# LOW-BIT QUANTIZATION FOR SEEING IN THE DARK

#### **Anonymous authors**

Paper under double-blind review

#### ABSTRACT

009 Several properties of raw data exhibit significant potential for enhancing images under extremely low-light conditions. Recently, many deep-learning methods for raw-based low-light image enhancement (LLIE) have demonstrated excellent performance. However, deploying them on resource-limited devices is restricted by high computational and storage demands. In this work, we propose a novel low-bit quantization method for raw-based LLIE model to improve their efficiency. Nevertheless, directly adopting existing quantizers for LLIE networks leads to 016 obvious performance drop due to two main reasons. *i*) The U-Net model, commonly employed in LLIE, faces challenges in identifying a suitable quantization range due to disparities in distribution between the encoder and decoder features. *ii*) Low-bit quantized LLIE networks struggle to restore clear details in low-light images because their features have a constraint capacity. We address these issues by introducing a novel low-bit quantization method, the Distribution-Separative Asymmetric Quantizer (DSAQ), designed specifically for U-Net architectures used in LLIE. In order to accurately determine the quantization intervals, DSAQ separates the distribution of encoder and decoder features before they are concatenated by the skip connection. We also make the quantizer asymmetric with trainable scale and offset parameters to suit skewed activation ranges caused by non-linear functions. To further enhance performance, we propose a uniform feature distillation technique, which allows the low-bit student model to effectively assimilate knowledge from the full-precision teacher model, bridging the gap in representation capability. Extensive experiments show that our approach not only greatly reduces the memory and computational requirements of raw-based LLIE models but also has a promising performance. Our low-bit quantized model can achieve comparable or superior results to full-precision counterparts.

032 033 034

000

001 002 003

004

006

008

010

011

012

013

014

015

017

018

019

021

024

025

026

028

029

031

#### 1 INTRODUCTION

Capturing high-quality images under extremely low-light conditions is important for night surveillance 037 and various downstream computer vision tasks Hong et al. (2021); Chen et al. (2023). However, the degradations including low signal-to-noise ratio (SNR) and obvious color cast caused by limited photon count make it challenging. Low-light image enhancement (LLIE) methods provide a postcapture solution that prevents noise amplification at high ISO settings and motion blur due to long 040 exposure time. 041

042 Recently, deep-learning LLIE methods trained with paired datasets Chen et al. (2018); Dong et al. 043 (2022); Wei et al. (2018) have shown outstanding performance. In this work, we focus on enhancing 044 low-light raw images because of their inherent advantages for LLIE Wei et al. (2022); Huang et al. (2022). On the one hand, they maintain a linear relationship with photon counts and a tractable noise distribution before passing through the image signal processing (ISP) pipeline. On the other hand, 046 they have a higher bit-depth that can distinguish subtle low-intensity details. Although deep-learning 047 models for raw-based LLIE Chen et al. (2018); Zhu et al. (2020); Jin et al. (2023) have achieved 048 promising results, the deployment of these neural networks on edge devices like mobile phones or embedded cameras is hindered by their high computational and storage demands.

A potential solution to this problem lies in the technique of network quantization Zhou et al. (2016); Li et al. (2020); Qin et al. (2023). Quantization involves converting the continuous weights and 052 activations (features) of a neural network into discrete low-bit representations, significantly reducing the model's memory footprint and accelerating its inference speed. In this work, we intend to quantize

the LLIE model to a range of 2-4 bits to achieve a higher compression ratio. Despite its benefits in terms of efficiency, low-bit network quantization may lead to a deterioration in model performance.

For raw-based LLIE methods, we recognize two main reasons for this performance degradation. 057 Firstly, many of raw LLIE networks are based on U-Net Ronneberger et al. (2015) structure Chen et al. (2018); Dong et al. (2022); Huang et al. (2022). In the U-Net architecture, we observe that the features from the encoder and decoder, concatenated through skip connections, show notable 060 differences in distribution. Moreover, the use of non-linear activation functions, such as LeakyReLU 061 Chen et al. (2018); Lamba & Mitra (2021), results in asymmetric distributions of positive and negative 062 values. These factors pose challenges in accurately determining the quantization range. Secondly, 063 features in low-bit quantized networks exhibit a representation capability gap compared to those in 064 full-precision networks. Therefore, existing knowledge distillation schemes Li et al. (2020); Zhong et al. (2022), which directly impose constraints on normalized features, cannot fully transfer intrinsic 065 semantic information from full-precision teacher model to the low-bit quantized student model. 066

067 In this paper, we present a novel low-bit quantization method for raw-based LLIE to solve the above 068 problems. Specifically, we propose a Distribution-Separative Asymmetric Quantizer (DSAQ) that is 069 tailored for U-Net based LLIE method. It quantizes the encoder and decoder features respectively before concatenation to facilitate the learning of quantization interval. In order to mitigate the 071 influence of non-linear functions on the distribution of activations, we introduce trainable scale and offset parameters to implement the asymmetric quantizer. We further introduce a uniform feature 072 distillation that maps features of quantized student model and full-precision teacher model into 073 a uniform latent feature space. So low-bit network can better obtain intrinsic information from 074 teacher model and restore clearer details from low-light images. Through extensive experiments, we 075 demonstrate that our quantization method surpasses previous quantizers in raw LLIE and the low-bit 076 network achieves comparable enhancement results to their full-precision counterparts. Our main 077 contributions can be summarized as follows:

- We propose a compact low-bit quantized model for low-light raw image enhancement, which can achieve satisfactory results with low memory and computation.
- We build a Distribution-Separative Asymmetric Quantizer (DSAQ) for U-Net structure. It separately determines the quantizer of different features before concatenation and introduces asymmetric quantization for activations with skewed distribution.
  - We design a uniformed feature distillation that reduces capacity difference between features in quantized and full-precision models. So the knowledge from teacher model can be easily transferred to student model.
- 2 RELATED WORK

079

081

082

084

085

087

089 090

091

092

In this section, we first review deep-learning methods for raw-based LLIE. We then review some quantization techniques for efficient neural network inference.

093 094 2.1 RAW-BASED LOW-LIGHT IMAGE ENHANCEMENT

Because of the merits of raw images discussed in Section 1, they are commonly used for LLIE in
extremely dark environments. The pioneering work Chen et al. (2018) builds a large-scale paired
short/long exposure raw image dataset for LLIE, dubbed See-in-the-Dark (SID). A U-Net is employed
for restoring noisy low-light raw input into bright RGB images. A parallel work DeepISP Schwartz
et al. (2019) uses an end-to-end neural network to process low-light raw images, which achieves
better visual quality than manufactured ISP. Based on the SID dataset, following work also introduces
residual learning Maharjan et al. (2019), self-guidance strategy Gu et al. (2019) and multi-criterion
loss Zamir et al. (2021) for single-stage raw to RGB LLIE.

Another line of methods decompose the problem of raw-based LLIE into different aspects and design
 multi-stage networks. EEMEFN Zhu et al. (2020) sequentially performs multi-exposure fusion and
 edge enhancement for LLIE. LDC Xu et al. (2020) enhances the low-frequency part and reconstructs
 the high-frequency details of the low-light images in two consecutive stages. MCR Dong et al.
 (2022) first learns to synthesize monochrome images with additional supervision. Then a dual-branch
 network is leveraged to fuse generated monochrome and color images to produce enhanced RGB

results. Huang *et al.* Huang et al. (2022) proposes a raw-guiding exposure enhancement network, which consist of three cascaded U-Nets for unprocessing, denoising and processing. DNF Jin et al. (2023) decouples raw-based LLIE into raw image denoising stage and RGB image color correction stage to mitigate the domain ambiguity. A feedback module enables feature interaction across two stages to reduce error accumulation.

The power of these neural networks to see in the dark relies on their model depth and computational complexity. Some work Lamba & Mitra (2021); Lamba et al. (2020) also improves efficiency of LLIE models by designing lightweight network architectures. In this work, we resort to network quantization to achieve efficient LLIE.

117 118

119

#### 2.2 NEURAL NETWORK QUANTIZATION

120 Neural network quantization involves reducing the precision of weights and activations in a neural network, representing them with a lower-bit (usually 2-8 bits) discrete representation Nagel et al. 121 (2021). This process can effectively reduce the model size and computation cost, and it can be 122 incorporated with other network compression techniques like parameter pruning Zhang et al. (2022); 123 Wang & Fu (2023) and knowledge distillation Zhu et al. (2023); Li et al. (2020). There are two 124 primary paradigms for network quantization: Post-Training Quantization (PTQ) Hubara et al. (2021); 125 Li et al. (2021) and Quantization Aware Training (QAT) Zhou et al. (2016); Choi et al. (2018); Li 126 et al. (2022; 2023); Esser et al. (2020). PTQ methods allow for the efficient quantization of pre-127 trained neural networks with minimal data and no retraining. However, they suffer from sub-optimal 128 performance due to fixed parameters and limited fine-tuning capabilities. In this paper, we adopt 129 QAT that retrain the network parameters with simulated quantization and full training data to achieve 130 a better performance in low-bit (*i.e.*, less than 4 bits) quantization. Additionally, 1-bit quantization 131 methods (also known as binary neural networks) are not discussed because they often rely on specific designs to avoid severe performance degradation Liu et al. (2018); Cai et al. (2023) and have different 132 hardware implementations Qin et al. (2023). 133

134 Low-bit network quantization with OAT has been widely applied to various computer vision tasks, 135 with early efforts primarily focusing on the quantization of classification models Choi et al. (2018); 136 Jacob et al. (2018); Gong et al. (2019); Jung et al. (2019); Esser et al. (2020); Bhalgat et al. (2020). 137 These methods incorporate either a learnable quantization interval Choi et al. (2018) or a learnable step size Esser et al. (2020) within the quantizer, optimizing these parameters along with network 138 weights to minimize task-specific loss Jung et al. (2019). In low-level vision, much work has explored 139 low-bit quantization for super-resolution networks, which typically consist of a head, main body, 140 and upsample tail Qin et al. (2023); Li et al. (2020); Wang et al. (2021). PAMS Li et al. (2020) 141 introduces a trainable clamp function and proposes a structured knowledge transfer strategy, enabling 142 the learning of high-level representations from the full-precision model. FQSR Wang et al. (2021) 143 fully quantizes all the layers including head and upsample tail in super-resolution networks. DDTB 144 Zhong et al. (2022) adopts trainable upper and lower bounds for the highly asymmetric activations. 145 DAQ Hong et al. (2022) uses a distribution-aware quantization that defines a quantize function for 146 each channel. QuantSR Qin et al. (2023) leverages a redistribution-driven learnable quantizer to 147 diversify the low-bit quantized representation. A depth-dynamic quantized architecture is designed to 148 achieve resource adaptive inference. In this work, we aim to quantize a U-Net-based model, which is widely used in raw LLIE. 149

150 151

152

# 3 Method

In this section, we first provide an overview of the process of low-bit quantization for U-Net style
 raw-based LLIE networks, along with the limitations of existing quantizers. Then, we present our
 Distribution-Separative Asymmetric Quantizer (DSAQ), which is specifically designed for U-Net structured LLIE models. Finally, we introduce uniform feature distillation for low-bit quantization.

157

158 3.1 LOW-BIT QUANTIZED LLIE U-NET 159

160 **Overall Network Architecture.** We utilize the U-Net architecture in SID Chen et al. (2018) as 161 the full-precision model and follow the same pipeline to process raw data. Given a low-light raw 162 image  $I^B \in \mathbb{R}^{H \times W}$  in Bayer array format, we pack each  $2 \times 2$  pattern into four channels to ensure



Figure 1: The architecture of our low-bit quantized LLIE model. The overall U-Net structure is illustrated on the left. The details of the distribution-separative asymmetric quantizer (DSAQ) and the uniform feature distillation are shown on the right.

each channel represents the same color. The packed raw image, denoted as  $I^P \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times 4}$ , is then multiplied with a pre-defined amplification ratio r. Finally, the packed and amplified image is fed into the network to restore a bright and clean RGB output image  $O \in \mathbb{R}^{H \times W \times 3}$ .

The overall U-Net architecture is shown in Figure 1, it contains four levels of encoders and decoders. 191 The encoder features are concatenated with upsampled decoder features from the previous level 192 through skip connection. Convolution blocks of the encoders and decoders consist of two convolutions 193 and are activated with LeakyReLU non-linear function. We apply quantization to all convolutional 194 layers in the encoders and decoders, except for the first and last convolutional layers, which are kept 195 at full precision. This approach helps prevent information loss in the input raw images and ensures a 196 higher fidelity in the final enhanced images. We use maxpooling for downsampling and quantized 197 transposed convolution for upsampling. 198

Formulation of Network Quantization. The common quantization scheme first maps real-valued vectors in the network into integer representation, then performs a de-quantization step to approximate the original value. The quantizer  $Q^b$  can be formulated as

$$\hat{\boldsymbol{x}} = \mathbf{Q}^{b}(\boldsymbol{x}) = \left\lfloor \operatorname{clip}(\frac{\boldsymbol{x}}{s}, Q_{n}, Q_{p}) \right\rceil \times s, \tag{1}$$

where x represents full-precision weights (*e.g.* kernels in convolution) or activations (*e.g.* feature maps in convolution), s denotes the scaling factor that converts real values to the quantization range and de-quantizes integers back to original value range.  $Q_n$  and  $Q_p$  represent the quantization range that  $Q_n = 0, Q_p = 2^b - 1$  for unsigned quantizers and  $Q_n = -2^{b-1}, Q_p = 2^{b-1} - 1$  for signed quantizers, where b is the bit-width of the quantizer. The function  $\operatorname{clip}(\cdot, Q_n, Q_p)$  limits the scaled values to the quantization range, and  $\lfloor \cdot \rceil$  rounds the real value to its nearest integer.  $\hat{x}$  is the low-bit discrete representation of the full-precision vector x.

211

202 203

183

185

186

## 212 3.2 DISTRIBUTION-SEPARATIVE ASYMMETRIC QUANTIZER

213

5.2 DISTRIBUTION SELARATIVE ASTMINETRIC QUANTLER

Asymmetric Activation Quantizer. The symmetric quantizer defined in Equation 1 allocates an
 equal number of bins for both positive and negative values. Despite its efficiency, this approach may
 exhibit suboptimal suitability for vectors with asymmetric distributions. In the LLIE model, two

230

231 232 233

234

235

236

237

238

239 240

241

249 250 251



Figure 2: (a) Skewness of activations and weights in each quantized convolution layer. (b) Unfitness of symmetric quantizer for asymmetric-distributed activations.

main factors result in the asymmetric distribution of activations. First, batch normalization layers are often removed in LLIE networks because they smooth the features, resulting in blurred enhancement images Li et al. (2020). Second, LeakyReLU is commonly used as the activation function Chen et al. (2018); Lamba & Mitra (2021), which compresses the range of negative values in features.

We analyze the skewness to measure the asymmetry of the activations and weights in quantized convolutions. The skewness of a vector  $\boldsymbol{x}$  with n values can be estimated by Joanes & Gill (1998)

Skewness
$$(\boldsymbol{x}) = \frac{\frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^3}{\sigma^3},$$
 (2)

where  $\mu$  is the sample mean and  $\sigma$  is the sample standard deviation. As illustrated in Figure 2(a), the skewness of convolution kernel weights is near zero as they follow a symmetric bell-shaped distribution Hong et al. (2022). However, a large positive skewness of activations shows their distributions are right-skewed. As shown in Figure 2(b), some quantization bins in the negative range are wasted when applying symmetric quantizer for right-skewed activations. For this reason, we use the symmetric quantizer defined in Equation 1 for the weights and the asymmetric quantizer for the activations. The asymmetric quantizer AQ<sup>b</sup> is defined as

$$\hat{\boldsymbol{a}} = \mathrm{AQ}^{b}(\boldsymbol{a}) = \left\lfloor \mathrm{clip}(\frac{\boldsymbol{a} - \beta}{s}, 0, 2^{b} - 1) \right\rfloor \times s + \beta, \tag{3}$$

where  $\hat{a}$  is the low-bit representation for asymmetric activation a, s and  $\beta$  are learnable parameters that represent scaling factor and offset respectively. We first initialize the scaling factor sto  $2\text{mean}(|a|)/\sqrt{Q_p}$ , which is calculated from the first batch of activations, then the offset  $\beta$  is initialized to  $sQ_n$ . In order to preserve more information during back-propagation, we utilize the soft gradient transformation function Qin et al. (2023) instead of the straight-through estimation (STE).

257 Distribution-Separative Quantization. The key of preserving the performance of full-precision 258 after quantization is to find a proper scaling factor s for each activation. Existing QAT methods Choi 259 et al. (2018); Esser et al. (2020) treat s as a learnable parameter and jointly optimize it with network 260 weights. However, we observe that the distribution range of features concatenated through the skip 261 connection exhibits significant differences. As illustrated in Figure 3, features upsampled from the former decoder have a larger value range than the encoder features. Therefore, learning a single 262 scaling factor for the concatenated features may lead to quantization unfitness for the activations. 263 Figure 3 shows two typical situations of the quantization unfitness in the quantizer. In the first row, the 264 quantizer learns a small scaling factor and the activations with large absolute values are scaled out of 265 the quantization range and clipped by the  $clip(\cdot, Q_n, Q_p)$  function, which causes the information loss 266 of decoder features. In the second row, the quantizer learns a large scaling factor and the activations 267 with small absolute values are quantized to zero, which leads to information loss of encoder features. 268

In order to preserve the information of both encoder and decoder features in the skip connection, we propose a simple yet effective distribution-separative quantization. Specifically, let



Figure 3: Distribution of activations after concatenation through skip connection.

 $a \in \mathbb{R}^{H \times W \times 2C}$  be the concatenated feature through skip connection and  $a = \text{Concat}(a^u, a^e)$ , where  $a^u, a^e \in \mathbb{R}^{H \times W \times C}$  are upsampled decoder feature and encoder feature. Our distributionseparative asymmetric quantization (DSAQ) for activations is defined as

$$DSAQ^{b}(\boldsymbol{a}) = Concat(AQ_{1}^{b}(\boldsymbol{a}^{u}), AQ_{2}^{b}(\boldsymbol{a}^{e})),$$
(4)

where  $AQ_1^b$  and  $AQ_2^b$  are quantizers that learn two different sets of quantization parameters (*i.e.* scaling factors and offsets) for  $a^u$  and  $a^e$  respectively. Compared with channel-wise quantizers Hong et al. (2022) that learn parameters for each channel of the activations, our DSAQ is a more efficient approach as only one additional set of quantization parameters is introduced to handle the distribution mismatch.

#### 3.3 UNIFORM FEATURE DISTILLATION

303 Inspired by previous work Li et al. (2020); Zhong et al. (2022), incorporating network quantization 304 with knowledge distillation can achieve a better performance. Full-precision networks can learn 305 more representative features, which provide abundant details and high-level semantic information for low-bit quantized networks. The structured knowledge transfer Li et al. (2020); Zhong et al. 306 (2022) used in previous quantized super-resolution networks directly minimize pixel-wise distance of 307 normalized features from full-precision teacher model and low-bit student model. However, there is 308 great capability gap between features from quantized models and their full-precision counterparts, 309 which makes it challenging for low-bit features to mimic float-point features. Existing work Zhu et al. 310 (2023) also leverages the quantized feature from the full-precision model for knowledge distillation 311 in the classification task. Although it makes low-bit models easier to learn the feature representation, 312 quantizatied float-point features lose the detailed information for enhancing low-light images. In 313 this work, we propose a uniform feature distillation for feature alignment and knowledge transfer. 314 Specifically, we introduce a full-precision feature uniform module (FUM) to process features from 315 the quantized network, which can be excluded during inference. The FUM projects the low-precision 316 feature to a uniform space with the full-precision features and mitigates the capability disparity. 317 Therefore, our uniform feature distillation facilitates the knowledge transfer from the teacher model to the low-bit student model without losing essential details in the full-precision features, which is 318 represented as 319

320

289 290 291

292 293

295

301

302

321

$$L_{distill} = \|\frac{F'_{US}}{\|F'_{US}\|_2} - \frac{F'_T}{\|F'_T\|_2}\|_2,$$
(5)

where  $F_{US} = \text{FUM}(F_S)$  is the processed uniformed feature,  $F_S, F_T \in \mathbb{R}^{H \times W \times C}$  are the features of student model and teacher model,  $F' = \sum_{i=1}^{C} |F_i|^2 \in \mathbb{R}^{H \times W}$  represents the spatial mapping Li

Mathad	Bits	SID-Sony		MCR		Params	FLOPs
Wiethou	(w/a)	PSNR	SSIM	PSNR	SSIM	(M)	(G)
SID Chen et al. (2018)	32/32	29.02	0.7866	29.43	0.9076	7.76	48.45
LLPack Lamba et al. (2020)	32/32	27.76	0.7675	24.53	0.8240	1.17	7.21
RRT Lamba & Mitra (2021)	32/32	28.54	0.7743	26.17	0.8438	0.78	5.17
Dorefa Zhou et al. (2016)	4/4	27.80	0.7677	27.18	0.8745		
PACT Choi et al. (2018)	4/4	27.65	0.7634	25.32	0.8558		
PAMS Li et al. (2020)	4/4	28.03	0.7527	25.20	0.8291		
LSQ Esser et al. (2020)	4/4	28.62	0.7790	28.61	0.8925	0.97	6.51
LLT Wang et al. (2022)	4/4	24.54	0.7170	20.61	0.5887		
QuantSR Qin et al. (2023)	4/4	28.73	0.7814	28.64	0.8923		
Ours	4/4	28.81	0.7823	29.00	0.8987		
Dorefa Zhou et al. (2016)	3/3	27.48	0.7502	25.76	0.8479		
PACT Choi et al. (2018)	3/3	26.82	0.7324	24.21	0.8257		
PAMS Li et al. (2020)	3/3	27.35	0.7437	22.26	0.7669		
LSQ Esser et al. (2020)	3/3	28.33	0.7722	27.45	0.8756	0.73	3.64
LLT Wang et al. (2022)	3/3	20.87	0.5870	20.25	0.6970		
QuantSR Qin et al. (2023)	3/3	28.53	0.7741	27.60	0.8810		
Ours	3/3	28.66	0.7772	28.39	0.8866		
Dorefa Zhou et al. (2016)	2/2	26.50	0.7173	23.67	0.7768		
PACT Choi et al. (2018)	2/2	25.96	0.7069	21.83	0.7335		
PAMS Li et al. (2020)	2/2	23.57	0.6008	18.65	0.6584		
LSQ Esser et al. (2020)	2/2	27.79	0.7586	25.02	0.8197	0.49	2.2
LLT Wang et al. (2022)	2/2	17.74	0.5518	-	-		
QuantSR Qin et al. (2023)	2/2	28.10	0.7617	25.62	0.8413		
Ours	2/2	28.14	0.7637	26.00	0.8430		

Table 1: Quantitative results on SID dataset and MCR dataset. LLT Wang et al. (2022) fails to converge on the MCR dataset in 2-bit setting so the results are denoted by '-'.

et al. (2020). We choose the output feature from convolution block of the last decoder for distillation. The overall training loss is defined as

$$L = \lambda_1 L_1 + \lambda_2 L_{distill},\tag{6}$$

where  $\lambda_1, \lambda_2$  are hyperparameters and we set  $\lambda_1 = 1, \lambda_2 = 100$ .

#### 4 EXPERIMENTS

354

355

356

357

358 359 360

361 362

363 364

365

366 367 In this section, we evaluate our low-bit quantized U-Net network on two raw-based LLIE datasets. We also provide a comprehensive analysis of our DSAQ and uniform feature distillation.

#### 368 4.1 EXPERIMENT SETTINGS

369 **Datasets.** We adopt two LLIE datasets with raw input images to evaluate our low-bit quantization 370 method. The SID Chen et al. (2018) dataset comprises 5094 RAW images captured in extremely 371 low-light conditions, along with their corresponding normal-light reference images. These images 372 were taken using two different cameras: Sony A7S2 (Bayer sensor with a resolution of  $4240 \times$ 373 2832) and Fuji X-T2 (Bayer sensor with a resolution of  $6000 \times 4000$ ). The exposure time for the 374 low-light images in the dataset ranges from 0.1s to 0.033s, which are 100 to 300 times shorter than 375 the corresponding reference images. The MCR Dong et al. (2022) dataset contains 3984 low-light 376 raw images with a resolution of  $1280 \times 1024$  captured from 498 indoor and outdoor scenes. Each scene includes one RGB reference image and 8 low-light raw images with exposure time ranging 377 from 1/4096s to 3/8s.



Figure 4: Visual comparison of different raw-based LLIE methods on SID datasets.

**Training Details.** We use the U-Net model in SID Chen et al. (2018) as the full-precision backbone for low-bit quantization. The weights in low-bit quantized model is initialized with the corresponding parameters in the pretrained full-precision U-Net. During training, the batch size is set to 1 and the size of input raw patch is set to  $1024 \times 1024$ . We train the low-bit quantized model for 300 epochs on these two raw LLIE datasets. We adopt the Adam optimizer Kingma & Ba (2015) with the learning rate set to  $10^{-4}$  and the cosine annealing scheduler for network optimization. All the networks are implemented with PyTorch Paszke et al. (2019) and trained on one NVIDIA RTX 3090 GPU.

**Evaluation Metrics.** We calculate average peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) with enhanced RGB output images and their reference images to evaluate the performance of all the methods. A higher PSNR and SSIM indicate a better restoration quality. We follow previous work Xu et al. (2023) to add  $\{\frac{1}{32}, \frac{1}{16}, \frac{1}{8}\}$  of the number of  $\{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}\}$ -bit operations with respective number of FLOPs to estimate the computational complexity of quantized neural networks.

#### 406 4.2 Compare with State-of-the-arts

Comparison Methods. We first give the performance of the full-precision pretrained SID Chen et al. (2018) U-Net. We then compare our low-bit quantization method with state-of-the-art quantization methods including Dorefa Zhou et al. (2016), PACT Choi et al. (2018), PAMS Li et al. (2020), LSQ Esser et al. (2020), LLT Wang et al. (2022) and QuantSR Qin et al. (2023). In addition, we also compare the low-bit quantization methods with some lightweight full-precision raw-based LLIE methods, including LLPack Lamba et al. (2020) and RRT Lamba & Mitra (2021).

414 Quantitative Results. As shown in Table 1, our low-bit quantized model achieves promising results 415 with low computational cost and memory overhead. Compared with the state-of-the-art quantization methods, our methods yields the best performance in all the 2-bit to 4-bit settings. On the MCR 416 dataset, our method outperforms LSQ Esser et al. (2020) in PSNR/SSIM metrics by 0.39dB/0.0062, 417 0.94dB/0.0110, and 0.98dB/0.0233 for 2-bit, 3-bit, and 4-bit network quantization. Compared 418 with lightweight raw-based LLIE methods, our 4-bit quantized model achieves 1.05dB/0.0148 and 419 4.47dB/0.0549 higher PSNR/SSIM than LLPack Lamba et al. (2020) on the SID and MCR datasets, 420 respectively. Additionally, our 3-bit quantized model surpasses RRT Lamba & Mitra (2021) on 421 both datasets with fewer parameters and computations. Regarding the compression ratio, our 4-422 bit quantized SID U-Net Chen et al. (2018) reduces the model size by 87.5% and the FLOPs by 423 86.6% relative to the full-precision counterpart, while maintaining comparable enhancement results. 424 The compression ratio can achieve 93.7% for parameters and 95.5% for computational costs when 425 quantize the full-precision model to 2-bit.

426

391 392

393

394

395

396

397

398

399

405

413

Visual Comparison. The qualitative results on the SID dataset and MCR dataset are illustrated in
 Figures 4 and 5, respectively. The input is amplifed with the ratio and post-processed for visualization.
 As shown in Figure 4, our 4-bit quantized U-Net yields enhanced images with high visual quality,
 comparable to those produced by the full-precision counterpart. Compared to other methods, our
 approach effectively suppresses severe noise while preserving clear details and textures in the
 enhanced image. Additionally, our method exhibits better color fidelity and consistency in flat areas.



Figure 5: Visual comparison of different raw-based LLIE methods on MCR datasets.

	Method		Bits (w/a)							
DSQ	Asym	UFD	4/4		3	/3	2/2			
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM		
X	X	X	28.62	0.7790	28.33	0.7722	27.79	0.7586		
$\checkmark$	×	×	28.73	0.7813	28.52	0.7741	27.77	0.7599		
×	$\checkmark$	×	28.62	0.7784	28.54	0.7765	27.98	0.7609		
$\checkmark$	$\checkmark$	×	28.72	0.7819	28.61	0.7761	28.14	0.7637		
$\checkmark$	$\checkmark$	$\checkmark$	28.81	0.7823	28.66	0.7772	27.90	0.7603		

Table 2: Ablation study of proposed DSAQ and UFD on SID-Sony dataset.

As shown in Figure 5, our quantization method also demonstrates better perceived quality than state-of-the-art quantization methods in the 2-bit and 3-bit settings.

477 4.3 Ablation Study

We conduct the ablation study to validate the effect of DASQ and uniform feature distillation on the SID dataset. The result is shown in Table 2, where DSQ, Asym and UFD represent whether to use distribution-separative quantization, asymmetric activation quantizer and uniform feature distillation respectively. We can observe from the fourth row that using DSAQ achieves better low-bit quantization performance on U-Net compared to the vanilla symmetric quantizer. It can also be found from the second and third rows that the distribution-separative strategy is more effective with relatively more quantization bins, while the asymmetric quantizer is more useful in lower-bit settings. From the last two rows, we find that the low capacity of 2-bit model limits knowledge transfer even with the feature uniformity module. So we empirically exclude the distillation loss in the 2-bit setting.

36	Distillation Scheme			Bits (w/a)							
37				4/4		3/3		2/2			
38			PS	NR S	SIM	PSNR	SSIM	PSNR	SSIM		
39	Feature Distillatio	Feature Distillation Li et al. (2020)			0.7804 2	28.61 0	0.7761	27.85	0.7600		
0	UFD		28	.81 0.	7823	28.66	0.7772	27.90	0.7603		
11											
2	Table 3	: Ablation stuc	ly of the	distillati	on sche	me on S	ID-Sony	dataset.			
3											
1		Device	CPU	G	PU	N	PU				
5		Bits (w/a)	32/32	32/32	16/16	6 16/16	5 4/8				
6		Time (ms)	190.4	56.3	18.3	3.7	1.7				
,											
3	Table 4: C	Comparison of	inferenc	e time o	n Qualc	comm Sn	apdragon	8 Gen 3.			
)		1					1 0				
	In order to prove the	effectiveness	of the ur	niform fe	eature d	istillatio	n (UFD)	scheme, v	we camp		
	it with the vanilla fea	ture distillation	n in PAN	AS Li et	al. (202	20), whic	h directly	uses the	normaliz		
	feature for distillation	n. The experim	nent resu	lts in Ta	ble 3 pr	oves that	t the prop	osed feat	ure unifo		
	module (FUM) can I	nitigrate the r	epresent	ation ga	p betwo	een the f	features f	rom low-	bit strud		
	models and full-preci-	sion teacher m	odel.								
,	4.4 ON-CHIP LATE	ENCY									
3	<b>W</b> 7				NT	. 1.1	1 1.	1.4			
)	we compare the later	icy of the float	ing-poir	IT SID U	-net m	odel With	n the low-	-bit quant	ized one		
	Qualcomin Snapurage	JII o Gen 3, Wr	nen supp	ons 4W/	oa quan	luzation	on the NP		widely u		

We compare the latency of the floating-point SID U-Net model with the low-bit quantized one on Qualcomm Snapdragon 8 Gen 3, which supports 4w/8a quantization on the NPU and is widely used in smartphones. The resolution of the testing image patch is set to  $256 \times 256$  and the inference time is shown in Table 4. The 4w/8a quantized U-Net model is about  $2.2 \times$  faster than the 16-bit floating-point model on NPU and  $33 \times$  faster than the 32-bit floating-point model running on GPU. Although most devices currently do not support 3-bit or 2-bit, we believe the lower-bit model will be more practical in the future. And our method may be useful for efficiently processing high-resolution images on smartphones or other edge devices.

516 517

518

538

## 5 CONCLUSION

519 In this paper, we propose a low-bit quantization method for raw-based LLIE networks. First, we 520 present a novel low-bit quantizer DSAQ for the U-Net architecture. In order to match the distribution 521 range of the features concatenated via skip connections, DSAQ employs two sets of quantization 522 parameters to separately quantize the two parts of the activations, thereby better fitting these two 523 different distributions. It also exploits the asymmetric activation quantizer for the skewed features activated by LeakyReLU non-linear function. Second, we introduce uniform feature distillation, 524 which employs a feature uniform module to reduce the capability gap between low-bit features 525 and full-precision features, facilitating knowledge transfer from the teacher model. However, it 526 shows limitations in the 2-bit setting, which is worth to explore in the future work. Extensive 527 experiments demonstrate that our low-bit quantized LLIE model can yield satisfactory results with 528 low computational and memory costs. 529

530 531 REFERENCES

Yash Bhalgat, Jinwon Lee, Markus Nagel, Tijmen Blankevoort, and Nojun Kwak. LSQ+: improving low-bit quantization through learnable offsets and better initialization. In *CVPRW*, pp. 2978–2985, 2020.

- Yuanhao Cai, Yuxin Zheng, Jing Lin, Xin Yuan, Yulun Zhang, and Haoqian Wang. Binarized spectral
   compressive imaging. In *NeurIPS*, 2023.
- Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *CVPR*, pp. 3291–3300, 2018.

549

551

555

577

582

- 540 Linwei Chen, Ying Fu, Kaixuan Wei, Dezhi Zheng, and Felix Heide. Instance segmentation in the 541 dark. IJCV, 131(8):2198-2218, 2023. 542
- Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, 543 and Kailash Gopalakrishnan. PACT: parameterized clipping activation for quantized neural 544 networks. arXiv preprint arXiv:1805.06085, 2018.
- 546 Xingbo Dong, Wanyan Xu, Zhihui Miao, Lan Ma, Chao Zhang, Jiewen Yang, Zhe Jin, Andrew 547 Beng Jin Teoh, and Jiajun Shen. Abandoning the bayer-filter to see in the dark. In CVPR, pp. 548 17410-17419, 2022.
- Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmen-550 dra S. Modha. Learned step size quantization. In ICLR, 2020.
- 552 Ruihao Gong, Xianglong Liu, Shenghu Jiang, Tianxiang Li, Peng Hu, Jiazhen Lin, Fengwei Yu, and 553 Junjie Yan. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. 554 In ICCV, pp. 4851–4860, 2019.
- Shuhang Gu, Yawei Li, Luc Van Gool, and Radu Timofte. Self-guided network for fast image 556 denoising. In ICCV, pp. 2511–2520, 2019.
- 558 Cheeun Hong, Heewon Kim, Sungyong Baik, Junghun Oh, and Kyoung Mu Lee. DAQ: channel-wise 559 distribution-aware quantization for deep image super-resolution networks. In WACV, pp. 913–922, 560 2022.
- 561 Yang Hong, Kaixuan Wei, Linwei Chen, and Ying Fu. Crafting object detection in very low light. In 562 BMVC, 2021. 563
- 564 Haofeng Huang, Wenhan Yang, Yueyu Hu, Jiaying Liu, and Ling-Yu Duan. Towards low light 565 enhancement with RAW images. IEEE TIP, 31:1391-1405, 2022.
- 566 Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training 567 quantization with small calibration sets. In ICML, pp. 4466-4475, 2021. 568
- 569 Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew G. Howard, 570 Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for 571 efficient integer-arithmetic-only inference. In CVPR, pp. 2704–2713, 2018.
- 572 Xin Jin, Linghao Han, Zhen Li, Chun-Le Guo, Zhi Chai, and Chongyi Li. DNF: decouple and 573 feedback network for seeing in the dark. In CVPR, pp. 18135–18144, 2023. 574
- 575 Derrick N Joanes and Christine A Gill. Comparing measures of sample skewness and kurtosis. 576 Journal of the Royal Statistical Society, 47(1):183–189, 1998.
- Sangil Jung, Changyong Son, Seohyung Lee, JinWoo Son, Jae-Joon Han, Youngjun Kwak, Sung Ju 578 Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing quantization 579 intervals with task loss. In CVPR, pp. 4350-4359, 2019. 580
- 581 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In ICLR, 2015.
- Mohit Lamba and Kaushik Mitra. Restoring extremely dark images in real time. In CVPR, pp. 583 3487-3497, 2021. 584
- 585 Mohit Lamba, Atul Balaji, and Kaushik Mitra. Towards fast and light-weight restoration of dark 586 images. In BMVC, 2020.
- Huixia Li, Chenqian Yan, Shaohui Lin, Xiawu Zheng, Baochang Zhang, Fan Yang, and Rongrong Ji. 588 PAMS: quantized super-resolution via parameterized max scale. In ECCV, pp. 564–580, 2020. 589
- Yanjing Li, Sheng Xu, Baochang Zhang, Xianbin Cao, Peng Gao, and Guodong Guo. Q-vit: Accurate and fully quantized low-bit vision transformer. In NeurIPS, 2022. 592
- Yanjing Li, Sheng Xu, Xianbin Cao, Xiao Sun, and Baochang Zhang. Q-DM: an efficient low-bit quantized diffusion model. In NeurIPS, 2023.

594 595 596	Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. BRECQ: pushing the limit of post-training quantization by block reconstruction. In <i>ICLR</i> , 2021.
597 598 599 600	Zechun Liu, Baoyuan Wu, Wenhan Luo, Xin Yang, Wei Liu, and Kwang-Ting Cheng. Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced training algorithm. In <i>ECCV</i> , pp. 747–763, 2018.
601 602	Paras Maharjan, Li Li, Zhu Li, Ning Xu, Chongyang Ma, and Yue Li. Improving extreme low-light image denoising via residual learning. In <i>ICME</i> , pp. 916–921, 2019.
603 604 605	Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tij- men Blankevoort. A white paper on neural network quantization. <i>arXiv preprint arXiv:2106.08295</i> , 2021.
607 608 609	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In <i>NeurIPS</i> , pp. 8024–8035, 2019.
610 611	Haotong Qin, Yulun Zhang, Yifu Ding, Yifan liu, Xianglong Liu, Martin Danelljan, and Fisher Yu. Quantsr: Accurate low-bit quantization for efficient image super-resolution. In <i>NeurIPS</i> , 2023.
612 613	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In <i>MICCAI</i> , pp. 234–241, 2015.
615 616	Eli Schwartz, Raja Giryes, and Alexander M. Bronstein. Deepisp: Toward learning an end-to-end image processing pipeline. <i>IEEE TIP</i> , 28:912–923, 2019.
617 618	Hu Wang, Peng Chen, Bohan Zhuang, and Chunhua Shen. Fully quantized image super-resolution networks. In <i>ACM MM</i> , pp. 639–647, 2021.
619 620	Huan Wang and Yun Fu. Trainability preserving neural structured pruning. In ICLR, 2023.
621 622	Longguang Wang, Xiaoyu Dong, Yingqian Wang, Li Liu, Wei An, and Yulan Guo. Learnable lookup table for neural network quantization. In <i>CVPR</i> , pp. 12413–12423, 2022.
623 624	Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. In <i>BMVC</i> , 2018.
625 626 627	Kaixuan Wei, Ying Fu, Yinqiang Zheng, and Jiaolong Yang. Physics-based noise modeling for extreme low-light photography. <i>IEEE TPAMI</i> , 44(11):8520–8537, 2022.
628 629	Ke Xu, Xin Yang, Baocai Yin, and Rynson W. H. Lau. Learning to restore low-light images via decomposition-and-enhancement. In <i>CVPR</i> , pp. 2278–2287, 2020.
630 631 632	Sheng Xu, Yanjing Li, Mingbao Lin, Peng Gao, Guodong Guo, Jinhu Lü, and Baochang Zhang. Q-DETR: an efficient low-bit quantized detection transformer. In <i>CVPR</i> , pp. 3842–3851, 2023.
633 634	Syed Waqas Zamir, Aditya Arora, Salman Khan, Fahad Shahbaz Khan, and Ling Shao. Learning digital camera pipeline for extreme low-light imaging. <i>Neurocomputing</i> , 452:37–47, 2021.
635 636	Yulun Zhang, Huan Wang, Can Qin, and Yun Fu. Learning efficient image super-resolution networks via structure-regularized pruning. In <i>ICLR</i> , 2022.
637 638 639 640	Yunshan Zhong, Mingbao Lin, Xunchao Li, Ke Li, Yunhang Shen, Fei Chao, Yongjian Wu, and Rongrong Ji. Dynamic dual trainable bounds for ultra-low precision super-resolution networks. In <i>ECCV</i> , pp. 1–18, 2022.
641 642 643	Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net: Train- ing low bitwidth convolutional neural networks with low bitwidth gradients. <i>arXiv preprint</i> <i>arXiv:1606.06160</i> , 2016.
644 645 646	Ke Zhu, Yin-Yin He, and Jianxin Wu. Quantized feature distillation for network quantization. In <i>AAAI</i> , pp. 11452–11460, 2023.
	Mintong Thu Dingha Dan Wai Chan and Vi Vang, EEMEEN, low light image anhangement via

<sup>647</sup> Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. EEMEFN: low-light image enhancement via edge-enhanced multi-exposure fusion network. In *AAAI*, pp. 13106–13113, 2020.