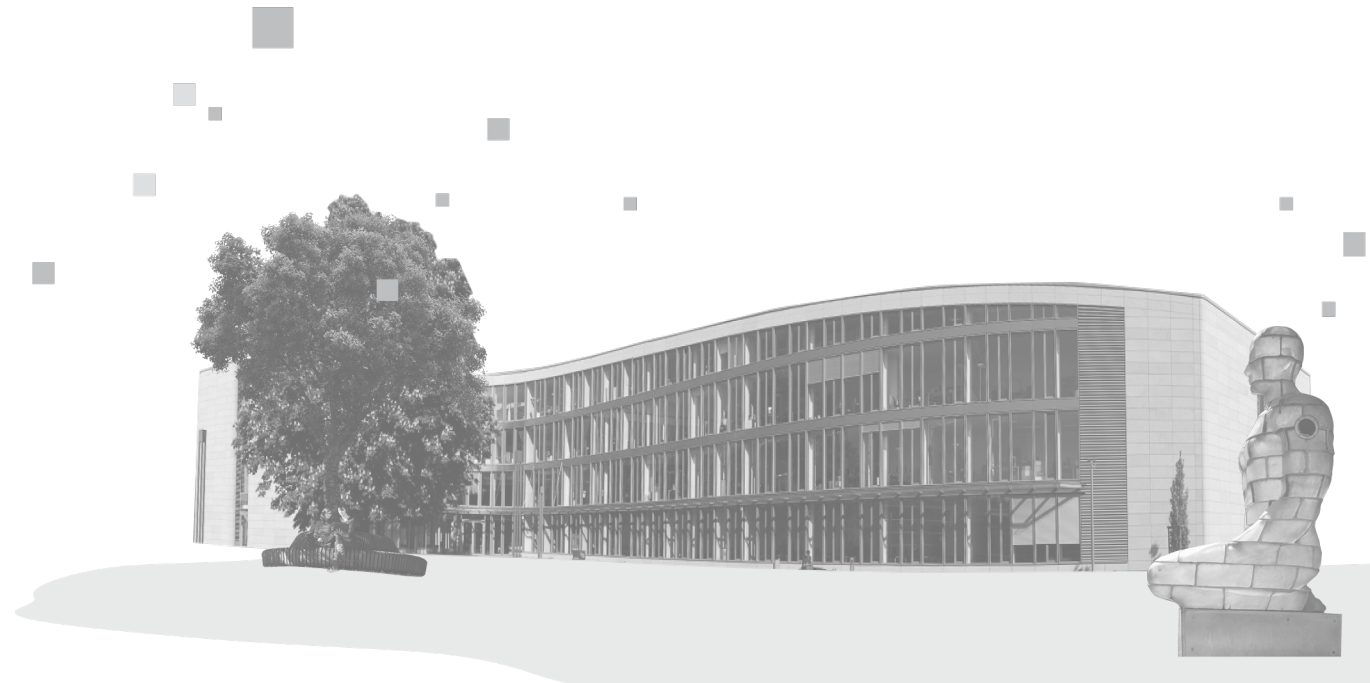# Building Local LLMs with Low-bit Quantization

## Reducing Model Size Without Sacrificing Performance

PD Dr. Haojin Yang,

Multimedia and Machine Learning Group,

Hasso Plattner Institute

**Design IT.**
**Create Knowledge.**

# Introduction to Low-Bit Quantization

- Low-bit quantization reduces model parameter precision to make LLMs smaller

- Converts 32-bit floating-point numbers to 8-bit, 4-bit, or even lower formats

- Enables running powerful LLMs on consumer hardware with limited resources

- Significantly reduces memory requirements and computational demands

- Trades minimal accuracy loss for substantial efficiency gains

- Essential technique for local deployment of AI models

# How Quantization Works

- Quantization reduces precision of model weights from 32-bit floating point to lower bit formats (8-bit, 4-bit, etc.)
- Process maps continuous values to a smaller, discrete set of values using scaling factors ($\alpha, \beta$)
- Post-training quantization (PTQ) applies quantization after model is fully trained
- Quantization-aware training (QAT) incorporates quantization during the training process
- **Weight-only quantization** focuses on model parameters while preserving activation precision
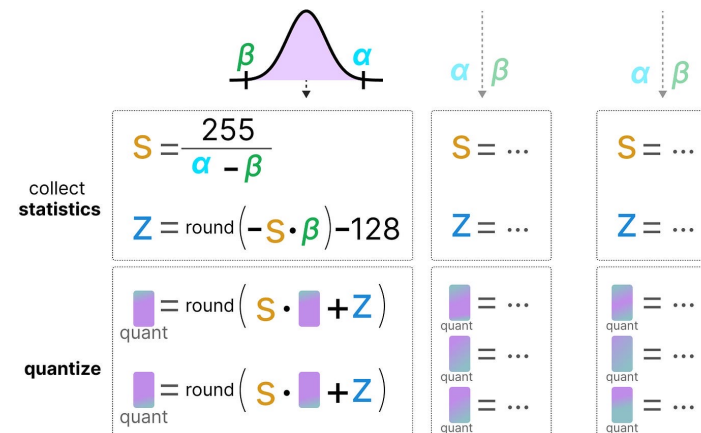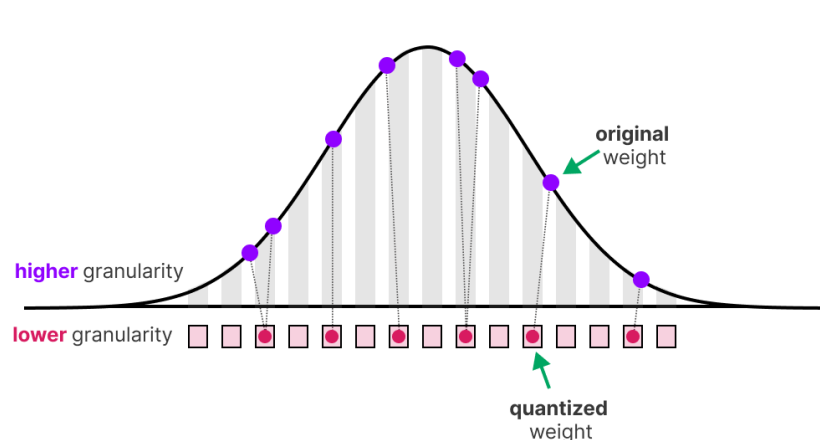- Different methods balance trade-offs between model size, inference speed, and accuracy



Image source: https://www.maartengrootendorst.com/blog/quantization/

# Quantization Techniques Compared

- 8-bit quantization: ~2x smaller models with minimal performance loss

- 4-bit quantization: ~8x smaller models with acceptable quality tradeoffs

- 2-bit/3-bit quantization: Experimental with significant quality degradation

- GPTQ vs AWQ: Popular algorithms with different optimization approaches

- INT4 vs NF4: Different number formats with varying performance characteristics

- Mixed-precision: Using different bit-widths for different model components

  - Layer-level mixed bit-width or channel-level mixed bit-width

**Model size estimation of a 70B Model with 70 billion parameters:**

$$64\text{-bit} = 64/8 \times 70B \approx 560 \text{ GB}$$
$$32\text{-bit} = 32/8 \times 70B \approx 280 \text{ GB}$$
$$16\text{-bit} = 16/8 \times 70B \approx 140 \text{ GB}$$
$$\mathbf{8\text{-bit} = 8/8 \times 70B \approx 70 \text{ GB}}$$
$$\mathbf{4\text{-bit} = 4/8 \times 70B \approx 35 \text{ GB}}$$
$$2\text{-bit} = 2/8 \times 70B \approx 17{,}7 \text{ GB}$$
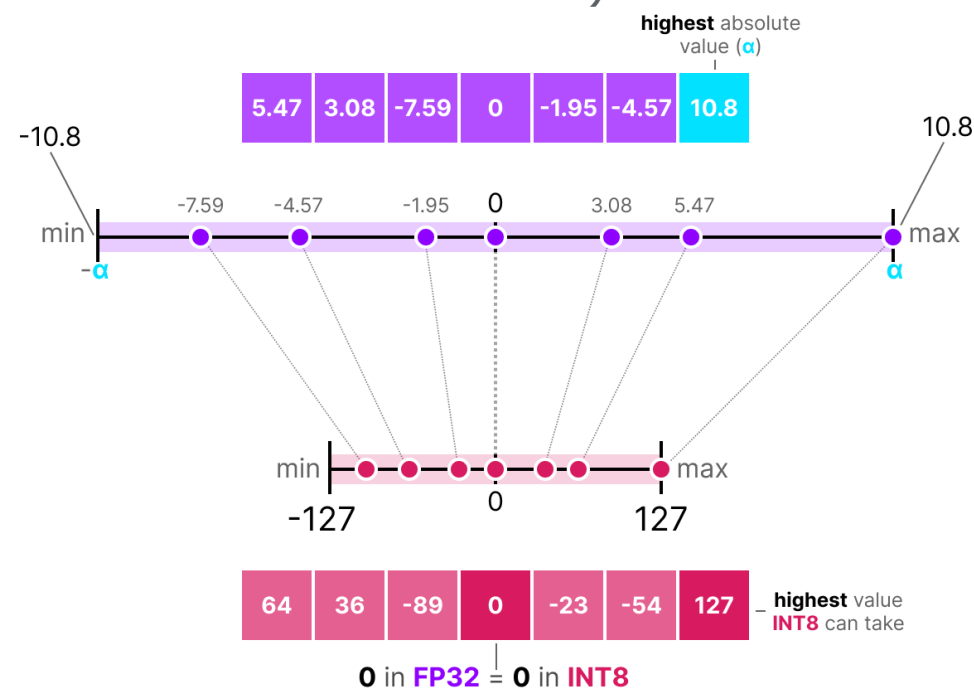
# Symmetric Quantization

- Maps floating-point values to a symmetric range around zero

- Zero in floating-point space maps exactly to zero in quantized space

- Uses a linear mapping centered around zero

- Scale factor: $s = \dfrac{(2^{b-1} - 1)}{\alpha}$ (where α is the maximum absolute value)

- Quantization: $q(x) = round(s \times x)$

- Dequantization: $\hat{x} = \dfrac{q(x)}{s}$



Image source: https://www.maartengrootendorst.com/blog/quantization/
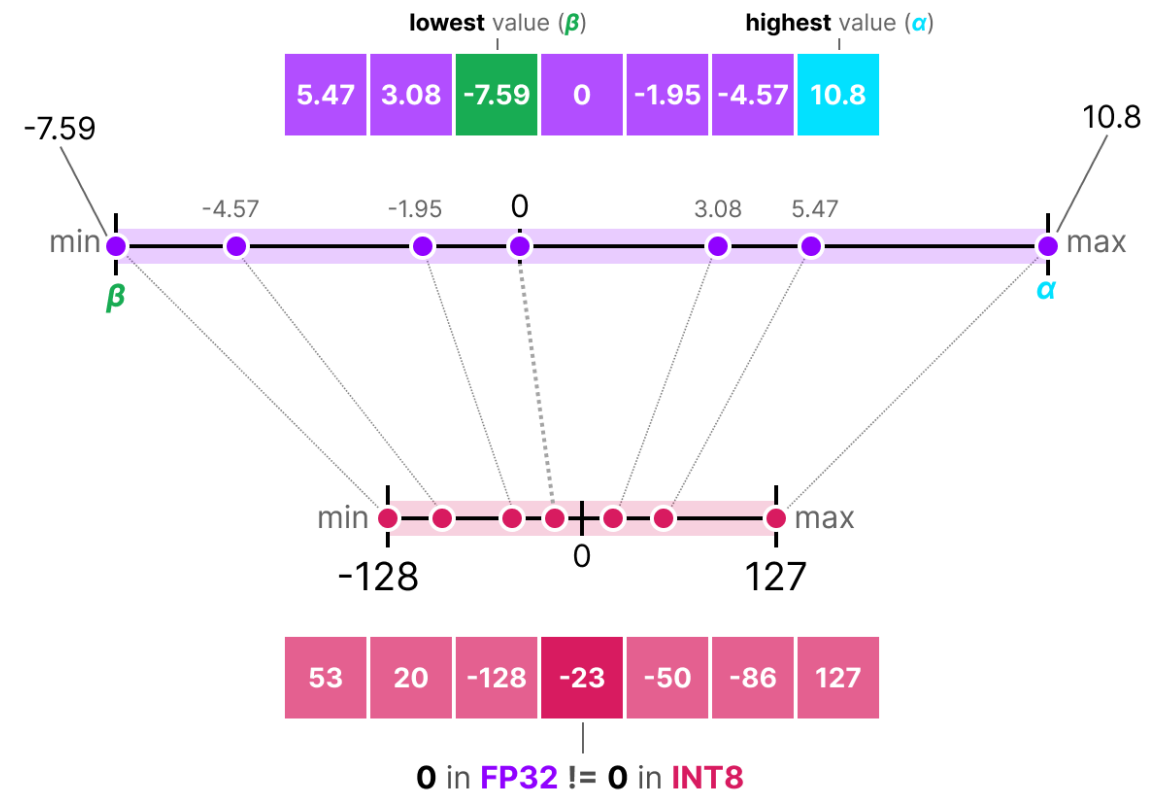
# Symmetric Quantization Example

HPI

- Example with original values: [8.5, 3.08, 3.02, -7.59, 10.8]

- Highest absolute value: 10.8

- Scale factor calculation: $S$ = 127 / 10.8 = 11.76

- Quantized values: [100, 36, 36, -89, 127]

- **Advantages:**

  - Simple implementation

  - Preserves zero exactly

  - No need to store zero-point offset

- **Limitations:**

  - Less efficient for data not centered around zero

  - Similar values may map to same quantized value (loss of precision)

# Asymmetric Quantization

- Not symmetric around zero

- Maps minimum (β) and maximum (α) values to full quantized range

- Uses zero-point quantization with shifted zero position

- Scale factor: $s = \dfrac{(2^b - 1)}{(\alpha - \beta)}$

- Zero point: $z = round(-\beta \times s)$

- Quantization: $q = round(x \times s + z)$

- Dequantization: $x = \dfrac{(q - z)}{s}$



Image source: https://www.maartengrootendorst.com/blog/quantization/

# Comparing Symmetric vs. Asymmetric Quantization

- Symmetric: Centered around zero, simpler calculations

- Asymmetric: Shifted zero position, better range utilization

- When to use Symmetric:

  - Data naturally centered around zero (like weights in neural networks)

  - When computational simplicity is important

  - When preserving zero exactly is critical

- When to use Asymmetric:

  - Data with skewed distribution (not centered around zero)

  - When maximizing representation accuracy is priority

  - For activations in neural networks (often positive-only)