

Technical Support Center: Mitigating Accuracy Drop in W4A8 FPTQ

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

Compound of Interest

Compound Name: FPTQ

Cat. No.: B2542558

[Get Quote](#)

This technical support center provides troubleshooting guides and frequently asked questions (FAQs) to help researchers, scientists, and drug development professionals mitigate accuracy drops during 4-bit weight and 8-bit activation (W4A8) Fine-grained Post-Training Quantization (FPTQ) of their models.

Troubleshooting Guides

Issue: Significant accuracy degradation after W4A8 Post-Training Quantization.

Root Cause: A common cause of accuracy drop in W4A8 quantization is the significant loss of precision, especially for weights, and the presence of outliers in activations that disrupt the quantization range.^{[1][2]} Naive quantization of both weights and activations to low bit-widths can lead to severe performance degradation.^{[2][3]}

Solution: Implement a fine-grained, layer-wise quantization strategy. Not all layers are equally sensitive to quantization. By analyzing the distribution of activations in each layer, you can apply more sophisticated quantization techniques only to the most problematic layers.

Experimental Protocol:

- **Layer-wise Analysis:** Profile the activation ranges for each layer of your model using a representative calibration dataset.

- **Identify Problematic Layers:** Identify layers with large activation outliers or asymmetric distributions. The **FPTQ** methodology suggests that layers with activation ranges between 15 and 150 are particularly challenging.[\[2\]](#)
- **Apply Conditional Quantization:**
 - For layers with significant outliers, apply Logarithmic Activation Equalization (LAE) to make the activation distribution more quantization-friendly.[\[2\]](#)
 - For other sensitive layers, consider alternative strategies like channel-wise shifting and scaling.[\[4\]](#)
 - For less sensitive layers, a standard per-token dynamic quantization may suffice.[\[2\]](#)
- **Fine-Grained Weight Quantization:** Utilize group-wise quantization for weights to provide more flexibility and reduce quantization error.[\[2\]](#)

Issue: Model accuracy is highly sensitive to the calibration dataset.

Root Cause: The calibration dataset is crucial for determining the quantization parameters (scale and zero-point). If the calibration data is not representative of the data the model will see during inference, the learned quantization parameters will be suboptimal, leading to a drop in accuracy.

Solution: Employ a Sequence-Length-Aware Calibration (SLAC) strategy. The variation in activation diversity can be related to the input sequence length. Calibrating the model with sequence lengths that are representative of the target task can mitigate accuracy losses.[\[5\]](#)

Experimental Protocol:

- **Analyze Target Task Sequence Lengths:** Determine the typical sequence lengths of inputs for your specific use case (e.g., drug-protein interaction prediction, scientific literature analysis).
- **Create a Representative Calibration Dataset:** Construct a calibration dataset with a distribution of sequence lengths that mirrors your target task.

- Perform Calibration: Use this tailored dataset to perform post-training quantization. This will ensure the quantization parameters are optimized for the expected inputs.

Frequently Asked Questions (FAQs)

Q1: What is W4A8 **FPTQ** and why is it challenging?

A1: W4A8 Fine-grained Post-Training Quantization is a technique to compress large models by representing their weights with 4-bit integers and activations with 8-bit integers.[1][2] This combination is advantageous as it reduces the memory footprint due to 4-bit weights and allows for faster computation using 8-bit matrix operations.[2][3] The primary challenge is the significant performance degradation that can occur due to the aggressive quantization of weights and the difficulty in quantizing activations without losing critical information, especially in the presence of outliers.[2][3]

Q2: What are "outliers" in activations and how do they impact quantization?

A2: Outliers are activation values that are significantly larger in magnitude than the majority of other activation values within a tensor. These outliers can skew the quantization range. When using a min-max quantization scheme, a single large outlier can force the vast majority of other values into a very small portion of the quantization grid, leading to a significant loss of precision for these values. Mitigating the effect of outliers is a key focus of advanced quantization methods.[4][6]

Q3: What is Logarithmic Activation Equalization (LAE)?

A3: Logarithmic Activation Equalization is a technique used in **FPTQ** to handle layers with challenging activation distributions.[2] It applies a logarithmic function to the activations to compress the range of values, making them more amenable to quantization. This is particularly effective for distributions with large outliers.

Q4: Should I use integer or floating-point formats for quantization?

A4: While integer (INT) quantization is common, floating-point (FP) quantization (e.g., FP8 for activations and FP4 for weights) can offer superior performance, especially for large language models.[7] FP formats can better represent values with a wide dynamic range, which can help

mitigate issues with outliers. However, hardware support (like NVIDIA's H100 GPUs) is a consideration for FP quantization.[7]

Q5: Can I improve accuracy without retraining the model?

A5: Yes, all the techniques discussed here are part of Post-Training Quantization (PTQ), which does not require retraining or fine-tuning.[1][2] Methods like **FPTQ**, Outlier Suppression+, and using a representative calibration dataset are designed to be applied to an already trained model.

Quantitative Data Summary

The following tables summarize the performance of W4A8 **FPTQ** compared to other quantization methods on common benchmarks for various Large Language Models (LLMs).

Table 1: Performance on the LAMBADA Dataset

Model	Method	Accuracy
LLaMA-7B	FP16 (Original)	75.28%
SmoothQuant W8A8	74.01%	
FPTQ W4A8	73.80%	
LLaMA-13B	FP16 (Original)	78.31%
SmoothQuant W8A8	77.83%	
FPTQ W4A8	77.74%	
LLaMA-30B	FP16 (Original)	80.01%
SmoothQuant W8A8	79.78%	
FPTQ W4A8	79.82%	

Data sourced from the **FPTQ** paper.[2]

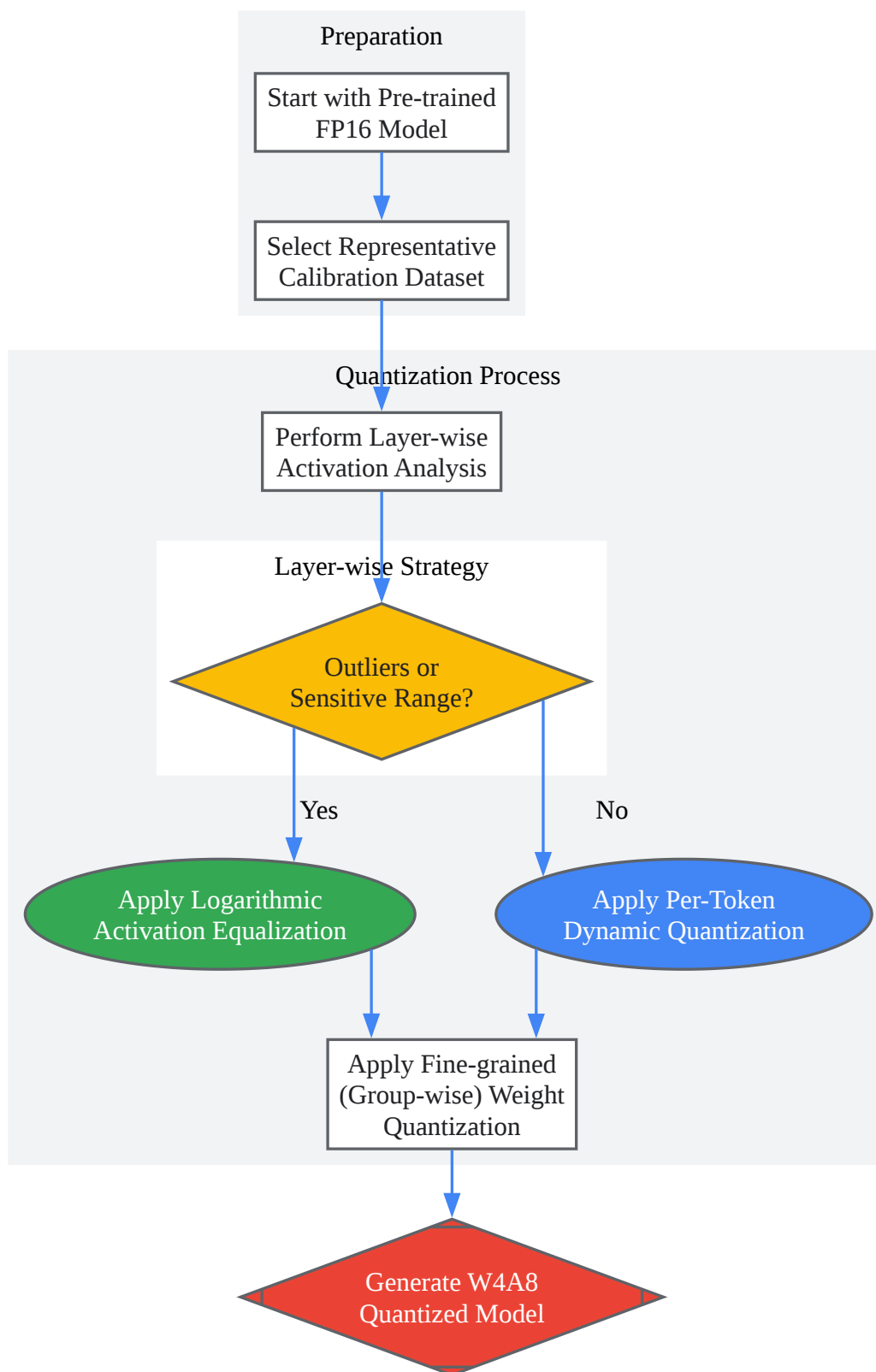
Table 2: Performance on Common Sense QA Datasets

Model	Method	PIQA	HS	ARCe	Avg.
LLaMA-7B	FP16 (Original)	78.4	78.8	53.4	70.2
LLM-QAT W4A8	77.2	77.3	51.9	68.8	
FPTQ W4A8	77.9	78.4	52.2	69.5	
LLaMA-13B	FP16 (Original)	79.8	81.0	58.1	73.0
LLM-QAT W4A8	79.1	80.1	55.1	71.4	
FPTQ W4A8	79.5	80.5	56.5	72.2	

Data sourced from the **FPTQ** paper, comparing against a Quantization-Aware Training (QAT) method.[\[8\]](#)

Visualizations

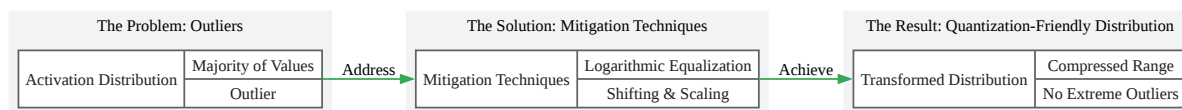
Experimental Workflow for W4A8 FPTQ



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

Caption: Workflow for applying Fine-grained Post-Training Quantization.

Logical Relationship of Outlier Mitigation



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

Caption: Conceptual diagram of mitigating activation outliers.

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. paperreading.club [paperreading.club]
- 2. FPTQ: FINE-GRAINED POST-TRAINING QUANTIZATION FOR LARGE LANGUAGE MODELS | OpenReview [openreview.net]
- 3. [2308.15987] FPTQ: Fine-grained Post-Training Quantization for Large Language Models [arxiv.org]
- 4. Outlier Suppression+: Accurate quantization of large language models by equivalent and effective shifting and scaling | OpenReview [openreview.net]
- 5. Enhancing Computation Efficiency in Large Language Models through Weight and Activation Quantization | OpenReview [openreview.net]
- 6. Rethinking the Outlier Distribution in Large Language Models: An In-depth Study [arxiv.org]
- 7. [2307.09782] ZeroQuant-FP: A Leap Forward in LLMs Post-Training W4A8 Quantization Using Floating-Point Formats [arxiv.org]
- 8. arxiv.org [arxiv.org]

- To cite this document: BenchChem. [Technical Support Center: Mitigating Accuracy Drop in W4A8 FPTQ]. BenchChem, [2025]. [Online PDF]. Available at: [\[https://www.benchchem.com/product/b2542558#mitigating-accuracy-drop-in-w4a8-fptq\]](https://www.benchchem.com/product/b2542558#mitigating-accuracy-drop-in-w4a8-fptq)

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