

# Understanding Deep Learning Model Compression on Heterogeneous Mobile Platforms

Qin Qing<sup>†</sup>, Jie Ren<sup>‡</sup>, Ling Gao<sup>†</sup>, Hai Wang<sup>†</sup>, Zheng Wang<sup>\*</sup>

<sup>†</sup>Northwest University, China, <sup>‡</sup>Shaanxi Normal University, China, <sup>\*</sup>Lancaster University, UK

<sup>†</sup> qq@stumail.nwu.edu.cn, <sup>‡</sup> renjie@snnu.edu.cn, <sup>†</sup>{gl, hwang}@nwu.edu.cn, <sup>\*</sup> z.wang@lancaster.ac.uk

## Abstract

## Keywords

*workload characterisation*

development. There are three mainstream compression methods, parameter pruning and sharing [FIX:\[\]](#), low rank decomposition [FIX:\[\]](#) and knowledge purification [FIX:\[\]](#). However, these methods have not tested on the mobile or embedded system. Our work focus on compressing the deep model and testing the performance on mobile platform. It is a necessary step to further verify the feasibility of the model's deployment on the mobile or embedded system.

## 1. Introduction

(Deep Learning has led an AI renaissance of sorts in recent years. Image is one of its most prominent success area. Results of the ImageNet Challenge points to a dramatic improvement in object detection, localization and classification.. So far, the accuracy of the image classification model has reached more than 95% (top5) [FIX:\[\]](#). While, the excellent performance of these deep networks depends on the powerful computing devices (i.e cloud platform) to deal with a huge number of operations for one inference (up to  $O(10^9)$ ) [FIX:\[\]](#). It is expected that the size of the deep neural networks (i.e., number of weights) will be larger for the more difficult task than the simpler task and thus require more energy to deal with the growing number of parameters and operations.

The current neural network based mobile applications all rely on the cloud server by putting the load on the server and waiting for the inference results back to the mobile, such as [FIX:\[\]](#) using the cloud-based APIs to recognize the xxx lables from an image, which leverages the power of cloud platform's machine learning technology to give a high level of accuracy. However, since the communication delay, the mobile applications which provide neural network based services can not satisfy the requirement for low latency. Therefore, these deep neural networks need to be migrated to mobile platforms to reduce the inference time. While the training and inference of an advanced deep learning model usually requires a large amount of computing resources, memory, and computing power, which becomes a huge obstacles for mobile devices and IoT devices.

Model compression techniques make it possible to run the inferences on mobile devices by reducing the model size and inference time. In recent years, this field has achieved great

In this paper, we compress 10 deep models which belong to four typical networks, mobilenet, vggnet, inception net, resnet by two typical method, prune and quantization. These models are meant to represent the state of the art in deep learning. Most are taken directly from top-tier research venues and have either set new accuracy records on competitive datasets. We focus on four areas: identifying the types of operations which dominate execution time, compressing the different kinds of models, retraining the compressed model to regain the accuracy, comparing the model size and inference time before and after the compression. We evaluate the compressed model on a representative heterogeneous mobile platform NVIDIA Jetson TX2. In summary, our contributions are as follows:

- We test 10 most widely used deep learning models in heterogeneous mobile platform. Testing and analyzing inference time, model size, and breakdown of execution time by operation type for each model. We are the first to do such work, which is the essential step in the migration of the deep learning model in the mobile platform.
- We compress the 10 models by quantization and pruning, and comparing the inference time, model size and accuracy before and after the compression for the first time..
- We find that the inference time varies for different models after using the quantification and prune. and it can be further proved that the quantification method has different effects on different models.

## 2. Background and Motivation

### 2.1. Workloads

**VGGNet:** VGGNet is a deep convolutional neural network developed by the Computer Vision Group at Oxford University and researchers at Google DeepMind. VGGNet explores the relationship between the depth of a convolutional neural network and its performance. By repeatedly stacking  $3 \times 3$  small convolutional kernels and  $2 \times 2$  largest pooled layers, VGGNet successfully builds a 16-19 layer deep convolution. [?] Neural Networks. Compared with the previous network structure, VGGNet significantly reduced the error rate, and achieved the second place of the ILSVRC 2014 competition classification project and the first place of the positioning project. At the same time, VGGNet's scalability is very strong, and the generalization performance of migrating to other image data is very good. VGGNet's structure is very simple, the entire network uses the same size of the convolution kernel size and maximum pool size. So far, VGGNet is still often used to extract image features.

**Inception Net:** Google Inception Net won first place in the ILSVRC 2014 competition. The Inception Net in the game is usually called Inception V1. Its greatest feature is that it controls the calculations and parameters while achieving very good classification performance. Inception V1 proposes Inception Module to increase the utilization of parameters. Although his depth reaches 22 levels, the parameter amount is much lower than AlexNet and VGGNet. Subsequently, Inception V2, Inception V3 and Inception V4 emerged successively in 2015 and 2016. Inception V2 refers to VGGNet and uses two  $3 \times 3$  convolutions instead of a  $5 \times 5$  convolution to reduce the amount of parameters. At the same time, the famous method of Batch Normalization was proposed to speed up the training of the network, and the accuracy of the classification after convergence can also be improved. Inception V3 splits the larger two-dimensional convolutional network into a smaller one-dimensional convolutional network. On the one hand, it saves a lot of parameters, speeds up calculations, and reduces over-fitting. It also adds a layer of non-linear expansion model. Ability to express so that the network can handle more and more abundant spatial features and increase the diversity of features. On the other hand, Inception V3 optimizes the structure of the Inception Module. Subsequent Inception V4 further integrated Resnet on the basis of Inception V3 to further improve accuracy.

**Resnet:** Microsoft Researcher Kaiming He and others proposed Resnet and won the championship in ILSVRC 2015. Resnet's TOP5 error rate was only 3.75%. Resnet uses a connection method called "shortcut connection" to reduce network parameters, and stack the input and output of the block. This simple addition will not add extra parameters and calculations to the network, but it can greatly increase the training speed of the model, improve the training effect, and when the number of layers of the model deepens, this simple structure can solve the problem of degradation well. There are

two different learning units in Resnet's structure: two levels of residual learning units and three levels of residual learning units. The two-level residual learning unit contains two  $3 \times 3$  convolutions with the same number of output channels. The three-level residual network uses a  $1 \times 1$  convolution before and after the middle  $3 \times 3$  convolution. The 50,101,152-layer networks we use are mainly three-tier residual structures.

**MobileNet:** MobileNet was proposed by Howard et al. in 2017. Its architecture consists of a standard convolution layer acting on the input image, a deep separable convolutional stack, and a final average pool and fully connected layer. MobileNet uses a deeply separable convolution to build a lightweight, deep neural network that can integrate standard volumes into a deep convolution and a dot convolution ( $1 \times 1$  convolution kernel). Depth convolution applies each convolution kernel to each channel, and  $1 \times 1$  convolution is used to combine the output of channel convolutions. This decomposition can effectively reduce the computation and reduce the model size. MobileNet has a clear advantage in terms of computational volume and model size.

### 2.2. Motivation

## 3. Experimental Setup

NVIDIA Jetson TX2 is NVIDIA's second-generation CUDA-capable edge device. TX2 runs Linux using a quad-core ARM CPU. It is equipped with an NVIDIA Pascal GPU, which contains specialized architecture for AI applications. TX2 has 8 GB of LPDDR4 RAM and supports PCIe2.0 and various peripherals such as UART, GPIOs, HDMI, USB 3.0 and 2.0, Ethernet, and 802.11ac WLAN. NVIDIA provides the required drivers and CUDA toolkits for TX2 via the JetPack SDK.

## 4. Experimental Results

### 4.1. Overall Results

### 4.2. Compare to Oracle

### 4.3. Analysis

## 5. Related Work

Our work lies at the intersection of multiple research areas: web browsing optimization, task scheduling, energy optimization and predictive modeling. There is no existing work that is similar to ours, in respect to optimizing web workloads across multiple optimization objectives on heterogeneous mobile platforms.

## 6. Conclusions

This paper has presented an automatic approach to optimize the mobile web