# LL-ICM: Image Compression for Low-level Machine Vision via Large Vision-Language Model

Yuan Xue[†], Qi Zhang[‡], Chuanmin Jia[♯], and Shiqi Wang[†]

[†] School of Computer Science, City University of Hong Kong, Hong Kong SAR, China
[‡] National Engineering Research Center of Visual Technology, Peking University, Beijing, China
[♯] Wangxuan Institute of Computer Technology, Peking University, Beijing, China

## Abstract

Image Compression for Machines (ICM) aims to compress images for machine vision tasks rather than human viewing. Current works predominantly concentrate on high-level tasks like object detection and semantic segmentation. However, the quality of original images is usually not guaranteed in the real world, leading to even worse perceptual quality or downstream task performance after compression. Low-level (LL) machine vision models, like image restoration models, can help improve such quality, and thereby their compression requirements should also be considered. In this paper, we propose a pioneered ICM framework for LL machine vision tasks, namely LL-ICM. By jointly optimizing compression and LL tasks, the proposed LL-ICM not only enriches its encoding ability in generalizing to versatile LL tasks but also optimizes the processing ability of down-stream LL task models, achieving mutual adaptation for image codecs and LL task models. Furthermore, we integrate large-scale vision-language models into the LL-ICM framework to generate more universal and distortion-robust feature embeddings for LL vision tasks. Therefore, one LL-ICM codec can generalize to multiple tasks. We establish a solid benchmark to evaluate LL-ICM, which includes extensive objective experiments by using both full and no-reference image quality assessments. Experimental results show that LL-ICM can achieve 22.65% BD-rate reductions over the state-of-the-art methods.

## Introduction

Image Coding for Machines (ICM) has become an emerging research topic that combines visual signal compression and understanding. Existing methods can be categorized according to the number of generated bitstreams: a single versatile bitstream for multiple tasks jointly, two bitstreams for human perception and machine task respectively and multiple bitstreams for diverse tasks separately. These ICM bitstreams can be compact representations of both original image [1] and extracted deep features [2], which usually have different levels of capability on texture reconstruction and semantic preservation. However, all current methods are limited by only optimizing for high-level (HL) vision tasks, while neglecting practical compression requirements for low-level (LL) ones.

Low-level vision tasks, such as denoising, debluring, inpainting, *etc.*, have been studied for decades. The main purpose of these tasks is to recover the details introduced by image capturing and delivering, and enhance the visual quality of the image content. Early works tend to propose one neural network for one specific LL vision task. Recently, considering the similarity of these tasks on pixel-level content analysis and processing, more and more works are trying to create a single model to solve diverse LL vision tasks simultaneously. By introducing large vision-language modeling, state-of-the-art all-in-one LL vision
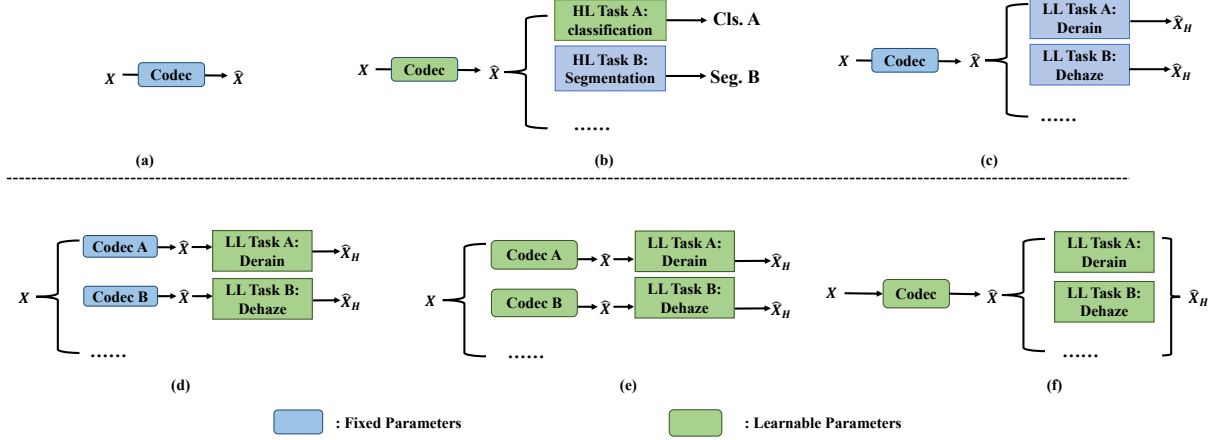
**(a)** ... **(b)** ... **(c)**

**(d)** ... **(e)** ... **(f)**

□ : Fixed Parameters    □ : Learnable Parameters

Figure 1: ICM frameworks from the simplest to the proposed. $X$ is the original image, $\hat{X}$ is the compressed image, and $\hat{X}_H$ is the enhanced image by downstream task model. $Cls.A$ and $Seg.B$ are the results of HL tasks A and B.

models can achieve impressive performance on important LL vision tasks. These methods are also more suitable for real-world applications ***where distortion types unknown or combined.***

Image compression is usually the subsequent processing step after image capturing, which is essential for content delivery with diverse network conditions. However, compression also leads to information lost and quality degradation. Since LL visual enhancement cannot always be performed before compression due to diverse restrictions, there are two approaches to improve its performance. First, we can design LL vision models that consider compression distortions or even have the capability of removing them. But this approach has two notable drawbacks: 1) There are many types of compression distortions from different codecs like JPEG, VVC, AVS, neural codecs, *etc.*, making such model hard to architect and train. 2) Without leveraging the power of LL vision models, the compression efficiency is suboptimal. For example, some pixels could have been compressed more severely and restored or enhanced appropriately in the downstream. Similarly, some signals are visually unpleasant or broken, which should not be preserved after compression. Second, and more ideally, we can create an end-to-end compression and LL vision model, which jointly optimizes the performance of enhancing the visual perceptual quality and saving the bit-rate cost at the same time.

In this work, we propose the first ICM framework for LL vision tasks, namely LL-ICM. Our goal is to achieve the rate-perception optimum by ***breaking the long lasting isolation between compression and LL visual enhancement via a unified model.*** Considering the diversity and complexity of distortion types in the real world, it is expensive or even impractical to customize LL-ICM model for each individually. Therefore, we further extend our framework so that it can adapt to versatile LL vision tasks using a single all-in-one model. Specifically, we introduce large-scale vision-language model for generalized feature embedding generation and exploit diffusion procedure to remove distortions efficiently. Our primary contributions are delineated as follows.

- We propose LL-ICM, the first compression framework for LL vision. By jointly optimizing compression and LL vision processing, LL-ICM can improve the performance of both compared to existing frameworks that only consider compression or LL task alone.

- We integrate large vision-language models (VLM) into LL-ICM. Based on the generalized feature representations learned by these models, we can use a diffusion model to reduce arbitrary type of distortions and increase the perceptual quality of compressed images. Therefore, one LL-ICM model can handle multiple LL vision tasks.

- We create a large-scale benchmark to evaluate the performance of LL-ICM through objective evaluations. Experimental results indicate that LL-ICM significantly outperforms existing compression methods on both full and no-reference quality measurements. Specifically, LL-ICM achieves 22.65% BD-rate reductions compared to state-of-the-art neural codecs.

## Method

*Exsiting ICM frameworks*

We first give a brief introduction to the existing ICM frameworks and analyze their drawbacks. Fig. 1(a) is a conventional image codec optimized for signal fidelity without considering the performance of downstream tasks on compression outputs. As machine intelligence evolves, the compression needs for high-level visual understanding are getting more and more attention, leading to the establishment of HL-ICM framework shown in Fig. 1(b). In this framework, the codec is optimized for HL task models, whether their weights are fixed or not.

Most ICM methods focus on HL vision tasks. However, LL tasks are also important for applications in the real world. The reason is that the quality of original images is usually not guaranteed. After compressed and transmitted, we need to enhance their quality by LL vision processing. The simplest ICM framework for LL-vision tasks is illustrated in Fig. 1(c), where the compressed image is processed by LL task models, but LL vision models do not consider the compression artifacts. In recent years, this issue has been addressed by the framework shown in Fig. 1(d), where LL vision models are trained with adaption to the compression distortion. However, the compression efficiency is neglected. Therefore, the LL-ICM framework in Fig. 1(e) offers a better solution that jointly improves the performance of compression and LL vision processing. Nevertheless, there are many LL vision tasks and models, and it is impractical to train a codec for each of them. In this paper, we take a step further to propose a unified one-for-all LL-ICM framework to optimize the codec for diverse LL vision tasks, which is shown in Fig. 1(f).

*Problem Definition*

LL-ICM aims to jointly optimize the performance of multiple LL vision tasks while minimizing compression costs. Inherently, LL vision tasks target at removing the artifacts from images causing quality degradations and improving the perceptual quality. Therefore, the objective of LL-ICM can be defined as rate-perception optimization (RPO), which is
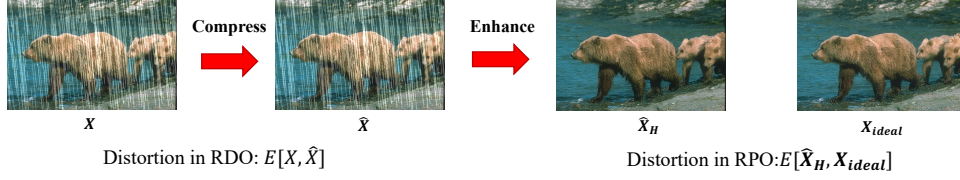
Figure 2: Difference between RPO and RDO Definition: RDO distortion measures the error of compressed image $\hat{\mathbf{X}}$ and the original image $\mathbf{X}$, while RPO measures the error of the generated enhanced image $\hat{\mathbf{X}}_{\mathbf{H}}$ and the high-quality ideal image $\mathbf{X}_{\mathbf{ideal}}$.

similar to the classical rate-distortion optimization (RDO) problem in image compression. There are several differences between them. Most importantly, RDO progress assumes that the original image has an ideal level of quality, so the compression is optimized to preserve signal fidelity. However, RPO progress admits that the original image can have several flaws, so the compression should be optimized to increase the perceptual quality with the assistance of LL vision models.

As illustrated in Fig. 2, the distortion in RDO, depends on the degradation of the original image $\mathbf{X}$, can be formulated as $\mathbb{E}(\mathbf{X}, \hat{\mathbf{X}})$, where $\hat{\mathbf{X}}$ is the compressed image, $\mathbb{E}$ is the error assessment function of compressed and original image. However, the perceptual distortion in RPO, depends on the enhancement of $\hat{\mathbf{X}}$, can be formulated as $\mathbb{E}(\tau(\hat{\mathbf{X}}), \mathbf{X}_{\mathbf{ideal}})$. We define $\mathbf{X}_{\mathbf{ideal}}$ as the ideal high-quality image, and $\tau(\cdot)$ is used to enhance $\hat{\mathbf{X}}$ so that its quality can be close to $\mathbf{X}_{\mathbf{ideal}}$, and $\hat{\mathbf{X}}_{\mathbf{H}}$ is the enhanced version of compressed image $\hat{\mathbf{X}}$: Generally, the RPO function $\mathbf{p}_{\text{opt}}$ can be described as follows:

$$\mathbf{p}_{\text{opt}} = \arg\min_{\mathbf{p}}\{\mathbb{E}(\hat{\mathbf{X}}_{\mathbf{H}}, \mathbf{X}_{\mathbf{ideal}}) + \lambda\mathcal{R}(\mathbf{X})\}, \tag{1}$$

where $\mathcal{R}$ represents the bit-rate of compressed image.

*LL-ICM framework*

We present a unified LL-ICM framework in Fig. 3, which integrates a neural image codec and a unified LL vision processing model. Notably, we incorporate a pre-trained VLM model into this framework to extract generalized features for handling different LL vision tasks.

In this study, we employ MLIC [3] as backbone image codec $\mathcal{C}$ due to its superior performance. The frozen VLM model $\mathcal{V}$ [4] is used to extract generalized features $\mathbf{F}$ as:

$$\mathbf{F} = \mathcal{V}(\hat{\mathbf{X}}) = \mathcal{V}(\mathcal{C}(\mathbf{X})). \tag{2}$$

After that, the LL-vision encoder recieve $\hat{\mathbf{X}}$ and encoded $\mathbf{F}$ by encoder $\mathcal{E}$ as input reference to generate $\hat{\mathbf{X}}_{\mathbf{H}}$:

$$\hat{\mathbf{X}}_{\mathbf{H}} = \tau(\hat{\mathbf{X}}, (\mathcal{E}(\mathbf{F})). \tag{3}$$

We adopt the image controller in DA-CLIP [5] as the encoder in the feature controller of our framework, which encodes the feature $\mathbf{F}$ to LL task type $\varphi$ and caption $\sigma$. For

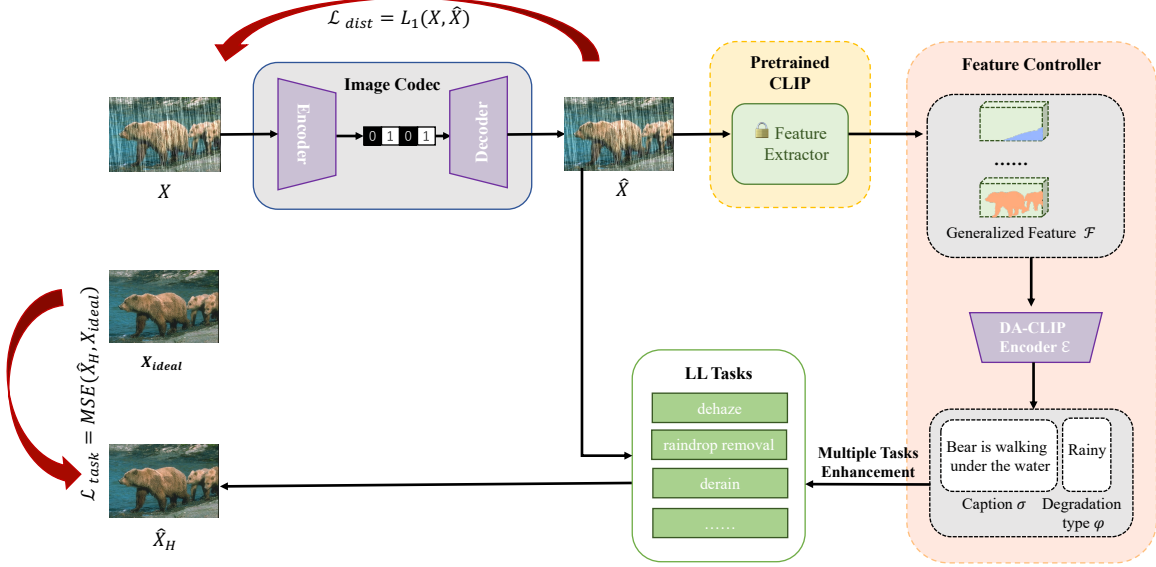Figure 3: Overview of the LL-ICM framework: $\hat{\mathbf{X}}$ is the compressed image. A pre-trained CLIP model extracts the generalized feature $\mathbf{F}$ from $\hat{\mathbf{X}}$. After the DA-CLIP encoder $\mathbf{E}$, the encoded $\mathbf{F}$ can be used for multiple LL tasks and reconstruct the enhanced image $\hat{\mathbf{X}}_{\mathbf{H}}$.

instance, a rainy image in Fig. 3 as input of DA-CLIP can be used to extract distortion type of "rainy" and a caption of the image as "bear is walking under the water". Thus, we have:

$$(\varphi, \sigma) = \mathcal{E}(\mathbf{F}). \tag{4}$$

This textual information, $\varphi$ and $\sigma$, guides the LL-vision enhancement network to generate $\hat{\mathbf{X}}_{\mathbf{H}}$. We employ the IR-SDE [6] diffusion network to handle multiple LL tasks and generate the $\hat{\mathbf{X}}_{\mathbf{H}}$.

Regarding learning objective of the proposed LL-ICM $\mathcal{L}$, we aim to optimize the distortion loss $\mathcal{L}_{dist}$ between $\mathbf{X}$ and $\hat{\mathbf{X}}$, the rate loss $\mathcal{L}_r$, and the loss of the downstream LL task $\mathcal{L}_{task}$ between $\hat{\mathbf{X}}_{\mathbf{H}}$ and $\mathbf{X}_{\mathbf{ideal}}$. The total loss $\mathcal{L}$ of the LL-ICM can be formulated as:

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{dist} + \beta \cdot \mathcal{L}_r + \gamma \cdot \mathcal{L}_{task}, \tag{5}$$

where the MSE loss are calculated for $\mathcal{L}{dist}$, and $\mathcal{L}r$ represents the bitrate of $\hat{\mathbf{X}}$. The $\mathcal{L}{task}$ is the diffusion loss, as described in IR-SDE [6], which ensures that $\hat{\mathbf{X}}_{\mathbf{H}}$ is as close as possible to $\mathbf{X}_{\mathbf{ideal}}$.

## Experiment

*Training Setting*

Our training procedure is divided into two stages. In the first stage, we only train the image codec to achieve high compression efficiency. We train the codec on two NVIDIA 4090 GPUs for 30 epochs, using the Adam optimizer with a batch size of 16. We randomly select $2 \times 10^5$ images from the COCO2017 and ImageNet datasets as the initial training set, and the images are randomly cropped to a patch size of 448x448. Following the settings of

CompressAI library, we set $\beta \in \{0.0002, 0.0008, 0.0018, 0.0035, 0.0130, 0.0350\}$ for $\mathcal{L}_{dist}$. The parameters $\alpha$ and $\gamma$ are set to be 1 and 0, respectively.

In the second stage, we load our pre-trained codec from the first stage. The LL task is then trained jointly with the codec and optimize it jointly for compression and LL vision tasks. In this work, we choose ***dehazing, raindrop removal, deraining, deshadowing, denoising, and inpainting*** as exemplar LL vision tasks. It should be noticed that the proposed LL-ICM framework can be extended to even more tasks because of the integrated large vision-language model. To further increase the generalizability of our models, we incorporate several in-the-wild datasets into our training, such as LOL [7] and RESIDE-6k [8]. The training datasets in this stage are listed in TABLE 1. Noticeably, we manually generate several distortions on datasets of the denoising task by adding gaussian noise to the original images with a level of 50.

Table 1: TRAINING AND TESTING DATASETS USED FOR EVALUATING LL-ICM

| LL Task Type | Dataset | Size(Train+Test) | LL Task Type | Dataset | Size(Train+Test) |
|---|---|---|---|---|---|
| deraining | Rain100H [9] | 1899+100 | dehazing | RESIDE-6k [8] | 6000+1000 |
| inpainting | CelebaHQ-256 [10] RePaint [10] | 29901+100 | denoising | DIV2K denoising Flick2k CBSD68 [11] | 3550+68 |
| deshadowing | SRD [12] | 2680+408 | raindrop removal | Raindrop [13] | 861+58 |

*Testing Setting*

To evaluate the performance of our proposed LL-ICM framework, we test the compression and LL vision task performance and compare our method with several state-of-the-art image codecs. The selected codecs include both traditional block-based codec like Enhance Enhanced Compression Model (ECM) [14], and deep learning-based ones like Balle2018 [15], Cheng2020 [16], and MLIC [3]. To ensure the fairness of the comparison, we also use IR-SDE [6] as the LL vision model to generate enhance outputs of the compressed images from these codecs. We conduct objective experiments by calculating both full-reference and no-reference image quality metrics on the generated $\hat{X}_H$.

Specifically, we use the popular full-reference image quality assessment (FR-IQA) metric LPIPS [17] and the state-of-the-art no-reference image quality assessment (NR-IQA), Q-Align [18] and LIQE [19], in objective test. As FR-IQA, we use the image pairs in the datasets as $X_{ideal}$ and $\hat{X}_H$, respectively. In contrast, we evaluate only the generated $\hat{X}_H$ in NR-IQA without using $X_{ideal}$ as the reference. Here, we introduce NR-IQA metrics because $X_{ideal}$ is often inaccessible or even does not exist in the real world. We also calculate the average coding gains and perceptual quality enhancements using the Bjontegaard Delta-rate (BD-rate) metrics. The BD-LPIPS and BD-Q-Align metrics are computed similarly to BD-PSNR in [20] by substituting the distortion metric with LPIPS or Q-Align. A negative BD-rate value indicates rate savings at equivalent quality levels. Negative BD-LPIPS and positive BD-Q-Align values signify quality improvements at the same bit rate.

*Rate-percetion Performance Result*

Fig. 4 illustrates the rate-perception (RP) performance results on different tasks. Fig. 4 (a)-(d) represents the results on dehazing, Raindrop removal, deraining, and deshadowing,
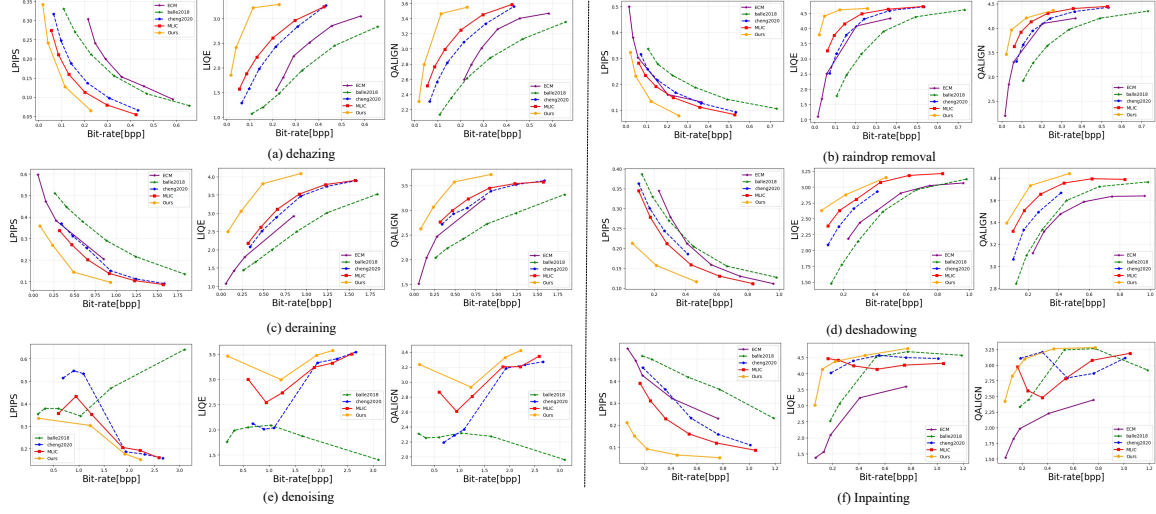
Figure 4: RP performance on 6 LL tasks: Our LL-ICM framework significantly outperforms all existing codecs. RP performance on (a)-(d) are the monotonic cases in which the perception quality increases as the bpp increases. The (e) and (f) are abnormal cases as the perception quality does not follow the regular law as (a)-(d).

respectively. The results show the average quality and bitrate in different datasets. It is easy to observe that the perceptual quality of all anchors in Fig. 4 (a)-(d) increases as the bpp increases. Our LL-ICM model achieves the best performance among the four tasks. To be specific, the LPIPS keeps the lowest and Q-Align situated at the top of the diagram.

However, not all LL tasks keep such regular results. Fig. 4 (g) shows the results of de-noise. As shown in the anchor "Balle2018", the perception performance abnormally decreases as the bpp increases. We visualize an example in Fig. 5. The noisy information is even more in $\hat{X}_H$ with higher bpp. Similar results are also shown in our LL-LCM framework. A turning point is shown in the Q-Align/LIQE diagram, representing the enhancement quality increases when bpp increases after the turning point. Thus, we can conclude that as the bpp increases, the codec keeps more detail and chaotic signal in CLQ, which will interfere with the de-noise downstream. Fortunately, our LL-ICM framework still performs best among the anchors. Similarly, the trend is shown in Fig. 4 (h).

Interestingly, not all LL tasks yield consistent results. Fig. 4 (e) shows the results for the denoising task. As shown in the result from method "Balle2018", the performance abnormally decreases as the bpp increases. We visualize a representative example in Fig. 5, where the noisy signal in $\hat{X}_H$ actually increases with higher bpp. Similar results are observed in our LL-ICM framework. A turning point is evident in the Q-Align diagram, indicating that the quality improves when the bpp increases beyond this point. Therefore, we can conclude that as bpp increases, the codec retains more detail and chaotic signals in $\hat{X}$, which can interfere with the denoising process. Nevertheless, our LL-ICM framework still outperforms other methods in this extreme condition. Besides, the turning point in MLIC is also obvious as shown in Fig. 4 (f). The quality decreases, and the inpainting result becomes worse when the bpp increases before this point.

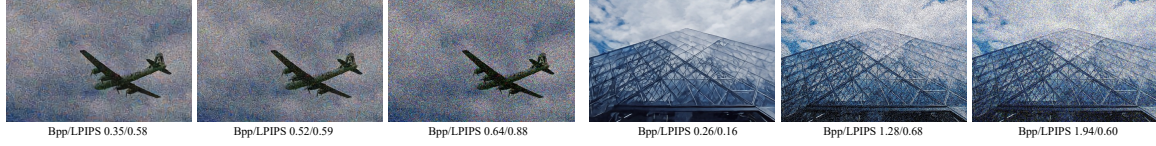| Bpp/LPIPS 0.35/0.58 | Bpp/LPIPS 0.52/0.59 | Bpp/LPIPS 0.64/0.88 | Bpp/LPIPS 0.26/0.16 | Bpp/LPIPS 1.28/0.68 | Bpp/LPIPS 1.94/0.60 |

Figure 5: An abnormal result on denoising: the noise information is more in higher bpp image compared to lower bpp image.

Table 2: BD-RATE, BD-LPIPS, AND BD-QALIGN RESULTS FOR DIFFERENT LL TASKS WITH BALLE2018 [15] BEING THE ANCHOR

| Tasks | dehazing | | | raindrop removal | | | denoising | | |
|---|---|---|---|---|---|---|---|---|---|
| | BD-Rate (%) ↓ | BD-LPIPS ↓ | BD-QAlign ↑ | BD-Rate (%) ↓ | BD-LPIPS ↓ | BD-QAlign ↑ | BD-Rate (%) ↓ | BD-LPIPS ↓ | BD-QAlign ↑ |
| Cheng2020 | -45.22 | -0.08 | 0.51 | -40.07 | -0.06 | 0.49 | -53.36 | -0.05 | 0.45 |
| ECM | -52.07 | -0.10 | 0.35 | -55.02 | -0.08 | 0.55 | -49.13 | -0.02 | 0.43 |
| MLIC | -59.72 | -0.11 | 0.69 | -54.37 | -0.08 | 0.63 | -91.89 | -0.10 | 0.63 |
| LL-ICM (Ours) | **-76.38** | **-0.17** | **1.18** | **-78.72** | **-0.16** | **0.93** | **-79.79** | **-0.03** | **0.43** |
| Tasks | deraining | | | deshadowing | | | inpainting | | |
| Cheng2020 | -41.53 | -0.11 | 0.58 | -25.32 | -0.04 | 0.23 | -52.65 | -0.13 | 0.14 |
| ECM | -48.96 | -0.12 | 0.52 | -15.07 | -0.01 | 0.13 | -55.44 | -0.12 | -0.63 |
| MLIC | -51.46 | -0.14 | 0.66 | -35.53 | -0.05 | 0.28 | -73.17 | 0.22 | -0.15 |
| LL-ICM (Ours) | **-77.17** | **-0.25** | **1.18** | **-74.08** | **-0.14** | **0.51** | **-96.07** | **-0.38** | **0.36** |

TABLE 2 summarizes the average coding gains and perceptual quality score enhancements on all evaluated LL tasks. LL-ICM outperforms all compared methods. The results demonstrate a substantial increase in perceptual quality score, with LL-ICM ***achieving 7.70%-49.23% coding gains*** across various LL tasks compared to state-of-the-art methods. At the same time, the LPIPS score is significantly improved by 0.06-0.16, and the Q-Align score is increased by 0.20-0.52.

*Qualitative Comparison*

Fig. 6 demonstrates the qualitative comparison of generated $\hat{X}_H$ images by LL-ICM against the other anchors. The compressed images are in similar bit rates. Two exemplary LL tasks are showcased to demonstrate the superior performance of our LL-ICM framework, which are image deraining and raindrop removal. Obviously, our method preserves more intricate details and texture in the reconstructed and enhanced images. For example, the result for the deraining task displayed in Fig. 6 reveals that MLIC introduces significant rippling artifacts, while our method keeps the image content clean. Similarly, for the raindrop removal task, the results of other codecs lack details and are not clear, while our method provides the most satisfying image quality.

## Conclusion

In this paper, we propose the first image compression framework for low level machine vision tasks, LL-ICM. LL-ICM aims to increase the performance of both compression and LL vision processing, addressing that many original images to be compressed may have quality flaws and need to be enhanced by LL vision models after compression. Built upon a capable neural image codec MLIC, we incorporate a large vision-language model into the proposed LL-ICM framework, extracting generalized features to support multiple downstream LL vision tasks simultaneously. The compression and LL vision processing

Figure 6: Qualitative comparisons between our method and other codecs on image deraining, inpainting and raindrop removal tasks. The first and the last lines show the original image $\mathbf{X}$ and the ideal high-quality image $\mathbf{X}_{\mathbf{ideal}}$. The other lines show the enhanced image $\hat{\mathbf{X}}_{\mathbf{H}}$ generated by different codecs.

are then jointly optimized during training. Extensive experiment results demonstrate that our LL-ICM framework significantly outperforms existing image codecs on different LL vision tasks.

## References

[1] Yi-Hsin Chen, Ying-Chieh Weng, Chia-Hao Kao, Cheng Chien, Wei-Chen Chiu, and Wen-Hsiao Peng, "Transtic: Transferring transformer-based image compression from human perception to machine perception," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 23297–23307.

[2] Chaoran Chen, Mai Xu, Shengxi Li, Tie Liu, Minglang Qiao, and Zhuoyi Lv, "Residual based hierarchical feature compression for multi-task machine vision," in *2023 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2023, pp. 1463–1468.

[3] Wei Jiang, Jiayu Yang, Yongqi Zhai, Peirong Ning, Feng Gao, and Ronggang Wang, "Mlic: Multi-reference entropy model for learned image compression," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 7618–7627.

[4] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al., "Learning transferable

visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.

[5] Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön, "Controlling vision-language models for universal image restoration," *arXiv preprint arXiv:2310.01018*, 2023.

[6] Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön, "Image restoration with mean-reverting stochastic differential equations," *arXiv preprint arXiv:2301.11699*, 2023.

[7] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu, "Deep retinex decomposition for low-light enhancement," *arXiv preprint arXiv:1808.04560*, 2018.

[8] Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia, "Ffa-net: Feature fusion attention network for single image dehazing," in *Proceedings of the AAAI conference on artificial intelligence*, 2020, vol. 34, pp. 11908–11915.

[9] Wenhan Yang, Robby T Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan, "Deep joint rain detection and removal from a single image," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 1357–1366.

[10] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool, "Repaint: Inpainting using denoising diffusion probabilistic models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11461–11471.

[11] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proceedings eighth IEEE International Conference on Computer Vision*. IEEE, 2001, vol. 2, pp. 416–423.

[12] Liangqiong Qu, Jiandong Tian, Shengfeng He, Yandong Tang, and Rynson WH Lau, "Deshadownet: A multi-context embedding deep network for shadow removal," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4067–4075.

[13] Rui Qian, Robby T Tan, Wenhan Yang, Jiajun Su, and Jiaying Liu, "Attentive generative adversarial network for raindrop removal from a single image," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 2482–2491.

[14] Mohsen Abdoli, Ramin G Youvalari, Karam Naser, Kevin Reuzé, and Fabrice Le Léannec, "Video compression beyond vvc: Quantitative analysis of intra coding tools in enhanced compression model (ecm)," *arXiv preprint arXiv:2404.07872*, 2024.

[15] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston, "Variational image compression with a scale hyperprior," *arXiv preprint arXiv:1802.01436*, 2018.

[16] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto, "Learned image compression with discretized gaussian mixture likelihoods and attention modules," in *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2020, pp. 7939–7948.

[17] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on Computer Vision and Oattern Recognition*, 2018, pp. 586–595.

[18] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al., "Q-align: Teaching lmms for visual scoring via discrete text-defined levels," *arXiv preprint arXiv:2312.17090*, 2023.

[19] Weixia Zhang, Guangtao Zhai, Ying Wei, Xiaokang Yang, and Kede Ma, "Blind image quality assessment via vision-language correspondence: A multitask learning perspective," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14071–14081.

[20] Gisle Bjontegaard, "Calculation of average psnr differences between rd-curves," *ITU SG16 Doc. VCEG-M33*, 2001.