

Foveation-based Deep Video Compression without Motion Search

Meixu Chen, *Student Member, IEEE*, Richard Webb, and Alan C. Bovik, *Fellow, IEEE*

Abstract—Virtual Reality (VR) and its applications have attracted significant and increasing attention. However, the requirements of much larger file sizes, different storage formats, and immersive viewing conditions pose significant challenges to the goals of acquiring, transmitting, compressing, and displaying high-quality VR content. At the same time, the great potential of deep learning to advance progress on the video compression problem has driven a significant research effort. Because of the high bandwidth requirements of VR, there has also been significant interest in the use of space-variant, foveated compression protocols. We have integrated these techniques to create an end-to-end deep learning video compression framework. A feature of our new compression model is that it dispenses with the need for expensive search-based motion prediction computations. This is accomplished by exploiting statistical regularities inherent in video motion expressed by displaced frame differences. Foveation protocols are desirable since, unlike traditional flat-panel displays, only a small portion of a video viewed in VR may be visible as a user gazes in any given direction. Moreover, even within a current field of view (FOV), the resolution of retinal neurons rapidly decreases with distance (eccentricity) from the projected point of gaze. In our learning based approach, we implement foveation by introducing a Foveation Generator Unit (FGU) that generates foveation masks which direct the allocation of bits, significantly increasing compression efficiency while making it possible to retain an impression of little to no additional visual loss given an appropriate viewing geometry. Our experiment results reveal that our new compression model, which we call the Foveated MOtionless VIdeo Codec (Foveated MOVI-Codec), is able to efficiently compress videos without computing motion, while outperforming foveated version of both H.264 and H.265 on the widely used UVG dataset and on the HEVC Standard Class B Test Sequences. The Foveated MOVI-Codec project page can be found at <https://github.com/Meixu-Chen/Foveated-MOVI-Codec>.

Index Terms—Foveation, motion, video compression.

I. INTRODUCTION

RECENT advances in Virtual Reality (VR) offered more immersive viewing experiences than even high-resolution flat-panel displays. However, VR presentations require much larger file sizes, high bandwidths, different storage formats, and spatial-purpose immersive viewing hardware, all of which pose significant challenges against the goals of acquiring, transmitting, compressing, and displaying high quality VR content. One advantage of VR, however, is that the two eyes have fixed positions, aside from eye movements, relative to the viewing screen. Because of this, the eye movements, and associated points of gaze on the displays can be measured. This makes it possible to exploit the fact that the density of retinal photosensors is highly non-uniform. The cone cells

used in photopic viewing achieve on peak density in the foveal region, which captures a circumscribed FOV of about 2.5° around gaze. This includes only 0.8% of all pixels on a flat panel display when viewed under typical conditions [1], and around 4% of pixels on a VR display [2], [3]. Since the density of photoreceptors falls away quite rapidly with increased eccentricity relative the fovea, much more efficient representations of what is perceived can be obtained by judiciously removing redundant information from peripheral regions.

While foveal processing protocols might be useful for many aspects of VR rendering and viewing, such as enhancement or brightening around the point of gaze, foveated compression may offer the most significant and obvious benefits. While this topic has been studied in the past [4]–[6], only recently has there been renewed interest in foveating modern codecs [7]. In this direction, because of their tremendous ability to learn efficient visual representations, deep learning models are viewed as promising vehicles for developing alternative video codecs. This also raises the possibility of creating very efficient, end-to-end deep foveated video compression engines, which is the main topic we study here. Much less work has been done on deep video compression than on learning-based image compression. However, a variety of uniform resolution (without foveation) deep video compression models have been devised [8]–[11]. For example, Wu *et al.* [8] proposed a deep video compression network utilizing the idea that video compression may be conceptualized as image interpolation. Chen *et al.* [9] enhanced spatial-temporal energy compaction in a learning-based video compression model by introducing a spatial energy compaction penalty into the loss function. The authors of [10] used bidirectional reference frames to compute motion maps in a hierarchical learning scheme. Lu *et al.* [11] developed a video compression network called DVC, where each component of a traditional hybrid codec is replaced by a deep network. In [12], a deep video compression without motion estimation and compensation is proposed, and the resulting codec is competitive to latest H.266 on high resolution videos against MS-SSIM.

In both traditional video codecs and recent deep learning-based models, motion estimation and compensation occupy a significant portion of compression system resources. Motion estimation requires expensive search processes that are amplified in very large format videos (e.g. 6K and above) that are needed to display naturalistic video in VR. In the model we proposed here, we avoid motion estimation via search, by instead training a specially-designed network that is able to efficiently represent the residuals between each frame and a set of spatially-displaced neighboring frames. Computing a set

M. Chen and A. C. Bovik are with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, USA (e-mail: chenmx@utexas.edu; bovik@ece.utexas.edu).

Richard webb is with Meta Reality Labs(email: rwebb@fb.com).

of frame differences, even over many displacement directions, is much less expensive than conducting matching-based search processes. Moreover, while the statistics of motion vectors are generally not regular, the intrinsic statistics of frame differences exhibit strong regularities [13], especially when the differences are taken between spatially displaced lying along the same local motion paths [14]. The strong internal structure of aligned, high-correlation frame differences makes them sparser and easier to efficiently represent in a deep architecture.

Our idea is inspired by the way the human visual system processes natural time-varying images. Many studies have produced strong evidence suggesting that the early stages of vision are primarily implicated in reducing redundancies in the sensed input visual signals [15], [16]. Indeed, much of early visual processing along the retino-cortical pathway appears to be devoted to processes of spatial and temporal decorrelation [17]–[22]. We have found that sets of spatially displaced frame differences, which are space-time processes, supply a rich and general way to exploit space-time redundancies [13], [14]. Our idea is also related to recent theories of the role of microsaccades in human visual information processing [18]–[21]. Microsaccades create small spatial displacements of the visual field from moment to moment. While microsaccades have been theorized to play roles in avoiding retinal saturation, maintaining accurate fixation in the presence of drifts, and preserving the perception of fine spatial details [21], it has been more recently theorized to play an important role in efficiently representing locally changing and shifting space-time visual information [18]–[21]. We believe that micro-saccadic eye movements may be an efficient adaptation to efficiently exploit local regularities induced by small spatial displacements over time, to achieve more efficient visual (neural) representations. This has inspired us to, in like manner, train a deep foveated coder-decoder network that compresses videos using regular displaced residual representations as inputs.

Because the compute bandwidths and data volumes involved in VR rendering are also unusually high, we have sought to reduce both data and computation in two ways: by perceptually relevant foveation, and by a perception-driven elimination of expensive motion computations.

The success of foveation based processing protocols involves several factors, including distribution of retinal ganglion cells [23], cortical magnification [24], and the steep grade of density of the photoreceptors [25]. The spacings of the photoreceptors and the receptive fields of the neurons they feed are smallest in the fovea [26]. The fovea covers an area in the approximate range of 0.8% to 4% of the pixels on a display, depending on the display size, resolution, and the assumed typical viewing distance [1], [2]. Recent advances in eye-tracking technology and their integration into consumer VR headsets have opened the possibility of using them to facilitate gaze-contingent video compression. Indeed, retinal foveation when combined with ballistic saccadic eye movements to direct visual resources, is a form of biological information compression. For example, the density of retinal ganglion cells (RGC) in the fovea is $325,000/mm^2$. If the entire retina had this output density, then about 350 million RGCs would

be implied. However, the number of axons carrying signals along the optic nerves of each eye is only around 1 million, hence foveation results in a 350-fold compression of data passed along the retino-cortical pathway [27]. In an analogous manner, considerable increases in digital video compression can be obtained by removing visual redundancies (relative to fixation) in the visual periphery.

Here we introduce a foveated deep video compression network that is efficient, statistically and perceptually motivated, and free to motion search computations. Our specific contributions may be summarized as follows:

- We incorporate foveation into a deep video compression model to achieve significant data reductions suitable for eye-tracked VR systems.
- We innovate the use of displaced frame differences to capture efficient representations of structures and temporal statistical redundancies induced by motion.
- The overall video compression system, which we call the Foveated MOVI-Codec is optimized using a single loss function.
- The new model is shown to obtain state-of-the-art compression performance on the widely used UVG dataset and on the HEVC Standard Test Sequence Class B dataset, outperforming H.264, H.265 and their foveated counterparts against the well-known perceptual video quality index Foveated Wavelet Image Quality Index (FWQI) [5], [28].

The rest of the paper is organized as follows. Section III details the architecture and training protocol used to create the Foveated MOVI-Codec model. Section IV explains the experiments on algorithm performance and comparisons that we conducted. Section V concludes the paper with a discussion of future research directions.

II. RELATED WORKS

A. Foveated Video Compression

Since the turn of the millennium, there has been a slowly growing interest in the use of foveation for such diverse image and video processing tasks as quality assessment [29], segmentation [30], and watermarking [31]. Methods of foveating visual content can be categorized into three ways: geometric transformations, space-varying filters, and space-variant multiresolution decompositions [32]. In the first of these, a foveated retinal sampling geometry is used to either apply a foveating coordinate transformation on an original uniform resolution image [33], or to average and map local pixel groups into superpixels [34], [35]. Filter-based methods process images with space-varying low-pass filter with cut-off frequencies determined by foveated resolution-reduction protocols [36], [37]. Multiresolution methods foveation involves decomposing images into bandpass scales, and only retaining scales specified by a foveal fall-off function defined relative to a measured or presumed fixation point [4], [38].

Recently, given significant advances in high resolution and immersive displays technologies, along with concurrent increases in VR content, interest of foveation as an efficient processing tool has quickened. Recent related models include

[39], where a neurobiological model of visual attention is used to predict high saliency regions and to generate saliency maps. A guidance map is also generated, using foveation to guide bit allocations when tuning quantization parameters in video compression system. Li *et al.* [40] trained a content-weighted CNN to conduct image compression, whereby the bitrates allocated to different parts of an image are adapted to the local content. Their system significantly outperforms JPEG and JPEG2000 in terms of SSIM when operating in a low bitrate regime. Mentzer *et al.* [41] proposed a similar but simpler model, by incorporating a second channel at the output of the encoder that is expanded into a mask which is used to modify the latent representations. DeepFovea [7] is a foveated reconstruction model, that employs a generative adversarial neural network. A peripheral video is reconstructed from a small fraction of pixels, by finding a closest matching video to the sparse input stream of pixels that lies on the learned manifold of natural videos. This method is fast enough to drive gaze-contingent head-mounted displays in real time.

B. Foveated Video Quality Assessment

When designing foveated compression systems, it is desirable to be able to access their perceptual efficiencies using quality measurement tools that account for the foveation. However, almost all available image quality measurement tools, such as SSIM [42], operate on spatially uniform resolution contents. However, there are a few foveated video quality assessment models, which can be conveniently divided into several types. One type of foveated VQA model uses purely static, spatial foveation, whereby measurement or prediction of the user's point of gaze guides the space variant measurement of quality as a function of eccentricity. For example, the Foveated Wavelet Image Quality Index (FWQI) utilizes wavelets to extract position-dependent spatial quality information [5], [28]. Several factors are taken into consideration, including the spatial contrast sensitivity function, which is used to determine local visual cutoff frequencies, which guides modeling of human visual sensitivity across the available wavelet subbands, when combined with assumption on viewing distance and the display resolution. Lee *et al.* [29] proposed a foveal signal-to-noise ratio (FSNR) to evaluate the quality of picture or video streams. In this method, a foveated image is obtained by a foveated coordinate transformation on the original image(s) to be quality-accessed.

A second type of foveated VQA model is based on retinal velocity. In addition to static foveation mechanisms, these kinds of models also take advantage of the fact that the contrast sensitivity of HVS to an object in a moving scene is influenced by the velocity of its map on the retina. Movement in a video may cause two effects: loss of acuity of the moving objects, modifications of perceived quality. Further, two factors can contribute to losses of acuity: increases of retinal image velocity, and increases of eccentricity relative to the foveal center. Based on these observations, Riomas-Drlje *et al.* [43] proposed a foveated mean squared error (FMSE) that models the effects of spatial acuity reduction due to motion. Another model called the foveation-based content Adaptive Structural SIMilarity index (FA-SSIM), which is based on the popular

IQA model SSIM [44] combines SSIM with a foveation-based sensitivity function.

You *et al.* [45] proposed a full reference attention-driven foveated video quality metric (AFViQ) that accounts for the localization of fixations in images and videos. All of the algorithms mentioned above assume that the point of fixation is the center of the image, which is not always true, and can lead to an invalid foveation model. As a result, algorithms based on automatic fixation detection have also been proposed. AFViQ attempted to solve this problem by integrating foveation into a wavelet-based distortion visibility model.

III. PROPOSED METHOD

A. Framework

Figure 1 illustrates the overall architecture of our network, which extends our previous MOVI-Codec [12]. The compression network is comprised of four components: a Displacement Calculation Unit (DCU), a Displacement Compression Network (DCN), a Foveation Generator Unit (FGU), and a Frame Reconstruction Network (FRN). The DCU computes displaced frame differences between the current frame and the previous reconstructed frame; the FGU generates foveation masks that later direct the allocation of bits in DCN; the DCN compresses displaced frame differences generated from DCU; and the FRN reconstructs input frames from the previous reconstructed frame and the reconstructed displaced frame differences.

The flow of our network is: Given an input video with frames x_1, x_2, \dots, x_T , for every frame x_t , calculated displaced frame differences between the current frame x_t and previous reconstructed frame \hat{x}_{t-1} via the DCU, after which the displaced frame differences d_t are input into the DCN. In the FGU, a perception-based foveation map P is generated from [4], [28] and used to generate a set of foveation masks $M(P)$. After the set of displaced frame differences d_t are encoded into latent representations y_t , the masks generated from the FGU direct the allocation of bits via element-wise multiplication of y_t and $M(P)$, producing a masked latent representation c_t , which is then quantized (via rounding) and decoded to \hat{d}_t . Finally, the FRN reconstructs the input frame \hat{x}_t from the reconstructed displaced frame differences \hat{d}_t and the previous reconstructed frames \hat{x}_{t-1} . The DCN and FRN are defined identically as in [46], so we do not further elaborate them here. We explain the DCU and FGU in the following.

B. Displacement Calculation Unit (DCU)

The DCU removes the need for any kind of motion vector search. Instead, the DCU learns to optimally represent time-varying images as sets of spatially displaced frame differences. Given a video with T frames x_1, x_2, \dots, x_T of width W and height H , two directional (spatially displaced) temporal differences are computed between each pair of adjacent frames, as shown in Figure 2. Assume that the inputs to the DCU a current frame x_t and a reconstructed previous frame \hat{x}_{t-1} . Then, at each spatial coordinate (i, j) , a set of spatially displaced differences is calculated as:

$$d_H(i, j)_t = x_t(i, j) - \hat{x}_{t-1}(i, j - s), \quad (1)$$

$$d_V(i, j)_t = x_t(i, j) - \hat{x}_{t-1}(i - s, j). \quad (2)$$

all periperal information. Figure 3 shows a quantized foveation map. Since the latent representations for the displaced frame differences d_t are 128 channels, the same mask is assigned for every 8 channels of latent representations. The quantized map in x axis is shown in Figure 4. After a set of n masks $M(P)$ are generated, we element-wise multiply $M(P)$ and the encoder output y_t to obtain quantized spatially variant (foveated) codes c_t which are then subjected to entropy coding and bitrate estimation, using the same procedure as [12].

While quantized foveation maps are used to train our model, during application (testing) we instead use isotropic 2D gaussian shaped foveation maps, where the gaussians are defined to follow the modified fall-off of visual acuity. This allows for smoother perceived changes of foveation, with the significant added benefit of making it possible to effect variable rate control by varying the widths (σ) of the gaussians. We define this parameter as foveation mask space constant (FMSC).

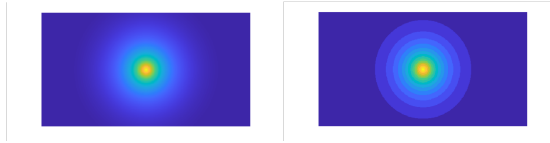


Fig. 3. Foveation map (left) and quantized foveation map (right), where brighter regions corresponds to larger value.

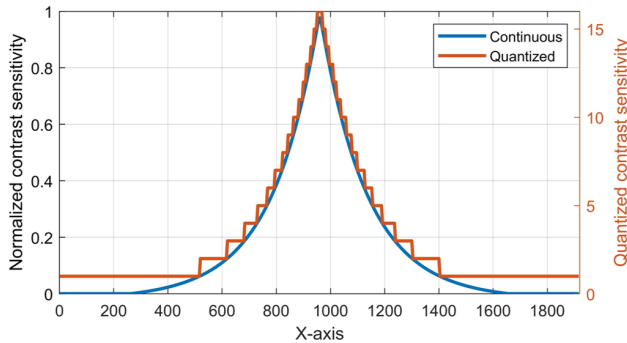


Fig. 4. Quantized contrast sensitivity function.

D. Bit Rate Allocation

Given an input frame x , let $y = E(x) \in R^{c \times h \times w}$ be the output of the encoder network, which includes c feature maps of sizes of $h \times w$. Also let $p = P(x)$ denote a $h \times w$ non-negative foveation map to be applied. The expand y using masks $\mathbf{m} \in R^{c \times h \times w}$ as follows:

$$m(i, j, k) = \begin{cases} 1 & \text{if } p(i, j) \geq \lfloor \frac{k}{c/L} \rfloor \cdot \frac{1}{L} \\ 0 & \text{others} \end{cases}, \quad (6)$$

where c is the number of channels in the latent representations y , and L is the number of desired compression levels across the foveation regions. In this way, more bits are allocated to the foveal region, preserving visual details with less sacrifice of bit rate. The sum of the foveation maps $\sum_{i,j} p_{i,j}$ naturally serves as a continuous estimate of compression rate, and can be directly adopted as a compression rate controller. Because

of the flexibility of this foveation map approach, it is not necessary to apply entropy rate estimation when training the encoder and decoder, using a simple binarizer for quantization of latent representations y .

E. Training Strategy

We are able to model the loss function considering only the distortion as follows:

$$D = [D_1(x_t, \hat{x}_t) + \beta D_2(d_t, \hat{d}_t)], \quad (7)$$

where D represents the distortion, and D_1 is the distortion between the input frame x_t and reconstructed frame \hat{x}_t , measured by foveation-weighted SSIM as detailed below at the end of this subsection. D_2 is the distortion between the displaced frame differences d_t and the reconstructed displaced frame differences \hat{d}_t , as measured by the MSE. The weight β controls the trade-off between the perceptual distortion D_1 and the pixel-to-pixel distortion D_2 .

To leverage multi-frame information in our RNN-based codec structure, we update the network parameters every set of N frames during model training, using the loss function in Equation 7, but modified to be a sum of losses over the k th set of N frames indexed $x_{t_k+1}, \dots, x_{t_k+N}$:

$$D_k = \frac{1}{N} \sum_{n=1}^N [D_1(x_{t_k+n}, \hat{x}_{t_k+n}) + \beta D_2(d_{t_k+n}, \hat{d}_{t_k+n})]. \quad (8)$$

During training, we selected a random $W \times W$ patch from each training video, and also randomly sampled a patch of the same size from the foveation map, to generate foveation masks from the patch. Foveation-weighted SSIM scores were calculated by applying a low-pass filter (Haar's filter) on the SSIM scores of each frame patch, then multiplying them by the foveation map patches. The overall workflow is shown in Figure 5.

IV. EXPERIMENTS

A. Settings

The Foveated MOVI-Codec networks that we experimented with were trained end-to-end on the Kinetics-600 dataset [47], [48] and on the Vimeo-90K dataset [49]. We used part of the testing set from Kinetics-600, which consists of around 10,000 videos, to conduct our experiments. From each video, a random 192×192 patch containing 49 frames was randomly selected for training, and normalized the values of each input video to $[-1, 1]$. Since Kinetics-600 dataset consist of YouTube videos of different resolutions, we randomly downsampled each original frames and extracted a 192×192 patches from the foveation maps to reduce any previously introduced compression artifacts. We randomly sampled 192×192 patches from the foveation maps to generate foveation masks for bitrate allocation. The Vimeo-90K dataset consists of 4,278 videos of fixed resolution 448×256 . Since the videos in this dataset each have 7 frames, we randomly selected patches from each of the same size as mentioned earlier (overall 7 frames) for training.

We fixed the mini-batch size to 8 for training, while the step length N of the recurrent network was set as 7. We used

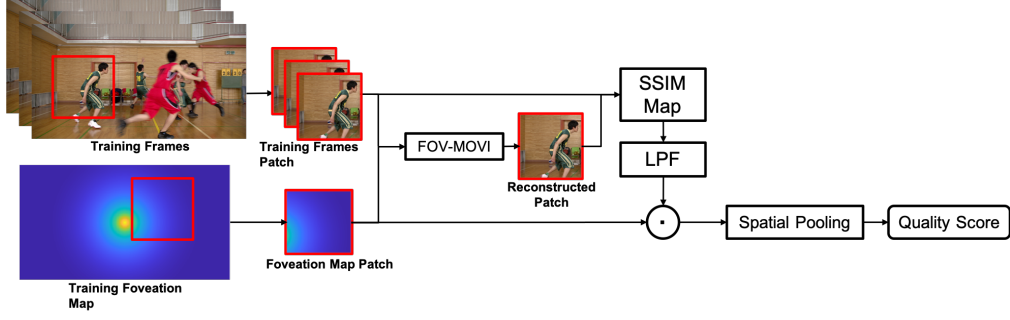


Fig. 5. Training strategy.

Adamax optimizer for training and set the initial learning rate to 0.0001. The whole system is implemented based on PyTorch and using one Titan RTX GPU. By training on both the Vimeo-90K and the Kinetics-600 datasets, we are able to generalize our model to a wider range of natural motions. We tested the Foveated MOVI-Codec on the JCT-VC Class B datasets [50] and the UVG datasets [51]. Both of these testing datasets have HD resolution contents (1920×1080).

In order to assess the reconstruction quality of the foveation compressed videos, we utilized the perceptually relevant FWQI foveated video quality measurement tool following the same settings in [7], with screen width being 0.02 meters and display distance being 0.012 meters. We also used the foveated SSIM model which deploys a fixed foveation map generated from the error sensitivity function from [28]. During testing, videos having different bitrates were generated using gaussian shape foveation maps with different foveation mask space constants, e.g. FMSCs of $\frac{H}{10}, \frac{H}{8}, \frac{H}{6}, \frac{H}{4}, \frac{H}{3}$, and $\frac{H}{2}$, where H is the height of the input frame. Exemplar 1D slices through the gaussians are shown in Figure 6, while the corresponding quantized maps are shown in Figure 7. We fixed $\beta = 1$.

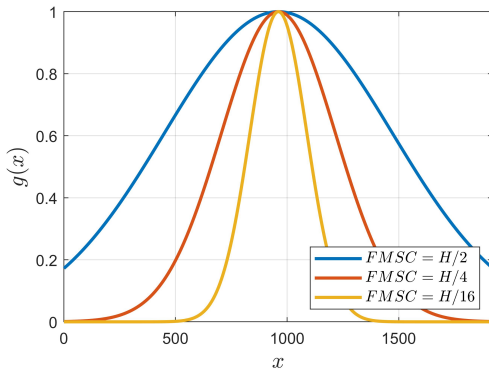


Fig. 6. Normalized sliced profiles of gaussian foveation masks.

B. Results

1) Rate-Distortion Curve

We compared our video compression engine against the standardized hybrid codecs H.264 and H.265, and also against our previous non-foveated model, the MOVI-Codec, on the UVG dataset and the HEVC Standard Test Sequences Class

B. In addition, we also implemented a foveated version of the hybrid codecs using the foveation method mentioned [4]. Both testing datasets have resolutions 1920×1080 .

Figure 8 shows the results obtained on the UVG and HEVC Class B datasets. Unsurprisingly, the foveated version of the hybrid codecs outperforms their foveated counterparts in terms of FWQI. These results also show that our foveated model outperformed the non-foveated MOVI-Codec on both datasets. Moreover, the Foveated MOVI-Codec outperformed both H.264 and H.265, as well as their foveated counterparts, on both datasets. It is worth noting that the measured qualities of the reconstructed videos produced by Foveated MOVI-Codec produced did not vary much with respect to bitrate, suggested that our model is able to maintain a high quality fovea, while decreasing the bitrate derived from the periphery without sacrificing perceptual video quality. Visualizations of example frames compressed using different levels of bitrates and qualities are shown in Figure 9. More exemplar reconstructed videos are included on our project page with link given in the Abstract. In these reconstructed frames, we selected three regions for detailed comparison: one in the foveal region and the other two others in peripheral. Our model is able to reconstruct videos having higher quality foveas and peripheral regions than the compared models, both visual and in terms of FWQI.

2) Latent Representations

As mentioned in Section III-C, the Foveated MOVI-Codec uses foveation maps to mediate bit allocations as a function of eccentricity relative to visual fixation. To visualize this process, we compared the latent representations (the encoded outputs) y_t in the Foveated MOVI-Codec against the encoded outputs y'_t of the original MOVI-Codec as shown in Figure 10. In the figure, the first row corresponds to reconstructed frames under different models, where the first column shows reconstructed frames from the MOVI-Codec, the second column contains reconstruction from the Foveated MOVI-Codec trained with a uniform (non-foveated) importance map with the masks of first N channels being one and zero elsewhere, and N is a random number during training. The remaining two columns show reconstructions from the Foveated MOVI-Codec with foveation mask space constants FMSCs equal to $H/2$ and $H/4$, respectively. The second row shows the corresponding accumulated feature maps, where brighter colors

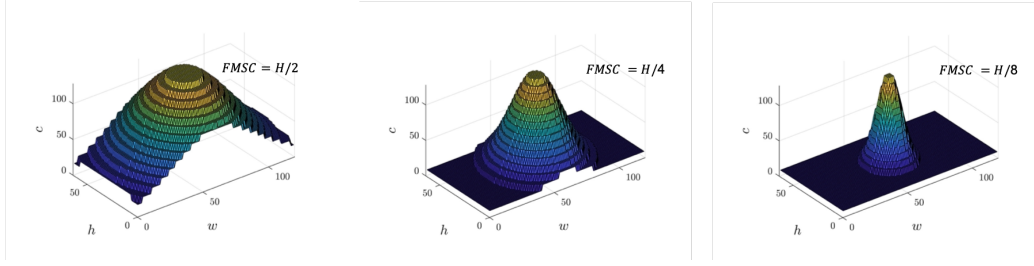


Fig. 7. Exemplar quantized gaussian foveation masks with different foveation mask space constants $FMSC$ s.

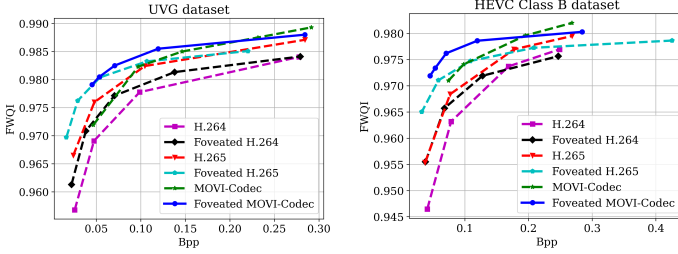


Fig. 8. FWQI of the compared models on the UVG dataset and HEVC B test sequences. All video resolutions are 1920×1080 .

correspond to larger numbers of more features. This shows that more features are used to represent the foveal region, as the foveation maps become narrower (smaller $FMSC$). The last row of Figure 10 shows the latent representations at each level (8 channels per level) of the reconstructed frames. Figure 11 also shows the sum of latent representations of y_t and y'_t (foveated and non-foveated, respectively). As shown in Figure 11, the sum is roughly flat for the non-foveated compressor, whereas the sum decreases as with the channel number for the foveated compressor. This suggests that the foveated network was able to learn more relevant features in the first few channels. From Figures 10-11, we may conclude that the model learned efficient features across channels and bit allocation, even without the masked multiplication.

3) Bit Allocation

Figure 12 shows the reconstructed frames from the Foveated MOVI-Codec, differenced frames between original frames and reconstructed frames, and bits and SSIM profiles, when using different foveation space constants. From the differenced frames, we can conclude that our model is able to reconstruct a foveal region similar to the original frame regardless of the mask used. The third row of Figure 12 shows the both a bit allocation plot and the SSIM map profile for different models. From the SSIM map profile, it may be observed that the lower number of bits allocated to the peripheral does not result in lower quality, since the quality of the reconstructed frames remain similar overall.

C. Discussion

Our experiments have shown that deploying foveation masks leads to much more efficient video compression for suitable environments, such as VR. Our new model outperformed H.264, H.265 and their foveated counterparts against FWQI across all testing sequences. Our model is best targeted at high

resolution, gaze contingent foveated compression applications in VR and AR. The hierarchical masks make it possible to transmit scalably, viz., the first levels of content when bandwidth is limited, followed by the other levels. Since foveation masks are used in our model, the first transmitted levels correspond to foveal regions which draw the attention, and are the most important, supplying additional efficiency related to traditional hybrid codecs. Further, the new method is faster than MOVI-Codec since it does not require arithmetic coding. In the current model, the foveation masks are fixed with respect to frame height. One future direction is to train sets of masks adaptive to contents. Another direction is to extend the framework to generate a foveation map based on frequency as well and use it to allocate the contents learnt in the latent representation.

V. CONCLUSION

We have proposed an end-to-end deep learning video compression framework that assigns bits according to a foveation protocol, assuming known visual fixations. We also achieve efficiency by training a deep space-time compression network to use displaced frame differences to compute efficient motion information by learning optimal between-frame interpolated representations. Our experimental results show that our approach, which we call FOV-MOVI-Codec, outperforms both H.264 and H.265 and foveated versions of them. The low complexity of our model, which avoids motion search, could make it amenable for implementations on resource-limited devices, such as smartphones, VR headsets, and AR glasses.

REFERENCES

- [1] B. Guenter, M. Finch, S. Drucker, D. Tan, and J. Snyder, "Foveated 3D graphics," *ACM Transactions on Graphics (TOG)*, vol. 31, no. 6, pp. 1–10, 2012.
- [2] A. Patney, M. Salvi, J. Kim, A. Kaplanyan, C. Wyman, N. Bentley, D. Luebke, and A. Lefohn, "Towards foveated rendering for gaze-tracked virtual reality," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 6, pp. 1–12, 2016.
- [3] M. Chen, Y. Jin, T. Goodall, X. Yu, and A. C. Bovik, "Study of 3d virtual reality picture quality," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 1, pp. 89–102, 2019.
- [4] W. S. Geisler and J. S. Perry, "Real-time foveated multiresolution system for low-bandwidth video communication," vol. 3299, pp. 294–305, 1998.
- [5] Z. Wang and A. C. Bovik, "Embedded foveation image coding," *IEEE Transactions on image processing*, vol. 10, no. 10, pp. 1397–1410, 2001.
- [6] Z. Wang, L. Lu, and A. C. Bovik, "Foveation scalable video coding with automatic fixation selection," *IEEE Transactions on Image Processing*, vol. 12, no. 2, pp. 243–254, 2003.
- [7] A. S. Kaplanyan, A. Sochenov, T. Leimkühler, M. Okunev, T. Goodall, and G. Rufo, "Deepfovea: Neural reconstruction for foveated rendering

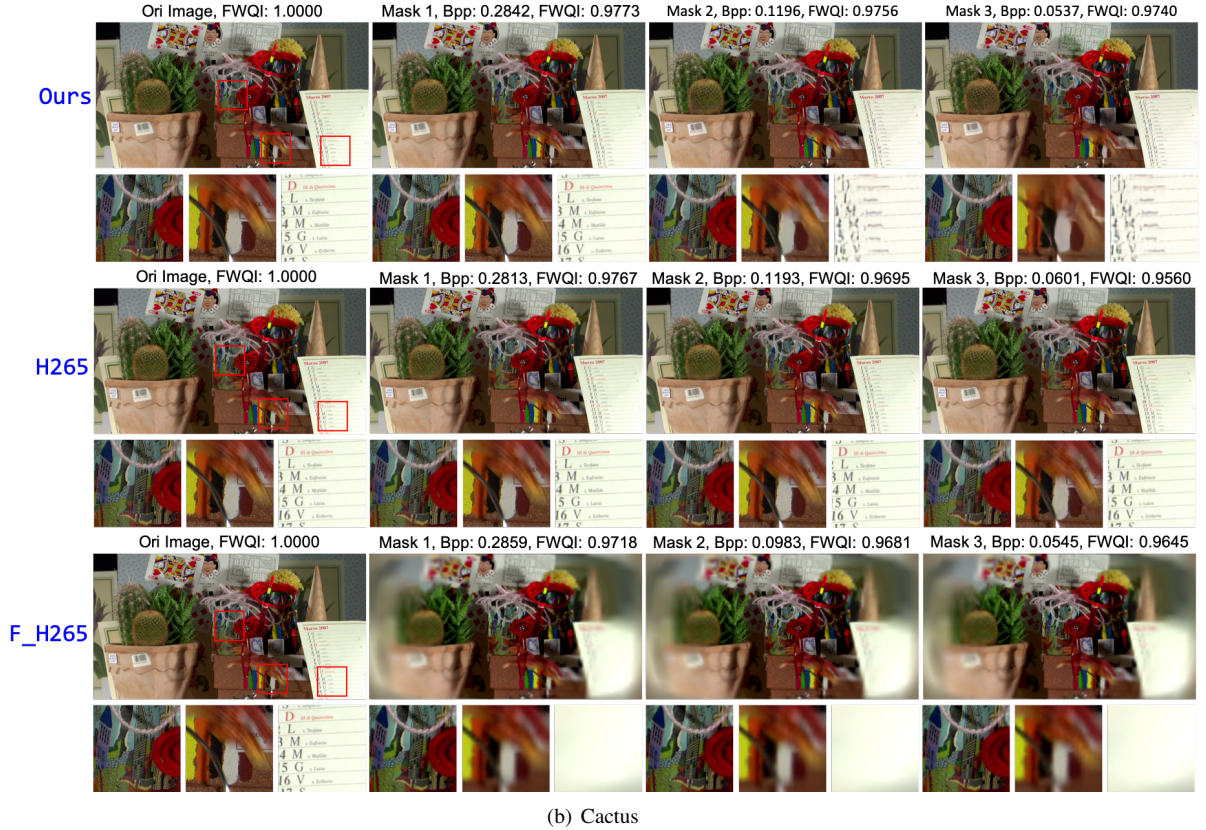
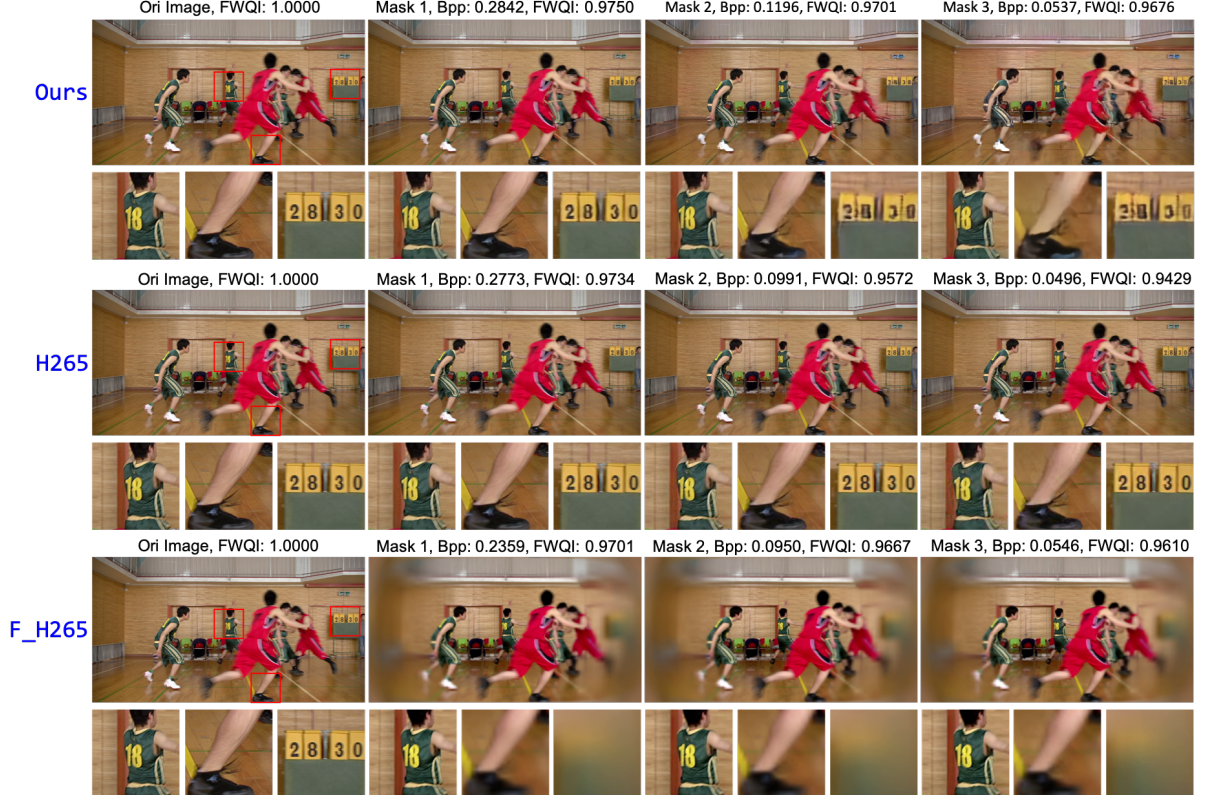


Fig. 9. Visualizations of exemplar foveated frames reconstructed by FOV-MOVI-Codec, H.265, and Foveated H.265 (denoted F_265) on the videos (a) Basketball drive and (b) Cactus.

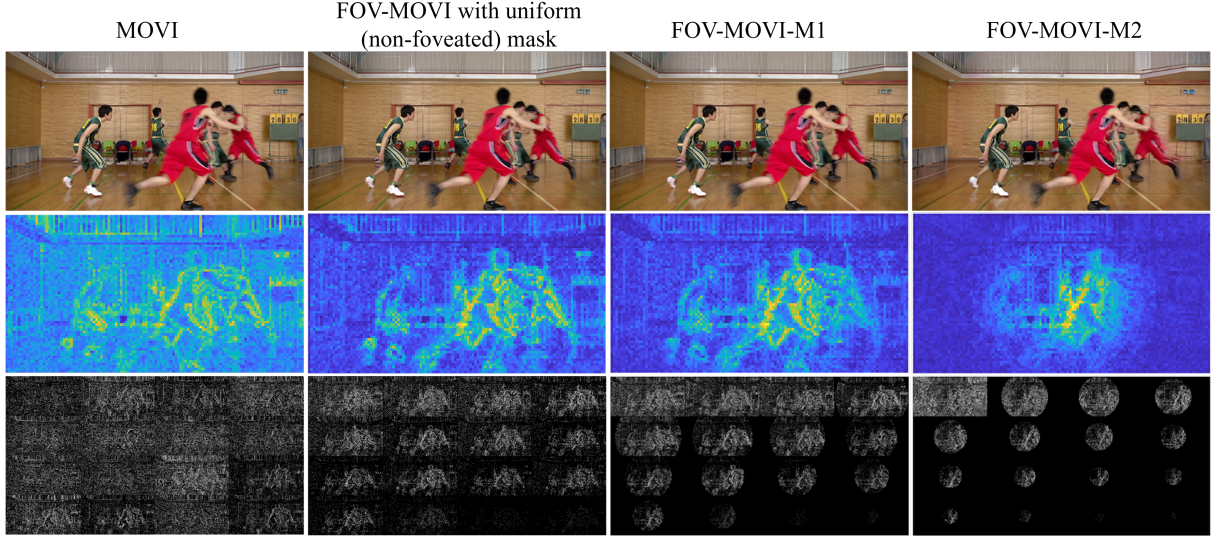


Fig. 10. Latent representations generated from four models. The first row corresponds to reconstructed frames from each model, the second row shows the cumulative latent representations, and the last row shows the latent representations at each compression level. FOV-MOVI-M1 is Foveated MOVI-Codec with foveation mask space constant $FMSC = H/2$ and FOV-MOVI-M2 is Foveated MOVI-Codec with $FMSC = H/4$, where H is the height of the frame.

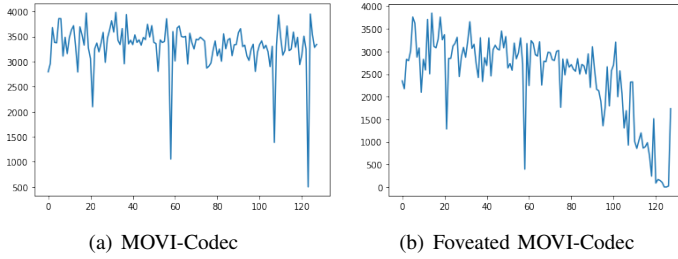


Fig. 11. Sum of latent representations for each channel, where the sum is decreasing in foveated version.

and video compression using learned statistics of natural videos,” *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–13, 2019.

- [8] C.-Y. Wu, N. Singhal, and P. Krahenbuhl, “Video compression through image interpolation,” pp. 416–431, 2018.
- [9] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, “Learning image and video compression through spatial-temporal energy compaction,” pp. 10071–10080, 2019.
- [10] R. Yang, F. Mentzer, L. V. Gool, and R. Timofte, “Learning for video compression with hierarchical quality and recurrent enhancement,” pp. 6628–6637, 2020.
- [11] G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao, “DVC: An end-to-end deep video compression framework,” pp. 11006–11015, 2019.
- [12] M. Chen, T. Goodall, A. Patney, and A. C. Bovik, “Learning to compress videos without computing motion,” *Signal Processing: Image Communication*, p. 116633, 2022.
- [13] R. Soundararajan and A. C. Bovik, “Video quality assessment by reduced reference spatio-temporal entropic differencing,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 4, pp. 684–694, 2012.
- [14] D. Y. Lee, H. Ko, J. Kim, and A. C. Bovik, “On the space-time statistics of motion pictures,” *JOSA A*, vol. 38, no. 7, pp. 908–923, 2021.
- [15] J. J. Atick and A. N. Redlich, “Towards a theory of early visual processing,” *Neural Computation*, vol. 2, no. 3, pp. 308–320, 1990.
- [16] F. Attneave, “Some informational aspects of visual perception,” *Psychological review*, vol. 61, no. 3, p. 183, 1954.
- [17] D. W. Dong and J. J. Atick, “Temporal decorrelation: a theory of lagged and nonlagged responses in the lateral geniculate nucleus,” *Network: Computation in Neural Systems*, vol. 6, no. 2, pp. 159–178, 1995.
- [18] M. Rucci and J. D. Victor, “The unsteady eye: an information-processing stage, not a bug,” *Trends in Neurosciences*, vol. 38, no. 4, pp. 195–206, 2015.
- [19] E. Chichilnisky and R. S. Kalmar, “Functional asymmetries in on and off ganglion cells of primate retina,” *Journal of Neuroscience*, vol. 22, no. 7, pp. 2737–2747, 2002.
- [20] R. Engbert, “Microsaccades: A microcosm for research on oculomotor control, attention, and visual perception,” *Progress in Brain Research*, vol. 154, pp. 177–192, 2006.
- [21] M. Poletti and M. Rucci, “A compact field guide to the study of microsaccades: Challenges and functions,” *Vision research*, vol. 118, pp. 83–97, 2016.
- [22] B. A. Olshausen and D. J. Field, “Emergence of simple-cell receptive field properties by learning a sparse code for natural images,” *Nature*, vol. 381, no. 6583, pp. 607–609, 1996.
- [23] D. Purves, G. Augustine, D. Fitzpatrick, L. Katz, A. LaMantia, J. McNamara, and S. Williams, “Functional specialization of the rod and cone systems,” *Neuroscience*, vol. 2, 2001.
- [24] B. M. Harvey and S. O. Dumoulin, “The relationship between cortical magnification factor and population receptive field size in human visual cortex: constancies in cortical architecture,” *Journal of Neuroscience*, vol. 31, no. 38, pp. 13604–13612, 2011.
- [25] H. Wässle, U. Grünert, J. Röhrenbeck, and B. B. Boycott, “Retinal ganglion cell density and cortical magnification factor in the primate,” *Vision research*, vol. 30, no. 11, pp. 1897–1911, 1990.
- [26] B. Cheung, E. Weiss, and B. Olshausen, “Emergence of foveal image sampling from learning to attend in visual scenes,” *arXiv preprint arXiv:1611.09430*, 2016.
- [27] C. Weber and J. Triesch, “Implementations and implications of foveated vision,” *Recent Patents on Computer Science*, vol. 2, no. 1, pp. 75–85, 2009.
- [28] Z. Wang, A. C. Bovik, L. Lu, and J. L. Koulouheris, “Foveated wavelet image quality index,” vol. 4472, pp. 42–53, 2001.
- [29] S. Lee, M. S. Pattichis, and A. C. Bovik, “Foveated video quality assessment,” *IEEE Transactions on Multimedia*, vol. 4, no. 1, pp. 129–132, 2002.
- [30] G. Boccignone, A. Chianese, V. Moscato, and A. Picariello, “Foveated shot detection for video segmentation,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 15, no. 3, pp. 365–377, 2005.
- [31] A. Koz and A. A. Alatan, “Foveated image watermarking,” vol. 3, pp. 657–660, 2002.
- [32] Z. Wang and A. C. Bovik, “Foveated image and video coding,” in *Digital Video Image Quality and Perceptual Coding*. CRC Press, 2005, pp. 431–458.
- [33] Z. Wang, *Rate-Scalable Foveated Image and Video Communications*. The University of Texas at Austin, 2001.
- [34] R. S. Wallace, P.-W. Ong, B. B. Bederson, and E. L. Schwartz, “Space variant image processing,” *International Journal of Computer Vision*, vol. 13, no. 1, pp. 71–90, 1994.

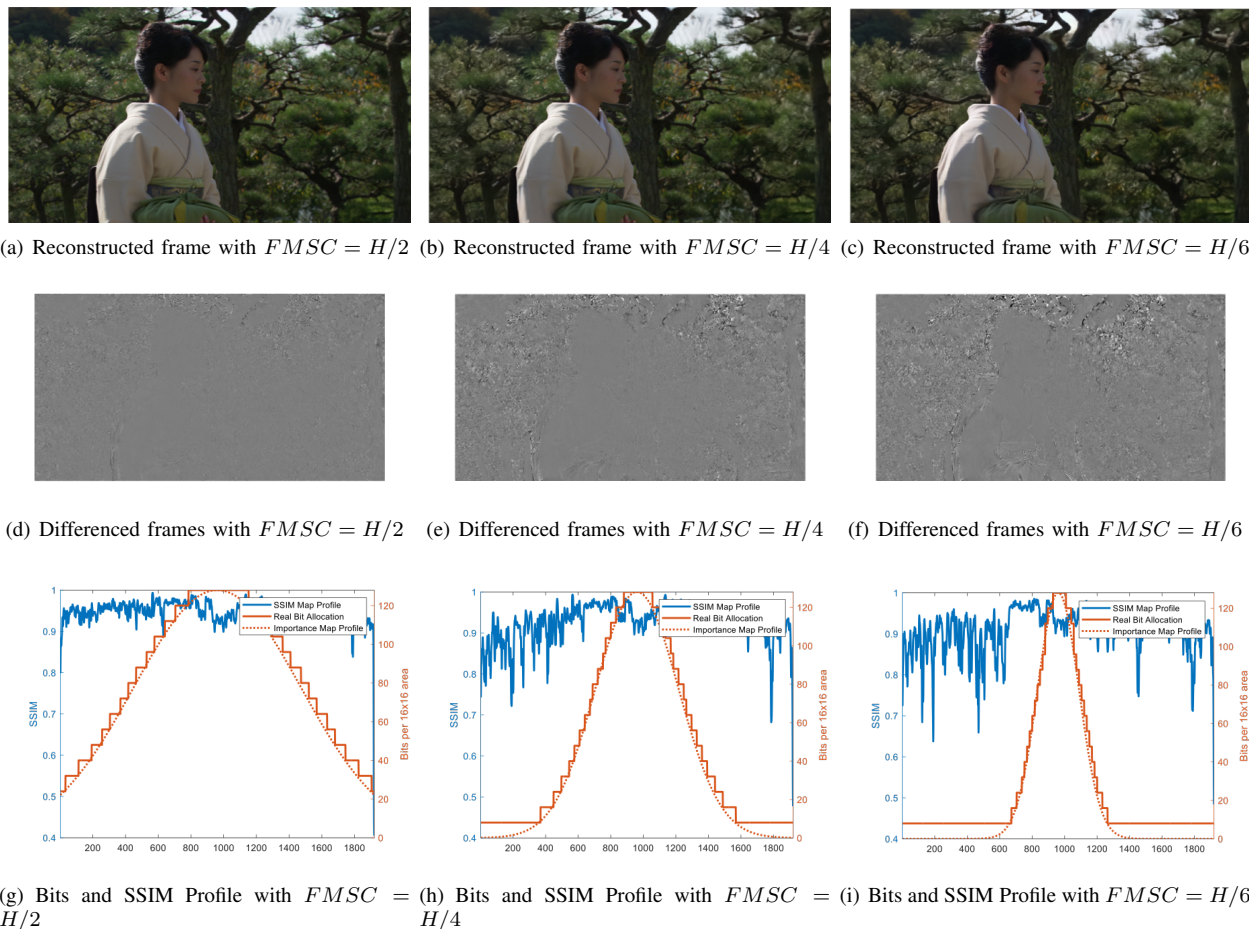


Fig. 12. Reconstructed frames, differenced frames and bit-SSIM profiles under different foveation space constants (FMSCs).

- [35] N. Tsumura, C. Endo, H. Haneishi, and Y. Miyake, "Image compression and decompression based on gazing area," vol. 2657, pp. 361–367, 1996.
- [36] S. Liu and A. C. Bovik, "Foveation embedded dct domain video transcoding," *Journal of Visual Communication and Image Representation*, vol. 16, no. 6, pp. 643–667, 2005.
- [37] H. R. Sheikh, S. Liu, Z. Wang, and A. C. Bovik, "Foveated multipoint videoconferencing at low bit rates," vol. 2, pp. II–2069, 2002.
- [38] P. J. Burt, "Smart sensing within a pyramid vision machine," *Proceedings of the IEEE*, vol. 76, no. 8, pp. 1006–1015, 1988.
- [39] Z. Li, S. Qin, and L. Itti, "Visual attention guided bit allocation in video compression," *Image and Vision Computing*, vol. 29, no. 1, pp. 1–14, 2011.
- [40] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, "Learning convolutional networks for content-weighted image compression," pp. 3214–3223, 2018.
- [41] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool, "Conditional probability models for deep image compression," pp. 4394–4402, 2018.
- [42] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [43] S. Rimac-Drlje, M. Vranješ, and D. Žagar, "Foveated mean squared error - a novel video quality metric," *Multimedia tools and applications*, vol. 49, no. 3, pp. 425–445, 2010.
- [44] S. Rimac-Drlje, G. Martinović, and B. Zovko-Cihlar, "Foveation-based content adaptive structural similarity index," pp. 1–4, 2011.
- [45] J. You, T. Ebrahimi, and A. Perkins, "Attention driven foveated video quality assessment," *IEEE Transactions on Image Processing*, vol. 23, no. 1, pp. 200–213, 2014.
- [46] M. Chen, T. Goodall, A. Patney, and A. C. Bovik, "Learning to compress videos without computing motion," *arXiv preprint arXiv:2009.14110*, 2020.
- [47] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, 2017.
- [48] J. Carreira, E. Noland, A. Banki-Horvath, C. Hillier, and A. Zisserman, "A short note about kinetics-600," *arXiv preprint arXiv:1808.01340*, 2018.
- [49] T. Xue, B. Chen, J. Wu, D. Wei, and W. T. Freeman, "Video enhancement with task-oriented flow," *International Journal of Computer Vision*, vol. 127, no. 8, pp. 1106–1125, 2019.
- [50] G. J. Sullivan, J.-R. Ohm, W.-J. Han, and T. Wiegand, "Overview of the high efficiency video coding (hevc) standard," *IEEE Transactions on circuits and systems for video technology*, vol. 22, no. 12, pp. 1649–1668, 2012.
- [51] Cisco, "Ultra video group test sequences," 2020. [Online]. Available: <http://ultravideo.cs.tut.fi>.



Fig. 13. Exemplar reconstructed foveated frames produced by FOV-MOVI-Codec, H.265, and Foveated H.265 (F_H265), from top to bottom.

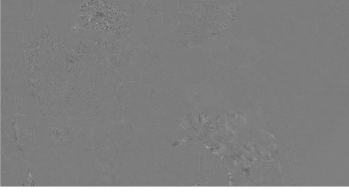
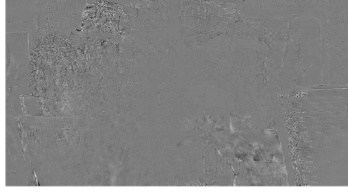
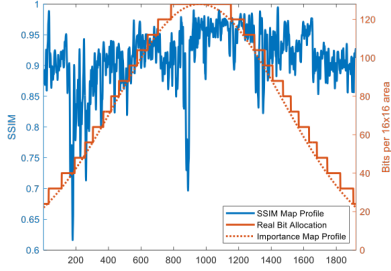
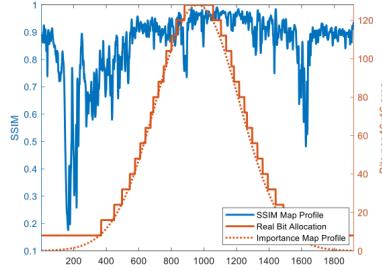
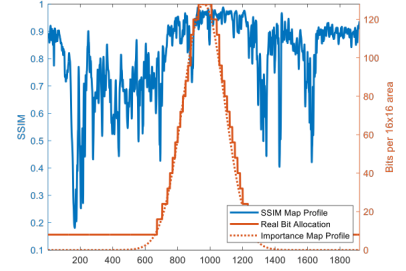
(a) Reconstructed frame with $FMSC = H/2$ (b) Reconstructed frame with $FMSC = H/4$ (c) Reconstructed frame with $FMSC = H/6$ (d) Difference image with $FMSC = H/2$ (e) [Difference image with $FMSC = H/4$ (f) [Difference image with $FMSC = H/6$ (g) Bits and SSIM Profile with $FMSC = H/2$ (h) Bits and SSIM Profile with $FMSC = H/4$ (i) Bits and SSIM Profile with $FMSC = H/6$

Fig. 14. Reconstructed frames, differenced frames and bit-SSIM profiles for different foveation mask space constants (FMSCs).