

MoVie: MULTIMODAL VIDEO COMPRESSION WITH TEXT GUIDANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advances in deep video compression have significantly improved rate-distortion performance. Compared to traditional codecs that rely on handcrafted motion estimation and block-based prediction, deep learning-based methods can learn more flexible and content-adaptive representations, leading to better compression efficiency. However, most existing approaches still focus primarily on low-level pixel motion modeling and lack semantic awareness, which limits their ability to preserve perceptual quality in complex scenes. In this paper, we propose **MoVie**, a **Multimodal Video** compression framework built upon a Text-guided Video Transformer–CNN Mixed block (*Text-VideoTCM*). Instead of relying on image-oriented feature extractors that ignore temporal cues, we design a video-focused network, jointly modeling local spatial structures and temporal dynamics, achieving a remarkable trade-off between computational cost and perceptual performance. To enhance semantic perception, a dual-stage text fusion mechanism is introduced: *Extractor* modules distill text-aware features at early layers, while *Injector* modules inject refined semantics in deeper stages. We also introduce a new recipe *history-conditioned coding* that adaptively leverages both previous and aggregated historical frames, alongside a spatial-channel factorized entropy model tailored for window-based Transformer, which jointly captures local spatial structures and inter-channel dependencies. Averaged over the UVG and MCL-JCV datasets, MoVie achieves substantial BD-rate reductions relative to HM: **−50.23%** for FID and **−14.64%** for LPIPS(VGGNet). While maintaining superior perceptual quality, our method substantially reduces computational cost, requiring only **55.76%** of the per-pixel kMACs of DCVC-FM.

1 INTRODUCTION

Image and video compression are essential for efficient storage, transmission, and streaming in computer vision and multimedia applications. Traditional codecs like JPEG Wallace (1991), H.265/HEVC Sullivan et al. (2012), and H.266/VVC Bross et al. (2021) rely on hand-crafted modules such as block transforms, quantization, and entropy coding to exploit statistical redundancies.

With the rise of deep learning, video compression has made significant strides. The DVC framework Lu et al. (2019) introduced end-to-end modeling of motion, residuals, and temporal priors, establishing a foundation for learned video coding. However, most existing models Sheng et al. (2022); Li et al. (2023) still adopt frame-by-frame encoding or rely on low-level motion (e.g., optical flow), limiting their ability to model long-term temporal and semantic dependencies. While Transformer-based methods like VCT Mentzer et al. (2022) offer improved temporal modeling, their use of global attention leads to prohibitive costs from the explosion of spatiotemporal tokens.

Although designed for video, many recent learned approaches still rely on image-based encoders, failing to fully exploit temporal structures. For example, VCT Mentzer et al. (2022) uses the spatial ELIC backbone He et al. (2022), and FLAVC Zhang et al. (2025) employs LIC-TCM Liu et al. (2023), lacking explicit modeling of inter-frame redundancy. Meanwhile, DCVC variants Li et al. (2021); Sheng et al. (2022); Li et al. (2023) rely heavily on CNNs for motion and residual prediction, which limits the ability to capture long-range dependencies and semantics in dynamic scenes.

Parallel progress in multimodal learning shows natural-language descriptions, as revealed by CLIP Radford et al. (2021), offer high-level semantics useful for vision tasks. In image compression, in-

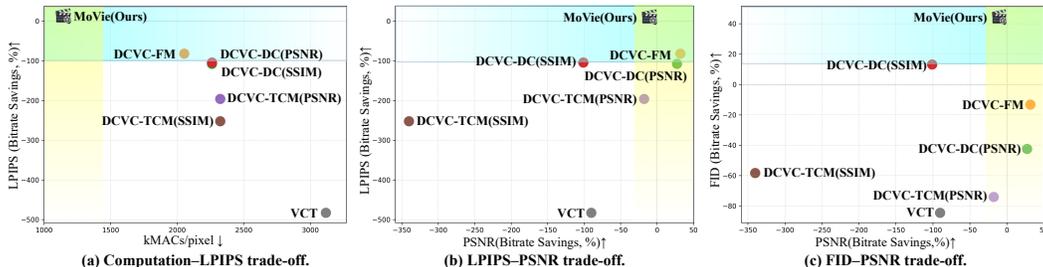


Figure 1: Trade-off performance of various LVC models in two aspects. Higher opacity indicates a better trade-off. We report bitrate savings ($-BD\text{-Rate}$, %) relative to HM-18.0 on UVG. Higher values indicate better performance. Note that LPIPS and FID are negated for $BD\text{-Rate}$ calculation.

tegrating textual priors improves perceptual quality Jiang et al. (2023); Qin et al. (2023); Lee et al. (2024). Although Zhang et al. (2024) have achieved strong results in video compression with multimodal large models, their high memory usage and computational cost limit practical application.

To address these limitations, we propose a Text-guided Video Transformer-CNN Mixed block (*Text-VideoTCM*), a unified backbone that efficiently integrates textual semantics with localized spatiotemporal modeling. Our design couples a Video Swin Transformer Liu et al. (2022) with a 3D CNN Tran et al. (2015) to jointly capture local spatial and temporal dependencies. The Swin Transformer restricts attention to 3D windows, enabling efficient modeling of short-range motion dynamics, while the 3D CNN focuses on preserving local textures and motion boundaries. The resulting fused representation strikes an effective balance between semantic abstraction and detail preservation. To incorporate semantic guidance from text, we further introduce a two-stage fusion mechanism: *Extractor* and *Injector*. Specifically, we design an Extractor to apply cross-attention between video tokens and CLIP text features, generating frame-specific semantics, and an Injector to embed these cues into latent representations, enhancing semantic consistency during compression.

Meanwhile, mainstream entropy models still treat latent features as short-range, missing richer dependencies spread across long sequences and channel groups. To bridge this gap, we design a *history-conditioned coding* strategy that caches and reuses long-term latent summaries, enabling the Transformer to condition on historical frames. We further propose a spatial-channel factorized entropy model that estimates probabilities across spatial neighborhoods and channel groups, enabling adaptive bit allocation based on motion and semantic saliency. These modules jointly improve rate control, reduce redundancy, and enhance temporal coherence.

Extensive experiments demonstrate that MoVi strikes a balance between perceptual quality, computational cost, and pixel-level fidelity. As shown in Fig. 1 (a), MoVi reduces computation by nearly half, owing to the efficient Text-VideoTCM block that models local spatiotemporal patterns. Our history-conditioned entropy model with spatial-channel factorization improves rate control by capturing long-range dependencies, while the text-guided fusion boosts perceptual quality by refining the distortion term. These components jointly improve overall $BD\text{-Rate}$ performance, as shown in Fig. 1 (b). MoVi consistently outperforms perceptual baselines on perceptual metrics while maintaining strong PSNR performance. Besides objective metrics such as $BD\text{-rate}$, MoVi also has good performance on perceptual metrics. Fig 1 (c) shows that MoVi consistently obtains high FID while maintaining strong PSNR.

Our key contributions are summarized as follows:

- We propose MoVi, a multimodal video compression framework that uses textual priors and cross-frame fusion for semantically guided and perceptually coherent compression.
- We propose the Text-VideoTCM block, which integrates textual semantics with localized spatiotemporal modeling to balance abstraction and detail preservation.
- We propose a history-conditioned coding strategy and a spatial-channel factorized entropy model to better capture long-term temporal dependencies and enhance entropy estimation.
- Extensive experiments demonstrate that MoVi achieves substantial $BD\text{-rate}$ reductions while excelling in both computation-perception and pixel-perception fidelity trade-offs.

2 RELATED WORK

2.1 CNN-BASED VIDEO COMPRESSION

Early works like DVC Lu et al. (2019) jointly optimize motion estimation, residual compression, and entropy modeling. Later methods adopt conditional coding Liu et al. (2020) to enhance context representation, with DCVC Li et al. (2021) introducing feature-domain motion compensation. DCVC-DC Li et al. (2023) eliminates explicit motion estimation using deformable alignment with spatiotemporal priors, whereas DCVC-FM Li et al. (2024b) adopts flow-matching for motion-free alignment. Recent studies Chen et al. (2024); Sheng et al. (2024) improve context diversity at the cost of increased complexity. Although effective at modeling temporal dependencies, 3D convolutions Pessoa et al. (2020) are rarely used in CNN-based video compression due to their high cost.

2.2 TRANSFORMER-BASED VIDEO COMPRESSION

VCT Mentzer et al. (2022) leverages Transformer-based temporal modeling in the entropy coding module via implicit motion reasoning, achieving competitive compression with a simplified design, though at the cost of high memory overhead. Nonetheless, its autoencoder retains a CNN-based design for spatial representation learning. Based on VCT, FLAVC Zhang et al. (2025) improves the encoder–decoder architecture using a Transformer-CNN mixed design Liu et al. (2023). However, the autoencoder processes frames independently, neglecting temporal dependencies. In contrast, our method employs a video-specific autoencoder that fully exploits inter-frame consistency.

2.3 TEXT-GUIDED PERCEPTUAL COMPRESSION

Following the success of vision-language models, text-guided image compression methods have recently emerged Bhowan et al. (2018); Weissman (2023), leveraging high-level semantics to enhance perceptual quality. One line of work Pan et al. (2022); Lei et al. (2023) explores using pre-trained text-guided generative models (e.g., diffusion models Rombach et al. (2022)) as decoders. Some approaches Wan et al. (2025); Zhang et al. (2024) employ pretrained generative models or multimodal large language models for video compression, but this incurs substantial computational cost and storage overhead, making them impractical for efficient or large-scale deployment.

Another line of work trains decoders from scratch, where Jiang et al. (2023) inject textual information into both encoder and decoder, while Qin et al. (2023) introduce text only into the decoder via semantic-spatial aware blocks. Recent methods like TACO Lee et al. (2024) show that adding text guidance only in the encoder can effectively boost perceptual quality without increasing compression cost. This opens a promising direction for end-to-end text-guided video compression, yet it remains unexplored in prior work. Our study bridges this gap by integrating efficient inter-frame modeling with semantic-aware compression into a unified, low-complexity framework.

3 METHOD

3.1 OVERALL ARCHITECTURE

The overall architecture of our proposed multimodal video compression framework is illustrated in Fig. 2. Inspired by the success of Transformer-based models Liu et al. (2022) in video understanding, we adapt the Video Swin Transformer—originally designed for recognition tasks—into a self-encoder structure optimized for video compression. Specifically, we augment the original architecture, which only contains a downsampling path, with a symmetric upsampling counterpart, forming a fully matched encoder-decoder design suitable for end-to-end video compression.

Our compression pipeline adopts a symmetric autoencoder architecture composed of an encoder $\mathbf{E}(\cdot)$ and a decoder $\mathbf{D}(\cdot)$, each consisting of three stages. In the encoder, we first partition the input video $x \in \mathbb{R}^{T \times H \times W \times 3}$ into non-overlapping patches along the spatial dimensions. Specifically, we divide the input into smaller spatial patches while preserving the full temporal resolution, unlike the original Video Swin Transformer which reduces the temporal dimension. This design choice ensures that every frame can be faithfully reconstructed in video compression tasks.

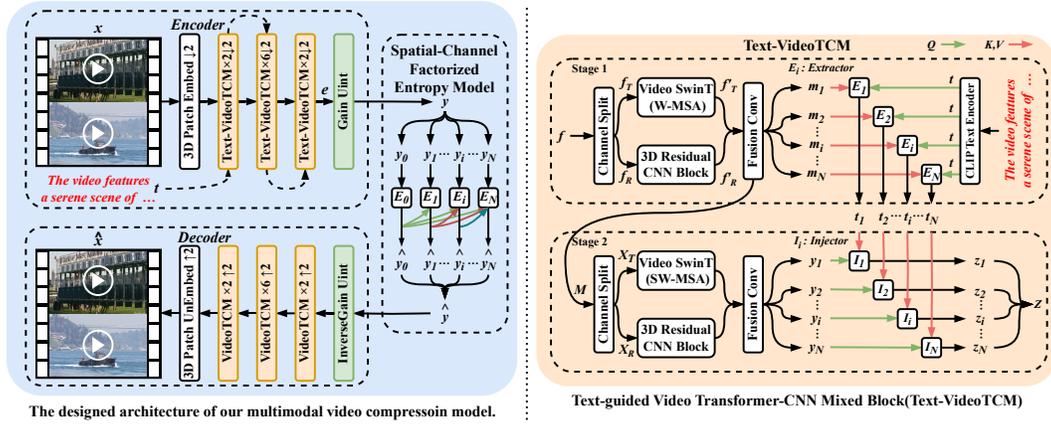


Figure 2: Overview of our multimodal video compression framework and the proposed *Text-guided Video Transformer-CNN Mixed Block (Text-VideoTCM)*. The encoder extracts latent features using a sequence of *Text-VideoTCM* or *VideoTCM* blocks, where *3D Patch Embed/UnEmbed* modules perform spatiotemporal tokenization and reconstruction. Unlike *Text-VideoTCM*, *VideoTCM* does not use text guidance during decoding and lacks the *Extractor (E)* and *Injector (I)* modules.

We then introduce a *Text-VideoTCM* block that captures spatial and temporal dependencies via parallel Video Swin Transformer Liu et al. (2022) and 3D CNN Tran et al. (2015) branches. The detailed architecture is described in the next section. Between stages, patch merging operations progressively reduce the spatial resolution while expanding the channel dimension, forming a hierarchical spatial pyramid. Notably, to ensure the reconstructability of compressed features, we perform spatial downsampling only four times, resulting in a final feature resolution of $\frac{H}{16} \times \frac{W}{16}$.

To support one single model operating at multiple bit-rates, we follow AG-VAE Cui et al. (2021) and adopt a gain-based rate-adjustment mechanism. Specifically, we incorporate a learnable Gain Unit in the encoder and its corresponding Inverse Gain Unit in the decoder. Given a user-specified rate-control index $s \in \{0, 1, \dots, S-1\}$ —which maps to a preset Lagrangian multiplier λ_s —the units scale the latent features before quantization. We refer to this process as f_{gain} , with the detailed algorithm provided in Appendix C.

The latent representation is processed by an entropy model $\mathbf{H}(\cdot)$, which estimates the distribution of y for accurate bitrate modeling. Before entropy coding, the features are quantized by $\mathbf{Q}(\cdot)$, which discretizes the continuous latent representations into a finite set of symbols. The resulting quantized representation \hat{y} is then entropy coded into a bitstream and passed to the decoder to reconstruct the video \hat{x} . Thus, the overall pipeline of our framework can be summarized as follows:

$$e = \mathbf{E}(x, t), \quad y = f_{\text{gain}}(s, e), \quad \hat{y} = \mathbf{Q}(\mathbf{H}(y)), \quad \hat{x} = \mathbf{D}(\hat{y}). \quad (1)$$

Global textual information is exploited only within the encoder during offline processing. Because this text offers a holistic description of the entire sequence, injecting it into entropy coding or decoding would potentially break frame-wise causality—the reconstruction of an early frame could then depend on knowledge of future content. In contrast, the encoder is allowed to inspect all frames beforehand, so integrating the global text at feature-extraction time introduces no causality issue nor additional bitrate overhead.

3.2 TEXT-VIDEO TCM

In this section, we provide a detailed description of the proposed *Text-VideoTCM* block. This architecture employs window-based self-attention (W-MSA) and shifted window-based self-attention (SW-MSA) in two successive stages to model intra-window and inter-window dependencies, thereby achieving better local feature representations. Within each stage, a dual-branch design integrates a Video Swin Transformer module and a 3D Residual CNN module to jointly capture both local details and global structure, leading to enhanced visual performance. Furthermore, to improve perceptual quality, we incorporate text features in both stages using different fusion strategies, aiming to en-

216 hance semantic consistency and cross-modal alignment. In the following, we elaborate on the design
 217 of each stage in detail.

218 Assuming that the input tensor is $f \in \mathbb{R}^{T \times H_f \times W_f \times C}$, we first evenly split it along the channel
 219 dimension into two sub-tensors, f_T and f_R , each with a shape of $\mathbb{R}^{T \times H_f \times W_f \times \frac{C}{2}}$. Then, f_T and f_R
 220 are respectively fed into a Video Swin Transformer (Video SwinT) block and a 3D Residual CNN
 221 Block, which are used to extract local and non-local features in parallel, resulting in the outputs f'_T
 222 and f'_R . Afterward, these two features are fused via a residual convolution module to obtain the
 223 final representation $M \in \mathbb{R}^{T \times H_f \times W_f \times C}$, and can be expressed as $\{m_1, m_2, \dots, m_i, \dots, m_T\} \in$
 224 $\mathbb{R}^{H_f \times W_f \times C}$. The entire fusion procedure is formally defined as:

$$\begin{aligned} 225 & f_T, f_R = \text{Split}(f), \\ 226 & f'_T = \text{VideoSwinT}(f_T), f'_R = \text{Res3D}(f_R), \\ 227 & M = f + \text{Conv3D}_{1 \times 1}(\text{Cat}(f'_T, f'_R)). \end{aligned} \quad (2)$$

230 The textual description is encoded by a frozen CLIP text encoder into $t \in \mathbb{R}^{L \times 512}$, where we fix the
 231 token length to $L = 38$ by truncating or padding the CLIP token sequence. Since it describes the
 232 entire video, t is tiled T times along the temporal axis and fused with image features X . Extractors
 233 $\{E_1, \dots, E_T\}$ then apply cross-attention to each token x_i with t : following CLIP, t serves as query
 234 Q and m_i provides keys/values (K, V) . This yields per-frame text features $\{t_i\}_{i=1}^T \in \mathbb{R}^{T \times L \times 512}$.
 235 The output of the extractor for the i -th frame is computed as:

$$236 \quad t_i = E_i(t, m_i) = t + \text{CrossAttention}(t, \text{Lin}(m_i)). \quad (3)$$

237 Here, $\text{Lin}(\cdot)$ projects image features to the text feature dimension, and $\text{CrossAttention}(\cdot)$ is a multi-
 238 head attention module with text as Q , and image as K, V .

239 Stage 2 re-encodes the feature map X with the same dual-branch design as Stage 1; the Video Swin
 240 Transformer uses shifted-window attention. To inject per-frame semantics into the visual stream,
 241 we introduce injectors $\{I_1, \dots, I_T\}$. Each I_i takes the second-stage image feature y_i and its paired
 242 text-guided token t_i , producing z_i as:

$$243 \quad z_i = I_i(y_i, t_i) = y_i + \text{CrossAttention}(y_i, \text{Lin}(t_i)). \quad (4)$$

244 The image feature serves as Q , and the text feature as K, V in the cross-attention module.

247 3.3 SPATIAL-CHANNEL FACTORIZED ENTROPY MODEL

248 Fig. 3 presents an overview of our proposed Spatial-Channel Factorized Entropy Model, which aims
 249 to enhance entropy estimation by jointly capturing spatial, temporal, and channel-wise dependencies
 250 in video latents. For clarity, we omit the batch and temporal indices in the following derivation and
 251 focus on a single frame latent $y_i \in \mathbb{R}^{H \times W \times C}$.

252 **Left: Windowed Spatial-Channel Factorization and Temporal Cache.** We first partition y_i into
 253 non-overlapping 8×8 spatial windows:

$$254 \quad y_i[p, q] \in \mathbb{R}^{8 \times 8 \times C}, \quad p = 1, \dots, \frac{H}{8}, \quad q = 1, \dots, \frac{W}{8}, \quad (5)$$

255 and denote the total number of windows by $N_w = (H/8)(W/8)$. Each window is then split along
 256 the channel dimension into 6 equal-sized groups:

$$257 \quad y_i[p, q] = (y_i^{(1)}[p, q], \dots, y_i^{(6)}[p, q]), \quad y_i^{(g)}[p, q] \in \mathbb{R}^{8 \times 8 \times C/6}. \quad (6)$$

261 We call $y_i^{(g)}[p, q]$ a *channel slice* and index it by the triple (p, q, g) . In the right part of Fig. 3, each
 262 slice is denoted by s_i and is processed sequentially along the channel dimension, i.e., at step g the
 263 input to the S-C Factorized Entropy Model encoder is $s_i = y_i^{(g)}[p, q]$. The corresponding decoded
 264 slice \hat{s}_i is then stored and concatenated with all previously decoded slices in the same window to
 265 form $\hat{s}_{<i}$, which serves as an additional conditional input when predicting the distribution of the
 266 next channel slice.

267 To effectively model temporal priors, we employ a VideoTCM Fusion module that aggregates in-
 268 formation from the previous latent feature \hat{y}_{i-1} and historical context latent feature $\hat{y}_{<i-1}$ to predict
 269 the mean F_{mean} and scale F_{scale} parameters for the current latent y_i .

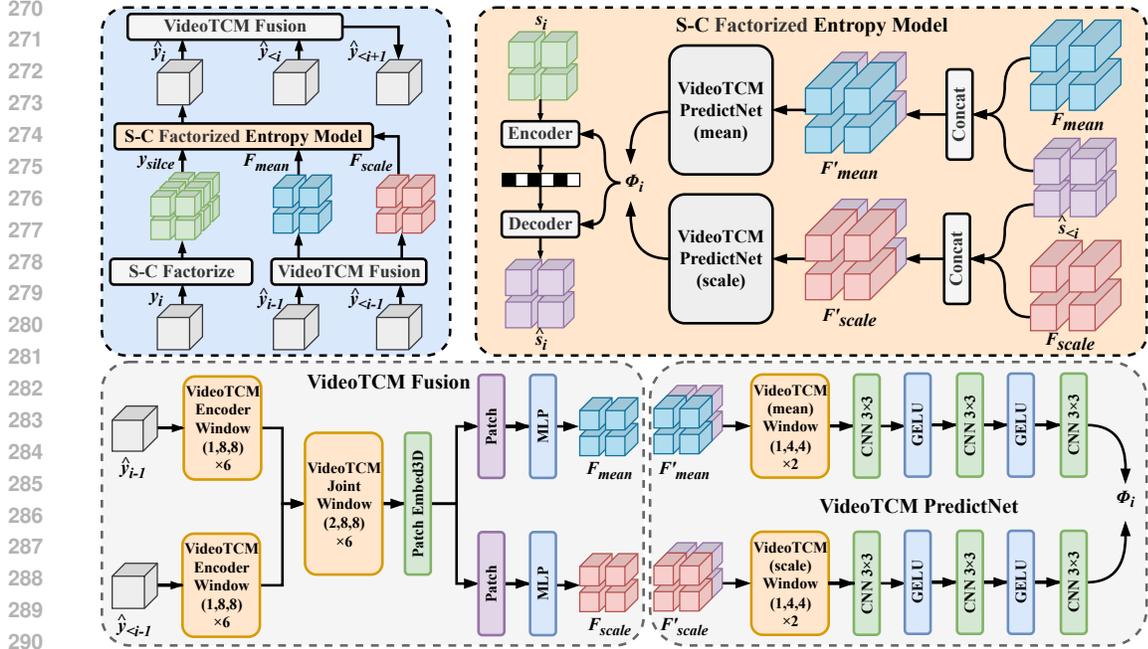


Figure 3: Overview of the proposed Spatial-Channel Factorized Entropy Model. Here, y_i denotes the current encoded frame, while \hat{y}_{i-1} and $\hat{y}_{<i-1}$ represent the previously and historical reconstructed frames. *S-C Factorized* refers to spatial-channel factorized.

Then, y_{slice} , along with the predicted F_{mean} and F_{scale} , are fed into the spatial-channel entropy coding block to produce the quantized latent \hat{y}_i . The output is further \hat{y}_i fused with $\hat{y}_{<i}$ to generate $\hat{y}_{<i+1}$ for the entropy modeling of y_{i+1} .

The recursive modeling process can be expressed as follows:

$$p_{\theta}(y_i | \hat{y}_{i-1}, \hat{y}_{<i-1}) = \text{Entropy}_{\theta}(\hat{y}_{i-1}, \hat{y}_{<i-1}), \quad \hat{y}_{<i} = f(\hat{y}_{<i-1}, \hat{y}_{i-1}), \quad \hat{y}_{<1} = y_0. \quad (7)$$

Right: S-C Factorized Entropy Model. The right part of the figure illustrates the proposed spatial-channel factorized entropy model, which refines the entropy parameters by learning spatial- and channel-aware distributions.

From the left part of the pipeline, the obtained latent tensor y_{slice} is first split along the channel dimension into a set of features s_i . Each channel slice s_i is then entropy coded sequentially to produce the quantized output \hat{s}_i .

Meanwhile, the conditional inputs F_{mean} and F_{scale} are concatenated with the encoded results of previous channels $s_{<i}$, yielding updated representations F'_{mean} and F'_{scale} . These are then fed into two independent prediction networks (predictNet) to produce the quantization parameters ϕ_i .

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

We optimize the model with four losses over the original video x , textual description t , reconstructed output \hat{x} , and quantized latent \hat{y} , the overall loss is defined as:

$$\mathcal{L} = \sum_{i=1}^T \left\{ \mathcal{R}(\hat{y}_i) + \lambda_s \mathcal{D}(x_i, \hat{x}_i) + k_p \text{LPIPS}(x_i, \hat{x}_i) + k_j [\mathcal{L}_{\text{con}}(f_I(\hat{x}_i), f_T(t)) + \beta \|f_I(x_i) - f_I(\hat{x}_i)\|_2^2] \right\}. \quad (8)$$

The first three terms correspond to standard losses: bitrate $\mathcal{R}(\cdot)$, distortion $\mathcal{D}(\cdot, \cdot)$ (MSE), and perceptual similarity (LPIPS). $f_I(\cdot)$ and $f_T(\cdot)$ denote CLIP’s image and text encoders, and \mathcal{L}_{con} is

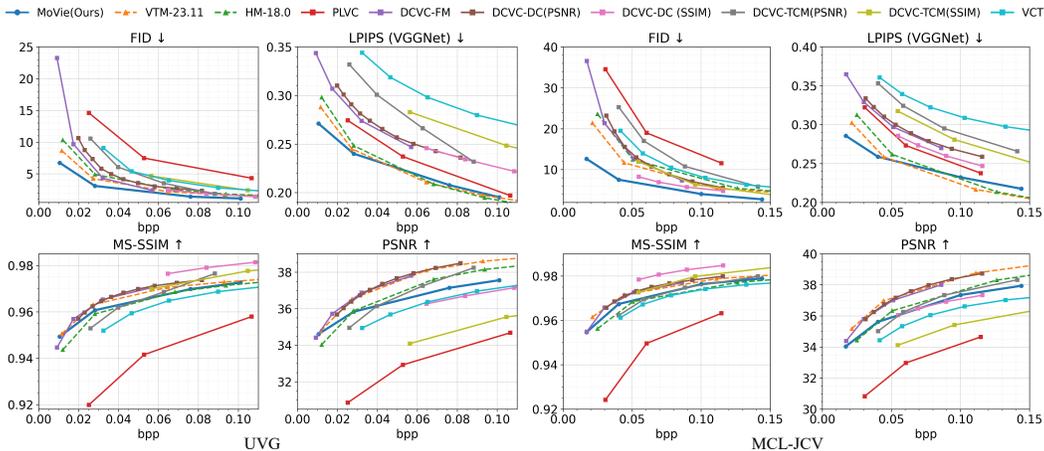


Figure 4: Overall compression performance with RD curves (\downarrow lower is better, \uparrow higher is better).

| Method | kMAC/pixel | UVG | | | | MCL-JCV | | | |
|---------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | | FID | LPIPS | MS-SSIM | PSNR | FID | LPIPS | MS-SSIM | PSNR |
| HM-18.0 | / | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| VTM-23.11 | / | -17.38 | -1.51 | -23.93 | -29.14 | -17.03 | -9.52 | -30.00 | -30.08 |
| PLVC | / | 219.19 | 41.19 | 356.35 | 565.80 | 98.93 | 42.90 | 229.91 | 323.97 |
| DCVC-FM | 2051.8 | 13.21 | 82.00 | -29.49 | -31.82 | 5.84 | 74.77 | -28.50 | -25.82 |
| DCVC-DC (PSNR) | 2260.8 | 42.42 | 107.60 | -30.25 | -27.30 | 12.70 | 93.48 | -31.70 | -26.58 |
| DCVC-DC (SSIM) | 2260.8 | -13.11 | 104.20 | -62.36 | 101.10 | -26.82 | 68.00 | -66.43 | 33.12 |
| DCVC-TCM (PSNR) | 2322.1 | 74.03 | 195.94 | -1.56 | 17.80 | 56.96 | 187.30 | 0.60 | 17.54 |
| DCVC-TCM (SSIM) | 2322.1 | 58.28 | 252.23 | -40.23 | 340.24 | -13.51 | 141.26 | -45.13 | 182.66 |
| VCT | 3115.9 | 84.43 | 482.72 | 47.61 | 90.23 | 29.66 | 349.99 | 23.96 | 79.13 |
| MoVie (Ours) | 1144.1 | -45.09 | -13.03 | -13.98 | 7.62 | -55.36 | -16.25 | -15.67 | 10.91 |

Table 1: Comparison on UVG and MCL-JCV in BD-Rate (%) with HM-18.0 as anchor. BD-Rate is computed as $-M$ for lower-is-better metrics. MAC is measured in kMACs/pixel.

the CLIP contrastive loss. We adopt the hyperparameters from TACO Lee et al. (2024), setting $k_p=1$, $k_j=0.005$, and $\beta=40$. Specifically, we follow the same schedule as CompressAI Bégin et al. (2020), assigning $\lambda_s = 0.00045, 0.0018, 0.0067, \text{ and } 0.0200$ for $s = 0, 1, 2, 3$, respectively. Following VCT Mentzer et al. (2022), we adopt a three-stage training: autoencoder pretraining with distortion/perceptual losses(w/o \mathcal{R}), entropy model training with \mathcal{R} , and joint fine-tuning. We train the three stages for 1M, 250k, and 500k steps respectively on two NVIDIA A6000 GPUs.

4.2 DATASETS AND EVALUATION

Following standard practice, we train our model on Vimeo90K Xue et al. (2019) using random 256×256 crops. For evaluation, we adopt UVG Mercat et al. (2020) and MCL-JCV Wang et al. (2016) to ensure scene diversity. Additional results on four HEVC datasets are provided in the Appendix J. All captions are generated by LLaVA Li et al. (2024a), with prompts included in the Appendix F. We evaluate the first 96 frames of each sequence under the standard protocol. Rate-distortion (RD) performance is assessed using four metrics: FID and LPIPS for perceptual quality, and PSNR and MS-SSIM for pixel-level fidelity. Although MS-SSIM models structural similarity across multiple scales, it still relies on direct pixel-wise comparisons with the reference, and thus is generally considered a fidelity rather than a perceptual metric.

4.3 RATE-DISTORTION PERFORMANCE

Our baselines include the Transformer model VCT Mentzer et al. (2022) and the DCVC family—DCVC-TCM Sheng et al. (2022), DCVC-DC Li et al. (2023), and DCVC-FM Li et al. (2024b). For DCVC-TCM/DCVC-DC we report both PSNR- and MS-SSIM-optimized variants;

| Method | Encoding time | Decoding time | Encoder kMACs/pixel | Decoder kMACs/pixel |
|---------|---------------|---------------|---------------------|---------------------|
| DCVC-DC | 458 | 404 | 1343.7 | 917.1 |
| DCVC-FM | 439 | 377 | 1132.4 | 919.4 |
| VCT | 614 | 531 | 1570.6 | 1515.1 |
| MoVie | 324 | 255 | 657.0 | 457.1 |

Table 2: Encoding/decoding time and MACs comparison. Times are averaged per-frame reconstruction on an NVIDIA A6000 GPU.

DCVC-FM is reported in its unified configuration. We also include traditional reference codecs HM-18.0 and VTM-23.11, using their latest releases; encoding commands are provided in Appendix D. Additionally, we compare with PLVC Yang et al. (2022), a representative work in perceptual video compression, to demonstrate the perceptual quality improvements achieved by our method.

As shown by the RD curves in Fig. 4 and the BD-Rate results in Table 1, MoVie delivers state-of-the-art perceptual quality, especially at low bitrates. It consistently attains the lowest FID across the full rate range, indicating that its text-guided, semantics-aware features better align with natural image statistics. Moreover, on LPIPS, MoVie leads at low bitrates and achieves the largest BD-Rate reductions overall.

Although MoVie is not the top performer in PSNR/MS-SSIM, it remains competitive with HM in PSNR BD-Rate and reduces MS-SSIM BD-Rate by 13.98%—despite not being explicitly optimized for SSIM. In the low-bitrate regime (leftmost points in Fig. 4), it achieves comparable fidelity to other methods. Compared with PSNR-oriented models, MoVie offers better perceptual quality (lower LPIPS) with only a minor PSNR trade-off. While MS-SSIM-optimized models excel on MS-SSIM, MoVie significantly outperforms them on the three metrics.

In addition, Table 2 reports separate encoding and decoding times together with the encoder/decoder complexities in kMACs/pixel. MoVie attains both the fastest codec runtime and the lowest computational cost, largely thanks to the efficiency of its window-based processing design.

| Method | kMACs/pixel | Inference time | PSNR \uparrow | MS-SSIM \uparrow | FID \downarrow | LPIPS \downarrow |
|------------------------------|---------------|----------------|-----------------|--------------------|------------------|--------------------|
| Enc-V / Dec-V (Vie) | 226.07 | 187ms | 46.78 | 0.99826 | 0.0062 | 0.0018 |
| Enc-T / Dec-V (MoVie) | 235.12 | 197ms | 46.49 | 0.99796 | 0.0041 | 0.0009 |
| Enc-V / Dec-T | 235.12 | 196ms | 45.08 | 0.99673 | 0.0053 | 0.0012 |
| Enc-T / Dec-T | 244.11 | 206ms | 45.69 | 0.99686 | 0.0057 | 0.0012 |
| LIC-TCM | 944.50 | 577ms | 45.58 | 0.99757 | 0.0065 | 0.0022 |
| ELIC | 479.52 | 183ms | 43.01 | 0.99411 | 0.0221 | 0.0098 |
| Global Text | 228.76 | 190ms | 40.04 | 0.99565 | 1.6403 | 0.0037 |
| Copy Text | 230.76 | 190ms | 44.59 | 0.99789 | 0.0212 | 0.0020 |
| Constant Text | 235.12 | 197ms | 40.93 | 0.99138 | 0.0593 | 0.0111 |

Table 3: Comparison of different autoencoders on UVG, **without** entropy coding. MAC is measured in kMACs/pixel. The reported inference time corresponds to the average time to reconstruct a single frame. All measurements are obtained on an NVIDIA RTX 3090 GPU.

4.4 EFFECT OF TEXT GUIDANCE

To better isolate the effect of text guidance, Table 3 reports variants with different text-injection locations: Enc-V/Dec-V (Vie, no text on either side), Enc-T/Dec-V (MoVie, text only in the encoder), Enc-V/Dec-T (text only in the decoder), and Enc-T/Dec-T (text on both encoder and decoder). The results show that MoVie (Enc-T/Dec-V) achieves the best perceptual scores (FID and LPIPS) with only a minor drop in PSNR compared to Vie. In contrast, adding text only to the decoder (Enc-V/Dec-T) or to both encoder and decoder (Enc-T/Dec-T) does not lead to further improvements and can even hurt reconstruction quality.

A plausible explanation is that, during autoencoder training, the encoder determines what information is preserved in the latent space. When text is injected at the encoder side, the latent representation is shaped to align with semantic cues, so the decoder can reconstruct text-relevant structures more faithfully from a coherent latent. If text is introduced only at the decoder, the latent remains video-only and the decoder must “correct” or hallucinate details from insufficient features, which can conflict with the latent and degrade both fidelity and perceptual quality. Injecting text on both sides over-conditions the model with similar guidance twice, making optimization harder without

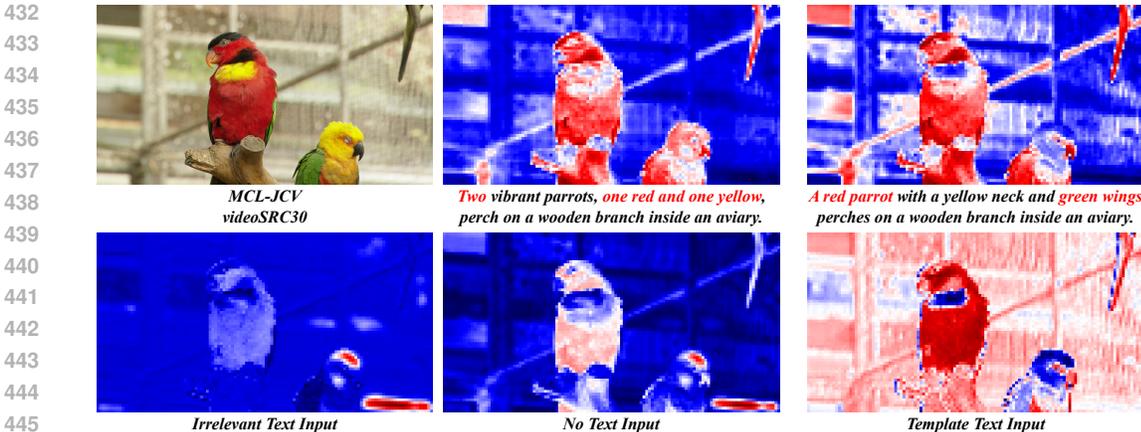


Figure 5: Visualization comparison of the stacked top 8 highest-entropy channels of encoder output y under different text inputs.

providing additional semantic information beyond the encoder-side guidance, which explains why Enc-T/Dec-T does not outperform the simpler Enc-T/Dec-V design.

In addition, we include two widely used *image*-compression backbones as autoencoder variants for video: LIC-TCM Liu et al. (2023) and ELIC He et al. (2022). These baselines illustrate the trade-off between reconstruction quality and efficiency in image-oriented designs. LIC-TCM delivers relatively strong reconstruction quality but incurs very high computational cost and long inference time, whereas ELIC is lightweight and fast but suffers from noticeably worse reconstruction quality. By contrast, our video-specific encoder–decoder is designed to explicitly exploit inter-frame relationships, which allows it to achieve both lower complexity and better reconstruction quality on video.

We also compare different text-integration strategies. As shown in Table 3, *Global Text* uses a single clip-level caption for all frames, and *Copy Text* simply replicates this caption at each time step. We also include a *Constant Text* variant, which keeps the MoVie-style architecture but replaces captions with the same fixed sentence across clips and frames to isolate meaningful text.

The results show that naively sharing one caption across frames is ineffective. Both Global Text (single-stage fusion with a clip-level embedding) and Constant Text (non-informative, fixed captions) exhibit clear drops in distortion and perceptual metrics, indicating that frame- or content-agnostic descriptions cannot match diverse frame content. Copy Text helps slightly but still lags behind our two-stage design. By first adapting the clip-level caption into *frame-specific* text features and then fusing them temporally, our approach exploits informative, content-dependent text and better captures frame-to-frame variations, yielding higher perceptual quality at similar complexity.

| Text | bpp | PSNR \uparrow | MS-SSIM \uparrow | FID \downarrow | LPIPS \downarrow |
|-----------------|--------------|-----------------|--------------------|------------------|--------------------|
| Ground Truth | 0.011 | 34.61 | 0.94930 | 6.731 | 0.19388 |
| Rephrased Text | 0.011 | 34.57 | 0.94928 | 6.731 | 0.19411 |
| Irrelevant Text | 0.035 | 29.35 | 0.92912 | 15.129 | 0.21280 |
| Without Text | 0.032 | 29.77 | 0.93068 | 14.694 | 0.21164 |
| Template Text | 0.011 | 34.26 | 0.93934 | 11.657 | 0.20340 |

Table 4: Robustness to textual perturbations. All results use the same model; bpp variations arise only from input text differences.

4.5 STABILITY EVALUATION

To assess the robustness of semantic guidance, we evaluate different textual inputs for the same video (Table 4). We use a factual caption as Ground Truth and GPT Hurst et al. (2024) to generate semantically similar alternatives (Rephrased Text). We also test two challenging settings: irrelevant descriptions (Irrelevant Text) and no descriptions (Without Text), where a simple fallback template is used as a generic substitute (Template Text).

As shown in Table 4, our model remains stable with semantically aligned text, showing robustness to paraphrasing. In contrast, irrelevant or missing text significantly degrades perceptual and fidelity

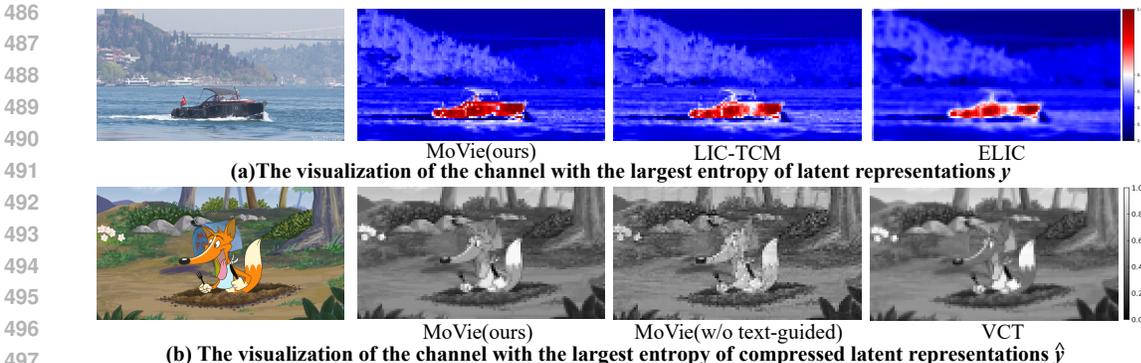


Figure 6: Visual comparison of different methods.

performance, highlighting the importance of semantic guidance. Our template serves as a fallback when accurate captions are unavailable, yielding better reconstructions than incorrect or missing text, though still worse than accurate descriptions. The specific template is provided in Appendix H.

As shown in Fig. 5, we visualize the stacked top-8 highest-entropy channels under different text inputs. When the description covers both birds, the model attends well to both targets; with a single-bird description, the other bird receives little attention. For irrelevant or no-text inputs, the attention distribution becomes diffuse and lacks focus. Template text provides some emphasis, but the contrast with the background remains limited.

4.6 ABLATION STUDIES

Fig. 6(a) shows the highest-entropy channel of y (mid-point 0.8). MoVie allocates more bits to the boat region, highlighting its text-guided focus on perceptually critical content, whereas LIC-TCM and ELIC exhibit more uniform, less semantic-aware distributions.

Fig. 6(b) provides a qualitative comparison of the compressed latent features \hat{y} . MoVie exhibits stronger, more localized activations around the fox, clearly outlining semantic regions such as the face and tail—demonstrating the benefit of text-guided fusion. In contrast, the w/o text-guided version of MoVie shows weaker and less focused responses, while VCT produces diffuse activations where the foreground blends into the background, indicating weaker semantic modeling.

Fig. 7 shows an ablation study of our modules. All variants are **fully trained** with entropy coding. “Baseline” denotes the Transformer-based VCT; “+VideoTCM” replaces the feature extractor with our Video Transformer-CNN Block; “+S-C” adds the Spatial-Channel Entropy Model; and “+Text” integrates the Extractor and Injector, yielding our final Text-guided VideoTCM.

5 CONCLUSION

We propose MoVie, a multimodal video compression framework that embeds textual semantics into a unified Video Transformer-CNN block. Through joint spatiotemporal modeling and two-stage text fusion, MoVie improves perceptual quality while reducing computation. A spatial-channel factorized entropy model and history-conditioned coding further enhance efficiency. Extensive experiments show that MoVie achieves an excellent trade-off between fidelity and perceptual quality, particularly at low bitrates, highlighting the potential of multimodal guidance. However, its performance is sensitive to text accuracy, and future work will focus on more robust semantic guidance to ensure stable compression under imperfect text conditions.

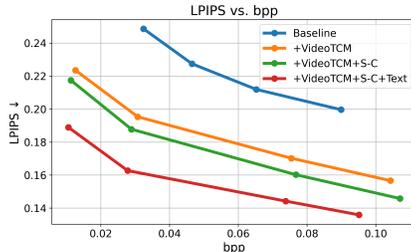


Figure 7: Ablation study on the UVG dataset, showing LPIPS across bitrates (bpp) to assess each component’s contribution to perceptual quality.

540 REPRODUCIBILITY STATEMENT

541

542 We have made every effort to ensure the reproducibility of our work. The complete source code,
543 including training and evaluation scripts, is provided in the supplementary material. Detailed de-
544 scriptions of the model architecture, training procedures, datasets, and hyperparameter settings are
545 included in the main text to facilitate replication. With these resources, we believe that all experi-
546 ments reported in this paper can be reliably reproduced.

547

548 REFERENCES

549

550 Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja. Compressai: a py-
551 torch library and evaluation platform for end-to-end compression research. *arXiv preprint*
552 *arXiv:2011.03029*, 2020.

553 Ashutosh Bhowan, Soham Mukherjee, Sean Yang, Shubham Chandak, Irena Fischer-Hwang, Kedar
554 Tatwawadi, Judith Fan, and Tsachy Weissman. Towards improved lossy image compression: Hu-
555 man image reconstruction with public-domain images. *arXiv preprint arXiv:1810.11137*, 2018.

556 Frank Bossen. Common test conditions and software reference configurations. In *3rd. JCT-VC*
557 *Meeting, Guangzhou, CN, October 2010*, 2010.

559 Benjamin Bross, Ye-Kui Wang, Yan Ye, Shan Liu, Jianle Chen, Gary J Sullivan, and Jens-Rainer
560 Ohm. Overview of the versatile video coding (vvc) standard and its applications. *IEEE Transac-*
561 *tions on Circuits and Systems for Video Technology*, 31(10):3736–3764, 2021.

562 Yi-Hsin Chen, Hong-Sheng Xie, Cheng-Wei Chen, Zong-Lin Gao, Martin Benjak, Wen-Hsiao Peng,
563 and Jörn Ostermann. Maskcrt: Masked conditional residual transformer for learned video com-
564 pression. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.

566 Ze Cui, Jing Wang, Shangyin Gao, Tiansheng Guo, Yihui Feng, and Bo Bai. Asymmetric gained
567 deep image compression with continuous rate adaptation. In *Proceedings of the IEEE/CVF Con-*
568 *ference on Computer Vision and Pattern Recognition*, pp. 10532–10541, 2021.

570 David Flynn, Detlev Marpe, Matteo Naccari, Tung Nguyen, Chris Rosewarne, Karl Sharman, Joel
571 Sole, and Jizheng Xu. Overview of the range extensions for the hevc standard: Tools, profiles,
572 and performance. *IEEE Transactions on Circuits and Systems for Video Technology*, 26(1):4–19,
573 2015.

574 Dailan He, Ziming Yang, Weikun Peng, Rui Ma, Hongwei Qin, and Yan Wang. Elic: Efficient
575 learned image compression with unevenly grouped space-channel contextual adaptive coding. In
576 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5718–
577 5727, 2022.

578 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-
579 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*
580 *arXiv:2410.21276*, 2024.

582 Xuhao Jiang, Weimin Tan, Tian Tan, Bo Yan, and Liquan Shen. Multi-modality deep network
583 for extreme learned image compression. In *Proceedings of the AAAI Conference on Artificial*
584 *Intelligence*, volume 37, pp. 1033–1041, 2023.

585 Hageong Lee, Minkyu Kim, Jun-Hyuk Kim, Seungeon Kim, Dokwan Oh, and Jaeho Lee. Neural
586 image compression with text-guided encoding for both pixel-level and perceptual fidelity. *arXiv*
587 *preprint arXiv:2403.02944*, 2024.

589 Eric Lei, Yiğit Berkay Uslu, Hamed Hassani, and Shirin Saeedi Bidokhti. Text+ sketch: Image
590 compression at ultra low rates. *arXiv preprint arXiv:2307.01944*, 2023.

591

592 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan
593 Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint*
arXiv:2408.03326, 2024a.

- 594 Jiahao Li, Bin Li, and Yan Lu. Deep contextual video compression. *Advances in Neural Information*
595 *Processing Systems*, 34:18114–18125, 2021.
- 596
- 597 Jiahao Li, Bin Li, and Yan Lu. Neural video compression with diverse contexts. In *Proceedings of*
598 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22616–22626, 2023.
- 599
- 600 Jiahao Li, Bin Li, and Yan Lu. Neural video compression with feature modulation. In *Proceedings*
601 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26099–26108,
602 2024b.
- 603 Jerry Liu, Shenlong Wang, Wei-Chiu Ma, Meet Shah, Rui Hu, Pranaab Dhawan, and Raquel Ur-
604 *tasun*. Conditional entropy coding for efficient video compression. In *European Conference on*
605 *Computer Vision*, pp. 453–468. Springer, 2020.
- 606
- 607 Jinming Liu, Heming Sun, and Jiro Katto. Learned image compression with mixed transformer-
608 *cnn architectures*. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
609 *recognition*, pp. 14388–14397, 2023.
- 610
- 611 Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin trans-
612 *former*. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
613 pp. 3202–3211, 2022.
- 614
- 615 Guo Lu, Wanli Ouyang, Dong Xu, Xiaoyun Zhang, Chunlei Cai, and Zhiyong Gao. Dvc: An
616 *end-to-end deep video compression framework*. In *Proceedings of the IEEE/CVF conference on*
617 *computer vision and pattern recognition*, pp. 11006–11015, 2019.
- 618
- 619 Fabian Mentzer, George Toderici, David Minnen, Sung-Jin Hwang, Sergi Caelles, Mario Lucic, and
620 *Eirikur Agustsson*. Vct: A video compression transformer. *arXiv preprint arXiv:2206.07307*,
621 2022.
- 622
- 623 Alexandre Mercat, Marko Viitanen, and Jarno Vanne. Uvg dataset: 50/120fps 4k sequences for
624 *video codec analysis and development*. In *Proceedings of the 11th ACM multimedia systems*
625 *conference*, pp. 297–302, 2020.
- 626
- 627 Zhihong Pan, Xin Zhou, and Hao Tian. Extreme generative image compression by learning text
628 *embedding from diffusion models*. *arXiv preprint arXiv:2211.07793*, 2022.
- 629
- 630 Jorge Pessoa, Helena Aidos, Pedro Tomás, and Mário AT Figueiredo. End-to-end learning of video
631 *compression using spatio-temporal autoencoders*. In *2020 IEEE Workshop on Signal Processing*
632 *Systems (SiPS)*, pp. 1–6. IEEE, 2020.
- 633
- 634 Shiyu Qin, Bin Chen, Yujun Huang, Baoyi An, Tao Dai, and Shu-Tao Xia. Perceptual image com-
635 *pression with cooperative cross-modal side information*. *arXiv preprint arXiv:2311.13847*, 2023.
- 636
- 637 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
638 *Girish Sastry, Amanda Aspell, Pamela Mishkin, Jack Clark, et al.* Learning transferable visual
639 *models from natural language supervision*. In *International conference on machine learning*, pp.
640 8748–8763. PmLR, 2021.
- 641
- 642 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
643 *resolution image synthesis with latent diffusion models*. In *Proceedings of the IEEE/CVF confer-*
644 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 645
- 646 Xihua Sheng, Jiahao Li, Bin Li, Li Li, Dong Liu, and Yan Lu. Temporal context mining for learned
647 *video compression*. *IEEE Transactions on Multimedia*, 25:7311–7322, 2022.
- 648
- 649 Xihua Sheng, Li Li, Dong Liu, and Houqiang Li. Spatial decomposition and temporal fusion based
650 *inter prediction for learned video compression*. *IEEE Transactions on Circuits and Systems for*
651 *Video Technology*, 34(7):6460–6473, 2024.
- 652
- 653 Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high
654 *efficiency video coding (hevc) standard*. *IEEE Transactions on circuits and systems for video*
655 *technology*, 22(12):1649–1668, 2012.

648 Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spa-
649 tiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international*
650 *conference on computer vision*, pp. 4489–4497, 2015.

651 Gregory K Wallace. The jpeg still picture compression standard. *Communications of the ACM*, 34
652 (4):30–44, 1991.

654 Rui Wan, Qi Zheng, and Yibo Fan. M3-cvc: Controllable video compression with multimodal
655 generative models. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech*
656 *and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2025.

657 Haiqiang Wang, Weihao Gan, Sudeng Hu, Joe Yuchieh Lin, Lina Jin, Longguang Song, Ping Wang,
658 Ioannis Katsavounidis, Anne Aaron, and C-C Jay Kuo. Mcl-jcv: a jnd-based h. 264/avc video
659 quality assessment dataset. In *2016 IEEE international conference on image processing (ICIP)*,
660 pp. 1509–1513. IEEE, 2016.

662 Tsachy Weissman. Toward textual transform coding. *IEEE BITS the Information Theory Magazine*,
663 3(2):32–40, 2023.

664 Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement
665 with task-oriented flow. *International Journal of Computer Vision*, 127(8):1106–1125, 2019.

667 Ren Yang, Radu Timofte, and Luc Van Gool. Perceptual learned video compression with recurrent
668 conditional gan. In *IJCAI*, pp. 1537–1544, 2022.

669 Chun Zhang, Heming Sun, and Jiro Katto. Flavc: Learned video compression with feature level
670 attention. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 28019–
671 28028, 2025.

672 Pingping Zhang, Jinlong Li, Kecheng Chen, Meng Wang, Long Xu, Haoliang Li, Nicu Sebe, Sam
673 Kwong, and Shiqi Wang. When video coding meets multimodal large language models: A unified
674 paradigm for video coding. *arXiv preprint arXiv:2408.08093*, 2024.

675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A APPENDIX

703
704 The appendix provides detailed descriptions of the network architecture, training strategy, dataset
705 usage, prompt construction, additional experimental results, comparisons under various traditional
706 codec settings, and the command-line configurations used for HM and VTM baselines.
707

708 B DATASET

709
710 We conduct training on the Vimeo-90k dataset Xue et al. (2019), which is distributed under the MIT
711 License. For evaluation, we use three publicly available and widely adopted datasets: the UVG
712 dataset Mercat et al. (2020) (BY-NC License), the MCL-JCV dataset Wang et al. (2016) (copyright
713 details are available online), and the standard HEVC test sequences Bossen (2010); Flynn et al.
714 (2015). All datasets are intended for academic research purposes and can be freely accessed on-
715 line. Furthermore, we have manually verified that none of the data contains personally identifiable
716 information or offensive content.
717

718 In the main paper, we report results on the UVG and MCL-JCV datasets. Specifically, Fig. 8 and
719 Table 6 present the rate-distortion curves and BD-rate results for the HEVC-B and HEVC-C datasets,
720 while Fig. 9 and Table 7 show the corresponding results on the HEVC-D and HEVC-E datasets.
721 We further present qualitative comparisons in Fig. 10 and Fig. 11 across different datasets. As
722 shown in Fig. 10, our method achieves better reconstruction in high-frequency regions, preserving
723 textures more faithfully. Moreover, Fig. 11 demonstrates that our method maintains visually pleasing
724 reconstruction quality even when operating at only half the bitrate of other methods.
725

726 Algorithm 1 GainUnit: Scale-conditioned Modulation

```

727 1: function GAIN INIT( $x, s$ )
728 2:    $g \leftarrow \exp(\log g_s)$ 
729 3:    $G \leftarrow \text{reshape}(g, [1, 1, C, 1, 1])$ 
730 4:   return  $x \cdot G$ 
731 5: end function
732 6: function INVERSE GAIN INIT( $x, s$ )
733 7:    $g \leftarrow \exp(-\log g_s)$ 
734 8:    $G \leftarrow \text{reshape}(g, [1, 1, C, 1, 1])$ 
735 9:   return  $x \cdot G$ 
736 10: end function

```

737 C RATE CONTROL

738
739 We implement a simple yet lightweight rate control mechanism by introducing a learnable Gain
740 Unit and Inverse Gain Unit as an additional modulation layer. The complete algorithm is provided
741 in Algorithm 1. We emphasize that this is a preliminary and straightforward implementation. Our
742 modular codec design allows for more advanced rate control strategies, which we leave for future
743 exploration.
744

745 Here, $\log g_s \in \mathbb{R}^C$ denotes a learnable gain vector indexed by scale s , and $x \in \mathbb{R}^{B \times D \times C \times H \times W}$.
746 The gain is broadcast over spatial and temporal dimensions after reshaping to $G \in \mathbb{R}^{1 \times 1 \times C \times 1 \times 1}$.
747

748 D TRADITIONAL CODEC

749
750 We used the official reference implementations of HM and VTM for traditional video compression:

- 751 • **HM**: version **18.0**, compiled on *Linux (GCC 13.2.0, 64-bit)* with *Range Extensions (RExt)*
752 enabled.
- 753 • **VTM**: version **23.11**, compiled on *Linux (GCC 13.2.0, 64-bit)* with *AVX2 SIMD* accelera-
754 tion.
755

756 These implementations are widely adopted as standard benchmarks for evaluating learned video
 757 compression frameworks, and all experiments were conducted under consistent compilation envi-
 758 ronments to ensure fair comparison.

759 We followed the *common test conditions (CTC)* specified by the JVET and JCT-VC standardization
 760 groups to ensure fair and reproducible comparison. The encoder configurations used for traditional
 761 codecs are as follows:

- 763 • **HM:**
 764 We used the `encoder_randomaccess_main.cfg` configuration file from HM-18.0.
- 765 • **VTM:**
 766 We used the `encoder_randomaccess_vtm.cfg` configuration file from VTM-23.11.

767 Each video was encoded with the following general command-line options:

```
768 EncoderAppStatic \
769   -c {config file}
770   --InputFile={input file}
771   --InputBitDepth=8
772   --OutputBitDepth=8
773   --FrameRate={frame rate}
774   --FramesToBeEncoded={frame number}
775   --SourceWidth={width} --SourceHeight={height}
776   --QP={qp}
```

778 **E TRAINING STRATEGY**

779 We adopt a three-stage training strategy for efficiency and stability.

- 782 • **Stage I:** Only the autoencoder is trained without the bitrate term, using a small number of
 783 frames to accelerate early convergence.
- 784 • **Stage II:** The entropy model is introduced and trained with all loss terms, using longer
 785 sequences for improved temporal modeling.
- 786 • **Stage III:** The entire model is jointly optimized with a reduced learning rate.

787 This strategy allows gradual integration of rate modeling while maintaining training stability.

| | Components trained | Loss | B | N_F | LR | Steps |
|-----|--------------------------|---------|-----|-------|--------|-------|
| 792 | Stage I Autoencoder | w/o r | 4 | 1 | $2E-4$ | 1M |
| 793 | Stage II Entropy Model | all | 4 | 7 | $2E-4$ | 250k |
| 794 | Stage III All Components | all | 4 | 7 | $5E-5$ | 500k |

795
 796 Table 5: We split training in three stages for training efficiency. r is bitrate, d is distortion, B is
 797 batch size, N_F the number of frames.

798
 799 **F PROMPT CONSTRUCTION**

800 We adopt LLaVA-Video Li et al. (2024a) (Qwen2 version) as our video-text generation backbone.
 801 For each short video clip from video dataset, we uniformly sample 64 frames and construct the
 802 following prompt:

```
803 The video lasts for  $X.XX$  seconds, and  $N$  frames are  

    804 uniformly sampled from it. These frames are located  

    805 at  $T_1s, \dots, T_Ns$ .  

    806 Please answer the following questions related to this  

    807 video.  

    808 Please describe this video in detail.
```

G SAMPLE TEXT DESCRIPTIONS FOR FIG. 6

To provide context for the visual results shown in Fig. 6, we include the corresponding automatically generated descriptions for each video segment:

(a) UVG – Bosphorus

The video features a serene scene of a dark-colored motorboat with a canopy, cruising on calm waters. The boat is adorned with a red flag on its stern. In the background, a large suspension bridge spans across the water, connecting two land masses covered in lush greenery and dotted with buildings. The sky is clear, suggesting a sunny day, and the water reflects the light, creating a tranquil atmosphere. Other boats are visible in the distance, adding to the sense of a bustling yet peaceful waterway.

(b) MCL-JCV – videoSRC20

The video features an animated orange fox with a white belly and a blue bib, sitting in a hole in the ground. The fox is holding a fork and knife, and appears to be looking around curiously. In the background, there is a mailbox labeled “MR FOX” and some greenery, including rocks and flowers. The fox then jumps out of the hole and onto the mailbox, causing it to wobble. The fox continues to jump around energetically, holding the utensils, and eventually runs off-screen. The scene transitions to the fox running through a forested area, still holding the fork and knife, and looking excited. The background includes trees, rocks, and patches of grass.

H CAPTION TEMPLATES FOR CONTROLLED EVALUATION

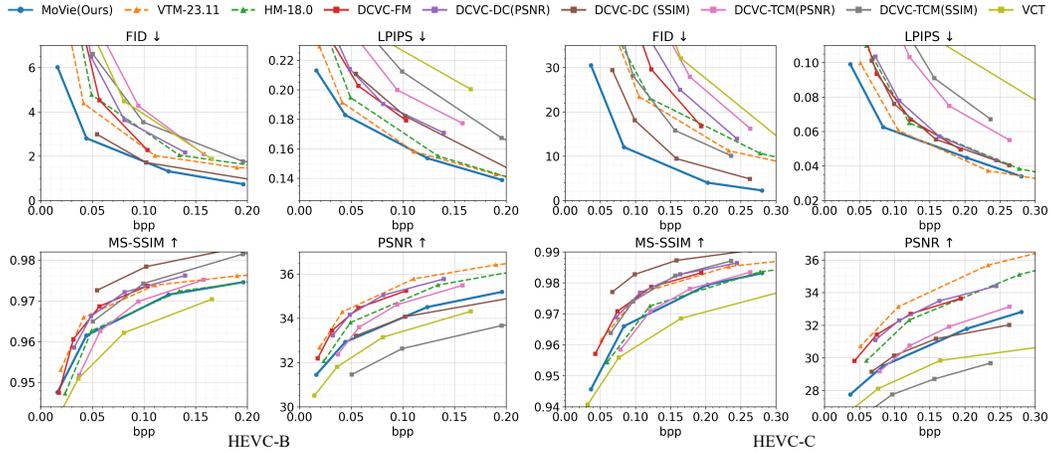
To evaluate the robustness of our method under standardized semantic guidance, we design a manually constructed caption template to simulate typical video descriptions in the absence of human annotations:

A broad, generic video scene showing one or more primary subjects (people, vehicles, or objects) moving through a typical indoor or outdoor environment under natural or artificial lighting, captured as a wide-to-medium shot with gentle handheld or stabilized camera motion (subtle pan/tilt/track), moderate global motion, and occasional local fast motion around the main subject; the composition centers the subject with clear separation from background, maintaining readable on-screen text, faces, logos, and motion boundaries as top visual priorities, while background textures and repetitive patterns are less critical; colors are neutral to slightly warm with balanced contrast, shadows are soft to moderate, and depth of field is shallow-to-medium to keep the subject prominent; the scene may include brief occlusions, partial reflections, or specular highlights, with mild motion blur on fast regions acceptable; overall pacing is steady, with a small number of salient moments (e.g., a gesture, a glance, a directional change) that should remain crisp and legible for downstream perception tasks.

I DETAILS ON LARGE LANGUAGE MODEL USAGE

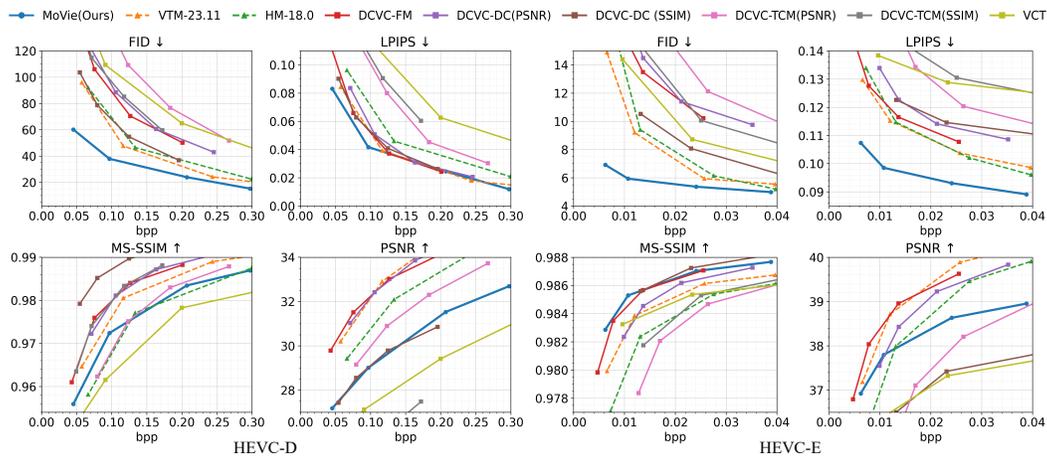
In this work, a Large Language Model (LLM) was used only as an auxiliary tool to polish writing, improve grammar, and refine expressions for clarity. It was not used to generate scientific content, design experiments, analyze results, or write substantive parts of the paper. All conceptual ideas, technical contributions, and experiment analyses were conducted entirely by the authors, with the LLM serving solely as a language aid rather than a content creator.

J VISUALIZATION RESULTS OF HEVC

Figure 8: Overall compression performance with RD curves (\downarrow lower is better, \uparrow higher is better).

| Method | MAC | HEVC-B | | | | HEVC-C | | | |
|---------------------|---------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| | | FID | LPIPS | MS-SSIM | PSNR | FID | LPIPS | MS-SSIM | PSNR |
| HM | / | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| VTM-23.11 | / | -23.07% | -14.83% | -34.08% | -32.13% | -9.86% | -20.86% | -31.38% | -30.33% |
| DCVC-FM | 2051.8 | 15.84% | 29.80% | -27.66% | -21.88% | 35.43% | 4.88% | -32.04% | -16.23% |
| DCVC-DC (PSNR) | 2260.8 | 33.13% | 44.20% | -29.81% | -17.19% | 50.74% | 13.56% | -32.98% | -7.53% |
| DCVC-DC (SSIM) | 2260.8 | -41.48% | 30.73% | -67.51% | 95.43% | -39.55% | 3.40% | -61.22% | 77.68% |
| DCVC-TCM (PSNR) | 2322.1 | 75.69% | 115.06% | 11.16% | 28.00% | 85.23% | 84.58% | 8.98% | 62.45% |
| DCVC-TCM (SSIM) | 2322.1 | 37.07% | 117.46% | -37.07% | NaN | -2.49% | 110.31% | -29.56% | NaN |
| VCT | 3115.9 | 41.06% | 239.06% | 57.41% | 122.59% | 78.81% | 205.07% | 57.96% | 232.82% |
| MoVie (Ours) | 1144.1 | -54.35% | -28.69% | -7.61% | 48.62% | -64.17% | -32.66% | -2.72% | 86.00% |

Table 6: Comparison on HEVC-B and HEVC-C in terms of BD-Rate (%) using HM-18.0 as the common anchor. BD-Rate is computed as $-M$ for lower-is-better metrics. MAC is measured in kMACs/pixel. NaN indicates failure in BD-rate computation due to excessive performance gap with the anchor.

Figure 9: Overall compression performance with RD curves (\downarrow lower is better, \uparrow higher is better).

| Method | MAC | HEVC-D | | | | HEVC-E | | | |
|---------------------|---------------|----------------|----------------|----------------|---------|----------------|----------------|----------------|---------|
| | | FID | LPIPS | MS-SSIM | PSNR | FID | LPIPS | MS-SSIM | PSNR |
| HM | / | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| VTM-23.11 | / | -10.25% | -25.56% | -29.53% | -28.87% | -11.75% | 0.27% | -32.42% | -33.45% |
| DCVC-FM | 2051.8 | 51.34% | -20.29% | -39.69% | -30.91% | 53.64% | 5.26% | -53.42% | -38.44% |
| DCVC-DC (PSNR) | 2260.8 | 63.80% | -12.01% | -39.19% | -24.57% | 82.33% | 48.50% | -37.97% | -11.30% |
| DCVC-DC (SSIM) | 2260.8 | 9.26% | -17.70% | -66.70% | 80.26% | 23.13% | 116.41% | -62.41% | 213.39% |
| DCVC-TCM (PSNR) | 2322.1 | 136.33% | 49.48% | 0.72% | 28.50% | 164.01% | 154.28% | 32.91% | 82.85% |
| DCVC-TCM (SSIM) | 2322.1 | 72.79% | 76.59% | -38.85% | NaN | 90.99% | 293.92% | -3.74% | NaN |
| VCT | 3115.9 | 103.36% | 104.01% | 39.25% | 193.70% | 48.52% | 277.70% | -18.25% | 133.63% |
| MoVie (Ours) | 1144.1 | -41.46% | -34.04% | -9.27% | 77.72% | -58.44% | -63.09% | -60.31% | -0.62% |

Table 7: Comparison on HEVC-D and HEVC-E in terms of BD-Rate (%) using HM-18.0 as the common anchor. BD-Rate is computed as $-M$ for lower-is-better metrics. MAC is measured in kMACs/pixel. NaN indicates failure in BD-rate computation due to excessive performance gap with the anchor.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

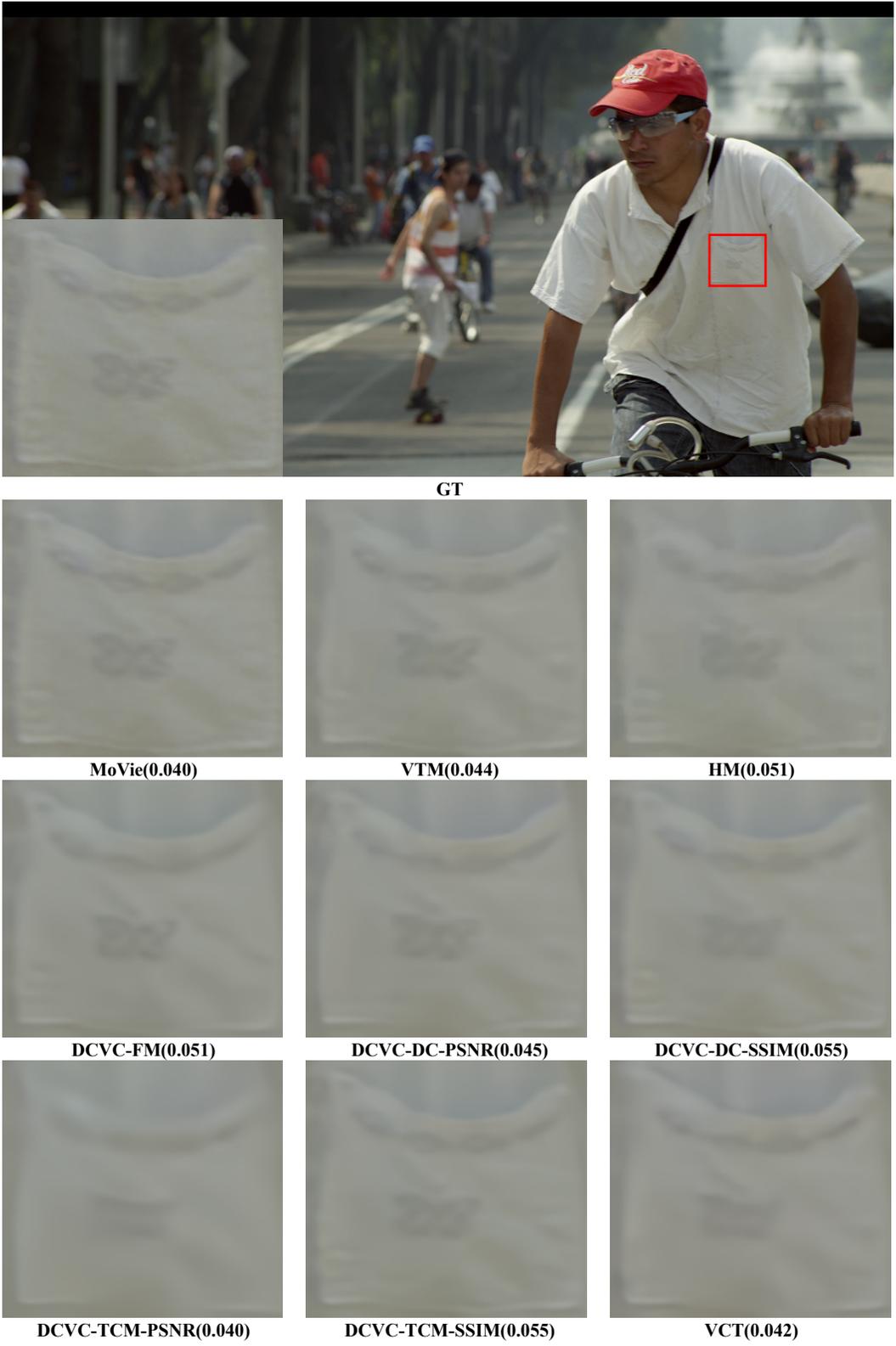


Figure 10: Subjective quality comparison on the MCL-JCV *videoSRC02* sequence.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

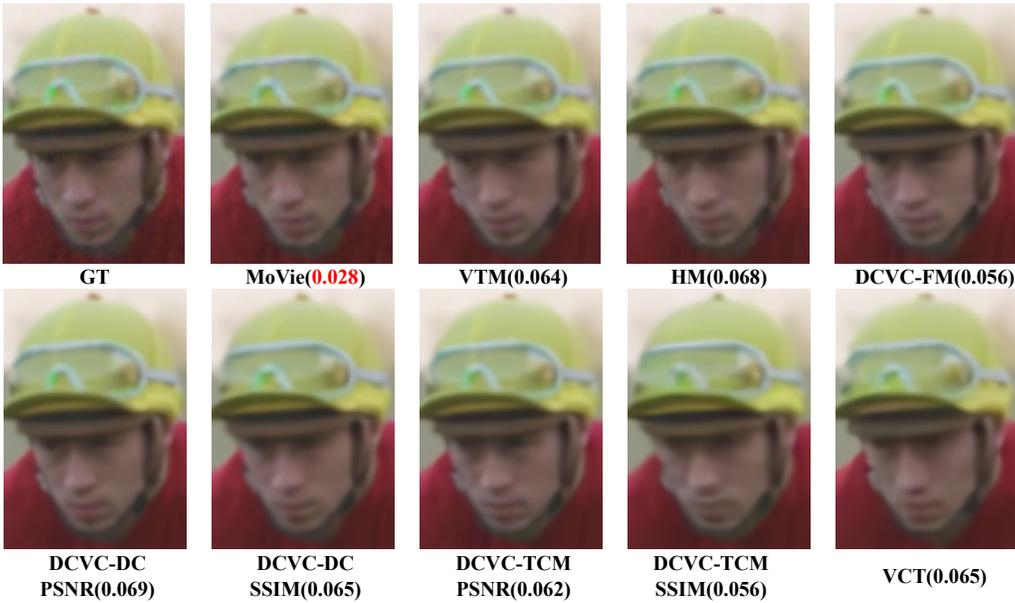


Figure 11: Subjective quality comparison on the UVG *Jockey* sequence.