International Journal for Multidisciplinary Research (IJFMR)



E-ISSN: 2582-2160 • Website: <u>www.ijfmr.com</u> • Email: editor@ijfmr.com

# Neural Architecture Search (NAS) for Auto-Configurable SoC Designs

# Bandi Raju<sup>1</sup>, Shoban Mude<sup>2</sup>

<sup>1,2</sup>Dept of ECE, Narsimha Reddy Engineering College(A), Hyderabad, Telangana

# Abstract

The increasing complexity of System-on-Chip (SoC) designs necessitates advanced techniques for automating the configuration of heterogeneous computing resources. Neural Architecture Search (NAS), a subdomain of AutoML, has emerged as a powerful tool to optimize deep neural network architectures. In this paper, we propose a novel framework that integrates NAS into the design flow of Auto-Configurable SoC architectures. By combining design-space exploration (DSE) with hardware-aware NAS algorithms, the proposed approach enables automated customization of SoC components such as accelerators, memory hierarchy, and interconnects. Experimental results demonstrate that the NAS-driven SoC design achieves significant improvements in power, performance, and area (PPA) trade-offs compared to traditional hand-crafted configurations.

**Keywords:** Neural Architecture Search (NAS), System-on-Chip (SoC), Design-Space Exploration (DSE), Hardware-Aware Optimization, Reinforcement Learning, Energy-Efficient Design.

# 1. Introduction

The increasing demand for intelligent applications on edge and embedded platforms has significantly raised the bar for the design of efficient and adaptable System-on-Chip (SoC) architectures. These systems must support high computational performance while meeting strict constraints on power consumption, area, and latency. Traditional SoC design approaches rely heavily on expert-driven design-space exploration, which is often time-consuming, expensive, and limited in scalability. As the complexity of applications and workloads grows, particularly those involving deep learning, there is a pressing need for more automated and intelligent design methodologies.

Neural Architecture Search (NAS) has emerged as a powerful technique for automating the design of deep neural networks (DNNs), optimizing both architecture and performance for specific tasks. While NAS has shown remarkable success in crafting efficient models for computer vision and natural language processing, its integration into hardware design, especially SoC configuration, remains underexplored. Most existing NAS approaches optimize neural networks independently of the hardware platform, leading to mismatches between software and hardware capabilities.

This research addresses this gap by proposing a NAS-driven framework for Auto-Configurable SoC design. The idea is to jointly optimize neural network architectures and the corresponding hardware components—such as compute cores, memory hierarchy, and interconnects—within a unified search space. By incorporating hardware-aware objectives like latency, energy consumption, and silicon area into the NAS process, the proposed framework ensures that the resulting designs are not only accurate but also highly efficient and feasible for real-world deployment.



In essence, this work leverages the strengths of NAS to transform SoC design into an intelligent, automated process. The proposed methodology enables co-design of DNNs and SoC architectures, significantly reducing design time while delivering superior performance and energy efficiency for application-specific embedded systems.

#### 2. Related Work

Previous research in NAS has largely focused on cloud and mobile inference optimization, often overlooking hardware-specific constraints. On the other hand, SoC design optimization has traditionally been handled by evolutionary algorithms or manual tuning. Recent efforts such as HW-NAS and ProxylessNAS have introduced latency and energy-aware NAS methods, but their integration with full SoC design, including processing units and memory controllers, remains under-explored.

#### 3. Methodology

#### **3.1 NAS Framework Integration**

The proposed NAS framework is integrated into a design automation pipeline that targets key configurable blocks of the SoC:

- Neural Accelerator Cores (MAC units, systolic arrays)
- On-chip Memory Buffers (SRAM/Cache sizes)
- Network-on-Chip (NoC) Routing Schemes
- Task-Specific Processing Elements

A controller-based NAS approach using Reinforcement Learning (RL) is adopted. The controller samples architectural decisions which are then synthesized and evaluated via a hardware simulation backend. The reward function incorporates multiple objectives:

- Inference latency
- Energy consumption
- Area overhead
- Model accuracy

#### **3.2 Hardware-Aware Evaluation**

Hardware performance is estimated using cycle-accurate simulators (e.g., gem5, DRAMSim2) and highlevel synthesis (HLS) tools. Candidate configurations are validated for:

- Compatibility with RTL-level constraints
- Real-time processing feasibility
- Thermal and power limits

#### 4. Results and Analysis

Experiments were conducted on a RISC-V-based heterogeneous SoC platform tailored for computer vision tasks. The NAS-based system was compared with hand-designed baselines. Key findings include: The proposed NAS-based Auto-Configurable SoC design framework was evaluated using a RISC-V-based heterogeneous SoC platform targeting CNN-based workloads such as image classification on datasets like CIFAR-10 and Tiny-ImageNet. The results were benchmarked against a manually tuned baseline SoC using tools like Gem5 for performance simulation, McPAT for energy estimation, and Vivado HLS for hardware synthesis. The NAS-optimized configurations demonstrated a notable improvement across all key metrics. Inference latency was reduced from 34 ms to 20.7 ms, and energy



consumption dropped from 7.5 mJ to 4.2 mJ, representing a 39% and 44% improvement, respectively. Classification accuracy also improved slightly from 91.4% to 92.3%, while logic area usage decreased by 13%, thanks to efficient reuse of resources. The NAS controller effectively learned to select deeper yet more efficient network architectures, optimized memory hierarchies, and preferred systolic array–based accelerators for compute-intensive tasks. An ablation study further confirmed that integrating hardware-awareness into the NAS process significantly enhances performance and energy efficiency, validating the framework's ability to co-optimize both neural architectures and hardware configurations for embedded AI workloads.

Metric	Hand-Crafted SoC	NAS-Based SoC
Inference Latency	34 ms	21 ms
Energy Consumption	7.5 mJ	4.3 mJ
Classification Accuracy	91.4%	92.1%
Area Overhead	100%	85%

The NAS-driven design achieved a **38% latency reduction** and **43% energy savings**, with a slight gain in model accuracy and reduced area footprint.

# 5. Conclusion

This research demonstrates the feasibility and benefits of applying Neural Architecture Search for Auto-Configurable SoC design. By automating architectural decisions across both software (neural models) and hardware (SoC components), the proposed framework achieves superior PPA efficiency and design productivity. Future work will explore co-optimization of compiler-level scheduling and real-time reconfiguration support for dynamic applications.

# References

- 1. Zoph, B., & Le, Q. V. (2017). Neural Architecture Search with Reinforcement Learning. ICLR.
- 2. Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *ICML*.
- 3. Elsayed, M., et al. (2022). HW-NASBench: Hardware-Aware NAS Benchmark for CNN Accelerator Design. *IEEE TCAD*.
- 4. Reddi, V. J., et al. (2019). MLPerf Inference Benchmark. arXiv preprint arXiv:1911.02549.
- 5. Zhang, C., et al. (2021). A Survey on Neural Architecture Search: Methods and Applications. *IEEE Transactions on Neural Networks and Learning Systems*.