
Versatile Learned Video Compression

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learned video compression methods have demonstrated great promise in catching
2 up with traditional video codecs in their rate-distortion (R-D) performance. How-
3 ever, existing learned video compression schemes are limited by the binding of
4 the prediction mode and the fixed network framework. They are unable to support
5 various inter prediction modes and thus inapplicable for various scenarios. In this
6 paper, to break this limitation, we propose a versatile learned video compression
7 (VLVC) framework that uses one model to support all possible prediction modes.
8 Specifically, to realize versatile compression, we first build a motion compensation
9 module that applies multiple 3D motion vector fields (*i.e.*, voxel flows) for weighted
10 trilinear warping in spatial-temporal space. The voxel flows convey the information
11 of temporal reference position that helps to decouple inter prediction modes away
12 from framework designing. Secondly, in case of multiple-reference-frame predic-
13 tion, we apply a flow prediction module to predict accurate motion trajectories with
14 a unified polynomial function. We show that the flow prediction module can largely
15 reduce the transmission cost of voxel flows. Experimental results demonstrate that
16 our proposed VLVC not only supports versatile compression in various settings but
17 also achieves comparable R-D performance with the latest Versatile Video Coding
18 (VVC) standard in terms of MS-SSIM.

19 1 Introduction

20 Video occupies more than 80% of network traffic and the amount of video data is increasing rapidly
21 [1]. Thus, the storage and transmission of video become more challenging. A series of hybrid video
22 coding standards have been proposed, such as AVC/H.264 [2], HEVC/H.265 [3] and the latest video
23 coding standard VVC/H.266 [4]. These traditional standards are manually designed, evolving for
24 decades. However, the development within the hybrid coding framework is gradually saturated.
25 Recently, the performance of video compression is mainly improved by designing more complex
26 prediction modes, leading to increased coding complexity.

27 Deep neural networks are currently promoting the development of data compression. Despite the
28 remarkable progress on the field of learned image compression [5–10], the area of learned video
29 compression is still in early stages. Existing methods for learned video compression can be grouped
30 into three categories, including frame interpolation-based methods [11, 12], 3D autoencoder-based
31 methods [13, 14], and predictive coding methods with optical flows [15–17]. So far, among them,
32 video compression with optical flow presents the best performance [18], where optical flow represents
33 pixel-wise motion vector (MV) fields utilized for inter frame prediction. In this paper, we also focus
34 on this predictive coding architecture. Previous works with optical flow are proposed to support
35 specific prediction mode, including unidirectional or bidirectional, single or multiple frame prediction.
36 They are too cumbersome to support versatile compression in various settings since they bind the inter
37 prediction mode with the fixed network framework. It is important to design a more flexible model to
38 handle all possible settings like traditional codecs. In this paper, we propose a versatile learned video

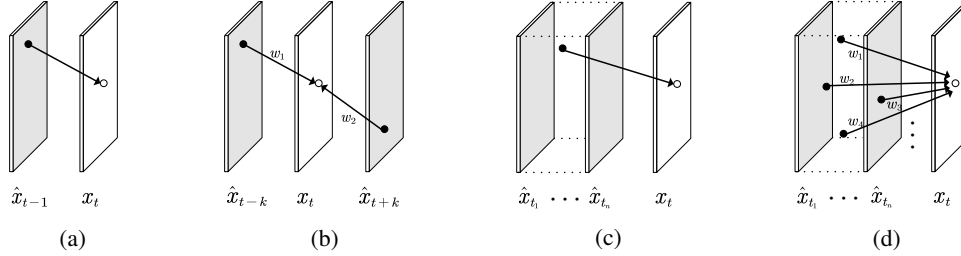


Figure 1: Different motion compensation (inter frame prediction) methods. (a) Unidirectional prediction with 2D optical flow [15]. (b) Bidirectional prediction with two optical flows and weight coefficients [12]. (c) Prediction with a single voxel flow, freely sampling the reference frames in space-time. (d) Prediction with multiple voxel flows via weighted trilinear warping.

39 compression (VLVC) framework that achieves coding flexibility as well as compression performance.
 40 A voxel flow based motion compensation module is adopted for higher flexibility, which is then
 41 extended into multiple voxel flows to perform weighted trilinear warping. In addition, in case of
 42 multiple-reference-frame prediction, a polynomial motion trajectories based flow prediction module
 43 is designed for better compression performance. Our motivations are described as follows.

44 **Motion compensation with multiple voxel flows.** Previous works such as [15] apply 2D optical
 45 flows for low-delay prediction using single reference frame (unidirectional prediction, see Fig. 1a).
 46 For the practical random access scenario, bidirectional reference frames are available for more
 47 accurate frame interpolation [12] (Fig. 1b). However, the reference positions in these works are
 48 determined by pre-defined prediction modes. They cannot adapt to various inter prediction modes
 49 where reference positions are various. In this paper, we apply 3D voxel flows to describe not only
 50 the spatial MVs, but also the information of temporal reference positions (Fig. 1c). We perform
 51 voxel flow based motion compensation via trilinear warping, which is applicable to single or multiple,
 52 unidirectional or bidirectional reference frames. Unlike [16] that adopts scale space flow with trilinear
 53 warping, we apply voxel flows for inter prediction in spatial-temporal space, which naturally renders
 54 our model more robust to different coding scenarios. Furthermore, beyond using single MV in every
 55 position of the current frame, we propose to use multiple voxel flows to describe multiple possible
 56 reference relationships (Fig. 1d). Then the target pixel is synthesized by weighted fusing of the
 57 warping results. We show that without increasing the coding cost of motion information, the motion
 58 compensation is thus more accurate, yielding less residuals and more efficient compression.

59 **Flow prediction with polynomial motion trajectories.** Exploiting multiple reference frames usu-
 60 ally achieves better compression performance since more reference information is provided. A
 61 versatile learned video compression model should cover this multi-reference case. While previous
 62 work [17] designs a complex flow prediction network to reduce the redundancies of 2D MV fields, the
 63 number and structure of reference frames are inherent and fixed within the framework. In this paper,
 64 we design a more intelligent method for flow prediction, *i.e.*, modelling the prediction modes with
 65 polynomial coefficients. We formulate different motion trajectories in a time interval by a unified
 66 polynomial function. The polynomial coefficients are solved by establishing a multivariate equation
 67 (see Section 3.2). Since this polynomial function models the accurate motion trajectories, it serves as
 68 a basic discipline that constrains the predicted motion to be reasonable. We show the transmission
 69 cost of voxel flows is reduced obviously with the help of additional motion trajectory information.

70 Thanks to the above two technical contributions, our proposed VLVC is not only applicable for various
 71 practical compression scenarios with different inter prediction modes, but also delivers impressive
 72 R-D performance on standard test sequences. Extensive experimental results demonstrate that our
 73 method is the first one to achieve comparable performance with VVC in terms of MS-SSIM in both
 74 low delay and random access configurations. Comprehensive ablation studies and discussions are
 75 provided to verify the effectiveness of our method.

76 2 Related Work

77 **Learned Image Compression** Recent advances in learned image compression [5–7], have shown
 78 the great success of nonlinear transform coding. Many existing methods are built upon hyperprior-

79 based coding framework [6], which are improved with more efficient entropy models [7, 8], variable-
 80 rate compression [19] and more effective quantization [9, 10]. While the widely used autoregressive
 81 entropy models provide significant performance gain in image coding, the high decoding complexity
 82 is not suitable for practical video compression. We thus employ the hyperprior model [6] without
 83 context models in our video compression framework.

84 **Learned Video Compression** Existing approaches [11–15, 17, 18, 20–24] can be roughly divided
 85 into three categories: frame interpolation-based methods [11, 12, 24], 3D autoencoder-based methods
 86 [13, 14], and predictive coding methods with optical flows [15–17]. Currently, researchers are more
 87 interested in the latter two methods. Although 3D autoencoder-based methods requires less time
 88 complexity, they barely achieve comparable performance with x265 in MS-SSIM [14]. Meanwhile,
 89 predictive coding methods with optical flows have outperformed HM in terms of PSNR [18].

90 Predictive-based video compression approaches [15, 20–24, 11, 12, 17] sequentially perform motion
 91 estimation, motion compression, motion compensation and residual compression. Chen et al. [24]
 92 first propose to predict block of pixels using learned neural network (DNN), and the residual is
 93 compressed by a RNN-based autoencoder. Wu et al. [11] propose a interpolation-based approach
 94 using traditional MVs. Lu et al. [15] propose an fully end-to-end trainable framework, where all key
 95 components in the classical video codec are implemented with neural networks. Rippel et al. [21]
 96 jointly compress the motion and residual information, and propose a latent state to memorize the
 97 information from the past. Djelouah et al. [12] perform interpolation by the decoded optical flow
 98 and blending coefficients. They reuse the same autoencoder of I-frame compression and directly
 99 quantize the corresponding latent space residual. Liu et al. [23] combine the optical flow estimation
 100 and motion compression networks into one-stage, and remove the redundancy of quantized flow
 101 representations using joint spatial-temporal priors. Yang et al. [20] propose a video compression
 102 framework with three hierarchical quality layers and recurrent enhancement. In [17], multiple frames
 103 motion prediction are introduced into the P-frame coding. Lu et al. [22] propose an content adaptive
 104 and error propagation aware method to reduce error accumulation and achieve adaptive coding.
 105 Agustsson et al. [16] replace the bilinear warping operation with scale-space flow which allows
 106 the model adaptively blur the reference content for better warping results. However, most existing
 107 methods are designed for particular prediction modes, resulting in inflexibility for different scenarios.

108 **Video Interpolation** The task of video interpolation is closely related to video compression. One
 109 pioneering work [25] proposes to use deep voxel flow to synthesize new video frames. Some works
 110 of video interpolation [26–28] directly generate the spatially-adaptive convolutional kernels for each
 111 motion vectors by neural networks. Most recently, [29, 30] proposed to relax the kernel shape,
 112 allowing the models to freely select multiple sampling points in space or space-time. In this paper,
 113 our employed multiple voxel flows is motivated by the accurate interpolation result in [30].

114 3 Versatile Learned Video Compression

115 To compress video, the original video sequence is first divided into groups of pictures (gop). Let
 116 $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ denote the frames in one gop unit where the gop size is T . To take advantage
 117 of previous decoded frames, our model predicts the current frame \mathbf{x}_t from n reference frame(s),
 118 *i.e.*, the lossy reconstruction results compared to the original frames. Here, we denote the reference
 119 frames as $\{\hat{\mathbf{x}}_{t_1}, \hat{\mathbf{x}}_{t_2}, \dots, \hat{\mathbf{x}}_{t_n}\}$, where $\{t_1, t_2, \dots, t_n\}$ is the index of temporal reference position. If
 120 multiple frames are taken as the reference (*i.e.*, $n > 1$), the reference frames are divided into two
 121 groups: one is used only for flow prediction, and the other is used for both flow prediction and motion
 122 compensation. In other words, the reference involved for motion compensation is only a sub-set of
 123 $\{\hat{\mathbf{x}}_{t_1}, \hat{\mathbf{x}}_{t_2}, \dots, \hat{\mathbf{x}}_{t_n}\}$, which could be concatenated into a volume denoted by $\hat{\mathbf{X}}_t$. If only one reference
 124 frame is available, the volume for warping $\hat{\mathbf{X}}_t = \{\hat{\mathbf{x}}_{t-1}\}$.

125 An overview of our video compression framework is shown in Fig. 2. Previous work [16] demon-
 126 strates that an implicit flow encoder can outperform a pre-trained optical flow network and simplify
 127 the network structures simultaneously. In our paper, we also abandon the use of a pre-trained optical
 128 flow network in motion encoder. The motion encoder and decoder are similar to image compression
 129 network [5]. While the work of [16] sends current frame and previous reconstruction into motion
 130 encoder, we make some modifications on the input of motion encoder. Specifically, in our framework,
 131 the motion encoder is fed with the current frame \mathbf{x}_t concatenated with predicted frames (represented

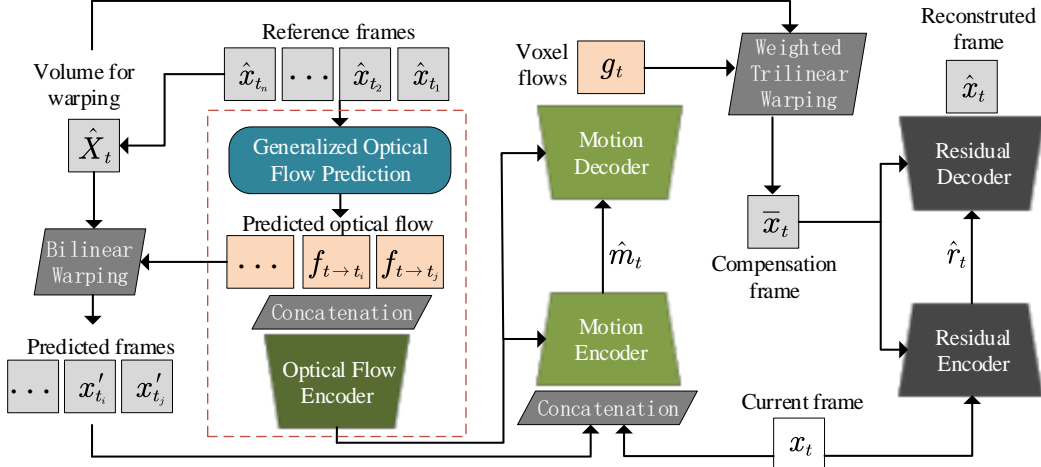


Figure 2: Overview of our inter-frame coding framework.

132 as x'_t in Fig. 2). Here, the predicted frame x'_t is an estimation of current frame x_t . The flow prediction
 133 module will predict 2D optical flows $f_{t \rightarrow t_i}$ to warp corresponding reference frames \hat{x}_{t_i} , each of
 134 which will generate a predicted frame x'_{t_i} . The predicted frames reveal how much information the
 135 decoder knows about current frame.

136 On the decoder side, a motion decoder will generate voxel flows which is used for motion compensa-
 137 tion via trilinear warping (details are explained in Section 3.1). In addition, when multiple reference
 138 frames are available, the flow prediction module is turned on (the red dashed box). As shown in
 139 Fig. 2, the flow prediction module is auxiliary for motion encoding and decoding, which reduces
 140 the transmission cost of the quantized motion latent \hat{m}_t . The specific mechanism of flow prediction
 141 can be found in Section 3.2. After motion compensation, we obtain a prediction of current frame as
 142 \bar{x}_t . The residual encoder and decoder are then used to compress the remaining residuals between the
 143 original frame x_t and the predicted frame \bar{x}_t , yielding the final reconstructed frame \hat{x}_t and quantized
 144 latent \hat{r}_t .

145 3.1 Motion compensation with multiple voxel flows

146 Voxel flow [25] is a per-pixel 3-D motion vector that describes relationships in spatial-temporal
 147 domain. Compared to 2D optical flow, voxel flow can inherently allow the codec to be aware of
 148 the sampling positions in the temporal dimension for various prediction modes. Given arbitrary
 149 number of reference frames, the model is expected to select the optimal reference frame for better
 150 reconstructing the current frame to be compressed. Such a 3-D motion descriptor helps to build a
 151 prediction-model-agnostic video compression framework, *i.e.*, versatile learned video codec.

152 In addition, single flow field is hard to represent complex motion (e.g. blurry motion), which may
 153 result in inaccurate motion compensation or high coding cost of motions. When reconstructing a
 154 local region, its reference information may not come from only one source. Considering a practical
 155 scene where multiple objects of the same types appear at the same time, more than one areas could be
 156 referred for reconstructing the local region. Thereby, in this work, we further propose to use multiple
 157 voxel flows to perform weighted trilinear warping by sampling in \mathbf{X}_t for multiple times. We remind
 158 our readers that \mathbf{X}_t is a volume consisting of some reference frames. Assume the dimension of
 159 \mathbf{X}_t is $D \times H \times W$ (usually reshaped into $H \times W \times D$ for warping), where D is the number of
 160 reference frames used for motion compensation. the motion decoder will generate multiple voxel
 161 flows by outputting a $(4M) \times H \times W$ tensor. Here, M refers to the number of flows. Therefore,
 162 every voxel flow is a 4-channel field that describes the 3-channel voxel flow $g^i = (g_x^i, g_y^i, g_z^i)$ with a
 163 corresponding weight channel g_w^i . Here, i ($1 \leq i \leq M$) is the index of voxel flow. To synthesize the
 164 target pixels in current frame, the weights g_w^i are normalized by a softmax function across M voxel

165 flows. We finally obtain the target pixel $\bar{x}[x, y]$ in spatial location $[x, y]$ by calculating the weighted
 166 sum of sampling results, formulated as:

$$\bar{x}[x, y] = \sum_{i=1}^M g_w^i(x, y) \mathbf{X}_t[x + g_x^i(x, y), y + g_y^i(x, y), g_z^i(x, y)]. \quad (1)$$

167 We experimentally find that compared with single voxel flow, the transmission cost of multiple voxel
 168 flows does not increase largely. The model is able to assign appropriate number of flows under the
 169 rate-distortion optimization goal. In other word, the model is optimized to avoid the transmission of
 170 unnecessary flows. Meanwhile, due to more accurate inter frame prediction, the transmission cost of
 171 residuals decreases obviously by using multiple voxel flows for weighted warping.

172 3.2 Generalized optical flow prediction

173 In our proposed VLVC framework, as illustrated in Fig. 2, we compress the spatial-temporal motion
 174 information via motion encoder and decoder. The concatenation of the predicted frames (*i.e.*, bilinear
 175 warping results using the predicted optical flow) and the current frame are fed into the motion encoder.
 176 The 3-D motion descriptor, *i.e.*, the voxel flows, are then decoded by the motion decoder given
 177 the quantized motion latent and the feature of predicted optical flow. In this process, the predicted
 178 optical flow reduces the spatial displacement need to be encoded, and also serves as the conditions to
 179 promote the generation of voxel flows. Thus, the optical flow prediction is clearly of great importance
 180 to reduce the redundancies of consecutive voxel flows in case of using multiple reference frames.

181 Specifically, there are two optical flow describe the motion between the reference frame \hat{x}_{t_j} and
 182 the target frame x_t : $f_{t_j \rightarrow t}$ and $f_{t \rightarrow t_j}$. The flow $f_{t \rightarrow t_j}$ describe the motion of each pixel from
 183 x_t , and therefore we can sample \hat{x}_{t_j} for each pixel in the target frame x_t via bilinear (backward)
 184 warping. However, $f_{t \rightarrow t_j}$ is unknown at decoder side because the pixels of target frame is unavailable.
 185 Fortunately, the pixels from reference frames are known at both encoder and decoder. We can first
 186 estimate the optical flow of pixels from a reference frame $bol \hat{x}_{t_j}$ to other reference frames, and then
 187 predict the flow $f_{t_j \rightarrow t}$. While we obtain $f_{t_j \rightarrow t}$, it cannot be directly used for motion compensation
 188 with bilinear warping.

189 Recently work [31] for video interpolation proposed a forward warping method to interpolate the
 190 target frame x_t by directly using the flow $f_{t_j \rightarrow t}$. For video compression, we aim to predict a
 191 approximation of the flow $f_{t \rightarrow t_j}$ to reduce the redundancies of the proposed voxel flows for better
 192 rate-distortion performance. We therefore employ the forward warping method [31] (named softmax
 193 splatting) to project the flow $f_{t_j \rightarrow t}$ to $f_{t \rightarrow t_j}$, which is a kind of flow reversal methods similar to [32].
 194 In the following part, we will describe a novel polynomial motion modeling method to predict $f_{t_j \rightarrow t}$
 195 given arbitrary reference frames and any target time stamp t . And a flow reversal layer based on
 196 softmax splatting is introduced for the final flow prediction.

197 **Polynomial motion modeling** For each pixel at t_j , we model the motion $f_{t_j \rightarrow t}$ by the k -order
 198 ($k < n$) polynomial functions:

$$f_{t_j \rightarrow t} = a_1 \times (t - t_j) + a_2 \times (t - t_j)^2 + \dots + a_k \times (t - t_j)^k, \quad (2)$$

199 where a_0, a_1, \dots, a_k are the polynomial coefficients. To solve the coefficients, we set t equals to the
 200 top- k nearest time stamp $\{t_{j_i}\}_{i=1}^k$ around t_j within the set of reference time stamp. Then we can
 201 obtain the following equation:

$$\begin{pmatrix} a_1 \\ a_2 \\ \dots \\ a_k \end{pmatrix} = \begin{pmatrix} (t_{j_1} - t_j) & (t_{j_1} - t_j)^2 & \dots & (t_{j_1} - t_j)^k \\ (t_{j_2} - t_j) & (t_{j_2} - t_j)^2 & \dots & (t_{j_2} - t_j)^k \\ \dots & \dots & \dots & \dots \\ (t_{j_k} - t_j) & (t_{j_k} - t_j)^2 & \dots & (t_{j_k} - t_j)^k \end{pmatrix}^{-1} \begin{pmatrix} f_{t_j \rightarrow t_{j_1}} \\ f_{t_j \rightarrow t_{j_2}} \\ \dots \\ f_{t_j \rightarrow t_{j_k}} \end{pmatrix} \quad (3)$$

202 where $f_{t_j \rightarrow t_{j_1}}, f_{t_j \rightarrow t_{j_2}}, \dots, f_{t_j \rightarrow t_{j_k}}$ can be obtained using off-the-shelf flow estimation network.
 203 Then we can derive the polynomial coefficients and apply them to Eq. (3) predict the forward flow
 204 from t_j to any time stamp t .

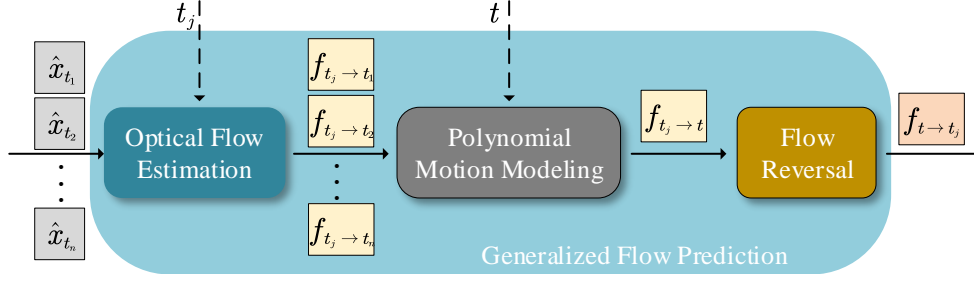


Figure 3: Generalized flow prediction module.

205 **Flow reversal via softmax splatting** While the forward flow $f_{t_j \rightarrow t}$ is predicted by the polynomial
 206 functions, it cannot be directly used for motion compensation. Therefore, we introduce a flow reversal
 207 layer to forward warping $-f_{t_j \rightarrow t}$ by softmax splatting [31]:

$$f_{t \rightarrow t_j} = \frac{\vec{\sum} (\exp(\mathbf{Z}) \cdot (-f_{t_j \rightarrow t}), f_{t_j \rightarrow t})}{\vec{\sum} (\exp(\mathbf{Z}), f_{t_j \rightarrow t})}, \quad (4)$$

208 where $\vec{\sum}$ is the summation splatting defined in [31], and \mathbf{Z} is an importance mask generated from a
 209 small network q as:

$$\mathbf{Z} = q(\hat{\mathbf{x}}_{t_j}, -\frac{1}{k} \sum_{i=1}^k \|\hat{\mathbf{x}}_{t_j} - \overleftarrow{w}(\hat{\mathbf{x}}_{t_i}, f_{t_j \rightarrow t_i})\|_1), \quad (5)$$

210 where \overleftarrow{w} is the bilinear backward warping operator.

211 3.3 Loss function

212 In previous works, the reference frames are determined according to pre-defined prediction modes.
 213 For example, the work of [17] applies four unidirectional reference frames, where the reference set
 214 is $\{\hat{\mathbf{x}}_{t-4}, \hat{\mathbf{x}}_{t-3}, \hat{\mathbf{x}}_{t-2}, \hat{\mathbf{x}}_{t-1}\}$. The work of [12] applies $\{(\hat{\mathbf{x}}_{t-1}, \hat{\mathbf{x}}_{t+1}), (\hat{\mathbf{x}}_{t-2}, \hat{\mathbf{x}}_{t+2}), (\hat{\mathbf{x}}_{t-3}, \hat{\mathbf{x}}_{t+3})\}$
 215 as the reference set for bilinear prediction. In this paper, to optimize a versatile video compression
 216 model, the model will have access to various reference structures during training to adapt to different
 217 prediction modes. Therefore, we apply the loss function to cover all the frames in the entire GOP as:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T [R_t(\hat{\mathbf{m}}_t, \hat{\mathbf{r}}_t | \hat{\mathbf{x}}_{t_i}, \dots, \hat{\mathbf{x}}_{t_j}) + \lambda \cdot \mathcal{D}(\mathbf{x}_t, \hat{\mathbf{x}}_t | \hat{\mathbf{x}}_{t_i}, \dots, \hat{\mathbf{x}}_{t_j})]. \quad (6)$$

218 Here, T is the GOP size during training. The maximum value of T is seven in our experiments since a
 219 7-frame GOP can cover most prediction modes. $\{\hat{\mathbf{x}}_{t_i}, \dots, \hat{\mathbf{x}}_{t_j}\}$ represents different reference set that
 220 may vary in different mini-batches. $R_t(\hat{\mathbf{m}}_t, \hat{\mathbf{r}}_t)$ is the rate of motion and residual. For simplicity, we
 221 omit the process of intra frame compression (at $t = 1$) in this loss function.

222 4 Experiments

223 4.1 Experimental setup

224 **Model details** The motion/residual compression modules are two auto-encoder style networks,
 225 where the bit-rate are estimated by the factorized and hyperprior entropy model [6, 7], respectively. We
 226 employ the off-the-shelf PWC-net [33] as the optical flow estimation network only in our generalized
 227 flow prediction module. We employ feature residual coding [34] instead of pixel residual coding for
 228 better performance. Detailed architecture can be found in supplementary.

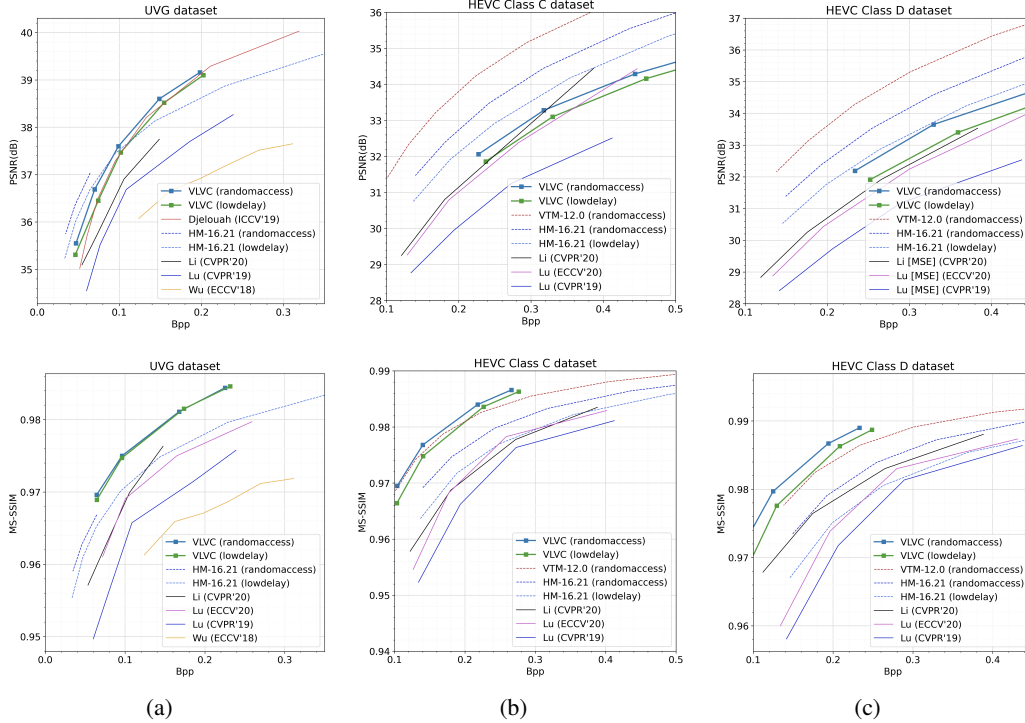


Figure 4: Rate-distortion Performance.

229 **Training sets** The models were trained on the Vimeo-90k septuplets dataset [35] which consists of
 230 89800 video clips with diverse content. The video clips are randomly cropped to 128×128 or $256 \times$
 231 256 pixel for training.

232 **Testing sets** The HEVC common test sequences [3] and the UVG dataset [36] are used for evalua-
 233 tion. The HEVC Classes B,C,D and E contain 16 videos with different resolution and content. The
 234 UVG dataset contains seven 1080p HD video sequences with 3900 frames in total.

235 **Implementation details** We optimize four models for MSE and four models for MS-SSIM [37].
 236 The video clip length T is set to 7 for training. We use the Adam optimizer [38] with batch size of 8
 237 and a initial learning rate of 5×10^{-5} . It is difficult to stably train the whole models from scratch. We
 238 first separately pre-train the intra-frame coding models and inter-frame coding models for MSE, with
 239 128×128 video crops and 1,200,000 training steps. Then we jointly optimize both the models with
 240 the gop loss Eq. (6) for 100,000 steps using different metrics and λ values. Finally, we fine-tuning all
 241 the models for 20,000 with a crop size of 256×256 and a reduced learning rate of 1×10^{-5}

242 **Evaluation Setting** We measure the quality of reconstructed frames using PSNR and MS-
 243 SSIM [37] in the RGB colorspace. The bits per pixel (bpp) is used to measure the average number
 244 of bits. We compare our method with the traditional video coding standards H.265/HEVC and
 245 H.266/VVC, as well as the state-of-art learning based methods including [15, 22, 11, 12, 17].

246 Recent works for learned video compression usually evaluate H.265 by using FFmpeg, with perform-
 247 ance is much lower than official implementation. In this paper, we evaluate H.265 and H.266 by
 248 using the implementation of the standard reference software HM 16.21[39] and VTM 12.0[40], respec-
 249 tively. We use the default low delay and random access configuration, and modify the gop structure
 250 and key frame interval for fair comparison. Detailed configuration can be found in supplementary.

251 4.2 Performance

252 We evaluate our model with the state-of-the-art learned video compression approaches, including the
 253 P-frame based methods of [15, 22, 23, 17], the interpolation based methods of [11, 12]. As shown

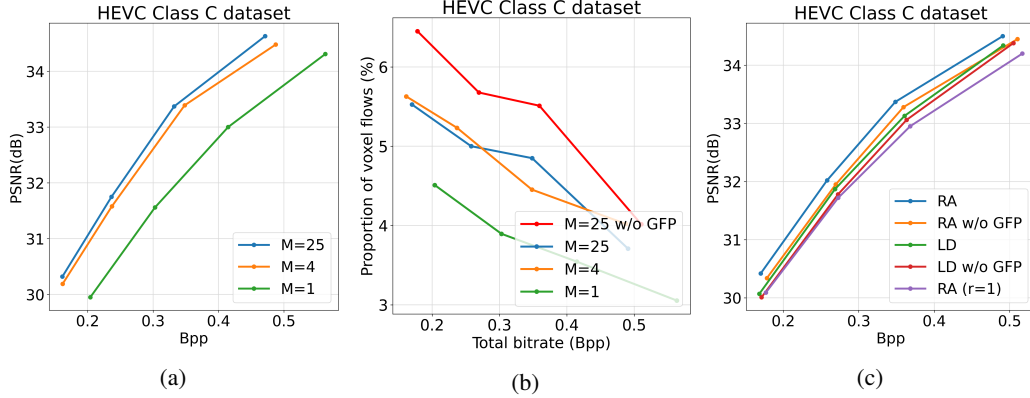


Figure 5: (a) Ablation on the number of voxel flows. (b) The Proportion of voxel flows in total bitrate. (c) Ablation on different coding configurations.

254 in Fig. 4, it can be observed that our proposed method significantly outperforms existing learned
 255 video compression methods in both PSNR and MS-SSIM. Note that the *VLVC (randomaccess)* and
 256 *VLVC (lowdelay)* are two different configurations from the same models. Besides, our model is the
 257 first end-to-end learned video compression method that achieves comparable R-D performance with
 258 H.266 in terms of MS-SSIM.

259 4.3 Ablation Study and Analysis

260 All the models reported in ablation studies are trained for MSE using 128×128 video clips. More
 261 ablation study results and visual results can be found in the supplementary.

262 **The effect of the voxel flow number** As shown in Fig. 5a, the number M of voxel flows signifi-
 263 cantly influence the overall rate-distortion performance. More voxel flows provide more possible
 264 sampling location for accurate motion compensation. Our proposed weighted trilinear warping with
 265 multiple voxel flows achieves about 1dB gain compared with the conventional trilinear warping with
 266 single voxel flow. Note that the performance gain is nearly saturated for $M = 25$, which is used as
 267 the default value in our models.

268 We also investigate the additional bitrate cost of multiple voxel flows. As shown in Fig. 5b, the
 269 proportion of multiple voxel flows in the total bitrate of video coding increases about $\frac{1}{3}$ at the same
 270 bitrate. In other words, our model can learn to improve the overall compression performance by
 271 transmitting a proper amount of additional motion information, which is represented as voxel flows.

272 **Versatile coding configurations** The proposed methods can deal with a various set of prediction
 273 modes. To evaluate the effectiveness of coding flexibility as well as the effectiveness of the proposed
 274 generalized flow prediction module, we simply change the input coding configurations of the same
 275 trained model at different bitrate points. Random access and low delay coding settings are denoted
 276 as *RA* and *LD*, respectively. As shown in Fig. 5c, the coding mode *RA* with bidirectional reference
 277 frames achieves a compression gain of about 0.4dB, compared with the unidirectional coding modes
 278 *LD*. Furthermore, the performance dropped about 0.1dB~0.3dB when we turn off the generalized
 279 flow prediction module for different coding settings, noted as *w/o GFP*. We also illustrate the bitrate
 280 reduction of the voxel flows shown in Fig. 5b, where $M=25$ reduce the bitrate of voxel flows about $\frac{1}{6}$
 281 compared to $M=25$ *w/o GFP*. Finally, we change the number of the reference frames for warping,
 282 which is set to 2 as default. We reduce the number to 1 in random access mode, noted as *RA (r=1)*,
 283 which performance is even worse than low delay setting.

284 **Visualization of voxel flows** The proposed voxel flows contain multiple 3-channel voxel flows
 285 $\{(g_x^i, g_y^i, g_z^i)\}_{i=1}^M$ and their weights $\{g_w^i\}_{i=1}^M$. We separately visualize the weighted temporal and
 286 spatial flow maps. The mean temporal flow map $\bar{g}_z = \sum_i g_w^i \cdot g_z^i$ describes the weighted centroid of
 287 voxel flows along the time axis. As shown in the fourth column of Fig. 6a, the \bar{g}_z performs like a
 288 occlusion map for bidirectional frame prediction. The pixels in black area cannot be found in the

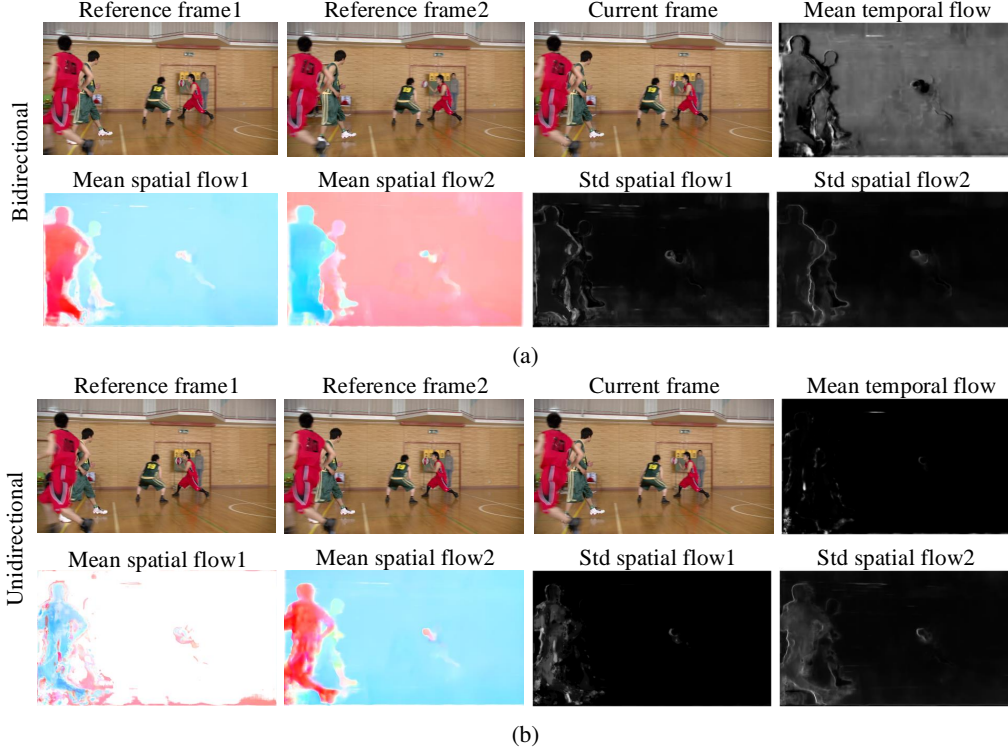


Figure 6: Visualization of the voxel flows for the same target frame with different reference frames, generated by the same model. (a) Bidirectional reference frames. (b) Unidirectional reference frames

289 first reference frame because the basketball player in red covers the background. Hence the voxel
 290 flows pay more attention on the second frame, resulting in large weights. The white area can be
 291 explained in the similar way for the first frame, and the gray area means that the voxel flows pay
 292 equally attention for both frames. For unidirectional frame prediction, the \bar{g}_z generated by the same
 293 model are almost black everywhere, demonstrating the flexibility of trilinear warping for different
 294 prediction mode.

295 We also visualize the weighted mean and weighted standard deviation of spatial flow maps (noted
 296 as mean spatial flow and std spatial flow) to investigate the spatial distribution of voxel flows. We
 297 first round and group the g_z^i to the nearest integer location of reference frames (e.g. 0 or 1), then
 298 separately calculate the mean flow map and std flow map for each group of voxel flows. As shown in
 299 the second and fourth rows of Fig. 6, the grouped spatial mean of voxel flows has similar distribution
 300 with optical flow. Different from optical flow, the voxel flows have large variance in the area of
 301 motion, occlusion and blur, shown in the std spatial flow map. Single optical flow is not able to find a
 302 accurate reference pixel and results in inefficient motion compensation. Multiple flow warping with
 303 weighted coefficients provide a choice to perform motion compensation using multiple reference
 304 pixels with better rate-distortion performance.

305 5 Conclusion

306 In this paper, we propose a versatile learned video coding (VLVC) framework that allows us to train
 307 one model to support various inter prediction modes. To this end, we apply voxel flows as a motion
 308 information descriptor along both spatial and temporal dimensions, and we perform reconstruction
 309 via proposed weighted trilinear warping using voxel flows for more effective motion compensation.
 310 Through formulating different inter prediction modes by a unified polynomial function, we design a
 311 novel flow prediction module to predict accurate motion trajectories. In this way, we significantly
 312 reduce the bits cost of encoding motion information. Thanks to above novel motion compensation
 313 and flow prediction, VLVC not only achieve the support of different inter prediction modes but also

314 yield competitive R-D performance compared to conventional VVC standard, which fosters practical
315 applications of learned video compression technologies.

316 **References**

- 317 [1] VNI Cisco. Cisco visual networking index: Forecast and trends, 2017–2022. White Paper, 1, 2018.
- 318 [2] Thomas Wiegand, Gary J Sullivan, Gisle Bjontegaard, and Ajay Luthra. Overview of the h. 264/avc video
319 coding standard. IEEE Transactions on circuits and systems for video technology, 13(7):560–576, 2003.
- 320 [3] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high efficiency
321 video coding (hevc) standard. IEEE Transactions on circuits and systems for video technology, 22(12):
322 1649–1668, 2012.
- 323 [4] Shan Liu Ye-Kui Wang Benjamin Bross, Jianle Chen. Versatile video coding (draft 9). Draft, JVET, 2020.
- 324 [5] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. arXiv
325 preprint arXiv:1611.01704, 2016.
- 326 [6] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image
327 compression with a scale hyperprior. arXiv preprint arXiv:1802.01436, 2018.
- 328 [7] David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors for
329 learned image compression. In Advances in Neural Information Processing Systems, pages 10771–10780,
330 2018.
- 331 [8] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression with dis-
332 cretized gaussian mixture likelihoods and attention modules. In Proceedings of the IEEE/CVF Conference
333 on Computer Vision and Pattern Recognition, pages 7939–7948, 2020.
- 334 [9] Eirikur Agustsson and Lucas Theis. Universally quantized neural compression. Advances in Neural
335 Information Processing Systems, 33, 2020.
- 336 [10] Zongyu Guo, Zhizheng Zhang, Runsen Feng, and Zhibo Chen. Soft then hard: Rethinking the quantization
337 in neural image compression, 2021.
- 338 [11] Chao-Yuan Wu, Nayan Singhal, and Philipp Krahenbuhl. Video compression through image interpolation.
339 In Proceedings of the European Conference on Computer Vision (ECCV), pages 416–431, 2018.
- 340 [12] Abdelaziz Djelouah, Joaquim Campos, Simone Schaub-Meyer, and Christopher Schroers. Neural inter-
341 frame compression for video coding. In Proceedings of the IEEE International Conference on Computer
342 Vision, pages 6421–6429, 2019.
- 343 [13] Amirhossein Habibian, Ties van Rozendaal, Jakub M Tomczak, and Taco S Cohen. Video compression
344 with rate-distortion autoencoders. In Proceedings of the IEEE International Conference on Computer
345 Vision, pages 7033–7042, 2019.
- 346 [14] Jerry Liu, Shenlong Wang, Wei-Chiu Ma, Meet Shah, Rui Hu, Pranaab Dhawan, and Raquel Urtasun.
347 Conditional entropy coding for efficient video compression, 2020.
- 348 [15] Guo Lu, Wanli Ouyang, Dong Xu, Xiaoyun Zhang, Chunlei Cai, and Zhiyong Gao. Dvc: An end-to-end
349 deep video compression framework. In Proceedings of the IEEE Conference on Computer Vision and
350 Pattern Recognition, pages 11006–11015, 2019.
- 351 [16] Eirikur Agustsson, David Minnen, Nick Johnston, Johannes Balle, Sung Jin Hwang, and George
352 Toderici. Scale-space flow for end-to-end optimized video compression. In Proceedings of the IEEE/CVF
353 Conference on Computer Vision and Pattern Recognition (CVPR), June 2020.
- 354 [17] Jianping Lin, Dong Liu, Houqiang Li, and Feng Wu. M-lvc: Multiple frames prediction for learned video
355 compression. arXiv preprint arXiv:2004.10290, 2020.
- 356 [18] Oren Rippel, Alexander G Anderson, Kedar Tatwawadi, Sanjay Nair, Craig Lytle, and Lubomir Bourdev.
357 Elf-vc: Efficient learned flexible-rate video coding. arXiv preprint arXiv:2104.14335, 2021.
- 358 [19] Ze Cui, Jing Wang, Bo Bai, Tiansheng Guo, and Yihui Feng. G-vae: A continuously variable rate deep
359 image compression framework. arXiv preprint arXiv:2003.02012, 2020.

- 360 [20] Ren Yang, Fabian Mentzer, Luc Van Gool, and Radu Timofte. Learning for video compression with
361 hierarchical quality and recurrent enhancement. [arXiv preprint arXiv:2003.01966](#), 2020.
- 362 [21] Oren Rippel, Sanjay Nair, Carissa Lew, Steve Branson, Alexander G Anderson, and Lubomir Bourdev.
363 Learned video compression. In [Proceedings of the IEEE International Conference on Computer Vision](#),
364 pages 3454–3463, 2019.
- 365 [22] Guo Lu, Chunlei Cai, Xiaoyun Zhang, Li Chen, Wanli Ouyang, Dong Xu, and Zhiyong Gao. Content
366 adaptive and error propagation aware deep video compression. [arXiv preprint arXiv:2003.11282](#), 2020.
- 367 [23] Haojie Liu, Lichao Huang, Ming Lu, Tong Chen, Zhan Ma, et al. Learned video compression via joint
368 spatial-temporal correlation exploration. [arXiv preprint arXiv:1912.06348](#), 2019.
- 369 [24] Zhibo Chen, Tianyu He, Xin Jin, and Feng Wu. Learning for video compression. [IEEE Transactions on](#)
370 [Circuits and Systems for Video Technology](#), 30(2):566–576, 2020.
- 371 [25] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis
372 using deep voxel flow. In [Proceedings of the IEEE International Conference on Computer Vision](#), pages
373 4463–4471, 2017.
- 374 [26] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In
375 [Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition](#), pages 670–679, 2017.
- 376 [27] Fitsum A Reda, Guilin Liu, Kevin J Shih, Robert Kirby, Jon Barker, David Tarjan, Andrew Tao, and
377 Bryan Catanzaro. Sdc-net: Video prediction using spatially-displaced convolution. In [Proceedings of the](#)
378 [European Conference on Computer Vision \(ECCV\)](#), pages 718–733, 2018.
- 379 [28] Wenbo Bao, Wei-Sheng Lai, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Memc-net: Motion
380 estimation and motion compensation driven neural network for video interpolation and enhancement. [IEEE](#)
381 [transactions on pattern analysis and machine intelligence](#), 2019.
- 382 [29] Hyeongmin Lee, Taeoh Kim, Tae-young Chung, Daehyun Pak, Yuseok Ban, and Sangyoung Lee. Adacof:
383 adaptive collaboration of flows for video frame interpolation. In [Proceedings of the IEEE/CVF Conference](#)
384 [on Computer Vision and Pattern Recognition](#), pages 5316–5325, 2020.
- 385 [30] Zhihao Shi, Xiaohong Liu, Kangdi Shi, Linhui Dai, and Jun Chen. Video frame interpolation via generalized
386 deformable convolution. [IEEE Transactions on Multimedia](#), 2021.
- 387 [31] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In [Proceedings of the](#)
388 [IEEE/CVF Conference on Computer Vision and Pattern Recognition](#), pages 5437–5446, 2020.
- 389 [32] Xiangyu Xu, Li Siyao, Wenxiu Sun, Qian Yin, and Ming-Hsuan Yang. Quadratic video interpolation.
390 [arXiv preprint arXiv:1911.00627](#), 2019.
- 391 [33] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid,
392 warping, and cost volume. In [Proceedings of the IEEE Conference on Computer Vision and Pattern](#)
393 [Recognition](#), pages 8934–8943, 2018.
- 394 [34] Runsen Feng, Yaojun Wu, Zongyu Guo, Zhizheng Zhang, and Zhibo Chen. Learned video compression
395 with feature-level residuals. In [Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern](#)
396 [Recognition Workshops](#), pages 120–121, 2020.
- 397 [35] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with
398 task-oriented flow. [International Journal of Computer Vision](#), 127(8):1106–1125, 2019.
- 399 [36] Alexandre Mercat, Marko Viitanen, and Jarmo Vanne. Uvg dataset: 50/120fps 4k sequences for video
400 codec analysis and development. In [Proceedings of the 11th ACM Multimedia Systems Conference](#), pages
401 297–302, 2020.
- 402 [37] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assess-
403 ment. In [The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers](#), 2003, volume 2,
404 pages 1398–1402. Ieee, 2003.
- 405 [38] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. [arXiv preprint](#)
406 [arXiv:1412.6980](#), 2014.
- 407 [39] HEVC Official Test Model HM. <https://hevc.hhi.fraunhofer.de>.
- 408 [40] VVC Official Test Model VTm. <https://jvet.hhi.fraunhofer.de>.

409 Checklist

410 The checklist follows the references. Please read the checklist guidelines carefully for information on
411 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
412 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
413 the appropriate section of your paper or providing a brief inline description. For example:

- 414 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 415 • Did you include the license to the code and datasets? **[No]** The code and the data are
416 proprietary.
- 417 • Did you include the license to the code and datasets? **[N/A]**

418 Please do not modify the questions and only use the provided macros for your answers. Note that the
419 Checklist section does not count towards the page limit. In your paper, please delete this instructions
420 block and only keep the Checklist section heading above along with the questions/answers below.

- 421 1. For all authors...
 - 422 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
423 contributions and scope? **[Yes]**
 - 424 (b) Did you describe the limitations of your work? **[No]**
 - 425 (c) Did you discuss any potential negative societal impacts of your work? **[No]**
 - 426 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
427 them? **[Yes]**
- 428 2. If you are including theoretical results...
 - 429 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
 - 430 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 431 3. If you ran experiments...
 - 432 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
433 mental results (either in the supplemental material or as a URL)? We describe datasets,
434 network architecture and more details in the supplementary.
 - 435 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
436 were chosen)? **[Yes]**
 - 437 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
438 ments multiple times)? **[No]**
 - 439 (d) Did you include the total amount of compute and the type of resources used (e.g., type
440 of GPUs, internal cluster, or cloud provider)? **[No]** In the supplementary.
- 441 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 442 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - 443 (b) Did you mention the license of the assets? **[Yes]**
 - 444 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
 - 445 (d) Did you discuss whether and how consent was obtained from people whose data you're
446 using/curating? **[No]**
 - 447 (e) Did you discuss whether the data you are using/curating contains personally identifiable
448 information or offensive content? **[No]**
- 449 5. If you used crowdsourcing or conducted research with human subjects...
 - 450 (a) Did you include the full text of instructions given to participants and screenshots, if
451 applicable? **[N/A]**
 - 452 (b) Did you describe any potential participant risks, with links to Institutional Review
453 Board (IRB) approvals, if applicable? **[N/A]**
 - 454 (c) Did you include the estimated hourly wage paid to participants and the total amount
455 spent on participant compensation? **[N/A]**