

QServe: W4A8KV4 Quantization and System **Co-design for Efficient LLM Serving**

Yujun Lin*, Haotian Tang*, Shang Yang*, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, Song Han MIT, NVIDIA, MIT-IBM Watson AI Lab, UMass Amherst. Accepted by MLSys 2025.

* indicates the equal contribution to this work.



Current status of LLM serving systems GEMM and Attention together dominate the LLM inference latency.



QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving



2

QoQ: W4A8KV4 Quantization QoQ stands for *quattuor-octo-quattuor*, which represents 4-8-4 in Latin.

- QoQ quantization uses \bullet
 - 8-bit activations to improve peak performance compared to FP16 activations

 - 4-bit weights to save memory bandwidth compared to INT8 weights \bullet • 4-bit KV Cache to save both storage and bandwidth compared to INT8 KV Cache.





Current status of LLM serving systems

W4A4 quantization leads to severe accuracy loss and even worse efficiency.

- Despite the existence of aggressive W4A4 quantization algorithms, they
 - Bring about significant accuracy loss;
 - Cannot run efficiently on current GPUs.

Wikitext2	Algorithm	Llama-2			
Perplexity ↓	Aigontinin	7B	70B		
W8A8	SmoothQuant	5.54	3.36		
W4A4	QuaRot	6.19	3.83		
W4A4 g128	Atom	6.12	3.73		
W4A8KV4	RTN	6.51	3.90		
W4A8KV4 g128	RTN	5.99	3.70		







QoQ: W4A8KV4 Quantization Algorithm for LLM Cloud Serving

Progressive Group Quantization - W4A8 Per-channel INT8 quantization followed by per-group INT4 quantization







Progressive Group Quantization - W4A8 Protective range [-119, 119] avoids dequantization overflow. We reconsider the dequantization process.



QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving

$$\hat{w}_{s8} = \left\lceil \frac{w_{s8}}{s_{u8}} \right\rfloor \cdot s_{u8} \le w_{s4} + \frac{1}{2} s_{u8}$$

$$\frac{q_{\text{max}} - w_{\text{s8min}}}{q_{\text{u4max}} - q_{\text{u4min}}} \le \frac{q_{\text{s8max}} - q_{\text{s8min}}}{q_{\text{u4max}} - q_{\text{u4min}}} = \frac{127 - (-127)}{15 - 0} = \frac{127 - (-127)}{15 - 0}$$

$$\hat{w}_{s8} \le 127$$

we only need

$$w_{s8} + \frac{1}{2}s_{u8} \le 127$$

which is

$$w_{s8} \le 127 - \frac{1}{2}s_{u8} \le 127 - \frac{1}{2} \cdot 17 = 119.5$$



SmoothAttention - KV4 Migrate quantization difficulty from K cache to Q matrix



QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving

(Original)

(SmoothAttention)

 $\mathbf{Z} = (\mathbf{Q}\mathbf{\Lambda}) \cdot (\mathbf{K}\mathbf{\Lambda}^{-1})^{T}, \quad \mathbf{\Lambda} = \operatorname{diag}(\lambda)$

Both Q and K are activations.



8

SmoothAttention - KV4 Migrate quantization difficulty from K cache to Q matrix



QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving

Both Q and K are activations.

$$\max\left(\left|\mathbf{K}_{i}\right|\right)^{\alpha}$$









QServe: Efficient LLM Serving System with W4A8KV4 Quantization

QServe System Overview







Quantized GEMM on GPUs Overhead in the GEMM main loop should be avoided.



(b) TensorRT-LLM (INT4 Weights and FP16 Activations)

QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving



(d) Ours (INT4 Weights and INT8 Activations)







Compute-Aware Weight Reorder Reducing pointer arithmetics overhead through offline weight reordering



(a) The Idmatrix instruction ensures that each thread gets what it needs for compute in W8A8 GEMM

H									<u>S</u>
—T0 →	TO	T1	T2	T3	TO	T1	T2	Т3	nne
—T1 →	T4	T5	T6	T7	T4	T5	T6	T7	Cha
<i>−…→</i>	T8	Т9	T10	T11	Т8	Т9	T10	T11	nt (
$-T7 \rightarrow$	T12	T13	T14	T15	T12	T13	T14	T15	utp
► 4xINT4					e Tile	s obtained by	y T0 from Ic	amatrix	0







(b) However, the Idmatrix instruction fails for W4A8 GEMM due to storage-compute mismatch



	— 19 →	T4	15	16	(С С
	<i>−…→</i>					U tr
	− T15 →	T28	T29	T30	T31	utpi
ipu	te Tile	es obtained	by T0 from	Idmatrix	abt ko	0

put Ch	anne	els —				<u>S</u>		
Т3	TO		T1	T2	Т3	nne		
T7	T4		T4		T5	T6	T7	Cha
T11	Т8		T8 T9		T10	T11	ut (
T15	Т	12	T13	T14	T15	utp		
compute Tiles obtained by T0 from Idmatrix								

(c) Our solution: **compute-aware weight reorder**



Efficient Dequantization with Reg-Level Parallelism Step 1. Dequantization from UINT4 to UINT8



QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving

Reorder every 32 weights to **minimize logical instructions**: 3 instructions to dequantize 32 numbers.





Efficient Dequantization with Reg-Level Parallelism Step 2. Dequantization from UINT8 to SINT8 (apply zero points and scales)



QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving

Subtraction after multiplication is more efficient than subtraction before multiplication.



16

Making Fused Attention Memory Bound Theoretical computation intensity = 1 does not hold in fused implementation

Technique

TRTLLM-KV8

Naive KV4-g128

+Efficient INT4-to-FP16 dequantization

+Control flow simplification

+QK product in half2

+Prefetch scales/zeros

+Tuning tiling parameters

+SV product in half2

	LLM-KV8	Speedup over TRT	Latency (ms)
		1.0x	0.42
	{	0.88x	0.48
quantization c speed up de		0.95x	0.44
stage atter automatic		1.08x	0.39
QServe ba		1.17x	0.36
brings fused back to the		1.31x	0.32
realizes 1.5x		1.40x	0.30
		1.50x	0.28







End-to-end Throughput Evaluation Results Up to 2.4x-3.5x faster than TensorRT-LLM on A100, L40S



	3.07						
1.01							



QServe outperforms the state-of-the-art W4A4

Wikitext2 Perplexity ↓			Llama-2		Llama			Mistral	Mixtral	Yi
Precision	Algorithm	7B	13B	70B	7B	13B	30B	7B	8x7B	34B
W8A8	SmoothQuant	5.54	4.95	3.36	5.73	5.13	4.23	5.29	3.89	4.69
W4A16 g128	AWQ	5.60 (+0.06)	4.97 (+0.02)	3.41 (+0.05)	5.78 (+0.05)	5.19 (+0.06)	4.21 (-0.02)	5.37 (+0.08)	4.02 (+0.13)	4.67 (-0.02
W4A4	QuaRot	6.19 (+0.65)	5.45 (+0.50)	3.83 (+ 0.47)	6.34 (+0.61)	5.58 (+ 0.45)	4.64 (+0.41)	5.77 (+ 0.48)	NaN	NaN
W4A4 g128	Atom	6.12 (+0.58)	5.31 (+0.36)	3.73 (+0.37)	6.25 (+0.52)	5.52 (+0.39)	4.61 (+0.38)	5.76 (+ 0.47)	4.48 (+0.59)	4.97 (+0.28
W4A8KV4	QServe	5.75 (+0.18)	5.12 (+0.17)	3.52 (+0.16)	5.93 (+0.20)	5.28 (+0.15)	4.34 (+0.11)	5.45 (+0.16)	4.18 (+0.29)	4.74 (+0.05
W4A8KV4 g128	QServe	5.70 (+0.13)	5.08 (+0.13)	3.47 (+0.11)	5.89 (+0.16)	5.25 (+0.12)	4.28 (+0.05)	5.42 (+0.13)	4.14 (+0.25)	4.76 (+0.07







Next Step: Long-context Serving with Sparse Attention LServe: Efficient Long-sequence LLM Serving with Unified Sparse Attention

- Unified static and dynamic sparse attention
- by 1.3-2.1x over vLLM



LServe: Efficient Long-sequence LLM Serving with Unified Sparse Attention [Yang et al, MLSys 2025]





QoQ Algorithm

QServe System

- Compute-aware Weight Reordering
- Efficient Dequantization (Mul \rightarrow Sub)
- Make Fused Attention Memory-bound

Q

()

TRT-LLM

on A100





