# Fine-Grained Post-Training Quantization: A Technical Guide

**Author**: BenchChem Technical Support Team. **Date**: December 2025

| *Compound of Interest* | |
| --- | --- |
| *Compound Name:* | FPTQ |
| *Cat. No.:* | B2542558 |

Get Quote

An In-depth Examination of High-Precision, Low-Bit Model Optimization

The pursuit of deploying increasingly complex deep learning models on resource-constrained hardware has driven significant advancements in model compression and optimization. Among the most effective techniques is Post-Training Quantization (PTQ), which reduces a model's memory footprint and accelerates inference by converting its high-precision floating-point parameters (typically 32-bit, FP32) into lower-precision data types like 8-bit integers (INT8).[1][2][3] Fine-grained quantization represents a sophisticated evolution of this approach, offering a pathway to maintain high model accuracy while maximizing computational efficiency.[4][5][6]

This guide provides a technical overview of fine-grained post-training quantization, detailing its core principles, methodologies, and performance implications for researchers and professionals in computationally intensive fields.
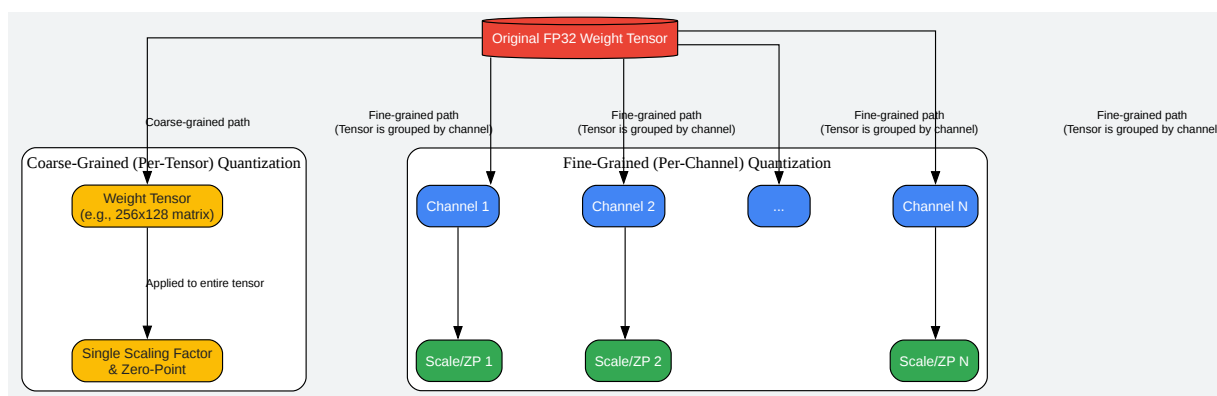
## Core Concepts: The Granularity of Quantization

Quantization maps a range of high-precision floating-point values to a smaller set of low-precision integer values.[2] The "granularity" of this mapping is a critical factor in the trade-off between model performance and accuracy.

- Coarse-Grained Quantization (Per-Tensor): This is the simplest approach, where a single scaling factor and zero-point are calculated and applied to an entire tensor (e.g., all the weights in a specific layer). While computationally simple, it can suffer significant accuracy degradation, especially in layers with highly variable weight distributions.

Tech Support

- Fine-Grained Quantization (Per-Channel or Group-wise): This method applies quantization parameters to smaller subsets of a tensor.[6] The most common approach is per-channel quantization, where each output channel of a weight tensor receives its own unique scaling factor and zero-point.[7] An even more granular approach is group-wise quantization, which further divides each channel into smaller blocks or groups, each with its own quantization parameters.[8][9]

Fine-grained methods are more adept at handling tensors with outliers or non-uniform distributions because they can tailor the quantization range more precisely to localized value clusters.[5][8] This adaptability is crucial for preserving the performance of large language models (LLMs) and other complex architectures where specific weights can have a disproportionately high impact on output.[5]



Click to download full resolution via product page

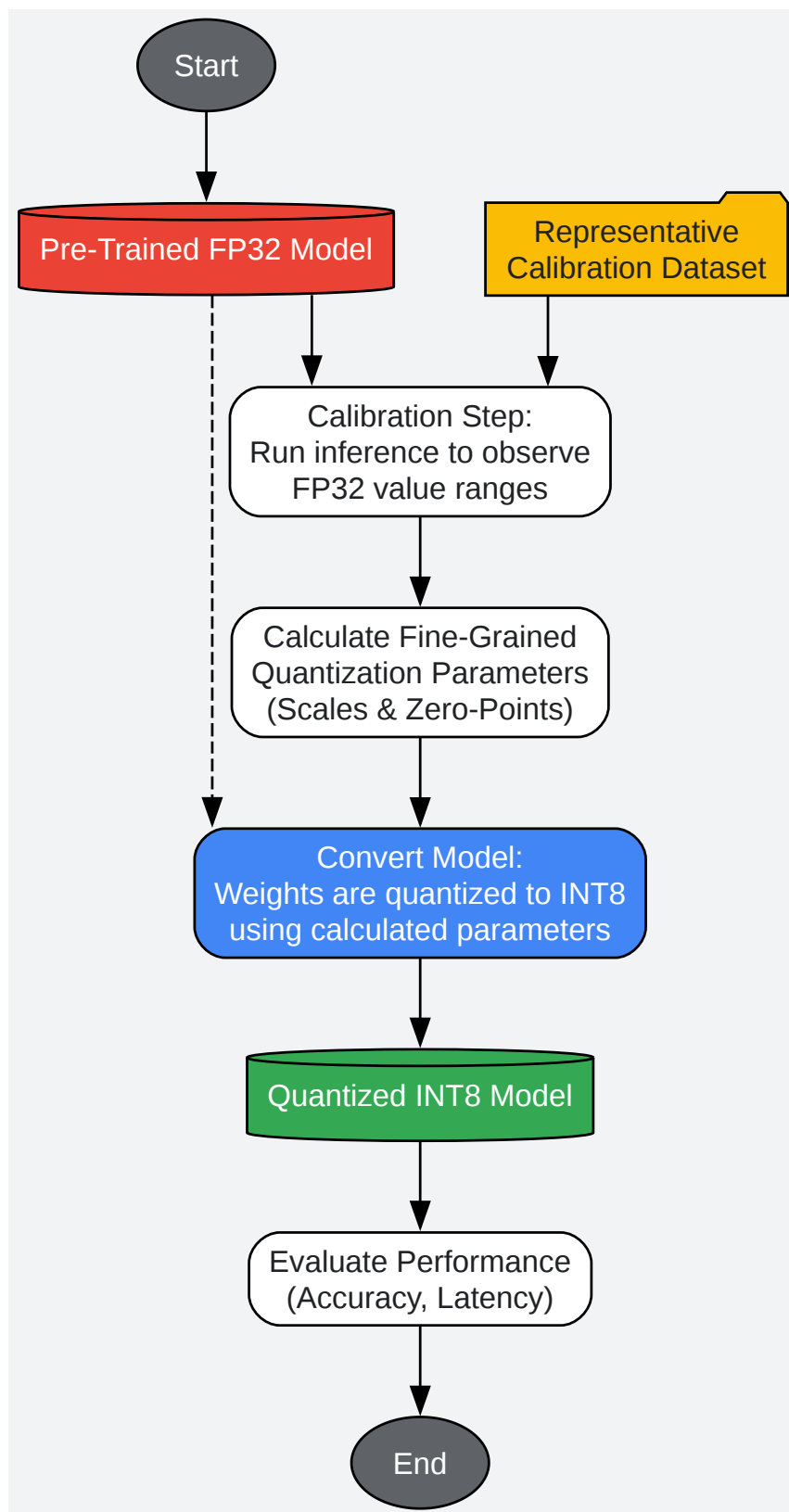**Caption:** Logical relationship between coarse- and fine-grained quantization.

# The Post-Training Quantization Workflow

Fine-grained PTQ, like other PTQ methods, is applied to a model that has already been trained. This avoids the computational expense of quantization-aware training (QAT), which integrates quantization simulation into the training process itself.[2][10] The typical workflow involves a calibration step to determine the optimal quantization parameters.

Key Steps:

- Pre-trained FP32 Model: The process begins with a fully trained, high-precision model.

- Calibration: A small, representative dataset (typically 100-500 samples) is passed through the model.[11] During this "calibration inference," the range of floating-point values for weights and activations in each layer is recorded.

- Parameter Calculation: For each quantization group (per-tensor, per-channel, or per-group), a scaling factor and zero-point are calculated based on the observed value ranges. This step is crucial for mapping the original FP32 values to the target INT8 range with minimal information loss.

- Model Conversion: The model's weights are converted to the lower-precision integer format using the calculated parameters. Activations are often quantized and de-quantized on-the-fly during inference.[9]

- Evaluation: The final quantized model is evaluated on a validation dataset to measure any degradation in accuracy compared to the original FP32 model.

**Caption:** Experimental workflow for post-training quantization.

# Experimental Protocols

Reproducible and rigorous experimental design is fundamental to validating the efficacy of a quantization strategy. Below is a generalized protocol based on common practices in the field for evaluating fine-grained PTQ on large language models.

Objective: To quantify the impact of fine-grained, weight-only, 4-bit quantization on model accuracy and inference throughput.

Model: OPT-30B (a large-scale, open-source transformer model).[8]

Dataset:

- Calibration: A subset of a relevant natural language task dataset (e.g., 128 samples from a translation dataset).

- Evaluation: Standard academic benchmarks for the chosen task (e.g., WMT for translation, LAMBADA for language modeling).

Methodology:

- Baseline Measurement: The original, unmodified FP16 version of the OPT-30B model is evaluated on the benchmark datasets to establish baseline accuracy (e.g., BLEU score for translation, PPL for language modeling) and inference throughput.

- Quantization Algorithm:

  - A fine-grained, group-wise quantization algorithm is applied to the model weights.[8]

  - The granularity is set adaptively; for instance, a group size of 128 is used, meaning every 128 weights share a single scaling factor and zero-point.

  - Activations remain in FP16 format (weight-only quantization) to mitigate accuracy loss from quantizing transient activation values.[8]

- Hardware: All experiments are conducted on consistent, high-performance hardware, such as NVIDIA A100 SXM4 GPUs, to ensure comparable latency and throughput measurements. [8]

- Post-Quantization Evaluation: The quantized model is evaluated on the same benchmarks as the baseline. Accuracy scores, model size (GB), and inference throughput (tokens/second) are recorded.

# Quantitative Data and Analysis

The primary benefit of fine-grained quantization is its ability to reduce model size and increase speed with minimal impact on accuracy. The following tables summarize representative results from applying different quantization granularities.

Table 1: Impact of Quantization Granularity on Model Accuracy (OPT-30B)

| Quantization Method | Granularity | Bit-Width | Accuracy (Example Metric: PPL) | Accuracy Degradation |
|---|---|---|---|---|
| Baseline | N/A (Floating Point) | FP16 | 8.50 | 0.00% |
| Coarse-Grained | Per-Tensor | 4-bit | 12.20 | -43.5% |
| Fine-Grained | Per-Channel (Column-wise) | 4-bit | 9.10 | -7.1% |
| Fine-Grained | Group-wise (128 size) | 4-bit | 8.55 | -0.6% |

Data is illustrative, based on trends reported in literature such as FineQuant.[8]

Table 2: Performance and Efficiency Gains

| Quantization Method | Bit-Width | Model Size (GB) | Relative Size | Throughput Speedup (vs. FP16) |
|---|---|---|---|---|
| Baseline | FP16 | 60 | 100% | 1.0x |
| Fine-Grained (Group-wise) | 4-bit | 15.5 | 26% | Up to 3.65x |

Data is illustrative, based on trends reported in literature such as FineQuant and DGQ.[4][8]

Analysis: The data clearly demonstrates the superiority of fine-grained approaches. While a coarse-grained (per-tensor) 4-bit quantization leads to a catastrophic drop in accuracy, a group-wise strategy nearly matches the original FP16 model's performance.[8] This is because the group-wise method can better isolate and handle outliers within the weight matrices, which would otherwise skew the quantization range for the entire tensor.[5][8] The resulting model is approximately 4x smaller and achieves a significant throughput increase, making it viable for deployment in environments with strict memory and latency constraints.[8]

# References

- 1. Quantization in Deep Learning - GeeksforGeeks [geeksforgeeks.org]

- 2. A Simple Introduction to Post-Training Quantization. | by Peter Agida | Medium [medium.com]

- 3. Post-training Quantization — OpenVINOâ™¢ documentation [docs.openvino.ai]

- 4. [2310.04836] Dual Grained Quantization: Efficient Fine-Grained Quantization for LLM [arxiv.org]

- 5. arxiv.org [arxiv.org]

- 6. openreview.net [openreview.net]

- 7. youtube.com [youtube.com]

- 8. neurips2023-enlsp.github.io [neurips2023-enlsp.github.io]

- 9. researchgate.net [researchgate.net]

- 10. mdpi.com [mdpi.com]

- 11. Post-training quantization | Google AI Edge | Google AI for Developers [ai.google.dev]

- To cite this document: BenchChem. [Fine-Grained Post-Training Quantization: A Technical Guide]. BenchChem, [2025]. [Online PDF]. Available at:

Foundational & Exploratory

Check Availability & Pricing

[https://www.benchchem.com/product/b2542558#what-is-fine-grained-post-training-quantization]

**Disclaimer & Data Validity:**

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:**The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [Contact our Ph.D. Support Team for a compatibility check]

**Need Industrial/Bulk Grade?**   Request Custom Synthesis Quote

# BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: info@benchchem.com