in 1, 2, \dots, |U|$$

© Midia Reshadi