

---

# Multi-Sample Training for Neural Image Compression

---

Tongda Xu<sup>1,2</sup>, Yan Wang<sup>1,2</sup>\*, Dailan He<sup>1</sup>, Chenjian Gao<sup>1,3</sup>, Han Gao<sup>1,4</sup>,  
Kunzan Liu<sup>1,5</sup>, Hongwei Qin<sup>1</sup>

<sup>1</sup>SenseTime Research, <sup>2</sup>Institute for AI Industry Research (AIR), Tsinghua University,

<sup>3</sup>Beihang University, <sup>4</sup>University of Electronic Science and Technology of China,

<sup>5</sup>Department of Electronic Engineering, Tsinghua University

{xutongda, wangyan}@air.tsinghua.edu.cn,

{hedailan, gaochenjian, gaohan1, liukunzan, qinhongwei}@sensetime.com

## Abstract

This paper considers the problem of lossy neural image compression (NIC). Current state-of-the-art (sota) methods adopt uniform posterior to approximate quantization noise, and single-sample pathwise estimator to approximate the gradient of evidence lower bound (ELBO). In this paper, we propose to train NIC with multiple-sample importance weighted autoencoder (IWAE) target, which is tighter than ELBO and converges to log likelihood as sample size increases. First, we identify that the uniform posterior of NIC has special properties, which affect the variance and bias of pathwise and score function estimators of the IWAE target. Moreover, we provide insights on a commonly adopted trick in NIC from gradient variance perspective. Based on those analysis, we further propose multiple-sample NIC (MS-NIC), an enhanced IWAE target for NIC. Experimental results demonstrate that it improves sota NIC methods. Our MS-NIC is plug-and-play, and can be easily extended to other neural compression tasks.

## 1 Introduction

Latent variable-based lossy neural image compression (NIC) has witnessed significant success. The majority of NIC follows the framework proposed by Ballé et al. [2017]: For encoding, the original image  $\mathbf{x}$  is transformed into  $\mathbf{y}$  by the encoder. Then  $\mathbf{y}$  is scalar-quantized into integer  $\tilde{\mathbf{y}}$ , estimated with an entropy model  $p(\tilde{\mathbf{y}})$  and coded. For decoding,  $\tilde{\mathbf{y}}$  is transformed back by the decoder to obtain reconstructed  $\tilde{\mathbf{x}}$ . The optimization target of NIC is R-D cost:  $R + \lambda D$ .  $R$  denotes the bitrate of  $\tilde{\mathbf{y}}$ ,  $D$  denotes the distortion between  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ , and  $\lambda$  denotes the hyper-parameter controlling their trade-off. During training, the quantization  $\tilde{\mathbf{y}} = \lfloor \mathbf{y} \rfloor$  is relaxed with  $\tilde{\mathbf{y}} = \mathbf{y} + \epsilon$  to simulate the quantization noise. And  $\epsilon$  is fully factorized uniform noise  $\epsilon \sim p(\epsilon) = \prod \mathcal{U}(-\frac{1}{2}, +\frac{1}{2})$ .

Ballé et al. [2017] further recognises that such training framework is closely related to variational inference. Indeed, the above process can be formulated as a graphic model  $\mathbf{x} \leftarrow \tilde{\mathbf{y}}$ . During encoding,  $\mathbf{x}$  is transformed into variational parameter  $\mathbf{y}$  by inference model (encoder), and  $\tilde{\mathbf{y}}$  is sampled from variational posterior  $q(\tilde{\mathbf{y}}|\mathbf{x})$ , which is a unit uniform distribution centered in  $\mathbf{y}$ . The prior likelihood  $p(\tilde{\mathbf{y}})$  is computed, and  $\tilde{\mathbf{y}}$  is transformed back by the generative model (decoder) to compute the likelihood  $p(\mathbf{x}|\tilde{\mathbf{y}})$ . Under such formulation, the prior is connected to the bitrate, the likelihood is connected to the distortion, and the posterior likelihood is connected to the bits-back bitrate (See Appendix. 2.3), which is 0 in NIC. Finally, the evidence lower bound (ELBO) is the negative  $R + \lambda D$  target (Eq. 1). Denote the transform function  $\tilde{\mathbf{y}}(\epsilon; \phi) = \mathbf{y} + \epsilon$ , and sampling  $\tilde{\mathbf{y}} \sim q(\tilde{\mathbf{y}}|\mathbf{x})$  is equivalent to transforming  $\epsilon$  through  $\tilde{\mathbf{y}}(\epsilon; \phi)$ . Then the gradient of ELBO is estimated via pathwise estimator with single-sample Monte Carlo (Eq. 2). This is the same as SGVB-1 [Kingma and Welling, 2013].

---

\*Yan Wang is the corresponding author.

$$\mathcal{L} = -(R + \lambda D) = \mathbb{E}_{q(\tilde{\mathbf{y}}|\mathbf{x})} \left[ \underbrace{\log p(\mathbf{x}|\tilde{\mathbf{y}})}_{\text{- distortion}} + \underbrace{\log p(\tilde{\mathbf{y}})}_{\text{- rate}} - \underbrace{\log q(\tilde{\mathbf{y}}|\mathbf{x})}_{\text{bits-back rate: 0}} \right] \quad (1)$$

$$\nabla_{\phi} \mathcal{L} = \mathbb{E}_{p(\epsilon)} [\nabla_{\phi} (\log \frac{p(\mathbf{x}, \tilde{\mathbf{y}}(\epsilon; \phi))}{q(\tilde{\mathbf{y}}(\epsilon; \phi)|\mathbf{x})})] \approx \nabla_{\phi} \log \frac{p(\mathbf{x}, \tilde{\mathbf{y}}(\epsilon; \phi))}{q(\tilde{\mathbf{y}}(\epsilon; \phi)|\mathbf{x})} \quad (2)$$

Ballé et al. [2018] further extends this framework into a two-level hierarchical structure, with graphic model  $\mathbf{x} \leftarrow \tilde{\mathbf{y}} \leftarrow \tilde{\mathbf{z}}$ . The variational posterior is fully factorized uniform distribution  $\mathcal{U}(\mathbf{y} - \frac{1}{2}, \mathbf{y} + \frac{1}{2})\mathcal{U}(z - \frac{1}{2}, z + \frac{1}{2})$  To simulate the quantization noise. And  $\mathbf{y}, z$  denote outputs of their inference networks.

$$\mathcal{L} = \mathbb{E}_{q(\tilde{\mathbf{y}}, \tilde{\mathbf{z}}|\mathbf{x})} \left[ \underbrace{\log p(\mathbf{x}|\tilde{\mathbf{y}})}_{\text{- distortion}} + \underbrace{\log p(\tilde{\mathbf{y}}|\tilde{\mathbf{z}})}_{\text{- rate}} + \underbrace{\log p(\tilde{\mathbf{z}})}_{\text{bits-back rate: 0}} - \log q(\tilde{\mathbf{y}}|\mathbf{x}) - \log q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}}) \right] \quad (3)$$

The majority of later NIC follows this hierarchical latent framework [Minnen et al., 2018, Cheng et al., 2020]. Some focus on more expressive network architectures [Zhu et al., 2021, Xie et al., 2021], some stress better context models [Minnen and Singh, 2020, He et al., 2021, Guo et al., 2021a], and some emphasize semi-amortization inference [Yang et al., 2020]. However, there is little research on multiple-sample methods, or other techniques for a tighter ELBO.

On the other hand, IWAE [Burda et al., 2016] has been successful in density estimation. Specifically, IWAE considers a multiple-sample lowerbound  $\mathcal{L}_k$  (Eq. 4), which is tighter than its single-sample counterpart. The benefit of such bound is that the implicit distribution defined by IWAE approaches true posterior as  $k$  increases [Cremer et al., 2017]. This suggests that its variational posterior is less likely to collapse to a single mode of true posterior, and the learned representation is richer. The gradient of  $\mathcal{L}_k$  is computed via pathwise estimator. Denote the exponential ELBO sample as  $w_i$ , its reparameterization as  $w(\epsilon_i; \phi)$ , and its weight  $\tilde{w}_i = \frac{w_i}{\sum w_j}$ . Then  $\nabla_{\phi} \mathcal{L}_k$  has the form of importance weighted sum (Eq. 5).

$$\mathcal{L}_k = \mathbb{E}_{q(\tilde{\mathbf{y}}_{1:k}|\mathbf{x})} \left[ \log \frac{1}{k} \sum_i^k \underbrace{\frac{p(\mathbf{x}, \tilde{\mathbf{y}}_i)}{q(\tilde{\mathbf{y}}_i|\mathbf{x})}}_{w_i} \right] = \mathbb{E}_{p(\epsilon_{1:k})} \left[ \log \frac{1}{k} \sum_i^k \underbrace{\frac{p(\mathbf{x}, \tilde{\mathbf{y}}(\epsilon_i; \phi))}{q(\tilde{\mathbf{y}}(\epsilon_i; \phi)|\mathbf{x})}}_{w(\epsilon_i; \phi)} \right] \quad (4)$$

$$\nabla_{\phi} \mathcal{L}_k = \mathbb{E}_{p(\epsilon_{1:k})} \left[ \sum_i^k \tilde{w}_i \nabla_{\phi} \log w(\epsilon_i; \phi) \right] \approx \sum_i^k \tilde{w}_i \nabla_{\phi} \log w(\epsilon_i; \phi) \quad (5)$$

In this paper, we consider the problem of training NIC with multiple-sample IWAE target (Eq. 4), which allows us to learn a richer latent space. First, we recognise that NIC’s factorized uniform variational posterior has impacts on variance and bias properties of gradient estimators. Specifically, we find NIC’s pathwise gradient estimator equivalent to an improved STL estimator [Roeder et al., 2017], which is unbiased even for the IWAE target. However, NIC’s IWAE-DReG estimator [Tucker et al., 2018] has extra bias, which causes performance decay. Moreover, we provide insights on a commonly adopted but little explained trick of training NIC from gradient variance perspective. Based on those analysis and observations, we further propose MS-NIC, a novel improvement of multiple-sample IWAE target for NIC. Experimental results show that it improves sota NIC methods [Ballé et al., 2018, Cheng et al., 2020] and learns richer latent representation. Our method is plug-and-play, and can be extended into neural video compression.

To wrap up, our contributions are as follows:

- We provide insights on the impact of the uniform variational posterior upon gradient estimators, bits-back coding and a commonly adopted but little discussed trick of NIC training from gradient variance perspective.
- We propose multiple-sample neural image compression (MS-NIC). It is a novel enhancement of hierarchical IWAE [Burda et al., 2016] for neural image compression. To the best of our knowledge, we are the first to consider a tighter ELBO for training neural image compression.

- We demonstrate the efficiency of MS-NIC through experimental results on sota NIC methods. Our method is plug-and-play for neural image compression and can be easily applied to neural video compression.

## 2 Gradient Estimation for Neural Image Compression

The common NIC framework (Eq. 1, Eq 3) adopts fully factorized uniform distribution  $q(\tilde{\mathbf{y}}, \tilde{\mathbf{z}}|\mathbf{x}) = \prod \mathcal{U}(y^i - \frac{1}{2}, y^i + \frac{1}{2}) \prod \mathcal{U}(z^j - \frac{1}{2}, z^j + \frac{1}{2})$  to simulate the quantization noise. Such formulation has the following special properties:

- Property I:  $q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$  and  $q(\tilde{\mathbf{y}}|\mathbf{x})$ 's support depends on the parameter.
- Property II:  $\log q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}}) = \log q(\tilde{\mathbf{y}}|\mathbf{x}) = 0$  on their support.

The impacts of these two properties are frequently neglected in previous works, which does not influence the results for single-sample pathwise gradient estimators (a.k.a. reparameterization trick in Kingma and Welling [2013]). In this section, we discuss the impacts of these two properties upon the variance and biasness of gradient estimators. Our analysis is based on single level latent (Eq. 1) instead of hierarchical latent (Eq. 3) to simplify notations.

### 2.1 Impact on Pathwise Gradient Estimators

First, let's consider the single-sample case. We can expand the pathwise gradient of ELBO in Eq. 2 into Eq. 6. As indicated in the equation,  $\phi$  contributes to  $\mathcal{L}$  in two ways. The first way is through the reparametrized  $\tilde{\mathbf{y}}(\epsilon; \phi)$  (pathwise term), and the other way is through the parameter of  $\log q(\tilde{\mathbf{y}}|\mathbf{x})$  (parameter score term). Generally, the parameter score term has higher variance than the pathwise term. The STL [Roeder et al., 2017] reduces the gradient by dropping the score. It is unbiased since the dropped term's expectation  $\mathbb{E}_{q(\tilde{\mathbf{y}}|\mathbf{x})}[\nabla_{\phi} \log q_{\phi}(\tilde{\mathbf{y}}|\mathbf{x})]$  is 0.

$$\nabla_{\phi} \mathcal{L} = \mathbb{E}_{p(\epsilon)} \left[ \underbrace{\nabla_{\tilde{\mathbf{y}}} \left( \log \frac{p(\mathbf{x}|\tilde{\mathbf{y}})p(\tilde{\mathbf{y}})}{q(\tilde{\mathbf{y}}|\mathbf{x})} \right)}_{\text{pathwise term}} \nabla_{\phi} \tilde{\mathbf{y}}(\epsilon; \phi) - \underbrace{\nabla_{\phi} \log q_{\phi}(\tilde{\mathbf{y}}|\mathbf{x})}_{\text{parameter score term}} \right] \quad (6)$$

Now let's consider the STL estimator of multiple-sample IWAE bound (Eq. 4). As shown in Tucker et al. [2018], the STL estimation of IWAE bound gradient is biased. To reveal the reason, consider expanding the gradient Eq. 5 into partial derivatives as we expand Eq. 2 into Eq. 6. Unlike single-sample case, the dropped parameter score term  $\mathbb{E}_{p(\epsilon_{1:k})}[\sum \tilde{w}_i (-\nabla_{\phi} \log q_{\phi}(\tilde{\mathbf{y}}|\mathbf{x}))]$  is no longer 0 due to the importance weight  $\tilde{w}_i$ . This means that STL loses its unbiasedness in general IWAE cases.

Regarding NIC, however, the direct pathwise gradient for IWAE bound is automatically an unbiased STL estimator. Property II means that variational posterior has constant entropy, which further means that the parameter score gradient is 0. So, NIC's pathwise gradient of IWAE bound is equivalent to an extended, unbiased STL estimator.

### 2.2 Impact on Score Function Gradient Estimators

In previous section, we show the blessing of NIC's special properties on pathwise gradient estimators. In this section, we show their curse on score function gradient estimators. Specifically, Property I implies that  $q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$  and  $q(\tilde{\mathbf{y}}|\mathbf{x})$  are not absolute continuous, and hence the score function gradient estimators of those distributions are biased.

For example, consider a univariate random variable  $x \sim p_{\theta}(x) = \mathcal{U}(\theta - \frac{1}{2}, \theta + \frac{1}{2})$ . Our task is to estimate the gradient of a differentiable function  $f(x)$ . And consider the  $\theta$ -independent random variable  $\epsilon \sim p(\epsilon) = \mathcal{U}(-\frac{1}{2}, +\frac{1}{2})$ , the transform  $x(\epsilon; \theta) = \theta + \epsilon$ . Under such conditions, the Monte Carlo estimated pathwise gradient and score function gradient are:

$$\text{pathwise gradient: } \nabla_{\theta} \mathbb{E}_{p_{\theta}(x)}[f(x)] = \nabla_{\theta} \mathbb{E}_{p(\epsilon)}[f(x(\epsilon; \theta))] \approx \frac{1}{N} \sum_i^N \nabla_{\theta} f(\theta + \epsilon_i) \quad (7)$$

$$\text{score function gradient: } \nabla_{\theta} \mathbb{E}_{p_{\theta}(x)}[f(x)] = \mathbb{E}_{p_{\theta}(x)}[\nabla_{\theta} \log p_{\theta}(x) f(x)] = 0 \quad (8)$$

Eq. 7 does not equal to Eq. 8, and Eq.8 is wrong. The score function gradient is only unbiased when the distribution satisfies the absolute continuity condition of [Mohamed et al., 2020]. This reflects that under the formulation of NIC, the equivalence between the score function gradient (a.k.a. REINFORCE [Williams, 1992]) and pathwise gradient (a.k.a reparameterization trick in [Kingma and Welling, 2013]) no longer holds.

Table 1: Effect of DReG gradient estimator in NIC.

	Sample Size	bpp	MSE	PSNR (db)	R-D cost
<i>Single-sample</i>					
Baseline [Ballé et al., 2018]	-	0.5273	32.61	33.28	1.017
<i>Multiple-sample</i>					
MS-NIC-MIX(pathwise gradient)	5	0.5259	31.84	33.38	1.003
MS-NIC-MIX(DReG gradient)	5	0.5316	35.09	32.90	1.058

Such equivalence is the cornerstone of many gradient estimators, and IWAE-DReG [Tucker et al., 2018] is one of them. IWAE-DReG is a popular gradient estimator for IWAE target (Eq. 4) as it resolves the vanish of inference network gradient SNR (signal to noise ratio). However, the correctness of IWAE-DReG depends on the equivalence between the score function gradient and pathwise gradient, which does not hold for NIC. Specifically, IWAE-DReG expand the total derivative of IWAE target as Eq. 9 and perform another round of reparameterization on the score function term as Eq. 10 to further reduce the gradient variance. However, Eq. 10 requires the equivalence of pathwise gradient and score function gradient.

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(\tilde{\mathbf{y}}_{1:k}|\mathbf{x})} [\log \frac{1}{k} \sum_{i=1}^k w_i] = \mathbb{E}_{p(\epsilon_{1:k})} [\underbrace{\sum_{i=1}^k \frac{w_i}{\sum_{j=1}^k w_j} \left( -\frac{\partial \log q_{\phi}(\tilde{\mathbf{y}}_i|\mathbf{x})}{\partial \phi} + \frac{\partial \log w(\epsilon_i; \phi)}{\partial \tilde{\mathbf{y}}_i} \frac{\partial \tilde{\mathbf{y}}(\epsilon_i; \phi)}{\partial \phi} \right)}_{\text{score function term}}] \quad (9)$$

$$\mathbb{E}_{q(\tilde{\mathbf{y}}_i|\mathbf{x})} \left[ \frac{w_i}{\sum_{j=1}^k w_j} \frac{\partial \log q_{\phi}(\tilde{\mathbf{y}}_i|\mathbf{x})}{\partial \phi} \right] = \mathbb{E}_{p(\epsilon_i)} \left[ \frac{\partial}{\partial \tilde{\mathbf{y}}_i} \left( \frac{w_i}{\sum_{j=1}^k w_j} \right) \frac{\partial \tilde{\mathbf{y}}(\epsilon_i; \phi)}{\partial \phi_i} \right] \quad (10)$$

As we show empirically in Tab. 1, blindly adopting IWAE-DReG estimator for multiple-sample NIC brings evident performance decay. Other than IWAE-DReG, many other gradient estimators such as NVIL [Mnih and Gregor, 2014], VIMCO [Mnih and Rezende, 2016] and GDReG [Bauer and Mnih, 2021] do not apply to NIC. They either bring some extra bias or are totally wrong.

### 2.3 Impact on Bits-Back Coding

It is well known that the ELBO  $\mathcal{L}$  is the minus overall bitrate for bits-back coding in compression [Hinton and Van Camp, 1993, Hinton et al., 1995, Chen et al., 2017], and the entropy of variational posterior is exactly the bits-back rate itself. For this reason, earlier works [Townsend et al., 2018, Yang et al., 2020] point out that [Ballé et al., 2018, Minnen et al., 2018] waste bits for not using bits-back coding on  $z$ . However, during training the differential entropy [Cover, 1999]  $\mathbb{E}_{q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})} [\log q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})]$  is constant. And this means that this term does not have impact on the optimization procedure. And due to the deterministic inference, the  $\log q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$  is 0, which means that the bitrate saved by bits-back coding is 0. In this sense, [Ballé et al., 2018, Minnen et al., 2018] is also optimal in bits-back coding perspective, although no actual bits-back coding is performed. In fact, there is no space for bits-back coding so long as encoder is deterministic. Since we can view deterministic encoder as a posterior distribution with mass 1 on a single point. And then the posterior’s entropy is always 0.

### 2.4 The *direct-y* Trick in Training NIC

In NIC, we feed deterministic parameter  $\mathbf{y}$  into  $z$  inference model instead of noisy samples  $\tilde{\mathbf{y}}$ . This implies that  $\tilde{\mathbf{z}}$  is sampled from  $q(\tilde{\mathbf{z}}|\mathbf{y})$  instead of  $q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$ . This trick is initially adopted in Ballé et al. [2018] and followed by most of the subsequent works. However, it is little discussed. In this

paper, we refer it to *direct-y* trick. Yang et al. [2020] observes that feeding  $\tilde{\mathbf{y}}$  instead of  $\mathbf{y}$  causes severe performance decay. We confirm this result in Tab. 2. Thus, *direct-y* trick is essential to train hierarchical NIC.

Table 2: Effects of *direct-y* on R-D performance. 2-level VAE is equivalent to Ballé et al. [2018] without *direct-y*.

	bpp	MSE	PSNR	R-D cost
2-level VAE	0.9968	33.08	33.22	1.493
[Ballé et al., 2018]	0.5273	32.61	33.28	1.017

Table 3: Effects of *direct-y* on gradient SNR of different parts of the model. 2-level VAE is equivalent to Ballé et al. [2018] without *direct-y*. "early" is  $5 \times 10^4$  iterations, "mid" is  $5 \times 10^5$  iterations and "late" is  $1 \times 10^6$  iterations. "infer" is the abbreviation for "inference model", and "gen" is the abbreviation for "generative model".

Iteration	Method	gradient SNR of #				
		y infer	y gen	z infer	z gen	z prior
early	2-level VAE	2.287	0.5343	0.3419	0.4099	0.9991
	Ballé et al. [2018]	2.174	0.5179	0.5341	0.3813	1.069
mid	2-level VAE	1.350	0.4793	0.2414	0.3583	0.8861
	Ballé et al. [2018]	1.334	0.4813	0.4879	0.3761	0.9693
late	2-level VAE	1.217	0.4746	0.2863	0.3439	0.8691
	Ballé et al. [2018]	1.206	0.4763	0.5506	0.3707	0.9339

One explanation is to view  $q(\tilde{\mathbf{z}}|\mathbf{y})$  as  $q(\tilde{\mathbf{z}}|\mathbf{x})$ , and  $q(\tilde{\mathbf{y}}, \tilde{\mathbf{z}}|\tilde{\mathbf{x}})$  factorized as  $q(\tilde{\mathbf{y}}|\mathbf{x})q(\tilde{\mathbf{z}}|\mathbf{x})$  (See Fig.1 (a)-(c)). A similar trick of feeding mean parameter can be traced back to the Helmholtz machine [Dayan et al., 1995]. However, this provides a rationale why *direct-y* is fine to be adopted but does not explain why sampling  $\tilde{\mathbf{z}}$  from  $q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$  fails. We provide an alternative explanation from the gradient variance perspective. Specifically,  $q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$  has two stochastic arguments that could cause high variance in the gradient of z inference model, and make its convergence difficult. To verify this, we follow Rainforth et al. [2018] to compare the gradient SNR, which is the absolute value of the empirical mean divided by standard deviation. We trace the gradient SNR during different training stages as model converges (See Sec. 5.1 for detailed setups).

As demonstrated in Tab. 3, the gradient SNR of z inference model of standard 2-level VAE (without *direct y*) is indeed significantly lower than Ballé et al. [2018] (with *direct y*) during all 3 stage of training. This result reveals that the z inference model is more difficult to train without *direct-y*. And such difficulty could be the source of the failure of NIC without *direct-y* trick.

### 3 Multiple-sample Neural Image Compression

In this section, we consider the multiple-sample approach based on the 2-level hierarchical framework by Ballé et al. [2018], which is the de facto NIC architecture adopted by many sota methods. To simplify notations,  $\log q(\tilde{\mathbf{z}}|\tilde{\mathbf{y}})$  and  $\log q(\tilde{\mathbf{y}}|\mathbf{x})$  in ELBO are omitted as they are 0. First, let's consider directly applying 2-level IWAE to NIC without *direct-y* trick (See Fig. 1 (d)). Regarding a  $k$  sample IWAE, we first compute parameter  $\mathbf{y}$  of  $q(\tilde{\mathbf{y}}|\mathbf{x})$  and sample  $\tilde{\mathbf{y}}_{1:k}$  from it. Then, we compute parameter  $\mathbf{z}_{1:k}$  of  $q(\tilde{\mathbf{z}}_{1:k}|\tilde{\mathbf{y}}_{1:k})$  and samples  $\tilde{\mathbf{z}}_{1:k}$  from it. Afterward,  $\tilde{\mathbf{y}}_{1:k}$  and  $\tilde{\mathbf{z}}_{1:k}$  are fed into the generative model and compute  $w_{1:k}$ . Finally, we follow Eq 5 to compute the gradient and update parameters. In fact, this is the standard 2-level IWAE in the original IWAE paper.

However, the vanilla 2-level IWAE becomes a problem for NIC with *direct-y* trick. Concerning a  $k$  sample IWAE, we sample  $\tilde{\mathbf{y}}_{1:k}$  from  $q(\tilde{\mathbf{y}}|\mathbf{x})$ . Due to the *direct-y* trick, we feed  $\mathbf{y}$  instead of  $\tilde{\mathbf{y}}_{1:k}$  into z inference network, and our  $q(\tilde{\mathbf{z}}|\mathbf{y})$  has only one parameter  $\mathbf{z}$  other than  $k$  parameter  $\mathbf{z}_{1:k}$ . If we follow the 2-level IWAE approach, only one sample  $\tilde{\mathbf{z}}$  is obtained, and  $w_{1:k}$  can not be computed.

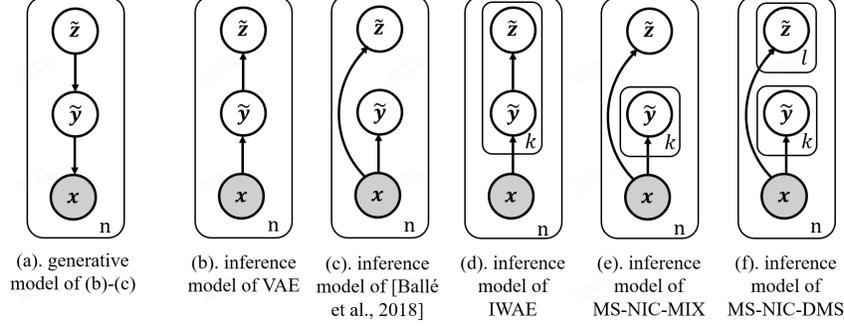


Figure 1: The plate notation of different NIC methods.  $x$  is the observed image,  $\tilde{y}$  and  $\tilde{z}$  are latent. The inference models show how we sample from variational posterior during training.  $n$  is the number of data points in dataset,  $k, l$  is the sample size of multiple-sample approaches. The generative model of (b), (c) is (a). The generative model of (d)-(f) is shown in Appendix. A.1. For clarity, we omit the parameters.

One method is to limit the multiple-sample part to  $\tilde{y}$  related term only and optimize other parts via single-sample SGVB-1, which produces our MS-NIC-MIX (See Fig 1 (e)). Another method is to sample another  $l$  samples of  $\tilde{z}_j$  from  $q(\tilde{z}|\tilde{y})$  and nest it with MS-NIC-MIX, which generates our MS-NIC-DMS (See Fig 1 (f)).

### 3.1 MS-NIC-MIX: Multiple-sample NIC with Mixture

One way to optimize multiple-sample IWAE target of NIC with *direct-y* trick is to sample  $\tilde{y}$   $k$  times to obtain  $\tilde{y}_{1:k}$  and  $\tilde{z}$  only 1 time. Then we perform  $k$  sample log mean of  $p(x|\tilde{y}_i)p(\tilde{y}_i|\tilde{z})$  to obtain a multiple-sample estimated  $\log p(x|\tilde{z})$ , add it with single-sample  $\log p(\tilde{z})$ . This brings a  $\mathcal{L}_k^{MIX}$  with the form of a mixture of 1-level VAE and 1-level IWAE ELBO:

$$\mathcal{L}_k^{MIX} = \mathbb{E}_{q_\phi(\tilde{z}|x)} [\mathbb{E}_{q_\phi(\tilde{y}_{1:k}|x)} [\log \frac{1}{k} \sum_i^k p(x|\tilde{y}_i)p(\tilde{y}_i|\tilde{z}) | \tilde{z}] + \log p(\tilde{z})] \quad (11)$$

Moreover,  $\mathcal{L}_k^{MIX}$  is a reasonably preferable target over ELBO as it satisfies the following properties (See Appendix. A.2 for proof):

1.  $\log p(x) \geq \mathcal{L}_k^{MIX}$
2.  $\mathcal{L}_k^{MIX} \geq \mathcal{L}_m^{MIX}$  for  $k \geq m$

Although  $\mathcal{L}_k^{MIX}$  does not converge to true  $\log p(x)$  as  $k$  grows, it is still a lower bound of  $\log p(x)$  and tighter than ELBO (as  $\mathcal{L}_1^{MIX} = \text{ELBO}$ ). Its gradient can be computed via pathwise estimator. Denote the per-sample integrand  $p(x|\tilde{y}_i)p(\tilde{y}_i|\tilde{z})$  as  $w_i^{MIX}$ , and its relative weight as  $\tilde{w}_i^{MIX}$ , then the gradient  $\nabla_\phi \mathcal{L}_k^{MIX}$  can be estimated as Eq. 13.

$$\begin{aligned} \mathcal{L}_k^{MIX} &= \mathbb{E}_{p(\epsilon_{1:k}^y, \epsilon^z)} [\log \frac{1}{k} \sum_i^k p(x|\tilde{y}(\epsilon_i^y; \phi))p(\tilde{y}(\epsilon_i^y; \phi)|\tilde{z}(\epsilon^z; \phi)) + \log p(\tilde{z}(\epsilon^z; \phi))] \\ &\approx \log \frac{1}{k} \sum_i^k \underbrace{p(x|\tilde{y}(\epsilon_i^y; \phi))p(\tilde{y}(\epsilon_i^y; \phi)|\tilde{z}(\epsilon^z; \phi))}_{w^{MIX}(\epsilon_{1:k}^y, \epsilon^z; \phi)} + \log p(\tilde{z}(\epsilon^z; \phi)) \end{aligned} \quad (12)$$

$$\nabla_\phi \mathcal{L}_k^{MIX} \approx \sum_i^k \tilde{w}_i^{MIX} \nabla_\phi \log w^{MIX}(\epsilon_{1:k}^y, \epsilon^z; \phi) + \nabla_\phi \log p(\tilde{z}(\epsilon^z; \phi)) \quad (13)$$

Another way to understand MS-NIC-MIX is to view the  $y$  inference/generative model as a single level IWAE, and the  $z$  inference/generative model as a large prior of  $\tilde{y}$  which is optimized via SGVB-1.

This perspective is often taken by works in NIC context model [Minnen et al., 2018, He et al., 2021], as the context model of NIC is often limited to  $\tilde{\mathbf{y}}$ .

### 3.2 MS-NIC-DMS: Multiple-sample NIC with Double Multiple Sampling

An intuitive improvement over MS-NIC-MIX is to add another round of multiple-sample over  $\tilde{\mathbf{z}}$ . Specifically, we sample  $\tilde{\mathbf{z}}$   $l$  times, nest it with  $\mathcal{L}_k^{MIX}$  to obtain  $\mathcal{L}_{k,l}^{DMS}$ :

$$\mathcal{L}_{k,l}^{DMS} = \mathbb{E}_{q_\phi(\tilde{\mathbf{z}}_{1:l}|\mathbf{x})} \left[ \log \frac{1}{l} \sum_j \exp \left( \mathbb{E}_{q_\phi(\tilde{\mathbf{y}}_{1:k}|\mathbf{x})} \left[ \log \frac{1}{k} \sum_i p(\mathbf{x}|\tilde{\mathbf{y}}_i) p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{z}}_j) \right] p(\tilde{\mathbf{z}}_j) \right) \right] \quad (14)$$

And we name it MS-NIC-DMS as it adopts multiple sampling twice. Moreover,  $\mathcal{L}_{k,l}^{DMS}$  is a reasonably better target for optimization over ELBO and  $\mathcal{L}_k^{MIX}$ , as it satisfies the following properties (See proof in Appendix. A.2):

1.  $\log p(\mathbf{x}) \geq \mathcal{L}_{k,l}^{DMS}$
2.  $\mathcal{L}_{k,l}^{DMS} \geq \mathcal{L}_{m,n}^{DMS}$  for  $k \geq m, l \geq n$
3.  $\mathcal{L}_{k,l}^{DMS} \geq \mathcal{L}_k^{MIX}$
4.  $\mathcal{L}_{k,l}^{DMS} \rightarrow \log p(\mathbf{x})$  as  $k, l \rightarrow \infty$ , under the assumption that  $\log \frac{p(\mathbf{x}|\tilde{\mathbf{y}}_i)p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{z}}_j)}{q(\tilde{\mathbf{y}}_i|\mathbf{x})}$  and  $\log \frac{p(\mathbf{x}|\tilde{\mathbf{z}}_j)p(\tilde{\mathbf{z}}_j)}{q(\tilde{\mathbf{z}}_j|\mathbf{x})}$  are bounded.

In other words, the target  $\mathcal{L}_{k,l}^{DMS}$  is a lowerbound of  $\log p(\mathbf{x})$ , converging to  $\log p(\mathbf{x})$  as  $k, l \rightarrow \infty$ , tighter than  $\mathcal{L}_k^{MIX}$  and tighter than ELBO (as  $\mathcal{L}_{1,1}^{DMS} = \text{ELBO}$ ). However, its Monte Carlo estimation is biased due to the nested transformation and expectation. Empirically, we find that directly adopting biased pathwise estimator works fine. And its gradient can be estimated by pathwise estimator similar to original IWAE target (See Eq. 5).

$$\begin{aligned} \mathcal{L}_{k,l}^{DMS} &= \mathbb{E}_{q(\epsilon_{1:l}^z)} \left[ \log \frac{1}{l} \sum_j \exp \left( \mathbb{E}_{q(\epsilon_{1:k}^y)} \left[ \log \frac{1}{k} \sum_i p(\mathbf{x}|\tilde{\mathbf{y}}(\epsilon_i^y; \phi)) p(\tilde{\mathbf{y}}(\epsilon_i^y; \phi)|\tilde{\mathbf{z}}(\epsilon_j^z; \phi)) \right] p(\tilde{\mathbf{z}}(\epsilon_j^z; \phi)) \right) \right] \\ &\approx \log \frac{1}{l} \sum_j \frac{1}{k} \sum_i \underbrace{p(\mathbf{x}|\tilde{\mathbf{y}}(\epsilon_i^y; \phi)) p(\tilde{\mathbf{y}}(\epsilon_i^y; \phi)|\tilde{\mathbf{z}}(\epsilon_j^z; \phi)) p(\tilde{\mathbf{z}}(\epsilon_j^z; \phi))}_{w^{DMS}(\epsilon_{1:k}^y, \epsilon_{1:l}^z; \phi)} \end{aligned} \quad (15)$$

Another interpretation of MS-NIC-DMS is to view it as a multiple level IWAE with repeated local samples. The  $\mathcal{L}_{k,l}^{DMS}$  Monte Carlo pathwise estimator has the form of IWAE with  $k \times l$  samples. However, there are multiple repeated samples that contain the same  $\tilde{\mathbf{y}}_i$  and  $\tilde{\mathbf{z}}_j$ . For example, the samples  $w_{1:6}^{IWAE}$  of 2 level IWAE with sample size 6 look like Eq. 16. While the samples  $w_{1:2,1:3}^{DMS}$  of MS-NIC-DMS with  $2 \times 3$  samples look like Eq. 17. We can see that in IWAE, we have 6 pairs of independently sampled  $\tilde{\mathbf{y}}$  and  $\tilde{\mathbf{z}}$ , while in MS-NIC-DMS, we have 2 independent  $\tilde{\mathbf{y}}$  and 3 independent  $\tilde{\mathbf{z}}$ , they are paired to generate 6 samples in total. Note that this is only applicable to NIC as  $\tilde{\mathbf{y}}$  and  $\tilde{\mathbf{z}}$  are conditionally independent given  $\tilde{\mathbf{x}}$  due to *direct-y* trick.

$$\begin{aligned} w_{1:6}^{IWAE} &= \{ p(\mathbf{x}|\tilde{\mathbf{y}}_1)p(\tilde{\mathbf{y}}_1|\tilde{\mathbf{z}}_1)p(\tilde{\mathbf{z}}_1), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_2)p(\tilde{\mathbf{y}}_2|\tilde{\mathbf{z}}_2)p(\tilde{\mathbf{z}}_2), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_3)p(\tilde{\mathbf{y}}_3|\tilde{\mathbf{z}}_3)p(\tilde{\mathbf{z}}_3), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_4)p(\tilde{\mathbf{y}}_4|\tilde{\mathbf{z}}_4)p(\tilde{\mathbf{z}}_4), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_5)p(\tilde{\mathbf{y}}_5|\tilde{\mathbf{z}}_5)p(\tilde{\mathbf{z}}_5), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_6)p(\tilde{\mathbf{y}}_6|\tilde{\mathbf{z}}_6)p(\tilde{\mathbf{z}}_6) \} \end{aligned} \quad (16)$$

$$\begin{aligned} w_{1:2,1:3}^{DMS} &= \{ p(\mathbf{x}|\tilde{\mathbf{y}}_1)p(\tilde{\mathbf{y}}_1|\tilde{\mathbf{z}}_1)p(\tilde{\mathbf{z}}_1), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_1)p(\tilde{\mathbf{y}}_1|\tilde{\mathbf{z}}_2)p(\tilde{\mathbf{z}}_2), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_1)p(\tilde{\mathbf{y}}_1|\tilde{\mathbf{z}}_3)p(\tilde{\mathbf{z}}_3), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_2)p(\tilde{\mathbf{y}}_2|\tilde{\mathbf{z}}_1)p(\tilde{\mathbf{z}}_1), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_2)p(\tilde{\mathbf{y}}_2|\tilde{\mathbf{z}}_2)p(\tilde{\mathbf{z}}_2), \\ &\quad p(\mathbf{x}|\tilde{\mathbf{y}}_2)p(\tilde{\mathbf{y}}_2|\tilde{\mathbf{z}}_3)p(\tilde{\mathbf{z}}_3) \} \end{aligned} \quad (17)$$

## 4 Related Work

### 4.1 Lossy Neural Image and Video Compression

Ballé et al. [2017] and Ballé et al. [2018] formulate lossy neural image compression as a variational inference problem, by interpreting the additive uniform noise (AUN) relaxed scalar quantization as a factorized uniform variational posterior. After that, the majority of sota lossy neural image compression methods adopt this formulation [Minnen et al., 2018, Minnen and Singh, 2020, Cheng et al., 2020, Guo et al., 2021a, Gao et al., 2021, He et al., 2022]. And Yang et al. [2020], Guo et al. [2021b] also require a AUN trained NIC as base. Moreover, the majority of neural video compression also adopts this formulation [Lu et al., 2019, 2020, Agustsson et al., 2020, Hu et al., 2021, Li et al., 2021], implying that MS-NIC can be extended to video compression without much pain.

Other approaches to train NIC include random rounding [Toderici et al., 2015, 2017] and straight through estimator (STE) [Theis et al., 2017]. Another promising approach is the VQ-VAE [Van Den Oord et al., 2017]. By the submission of this manuscript, one unarchived work [Zhu et al., 2022] has shown the potential of VQ-VAE in practical NIC. Our MS-NIC does not apply to the approaches mentioned in this paragraph, as the formulation of variational posterior is different.

### 4.2 Tighter Lower Bound for VAE

IWAE [Burda et al., 2016] stirs up the discussion of adopting tighter lower bound for training VAEs. However, at the first glance it is not straightforward why it might works. Cremer et al. [2018] decomposes the inference suboptimality of VAE into two parts: 1) The limited expressiveness of inference model. 2) The gap between ELBO and log likelihood. However, this gap refers to inference not training. The original IWAE paper empirically shows that IWAE can learn a richer latent representation. And Cremer et al. [2017] shows that the IWAE target converges to ELBO under the expectation of true posterior. And thus the posterior collapse is avoided.

From the information preference [Chen et al., 2017] perspective, VAE prefers to distribute information in generative distribution than autoencoding information in the latent. This preference formulates another view of posterior collapse. And it stems from the gap between ELBO and true log likelihood. There are various approaches alleviating it, including *soft free bits* [Theis et al., 2017] and *KL annealing* [Serban et al., 2017]. In our opinion, IWAE also belongs to those methods, and it is asymptotically optimal. However, we have not found many works comparing IWAE with those methods. Moreover, those approaches are rarely adopted in NIC community.

Many follow-ups of IWAE stress gradient variance reduction [Roeder et al., 2017, Tucker et al., 2018, Rainforth et al., 2018], discrete latent [Mnih and Rezende, 2016] and debiasing IWAE target [Nowozin, 2018]. Although the idea of tighter low bound training has been applied to the field of neural joint source channel coding [Choi et al., 2018, Song et al., 2020], to the best of our knowledge, no work in NIC consider it yet.

### 4.3 Multi-Sample Inference for Neural Image Compression

Theis and Ho [2021] considers the similar topic of importance weighted NIC. However, it does not consider training of NIC. Instead, it focuses on achieving IWAE target with an entropy coding technique named *softmin*, just like BB-ANS [Townsend et al., 2018] achieving ELBO. It is alluring to apply *softmin* to MS-NIC, as it closes the multiple-sample training and inference gap. However, it requires large number of samples (e.g. 4096) to achieve slight improvement for  $64 \times 64$  images. The potential sample size required for practical NIC is forbidding. Moreover, we believe the stochastic lossy encoding scheme [Agustsson and Theis, 2020] that Theis and Ho [2021] is not yet ready to be applied (See Appendix. A.7 for details).

## 5 Experimental Results

### 5.1 Experimental Settings

Following He et al. [2022], we train all the models on the largest 8000 images of ImageNet [Deng et al., 2009], followed by a downsampling according to Ballé et al. [2018]. And we use Kodak [Kodak,

1993] for evaluation. For the experiments based on Ballé et al. [2018] (include Tab. 1, Tab. 2), we follows the setting of the original paper except for the selection of  $\lambda_s$ , For the selection of  $\lambda_s$ , we set  $\lambda \in \{0.0016, 0.0032, 0.0075, 0.015, 0.03, 0.045, 0.08\}$  as suggested in Cheng et al. [2020]. And for the experiments based on Cheng et al. [2020], we follows the setting of original paper. More detailed experimental settings can be found in Appendix. A.4.

And when comparing the R-D performance of models trained on multiple  $\lambda_s$ , we use Bjontegaard metric (BD-Metric) and Bjontegaard bitrate (BD-BR) [Bjontegaard, 2001], which is widely applied when comparing codecs. More detailed experimental results can be found in Appendix. A.5.

Table 4: Results based on Ballé et al. [2018].

	PSNR		MS-SSIM	
	BD-BR (%)	BD-Metric	BD-BR (%)	BD-Metric
<i>Single-sample</i>				
Baseline [Ballé et al., 2018]	0.000	0.000	0.000	0.0000
<i>Multiple-sample</i>				
IWAE [Burda et al., 2016]	64.23	-2.318	68.67	-0.01648
MS-NIC-MIX	-3.847	0.1877	-4.743	0.001618
MS-NIC-DMS	-4.929	0.2405	-5.617	0.001976

Table 5: Results based on Cheng et al. [2020]. The BD Metrics of IWAE can not be computed as its R-D is not monotonously increasing.

	PSNR		MS-SSIM	
	BD-BR (%)	BD-Metric	BD-BR (%)	BD-Metric
<i>Single-sample</i>				
Baseline [Cheng et al., 2020]	0.0000	0.0000	0.0000	0.0000
<i>Multiple-sample</i>				
IWAE [Burda et al., 2016]	-	-	-	-
MS-NIC-MIX	-1.852	0.0805	2.238	-0.0006764
MS-NIC-DMS	-2.378	0.1046	1.998	-0.0006054

## 5.2 R-D Performance

We evaluate the performance of MS-NIC-MIX and MS-NIC-DMS based on sota NIC methods [Ballé et al., 2018, Cheng et al., 2020]. Empirically, we find that MS-NIC-MIX works best with sample size 8, and MS-NIC-DMS with sample size 16. The experimental results on sample size selection can be found in Appendix. A.3. Without special mention, we set the sample size of MS-NIC-MIX to 8 and MS-NIC-DMS to 16.

For Ballé et al. [2018], MS-NIC-MIX saves around 4% of bitrate compared with single-sample baseline (See Tab. 4). And MS-NIC-DMS saves around 5% of bitrate. On the other hand, the original IWAE suffers performance decay as it is not compatible with *direct-y* trick. For Cheng et al. [2020], we find that both MS-NIC-MIX and MS-NIC-DMS suppress baseline in PSNR. However, it is not as evident as Ballé et al. [2018]. Moreover, the MS-SSIM is slightly lower than the baseline. This is probably due to the auto-regressive context model. Besides, the original IWAE without *direct-y* trick suffers from severe performance decay in both cases. The BD metric of IWAE on Cheng et al. [2020] can not be computed as its R-D is not monotonous increasing, we refer interested readers to Appendix. A.5 for details.

## 5.3 Latent Space Representation of MS-NIC

To better understand the latent learned by MS-NIC, we evaluate the variance and coefficient of variation (Cov) of per-dimension latent distribution mean parameter  $\mathbf{y}^{(i)}$ ,  $\mathbf{z}^{(i)}$ , with regard to input

distribution  $p(\mathbf{x})$ . As we are also interested in the discrete representation, we provide statistics of rounded mean  $\bar{\mathbf{y}}^{(i)}, \bar{\mathbf{z}}^{(i)}$ . These metrics show how much do latents vary when input changes, and a large variation in latents means that there are useful information encoded. A really small variation indicates that the latent is "dead" in that dimension.

As shown in Tab. 10 of Appendix. A.6, the latent of multiple-sample approaches has higher variance than those of single-sample approach. Moreover, the  $\text{Cov}(\mathbf{y})$  of multiple-sample approaches is around 4 – 5 times higher than single-sample approach. Although the  $\text{Cov}(\mathbf{z})$  of multiple-sample approaches is around 2 times lower, the main contributor of image reconstruction is  $\mathbf{y}$ , and  $\mathbf{z}$  only serves to predict  $\mathbf{y}$ 's distribution. Similar trend can be concluded from quantized latents  $\bar{\mathbf{y}}, \bar{\mathbf{z}}$ . From the variance and Cov perspective, the latent learned by MS-NIC is richer than single-sample approach. It is also noteworthy that although the variance and Cov of  $\mathbf{y}, \bar{\mathbf{y}}$  of MS-NIC is significantly higher than single-sample approach, the bpp only varies slightly.

Table 6: The average of per-dimension latent variance and Cov across Kodak test images. The model is trained with  $\lambda = 0.015$ .

Method	Var(#)		Cov(#)		bpp of #	
	$\bar{\mathbf{y}}$	$\bar{\mathbf{z}}$	$\bar{\mathbf{y}}$	$\bar{\mathbf{z}}$	$\bar{\mathbf{y}}$	$\bar{\mathbf{z}}$
<i>Single-sample</i>						
Ballé et al. [2018]	1.499	0.3255	19.70	9.944	0.5136	0.01342
<i>Multiple-sample</i>						
MS-NIC-MIX	1.906	0.7594	111.1	7.425	0.5108	0.01521
MS-NIC-DMS	1.919	0.7648	95.51	7.243	0.5092	0.01634

## 6 Limitation & Discussion

A major limitation of our method is that the improvement in R-D performance is marginal, especially when based on Cheng et al. [2020]. Moreover, evaluations on more recent sota methods are also helpful to strengthen the claims of this paper. In general, we think that the performance improvement of our approach is bounded by how severe the posterior collapse is in neural image compression. We measure the variance in latent dimension according to data in Fig. A.6. And from that figure it might be observed that the major divergence of IWAE and VAE happens when the variance is very small. And for the area where variance is reasonably large, the gain of IWAE is not that large. This probably indicates that the posterior collapse in neural image compression is only alleviated to a limited extent.

See more discussion in why the result on Cheng et al. [2020] is negative in Appendix. A.9

## 7 Conclusion

In this paper we propose MS-NIC, a multiple-sample importance weighted target for training NIC. It improves sota NIC methods and learns richer latent representation. A known limitation is that its R-D performance improvement is limited when applied to models with spatial context models (e.g. Cheng et al. [2020]). Despite the somewhat negative result, this paper provides insights to the training of NIC models from VAE perspective. Further work could consider improving the performance and extend it into neural video compression.

## Acknowledgments and Disclosure of Funding

This work is supported by SenseTime Research. The content is solely the responsibility of the authors and does not necessarily represent the official views of SenseTime Research.

## References

- E. Agustsson and L. Theis. Universally quantized neural compression. *Advances in neural information processing systems*, 33:12367–12376, 2020.
- E. Agustsson, D. Minnen, N. Johnston, J. Balle, S. J. Hwang, and G. Toderici. Scale-space flow for end-to-end optimized video compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8503–8512, 2020.
- J. Ballé, V. Laparra, and E. P. Simoncelli. End-to-end optimized image compression. In *International Conference on Learning Representations*, 2017.
- J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston. Variational image compression with a scale hyperprior. In *International Conference on Learning Representations*, 2018.
- M. Bauer and A. Mnih. Generalized doubly reparameterized gradient estimators. In *International Conference on Machine Learning*, pages 738–747. PMLR, 2021.
- G. Bjontegaard. Calculation of average psnr differences between rd-curves. *VCEG-M33*, 2001.
- B. Bross, J. Chen, J.-R. Ohm, G. J. Sullivan, and Y.-K. Wang. Developments in international video coding standardization after avc, with an overview of versatile video coding (vvc). *Proceedings of the IEEE*, 109(9):1463–1493, 2021.
- Y. Burda, R. B. Grosse, and R. Salakhutdinov. Importance weighted autoencoders. In *ICLR (Poster)*, 2016.
- X. Chen, D. P. Kingma, T. Salimans, Y. Duan, P. Dhariwal, J. Schulman, I. Sutskever, and P. Abbeel. Variational lossy autoencoder. In *International Conference on Learning Representations*, 2017.
- Z. Cheng, H. Sun, M. Takeuchi, and J. 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*, pages 7939–7948, 2020.
- K. Choi, K. Tatwawadi, T. Weissman, and S. Ermon. Necst: neural joint source-channel coding. 2018.
- T. M. Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- C. Cremer, Q. Morris, and D. Duvenaud. Reinterpreting importance-weighted autoencoders. 2017.
- C. Cremer, X. Li, and D. Duvenaud. Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning*, pages 1078–1086. PMLR, 2018.
- P. Dayan, G. E. Hinton, R. M. Neal, and R. S. Zemel. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995.
- J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- G. Flamich, M. Havasi, and J. M. Hernández-Lobato. Compressing images by encoding their latent representations with relative entropy coding. *Advances in Neural Information Processing Systems*, 33:16131–16141, 2020.
- G. Gao, P. You, R. Pan, S. Han, Y. Zhang, Y. Dai, and H. Lee. Neural image compression via attentional multi-scale back projection and frequency decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14677–14686, 2021.
- Z. Guo, Z. Zhang, R. Feng, and Z. Chen. Causal contextual prediction for learned image compression. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021a.
- Z. Guo, Z. Zhang, R. Feng, and Z. Chen. Soft then hard: Rethinking the quantization in neural image compression. In *International Conference on Machine Learning*, pages 3920–3929. PMLR, 2021b.

- D. He, Y. Zheng, B. Sun, Y. Wang, and H. Qin. Checkerboard context model for efficient learned image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14771–14780, 2021.
- D. He, Z. Yang, W. Peng, R. Ma, H. Qin, and Y. Wang. Elic: Efficient learned image compression with unevenly grouped space-channel contextual adaptive coding. *arXiv preprint arXiv:2203.10886*, 2022.
- G. E. Hinton and D. Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pages 5–13, 1993.
- G. E. Hinton, P. Dayan, B. J. Frey, and R. M. Neal. The "wake-sleep" algorithm for unsupervised neural networks. *Science*, 268(5214):1158–1161, 1995.
- Z. Hu, G. Lu, and D. Xu. Fvc: A new framework towards deep video compression in feature space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1502–1511, 2021.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- E. Kodak. Kodak lossless true color image suite. <http://r0k.us/graphics/kodak/>, 1993.
- J. Li, B. Li, and Y. Lu. Deep contextual video compression. *Advances in Neural Information Processing Systems*, 34, 2021.
- I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao. Dvc: An end-to-end deep video compression framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11006–11015, 2019.
- G. Lu, X. Zhang, W. Ouyang, L. Chen, Z. Gao, and D. Xu. An end-to-end learning framework for video compression. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3292–3308, 2020.
- D. Minnen and S. Singh. Channel-wise autoregressive entropy models for learned image compression. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 3339–3343. IEEE, 2020.
- D. Minnen, J. Ballé, and G. D. Toderici. Joint autoregressive and hierarchical priors for learned image compression. *Advances in neural information processing systems*, 31, 2018.
- A. Mnih and K. Gregor. Neural variational inference and learning in belief networks. In *International Conference on Machine Learning*, pages 1791–1799. PMLR, 2014.
- A. Mnih and D. Rezende. Variational inference for monte carlo objectives. In *International Conference on Machine Learning*, pages 2188–2196. PMLR, 2016.
- S. Mohamed, M. Rosca, M. Figurnov, and A. Mnih. Monte carlo gradient estimation in machine learning. *J. Mach. Learn. Res.*, 21(132):1–62, 2020.
- S. Nowozin. Debiasing evidence approximations: On importance-weighted autoencoders and jackknife variational inference. In *International conference on learning representations*, 2018.
- M. B. Paulus, C. J. Maddison, and A. Krause. Rao-blackwellizing the straight-through gumbel-softmax gradient estimator. *arXiv preprint arXiv:2010.04838*, 2020.
- T. Rainforth, A. Kosiorek, T. A. Le, C. Maddison, M. Igl, F. Wood, and Y. W. Teh. Tighter variational bounds are not necessarily better. In *International Conference on Machine Learning*, pages 4277–4285. PMLR, 2018.

- G. Roeder, Y. Wu, and D. K. Duvenaud. Sticking the landing: Simple, lower-variance gradient estimators for variational inference. *Advances in Neural Information Processing Systems*, 30, 2017.
- T. Ryder, C. Zhang, N. Kang, and S. Zhang. Split hierarchical variational compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 386–395, 2022.
- I. Serban, A. Sordoni, R. Lowe, L. Charlin, J. Pineau, A. Courville, and Y. Bengio. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Y. Song, M. Xu, L. Yu, H. Zhou, S. Shao, and Y. Yu. Infomax neural joint source-channel coding via adversarial bit flip. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5834–5841, 2020.
- L. Theis and E. Agustsson. On the advantages of stochastic encoders. *arXiv preprint arXiv:2102.09270*, 2021.
- L. Theis and J. Ho. Importance weighted compression. In *Neural Compression: From Information Theory to Applications—Workshop@ ICLR 2021*, 2021.
- L. Theis, W. Shi, A. Cunningham, and F. Huszár. Lossy image compression with compressive autoencoders. *arXiv preprint arXiv:1703.00395*, 2017.
- G. Toderici, S. M. O’Malley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, and R. Sukthankar. Variable rate image compression with recurrent neural networks. *arXiv preprint arXiv:1511.06085*, 2015.
- G. Toderici, D. Vincent, N. Johnston, S. Jin Hwang, D. Minnen, J. Shor, and M. Covell. Full resolution image compression with recurrent neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 5306–5314, 2017.
- J. Townsend, T. Bird, and D. Barber. Practical lossless compression with latent variables using bits back coding. In *International Conference on Learning Representations*, 2018.
- G. Tucker, D. Lawson, S. Gu, and C. J. Maddison. Doubly reparameterized gradient estimators for monte carlo objectives. 2018.
- A. Van Den Oord, O. Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Y. Xie, K. L. Cheng, and Q. Chen. Enhanced invertible encoding for learned image compression. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 162–170, 2021.
- Y. Yang, R. Bamler, and S. Mandt. Improving inference for neural image compression. *Advances in Neural Information Processing Systems*, 33:573–584, 2020.
- X. Zhu, J. Song, L. Gao, F. Zheng, and H. T. Shen. Unified multivariate gaussian mixture for efficient neural image compression. *arXiv preprint arXiv:2203.10897*, 2022.
- Y. Zhu, Y. Yang, and T. Cohen. Transformer-based transform coding. In *International Conference on Learning Representations*, 2021.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
  - (b) Did you describe the limitations of your work? [\[Yes\]](#)
  - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Appendix

- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
  - (b) Did you include complete proofs of all theoretical results? [Yes] In Appendix
- 3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [No]
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] It is really expensive to train a NIC model. The error bar is rarely reported out of toy size model and dataset.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes]
  - (b) Did you mention the license of the assets? [N/A]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
  
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
- 5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]