

TVC: Tokenized Video Compression with Ultra-Low Bitrate

Lebin Zhou
Santa Clara University
Santa Clara, CA, USA
lzhou@scu.edu

Cihan Ruan
Santa Clara University
Santa Clara, CA, USA
cruan@scu.edu

Nam Ling
Santa Clara University
Santa Clara, CA, USA
nling@scu.edu

Wei Wang
Futurewei Technologies, Inc.
San Jose, CA, USA
rickweiwang@futurewei.com

Wei Jiang
Futurewei Technologies, Inc.
San Jose, CA, USA
wjiang@futurewei.com

Abstract

Tokenized visual representations have shown great promise in image compression, yet their extension to video remains under-explored due to the challenges posed by complex temporal dynamics and stringent bitrate constraints. In this paper, we propose **Tokenized Video Compression (TVC)**, the first token-based dual-stream video compression framework designed to operate effectively at ultra-low bitrates. TVC leverages the powerful Cosmos video tokenizer to extract both discrete and continuous token streams. The discrete tokens (*i.e.*, code maps generated by FSQ) are partially masked using a strategic masking scheme, then compressed losslessly with a discrete checkerboard context model to reduce transmission overhead. The masked tokens are reconstructed by a decoder-only transformer with spatiotemporal token prediction. Meanwhile, the continuous tokens, produced via an autoencoder (AE), are quantized and compressed using a continuous checkerboard context model, providing complementary continuous information at ultra-low bitrate. At the Decoder side, both streams are fused using ControlNet, with multi-scale hierarchical integration to ensure high perceptual quality alongside strong fidelity in reconstruction. This work mitigates the long-standing skepticism about the practicality of tokenized video compression and opens up new avenues for semantics-aware, token-native video compression.

CCS Concepts

• **Computing methodologies** → **Video compression**; *Neural networks*; Computer vision; • **Information systems** → *Multimedia information systems*.

Keywords

Video Compression, Dual-Stream Architecture, Discrete-Continuous Tokenization, Neural Codecs, Deep Learning

1 Introduction

Generative visual priors, *a.k.a.* visual tokens, learned from a massive amount of images and videos to model distributions of the entire visual space, have shown great success in various restoration tasks, such as super-resolution [5, 21, 28], quality enhancement [35], image compression [8, 24], *etc.* The key is the learned Tokenized Visual Representation (TVR), into which the input visual signals are transformed by a tokenizer, and based on which a pixel decoder reconstructs the output visual signals. The tokenizer, TVR, and

pixel decoder are optimized end-to-end to balance the efficiency of representation and the quality of the reconstruction.

There are two types of TVR: continuous TVR (C-TVR) and discrete TVR (D-TVR). The discrete tokenizers, *e.g.*, VQGAN (Vector Quantized Generative Adversarial Network) [8, 19] and FSQ (Finite Scalar Quantization) [26, 29], transform inputs into discretized sequences of latent codes by learning a latent visual space that is partitioned into cells with unequal volumes, as illustrated in Figure 2 (A-B). Continuous tokenizers, *e.g.*, Autoencoder (AE), Variational Autoencoder (VAE) [16, 29], transform inputs into a learned latent space as continuous latent features, which are then quantized for transmission. The quantization process partitions the latent space of C-TVR into equal cells, as illustrated in Figure 2 (C).

C-TVR has been largely studied for Learned Image Compression (LIC) [6, 13], since the equal partition of the visual space preserves the local sensitivity of the visual features. Similar/different visual appearances are assigned to similar/different cells in latent space, and the granularity of partition (*i.e.*, scale of quantization) determines the mapping sensitivity, hence the fidelity of the reconstruction to the input. This aligns with the conventional compression target of balancing rate-distortion (RD). However, at ultra-low bitrates, the performance suffers from severe information loss caused by heavy quantization, as is the case for conventional compression methods.

D-TVR has gained increasing attention in recent years for low-bitrate LIC [8, 24] due to its ability to generate reconstruction with high perceptual quality against input degradations. The learned unequal partitions focus on modeling the general salient visual cues described by dense areas in the visual space, while allocating sparse cells over less representative areas. Such a “smart” token allocation is especially useful for ultra-low-bitrate scenarios, where outputs can be reconstructed with high perceptual quality by using only a small number of salient tokens.

Neither D-TVR nor C-TVR has been applied to compress videos, although there is a much stronger need for ultra-low-bitrate video compression compared to images, since efficient video transmission has become a performance bottleneck for a vast amount of applications in entertainment, gaming, AR/VR, *etc.* Previous video compression methods, including Learned Video Compression (LVC) based on neural networks [10, 22], follow a pipeline based on motion prediction and residual coding, which was designed decades ago, particularly for traditional video coding [3, 17, 31]. We strongly believe that the overall performance of LVC, in terms of both computation and bitrate, can be largely improved by redesigning a

holistic framework based on TVRs. However, there is a fundamental dilemma about using TVR for LVC. On the one hand, to reliably represent the complex spatial-temporal visual content in general videos, D-TVR usually requires a very large number of spatial-temporal tokens (tens or even hundreds of millions), and C-TVR usually requires a latent space with very high dimensionality. The number of bits to represent such token indices or high-dimensional latent features can be prohibitively high, defeating the purpose of video compression.

In this paper, we introduce a novel dual-stream Tokenized Video Compression (TVC) framework, systematically designed for TVR-based LVC, as illustrated in Figure 3. Our approach addresses the aforementioned challenges from several aspects and integrates both D-TVR and C-TVR to leverage their complementary strengths. The key contributions of our work are summarized as follows.

- We aim for both high fidelity and high perceptual quality in reconstructed videos at ultra-low bitrates. Although D-TVR excels at high-quality visual generation, it can result in loss of fine-fidelity details. Visual differences in the sparse regions of the latent space can be overlooked, while inauthentic details can be hallucinated in denser areas (as illustrated in Figure 1). By integrating the fidelity-preserving C-TVR stream, we achieve a balanced outcome that ensures both high fidelity and enhanced perceptual quality.
- We take advantage of the high redundancy in videos to largely reduce the number of transferred tokens by employing the token prediction mechanism. As shown by MAGE [19], the token space creates an efficient probability space for prediction, where masked token prediction can effectively recover image content using spatial contextual information. Compared to images, token-based modeling is even more effective for videos, since semantic patterns tend to persist between frames [7]. Tokens capture content at a conceptual level where semantic consistency naturally emerges, *i.e.*, similar visual concepts in adjacent frames often map to identical or closely related tokens, creating temporal predictable patterns. Our TVC exploits this semantic redundancy through spatiotemporal token prediction to achieve efficient compression.
- We build upon Nvidia’s Cosmos Tokenizer [29], which extracts compact discrete and continuous token streams through wavelet transforms, spatiotemporal 3D convolution, and spatiotemporal causal self-attention. To effectively further compress these dual streams, we introduce two checkerboard context models (CCM). For discrete tokens, we apply 3D masking to drop off masked tokens and losslessly compress the visible tokens using a discrete CCM. For continuous tokens, we perform quantization with a continuous CCM for effective lossy compression. This design significantly reduces both spatial and temporal redundancy to achieve ultra-low bitrates.
- Inspired by the residual conditioning mechanism of ControlNet [34], we design a hierarchical pixel decoder that fuses D-TVR and C-TVR streams through multi-scale residual injection. Specifically, the C-TVR latents are decoded

via a cascade of Cosmos pixel decoder layers, whose intermediate features are injected into the corresponding layers of the D-TVR pixel decoder in a residual manner. This dual-stream fusion architecture enhances both pixel-level fidelity and semantic consistency across video frames.

Compared to conventional LVC methods like [10, 22], our TVC approach significantly improves ultra-low bitrate performance with considerably reduced computational overhead. Likewise, compared to the straightforward extension of dual-stream LIC methods such as [23], our TVC consistently outperforms by leveraging effective video tokenizers and enhancing the reduction of spatio-temporal redundancy through advanced spatio-temporal token prediction.

Operating in the token space, we formulate video compression as a token selection and prediction problem. Our key insight is that, with a structured and temporally causal token space like Cosmos, ultra-low bitrate video compression becomes not only feasible but highly effective. Unlike prior methods that focus on compressing code indices or residuals, we directly compress the tokens themselves, yielding a highly compact and semantically consistent representation across frames. Our approach mitigates the long-standing dilemma about the practicality of tokenized video compression and opens up new opportunities for learned video compression.

It is worth mentioning that our TVC framework is designed with flexibility in mind, allowing seamless integration with different video tokenizers beyond Cosmos. This adaptability stems from our generic token prediction and fusion mechanisms, which are agnostic to any specific tokenizer architecture.

We conduct extensive experiments and ablation studies on multiple benchmark datasets. As the first tokenized video compression system operating at ultra-low bitrates, TVC achieves an LPIPS of around 0.30 without relying on conventional optical flow or MEM (motion estimation and motion compensation). At such extreme bitrate levels, pixel-level distortion metrics like PSNR and perceptual quality metrics like LPIPS often diverge, highlighting the trade-off between structural fidelity and visual realism. In this context, LPIPS serves as a more reliable indicator of perceptual quality, aligning more closely with human visual preferences. Our results exhibit strong perceptual consistency across frames, demonstrating the potential of token-based modeling for semantic-level video reconstruction under extreme compression constraints.

2 Related Works

2.1 Tokenized Visual Representation

The core of image/video generation is the learned TVR by modeling the probability distributions of the visual space. For D-TVR, the discrete tokenizer, *e.g.*, VQGAN [8] or FSQ [26], transforms input visual signals into discretized sequences of latent codes. For C-TVR, the continuous tokenizer, *e.g.*, VAE [16, 29], transforms inputs into a learned latent space as continuous latent features. Then a pixel decoder transforms TVR back to the pixel space to reconstruct the visual signal. The tokenizer, TVR, and pixel decoder are optimized end-to-end to balance the efficiency of representation and the reconstruction quality.

The spatiotemporal visual content in videos is inherently complex, and most existing tokenizers focus solely on spatial image content [8, 16, 26]. Even for static images, developing an efficient



Figure 1: Visual comparison between single-stream and dual-stream reconstruction at extremely low bitrates. C-TVR produces overly smooth and blurry outputs as a result of heavy quantization, while D-TVR often hallucinates inauthentic details derived from learned salient codewords. Our dual-stream design combines their complementary strengths and faithfully recovers both rich textures and structural details.

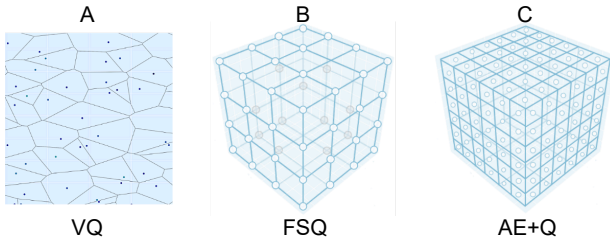


Figure 2: Visualization of quantization strategies of different TVRs. Each cell represents a token, with the dot being its centroid (latent feature). The uniform quantization of AE+Q results in equal partition of the latent space. In contrast, the non-uniform partition of VQ and FSQ emphasizes frequent visual patterns and enables adaptive representation, leading to higher perceptual quality under extreme compression.

class-agnostic tokenizer capable of handling diverse image categories remains challenging. Improvements have been made through class-guided tokenizers. For instance, AdaCode [21] introduces a set of basis codebooks to capture semantic-specific visual details and employs a weight map to effectively combine these basis codebooks.

However, directly applying image-based TVR to tokenize individual video frames introduces temporal inconsistencies, resulting in jittering and flickering artifacts. To address this, video-based TVRs [29, 33] have been developed by training on extensive datasets comprising tens of thousands or even millions of hours of video content. Given the heightened complexity of spatiotemporal visual content, video-based TVRs typically require a significantly larger number of spatiotemporal tokens or a highly dimensional latent space for accurate modeling compared to their image-based counterparts.

2.2 Tokenization for Image Compression

Since the pioneer work of Ballé et al. [1], C-TVR based on VAE has emerged as a prominent approach in LIC. In this method, a continuous tokenizer transforms an image into a latent feature, which is then processed through traditional quantization and entropy coding to produce a compact bitstream with continuous values. The

Decoder reconstructs the image by applying conventional entropy decoding and dequantization to recover the degraded latent feature. While significant advancements have been made in improving the entropy model to mitigate information loss during quantization, achieving high-quality reconstruction at ultra-low bitrates remains challenging, because aggressive quantization severely degrades the recovered latent feature, often resulting in artifacts like blurring, blocky effects, etc.

The idea of D-TVR naturally aligns with compression tasks, where an image is mapped to a sequence of discretized latent codes. Each code corresponds to the index of a partitioned cell, *i.e.*, a visual token, in the learned latent space. Using the index, the Decoder retrieves the latent feature of the corresponding cell for reconstruction. Compared to the original embedded feature, the retrieved latent feature of the mapped token is a quantized approximation. The granularity of this partition, *i.e.* the number of tokens, directly determines the tradeoff between bitrate and reconstruction quality. A finer partition with more tokens allows for higher reconstruction quality at the cost of increased bitrate.

D-TVR offers two notable advantages for compression. First, the use of integer indices is highly efficient for data transmission, ensuring robustness across different platforms and minimizing computation mismatches between sender and receiver devices. Second, D-TVR can achieve high perceptual quality in reconstructed content by leveraging high-quality tokens. Even when the input quality is degraded, retrieving a set of visual tokens similar to those representing a high-quality input can significantly enhance the reconstruction outcome. This capability is particularly advantageous in ultra-low bitrate scenarios, where C-TVR or traditional methods often struggle to deliver high-quality results.

Therefore, D-TVR has been used recently for ultra-low bitrate LIC. Since learning a universally rich visual codebook capable of capturing diverse image content remains challenging, the direct application of VQGAN to LIC [25] has demonstrated improvements in perceptual quality at extremely low bitrates. To enhance performance across a broader range of bitrates, [12] introduced multiple visual codebooks that capture class-specific visual details [21], combined with a weight masking mechanism to mitigate the increased

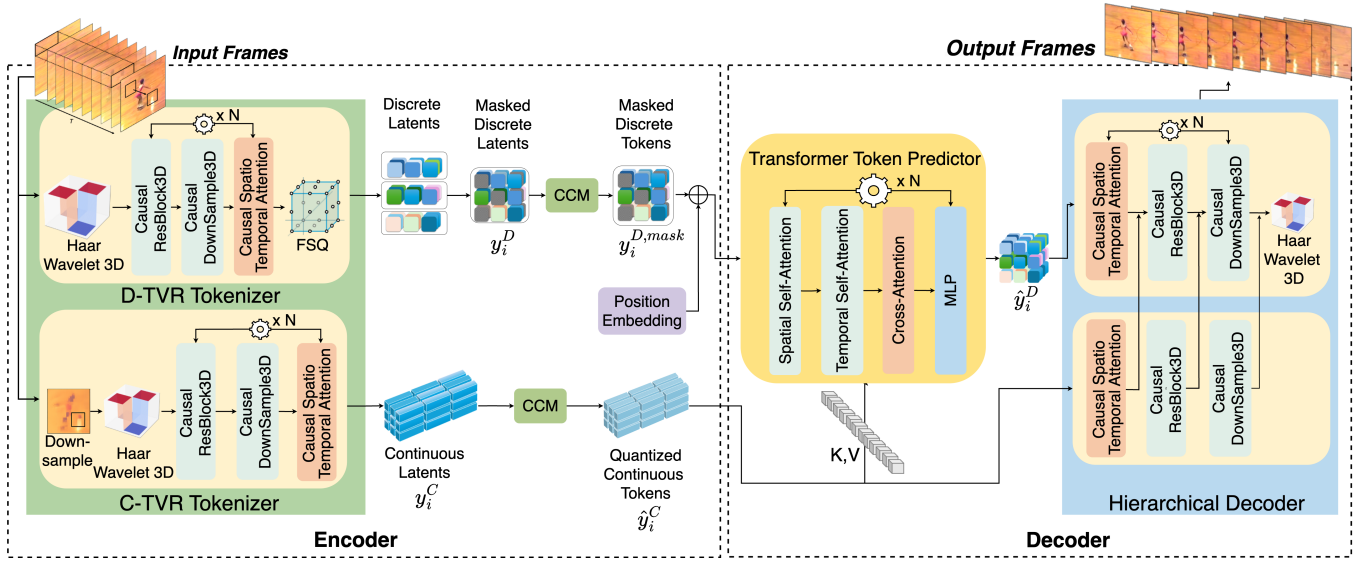


Figure 3: Overview of the proposed TVC framework. The tokenizer encoder extracts discrete and continuous token streams from each input GoP. Discrete tokens are masked and compressed using a CCM (checkerboard context model), and continuous tokens are quantized and also compressed with CCM. At the Decoder side, masked discrete tokens are reconstructed via a transformer with cross-attention to the continuous stream. Both streams are fused via a ControlNet-style hierarchical pixel decoder for final frame reconstruction.

bitrate associated with transmitting the additional weights for fusing multiple codebooks. Despite the advancements, D-TVR alone still faces the inherent fidelity loss issue, particularly at ultra-low bitrates. In such cases, a smaller codebook excessively compresses the visual space, causing different images to be treated as variations of one another and mapped to the same set of codewords. This results in reconstructions that, while potentially visually appealing, may lack pixel-level — and even semantic-level — accuracy, failing to faithfully reproduce the original input.

To overcome the inherent limitations of both the C-TVR and D-TVR methods, researchers have proposed dual-stream frameworks that leverage their complementary strengths. HybridFlow [24] exemplified this approach by integrating a discrete codebook-based stream with a continuous feature stream, effectively balancing perceptual quality and fidelity. HDCompression [23] further advanced this idea by introducing diffusion models into the dual-stream architecture, enhancing the performance in image compression.

Building upon these advancements, we extend the dual-stream framework to the video domain, addressing the unique challenges introduced by the temporal dimension in our tokenization-based video compression approach. This adaptation not only improves spatiotemporal redundancy reduction, but also ensures robust, balanced performance in ultra-low bitrate scenarios.

2.3 Tokenization for Video Compression

In contrast to static images, the application of tokenization to video data is relatively underexplored. Although methods like MagViT [19] and TimeSformer [2] have extended token-based architectures to video generation and high-level recognition tasks, no prior work has investigated token-based video compression.

This may be due to the inherent challenges mentioned above in applying tokenization to video data in Section 2.1. Compared to static images, the complex spatiotemporal visual content in videos demands either a significantly larger number of spatiotemporal tokens or a highly dimensional latent space for accurate modeling. This increased complexity substantially increases the number of bits required to effectively represent video-based TVR.

To effectively address this challenge, we integrate the state-of-the-art (SOTA) Cosmos tokenizer [29], which features an efficient temporally causal design. Furthermore, we introduce a masked token and prediction mechanism that aligns with Cosmos’s temporally causal structure, effectively exploiting visual redundancy in videos for improved compression efficiency.

3 Methodology: Tokenized Video Compression

An overview of the full TVC pipeline is shown in Figure 3. We formulate video compression as a token-level selection and prediction problem and tackle the challenges of ultra-low bitrate scenarios through four key components:

- (A) A dual-stream tokenizer that leverages Cosmos’s spatiotemporal tokenizer to extract discrete and continuous token representations;
- (B) Dedicated, entropy-efficient token compression pipelines for each stream, utilizing checkerboard context models (CCM);
- (C) A masked token prediction mechanism based on spatiotemporal context, implemented through a lightweight Transformer;
- (D) A multi-scale hierarchical pixel decoder that fuses the discrete and continuous streams to harness their complementary strengths.

This architecture enables the reconstruction of high-quality videos from highly sparse token representations, delivering efficient compression while preserving visual integrity and temporal consistency.

3.1 Discrete and Continuous Tokenizers

As shown in Figure 3, our TVC framework tokenizes each video at the granularity of a Group of Pictures (GoP), allowing the model to capture both intra-frame and inter-frame dependencies. To represent videos compactly yet expressively, we extract two types of token streams from each GoP: a discrete token stream that encodes high-level semantics, and a continuous stream that preserves fine-grained spatial fidelity. Both are derived from the Cosmos tokenizer, which features a pretrained encoder-decoder architecture capable of producing both discrete and continuous token outputs. This dual-stream representation forms the foundation of our efficient and semantically rich video compression pipeline.

Following the notations in Cosmos [29], let $\mathbf{x}_i \in \mathbb{R}^{(1+T) \times H \times W \times 3}$ denote the i -th GoP, where H , W , T are height, width, and frame number. The tokenizer encoder utilizes wavelet transforms and spatiotemporal factorized 3D convolution to transform \mathbf{x}_i into a compact latent space. The spatiotemporal factorized causal self-attention mechanism effectively captures long-range spatiotemporal dependencies, resulting in a continuous latent token $y_i^C \in \mathbb{R}^{(HT') \times H' \times W' \times d_C}$ retaining rich contextual information of \mathbf{x}_i . To achieve a more compact representation in the discrete space, the discrete tokenizer incorporates the FSQ to further quantize the continuous feature into a discrete token map $y_i^D \in \mathbb{R}^{(HT') \times H' \times W' \times d_D}$. The pixel decoder leverages both discrete and continuous tokens to reconstruct a high-quality GoP $\hat{\mathbf{x}}_i \in \mathbb{R}^{(1+T) \times H \times W \times 3}$ that closely approximates the original \mathbf{x}_i .

3.2 Token Compression and Decompression

To enable ultra-low bitrate transmission, both the continuous latents y_i^C and discrete code maps y_i^D are compressed via CCMs and entropy coding.

3.2.1 Continuous Stream. We adopt a quantization-plus-entropy model pipeline based on the continuous CCM [9]. First, the i -th GoP \mathbf{x}_i is spatially downsampled to further reduce the size of the resulting continuous latent token y_i^C . It is quantized to \hat{y}_i^C via uniform quantization, and \hat{y}_i^C is compressed by an arithmetic encoder. The CCM predicts the Gaussian distribution parameters μ, σ for each token using both decoded context and a learned hyperprior \hat{z}_i , enabling efficient entropy modeling.

3.2.2 Discrete Stream. For the discrete token stream y_i^D , we compress the code map generated by FSQ rather than raw codebook indices, which offers a narrower dynamic range and thus better entropy coding precision. Discrete stream compression employs a CCM with a 3D spatiotemporal masking strategy to enable partial token transmission, *i.e.*, only visible tokens are transferred.

During training, random masks are sampled as in [20]; at inference, a deterministic fixed mask is used to ensure inference consistency. The fixed mask is generated by unmasking the first n_{visible} tokens in every predefined interval along the flattened token sequence, with the overall mask rate defined as:

$$\text{Mask Rate} = 1 - (n_{\text{visible}} / \text{Mask Interval}). \quad (1)$$

This 1D mask is reshaped into a 3D binary map $m \in \{0, 1\}^{(1+T') \times H' \times W'}$. The masked tokens are dropped and the visible ones are encoded via a CCM-based arithmetic encoder. As y_i^D are integers from FSQ, this

compression process of the visible tokens is lossless. The discarded tokens are later predicted in the Decoder.

3.2.3 Masking Consistency and Prediction Alignment. To maintain semantic and structural alignment between the Encoder and Decoder, we adopt two consistency strategies. First, instead of using fixed constants as mask labels [19, 24], we dynamically initialize masked positions using values from the continuous stream \hat{y}_i^C . This improves both entropy modeling and Transformer-based prediction. Second, we apply the same 3D mask to the Encoder-side input of the CCM, preserving spatial layout even for discarded tokens. This ensures alignment across the compression and reconstruction pipelines. Together, these techniques improve prediction fidelity and reduce mismatch in downstream token recovery.

3.3 Masked Token Prediction

The prediction of masked discrete tokens is performed by a decoder-only Transformer, as shown in Figure 3. To facilitate consistent spatiotemporal representation across token streams, we deliberately reapply the Cosmos tokenizer encoder to the discrete tokens prior to the Transformer prediction. This shallow encoder pass acts as a unified projection layer, aligning the structural characteristics of discrete tokens with those of the continuous stream. While seemingly redundant, this step ensures architectural symmetry, stabilizes early-stage Transformer training, and retains the inductive bias embedded in the original tokenization process.

For compatibility with standard token prediction tasks, we convert the FSQ-generated code map back to its corresponding index map using the implicit codebook. Given a partially observed index map $y_i^{D(\text{index}), \text{mask}}$, where masked tokens have been dropped during compression, the Transformer reconstructs the complete token sequence $\hat{y}_i^{D(\text{index})}$ by attending to both intra-stream and cross-stream contexts:

$$\hat{y}_i^{D(\text{index})} = \text{Transformer}(Q = y_i^{D(\text{index}), \text{mask}}, K = \hat{y}_i^C, V = \hat{y}_i^C). \quad (2)$$

The Transformer comprises stacked spatiotemporal attention blocks, with learnable spatial and temporal positional embeddings to preserve the structure of token layouts. These layers capture both intra-frame dependencies and inter-frame correlations, enabling context-aware token reconstruction. Additionally, the cross-attention module leverages the continuous latent \hat{y}_i^C as Key/Value (K,V) inputs, providing guidance to improve prediction accuracy.

3.4 Hierarchical Decoder Fusion

Inspired by the residual conditioning mechanism of ControlNet [34], our pixel decoder exploits the continuous token stream as a conditional control branch to guide the Cosmos tokenizer’s pixel decoder for video reconstruction. In particular, features extracted from the continuous stream are injected into the main decoding path in a residual form at multiple spatial resolutions and are progressively fused with the discrete stream pixel decoder. This multi-scale residual fusion mechanism enables effective cross-stream information interaction, enhancing the pixel decoder’s contextual awareness. As a result, the model can better leverage the structurally preserving guidance from the continuous stream, leading to improved reconstruction accuracy and structural consistency.

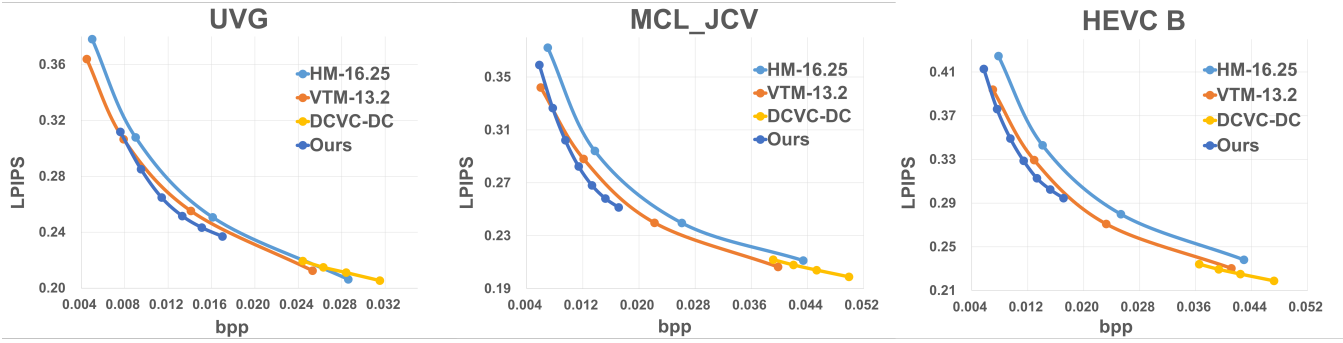


Figure 4: Rate-distortion (R-D) comparison (in terms of LPIPS; lower is better) across the UVG, MCL-JCV, and HEVC Class B datasets. Our TVC clearly outperforms conventional codecs such as VTM and HM in most cases. While existing neural codecs like DCVC-DC operate in a higher bitrate range, TVC is the first neural codec to achieve effective video reconstruction at ultra-low bitrates in the RGB colorspace, setting a new benchmark in this challenging scenario.

3.5 Training Pipeline

We adopt a four-stage training procedure to progressively optimize different components of the TVC framework.

1. Pretraining. We use the pretrained Cosmos discrete and continuous tokenizers [29], operating with $T=8$, $H/H'=16$, $W/W'=16$, to extract y_i^D and y_i^C . These tokenizers are kept frozen throughout the entire training process.

2. Context model training. The CCMs for both streams are trained using an entropy-only loss:

$$\mathcal{L}_{\text{bpp}} = \mathbb{E}_{x \sim p_x} \left[-\log_2 p_{\hat{y}|\hat{z}}(\hat{y}_i|\hat{z}_i) - \log_2 p_{\hat{z}}(\hat{z}_i) \right]. \quad (3)$$

This objective corresponds to the expected bit cost under the predicted Gaussian distribution. No distortion loss is applied at this stage, as the tokenizers are frozen.

3. Transformer training. The Transformer decoder for reconstructing masked discrete tokens is trained based on a negative log-likelihood loss that is computed only over masked positions:

$$\mathcal{L}_{\text{recon}} = -\frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \log p_{\theta}(y_i^D | y_i^{D,\text{mask}}, \hat{y}_i^C), \quad (4)$$

where \mathcal{M} contains masked token indices. Dropped tokens are recovered using both intra-stream context and cross-stream guidance.

4. Pixel decoder finetuning. Finally, we fine-tune the ControlNet-style pixel decoder for video reconstruction. Features from the C-TVR stream are injected into the D-TVR decoder in a residual manner at multiple scales. The training loss combines pixel-wise and perceptual terms:

$$\mathcal{L}_{\text{pixel}} = \lambda_1 \cdot \mathcal{L}_{\text{L1}}(x, \hat{x}) + \lambda_2 \cdot \mathcal{L}_{\text{perc}}(x, \hat{x}). \quad (5)$$

This stage enhances visual fidelity by aligning reconstructed frames with both low-level structure and high-level semantics.

3.6 Complexity Analysis

The proposed TVC is a holistic framework built upon TVR, which is inherently computationally efficient. It avoids redundant computation across modules and employs a single-pass inference process to handle each GoP for both the Encoder and Decoder. In contrast, traditional video codecs based on motion estimation and residual coding require replicating the entire Decoder-side process within

the Encoder, increasing complexity. Existing neural video codecs adopt a piecemeal replacement strategy, substituting individual components with separate networks – resulting in significant redundancy and high computational overhead.

In addition, every module in TVC is carefully designed for computational efficiency. The Cosmos discrete and continuous tokenizers, along with the pixel decoder, are up to 12× faster than the best available SOTA tokenizers [29]. The CCM method is specifically optimized for fast learned image compression, offering SOTA computational efficiency. In the D-TVR stream, our masking mechanism significantly reduces the number of discrete tokens to be processed, while the C-TVR stream benefits from input downsampling, further minimizing its computational load. Additionally, the transformer-based token predictor is lightweight by design. As a result, the integration of these efficient components within a single-pass inference process enables TVC to run efficiently.

Through empirical evaluation, our TVC runs about 2× faster than SOTA neural codecs like DCVC-DC [18], and runs more than 150× faster than conventional codecs like VTM [14]. On average, the complete Encoder–Decoder pipeline of TVC runs at 0.3 sec/frame, compared to 50 sec/frame for VTM and 0.6 sec/frame for DCVC-DC, when evaluated on a single NVIDIA L40s GPU – demonstrating TVC’s practical efficiency. Notably, TVC’s concise architecture ensures stable and consistent computational complexity across different videos. This contrasts with traditional codecs, which exhibit significant variability in computation due to complex module selection and tool-switching mechanisms affected by video content.

4 Experiments

4.1 Experimental Settings

Datasets. We trained our model on the Kinetics-600 training dataset [4, 15]. Specifically, video clips containing more than 300 frames and a resolution greater than 256×256 were selected from the training split. A centered 256×256 crop was extracted from each selected video to construct the training set. For evaluation, we followed prior works and conducted experiments on standard benchmarks, including HEVC Class B [31], UVG [27], and MCL-JCV [32].

Compared baselines. We compared our method against both traditional and learned video compression approaches. For traditional



Figure 5: Qualitative comparison of reconstruction results across multiple test videos. TVC consistently outperforms conventional codecs like HM and VTM in preserving both semantic structures and visual sharpness, even at significantly lower bitrates. The actual bpp values are shown for each method. Our method achieves high-fidelity reconstruction with reduced visual artifacts, such as blurring, ghosting, and blocking.

codecs, we included HM [11] and VTM [14], configured with the Low Delay (LD) setting. Among learned methods, we evaluated against DCVC-DC [18], a representative SOTA neural video codec. However, due to our focus on ultra-low bitrate settings, direct comparisons are limited, as existing methods like DCVC-DC do not operate effectively at such extreme bitrate levels.

Metrics. In the ultra-low bitrate regime, conventional metrics such as PSNR and MS-SSIM tend to prioritize structural fidelity while failing to capture perceptual realism. To better reflect human visual perception under extreme compression, we adopt LPIPS as our primary evaluation metric, due to its stronger alignment with perceptual quality and subjective visual preference.

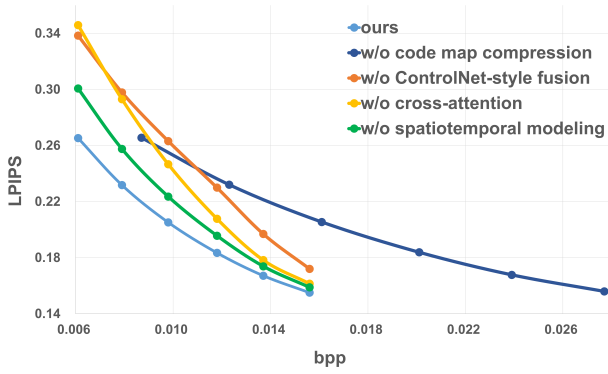


Figure 6: Ablation study evaluating the impact of key design components. Each variant disables a specific module from the full model, including code map compression, dual-stream fusion, cross-attention, and spatiotemporal modeling. Results show consistent performance drops across all ablations, highlighting the importance of each component.

4.2 Quantitative Performance

Figure 4 presents the quantitative rate-distortion (R-D) performance of various methods in terms of **LPIPS** on the UVG, MCL-JCV, and HEVC Class B datasets under ultra-low bitrate conditions. Compared to both traditional codecs and prior neural video compression approaches, our TVC consistently achieves superior perceptual quality, especially in the low-bitrate regime.

Notably, TVC outperforms conventional codecs such as HM and VTM, and surpasses neural baselines like DCVC-DC across all datasets in most bitrate ranges, particularly in the sub-0.02 bpp regime. This performance gain stems from its structured dual-stream architecture and token prediction pipeline. The integration of the checkerboard entropy model and cross-attention transformer decoding enables more accurate and compact reconstructions, particularly in texture-rich or fast-motion scenes where other methods tend to degrade.

4.3 Qualitative Performance

Figure 5 presents visual comparisons between TVC and both traditional and neural codecs across diverse video scenes. Despite operating at significantly lower bitrates (e.g., 0.013x vs. 0.018x), TVC retains fine details such as edges, textures, and semantic boundaries that are often smoothed out or distorted in HM and VTM outputs.

Notably, the ControlNet-inspired fusion of continuous and discrete streams enables fine-grained texture reconstruction, while the transformer-based token prediction improves global coherence and temporal consistency. These qualitative results validate our token-centric design, which effectively preserves both structure and semantics even under extreme compression.

4.4 Ablation Study

We conduct ablation experiments on the Kinetics-600 test set to assess the contribution of each architectural component. As shown in Figure 6, removing any of the core modules leads to noticeable degradation in LPIPS, confirming their complementary roles:

- **Code map compression.** Disabling our code map compression (CCM) leads to a consistent drop in rate-distortion performance—LPIPS increases by up to 0.03 at 0.01 bpp—due to inefficient entropy modeling over the high-cardinality index space.
- **Cross-attention.** Without the cross-attention mechanism, the perceptual quality degrades significantly, especially in semantically complex scenes. LPIPS worsens by 0.04 on average, indicating that continuous-token guidance is essential for accurate reconstruction of masked discrete tokens.
- **Spatiotemporal modeling.** Removing temporal attention and position encoding weakens the model’s ability to capture inter-frame dependencies, resulting in flickering artifacts and unstable predictions across frames. The LPIPS degradation becomes especially severe at lower bitrates.
- **ControlNet-style fusion.** The absence of dual-stream fusion causes the sharpest degradation in perceptual quality across all settings. LPIPS rises by over 0.05 at low bitrates, highlighting the critical role of stream interaction in preserving fine-grained structure.

4.5 Further Discussions

Our TVC framework operates entirely in token space and is inherently modular, allowing for seamless integration with a wide range of tokenization strategies. While our current implementation builds upon Cosmos, recent work such as [30] introduces potential alternatives. The modular nature of TVC enables smooth integration with such emerging tokenizers. TVC can adaptively evolve in tandem with future advances in visual tokenization techniques.

By reformulating compression as a token prediction and selection task, TVC represents a paradigm shift from traditional signal-based coding to semantic-aware modeling. This perspective encourages the co-design of tokenizers and compressors, paving the way for end-to-end trainable systems that are jointly optimized for compression efficiency, perceptual quality, and controllability.

Looking ahead, we envision extending the compression pipeline with class-conditional tokenization, multimodal alignment, and instruction-driven generation backbones. These advancements could enable personalized and task-specific compression solutions, tailored to applications like video editing and synthesis.

5 Conclusion

We introduce TVC, a novel video compression framework that implements a fully tokenized pipeline capable of operating at ultra-low bitrates. This marks a significant departure from conventional video compression approaches. By reframing compression as a task of token selection and prediction, TVC transcends traditional paradigms based on motion estimation and residual coding, enabling a more semantic, flexible, and efficient representation of video content.

The dual-stream architecture effectively separates semantic content from fidelity details, creating a more flexible representation space. The integration of causal attention mechanisms, checkerboard modeling, and ControlNet-inspired fusion demonstrates that a semantic-first, perception-driven approach can achieve remarkable compression efficiency at rates as low as 0.01 bpp.

TVC’s contribution goes beyond technical innovation by introducing a paradigm shift from pixel-based to token-based representation, aligning video compression more closely with modern generative modeling techniques. This perspective opens promising opportunities for integrating compression with a wide range of downstream applications, including content retrieval, synthesis, and interactive systems.

References

- [1] Johannes Ballé, David Minnen, Saurabh Singh, Sung-Jin Hwang, and Nick Johnston. 2018. Variational image compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436* (2018).
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is Space-Time Attention All You Need for Video Understanding?. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- [3] Benjamin Bross, Jianle Chen, Shan Liu, Jens-Rainer Ohm, Andrew Segall, and Gary J Sullivan. 2021. Overview of the versatile video coding (VVC) standard. *IEEE Transactions on Circuits and Systems for Video Technology* 31, 10 (2021), 3736–3764.
- [4] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. 2018. A short note about kinetics-600. *arXiv preprint arXiv:1808.01340* (2018).
- [5] Chaofeng Chen, Xinyu Shi, Yipeng Qin, Xiaoming Li, Xiaoguang Han, Tao Yang, and Shihui Guo. 2022. Real-world blind super-resolution via feature matching with implicit high-resolution priors. In *Proceedings of the 30th ACM International Conference on Multimedia*. 1329–1338.
- [6] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. 2020. Learned image compression with discretized gaussian mixture likelihoods and attention modules. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 7939–7948.
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [8] Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12873–12883.
- [9] Dailan He, Yaoyan Zheng, Baocheng Sun, Yan Wang, and Hongwei Qin. 2021. Checkerboard context model for efficient learned image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14771–14780.
- [10] Zhong Hu, Ruili Yang, Wangmeng Li, and Shuaicheng Liu. 2021. FVC: A new framework towards deep video compression in feature space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1502–1511.
- [11] JCT-VC. 2016. HM: HEVC Reference Software. <https://vcgit.hhi.fraunhofer.de/jvet/HM/>.
- [12] Wei Jiang, Wei Wang, and Yue Chen. 2024. Neural Image Compression Using Masked Sparse Visual Representation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 4189–4197.
- [13] Wei Jiang, Jiayu Yang, Yongqi Zhai, Peirong Ning, Feng Gao, and Ronggang Wang. 2023. MLIC: Multi-reference entropy model for learned image compression. In *Proceedings of the 31st ACM International Conference on Multimedia*. 7618–7627.
- [14] JVT. 2020. VTM: VVC Reference Software. https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware_VTM/.
- [15] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. 2017. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950* (2017).
- [16] Diederik P Kingma, Max Welling, et al. 2013. Auto-encoding variational bayes.
- [17] Didier Le Gall. 1991. MPEG: A video compression standard for multimedia applications. *Commun. ACM* 34, 4 (1991), 46–58.
- [18] Jiahao Li, Bin Li, and Yan Lu. 2023. Neural video compression with diverse contexts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 22616–22626.
- [19] Tianhong Li, Huiwen Chang, Shlok Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan. 2023. MAGE: Masked generative encoder to unify representation learning and image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2142–2152.
- [20] Tianhong Li Li, Huiwen Chang, Shlok Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan. 2023. MAGE: Masked generative encoder to unify representation learning and image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2142–2152.
- [21] Kechun Liu, Yitong Jiang, Inchang Choi, and Jinwei Gu. 2023. Learning image-adaptive codebooks for class-agnostic image restoration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5373–5383.
- [22] Guo Lu, Wanli Ouyang, Dong Xu, Xin Zhang, and Zhiwei Gao. 2019. DVC: An end-to-end deep video compression framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11006–11015.
- [23] Lei Lu, Yize Li, Yanzhi Wang, Wei Wang, and Wei Jiang. 2025. HDCompression: Hybrid-Diffusion Image Compression for Ultra-Low Bitrates. *arXiv preprint arXiv:2502.07160* (2025).
- [24] Lei Lu, Yanyue Xie, Wei Jiang, Wei Wang, Xue Lin, and Yanzhi Wang. 2024. HybridFlow: Infusing Continuity into Masked Codebook for Extreme Low-Bitrate Image Compression. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 3010–3018.
- [25] Qi Mao, Tinghan Yang, Yinu Zhang, Shuyin Pan, Meng Wang, Shiqi Wang, and Siwei Ma. 2023. Extreme Image Compression using Fine-tuned VQGAN Models. *arXiv preprint arXiv:2307.08265* (2023).
- [26] Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. 2023. Finite scalar quantization: VQ-VAE made simple. *arXiv preprint arXiv:2309.15505* (2023).
- [27] Alexandre Mercat, Marko Viitanen, and Jarno Vanne. 2020. UVG dataset: 50/120fps 4K sequences for video codec analysis and development. In *Proceedings of the 11th ACM multimedia systems conference*. 297–302.
- [28] NVIDIA Research. 2025. Cosmos World Foundation Model Platform for Physical AI. *arXiv preprint arXiv:2501.03575* (2025). <https://arxiv.org/abs/2501.03575>
- [29] Felix Reda et al. 2023. Cosmos Tokenizer. <https://github.com/NVIDIA/Cosmos-Tokenizer>. Accessed: 2025-02-24.
- [30] Kyle Sargent, Kyle Hsu, Justin Johnson, Li Fei-Fei, and Jiajun Wu. 2025. Flow to the Mode: Mode-Seeking Diffusion Autoencoders for State-of-the-Art Image Tokenization. *arXiv preprint arXiv:2503.11056* (2025).
- [31] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. 2012. Overview of the high efficiency video coding (HEVC) standard. *IEEE Transactions on Circuits and Systems for Video Technology* 22, 12 (2012), 1649–1668.
- [32] Haiqiang Wang, Weihao Gan, Sudeng Hu, Joe Yuchieh Lin, Lina Jin, Longguang Song, Ping Wang, Ioannis Katsavounidis, Anne Aaron, and C-C Jay Kuo. 2016. MCL-JCV: a JND-based H. 264/AVC video quality assessment dataset. In *2016 IEEE international conference on image processing (ICIP)*. IEEE, 1509–1513.
- [33] Junke Wang, Yi Jiang, Zehuan Yuan, Binyue Peng, Zuxuan Wu, and Yu-Gang Jiang. 2024. OmniTokenizer: A Joint Image-Video Tokenizer for Visual Generation. *arXiv preprint arXiv:2406.09399* (2024).
- [34] Lvmin Zhang and Maneesh Agrawala. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. *arXiv:2302.05543 [cs.CV]*
- [35] Wenbin Zou, Hongxia Gao, Tian Ye, Liang Chen, Weipeng Yang, Shasha Huang, Hongshen Chen, and Sixiang Chen. 2023. VQCNIR: Clearer Night Image Restoration with Vector-Quantized Codebook. *arXiv preprint arXiv:2312.08606* (2023).

Received XX April 2025; revised XX XX 2025; accepted XX XX 2025