



MLOps Meet-up Helsinki

# Unlocking NVFP4: Low Precision Numerics on NVIDIA Blackwell

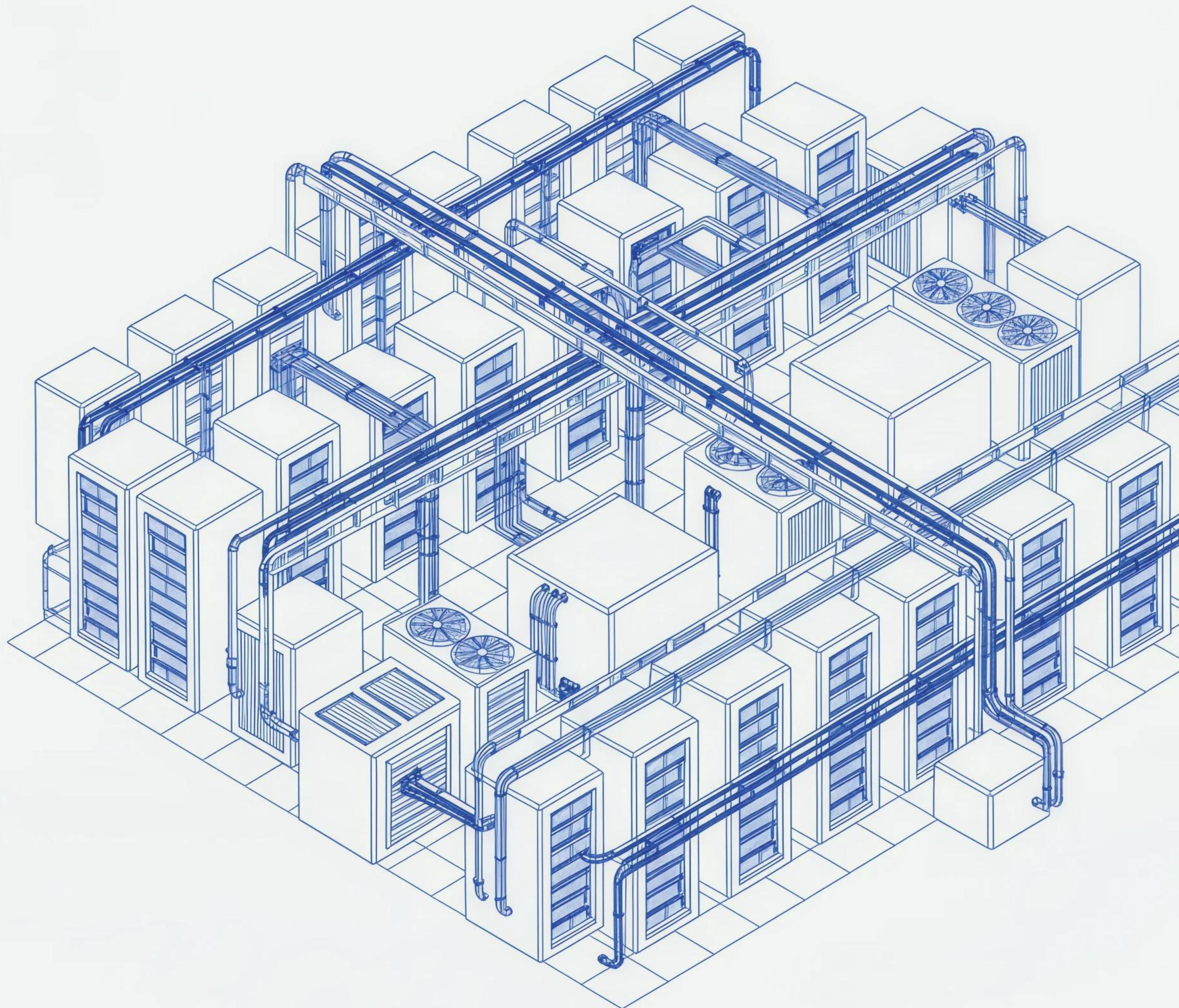
Speaker:



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

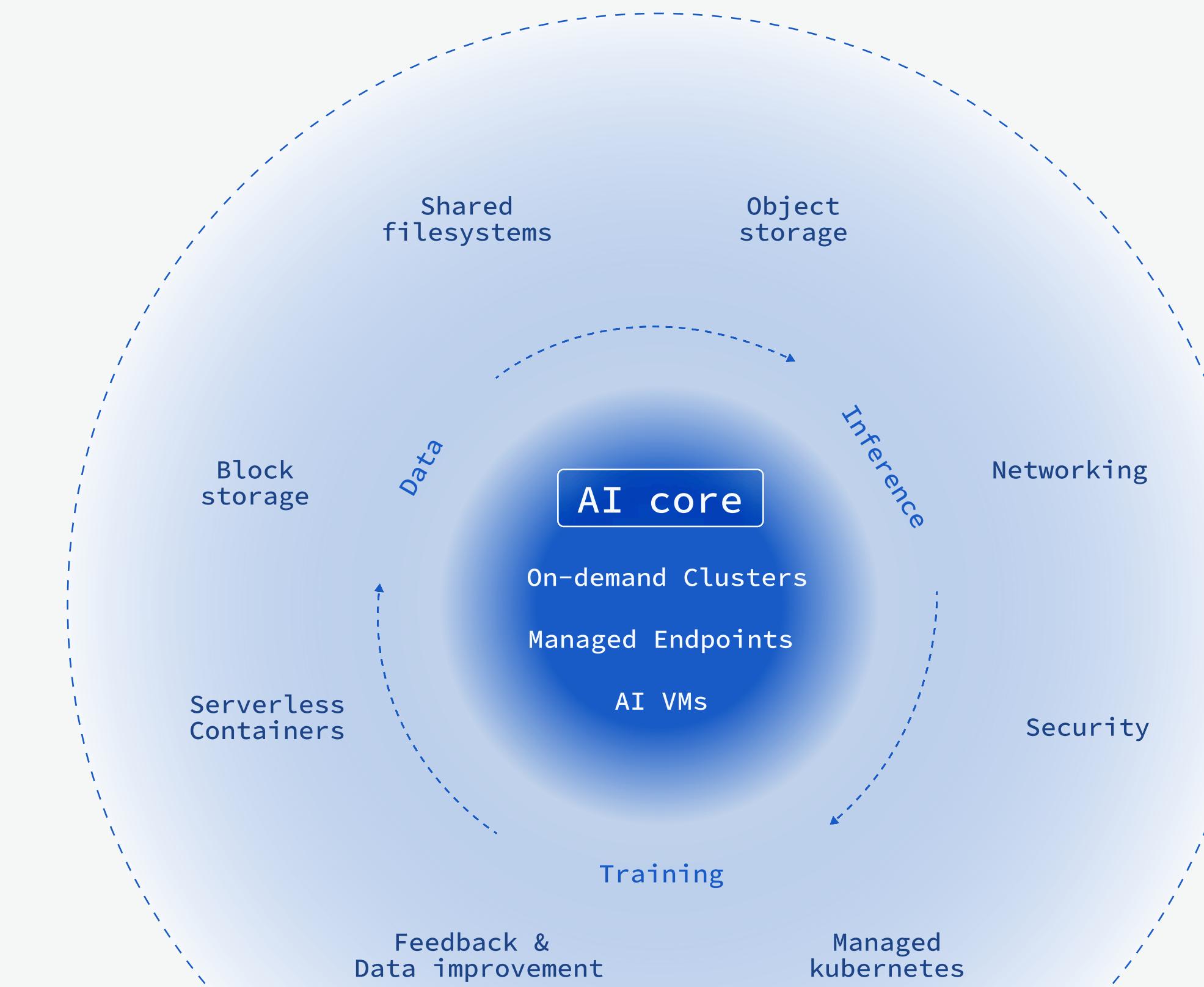
Redesigned from scratch for GPU workloads, beating legacy piecemeal cloud architecture

## Unified control plane

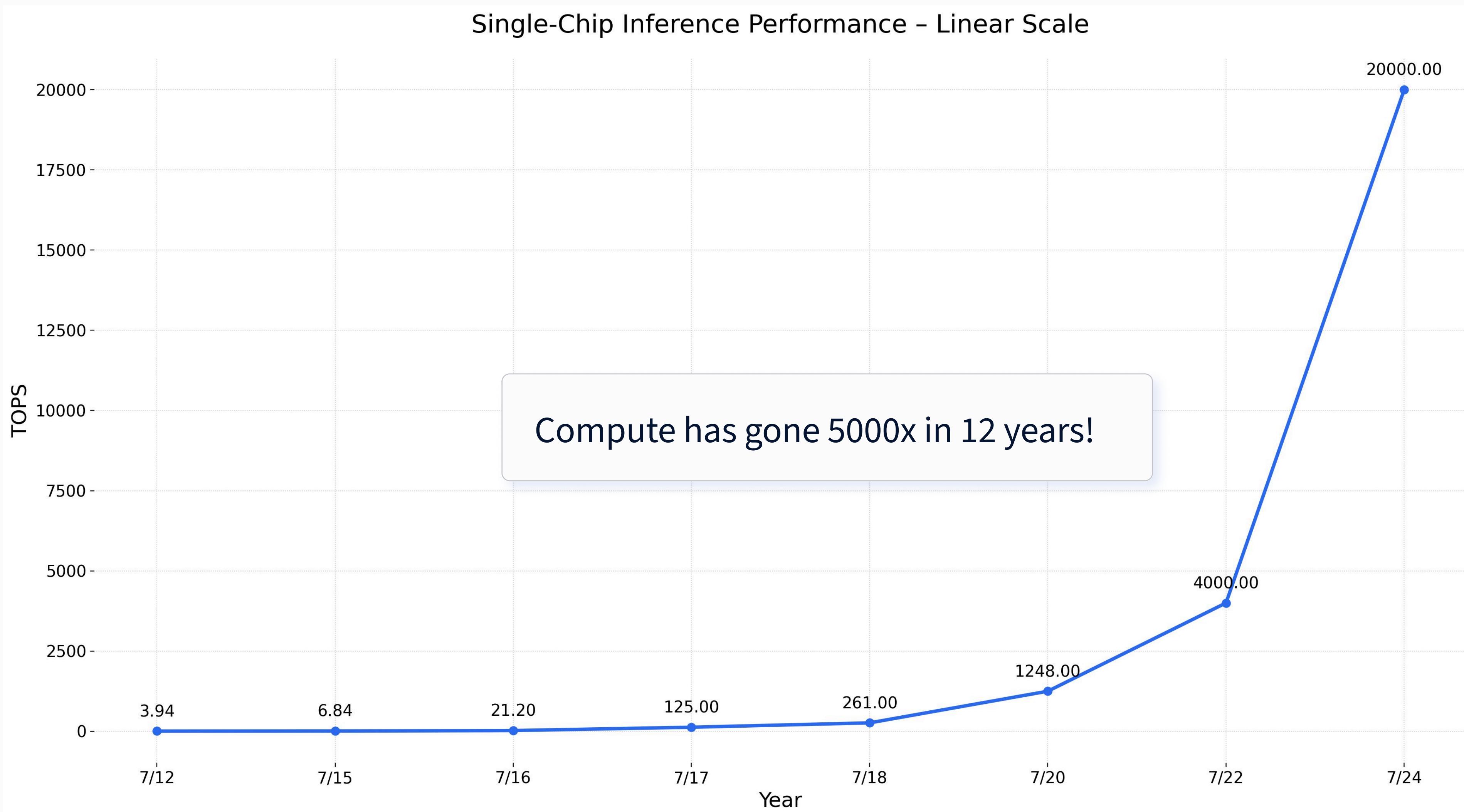
One cohesive compute layer, beating legacy fragmented layer system on utilisation and developer experience

## Deep AI framework integration

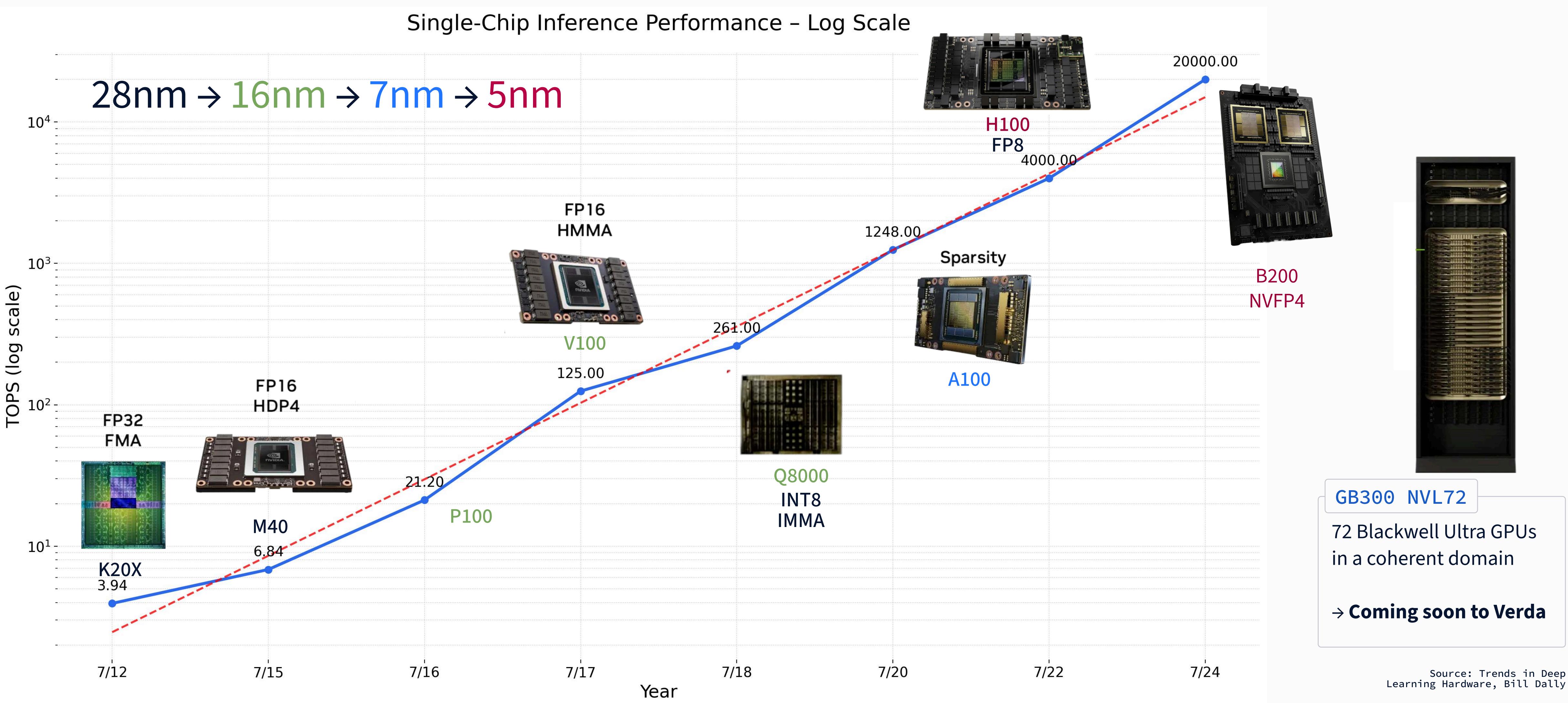
Natively optimised for PyTorch, JAX and more, enabling faster adoption of new AI tooling



# Why low precision matters for a cloud company?



# How did compute go 5000x?



# Where does the increases come from?

## Number Representation

FP32 → BF16 → FP8 → **FP4**

**32x (64x w/ sparsity)**

## Amortised Instructions (Tensor cores)

FMA → HMMA → WGMMA → UMMA

**12.5x**

## Process (fab size)

28nm → 16nm → 7nm → 5nm

**~2.5x (more transistors per mm<sup>2</sup>)**

## HW-SW co-design

FA, FA2, FA3, FA4

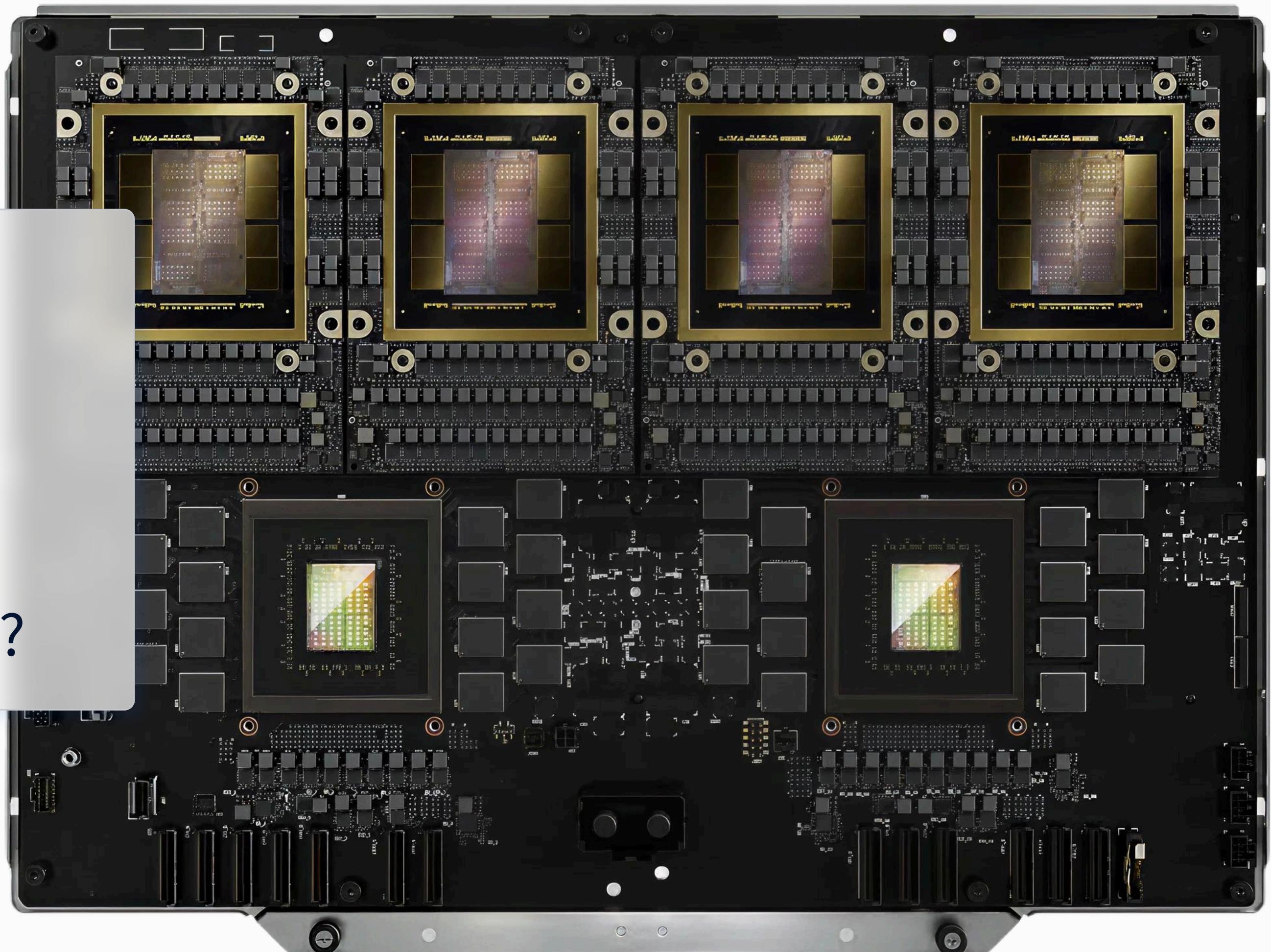
Section

# Low Precision Numerics

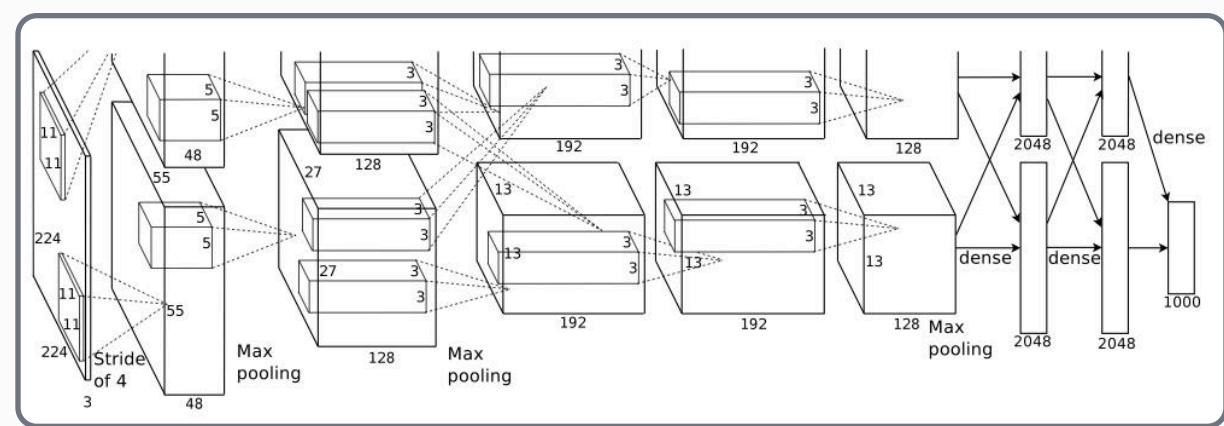
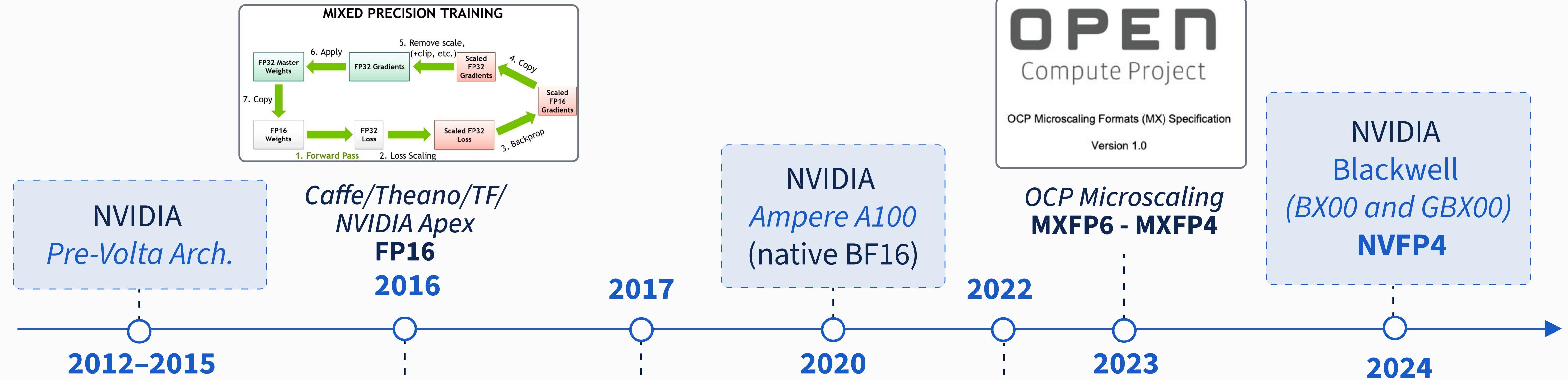
# Low Precision Numerics on NVIDIA Blackwell

## Outline

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



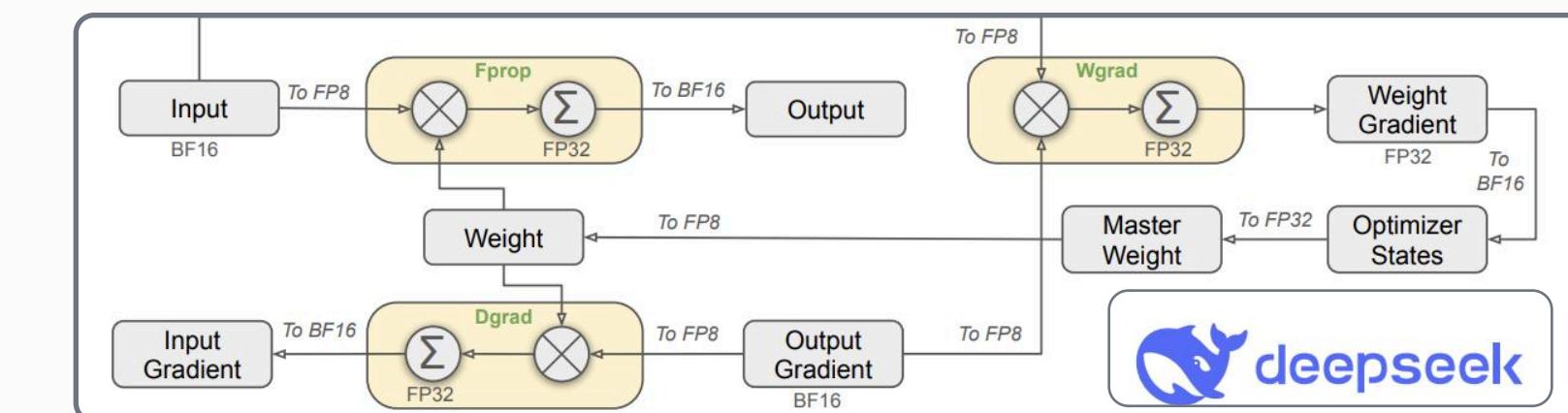
# Timeline: Floating Point in AI



**Jeff Dean** @JeffDean

The bfloat16 format was something we identified and started using during the early days of TensorFlow and also in its previous system, DistBelief. The format was basically the IEEE fp32 format with 16 bits of mantissa removed.

9:12 PM · May 31, 2019



# Why FP4 Training?

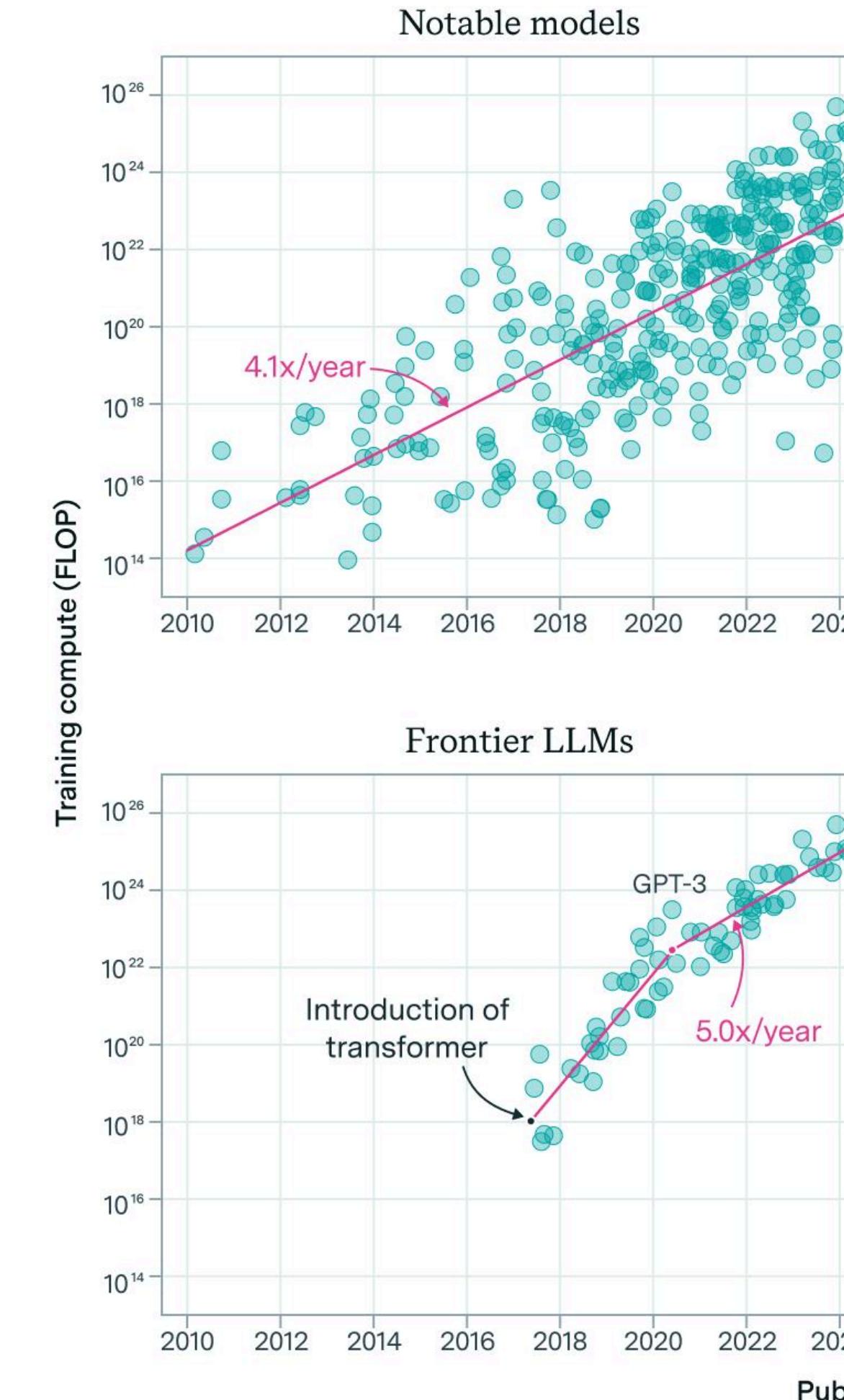
**Training compute for frontier LLMs is growing at **4-5x/year**, demanding:**

- smaller data formats
- high throughput

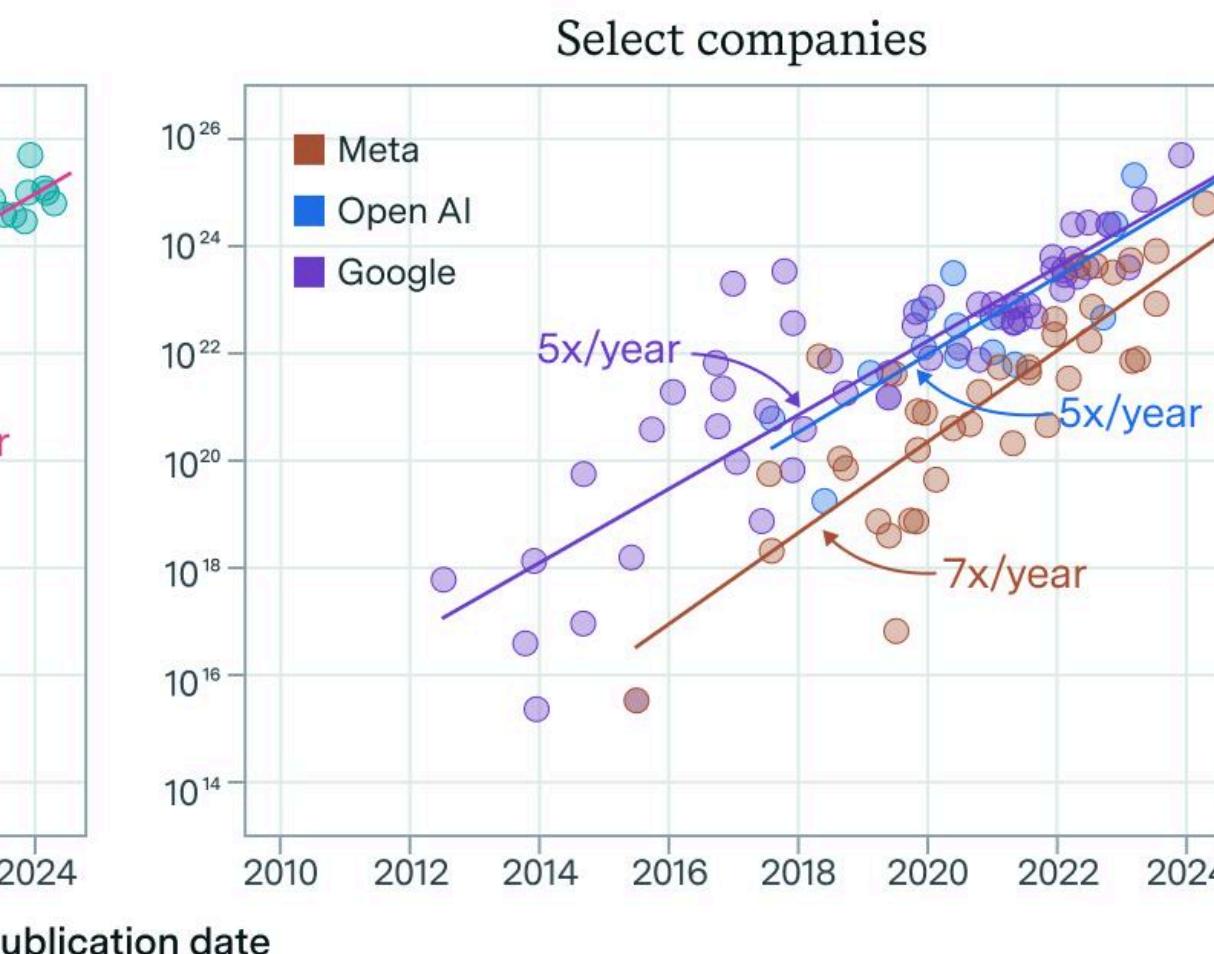
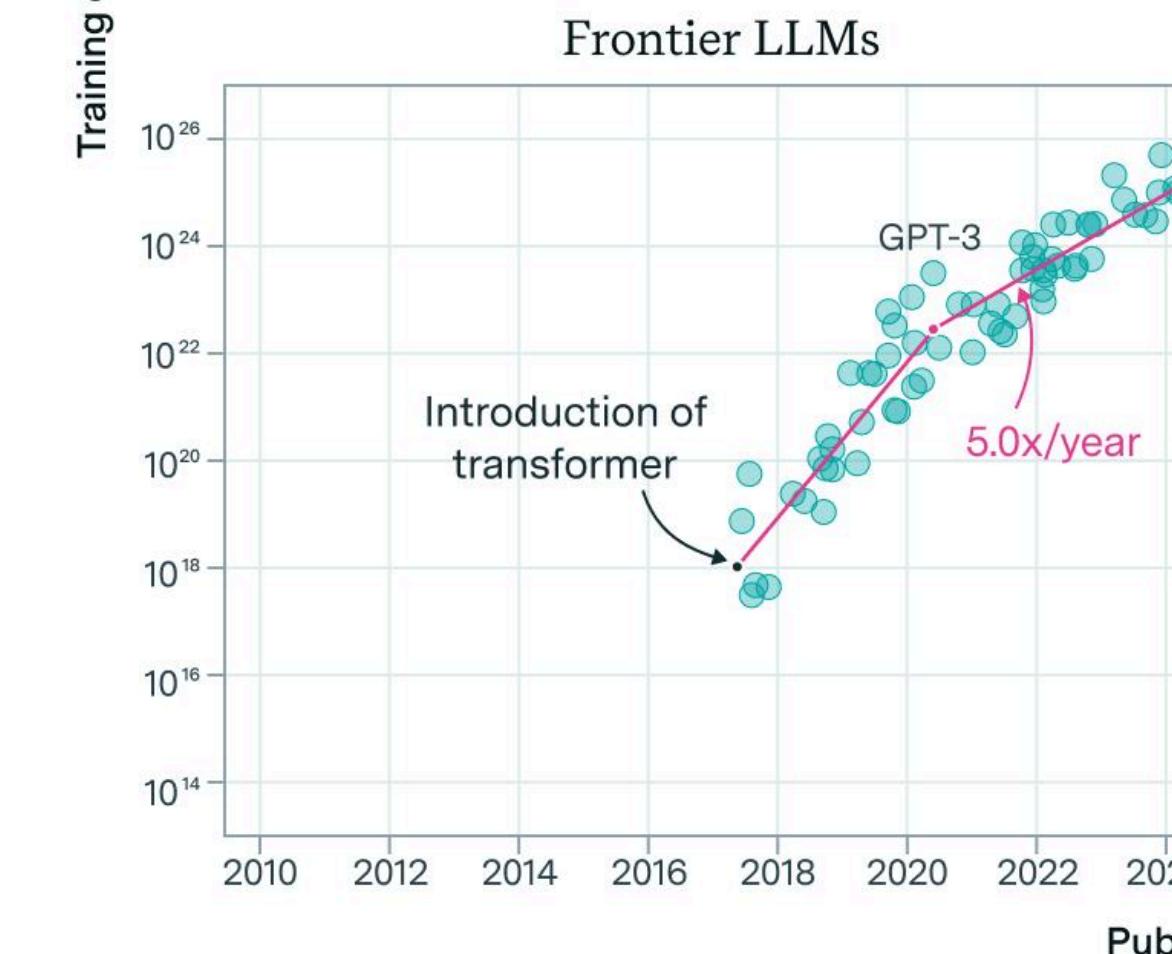
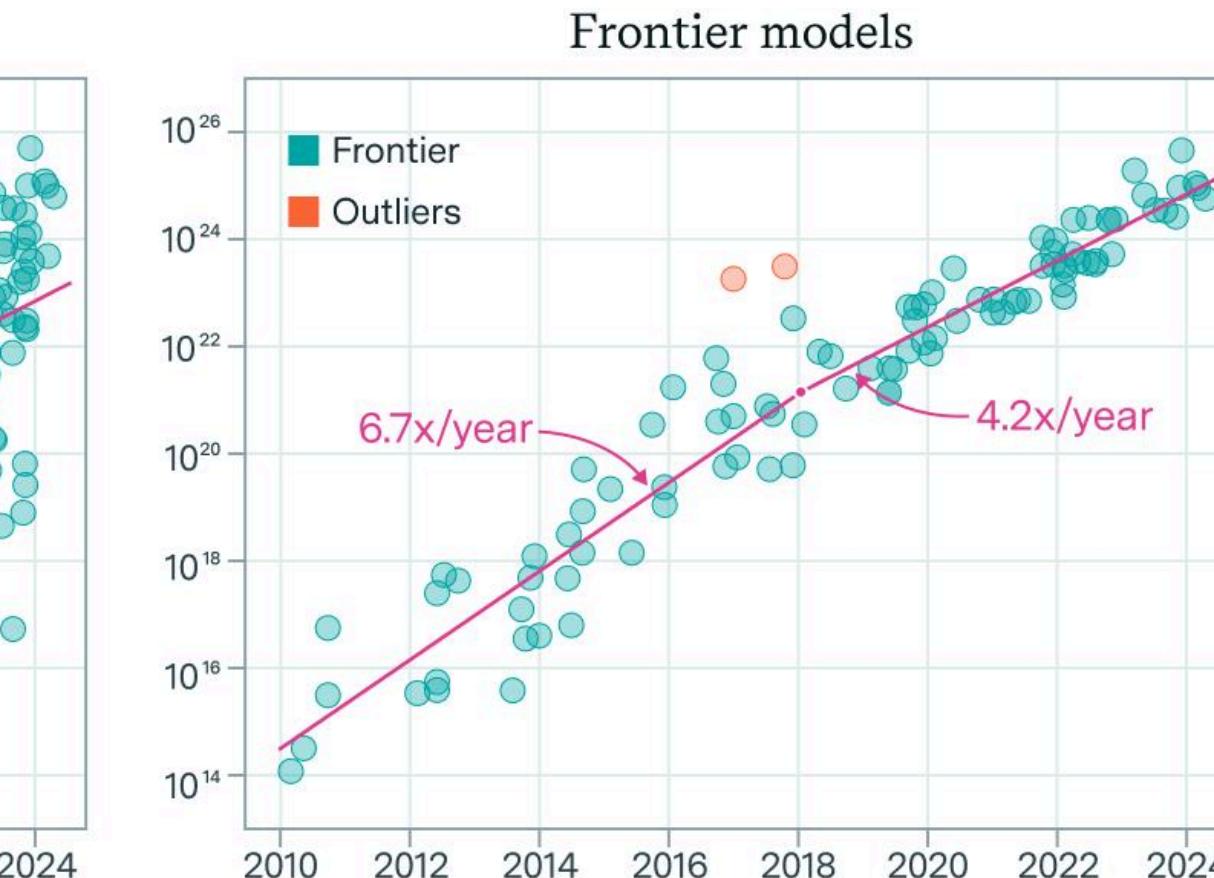
**HW evolution:**

- FP8 (4th gen tensor cores):  
~1.5-2x throughput vs. FP16
- FP4 (5th gen tensor cores):  
4x throughput vs. FP16

Summary of compute trends in AI



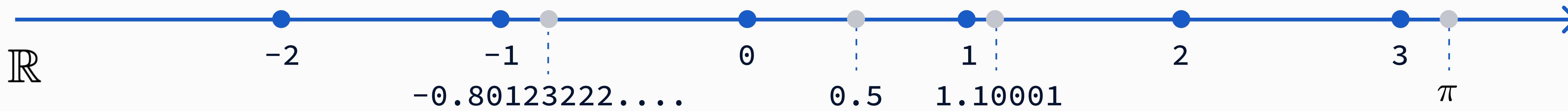
EPOCH AI



# 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)

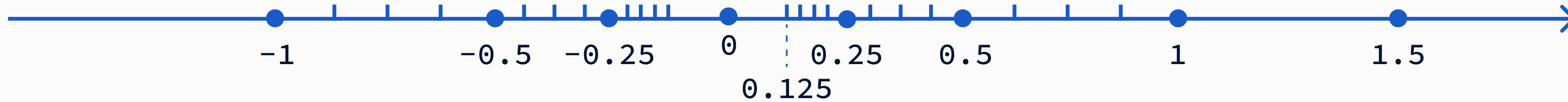
## Need for number representation to sample the real line:

- 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
- **Mantissa:** 10 bits
- **Exponent bias:**  $2^{5-1} - 1 = 15$

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

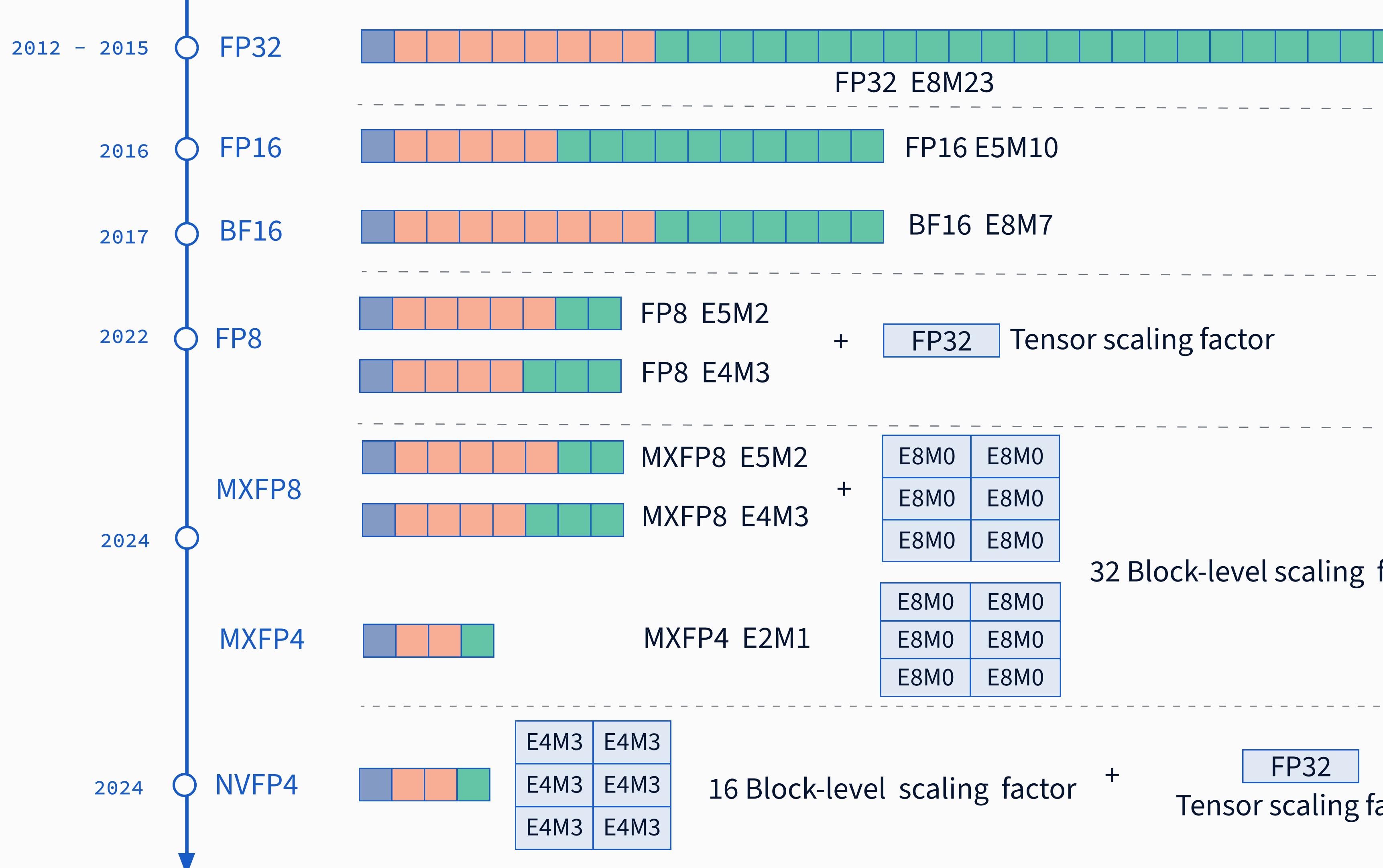
## FP 16 bits:

$$0.10001.101000000 = (6.510)_{10}$$

$$\begin{aligned}(-1)^0 \times 1.101_2 \times 2^{10001_2 - 15} &= (1 \times (1 + 2^{-1} + 2^{-3}) \times 2^2)_{10} = \\&= (4 \times 1.625)_{10} = (6.510)_{10}\end{aligned}$$

# Floating Point in AI

14  
█ Sign   █ Exponent   █ Mantissa



# 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

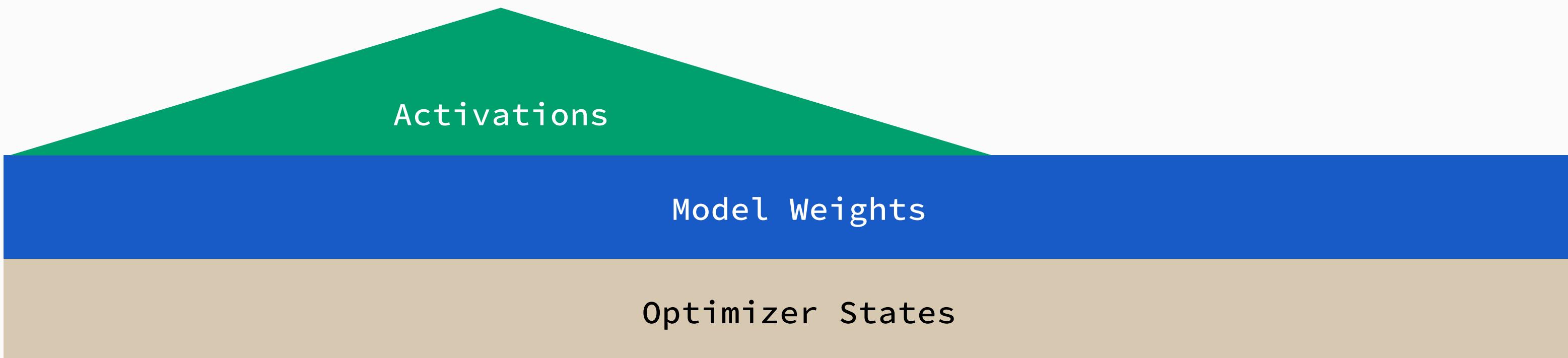


# Model Quantization Algorithms

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

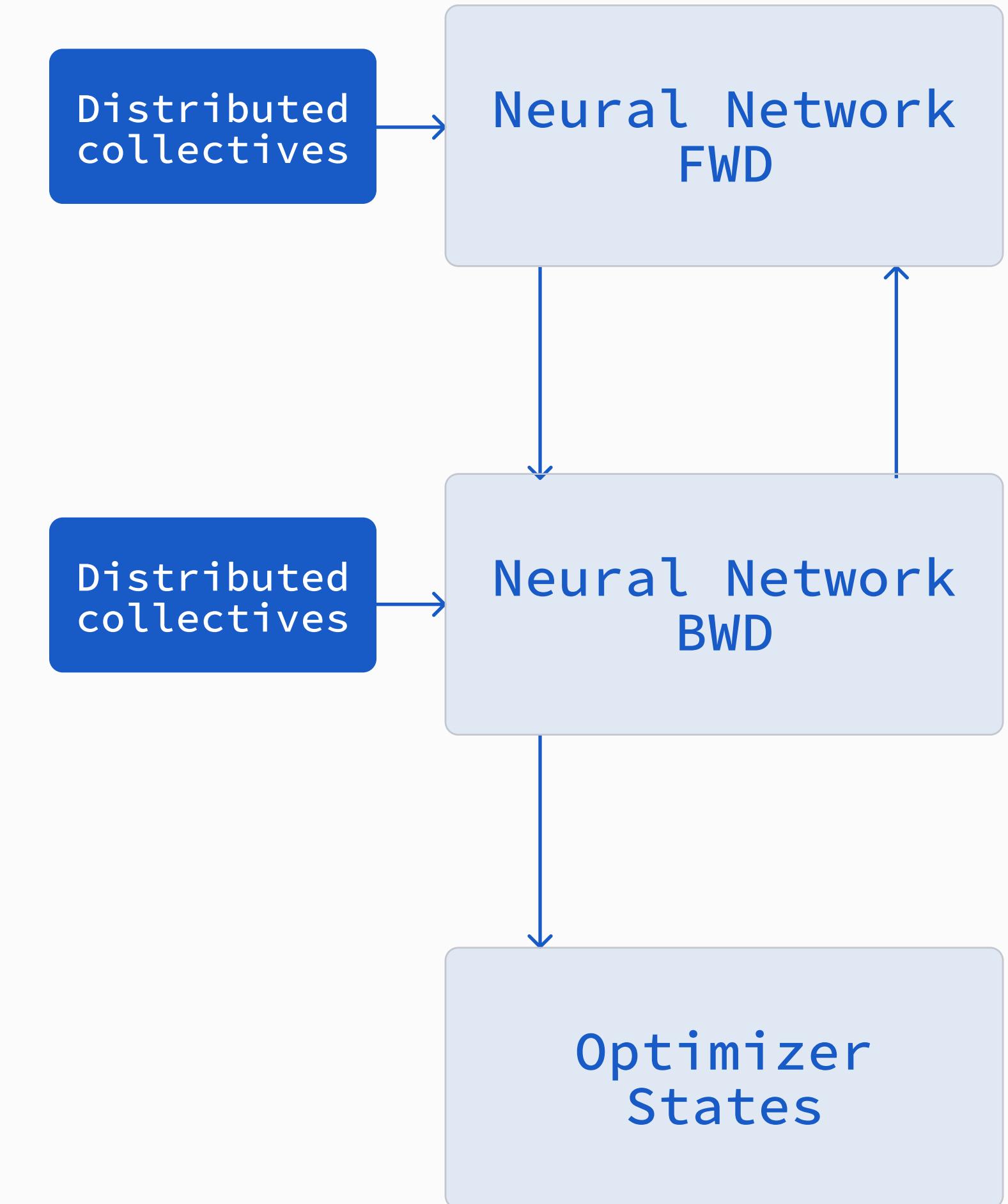
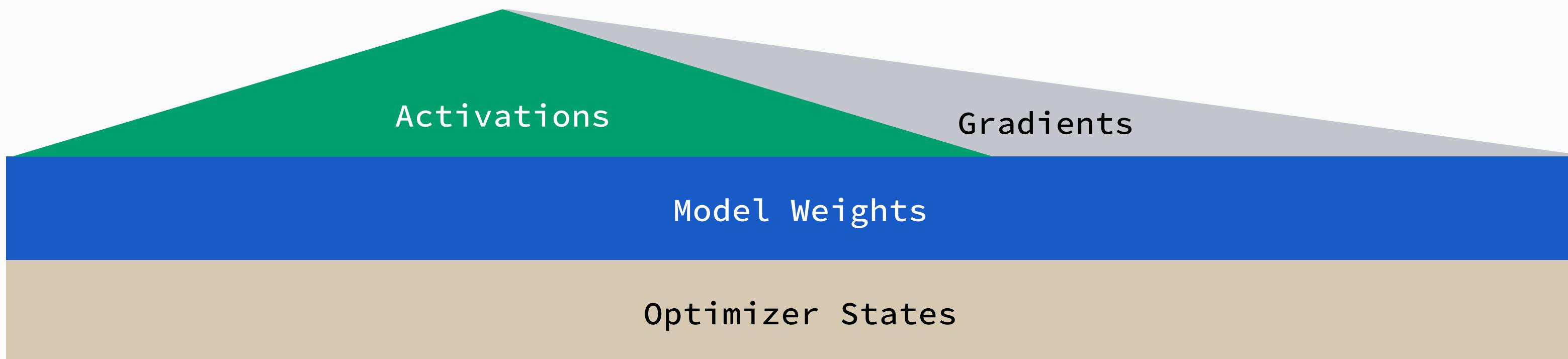
Neural Network  
FWD



# Model Quantization Algorithms

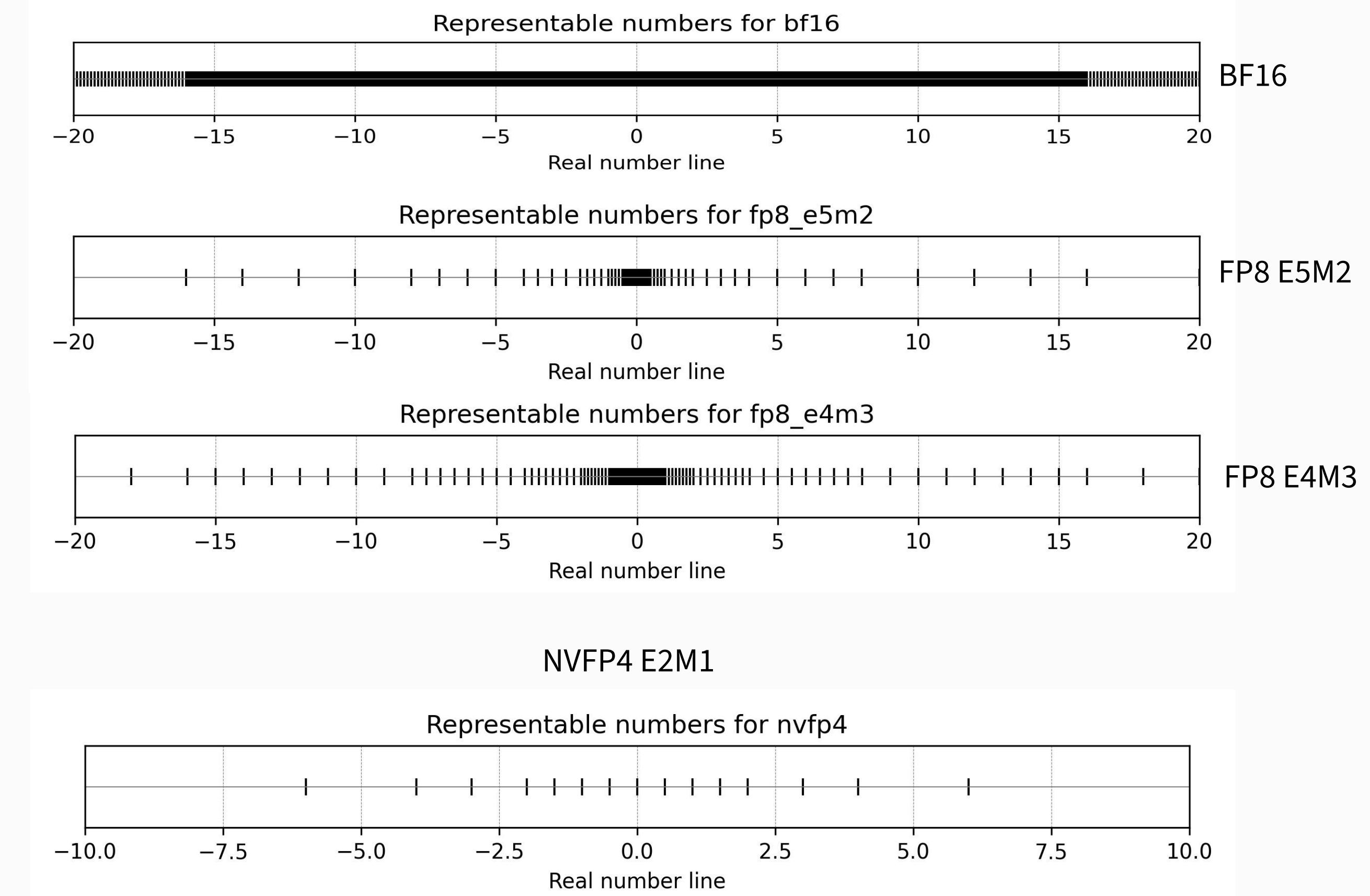
There are several strategies to quantize neural networks:

- **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



# NVFP4: Are 16 values enough?

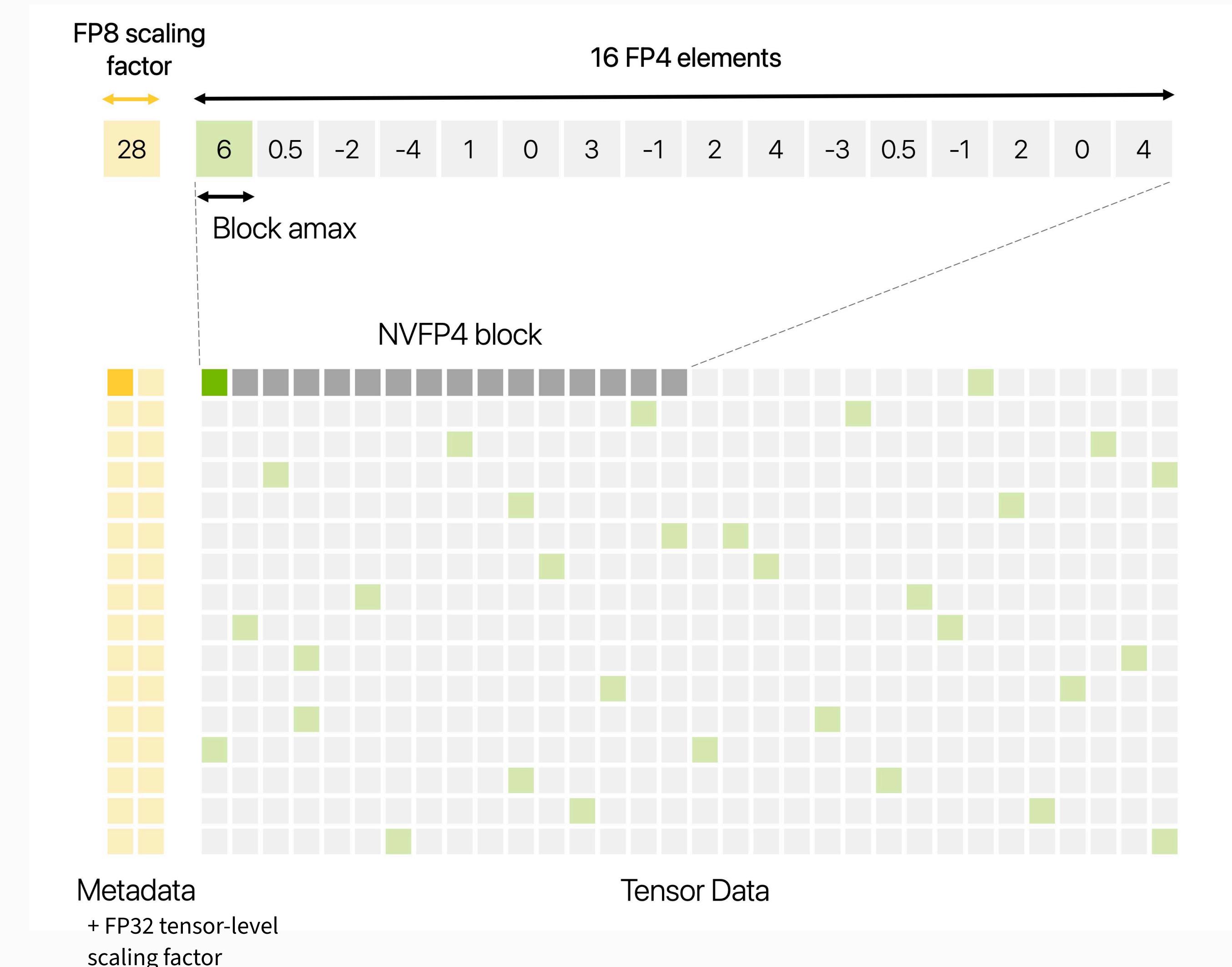
sign	exponent	mantissa		
0	0	0	0	= 0.0
0	0	0	1	= 0.5
0	0	1	0	= 1.0
0	0	1	1	= 1.5
0	1	0	0	= 2.0
0	1	0	1	= 3.0
0	1	1	0	= 4.0
0	1	1	1	= 6.0
1	0	0	0	= -0.0
1	0	0	1	= -0.5
1	0	1	0	= -1.0
1	0	1	1	= -1.5
1	1	0	0	= -2.0
1	1	0	1	= -3.0
1	1	1	0	= -4.0
1	1	1	1	= -6.0



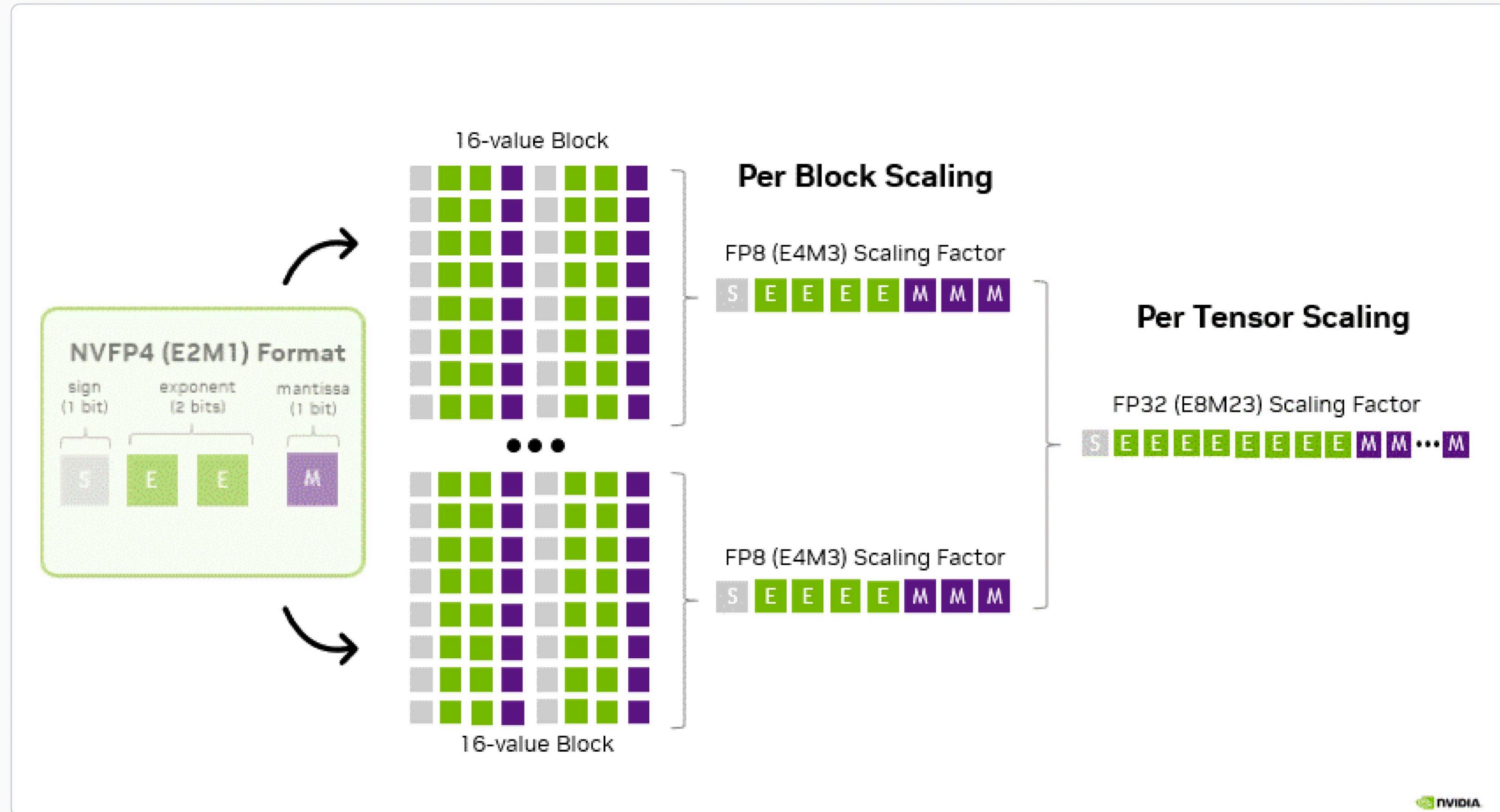
# NVFP4 vs. MXFP4

The smaller block size and more precise scaling allows:

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



# NVIDIA NVFP4

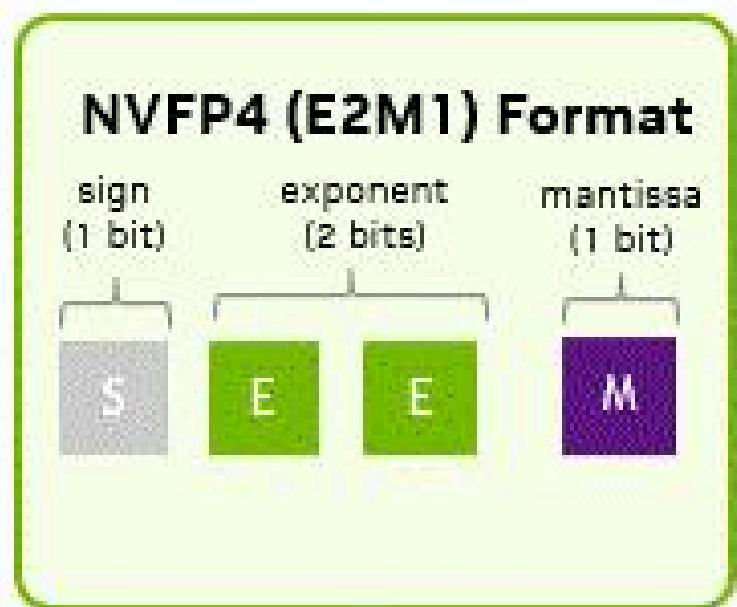


## Why two levels?

- E4M3 block scales provide fractional precision for accurate scaling, but max out at 448.
- FP32 tensor scale first shrinks values into E4M3's limited range to prevent overflow.



# NVIDIA NVFP4



## Why two levels?

- E4M3 block scales provide fractional precision for accurate scaling, but max out at 448.
- FP32 tensor scale first shrinks values into E4M3's limited range to prevent overflow.



# NVFP4 Recipe

- 2D Scaling

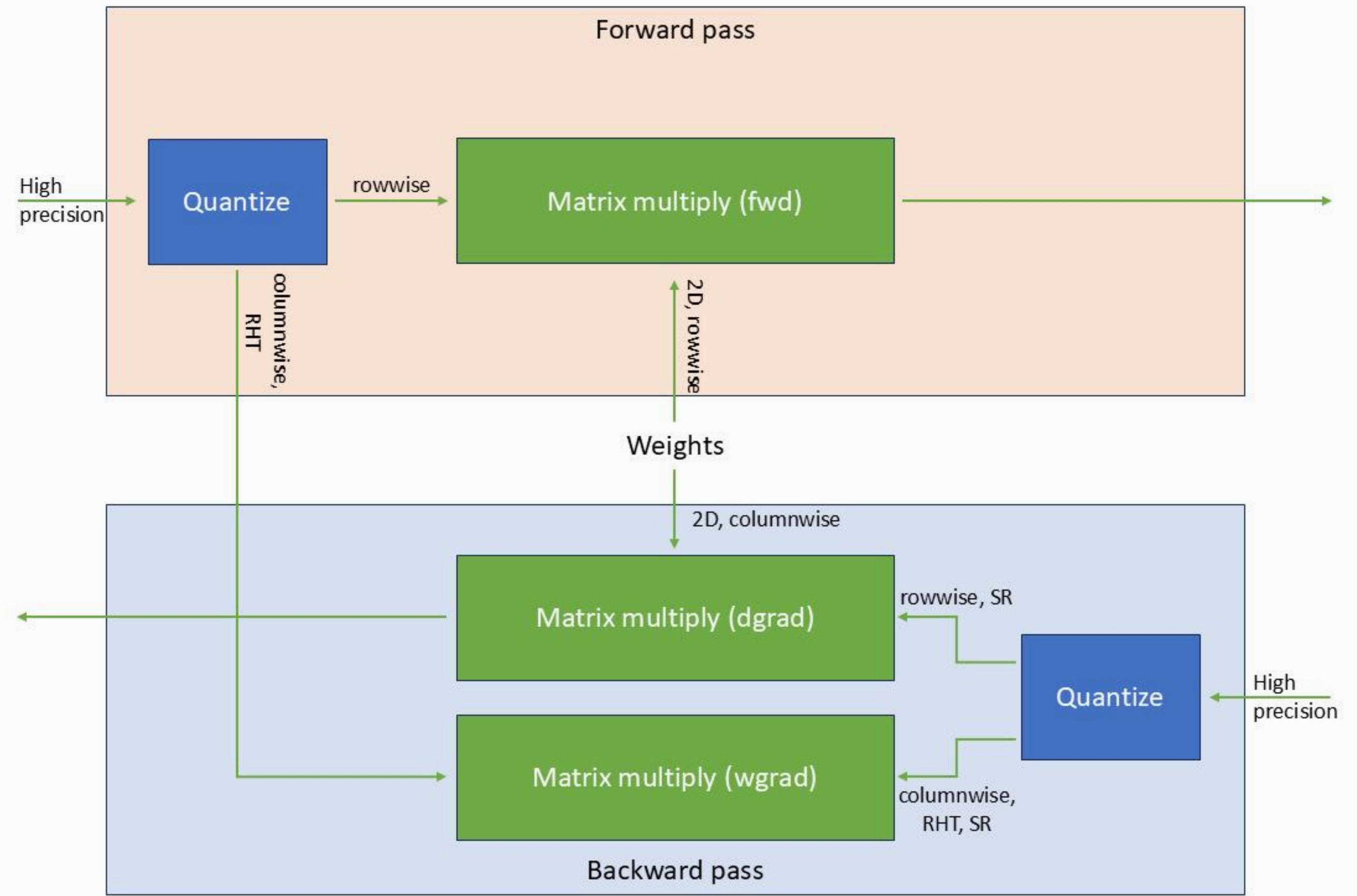
$$y = Q(x) \cdot Q(W^T) \text{ fwd pass}$$

$$\partial x = Q(\partial y) \cdot Q(W) \text{ bwd dgrad}$$

**Issue:**  $Q(W) \neq Q(W^T)$

**Solution:**

2D **16 x 16** blocks when quantizing weights



Linear Layer flow w/ NVFP4

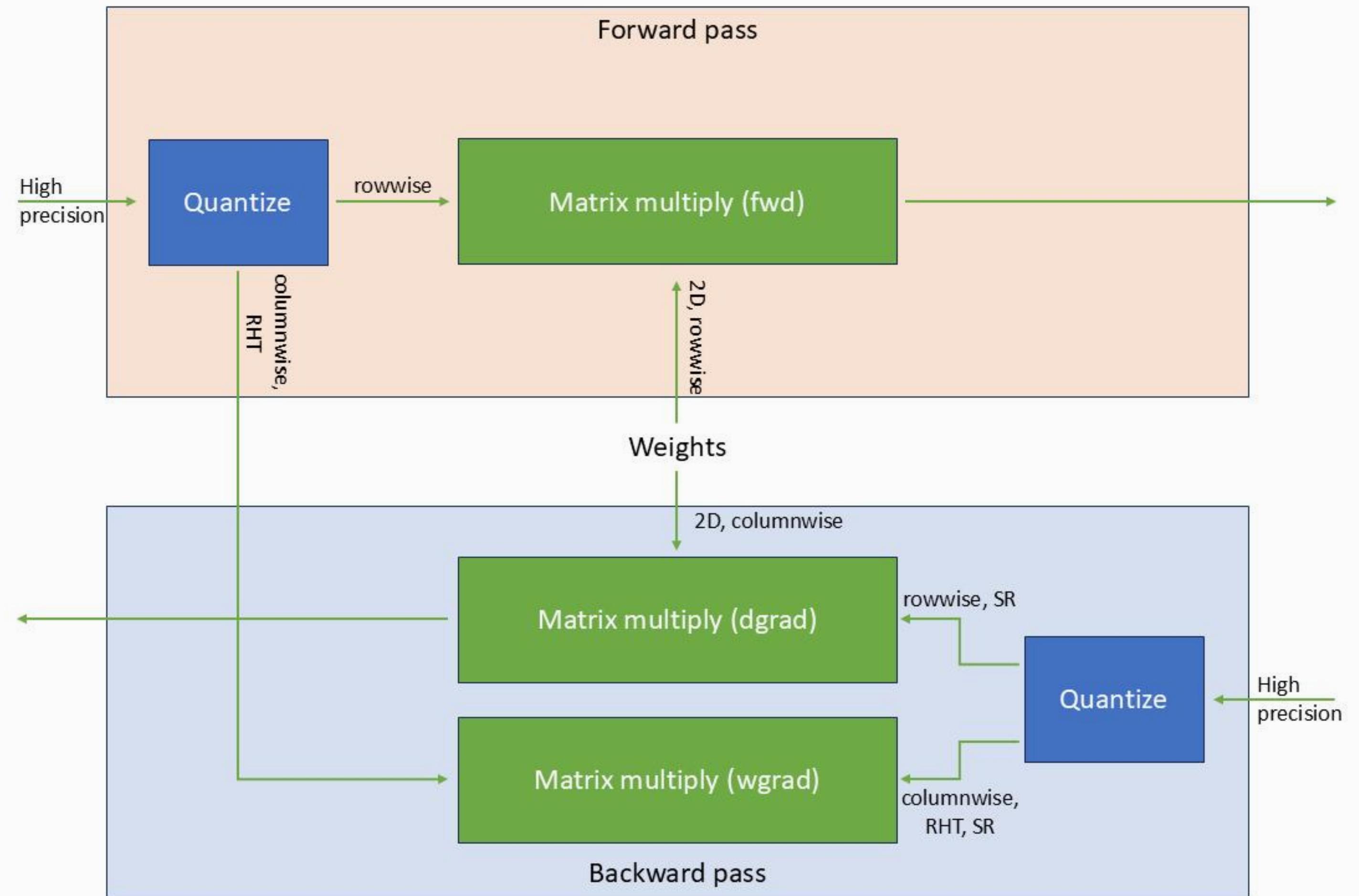
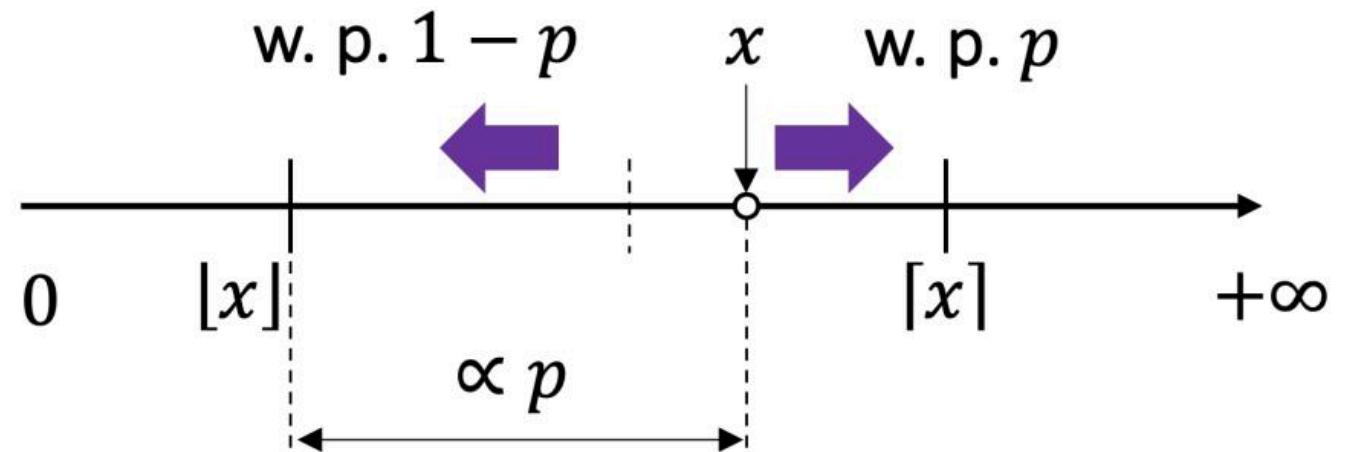
# NVFP4 Recipe

- 2D Scaling
- Stochastic Rounding

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

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

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



Linear Layer flow w/ NVFP4

# NVFP4 Recipe

- 2D Scaling
- Stochastic Rounding
- Random Hadamard Transforms:
  - Activation outliers (Wgrad inputs)

$$H^T H = I, \quad (xH)(WH)^T = XW^T$$

$$\text{RHT}(x) = H \cdot \text{diag}(s) \cdot x$$

where  $s \in \{-1, +1\}^n$

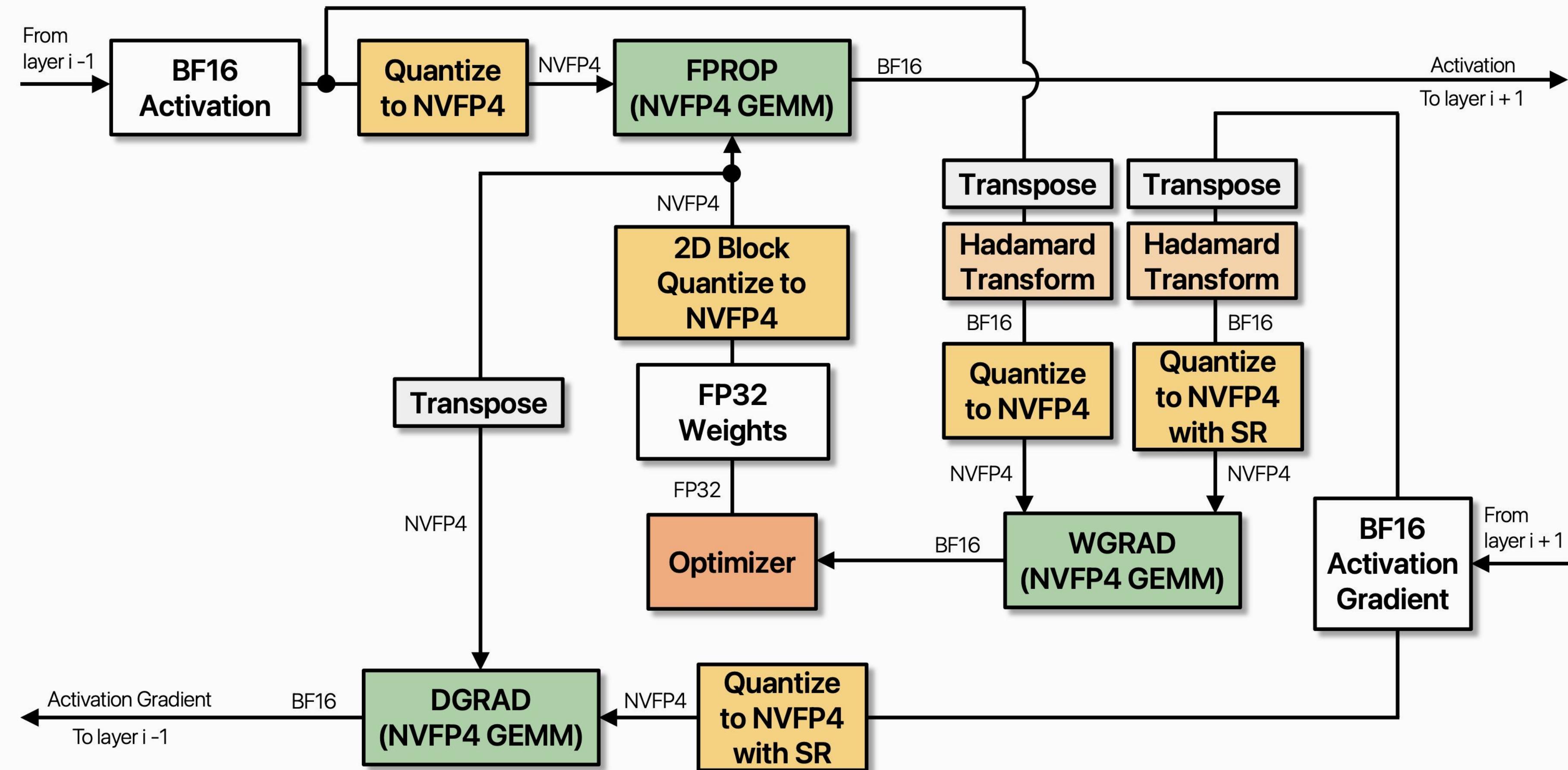


```

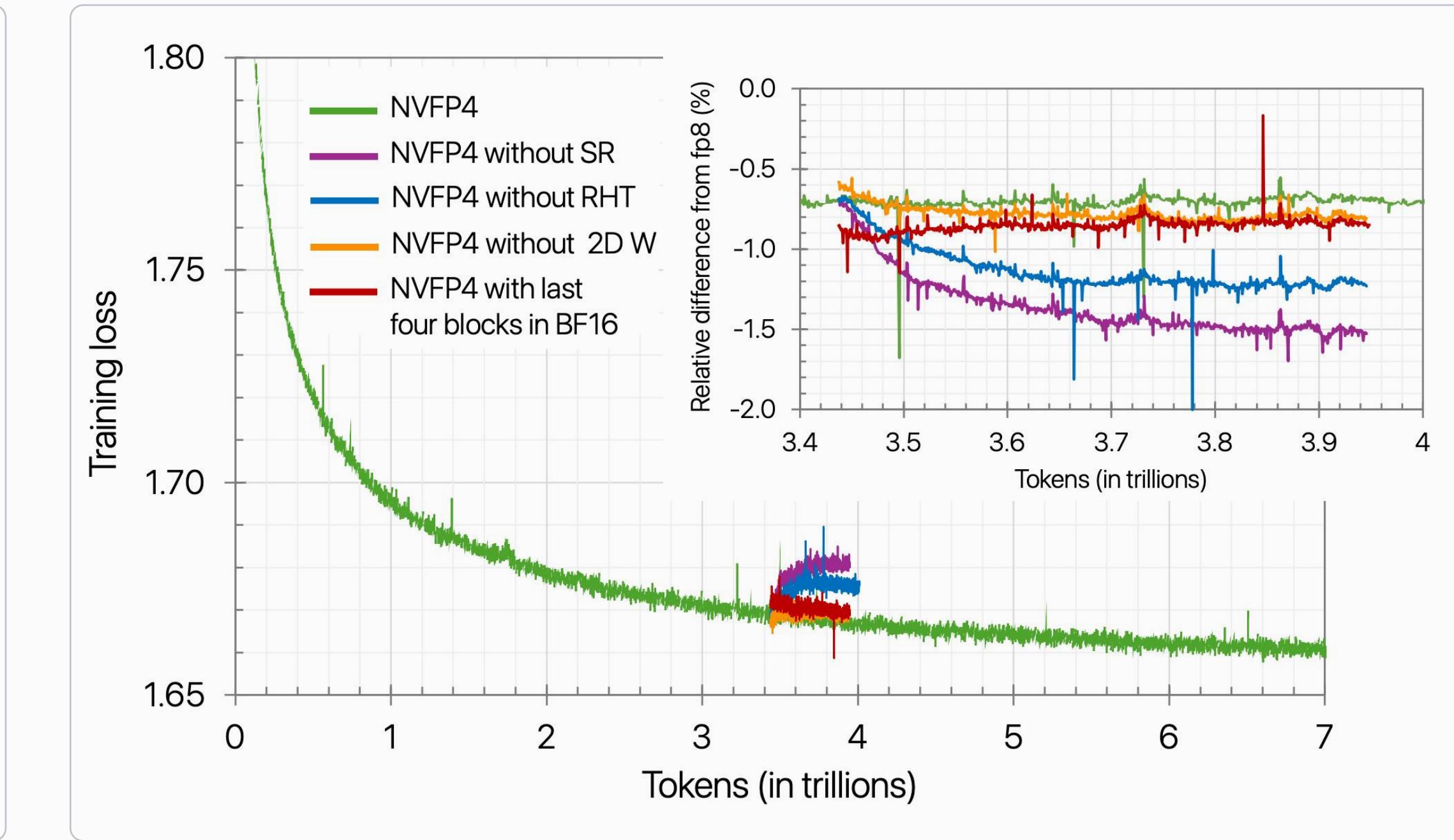
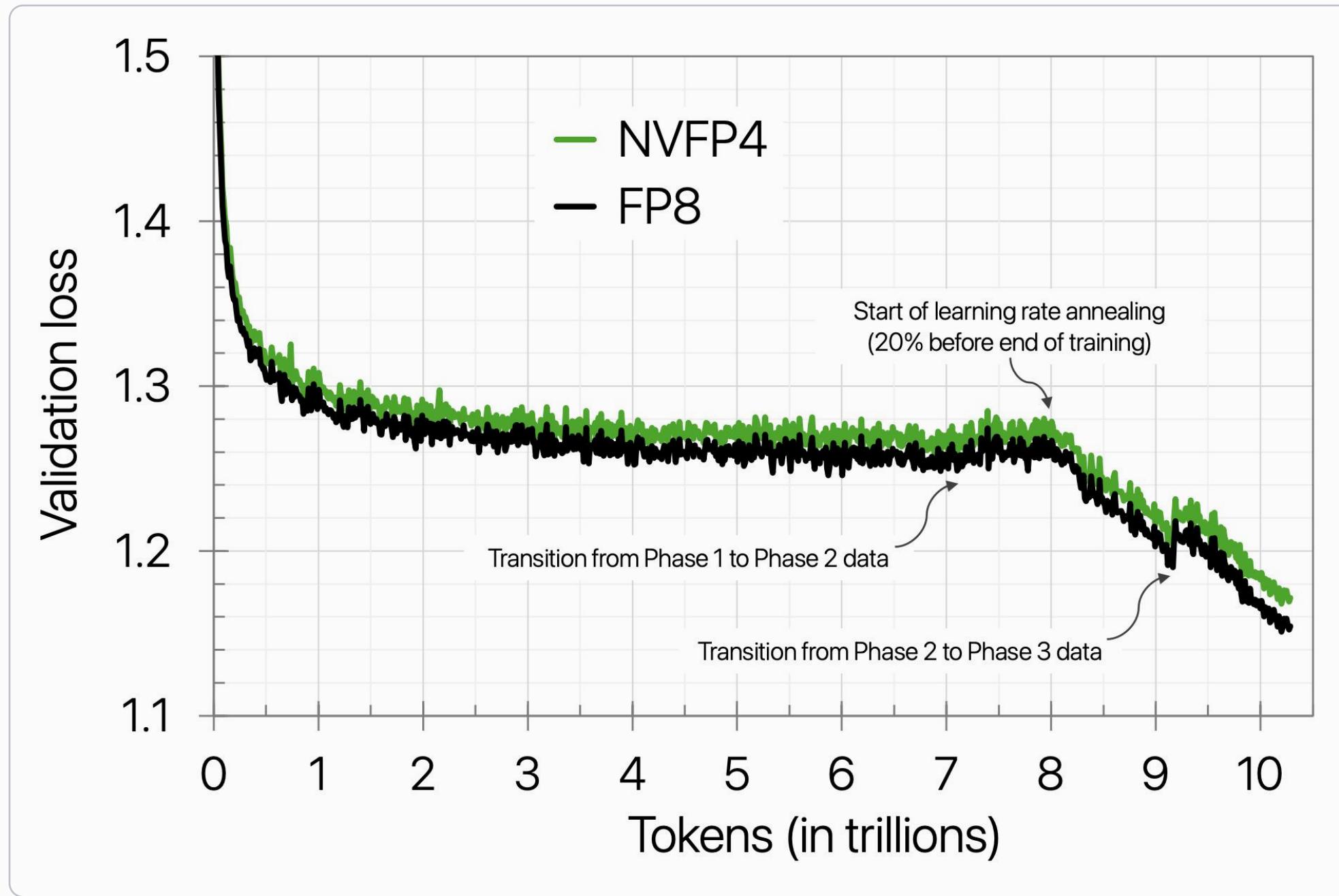
48 ... def get_wgrad_sign_vector(device: int) -> torch.Tensor:
49     """Hard-coded random signs for Hadamard transform.
50
51     https://xkcd.com/221/
52
53     """
54     return torch.tensor(
55         [1, 1, 1, -1, 1, -1, -1, -1, -1, -1, 1, -1, 1, -1, -1],
56         dtype=torch.float32,
57         device=device,
58     )

```

# NVIDIA NVFP4 Recipe



# NVIDIA NVFP4 Results



- **NVFP4 ≈ FP8:**  
Loss curves nearly identical throughout 10T token training run
- **Throughput : 2x on GB200 / 3x on GB300** vs. FP8
- **Memory: ~2x** (half the memory)

# Conclusions

Use NVFP4 when:

- You have Blackwell hardware (B300, GB300, B200, GB200, RTX PRO 6000)
- Serving large models where memory is the bottleneck
- Training frontier models and achieving 2-3x throughput

Summary of NVFP4 recipe:

GEMM	RHT	Scaling	Rounding
Fprop	N	$x: 1D, W^T: 2D$	RTN
Dgrad	N	$\partial y: 1D, W: 2D$	$\partial y: SR$
Wgrad	Y	$x^T: 1D, \partial y: 1D$	$x: RTN, \partial y: SR$

Things to keep in mind:

- Still an active area of research
- Requires complex training recipes
- Attention kernels still in development



# Thanks for listening!

Use the code:

**MLOPS-COMMUNITY-FEB26**

(200€ free credits)

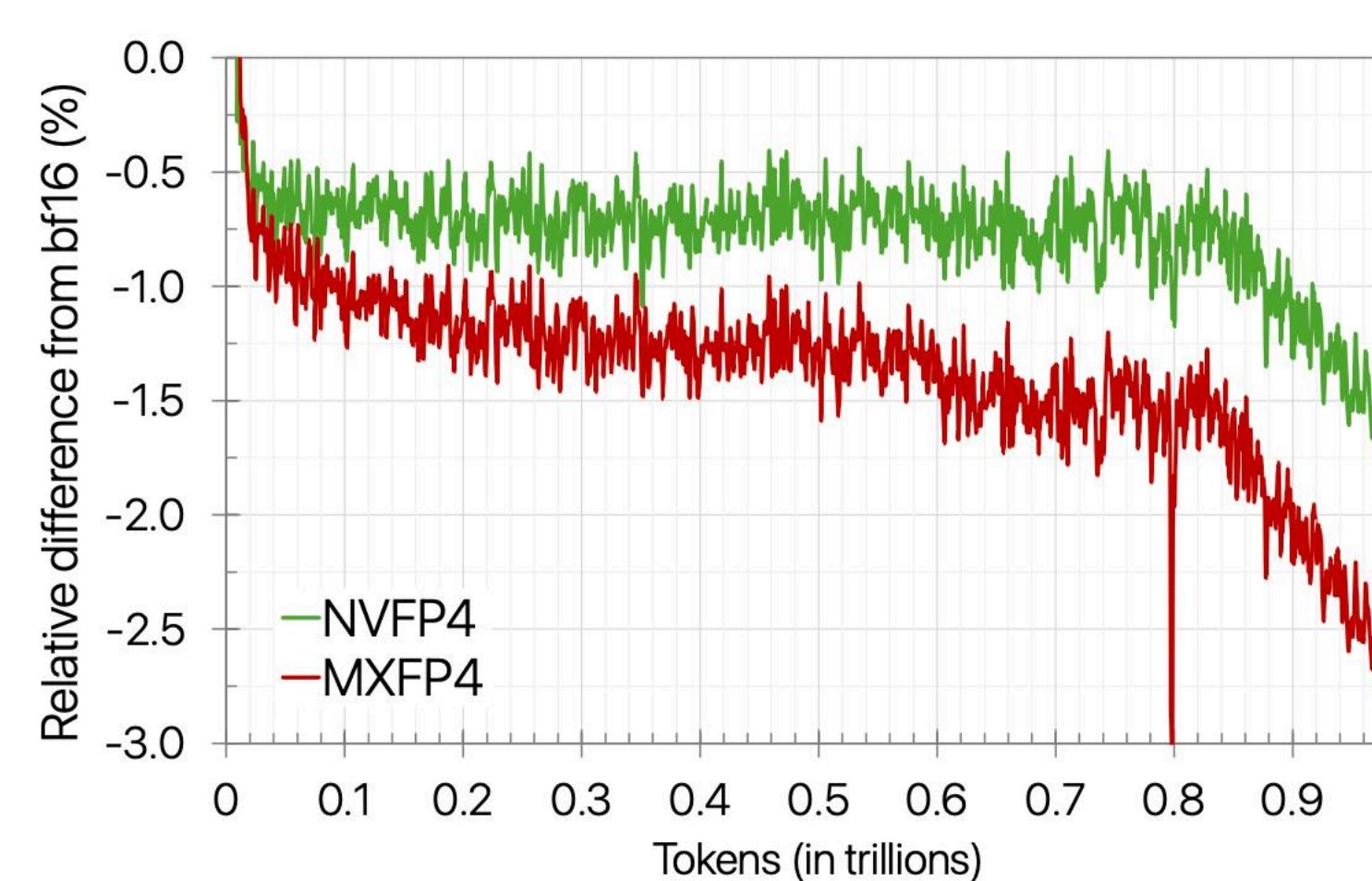
Contacts:



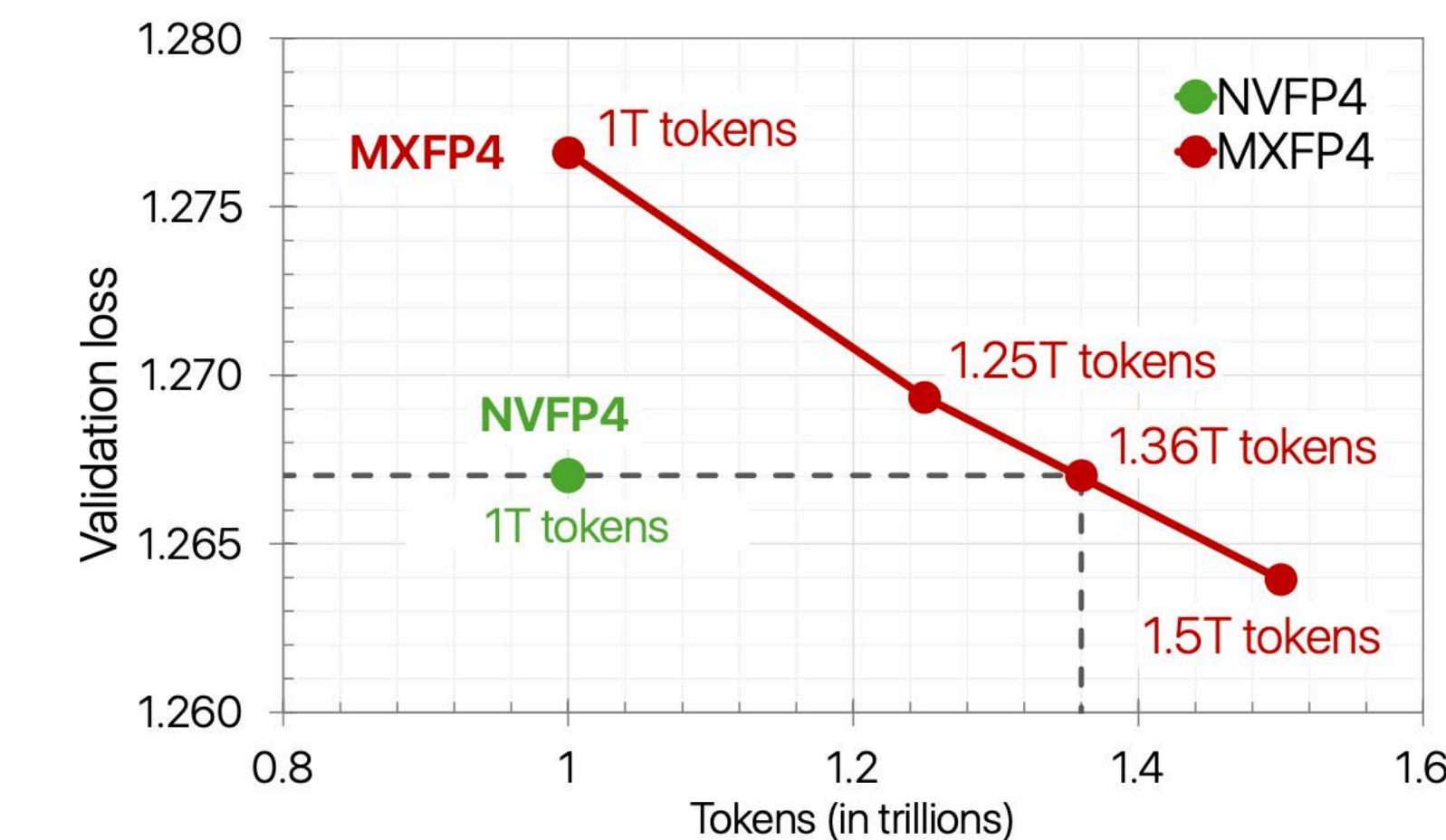
Riccardo Mereu  
[riccardo@verda.com](mailto:riccardo@verda.com)

# NVIDIA NVFP4 Results pt. 2

Format	Element	Scale	Block	Speedup vs. BF16	
				GB200	GB300
MXFP8	E5M2/E4M3	UE8M0	32	2×	2×
MXFP6	E3M2/E2M3	UE8M0	32	2×	2×
MXFP4	E2M1	UE8M0	32	4×	6×
NVFP4	E2M1	E4M3	16	4×	6×



(a) Relative difference between training loss of BF16 (baseline) and NVFP4 and MXFP4 pretraining.



(b) Final validation loss for NVFP4 and MXFP4 pretraining with different number of tokens.

Figure 6 | NVFP4 vs MXFP4 comparisons: (a) training-loss difference; (b) validation perplexity across token budgets.

# DeepSeek-V3 Mixed Precision

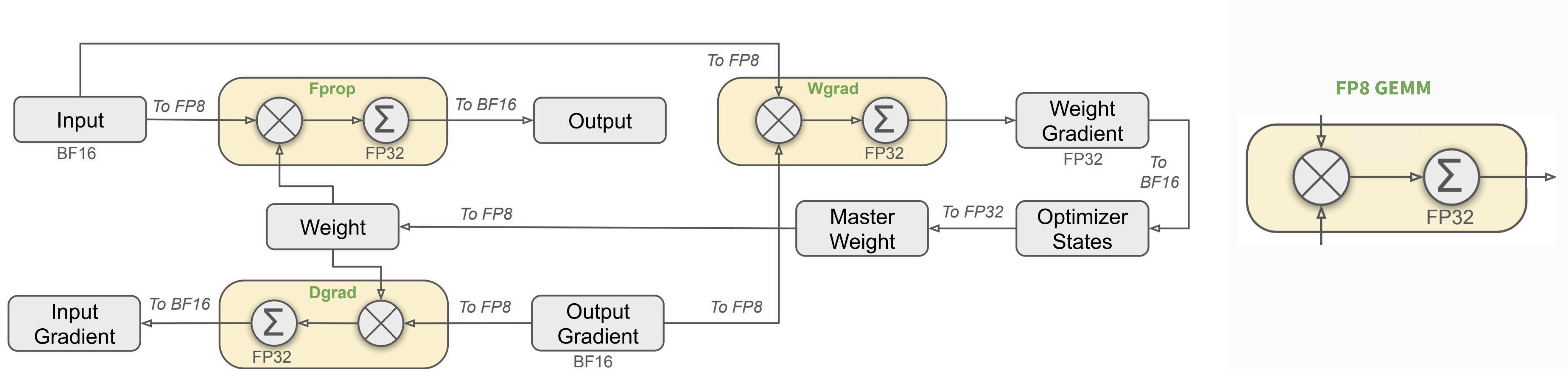


Figure 6 | The overall mixed precision framework with FP8 data format. For clarification, only the Linear operator is illustrated.

# DeepSeek V3

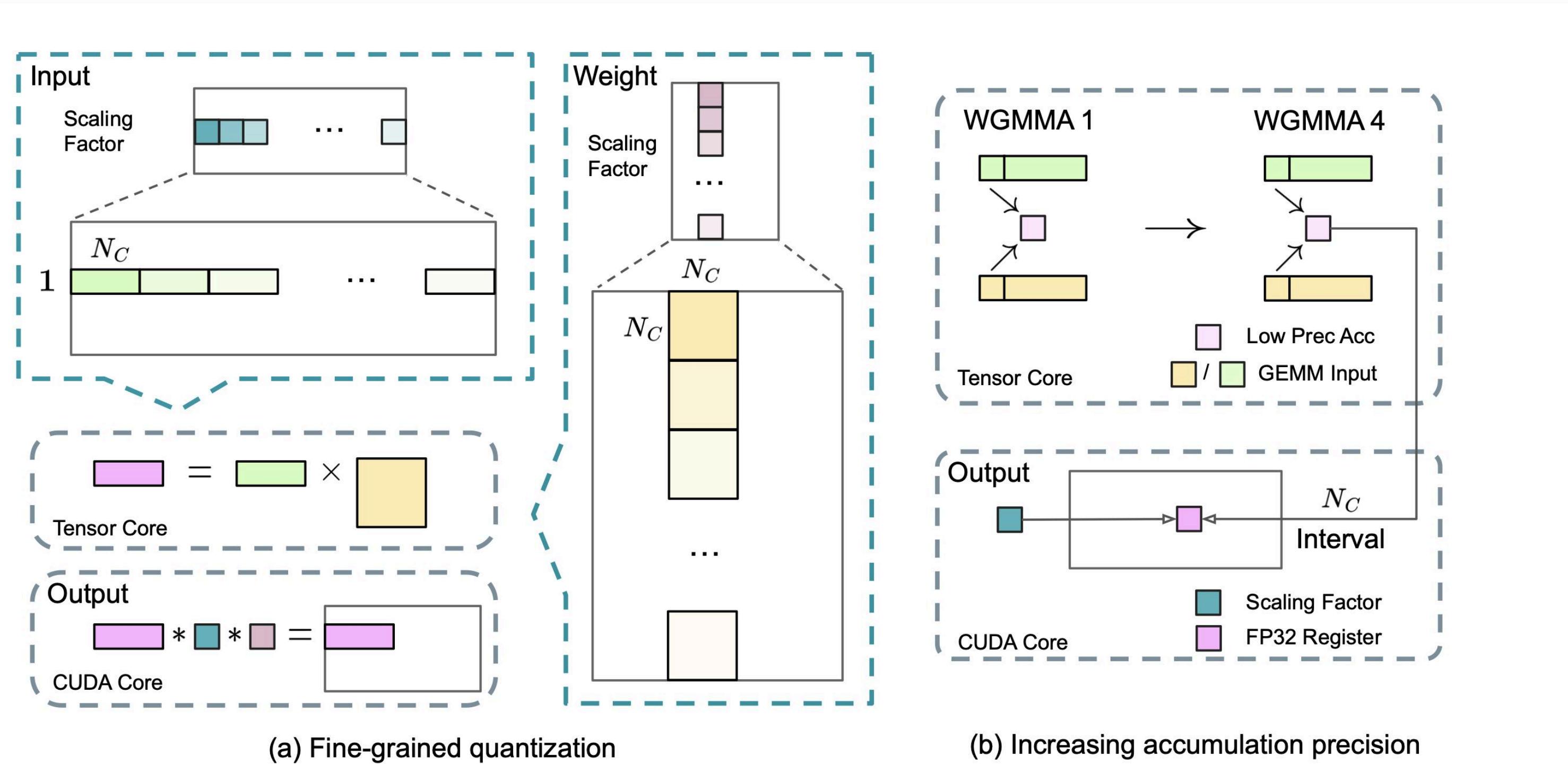


Figure 7 | (a) We propose a fine-grained quantization method to mitigate quantization errors caused by feature outliers; for illustration simplicity, only Fprop is illustrated. (b) In conjunction with our quantization strategy, we improve the FP8 GEMM precision by promoting to CUDA Cores at an interval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