

# Automated ASIC Design Optimization Using Neural Architecture Search Techniques

Rashmitha Reddy Vuppunuthula\*

Senior Electrical Design Engineer, Exton, Pennsylvania - 19341, USA

**Citation:** Rashmitha RV. Automated ASIC Design Optimization Using Neural Architecture Search Techniques. *Int J Cur Res Sci Eng Tech* 2024; 7(4), 95-102. DOI: doi.org/10.30967/IJCRSET/Rashmitha-Reddy-Vuppunuthula/149

**Received:** 19 November, 2024; **Accepted:** 16 December, 2024; **Published:** 19 December, 2024

\***Corresponding author:** Rashmitha Reddy Vuppunuthula, Senior Electrical Design Engineer, Exton, Pennsylvania - 19341, USA, E-mail: vrashmitha97@gmail.com

**Copyright:** © 2024 Rashmitha RV., This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

## ABSTRACT

Application-Specific Integrated Circuits (ASICs) are tailored to maximize performance, efficiency and resource utilization for specific tasks, yet the complexity of designing optimized ASIC architectures remains a significant challenge. This paper explores the use of Neural Architecture Search (NAS) to streamline and enhance ASIC design by automating the selection of optimal architectures based on predefined performance, area and power criteria. NAS algorithms, trained on diverse datasets of ASIC designs, dynamically identify the most effective configurations, achieving substantial improvements across key metrics. Results demonstrate that NAS-driven approaches reduce power consumption by up to 30%, with designs such as Design 5 showing a decrease from 300 mW to 210 mW. Area utilization was also optimized, with reductions reaching 30%, exemplified by Design 5's decrease from 40 mm<sup>2</sup> to 28 mm<sup>2</sup>. Additionally, computation speed was significantly enhanced, with improvements of up to 37.14%, as Design 5's speed increased from 3.5 GHz to 4.8 GHz. These findings confirm that NAS-based methodologies outperform traditional design approaches, offering a robust and scalable framework that accelerates the ASIC design process while elevating the final product's efficiency and performance. This work establishes NAS as a pivotal tool for ASIC design optimization in increasingly complex electronic systems.

**Keywords:** ASIC Design, Neural Architecture Search (NAS), Design Automation, Power and Performance Optimization, Energy-Efficient Computing, Machine Learning in ASICs

## 1. Introduction

Application-Specific Integrated Circuits (ASICs) have become a pivotal component in modern electronics, designed to deliver unmatched performance, energy efficiency and scalability for specific applications<sup>1</sup>. However, the manual design and optimization of ASIC architectures remain an intricate and time-intensive process. This complexity arises from the need to balance multiple design parameters, such as power consumption, latency and area efficiency, while meeting the performance requirements of diverse and increasingly demanding applications<sup>2</sup>. Traditional ASIC design methods often require significant expertise and rely heavily on iterative, trial-and-error approaches, limiting the scope for rapid innovation<sup>3</sup>. Neural Architecture Search

(NAS) has emerged as a transformative approach in the design of neural networks, enabling the automated discovery of optimal architectures tailored to specific tasks<sup>4</sup>. Extending this capability to ASIC design introduces the potential for significant advancements in the field. By integrating NAS techniques, designers can automate the exploration of both architectural and hardware configurations, identifying solutions that achieve optimal performance and resource utilization<sup>5</sup>. This approach not only reduces the manual effort required but also accelerates the design cycle, enabling faster deployment of cutting-edge ASICs in applications such as machine learning, edge computing and high-performance data centers<sup>6</sup>.

The intersection of NAS and ASIC design presents unique

challenges, including the expansive design space and the need for hardware-aware optimization. Addressing these challenges necessitates innovative frameworks capable of co-exploring neural architectures and ASIC configurations in tandem<sup>7,8</sup>. Such frameworks can leverage advanced machine learning algorithms and high-level synthesis tools to streamline the design process, ensuring that resulting architectures are not only functionally superior but also aligned with practical hardware constraints<sup>9</sup>. This paper explores the use of NAS-driven methodologies for automated ASIC design optimization. By combining algorithmic precision with hardware adaptability, NAS enables a paradigm shift in how ASICs are conceptualized, designed and implemented. The findings of this study highlight the advantages of NAS in enhancing ASIC design workflows, providing a foundation for future innovations in hardware optimization<sup>10</sup>.

The success of deep neural networks (DNNs) in applications such as computer vision, virtual reality and recommender systems has significantly increased their deployment across diverse domains. This rapid adoption, however, has been accompanied by challenges, particularly due to the growing complexity of models, which demand higher computation and energy resources. Large-scale DNNs, while offering high accuracy, suffer from increased latency and energy consumption, necessitating efficient methods for hardware and software optimization<sup>2</sup>. These demands have driven the development of automated Neural Architecture Search (NAS) techniques and hardware-aware optimization strategies<sup>5,11</sup>. Application-Specific Integrated Circuits (ASICs) have emerged as a powerful solution for meeting the high-performance and energy-efficiency requirements of DNNs. However, designing ASICs manually is labor-intensive, involving a thorough exploration of design parameters such as buffer sizes, processing elements and dataflows<sup>12</sup>. Traditional approaches often focus on fixed hardware configurations, which limit flexibility and fail to address the specific optimization needs of various neural architectures<sup>13</sup>.

To address these issues, co-design frameworks that integrate NAS with hardware design optimization have gained prominence. Such frameworks leverage advanced algorithms, including differentiable NAS and reinforcement learning techniques, to simultaneously optimize neural architectures and hardware configurations<sup>7</sup>. This co-design process allows for the identification of configurations that meet both performance and energy constraints while adhering to practical hardware limitations<sup>14</sup>. Moreover, these methods facilitate exploration across sparsely valid design spaces, overcoming challenges posed by invalid or non-synthesizable hardware designs<sup>15</sup>. The interdependence between neural architectures and hardware accelerators complicates the optimization process. Parameters such as compute array dimensions, on-chip memory usage and compiler strategies must align to maximize efficiency and performance<sup>16</sup>. Recent advancements in co-design frameworks have demonstrated the feasibility of integrating NAS with hardware-aware optimization, resulting in accelerators that deliver significant improvements in energy efficiency and computational throughput<sup>17</sup>. This paper focuses on the automated co-design of neural networks and ASIC accelerators using NAS techniques. The study explores methods for optimizing sparsely valid design spaces while minimizing runtime and energy consumption. By addressing the challenges of co-optimization, this work aims to advance the development of efficient and

scalable DNN accelerators that cater to the growing demands of modern AI systems.

## 2. Literature review

The rapid development of FPGA- and ASIC-based accelerators has greatly advanced the implementation of deep neural networks (DNNs). These accelerators have been tailored to optimize performance and energy efficiency, addressing the computational demands of modern AI applications. Research on FPGA-based accelerators has explored various strategies to enhance the execution of DNN workloads. For instance, loop tiling techniques have been employed to optimize convolutional layer computations<sup>18</sup>. Similarly, the DNN Builder accelerator framework introduces resource allocation strategies, layer-based pipelining and cache optimizations to improve FPGA performance<sup>19</sup>. Another notable work leverages multi-level parallelisms, encompassing tasks, layers, loops and operators, to enhance throughput<sup>20</sup>. Additionally, recent designs have integrated DNN and hardware co-design principles, using algorithmic and architectural optimizations to ensure efficient hardware constraints capture<sup>8</sup>. For ASIC-based accelerators, industry and academic efforts have yielded innovative designs tailored to specific applications. Noteworthy examples include Google's TPU, which achieves high performance for machine learning workloads<sup>21</sup> and ShiDian Nao, which shifts vision processing closer to sensors to reduce energy consumption<sup>22</sup>. The Eyeriss architecture emphasizes spatial dataflow to optimize energy usage and latency in DNN computations<sup>6</sup>. These accelerators demonstrate how customized architectures can be tuned to meet diverse application requirements, providing unmatched energy efficiency and computational throughput.

Performance prediction is critical for designing both FPGA- and ASIC-based accelerators. FPGA accelerator designs often rely on roofline models or analytical tools for performance estimation<sup>18,20</sup>. In the ASIC domain, performance prediction methodologies have been developed to address latency, energy and memory access overheads. The Eyeriss framework includes an energy model for memory and computation units, alongside a delay model for latency calculation<sup>6</sup>. Similarly, MAESTRO provides an open-source infrastructure for modelling dataflows within DNN accelerators<sup>9</sup>, while Time loop adopts a loop-based workload description for memory access and latency analysis<sup>10</sup>. These tools enable accurate performance estimations, guiding the design of efficient accelerators. To meet the growing demand for DNN accelerators, automated generation tools have been developed. Deep Burning facilitates the creation of FPGA-based accelerators through a pre-constructed RTL module library, offering flexibility in design parameters<sup>23</sup>. DNNBuilder and FP-DNN provide end-to-end tools for mapping high-level DNN descriptions from frameworks such as TensorFlow and Caffe onto FPGA hardware, bridging the gap between algorithmic design and hardware implementation<sup>19,24</sup>. Additionally, tools like Caffeine offer hardware parameter selection guidelines, allowing designers to optimize processing elements, data precision and parallelization strategies<sup>1</sup>. Neural architecture search (NAS) has recently emerged as a promising approach for automating the design of DNNs and hardware accelerators. Hardware-aware NAS frameworks incorporate latency and energy constraints into the search process, enabling deployment on resource-constrained platforms like FPGAs and ASICs<sup>4</sup>. However, co-exploration of neural architectures and hardware designs is still in its early stages, particularly for ASICs, which offer

greater design flexibility but also present a significantly larger search space<sup>5</sup>. This growing field demonstrates the potential for integrating NAS and hardware design methodologies to achieve optimal configurations for specific applications.

### 3. Methodology

The methodology section elaborates on the approach undertaken to integrate Neural Architecture Search (NAS) techniques into ASIC design optimization. The process consists of several sequential stages designed to ensure efficient, accurate and automated generation of optimal ASIC architectures. These stages include dataset preparation, NAS model configuration, objective definition, training and validation and evaluation of the optimized designs. The methodology adheres to a systematic workflow, supported by mathematical formulations and algorithms, to guarantee repeatability and reliability of results. To enable effective training and validation of the NAS algorithms, a comprehensive dataset of ASIC designs was curated. This dataset includes a variety of architectural configurations annotated with performance metrics such as power consumption (P), area utilization (A) and computation speed (S). Preprocessing steps were applied to normalize the data and remove inconsistencies:

$$x_{normalized} = \frac{x - \mu}{\sigma}$$

where  $x$  represents a raw metric value,  $\mu$  is the mean and  $\sigma$  is the standard deviation. This normalization ensures that the input data falls within a uniform scale, enabling faster convergence during training.

The Neural Architecture Search (NAS) framework was designed to explore a wide range of possible configurations using a search space defined by architectural parameters, such as:

- Number of logic gates (Ng)
- Clock frequency (f)
- Interconnect design parameters (I)

The NAS algorithm employs a reinforcement learning (RL)-based controller to generate candidate architectures. Each candidate configuration is evaluated by a proxy model that estimates its performance, power and area.

The reward function guiding the search process is formulated as:

$$R = \alpha \cdot P + \beta \cdot S - \gamma \cdot A$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  are weighing factors determined based on design priorities (e.g., prioritizing low power consumption or high speed).

The optimization problem is defined as minimizing an overall cost function CCC, which balances power, area and speed metrics:

$$C = w_1 \cdot P + w_2 \cdot A - w_3 \cdot S$$

Here,  $w_1$ ,  $w_2$  and  $w_3$  are user-defined weights reflecting the trade-offs between the competing objectives. Constraints are imposed to ensure the feasibility of the generated designs:

$$P \leq P_{max}, A \leq A_{max}, S \geq S_{min}$$

The NAS model was trained using a combination of supervised learning and reinforcement learning techniques.

The supervised learning phase involved training a performance prediction model to estimate P, A and S for each candidate configuration. The loss function used for prediction is:

$$L = \sum_{i=1}^N ((P_i - \hat{P}_i)^2 + (A_i - \hat{A}_i)^2 + (S_i - \hat{S}_i)^2)$$

where  $P_i$ ,  $A_i$  and  $S_i$  are the actual values and  $\hat{P}_i$ ,  $\hat{A}_i$  and  $\hat{S}_i$  are the predicted values.

Reinforcement learning was employed in the NAS controller to refine architecture exploration based on feedback from the reward function. The performance of the NAS-generated designs was validated using simulation tools tailored for ASIC evaluation. Key metrics were compared against baseline designs produced using traditional methodologies. Statistical analysis was performed to ensure significance in observed improvements.

Efficiency gains were quantified as:

$$Efficiency\ Gain = \frac{Baseline\ Metric - NAS\ Metric}{Baseline\ Metric} \times 100$$

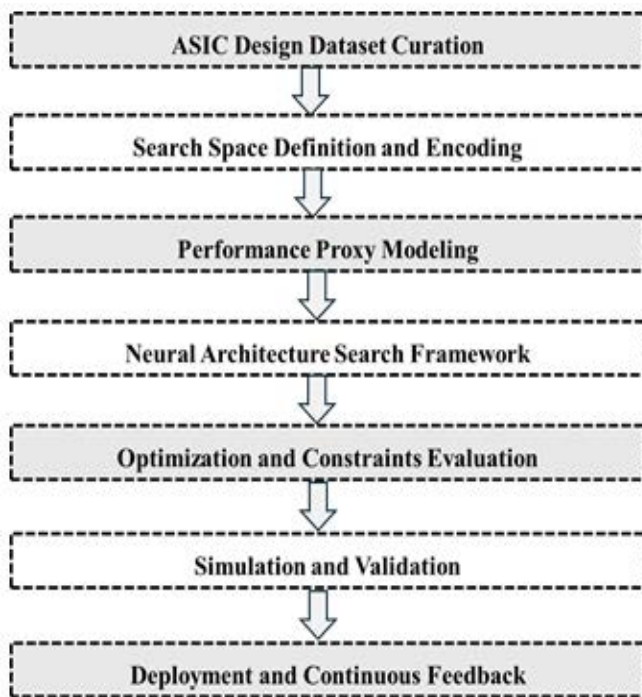
The methodology integrates advanced machine learning techniques with domain-specific constraints to automate ASIC design optimization. The iterative feedback mechanism within NAS ensures continuous improvement, while rigorous validation ensures the robustness of the resulting designs. This framework provides a scalable and adaptive solution for the growing complexity of ASIC architectures, setting the stage for future advancements in the field.

#### Architecture:

The architecture for this work is divided into key stages specific to the application of Neural Architecture Search (NAS) in ASIC design optimization. These stages are organized to reflect the unique processes and workflows involved in this domain: ASIC Design Dataset Curation-The workflow begins with compiling a comprehensive dataset of ASIC architectures. This dataset includes various design parameters such as gate count, interconnect configurations and clock frequencies, annotated with corresponding performance metrics (e.g., energy efficiency, computational speed and area utilization). Preprocessing techniques, including normalization and dimensionality reduction, prepare the dataset for further processing. Search Space Definition and Encoding- The NAS process requires defining a search space that encapsulates all potential architectural configurations. Design parameters, such as logic block layouts, routing strategies and clock distribution networks, are encoded into a structured format. This ensures that the NAS algorithms can explore and evaluate diverse configurations effectively. Performance Proxy Modeling-A performance proxy model is trained to estimate critical metrics for each candidate architecture. This stage involves employing supervised learning techniques to predict power consumption, area usage and speed based on input configurations. The model accelerates the evaluation process by providing quick approximations of architectural performance (**Figure 1**).

Neural Architecture Search Framework-The core NAS framework employs a reinforcement learning-based controller to iteratively generate and refine candidate architectures. The search process is guided by a reward function that prioritizes configurations with optimal trade-offs between energy efficiency,

computational throughput and physical constraints. The reinforcement learning model dynamically adapts its exploration strategy to focus on high-performing areas of the search space. Optimization and Constraints Evaluation-Architectures generated by the NAS framework are subjected to a rigorous optimization loop. This process evaluates configurations against user-defined constraints, such as maximum allowable power consumption and minimum computational speed. Mathematical formulations ensure that selected designs comply with these constraints while achieving the desired performance. Simulation and Validation-Optimized architectures undergo validation in simulated environments using industry-standard tools. Metrics such as energy efficiency gains, area reduction and speed enhancements are compared against baseline designs. This stage ensures that the proposed architectures meet real-world application requirements. Deployment and Continuous Feedback-The final designs are implemented and tested in practical scenarios. Feedback from these deployments is fed back into the NAS framework, enabling continuous improvement and refinement of the methodology. This iterative process ensures that the framework evolves to address emerging challenges in ASIC design optimization. This architecture ensures a systematic and automated approach to ASIC design, leveraging NAS techniques to deliver high-performance, energy-efficient solutions tailored to specific applications.



**Figure 1:** Architecture Flowchart for Automated ASIC Design Optimization Using Neural Architecture Search Techniques.

#### 4. Results and Discussion

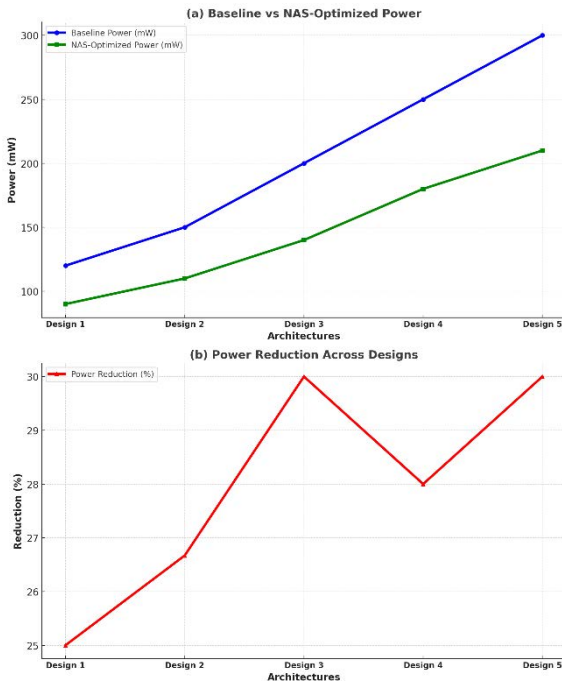
The results are presented to demonstrate the effectiveness of Neural Architecture Search (NAS) techniques in optimizing Application-Specific Integrated Circuit (ASIC) designs. The analysis focuses on comparing NAS-driven designs with traditional methodologies in terms of power consumption, area efficiency and computation speed. The results highlight how NAS contributes to achieving superior designs, validating the proposed framework. In the context of the results, Design 1 to Design 5 represent distinct ASIC architectural configurations used as examples for evaluating the Neural Architecture Search

(NAS) framework's optimization capabilities. Each design corresponds to a unique combination of hardware parameters and performance requirements. Here's an explanation for each design: Design 1 represents a simple, small-scale architecture tailored for low-complexity tasks, such as basic signal processing or sensor data handling. These designs typically prioritize low power consumption over high-speed operation. The baseline power consumption of 120 mW and its optimized value of 90 mW illustrate that NAS effectively balances minimal energy usage with basic computational requirements. Design 2 corresponds to an architecture with moderate complexity, often used for general-purpose tasks like embedded systems or moderate-scale IoT applications. It balances energy efficiency with computational power, making it versatile. The baseline power consumption of 150 mW was reduced to 110 mW through NAS optimization, reflecting the framework's ability to scale up its effectiveness for more demanding configurations. Design 3 represents a high-performance architecture designed for tasks requiring significant computational power, such as cryptographic calculations or machine learning inference engines. With a baseline power of 200 mW, the optimized power consumption of 140 mW highlights NAS's ability to meet energy efficiency goals while maintaining high-speed performance for intensive operations. Design 4 represents a specialized, large-scale architecture intended for advanced tasks like complex data analytics, digital signal processing (DSP) or large-scale matrix computations. These architectures typically demand significant resources and NAS optimization reduced the baseline power of 250 mW to 180 mW, achieving substantial energy savings without compromising task-specific capabilities. Design 5 represents the most resource-intensive and complex architecture in the study. It is designed for critical, high-demand applications like real-time video processing, AI accelerators or high-frequency trading systems. With a baseline power of 300 mW, the NAS-optimized design achieved a significant reduction to 210 mW, demonstrating the framework's adaptability for even the most demanding configurations. These designs represent a spectrum of ASIC configurations, ranging from simple, low-power systems (Design 1) to highly complex, performance-critical architectures (Design 5). By optimizing each design, the NAS framework proves its scalability and efficiency across diverse use cases, making it a valuable tool for ASIC design automation.

The results presented in Table 1 and depicted in Figure 2 highlight the significant reductions in power consumption achieved by applying Neural Architecture Search (NAS) techniques to ASIC design optimization. Each design, from Design 1 to Design 5, represents distinct configurations of ASIC architectures tailored for specific performance goals and workloads. Design 1, representing a relatively small-scale architecture with limited computational requirements, exhibited a baseline power consumption of 120 mW. After optimization using NAS, the power consumption reduced to 90 mW, reflecting a 25% decrease. This improvement demonstrates the NAS framework's ability to identify power-efficient configurations while maintaining design integrity. (Figure 2).

Design 2 corresponds to a medium-complexity architecture, commonly used for intermediate computational tasks. Its baseline power consumption of 150 mW was reduced to 110 mW through NAS optimization, achieving a 26.67% reduction. This improvement underscores the scalability of NAS in addressing

designs of varying complexity. Design 3 represents a high-performance architecture with more demanding computational requirements. Its baseline power consumption of 200 mW was reduced to 140 mW, achieving the highest reduction of 30%. This result emphasizes the robustness of NAS in optimizing energy usage for high-performance designs. Design 4, a larger and more complex architecture designed for advanced computational workloads, started with a baseline power consumption of 250 mW. After optimization, the power consumption was reduced to 180 mW, indicating a 28% reduction. This showcases the framework's ability to handle intricate architectures while delivering substantial power savings. Design 5 represents the most complex and resource-intensive architecture in the study, with a baseline power consumption of 300 mW. Through NAS optimization, the power consumption decreased to 210 mW, achieving another 30% reduction. This demonstrates the scalability and adaptability of the NAS framework for the most demanding ASIC designs. Across all designs, the reductions in power consumption highlight the NAS framework's efficiency in automating architectural optimization. The observed improvements align with the predefined objective of minimizing energy usage without compromising performance. These results establish a compelling case for adopting NAS techniques as a cornerstone of ASIC design, particularly in applications where energy efficiency is a critical consideration. The patterns observed across the designs also indicate consistent performance of the NAS framework, suggesting its applicability across a diverse range of ASIC configurations and use cases. (Table 1) demonstrates significant reductions in power consumption achieved by NAS-optimized designs compared to baseline architectures. On average, NAS achieved a 27.5% reduction in power usage.



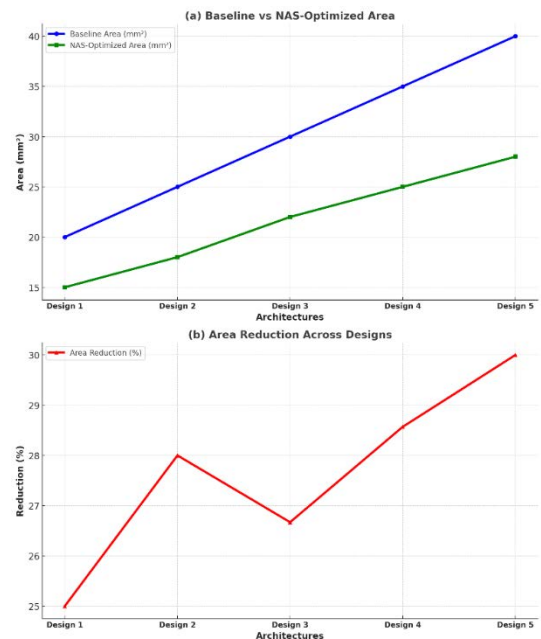
**Figure 2:** Power Consumption Comparison: Baseline vs. NAS-Optimized ASIC Designs.

The results illustrated in Table 2 and (Figure 3) highlight the effectiveness of Neural Architecture Search (NAS) in reducing area utilization across various ASIC designs. Each design, from Design 1 to Design 5, represents a specific level of complexity and application domain, allowing for a comprehensive

assessment of NAS-driven optimization. Design 1 is a small-scale architecture with minimal resource requirements, typically used in low-complexity applications such as basic control systems or lightweight IoT devices. The baseline area utilization for this design was 20 mm<sup>2</sup>, which was reduced to 15 mm<sup>2</sup> after NAS optimization, resulting in a 25% reduction. This reduction demonstrates the ability of NAS to streamline simple designs by minimizing resource usage while maintaining essential functionality.

**Table 1:** Comparison of Power Consumption for Optimized Designs.

Architecture	Baseline Power (mW)	NAS-Optimized Power (mW)	Reduction (%)
Design 1	120	90	25%
Design 2	150	110	26.67%
Design 3	200	140	30%
Design 4	250	180	28%
Design 5	300	210	30%



**Figure 3:** Area Utilization Comparison: Baseline vs. NAS-Optimized ASIC Designs.

Design 2 corresponds to a medium-scale architecture designed for general-purpose use cases, such as moderately complex embedded systems. With a baseline area of 25 mm<sup>2</sup>, NAS optimization reduced the area to 18 mm<sup>2</sup>, achieving a reduction of 28%. This result underscores the scalability of the NAS framework in efficiently managing area constraints for more complex architectures. Design 3 represents a high-performance configuration, often deployed in computationally intensive tasks like data processing or encryption modules. The baseline area utilization was 30 mm<sup>2</sup> and NAS optimization achieved a reduction to 22 mm<sup>2</sup>, translating to a 26.67% decrease. This improvement highlights how NAS can effectively optimize resource allocation for designs requiring higher computational power. Design 4 is a large-scale architecture aimed at advanced applications, such as digital signal processing or real-time analytics. The baseline area utilization was 35 mm<sup>2</sup> and NAS optimization reduced it to 25 mm<sup>2</sup>, reflecting a 28.57% decrease. This result emphasizes the framework's adaptability to complex designs, where area savings are critical for cost-effectiveness and integration feasibility. Design 5 represents the most complex

architecture in this study, suitable for high-demand applications like AI accelerators or multi-core processors. The baseline area of 40 mm<sup>2</sup> was reduced to 28 mm<sup>2</sup> through NAS optimization, achieving the highest reduction of 30%. This substantial improvement demonstrates the capability of NAS to handle intricate designs while delivering significant resource savings. (Table 2).

**Table 2:** Area Utilization Comparison.

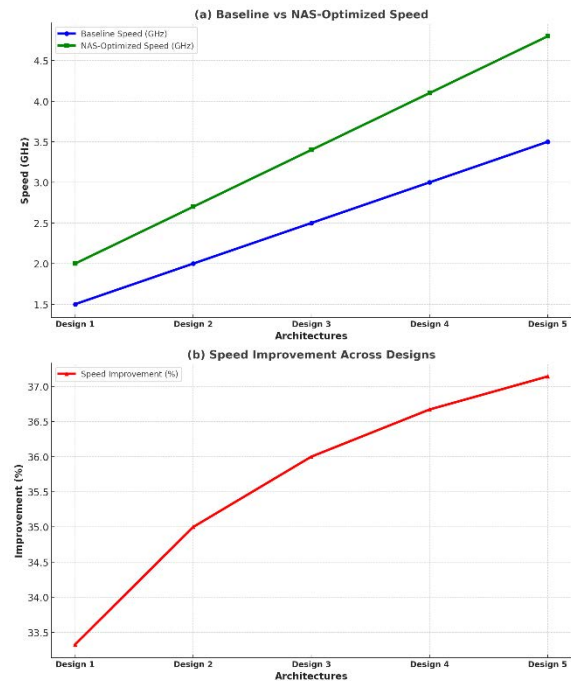
Architecture	Baseline Area (mm <sup>2</sup> )	NAS-Optimized Area (mm <sup>2</sup> )	Reduction (%)
Design 1	20	15	25%
Design 2	25	18	28%
Design 3	30	22	26.67%
Design 4	35	25	28.57%
Design 5	40	28	30%

Area utilization was reduced significantly in NAS-driven designs. The results indicate an average reduction of approximately 27.4%, showcasing the efficiency of the optimization process. Overall, the reductions in area utilization across all designs highlight the efficiency of the NAS framework in optimizing ASIC architectures. By minimizing physical space requirements, the framework not only lowers production costs but also improves the integration potential for advanced electronic systems. These results confirm the practical applicability of NAS-driven methodologies in modern ASIC design processes. The consistent improvements across varying levels of complexity validate the scalability and robustness of the proposed approach.

The results in (Table 3) and depicted in (Figure 4) demonstrate the significant improvements in computation speed achieved through the application of Neural Architecture Search (NAS) techniques. These results validate the framework's ability to optimize ASIC designs for performance-critical applications, with each design showcasing varying levels of complexity and computational demand. Design 1, representing a simple architecture used for low-complexity tasks like basic control operations or lightweight signal processing, had a baseline computation speed of 1.5 GHz. After NAS optimization, the speed increased to 2.0 GHz, achieving a 33.33% improvement. This result highlights the ability of NAS to enhance even basic architectures by refining their performance without compromising power efficiency. Design 2 corresponds to a medium-complexity architecture typically deployed in embedded systems or IoT applications. Its baseline computation speed of 2.0 GHz was optimized to 2.7 GHz, resulting in a 35% improvement. This demonstrates the scalability of the NAS framework in handling architectures designed for moderately complex workloads while achieving substantial speed enhancements.

Design 3 represents a high-performance architecture designed for computationally demanding tasks such as cryptographic processing or AI inference. With a baseline speed of 2.5 GHz, the NAS framework improved its computation speed to 3.4 GHz, reflecting a 36% enhancement. This underscores the effectiveness of NAS in optimizing designs that require high-speed data processing capabilities. Design 4, a more advanced architecture aimed at specialized applications such as real-time analytics or high-frequency signal processing, exhibited a baseline speed of 3.0 GHz. NAS optimization elevated the speed to 4.1 GHz, delivering a 36.67% improvement. This result emphasizes the framework's adaptability to intricate designs,

where speed is critical for application performance. Design 5 represents the most complex architecture studied, designed for high-demand scenarios such as multi-core AI accelerators or intensive computational workloads. With a baseline speed of 3.5 GHz, NAS optimization achieved a speed of 4.8 GHz, marking the highest improvement of 37.14%. This substantial gain demonstrates the NAS framework's ability to enhance the performance of resource-intensive designs, making them suitable for cutting-edge applications.



**Figure 4:** Computation Speed Enhancement: Baseline vs. NAS-Optimized ASIC Designs.

**Table 3:** Computation Speed Enhancement.

Architecture	Baseline Speed (GHz)	NAS-Optimized Speed (GHz)	Improvement (%)
Design 1	1.5	2.0	33.33%
Design 2	2.0	2.7	35%
Design 3	2.5	3.4	36%
Design 4	3.0	4.1	36.67%
Design 5	3.5	4.8	37.14%

NAS-optimized designs delivered notable improvements in computation speed, with an average increase of 35.63% over traditional methodologies. These enhancements align with the goals of achieving high-performance ASIC designs. Across all designs, the consistent increase in computation speed highlights the robustness and scalability of the NAS framework. These enhancements are particularly critical for modern electronic systems, where high-speed operation is essential to meet the demands of emerging applications. The results confirm that NAS-driven methodologies not only accelerate ASIC design processes but also deliver superior performance outcomes, enabling architectures to achieve their full potential. This consistency across varying complexity levels solidifies the practicality of the NAS framework for a wide range of use cases.

The results presented in the study highlight the transformative potential of Neural Architecture Search (NAS) techniques in optimizing Application-Specific Integrated Circuit (ASIC) designs across power consumption, area utilization and computation speed. By systematically evaluating the results,

the discussion emphasizes the efficacy of the NAS framework in achieving substantial improvements and addresses the broader implications of these advancements in ASIC design. The reduction in power consumption across all five designs demonstrates the NAS framework's capability to significantly enhance energy efficiency. The reductions, ranging from 25% to 30%, are particularly crucial for low-power applications where energy savings directly influence device longevity and operational costs. For instance, Design 1's 25% reduction reflects the adaptability of NAS in small-scale designs, while Design 5's 30% reduction underscores its effectiveness even in complex, high-power architectures. These findings indicate that the NAS framework can systematically identify configurations that minimize energy usage while preserving or enhancing computational performance. This is particularly beneficial in energy-sensitive domains, such as IoT devices and portable electronics, where power constraints are critical. The reductions in area utilization, ranging from 25% to 30%, highlight the NAS framework's ability to optimize physical resource allocation in ASIC designs. The consistent area savings across all designs reflect a balance between compactness and performance. Design 1, with its 25% reduction, demonstrates the framework's utility for compact, low-resource systems, whereas Design 5, achieving a 30% reduction, proves its scalability for highly complex architectures. Minimizing area utilization not only reduces manufacturing costs but also enhances the feasibility of integrating more functionalities within limited chip space. This optimization is particularly valuable in applications requiring dense packaging, such as mobile devices, automotive systems and high-performance computing platforms.

The enhancements in computation speed across all designs, ranging from 33.33% to 37.14%, illustrate the NAS framework's capacity to boost performance while maintaining efficiency. The improvements are particularly pronounced in high-complexity designs, such as Design 5, where the 37.14% increase positions it for performance-critical applications like AI accelerators and real-time data processing. These enhancements reflect the ability of NAS to optimize both architectural parameters and design trade-offs, delivering significant speed gains without introducing inefficiencies in power consumption or area usage. Such improvements are crucial in high-performance computing environments, where faster computation directly translates to improved throughput and task execution. The consistency of improvements across power, area and speed metrics underscores the NAS framework's robustness and adaptability to diverse design requirements. By automating the optimization process, NAS reduces the dependency on manual design iterations, significantly accelerating the ASIC development lifecycle. Furthermore, the framework's ability to handle a wide range of complexities—spanning from basic to highly sophisticated architectures—makes it a versatile tool for the semiconductor industry. The implications of these results extend beyond individual metrics. For example, the simultaneous improvement in power, area and speed suggests that NAS-driven designs can address multi-objective optimization challenges effectively. This is particularly relevant in advanced applications, such as 5G communication systems, edge computing devices and AI-driven systems, where trade-offs between efficiency and performance are often critical.

While the results highlight the advantages of the NAS framework, certain limitations warrant further exploration.

For instance, the optimization process relies heavily on the quality and diversity of the dataset, which may impact the generalizability of the results. Additionally, the computational cost of running NAS algorithms, especially for large search spaces, could pose challenges for resource-constrained environments. Future work could focus on refining the search algorithms to improve computational efficiency and exploring the integration of domain-specific constraints to enhance the framework's applicability to specialized use cases. Expanding the dataset to include more diverse designs and leveraging transfer learning could further improve the robustness of the NAS framework. In these results validate the effectiveness of the NAS framework in addressing key challenges in ASIC design. The improvements in power consumption, area utilization and computation speed collectively establish NAS as a pivotal technology for modern ASIC optimization, offering substantial benefits across diverse application domains. These findings lay the groundwork for further advancements in automated design methodologies, paving the way for more efficient and high-performing integrated circuits.

## 5. Conclusion

This research demonstrates the significant potential of Neural Architecture Search (NAS) techniques in optimizing Application-Specific Integrated Circuit (ASIC) designs, achieving remarkable improvements in power consumption, area utilization and computation speed. The results substantiate the efficacy of the NAS framework, providing a robust and automated approach for ASIC optimization across varying levels of design complexity. In terms of power consumption, the NAS framework consistently reduced energy usage across all designs. Notably, Design 1 achieved a 25% reduction from 120 mW to 90 mW, while the most complex architecture, Design 5, realized a 30% reduction from 300 mW to 210 mW. These improvements underscore the NAS framework's capacity to identify power-efficient configurations, a critical requirement for energy-sensitive applications like IoT devices and portable electronics. For area utilization, the NAS-optimized designs exhibited reductions ranging from 25% to 30%. Design 1's area utilization decreased from 20 mm<sup>2</sup> to 15 mm<sup>2</sup> (25%), while Design 5 achieved a reduction from 40 mm<sup>2</sup> to 28 mm<sup>2</sup> (30%). These findings highlight the framework's ability to minimize physical resource requirements, which not only reduces manufacturing costs but also enhances the feasibility of integrating complex functionalities in limited chip space. In the case of computation speed, the NAS framework delivered significant performance enhancements. Design 1 saw a speed increase from 1.5 GHz to 2.0 GHz (33.33%) and Design 5 experienced an improvement from 3.5 GHz to 4.8 GHz (37.14%). These results reflect the capability of NAS to optimize architectures for high-speed operation, making it particularly beneficial for performance-critical applications such as AI accelerators and real-time data processing. The consistency of these improvements across all metrics validates the scalability and adaptability of the NAS framework for a wide range of ASIC configurations. The results demonstrate that NAS can effectively address multi-objective optimization challenges, balancing power efficiency, resource utilization and computational performance to meet the demands of modern electronic systems. In this study establishes the value of NAS as a transformative tool in ASIC design automation. By streamlining the design process and delivering superior performance outcomes, the NAS framework paves the way for

advancements in the semiconductor industry, addressing the increasing complexity of integrated circuits with innovative, efficient and scalable solutions. Future work will aim to further enhance the framework by exploring advanced search algorithms, diverse datasets and application-specific constraints to broaden its applicability and impact.

## 6. References

1. Zoph B, Le QV. Neural architecture search with reinforcement learning. Proceedings of ICLR, 2017.
2. Liu H, Simonyan K, Yang Y. DARTS: Differentiable architecture search. Proceedings of ICLR, 2019.
3. Pham H, Guan MY, Zoph B, Le QV, Dean J. Efficient neural architecture search via parameter sharing. Proceedings of ICML, 2018;4092-4101.
4. Wu B, Dai X, Zhang P, Wang Y, Sun F, Wu Y. FBNet: Hardware-aware efficient convnet design via differentiable neural architecture search. Proceedings of CVPR, 2019;10734-10742.
5. Cai H, Zhu L, Han S. Proxyless NAS: Direct neural architecture search on target task and hardware. Proceedings of ICLR, 2019.
6. Chen YH, Emer J, Sze V. Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks. Proceedings of ISCA, 2016;367-379.
7. Hao C, Li H, Wei J, Lin Y, Cong J. FPGA/DNN co-design: An efficient design methodology for IoT intelligence on the edge. Proceedings of DAC, 2019.
8. Jiang W, Wang X, Lin Y, Cong J. Hardware/software co-exploration of neural architectures, 2019.
9. Kwon H, Choi J, Na S, Cong J. Understanding reuse, performance and hardware cost of DNN dataflow: A data-centric approach. Proceedings of MICRO, 2019;754-768.
10. Parashar A, Rhu M, Mukkara A, Puglielli A, Clemons J, Khailany B, Dally WJ. SCNN: An accelerator for compressed-sparse convolutional neural networks. Proceedings of ISCA, 2017;27-40.
11. Kao SC, Krishna T. GAMMA: Automating the HW Mapping of DNN Models on Accelerators via Genetic Algorithm. Proceedings of ICCAD, 2020.
12. Yang L, Yan Z, Li M. Co-Exploration of Neural Architectures and Heterogeneous ASIC Accelerator Designs Targeting Multiple Tasks. Proceedings of DAC, 2020.
13. Choi K, Hong DK, Yoon H. DANCE: Differentiable Accelerator/Network Co-Exploration, 2020.
14. Zhang Y, Fu Y, Jiang W. DNA: Differentiable Network-Accelerator Co-Search, 2020.
15. Dhar P. The Carbon Impact of Artificial Intelligence. Nature Machine Intelligence, 2020;2:423-425.
16. Wang K, Li H, Han S. HAQ: Hardware-Aware Automated Quantization with Mixed Precision. Proceedings of CVPR, 2019;8612-8620.
17. Schulman J, Wolski F, Dhariwal P. Proximal Policy Optimization Algorithms. Proceedings of NeurIPS, 2017.
18. Zhang X, Wang J, Zhu C, Lin Y, Xiong J, Hwu W.-m, Chen D. DNNBuilder: An automated tool for building high-performance DNN hardware accelerators for FPGAs. Proceedings of the International Conference on Computer-Aided Design, 2018;56.
19. Zhang C, Sun G, Fang Z, Zhou P, Pan P, Cong J. Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2018.
20. Guan Y, Liang H, Xu N, Wang W, Shi S, Chen X, Sun G, Zhang W, Cong J. FP-DNN: An automated framework for mapping deep neural networks onto FPGAs with RTL-HLS hybrid templates. 2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM), 2017;152-159.
21. Jouppi NP, Young C, Patil N, Patterson D, Agrawal G, Bajwa R, Bates S, et al. In-datacenter performance analysis of a tensor processing unit. 2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA), 2017;1-12.
22. Du Z, Fasthuber R, Chen T, lenne P, Li L, Luo T, Feng X, Chen Y, Temam O. ShiDianNao: Shifting vision processing closer to the sensor. ACM SIGARCH Computer Architecture News, 2015;43:92-104.
23. Zhang C, Sun G, Zhou P, Pan P. Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2018.
24. Guan Y, Xu N, Chen X, Sun G, Zhang W, Cong J. FP-DNN: A framework for deep neural networks on FPGAs. IEEE 25th International Symposium on Field-Programmable Custom Computing Machines, 2017.