# PTQTP: Post-Training Quantization to Trit-Planes for Large Language Models

**He Xiao[1]\*, Runming Yang[1]\*, Qingyao Yang[1]\*, Wendong Xu[1],**
**Zhen Li[2], Yupeng Su[3], Zhengwu Liu[1], Hongxia Yang[2], Ngai Wong[1]†**

[1] The University of Hong Kong [2] The Hong Kong Polytechnic University
[3] University of California, Santa Barbara

\* Equal Contribution † Corresponding author: nwong@eee.hku.hk

## Abstract

Post-training quantization (PTQ) of large language models (LLMs) to extremely low bit-widths remains challenging due to the fundamental trade-off between computational efficiency and representational capacity. While existing ultra-low-bit methods rely on binary approximations or quantization-aware training(QAT), they often suffer from either limited representational capacity or huge training resource overhead. We introduce **PTQ** to **T**rit-**P**lanes (PTQTP), a structured PTQ framework that decomposes weight matrices into dual ternary $\{-1, 0, 1\}$ trit-planes. This approach achieves multiplication-free additive inference by decoupling weights into discrete topology (trit-planes) and continuous magnitude (scales), effectively enabling high-fidelity sparse approximation. PTQTP provides: (1) a theoretically grounded progressive approximation algorithm ensuring global weight consistency; (2) model-agnostic deployment without architectural modifications; and (3) uniform ternary operations that eliminate mixed-precision overhead. Comprehensive experiments on LLaMA3.x and Qwen3 (0.6B-70B) demonstrate that PTQTP significantly outperforms sub-4bit PTQ methods on both language reasoning tasks and mathematical reasoning as well as coding. PTQTP rivals the 1.58-bit QAT performance while requiring only single-hour quantization compared to 10-14 GPU days for training-based methods, and the end-to-end inference speed achieves $4.63\times$ faster than the FP16 baseline model, establishing a new and practical solution for efficient LLM deployment in resource-constrained environments. Code will available at https://github.com/HeXiao-55/PTQTP.

## 1 Introduction

The unprecedented success of large language models (LLMs) has revolutionized natural language processing, but their deployment is hindered by exorbitant computational and memory
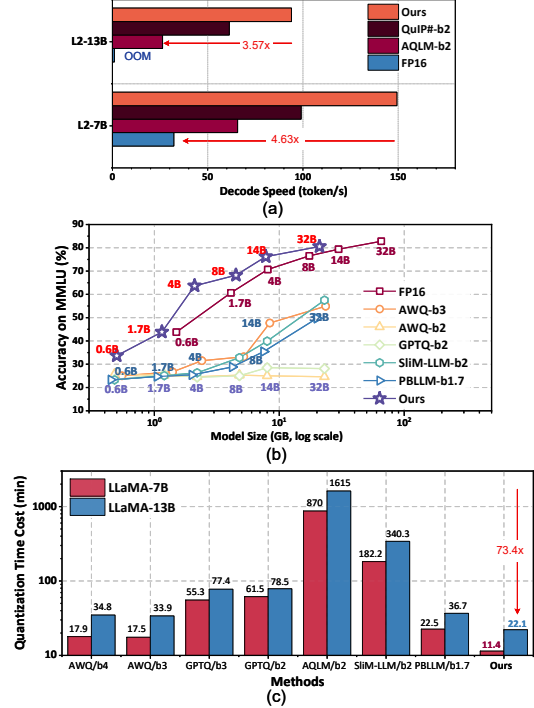


Figure 1: PTQTP minimizes compression costs while maintaining excellent performance. (a) PTQTP can deliver $4.63\times$ end-to-end decode speed on NVIDIA RTX 3090 GPU. (b) The language reasoning evaluation on Qwen3 demonstrates that PTQTP outperforms existing low-bit quantization methods. (c) Quantization runtime comparison shows PTQTP achieves $73.4\times$ speedup over AQLM and $1.57\times$ over AWQ on LLaMA-7B.

demands (Brown et al., 2020). Models such as LLaMA3 (Grattafiori et al., 2024), Qwen3 (Qwen-Team, 2025), DeepSeek-R1-671B (Guo et al., 2025), with hundreds of billions of parameters, require specialized hardware and massive energy consumption, limiting accessibility on edge devices and raising environmental concerns (Patterson et al., 2021). To address this, sub-4 bit low-bit quantization including binary (1-bit) and ternary (1.58-bit), has emerged as a viable approach, facilitating affordable memory consumption aligned with efficient and edge implementation demands.

Post-training quantization (PTQ) offers a practical pathway for compressing pretrained LLMs without retraining, with recent advancements like GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2024), and many exquisite PTQ methods (Shao et al., 2023; Ashkboos et al., 2024; Xiao et al., 2023; Tseng et al., 2024b) achieving effective 4-bit quantization with near full-precision accuracy; Near-2bit PTQ methods (Egiazarian et al., 2024a; Huang et al., 2024b; Xiao et al., 2025; Tseng et al., 2024a; Chee et al., 2023; Zhao et al., 2025) concentrated on solving the outlier challenge by using pre-processing or high-precision protection methods before PTQ to smooth outliers and extra fine-tuning strategies to help refine the performance of the quantized model.

However, 2-4 bit operations still rely on costly multiply-accumulate (MAC) operations on existing hardware. Binary (1-bit) and ternary (1.58-bit) PTQ can effectively reduce algebraic multiplication to addition to save the inference consumption, but push PTQ to extreme bit-widths coexist with challenges and opportunities. ***Two main Challenges***: Binary PTQ that achieves 1-bit quantization through unstructured weight categorization, but sacrifices representational capacity (Huang et al., 2024a; Li et al., 2025; Shang et al., 2023; Zhao et al., 2025); Quantization-aware training (QAT) with extremely low-bit widths, which requires costly retraining, e.g., BitNet (Wang et al., 2023; Ma et al., 2024) for pretraining binary/ternary models. ***Opportunities for Ternary***: Ternary operations further enhance the logic selection capabilities compared to binary operations. Moreover, from the hardware perspective, ternary PTQ is not merely quantization but a step towards transforming LLM from an *arithmetic-bound* to a *logic-bound* operation. This represents a qualitative leap for future neuromorphic computing or edge AI. Structured ternary PTQ, offering higher expressivity than binary while avoiding QAT's pre-training overhead and extra fine-tuning, remains underexplored (if not unexplored).

We introduce PTQTP (**P**ost-**T**raining **Q**uantization to **T**rit-**P**lanes), a novel structured PTQ method that bridges the gap between binary efficiency and ternary expressiveness while avoiding costly training overhead and time-consumed fine-tuning step. Unlike binary quantization, which *forces weights to* $\pm 1$, *introducing significant quantization noise that disrupts precise logic flows*, causing sophisticated reasoning

(such as mathematics) to collapse. PTQTP leverages a magnitude-topology decoupling strategy, it decomposes full-precision weights into a linear superposition of dual ternary trit-planes $\{-1, 0, 1\}$ modulated by a column of row-wise scalars. This essentially performs a constructive interference: the "0" state allows the model to selectively silence noise features. Furthermore, by converting heavy multiplications into lightweight ternary additions, PTQTP offers a practical, uniform, and mathematically robust solution for the deployment of powerful LLMs, opening new frontiers for efficient inference in real-world applications.

1. PTQTP captures both the principal skeleton and the residual details of the weight distribution using a collaborative dual trit-planes structure, producing an expressive space that is more flexible and contains richer sparsity than current 1-3 bit methods. Our work demonstrates that structured ternary quantization strikes a sweet spot between computational simplicity and representational power, as shown in Fig.1.

2. PTQTP preserves complex inference topology at extremely low-bit precision, addressing a critical reasoning performance degradation in the literature. PTQTP offers a robust, model-agnostic solution for extremely low-bit PTQ. By eliminating the need for architecture-specific adjustments, it supports seamless deployment on quantization-sensitive models (e.g., LLaMA3.x, Qwen3), while consistently preserving performance and surpassing state-of-the-art architecture-dependent approaches in scalability.

3. Extensive experiments demonstrate that PTQTP consistently outperforms state-of-the-art low-bit methods, surpassing most 1–3 bit PTQ approaches in language benchmarks while enabling faster quantization and hardware-efficient multiplication-free operations. Remarkably, it rivals or even exceeds 1.58-bit QAT, despite requiring over $10^4 \times$ fewer GPU hours, without any retraining or post-PTQ fine-tuning, underscoring its efficiency and generalizability across advanced models.

## 2 Methodology

### 2.1 Magnitude-Topology Decoupling

Unlike conventional scalar quantization which maps weights directly to a discrete grid, we conceptualize the quantization of LLM weights as a magnitude-topology decoupling problem. Let $W \in \mathbb{R}^{n \times d}$ be the weight matrix. We assume
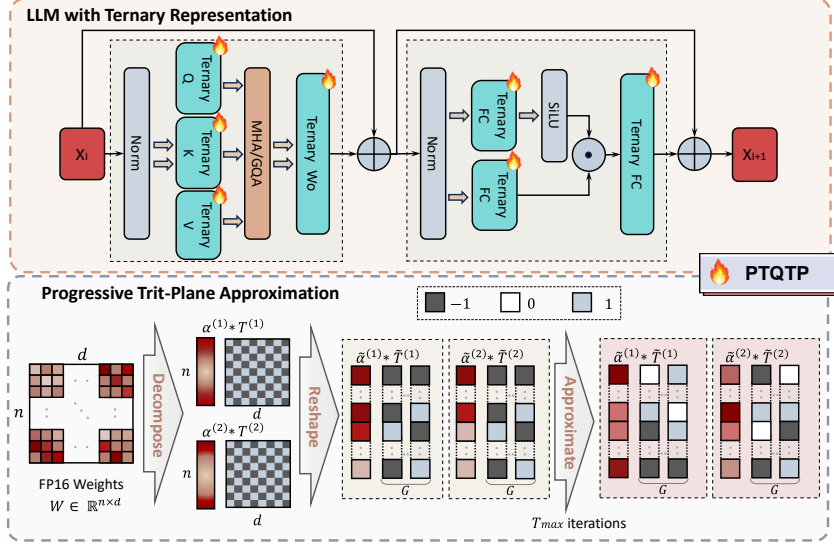
Figure 2: PTQTP workflow overview: (top) Linear layer transformation pathway for ternary quantization in LLaMA architecture; (bottom) Group-wise progressive trit-plane approximation process, where $G$ represents group size and $T_{max}$ indicates maximum iteration count.

that $W$ carries two distinct types of information: (1) *Magnitude:* The amplitude or energy of these connections, which varies across channels. (2) *Topology:* The structural connectivity indicating which input features positively or negatively contribute to the output.

Existing binary quantization $W \in \{-1, +1\}$ suffers from the *forced activation* problem, where near-zero weights, *i.e., noise*, are amplified to $\pm 1$, destroying the delicate logical inference paths essential for **complex reasoning tasks like *MATH* and *Coding***. PTQTP addresses this by decomposing $W$ into a linear superposition of dual sparse trit-planes:

$$W \approx \hat{W} = \sum_{k=1}^{K=2} \underbrace{\text{diag}(\alpha^{(k)})}_{Magnitude} \cdot \underbrace{T^{(k)}}_{Topology} \quad (1)$$

$$W \approx \underbrace{\alpha^{(1)} T^{(1)}}_{Coarse\ Structure} + \underbrace{\alpha^{(2)} T^{(2)}}_{Fine\text{-}grained\ Correction} \quad (2)$$

Where $T^{(k)} \in \{-1, 0, 1\}^{n \times d}$ represents the discrete routing topology, and $\alpha^{(k)} \in \mathbb{R}^n$ represents the continuous channel gain. Crucially, the inclusion of the "0" state allows PTQTP to perform implicit *denoising*, filtering out irrelevant features rather than forcing them into binary states. The superposition of dual trit-planes enables *constructive interference*, where the second plane $T^{(2)}$ acts as a spectral residual compensator, capturing the high-frequency details missed by the principal structural plane $T^{(1)}$ while ignore the parts is reconstructed.
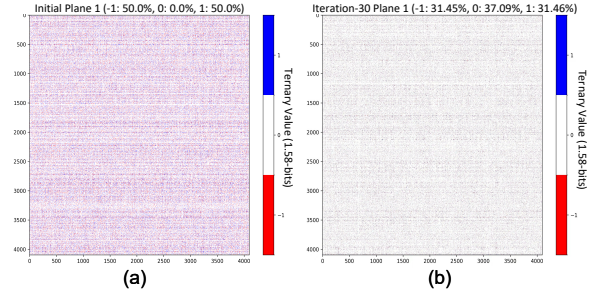


Figure 3: Iterations on LLaMA3.1-8B. (a) Initialized single trit-plane. (b)After 30 iterations.

## 2.2 Dual Trit-planes Approximation

The general process of PTQTP is illustrated in Fig.2. We use $T_i^{(k)}$ to denote the $i$th row of the $k$th trit-plane. First, we fix $T_i^{(k)}$ and optimize the scaling coefficients $\alpha_i^{(k)}$ without any bias terms. To solve for optimal scaling coefficients $\alpha_i^{(k)}$, we use the adaptive ridge regression, i.e., linear regression with a regularization term to the least-squares function. We first define the local basis matrix $S_i$:

$$S_i = \left[ (T_i^{(1)})^T \ (T_i^{(2)})^T \right] \in \{-1, 0, 1\}^{d \times 2} \quad (3)$$

$$A_i = (S_i)^T S_i + \lambda_i I_2 \in \mathbb{R}^{2 \times 2} \quad (4)$$

$$b_i = (S_i)^T W_i^T \in \mathbb{R}^2 \quad (5)$$

Here $\lambda_i \in \mathbb{R}$ is a regularization parameter to improve numerical stability, $I_2$ is the $2 \times 2$ identity matrix. By constructing the local basis matrix row by row, we can find a closed-form solution

3

$\alpha_i(\in \mathbb{R}^2) := [\alpha_i^{(1)} \alpha_i^{(2)}]^T = A_i^{-1} b_i$ for each $i$ independently. However, due to the discrete nature of trit-planes, the method may converge to different minima depending on the regularization parameter $\lambda_i$, which is crucial for the regression performance: if it is too small, the regularization effect may be negligible and the solution may be unstable. If it is too large, the coefficients $\alpha_i$ can be overly diminished, leading to a poor approximation of the original weight matrix.

### 2.3 Progressive Optimization

**Adaptive Regularization.** To achieve a robust decomposes process, we adaptively regularize the approximation and further optimize it through progressive optimization. First, we estimate the condition number of the $2 \times 2$ system:

$$\kappa_{i,approx} = \|A_i\|_F \cdot \|A_i^{-1}\|_F \qquad (6)$$

This measure guides dynamic updates to $\lambda_i$ with the constraint $\lambda_{i,new} = \lambda_i \leq \lambda_{\max}(:= 1.0)$ when $\kappa_{i,approx} < 10^6$. This adaptation mitigates under-regularization (leading to singularity and unstable solutions) and over-regularization (causing excessive coefficient shrinkage), ensuring robustness across weight blocks with varying numerical conditions. Furthermore, the optimal scaling coefficients $\alpha_i$ are then found by solving the following optimization problem:

$$\alpha_i = \arg\min_{\theta_i \in \mathbb{R}^2} \|W_i - S_i \theta_i\|_F^2 + \lambda_i \|\theta_i\|_F^2 \qquad (7)$$

Where $\theta_i$ is a crucial variable that represents the intermediate solution to the scaling coefficients in the process of updating $\alpha_i$. In addition, $\lambda_i \|\theta_i\|_F^2$ is the regularization term that penalizes large values of coefficients $\theta_i$. The Frobenius norm term represents the squared error between the linear combination of the $W_i$ and $\hat{W}_i$, where $S_i$ incorporates the current iteration's trit-plane values. Specifically, we perform the optimization steps $\leq T_{max}$. Each update step is designed to not increase the Frobenius norm $\|W - \hat{W}\|_F^2$. By ensuring that this norm monotonically decreases, we can guarantee that the algorithm will converge to a local minimum. Therefore, we refine trit-plane elements through local exhaustive search, updating $T_{i,j}^{(k)}$ to the value in $c_m \in \{-1, 0, 1\}$ that minimizes squared error:

$$T_{i,j}^{(k)} = \arg\min_{c_m^{(k)} \in \{-1,0,1\}} \left( W_{ij} - \sum_{k=1}^2 \alpha_i^{(k)} c_m^{(k)} \right)^2 \qquad (8)$$

---

**Algorithm 1** Approximation Process of PTQTP

**Require:** Weight matrix $W$, group size $G$, max iterations $T_{\max}$, tolerance $\epsilon$
**Ensure:** Quantized weight approximation $\hat{W}$
1: Divide $W$ into groups $\tilde{W}_i \in \mathbb{R}^G$
2: Initialize $\tilde{T}_{i,(0)}^{(k)} \leftarrow \text{sign}(\tilde{W}_i)$, $\tilde{\alpha}_{i,(0)}^{(k)} \leftarrow [1,1]$
3: **for** $t \leftarrow 1$ to $T_{\max}$ **do**
4:     **for** each row $i$ **do** row-wise decomposition
5:         $\tilde{S}_i \leftarrow [(\tilde{T}_{i,(t-1)}^{(1)})^T \quad (\tilde{T}_{i,(t-1)}^{(2)})^T]$
6:         $\tilde{\alpha}_{i,(t)}^{(k)} \leftarrow \text{RidgeRegression}(\tilde{S}_i, \tilde{W}_i, \lambda_i)$
7:     **end for**
8:     **for** each row $i$ **do**
9:         $\tilde{T}_{i,(t)}^{(1)}, \tilde{T}_{i,(t)}^{(2)} \leftarrow \arg\min_{c_m^{(1)}, c_m^{(2)}} \|\tilde{W}_i - \tilde{\alpha}_{i,(t)}^{(1)} c_m^{(1)} - \tilde{\alpha}_{i,(t)}^{(2)} c_m^{(2)}\|_F^2$
10:     **end for**
11:     **if** $\max_i \|\alpha_{i,(t)}^{(k)} - \alpha_{i,(t-1)}^{(k)}\|_F < \epsilon$ **then break**
12:     **end if**
13: **end for**

---

Therefore, PTQTP achieves $\mathbb{E}[W_{ij}] \approx \alpha_i^T \cdot \mathbb{E}\left[[T_{ij}^{(1)} \; T_{ij}^{(2)}]^T\right]$ and is *bias-free* and *mask-free*, its uniform model architecture preserves hardware-friendly design and operations. Fig. 3 illustrates the example process of a single trit-plane update.

**Batch Processing.** We further introduce group-wise (aka block-wise) processing (Frantar et al., 2022; Lin et al., 2024) in PTQTP, namely, by reshaping $W$ from $n \times d$ to $\frac{nd}{G} \times G$, or equivalently grouping into $G$ columns. Specifically, we set $G = 128$ similar to other related works, and we denote such group-wise operation with the tilde ($\tilde{\circ}$) notation, then we have

$$\tilde{A}_i = \tilde{S}_i^T \tilde{S}_i + \lambda_i I_2, \qquad (9)$$

$$\tilde{\alpha}_i = \tilde{A}_i^{-1} \tilde{b}_i, \quad \tilde{b}_i = \tilde{S}_i^T \tilde{W}_i^T \qquad (10)$$

Such grouping of columns generally reshapes $W$ into a taller $\tilde{W}$, alongside with lengthened $\tilde{\alpha}^{(k)}$'s. Nonetheless, the increase in the latter incurs negligible parameter overhead compared to the model size, which is far outweighed by an improved approximation and performance.

As shown in **Algorithm 1**, we first decompose all the FP16 linear projection weights into trit-planes for magnitude-topology decoupling and approximation. Then we apply progressive optimization and adaptive regulation to automatically save optimal parameters. Furthermore, the unimportant information is filtered out, while the salient feature information is retained. *Experiment results*

Table 1: Perplexity (↓) comparison across SOTA model families on WikiText2 and C4 datasets (*group size=128*). Fine-Tuning means the further optimization step in the quantization method; Salience detect means there are extra search and protect salience weight step in the method. L2: LLaMA2; L3: LLaMA3; Q3: Qwen3; N/A: Model size variant not available; D†: Dual; **Bold**: best result; underlined: second-best.

| Method | Fine Tuning | Salience Detect | WikiText2 | | | | | | | | | C4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | L2-7 | L3-1 | L3-3 | L3-8 | L3-70 | Q3-1.7 | Q3-4 | Q3-8 | Q3-32 | L2-7 | L3-8 | L3-70 |
| FP16 | N/A | N/A | 5.47 | 9.75 | 7.80 | 6.23 | 2.81 | 16.70 | 13.64 | 9.71 | 8.64 | 7.26 | 9.54 | 7.17 |
| QTIP-b4 | × | × | **5.69** | N/A | N/A | **5.69** | **2.75** | **8.46** | **7.09** | **6.28** | N/A | **6.63** | **7.22** | **5.83** |
| GPTQ-b3 | × | × | 6.44 | <u>18.99</u> | <u>16.34</u> | 9.12 | 8.14 | 6.17e4 | 1.34e5 | 5.02e5 | <u>37.10</u> | 7.95 | 17.68 | 10.04 |
| AWQ-b2 | × | × | 2.22e5 | 1.64e5 | 1.11e5 | 7.98e5 | 4.52e4 | 7.52e6 | 1.38e7 | 1.21e7 | N/A | N/A | N/A | N/A |
| GPTQ-b2 | × | × | 52.22 | 4.97e3 | 2.42e3 | 717.54 | 18.79 | 655.10 | 5.92e5 | 7.77e4 | 38.40 | 35.27 | 394.74 | 122.55 |
| OmniQuant-b2 | ✓ | × | 11.06 | 771.05 | 538.64 | 2.08e3 | N/A | N/A | 5.30e4 | 7.31e4 | N/A | 15.02 | N/A | N/A |
| QUIP#-b2 | ✓ | ✓ | 39.73 | N/A | N/A | 84.97 | N/A | N/A | N/A | N/A | N/A | 31.94 | 130.00 | N/A |
| SliM-LLM-b2 | ✓ | ✓ | 15.84 | 3.31e2 | 8.19e2 | 31.52 | N/A | 4.41e4 | 39.71 | N/A | N/A | 84.92 | 390.02 | N/A |
| PBLLM-b1.7 | ✓ | ✓ | 66.41 | 3.87e3 | 1.01e4 | 1.89e3 | N/A | 9.27e5 | 4.68e3 | 1.80e3 | N/A | 80.69 | 104.15 | N/A |
| PTQ1.61-b1.61 | ✓ | ✓ | 12.70 | N/A | N/A | 22.90 | N/A | N/A | N/A | N/A | N/A | N/A | 33.82 | N/A |
| **PTQTP-b1.58-D†** | × | × | <u>6.30</u> | **17.15** | **10.24** | <u>8.53</u> | <u>7.76</u> | <u>32.46</u> | <u>18.25</u> | <u>11.80</u> | **10.06** | <u>6.76</u> | <u>13.24</u> | <u>12.64</u> |

*demonstrate that PTQTP always converges within 50 iterations*, enabling a stable and efficient compression method.

## 3 Experiments

**Quantization Configuration.** We implemented PTQTP on PyTorch platform using models from Huggingface (Paszke et al., 2019). All the weights in linear projection were quantized with tolerance $\epsilon = 10^{-4}$ and maximum iterations $T_{max} = 50$. For numerical stability, we employed dynamic regularization with $\lambda \in [10^{-8}, 1]$. No task-specific calibration, tuning, or fine-tuning was applied in any experiment. All evaluations were conducted on a single NVIDIA A100 80GB GPU.

**Model Architectures.** We evaluated PTQTP across multiple mainstream LLM families including Qwen3 (QwenTeam, 2025), LLaMA3.x (Grattafiori et al., 2024) and LLaMA series (Touvron et al., 2023a,b). To assess cross-domain generalization capabilities, we also tested instruction-tuned variants of these models.

**Baseline Methods.** We compared PTQTP with three categories of methods: (1) Extremely low-bit PTQ methods including: PBLLM (Shang et al., 2023), SliM-LLM (Huang et al., 2024a), QUIP# (Tseng et al., 2024a), QUIP (Chee et al., 2023), AQLM (Egiazarian et al., 2024a); (2) Popular PTQ Methods including: GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2024), QTIP (Tseng et al., 2024b), OmniQuant (Shao et al., 2023); and (3) 1.58-bit QAT approaches (Ma et al., 2024, 2025)

to demonstrate PTQTP's effectiveness even against quantization-aware training-based methods. For fair comparison with smaller models, we included common 1B-3B LLMs such as SmolLM2 (Allal et al., 2025) and MiniCPM (Hu et al., 2024).

**Evaluation Protocol.** We assessed language modeling capability through perplexity measurements on WikiText-2 (Merity et al., 2016) and C4 (Raffel et al., 2020). To evaluate reasoning abilities, we utilized ARC-Challenge, ARC-Easy (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020), Winogrande (Sakaguchi et al., 2021), and MMLU (Hendrycks et al., 2021). For comprehensive comparison with SOTA methods, we further evaluated coding and mathematical reasoning performance on LLaMA3.x and Qwen3 models. All evaluations were conducted using the standard lm-eval benchmarking tool (Gao et al., 2024).

### 3.1 Main Results

**Perplexity on Mainstream Model Backbones.** Table 1 compares the perplexity of PTQTP on WikiText2 and C4 datasets with baselines across LLaMA2 , LLaMA3.x and Qwen 3 variants (0.6B to 70B). The results demonstrate that PTQTP consistently outperforms existing extremely low-bit (1-3 bit) quantization schemes and approaches or exceeds 4-bit methods across diverse architectures. ***This robustness is particularly pronounced in SOTA small LLMs (0.6B-3B), which are typically more vulnerable to quantization due to their higher information density from advanced pre-***

Table 2: Zero-shot Reasoning Tasks Results on LLaMA2-7B and Qwen3-14B models (accuracy %). D[†]: Dual; **Bold**: best result; underlined: second-best. OBQA: OpenBookQA.

| Method | ARC-C | ARC-E | BoolQ | HellaSwag | OBQA | PIQA | WinoG | Average |
|---|---|---|---|---|---|---|---|---|
| | | | | LLaMA2-7B | | | | |
| FP16 | 43.09 | 75.46 | 79.36 | 57.11 | 33.2 | 78.13 | 69.38 | 62.25 |
| GPTQ-b4 | **42.41** | **75.8** | **77.58** | **56.45** | **32.8** | **78.67** | **70.32** | **62.01** |
| GPTQ-b3 | 39.68 | 72.56 | 73.55 | 53.97 | 28.8 | 76.06 | 66.77 | 58.77 |
| GPTQ-b2 | 20.31 | 29.08 | 50.67 | 27.45 | 12.4 | 54.9 | 51.07 | 35.13 |
| AWQ-b2 | 22.27 | 26.68 | 39.94 | 26.17 | 15.4 | 52.94 | 49.41 | 33.26 |
| OmniQuant-b2 | 20.73 | 38.8 | N/A | 30.11 | N/A | 57.34 | 51.78 | N/A |
| QUIP-b2 | 18.94 | 28.7 | N/A | 27.52 | N/A | 56.53 | 48.78 | N/A |
| QUIP#-b2 | 28.84 | 55.56 | N/A | 42.94 | N/A | 71.38 | 62.43 | N/A |
| AQLM-b2 | 32.76 | 63.68 | N/A | 49.55 | N/A | 74.76 | 65.67 | N/A |
| AQLM*-b2 | 34.9 | 66.5 | N/A | 50.88 | N/A | 74.92 | 62.43 | N/A |
| SliM-LLM-b2 | 23.38 | 47.81 | 59.97 | 33.76 | 17.8 | 63.82 | 56.91 | 43.35 |
| PB-LLM-b1.7 | 19.11 | 28.2 | 41.25 | 27.03 | 13.2 | 54.46 | 49.09 | 33.19 |
| PTQ1.61-b1.61 | 22.27 | 47.18 | N/A | 35.78 | N/A | 63.22 | 52.25 | N/A |
| **PTQTP-b1.58-D[†]** | 42.24 | 74.45 | 74.46 | 54.89 | 31.4 | 77.26 | 68.43 | 60.45 |
| | | | | Qwen3-14B-instruct | | | | |
| FP16 | 55.8 | 83.5 | 86.67 | 61.75 | 35.2 | 80.52 | 74.11 | 68.22 |
| AWQ-b4 | **53.84** | **81.61** | **85.6** | 59.13 | **33.2** | **79.71** | **72.38** | **66.49** |
| GPTQ-b4 | 20.99 | 25.38 | 38.44 | 26.01 | 14.4 | 52.45 | 50.91 | 32.65 |
| AWQ-b3 | 40.19 | 65.36 | 71.13 | 42.5 | 31.8 | 69.75 | 60.85 | 54.51 |
| AWQ-b2 | 21.59 | 25.17 | 42.14 | 25.73 | 16.2 | 53.81 | 52.01 | 33.81 |
| GPTQ-b2 | 22.61 | 24.96 | 38.17 | 25.57 | 13.2 | 52.07 | 50.28 | 32.41 |
| SliM-LLM-b2 | 29.35 | 52.54 | 61.2 | 31.52 | 20.4 | 61.83 | 52.04 | 44.13 |
| PB-LLM-b1.7 | 20.73 | 25.93 | 38.04 | 25.76 | 15 | 54.08 | 47.99 | 32.5 |
| PT2-LLM | 23.63 | 53.03 | 62.17 | 33.65 | 20.6 | 62.95 | 59.75 | 45.11 |
| **PTQTP-b1.58-D[†]** | 53.5 | 81.44 | 75.5 | **60.27** | 32.2 | 78.73 | 71.35 | 64.71 |

**training recipes.** The exceptional performance retention of PTQTP establishes it as a robust, generalizable, and efficient quantization solution for both current and future models.

**Language Reasoning Tasks.** Models with larger parameters generally exhibit greater resistance to quantization-induced performance loss, making them ideal for evaluating extreme quantization effectiveness. Results in Table 2 reveal dramatic disparities in capability retention across PTQ methods. Existing solutions exhibit significant performance degradation when performing extremely low bit quantization (less than 2 bits), with a performance gap of almost half. In contrast, PTQTP demonstrates its consistent performance across diverse benchmarks (*average performance accuracy retention near to 95%, and the language reasoning ability is very close to 4-bit methods though PTQTP only use ternary representation precision*) without requiring dynamic bit allocation, salient weight protection, or sensitivity-aware quantization fundamentally challenges the trade-off between quantization and reasoning ability.

**Complex Reasoning *vs* FP16/1.58-bit QAT Language Models.** PTQTP's versatility enables extremely low-bit quantization of advanced models while maintaining near-baseline performance. Unlike QAT methods that require extensive pretraining and fine-tuning, PTQTP applies uniform treatment across model layers without post-quantization adjustments. Table 3 demonstrates PTQTP's minimal performance degradation compared to baselines, and its ability to match or exceed 1.58-bit QAT BitNet schemes (Ma et al., 2024, 2025) at similar model sizes. *These results confirm PTQTP's exceptional stability and establish it as a true plug-and-play solution for model-agnostic 1.58-bit quantization.* Furthermore, compared to other existing PTQ schemes in extreme cases (2 bits), we found that they also perform poorly in terms of mathematical capabilities, especially without using an additional codebook, such as AWQ, which almost completely loses its mathematical and coding capabilities. On the other hand, PTQTP can better protects overall performance including complex reasoning, even compared to the AQLM (with additional codebook).

6

Table 3: Complex tasks comparison between PTQTP-quantized models and leading instruction-tuned LLMs (1B-4B parameters)/ SOTA PTQ-quantized models across efficiency metrics and benchmark performance (accuracy %). D[†] : Dual. **Bold**: best result; <u>underlined</u>: second-best. * Results claimed by BitNet-b1.58 paper.

| Model (Params) | Math-500 | GSM8K | HumanEval | MBPP | MMLU | Storage/GB |
|---|---|---|---|---|---|---|
| LLaMA3.2 (1B) | 14.40 | 46.47 | 41.46 | 52.14 | 46.81 | 2.74 |
| Qwen3 (1.7B) | 72.20 | 74.60 | 65.24 | 63.42 | 60.63 | 4.04 |
| SmolLM2 (1.7B) | 21.00 | 50.19 | 35.98 | 49.81 | 49.73 | 3.42 |
| MiniCPM (2B) | 0.60 | 56.48 | 39.02 | 55.64 | 52.08 | 5.45 |
| BitNet-b1.58 (2B)* | 43.40 | 58.38 | N/A | N/A | 53.17 | 1.18 |
| LLaMA3.2 (3B) | 45.20 | 77.56 | 62.20 | 63.81 | 62.24 | 5.98 |
| Qwen3 (4B) | **83.00** | <u>86.96</u> | **77.44** | **79.38** | **70.67** | 7.49 |
| Qwen3-PTQTP-b1.58-D[†] (1.7B) | 43.80 | 54.21 | 45.53 | 48.78 | 43.82 | 1.90 |
| Qwen2.5-PTQTP-b1.58-D[†] (3B) | 42.80 | 62.47 | 25.61 | 44.75 | 56.80 | 2.61 |
| LLaMA3.2-PTQTP-b1.58-D[†] (3B) | 34.20 | 65.81 | 17.68 | 52.92 | 55.08 | 2.94 |
| Qwen3-PTQTP-b1.58-D[†] (4B) | <u>82.60</u> | **87.79** | <u>71.34</u> | <u>69.26</u> | <u>63.65</u> | 3.35 |
| LLaMA3 (8B) | 27.80 | 79.23 | 59.15 | 68.48 | 68.28 | 14.96 |
| LLaMA3-AWQ-b2 (8B) | 0.00 | 0.15 | 0.00 | 0.00 | 24.82 | 2.86 |
| LLaMA3-AQLM-b2.07 (8B) | 12.40 | 60.42 | 28.66 | 29.27 | 47.86 | 4.08 |
| LLaMA3-PTQTP-b1.58-D[†] (8B) | 15.40 | 67.02 | 42.07 | 45.91 | 54.54 | 5.61 |

Table 4: Perplexity (PPL) of LLaMA-3.1 8B and LLaMA-3.2 3B on WikiText-2, PTB, and C4 under different regularisation strengths. Lower is better.

| Condition | LLaMA-3.1 8B | | | LLaMA-3.2 3B | | |
|---|---|---|---|---|---|---|
| | WikiText2 | PTB | C4 | WikiText2 | PTB | C4 |
| 1 | 12.45 | 17.52 | 17.76 | 12.52 | 21.20 | 18.09 |
| 5 | 10.86 | 16.11 | 16.17 | 11.52 | 19.98 | 17.13 |
| $10^1$ | 9.58 | 14.92 | 14.65 | 10.92 | 18.64 | 16.18 |
| $10^2$ | 8.97 | 14.20 | 13.81 | 10.52 | 17.76 | 15.56 |
| $10^4$ | 8.62 | 13.84 | 13.40 | 10.26 | 17.32 | 15.26 |
| $10^6$ | 8.53 | 13.72 | 13.22 | 10.24 | 17.41 | 15.28 |
| $10^{12}$ | 8.53 | 13.72 | 13.22 | 10.24 | 17.41 | 15.28 |

## 3.2 Ablation Studies

**Condition Numbers.** Table 4 summarizes the validation perplexities obtained when sweeping a regularization coefficient (in log-space) across two model scales: LLaMA-3.1 8B and LLaMA-3.2 3B. All PPLs are reported on WikiText-2, PTB, and the C4 validation split. Across corpora, both models exhibit monotonically decreasing perplexity (illustrates robustness) as the condition numbers increases from $10^0$ to $10^2$, after which the metric saturates—indicating a regime where further rise the condition bounds offers negligible improvement.

**Progressive Search Iterations.** Fig. 4(a)-(c) identify a critical 30-iteration threshold that optimally balances perplexity and efficiency across model scales. Within this range, both test models

achieve remarkable convergence. Smaller models converge faster due to simpler weight structures, while larger models leverage their expressivity for steeper initial gains. Beyond 30 iterations, perplexity stabilizes while quantization time increases linearly, indicating diminishing returns from additional refinement.

**Bound of Tolerance.** Fig. 4(d)-(f) reveal the crucial trade-off between tolerance $\epsilon$ and quantization performance. Tighter $\epsilon$ values improve perplexity but at significant computational cost: LLaMA3.1-8B achieves 9.1% perplexity reduction with a 172% runtime increase, while LLaMA3.2-3B shows 4.5% improvement with 137% longer runtime. Smaller models show better stability at low $\epsilon$, while larger models require stricter tolerance to capture complex weight relationships. Our experiments identify $\epsilon \in [10^{-3}, 10^{-4}]$ as the optimal range for balancing accuracy and efficiency, particularly important for resource-constrained deployment scenarios.

## 4 Inference Efficiency

**Hardware Efficiency.** The PTQTP framework leverages ternary operations $\{-1, 0, 1\}$ to achieve hardware-efficient computations. For a ternary element $c_m \in \{-1, 0, 1\}$ and scaling coefficient $\alpha$, the multiplication operation with $c_m$ in dequantization or dequantization-free forward can be implemented as an addition, eliminating floating-point
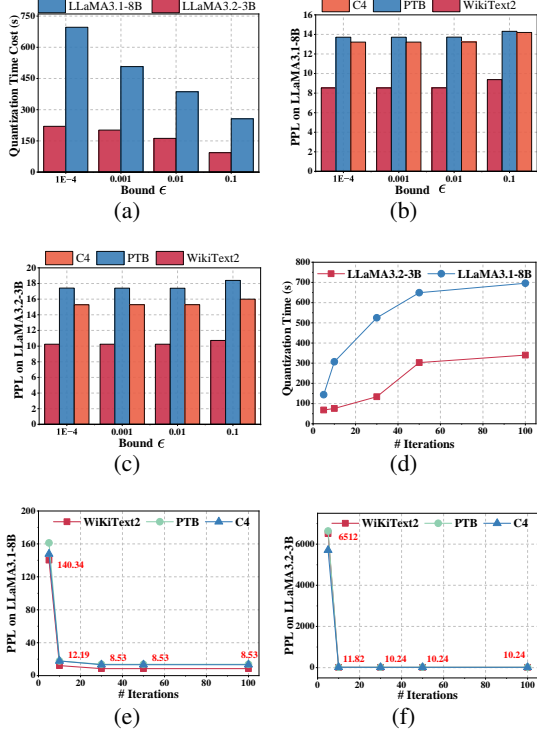
Figure 4: (a)-(c). Effect of progressive search iterations on quantization time (left sub-figure) and perplexity (PPL) (middle and right sub-figures) for LLaMA3.1-8B and LLaMA3.2-3B models. (d)-(f). Trade-off between tolerance bounds ($\epsilon$), quantization time (left sub-figure), and model perplexity (PPL) (middle and right sub-figures) for LLaMA3.1-8B and LLaMA3.2-3B architectures.

multipliers and replacing them with sign flips, identity mappings, or zeroing operations. This reduces the arithmetic intensity to $\mathcal{O}(1)$ per element in the matrices. The far lower multiplication and the inherent sparsity in PTQTP than other low-bit methods and the decomposition into two trit-planes $T^{(1)}, T^{(2)}$, allowing parallelized additive superposition of dual trit-planes, where each plane can be processed independently on parallelized architectures, exploiting data-level parallelism and minimizing latency in inference acceleration. This bandwidth reduction is critical for latency-bound applications, as memory access often dominates inference time in modern hardware architectures, and the final end-to-end experiments also verify the efficiency of PTQTP.

**End-to-End Speed.** To fully exploit the hardware-friendly properties of ternary weight quantization in PTQTP, we design a ternary matrix multiplication CUDA kernel based on LUT-GEMM (Park et al., 2022). Since the weight values are restricted to {-1,0,1}, all possible

Table 5: The end-to-end generation speed (tokens/s) comparison. *: b1.58-Dual.

|        | FP16 | AQLM-b2 | QuIP#-b2 | PTQTP* |
|--------|------|---------|----------|--------|
| L2-7B  | 32.2 | 65.7    | 99.04    | 149.35 |
| L2-13B | OOM  | 26.3    | 61.32    | 93.98  |

dot products between activation vectors and low-bit weight patterns can be precomputed and stored in lookup tables (LUTs) in the shared memory. Notably, these dot products involve only addition operations. As a result, subsequent matrix multiplications can be performed via direct LUT lookups without any dequantization, and the computational overhead is limited to a one-time additive cost during LUT construction and shared memory read.

In our experiments, we evaluate the end-to-end generation throughput of PTQTP and other baseline methods (AQLM (Egiazarian et al., 2024b) and QuIP#) on a single 24GB RTX 3090 GPU. The input prompt length is set to 512 tokens, the generated sequence length is 512 tokens, and the batch size is 1. The results are summarized in Table 5. For LLaMA-2 7B inference, PTQTP achieves a 4.63× speedup compared to the FP16 baseline. For LLaMA-2 13B, PTQTP delivers 3.57× and 1.53× speedups over AQLM and QuIP#, respectively.

## 5 Conclusion

We introduced PTQTP, a structured PTQ framework that achieves ternary quantization through systematic trit-plane decomposition without retraining or fine-tuning. Extensive experiments demonstrate that PTQTP outperforms existing low-bit PTQ methods, and approaches or even surpasses the accuracy of 4-bit PTQ and 1.58-bit QAT techniques. Remarkably, PTQTP preserves mathematical reasoning and coding capabilities that catastrophically degrade or collapse in other extremely low-bit PTQ schemes, challenging the conventional wisdom that such tasks inherently require higher precision. In addition, PTQTP's model-agnostic design ensures robust deployment across diverse architectures from 1B to 70B models while providing a impressive end-to-end speed gains. By addressing the core challenges of expressiveness, stability, scalability, and hardware efficiency, PTQTP sets a new yardstick for compressing high-performance LLMs for resource-constrained platforms.

## Limitations

PTQTP's memory layout can be further optimized through bit-packing, where 8 ternary elements are stored in a single byte, improving cache locality and reducing cache misses by 20%-30% compared to non-packed formats. In addition, PTQTP's multiplication-free operations align perfectly with hardware accelerators designed for binary/ternary neural networks, such as customized ASICs. These architectures can execute sign flips and additions in single clock cycles, achieving peak throughput with minimal energy consumption. For instance, ternary operations can leverage bitwise XOR for sign inversion, enabling sub-nanosecond latency per operation on specialized hardware. However, currently, technical support is still not widespread, and PTQTP's storage and technical computing still rely on modern computing devices, which limits the performance of PTQTP to some extent. Future works can focus on software-hardware co-design and hardware framework design for 1.58-bit inference acceleration.

## LLM Usage Disclosure

This paper is primarily the work of the human authors, but we also made use of several advanced LLMs, including ChatGPT-5 and Gemini 3. They were employed to assist with code development and debugging, support result analysis, and help with formatting and language polishing. We acknowledge the contributions of these LLMs, while fully recognizing their limitations, and take full responsibility for all content presented under our names. We did not include hidden prompt-injection text in the submission, and all external data and code comply with their respective licenses.

## References

Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, and 1 others. 2025. Smollm2: When smol goes big–data-centric training of a small language model. *arXiv preprint arXiv:2502.02737*.

Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Pashmina Cameron, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. 2024. Quarot: Outlier-free 4-bit inference in rotated llms. *Advances in Neural Information Processing Systems*, 37:100213–100240.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.

Tom Brown and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. 2023. Quip: 2-bit quantization of large language models with guarantees. *Advances in Neural Information Processing Systems*, 36:4396–4429.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.

Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. 2024a. Extreme compression of large language models via additive quantization. *arXiv preprint arXiv:2401.06118*.

Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. 2024b. Extreme compression of large language models via additive quantization. *ArXiv*, abs/2401.06118.

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. 2022. GPTQ: Accurate post-training compression for generative pretrained transformers. *arXiv preprint arXiv:2210.17323*.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. The language model evaluation harness.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and Aiesha Letman. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.

Seongmin Hong, Seungjae Moon, Junsoo Kim, Sungjae Lee, Minsub Kim, Dongsoo Lee, and Joo-Young Kim. 2023. Dfx: A low-latency multi-fpga appliance for accelerating transformer-based text generation. In *Proceedings of the 55th Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO '22, page 616–630.

Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, and 1 others. 2024. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*.

Wei Huang, Yangdong Liu, Haotong Qin, Ying Li, Shiming Zhang, Xianglong Liu, Michele Magno, and XIAOJUAN QI. 2024a. BiLLM: Pushing the limit of post-training quantization for LLMs. In *Forty-first International Conference on Machine Learning*.

Wei Huang, Haotong Qin, Yangdong Liu, Yawei Li, Qinshuo Liu, Xianglong Liu, Luca Benini, Michele Magno, Shiming Zhang, and Xiaojuan Qi. 2024b. Slim-llm: Salience-driven mixed-precision quantization for large language models. *arXiv preprint arXiv:2405.14917*.

Zhiteng Li, Xianglong Yan, Tianao Zhang, Haotong Qin, Dong Xie, Jiang Tian, zhongchao shi, Linghe Kong, Yulun Zhang, and Xiaokang Yang. 2025. ARB-LLM: Alternating refined binarizations for large language models. In *The Thirteenth International Conference on Learning Representations*.

Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. In *MLSys*.

Shuming Ma, Hongyu Wang, Shaohan Huang, Xingxing Zhang, Ying Hu, Ting Song, Yan Xia, and Furu Wei. 2025. Bitnet b1.58 2b4t technical report. *Preprint*, arXiv:2504.12285.

Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. 2024. The era of 1-bit llms: All large language models are in 1.58 bits. *Preprint*, arXiv:2402.17764.

Mitch Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. 1994. The penn treebank: Annotating predicate argument structure. In *Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*.

Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *Preprint*, arXiv:1609.07843.

Gunho Park, Baeseong Park, Minsub Kim, Sungjae Lee, Jeonghoon Kim, Beomseok Kwon, Se Jung Kwon, Byeongwook Kim, Youngjoo Lee, and Dongsoo Lee. 2022. Lut-gemm: Quantized matrix multiplication based on luts for efficient inference in large-scale generative language models. In *International Conference on Learning Representations*.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. Pytorch: An imperative style, high-performance deep learning library. *Preprint*, arXiv:1912.01703.

David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluis-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. 2021. Carbon emissions and large neural network training. *arXiv preprint arXiv:2104.10350*.

QwenTeam. 2025. Qwen3.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.

Yuzhang Shang, Zhihang Yuan, Qiang Wu, and Zhen Dong. 2023. Pb-llm: Partially binarized large language models. *Preprint*, arXiv:2310.00034.

Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. 2023. Omniquant: Omnidirectionally calibrated quantization for large language models. *arXiv preprint arXiv:2308.13137*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, and Nikolay Bashlykov. 2023b. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

10

Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. 2024a. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks. *arXiv preprint arXiv:2402.04396*.

Albert Tseng, Qingyao Sun, David Hou, and Christopher M De Sa. 2024b. Qtip: Quantization with trellises and incoherence processing. *Advances in Neural Information Processing Systems*, 37:59597–59620.

Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. 2023. Bitnet: Scaling 1-bit transformers for large language models. *Preprint*, arXiv:2310.11453.

Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International conference on machine learning*, pages 38087–38099. PMLR.

He Xiao, Qingyao Yang, Dirui Xie, Wendong Xu, Wenyong Zhou, Haobo Liu, Zhengwu Liu, and Ngai Wong. 2025. Exploring layer-wise information effectiveness for post-training quantization in small language models. *arXiv preprint arXiv:2508.03332*.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.

Shulin Zeng, Jun Liu, Guohao Dai, Xinhao Yang, Tianyu Fu, Hongyi Wang, Wenheng Ma, Hanbo Sun, Shiyao Li, Zixiao Huang, Yadong Dai, Jintao Li, Zehao Wang, Ruoyu Zhang, Kairui Wen, Xuefei Ning, and Yu Wang. 2024. Flightllm: Efficient large language model inference with a complete mapping flow on fpgas. *Preprint*, arXiv:2401.03868.

Jiaqi Zhao, Miao Zhang, Ming Wang, Yuzhang Shang, Kaihao Zhang, Weili Guan, Yaowei Wang, and Min Zhang. 2025. Ptq1. 61: Push the real limit of extremely low-bit post-training quantization methods for large language models. *arXiv preprint arXiv:2502.13179*.

## A  Related Works

**Recent Progress in PTQ for LLMs.**  PTQ has emerged as a promising approach for compressing LLMs without the computational overhead of retraining. GPTQ (Frantar et al., 2022) pioneered efficient 4-bit quantization through layer-wise optimization. AWQ (Lin et al., 2024) further improved this by incorporating activation-aware scaling, demonstrating that the consideration of activations during quantization leads to better performance. Recent advancements like QuIP# (Tseng et al., 2024a), AQLM (Egiazarian et al., 2024a),

and (Huang et al., 2024b; Chee et al., 2023; Zhao et al., 2025) have pushed the boundaries of PTQ by utilizing vector quantization or learnable equivalent transformations to handle outliers. While these methods achieve impressive accuracy, they often introduce non-trivial decoding overhead (e.g., codebook lookups) or require complex calibration optimization. LieQ (Xiao et al., 2025) attempts to do the structured and explainable quantization using layer-wise mixed precision scheme,achieving extremely low-bit PTQ with well-reasoned and hardware-friendly properties. In addition, other works have pushed towards lower bit width PTQ. PBLLM (Shang et al., 2023) and (Huang et al., 2024a; Li et al., 2025) achieved 1-bit quantization by identifying and preserving structurally salient weights while applying different quantization strategies to different weight groups. However, these methods rely on complex, unstructured weight categorizations that complicate hardware implementation.

**Ternary Quantization.**  The transition from binary to ternary quantization represents a leap of model expressiveness while still preserving multiplier-less arithmetic. BitNet's (Wang et al., 2023) extension to 1.58-bit (ternary) weights (Ma et al., 2024) demonstrated superior performance compared to binary models, suggesting the value of the additional state in the ternary representation. However. existing ternary quantization methods mainly (if not all) rely on QAT approaches (such as BitNet), requiring extensive pretraining and substantial computational resources. In addition, previous methods use unstructured weight quantization to compensate for the degenerate representation ability. Unlike these works that optimize progressive reconstruction via retrained context models, PTQTP is a post-training, structured and model-agnostic framework that decomposes LLM weights into uniform and collaborative tritplanes with adaptive ridge-scaled coefficients to yield 1.58-bit quantization without retraining or fine-tuning.

**Hardware Efficiency.**  A key advantage of binary and ternary quantization lies in their suitability for FPGA or ASIC deployment (e.g., FlightLLM (Zeng et al., 2024), DFX (Hong et al., 2023)), where multiplier-less feature–weight products can be efficiently implemented in hardware. Specifically, binary operations can be realized with XNOR gates and bit-counting, while ternary opera-
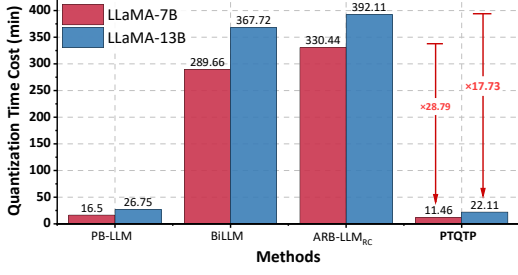
Figure 5: Runtime performance evaluation. PTQTP outperforms existing extremely low-bit quantization methods while achieving $17.73\times - 28.79\times$ speedup over ARB-LLM$_{RC}$.

tions leverage simple adders, preserving computational efficiency while offering greater expressivity. Therefore, PTQTP focuses on the specific niche of *multiplication-free* inference via structured ternary decomposition, prioritizing hardware efficiency on edge devices where adders are significantly cheaper than multipliers.

## B More Perspective on Efficiency Analysis

**Computational Complexity.** For row $i$, the normal equation matrix $A_i = (S_i)^T S_i + \lambda_i I_2$ in PTQTP is $2 \times 2$, with inverse computed in constant time $\mathcal{O}(1)$ using the formula:

$$A_i^{-1} = \frac{1}{\det(A_i)} \begin{bmatrix} A_i(2,2) & -A_i(1,2) \\ -A_i(2,1) & A_i(1,1) \end{bmatrix} \quad (11)$$

Where $\det(A_i) = A_i(1,1)A_i(2,2) - A_i(1,2)^2$. The vector $b_i = S_i^T W_i^T$ requires $2d$ operations (dot products for each column of $S_i$), leading to $\mathcal{O}(d)$ per row. Adding up to $n$ rows, this step is $\mathcal{O}(nd)$ per iteration. In addition, for each element $(i,j)$ in $W$, evaluating $(T^{(1)}, T^{(2)}) \in \{-1, 0, 1\}^2$ involves 9 arithmetic operations per element, resulting in $\mathcal{O}(1)$ complexity per element. For $d$ elements per row, this is $\mathcal{O}(nd)$ in $n$ rows per iteration. Therefore, with $T_{\max}$ iterations, the total computational complexity is $\mathcal{O}(T_{\max} \cdot nd)$.

**Time Complexity.** For each iteration, the time complexity can be presented as $\mathcal{O}(nd)$. With empirical convergence in $T_{\max} \leq 50$ iterations, the total time complexity is $\mathcal{O}(T_{\max} \cdot nd)$, exhibiting linear scalability with the dimension of weight matrix $n, d$, critical for deployment on LLMs with billions of parameters.

$\mathcal{O}(T_{\max} \cdot nd)$ complexity is optimal for low-bit quantization when considering both approximation

quality and computational efficiency. Compared to existing PTQ methods, *PTQTP achieves lower per-iteration complexity due to its closed-form solutions and fixed-dimensional ridge regression.*

**Run Time Consumption.** The experiment results in Fig. 4 show that PTQTP can stably converges less than 50 iterations with different model structures, resulting in 28.79 $\times$ speedup during quantization on modern GPUs, which are shown in Fig. 5(a). In other words, PTQTP can quickly find the optimal trit-planes to achieve model compression. Compared to existing extremely low-bit schemes, it significantly reduces runtime overhead, which we believe is mainly due to the absence of additional outlier handling and protection operations. Furthermore, PTQTP also boasts the hardware efficiency advantages of binary compression, and its performance is stronger and more stable. More performance comparison are shown in Appendix D.

**Memory Saving.** For $W \in \mathbb{R}^{n \times d}$, group size $k$, the memory of $\hat{W}$ after standard quantization of bits $m$ is

$$\mathcal{M} = \underbrace{n \times d \times m}_{B} + \underbrace{[d/k]}_{\text{multiple groups}} \times \underbrace{n \times 16}_{\text{row-wise FP16 } \alpha} \quad (12)$$

The memory required by BiLLM can be formulated as follows, where $c$ is the number of salient columns for $W$.

$$\mathcal{M}_{\text{BiLLM}} = \underbrace{n \times d}_{\text{group bitmap}} + \underbrace{d}_{\text{salient column bitmap}}$$
$$+ \underbrace{2 \times n \times c + [d/k] \times 3n \times 16}_{\text{second-order binarization}}$$
$$+ \underbrace{n \times (d-c) + \overbrace{[d/k] \times 2n \times 32}^{\text{2 groups}}}_{\text{first-order binarization}} \quad (13)$$

Previous work (Li et al., 2025) derived the total memory occupation ARB-RC and ARB-RC + CGB (grouped column bitmap), which can be formulated

as:

$$
\mathcal{M}_{\text{ARB-RC}} = \underbrace{n \times d}_{\text{group bitmap}} + \underbrace{d}_{\text{salient column bitmap}}
$$

$$
+ \underbrace{2 \times n \times c + ([d/k]] \times 2n + 2c) \times 16}_{\text{second-order binarization}}
$$

$$
+ \underbrace{n \times (d-c) + \overbrace{([d/k] \times n + (d-c)) \times 16 \times 2}^{\text{2 groups}}}_{\text{first-order binarization}}
$$

$$(14)$$

$$
\mathcal{M}_{\text{ARB-RC+CGB}} = \underbrace{n \times d}_{\text{group bitmap}} + \underbrace{d}_{\text{salient column bitmap}}
$$

$$
+ \underbrace{2 \times n \times c + ([d/k] \times 2n + 2c) \times 16 \times 2}_{\text{second-order binarization}}
$$

$$
+ \underbrace{n \times (d-c) + \overbrace{([d/k] \times n + (d-c)) \times 16 \times 2}^{\text{2 groups}}}_{\text{first-order binarization}}
$$

$$(15)$$

In PTQTP, each trit-plane containing 3 states has to be stored as a 2bit datatype due to the hardware constraint. The total memory of PTQTP is

$$
\mathcal{M}_{\text{PTQTP}} = \underbrace{2 \times n \times d \times 2}_{\text{2 Trit-Plane}} + \underbrace{[d/k] \times 2n \times 16}_{\text{row-wise FP16 } \alpha}
$$

$$(16)$$

Figure 5(b) demonstrates the estimated memory demand of PTQTP with other binary quantization methods, derived from the above formulas. The proposed PTQTP slightly increased the memory consumption to other binary quantization approaches. This is a trade-off between storage and representational capacity. However, methods such as BiLLM and ARB-LLM$_{\text{RC}}$ explicitly divide columns into first-order and second-order groups based on their saliency (as shown in Eqs.13, 14, 15), assigning more bit planes and FP16 parameters to the more salient second-order columns. As a result, *although PTQTP uses trit-planes to represent quantized weights, it does not incur significant memory overhead compared to binary-based methods*.

## C Implementation Details About PTQTP

**Progressive and Adaptive Regularization.** Algorithm 2 illustrates more details about the process of PTQTP. In the beginning, the trit-planes $T^{(k)}$ are initialized using the sign function of $W$, with

---

**Algorithm 2** Progressive and Adaptive Regularization

**Require:** Weight matrix $W \in \mathbb{R}^{n \times d}$, max iterations $T_{\max}$, tolerance $\epsilon$
**Ensure:** Optimized parameters $\alpha, T^{(1)}, T^{(2)}$
1: **Initialize:**
2: $T_{(0)}^{(k)} \leftarrow \text{sign}(W)$ with $0 \rightarrow 1$ replacement for $k = 1, 2$
3: Let $\alpha_{(0)} \leftarrow [1, 1] \otimes \mathbf{1}_n$
4: Let $\lambda_{(0)} \leftarrow 10^{-8} \cdot \mathbf{1}_n$
5: **for** $t = 1$ **to** $T_{\max}$ **do**
6:     **for** row $i = 1$ **to** $n$ **do**
7:         $S_i \leftarrow [T_{i,(t-1)}^{(1)T} \quad T_{i,(t-1)}^{(2)T}]$
8:         $A_i \leftarrow (S_i)^T S_i + \lambda_i I_2$
9:         $b_i \leftarrow (S_i)^T W_i^T$
10:         Solve $\alpha_{(i,t)} \leftarrow (A_i)^{-1} b_i$
11:         $\lambda_{i,new} \leftarrow \lambda_i \sqrt{\kappa_{i,approx}/10^6}$,
12:         When $\kappa_{i,approx} \geq 10^6$
13:     **end for**
14:     $\mathcal{C} \leftarrow \{-1, 0, 1\}^2$
15:     **for** element $(i, j) \in W$ **do**
16:         Evaluate error for all $\mathbf{c}_m \in \mathcal{C}$:
17:         $e_m \leftarrow \left( W_{ij} - \sum_{k=1}^{2} \alpha_{i,(t)}^{(k)} c_m^{(k)} \right)^2$
18:         $m^* \leftarrow \arg\min_m e_m$
19:         $T_{ij,(t)}^{(1)}, T_{ij,(t)}^{(2)} \leftarrow \mathbf{c}_m^*$
20:     **end for**
21:     **if** $\|\alpha_{(t)} - \alpha_{(t-1)}\|_F < \epsilon$ **then**
22:         **break**
23:     **end if**
24: **end for**
25: **Return** $\alpha_{(t)}, \{T_{(t)}^{(1)}, T_{(t)}^{(2)}\}$

---

zero entries replaced by 1 to ensure valid ternary values (later expanded to include 0 through optimization); scaling coefficients $\alpha$ are initialized as uniform vectors $[1, 1]$ replicated across all rows; and the regularization parameter $\lambda$ starts at a small value $10^{-8}$ to promote numerical stability.

In each iteration, the algorithm first updates the continuous scaling coefficients $\alpha$ row by row. If the condition number $\kappa_{i,approx} \geq 10^6$, indicating ill-conditioning, the regularization parameter $\lambda_i$ is adaptively increased $\leq \lambda_{max} = 1.0$ to stabilize the solution. This adaptive regularization mitigates the sensitivity to small input perturbations. The scaling coefficients are then solved via closed-form ridge regression, ensuring efficient computation without iterative solvers. Next, the algorithm optimizes the discrete trit-planes. It generates all 9 possible

ternary value combinations $\mathcal{C} = \{-1, 0, 1\}^2$ for the pair $(T^{(1)}, T^{(2)})$, and for each matrix element $(i, j)$, computes the approximation error for each combination. The combination $\mathbf{c}_{m^*}$ that minimizes the squared error $e_m$ is selected to update $T_{ij}^{(1)}$ and $T_{ij}^{(2)}$, effectively conducting a local exhaustive search to find the best ternary representation for each weight entry. Finally, convergence is checked by monitoring the Frobenius norm difference between consecutive scaling coefficient updates. If $\|\alpha_{(t)} - \alpha_{(t-1)}\|_F < \epsilon$, the algorithm terminates early, balancing optimization quality with computational efficiency.

**Regularized Initialization.** We adopt a progressive initialization strategy to kick-start the optimization process. The first step of our initialization focuses on approximating the weight matrix $\tilde{W}_i$ group by group. We start by initializing the trit-planes $\tilde{T}_{i,(0)}^{(1)}$ $\tilde{T}_{i,(0)}^{(2)}$ using the sign function of $\tilde{W}_i$:

$$\tilde{T}_{i,(0)}^{(k)} = \text{sign}(\tilde{W}_i) \qquad (17)$$

After obtaining the group-wise initialized trit-planes $\tilde{T}_{i,(0)}^{(1)}$ and $\tilde{T}_{i,(0)}^{(2)}$, we construct the matrix $\tilde{S}_{i,(0)}$ where each row contains elements from $\tilde{T}_{i,(0)}^{(1)}$ and $\tilde{T}_{i,(0)}^{(2)}$. The initial scaling coefficients $\tilde{\alpha}_{i,(0)}$ are obtained by solving:

$$\tilde{\alpha}_{i,(0)} = \arg\min_{\alpha} \left( \|\tilde{W}_i - \tilde{S}_{i,(0)}\theta\|_F^2 + \lambda\|\theta\|_F^2 \right) \qquad (18)$$

The closed-form solution with dimension-compatible regularization becomes:

$$\tilde{\alpha}_{i,(0)} = \left( \tilde{S}_{i,(0)}^T \tilde{S}_{i,(0)} + \lambda I \right)^{-1} \tilde{S}_{i,(0)}^T \tilde{W}_i^T \qquad (19)$$

The coefficient bound should be modified as follows, where $\sigma_{\max}(\cdot)$ denotes the maximum singular value.

$$\|\tilde{\alpha}_{i,(0)}\|_2 \le \frac{\sigma_{\max}(\tilde{S}_{i,(0)})}{\sigma_{\min}^2(\tilde{S}_{i,(0)}) + \lambda} \|\tilde{W}_i^T\|_2 \qquad (20)$$

## D  More Performance Illustration

**Full experiments results of PTQTP.** In this supplementary chapter, we listed all our test results. Table 6 displays the results of the benchmark experiment with a certain level of difficulty, showing the quantization stability of our PTQTP for the Qwen3-14B quantized model. Table 7 illustrated the specific accuracy comparison about PTQTP on 6 zero-shot reasoning tasks.

Table 6: Accuracy of 7 language reasoning tasks on Qwen3 family with PTQTP. We compare the results among Qwen-3 Models FP16 and PTQTP across model size from 0.6B to 32B. HellaS: HellaSwag. **Bold**: Keep original performance more than 95% after using PTQTP.

| Metrics | Method | Size | | | | | |
|---|---|---|---|---|---|---|---|
| | | 0.6B | 1.7B | 4B | 8B | 14B | 32B |
| ARC-C | FP16 | 52.88 | 65.08 | 79.66 | 82.37 | 82.37 | 88.81 |
| | PTQTP | 44.41 | 51.53 | 76.95 | 68.81 | 80.68 | 77.97 |
| | (%) | 84.0 | 79.2 | **96.6** | 83.5 | **97.9** | 87.8 |
| ARC-E | FP16 | 73.72 | 80.25 | 87.13 | 89.07 | 90.83 | 92.24 |
| | PTQTP | 56.44 | 70.02 | 86.60 | 81.13 | 88.71 | 86.24 |
| | (%) | 76.6 | 87.3 | **99.4** | 91.1 | **97.7** | 93.5 |
| BoolQ | FP16 | 66.42 | 78.87 | 86.36 | 86.39 | 89.45 | 87.98 |
| | PTQTP | 62.57 | 27.55 | 84.13 | 85.44 | 89.02 | 86.27 |
| | (%) | 94.2 | 34.9 | **97.4** | **98.9** | **99.5** | **98.1** |
| HellaS | FP16 | 43.60 | 58.03 | 78.19 | 84.67 | 88.09 | 90.80 |
| | PTQTP | 32.85 | 37.92 | 71.02 | 79.75 | 85.13 | 89.56 |
| | (%) | 75.3 | 65.3 | 90.8 | 94.2 | **96.6** | **98.6** |
| PIQA | FP16 | 58.98 | 67.57 | 77.26 | 80.20 | 81.07 | 83.51 |
| | PTQTP | 57.67 | 57.89 | 73.18 | 73.78 | 79.33 | 80.41 |
| | (%) | **97.8** | 85.7 | **94.7** | 92.0 | **97.9** | **96.3** |
| WinoG | FP16 | 50.51 | 50.99 | 58.96 | 54.14 | 68.90 | 77.82 |
| | PTQTP | 50.43 | 51.07 | 53.67 | 59.67 | 67.72 | 73.64 |
| | (%) | **99.8** | **100.2** | 91.0 | **110.2** | **98.3** | 94.6 |
| MMLU | FP16 | 43.76 | 60.67 | 70.68 | 76.55 | 79.38 | 82.81 |
| | PTQTP | 33.64 | 43.82 | 63.65 | 68.23 | 76.20 | 80.56 |
| | (%) | 76.9 | 72.2 | 90.1 | 89.1 | **96.0** | **97.3** |

**Performance vs Binary PTQ.** We list all of our experiment results here for reference. We further compared PTQTP performance on different datasets with various model size and backbones, the results are illustrated in Table 8 and Table 9. According to these results, we find specialized binary schemes (PBLLM, BiLLM) show catastrophic failure on Math-500 (0%) and GSM8K (<2%). Even ARB-LLM$_{RC}$, despite architectural modifications, achieves only 1.80% and 32.22% on these benchmarks—significantly below baseline. ***This evidence confirms that mathematical reasoning is uniquely vulnerable to precision loss in conventional PTQ approaches***, exhibiting $19\times$ greater performance degradation compared to linguistic tasks, with non-negligible impacts (12.7% mean reduction) on general reasoning benchmarks. These findings suggest that architectural modifications alone—such as salient weight protection—cannot effectively preserve mathematical reasoning capabilities at ultra-low precision. In contrast, PTQTP-b1.58 achieves 82.40% on Math-500 and 85.44% on GSM8K—less than 5% degradation from the baseline. This preservation of mathematical reasoning at ultra-low bit-widths suggests PTQTP effectively decouples numerical precision requirements from model parameterization.

Table 7: Comparison between PTQTP-quantized models and leading instruction-tuned LLMs (1B-4B parameters) across efficiency metrics and benchmark performance (accuracy %). $D^\dagger$ : Dual. **Bold**: best result; underlined: second-best. * Results claimed by BitNet-b1.58 paper.

| Model (Params) | ARC-C | ARC-E | BoolQ | HellaSwag | PIQA | WinoG |
|---|---|---|---|---|---|---|
| LLaMA3.2 (1B) | 55.93 | 69.84 | 62.20 | 40.12 | 58.16 | 50.04 |
| Qwen3 (1.7B) | **82.71** | 80.07 | 78.59 | 57.76 | 67.36 | 50.75 |
| SmolLM2 (1.7B) | 49.15 | 76.01 | 67.25 | 50.41 | 67.57 | 50.04 |
| MiniCPM (2B) | 62.03 | 26.46 | 79.76 | 26.49 | 53.75 | 54.22 |
| BitNet-b1.58 (2B)* | 49.91 | 74.79 | 80.18 | <u>68.44</u> | **77.09** | **71.90** |
| LLaMA3.2 (3B) | <u>80.68</u> | **88.36** | <u>80.89</u> | 63.75 | 71.38 | 53.75 |
| Qwen3-PTQTP-b1.58-D$^\dagger$ (1.7B) | 51.53 | 70.02 | 27.55 | 37.92 | 57.89 | 51.07 |
| Qwen2.5-PTQTP-b1.58-D$^\dagger$ (3B) | 75.93 | 85.54 | 78.38 | 60.87 | 66.32 | <u>56.27</u> |
| LLaMA3.2-PTQTP-b1.58-D$^\dagger$ (3B) | 74.92 | 82.89 | 74.34 | 55.61 | 64.42 | 51.85 |
| Qwen3-PTQTP-b1.58-D$^\dagger$ (4B) | 76.95 | <u>86.60</u> | **84.13** | 71.02 | <u>73.18</u> | 53.67 |

Table 8: Comparison between PTQTP and binary PTQ methods on LLaMA family across multiple datasets: WikiText2, PTB ([Marcus et al., 1994](#)), C4 (Group Size is 128). N/A: No corresponding result; D$^\dagger$: Dual. **Bold**: best result.

| Model | Method/#Bits | WikiText2 | | | PTB | | | C4 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 7B/8B | 13B | 65B/70B | 7B/8B | 13B | 65B/70B | 7B/8B | 13B | 65B/70B |
| | FP16 | 5.68 | 5.09 | 3.53 | 41.15 | 28.09 | 25.05 | 7.34 | 6.79 | 5.81 |
| LLaMA | GPTQ/2.00 | 129.19 | 20.46 | 8.66 | 1421.47 | 224.45 | 47.70 | 79.06 | 18.97 | 10.23 |
| | PB-LLM/1.70 | 82.76 | 44.93 | 12.81 | 603.57 | 237.22 | 119.19 | 76.63 | 40.64 | 15.30 |
| | BiLLM/1.09 | 49.79 | 14.58 | 8.37 | 373.81 | 84.87 | 44.68 | 46.96 | 16.83 | 11.09 |
| | ARB-LLM$_{RC}$/1.09 | 14.03 | 10.18 | 6.56 | 195.94 | 54.38 | 32.20 | 17.38 | 12.48 | 8.91 |
| | PTQTP-b1.58-D$^\dagger$ | **6.40** | **5.66** | **4.04** | **50.41** | **30.25** | **24.67** | **8.32** | **7.41** | **6.28** |
| | FP16 | 5.47 | 4.88 | 3.32 | 37.9 | 50.93 | 24.25 | 7.26 | 6.73 | 5.71 |
| LLaMA-2 | GPTQ/2.00 | 52.22 | 23.63 | 8.18 | 5583.96 | 419.07 | 50.51 | 35.27 | 19.66 | 9.55 |
| | PB-LLM/1.70 | 66.41 | 236.40 | 28.37 | 657.24 | 816.31 | N/A | 80.69 | 184.67 | N/A |
| | BiLLM/1.08 | 32.31 | 21.35 | 13.32 | 5243.01 | 309.12 | 72.02 | 39.38 | 25.87 | 17.30 |
| | ARB-LLM$_{RC}$/1.08 | 16.44 | 11.85 | 6.16 | 389.59 | 198.17 | 32.79 | 20.38 | 14.36 | 8.65 |
| | PTQTP-b1.58-D$^\dagger$ | **6.32** | **5.29** | **3.95** | **42.53** | **46.9** | **28.15** | **8.41** | **7.34** | **6.3** |
| | FP16 | 6.14 | N/A | 2.86 | 10.05 | N/A | 8.53 | 9.54 | N/A | 7.17 |
| LLaMA-3 | GPTQ/2.00 | 1480.43 | N/A | 82.23 | 717.24 | N/A | 79.20 | 394.74 | N/A | 122.55 |
| | PB-LLM/1.70 | 73.08 | N/A | 22.96 | 106.25 | N/A | 45.13 | 104.15 | N/A | 40.69 |
| | BiLLM/1.06 | 55.80 | N/A | 66.30 | 87.25 | N/A | 97.13 | 61.04 | N/A | 198.86 |
| | ARB-LLM$_{RC}$/1.06 | 27.42 | N/A | 11.10 | 45.49 | N/A | 15.34 | 35.70 | N/A | 15.44 |
| | PTQTP-b1.58-D$^\dagger$ | **8.51** | N/A | **6.12** | **14.02** | N/A | **10.91** | **13.24** | N/A | **12.64** |

Table 9: Performance comparison between PTQTP and binary PTQ methods on Qwen3-14B across mathematical reasoning and general benchmarks (accuracy %). D$^\dagger$: Dual. **Bold**: best result.

| Model | Math-500 | GSM8K | ARC-C | ARC-E | BoolQ | HellaSwag | PIQA | MMLU | WinoG |
|---|---|---|---|---|---|---|---|---|---|
| Baseline-FP16 | 86.60 | 89.39 | 95.25 | 90.83 | 89.45 | 88.09 | 81.07 | 79.38 | 68.90 |
| PBLLM-b1.07 | 0.00 | 1.21 | 4.75 | 7.76 | 29.14 | 7.38 | 25.46 | 12.11 | 49.57 |
| BiLLM-b1.05 | 0.00 | 0.83 | 50.17 | 4.76 | 49.39 | 22.15 | 7.40 | 24.40 | 49.96 |
| ARB-LLM$_{RC}$-b1.05 | 1.80 | 32.22 | 76.95 | 82.01 | 80.00 | 50.01 | 69.91 | 54.08 | 50.36 |
| PTQTP-b1.58-D$^\dagger$ | **82.40** | **85.44** | **93.56** | **88.54** | **89.02** | **85.12** | **79.05** | **76.18** | **67.88** |