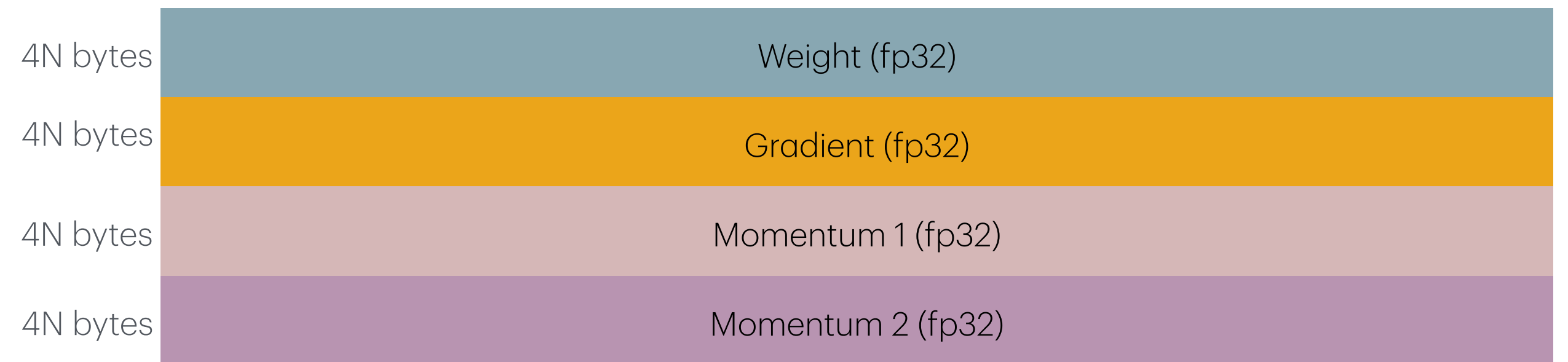


Quantization

Training large models

Memory requirements

- Without optimization
 - Model parameters: N
 - Weights: N floats
 - Gradients: N floats
 - Momentum: N floats
 - 2nd momentum (ADAM): N floats
- $16N$ bytes without counting activations



Inference in large models

Memory requirements

- Without optimization
 - Model parameters: N
 - Weights: N floats
- 4N bytes without counting activations
- 8B parameter model (i.e. Llama 3.1)
 - 32GB of memory for just weights

4N bytes

Weight (fp32)

Model	Micro-architecture	Launch	Core	Memory				
				Bus type	Bus width (bit)	Size (GB)	Clock (MT/s)	Bandwidth (GB/s)
A100 GPU accelerator (PCIe card) ^{[44][45]}	Ampere	May 14, 2020 ^[46]	1× GA100-883AA-A1	HBM2	5,120	40 or 80	1,215	1,555
A40 GPU accelerator (PCIe card) ^[43]	Ampere	October 5, 2020	1× GA102	GDDR6	384	48	7,248	695.8
L40 GPU accelerator ^[50]	Ada Lovelace	October 13, 2022	1× AD102 ^[51]	GDDR6	384	48	2,250	864
H100 GPU accelerator (PCIe card) ^[47]	Hopper	March 22, 2022 ^[48]	1× GH100 ^[49]	HBM2E	5120	80	1,000	2,039
H100 GPU accelerator (SXM card)	Hopper	March 22, 2022 ^[48]	1× GH100 ^[49]	HBM3	5,120	80	1,500	3,352

[1] https://en.wikipedia.org/wiki/List_of_Nvidia_graphics_processing_units#Data_Center_GPUs

Inference in large models

Memory requirements

- Without optimization
 - Model parameters: N
 - Weights: N bfloat16
- $4N$ bytes without counting activations
- 8B parameter model (i.e. Llama 3.1)
 - 16GB of memory for just weights

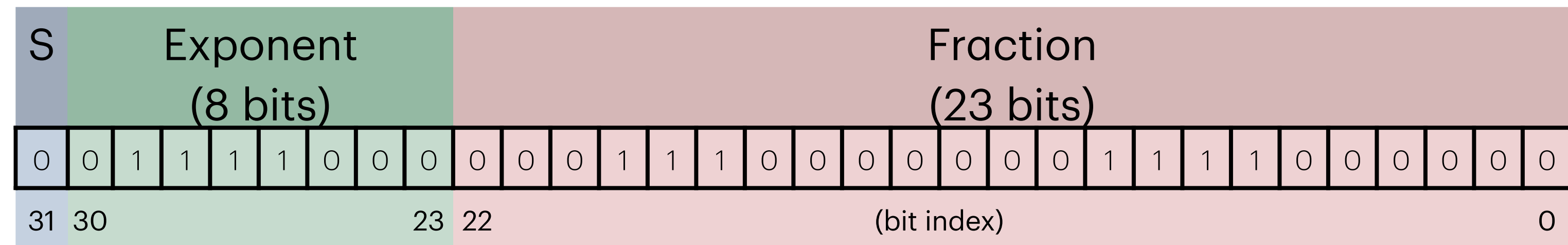
2N bytes

Weight (bf16)



float32

float32	
Sign	1 bit
Exponent	8 bit
Mantissa	23 bit
Precision (relative)	1E-07
Max value	3E+38
Min value (normal)	1.7E-38

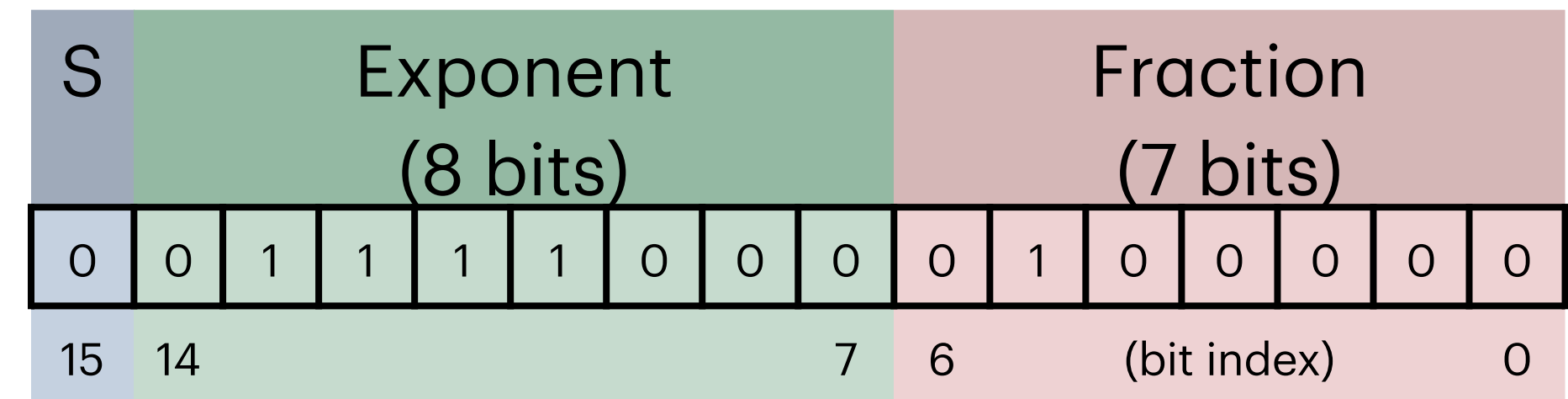


$$\text{Precision (relative)} = \frac{x_2 - x_1}{x_1}$$

where x_2 is smallest value $x_2 > x_1$

Bfloat16

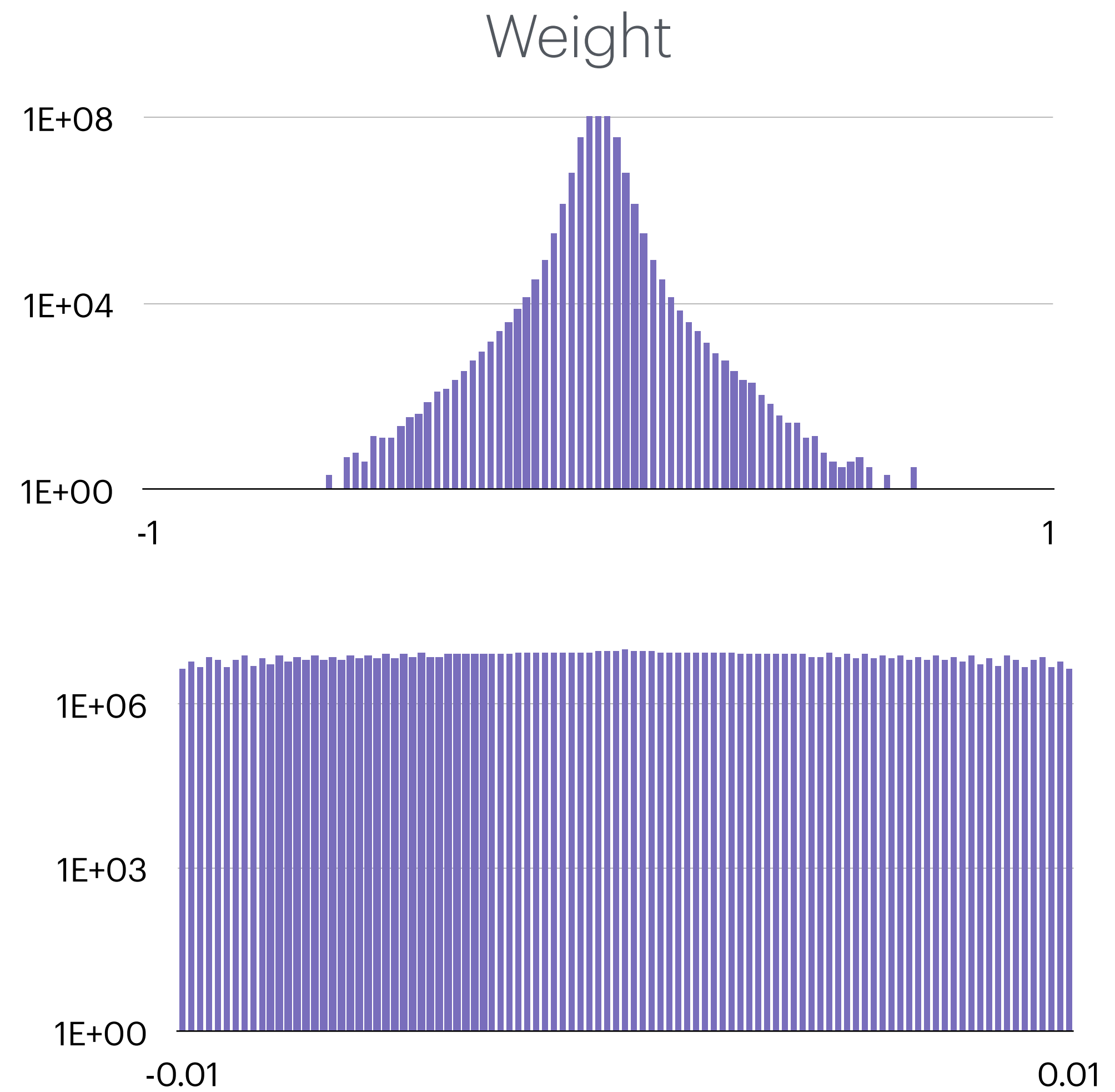
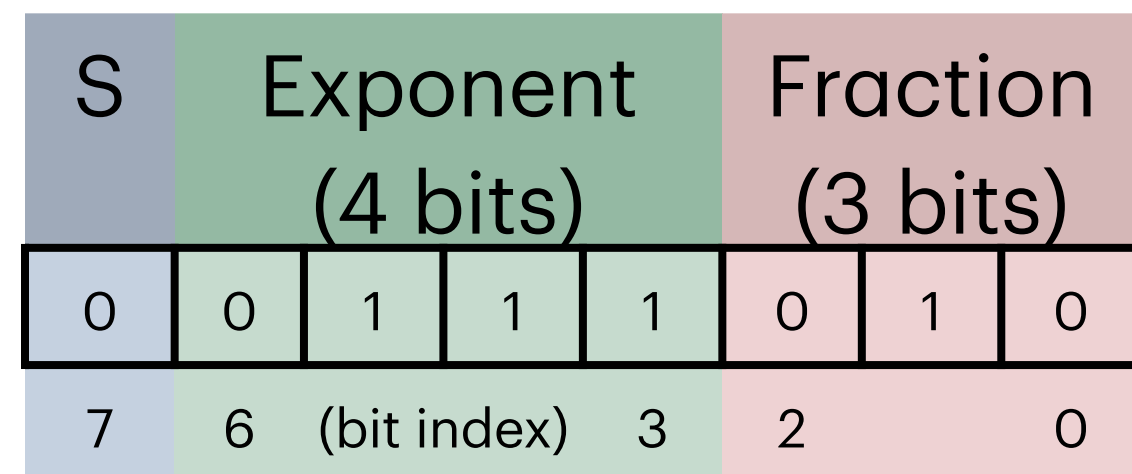
Bfloat16	
Sign	1 bit
Exponent	8 bit
Mantissa	7 bit
Precision (relative)	7.8E-03
Max value	3E+38
Min value (normal)	1.7E-38



Can we go lower?

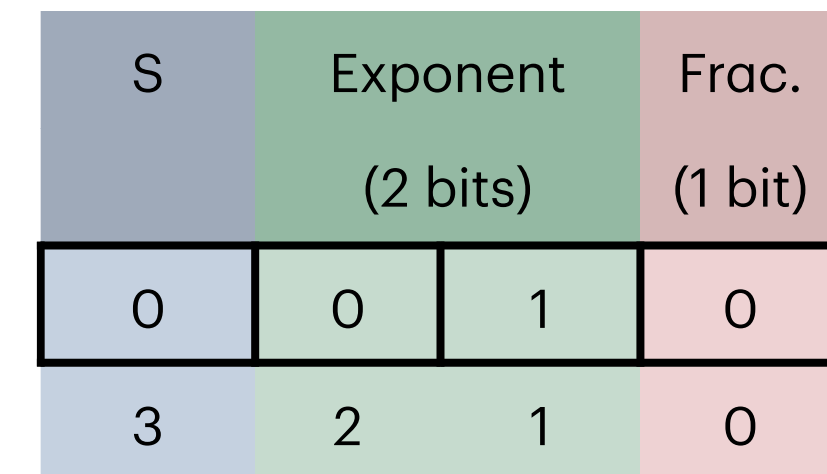
Float8	
Sign	1 bit
Exponent	4 bit
Mantissa	3 bit
Precision (relative)	0.125
Max value	240
Min value (normal)	0.015625

Llama 3.1 8b	
Largest weight	3.21
Smallest weight	-0.83
Largest gradient	5.09
Smallest gradient	-4.41
Largest activation	227.42
Smallest activation	-206.84



Can we go lower?

Float4	
Sign	1 bit
Exponent	4 bit
Mantissa	3 bit
Precision (relative)	0.125
Max value	240
Min value (normal)	0.015625



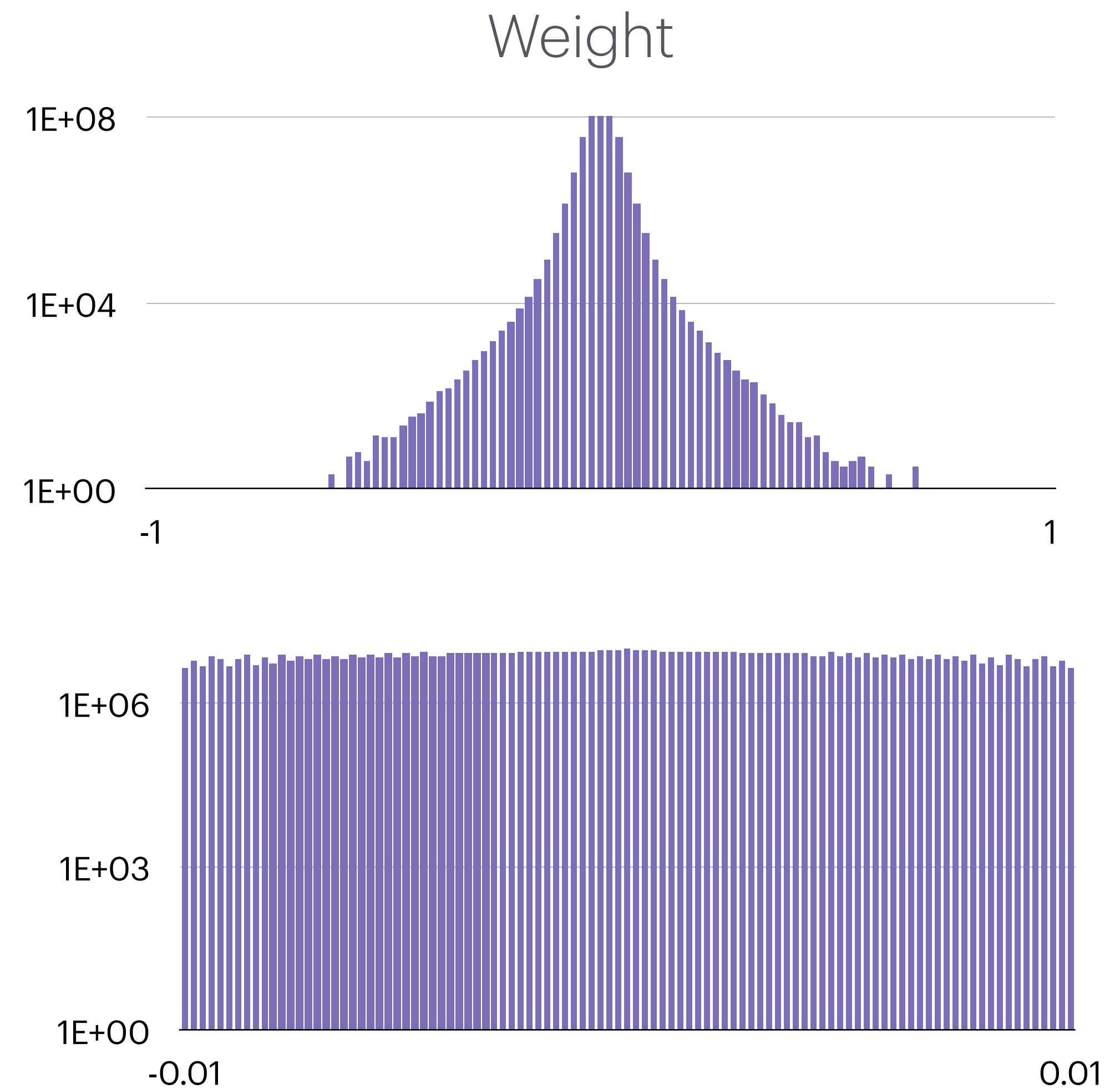
	0 ... 0	0 ... 1	1 ... 0	1 ... 1
... 00 ...	0	0.5	0	-0.5
... 01 ...	1	1.5	-1	-1.5
... 10 ...	2	3	-2	-3
... 11 ...	Inf	NaN	-Inf	NaN

Can we go lower?

Float4	
Sign	1 bit
Exponent	4 bit
Mantissa	3 bit
Precision (relative)	0.125
Max value	240
Min value (normal)	0.015625

Llama 3.1 8b	
Largest weight	3.21
Smallest weight	-0.83
Largest gradient	5.09
Smallest gradient	-4.41
Largest activation	227.42
Smallest activation	-206.84

S	Exponent		Frac.
	(2 bits)		(1 bit)
0	0	1	0
3	2	1	0



Integer scale quantization

- Find $T = \max_i |W_i|$

- For K-bit integer

- $Q_i = \operatorname{argmin}_q \left| \frac{W_i}{T}(2^{K-1} - 1) - q \right|$

for any K-bit signed integer q

- $Q_i = \operatorname{round} \left(\frac{W_i}{T}(2^{K-1} - 1) \right)$

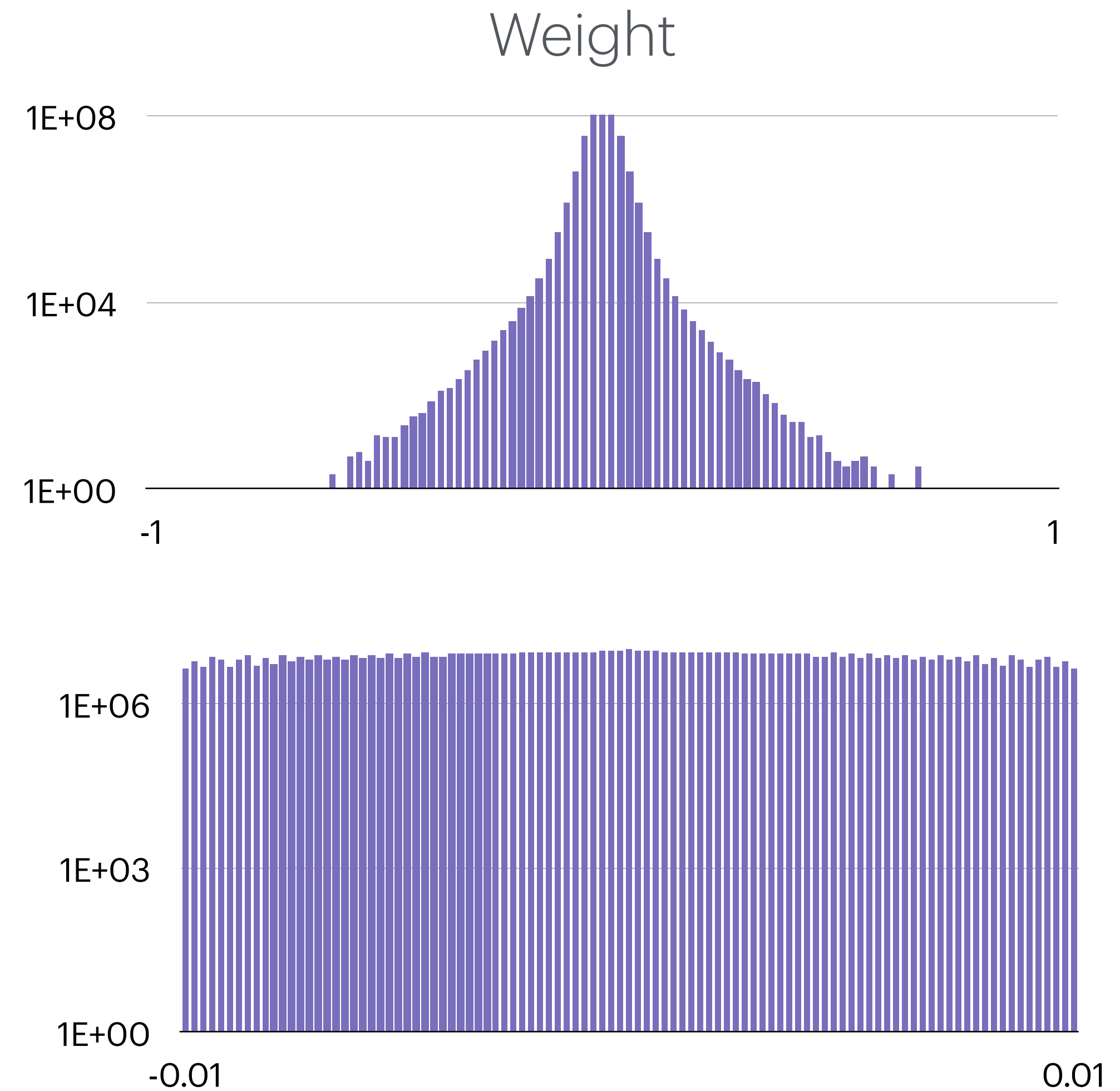
Integer scale quantization

- Quantization error

$$\frac{T}{2^{K-1} - 1}$$

- Weights: $\frac{NK}{8} + 2$ bytes

Llama 3.1 8b	
Largest weight	3.21
Smallest weight	-0.83
Largest gradient	5.09
Smallest gradient	-4.41
Largest activation	227.42
Smallest activation	-206.84



Integer affine quantization

- Find $A = \min_i |W_i|$, $B = \max_i |W_i|$

- For K-bit integer

- $Q_i = \operatorname{argmin}_q \left| \frac{W_i - A}{B - A} (2^K - 1) - q \right|$

for any K-bit unsigned integer q

- $Q_i = \operatorname{round} \left(\frac{W_i - A}{B - A} (2^K - 1) \right)$

Integer affine quantization

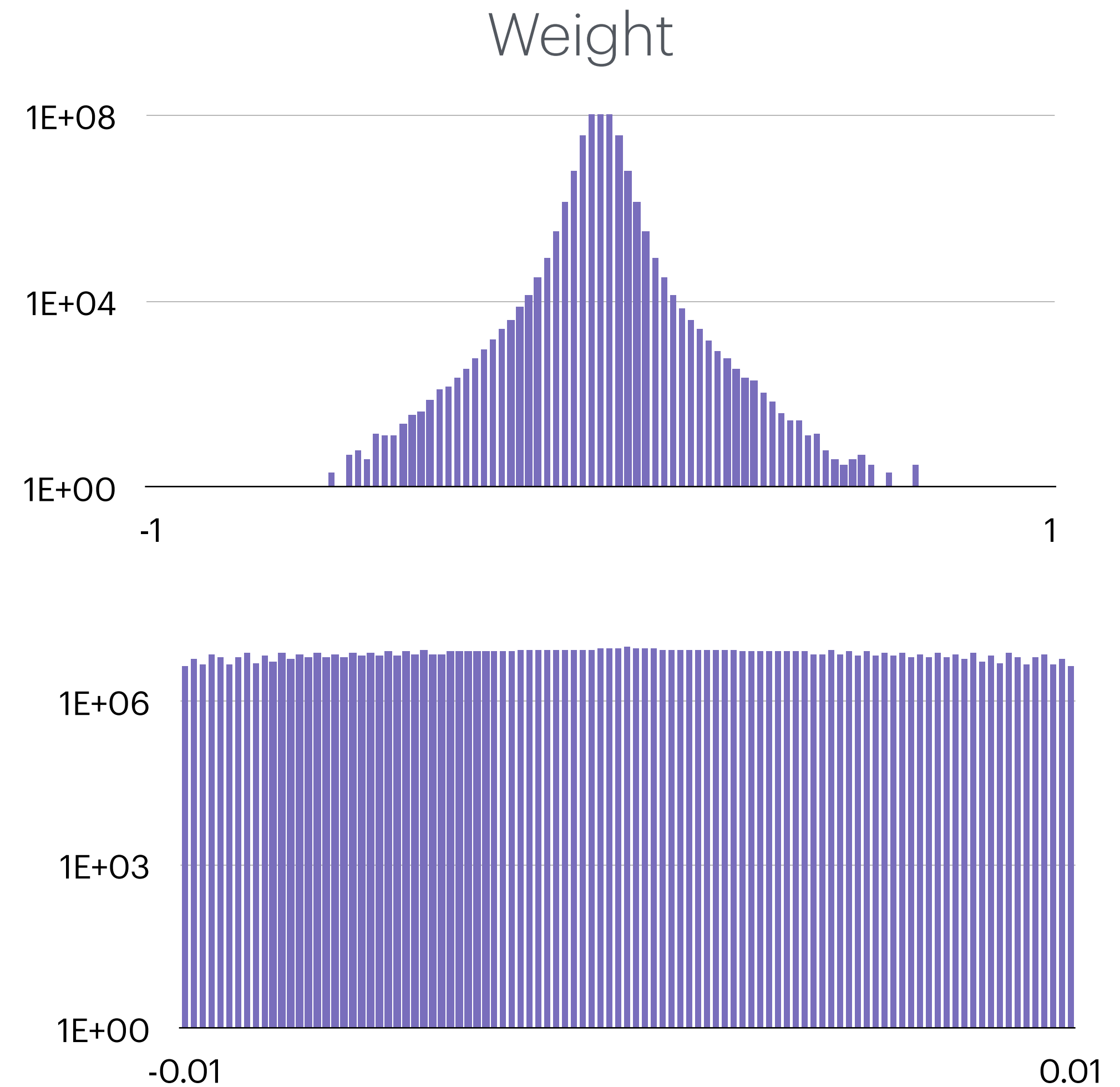
- Quantization error

$$\frac{B - A}{2^K - 1}$$

- $2^K - 1$

- Weights: $\frac{NK}{8} + 4$ bytes

Llama 3.1 8b	
Largest weight	3.21
Smallest weight	-0.83
Largest gradient	5.09
Smallest gradient	-4.41
Largest activation	227.42
Smallest activation	-206.84



Blockwise Quantization

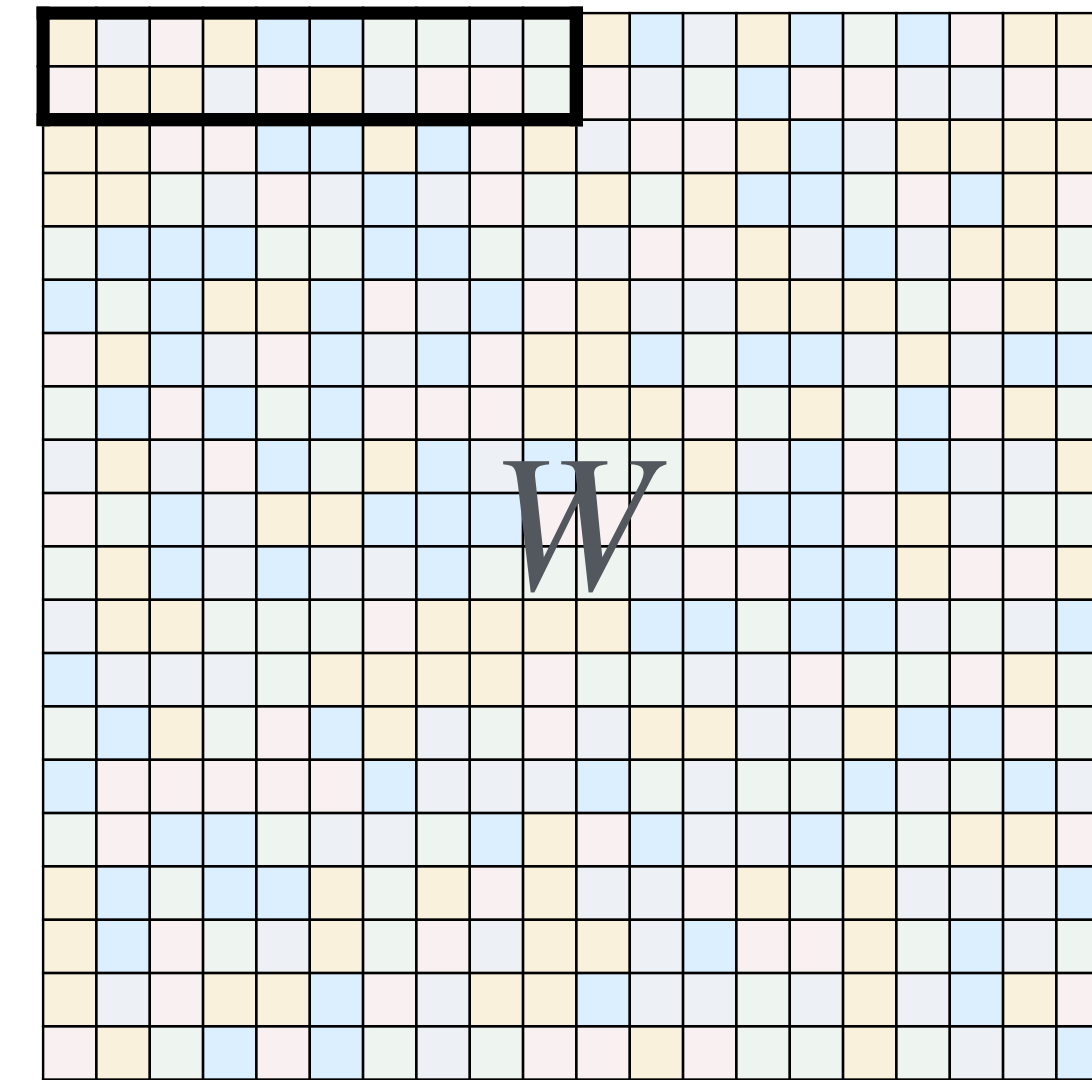
- Separate quantization constant T, A, B for a block of S weights

- Scale: $\frac{NK}{8} + 2\frac{N}{S}$ bytes

- $\frac{16}{S}$ extra bits per parameters

- Affine: $\frac{NK}{8} + 4\frac{N}{S}$ bytes

- $\frac{32}{S}$ extra bits per parameters



Double Quantization

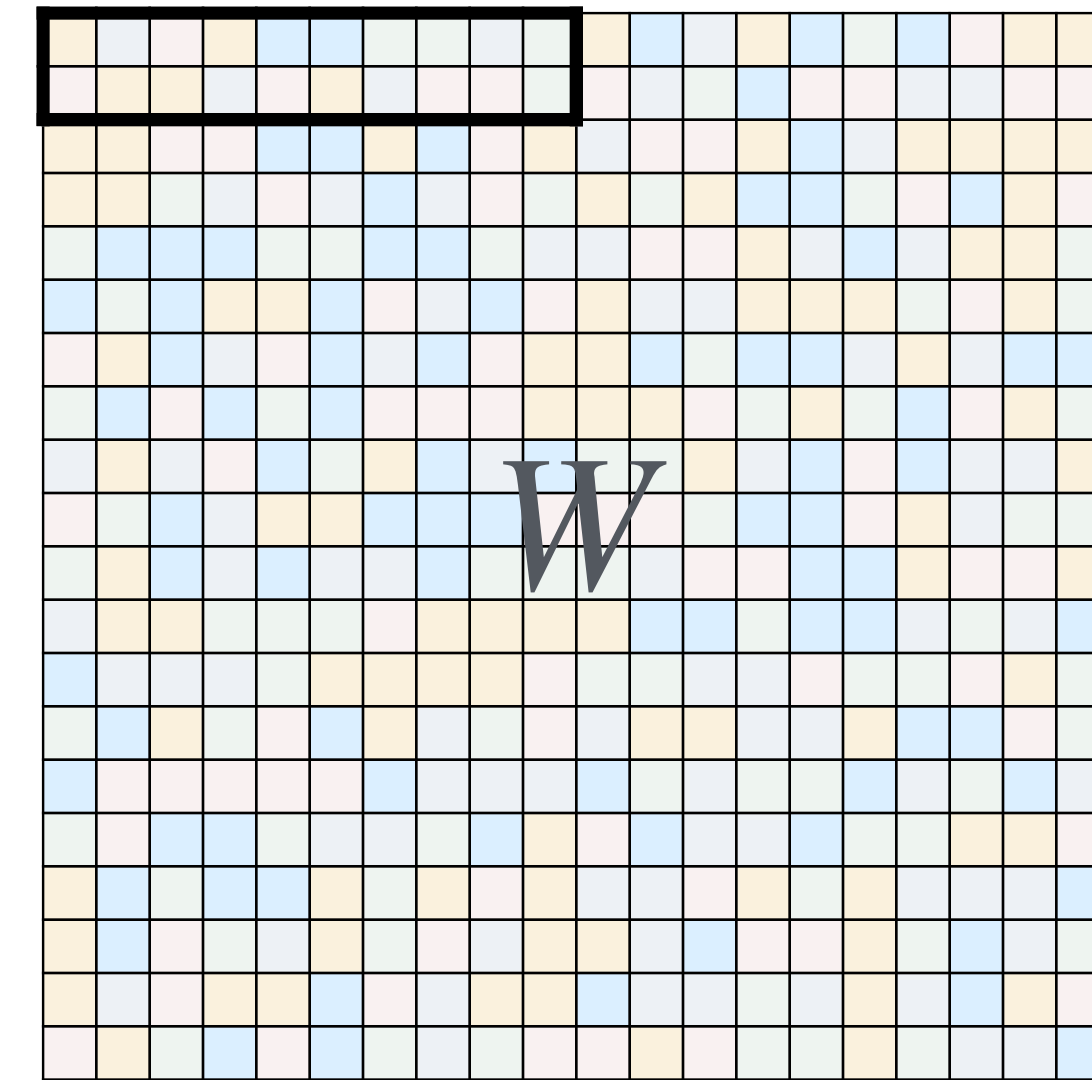
- Let's quantize the quantization factors T, A, B

- Scale: $\frac{NK}{8} + \frac{NK}{8S} + 2$ bytes

- $\frac{K}{S}$ extra bits per parameters

- Affine: $\frac{NK}{8} + \frac{2NK}{8S} + 4$ bytes

- $\frac{2K}{S}$ extra bits per parameters

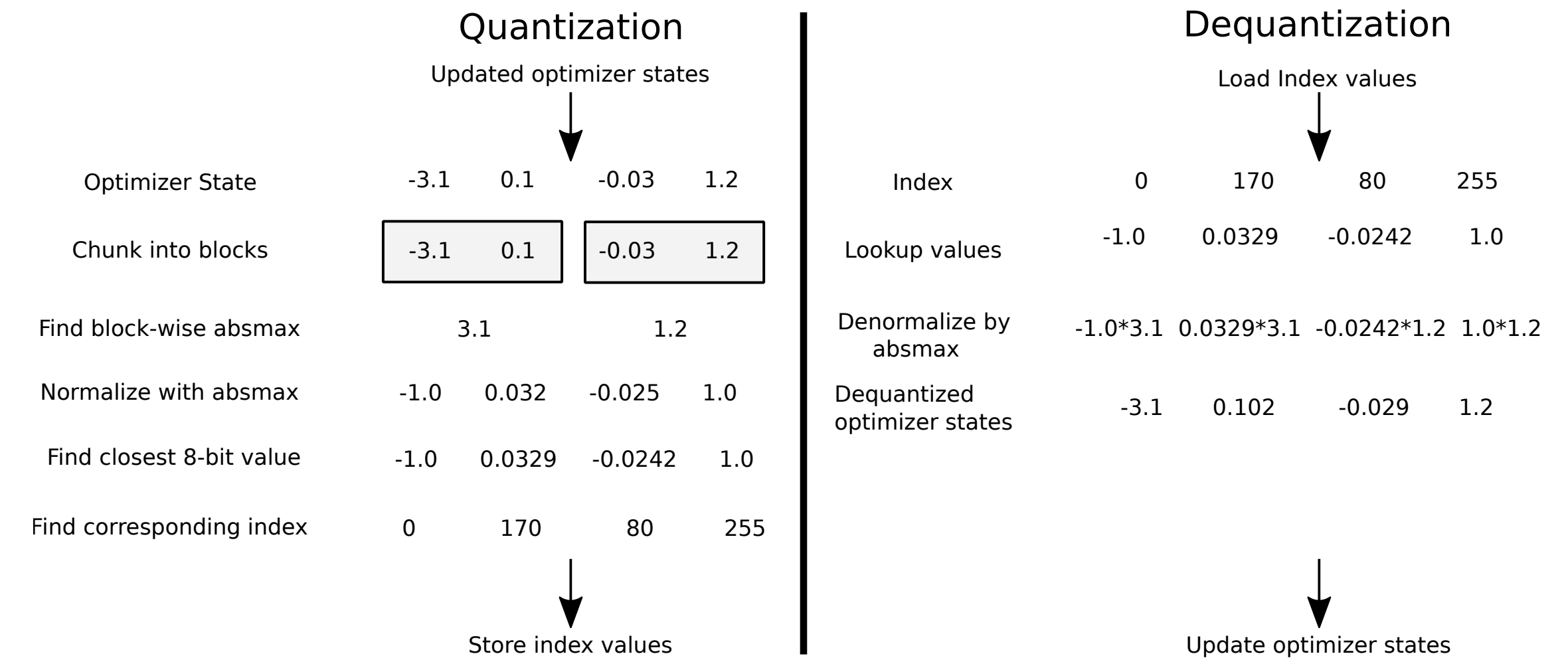


Beyond Linear Quantization

- 8-Bit Approximations for Parallelism in Deep Learning : <https://arxiv.org/abs/1511.04561>
- SqueezeLLM: <https://arxiv.org/abs/2306.07629>
- Extreme Compression of Large Language Models via Additive Quantization: <https://arxiv.org/abs/2401.06118>
- PV-Tuning: <https://arxiv.org/abs/2405.14852>

8-bit Adam

- Quantize 1st and 2nd momentum in Adam
 - 1st momentum: int8
 - 2nd momentum: uint8
 - Non-linear quantization
- Requires “stable” embeddings for LLMs
 - 32-bit optimizer states, normalizations



Stochastic rounding

How to train with quantized weights?

- Deterministic rounding

$$Q_i = \mathit{round} \left(\frac{W_i}{T} (2^{K-1} - 1) \right)$$

- Works only if update is largest than quantization

- Stochastic rounding

$$Q_i = \mathit{sround} \left(\frac{W_i}{T} (2^{K-1} - 1) \right)$$

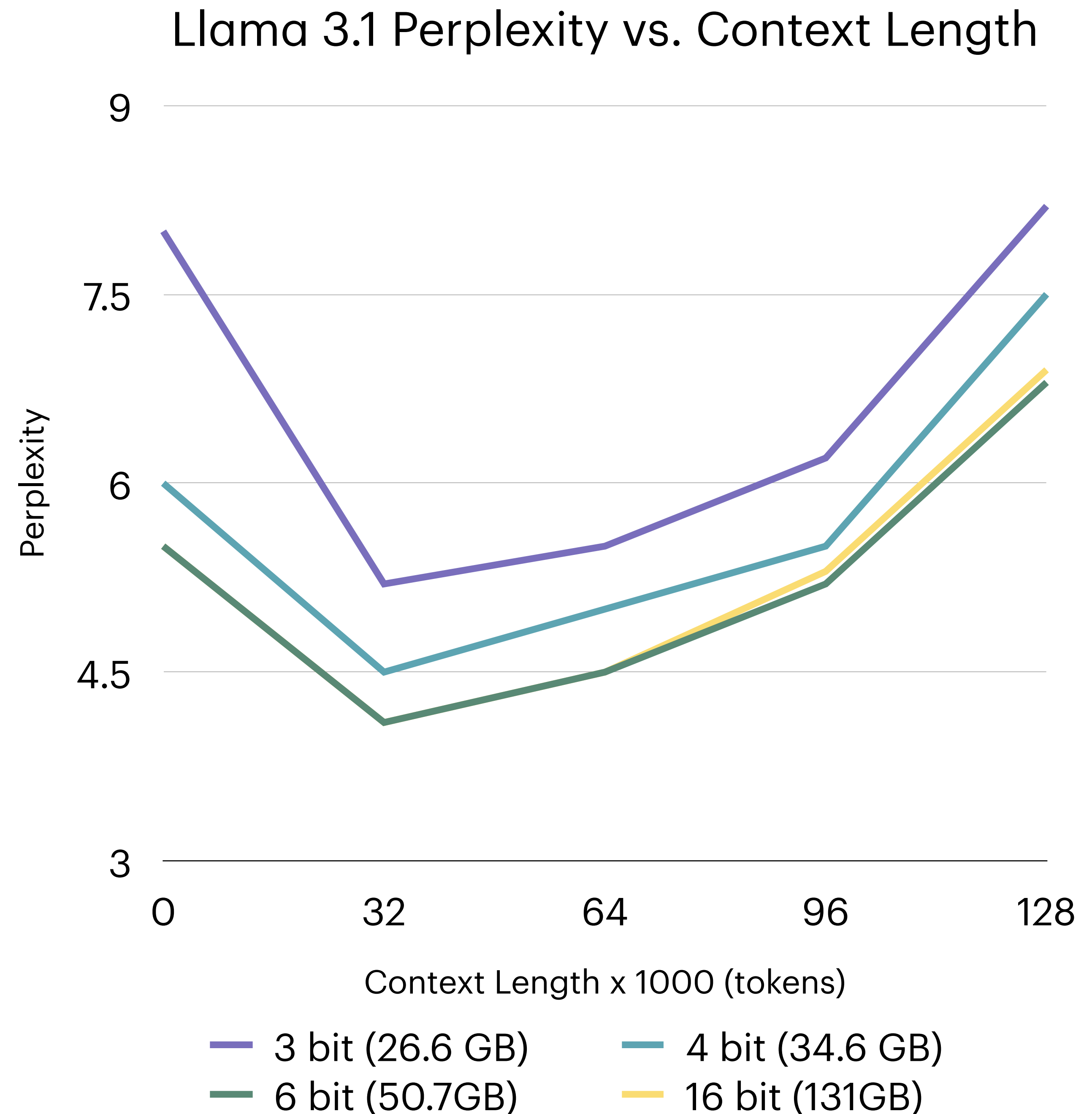
$$\mathit{sround}(x) = \begin{cases} [x] + 1 & \text{with } p \leq x - [x] \\ [x] & \text{otherwise} \end{cases}$$

$$E[\mathit{round}(x)] = \mathit{round}(x)$$

$$E[\mathit{sround}(x)] = x$$

How low can we go?

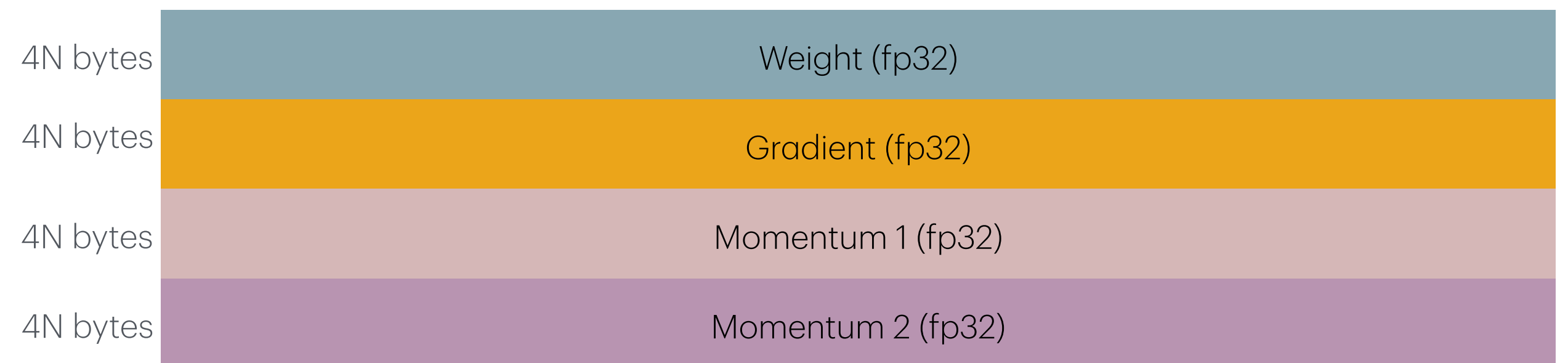
- GPT-style LLM can store about 2 bits of information per parameter
- Under ideal conditions
- 4 bits in practice
- Only after training!



Training large models

Memory requirements

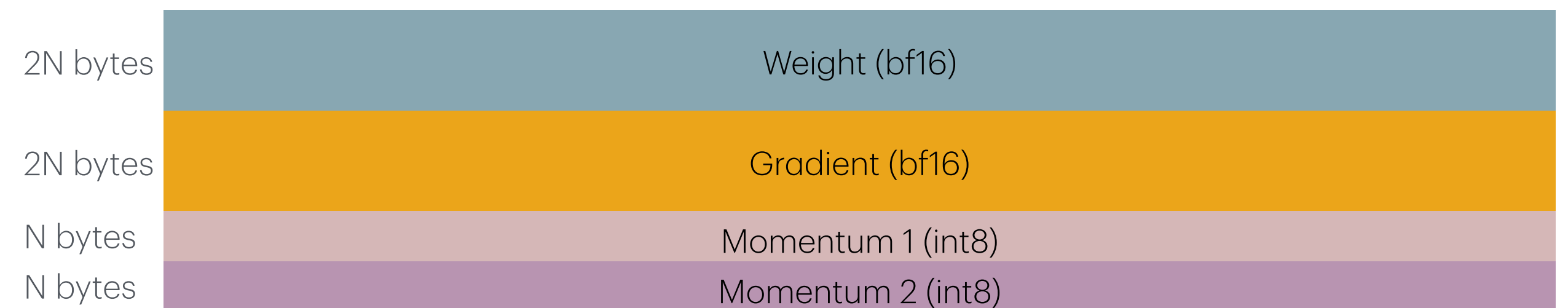
- Without optimization:
 - Model parameters: N
 - Weights: N floats
 - Gradients: N floats
 - Momentum: N floats
 - 2nd momentum (ADAM): N floats
- $16N$ bytes without counting activations



Training large models

Memory requirements

- Quantization:
 - Model parameters: N
 - Weights: N floats/bfloat16
 - Gradients: N floats/bfloat16
 - Momentum: N uint8
 - 2nd momentum (ADAM): N uint8
- $6N$ bytes without counting activations

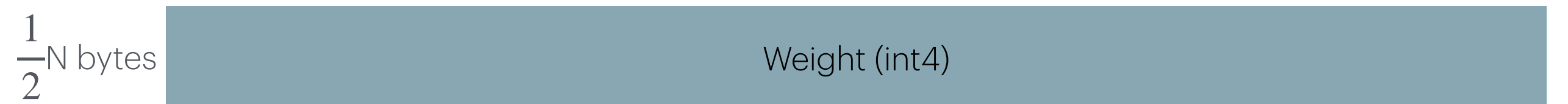


Inference in large models

Memory requirements

- Quantization:
 - Model parameters: N
 - Weights: N int4/int8
- $\frac{1}{2}$ N bytes without counting activations

$\frac{1}{2}$ N bytes



References

- [1] Tim Dettmers, et al. 8-bit optimizers via block-wise quantization. 2022. ([link](#))
- [2] Tim Dettmers, et al. QLoRA: Efficient Finetuning of Quantized LLMs. 2023. ([link](#))
- [3] Hao Li, et al. Training Quantized Nets: A Deeper Understanding. 2017. ([link](#))
- [4] Zeyuan Allen-Zhu, et al. Physics of Language Models: Part 3.3, Knowledge Capacity Scaling Laws. 2024. ([link](#))
- [5] Tim Dettmers. 8-bit approximations for parallelism in deep learning. 2015 ([link](#))
- [6] Sehoon Kim, et al. SqueezeLLM: Dense-and-Sparse Quantization. 2024 ([link](#))
- [7] Vage Egiazarian, et al. Extreme Compression of Large Language Models via Additive Quantization. 2024 ([link](#))
- [8] Vladimir Malinovskii, et al. PV-Tuning: Beyond Straight-Through Estimation for Extreme LLM Compression. 2024 ([link](#))