



Instant Quantization of Neural Networks using Monte Carlo Methods

EMC2 Workshop
@ NeurIPS 2019

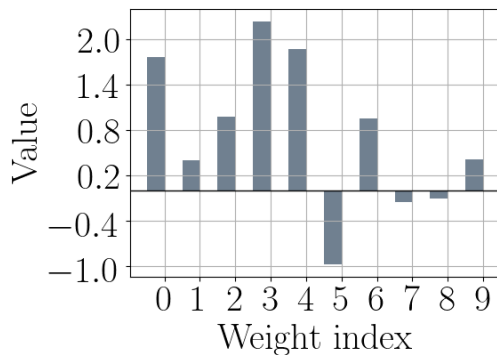
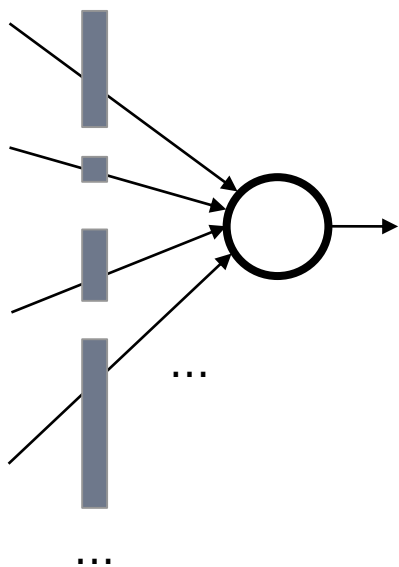
Gonçalo Mordido
Matthijs Van Keirsbilck
Alexander Keller

Hasso Plattner Institute
NVIDIA
NVIDIA

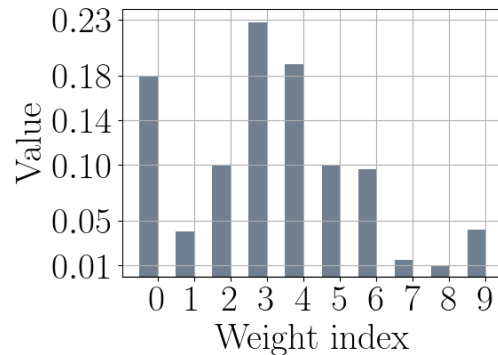
Motivation and idea

- neural network quantization/sparsity
 - lower cost: compute, memory, power, bandwidth, ...
- quantization usually requires retraining
- idea: use importance sampling
 - fast and efficient due to stratified sampling
 - sparsity and bit-width adjustable by the number of samples
 - **no additional training**

Monte Carlo Quantization (MCQ)

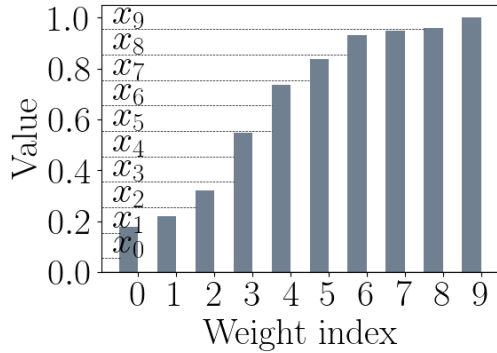


full precision values

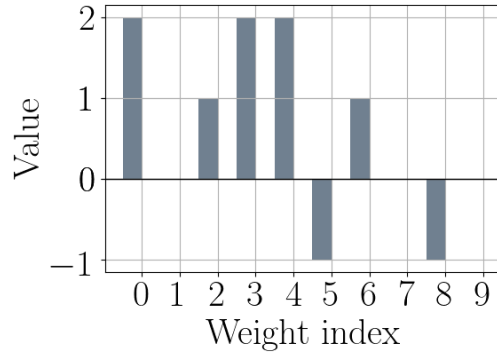


PDF

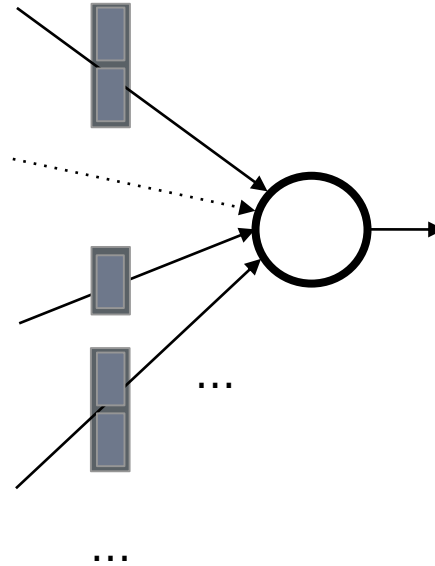
Monte Carlo Quantization (MCQ)



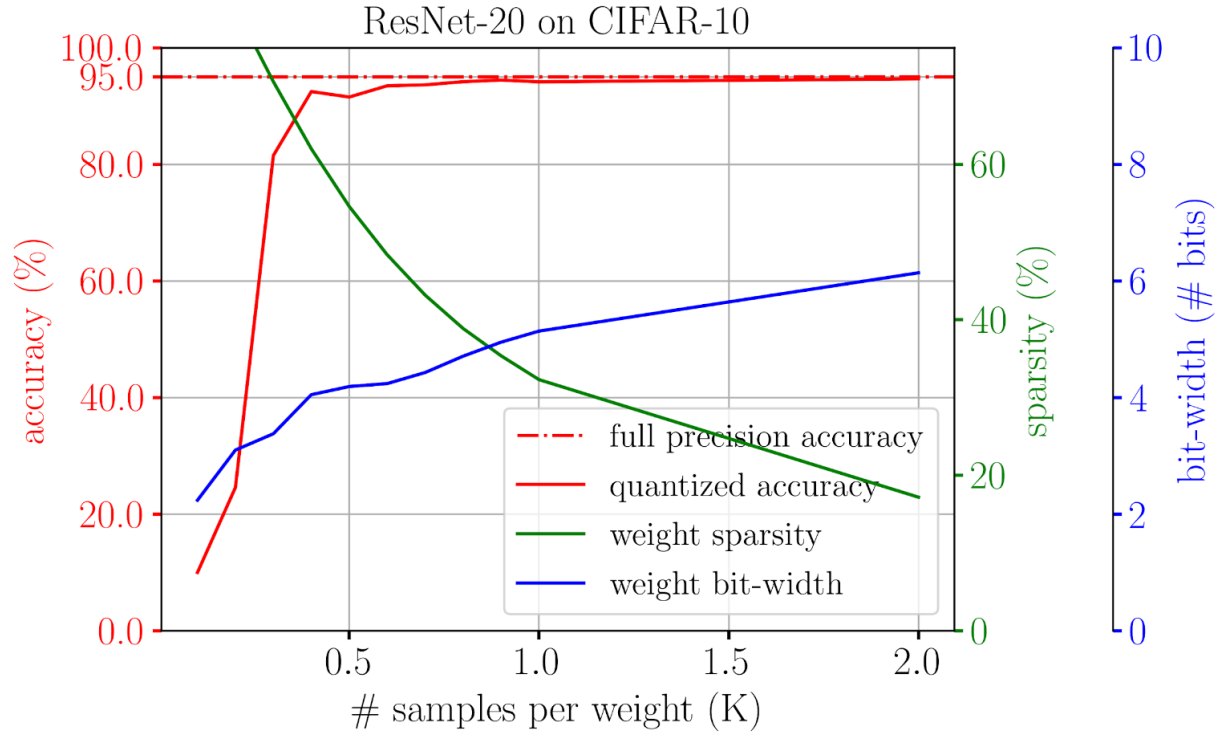
CDF



integer values



Results



Monte Carlo Neural Networks

- simple method to quantize/sparsify models
 - low accuracy loss
 - no retraining
- general applicability
 - weights and/or activations
 - related to random walks
- future work
 - quantized gradients
 - integer neural networks

