LLM Quantization - Basic methods

Lei Li



Carnegie Mellon University Language Technologies Institute

Ack: Kath Choi, Aashiq Muhamed, Zhen Wu, Hongyi Jin, Wayne Wang

Summary

- Pipeline Parallelism

 split by layers (horizonal split)
 eliminate the bubbles (idle)
 interleaving forward/backward
- Tensor Parallelism
 o split the matrix computation

Today's Topic

- Low precision numbers in computer
- Basic Quantization Methods

LLM Training and Inference are Costly!

Llama-70B

39.3M H100-80GB GPU hours to train requires 140GB GPU memory for inference

<u>Deepseek V3</u> (671GB) 2.8M H800 GPU hours to train requires > 400GB GPU memory for inference



Model Quantization

- Use low-bit precision to store parameters and layer outputs
- Quantization can

o reduce memory → larger batch size
o speed up calculation, more operations in one cycle

• Cons: potentially reduce accuracy

Precision Formats

$$\pm 1.(mantissa) * 2^{\{a-2^{\{exponent-1\}}\}}$$



BF16 favors Dynamic Range than Precision.

INT8 range: -128~127, Precision: integer



Image from https://deeprec.readthedocs.io/en/latest/BFloat16.html

BF16/FP16 Calculations are faster!

HFMA2: Half-precision Fused Multiply-Add for 2 elements in one cycle (2x speedup)



CUDA APIs for Half Precision

Performs half2 vector addition in round-to-nearest-even mode.

___device____half2 ___hadd2(const ___half2 a, const ___half2 b)

Performs half2 vector fused multiply-add in round-to-nearesteven mode.

__device___half2 __hfma2(const __half2 a, const __half2 b, const __half2 c)

Direct Quantization Approach

• Using lower precision

o converting parameters from FP32 to INT8 or INT4
o perform all computation in lower prevision.

• Reduce model accuracy:

 \circ Loss of Precision \rightarrow accumulate quantization noise

 o Range mismatch → values are clipped and lead to information loss

 \circ Quantization error \rightarrow rounding errors

Quantize a number



Code implementation for quant(.)

import torch

```
def absmax_quantize(X):
```

Calculate scale

scale = 127 / torch.max(torch.abs(X))

Quantize

```
X_quant = (scale * X).round()
```

Dequantize

X_dequant = X_quant / scale
return X_quant.to(torch.int8), X_dequant

def zeropoint_quantize(X): # Calculate value range (denominator) $x_range = torch.max(X) - torch.min(X)$ x_range = 1 if x_range == 0 else x_range # Calculate scale scale = 255 / x_range # Shift by zero-point zeropoint = (-scale * torch.min(X) - 128).round() # Scale and round the inputs $X_quant = torch.clip((X * scale + zeropoint).round(), -128,$ 127) # Dequantize X_dequant = (X_quant - zeropoint) / scale return X_quant.to(torch.int8), X_dequant

Direct Quantization Colab

https://colab.research.google.com/drive/1DPr4mUQ92Ccxf4GgAaB6dFcFnWIvqYi?usp=sharing

Today's Topic

- Low precision numbers in computer
- Basic Quantization Methods







Basic Approach: Layer-Wise Quantization

- Find the quantized matrix \widehat{W} that minimizes the linear layer's output.
- Performs quantization layer-by-layer $\operatorname{argmin}_{\widehat{W}} \| WX - \widehat{W}X \|_{2}^{2}$

W: linear projection weights (e.g. in FFN and attention) X: layer input

• Limitation: could still lead to accumulation of quantization error

AdaQuant

- minimize the error between the quantized / full-precision layer outputs for each layer
- adding continuous V to W and quantize

$$(\widehat{\Delta_{w}}, \widehat{\Delta_{x}}, \widehat{V}) = \underset{\Delta_{w} \Delta_{x}, V}{\operatorname{argmin}} \| WX - Q_{\Delta_{w}}(W + V) \cdot Q_{\Delta_{x}}(X) \|^{2})$$
quantized results

Hubara et al. Improving Post Training Neural Quantization: Layer-wise Calibration and Integer Programming. 2020. 18

Is Quantization Accurate?

MobileNet-v2	Float Model	Direct Quantized Model	Quantized Model with ADAQUANT
Model Size	8.4 MB	2.3 MB	2.4 MB
P rec@1	65.4 %	1.7 %	52.3 %
P rec@5	85.7 %	5.6 %	75.7 %



ZeroQuant

- Layer-by-layer knowledge distillation
 - Use the original model as Teacher
 - the quantized model is student

$$\mathcal{L}_{LKD,k} = \left| L_k \cdot L_{k-1} \cdot L_{k-2} \cdot \dots \cdot L_1(X) - \hat{L}_k \cdot L_{k-1} \cdot L_{k-2} \cdot \dots \cdot L_1(X) \right|^2$$
Layer L_{k-1}

Update quantization scheme for L_k

Layer L_k

Update quantization bistillation Loss

Yao et al. ZeroQuant: Efficient and Affordable Post-Training Quantization for Large-Scale Transformers. Neurips 2022. ²¹

ZeroQuant

• Quantization-Optimized Transformer Kernels (fusion)



Figure 2: The illustration of normal (left) and our fused (right) INT8 GeMM.

- The scalability is verified up to 20B models (GPT-NeoX_{20B})
- At 1.3B scale, computation time is ~3 hours
 - but slower than GPTQ (x100 larger in ~4 hours)
- integrated in Deepspeed

Yao et al. ZeroQuant: Efficient and Affordable Post-Training Quantization for Large-Scale Transformers. Neurips 2022.

22

LLM.int8()

- Using 8-bit quantization for matrix multiplications
- But, extreme outliers in features (activation values)

 need for wider numerical ranges
 Quantize all parameters without distinguishing them separately can result in accuracy degradation



LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale. Dettmers et al. Neurips 2022.

LLM.int8()

- Keep outliers in higher precision while quantizing the rest
- Outliers: large magnitude (>= 6.0), affects >= 25% layers, and affects >= 6% sequence dimensions



I M int8(): 8-bit Matrix Multiplication for Transformers at Scale. Dettmers et al. Neurips 2022.

Summary

- Iow-bit number representation in computer

 BF16: 16-bit half precision floating point numbers, better for ML
 int8
- Direct quantization

o absmax: linearly scale according to max abs value
o zero-point: finding zero-point and scale

- Layer-wise quantization approaches

 AdaQuant
 - KD: ZeroQuant
 - o LLM.int8()

Next

• Scaling Quantization for large models: GPTQ