

# POWER CONSUMPTION OPTIMIZATION AND DELAY BASED ON ANT COLONY ALGORITHM IN NETWORK-ON-CHIP

T. He<sup>1\*</sup> – Y. Guo<sup>1</sup>

<sup>1</sup>National Digital Switching System Engineering & Technological Research Center, Zhengzhou 450002, Henan, China

---

## ARTICLE INFO

### Article history:

Received: 11.9.2013.

Received in revised form: 13.10.2013.

Accepted: 13.10.2013.

### Keywords:

Network-on-chip

Optimization

Power consumption

Delay on chip

## Abstract:

*With a further increase of the number of on-chip devices, the bus structure has not met the requirements. In order to make better communication between each part, the chip designers need to explore a new NoC structure to solve the interconnection of an on-chip device. For the purpose of improving the performance of a network-on-chip without a significant increase in power consumption, the paper proposes a network-on-chip that selects NoC (Network-On-Chip) platform with 2-dimension mesh as the carrier and incorporates communication power consumption and delay into a unified cost function. The paper uses ant colony optimization for the realization of NoC map facing power consumption and delay potential. The experiment indicates that in comparison with a random map, single objective optimization can separately account for (30%~47%) and (20%~39%) of communication power consumption and execution time, and joint objective optimization can further excavate the potential of time dimension in a mapping scheme dominated by the power.*

---

## 1 Introduction

With an increase in the number of CPU (Central Processing Unit), various complicated buses with a topological structure and communication means emerge to support multi-CPU system [1]. One of the branches enters the field of the network structure and routing communication, and develops a new NoC architecture [2]. Compared with SoC (System on Chip), NoC is a system on chip technology with higher level and larger scale, and a network system on chip [3-5]. The network structure based on routing technology solves the problem of series fundamentally. The communication between each node is not limited to a path, and so the multi-task and multi-progress communication enable parallel operations in time shaft [6-8]. The parallel

computation of nodes and parallel communication of nodes enable a true parallel operation. Instead of calculation, communication has now emphasized a NoC design. . However, the performance and power consumption of NoC have become the bottleneck of the entire system in further optimization.

Performance parameters of NoC include delay, jitter, throughput, etc. The changes of network status may cause the differences among each packet transmission delay and result in jitter. Increasing performance of NoC is required with real-time services emerging. For example, not only real-time but also the smaller jitter is required for audio and video services to playback smoothly. The researchers mainly aim at routing unit structure and algorithms in performance optimization. [9] proposed a low-latency low power fault-tolerant routing unit

---

\* Corresponding author. Tel.: 08613700882262 ; fax:086037165903522  
E-mail address: taohe163@163.com.

structure, named RoCo. [10] proposed DAMQ (Dynamically Allocated Multiple Queue) to achieve a virtual channel mechanism that allows different direction port sharing the same cache. [11] then proposed a "Neighbors-on-Path" selection strategy in a routing process to reduce routing delay. [12] proposed a dynamic routing algorithm which can avoid deadlock and livelock in routing process. In order to reduce network delay, [13] and [14] propose the dynamic routing algorithms with fault-tolerant routing. As the above documents do not consider congestion waiting time, they are not very reasonable. The power consumption is another important constraint of NoC design because the power consumption of communication accounts for a relatively large part of the entire system consumption. e.g., in MIT Raw, the communication interconnection network consumed 36% of the total power consumption in the entire system [15]. The communication structure of the Alpha processor consumed 20% of the power consumption of the whole chip [16]. The GALS (Global Asynchronous Local Synchronous) mechanism used in network-on-chip can effectively reduce power consumption. In [17], the paper analyzes performance and power consumption of the system using GALS, and the experimental results show that the loss of a system performance is less than 1% of the original, but the power consumption is reduced by 40% in 0.18  $\mu$  m process. According to the change of processor utilization, DVS (Dynamic Voltage Scaling), proposed in [18], reduces both chip frequency and voltage to decrease power consumption. In the paper, the author used an existing network status to predict the future load, connection frequency and voltage dynamically. In [19] the research scientists proposed a variable frequency link to reduce power consumption by adjusting the link voltage. However, above literatures do not consider the power consumption and delay simultaneously.

From the above analysis it is quite evident that this article proposes a network-on-chip that selects NoC platform with 2-dimension mesh as the carrier. The paper incorporates communication delay and power consumption into a unified cost function and uses ant colony optimization to realize NoC map facing power consumption and delay potential. This method takes into account the objective function of the power consumption and execution time and the delay is optimized indirectly by optimizing link load distribution, which avoids the problem of accurate modeling of NoC waiting delay.

## 1.1 Power consumption and delay model

This NoC platform is a parameterized and two-dimensional grid structured model. NoC consists of a router, links ( $4 \times 4$  scale, as shown in Fig. 1). The router is a core component of networks and constitutes the switch node.

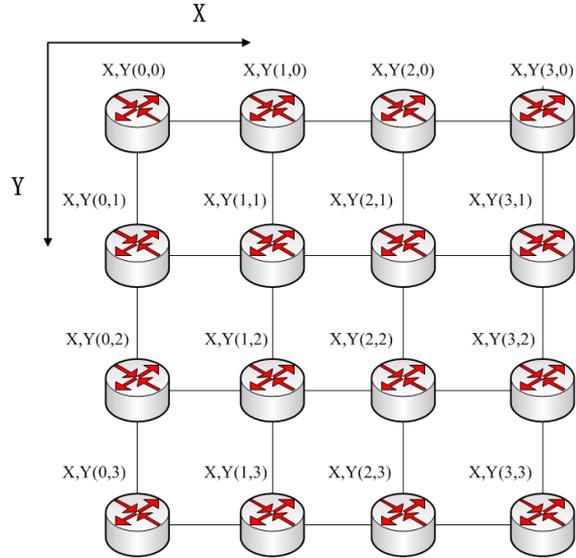


Figure 1. NoC platform structure.

The platform adopts a wormhole switching mode, input channel cycle FIFO (First in First out) caching strategy, rotary priority arbitration mechanism and static XY routing algorithm. Static routing can ensure the prediction of communication path so as to develop a new research resource on NoC map.

## 1.2 Power consumption model

In NoC model, the average power consumption of node  $r_i$  transmitting to node  $r_j$  for single bit data is shown in Equation (1).

$$E_{bit}^{i,j} = (h_{i,j} + 1) \times E_{Sbit} + h_{i,j} \times E_{Lbit} \quad (1)$$

In the equation,  $E_{Sbit}$  and  $E_{Lbit}$  respectively represents the power consumption on switch and links.  $h_{i,j}$  represents the number from node  $r_i$  to node  $r_j$ , that is, the Manhattan distance between two nodes. Equation (1) is a linear equation in which  $E_{Sbit}$  and  $E_{Lbit}$  are constant coefficients, and we can see that the maximal power consumption is equivalent to the maximal  $h_{i,j}$ . Therefore, the Manhattan distance between two nodes can be regarded as their communication power consumption index.

### 1.3 Delay model

In two-dimensional NoC grid adopting wormhole switching and static XY routing, the delay  $T_{i,j}$  of transmitting a data packet from the source node  $r_i$  to the goal node  $r_j$  is defined as a time slot, the packet begins in/at  $r_i$ , and packet ends when reaching  $r_j$ .  $T_{i,j}$  model is shown in Equation (2).

$$T_{i,j} = (T_b + T_w) \times h_{i,j} + T_b \times (B - 1) \quad (2)$$

In Equation (2),  $T_b$  is the time of a frame data (flit) passing through a switch and a link when there is no congestion,  $T_w$  is the average waiting time of packet head in switching node when there is congestion,  $h_{i,j}$  is the Manhattan distance between  $r_i$  and  $r_j$ ,  $B$  is the number of data frame included in data packet and finally  $T_b$  and  $B$  are constant coefficients. In Equation (2), the first item is the time of packet needed to reach the destination, and the second item is the time of the subsequent frame needed to reach the destination. Both of them are constant terms. Therefore, the optimizing space of  $T_{i,j}$  lies in the first item and depends on  $T_w$  and  $h_{i,j}$ . Thus,  $h_{i,j}$  is defined by the corresponding position relationship of  $r_i$  and  $r_j$ , and subsequently  $T_w$  relies on the congestion status of the network.

Transmission delay of NoC has a dynamic feature. It is more complex than the delay of shared bus, and therefore it could not be accurately predicted. The map strategies of estimating delay by ignoring or setting  $T_w$  artificially are not accurate. The difficulty and complexity of precise modeling  $T_w$  are taken into consideration, so to consider these issues, we turn to the method of optimizing the key factors influencing  $T_w$  for indirect optimization of  $T_w$ . Only a rational position of key factors can guarantee the optimization effect, and it may be better than the method for optimizing  $T_w$  directly. The most effective way of reducing  $T_w$  is to remit the network congestion, and consequently the key to remitting the congestion is to balance the link load. Based on the ideas, the paper selects link load balancing to remit the congestion and optimize  $T_w$ . The transmission delay is thus reduced, and the objective to reduce the time of implementing tasks achieved.

### 1.4 Objective function

An optimal object of NoC is positioned in communication power consumption and execution

time. According to the ideas in chapter 2.2 and 2.3, the definition of power consumption index is shown in Equation (3).

$$E(C) = \sum_{i=1}^N \sum_{j=1}^N w_{i,j} \times h_{i,j} \quad (3)$$

In the equation,  $w_{i,j}$  is the communication traffic from  $r_i$  to  $r_j$ ,  $N$  is the number of nodes. It is clear that the optimum objective of communication power consumption is the sum of the weighted Manhattan distance between each minimized nodes, minimize ( $E(C)$ ).

As shown in the first item of Equation (2), the optimizable transmission delay item can essentially be divided into two parts, latency time relating to congestion and fixed transmission time relating to distance. The fixed transmission time is decided by  $h_{i,j}$ , and latency time is decided by the product of  $T_w$  and  $h_{i,j}$ .  $h_{i,j}$  has been considered in optimizing communication power consumption, whereas  $T_w$  is not incorporated into the optimization scope.  $T_w$  has a great influence on the delay, and the optimization space is considerable. According to the idea expounded in chapter 2.2, the goal of optimizing delay has been achieved by balancing link load. The variance of link load is the imbalance index of link load on optimization, and the definition is shown in Equation (4).

$$VAR(L) = \sum_{i=1}^M \left[ Load(l_i) - Load(l)_{avg} \right]^2 / M \quad (4)$$

In the equation,  $l_i$  is the  $i^{th}$  link of NoC,  $M$  is the total number of links,  $Load(l_i)$  is the load amount of link  $l_i$ ,  $Load(l)_{avg}$  is the average link load amount, the variance represents the dispersion degree of the distribution, the larger the variance, the more uneven the distribution. So the optimization goal of communication delay is transferred into the minimized link load variance, minimize ( $VAR(L)$ ). All these facts show that optimizing power consumption needs to optimize  $E(C)$ , and optimizing delay needs to optimize  $VAR(L)$  and  $E(C)$  simultaneously, in which  $VAR(L)$  is the dominance. The definition of unified cost function is as follows:

$$cos = \lambda \times E(C) + (1 - \lambda) \times VAR(L) \quad (5)$$

In the equation,  $\lambda$  is a proportionality coefficient that is used to adjust the proportion of communication power consumption and delay in cost function, and the value range is  $0 \sim 1$ . So the

optimization objective of NoC map method in the paper is minimize ( $cos$ ). When  $\lambda = 1$ , it represents the optimization of communication, and when  $\lambda = 0$ , it represents the optimization of communication delay. When  $\lambda$  takes the values between 0 and 1, the optimization embodies the compromise of communication power consumption and delay.

## 2 Mapping approach for power consumption and delay

Since V. Maniezzo, A. Colorini and M. Dorigo applied Ant System algorithm for the first time to solve quadratic assignment problem, many researchers have proposed improved algorithms, and the relevant introduction can refer to the literatures. NoC map belongs to quadratic assignment problem. The paper adopts AS algorithm and introduces pheromone bound in MMAS algorithm to inhibit search sinking in local optimal solution untimely.

### 2.1 Objective function

The definition and the settings for basic parameters are shown in Table 1.

Table 1. Parameter definition and initialization

Parameters	Definition	Values
M	Number of ants in ant colony	N (NoC scale)
$\alpha$	Degree of importance of information in path	1
$\beta$	Degree of importance for directive factor	3,5
$\rho$	Degree of information reducing with time lapse	0,7
$\tau(0)$	Initial value of pheromone intensity	1

The parameters flows are adjusted according to the task characteristics.  $\eta_{i,j}$  is called prior knowledge or visibility, which means inspiration information of IP core assigning to resource node  $r_j$ , and it determines the probability distribution with pheromone strength. Inspiration information has a great influence on convergence speed of algorithm, which is considered deliberately. Two functions are defined here:

$$center(i) = \sum_{j=1}^N h_{i,j} \tag{6}$$

$$vip(i) = \sum_{j=1}^N w_{i,j} + \sum_{j=1}^n w_{j,j} \tag{7}$$

$center(i)$  represents the central degree of  $r_i$  in NoC and means communication capability; the less the value, the stronger the communication capability.  $vip(i)$  means the degree of importance for  $v_i$  in the application character graph; the greater the value, the more important it is. It is obvious that an optimal map always assigns the most important IP core to the resource node with the strongest communication capability. Therefore, the definition of inspiration information is as shown in Equation (8).

$$\eta_{i,j} = vip(j)center(i) \tag{8}$$

$\eta_{i,j}$  embodies the reasonableness of IP core assigned to the resource node  $r_i$ .

### 2.2 Construction of solution

$p_{i,j}^k(t)$  means the probability that the ant k assigns  $v_j$  to  $r_i$  in the cycle :

$$p_{i,j}^k(t) = \begin{cases} [\tau_{i,j}(t)]^\alpha \times [\eta_{i,j}]^\beta, & j \notin tabu\ k \\ 0, & j \in tabu\ k \end{cases} \tag{9}$$

$tabu\ k$  ( $k=1, 2, \dots, M$ ) is used to record the assigned IP core of the ant k and is called tabu list. Also,  $\tau_{i,j}(t)$  means the pheromone intensity of the path from  $v_j$  transmitted to  $r_i$  on the cycle of the t time. The solution procedure is as follows. A core is selected from the optional set according to the probability  $p_{i,j}^k(t)$  assigned to  $r_1$ , and the core is added to tabu. And then an unassigned core is selected according to the probability  $p_{2,j}^k(t)$  for assigning to  $r_2$ , and the core is added to tabu. The step has been implemented N times until all cores have been assigned to the corresponding resources, and tabu has been full. In a cycle, every ant only implements the above process for one time and M solutions for mapping problem can be obtained.

### 2.3 Pheromone update

With the advancement of evolution process, the left pheromone vanishes gradually. The parameter  $\rho$  is used to represent the persistence degree of

pheromone ( $0 < \rho < 1$ ). When all ants have finished one cycle, the pheromone should be updated. And information content of each assignment path should be updated according to Equation (10).

$$\tau_{i,j}(t+1) = \rho \times \tau_{i,j}(t) + \Delta\tau_{i,j} \quad (10)$$

In the equation,  $\Delta\tau_{i,j}$  means information gain and is represented as,

$$\Delta\tau_{i,j} = \sum_{k=1}^M \Delta\tau_{i,j}^k \quad (11)$$

$\Delta\tau_{i,j}^k$  means the left information content of the ant k in assignment path ( $v_j \rightarrow r_i$ ) in this cycle, and the calculation equation is,

$$\Delta\tau_{i,j}^k = \begin{cases} \frac{1}{\cos(k)} & \text{map}(k) \text{ includes } (v_j \rightarrow r_i) \\ 0 & \text{else} \end{cases} \quad (12)$$

In the equation,  $\cos(k)$  is the cost of the ant k completing distribution program, which can be referred to the destination in chapter 2.3. It is clear that the optimum solution has the minimum cost, and therefore it has the greatest contribution to information update. In order to avoid search stagnation, the pheromone intensity is provided with the bound  $T_{\max}$  and  $T_{\min} = T_{\max} / 5$ . And the calculation equation is Equation (13) and Equation (14).

$$T_{\max} = \frac{1}{1-\rho} \times \frac{1}{\text{BestSolution}} \quad (13)$$

$$T_{\min} = T_{\max} / 5 \quad (14)$$

In the equation, Best Solution is the cost of the optimal solution for the present moment.

### 3 Algorithm simulation

#### 3.1 Realization of mapping algorithm

Program source code of mapping algorithm is written by C++, and compilation and simulation are completed under Microsoft Visual C++ 6.0. The paper adopts the task graph needed for the experiment and is generated randomly from self-compiled program. The experiment adopts three different scales of the application characterization graph. The property of each application characterization graph is shown in Table 2. Additionally, the NoC target platform corresponding to three applications is respectively two dimensional network structure with the rule of  $3 \times 3$ ,  $4 \times 4$  and  $5 \times 5$ .

Table 2. Property of application characterization

	Tasks	Edges	Communication quantity
Application 1	9	64	505
Application 2	16	176	751
Application 3	25	457	1804

The paper uses the mapping scheme with the order ( $v_j \rightarrow r_i$ ) as reference scheme (Ref). As the application characterization graph is generated randomly, the reference is equivalent to the random map. In the experiment, the applications of three scales for mapping are optimized according to the optimal object with different emphases ( $\lambda = 1, 0.5, 0$ ).

Communication power consumption and link load cost variance of each mapping scheme is respectively shown in Fig. 2 and Fig. 3.

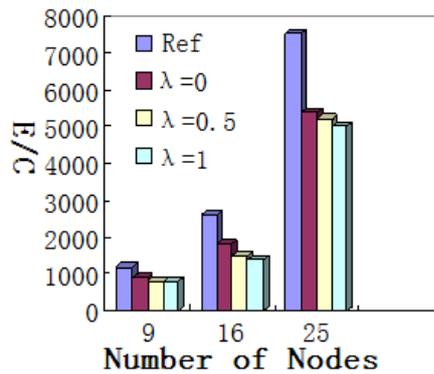


Figure 2. Communication power consumption.

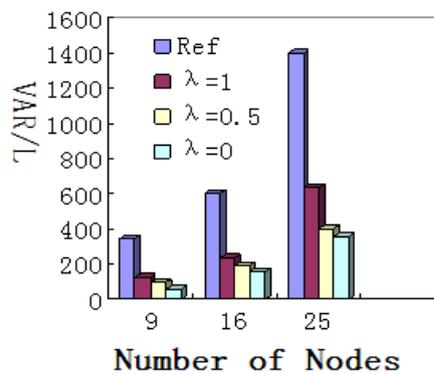


Figure 3. Results of link load variance.

From Fig. 2 and Fig. 3, we can see that compared with a reference scheme,  $\lambda = 1$  mapping scheme  $E(C)$  and  $VAR(L)$  is respectively reduced by 30~47%

and 55~65%,  $\lambda = 0.5$  mapping scheme is respectively reduced by 30~45% and 70~77%,  $\lambda = 0$  mapping scheme is respectively reduced by 24~30% and 75-83%. The experiment results are consistent with the algorithm goal. The optimization of a single goal can minimize the corresponding indexes ( $E(C)$  or  $VAR(L)$ ) in a great degree, and joint objective optimization can consider two indexes, which can make a good compromise.

### 3.2 Execution time simulation

The above experiment results show that the introduction of link load variance into objective function can evidently balance the distribution of link load. But it needs further verification if link load balancing can improve the network congestion, reduce delay and reduce application execution time. Therefore, we convert the communication content between IP cores into the data packet with the corresponding quantity (the communication content of  $v_j \rightarrow r_i$  is converted into the data packet of  $w_{i,j}$ ). IP core is used to act as a sending and receiving unit for making real-time simulation on communication process of a specific mapping scheme. Other parameters of a NoC platform are that the buffering depth is 8, the data packet length is 3 frames, and the width of each frame is 35 bits. The execution time of each mapping scheme is shown in Fig. 4.

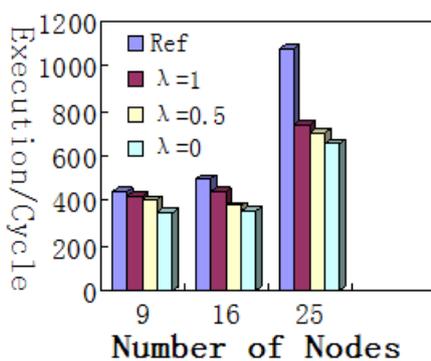


Figure 4. Results of the execution time.

We can know from Fig. 4 that compared with a reference scheme, the execution time of the mapping schemes for  $\lambda = 1, 0.5$  and  $0$  is respectively reduced by 4~32%, 7~35% and 20~39%. Besides, the experimental results of execution time are consistent with the link load variance, which corroborates the analysis results in Chapter 2.3.

### Conclusions

The paper has presented a brief introduction into the basic and relevant knowledge about network on chip technology. Employing research and analysis of the existing communication power consumption and transmission delay models, the paper has incorporated communication power consumption and delay into a unified cost function and used ant colony optimization to realize a NoC map facing power consumption and delay potential. The experiment has indicated that single objective optimization can separately get a certain degree of improvement compared with the random map. Through appropriate parameters adjustment, a balance between power consumption and delay can be achieved, and also the overall performance of the system in the joint objective optimization enhanced. This paper has provided a referable design solution for LowPower and High Performance NoC.

### 4 Acknowledgement

Main project of National 863 Program "2009AA012201", The key scientific and technological projects of Science and Technology Commission of Shanghai Municipality "08dz501600".

### References

- [1] Greenberg, R. I, Oh. H. C.: *Universal wormhole routing*, IEEE Transactions on Parallel and Distributed Systems, 8 (1997) 3, 254-262.
- [2] Ni, L. M., McKinley, P. K.: *A survey of wormhole routing techniques in direct networks*, IEEE Tran.-on-Computers, 23 (1993) 2, 62-76.
- [3] Tingqiang S., Chuanlai L., Sikun L.: *Research on IP core application technology in SoC design*, Journal of Qingdao Institute of Chemical, 23 (2003) 3, 260-263.
- [4] Chenyang G., Weipu X., Fei S.: *Study on IP multiplex technique*, Microelectronics Journal, 32 (2002) 4, 257-260.
- [5] Kaiyu W.: *Improved DyAD Algorithm for Network-on-Chip*, Journal of Convergence Information Technology, 7 (2012) 8, 62- 72.
- [6] Yi L., Gang L., Yintang Y., Zijin L.: *A 9Gb/s/ch Low-Swing Transceiver for Interconnection between NoC routers*, Advances in Information Sciences and Service Sciences, 4 (2012) 7, 12 - 22.
- [7] Dally, W. J, Towles, B.: *Route Packets Not Wires: On-Chip Interconnection Networks*, Proc of DAC, 2001, 684 -689.

- [8] Shashi Kumar et al: *A network on chip architecture and design methodology*, Proceedings of IEEE Computer Society Annual Symposium, 2002, 105 -112.
- [9] Kim, J., Nicopoulos, C., Park, D.: *A Gracefully Degrading and Energy-Efficient Modular Router Architecture for On-Chip Networks*, Proceedings of the 33rd International Symposium on Computer Architecture, 2006, 4 - 15
- [10] Jin L., Jose. G., Delgado F.: *A Shared Self-Compacting Buffer for Network-on-Chip Systems*, MWSCAS'06. 49th IEEE International Midwest Symposium, 2 (2006), 22-30.
- [11] Ascia, G., Catania, V., Palesi, M., Patti, D.: *Neighbors-on-Path: A new selection strategy for On-Chip Networks*, Proceedings of the 2006 IEEE/ACM/IFIP Workshop, Seoul Korea, 2006, 79-84.
- [12] Sabry, M.M., El-Kharashi, M.W., Bedor, H.S.: *A new Dynamic Routing Algorithm for Networks-on-Chips*, IEEE Pacific Rim Conference on Communications, Computers and Signal Processing, VictoriaBC Canada, 2007, 22-24.
- [13] Chen, K. H., Chiu, G. M.: *Fault-Tolerant Routing Algorithm for Meshes without Using Virtual Channels*, Journal of Information Science and Engineering, 14 (1998) 4, 65-783.
- [14] Schonwald, T. Z., Jochen B.: *Full Adaptive Fault-Tolerant Routing Algorithm for Network-on-Chip Architectures*, 10th Euro micro Conference on Digital System Design Architectures, Methods and Tools, Pisa Italy, 2007, 527-34.
- [15] Hangsheng, W., Li-Shiuan, P., Malik: *Power-driven design of router microarchitectures in NoC*, In Proceedings of the 36th Annual IEEE/ACM International Symposium on Microarchitecture, Washington USA, 2003, 105-116.
- [16] Zhiyi Y., Bevan M. B.: *Performance and Power Analysis of Globally Asynchronous Locally Synchronous Multi-Processor Systems*, Proceedings of the 2006 Emerging VLSI Technologies and Architectures, Karlsruhe Germany, 2006, 343-349.
- [17] Li S., Li-Shiuan P., Niraj K. J.: *Dynamic Voltage Scaling with Links for Power Optimization of Interconnection Networks*, Proceedings of the 9th International Symposium on High-Performance Computer Architecture, California USA, 2003, 91-102.
- [18] Kim, J., Horowitz, M.: *Adaptive supply serial links with sub-1V operation and per-pin clock recovery*, Digest of Technical Papers. ISSCC. 2002 IEEE International, San Francisco USA, 8 (2002), 216-480.

