# FPTQ for large language models explained

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

| Compound of Interest | |
| --- | --- |
| Compound Name: | FPTQ |
| Cat. No.: | B2542558 |

Get Quote

An In-depth Technical Guide to Fine-grained Post-Training Quantization (**FPTQ**) for Large Language Models

# Introduction

The deployment of large language models (LLMs) is often hindered by their substantial size, which demands significant storage and computational resources.[1][2] Quantization has become a mainstream technique to compress these models and accelerate inference.[3][4] This process primarily revolves around two main strategies: W8A8 (8-bit weights and 8-bit activations) and W4A16 (4-bit weights and 16-bit activations).[5]

This technical guide delves into Fine-grained Post-Training Quantization (**FPTQ**), a novel W4A8 post-training quantization method that synergistically combines the advantages of both popular recipes.[1][2] **FPTQ** leverages the reduced memory input/output (I/O) of 4-bit weight quantization and the computational acceleration of 8-bit matrix operations.[4][6] The primary challenge with a naive W4A8 approach is a significant degradation in model performance.[2][5] **FPTQ** addresses this by employing layer-wise activation quantization strategies, featuring a unique logarithmic equalization for more challenging layers, combined with fine-grained weight quantization.[5][6] This method has demonstrated state-of-the-art performance for W4A8 quantized models like BLOOM, LLaMA, and LLaMA-2 without the need for extensive fine-tuning.[2][4]
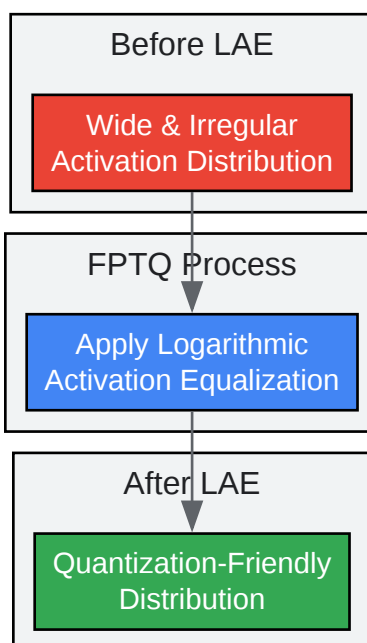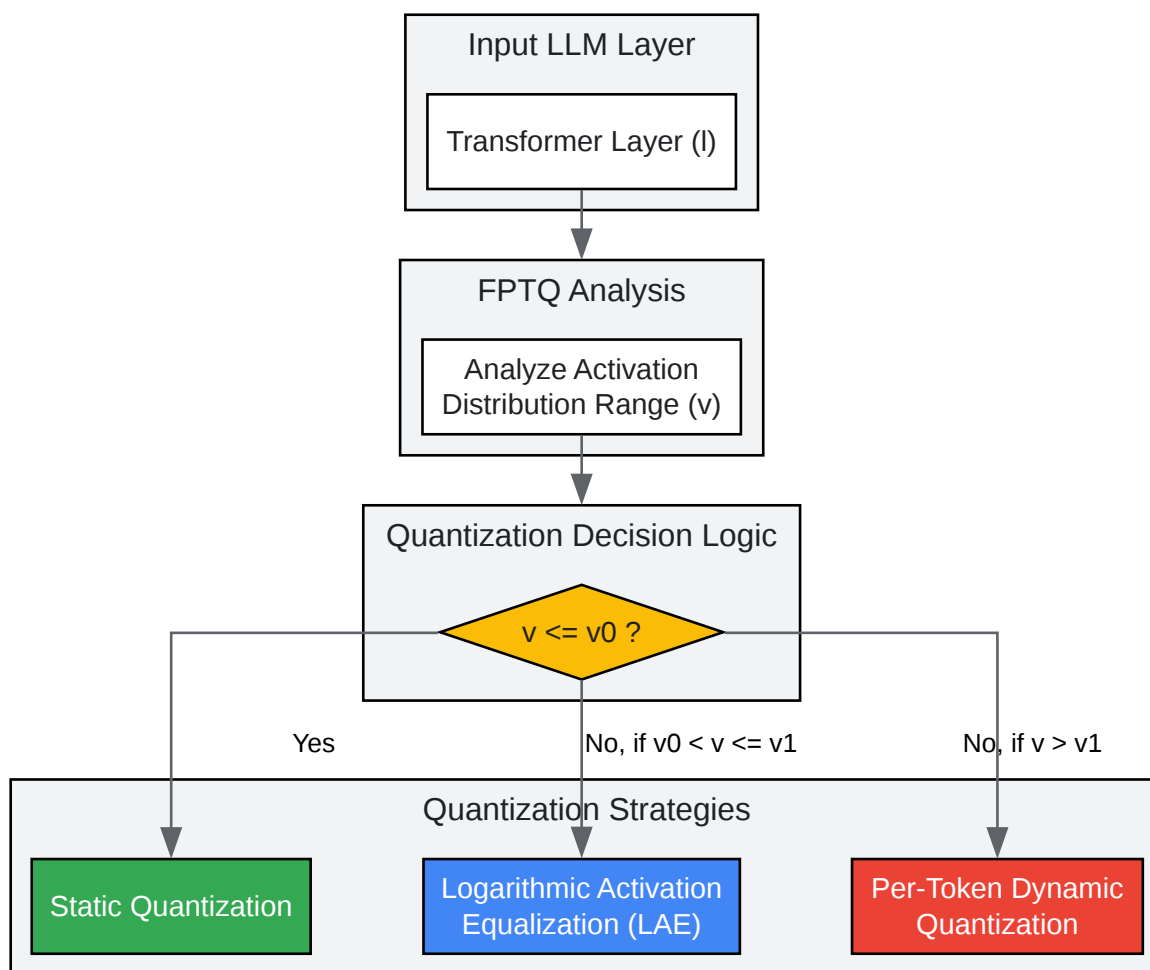
# Core Concepts of **FPTQ**

The fundamental innovation of **FPTQ** is its hybrid approach to quantization that adapts to the different characteristics of layers within a transformer architecture. It recognizes that a one-
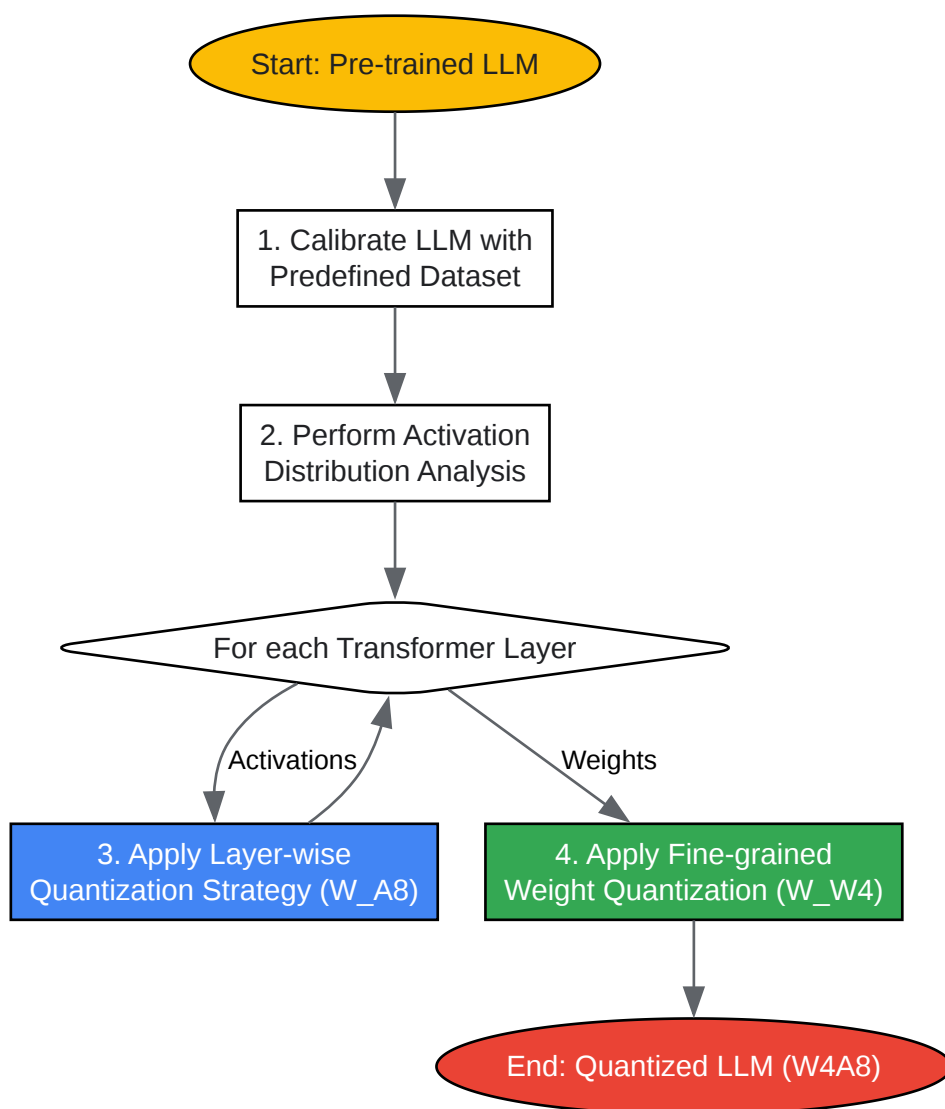
size-fits-all quantization strategy is suboptimal.

## Layer-wise Quantization Strategy

A key observation is that activation distributions vary significantly across different layers of an LLM. Some layers are amenable to simple static quantization, while others exhibit activation ranges that are challenging to quantize without significant error.[1] Applying per-tensor static quantization across all layers can lead to substantial performance loss, whereas using per-token dynamic quantization for all layers introduces computational overhead that can negate the benefits of quantization.[1][5]

**FPTQ** resolves this by implementing a layer-specific policy. It analyzes the activation distribution for each layer and selects the most appropriate quantization granularity, creating a more balanced and efficient model.[1]

Tech Support

## Input LLM Layer

Transformer Layer (l)

## FPTQ Analysis

Analyze Activation Distribution Range (v)

## Quantization Decision Logic

$v <= v0$ ?

Yes

No, if $v0 < v <= v1$

No, if $v > v1$

## Quantization Strategies

Static Quantization

Logarithmic Activation Equalization (LAE)

Per-Token Dynamic Quantization

## Before LAE

Wide & Irregular Activation Distribution

## FPTQ Process

Apply Logarithmic Activation Equalization

## After LAE

Quantization-Friendly Distribution

```
           ┌─────────────────────┐
           │ Start: Pre-trained  │
           │        LLM          │
           └─────────────────────┘
                     │
                     ▼
           ┌─────────────────────┐
           │ 1. Calibrate LLM    │
           │ with Predefined     │
           │      Dataset        │
           └─────────────────────┘
                     │
                     ▼
           ┌─────────────────────┐
           │ 2. Perform          │
           │ Activation          │
           │ Distribution        │
           │ Analysis            │
           └─────────────────────┘
                     │
                     ▼
           ◇ For each Transformer Layer ◇
           Activations              Weights
              │                        │
              ▼                        ▼
   ┌──────────────────┐      ┌──────────────────┐
   │ 3. Apply Layer-  │      │ 4. Apply Fine-   │
   │ wise Quantization│      │ grained Weight   │
   │ Strategy (W_A8)  │      │ Quantization     │
   │                  │      │ (W_W4)           │
   └──────────────────┘      └──────────────────┘
                                      │
                                      ▼
                          ┌──────────────────────┐
                          │ End: Quantized LLM   │
                          │       (W4A8)         │
                          └──────────────────────┘
```

Click to download full resolution via product page

---

***Need Custom Synthesis?***

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

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

---

# References

- 1. openreview.net [openreview.net]

- 2. [2308.15987] FPTQ: Fine-grained Post-Training Quantization for Large Language Models [arxiv.org]

- 3. FPTQ: Fine-grained Post-Training Quantization for Large Language Models | DeepAI [deepai.org]

- 4. researchgate.net [researchgate.net]

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

- 6. Ribbit Ribbit â°ᴹᴰˢᵀˢ Discover Research the Fun Way [ribbitribbit.co]

- To cite this document: BenchChem. [FPTQ for large language models explained]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b2542558#fptq-for-large-language-models-explained]

**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