

This is a pre print version of the following article:

DeepWiVe: Deep-Learning-Aided Wireless Video Transmission / Tung, T.; Gunduz, D.. - In: IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS. - ISSN 0733-8716. - 40:9(2022), pp. 1-1.  
[10.1109/JSAC.2022.3191354]

*Terms of use:*

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

05/05/2026 19:09

(Article begins on next page)

# DeepWiVe: Deep-Learning-Aided Wireless Video Transmission

Tze-Yang Tung and Deniz Gündüz

Information Processing and Communications Lab (IPC-Lab),  
Imperial College London, UK

{tze-yang.tung14, d.gunduz}@imperial.ac.uk

**Abstract**—We present *DeepWiVe*, the first-ever end-to-end joint source-channel coding (JSCC) video transmission scheme that leverages the power of deep neural networks (DNNs) to directly map video signals to channel symbols, combining video compression, channel coding, and modulation steps into a single neural transform. Our DNN decoder predicts residuals without distortion feedback, which improves video quality by accounting for occlusion/disocclusion and camera movements. We simultaneously train different bandwidth allocation networks for the frames to allow variable bandwidth transmission. Then, we train a bandwidth allocation network using reinforcement learning (RL) that optimizes the allocation of limited available channel bandwidth among video frames to maximize overall visual quality. Our results show that *DeepWiVe* can overcome the *cliff-effect*, which is prevalent in conventional separation-based digital communication schemes, and achieve graceful degradation with the mismatch between the estimated and actual channel qualities. *DeepWiVe* outperforms H.264 video compression followed by low-density parity check (LDPC) codes in all channel conditions by up to 0.0462 on average in terms of the multi-scale structural similarity index measure (MS-SSIM), while beating H.265 + LDPC by up to 0.0058 on average. We also illustrate the importance of optimizing bandwidth allocation in JSCC video transmission by showing that our optimal bandwidth allocation policy is superior to the naïve uniform allocation. We believe this is an important step towards fulfilling the potential of an end-to-end optimized JSCC wireless video transmission system that is superior to the current separation-based designs.

## I. INTRODUCTION

Video content contributes to more than 80% of Internet traffic and the percentage is only expected to increase [1]. Video compression is widely used to reduce the bandwidth requirement when transmitting video signals wirelessly. This follows the modular approach employed in almost all wireless video transmission systems, where the end-to-end transmission problem is divided into two: (1) a source encoder that compresses the video into a sequence of bits of the shortest possible length such that a reconstruction of the original video is possible within an allowable distortion; and (2) a channel encoder that introduces redundancies such that the compressed bits are protected against channel errors and interference. A diagram of this modular-based approach is shown in Fig. 1.

This separate source and channel coding design provides modularity and allows independent optimization of each component. It has been applied successfully in a large variety

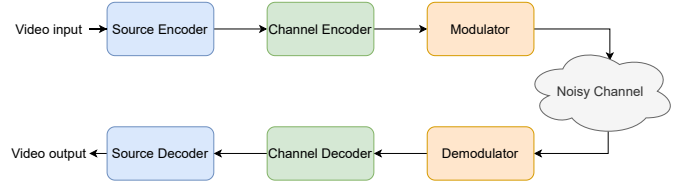


Fig. 1. Diagram of a typical separation-based digital video delivery system employed by almost all communication systems today.

of applications from on-demand mobile video streaming to video conferencing and digital TV broadcasting. However, the limits of the separation-based designs are beginning to rear, with the emergence of more demanding and challenging video delivery applications, such as wireless virtual reality (VR) and drone-based surveillance systems, which have ultra-low latency requirements, suffer from highly unpredictable channel conditions, and need to be implemented on energy limited mobile devices.

In the context of wireless video transmission, separation-based designs lead to what is known as the *cliff-effect*. That is, when the channel condition deteriorates below the level anticipated by the channel encoder, the source information becomes irrecoverable. This leads to a cliff edge deterioration of the system performance. As a result, most current systems operate at a much more conservative transmission rate than that is suggested by the instantaneous channel capacity, and employ additional error correction mechanisms through automatic repeat requests (ARQ).

An alternative to the separation-based architecture is joint source-channel coding (JSCC). The most natural JSCC approach continues to employ separate modules for compression and communication, but jointly optimizes various parameters of these modules in a cross-layer framework. While there have been many such proposals over the years [2]–[9], these techniques typically do not provide sufficient gains to justify the significant increase in system complexity.

A more fundamental approach is to design the transmission system from scratch, without considering any digital interface in between. The best example for such an approach is analog communications, such as AM/FM radio or analog TV, where the information is directly modulated onto the carrier waveform without any compression. Analog communication can overcome the cliff-effect problem by showing graceful degradation with the channel parameters. From a fundamental

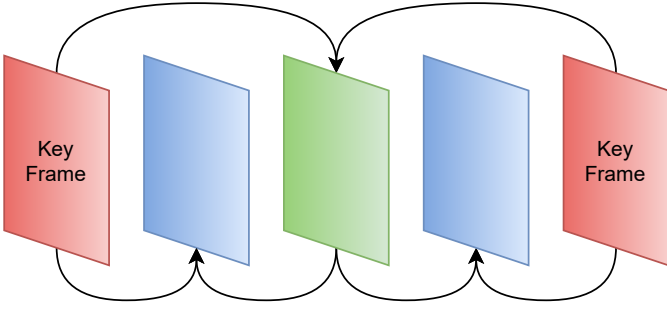


Fig. 2. Diagram of a typical interpolation structure used in video compression algorithm.

information theoretic perspective, when transmitting independent Gaussian samples over an additive white Gaussian noise (AWGN) channel, with one sample per channel use on average, uncoded transmission, where each sample is simply scaled and transmitted, meets the theoretical Shannon bound [10]. With digital transmission, the same performance can only be achieved by vector-quantizing an arbitrarily long sequence of source samples, followed by a capacity achieving channel code. Benefits of analog transmission has also been shown in various multi-user scenarios [11], [12]. However, analog modulation cannot exploit the available bandwidth efficiently, and the optimality of simple uncoded transmission does not generalize to bandwidth-mismatched scenarios. However, despite this theoretical suboptimality, analog modulation approaches to image and video transmission have gained recent popularity [13]–[16], mainly due to their low computational complexity and the graceful degradation with channel quality.

Recently, it has been shown in [17], [18] that deep neural networks (DNNs) can break the complexity barrier in designing effective JSCC schemes, focusing particularly on the image transmission problem. The approach there is to train DNN models in an autoencoder architecture with a non-trainable channel layer between the encoder and decoder. The authors showed that this approach not only provides graceful degradation with channel quality, but can also achieve results similar to or superior than state-of-the-art separation-based digital designs. An extension to this work [19] further shows that JSCC is successively refineable; that is, an image can be transmitted in stages, and each additional channel block further refines the quality of the decoded image, with almost no additional cost.

Herein, we propose a deep learning-based JSCC solution for wireless video transmission, called *DeepWiVe*, which is trained to optimize the reconstructed video quality in an end-to-end fashion. Although the problem of video compression using deep learning has received significant attention [20]–[28], no prior work has considered deep learning aided wireless video delivery.

The core idea behind most video compression algorithms is to exploit the temporal correlations across video frames. In a standard video sequence, motion differences between successive frames is typically very small, and video compression algorithms can be very efficient by identifying intermittent key frames, which are compressed and decoded independently,

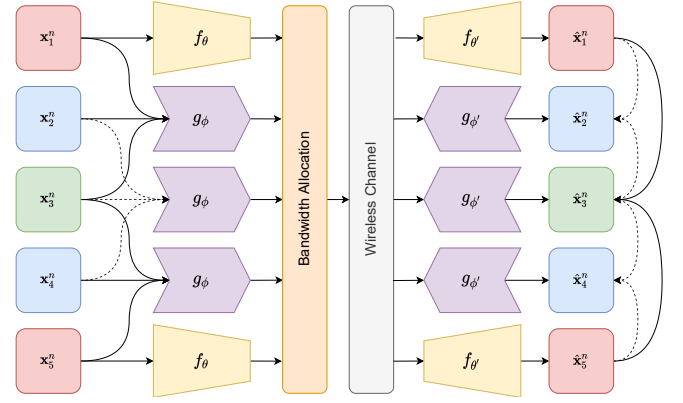


Fig. 3. *DeepWiVe* system overview.

and conveying only the motion and residual information for the remaining frames, thereby exploiting temporal redundancy. Here, the residual information refers to the difference between the true frame and the motion compensated key frame. Fig. 2 shows a typical interpolation structure in a video compression algorithm. For compression algorithms, the residual information is available at the encoder as the encoder can simply decode the compressed frames and observe the difference between the source and its reconstruction. On the other hand, for JSCC, the residual depends on the channel condition during transmission. Therefore, the residual is not known at the encoder. To overcome this, we propose to use a DNN to predict the residual, without the need for distortion feedback.

In *DeepWiVe* we directly map the video sequence into the channel vector under a channel bandwidth constraint for the transmission of a group-of-pictures (GoP). Similarly to [19], we employ variable bandwidth transmission and allocate the bandwidth dynamically using reinforcement learning (RL) to train a bandwidth allocation network that optimizes the bandwidth utilization of each frame. Our results show that *DeepWiVe* can meet or beat industry standard video compression codecs, such as H.264, combined with state-of-the-art channel codes, such as low density parity check (LDPC) codes, in almost all channel conditions tested, while achieving graceful degradation of video quality with respect to channel quality, thereby avoiding the *cliff-effect*.

The contributions of this paper are summarized as follows:

- 1) We propose *DeepWiVe*, a JSCC-based wireless video transmission scheme leveraging DNNs to jointly compress and channel code video frames in an end-to-end manner to maximize the end video quality.
- 2) We train our DNN decoder to predict residuals without the need for distortion feedback.
- 3) We optimize the bandwidth allocation among video frames by training a bandwidth allocation network using RL, and show that doing so achieves superior performance than naïve uniform allocation.
- 4) We show that it is possible to achieve variable bandwidth transmission by simply training different bandwidth allocation networks.
- 5) Numerical results show that our proposed *DeepWiVe* is superior to industry standard H.264 [29] codec with state-of-the-art LDPC channel codes [30] in all channel

conditions tested and can avoid the cliff-effect. It also beats H.265 [31] using the same channel codes when evaluated in terms of the multi-scale structural similarity index measure (MS-SSIM).

## II. RELATED WORK

JSCC for video delivery has consistently received attention over the years. The earliest work we could find is [6], which studies the problem of video multicast to heterogeneous receivers. They approached the problem from the receivers' perspective, where the source video is encoded in a hierarchical manner, with each layer of the hierarchy distributed on a separate network channel. Each receiver then adapts to its local channel capacity by adjusting the number of layers it decodes. In a similar line of work, [32] uses scalable video coding (SVC), which encodes the source video into multiple bitstreams, with a base layer that represents the lowest supported quality and a set of enhancement layers representing versions of the video at different qualities. Since the base layer and lower quality layers require more error protection than higher quality layers, unequal error protection (UEP) of packets is used. To determine the optimal channel code rate for each layer such that the average distortion is minimized, they devised a low complexity search algorithm to find the optimal choice of channel code rates among a set of available rates. There is a large body of works that use source coding schemes similarly to SVC and minimize the end-to-end distortion by jointly optimizing various parameters of the source and channel codes [4], [5], [7]–[9], [33]. However, none of the scheme proposed were able to achieve adequate performance gains for the increased complexity they introduce as a result of the optimization of the parameters. Moreover, with the exception of [6] and [8], which model the channel as a physical link with a certain error probability, the above works mainly consider time-varying channels with proposed schemes aimed at adapting to channel variations. In this work, we will focus on AWGN channels.

A completely refreshed approach to JSCC video delivery, called SoftCast, utilizing low complexity methods to map videos or images from the pixel domain to channel symbols directly was first introduced in [13]. Their scheme involves a hybrid digital and analog design by leveraging frequency domain sparsity. Since then, various works have improved upon [13] by optimizing different aspects of the hybrid digital and analog design. In [13]–[16], frequency domain sparsity is exploited by utilizing compressed sensing to reduce the bandwidth requirement. In [34], the power allocation problem is addressed by optimizing the division of frequency domain coefficients into chunks. In [35], the proposed scheme separates a base layer (low frequency coefficients) from the enhancement layer (high frequency coefficients) and sends the base layer using a digital separation-based system, while the enhancement layer is sent using a scheme similar to Soft-Cast. This method also improves power allocation since low frequency coefficients have larger values than high frequency coefficients in natural images so the majority of the signal energy is sent through the digital system. In [36], [37], a

similar approach is utilized, however, instead of the high frequency components, they send the gradient of the image as the enhancement information. In [38], the prior used for the estimation is modified to obtain better reconstruction quality (SoftCast assumes a Gaussian prior). In [39], the variability of temporal redundancies between frames is addressed by introducing adaptive GoP sizes. Although these methods have been shown to overcome the *cliff-effect*, they are not competitive to separation-based schemes when it comes to video quality and cannot exploit the available bandwidth, or adapt to channel and network conditions dynamically.

Following the success of DNNs in various image and video intelligence tasks, there has been a growing interest in employing them for video compression [20]–[28]. Although some recent works have reported competitive or superior performance with respect to state-of-the-art video compression standards H.264/5, none of the works have considered the wireless video delivery problem, taking into account the channel conditions. In principle, the works treat the communication layer simply as a perfect bit pipe. As a result, the shortcomings of separation-based schemes, such as the *cliff-effect*, are inherent to those works if applied in a wireless communication scenario.

The closest prior art to our work are [17]–[19], [40], [41], which explore the JSCC problem, but they focus on image transmission. In the context of video transmission, there are unique challenges that sets it apart from simple image transmission. Namely, exploiting the inter-frame redundancies to improve coding efficiency and to optimize resource allocation across the frames. Therefore, extending the problem from image to video transmission is not a trivial task.

## III. PROBLEM FORMULATION

We consider the problem of wireless video transmission in a constrained bandwidth setting. Consider a video sequence  $\mathbf{X} = \{\mathbf{X}^n\}_{n=1}^T$ , where  $\mathbf{X}^n = \{\mathbf{x}_1^n, \dots, \mathbf{x}_N^n\}$ ,  $\mathbf{x}_i^n \in \mathbb{R}^{H \times W \times 3}$ ,  $\forall i \in [1, N]$ , represents the  $n$ th GoP in the video sequence. Each frame  $\mathbf{x}_i^n$  is represented as a 24bit RGB image. We wish to design an encoding function  $E : \mathbb{R}^{TN \times H \times W \times 3} \mapsto \mathbb{C}^{Tk}$ , which maps the video sequence  $\mathbf{X}$  to a set of complex symbols  $\mathbf{z} = E(\mathbf{X}) \in \mathbb{C}^{Tk}$ , and a decoding function  $D : \mathbb{C}^{Tk} \mapsto \mathbb{R}^{TN \times H \times W \times 3}$ , which maps a noise corrupted version of the encoder output  $\mathbf{y} = \mathbf{z} + \mathbf{n}$ , caused by the wireless channel, to an approximate reconstruction of the original video sequence  $\hat{\mathbf{X}} = D(\mathbf{y})$ .

In this setting, we restrict the number of channel uses to  $k$  per GoP, which can be considered as a bandwidth constraint and we define the *bandwidth compression ratio* as

$$\rho = \frac{k}{3HWN}. \quad (1)$$

We consider an additive white Gaussian noise (AWGN) channel  $\mathbf{n} \sim \text{CN}(0, \sigma^2 \mathbf{I})$  that follows a complex Gaussian noise distribution with zero mean and covariance  $\sigma^2 \mathbf{I}$  ( $\mathbf{I}$  is the identity matrix). We impose an average power constraint at the transmitter, such that

$$\frac{1}{Tk} \mathbb{E}_{\mathbf{z}} [||\mathbf{z}||_2^2] \leq P, \quad (2)$$

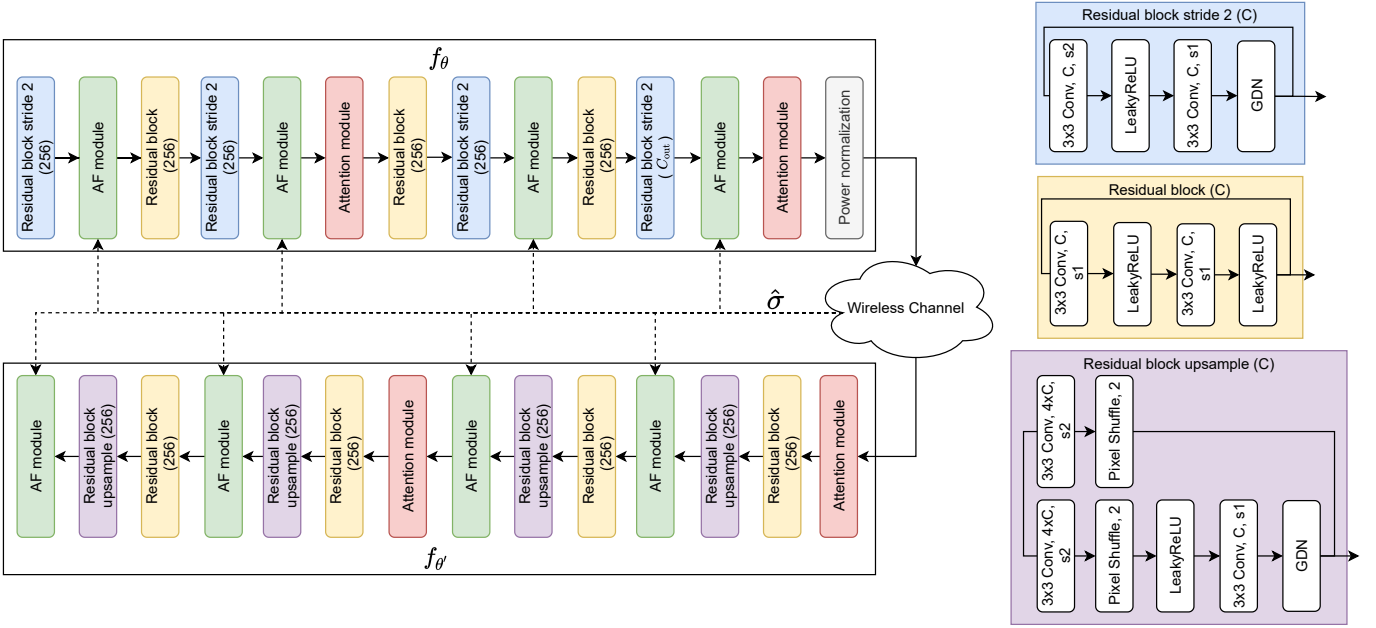


Fig. 4. Key frame encoder/decoder ( $f_\theta, f_{\theta'}$ ) network architectures.

where the expectation is taken over the distribution of the encoder output. Accordingly, the channel signal-to-noise ratio (SNR) is defined as

$$\text{SNR} = 10 \log_{10} \left( \frac{P}{\sigma^2} \right) \text{ dB}. \quad (3)$$

We measure the average quality of the reconstructed video using two metrics: peak signal-to-noise ratio (PSNR) and MS-SSIM. They are defined as

$$\text{PSNR}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{1}{TN} \sum_{n=1}^T \sum_{i=1}^N 10 \log_{10} \left( \frac{255^2}{l_{\text{PSNR}}(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n)} \right) \text{ dB}, \quad (4)$$

and

$$\text{MS-SSIM}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{1}{TN} \sum_{n=1}^T \sum_{i=1}^N 1 - l_{\text{MS-SSIM}}(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n), \quad (5)$$

where

$$l_{\text{PSNR}}(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n) = \frac{1}{3HW} \|\mathbf{x}_i^n - \hat{\mathbf{x}}_i^n\|_2^2, \quad (6)$$

and

$$l_{\text{MS-SSIM}}(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n) = 1 - \text{MS-SSIM}(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n). \quad (7)$$

MS-SSIM is defined between two frames as

$$\text{MS-SSIM}(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n) \quad (8)$$

$$= [l_M(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n)]^{\alpha_M} \prod_{j=1}^M [c_j(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n)]^{\beta_j} [s_j(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n)]^{\gamma_j}, \quad (9)$$

where

$$l_M(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n) = \frac{2\mu_{\mathbf{x}_i^n} \mu_{\hat{\mathbf{x}}_i^n} + c_1}{\mu_{\mathbf{x}_i^n}^2 + \mu_{\hat{\mathbf{x}}_i^n}^2 + c_1}, \quad (10)$$

$$c_j(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n) = \frac{2\sigma_{\mathbf{x}_i^n} \sigma_{\hat{\mathbf{x}}_i^n} + c_2}{\sigma_{\mathbf{x}_i^n}^2 + \sigma_{\hat{\mathbf{x}}_i^n}^2 + c_2}, \quad (11)$$

$$s_j(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n) = \frac{\sigma_{\mathbf{x}_i^n \hat{\mathbf{x}}_i^n} + c_3}{\sigma_{\mathbf{x}_i^n} \sigma_{\hat{\mathbf{x}}_i^n} + c_3}. \quad (12)$$

Here,  $\mu_{\mathbf{x}_i^n}$ ,  $\sigma_{\mathbf{x}_i^n}^2$ ,  $\sigma_{\mathbf{x}_i^n \hat{\mathbf{x}}_i^n}^2$  are the mean and variance of  $\mathbf{x}_i^n$ , and the covariance between  $\mathbf{x}_i^n$  and  $\hat{\mathbf{x}}_i^n$ , respectively.  $c_1$ ,  $c_2$ , and  $c_3$  are coefficients for numeric stability;  $\alpha_M$ ,  $\beta_j$ , and  $\gamma_j$  are the weights for each of the components. Each  $c_j(\cdot)$  and  $s_j(\cdot)$  are computed at a different downsampled scale of  $(\mathbf{x}_i^n, \hat{\mathbf{x}}_i^n)$ . We use the default parameter values of  $(\alpha_M, \beta_j, \gamma_j)$  provided by the original paper [42]. MS-SSIM has been shown to be as good and better at approximating human visual perception than the more simplistic structural similarity index (SSIM) on different subjective image and video databases.

The overall goal of our design is to maximize the video quality, measured by either Eqn. (4) or (5), between the input video  $\mathbf{X}$  and its reconstruction at the decoder  $\hat{\mathbf{X}}$ , under the given constraints on the available bandwidth ratio  $\rho$  and the average power  $P$ .

#### A. Joint Source-Channel Video Coding

In this section, we present our proposed DNN-based joint source-channel video encoding and decoding scheme. We will deconstruct the design of the encoder ( $E$ ) and decoder ( $D$ ) into three parts: the key frame encoder/decoder ( $f_\theta, f_{\theta'}$ ), parameterized by  $(\theta, \theta')$ , the interpolation encoder/decoder ( $g_\phi, g_{\phi'}$ ), parameterized by  $(\phi, \phi')$ , and the bandwidth allocation function  $q_\psi$ , parameterized by  $\psi$ . We will represent all these functions with DNNs, where the parameters of the functions

correspond to the weights of these DNNs. An overview of our scheme is shown in Fig. 3.

One of the drawbacks in prior works was the need to train multiple networks, one for each channel condition. To address this issue, [40] proposed an attention feature (AF) module, motivated by resource assignment strategies in traditional JSCC schemes [43], which allows the network to learn to assign different weights to different features for a given SNR. This is achieved by deliberately randomizing the channel SNR during training, and providing the AF modules with the current SNR. By doing so, the results in [40] show that a single model whose parameters are adjusted to the channel SNR with the help of the AF modules perform at least as well as the models trained for each SNR individually. We adopt the AF module proposed by [40] in *DeepWive* to obtain a single model that can work over a range of SNRs.

The encoding and decoding procedures are described herein. Consider the  $n$ th GoP,  $\mathbf{X}^n = \{\mathbf{x}_1^n, \dots, \mathbf{x}_N^n\}$ . The last ( $\mathbf{x}_N^n$ ) frame is called the key frame and is compressed and transmitted using the key frame encoder  $f_\theta : \mathbb{R}^{H \times W \times 3} \mapsto \mathbb{C}^k$ ,

$$\mathbf{z}_i^n = f_\theta(\mathbf{x}_i^n, \hat{\sigma}^2), \quad i = N, \quad (13)$$

where  $\hat{\sigma}^2$  is the estimated channel noise power at the transmitter. Here, each element of  $\mathbf{z}_i^n$ , denoted by  $z_{i,j}^n$ , represents the in-phase (I) and quadrature (Q) components of a complex channel symbol. We note that, while the channel input/output values are complex, we employ real-valued DNN architecture. The mapping between real network outputs and complex channel inputs (or vice versa) is achieved by pairing consecutive real values at the output of the encoder DNN.

The values in the complex latent vector are first normalized according to

$$\hat{z}_{i,j}^n = \sqrt{kP} \frac{z_{i,j}^n}{\sqrt{(\mathbf{z}_i^n)^H \mathbf{z}_i^n}}, \quad j = 1, \dots, k, \quad (14)$$

where  $P$  is the power constraint,  $k$  is the bandwidth constraint, and  $H$  refers to the Hermitian transpose.

These values are then directly sent through the channel as,

$$\hat{\mathbf{y}}_i^n = \hat{\mathbf{z}}_i^n + \mathbf{n}. \quad (15)$$

Consequently, the key frame decoder  $f_{\theta'} : \mathbb{C}^k \mapsto \mathbb{R}^{H \times W \times 3}$  that maps the noisy latent vector  $\hat{\mathbf{y}}_i^n \in \mathbb{C}^k$  observed at the receiver back to the original frame domain  $\hat{\mathbf{x}}_i^n \in \mathbb{R}^{H \times W \times 3}$  is defined as:

$$\hat{\mathbf{x}}_i^n = f_{\theta'}(\hat{\mathbf{y}}_i^n, \hat{\sigma}^2), \quad i = N. \quad (16)$$

The loss between the original frame  $\mathbf{x}_i^n$  and the reconstructed frame  $\hat{\mathbf{x}}_i^n$  is computed using Eqn. (6) or (7) depending on which performance measure is being used. The network weights  $(\theta, \theta')$  are then updated via backpropagation with respect to the gradient of the loss.

The network architectures of the key frame encoder and decoder are shown in Fig. 4. The Pixel Shuffle module, first proposed in [44], increases the height and width of the input tensor dimensions efficiently by exchanging channel dimensions for height and width dimensions. The GDN layer refers to Generalised Divisive Normalization, initially proposed in [45], and has been shown to be effective in density

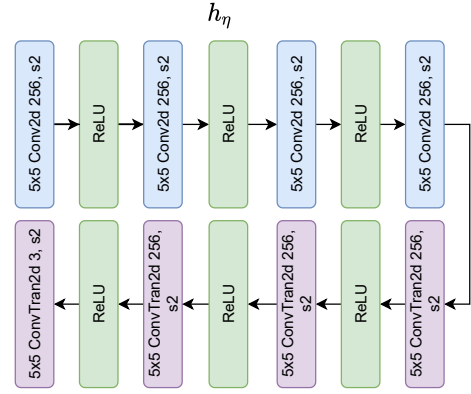


Fig. 5. Architecture of the SSF estimator network  $h_\eta$ .

modeling and compression of images. The Attention layer refers to the simplified attention module proposed in [46], which reduces the computation cost of the attention module originally proposed in [47]. The attention mechanism has been used in both [46] and [47] to improve the compression efficiency by focusing the neural network on regions in the image that require higher bit rate. The network in Fig. 4 is fully convolutional, therefore it can accept input of any height ( $H$ ) and width ( $W$ ), making it versatile to any video resolution.

For the remaining frames, i.e.,  $\mathbf{x}_i^n$ ,  $i = 1, 2, \dots, N - 1$ , we use the interpolation encoder  $g_\phi(\cdot)$  to encode the motion information  $(\mathbf{f}_{i-t}^n, \mathbf{f}_{i+t}^n)$  and residual information  $(\mathbf{r}_{i-t}^n, \mathbf{r}_{i+t}^n)$  of  $\mathbf{x}_i^n$  with respect to two reference frames  $(\bar{\mathbf{x}}_{i-t}^n, \bar{\mathbf{x}}_{i+t}^n)$  that are  $t$  frames away from the current frame. The reference frames  $\bar{\mathbf{x}}_{i-t}^n, \bar{\mathbf{x}}_{i+t}^n$  are what the encoder expects the corresponding reconstructed frames  $\hat{\mathbf{x}}_{i-t}^n, \hat{\mathbf{x}}_{i+t}^n$  to be. This is done via channel emulation and decoding at the transmitter side to obtain an approximation of what the transmitter expects the receiver to reconstruct.

We follow the same interpolation structure as the one presented in Fig. 2 for a GoP size  $N = 4$ . That is, for  $i = 2$ ,  $t = 2$ , while for  $i = 1, 3$ ,  $t = 1$ . We define  $\bar{\mathbf{x}}_0^n = \bar{\mathbf{x}}_{N-1}^n$ , and assume that the GoPs are encoded and decoded sequentially, such that the frames from the previous decoded GoP are available as reference for the current GoP. To interpolate  $\mathbf{x}_i^n$  from  $\bar{\mathbf{x}}_{i-t}^n$  and  $\bar{\mathbf{x}}_{i+t}^n$ , motion information, such as optical flow [48], is usually used to warp a reference image by translating the pixels in the reference image according to the optical flow vectors. These optical flow vectors describe the horizontal and vertical translations of each pixel in a reference image in order to transform it into the target image. The difference between the optical flow transformed image and the true target image is called the *residual*, which is used to capture information that cannot be described by optical flow, such as occlusion/disocclusion and camera movements.

To estimate the motion information  $(\mathbf{f}_{i-t}^n, \mathbf{f}_{i+t}^n)$ , instead of using optical flow, we use scaled space flow (SSF), which was first proposed by [26] as a more general description of pixel warping than optical flow. The idea is to blur regions of the frame where the motion is difficult to model using traditional pixel warping and instead compensate those regions using the residual. To that end, in scale-space warping (SSW), a frame

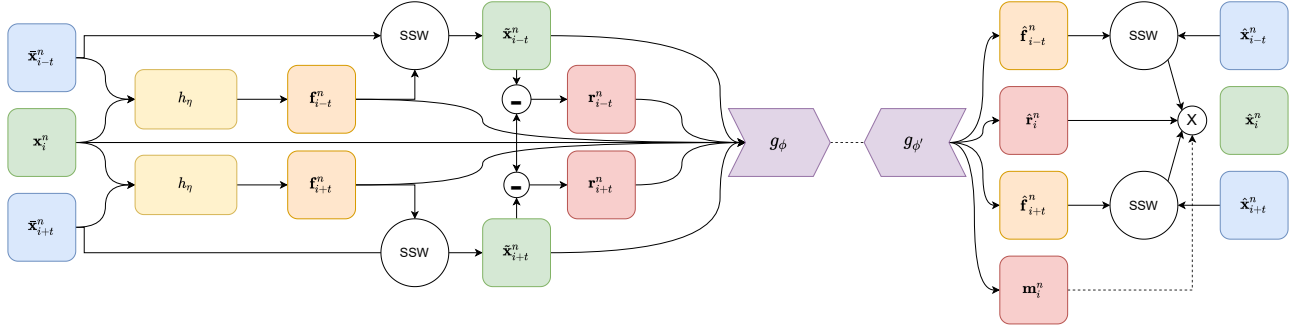


Fig. 6. Information flow over the interpolation network.

is first transformed into a fixed-resolution volume  $\bar{\mathbf{X}}_{i+t}^n = [\bar{\mathbf{x}}_{i+t}^n, \bar{\mathbf{x}}_{i+t}^n \otimes G(\sigma_0), \bar{\mathbf{x}}_{i+t}^n \otimes G(2\sigma_0), \dots, \bar{\mathbf{x}}_{i+t}^n \otimes G(2^{V-1}\sigma_0)]$ , where  $\bar{\mathbf{x}}_{i+t}^n \otimes G(\sigma_0)$  denotes Gaussian blurring of the frame  $\bar{\mathbf{x}}_{i+t}^n$  by convolving  $\bar{\mathbf{x}}_{i+t}^n$  with a Gaussian kernel  $G(\sigma_0)$  with standard deviation  $\sigma_0$  and  $\otimes$  is the convolution operation.  $\bar{\mathbf{X}}_{i+t}^n \in \mathbb{R}^{H \times W \times (V+1)}$  represents a progressively blurred version of  $\bar{\mathbf{x}}_{i+t}^n$ , which can be sampled at continuous points via trilinear interpolation. The scaled space flow  $\mathbf{f}_{i+t}^n \in \mathbb{R}^{H \times W \times 3}$  that warps frame  $\bar{\mathbf{x}}_{i+t}^n$  to an approximation of  $\mathbf{x}_i^n$  denoted by  $\tilde{\mathbf{x}}_{i+t}^n$  is then defined as

$$\tilde{\mathbf{x}}_{i+t}^n = \text{SSW}(\bar{\mathbf{x}}_{i+t}^n, \mathbf{f}_{i+t}^n) \quad (17)$$

$$\begin{aligned} \text{s.t. } \tilde{\mathbf{x}}_{i+t}^n[x, y] \\ = \bar{\mathbf{X}}_{i+t}^n[x + \mathbf{f}_{i+t}^n[x, y, 1], y + \mathbf{f}_{i+t}^n[x, y, 2], \mathbf{f}_{i+t}^n[x, y, 3]] \end{aligned} \quad (18)$$

To estimate the scaled space flow  $\mathbf{f}_{i+t}^n$ , we use the network architecture  $h_\eta: \mathbb{R}^{H \times W \times 6} \mapsto \mathbb{R}^{H \times W \times 3}$  proposed in [26],

$$\mathbf{f}_{i+t}^n = h_\eta(\mathbf{x}_i^n, \hat{\mathbf{x}}_{i+t}^n). \quad (19)$$

Fig. 5 shows the architecture of the SSF estimator network  $h_\eta$ . This architecture is similar to U-Net [49], which was first proposed for bioimage segmentation. The architecture progressively downsamples the input using convolutional layers before upsampling it back to the frame dimensions. The channel dimension of the output SSF  $\mathbf{f}_{i+t}^n$ , instead of representing the three color channels, represent the  $(x, y, z)$  sampling points in the SSF volume  $\bar{\mathbf{X}}_{i+t}^n$ .

Given the above definition of SSF, the residual  $\mathbf{r}_{i+t}^n$  is defined as

$$\mathbf{r}_{i+t}^n = \mathbf{x}_i^n - \tilde{\mathbf{x}}_{i+t}^n. \quad (20)$$

The interpolation encoder  $g_\phi: \mathbb{R}^{H \times W \times 21} \mapsto \mathbb{C}^k$  defines the mapping

$$\begin{aligned} \mathbf{z}_i^n = g_\phi(\mathbf{x}_i^n, \bar{\mathbf{x}}_{i-t}^n, \bar{\mathbf{x}}_{i+t}^n, \mathbf{r}_{i-t}^n, \mathbf{r}_{i+t}^n, \mathbf{f}_{i-t}^n, \mathbf{f}_{i+t}^n, \hat{\sigma}^2), \\ i = 1, 2, \dots, N-1. \end{aligned} \quad (21)$$

The vector  $\mathbf{z}_i^n \in \mathbb{C}^k$  is power normalized according to Eqn. (14) and sent across the channel according to Eqn. (15).

Given the noisy  $\hat{\mathbf{y}}_i^n$ , the decoder first estimates the SSF, the residual, and a mask. That is, the interpolation decoder  $g_{\phi'}: \mathbb{C}^k \mapsto \mathbb{R}^{H \times W \times 12}$  defines the mapping

$$(\hat{\mathbf{f}}_{i-t}^n, \hat{\mathbf{f}}_{i+t}^n, \hat{\mathbf{r}}_i^n, \mathbf{m}_i^n) = g_{\phi'}(\hat{\mathbf{y}}_i^n, \hat{\sigma}^2), \quad (22)$$

where  $\hat{\mathbf{f}}_{i\pm t}^n \in \mathbb{R}^{H \times W \times 3}$ ,  $\hat{\mathbf{r}}_i^n \in \mathbb{R}^{H \times W \times 3}$ , and  $\mathbf{m}_i^n \in \mathbb{R}^{H \times W \times 3}$ .  $\mathbf{m}_{i,c}^n \in \mathbb{R}^{H \times W}$ ,  $c = 1, 2, 3$ , a 2D matrix in the third dimension of  $\mathbf{m}_i^n$ , satisfies:

$$\sum_{c=1}^3 \mathbf{m}_{i,c}^n = \mathbf{1}_{H \times W}. \quad (23)$$

That is, for each  $H$  and  $W$  index of the mask  $\mathbf{m}_i^n$ , the sum of values along the channel dimension is equal to 1, which is achieved by using the softmax activation. The reconstructed frame is then defined as:

$$\begin{aligned} \hat{\mathbf{x}}_i^n = (\mathbf{m}_i^n)_1 * \text{SSW}(\hat{\mathbf{x}}_{i-t}^n, \hat{\mathbf{f}}_{i-t}^n) + \\ (\mathbf{m}_i^n)_2 * \text{SSW}(\hat{\mathbf{x}}_{i+t}^n, \hat{\mathbf{f}}_{i+t}^n) + (\mathbf{m}_i^n)_3 * \hat{\mathbf{r}}_i^n, \end{aligned} \quad (24)$$

where  $*$  refers to element-wise multiplication. The predicted mask  $\mathbf{m}_i^n$  acts as a convex set of weights to sum the two motion compensated predictions and the predicted residual, such that the resultant prediction  $\hat{\mathbf{x}}_i^n$  remains within  $[0, 255]$ . A diagram of the interpolation structure described herein can be seen in Fig. 6. The architectures of  $g_\phi$  and  $g_{\phi'}$  are functionally the same as  $f_\theta$  and  $f_{\theta'}$ , except the size of the input tensor, which is the concatenation of  $(\mathbf{x}_i^n, \bar{\mathbf{x}}_{i-t}^n, \bar{\mathbf{x}}_{i+t}^n, \mathbf{r}_{i-t}^n, \mathbf{r}_{i+t}^n, \mathbf{f}_{i-t}^n, \mathbf{f}_{i+t}^n)$  along the channel dimension. As in the key frame, the network weights  $(\phi, \phi', \eta)$  are updated via backpropagation with respect to the gradient of the loss between the original and the reconstructed frame.

## B. Bandwidth Allocation

In the previous section, we have assumed that each frame utilizes the full bandwidth of  $k$  channel uses allowed for each GoP. In order to satisfy the bandwidth constraint defined in Section III, the encoder must decide how to allocate  $k$  channel uses to the  $N$  frames in a GoP. Intuitively, if the frame in consideration  $\mathbf{x}_i^n$  is exactly the same with respect to the reference frames  $(\bar{\mathbf{x}}_{i-t}^n, \bar{\mathbf{x}}_{i+t}^n)$  that it is interpolated from, then no information needs to be transmitted. On the other hand, if there is significant differences with respect to the reference frames, then more information needs to be sent in order to accurately interpolate the frame. Since the last frame of a previous GoP becomes the reference frame of the next GoP ( $\bar{\mathbf{x}}_N^n = \bar{\mathbf{x}}_0^{n+1}$ ), we formulate the problem of allocating available bandwidth in each GoP as a Markov decision process (MDP) and solve the optimal bandwidth allocation policy using reinforcement learning.

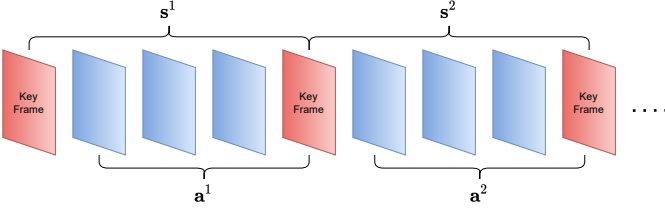


Fig. 7. Diagram illustrating the bandwidth allocation problem formulated as a MDP.

A MDP is defined by the tuple  $(\mathcal{S}, \mathcal{A}, P, r)$ , where  $\mathcal{S}$  is the set of states,  $\mathcal{A}$  is the action set,  $P$  is the probability transition kernel that defines the probability of one state transitioning to another state given an action, and  $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$  is the reward function. At each time step  $n$ , an agent observes state  $\mathbf{s}^n \in \mathcal{S}$  and takes an action  $\mathbf{a}^n \in \mathcal{A}$  based on its policy  $\pi : \mathcal{S} \mapsto \mathcal{A}$ . The state then transitions to  $\mathbf{s}^{n+1}$  according to the probability  $P(\mathbf{s}^{n+1} | \mathbf{s}^n, \mathbf{a}^n)$ , and the agent receives a reward  $r^n(\mathbf{s}^n, \mathbf{a}^n)$ . The objective is to maximize the expected sum of rewards  $J(\pi) = \mathbb{E}_{\mathbf{s}_1 \sim \omega_{\mathbf{s}_1}, \pi} [\sum_{i=1}^{\infty} \gamma^{(i-1)} r^i(\mathbf{s}^i, \mathbf{a}^i)]$ , where  $\omega_{\mathbf{s}_1}$  is the initial state distribution and  $\gamma \in (0, 1)$  is the reward discount factor to ensure convergence.

In the dynamic bandwidth allocation problem, we define the state at time step  $n$  as  $\mathbf{s}^n = \{\mathbf{M}^n, \mathbf{R}^n, \mathbf{F}^n, \hat{\sigma}^2\}$ , where

$$\mathbf{M}^n = \{\bar{\mathbf{x}}_i^n\}_{i=0}^N, \quad (25)$$

$$\mathbf{R}^n = \{\mathbf{r}_{i \pm t}^n\}_{i=1}^{N-1}, \quad (26)$$

$$\mathbf{F}^n = \{\mathbf{f}_{i \pm t}^n\}_{i=1}^{N-1}, \quad (27)$$

where  $\bar{\mathbf{x}}_0^n = \bar{\mathbf{x}}_N^{n-1}$ . The action set  $\mathcal{A}$  is the set of all the different ways the available bandwidth  $k$  can be allocated to each frame in the GoP. In order for the decoder functions  $f_{\theta'}$  and  $g_{\phi'}$  to be able to decode each frame that has been given different amounts of bandwidth, we use a result in [19], which showed that joint source-channel encoded images can be successively refined by sending increasingly more information. This is achieved by dividing the latent vectors  $\mathbf{z}_i^n$  into  $V$  equal sized blocks (i.e.,  $\mathbf{z}_i^n = \{\mathbf{z}_{i,1}^n, \dots, \mathbf{z}_{i,V}^n\}$ ,  $\mathbf{z}_{i,v}^n \in \mathbb{C}^{\frac{k}{V}}$ ,  $v = 1, \dots, V$ ), while randomly varying the number of blocks  $v_i^n$  of the latent code transmitted in each batch  $\mathbf{z}_i^n(v_i^n) = \{\mathbf{z}_{i,1}^n, \dots, \mathbf{z}_{i,v_i^n}^n\}$ ,  $v_i^n \leq V$ . This training process leads to the descending ordering of information from  $\mathbf{z}_{i,1}^n$  to  $\mathbf{z}_{i,V}^n$ . As such, each action represents  $\mathbf{a}^n = [v_1^n, \dots, v_N^n]$ . We implement this training process in the algorithm described in Section III-A by zeroing out the blocks in the latent vector not transmitted. As such, the action set is all the ways to assign  $V$  blocks to the  $N$  frames in the GoP; that is,  $\sum_{i=1}^N v_i^n = k$ . Consequently, it can be shown that the number of ways to assign  $V$  blocks to  $N$  sets without replacement is

$$|\mathcal{A}| = \frac{(V + N - 1)!}{V!(N - 1)!}. \quad (28)$$

Since we are concerned with maximizing the visual quality of the final video, we define the reward function  $r^n$  as

$$r^n = -\log_{10} \left( l(\mathbf{X}^n, \hat{\mathbf{X}}^n) \right), \quad (29)$$

where  $l(\cdot, \cdot)$  is either  $l_{\text{PSNR}}$  or  $l_{\text{MS-SSIM}}$  depending on the video quality metric used. Note that in the previous section,

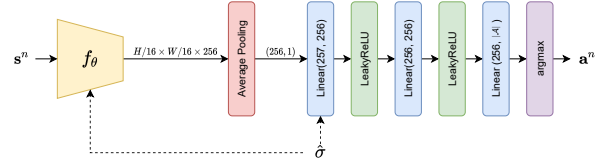


Fig. 8. Architecture of the bandwidth allocation network  $q_\psi$ . The convolutional part of the network for feature extraction is functionally the same as the key encoder network ( $f_\theta$ ) but with  $21(N - 1) + 6$  input dimensions to account for all the tensors in the state  $\mathbf{s}^n$ .

we said that the transmitter performs channel emulation in order to obtain the reference frames  $(\bar{\mathbf{x}}_{i-t}^n, \bar{\mathbf{x}}_{i+t}^n)$ . Since the precise estimate of the reference frames is bootstrapped to the amount of bandwidth allocated, we initially assume a uniform bandwidth allocation (i.e.,  $v_i^n = V/N$ ) when computing the SSF and residuals; but once the allocation has been done, the reference frame in the next state (i.e.,  $\bar{\mathbf{x}}_0^{n+1} \in \mathbb{M}^{n+1}$ ) is estimated based on the bandwidth allocated.

To solve the MDP described herein and to learn the optimal allocation policy, we use deep Q-learning [50], where the network  $q_\psi$  seeks to approximate the Q-function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . The purpose of the Q-function is to map each state and action pair to a Q value, which represents the total discounted reward from step  $n$  given the state and action pair  $(\mathbf{s}^n, \mathbf{a}^n)$ . That is,

$$Q(\mathbf{s}^n, \mathbf{a}^n) = E \left[ \sum_{i=n}^{\infty} \gamma^{i-n} r^i \mid \mathbf{s}^n, \mathbf{a}^n \right], \quad \forall (\mathbf{s}^n, \mathbf{a}^n) \in \mathcal{S} \times \mathcal{A}. \quad (30)$$

As is typical in DQN methods, we use *replay buffer*, *target network*, and  $\epsilon$ -*greedy* to aid the learning of the Q-function. The replay buffer  $\mathcal{R}$  stores experiences  $(\mathbf{s}^n, \mathbf{a}^n, r^n, \mathbf{s}^{n+1})$  and are sampled uniformly to update the parameters  $\psi$ . This prevents the states from being correlated, which would break the assumption in most optimization algorithms that the samples are independent. We use target parameters  $\psi^-$ , which are copies of  $\psi$ , to compute the DQN loss function:

$$L_{\text{DQN}}(\psi) = \left( r^n + \gamma \max_{\mathbf{a}} \{q_{\psi^-}(\mathbf{s}^{n+1}, \mathbf{a})\} - q_{\psi}(\mathbf{s}^n, \mathbf{a}^n) \right)^2. \quad (31)$$

The parameters  $\psi$  are then updated via gradient descent according to the gradient  $\nabla_{\psi} L_{\text{DQN}}(\psi)$ . In practice, the DQN loss is approximated with a batch of samples from the replay buffer  $\mathcal{B} \subset \mathcal{R}$ .

The target network parameters are updated via

$$\psi^- \leftarrow \tau \psi + (1 - \tau) \psi^-, \quad (32)$$

where  $0 \leq \tau \leq 1$ . The target networks here stabilize the updates. Due to Q-learning being bootstrapped, if the same  $q_\psi$  is used to estimate the state-action value of GoP number  $n$  and  $n + 1$ , both values would move at the same time, which may lead to the updates to never converge. By introducing the target networks, this effect is reduced due to the much slower updates of the target network, as done in Eqn. (32).

To promote exploration, we use  $\epsilon$ -greedy, which chooses a random action with probability  $\epsilon$  at each GoP. That is,

$$\mathbf{a}^n = \begin{cases} \arg \max_{\mathbf{a}} q_{\psi}(\mathbf{s}^n, \mathbf{a}), & \text{w.p. } 1 - \epsilon \\ \mathbf{a} \sim \text{Uniform}(\mathcal{A}), & \text{w.p. } \epsilon \end{cases}, \quad (33)$$

where  $\mathbf{a} \sim \text{Uniform}(\mathcal{A})$  denotes an action that is sampled uniformly from the action set  $\mathcal{A}$ . A diagram of the architecture used for  $q_\psi$  is shown in Fig. 8.

Upon initialization, we send the first frame  $\mathbf{x}_1^1$  using full bandwidth  $k$ . The first frame can be considered as a GoP on its own. For all subsequent GoPs, we perform optimal bandwidth allocation as described in this section. A diagram illustrating the bandwidth allocation problem formulated herein is shown in Fig. 7.

#### IV. NUMERICAL RESULTS

##### A. Training Details

We train our models on the UCF101 dataset [51] using Pytorch [52], with the Adam optimizer [53] at learning rate  $1e^{-5}$ . We split the dataset with 0.8 : 0.1 : 0.1 ratio between training, validation and evaluation. We train the JSCC ( $f_\theta, f_{\theta'}, g_\phi, g_{\phi'}, h_\eta$ ) networks first, randomizing the latent vector block sizes as described in Section III-B, before we train the bandwidth allocation network  $q_\psi$  to find the optimal bandwidth allocation policy. We also assume that during this phase, the transmitter and receiver can estimate the channel SNR accurately  $\hat{\sigma}^2 = \sigma^2$ . We use a training batch size of 4 when training the JSCC networks and a batch size of 8 when training the bandwidth allocation network. The batch sizes are relatively small due to memory restrictions of the available GPU (11GB Nvidia RTX 2080Ti). We use early stopping based on the validation error with a patience of 8 epochs. We adjust the learning rate based on the number of bad epochs: if the validation error does not improve for 4 epochs in a row, the learning rate is multiplied by 0.8.

We define the SNR estimated by the transmitter and receiver to be

$$\text{SNR}_{\text{Est}} = 10 \log_{10} \left( \frac{P}{\hat{\sigma}^2} \right). \quad (34)$$

For training the bandwidth allocation network, we choose DQN hyper-parameters  $\gamma = 0.99$ ,  $\tau = 0.005$ , and a replay buffer size  $|\mathcal{R}| = 1000$ . The function used for  $\epsilon$ -greedy exploration is

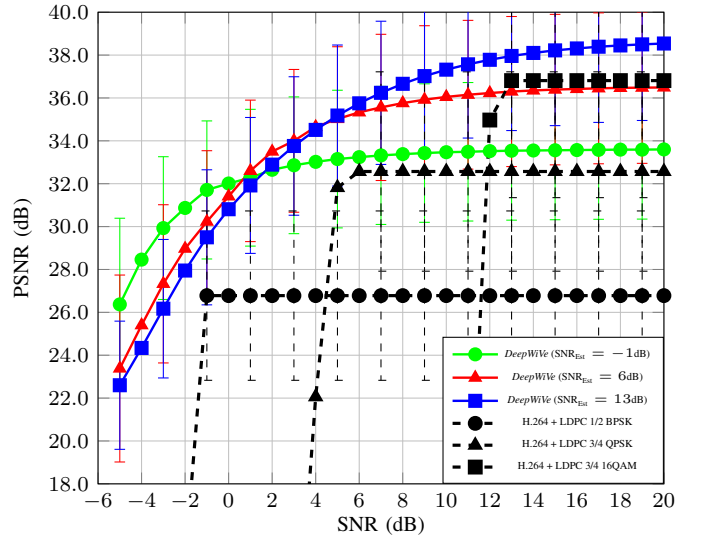
$$\epsilon = \epsilon_{\text{end}} + (\epsilon_0 - \epsilon_{\text{end}}) e^{\left( -\frac{\text{episode}}{\lambda} \right)}, \quad (35)$$

where  $\lambda$  controls the decay rate of  $\epsilon$ , and in each episode the bandwidth allocation network allocates the bandwidth resource within one video sequence. We choose  $\epsilon_0 = 0.9$ ,  $\epsilon_{\text{end}} = 0.05$ , and  $\lambda = 1000$  for these parameters.

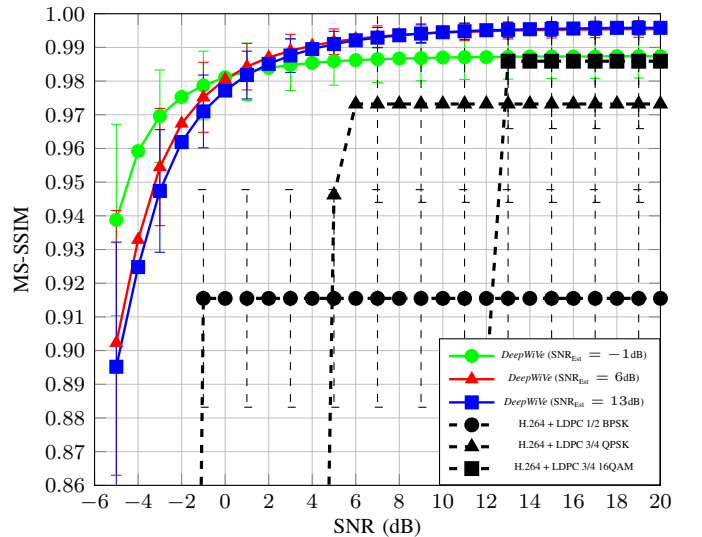
We train our model at different channel SNRs and evaluate each model at the same range of SNRs. In each batch, the training SNR is sampled uniformly from the range  $[-5, 20]$  dB. We chose  $N = 4$  and  $V = 20$  to train our models. Note that although we can choose  $V = k$  to achieve very fine grained bandwidth control, this would lead to a very large action space, which makes Q-learning difficult. We let  $P = 1$  and compute the necessary  $\sigma^2$  to achieve the desired SNR.

##### B. Simulation Results

We compare the performance of our model with that of the conventional separation-based schemes. In particular, we



(a) PSNR



(b) MS-SSIM

Fig. 9. Comparison of DeepWiVe using different  $\text{SNR}_{\text{Est}}$  to H.264 paired with LDPC codes ( $\rho = 0.031$ ).

use the H.264 [29] and H.265 [31] video compression codecs for source coding, LDPC codes [30] for channel coding, and QAM modulation. We plot the average video quality across the test dataset using each of the schemes considered herein and error bars representing the standard deviation of the video qualities.

In Fig. 9, we show the effect of channel estimation error on the performance of *DeepWiVe*. We specifically compare models using  $\text{SNR}_{\text{Est}}$  where the H.264 codec paired with a specific LDPC code rate experiences the *cliff-effect*. It is clear that *DeepWiVe* is able to overcome the *cliff-effect*, as video quality gracefully degrades as the SNR decreases even as  $\text{SNR}_{\text{Est}}$  remains the same. This is in contrast to the cliff edge drop off that separation-based designs suffer from. The *cliff-effect* is due to the fact that as the SNR decreases, the capacity of the channel decreases, leading to certain LDPC rate and modulation order pairs to communicate at a rate greater than

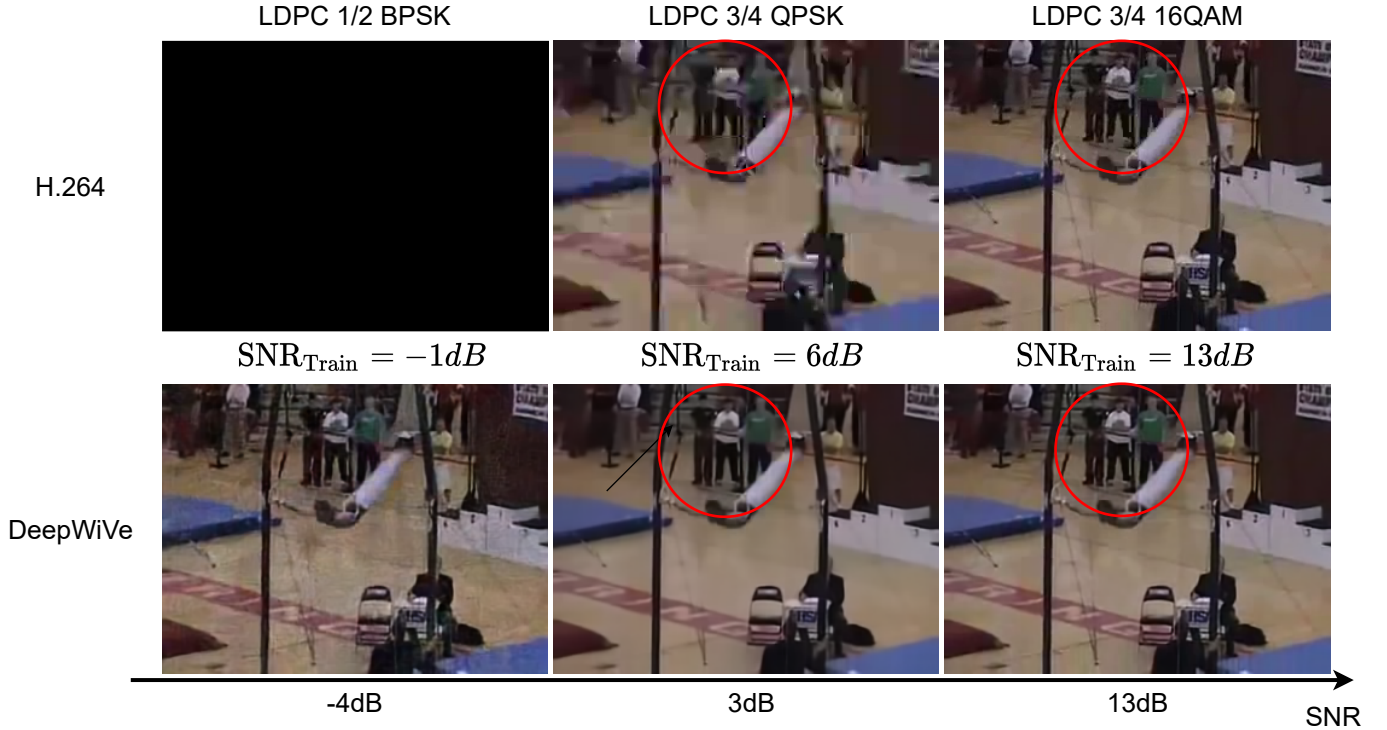


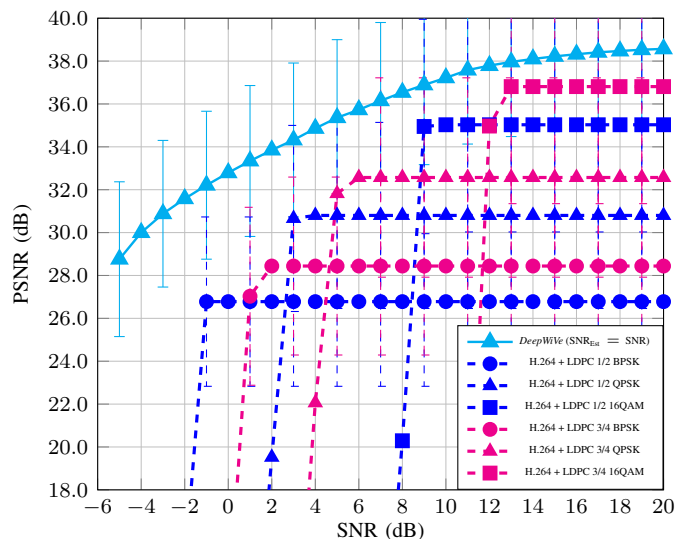
Fig. 10. Visual examples of *DeepWiVe* v.s. H.264. The difference in video quality can be seen most clearly in the 3 people standing at the back of the scene (encircled).

the capacity of the channel, which leads to highly unreliable communication, manifested as a cliff edge deterioration of the video quality. We can also see that the variations in the video quality using *DeepWiVe* is lower than those produced by the separation scheme as indicated by the smaller error bars. The error bars represent the standard deviation of the video quality at the receiver side. This is likely due to the fact that the H.264 codec does not have a continuous range of compression rates available but rather a set of discrete levels it can compress. Depending on the complexity of the video, a given target distortion may lead to a larger rate than is allowed by the instantaneous channel condition and it must reduce the target distortion level. Since the allowed target distortion levels can be far apart, this may mean the video is compressed more conservatively than is suggested by the channel in order to meet the channel condition, leading to a large variation in the resultant video quality. *DeepWiVe*, on the other hand, does not have this issue as we do not define a set of possible compression rates. Instead, the weights are adjusted by the AF modules based on the current channel condition to meet the rate-distortion curve as closely as possible. Fig. 10 shows the visual results of the plots shown in Fig. 9 for a specific video in the test dataset. It can be seen that at  $\text{SNR} = 13\text{dB}$ , the visual qualities of the videos produced by H.264 and *DeepWiVe* are similar. However, at  $\text{SNR} = 3\text{dB}$ , the video produced by H.264 starts to look very pixelated, while *DeepWiVe* is still able to retain a smooth looking frame. At  $\text{SNR} = -4\text{dB}$ , the capacity of the channel is too low for H.264 to compress the video sufficiently, therefore the output is simply black, while *DeepWiVe* is still able to achieve a reasonable video quality

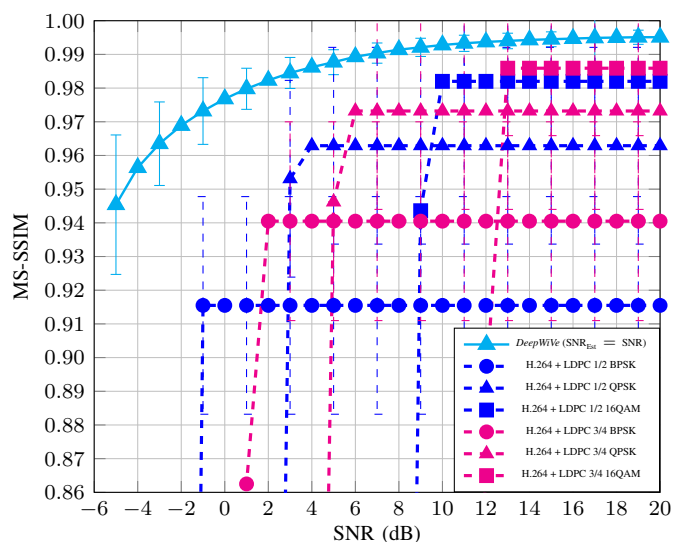
despite the very low channel SNR.

In Fig. 11, we present the comparison of *DeepWiVe* using the accurate estimate of the channel SNR (i.e.,  $\text{SNR}_{\text{Est}} = \text{SNR}$ ) with separation employing H.264, and in Fig. 12 employing H.265. In Fig. 11, we see that at  $\rho = 0.031$ , *DeepWiVe* is superior to the separation based scheme using H.264 in all the SNRs tested. This is due to the end-to-end optimization of the JSCC encoder and decoder producing a superior compression, channel code, and modulation scheme than the separation-based scheme. When compared to H.265, we see in Fig. 12a that H.265 outperforms *DeepWiVe* in terms of the PSNR metric. However, when compared with the more perceptually aligned MS-SSIM metric in Fig. 12b, we see that *DeepWiVe* can also outperform separation-based transmission with H.265. We also highlight that, in the very low SNR regime (i.e.,  $\text{SNR} < -1\text{dB}$ ), H.265 was unable to meet the compression rate required, and therefore did not produce results in that range. *DeepWiVe* on the other hand, did not have this problem. We believe that further optimization of the network architecture can bring *DeepWiVe* on par or surpass H.265 evaluated using the PSNR metric for higher SNR values as well. On average, *DeepWiVe* is 1.51dB better in PSNR and 0.0088 better in MS-SSIM than H.264 for  $\text{SNR} \in [13, 20]\text{dB}$ , 3.61dB better in PSNR and 0.0281 better in MS-SSIM for  $\text{SNR} \in [3, 6]\text{dB}$ . For H.265, *DeepWiVe* is 2.69dB worse in PSNR and 0.0056 better in MS-SSIM for  $\text{SNR} \in [13, 20]\text{dB}$ , 1.59dB worse in PSNR and 0.0109 better in MS-SSIM for  $\text{SNR} \in [3, 6]\text{dB}$ .

Next, we investigate the variable bandwidth transmission capability of *DeepWiVe* by decreasing the bandwidth compression ratio  $\rho$ , thereby increasing the compression of the



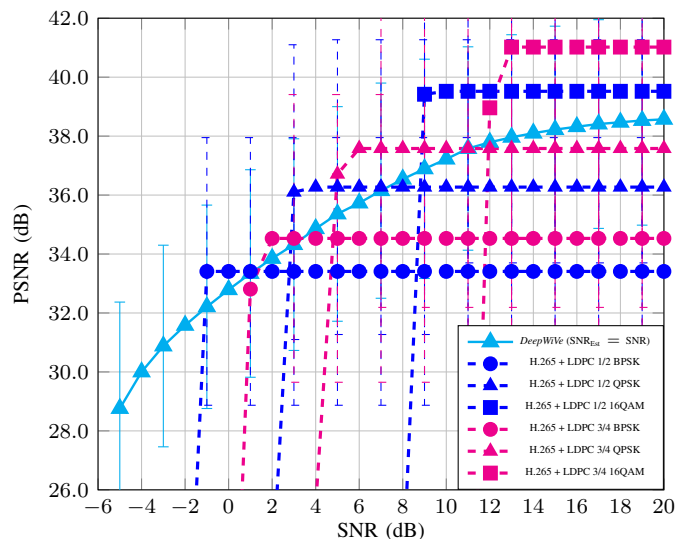
(a) PSNR



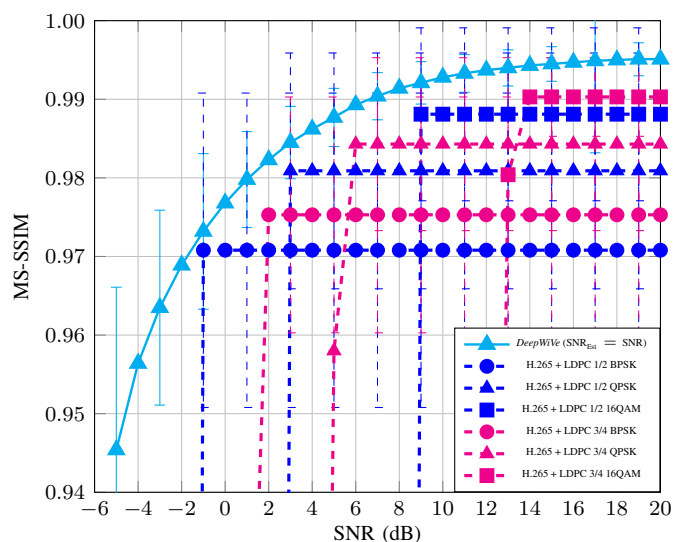
(b) MS-SSIM

Fig. 11. Performance comparison of DeepWiVe to H.264 paired with LDPC codes ( $\rho = 0.031$ ).

video. To decrease  $\rho$ , we do not require the retraining of the autoencoder networks ( $f_\theta, f_{\theta'}, g_\phi, g_{\phi'}, h_\eta$ ); we only need to retrain the bandwidth allocator  $q_\psi$ . As shown in Fig. 13, we see that *DeepWiVe* beats H.264 with LDPC coding for all SNRs tested in terms of both the PSNR and MS-SSIM metrics. It also beats H.265 using the MS-SSIM metric as shown in Fig. 14b, although again, it falls short of H.265 in terms of the PSNR metric for  $\text{SNR} > -1$  dB (Fig. 14a). This shows that *DeepWiVe* can achieve variable bandwidth transmission using RL to allocate an arbitrary number of blocks to meet the desired transmission bandwidth, as outlined in Section III-B. On average, for  $\rho = 0.018$ , *DeepWiVe* is 2.05 dB better in PSNR and 0.0129 better in MS-SSIM than H.264 for  $\text{SNR} \in [13, 20]$  dB, 4.72 dB better in PSNR and 0.0462 better in MS-SSIM for  $\text{SNR} \in [3, 6]$  dB. For H.265, *DeepWiVe* is 3.37 dB worse in PSNR and 0.0053 better in MS-SSIM for



(a) PSNR

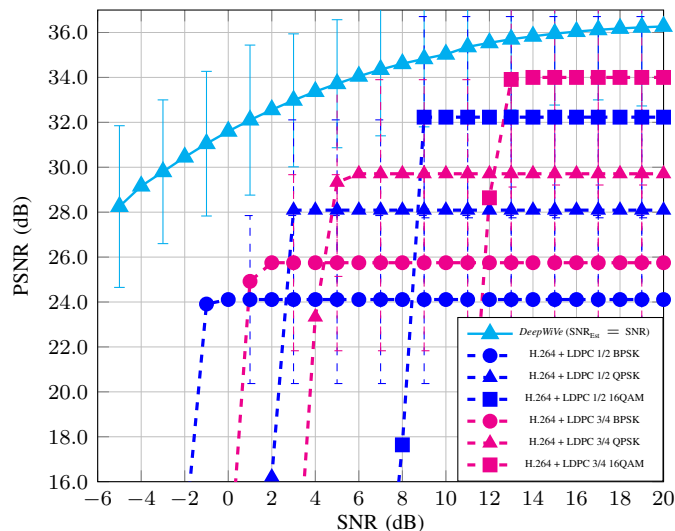


(b) MS-SSIM

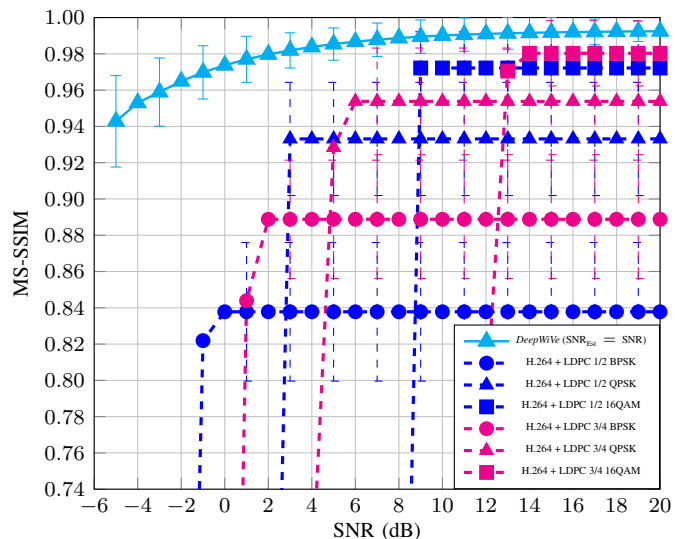
Fig. 12. Performance comparison of DeepWiVe to H.265 paired with LDPC codes ( $\rho = 0.031$ ).

$\text{SNR} \in [13, 20]$  dB, 3.16 dB worse in PSNR and 0.0058 better in MS-SSIM for  $\text{SNR} \in [3, 6]$  dB.

Lastly, we evaluate the performance of our models with and without the optimal allocation of bandwidth, as an ablation study. We compare the results obtained by using the allocation network  $q_\psi$  with that of uniform allocation (i.e.,  $v_i^n = 5, \forall i, n$  for  $\rho = 0.031$  and  $v_i^n = 3, \forall i, n$  for  $\rho = 0.018$ ). In Fig. 15, it can be seen that there is a clear and significant improvement in performance when the allocation is optimized by our allocation network  $q_\psi$ , by 0.67 dB and 0.0008 in PSNR and MS-SSIM, respectively, for  $\rho = 0.031$ . For  $\rho = 0.018$ , the gains are 0.99 dB and 0.0023 in PSNR and MS-SSIM, respectively. We observe that the gain from bandwidth allocation is more significant when  $\rho$  is smaller; that is, when the available channel bandwidth is more limited.



(a) PSNR

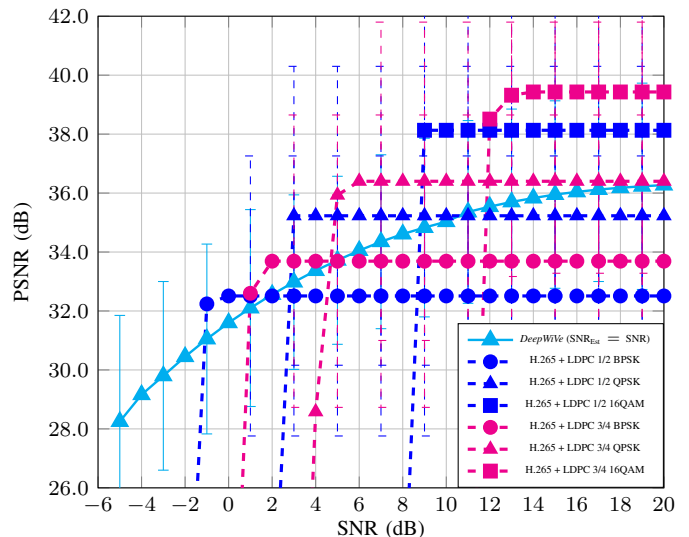


(b) MS-SSIM

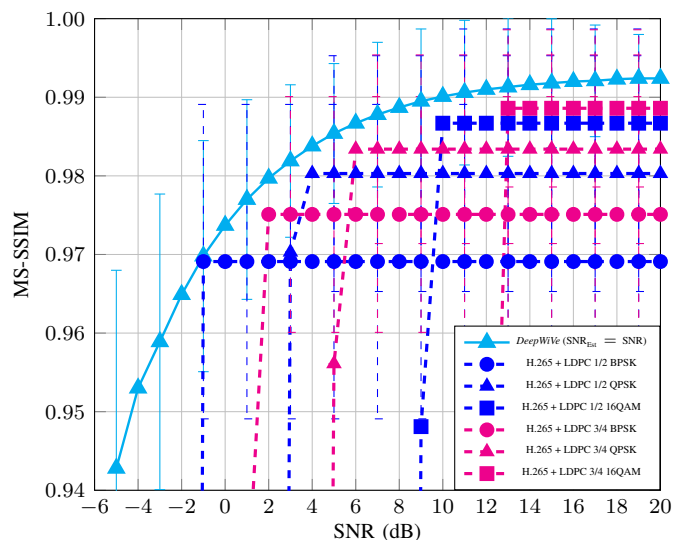
Fig. 13. Performance comparison of DeepWiVe to H.264 paired with LDPC codes ( $\rho = 0.018$ ).

## V. CONCLUSION

We presented the first ever DNN-aided joint source-channel wireless video transmission scheme in the literature. Our novel architecture, called *DeepWiVe*, is capable of dynamic bandwidth allocation and residual estimation without the need for distortion feedback. Additionally, it utilizes RL to learn a bandwidth allocation network that optimizes the allocation of available bandwidth within a given GoP in a dynamic fashion with the goal of maximizing the visual quality of the video under the given bandwidth constraint. Our results show that *DeepWiVe* overcomes the *cliff-effect* that all separation-based schemes suffer from, and achieves a graceful degradation with channel quality. In highly bandwidth constrained scenarios, *DeepWiVe* produces far superior video quality compared to both H.264 and H.265. We also show that our bandwidth allocation strategy is effective, improving upon the naïve uniform allocation by up to 0.87dB in PSNR. Our overall



(a) PSNR



(b) MS-SSIM

Fig. 14. Performance comparison of DeepWiVe to H.265 paired with LDPC codes ( $\rho = 0.018$ ).

results also show that *DeepWiVe* is better than the separation-based schemes using industry standard H.264 codec and LDPC channel codes in all the channel conditions considered by up to 4.72dB. Although, H.265 performs better than *DeepWiVe* in terms of the PSNR metric, *DeepWiVe* outperforms H.265 when compared in terms of MS-SSIM, which is widely accepted as a performance measure that better represents human perceptual quality.

## REFERENCES

- [1] "Cisco visual networking index: forecast and methodology 2016-2021.," 2017.
- [2] G. Cheung and A. Zakhour, "Bit allocation for joint source/channel coding of scalable video," *IEEE Transactions on Image Processing*, vol. 9, pp. 340–356, Mar. 2000.
- [3] I. Kozintsev and K. Ramchandran, "Robust image transmission over energy-constrained time-varying channels using multiresolution joint source-channel coding," *IEEE Transactions on Signal Processing*, vol. 46, pp. 1012–1026, Apr. 1998.

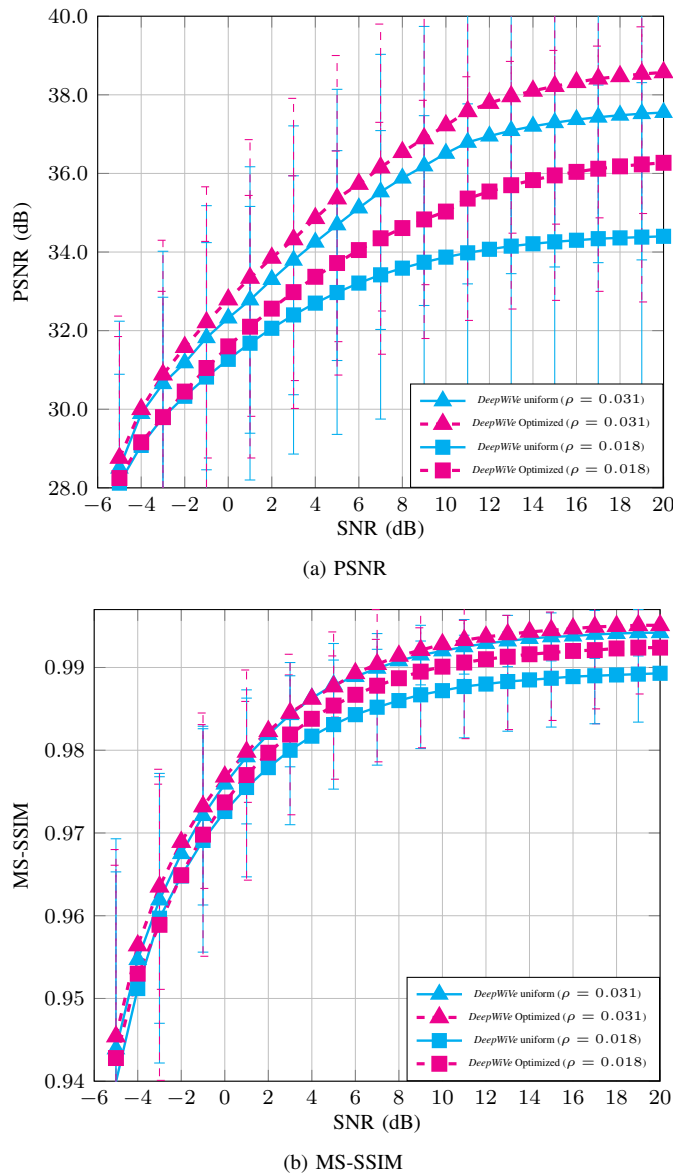


Fig. 15. Comparison of uniform bandwidth allocation versus optimal bandwidth allocation via auxiliary bandwidth allocation network  $q_\phi$ .

- [4] L. Kondi and A. Katsaggelos, "Joint source-channel coding for scalable video using models of rate-distortion functions," in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings*, vol. 3, pp. 1377–1380 vol.3, May 2001.
- [5] W. Ji, Z. Li, and Y. Chen, "Joint source-channel coding and optimization for layered video broadcasting to heterogeneous devices," *IEEE Transactions on Multimedia*, vol. 14, pp. 443–455, Apr. 2012.
- [6] G. Cheung and A. Zakhor, "Joint source/channel coding of scalable video over noisy channels," *AIP Conference Proceedings*, vol. 387, pp. 957–962, Jan. 1997.
- [7] L. Kondi, F. Ishtiaq, and A. Katsaggelos, "Joint source-channel coding for motion-compensated DCT-based SNR scalable video," *IEEE Transactions on Image Processing*, vol. 11, pp. 1043–1052, Sept. 2002.
- [8] T. P.-c. Chen and T. Chen, "Adaptive joint source-channel coding using rate shaping," in *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 2, pp. II–1985–II–1988, May 2002.
- [9] J. Wu, Y. Shang, J. Huang, X. Zhang, B. Cheng, and J. Chen, "Joint source-channel coding and optimization for mobile video streaming in heterogeneous wireless networks," *EURASIP Journal on Wireless Communications and Networking*, vol. 2013, p. 283, Dec. 2013.
- [10] T. Goblick, "A coding theorem for time-discrete analog data sources," *IEEE Transactions on Information Theory*, vol. 15, pp. 401–407, May 1969.
- [11] I. E. Aguerri and D. Gündüz, "Joint source-channel coding with time-varying channel and side-information," *IEEE Transactions on Information Theory*, vol. 62, pp. 736–753, Feb. 2016.
- [12] A. Lapidith and S. Tinguely, "Sending a bivariate Gaussian source over a Gaussian MAC with feedback," vol. 56, pp. 2246–2250, July 2007.
- [13] S. Jakubczak and D. Katabi, "SoftCast: One-size-fits-all wireless video," in *Proceedings of the ACM SIGCOMM 2010 Conference, SIGCOMM '10*, pp. 449–450, 2010.
- [14] T. Tung and D. Gündüz, "SparseCast: Hybrid digital-analog wireless image transmission exploiting frequency domain sparsity," *IEEE Communications Letters*, pp. 1–1, 2018.
- [15] W. Yin, X. Fan, Y. Shi, R. Xiong, and D. Zhao, "Compressive sensing based soft video broadcast using spatial and temporal sparsity," *Mobile Networks and Applications*, vol. 21, pp. 1002–1012, Dec. 2016.
- [16] A. Wang, B. Zeng, and H. Chen, "Wireless multicasting of video signals based on distributed compressed sensing," *Signal Processing: Image Communication*, vol. 29, pp. 599–606, May 2014.
- [17] E. Boursoulatze, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4774–4778, May 2019.
- [18] E. Boursoulatze, D. Burth Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, 2019.
- [19] D. B. Kurka and D. Gündüz, "Bandwidth-agile image transmission with deep joint source-channel coding," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2021.
- [20] C.-Y. Wu, N. Singhal, and P. Krahenbuhl, "Video Compression through Image Interpolation," in *2018 European Conference on Computer Vision (ECCV)*, pp. 416–431, Sept. 2018.
- [21] G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao, "DVC: an end-to-end deep video compression framework," in *2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11006–11015, Sept. 2019.
- [22] O. Rippel, S. Nair, C. Lew, S. Branson, A. G. Anderson, and L. Bourdev, "Learned Video Compression," pp. 3454–3463, Oct. 2019.
- [23] A. Djelouah, J. Campos, S. Schaub-Meyer, and C. Schroers, "Neural inter-frame compression for video coding," in *2019 International Conference on Computer Vision (ICCV)*, pp. 6420–6428, Oct. 2019.
- [24] A. Habibiyan, T. v. Rozendaal, J. M. Tomczak, and T. S. Cohen, "Video compression with rate-distortion autoencoders," in *2019 International Conference on Computer Vision (ICCV)*, pp. 7033–7042, Oct. 2019.
- [25] J. Han, S. Lombardo, C. Schroers, and S. Mandt, "Deep generative video compression," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, Nov. 2019.
- [26] E. Agustsson, D. Minnen, N. Johnston, J. Balle, S. J. Hwang, and G. Toderici, "Scale-space flow for end-to-end optimized video compression," in *2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8500–8509, June 2020.
- [27] B. Liu, Y. Chen, S. Liu, and H.-S. Kim, "Deep learning in latent space for video prediction and compression," in *2021 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 701–710, 2021.
- [28] Z. Hu, G. Lu, and D. Xu, "FVC: A new framework towards deep video compression in feature space," in *2021 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1502–1511, 2021.
- [29] T. Wiegand, G. Sullivan, G. Bjontegaard, and A. Luthra, "Overview of the H.264/AVC video coding standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, pp. 560–576, July 2003.
- [30] R. Gallager, "Low-density parity-check codes," *IRE Transactions on Information Theory*, vol. 8, pp. 21–28, Jan. 1962.
- [31] J.-R. Ohm and G. J. Sullivan, "High efficiency video coding: The next frontier in video compression [Standards in a Nutshell]," *IEEE Signal Processing Magazine*, vol. 30, pp. 152–158, Jan. 2013.
- [32] M. Stoufs, A. Munteanu, J. Cornelis, and P. Schelkens, "Scalable joint source-channel coding for the scalable extension of H.264/AVC," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 18, pp. 1657–1670, Dec. 2008.
- [33] H. Y. Shutoy, D. Gunduz, E. Erkip, and Y. Wang, "Cooperative source and channel coding for wireless multimedia communications," *IEEE Journal of Selected Topics in Signal Processing*, vol. 1, no. 2, pp. 295–307, 2007.
- [34] R. Xiong, F. Wu, X. Fan, C. Luo, S. Ma, and W. Gao, "Power-distortion optimization for wireless image/video SoftCast by transform coefficients energy modeling with adaptive chunk division," in *2013*

- Visual Communications and Image Processing (VCIP)*, pp. 1–6, Nov. 2013.
- [35] Z. Song, R. Xiong, S. Ma, and W. Gao, “Hybridcast: A wireless image/video SoftCast scheme using layered representation and hybrid digital-analog modulation,” in *2014 IEEE International Conference on Image Processing (ICIP)*, pp. 6001–6005, Oct. 2014.
- [36] R. Xiong, H. Liu, S. Ma, X. Fan, F. Wu, and W. Gao, “G-CAST: Gradient based image SoftCast for perception-friendly wireless visual communication,” in *2014 Data Compression Conference*, pp. 133–142, Mar. 2014.
- [37] H. Liu, R. Xiong, X. Fan, D. Zhao, Y. Zhang, and W. Gao, “CG-Cast: Scalable wireless image SoftCast using compressive gradient,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, pp. 1832–1843, June 2019.
- [38] R. Xiong, J. Zhang, F. Wu, and W. Gao, “High quality image reconstruction via non-local collaborative estimation for wireless image/video softcast,” in *2014 IEEE International Conference on Image Processing (ICIP)*, pp. 2542–2546, Oct. 2014.
- [39] A. Trioux, F.-X. Coudoux, P. Corlay, and M. Gharbi, “Performance assessment of the adaptive GoP-size extension of the wireless SoftCast video scheme,” in *2020 10th International Symposium on Signal, Image, Video and Communications (ISIVC)*, pp. 1–6, Apr. 2021.
- [40] J. Xu, B. Ai, W. Chen, A. Yang, and P. Sun, “Wireless image transmission using deep source channel coding with attention modules,” *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2021.
- [41] D. B. Kurka and D. Gündüz, “DeepJSCC-f: Deep joint source-channel coding of images with feedback,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, pp. 178–193, May 2020.
- [42] Z. Wang, E. P. Simoncelli, and A. C. Bovik, “Multiscale structural similarity for image quality assessment,” in *The Thirty-Seventh Asilomar Conference on Signals, Systems Computers, 2003*, vol. 2, pp. 1398–1402 Vol.2, Nov. 2003.
- [43] K. Sayood, H. H. Otu, and N. Demir, “Joint source/channel coding for variable length codes,” *IEEE Transactions on Communications*, vol. 48, pp. 787–794, May 2000.
- [44] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1874–1883, June 2016. ISSN: 1063-6919.
- [45] J. Ballé, V. Laparra, and E. P. Simoncelli, “Density Modeling of Images using a Generalized Normalization Transformation,” *arXiv:1511.06281 [cs]*, Feb. 2016. arXiv: 1511.06281.
- [46] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, “Learned image compression with discretized Gaussian mixture likelihoods and attention modules,” in *2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [47] X. Wang, R. Girshick, A. Gupta, and K. He, “Non-local neural networks,” in *2018 Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7794–7803, June 2018.
- [48] S. S. Beauchemin and J. L. Barron, “The computation of optical flow,” *ACM Computing Surveys*, vol. 27, pp. 433–466, Sept. 1995.
- [49] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3431–3440, June 2015.
- [50] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, pp. 529–533, Feb. 2015.
- [51] K. Soomro, A. R. Zamir, and M. Shah, “UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild,” *arXiv:1212.0402 [cs]*, Dec. 2012. arXiv: 1212.0402.
- [52] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in PyTorch,” Oct. 2017.
- [53] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” *arXiv:1412.6980 [cs]*, Jan. 2017. arXiv: 1412.6980.