# Deployment of Frequency Regularization on Android Mobile Devices

1<sup>st</sup> Wenhao You

Department of Computing Science

University of Alberta

Edmonton, Canada

wyou1@ualberta.ca

2<sup>nd</sup> Leo Chang

Department of Computing Science

University of Alberta

Edmonton, Canada

basu@ualberta.ca

Index Terms—convolutional neural networks, frequency regularization.

#### I. INTRODUCTION

Currently, people can not live without mobile devices. These compact yet powerful gadgets have become indispensable tools for communication, entertainment, and information. Their portability and versatility make them a constant companion. Meanwhile, Convolutional neural networks play an important role in computer vision applications. However, these neural networks are usually implemented on high-specification hardware. There are several advantages of running convolutional neural networks on mobile devices: privacy, internet, and runtime. For enhancing privacy, personal information does not need to be uploaded or transmitted to the cloud servers anymore. For reading the dependence on the internet connection, the functionality on local devices can replace some internet services. For the runtime, especially some applications that need real-time feedback, without connecting to the cloud server can shorten the processing time. In all, convolutional neural networks can totally replace the usage of many applications on mobile devices, ensuring personal data security.

According to the popularity of mobile devices and the benefits of convolutional neural networks, we want to find a way to deploy some large and complex convolutional neural networks on mobile devices, leading to the question: "How can we deploy large convolutional neural networks on mobile devices?"

We found five methods to achieve our goal: upgrade the hardware of mobile devices; use Extreme Learning Machine (ELM) [1] to allocate the weight of the hidden layers randomly in order to train large models on mobile devices faster; use NestDNN [2] dynamically adjust the size and computational complexity of the network based on available resources on mobile devices; use "One-shot Whole Network Compression" [3] to prune, quantize, and compress the neural networks; use Frequency Regularization (FR) [4] to reduce parameters by removing high-frequency component. We make a more

detailed introduction to their drawbacks and limitations in Section II.

After conducting a thorough literature review, considering all the limitations, accuracy, complexity, and future potential, we choose Frequency Regularization (FR) as our target algorithm, deploying it on one of the most popular mobile devices - Android mobile. Our main idea uses FR to compress a convolutional neural network U-Net and then decompresses the uploaded compressed model on an Android mobile device in a short time. After that, use the decompressed model to do image segmentation for the Carvana Image Masking Challenge Dataset [5].

The proposed main idea can be divided into three directions to achieve the final goal respectively:

- Direct deployment and tuning of FR code [6] on Android: Instead of introducing an additional operating system layer, this method focuses on deploying the Frequency Regularization (FR) algorithm within the Android ecosystem directedly. The core idea of this direction is to optimize and adjust the FR parameters specifically for the Android hardware and software architecture.
- Deployment of Linux environment on Android system:
   It aims to add another layer to the Android system, providing a more controlled and flexible development environment for deploying and testing deep learning models. We choose Termux [7], [8], [9] which is an Android terminal application and Linux environment to deploy.
- Deployment of FR code on Android Studio: This method aims to optimize the FR algorithm for Android's native architecture, ensuring seamless integration and operation within Android. The main idea of this method is similar to the first direction above.

## II. RELATED WORK

Extreme learning machine (ELM) has been widely used in artificial intelligence field over the last decades [1], [10], [11], [12], [13]. Although this algorithm has seen significant development, it also has several drawbacks:

 Poor tunability: It has poor controllability of the model since ELM calculates the least squares solution directly, and users cannot analyze the characteristics of the

- datasets to fine-tune. Adjusting models based on specific performance of mobile devices is important to mobile development.
- Lack of robustness: The performance of the model can be affected significantly while including certain outliers in different datasets, indicating poor robustness. Deployment on mobile devices needs to handle various inputs, including every potential outlier. Although there are many advanced versions of ELM [14], [15], [16], [17] they lacks universality and are not as easy as other algorithms to deploy.
- Overfitting issues: While deploying large convolutional neural networks on mobile devices, model generalization is crucial since overfitting can result in poor performance on unseen data. ELM easily leads to overfitting issues because it is based on empirical risk minimization without considering structural risk. Xue et al. [18]pointed to a regularization strategy to solve this problem by featureselection.

NestDNN is a framework that takes the dynamics of runtime resources into account [2]. The experiment of Fang et al. [2] achieves as much as 4.2% increase in inference accuracy, 2.0x increase in video frame processing rate and 1.7x reduction in energy consumption. However, NestDNN also comes with some limitations. Its computational cost is significantly higher by using filter pruning method Triplet Response Residual (TRR). The high computational cost could probably exceed the processing capabilities of existing mobile devices and the runtime of model generation may be too long, which is not suitable for our deployment.

"One-shot Whole Network Compression" [3] includes removing certain parameters or structures, which is irreversible. Moreover, by using this compression method, the accuracy is too low. For example, in the experiment of Kim et al., by using AlexNet, the accuracy of the compressed model can drop by more than 50%. In order to increase its accuracy, we have to make fine-tuning on the compressed model. Increasing accuracy requires at least more than 10 training epochs, which wastes too much time.

The proposed frequency regularization (FR) [4] works by restricting the non-zero elements of network parameters in the frequency domain, thus reducing information redundancy. Table I illustrates the evaluation of the proposed frequency regularization on UNet, according to compression rate, number of parameters, and dice score. Dice score is a metric for assessing the quality of image segmentation and ranges from 0 to 1, where 0 indicates no overlap and 1 indicates perfect overlap. The data under the dashed line represents the result under the most extreme condition in which only 759 float16 parameters are kept in UNet-v4. Thus, according to the surprised and satisfying experiment outcomes, we choose the frequency regularization as our compression method, to deploy it on mobile devices (i.e. Android system).

TABLE I
EVALUATION OF THE PROPOSED FREQUENCY REGULARIZATION ON UNET
FOR IMAGE SEGMENTATION TASKS USING CARVANA IMAGE MASKING
CHALLENGE DATASET [4], [5].

	Dice Score	Compression Rate	# of Parameters
UNet-ref	99.13%	100%(1×)	31,043,586
UNet-v1	99.51%	1%(100x)	310,964
UNet-v2	99.37%	$0.1\%(1000\times)$	31,096
UNet-v3	98.86%	$0.0094\%(10573\times)$	2,936
ŪNet-v4	97.19%	$\overline{0.0012\%(81801\times)}^{-}$	759(float16)

Termux [7], [8], [9] is an Android terminal application and Linux environment. It works directly with no rooting or setup required. To use Termux, the system needs to meet some requirements: Android 5.0 - 12.0; CPU: AArch64, ARM, i686, x86\_64; at least 300 MB of disk space. It is open source and can be accessed at https://github.com/termux/termux-app. The instruction of deploy Termux on Android devices is available at https://github.com/btxcy/NeuralOnMobile#readme.

#### III. METHODOLOGY

### A. Deployment on Android Devices

To deploy the Frequency Regularization Compression Algorithm on Android devices, we have methodologies: using Termux to implement Linux environment on Android system and developing an Android Application with FR image segmentation by using Android Studio.

- Termux: "why we need termux, use for doing what"
- Android Studio: The reason for applying the algorithm by the Android Studio is because there is less memory consumption during the run. The logic of running on Termux is to create a virtual environment and run an operating system, in our experiment, Ubuntu. Although we are running the compact version of Ubuntu, the limitation is that Ubuntu does not recognize the type of device we are running on. Therefore, it brings up the problem of installing the incapable PyTorch or other libraries. Whereas on an Android application, there is no virtual environment involved.
  - Normal FR on Android Application:
  - 4-parts FR on Android Application:

#### B. Qualitative and Quantitative Metrics

To analyze the results of our experiment, we have three qualitative and quantitative metrices: usage of random access memory, dice score, and visual perception.

- Usage of Random Access Memory (RAM):
- Dice Score:
- Visual Perception:

#### IV. EXPERIMENTS

#### A. Packaging frequency regularization Project

As current and future plans of Zhao et al. [6], [4], we have accomplished the development of a pip repository for their Frequency Regularization technique and committed to their

original repository. We can now integrate Frequency Regularization into our project by simply running the command line in Linux environment:

#### \$ pip install frereg

This step is instrumental in simplifying the deployment of condensed yet potent models in pragmatic applications. We will use this python library in Section IV-B.

## B. Build Linux environment and install the frequency regularization package

We utilized an Android device which ran version 12.0.1 for this section. In order to facilitate the deployment of an Ubuntu environment within the Android system, we downloaded Termux [7], [8], [9], which version is v0.118.0. Upon accessing Termux, we employed a suite of basic tools, such as wget, proot, and git, to establish the Ubuntu environment. The Ubuntu package [19] we used is quite different from the conventional Ubuntu installations on normal personal computers. For more detailed steps of setting up the environment, please check our source code repository [20]. After initiating the Ubuntu environment and installing both Python and the Python-pip tool, we used the command line we mentioned in Section IV-A, installing the frequency regularization project successfully.

#### C. Deployment of Android Studio

As mentioned in Section III-A, the system can install an incapable PyTorch function through the virtual environment with mobile devices. Therefore, running the Python Script on the Android Studio required the python library Chaquopy [21] in this research. It helps us to implement the pip library from Python. However, in some situations, some PyTorch features may still not work when running the code because of the difference in hardware structure between the computer and mobile devices [22]. For instance, the Discrete Fast Fourier Transforms, mobile devices do not support this operation due to the requirement of GPU involvement. Therefore, an alternative solution is to convert the tensor data to the NumPy array, which runs through the CPU. After the transformation, we then convert back to the PyTorch tensor array. Many functions require conversion to make it work in the Frequency Regularization algorithm.

We show the outputs from the algorithm running on an Android phone by making an application using Android Studio to run the Frequency Regularization. The only feature the application provides currently is importing the images from the local machine directory. After that, we get the output from the compressed neural network.

The left-hand column of the IV-C shows the original images we want to implement the image segmentation. From left to right, the three columns represent three outputs generated by an Android phone respectively: output from non-compressed model, output from FR compressed model, output from FR compressed model with 4 parts. Table II illustrates the results

among three ways of image segmentation. analyze the avg. values here!

#### TABLE II

COMPARISON OF RAM USAGE FOR IMAGE SEGMENTATION FOR CARVANA IMAGE MASKING CHALLENGE DATASET [5] ON AN ANDROID PHONE: UNCOMPRESSED U-NET, FR COMPRESSED AND DECOMPRESSED U-NET, AND FR COMPRESSED AND DECOMPRESSED U-NET PROCESSING FOUR PARTS OF AN IMAGE SIMULTANEOUSLY.

	Non-compressed	Compressed	Compressed(4-parts)
1 <sup>st</sup> Avg.	4.2966GB	1.2324GB	1.2553GB
$2^{nd}$ Avg.	4.2966GB	1.2858GB	1.1073GB
$-\frac{3^{rd}}{\text{Total Avg.}}$ Avg	<u>4.2966GB</u> - <u>4.2966GB</u>	<del>xxxxx</del> -	$\frac{1.2546GB}{1.20\overline{5}7\overline{G}\overline{B}}$

TABLE III DICE SCORES

	Non-compressed	Compressed	Compressed(4-parts)
$1^{st}$ Avg. $2^{nd}$ Avg. $3^{rd}$ Avg.	XX	XX	XX
$2^{nd}$ Avg.	XX	XX	XX
$3^{rd}$ Avg.	XX	XX	XX
Total Avg.	XX	XX	XX

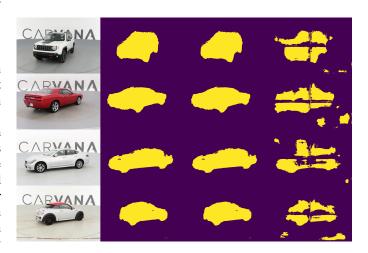


Fig. 1. Comparison between Original Images and Outputs from Image Segmentation for Carvana Image Masking Challenge Dataset [5] on an Android Phone: Uncompressed U-Net, FR Compressed and Decompressed U-Net, and FR Compressed and Decompressed U-Net Processing Four Parts of an Image Simultaneously.

#### V. LIMITATION AND FUTURE WORK

XXXX

#### VI. CONCLUSION

XXXX

#### REFERENCES

[1] A. Akusok, L. Leal, K.-M. Björk, and A. Lendasse, "High-performance elm for memory constrained edge computing devices with metal performance shaders," in *Proceedings of ELM2019*, ser. Proceedings in Adaptation, Learning and Optimization, J. Cao, C. Vong, Y. Miche, and A. Lendasse, Eds. International: Springer, 2021, pp. 79–88.

- [2] B. Fang, X. Zeng, and M. Zhang, "Nestdnn: Resource-aware multi-tenant on-device deep learning for continuous mobile vision," in *Proceedings of the 24th Annual International Conference on Mobile Computing and Networking*, ser. MobiCom '18. ACM, Oct. 2018. [Online]. Available: http://dx.doi.org/10.1145/3241539.3241559
- [3] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin, "Compression of deep convolutional neural networks for fast and low power mobile applications," 2016.
- [4] C. Zhao, G. Dong, S. Zhang, Z. Tan, and A. Basu, "Frequency regularization: Reducing information redundancy in convolutional neural networks," *IEEE Access*, vol. 11, pp. 106793–106802, 2023.
- [5] M. M. M. P. W. C. Brian Shaler, DanGill, "Carvana image masking challenge," 2017. [Online]. Available: https://kaggle.com/competitions/ carvana-image-masking-challenge
- [6] G. Dong, "Frequency regularization," https://github.com/guanfangdong/ pytorch-frequency-regularization, 2023.
- [7] termux, "Termux application," 2023. [Online]. Available: https://github.com/termux/termux-packages
- [8] M. Rana, "Open and accessible education with virtual reality," 2022.[Online]. Available: https://comserv.cs.ut.ee/ati\_thesis/datasheet.php?id=75298.
- [9] termux, "The termux wiki," 2023. [Online]. Available: https://wiki. termux.com/wiki/Main\_Page
- [10] R. N. Shifei Ding, Xinzheng Xu, "Extreme learning machine and its applications," *Neural Computing and Applications*, vol. 25, no. 3, pp. 549–556, 2014. [Online]. Available: https://doi.org/10.1007/s00521-013-1522-8
- [11] S.-H. W. Y.-D. Z. Jian Wang, Siyuan Lu, "A review on extreme learning machine," *Multimedia Tools and Applications*, vol. 81, no. 29, pp. 41611–41660, 2022. [Online]. Available: https://doi.org/10.1007/ s11042-021-11007-7
- [12] C. Deng, G. Huang, J. Xu, and J. Tang, "Extreme learning machines: new trends and applications," *Science China Information Sciences*, vol. 58, no. 2, pp. 1–16, 2015.
- [13] G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, "Extreme learning machine: Theory and applications," *Neurocomputing*, vol. 70, no. 1, pp. 489–501, 2006, neural Networks. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0925231206000385
- [14] J. M. Fossaceca, T. A. Mazzuchi, and S. Sarkani, "Mark-elm: Application of a novel multiple kernel learning framework for improving the robustness of network intrusion detection," *Expert Systems with Applications*, vol. 42, no. 8, pp. 4062–4080, 2015. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0957417414008197
- [15] K. Zhang and M. Luo, "Outlier-robust extreme learning machine for regression problems," *Neurocomputing*, vol. 151, pp. 1519–1527, 2015. [Online]. Available: https://www.sciencedirect.com/science/article/pii/ S0925231214012053
- [16] Q.-Y. Zhu, A. Qin, P. Suganthan, and G.-B. Huang, "Evolutionary extreme learning machine," *Pattern Recognition*, vol. 38, no. 10, pp. 1759–1763, 2005. [Online]. Available: https://www.sciencedirect.com/ science/article/pii/S0031320305001809
- [17] K. Sun, J. Zhang, C. Zhang, and J. Hu, "Generalized extreme learning machine autoencoder and a new deep neural network," *Neurocomputing*, vol. 230, pp. 374–381, 2017. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S092523121631503X
- [18] X. Ying, "An overview of overfitting and its solutions," Journal of Physics: Conference Series, vol. 1168, no. 2, p. 022022, feb 2019. [Online]. Available: https://dx.doi.org/10.1088/1742-6596/1168/ 2/022022
- [19] A. Argentakis, "ubuntu-in-termux," 2023. [Online]. Available: https://github.com/MFDGaming/ubuntu-in-termux
- [20] W. You and L. Chang, "Nerual networks on mobile devices," 2023. [Online]. Available: https://github.com/btxcy/NeuralOnMobile
- [21] Chaquopy, "Chaquopy 14.0 documentation," https://chaquo.com/ chaquopy/doc/current/, 2023, accessed: 2023-12-05.
- [22] R. Fojtik, "New processor architecture and its use in mobile application development," in *Digital Science*, T. Antipova, Ed. Cham: Springer International Publishing, 2022, pp. 545–556.