

# CORRECTING THE SUB-OPTIMAL BIT ALLOCATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In this paper, we investigate the problem of bit allocation in Neural Video Compression (NVC). First, we reveal that a recent bit allocation approach claimed to be optimal is, in fact, sub-optimal due to its implementation. Specifically, we find that its sub-optimality lies in the improper application of semi-amortized variational inference (SAVI) on latent with non-factorized variational posterior. Then, we show that the corrected version of SAVI on non-factorized latent requires recursively applying back-propagating through gradient ascent, based on which we derive the corrected optimal bit allocation algorithm. Due to the computational in-feasibility of the corrected bit allocation, we design an efficient approximation to make it [tractable](#). Empirical results show that our proposed correction significantly improves the incorrect bit allocation in terms of R-D performance and bitrate error, and outperforms all other bit allocation methods by a large margin. The source code is provided in the supplementary material.

## 1 INTRODUCTION

Recently, bit allocation for Neural Video Compression (NVC) has drawn growing attention thanks to its great potential in boosting compression performance. Due to the frame reference structure in video coding, it is sub-optimal to use the same R-D (Rate-Distortion) trade-off parameter  $\lambda$  for all frames. In bit allocation task, bitrate is allocated to different frames/regions to minimize R-D cost  $R + \lambda D$ , where  $R$  is total bitrate,  $D$  is total distortion, and  $\lambda$  is the Lagrangian multiplier controlling R-D trade-off. Li et al. (2022) are the pioneer of bit allocation for NVC, who improve the empirical R-D (Rate-Distortion) model from traditional video codec (Li et al., 2014; 2016) and solve the per-frame Lagrangian multiplier  $\lambda$ . Other concurrent works adopt simple heuristics for coarse bit allocation (Cetin et al., 2022; Hu et al., 2022).

Most recently, BAO (Bit Allocation using Optimization) (Xu et al., 2022) proposes to formulate bit allocation as semi-amortized variational inference (SAVI) (Kim et al., 2018; Marino et al., 2018) and solves it by gradient-based optimization. Specifically, it directly optimizes the variational posterior parameter to be quantized and encoded by gradient ascent, aiming at maximizing the minus overall R-D cost, which is also the evident lowerbound (ELBO). BAO does not rely on any empirical R-D model and thus outperforms previous work. Further, BAO shows its optimality by proving its equivalence to bit allocation with precise R-D model.

In this paper, we first show that BAO (Xu et al., 2022) is in fact, sub-optimal due to its implementation. Specifically, we find that it abuses SAVI (Kim et al., 2018; Marino et al., 2018) on latent with non-factorized variational posterior, which brings incorrect gradient signal during optimization. To solve this problem, we first extend SAVI to non-factorized latent by back-propagating through gradient ascent (Domke, 2012). Then based on that, we correct the sub-optimal bit allocation in BAO to produce true optimal bit allocation for NVC. Furthermore, we propose a computational feasible approximation to such correct but intractable bit allocation method. And we show that our approximation outperforms the incorrect bit allocation (BAO) in terms of R-D performance and bitrate error, and performs better than all other bit allocation methods.

To summarize, our contributions are as follows:

- We demonstrate that a previously claimed optimal bit allocation method is actually sub-optimal. We find that its sub-optimality comes from the improper application of SAVI to non-factorized latent.

- We present the correct way to conduct SAVI on non-factorized latent by recursively applying back-propagation through gradient ascent. Based on this, we derive the corrected optimal bit allocation algorithm for NVC.
- Furthermore, we propose a computational efficient approximation of the optimal bit allocation to make it feasible. Our proposed approach improves the R-D performance and bitrate error over the incorrect bit allocation (BAO), and outperforms all other bit allocation methods for NVC.

## 2 PRELIMINARIES

### 2.1 NEURAL VIDEO COMPRESSION

The input of NVC is a GoP (Group of Picture)  $\mathbf{x}_{1:T}$ , where  $\mathbf{x}_i \in R^{H \times W}$  is the  $i^{th}$  frame with  $H \times W$  pixels, and  $T$  is the number of frame inside the GoP. Most of the works in NVC follow a latent variable model with temporal autoregressive relationship (Yang et al., 2020a). Specifically, to encode  $\mathbf{x}_i$ , we first extract the motion latent  $\mathbf{w}_i = f_\phi^w(\mathbf{x}_i, \mathbf{x}'_i)$  from current frame  $\mathbf{x}_i$  and previous reconstructed frame  $\mathbf{x}'_{i-1}$ , where  $f_\phi^w(\cdot)$  is the motion encoder parameterized by  $\phi^1$ . Then, we encode the quantized latent  $\tilde{\mathbf{w}}_i = \lfloor \mathbf{w}_i \rfloor$  with the probability mass function (pmf) estimator  $P_\theta(\tilde{\mathbf{w}}_i | \tilde{\mathbf{w}}_{<i}, \tilde{\mathbf{y}}_{<i})$  parameterized by  $\theta$ , where  $\lfloor \cdot \rfloor$  is the rounding. Then, we obtain the residual latent  $\mathbf{y}_i = f_\phi^y(\mathbf{x}_i, \mathbf{x}'_i, \tilde{\mathbf{w}}_i)$ , where  $f_\phi^y(\cdot)$  is the residual encoder. Then, similar to how we treat  $\mathbf{w}_i$ , we encode the quantized latent  $\tilde{\mathbf{y}}_i = \lfloor \mathbf{y}_i \rfloor$  with pmf  $P_\theta(\tilde{\mathbf{y}}_i | \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{<i})$ . Finally, we obtain the reconstructed frame  $\mathbf{x}'_i = g_\theta^x(\mathbf{x}'_{i-1}, \tilde{\mathbf{w}}_i, \tilde{\mathbf{y}}_i)$ , where  $g_\theta^x(\cdot)$  is the decoder parameterized by  $\theta$ .

As only the motion latent  $\tilde{\mathbf{w}}_i$  and residual latent  $\tilde{\mathbf{y}}_i$  exist in the bitstream, the above process can be simplified as Eq. 1 and Eq. 2, where  $f_\phi(\cdot)$  is the generalized encoder and  $g_\theta(\cdot)$  is the generalized decoder. The target of NVC is to minimize the per-frame R-D cost  $R_i + \lambda_i D_i$  (Eq. 3), where  $R_i$  is the bitrate,  $D_i$  is the distortion and  $\lambda_i$  is the Lagrangian multiplier controlling R-D trade-off. The bitrate  $R_i$  and distortion  $D_i$  is computed as Eq. 2, where  $d(\cdot, \cdot)$  is the distortion metric. And  $\lambda_i D_i$  can be further interpreted as the data likelihood term  $-\log p_\theta(\mathbf{x}_i | \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{\leq i})$  so long as we treat  $\lambda_i D_i$  as the energy function of a Gibbs distribution (Minnen et al., 2018). Specifically, when  $d(\cdot, \cdot)$  is MSE, we can interpret  $\lambda_i D_i = -\log p_\theta(\mathbf{x}_i | \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{\leq i}) + const$ , where  $p_\theta(\mathbf{x}_i | \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{\leq i})$  is a Gaussian distribution  $\mathcal{N}(\hat{\mathbf{x}}_i, 1/2\lambda_i T)$ .

$$\mathbf{w}_i = f_\phi(\mathbf{x}_i, \tilde{\mathbf{w}}_{<i}, \tilde{\mathbf{y}}_{<i}), \mathbf{y}_i = f_\phi(\mathbf{x}_i, \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{<i}), \text{ where } \tilde{\mathbf{w}}_i = \lfloor \mathbf{w}_i \rfloor, \tilde{\mathbf{y}}_i = \lfloor \mathbf{y}_i \rfloor \quad (1)$$

$$R_i = \log P_\theta(\tilde{\mathbf{w}}_i, \tilde{\mathbf{y}}_i | \tilde{\mathbf{w}}_{<i}, \tilde{\mathbf{y}}_{<i}), D_i = d(\mathbf{x}_i, g_\theta(\tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{\leq i})) \quad (2)$$

$$\max -(R_i + \lambda_i D_i) \quad (3)$$

On the other hand, NVC is also closely related to Variational Autoencoder (VAE) (Kingma & Welling, 2013). As the rounding  $\lfloor \cdot \rfloor$  is not differentiable, Ballé et al. (2016); Theis et al. (2017) propose to relax it by additive uniform noise (AUN), and replace  $\tilde{\mathbf{w}}_i = \lfloor \mathbf{w}_i \rfloor$ ,  $\tilde{\mathbf{y}}_i = \lfloor \mathbf{y}_i \rfloor$  with  $\tilde{\mathbf{w}}_i = \mathbf{w}_i + \mathcal{U}(-0.5, 0.5)$ ,  $\tilde{\mathbf{y}}_i = \mathbf{y}_i + \mathcal{U}(-0.5, 0.5)$ . Under such formulation, the above encoding-decoding process becomes a VAE on graphic model  $\tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{\leq i} \rightarrow \mathbf{x}_i$  with variational posterior as Eq. 4, where  $\mathbf{w}_i, \mathbf{y}_i$  plays the role of variational posterior parameter. Then, minimizing the overall R-D cost (Eq. 3) is equivalent to maximizing the evident lowerbound (ELBO) (Eq. 5).

$$q_\phi(\tilde{\mathbf{w}}_i | \mathbf{x}_i, \tilde{\mathbf{w}}_{<i}, \tilde{\mathbf{y}}_{<i}) = \mathcal{U}(\mathbf{w}_i - 0.5, \mathbf{w}_i + 0.5), q_\phi(\tilde{\mathbf{y}}_i | \mathbf{x}_i, \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{<i}) = \mathcal{U}(\mathbf{y}_i - 0.5, \mathbf{y}_i + 0.5) \quad (4)$$

$$-(R_i + \lambda_i D_i) = \mathbb{E}_{q_\phi} \left[ \underbrace{\log P_\theta(\tilde{\mathbf{w}}_i, \tilde{\mathbf{y}}_i | \tilde{\mathbf{w}}_{<i}, \tilde{\mathbf{y}}_{<i})}_{-R_i} + \underbrace{\log p_\theta(\mathbf{x}_i | \tilde{\mathbf{w}}_{\leq i}, \tilde{\mathbf{y}}_{\leq i})}_{-\lambda_i D_i} \right] \quad (5)$$

bits-back bitrate: 0

### 2.2 BIT ALLOCATION FOR NEURAL VIDEO COMPRESSION

It is well known to video coding community that using the same R-D trade-off parameter  $\lambda_i$  to optimize R-D cost in Eq. 3 for all  $T$  frames inside a GoP is suboptimal (Li et al., 2014; 2016). This sub-optimality comes from the frame reference structure and is explained in detail by Li et al. (2022); Xu et al. (2022). The target of bit allocation is to maximize the minus of overall R-D cost

<sup>1</sup>Following previous works in deep generative modeling (Kingma & Welling, 2013; Kim et al., 2018), we denote all parameters related to encoder as  $\phi$ , and all parameters related to decoder and prior as  $\theta$ .

(ELBO)  $\mathcal{L}$  as Eq. 6 given the overall R-D trade-off parameter  $\lambda_0$ , instead of maximizing  $\mathcal{L}_i$  of each frame  $i$  separately.

The pioneer work of bit allocation in NVC (Li et al., 2022) follows bit allocation for traditional video codec (Li et al., 2016). Specifically, it adopts empirical models to approximate the relationship of the rate dependency  $\partial R_{i+1}/\partial R_i$  and distortion dependency  $\partial D_{i+1}/\partial D_i$  between frames. Then it takes those models into Eq. 6 to solve  $\lambda_{1:T}^*$  explicitly as Eq. 7.left. However, its performance heavily relies on the accuracy of empirical models.

$$\max \mathcal{L} = \sum_{i=1}^T \mathcal{L}_i, \text{ where } \mathcal{L}_i = -(R_i + \lambda_0 D_i) \quad (6)$$

$$\lambda_{1:T}^* \leftarrow \arg \max_{\lambda_{1:T}} \mathcal{L}(\lambda_{1:T}), \text{ versus } \mathbf{w}_{1:T}^*, \mathbf{y}_{1:T}^* \leftarrow \arg \max_{\mathbf{w}_{1:T}, \mathbf{y}_{1:T}} \mathcal{L}(\mathbf{w}_{1:T}, \mathbf{y}_{1:T}) \quad (7)$$

On the other hand, BAO (Xu et al., 2022) does not solve  $\lambda_{1:T}^*$  explicitly. Instead, it adopts SAVI (Kim et al., 2018; Marino et al., 2018) to achieve implicit bit allocation. To be specific, it initializes the variational posterior parameter  $\mathbf{w}_{1:T}^0, \mathbf{y}_{1:T}^0$  from fully amortized variational inference (FAVI) as Eq. 1. Then, it optimizes  $\mathbf{w}_{1:T}, \mathbf{y}_{1:T}$  via gradient ascent to maximize  $\mathcal{L}$  as Eq. 7.right. During this procedure, no empirical model is required. BAO further proves that optimizing Eq. 7.right is equivalent to optimizing Eq. 7.left with precise rate and distortion dependency model  $\partial R_{i+1}/\partial R_i, \partial D_{i+1}/\partial D_i$  (See Thm. 1, Thm. 2 in Xu et al. (2022)). Thus, BAO claims that it is optimal assuming gradient ascent achieves global maximum. However, in next section, we show that BAO (Xu et al., 2022) is in fact suboptimal due to its implementation.

### 3 WHY BAO IS SUP-OPTIMAL

BAO (Xu et al., 2022) achieves the SAVI (Kim et al., 2018; Marino et al., 2018) target in Eq. 7.right by gradient-based optimization. More specifically, its update rule is described as Eq. 8 and Eq. 9, where  $K$  is the total number of gradient ascent steps, and  $\mathbf{w}_i^k, \mathbf{y}_i^k$  is the posterior parameter  $\mathbf{w}_i, \mathbf{y}_i$  after  $k$  steps of gradient ascent. In the original paper of BAO, the authors also find that directly optimizing  $\mathbf{w}_i, \mathbf{y}_i$  simultaneously by Eq. 8 and Eq. 9 performs worse than optimizing  $\mathbf{y}_i$  alone using Eq. 9, but they have not offered any explanation. It is obvious that optimizing  $\mathbf{y}_i$  alone is sub-optimal. However, it is not obvious why jointly optimizing  $\mathbf{w}_i, \mathbf{y}_i$  with Eq. 8 and Eq. 9 fails.

$$\mathbf{w}_i^{k+1} \leftarrow \mathbf{w}_i^k + \alpha \frac{d\mathcal{L}(\mathbf{w}_{1:T}^k, \mathbf{y}_{1:T}^k)}{d\mathbf{w}_i^k}, \text{ where } \frac{d\mathcal{L}(\mathbf{w}_{1:T}^k, \mathbf{y}_{1:T}^k)}{d\mathbf{w}_i^k} = \sum_{j=i}^T \frac{\partial \mathcal{L}_j(\mathbf{w}_{1:j}^k, \mathbf{y}_{1:j}^k)}{\partial \mathbf{w}_i^k} \quad (8)$$

$$\mathbf{y}_i^{k+1} \leftarrow \mathbf{y}_i^k + \alpha \frac{d\mathcal{L}(\mathbf{w}_{1:T}^k, \mathbf{y}_{1:T}^k)}{d\mathbf{y}_i^k}, \text{ where } \frac{d\mathcal{L}(\mathbf{w}_{1:T}^k, \mathbf{y}_{1:T}^k)}{d\mathbf{y}_i^k} = \sum_{j=i}^T \frac{\partial \mathcal{L}_j(\mathbf{w}_{1:j}^k, \mathbf{y}_{1:j}^k)}{\partial \mathbf{y}_i^k} \quad (9)$$

In fact, the update rule in Eq. 8 and Eq. 9 is exactly the SAVI (Kim et al., 2018; Marino et al., 2018) when  $\mathbf{w}_i, \mathbf{y}_i$  fully factorizes (e.g. the full factorization used in mean-field (Blei et al., 2017)). However, in NVC the  $\mathbf{w}_i, \mathbf{y}_i$  has complicated auto-regressive relationships (See Eq. 1 and Fig. 1.(a)). Abusing SAVI on non-factorized latent causes gradient error in two aspects: (1). The total derivative  $d\mathcal{L}/d\mathbf{w}_i, d\mathcal{L}/d\mathbf{y}_i$  is incomplete. (2). The total derivative  $d\mathcal{L}/d\mathbf{w}_i, d\mathcal{L}/d\mathbf{y}_i$  and partial derivative  $\partial \mathcal{L}_j/\partial \mathbf{w}_i, \partial \mathcal{L}_j/\partial \mathbf{y}_i$  is evaluated at wrong value. In next two sections, we elaborate those two issues with  $\mathbf{w}_i$  related equations in main text and  $\mathbf{y}_i$  related equations in Appendix. A.2.

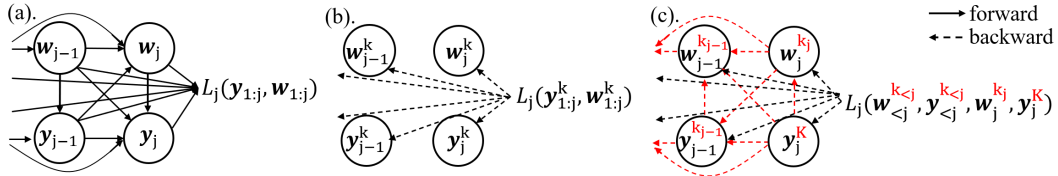


Figure 1: (a). The gradient structure of NVC without SAVI. (b). After  $k$  step of SAVI/gradient ascent, the gradient structure of NVC is broken. (c). The proposed approach using back-propagating through gradient ascent. We mark the difference between (b) and (c) in red.

### 3.1 INCOMPLETE TOTAL DERIVATIVE EVALUATION

According to the latent generation procedure described by Eq. 1 and Eq. 2, we draw the computational graph to describe the latent dependency as Fig. 1.(a). Based on that, we expand the total derivative  $d\mathcal{L}/dw_i, d\mathcal{L}/dy_i$  as Eq. 10 and Eq. 22.

$$\begin{aligned} \frac{d\mathcal{L}(\mathbf{w}_{1:T}, \mathbf{y}_{1:T})}{dw_i} &= \sum_{j=i}^T \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{dw_i} \\ \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{dw_i} &= \underbrace{\sum_{l=i+1}^j \frac{\partial \mathbf{w}_l}{\partial w_i} \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{d\mathbf{w}_l}}_{\text{ignored by BAO}} + \underbrace{\sum_{l=i}^j \frac{\partial \mathbf{y}_l}{\partial w_i} \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{d\mathbf{y}_l}}_{\text{considered by BAO}} + \underbrace{\frac{\partial \mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{\partial w_i}}_{\text{considered by BAO}} \end{aligned} \quad (10)$$

As shown in Eq. 8, Eq. 9 and Fig. 1.(b), BAO (Xu et al., 2022) treats the total derivative  $d\mathcal{L}/dw_i, d\mathcal{L}/dy_i$  as the sum of the frame level partial derivative  $\partial \mathcal{L}_j/\partial w_i, \partial \mathcal{L}_j/\partial y_i$ , which is the direct contribution of frame  $i^{\text{th}}$  latent  $w_i, y_i$  to  $j^{\text{th}}$  frame's R-D cost  $\mathcal{L}_j$  (as marked in Eq. 10 and Eq. 22). This incomplete evaluation of gradient signal brings sub-optimality. Further, it is not possible to correct BAO by simply including other parts of gradient into consideration. As BAO jointly updates all the latent  $\mathbf{w}_{1:T}, \mathbf{y}_{1:T}$ , the relationship of Eq. 2 only holds for the initial latent parameters  $\mathbf{w}_{1:T}^0, \mathbf{y}_{1:T}^0$  produced by FAVI. And this important relationship is broken for parameters  $\mathbf{w}_{1:T}^k, \mathbf{y}_{1:T}^k$  after  $k \geq 1$  steps of update.

### 3.2 INCORRECT VALUE TO EVALUATE GRADIENT

As shown in Eq. 8 and Eq. 9, BAO (Xu et al., 2022) simultaneously updates all the posterior parameter  $\mathbf{w}_{1:T}, \mathbf{y}_{1:T}$  with gradient evaluated at the same gradient ascent step  $\mathbf{w}_{1:T}^k, \mathbf{y}_{1:T}^k$ . However, as we show later in Sec. 4.1 and Fig. 1.(c), this is sub-optimal as all the descendant latent  $w_{>i}, y_{>i}$  of  $w_i$  should already complete all  $K$  steps of gradient ascent before the gradient of  $w_i$  is evaluated. Moreover,  $w_{>i}, y_{>i}$  should be initialized by FAVI using precedents latent. Similar rule applies to  $y_i$ . Specifically, the correct value to evaluate the gradient is as Eq. 11 and Eq. 23, where  $w_i^{k_i}$  denotes the latent  $w_i$  after  $k_i$  steps of update, and  $y_i^{k'_i}$  denotes the latent  $y_i$  after  $k'_i$  steps of update.

$$\begin{aligned} \mathbf{w}_i^{k_i+1} &\leftarrow \mathbf{w}_i^{k_i} + \alpha \frac{d\mathcal{L}(\mathbf{w}_1^{k_1}, \dots, \mathbf{w}_i^{k_i}, \mathbf{w}_{>i}^K, \mathbf{y}_1^{k'_1}, \dots, \mathbf{y}_{i-1}^{k'_{i-1}}, \mathbf{y}_{\geq i}^K)}{d\mathbf{w}_i^{k_i}}, \\ &\text{where } \mathbf{w}_{>i}^0, \mathbf{y}_{\geq i}^0 = f(\mathbf{x}, \mathbf{w}_1^{k_1}, \dots, \mathbf{w}_i^{k_i}, \mathbf{y}_1^{k'_1}, \dots, \mathbf{y}_{i-1}^{k'_{i-1}}) \end{aligned} \quad (11)$$

Similar to the incomplete total derivative evaluation, this problem does not have a simple solution. In next section, we show how to correct both of the above-mentioned issues by recursively applying back-propagating through gradient ascent (Domke, 2012).

## 4 CORRECTING THE SUB-OPTIMAL BIT ALLOCATION

In this section, we first extend the generic SAVI Kim et al. (2018); Marino et al. (2018) to 2-level non-factorized latent. Then we further extend this result to latent with any dependency that can be described by a DAG (Directed Acyclic Graph). And finally, we correct the sub-optimal bit allocation by applying the result in DAG latent to NVC.

### 4.1 SAVI ON 2-LEVEL NON-FACTORIZED LATENT

In this section, we extend the SAVI on 1-level latent (Kim et al., 2018) to 2-level non-factorized latent. We denote  $\mathbf{x}$  as evidence,  $\mathbf{a}$  as the variational posterior parameter of the first level latent  $\tilde{\mathbf{a}}$ ,  $\mathbf{b}$  as the variational posterior parameter of the second level latent  $\tilde{\mathbf{b}}$ , and the ELBO to maximize as  $\mathcal{L}(\mathbf{a}, \mathbf{b})$ . The posterior  $q(\tilde{\mathbf{a}}, \tilde{\mathbf{b}}|\mathbf{x})$  factorizes as  $q(\tilde{\mathbf{a}}|\mathbf{x})q(\tilde{\mathbf{b}}|\tilde{\mathbf{a}}, \mathbf{x})$ , which means that  $\mathbf{b}$  depends on  $\mathbf{a}$ . Given  $\mathbf{a}$  is fixed, we can directly follow Kim et al. (2018); Marino et al. (2018) to optimize  $\mathbf{b}$  to maximize ELBO by SAVI. However, it requires some tricks to optimize  $\mathbf{a}$ .

**Algorithm 1:** SAVI on 2-level Latent

```

1 procedure solve-2-level( $\mathbf{x}, \mathbf{a}^k$ )
2   initialize  $\mathbf{a}^0 \leftarrow f(\mathbf{x})$  from FAVI
3   for  $k = 0, \dots, K - 1$  do
4      $\frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{a}^k} = \text{grad-2-level}(\mathbf{x}, \mathbf{a}^k)$ 
5      $\mathbf{a}^{k+1} \leftarrow \mathbf{a}^k + \alpha \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{a}^k}$ 
6   return  $\mathbf{a}^K, \mathbf{b}^K$ 

7 procedure grad-2-level( $\mathbf{x}, \mathbf{a}^k$ )
8    $\mathbf{b}^0 \leftarrow f(\mathbf{x}, \mathbf{a}^k)$  from FAVI
9   for  $k' = 0, \dots, K - 1$  do
10     $\mathbf{b}^{k'+1} \leftarrow \mathbf{b}^{k'} + \alpha \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{d\mathbf{b}^{k'}}$ 
11     $\overleftarrow{\mathbf{a}} \leftarrow \frac{\partial \mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{\partial \mathbf{a}^k}$ 
12     $\overleftarrow{\mathbf{b}}^K \leftarrow \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{b}^K}$ 
13    for  $k' = K - 1, \dots, 0$  do
14       $\overleftarrow{\mathbf{a}} \leftarrow \overleftarrow{\mathbf{a}} + \alpha \frac{\partial^2 \mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{\partial \mathbf{a}^k \partial \mathbf{b}^{k'}} \overleftarrow{\mathbf{b}}^{k'+1}$ 
15       $\overleftarrow{\mathbf{b}}^{k'} \leftarrow \overleftarrow{\mathbf{b}}^{k'} + \alpha \frac{\partial^2 \mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{\partial \mathbf{b}^{k'} \partial \mathbf{b}^{k'}} \overleftarrow{\mathbf{b}}^{k'+1}$ 
16       $\overleftarrow{\mathbf{a}} = \overleftarrow{\mathbf{a}} + \frac{\partial \mathbf{b}^0}{\partial \mathbf{a}^k} \overleftarrow{\mathbf{b}}^0$ 
17    return  $\frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{a}^k} = \overleftarrow{\mathbf{a}}$ 

```

**Algorithm 2:** SAVI on DAG Latent

```

1 procedure solve-dag( $\mathbf{x}$ )
2   sort  $\mathbf{a}_1, \dots, \mathbf{a}_N$  in topological order
3   for  $\mathbf{a}_j$  with parent  $\mathcal{P}(\mathbf{a}_j) = \emptyset$ 
4     add  $\mathbf{a}_j$  to fake node  $\mathbf{a}_0$ 's children  $\mathcal{C}(\mathbf{a}_0)$ 
5   grad-dag( $\mathbf{x}, \mathbf{a}_0^0$ )
6   return  $\mathbf{a}_1^K, \dots, \mathbf{a}_N^K$ 

7 procedure grad-dag( $\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}$ )
8   for  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)$  in topological order do
9      $\mathbf{a}_j^0 \leftarrow f(\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_{<j}^{k_{<j}})$  from FAVI
10    for  $k_j = 0, \dots, K - 1$  do
11       $\frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_j^{k_j}, \mathbf{a}_{>j}^K)}{d\mathbf{a}_j^{k_j}} \leftarrow \text{grad-dag}(\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_j^{k_j})$ 
12       $\mathbf{a}_j^{k_j+1} \leftarrow \mathbf{a}_j^{k_j} + \alpha \frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_j^{k_j}, \mathbf{a}_{>j}^K)}{d\mathbf{a}_j^{k_j}}$ 
13     $\overleftarrow{\mathbf{a}}_i \leftarrow \frac{\partial \mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{\partial \mathbf{a}_i^{k_i}}$ 
14    for  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)$  do
15       $\overleftarrow{\mathbf{a}}_j \leftarrow \mathbf{0}, \mathbf{a}_j^K \leftarrow \frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>j}^K)}{d\mathbf{a}_j^K}$ 
16      for  $k_j = K - 1, \dots, 0$  do
17         $\overleftarrow{\mathbf{a}}_j \leftarrow \overleftarrow{\mathbf{a}}_j + \alpha \frac{\partial^2 \mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_j^{k_j}, \mathbf{a}_{>j}^K)}{\partial \mathbf{a}_i^{k_i} \partial \mathbf{a}_j^{k_j}} \overleftarrow{\mathbf{a}}_j^{k_j+1}$ 
18         $\overleftarrow{\mathbf{a}}_j^{k_j} \leftarrow \overleftarrow{\mathbf{a}}_j^{k_j+1} + \alpha \frac{\partial^2 \mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_j^{k_j}, \mathbf{a}_{>j}^K)}{\partial \mathbf{a}_j^{k_j} \partial \mathbf{a}_j^{k_j}} \overleftarrow{\mathbf{a}}_j^{k_j+1}$ 
19       $\overleftarrow{\mathbf{a}}_i \leftarrow \overleftarrow{\mathbf{a}}_i + \overleftarrow{\mathbf{a}}_j + \frac{\partial \mathbf{a}_j^0}{\partial \mathbf{a}_i^{k_i}} \overleftarrow{\mathbf{a}}_j^0$ 
20    return  $\frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{d\mathbf{a}_i^{k_i}} = \overleftarrow{\mathbf{a}}_i$ 

```

The intuition is, we do not want to find a  $\mathbf{a}$  that maximizes  $\mathcal{L}(\mathbf{a}, \mathbf{b})$  given a fixed  $\mathbf{b}$  (or we have the gradient issue described in Sec. 3). Instead, we want to find a  $\mathbf{a}$ , whose  $\max_{\mathbf{b}} \mathcal{L}(\mathbf{a}, \mathbf{b})$  is maximum. This translates to the optimization problem as Eq. 12. In fact, Eq. 12 is a variant of setup in back-propagating through gradient ascent (Samuel & Tappen, 2009; Domke, 2012). The difference is, our  $\mathbf{a}$  also contributes directly to optimization target  $\mathcal{L}(\mathbf{a}, \mathbf{b})$ . From this perspective, Eq. 12 is more closely connected to Kim et al. (2018), if we treat  $\mathbf{a}$  as the model parameter and  $\mathbf{b}$  as latent.

$$\mathbf{a} \leftarrow \arg \max_{\mathbf{a}} \mathcal{L}(\mathbf{a}, \mathbf{b}^*(\mathbf{a})), \text{ where } \mathbf{b}^*(\mathbf{a}) \leftarrow \arg \max_{\mathbf{b}} \mathcal{L}(\mathbf{a}, \mathbf{b}) \quad (12)$$

And as SAVI on 1-level latent (Kim et al., 2018; Marino et al., 2018), we need to solve Eq. 12 using gradient ascent. Specifically, denote  $\alpha$  as step size (learning rate),  $K$  as the total gradient ascent steps,  $\mathbf{a}^k$  as the  $\mathbf{a}$  after  $k$  step update,  $\mathbf{b}^{k'}$  as the  $\mathbf{b}$  after  $k'$  step update, and  $f(\cdot)$  as FAVI procedure generating initial posterior parameters  $\mathbf{a}^0, \mathbf{b}^0$ , the optimization problem as Eq. 12 translates into the update rule as Eq. 13. Eq. 13 is the guidance for designing optimization algorithm, and it also explains why the gradient of BAO (Xu et al., 2022) is evaluated at wrong value (See Sec. 3.2).

$$\mathbf{a}^{k+1} \leftarrow \mathbf{a}^k + \alpha \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{a}^k}, \mathbf{b}^{k'+1} \leftarrow \mathbf{b}^{k'} + \alpha \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{d\mathbf{b}^{k'}}, \text{ where } \mathbf{b}^0 = f(\mathbf{x}, \mathbf{a}^k) \quad (13)$$

To solve Eq. 13, we note that although  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})/d\mathbf{b}^{k'}$  is directly computed,  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k$  is not straightforward. Resorting to previous works (Samuel & Tappen, 2009; Domke, 2012) in implicit differentiation and extending the results in Kim et al. (2018) from model parameters to variational posterior parameters, we implement Eq. 13 as Alg. 1. Specifically, we first initialize  $\mathbf{a}^0$  from FAVI. Then we conduct gradient ascent on  $\mathbf{a}$  with gradient  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k$  computed from the procedure grad-2-level( $\mathbf{x}, \mathbf{a}^k$ ). And inside grad-2-level( $\mathbf{x}, \mathbf{a}^k$ ),  $\mathbf{b}$  is also updated by gradient

ascent, the above procedure corresponds to Eq. 13. The key of Alg. 1 is the evaluation of gradient  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k$ . Formally, we have:

**Theorem 1.** *After  $\text{grad-2-level}(\mathbf{x}, \mathbf{a}^k)$  of Alg. 1 executes, we have the return value  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k = \overleftarrow{\mathbf{a}}$ . (See proof in Appendix. A.1.)*

#### 4.2 SAVI ON DAG-DEFINED NON-FACTORIZED LATENT

In this section, we extend the result from previous section to SAVI on general non-factorized latent with dependency described by any DAG. This DAG is the computational graph during network inference, and it is also the directed graphical model (DGM) (Koller & Friedman, 2009) defining the factorization of latent variables during inference. This is the general case covering all dependency that can be described by DGM. This extension is necessary to perform SAVI on latent with complicated dependency (e.g. bit allocation of NVC).

Similar to the 2-level latent setup, we consider performing SAVI on  $N$  variational posterior parameter  $\mathbf{a}_1, \dots, \mathbf{a}_N$  with their dependency defined by a computational graph  $\mathcal{G}$ , i.e., their corresponding latent variable  $\tilde{\mathbf{a}}_1, \dots, \tilde{\mathbf{a}}_N$ 's posterior distribution factorizes as  $\mathcal{G}$ . Specifically, we denote  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i), \mathbf{a}_i \in \mathcal{P}(\mathbf{a}_j)$  if an edge exists from  $\mathbf{a}_i$  to  $\mathbf{a}_j$ . This indicates that  $\tilde{\mathbf{a}}_j$  conditions on  $\tilde{\mathbf{a}}_i$ . Without loss of generality, we assume  $\mathbf{a}_1, \dots, \mathbf{a}_N$  is sorted in topological order. This means that if  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i), \mathbf{a}_i \in \mathcal{P}(\mathbf{a}_j)$ , then  $i < j$ . Each latent is optimized by  $K$ -step gradient ascent, and  $\mathbf{a}_i^{k_i}$  denotes the latent  $\mathbf{a}_i$  after  $k_i$  steps of update. Then, similar to 2-level latent, we have the update rule as Eq. 14:

$$\mathbf{a}_i^{k_i+1} \leftarrow \mathbf{a}_i^{k_i} + \alpha \frac{d\mathcal{L}(\mathbf{a}_1^{k_1}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{d\mathbf{a}_i^{k_i}}, \text{ where } \mathbf{a}_{>i}^0 = f(\mathbf{x}, \mathbf{a}_1^{k_1}, \dots, \mathbf{a}_i^{k_i}) \quad (14)$$

, which can be translated into Alg. 2. Specifically, we first sort the latent in topological order. Then, we add a fake latent  $\mathbf{a}_0$  to the front of all  $\mathbf{a}$ s. Its children are all the  $\mathbf{a}$ s with 0 in-degree. Then, we can solve the SAVI on  $\mathbf{a}_1, \dots, \mathbf{a}_N$  using gradient ascent by executing the procedure  $\text{grad-dag}(\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i})$  in Alg. 2 recursively. Inside procedure  $\text{grad-dag}(\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i})$ , the gradient to update  $\mathbf{a}_i$  relies on the convergence of its children  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)$ , which is implemented by the recursive depth-first search (DFS) in line 11. And upon the completion of procedure  $\text{grad-dag}(\mathbf{x}, \mathbf{a}_0^0)$ , all the latent converges to  $\mathbf{a}_1^K, \dots, \mathbf{a}_N^K$ . Similar to the 2-level latent case, the key of Alg. 2 is the evaluation of gradient  $d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)/d\mathbf{a}_i^{k_i}$ . Formally, we have:

**Theorem 2.** *After the procedure  $\text{grad-dag}(\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i})$  in Alg. 2 executes, we have the return value  $d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)/d\mathbf{a}_i^{k_i} = \overleftarrow{\mathbf{a}}_i$ . (See proof in Appendix. A.1.)*

To better understand how Alg. 2 works, we provide a detailed example in Fig. 5 of Appendix. A.3.

#### 4.3 CORRECTING THE SUB-OPTIMAL BIT ALLOCATION USING SAVI ON DAG

With the result in previous section, correcting BAO (Xu et al., 2022) seems to be trivial. We only need to sort the latent in topological order as  $\mathbf{w}_1, \mathbf{y}_1, \dots, \mathbf{w}_T, \mathbf{y}_T$ , treat them as  $\mathbf{a}_1, \dots, \mathbf{a}_{2T+1}$  and run Alg. 2 to obtain the optimized latent parameters  $\mathbf{w}_1^K, \mathbf{y}_1^K, \dots, \mathbf{w}_T^K, \mathbf{y}_T^K$ . And the gradient  $d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)/d\mathbf{a}_i^{k_i}$  computed in Alg. 2 resolves the issue of BAO described in Sec. 3.1 and Sec. 3.2. However, an evident problem is the temporal complexity. Given the latent number  $N$  and gradient ascent step number  $K$ , Alg. 2 has temporal complexity of  $\Theta(K^N)$ . NVC with GoP size 10 has approximately  $N = 20$  latent, and the SAVI on NVC (Xu et al., 2022) takes around  $K = 2000$  step to converge. For bit allocation, the complexity of Alg. 2 is  $\approx 2000^{20}$ , which is intractable. On the other hand, BAO's complexity is reasonable ( $\Theta(KN) \approx 4 \times 10^4$ ). Thus, in next section, we provide a feasible approximation to such intractable corrected bit allocation.

#### 4.4 FEASIBLE APPROXIMATION TO THE CORRECTED BIT ALLOCATION

In order to solve problem with practical size such as bit allocation on NVC, we provide an approximation to the SAVI (Kim et al., 2018; Marino et al., 2018) on DAG described in Sec. 4.2. The general idea is that, when being applied to bit allocation of NVC, the accurate SAVI on DAG (Alg. 2) satisfies both requirement on gradient signal described in Sec. 3.1 and Sec. 3.2. We can not make it

tractable without breaking them. Thus, we break one of them and achieve a reasonable complexity, while maintain a superior performance compared with BAO (Xu et al., 2022).

We consider the approximation in Eq. 15 which breaks the requirement for gradient evaluation in Sec. 3.2. Based on Eq. 15 and the requirement in Sec. 3.1, we design an approximation of accurate SAVI as Alg. 4. When being applied to bit allocation in NVC, it satisfies the gradient requirement in Sec. 3.1 while maintaining a temporal complexity of  $\Theta(KN)$  as BAO.

$$\frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^{K_i})}{d\mathbf{a}_i^{k_i}} \approx \frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^0)}{d\mathbf{a}_i^{k_i}} \quad (15)$$

Specifically, with the approximation in Eq. 15, the recurrent gradient computation in Alg. 2 becomes unnecessary as the right hand side of Eq. 15 does not require  $\mathbf{a}_{>i}^K$ . However, to maintain the dependency of latent described in Sec. 3.1, as Alg. 2, we still need to ensure that the children node  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)$  are re-initialized by FAVI every-time when  $\mathbf{a}_i$  is updated. Therefore, a reasonable approach is to traverse the graph in topological order. We keep the children node  $\mathbf{a}_j$  untouched until all its parent node  $\mathbf{a}_i \in \mathcal{P}(\mathbf{a}_j)$ 's gradient ascent is completed and  $\mathbf{a}_i^K$  is known. And the resulting approximate SAVI algorithm is as Alg. 4. When applied to bit allocation, it satisfies the gradient requirement in Sec. 3.1, and as BAO, its temporal complexity is  $\Theta(KN)$ .

Algorithm 3: BAO on DAG Latent	Algorithm 4: Approximate SAVI on DAG latent
<pre> 1 <b>procedure</b> solve-bao(<math>\mathbf{x}</math>) 2   <math>\mathbf{a}_1^0, \dots, \mathbf{a}_N^0 \leftarrow f(\mathbf{x})</math> from FAVI 3   <b>for</b> <math>k = 0, \dots, K - 1</math> <b>do</b> 4     <b>for</b> <math>i = 1, \dots, N</math> <b>do</b> 5       <math>\mathbf{a}_i^{k+1} \leftarrow \mathbf{a}_i^k + \alpha \frac{\partial \mathcal{L}(\mathbf{a}_1^k, \dots, \mathbf{a}_N^k)}{\partial \mathbf{a}_i^k}</math> 6   <b>return</b> <math>\mathbf{a}_1^K, \dots, \mathbf{a}_N^K</math> </pre>	<pre> 1 <b>procedure</b> solve-approx-dag(<math>\mathbf{x}</math>) 2   sort <math>\mathbf{a}_1, \dots, \mathbf{a}_N</math> in topological order 3   <b>for</b> <math>i = 1, \dots, N</math> <b>do</b> 4     <math>\mathbf{a}_i^0, \dots, \mathbf{a}_N^0 \leftarrow f(\mathbf{x}, \mathbf{a}_{&lt;i}^K)</math> from FAVI 5     <b>for</b> <math>k = 0, \dots, K - 1</math> <b>do</b> 6       <math>\frac{d\mathcal{L}(\mathbf{a}_{&lt;i}^K, \mathbf{a}_i^k, \mathbf{a}_{&gt;i}^K)}{d\mathbf{a}_i^k} \approx \frac{d\mathcal{L}(\mathbf{a}_{&lt;i}^K, \mathbf{a}_i^k, \mathbf{a}_{&gt;i}^0)}{d\mathbf{a}_i^k}</math> 7       <math>\mathbf{a}_i^{k+1} \leftarrow \mathbf{a}_i^k + \alpha \frac{d\mathcal{L}(\mathbf{a}_{&lt;i}^K, \mathbf{a}_i^k, \mathbf{a}_{&gt;i}^K)}{d\mathbf{a}_i^k}</math> 8   <b>return</b> <math>\mathbf{a}_1^K, \dots, \mathbf{a}_N^K</math> </pre>

To better understand BAO (Xu et al., 2022) in SAVI context, we rewrite it by general SAVI notation instead of NVC notation in Alg. 3. We highlight the difference between BAO (Alg. 3) (Xu et al., 2022), the accurate SAVI on DAG latent (Alg. 2) and the approximate SAVI on DAG latent (Alg. 4) from several aspects:

- **Graph Traversal Order:** BAO performs gradient ascent on  $\mathbf{a}_{1:T}$  all together. The accurate SAVI only updates  $\mathbf{a}_i$  when  $\mathbf{a}_{>i}$ 's update is complete and  $\mathbf{a}_{<i}^K$  is known. The approximate SAVI only updates  $\mathbf{a}_i$  when  $\mathbf{a}_{<i}$ 's update is complete and  $\mathbf{a}_{<i}^K$  is known.
- **Gradient Correctness:** When being applied to bit allocation in NVC, BAO violates the gradient rule in Sec. 3.1 and Sec. 3.2, accurate SAVI satisfies both rules, approximate SAVI satisfies Sec. 3.1 and violates Sec. 3.2.
- **Temporal Complexity:** With the latent number  $N$  and steps of gradient ascent  $K$ , the complexity of BAO is  $\Theta(KN)$ , the complexity of accurate SAVI is  $\Theta(K^N)$  and the complexity of approximate SAVI is  $\Theta(KN)$ .

Then we can simply apply Alg. 4 to bit allocation in NVC to obtain a feasible approximation of the corrected optimal bit allocation. And in Sec. 6.2, we empirically show that our approximation improves the R-D performance over BAO (Xu et al., 2022) with even smaller number of updates.

## 5 RELATED WORK: BIT ALLOCATION & SAVI FOR NEURAL COMPRESSION

Li et al. (2022) are the pioneer of bit allocation for NVC and their work is elaborated in Sec. 2.2. Other recent works that consider bit allocation for NVC only adopt simple heuristic such as inserting 1 high quality frame per 4 frames (Hu et al., 2022; Cetin et al., 2022). On the other hand, OEU (Lu et al., 2020) is also recognised as frame-level bit allocation while its performance is inferior than BAO (Xu et al., 2022). BAO is the most recent work with best R-D performance. It is elaborated in Sec. 2.2 and Sec. 3, and corrected in the previous section.

Semi-Amortized Variational Inference (SAVI) is proposed by Kim et al. (2018); Marino et al. (2018). The idea is that works following Kingma & Welling (2013) use fully amortized inference parameter  $\phi$  for all data, which leads to the amortization gap (Cremer et al., 2018). SAVI reduces this gap by optimizing the variational posterior parameter after initializing it with inference network. It adopts back-propagating through gradient ascent (Domke, 2012) to evaluate the gradient of model parameters. We adopt a similar method to extend SAVI to non-factorized latent. When applying SAVI to practical neural codec, researchers abandon the nested model parameter update for efficiency. Prior works (Djelouah & Schroers, 2019; Yang et al., 2020b; Zhao et al., 2021; Gao et al., 2022) adopt SAVI to boost R-D performance and achieve variable bitrate in image compression. And BAO (Xu et al., 2022) is the first to consider SAVI for bit allocation.

## 6 EXPERIMENTS

### 6.1 EXPERIMENTAL SETTINGS

We implement our approach in PyTorch 1.9 with CUDA 11.2, and run the experiments on NVIDIA(R) A100 GPU. Most of the other settings are intentionally kept the same as BAO (Xu et al., 2022). Specifically, we adopt HEVC Common Testing Condition (CTC) (Bossen et al., 2013) and UVG dataset (Mercat et al., 2020). And we measure the R-D performance in Bjontegaard-Bitrate (BD-BR) and BD-PSNR (Bjontegaard, 2001). For baseline NVC (Lu et al., 2019; Li et al., 2021), we adopt the official pre-trained models. And we select target  $\lambda_0 = \{256, 512, 1024, 2048\}$ . For gradient ascent, we adopt Adam (Kingma & Ba, 2014) optimizer with  $lr = 1 \times 10^{-3}$ . We set the gradient ascent step  $K = 2000$  for the first frame and  $K = 400$  for other frames. More details are presented in Appendix. A.5.

### 6.2 QUANTITATIVE RESULTS

As shown in Tab. 1, our method consistently improves the R-D performance in terms of BD-BR over BAO (Xu et al., 2022) on both baseline methods and all datasets. Moreover, this improvement is especially significant (more than 10% in BD-BR) when the baseline is DCVC (Li et al., 2021). And both BAO and our proposed correction significantly outperform other approaches. It is also noteworthy that with our bit allocation, DVC (the SOTA method in 2019) already outperforms DCVC (the SOTA method in 2021) by large margin (See the red solid line and black dash line in Fig. 2).

Method	BD-BR (%) ↓				
	Class B	Class C	Class D	Class E	UVG
<i>DVC (Lu et al., 2019) as Baseline</i>					
Li et al. (2016) <sup>1</sup>	20.21	17.13	13.71	10.32	16.69
Li et al. (2022) <sup>1</sup>	-6.80	-2.96	0.48	-6.85	-4.12
OEU (Lu et al., 2020) <sup>2</sup>	-13.57	-11.29	-18.97	-12.43	-13.78
BAO (Xu et al., 2022) <sup>2</sup>	-28.55	-26.82	-25.37	-32.54	-27.68
Proposed	-32.10	-31.71	-35.86	-32.93	-30.92
<i>DCVC (Li et al., 2021) as Baseline</i>					
OEU (Lu et al., 2020) <sup>2</sup>	-10.75	-14.34	-16.30	-7.15	-16.07
BAO (Xu et al., 2022) <sup>2</sup>	-20.59	-19.69	-20.60	-23.33	-25.22
Proposed	-32.89	-33.10	-32.01	-36.88	-39.66

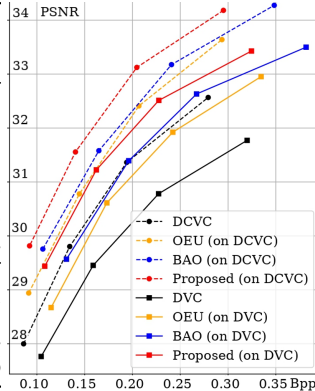


Table 1: The BD-BR of our approach compared with others. <sup>1</sup> comes from Li et al. (2022). <sup>2</sup> comes from Xu et al. (2022).

Figure 2: The R-D curve on HEVC Class D.

Other than R-D performance, the bitrate error of our approach is also significantly smaller than BAO (Xu et al., 2022) (See Tab. 2). The bitrate error is measured as the relative bitrate difference before and after bit allocation. The smaller it is, the easier it is to achieve the desired bitrate accurately. For complexity, our approach only performs 920 steps of gradient ascent per-frame, while BAO requires 2000 steps.

See more quantitative results (BD-PSNR & R-D curves) in Appendix. A.6.



### 6.3 ABLATION STUDY, ANALYSIS & QUALITATIVE RESULTS

Tab. 3 shows that for BAO (Xu et al., 2022), jointly optimizing  $w_{1:T}, y_{1:T}$  performs worse than optimizing  $y_{1:T}$  or  $w_{1:T}$  alone. This counter-intuitive phenomena comes from its incorrect estimation of gradient signal. For the proposed approach that corrects this, jointly optimizing  $w_{1:T}, y_{1:T}$  performs better than optimizing  $y_{1:T}$  or  $w_{1:T}$  alone, which is aligned with our intuition.

Method	Bitrate-Error (%) ↓					Method	BD-BR (%) ↓
	Class B	Class C	Class D	Class E	UVG		
<i>DVC (Lu et al., 2019) as Baseline</i>						BAO ( $y$ )	-25.37
BAO (Xu et al., 2022) <sup>2</sup>						BAO ( $w$ )	-22.24
Proposed						BAO ( $y, w$ )	-14.76
DCVC (Li et al., 2021) as Baseline						Proposed ( $y$ )	-32.60
BAO (Xu et al., 2022) <sup>2</sup>						Proposed ( $w$ )	-31.56
Proposed						Proposed ( $y, w$ )	-35.86

Table 2: The bitrate error of our approach compared with BAO.

Table 3: Ablation study with HEVC Class D and DVC (Lu et al., 2019).

To better understand why our method works, we present the R-D cost, distortion and rate versus frame/latent index for different methods in Fig. 3: *top-left* shows that the R-D cost of our approach consistently decreases according to SAVI stage. Moreover, it outperforms BAO after 4<sup>th</sup> frame; *top-right* shows that for each frame the R-D cost of our method is lower than BAO; *bottom-left* shows that the distortion part of R-D cost of our approach is approximately the same as BAO. While *bottom-right* shows that the advantage of our approach over BAO lies in the bitrate. More specifically, BAO increases the bitrate of  $y_i$ s after SAVI, while our correction decreases it.

See more analysis in Appendix. A.9 and qualitative results in Appendix. A.10.

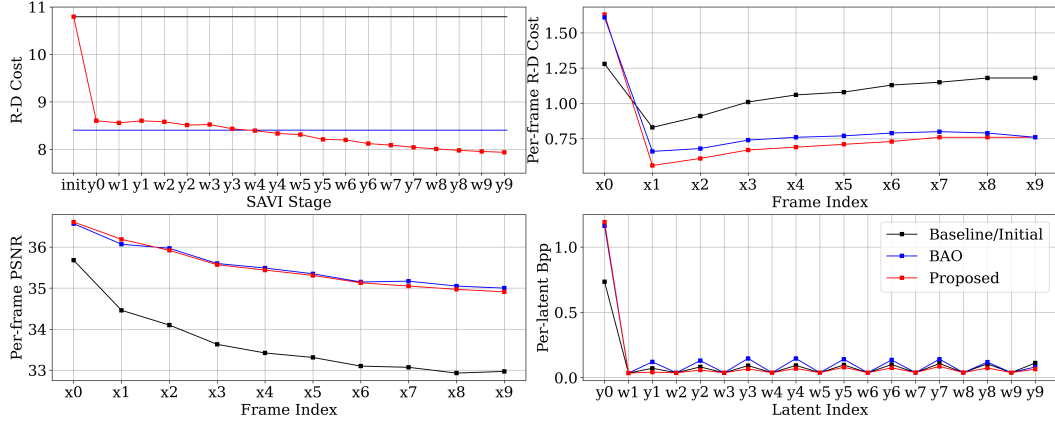


Figure 3: *top-left*. R-D cost vs. SAVI stage. *top-right*. R-D cost vs. frame index. *bottom-left*. PSNR vs. frame index. *bottom-right*. bpp vs. latent index. See enlarged-version in Appendix. A.9.

## 7 DISCUSSION & CONCLUSION

Despite our correction is already more efficient than original BAO (Xu et al., 2022), its encoding speed remains far from real-time. Thus, it is limited to scenarios where R-D performance matters much more than encoding time (e.g. video on demand). See more discussion in Appendix. A.11.

To conclude, we show that a previous bit allocation method for NVC is sub-optimal as it abuses SAVI on non-factorized latent. Then, we propose the correct SAVI on general non-factorized latent by back-propagating through gradient ascent, and we further propose a feasible approximation to make it tractable for bit allocation. Experimental results show that our correction significantly improves the R-D performance.

## ETHICS STATEMENT

Improving the R-D performance of NVC has positive social value, in terms of reducing carbon emission by saving the resources required to transfer and store videos. Moreover, unlike traditional codecs such as H.266 (Bross et al., 2021), neural video codec does not require dedicated hardware. Instead, it can be deployed with general neural accelerators. Improving the R-D performance of NVC prompts the practical deployment of video codecs that are independent of dedicated hardware, and lowers the hardware-barrier of playing multi-media contents.

## REPRODUCIBILITY STATEMENT

For theoretical results, both of the two theorems are followed by proof in Appendix. A.1. For a relatively complicated novel algorithm (Alg. 2), we provide an illustration of the step by step execution procedure in Appendix. A.3. For experiment, both of the two datasets are publicly accessible. In Appendix. A.5, we provide more implementation details including all the hyper-parameters. Moreover, we provide our source code for reproducing the empirical results in supplementary material.

## REFERENCES

- Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016.
- Gisle Bjontegaard. Calculation of average psnr differences between rd-curves. *VCEG-M33*, 2001.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Frank Bossen et al. Common test conditions and software reference configurations. *JCTVC-L1100*, 12(7), 2013.
- Benjamin Bross, Jianle Chen, Jens-Rainer Ohm, Gary J Sullivan, and Ye-Kui Wang. Developments in international video coding standardization after avc, with an overview of versatile video coding (vvc). *Proc. IEEE*, 109(9):1463–1493, 2021.
- Eren Cetin, M Akın Yılmaz, and A Murat Tekalp. Flexible-rate learned hierarchical bi-directional video compression with motion refinement and frame-level bit allocation. *arXiv e-prints*, pp. arXiv–2206, 2022.
- Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders. In *Int. Conf. on Machine Learning*, pp. 1078–1086. PMLR, 2018.
- Joaquim Campos Meierhans Simon Abdelaziz Djelouah and Christopher Schroers. Content adaptive optimization for neural image compression. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit*, 2019.
- Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pp. 318–326. PMLR, 2012.
- Chenjian Gao, Tongda Xu, Dailan He, Hongwei Qin, and Yan Wang. Flexible neural image compression via code editing. *arXiv preprint arXiv:2209.09422*, 2022.
- Zhihao Hu, Guo Lu, Jinyang Guo, Shan Liu, Wei Jiang, and Dong Xu. Coarse-to-fine deep video coding with hyperprior-guided mode prediction. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit*, pp. 5921–5930, 2022.
- Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized variational autoencoders. In *Int. Conf. on Machine Learning*, pp. 2678–2687. PMLR, 2018.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Bin Li, Houqiang Li, Li Li, and Jinlei Zhang.  $\lambda$  domain rate control algorithm for high efficiency video coding. *IEEE Trans. Image Process.*, 23(9):3841–3854, 2014.
- Jiahao Li, Bin Li, and Yan Lu. Deep contextual video compression. *Advances in Neural Information Processing Systems*, 34, 2021.
- Li Li, Bin Li, Houqiang Li, and Chang Wen Chen.  $\lambda$ -domain optimal bit allocation algorithm for high efficiency video coding. *IEEE Trans. Circuits Syst. Video Technol.*, 28(1):130–142, 2016.
- Yanghao Li, Xinyao Chen, Jisheng Li, Jiangtao Wen, Yuxing Han, Shan Liu, and Xiaozhong Xu. Rate control for learned video compression. In *ICASSP 2022-2022 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2829–2833. IEEE, 2022.
- Guo Lu, Wanli Ouyang, Dong Xu, Xiaoyun Zhang, Chunlei Cai, and Zhiyong Gao. Dvc: An end-to-end deep video compression framework. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit*, pp. 11006–11015, 2019.
- Guo Lu, Chunlei Cai, Xiaoyun Zhang, Li Chen, Wanli Ouyang, Dong Xu, and Zhiyong Gao. Content adaptive and error propagation aware deep video compression. In *European Conference on Computer Vision*, pp. 456–472. Springer, 2020.
- Joe Marino, Yisong Yue, and Stephan Mandt. Iterative amortized inference. In *Int. Conf. on Machine Learning*, pp. 3403–3412. PMLR, 2018.
- Alexandre Mercat, Marko Viitanen, and Jarno Vanne. Uvg dataset: 50/120fps 4k sequences for video codec analysis and development. In *Proceedings of the 11th ACM Multimedia Systems Conference*, pp. 297–302, 2020.
- David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors for learned image compression. *Advances in neural information processing systems*, 31, 2018.
- Kegan GG Samuel and Marshall F Tappen. Learning optimized map estimates in continuously-valued mrf models. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 477–484. IEEE, 2009.
- Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár. Lossy image compression with compressive autoencoders. In *5th Int. Conf. on Learning Representations, ICLR 2017*, 2017.
- Tongda Xu, Han Gao, Chenjian Gao, Jinyong Pi, Yanghao Li, Yuanyuan Wang, Ziyu Zhu, Dailan He, Mao Ye, Hongwei Qin, and Yan Wang. Bit allocation using optimization. *arXiv preprint arXiv:2209.09422*, 2022.
- Ruihan Yang, Yibo Yang, Joseph Marino, and Stephan Mandt. Hierarchical autoregressive modeling for neural video compression. *arXiv preprint arXiv:2010.10258*, 2020a.
- Yibo Yang, Robert Bamler, and Stephan Mandt. Improving inference for neural image compression. *Advances in Neural Information Processing Systems*, 33:573–584, 2020b.
- Jing Zhao, Bin Li, Jiahao Li, Ruiqin Xiong, and Yan Lu. A universal encoder rate distortion optimization framework for learned compression. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit*, pp. 1880–1884, 2021.

## A APPENDIX

## A.1 PROOF OF THM 1 AND THM 2

**Theorem 1.** *After the procedure  $\text{grad-2-level}(\mathbf{x}, \mathbf{a}^k)$  of Alg. 1 executes, we have the return value  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k = \overleftarrow{\mathbf{a}}$ .*

*Proof.* This proof extends the proof of Thm. 1 in Domke (2012), and it also serves as a formal justification of Alg. 1 in Kim et al. (2018). Note that our paper and Kim et al. (2018) are subtly different from Samuel & Tappen (2009); Domke (2012) as our high level parameter  $\mathbf{w}$  not only generate low level parameter  $\mathbf{y}$ , but also directly contributes to optimization target (See Fig. 4).

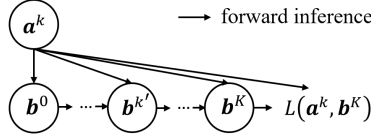


Figure 4: The computational graph corresponding to Eq. 16

As the computational graph in Fig. 4 shows, we can expand  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k$  as Eq. 16, with each term solved in Eq. 18 and Eq. 19.

$$\frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{a}^k} = \underbrace{\frac{\partial \mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{\partial \mathbf{a}^k}}_{\text{known}} + \sum_{k'=0}^K \underbrace{\frac{\partial \mathbf{b}^{k'}}{\partial \mathbf{a}^k}}_{\text{Eq. 18}} \underbrace{\frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{b}^{k'}}}_{\text{Eq. 19}} \quad (16)$$

To solve Eq. 16, we first note that  $\partial \mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/\partial \mathbf{a}^k$ ,  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{b}^K$ ,  $\partial \mathbf{b}^0/\partial \mathbf{a}^k$  is naturally known. Then, by taking partial derivative of the update rule of gradient ascent  $\mathbf{b}^{k'+1} \leftarrow \mathbf{b}^{k'} + \alpha d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})/d\mathbf{b}^{k'}$  with regard to  $\mathbf{a}^k, \mathbf{b}^{k'}$ , we have Eq. 17 and Eq. 18. Note that Eq. 18 is the partial derivative  $\partial \mathbf{b}^{k'+1}/\partial \mathbf{a}^k$  instead of total derivative  $d\mathbf{b}^{k'+1}/d\mathbf{a}^k = (\partial \mathbf{b}^{k'+1}/\partial \mathbf{b}^{k'}) (d\mathbf{b}^{k'}/d\mathbf{a}^k) + \partial \mathbf{b}^{k'+1}/\partial \mathbf{a}^k$ .

$$\frac{\partial \mathbf{b}^{k'+1}}{\partial \mathbf{b}^{k'}} = I + \alpha \frac{\partial^2 \mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{\partial \mathbf{b}^{k'} \partial \mathbf{b}^{k'}} \quad (17)$$

$$\frac{\partial \mathbf{b}^{k'+1}}{\partial \mathbf{a}^k} = \alpha \frac{\partial^2 \mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{\partial \mathbf{a}^k \partial \mathbf{b}^{k'}} \quad (18)$$

And those second order terms can either be directly evaluated or approximated via finite difference as Eq. 20. As Eq. 18 already solves the first term on the right hand side of Eq. 16, the remaining issue is  $d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{b}^{k'}$ . To solve this term, we expand it recursively as Eq. 19 and take Eq. 17 into it.

$$\frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{b}^{k'}} = \frac{\partial \mathbf{b}^{k'+1}}{\partial \mathbf{b}^{k'}} \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)}{d\mathbf{b}^{k'+1}} \quad (19)$$

And the above solving process can be described by the procedure  $\text{grad-2-level}(\mathbf{x}, \mathbf{a}^k)$  of Alg. 1.

Specifically, the iterative update of  $\mathbf{b}^{k'+1}$  in line 15 corresponds to recursively expanding Eq. 19 with Eq. 17, and the iterative update of  $\overleftarrow{\mathbf{a}}$  in line 14 corresponds to recursively expanding Eq. 16 with Eq. 18 and Eq. 19. Upon the return of  $\text{grad-2-level}(\mathbf{x}, \mathbf{a}^k)$  of Alg. 1, we have  $\overleftarrow{\mathbf{a}} = d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^K)/d\mathbf{a}^k$ .

The complexity of the Hessian-vector product in line 14 and 15 of Alg. 1 may be reduced using finite difference following (Domke, 2012) as Eq. 20.

$$\begin{aligned} \frac{\partial^2 \mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{\partial \mathbf{a}^k \partial \mathbf{b}^{k'}} \mathbf{v} &= \lim_{r \rightarrow 0} \frac{1}{r} \left( \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'} + r\mathbf{v})}{d\mathbf{a}^k} - \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{d\mathbf{a}^k} \right) \\ \frac{\partial^2 \mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{\partial \mathbf{b}^{k'} \partial \mathbf{b}^{k'}} \mathbf{v} &= \lim_{r \rightarrow 0} \frac{1}{r} \left( \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'} + r\mathbf{v})}{d\mathbf{b}^{k'}} - \frac{d\mathcal{L}(\mathbf{a}^k, \mathbf{b}^{k'})}{d\mathbf{b}^{k'}} \right) \end{aligned} \quad (20)$$

□

**Theorem 2.** *After the procedure `grad-dag`( $\mathbf{x}, \mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}$ ) in Alg. 2 executes, we have the return value  $d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)/d\mathbf{a}_i^{k_i} = \overleftarrow{\mathbf{a}}_i$ .*

*Proof.* Consider computing the target gradient with DAG  $\mathcal{G}$ . The  $\mathbf{a}_i^{k_i}$ 's gradient is composed of its own contribution to  $\mathcal{L}$  in addition to the gradient from its children  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)$ . Further, as we are considering the optimized children  $\mathbf{a}_j^K$ , we expand the children node  $\mathbf{a}_j$  as Fig. 4. Then, we have:

$$\frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{d\mathbf{a}_i^{k_i}} = \underbrace{\frac{\partial \mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{\partial \mathbf{a}_i^{k_i}}}_{\text{known}} + \sum_{\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)} \left( \underbrace{\sum_{k_j=0}^K \frac{\partial \mathbf{a}_j^{k_j}}{\partial \mathbf{a}_i^{k_i}}}_{\text{Eq. 18}} \underbrace{\frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_{j-1}^{k_{j-1}}, \mathbf{a}_{\geq j}^K)}{d\mathbf{a}_j^{k_j}}}_{\text{Eq. 19}} \right) \quad (21)$$

The first term on the right-hand side of Eq. 21 can be trivially evaluated. The  $\partial \mathbf{a}_j^{k_j} / \partial \mathbf{a}_i^{k_i}$  can be evaluated as Eq. 18. And the  $d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_{j-1}^{k_{j-1}}, \mathbf{a}_{\geq j}^K) / d\mathbf{a}_j^{k_j}$  can be iteratively expanded as Eq. 19. We highlight several key differences between Alg. 2 and Alg. 1 which are reflected in the implementation of Alg. 2:

- The gradient evaluation of current node  $\mathbf{y}_i$  requires gradient of its plural direct children  $\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)$ , instead of the single child in 2-level case. The children traversal part of Eq. 19 corresponds to the two extra for loop in line 8 and 14 of Alg. 2.
- The gradient ascent update of child latent parameter  $\mathbf{a}_j^{k_j+1} \leftarrow \mathbf{a}_j^{k_j} + \alpha d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_j^{k_j}, \mathbf{a}_{>j}^K) / d\mathbf{a}_j^{k_j}$  can be conducted trivially only if  $\mathcal{C}(\mathbf{a}_j)$  is empty, otherwise the gradient has to be evaluated recursively using Eq. 21. And this part corresponds to the recursive call in line 11 of Alg. 2.

And the other part of Alg. 2 is the same as Alg. 1. So the rest of the proof follows Thm. 1. Similarly, the Hessian-vector product in line 17 and 18 of Alg. 2 may be approximated as Eq. 20. However, this does not save Alg. 2 from an overall complexity of  $\Theta(K^N)$ .  $\square$

## A.2 THE COMPLETE FORMULA FOR SEC. 3.1 AND SEC. 3.2

In this section, we provide the complete formula on  $\mathbf{y}_i$  related gradient for Sec. 3.1 and Sec. 3.2. Specifically, Eq. 22 is paired with Eq. 10, and Eq. 23 is paired with Eq. 11.

$$\begin{aligned} \frac{d\mathcal{L}(\mathbf{w}_{1:T}, \mathbf{y}_{1:T})}{d\mathbf{y}_i} &= \sum_{j=i}^T \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{d\mathbf{y}_i} \\ \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{d\mathbf{y}_i} &= \underbrace{\sum_{l=i+1}^j \left( \frac{\partial \mathbf{y}_l}{\partial \mathbf{y}_i} \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{d\mathbf{y}_l} + \frac{\partial \mathbf{w}_l}{\partial \mathbf{y}_i} \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{d\mathbf{w}_l} \right)}_{\text{ignored by BAO}} + \underbrace{\frac{\partial \mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{\partial \mathbf{y}_i}}_{\text{considered by BAO}} \quad (22) \end{aligned}$$

$$\begin{aligned} \mathbf{y}_i^{k'_i+1} &\leftarrow \mathbf{y}_i^{k'_i} + \alpha \frac{d\mathcal{L}(\mathbf{w}_1^{k_1}, \dots, \mathbf{w}_i^{k_i}, \mathbf{w}_{>i}^K, \mathbf{y}_1^{k'_1}, \dots, \mathbf{y}_i^{k'_i}, \mathbf{y}_{>i}^K)}{d\mathbf{y}_i^{k'_i}}, \\ &\text{where } \mathbf{w}_{>i}^0, \mathbf{y}_{>i}^0 = f(\mathbf{x}, \mathbf{w}_1^{k_1}, \dots, \mathbf{w}_i^{k_i}, \mathbf{y}_1^{k'_1}, \dots, \mathbf{y}_i^{k'_i}) \quad (23) \end{aligned}$$

## A.3 AN EXAMPLE OF EXECUTION OF ALG. 2

In this section, we provide an example of the full `procedure` of execution of Alg. 2 in Fig. 5. The setup is as Fig. 5.(0): we have  $N = 3$  latent  $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$  and gradient ascent step  $K = 2$ , connected by a DAG shown in the figure.

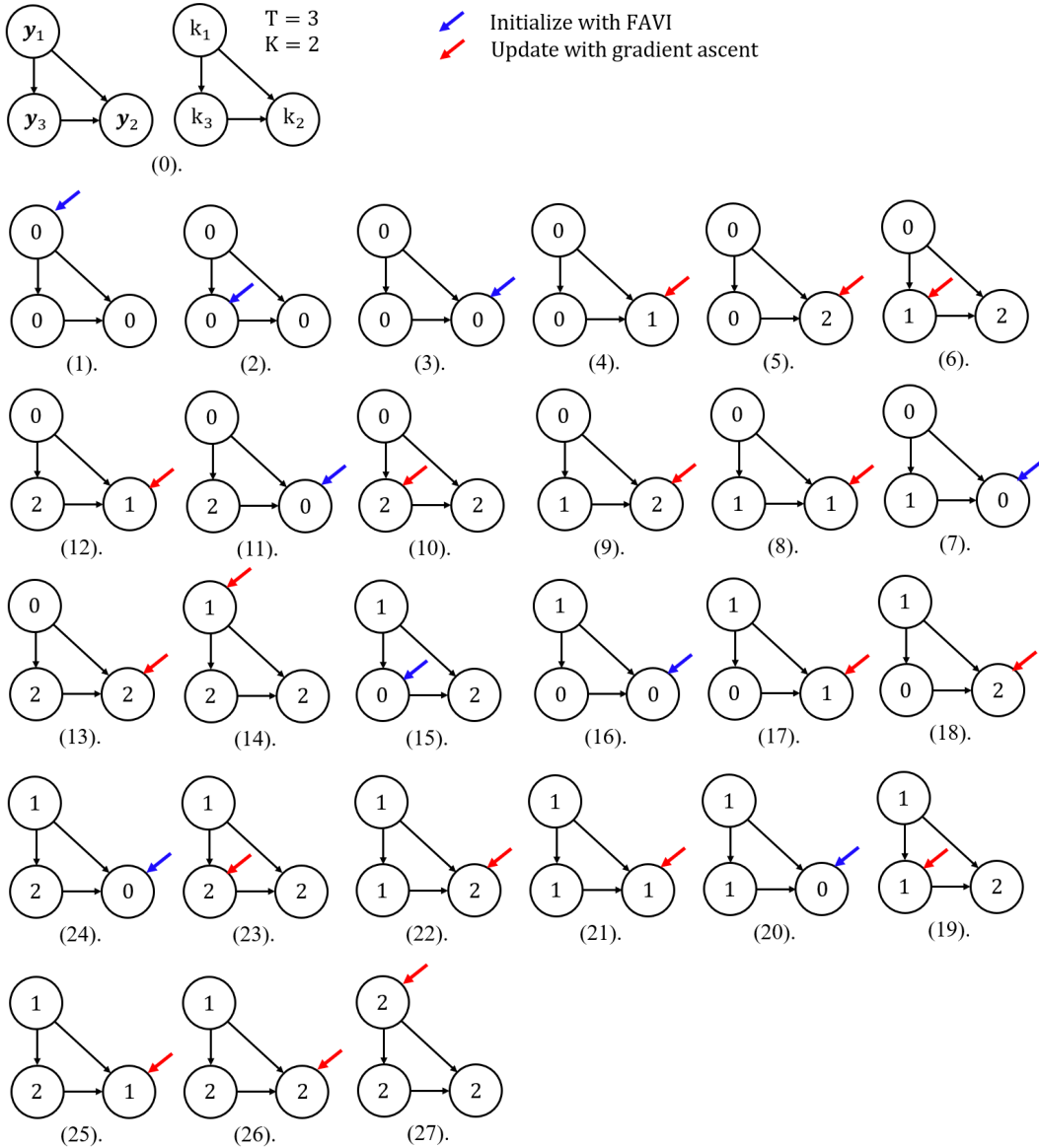


Figure 5: (0). The example setup. (1)- (27). The execution procedure. The number in the circle indicates the gradient step  $k_i$  of each node  $y_i$ . The bold blue/red arrow indicates that the current node is under initialization/gradient ascent.

#### A.4 EXTENDING THE ANALYSIS SEC. 3 TO GENERAL DAG CASE

As the Alg. 2 and Alg. 4 are applicable to general SAVI (Kim et al., 2018; Marino et al., 2018) beyond bit allocation, it is helpful to understand their merit to extend the analysis in Sec. 3 from bit allocation to general DAG scenario. In this section, we consider the same problem setup as Sec. 4.2. Similar to bit allocation case, BAO has the gradient incomplete and gradient value incorrect problem. The gradient incomplete issue is presented as Eq. 24, and gradient value incorrect issue is presented

as Eq. 25.

$$\frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{d\mathbf{a}_i^{k_i}} = \underbrace{\frac{\partial\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{\partial\mathbf{a}_i^{k_i}}}_{\text{considered by BAO}} + \underbrace{\sum_{\mathbf{a}_j \in \mathcal{C}(\mathbf{a}_i)} \left( \sum_{k_j=0}^K \frac{\partial\mathbf{a}_j^{k_j}}{\partial\mathbf{a}_i^{k_i}} \frac{d\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_{j-1}^{k_{j-1}}, \mathbf{a}_{\geq j}^K)}{d\mathbf{a}_j^{k_j}} \right)}_{\text{ignored by BAO}} \quad (24)$$

$$\underbrace{\frac{\partial\mathcal{L}(\mathbf{a}_0^{k_0}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^K)}{\partial\mathbf{a}_i^{k_i}}}_{\text{approximation of BAO in gradient value}} \approx \frac{\partial\mathcal{L}(\mathbf{a}_0^{k_i}, \dots, \mathbf{a}_i^{k_i}, \mathbf{a}_{>i}^{k_i})}{\partial\mathbf{a}_i^{k_i}} \quad (25)$$

## A.5 MORE IMPLEMENTATION DETAILS

In the main text, we use  $\mathbf{y}_i$  as all the latent variable related to residual. In practice, it is divided into  $\mathbf{y}_i, \mathbf{z}_i, \Delta_i^y$ , which refer to the first level latent of residual, second level latent of residual and quantization step size of first level latent of residual respectively. In practice, as BAO (Xu et al., 2022), all of the 3 parts are involved in SAVI jointly. We note that this is not a problem as they fully factorize. And for DVC (Lu et al., 2019),  $\mathbf{w}_i$  indeed represent the latent of motion. As for DVC, the motion has only one level of latent. However for DCVC (Li et al., 2021),  $\mathbf{w}_i$  is divided into  $\mathbf{w}_i, \mathbf{v}_i, \Delta_i^w$ , which refer to the first level latent of motion, second level latent of motion and quantization step size of first level latent of motion respectively. Similar to  $\mathbf{y}_i$ , all of the 3 parts are involved in SAVI jointly, and this is not a problem as they fully factorize.

Following BAO (Xu et al., 2022), we set the target  $\lambda_0 = \{256, 512, 1024, 2048\}$ , which also follows the baselines (Lu et al., 2019; Li et al., 2021). We adopt the official pre-train models for both of the baseline methods (Lu et al., 2019; Li et al., 2021). We do not have a training dataset or implementation details for training amortized encoder / decoder as all the experiments are performed on official pre-trained models. For gradient ascent, we set  $K = 2000$  for the first I frame and  $K = 400$  for all other P frames. On average, the gradient ascent steps for each frame is 920, which is smaller than 2000 in BAO.

## A.6 MORE QUANTITATIVE RESULTS

In this section we present more quantitative results. In Tab. 4 we show the BD-PSNR of our proposed method and other methods as a supplementary to the BD-BR results (Tab. 1). Furthermore, in Fig. 6, we present R-D curve on all classes of HEVC CTC and UVG dataset as a supplementary to the HEVC Class D plot (Fig. 2).

Method	BD-PSNR (dB) $\uparrow$				
	Cls B	Cls C	Cls D	Cls E	UVG
<i>DVC (Lu et al., 2019) as Baseline</i>					
Li et al. (2016) <sup>1</sup>	-0.54	-	-	-0.32	-0.47
Li et al. (2022) <sup>1</sup>	0.19	-	-	0.28	0.14
OEU (Lu et al., 2020)	0.39	0.49	0.83	0.48	0.48
BAO (Xu et al., 2022) <sup>2</sup>	0.87	1.11	1.17	1.35	0.98
Proposed	1.03	1.38	1.67	1.41	1.15
<i>DCVC (Li et al., 2021) as Baseline</i>					
OEU (Lu et al., 2020)	0.30	0.58	0.74	0.29	0.50
BAO (Xu et al., 2022) <sup>2</sup>	0.52	0.76	0.96	0.96	0.69
Proposed	0.91	1.37	1.55	1.58	1.18

Table 4: The BD-PSNR of our approach compared with baselines (w/o bit allocation) and other bit allocation approaches. <sup>1</sup> comes from Li et al. (2022). <sup>2</sup> comes from (Xu et al., 2022).

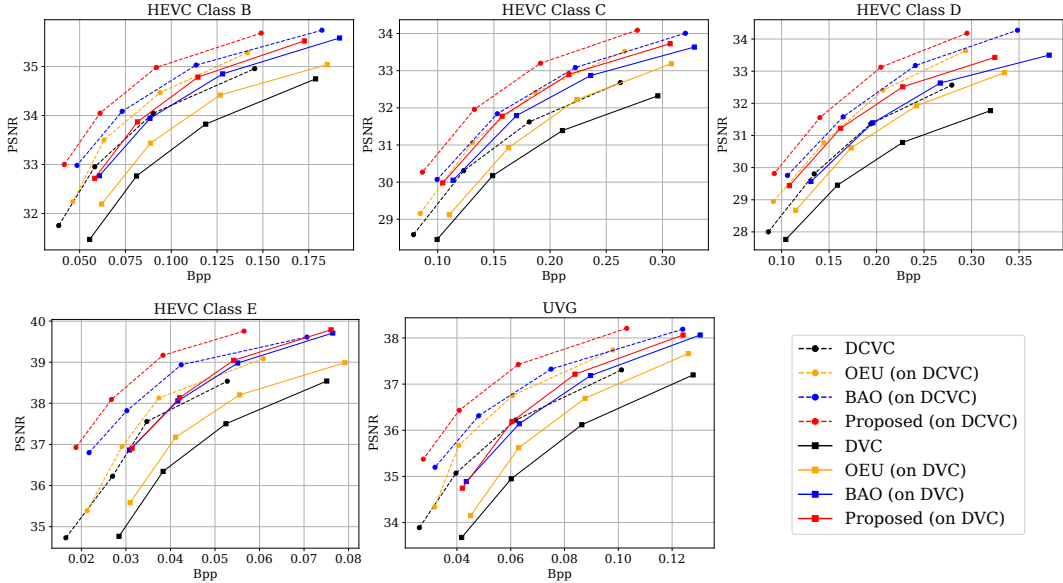


Figure 6: The R-D performance of our approach compared with baselines (w/o bit allocation) and other bit allocation approaches.

### A.7 COMPLEXITY & SCALABILITY

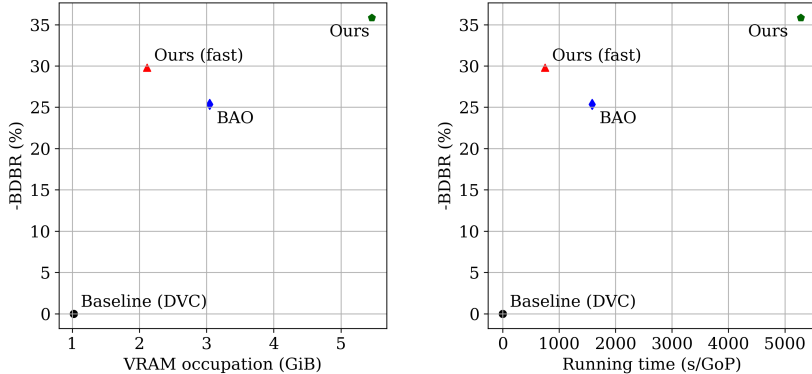


Figure 7: Spatial temporal complexity analysis comparing BAO (Xu et al., 2022), the proposed approach and a fast approximation of the proposed approach. The analysis is done on DVC baseline and HEVC Class D dataset.

We perform additional evaluation to compare the proposed method with BAO (Xu et al., 2022) in terms of temporal complexity and memory cost. The evaluation result can be found in Fig. 7. The general result is that our approach is  $\approx 2.8$  times slower and cost  $\approx 2.0$  times memory than BAO, despite the optimization stepsize is smaller. This extra complexity comes from the cost of sequential optimization of latent. And our current method in its naïve form is slower than BAO while performs better. Jointly consider RD performance, time and memory, our method does not dominate BAO.

However, as our approach enables a sequential style semi-amortized variational inference (SAVI) (Kim et al., 2018; Marino et al., 2018) on latents, there exists a very simple trick to speed it up. Moreover, this trick also resolves the scalability issue. Specifically, to optimize the  $i^{th}$  frame’s latent, we do not compute the R-D cost of all the frames after it as we do now. Instead, we limit the



R-D cost computation to a small fixed size of frames. Formally, we approximate the gradient as:

$$\frac{d\mathcal{L}(\mathbf{w}_{1:T}, \mathbf{y}_{1:T})}{dw_i} \approx \sum_{j=i}^{i+C} \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{dw_i}, \quad \frac{d\mathcal{L}(\mathbf{w}_{1:T}, \mathbf{y}_{1:T})}{dy_i} = \sum_{j=i}^{i+C} \frac{d\mathcal{L}_j(\mathbf{w}_{1:j}, \mathbf{y}_{1:j})}{dy_i} \quad (26)$$

, where  $C$  is a preset constant indicating the number of future frames we included for consideration. With this trick, our algorithm approach cost only  $\approx 50\%$  of time and  $\approx 60\%$  of memory compared with BAO, while remains a superior performance ( $\approx 5\%$  better in BDBR) (Ours (fast) in Fig. 7, the results are based on DVC Class D  $c = 2$ ). With this trick, jointly consider RD performance, time and memory, our approach clearly dominates BAO. Furthermore, with this trick, the scalability issue of

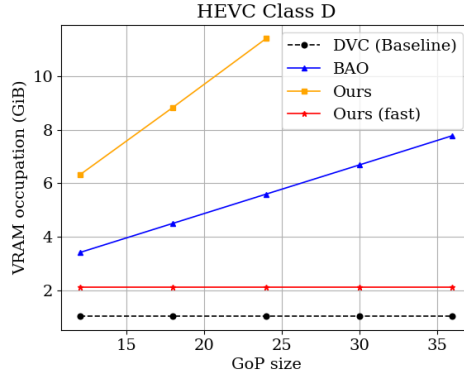


Figure 8: Scalability analysis comparing BAO (Xu et al., 2022), the proposed approach and the fast approximation of the proposed approach. The analysis is done on DVC baseline and HEVC Class D dataset.

our approach is significantly alleviated. As shown in Fig. 8, the memory cost our approach with this trick is constant to GoP size, while that of BAO and our approach without this trick grows linearly with GoP size. This means that with this trick, our approach becomes scalable to any GoP, which is superior than BAO.

#### A.8 IMPACT ON OEU

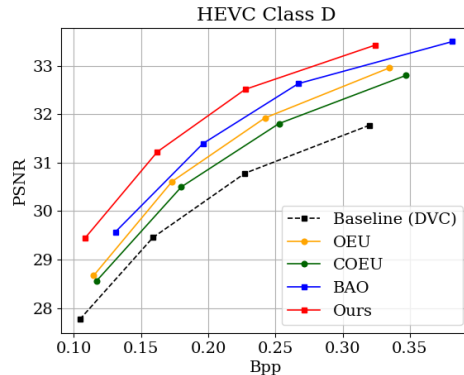


Figure 9: Replace the joint update of OEU (Lu et al., 2020) by our sequential update (Alg. 4). The analysis is done on DVC baseline and HEVC Class D dataset.

Another interesting question to ask is whether the sequential updating algorithm (Alg. 4) benefits the OEU (Lu et al., 2020). Indeed, OEU (Lu et al., 2020) and BAO (Xu et al., 2022) are quite similar at the first glance. However, it is important to note that the theoretical foundation of BAO and this paper is SAVI (Kim et al., 2018; Marino et al., 2018). However, OEU does not fit into SAVI

framework. More specifically, its encoder parameter to be updated does not factorizes as the DAG defined by variational posterior. Thus, applying Alg. 4 is incorrect. To verify this empirically, we change the OEU from BAO’s joint optimization to ours sequential optimization (Alg. 4), and the results show that this change degrades R-D performance (See COEU line in Fig. 9).

#### A.9 MORE ANALYSIS

In this section, we extend the analysis on why the proposed approach works and what is the difference between the proposed approach and BAO (Xu et al., 2022). In the approximate SAVI on DAG latent (Alg. 4), we solve SAVI approximately latent by latent in topological order. For bit allocation of NVC with 10 frames, this topological order is  $\mathbf{y}_0, \mathbf{w}_1, \mathbf{y}_1, \dots, \mathbf{w}_9, \mathbf{y}_9$ , where  $\mathbf{y}_0$  is the latent of I frame,  $\mathbf{w}_i$  is the motion latent of  $i^{th}$  P frame and  $\mathbf{y}_i$  is the residual latent of  $i^{th}$  P frame. In Fig. 10, we show the relationship between R-D cost and the stage of approximate SAVI. We can see that the R-D cost reduces almost consistently with the growing of SAVI stage, which indicates that our approximate SAVI on DAG (Alg. 4) is successful. Specifically, despite our approach is inferior to BAO (Xu et al., 2022) upon the convergence of  $\mathbf{y}_3$ , it attains significant advantage over BAO after  $\mathbf{y}_9$  converges.

In Fig. 11, we compare the distribution of R-D cost, PSNR and Bpp across frame and latent of the baseline DVC (Lu et al., 2019), BAO Xu et al. (2022) and the proposed approach. For R-D cost, it is obvious that our proposed approach’s R-D cost is lower than BAO and baseline, which indicates a better R-D performance. For bpp, it is interesting to observe that despite all three methods have similar bpp of motion related latent  $\mathbf{w}_{1:T}$ , the bpp of residual related latent  $\mathbf{y}_{1:T}$  is quite different. Specifically, BAO increases the bpp of  $\mathbf{y}_{1:T}$  compared with baseline, while our approach decreases the bpp of  $\mathbf{y}_{1:T}$  compared with baseline. This explains why our approach has lower bitrate compared with BAO, and also explains why our approach has significantly less bitrate error. For the PSNR metric, both our approach and BAO significantly improve the baseline. And the difference between proposed approach and BAO is not obvious. We can conclude that the benefits of the proposed approach over BAO comes from the bitrate saving instead of quality enhancing.

#### A.10 QUALITATIVE RESULTS

In Fig. 12, Fig. 13, Fig. 14 and Fig. 15, we present the qualitative result of our approach compared with the baseline approach. We note that compared with the reconstruction frame of baseline approach, the reconstruction frame of our proposed approach preserves significantly more details with lower bitrate, and looks much more similar to the original frame. We intentionally omit the qualitative comparison with BAO (Xu et al., 2022) as it is not quite informative. Specifically, from Fig. 2 we can observe that the PSNR difference of BAO and our approach is very small (within  $\pm 0.1\text{dB}$ ). And our main advantage over BAO comes from bitrate saving instead of quality improvement. Thus the qualitative difference between the proposed method and BAO is likely to fall below just noticeable difference (JND).

#### A.11 MORE DISCUSSION

Other weakness includes scalability. Our method requires jointly considering all the frame inside the GoP, which is impossible when the GoP size is large or when GoP size is unknown for live streaming tasks. Furthermore, currently the gradient ascent step number is merely chosen as an empirical sweet spot between speed and performance. A thorough grid search is desired to better understand its effect on performance.

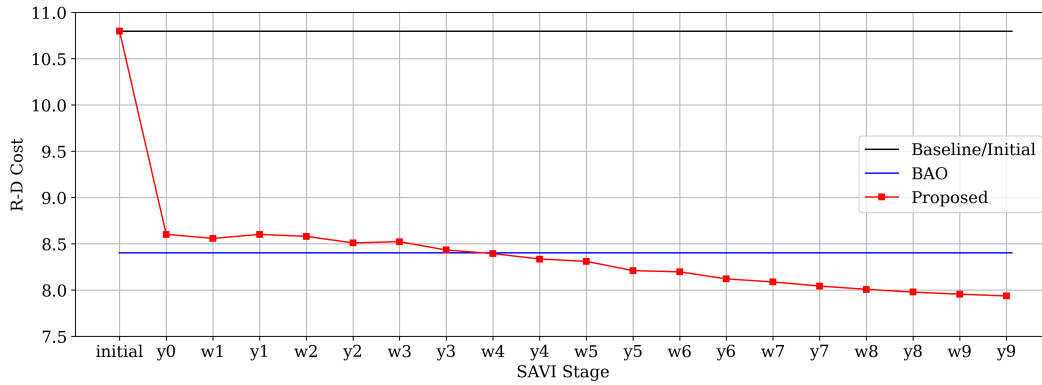


Figure 10: The R-D cost versus SAVI Stage. Experiment is conducted with DVC (Lu et al., 2019) as baseline and *BasketballPass* of HEVC Class D as data.

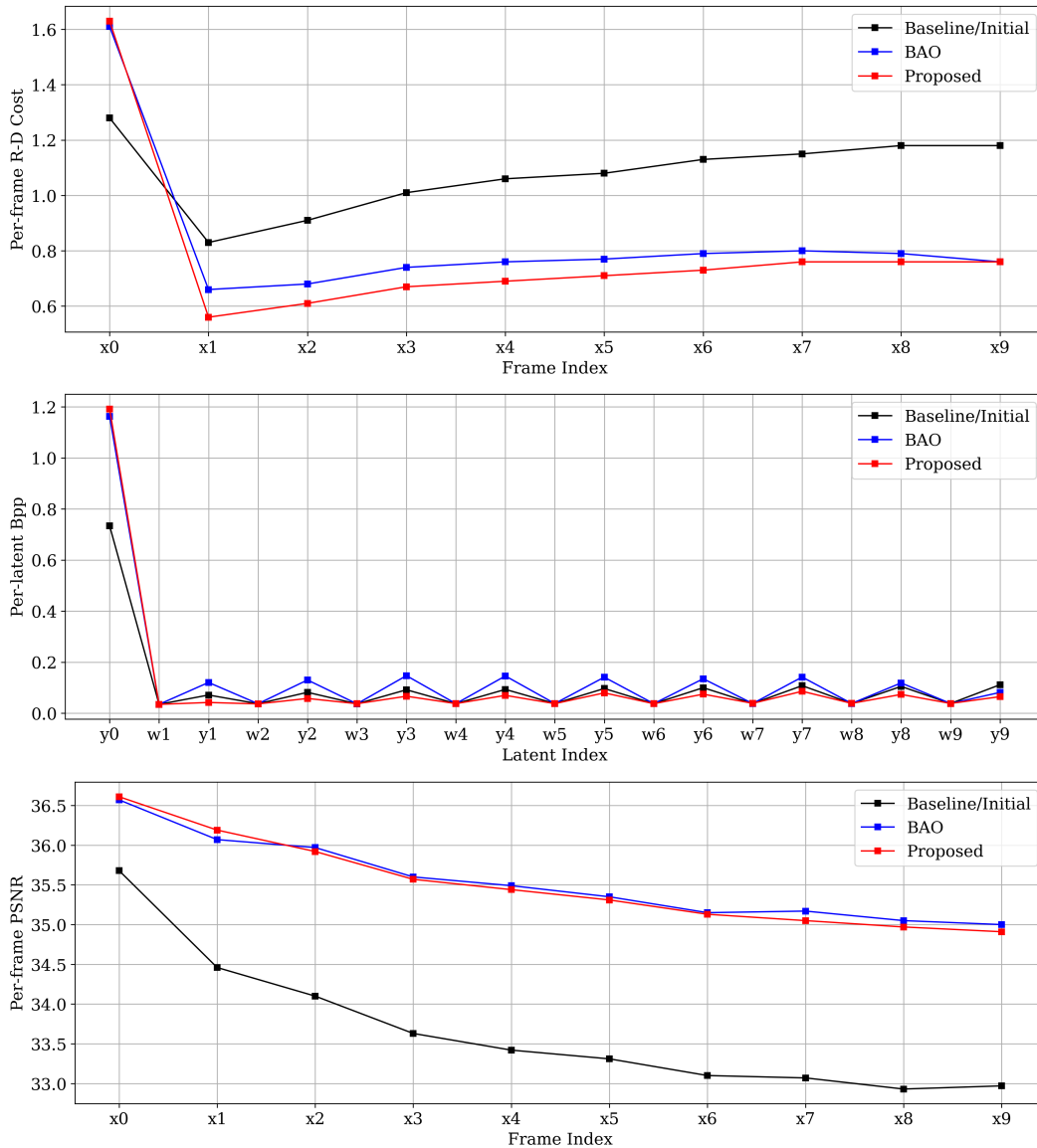


Figure 11: The R-D cost, bpp and PSNR versus frame/latent index. Experiment is conducted with DVC (Lu et al., 2019) as baseline and *BasketballPass* of HEVC Class D as data.



Figure 12: Qualitative results using *BasketballPass* of HEVC Class D. *Top*. Original frame. *Middle*. Baseline codec (DVC)’s reconstruction result with  $\text{bpp} = 0.148$  and  $\text{PSNR} = 33.06\text{dB}$ . *Bottom*. Proposed method’s reconstruction result with  $\text{bpp} = 0.103$  and  $\text{PSNR} = 34.91\text{dB}$ .



Figure 13: Qualitative results using *BlowingBubbles* of HEVC Class D. *Top*. Original frame. *Middle*. Baseline codec (DVC)’s reconstruction result with  $\text{bpp} = 0.206$  and  $\text{PSNR} = 30.71\text{dB}$ . *Bottom*. Proposed method’s reconstruction result with  $\text{bpp} = 0.129$  and  $\text{PSNR} = 32.34\text{dB}$ .



Figure 14: Qualitative results using *BQSquare* of HEVC Class D. *Top*. Original frame. *Middle*. Baseline codec (DVC)’s reconstruction result with  $\text{bpp} = 0.232$  and  $\text{PSNR} = 28.72\text{dB}$ . *Bottom*. Proposed method’s reconstruction result with  $\text{bpp} = 0.128$  and  $\text{PSNR} = 30.87\text{dB}$ .



Figure 15: Qualitative results using *RaceHorses* of HEVC Class D. *Top*. Original frame. *Middle*. Baseline codec (DVC)’s reconstruction result with  $\text{bpp} = 0.448$  and  $\text{PSNR} = 30.48\text{dB}$ . *Bottom*. Proposed method’s reconstruction result with  $\text{bpp} = 0.379$  and  $\text{PSNR} = 31.92\text{dB}$ .