# TOWARDS QUANTIZATION-AWARE TRAINING FOR ULTRA-LOW-BIT REASONING LLMS

**Anonymous authors** 

000

001

002003004

010

011

012

013

014

016

018

019

021

024

025

026

027

028

029

031

033

034

037

038

040

041

043

044

046

047

048

049

050 051

052

Paper under double-blind review

#### **ABSTRACT**

Large language models (LLMs) have achieved remarkable performance across diverse reasoning tasks, yet their deployment is hindered by prohibitive computational and memory costs. Quantization-aware training (QAT) enables ultra-lowbit compression (< 4 bits per weight), but existing QAT methods often degrade reasoning capability, partly because complex knowledge structures are introduced during the post-training process in LLMs. In this paper, through a systematic investigation of how quantization affects different data domains, we find that its impact on pre-training and reasoning capabilities differs. Building on this insight, we propose a novel two-stage QAT pipeline specifically designed for reasoning LLMs. In the first stage, we quantize the model using mixed-domain calibration data to preserve essential capabilities across domains; in the second stage, we fine-tune the quantized model with a teacher-guided reward-rectification loss to restore reasoning capability. We first demonstrate that mixed-domain calibration outperforms single-domain calibration at maximum 2.74% improvement on average over six tasks including reasoning and pre-trained tasks. Following experiments on five reasoning benchmarks show that our 2-bit-quantized Qwen3-8B outperforms post-training quantization (PTQ) baselines by 50.45% on average. Moreover, compared to ultra-low-bit-specialized models such as BitNet-2B4T, our pipeline achieves about 2\% higher mathematical-reasoning accuracy using only 40K training sequences.

# 1 Introduction

Large language models (LLMs) have demonstrated remarkable performance across various tasks, including mathematics (Shao et al., 2024; Wang et al.; Yang et al., 2024), coding (Hui et al., 2024; Roziere et al., 2023), and knowledge-intensive question answering (Lu et al., 2022). However, their prohibitive computational and memory requirements pose significant challenges for deployment in inference. One promising direction for reducing these inference costs is weight quantization (Zhou et al., 2024; Lang et al., 2024), which employs low-bit widths for model weights. Among various quantization methods, *quantization-aware training* (QAT), which fine-tunes the model with quantized weights, is especially effective for *ultra-low-bit widths* (< 4 bits) (Wang et al., 2023; Ma et al., 2024; Xu et al., 2024), enabling us to deploy lightweight and fast LLMs. For example, 2-bit quantized LLMs via QAT can achieve performance comparable to their pre-quantized fp16 counterparts (Ma et al., 2024; Kaushal et al., 2024; Liu et al., 2025c).

Despite the promising performance of QAT, existing approaches suffer from severe performance degradation on reasoning benchmarks (Du et al., 2024), such as mathematics, and instruction-following tasks (Lee et al., 2025). We hypothesize that this degradation arises from the complex knowledge structures introduced during post-training. The post-training process is an extensive process that includes supervised fine-tuning (Wei et al., 2021) and preference optimization (Ouyang et al., 2022; Rafailov et al., 2023), introducing new reasoning capabilities with existing commonsense knowledge acquired during pre-training. While it creates heterogeneous knowledge structures, it remains unclear how quantization affects the model's performance on reasoning capabilities and pre-trained commonsense knowledge.

To address this gap, we conduct a systematic investigation of how quantization impacts different knowledge domains in post-training LLMs. Our analysis reveals that quantization creates inherent trade-offs between commonsense knowledge preservation and reasoning capability retention, where different domains exhibit varying sensitivity to quantization. Specifically, while performance on

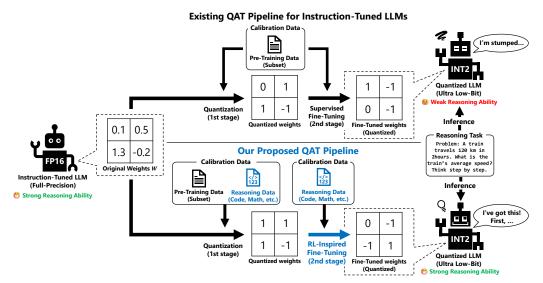


Figure 1: Comparison of the existing QAT pipeline with the proposed pipeline.

commonsense knowledge benchmarks remains relatively stable even with quantization using outof-domain data, reasoning capabilities exhibit significant sensitivity to quantization data, suggesting that different domains have distinct requirements for effective quantization.

Based on this analysis, we introduce a quantization framework specifically designed for post-trained LLMs that address diverse knowledge domains through a novel two-stage pipeline. Following our observation, our quantization framework is designed to dedicate computational resources to maintain reasoning capability, with minimal efforts to preserve general knowledge. Specifically, the first stage carries out block-wise quantization with mixed-domain calibration. This mixed-domain calibration preserves essential reasoning capabilities that are difficult to restore, while also maintaining commonsense knowledge. Subsequently, we perform end-to-end fine-tuning with reinforcement learning inspired objectives to enhance reasoning capability. This unified framework enables extremely low-bit quantization of post-trained LLMs with minimal reasoning performance degradation.

Extensive experiments on five reasoning benchmarks demonstrate the effectiveness of our approach. Our method achieves significant improvements over existing post-training quantization methods for reasoning LLMs. Specifically, our 2-bit quantized Qwen3-8B outperforms other quantization methods by 50.45% on average. Notably, even when compared to specialized ternary LLMs like BitNet-2B4T, our 2-bit model with 1.7B parameters demonstrates superior mathematical reasoning performance with substantially reduced training costs—achieving 2.5% improvement using only 40 K training sequences.

Our contributions can be summarized as follows:

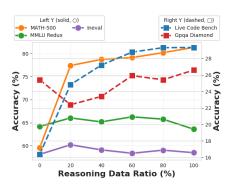
- We empirically demonstrate how quantization differently affects commonsense knowledge acquired during pre-training and reasoning capabilities developed in post-training. Our results highlight the importance of designing mixed calibration data to effectively preserve both of them.
- We propose two-stage quantization pipeline for post-trained LLMs that combines mixed-domain calibration and RL-inspired fine-tuning to preserve reasoning capabilities while achieving extremely low-bit quantization.
- We demonstrate that our approach achieves state-of-the-art (SOTA) performance on multiple reasoning benchmarks with both 2-bit and 3-bit quantization.

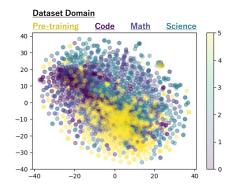
#### 2 Preliminaries

This section outlines the weight quantization and recent quantization-aware training (QAT) pipeline.

**Weight Quantization:** Weight quantization maps the model weights to low-bit width counterparts. Given a full-precision weight  $w \in \mathbb{R}$ , we obtain its dequantized approximation w' via

$$q := \text{clamp}(|w/s| + z, 0, 2^N - 1), \quad w' := s(q - z),$$





- (a) Performance for general (green and purple •) and mathematical (orange •, blue ■, and red ■) tasks on the reasoning data ratio.
- (b) t-SNE visualization of the activations in the 14th Transformer block, where colors denote dataset domains

Figure 2: Impact of data-domain composition on Qwen3-1.7B. In the ultra-low-bit quantized model, replacing a part of the calibration data drawn from the pre-training data (FineWeb-Edu) with reasoning data (OpenThoughts3-1.2M) leaves general-task accuracy unchanged while improving reasoning accuracy (left). In the full-precision model, t-SNE of the model's activations shows tight clusters for pre-training inputs but wide dispersion for reasoning inputs (right).

where s>0 is the scale factor,  $z\in\mathbb{R}$  is the zero point, and N is the target bit width.  $\lfloor\cdot\rceil$  represents the nearest integer function (i.e., the round function), and  $\operatorname{clamp}(x,a,b)$  clamps input x to the interval [a,b]. Since the scale s and zero point z are shared across groups of weights (e.g., an entire matrix, a channel, or a block), each weight is represented only by an N-bit code q, with s and z stored once per group, achieving a low-bit width per-weight representation.

Weight quantization approaches are mainly categorized into two strategies: 1) post-training quantization (PTQ), which converts pre-trained model weights to low-bit widths without retraining; and 2) quantization aware training (QAT), which quantizes and fine-tunes the weights simultaneously. PTQ can quickly quantize weights with small or even without calibration data, while it struggles with ultra-low-bit quantization. In contrast, QAT can flexibly fine-tune the full-precision weights w, scaling factor s, and zero point s, achieving performance comparable to the full-precision model in ultra-low-bit scenarios.

Quantization-Aware Training Pipeline: As illustrated in the top panel of Figure 1, existing QAT pipelines generally comprise two stages: 1) an initial quantization stage; and 2) a fine-tuning stage. The first stage initializes the quantized weights that serve as the starting point for the subsequent stage. Some methods omit this step, whereas the latest state-of-the-art (SOTA) QAT approaches (Chen et al., 2024; Du et al., 2024) have demonstrated that quantizing weights using a subset of the pre-training dataset as calibration data enables stable fine-tuning during the subsequent stage. In the second stage, the weights quantized in the first stage are fine-tuned by minimizing a training objective. All parameters (i.e., weights, scale, zero point) or some of them are fine-tuned, and this paper fine-tunes only the scale, following one of the SOTA QAT approaches, EfficientQAT. Also, the training objective is typically either a self-supervised pre-training loss (Liu et al., 2025c; Chen et al., 2024) or a knowledge-distillation loss (Du et al., 2024; Lee et al., 2025).

#### 3 REASONING-ORIENTED TWO-STAGE QUANTIZATION AWARE TRAINING

This section proposes a novel QAT pipeline that enables the preservation of reasoning capabilities after ultra-low-bit quantization. We first analyze the impact of quantization on various knowledge domains, and based on these findings, we introduce the reasoning-oriented QAT pipeline.

# 3.1 QUANTIZATION IMPACTS ACROSS KNOWLEDGE DOMAINS

This section analyzes how domain selection for calibration data affects the overall model performance. Existing quantization approaches mainly perform quantization with either pre-training data (Liu et al., 2025c; Chen et al., 2024), or domain-specific data, such as mathematics (Liu et al., 2025a). While previous work has selected calibration data tailored to specific target tasks, the crosstask implications of such task-specific calibration choices remain largely unexplored.

To investigate the effect of domain selection for calibration data, we analyze the impact of selecting different calibration datasets on performance across multiple tasks and knowledge domains. As shown in the results of Qwen3-1.7B quantized to 3-bits by EfficientQAT (Chen et al., 2024) (Figure 2a), the tasks can be broadly grouped into two trends: 1) tasks for which performance improves as the amount of reasoning data increases; and 2) tasks for which performance remains almost constant regardless of the amount of reasoning data. Notably, all tasks in the first category are represented in the reasoning dataset, i.e., code (blue), mathematical tasks (orange), scientific questions and answers (red). These results indicate that reasoning data tends to suffer from domain shift, while tasks related to commonsense knowledge are less sensitive to calibration datasets. On the other hand, common tasks also demonstrate performance degradation when calibrated with pure reasoning data, suggesting the importance of dataset diversity even for tasks that appear less sensitive to calibration choices.

These distinct trends happen as the intermediate distributions between pre-trained data and reasoning data differ, as shown in Figure 2b. The distributional mismatch leads to suboptimal quantization performance when calibrating on single-domain data, resulting in higher quantization errors for tasks that require domain-specific representations.

#### 3.2 Proposed Method

We now introduce the novel QAT pipeline for ultra-low-bit reasoning LLMs.

**Knowledge Domain Selection in Calibration data:** We first focus on the mixing ratio of the knowledge domain in calibration data. The results in Section 3.1 illustrate the importance of selecting appropriate calibration data when a QAT pipeline is applied to post-trained LLMs. In particular, it is important to mix pre-training data and reasoning data in an appropriate ratio. Building on these findings, we propose using novel calibration data in the first stage of the QAT pipeline. This data is composed of 80% reasoning-focused data and 20% pre-training data, designed to bias the calibration process toward reasoning while retaining coverage of pre-training distributions.

**Supervised Fine-Tuning With Reward Rectification Loss:** We secondly aim at the fine-tuning stage. Quantization of the first stage using proposed calibration data mixed with pre-training data preserves the fundamental capabilities of the LLM, enabling us to focus on enhancing reasoning capabilities during the fine-tuning stage. A straightforward approach to enhance reasoning ability is to perform supervised fine-tuning using reasoning data, but such training does not effectively generalize into reasoning data (Chu et al., 2025). Employing reinforcement learning could improve generalization on reasoning data, but online text generation incurs auto-regressive text generations, resulting in huge training overhead. To balance training efficiency and generalization on unseen data, we employ reweighted rectification (Wu et al., 2025) for supervised fine-tuning to make the objective function reinforcement-like.

Reward rectification is scaling factors for the loss function in supervised fine-tuning. Given the datasets  $\mathcal{D} = \{x, y^*\}$  and the supervised fine-tuning loss  $\mathcal{L}_{SFT}(\theta)$ , reward rectification loss  $\mathcal{L}(\theta)$  dynamically reweights the supervised loss as follows:

$$\mathcal{L}(\theta) = \mathcal{L}_{SFT}(\theta) \cdot sg(1/w),$$

where w is the dynamic reweighting factor and  $sg(\cdot)$  denotes the stop-gradient operator.

This formulation can be viewed as bridging supervised fine-tuning and reinforcement learning. In particular, choosing  $w=1/\pi_\theta(y\mid x)$  yields a gradient equivalent to an on-policy policy-gradient update with the reward function:

$$r(x,y) = \mathbf{1}[y = y^*],$$

where  $\pi_{\theta}(y \mid x)$  is the model's conditional probability of generating an output y given an input x under parameters  $\theta$ . This dynamic re-weighting can avoid over-concentration on low-probability reference tokens, improving generalization despite not using additional sampling or reward functions.

Table 1: Accuracy comparison for different calibration data on 6 benchmarks. Higher values are better. We define the group size as 128. Mixed data contains 80% of reasoning data and 20% of pre-training data.

		Re	Reasoning Tasks			Pre-trained Tasks			
Model (Qwen3)	Bit Width	Dataset Type	MATH- 500	Live Code Bench	GPQA- Diamond	MMLU- Redux	CSR	IFEVAL	Avg.
1.7B		Pre-training	59.53	16.36	25.42	64.16	57.96	58.60	47.01
	w3	Reasoning	81.33	29.35	26.60	63.58	55.82	59.52	52.70
		Mixed	80.20	29.32	25.42	65.76	56.71	59.70	52.85
	w2	Pre-training	0.13	0.00	0.00	0.00	50.75	13.86	10.79
		Reasoning	20.60	0.09	7.58	19.93	44.77	25.51	19.75
		Mixed	18.68	0.47	7.58	28.92	49.11	24.58	21.50
4B	w3	Pre-training	81.80	35.42	41.92	77.98	63.94	71.16	62.04
		Reasoning	90.90	46.89	42.76	78.55	63.71	71.72	65.76
		Mixed	90.70	46.85	45.45	78.91	63.97	74.86	66.79
40		Pre-training	2.73	0.00	6.23	26.89	59.34	17.74	18.82
	w2	Reasoning	33.80	5.78	11.62	46.26	53.41	32.90	30.63
		Mixed	22.60	5.12	14.14	51.20	55.89	31.05	30.00
8B -		Pre-training	87.00	37.25	41.66	82.60	69.19	76.34	65.6
	w3	Reasoning	91.80	53.84	51.52	81.75	68.14	75.79	70.4
		Mixed	92.40	51.75	48.99	83.33	69.16	81.15	71.1
		Pre-training	5.33	0.28	4.55	41.72	63.35	19.41	22.4
	w2	Reasoning	42.27	8.15	10.44	51.40	55.05	35.86	33.8
		Mixed	40.00	7.30	11.11	57.11	60.81	43.25	36.6

While the original reward rectification uses the student model's own probability  $\pi_{\theta}(y \mid x)$  for reweighting, in the QAT, the quantized model's distribution becomes less reliable due to precision loss. Using the quantized model's own probabilities for reweighting could amplify these errors.

Therefore, we leverage the teacher model's probability  $\pi_t(y^*|x)$  as a more reliable reference for the reweighting factor. This teacher-guided approach ensures that the reweighting process is based on the target distribution we aim to recover, rather than the potentially corrupted distribution of the quantized model.

Thus, we introduce teacher-guided reward rectification loss  $\mathcal{L}(\theta)$ , where the teacher model  $\pi_t$  controls the scale of supervised loss function. Given the teacher probability with labeled data  $\pi_t(y^* \mid x)$ , teacher guided reward rectification loss can be represented as:

$$\mathcal{L}_t(\theta) = \mathcal{L}_{SFT}(\theta) \cdot sg(\pi_t(y^*|x)).$$

Intuitively, this formulation represents that the supervised loss values are amplified when the probability of the quantized model for the label is smaller than that of the teacher probability. When the distribution of the quantized model becomes close to the original distribution, this scaling factor acts as the original reward rectification.

To align the overall probabilistic distribution of the quantized model with original LLMs, we further introduce an additional KL divergence loss. Finally, our training loss function can be represented as:

$$\mathcal{L}(\theta) = \alpha \mathcal{L}_t(\theta) + \beta D_{\text{KL}}(\pi_{\text{T}}(\cdot|x)||\pi_{\text{S}}(\cdot|x)), \tag{1}$$

where  $D_{\mathrm{KL}} \big( \pi_{\mathrm{T}}(\cdot|x) || \pi_{\mathrm{S}}(\cdot|x) \big) = \sum_{y} \pi_{\mathrm{T}}(y|x) \log \frac{\pi_{\mathrm{T}}(y|x)}{\pi_{\mathrm{S}}(y|x)}$  is the KL divergence between the fp16 model and the quantized model, and  $\alpha, \beta$  is hyperparameters that control the effects of teacher-guided reward rectification loss and kl divergence loss.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETTINGS

**Training:** We conduct experiments on Qwen3 instruction-tuned models (Yang et al., 2025). For block-wise calibration, we use a total 4,096 samples with a context length of 2,048. Calibration

datasets consist of 80% of sequences sampled from OpenThoughts-1.2M (Guha et al., 2025) and the remaining 20% sampled from FineWeb-Edu (Lozhkov et al., 2024). We use different learning rates for quantization parameters (1e-4) and weight parameters (1e-5). For 2-bit quantization, we use a larger learning rate of 2e-5 for weights.

During supervised fine-tuning, models are with 32,768 samples from OpenThoughts-1.2M. We optimize all trainable parameters with the same learning rate. The learning rate for 3-bit quantization is 1e-6, while we use a larger learning rate for 2-bit quantization, 5e-6 for the 1.7B parameter, and 1e-4 for other parameters. We use the AdamW optimizers (Loshchilov & Hutter, 2019) with the cosine annealing learning rate decay (Loshchilov & Hutter, 2017). Models are fine-tuned with a batch size of 64 and one epoch for 3-bit models and 3 epochs for 2-bit models. We filter out the top-20 probabilities for the KL loss. We set  $\alpha=0.2$  and  $\beta=1.0$  in Equation (1) unless explicitly stated otherwise.

**Evaluation:** We evaluate the zero-shot accuracy on five benchmarks including We evaluate the zero-shot accuracy on five benchmarks, including MATH-500 (Lightman et al., 2023), Live Code Bench (White et al., 2024), MMLU-Redux (Gema et al., 2024), GPQA- Diamond (Rein et al., 2024), and IFEval (Zhou et al., 2023), using the evalscope (Team, 2024). These tasks are evaluated in open-ended text generation. We use token-level sampling, whose tokens are sampled from the top 20 highest tokens with a temperature of 0.6. We basically use a maximum sequence length of 32K for all benchmarks. However, we reduce the maximum sequence length to 8K on the lower-performance models to avoid excessive text generation due to the absence of a stop token. All evaluations are conducted three times, and we report the average accuracy.

Quantization baselines to post-trained LLMs: We compare our method with two PTQ quantization baselines, GPTQ (Frantar et al., 2022) and AWQ (Lin et al., 2024), both of which are evaluated on reasoning benchmarks (Liu et al., 2025a). To quantize these two baselines, we follow a similar strategy as conducted by Liu et al. (2025a). Specifically, we perform quantization using 128 samples from NuminaMath (LI et al., 2024). We reproduce these quantized models locally, except for 3 and 4 bits AWQ quantization for Qwen-8B, as reproduced performance is much inferior to the performance claimed in the paper (Liu et al., 2025a).

# 4.2 Dataset effects on overall performance

This section evaluates our proposed calibration datasets against single-domain calibrations using either pre-training data or reasoning data. We evaluate five benchmarks described in Section 4.1 with the additional commonsense reasoning tasks (CSR) to evaluate general knowledge of quantized models. The CSR includes five subset tasks: ARC-e, ARC-c (Clark et al., 2018), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), and WinoGrande Sakaguchi et al. (2021). Table 1 demonstrates that mixed domain calibration outperforms single-domain calibration across various parameters in both 2-bit and 3-bit quantization settings. Compared with pre-training datasets, its performance improvements on reasoning benchmarks, including mathematical and coding tasks, are particularly notable. In addition, mixed data achieves comparable performance to pre-trained data and other benchmarks in reasoning tasks. When compared with reasoning-only datasets, the performance of mixed datasets on reasoning benchmarks is close, while performance in pre-trained domain benchmarks tends to be superior. These results suggest that including a small portion of pre-training data can be effective in maintaining commonsense knowledge. resulting in better generalization results across a wide range of tasks.

## 4.3 COMPARISON WITH PRIOR QUANTIZATION APPROACHES.

This section compares our proposal with other quantization methods, including GPTQ and AWQ. Table 2 shows that our proposal significantly outperforms existing quantization approaches in both 2-bit and 3-bit quantization settings. The performance improvements are particularly notable in 2-bit quantization settings. While the performance of GPTQ or AWQ quantized models is extremely low, our quantized models not only achieve solid performance but also exhibit steadily increasing accuracy as the number of parameters grows. This trend suggests that our quantization method scales effectively with model size even for extremely low-bit quantization. For 3-bit quantization, our approach dramatically improves performance, particularly on smaller models. For example, our

Table 2: Comparison of quantization methods and bit-widths on Qwen3 models. Values are % (higher is better). We define the group size as 128.

Settings			Benchmarks (%)					
Model	Method	Bit-width (W/A)	MATH500	LiveCodeBench	MMLU-Redux	GPQA-Diamond	IFEval	Avg.
Qwen3	-1.7B							
	FP (Baseline)	bfloat16	89.00	53.60	74.70	38.38	70.43	65.22
	GPTQ	4/16	86.33	36.87	70.23	34.34	66.17	58.79
	AWQ	4/16	87.40	44.39	71.61	35.86	65.25	60.90
	GPTQ	3/16	58.20	0.03	42.38	9.26	31.79	28.3
	AWQ	3/16	57.93	5.81	53.51	17.51	47.69	36.49
	Proposal	3/16	82.67	32.96	67.70	31.65	61.00	55.20
	GPTQ	2/16	2.07	0.00	5.76	4.71	8.50	4.21
	AWQ	2/16	0.00	0.00	27.25	8.08	12.26	9.52
	Proposal	2/16	48.60	6.54	40.09	14.48	32.22	28.39
Owen3	-4B							
	FP (Baseline)	bfloat16	93.60	71.19	84.25	51.52	83.55	76.82
	GPTQ	4/16	93.40	66.16	82.06	50.67	81.15	74.69
	AWQ	4/16	93.00	65.69	82.95	50.17	80.96	74.55
	GPTO	3/16	84.93	21.17	65.59	24.92	54.71	50.26
	AWQ	3/16	88.33	37.19	74.16	33.67	71.90	61.05
	Proposal	3/16	89.53	50.20	79.79	46.80	75.29	68.32
	GPTO	2/16	3.67	0.00	7.50	8.92	8.50	5.71
	AWQ	2/16	0.00	0.00	0.0	0.00	11.83	2.37
	Proposal	2/16	77.13	19.49	61.69	26.94	50.09	47.07
Qwen3	-8B							
•	FP (Baseline)	bfloat16	94.00	73.00	87.30	61.62	86.51	80.49
	GPTQ	4/16	94.60	69.57	86.84	58.08	86.88	79.19
	AWQ	4/16	97.0	54.7	N/A	59.6	N/A	N/A
	GPTQ	3/16	92.07	39.39	79.04	46.80	76.89	66.83
	AWQ	3/16	92.9	35.3	N/A	46.8	N/A	N/A
	Proposal	3/16	91.47	60.03	84.50	47.47	78.80	72.45
	GPTQ	2/16	2.80	0.00	6.38	5.22	8.69	4.62
	AWQ	2/16	0.00	0.00	5.93	3.87	10.17	3.99
	Proposal	2/16	80.40	28.50	72.63	34.51	59.33	55.07

3-bit quantization of Qwen3-1.7B achieves an average accuracy of 55.20% across five tasks, which is 18.71% higher than existing PTQ methods. These results highlight that our approach is especially effective when model capacity is constrained, such as in cases of ultra-low bit widths or limited parameter counts.

## 4.4 ABLATION STUDY

There are very few quantization-aware (QAT) training approaches that can be directly compared to ours, as most existing methods target different evaluation settings. Instead, this section presents ablation studies that compare our approach with key components derived from existing QAT methods.

Effectiveness of Block-wise Calibration: The main differences between our approach and existing methods are the use of block-wise calibration before fine-tuning and the loss function. To analyze the effects of these two components, we either replace the loss function with conventional cross-entropy loss, which is basically used in QAT (Liu et al., 2025c). In these experiments, we fine-tune the 2-bit quantized Qwen-3 1.7B model for one epoch using 32K sequences with a learning rate of 5e-6. Table 3 summarizes the results. Here, "S" denotes the conventional supervised fine-tuning with cross entropy loss function, "R" denotes the teacher-guided reward rectification loss in Section 3.2, and "C" means the existence of a calibration stage. Therefore, "C+S" denotes supervised fine-tuning after calibration. As shown in Table 3, both calibration data and proposed loss function significantly enhance model performance. We also find that modifying the supervised loss led to substantial improvements on reasoning benchmarks. In particular, on MATH-500, the accuracy increased by 15.43% when moving from cross entropy loss to our proposed loss, and on Live Code Bench, it increased by 5.75%. More importantly, the performance on reasoning benchmarks degrades after conducting supervised fine-tuning with cross-entropy loss. These results indicate that conventional QAT approaches, which rely primarily on cross-entropy loss for supervised fine-tuning, are insufficient for post-trained LLMs.

Table 3: Ablation on block-wise calibration and loss choice: "S" denotes conventional supervised fine-tuning; "R" denotes our proposed loss; and "C" indicates the use of block-wise calibration.

Training	S	R	С	C+S	C+R
MATH-500	1.4	1.60	28.57	22.70	38.13
Live Code Bench	0.0	0.0	0.47	0.00	5.75
MMLU Redux	3.54	3.53	28.57	35.78	36.64
<b>GPQA Diamond</b>	6.06	6.06	7.58	14.31	14.14
IFEval	10.72	12.20	24.58	23.66	31.61

Table 4: Effect of the loss function on reasoning performance. Methods evaluated with 3-bit quantization on Qwen3 1.7B.

$(\alpha, \beta)$	(1.0, 0.0)	(0.0, 1.0)	(0.2, 1.0)
MATH-500	78.20	82.80	82.67
Live Code Bench	28.47	32.35	32.96
MMLU Redux	66.14	66.95	67.70
<b>GPQA Diamond</b>	21.72	30.30	31.65
IFEval	60.63	63.03	62.11

Table 5: Benchmark comparison of our proposal with INT2 Qwen3 family and BitNet b1.58 2B.

Benchmark (Metric)	Qwen3	Qwen3	Qwen3	BitNet b1.58
	8B-int2	4B-int2	1.7B-int2	2B
Training Tokens	0.2M	0.2M	0.8M	4T
Activation	bf16	bf16	bf16	int8
MATH-500 (0-shot; EM)	80.13	77.13	48.60	43.40
GSM8K (4-shot; EM)	88.93	81.71	57.47	58.38
IFEval (0-shot; Instruct-Strict)	59.33	50.09	45.29	53.48
Average	76.13	69.91	50.75	51.75

**Effectiveness of Loss Weighting:** This section studies the contribution of two loss terms, the teacher-guided reward-rectification loss and the KL-divergence loss, to overall performance. We fine-tune a 3-bit quantized Qwen3-1.7B for a single epoch with different  $\alpha$  and  $\beta$ . We evaluate three different weighting schemes:  $(\alpha, \beta) = (1, 0)$ , which applies only the reward rectification loss;  $(\alpha, \beta) = (0, 1)$ , which applies only the KL divergence loss; and  $(\alpha, \beta) = (0.2, 1)$  (i.e., our proposed configuration), which combines both losses with the specified weights. Table 4 demonstrates that combining the two losses improves overall model performance.

# 4.5 Comparison with BitNet1.58 2B4T

This section compares our quantized models with BitNet1.58 2B4T, a native ternary LLM trained from scratch. To align the bit-width, our QAT pipeline quantizes Qwen3 models into 2 bits.

Table 5 describes the accuracies on two mathematical benchmarks including MATH-500 and GSM8K, and IFEval. We referred to the results of BitNet1.58 2B4T from (Ma et al., 2025). As shown in Table 5, our INT2 quantized model achieves superior mathematical performance with lower parameter requirements and significantly fewer tokens required for the quantization process. These results demonstrate that by designing an appropriate QAT pipeline, it is possible to leverage pre-trained features, leading to promising reasoning performance train high-accuracy 2-bit models at a fraction of the training costs.

In addition, our approach demonstrates superior scalability compared to BitNet 1.58 2B4T. Because we fine-tune pre-trained LLMs using only a limited number of sequences, we can easily produce models with different parameter counts. This enables a flexible trade-off between performance and resource usage, as illustrated in Table 5, which demonstrates results of several parameter variations of our quantized models.

# 5 RELATED WORKS

In this section, we briefly summarize the quantization approaches. Quantization approaches can be categorized into post-training quantization (PTQ) and quantization-aware training (QAT) depending on whether fine-tuning is performed or not. This section deals with weight-only quantization of large language models (LLMs) addressed in this work.

**Post-training quantization** (PTQ) converts full-precision weights into lower-bit counterparts without relying on fine-tuning. To obtain better quantization parameters, recent methods optimize the reconstruction problem either at the linear projection level (Frantar et al., 2022; Lin et al., 2024)

or at the transformer block level (Lee et al., 2023; Shao et al., 2023). While PTQ has achieved strong initial success in LLMs, initial approaches still face limitations in achieving extremely low-bit quantization without losing their performance. To overcome these challenges, research has shifted toward more aggressive quantization, such as 3-bit or 2-bit. Some approaches target such low-bit quantization with integer representation (Shao et al., 2023; Zhao et al., 2024; Chee et al., 2023), demonstrating noticeable performance at these bit-widths. To further improve the trade-offs between accuracy and model size, recent approaches introduce vector quantization (Egiazarian et al., 2024; Tseng et al., 2024; Malinovskii et al., 2024). Despite their promising performance, vector quantization introduces substantial overhead in inference (Gong et al., 2024).

Quantization-aware training (QAT), in contrast, can enhance quantized model performance by incorporating fine-tuning. With the additional computational cost for fine-tuning, QAT enables the use of hardware-friendly numerical representations, such as integers, for low-bit quantization, resulting in minimal overhead at inference time. There are several choices for optimization targets for fine-tuning. LLM-QAT (Liu et al., 2023) and BitDistiller (Du et al., 2024) explore knowledge distillation within QAT literature. BitNet b1.58 (Ma et al., 2024), Spectra (Kaushal et al., 2024), and ParetoQ (Liu et al., 2025c) employ fine-tuning in a self-supervised manner using pre-training data. By spending billions of tokens for fine-tuning, these approaches realize promising performance with ternary or 2-bit. Given the substantial training costs of these approaches, recent work has focused on improving the training efficiency of QAT approaches. EfficientQAT (Chen et al., 2024) introduces two two-stage pipeline that perform end-to-end backpropagation following block-wise calibration. UPQ (Lee et al., 2025) modifies the two-stage QAT pipeline to use knowledge distillation and progressive quantization, demonstrating the promising performance on instruction-tuned LLMs. However, most existing quantization approaches have primarily focused on pre-training LLMs, with limited exploration of their effectiveness on complex reasoning capabilities that are crucial for modern LLM applications. In this paper, we investigate how quantization affects reasoning performance and propose methods to preserve reasoning capabilities in quantized LLMs.

Quantization and Reasoning. Several comprehensive analyses have explored the effects of quantization on reasoning capability. Li et al. (2025) and Liu et al. (2025b) demonstrate that ultra-low-bit quantization leads to severe performance drops on reasoning benchmarks such as mathematical tasks. Liu et al. (2025b) demonstrates that less-than-4-bit quantization leads to severe performance drops on reasoning benchmarks such as mathematical tasks. Mekala et al. (2025) systematically analyzes the effects of quantization on long-context reasoning tasks, demonstrating that even 4-bit models incur substantial losses. Although these analyses reveal critical challenges for existing quantization approaches, few studies have explored effective strategies to maintain reasoning performance under such aggressive settings. However, one notable example is BitNet 2B4T (Ma et al., 2025), which demonstrates strong performance on mathematical benchmarks with ternary LLMs by performing quantization-aware training over four trillion tokens. In this paper, we explore a more efficient approach for reasoning-oriented LLMs. By combining block-wise quantization with RL-inspired fine-tuning using limited tokens, we obtain highly accurate 2- and 3-bit LLMs with significantly fewer fine-tuning sequences.

#### 6 Conclusion

This paper addresses the critical challenge of maintaining reasoning capabilities in ultra-low-bit quantized large language models (LLMs). Through systematic analysis, we demonstrate that quantization affects different knowledge domains unevenly—while pre-training knowledge remains robust, reasoning capabilities show severe degradation. Building on this insight, we develop a novel two-stage quantization-aware training pipeline specifically designed for post-trained reasoning LLMs. Our approach combines mixed-domain calibration with teacher-guided reward rectification loss to preserve and restore reasoning abilities under aggressive quantization. Experiments across five reasoning benchmarks validate our method, with 2-bit quantized Qwen3-8B achieving 50.45% average improvement over existing approaches. Notably, our method outperforms BitNet 2B4T on mathematical reasoning while requiring dramatically fewer training resources—40K sequences versus 4 trillion tokens. We establish the first quantization framework specifically targeting reasoning-oriented LLMs, providing practical solutions for efficient model compression without sacrificing cognitive capabilities. This work provides a foundation for future developments in efficient, high-performance quantized reasoning models, enabling broader deployment of sophisticated AI systems.

## REFERENCES

- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of large language models with guarantees. *Advances in Neural Information Processing Systems*, 36:4396–4429, 2023.
- Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, and Ping Luo. Efficientqat: Efficient quantization-aware training for large language models. *arXiv preprint arXiv:2407.11062*, 2024.
- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*, 2025.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Dayou Du, Yijia Zhang, Shijie Cao, Jiaqi Guo, Ting Cao, Xiaowen Chu, and Ningyi Xu. Bitdistiller: Unleashing the potential of sub-4-bit llms via self-distillation. *arXiv preprint arXiv:2402.10631*, 2024.
- Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. Extreme compression of large language models via additive quantization. *arXiv* preprint arXiv:2401.06118, 2024.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv* preprint arXiv:2210.17323, 2022.
- Ruihao Gong, Yang Yong, Shiqiao Gu, Yushi Huang, Yunchen Zhang, Xianglong Liu, and Dacheng Tao. Llm-qbench: A benchmark towards the best practice for post-training quantization of large language models. *CoRR*, 2024.
- Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer, Liangyu Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff, Shiye Su, Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie Ji, Yichuan Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave, Alon Albalak, Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit Bansal, Saadia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan, Mike A. Merrill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran Sathiamoorthy, Alexandros G. Dimakis, and Ludwig Schmidt. Openthoughts: Data recipes for reasoning models, 2025. URL https://arxiv.org/abs/2506.04178.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- Ayush Kaushal, Tejas Vaidhya, Arnab Kumar Mondal, Tejas Pandey, Aaryan Bhagat, and Irina Rish. Spectra: Surprising effectiveness of pretraining ternary language models at scale. *arXiv* preprint arXiv:2407.12327, 2024.
- Jiedong Lang, Zhehao Guo, and Shuyu Huang. A comprehensive study on quantization techniques for large language models. In 2024 4th International Conference on Artificial Intelligence, Robotics, and Communication (ICAIRC), pp. 224–231. IEEE, 2024.
- Jung Hyun Lee, Jeonghoon Kim, Se Jung Kwon, and Dongsoo Lee. Flexround: Learnable rounding based on element-wise division for post-training quantization. In *International Conference on Machine Learning*, pp. 18913–18939. PMLR, 2023.

- Jung Hyun Lee, Seungjae Shin, Vinnam Kim, Jaeseong You, and An Chen. Unifying block-wise ptq and distillation-based qat for progressive quantization toward 2-bit instruction-tuned llms. *arXiv* preprint arXiv:2506.09104, 2025.
  - Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. Numinamath. [https://huggingface.co/AI-MO/NuminaMath-CoT] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina\_dataset.pdf), 2024.
  - Zhen Li, Yupeng Su, Runming Yang, Congkai Xie, Zheng Wang, Zhongwei Xie, Ngai Wong, and Hongxia Yang. Quantization meets reasoning: Exploring llm low-bit quantization degradation for mathematical reasoning. *arXiv* preprint arXiv:2501.03035, 2025.
  - Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device Ilm compression and acceleration. *Proceedings of machine learning and systems*, 6:87–100, 2024.
  - Ruikang Liu, Yuxuan Sun, Manyi Zhang, Haoli Bai, Xianzhi Yu, Tiezheng YU, Chun Yuan, and Lu Hou. Quantization hurts reasoning? an empirical study on quantized reasoning models. In Second Conference on Language Modeling, 2025a. URL https://openreview.net/forum?id=BM192Ps5Nv.
  - Ruikang Liu, Yuxuan Sun, Manyi Zhang, Haoli Bai, Xianzhi Yu, Tiezheng Yu, Chun Yuan, and Lu Hou. Quantization hurts reasoning? an empirical study on quantized reasoning models. *arXiv* preprint arXiv:2504.04823, 2025b.
  - Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023.
  - Zechun Liu, Changsheng Zhao, Hanxian Huang, Sijia Chen, Jing Zhang, Jiawei Zhao, Scott Roy, Lisa Jin, Yunyang Xiong, Yangyang Shi, et al. Paretoq: Scaling laws in extremely low-bit llm quantization. *arXiv preprint arXiv:2502.02631*, 2025c.
  - Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=Skq89Scxx.
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=Bkq6RiCqY7.
  - Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest collection of educational content, 2024. URL https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu.
  - Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022.
  - Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Lifeng Dong, Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-bit llms: All large language models are in 1.58 bits. *arXiv preprint arXiv:2402.17764*, 1(4), 2024.
  - Shuming Ma, Hongyu Wang, Shaohan Huang, Xingxing Zhang, Ying Hu, Ting Song, Yan Xia, and Furu Wei. Bitnet b1. 58 2b4t technical report. *arXiv preprint arXiv:2504.12285*, 2025.
  - Vladimir Malinovskii, Denis Mazur, Ivan Ilin, Denis Kuznedelev, Konstantin Burlachenko, Kai Yi, Dan Alistarh, and Peter Richtarik. Pv-tuning: Beyond straight-through estimation for extreme llm compression. *Advances in Neural Information Processing Systems*, 37:5074–5121, 2024.

- Anmol Mekala, Anirudh Atmakuru, Yixiao Song, Marzena Karpinska, and Mohit Iyyer. Does quantization affect models' performance on long-context tasks? *arXiv preprint arXiv:2505.20276*, 2025.
  - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
  - Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
  - Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
  - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
  - Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language models. *arXiv* preprint arXiv:2308.13137, 2023.
  - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. *URL https://arxiv. org/abs/2402.03300*, 2(3):5, 2024.
  - Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks. *arXiv* preprint *arXiv*:2402.04396, 2024.
  - Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. Bitnet: Scaling 1-bit transformers for large language models. *arXiv preprint arXiv:2310.11453*, 2023.
  - Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024. *URL https://arxiv. org/abs/2312.08935*.
  - Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
  - Yongliang Wu, Yizhou Zhou, Zhou Ziheng, Yingzhe Peng, Xinyu Ye, Xinting Hu, Wenbo Zhu, Lu Qi, Ming-Hsuan Yang, and Xu Yang. On the generalization of sft: A reinforcement learning perspective with reward rectification. *arXiv preprint arXiv:2508.05629*, 2025.
  - Yuzhuang Xu, Xu Han, Zonghan Yang, Shuo Wang, Qingfu Zhu, Zhiyuan Liu, Weidong Liu, and Wanxiang Che. Onebit: Towards extremely low-bit large language models. *Advances in Neural Information Processing Systems*, 37:66357–66382, 2024.
  - An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
  - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv* preprint arXiv:2505.09388, 2025.
    - Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and accurate llm serving. *Proceedings of Machine Learning and Systems*, 6:196–209, 2024.

Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, et al. A survey on efficient inference for large language models. *arXiv preprint arXiv:2404.14294*, 2024.