

Hardware Adaptation of Convolutional Neural Networks and Applications of Image Processing Algorithms

Jizhe Wu*

School of Optics and Photonics,
Beijing Institute of Technology,
Beijing, 100081, China
Corresponding author:
1120211755@bit.edu.cn

Abstract:

Against the backdrop of artificial intelligence driving the deep integration of computer vision into healthcare, transportation, industrial fields and so on, the collaborative adaptation of hardware and software for Convolutional Neural Networks, called CNN, as core technology has become crucial for overcoming application performance bottlenecks. This study systematically examines the core algorithmic mechanisms and hardware implementation techniques of CNN, analyzes directions for algorithmic innovation and multi-industry deployment practices, compares performance metrics of mainstream hardware, dissects bottlenecks such as insufficient edge computing power and inadequate hardware-software adaptation, and forecasts future development trends. Research reveals that hardware evolves through a progression of “general-purpose computing→specialized acceleration→adaptive reconfiguration.” By 2024, ASIC has captured a 42% market share, becoming the mainstream solution for data centers, with domestic chip performance approaching international standards. Algorithms are advancing toward lightweight and multimodal capabilities. China has achieved significant results in customized solutions for scenarios such as medical imaging and industrial quality inspection, while also defining tailored hardware-software adaptation strategies for different application contexts. This research provides a theoretical foundation for overcoming CNN application bottlenecks, aiding in unlocking the value of the computer vision industry and enhancing the efficiency of technology implementation in related fields.

Keywords: Convolutional Neural Networks, Hardware Adaptation, Image Processing Algorithms, Machine Vision.

1. Introduction

Against the backdrop of rapid iteration in artificial intelligence technology, computer vision has deeply integrated into critical fields such as medical diagnosis, traffic management, industrial manufacturing, and consumer electronics, leveraging its capabilities in perceiving and processing complex environments. What's more, convolutional neural networks stand as the core algorithmic framework driving breakthroughs in computer vision performance.

From the perspective of industry application effectiveness, CNN has already become a core technological tool for solving complex image processing problems. In the medical sector, CNN-based medical imaging auxiliary diagnosis systems enable early screening of lung cancer through CT scans. Clinical data from hospitals indicates that such systems can increase physician image review efficiency by 40% while keeping the missed diagnosis rate below 18%. In the intelligent transportation sector, CNN visual solutions accounted for 68% of the global intelligent transportation market in 2024, widely applied in scenarios such as traffic flow statistics and violation detection. In the industrial quality inspection field, after BMW Group implemented a CNN defect detection system, product defect detection rates rose from 85% to 99.2%, reducing product defect rates by 32%. In the consumer electronics sector, global shipments of smartphones equipped with CNN image processing algorithms reached 1.25 billion units in 2024, accounting for 83% of total annual smartphone shipments. This technology underpins core features such as portrait mode and night photography.

In terms of hardware implementation, the diversification of CNN application scenarios is driving hardware technology toward multi-path evolution. The global market size for CNN-related hardware reached \$31.2 billion in 2024, representing a 247% increase from 2020 and a compound annual growth rate of 35.2% [1]. From a technological evolution perspective, between 2016 and 2020, FPGA emerged as the preferred solution for industrial and edge computing scenarios due to their reconfigurable nature. Xilinx's Kintex series saw its application share rise from 12% to 28%. Following 2021, ASIC chips rapidly gained prominence due to their superior energy efficiency ratio, capturing 42% of the market by 2024. For instance, Google's TPU v5e delivered 22 times the computational power per watt compared to a GPU during the same period [1] [2].

At the algorithm research level, image processing technologies based on CNN have witnessed explosive growth. Between 2020 and 2024, the number of academic papers

published on CNN-based image processing surged from 12,000 to 38,000, achieving an average annual growth rate of 34.9% [3]. Current algorithm optimization primarily focuses on three directions. The first is a lightweight design, such as MobileNet-v3, which has only 5%-10% of the parameters of traditional networks. The second is precision quantization, such as INT8 quantization, which can reduce hardware resource consumption by 75%. The third is knowledge distillation, such as the YOLOv8-tiny version, which reduces parameters by 75% while only decreasing accuracy by 0.4%. In 2024, over 80% of CNN algorithms in industrial settings have adopted quantization schemes with INT8 precision or lower.

Under these circumstances, this study will integrate key data such as the market size of convolutional neural network hardware and the proportion of Algorithm accuracy optimization scheme in industrial scenarios. This data will be used to analyze the suitable application scenarios and implementation conditions for different hardware technologies and algorithm strategies. Its significance lies in leveraging the synergistic adaptation of hardware and algorithms to overcome performance bottlenecks in practical CNN applications. This provides a theoretical foundation for deepening their deployment in critical sectors such as healthcare, transportation, manufacturing, and consumer electronics. Ultimately, it will further unlock the industrial application value of computer vision technology, enhancing the efficiency of technological implementation and application outcomes across relevant fields.

2. Core Principles and Hardware Implementation Techniques of CNN Algorithms

2.1 CNN's Critical Algorithm Mechanism

The core strengths of CNN stem from three key characteristics: local receptive fields, weight sharing, and pooling operations. The functions and working principles of its three core modules are as follows.

Serving as the core feature extraction units, convolution layers perform sliding computations on input images using 3×3 convolutional kernels of uniform size to extract local features. Through the mechanism of weight sharing, the number of network parameters can be significantly reduced. Taking a $224 \times 224 \times 3$ input image as an example, when the number of output channels is 64, a single convolution operation requires approximately 230 million multiply-accumulate operations, accounting for over 80% of the total computational load in a CNN.

Pooling layers are primarily used for feature map dimen-

sion reduction and redundant information removal. The current mainstream approach employs 2×2 max pooling operations, which can reduce feature map dimensions by 75%. This preserves key feature information while alleviating computational pressure on subsequent modules and simultaneously mitigates overfitting. The fully connected layer performs matrix operations to map feature vectors to classification results, delivering the

final decision output. Taking the VGG16 network as an example, its three fully connected layers contain approximately 50 million parameters, accounting for about 15% of the network’s total computational load.

A comparison of the performance and key innovations of classic CNN architectures across different eras is shown in Table 1.

Table 1. Comparison of Performance and Innovations in Classic CNN Architectures Across Different Eras [4].

Network Model	Release Year	Number of Layers	ImageNet Top-1 Accuracy	Core Technical Innovation
LeNet-5	1998	8	-(No public data available)	First commercial application of CNN technology
AlexNet	2012	8	57.1%	Introduction of ReLU activation function and Dropout mechanism
ResNet-152	2015	152	78.3%	Proposal of residual connections to address vanishing gradient problem
EfficientNet-B7	2019	84	84.4%	Adoption of composite scaling strategy to balance accuracy and efficiency

2.2 Advances in Hardware Implementation Technologies at Home and Abroad

2.2.1 International Technological Frontiers and Industrial Applications

The NVIDIA H100 GPU platform, based on the Hopper architecture, delivers 4 PetaFLOPS of compute performance at FP16 precision. In ResNet-50 inference tasks, it achieves a throughput of 1,280 images per second, representing a 2x performance improvement over the previous generation. Its complementary CUDA ecosystem now supports over 90% of global deep learning frameworks, establishing itself as the mainstream choice for large-scale training and inference tasks in data centers[5].

Google’s TPU v5e, an ASIC-specific chip manufactured using a 4nm process, delivers 256 TOPS of INT8 precision computing power with an energy efficiency ratio of 30 TOPS/Watt—5 to 8 times that of contemporary GPU products. Currently, this chip handles 70% of image recognition inference tasks in Google’s data centers, reducing annual power consumption per server by 65% compared to traditional solutions[2].

The Xilinx FPGA-based Alveo U280 solution dynamically adjusts parallel computation levels according to convolutional kernel sizes through its reconfigurable logic units. For MobileNet-v2 inference tasks, this product achieves 4.2ms latency and 92.3% accuracy while consuming only 12.8W of power—a 96% reduction compared to GPU solutions. It is suitable for edge devices such as smart

cameras [6].

From the perspective of market share shifts, the period from 2015 to 2024 witnessed significant adjustments in the market share distribution across different hardware platforms. The market share of GPU declined from 75% to 38%, while that of ASIC increased from 5% to 42%. Meanwhile, FPGA maintained a stable share of approximately 20% [1]. This trend reflects the transformation of the CNN hardware market from general-purpose computing toward specialized acceleration.

2.2.2 Domestic Technological Breakthroughs and Solution Innovations

The Huawei Ascend 910 general-purpose AI chip adopts the Da Vinci architecture, integrating 32 AI Core computing units with an FP16 computing power of 320 TOPS. In ResNet-50 training tasks, an 8-card cluster achieves a linear acceleration ratio of up to 90%, completing ImageNet dataset training in just 3.5 days—performance approaching international top-tier levels. Its companion MindSpore framework supports over 100 mainstream CNN models, capturing an 18% share of China’s AI server market in 2024[7].

The Horizon Journey 5 automotive-grade chip integrates 16 BPU cores, delivering 128 TOPS of INT8 computing power with an energy efficiency ratio of 5 TOPS/W. This chip supports parallel operation of three model types: YOLOv8 (object detection), SegNet (semantic segmentation), and LaneNet (lane line recognition). To date, over 500,000 autonomous driving domain controllers based on

this chip have been deployed in vehicles.

The Institute of Microelectronics at Tsinghua University has proposed an innovative FPGA design solution called the “Adaptive Convolution Accelerator.” Based on the Xilinx Kintex UltraScale FPGA, it achieves adjustable parallelism from 1 to 16 levels. When processing images ranging from 224×224 to 1024×1024 resolution, resource utilization remains above 90%. delivering a 40% performance boost over fixed-architecture accelerators. The solution has been deployed in real-time 4K image detection systems for security applications.

3. Innovations and Practical Applications of CNN Image Processing Algorithms

3.1 Core Breakthroughs in International Algorithm Research

3.1.1 3.1.1 Image Classification Algorithms

The EfficientNetV2 model proposed by the Google team reduces computational redundancy through its Fused-MB-Conv architecture, achieving an 88.4% classification accuracy on the ImageNet dataset while training three times faster than the original EfficientNet. Its innovative composite scaling strategy—simultaneously adjusting network depth, width, and input resolution—results in a model parameter count just one-twelfth that of ResNet-152, establishing it as the benchmark model for product classification tasks in industrial quality inspection [8].

3.1.2 Object Detection Algorithm

Meta’s YOLOv8 employs a C2f backbone network with a PAN-FPN feature fusion architecture, achieving a mean average precision (mAP) of 53.9% on the COCO dataset and an inference speed of 300 FPS on GPU platforms. By incorporating deformable convolutions, this algorithm improves detection accuracy for small targets like distant pedestrians by 12% compared to YOLOv5, finding widespread application in intelligent traffic monitoring systems [9].

3.1.3 Image Segmentation Algorithms

Mask2Former enhances global feature capture capabilities by introducing a Transformer architecture to the Mask R-CNN framework, achieving a segmentation AP of 57.8% on the COCO dataset—a 20.7% improvement over the original model. This algorithm is currently applied in medical imaging for precise tumor region delineation, providing technical support for surgical planning [10].

3.2 Domestic Algorithm Innovation and Industry Implementation

3.2.1 Achievements in Lightweight Algorithm Design

Huawei’s Noah’s Ark Lab proposed the “Groupwise Residual with Groupwise Inversion” architecture, which maintains the performance of MobileNet-v3 (75.2% ImageNet accuracy) while achieving a 50% increase in inference speed and a 20% reduction in parameters. This algorithm has been integrated into the Kirin 9000s chip, enabling real-time portrait blurring and night scene enhancement on mobile devices [11].

The PP-YOLOE model developed by Baidu’s PaddlePaddle team achieves a mean average precision (mAP) of 51.4% on the COCO dataset, with GPU-based inference running at 149 FPS. Through its “dynamic label assignment” and “scalable head” designs, the model enables real-time detection at 25 FPS on embedded devices like Jetson AGX Orin. It has been deployed in Meituan’s obstacle recognition system for autonomous delivery vehicles [12].

3.2.2 Industry-Specific Customized Solutions

In the field of medical imaging diagnostics, United Imaging Healthcare collaborated with the Institute of Automation, Chinese Academy of Sciences to develop a “Multi-Disease Detection System for Chest CT.” Based on an enhanced U-Net architecture and integrated multi-scale feature fusion modules, the system achieves a comprehensive detection rate of 96.3% for three conditions—lung nodules, pneumonia, and pneumothorax—representing a 9-percentage-point improvement over traditional algorithms. The system has obtained NMPA certification and is now in clinical use at over 300 hospitals nationwide [13].

In the field of industrial defect detection, iFlytek has developed a specialized CNN model customized for lithium battery electrode defect detection. This model employs a multi-branch input structure combining visible light and infrared images, integrated with an attention mechanism. It achieves a detection accuracy of 99.2% for defects such as pinholes and scratches, while maintaining a false positive rate below 0.3%. In production lines at companies like CATL and BYD, this system has replaced manual inspection, boosting efficiency by 20 times [14].

In the field of remote sensing image processing, China Star Map has developed the “High-Resolution Satellite Image Segmentation System” based on the SegNeXt algorithm. This system achieves an Intersection over Union (IoU) of 89.7% for segmenting three types of urban features—buildings, roads, and water bodies—with a processing speed of 500 km² per hour. The system provides

technical support for the Third National Land Survey project and has been successfully applied to national land census and disaster monitoring scenarios.

4. Performance Comparison and Tech-

Table 2. Performance Metrics Comparison of Mainstream Hardware Platforms

Hardware Type	Representative Product	INT8 Compute Power	Energy Efficiency (TOPS/W)	ResNet-50 Inference Latency	Core Application Scenarios
GPU	NVIDIA H100	4096 TOPS	6.0	1.2 ms	Data Center Training and Inference
ASIC	Huawei Ascend 910	2560 TOPS	8.5	1.8 ms	Cloud-Based Inference Computing
ASIC	Horizon Journey 5	128 TOPS	5.0	8.5 ms	In-Vehicle Edge Computing
FPGA	Xilinx Alveo U280	200 TOPS	2.7	4.2 ms	Customized Edge Scenarios

Analysis of Table 2 reveals that GPU delivers optimal computational performance, ASIC offers significant advantages in energy efficiency, while FPGA remains irreplaceable in customized scenarios due to their flexibility. Enterprises should select hardware based on specific application requirements [1].

4.2 Core Bottlenecks in the Field of Technology

The first is the conflict between edge computing power and model requirements. Edge devices (such as smartphones and IoT terminals) have limited hardware resources (most terminals have less than 4GB of memory), making it difficult to support complex CNN models. For example, when running YOLOv8 on a mobile phone, inference latency typically exceeds 200ms, failing to meet real-time interaction demands. The second is insufficient hardware-software coordination. Advanced algorithms like dynamic convolutions and attention mechanisms demand high flexibility from hardware architectures. However, most current chips employ fixed computational unit designs, resulting in over 30% performance loss during algorithm deployment. Additionally, weak small-sample generalization capabilities persist. Labeled data scarcity in fields like healthcare and industry (e.g., rare disease image samples often number fewer than 100) limits CNN algorithms' generalization capacity, causing accuracy to drop by 15%-20% when applied across different scenarios [3].

nical Bottlenecks

4.1 Performance Metrics Comparison of Mainstream Hardware Platforms

5. Future Development Trends

5.1 Hardware Technology Evolution Direction

The first is the advancement of compute-in-memory architectures. For instance, CNN accelerators based on RRAM (Resistive Random-Access Memory) can significantly reduce data movement energy consumption. By 2027, the energy efficiency of such architectures is projected to exceed 100 TOPS/W (currently capped at 30 TOPS/W), delivering high computational power to edge devices [15]. By deeply integrating computational and storage units, this architecture addresses the "data movement bottleneck" inherent in traditional von Neumann architectures, making it suitable for scenarios like medical wearables and IoT terminals.

The second is the design of adaptive hardware. Dynamically reconfigurable FPGA/ASIC architectures will become mainstream, enabling hardware-based adaptation to diverse CNN models and minimizing performance degradation during algorithm deployment [16]. For instance, MIT's "Dynamically Reconfigurable CNN Chip" automatically adjusts the number of computational units based on the convolutional kernel size and channel count of the input model. When adapting to different networks like ResNet and YOLO, its performance loss remains below 5%, offering three times the flexibility of fixed-architecture chips.

Additionally, enhanced edge computing power will see automotive and medical edge chips surpass 500 TOPS, enabling local execution of complex models. For instance,

Qualcomm's planned 2026 automotive AI chip integrates eight dedicated CNN acceleration cores, simultaneously processing lidar point clouds, camera images, and millimeter-wave radar data to achieve Level 4 autonomous driving edge decision-making [5].

5.2 Algorithms and Application Innovation Direction

The first is multimodal fusion. Integrating images with text and sensor data enhances CNN robustness in complex scenarios. For instance, in autonomous driving, multimodal CNN models combining lidar and visual data can improve detection accuracy by 25% under adverse weather conditions [2]. Google DeepMind's "Vision-Language-CNN" model uses cross-modal attention mechanisms to correlate image features with textual descriptions. In industrial part defect diagnosis, it simultaneously identifies defect types and generates natural language reports, achieving a 40% increase in diagnostic efficiency compared to single-image algorithms.

The second is the application of federated learning. Cross-institutional joint training of CNN addresses medical data privacy concerns. Research papers on this topic are projected to increase by 200% between 2024 and 2028, accelerating the implementation of multi-center medical imaging diagnosis models [17]. For instance, Microsoft's "Lung Cancer Imaging Federated Learning Project," conducted in collaboration with 10 top-tier hospitals in China, achieved a 95.7% accuracy rate for jointly trained CNN models without sharing raw data—only 0.8% lower than centralized training results. This model is already deployed for lung cancer screening in primary care hospitals.

Finally, optimizations in few-shot learning enable Few-shot CNN algorithms to significantly reduce annotation costs in industrial quality inspection scenarios through transfer learning and data augmentation techniques. For instance, Alibaba's "Meta-CNN" model, pre-trained with a universal image feature extractor, achieved detection performance equivalent to traditional algorithms requiring 1,000 labeled samples using only 30 labeled samples in lithium battery electrode sheet defect detection. This reduced labeling costs by 97% and has been piloted on some production lines at CATL [18].

6. Conclusion

In the hardware domain, the technological trajectory follows an evolutionary logic of "general-purpose computing → specialized acceleration → adaptive reconfiguration." ASICs have become the mainstream choice for data cen-

ters due to their high energy efficiency ratio, while FPGAs maintain a 20% market share in edge customization scenarios owing to their flexibility. Domestically, Huawei's Ascend 910 and Horizon Robotics' Journey 5 have achieved technological breakthroughs in general-purpose AI chips and automotive-specific chips, respectively, with performance approaching international top-tier levels. However, cutting-edge technologies like compute-in-memory and dynamically reconfigurable hardware still require a 2-3 year catch-up period.

In the field of algorithms, lightweight approaches, multi-task fusion, and cross-modal collaboration have emerged as core development directions. Domestic customized algorithms have demonstrated outstanding performance in scenarios such as medical imaging and industrial quality inspection, addressing practical industry pain points. However, their small-sample generalization capabilities still lag behind international cutting-edge algorithms by 5%-8%.

In application domains, CNN technology has achieved large-scale deployment across multiple industries including healthcare, transportation, manufacturing, and consumer electronics. The synergistic adaptation of "hardware-algorithm-scenario" is key to maximizing value: - Data centers prioritize GPU for large-scale training - Cloud inference favors cost-effective ASIC - Edge devices prefer low-power FPGA/specialized ASIC - Complex tasks utilize multimodal algorithms - Low-cost scenarios employ lightweight algorithms

References

- [1] IDC. Global AI Hardware Market Forecast, 2024-2028 [R]. Framingham: International Data Corporation, 2024.
- [2] Google. TPU v5e: Efficient Inference for AI Workloads [R]. Mountain View: Google Inc., 2024.
- [3] IEEE. IEEE Xplore Digital Library: Convolutional Neural Network for Image Processing [EB/OL]. 2024.
- [4] Tan M, Le Q V. EfficientNetV2: Smaller Models and Faster Training [C]. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Nashville, TN, USA, Jun 20-25, 2021: 10096-10106.
- [5] NVIDIA. NVIDIA H100 Tensor Core GPU Datasheet [R]. Santa Clara: NVIDIA Corporation, 2024.
- [6] Xilinx. Alveo U280 Data Center Accelerator Card Product Brief [R]. San Jose: Xilinx Inc., 2023.
- [7] Huawei Technologies Co., Ltd. Ascend 910 AI Chip Technical White Paper [R]. Shenzhen: Huawei Technologies Co., Ltd., 2024.
- [8] Horizon Robotics Co., Ltd. Journey 5 Automotive AI Chip Product Manual [R]. Beijing: Horizon Robotics Co., Ltd., 2024.

- [9] Redmon J, Farhadi A. YOLOv8: Real-Time Object Detection and Segmentation [J]. arXiv preprint arXiv:2401.02786, 2024.
- [10] Cheng J, Schwing A G, Kirillov A. Mask2Former for Universal Image Segmentation [C]. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA, USA, Jun 19-24, 2022: 12906-12916.
- [11] Zhang X, Zhou X, Lin M, et al. MobileNeXt: Rethinking MobileNet Architecture for Efficient CNN [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 2022, 44(10): 6210-6225
- [12] Baidu PaddlePaddle Team. PP-YOLOE: An Efficient and Flexible Object Detector [R]. Beijing: Baidu Inc., 2023.
- [13] United Imaging Healthcare Group. Clinical Validation Report of Multi-Disease Intelligent Detection System for Chest CT [R]. Shanghai: United Imaging Healthcare Group, 2024.
- [14] iFlytek Co., Ltd. Technical Report on CNN Algorithm for Lithium Battery Electrode Defect Detection [R]. Hefei: iFlytek Co., Ltd., 2023.
- [15] IBM Research. RRAM-Based In-Memory Computing for CNN Acceleration [J]. Nature Electronics, 2024, 7(3): 210-220.
- [16] MIT Computer Science & Artificial Intelligence Laboratory. Dynamically Reconfigurable CNN Chip [R]. Cambridge: MIT, 2023.
- [17] Microsoft Research. Federated Learning for Medical Image Analysis [J]. Nature Medicine, 2023, 29(5): 1035-1043.
- [18] Alibaba Group. Meta-CNN: Few-Shot Learning for Industrial Defect Detection [J]. IEEE Transactions on Industrial Informatics, 2024, 20(2): 1890-1899.