

# An Information-Theoretic Regularizer for Lossy Neural Image Compression

## Supplementary Material

### 6. Proof of Lemma 1

For a deterministic quantization process  $Q(\cdot)$ , the conditional probability  $p(U|X)$  can only take values of 0 or 1, i.e.,

$$p(U|X) = \begin{cases} 1, & \text{if } U = Q(X) \\ 0, & \text{if } U \neq Q(X) \end{cases}. \quad (17)$$

Similarly, for a deterministic dequantization process  $Q^{-1}(\cdot)$ , it is by definition a bijection function, i.e., given the index  $U$  or the reconstruction  $\hat{X}$ , we can uniquely determine the corresponding value of  $\hat{X}$  or  $U$ , respectively. Therefore, we have

$$p(\hat{X}|U) = \begin{cases} 1, & \text{if } \hat{X} = Q^{-1}(U) \\ 0, & \text{if } \hat{X} \neq Q^{-1}(U) \end{cases}, \quad (18)$$

$$p(U|\hat{X}) = \begin{cases} 1, & \text{if } U = Q(\hat{X}) \\ 0, & \text{if } U \neq Q(\hat{X}) \end{cases}, \quad (19)$$

and overall,

$$p(\hat{X}|X) = \begin{cases} 1, & \hat{X} = Q^{-1}(Q(X)) \\ 0, & \hat{X} \neq Q^{-1}(Q(X)) \end{cases}. \quad (20)$$

Thus, given Eqn.(18) to Eqn.(20), the following conditional entropy is by definition equal to 0:

$$H(\hat{X}|U) = H(U|\hat{X}) = H(\hat{X}|X) = 0. \quad (21)$$

Furthermore, from the equality of mutual information  $I(X; \hat{X})$  and the chain rule of joint entropy  $H(\hat{X}, U)$ , we have

$$\begin{cases} I(X; \hat{X}) = H(\hat{X}) - H(\hat{X}|X) \\ H(U|\hat{X}) + H(\hat{X}) = H(\hat{X}|U) + H(U) \end{cases}. \quad (22)$$

Substituting Eqn.(21) into Eqn.(22), we conclude

$$\begin{cases} I(X; \hat{X}) = H(\hat{X}) \\ H(\hat{X}) = H(U) \end{cases}, \quad (23)$$

and thus

$$H(U) = I(X; \hat{X}). \quad (24)$$

### 7. Proof of Lemma 2

Recalling the proof in Sec. 6, since both the analysis transform  $T_A(\cdot)$  and synthesis transform  $T_S(\cdot)$  are deterministic, the following holds:

$$p(U|X) = \begin{cases} 1, & \text{if } U = Q(T_A(X)) \\ 0, & \text{if } U \neq Q(T_A(X)) \end{cases}, \quad (25)$$

$$p(\hat{X}|U) = \begin{cases} 1, & \text{if } \hat{X} = T_S(Q^{-1}(U)) \\ 0, & \text{if } \hat{X} \neq T_S(Q^{-1}(U)) \end{cases}, \quad (26)$$

$$p(\hat{X}|X) = \begin{cases} 1, & \text{if } \hat{X} = T_S(Q^{-1}(Q(T_A(X)))) \\ 0, & \text{if } \hat{X} \neq T_S(Q^{-1}(Q(T_A(X)))) \end{cases}, \quad (27)$$

and therefore

$$H(\hat{X}|U) = H(\hat{X}|X) = 0. \quad (28)$$

Herein, the core distinction between the transform coding and direct coding models lies in the fact that the synthesis transform  $T_S(\cdot)$ , unlike the dequantization function  $Q^{-1}(\cdot)$ , does not inherently guarantee a bijective mapping, particularly in the context of neural transforms. Consequently, given the reconstruction  $\hat{X}$ , there may be uncertainty in  $\hat{Y}$  and thus  $U$ , i.e.,

$$H(U|\hat{X}) \neq 0. \quad (29)$$

Substituting Eqn.(28) and Eqn.(29) into Eqn.(22), we can conclude

$$\begin{cases} I(X; \hat{X}) = H(\hat{X}) \\ H(\hat{X}) = H(U) - H(U|\hat{X}) \end{cases}, \quad (30)$$

and thus

$$H(U) = I(X; \hat{X}) + H(U|\hat{X}). \quad (31)$$

Table 3. BD-Rate comparison between our reproduction and the pre-trained models [44].

	Kodak	Tecnick	CLIC	Average
<i>hyperprior</i>	0.52%	1.88%	-1.24%	0.39%
<i>autoregressive</i>	0.56%	2.25%	-1.56%	0.41%
<i>attention</i>	-3.32%	-3.33%	-5.68%	-4.11%

### 8. Reproduced baselines

We retrain the *hyperprior* [4], *autoregressive* [8], and *attention* [10] models from scratch, adhering to the default implementation and training configurations of CompressAI's [44]. Four bit-rate points, i.e.,  $\lambda \in \{0.0018, 0.0035, 0.0067, 0.0130\}$  are trained with  $2 \times 10^6$  steps. The evaluation results are summarized in Table 3. On

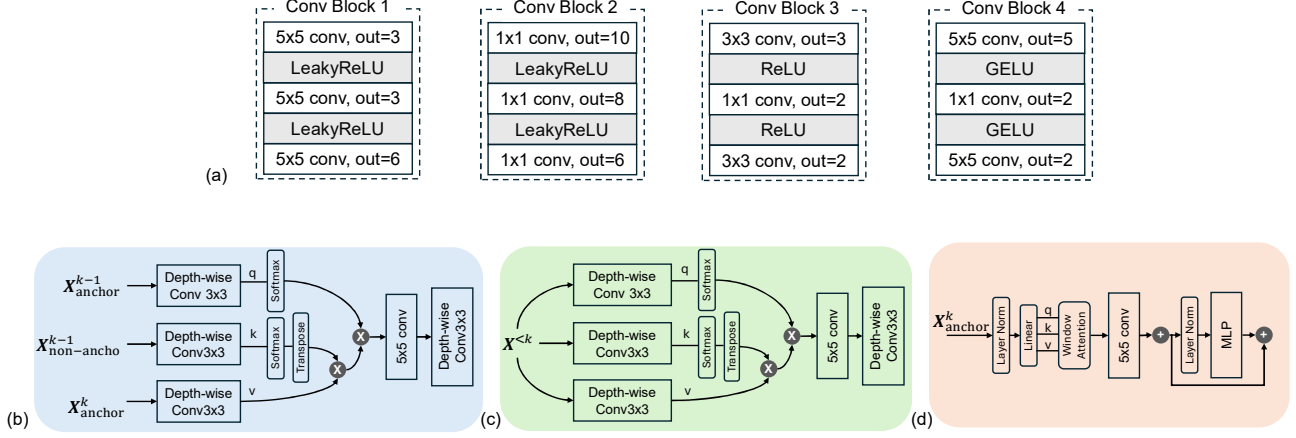


Figure 7. Details of (a) Four convolutional modules. (b) Inter attention module, (c) Intra attention module, and (d) Checkerboard attention module from MLIC++ [12].

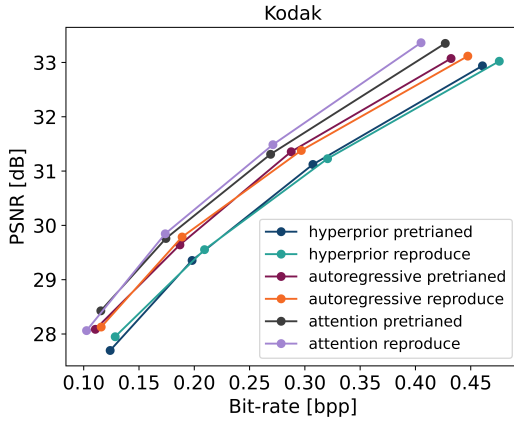


Figure 8. Comparison between the pre-trained models from CompressAI's and our reproduction on the Kodak dataset. Our reproduction achieves 0.52%, 0.56%, and -3.32% BD-Rates for the *hyperprior*, *autoregressive*, and *attention* models, respectively.

average, when compared to the pre-trained models of CompressAI's, our reproduction models yield 0.39%, 0.41%, and -4.11% BD-Rates [57] for the *hyperprior*, *autoregressive*, and *attention* models, respectively. Notably, on the CLIC dataset, our reproduction models outperform their pre-trained counterparts by -1.24%, -1.56%, and -5.68% BD-Rates for the *hyperprior*, *autoregressive*, and *attention* models, respectively. For the *attention* model, our reproduction model outperforms the pre-trained model on all three evaluation datasets of Kodak, CLIC, and Tecnick by -3.32%, -3.33%, and -5.68% BD-Rates, respectively. Note that in the original *attention* network [10], the latent is modeled as a mixture of Gaussians with 3 clusters. In CompressAI's default implementation, this cluster number is simplified to 1. We follow this latent implementation. The de-

tailed rate-distortion curves for the Kodak dataset are visualized in Fig. 8.

## 9. More details on source entropy models

Herein, the module designs are identical to the latent designs [4, 8, 10, 12, 43], with only minor dimension adjustment to fit  $X$ . The details are depicted in Fig. 7. For the attention modules from MLIC++, the depthwise separable convolution is adopted [56].

## 10. Regularization with $H(U|\hat{X})$ estimation

As shown in Fig. 9(a), we further examine the regularization performance by training compression networks with the following minimization objective:

$$R + \lambda D + \alpha \underbrace{\mathbb{E}_{\mathbf{X}} [\log q_{\theta}(\mathbf{X}|\hat{\mathbf{X}}) - \log q_{\varphi}(\mathbf{U}|\hat{\mathbf{X}})]}_{\approx -(H(\mathbf{X}|\hat{\mathbf{X}}) - H(\mathbf{U}|\hat{\mathbf{X}}))}, \quad (32)$$

where an additional entropy model  $q_{\varphi}(\mathbf{U}|\hat{\mathbf{X}})$  is jointly estimated with  $q_{\theta}(\mathbf{X}|\hat{\mathbf{X}})$  by maximizing

$$\max \mathbb{E}_{\mathbf{X}} [\log q_{\theta}(\mathbf{X}|\hat{\mathbf{X}}) + \log q_{\varphi}(\mathbf{U}|\hat{\mathbf{X}})]. \quad (33)$$

The training procedure follows Algorithm 1, consisting of alternating steps: one step updates the compression network using Eqn.(32), while another step updates the regularizer using Eqn.(33). The distribution  $q_{\varphi}(\mathbf{U}|\hat{\mathbf{X}})$  is modeled as a factorized Gaussian, with its neural architecture identical to that of the analysis transform  $T_A$ . We evaluate this approach on the *hyperprior* model, with the regularization factor  $\alpha$  set to 0.1 and all other training settings consistent with Sec. 4.1. The corresponding results are presented in Fig. 9(b). As shown in Fig. 9(b), introducing the  $H(\mathbf{U}|\hat{\mathbf{X}})$  term only provides a marginal gain (approximately -0.2% BD-Rate) on the CLIC dataset at  $0.5 \times 10^6$

steps, while no significant improvement is observed in other cases. One possible explanation is that learning the distribution  $q_\varphi(\mathbf{U}|\hat{\mathbf{X}})$  is highly challenging, as it requires training an additional “encoder” function to accurately map  $\hat{\mathbf{X}}$  to  $\mathbf{U}$ . Moreover, due to computational complexity considerations, the proposed Algorithm 1 updates the regularization terms with only one gradient step per iteration. Experiments reveal that even with this single-step update strategy, the overall training complexity has already increased from approximately 28% to 45%. One might suggest reusing the entire analysis transform  $T_A$  instead of learning a new mapping from  $\hat{\mathbf{X}}$  to  $\mathbf{U}$  to reduce complexity. We have explored this alternative, but its performance is inferior to the results presented here.

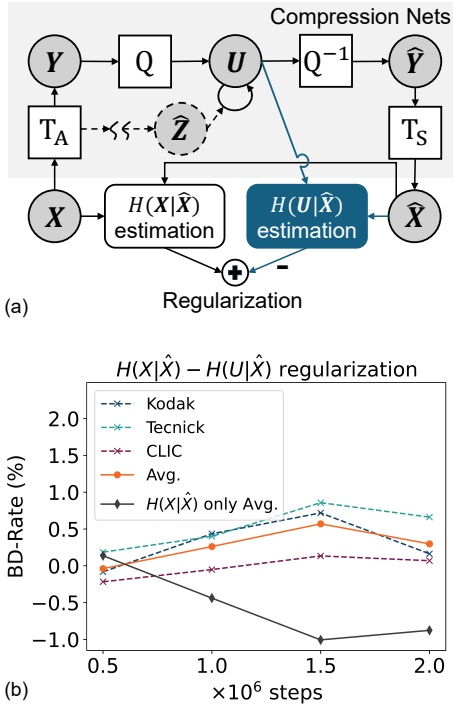


Figure 9. (a) Illustration of regularization with an additional  $H(\mathbf{U}|\hat{\mathbf{X}})$  term; (b) Performance evaluation on the *hyperprior* model.