

Quantization-Aware Training: A Comprehensive Survey

Anonymous ACL submission

Abstract

With the unprecedentedly rapid development of LLMs, the prohibitive memory footprints and intensive computational demands of models have emerged as critical bottlenecks. Quantization-Aware Training (QAT) techniques have emerged as a primary solution to address these challenges by explicitly simulating quantization effects within the training loop, enabling low-bit models to achieve accuracy comparable to their full-precision counterparts. In this work, we strive to provide a comprehensive survey on QAT, serving as a valuable resource for researchers aiming to understand the theory of QAT and its evolving implementation landscape. To the best of our knowledge, this is the first systematic survey dedicated to reviewing the recent developments of QAT. In this paper, we systematically review existing QAT methodologies based on a taxonomy organized by quantization targets. We provide an in-depth analysis of the technical connections and distinctions among these methods and summarize their evaluation paradigms. Furthermore, we discuss persistent challenges and outline potential directions for future research.

1 Introduction

Recent years have witnessed rapid progress in large-scale foundation models across language, vision, and generative modeling, including large language models (LLMs), vision transformers (ViTs), and diffusion-based generative models. However, as the scale of these models grows exponentially, their demands on memory and computational resources far exceed the capacity of most edge devices and even consumer-grade GPUs (Nagel et al., 2021; Wei et al., 2024; Yang et al., 2025; Liu et al., 2025a). To mitigate these resource bottlenecks, model quantization serves as a fundamental technique, compressing the numerical representation of these massive architectures by lowering the bit-width of their

parameters (Jacob et al., 2018; Nagel et al., 2021; Wei et al., 2024).

Model quantization methods can be broadly divided into Post-training quantization (PTQ) and Quantization-aware training (QAT) (Nagel et al., 2021; Wei et al., 2024; Yang et al., 2025). PTQ applies low-precision quantization to a pre-trained language model without requiring further training, thereby reducing model size while improving inference efficiency (Xiao et al., 2023; Liu et al., 2024b; Zhang and Shrivastava, 2024). However, directly quantizing the model can result in a significant performance degradation due to numerical overflow and limited representational range, which has been repeatedly observed across numerous quantization methods, making researchers cautious about directly applying PTQ. For instance, Bondarenko et al. (2021) identified that activations in residual connections contain structured outliers with high dynamic ranges, which are difficult to represent in low-bit formats, thereby causing significant quantization error.

In contrast to PTQ, QAT simulates quantization during training using fake quantization modules, optimizing the model directly under quantized inference constraints and often recovering the accuracy loss of PTQ, especially at ultra-low bit-widths such as INT4 (Jacob et al., 2018; Bengio et al., 2013). When combined with other parameter-efficient adaptation techniques such as low-rank adaptation (Ke et al., 2024; Xu et al., 2024) and knowledge distillation (Hinton et al., 2015; Kim et al., 2022), QAT can further improve ultra-low precision training and quantized performance. Specifically, Xu et al. (2024) proposed QALoRA to integrate group-wise quantization into low-rank adapters, effectively closing the accuracy gap between PTQ and QAT with minimal overhead. Furthermore, Chen et al. (2025a) demonstrated that such efficient QAT methods can successfully retain the performance of LLMs even in ultra-low

084 precision (e.g., 2-bit) settings.

085 Despite the growing interest and advancements,
086 there is a lack of comprehensive surveys that thor-
087 oughly examine the various methods, challenges,
088 and applications specific to QAT. Most existing
089 surveys on neural network quantization adopt a
090 broad perspective, and consequently provide only a
091 coarse-grained analysis of QAT. For instance, [Wei
092 et al. \(2024\)](#) present a broad taxonomy of neural
093 network quantization, but their discussion of QAT
094 remains largely introductory and does not offer
095 a detailed analysis of key challenges. Likewise,
096 [Nagel et al. \(2021\)](#) focus on deployment-oriented
097 industrial pipelines, with limited theoretical or an-
098 alytical depth of QAT. More recently, [Yang et al.
099 \(2025\)](#) organize quantization methods primarily by
100 target architectures, but this coarse-grained categor-
101 ization offers limited insight into module-specific
102 QAT strategies. In contrast, this survey focuses pri-
103 marily on QAT and organizes the literature accord-
104 ingly. This perspective highlights how the numeri-
105 cal characteristics and access patterns of different
106 targets shape the design of QAT schemes, and why
107 the benefits and costs of QAT differ across targets.

108 We organize recent progress in QAT from four
109 aspects: theoretical foundations, quantization tar-
110 gets, evaluation paradigms, and future horizons.
111 Firstly, we establish the theoretical foundations,
112 focusing on the mathematical formulation of uni-
113 form quantization, the fake quantization mecha-
114 nism, and gradient approximation techniques such
115 as the Straight-Through Estimator (STE). As QAT
116 adapts to diverse model architectures, we explore
117 the quantization targets involved in compressing
118 them, emphasizing not only standard weights and
119 activations but also emerging bottlenecks like the
120 KV cache and gradients. We provide a comprehen-
121 sive overview of evaluation methods, distinguish-
122 ing between structural metrics and functional met-
123 rics. Lastly, we discuss persistent limitations such
124 as optimization instability and explore potential
125 directions for future research.

126 The remainder of this survey is organized as fol-
127 lows: Section 2 establishes the preliminary knowl-
128 edge and theoretical framework. Section 3 cat-
129 egorizes existing research based on quantization
130 targets. Finally, Section 4 identifies key challenges
131 and outlines future opportunities, while evaluation
132 paradigms and metrics are detailed in Appendix A.
133 Figure 1 provides a chronological overview of rep-
134 resentative QAT methods (2020–2025), comple-
135 menting our target-based taxonomy in Section 3.

2 Preliminary Knowledge 136

137 This section introduces the notation and mathemati-
138 cal formulation of uniform affine quantization, and
139 summarizes the core mechanism of QAT ([Jacob
140 et al., 2018; Nagel et al., 2021](#)). We focus on the
141 quantize–dequantize (fake quantization) operator
142 used during training and the surrogate gradients
143 (e.g., STE) needed to backpropagate through round-
144 ing and clipping.

2.1 Uniform Quantization 145

146 Quantization maps high-precision real-valued ten-
147 sors (typically FP32) to a low-precision discrete
148 set amenable to efficient integer arithmetic and re-
149 duced memory ([Jacob et al., 2018; Nagel et al.,
150 2021](#)). Different precisions (bit-widths b) offer a
151 trade-off between model efficiency and representa-
152 tional capacity: while FP32 provides high dynamic
153 range and precision, lower bit-widths (e.g., INT8,
154 INT4) significantly reduce memory footprint and
155 computational cost, yet at the risk of information
156 loss. Uniform quantization is a de facto standard
157 due to its simplicity and compatibility with integer
158 arithmetic ([Jacob et al., 2018; Nagel et al., 2021](#)).

159 Let $x \in \mathbb{R}^{d_1 \times \dots \times d_k}$ denote a real-valued tensor
160 (e.g., weights or activations). Uniform quantization
161 is parameterized by a positive step size $s > 0$ and
162 an integer zero-point $z \in \mathbb{Z}$. The step size sets the
163 grid resolution, and the zero-point makes real zero
164 exactly representable in the integer domain.

165 Let b be the bit-width and let $[q_{\min}, q_{\max}]$ be the
166 integer representable range (e.g., $q_{\min} = -2^{b-1}$
167 and $q_{\max} = 2^{b-1} - 1$ for signed integers). The
168 quantization operator first maps x to an integer
169 tensor $x_{\text{int}} \in \mathbb{Z}^{d_1 \times \dots \times d_k}$:

$$x_{\text{int}} = \text{clamp}\left(\left\lfloor \frac{x}{s} \right\rfloor + z, q_{\min}, q_{\max}\right), \quad (1) \quad 170$$

171 where $\lfloor \cdot \rfloor$ denotes deterministic rounding-to-
172 nearest, and $\text{clamp}(\cdot, q_{\min}, q_{\max})$ clips values to
173 the integer range ([Jacob et al., 2018](#)).

174 During training, QAT typically performs compu-
175 tations in floating point while emulating quantiza-
176 tion effects via a quantize–dequantize (fake quanti-
177 zation) step ([Jacob et al., 2018](#)):

$$\hat{x} = s \cdot (x_{\text{int}} - z), \quad (2) \quad 178$$

179 where \hat{x} is the fake-quantized proxy of x used in the
180 forward pass. The quantization error (or injected
181 noise) is $\epsilon = x - \hat{x}$.

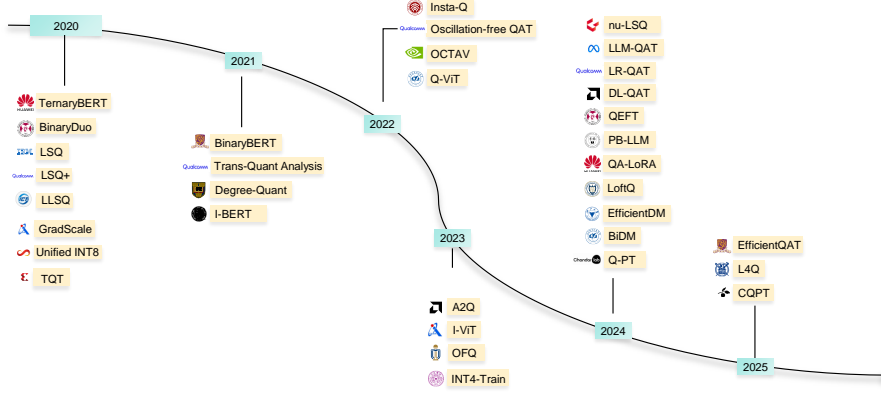


Figure 1: Timeline of representative QAT methods from 2020 to 2025. Works are grouped by publication year along the timeline to provide a high-level view of the field’s evolution (non-exhaustive).

A common choice is to derive (s, z) from a chosen clipping range $[x_{\min}, x_{\max}]$:

$$s = \frac{x_{\max} - x_{\min}}{q_{\max} - q_{\min}}, \quad z = \left[q_{\min} - \frac{x_{\min}}{s} \right]. \quad (3)$$

In practice, $[x_{\min}, x_{\max}]$ can be estimated from calibration statistics (e.g., min–max) (Nagel et al., 2021), optimized by minimizing local reconstruction error (Sakr et al., 2022), or learned jointly with model parameters (Esser et al., 2020).

Symmetric vs. Asymmetric Quantization. Uniform affine quantization is commonly instantiated in two forms: *Asymmetric quantization* uses a nonzero z (Eq. 3), which is beneficial for non-zero-centered or non-negative activations (e.g., ReLU). It reduces wasted quantization levels on skewed ranges, at the cost of extra zero-point terms in integer kernels (Jacob et al., 2018; Nagel et al., 2021). *Symmetric quantization* typically enforces $z = 0$ and uses a symmetric clipping range (e.g., $x_{\min} = -x_{\max}$):

$$\hat{x} = s \cdot x_{\text{int}}. \quad (4)$$

Consequently, many deployment pipelines apply symmetric quantization to network weights to simplify integer kernels, while network activations are often quantized asymmetrically due to their skewed ranges.

2.2 QAT Mechanism

Unlike PTQ, QAT incorporates quantization effects during training by inserting quantization operators into the training graph (Jacob et al., 2018; Nagel et al., 2021). Specifically, QAT inserts fake quantization operators into the forward pass so that the training process can adapt the network parameters to minimize the task loss under discretization. This narrows down the training–inference discrepancy, resulting in models whose predictions are robust to the low-precision scenarios during inference.

Fake Quantization. In QAT, quantization is simulated by replacing some tensors with their fake-quantized counterparts. For a layer that computes $y = f(Wx + b)$, QAT computes

$$\begin{aligned} \hat{W} &= \mathcal{Q}(W), \quad \hat{x} = \mathcal{Q}(x), \\ y &= f(\hat{W} \hat{x} + b), \end{aligned} \quad (5)$$

where $\mathcal{Q}(\cdot)$ denotes the quantize–dequantize operator (Eqs. 1–2). Model parameters are often not stored as integers during training; instead, QAT maintains full-precision master weights W and uses \hat{W} only for forward computation. This preserves small gradient updates that would otherwise vanish after discretization and stabilizes training.

Gradient Mismatch. A critical challenge is that Eq. 1 involves a non-differentiable rounding operator. Formally, the derivative of rounding is zero

almost everywhere and undefined at integers:

$$\frac{\partial \lfloor u \rfloor}{\partial u} = 0 \quad \text{a.e.} \quad (6)$$

Therefore, naive backpropagation through the quantization path yields vanishing or ill-defined gradients, preventing effective optimization of W (Nagel et al., 2022).

Gradient Approximation. To address the gradient mismatch issue, QAT relies on surrogate gradients that provide a usable gradient through rounding and clipping. Straight-Through Estimator (STE) Bengio et al. (2013) is commonly used, treating rounding as an identity operation in the backward pass (optionally combined with clipping). A common formulation is

$$\frac{\partial \mathcal{L}}{\partial x} \approx \frac{\partial \mathcal{L}}{\partial \hat{x}} \cdot \mathbb{I}(x \in [x_{\min}, x_{\max}]), \quad (7)$$

where $[x_{\min}, x_{\max}]$ is the chosen clipping range and $\mathbb{I}(\cdot)$ is the indicator function. Intuitively, STE allows gradients to pass through quantization as if $\hat{x} \approx x$ within the valid range, while saturating gradients when inputs are clipped.

Optimizing Quantization Parameters. Beyond model weights, some QAT methods also learn quantization parameters such as step sizes, offsets (zero-points, Bhalgat et al. (2020)), and clipping thresholds (Jain et al., 2020). LSQ (Esser et al., 2020) treats s c.f. Eqs. 1 as a learnable variable and updates it via gradient descent using an STE-style surrogate. For example, consider the symmetric quantization case ($z = 0$, c.f. Eq. 4) with

$$\hat{x} = s \cdot \text{clamp}\left(\left\lfloor \frac{x}{s} \right\rfloor, q_{\min}, q_{\max}\right). \quad (8)$$

According to LSQ, the gradient w.r.t. s can be written piecewise as

$$\frac{\partial \hat{x}}{\partial s} = \begin{cases} \text{clamp}\left(\left\lfloor \frac{x}{s} \right\rfloor, q_{\min}, q_{\max}\right) - \frac{x}{s}, & q_{\min} < \frac{x}{s} < q_{\max} \\ q_{\min}, & \frac{x}{s} \leq q_{\min}, \\ q_{\max}, & \frac{x}{s} \geq q_{\max}. \end{cases} \quad (9)$$

Learning s enables the model to balance clipping-induced distortion and rounding-induced error in a task-driven manner. In practice, QAT could further incorporate distillation or regularization terms to improve stability under aggressive low-bit settings (Hinton et al., 2015; Kim et al., 2022).

3 Taxonomy

In this section, we firstly structure the QAT literature by the *quantization target*, e.g., weights, activations, KV cache, and gradients. While many methods quantize multiple targets simultaneously, we group each work under the targets where its key ideas, algorithmic designs, and ablations are centered. A taxonomy of representative QAT frameworks under this view is summarized in Table 1.

3.1 Quantization Target: Model Weights

This section surveys works in which the QAT method primarily quantizes the model weights. However, it should be noted that QAT methods often simultaneously quantize the activations — often at a fixed precision such as 8-bit — as an auxiliary design choice. The goal of weight-centric QAT is to compress model parameters while keeping activations in high precision. The key obstacles include unstable optimization under STE, limited expressivity of low-bit grids, and deployment constraints such as accumulator overflow and supported granularities. We summarize recent progress along two themes covering PEFT integration and ultra-low bit.

3.1.1 PEFT Integration

Conventional end-to-end QAT quantizes all parameters, which is infeasible for billion-scale models. Recent work therefore focuses on parameter-efficient weight-QAT, where quantization is restrictively applied to structured or low-rank subspaces. One line of research reduces training memory by updating only quantization-related or auxiliary parameters. EfficientQAT (Chen et al., 2025a) progressively learns only quantization parameters, while DL-QAT (Ke et al., 2024) formulates QAT as low-rank adaptation in a decomposed quantization space, enabling effective ultra-low-bit adaptation. Another line of work integrates QAT with low-rank adaptation, inspired by LoRA (Hu et al., 2022). L4Q (Jeon et al., 2025) achieves LoRA-like training cost with a fully quantized final model, whereas QA-LoRA (Xu et al., 2024) increases quantization flexibility via group-wise operators to enable lossless adapter merging.

Initialization and precision allocation are also crucial. LoftQ (Li et al., 2024b) reduces quantization error through joint backbone quantization and low-rank initialization, while QEFT (Lee et al., 2024) introduces sensitivity-aware mixed precision with hardware-friendly reordering.

322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371

3.1.2 Ultra-Low Bit

Pushing the precision of model weights below 3 bits often leads to severe optimization instability and capacity degradation (Bai et al., 2021; Zhang et al., 2020). Prior work mitigates these issues through sensitivity-aware precision allocation (Shen et al., 2019) and strong supervision, e.g., knowledge distillation and progressive training (Bai et al., 2021). For decoder-only LLMs, preserving a small subset of salient weights becomes particularly important in ultra-low-bit regimes. PB-LLM (Yuan et al., 2024) shows that full binarization harms reasoning ability and instead advocates partial binarization, where salient weights remain in higher precision. Similarly, Nielsen et al. (2025) demonstrates that gradually transitioning from high to ultra-low precision, together with optimizer-state preservation, can alleviate loss spikes under extreme quantization.

Overall, ultra-low-bit weight-centric QAT can reduce model memory footprint by an order of magnitude, enabling integer-only inference and on-device deployment. A key takeaway from these works is that success in ultra-low-bit regimes depends less on increasingly sophisticated low-bit optimizers (Xu et al., 2025) and more on selectively preserving critical weights and stabilizing training with strong supervision.

3.2 Quantization Target: Activations

Activation-centric QAT maps intermediate features (e.g., token embeddings and residual streams) from high-precision floating-point formats to low-bit representations during the forward pass. Activation-centric quantization enables integer-only inference, reduces memory overhead in long-sequence processing, and improves energy efficiency on resource-constrained hardware. However, activation QAT faces several challenges: (i) structural outliers that dominate dynamic ranges, (ii) gradient instability caused by discrete mappings, and (iii) non-stationary activation statistics, particularly in diffusion models.

To mitigate the performance degradation caused by discretizing activations, recent work in activation QAT focuses on introducing learnable quantization parameters and enhanced gradient estimation tailored for activation distributions. LSQ+ (Bhalgat et al., 2020) introduces asymmetric activation quantization with learnable offsets to better fit asymmetric activation distributions. Under ex-

treme compression such as binarization, gradient mismatch becomes the primary bottleneck. BinaryDuo (Kim et al., 2020) formalizes this issue via Coordinate Discrete Gradient (CDG) and mitigates it through training–inference decoupling: ternary activations are used during training to stabilize gradients, while binary activations are adopted at inference to maintain deployment efficiency.

3.3 Quantization Target: Weight + Activation

Joint weight and activation quantization aligns training with the coupled distortion of the genuine inference graph used in practical deployment, enabling weights, activations, and operators to co-adapt to the quantization noise. However, it introduces coupled challenges including non-stationary activation distributions, global range dependencies, and severe gradient mismatch at ultra-low bit widths, alongside the need for integer-friendly operator redesigns. We organized existing work into four directions: representational integer Transformers, dynamic adaptation, quantized generators.

3.3.1 Quantization Range

A key challenge in QAT is accurately modeling the quantization range of activations. Under joint quantization, upstream quantization continuously alters activation distributions, causing the effective dynamic range to drift during training. Consequently, static, layer-wise range calibration is often inadequate. To address this issue, many methods treat the quantization range as a learnable parameter and optimize it jointly with model weights. The range can be adapted in different forms. Some approaches adjust it implicitly by learning non-uniform step sizes to better accommodate heavy-tailed or multimodal activations (Gongyo et al., 2024). Others explicitly constrain the range through symmetric linear quantization, as in LLSQ (Zhao et al., 2020), to maintain compatibility with integer-only hardware. The quantization range can also be controlled via learnable clipping thresholds. TQT (Jain et al., 2020) directly optimizes clipping boundaries, while OCTAV (Sakr et al., 2022) further incorporates gradient-aware differentiation to better handle forward–backward interactions in low-precision training.

3.3.2 Integer Transformers

Research on Integer Transformers (Kim et al., 2021) views Transformer quantization as a system-level problem, where numerical fragility arises

372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420

not only from linear projections but also from core operators such as softmax, normalization, and smooth nonlinearities. Under strict integer arithmetic constraints, joint quantization therefore requires operator-aware approximations and careful scale management.

Representative Integer Transformer approaches, including IBERT (Kim et al., 2021) and IViT (Li and Gu, 2023), pursue integer-only inference pipelines by reformulating or approximating non-integer-friendly operations and explicitly controlling intermediate activation scales to maintain stability across Transformer blocks. Subsequent studies further reveal that, within this integer computation regime, weight-only quantization is insufficient: the dominant quantization errors and memory costs often originate from the activations inside attention and feed-forward sublayers. Accordingly, methods such as QViT (Li et al., 2022) emphasize joint quantization of weights and activations as a prerequisite for preserving accuracy in Integer Transformers.

3.3.3 Dynamic Adaptation

Because activation distributions vary with input, topology, or context, joint quantization can fail when treated as a globally fixed configuration. This motivates heterogeneity-aware strategies that adapt quantization behavior to structural or instance-specific variation. DegreeQuant (Tailor et al., 2021) targets graph neural networks, where degree-dependent aggregation induces systematic differences in activation magnitude and variance across nodes; it thus incorporates degree awareness into QAT to avoid overly conservative global ranges. Similarly, to address instance-specific variation in vision tasks, Liu et al. (2022) propose a dynamic framework adjusting quantization precision conditioned on individual input samples. Their method employs a lightweight bit-controller to predict optimal layer-wise bit-widths for both weights and activations, thereby reducing computational redundancy for simple inputs while preserving high precision for complex ones.

3.3.4 Quantized Generators

For generative models such as diffusion models, quantization introduces small errors at each step, which can accumulate along the diffusion time steps. For example, BiDM (Zheng et al., 2024) shows that degradation under joint quantization is not only per-layer but also per-step, because early-

step errors can change the denoising trajectory and cannot be fully corrected later. EfficientDM (He et al., 2024) focuses on efficient quantization-aware fine-tuning to adapt pretrained diffusion models to low precision under constrained training budgets. Together, these works suggest that joint QAT for generative models must account for timestep distribution shifts and compounding errors, making robustness and training efficiency jointly critical.

3.4 Quantization Target: KV Cache

As the Transformer context window expands, the KV cache becomes a primary bottleneck for memory footprint and bandwidth. Quantization of the KV cache has emerged as a critical compression target, distinct from but complementary to weight or/and activation quantization. Unlike static storage compression, KV cache quantization introduces errors that recursively participate in attention, leading to accumulated perturbations along decoding steps. and training cost. To address temporal variability, outlier amplification, and training cost, we group them into two types: distillation, and efficient end-to-end QAT.

3.4.1 Distillation

This category addresses the challenge of data availability, particularly when the pretraining data are inaccessible. Data-Free Distillation (Liu et al., 2024a) replaces missing pretraining data with synthetic sequences sampled from a full-precision teacher model, and uses distillation objectives to guide a quantized student model. A key principle is to treat the KV cache as part of the quantized computational graph rather than a post-hoc target: per-token quantization for cached keys and values is enabled during the forward pass so gradients can propagate through the quantized KV path. This alignment between training and inference-time behavior facilitates simultaneous low-bit quantization of weights, activations, and the KV cache, though it does not inherently eliminate the training overhead of QAT.

3.4.2 Efficient End-to-End QAT

The second framework prioritizes training efficiency, aiming to make end-to-end QAT feasible for large-scale models without sacrificing KV cache quantization. Bondarenko et al. (2024) propose LR-QAT, which combines low-rank auxiliary weights compatible with the quantization grid, downcasting of frozen pretrained parameters, and gradi-

ent checkpointing to reduce memory overhead. Despite focusing on training feasibility, LR-QAT maintains rigorous KV-cache quantization via per-token granularity for the keys and values, showing that KV compression can be integrated into memory-constrained training pipelines without re-designing inference-time caching. This approach highlights a practical balance between memory efficiency and quantization fidelity, making it particularly suitable for large-scale transformer models deployed in resource-limited environments.

3.5 Quantization Target: Gradient

Gradient quantization discretizes tensors computed during backpropagation, including layer output gradients, activation gradients, and parameter gradients. Weight or/and activation quantization primarily affect the forward pass, while gradient quantization directly perturbs the optimization trajectory. Since gradients determine the update directions and step sizes, errors introduced here can compound through deep gradient chains. Despite these stability challenges, gradient quantization can reduce training-time memory, lower bandwidth costs, and enable integer-friendly backward propagation.

Early work on gradient QAT focuses on controlling gradient magnitudes to avoid saturation and quantization noise at very low precision. Sun et al. (2020) propose GradScale, which uses layer-wise dynamic scaling with overflow feedback to fit gradients into the FP4 range, and introduce Two-Phase Rounding (TPR) with complementary FP4 grids to reduce rounding bias and improve stability. Beyond magnitude errors, gradient direction is also critical for stable training. Zhu et al. (2020) measure directional deviation using a cosine-based metric and propose direction-sensitive gradient clipping and direction-aware learning rate scaling to mitigate large direction errors due to quantization. For Transformers, aggressive low-bit gradient quantization is more fragile than in CNNs. Xi et al. (2023) exploit structural sparsity by splitting gradients into high- and low-bit components and computing gradients only for important ones. Chitsaz et al. (2024) further show that gradients are highly noise-sensitive: 4-bit gradient quantization is often unstable, while 8-bit quantization can work when applied selectively with per-token scaling. Overall, these studies show that effective gradient quantization requires careful range control, selective application, and awareness of model structure to maintain training stability.

4 Challenges and Future Directions

We have systematically reviewed current quantization-aware training strategies, based on the quantization targets. QAT alleviates the intensive computational demands of utilizing large-scale models. It further mitigates accuracy degradation in low-precision scenarios by explicitly injecting quantization effects (e.g., clipping, rounding, and rescaling) into model training, so that model parameters and activations can easily adapt to the discrete inference arithmetic in real-world deployment (Jacob et al., 2018). QAT commonly relies on surrogate gradients, most notably the straight-through estimator (STE), to enable end-to-end optimization (Bengio et al., 2013) to escape from non-differentiable operations like rounding. Despite strong empirical success, QAT still faces the following persistent challenges.

Optimization Instability QAT is fundamentally challenged by the mismatch between the discrete, non-differentiable nature of quantization and gradient-based optimization. Since quantization functions are piecewise constant, QAT relies on surrogate gradients such as STE to enable back-propagation. However, these heuristics do not reflect the true derivatives of the quantized objective, introducing biased gradient estimations and gradient mismatch (Bengio et al., 2013; Sakr et al., 2022). This mismatch naturally degrades optimization stability and convergence. Inaccurate gradients may fail to consistently reduce the quantized loss: near the quantization thresholds, small weight updates can cause abrupt changes in quantized values, leading to weight oscillation (Nagel et al., 2022; Liu et al., 2023). These problems are exacerbated in ultra-low-bit settings, where coarse quantization produces a highly non-smooth optimization landscape. In such regimes, surrogate gradients may be too weak or misaligned to move parameters across quantization boundaries, causing frozen weights that remain trapped within a single quantization bin (Lee et al., 2021). Overall, these limitations motivate the development of optimization methods that better align backward gradients with the true discrete objective.

Outliers in Activation. While weight quantization has become relatively mature, activation quantization remains a major bottleneck, especially for Transformers and LLMs. A key challenge is the presence of strong, input-dependent outliers in acti-

621 vations—commonly observed in attention modules
622 and residual streams—which significantly skew
623 the quantization range and degrade accuracy (Bon-
624 darenko et al., 2021; Xiao et al., 2023; Zhang and
625 Shrivastava, 2024). Moreover, under aggressive
626 joint quantization of weights and activations (e.g.,
627 W4A4), quantization noise accumulates along the
628 layers. This accumulated noise interacts unfavor-
629 ably with LayerNorm and Softmax, further ampli-
630 fying performance degradation. As a result, re-
631 covering accuracy often relies on complex algo-
632 rithmic techniques, such as knowledge distillation
633 or rotation-based transformations, which introduce
634 additional optimization difficulty and engineering
635 overhead (Zhang et al., 2020; Liu et al., 2024b).

636 **The Simulation-to-Deployment Gap.** A persis-
637 tent practical challenge lies in the discrepancy be-
638 tween fake quantization—typically simulated in
639 FP32 during training—and true integer-only in-
640 ference at deployment time. QAT pipelines often
641 overlook low-level hardware constraints, including
642 accumulator overflow, hardware-specific rounding
643 modes, and operator fusion behaviors (Jacob et al.,
644 2018; Colbert et al., 2023). As a result, models con-
645 verged successfully under QAT simulations may
646 experience substantial accuracy degradation—or
647 even fail to execute efficiently—when deployed on
648 real-world integer-only hardware accelerators such
649 as DSPs or NPUs (Kim et al., 2021; Li and Gu,
650 2023).

651 **High Training Overhead.** QAT requires a full
652 training pipeline, including access to labeled
653 datasets and substantial computational resources
654 for gradient-based optimization. This overhead be-
655 comes prohibitive for billion-parameter LLMs or
656 diffusion models (Liu et al., 2024a; Chen et al.,
657 2025a). Moreover, QAT typically maintains full-
658 precision master weights and optimizer states dur-
659 ing training, which significantly increases memory
660 consumption. As a result, there is an inherent ten-
661 sion trade-off between the objective of model com-
662 pression and the high computational and memory
663 costs incurred by the compression process itself
664 (Dremov et al., 2025; Huang et al., 2024).

665 **Homogenization.** Model homogenization is an
666 emerging challenge in model quantization. Most
667 existing quantization methods rely on uniform
668 quantizers, fixed bit-width allocation, and layer-
669 agnostic design choices to simplify deployment
670 and ensure hardware efficiency. While effective

671 for reducing model size and inference cost, these
672 choices can significantly constrain the expressive
673 capacity of quantized models. As a result, models
674 with different architectures, training objectives, or
675 task specializations often converge to highly simi-
676 lar weight distributions and activation patterns af-
677 ter quantization, leading to diminished representa-
678 tional diversity (Nagel et al., 2021; Li et al., 2024a).
679 This quantization-induced homogenization effect
680 becomes more pronounced in large-scale models.
681 As model size increases, quantization noise accu-
682 mulates across layers, encouraging “averaged” in-
683 ternal representations that suppress fine-grained
684 features and task-specific behaviors. This phe-
685 nomenon can ultimately limit the performance ceil-
686 ings and reduce the adaptability to downstream
687 tasks, particularly in LLMs (Yu et al., 2024). Ad-
688 dressing this trade-off between compression effi-
689 ciency and representational diversity remains an
690 open research challenge in model quantization.

691 Evaluating quantization methods is crucial, as
692 reductions in bit-width or model size do not neces-
693 sarily translate into faster inference or consistently
694 accurate performance. Quantization can introduce
695 subtle changes that may not be captured by simple
696 metrics such as compression ratio or theoretical
697 speedup. A thorough assessment is needed to un-
698 derstand its real impact. In particular, evaluating
699 across multiple dimensions helps identify trade-
700 offs that may not be apparent when considering ac-
701 curacy alone. We survey QAT evaluation methods
702 from three key perspectives: performance, which
703 measures task-specific accuracy; robustness, which
704 examines the model’s resilience to perturbations;
705 and efficiency, which considers both computational
706 cost and memory usage. For a more detailed dis-
707 cussion, please refer to Appendix A.

708 5 Conclusion

709 This survey reviews existing quantization-aware
710 training (QAT) methods as an effective method of
711 enabling efficient inference of large-scale models
712 by simulating quantization during training. We
713 present a unified taxonomy of the QAT methods
714 based on various quantization targets, ranging from
715 weight quantization to gradient quantization. We
716 also highlight the key challenges faced by QAT,
717 e.g., their noisy gradient estimations, high training
718 overhead, and model homogenization. Solving the
719 challenges is non-trivial, and we share promising
720 directions to address them.

721 Limitations

722 This review focuses primarily on quantization-
723 aware training, emphasizing quantization methods
724 and the challenges they entail. It does not cover rel-
725 evant datasets and benchmarks. To maintain a clear
726 and focused discussion, we present detailed evalua-
727 tions of the QAT methods in Appendix A due to the
728 page limitations. We intend to incorporate a com-
729 prehensive discussion of relevant datasets, bench-
730 marks, and evaluations in the final manuscript.

731 References

732 Haoli Bai, Wei Zhang, Lu Hou, Lifeng Shang, Jin Jin,
733 Xin Jiang, Qun Liu, Michael Lyu, and Irwin King.
734 2021. [BinaryBERT: Pushing the limit of BERT quan-](#)
735 [tization](#). In *Proceedings of the 59th Annual Meet-*
736 *ing of the Association for Computational Linguistics*
737 *and the 11th International Joint Conference on Natu-*
738 *ral Language Processing (Volume 1: Long Papers)*,
739 pages 4334–4348, Online. Association for Computa-
740 tional Linguistics.

741 Yoshua Bengio, Nicholas Léonard, and Aaron Courville.
742 2013. Estimating or propagating gradients through
743 stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*.

745 Yash Bhalgat, Jinwon Lee, Markus Nagel, Tijmen
746 Blankevoort, and Nojun Kwak. 2020. Lsq+: Improv-
747 ing low-bit quantization through learnable offsets and
748 better initialization. In *Proceedings of the IEEE/CVF*
749 *conference on computer vision and pattern recogni-*
750 *tion workshops*, pages 696–697.

751 Yelysei Bondarenko, Riccardo Del Chiaro, and Markus
752 Nagel. 2024. Low-rank quantization-aware training
753 for llms. *arXiv preprint arXiv:2406.06385*.

754 Yelysei Bondarenko, Markus Nagel, and Tijmen
755 Blankevoort. 2021. [Understanding and overcoming](#)
756 [the challenges of efficient transformer quantization](#).
757 In *Proceedings of the 2021 Conference on Empiri-*
758 *cal Methods in Natural Language Processing*, pages
759 7947–7969, Online and Punta Cana, Dominican Re-
760 public. Association for Computational Linguistics.

761 Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang,
762 Peng Gao, Kaipeng Zhang, and Ping Luo. 2025a. [Efficientqat: Efficient quantization-aware training for](#)
763 [large language models](#). In *Proceedings of the 63rd*
764 *Annual Meeting of the Association for Computational*
765 *Linguistics (Volume 1: Long Papers)*, pages 10081–
766 10100.

768 Mengzhao Chen, Chaoyi Zhang, Jing Liu, Yutao Zeng,
769 Zeyue Xue, Zhiheng Liu, Yunshui Li, Jin Ma, Jie
770 Huang, Xun Zhou, and 1 others. 2025b. Scaling
771 law for quantization-aware training. *arXiv preprint*
772 *arXiv:2505.14302*.

Kamran Chitsaz, Quentin Fournier, Goncalo Mordido,
and Sarath Chandar. 2024. [Exploring quantization](#)
[for efficient pre-training of transformer language](#)
[models](#). In *Findings of the Association for Com-*
putational Linguistics: EMNLP 2024, pages 13473–
13487, Miami, Florida, USA. Association for Com-
putational Linguistics.

Ian Colbert, Alessandro Pappalardo, and Jakoba Petri-
Koenig. 2023. A2q: Accumulator-aware quantiza-
tion with guaranteed overflow avoidance. In *Proceed-*
ings of the IEEE/CVF International Conference on
Computer Vision, pages 16989–16998.

Ian Colbert, Alessandro Pappalardo, Jakoba Petri-
Koenig, and Yaman Umuroglu. 2024. A2q+: im-
proving accumulator-aware weight quantization. In
Proceedings of the 41st International Conference on
Machine Learning, ICML’24. JMLR.org.

Peiyan Dong, Lei Lu, Chao Wu, Cheng Lyu, Geng
Yuan, Hao Tang, and Yanzhi Wang. 2023. Pack-
qvit: Faster sub-8-bit vision transformers via full and
packed quantization on the mobile. *Advances in Neu-*
ral Information Processing Systems, 36:9015–9028.

Aleksandr Dremov, David Grangier, Angelos
Katharopoulos, and Awni Hannun. 2025. Compute-
optimal quantization-aware training. *arXiv preprint*
arXiv:2509.22935.

Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani,
Rathinakumar Appuswamy, and Dharmendra S.
Modha. 2020. [Learned step size quantization](#). In
8th International Conference on Learning Represen-
tations, ICLR 2020, Addis Ababa, Ethiopia, April
26-30, 2020. OpenReview.net.

Shinya Gongyo, Jinrong Liang, Mitsuru Ambai, Rei
Kawakami, and Ikuro Sato. 2024. Learning non-
uniform step sizes for neural network quantization.
In *Proceedings of the Asian Conference on Computer*
Vision, pages 4385–4402.

Yefei He, Jing Liu, Weijia Wu, Hong Zhou, and Bohan
Zhuang. 2024. [Efficientdm: Efficient quantization-](#)
[aware fine-tuning of low-bit diffusion models](#). In
The Twelfth International Conference on Learning
Representations, ICLR 2024, Vienna, Austria, May
7-11, 2024. OpenReview.net.

Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean.
2015. [Distilling the knowledge in a neural network](#).
ArXiv, abs/1503.02531.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
Weizhu Chen, and 1 others. 2022. Lora: Low-rank
adaptation of large language models. *ICLR*, 1(2):3.

Yuxuan Hu, Xiaodong Chen, Cuiping Li, Hong Chen,
and Jing Zhang. 2025. Quad: Quantization and
parameter-efficient tuning of llm with activation de-
composition. *arXiv preprint arXiv:2503.19353*.

827	Xijie Huang, Zechun Liu, Shih-Yang Liu, and Kwang-Ting Cheng. 2024. Robust and efficient quantization-aware training via coresets selection . <i>Transactions on Machine Learning Research</i> .	882
828		883
829		884
830		885
831	Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. 2018. Quantization and training of neural networks for efficient integer-arithmic-only inference . In <i>2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 2704–2713.	886
832		887
833		888
834		889
835		890
836		891
837		
838	Sambhav Jain, Albert Gural, Michael Wu, and Chris Dick. 2020. Trained quantization thresholds for accurate and efficient fixed-point inference of deep neural networks. <i>Proceedings of Machine Learning and Systems</i> , 2:112–128.	892
839		893
840		894
841		895
842		896
843	Hyesung Jeon, Yulhwa Kim, and Jae-Joon Kim. 2025. L4q: parameter efficient quantization-aware fine-tuning on large language models. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2002–2024.	897
844		898
845		899
846		900
847		901
848		902
849	Wenjing Ke, Zhe Li, Dong Li, Lu Tian, and Emad Barsoum. 2024. DL-qat: Weight-decomposed low-rank quantization-aware training for large language models . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 113–119.	903
850		904
851		905
852		906
853		907
854		908
855	Hyunjun Kim, Kyungsu Kim, Jinseok Kim, and Jae-Joon Kim. 2020. Binaryduo: Reducing gradient mismatch in binary activation network by coupling binary activations . <i>ArXiv</i> , abs/2002.06517.	909
856		910
857		911
858		912
859	Minsoo Kim, Sihwa Lee, Suk-Jin Hong, Du-Seong Chang, and Jungwook Choi. 2022. Understanding and improving knowledge distillation for quantization aware training of large transformer encoders . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 6713–6725, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	913
860		914
861		915
862		916
863		917
864		918
865		
866		
867	Sehoon Kim, Amir Gholami, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. 2021. I-bert: Integer-only bert quantization . In <i>International conference on machine learning</i> , pages 5506–5518. PMLR.	919
868		920
869		921
870		922
871	Changhun Lee, Jun-gyu Jin, YoungHyun Cho, and Eunhyeok Park. 2024. QEFT: Quantization for efficient fine-tuning of LLMs . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 13823–13837, Miami, Florida, USA. Association for Computational Linguistics.	923
872		924
873		925
874		926
875		927
876		928
877	Junghyup Lee, Dohyung Kim, and Bumsub Ham. 2021. Network quantization with element-wise gradient scaling. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 6448–6457.	929
878		930
879		931
880		932
881		933
		934
		935
		936
	Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. 2024a. Autoregressive image generation without vector quantization . <i>ArXiv</i> , abs/2406.11838.	
	Yanjing Li, Sheng Xu, Xianbin Cao, Xiao Sun, and Baochang Zhang. 2023. Q-dm: an efficient low-bit quantized diffusion model . In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc.	
	Yanjing Li, Sheng Xu, Baochang Zhang, Xianbin Cao, Peng Gao, and Guodong Guo. 2022. Q-vit: Accurate and fully quantized low-bit vision transformer . <i>Advances in neural information processing systems</i> , 35:34451–34463.	
	Yixiao Li, Yifan Yu, Chen Liang, Nikos Karampatziakis, Pengcheng He, Weizhu Chen, and Tuo Zhao. 2024b. Loftq: Lora-fine-tuning-aware quantization for large language models . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
	Yuhang Li, Xin Dong, and Wei Wang. 2020. Additive powers-of-two quantization: An efficient non-uniform discretization for neural networks . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	
	Zhikai Li and Qingyi Gu. 2023. I-vit: Integer-only quantization for efficient vision transformer inference . In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 17065–17075.	
	Kai Liu, Qian Zheng, Kaiwen Tao, Zhiteng Li, Haotong Qin, Wenbo Li, Yong Guo, Xianglong Liu, Linghe Kong, Guihai Chen, Yulun Zhang, and Xiaokang Yang. 2025a. Low bit model quantization for deep neural networks: A survey . <i>arXiv preprint arXiv:2505.05530</i> .	
	Shih-Yang Liu, Zechun Liu, and Kwang-Ting Cheng. 2023. Oscillation-free quantization for low-bit vision transformers . In <i>International conference on machine learning</i> , pages 21813–21824. PMLR.	
	Xuwen Liu, Zhikai Li, Minghao Jiang, Mengjuan Chen, Jianquan Li, and Qingyi Gu. 2025b. Dilatequant: Accurate and efficient quantization-aware training for diffusion models via weight dilation . In <i>Proceedings of the 33rd ACM International Conference on Multimedia, MM '25</i> , page 8399–8408, New York, NY, USA. Association for Computing Machinery.	
	Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. 2024a. Llm-qat: Data-free quantization aware training for large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 467–484.	

937	Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. 2024b. Spinquant: Llm quantization with learned rotations. <i>arXiv preprint arXiv:2405.16406</i> .	993	Haoxuan Wang, Yuzhang Shang, Zhihang Yuan, Junyi Wu, Junchi Yan, and Yan Yan. 2025. Quest: Low-bit diffusion model quantization via efficient selective finetuning. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 15542–15551.	994		995		996		997		998	
943	Zhenhua Liu, Yunhe Wang, Kai Han, Siwei Ma, and Wen Gao. 2022. Instance-aware dynamic neural network quantization. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 12434–12443.	999	Zheng Wang, Juncheng B Li, Shuhui Qu, Florian Metze, and Emma Strubell. 2022. Squat: Sharpness-and quantization-aware training for bert. <i>arXiv preprint arXiv:2210.07171</i> .	1000		1001		1002					
948	Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. 2021. A white paper on neural network quantization. <i>ArXiv</i> , abs/2106.08295.	1003	Lu Wei, Zhong Ma, Chaojie Yang, and Qin Yao. 2024. Advances in the neural network quantization: A comprehensive review. <i>Applied Sciences</i> , 14(17).	1004		1005							
952	Markus Nagel, Marios Fournarakis, Yelysei Bondarenko, and Tijmen Blankevoort. 2022. Overcoming oscillations in quantization-aware training. In <i>International Conference on Machine Learning</i> , pages 16318–16330. PMLR.	1006	Quan Wei, Chung-Yiu Yau, Hoi To Wai, Yang Zhao, Dongyeop Kang, Youngsuk Park, and Mingyi Hong. 2025. RoSTE: An efficient quantization-aware supervised fine-tuning approach for large language models. In <i>Forty-second International Conference on Machine Learning</i> .	1007		1008		1009		1010		1011	
957	Jacob Nielsen, Peter Schneider-Kamp, and Lukas Galke. 2025. Continual quantization-aware pre-training: When to transition from 16-bit to 1.58-bit pre-training for BitNet language models? In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 13483–13493, Vienna, Austria. Association for Computational Linguistics.	1012	Haocheng Xi, ChangHao Li, Jianfei Chen, and Jun Zhu. 2023. Training transformers with 4-bit integers. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1013		1014		1015					
964	Charbel Sakr, Steve Dai, Rangharajan Venkatesan, Brian Zimmer, William J. Dally, and Brucec Khailany. 2022. Optimal clipping and magnitude-aware differentiation for improved quantization-aware training. In <i>International Conference on Machine Learning</i> .	1016	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 <i>International conference on machine learning</i> , pages 38087–38099. PMLR.	1017		1018		1019		1020			
968	Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. 2019. Q-bert: Hessian based ultra low precision quantization of bert. In <i>AAAI Conference on Artificial Intelligence</i> .	1021	Cong Xu, Wenbin Liang, Mo Yu, Anan Liu, Keqin Zhang, Lizhuang Ma, Jianyong Wang, Jun Wang, and Wei Zhang. 2025. Pushing the limits of low-bit optimizers: A focus on ema dynamics. <i>ArXiv</i> , abs/2505.00347.	1022		1023		1024		1025			
975	Junhyuk So, Jungwon Lee, Daehyun Ahn, Hyungjun Kim, and Eunhyeok Park. 2023. Temporal dynamic quantization for diffusion models. <i>Advances in neural information processing systems</i> , 36:48686–48698.	1026	Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. 2024. Qa-lora: Quantization-aware low-rank adaptation of large language models. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1027		1028		1029		1030		1031	
980	Xiao Sun, Naigang Wang, Chia-Yu Chen, Jiamin Ni, Ankur Agrawal, Xiaodong Cui, Swagath Venkataramani, Kaoutar El Maghraoui, Vijayalakshmi Viji Srinivasan, and Kailash Gopalakrishnan. 2020. Ultra-low precision 4-bit training of deep neural networks. <i>Advances in Neural Information Processing Systems</i> , 33:1796–1807.	1032	Jiawei Yang, Zhongbo Li, Zeqin Feng, and Yongqiang Xie. 2025. A survey on neural network quantization. In <i>Proceedings of the 2025 6th International Conference on Computer Information and Big Data Applications, CIBDA '25</i> , page 384–394, New York, NY, USA. Association for Computing Machinery.	1033		1034		1035		1036		1037	
987	Shyam A Tailor, Javier Fernandez-Marques, and Nicholas D Lane. 2021. Degree-quant: Quantization-aware training for graph neural networks. In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	1038	Jeffrey Yu, Kartik Prabhu, Yonatan Urman, Robert M. Radway, Eric Han, and Priyanka Raina. 2024. 8-bit transformer inference and fine-tuning for edge accelerators. <i>Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3</i> .	1039		1040		1041		1042		1043	
992		1044	Zhihang Yuan, Yuzhang Shang, and Zhen Dong. 2024. PB-LLM: partially binarized large language models. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1045		1046		1047		1048		1049	

1050
1051
1052
1053
1054
1055

1056
1057
1058
1059

1060
1061
1062
1063
1064
1065
1066

1067
1068
1069
1070
1071
1072

1073
1074
1075
1076
1077

1078
1079
1080
1081
1082
1083
1084
1085

1086
1087
1088
1089
1090
1091
1092

1093
1094
1095
1096
1097
1098

1099

1100
1101
1102
1103

Bonan Zhang, Chia-Yu Chen, and Naveen Verma. 2024. Reshape and adapt for output quantization (raoq): quantization-aware training for in-memory computing systems. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org.

Tianyi Zhang and Anshumali Shrivastava. 2024. Leanquant: Accurate and scalable large language model quantization with loss-error-aware grid. *arXiv preprint arXiv:2407.10032*.

Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. 2020. TernaryBERT: Distillation-aware ultra-low bit BERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 509–521, Online. Association for Computational Linguistics.

Maosen Zhao, Pengtao Chen, Chong Yu, Yan Wen, Xudong Tan, and Tao Chen. 2025. Pioneering 4-bit fp quantization for diffusion models: Mixup-sign quantization and timestep-aware fine-tuning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 18134–18143.

Xiandong Zhao, Ying Wang, Xuyi Cai, Cheng Liu, and Lei Zhang. 2020. Linear symmetric quantization of neural networks for low-precision integer hardware. In *International Conference on Learning Representations*.

Xingyu Zheng, Xianglong Liu, Yichen Bian, Xudong Ma, Yulun Zhang, Jiakai Wang, Jinyang Guo, and Haotong Qin. 2024. Bidm: Pushing the limit of quantization for diffusion models. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

Qihua Zhou, Song Guo, Zhihao Qu, Jingcai Guo, Zhenda Xu, Jiewei Zhang, Tao Guo, Boyuan Luo, and Jingren Zhou. 2021. Octo:{INT8} training with loss-aware compensation and backward quantization for tiny on-device learning. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*, pages 177–191.

Feng Zhu, Ruihao Gong, Fengwei Yu, Xianglong Liu, Yanfei Wang, Zhelong Li, Xiuqi Yang, and Junjie Yan. 2020. Towards unified int8 training for convolutional neural network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1969–1979.

A Evaluation

Quantization-Aware Training (QAT) differs from post-training quantization (PTQ) in that quantization operators are inserted into the optimization loop. As a result, evaluating QAT is inherently

multi-dimensional: a method must (i) achieve the desired task quality at a target precision, (ii) exhibit stable and well-behaved training dynamics despite non-smooth quantization and straight-through estimators (STE), and (iii) satisfy deployment constraints (e.g., integer overflow limits, accumulator precision, or analog non-idealities). To make these objectives explicit, we organize QAT evaluation into structural metrics and functional metrics, and then discuss a cross-cutting, boundary-based view enabled by recent scaling-law analyses (Chen et al., 2025b; Dremov et al., 2025).

Evaluation setup. Many metrics are only meaningful once the quantization configuration is fixed. We thus recommend stating (a) **targets** (weights/activations/gradients/KV-cache and any nonlinear approximations), (b) **granularity** (per-tensor/per-channel/group, and group size if applicable), (c) **integer format and rounding** (symmetric/asymmetric, stochastic/deterministic rounding, clipping rule), and (d) **hardware model vs. measurement** (simulator assumptions, accumulator bit-width, ADC precision, kernel backend, or on-device profiling). When reporting efficiency, also specify sequence length/batch size (for transformers), compiler/kernel stack, and whether comparisons share identical implementations.

Structural metrics vs. functional metrics. Structural metrics diagnose whether QAT behaves as intended during training: (1) quantization distortion (tensor-level and its propagation), (2) optimization mismatch induced by STE and quantization (e.g., oscillations and gradient misalignment), and (3) feasibility under numeric/device constraints (overflow/saturation/non-ideality margins). Functional metrics measure usefulness in deployment: task quality at target precision, training and inference efficiency on the intended system, and robustness to precision, data and system shifts.

A.1 Structural Evaluation

A.1.1 Quantization Error

Structural evaluation often begins with direct measures of how closely quantized tensors approximate their full-precision counterparts.

Tensor-level distortion. Given a full-precision tensor X (weights or activations) and a quantizer $Q(\cdot)$, common proxies include normalized MSE,

$$\text{nMSE}(X) = \frac{\|X - Q(X)\|_F^2}{\|X\|_F^2}, \quad (10)$$

1250	Comparability across configurations.	To avoid cherry-picking, we recommend reporting a small grid of standard settings: (i) multiple bit-widths (e.g., INT8, W4A4, W3A3, W2A2 as applicable), (ii) ablations on which tensors are quantized (W-only, A-only, KV-cache), and (iii) granularity choices (per-channel vs. per-tensor vs. group). When distillation or data-free supervision is used, report the teacher metric and the student–teacher gap explicitly (Liu et al., 2024a).	1299
1251			1300
1252			1301
1253			
1254			
1255			
1256			
1257			
1258			
1259			
1260	A.2.2 Efficiency		
1261	Algorithmic cost proxies.	Common proxies include peak memory, wall-clock time per step, GPU hours, and total training FLOPs. For QAT schedules that mix precisions, report compute and the schedule itself (transition points, phase lengths). Studies on scaling and compute-optimality treat compute budget as a first-class metric and report validation loss/perplexity as a function of total FLOPs under alternative precision allocations (Chen et al., 2025b; Dremov et al., 2025; Nielsen et al., 2025).	
1262			
1263			
1264			
1265			
1266			
1267			
1268			
1269			
1270			
1271			
1272	System measurements.	Ultimately, inference benefits must be measured on the intended system: latency, throughput, memory footprint, and energy (when possible). For fair comparison, report hardware model, kernel backend, and input shapes. Mobile/edge evaluations often quantify speedups from packing or integerizing operators (Dong et al., 2023; Li and Gu, 2023), while hardware-aware QAT in IMC or limited-accumulator settings reports energy/throughput under varying device constraints (Zhang et al., 2024; Colbert et al., 2023, 2024). On-device learning evaluates training-time energy and wall-clock feasibility under microcontroller constraints (Zhou et al., 2021).	
1273			
1274			
1275			
1276			
1277			
1278			
1279			
1280			
1281			
1282			
1283			
1284			
1285			
1286	A.2.3 Robustness		
1287	Precision robustness.	A basic robustness check is performance across multiple precisions and configurations, revealing whether a recipe overfits to a particular bit-width or quantization target (e.g., W-only vs. W + A). Dynamic quantization methods additionally test whether conditional policies generalize across datasets (Liu et al., 2022).	
1288			
1289			
1290			
1291			
1292			
1293			
1294	Data and task shifts.	Robustness should include evaluation on distribution shifts or unseen tasks/domains, especially for LLM and PEFT-based QAT where extreme low-bit regimes may disproportionately affect reasoning or safety behavior (Xu et al., 2024; Li et al., 2024b; Jeon et al., 2025; Ke et al., 2024; Lee et al., 2024; Wei et al., 2025; Hu et al., 2025; Bondarenko et al., 2024).	1299
1295			1300
1296			1301
1297			
1298			
		System shifts and implementation sensitivity.	1302
		Finally, robustness must consider deployment variation: kernel/back-end changes, rounding modes, accumulator precision differences, or analog non-idealities. Reporting sensitivity to such factors helps distinguish recipes that generalize to real deployments from those that overfit a particular simulator or kernel stack (Zhang et al., 2024; Colbert et al., 2023).	1303
			1304
			1305
			1306
			1307
			1308
			1309
			1310
		Recommended minimal functional protocol.	1311
		At a minimum, we recommend reporting: (i) the task performance at the target precision, (ii) at least one efficiency metric (either a proxy metric and/or a system-level metric), and (iii) one robustness slice (e.g., across different bit-widths or under a deployment-time distribution shift), with all quantization configuration details clearly specified upfront. Adhering to this reporting protocol is critical for ensuring reproducibility and enables fair and meaningful comparisons across diverse quantization methods, reducing ambiguity in evaluation and strengthening the reliability of empirical conclusions.	1312
			1313
			1314
			1315
			1316
			1317
			1318
			1319
			1320
			1321
			1322
			1323
			1324

Approach	Venue	Centric-Target	Bits	Gran.	Model
TernaryBERT	2020 EMNLP	W	2	Layer, Row	BERT
BinaryBERT	2021 ACL	W	1	Layer, Tensor	BERT
Nagel et al.	2022 ICML	W	3, 4	Tensor	General
A2Q	2023 ICCV	W	5–8	Channel	General
OFQ	2023 ICML	W	2, 3, 4	Row	ViT
DL-QAT	2024 EMNLP	W	3, 4	Group	LLM
QEFT	2024 EMNLP Findings	W	4	Group	LLM
PB-LLM	2024 ICLR	W	1	Column	LLM
QA-LoRA	2024 ICLR	W	4	Group	LLM
LoftQ	2024 ICLR	W	2, 4	Tensor	LLM
EfficientQAT	2025 ACL	W	2	Group	LLM
L4Q	2025 ACL	W	3, 4	Group	LLM
CQPT	2025 ACL Findings	W	1.58	Layer, Token	LLM
LSQ+	2020 CVPR	A	2, 4	Layer	General
BinaryDuo	2020 ICLR	A	1	/	BNN
LLSQ	2020 ICLR	W + A	W4A4, W3	Channel, Layer	General
Jain et al.	2020 MLSys	W + A	W8A8, W4	Tensor	General
Degree-Quant	2021 ICLR	W + A	W8A8, W4	Tensor	GNN
Liu et al.	2022 CVPR	W + A	2–6	Layer, Token	General
OCTAV	2022 ICML	W + A	4–6	Tensor	General
Q-ViT	2022 NeurIPS	W + A	W4A4, W3	/	ViT
nu-LSQ	2024 ACCV	W + A	2, 3, 4	Layer	General
EfficientDM	2024 ICLR Spotlight	W + A	W4A4	Channel, Layer	DM
BiDM	2024 NeurIPS	W + A	W/A 1	Layer	DM
per-embed	2021 EMNLP	W + A	W8A8	Tensor	Transformer
I-BERT	2021 ICML	W + A	8	/	BERT
I-ViT	2023 ICCV	W + A	8	/	ViT
Zhu et al.	2020 CVPR	G	8	Layer	CNN
GradScale	2020 NeurIPS	G	4	Layer	General
Xi et al.	2023 NeurIPS	G	4	Token	Transformer
Chitsaz et al.	2024 EMNLP Findings	G	4, 8	Tensor, Channel, Token	LLM
LLM-QAT	2024 ACL Findings	KV	8, 16	Token	LLM
LR-QAT	2024 arXiv	KV	4, 8	/	LLM

Table 1: Taxonomy of quantization-aware training (QAT) frameworks. Methods are grouped by quantization centric-target (weights, activations, weights+activations, gradients, and KV cache) and summarized by bit-width, granularity, and model.