# Implementing Logarithmic Equalization in Fine-grained Post-Training Quantization (FPTQ)

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

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

Get Quote

Application Notes and Protocols for Researchers, Scientists, and Drug Development Professionals
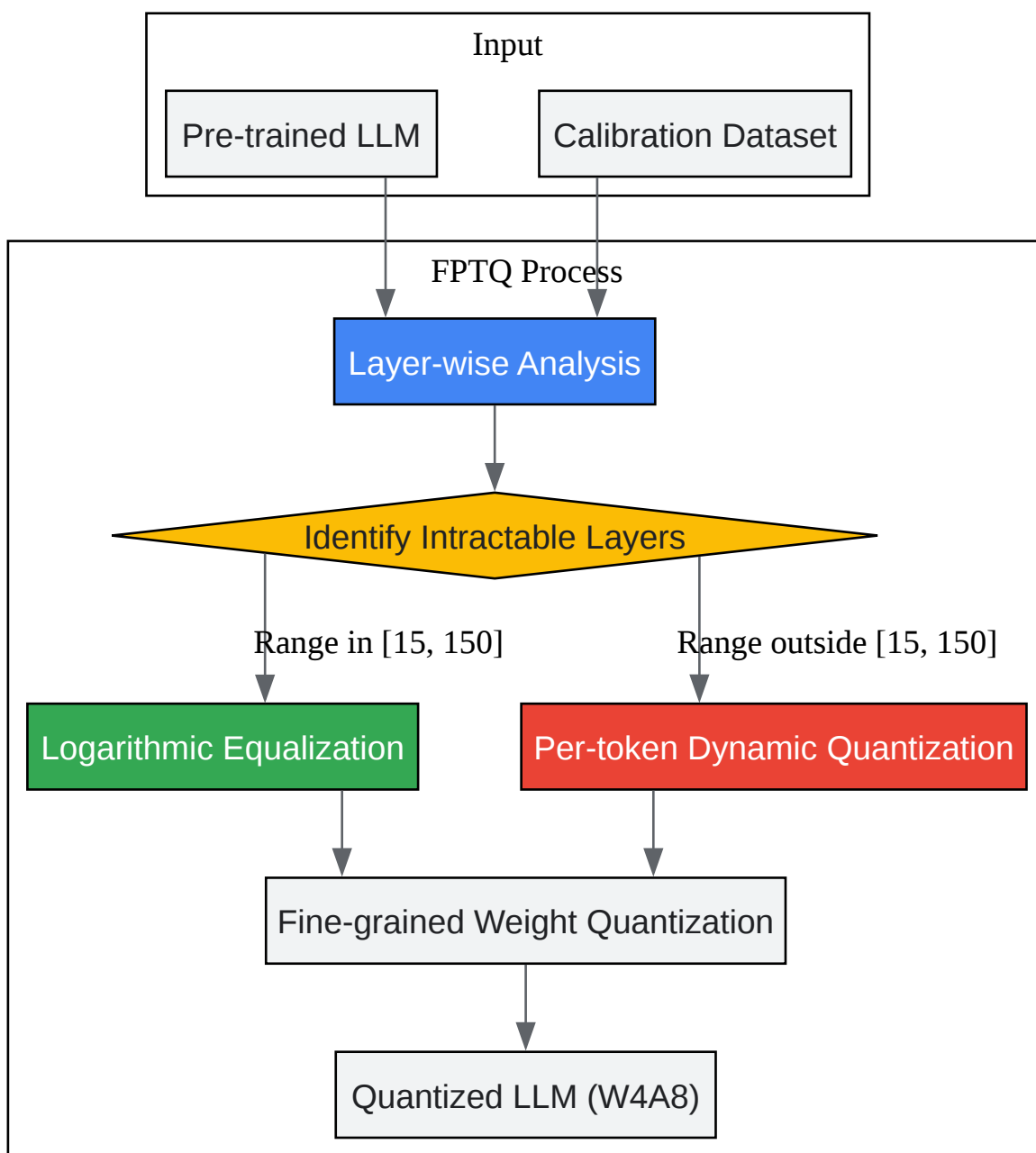
## Introduction

Fine-grained Post-Training Quantization (**FPTQ**) is a novel method for compressing large language models (LLMs), enabling their deployment in resource-constrained environments.[1][2][3][4] A key innovation within **FPTQ** is the application of logarithmic equalization to the activations of specific, challenging layers within the model.[1][2][3][4] This technique, combined with fine-grained weight quantization, allows for a W4A8 (4-bit weights, 8-bit activations) quantization scheme that maintains high model performance without the need for costly retraining.[2][3][5] These application notes provide a detailed overview of the principles behind logarithmic equalization in **FPTQ**, protocols for its implementation, and quantitative data from relevant studies.

## Principle of Logarithmic Equalization in **FPTQ**

In the context of post-training quantization, a significant challenge arises from "intractable layers" where activation values have a wide dynamic range. Standard quantization methods struggle with these layers, leading to substantial performance degradation. **FPTQ** addresses this by employing a layer-wise activation quantization strategy that includes a novel logarithmic equalization for these problematic layers.[2][3]

The core idea is to apply a logarithmic function to the activation values, which compresses the range of high-magnitude values while expanding the range of low-magnitude values. This makes the distribution of activations more uniform and amenable to quantization. The **FPTQ** method specifically applies this logarithmic equalization when the range of activations in a given layer falls between 15 and 150.[5] For layers with activation ranges outside this window, a per-token dynamic quantization approach is used as a fallback.[5] The quantization scaling factor for the equalized activations is determined based on a logarithmic function of the maximum activation values.[5]

## Logical Flow of **FPTQ** with Logarithmic Equalization

Caption: Workflow of the **FPTQ** process, highlighting the conditional application of logarithmic equalization.

# Experimental Protocols

The following protocols are synthesized from the methodologies described in the **FPTQ** research papers and general post-training quantization workflows.

# Environment Setup

Objective: To prepare the necessary computational environment and libraries.

Protocol:

- Install a deep learning framework such as PyTorch.

- Install the Hugging Face transformers library for model loading and manipulation.

- Install libraries for quantization, such as bitsandbytes (for baseline comparisons) and any available **FPTQ**-related packages.

- Ensure access to a CUDA-enabled GPU for efficient model processing.

# Model and Calibration Dataset Preparation

Objective: To load a pre-trained LLM and a representative dataset for calibration.

Protocol:

- Load a pre-trained large language model (e.g., LLaMA, BLOOM) using the transformers library.

- Select a calibration dataset that is representative of the data the model will encounter in the target application. A subset of a common dataset like C4 or WikiText is often used.

- The calibration process involves feeding a small, carefully selected set of input data through the original model to observe the resulting activation values.[6] The goal is to gather statistics about the activation distributions to inform the quantization process.[7]
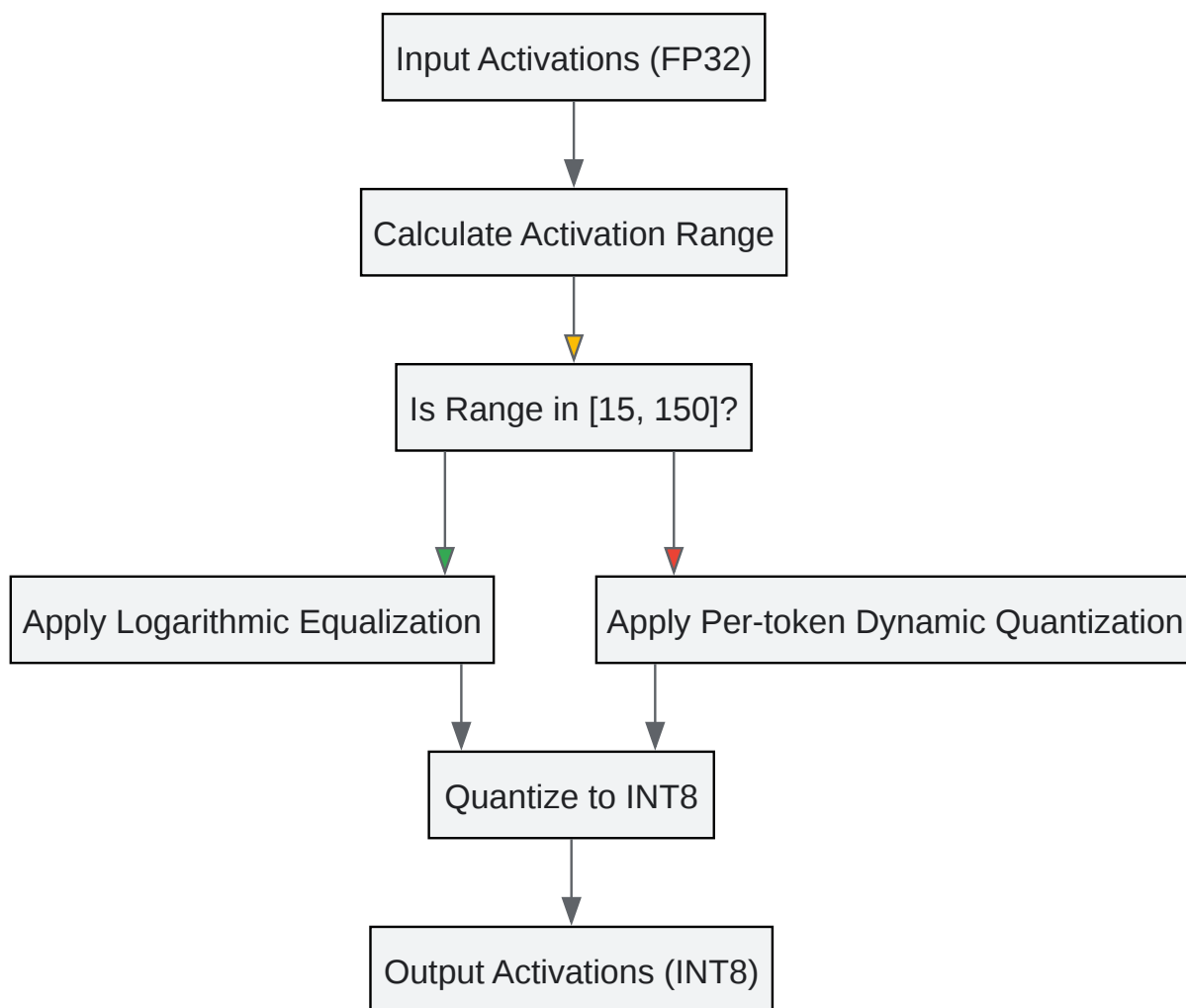
# FPTQ Implementation Protocol

Objective: To apply the **FPTQ** method, including logarithmic equalization, to the loaded model.

Protocol:

- Iterate Through Model Layers: Process the model layer by layer to apply a tailored quantization strategy.

 Tech Support

- Activation Range Analysis: For each layer, pass the calibration dataset through the model and record the minimum and maximum activation values to determine the dynamic range.

- Conditional Logarithmic Equalization:

  - If the activation range for a layer is between 15 and 150, apply logarithmic equalization to the activations. The scaling factor for quantization is then calculated based on a log function of the maximum activation values.

  - If the activation range is outside of this window, apply a standard per-token dynamic quantization to the activations.

- Fine-grained Weight Quantization: Apply group-wise quantization to the model's weights. This involves dividing the weight tensors into smaller groups and quantizing each group independently to 4-bit integers.

- Generate Quantized Model: Save the modified model with the quantized weights and the necessary information for on-the-fly activation quantization.

```
┌─────────────────────────────┐
│   Input Activations (FP32)  │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│   Calculate Activation Range│
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│     Is Range in [15, 150]?  │
└─────────────────────────────┘
       │                │
       ▼                ▼
┌──────────────────┐  ┌────────────────────────────────┐
│ Apply Logarithmic│  │ Apply Per-token Dynamic        │
│ Equalization     │  │ Quantization                   │
└──────────────────┘  └────────────────────────────────┘
         │                     │
         ▼                     ▼
      ┌─────────────────────────┐
      │    Quantize to INT8     │
      └─────────────────────────┘
                  │
                  ▼
      ┌─────────────────────────┐
      │ Output Activations(INT8)│
      └─────────────────────────┘
```

Click to download full resolution via product page

Caption: Decision logic for applying logarithmic equalization within a layer.

# Quantitative Data and Performance Benchmarks

The effectiveness of **FPTQ** has been demonstrated on various LLMs and benchmarked against other quantization methods. The following tables summarize the performance of **FPTQ** in terms of perplexity (a measure of language model quality; lower is better) and task-specific accuracy.

# Table 1: Perplexity on WikiText2 for LLaMA Models

| Model | Original (FP16) | SmoothQuant (W8A8) | FPTQ (W4A8) |
|---|---|---|---|
| LLaMA-7B | 5.34 | 5.35 | 5.34 |
| LLaMA-13B | 4.75 | 4.76 | 4.75 |
| LLaMA-30B | 4.02 | 4.03 | 4.02 |
| LLaMA-65B | 3.65 | 3.67 | 3.66 |

Data sourced from **FPTQ** research submissions.

## Table 2: Zero-shot Accuracy on Common Sense Reasoning Tasks

| Model | Task | Original (FP16) | SmoothQuant (W8A8) | FPTQ (W4A8) |
|---|---|---|---|---|
| LLaMA-7B | PIQA | 78.4 | 78.3 | 78.5 |
| | HellaSwag | 78.8 | 78.6 | 78.7 |
| | WinoGrande | 72.4 | 72.0 | 72.2 |
| LLaMA-13B | PIQA | 80.0 | 79.9 | 80.0 |
| | HellaSwag | 81.3 | 81.1 | 81.2 |
| | WinoGrande | 75.8 | 75.3 | 75.6 |

Data sourced from **FPTQ** research submissions.

## Conclusion

**FPTQ** with logarithmic equalization presents a compelling solution for the efficient deployment of large language models. By strategically applying logarithmic compression to the activations of challenging layers, **FPTQ** achieves a W4A8 quantization scheme with minimal to no performance degradation. The provided protocols and quantitative data serve as a guide for researchers and developers looking to implement and evaluate this advanced quantization technique. The ability to significantly reduce the memory and computational footprint of LLMs,

Tech Support

as demonstrated by **FPTQ**, is a critical step towards their broader application in diverse scientific and industrial domains.

> **Need Custom Synthesis?**
>
> BenchChem offers custom synthesis for rare earth carbides and specific isotopiclabeling.
>
> Email: info@benchchem.com or Request Quote Online.

# References

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

- 2. FPTQ: Fine-grained Post-Training Quantization for Large Language Models [paperreading.club]

- 3. researchgate.net [researchgate.net]

- 4. [2308.15987] FPTQ: Fine-grained Post-Training Quantization for Large Language Models [arxiv.org]

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

- 6. apxml.com [apxml.com]

- 7. medium.com [medium.com]

- To cite this document: BenchChem. [Implementing Logarithmic Equalization in Fine-grained Post-Training Quantization (FPTQ)]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15621169#implementing-logarithmic-equalization-in-fptq]

**Disclaimer & Data Validity:**

**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