



# Systematic Neural Network Quantization and Pruning

Amir Gholami, Z. Yao, Z. Dong, S. Kim, Z. Zheng, E. Tan, Q. Huang, S. Yu, L. Wang, Y. Wang, Michael Mahoney, Kurt Keutzer

University of California Berkeley

Microsoft Research Summit, Oct 2021







### Outline

## >Introduction

- Quantization
- ➢Pruning
- Conclusions and Future Work

#### Al and Memory Wall



© Pallas Group, UCB

4

#### Al and Memory Wall



© Pallas Group, UCB

5

### Outline

>Introduction

## Quantization

## ➢Pruning

Conclusions and Future Work

#### **Quantization Enables Fewer Memory Accesses**

- Memory accesses are the principal cost in both latency and energy
- Lower precision weights in DNN mean each memory access brings more data values
- More data values few accesses overall



Operation:	Energy (pJ)	Relative Energy Cost
8b Add	0.03	
16b Add	0.05	
32b Add	0.1	
16b FP Add	0.4	
32b FP Add	0.9	
8b Multiply	0.2	
32b Multiply	3.1	
16b FP Multiply	1.1	
32b FP Multiply	3.7	
32b SRAM Read (8KB)	5	
32b DRAM Read	640	
[Horowitz, ISSCC 20	)14]	1 10 10 <sup>2</sup> 10 <sup>3</sup> 10 <sup>4</sup>

Image Credit: Sdxcentral, nvidia Table credit: Mark Horowitz

### Simulated Quantization!



- Simulated quantization performs arithmetic in software using FP32 arithmetic, but this can lead to discrepancy when deployed on integer-only hardware
- Also many components such as BN or residual connection are computed using FP32 which is not executable on integer-only hardware

Z. Yao\*, Z. Dong\*, Z. Zheng\*, A. Gholami\*, E. Tan, J. Li, L. Yuan, Q. Huang, Y. Wang, M. W. Mahoney, K. Keutzer, HAWQ-V3: Dyadic Neural Network Quantization in Mixed Precision, ICML'21.

#### Simulated/Fake Quantization Error

 Simulated/Fake quantization can create O(1) error for simple operations such as residual connection. This is because rounding operation is not linear:

 $Int(a+b) \neq Int(a) + Int(b)$ 

Example:

$$Int(2.4 + 1.3) = 4 \neq Int(2.4) + Int(1.3) = 3$$

 This difference is O(1) and will propagate throughout the network, especially for low precision quantization as shown in the figure

Z. Yao\*, Z. Dong\*, Z. Zheng\*, A. Gholami\*, E. Tan, J. Li, L. Yuan, Q. Huang, Y. Wang, M. W. Mahoney, K. Keutzer, HAWQ-V3: Dyadic Neural Network Quantization in Mixed Precision, ICML, 2021.



The normalized difference between fake quantization in PyTorch and corresponding values when deployed in hardware for ResNet50 on ImageNet. As one can see the error becomes quite significant, especially for low precision quantization

### Integer-only Quantization!



 Key idea is to not fold BN during training and only do BN folding at the end, otherwise you will effectively be training a model without BN which is known to be hard to optimize



10

### Integer-only Quantization Works: CV

Method	Int	Uni	BL	Precision	Size	BOPS	Top-1
Baseline	X	✓	73.03	W32A32	13.2	322	73.03
RVQuant (Park et al., 2018) CalibTIB(Hubara et al., 2020)	X X	×	70.10 71.90	W8A8 W8A8	3.3 3.3	20 20	70.29 71.60
HAWQV3	<ul> <li>Image: A start of the start of</li></ul>	<ul> <li>Image: A start of the start of</li></ul>	73.03	W8A8	3.3	20	73.01

#### (a) *MobileNetV2*

#### (c) *InceptionV3*

Method	Int	Uni	BL	Precision	Size	BOPS	Top-1
Baseline	X	1	78.88	W32A32	90.9	5850	78.88
Integer Only (Jacob et al., 2018)	<	<ul> <li></li> <li></li> </ul>	78.30	W8A8	22.7	366	74.20
HAWQV3	~	✓	78.88	W8A8 W8A8	22.7	366	<b>78.76</b>

+4%

Z. Yao\*, Z. Dong\*, Z. Zheng\*, A. Gholami\*, E. Tan, J. Li, L. Yuan, Q. Huang, Y. Wang, M. W. Mahoney, K. Keutzer, HAWQ-V3: Dyadic Neural Network Quantization in Mixed Precision, ICML'21. Online Code: <a href="https://github.com/Zhen-Dong/HAWQ">https://github.com/Zhen-Dong/HAWQ</a>

11

#### Integer-only Quantization Works: Transformers

(a) RoBERTa-Base											
	Int-only	MNLI-m	MNLI-mm	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.
Baseline	×	87.8	87.4	90.4	92.8	94.6	61.2	91.1	90.9	78.0	86.0
I-BERT	1	87.5	87.4	90.2	92.8	95.2	62.5	90.8	91.1	79.4	86.3
Diff		-0.3	0.0	-0.2	0.0	+0.6	+1.3	-0.3	+0.2	+1.4	+0.3
(b) RoBERTa-Large											
	Int-only	MNLI-m	MNLI-mm	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.
Baseline	×	90.0	89.9	92.8	94.1	96.3	68.0	92.2	91.8	86.3	89.0
I-BERT	1	90.4	90.3	93.0	94.5	96.4	69.0	92.2	93.0	87.0	89.5
Diff		+0.4	+0.4	+0.2	+0.4	+0.1	+1.0	0.0	+1.2	+0.7	+0.5

#### > I-BERT is integrated in HuggingFace and soon will also be available in Google's Model Garden

S. Kim\*, A. Gholami\*, Z. Yao\*, M. Mahoney, K. Keutzer, I-BERT: Integer-only BERT Quantization, ICML, 2021. [Online Code] S. Kim, A. Gholami, Z. Yao, A. Nrusimha, B. Zhai, T. Gao, M. Mahoney, K. Keutzer, Q-ASR: Integer-only Zero-shot Quantization for Efficient Speech Recognition arxiv: 2103.16827, 2021. [Online Code]

12

#### What about Sub-INT8 Quantization?

Can we go even further and perform lower precision quantization?



Uniform low precision does not work as it can significantly degrade accuracy => Use mixed-precision

- But how should we set the precision for each kernel?

#### **Hessian Aware Quantization**

This is somewhat similar to the **Jenga** game. We only remove blocks that are not sensitive.

- Only use low precision quantization  $\geq$ for insensitive parameters (flat loss landscape)
- Use high precision quantization for sensitive parameters (sharp loss landscape)

This sensitivity can be calculated through Hessian which quantifies the relative sharpness/flatness of the loss landscape.



Image from UniversityCoop

Dong Z, Yao Z, Arfeen D, Gholami A, Mahoney MW, Keutzer K. Hawq-v2: Hessian aware trace-weighted quantization of neural networks. NeurIPS, 2020. Yu S, Yao Z, Gholami A, Dong Z, Mahoney MW, Keutzer K. Hessian-Aware Pruning and Optimal Neural Implant. WACV, 2022(Accepted).

### HAWQ Also works on NLP



#### Named Entity Recognition Task

Method	w-bits	e-bits	$F_1$	Size
Baseline	32	32	95.00	410.9
Q-BERT	8	8	94.79	102.8
DirectQ	4	8	89.86	62.2
Q-BERT	4	8	94.90	62.2
DirectQ	3	8	84.92	52.1
Q-BERT	3	8	94.78	52.1
Q-BERT <sub>MP</sub>	2/4 мр	8	94.55	52.1
DirectQ	2	8	54.50	42.0
Q-BERT	2	8	91.06	42.0
$Q\text{-}BERT_{\text{MP}}$	2/3 мр	8	94.37	45.0

S. Sheng, Z. Dong, J. Ye, L. Ma, Z. Yao, A. Gholami, M. Mahoney, K. Keutzer, Q-BERT: Hessian Based Ultra Low Precision Quantization of BERT, AAAI, 2020.

### Outline

>Introduction

➢Quantization

## ➢Pruning

Conclusions and Future Work

### Pruning is Quantization with 0 bits

• Pruning is a special case of quantization with 0 bits, and all the previous results apply to pruning as well



LeCun Y, Denker J, Solla S. Optimal brain damage. Advances in neural information processing systems. 1989;2:598-605.	
Yu S*, Yao Z*, Gholami A*, Dong Z*, Kim S, Mahoney MW, Keutzer K. Hessian-Aware Pruning and Optimal Neural Implant. WACV, 2022 (accepted).	17

### Neural Implant



- Completely removing some of the neuron results in irrecoverable accuracy loss
- => We can insert a small low dimensional implant instead of complete removal

Yu S\*, Yao Z\*, Gholami A\*, Dong Z\*, Kim S, Mahoney MW, Keutzer K. Hessian-Aware Pruning and Optimal Neural Implant. WACV, 2022 (accepted).

#### HAP with Neural Implant Achieves SOTA

• HAP with Neural Implant consistently achieves higher performance



Pruning results on Cifar-10. Y-axis is accuracy, and x-axis is remaining parameters after pruning.

Yu S*, Yao Z*, Gholami A*, Dong Z*, Kim S, Mahoney MW, Keutzer K. Hessian-Aware Pruning and Optimal Neural Implant. WACV 2022 (accepted).	
Code: https://github.com/yaozhewei/HAP	19

#### Summary: One Size Does Not Fit All

- Quantization and Pruning of a model depends on the loss landscape which is a function of:
  - Model Architecture
  - Dataset
- We can quantify the sensitivity to quantization and pruning efficiently using Hessian



Less Robust to Quantization



More Robust to Quantization

### Widely Adopted by Industry

• Incorporated in Intel OpenVino Library and used for low precision quantization on FPGAs







## Thank You for Listening! amirgh@berkeley.edu

## Slides will be available on amirgholami.org

Hessian Code: HAWQV3 Code: I-BERT Code: Neural Implant Code: https://github.com/amirgholami/PyHessian https://github.com/Zhen-Dong/HAWQ https://github.com/kssteven418/I-BERT https://github.com/yaozhewei/HAP





