# LoRA-C: Parameter-Efficient Fine-Tuning of Robust CNN for IoT Devices

Chuntao Ding, Xu Cao, Jianhang Xie, Linlin Fan, Shangguang Wang, Zhichao Lu

**Abstract**—Efficient fine-tuning of pre-trained convolutional neural network (CNN) models using local data is essential for providing high-quality services to users using ubiquitous and resource-limited Internet of Things (IoT) devices. Low-Rank Adaptation (LoRA) fine-tuning has attracted widespread attention from industry and academia because it is simple, efficient, and does not incur any additional reasoning burden. However, most of the existing advanced methods use LoRA to fine-tune Transformer, and there are few studies on using LoRA to fine-tune CNN. The CNN model is widely deployed on IoT devices for application due to its advantages in comprehensive resource occupancy and performance. Moreover, IoT devices are widely deployed outdoors and usually process data affected by the environment (such as fog, snow, rain, etc.). The goal of this paper is to use LoRA technology to efficiently improve the robustness of the CNN model. To this end, this paper first proposes a strong, robust CNN fine-tuning method for IoT devices, LoRA-C, which performs low-rank decomposition in convolutional layers rather than kernel units to reduce the number of fine-tuning parameters. Then, this paper analyzes two different rank settings in detail and observes that the best performance is usually achieved when $\alpha/r$ is a constant in either standard data or corrupted data. This discovery provides experience for the widespread application of LoRA-C. Finally, this paper conducts many experiments based on pre-trained models. Experimental results on CIFAR-10, CIFAR-100, CIFAR-10-C, and Icons50 datasets show that the proposed LoRA-Cs outperforms standard ResNets. Specifically, on the CIFAR-10-C dataset, the accuracy of LoRA-C-ResNet-101 achieves 83.44% accuracy, surpassing the standard ResNet-101 result by +9.5%. On the Icons-50 dataset, the accuracy of LoRA-C-ResNet-34 achieves 96.9% accuracy, surpassing the standard ResNet-34 result by +8.48%. In addition, compared with full parameter fine-tuning, LoRA-C can reduce the amount of updated parameters by more than 99%.

**Index Terms**—Internet of things, LoRA, Cloud-Device Collaboration, CNN.

✦

## 1 INTRODUCTION

TRAINING large models is time-consuming and expensive. For example, training Meta's LLaMA-65B model requires 2,048 NVIDIA A100 GPUs, takes about 21 days, and costs more than $ 2.4 million. Larger models, such as OpenAI's GPT-3, contain 175 billion parameters and cost more than $ 4 million to train[1]. Moreover, these large models [1], [2] usually have broad general capabilities, but their performance on specific tasks is insufficient to meet practical needs. Using task-specific data to fine-tune these large models to meet the needs of specific tasks has attracted widespread attention and spawned many effective fine-tuning methods [3]–[12]. Among these fine-tuning methods, Low-Rank Adaptation (LoRA) [9], [13]–[18] is widely used

- *Chuntao Ding is with School of Artificial Intelligence, Beijing Normal University, Beijing, China. E-mail: ctding@bnu.edu.cn*
- *Xu Cao and Jianhang Xie are with Key Laboratory of Big Data & Artificial Intelligence in Transportation, Ministry of Education, with School of Computer Science and Technology, Beijing Jiaotong University, Beijing 100044, China. E-mail: {caoxuu; xiejianhang}@bjtu.edu.cn.*
- *Linlin Fan is with the School of Computer and Information Engineering, Henan Normal University, Henan 453007, China. E-mail: huazheng1106@gmail.com.*
- *Shangguang Wang is with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing, China. E-mail: sgwang@bupt.edu.cn.*
- *Zhichao Lu is with the Department of Computer Science, City University of Hong Kong, HKSAR. E-mail: luzhichaocn@gmail.com.*

1. https://www.cnbc.com/2023/03/13/chatgpt-and-generative-ai-are-booming-but-at-a-very-expensive-price.html



Fig. 1: Overview of system architecture.

in the Transformer architecture because of its simplicity, effectiveness, and no additional reasoning burden.

On the other hand, ABI Research reports that the global number of smart cameras using artificial intelligence (AI) chips will reach 350 million in 2025 [2]. In addition, Internet of Things (IoT) devices are widely deployed outdoors and usually process data affected by the environment (such as fog, snow, rain, etc.). Compared with Transformer [19]–[21], Convolutional Neural Networks (CNN) [22]–[25] have more advantages in balancing resource consumption and performance in visual tasks and has been widely deployed on resource-constrained IoT devices. Hence, it has become mainstream to use local data to fine-tune pre-trained CNN models transmitted from the cloud to achieve high performance on ubiquitous smart cameras containing AI chips. As shown in Fig. 1, the pre-trained model is fine-tuned using local data on IoT devices to improve the model's perfor-

2. https://www.abiresearch.com/press/global-installed-base-smart-city-cameras-ai-chipset-reach-over-350-million-2025/

mance in processing images affected by the environment and provide high-quality services.

However, there are few studies on fine-tuning CNN using LoRA. The goal of this paper is to use LoRA to fine-tune the CNN model on IoT devices to improve its robustness and provide high-quality services.

To this end, this paper first proposes the LoRA-C, a layer-wise LoRA for CNN fine-tuning. Inspired by the application of LoRA in the Transformer structure, a natural approach is to use the low-rank decomposition of each convolution kernel as the newly added weights. However, due to the kernel parameters sharing characteristics of CNN, low-rank decomposition based on convolution kernels will result in too many updated parameters, making it difficult to deploy on resource-constrained IoT devices. Therefore, instead of adding low-rank decomposition kernel-wise, this paper proposes LoRA-C, which adds low-rank decomposition in convolutional layer-wise to efficiently fine-tune CNN models. This paper also experimentally verifies that LoRA-C achieves significant performance improvements with only a few parameters updated.

Then, this paper discusses in detail the relationship between $\alpha$, which controls the weight of the newly added parameters, and the rank of $r$ of the matrix. Extensive experiments show that the proposed LoRA-C usually performs best when $\alpha/r$ is a constant. This finding provides experience for efficiently fine-tuning CNN models on IoT devices. This paper also studies two settings of matrix rank $r$ and gives rank $r$ settings suitable for different models through many experiments.

Finally, this paper conducts extensive experiments on the CIFAR-10, CIFAR-100, CIAFR-10-C, and Icons-50 datasets, and the experimental results show that our proposed LoRA-C achieves strong performance. For example, our proposed LoRA-C-ResNet-50 achieves 82.98% accuracy on the CIFAR-100 dataset, which surpasses the standard ResNet-50 result by +5.29%. On the CIFAR-10-C and Icons-50 datasets, our proposed LoRA-C-ResNet-34 achieves an average accuracy of 78.45% and 96.9%, respectively, which is +5.9% and +8.48% higher than the results of the standard ResNet-34. In addition, we also show that LoRA-C has the advantage in the number of updated parameters, making it more suitable for deployment on resource-constrained IoT devices.

In summary, our main contributions are as follows:

- To our knowledge, this is the first work that seeks to improve CNN robustness by incorporating LoRA. Specifically, according to the characteristics of the CNN model, this paper proposes LoRA-C, a convolutional layer-wise low-rank decomposition method for effective fine-tuning.
- This paper observes that the proposed LoRA-C generally performs best when $\alpha/r$ is a constant, which provides experience for efficiently fine-tuning CNN models on IoT devices.
- This paper compares the impact of rank settings on model robustness and analyzes its relationship with model capacity in detail.
- Extensive experiments on the CIFAR-10, CIFAR-100, CIFAR-10-C, and Icons-50 datasets demonstrate that our proposed LoRA-C achieves higher accuracy.

The remainder of the paper is organized as follows. Section 2 reviews the work on convolutional neural networks, cloud-device collaboration, and parameter-efficient fine-tuning technologies. Section 3 describes the proposed LoRA-C in detail. Section 4 presents our evaluation results, and Section 5 concludes the paper.

## 2 RELATED WORK

This section briefly introduces the Convolutional Neural Networks (CNN), cloud-device collaboration, and Parameter-Efficient Fine-Tuning (PEFT) techniques.

### 2.1 Convolutional Neural Networks

The great success of convolution operations in visual tasks has spawned many classic convolutional neural network (CNN) models in recent years, such as AlexNet [26], VGG [22], Inception and its variants [27]–[30], ResNet and its variants [25], [31], Xception [32], DenseNet [24], SENet [33], etc. Among them, ResNet [25] turns the CNN model into a deep CNN, introducing the residual module. SENet [33] further improves the performance of the CNN model by introducing the idea of attention in the channel dimension.

The lightweight CNN models include the MobileNets series [34]–[36], which separates spatial and channel dimensions; the ShuffleNet series [23], [37], which uses group convolution; and the GhostNet series [38], [39], which uses linear transformation to generate other feature maps based on a few feature maps. Recently, RepVGG [40] introduces the re-parameterization to speed up the inference of the CNN model. MonoCNN [41], [42] and SineFM [43] have been proposed to improve the robustness of CNN models. On the one hand, they reduce the number of learning parameters in the CNN model to alleviate the parameter update from being completely dependent on training data. On the other hand, they use nonlinear transformation rules to regularize the model to improve the robustness of the CNN model.

Different from the above methods, this study provides another idea, which aims to explore the use of LoRA to improve the robustness of the CNN model, so as to help CNN provide high-quality services on IoT devices.

### 2.2 Cloud-Device Collaboration

Due to the limited resources of IoT devices, cloud-device collaborative training and deployment of CNN models have become mainstream. For example, Kang *et al.* [44] propose to divide the CNN into a head running on the device and a tail running on the cloud and determine the split point based on the device load and the amount of cloud-to-device parameter transmission. Zhang *et al.* [45] train the CNN model through cloud-edge collaboration and prune the CNN in the cloud to minimize the number of model transmission parameters while retaining the model performance. Stefanos *et al.* [46] proposed a progressive inference method for collaborative device and cloud computing and used compression [47] to reduce the amount of parameter exchange between the device and the cloud. However, the
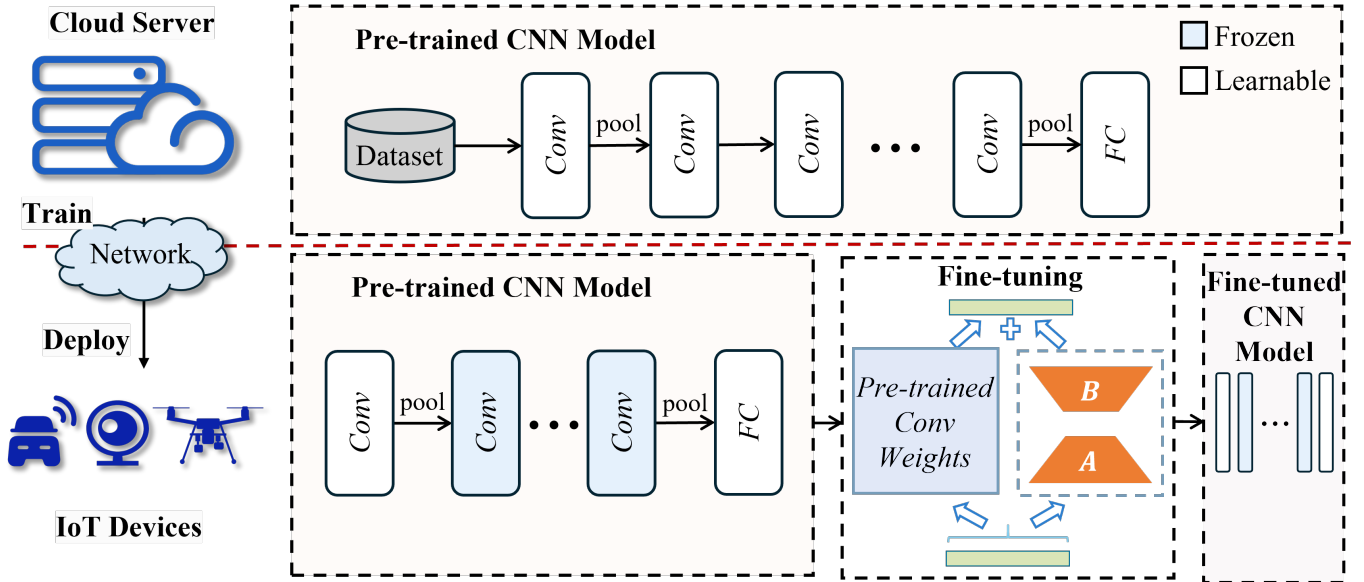
Fig. 2: Overview of the proposed framework.

above methods are highly dependent on network conditions. When the network condition is unstable or unavailable, the service quality will be reduced or even unavailable.

To decouple service quality from network conditions, many works proposed that the cloud server first trains the model and then sends the trained model to IoT devices [41]–[43], [48], [49]. For example, Ding *et al.* [41] propose to train a CNN model containing a small number of learning parameters, namely MonoCNN, on the cloud, and then send the trained MonoCNN model to IoT devices.

This paper studies the cloud training of CNN models first and then sends the trained CNN models to IoT devices. The difference is that IoT devices efficiently fine-tune the CNN models sent from the cloud based on their data rather than directly using these CNN models to provide services.

### 2.3 Parameter-Efficient Fine-Tuning Techniques

In recent years, many excellent parameter-efficient fine-tuning techniques [3], [4], [50] have been proposed, aiming to improve the performance of pre-trained models on new tasks by minimizing the number of fine-tuning parameters and computational complexity, thereby reducing the cost of training models on new tasks. Existing methods can be divided into three categories.

The first category is called the Adapter [4]–[7], which inserts a module with a small number of parameters into the pre-trained model for each downstream task, and the inserted module with a small number of parameters is called the Adapter module. In this category, only the Adapter module parameters are updated during training for downstream tasks, thereby efficiently migrating the capabilities of the powerful foundation models to many downstream tasks while ensuring the model's performance in downstream tasks. However, it increases the depth of the pre-trained model, thereby increasing the latency of model reasoning.

The second category is called the P-tuning [8], [51], [52], which constructs a task-related virtual token as a prefix

before inputting the token. It then only updates the parameters of the prefix part during training, while the other parameters in the Transformer are fixed. However, P-tuning methods are challenging to train and reduce the available sequence length of the model.

The third category is called LoRA and its variants [9], [14]–[18]. The full name of LoRA is low-rank adaptation, and is first proposed by Hu *et al.* [9]. Its core idea is to simulate the change of model parameters when the model adapts to new tasks through low-rank decomposition to realize indirect training of large models with an extremely small number of parameters. LoRA assumes that the pre-trained foundation models are over-parameterized and have a small "intrinsic dimension" [53], [54], that is, there is an extremely low dimensional. Fine-tuning the parameters of the extremely low dimension can have the same effect as fine-tuning in the full parameter space. LoRA does not change the depth of the pre-trained model, and can merge parameters added for new tasks during inference, thus generating no additional inference delay. Given the above advantages of LoRA, LoRA and its variants have been extensively studied in many applications. For example, Zhang *et al.* [17] combine LoRA with pruning and propose LoRAPrune. Dettmers *et al.* [14] and Li *et al.* [16] combine LoRA with quantization and propose QLoRA and LoftQ, respectively. Liu *et al.* [18] propose the DoRA based on LoRA. It first decomposes the pre-trained weights into their amplitude and direction components and then fine-tunes the two, which enhances the learning ability and training stability of LoRA.

Besides, the research on LoRA fine-tuning with CNN model is rare currently [55]–[58]. The former two works [55], [56] attempted to combine LoRA with convolution but did not analyze it further and deeper. Zhong *et al.* [57] combined the lightweight convolution into the mixture of expert network to improve LoRA, but did not fine-tune the convolution itself. Yeh *et al.* [58] proposed the LoCon which has two convolutions for dimension-down (with $k \times k$

convolution) and dimension-up (with $1 \times 1$ convolution) to simulate the matrices decomposition to reduce the number of fine-tuning parameters.

Our research falls within the third category. We aim to use LoRA to improve the robustness of CNN models, a direction with high practical application value but little research. This paper will study in depth how to apply LoRA to CNN, so that it can effectively improve the performance of the CNN model while only updating a very small number of parameters.

## 3 DESIGN OF THE PROPOSED APPROACH

### 3.1 Overview

Fig. 2 illustrates the overview of the proposed approach. IoT devices are usually deployed outdoors and collect a lot of data affected by the environment (such as fog, rain, and snow). We refer to the performance of the CNN model when processing this type of data as the robustness of the model. The standard CNN model is low-robust when processing this type of data. We aim to efficiently fine-tune the pre-trained model on resource-constrained IoT devices using local data to improve model robustness and handle the environmental-affected data collected by IoT devices.

To this end, we propose parameter-efficient fine-tuning of robust CNN method, LoRA-C. We first use the cloud to train the complex CNN model, known as the pre-trained model, and then send the pre-trained model to IoT devices. On the IoT device, we freeze all convolutions of the pre-trained model, excluding its first and last layers, and add LoRA-C branches to each frozen convolution layer. Then, the IoT device uses a small amount of local data to fine-tune this pre-trained model based on the LoRA-C. When providing CNN-powered IoT smart services, the parameters of the LoRA-C branch can be fused into the pre-trained convolutional weights, which will not introduce additional inference latency.

Next, we will conduct a fine-grained analysis from motivation and design of LoRA-C.

### 3.2 Motivation

This section will illustrate our motivation by analyzing the two most relevant techniques, namely low-rank adaptation and the convolutional layer.

**Low-Rank Adaptation** [9]: Low-rank adaptation (LoRA) uses low-rank decomposition to simulate the change in model parameters when the model adapts to new tasks, thereby achieving indirect training of large models by updating very few parameters. Formally, for a pre-trained weight matrix $W_0 \in \mathbb{R}^{d \times k}$, the newly added weight matrix $\Delta W = BA \in \mathbb{R}^{d \times k}$, matrix $A \in \mathbb{R}^{r \times k}$ and matrix $B \in \mathbb{R}^{d \times r}$, and the rank $r \ll min(d, k)$. When fine-tuning the model, only the parameters of $\Delta W$ are updated, while the parameters of $W_0$ are frozen. $W_0$ and $\Delta W$ are multiplied with the same input, and their respective output vectors are summed coordinate-wise. Given input $x$, the forward pass yields:

$$h = W_0 x + \alpha \Delta W x \ , \tag{1}$$

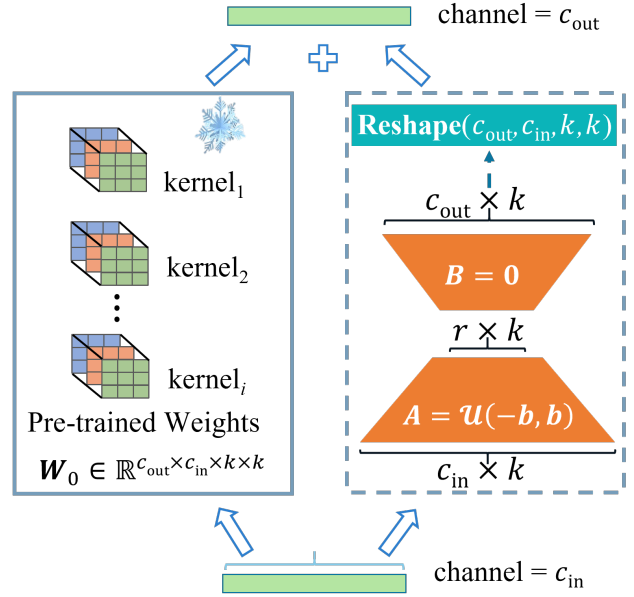where $\alpha$ is a constant. In the $\Delta W$, the settings of LoRA initialize $A$ with the random Gaussian distribution and



Fig. 3: The proposed LoRA-C.

set $B$ to all zeros. Since $r \ll min(d, k)$, the number of parameters of $\Delta W$ is much smaller than that of $W_0$. In addition, during inference, $\Delta W$ can be completely added to $W_0$ without changing the model architecture and without incurring any additional inference overhead.

**Convolutional Layer** [26], [59], [60]: The convolution uses sliding kernels to extract features. Let $x \in \mathbb{R}^{c_{\text{in}} \times w \times h}$ denote the input feature maps, where $w$, $h$, and $c_{\text{in}}$ are the input width, input height, and input channels, respectively; and $W \in \mathbb{R}^{c_{\text{out}} \times c_{\text{in}} \times k \times k}$ denote the convolution weights, where a kernel with size of $k \times k$. Then the generated output feature maps is $y \in \mathbb{R}^{c_{\text{out}} \times w' \times h'}$, where $w'$, $h'$ and $c_{\text{out}}$ are the output width, output height, and output channels, respectively. So the forward of the convolutional layer can be described by:

$$y = W \otimes x \ , \tag{2}$$

where the $\otimes$ is a convolution operation. For a kernel $W^{(n)} \in \mathbb{R}^{c_{\text{in}} \times k \times k}, (0 < n < c_{\text{out}})$, the output feature map $y^{(n)}$ in output channel $n$ is:

$$y^{(n)} = W^{(n)} \otimes x \ , \tag{3}$$

where a value of feature map $y^{(n)}$ in the position $(i, j)$ can be calculated by:

$$y_{i,j}^{(n)} = \sum_{m=1}^{c_{\text{in}}} \sum_{u=1}^{k} \sum_{v=1}^{k} W_{m,u,v}^{(n)} \cdot x_{m,i-u+1,j-v+1} \ . \tag{4}$$

Most existing studies use LoRA to fine-tune models based on the Transformer architecture, and there are very few LoRA studies targeting CNN. Given the wide application of CNN on IoT devices, this paper studies using LoRA to fine-tune the CNN model. When fine-tuning the Transformer-based architecture, the $BA$ in LoRA acts as the $Q$, $K$, and $V$ matrices. *So how can LoRA be efficiently combined with CNN?*

## 3.3 Design of LoRA-C

This section first presents LoRA-C and then explains the reasons for this design.

**LoRA-C**: Fig. 3 illustrates the proposed LoRA-C, which is a convolutional layer-wise low-rank decomposition method for CNN fine-tuning. For a convolutional layer with the input channel $c_{\text{in}}$ and output channel $c_{\text{out}}$, the pre-trained convolutional weights $\boldsymbol{W}_0 \in \mathbb{R}^{c_{\text{out}} \times c_{\text{in}} \times k \times k}$ are froze by removing the gradients, which are not updated during fine-tuning.

We introduce a LoRA-C branch with layer-wise decomposed matrices for updating the pre-trained weights $\boldsymbol{W}_0$, that is, decomposing the updating increment of weights $\Delta \boldsymbol{W}$ into two matrices, $\boldsymbol{A}$ and $\boldsymbol{B}$, where $\boldsymbol{A} \in \mathbb{R}^{r \times c_{\text{in}} \times k}$ and $\boldsymbol{B} \in \mathbb{R}^{c_{\text{out}} \times k \times r}$.

As with LoRA, we also introduce a constant, $\alpha$, to scale the weight of the newly added parameters in LoRA-C, and the updating of fine-tuning can be described by Equation 5:

$$\boldsymbol{W} = \boldsymbol{W}_0 + \alpha \Delta \boldsymbol{W} = \boldsymbol{W}_0 + \alpha \boldsymbol{B} \boldsymbol{A} \ . \tag{5}$$

The increments of pre-trained weights $\boldsymbol{W}_0$ brought by fine-tuning is equivalent to $\boldsymbol{B} \boldsymbol{A}$, where $\boldsymbol{B} \boldsymbol{A} \in \mathbb{R}^{c_{\text{out}} \times k \times c_{\text{in}} \times k}$. The $\boldsymbol{B} \boldsymbol{A}$ matrix can be reshape into a tensor whose dimensions are equal to the pre-trained weights $\boldsymbol{W}_0$, i.e., $\mathbf{Reshape}\,(\boldsymbol{B}\boldsymbol{A}) \in \mathbb{R}^{c_{\text{out}} \times c_{\text{in}} \times k \times k}$.

For the initialization of $\boldsymbol{A}$ and $\boldsymbol{B}$ matrices, we follow the settings of LoRA, set the matrix $\boldsymbol{B}$ with all zeros $\mathcal{Z}$ and initialize matrix $\boldsymbol{A}$ with Kaiming uniform distribution $\mathcal{U}(-b, b)$, where the $b$ is the bounding of uniform distribution [61]. Let $\boldsymbol{T}$ denote a random initial tensor, the initialization of $\boldsymbol{A}$ and $\boldsymbol{B}$ is as follows:

$$\begin{cases} \boldsymbol{A} = \mathcal{U}(\boldsymbol{T}), \ \boldsymbol{T} \in \mathbb{R}^{r \times c_{\text{in}} \times k} \\ \boldsymbol{B} = \mathcal{Z}(\boldsymbol{T}), \ \boldsymbol{T} \in \mathbb{R}^{c_{\text{out}} \times k \times r} \end{cases} . \tag{6}$$

The forward of convolution operation in LoRA-C is:

$$y = (\boldsymbol{W}_0 + \alpha \boldsymbol{B} \boldsymbol{A}) \otimes x \ . \tag{7}$$

For the rank of $\Delta \boldsymbol{W}$ (or $\boldsymbol{B}\boldsymbol{A}$), we provide two settings of $r$ and $r * k$, while the impact of the two different ranks will be analyzed in Section 3.3.2.

### 3.3.1 Decomposed Matrices Granularity

To apply the mechanism of LoRA fine-tuning to CNN model, we made modifications in layer-wise granularity to adapt to the structure of convolutional layer. *So why choose a laye-wise decomposition matrices instead of the more fine-grained kernel-wise?*

We will analyze the choice of layer-wise decomposition matrices granularity in terms of parametric efficiency. For the full fine-tuning, the number of updated parameters for pre-trianed weights $\boldsymbol{W}_0$ is:

$$\mathcal{P}_{\text{full\_ft}} = c_{\text{out}} \cdot c_{\text{in}} \cdot k^2 \ . \tag{8}$$

We aim to reduce the number of updated parameters in the convolutional weights by utilizing the LoRA technique. *Kernel-Wise Decomposed Matrices*: Based on the previous experience of LoRA in attention, we first consider assigning the $\boldsymbol{B}\boldsymbol{A}$ matrix for each convolutional kernel to fine-tune, i.e., kernel-wise decomposed matrices in LoRA-C.

Let $\boldsymbol{W}^{(i)} \in \mathbb{R}^{c_{\text{in}} \times k \times k}$ denote the $i$-th convolutional kernel in output channel $i$, where $0 < i < c_{\text{out}}$. So, for each $\boldsymbol{W}^{(i)}$, it can introduce a pair of low-rank matrices $\boldsymbol{A}^{(i)}$ and $\boldsymbol{B}^{(i)}$ with rank $r$, where $\boldsymbol{A}^{(i)} \in \mathbb{R}^{r \times c_{\text{in}} \times k}$ and $\boldsymbol{B}^{(i)} \in \mathbb{R}^{k \times r}$. So the updating of one kernel can be given by Equation 9:

$$\boldsymbol{W}^{(i)} = \boldsymbol{W}_0^{(i)} + \alpha \Delta \boldsymbol{W}^{(i)} = \boldsymbol{W}_0^{(i)} + \alpha \boldsymbol{B}^{(i)} \boldsymbol{A}^{(i)} \ , \tag{9}$$

so the forward of convolution in kernel-wise is:

$$y^{(i)} = \left( \boldsymbol{W}_0^{(i)} + \alpha \boldsymbol{B}^{(i)} \boldsymbol{A}^{(i)} \right) \otimes x \ . \tag{10}$$

For kernel-wise decomposition, the number of parameters for $\boldsymbol{A}^{(i)}$ and $\boldsymbol{B}^{(i)}$ is $c_{\text{in}} \cdot r \cdot k + r \cdot k$, and there are $c_{\text{out}}$ kernels in a convolution layer. Thus, the number of updated parameters for kernel-wise decomposed matrices is:

$$\mathcal{P}_{\text{kernel\_wise}} = c_{\text{out}} \cdot (c_{\text{in}} \cdot r \cdot k + r \cdot k) \ . \tag{11}$$

*Layer-Wise Decomposed Matrices*: For layer-wise decomposition using a pair of decomposed matrices with rank $r$ to fine-tune the convolutional weights in convolutional layer, the corresponding number of updated parameters can be described in Equation 12:

$$\mathcal{P}_{\text{layer\_wise}} = (c_{\text{out}} + c_{\text{in}}) \cdot r \cdot k \ . \tag{12}$$

*Updated Parameters Reduction*: For the two different LoRA-C granularities of low-rank decomposition, kernel-wise and layer-wise, we analyze the number of updated parameters compared to full fine-tuning, respectively.

For an intuitive comparison, we choose the ratio of the amounts of updating parameter of two LoRA-C granularities to the full fine-tuning number of parameters. The ratio of kernel-wise decomposition can be calculated via Equation 13.

$$\begin{aligned} \mathcal{R}_{\text{kernel\_wise}} &= \frac{\mathcal{P}_{\text{kernel\_wise}}}{\mathcal{P}_{\text{full\_ft}}} = \frac{c_{\text{out}} \cdot (c_{\text{in}} \cdot r \cdot k + r \cdot k)}{c_{\text{out}} \cdot c_{\text{in}} \cdot k^2} \\ &= \frac{c_{\text{in}} \cdot r + r}{c_{\text{in}} \cdot k} \end{aligned} . \tag{13}$$

We can find that the number of parameters in kernel-wise will be greater than the full fine-tuning one when $r \geq k$ as $\mathcal{R}_{\text{kernel\_wise}} \geq \frac{c_{\text{in}} \cdot k + k}{c_{\text{in}} \cdot k} > 1$. For an example, the representative CNN, ResNet, the kernel size $k$ typically takes values of 1 or 3, so the $r$ should not be greater than 3 for a positive gain. However, the model usually has ill perform when $r$ is too small as our ablation studies in Section 4.6. Therefore, the design of kernel-wise LoRA-C is difficult to achieve a balance between performance and parametric efficiency.

The ratio of layer-wise decomposition is:

$$\begin{aligned} \mathcal{R}_{\text{layer\_wise}} &= \frac{\mathcal{P}_{\text{layer\_wise}}}{\mathcal{P}_{\text{full\_ft}}} = \frac{(c_{\text{out}} + c_{\text{in}}) \cdot r \cdot k}{c_{\text{out}} \cdot c_{\text{in}} \cdot k^2} \\ &= \frac{(c_{\text{out}} + c_{\text{in}}) \cdot r}{c_{\text{out}} \cdot c_{\text{in}} \cdot k} \end{aligned} . \tag{14}$$

In the layer-wise granularity, the number of parameters mainly depends on $c_{\text{out}}, c_{\text{in}}, r$ and $k$. In typical CNN for computer vision, $c_{\text{out}}$ and $c_{\text{in}}$ are generally greater than 3 (or even larger, e.g., 64, 128, 256, ...), so $c_{\text{out}} + c_{\text{in}}$ is usually smaller than $c_{\text{out}} \cdot c_{\text{in}}$. For example, for a $k = 1$ convolution with $c_{\text{in}} = 256$, $c_{\text{out}} = 512$, and rank $r = 128$, the $\mathcal{R}_{\text{layer\_wise}} = \frac{(512+256) \cdot 128}{512 \cdot 256 \cdot 1} = 0.75 < 1$ can also maintain
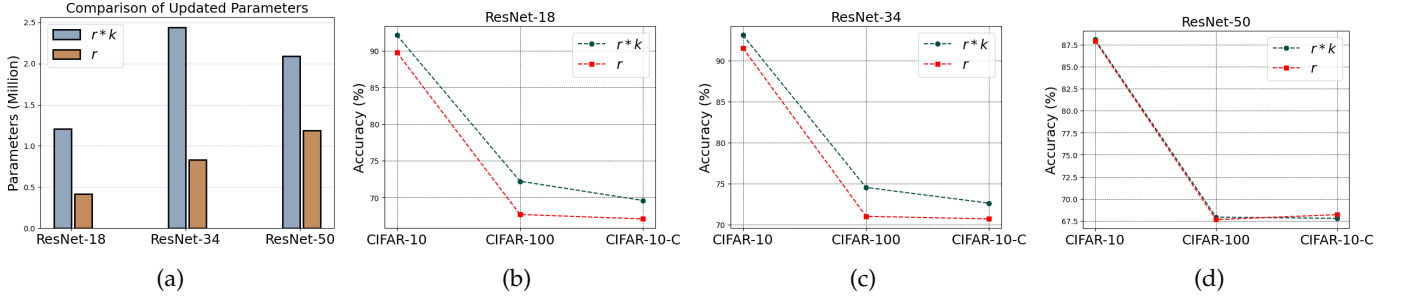
Fig. 4: Comparison of updated parameters, performance, and robustness. (a) Comparison of updated parameters under two settings. (b) Performance comparison of the two settings when using ResNet-18. (c) Comparison of two settings when using ResNet-34. (d) Performance comparison of the two settings when using ResNet-50. We conclude that the second setting is used on models with model parameters greater than or equal to ResNet-50, and the first setting is used on models with model parameters less than ResNet-50.

a positive gain of updated parameters reduction in a higher rank $r$ ($r = 128 \gg k = 1$).

Compared to kernel-wise decomposition which causes parameter inflation, the layer-wise one has fewer fine-tuning parameters for updating. Therefore, we choose layer-wise decomposed matrices as the decomposition in LoRA-C.

### 3.3.2 Impact of Decomposed Matrices Rank

There are also two settings of the rank of $\Delta \boldsymbol{W}$ when fine-tuning the CNN model using LoRA. The first is to associate the rank with the size of the convolutional kernel, that is, the rank is $r * k$. The second is to directly set rank to $r$. The number of updated parameters for the latter is $1/k$ of the former. *How to choose between these two settings?*

To this end, on CIFAR-10/-100 datasets, based on ResNet-18/-34/-50, we set $r$ to $\{1, 2, 4, 8, 16, 32, 64\}$ and $\alpha$ to $\{1, 2, 4, 8, 16, 32, 64, 128\}$ for traversal. The average amount of updated parameters for each model under different settings, the average performance of each model under each setting on each dataset, and the robustness on CIFAR-10-C dataset are shown in Fig. 4.

As shown in Fig. 4a, the amount of updated parameters under the first setting is more than twice that of the second. Correspondingly, as shown in Fig. 4b and Fig. 4c, the model accuracy under the first setting is higher than that under the second setting. However, in the case of the ResNet-50, as shown in Fig. 4d, the model accuracy under the two settings is essentially equal, and even on CIFAR-10-C, the model under the second setting is more robust.

The two settings seem to be impossible to choose. However, we see from Fig. 4 that the gap in accuracy between the two settings is narrowing when we go from ResNet-18 to ResNet-34 and then to ResNet-50. For example, on the CIFAR-10 dataset, when using ResNet-18, the model accuracy of the first setting is 2.33% higher than that of the second one. When using ResNet-34, the model accuracy of the first setting is 1.55% higher than that of the second. When using ResNet-50, the accuracy of the two settings is roughly the same.

It is worth mentioning that the second setting outperforms the first one in terms of improving model robustness as the model becomes larger in CIFAR-10-C. For example, when using ResNet-18, the model accuracy of the first

setting is 2.51% higher than the model accuracy of the second one. When using ResNet-34, the model accuracy of the first setting is 1.92% higher than the model accuracy of the second one. While using ResNet-50, the model accuracy of the second setting exceeds that of the first.

To this end, we comprehensively consider the number of updated parameters and model accuracy and conclude that the second setting is used on models with model parameters greater than or equal to ResNet-50, and the first setting is used on models with model parameters less than ResNet-50. The experiments in this paper follow this setting.

### 3.4 Procedures of LoRA-C on IoT Devices

#### 3.4.1 LoRA-C Fine-Tuning

The CNN fine-tuning procedure of LoRA-C consist of two stages: preliminary of low-rank decomposition and tuning the pre-trained CNN model, as depicted in Algorithm 1.

---

**Algorithm 1:** LoRA-C Fine-Tuning

**Input:** The pre-trained CNN model $\mathcal{M}$, local corrupted datasets $\mathcal{D}_{\text{local}}$

**Output:** The robust CNN model $\mathcal{M}_{\text{robust}}$

1 **Preliminary of Low-Rank Decomposition**:
2 Local corrupted datasets $\mathcal{D}_{\text{local}}$ initialization;
3 **for** *layer $\ell$ in $\mathcal{M}$* **do**
4     Set the rank $r$ of decomposed matrix $\Delta \boldsymbol{W}_{\ell}$;
5     Register LoRA-C branch with $\boldsymbol{A}_{\ell}$, $\boldsymbol{B}_{\ell}$ matrices as learnable parameters in layer-wise;
6     Freeze pre-trained convolutional weights $\boldsymbol{W}_{\ell}$ by removing gradients;

7 **Tuning the pre-trained CNN model $\mathcal{M}$**:
8 **for** *epoch in $1, 2, \ldots, epoch_{\max}$* **do**
9     **for** *batch in $\mathcal{D}_{\text{local}}$* **do**
10         *batch* is a mini-batch from $\mathcal{D}_{\text{local}}$;
11         Activate LoRA-C branches with decomposed matrices $\boldsymbol{A}$, $\boldsymbol{B}$;
12         Compute predictions by $\mathcal{M}$.forward();
13         Compute gradients of matrices $\boldsymbol{A}$, $\boldsymbol{B}$;
14         Update parameters in LoRA-C branches;

15 **Return** $\mathcal{M}_{\text{robust}}$;

---

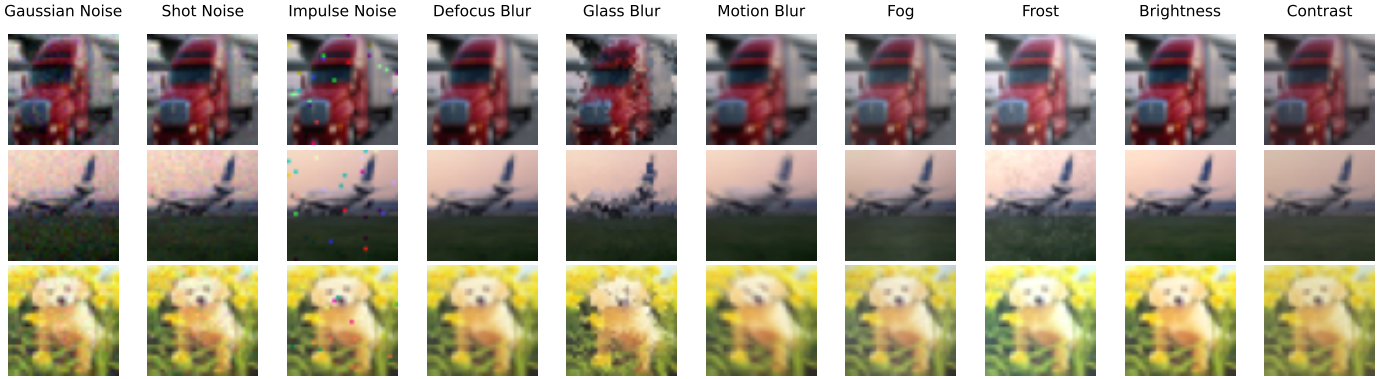| Gaussian Noise | Shot Noise | Impulse Noise | Defocus Blur | Glass Blur | Motion Blur | Fog | Frost | Brightness | Contrast |

Fig. 5: Visualization examples of commonly observable corruptions in CIFAR-10-C dataset.



Fig. 6: Visualization examples of the Icons-50 dataset. For each category, the image in the upper left corner is from the test set, and the other three are from the training set. As shown, the style of training data is very different from that of test data.

*Preliminary of Low-Rank Decomposition*: Before tuning, we have to prepare the settings of the pre-trained CNN model sent by the cloud server. We need to add new branches to each convolution layer, referred to as the LoRA-C branch. The LoRA-C branch initializes the $A_\ell$ and $B_\ell$ matrices at each layer $\ell$. At the same time, the gradients need to be removed for each pre-trained convolutional weights and retained only for the LoRA-C branches to reduce the parameter burden of fine-tuning.

*Tuning the Pre-trained CNN Model*: The on-device fine-tuning process is similar to model training, but usually using with fewer epochs. The difference is that only the $A$ and $B$ matrices compute the gradients and update the parameters during tuning iteration. Therefore, the small amount of computational overhead bringing by parametric efficiency enables the CNN fine-tuning of LoRA-C to adapt to the computational resources of IoT devices.

### 3.4.2 LoRA-C CNN Inference

After fine-tuning with local datasets, the IoT device can get a robust CNN model. The pre-trained model can be removed if storage resources are insufficient, and only storing the updated LoRA-C branches. This is because the pre-trained

model can be sent by cloud assist and LoRA-C does not change pre-trained weights. The IoT device only needs to maintain the portion of the incremental updates.

---

**Algorithm 2:** LoRA-C Inference

---
**Input:** The robust CNN model $\mathcal{M}_{\text{robust}}$
**Output:** The robust CNN inference model $\mathcal{M}_{\text{inference}}$

1 **Convert the $\mathcal{M}_{\text{robust}}$ to the $\mathcal{M}_{\text{inference}}$:**
2 **for** *layer $\ell$ in $\mathcal{M}_{\text{robust}}$* **do**
3      Recompose $A_\ell$, $B_\ell$ matrices to $\Delta W_\ell$;
4      Merge the updated incremental matrix $\Delta W_\ell$ into frozen pre-trained conv weights $W_\ell$ by coordinate-wise summation;
5      Streamline this layer further with other re-parameterization techniques, e.g., BN fusion;
6 **Return $\mathcal{M}_{\text{inference}}$;**

---

For on-device CNN model inference, as shown in Algorithm 2, we only need to recompose the decomposed matrices $A_\ell$ and $B_\ell$ of the LoRA-C branches into the increments $\Delta W_\ell$ in each layer $\ell$, then fuse them to pre-trained weights by coordinate-wise summation to eliminate the branches. Finally, we can streamline the model further by some other re-parameterization techniques [40], e.g., batch normalization fusion. In this case, the neural architecture of LoRA-C CNN model for inference is exactly the same as the pre-trained one. Therefore, there is not any degradation of IoT quality of services in terms of inference latency.

## 4 EXPERIMENTS

This section introduces our experimental setup, including this work's datasets, baselines, and implementation details. Then, we conduct experimental comparisons and present analyses.

### 4.1 Experimental Setup

**Datasets.** We conduct experiments on four benchmark datasets, i.e., CIFAR-10, CIFAR-100, CIFAR-10-C and Icons-50 datasets. CIFAR-10 and CIFAR-100 are two standard datasets, and CIFAR-10-C and Icons-50 are used as corrupted datasets. Among them, noise is introduced into the validation set of CIFAR-10-C (see Fig. 5 for a visualization),

TABLE 1: Comparison between training from scratch (SCR), full fine-tuning (FT), and LoRA-C on CIFAR-10/-100 datasets. Our results are highlighted with shading. #P and #P (LoRA-C.) refer to the number of parameters that need to be updated when training models and when using our proposed method, respectively. Our method achieves better performance with less number of updated parameters.

| Models | #P | CIFAR-10 | | | | CIFAR-100 | | | |
| | | #P (LoRA-C.) | Acc. (SCR/FT) | Acc. (LoRA-C.) | $\Delta_{Acc}(\uparrow)$ | #P (LoRA-C.) | Acc. (SCR/FT) | Acc. (LoRA-C.) | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|---|---|---|
| ResNet-18 | 11.23 | 4.26 | 95.45/94.96 | 95.69 | **0.24** | 4.26 | 76.88/76.55 | 79.93 | **3.05** |
| ResNet-34 | 21.29 | 0.29 | 95.08/95.48 | 96.06 | **0.98** | 2.16 | 78.38/78.12 | 82.49 | **4.11** |
| ResNet-50 | 23.52 | 3.61 | 95.34/95.70 | 96.59 | **1.25** | 2.26 | 77.69/78.68 | 82.98 | **5.29** |
| ResNet-101 | 42.51 | 1.95 | 95.77/95.25 | 97.10 | **1.33** | 4.49 | 79.64/78.65 | 84.52 | **4.88** |

and the styles of the training data and test data of Icons-50 are pretty different (see Fig. 6 for a visualization).

- CIFAR-10 and CIFAR-100 [62] are multi-class natural object datasets for image classification. They consist of 50, 000 training images and 10,000 test images in 10/100 categories. Each image has a resolution of $32 \times 32$ pixels.
- CIFAR-10-C [63] is a test dataset after using synthetic common perturbations and noise corruptions on the CIFAR-10 test set. It consists of 10, 000 test images of 19 types of damage in 4 categories. The resolution of each image is $32 \times 32$ pixels.
- Icons-50 [63] consists of 10,000 images in 50 categories collected from different companies. We set the resolution of all images to $32 \times 32$ pixels. For training, the data of one company is retained as test data, and the data of the rest of the companies is used as training data.

Please refer to [41] for more explanation of the CIFAR-10-C and Icons-50 datasets.

**Baselines.** We use RestNet series methods as the baselines, including ResNet-18/-34/-50/-101 [25].

**Implementation Details.** We implement our method in Python 3.8 and PyTorch 2.2 with CUDA 12.1. All experiments are performed on NVIDIA RTX 4090 GPUs. We use the pre-trained weights of the ResNet series (i.e., ResNet-18/-34/-50/-101) on ImageNet-1K provided by PyTorch and replace the first $7 \times 7$ convolution in ResNet with a $3 \times 3$ convolution to process $32 \times 32$ image data. In addition, the output head is modified for different datasets. When fine-tuning the model, the updating parameters include the replaced first layer $3 \times 3$ convolution, the last fully connected layer, the added $A$ and $B$ matrices, and other parameters are frozen by removing gradients. We set $r$ to $\{1, 2, 4, 8, 16, 32, 64\}$ and $\alpha$ to $\{1, 2, 4, 8, 16, 32, 64, 128\}$ for traversal. We use the SGD optimizer for all datasets and set the weight decay to $5e^{-4}$ and the initial learning rate to $0.1$.

### 4.2 Experimental Results

### 4.3 The amount of Fine-Tuning Parameters

The proposed LoRA-C reduces the number of model updated parameters by one order of magnitude compared to fine-tuning the full model parameters. As shown in Fig. 7, compared with the full parameter fine-tuning ResNet-18, the amount of fine-tuning parameters of LoRA-C-ResNet-18 is reduced by up to about 99.35%. Compared with the
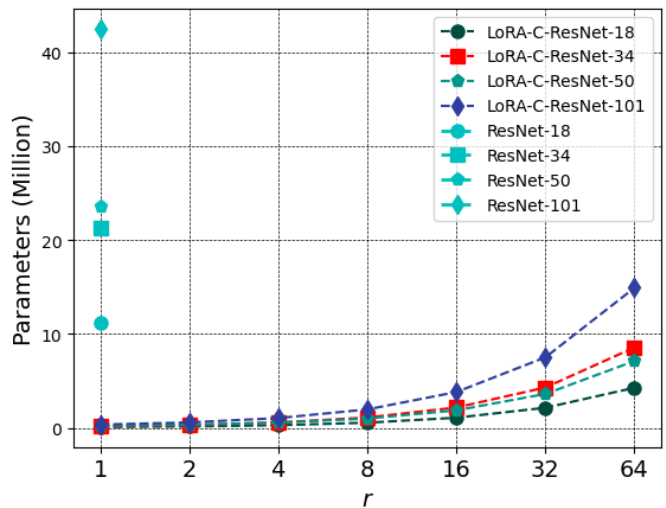


Fig. 7: Comparison of the updated parameters of LoRA-Cs and the standard models.

full parameter fine-tuning ResNet-101, the amount of fine-tuning parameters of LoRA-C-ResNet-101 is reduced by up to about 99.21%. LoRA-C-based fine-tuning freezes the backbone model parameters and only fine-tunes the newly added parameters. The number of newly added parameters is small because LoRA-C performs low-rank decomposition in units of convolutional layers, thus effectively reducing the number of model fine-tuning parameters.

### 4.4 Results on Standard Benchmarks

The proposed LoRA-C achieves improved model performance with only a few parameters fine-tuned. As shown in Table 1, on the CIFAR-10 dataset, based on ResNet-18, our proposed LoRA-C achieves 95.69% accuracy, surpassing the training from scratch (SCR) result by +0.24%. On the CIFAR-100 dataset, our proposed LoRA-C achieves 79.93% accuracy, surpassing the SCR result by +3.05%. This is mainly because LoRA-C relies on the existing knowledge learned by the model on ImageNet. When the main model is frozen, and only the newly added parameters are updated, it is equivalent to allowing the new branches to learn specific features based on the knowledge learned by the model on ImageNet. Therefore, the LoRA-C can achieve higher accuracy than training the model from scratch.

The proposed LoRA-C can improve the model performance when the model accuracy is relatively low. As shown in Table 1, our proposed LoRA-C significantly outperforms

TABLE 2: Robustness to common observable corruptions. SCR means the model is trained from scratch. FT means fine-tuning all parameters using new data based on the pre-trained model. Our results are highlighted with shading.

(a) ResNet18

| Method | Noise | Blur | Weather | Digital | mean | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|
| SCR | 56.48 | 71.88 | 85.20 | 80.49 | 73.51 | – |
| FT | 55.51 | 70.29 | 84.14 | 80.26 | 72.55 | -0.96% |
| LoRA-C (Ours) | 56.07 | 76.89 | 85.47 | 82.11 | **75.14** | +1.63% |

(b) ResNet34

| Method | Noise | Blur | Weather | Digital | mean | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|
| SCR | 58.11 | 71.32 | 85.49 | 80.53 | 73.86 | – |
| FT | 59.1 | 72.33 | 85.82 | 81.18 | 74.61 | +0.75% |
| LoRA-C (Ours) | 66.08 | 81.06 | 87.69 | 84.22 | **79.76** | +5.9% |

(c) ResNet50

| Method | Noise | Blur | Weather | Digital | mean | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|
| SCR | 53.81 | 70.56 | 85.70 | 81.63 | 72.93 | – |
| FT | 54.92 | 71.35 | 85.87 | 81.17 | 73.33 | +0.4% |
| LoRA-C (Ours) | 66.76 | 81.94 | 89.82 | 84.42 | **80.74** | +7.81% |

(d) ResNet101

| Method | Noise | Blur | Weather | Digital | mean | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|
| SCR | 55.10 | 73.03 | 85.30 | 82.35 | 73.95 | – |
| FT | 56.22 | 73.48 | 85.61 | 82.05 | 74.34 | +0.39% |
| LoRA-C (Ours) | 73.78 | 82.87 | 90.58 | 86.53 | **83.44** | +9.50% |

SCR. On the CIFAR-100 dataset, based on ResNet-34, our proposed LoRA-C surpasses the SRC by +4.11%. Based on ResNet-50, our proposed LoRA-C surpasses the SCR by +5.29%. The low performance of the model indicates that it is difficult for the model to extract effective features from the current dataset. LoRA-C can achieve high accuracy because it uses a pre-trained model, that is, it extracts effective features based on the existing knowledge of the pre-trained model. It can also be understood that based on the pre-trained model, the newly added network branch is easier to update parameters through gradient descent.

As the model capacity increases, the model accuracy improves more significantly. As shown in Table 1, on the CIFAR-10 dataset, based on ResNet-18, LoRA-C improves accuracy by 0.24% compared to SCR. Based on ResNet-34, LoRA-C improves accuracy by 0.98% compared to SCR. Based on ResNet-101, LoRA-C improves accuracy by 1.33% compared to SCR. Compared with simple-structured models processing complex data, complex-structured models can achieve higher performance when processing complex data. Complex models are more likely to fall into the gradient vanishing problem, which makes it impossible to effectively update model parameters. However, once the parameters can be effectively updated through gradient descent, their performance often exceeds that of simple models. LoRA-C uses the pre-trained model to update parameters so that the newly added branches can effectively use gradient descent to update their parameters, thereby significantly improving the performance of large-capacity CNN models.

The full fine-tuning (FT) and SCR have similar accuracy. In some cases, the accuracy of FT is even lower than that
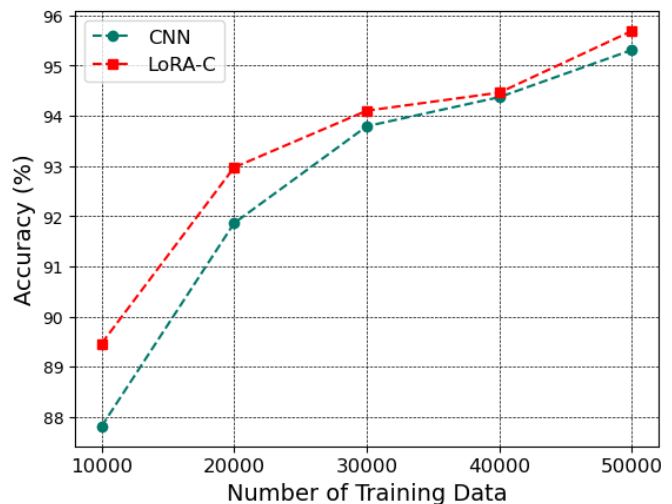


Fig. 8: Robustness to limited training data.

of SCR. For example, based on ResNet-18, on CIFAR-10 dataset, The accuracy of SCR is higher than that of FT. Based on ResNet-34, the accuracy of SCR is lower than that of FT. This may be due to the overfitting problem of FT.

## 4.5 Results on Robustness

### 4.5.1 Limited Training data

LoRA-C outperforms standard CNN with limited training data, as shown in Fig. 8. The backbone network is learned on the ImageNet dataset. According to the characteristics of the CNN model, learning on this basis is conducive to

the gradient transfer of the model. It is more conducive to obtaining parameters enabling the model to perform well.

### 4.5.2 Corrupted Data

The proposed LoRA-C significantly improves the performance of handling corrupted data. As shown in Table 2, based on ResNet-34, compared to the accuracy of SCR, LoRA-C has an accuracy improvement of 5.9%. Based on ResNet-50, compared to the accuracy of SCR, LoRA-C has an accuracy improvement of 7.81%. It is worth mentioning that the LoRA-C-based fine-tuning method achieves better performance in all categories, as shown in Table 2. In LoRA-C, the backbone network parameters are frozen and will not be updated with local data. The parameters of the backbone network are obtained based on ImageNet training, which means that the backbone network retains the knowledge learned on ImageNet. On the one hand, based on the backbone network, the knowledge learned by the backbone network on ImageNet can improve performance by using local data to fine-tune the newly added model parameters. On the other hand, the parameters of the backbone model are frozen, which can play a regularization role. Therefore, the proposed LoRA-C can achieve better results on the corrupted dataset that is, LoRA-C has strong robustness.

### 4.5.3 Data Under Different Styles

The proposed LoRA-C achieves high performance when dealing with different training and test data styles. As shown in Table 3, based on ResNet-18, compared to the accuracy of SCR, LoRA-C has an accuracy improvement of 6.57%. Based on ResNet-34, compared to the accuracy of SCR, LoRA-C has an accuracy improvement 8.48%. In addition, similar to those obtained with corrupted data, the LoRA-C-based fine-tuning method achieves better performance in all categories. The reasons for this are the same as in Section 4.5.2: (i) Learning is performed based on the existing knowledge of the backbone network. (ii) The parameters of the backbone model are frozen, regularizing the model with the newly added parameters. The above results prove that the proposed LoRA-C is highly robust.

## 4.6 Hyperparameter Study

Two hyperparameters have a great impact on the LoRA-C, namely $\alpha$ and $r$. The $\alpha$ measures the proportion of newly added branches compared to the backbone network. The larger the $\alpha$, the greater the proportion of newly added branches, and vice versa. The $r$ represents the rank of $\Delta W$. The larger $r$ is, the more parameters are fine-tuned, and vice versa.

Figs. 9 and 10 illustrate the relationship between model accuracy and $\alpha$ for a given $r$. We observe that the performance of the LoRA-C does not always improve with the increase of $\alpha$. For example, when $r = 2$, LoRA-C achieves the highest accuracy when $\alpha = 4$. When $r = 4$, LoRA-C achieves the highest accuracy when $\alpha = 8$. When $r = 64$, LoRA-C achieves the highest accuracy when $\alpha = 128$. We also observe that the best performance is usually obtained when $\frac{\alpha}{r} = 2$, as shown in Figs. 9 and 10.

We also illustrate the relationship between model accuracy and $r$ for a given $\alpha$, as shown in Fig. 11 and Fig. 12.

We observe that (i) the best performance is usually obtained when $\frac{\alpha}{r} = 2$. (ii) When fine-tuning, it is not the case that the more parameters updated, the better.

In addition, We also illustrate the relationship between model accuracy and $r$ for a given $\alpha$ on CIFAR-10-C dataset, as shown in Fig. 13. The best performance is usually obtained when $\frac{\alpha}{r} = 2$. This discovery provides experience for the widespread application of LoRA-C.

## 5 CONCLUSION AND FUTURE WORK

This paper proposes a fine-tuning method for robust CNNs for IoT devices, LoRA-C, which performs low-rank decomposition in convolutional layers to reduce the number of fine-tuning parameters. By setting the ratio of $\alpha$, which controls the proportion of newly added branches, to the rank of the weight matrix to a constant, the accuracy of the fine-tuned model significantly exceeds that of the fully trained model. Experimental results on CIAFR-10, CIFAR-100, CIFAR-10-C, and Icons50 datasets demonstrate the effectiveness of the proposed LoRA-C.

Given that LoRA-C effectively improves the robustness of the CNN model, this motivates our future work to start from the following two points: (i) Set $\alpha$ to be learnable, and it will learn the optimal value according to the value of $r$ during model training. (ii) Apply LoRA to MonoCNN or SineFM to explore how to use LoRA to fine-tune the CNN model based on nonlinear mapping.

## REFERENCES

[1] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," in *CoRR abs/2302.13971*, 2023, pp. 1–27.

[2] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. Rao, P. Barnes, Y. Tay, N. Shazeer, V. Prabhakaran, E. Reif, N. Du, B. Hutchinson, R. Pope, J. Bradbury, J. Austin, M. Isard, G. Gur-Ari, P. Yin, T. Duke, A. Levskaya, S. Ghemawat, S. Dev, H. Michalewski, X. Garcia, V. Misra, K. Robinson, L. Fedus, D. Zhou, D. Ippolito, D. Luan, H. Lim, B. Zoph, A. Spiridonov, R. Sepassi, D. Dohan, S. Agrawal, M. Omernick, A. M. Dai, T. S. Pillai, M. Pellat, A. Lewkowycz, E. Moreira, R. Child, O. Polozov, K. Lee, Z. Zhou, X. Wang, B. Saeta, M. Diaz, O. Firat, M. Catasta, J. Wei, K. Meier-Hellstern, D. Eck, J. Dean, S. Petrov, and N. Fiedel, "Palm: Scaling language modeling with pathways," *Journal of Machine Learning Research*, vol. 24, no. 2023, pp. 1–113, 2023.

[3] N. Ding, Y. Qin, G. Yang, F. Wei, Z. Yang, Y. Su, S. Hu, Y. Chen, C.-M. Chan, W. Chen, J. Yi, W. Zhao, X. Wang, Z. Liu, H.-T. Zheng, J. Chen, Y. Liu, J. Tang, J. Li, and M. Sun, "Parameter-efficient fine-tuning of large-scale pre-trained language models," *Nature Machine Intelligence*, vol. 5, no. 3, pp. 220–235, 2023.

[4] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. de Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for nlp," in *Proceedings of the International Conference on Machine Learning*, 2019, pp. 2790–2799.

[5] J. Pfeiffer, A. Kamath, A. Rücklé, K. Cho, and I. Gurevych, "Adapterfusion: Non-destructive task composition for transfer learning," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*, 2021, pp. 487–503.

[6] A. Rücklé, G. Geigle, M. Glockner, T. Beck, J. Pfeiffer, N. Reimers, and I. Gurevych, "Adapterdrop: On the efficiency of adapters in transformers," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 7930–7946.
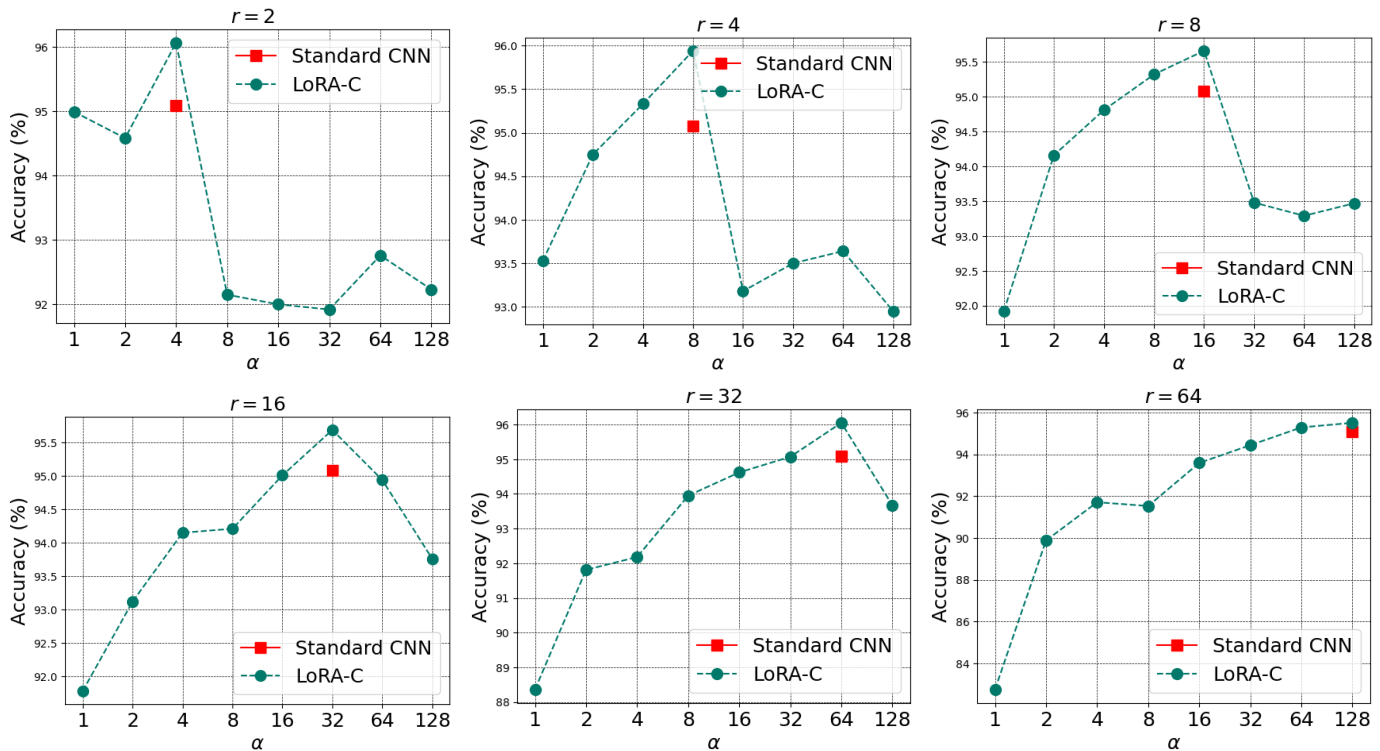
TABLE 3: Robustness to different styles. Our results are highlighted with shading.

(a) ResNet18

| Method | Apple | Facebook | Google | Samsung | mean | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|
| SCR | 93.99 | 90.82 | 85.75 | 85.71 | 89.07 | – |
| FT | 95.32 | 91.11 | 86.85 | 86.52 | 89.95 | +0.88% |
| LoRA-C (Ours) | 98.57 | 95.98 | 94.49 | 93.51 | **95.64** | +6.57% |

(b) ResNet34

| Method | Apple | Facebook | Google | Samsung | mean | $\Delta_{Acc}(\uparrow)$ |
|---|---|---|---|---|---|---|
| SCR | 93.79 | 88.19 | 86.61 | 85.11 | 88.42 | – |
| FT | 95.70 | 90.44 | 86.38 | 83.79 | 89.08 | +0.66% |
| LoRA-C (Ours) | 99.14 | 97.32 | 96.69 | 94.43 | **96.90** | +8.48% |



Fig. 9: Given $r$, the impact of $\alpha$ on the CIFAR-10 dataset.

[7] R. K. Mahabadi, J. Henderson, and S. Ruder, "Compacter: Efficient low-rank hypercomplex adapter layers," in *Advances in Neural Information Processing Systems*, 2021, pp. 1022–1035.

[8] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021, pp. 4582–4597.

[9] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," in *Proceedings of the International Conference on Learning Representations*, 2022.

[10] S. Chen, C. Ge, Z. Tong, J. Wang, Y. Song, J. Wang, and P. Luo, "Adaptformer: Adapting vision transformers for scalable visual recognition," in *Advances in Neural Information Processing Systems*, 2022, pp. 16 664–16 678.

[11] D. Lian, D. Zhou, J. Feng, and X. Wang, "Scaling & shifting your features: A new baseline for efficient model tuning," in *Advances in Neural Information Processing Systems*, 2022, pp. 109–123.

[12] Y. Fan, O. Watkins, Y. Du, H. Liu, M. Ryu, C. Boutilier, P. Abbeel, M. Ghavamzadeh, K. Lee, and K. Lee, "Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models," in *Advances in Neural Information Processing Systems*, 2024, pp. 1–28.

[13] D. J. Kopiczko, T. Blankevoort, and Y. M. Asano, "Vera: Vector-based random matrix adaptation," in *Proceedings of the International Conference on Learning Representations*, 2024, pp. 1–21.

[14] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "Qlora: Efficient finetuning of quantized llms," in *Advances in Neural Information Processing Systems*, 2023, pp. 10 088–10 115.

[15] Q. Zhang, M. Chen, A. Bukharin, P. He, Y. Cheng, W. Chen, and T. Zhao, "Adaptive budget allocation for parameter-efficient fine-tuning," in *Proceedings of the International Conference on Learning Representations*, 2023, pp. 1–17.

[16] Y. Li, Y. Yu, C. Liang, N. Karampatziakis, P. He, W. Chen, and T. Zhao, "Loftq: Lora-fine-tuning-aware quantization for large language models," in *Proceedings of the International Conference on Learning Representations*, 2024, pp. 1–23.

[17] M. Zhang, H. Chen, C. Shen, Z. Yang, L. Ou, X. Yu, and B. Zhuang, "Loraprune: Structured pruning meets low-rank parameter-efficient fine-tuning," in *Proceedings of the Findings of the Association for Computational Linguistics*, 2024, pp. 3013–3026.

[18] S.-Y. Liu, C.-Y. Wang, H. Yin, P. Molchanov, Y.-C. F. Wang, K.-T. Cheng, and M.-H. Chen, "Dora: Weight-decomposed low-rank adaptation," in *Proceedings of the International Conference on Machine*
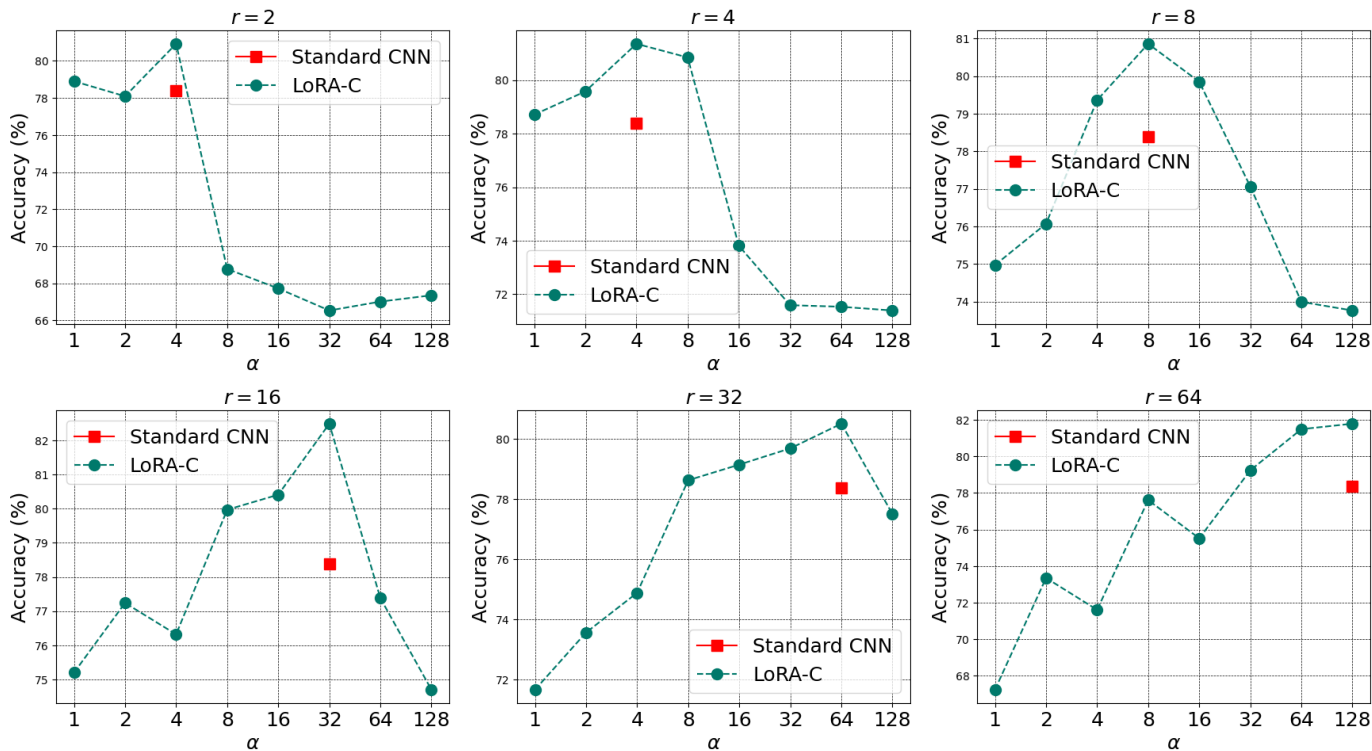
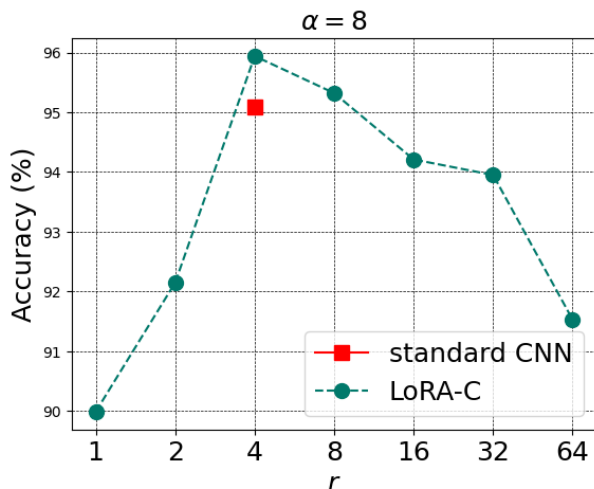Fig. 10: Given $r$, the impact of $\alpha$ on the CIFAR-100 dataset.



Fig. 11: Given $\alpha = 8$, the impact of $r$ on CIFAR-10 dataset.


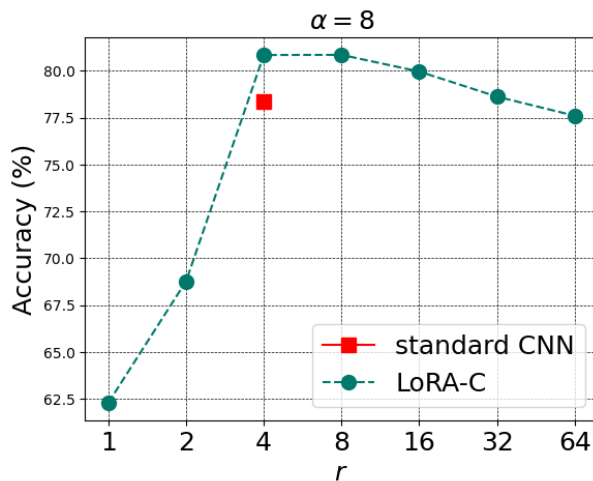
Fig. 12: Given $\alpha = 8$, the impact of $r$ on CIFAR-100 dataset.

*Learning*, 2024, pp. 1–23.

[19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.

[20] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 012–10 022.

[21] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, F. Wei, and B. Guo, "Swin transformer v2: Scaling up capacity and resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 009–12 019.

[22] K. Simonyan and A. Zisserman, "Very deep convolutional net-

works for large-scale image recognition," in *Proceedings of the International Conference on Learning Representations*, 2015, pp. 1–14.

[23] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6848–6856.

[24] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2261–2269.

[25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[26] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1106–1114.
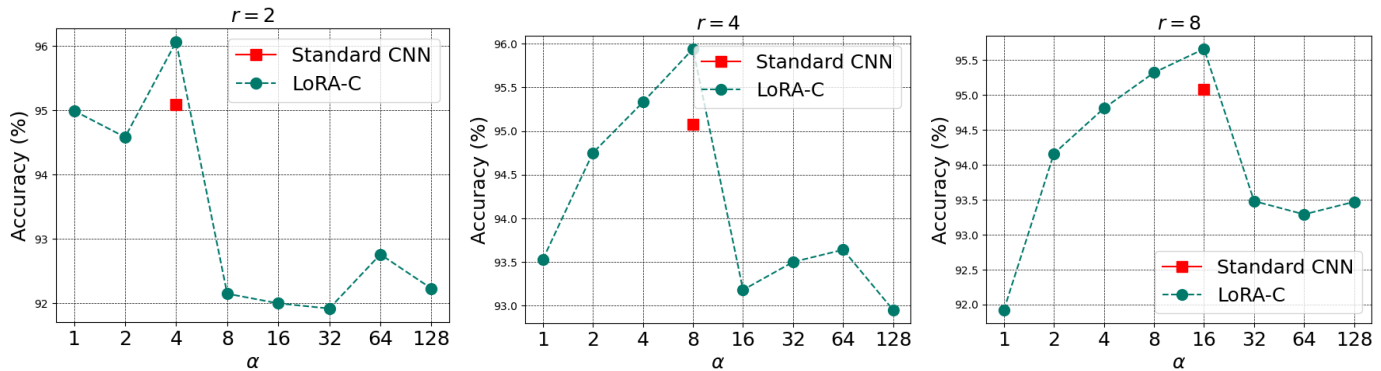
Fig. 13: Given $\alpha$, the impact of $r$ on the CIFAR-10-C dataset.

[27] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.

[28] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proceedings of the International Conference on Machine Learning*, 2015, pp. 448–456.

[29] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.

[30] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017, pp. 4278–4284.

[31] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Proceedings of the European Conference on Computer Vision*, 2016, pp. 630–645.

[32] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1251–1258.

[33] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7132–7141.

[34] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," in *CoRR abs/1704.04861*, 2017, pp. 1–9.

[35] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

[36] A. Howard, R. Pang, H. Adam, Q. V. Le, M. Sandler, B. Chen, W. Wang, L.-C. Chen, M. Tan, G. Chu, V. Vasudevan, and Y. Zhu, "Searching for mobilenetv3," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1314–1324.

[37] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "Shufflenet v2: Practical guidelines for efficient cnn architecture design," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 122–138.

[38] K. Han, Y. Wang, Q. Tian, J. Guo, C. Xu, and C. Xu, "Ghostnet: More features from cheap operations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1577–1586.

[39] Y. Tang, K. Han, J. Guo, C. Xu, C. Xu, and Y. Wang, "Ghostnetv2: Enhance cheap operation with long-range attention," in *Advances in Neural Information Processing Systems*, 2022, pp. 9969–9982.

[40] X. Ding, X. Zhang, N. Ma, J. Han, G. Ding, and J. Sun, "Repvgg: Making vgg-style convnets great again," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13733–13742.

[41] C. Ding, Z. Lu, F. Juefei-Xu, V. N. Boddeti, Y. Li, and J. Cao, "Towards transmission-friendly and robust cnn models over cloud and device," *IEEE Transactions on Mobile Computing*, vol. 22, no. 10, pp. 6176–6189, 2023.

[42] J. Yu, C. Ding, H. Li, Y. Jin, and Y. Li, "Localized knowledge distillation helps iot devices provide high-performance visual services," in *Proceedings of the IEEE International Conference on Web Services*, 2023, pp. 170–178.

[43] Z. Lu, C. Ding, S. Wang, R. Cheng, F. Juefei-Xu, and V. N. Boddeti, "Seed feature maps-based cnn models for leo satellite remote sensing services," in *Proceedings of the IEEE International Conference on Web Services*, 2023, pp. 415–425.

[44] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. N. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," in *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems*, 2017, pp. 615–629.

[45] X. Zhang, M. Qiao, L. Liu, Y. Xu, and W. Shi, "Collaborative cloud-edge computation for personalized driving behavior modeling," in *Proceedings of 4th ACM/IEEE Symposium on Edge Computing*, 2019, pp. 209–221.

[46] S. Laskaridis, S. I. Venieris, M. Almeida, I. Leontiadis, and N. D. Lane, "Spinn: Synergistic progressive inference of neural networks over device and cloud," in *Proceedings of 26th Annual International Conference on Mobile Computing and Networking*, 2020, pp. 1–15.

[47] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," in *Proceedings of International Conference on Learning Representations*, 2016, pp. 1–14.

[48] Z. Lu, C. Ding, F. Juefei-Xu, V. N. Boddeti, S. Wang, and Y. Yang, "Tformer: A transmission-friendly vit model for iot devices," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 2, pp. 598–610, 2023.

[49] C. Ding, A. Zhou, Y. Liu, R. N. Chang, C.-H. Hsu, and S. Wang, "A cloud-edge collaboration framework for cognitive service," *IEEE Transactions on Cloud Computing*, vol. 10, no. 3, pp. 1489–1499, 2022.

[50] Y. Mao, L. Mathias, R. Hou, A. Almahairi, H. Ma, J. Han, S. Yih, and M. Khabsa, "Unipelt: A unified framework for parameter-efficient language model tuning," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022, pp. 6253–6264.

[51] B. Lester, R. Al-Rfou, and N. Constant, "The power of scale for parameter-efficient prompt tuning," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 3045–3059.

[52] X. Liu, K. Ji, Y. Fu, Z. Du, Z. Yang, and J. Tang, "P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks," in *CoRR abs/2110.07602*, 2021, pp. 1–8.

[53] C. Li, H. Farkhoor, R. Liu, and J. Yosinski, "Measuring the intrinsic dimension of objective landscapes," in *Proceedings of the International Conference on Learning Representations*, 2018, pp. 1–23.

[54] A. Aghajanyan, S. Gupta, and L. Zettlemoyer, "Intrinsic dimensionality explains the effectiveness of language model fine-tuning," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021, pp. 7319–7328.

[55] S. Aleem, J. Dietlmeier, and S. L. Eric Arazo, "Convlora and adabn based domain adaptation via self-training," in *Proceedings of the IEEE International Symposium on Biomedical Imaging*, 2024, pp. 1–5.

[56] D. Aggarwal, Y. Mittal, and U. Kumar, "Advancing image classification through parameter-efficient fine-tuning: A study on loRA with plant disease detection datasets," in *The Second Tiny Papers Track at ICLR 2024*, 2024.

[57] Z. Zhong, Z. Tang, T. He, H. Fang, and C. Yuan, "Convolution meets loRA: Parameter efficient finetuning for segment anything model," in *Proceedings of the International Conference on Learning Representations*, 2024.

[58] S.-Y. YEH, Y.-G. Hsieh, Z. Gao, B. B. W. Yang, G. Oh, and Y. Gong, "Navigating text-to-image customization: From lyCORIS fine-tuning to model evaluation," in *Proceedings of the International Conference on Learning Representations*, 2024.

[59] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.

[60] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[61] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1026–1034.

[62] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[63] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," in *Proceedings of the International Conference on Learning Representations*, 2019.