#### Jetfire : Efficient and Accurate Transformer Pretrining with INT8 Data Flow and Per-Block Quantization

Haocheng Xi, Yuxiang Chen, Kang Zhao, KAI JUN TEH, Jianfei Chen, Jun Zhu

Presenter : Taehyeon Kim, Yongha Shin, Sanghyeon Cho

# Contents

- Introduction
- Related Work
- INT8 Data Flow
- Per Block Quantization
- Linear Layer Operator
- Non-Linear Operator
- Experiments
- Conclusion

• Large-scale pre-trained transformer-based models e.g. GPT-4, LLAMA, PaLM



• Pre-training transformers require numerous **computations** & high-bandwidth **memory** e.g. sub-layer of FP16 FFN



Pre-training transformers require numerous computations & high-bandwidth memory
e.g. sub-layer of FP16 FFN



Pre-training transformers require numerous computations & high-bandwidth memory
e.g. sub-layer of FP16 FFN



Pre-training transformers require numerous computations & high-bandwidth memory
e.g. sub-layer of FP16 FFN



- Fully Quantized Training (FQT) with INT8
  - Computation : FLOP  $\Rightarrow$  Integer operation
  - Memory bandwidth reduce



- Post-Training Quantization
  - No speedup for training
  - Quantization error for output

#### • Quantization-Aware Training

- $\circ~$  focus on training accuracy, not speedup
- Weight still full-precision



- Fully Quantized Training : INT8 training for CNNs
- Direction Sensitive Gradient Clipping
- Deviation Counteractive Learning Rate Scaling



F. Zhu, arxiv:1912.12607v1

• Gradient Vectorized Quantization

$$\widehat{G}_{W} = \left[\widehat{G}_{W_{1}}, \cdots \widehat{G}_{W_{C_{out}}}\right]^{T}, \widehat{G}_{W_{i}} = q(G_{W_{i}}) \cdot \frac{S_{x}}{127} \cdot \frac{S_{i}}{127}$$

• Magnitude-aware Clipping Strategy





K. Zhao, "Distribution Adaptive INT8 Quantization for training CNNs"

#### • Fully Quantized Training : SwitchBack



• Fully Quantized Training : FP8 training with Hopper architecture



- Limitations on existing FQT
- 1. Previous FQT methods designed for CNN, low accuracy for Transformer
- 2. Most FQT focus on computation reduction, not data access overhead
- 3. Some FQT techniques (FP8) requires specialized hardware



• Previous







## **3 INT8 Data Flow**

- Reduce computation and amount of memory access
- Accelerate nonlinear operators
- Activation memory consumption and amount of communication reduced





- 2) Linear Layer Operator
- 3) Non-Linear Layer Operator



Per Channel Quantization is advantageous for activations

due to the presence of channel-wise outliers



Layer Norm

Linear

Per Token Quantization is advantageous for gradient due to the presence of token-wise outliers



INT8 numerical format must accurately support the following three MMs of a linear layer

$$\mathbf{Y} = \mathbf{X}\mathbf{W}^{\mathrm{T}}$$
$$\nabla_{\mathbf{X}} = \nabla_{\mathbf{Y}}\mathbf{W}$$
$$\nabla_{\mathbf{W}} = \nabla_{\mathbf{Y}}^{\mathrm{T}}\mathbf{X}$$



But, using per-channel for forward pass and per-token quantization for backward pass is not hardwarefriendly



 $S_{\mathbf{X}\mathbf{k}}$  prevents the utilization of high-speed matrix multiplication (MM) in hardware

#### **Per-block quantization** can achieve computational efficiency and preserve accuracy at the same time



Per-block quantization achieved an error rate that falls between that of per-channel and per-token quantization

Final layer activation quant error

17

## 4 Method

- 1) Per Block Quantization
- 2) Linear Layer Operator
- 3) Non-Linear Layer Operator



### 4 Method Linear Layer Operator

Jetfire uses **3-Level Tiling of Matrix Multiplication** based on the GPU architecture

> Level 1: CUDA thread block Level 2: Quantization block Level 3: WMMA operation

## 4 Method Linear Layer Operator

Level 1: CUDA thread block

 $\mathbf{Y} = \mathbf{X}\mathbf{W}^{\mathrm{T}}$ , where  $\mathbf{X} \in \mathbb{R}^{N \times C}$ ,  $\mathbf{W} \in \mathbb{R}^{D \times C}$ ,  $\mathbf{Y} \in \mathbb{R}^{N \times D}$ ,





#### Level 2: Quantization block





#### Level 3: WMMA operation



$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \cdot \begin{bmatrix} A' & B' \\ C' & D' \end{bmatrix} = \begin{bmatrix} AA' + BC' & AB' + BD' \\ A'C + C'D & B'C + DD' \end{bmatrix}$$

Block 1 Shared  $\mathbf{W}_{jk}$ memory  $\mathbf{X}_{ik}$ Register  $\mathbf{X}_{ik,p}$   $\mathbf{W}_{ik,q}$ ... Thread Global memory X W

#### GPU memory hierarchy

### 4 Method Linear Layer Operator



### 4 Method

- 1) Per Block Quantization
- 2) Linear Layer Operator
- 3) Non-Linear Layer Operator



#### 4 Method Non-Linear Layer Operator

Authors observed that non-linear operation are memorybounded





Precision differences in global memory read and write operations on GLEU

### 4 Method Non-Linear Layer Operator

GELU, LayerNorm, Dropout, Add



All FP32 values are stored only in shared memory due to kernel fusion, while global memory handles only INT8 read/write operations

# 5 Experimental Settings

- Task
  - Machine translation, image classification, generative model pretraining
- Quantization
  - INT 8: Linear layers of MLP, attention, non-linear layers (GELU, LayerNorm, Dropout)
  - FP 16: Multi-head attention (used FlashAttention[1])
  - FP 32: Master copy of the weights
- Comparison method
  - FP16: Floating point training baseline
  - Per-tensor quantization
  - SwitchBack

[1] Dao, T., Fu, D., Ermon, S., Rudra, A., and R´e, C. Flashattention: Fast and memory-efficient exact attention with io-awareness. Advances in Neural Information Process- ing Systems, 35:16344–16359, 2022.

#### Converged model accuracy

- Machine translation transformer
  - Train a transformer base model on WMT 14 En-De dataset

*Table 4.* Results on machine translation, deit pretraining, GPT2 pretraining, and GLUE fine-tuning result based on the pretrained model. FP refers to floating-point, SwitchBack refers to per-token quantization. '-' means the model does not converge.

			BASELINE			OURS
MODEL	#PARAMS(M)	METRIC	FP	<b>SWITCHBACK</b>	PER-TENSOR	JETFIRE
TRANSFORMER-BASE	61	BLEU	26.49	26.46	26.04	26.49
DEIT-TINY	5		64.08	63.55	_	63.95
DEIT-SMALL	22	TOP1 ACC	73.43	72.80	_	73.31
DEIT-BASE	86		75.67	75.62	_	76.03
GPT2-BASE	124		2.9074	3.0796	3.1638	2.8597
GPT2-MEDIUM	350	VALID LOSS	2.6612	2.9141	3.1795	2.4195
GPT2-LARGE	774		2.5993	3.0512	2.9775	2.4696
GPT2-BASE	124		$78.50_{0.45}$	$78.15_{0.15}$	$76.33_{0.94}$	$78.18_{0.34}$
GPT2-MEDIUM	350	GLUE SCORE	$82.00_{0.50}$	$80.03_{0.15}$	$79.19_{0.22}$	$81.60_{0.26}$
GPT2-LARGE	774		$83.01_{0.24}$	$78.74_{0.24}$	$75.88_{0.35}$	$82.94_{0.70}$

#### Converged model accuracy

- Image classification Deit
  - Pretraining for Deit (Tiny, Small, Base) model on ImageNet1K

*Table 4.* Results on machine translation, deit pretraining, GPT2 pretraining, and GLUE fine-tuning result based on the pretrained model. FP refers to floating-point, SwitchBack refers to per-token quantization. '-' means the model does not converge.

			BASELINE			OURS
MODEL	#PARAMS(M)	METRIC	FP	SWITCHBACK	PER-TENSOR	JETFIRE
TRANSFORMER-BASE	61	BLEU	26.49	26.46	26.04	26.49
DEIT-TINY	5		64.08	63.55	_	63.95
DEIT-SMALL	22	TOP1 ACC	73.43	72.80	_	73.31
DEIT-BASE	86		75.67	75.62	_	76.03
GPT2-BASE	124		2.9074	3.0796	3.1638	2.8597
GPT2-MEDIUM	350	VALID LOSS	2.6612	2.9141	3.1795	2.4195
GPT2-LARGE	774		2.5993	3.0512	2.9775	2.4696
GPT2-BASE	124		$78.50_{0.45}$	$78.15_{0.15}$	$76.33_{0.94}$	$78.18_{0.34}$
GPT2-MEDIUM	350	GLUE SCORE	$82.00_{0.50}$	$80.03_{0.15}$	$79.19_{0.22}$	$81.60_{0.26}$
GPT2-LARGE	774		$83.01_{0.24}$	$78.74_{0.24}$	$75.88_{0.35}$	$82.94_{0.70}$

#### Converged model accuracy

- Image classification Swin transformer and ViT
  - Pretraining for Swin (Tiny, Small, Base) model on ImageNet1K
  - Fine-tune ViT (Base, Large) model on ImageNet1K

MODEL	Swin-tiny	SWIN-SMALL	Swin-base
FP	77.55	80.39	80.45
Jetfire	77.51	80.39	80.37
	VIT-BASE	VIT-LARGE	
FP	83.45	85.72	
Jetfire	83.48	85.67	

Table 5. Comparison of FP and Jetfire

#### Converged model accuracy

- Generative model pretraining- GPT2
  - Training for GPT2 (base, medium, large) model on OpenWebText

*Table 4.* Results on machine translation, deit pretraining, GPT2 pretraining, and GLUE fine-tuning result based on the pretrained model. FP refers to floating-point, SwitchBack refers to per-token quantization. '-' means the model does not converge.

			BASELINE			OURS
MODEL	#PARAMS(M)	METRIC	FP	SWITCHBACK	PER-TENSOR	JETFIRE
TRANSFORMER-BASE	61	BLEU	26.49	26.46	26.04	26.49
DEIT-TINY	5		64.08	63.55	_	63.95
DEIT-SMALL	22	TOP1 ACC	73.43	72.80	_	73.31
DEIT-BASE	86		75.67	75.62	_	76.03
GPT2-BASE	124		2.9074	3.0796	3.1638	2.8597
GPT2-MEDIUM	350	VALID LOSS	2.6612	2.9141	3.1795	2.4195
GPT2-LARGE	774		2.5993	3.0512	2.9775	2.4696
GPT2-BASE	124		$78.50_{0.45}$	$78.15_{0.15}$	$76.33_{0.94}$	$78.18_{0.34}$
GPT2-MEDIUM	350	GLUE SCORE	$82.00_{0.50}$	$80.03_{0.15}$	$79.19_{0.22}$	$81.60_{0.26}$
GPT2-LARGE	774		$83.01_{0.24}$	$78.74_{0.24}$	$75.88_{0.35}$	$82.94_{0.70}$

# 7 Ablation Study

#### CUDA kernel and Triton kernel block size

- Block size for Triton and CUDA kernels is crucial
  - Large: Decrease in parallelism
  - Small: Low utilization of bandwidth and computational resources
- Block size of Triton:  $64 \times 64$
- Block size of CUDA:  $128 \times 32 \times 128$



*Figure 6.* Speed test of GELU and GEMM operator. (a) Triton kernel speedup with different Triton block sizes.(b) GEMM CUDA kernel speed with different CUDA block sizes.

#### **Operator and End-to-End Experiments** 8

#### Linear layer & non-linear operator speedup



Figure 8. Matrix Multiplication Speed test for Different Methods in different settings (B=Batch Size, N=Sequence Length).



Figure 7. Speed comparision between our INT8 non-linear operator and pytorch FP16 implementation.

# Operator and End-to-End Experiments

#### End-to-end speedup

• Show the end-to-end speedup for Jetfire method over PyTorch's FP16

*Table 6.* Acceleration ratios for End-to-end comparison (SB refers to SwitchBack basic version) on GPT2 model.

	Forward		BACKWARD		OVERALL	
HIDDEN SIZE	SB	OURS	SB	OURS	SB	OURS
4096	1.50	1.32	1.18	1.46	1.27	1.42
2048	1.53	1.29	1.24	1.41	1.32	1.37
1024	0.94	0.97	1.14	1.11	1.07	1.07

## Sector Conclusion

- Jetfire proves that full-pipeline INT8 pre-training can match FP baselines without accuracy loss
- By storing activations, weights, and gradients in INT8, it halves computation, memory traffic, and GPU usage

## 10 Limitation

- Multi-Head Attention remains in FP16 via the FlashAttention kernel, so Jetfire is not a fully INT8 end-to-end pipeline
- The experiments lack INT4 Transformer[2] baselines
- The per-block quantization scale matrix is large, leading to increased communication overhead for scale values.
- 4-bit per block quantization may not be accurately performed in a Transformer.

[2] Xi, H., Li, C., Chen, J., & Zhu, J. (2023). Training transformers with 4-bit integers. Advances in Neural Information Processing Systems, 36, 49146-49168.

# **Thank You!**

# Appendix

Table 2. Meaning of Key Constants.				
$B_N/B_C/B_D$	CUDA block size in MM			
$T_N/T_C/T_D$	Number of CUDA blocks along each axis			
$\overline{B}$	Quantization block size			
$\overline{R_N/R_C/R_D}$	Number of quantization blocks			
	in a CUDA block along each axis			

Table 3. Ti	me complexity	of different	operations	in	MM.
-------------	---------------	--------------	------------	----	-----

		Method	
OPERATION	BASIC INT8	SWITCHBACK	OURS
MM	$B_N B_D C$	$B_N B_D C$	$B_N B_D C$
16-bit Load/Store	$(B_N + B_D)C + B_N B_D$	$\frac{B_N B_C}{T_D} + \frac{B_D B_C}{T_N} + B_N B_D$	-
8-bit Load/Store	-	$(B_N + B_D)C$	$(B_N + B_D)C + B_N B_D$
DEQUANTIZE	-	$B_N B_D$	$B_N B_D T_C$
QUANTIZE	_	$\frac{B_N B_C}{T_D} + \frac{B_D B_C}{T_N}$	$B_N B_D$

# Appendix

#### Algorithm 2 INT8 Non-Linear Operator

**Require:** INT8 Matrix  $\mathbf{X} \in \mathbb{R}^{N \times C}$ , FP16 scale matrix  $\mathbf{S}_{\mathbf{X}} \in \mathbb{R}^{L_N \times L_C}$ , element-wise function f 1: Define  $T_N = \left| \frac{N}{B_N} \right|, T_C = \left| \frac{C}{B_C} \right|$ 2: Define  $R_N = \left\lceil \frac{B_N}{B} \right\rceil, R_C = \left\lceil \frac{B_C}{B} \right\rceil$ 3: for  $1 < i < T_N$  do for  $1 \leq j \leq T_C$  do 4: Load INT8 block  $\mathbf{X}_{ij} \in \mathbb{R}^{B_N \times B_C}$ ,  $\mathbf{S}_{\mathbf{X}_{ij}} \in \mathbb{R}^{R_N \times R_C}$ 5: Dequantize  $X_{ij}$  and  $S_{X_{ij}}$  to get  $X_{ij}^{FP32}$ 6: Operate:  $\mathbf{Y}_{ii}^{\mathbf{FP32}} = f(\mathbf{X}_{ii}^{\mathbf{FP32}})$ 7: Quantize  $\mathbf{Y}_{ij}^{\mathbf{FP32}}$  to get  $\mathbf{Y}_{ij}^{\mathbf{INT32}} \in \mathbb{R}^{B_N \times B_C}$  and scale factor  $\mathbf{S}_{\mathbf{Y}_{ij}} \in \mathbb{R}^{R_N \times R_C}$ 8: Save  $\mathbf{Y}_{ij}^{\mathbf{INT32}}$  and  $\mathbf{S}_{\mathbf{Y}_{ij}}$  to global memory. 9: 10: end for 11: **end for**