

Enhancing Trustworthiness with Mixed Precision: Benchmarks, Opportunities, and Challenges

Guanxi Lu

Department of Electrical and
Electronic Engineering
Imperial College London
guanxi.lu22@imperial.ac.uk

Hao (Mark) Chen

Department of Computing
Imperial College London
hao.chen20@imperial.ac.uk

Zhiqiang Que

Department of Computing
Imperial College London
z.que@imperial.ac.uk

Wayne Luk

Department of Computing
Imperial College London
w.luk@imperial.ac.uk

Hongxiang Fan

Department of Computing
Imperial College London
hongxiang.fan@imperial.ac.uk

Abstract—Large language models (LLMs) have shown promising performance across various tasks. However, their autoregressive decoding process poses significant challenges for efficient deployment on existing AI hardware. Quantization alleviates memory and compute pressure by compressing weights, activations, and KV caches to low precisions while preserving generation quality. However, existing quantization frameworks typically focus on perplexity or classification accuracy, often omitting critical trustworthiness metrics. This gap introduces risks when applying quantized LLMs to downstream high-stakes domains such as finance and healthcare. In this work, we systematically investigate the impact of quantization on four trustworthiness metrics (adversarial robustness, fairness, machine ethics, and out-of-distribution robustness) and identify the instability across compression ratios and quantization methods. Building on these observations, we develop a novel precision-ensemble voting approach that leverages predictions from mixed-precision variants of the same model and consistently improves performance by up to 5.8% on trustworthiness metrics. Our results highlight the importance of considering trustworthiness when developing model compression techniques and point to research opportunities at the intersection of compression and trustworthiness for safety-critical applications.

Index Terms—large language models, model quantization, mixed precision, model compression, natural language processing, low-bit inference, post-training quantization.

I. INTRODUCTION

Large language models (LLMs) [1], [2] have witnessed rapid advancements, demonstrating remarkable capabilities across a broad range of natural language processing tasks. However, these capabilities come with a huge demand for memory and compute, posing significant challenges in resource-constrained settings. Low-bit quantization entails reducing the bit-width of tensors, thereby decreasing memory footprint and easing computational requirements, while preserving generation quality and emergent capabilities such as in-context learning and instruction-following. Quantization methods are commonly grouped into post-training quantization (PTQ) and quantization-aware training (QAT), with the former widely adopted when further training is infeasible.

Existing PTQ frameworks focus on reducing the precision of weights [3], activations [4], and key–value (KV) caches [5]. Although these frameworks often preserve perplexity and accuracy on multi-domain tasks, they typically overlook trustworthiness metrics. This neglect introduces risks when deploying quantized LLMs to downstream applications, potentially leading to unfair, non-robust, or even harmful behaviors.

This work highlights the necessity to consider trustworthiness when compressing LLMs for deployment, using weight quantization as a representative example. We begin by investigating quantized models on both multi-domain tasks and four trustworthiness-focused metrics (adversarial robustness, fairness, machine ethics, and out-of-distribution robustness). Consistent with prior work [6], we find that quantization frameworks typically preserve performance at 8-bit. When further compressed to 3-bit and 4-bit, models often maintain accuracy on multi-domain tasks, but perform divergently across quantization methods and trustworthiness metrics. We further observe that although low-precision models can outperform non-compressed dense models on certain dimensions, their performance is less stable, suffers from high refusal rates, and can exhibit abrupt failures.

To improve the robustness of low-precision models, we introduce a novel precision-ensemble voting approach utilizing multi-precision LLMs. Featuring refusal filtering and majority voting, our approach achieves stable and desirable performance using low-precision models, obtaining superior performance to large dense models by up to 5.8%. Our study underscores the importance of considering trustworthiness metrics under model compression and outlines challenges and opportunities for future research on robust, efficient LLMs.

II. BACKGROUND

A. Low-Precision LLM Inference

Quantization is a widely adopted model compression technique that represents tensors in low-precision number formats, thereby reducing both computational cost and memory

footprint. Quantization has been extensively applied to traditional neural networks (e.g., CNNs/RNNs) [7]–[10], but in the era of transformer-based LLMs [11], self-attention and layer normalization pose new challenges. Contemporary quantization workflows comprise post-training quantization (PTQ) and quantization-aware training (QAT) [12]. PTQ compresses a pre-trained model without retraining and thus introduces minimal overhead. QAT considers quantization error during the training, optimizes parameters for low-bit representations, and typically achieves higher accuracy than PTQ. For LLMs, research has applied quantization to weights, activations, and key–value (KV) caches.

1) *Weight-Only Quantization*: Weight-only quantization applies lower precision to weights while keeping activations in their original precision. In this setting, GPTQ [13] conducts blockwise second-order optimization, adjusting per-weight rounding to efficiently minimize layer-output reconstruction error. AWQ [3] rescales activation-informed channels to preserve salient directions before applying symmetric per-channel weight quantization offline. SqueezeLLM [14] employs a dense–sparse decomposition, isolating outlier channels while quantizing the remaining dense weights more aggressively. AnyPrecisionLLM [15] uses bit-sliced weights to enable runtime-selectable precisions, adapting to diverse hardware budgets and workloads. For ultra-low bitwidths, QTIP and AQLM propose codebook- or entropy-aware schemes to preserve accuracy and stability. Furthermore, PDMD [16] adopts an adaptive decoding strategy that progressively reduces precision as generation proceeds. We focus on how weight-only PTQ frameworks affect model trustworthiness in this work.

2) *Weight-Activation Quantization*: Weight–activation quantization compresses both weights and activations, enabling low-bit matrix multiplications (e.g., W8A8: 8-bit weights and 8-bit activations). In this setting, ZeroQuant [4] first explores weight–activation quantization for LLMs using group-wise weight quantization and token-wise activation quantization to enable W8A8 inference. SmoothQuant [17] migrates activation quantization difficulty to the weights via per-channel rescaling, smoothing activation outliers for training-free W8A8 quantization. RPTQ [18] reorders activation channels into range-homogeneous clusters and fuses the resulting permutations into adjacent layers, enabling robust low-bit activation quantization.

3) *KV Cache Quantization*: In LLM inference, the size of the KV cache scales rapidly as batch size and sequence length increase, motivating the compression of the stored KV pairs. In this setting, KVQuant [5] employs per-channel key quantization, pre-RoPE key quantization, layer-sensitive non-uniform datatypes, and per-vector dense–sparse handling to achieve sub-4-bit KV caches. KIVI [19] provides a tuning-free asymmetric 2-bit scheme that quantizes keys per-channel and values per-token with a streaming- and hardware-friendly implementation. WKVQuant [20] jointly quantizes model weights and the past-only KV cache to low bitwidths, improving attention efficiency while preserving stability.

Type	Framework (bits)
Weight-Only	GPTQ (3–8); AWQ (3–8); QTIP (2–4); AQLM (2–4); SqueezeLLM (2–8)
Weight+Activation	ZeroQuant (W8A8); SmoothQuant (W8A8); RPTQ (W4A16 / W4A8 / W4A4)
KV Cache	KVQuant (2–4); KIVI (2); WKVQuant (4)

TABLE I: Summary of low-precision LLM post-training quantization frameworks with typical bit settings.

B. Trustworthiness of Low-Precision LLMs

LLMs are increasingly integrated into daily applications. However, they pose risks to users, including generating biased content and disclosing sensitive information. These risks are particularly consequential in safety-critical domains such as healthcare [10], [21]–[23] and finance [24]. Numerous studies have evaluated the trustworthiness of LLMs [6], [25], [26]: DecodingTrust [6] evaluates the trustworthiness of LLMs in eight different trustworthiness metrics; TrustLLM [25] provides an open evaluation suite and taxonomy that measure robustness, fairness, privacy, and transparency. As LLMs are adopted in broader applications, trustworthiness research also becomes application-specific. In agentic systems, [27] proposes governance controls and runtime monitors for autonomous agents, addressing provenance, policy compliance, and oversight. In healthcare applications, [28] synthesizes evaluation protocols and safeguards for clinical LLMs, emphasizing reliability and harm mitigation.

Current quantization frameworks focus on evaluating perplexity on pretraining datasets, as well as zero-shot and few-shot accuracy on classification and reasoning tasks. However, these metrics may not fully capture model capabilities such as instruction-following, nor behaviors related to hallucination [29]. In the trustworthiness domain, [30] highlights the impact of model compression and reports that quantization typically results in less degradation across multiple trustworthiness metrics compared with pruning. For specific dimensions, [31] benchmarks the robustness of quantized models in code generation, [32] evaluates the harmfulness of quantized models, and [33] assesses the truthfulness of quantized models. Existing studies [6], [31] further observe that dense model size, the chosen quantization framework, the degree of quantization, and the selected trustworthiness metrics all influence the performance of quantized models. Building on these insights, we investigate opportunities to improve the trustworthiness of low-precision LLMs.

III. BENCHMARKING TRUSTWORTHINESS OF QUANTIZED LLMs

Before investigating how trustworthiness can be enhanced in mixed-precision settings, we first evaluate the effect of quantization on model trustworthiness. In this section, we explore 1) how quantization impacts trustworthiness across multiple dimensions; 2) how trustworthiness varies with the chosen quantization framework; and 3) how the compression ratio influences robustness of the trustworthiness.

A. Models and Quantization Frameworks

We evaluate trustworthiness on LLaMA-2-Chat [34], an instruction-tuned variant of LLaMA-2 optimized for multi-turn dialogue. We consider two dense LLMs, 7B and 13B, as baselines for evaluation. Both configurations are popular for latency-constrained or on-device deployments.

Following prior work [30], [35], we apply AWQ [3] and GPTQ [13], two representative post-training weight-only quantization frameworks to LLaMA-2-13B-Chat. AWQ performs activation-aware weight quantization by estimating channel salience from calibration activations and selecting scales to preserve layer responses. GPTQ formulates weight quantization as a blockwise quadratic reconstruction problem, using an approximate Hessian to compute error-compensated rounding that minimizes output distortion. We adopt three different bit-widths: 8-bit, 4-bit, and 3-bit. Prior work [30] also investigates structured pruning, another popular model compression technique, reporting that quantization constitutes a more reliable method to preserve trustworthiness than pruning at similar compression ratios.

B. Evaluation Metrics

Following the configuration in [30], we leverage three metrics to investigate the questions in Sec. III: 1) *Accuracy on multi-domain tasks*. We report average accuracy on **Massive Multitask Language Understanding (MMLU)** [36], a benchmark covering 57 tasks including elementary mathematics, U.S. history, computer science, and law; 2) *Trustworthiness*. We evaluate four trustworthiness metrics, introduced below; and 3) *Refusal rate*. For each trustworthiness metric, we measure the refusal rate. Refusal rate characterizes the frequency that LLMs refuse to provide an explicit answer, by answering “I don’t know”, or giving neutral responses. In most trustworthiness metrics, refusals contribute to inaccuracy.

Trustworthiness is multifaceted and resists a single quantified definition. Following TrustLLM [25] which describes trustworthiness in LLMs as the accurate representation of information, facts, and results, prior work has proposed taxonomies spanning safety, robustness, fairness, toxicity, reliability, privacy, machine ethics, and explainability [6], [25], [26]. In this work, we evaluate four representative dimensions: adversarial robustness, fairness, machine ethics, and out-of-distribution (OOD) robustness. We leverage the evaluation framework of DecodingTrust [6], and score each dimension on a 0 ~ 100 scale.

1) *Adversarial Robustness*: Adversarial robustness characterizes a model’s stability under adversarially perturbed inputs, using **AdvGLUE** [37] and **AdvGLUE++** [6]. AdvGLUE applies 14 textual adversarial attack methods to GLUE tasks; DecodingTrust further introduces AdvGLUE++, an extension that generates adversarial texts using LLMs. We report accuracy on three representative and challenging tasks: Sentiment Analysis (SST-2), Duplicate Question Detection (QQP), and Multi-Genre Natural Language Inference (MNLI), as well as the refusal rate.

2) *Fairness*: Fairness characterizes the robustness of model predictions with respect to sensitive attributes. To evaluate fairness, we use the Adult dataset [38], which includes attributes such as age, sex, and race to predict whether a person’s income exceeds \$50k per year. The fairness score is aggregated using demographic parity difference (DPD) and equalized odds difference (EOD). DPD captures disparities in the rate of positive predictions between different sensitive attribute values; EOD incorporates ground-truth labels by comparing groups on both true positive rate and false positive rate. For both metrics, larger values indicate greater unfairness, and zero indicates parity.

3) *Machine Ethics*: Machine ethics characterizes common-sense moral judgments aligned with principles that humans intuitively accept. We evaluate using 2,109 short samples from the ETHICS dataset [39] under zero-shot and few-shot settings, strengthened with jailbreaking and evasive prompts. Models are expected to recognize immoral actions and remain robust to evasive wording or jailbreak attempts. The final score represents the model’s ability to identify immoral actions.

4) *Out-of-Distribution Robustness*: Out-of-distribution (OOD) robustness characterizes performance on inputs that deviate substantially from the training distribution. The evaluation covers style-transformed paraphrases and out-of-scope knowledge, assessing whether the model identifies OOD scenarios and attains high accuracy on inputs it does not refuse. The final score aggregates the refusal rate and meaningful accuracy, where meaningful accuracy (MAcc) denotes accuracy conditional on non-refusal.

C. Performance on Multi-domain Tasks and Trustworthiness

We evaluate accuracy on multi-domain tasks and on multiple trustworthiness metrics, as depicted in Fig. 1. Although all evaluations can be cast as classification after post-processing, quantization exhibits a different impact on trustworthiness metrics than on general-task accuracy. We analyze aggregated performance score and refusal rate across all metrics, thereby addressing the three questions in Sec. III.

1) *MMLU*: We evaluate the models on MMLU, which spans 57 subjects in a multiple-choice format to represent the performance in multi-domain tasks. All questions were answered, yielding a zero refusal rate. In terms of accuracy, the change relative to the 13B dense baseline is negligible for both 8-bit and 4-bit quantization. Further compressing the model to 3-bit leads to a noticeable performance drop; however, accuracy for most configurations remains higher than that of the 7B dense baseline. These results indicate that low-precision quantization can largely preserve multi-domain task accuracy compared to dense models and can outperform smaller-parameter dense models. Comparing the two quantization methods, both AWQ and GPTQ preserve accuracy at 8-bit, while AWQ is more robust at 4-bit and 3-bit.

2) *Adversarial Robustness*: Adversarial robustness is assessed on the AdvGLUE++ variants of three GLUE tasks (SST-2, QQP, and MNLI). At 8-bit, both AWQ and GPTQ achieve accuracy close to the 13B dense baseline. Under lower

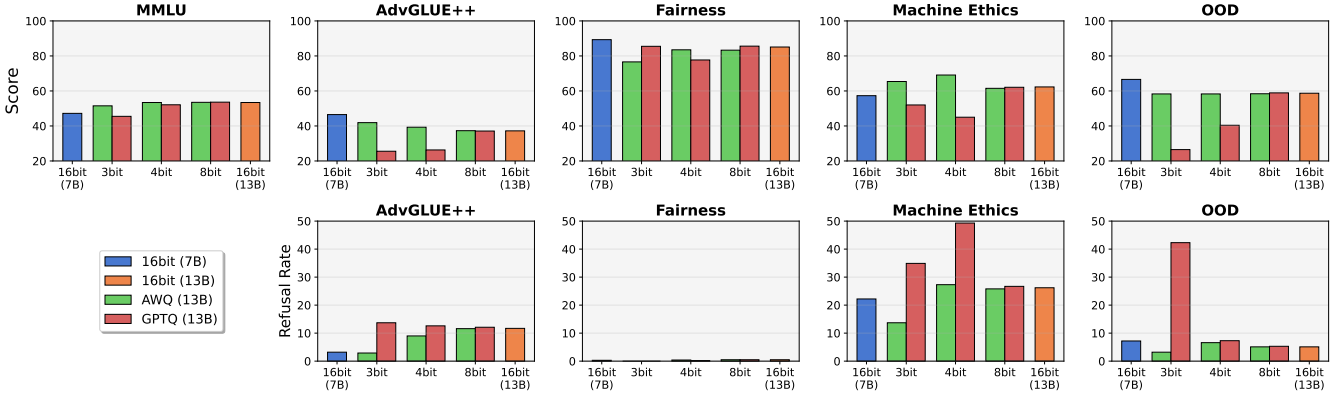


Fig. 1: Accuracy and refusal rate for multi-domain tasks (MMLU) and trustworthiness metrics. For MMLU, the refusal rate is zero and therefore omitted. On multi-domain tasks, AWQ and GPTQ match the 13B dense baseline at 8-bit and remain close at 4-bit, with a larger degradation at 3-bit. On trustworthiness metrics, both methods are comparable to the dense baseline at 8-bit; at lower precisions (4- and 3-bit), AWQ is more robust than GPTQ.

precisions, AWQ exhibits improved robustness at 4-bit and 3-bit, whereas GPTQ degrades by more than 10%, at 4-bit and 3-bit. Accuracy is computed over the valid label set; refusals and outputs that cannot be mapped to task labels are counted as errors. Consistent with this definition, we observe that AWQ low-bit models have low refusal rates, while GPTQ at higher compression yields more label-inconsistent (“contradictory”) answers rather than explicit abstentions, resulting in refusals. This pattern suggests that aggressive GPTQ compression may impair calibration and label consistency.

Comparing dense baselines, the 7B model outperforms the 13B model on AdvGLUE++ by 9.3% and yields a lower refusal rate. One possible explanation is that smaller models are less sensitive to spurious lexical cues introduced by adversarial perturbations, whereas larger models are more easily misled by minor word changes. The observed correlation between lower refusal rates and higher accuracy holds for both dense and quantized variants.

3) *Fairness*: Fairness is assessed on the Adult dataset [38], with scores aggregated from two metrics: demographic parity difference (DPD) and equalized odds difference (EOD). Unlike the other trustworthiness metrics, the fairness score does not completely rely on correctly predicting the category. Both dense and quantized models achieve high fairness, with scores close to or above 80% (higher is better under our normalized fairness score); all configurations also achieve low refusal rates, since the requests are not ambiguous or misleading. For AWQ, the fairness score remains similar at 8-bit and 4-bit, while at 3-bit the performance decreases; GPTQ attains high fairness at 8-bit and 3-bit, with a decrease at 4-bit. Additionally, the breakdown results show that the magnitude of DPD is typically lower than that of EOD, indicating that EOD contributes more to the overall unfairness; for high-fairness configurations, the difference between EOD and DPD is small.

Comparing dense baselines, the 7B model also outperforms the 13B model, as well as the quantized 13B variants. The performance gap is mainly associated with cases showing extreme

base-rate imbalance, where the label distribution is highly skewed. In such cases, larger models appear more sensitive to the imbalance and make more biased classifications.

4) *Machine Ethics*: Machine ethics is assessed on the ETHICS dataset [39], reporting accuracy together with a false positive ratio (FPR) that penalizes false positive responses to jailbreaking and evasive inputs. At 8-bit, both quantization methods achieve performance similar to the 13B dense baseline. At lower bit-widths (4- and 3-bit), AWQ shows improved performance relative to the dense baseline, whereas GPTQ exhibits a significant drop. A breakdown indicates that the true positive rate on the immoral class is similar for all models; the 4- and 3-bit AWQ-quantized models benefit from a lower FPR on evasive sentences, while low-bit GPTQ-quantized models frequently produce invalid or abstaining outputs under few-shot, yielding a higher refusal rate. This highlights the model’s degradation of in-context learning capability.

Comparing dense baselines, the 7B dense model scores lower than the 13B dense model and lags behind the AWQ-quantized variants. Although the 7B model’s accuracy is approximately 10% lower, its lower refusal rate narrows the gap with the 13B dense model.

5) *Out-of-Distribution Robustness*: OOD robustness is assessed on both style-shifted inputs and out-of-scope knowledge queries. For style-shifted inputs (e.g., Shakespearean phrasing), robustness means maintaining high accuracy with a low refusal rate. For out-of-scope queries, robustness requires detecting that the query is beyond the model’s knowledge and abstaining appropriately, while answering in-scope queries correctly, achieving high meaningful accuracy (MAcc). AWQ-quantized models maintain OOD performance relative to the 13B dense baseline, whereas GPTQ-quantized models degrade significantly at 4- and 3-bit. A breakdown attributes the 3-bit degradation primarily to elevated refusal rates, and the 4-bit degradation to low accuracy on out-of-scope knowledge tasks, indicating a failure mode similar to that observed in the machine-ethics dimension.

Comparing dense baselines, the 7B model outperforms the

13B model and the quantized variants on OOD robustness. Specifically, in the zero-shot out-of-scope evaluation, the 7B model achieves nearly double the accuracy while exhibiting a higher refusal rate, indicating a more conservative behavior on knowledge-unknown queries.

6) *Observations and Insights:* The above observations provide insights into the three questions in Sec. III. 1) Quantization methods have heterogeneous impacts on multi-domain tasks and on different trustworthiness metrics. This difference is often missed by standard evaluations of quantized models. 2) For both AWQ and GPTQ, 8-bit quantization largely preserves performance. AWQ is more robust at 4- and 3-bit, whereas low-precision GPTQ can substantially degrade few-shot adherence and in-context learning in the LLaMA-2 series. 3) Smaller models, whether quantized or trained with fewer parameters, can outperform on certain trustworthiness dimensions by being less sensitive to ambiguous phrasing, where larger models are more likely to err. These insights motivate our attempt to enhance quantized models’ robustness in trustworthiness, as described in Sec. IV.

IV. ENHANCING TRUSTWORTHINESS USING PRECISION-ENSEMBLE VOTING

A. Motivation

In Sec. III, we examine the performance of quantized models at different compression ratios on multi-domain tasks and trustworthiness metrics. As shown in Fig. 1, models quantized to low bit-widths can surpass dense models on certain trustworthiness metrics while maintaining comparable accuracy on multi-domain tasks. However, a critical bottleneck of low-precision quantization is reduced robustness: quantized models are more vulnerable to high refusal rates and can suffer abrupt performance drops in specific scenarios. Prior work [30] also reports performance instability across model families and quantization methods. These observations motivate mechanisms that enable low-precision models to provide robust and consistent predictions.

We address this bottleneck by leveraging the idea of test-time optimization [40], [41], which invests additional inference compute to enhance performance. To this end, we propose a simple voting-based precision-ensemble that aggregates the predictions of multi-precision variants quantized from the same dense LLM.

B. Precision-Ensemble Voting

We pursue robust quantized models via **precision-ensemble voting**, illustrated in Fig. 2. The procedure comprises four stages, quantization, generation, filtering, and voting, as detailed in Algorithm 1. In the quantization stage (lines 1–3), the dense backbone (e.g., 13B) is quantized to multiple bit-widths or the same bit-width using different seeds. In the generation stage (lines 4–5), each quantized LLM generates a response in parallel, and predictions are mapped to discrete labels. Many trustworthiness benchmarks can be cast as classification, enabling the use of voting to aggregate results. Labels can be extracted using heuristics or an LLM-as-a-Judge [42].

In the filtering stage (lines 6–9), we discard invalid or contradictory outputs (e.g., empty, unparseable, or multi-label generations) and mark refusals. For standard classification metrics, refusals are removed. For OOD robustness tasks that evaluate out-of-scope queries, explicit refusals (e.g., “I don’t know”) are retained as a dedicated label. In the final voting stage, the remaining candidates are aggregated via unweighted majority voting. In the event of a tie, we select the highest-precision model’s prediction.

Algorithm 1 Precision-ensemble voting

Require: Dense model \mathcal{M}_{fp} , precisions $\mathcal{B} = \{N_1, \dots, N_m\}$, prompt x , decoding params θ , refusal filter f

Ensure: Final label \hat{y}

```

1: for all  $N_i \in \mathcal{B}$  do
2:    $\mathcal{M}_i \leftarrow \text{Quantize}(\mathcal{M}_{fp}, N_i)$ 
3: end for
4: Generate:  $G \leftarrow \{\text{Generate}(\mathcal{M}_i, x; \theta)\}_{i=1}^m$ 
5: Map to labels:  $C \leftarrow \{\text{Postprocess}(g) \mid g \in G\}$ 
6: Filter refused candidates:  $C' \leftarrow \{c \in C \mid \neg f(c)\}$ 
7: if  $C' = \emptyset$  then
8:   return REFUSED
9: end if
10: Aggregation:  $\hat{y} \leftarrow \text{Maj.Voting}(C')$ 
11: return  $\hat{y}$ 

```

C. Effectiveness of Precision-Ensemble Voting

1) *Evaluation Setup:* We quantize the dense LLaMA-2-13B-Chat model to three precisions: 3-, 4-, and 8-bit. Each quantized model and the precision-ensemble are evaluated on multi-domain tasks (MMLU) and on three trustworthiness metrics: adversarial robustness, machine ethics, and out-of-distribution robustness. The fairness dimension is excluded in this evaluation because it quantifies inter-group disparities between predictions rather than absolute accuracy, and is therefore not directly comparable under this setup. For each metric, we compare the ensemble prediction against the 7B and 13B dense baselines and against the best-performing single-precision result for each quantization method.

2) *Effectiveness on Multi-Domain Tasks:* According to Fig. 3, precision-ensemble voting achieves accuracy that is close to both the best single-precision results and the 13B dense baseline, and it outperforms the 7B dense model. The gains on MMLU are modest because low-bit quantization already preserves accuracy for these models. Nevertheless, the ensemble effectively mitigates instability observed in low-precision GPTQ-quantized models.

3) *Effectiveness on Trustworthiness Metrics:* Fig. 3 shows that precision-ensemble voting delivers robust, desirable performance. Compared with dense baselines, the ensemble results consistently outperform the 13B dense model by up to 5.8%, but fail to match the 7B dense baseline, as small models often benefit from being less sensitive to spurious lexical cues. Relative to the best single-precision results, the

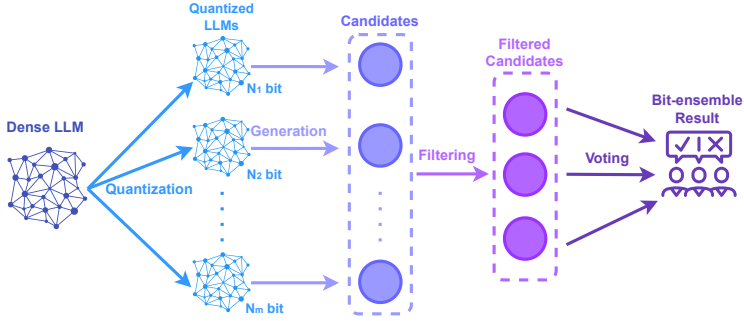


Fig. 2: Workflow of precision-ensemble voting. A dense LLM is quantized to multiple precisions; each quantized model generates its own response. After response filtering, the remaining responses are aggregated via unweighted majority voting.

precision-ensemble voting approach performs better in most trustworthiness metrics and for both quantization frameworks.

We attribute these gains to two factors: 1) the filtering stage (for standard classification metrics) removes refusals, thereby lowering the measured refusal rate and preserving meaningful labels for aggregation; and 2) unweighted majority voting improves stability by reducing variance across bit-widths: when a quantized model fails on a specific instance, other bit-widths do not tend to fail as well, and the aggregation exploits this partial error diversity.

V. CHALLENGES AND OPPORTUNITIES

A. Mixed Precision for Multi-Modal Trustworthiness

Challenges: The rapid development and increasing deployment of multi-modal LLMs, integrating vision, speech, action, and language, introduces a more complex trustworthiness problem than in text-only models. In particular, embodied AI systems, which rely on cross-modal reasoning to interact with the real world, demand heightened attention to trustworthiness across all modalities.

Opportunities: Multi-modal setting opens new opportunities for modality-aware mixed-precision quantization that jointly considers efficiency and trustworthiness. Different bit-ensemble strategies and precision scheduling can be adaptively tuned based on modality-specific information collected at runtime. Our findings in the single-modality case lay the groundwork for future exploration in multi-modal scenarios.

B. Joint Compression and Trust-Aware Optimization

Challenges: While mixed-precision quantization has proven effective in preserving or even enhancing trustworthiness, real-world deployments often combine it with other compression techniques such as pruning, tensor decomposition, and knowledge distillation. However, the interplay between these methods and their joint effect on trustworthiness remains underexplored. In particular, naive combinations may introduce unpredictable degradation in ethical behavior, adversarial robustness, or fairness, especially under low-bit regimes.

Opportunities: We identify a significant opportunity to unify compression design with trustworthiness objectives. One

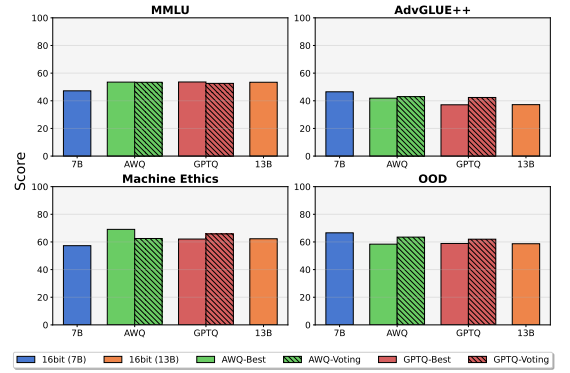


Fig. 3: The precision-ensemble voting mechanism maintains MMLU accuracy while consistently improving performance on the trustworthiness benchmarks.

promising direction is to formulate model compression as an automated sparsity search problem, where trustworthiness metrics are directly embedded in the optimization objective. This would enable the development of trust-aware auto-compression pipelines capable of jointly tuning bit-width, sparsity, and decomposition rank for maximal efficiency while enhancing trustworthiness.

C. Algorithm-System-Hardware Co-Design

Challenges: Efficient system and hardware support for mixed-precision and bit-ensemble execution is critical for real-world deployment. At the hardware level, designing a unified compute unit that can efficiently and concurrently support operations at multiple precisions remains a major challenge, due to inherent trade-offs among performance, power, and area (PPA). At the system level, effective scheduling and execution policies are required to manage precision diversity, particularly in multi-batch and multi-tenant deployment settings.

Opportunities: The hardware and system design choices span a large configuration space. When jointly considered with algorithmic parameters, this forms a vast co-design space where algorithmic performance and hardware/system efficiency can be optimized together [43]. This joint optimization can be tackled using techniques such as Bayesian optimization or reinforcement learning to efficiently explore the co-design space and derive Pareto-optimal solutions.

VI. CONCLUSION

This work analyzes the impact of quantization on trustworthiness, highlights the stability bottleneck with low precision, and proposes a precision-ensemble voting approach to improve robustness. We underscore the importance of considering trustworthiness when compressing LLMs for deployment, using weight quantization as a representative case study with quantitative experimental results. Future directions include extending trustworthiness analysis to multi-modal settings, exploring joint compression and trust-aware optimization frameworks, and pursuing end-to-end algorithm–system–hardware co-design for trustworthy and efficient LLM deployment.

REFERENCES

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [2] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv *et al.*, “Qwen3 technical report,” *arXiv preprint arXiv:2505.09388*, 2025.
- [3] J. Lin, J. Tang, H. Tang, S. Yang, W.-M. Chen, W.-C. Wang, G. Xiao, X. Dang, C. Gan, and S. Han, “Awq: Activation-aware weight quantization for on-device llm compression and acceleration,” *Proceedings of machine learning and systems*, vol. 6, pp. 87–100, 2024.
- [4] Z. Yao, R. Y. Aminabadi, M. Zhang, X. Wu, C. Li, and Y. He, “Zeroquant: Efficient and affordable post-training quantization for large-scale transformers,”
- [5] C. Hooper, S. Kim, H. Mohammadzadeh, M. W. Mahoney, Y. S. Shao, K. Keutzer, and A. Gholami, “Kvquant: Towards 10 million context length llm inference with kv cache quantization,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 1270–1303, 2024.
- [6] B. Wang, W. Chen, H. Pei, C. Xie, M. Kang, C. Zhang, C. Xu, Z. Xiong, R. Dutta, R. Schaeffer *et al.*, “Decodingtrust: A comprehensive assessment of trustworthiness in gpt models,” 2023.
- [7] M. Nagel, M. Fournarakis, R. A. Amjad, Y. Bondarenko, M. Van Baalen, and T. Blankevoort, “A white paper on neural network quantization,” *arXiv preprint arXiv:2106.08295*, 2021.
- [8] H. Fan, H.-C. Ng, S. Liu, Z. Que, X. Niu, and W. Luk, “Reconfigurable acceleration of 3d-cnns for human action recognition with block floating-point representation,” in *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2018, pp. 287–2877.
- [9] H. Fan, G. Wang, M. Ferianc, X. Niu, and W. Luk, “Static block floating-point quantization for convolutional neural networks on fpga,” in *2019 International Conference on Field-Programmable Technology (ICFPT)*. IEEE, 2019, pp. 28–35.
- [10] H. Fan, S. Liu, Z. Que, X. Niu, and W. Luk, “High-performance acceleration of 2-d and 3-d cnns on fpgas using static block floating point,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 8, pp. 4473–4487, 2021.
- [11] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [12] A. Gholami, S. Kim, Z. Dong, Z. Yao, M. W. Mahoney, and K. Keutzer, “A survey of quantization methods for efficient neural network inference,” in *Low-power computer vision*. Chapman and Hall/CRC, 2022, pp. 291–326.
- [13] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, “Gptq: Accurate post-training quantization for generative pre-trained transformers,” *arXiv preprint arXiv:2210.17323*, 2022.
- [14] S. Kim, C. Hooper, A. Gholami, Z. Dong, X. Li, S. Shen, M. W. Mahoney, and K. Keutzer, “Squeezellm: dense-and-sparse quantization,” in *Proceedings of the 41st International Conference on Machine Learning*, 2024, pp. 23 901–23 923.
- [15] Y. Park, J. Hyun, S. Cho, B. Sim, and J. W. Lee, “Any-precision llm: low-cost deployment of multiple, different-sized llms,” in *Proceedings of the 41st International Conference on Machine Learning*, 2024, pp. 39 682–39 701.
- [16] H. M. Chen, F. Tan, A. Kouris, R. Lee, H. Fan, and S. Venieris, “Progressive mixed-precision decoding for efficient llm inference,” in *The Thirteenth International Conference on Learning Representations*.
- [17] G. Xiao, J. Lin, M. Seznec, H. Wu, J. Demouth, and S. Han, “Smoothquant: accurate and efficient post-training quantization for large language models,” in *Proceedings of the 40th International Conference on Machine Learning*, 2023, pp. 38 087–38 099.
- [18] Z. Yuan, L. Niu, J. Liu, W. Liu, X. Wang, Y. Shang, G. Sun, Q. Wu, J. Wu, and B. Wu, “Rptq: Reorder-based post-training quantization for large language models,” *arXiv preprint arXiv:2304.01089*, 2023.
- [19] Z. Liu, J. Yuan, H. Jin, S. Zhong, Z. Xu, V. Braverman, B. Chen, and X. Hu, “Kivi: a tuning-free asymmetric 2bit quantization for kv cache,” in *Proceedings of the 41st International Conference on Machine Learning*, 2024, pp. 32 332–32 344.
- [20] Y. Yue, Z. Yuan, H. Duanmu, S. Zhou, J. Wu, and L. Nie, “Wkvquant: Quantizing weight and key/value cache for large language models gains more,” *CoRR*, 2024.
- [21] Asgari *et al.*, “A framework to assess clinical safety and hallucination rates of llms for medical text summarisation,” *npj Digital Medicine*, vol. 8, no. 1, p. 274, 2025.
- [22] H. Fan, M. Chen, L. Castelli, Z. Que, H. Li, K. Long, and W. Luk, “When monte-carlo dropout meets multi-exit: Optimizing bayesian neural networks on fpga,” in *2023 60th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2023, pp. 1–6.
- [23] H. Fan, M. Ferianc, and W. Luk, “Enabling fast uncertainty estimation: accelerating bayesian transformers via algorithmic and hardware optimizations,” in *Proceedings of the 59th ACM/IEEE Design Automation Conference*, 2022, pp. 325–330.
- [24] T. Hu, T. Hu, L. Bai, Y. Zhao, A. Cohan, and C. Zhao, “Fintrust: A comprehensive benchmark of trustworthiness evaluation in finance domain,” in *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 2025, pp. 10 110–10 139.
- [25] Y. Huang *et al.*, “Trustllm: Trustworthiness in large language models,” 2024. [Online]. Available: <https://openreview.net/forum?id=bWUULwwMp>
- [26] Y. Liu *et al.*, “Trustworthy llms: a survey and guideline for evaluating large language models’ alignment,” *arXiv preprint arXiv:2308.05374*, 2023.
- [27] S. Raza, R. Sapkota, M. Karkee, and C. Emmanouilidis, “Trism for agentic ai: A review of trust, risk, and security management in llm-based agentic multi-agent systems,” *arXiv preprint arXiv:2506.04133*, 2025.
- [28] M. Aljohani, J. Hou, S. Kommu, and X. Wang, “A comprehensive survey on the trustworthiness of large language models in healthcare,” *arXiv preprint arXiv:2502.15871*, 2025.
- [29] J. Lee, S. Park, J. Kwon, J. Oh, and Y. Kwon, “Exploring the trade-offs: Quantization methods, task difficulty, and model size in large language models from edge to giant,” *arXiv preprint arXiv:2409.11055*, 2024.
- [30] J. Hong, J. Duan, C. Zhang, Z. Li, C. Xie, K. Lieberman *et al.*, “Decoding compressed trust: Scrutinizing the trustworthiness of efficient llms under compression,” *arXiv:2403.15447*, 2024.
- [31] S. Fang, W. Ding, A. Mastropaolo, and B. Xu, “Smaller= weaker? benchmarking robustness of quantized llms in code generation,” *arXiv preprint arXiv:2506.22776*, 2025.
- [32] Y. Belkhirer, G. Zizzo, and S. Maffei, “Harmlevelbench: Evaluating harm-level compliance and the impact of quantization on model alignment,” *arXiv preprint arXiv:2411.06835*, 2024.
- [33] Y. Fu, X. Long, R. Li, H. Yu, M. Sheng, X. Han, Y. Yin, and P. Li, “Quantized but deceptive? a multi-dimensional truthfulness evaluation of quantized llms,” *arXiv preprint arXiv:2508.19432*, 2025.
- [34] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [35] A. Kharinaev, V. Moskvoretiskii, E. Shvetsov, K. Studenikina, B. Mikhail, and E. Burnaev, “Investigating the impact of quantization methods on the safety and reliability of large language models,” *arXiv preprint arXiv:2502.15799*, 2025.
- [36] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, “Measuring massive multitask language understanding,” *arXiv preprint arXiv:2009.03300*, 2020.
- [37] B. Wang, C. Xu, S. Wang, Z. Gan, Y. Cheng, J. Gao *et al.*, “Adversarial glue: A multi-task benchmark for robustness evaluation of language models,” in *Advances in Neural Information Processing Systems*, 2021.
- [38] B. Becker and R. Kohavi, “Adult,” UCI Machine Learning Repository, 1996, DOI: <https://doi.org/10.24432/C5XW20>.
- [39] D. Hendrycks, C. Burns, S. Basart, A. Critch, J. Li, D. Song, and J. Steinhardt, “Aligning ai with shared human values,” in *International Conference on Learning Representations*.
- [40] C. Snell, J. Lee, K. Xu, and A. Kumar, “Scaling llm test-time compute optimally can be more effective than scaling model parameters,” *arXiv preprint arXiv:2408.03314*, 2024.
- [41] H. M. Chen, G. Lu, Y. Okoshi, Z. Mo, M. Motomura, and H. Fan, “Rethinking optimal verification granularity for compute-efficient test-time scaling,” *arXiv preprint arXiv:2505.11730*, 2025.
- [42] D. Li *et al.*, “From generation to judgment: Opportunities and challenges of llm-as-a-judge,” in *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 2025, pp. 2757–2791.
- [43] H. Fan, M. Ferianc, Z. Que, S. Liu, X. Niu, M. R. D. Rodrigues, and W. Luk, “Fpga-based acceleration for bayesian convolutional neural networks,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 41, no. 12, pp. 5343–5356, 2022.