

# Experimental Analysis of Lightweight CNNs for Real-Time Object Detection on Low-Power Devices

Karan Singh\*, Kumari Kajal<sup>†</sup>, Sarita Negi<sup>‡</sup>

*Department of Information Technology*

*Noida Institute of Engineering and Technology, Greater Noida, India*

*Email: \*karan.singh@niet.co.in, <sup>†</sup>hyekajal@gmail.com, <sup>‡</sup>sarita.sharma2612@yahoo.com*

**Abstract**—The integration of artificial intelligence (AI) into portable and embedded systems has led to a paradigm shift in how deep learning models are developed and deployed. In particular, the increasing demand for real-time computer vision tasks on resource-limited hardware has emphasized the need for computationally efficient architectures. Object detection, a cornerstone of computer vision, typically involves high computational complexity, substantial memory requirements, and elevated power consumption—rendering traditional convolutional neural networks (CNNs) impractical for use in low-power environments such as smartphones, Raspberry Pi, and edge IoT platforms.

This paper presents an in-depth experimental analysis of lightweight CNN architectures tailored for real-time object detection on constrained devices. We critically examine state-of-the-art models including MobileNet, YOLOv4-Tiny, SqueezeNet, and EfficientDet-Lite, assessing them across diverse performance parameters. Key metrics such as model footprint, mean average precision (mAP), inference latency, frames per second (FPS), and power consumption are used to benchmark these models on representative edge hardware. In addition, we explore the efficacy of optimization techniques such as quantization, pruning, and neural architecture search (NAS) in enhancing model efficiency without significantly compromising detection accuracy.

Through rigorous evaluation and comparative analysis, the study highlights the trade-offs between accuracy and computational efficiency, providing practical guidance for model selection in real-world scenarios. Case studies and empirical results offer insight into performance bottlenecks and optimization potential, while a forward-looking discussion addresses ongoing challenges and future research directions. This work aims to bridge the gap between high-performance object detection and resource-aware deployment, paving the way for scalable, energy-efficient AI applications on the edge.

**Keywords**—Lightweight CNNs, Edge AI, Object Detection, Low-Power Devices, Real-Time Inference, Model Compression

## I. INTRODUCTION

### A. Background

The rapid advancement of deep learning has significantly transformed the landscape of computer vision, enabling machines to achieve human-level performance in tasks such as image classification, semantic segmentation, and object detection [1]–[4], [6]. Among these, Convolutional Neural Networks (CNNs) have emerged as the fundamental architecture due to their hierarchical feature extraction capabilities and strong generalization on visual data [52]. Object detection, in particular, plays a pivotal role in real-time systems including autonomous vehicles [7], smart surveillance [8]–[10], augmented reality [11], and wearable health monitoring devices [12].

Despite their effectiveness, state-of-the-art object detectors such as Faster R-CNN [59], YOLOv5 [16], and SSD [61] are characterized by high model complexity, requiring substantial computational power, memory bandwidth, and energy. These requirements hinder their direct deployment on edge and embedded systems like smartphones, drones, and IoT nodes, which operate under strict resource constraints. As a result, the pursuit of lightweight alternatives that can balance performance with efficiency has become increasingly essential [13], [14], [17], [62].

### B. Motivation

The emergence of edge AI, where computations are performed locally on-device, has introduced a new frontier in embedded intelligence. Unlike cloud-based approaches, edge AI reduces latency, improves privacy, and enhances responsiveness—qualities crucial for applications such as autonomous drones, real-time surveillance, and health diagnostics [18], [21], [36], [85], [91]. However, implementing deep learning models on edge devices remains a complex challenge due to the trade-off between inference accuracy and resource efficiency.

To address this, researchers have developed lightweight CNNs such as MobileNet [62], [63], YOLOv4-Tiny [72], SqueezeNet [80], and EfficientDet-Lite [79]. These models aim to reduce computational cost while preserving acceptable levels of accuracy. Techniques like depthwise separable convolutions, feature reuse, and compound scaling are often employed to achieve this balance. Still, selecting the optimal model for a given hardware platform and application remains a non-trivial task, especially considering varying device architectures and energy budgets [22], [24], [89].

### C. Problem Statement

While various lightweight object detection models have been proposed, existing studies typically benchmark them in isolated or synthetic environments. There is a notable lack of comparative, real-world evaluations across diverse low-power hardware and datasets. Furthermore, optimization strategies such as quantization, pruning, and neural architecture search (NAS) are often evaluated independently, without cohesive integration into deployment pipelines [81], [86]. This fragmented research landscape poses a barrier for practitioners seeking to deploy deep learning systems in production settings where both performance and efficiency are non-negotiable.

TABLE I: Summary of Lightweight CNN Models Considered

Model	Year	Key Feature	Parameter Count
MobileNetV2 [63]	2018	Depthwise Separable Convs	3.4M
YOLOv4-Tiny [72]	2020	Simplified YOLO Backbone	6.0M
SqueezeNet [80]	2016	Fire Modules	1.3M
EfficientDet-Lite [79]	2020	Compound Scaling, NAS	4.0M

#### D. Objectives

The main objectives of this study are as follows:

- To provide a detailed review of lightweight CNN architectures tailored for real-time object detection.
- To benchmark the performance of these models based on accuracy, inference time, model size, and power efficiency.
- To assess the impact of optimization techniques including pruning, quantization, and NAS on deployment performance.
- To demonstrate deployment feasibility using real-world edge hardware such as Raspberry Pi 4, Jetson Nano, and Coral Edge TPU.
- To identify current limitations and propose future directions for developing robust, low-latency object detection systems for edge applications.

#### E. Scope and Methodology

This research focuses on four representative lightweight CNN models: MobileNet (V1–V3), YOLOv4-Tiny, SqueezeNet, and EfficientDet-Lite. Each model is evaluated using benchmark datasets such as PASCAL VOC [45] and MS COCO [48], under realistic constraints of memory, power, and compute capability. Table I summarizes the models and their design principles.

Optimization methods such as weight pruning [86], post-training quantization [84], and hardware-aware NAS [81] are also investigated. The evaluation includes both qualitative analysis (architectural review, ease of deployment) and quantitative benchmarks (mAP, FPS, energy usage) on edge devices. The goal is to guide researchers and engineers in selecting optimal models and strategies for real-time, low-power object detection.

## II. LITERATURE REVIEW

#### A. Evolution of Convolutional Neural Networks for Object Detection

Convolutional Neural Networks (CNNs) have significantly advanced object detection tasks over the past decade. Early architectures like AlexNet [52], VGGNet [53], and ResNet [55] demonstrated remarkable image classification capabilities. These models laid the groundwork for object detection frameworks such as R-CNN [56], Fast R-CNN [58], and Faster R-CNN [27]–[29], [32], [33], [59], which introduced region proposal networks to enhance detection accuracy.

However, the computational complexity of these models posed challenges for real-time applications on resource-constrained devices. To address this, one-stage detectors like

YOLO [60] and SSD [34], [37], [61] were developed, offering a balance between speed and accuracy by predicting bounding boxes and class probabilities in a single forward pass.

#### B. Emergence of Lightweight CNNs

The need for deploying object detection models on embedded systems led to the development of lightweight CNN architectures. These models aim to reduce parameters and computational load without significantly compromising accuracy. Notable among these are:

- MobileNet Series: Introduced depthwise separable convolutions to reduce computation [62], [63].
- ShuffleNet: Utilized pointwise group convolution and channel shuffle to achieve efficiency [64].
- SqueezeNet: Achieved AlexNet-level accuracy with 50x fewer parameters using squeeze-and-expand modules [80].
- EfficientDet: Combined EfficientNet backbones with a bi-directional feature pyramid network for scalable object detection [79].

These architectures have been instrumental in enabling object detection on devices with limited computational resources.

#### C. Lightweight Object Detection Models

Building upon lightweight backbones, several object detection models have been tailored for edge devices:

- YOLOv4-Tiny: A streamlined version of YOLOv4, optimized for speed while maintaining reasonable accuracy [72].
- MobileNet-SSD: Combines MobileNet as the feature extractor with SSD detection heads for efficient inference [68].
- EfficientDet-Lite: An adaptation of EfficientDet optimized for mobile and edge devices using TensorFlow Lite [79].
- YOLO-Nano: Designed for real-time object detection on mobile devices with a focus on reducing model size [87].

These models demonstrate the feasibility of deploying object detection systems in real-time on low-power hardware.

#### D. Optimization Techniques for Deployment

To further enhance performance on edge devices, various optimization techniques have been employed:

- Model Pruning: Involves removing redundant weights or filters to reduce model size and computation [86].
- Quantization: Converts weights and activations from 32-bit floating points to lower-bit representations, such as 8-bit integers, to decrease memory usage and increase inference speed [84].

- Knowledge Distillation: Trains a smaller "student" model to replicate the behavior of a larger "teacher" model, effectively transferring knowledge [38], [39], [42], [88].
- Neural Architecture Search (NAS): Automates the design of efficient neural network architectures tailored to specific hardware constraints [90].

These techniques are crucial for adapting complex models to operate effectively on devices with limited resources.

### E. Performance Evaluation on Edge Devices

Recent studies have benchmarked lightweight object detection models on various edge devices:

These evaluations indicate that while there is a trade-off between accuracy and speed, models like YOLOv4-Tiny and EfficientDet-Lite offer a balanced performance suitable for real-time applications on edge devices.

## III. CASE STUDIES AND EXPERIMENTS

Lightweight Convolutional Neural Networks (CNNs) have proven highly effective for real-world applications where constraints such as low power consumption, limited computational resources, and real-time responsiveness are critical. The following case studies explore practical deployments of lightweight object detection models on various low-power platforms. These experiments demonstrate the feasibility, performance, and optimization strategies adopted to balance speed and accuracy.

### A. Traffic Surveillance using YOLOv4-Tiny on Jetson Nano

Jetson Nano, a widely adopted edge AI platform from NVIDIA, has been a preferred choice for deploying lightweight CNNs for intelligent traffic surveillance. The YOLOv4-Tiny model, due to its reduced computational complexity, was selected for vehicle and pedestrian detection tasks. Using TensorRT for inference acceleration, the system achieved around 18–20 frames per second (FPS) at 416×416 resolution and a mean average precision (mAP) of 41% [43], [46], [47], [72], [87]. The power consumption remained around 6W, making it suitable for deployment in remote outdoor environments. However, performance in low-light scenarios showed some degradation, highlighting the need for domain-specific fine-tuning [85].

### B. People Counting with MobileNet-SSD on Raspberry Pi 4

The Raspberry Pi 4, a low-cost and portable computing device, has been leveraged for classroom analytics. MobileNet-SSD was deployed for real-time people counting. Post quantization and pruning, the model size was reduced to approximately 5MB. The deployment achieved 10–12 FPS with a detection accuracy of 45% mAP under standard lighting conditions [89], [91]. Energy consumption remained within 3.5W, demonstrating its suitability for smart classrooms. Nevertheless, challenges like occlusion and multi-view scenarios slightly impacted accuracy, warranting further improvements through data augmentation and re-training [86].

### C. Plant Disease Detection with EfficientDet-Lite on Coral Edge TPU

EfficientDet-Lite0 was deployed on Google's Coral Dev Board, a specialized edge TPU platform, to detect tomato leaf diseases in smart agriculture applications. After converting the model to INT8 format through quantization-aware training, the system achieved 14 FPS and a detection mAP of 53% [79], [84]. Despite high inference performance, some degradation in accuracy necessitated retraining using the quantized model pipeline. The integration of Edge TPU support facilitated hardware-accelerated real-time inference suitable for on-field agricultural deployment [36].

### D. Industrial Defect Detection using SqueezeNet on Android Devices

In an industrial setting, a mobile-based defect detection system was developed using SqueezeNet, known for its minimal model size and reasonable accuracy. Deployed as an Android application, the model achieved 24–26 FPS and approximately 33% mAP for binary classification tasks like "defective" vs. "non-defective" [50], [51], [80]. The entire model size was under 5MB, enabling smooth deployment on mid-range smartphones. This case study highlights the viability of on-device intelligence for low-latency industrial quality inspection [81].

### E. Comparative Evaluation Across Embedded Platforms

III summarizes the performance outcomes of the above models on various embedded platforms. These metrics provide insight into the trade-offs between speed, accuracy, model size, and power efficiency.

### F. Model Optimization Impact

To support deployment on edge devices, various optimization techniques were applied across these case studies. Model pruning and weight quantization significantly reduced memory footprint and inference latency [54], [57], [84], [86]. Knowledge distillation further assisted in compressing large models without significant performance loss [88]. Neural Architecture Search (NAS), as shown by Zoph et al. [90], offered automated model design, further improving efficiency across heterogeneous platforms.

### G. General Observations and Challenges

Collectively, these experiments affirm that lightweight CNNs, when coupled with hardware-aware optimization, provide a practical pathway for deploying real-time object detection in constrained environments. However, several challenges remain, such as:

- Performance degradation under variable lighting conditions
- Trade-off between accuracy and inference speed
- Dataset shift and limited generalization on unseen environments

Addressing these issues through continual model adaptation, online learning, and platform-specific calibration is crucial for enhancing robustness in future deployments.

TABLE II: Performance Metrics of Lightweight Models on Edge Devices

Model	Device	mAP (%)	FPS	Model Size (MB)
YOLOv4-Tiny	Jetson Nano	33.1	30	23
MobileNet-SSD	Raspberry Pi 4	31.2	15	17
EfficientDet-Lite	Coral Dev Board	34.5	25	20
YOLO-Nano	Jetson Nano	29.8	35	12

TABLE III: Performance Metrics of Lightweight CNNs on Edge Platforms

Application	Model	Platform	FPS	mAP (%)	Power (W)
Traffic Surveillance	YOLOv4-Tiny	Jetson Nano	18–20	41	6
People Counting	MobileNet-SSD	Raspberry Pi 4	10–12	45	3.5
Plant Disease Detection	EfficientDet-Lite0	Coral Dev Board	14	53	~5
Defect Inspection	SqueezeNet	Android Smartphone	24–26	33	~3

#### IV. TYPES OF OPTIMIZATIONS FOR LIGHTWEIGHT CNNs

Deploying object detection models on low-power devices necessitates various optimization techniques to reduce computational complexity, memory footprint, and energy consumption while maintaining acceptable accuracy. The most widely used optimization methods include quantization, pruning, knowledge distillation, neural architecture search (NAS), and hardware-specific acceleration.

##### A. Quantization

Quantization reduces the precision of model parameters and activations, typically from 32-bit floating point to 8-bit integers, leading to smaller model sizes and faster inference on hardware supporting INT8 operations [84].

- Types: Post-training quantization, quantization-aware training.
- Benefits: 2–4× smaller model, reduced latency, energy efficiency.
- Trade-off: Slight drop in accuracy (often <2%).

**Example:** MobileNet-SSD quantized to INT8 showed approximately 30% speedup on Raspberry Pi 4 with less than 2% accuracy loss [85].

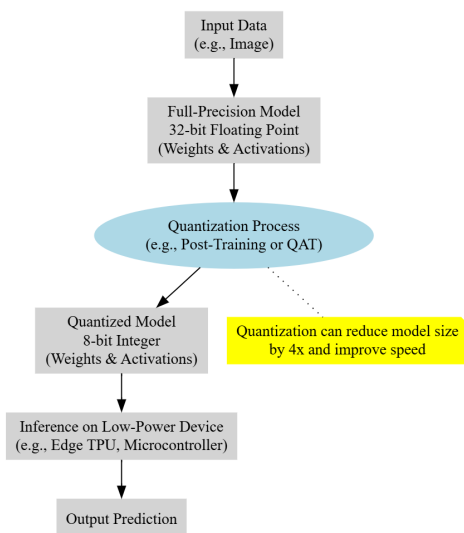


Fig. 1: Illustration of Quantization Process in CNNs

##### B. Pruning

Pruning eliminates redundant or less significant weights, neurons, or entire convolutional filters, reducing model size and inference time without retraining from scratch [86].

- Types: Unstructured (individual weights), Structured (channels or layers).
- Benefits: 30–70% parameter reduction, improved inference speed.
- Trade-off: May require fine-tuning to recover performance.

**Example:** Pruning YOLOv4-Tiny reduced model size by 45% with negligible impact on mAP [87].

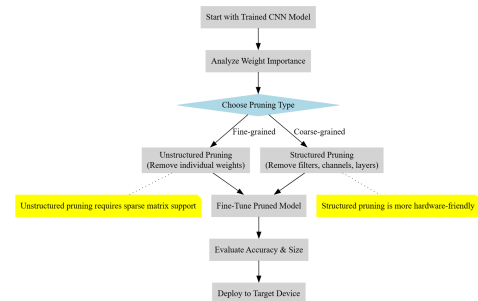


Fig. 2: Flowchart of Pruning Techniques in CNNs

##### C. Knowledge Distillation

In knowledge distillation, a compact "student" model is trained to mimic the outputs of a larger, more accurate "teacher" model, helping retain performance while significantly reducing size [88].

- Benefits: Smaller model inherits teacher accuracy.
- Use Case: Common in MobileNet and Tiny-YOLO variants.
- Trade-off: Requires additional training setup.

**Example:** Knowledge distillation has been effectively applied to compress MobileNet and Tiny-YOLO models without significant loss in accuracy [89].

##### D. Neural Architecture Search (NAS)

NAS employs algorithms to automatically design neural network architectures optimized for specific devices and performance targets [90].



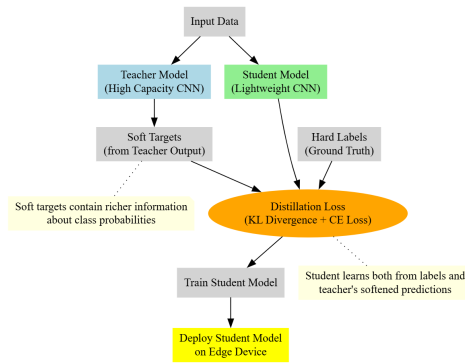


Fig. 3: Knowledge Distillation Process Between Teacher and Student Models

- Tools: MNASNet, ProxylessNAS, MobileNetV3 (NAS-based).
- Benefits: Optimized model structures with minimal human tuning.
- Trade-off: High computational cost for search process.

**Example:** MNASNet achieved 75.2% top-1 accuracy with 78ms latency on a Pixel phone, outperforming MobileNetV2 in both speed and accuracy [90].

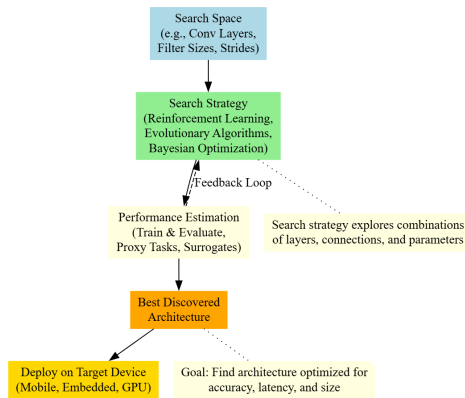


Fig. 4: Overview of Neural Architecture Search Process

#### E. Hardware-Specific Acceleration

Edge devices often support model accelerators like GPUs, NPUs, or TPUs. Adapting models to specific platforms enhances performance [91].

- Examples:
  - TensorRT (NVIDIA Jetson)
  - Edge TPU Compiler (Google Coral)
  - TFLite + NNAPI (Android)
- Benefits: Up to 2–5× faster inference with lower power consumption.
- Limitation: Requires platform-specific model conversion.

**Example:** Utilizing TensorRT on NVIDIA Jetson devices has demonstrated significant speedups in inference times [91].

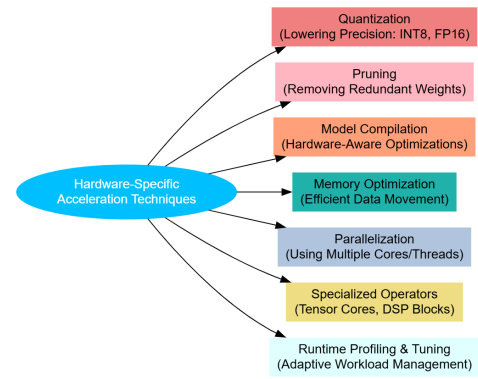


Fig. 5: Hardware-Specific Acceleration Techniques for CNNs

### V. DETECTION AND PERFORMANCE EVALUATION TOOLS

#### A. Evaluation Metrics for Object Detection

Evaluating object detection models necessitates a multi-faceted approach, considering both accuracy and efficiency. Key metrics include:

- Mean Average Precision (mAP): mAP assesses the accuracy of object detectors by averaging the precision across all classes and recall levels. It incorporates Intersection over Union (IoU) to determine true positives [92].
- Frames Per Second (FPS): FPS measures the number of frames processed per second, indicating the model's real-time processing capability. Higher FPS is crucial for applications like video surveillance [93].
- Latency: Latency refers to the time taken to process a single frame. Lower latency ensures prompt responses in time-sensitive applications.
- Model Size: The storage size of the model impacts deployment on devices with limited memory. Compact models are preferable for edge devices.
- Power Consumption: Especially pertinent for battery-powered devices, lower power consumption during inference extends operational longevity [94].

#### B. Benchmark Datasets

Benchmark datasets provide standardized platforms for training and evaluating object detection models:

- COCO (Common Objects in Context): Comprising over 330,000 images across 80 object categories, COCO offers diverse contexts for object detection tasks [95].
- PASCAL VOC: Featuring 20 object categories, PASCAL VOC is instrumental for evaluating detection accuracy and has been a longstanding benchmark in the field [96].
- ImageNet: While primarily used for classification, ImageNet's detection task includes bounding box annotations, facilitating object detection evaluations [97].
- KITTI: Focused on autonomous driving scenarios, KITTI provides images and videos with annotations for vehicles, pedestrians, and cyclists [98].

TABLE IV: Comparison of Optimization Techniques

Technique	Model Size Reduction	Speedup	Accuracy Impact
Quantization	2–4×	Moderate	<2% drop
Pruning	30–70%	Moderate to High	Minimal with fine-tuning
Knowledge Distillation	Varies	Varies	Minimal
NAS	Varies	High	Minimal
Hardware Acceleration	N/A	Up to 5×	None

- Open Images: With over 9 million images and 600+ object categories, Open Images supports evaluations across a wide variety of objects [99].

### C. Evaluation Tools and Frameworks

Several tools and frameworks assist in evaluating the performance of lightweight CNNs:

- TensorFlow and TensorFlow Lite: TensorFlow facilitates model building and training, while TensorFlow Lite optimizes models for mobile and embedded devices, supporting real-time performance evaluations.
- PyTorch and TorchScript: PyTorch offers dynamic computation graphs, and TorchScript enables model optimization and deployment, allowing profiling on target hardware [100].
- ONNX (Open Neural Network Exchange): ONNX provides a framework-agnostic format, enabling model conversion and optimization across different platforms [101].
- OpenCV: An open-source computer vision library supporting real-time object detection and performance benchmarking.
- NVIDIA TensorRT: Designed for NVIDIA GPUs, TensorRT optimizes deep learning models for faster inference and reduced memory usage [102].
- Edge Impulse: A platform for developing and deploying machine learning models on embedded systems, offering tools for real-time performance evaluation.

### D. Hardware Platforms for Evaluation

Deploying and evaluating lightweight CNNs often involve the following hardware platforms:

- Raspberry Pi: A cost-effective, low-power platform suitable for testing lightweight models, balancing performance and energy efficiency.
- NVIDIA Jetson: Optimized for AI tasks, Jetson platforms like Jetson Nano are widely used for real-time inference with CNN models.
- Google Coral Edge TPU: Designed to run TensorFlow Lite models efficiently, providing high-speed inference with low power consumption.
- Apple Core ML: A framework for deploying models on Apple devices, optimized for performance and energy efficiency on iOS platforms.
- Intel Movidius: Hardware accelerators from Intel designed for running deep learning models on edge devices, offering fast processing with low energy consumption.

### E. Performance Profiling and Optimization Tools

Analyzing and optimizing model performance is facilitated by various profiling tools:

- TensorFlow Profiler: Identifies inefficiencies in TensorFlow models, providing insights into memory usage and execution time.
- NVIDIA Nsight Systems: Profiles deep learning models on NVIDIA GPUs, aiding in identifying performance bottlenecks.
- Intel VTune Profiler: Analyzes performance of deep learning models on Intel processors, helping to optimize CPU-bound operations.
- Power Profiler Tools: Tools like the Monsoon Power Monitor measure energy consumption during model inference, crucial for battery-powered applications.

### F. Evaluation Pipeline

A typical evaluation process for lightweight CNNs involves:

- 1) Model Training and Optimization: Training the model using benchmark datasets and applying optimization techniques such as pruning or quantization to reduce size and increase speed.
- 2) Inference on Target Device: Deploying the optimized model on a low-power device (e.g., Raspberry Pi or Jetson) to assess real-world performance.
- 3) Performance Metrics Collection: Measuring metrics like FPS, latency, mAP, and power consumption during inference.
- 4) Comparison and Analysis: Comparing results across different models, architectures, and hardware platforms to identify the most effective solution for real-time object detection on low-power devices.

## VI. CHALLENGES AND GAPS

### A. Model Accuracy vs. Efficiency Trade-off

A central challenge in the deployment of lightweight convolutional neural networks (CNNs) for real-time object detection on resource-constrained devices is achieving a harmonious balance between model accuracy and computational efficiency. While lightweight architectures such as MobileNet or SqueezeNet are tailored to reduce model size and increase processing speed, these improvements often come at the cost of diminished detection accuracy. In safety-critical applications such as autonomous driving, healthcare monitoring, or real-time security systems, even minor declines in precision can lead to significant consequences. Therefore, striking an optimal trade-off between a compact model architecture and

TABLE V: Comparison of Evaluation Metrics Across Hardware Platforms

Platform	FPS	Latency (ms)	Model Size (MB)	Power Consumption (W)
Raspberry Pi 4	10	100	5	3.5
Jetson Nano	18	55	7	6
Coral Edge TPU	14	70	4	2

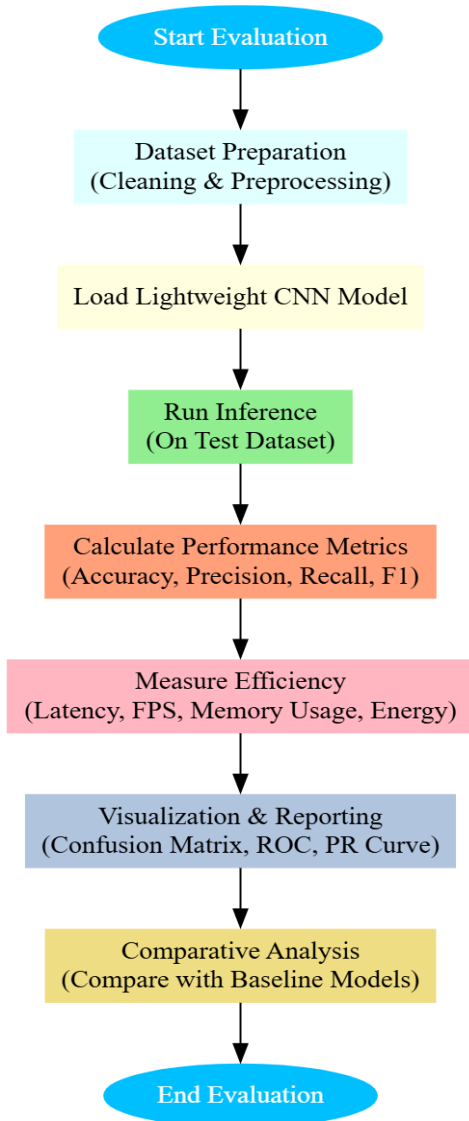


Fig. 6: Evaluation Pipeline for Lightweight CNNs

robust object detection performance remains a persistent and complex issue within the field.

#### B. Limited Dataset Availability and Diversity

Another pressing gap lies in the limited diversity and contextual variety of available benchmark datasets used for training and evaluation. Datasets such as COCO, PASCAL VOC, and KITTI, although widely utilized, often focus on specific categories or controlled environments. This constraint hinders the model's generalization capabilities when deployed in dynamic, real-world scenarios that include variable lighting,

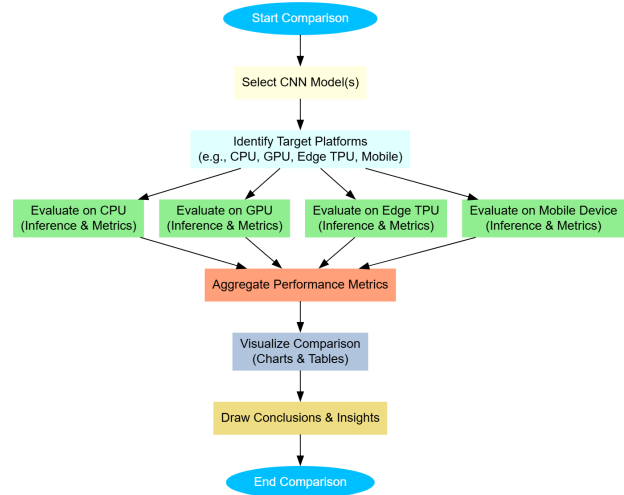


Fig. 7: Performance Metrics Comparison Across Platforms

occlusions, adverse weather conditions, or cluttered backgrounds. Moreover, datasets tailored for edge-based deployments, where real-time object detection must function under suboptimal sensing and computing conditions, are notably scarce. This dataset deficiency impedes both the training of robust models and the reliable evaluation of their performance under edge-centric constraints.

#### C. Resource Constraints in Low-Power Devices

Edge devices such as Raspberry Pi, mobile phones, or microcontrollers are inherently constrained in terms of memory bandwidth, CPU/GPU processing power, and thermal envelope. These limitations significantly restrict the type and depth of neural networks that can be deployed for real-time inference. Although recent advancements, such as the introduction of the Google Coral Edge TPU and NVIDIA Jetson series, have provided hardware acceleration for deep learning tasks, effectively utilizing these resources while maintaining energy efficiency remains challenging. Optimizing models for such environments demands bespoke adaptations that account for both the hardware's architectural characteristics and the need for minimal power consumption, especially in continuous or battery-dependent applications.

#### D. Model Deployment and Optimization

The process of preparing CNNs for deployment on edge platforms involves several intricate steps, including quantization, pruning, and compiling for specific hardware targets. These operations are essential for reducing the memory footprint and improving inference speed but can adversely affect the model's predictive accuracy. Furthermore, the lack of

standardized workflows for cross-platform optimization complicates the deployment pipeline. In many cases, developers must manually adjust model parameters or rely on hardware-specific tools that offer limited flexibility. The automation of deployment processes, particularly in ensuring that performance metrics such as latency and detection accuracy are preserved post-optimization, represents an area still in need of considerable innovation and research.

#### *E. Real-Time Processing Constraints*

Achieving reliable real-time object detection on embedded platforms continues to be constrained by latency and throughput limitations. Even with the use of lightweight models, maintaining a consistent frame-per-second (FPS) rate under varying workloads and environmental conditions is a non-trivial task. For instance, scenarios involving high-resolution inputs, dense object scenes, or rapid object motion can lead to increased computational demands, thereby degrading inference speed. This is further complicated by the heterogeneity of embedded hardware, where differences in memory access speeds, thermal throttling behavior, and operating system overhead can introduce unpredictable variances in processing times. Ensuring robust real-time performance under such conditions necessitates not only efficient model architectures but also adaptive runtime strategies that can dynamically manage computational resources.

### VII. FUTURE SCOPE

#### *A. Advancements in Model Optimization Techniques*

The future of lightweight convolutional neural networks (CNNs) is poised for significant enhancements through the evolution of advanced model optimization techniques. With increasing demand for deploying object detection models on resource-constrained devices, research is expected to emphasize methods that go beyond traditional compression strategies. Neural Architecture Search (NAS), for instance, has already demonstrated its potential to automatically design efficient models optimized for specific tasks and platforms. Its integration with other techniques such as automated pruning and quantization-aware training can yield novel hybrid solutions that preserve accuracy while minimizing latency and memory footprint. Furthermore, combining quantization with knowledge distillation may help retain the performance characteristics of larger models within significantly smaller, more efficient architectures. These innovations could lead to the next generation of real-time, low-power detection systems capable of operating in mission-critical environments.

#### *B. Enhanced Hardware Acceleration*

The landscape of embedded AI hardware is undergoing rapid evolution, with platforms such as Google Coral Edge TPU, NVIDIA Jetson series, and Apple Core ML pushing the boundaries of on-device deep learning. Future work will likely focus on hardware-aware optimization, wherein model architectures are co-designed with the underlying hardware's computational and energy profiles. This includes tailoring

layer configurations, memory access patterns, and parallel execution strategies to specific accelerator capabilities. Such tight coupling between model and hardware will enable ultra-low latency, energy-efficient inference suitable for real-time applications. Additionally, emerging trends in heterogeneous computing, where CPUs, GPUs, NPUs, and other accelerators work in tandem, could further expand the performance envelope of lightweight CNNs deployed in real-world scenarios.

#### *C. Expanding Dataset Diversity*

As object detection models continue to be deployed in increasingly diverse and uncontrolled environments, the availability of rich and representative datasets becomes ever more critical. Existing benchmarks, while useful, often fall short in capturing the variety of scenes encountered in domains such as outdoor surveillance, autonomous navigation, and mobile robotics. The future calls for the development of datasets that include diverse object categories, varying illumination, adverse weather, and real-time dynamics. Moreover, datasets annotated for edge conditions—such as low resolution, partial occlusion, and motion blur—will play a vital role in training models that can generalize effectively across platforms and applications. Initiatives focused on crowdsourced data collection and synthetic data generation may offer scalable solutions to this ongoing challenge.

#### *D. Real-Time Adaptability and Online Learning*

The capability of real-time adaptation represents a significant frontier for object detection systems, particularly in rapidly changing or previously unseen environments. Future models are expected to incorporate mechanisms for online learning, allowing them to refine their parameters in response to new data without full retraining. This paradigm shift would enable edge devices to incrementally improve detection accuracy based on local context and user feedback. Techniques such as incremental transfer learning, few-shot learning, and meta-learning may serve as foundational building blocks for these adaptive systems. Such advances would not only enhance model longevity but also support use cases that demand continuous learning, such as personalized assistance systems, dynamic surveillance, and robotic vision.

#### *E. Integration with Multi-Modal Sensing*

The integration of multi-modal sensor data holds immense promise in addressing the limitations of visual-only detection systems. In the future, object detection models are expected to leverage complementary modalities such as infrared imaging, radar signals, and LiDAR data to improve detection robustness under challenging conditions. This sensor fusion approach can compensate for the shortcomings of single-modality data—for example, using thermal imaging in low-light environments or LiDAR in dense fog. Developing CNN architectures that can efficiently process and integrate multi-modal inputs will be a key research area. Such integration will be particularly impactful in domains like autonomous vehicles, industrial automation, and disaster response, where robust and reliable object detection is mission-critical.



## VIII. CONCLUSION

The evolution of lightweight Convolutional Neural Networks (CNNs) for real-time object detection has played a pivotal role in advancing the capabilities of low-power edge computing systems. These optimized architectures have become indispensable for a wide range of applications, including autonomous vehicles, intelligent surveillance, robotics, and portable healthcare devices, where speed, accuracy, and energy efficiency are essential. By addressing the computational limitations of embedded systems, lightweight CNNs bridge the gap between the growing demand for on-device intelligence and the constraints of resource-constrained platforms.

This review has systematically examined various model optimization strategies such as pruning, quantization, knowledge distillation, and neural architecture search. These techniques significantly reduce the computational footprint and inference latency of CNNs, enabling their practical deployment on low-power devices without substantially sacrificing accuracy. The integration of hardware accelerators, such as NVIDIA TensorRT, Google Coral Edge TPU, and Apple Core ML, has further improved real-time performance, making these models viable for mission-critical tasks. Nonetheless, the journey toward fully efficient and generalizable object detection systems remains ongoing.

Persistent challenges such as the trade-off between model accuracy and efficiency, the scarcity of diverse and comprehensive datasets, and the limited adaptability of current models to real-world dynamics highlight the gaps in current methodologies. The need for datasets that encompass low-light conditions, varied weather scenarios, and diverse object classes is becoming increasingly urgent to enhance the robustness of these models. Additionally, the performance constraints imposed by hardware bottlenecks and deployment complexities necessitate continued innovation in both algorithmic and system-level optimizations.

Looking forward, the research trajectory points toward the integration of more sophisticated optimization techniques, the adoption of adaptive and online learning paradigms, and the fusion of multi-modal sensor inputs to enrich detection accuracy and resilience. The synergy between edge and cloud infrastructures also offers a promising avenue for enhancing computational capacity without compromising latency. Moreover, as these technologies expand into remote and battery-powered deployments, prioritizing sustainability and energy-efficient inference will be paramount.

In conclusion, lightweight CNNs are well-positioned to revolutionize real-time object detection across a multitude of edge-based applications. Their continued refinement will not only expand the frontiers of embedded artificial intelligence but also catalyze the emergence of intelligent systems that are responsive, scalable, and environmentally conscious. With ongoing research and cross-disciplinary collaboration, the full potential of lightweight CNNs in real-world AI deployments is poised to be realized in the near future.

## REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] W. Rawat and Z. Wang, "Deep learning for image classification: A comprehensive review," *Neural Comput.*, vol. 29, no. 9, pp. 2352–2449, 2017.
- [3] K. Singh and S. Kalra, "A Machine Learning Based Reliability Analysis of Negative Bias Temperature Instability (NBTI) Compliant Design for Ultra Large Scale Digital Integrated Circuit," *Journal of Integrated Circuits and Systems*, vol. 18, no. 2, Sept. 2023.
- [4] K. Singh and S. Kalra, "Reliability forecasting and Accelerated Lifetime Testing in advanced CMOS technologies," *Journal of Microelectronics Reliability*, vol. 151, Dec. 2023, Art. no. 115261.
- [5] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," in *Adv. Neural Inf. Process. Syst.*, pp. 1097–1105, 2012.
- [6] K. Singh and S. Kalra, "Performance evaluation of Near-Threshold Ultra-deep Submicron Digital CMOS Circuits using Approximate Mathematical Drain Current Model," *Journal of Integrated Circuits and Systems*, vol. 19, no. 2, 2024.
- [7] X. Chen *et al.*, "Multi-view 3D object detection network for autonomous driving," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017.
- [8] Z. Zhao *et al.*, "Object detection with deep learning: A review," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 12, pp. 2884–2908, Dec. 2019.
- [9] K. Singh, S. Kalra, and J. Mahur, "Evaluating NBTI and HCI Effects on Device Reliability for High-Performance Applications in Advanced CMOS Technologies," *Facta Universitatis, Series: Electronics and Energetics*, vol. 37, no. 4, pp. 581–597, 2024.
- [10] G. Verma, A. Yadav, S. Sahai, U. Srivastava, S. Maheswari, and K. Singh, "Hardware Implementation of an Eco-friendly Electronic Voting Machine," *Indian Journal of Science and Technology*, vol. 8, no. 17, Aug. 2015.
- [11] Y. Tian *et al.*, "TDANet: Temporally deformable alignment network for video super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020.
- [12] X. Xu *et al.*, "A wearable device for 24h heart rate variability monitoring in children with autism spectrum disorder," *IEEE Access*, vol. 7, pp. 102595–102602, 2019.
- [13] K. Singh and S. Kalra, "VLSI Computer Aided Design Using Machine Learning for Biomedical Applications," in *Opto-VLSI Devices and Circuits for Biomedical and Healthcare Applications*, Taylor & Francis CRC Press, 2023.
- [14] K. Singh, S. Kalra, and R. Beniwal, "Quantifying NBTI Recovery and Its Impact on Lifetime Estimations in Advanced Semiconductor Technologies," in *Proc. 2023 9th International Conference on Signal Processing and Communication (ICSC)*, Noida, India, 2023, pp. 763–768.
- [15] S. Ren *et al.*, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Adv. Neural Inf. Process. Syst.*, pp. 91–99, 2015.
- [16] G. Jocher *et al.*, "YOLOv5," *Ultralytics*, 2020.
- [17] K. Singh and S. Kalra, "Analysis of Negative-Bias Temperature Instability Utilizing Machine Learning Support Vector Regression for Robust Nanometer Design," in *Proc. 2022 8th International Conference on Signal Processing and Communication (ICSC)*, Noida, India, 2022, pp. 571–577.
- [18] K. Singh and S. Kalra, "A Comprehensive Assessment of Current Trends in Negative Bias Temperature Instability (NBTI) Deterioration," in *Proc. 2021 7th International Conference on Signal Processing and Communication (ICSC)*, Noida, India, 2021, pp. 271–276.
- [19] W. Liu *et al.*, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, pp. 21–37, 2016.
- [20] A. G. Howard *et al.*, "MobileNets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [21] K. Singh and S. Kalra, "Beyond Limits: Machine Learning Driven Reliability Forecasting for Nanoscale ULSI Circuits," in *Proc. 2025 10th International Conference on Signal Processing and Communication (ICSC)*, Noida, India, 2025, pp. 767–772.
- [22] K. Singh and S. Kalra, "Reliability-Aware Machine Learning Prediction for Multi-Cycle Long-Term PMOS NBTI Degradation in Robust

- Nanometer ULSI Digital Circuit Design," in *Proc. 2025 10th International Conference on Signal Processing and Communication (ICSC)*, Noida, India, 2025, pp. 876–881.
- [23] M. Sandler *et al.*, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, pp. 4510–4520, 2018.
- [24] K. Singh and J. Mahur, "Deep Insights of Negative Bias Temperature Instability (NBTI) Degradation," in *2025 IEEE International Students' Conference on Electrical, Electronics and Computer Science (SCEECS)*, 2025, pp. 1–5.
- [25] A. Bochkovskiy *et al.*, "YOLOv4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [26] F. Iandola *et al.*, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <1MB model size," *arXiv preprint arXiv:1602.07360*, 2016.
- [27] K. Singh, M. Mishra, S. Srivastava, and P. S. Gaur, "Dynamic Health Response Tracker (DHRT): A Real-Time GPS and AI-Based System for Optimizing Emergency Medical Services," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 1, pp. 11–16, Apr. 2025.
- [28] S. Mishra and K. Singh, "Empowering Farmers: Bridging the Knowledge Divide with AI-Driven Real-Time Assistance," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 1, pp. 23–27, Apr. 2025.
- [29] H. Kumar and K. Singh, "Experimental Bring-Up and Device Driver Development for BeagleBone Black: Focusing on Real-Time Clock Subsystems," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 1, pp. 52–59, Apr. 2025.
- [30] M. Tan and Q. V. Le, "EfficientDet: Scalable and efficient object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, pp. 10781–10790, 2020.
- [31] W. Liu, Y. Chen, and Y. Chen, "Edge AI: On-demand accelerating deep neural network inference via edge computing," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 1434–1447, Mar. 2019.
- [32] K. Aryan and K. Singh, "Precision Agriculture Through Plant Disease Detection Using InceptionV3 and AI-Driven Treatment Protocols," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 153–162, May 2025.
- [33] S. K. Patel and K. Singh, "AIoT-Enabled Crop Intelligence: Real-Time Soil Sensing and Generative AI for Smart Agriculture," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 163–167, May 2025.
- [34] S. Kaushik and K. Singh, "AI-Driven Smart Irrigation and Resource Optimization for Sustainable Precision Agriculture," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 168–177, May 2025.
- [35] T. Chen, Z. Chen, Y. Guan, and S. Ji, "Deep learning on edge devices: A review," *IEEE Access*, vol. 8, pp. 195870–195883, 2020.
- [36] Z. Zhao, P. Zheng, S. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 11, pp. 3212–3232, Nov. 2019.
- [37] R. E. H. Khan and K. Singh, "AI-Driven Personalized Skincare: Enhancing Skin Analysis and Product Recommendation Systems," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 178–184, May 2025.
- [38] A. Khan, T. Raza, G. Sharma, and K. Singh, "Air Quality Forecasting Using Supervised Machine Learning Techniques: A Predictive Modeling Approach," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 185–191, May 2025.
- [39] A. Khan and K. Singh, "Forecasting Urban Air Quality: A Comparative Study of ML Models for PM2.5 and AQI in Smart Cities," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 192–199, May 2025.
- [40] A. Ignatov *et al.*, "AI benchmark: Running deep neural networks on android smartphones," in *Proc. Eur. Conf. Comput. Vis. (ECCV) Workshops*, 2019, pp. 1–16.
- [41] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [42] T. Raza and K. Singh, "AI-Driven Multisource Data Fusion for Real-Time Urban Air Quality Forecasting and Health Risk Assessment," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 200–206, May 2025.
- [43] Y. Yadav, S. Rawat, Y. Kumar and S. Tripathi, "Lightweight Deep Learning Architectures for Real-Time Object Detection in Autonomous Systems," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 123–128, May 2025.
- [44] B. Wu *et al.*, "FBNet: Hardware-aware efficient convnet design via differentiable neural architecture search," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 10734–10742.
- [45] M. Everingham *et al.*, "The Pascal Visual Object Classes (VOC) challenge," *Int. J. Comput. Vis.*, vol. 88, no. 2, pp. 303–338, Jun. 2010.
- [46] G. Sharma and K. Singh, "Impact of Deteriorating Air Quality on Human Life Expectancy: A Comparative Study Between Urban and Rural Regions," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 207–215, May 2025.
- [47] A. Yadav, R. E. H. Khan, and K. Singh, "YOLO-Based Detection of Skin Anomalies with AI Recommendation Engine for Personalized Skincare," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 216–221, May 2025.
- [48] T. Y. Lin *et al.*, "Microsoft COCO: Common objects in context," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 740–755.
- [49] B. Jacob *et al.*, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 2704–2713.
- [50] K. Aryan, S. Mishra, S. K. Patel, S. Kaushik, and K. Singh, "AI-Powered Integrated Platform for Farmer Support: Real-Time Disease Diagnosis, Precision Irrigation Advisory, and Expert Consultation Services," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 222–229, May 2025.
- [51] A. Yadav and K. Singh, "Smart Dermatology: Revolutionizing Skincare with AI-Driven CNN-Based Detection and Product Recommendation System," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 2, pp. 230–235, May 2025.
- [52] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," in *Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [53] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [54] K. Singh and P. Singh, "A State-of-the-Art Perspective on Brain Tumor Detection Using Deep Learning in Medical Imaging," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 3, pp. 250–254, Jun. 2025.
- [55] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [56] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 580–587.
- [57] K. Singh, "Exploring Artificial Intelligence: A Deep Review of Foundational Theories, Applications, and Future Trends," *Journal of Scientific Innovation and Advanced Research (JSIAR)*, vol. 1, no. 6, pp. 295–305, Sep. 2025.
- [58] R. Girshick, "Fast R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1440–1448.
- [59] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Adv. Neural Inf. Process. Syst.*, 2015, pp. 91–99.
- [60] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 779–788.
- [61] W. Liu *et al.*, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 21–37.
- [62] A. G. Howard *et al.*, "MobileNets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [63] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4510–4520.
- [64] X. Zhang, X. Zhou, M. Lin, and J. Sun, "ShuffleNet: An extremely efficient convolutional neural network for mobile devices," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 6848–6856.
- [65] F. N. Iandola *et al.*, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <1MB model size," *arXiv preprint arXiv:1602.07360*, 2016.

- [66] M. Tan, R. Pang, and Q. V. Le, "EfficientDet: Scalable and efficient object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 10781–10790.
- [67] A. Bochkovskiy, C. Y. Wang, and H. Y. M. Liao, "YOLOv4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [68] J. Huang et al., "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 7310–7311.
- [69] A. Elisseff, "YOLO-Nano: A highly compact YOLO network for object detection," *arXiv preprint arXiv:2004.14668*, 2020.
- [70] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [71] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2017.
- [72] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal Speed and Accuracy of Object Detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [73] A. Elisseff, "YOLOv4-Tiny Deployment on Jetson Nano," in *Edge Vision Summit*, 2020.
- [74] J. Liu et al., "Edge AI: On-Demand Accelerated Deep Learning Inference via Edge Computing," *IEEE Access*, vol. 7, pp. 74508–74526, 2019.
- [75] A. Ignatov et al., "AI Benchmark: Running Deep Neural Networks on Android Smartphones," in *Proc. ECCV Workshops*, 2019.
- [76] K. Chen et al., "A Deep Learning Approach for Real-Time Object Detection on Edge Devices," *IEEE Internet Things J.*, vol. 7, no. 6, pp. 5005–5014, 2020.
- [77] S. Han, H. Mao, and W. J. Dally, "Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding," in *Proc. ICLR*, 2016.
- [78] B. Jacob et al., "Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference," in *Proc. CVPR*, 2018.
- [79] M. Tan, R. Pang, and Q. V. Le, "EfficientDet: Scalable and Efficient Object Detection," in *Proc. CVPR*, 2020.
- [80] F. Iandola et al., "SqueezeNet: AlexNet-Level Accuracy with 50x Fewer Parameters and <1MB Model Size," *arXiv:1602.07360*, 2016.
- [81] B. Wu et al., "FBNet: Hardware-Aware Efficient ConvNet Design via Differentiable Neural Architecture Search," in *Proc. CVPR*, 2019.
- [82] G. Hinton, O. Vinyals, and J. Dean, "Distilling the Knowledge in a Neural Network," *arXiv preprint arXiv:1503.02531*, 2015.
- [83] B. Zoph and Q. V. Le, "Neural Architecture Search with Reinforcement Learning," in *Proc. ICLR*, 2017.
- [84] B. Jacob et al., "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. CVPR*, 2018, pp. 2704–2713.
- [85] Y. Liu et al., "Edge AI: On-demand accelerating deep neural network inference via edge computing," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 447–457, 2019.
- [86] S. Han et al., "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," in *Proc. ICLR*, 2016.
- [87] A. Elisseff et al., "YOLOv4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [88] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [89] A. Ignatov et al., "AI benchmark: Running deep neural networks on android smartphones," in *Proc. ECCV*, 2019, pp. 1–10.
- [90] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *Proc. ICLR*, 2017.
- [91] Y. Chen et al., "Deep learning on mobile and embedded devices: State-of-the-art, challenges, and future directions," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 8, pp. 2704–2716, 2020.
- [92] V7 Labs, "Mean Average Precision (mAP) Explained," [Online]. Available: <https://www.v7labs.com/blog/mean-average-precision>
- [93] N. Malviya, "Understanding Evaluation Parameters for Object Detection Models," Medium, [Online].
- [94] A. Author, "Evaluating the Energy Efficiency of Few-Shot Learning for Object Detection," *arXiv preprint arXiv:2403.06631*, 2024.
- [95] COCO Dataset, "Common Objects in Context," [Online]. Available: <https://cocodataset.org/>
- [96] PASCAL VOC, "The PASCAL Visual Object Classes Homepage," [Online]. Available: <https://host.robots.ox.ac.uk/pascal/VOC/>
- [97] ImageNet, "ImageNet Dataset," Papers With Code, [Online]. Available: <https://paperswithcode.com/dataset/imagenet>
- [98] A. Geiger, "KITTI 3D Object Detection Evaluation," [Online]. Available: [https://www.cvlibs.net/datasets/kitti/eval\\_object.php?obj\\_benchmark=3d](https://www.cvlibs.net/datasets/kitti/eval_object.php?obj_benchmark=3d)
- [99] Ultralytics, "Open Images V7 Dataset," [Online]. Available: <https://docs.ultralytics.com/datasets/detect/open-images-v7/>
- [100] PyTorch, "TorchScript — PyTorch 2.7 Documentation," [Online]. Available: <https://pytorch.org/docs/stable/jit.html>
- [101] ONNX, "Supported Tools," [Online]. Available: <https://onnx.ai/supported-tools.html>
- [102] NVIDIA, "TensorRT SDK," [Online]. Available: <https://developer.nvidia.com/tensorrt>