

Frequency Regularization: Efficient Neural Network Compression for Mobile Device Deployment

Wenhao You(wyou1@ualberta.ca) Leo Chang(chewei@ualberta.ca)

Abstract

In general, neural networks require high-specification hardware because of their complexity. However, during the daily life, people have the requirements for running neural networks applications on their own mobile devices but have limited hardware. Even the researchers have to use the computer in their lab instead of running the pre-trained models by using mobile devices. In this proposal, we focus on how to implement a simple and creative compression algorithm that is applied to a basic neural network in an Android mobile device, where the algorithm is called Frequency Regularization (FR). If the algorithm can be implemented on the mobile device, the general neural networks can also be attempted to run on limited hardware devices. As this algorithm does not reduce lots of accuracy compared to the un-compressed neural network, it has a high possibility that we can overcome to run neural networks on the limited hardware devices and without the loss of accuracy.

Introduction

Neural networks have been implemented on many applications and performed a great results. However, these neural networks are mostly implemented on a high-specification hardware such as computers since there are many parameters in different neural networks. The more parameters in the neural networks, the more storages and internet transmission rates are occupied. In this research, we try to look for a method that can make the neural networks run on a limited hardware devices, such as mobile devices. Although mobile devices have more enhanced hardware nowadays, but because of its limited size, the efficacy is still not as powerful as the computers. There are several benefits of running the neural networks on limited hardware devices, specifically, mobile devices. The most important part is the privacy and the data security. Processing data on-device ensures that personal information does not need to upload or transmit to the cloud servers, which enhanced privacy.

Running neural networks on mobile devices can also bring other advantages. If we are able to implement the neural network models directly on the mobile devices, it can reduce our dependence on the internet connection to some extent and since mobile devices play a significant role in our daily lives, it can be easily access. Moreover, for some applications which need real-time feedback, such as object recognition and image segmentation, neural networks on mobile devices can also make the process faster to improve user experience. The most important part of modern life is that we can keep the data processing and inference in our own internal system of the mobile devices rather than sending it to a remote server, which can increase the processing speed.

The new method called Frequency regularization (FR) [1] is based on the frequency domain to compress the neural network. It can be divided into two steps: dynamic tail-truncation and inverse discrete cosine transform (IDCT). The main idea of this method is that some of the neural networks parameters may be redundant in their representation. The algorithm focus on using the frequency domain transform for network regularization to restrict information redundancy during the training process and introduce FR as shown in Fig. 1. Therefore, the redundant part can be removed without any side effect on the performance of the network.

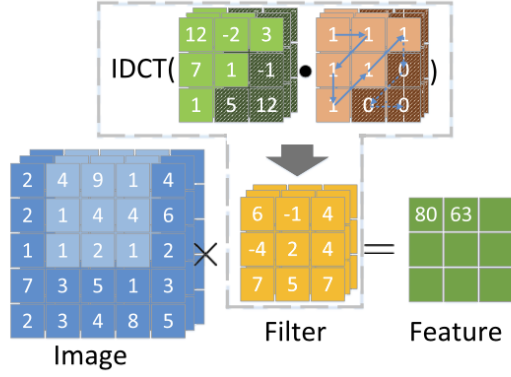


Figure 1: Illustration of the proposed frequency regularization. The tail elements of tensors in the frequency domain are truncated first, then input into the inverse of the discrete cosine transform to reconstruct the spatial tensor for learning features [1].

simple, efficient, and creative method to run the large neural network on mobile devices based on Frequency regularization algorithm. Our purpose is not only to accomplish this implementation but also to make sure the performance is similar to the result running on the high-specification hardware such as computers. We hope to provide a fast and secure neural network application experience on mobile devices.

In order to implement the neural networks on mobile devices, we have to think about the limitations of hardware and its characteristics. Many popular deep learning libraries, such as PyTorch and TensorFlow, are not fully complete and perfect to implement directly on the mobile devices. Thus, we need to edit some configurations in order to optimize specifically for the mobile devices. By using this method, we can convert our standard neural network models into a lighter version that can fit the mobile devices hardware better. Moreover, beyond this idea, we can also consider improving the hardware parts such as creating special neural network processors for the mobile. In this paper, we will only explore the method in the software portion instead of the hardware.

In conclusion, this paper explores a simple, efficient, and creative method to run the large neural network on mobile devices based on Frequency regularization algorithm.

Brief Summary of Existing Work

In this section, we cover the methods about implementing neural networks on the mobile over three sections. In Section 0.1, upgrading hardware on mobile devices and which part are introduced. At the algorithm level, there is an algorithm called NestDNN [2] based on the dynamics or runtime resources to enable resource-aware multi-tenant on-device deep learning for mobile vision systems as shown in Section 0.2. Moreover, there is an algorithm called “one-shot whole network compression” [3] to compress convolutional neural networks (CNNs) for mobile deployment as shown in Section 0.3.

0.1 Upgrade Hardware

Hardware accelerators such as the digital signal processors offer solutions to overcome these challenges. Many manufacturers have already produced the better hardware. For example, Apple presented the Metal Performance Shaders with support of CNNs accelerated by GPU. This is a solution built on top of the Metal API and allows custom operations. Since there is an experiment [4] which confirmed the feasibility of Big Data analysis on modern mobile devices. Moreover, the hardware on Android devices are upgraded by four main companies, such as Qualcomm, HiSilicon, MediaTek, and Samsung [5], which can be used for neural networks as well.

0.2 Algorithm NestDNN

NestDNN is a framework that takes the dynamics of runtime resources into account to enable resource-aware multi-tenant on-device deep learning for mobile vision systems. NestDNN enables each deep learning model to offer flexible resource-accuracy trade-offs [2]. Fig. 2 illustrates the architecture of NestDNN detailedly, which is split into an offline stage and an online stage. The offline stage consists of three phases: model pruning, model recovery, and model profiling [2]. There are five models involved in NestDNN as shown below:

- **Vanilla Model:** Off-the-shelf deep learning model (e.g., ResNet) trained on a given dataset (e.g., ImageNet).

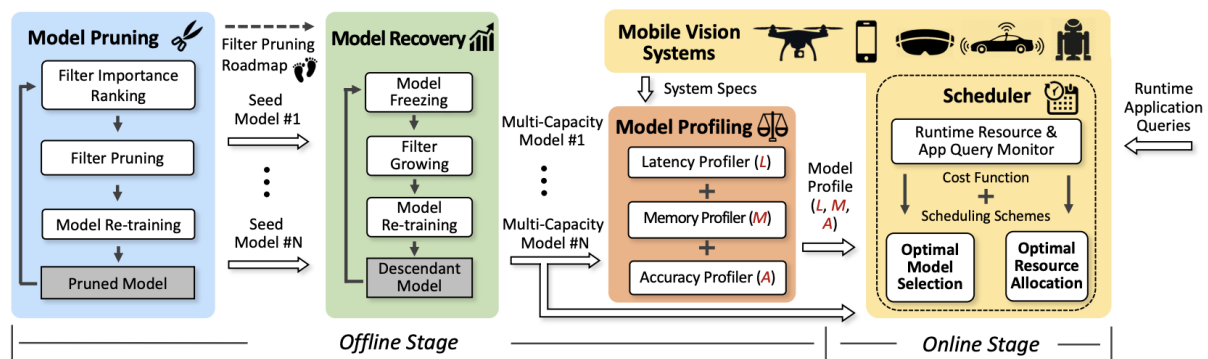


Figure 2: NestDNN architecture [2].

- **Pruned Model:** Intermediate result obtained in model pruning stage.
- **Seed Model:** The smallest pruned model generated in model pruning which meets the minimum accuracy goal set by the user. It is also the starting point of model recovery stage.
- **Descendant Model:** A model grown upon the seed model in model recovery stage. It has a unique resource-accuracy trade-off.
- **Multi-Capacity Model:** The final descendant model that has the capacities of all the previously generated descendant models.

0.3 Algorithm One-shot Whole Network Compression

"One-shot Whole Network Compression" is a simple way to compress CNNs to make them more suitable for fast and low-power applications on the mobile. This algorithm can be divided into three steps: rank selection, tensor decomposition, and fine-tuning. Each step can be implemented easily by using public tools such as VBMF for rank selection, Tucker decomposition for tensor decomposition, and Caffe for fine-tuning [3].

What you Plan To Do

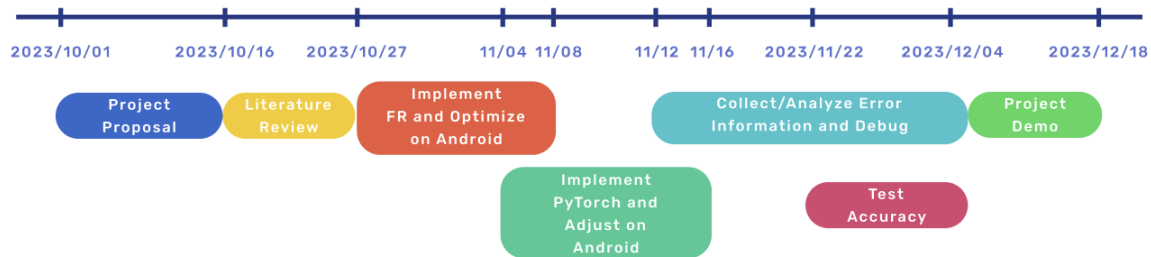
There are five sub-topics for implementing Frequency Regularization on the Android devices. The first one "Implement Frequency Regularization" have done in Basu's lab in University of Alberta (2023). The remaining topics will be finished before the lecture ends.

- **Implement Frequency Regularization:** Implement the Frequency Regularization in order to compress the redundant information on neural network. Make sure that the accuracy and velocity is acceptable for the neural network.
- **Configure Emulator:** In this paper, we choose Android Studio as our emulator. Install it successfully on our own local machines.
- **Implement PyTorch Library:** Implement PyTorch Android library to the emulator. Collect and analyze the error information if we get.
- **Adjust Dependencies:** Potential errors may probably caused by Android PyTorch Library. Since the library on android lacks of improvement not as much as laptop.
- **Adjust Frequency Regularization:** We can also try to change something on the algorithm to let the new one fitting the Android system.

How You Plan to Implement Your Ideas

In this research, we are only going to put the compressed neural network to the mobile devices and make it run on the emulator.

Project Timeline



Short Description of 5 LABS

In this paper, we will divide the project into five labs (milestones), which can help us to double check our progress and adjust it timely.

- Deploy the Frequency Regularization algorithm with the UNet (one basic Neural Network) on the computer and adjust some parameters for the Android emulator.
- Attempt to implement PyTorch on the emulator and summarize the error messages.
- Analyze the error messages and try to do the debug process.
- xx
- Test all functions and neural networks to check the effectiveness and accuracy on the emulator after importing the compressed neural network. We compare the result between the computer and mobile devices.

References

- [1] C. Zhao, G. Dong, S. Zhang, Z. Tan, and A. Basu. Frequency regularization: Reducing information redundancy in convolutional neural networks. *IEEE Access*, September 2023. Department of Computing Science, University of Alberta, Edmonton, AB.
- [2] Biyi Fang, Xiao Zeng, and Mi Zhang. Nestdnn: Resource-aware multi-tenant on-device deep learning for continuous mobile vision. pages 115–127, 10 2018.
- [3] Yong-Deok Kim, Eunhyeok Park, Sungjoo Yoo, Taelim Choi, Lu Yang, and Dongjun Shin. Compression of deep convolutional neural networks for fast and low power mobile applications. 11 2015.
- [4] Anton Akusok, Leonardo Espinosa Leal, Kaj-Mikael Björk, and Amaury Lendasse. High-performance elm for memory constrained edge computing devices with metal performance shaders. In Jiuwen Cao, Chi Man Vong, Yoan Miche, and Amaury Lendasse, editors, *Proceedings of ELM2019*, pages 79–88, Cham, 2021. Springer International Publishing.
- [5] Andrey Ignatov, Radu Timofte, William Chou, Ke Wang, Max Wu, Tim Hartley, and Luc Van Gool. Ai benchmark: Running deep neural networks on android smartphones. In Laura Leal-Taixé and Stefan Roth, editors, *Computer Vision – ECCV 2018 Workshops*, pages 288–314, Cham, 2019. Springer International Publishing.