

Large Language Model for Lossless Image Compression with Visual Prompts

Junhao Du¹ Chuqin Zhou¹ Ning Cao² Gang Chen² Yunuo Chen¹
 Zhengxue Cheng¹ Li Song¹ Guo Lu^{1*} Wenjun Zhang¹

¹Shanghai Jiao Tong University

²China Telecommunications Corporation

Abstract

Recent advancements in deep learning have driven significant progress in lossless image compression. With the emergence of Large Language Models (LLMs), preliminary attempts have been made to leverage the extensive prior knowledge embedded in these pretrained models to enhance lossless image compression, particularly by improving the entropy model. However, a significant challenge remains in bridging the gap between the textual prior knowledge within LLMs and lossless image compression. To tackle this challenge and unlock the potential of LLMs, this paper introduces a novel paradigm for lossless image compression that incorporates LLMs with visual prompts. Specifically, we first generate a lossy reconstruction of the input image as visual prompts, from which we extract features to serve as visual embeddings for the LLM. The residual between the original image and the lossy reconstruction is then fed into the LLM along with these visual embeddings, enabling the LLM to function as an entropy model to predict the probability distribution of the residual. Extensive experiments on multiple benchmark datasets demonstrate our method achieves state-of-the-art compression performance, surpassing both traditional and learning-based lossless image codecs. Furthermore, our approach can be easily extended to images from other domains, such as medical and screen content images, achieving impressive performance. These results highlight the potential of LLMs for lossless image compression and may inspire further research in related directions.

1. Introduction

Lossless image compression aims to reduce image size as much as possible without introducing any distortion, making it essential for high-quality data storage and transmission. Furthermore, the techniques used in lossless compression often play a key role in lossy compression meth-

ods. Over the past few decades, numerous effective lossless image codecs have been developed. Among these, traditional codecs such as PNG [6], WebP [15], FLIF [39], and JPEG-XL [1] have achieved strong compression performance through hand-crafted coding algorithms. For example, JPEG-XL employs invertible transforms and a sophisticated context model, including tree structure and pre-context predictor selection, to compress images effectively. In recent years, learning-based lossless image codecs [2, 28, 29, 54, 55] become increasingly popular. L3C [28], for instance, utilizes a hierarchical probability prediction framework and introduces auxiliary latent representations to model the probability distribution of image data. These state-of-the-art (SOTA) methods typically rely on empirical knowledge in image compression and employ meticulously designed models to achieve better compression performance.

Recently, Large Language Models (LLMs) have achieved significant breakthroughs in Natural Language Processing tasks, and their applications have extended to vision tasks, driving substantial progress in areas such as image generation [14, 34] and image restoration [56]. The primary objective of LLMs is to predict the probability distribution of the next token in a sequence. Consequently, more advanced LLM results in more precise modeling of data distribution. Similarly, entropy coding in lossless compression seeks to accurately model data distribution to minimize the coding bitrate. This parallel suggests that LLMs could potentially serve as powerful tools for entropy coding.

Recent work by Delétang et al. [11] supports this perspective, demonstrating that LLMs not only achieve impressive results in text compression but also demonstrate strong potential for lossless image compression. This highlights the advantages of leveraging LLMs in the compression domain. However, pretrained LLMs primarily encapsulate textual prior knowledge, whereas image compression relies more on visual information for optimal performance. Therefore, it is crucial to bridge the gap between the textual nature of LLMs and visual data compression tasks. Unfortunately, the existing approach [11] directly treats the pixel

*Corresponding author.

values of input images as indexes for LLMs, overlooking the inherent spatial relationships within the images. Consequently, the compression efficiency of this method is sub-optimal. For instance, the model proposed by Delétang et al. [11] with 7B parameters performs only slightly better than PNG [6]. Thus, how to effectively unlock the prior knowledge of LLMs and activate their potential for lossless image compression remains a critical issue that deserves in-depth exploration.

In this work, we propose a novel framework for lossless image compression that leverages the LLM with visual prompts. Specifically, the image is initially compressed using a lossy codec, and this lossy reconstruction is then employed as visual prompts for the LLM. Subsequently, the LLM is used to predict the probability distribution of the residual between the lossy reconstruction and the original image. Finally, the probability distributions of the residual pixels are modeled using the Gaussian Mixture Model (GMM), where the parameters are predicted from the output features generated by the LLM. Furthermore, by finetuning the pretrained LLM with Low-Rank Adaption (LoRA) [20], we further enhance our compression performance. Our approach has been evaluated on several benchmark datasets, including Kodak, CLIC, and DIV2K. The results demonstrate that our method achieves SOTA performance, comparable to other well-designed codecs. Our research provides novel insights into lossless image compression and highlights the potential of LLMs for this task.

Our main contributions can be summarized as follows:

- By employing the lossy reconstruction as visual prompts for the LLM, we guide the LLM for more efficient lossless data compression.
- The extensive experimental results demonstrate the SOTA performance of our approach on benchmark datasets. Moreover, our approach can be readily applied to images from other domains, such as screen content images and medical images.

2. Related Work

2.1. Lossy Image Compression

Lossy image compression methods aim to minimize coding distortion at a given bitrate. Traditional lossy image coding standards, such as JPEG [45] and BPG [4], employ manually designed modules to improve the compression performance. For instance, the widely-used JPEG codec leverages the discrete cosine transform (DCT) to reduce spatial redundancy and employs Huffman coding to further reduce bitrates losslessly. Most lossy codecs adhere to the rate-distortion principle, selecting optimal coding modes to achieve better compression performance.

Recent advancements in learning-based lossy image compression [22, 24, 25] have surpassed the SOTA tradi-

tional codecs like VVC [7]. The hyperprior model by Ballé et al. [3] has been studied as a powerful paradigm, applying lossy transforms, quantization, and efficient lossless encoding of latent representations. Some works [10, 57, 58] employ advanced architectures, such as attention mechanism [44] and Swin-Transformer [26], to improve information retention during lossy transforms. Additionally, studies like [31] have optimized the lossless latent coding, incorporating autoregressive components with the hyperprior to capture causal context. Refinements of the context model have led to further improvements in compression [17, 18, 30].

Many advancements in hyperprior-based methods focus on enhancing the lossless compression of latent representations by achieving more accurate distribution estimation. Consequently, lossy and lossless image compression are closely related, with lossless compression techniques often contributing to greater efficiency in lossy compression.

2.2. Lossless Image Compression

Traditional lossless image codecs, such as PNG [6], WebP [15], FLIF [39], and JPEG-XL [1], typically utilize hand-crafted techniques to reduce intra-image redundancy. These methods typically follow a process of filtering, transforming, quantizing, and applying entropy coding to generate the final bitstream. Recently, learning-based lossless image compression has gained significant attention, typically consisting of two stages: 1) constructing a statistical model to capture the probability distribution of image data. 2) utilizing this statistical model to encode the image into a bitstream using entropy tools such as arithmetic coding (AC) or asymmetric numerical systems (ANS) [13]. We employ AC as the lossless data compression technique, due to its widespread use in coding systems and its ability to generate nearly optimal-length codes based on a given probability distribution and input sequence. It encodes an entire message as a single number within the interval $[0, 1)$ (represented in binary), using a probabilistic model to subdivide the interval into subintervals proportional according to each symbol's probability.

To enhance statistical models for lossless image compression, deep generative models have been introduced and can be broadly categorized into three types: 1) *Autoregressive models*, such as PixelRNN [43], PixelCNN [42] and PixelCNN++ [37], which predict pixel distributions based on conditional dependencies with previously obtained pixels via masked convolutions. 2) *Flow models*, such as iVPF [55] and iFlow [54], which leverage invertible transforms to simplify latent distributions for efficient entropy coding. 3) *Variational Auto-Encoder (VAE) models*, like L3C [28], which employ VAE architectures to model image distributions. It is noteworthy that some studies have managed to achieve lossless compression by first compressing

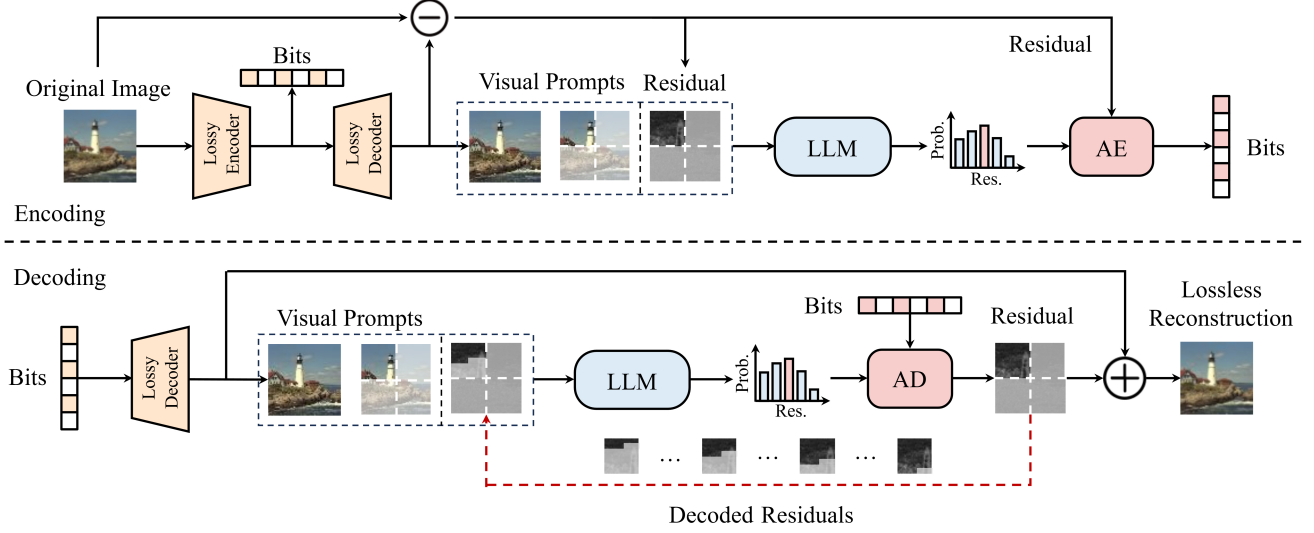


Figure 1. Overview of the encoding and decoding process. A lossy reconstruction \mathbf{x}_l and its patch \mathbf{x}_l^n serve as visual prompts for the LLM to predict the residual’s probability distribution, with the decoding process mirroring encoding by generating residual tokens autoregressively. The red dashed line represents the autoregressive process, where the decoded residuals serve as input to the LLM to predict the probability distribution of the next residual. This process continues until all residuals are decoded. (AE: Arithmetic Encoder. AD: Arithmetic Decoder. LLM: Large Language Model.)

the image using a lossy encoder, and then compressing the residuals. For example, RC [29] integrates BPG for image compression and a CNN for residual compression, whereas DLPR [2] combines VAE with autoregressive models to enhance performance.

However, these methods typically rely on complex network designs and are constrained by limited training datasets, especially in the fields like medical images where data is scarce. This highlights the need for a simple pipeline that leverages the extensive prior knowledge embedded in pretrained models from other datasets to enhance compression efficiency.

2.3. Large Language Models

Large language models (LLMs) have gained significant attention in natural language processing (NLP) and artificial general intelligence (AGI) for their impressive abilities in language generation, in-context learning, world knowledge, and reasoning [47]. LLMs can quickly adapt to specific tasks using techniques like Adapters [19] and Low-Rank Adaptation (LoRA) [20]. Recent research has extended the potential of LLMs to computer vision tasks, such as image classification and segmentation [16, 53]. However, these studies primarily focus on aligning textual and visual semantics while overlooking low-level visual features. Addressing this gap, LM4LV [56] employs LLMs for image restoration, emphasizing their understanding of low-level visual features. Additionally, Delétang et al.[11] demonstrates that LLMs, when viewed as compressors, can out-

perform traditional codecs like PNG in lossless compression for grayscale images, highlighting their potential in this field.

3. Methodology

The overall framework of our proposed lossless image compression pipeline is illustrated in Fig. 1. The original image \mathbf{x} is first compressed using a lossy codec, producing a lossy reconstructed image \mathbf{x}_l . Then we divide \mathbf{x}_l and the residual image \mathbf{r} into non-overlapping patches of size $p \times p$, denoted as $\{\mathbf{x}_l^1, \dots, \mathbf{x}_l^N\}$ and $\{\mathbf{r}^1, \dots, \mathbf{r}^N\}$, where N represents the total number of patches. During the encoding process, each patch is processed independently. We predict the probability distribution of each pixel within a residual patch in an autoregressive manner and encode these pixels using arithmetic coding. For instance, when encoding patch \mathbf{r}^n (where $n = 1, 2, \dots, N$), the entire lossy reconstruction \mathbf{x}_l and its corresponding lossy reconstructed patch \mathbf{x}_l^n are used as visual prompts to extract visual embeddings for the LLM. The pixels in residual patch \mathbf{r}^n are then autoregressively fed into the LLM to estimate the probability distribution. Given the estimated distributions, we losslessly encode \mathbf{r}^n into a bitstream using arithmetic encoding. The final bitstream comprises the lossy reconstruction \mathbf{x}_l and its corresponding residual image \mathbf{r} .

During the decoding procedure, the lossy reconstructed image \mathbf{x}_l is first decoded. Both \mathbf{x}_l and its patch \mathbf{x}_l^n are then utilized as visual prompts to autoregressively obtain the dis-

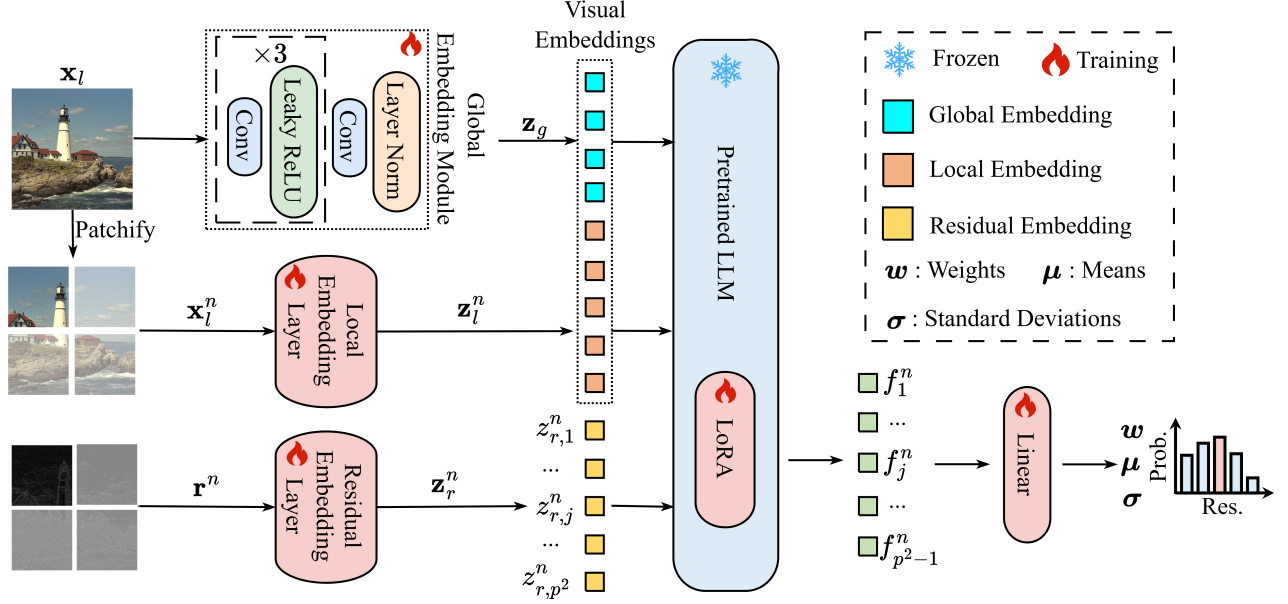


Figure 2. Our distribution estimation framework based on LLM. Visual embeddings, including the global embeddings \mathbf{z}_g and local embeddings \mathbf{z}_l^n , enhance the inference. The output feature of LLM f^n are projected onto a Gaussian Mixture Model (GMM) to estimate the residual's probability distribution.

tribution for each pixel in the residual patch \mathbf{r}^n . Finally, the full residual image is decoded, and the original image is reconstructed by combining the lossy reconstruction \mathbf{x}_l with the residual image \mathbf{r} . It is important to note that for the lossy codecs, we can either choose a traditional compression method or employ an end-to-end learned compression method. Here we use BPG [4] as the default lossy codec.

3.1. Input Embeddings

In existing LLMs, the tokenizer converts text into corresponding indexes, which are then used to obtain embeddings through an embedding layer. For image compression task, Delétang et al. [11] proposes using pixel values directly as indexes and reusing the embeddings originally trained for text dataset. However, this approach may not fully capture the relationships within the image domain, and the mismatch between textual embeddings and image pixel values may lead to poor performance. Moreover, the prompt technique, which is crucial for large language models, has been overlooked in [11].

To address the aforementioned challenges, we introduce visual prompts and visual embeddings as illustrated in Fig. 2. For compressing a residual patch, the visual prompts consist of two components: global lossy image and local lossy patch. To extract global embeddings $\mathbf{z}_g \in \mathbb{R}^{k_g \times d}$, we design a simple Global Embedding Module that utilizes several convolutional layers to capture pixel relationships from \mathbf{x}_l . For the local embeddings $\mathbf{z}_l^n \in \mathbb{R}^{p^2 \times d}$ of patch n ,

we directly use pixel values as indexes, with the embedding layer jointly optimized with the entire framework. These global and local embeddings together form visual embeddings, supplying the LLM with both global and local visual information about the image. For compressing the residual patch \mathbf{r}^n , the learnable Residual Embedding Layer extracts residual embeddings $\mathbf{z}_r^n \in \mathbb{R}^{p^2 \times d}$. These elements allow us to integrate image information with the LLM's prior knowledge, bridging the gap between image and text tasks, ultimately enhancing compression efficiency.

3.2. Distribution Estimation Using LLM

In our proposed framework, we utilize the LLM as a conditional probability estimator, leveraging the lossy reconstruction as visual prompts to predict the probability distribution of the residual image. The estimated distribution is then applied to losslessly encode the residual patch via arithmetic coding.

As illustrated in Fig. 2, the visual embeddings for the LLM consist of global embeddings \mathbf{z}_g and local embeddings \mathbf{z}_l^n . The residual is compressed in a pixel-by-pixel manner. For each pixel in the residual patch \mathbf{r}^n , we employ an autoregressive approach to estimate its probability distribution. Specifically, to predict the probability distribution of residual pixel r_j^n at position j (where $j = 1, 2, \dots, p^2$), the visual embeddings, along with previously obtained residual embeddings $\{z_{r,1}^n, \dots, z_{r,j-1}^n\}$, are concatenated into a sequence and fed into the LLM. The LLM

then outputs the corresponding prediction, calculated as follows:

$$f_j^n = F(\mathbf{z}_g, \mathbf{z}_l^n, z_{r,1}^n, \dots, z_{r,j-1}^n), \quad (1)$$

where $f_j^n \in \mathbb{R}^d$ is the output feature of the LLM for the pixel at position j .

To estimate the distribution more accurately, we go beyond directly outputting probabilities and instead predict the parameters of the probability distribution. Specifically, we introduce a Gaussian Mixture Model (GMM) [10] for effective distribution modeling. The parameters of the GMM are derived by linearly projecting the LLM output feature f_j^n . These parameters include the weights w_j^n , means μ_j^n , and standard deviations σ_j^n . Consequently, the probability distribution of the residual values can be expressed as follows:

$$p(r_j^n | \mathbf{x}_l, \mathbf{x}_l^n, r_{<j}^n) = p(r_j^n | f_j^n) \\ \sim \sum_{k=1}^K w_j^{n,(k)} \mathcal{N}(\mu_j^{n,(k)}, \sigma_j^{2n,(k)}), \quad (2)$$

where k denotes the index of mixtures, K denotes the total number of mixtures, and $\mathcal{N}(\mu, \sigma^2)$ denotes a Gaussian distribution with mean μ and standard deviation σ .

3.3. Loss Function

In our proposed method, the primary objective is to minimize the discrepancy between the estimated distribution $p(r)$ and the real distribution $q(r)$. We quantify this discrepancy using cross-entropy: the lower the cross-entropy, the closer $p(r)$ approximates $q(r)$, resulting in fewer bits required by the entropy coder to encode r . Specifically, we train our model by optimizing the following loss function:

$$\mathcal{L} = H(q, p) = \mathbb{E}_{r \sim q}[-\log p(r)] \\ = - \sum_r q(r) \log p(r) \\ = - \sum_{n=1}^N \sum_{j=1}^{p^2} \log \left\{ \sum_{k=1}^K w_j^{n,(k)} \left[c^{(k)}(r_j^n + \frac{1}{2}) \right. \right. \\ \left. \left. - c^{(k)}(r_j^n - \frac{1}{2}) \right] \right\}, \quad (3)$$

where $c^{(k)}(\cdot)$ is the cumulative distribution function of a Gaussian distribution defined by the mean $\mu_j^{n,(k)}$ and the standard deviation $\sigma_j^{n,(k)}$.

4. Experimental Results

4.1. Experimental Settings

Training Details. We train the entire framework in two stages. In the first stage, we freeze the LLM and optimize all other modules. This stage is trained on the ImageNet2012

dataset [36] using the AdamW optimizer [27] with a learning rate of 1×10^{-4} . In the second stage, we apply the LoRA [20] to finetune the LLM. For this, we utilize the DIV2K training dataset [21] to finetune the entire framework.

In this paper, we use LLaMA3-8B [12] as the default LLM, unless otherwise specified. The original images are lossy compressed using BPG [4] with the compression parameter of $Q = 28$. The lossy reconstructions and the original images are then randomly cropped into patch pairs of size 16×16 , which serve as inputs to the model. To model the distribution, we employ a Gaussian Mixture Model with $K = 5$.

Our method is implemented using the PyTorch framework [35] and requires 3 days to train the entire model on 4 NVIDIA A100 GPUs. Additionally, the arithmetic coding is implemented using the yaecl tool library [51].

Datasets. To evaluate the performance of the model, we select four different datasets. 1) *DIV2K* [21]: This dataset contains 100 high-resolution color images. 2) *CLIC.mobile* [41]: The CLIC mobile validation dataset consists of 61 color images taken with mobile phones, with most images in 2K resolution. 3) *CLIC.pro* [41]: The CLIC professional validation dataset includes 41 color images captured by professional photographers, with the majority of images in 2K resolution. 4) *Kodak* [23]: This dataset contains 24 uncompressed 768×512 color images and is widely used as a benchmark for lossy image compression.

Baseline Codecs. To validate the effectiveness of our method, we compare it against eight traditional lossless image encoders: PNG [6], JPEG-LS [48], CALIC [49], JPEG2000 [38], WebP [15], BPG [4], FLIF [39], and JPEG-XL [1]. In addition, we include five representative learning-based lossless image compression methods for comparison: L3C [28], RC [29], iVPF [55], iFlow [54], and DLPR [2]. We also reproduce the LLM-based lossless image codec [11] in our experiments. Since the LLM used in their approach is not open-source, we substitute it with LLaMA3-8B as the default model while following their other settings.

Metric. We use bits per subpixel (bpsp) as the metric to evaluate the compression ratios. The bpsp is calculated by dividing the total bits in the compressed file by the number of subpixels, where each RGB pixel consists of three subpixels.

4.2. Main Results

As shown in Tab. 1, our proposed method achieves state-of-the-art lossless compression performance across all test datasets. On the high-resolution DIV2K and CLIC datasets, our approach further reduces file size by 12.3%-17.9% compared to the best traditional lossless compression scheme JPEG-XL. When compared to SOTA learning-based meth-

Category	Codec	DIV2K	CLIC.pro	CLIC.mobile	Kodak
Traditional	PNG [6]	4.23	3.93	3.93	4.35
	JPEG-LS [48]	2.99	2.82	2.53	3.16
	CALIC [49]	3.07	2.87	2.59	3.18
	JPEG2000 [38]	3.12	2.93	2.71	3.19
	WebP [15]	3.11	2.90	2.73	3.18
	BPG [4]	3.28	3.08	2.84	3.38
	FLIF [39]	2.91	2.72	2.48	2.90
	JPEG-XL [1]	2.79	2.63	2.36	2.87
Learning-based	L3C [28]	3.09	2.94	2.64	3.26
	RC [29]	3.08	2.93	2.54	-
	iVPF [55]	2.68	2.54	2.39	-
	iFlow [54]	2.57	2.44	2.26	-
	DLPR [2]	2.55	2.38	2.16	2.86
LLM-based	Delétang et al. [11]	4.25	3.99	4.12	4.84
	Ours	2.29	2.25	2.07	2.83

Table 1. Lossless image compression performance (bpsp) of our proposed method compared to other lossless image codecs on DIV2K, CLIC.pro, CLIC.mobile and Kodak datasets.

ods such as DLPR [2] and iFlow [54], our approach also demonstrates superior results. For example, the bpsp of DLPR is 2.55, while our method achieves 2.29, reflecting a 10.2% improvement. Additionally, in comparison with a LLM-based codec [11], our method reduces the bpsp from 4.84 to 2.83 on the Kodak dataset. These results clearly demonstrate that LLMs can be effectively applied to lossless image compression, surpassing even the latest SOTA compression methods. Moreover, these results underscore how our architecture, enhanced with visual prompts, significantly improves the performance of LLM-based codecs in the lossless image compression task.

4.3. Ablation Studies

To further analyze our architecture, we conduct ablation studies as shown in Tabs. 2 to 4.

Visual Prompts. We begin by establishing a simple baseline where the LLM is fixed, without the use of visual prompts. Experimental results show that introducing visual prompts, i.e. the information from lossy reconstruction, reduces the bpsp from 4.84 to 3.19, underscoring the effectiveness of visual prompts in enhancing the LLM-based compression framework.

To explore the role of visual prompts in conjunction with LoRA, we conduct the finetuning experiments based on the method of Delétang et al. [11], and the experimental results are presented in Tab. 2. It is evident that, after applying the LoRA finetuning, our visual prompts continue to achieve a performance gain of 9.8% to 12.7% on high-resolution DIV2K and CLIC datasets.

Patch Size. In our main experiments, we use a patch size of 16×16 and then extend our evaluation to 24×24 . Increasing the patch size results in a slight performance improvement, with the bpsp decreasing from 3.19 to 3.16. This enhancement can be attributed to the larger patch sizes, which allow for longer contexts that provide more information for the model to process. This additional information enhances the model’s ability to capture intricate details and relationships within the image data, ultimately facilitating better compression.

LLM Size. We conduct experiments utilizing three LLaMA models with varying parameters and test them on the Kodak dataset to evaluate the impact of LLM size on compression performance. As shown in Tab. 3, the results indicate that compression performance decreases as the model size decreases; however, the degradation in performance is not significant, as smaller models can still achieve acceptable performance.

Lossy Image Codec. In this experiment, we evaluate the impact of the quantization parameter (QP) in the BPG codec on the performance of our proposed framework. We train our framework using different QP values, with the corresponding results presented in Tab. 4. While a lower QP increases the bpsp for lossy compression, it decreases the bpsp for lossless residual compression. Experiments show the final bpsp results are similar in range between [22, 34] and the QP value has a limited influence. Based on these findings, we select BPG with a QP value of 28 as the default lossy codec in our experiment.

Additionally, our framework accommodates other lossy

Codec	DIV2K	CLIC.pro	CLIC.mobile	Kodak
Delétang et al.	4.25	3.99	4.12	4.84
Ours	2.81(-33.9%)	2.71(-32.1%)	2.50(-39.3%)	3.19(-34.1%)
Delétang et al. (after LoRA)	2.54	2.50	2.34	3.00
Ours (after LoRA)	2.29(-9.8%)	2.25(-10.0%)	2.07(-11.5%)	2.83(-5.7%)

Table 2. Performance comparison for Delétang et al. (i.e., without visual prompts) and Ours (i.e., with visual prompts).

Method	bpsp	Loss
Ours (1B)	3.24	1.6%
Ours (3B)	3.21	0.6%
Ours (8B)	3.19	-

Table 3. Results comparison by LLM size on Kodak dataset.

Lossy Codec	Lossy	Residual	Total
BPG (QP=14)	0.95	2.43	3.38
BPG (QP=22)	0.48	2.72	3.20
BPG (QP=28)	0.27	2.92	3.19
BPG (QP=34)	0.13	3.13	3.26
BPG (QP=42)	0.04	3.38	3.42
JPEG (quality=30)	0.20	3.30	3.50
JPEG (quality=50)	0.29	2.99	3.28
JPEG (quality=70)	0.40	2.96	3.36

Table 4. Ablation experiments for lossy image codecs, test results on the Kodak dataset, using bpsp as a metric.

codecs, such as JPEG, which also demonstrates similar performance trends as presented in Tab. 4. When quality settings are appropriate, both BPG and JPEG contribute positively to our architecture. However, extreme QP settings can lead to performance degradation; setting the QP too low increases the bitrate for lossy coding, while setting it too high results in excessive residuals needing lossless compression.

Orders of GMM. GMM is commonly used in image compression[2, 10]. Residual image samples often exhibit complex distributions due to their high-frequency nature, making them challenging to model. Compared to the Gaussian Single Model (GSM), where $K = 1$, GMM incorporates a minimal increase in parameters while providing significantly improved modeling capabilities. Our ablation study on the orders K in GMM indicates that $K = 5$ significantly outperforms $K = 1$, leading to a reduction in bpsp from 3.29 to 3.19. This highlights its superior ability to capture complex distributions.

4.4. Computational Complexity

Although our LLM-based codec demonstrates superior performance, surpassing classical and other learned-based codecs through its advanced intelligence, its decoding time, as shown in Tab. 5, is considerably slower than other baselines. This is primarily due to the inherent limitations of autoregressive models and the large number of parameters in LLMs.

4.5. Lossless Compression for Images Across Diverse Domains

In this section, we apply our proposed pipeline to images from various domains, including screen content images (SCIs) and medical images. Traditional codecs often require specialized tools, such as the intra block copy technique for SCIs, to improve compression performance, which introduces additional design complexity [50]. In contrast, learning-based codecs can adapt to these diverse image types through training on sufficiently large datasets. Our proposed pipeline further advances by leveraging the extensive prior information embedded in the LLM, resulting in enhanced compression performance across these diverse image types.

Screen Content Image Compression. Screen Content Images (SCIs) typically contain text and graphics, with computer-generated elements constituting over 90% of SCIs. Compared to natural images, SCIs are characterized by sharp edges, a limited color palette, high contrast, and markedly different regional complexity, often exhibiting little to no noise [32].

In this experiment, we utilize HM-SCC [50] as the default lossy codec (QP=28). We evaluate performance on the SCID dataset [33], with the results presented in Tab. 6. The results indicate that our method, finetuned on the natural image dataset (i.e., DIV2K), demonstrates competitive generalization ability and can be effectively applied to the SCI domain, achieving a 5.1% improvement over DLPR. Furthermore, finetuning on the SCI dataset DSCIC [46] significantly enhances the model’s performance within the SCI domain, reaching a SOTA level with a bpsp of 1.11, representing a substantial improvement of 10.5% compared to JPEG-XL.

Codec	Params	Enc/Dec kMACs/pixel	Enc/Dec Times (second/image)
L3C [28]	5M	252.59/431.31	8.17/7.89
DLPR [2]	37M	$1.8 \times 10^4 / 1.3 \times 10^4$	1.26/1.80
Delétang et al. [11]	8B	2.1×10^7	10.44/288.0
Ours (1B)	1B+2M	5.9×10^6	3.84/141.6
Ours (3B)	3B+3M	1.7×10^7	10.08/338.4
Ours (8B)	8B+4M	4.2×10^7	21.12/495.6

Table 5. Comparison of runtimes and kMACs on Kodak dataset.

Codec	bpsp	Gain
PNG [6]	1.79	+14.0%
BPG [4]	1.57	-
WebP [15]	1.28	-18.5%
JPEG-XL [1]	1.24	-21.0%
HM-SCC [50]	1.18	-24.8%
L3C [28]	2.67	+70.1%
DLPR [2]	1.58	+0.6%
Ours(DIV2K)	1.50	-4.5%
Ours(SCI)	1.11	-29.3%

Table 6. Applying our model to screen content image compression, test results on the SCID dataset, using bpsp as a metric.

Codec	Axial	Coronal	Sagittal
PNG [6]	5.36	4.58	5.58
JP3D [8]	4.98	4.15	5.28
JPEG-XL [1]	4.72	3.89	5.09
HEVC [40]	5.19	4.47	5.58
VVC [7]	4.96	4.10	5.32
L3C [28]	5.16	4.45	5.52
ICEC [9]	4.64	3.84	4.97
aiWave [52]	4.55	3.80	4.83
Ours	4.46	3.57	4.83

Table 7. Applying our model to medical image compression, test results on the MRNet dataset, using bpsp as a metric.

Medical Image Compression. Most medical images are 3D, producing large data volumes that challenge storage and transmission. Effective compression is crucial. Although lossy compression offers higher ratios, it risks distorting images and compromising diagnostic accuracy, potentially leading to medical errors. Thus, lossless compression is preferred for maintaining data integrity and meeting strict standards.

Traditional lossless image compression methods, such as PNG [6] and JPEG-XL [1], individually encode each slice of 3D medical images. In addition, video coding techniques like HEVC [40] and VVC [7], along with traditional medical image compression method JP3D [8], treat 3D medical images as video sequences or volumetric data. The latest learned lossless compression methods, including L3C [28], ICEC [9], and aiWave [52], are also used as baselines.

Given that medical images are three-dimensional, we split the input medical images into 3-channel slices for processing. In this experiment, we use JPEG-XL as our lossy codec, empirically setting the corresponding quality to 68. Following prior work [9], our framework is finetuned on the MRNet training dataset [5] and tested on the MRNet validation dataset. The test results are presented in Tab. 7.

Our model demonstrates superior compression performance for lossless medical image compression. For the

Axial subset, the average bpsp of the proposed method is 4.46, compared to 4.72 for JPEG-XL. Moreover, when compared to the learning-based lossless codec L3C [28], which is also finetuned on medical images in this experiment, our approach shows significantly better compression performance. On the Coronal subset, our method further saves 6.1% bit consumption compared with aiWave [52]. This improvement can be attributed to our method’s utilization of the extensive prior information embedded in LLMs, enhancing overall performance.

5. Conclusion

Our work demonstrates that LLMs hold significant potential for lossless image compression. By designing embeddings tailored for image data and incorporating visual prompts, we achieve state-of-the-art lossless compression performance. Additionally, this framework can be effectively adapted to other image compression domains, such as screen content images and medical images. While our exploration of this framework is still in its early stages, we believe that this LLM-based method has the potential to become a new paradigm for image compression in the near future.

References

- [1] Jyrki Alakuijala, Ruud Van Asseldonk, Sami Boukourt, Martin Bruse, Iulia-Maria Comşa, Moritz Firsching, Thomas Fischbacher, Evgenii Kliuchnikov, Sebastian Gomez, Robert Obryk, et al. JPEG XL next-generation image compression architecture and coding tools. In *Applications of digital image processing XLII*, pages 112–124. SPIE, 2019. 1, 2, 5, 6, 8
- [2] Yuanchao Bai, Xianming Liu, Kai Wang, Xiangyang Ji, Xiaolin Wu, and Wen Gao. Deep lossy plus residual coding for lossless and near-lossless image compression. *IEEE TPAMI*, 2024. 1, 3, 5, 6, 7, 8
- [3] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image compression with a scale hyperprior. In *ICLR*, 2018. 2
- [4] Fabrice Bellard. BPG Image format, 2018. 2, 4, 5, 6, 8
- [5] Nicholas Bien, Pranav Rajpurkar, Robyn L Ball, Jeremy Irvin, Allison Park, Erik Jones, Michael Bereket, Bhavik N Patel, Kristen W Yeom, Katie Shpanskaya, et al. Deep-learning-assisted diagnosis for knee magnetic resonance imaging: development and retrospective validation of MR-Net. *PLoS medicine*, 15(11):e1002699, 2018. 8
- [6] Thomas Boutell. Png (portable network graphics) specification version 1.0. Technical report, 1997. 1, 2, 5, 6, 8
- [7] Benjamin Bross, Ye-Kui Wang, Yan Ye, Shan Liu, Jianle Chen, Gary J. Sullivan, and Jens-Rainer Ohm. Overview of the Versatile Video Coding (VVC) Standard and its Applications. *IEEE TCSVT*, 31(10):3736–3764, 2021. 2, 8
- [8] Tim Bruylants, Peter Schelkens, and Alexis Tzannes. JP3D—extensions for three-dimensional data (part 10). *The JPEG 2000 Suite*, pages 199–227, 2009. 8
- [9] Zhenghao Chen, Shuhang Gu, Guo Lu, and Dong Xu. Exploiting intra-slice and inter-slice redundancy for learning-based lossless volumetric image compression. *IEEE TIP*, 31:1697–1707, 2022. 8
- [10] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned Image Compression With Discretized Gaussian Mixture Likelihoods and Attention Modules. In *CVPR*, 2020. 2, 5, 7
- [11] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, et al. Language modeling is compression. *arXiv preprint arXiv:2309.10668*, 2023. 1, 2, 3, 4, 5, 6, 8
- [12] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 5
- [13] Jarek Duda. Asymmetric numeral systems: entropy coding combining speed of Huffman coding with compression rate of arithmetic coding. *arXiv preprint arXiv:1311.2540*, 2013. 2
- [14] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. SEED-X: Multimodal Models with Unified Multi-granularity Compression and Generation. *arXiv preprint arXiv: 2404.14396*, 2024. 1
- [15] Google. WebP Compression Techniques. Technical Report TR-2010-1, Google, 2010. 1, 2, 5, 6, 8
- [16] Chenhui Gou, Abdulwahab Felemban, Faizan Farooq Khan, Deyao Zhu, Jianfei Cai, Hamid Rezaatofighi, and Mohamed Elhoseiny. How Well Can Vision Language Models See Image Details? *arXiv preprint arXiv: 2408.03940*, 2024. 3
- [17] Dailan He, Yaoyan Zheng, Baocheng Sun, Yan Wang, and Hongwei Qin. Checkerboard Context Model for Efficient Learned Image Compression. In *CVPR*, pages 14771–14780, 2021. 2
- [18] Dailan He, Ziming Yang, Weikun Peng, Rui Ma, Hongwei Qin, and Yan Wang. ELIC: Efficient Learned Image Compression with Unevenly Grouped Space-Channel Contextual Adaptive Coding. In *CVPR*, pages 5708–5717, 2022. 2
- [19] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In *ICML*, pages 2790–2799. PMLR, 2019. 3
- [20] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 3, 5
- [21] Andrey Ignatov, Radu Timofte, et al. PIRM challenge on perceptual image enhancement on smartphones: report. In *ECCV*, 2019. 5
- [22] Wei Jiang and Ronggang Wang. MLIC++: Linear Complexity Multi-Reference Entropy Modeling for Learned Image Compression. In *Proc. Int. Conf. Mach. Learn. Workshops*, 2023. 2
- [23] Eastman Kodak. Kodak lossless true color image suite (PhotoCD PCD0992). URL <http://r0k.us/graphics/kodak/>, 6:2, 1993. 5
- [24] Han Li, Shaohui Li, Wenrui Dai, Chenglin Li, Junni Zou, and Hongkai Xiong. Frequency-Aware Transformer for Learned Image Compression. In *ICLR*, 2024. 2
- [25] Jinming Liu, Heming Sun, and Jiro Katto. Learned Image Compression with Mixed Transformer-CNN Architectures. In *CVPR*, pages 14388–14397, 2023. 2
- [26] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *ICCV*, pages 9992–10002, 2021. 2
- [27] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5
- [28] Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. Practical full resolution learned lossless image compression. In *CVPR*, pages 10629–10638, 2019. 1, 2, 5, 6, 8
- [29] Fabian Mentzer, Luc Van Gool, and Michael Tschannen. Learning better lossless compression using lossy compression. In *CVPR*, pages 6638–6647, 2020. 1, 3, 5, 6
- [30] David Minnen and Saurabh Singh. Channel-Wise Autoregressive Entropy Models for Learned Image Compression. In *ICIP*, pages 3339–3343, 2020. 2

- [31] David Minnen, Johannes Ballé, and George Toderici. Joint Autoregressive and Hierarchical Priors for Learned Image Compression. In *NeurIPS*, pages 10794–10803, 2018. 2
- [32] Tung Nguyen, Xiaozhong Xu, Felix Henry, Ru-Ling Liao, Mohammed Golam Sarwer, Marta Karczewicz, Yung-Hsuan Chao, Jizheng Xu, Shan Liu, Detlev Marpe, et al. Overview of the screen content support in VVC: Applications, coding tools, and performance. *IEEE TCSVT*, 31(10):3801–3817, 2021. 7
- [33] Zhangkai Ni, Lin Ma, Huanqiang Zeng, Ying Fu, Lu Xing, and Kai-Kuang Ma. SCID: A database for screen content images quality assessment. In *2017 International Symposium on Intelligent Signal Processing and Communication Systems (ISPACS)*, pages 774–779, 2017. 7
- [34] Ziqi Pang, Ziyang Xie, Yunze Man, and Yu-Xiong Wang. Frozen Transformers in Language Models Are Effective Visual Encoder Layers. In *ICLR*, 2024. 1
- [35] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. 2017. 5
- [36] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 115(3):211–252, 2015. 5
- [37] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. Pixelcnn++: A pixelcnn implementation with discretized logistic mixture likelihood and other modifications. In *ICLR*, 2017. 2
- [38] Athanassios Skodras, Charilaos Christopoulos, and Touradj Ebrahimi. The JPEG 2000 still image compression standard. *IEEE Signal processing magazine*, 18(5):36–58, 2001. 5, 6
- [39] Jon Sneyers and Pieter Wuille. FLIF: Free lossless image format based on MANIAC compression. In *ICIP*, pages 66–70. IEEE, 2016. 1, 2, 5, 6
- [40] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high efficiency video coding (HEVC) standard. *IEEE TCSVT*, 22(12):1649–1668, 2012. 8
- [41] George Toderici, Wenzhe Shi, Radu Timofte, Lucas Theis, Johannes Balle, Eirikur Agustsson, Nick Johnston, and Fabian Mentzer. Workshop and Challenge on Learned Image Compression (CLIC2020), 2020. 5
- [42] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. In *NeurIPS*, 2016. 2
- [43] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, pages 1747–1756. PMLR, 2016. 2
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *NeurIPS*, pages 5998–6008, 2017. 2
- [45] Gregory K. Wallace. The JPEG Still Picture Compression Standard. *Communication ACM*, 34(4):30–44, 1991. 2
- [46] Feifeng Wang, Liquan Shen, Qi Teng, and Zhaoyi Tian. DSCIC: Deep Screen Content Image Compression. *IEEE TCSVT*, pages 1–1, 2024. 7
- [47] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, and Jifeng Dai. VisionLLM: Large Language Model is also an Open-Ended Decoder for Vision-Centric Tasks. In *NeurIPS*, 2023. 3
- [48] Marcelo J Weinberger, Gadiel Seroussi, and Guillermo Sapiro. The LOCO-I lossless image compression algorithm: Principles and standardization into JPEG-LS. *IEEE TIP*, 9(8):1309–1324, 2000. 5, 6
- [49] Xiaolin Wu and Nasir Memon. Context-based, adaptive, lossless image coding. *IEEE transactions on Communications*, 45(4):437–444, 1997. 5, 6
- [50] Jizheng Xu, Rajan Joshi, and Robert A. Cohen. Overview of the Emerging HEVC Screen Content Coding Extension. *IEEE TCSVT*, 26(1):50–62, 2016. 7, 8
- [51] Tongda Xu, Han Gao, Chenjian Gao, Jinyong Pi, Yanghao Li, Yuanyuan Wang, Ziyu Zhu, Dailan He, Mao Ye, Hongwei Qin, et al. Bit allocation using optimization. *arXiv preprint arXiv:2209.09422*, 2022. 5
- [52] Dongmei Xue, Haichuan Ma, Li Li, Dong Liu, and Zhiwei Xiong. aiWave: Volumetric image compression with 3-D trained affine wavelet-like transform. *IEEE Transactions on Medical Imaging*, 42(3):606–618, 2022. 8
- [53] Senqiao Yang, Tianyuan Qu, Xin Lai, Zhuotao Tian, Bohao Peng, Shu Liu, and Jiaya Jia. An Improved Baseline for Reasoning Segmentation with Large Language Model. *arXiv preprint arXiv: 2312.17240*, 2023. 3
- [54] Shifeng Zhang, Ning Kang, Tom Ryder, and Zhenguo Li. iflow: Numerically invertible flows for efficient lossless compression via a uniform coder. In *NeurIPS*, pages 5822–5833, 2021. 1, 2, 5, 6
- [55] Shifeng Zhang, Chen Zhang, Ning Kang, and Zhenguo Li. iVPF: Numerical invertible volume preserving flow for efficient lossless compression. In *CVPR*, pages 620–629, 2021. 1, 2, 5, 6
- [56] Boyang Zheng, Jinjin Gu, Shijun Li, and Chao Dong. LM4LV: A Frozen Large Language Model for Low-level Vision Tasks. *arXiv preprint arXiv:2405.15734*, 2024. 1, 3
- [57] Yinhao Zhu, Yang Yang, and Taco Cohen. Transformer-based Transform Coding. In *ICLR*, 2022. 2
- [58] Renjie Zou, Chunfeng Song, and Zhaoxiang Zhang. The Devil Is in the Details: Window-based Attention for Image Compression. In *CVPR*, pages 17471–17480, 2022. 2