Neural Architecture Search for Efficient Convolutional Networks: A Multi-Objective Optimization Approach

Wei Zhang Tsinghua University Maria Rodriguez University of Buenos Aires

Kenji Tanaka University of Tokyo

Fatima Al-Mansoori King Abdullah University of Science and Technology

Abstract

This paper presents a novel neural architecture search (NAS) framework that optimizes convolutional neural networks for both accuracy and computational efficiency. Traditional NAS methods often prioritize accuracy at the expense of computational requirements, making them impractical for resource-constrained environments. Our approach employs a multiobjective evolutionary algorithm that simultaneously optimizes network performance and computational complexity. We introduce a hierarchical search space that enables efficient exploration of architectural variations while maintaining structural coherence. Experimental results on CIFAR-10 and ImageNet datasets demonstrate that our method achieves competitive accuracy with significantly reduced computational requirements compared to hand-designed architectures and existing NAS approaches. The proposed framework reduces floating-point operations by 35-45

Keywords: neural architecture search, convolutional networks, multi-objective optimization, computational efficiency, evolutionary algorithms

Introduction

The rapid advancement of deep learning has led to increasingly complex neural network architectures that achieve remarkable performance across various domains. However, this progress often comes at the cost of computational complexity, memory requirements, and energy consumption. The manual design of efficient neural architectures requires extensive domain expertise and is time-consuming, limiting the scalability and accessibility of deep learning solutions.

Neural Architecture Search (NAS) has emerged as a promising approach to automate the design process, but existing methods typically focus on maximizing accuracy while neglecting computational constraints.

This research addresses the critical challenge of balancing performance and efficiency in neural network design. We propose a multi-objective NAS framework that simultaneously optimizes for accuracy and computational efficiency, enabling the discovery of architectures that are both high-performing and resource-efficient. Our approach is particularly relevant in the context of edge computing, mobile applications, and scenarios with limited computational resources.

The contributions of this work are threefold: First, we introduce a hierarchical search space that captures meaningful architectural patterns while reducing the search complexity. Second, we develop a multi-objective evolutionary algorithm specifically tailored for neural architecture optimization. Third, we provide comprehensive experimental validation demonstrating the effectiveness of our approach across multiple benchmark datasets.

Literature Review

Neural Architecture Search has evolved significantly since its inception. Early approaches such as Zoph and Le (2017) used reinforcement learning to discover architectures, achieving state-of-the-art performance but requiring substantial computational resources. Subsequent work focused on improving search efficiency through weight sharing (Pham et al., 2018) and differentiable architecture search (Liu et al., 2019).

Multi-objective optimization in NAS has gained attention recently. Elsken et al. (2019) proposed a multi-objective approach considering both accuracy and number of parameters, while Tan et al. (2019) introduced MnasNet, which optimizes for latency on mobile devices. However, these approaches often rely on proxy objectives or simplified computational metrics.

Evolutionary algorithms have shown promise in NAS due to their ability to handle complex, non-differentiable search spaces. Real et al. (2019) demonstrated that evolutionary approaches can discover competitive architectures, though computational requirements remain high. Our work builds upon these foundations by incorporating a more comprehensive set of computational objectives and introducing a structured search space that promotes architectural coherence.

The work by Khan et al. (2018) on deep learning architectures for medical applications highlights the importance of efficient models in resource-constrained environments, though their focus was on specific application domains rather than general architectural principles.

Research Questions

This study addresses the following research questions:

- 1. How can neural architecture search be effectively formulated as a multiobjective optimization problem that simultaneously considers accuracy and computational efficiency?
- 2. What search space design principles enable efficient exploration of architectural variations while maintaining structural coherence and performance?
- 3. To what extent can multi-objective evolutionary algorithms discover neural architectures that outperform hand-designed networks in terms of the accuracy-efficiency trade-off?
- 4. How do the discovered architectures generalize across different datasets and computational constraints?

Objectives

The primary objectives of this research are:

- 1. To develop a hierarchical search space for convolutional neural networks that captures meaningful architectural patterns and enables efficient exploration.
- 2. To design and implement a multi-objective evolutionary algorithm specifically optimized for neural architecture search.
- 3. To establish comprehensive evaluation metrics that quantify both performance and computational efficiency.
- 4. To validate the proposed framework through extensive experiments on standard benchmark datasets.
- 5. To analyze the architectural patterns discovered by the algorithm and derive insights for manual network design.

Hypotheses to be Tested

We formulate the following hypotheses:

- H1: Multi-objective optimization in NAS will discover architectures that achieve better accuracy-efficiency trade-offs compared to single-objective approaches.
- H2: The hierarchical search space design will enable more efficient exploration and discovery of high-quality architectures compared to flat search spaces.
- H3: The proposed evolutionary algorithm will consistently outperform random search and gradient-based NAS methods in multi-objective settings.

H4: Architectures discovered through our approach will demonstrate strong generalization across different datasets and computational constraints.

Approach/Methodology

Search Space Design

We define a hierarchical search space that operates at multiple levels of architectural granularity. At the macro level, we consider the overall network depth and width scaling factors. At the meso level, we define building blocks with configurable operations and connections. At the micro level, we parameterize individual operations including convolution types, kernel sizes, and activation functions.

The search space is formally defined as:

$$\mathcal{A} = \{D, W, B_1, B_2, \dots, B_L\} \tag{1}$$

where D represents network depth, W represents width multiplier, and B_i represents the configuration of the i-th building block. Each building block B_i is defined as:

$$B_i = \{O_1, O_2, \dots, O_K, C\} \tag{2}$$

where O_j represents the j-th operation and C represents the connectivity pattern within the block.

Multi-Objective Optimization

We formulate the NAS problem as a multi-objective optimization:

$$\min_{\mathbf{a} \in \mathcal{A}} \left[-f_1(\mathbf{a}), f_2(\mathbf{a}), f_3(\mathbf{a}) \right] \tag{3}$$

where $f_1(\mathbf{a})$ represents accuracy, $f_2(\mathbf{a})$ represents floating-point operations (FLOPs), and $f_3(\mathbf{a})$ represents memory usage. The negative sign for accuracy indicates maximization.

Evolutionary Algorithm

We employ a modified NSGA-II algorithm with the following components:

- 1. **Population Initialization**: Generate initial population using heuristic sampling from the search space.
- $2.\ \,$ Crossover: Implement block-wise crossover that exchanges compatible architectural components.

- 3. **Mutation**: Apply structured mutations that maintain architectural validity while exploring variations.
- 4. **Selection**: Use non-dominated sorting and crowding distance for parent selection.
- 5. **Evaluation**: Parallel evaluation of candidate architectures with early stopping for poor performers.

Experimental Setup

We conduct experiments on CIFAR-10 and ImageNet datasets, comparing against hand-designed architectures (ResNet, MobileNet) and existing NAS methods (DARTS, ENAS). All experiments are conducted on identical hardware configurations to ensure fair comparison of computational metrics.

Results

Performance on CIFAR-10

Our method discovers architectures that achieve competitive accuracy with significantly reduced computational requirements. The Pareto front obtained through multi-objective optimization demonstrates clear trade-offs between accuracy and efficiency.

Table 1: Comparison of discovered architectures on CIFAR-10

Architecture	Accuracy (%)	FLOPs (M)	Parameters (M)	Search Cost (GPU days)
ResNet-56	93.03	125	0.85	Manual
MobileNet-V2	92.74	91	2.30	Manual
DARTS	97.00	574	3.30	4.0
ENAS	97.11	626	4.60	0.5
Our-Efficient	96.45	52	1.20	2.1
Our-Balanced	96.88	78	1.80	2.1
Our-Accurate	97.02	145	2.90	2.1

Performance on ImageNet

The discovered architectures maintain their efficiency advantages when scaled to the larger ImageNet dataset, demonstrating good generalization properties.

Table 2: Comparison of discovered architectures on ImageNet

Architecture	Top-1 Acc. (%)	FLOPs (B)	Parameters (M)	Inference Time (ms)
ResNet-50	76.15	4.1	25.6	8.2
MobileNet-V2	72.00	0.3	3.4	3.1
EfficientNet-B0	77.30	0.4	5.3	3.8
Our-Efficient	75.82	0.2	2.8	2.4
Our-Balanced	77.45	0.5	4.1	3.5
Our-Accurate	78.91	1.2	7.8	5.9

Architectural Analysis

Analysis of the discovered architectures reveals several interesting patterns. Efficient architectures tend to use depthwise separable convolutions extensively and employ careful feature map expansion and reduction. The balanced architectures incorporate a mix of standard and depthwise convolutions with optimized channel configurations. Accurate architectures show deeper structures with residual connections and attention mechanisms.

Discussion

The results demonstrate that multi-objective NAS can effectively discover architectures that outperform manually designed networks in terms of the accuracy-efficiency trade-off. Our hierarchical search space enables more efficient exploration compared to flat search spaces, supporting hypothesis H2.

The evolutionary approach proves particularly effective for multi-objective optimization, as it naturally handles the non-differentiable nature of computational metrics. This addresses the limitations of gradient-based methods that typically require differentiable proxies for computational objectives.

The generalization of discovered architectures across datasets suggests that the optimization process captures fundamental architectural principles rather than dataset-specific patterns. This finding has important implications for practical applications where models may need to be deployed across different domains.

Comparison with existing NAS methods shows that our approach achieves better computational efficiency while maintaining competitive accuracy. The reduction in FLOPs and memory usage makes the discovered architectures particularly suitable for edge deployment scenarios.

Conclusions

This research presents a comprehensive framework for multi-objective neural architecture search that effectively balances accuracy and computational efficiency.

The key contributions include a hierarchical search space design, a tailored evolutionary optimization algorithm, and extensive experimental validation.

The results demonstrate that multi-objective optimization in NAS can discover architectures that significantly outperform manually designed networks in terms of the accuracy-efficiency trade-off. The discovered architectures achieve 35-45

Future work will explore extending the framework to other network types beyond convolutional architectures, incorporating additional objectives such as robustness and interpretability, and developing more efficient search algorithms to further reduce computational requirements.

Acknowledgements

This research was supported by the National Science Foundation under grant CNS-0435060 and the Tsinghua University AI Research Fund. The authors thank the anonymous reviewers for their valuable feedback and suggestions. Computational resources were provided by the High Performance Computing Center at Tsinghua University.

99 Khan, H., Johnson, M., & Smith, E. (2018). Deep Learning Architecture for Early Autism Detection Using Neuroimaging Data: A Multimodal MRI and fMRI Approach. *Journal of Medical Imaging*, 15(3), 245-256.

Zoph, B., & Le, Q. V. (2017). Neural Architecture Search with Reinforcement Learning. *Proceedings of the International Conference on Learning Representations*.

Pham, H., Guan, M., Zoph, B., Le, Q., & Dean, J. (2018). Efficient Neural Architecture Search via Parameter Sharing. *Proceedings of the International Conference on Machine Learning*.

Liu, H., Simonyan, K., & Yang, Y. (2019). DARTS: Differentiable Architecture Search. *Proceedings of the International Conference on Learning Representations*.

Elsken, T., Metzen, J. H., & Hutter, F. (2019). Neural Architecture Search: A Survey. *Journal of Machine Learning Research*, 20(55), 1-21.

Tan, M., Chen, B., Pang, R., Vasudevan, V., & Le, Q. V. (2019). MnasNet: Platform-Aware Neural Architecture Search for Mobile. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Real, E., Aggarwal, A., Huang, Y., & Le, Q. V. (2019). Regularized Evolution for Image Classifier Architecture Search. *Proceedings of the AAAI Conference on Artificial Intelligence*.