



Yield Enhancement by Robust Application-Specific Mapping on Networks-on-Chip



A. Dutta Choudhury
ALaRI - University of Lugano
Lugano, Switzerland
anirban.dutta.choudhury@lu.unisi.ch



G. Palermo, C. Silvano, V. Zaccaria
Politecnico di Milano, Milano (ITALY)
Dipartimento di Elettronica e Informazione
silvano@elet.polimi.it



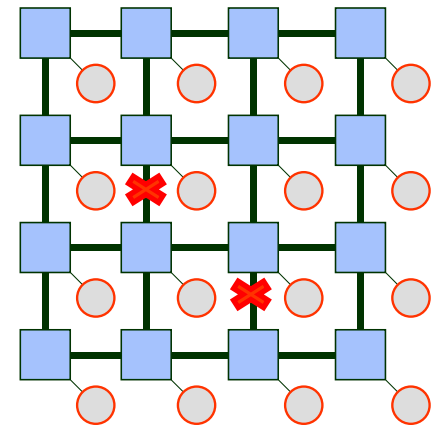
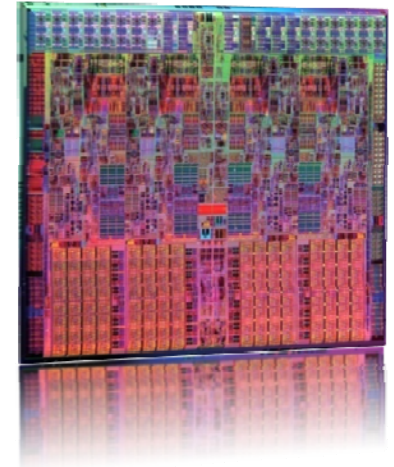
Outline

- Introduction and Motivations
- Methodology for Fault-Tolerant NoCs
 - Target Problem: Robust IP-to-NoC Mapping Problem
 - Problem Reformulation as Multi-Objective Mapping Problem
 - Proposed Multi-Objective Mapping Heuristic
- Experimental Results
- Conclusions



Introduction

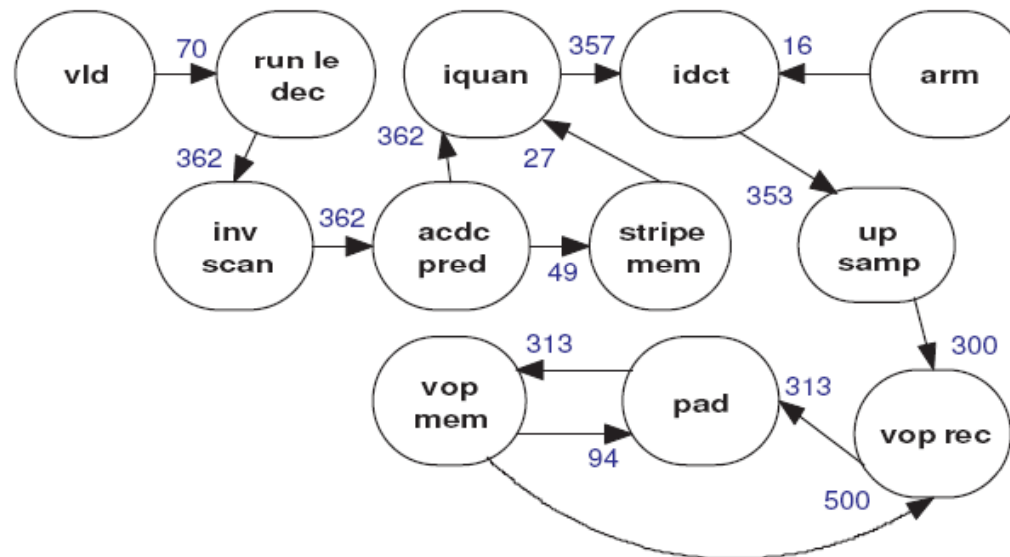
- Energy-aware fault-tolerant Network-on-Chip designs are crucial for the design of complex embedded systems.
- Current technological defect densities and production yield are the motivating factors to introduce design-for-manufacturability techniques during the high-level design of complex embedded systems based on NoCs.
- *Given the problem of mapping an IP graph to a NoC topology graph, our main goal is to increase the robustness of the NoC-based system with respect to link failures due to manufacturing defects.*





Definition of IP Graph

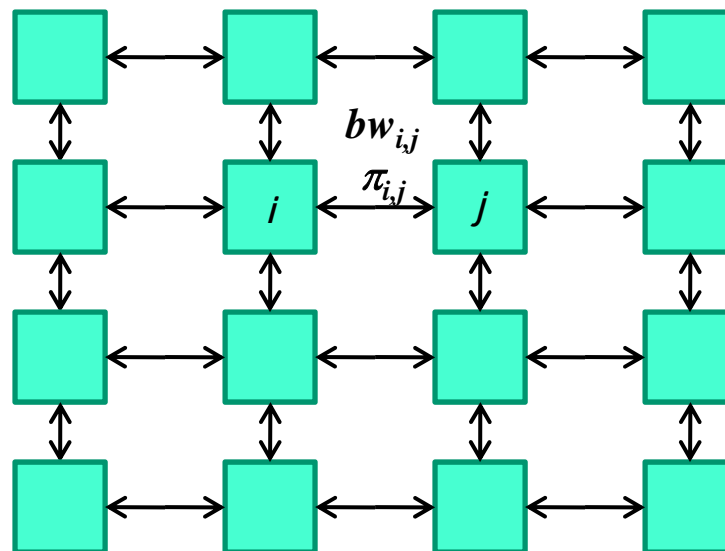
- The **IP Graph** is a direct graph $G(V,E)$ where V is the set of IPs of the target System-on-chip and E is the set of edges representing the communication between the IPs $v_i \in V$ and $v_j \in V$. The weight $w_{i,j}$ of the edge $e_{i,j} \in E$ represents the bandwidth of the direct communication from v_i to v_j
- An example: IP Graph of Video Object Plane Decoder (part of MPEG4 decoding algorithm)





Definition of NoC Topology Graph

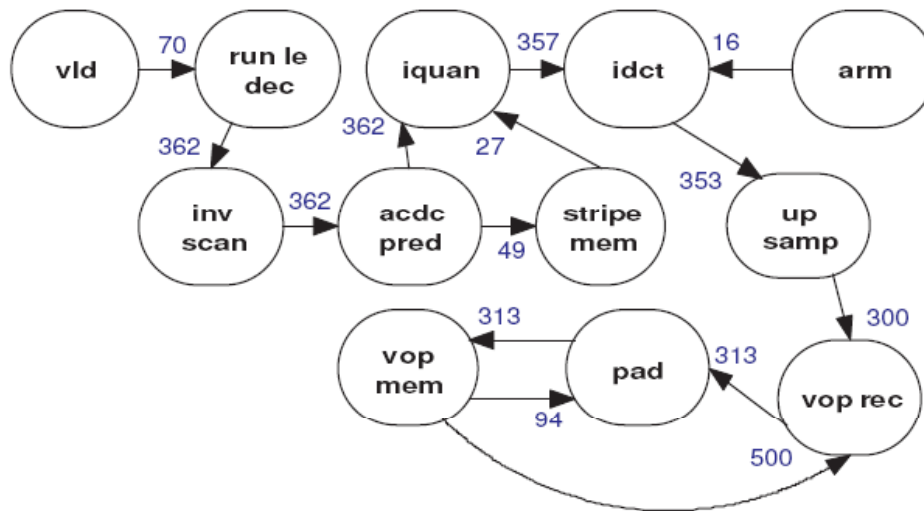
- The **NoC Topology Graph** is a direct graph $P(U, F)$ where U is the set of network nodes and F is the set of direct edges representing a link between the network nodes $u_i \in U$ and $u_j \in U$.
The weight of the edge $f_{i,j} \in F$ represents the bandwidth $bw_{i,j}$ available across $f_{i,j}$ and a probability $\pi_{i,j}$ of becoming unavailable due to manufacturing defects.
- An example: 2D-mesh with 16-nodes



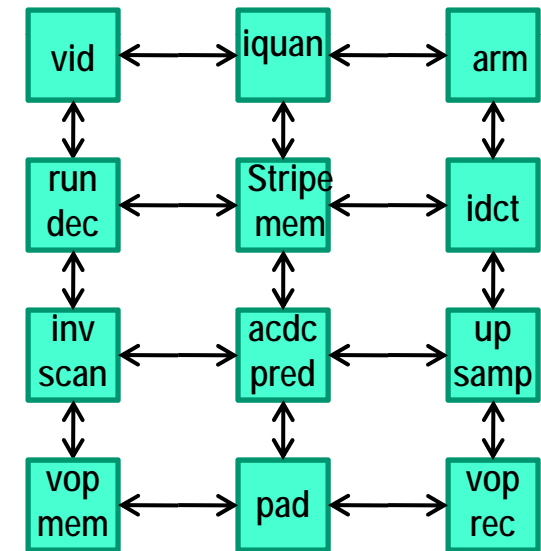
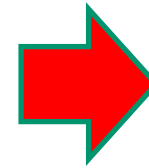


Mapping from IP Graph to NoC Topology Graph

- The **IP-to-Node mapping** function $M : V \rightarrow U$ is defined as the set IP-to-Node mappings (v_i, u_j) , representing the IP $v_i \in V$ mapped to the network node $u_j \in U$.
- An example: IP Graph of VOPD mapped to 12-node 2D-mesh



G



P

$$M \in \mathcal{M}(P, G)$$



Motivations

- A network fault makes some of the topology links unavailable and it generates a change on the fault-free NoC Topology Graph P^o , resulting into a new P^* .
- Assuming a dynamic routing policy, when the fault is detected, the new P^* can, in principle, still be working with sub-optimal power and delay figures
- However it could happen that the new P^* cannot enable any routing policy without incurring in a deadlocked/not-working network.

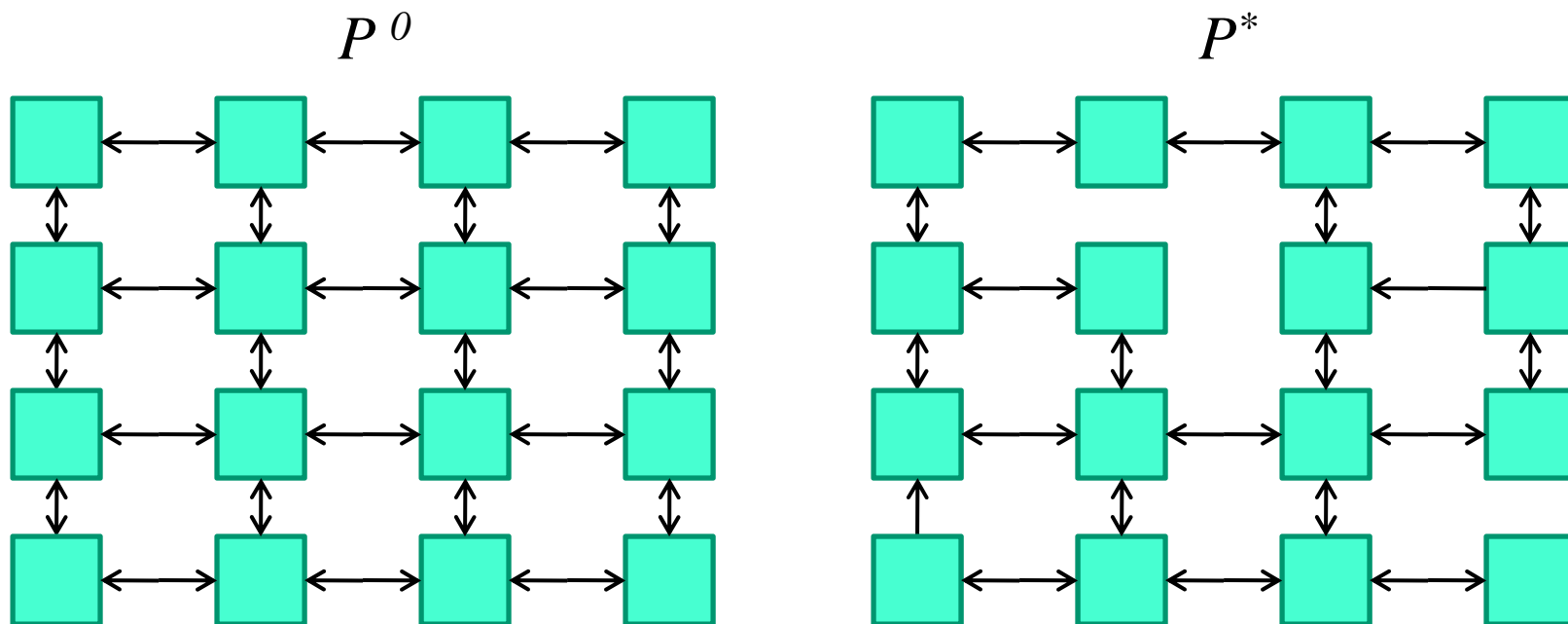
- *The motivating question behind this paper is:*

Is it possible to devise an IP mapping to NoC topology graph such that, given a network fault probability distribution, the power efficiency and the performance of the new topology P^ are optimized while minimizing the probability of incurring into a deadlocked/not-working network?*



Network Fault Model

- The actual network topology graph P^* can be different from the original fault-free topology graph P^0 because of some of the original links could be unavailable due to manufacturing defects. Each scenario P^* is associated with a probability of occurrence $\pi(P^*)$





Target Problem: Robust Optimization Problem

- To find an optimal IP-to-Node mapping M such that the estimated power consumption η and application execution delay τ of NoC-based system are minimized for all possible scenarios P^* derived from the original P^o . Moreover, we want to minimize the probability of deadlocked combinations (M, P^*) :

$$\min_{M \in \mathcal{M}(\mathcal{P}^o, \mathcal{G})} \left[\begin{array}{c} \eta(M, \mathcal{P}^*) \\ \tau(M, \mathcal{P}^*) \\ \text{prob}(\perp (M, \mathcal{P}^*)) \end{array} \right] \forall \mathcal{P}^* \in \Pi(\mathcal{P}^o)$$

where \perp is a predicate which is true whenever the combination of the actual mapping M and the current scenario P^* results into a network deadlock (*when using a minimum path routing, or when the communication requirements are not satisfied*).



Problem Reformulation as Multi-Objective Mapping

- Optimization of the mean μ and variance σ^2 of the target system metrics y (power η and delay τ) \Rightarrow definition of an aggregate quality measure Q_y of each system metric y given a set of N samples y_i corresponding to deadlock-free scenarios:

$$Q_y = \frac{1}{\left(\frac{1}{N} \sum_{i=1}^N y_i^2\right)}$$

- Maximization of the yield with respect to non-deadlocked combinations:

$$Y(M) = 1 - \text{prob}(\perp (M, \mathcal{P}^*) | \pi(\mathcal{P}^*))$$

- Multi-Objective optimization problem to find a Pareto set H of mapping solutions with respect to the objective functions: power quality, latency quality and yield

$$\max_{M \in \mathcal{M}(\mathcal{P}^o, \mathcal{G})} \begin{bmatrix} Q_\eta(M) \\ Q_\tau(M) \\ Y(M) \end{bmatrix}$$



Proposed Robust Mapping Algorithm

Require: $\mathcal{G}(V, E), \mathcal{P}^o(U, F), R, \rho$

Ensure: $|U| \geq |V|$

1: $H = \{ \}$

2: $cov = \infty$

3: **while** $cov > 0$ **do**

4: $H_R =$ generate R random initial mappings from $\mathcal{M}(\mathcal{P}^o, \mathcal{G})$

5: $cov = \chi(\psi(H \cup H_R), H)$

6: $H = \psi(H \cup H_R)$

7: **end while**

8: $C = \{C_{min}, C_{avg}, C_{max}\} =$ k-means clustering of H into 3 sets, considering $Y(M), \forall M \in H$

9: **for** $\forall C_i \in C$ **do**

10: $k_0 = \arg \max \Phi(M), \forall M \in C_i$

11: $N(k_0) = \nu(k_0, \rho)$

12: $k_1 = \arg \max \Phi(M), \forall M \in N(k_0)$

13: $f_0 = \Phi(k_0), f_1 = \Phi(k_1)$

14: **while** $f_1 > f_0$ **do**

15: $k_0 = k_1, f_0 = f_1$

16: $N(k_0) = \nu(k_0, \rho).$

17: $k_1 = \arg \max \Phi(M), \forall M \in N(k_0)$

18: $f_1 = \Phi(k_1)$

19: **end while**

20: $h_i = k_0$

21: **end for**

22: **return** $\arg \max \Phi(M), \forall M \in (\{h_{min}, h_{avg}, h_{max}\})$

Iterative Pareto filtering of a set of random mappings

$H =$ Pareto front of mapping solutions

K-means clustering to partition the Pareto-set H in 3 yield classes

Steepest climbing neighborhood search to optimize the geometric average Φ of objective functions

Return best mapping w.r.t. Φ



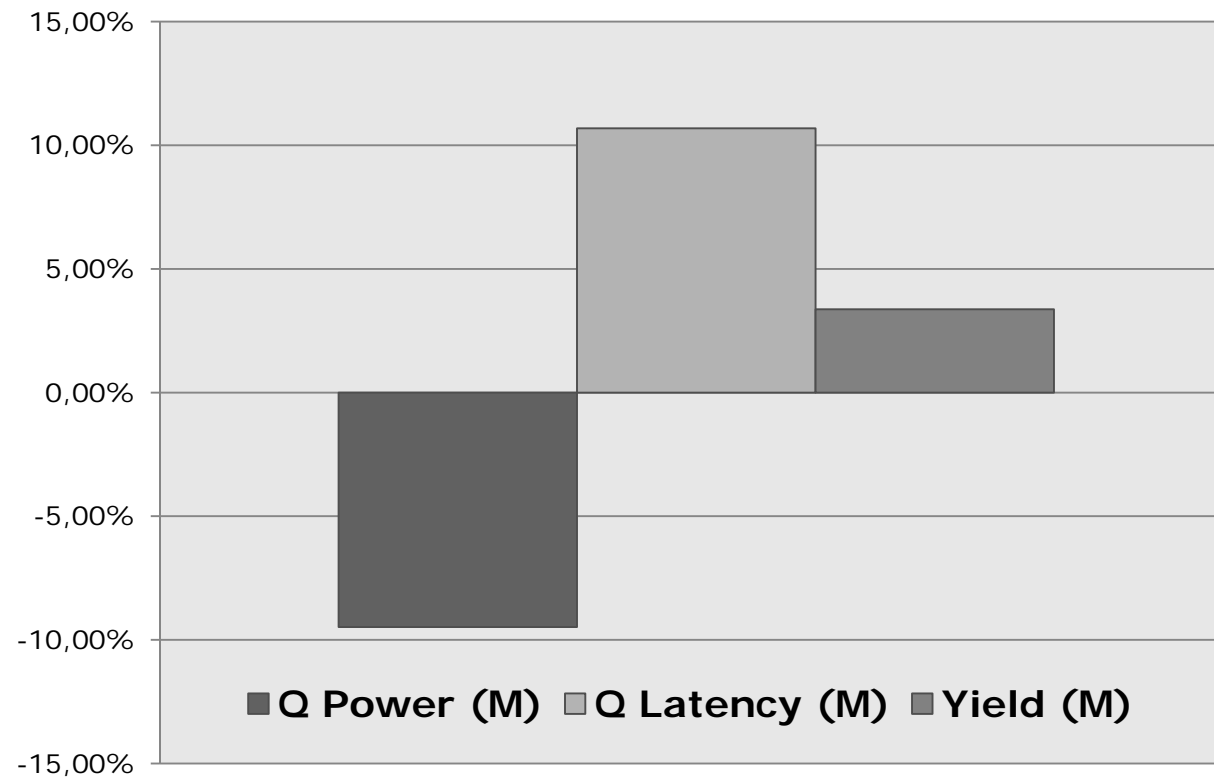
Experimental Results

- Mapping of IP-Graph of VODP application to a (4 x 3) mesh NoC topology.
 - Power model characterized with 90 nm STMicro library
 - Latency model derived from PIRATE NoC simulation framework



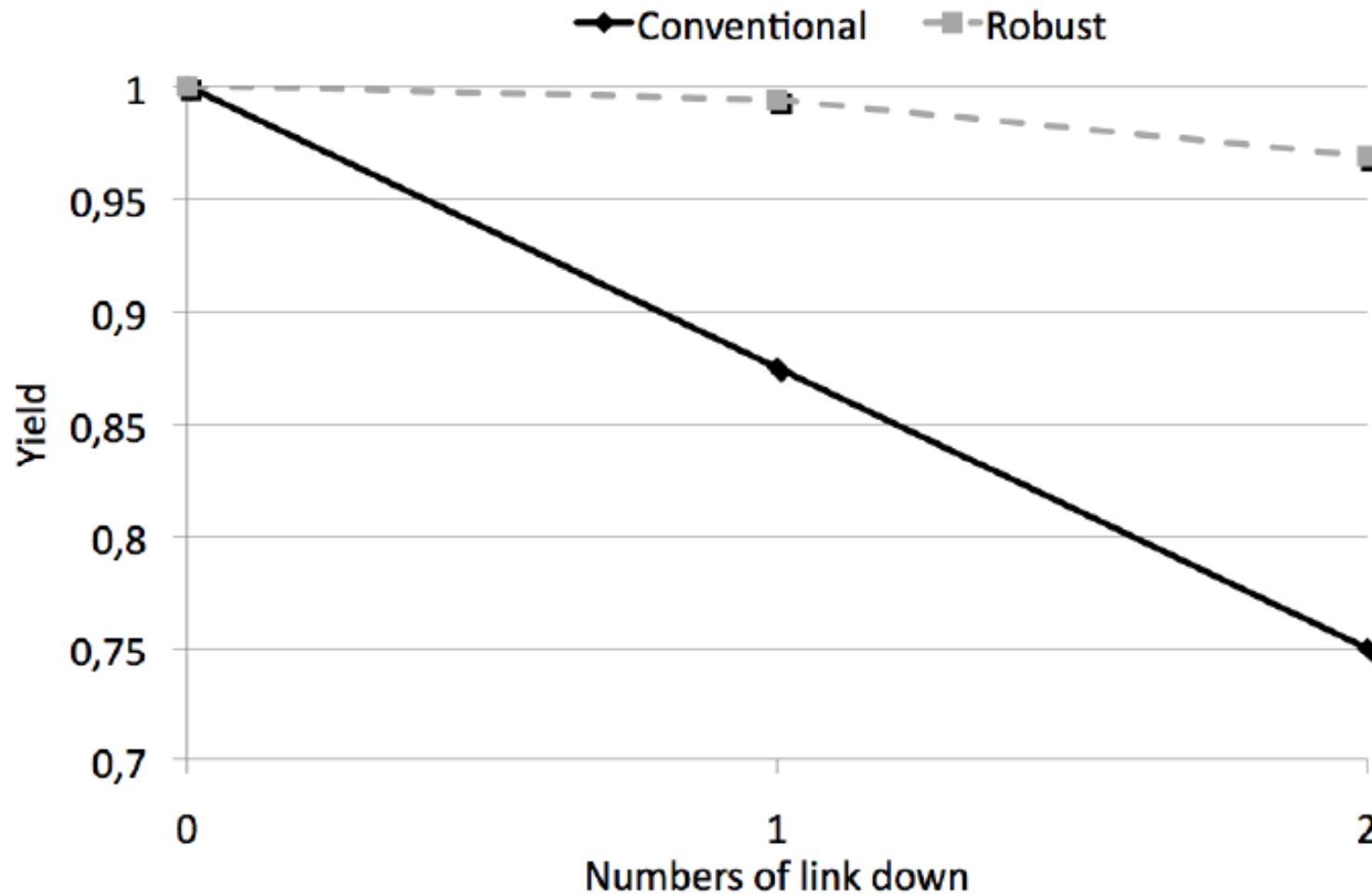
Comparison: Robust vs Conventional Mapping

- Comparison of the proposed robust approach with respect to a conventional mapping (SUNMAP) minimizing the aggregate NoC bandwidth considering non-faulty topology



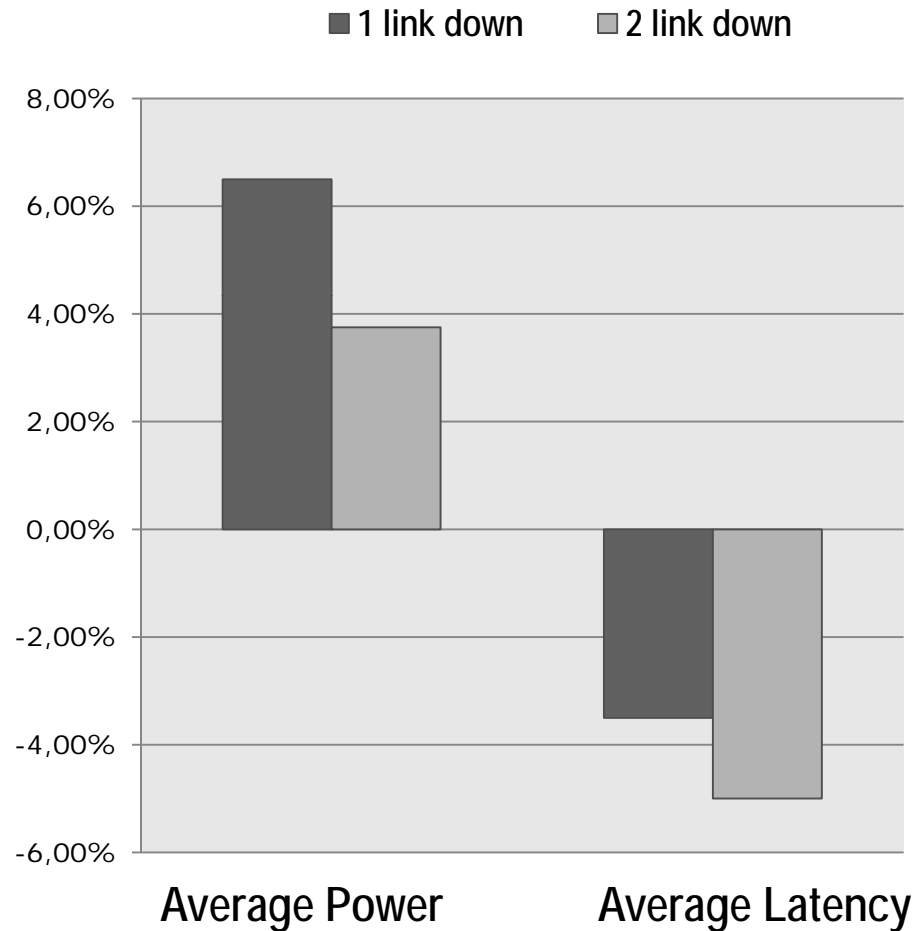


Yield of robust and conventional approaches for 1 link down and 2 links down





Percentage variation of robust vs conventional approach in terms of average power and latency





Conclusions

- A robust application-specific methodology has been proposed for identifying optimal IP-to-NoC mappings in Chip Multi Processor architectures.
- The proposed robust mapping approach increases the probability to derive a feasible routing even in the case of faulty links.
- Future research is directed towards the analysis and optimization of the influence of more complex probability distributions to the overall system metrics.
- This work is part of the ICT-FP7 EU project MULTICUBE on Automatic Design Space Exploration for CMPs.

www.multicube.eu

