# what is Fine-grained Post-Training Quantization

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

| Compound of Interest | |
|---|---|
| Compound Name: | FPTQ |
| Cat. No.: | B15621169 |

Get Quote

An In-depth Technical Guide to Fine-grained Post-Training Quantization

For Researchers, Scientists, and Drug Development Professionals

## Introduction to Post-Training Quantization

In the era of large-scale neural networks, particularly in fields like drug discovery and scientific research where model complexity is ever-increasing, deploying these models efficiently presents a significant challenge. Post-Training Quantization (PTQ) has emerged as a critical optimization technique to reduce the computational and memory demands of these models without the need for costly retraining.[1][2] PTQ converts the high-precision floating-point parameters (typically 32-bit) of a trained model to lower-precision data types, such as 8-bit integers (INT8) or even 4-bit integers (INT4).[3] This reduction in precision leads to a smaller memory footprint, faster inference speeds, and lower power consumption, making it feasible to deploy large models on resource-constrained environments.[3]

## The Core Principles of Fine-grained Quantization

While basic PTQ applies a uniform quantization scale to an entire tensor (per-tensor quantization), this "coarse-grained" approach can lead to significant accuracy degradation, especially for models with diverse weight distributions.[4] Fine-grained quantization addresses this limitation by applying quantization parameters at a more granular level. This approach recognizes that different parts of a neural network have varying sensitivities to quantization.[5]
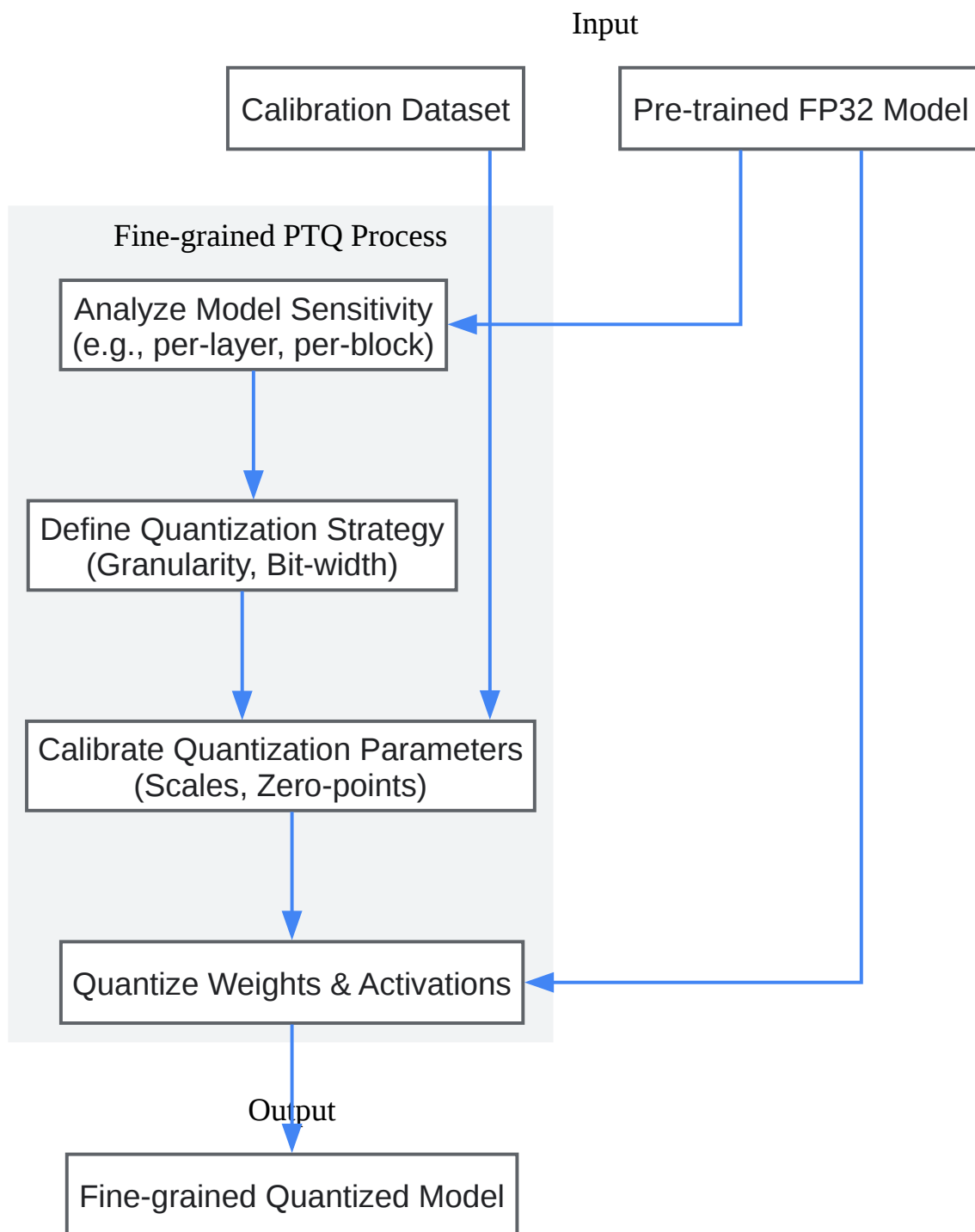
The primary granularities in fine-grained quantization include:

- Per-channel quantization: This method applies a unique scaling factor for each channel in a convolutional layer's weight tensor.[6] It is a widely supported and effective technique for preserving accuracy in convolutional neural networks.[7]

- Group-wise/Block-wise quantization: Here, the tensor is divided into smaller blocks or groups, and a separate scaling factor is applied to each. This is particularly effective for quantizing large language models (LLMs) to very low bit-widths (e.g., 4-bit), as it can better handle the presence of outlier values within the weights.[8][9][10]

Fine-grained quantization often employs a mixed-precision strategy, where different layers or even different parts of a single layer are quantized to different bit-widths based on their sensitivity.[5][11][12] More sensitive components might be kept at a higher precision (e.g., 8-bit or even 16-bit floating-point), while less sensitive parts can be aggressively quantized to lower bit-widths (e.g., 4-bit or 3-bit).[5][13] This selective application of quantization strength allows for a better trade-off between model compression and accuracy.

## Logical Flow of Fine-grained Post-Training Quantization

The general workflow for applying fine-grained post-training quantization involves several key steps. The following diagram illustrates this process, from analyzing the pre-trained model to deploying the optimized version.

Input

Calibration Dataset

Pre-trained FP32 Model

Fine-grained PTQ Process

Analyze Model Sensitivity
(e.g., per-layer, per-block)

Define Quantization Strategy
(Granularity, Bit-width)

Calibrate Quantization Parameters
(Scales, Zero-points)

Quantize Weights & Activations

Output

Fine-grained Quantized Model

Click to download full resolution via product page

Caption: A generalized workflow for fine-grained post-training quantization.
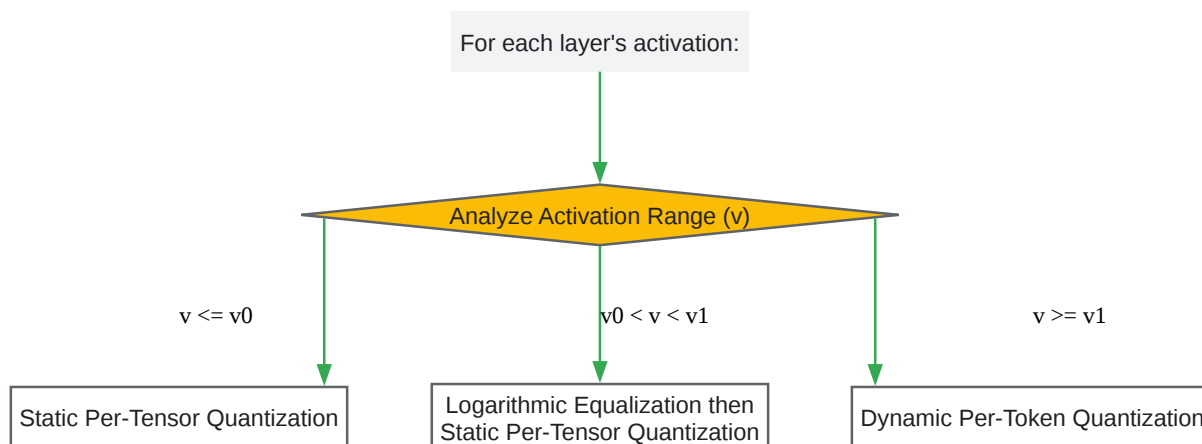
# Key Methodologies in Fine-grained Post-Training Quantization

Several advanced methodologies have been developed to implement fine-grained PTQ effectively. These methods often introduce novel techniques to mitigate the accuracy loss associated with aggressive quantization.

## FPTQ: Fine-grained Post-Training Quantization

**FPTQ** is a method that enables W4A8 quantization (4-bit weights and 8-bit activations) for large language models.[11][14][15] A key challenge in W4A8 quantization is the performance degradation that can occur. **FPTQ** addresses this by employing layer-wise activation quantization strategies, including a novel logarithmic equalization technique for layers that are difficult to quantize, combined with fine-grained weight quantization.[14][16][17] This approach avoids the need for fine-tuning after quantization.[11][14]

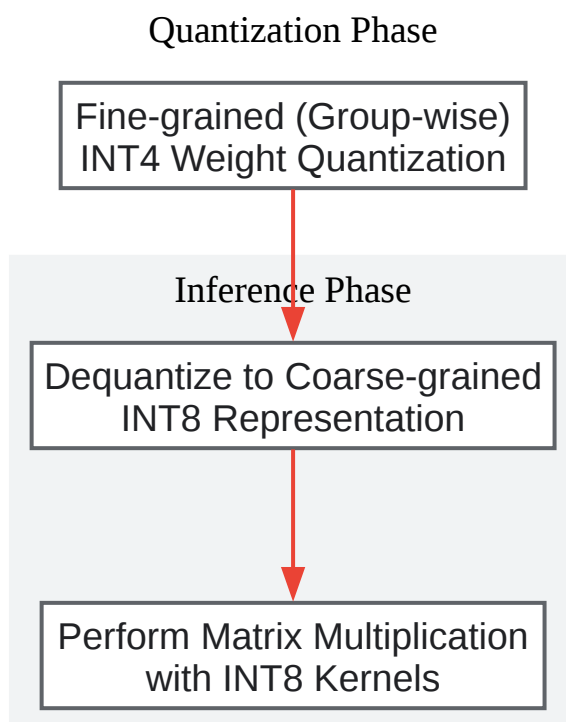The decision-making process within **FPTQ** for handling activations is illustrated below:

For each layer's activation:

Analyze Activation Range (v)

v <= v0

v0 < v < v1

v >= v1

Static Per-Tensor Quantization

Logarithmic Equalization then Static Per-Tensor Quantization

Dynamic Per-Token Quantization

Click to download full resolution via product page

Caption: **FPTQ**'s layer-wise activation quantization decision logic.

# Dual Grained Quantization (DGQ)

DGQ is another novel A8W4 quantization scheme for LLMs that aims to balance performance and inference speed.[9][18] It addresses a key drawback of fine-grained quantization: the disruption of continuous integer matrix multiplication, which can lead to inefficient inference.[9] DGQ dequantizes the fine-grained INT4 weights into a coarse-grained INT8 representation to perform matrix multiplication using efficient INT8 kernels.[9][18][19] The method also includes a two-phase grid search algorithm to determine the optimal fine-grained and coarse-grained quantization scales and a percentile clipping mechanism to handle activation outliers.[9][18]

The core workflow of DGQ is as follows:

Quantization Phase

Fine-grained (Group-wise)
INT4 Weight Quantization

Inference Phase

Dequantize to Coarse-grained
INT8 Representation

Perform Matrix Multiplication
with INT8 Kernels

Click to download full resolution via product page

Caption: The dual-grained quantization and inference process of DGQ.

# FGMP: Fine-Grained Mixed-Precision Quantization

FGMP is a hardware-software co-design methodology for mixed-precision PTQ.[20] It aims to maintain high accuracy by quantizing the majority of weights and activations to a lower

Tech Support

precision while keeping more sensitive parts at a higher precision.[16][20] A key innovation in FGMP is a policy that uses Fisher information to weight the perturbation of each value, thereby identifying which blocks of weights and activations should be kept in higher precision to minimize the impact on the model's loss.[20] FGMP also introduces a sensitivity-weighted clipping approach for the blocks that are quantized to lower precision.[20]

# FineQ: Software-Hardware Co-Design for Low-Bit Quantization

FineQ is another software-hardware co-design approach that focuses on low-bit, fine-grained mixed-precision quantization.[5] It partitions weights into very fine-grained clusters and considers the distribution of outliers within these clusters to strike a balance between model accuracy and memory overhead.[5] A notable feature of FineQ is its outlier protection mechanism, which uses 3 bits to represent outliers within a cluster.[5] The hardware component of FineQ is an accelerator that uses temporal coding to efficiently support the quantization algorithm.[5]

# Experimental Protocols and Data Presentation

The effectiveness of fine-grained PTQ methods is typically evaluated on large language models using standard academic benchmarks.

## Experimental Setup

- Models: The experiments often utilize a range of open-source large language models from families like LLaMA, BLOOM, and OPT of varying sizes (e.g., 7B, 13B, 70B parameters).[21][22]

- Datasets: For performance evaluation, standard datasets are used. Perplexity, a measure of how well a probability model predicts a sample, is often calculated on datasets like WikiText-2 and C4.[23][24] Zero-shot performance on various downstream tasks is also a common evaluation metric.[8]

- Calibration: A small, representative dataset is used to calibrate the quantization parameters (scales and zero-points). This calibration set typically consists of a few hundred to a thousand samples.

- Evaluation Metrics:

  - Perplexity (PPL): Lower perplexity indicates better model performance.

  - Zero-shot Task Accuracy: The accuracy of the quantized model on various downstream tasks without any task-specific fine-tuning.

  - Memory Footprint: The reduction in model size after quantization.

  - Inference Speedup: The improvement in inference latency and throughput.[4]

  - Energy Efficiency: The reduction in energy consumption during inference.[7]

# Quantitative Data Summary

The following tables summarize the performance of different fine-grained PTQ methods on various models and datasets as reported in the literature.

Table 1: Perplexity on WikiText-2 for LLaMA Models

| Model | Original FP16 PPL | Quantization Method | Average Bit-width | Quantized PPL | Perplexity Degradation |
|-------|-------------------|---------------------|-------------------|---------------|------------------------|
| LLaMA-2-7B | - | FGMP (NVFP4/FP8) | - | <1% degradation vs FP8 | <1% vs FP8 |
| LLaMA-2-7B | 5.41 | FineQ | 2.33 | 5.89 | 0.48 |
| LLaMA-2-13B | 4.88 | FineQ | 2.33 | 5.34 | 0.46 |
| LLaMA-2-70B | 3.53 | GGUF (4-bit) | 4 | 3.61 | 0.08 |

Note: Direct comparison between methods can be challenging due to variations in experimental setups. The data is aggregated from multiple sources.[20][22][25]

Table 2: Performance of DGQ on LLaMA-2-7B

| Metric | A16W16 (Baseline) | A16W4 | DGQ (A8W4) |
|---|---|---|---|
| Perplexity (WikiText-2) | 5.41 | - | ~5.71 |
| Memory Reduction | 1x | - | 1.12x vs A16W4 |
| Speedup | 1x | - | 3.24x vs A16W4 |

Source: Data synthesized from the DGQ papers.[9][18][19]

Table 3: FineQ Performance on C4 Dataset

| Model | Original FP16 PPL | FineQ (2.33-bit) PPL | Perplexity Degradation |
|---|---|---|---|
| LLaMA-2-7B | 7.92 | 8.35 | 0.43 |
| LLaMA-2-13B | 7.18 | 7.59 | 0.41 |

Source: Data from the FineQ paper.[25]

# Conclusion

Fine-grained post-training quantization represents a significant advancement in the efficient deployment of large-scale neural networks. By moving beyond coarse, per-tensor quantization and adopting more granular strategies like per-channel, group-wise, and mixed-precision quantization, researchers and practitioners can achieve substantial reductions in model size, memory bandwidth, and power consumption with minimal degradation in model accuracy. Methodologies such as **FPTQ**, DGQ, FGMP, and FineQ demonstrate the ongoing innovation in this field, pushing the boundaries of what is possible with low-bit quantization. For professionals in computationally intensive domains like drug development, leveraging these techniques can unlock the potential of state-of-the-art models in resource-constrained environments, accelerating research and discovery.

---

***Need Custom Synthesis?***

*BenchChem offers custom synthesis for rare earth carbides and specific isotopiclabeling.*

*Email: info@benchchem.com or Request Quote Online.*

---

# References

- 1. jshapira.com [jshapira.com]

- 2. [2504.19746] FineQ: Software-Hardware Co-Design for Low-Bit Fine-Grained Mixed-Precision Quantization of LLMs [arxiv.org]

- 3. Optimizing LLMs for Performance and Accuracy with Post-Training Quantization | NVIDIA Technical Blog [developer.nvidia.com]

- 4. apxml.com [apxml.com]

- 5. FineQ: Software-Hardware Co-Design for Low-Bit Fine-Grained Mixed-Precision Quantization of LLMs [arxiv.org]

- 6. [PDF] FGMP: Fine-Grained Mixed-Precision Weight and Activation Quantization for Hardware-Accelerated LLM Inference | Semantic Scholar [semanticscholar.org]

- 7. FineQ: Software-Hardware Co-Design for Low-Bit Fine-Grained Mixed-Precision Quantization of LLMs | IEEE Conference Publication | IEEE Xplore [ieeexplore.ieee.org]

- 8. [2402.16775] A Comprehensive Evaluation of Quantization Strategies for Large Language Models [arxiv.org]

- 9. Dual Grained Quantization: Efficient Fine-grained Quantization for LLM | OpenReview [openreview.net]

- 10. m.youtube.com [m.youtube.com]

- 11. [PDF] FPTQ: Fine-grained Post-Training Quantization for Large Language Models | Semantic Scholar [semanticscholar.org]

- 12. researchgate.net [researchgate.net]

- 13. ojs.aaai.org [ojs.aaai.org]

- 14. researchgate.net [researchgate.net]

- 15. [2308.15987] FPTQ: Fine-grained Post-Training Quantization for Large Language Models [arxiv.org]

- 16. researchgate.net [researchgate.net]

- 17. [2502.13178] Benchmarking Post-Training Quantization in LLMs: Comprehensive Taxonomy, Unified Evaluation, and Comparative Analysis [arxiv.org]

- 18. [2310.04836] Dual Grained Quantization: Efficient Fine-Grained Quantization for LLM [arxiv.org]

- 19. arxiv.org [arxiv.org]

- 20. [2504.14152] FGMP: Fine-Grained Mixed-Precision Weight and Activation Quantization for Hardware-Accelerated LLM Inference [arxiv.org]

- 21. arxiv.org [arxiv.org]

- 22. lesswrong.com [lesswrong.com]

- 23. Daily Papers - Hugging Face [huggingface.co]

- 24. arxiv.org [arxiv.org]

- 25. themoonlight.io [themoonlight.io]

- To cite this document: BenchChem. [what is Fine-grained Post-Training Quantization]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15621169#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