# FPTQ Performance on MMLU and Other Benchmarks: A Comparative Guide

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

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

For Researchers, Scientists, and Drug Development Professionals

This guide provides an objective comparison of Few-shot Parameter-efficient Transfer learning via Quantization (**FPTQ**) performance against other state-of-the-art methods on the Massive Multask Language Understanding (MMLU) benchmark and other relevant evaluations. Experimental data, detailed methodologies, and workflow visualizations are presented to offer a comprehensive overview for researchers and professionals in computationally intensive fields.

## Executive Summary

**FPTQ** is a post-training quantization (PTQ) method designed to reduce the memory footprint and computational cost of large language models (LLMs) by quantizing weights to 4-bit integers and activations to 8-bit integers (W4A8). This approach is particularly advantageous for deploying large models in resource-constrained environments. While direct MMLU benchmark scores were not extensively detailed in the primary **FPTQ** publication, subsequent discussions and comparisons have provided insights into its performance. This guide synthesizes available data to compare **FPTQ** with other prominent quantization and parameter-efficient fine-tuning (PEFT) techniques.

## Performance Comparison

The following tables summarize the performance of **FPTQ** and other methods on the MMLU benchmark and other common sense reasoning tasks. MMLU is a widely recognized

benchmark designed to evaluate the knowledge and problem-solving abilities of LLMs across 57 diverse subjects.[1][2][3][4]

Table 1: **FPTQ** Performance on MMLU (LLaMA Models)

In a direct comparison with a Quantization-Aware Training (QAT) method, **FPTQ**, a post-training quantization technique, demonstrated superior performance on the MMLU benchmark for LLaMA models under a W4A8 setting.[5]

| Model | Method | Humanities | STEM | Social Sciences | Other | Average |
|---|---|---|---|---|---|---|
| LLaMA-7B | FPTQ (W4A8) | 45.2 | 35.7 | 47.9 | 44.9 | 43.4 |
| | LLM-QAT (W4A8) | 41.5 | 34.1 | 44.5 | 41.7 | 40.4 | |
| LLaMA-13B | FPTQ (W4A8) | 52.8 | 41.1 | 57.1 | 52.2 | 50.8 |
| | LLM-QAT (W4A8) | 48.1 | 38.2 | 52.3 | 48.1 | 46.7 | |

Table 2: MMLU Performance of Various Quantization Methods on Llama 3 Models

This table provides a broader comparison of different quantization techniques on Llama 3 models, highlighting the trade-offs between model size and accuracy. The results indicate that FP8 quantization maintains the highest fidelity to the baseline FP16 performance.[6]

| Model | Method | MMLU Score |
|---|---|---|
| Llama 3 8B | FP16 (Baseline) | 68.9 |
| FP8 | 68.7 | |
| INT8 SQ | 68.5 | |
| INT4 AWQ | 67.2 | |
| Llama 3 70B | FP16 (Baseline) | 79.9 |
| FP8 | 79.8 | |
| INT8 SQ | 79.6 | |
| INT4 AWQ | 78.5 | |

Table 3: Comparison of Parameter-Efficient Fine-Tuning (PEFT) Methods

PEFT methods offer an alternative to full fine-tuning by only updating a small subset of a model's parameters. This significantly reduces computational and storage costs. Low-Rank Adaptation (LoRA) is a popular PEFT technique that has been shown to achieve performance comparable to full fine-tuning with a fraction of the trainable parameters.[7][8][9][10]

Tech Support

| Method | Trainable Parameters (vs. Full Fine-Tuning) | Performance vs. Full Fine-Tuning | Key Advantages |
|---|---|---|---|
| Full Fine-Tuning | 100% | Baseline | Highest potential performance, but computationally expensive. |
| LoRA | ~0.01% - 1% | 97-99% of Full Fine-Tuning performance | Drastically reduced memory and storage footprint; faster training.[7] |
| Adapters | ~0.1% - 5% | Comparable to Full Fine-Tuning | Modular and composable; can be added to models without altering original weights. |

# Experimental Protocols

**FPTQ** Evaluation on MMLU

The evaluation of **FPTQ** on the MMLU benchmark, as inferred from standard practices and the authors' descriptions, likely followed a 5-shot evaluation protocol.[11] This involves providing the model with five examples from a specific task within the MMLU benchmark before presenting the actual question. This few-shot approach assesses the model's ability to generalize from a small number of examples.

The core of the **FPTQ** methodology involves a layer-wise quantization strategy that adapts to the distribution of activations in different layers of the neural network. A key component is logarithmic equalization for layers with challenging activation distributions. The process is a post-training approach, meaning it is applied to an already trained LLM without the need for costly retraining.[12][13]
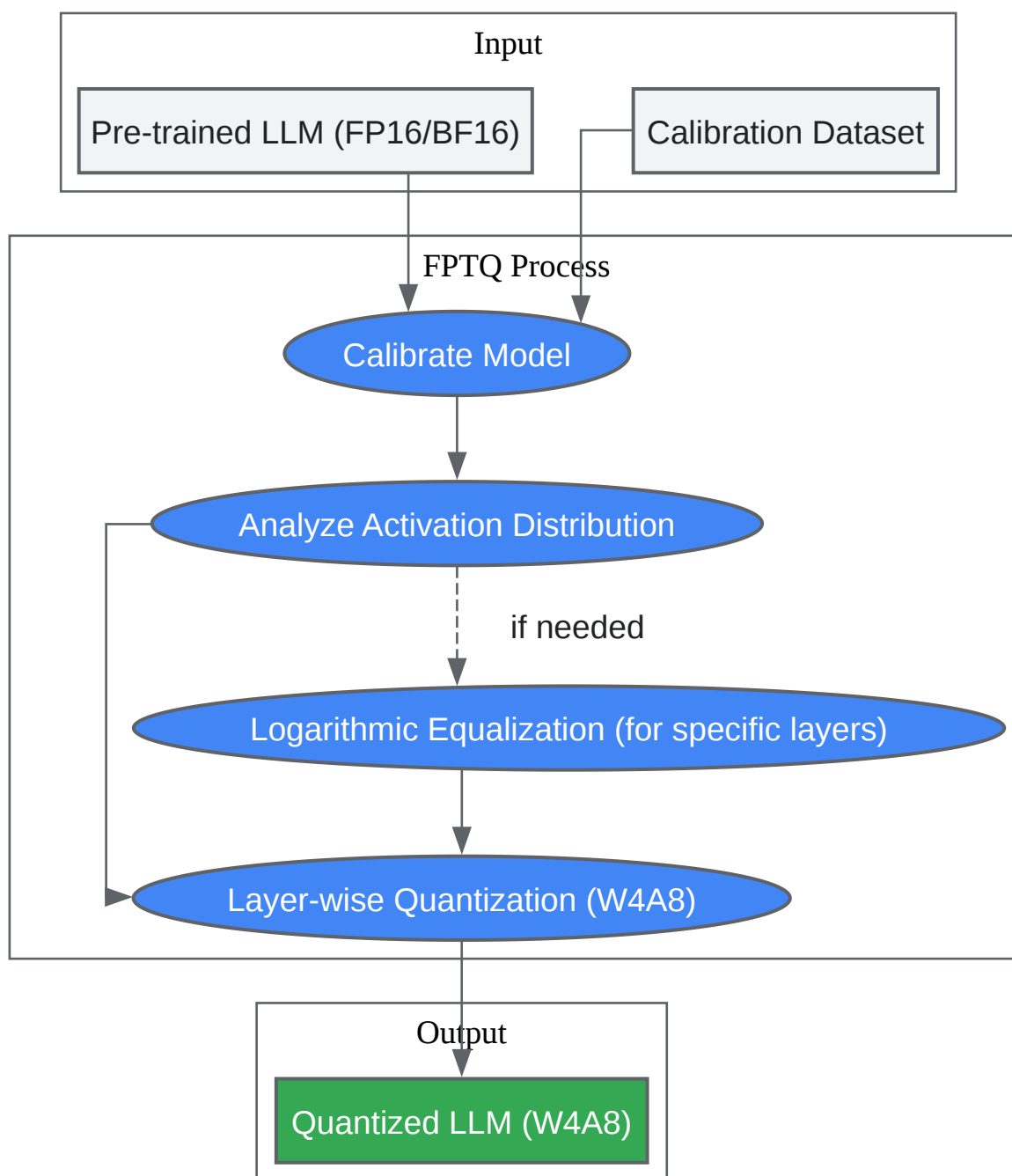
General MMLU Evaluation

The MMLU benchmark consists of multiple-choice questions across 57 subjects.[2][4] Performance is typically measured by the model's accuracy in selecting the correct answer. Evaluations are conducted in either a zero-shot or few-shot setting.[4] In a few-shot setting, a small number of example questions and their correct answers are provided in the prompt to give the model context for the task.

# Visualizing **FPTQ** and its Application

**FPTQ** Workflow

The following diagram illustrates the general workflow of the Fine-grained Post-Training Quantization (**FPTQ**) process.

**Input**

Pre-trained LLM (FP16/BF16)

Calibration Dataset

**FPTQ Process**

Calibrate Model

Analyze Activation Distribution

*if needed*

Logarithmic Equalization (for specific layers)

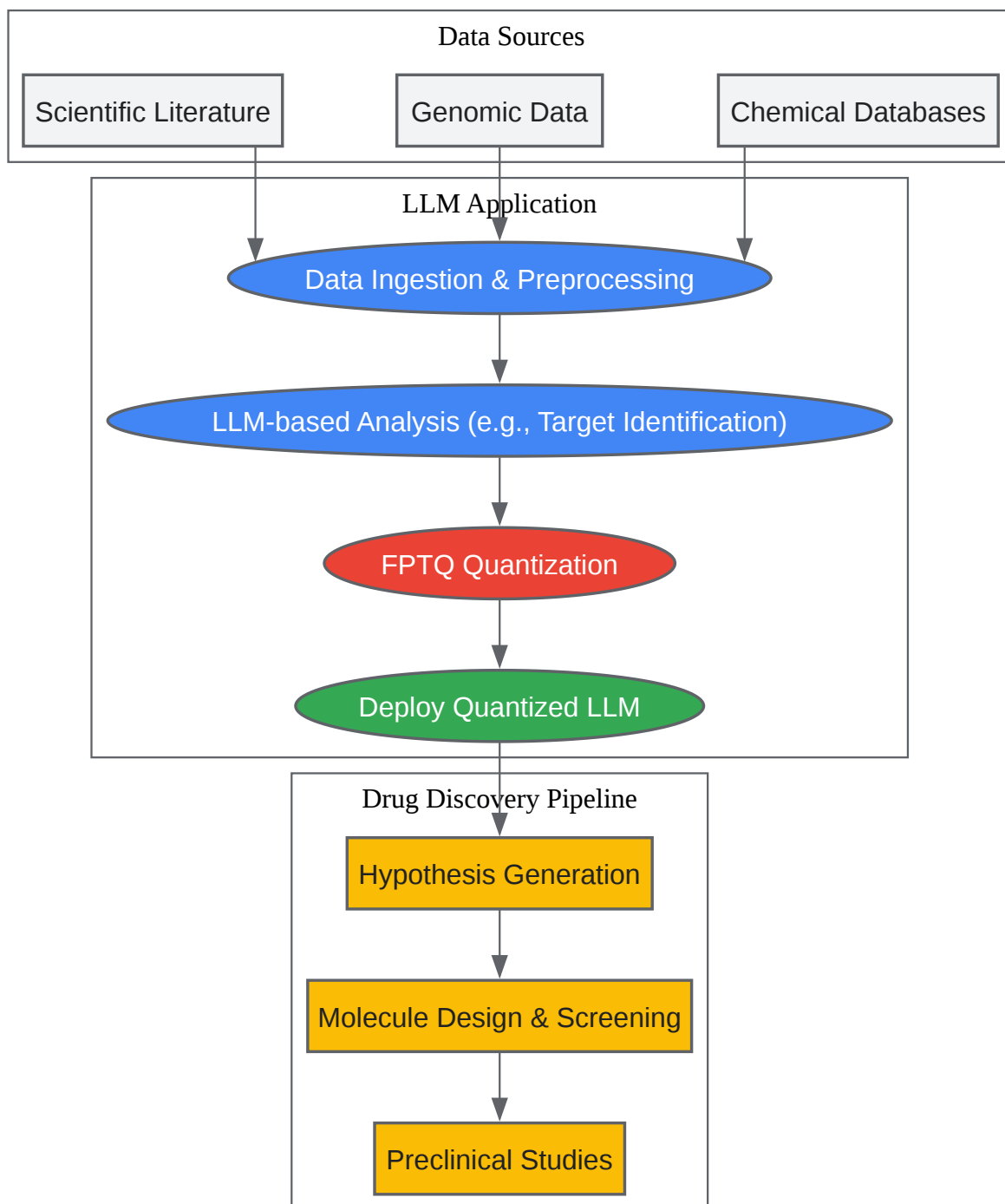Layer-wise Quantization (W4A8)

**Output**

Quantized LLM (W4A8)

Click to download full resolution via product page

Caption: A diagram illustrating the **FPTQ** workflow.

Conceptual Application in Drug Development

While **FPTQ** is a general-purpose model optimization technique, its application in drug development can be conceptualized in workflows that leverage large language models for tasks like scientific literature analysis, target identification, and molecule property prediction. By quantizing these LLMs, researchers can deploy them more efficiently, for instance, on local servers for proprietary data analysis without relying on cloud APIs.

Caption: Conceptual workflow for **FPTQ** in drug development.

***Need Custom Synthesis?***

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

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

# References

- 1. Why MMLU Remains the Most Honest AI Benchmark in 2025 [graphlogic.ai]

- 2. MMLU - Wikipedia [en.wikipedia.org]

- 3. Paper Breakdown #1: MMLU — LLMs Have Exams Too! A Post on Benchmarking | by Alakarthika Ulaganathan | Medium [medium.com]

- 4. galileo.ai [galileo.ai]

- 5. FPTQ: FINE-GRAINED POST-TRAINING QUANTIZATION FOR LARGE LANGUAGE MODELS | OpenReview [openreview.net]

- 6. Post-Training Quantization of LLMs with NVIDIA NeMo and NVIDIA TensorRT Model Optimizer | NVIDIA Technical Blog [developer.nvidia.com]

- 7. beautifulcode.co [beautifulcode.co]

- 8. Gen AI Fine-Tuning Techniques: LoRA, QLoRA, and Adapters Compared [digitaldividedata.com]

- 9. Comparison between parameter-efficient techniques and full fine-tuning: A case study on multilingual news article classification - PMC [pmc.ncbi.nlm.nih.gov]

- 10. predibase.com [predibase.com]

- 11. openreview.net [openreview.net]

- 12. arxiv.org [arxiv.org]

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

- To cite this document: BenchChem. [FPTQ Performance on MMLU and Other Benchmarks: A Comparative Guide]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15621169#fptq-performance-on-mmlu-and-other-benchmarks]

**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