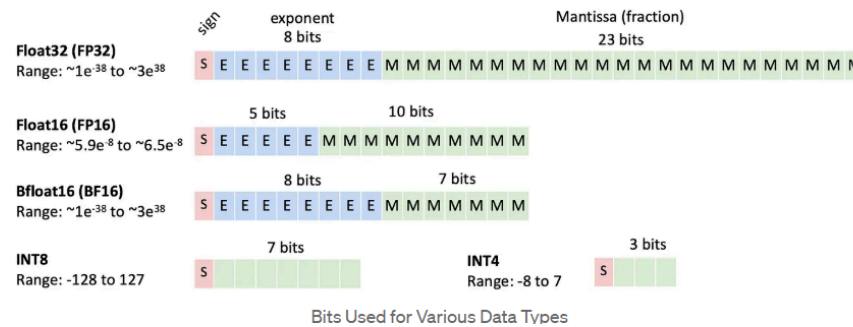


# Quantization Techniques in Deep Learning

In the era of large-scale deep learning models, optimizing inference efficiency without compromising performance is critical for real-world deployments. Quantization has emerged as a fundamental approach to achieving this optimization, particularly for edge devices, GPUs, and custom hardware accelerators.

This guide provides an in-depth understanding of quantization techniques, focusing on five key methods—**FP32, Dynamic Quantization, Static Quantization, Quantization-Aware Training (QAT), and Mixed Precision (FP16)**. Each section includes mathematical formulations, practical trade-offs, and experimental results to give you a well-rounded understanding of their applications.

For code, follow [Google Colab](#) or Github



## 1. Why Quantization?

Quantization reduces model size, memory consumption, and computational load by representing weights and activations in a lower-precision format, such as 8-bit integers (INT8) instead of 32-bit floating-point (FP32) values.

Key advantages:

- **Faster Inference:** Reduces computational complexity.
- **Smaller Model Sizes:** Enables deployment on resource-constrained devices.
- **Energy Efficiency:** Ideal for real-time applications and battery-powered devices.

Mathematically, the transformation from high-precision to quantized values can be expressed as:

$$\hat{x} = \text{round} \left( \frac{x}{s} \right)$$

where:

- $x$ : High-precision value.
- $s$ : Scaling factor that maps the floating-point range to the quantized range.
- $\hat{x}$ : Quantized value, often stored as INT8 or FP16.

The reverse transformation is used to reconstruct the approximate value:

$$x' = \hat{x} \cdot s$$

where  $x'$  is the reconstructed value.

## 2. Metrics to Evaluate Quantization

Quantization affects model accuracy and runtime efficiency. To quantify these trade-offs, the following metrics are commonly evaluated:

### 2.1 Perplexity (P)

Perplexity measures the quality of language models and is defined as:

$$\mathcal{P} = \exp \left( -\frac{1}{N} \sum_{i=1}^N \log p(t_i | \text{context}) \right)$$

where:

- N: Number of tokens.
- $p(t_i | \text{context})$ : Probability of token  $t_i$  given its context.

Lower perplexity indicates better model performance.

## 2.2 Latency (T)

Latency represents the average time required for a single inference pass. For MMM repetitions:

$$T = \frac{1}{M} \sum_{j=1}^M t_j$$

where  $t_j$  is the time for the  $j$ -th inference run.

## 2.3 Model Size (S)

The total model size in megabytes (MB) is:

$$S = \sum_{i=1}^K \text{elements}(w_i) \cdot \text{element\_size}(w_i)$$

where:

- K: Total number of weight tensors.
- $w_i$ : i-th weight tensor.

. . .

## 3. Quantization Techniques

### 3.1 Full Precision (FP32)

The baseline representation uses 32-bit floating-point (FP32) numbers for weights and activations. It provides the highest precision but is computationally and memory-intensive.

- **Mathematics:**

For weights  $w$  and activations  $a$ :

$$w, a \in \mathbb{R}^{32}$$

- **Advantages:**

- Maximum precision.
- No information loss during computations.

- **Drawbacks:**

- High memory usage.
- Slower inference on low-power devices.

- **Results:**

- Perplexity: 53.63
- Latency: 125.88 ms
- Model Size: 497.76 MB

. . .

## 3.2 Dynamic Quantization

Dynamic quantization reduces memory usage by quantizing weights at runtime while keeping activations in FP32. This method is simple and works well for linear layers.

- **Mathematics:** Weights are quantized dynamically using:

$$\hat{w} = \text{round} \left( \frac{w}{s} \right)$$

where:

$$s = \frac{\max(w) - \min(w)}{2^n - 1}$$

n is the quantization band-width

- **Advantages:**
  - Minimal overhead.
  - Suitable for low-complexity inference.
- **Drawbacks:**
  - Activations remain unquantized, limiting memory benefits.
  - Suboptimal accuracy for non-linear operations.

- **Results:**
- Perplexity: 3896.48
- Latency: 133.59 ms
- Model Size: 497.76 MB

. . .

### 3.3 Static Quantization

Static quantization converts both weights and activations to a lower precision ahead of runtime using a calibration dataset to compute scales and zero-points.

- **Mathematics:** For activations  $a$ :

$$\hat{a} = \text{round} \left( \frac{a - z}{s} \right)$$

where:

- $s$ : Scale factor.
- $z$ : Zero-point, representing the quantized value of zero.
- Reconstructed activation:



$$a' = s \cdot \hat{a} + z$$

- **Advantages:**
- Significant reduction in latency.
- Optimized for deployment on INT8-optimized hardware.
- **Drawbacks:**
- Requires calibration datasets.
- Limited accuracy for complex models.
- **Results:**
- Perplexity: 3896.48
- Latency: 83.17 ms
- Model Size: 497.76 MB

. . .

### 3.4 Quantization-Aware Training (QAT)

QAT simulates quantization during training, allowing the model to adapt to quantization-induced errors.

- **Mathematics:** During training:

$$\hat{w}_{\text{QAT}} = \text{round} \left( \frac{w}{s} \right)$$

$$w_{\text{QAT}} = s \cdot \hat{w}_{\text{QAT}}$$

- Inference uses the quantized weights directly.
- **Advantages:**
- High accuracy, close to FP32.
- Suitable for complex networks like transformers.
- **Drawbacks:**
- Computationally expensive training process.
- Requires retraining from scratch or fine-tuning.
- **Results:**
- Perplexity: 53.89
- Latency: 1571.14 ms
- Model Size: 248.88 MB

. . .

### 3.5 Mixed Precision (FP16)

Mixed precision uses 16-bit floating-point (FP16) for both weights and activations, striking a balance between performance and precision.

- **Mathematics:** Weights and activations are represented as:

$$w_{\text{FP16}}, a_{\text{FP16}} \in \mathbb{R}^{16}$$

- **Advantages:**
  - Significant size reduction.
  - Suitable for modern GPUs with native FP16 support.
- **Drawbacks:**
  - Possible precision loss for small values.
  - Requires specialized hardware.
- **Results:**
  - Perplexity: 53.89
  - Latency: 1547.32 ms
  - Model Size: 248.88 MB

. . .

## 4. Trade-offs Between Techniques

The table below summarizes the trade-offs between different quantization methods:

Technique	Perplexity	Latency
FP32	53.63	125
Dynamic Quantization	3896.48	60
Static Quantization	3896.48	5
QAT	53.89	1571
Mixed Precision (FP16)	53.89	154

For code, follow [Google Colab](#) or [Github](#)

. . .

## 5. Recommendations

1. Use **FP32** for high accuracy in research settings.
2. Apply **Dynamic Quantization** for simple linear models.
3. Opt for **Static Quantization** in resource-constrained devices.
4. Leverage **QAT** for high-stakes applications where accuracy matters.
5. Deploy **FP16** for state-of-the-art GPUs.

. . .

## References

1. Quantization in PyTorch
2. TensorFlow Model Optimization
3. Research papers:
  - “Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference”
  - “Mixed Precision Training”

