

---

# Extended Abstract for Research Project

---

**József Konczer**  
Imagination Technologies Ltd  
konczer.j@gmail.com

## Abstract

1       The proposed project is the continuation and generalization of the paper: Self-  
2       Compressing Neural Networks (1). The main contribution of the paper is a novel  
3       quantization-aware training (QAT) method. The approach can be used both to  
4       compress models during training, and to compress pretrained models. The ex-  
5       periments of the original paper were conducted on visual models (classification  
6       on CIFAR10 dataset), which I would extend and generalize to natural language  
7       processing models and tasks.

## 8   1   Introduction

9       Deep neural networks (DNNs) are a powerful tool that have shown unmatched performance in various  
10      tasks in computer vision, natural language processing and optimal control, to mention only a few. The  
11      high computational resource requirements, however, constitute one of the main drawbacks of DNNs,  
12      hindering their massive adoption on edge devices. With the growing number of tasks performed  
13      on edge devices, e.g., smartphones or embedded systems, and the availability of dedicated custom  
14      hardware for DNN inference, the subject of DNN compression has gained popularity. (2)

15     The objective in the paper (1) was threefold: (1) to compress networks during training to realize  
16     benefits in training time; (2) to reduce the size of weight and activation tensors by eliminating  
17     redundant channels; and (3) to reduce the number of bits required to represent weights. The second  
18     and third points produce a smaller network expected to execute more efficiently on devices supporting  
19     variable bit depth weight formats. Despite being conceptually simple, the approach was effective and  
20     the authors demonstrate high compression rates on an example classification network.

21     The main novelty in paper (1) was a differentiable quantization scheme, which is a representation of  
22     model weights  $x$  on a tunable but finite bit depth  $b$  and exponent of a floating-point representation  $e$ :

$$q(x, b, e) = 2^e \left[ \min(\max(2^{-e}x, -2^b), 2^b - 1) \right] \quad (1)$$

23     Where  $[\cdot]$  is the rounding function which rounds to nearest integer with ties to nearest even. The  
24     quantization works for continuous  $b, e$  parameters, which allows to apply standard learning methods.

## 25   2   Related Work

26     Quantization-aware training (QAT) is both a relatively old concept (3), and an active research area (4).  
27     While the proposed work aims to bridge multiple research areas: low bit depth neural networks(5),  
28     QAT(6; 4; 3), and induced sparsity (particularly channel pruning)(7).

29 **3 Visual DNN Results**

30 The main results of the method on visual tasks are demonstrated in figures 1, 2:

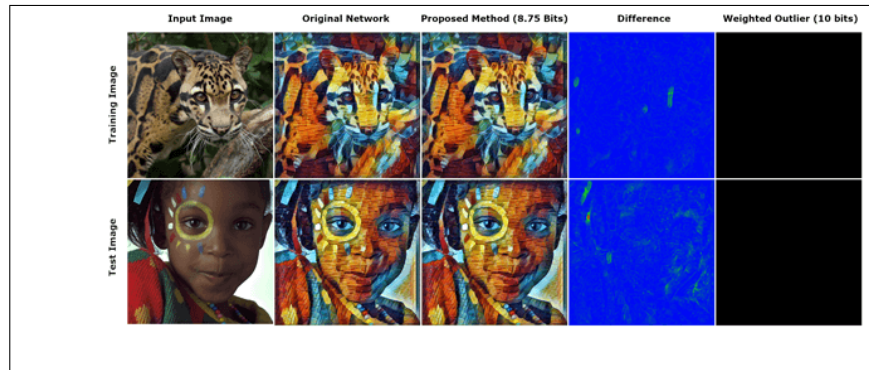
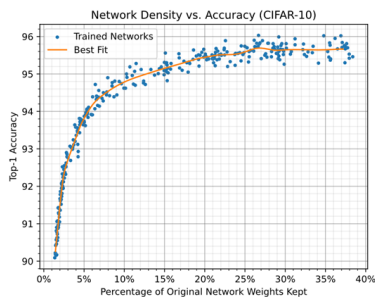
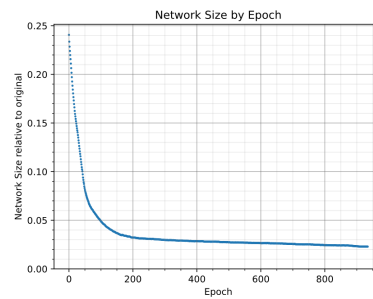


Figure 1: Style transfer results on pretrained model: the results of the original and proposed compressed model are hardly distinguishable. Source: Szabolcs Cséfalvay's blog post



(a) Accuracy-Network size relation



(b) Network size shrinking during training

Figure 2: Compression during training. Source: Szabolcs Cséfalvay's blog post

31 **4 Language Models**

32 The proposed self compressing method is model agnostic, therefore it could be applied to language  
 33 models as well. (An example for an already existing technique is Alpha Tuning (8))

34 **5 Outline of the Proposed Project**

- 35 • **Reproduction:**
  - 36 – Reproducing the results on a ResNet-v2 50(9; 10)(Implemented in PyTorch)
  - 37 – trained and tested on CIFAR 10 (11)
- 38 • **Extension:**
  - 39 – **Possibility I:**
    - 40 \* Self compressed BERT (12) (Implemented in PyTorch)
    - 41 \* trained on Wikipedia data set e.g. WikiText-2
  - 42 – **Possibility II:**
    - 43 \* Self compressed RoBERTa (13) (Implemented in PyTorch)
    - 44 \* trained on Goodreads Dataset (14; 15)

## 45 References

- 46 [1] S. Cséfalvay and J. Imber, “Self-compressing neural networks,” 2023. [Online]. Available:  
47 <https://arxiv.org/abs/2301.13142>
- 48 [2] Y. Nahshan, B. Chmiel, C. Baskin, E. Zheltonozhskii, R. Banner, A. M. Bronstein, and  
49 A. Mendelson, “Loss aware post-training quantization,” *Machine Learning*, vol. 110, no. 11-12,  
50 pp. 3245–3262, Oct. 2021. [Online]. Available: <https://doi.org/10.1007/s10994-021-06053-z>
- 51 [3] A. Gholami, S. Kim, Z. Dong, Z. Yao, M. W. Mahoney, and K. Keutzer, “A survey of  
52 quantization methods for efficient neural network inference,” 2021. [Online]. Available:  
53 <https://arxiv.org/abs/2103.13630>
- 54 [4] A. D’efossez, Y. Adi, and G. Synnaeve, “Differentiable model compression via pseudo quanti-  
55 zation noise,” *ArXiv*, vol. abs/2104.09987, 2021.
- 56 [5] E. Wang, J. J. Davis, P. Y. K. Cheung, and G. A. Constantinides, “Lutnet: Learning fpga  
57 configurations for highly efficient neural network inference,” *IEEE Transactions on Computers*,  
58 vol. 69, no. 12, pp. 1795–1808, 2020.
- 59 [6] R. Gong, X. Liu, S. Jiang, T. Li, P. Hu, J. Lin, F. Yu, and J. Yan, “Differentiable soft  
60 quantization: Bridging full-precision and low-bit neural networks,” 2019. [Online]. Available:  
61 <https://arxiv.org/abs/1908.05033>
- 62 [7] Y. He, X. Zhang, and J. Sun, “Channel pruning for accelerating very deep neural networks,” in  
63 *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- 64 [8] S. J. Kwon, J. Kim, J. Bae, K. M. Yoo, J.-H. Kim, B. Park, B. Kim, J.-W. Ha, N. Sung,  
65 and D. Lee, “Alphatuning: Quantization-aware parameter-efficient adaptation of large-scale  
66 pre-trained language models,” *ArXiv*, vol. abs/2210.03858, 2022.
- 67 [9] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” 2015.  
68 [Online]. Available: <https://arxiv.org/abs/1512.03385>
- 69 [10] —, “Identity mappings in deep residual networks,” in *Computer Vision – ECCV 2016*,  
70 B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing,  
71 2016, pp. 630–645.
- 72 [11] A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny  
73 images,” *Computer Science Department, University of Toronto*, 2009. [Online]. Available:  
74 <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>
- 75 [12] —, “Learning multiple layers of features from tiny images,” *Computer Science  
76 Department, University of Toronto*, 2009. [Online]. Available: [https://www.cs.toronto.edu/  
77 ~kriz/learning-features-2009-TR.pdf](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf)
- 78 [13] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer,  
79 and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” 2019. [Online].  
80 Available: <https://arxiv.org/abs/1907.11692>
- 81 [14] M. Wan and J. J. McAuley, “Item recommendation on monotonic behavior chains,” in  
82 *Proceedings of the 12th ACM Conference on Recommender Systems, RecSys 2018, Vancouver,  
83 BC, Canada, October 2-7, 2018*, S. Pera, M. D. Ekstrand, X. Amatriain, and J. O’Donovan, Eds.  
84 ACM, 2018, pp. 86–94. [Online]. Available: <https://doi.org/10.1145/3240323.3240369>
- 85 [15] M. Wan, R. Misra, N. Nakashole, and J. J. McAuley, “Fine-grained spoiler detection from  
86 large-scale review corpora,” in *Proceedings of the 57th Conference of the Association for  
87 Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long  
88 Papers*, A. Korhonen, D. R. Traum, and L. Màrquez, Eds. Association for Computational  
89 Linguistics, 2019, pp. 2605–2610. [Online]. Available: <https://doi.org/10.18653/v1/p19-1248>