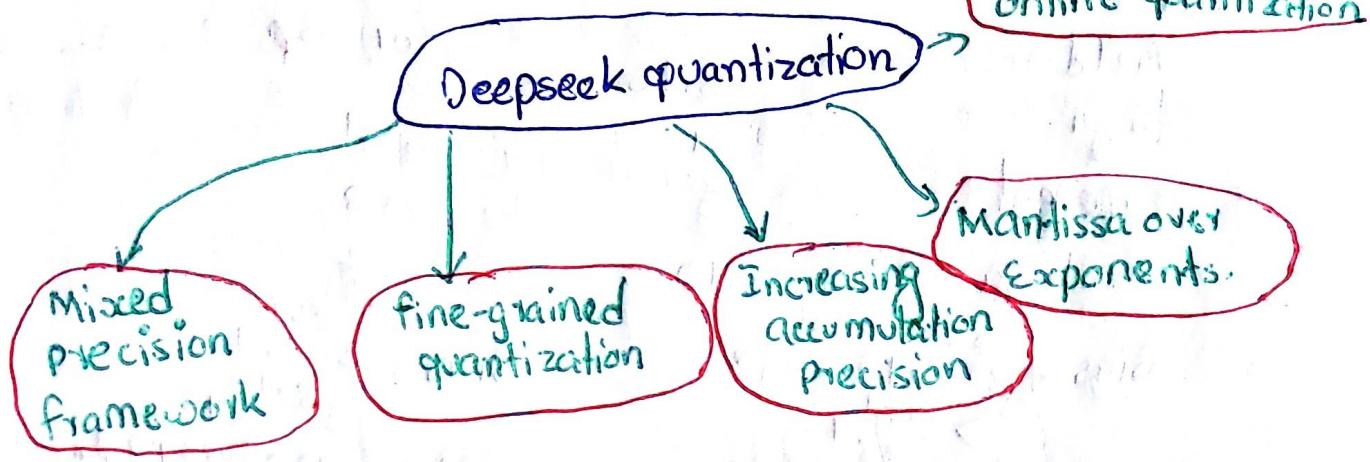


# Quantization



How a number is stored in FP32??

→ According to IEEE-754 standard.

field      Bits  
 Sign      1-bit  
 Exponent( $e$ )      8 bits  
 Mantissa      23 bits  
 suppose  $x = 3.1415927410125732$ .

Binary(approx) =  $11.001001000011111$

Normalized =  $1.100100100011111 \times 2^e$

=  $(\epsilon)^s \times (1 + \text{mantissa}) \times 2^{e-127}$  (formulae)

stored 32bit

0    1 0 0 0 0 0 0 0  
 1bit      8bit

1 0 0 1 0 0 1 0 0 0 0 1 1 1 1  
 23 bits

To find the exponent bit we compare the formulae with the Normalized binary of the number.

$$\text{we get, } e-127 = 1$$

$$e = 128$$

we add bias to minimize the negative to occur.

for FP16

field

Sign

Exponent

Mantissa

$$\text{Bias} = 2^{e-1} - 1 \Rightarrow 2^5 - 1 = 15$$

Bits

1

5

10

1

8

7

for Bf16

$$\text{Bias} = 2^{e-1} - 1 = 2^{8-1} - 1 = 128 - 1 = 127$$

## Mixed Precision Framework

↪ suppose

$$y = Wx$$

computed initially  
in FP32 for numerical  
stability then converted  
back to BF16.

↪ converted to FP8 from BF16.

↪ maintained in higher precision (BF16 or FP32),  
then converted to FP8 on the fly.

(Fprop in paper/figure).

$$y_{\text{FP32}} = W_{\text{FP8}} \times x_{\text{FP8}}$$

↪ stored in BF16 to optimize memory.

In the middle of a neural network,

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} * W^T$$

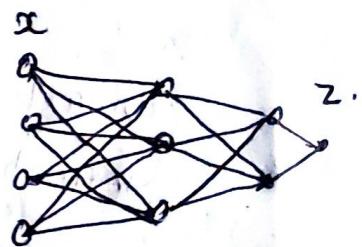
↪ it comes as,  
 $z = Wx$

$$\frac{\partial z}{\partial x}$$

computed  
initially as FP32, then  
converted to BF16.

stored as  
BF16 but  
converted to

FP8 for  
computation.



we know,

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} * \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z}$$

To maintain  
dimension in  
back propagation

Dgrad  
in paper

$$\frac{\partial L}{\partial x_{\text{FP32}}} = \frac{\partial L}{\partial z_{\text{FP8}}} * W^T_{\text{FP8}}$$

↪ stored as BF16.

weights are stored at higher  
precision so that we can confirm,  
there's update in it. it's learning!!

## W grad in paper

$$\frac{dL}{dW} = X^T \frac{\delta L}{\delta Z}$$

↓  
stored in  
FP8.

computed  
+ stored in  
FP32 to ensure  
numerical  
stability for weight  
updates.

→ converted from BF16 to FP8

why to do in such approaches?

### FP8 Precision

- ↳ significantly reduces computational & memory cost.
- ↳ provides up to 2x speed improvement over BF16 operations.
- ↳ ideal for computationally intensive GEMM operations.

### Higher Precision Retained (FP32/BF16):

- 1) Embedding modules, output heads, gating modules, normalization operators, attention operators are sensitive. Hence, these are retained in higher precision.
- 2) Weights & optimizer states also stored in FP32 ensure training stability.

This balance of high-precision storage with low-precision calculation allows the best trade-off between computational efficiency, memory reduction & numerical stability.

## b) Fine grained quantization

↳ This method was introduced to protect the precision of weight dimension when they are quantized to lower precision.

\* If any outlier occur, it makes the numbers with different values in quantization & dequantization.

$$\text{Suppose: } x = [2, 3, 4, 500] \times \frac{127}{500} = [0.508, 0.762, 1.016, 127.0]$$

Now, quantize to FP8. Assume EAM3 quantization

$$= [0.5, 0.75, 1.0, 127]$$

Dequantize back

$$\text{Recovered Tensor} = [0.5, 0.75, 1.0, 127] \times 500 \approx [1.97, 2.95, 3.94, 500]$$

This occurs due to outlier.

So, what to do??

We have billion parameters in weight, in h-dim also, we have

large no. of values, so,

we will group them in groups & then do quantization,

so, that outlier in one group didn't affect the precision of

another group.

$$x = [1.2, 2.1, 2.001, 3, 76, 0.1, 1.2, 3]$$

group 1    group 2

scaled by factor 1                              scaled by factor 2.

$$\Rightarrow = \boxed{x} \times \boxed{W}$$

tensor core. (General matrix mul occurs in FP8 inside NVIDIA Tensor core.)

In CUDA core.

$$\text{FP8} * \text{upscale factor} = \text{FP32}$$



w11 → scaling factor 1

w12 → scaling factor 2

w21 → scaling factor 3

w22 → scaling factor 4.

## Increasing accumulation Precision

→ Using higher-precision arithmetic only for the accumulation (summing) step during neural-network computations. Even though the inputs, weights, & outputs are stored in lower precision.

We know, in almost every neural network layer, we have

$$y_i = \sum_{j=1}^d x_j \cdot w_{ji}$$

This has two different operations.

- i) Multiply  $\rightarrow x_j \cdot w_{ji}$
- ii) Accumulate (sum)  $\rightarrow \sum_j (\cdot)$ .

Deepseek-V3 keeps (1) in low precision & (2) in higher precision.

why low-precision accumulation is dangerous??

→ suppose Deepseek uses FP8/BF16 for speed

problem:

→ FP8/BF16 has few mantissa bits.

→ when summing thousands of terms, small values get

- rounded away

- lost entirely.

This causes:

- large numerical error
- unstable training
- exploding / vanishing gradients.

## WGMM

Low-precision matrix Multiply Accumulate.

→ It's a local-level GPU instruction family introduced with NVIDIA Volta (H100) & it's central to how models like V3 run FP8 efficiently.

Key: multiplication in low precision in tensor core & accumulation is done in FP32 in CUDA core.

Mantissa over Exponents → E4M3 → 3 mantissa bit higher precision  
in FP8 there are two format → E5M2.

Deepseek prefer S exponent bits high range  
Always use floating point number E4M3 format.

but prior works prefers floating point number E4M3 format.

\* Forward pass: E4M3

\* Backward pass: E5M2.

## Online quantization

↳ Scaling factor is decided based on the numbers that actually come in calibration (Delayed quantization)

main problem of offline quantization  
↳ we run a calibration dataset, and decide the scaling factor at once & after that we do quantization based on that.

\* If a rare token appears, overflow or loss.

What deepseek does?  
↳ it calculate the scaling factor for the fly

## Tensors

\* In inference, each weight tensor has one (or a small set of) scale value(s), computed ONCE, & reused for all tokens.

so, in inference, we use offline quantization.