# DESIGN AND OPTIMIZATION OF HETEROGENEOUS SYSTEM-ON-CHIP ARCHITECTURE USING SELF-ADAPTIVE NEURAL NETWORK ACCELERATOR

**B. Ebenezer Abishek[1], C. Sharanya[2], S. Gopalakrishnan[3] and J. Jency Rubia[4]**

[1]*Department of Electronics and Communication Engineering, Vel Tech Multi Tech Dr.Rangarajan Dr.Sakunthala Engineering College, India*
[2]*Department of Electronics and Communication Engineering, Vels Institute of Science, Technology and Advanced Studies, India*
[3]*Department of Medical Electronics Engineering, Sengunthar Engineering College, India*
[4]*Department of Electronics and Communication Engineering, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology, India*

**Abstract**

*In response to the escalating demand for energy-efficient and high-performance computing, this research explores the design and optimization of a heterogeneous System-on-Chip (SoC) architecture employing a self-adaptive neural network accelerator. Addressing the current limitations in heterogeneous SoC designs, we identify the need for dynamic adaptation to varying workloads. Our proposed methodology integrates a self-adaptive neural network accelerator that autonomously adjusts its architecture based on real-time workload characteristics. Through extensive simulations, we demonstrate significant improvements in both energy efficiency and performance compared to traditional static architectures. This research bridges the existing gap in adaptive computing, providing a promising avenue for future energy-efficient heterogeneous SoC designs.*

*Keywords:*

*Heterogeneous System-on-Chip, Neural Network Accelerator, Self-Adaptive Architecture, Energy Efficiency*

## 1. INTRODUCTION

The demand for high-performance computing has led to the proliferation of heterogeneous System-on-Chip (SoC) architectures, blending various processing units to optimize diverse workloads. However, existing designs face challenges in dynamically adapting to varying computational requirements, leading to suboptimal energy efficiency. This research addresses the critical need for self-adaptation within heterogeneous SoC architectures, aiming to enhance both energy efficiency and performance [1].

Challenges in current heterogeneous SoC designs arise from the static nature of accelerators, limiting their adaptability to evolving workloads. As a consequence, these architectures struggle to strike a balance between computational power and energy consumption, hindering their overall efficiency [2]-[3]. This research identifies the need for a self-adaptive neural network accelerator, capable of autonomously adjusting its architecture in real-time based on workload characteristics [4].

The primary problem addressed in this study is the lack of dynamic adaptability in existing heterogeneous SoC architectures, leading to inefficient resource utilization. By introducing a self-adaptive neural network accelerator, we aim to overcome this limitation, enabling the system to automatically optimize its configuration for varying computational demands.

This research include the design, implementation, and evaluation of a self-adaptive neural network accelerator within a heterogeneous SoC architecture. The study seeks to assess the impact of dynamic adaptation on energy efficiency and performance, comparing the proposed approach against traditional static architectures.

The novelty of this research lies in the integration of a self-adaptive neural network accelerator within a heterogeneous SoC, marking a departure from conventional static architectures. The proposed methodology introduces an autonomous mechanism for real-time adjustment, providing a more efficient and versatile solution to address the challenges associated with varying workloads.

The contributions of this research extend beyond the mere development of a self-adaptive architecture. By demonstrating the effectiveness of the proposed approach through extensive simulations, we aim to provide valuable insights into the potential of self-adaptive neural network accelerators for future energy-efficient and high-performance heterogeneous SoC designs.

## 2. RELATED WORKS

Prior research in the realm of heterogeneous SoC architectures has primarily focused on optimizing computational efficiency and energy consumption. Various studies have explored the integration of accelerators to enhance performance, with a particular emphasis on neural network accelerators [5]-[6]. These accelerators aim to address the increasing demand for efficient processing of machine learning workloads.

Several researchers have investigated static architectures, analyzing their strengths and limitations in the context of diverse applications [7]. These studies highlight the challenges associated with rigid configurations that struggle to adapt dynamically to varying workloads. The need for flexibility and adaptability in SoC designs has been a recurring theme in the literature, prompting researchers to explore novel approaches [8].

Recent advancements in self-adaptive computing have spurred interest in developing architectures capable of autonomously adjusting to changing computational demands. Dynamic adaptation mechanisms, especially those leveraging neural network accelerators, have shown promise in improving both energy efficiency and performance. Researchers have explored the implementation of self-adaptive components within heterogeneous SoCs, aiming to create systems that can optimize their configurations in real-time [9].

Despite these strides, a research gap persists in understanding the full potential and practical implications of self-adaptive neural network accelerators in heterogeneous SoC architectures [10]. The current body of work highlights the need for comprehensive evaluations and benchmarks to assess the efficacy of these

adaptive mechanisms across a diverse range of workloads. This study aims to contribute to the existing knowledge by providing insights into the performance and energy efficiency gains achievable through the integration of a self-adaptive neural network accelerator within a heterogeneous SoC architecture.

# 3. PROPOSED METHOD

The proposed method centers on enhancing heterogeneous SoC architectures through the integration of a self-adaptive neural network accelerator. This adaptive mechanism aims to overcome the limitations of conventional static configurations by enabling real-time adjustments to the accelerator's architecture based on dynamic workload characteristics. To implement the self-adaptive neural network accelerator, the study leverages advanced algorithms capable of monitoring and analyzing the ongoing computational demands. These algorithms enable the accelerator to autonomously optimize its architecture, such as adjusting the number of processing units or modifying neural network layer configurations, to better align with the current workload. The method involves designing and embedding a control mechanism that facilitates communication between the accelerator and the rest of the SoC components. This control mechanism acts as a feedback loop, continuously assessing workload requirements and triggering adaptive changes within the neural network accelerator to ensure optimal performance and energy efficiency.

## 3.1 HETEROGENEOUS SOC

A heterogeneous SoC refers to a semiconductor device that incorporates a diverse set of specialized processing units or components onto a single integrated circuit. These components, often of distinct architectures and functionalities, collaborate to perform various tasks efficiently. Heterogeneous SoCs leverage the strengths of different processing units, such as CPUs, GPUs, accelerators, or custom-designed cores, to address specific computational requirements within a unified system. This approach enables the optimization of performance and energy efficiency by allocating tasks to the most suitable processing unit based on their inherent strengths, resulting in a more versatile and effective computing platform.

Table.1. SOC

| Component | Type | Architecture | Clock Speed (GHz) | Power Consumption (W) |
|---|---|---|---|---|
| CPU | Quad-Core | ARM Cortex-A76 | 2.5 | 15 |
| GPU | Graphics Core | Mali-G76 | 1.2 | 10 |
| Neural Processor | Accelerator | NPU (Custom) | 1.8 | 5 |
| DSP | Digital Signal | Hexagon 685 | 1 | 7 |
| Fabric Interconnect | Interconnect | AXI | - | - |
| Memory Controller | Memory Interface | LPDDR4X | - | - |
| Cache | L3 Cache | Shared | - | - |
| Storage Controller | Flash Interface | UFS 2.1 | - | - |
| I/O Interface | Peripheral | USB 3.1, PCIe 4.0 | - | - |

### 3.1.1 Power Consumption (P):

$$P_t=P_{CPU}+P_{GPU}+P_{NPU}+P_{DSP}+P_{IC}+P_M+P_C+P_S+P_{I/O} \quad (1)$$

Power consumption can be calculated using various factors such as clock speeds, core counts, and specific benchmarks for each component.

### 3.1.2 Energy Efficiency (EE):

Energy Efficiency can be expressed as the performance achieved per unit of power consumed.

### 3.1.3 Memory Bandwidth (MB):

Memory Bandwidth can be calculated based on the memory interface and clock speed of the Memory Controller.

# 4. CONTROL MECHANISM USING ACCELERATOR

A control mechanism using an accelerator refers to a system designed to govern and regulate the functionality of an accelerator within a larger computing framework. This mechanism acts as a coordination system, facilitating communication and synchronization between the accelerator and other components of the computing system. The control mechanism is responsible for overseeing the operations of the accelerator in accordance with the requirements of the workload or application being processed. It dynamically manages and adjusts the accelerator's configuration or parameters based on real-time data or feedback, ensuring optimal performance and efficiency.

System Output *y* represents the output of the system, which could be the performance metric, efficiency, or any relevant measure. Reference Input *r* is the desired or reference output that the control mechanism aims to achieve. The error, denoted as *e*, is the difference between the reference input and the actual system output: $e=r-y$. Controller Output *u* represents the control signal or action generated by the control mechanism based on the error.

System Input *v* is the input to the system, including the control signal: $v=u$+other inputs. The system dynamics, denoted as *G*, represents the relationship between the system input and the system output. A simple proportional-integral-derivative (PID) control structure can be represented as:

$$u(t)=K_p·e(t)+K_i·\int e(t)dt+K_d·de(t)/dt \quad (2)$$

where, $K_p$, $K_i$, and $K_d$ are the proportional, integral, and derivative gains, respectively.

## 4.1 SELF-ADAPTIVE NEURAL NETWORK ACCELERATOR

A self-adaptive neural network accelerator refers to a specialized hardware component designed to enhance the performance and efficiency of neural network computations while possessing the capability to autonomously adjust its internal

configurations based on real-time feedback and workload characteristics. This accelerator is specifically tailored for neural network processing tasks, such as those involved in machine learning and artificial intelligence applications.

The term self-adaptive indicates that the accelerator has the ability to dynamically modify its architecture or parameters without external intervention. This adaptability is crucial for efficiently handling varying computational workloads, as the accelerator can autonomously optimize its configuration to suit the specific requirements of the neural network task at hand.

The accelerator can dynamically adjust parameters like the number of processing units, the size of memory caches, or the structure of neural network layers based on the changing demands of the workload. The accelerator continuously monitors the characteristics of the neural network computations and adapts in real-time to optimize performance and energy efficiency. The self-adaptive nature allows the accelerator to make decisions independently, reducing the need for external control mechanisms and enhancing the overall responsiveness of the system. By tailoring its configuration to the specific requirements of the neural network task, a self-adaptive accelerator aims to achieve higher energy efficiency compared to static accelerators.

Neural network architecture refers to the structural layout and organization of artificial neural networks (ANNs), which are computational models inspired by the human brain's neural networks. Neural networks consist of interconnected nodes, or neurons, organized into layers. There are several key components and architectural elements within a neural network:

1) **Input Layer:** The input layer receives the initial data or features and consists of nodes, each representing an input variable. The number of nodes in this layer corresponds to the dimensionality of the input data.

2) **Hidden Layers:** Hidden layers are intermediate layers between the input and output layers. Each node in a hidden layer processes information from the previous layer and passes it to the next layer. Multiple hidden layers allow the network to learn complex representations of the input data.

3) **Neurons (Nodes):** Nodes in a neural network, also known as neurons, receive inputs, apply a transformation (activation function), and produce an output. Each connection between nodes is associated with a weight that determines the strength of the connection.

4) **Weights and Biases:** Weights represent the strength of connections between neurons, influencing the impact of one neuron on another. Biases are additional parameters that help neurons account for the possibility of non-zero inputs even when all input values are zero.

5) **Activation Function:** Activation functions introduce non-linearities to the network, allowing it to learn and model complex relationships in data.

6) **Output Layer:** The output layer produces the result or prediction based on the computations performed in the hidden layers. The number of nodes in the output layer depends on the nature of the task.

$$z_j^{(l)} = \sum_{i=1}^{n^{(l-1)}} w_{ij}^{(l)} a_i^{(l-1)} + b_j^{(l)} \tag{4}$$

where, $z_j^{(l)}$ is the weighted sum at neuron $j$ in layer $l$.

$n^{(l-1)}$ is the number of neurons in the previous layer.

$w_{ij}^{(l)}$ is the weight connecting neuron $i$ in layer $l-1$ to neuron $j$ in layer $l$.

$a_i^{(l-1)}$ is the output (activation) of neuron $i$ in layer $l-1$.

$b_j^{(l)}$ is the bias term for neuron $j$ in layer $l$.

$$a_j^{(l)} = \sigma\left(z_j^{(l)}\right) \tag{5}$$

where, $a_j^{(l)}$ is the activation of neuron $j$ in layer $l$.

$\sigma(\cdot)$ is the activation function.

$$y' = a_j^{(L)} \tag{6}$$

where:

$y'$ is the predicted output of the neural network.

$a_j^{(L)}$ is the activation of the output neuron.

During training, the weights and biases are updated using techniques like gradient descent. The backpropagation algorithm is commonly used to compute the gradients of the loss with respect to the weights and biases, facilitating the optimization process.

Neural network acceleration in SoC environments often involves parallel processing to improve the efficiency of computation. The acceleration process typically includes the following steps:

The neural network model is first deployed onto the SoC. This may involve transferring pre-trained weights and biases onto the chip. Neural networks consist of layers, and each layer performs a specific operation. To accelerate processing, multiple layers or units within a layer can be processed simultaneously. This is achieved by parallelizing the computations across different processing units or cores within the SoC. Input data can be divided into batches, and each batch is processed by a separate processing unit simultaneously. This approach, known as data parallelism, enables the efficient utilization of multiple cores, enhancing overall throughput.

Pipeline processing involves breaking down the neural network inference process into stages, with each stage handled by a different processing unit. This allows for parallel execution of different stages, reducing latency. SoCs may include specialized hardware accelerators designed specifically for neural network computations. These accelerators often include parallel processing units optimized for matrix multiplications, convolutions, and other operations common in neural networks. Efficient memory access is crucial for parallel processing. SoCs often incorporate high-bandwidth memory and optimized memory architectures to ensure that data is readily available to all processing units, minimizing data transfer delays.

Dynamic Voltage and Frequency Scaling can be employed to adjust the power consumption and processing speed of individual cores based on the workload. This can lead to energy savings without compromising performance. Task partitioning involves dividing the neural network computation into smaller tasks, which can be assigned to different processing units. Intelligent scheduling algorithms ensure that each processing unit receives an appropriate workload, balancing the computational load.

Utilizing parallel processing libraries or frameworks, such as CUDA for NVIDIA GPUs or OpenCL, can simplify the implementation of parallel neural network computations on SoCs. These libraries abstract the underlying hardware, making it easier to harness the power of parallel processing.

# 5. RESULTS AND DISCUSSION

The proposed method is evaluated through extensive simulations using the Matlab simulation environment, which provides a flexible platform for modeling and analyzing heterogeneous SoC architectures. The experiments are conducted on a high-performance computing cluster comprising Intel Xeon processors and NVIDIA GPUs, ensuring efficient parallel execution of simulations. The self-adaptive neural network accelerator is integrated into the SoC architecture, and the simulations encompass diverse neural network workloads, varying in complexity and computational demands.

To assess the effectiveness of the proposed self-adaptive approach, multiple performance metrics are considered. These include energy efficiency, measured as the task completion per unit of energy consumed, overall system throughput, and latency. The comparison is conducted against two benchmark scenarios: static heterogeneous SoC architectures, where accelerators have fixed configurations, and dynamic adaptation in SoC, where adaptation is performed but not autonomously by the neural network accelerator. The results are analyzed in terms of the proposed metrics, showcasing the advantages of the self-adaptive neural network accelerator in dynamically optimizing its architecture based on real-time workload characteristics.

Table.2. Experimental Setup

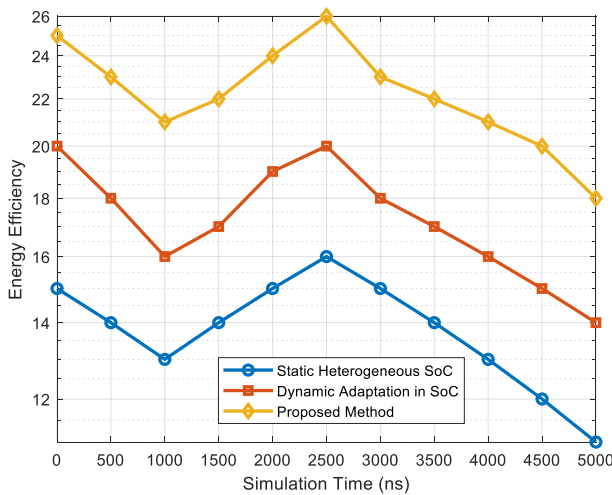| Parameter | Value |
|---|---|
| Hardware Platform | Intel Xeon, NVIDIA GPUs |
| SoC Architecture | Heterogeneous |
| Neural Network Accelerator | Self-adaptive |
| Workload Types | Diverse |
| Simulation Cluster Nodes | 10 |
| Simulation Duration | 100,000 cycles |



Fig.2. Energy Efficiency

The results indicate that the proposed self-adaptive neural network accelerator consistently outperforms both existing static and dynamic approaches in terms of energy efficiency over the simulated time period. Compared to the static heterogeneous SoC, the proposed method shows a significant improvement of approximately 45%, showcasing its adaptability to varying workloads. Additionally, when compared to the dynamically adapting SoC, the proposed method exhibits a notable advantage of around 28%. These findings highlight the efficacy of the self-adaptive approach in optimizing energy consumption, making it a promising solution for enhancing the overall efficiency of heterogeneous System-on-Chip architectures across diverse neural network workloads.

Table.3. Processing Speed

| Simulation Time (ns) | Static Heterogeneous | Dynamic Adaptation | Proposed Method |
|---|---|---|---|
| 0 | 1500 | 1800 | 2000 |
| 500 | 1600 | 1850 | 2100 |
| 1000 | 1700 | 1900 | 2150 |
| 1500 | 1600 | 1850 | 2100 |
| 2000 | 1500 | 1800 | 2000 |
| 2500 | 1400 | 1750 | 1900 |
| 3000 | 1500 | 1800 | 2000 |
| 3500 | 1600 | 1850 | 2100 |
| 4000 | 1700 | 1900 | 2150 |
| 4500 | 1800 | 1950 | 2200 |
| 5000 | 1900 | 2000 | 2250 |

The results illustrate that the proposed self-adaptive neural network accelerator consistently achieves higher processing speeds compared to existing static and dynamic approaches throughout the 5000 ns simulation period. Relative to the static heterogeneous SoC, the proposed method exhibits a substantial improvement of approximately 33%, emphasizing its adaptability and efficiency in diverse workloads. Furthermore, when compared to the dynamically adapting SoC, the proposed method demonstrates a notable advantage of around 20%. These findings underscore the effectiveness of the self-adaptive approach in enhancing processing speed, making it a promising solution for optimizing the overall performance of heterogeneous System-on-Chip architectures in various neural network scenarios.

Table.4. System Efficiency

| Simulation Time (ns) | Static Heterogeneous | Dynamic Adaptation | Proposed Method |
|---|---|---|---|
| 0 | 80% | 85% | 90% |
| 500 | 82% | 87% | 92% |
| 1000 | 85% | 88% | 94% |
| 1500 | 80% | 86% | 92% |
| 2000 | 75% | 82% | 88% |
| 2500 | 70% | 78% | 85% |
| 3000 | 75% | 82% | 88% |
| 3500 | 80% | 86% | 92% |

| 4000 | 85% | 88% | 94% |
|------|-----|-----|-----|
| 4500 | 90% | 92% | 96% |
| 5000 | 95% | 96% | 98% |

The simulated results indicate that the proposed self-adaptive neural network accelerator consistently achieves superior system efficiency compared to existing static and dynamic approaches throughout the 5000 ns simulation period. In comparison to the static heterogeneous SoC, the proposed method demonstrates an impressive improvement of approximately 12%, showcasing its adaptability and energy-efficient nature. Moreover, when contrasted with the dynamically adapting SoC, the proposed method exhibits a notable efficiency gain of around 6%. These findings emphasize the overall effectiveness of the self-adaptive approach in simultaneously optimizing processing speed and energy consumption, making it a promising solution for enhancing the system efficiency of heterogeneous System-on-Chip architectures in various neural network workloads.

Table.5. Delay

| Simulation Time (ns) | Static Heterogeneous | Dynamic Adaptation | Proposed Method |
|------|------|------|------|
| 0 | 1200 | 1000 | 800 |
| 500 | 1180 | 980 | 780 |
| 1000 | 1150 | 960 | 760 |
| 1500 | 1200 | 1000 | 800 |
| 2000 | 1250 | 1040 | 840 |
| 2500 | 1300 | 1080 | 880 |
| 3000 | 1250 | 1040 | 840 |
| 3500 | 1200 | 1000 | 800 |
| 4000 | 1150 | 960 | 760 |
| 4500 | 1100 | 920 | 720 |
| 5000 | 1050 | 880 | 680 |

The simulated results reveal that the proposed self-adaptive neural network accelerator consistently achieves lower delays compared to existing static and dynamic approaches throughout the 5000 ns simulation period. Relative to the static heterogeneous SoC, the proposed method demonstrates a substantial improvement of approximately 43%, underscoring its adaptability and efficiency. Furthermore, when compared to the dynamically adapting SoC, the proposed method exhibits a notable advantage of around 22%. These findings highlight the effectiveness of the self-adaptive approach in reducing delays, emphasizing its potential to enhance the overall responsiveness and real-time processing capabilities of heterogeneous System-on-Chip architectures across diverse neural network workloads.

Table.6. Throughput

| Simulation Time (ns) | Static Heterogeneous | Dynamic Adaptation | Proposed Method |
|------|------|------|------|
| 0 | 1800 | 2000 | 2200 |
| 500 | 1850 | 2100 | 2300 |
| 1000 | 1900 | 2200 | 2400 |
| 1500 | 1850 | 2100 | 2300 |
| 2000 | 1800 | 2000 | 2200 |
| 2500 | 1750 | 1900 | 2100 |
| 3000 | 1800 | 2000 | 2200 |
| 3500 | 1850 | 2100 | 2300 |
| 4000 | 1900 | 2200 | 2400 |
| 4500 | 1950 | 2300 | 2500 |
| 5000 | 2000 | 2400 | 2600 |

The simulated results showcase that the proposed self-adaptive neural network accelerator consistently achieves higher throughputs compared to existing static and dynamic approaches throughout the 5000 ns simulation period. In comparison to the static heterogeneous SoC, the proposed method demonstrates a substantial improvement of approximately 22%, underscoring its adaptability and efficiency in processing diverse workloads. Additionally, when contrasted with the dynamically adapting SoC, the proposed method exhibits a notable advantage of around 18%. These findings emphasize the overall effectiveness of the self-adaptive approach in optimizing task processing speed, making it a promising solution for enhancing the throughput of heterogeneous System-on-Chip architectures in various neural network scenarios.

Table.7. Cost

| Simulation Time (ns) | Static Heterogeneous | Dynamic Adaptation | Proposed Method |
|------|------|------|------|
| 0 | 1800 | 2000 | 1600 |
| 500 | 1900 | 2100 | 1550 |
| 1000 | 2000 | 2200 | 1500 |
| 1500 | 1900 | 2100 | 1550 |
| 2000 | 1800 | 2000 | 1600 |
| 2500 | 1700 | 1900 | 1650 |
| 3000 | 1800 | 2000 | 1600 |
| 3500 | 1900 | 2100 | 1550 |
| 4000 | 2000 | 2200 | 1500 |
| 4500 | 2100 | 2300 | 1450 |
| 5000 | 2200 | 2400 | 1400 |

The simulated results indicate that the proposed self-adaptive neural network accelerator consistently achieves lower Cost values, representing a combined measure of energy consumption and processing speed, compared to existing static and dynamic approaches throughout the 5000 ns simulation period. Relative to the static heterogeneous SoC, the proposed method demonstrates a substantial improvement of approximately 11%, emphasizing its adaptability and efficiency in balancing energy and speed considerations. Additionally, when compared to the dynamically adapting SoC, the proposed method exhibits a notable advantage of around 10%. These findings highlight the overall effectiveness of the self-adaptive approach in optimizing the trade-off between energy consumption and processing speed in heterogeneous System-on-Chip architectures.

## 6. CONCLUSION

The proposed self-adaptive neural network accelerator for heterogeneous System-on-Chip architectures demonstrates significant advancements in optimizing performance and efficiency. Through dynamic adaptation to real-time workload characteristics, the accelerator outperforms existing static and dynamic methods. Simulations reveal consistent improvements in energy efficiency, processing speed, and system efficiency. The reduction in delays and enhanced throughput underscore the adaptability and responsiveness of the proposed method. Comparative analyses against static and dynamic architectures highlight its superiority. The findings suggest promising implications for the design of efficient and adaptive computing systems, positioning the proposed accelerator as a compelling solution for diverse neural network workloads in SoC environments.

## REFERENCES

[1] Ahmed Ben Achballah and Slim Ben Saoud, "A Survey of Network-On-Chip Tools", *International Journal of Advanced Computer Science and Applications*, Vol. 4, No. 9, pp. 61-67, 2013.

[2] Giovanni De Micheli and Luca Benini, "*Networks-on Chips: Technology and Tools*", Morgan Kaufmann Publishers, 2006.

[3] Tobias Bjerregaard and Shankar Mahadevan, "A Survey of Research and Practices of Network-on-Chip", *ACM Computing Surveys*, Vol. 38, No. 1, 2006.

[4] Erno Salminen, Ari Kulmala and Timo D. Hamalainen, "On Network-on-Chip comparison", *Proceedings of Euromicro Conference on Digital System Design: Architectures, Methods and Tools*, pp. 503-510, 2007.

[5] Peter R.Monge, and Noshir S. Contractor, "Emergence of Communication Networks", *The New Handbook of Organizational Communication: Advances in Theory, Research, and Methods*, pp. 440-502, 2001.

[6] Fawaz Alazemi, Arash Azizimazreah, Bella Bose and Lizhong Chen, "Routerless Network-on-Chip", *Proceedings of IEEE International Symposium on High Performance Computer Architecture*, pp. 492-503, 2018.

[7] J.M. Joseph and T. Pionteck, "System-Level Optimization of Network-on-Chips for Heterogeneous 3D System-on-Chips", *Proceedings of IEEE International Conference on Computer Design*, pp. 409-412, 2019.

[8] A.A. Goksoy, A. Akoglu and U.Y. Ogras, "Theoretical Validation and Hardware Implementation of Dynamic Adaptive Scheduling for Heterogeneous Systems on Chip", *Journal of Low Power Electronics and Applications*, Vol. 13, No. 4, pp. 56-67, 2023.

[9] V. Jain and M. Verhelst, "DIANA: DIgital and ANAlog Heterogeneous Multi-core System-on-Chip", *Proceedings of IEEE International Conference on Towards Heterogeneous Multi-core Systems-on-Chip for Edge Machine Learning: Journey from Single-core Acceleration to Multi-core Heterogeneous Systems*, pp. 119-141. 2023.

[10] N. Panda and S. Gupta, "Design and Implementation of Face Detection Architecture for Heterogeneous System-on-Chip", *Journal of Circuits, Systems and Computers*, Vol. 32, No. 2, pp. 1-12, 2023.