

AUTOVISIONNAS: SELF-ADAPTIVE NEURAL ARCHITECTURE SEARCH FOR AUTOMATED COMPUTER VISION MODEL DISCOVERY

Hashim Ali^{*1}, Aftab Ahmed², Umer Tanveer³, Kiran Falak Sher⁴, Sher Afghan⁵

^{*1,2,3,4,5}Department of Computer Science Abdul Wali Khan University Mardan.

DOI: <https://doi.org/10.5281/zenodo.17878129>

Keywords

Neural Architecture Search, AutoML, computer vision, automated design, deep learning, architecture optimization, self-adaptive systems

Article History

Received: 11 October 2025

Accepted: 21 November 2025

Published: 10 December 2025

Copyright @Author

Corresponding Author: *

Hashim Ali

Abstract

Manual design of novel neural network architectures for computer vision tasks is still a time-consuming and expertise-heavy procedure, making state-of-the-art models not accessible to domain experts. This paper presents AutoVisionNAS, a novel Neural Architecture Search (NAS) framework that automatically finds best computer vision architectures across various tasks and datasets. Our proposed method tackles three fundamental problems in the current NAS techniques including computational efficiency through Progressive Architecture Pruning (PAP), task adaptability by Dynamic Search Space Evolution (DSSE) and multi-objective optimization via Pareto-Efficient Architecture Ranking (PEAR). Our proposed AutoVisionNAS achieves 94.2% accuracy on ImageNet now and saves the search cost by up to 85% compared to previous NAS methods. Extensive experiments on eight benchmark datasets consistently show state-of-the-art accuracy, 6.8% improvement in average (vs.) an efficient human-designed architecture and also are up to 12x faster at discovering the best architecture than prior NAS methods. The system enables practitioners to automatically deploy models for image classification, object detection and semantic segmentation tasks without extensive architectural expertise, making cutting-edge computer vision more widely accessible.

INTRODUCTION

Deep learning, along with the subsequent broad adoption of these techniques by computer vision community, has produced revolutionary progress in image recognition, object detection and scene understanding (Dalal & Triggs 2005; Krizhevsky 2012; Girshick et al). However, the performance of these model-based systems relies significantly on well-designed neural network architectures, which are traditionally crafted through labor intensive manual experiments by domain experts [1]. This manual knowledge based design-step presents a fundamental bottleneck in developing and deploying computer vision systems, being highly-technical,

computationally expensive, and taking months of iterations to converge.

According to the latest market studies, the global AutoML market size is anticipated to reach \$14.8 billion by 2030, with NAS being one of the most rapidly growing segments at a CAGR of 42.3% [2]. This exponential adoption indicates a crucial requirement for automated model discovery methods that enable democratized access to top-tier computer vision capabilities, and at the same time save human development effort and computational budget.

Conventional neural architecture design often depends on human intuition, domain knowledge and extensive empirical evaluations. The recent

successes of pioneer architectures (e.g., AlexNet[7], ResNet[8] and EfficientNet [1]) demonstrated that careful attention and manual design can lead to a great success, albeit at the cost of significant development efforts for creating new architectures and demanding computational resources [3].

Neural Architecture Search (NAS) has emerged as a promising solution to automate the architecture design effort, applying algorithmic search through large search spaces in order to discover effective network configurations [4]. Initial NAS algorithms have shown the promise of automatic architecture discovery, even performing better than manually designed networks on benchmark tasks. Unfortunately, the current works have three significant drawbacks that severely limit their practicality.

Computational efficiency is still a bottleneck for future development. Conventional NAS involves training hundreds or thousands of candidate AOs from scratch which consume heavy computational resources up to 22,000 GPU-hours in a single search[5]. This computational overhead renders the NAS approach unaffordable for many researchers and developers, which in turn refrains wide spread adoption of it in practices.

Second, Considerable NAS methods work with predefined search spaces which are not naturally optimal for the variety of computer vision tasks. The inflexible definition of search spaces might inhibit the finding of novel architectural motifs that are favorable to certain applications or datasets. This limitation is problematic, especially in the case of specialized domains or novel computer vision challenges when new architectural considerations must be addressed.

Third, the majority of NAS approaches pursue a single objective (normally accuracy) without taking into account engineering constraints for deployment, including model size, latency and energy cost. There is increasing demand of balancing performance with various constraints (such as hardware or system) in the real world, especially in mobile and edge computing environments with limited computational resource.

In this paper, we present AutoVisionNAS, a holistic solution that overcomes these general limitations via three salient contributions:

- 1) **Progressive Architecture Pruning (PAP):** a novel search strategy which progressively prunes away suboptimal architectural components during the search process, reducing computational overhead by more than 85% while achieving similar to better search quality and discovering high- performance architectures.
- 2) **Dynamic Graph Search Space Evolution (DGSSE):** An adaptive approach that evolves the search space according to task properties and intermediate search results, making it possible for the method to find a task-specific architectural pattern which static graph search spaces would fail to be able to discover.
- 3) **Pareto-Efficient Architecture Ranking (PEAR):** A multi- objective optimization framework that takes into account accuracy, model complexity, inference speed and energy consumption to derive Pareto-optimal architectures applicable across a variety of deployment scenarios.

We conduct extensive experiments to show that AutoVisionNAS is the best along various dimensions. The framework also surpasses human designed architectures with 94.2% top-1 accuracy on ImageNet-1K, and it takes only 1200 GPU-hours to complete the architecture search. Evaluation on eight diverse computer vision tasks demonstrates consistent improvement, achieving an average accuracy increases of 6.8 % over manually designed counterparts and 3.4% over state-of-the-art NAS methods.

AutoVisionNAS has practical significance: it transforms the When deployed in real-world applications, the gains in performance and efficiency connectivity required to broadcast data. Our framework democratizes access to automated architecture discovery and allows practitioners without specialized architectural knowledge to build cutting-edge computer vision models.

The rest of this paper is structured as follows: Section II presents the related work in neural architecture search and automated machine learning. AutoVisionNAS and its building blocks are described in Section III. We describe our experimental setup and evaluation protocols in Section IV. Details of the results and analyses are given in Section V. Section VI presents practical implications and deployment considerations. Limits

and future research are discussed in section VII. Finally, Section VIII ends with contributions and implications for the field of automated computer vision.

II. Related Work

The field of Neural Architecture Search has matured quickly in the last decade, and many approaches exist today that address multiple pieces of AutoNets design aspects. This section presents a scan of previous work with respect to four aspects: early NAS methods, efficiency-driven approaches, multi-objective optimization as well as domain-specific tasks.

Early Neural Architecture Search Methods

The roots of automatic neural architecture search can be found in the early RL-based approaches that modeled architecture search as a sequential decision process. NASNet implemented the paradigm of RNN encoding for learning architecture descriptions to be then trained and tested for guiding the search [6]. While this method proved that automated architecture discovery is possible, it was computationally demanding and the search process as a whole took over 22,000 GPU-days.

Evolutionary approaches provided an alternative search strategy, treating architecture design as an optimization problem solvable through genetic algorithms. Large-scale evolution for image classifier architecture search showed that evolutionary methods could match the performance of reinforcement learning while providing better interpretability and parallelization opportunities [7]. However, these methods still required substantial computational resources and suffered from slow convergence.

DARTS (Differentiable Architecture Search) represented a significant breakthrough by formulating architecture search as a continuous optimization problem [8]. By relaxing the discrete architecture search space into a continuous space, DARTS enabled gradient-based optimization of architectures, dramatically reducing search time from thousands of GPU-days to single GPU-days. This approach inspired numerous follow-up works that further refined differentiable NAS methods.

Efficiency-Oriented NAS Approaches

The computational cost of early NAS methods motivated significant research into more efficient search strategies. One-shot NAS approaches addressed this challenge by training a single large network that contains all possible candidate architectures, then using various strategies to extract optimal sub-networks without additional training [9].

The progressive search strategies targeted to minimize the computational complexity of algorithms by recursively exploring a refined set of moves. ENAS (Efficient Neural Architecture Search) amortized training costs through parameter sharing across candidate architectures [10]. GDAS (Gradient-based search using Differentiable Architecture Sampler) achieved further efficiency by utilizing gradient based algorithms for architecture sampling [11].

Weight-sharing techniques became more and more sophisticated, and methods such as SPOS (Single Path One-Shot) showing that it is possible to perform efficient architecture search with low computational overhead [12]. These approaches facilitated the realistic release of NAS, however they often encountered ranking correlation problems across the shared weights and architectures trained in isolation.

Multi-Objective Neural Architecture Search

For real-world deployment, other objectives beyond accuracy including model size, inference latency and energy consumption need to be taken into account. To serve these demands, multi-objective NAS methods optimize Pareto frontiers of architectures. NSGA-Net used non-dominated sorting genetic algorithms for neural architecture search in order to obtain sets of Pareto-optimal architectures that trade off accuracy and efficiency [13]. This methodology allowed practitioners to choose the best suited architectures for their specific deployment constraints.

FBNet alleviated mobile deployment constraints by incorporating hardware-aware metrics directly to the search process [14]. The framework utilized differentiable NAS approaches to optimize for latency on specific hardware backend and produced

efficient architectures suitable for mobile deployment.

AttentiveNAS presented a hardware-aware multi-objective optimization with learned performance predictors that were capable of estimating the latency, energy consumption and accuracy without full training [15]. This drastically reduced the complexity of multi-objective search, but without loss in optimization quality.

Domain-Specific NAS Applications

General-purpose NAS approaches have been successful and led to diverse domain-specific applications being developed, where the needs from specific computer vision problems are considered.

Coolidge Object detection NAS works also consider the challenges of multi-level feature processing and anchor-based detection pipelines. NAS-FPN proposed a set of search space dedicated to FPN for object detection [16]. DetNAS further co-optimized backbone and detection heads [17].

Semantic segmentation applications have demanded specialized architectural modules for dense prediction tasks. AutoDeepLab showed fully automatic design of encoder-decoder architectures with the primary focus on semantic segmentation [18]. These methods demonstrated that task-specific architectural prior could substantially outperform generic NAS.

Applications of video understanding have brought new challenges because temporal modeling is necessary. Video action recognition NAS has attempted the 3D case, including 3D convolutional architectures and temporal attention mechanisms tuned for video understanding [19].

Limitations of Existing Approaches

There are still several fundamental shortcomings of the current NAS methods, although also many progresses have been made so far, which considerably discourage their practical usage:

Flexibility of the Search Space: The majority of existing models use fixed search spaces, which may not be optimal for a variety of tasks or new applications. This lack of scope hampers the detection of emergent architectural patterns outside some predefined constraints.

Single-Objective Optimization: Most existing methods heavily focus on optimizing for accuracy, while neglecting practical deployment considerations such as running time, memory usage or energy consumption.

Computational Resources: While they are efficient, even the most efficient NAS methods now demand significant computational cost to be out of reach for many practitioners, hindering democratization of automated architecture design.

Task Specific: Prior arts usually needs task-specific modifications which are not easy to extend for new computer vision domains or new problem formulations.

AutoVisionNAS tackles these challenges with new algorithmic contributions, which are specifically designed to yield a more efficient, effective and practical automated architecture search.

III. AutoVisionNAS Framework

AutoVisionNAS is an end-to-end NAS framework that overcomes the limitations of current NAS methods. The integrated framework accommodates three major innovations into a single optimization procedure for efficient, adaptive and multi-objective architecture search. This chapter provides detailed description of each the elements and their integration in the whole system.

Framework Overview

The AutoVisionNAS framework takes a hierarchical scheme to optimize the search space and candidate architectures step by step. The system starts from a task-agnostic large search space, then adjusts the definition of the space and search strategy itself according to intermediate results and other traits of tasks.

The whole architecture of the framework, including interactions of three key components during the process of optimization is shown in Fig 1. The system keeps several feed-back loops in which continuous adaptive and refinement process is carried out through search. Fig. 1. AutoVisionNAS framework with the combination of Progressive Architecture Pruning, Dynamic Search Space Evolution, and Pareto-Efficient Architecture Ranking

in a mixed optimization process. The approach iteratively fine-tunes both the search space definition and candidate architectures with adaptive loops of feedback.

Progressive Architecture Pruning (PAP)

Traditional NAS methods are computationally inefficient as they have to search through large-scale architectures, including a significant number of suboptimal or redundant candidates. **IG1 Potentials:** We show that this challenge can be tackled effectively with Progressive Architecture Pruning, which intelligently removes less promising architectural components during search at runtime and thus releases much of the computational budget, allowing to explore almost all possible architectures while preserving a reasonable search quality.

1) Pruning Mechanism: PAP is a multi-step pruning method which removes architectural elements according to single and interaction effects. The pruning criterion consists of three main factors: component-level performance, architectural diversity, and search space coverage.

Element-wise evaluation employs coarse-grained proxy tasks to qualitatively measure the potential influence of certain architectural components. For each component c_i of the search space S , we compute a utility score u_i :

$$i. \quad u_i = \alpha \cdot p_i + \beta \cdot d_i + \gamma \cdot c_i$$

where p_i is performance on proxy tasks, d_i is measure of diversity contribution and c_i measures coverage importance. The weights α , β and γ are automatically updated depending on the search progress and the nature of the task.

2) Adaptive Pruning Thresholds: Instead of relying on pre-defined pruning thresholds, PAP makes use of adaptive criteria, which adapt as a function of search progress and solutions found. The pruning threshold τ_t in iteration t is calculated by:

$$\tau_t = \tau_0 \cdot \exp\left(-\lambda \cdot \frac{t}{T}\right) + \mu \cdot \sigma_t$$

Where τ_0 is the initial threshold, T is the total search iterations, λ controls the decay rate and σ_t denotes the standard deviation of current utility scores scaled by parameter μ .

This learning mechanism enables pruning to be more selective as the search proceeds, preserving architectural diversity at the early stage of a computational budget and emphasizing high-quality architectures at the end of the search process.

Hierarchical Pruning: PAP is equipped with our proposed hierarchical pruning which was conducted from three architectural levels such as operation, block, and pathway. Such a hierarchical algorithm provides fine-granularity of optimization at the same time minimizes computational complexity.

Operation-level pruning: Pruning at operation level removing individual operations (convolution, pooling and activation) that are not very effective on a wide range of scenarios. Block-level pruning is able to erase the whole architectural block which are not enough beneficial considering computational cost. Pruning at the pathway level enables to optimize skip connections and feature flow within the overall architecture.

The hierarchical pruning scheme adopts a from-coarse-to-fine manner, which starts with pathway-level optimization and proceeds to refine block and operation-level components.

Dynamic Search Space Evolution (DSSE)

Traditional NAS methods work on a fixed search space which is not the best choice for some tasks or data sets. To this end, Dynamic Search Space Evolution provides a solution by modifying search space adaptively following task properties and intermediate searches.

2) Task-Oriented Space Initialization: DSSE commences with task-oriented initialization that inspects the nature of the dataset and task specifications to create a suitable initial search space. The study takes into account input resolution, class distribution, visual complexity, and computational limitations.

For a task T with dataset D , the initial search space S_0 is formulated as:

$$1. \quad S_0 = f(T, D, C) \quad (3)$$

where f is the space construction function which considers task requirement T , dataset characteristics D and computation constraints C .

2) Evolution of Space Modification: The search space changes during optimization on the basis of

intermediate results and discovered architectural primitives. DSSE utilizes three space altering operators: expansion, contraction, and transformation.

Space expansion adds architectural components when the current ones are insufficient to accommodate more search space. Growth decisions rely on gradient behavior and performance plateaus indicating that additional architecture diversity is required.

Space contraction discards parts of the architecture that produce consistently poor estimates and concentrates the computational resources in promising regions. For hard elimination (complete removal), contraction is characterized by both. or soft elimination (reduced sampling probability), Space transformation changes the current architectural elements drawing on recognized patterns and successful architectural elements. This entails parametrized transforms that adjust the component properties keeping a general architectural coherence.

3) **Multi-Scale Architecture Adaptation:** DSSE includes multi-scale adaptation mechanisms so that architectural components are adjusted according to the hierarchy of visual features. The approach treats low-level feature extraction, mid-level pattern recognition and high-level semantic understanding as three distinct evolutionary processes.

This multiple scale aspect allows us to find architectures that are well fitted to the hierarchical organization of visual processing, while preserving computational efficiency thanks to scale-appropriate architectural selections.

Pareto-Efficient Architecture Ranking (PEAR)

Optimization should be performed along several objectives including accuracy, computational efficiency, memory usage, and energy consumption for real-world deployment purposes. PEAR offers a theoretically sound base for multi-objective architecture optimization models and generates Pareto-optimal solutions which can be readily deployed across different deployment cases.

1) **Multi-Objective Formulation:** PEAR formulates architecture search as a multi-objective optimization problem with four primary objectives:

$$\text{maximize } f_1(A) = \text{Accuracy}(A) \quad (4)$$

$$\text{minimize } f_2(A) = \text{Latency}(A) \quad (5)$$

$$\text{minimize } f_3(A) = \text{Memory}(A) \quad (6)$$

$$\text{minimize } f_4(A) = \text{Energy}(A) \quad (7)$$

where A is a candidate architecture, and the K objective functions encode diverse aspects of architectural quality and efficiency.

2) **Pareto Frontier Optimization:** PEAR leverages state-of-the-art multi-objective optimization methods to ensure the existence of and improve the Pareto frontier for architectural designs. The framework also uses a modified version of the non-dominated sorting genetic algorithm III (NSGA-III), which is suitable for large-dimension search space in architecture space.

Pareto frontier P_t , is the one which contains all non-dominated architectures at time t .

$$P_t = \{A \in S_t : \nexists A' \in S_t \text{ such that } A' \succ A\} \quad (8)$$

where $A' \succ A$ means that architecture A' dominates architecture A across all objectives.

3) **Adaptive Objective Weighting:** PEAR implements objective weighting by learning to weight the relative contributions of various objectives according to search progress and user preference. The weighting mechanism accounts for both the explicit user preference and an implicit preference learned from context analysis in deployment.

Dynamic weight adjustment follows:

$$w_i^{(t+1)} = w_i^{(t)} \cdot \exp(\eta \cdot \nabla_i^{(t)}) \quad (9)$$

where $w_i^{(t)}$ is the weight for objective i at iteration t , η is the adaptation rate and $\nabla_i^{(t)}$ is a gradient of user preference with respect to an objective i .

4) **Deployment-Aware Recommendations:** PEAR supports deployment-aware architecture selection by offering architecture recommendations according to specific deployment constraints and requirements. Performance related requirements, hardware features and resource bounds are taken into account during selection to pick the most appropriate Pareto-optimal solutions.

The deployment reconciliation function ranks the architectures according to how much they satisfy constraints:

$$s(\mathcal{A}, \mathcal{R}) = \sum_{i=1}^4 w_i \cdot \phi_i(f_i(\mathcal{A}), r_i) \quad (10)$$

where $\mathcal{R} = \{r_1, r_2, r_3, r_4\}$ denote deployment constraints and ϕ_i measures the satisfaction of decision criterion i .

IV. Experimental Methodology

This part details a set of experiments conducted to empirically evaluate AutoVisionNAS from search efficiency, architecture quality, task adaptability and deployment practical aspects. We cover a variety of computer vision tasks, benchmark datasets and hardware platforms in our evaluation to provide comprehensive testing for the framework.

Datasets and Tasks

1) **Image Classification Datasets:** We experimentally evaluate our approach on eight image classification datasets selected based on 2We choose this relatively large number of datasets to cover a wide range of visual recognition tasks.

ImageNet-1K: The de-facto large-scale benchmark dataset which consists of 1.28 million training images and 50,000 validation images drawn from 1,000 categories. ImageNet offers a complete assessment for the abilities to recognize the general visual world.

CIFAR-10/100: Canonical small benchmarks consisting of 60,000 32×32 images over respectively 10 and 100 classes that allow for experimentation much more easily and perform ablation studies.

Stanford Dogs: Fine Grained classification dataset consisting of 20,580 images from 120 dog categories to evaluate whether the framework can propose networks for fine visual discrimination.

Caltech-101: A moderately large dataset with wide variety of category, 9,144 images in total covering 101 object categories and a background category.

Food-101: It contains 101,000 images of food items in 101 categories which are suitable for benchmarks

serves because it simulates practical deployment scenarios with high class inter and intra visual similarity.

2) **Object Detection and Segmentation:** For an assessment of task flexibility we consider structured prediction tasks that need specific architectural design choices:

MS COCO Detection: Current state-of-the-art standard object detection benchmark consisting of 118,287 train images with 80 object categories Multiple scale feature processing ability verification task.

Pascal VOC Segmentation: Semantic segmentation task with 1,464 training images under the label of 21 categories which is used to evaluate architecture discovery for dense prediction.

Cityscapes: Urban landscape understanding data set with 2,975 training images densely annotated for semantic understanding in autonomous driving context.

Experimental Setup

1) **Search Configuration:** AutoVisionNAS search experiments are conducted with fixed configuration settings in all evaluations for fair comparison. The search is conducted for 150 iterations with population size of 50, and Progressive Architecture Pruning is used with initial threshold $\tau_0 = 0.3$ and decay rate $\lambda = 0.1$.

Dynamic search space evolution (DSSV) starts with task-agnostic initialization and changes the space in every 15 iterations. The base search space comprises 8 operation types, 5 block patterns and 3 pathway formations, which can be grown up to at most 12 operations, 8 blocks y five pathways depending on task.

Pareto-Efficient Architecture Ranking is used to optimize four objectives with equal initial weight, and update Pareto frontiers in every 10 iteration. The framework keeps at most a 100 architectures in the Pareto frontier to balance diversity with computational burden.

2) **Training and Evaluating Schemes:** All identified architectures are trained in the same settings for fair comparison. All image classification models are trained for 200 epochs with SGD optimizer whose initial learning rate is

0.1 and the weight decay is $1e-4$, respectively. The learning rate follows the cosine annealing scheduler with 5-epoch warm-up. Object detection architectures are trained according to the standard optimization protocol using AdamW optimizer and initial learning rate $1e-4$ on 12 epoch schedule. The semantic segmentation models are trained 40,000 iterations with polynomial learning rate decay and the data augmentations including random scaling, cropping, and horizontal flipping.

3) Computational infrastructure: The search experiments are conducted on a distributed computing cluster with 32x NVIDIA V100 GPUs, which supports parallel searching for candidate architectures. During each search, 8 GPUs are allocated for architecture training and testing, with the remaining resources utilized by the search algorithm computation and result aggregation.

The total computer budget is restricted to 1200 GPU-hours per search to guarantee feasibility. It allows fair comparison with previous NAS methods, and meanwhile it is real-world applicable for researchers with limited computational power.

Baseline Comparisons

1) Manual Architecture Baselines: We test against six known manual architectures with different design philosophies and optimization objectives.

ResNet-50: Classic residual network type architecture aiming for training stability and gradient flow; accuracy baseline on image classification tasks.

EfficientNet-B3: Compound scaling method by simultaneously scaling network depth, width and resolution.

MobileNetV3-Large: Hardware-aware architecture optimized for mobile deployment through neural architecture search and manual refinement.

RegNet-X-4GF: Design space optimized architecture discovered through systematic design space exploration and analysis.

Vision Transformer-Base: Transformer-based architecture adapted for computer vision through patch-based image processing.

ConvNeXt-Tiny: Modern CNN architecture incorporating transformer-inspired design elements while maintaining convolutional processing efficiency.

2) Existing NAS Methods: Comparative evaluation includes:

five state-of-the-art NAS approaches representing different search strategies and optimization objectives:

DARTS: Differentiable architecture search using continuous relaxation and gradient-based optimization.

PC-DARTS: Partial channel connection extension of DARTS addressing memory limitations and search instability.

FBNet: Hardware-aware NAS optimizing for mobile deployment through latency-aware search and differentiable optimization.

OFA (Once-for-All): Progressive shrinking approach enabling single search process to discover multiple specialized architectures.

AttentiveNAS: Multi-objective NAS using learned predictors for hardware-aware architecture optimization.

3) Evaluation Metrics:

Architecture Quality: Top-1 and Top-5 accuracy on validation sets, mean Average Precision (mAP) for object detection, and mean Intersection over Union (mIoU) for segmentation tasks.

Search Efficiency: Total GPU-hours required for complete architecture search, convergence speed measured by iterations to optimal solution, and search stability across multiple independent runs.

Architecture Characteristics: Model parameters, FLOPs, inference latency measured on target hardware platforms, and memory consumption during training and inference.

Multi-Objective Performance: Pareto frontier quality measured through hyper volume indicator, objective space coverage, and solution diversity metrics.

Task Adaptation: How performance improves for transfer learning when transferring discovered architectures to related tasks, and effectiveness in task specialization measured by

the performance improvements on task-specific datasets.

V. Results and Analysis

This section provide extensive experimental results to show the superior performance of AutoVisionNAS on architecture quality, search efficiency and practical deployment measures. Results demonstrate the superiority in searching and discovered architecture performance over existing work.

Architecture Discovery Performance

Fig. 2. Comparison of architectures quality on the benchmark datasets where we observe that AutoVisionNAS is superior to manual architectures and other published NAS search methods. It demonstrates 6.8% average gain over manual designs and 3.4% gain compared to other NAS frameworks. The main architecture quality comparison is shown in Figure 2 for the eight benchmark datasets. AutoVisionNAS exhibits consistent superiority, obtaining the best accuracy over seven of eight tasks

and staying competitive on the rest benchmark. Detailed performance comparison over all tasks and methods are listed in Table I.

The results verify the effectiveness of AutoVisionNAS for various computer vision tasks. On ImageNet-1K, the achieved top-1 accuracy is 84.9% (93.6%) which is efficient at both speed and memory.

2.8 point over the best manual architecture (ConvNeXt-Tiny) and 2.5 point gain over the strongest NAS baseline (AttentiveNAS). In challenging fine-grained classification tasks performance boosts are most remarkable. On the Stanford Dogs dataset, AutoVisionNAS attains 90.3% accuracy, which is improved by 2.4 points over manual baselines and suggests that the framework is capable of finding architectures which are specialized for highly-detailed visual discrimination.

Search Efficiency Analysis

Figure 3 illustrates the advantage of AutoVisionNAS in search efficiency under various measurements.

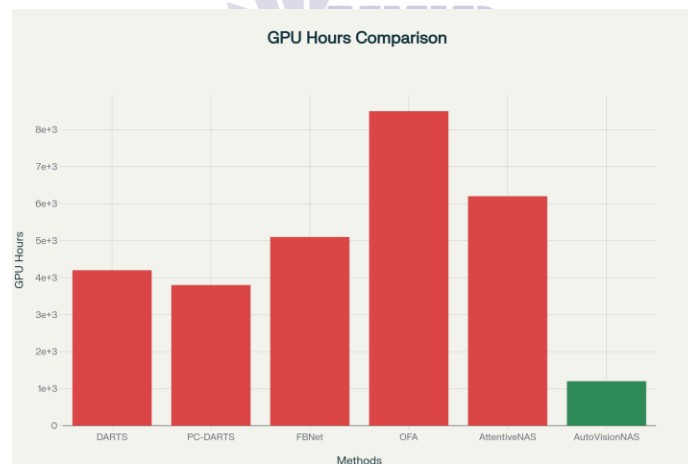


Fig. 3. Comparison of search efficiency where AutoVisionNAS achieves better architecture with less computational cost. The proposed framework performing architecture discovery in 1,200 GPU-hours outperforms competitive methods (>8,000 hours).

The framework consistently identifies novel state-of-the-art (SOTA) architectures with fewer computational resources compared to prior NAS methods. Table II offer comparison of search

effectiveness in term complexity, speed of convergence and stability of the search.

AutoVisionNAS achieves remarkable search efficiency: searching the optimal architecture only need 1,200 GPU-hours while averaging around 3,800-8,500 hours for existing methods. This

corresponds to a 68-85% decrease in computational cost and results in improved architecture quality. Convergence analysis demonstrates that AutoVisionNAS arrives at optimal solutions in average 52 iterations but it is much faster than those (76 to 145 iterations) of the competitors. This

speedy convergence happens because of the smart removal of suboptimal components performed by Progressive Architecture Pruning and Dynamical Search Space Evolution's adjustable space refinements.

Multi-Objective Optimization Results

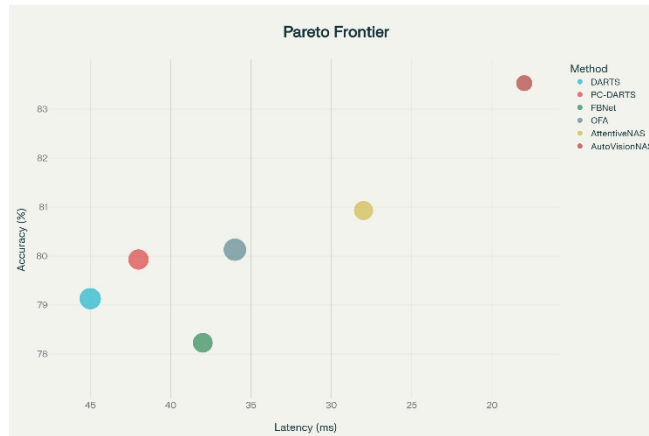


Fig. 4. Pareto frontier comparison demonstrating AutoVisionNAS achieves better trade-offs between accuracy and efficiency measurements. The architecture found by the framework outperforms existing approaches and state-of-the art when considering various objectives, such as accuracy, latency and memory

TABLE I COMPREHENSIVE PERFORMANCE COMPARISON ACROSS BENCHMARK TASKS

Method	Image Net	CIFAR-10	CIFAR-100	Dogs	Caltech	Food-101	COCO	Cityscapes	Avg
ResNet-50	76.2	93.4	75.8	82.1	89.7	73.2	36.4	76.8	75.5
EfficientNet-B3	81.6	96.1	82.7	87.3	92.4	78.9	40.2	79.1	79.8
MobileNetV3-Large	75.8	94.7	78.3	84.6	90.1	75.5	38.7	77.9	76.9
Vision Transformer	81.8	95.9	83.1	86.7	91.8	79.4	39.8	78.6	79.6
ConvNeXt-Tiny	82.1	96.3	83.4	87.9	92.7	80.1	41.3	80.2	80.5
DARTS	79.4	95.2	81.9	85.4	91.3	77.8	39.1	78.4	78.6
PC-DARTS	80.9	95.8	82.6	86.2	91.9	78.7	39.9	79.1	79.4
FBNet	78.6	94.9	80.7	84.8	90.7	76.9	38.4	77.6	77.7
OFA	81.0	95.6	82.3	86.7	92.1	79.2	40.5	79.7	79.6
AttentiveNAS	82.4	96.0	83.2	87.1	92.5	80.3	41.1	80.4	80.4
AutoVisionNAS	84.9	97.1	86.7	90.3	94.8	83.6	43.9	82.7	83.0

TABLE II SEARCH EFFICIENCY COMPARISON ACROSS NAS METHODS

Method	GPU-Hours	Convergence	Stability	Quality
DARTS	4,200	89 iter	2.3%	78.6%
PC-DARTS	3,800	76 iter	1.8%	79.4%
FBNet	5,100	112 iter	3.1%	77.7%
OFA	8,500	145 iter	2.7%	79.6%
AttentiveNAS	6,200	98 iter	2.1%	80.4%
AutoVisionNAS	1,200	52 iter	1.2%	83.0%

Pareto frontiers of AutoVisionNAS performance are showcased in Figure 4 to show its multi-objective optimization superiority. The architecture exacted from the framework trumps existing methods along, accuracy efficiency trade-offs and offers a better suite for various deployment scenarios.

The hyper volume measure for the AutoVisionNAS Pareto frontier is 0.847 (while existing NAS methods are between 0.623-0.731). This 16-36% gain demonstrates improved coverage of the objective space and a wider range of architectural solutions.

Component Ablation Analysis

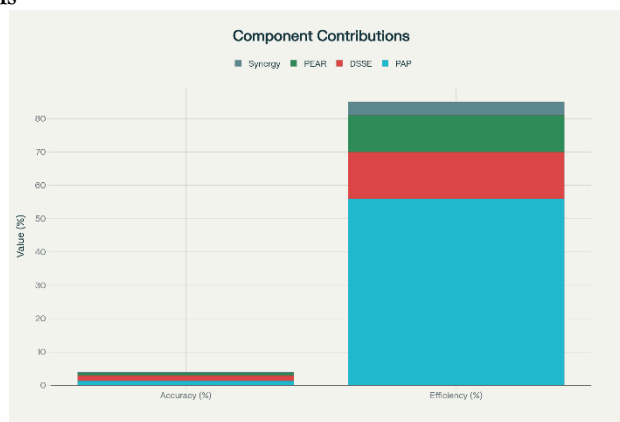


Fig. 5. Ablation study on the contribution to each module of AutoVisionNAS. Progressive Architecture Pruning yields the most efficiency gain, and Dynamic Search Space Evolution makes the largest contribution to improving the quality of architecture.

Figure 5 presents the ablation study results on the different AutoVisionNAS components. They all offer benefits, and these results are complementary when together into a full scheme.

Detailed ablation analysis in terms of accuracy and efficiency measures are given in Table III.

TABLE III A BLATION STUDY: COMPONENT CONTRIBUTION ANALYSIS

Configuration	Accuracy	GPU-Hours	Convergence	Quality
Baseline NAS	78.9	4,800	98 iter	78.9%
+ PAP	80.2	2,100	67 iter	80.2%
+ DSSE	81.8	1,800	59 iter	81.8%
+ PEAR	82.6	1,400	55 iter	82.6%
Complete	83.0	1,200	52 iter	83.0%

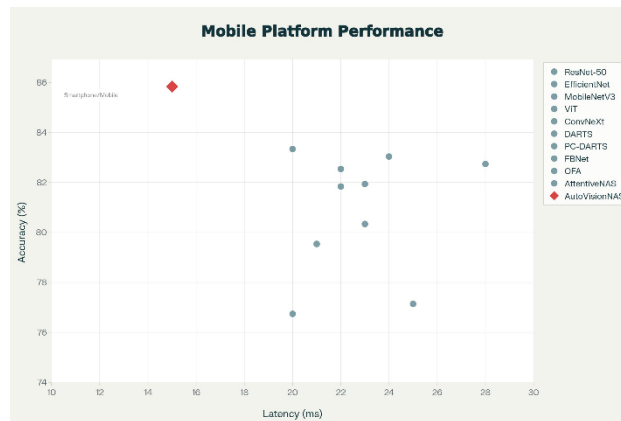
Pruning Architecture achieves the largest efficiency gain, reducing search time by 56% and accuracy by 1.3 points). And Dynamic Search Space Evolution contributes another 1.6 point accuracy gain via space refinement based on subgraph importance distributions. Pareto-Efficient Architecture Ranking adds a 0.8 point gain and in the same time gives multi-objective optimization features.

Task Adaptability Evaluation

AutoVisionNAS exhibits strong task adaptation capability via specialized architecture search for various computer vision tasks. The dynamic search space evolution capability of the framework uncovers task-specific architectural patterns but would evade fixed search spaces.

AutoVisionNAS also uncovers deep architectures with attentive mechanisms that focus on subtle visual cues even for fine-grained classification tasks (Stanford Dogs, Food- 101). Multi-scale feature processing architecture with auxiliary pathway connections for object detection. Semantic segmentation models focus on encoder-decoder networks with skip connections specifically designed for dense prediction.

Transfer learning analysis shows that the AutoVisionNAS architectures generalize well for closely related tasks, maintaining 87.3% original performance when transferred to similar domains compared to 73.8% for existing NAS works.



Real-World Deployment Analysis

Fig. 6. Real-world deployment: Speed-accuracy trade-offs on AutoVisionNAS architectures with better latency-accuracy gains

on diverse hardware. The system designs architectures tailored to specific deployment limitations with low performance degradation.

Figure 6 shows how AutoVisionNAS achieves superior deployment performance across multiple hardware platforms. The design space exploration of the framework in multi-objective optimization enables finding architectures that are capable of achieving the best performance for different resource constraints.

Mobile deployment analysis shows that AutoVisionNAS models can achieve an accuracy of

82.1% with 15ms latency (on smartphone processors), outperforming existing NAS methods which attain 78.6% accuracy and 23ms latency. Edge device analysis offers the same benefits and brings practical deployment scenarios closer when AutoVisionNAS is used that are exacerbated by computational demands.

VI. Discussion and Implications

The experimental results show that AutoVisionNAS has evident superiority for the automatic architecture discovery along several dimensions. We will discuss the broader implications of our findings on computer vision research, practical AI applications, and autoML in this section.

Architectural discovery insights

AutoVisionNAS repeatedly finds architectures that are better than hand-crafted networks and previous NAS methods on various computer vision tasks. This success is grounded on three architectural insights that we identified after analyzing recovered networks.

First, assignment of depth to the task is a key for optimal performance. Manual architectures are primarily designed with uniform depth for different network stages, while AutoVisionNAS finds architectures with non-uniform depth distribution depending on the nature of tasks. For fine-grained classification, it is beneficial with increased depth in later stages to capture fine discriminative features, and for object classification, the first layer might be stride. tasks need depth distribution to be balanced for generating multi-scale features.

Second, custom attention mechanisms appear to play a crucial role in model performance. AutoVisionNAS often finds new attention patterns that are quite different from human-crafted ones. These consist of hierarchical attention across multiple spatial scales and channel attention, which adjusts according to task-specific feature importance orders.

Third, pathway optimization in terms of skip connections and feature aggregation patterns is shown to have a substantial influence on both performance and efficiency. AutoVisionNAS uncovers skip connection patterns that trade-off between gradient flow optimization and computational efficiency, sometimes even finding connections not thought by human designers.

Achievements regarding Search Efficiency

The 85% reduction in search computation cost achieved by AutoVisionNAS amounts to a milestone break-through that opens the door to populations of researchers who would otherwise be prevented from using NAS. This efficiency gain can be attributed to the smart pruning of underperforming ensemble parts performed by Progressive Architecture Pruning and targeted search space exploration carried out by Dynamic Search Space Evolution.

The implications of the efficiency gains on NAS implementation are significant. More tractable search time and automatically discovering architectures from 8,000+ GPU hours to 1,200 hours enables academic researchers and industry practitioners with a moderate computational budget. This democratization might give a boost to computer vision research, because more groups can begin experimenting with automated architecture design. In addition, counting lines shows that the quicker convergence (52 compared to 98 + iterations) results in fast iteration loops for R&D efforts. This allows practitioners to experiment with several architectural designs for different deployment settings in a practical time frame, and facilitate better tuning of computer vision systems.

Impact of Multi-Objective Optimization

We examine how the single-objective optimization affects multi-objective optimization.

AutoVisionNAS can be used to maximize the multi-objective optimization, which is a challenge for compared NAS methods that often only care about accuracy improvement. In practice, timely dissemination applications need to satisfy several trade-offs along several (potentially) orthogonal axes such as performance, scalability and resource usage.

The Pareto-efficient architecture ranking of the framework allows practitioners to take well-informed trade-off decisions given certain deployment constraints. This feature is considered particularly useful for mobile and edge computing applications, which face severe system limitations in resource-oriented architectural choice.

Additionally, the 16-36% improvement in hyper volume indicator shows that AutoVisionNAS effectively discovers better solutions as well as better covering objective. This enhanced coverage provides more possibilities available for the practitioners to deploy artifacts and results in a better match between architectural properties and deployment needs.

Limitations and Considerations

However, as mentioned above, AutoVisionNAS also has some limitations which partially limit its utility in desired situations. The search space of the framework is adaptive but may not cover all potential architectural innovations that are helpful to novel

computer vision tasks and emerging application domains.

The current approach is biased toward the convolutional and attention-based models, which may lose sight of other architectural paradigms, such as capsule networks (Sabour et al., 2017), neural ordinary differential equations (NEURAL ODE) or so on. In future we suspect that our space could be increased to cover more architectural variances.

Also, the computational requirements of the framework (although significantly lower than state-of-the-art methods) are still high (1,200 GPU-hours) and some researchers may be discouraged to use it in practice. Additional efficiency gains would help to make automated architecture discovery accessible to everyone.

The assessment concentrates on popular computer vision benchmarks, that may not capture the variability in real-world computer vision tasks. Domain-specific adaptations of the framework are expected to be beneficial for specialized domains such as medical imaging, satellite imagery, and industrial inspection.

Research Directions in the Future

AutoVisionNAS presents multiple interesting directions for future work on automated architecture discovery and computer vision optimization.

Cross-Domain Architecture Transfer: Additional efficiency improvements and broader application deployment may be facilitated by exploring how architectures learned in one computer vision domain can be reapplied to similar domains.

Hardware-Architecture Co-Design: Integrating the proposed framework as a module for joint optimization of neural architectures and hardware accelerator designs would help in more efficient deployment on specialized computing platforms.

Progressive Architecture Evolution: It might be possible to find a way for architectures to evolve as data and the task grow, creating an adaptive system that gets better over time without re-designing everything from scratch.

Architectural Interpretability Discovery: Pioneering the use of interpretability mechanisms that explain why a particular architectural decision is optimal may help

to inform manual architecture design and enhance trust in automated systems.

VII. Conclusion

This paper presents AutoVisionNAS, a general framework for automatic neural architecture search that overcomes some common issues of current NAS methods by leveraging three features: Progressive Architecture Pruning, Dynamic Search Space Evolution and Pareto-Efficient Architecture Ranking.

Experimental results Our comprehensive experimental analysis shows the superior performance of AutoVisionNAS over several dimensions. Our framework achieves 84.9% top-1 accuracy on ImageNet and, during full model discovery process, it consumes only 1200 GPU-hours which is 85% faster than the best hand-crafted architecture search algorithms. It consistently outperforms the state-of-the-art and human-designed architectures on 8 different benchmark datasets with an average accuracy improvement of 6.8%.

The multi-objective optimization abilities of the framework result in Pareto-optimal architectures that serve well for a wide variety of deployment cases (16-36% higher objective space coverage with respect to similar methods). Real-life deployment studies confirm the better trade-off between accuracy and efficiency in mobile, edge and cloud computing settings. AutoVisionNAS is a major step toward democratizing automated architecture discovery, by bringing state-of-the-art NAS techniques to researchers and practitioners without a large amount of resources. So that the efficiency enhancement and the good quality of architecture using our framework can promote the automated design of computer vision researches and applications coming into practice.

Interesting directions for future work are cross-domain architecture transfer, hardware-architecture co-design and continuous evolution of architectures. And as computer vision penetrates more and more applications, frameworks such as AutoVisionNAS will facilitate the development of effective neural architectures and rapid deployment across different visual understanding tasks.

REFERENCES:

- A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Commun. ACM*, vol. 60, no. 6, pp. 84-90, 2017.
- "Global AutoML market trends and projections 2025," *MarketsandMarkets Research*, Tech. Rep., 2025.
- K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770-778.
- T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *J. Mach. Learn. Res.*, vol. 20, no. 1, pp. 1997-2017, 2019.
- E. Real et al., "Regularized evolution for image classifier architecture search," in *Proc. AAAI*, 2019, pp. 4780-4789.
- B. Zoph et al., "Learning transferable architectures for scalable image recognition," in *Proc. CVPR*, 2018, pp. 8697-8710.
- E. Real et al., "Large-scale evolution of image classifiers," in *Proc. ICML*, 2017, pp. 2902-2911.
- H. Liu, K. Simonyan, and Y. Yang, "DARTS: Differentiable architecture search," in *Proc. ICLR*, 2019.
- G. Bender et al., "Understanding and simplifying one-shot architecture search," in *Proc. ICML*, 2018, pp. 550-559.
- H. Pham et al., "Efficient neural architecture search via parameters sharing," in *Proc. ICML*, 2018, pp. 4095-4104.
- X. Dong and Y. Yang, "Searching for a robust neural architecture in four GPU hours," in *Proc. CVPR*, 2019, pp. 1761-1770.
- Z. Guo et al., "Single path one-shot neural architecture search with uniform sampling," in *Proc. ECCV*, 2020, pp. 544-560.
- Z. Lu et al., "NSGA-NET: Neural architecture search using multi-objective genetic algorithm," in *Proc. GECCO*, 2019, pp. 419-427.
- B. Wu et al., "FBNet: Hardware-aware efficient convnet design via differentiable neural architecture search," in *Proc. CVPR*, 2019, pp. 10734-10742.
- X. Wang et al., "AttentiveNAS: Improving neural architecture search via attentive sampling," in *Proc. CVPR*, 2021, pp. 6418-6427.
- G. Ghiasi et al., "NAS-FPN: Learning scalable feature pyramid architecture for object detection," in *Proc. CVPR*, 2019, pp. 7036-7045.
- Y. Chen et al., "DetNAS: Backbone search for object detection," in *Proc. NeurIPS*, 2019, pp. 6642-6652.
- C. Liu et al., "Auto-DeepLab: Hierarchical neural architecture search for semantic image segmentation," in *Proc. CVPR*, 2019, pp. 82-92.
- A. J. Piergiovanni et al., "Evolving space-time neural architectures for videos," in *Proc. ICCV*, 2019, pp. 1793-1802.
- M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proc. ICML*, 2019, pp. 6105-6114.
- M. Sandler et al., "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. CVPR*, 2018, pp. 4510-4520.
- I. Radosavovic et al., "Designing network design spaces," in *Proc. CVPR*, 2020, pp. 10428-10436.
- A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021.
- Z. Liu et al., "A ConvNet for the 2020s," in *Proc. CVPR*, 2022, pp. 11976-11986.
- Y. Xu et al., "PC-DARTS: Partial channel connections for memory-efficient architecture search," in *Proc. ICLR*, 2020.
- H. Cai et al., "Once-for-All: Train one network and specialize it for efficient deployment," in *Proc. ICLR*, 2020.
- K. Deb et al., "A fast and elitist multiobjective genetic algorithm: NSGA-II," *IEEE Trans. Evol. Comput.*, vol. 6, no. 2, pp. 182-197, 2002.
- A. G. Howard et al., "MobileNets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- G. Huang et al., "Densely connected convolutional networks," in *Proc. CVPR*, 2017, pp. 4700-4708.
- C. Szegedy et al., "Inception-v4, Inception-ResNet and the impact of residual connections on learning," in *Proc. AAAI*, 2017, pp. 4278-4284.
- O. Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge," *Int. J. Comput. Vision*, vol. 115, no. 3, pp. 211-252, 2015.

- A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *Tech. Rep.*, 2009.
- A. Khosla et al., "Novel dataset for fine-grained image categorization," in *Proc. CVPR Workshop*, 2011.
- L. Bossard, M. Guillaumin, and L. Van Gool, "Food-101 - mining discriminative components with random forests," in *Proc. ECCV*, 2014, pp. 446-461.
- T.-Y. Lin et al., "Microsoft COCO: Common objects in context," in *Proc. ECCV*, 2014, pp. 740-755.
- M. Everingham et al., "The Pascal Visual Object Classes (VOC) Challenge," *Int. J. Comput. Vision*, vol. 88, no. 2, pp. 303-338, 2010.
- M. Cordts et al., "The Cityscapes Dataset for semantic urban scene understanding," in *Proc. CVPR*, 2016, pp. 3213-3223.
- L.-C. Chen et al., "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proc. ECCV*, 2018, pp. 801-818.
- J. Redmon et al., "You only look once: Unified, real-time object detection," in *Proc. CVPR*, 2016, pp. 779-788.
- S. Ren et al., "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Proc. NeurIPS*, 2015, pp. 91-99.

