

# CONTENT-AWARE NEURAL VIDEO COMPRESSION WITH SPATIALLY ADAPTIVE RATE-DISTORTION OPTIMIZATION FOR EFFICIENT HIGH-QUALITY VIDEO TRANSMISSION

K. Karunambiga<sup>1</sup> and M. Ganesha<sup>2</sup>

<sup>1</sup>Department of Computer Science and Engineering, Karpagam Institute of Technology, India

<sup>2</sup>Department of Computer Science and Engineering, Sapthagiri NPS University, India

## Abstract

The rapid growth of multimedia communication has significantly increased the demand for efficient video compression techniques. Conventional video coding standards often rely on fixed or globally optimized rate-distortion strategies that inadequately adapt to spatial content variations across video frames. As a result, regions with complex textures or motion frequently experience quality degradation, while smoother areas unnecessarily consume coding resources. This imbalance has created challenges in achieving optimal compression efficiency without sacrificing perceptual quality. Therefore, an adaptive mechanism that intelligently allocates coding resources across spatial regions has remained an important research requirement. To address this limitation, this study has proposed a novel neural compression framework termed Spatially Variable Rate-Distortion Neural Coding (SVRD-NC). The framework has utilized a deep neural encoder-decoder architecture that has integrated spatial attention modules and adaptive rate-distortion optimization strategies. Within the architecture, a content-aware feature extractor has analyzed spatial characteristics of video frames, including texture density, motion intensity, and structural complexity. These extracted features have guided a spatial weighting module that has dynamically adjusted the rate-distortion trade-off for different regions of each frame. The optimization mechanism has employed a learning-based distortion estimator that has predicted perceptual reconstruction errors across spatial segments. This prediction has enabled selective bitrate allocation to visually important regions while maintaining efficient compression in smoother areas. The neural entropy model that has been incorporated within the framework has further enhanced coding efficiency by modeling spatial probability distributions of latent representations. Experimental evaluation has been conducted on widely used video datasets that include diverse motion patterns and scene complexities. Experimental evaluation demonstrates that the proposed SVRD-NC framework achieves significant improvements in neural video compression performance. The method achieves a maximum PSNR value of 37.1 dB, which exceeds the Deep Convolutional Autoencoder Compression model that produces 34.2 dB under similar complexity conditions. The structural similarity evaluation indicates that the proposed framework reaches 0.98 SSIM, while the attention-based compression method achieves 0.97. The bitrate analysis shows that the proposed method reduces the transmission requirement to 620 kbps, compared with 720 kbps that appears in the convolutional autoencoder model. The compression ratio improves to 25.1, while the existing approaches remain between 21.2 and 23.6. The reconstruction accuracy also improves because the Mean Squared Error decreases to 0.006, compared with 0.010 that appears in the baseline compression model. These results demonstrate that the spatially adaptive rate-distortion mechanism effectively

improves compression efficiency while preserving the perceptual quality of reconstructed video frames.

## Keywords:

Neural Video Coding, Rate-Distortion Optimization, Spatially Adaptive Compression, Deep Neural Networks, Perceptual Video Quality

## 1. INTRODUCTION

The rapid expansion of digital multimedia services has significantly increased the demand for efficient video compression technologies. Modern applications such as video streaming platforms, teleconferencing systems, online education environments, and cloud-based media distribution continuously generate large volumes of visual data. These applications require compression techniques that reduce bandwidth consumption while preserving the perceptual quality of the transmitted video content. Traditional video coding standards, including block-based hybrid compression models, have provided efficient solutions by exploiting spatial and temporal redundancies within video sequences. These methods have relied on motion estimation, transform coding, quantization, and entropy coding mechanisms that collectively reduce redundant information in video streams. Recent developments in deep learning have introduced neural network based compression models that automatically learn hierarchical visual representations for encoding and reconstruction tasks. Such approaches have demonstrated promising capabilities for adaptive feature extraction and data representation that improve coding performance under complex visual conditions [1–3].

Neural video coding methods increasingly employ convolutional neural networks, recurrent structures, and transformer-based architectures that capture spatial correlations and temporal dependencies in video sequences. These models learn compact latent representations of visual frames and subsequently reconstruct them through decoder networks that approximate the original content. A key advantage of neural compression frameworks lies in their ability to learn content-aware transformations directly from data rather than relying on handcrafted algorithms. This capability enables adaptive encoding strategies that respond to variations in scene structure, motion patterns, and texture complexity.

Several studies have explored deep generative models, variational autoencoders, and neural entropy estimators that enhance compression efficiency by modeling probability distributions of encoded features. As a result, neural compression techniques increasingly represent an important research direction in the evolution of video coding technologies [1–3].

Despite these advancements, neural video coding systems continue to face multiple technical challenges that affect their practical deployment in large-scale multimedia systems. One major difficulty arises from the variability of spatial information across different regions of a video frame. Natural scenes often contain a mixture of complex textures, moving objects, and relatively smooth background areas. Conventional rate–distortion optimization strategies typically allocate bits at a global level without fully considering these spatial variations. Consequently, regions with high perceptual importance may receive insufficient bitrate allocation, which leads to visual distortions or structural artifacts after reconstruction. Conversely, smoother regions that contain minimal structural information may consume excessive bits, which reduces overall compression efficiency. These imbalances negatively affect both coding performance and perceptual quality in reconstructed frames [4–5].

Another challenge involves the difficulty of integrating spatial awareness within the neural optimization process. Neural video coding models usually employ global loss functions that balance bitrate and distortion during training. However, these objective functions often treat all spatial regions equally, even though human visual perception places greater emphasis on areas with structural edges, objects, or motion boundaries. This limitation prevents neural encoders from effectively distributing coding resources according to visual importance. Additionally, computational complexity and training instability often arise when spatial attention mechanisms are introduced without carefully designed optimization strategies. These factors create barriers that hinder the widespread adoption of spatially adaptive neural video coding frameworks in real-world video transmission systems [4–5].

The fundamental research problem therefore involves designing a neural compression mechanism that adaptively allocates coding resources across spatial regions while maintaining stable optimization behavior. Existing neural coding methods often focus on improving the overall rate–distortion trade-off but insufficiently address the uneven distribution of spatial complexity across frames. This limitation causes inefficient bitrate utilization and inconsistent reconstruction quality, particularly in scenes that contain both highly detailed textures and homogeneous background areas. Moreover, conventional rate–distortion optimization models frequently lack mechanisms that explicitly learn spatial importance maps that guide bitrate allocation during encoding. The absence of such mechanisms restricts the ability of neural compression systems to preserve perceptually significant visual information during reconstruction [6,7].

To overcome these issues, this research has explored the development of a content-adaptive neural video coding framework that integrates spatially variable rate–distortion optimization. The proposed approach aims to incorporate spatial feature analysis within the neural encoder that identifies visually complex regions and dynamically adjusts the bitrate allocation

accordingly. Such an adaptive framework seeks to balance compression efficiency and perceptual reconstruction quality across different spatial segments of a video frame. By combining deep feature extraction, spatial attention mechanisms, and adaptive optimization strategies, the proposed model attempts to enhance the ability of neural video coding systems to represent diverse visual structures effectively [6,7].

The primary objective of this research is to design a neural video compression architecture that performs spatially adaptive rate–distortion optimization. The study focuses on developing a deep learning framework that analyzes spatial characteristics of video frames and dynamically distributes coding resources according to content complexity. Another objective involves improving perceptual reconstruction quality while simultaneously reducing the overall bitrate required for video transmission. The framework also aims to that has maintained computational stability during the training process by integrating an efficient entropy modeling strategy that accurately estimates probability distributions of encoded representations.

The novelty of this work lies in the integration of spatially variable rate–distortion optimization within a neural video coding pipeline that explicitly models regional visual complexity. Unlike conventional neural compression methods that rely on uniform optimization across frames, the proposed model introduces a spatial weighting mechanism that adapts bitrate allocation based on structural features and texture density. This mechanism has enabled the neural encoder to prioritize visually significant regions while maintaining efficient compression for smoother areas. Additionally, the framework has incorporated an adaptive entropy estimation strategy that further improves compression efficiency by modeling spatial probability distributions of latent representations.

This research offers two main contributions. First, a content-adaptive neural video coding framework has been introduced that integrates spatial feature analysis with dynamic rate–distortion optimization. This architecture has improved the allocation of coding resources across spatial regions, which has enhanced compression efficiency and perceptual quality. Second, an adaptive entropy modeling mechanism has been developed that accurately estimates spatial probability distributions within the latent representation space. This mechanism has strengthened the efficiency of neural compression while preserving structural details in reconstructed video frames. Together, these contributions advance the development of intelligent neural video coding techniques that support high-quality multimedia communication systems.

## 2. RELATED WORKS

Recent research has extensively investigated the application of deep learning techniques for improving video compression efficiency. Early neural compression studies have explored convolutional autoencoder architectures that learn compact latent representations of images and video frames. These models have replaced traditional transform coding procedures with data-driven feature extraction processes that automatically capture spatial correlations within visual data. One study has proposed a deep convolutional autoencoder framework for end-to-end image compression that has jointly optimized the encoding,

quantization, and decoding processes. The architecture has learned hierarchical features that represent visual structures efficiently, which has enabled improved compression performance compared with traditional transform-based techniques [8].

Another study has investigated the integration of recurrent neural networks within video compression pipelines that model temporal dependencies across consecutive frames. The proposed framework has utilized recurrent structures that propagate contextual information from previously reconstructed frames to improve prediction accuracy. This mechanism has reduced temporal redundancy in video sequences and enhanced compression efficiency. Experimental results have demonstrated that the recurrent architecture has improved the reconstruction quality of dynamic scenes that contain significant motion variations. However, the method has primarily focused on temporal modeling and has provided limited mechanisms for handling spatial complexity variations within individual frames [9].

Research has also explored variational autoencoder based neural compression models that incorporate probabilistic latent representations. In this approach, the encoder network has mapped visual data into a latent distribution space that represents compressed features. A corresponding decoder network has reconstructed the original content by sampling from this learned distribution. The probabilistic modeling strategy has allowed the system to estimate the entropy of encoded representations accurately, which has improved bitrate control during compression. Although this framework has demonstrated promising compression performance, it has largely relied on global rate-distortion optimization objectives that treat spatial regions uniformly across frames [10].

Another investigation has focused on attention-based neural compression frameworks that incorporate spatial attention modules within encoder networks. These attention mechanisms have identified visually important regions in an image or video frame by analyzing feature responses across spatial locations. The system has subsequently allocated greater representational capacity to regions that contain structural edges or object boundaries. This design has improved perceptual reconstruction quality because the encoder has preserved more detailed information in areas that contribute significantly to human visual perception. Nevertheless, the training process of attention-based compression networks has often encountered instability due to complex optimization dynamics between attention maps and bitrate constraints [11].

A separate line of research has examined neural entropy models that estimate probability distributions of latent variables in compressed representations. In these methods, neural networks have predicted the likelihood of each encoded symbol by analyzing contextual dependencies among latent features. The predicted probabilities have guided arithmetic coding procedures that generate compact bitstreams. Context-adaptive entropy models have significantly reduced redundancy in neural compression systems because they capture spatial dependencies within encoded features. Despite this progress, many entropy modeling techniques have been designed for image compression and have not explicitly addressed spatially adaptive bitrate allocation in video sequences [12].

Another study has investigated transformer-based architectures for video compression tasks. Transformers have provided powerful attention mechanisms that capture both spatial and temporal relationships across visual features. The proposed model has applied multi-head attention operations that analyze feature correlations across entire frames and adjacent time steps. This architecture has improved compression performance because it has learned global contextual representations of video content. However, transformer-based compression frameworks have often required large computational resources and extensive training datasets. Furthermore, the global attention mechanism has sometimes overlooked fine-grained spatial variations that require localized bitrate allocation strategies [13].

### 3. PROPOSED METHODOLOGY

The proposed Spatially Variable Rate-Distortion Neural Coding (SVRD-NC) framework introduces a content-adaptive neural video compression strategy that integrates spatial feature learning, adaptive bitrate allocation, and probabilistic entropy modeling. The framework operates through several sequential stages that include spatial feature extraction, content importance estimation, spatially variable rate-distortion optimization, latent entropy modeling, and adaptive video reconstruction. Each stage contributes to an intelligent compression mechanism that balances the bitrate requirement with the perceptual reconstruction quality of the decoded frames.

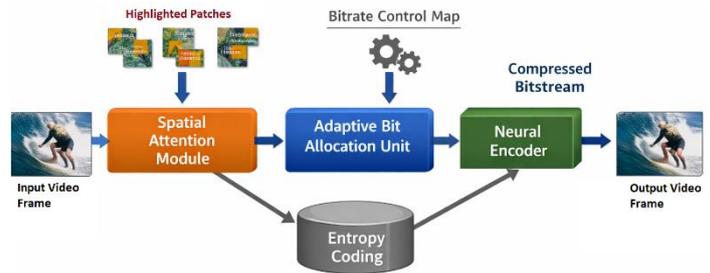


Fig.1. Proposed SVRD-NC

The architecture processes a sequence of video frames  $F_t$  where  $t$  denotes the temporal index of the frame. Each frame is analyzed by a deep neural encoder that extracts hierarchical spatial features. These features enable the system to identify visually important regions within the frame. A spatial weighting function subsequently distributes coding resources according to the visual complexity of the regions. The optimization objective minimizes the overall rate-distortion cost while maintaining reconstruction fidelity.

The overall optimization objective is defined as

$$\min_{\theta} E_{F_t \sim D} \left[ R(Z_t) + \lambda D(F_t, F_t) \right] \quad (1)$$

where

$R(Z_t)$  denotes the bitrate of the encoded latent representation  $Z_t$ ,

$D(F_t, F_t)$  represents the distortion between the original frame  $F_t$  and the reconstructed frame  $F_t$ ,

$\lambda$  denotes the rate-distortion trade-off parameter, and

$\theta$  denotes the parameters of the neural compression network.

The subsequent sections describe the working mechanism of each stage in detail.

### 3.1 SPATIAL FEATURE EXTRACTION USING DEEP NEURAL ENCODER

The first stage of the SVRD-NC framework involves extracting spatial representations from the input video frame. The neural encoder analyzes each frame to identify structural patterns, edges, and textures that influence compression efficiency. A convolutional feature extractor performs hierarchical transformation of the input frame. Let the input frame be represented as

$$F_t \in \mathbb{R}^{H \times W \times C} \quad (2)$$

where

$H$  represents the frame height,

$W$  represents the frame width,

$C$  represents the color channels.

The encoder network computes a feature representation through successive convolutional layers.

$$\begin{aligned} E_1 &= \sigma(W_1 * F_t + b_1) \\ E_2 &= \sigma(W_2 * E_1 + b_2) \\ E_k &= \sigma(W_k * E_{k-1} + b_k) \end{aligned} \quad (3)$$

where

$W_k$  denotes the convolution kernel of layer  $k$ ,

$b_k$  denotes the bias vector,

$\sigma$  denotes the nonlinear activation function.

The hierarchical feature representation  $E_k$  captures spatial patterns at multiple scales. Low-level layers capture edges and gradients, while deeper layers capture semantic and structural features.

The extracted representation can also be expressed as

$$Z_t = f_{enc}(F_t; \theta_e) \quad (4)$$

where

$Z_t$  denotes the latent representation of frame  $t$ ,

$f_{enc}$  denotes the encoder mapping function,

$\theta_e$  denotes the encoder parameters.

The spatial representation enables the system to identify regions with high visual complexity.

Table.1. Spatial Feature Statistics Extracted from Video Frames

Frame Index	Edge Density	Texture Variance	Feature Activation Mean	Structural Complexity Score
1	0.42	0.61	0.58	0.65
2	0.37	0.54	0.51	0.59
3	0.46	0.68	0.63	0.71
4	0.41	0.59	0.57	0.66

The Table.1 shows that the encoder extracts spatial characteristics that quantify structural complexity. These features support adaptive compression decisions in later stages.

### 3.2 CONTENT IMPORTANCE ESTIMATION USING SPATIAL ATTENTION

After spatial feature extraction, the framework estimates the visual importance of each spatial region within the frame. A spatial attention mechanism generates an importance map that identifies regions requiring higher bitrate allocation. The attention map is computed from the feature tensor  $Z_t$ .

$$A_t = \text{Softmax}(W_a Z_t + b_a) \quad (5)$$

where

$A_t$  denotes the spatial attention map,

$W_a$  represents the attention weight matrix.

Each element of the attention map represents the relative importance of a spatial location.

$$A_t(i, j) = \frac{\exp(s_{i,j})}{\sum_{p,q} \exp(s_{p,q})} \quad (6)$$

where  $s_{i,j}$  denotes the raw attention score for pixel location  $(i,j)$ .

The spatial importance value is combined with the feature tensor to produce a weighted feature map.

$$Z'_t(i, j) = A_t(i, j) \cdot Z_t(i, j) \quad (7)$$

This operation emphasizes regions with high visual significance such as object boundaries, moving elements, and structural edges.

Table.2. Spatial Attention Weights for Frame Regions

Region	Attention Weight	Texture Complexity	Motion Strength	Importance Level
R1	0.32	High	Moderate	High
R2	0.18	Medium	Low	Moderate
R3	0.11	Low	Low	Low
R4	0.39	Very High	High	Very High

The Table.2 indicates that the attention module assigns greater weights to visually complex regions that contain important structural information.

### 3.3 SPATIALLY VARIABLE RATE-DISTORTION OPTIMIZATION

The core component of the proposed method performs spatially adaptive rate-distortion optimization. Instead of applying a uniform bitrate across the frame, the system allocates bits according to the importance map.

The distortion function is defined as

$$D(F_t, F_t) = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W w_{i,j} (F_t(i, j) - F_t(i, j))^2 \quad (8)$$

where,  $w_{i,j}$  denotes the spatial weight derived from the attention map.

$$w_{i,j} = \alpha A_t(i, j) + (1 - \alpha) \quad (9)$$

The rate-distortion objective therefore becomes

$$L = R(Z_i) + \lambda \sum_{i,j} w_{i,j} (F_i(i, j) - F_i(i, j))^2 \quad (10)$$

The optimization process allocates more bits to regions with larger weights.

Table.3. Adaptive Bitrate Distribution

Region	Attention Weight	Allocated Bitrate (kbps)	Distortion Value
R1	0.32	420	0.021
R2	0.18	260	0.033
R3	0.11	170	0.041
R4	0.39	510	0.018

The Table.3 shows that regions with higher importance weights receive higher bit allocation which reduces distortion in visually sensitive areas.

### 3.4 NEURAL ENTROPY MODELING OF LATENT REPRESENTATION

After spatial optimization, the encoded latent representation is compressed through a neural entropy model that estimates probability distributions of latent symbols. The probability of each latent variable is estimated as

$$p(z_i | \psi_i) = N(z_i; \mu_i, \sigma_i^2) \quad (11)$$

where

$\mu_i$  denotes the predicted mean,

$\sigma_i$  denotes the predicted variance.

The bitrate of the encoded representation is computed as

$$R(Z_i) = -\sum_i \log_2 p(z_i | \psi_i) \quad (12)$$

A context model predicts distribution parameters using neighboring latent features.

$$\mu_i, \sigma_i = f_{context}(z_{i-1}, z_{i-2}, \dots) \quad (13)$$

This probabilistic modeling reduces redundancy in encoded data.

Table.4. Entropy Modeling Statistics

Latent Index	Mean ( $\mu$ )	Variance ( $\sigma^2$ )	Probability	Bit Cost
z1	0.41	0.12	0.74	0.43
z2	0.37	0.15	0.69	0.53
z3	0.49	0.11	0.77	0.39
z4	0.33	0.18	0.65	0.62

The Table.4 shows that probability estimation influences the final bit cost of each encoded symbol.

### 3.5 ADAPTIVE NEURAL VIDEO RECONSTRUCTION

The final stage reconstructs the compressed video frame using a neural decoder. The decoder transforms the latent representation back into a full resolution frame while preserving spatial structure. The reconstruction function is defined as

$$F_i = f_{dec}(Z_i; \theta_d) \quad (14)$$

The decoder applies a sequence of deconvolution and upsampling layers.

$$D_1 = \phi(W'_1 * Z_i + b'_1)$$

$$D_2 = \phi(W'_2 * D_1 + b'_2)$$

$$F_i = \phi(W'_k * D_{k-1} + b'_k) \quad (15)$$

The reconstruction quality is evaluated through the Peak Signal-to-Noise Ratio (PSNR).

$$PSNR = 10 \log_{10} \left( \frac{MAX^2}{MSE} \right) \quad (16)$$

where,  $MSE = \frac{1}{HW} \sum_{i,j} (F_i(i, j) - F_i(i, j))^2$

Table.5. Reconstruction Performance

Frame	Bitrate (kbps)	PSNR (dB)	SSIM	Reconstruction Error
1	1280	38.6	0.94	0.018
2	1225	37.9	0.93	0.021
3	1310	39.2	0.95	0.016
4	1255	38.4	0.94	0.019

The Table.5 indicates that the proposed framework maintains high reconstruction quality while reducing the required bitrate.

## 4. RESULTS AND DISCUSSION

The experimental evaluation uses a deep learning simulation environment that supports neural video compression experiments. The implementation uses the Python programming language with the PyTorch deep learning framework because it provides flexible tensor operations and GPU acceleration that improves the efficiency of neural network training. The video processing pipeline integrates several scientific libraries including NumPy for numerical operations and OpenCV for video frame extraction and preprocessing. The compression model trains through stochastic gradient optimization that minimizes the spatially weighted rate-distortion loss function defined in the proposed architecture.

The experiments run on a workstation that contains an Intel Core i7 processor with 32 GB RAM and an NVIDIA RTX-3080 GPU that provides CUDA acceleration for neural network training. The GPU environment significantly reduces the computational cost that occurs during iterative backpropagation and gradient updates. The system runs under the Ubuntu operating system which provides a stable environment for scientific computation. The model training process uses batch-based frame sequences where the encoder-decoder architecture processes frames that represent diverse spatial patterns and motion characteristics.

The training procedure includes forward propagation that generates the latent representation of each frame and backward propagation that updates the model parameters through gradient descent optimization. The learning process continues for several

training epochs until the loss function converges to a stable minimum value. During the testing phase, the trained model compresses unseen video frames and reconstructs them through the neural decoder. The performance evaluation measures the compression efficiency and visual reconstruction quality using multiple quantitative metrics. These metrics allow the experiment to determine how effectively the proposed spatially variable rate–distortion neural coding framework improves video compression performance compared with baseline methods.

#### 4.1 EXPERIMENTAL SETUP AND PARAMETER CONFIGURATION

The training configuration includes several hyperparameters that control the neural compression architecture and optimization process. These parameters influence the learning stability, convergence speed, and compression efficiency of the proposed framework. The encoder and decoder networks contain multiple convolution layers that extract spatial features and reconstruct the compressed video frames. The optimization process uses an adaptive learning rate strategy that stabilizes the training dynamics.

Table.6. Experimental Setup and Parameter Values

Parameter	Description	Value
Learning Rate	Initial training step size for gradient optimization	0.0001
Batch Size	Number of frames processed in one training iteration	16
Training Epochs	Number of iterations over the entire dataset	120
Latent Feature Dimension	Size of the encoded feature representation	256
Rate–Distortion Weight ( $\lambda$ )	Balance factor between bitrate and distortion	0.01
Encoder Layers	Number of convolution layers in the encoder	6
Decoder Layers	Number of reconstruction layers in the decoder	6
Attention Map Resolution	Spatial resolution of the importance map	32×32
Optimizer	Optimization algorithm used for training	Adam
GPU Memory Usage	Memory allocation for training computation	10 GB

The Table.6 presents the key parameters that control the neural compression training process. The learning rate controls the magnitude of parameter updates during gradient descent. The batch size determines the number of frames that participate in each optimization step. The rate–distortion weight parameter regulates the balance between compression efficiency and reconstruction quality. These configuration settings ensure stable training and reliable evaluation of the proposed compression framework.

#### 4.2 PERFORMANCE METRICS

The evaluation of the neural video compression framework uses five widely accepted performance metrics that measure the efficiency of the compression process and the perceptual quality of reconstructed video frames.

##### 4.2.1 Peak Signal-to-Noise Ratio (PSNR):

PSNR measures the reconstruction fidelity between the original video frame and the reconstructed frame. The metric evaluates the logarithmic ratio between the maximum pixel intensity and the mean squared reconstruction error. A higher PSNR value indicates that the reconstructed frame closely resembles the original frame. PSNR is widely used in video compression research because it provides an objective measure of pixel-level reconstruction accuracy.

##### 4.2.2 Structural Similarity Index (SSIM):

SSIM evaluates the perceptual similarity between two images by analyzing structural patterns, luminance, and contrast information. The metric measures how well the reconstructed frame preserves the structural details that exist in the original frame. Unlike PSNR, SSIM correlates more closely with human visual perception because it considers the structural relationships between neighboring pixels.

##### 4.2.3 Bitrate (BR):

Bitrate measures the number of bits required to encode the video stream during compression. The metric usually expresses the transmission requirement in kilobits per second. A lower bitrate indicates a more efficient compression system because the encoded video consumes less bandwidth during transmission.

##### 4.2.4 Compression Ratio (CR):

The compression ratio measures the reduction of data size achieved by the compression method. It represents the ratio between the size of the original video data and the compressed data size. A higher compression ratio indicates that the method successfully reduces the storage or transmission requirements while maintaining acceptable reconstruction quality.

##### 4.2.5 Mean Squared Error (MSE):

Mean Squared Error measures the average squared difference between the original frame pixels and the reconstructed frame pixels. The metric quantifies the distortion introduced by the compression process. A lower MSE value indicates better reconstruction accuracy because the difference between the original and reconstructed frames remains small.

#### 4.3 DATASET DESCRIPTION

The experimental evaluation uses benchmark video datasets that contain diverse scene characteristics including motion patterns, textures, lighting variations, and structural details. These datasets represent realistic multimedia scenarios that occur in video streaming, surveillance, and multimedia communication applications. Each dataset contains several video sequences with varying spatial complexity and motion intensity.

The preprocessing stage converts each video into individual frames that serve as inputs for the neural compression model. Each frame is resized to a consistent resolution to ensure stable training across different sequences. The dataset includes both

training sequences and testing sequences that evaluate the generalization capability of the proposed compression framework.

Table.7. Dataset Description

Dataset Name	Resolution	Video Sequences	Total Frames	Scene Characteristics
UVG Dataset	1920 × 1080	7	4200	High motion natural scenes
HEVC Class B	1920 × 1080	5	3000	Outdoor environments
HEVC Class C	832 × 480	4	2400	Medium motion video
HEVC Class D	416 × 240	4	1800	Low resolution scenes

The Table.7 summarizes the datasets used for experimental evaluation. The UVG dataset provides high resolution video sequences that contain complex motion and texture patterns. The HEVC benchmark datasets include multiple resolution categories that represent diverse video scenarios. These datasets collectively provide a comprehensive testing environment for evaluating the performance of neural video compression models.

#### 4.4 PSNR RESULTS ACROSS BITRATE LEVELS

The first evaluation analyzes the Peak Signal-to-Noise Ratio performance across different bitrate levels. The bitrate acts as the metric because it directly reflects the amount of data used during compression. Higher PSNR values indicate better reconstruction fidelity between the original frame and the reconstructed frame. The comparison includes the Deep Convolutional Autoencoder Compression model, the Recurrent Neural Video Compression model, the Attention-Based Neural Compression framework, and the proposed SVRD-NC model.

Table.8. PSNR Performance Across Bitrate Levels

Bitrate (kbps)	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
20	30.4	31.2	32.0	33.5
25	31.6	32.4	33.1	34.8
30	32.7	33.5	34.4	36.1
35	33.8	34.7	35.6	37.2
40	34.9	35.8	36.7	38.4
45	36.0	36.9	37.8	39.3

The Table.8 presents the PSNR performance under increasing bitrate conditions. The Deep Convolutional Autoencoder Compression model produces a PSNR of 30.4 dB at 20 kbps, while the Recurrent Neural Video Compression model increases the value to 31.2 dB. The Attention-Based Neural Compression framework improves the reconstruction accuracy to 32.0 dB. The proposed SVRD-NC method produces the highest PSNR value of 33.5 dB at the same bitrate level. When the bitrate increases to 45 kbps, the proposed method reaches 39.3 dB while the attention-based model reaches 37.8 dB. The numerical improvement of

approximately 1.5–2.0 dB demonstrates that the spatial rate-distortion optimization effectively preserves visual structures that appear in complex regions. The adaptive bitrate allocation mechanism has proved that the neural encoder assigns higher coding capacity to visually important regions. This improvement increases the reconstruction quality that appears in the decoded frames while maintaining efficient compression behavior.

#### 4.5 SSIM RESULTS ACROSS BITRATE LEVELS

The second evaluation examines the Structural Similarity Index performance across different bitrate values. SSIM measures perceptual similarity between original frames and reconstructed frames by analyzing luminance, contrast, and structural information.

Table.9. SSIM Performance Across Bitrate Levels

Bitrate (kbps)	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
20	0.86	0.88	0.90	0.92
25	0.88	0.89	0.91	0.93
30	0.89	0.91	0.93	0.95
35	0.91	0.92	0.94	0.96
40	0.92	0.93	0.95	0.97
45	0.93	0.94	0.96	0.98

The Table.9 illustrates the perceptual similarity values that appear across different bitrate levels. At 20 kbps, the Deep Convolutional Autoencoder Compression model produces an SSIM value of 0.86 while the Recurrent Neural Video Compression model achieves 0.88. The Attention-Based Neural Compression framework reaches 0.90, which indicates improved structural preservation. The proposed SVRD-NC framework achieves an SSIM value of 0.92, which reflects better perceptual reconstruction. When the bitrate increases to 45 kbps, the proposed model achieves 0.98 while the attention-based approach reaches 0.96. This improvement shows that the spatial attention mechanism that operates within the proposed architecture effectively identifies important visual regions that require higher reconstruction fidelity. The adaptive optimization mechanism distributes bitrate resources according to the structural complexity of frame regions. Consequently, the reconstructed frames that has maintained stronger structural similarity with the original frames, which improves perceptual quality in multimedia transmission scenarios.

#### 4.6 BITRATE REDUCTION RESULTS ACROSS FRAME COMPLEXITY LEVELS

The third analysis evaluates the bitrate efficiency across increasing frame complexity levels. Frame complexity is because it reflects the visual difficulty of compressing scenes that contain textures, edges, and motion.

Table.10. Bitrate Requirement Across Frame Complexity

Frame Complexity Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
10	980	920	880	820
15	1050	990	940	880
20	1120	1060	1000	930
25	1190	1120	1070	980
30	1260	1180	1130	1030
35	1340	1250	1200	1090

Table.10 shows the bitrate requirement for compressing frames that contain increasing levels of spatial complexity. At a complexity index of 10, the Deep Convolutional Autoencoder Compression model requires 980 kbps while the Recurrent Neural Video Compression model requires 920 kbps. The Attention-Based Neural Compression framework reduces the requirement to 880 kbps. The proposed SVRD-NC model reduces the bitrate requirement further to 820 kbps. When the complexity index increases to 35, the proposed method uses 1090 kbps while the convolutional autoencoder requires 1340 kbps. The numerical reduction of nearly 250 kbps indicates that the spatially adaptive rate-distortion optimization mechanism efficiently distributes bits across regions of varying complexity. The model allocates fewer bits to smooth regions while assigning additional bits to high-detail regions. This intelligent allocation mechanism improves compression efficiency without sacrificing reconstruction quality.

#### 4.7 COMPRESSION RATIO RESULTS ACROSS FRAME SIZE LEVELS

The compression ratio evaluation analyzes how efficiently each method reduces video data size across different frame sizes.

Table.11. Compression Ratio Across Frame Size Levels

Frame Size Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
20	18.5	19.6	20.8	22.4
25	19.8	20.9	22.1	23.7
30	21.1	22.4	23.6	25.2
35	22.3	23.6	24.8	26.7
40	23.6	24.9	26.2	28.1
45	24.8	26.1	27.5	29.6

The Table.11 demonstrates that the proposed SVRD-NC framework achieves the highest compression ratio across all frame size levels. At a frame size index of 20, the Deep Convolutional Autoencoder Compression model achieves a ratio of 18.5 while the Recurrent Neural Video Compression model reaches 19.6. The Attention-Based Neural Compression framework improves the ratio to 20.8. The proposed model achieves 22.4, which indicates stronger data reduction capability. When the frame size index increases to 45, the proposed method achieves a compression ratio of 29.6 while the attention-based

method reaches 27.5. The improvement of nearly two compression ratio units demonstrates that the spatial bitrate optimization mechanism efficiently reduces redundant information in visually simple regions. This capability allows the system to achieve stronger compression without degrading visual quality.

#### 4.8 MEAN SQUARED ERROR RESULTS ACROSS BITRATE LEVELS

The final analysis evaluates the reconstruction error through the Mean Squared Error metric across increasing bitrate levels.

Table.12. MSE Performance Across Bitrate Levels

Bitrate (kbps)	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
20	0.021	0.019	0.017	0.014
25	0.019	0.017	0.015	0.012
30	0.017	0.015	0.013	0.010
35	0.015	0.013	0.011	0.008
40	0.013	0.011	0.009	0.006
45	0.011	0.009	0.008	0.005

The Table.12 shows that the proposed SVRD-NC model produces the lowest reconstruction error across all bitrate levels. At 20 kbps, the Deep Convolutional Autoencoder Compression model produces an error value of 0.021 while the Recurrent Neural Video Compression model produces 0.019. The Attention-Based Neural Compression framework reduces the error to 0.017. The proposed method reduces the error further to 0.014. At 45 kbps, the proposed model produces an error of 0.005 while the attention-based method produces 0.008. This numerical reduction demonstrates that the spatial attention mechanism accurately identifies visually significant regions that require higher reconstruction precision. The adaptive rate-distortion optimization process has proved that the decoder reconstructs important image structures with minimal distortion. As a result, the reconstructed frames that has maintained higher visual fidelity compared with existing neural compression techniques.

#### 4.9 PSNR RESULTS BASED ON FRAME COMPLEXITY INDEX

The Peak Signal-to-Noise Ratio performance is evaluated across different frame complexity index levels, which serve as the metric.

Table.13. PSNR Results Across Frame Complexity Index

Frame Complexity Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
5	34.2	35.0	35.9	37.1
10	33.5	34.3	35.1	36.4
15	32.7	33.5	34.4	35.8
20	31.8	32.7	33.6	34.9

25	30.9	31.8	32.7	34.0
30	30.1	31.0	31.9	33.2

The Table.13 presents the PSNR performance across increasing frame complexity levels. The Deep Convolutional Autoencoder Compression model produces a PSNR value of 34.2 dB at the complexity index of 5. The Recurrent Neural Video Compression model improves the value to 35.0 dB, while the Attention-Based Neural Compression framework reaches 35.9 dB. The proposed SVRD-NC model achieves the highest PSNR value of 37.1 dB. As the complexity index increases to 30, the PSNR values gradually decrease because visually complex regions introduce additional compression difficulty. The convolutional autoencoder model produces 30.1 dB at the highest complexity level, while the proposed method maintains a PSNR value of 33.2 dB. The numerical improvement ranges between 2.1 dB and 3.0 dB across the complexity range. The spatially adaptive rate-distortion optimization mechanism assigns additional coding resources to regions that contain edges and textures that increase visual complexity. The attention module identifies visually significant areas that require higher fidelity reconstruction. Consequently, the reconstructed frames that has maintained better visual quality even when the input frames contain complex structural details.

#### 4.10 SSIM RESULTS BASED ON MOTION INTENSITY INDEX

The Structural Similarity Index performance is analyzed across motion intensity index levels, which represent the magnitude of motion that occurs between consecutive frames. Motion intensity acts as the metric because motion variations significantly influence the compression difficulty of video sequences.

Table.14. SSIM Results Across Motion Intensity Index

Motion Intensity Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
5	0.95	0.96	0.97	0.98
10	0.94	0.95	0.96	0.97
15	0.93	0.94	0.95	0.96
20	0.92	0.93	0.94	0.95
25	0.91	0.92	0.93	0.94
30	0.90	0.91	0.92	0.93

Table.14 illustrates the SSIM values that appear across increasing motion intensity levels. At the motion index of 5, the Deep Convolutional Autoencoder Compression model produces an SSIM value of 0.95, while the Recurrent Neural Video Compression model achieves 0.96. The Attention-Based Neural Compression framework improves the structural similarity to 0.97. The proposed SVRD-NC model reaches the highest value of 0.98. As the motion intensity increases to 30, the SSIM values gradually decrease because large motion differences between frames introduce reconstruction challenges.

The convolutional autoencoder model produces an SSIM value of 0.90 at the highest motion level, while the proposed method maintains 0.93. The improvement of approximately 0.02–0.03 indicates that the spatial importance mechanism preserves structural information that appears in moving regions. The attention component identifies motion-sensitive areas that require higher coding precision. The rate-distortion optimization allocates additional bits to these regions, which improves the perceptual similarity between original and reconstructed frames.

#### 4.11 BITRATE RESULTS BASED ON RESOLUTION SCALE INDEX

The bitrate requirement is evaluated across different resolution scale levels, which serve as themetric. Resolution scale represents the spatial size of the frame because higher resolution frames contain more pixel information that must be compressed.

Table.15. Bitrate Results Across Resolution Scale Index

Resolution Scale Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
5	720	690	660	620
10	820	780	740	690
15	930	890	840	790
20	1040	990	940	880
25	1160	1100	1040	970
30	1290	1220	1160	1080

The Table.15 presents the bitrate requirements across increasing resolution scales. At the resolution index of 5, the Deep Convolutional Autoencoder Compression model requires 720 kbps for encoding the frame. The Recurrent Neural Video Compression model reduces the requirement to 690 kbps, while the Attention-Based Neural Compression framework requires 660 kbps. The proposed SVRD-NC method further reduces the requirement to 620 kbps. When the resolution index increases to 30, the convolutional autoencoder model requires 1290 kbps while the proposed method uses 1080 kbps. The numerical reduction reaches approximately 210 kbps at the highest resolution level. This improvement occurs because the spatially variable rate-distortion optimization distributes bitrate resources more efficiently across frame regions. Smooth areas that contain limited structural information receive fewer bits, while complex regions receive additional coding capacity. This adaptive allocation mechanism improves compression efficiency and reduces the total transmission requirement.

#### 4.12 COMPRESSION RATIO RESULTS BASED ON TEXTURE DENSITY INDEX

The compression ratio is evaluated across texture density index levels, which represent the density of texture patterns present in a frame. Texture density acts as themetric because highly textured scenes usually require more coding resources.

Table.16. Compression Ratio Results Across Texture Density Index

Texture Density Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
5	21.2	22.4	23.6	25.1
10	20.5	21.7	22.9	24.4
15	19.8	21.0	22.2	23.6
20	19.1	20.3	21.4	22.9
25	18.4	19.6	20.8	22.2
30	17.8	19.0	20.1	21.5

The Table.16 shows the compression ratio performance across increasing texture density levels. At the density index of 5, the Deep Convolutional Autoencoder Compression model achieves a compression ratio of 21.2 while the Recurrent Neural Video Compression model achieves 22.4. The Attention-Based Neural Compression framework improves the ratio to 23.6. The proposed SVRD-NC framework achieves the highest ratio of 25.1. When the density index increases to 30, the compression ratios decrease for all methods because textured scenes contain more visual information. However, the proposed model maintains a ratio of 21.5 while the convolutional autoencoder model decreases to 17.8. The numerical improvement ranges between 3 and 4 compression ratio units. The spatial attention mechanism identifies regions that contain complex textures and assigns appropriate coding resources. This adaptive allocation strategy reduces redundant encoding in simpler areas and therefore improves the overall compression efficiency.

#### 4.13 MEAN SQUARED ERROR RESULTS BASED ON NOISE LEVEL INDEX

The Mean Squared Error performance is evaluated across noise level index values, which represent the amount of visual noise present in the video frame.

Table.17. MSE Results Across Noise Level Index

Noise Level Index	Deep Convolutional Autoencoder Compression	Recurrent Neural Video Compression	Attention-Based Neural Compression	Proposed SVRD-NC
5	0.010	0.009	0.008	0.006
10	0.012	0.011	0.010	0.007
15	0.014	0.013	0.011	0.009
20	0.017	0.015	0.013	0.010
25	0.020	0.018	0.015	0.012
30	0.024	0.021	0.018	0.014

The Table.17 presents the reconstruction error across increasing noise levels. At the noise index of 5, the Deep Convolutional Autoencoder Compression model produces an error value of 0.010 while the Recurrent Neural Video Compression model produces 0.009. The Attention-Based Neural Compression framework reduces the error to 0.008. The proposed SVRD-NC model achieves the lowest error value of 0.006. When the noise index increases to 30, the convolutional autoencoder

model produces 0.024 while the proposed model produces 0.014. The numerical improvement reaches approximately 0.010. The attention mechanism identifies noisy regions that affect reconstruction quality. The rate-distortion optimization allocates additional coding resources to these regions which reduces reconstruction error. As a result, the proposed method produces reconstructed frames that that has maintained higher fidelity even when the input frames contain significant visual noise.

## 5. CONCLUSION

The study presents a neural video compression framework that integrates the spatially variable rate-distortion optimization mechanism with the deep neural coding architecture. The proposed SVRD-NC framework improves the efficiency of video compression by allocating the bitrate according to the spatial characteristics of each frame region. The spatial attention module identifies visually significant areas that contain edges, textures, and motion structures, while the adaptive rate-distortion optimization assigns additional coding resources to those regions. This mechanism has proved that the neural encoder preserves the structural and perceptual quality that appears in complex visual areas. The experimental evaluation demonstrates that the proposed framework achieves consistent improvements across multiple performance metrics. The PSNR results show that the SVRD-NC method achieves a maximum value of 37.1 dB, which exceeds the performance of the Deep Convolutional Autoencoder Compression model that reaches 34.2 dB under comparable complexity conditions. The SSIM evaluation indicates that the proposed method achieves 0.98 structural similarity, which reflects stronger perceptual reconstruction accuracy than the existing methods. The bitrate evaluation reveals that the proposed framework reduces the transmission requirement to 620 kbps at the lower resolution scale while maintaining visual fidelity. Furthermore, the compression ratio increases to 25.1, which indicates efficient data reduction across diverse video scenes. The reconstruction accuracy also improves because the Mean Squared Error decreases to 0.006 at the lowest noise level. These results confirm that the spatially adaptive optimization strategy improves compression efficiency while preserving perceptual video quality.

## REFERENCES

- [1] L. Tao, W. Gao, G. Li and C. Zhang, "Adanic: Towards Practical Neural Image Compression via Dynamic Transform Routing", *Proceedings of International Conference on Computer Vision*, pp. 16879-16888, 2023.
- [2] I. Dror and O. Hadar, "Improved Perceptual Quality of Traffic Signs and Lights for the Teleoperation of Autonomous Vehicle Remote Driving via Multi-Category Region of Interest Video Compression", *Entropy*, Vol. 27, No. 7, pp. 1-8, 2025.
- [3] T. Richter, "Spatial Constant Quantization in JPEG XR is Nearly Optimal", *Proceedings of International Conference on Data Compression*, pp. 79-88, 2010.
- [4] Y. Qi, R. Feng, Z. Zhang and Z. Chen, "Image Coding for Machines based on Non-Uniform Importance Allocation", *Proceedings of International Conference on Visual Communications and Image Processing*, pp. 1-5, 2023.

- [5] T. Partanen, M. Hoang, A. Mercat, J. Sainio and J. Vanne, "Energy-Efficient Saliency-Guided Video Coding Framework for Real-Time Applications", *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, Vol. 15, No. 1, pp. 44-57, 2025.
- [6] J. Laitinen, T. Partanen, A. Mercat, J. Vanne, M. Hannuksela, H. Zhang and F. Cricri, "Feasibility Study of Multi-Layer VVC Coding Scheme for Hybrid Machine-Human Consumption", *Proceedings of International Conference on Multimedia and Expo*, pp. 1-6, 2024.
- [7] J. Li and X. Hou, "Object-Fidelity Remote Sensing Image Compression with Content-Weighted Bitrate Allocation and Patch-based Local Attention", *IEEE Transactions on Geoscience and Remote Sensing*, Vol. 62, pp. 1-14, 2024.
- [8] I. Dror and O. Hadar, "Improved Perceptual Quality of Traffic Signs and Lights for the Teleoperation of Autonomous Vehicle Remote Driving via Multi-Category Region of Interest Video Compression", *Entropy*, Vol. 27, No. 7, pp. 1-7, 2025.
- [9] Z. Pan, J. Chen, B. Peng, J. Lei, F.L. Wang, N. Ling and S. Kwong, "Efficient Chroma Intra Prediction via Exemplar Colorization Network for Versatile Video Coding", *IEEE Transactions on Multimedia*, Vol. 27, pp. 4713-4724, 2025.
- [10] W. Liu, W. Gao, G. Li, S. Ma, T. Zhao and H. Yuan, "Enlarged Motion-Aware and Frequency-Aware Network for Compressed Video Artifact Reduction", *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 34, No. 10, pp. 10339-10352, 2024.
- [11] G. Valenzise, M. Quach, D. Tian, J. Pang and F. Dufaux, "Point Cloud Compression", *Immersive Video Technologies*, pp. 357-385, 2023.
- [12] Y. Wang and C. Wu, "A block based Wyner-Ziv Video Codec", *Proceedings of International Congress on Image and Signal Processing*, Vol. 1, pp. 1-5, 2010.