

QUANTIZED NEURAL NETWORKS FOR ON-BOARD PROCESSING IN EARTH OBSERVATION BY SATELLITE

Cédric Gernigon¹, Silviu Filip¹, Olivier Sentieys¹, Clément Coggiola²

cedric.gernigon@inria.fr

¹ Univ Rennes, INRIA ² CNES



Data transmission bottleneck

- The increase in the number of high-resolution sensors leads to higher downlink bandwidth constraints
- Deep learning (DL) has made big breakthroughs in image processing



Efficient data transmission

- On-board data processing reduces the amount of transmission
- Most DL algorithms are incompatible with embedded constraints
→ Need for efficient DL

Background

What is Quantization?

- Reduce the arithmetic bit-width (e.g., 32-bit floating point → 8-bit integer)

Why is it interesting?

- Lower power consumption
- Enables the use of FPGAs

Are there any disadvantages?

- Aggressive quantization leads to loss of accuracy

The impact of quantization on accuracy

Post-Training Quantization (PTQ): quantize a pre-trained network → lower accuracy (-), data free (+), low computational cost (+)

Quantization aware training (QAT): train a network with quantized parameters → better accuracy (+), computationally expensive (-)

How to find a suitable bit-width?

Precision tuning during training

Learning the bit-width during training

- Find a uniform bit-width for both weights and activations

Hardware Loss:

$N_t \in \mathbb{R}_+^*$: bit-width, a new parameter to optimize

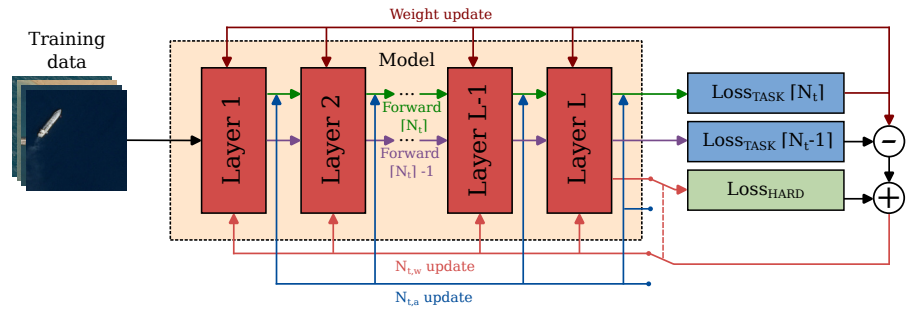
$$L_{\text{Hard}} = \alpha \cdot \sum_{l=1}^{L-1} [N_t^{(l)}]$$

Total loss:

$$L_{\text{Total}} = \lambda \cdot L_{\text{Task}} [N_t] + (1 - \lambda) \cdot L_{\text{Hard}} [N_t]$$

Bit-width gradient approximation:

$$\lambda \cdot (L_{\text{Task}} [N_t] - L_{\text{Task}} [N_t - 1]) + (1 - \lambda) \cdot \frac{\partial L_{\text{Hard}} [N_t]}{\partial [N_t]}$$



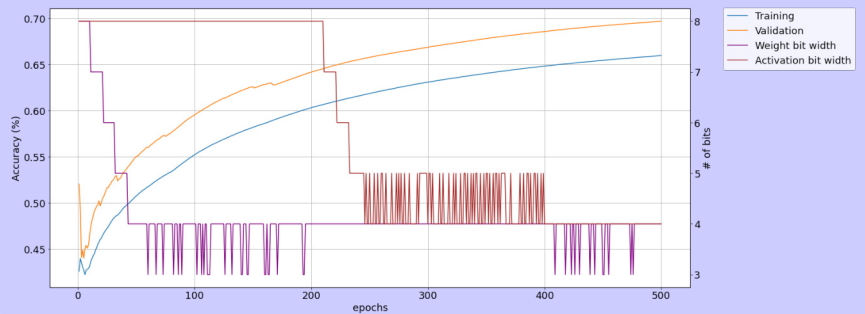
Training settings for image segmentation

Model: U-NET

Dataset: Airbus Ship Kaggle

Arithmetic: Integer

Hyper-parameters: $\lambda = 0.97$, $\alpha = 4$

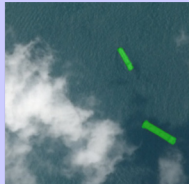


Experimental Results

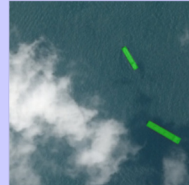
Original image



FP32 model



Expected result



W4A4 model



Accuracy
Size

75.6%
288MB

70.1%
36MB

Conclusion & future works

- Promising initial results for learning the bit-width of weights and activations
- Extend the approach to mixed precision and other arithmetic
- Model a more realistic hardware loss function
- Develop an FPGA-based accelerator