

On the Combination of Tensor Decomposition and Quantization for CNN Compression

Clément LARODIE

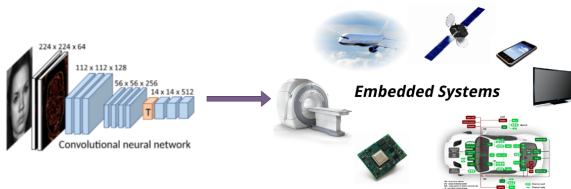
Encadrant: Mohamed Ouerfelli

CEA List/LVML

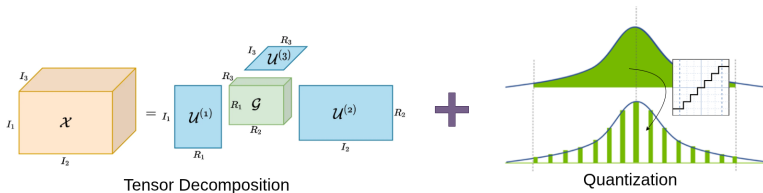
JDSE 2025 – September 25, 2025



CNN Compression



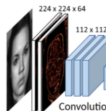
Practical study of an hybrid approach



CNN Compression

On the Combination of Tensor Decomposition and Quantization for CNN Compression

Clément Laroudie (clement.laroudie@cea.fr), Mohamed Oukrif (mohamed.oukri@cea.fr)
CEA LIST, GIGASCALEXCEL



Motivation

- CNNs achieve strong predictive performance but require substantial computation, limiting their deployment on resource-constrained devices such as smartphones or edge AI systems.
- Recently, post-training compression methods can reduce its footprint. Tensor decomposition (TD) and quantization (Q) are two widely used techniques for reducing the computational burden of CNNs.
- Although extensively studied individually, their joint behavior remains understudied. We provide a preliminary empirical study.



Figure 1: The layers of a CNN, where each feature map is a 3-way tensor.

- A convolution is parameterized by a kernel tensor $\mathcal{K} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ mapping an input image $2 \times \mathbb{R}^{I_1 \times I_2}$ to an output image $2 \times \mathbb{R}^{I_1 \times I_2}$.



Figure 2: A 2D convolution.

Pra

Methodology

Study of the two major decompositions:

- CP decomposition [1]: Factoring along the four modes (Fig. 3).
- Tucker decomposition [2]: Factoring along the first two modes only (Fig. 4).



Figure 3: CP decomposition approximates a 3-way tensor \mathcal{X} as a sum of rank-1 tensors — a generalization of matrix-vector tensors.



Figure 4: Tucker decomposition of a 3-way tensor \mathcal{X} as a core tensor \mathcal{G} and factors $\mathcal{U}^{(i)}$.

- **Models:** ResNet101/152, GoogLeNet, AlexNet.
- **Datasets:** CIFAR-10 and ImageNet.
- **Rank Selection:** Both parameter ratio and automatic rank selection via VAMP [3].
- **Quantization:** FP16/INT8 with Pytorch and QNNAP.

Preliminary Results

Tucker2

- Accuracy drops from TD+Q closely match the sum of the individual degradations, suggesting near-independence (Fig. 5).
- The results also generalize with the Signa Data-Aware Tucker Decomposition.
- **Practical pipeline:** Decompose layers with Tucker 2 (keeping later ones) to reach a desired ratio, then apply INT8 quantization to achieve an additional constant reduction.

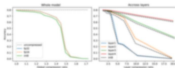


Figure 5: Realistic compression on ImageNet using Tucker 2 decomposition. Left: global compression with all convolutions compressed (except the last). Right: per-layer compression where only one layer is decomposed.

CP4

- CP + INT8 fails (Fig. 6) due to exploding component ranges.
- Corrective methods as EPC [5] (Fig. 7) seem to mitigate this effect (not perfect yet).

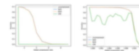


Figure 6: ResNet-101, global CP4 compression. Figure 7: ResNet-101, first layer CP4 with EPC correction.

Current & Future work

- **The CP irregularity:** We will continue investigating stable CP methods as EPC [5].
- **Lower precisions:** We need to investigate lower bit quantization (e.g. INT4).
- **Joint optimization:** We only evaluate post-training pipelines. Quantization-Aware Decomposition [4] and lightweight fine-tuning remain unexplored.
- **Limited to standard CNNs:** Extensions to transformers and attention models are planned.

References

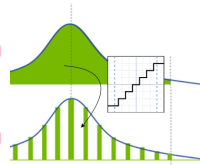
- [1] CP decomposition.
- [2] Tucker decomposition.
- [3] VAMP.
- [4] Quantization-Aware Decomposition.
- [5] EPC.



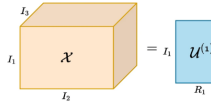
ems



ach



Quantization



Tensor Deci

