

11/8

Quantization

sign
 (-1) \times (1 + fraction)
 Exponent \rightarrow
 $\times 2$

message

\rightarrow It turns out that FP8 representation has many advantages over INT8 for inference, but we need to carefully tune.

$x^{(q)} = s \cdot \left\lfloor \frac{x}{s} \right\rfloor$: express the quantized input.

Scale factor s

HRNet, DeepLab V3, ViT,

BERT-base, ResNet18

E4M3, E5M2

E4M3
 NaN = 5.1111111₂

GPT3, int8()

Model = GPT2, PPL = 33.5, outliers
 = countil, 1-sided = 1. (\downarrow)

+ Iteration 비구하기

11/8 reproduce challenge 한거

① LDPC (64 / 120) 시작

BER 7844 sum product

algorithm / min-sum algorithm

과 비교. / 43 331 (MSA, SA)

③ FLOPS

② sparsity 비율

계산가능?

Self-attention

Model Size vs Quantization

\rightarrow quantized model의 성능
 = model 커질수록 (tend to
 degrade as model get
 larger) 성능 \downarrow .

\therefore why? emergence of outliers.
 Emergence of 'outliers
 features.'

$$1.01 = 1 + \frac{1}{2} \times 0 + \frac{1}{2^2} \times 1$$

$$= \frac{5}{4}$$

Quartiles

= (-8, -7, -6) : GPT2

w/outliers

with

LLM.int8()

- ① 8-bit vector-wise quantization
- ② 16-bit decomposition.

$$\langle 2 \rangle \|W \cdot X - \hat{W} \cdot X\|_F^2$$

layer by layer
optimization.

3) zero quant = 3번째 주시

KD = knowledge distillation.

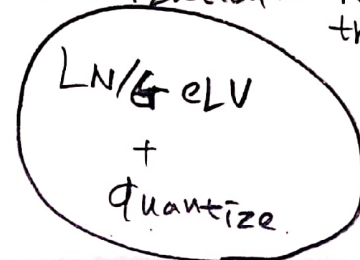
↳ low precision quantization,
KD requires you to load both teacher and student.

- ①
- ↳ 이전) Not really consider GPT
- ↳ state models (decoder-only).
- ② Focusing only on high-precision quantization. (FP16/INT8)

quantization - optimized transformer kernels.

③ KD = knowledge distillation.

Who : teacher model.



1) fine-grained hardware friendly quantization.

↳ INT8 cannot fully capture different numerical ranges.

(multiplier precision)
real(2) / 16-bit.

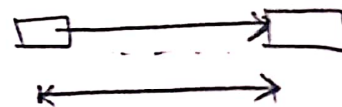
$\frac{L_{KD}}{L}$
Large knowledge distillation

end-to-end average latency

↳ architecture: Volta, Turing, NVIDIA, Ampere.

sequential layer-by-layer KD

$$\min_{L_i} \|L_i(X_{i-1}) - L_i^q(X_{i-1}^q)\|_2^2$$



W8A5 \xleftrightarrow{VS} W16A16
54M

256

↳ 128

= 28.13

128

↳ 140

↳ 「milliseconds」

BERT model →



비교 가능한 모델

