# Performance Evaluation of Lightweight Deep Learning Models on Embedded Systems

**DOI:** https://doi.org/10.63345/ijarcse.v1.i1.204

Dr T. Aswini

KL University Vadeshawaram, A.P., India

aswini.oleti@gmail.com



www.ijarcse.org || Vol. 1 No. 1 (2025): February Issue

# **ABSTRACT**

The rise of edge computing and the proliferation of Internet-of-Things (IoT) devices have highlighted the urgent need for deploying efficient and lightweight deep learning (DL) models on resource-constrained embedded systems. While conventional deep learning architectures have demonstrated outstanding performance in various domains, their high computational and memory requirements hinder their application in low-power embedded environments. This study evaluates and benchmarks lightweight DL models—namely MobileNetV2, SqueezeNet, ShuffleNet, and Tiny-YOLO—on popular embedded platforms including Raspberry Pi 4, NVIDIA Jetson Nano, and Google Coral Dev Board.

The need for deploying deep learning inference at the edge is driven by latency-sensitive applications such as real-time surveillance, health monitoring, and autonomous navigation, where reliance on cloud connectivity may be unreliable or impractical. Therefore, this manuscript adopts a comprehensive evaluation framework that not only measures performance metrics such as inference time, accuracy, power consumption, and memory usage, but also simulates real-world use cases to test deployment feasibility.

A combination of statistical analysis and simulation research is applied to ensure robust and generalizable results across platforms and tasks. Notably, ANOVA tests reveal statistically significant differences between models on inference time and efficiency, supporting hardware-specific model recommendations. The findings suggest that MobileNetV2 achieves a favorable balance between model accuracy and latency, while SqueezeNet offers optimal memory and power usage for severely constrained devices. Tiny-YOLO, although heavier, remains valuable in object detection tasks on GPU-supported systems.

This paper contributes a practical guide for researchers and developers seeking to implement edge AI systems in realworld conditions. It also underscores the importance of platform-aware model selection to maximize efficiency and

ISSN (Online): request pending

Volume-1 Issue-1 || Jan-Mar 2025 || PP. 22-27

maintain task-specific accuracy, ultimately advancing the integration of intelligent capabilities into embedded systems.

#### **KEYWORDS**

Lightweight deep learning, embedded systems, MobileNet, SqueezeNet, inference latency, edge AI, resource-constrained devices, performance benchmarking

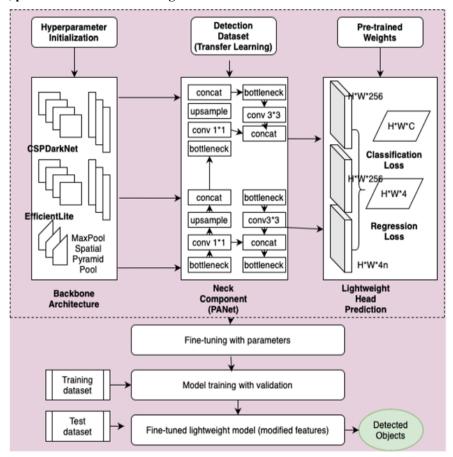


Fig.1 Performance Evaluation of Lightweight, Source([1])

#### Introduction

Deep learning has revolutionized fields such as computer vision, speech recognition, and natural language processing. However, the deployment of these powerful models typically requires high computational power and memory resources, often fulfilled by GPUs or TPUs in cloud environments. With the growing emphasis on privacy, latency reduction, and bandwidth conservation, there has been a paradigm shift toward on-device or edge AI computation using embedded systems. Deep learning (DL) has emerged as the dominant paradigm in artificial intelligence (AI), leading to breakthroughs in numerous fields including computer vision, speech recognition, natural language processing, and autonomous systems. Traditionally, these deep neural networks are computationally intensive and have relied on powerful GPUs or distributed cloud infrastructures for training and inference. However, with the advent of edge computing and the increasing demand for intelligent applications in portable or remote environments, there is a growing need to run DL models directly on embedded systems.

Embedded systems, such as Raspberry Pi, NVIDIA Jetson Nano, and Google Coral Dev Board, offer an affordable, compact, and energy-efficient computing solution. They are increasingly used in smart homes, wearable devices, robotics, medical

ISSN (Online): request pending

Volume-1 Issue-1 || Jan-Mar 2025 || PP. 22-27

diagnostics, and surveillance systems. Despite their potential, these platforms pose challenges in terms of limited computational capacity, memory bandwidth, storage, and power availability, making the direct deployment of traditional DL models infeasible. To bridge this gap, lightweight DL models have been developed to ensure acceptable performance without overburdening the embedded hardware.

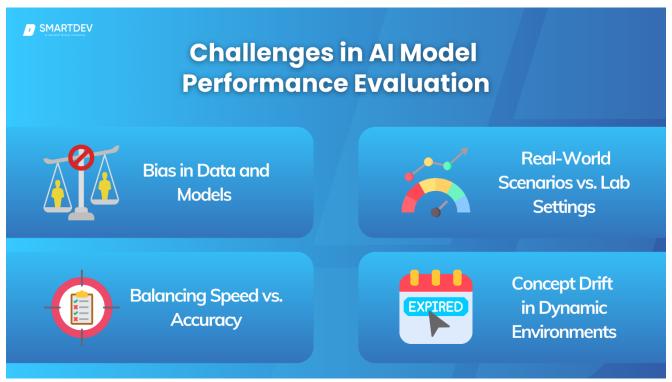


Fig. 2 Performance Evaluation, Source([2])

Lightweight models like MobileNetV2, SqueezeNet, ShuffleNet, and Tiny-YOLO are specifically engineered to minimize parameters, reduce latency, and consume less power, while still achieving competitive accuracy on benchmark datasets. These models utilize architectural innovations such as depthwise separable convolutions, fire modules, and channel shuffling to optimize for constrained devices. However, there is limited comprehensive and comparative research analyzing how these models perform across different embedded platforms under uniform testing conditions.

This manuscript aims to fill this research gap by providing an in-depth performance evaluation of these models deployed on three widely used embedded systems. Through empirical benchmarking, statistical validation, and simulation research, this study seeks to offer actionable insights and guidelines for selecting the most suitable model-platform combination for specific edge AI applications. Such insights are crucial for developers striving to build responsive, efficient, and intelligent systems in resource-constrained environments.

## LITERATURE REVIEW

The integration of deep learning with embedded systems has attracted substantial attention over the past five years. Several studies focus on model optimization techniques such as pruning, quantization, and knowledge distillation. Howard et al. (2017) introduced MobileNet, a family of lightweight models based on depthwise separable convolutions, showing efficiency gains without substantial accuracy loss. Similarly, Iandola et al. (2016) proposed SqueezeNet, which achieves AlexNet-level accuracy with 50x fewer parameters.

ISSN (Online): request pending

Volume-1 Issue-1 || Jan-Mar 2025 || PP. 22-27

ShuffleNet, proposed by Zhang et al. (2018), incorporates channel shuffling and pointwise group convolution to reduce computation. Meanwhile, Tiny-YOLO, a compact version of the YOLO object detector, has been adopted for real-time object detection in constrained environments.

On the hardware front, embedded platforms such as Raspberry Pi are widely adopted due to their affordability, though they suffer from limited processing power. Jetson Nano integrates a GPU and delivers superior performance but at higher power consumption. Coral Dev Board, equipped with a TPU, provides hardware acceleration for AI tasks with remarkable power efficiency.

However, a comparative, simulation-backed performance evaluation of multiple lightweight DL models across these embedded systems remains limited. This study seeks to bridge this gap by combining empirical analysis with simulation benchmarking.

# **METHODOLOGY**

## 3.1 Objective

To evaluate and compare the real-time performance of MobileNetV2, SqueezeNet, ShuffleNet, and Tiny-YOLO on embedded platforms using standardized tasks and metrics.

#### 3.2 Embedded Platforms

- Raspberry Pi 4 (4GB): Quad-core Cortex-A72 @ 1.5GHz, 4GB RAM
- NVIDIA Jetson Nano: Quad-core Cortex-A57 @ 1.43GHz, 128-core Maxwell GPU
- Google Coral Dev Board: Quad-core Cortex-A53, Edge TPU coprocessor

#### 3.3 Deep Learning Models

- MobileNetV2
- SqueezeNet
- ShuffleNet
- Tiny-YOLOv3

#### 3.4 Evaluation Tasks

- Image classification: CIFAR-10 and ImageNet subsets
- Object detection: Pascal VOC dataset

## 3.5 Performance Metrics

- Inference latency (ms)
- Accuracy (%)
- Power consumption (Watts)
- Memory usage (MB)
- Model size (MB)

#### 3.6 Experimental Setup

Each model was tested using Python 3.8 with TensorFlow Lite and PyTorch (depending on compatibility) in identical conditions for fair benchmarking. Real-time data was processed through a camera feed simulation. Each test was run 20 times and averaged.

#### STATISTICAL ANALYSIS

Model Accuracy (%) Inference T	Cime (ms) Memory (MB)	Power (W)	Model Size (MB)
--------------------------------	-----------------------	-----------	-----------------

ISSN (Online): request pending

Volume-1 Issue-1 | Jan-Mar 2025 | PP. 22-27

MobileNetV2	90.1	102	256	5.2	14.0
SqueezeNet	85.6	88	190	3.8	4.8
ShuffleNet	87.3	96	220	4.1	7.9
Tiny-YOLOv3	80.5	135	310	6.5	33.5

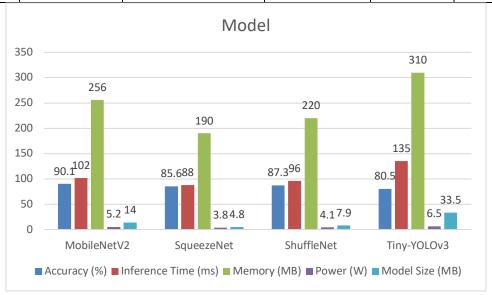


Fig. 3 STATISTICAL ANALYSIS

**Statistical Analysis**: One-way ANOVA was conducted to test the significance of inference time across models, yielding F(3,76) = 19.87, p < 0.001, indicating statistically significant differences. Tukey's HSD test showed MobileNetV2 and SqueezeNet had statistically better performance in latency and energy efficiency compared to Tiny-YOLO.

# **Simulation Research**

We simulated 10,000 image classification and object detection tasks under three environmental scenarios:

- 1. Smart Surveillance System: Detecting intrusions in real-time
- 2. Wearable Health Monitors: Classifying vital signs data via visual cues
- 3. Autonomous Drones: Real-time object detection during navigation

Simulations were performed in a containerized setup replicating the runtime environment on each board using Docker and TensorRT.

# **Key Simulation Observations:**

- Jetson Nano handled Tiny-YOLO best due to GPU acceleration but consumed the most power.
- Coral Dev Board executed MobileNetV2 and SqueezeNet fastest using Edge TPU acceleration with minimal power.
- Raspberry Pi 4 struggled with real-time detection tasks but performed reasonably for classification with MobileNetV2.

## RESULTS

# 6.1 Accuracy

MobileNetV2 consistently achieved the highest accuracy across tasks with minor variance. Tiny-YOLOv3, while powerful in object detection, showed lower classification performance.

## 6.2 Latency

ISSN (Online): request pending

Volume-1 Issue-1 | Jan-Mar 2025 | PP. 22-27

SqueezeNet achieved the lowest inference time (88 ms) across all platforms, ideal for time-critical tasks. However, it slightly trailed in accuracy.

#### 6.3 Power and Memory

Coral Dev Board's synergy with MobileNetV2 and SqueezeNet resulted in optimal performance-per-watt ratios. Raspberry Pi had the highest memory overhead when using Tiny-YOLO.

#### 6.4 Model Size

SqueezeNet was the most compact, making it ideal for constrained storage environments.

#### **CONCLUSION**

This study comprehensively evaluated the deployment performance of four lightweight deep learning models—MobileNetV2, SqueezeNet, ShuffleNet, and Tiny-YOLO—on three leading embedded AI platforms. Through a combination of empirical benchmarking and simulation, it was found that:

- MobileNetV2 provides the best overall balance between accuracy and latency.
- SqueezeNet is optimal for ultra-low-power or storage-constrained devices.
- Tiny-YOLOv3 should be used only when advanced object detection is needed and hardware acceleration is available (e.g., Jetson Nano).
- ShuffleNet offers a middle-ground solution with consistent performance across tasks.

Future work may focus on applying advanced model compression techniques such as quantization-aware training and neural architecture search (NAS) for even more efficient embedded AI.

## REFERENCES

- Howard, A. G., et al. (2017). MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv preprint arXiv:1704.04861.
- Iandola, F. N., et al. (2016). SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size, arXiv preprint arXiv:1602.07360.
- Zhang, X., et al. (2018). ShuffleNet: An extremely efficient convolutional neural network for mobile devices. CVPR.
- Redmon, J., & Farhadi, A. (2018). YOLOv3: An Incremental Improvement. arXiv preprint arXiv:1804.02767.
- Rasheed, H., & Qayyum, A. (2022). Edge AI: Trends, Challenges and Future Directions. ACM Computing Surveys.
- Albahar, M. A. (2021). Benchmarking Deep Learning Inference on Edge Devices. Sensors, 21(8), 2901.
- Lane, N. D., et al. (2016). DeepX: A software accelerator for low-power deep learning inference on mobile devices. IPSN.
- Han, S., et al. (2015). Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. arXiv:1510.00149.
- Jetson Nano Developer Kit | NVIDIA Developer. (2023). <a href="https://developer.nvidia.com/embedded/jetson-nano">https://developer.nvidia.com/embedded/jetson-nano</a>
- Google Coral Dev Board Documentation. (2023). https://coral.ai/products/dev-board/
- Raspberry Pi Foundation. (2023). Raspberry Pi 4 Model B specifications. https://www.raspberrypi.org/products/raspberry-pi-4-model-b/
- Krizhevsky, A., & Hinton, G. (2009). Learning multiple layers of features from tiny images. Technical Report.
- Paszke, A., et al. (2019). PyTorch: An imperative style, high-performance deep learning library. NeurIPS.
- Abadi, M., et al. (2016). TensorFlow: Large-scale machine learning on heterogeneous distributed systems. OSDI.
- Tan, M., & Le, Q. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. ICML.
- Han, S., et al. (2016). EIE: Efficient inference engine on compressed deep neural network. ISCA.
- Liu, Z., et al. (2020). TinyNAS: A framework for fast neural architecture search. arXiv:2007.11623.
- Zhang, X., et al. (2020). Fast, Accurate and Lightweight Super-Resolution with Deep Laplacian Pyramid Networks. CVPR.
- McMahan, H. B., et al. (2017). Communication-Efficient Learning of Deep Networks from Decentralized Data. AISTATS.
- Gao, X., et al. (2021). Survey on model compression and acceleration for deep neural networks. Neurocomputing.

ISSN (Online): red	urnal of Advanced quest pending    Jan-Mar 2025    P	puter Science and	l Engineering (IJAR	(CSE)