

Content-Aware Radiance Fields: Aligning Model Complexity with Scene Intricacy Through Learned Bitwidth Quantization

Weihang Liu^{1,*}, Xue Xian Zheng^{2,*}, Jingyi Yu^{1,3}, and Xin Lou^{1,3,†}

¹ ShanghaiTech University

² King Abdullah University of Science and Technology

³ Key Laboratory of Intelligent Perception and Human-Machine Collaboration

Abstract. The recent popular radiance field models, exemplified by Neural Radiance Fields (NeRF), Instant-NGP and 3D Gaussian Splatting, are designed to represent 3D content by that training models for each individual scene. This unique characteristic of scene representation and per-scene training distinguishes radiance field models from other neural models, because complex scenes necessitate models with higher representational capacity and vice versa. In this paper, we propose content-aware radiance fields, aligning the model complexity with the scene intricacies through Adversarial Content-Aware Quantization (A-CAQ). Specifically, we make the bitwidth of parameters differentiable and trainable, tailored to the unique characteristics of specific scenes and requirements. The proposed framework has been assessed on Instant-NGP, a well-known NeRF variant and evaluated using various datasets. Experimental results demonstrate a notable reduction in computational complexity, while preserving the requisite reconstruction and rendering quality, making it beneficial for practical deployment of radiance fields models. Codes are available at https://github.com/WeihangLiu2024/Content_Aware_NeRF.

Keywords: Radiance fields · Content-aware · Quantization · Model complexity

1 Introduction

Triggered by the phenomenal Neural Radiance Fields (NeRF) [23], the idea of representing 3D scenes using trainable models has been widely adopted in many reconstruction and rendering applications. Different from traditional explicit representations like meshes and point clouds, radiance field techniques encode a 3D scene with learnable parameters, which are obtained by training models based on sparse samples of the scene. Despite its advantages, the high-quality representation and rendering offered by radiance fields come at the expense of significant

* Equal Contribution.

† Corresponding author.

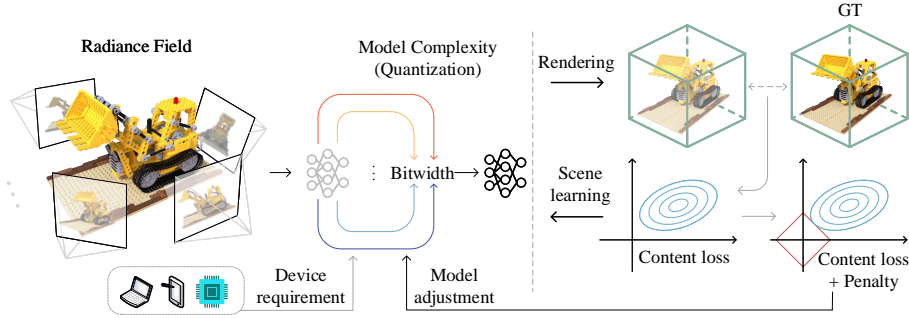


Fig. 1: Overview of content-aware radiance fields. In this work, model complexity is aligned with scene intricacy through learned bitwidth quantization.

computational complexity. This challenge has been extensively studied, leading to numerous research endeavors aimed at exploring more computationally efficient radiance field models [12, 22, 28, 33].

While many variants of NeRF have significantly improved efficiency, they typically compress all scenes to a single fixed scale. Very few of them consider dealing with the scene contents differently, i.e., employing a *content-aware* strategy to align model complexity with scene intricacy. This principle, as illustrated in Fig. 1, forms the cornerstone of our work in this paper. It is unique to radiance fields due to their distinctive attributes of scene representation and per-scene training. Intuitively, detailed scenes rich in geometry and texture necessitate more sophisticated radiance field models to capture their nuances. While on the other hand, encoding less complex scenes with simpler models not only suffices for adequate representation but also enhances the efficiency by reducing computational and memory requirements.

Quantization has been proven a fundamental yet effective technique for reducing model complexity. The selection of parameter bitwidth is crucial for balancing computation and performance, bridging the gap between content information loss and computational efficiency. Recent works have investigated this technique on radiance field models by quantizing them to a pre-defined fixed bitwidth [27] as well as learning quantization range [40]. Nevertheless, these methods often rely on extensive human expertise to select appropriate quantization parameters, resorting to trial and error approaches. Moreover, reducing the bitwidth with controllable performance loss is impractical within their frameworks.

In this work, we introduce the concept of content-aware radiance fields, which adaptively quantize models by exploiting the content differences among scenes and the features of each layer. The quantization bitwidth of radiance models, which directly correlates with computational complexity, are therefore content-aware, i.e., intricate scenes are accommodated with higher bitwidth models, while simpler scenes utilize lower bitwidth counterparts to reduce computational complexity. The key insights for our method are illustrated in Fig. 2, where average

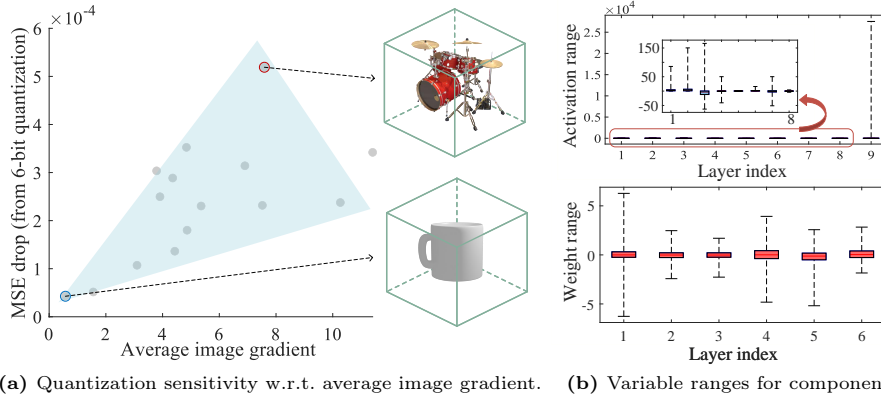


Fig. 2: The insights of the proposed LBQ. (a) The correlation between quantization sensitivity and scene complexity measured using average image gradient of the training set. Scenes with complex (simple) geometry and texture suffer more (less) accuracy degradation from quantization. (b) Exhibits the notable distinction of variables' distribution among different components. Those distributed in large (small) range is required to be quantized with high (low) bitwidth.

image gradient is used as an estimator for the complexity of the scenes. Results in Fig. 2a indicate that scenes with complex structure and texture suffer more from quantization than those with lower complexity. Besides, to achieve integer-only inference in rendering, all layers are required to be quantized by considering the significant statistical distinction of output features through the entire pipeline (shown in Fig. 2b). Moreover, the results in Fig. 3 demonstrate that scene-wise as well as layer-wise quantization is beneficial and quantization-aware training (QAT) can effectively combat quantization degradation.

Specifically, to establish the connection between quantization sensitivity and bitwidth, we propose Learned Bitwidth Quantization (LBQ) framework illustrated in Fig. 4, facilitating learnable bitwidth based on reconstructed contents. To fully extract the representation capability of quantized models, Adversarial Content-Aware Quantization (A-CAQ) algorithm is further proposed, which searches scene-dependent bitwidth and optimizes radiance fields simultaneously. The main contributions of this work, as illustrated in Fig. 1, are as follows:

- We introduce content-aware radiance fields that aligns the complexity of models with the intricacies of scenes through quantization, showing that the parameter bitwidth bridges the gap between the content information loss and computational efficiency.
- We introduce the LBQ framework which models the integer bitwidth with "soft bitwidth". This approach makes bitwidth differentiable from the content information loss during training, thereby obviating extensive human expertise to select bitwidths through a trial and error approach.

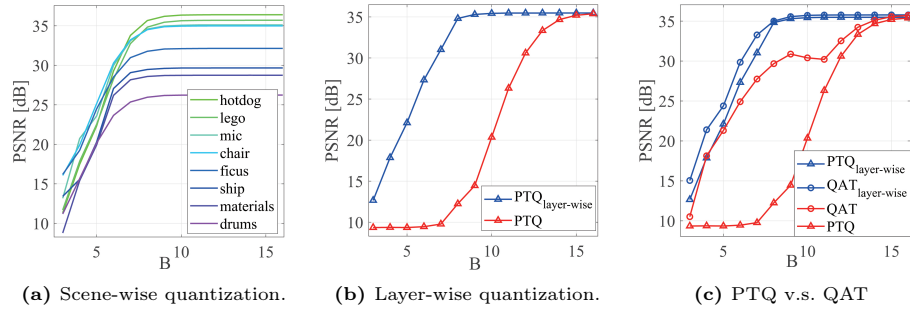


Fig. 3: The insights of the proposed A-CAQ. (a) Layer-wise quantization results for different scenes with various bitwidths, which reveals content-aware characteristics of quantization effects. (b) Results of layer-wise and non layer-wise quantization for the "lego" scene. The huge accuracy gap verifies the significance of mixed-precision models. (c) QAT alleviates performance degradation as quantization error expands.

- We propose A-CAQ that integrates LBQ to penalize layer-wise bitwidth, enabling dynamic learning of lower bitwidth with negligible reconstruction and rendering quality loss. Experimental results on various datasets confirm the superiority of the proposed method.

2 Related Work and Motivation

Scene representations and radiance fields. Various follow-up works on NeRF have focused specifically on improving computational efficiency. Rendering based on simplified representations is much more efficient while maintaining quality, which motivates studies of compression for scene primitives. A number of recent works achieve this goal from different perspectives [6, 9, 11, 14, 21, 24, 29]. Specifically, by employing a learnable codebook along with a compact MLP [24, 29], Instant-NGP mitigates representation workload of network parameters, resulting in superior results in terms of both rendering quality and speed among NeRF variants.

While content information is easily acquired by explicit primitives, it is challenging for implicit representations. VQ-AD [29] and SHACIRA [14] introduced scalable compression of trainable feature grid which render content with different qualities. Recursive NeRF [37] proposes to dynamically grow the network, considering complexity variations of different patches within one scene. However, none of these works are metric-oriented or consider the complexity of different scenes. The search for the most suitable compression, considering both content and available resources, remains an empirical task.

Model quantization. Network quantization has been demonstrated as a potent technique for compressing neural models [4, 7, 13, 26, 38]. Generally, Post-Training

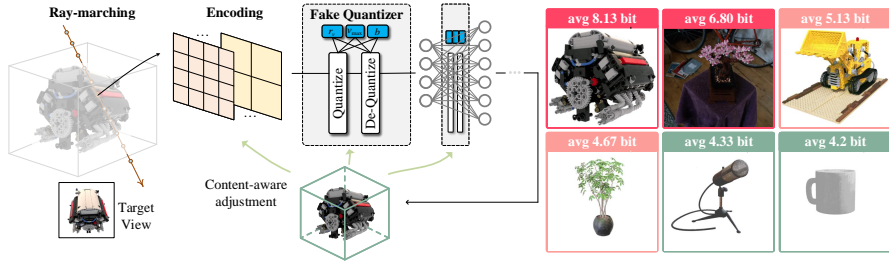


Fig. 4: Overview of the LBQ quantization framework. The fake quantizers are inserted into different components including encoding and MLPs, which are parameterized with variable range scale, upper bound and bitwidth. Quantization error is simulated by quantize and de-quantize procedure and quantization parameters are updated with direction indicated by gradient descent. Examples with different complexity are given on the right.

Quantization (PTQ) [8, 18] and Quantization-Aware Training (QAT) [41] are two common methods for model quantization. Some recent works endeavor to dynamically allocate bitwidths based on the features of each layer, resulting in mixed-precision quantization [5, 15, 34]. These mixed-precision methods have succeeded in obtaining efficient network for various applications [15, 17].

One significant limitation of existing PTQ and QAT is that they quantize pre-trained models using hyper-defined bitwidths. Re-quantizing models to other bitwidths can be challenging, as it suffer from differences in statistical characteristics of weights and activations [19]. One straightforward solution is to find the optimal bitwidth by considering specific application requirements. CADyQ [17] employ this idea to quantize super-resolution (SR) image network based on image gradient and standard deviation of features. Nevertheless, different from SR net based on convolutional neural networks (CNNs), there is no explicit pattern for MLPs features or 3D scene complexity measurement for neural fields. Therefore, we propose a novel framework that enables bitwidth differentiation, leveraging content-related Mean Squared Error (MSE) to adeptly select optimal scene-dependent layer-wise bitwidth.

3 Proposed Method

3.1 Preliminary

Radiance fields with learnable codebook. Instant-NGP, a well-known NeRF variants incorporating a learnable multi-resolution hash table, is selected as our baseline model due to its inclusion of prevalent feature modalities in radiance fields, namely spatial levels-of-detail (LOD) feature grid and MLPs. Specifically, it approximates the mapping from continuous 5D vectors to opacities σ and view-dependent colors \mathbf{c} , denoted as

$$F : (\mathbf{x}, \mathbf{d}) \rightarrow (\sigma, \mathbf{c}), \quad (1)$$

where \mathbf{x} and \mathbf{d} are 3D location coordinate (of sample points) and 2D view direction respectively. Given a multi-resolution feature table with trainable parameters Θ , the location coordinates are encoded as $\mathbf{y}_x = \text{enc}(\mathbf{x}, \Theta)$, and the view directions are encoded as \mathbf{y}_d with spherical harmonics [24]. The encoded vector is then fed to the MLP to produce the color and opacity as

$$(\sigma, \mathbf{c}) = m(\mathbf{y}, \Phi), \quad (2)$$

where $\mathbf{y} = (\mathbf{y}_x, \mathbf{y}_d)$, Φ denotes the weights of MLPs. With the colors and opacities of all sample points, the expected color along a ray \mathbf{r} can be calculated using volume rendering as

$$\mathbf{C}(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right), \quad (3)$$

$\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ are sampled points along the ray defined by origin \mathbf{o} and direction \mathbf{d} .

Quantization. To attain integer-only inference, quantization of both parameters (*i.e.* $\Omega = \{\Theta, \Phi\}$) and activations for each layer is imperative. As shown in Fig. 4, fake quantizers can be inserted to simulate quantization error in training procedure. Given the data to be quantized \mathbf{v} , along with the step size s and zero-point offset Z , the function of a fake quantizer is to dequantize the quantized data $\hat{\mathbf{V}}$. This operation can be expressed as

$$\hat{\mathbf{v}} = s(\hat{\mathbf{V}} - Z) = s \left[\text{clamp} \left(\left\lfloor \frac{\mathbf{v}}{s} \right\rfloor + Z; q_{\min}, q_{\max} \right) - Z \right], \quad (4)$$

where $\lfloor \cdot \rfloor$ is the round to nearest operator, and $\text{clamp}(\cdot)$ is defined as

$$\text{clamp}(x; a, c) = \begin{cases} a & x < a \\ x & a \leq x \leq c \\ c & x > c \end{cases}. \quad (5)$$

Based on Eq. (4), quantization training can be conducted to mitigate introduced quantization noise.

3.2 Differentiable Quantization for Radiance Fields

Learned Bitwidth Quantization (LBQ). The quantization training model has been established in Sec. 3.1. In this section, we initially reveal that the prevalent QAT-based quantization technique is inherently limited to fixed-bitwidth scenarios, followed by the introduction of the innovative LBQ scheme used in our content-aware quantization framework.

Existing QAT schemes, such as LSQ [10], LSQ+ [3], aim to directly train the step size and zero-point offset as

$$s = \frac{v_{\max} - v_{\min}}{q_{\max} - q_{\min}} = \frac{r_v}{r_q}, \quad (6)$$

Table 1: Derivatives to key parameters.

Variable range	$\partial\hat{v}/\partial r_v$	$\partial\hat{v}/\partial b$	$\partial\hat{v}/\partial v_{\max}$
$v_{\min} \leq v \leq v_{\max}$	$(s \cdot \lfloor v/s \rfloor - v) / r_v$	$(v - s \cdot \lfloor v/s \rfloor) \cdot 2^B \ln 2 / r_q$	0
$v > v_{\max}$	$1 - v_{\max}/r_v - Z/r_q$	$(v_{\max} - r_v + sZ) \cdot 2^B \ln 2 / r_q$	1
$v < v_{\min}$	$-v_{\max}/r_v - Z/r_q$		

$$Z = \left\lfloor q_{\max} - \frac{v_{\max}}{s} \right\rfloor = \left\lfloor q_{\max} - \frac{v_{\max}}{r_v} r_q \right\rfloor, \quad (7)$$

where r_v and $r_q = 2^B - 1$ are the range scales of \mathbf{v} and $\hat{\mathbf{V}}$, respectively. To make the "round" operation differentiable, the Straight-Trough-Estimator (STE) [2] $\partial \lfloor x \rfloor / \partial x = 1$ is assumed. It indicates that s and Z can only be regarded as leaf nodes with fixed bitwidth, where r_q is considered a constant.

To further make bitwidth scalable, we introduce floating point "soft bitwidth" b , where $B = \lfloor b \rfloor$. Moreover, r_v , v_{\max} and b are selected as trainable quantization parameters, where r_v and v_{\max} serve as replacements of step size and offset, respectively. The partial derivatives of Eq. (4) can then be derived as in Tab. 1.

The foregoing equations elucidate the operational mechanism of the proposed differentiable quantization:

- The parameter v_{\max} represents the offsets of the quantization range, calibrated solely upon the occurrence of overflow, as explicated in Tab. 1.
- Both r_v and b are updated in response to observed quantization errors calculated with the specific term $(v - s \cdot \lfloor v/s \rfloor)$ in Tab. 1.

This proposed method makes bitwidth learnable during the training procedure.

Quantization schemes for radiance field models. The criterion for selecting quantization parameters is to minimize the error summation generated from rounding and overflow [25]. Considering statistical properties and inference efficiency, we establish three different quantization schemes for different components within the radiance field pipeline.

Neural weights. [3] indicates that symmetric signed quantization is highly recommended for neural weights which are empirically distributed symmetrically around zero. More importantly, quantizing weights in this manner introduces no additional computational overhead during inference.

ReLU and exponential activations. As these functions always produce positive output, we use asymmetric unsigned quantization. Similar settings are found in [10]. Based on this configuration, lower bitwidth can be reached.

Positional encoding (PE) and others. To deal with components lacking significant statistical features, we introduce v_{\max} to represent trainable offset. [25]

Table 2: Quantization schemes for different components in radiance field pipelines.

Module name	v_{\max}	r_v	b	$[q_{\min}, q_{\max}]$
Neural weights	N/A	trainable	trainable	$[-2^{B-1}, 2^{B-1} - 1]$
ReLU and exponential	N/A	trainable	trainable	$[0, 2^B - 1]$
PE and others	trainable	trainable	trainable	$[0, 2^B - 1]$

proves that this configuration will not introduce any additional computational overhead during inference. The learnable codebook is quantized in this manner.

The training details about these schemes are provided in Tab. 2. Calculations of gradients w.r.t. three different sets of trainable quantization parameters can be found in *Supplementary materials*.

3.3 Adversarial Content-Aware Quantization (A-CAQ)

As we have illustrated in Fig. 2, the impact of quantization varies depending on the contents of scenes. To leverage this characteristic effectively, it is crucial to extract content-related information from radiance fields. While it is straightforward for explicit representations such as 3D Gaussian Splatting (3DGS) [20], it presents a formidable challenge for neural fields, primarily owing to the implicit nature of the patterns encapsulated within MLPs. Nevertheless, MSE between rendered view and ground truth emerges as a fortuitous revelation, serving as a compelling indicator of the precision achieved in reconstructing the 3D content. Therefore, the MSE loss assumes a pivotal role for ascertaining content-awareness of radiance field modeling.

Utilizing LBQ proposed in Sec. 3.2, scene-dependent layer-wise quantization schemes can therefore be learned from MSE. However, bitwidth cannot adhere to the same objective function that guides other parameters. The discrepancy arises from the accuracy degradation inevitable introduced by reducing bitwidth, which operates in adversarial manner compared to optimizing MSE. To address this problem, we proposed to define bitwidth learning loss as

$$\mathcal{L}^{\text{bit}} = \sqrt{\|\mathcal{L}^{\text{NeRF}} - \mathcal{L}^{\text{metric}}\|} + \sum_{i \in \mathcal{M}} \epsilon_i B_i, \quad (8)$$

where $\mathcal{M} = \{1, 2, \dots, M\}$ is the indexes of layers or components, $\mathcal{L}^{\text{NeRF}}$, defined as

$$\mathcal{L}^{\text{NeRF}} = \sum_{\ell \in \mathcal{R}} \|\hat{C}(\ell) - C(\ell)\|_{\mathbb{F}}^2, \quad (9)$$

is the MSE loss function used for training radiance field [23], \mathcal{R} is the set of rays in one batch, and the term $\mathcal{L}^{\text{metric}}$ is a hyper-defined metric representing the minimal accuracy requirement. By employing various loss metrics, rendering quality can be controlled, transferring redundant accuracy to efficiency. Weighted bitwidth penalties are further introduced to remove redundant bitwidth that

have minimal impact on the overall loss. In addition, the weights ϵ_i offer greater flexibility, enabling assignment according to specific requirements. For example, higher penalty could be allocated to the bitwidth of the learnable codebook to obtain memory-efficient quantization schemes.

Dynamic bitwidth search can be achieved by solely optimizing Eq. (8). However, this approach results in learning bitwidth in a PTQ manner. Directly searching for bitwidth in the QAT space is impractical, as QAT accuracy can only be obtained through trial and error. As illustrated in Fig. 3c, the performance gap between PTQ and QAT grows when quantizing to lower bitwidth. To mitigate this degradation, we propose A-CAQ, a multi-task learning-based method expressed as

$$\mathcal{L}^{\text{A-CAQ}} = \min_{\mathcal{Q}} \mathcal{L}^{\text{NeRF}} + \min_{\mathbf{b}} \mathcal{L}^{\text{bit}}. \quad (10)$$

where $\mathcal{Q} = \{\Omega, \mathbf{v}_{\max}, \mathbf{r}_v\}$. This task can be effectively solved by optimizing Eq. (8) and Eq. (9) alternatively. Detailed pseudo codes are found in *Supplementary Materials*.

The interpretation of Eq. (10) primarily pertains to adversarial learning: minimizing Eq. (8) leads to a lower bitwidth solution as well as higher $\mathcal{L}^{\text{NeRF}}$. And a lower bitwidth provides $\mathcal{L}^{\text{NeRF}}$ with more potential to alleviate accuracy degradation, thereby creating more search space to achieve a lower bitwidth solution.

4 Experiments

Experiments are conducted to demonstrate the effectiveness and versatility of our content-aware quantization framework for radiance fields. We first provide the implementation details, including model and training configurations, in Sec. 4.1, followed by quantitative and qualitative results in Sec. 4.2 and Sec. 4.3, respectively. Ablation studies and complexity analysis are presented to validate the effectiveness of the proposed algorithm in Sec. 4.4 and Sec. 4.5, respectively.

4.1 Implementation details

Models. As mentioned in Sec. 3.1, our proposed quantization framework is evaluated on Instant-NGP [31], a well-known NeRF variant. All activations and parameters are quantized to facilitate integer-only inference. The bitwidth is constrained within the range of [2, 32], as binary quantization consistently yields nonsensical rendering results. Quantization schemes for each component are presented in Tab. 2.

Training details For layer-wise quantization where bitwidth is learned individually for each component, feature quantization rate (FQR) [17], defined as $\text{FQR} = \frac{\sum_{i \in \mathcal{M}} B_i}{M}$, is introduced to measure the quantization performance. To

Table 3: A-CAQ results for metric-guided bitwidth learning on different datasets including Synthetic-NeRF [23], RTMV [32] and Mip-NeRF360 [1].

Dataset	Full precision		MDL		MGL ($10^{-3.2}$)		MGL (10^{-3})	
	PSNR $_{\uparrow}$	FQR $_{\downarrow}$	PSNR $_{\uparrow}$	FQR $_{\downarrow}$	PSNR $_{\uparrow}$	FQR $_{\downarrow}$	PSNR $_{\uparrow}$	FQR $_{\downarrow}$
chair	35.06	32.00	34.57	7.60	29.63	4.80	27.23	4.60
V8	27.68	32.00	27.39	8.00	27.39	7.40	26.16	5.67
bonsai	28.13	32.00	27.49	7.00	27.73	7.07	26.50	5.73

perform training in quantization mode, quantization parameters are initialized using the simple PTQ method. All components are initially quantized to 8-bit except for exponential activation, which is quantized to 32-bit due to its extensive range scale (see Fig. 2b).

As we introduce per-defined $\mathcal{L}^{\text{metric}}$, the compression rate can be manipulated according to specific requirements. Based on the average full precision loss among training set $\mathcal{L}_{\text{fp}}^{\text{NeRF}}$, we conduct experiments with two scenarios:

Minimal Degradation bitwidth Learning (MDL). For application with high-fidelity requirements, we need to maintain accuracy while minimizing bitwidth. The metric is thus defined as

$$\mathcal{L}^{\text{metric}} = \mathcal{L}_{\text{fp}}^{\text{NeRF}}. \quad (11)$$

Metric-Guided bitwidth Learning (MGL). As numerous studies strive for ever-higher accuracy, there inevitably arises surplus accuracy for applications with varying precision requirements, such as LOD and rendering on resource-constrained edge devices. The surplus accuracy can be traded off for efficiency by quantizing into lower bitwidth. In this case, the metric is defined as

$$\mathcal{L}^{\text{metric}} > \mathcal{L}_{\text{fp}}^{\text{NeRF}}. \quad (12)$$

4.2 Quantitative Results

In this section, we evaluate the performance of our A-CAQ algorithm on different datasets, including Synthetic-NeRF [23], RTMV [32], and Mip-NeRF360 [1]. These datasets comprise scenes with varying levels of complexity.

The experiments begin with various accuracy requirements to demonstrate the capability and flexibility of the metric-guided feature. The experimental results are presented in Tab. 3. Quantization schemes of different accuracy can be effectively learned with different $\mathcal{L}^{\text{metric}}$. MSE loss is selected as a general metric here. Note that other quality metrics such as PSNR, Structural Similarity Index (SSIM [35]) and Learned Perceptual Image Path Similarity (LPIPS [42]) can also be used according to the specific requirements of different applications.

To validate the effectiveness of our A-CAQ in consistently producing efficient quantization schemes while maintaining requisite quality, we compare it with PTQ and LSQ+ schemes with different bitwidth configurations, which are

Table 4: Quantitative comparisons. Instant-NGP quantized with PTQ [27], LSQ+ [3,40,43], and A-CAQ are compared. The FQR and PSNR are reported to measure the complexity and accuracy, respectively. The results demonstrate that proposed methods succeed in reducing FQR while minimizing accuracy degradation.

Method		Synthetic-NeRF		Mip-NeRF360		RTMV (V8)	
		FQR \downarrow	PSNR \uparrow	FQR \downarrow	PSNR \uparrow	FQR \downarrow	PSNR \uparrow
	NGP	32.00	32.42	32.00	25.55	32.00	27.68
MDL	NGP-PTQ [27]	9.60	31.98	9.60	25.38	9.60	27.29
	NGP-LSQ+ [3,40,43]	9.60	32.11	9.60	25.48	9.60	27.40
	NGP-A-CAQ (Ours)	7.76	32.00	7.11	25.30	8.00	27.39
MGL ($10^{-3.2}$)	NGP-PTQ	6.60	22.29	7.60	22.62	7.60	22.96
	NGP-LSQ+	6.60	25.06	7.60	23.84	7.60	24.50
	NGP-A-CAQ (Ours)	5.33	27.58	6.86	25.18	7.40	27.39
MGL (10^{-3})	NGP-PTQ	5.60	17.75	6.60	17.54	5.60	10.54
	NGP-LSQ+	5.60	20.26	6.60	22.44	5.60	14.36
	NGP-A-CAQ (Ours)	4.74	25.96	6.58	24.74	5.67	26.16

widely used for radiance field model quantization [27,40,43]. As shown in Tab. 4, the proposed methods can identify optimal quantization schemes for various scenarios. For high-fidelity applications (labeled as MDL), integer-only rendering is achieved with negligible PSNR loss (< 0.5 dB), while the FQR is significantly reduced. In resource-constrained scenarios (labeled as MGL), our method achieves state-of-the-art results for both accuracy and efficiency.

4.3 Qualitative Results

We also present qualitative results and comparisons with other widely used methods. As depicted in Fig. 5, our A-CAQ consistently yields visually clean results across various loss metrics while requiring fewer bits. In contrast, existing quantization paradigms introduce significant distortion with comparable bitwidth configurations, as they are unable to discern quantization sensitivities among different scenes and layers.

To demonstrate the effectiveness and flexibility of our method, LOD rendering results are presented in Fig. 6. Our dynamic quantization scheme can simultaneously filter and compress radiance field models, proving advantageous for LOD tasks [29]. Furthermore, quantization results for scenes with varying complexities are illustrated in Fig. 7. Content complexity influences the quantization sensitivity for output features of each layer, which can be distinguished with the proposed A-CAQ algorithm.

4.4 Ablation study

Scene-wise and layer-wise quantization. To verify the significance of scene-wise and layer-wise quantization, we compare our proposed quantization framework with fixed-bitwidth QAT schemes, whose bitwidths are determined through a trial and error approach. Results are presented in Tab. 5. Quantizing all scenes

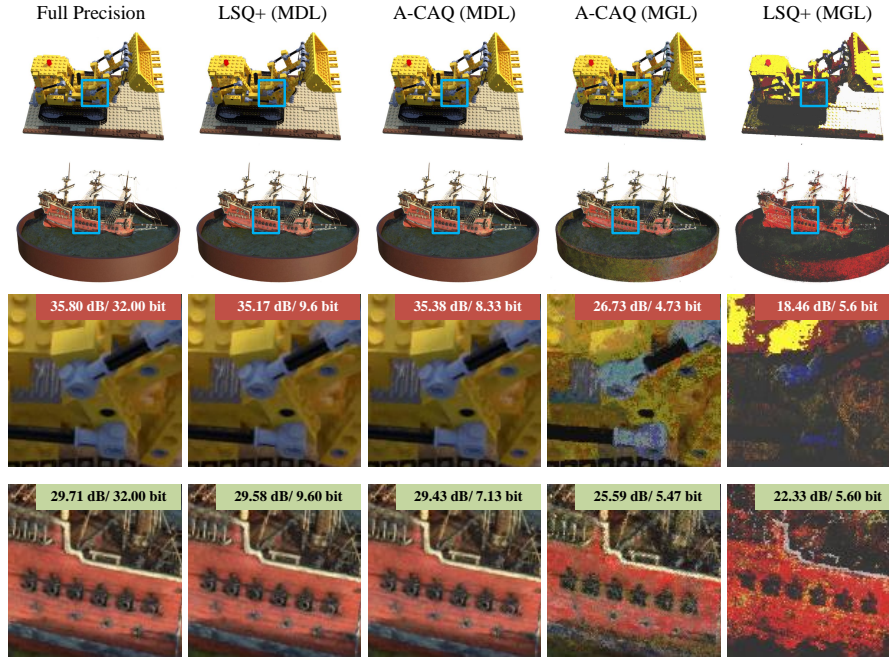


Fig. 5: Qualitative results of the "lego" and "ship" datasets from Synthetic-NeRF. Proposed A-CAQ outputs visual clean results for both MDL and MGL scenarios.

with a fixed bitwidth throughout the entire pipeline results in performance loss and resource wastage (model **(1a)**). To accommodate the varying complexities observed in different scenes, a manual adjustment of scene-specific bitwidth is conducted, which mitigates the degradation in quantization (model **(1b)**). Utilizing A-CAQ, scene-dependent layer-wise quantization bitwidths are effectively learned. Following the MDL approach, the lowest average bitwidth schemes are learned while maintaining accuracy (model **(Ours, 1e)**). As a comparison, we quantize all scenes with the same conservative scheme: quantizing each layer with the highest layer-wise bitwidth learned by A-CAQ among all scenes (model **(Ours, 1d)**). This achieves the highest PSNR among all models while sacrificing extra model complexity.

Adversarial learning. Our A-CAQ presents a multi-task learning problem. Here, we analyze the effects of two optimization tasks individually, which are reported in Tab. 6. Quantization without any post-processing results in the worst accuracy drop with the highest computational consumption (model **(2a)**). Post-training can, in turn, alleviate the performance degradation (model **(2b)**). Minimizing \mathcal{L}^{bit} allows for learning the layer-wise quantization bitwidth, which searches for the lowest bitwidth in the PTQ performance space (model **(2c)**).

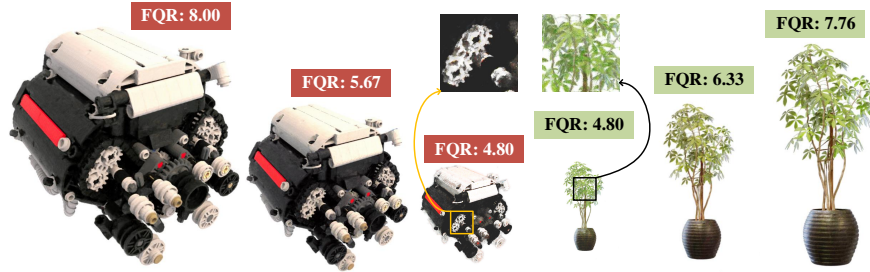


Fig. 6: Quantized LOD. Our proposed method can compress scenes in LOD style. The detail levels can be easily manipulated considering available resources and required accuracy.

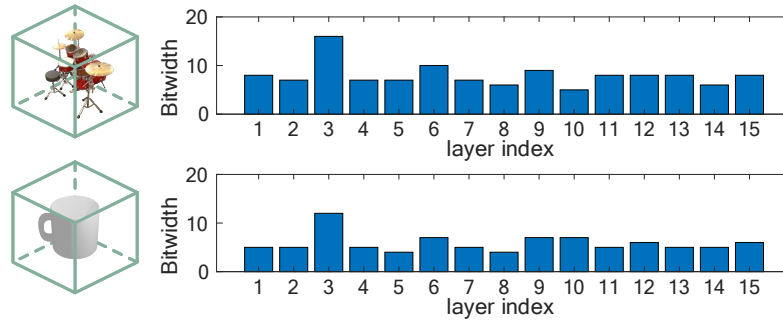


Fig. 7: Visualization of content-aware layer-wise quantization. The total number of components in our radiance field pipeline $M = 15$ including codebook and MLPs.

Our proposed method can achieve much lower bitwidths according to the required accuracy as we search for bitwidth in the QAT performance space.

4.5 Complexity analysis

Our proposed quantization method can determine the optimal content-dependent bitwidth for both high-fidelity (MDL) and resource-constrained (MGL) scenarios. To demonstrate how bitwidth influences resource overhead, we use the number of operations, weighted by bitwidth (BitOps) [36], involved in multiplications as a metric for computational complexity which has been widely employed for mixed-precision network systems [16, 30, 39]. Additionally, model sizes are listed to reflect memory consumption.

The results presented in Tab. 7 demonstrate the efficiency of our method. In the MDL scenario, our framework achieves a reduction of $\sim 90.78\%$ in BitOps compared to the baseline, and $\sim 11.66\%$ reduction compared to quantization with LSQ+. Regarding memory consumption, our proposed A-CAQ achieves a

Table 5: Ablation study on layer-wise and scene-wise quantization.

	layer	scene	FQR _↓	PSNR _↑
(1a)	✗	✗	14.00	31.56
(1b)	✗	✓	14.27	32.04
(LSQ+, 1c)	✓	✗	9.60	32.11
(Ours, 1d)	✓	✗	8.93	32.30
(Ours, 1e)	✓	✓	7.76	32.00

Table 6: Ablation study on the adversarial losses: $\mathcal{L}^{\text{NeRF}}$ and \mathcal{L}^{bit} .

	loss		MDL		MGL	
	$\mathcal{L}^{\text{NeRF}}$	\mathcal{L}^{bit}	FQR _↓	PSNR _↑	FQR _↓	PSNR _↑
(2a)	✗	✗	9.60	31.98	6.60	22.29
(2b)	✓	✗	9.60	32.11	6.60	25.06
(2c)	✗	✓	8.07	32.00	6.43	27.95
Ours	✓	✓	7.76	32.00	5.33	27.58

Table 7: Complexity analysis of A-CAQ for inference. BitOps [36] for rendering one 800×800 image and model storage are measured for time and space consumption, respectively.

	MDL			MGL		
	PSNR _↑	BitOps _↓ [T]	Storage _↓ [MB]	PSNR _↑	BitOps _↓	Storage _↓
NGP [24]	32.42	71.01	46.56	-	-	-
NGP-LSQ+ [3,40,43]	32.11	7.41	11.64	25.06	4.70	7.28
NGP-A-CAQ (Ours)	32.00	6.55	10.39	25.96	3.94	4.92

reduction of $\sim 77.68\%$ compared to the baseline, and $\sim 10.74\%$ reduction compared to LSQ+. In the MGL scenario, the reduction expands to $\sim 94.45\%$ and $\sim 16.17\%$ in BitOps compared to NGP and NGP-LSQ+ respectively. Moreover, there is a reduction of $\sim 89.43\%$ and $\sim 32.34\%$ in storage compared to NGP and NGP-LSQ+ respectively. Notably, our method also achieves a 0.9 dB increase in PSNR compared to NGP-LSQ+.

5 Conclusion

In this work, we introduce the concept of content-aware radiance field and explore the relationship between 3D scene complexity and quantization schemes. This motivates us to propose the A-CAQ algorithm, which learns bitwidth using gradients obtained by automatically perceiving the scene. Our method dynamically allocates layer-wise and scene-wise bitwidths. Experimental results demonstrate that the proposed algorithm significantly reduces model complexity for various scenarios through quantization, under different requirements.

Limitations. This study explores the novel concept of content-aware radiance fields, deftly integrating mixed-precision quantization into its framework. Despite the progress made, the realm of content-aware radiance fields beckons with unexplored territories ripe for detailed examination. For instance, by conducting a targeted search of network architecture aligned with reliable indicators that reflect the complexity of 3D scenes, it becomes feasible to construct a more comprehensive content-aware radiance field framework.

Acknowledgements

This work was supported by the Central Guided Local Science and Technology Foundation of China (YDZX20223100001001).

References

1. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. In: CVPR. pp. 5855–5864 (2022). <https://doi.org/10.1109/CVPR52688.2022.00539> 10
2. Bengio, Y., Léonard, N., Courville, A.: Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation. arXiv preprint arXiv:1308.3432 (2013) 7
3. Bhalgat, Y., Lee, J., Nagel, M., Blankevoort, T., Kwak, N.: LSQ+: Improving low-bit quantization through learnable offsets and better initialization. In: CVPR. pp. 696–697 (2020). <https://doi.org/10.1109/CVPRW50498.2020.00356> 6, 7, 11, 14
4. Cai, Z., He, X., Sun, J., Vasconcelos, N.: Deep Learning with Low Precision by Half-Wave Gaussian Quantization. In: CVPR. pp. 5918–5926 (2017). <https://doi.org/10.1109/CVPR.2017.574> 4
5. Cai, Z., Vasconcelos, N.: Rethinking Differentiable Search for Mixed-Precision Neural Networks. In: CVPR. pp. 2349–2358 (2020). <https://doi.org/10.1109/CVPR42600.2020.00242> 5
6. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: TensorRF: Tensorial Radiance Fields. In: ECCV. pp. 333–350 (2022). https://doi.org/10.1007/978-3-031-19824-3_20 4
7. Choi, Y., El-Khamy, M., Lee, J.: Towards the Limit of Network Quantization. arXiv preprint arXiv:1612.01543 (2016) 4
8. Choukroun, Y., Kravchik, E., Yang, F., Kisilev, P.: Low-bit Quantization of Neural Networks for Efficient Inference. In: 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW). pp. 3009–3018 (2019). <https://doi.org/10.1109/ICCVW.2019.00363> 5
9. Deng, C.L., Tartaglione, E.: Compressing Explicit Voxel Grid Representations: fast NeRFs become also small. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1236–1245 (2023). <https://doi.org/10.1109/WACV56688.2023.00129> 4
10. Esser, S.K., McKinstry, J.L., Bablani, D., Appuswamy, R., Modha, D.S.: Learned Step Size Quantization. In: ICLR (2020) 6, 7
11. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR. pp. 5501–5510 (2022). <https://doi.org/10.1109/CVPR52688.2022.00542> 4
12. Garbin, S.J., Kowalski, M., Johnson, M., Shotton, J., Valentin, J.: FastNeRF: High-Fidelity Neural Rendering at 200FPS. In: ICCV. pp. 14346–14355 (2021). <https://doi.org/10.1109/ICCV48922.2021.01408> 2
13. Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M.W., Keutzer, K.: A Survey of Quantization Methods for Efficient Neural Network Inference. In: Low-Power Computer Vision, pp. 291–326. Chapman and Hall/CRC (2022) 4
14. Girish, S., Shrivastava, A., Gupta, K.: SHACIRA: Scalable HAsH-grid Compression for Implicit Neural Representations. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 17513–17524 (2023). <https://doi.org/10.1109/ICCV51070.2023.01606> 4

15. Guo, L., Fei, W., Dai, W., Li, C., Zou, J., Xiong, H.: Mixed-Precision Quantization of U-Net for Medical Image Segmentation. In: 2022 IEEE International Symposium on Circuits and Systems (ISCAS). pp. 2871–2875 (2022). <https://doi.org/10.1109/ISCAS48785.2022.9937283> 5
16. Guo, Z., Zhang, X., Mu, H., Heng, W., Liu, Z., Wei, Y., Sun, J.: Single Path One-Shot Neural Architecture Search with Uniform Sampling. In: ECCV. pp. 544–560 (2020). https://doi.org/10.1007/978-3-030-58517-4_32 13
17. Hong, C., Baik, S., Kim, H., Nah, S., Lee, K.M.: CADyQ: Content-Aware Dynamic Quantization for Image Super-Resolution. In: ECCV. pp. 367–383 (2022). https://doi.org/10.1007/978-3-031-20071-7_22 5, 9
18. Hubara, I., Nahshan, Y., Hanani, Y., Banner, R., Soudry, D.: Accurate Post Training Quantization with Small Calibration Sets. In: International Conference on Machine Learning. pp. 4466–4475. PMLR (2021) 5
19. Jin, Q., Yang, L., Liao, Z.: AdaBits: Neural Network Quantization With Adaptive Bit-Widths. In: CVPR. pp. 2146–2156 (2020). <https://doi.org/10.1109/CVPR42600.2020.00222> 5
20. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* **42**(4) (2023). <https://doi.org/10.1145/3592433> 8
21. Li, L., Shen, Z., Wang, Z., Shen, L., Bo, L.: Compressing Volumetric Radiance Fields to 1 MB. In: CVPR. pp. 4222–4231 (2023). <https://doi.org/10.1109/CVPR52729.2023.00411> 4
22. Luo, H., Chen, A., Zhang, Q., Pang, B., Wu, M., Xu, L., Yu, J.: Convolutional Neural Opacity Radiance Fields. In: 2021 IEEE International Conference on Computational Photography (ICCP). pp. 1–12 (2021). <https://doi.org/10.1109/ICCP51581.2021.9466273> 2
23. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In: ECCV (2020). https://doi.org/10.1007/978-3-030-58452-8_24 1, 8, 10
24. Müller, T., Evans, A., Schied, C., Keller, A.: Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* **41**(4) (2022). <https://doi.org/10.1145/3528223.3530127> 4, 6, 14
25. Nagel, M., Fournarakis, M., Amjad, R.A., Bondarenko, Y., Van Baalen, M., Blankevoort, T.: A White Paper on Neural Network Quantization. arXiv preprint arXiv:2106.08295 (2021) 7
26. Polino, A., Pascanu, R., Alistarh, D.: Model Compression via Distillation and Quantization. arXiv preprint arXiv:1802.05668 (2018). <https://doi.org/10.48550/arXiv.1802.05668> 4
27. Rao, C., Yu, H., Wan, H., Zhou, J., Zheng, Y., Wu, M., Ma, Y., Chen, A., Yuan, B., Zhou, P., Lou, X., Yu, J.: ICARUS: A Specialized Architecture for Neural Radiance Fields Rendering. *ACM Trans. Graph.* **41**(6), 1–14 (2022). <https://doi.org/10.1145/3550454.3555505> 2, 11
28. Reiser, C., Peng, S., Liao, Y., Geiger, A.: KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14335–14345 (2021). <https://doi.org/10.1109/ICCV48922.2021.01407> 2
29. Takikawa, T., Evans, A., Tremblay, J., Müller, T., McGuire, M., Jacobson, A., Fidler, S.: Variable Bitrate Neural Fields. In: ACM SIGGRAPH 2022 Conference Proceedings. pp. 1–9 (2022). <https://doi.org/10.1145/3528233.3530727> 4, 11

30. Tang, C., Ouyang, K., Wang, Z., Zhu, Y., Ji, W., Wang, Y., Zhu, W.: Mixed-Precision Neural Network Quantization via Learned Layer-Wise Importance. In: ECCV. pp. 259–275 (2022). https://doi.org/10.1007/978-3-031-20083-0_16 13
31. Tang, J.: Torch-ngp: a PyTorch implementation of instant-ngp (2022), <https://github.com/ashawkey/torch-ngp> 9
32. Tremblay, J., Meshry, M., Evans, A., Kautz, J., Keller, A., Khamis, S., Müller, T., Loop, C., Morrical, N., Nagano, K., Takikawa, T., Birchfield, S.: RTMV: A Ray-Traced Multi-View Synthetic Dataset for Novel View Synthesis. arXiv preprint arXiv:2205.07058 (2022) 10
33. Wadhvani, K., Kojima, T.: Squeezenerf: Further factorized fastnerf for memory-efficient inference. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2717–2725 (2022). <https://doi.org/10.1109/CVPRW56347.2022.00307> 2
34. Wang, K., Liu, Z., Lin, Y., Lin, J., Han, S.: HAQ: Hardware-Aware Automated Quantization With Mixed Precision. In: CVPR. pp. 8612–8620 (2019). <https://doi.org/10.1109/CVPR.2019.00881> 5
35. Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E.: Image Quality Assessment: From Error Visibility to Structural Similarity. IEEE Transactions on Image Processing 13(4), 600–612 (2004). <https://doi.org/10.1109/TIP.2003.819861> 10
36. Wu, B., Wang, Y., Zhang, P., Tian, Y., Vajda, P., Keutzer, K.: Mixed Precision Quantization of ConvNets via Differentiable Neural Architecture Search. arXiv preprint arXiv:1812.00090 (2018) 13, 14
37. Yang, G.W., Zhou, W.Y., Peng, H.Y., Liang, D., Mu, T.J., Hu, S.M.: Recursive-NeRF: An Efficient and Dynamically Growing NeRF. IEEE Transactions on Visualization and Computer Graphics (2022). <https://doi.org/10.1109/TVCG.2022.3204608> 4
38. Yang, J., Shen, X., Xing, J., Tian, X., Li, H., Deng, B., Huang, J., Hua, X.s.: Quantization Networks. In: CVPR (June 2019) 4
39. Yang, L., Jin, Q.: FracBits: Mixed Precision Quantization via Fractional Bit-Widths. In: Proceedings of the AAAI Conference on Artificial Intelligence. pp. 10612–10620 (2021). <https://doi.org/10.1609/aaai.v35i12.17269> 13
40. Ye, Z., Hu, Q., Zhao, T., Zhou, W., Cheng, J.: MCUNeRF: Packing NeRF into an MCU with 1MB Memory. In: Proceedings of the 31st ACM International Conference on Multimedia. pp. 9082–9092 (2023). <https://doi.org/10.1145/3581783.3612109> 2, 11, 14
41. Youn, J., Song, J., Kim, H.S., Bahk, S.: Bitwidth-Adaptive Quantization-Aware Neural Network Training: A Meta-Learning Approach. In: ECCV. pp. 208–224 (2022). https://doi.org/10.1007/978-3-031-19775-8_13 5
42. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In: CVPR. pp. 586–595 (2018). <https://doi.org/10.1109/CVPR.2018.00068> 10
43. Zhao, T., Chen, J., Leng, C., Cheng, J.: TinyNeRF: Towards 100x Compression of Voxel Radiance Fields. In: Proceedings of the AAAI Conference on Artificial Intelligence. pp. 3588–3596 (2023). <https://doi.org/10.1609/aaai.v37i3.25469> 11, 14

A Additional details of LBQ and A-CAQ

A.1 Gradients for differentiable quantization

We introduce three quantization schemes for three different components within the radiance field pipeline in our proposed LBQ framework. To efficiently back-propagate through the simulated quantizer block, derivatives to three sets of key parameters are derived as follows

Neural weights. The weights are quantized with signed symmetric quantization, and the partial derivatives to the key parameters r_v and b are given as

$$\frac{\partial \hat{v}}{\partial r_v} = \begin{cases} \frac{1}{r_v} (s \cdot \lfloor v/s \rfloor - v) & v_{\min} \leq v \leq v_{\max} \\ \frac{q_{\max}}{r_q} & v > v_{\max} \\ \frac{q_{\min}}{r_q} & v < v_{\min} \end{cases}, \quad (13)$$

$$\frac{\partial \hat{v}}{\partial b} = \begin{cases} \frac{2^B \ln 2}{r_q} (v - s \cdot \lfloor v/s \rfloor) & v_{\min} \leq v \leq v_{\max} \\ \frac{2^B \ln 2}{r_q} (v_{\max} - q_{\max} \cdot s) & v > v_{\max} \\ \frac{2^B \ln 2}{r_q} (v_{\min} - q_{\min} \cdot s) & v < v_{\min} \end{cases}. \quad (14)$$

ReLU and exponential activations. The ReLU and exponential activations are quantized with unsigned symmetric quantization, the partial derivatives are derived as

$$\frac{\partial \hat{v}}{\partial r_v} = \begin{cases} \frac{1}{r_v} (s \cdot \lfloor v/s \rfloor - v) & v_{\min} \leq v \leq v_{\max} \\ 1 & v > v_{\max} \\ 0 & v < v_{\min} \end{cases}, \quad (15)$$

$$\frac{\partial \hat{v}}{\partial b} = \begin{cases} \frac{2^B \ln 2}{r_q} (v - s \cdot \lfloor v/s \rfloor) & v_{\min} \leq v \leq v_{\max} \\ 0 & \text{otherwise} \end{cases}. \quad (16)$$

PE and others. PE and other components are quantized with asymmetric quantization which is regarded as a general form. Considering the offset is equivalently replace by v_{\max} , the derivatives w.r.t. key parameters are given as

$$\frac{\partial \hat{v}}{\partial v_{\max}} = \begin{cases} 0 & v_{\min} \leq v \leq v_{\max} \\ 1 & \text{otherwise} \end{cases}, \quad (17)$$

$$\frac{\partial \hat{v}}{\partial r_v} = \begin{cases} \frac{1}{r_v} (s \cdot \lfloor v/s \rfloor - v) & v_{\min} \leq v \leq v_{\max} \\ 1 - \frac{v_{\max}}{r_v} - \frac{z}{r_q} & v > v_{\max} \\ -\frac{v_{\max}}{r_v} - \frac{z}{r_q} & v < v_{\min} \end{cases}, \quad (18)$$

$$\frac{\partial \hat{v}}{\partial b} = \begin{cases} \frac{2^B \ln 2}{r_q} (v - s \cdot \lfloor v/s \rfloor) & v_{\min} \leq v \leq v_{\max} \\ \frac{2^B \ln 2}{r_q} (v_{\min} + s \cdot z) & \text{otherwise} \end{cases}. \quad (19)$$

A.2 Pseudo code for A-CAQ

Our proposed A-CAQ is a multi-task learning-based algorithm, and the pseudo codes are shown in Algorithm 1.

Algorithm 1 A-CAQ

Input: Trained model with initial quantization parameters $\mathcal{Q}^{0,0} = \{\Omega^{0,0}, \mathbf{v}_{\max}^{0,0}, \mathbf{r}_v^{0,0}\}$; initial soft bitwidth $\mathbf{b}^{0,0}$; training set \mathcal{D} ; learning rate β ; hyper-defined metric $\mathcal{L}^{\text{metric}}$

Output: The convergent model parameters $\mathcal{Q}^{E-1,I}, \mathbf{b}^{E-1,I}$

```

for epoch  $i = 0$  to  $E - 1$  do
  for iteration  $j = 1$  to  $I$  do
     $\mathcal{B}_j \leftarrow \mathcal{D}$  ▷ Get Batch using ray-marching
     $(\sigma, \mathbf{c}) \leftarrow F(\mathbf{x}, \mathbf{d} \mid \mathcal{Q}^{i,j-1}, b^{i,j-1})$  for all  $(\mathbf{x}, \mathbf{d}) \in \mathcal{B}_j$ 
     $\hat{C}(\ell) \leftarrow$  volume rendering for all  $\ell \in \mathcal{R}$ 
     $\nabla_j^{(1)} \leftarrow \nabla_{\mathcal{Q}^{i,j-1}} \mathcal{L}^{\text{NeRF}}(\hat{C}(\ell) \mid \mathcal{Q}^{i,j-1}, b^{i,j-1})$ 
     $\nabla_j^{(2)} \leftarrow \nabla_{b^{i,j-1}} \mathcal{L}^{\text{bit}}(\hat{C}(\ell) \mid \mathcal{Q}^{i,j-1}, b^{i,j-1})$ 
     $\mathcal{Q}^{i,j} \leftarrow \mathcal{Q}^{i,j-1} - \frac{\beta}{|\mathcal{B}_j|} \nabla_j^{(1)}$ 
     $\mathbf{b}^{i,j} \leftarrow \mathbf{b}^{i,j-1} - \frac{\beta}{|\mathcal{B}_j|} \nabla_j^{(2)}$  ▷ Update with gradients averaged on batch size
  end for
end for

```

B More Implementation Details

Given that our framework is optimized based on a well-trained full-precision model, the learning rates for the two tasks are differently configured. Specifically, as the bitwidths are trained from scratch, the initial bitwidth learning rates are set to 1×10^{-2} , while those for other parameters are 1×10^{-3} . The summation of bitwidth penalties is fixed to $\sum_i \epsilon_i = 1 \times 10^{-3}$. And training of A-CAQ finishes after 3000 iterations. Other training configurations, such as batch size, learning rate schedule and optimizer, follow the settings of the baseline model.

C Details on penalty

The summation of bitwidth penalties serves as a per-defined parameters governs the balance between the accuracy and efficiency. To illustrate its impact, experimental results on Synthetic-NeRF dataset are shown in Fig. 8. Lower penalties result in higher accuracy, albeit at the expense of efficiency. Moreover, our proposed A-CAQ outperforms LSQ+ in terms of both accuracy and efficiency, offering a straightforward means to achieve higher levels of both metrics.

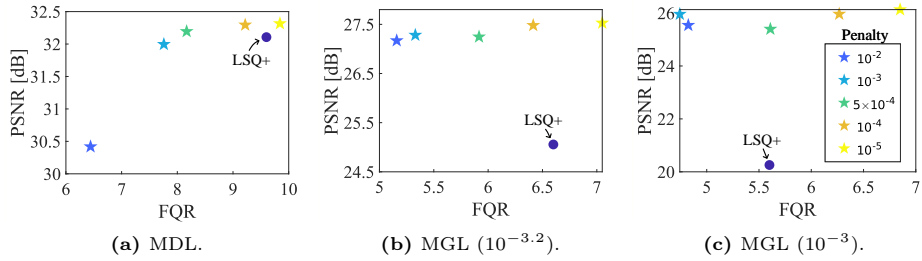


Fig. 8: Penalty effects on A-CAQ.

D More qualitative results

Fig. 9 and Fig. 10 provide additional qualitative results in Mip-NeRF360 and Synthetic-MeRF dataset.



Fig. 9: Rendering results on Mip-NeRF360 dataset.



Fig. 10: Rendering results on Synthetic-NeRF dataset.