

# Evaluating FPTQ Quantized LLMs: A Comparative Guide

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

## Compound of Interest

Compound Name: FPTQ

Cat. No.: B2542558

[Get Quote](#)

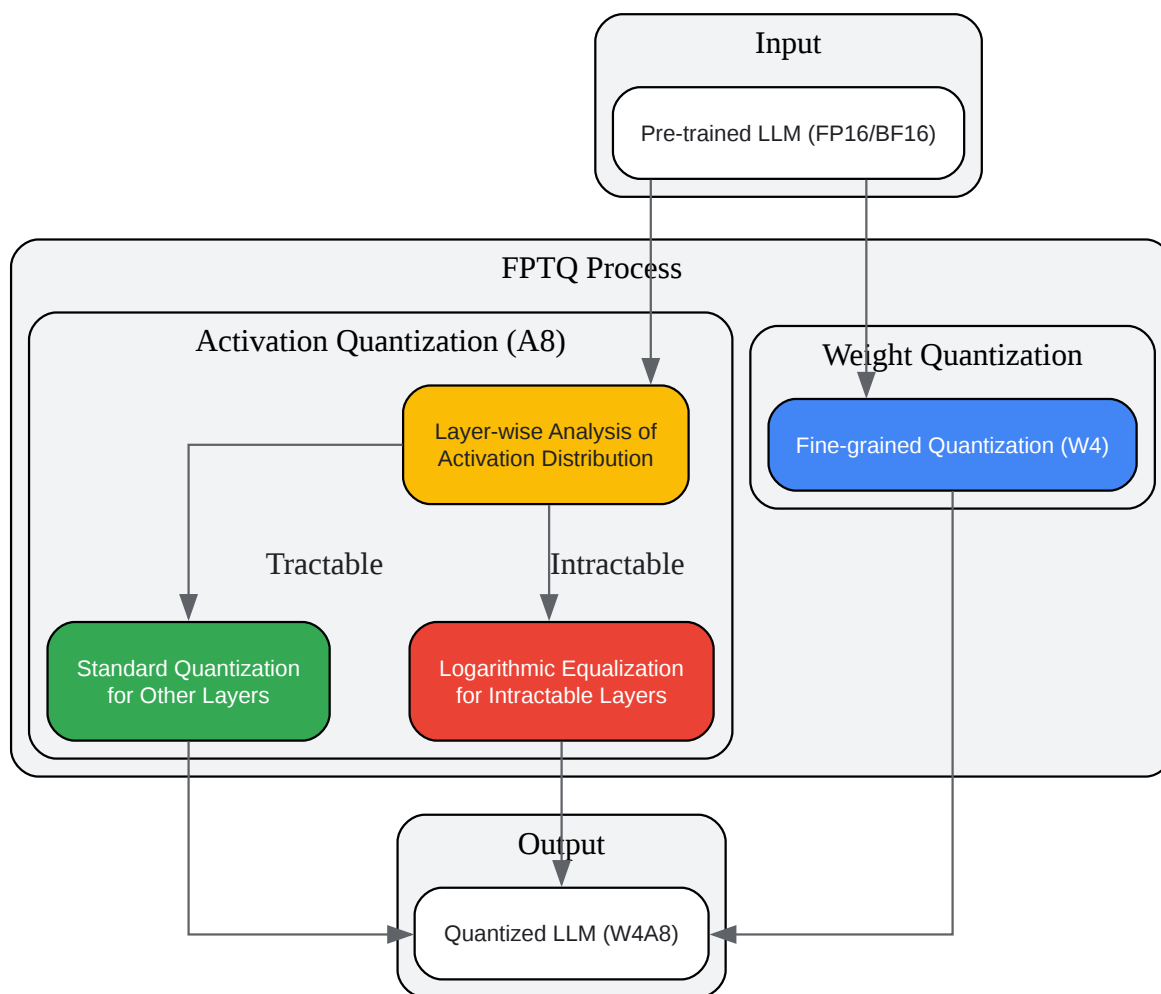
In the rapidly evolving landscape of Large Language Models (LLMs), post-training quantization (PTQ) has emerged as a critical technique for compressing these massive models, enabling their deployment in resource-constrained environments. This guide provides a comparative analysis of the Fine-grained Post-Training Quantization (**FPTQ**) method against other prominent PTQ alternatives, supported by experimental data and detailed methodologies.

## Introduction to FPTQ

**FPTQ** is a post-training quantization method designed to address the performance degradation often seen in low-bit quantization of LLMs.<sup>[1]</sup> It introduces a novel W4A8 (4-bit weights, 8-bit activations) quantization scheme that combines fine-grained weight quantization with a layer-wise activation quantization strategy.<sup>[1][2]</sup> A key innovation in **FPTQ** is the use of logarithmic equalization for layers that are difficult to quantize, which helps to create a more quantization-friendly distribution of activation values.<sup>[3]</sup> This approach aims to leverage the I/O utilization benefits of 4-bit weight quantization and the computational acceleration of 8-bit matrix operations without the need for extensive retraining.<sup>[1]</sup>

## FPTQ Workflow

The **FPTQ** process can be visualized as a multi-step workflow that strategically applies different quantization techniques to the weights and activations of a pre-trained LLM.



[Click to download full resolution via product page](#)

FPTQ Workflow Diagram

## Experimental Protocols

The evaluation of **FPTQ** and other PTQ methods typically involves a standardized set of protocols to ensure fair and reproducible comparisons.

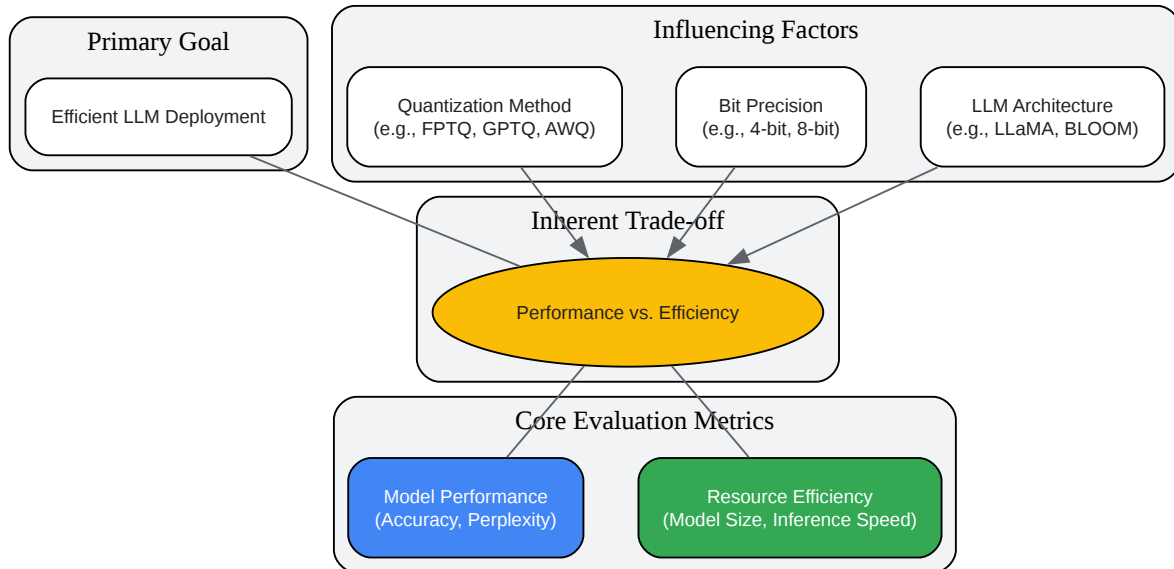
Models and Datasets:

- **Models:** Experiments are conducted on a range of open-source LLMs of varying sizes, such as the BLOOM and LLaMA series.[\[3\]](#)
- **Calibration Data:** A small, representative dataset is used to determine the quantization parameters. For instance, a subset of the C4 dataset is often utilized for this purpose.
- **Evaluation Benchmarks:** The performance of the quantized models is assessed on various downstream tasks, including:
  - **Common Sense Reasoning:** Datasets like Common Sense QA are used to evaluate the model's reasoning capabilities.[\[3\]](#)
  - **Massive Multitask Language Understanding (MMLU):** This benchmark tests the model's knowledge across a wide range of subjects.[\[3\]](#)
  - **Perplexity:** Measured on datasets like WikiText2 to assess the model's language modeling fluency.

**Quantization Procedure:** The core of the **FPTQ** method involves a layer-wise approach to activation quantization. An analysis is performed to identify layers with activation distributions that are challenging to quantize. These "intractable" layers then undergo logarithmic equalization to reshape their distribution, making them more amenable to quantization. The remaining layers are quantized using standard techniques. This is combined with a fine-grained, group-wise quantization of the model's weights.

## Evaluation Criteria for Quantized LLMs

The effectiveness of a quantization method is evaluated based on a trade-off between several key metrics. The logical relationship between these criteria is illustrated in the diagram below.



[Click to download full resolution via product page](#)

### Evaluation Criteria Relationship

## Performance Comparison

The following tables summarize the performance of **FPTQ** on various LLMs and provide a comparative context with other popular PTQ methods.

Disclaimer: The results presented below are compiled from different research papers. Direct comparison may not be entirely fair due to potential variations in the experimental setups, including the specific calibration datasets and evaluation frameworks used.

## FPTQ Performance on LLaMA and BLOOM Models

| Model       | Method | MMLU | Common Sense QA |
|-------------|--------|------|-----------------|
| LLaMA-7B    | FP16   | 63.4 | 75.1            |
| FPTQ (W4A8) | 63.1   | 74.8 |                 |
| LLaMA-13B   | FP16   | 68.9 | 77.3            |
| FPTQ (W4A8) | 68.5   | 77.0 |                 |
| LLaMA-30B   | FP16   | 74.8 | 79.2            |
| FPTQ (W4A8) | 74.3   | 78.9 |                 |
| LLaMA-65B   | FP16   | 77.6 | 80.5            |
| FPTQ (W4A8) | 77.1   | 80.1 |                 |
| BLOOM-7B1   | FP16   | -    | 71.2            |
| FPTQ (W4A8) | -      | 70.9 |                 |

Data sourced from the **FPTQ** paper.[\[3\]](#)

## Comparative Performance of Other PTQ Methods

This table presents results from a broader benchmark study on various PTQ methods.

| Model       | Method | Bit-width | MMLU |
|-------------|--------|-----------|------|
| LLaMA-2-7B  | FP16   | 16        | 63.9 |
| GPTQ        | 4      | 62.8      | 63.9 |
| AWQ         | 4      | 63.1      |      |
| QuIP        | 4      | 62.5      |      |
| LLaMA-2-13B | FP16   | 16        | 69.8 |
| GPTQ        | 4      | 68.7      | 69.8 |
| AWQ         | 4      | 69.0      |      |
| QuIP        | 4      | 68.3      |      |
| LLaMA-2-70B | FP16   | 16        | 77.4 |
| GPTQ        | 4      | 76.5      | 77.4 |
| AWQ         | 4      | 76.8      |      |
| QuIP        | 4      | 76.2      |      |

Note: These results are from a general PTQ benchmark study and may not be directly comparable to the **FPTQ** results due to differences in the LLaMA model versions and evaluation setups.

## Conclusion

**FPTQ** presents a promising approach to W4A8 quantization, demonstrating minimal performance degradation on several large language models.[3] Its innovative use of logarithmic equalization for handling activation outliers appears to be effective in preserving model accuracy.[3] While direct, controlled comparisons with other state-of-the-art methods like GPTQ and AWQ are not readily available in existing literature, the reported results for **FPTQ** are highly competitive.

For researchers and drug development professionals, the ability to deploy powerful LLMs on local or specialized hardware without significant performance loss is a compelling advantage. **FPTQ**, along with other advanced PTQ techniques, is a critical area of research that promises

to make these powerful AI tools more accessible and efficient for a wide range of scientific applications. Further research with comprehensive, head-to-head benchmark comparisons will be crucial for a definitive assessment of the relative strengths and weaknesses of each quantization method.

#### Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: [info@benchchem.com](mailto:info@benchchem.com) or [Request Quote Online](#).

## References

- 1. [2405.04532] QServe: W4A8KV4 Quantization and System Co-design for Efficient LLM Serving [arxiv.org]
- 2. researchgate.net [researchgate.net]
- 3. What Makes Quantization for Large Language Models Hard? An Empirical Study from the Lens of Perturbation [arxiv.org]
- To cite this document: BenchChem. [Evaluating FPTQ Quantized LLMs: A Comparative Guide]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b2542558#evaluating-fptq-quantized-llms]

#### 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](mailto:info@benchchem.com)