

A Comparative Guide to Post-Training Quantization: FPTQ vs. GPTQ and AWQ

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

Compound of Interest

Compound Name: FPTQ

Cat. No.: B15621169

[Get Quote](#)

In the rapidly evolving landscape of large language models (LLMs), post-training quantization (PTQ) has emerged as a critical technique for reducing model size and accelerating inference without the need for costly retraining. This guide provides a detailed comparison of a novel PTQ method, Fine-grained Post-training Quantization (**FPTQ**), against two other prominent methods: Generative Pre-trained Transformer Quantization (GPTQ) and Activation-aware Weight Quantization (AWQ). This document is intended for researchers, scientists, and drug development professionals who are leveraging LLMs in their work and require a deeper understanding of the trade-offs associated with different quantization techniques.

Core Methodologies

Post-training quantization methods aim to convert the high-precision floating-point weights and, in some cases, activations of a neural network to lower-precision integer formats. This reduction in bit-width leads to a smaller memory footprint and can enable faster computations on compatible hardware. The primary challenge lies in minimizing the accuracy degradation that can occur during this conversion. **FPTQ**, GPTQ, and AWQ employ distinct strategies to address this challenge.

FPTQ (Fine-grained Post-training Quantization) is a method that focuses on a W4A8 quantization scheme, meaning it quantizes weights to 4-bit integers and activations to 8-bit integers.^{[1][2][3]} This approach seeks to combine the benefits of reduced memory bandwidth from 4-bit weights with the computational efficiency of 8-bit matrix operations for activations.^[1] ^[2] A key innovation in **FPTQ** is its layer-wise approach to activation quantization, which

employs a novel logarithmic equalization technique for layers that are more difficult to quantize. [1][2][3] For the most challenging layers, **FPTQ** resorts to a per-token dynamic quantization approach.[4]

GPTQ (Generative Pre-trained Transformer Quantization) is a one-shot, layer-wise weight quantization method that aims to minimize the quantization error by using approximate second-order (Hessian) information.[5][6] This allows GPTQ to achieve very low bit-widths, such as 3 or 4 bits, with minimal loss of accuracy.[5] GPTQ is a weight-only quantization method, meaning it primarily focuses on compressing the model's weights while activations are typically kept at a higher precision (e.g., 16-bit floating point).

AWQ (Activation-aware Weight Quantization) is another weight-only quantization method that operates on the principle that not all weights are equally important.[5][6] AWQ identifies salient weights by analyzing the distribution of activations and protects these important weights from the full impact of quantization.[5][6] This is achieved by applying a scaling factor to the weights based on the activation magnitudes, which helps to preserve the model's performance, especially for instruction-tuned LLMs.[7]

Quantitative Performance Comparison

The following tables summarize the performance of **FPTQ**, GPTQ, and AWQ on various benchmarks. It is important to note that the experimental setups for these results may differ across the original sources. Please refer to the experimental protocols section for more details.

Performance on LLaMA Models

Model	Method	WikiText-2 (PPL)	C4 (PPL)
LLaMA-7B	FP16 (Baseline)	5.86	8.35
FPTQ (W4A8)	6.15	8.65	
SmoothQuant (W8A8)	6.25	8.78	
GPTQ (W4A16)	6.07	8.57	
LLaMA-13B	FP16 (Baseline)	5.09	7.37
FPTQ (W4A8)	5.38	7.69	
SmoothQuant (W8A8)	5.34	7.64	
GPTQ (W4A16)	5.26	7.57	
LLaMA-30B	FP16 (Baseline)	4.18	6.22
FPTQ (W4A8)	4.31	6.36	
SmoothQuant (W8A8)	4.32	6.38	
GPTQ (W4A16)	4.29	6.35	
LLaMA-65B	FP16 (Baseline)	3.65	5.56
FPTQ (W4A8)	3.73	5.66	
SmoothQuant (W8A8)	3.75	5.67	
GPTQ (W4A16)	3.71	5.63	

Note: The **FPTQ**, SmoothQuant, and GPTQ results for LLaMA models are from the **FPTQ** paper. The GPTQ variant used is W4A16.

Performance on BLOOM-7B1

Method	PIQA (acc)	ARC-e (acc)	BoolQ (acc)
FP16 (Baseline)	79.2	69.9	71.8
FPTQ (W4A8)	79.1	69.7	71.2
SmoothQuant (W8A8)	78.8	69.5	70.9
GPTQ (W4A16)	79.0	69.6	71.1

Note: The **FPTQ**, SmoothQuant, and GPTQ results for the BLOOM-7B1 model are from the **FPTQ** paper. The GPTQ variant used is W4A16.

MMLU Benchmark (LLaMA Models)

Model	Method	MMLU (5-shot acc)
LLaMA-7B	FP16 (Baseline)	45.3
FPTQ (W4A8)	43.8	
SmoothQuant (W8A8)	41.2	
GPTQ (W4A16)	44.2	
LLaMA-13B	FP16 (Baseline)	54.8
FPTQ (W4A8)	52.1	
SmoothQuant (W8A8)	51.5	
GPTQ (W4A16)	53.9	
LLaMA-30B	FP16 (Baseline)	62.4
FPTQ (W4A8)	60.2	
SmoothQuant (W8A8)	59.8	
GPTQ (W4A16)	61.5	
LLaMA-65B	FP16 (Baseline)	67.2
FPTQ (W4A8)	65.8	
SmoothQuant (W8A8)	65.1	
GPTQ (W4A16)	66.8	

Note: The **FPTQ**, SmoothQuant, and GPTQ results for the MMLU benchmark are from the **FPTQ** paper. The GPTQ variant used is W4A16.

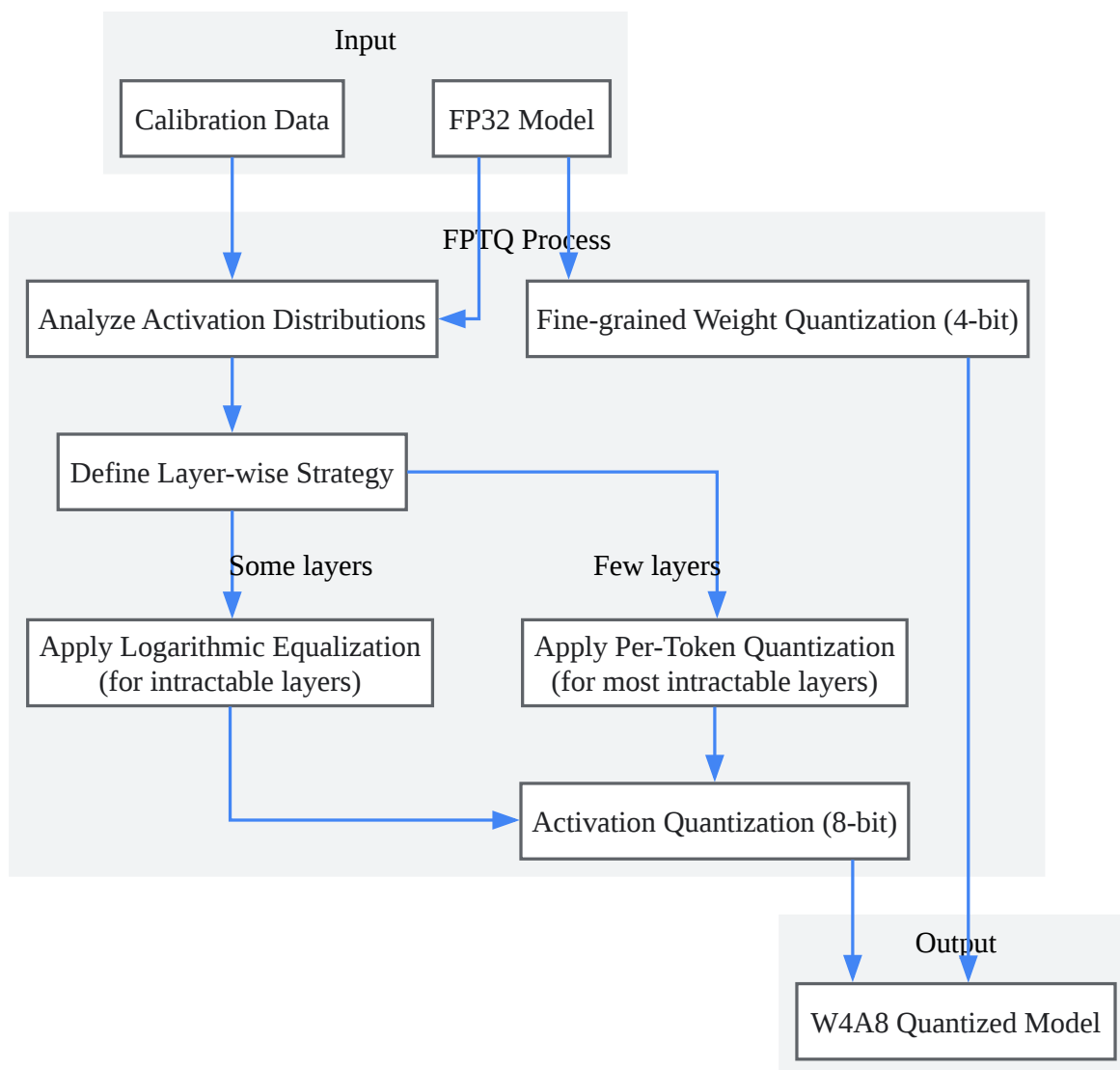
Experimental Protocols

FPTQ: The experiments for **FPTQ** were conducted using a calibration set of 128 random samples from the C4 dataset, with each sample having a sequence length of 2048. The evaluation was performed on WikiText-2 and C4 for perplexity, and on commonsense reasoning benchmarks (PIQA, ARC-e, BoolQ) and MMLU for accuracy.

GPTQ and AWQ (General): For GPTQ and AWQ, the calibration process also typically involves a small number of samples (e.g., 128) from a representative dataset like C4 or WikiText. The group size for quantization is a crucial hyperparameter, with a common choice being 128. The performance of GPTQ and AWQ can be influenced by the choice of calibration data, and some studies have shown that GPTQ can be prone to overfitting to the calibration set.[8]

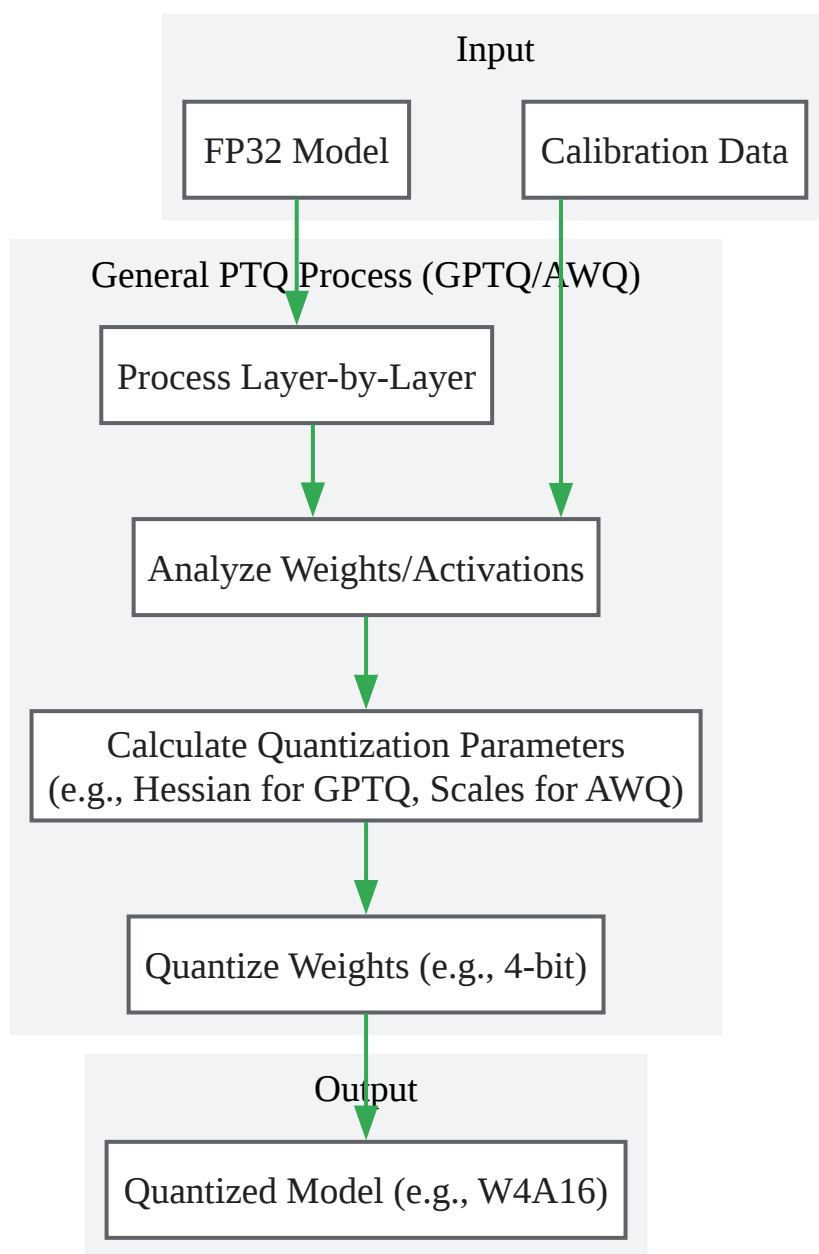
Visualizing the Workflows

To better understand the underlying processes of these quantization methods, the following diagrams, generated using the DOT language, illustrate their respective workflows.



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

Caption: **FPTQ** workflow, highlighting the layer-wise strategy for activation quantization.



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

Caption: General workflow for weight-only PTQ methods like GPTQ and AWQ.

Conclusion

FPTQ presents a compelling W4A8 quantization strategy that aims to balance memory savings and computational efficiency. The available data suggests that **FPTQ** is competitive with, and in

some cases outperforms, other PTQ methods like SmoothQuant and a W4A16 variant of GPTQ, particularly in maintaining performance on commonsense reasoning tasks.

GPTQ and AWQ are powerful weight-only quantization techniques that can achieve significant model compression with minimal accuracy loss. The choice between them may depend on the specific model, the importance of quantization speed (AWQ is generally faster to apply than GPTQ), and the nature of the downstream tasks.[9]

For researchers and professionals in drug development and other scientific fields, the choice of a PTQ method will depend on a careful consideration of the trade-offs between model accuracy, inference speed, and memory footprint. While **FPTQ** shows promise as a W4A8 solution, the maturity and widespread adoption of GPTQ and AWQ make them strong contenders for weight-only quantization. As the field of LLM quantization continues to evolve, it is crucial to stay informed about new methods and to perform in-house evaluations on specific use cases to determine the optimal approach.

Need Custom Synthesis?

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

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

References

- 1. aclanthology.org [aclanthology.org]
- 2. GitHub - pprp/Awesome-LLM-Quantization: Awesome list for LLM quantization [github.com]
- 3. [2303.13003] Benchmarking the Reliability of Post-training Quantization: a Particular Focus on Worst-case Performance [arxiv.org]
- 4. esann.org [esann.org]
- 5. themoonlight.io [themoonlight.io]
- 6. medium.com [medium.com]
- 7. Which Quantization Method Is Best for You?: GGUF, GPTQ, or AWQ... | E2E Networks [e2enetworks.com]
- 8. Why LLM Benchmarks Can Be Misleading - AWQ vs. GPTQ [bitbasti.com]

- 9. apxml.com [apxml.com]
- To cite this document: BenchChem. [A Comparative Guide to Post-Training Quantization: FPTQ vs. GPTQ and AWQ]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15621169#benchmarking-fptq-against-other-ptq-methods]

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