

EurIPS

Unlocking FP4: Low Precision AI on NVIDIA Blackwell

Speakers

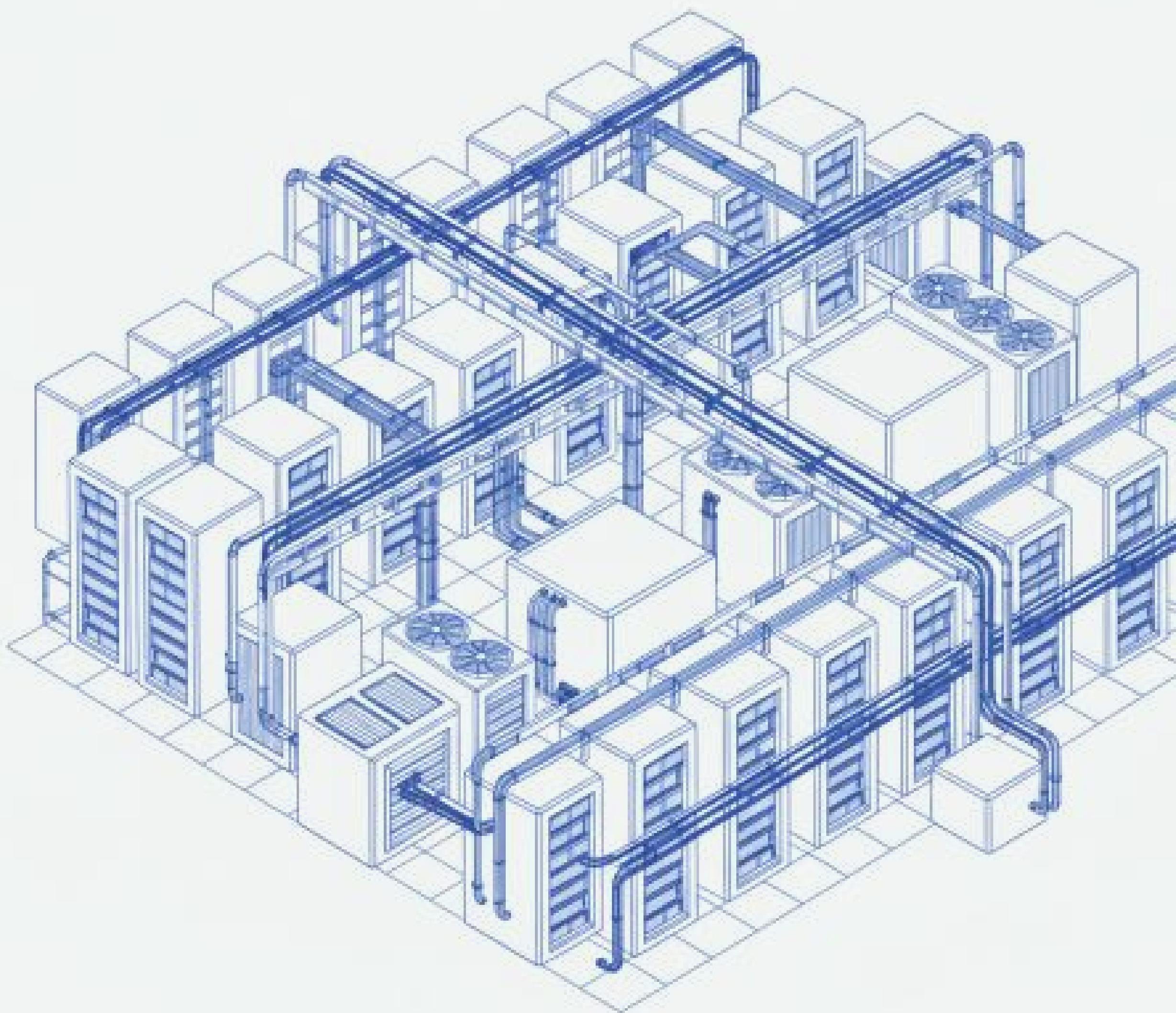


Paul Chang
ML Engineer



Riccardo Mereu
ML Engineer

Hyperscalers were built for CPU - AI cloud is different



The AI Cloud



Compute cost matters

Legacy cloud was built around saving developer time.
AI cloud optimizes for performance and compute costs.



GPU Data Centers need power

Data center build out starts from zero as energy density, power supply and cooling innovations make old footprint obsolete



Sovereignty is critical

Geopolitical shift towards national sovereignty especially in Europe requires an AI-native Hyperscaler built in Europe

Our AI-Native cloud beats hyperscalers for AI workloads

Clean modern architecture

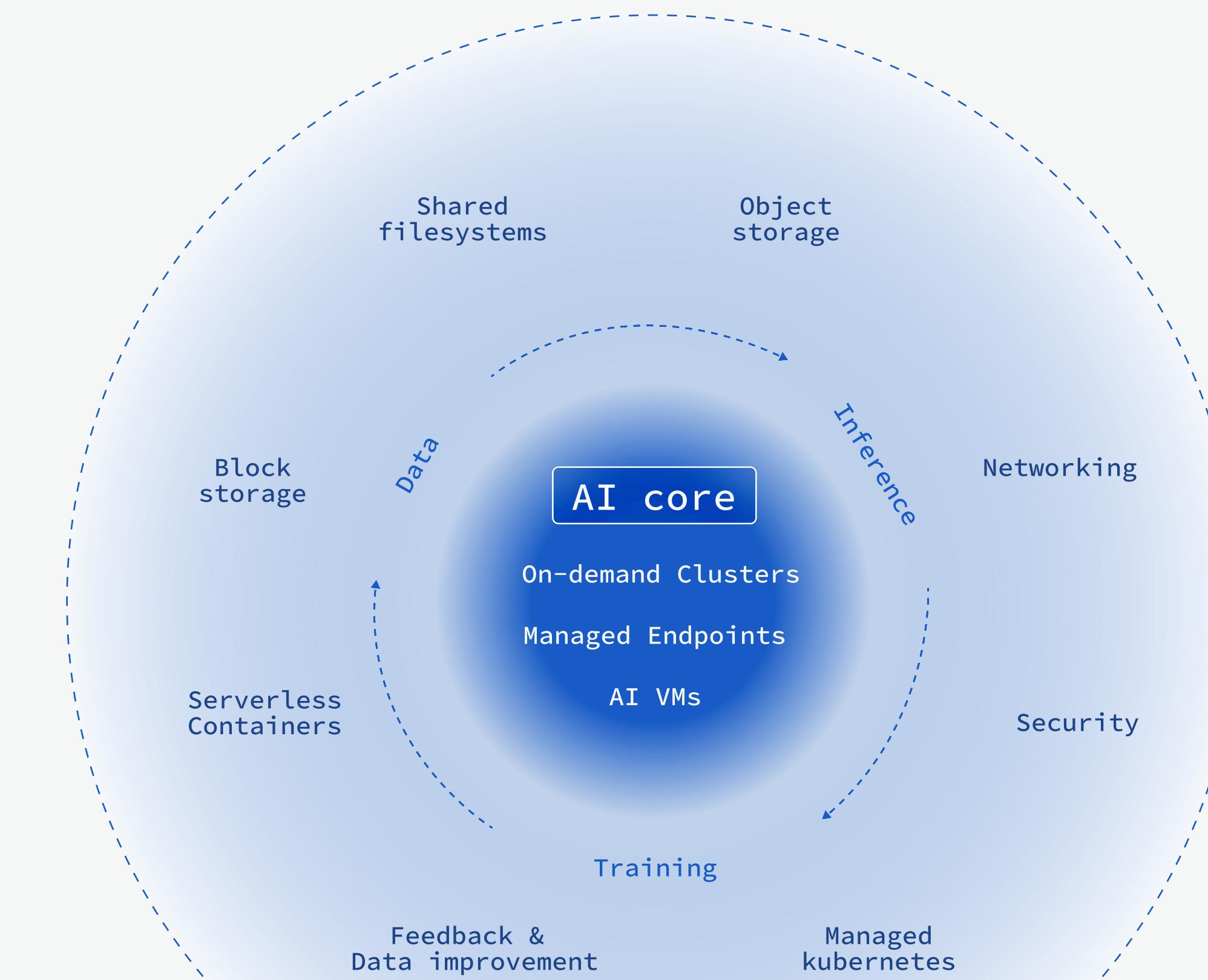
Redesigned from scratch for GPU workloads, beating legacy piecemeal cloud architecture

Unified control plane

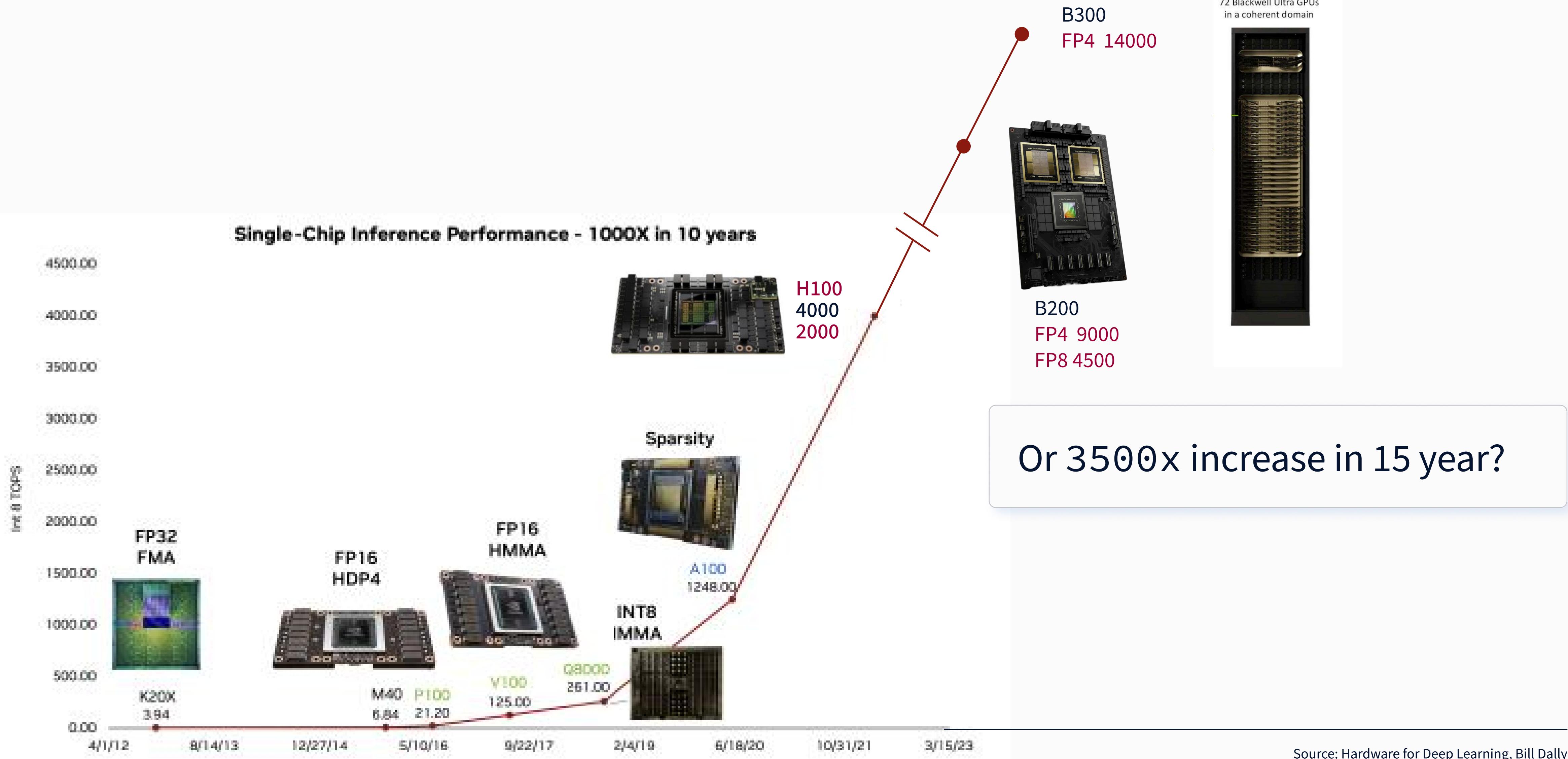
One cohesive compute layer, beating legacy fragmented layer system on utilisation and developer experience

Deep AI framework integration

Natively optimised for PyTorch, JAX and more, enabling faster adoption of new AI tooling



How compute has gone 1000x?



Where does the increases come from?

- Number Representation:
16x (64x) FP32 → BF16 → FP8 → FP4.
- Amoritized Instructions (Tensor cores):
12.5x FMA → HMMA → WGMMA → UMMA
- Process (fab size):
2.5x 28nm → 16nm → 7nm → 5nm (more transistors per mm²)
- Model efficiency has also boosted performance, FA, FA2, FA3, FA4.

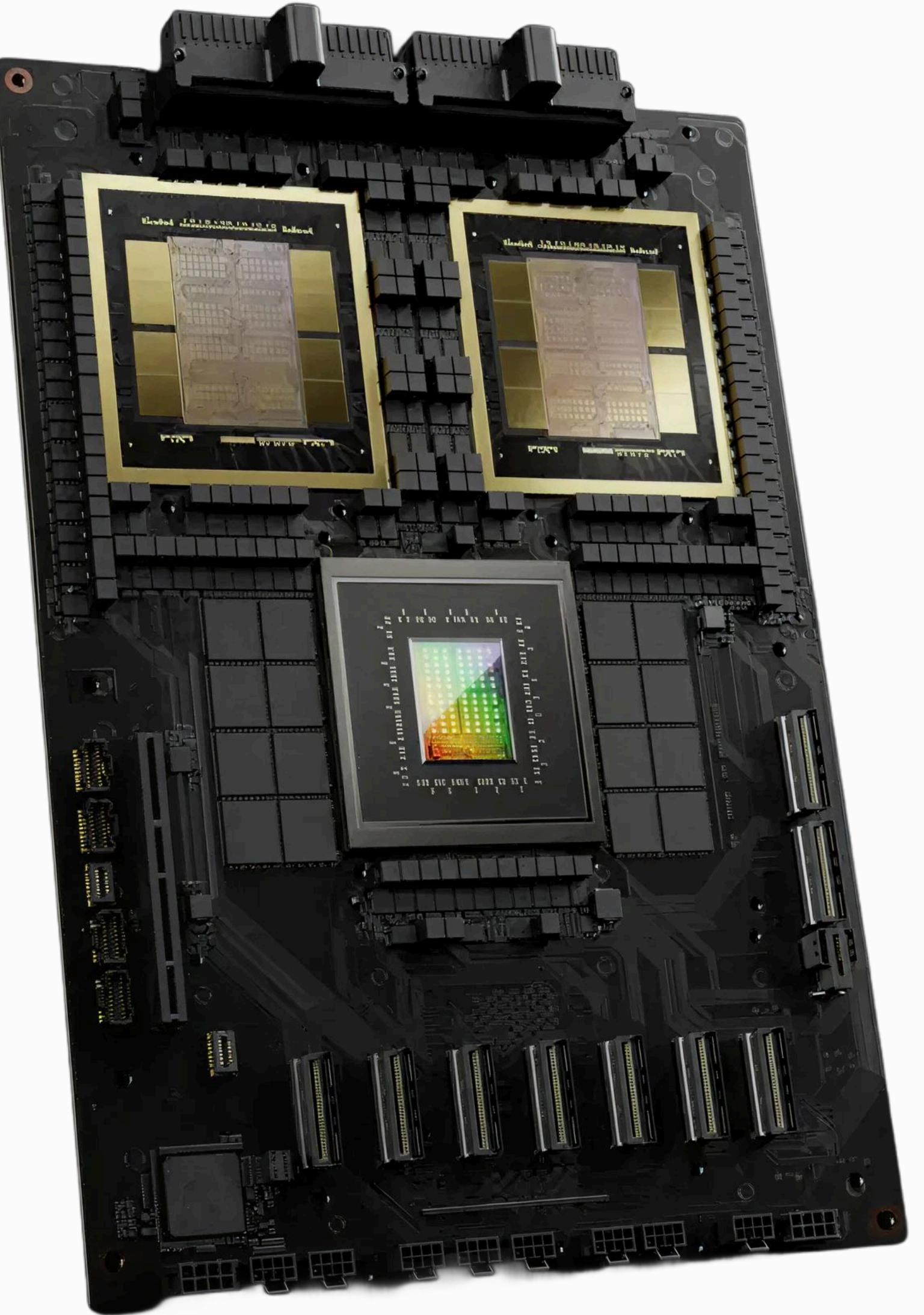
Section

Low Precision AI

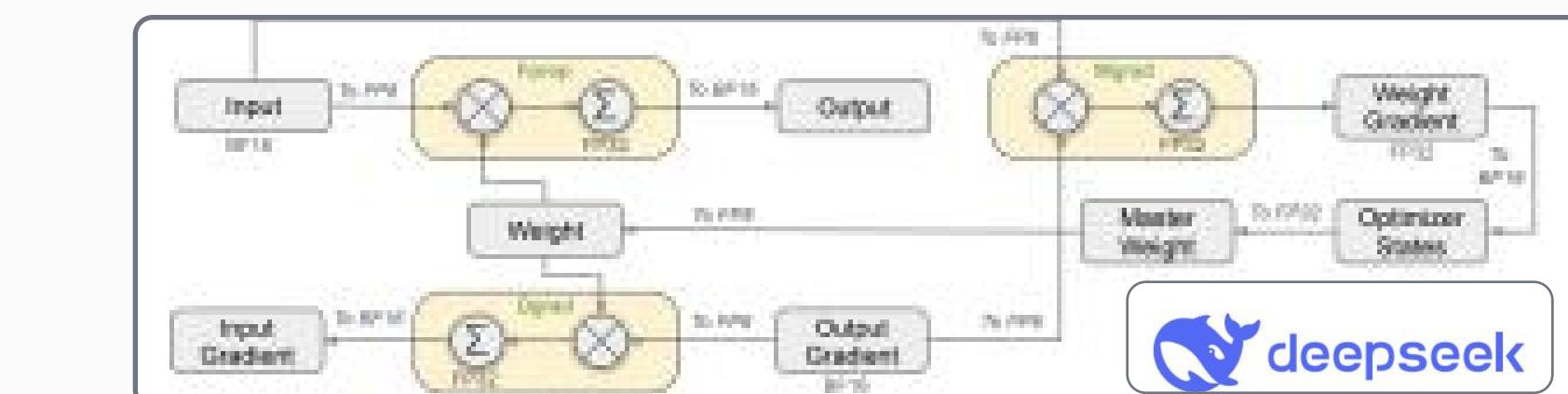
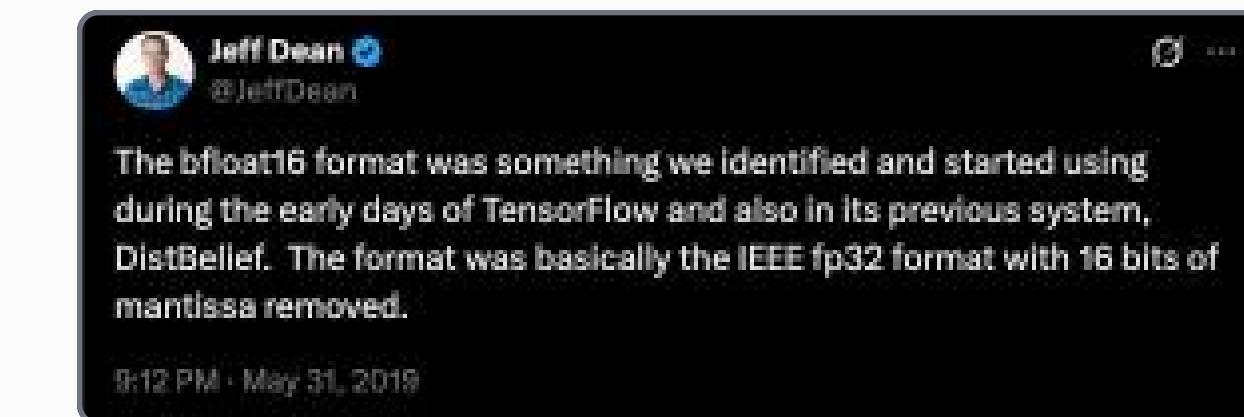
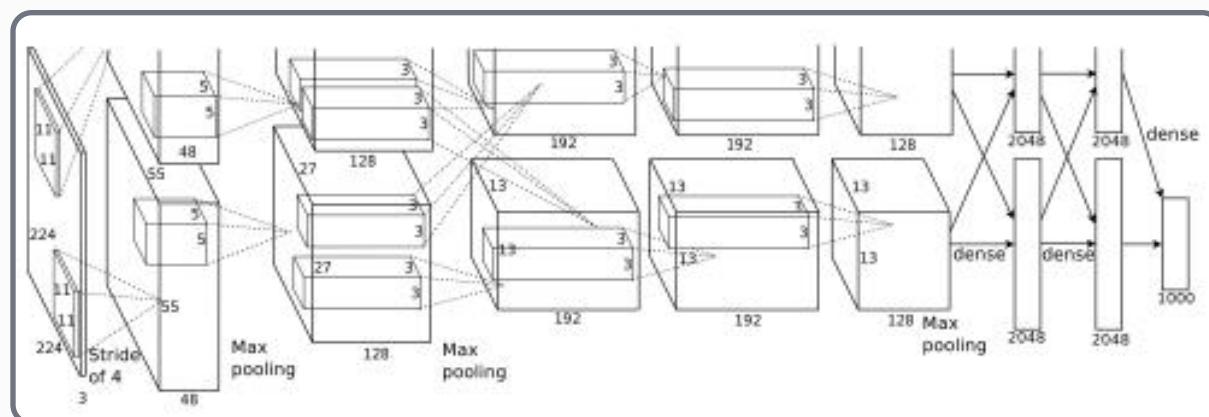
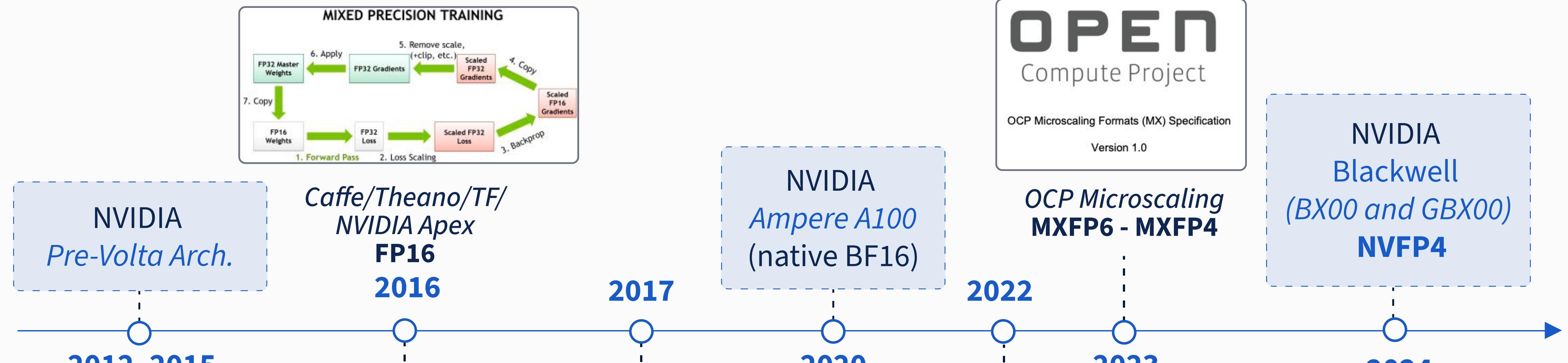
Unlocking FP4: Low Precision AI on NVIDIA Blackwell

Outline

1. Why do we need Low Precision AI?
2. What is NVFP4?
3. How NVIDIA Blackwell enable this?



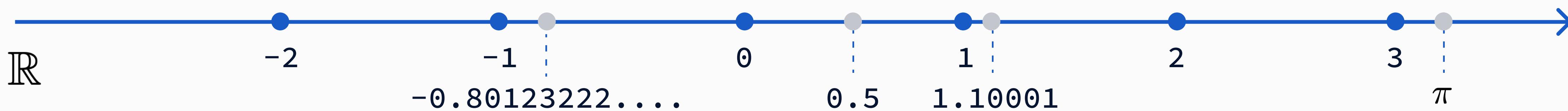
Timeline: Floating Point in AI



Floating Point Basics

Real Number Line:

- Arbitrarily large and small values → **Infinite range**
- Also between any pair of numbers → **Infinite precision**



We cannot do infinite:

- Finite memory
- Finite silicon area (or time to complete operations)

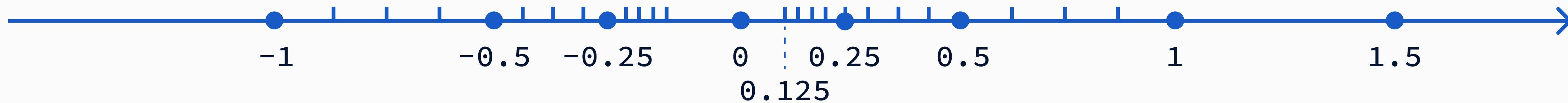
Need for number representation to sample the real line:

- Integer variants
- Floating Point

Floating Point Basics

FP Line:

- FP samples the real number line with equal number of samples between each pair of adjacent powers of 2
- *Example: E5M10 has 1024 samples between 0.25 and 0.5, between 128 and 256*



Bit Fields (ExMy):

- **Sign:** 1 bit (0 positive, 1 negative)
- **Exponent:** E bits, which power of 2 is sampled → *dynamic range*
- **Mantissa:** M bits, samples between powers of two → *precision*

$$N_{10} = (-1)^S \times 1.M_{10} \times 2^{E_{10} - \text{bias}}$$

$$\text{Exponent bias} = 2^{E_N - 1} - 1$$

Floating Point Basics

FP16 Example (E5M10)

- **Sign:** 1 bit (0 positive, 1 negative)
- **Exponent:** 5 bits $2^{5-1} - 1 = 15$
- **Mantissa:** 10 bits
- **Exponent bias:**

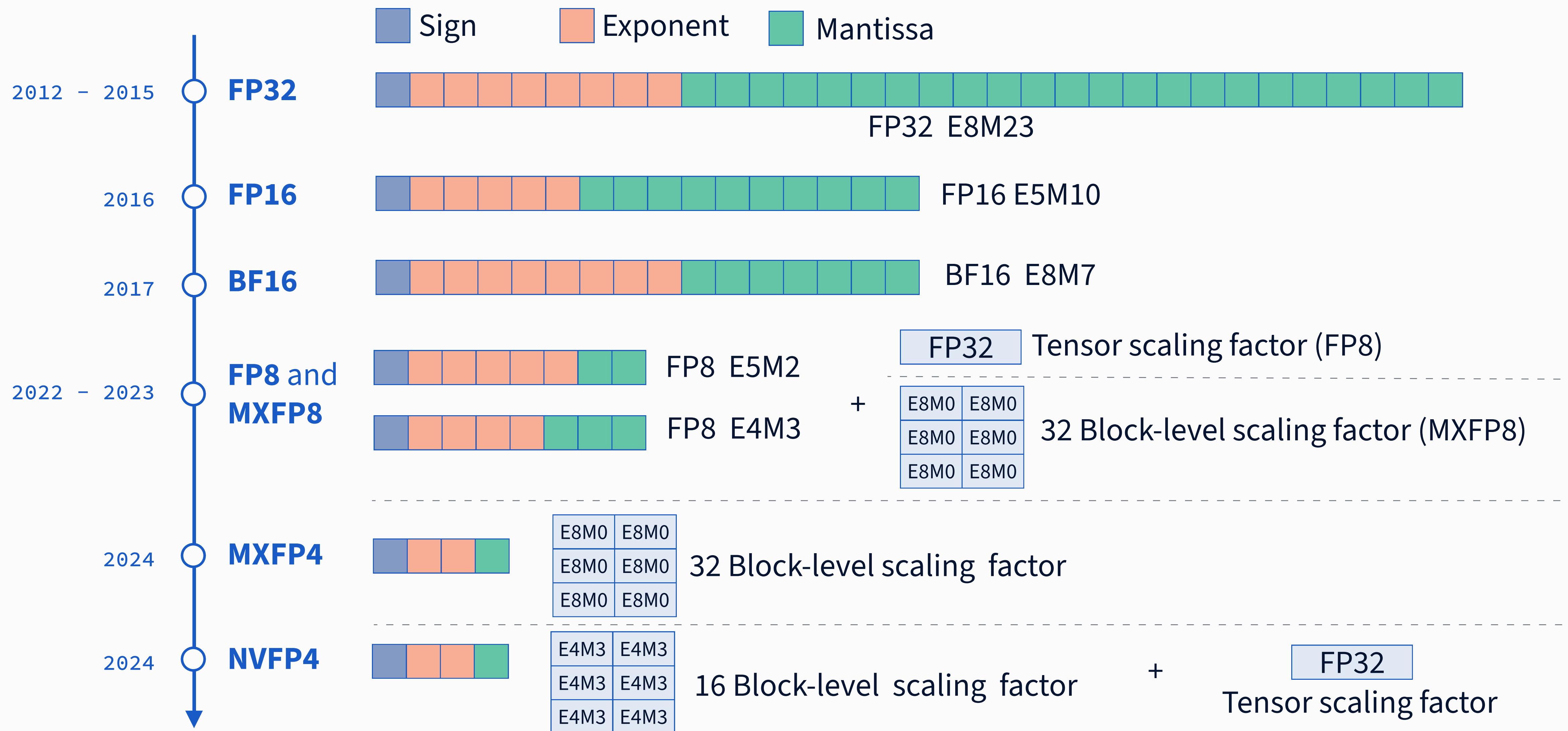
$$N_{10} = (-1)^S \times 1.M_{10} \times 2^{E_{10}-\text{bias}}$$

FP 16 bits:

$$0.\textcolor{blue}{10001}.\textcolor{red}{1010000000} = (6.510)_{10}$$

$$\begin{aligned}(-1)^{\textcolor{blue}{0}} \times 1.\textcolor{red}{101}_2 \times 2^{\textcolor{brown}{10001}_2 - 15} &= (1 \times (1 + 2^{-1} + 2^{-3}) \times 2^2)_{10} = \\&= (4 \times 1.625)_{10} = (6.510)_{10}\end{aligned}$$

Floating Point in AI



Model Quantization Algorithms

There are several strategies to quantize neural networks:

- **Post Training Quantization (PTQ)**

- Turn pre-trained **weights** and/or **activations** into low-precision formats
- No need for training data or training hardware
- May lead to a lower quality model

Neural Network
Inference

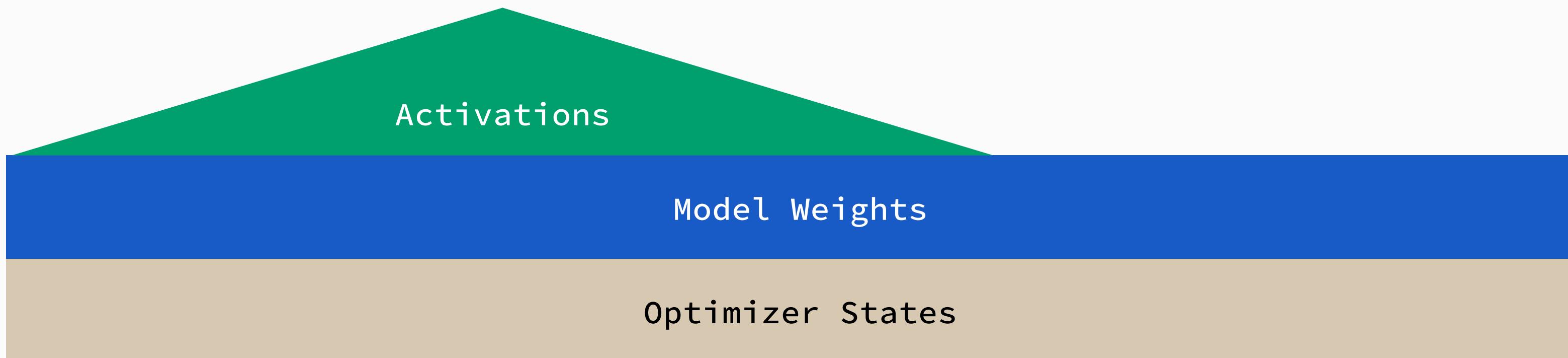


Model Quantization Algorithms

There are several strategies to quantize neural networks:

- **Post Training Quantization (PTQ)**
 - Turn pre-trained **weights** and/or **activations** into low-precision formats
 - No need for training data or training hardware
 - May lead to a lower quality model
- **Quantization Aware Training (QAT)**
 - Adds quantization to the forward for training/fine-tuning

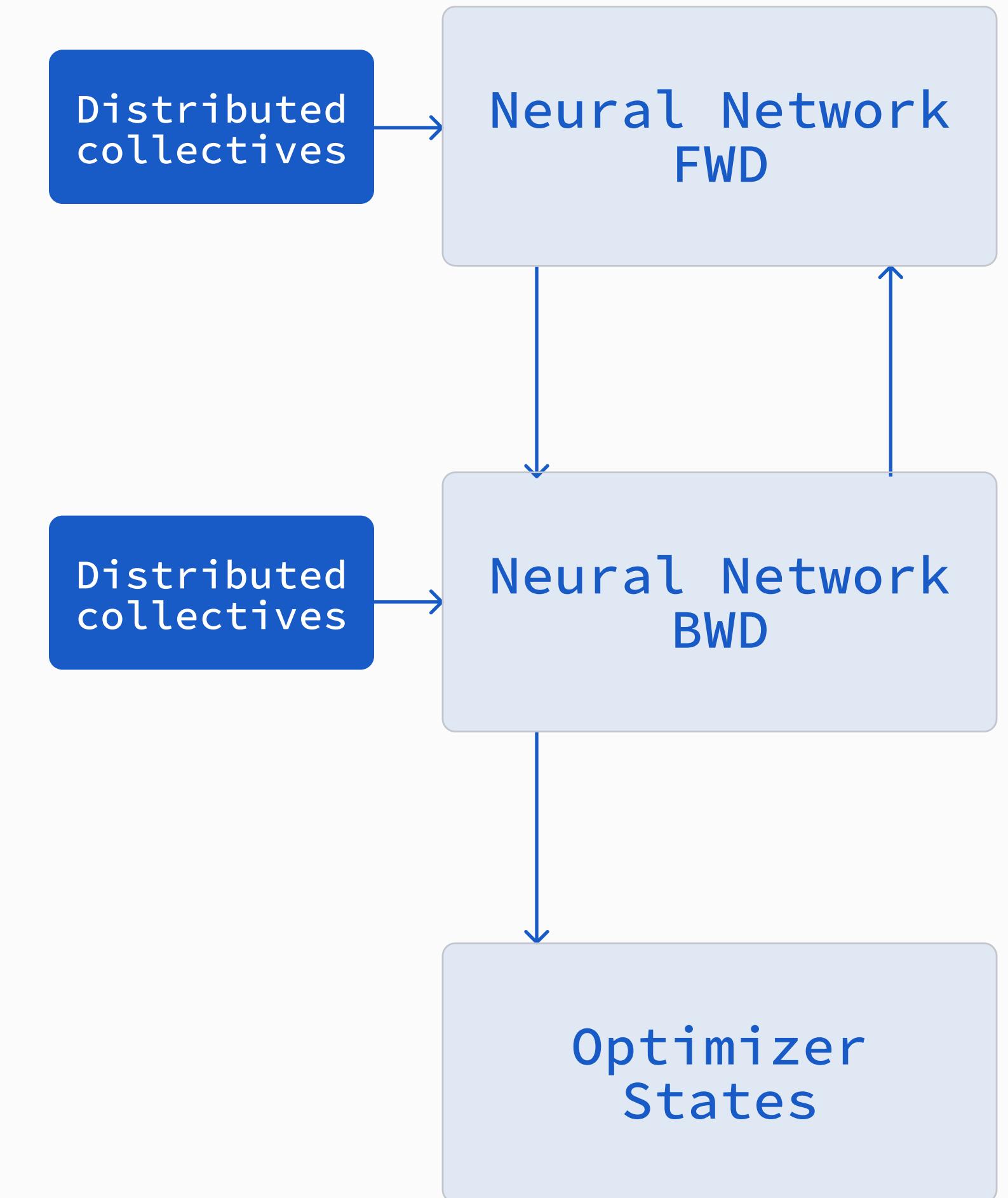
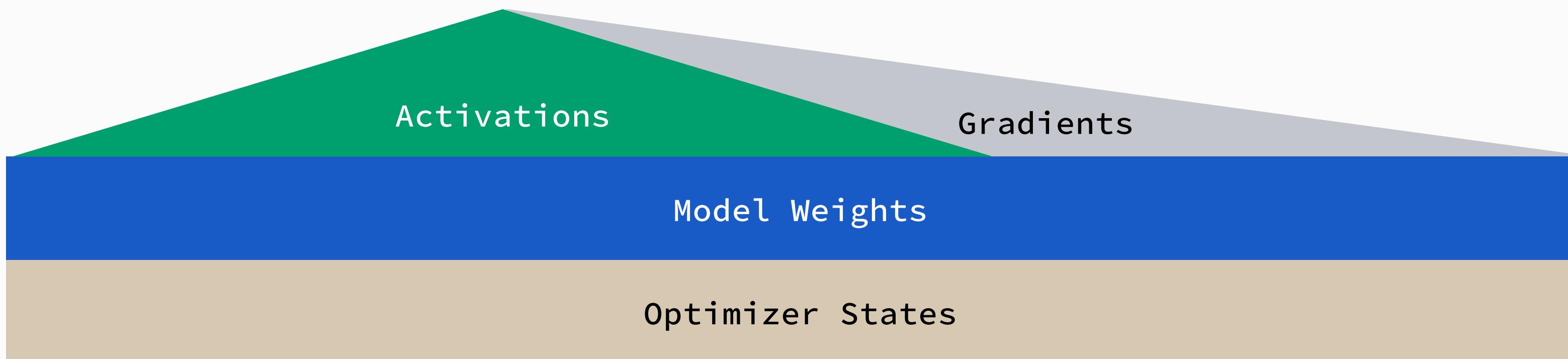
Neural Network
FWD



Model Quantization Algorithms

There are several strategies to quantize neural networks:

- **Post Training Quantization (PTQ)**
 - Turn pre-trained **weights** and/or **activations** into low-precision formats
 - No need for training data or training hardware
 - May lead to a lower quality model
- **Quantization Aware Training (QAT)**
 - Adds quantization to the forward for training/fine-tuning
- **Quantized Training (QT)**
 - Adds quantization also to the backward pass improving over QAT



DeepSeek-V3 Mixed Precision

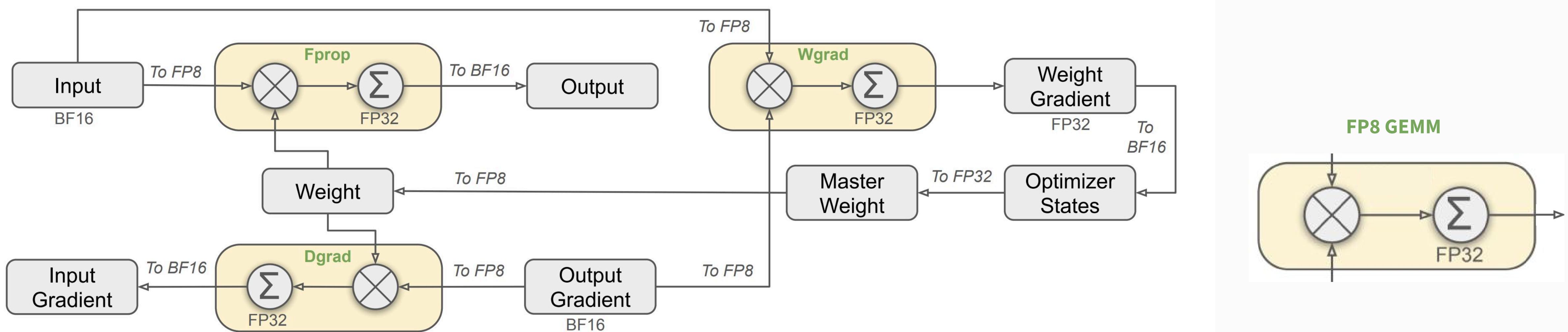
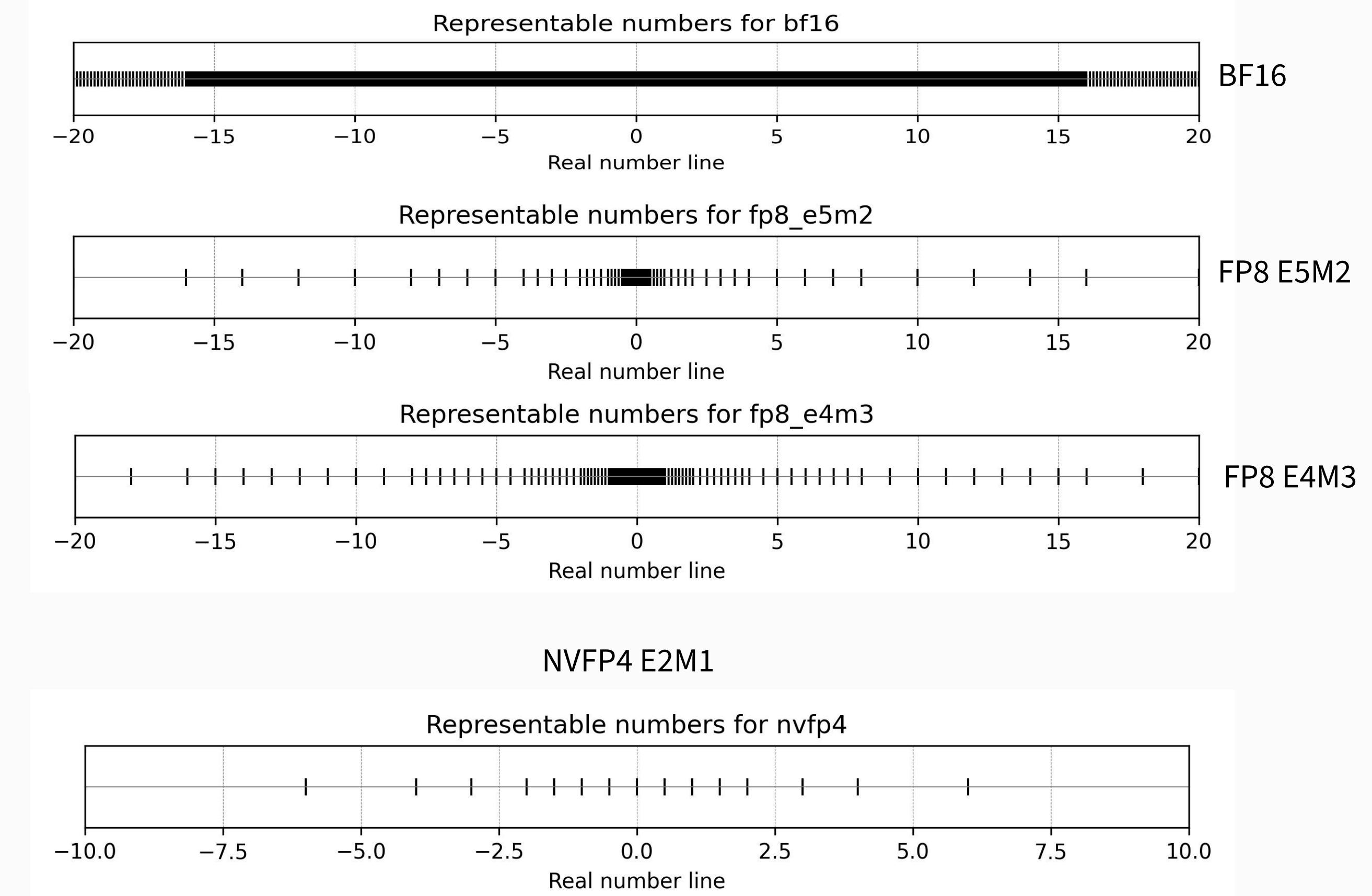


Figure 6 | The overall mixed precision framework with FP8 data format. For clarification, only the Linear operator is illustrated.

NVFP4: Are 16 values enough?

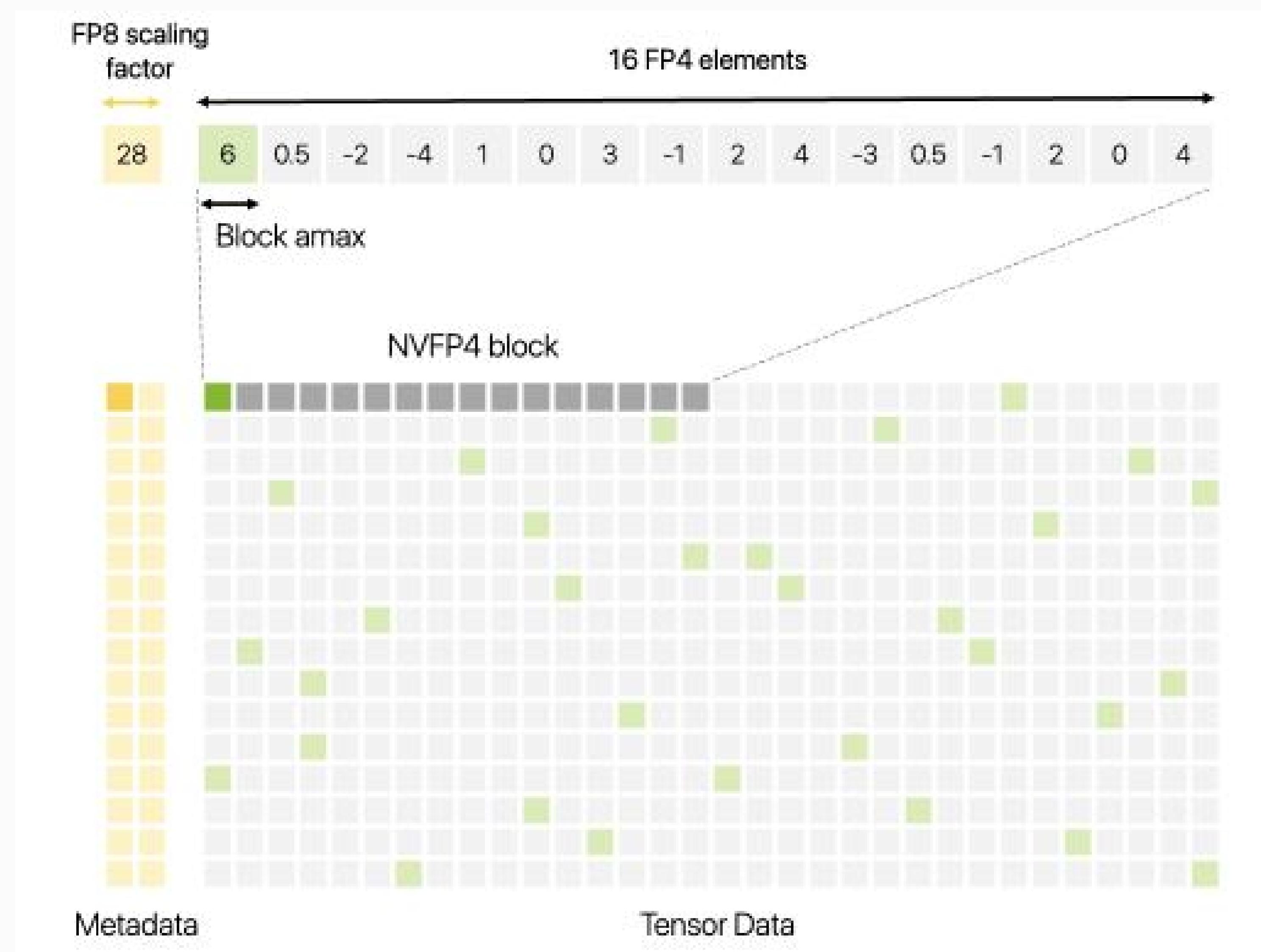
sign	exponent	mantissa		
0	0 0 0	0	= 0.0	
0	0 0 0	1	= 0.5	
0	0 0 1	0	= 1.0	
0	0 1 1	1	= 1.5	
0	1 0 0	0	= 2.0	
0	1 0 1	1	= 3.0	
0	1 1 0	0	= 4.0	
0	1 1 1	1	= 6.0	
1	0 0 0	0	= -0.0	
1	0 0 0	1	= -0.5	
1	0 1 0	0	= -1.0	
1	0 1 1	1	= -1.5	
1	1 0 0	0	= -2.0	
1	1 0 1	1	= -3.0	
1	1 1 0	0	= -4.0	
1	1 1 1	1	= -6.0	



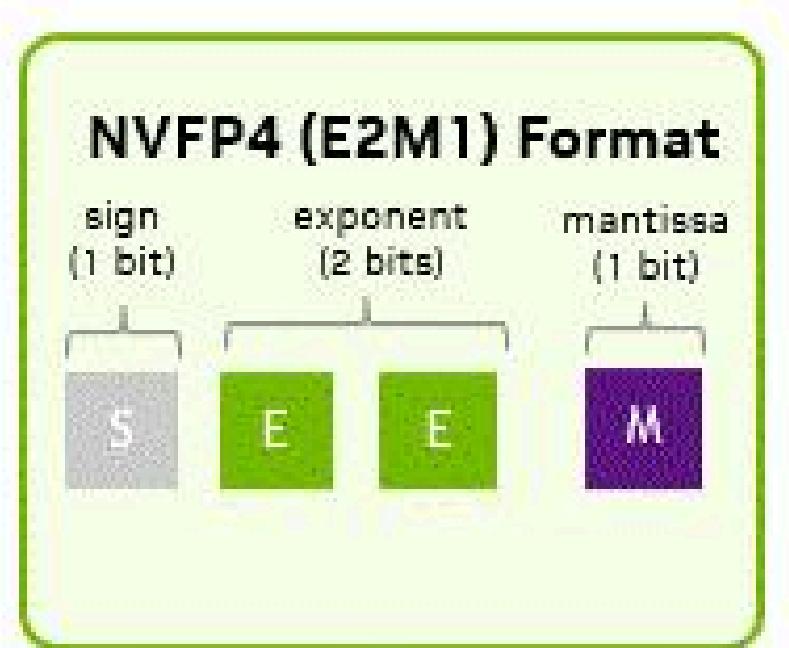
NVFP4 vs. MXFP4

The smaller block size and more precise scaling allows:

- to increase accuracy of outliers
- to minimize the values quantized to zero



NVIDIA NVFP4



NVFP4 Recipe

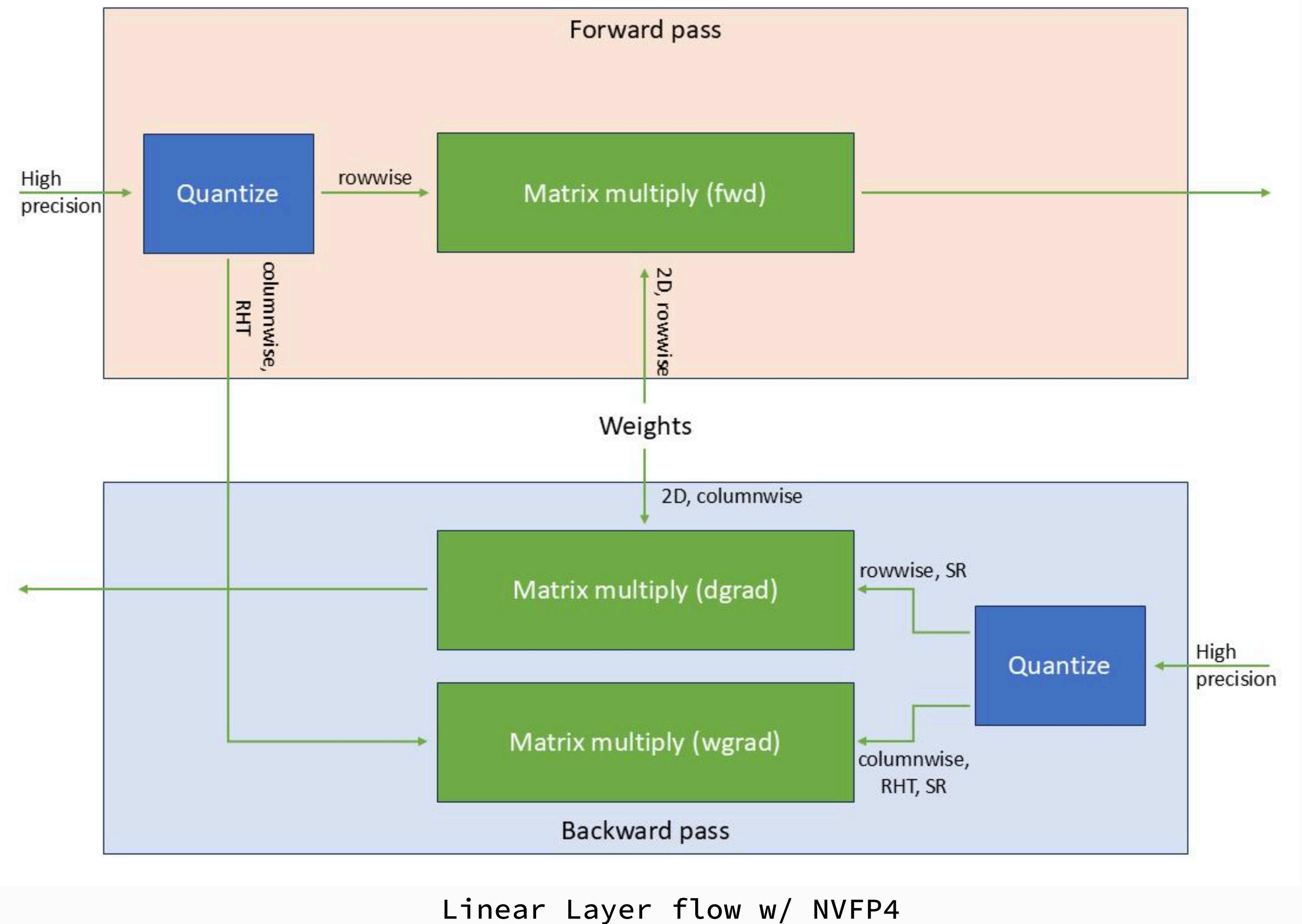
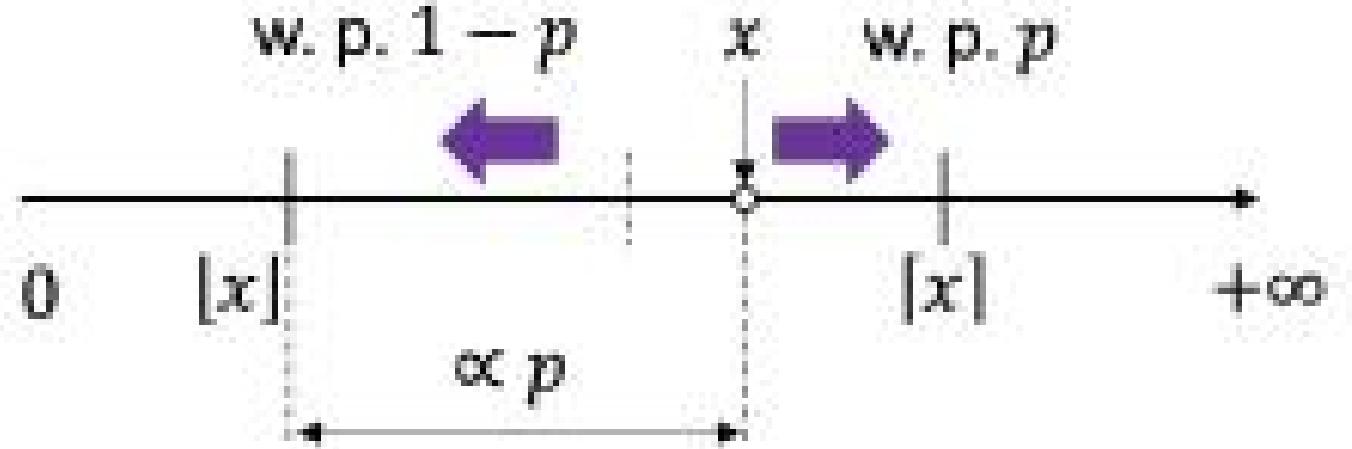
NVFP4 Recipe:

- 2D Scaling
- Stochastic Rounding

$$\text{Round}(x) = \begin{cases} \lfloor x \rfloor, & \text{w/ prob. } 1 - p, \\ \lceil x \rceil, & \text{w/ prob. } p. \end{cases}$$

$$p = (x - \lfloor x \rfloor)/(\lceil x \rceil - \lfloor x \rfloor)$$

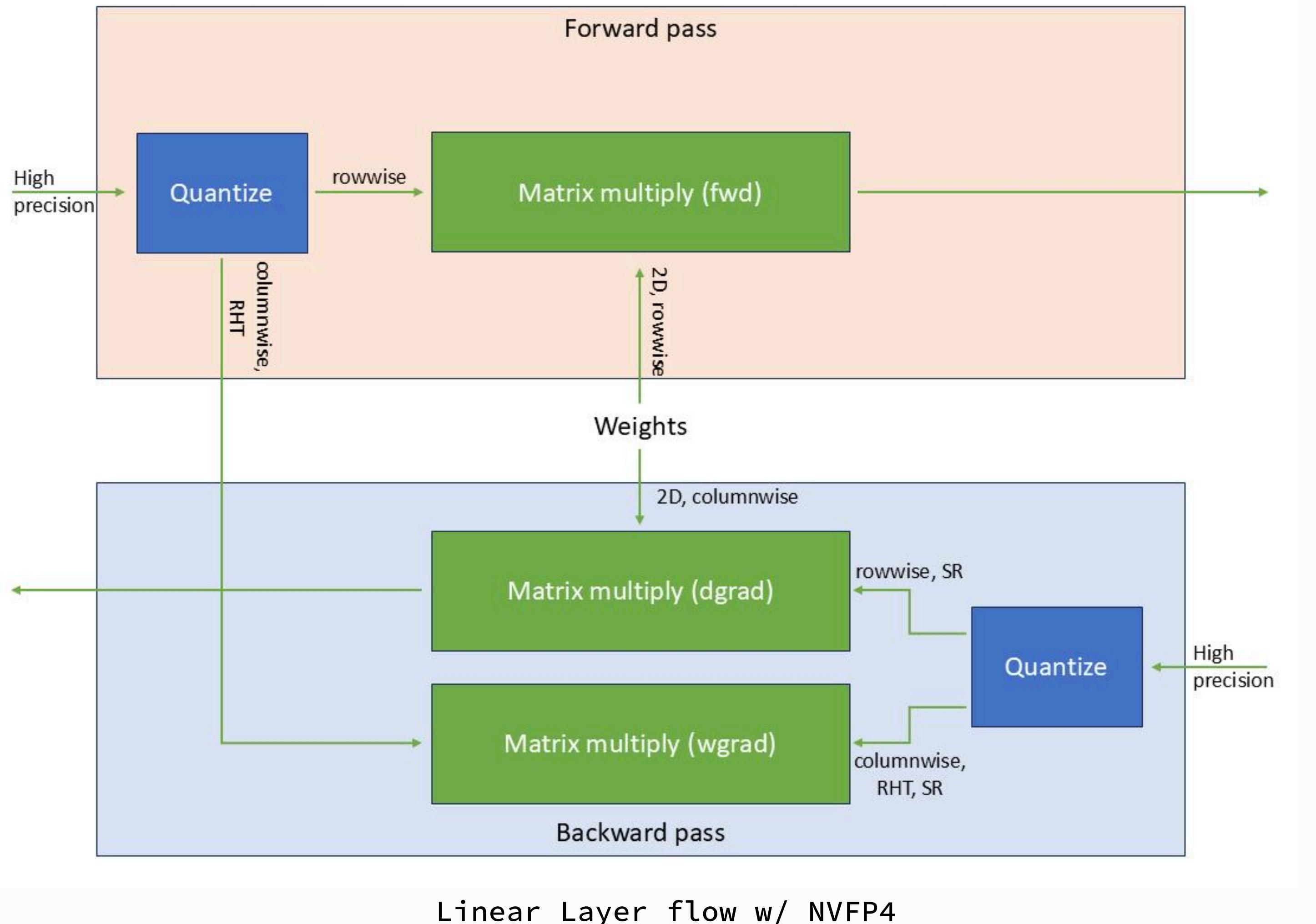
$$\mathbb{E}[\text{Round}(x)] = x$$



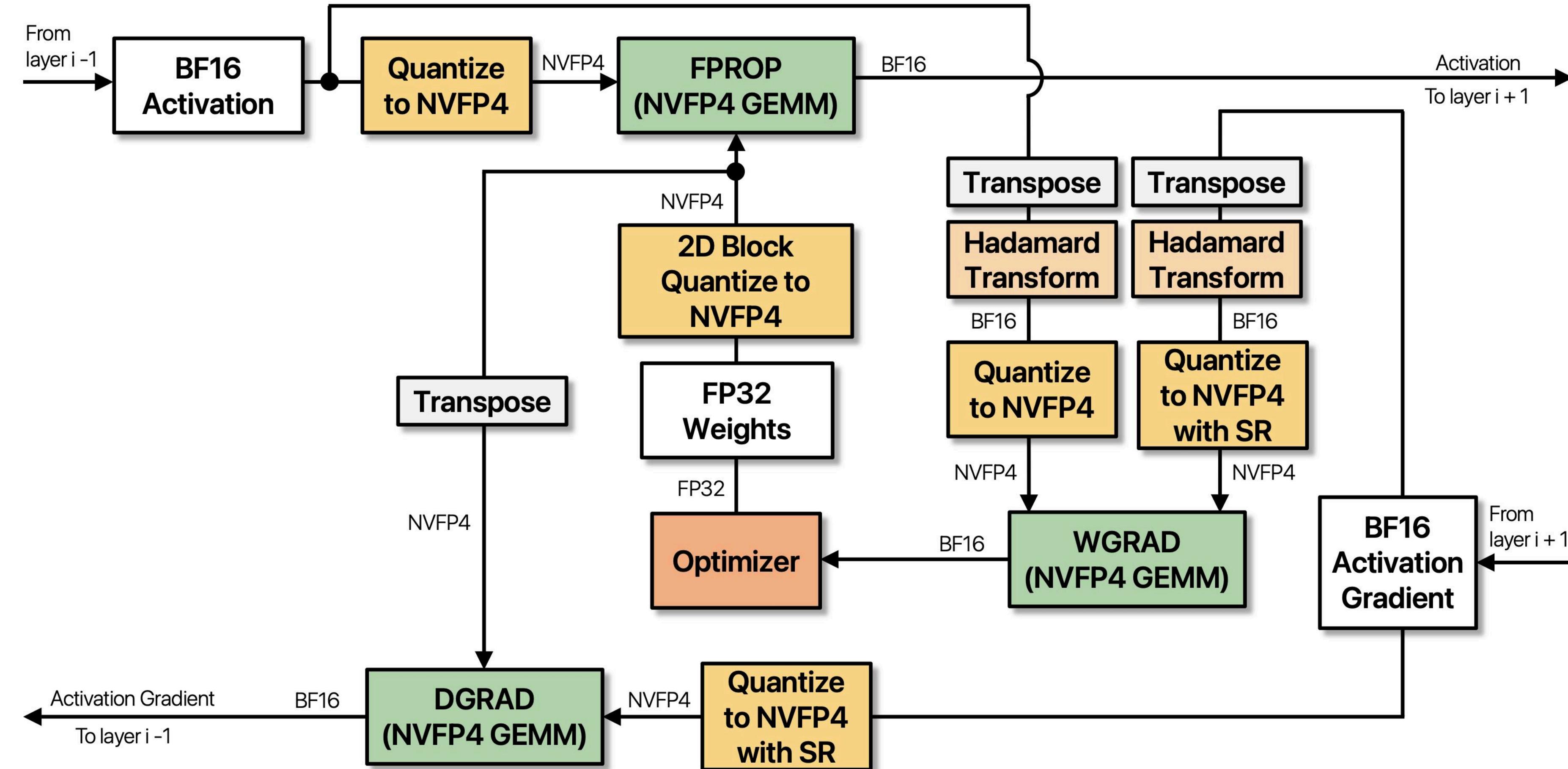
NVFP4 Recipe

NVFP4 Recipe:

- 2D Scaling
- Stochastic Rounding
- Random Hadamard Transforms



NVIDIA NVFP4 Recipe



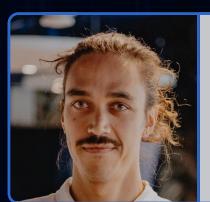
Thanks for listening!

Use the code:

EURIPS-2025

(50€ credit)

Contacts:



Paul Chang
ML Engineer



Riccardo Mereu
ML Engineer

DeepSeek V3

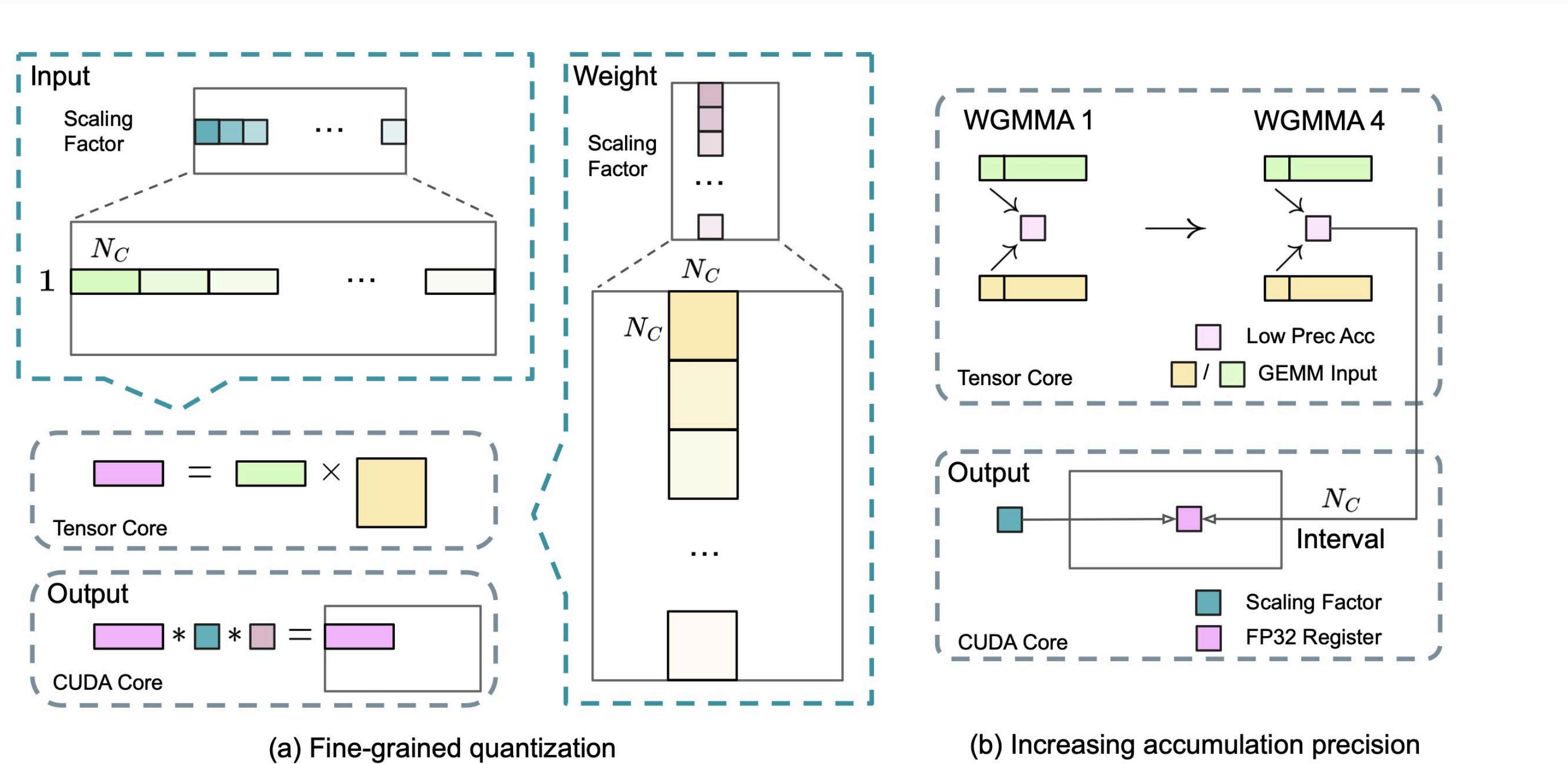


Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval of $N_C = 128$ elements MMA for the high-precision accumulation.