



# Optimization Research of Deep Learning Algorithms in Real-time Image Processing

Genjuan Ma, Yan Li

Communication University of China, Nanjing, Nanjing 210000, Jiangsu, China

DOI: 10.32629/jher.v5i5.3065

---

**Abstract:** This paper studies the optimization method of deep learning algorithm in real-time image processing, focusing on the optimization strategies of model compression, hardware acceleration and data processing. Through these technologies, the reasoning speed of the image processing system is significantly improved, while the high accuracy of the model is maintained. The research shows that these optimization strategies have wide application potential in resource-constrained environments, providing efficient solutions for practical applications.

**Keywords:** real-time image processing; deep learning optimization

---

## 1. Introduction

With the rapid development of computer vision technology, image processing has been widely used in many fields, including intelligent security, autonomous driving, augmented reality, medical image analysis, etc. However, the demand for real-time image processing is constantly increasing, which requires the algorithm to not only to have the ability of high precision image analysis, but also to meet the requirements of real-time response under limited computing resources. Traditional image processing algorithms are often difficult to balance between accuracy and speed, and the rise of deep learning has brought new hope for solving this problem. Through its powerful feature extraction capability, the deep neural network has already achieved remarkable results in image classification, object detection, image segmentation and other tasks, but its computational complexity also brings new challenges for real-time processing. Therefore, how to optimize the deep learning model while maintaining the high efficiency of the algorithm has become a hot issue in the current image processing research.

## 2. The challenge of real-time image processing

### 2.1 Computational complexity and time delay

Deep learning algorithms, especially convolutional neural networks, rely on a large number of convolution operations and fully connected layer computations, which leads to a significant increase in their computational complexity when processing high-resolution images. Real-time image processing requires the system to complete the whole process from image acquisition to decision output in milliseconds, but standard deep learning models, especially those deep networks with large numbers of parameters and levels, usually require long inference time. The time delay caused by this calculation burden directly affects the real-time performance of the system, and cannot meet the scenarios such as autonomous driving and video surveillance, which require immediate response.

### 2.2 Hardware limitations and resource bottlenecks

Deep learning models often rely on strong hardware support, such as graphics processing units, tensor processing units, etc. While high-performance hardware can significantly improve the execution efficiency of algorithms, computing resources are often very limited in practice, especially on edge devices. These devices often do not have the computing power of high-end hardware, and are limited by power consumption and storage space, making it difficult to support the real-time operation of large-scale deep learning models. Therefore, the effective deployment of deep learning algorithms in a resource-constrained environment becomes a major challenge, requiring a balance between model design and hardware utilization.

### 2.3 Data size and bandwidth limitations

Real-time image processing typically involves processing large-scale, high-resolution image or video data. Especially in scenarios such as video surveillance and intelligent driving, the system needs to continuously process real-time video streams, which puts great pressure on data transmission and processing. With the improvement of image resolution and

data frame rate, network bandwidth and storage devices also need to be able to cope with massive data streams. In addition, the real-time nature of data requires the equipment to quickly obtain, process and feedback the results, and any delay in data transmission will lead to the degradation of the overall performance of the system. Optimizing data transmission and processing links to reduce delay and bandwidth requirements is one of the important links to realize real-time processing.

## **2.4 Accuracy and speed balance in real-time processing**

In real-time image processing, there is often a mutual restriction relationship between accuracy and speed. The high accuracy of deep learning models usually relies on deeper network structures and more complex computational processes, but this will inevitably increase inference time and lead to a decrease in processing speed. In some applications, speed may be more important than accuracy, such as drone visual navigation or augmented reality scenes in video games, while in others, such as medical image analysis or autonomous driving, accuracy is critical, and any subtle errors can have serious consequences.

## **2.5 Complexity of multitasking parallel processing**

In many real-time image processing applications, the system not only needs to complete a single task, but also needs to handle multiple tasks at the same time. For example, in autonomous driving scenarios, vehicles need to identify pedestrians, vehicles, traffic signs, etc. on the road in real time, and perform path planning and environmental awareness at the same time. This kind of multi-task parallel processing increases the complexity of the system, and puts forward higher requirements for the design of computing resources and algorithms.

## **2.6 Influence of environmental changes on algorithm robustness**

Real-time image processing systems usually run in complex and changeable environments. For example, self-driving vehicles need to process image data under different weather, lighting conditions and road conditions, while monitoring systems need to capture clear images under different time periods and different background noises. target image. The uncertainty of this environment requires the algorithm to be highly robust, that is, it can still maintain high recognition accuracy in the face of noise, occlusion, illumination changes and other interferences. Although deep learning models perform well under ideal conditions, they often face the problems of model instability and decreased inference accuracy when dealing with changeable actual scenarios.

# **3. Application of deep learning algorithms in image processing**

## **3.1 Basic application of convolutional neural network in image processing**

Convolutional neural network is the most basic and widely used model structure of deep learning in the field of image processing. CNN effectively reduces the number of parameters and improves the learning ability of the network through local receptive fields, weight sharing and pooling operations. In tasks such as image classification and object recognition, CNN can extract different levels of features from images through multi-layer convolution operations, such as edges, corners, textures, etc., and finally form a global representation. Classic CNN models such as LeNet, AlexNet and VGG have achieved breakthrough results on large-scale data sets such as ImageNet, which has greatly promoted the development of image processing.

In real-time image processing , CNN is the foundation of many core tasks. In video surveillance systems, CNN is used for real-time face detection and recognition to help improve the efficiency and reliability of security systems. In the vision system of driverless cars, CNN is used to detect pedestrians, vehicles and other obstacles on the road in real time to ensure driving safety. As the depth and complexity of CNN increase, its computational cost also increases. How to further optimize the CNN structure to meet the high efficiency requirements in real-time applications is one of the important research directions at present.

## **3.2 Development and application of object detection algorithm**

Object detection is one of the key tasks in image processing, which requires the model to not only recognize the object in the image, but also accurately locate the bounding box of the object in the image. Before deep learning, traditional object detection methods such as Haar cascade and HOG + SVM worked well on small-scale data sets, but could not handle complex scenes and diverse objects. Object detection algorithms based on deep learning, such as YOLO and SSD, combine the feature extraction capabilities of convolutional neural networks and the end-to-end learning mode, significantly improving detection accuracy and speed.

Among them, algorithms such as YOLO and SSD are specially designed for real-time processing. YOLO simplifies

the detection task to a regression problem, and simultaneously predicts the object category and position in one forward propagation, which greatly improves the detection speed. SSD realizes efficient multi-target detection through multi-scale feature maps, which is suitable for image data with high resolution. In real-time video analysis, this kind of algorithm can be used in monitoring, intelligent transportation system and other scenarios to realize real-time detection and tracking of pedestrians and vehicles.

### **3.3 Application of image segmentation algorithm**

Image segmentation is an advanced task in image processing that requires each pixel in an image to be assigned to a corresponding category label. Different from object detection, image segmentation requires classification accurate to the pixel level, which puts forward higher requirements for the expressive ability and computational complexity of the model. Deep learning algorithms have also demonstrated powerful capabilities in the field of image segmentation, especially models such as Fully Convolutional Networks and Mask R-CNN, which have been widely used in medical image analysis, autonomous driving and other fields.

U-Net is a classic model in medical image segmentation. Through the symmetrical encoder-decoder structure, it can effectively capture the global and local information in the image and achieve high-precision pixel-level segmentation. Mask R-CNN adds a segmentation branch on the basis of Faster R-CNN, which can perform target detection and instance segmentation at the same time. It is very suitable for real-time processing tasks with complex scenes and various types of targets. In the field of autonomous driving, image segmentation algorithms help vehicles identify key information such as road boundaries and lane lines to ensure that they can make accurate driving decisions in complex environments.

### **3.4 Image super-resolution and enhancement technology**

Image super-resolution techniques aim to reconstruct high-resolution versions from low-resolution images, which are important in many applications of image processing. Deep learning-based super-resolution algorithms, such as SRCNN and EDSR, achieve high-quality image reconstruction by learning the mapping relationship between high-resolution and low-resolution images.

In the fields of video surveillance and satellite image processing, image super-resolution technology can improve the clarity of image details and help the system to better identify the target object. At the same time, the image enhancement technology improves the visual effect of images through denoising, contrast enhancement, image correction and other operations. Such techniques are particularly important for the processing of low-quality image data, especially real-time video streaming data.

## **4. Optimization strategy of deep learning algorithm in real-time image processing**

### **4.1 Model compression technology**

The size and complexity of deep neural networks is one of the main reasons for the computational burden, especially in convolutional neural networks, which involve a large number of parameters and computational resources. Model compression technology aims to improve the reasoning speed by reducing the number of model parameters and calculation requirements, while maintaining the accuracy of the model as much as possible. The main compression methods include model pruning, quantification, and knowledge distillation.

#### **4.1.1 Model pruning and quantification**

Model pruning is a technique that reduces the size of a network by removing redundant or unimportant network connections. Pruning can be done after training, analyzing which weights have less impact on the final output, and subsequently removing these weights. The pruned network not only reduces the amount of parameters, but also significantly improves the inference speed, especially when running on edge devices or mobile devices. Common pruning methods include global pruning, structured pruning, and unstructured pruning.

Quantization reduces the storage and computational costs of the model by compressing the model parameters from high-precision floating-point numbers to low-precision representations. The quantized model can greatly reduce the consumption of hardware resources while maintaining high accuracy. This technology is particularly suitable for real-time processing scenarios in mobile devices and embedded systems.

#### **4.1.2 Knowledge distillation**

Knowledge distillation is a technique by transferring knowledge from a large "teacher" model to a smaller "student" model. The teacher model is first trained on a large amount of data to achieve high accuracy, and then its knowledge is transferred to the student model. The student model learns the output distribution of the teacher model during the training

process, so that it has fewer parameters and lower computational complexity while maintaining high accuracy. This method can effectively balance model accuracy and real-time reasoning speed, and is suitable for scenarios with limited hardware resources.

#### **4.1.3 Network structure search**

Network structure search is an automatic optimization method, which automatically designs efficient neural network structures through algorithms. NAS can explore different network topologies and find network models that can not only meet real-time requirements, but also maintain high accuracy. In recent years, based on reinforcement learning and evolutionary algorithms The NAS method has gradually matured and become one of the important tools for model optimization. In real-time image processing scenarios, NAS can help find the best network structure suitable for specific tasks and devices, significantly improving processing efficiency.

### **4.2 Optimization methods for accelerated reasoning**

The acceleration of inference process is a key requirement in real-time image processing. Aiming at inference efficiency, researchers have proposed a variety of optimization strategies, especially in the fields of hardware acceleration and parallel computing. These technologies greatly reduce inference delay.

#### **4.2.1 Hardware acceleration**

Hardware accelerator is an important means to improve the speed of deep learning inference, mainly including graphics processing unit, tensor processing unit and field programmable gate array. GPU can significantly accelerate the forward and backward propagation of convolutional neural networks through massively parallel computing capabilities, and has a wide range of applications in high-performance computing tasks. As hardware specially designed by Google for deep learning tasks, TPU has more advantages in matrix operation and tensor calculation, further improving the reasoning efficiency in real-time image processing.

An FPGA is a programmable hardware with a high degree of customizability. In real-time image processing tasks, FPGA can optimize hardware logic circuits according to specific application scenarios, thereby achieving low-latency, high-efficiency deep learning inference. Although FPGA is difficult to design and develop, its application prospects in embedded and edge computing devices are broad.

#### **4.2.2 Parallel computing and heterogeneous computing**

Parallel computing is a technique that accelerates overall processing by performing multiple computing tasks simultaneously. In deep learning reasoning, especially the convolution operation of convolutional neural network, it can be divided into multiple subtasks for parallel processing, thus greatly reducing the computation time. In addition to implementing parallel computing on a single device, inference tasks can also be performed in parallel on multiple devices through a distributed computing framework.

#### **4.2.3 Online learning and incremental update**

Online learning is a technique that dynamically adjusts model parameters in a real-time environment, allowing the model to be updated and optimized progressively as new data arrives. For real-time image processing systems, online learning allows the model to continuously adapt to new environments and new tasks, avoiding the need for frequent retraining. Incremental update technology allows the model to be optimized on a small scale without complete retraining, reducing system downtime and improving processing efficiency.

### **4.3 Optimization of real-time data processing**

#### **4.3.1 Acceleration strategies for data preprocessing and enhancement**

Data preprocessing is a necessary step in image processing system, including image scaling, normalization, denoising and other operations. In real-time scenarios, data preprocessing may become the bottleneck of system performance. Therefore, optimizing the data preprocessing process, such as using fast image transformation algorithms and parallelizing the data enhancement process, can effectively reduce the latency. For video stream processing, inter-sampling or resolution dynamic adjustment strategies may be used to reduce unnecessary computational burden.

#### **4.3.2 Improvement of real-time image acquisition and transmission**

Data acquisition and transmission in real-time image processing system are also very important. In this process, using efficient image compression and transmission protocols, optimizing network bandwidth allocation and other measures can effectively reduce the delay caused by image data transmission and ensure the timeliness of real-time processing.

## **5. Performance evaluation of deep learning optimization algorithms**

The trade-off between accuracy and efficiency is a core consideration when evaluating the performance of deep learning

optimization algorithms. Although the optimization algorithm aims to improve the inference speed and resource utilization of the model, it must ensure that the prediction accuracy of the model is not significantly reduced. In specific applications, the optimization effect is usually measured by comparing the inference time, computing resource occupation and prediction accuracy before and after optimization. In addition, detailed experimental comparison is needed for different optimization strategies. On resource-constrained devices, whether the optimized algorithm can run smoothly is also a part that cannot be ignored in the evaluation. For real-time application scenarios, the evaluation also needs to pay attention to the delay performance of the algorithm, that is, the total time from data input to result output.

## 6. Conclusion

This study explores the optimization strategy of deep learning algorithm in real-time image processing. Through model compression, hardware acceleration, and data processing optimization, inference efficiency is significantly improved while maintaining high accuracy. These optimization methods provide an effective way to solve the computational bottleneck in real-time image processing, especially in resource-limited devices. Future research can further combine more efficient algorithms and hardware innovations to continuously improve the performance of real-time processing and provide support for more practical scenarios.

## References

---

- [1] Xing Zhihua. Research on image real-time edge detection system based on FPGA [D]. Donghua University, 2023.
- [2] Wang Hai, Zhang Cheng, Zhang Suyi, et al. Research on an ingredient ratio detection method based on real-time image technology [J]. *China Food Industry*, 2023, (11): 84-88 +93.
- [3] Zhou Feng. Research on the Application of Image Analysis and Processing Technology Based on Deep Learning in Detecting Cracks in Building Exterior Walls [J]. *Journal of Hebei Software Vocational and Technical College*, 2024, Vol. 26(3): 31-33.
- [4] Zhang Kezhi, Wei Guoqiang, Feng Ze, et al. Research on the Application of Deep Learning Technology in Intelligent Image Processing [J]. *Modern Information Technology*, 2021, Vol. 5(10): 15-19, 26.
- [5] Shi Zhibin, Luo Wangchun, Mo Bingbing, et al. Research on Real-time Image Processing and Recognition Algorithms for Power Inspection Based on Deep Learning [J]. *Electronic Design Engineering*, 2022, Vol. 30(23): 189-193.

## Author Bio

Genjuan Ma: born in 1981, female, Yancheng, Jiangsu Province, master, senior engineer, research direction: Software Engineering, Big Data and artificial intelligence.

Yan Li: born in 1986, female, native of Yangzhou, Jiangsu province, Han nationality, undergraduate, research interests: senior engineer, big data, artificial intelligence, cloud computing, Java.