

NIRVANA: Neural Implicit Representations of Videos with Adaptive Networks and Autoregressive Patch-wise Modeling

Shishira R Maiya^{1,2,*}, Sharath Girish^{1,*}, Max Ehrlich¹, Hanyu Wang¹, Kwot Sin Lee²,

Patrick Poirson², Pengxiang Wu², Chen Wang², Abhinav Shrivastava¹

¹University of Maryland ²Snap Inc

{maxehr, sgirish, shishira, hwang}@umd.edu

{klee6, ppoirson, pwu, cwang6}@snapchat.com abhinav@cs.umd.edu

Abstract

Implicit Neural Representations (INR) have recently shown to be powerful tool for high-quality video compression. However, existing works are limiting as they do not explicitly exploit the temporal redundancy in videos, leading to a long encoding time. Additionally, these methods have fixed architectures which do not scale to longer videos or higher resolutions. To address these issues, we propose NIRVANA, which treats videos as groups of frames and fits separate networks to each group performing patch-wise prediction. The video representation is modeled autoregressively, with networks fit on a current group initialized using weights from the previous group’s model. To further enhance efficiency, we perform quantization of the network parameters during training, requiring no post-hoc pruning or quantization. When compared with previous works on the benchmark UVG dataset, NIRVANA improves encoding quality from 37.36 to 37.70 (in terms of PSNR) and the encoding speed by 12×, while maintaining the same compression rate. In contrast to prior video INR works which struggle with larger resolution and longer videos, we show that our algorithm is highly flexible and scales naturally due to its patch-wise and autoregressive designs. Moreover, our method achieves variable bitrate compression by adapting to videos with varying inter-frame motion. NIRVANA achieves 6× decoding speed and scales well with more GPUs, making it practical for various deployment scenarios.

1. Introduction

In the information age today, where petabytes of content is generated and consumed every hour, the ability to compress data fast and reliably is important. Not only

*First two authors contributed equally

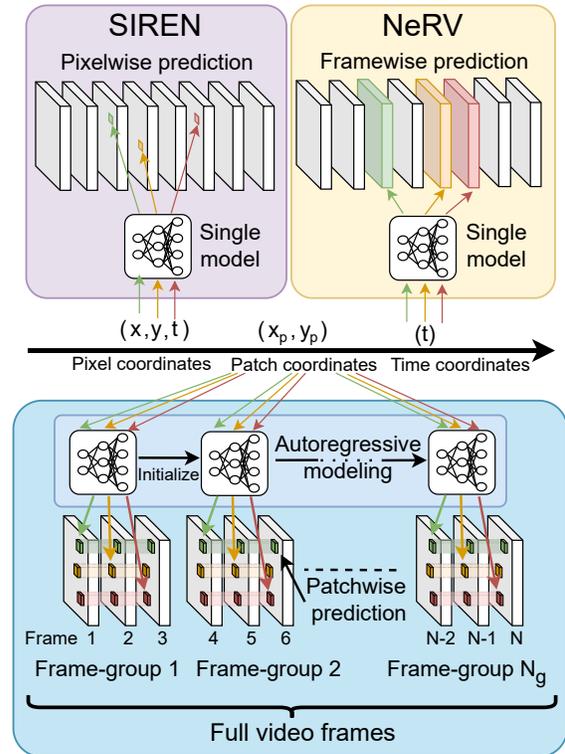


Figure 1. **Overview of NIRVANA:** Prior video INR works perform either pixel-wise or frame-wise prediction. We instead perform spatio-temporal patch-wise prediction and fit individual neural networks to groups of frames (clips) which are initialized using networks trained on the previous group. Such an autoregressive patch-wise approach exploits both spatial and temporal redundancies present in videos while promoting scalability and adaptability to varying video content, resolution or duration.

does compression make data cheaper for server hosting, it makes content accessible to population/regions with low-bandwidth. Conventionally, such compression is achieved through codecs like JPEG [45] for images and HEVC [40], AV1 [10] for videos, each of which compresses data via tar-

getted hand-crafted algorithms. These techniques achieve acceptable trade-offs, leading to their widespread usage.

With the rise of deep learning, machine learning-based codecs [4,5,31,43] showed that it is possible to achieve better performance in some aspects than conventional codecs. However, these techniques often require large networks as they attempt to generalize to compress all data from the distribution. Furthermore, such generalization is contingent on the training dataset used by these models, leading to poor performance for Out-of-Distribution (OOD) data across different domains [47] or even when the resolution changes [7]. This greatly limits its *real-world practicality* especially if the input data to be compressed is significantly different from what the model has seen during training. In recent years, a new paradigm, Implicit Neural Representations (INR), emerged to solve the drawbacks of model-learned compression methods. Rather than attempting to generalize to all data from a particular distribution, its key idea is to train a network that specifically fits to a signal, which can be an image [37], video [8], or even a 3D scene [30]. With this conceptual shift, a neural network is no longer just a predictive tool, rather it is now an efficient storage of data. Treating the *neural network as the data*, INR translate the data compression task to that of model compression. Such a continuous function mapping further benefits downstream tasks such as image super-resolution [23], denoising [32], and inpainting [37].

Despite these advances, videos vary widely in both spatial resolutions and temporal lengths, making it challenging for networks to encode videos in a practical setting. Towards solving this task, an early method, SIREN [37], attempted to learn a direct mapping from 3D spatio-temporal coordinates of a video to each pixel’s color values. While simple, this is computationally inefficient and does not factor in the spatio-temporal redundancies within the video. Later, NeRV [8] proposed to instead map 1D temporal coordinates in the form of frame indices directly to generate a whole frame. While this improves the reconstruction quality, such a mapping still does not capture the temporal redundancies between frames as it treats each frame individually as a separate image encoding task. Finally, mapping only the temporal coordinate also means one would need to modify the architecture in order to encode videos of different spatial resolutions.

To address the above issues, we propose NIRVANA, a method that exploits spatio-temporal redundancies to encode videos of arbitrary lengths and resolutions. Rather than performing a pixel-wise prediction (*e.g.*, SIREN) or a whole-frame prediction (*e.g.*, NeRV), we predict *patches*, which allows our model to adapt to videos of different spatial resolutions without modifying the architecture. Our method takes in the centroids of patches (patch coordinates) (x_p, y_p) as inputs and outputs a corresponding patch vol-

ume. Since patches can be arranged for different resolutions, we do not require any architectural modification when the input video resolution changes. Furthermore, to exploit the temporal nature of videos, we propose to train individual, small models for each group of video frames (“clips”) in an autoregressive manner: the model weights for predicting each frame group is initialized from the weights of the model for the previous frame group. Apart from the obvious advantage that we can scale to longer sequences without changing the model architecture, this design exploits the temporal nature of videos that, intuitively, frame groups with similar visual information (*e.g.*, static video frames) would have similar weights, allowing us to further perform residual encoding for greater compression gains. This *adaptive* nature, that static videos gain greater compression than dynamic ones, is a big advantage over NeRV where the compression for identical frames remain fixed as it models each frame individually. To obtain further compression gains, we employ recent advances in the literature to add entropy regularization for quantization, and encode model weights for each video during training [19]. This further adapts the compression level to the complexity of each video, and avoids any post-hoc pruning and fine-tuning as in NeRV, which can be slow.

Finally, we show that despite its autoregressive nature, our model is linearly parallelizable with the number of GPUs by chunking each video into disjoint groups to be processed. This strategy largely improves the speed while maintaining the superior compression gains of our method, making it practical for various deployment scenarios.

We evaluate NIRVANA on the benchmark UVG dataset [29]. We show that NIRVANA reaches the same levels of PSNR and Bits-per-Pixel (BPP) compression rate with almost $12\times$ the encoding speed of NeRV. We verify that our algorithm adapts to varying spatial and temporal scales by providing results on videos in the UVG dataset with 4K resolution at 120fps, as well as on significantly longer videos from the YouTube-8M dataset, both of which are challenging extensions which have not been attempted on for this task. We show that our algorithm outperforms the baseline with much smaller encoding times and that it naturally scales with no performance degradation. We conduct ablation studies to show the effectiveness of various components of our algorithm in achieving high levels of reconstruction quality and understand the sources of improvements.

Our contributions are summarized below:

- We present NIRVANA, a patch-wise autoregressive video INR framework which exploits both spatial and temporal redundancies in videos to achieve high levels of encoding speedups ($12\times$) at similar reconstruction quality and compression rates.
- We achieve a $6\times$ speedup in decoding times and scale well with increasing number of GPUs, making it prac-

tical in various deployment scenarios.

- Our framework adapts to varying video resolutions and lengths without performance degradations. Different from prior works, it achieves adaptive bitrate compression based on the amount of inter-frame motion for each video.

2. Related Work

Implicit Neural Representations (INR) are a novel family of methods designed to map a set of coordinates to a specific signal - such as a single image or video - using a neural network as a function for such mappings. SIREN [37] demonstrates that by utilizing periodic activation functions in MLPs, such a function can be fit and used to map a wide array of signals, including images, 3D shapes and videos. As an alternative, [42] shows that an INR network with standard activations can be trained by utilizing random Fourier features. [28] and [35] propose adaptive block-based approaches whose complexity mirrors the underlying signals. Frequency-based approaches are proposed in [15, 27, 36] that enable multi-scale representations. The first image specific INR method is COIN [12], which is extended to encode multiple images through network modulations in COIN++ [13]. A method for learning local implicit functions is proposed in [9] that results in smoother super-resolution outputs. Several methods [11, 39, 41] have explored meta-learning approaches to reduce the long encoding times of image INRs; [22] further shows that directly initializing a meta-sparse network not only gives a good initialization but also helps with model compression.

INR for videos. Despite the significant advances in INR methods for image compression, videos present a more challenging task for INR methods. For example, if we naively add time as an extra dimension to image-based methods, such as in SIREN [37], the resulting outputs are grainy. In [34, 48], INR methods that utilize flow-based information to encode videos are introduced; however, they cannot scale beyond short low-resolution videos. NeRV [8] is the first method to scale video compression using INRs. They modify the implicit mapping function to learn a direct mapping from a video frame index to an entire frame. Further extensions of NeRV, such as patch-based versions [3, 25], provide minor improvements over the original architecture. Despite good reconstruction, the problems of long encoding times, the lack of inter-frame encoding, and the inability to adapt to video content act as major drawbacks for widespread adoptions.

Model compression is typically achieved by pruning or quantization of network weights. A plethora of works perform model pruning with minimal loss of performance [16–18, 20]. Pruned models contain a majority of zeros and can be stored in sparse matrix formats [21] for reduction in model size. Alternatively, quantization works [14, 26, 38, 44]

to reduce the number of bits needed to store each model weight, resulting in reduced disk space. As implicit neural networks represent the data using their model weights, data compression translates to model compression. In this work, we adopt the works of [19, 33] for model compression by maintaining a set of quantized parameters during training which are then stored on disk. These recent methods have shown to achieve high levels of compression through entropy regularization without sacrificing downstream network performance. We perform quantization and model compression *during training*, unlike the post-hoc pruning and quantization in NeRV [8].

3. Approach

3.1. Background

Given a video $V \in \mathbb{R}^{N \times H \times W \times 3}$ consisting of N frames, each with spatial resolution $H \times W$, an INR defines a mapping from the spatio-temporal coordinate $X = (x, y, t); X \in \mathbb{R}^3$ to the pixel value $p \in \mathbb{R}^3$. Thus, it implicitly represents a parameterized function $h_\theta : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ parameterized by θ . The function is typically trained by minimizing the MSE-loss $\|h_\theta(X) - V\|_2$. While this is a straightforward extension of image-based INR methods to the spatio-temporal domain, it fails to exploit the spatial and temporal consistency in videos. Pixel-wise prediction leads to redundant computation and long encoding times while also producing blurry outputs [37], [12]. To mitigate this issue, NeRV [8] proposes to directly encode the frame index $t \in \mathbb{R}$ as a positional embedding input to a model which outputs the entire image frame $\mathbb{R}^{H \times W \times 3}$. NeRV consists of several MLP layers followed by convolution layers which upsamples the latent representation to the target frame’s spatial resolution. While this formulation improves upon the naive pixel-based formulation, it does not adapt to arbitrary resolutions, and does not capture the temporal dependencies between frames as it effectively acts as only an image encoder for each frame.

3.2. Autoregressive Patch-wise Modeling

Patch prediction. The two dominant INR paradigms for video encodings, SIREN [37] and NeRV [8] represents two extreme ends of a spectrum: the former predicts every pixel in a video volume independently, while the latter predicts the pixels for a single frame simultaneously. Through pixel-wise prediction, SIREN accommodates varying the output image’s spatial resolution but does not exploit the spatial consistency of the image. In contrast, while NeRV exploits the spatial consistency of an image through convolutions, it cannot vary the image’s spatial resolution. We adopt a middle ground between these two extremes by adopting a modeling approach that instead predicts *patches* of an image. This gives us the best of both approaches, since our

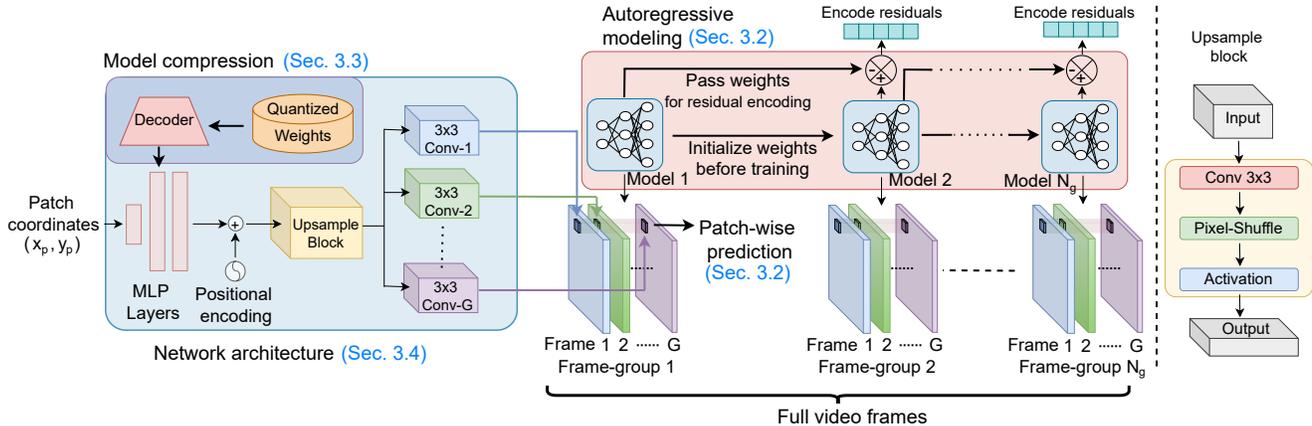


Figure 2. **Overview of NIRVANA:** We propose an autoregressive video INR framework which performs patch-wise prediction of groups of frames by fitting separate networks to each group. Each network is initialized with the previous group’s network weights. Our architecture consists of several SIREN MLP layers followed by an upsampling block (right). It takes patch coordinates as inputs and outputs patches across a group of G frames. We maintain a set of quantized weights which are decoded to obtain the network weights. Post training, we encode weight residuals which are the difference between the quantized weights of the current and previous group’s network.

model utilizes the spatial consistency of an image while still naturally scaling to varying image resolutions.

We push further in this direction by exploiting the temporal consistency in videos to predict a *volume of patches* across neighboring frames. We thus predict a patch group $\mathbf{P} \in \mathbb{R}^{G \times H_p \times W_p \times 3}$, where G is the number of frames in a group and (H_p, W_p) is the patch size. This method enables us to reduce the amount of redundant computation in both the spatial and temporal dimensions, leading to significantly shorter encoding times compared to NeRV or SIREN.

Autoregressive networks. While it is straightforward to input a 3D patch coordinate (x_p, y_p, g) (where (x_p, y_p) are the patch centroids within a frame and g is the frame index within each group) and output a patch volume using a single network, it still suffers from adaptability to varying video content, resolutions, or durations as mentioned in Section 3.1. To overcome these challenges, we propose to autoregressively fit separate networks to each frame group. Each network is fed with the same set of inputs, namely the centroids of patches (x_p, y_p) , and it outputs the corresponding 3D patch volume of the group. Every subsequent network is finetuned from the previous one, leading to shorter encoding times. As video content does not change much over a short time period of multiple frames, the difference in weights (or weight residuals) after fine-tuning is small. Thus, encoding the weight residuals instead of the weights themselves leads to higher compression rates. The design of the algorithm 1 allows us to encode multiple chunks of a long video in a parallel manner, a key feature missing from existing methods. This means NIRVANA can scale linearly with the number of GPUs without any drop in performance (see Section 4.6).

3.3. Model Compression and Weight Storage

Since the network weights are the latent representations for the video, network size directly translates to bitrate of the video encoding. To reduce network size, we adapt existing works which perform model weights/latent representation compression [19, 33]. For a weight parameter $\mathbf{W} \in \theta$, where θ is the set of model parameters, we maintain a corresponding quantized latent weight $\tilde{\mathbf{W}}$. The continuous weight parameter \mathbf{W} is then obtained as $\mathbf{W} = f_\phi(\tilde{\mathbf{W}})$ where f_ϕ is a learned linear transform. The entire setup is trained end to end, without any post-hoc quantization and fine-tuning. As previously explained, we encode the quantized latent residuals instead of the latents themselves to achieve higher levels of compression. We encode these residuals using arithmetic coding [46], a lossless entropy-based compression algorithm. In order to encourage the latents to have lower entropy, we add an entropy regularization term to our loss function. This term encourages the network to have a lower entropy and hence a lower bitrate. When decoding, each weight is obtained sequentially by cumulatively adding residuals. This approach helps in making the bitrate of NIRVANA adaptive to the video content: for frame-groups that have little motion, they are already closer to convergence and thus have small differences in their network weights, leading to sparser residuals and subsequently, lower bitrates. This feature is missing in other methods as the models are fixed for a given video.

3.4. Network Architecture

In this section, we describe the network architecture which is used for each frame group, as illustrated in Fig. 2. For a group of G frames, we segment patch volumes

Algorithm 1 Sequential Video INR

- 1: **Init** Randomly initialize network h_{θ_0} with initial weights θ_0 and training iterations T . The video contains N frames segmented into N_g frame groups of size G each.
 - 2: **for** g in $0, 1, 2, \dots, N_g - 1$ **do**
 - 3: **if** $g = 0$ **then**
 - 4: Train h_{θ_0} for T iterations for all patches
 - 5: Store weights h_{θ_0}
 - 6: **else**
 - 7: Initialize weights $h_{\theta_k} \leftarrow h_{\theta_{k-1}}$
 - 8: Finetune h_{θ_k} for T_r iterations for all patches
 - 9: Store quantized latent of residuals $h_{\theta_k} - h_{\theta_{k-1}}$
 - 10: **end if**
 - 11: **end for**
-

of shape (H_p, W_p, G) . The input to the network is the 2D patch coordinate $(x_p, y_p) \in \mathbb{R}^2$ and the output is the corresponding RGB patch volume $\mathbf{P} \in \mathbb{R}^{G \times H_p \times W_p \times 3}$. We stack multiple MLP layers with SIREN activations to obtain an output feature representation vector $s_p \in \mathbb{R}^d$ of dimension d . We replicate s_p by G times, and add positional encoding vectors based on the position of the frame within the group, using the following embedding function:

$$\tau(t, 2i) = \sin\left(\frac{t}{f^{2i}}\right) \tau(t, 2i + 1) = \cos\left(\frac{t}{f^{2i}}\right), i \in [0, d) \quad (1)$$

where $t \in [0, G)$ represents the position of the frame within the group of G frames. We then add a decoder block as in [8] followed by a 3×3 convolutional layer to output G patches. For a video with N frames, we segment it into N_g frame-groups with each group consisting of G frames ($N = N_g \times G$). For the g^{th} frame-group ($g \in [0, N_g)$), the corresponding network is represented as h_{θ_g} consisting of parameters θ_g . The overall loss objective for training the network for the g^{th} frame-group is therefore

$$\mathcal{L}_g = \mathcal{L}_{\text{mse}}(h_{\theta_g}(p), v_g) + \lambda_I \mathcal{L}_{\text{ent}}(\theta_g) \quad (2)$$

where p represents patch grid coordinates and v_g means the corresponding ground-truth frame-group pixel values. $\mathcal{L}_{\text{ent}}(\theta_g)$ represents the entropy regularization loss on the model parameters θ_g . It is weighed by the coefficient λ_I controlling the rate-distortion trade-off for reconstructing the frame groups.

4. Experiments

4.1. Datasets and Implementation Details

The standard benchmark UVG dataset [29] is used to compare our approach NIRVANA with prior video INR

works. Following similar setups [8], our approach is evaluated on 7 videos from the dataset at 1080p resolution (UVG-HD) and 120 fps with 6 videos consisting of 600 frames and 1 with 300 frames. To show our scalability for higher resolution videos, we show results for the 7 videos at 4K resolution (UVG-4K) as well. We additionally include a video from the Youtube-8M (see Appendix) [1] dataset at 1080p resolution and 60 fps with 3 separate versions segmented at 2000/3000/4000 frames to demonstrate our model’s capability for long videos. We use the standard PSNR (in dB) to measure the reconstruction quality and bits-per-pixel (BPP) to measure the compression rate. We also include encoding times as well as their decoding speed in fps.

The MLP network consists of 5 SIREN layers with a layer size of 512 and sine activation. The network predicts 32×32 patches for 3 frames ($G = 3$) in all our experiments unless mentioned otherwise. The number of iterations is set to 16000 for the first group in order to obtain a good initialization and 2000 iterations for subsequent groups. We set the learning rate to be $5e-4$ and optimize the network with the MSE and entropy regularization loss. The entropy loss weight coefficient λ_I as defined in Equation 2 is set to $1e-4$. In practice, the coefficient can be varied to control the trade-off between PSNR and BPP. We use the torchac library to perform arithmetic encoding of the quantized weight residuals. Since the convolutional layers typically contain only a fraction of the total network parameters ($\sim 10\%$), we do not quantize their weights and use the LZMA compression method for storing their residuals.

We use pixel-wise method SIREN [37] and frame-wise method NeRV [8] as our baselines. For SIREN, we use a 5-layer MLP with hidden dimension of 2048. For NeRV, we use the NeRV-L configuration as specified in the paper. Encoding times reported are equivalent to when fully run on a single NVIDIA RTX 2080 GPU. For NeRV, we fit separate models to each video and remove 40% of the parameters during the pruning stage with the remaining weights quantized to lowest possible bit-width without significant performance drop. Further implementation details can be found in the Appendix.

4.2. UVG-HD

Comparisons on the UVG-HD dataset are summarized in Table 1. By varying the patch size, we let NIRVANA achieve similar BPP to SIREN and NeRV respectively. NIRVANA outperforms SIREN by a significant margin in terms of PSNR while having $6\times$ shorter encoding times. Similarly, our approach obtains speedups of $\sim 12\times$ compared to NeRV while still achieving marginally higher PSNR (+0.34dB) and lower BPP. Additionally, we obtain a decoding speed of ~ 65 FPS which is nearly $65\times$ and $6\times$ speedup in inference time/decoding compared to SIREN and NeRV respectively. This shows the efficacy of our

Dataset	Method	Encoding Time (Hours) ↓	Decoding Speed (FPS) ↑	PSNR ↑	BPP ↓
UVG-HD	SIREN	~30	15.62	27.20	0.28
	NIRVANA (Ours)	5.44	87.91	34.71	0.32
	NeRV	~80	11.01	37.36	0.92
	NIRVANA (Ours)	6.71	65.42	37.70	0.86
UVG-4K	NeRV	~134	8.27	35.24	0.28
	NIRVANA (Ours)	20.89	50.83	35.18	0.27

Table 1. **Comparison with video INR approaches on UVG benchmarks.** We vary patch size of NIRVANA on UVG-HD to match the BPP of SIREN and NeRV respectively. NIRVANA achieved much faster encoding and decoding speed, while maintaining better or on-par quality at comparable BPP.

framework to reduce redundant computation in both the spatial and temporal domains. We show our qualitative results in Fig. 6.

4.3. UVG-4K

To analyze the spatial adaptability of our approach, we test our method on the UVG-4K dataset, with results shown in Table 1. Compared to NeRV, we achieve a $\sim 6\times$ speedup in both encoding and decoding times, while maintaining similar PSNR (35.24 vs 35.18) and slightly better BPP (0.28 vs 0.27). Furthermore, to adapt to such higher resolution data, our method does not require any architectural modifications. In contrast, NeRV requires architectural modifications by adding a $2\times$ upsampling block to increase the resolution of the predicted image frames. Note that a higher PSNR can be achieved with longer training schedules as we show in Section 5.4, but we limit to 2000 iterations to maintain consistency across datasets.

4.4. Long Videos

We now analyze the effect of increasing video duration for our approach. We utilize a video from the Youtube-8M dataset and evaluate on 3 separate segments consisting of the first 2000/3000/4000 frames. Results are summarized in Table 2. Our approach maintains a similar PSNR (< 0.3 drop) with increased number of frames while still encoding at a similar bitrate (< 0.04 increase). In contrast, even with higher encoding times ($4\times$ slower), NeRV suffers from significant degradation on longer videos with PSNR dropping from 33.38 \rightarrow 31.6 \rightarrow 30.53. Since NeRV’s model size remains constant, its BPP reduces with increased number of frames. However, the fixed number of network parameters limits its ability to fit to a larger set of frames, leading to performance drops.

Additionally, since our approach is autoregressive, it needs to be trained only once even with increasing number of frames. Networks for future frames are simply initialized with the weights of the previous networks and trained before encoding the weight residuals. Such a modeling makes it suitable for online scenarios as well with a constant

Num Frames	Method	Encoding Time (Hours) ↓	PSNR ↑	BPP ↓
2000	NeRV	84.44	33.38	0.22
	NIRVANA (Ours)	20.85	35.43	0.62
3000	NeRV	134.58	31.6	0.16
	NIRVANA (Ours)	31.37	35.21	0.64
4000	NeRV	190.30	30.53	0.12
	NIRVANA (Ours)	41.84	35.15	0.65

Table 2. **Video duration adaptability:** For longer videos, we maintain similar reconstruction quality (~ 35 PSNR) and compression rate (~ 0.62 BPP). We retain a significantly faster encoding speed than NeRV which suffers from significant degradation with increased number of frames.

stream of frames. In contrast, NeRV requires separate models to be trained for different video durations as each training epoch consists of training on all frames. Note that both approaches scale linearly with increased video duration, but NeRV fits the same model to larger video signals, leading to performance drops. Specific architectural modifications are necessary to improve the PSNR for longer videos which comes at the cost of even higher encoding times.

4.5. Adaptive Compression

Videos can consist of different levels of inter-frame motion with more static videos containing higher levels of temporal redundancy in comparison to dynamic ones. To illustrate the capability of our approach to exploit such redundancies, we evaluate the compression rate on 6 videos in the UVG-HD dataset which consist of 600 frames. We sort each video according to their average MSE between subsequent frames, which serves as a proxy to the amount of temporal redundancy in the video. More static scenes like Honeybee have a lower MSE compared to highly dynamic scenes like Jockey. We plot the PSNR and BPP of NIRVANA and NeRV with increasingly dynamic video content and show the results in Fig. 3. Note that the average PSNR/BPP over the 6 videos can be increased or decreased by varying other hyperparameters such as patch size, number of groups, entropy loss coefficient etc. (as shown in Section 5), but we focus on adaptability to videos for a given hyperparameter configuration. We see that our approach has an adaptive bitrate compression with more static scenes like Honeybee (MSE $2.2e-4$) that has a lower bitrate (0.51), compared to dynamic scenes such as Jockey (MSE $0.9e-3$) which are allocated more bits (0.96). We maintain similar PSNR as NeRV which has a constant BPP due to the same model applied to all videos. While NeRV’s quantization bit width can be reduced further for lower BPP, it is a post-hoc approach which comes at the cost of PSNR and requires tuning for each video. In contrast, our approach adaptively varies the BPP during training with no change in hyperparameters.

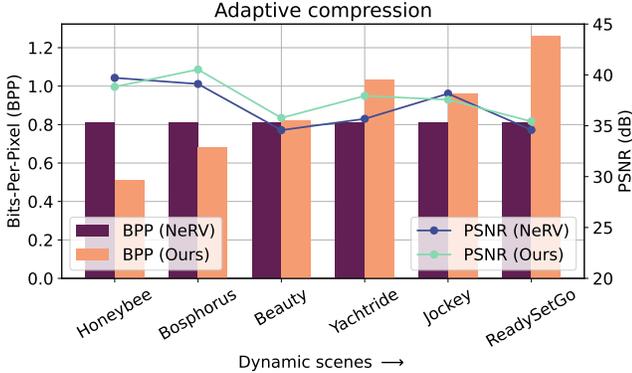


Figure 3. **Video content adaptability:** 6 videos are sorted in increasing order of variation between subsequent frames. Our approach shows adaptive bitrate compression, with more static scenes exhibiting lower BPP, while highly dynamic ones being allocated more bits while maintaining a similar PSNR as NeRV (and $12\times$ encoding speed).

4.6. GPU Parallelization

We now analyze the scalability of our approach with a larger number of GPUs. In Fig. 4 we plot the encoding times for “Jockey” (both 1080p and 4K versions) from UVG dataset for NeRV (using distributed training) and our methods. The design in Algorithm 1 allows different *chunks* of the source video to be processed autoregressively on separate GPUs. As the number of GPUs are scaled with a factor of $2\times$, our approach achieves close to linear scaling with very little overhead for the case of UVG-4K ($1.0\times \rightarrow 2.0\times \rightarrow 3.8\times \rightarrow 7.3\times$) compared to a weaker scaling of NeRV ($1.0\times \rightarrow 1.7\times \rightarrow 2.7\times \rightarrow 4.3\times$). UVG-HD shows a higher amount of overhead but still scales fairly well with increased GPUs compared to NeRV. Thus, we see the capability of parallelization of our approach with higher number of GPUs. Also note that the time taken by NeRV for HD videos on 8 GPUs is still higher than the time taken by NIRVANA on a single GPU.

5. Ablation Studies

In this section, we study the impact of various parameters of our approach on the PSNR-BPP tradeoff curves as seen in Fig. 5. By varying the entropy loss coefficient, we obtain different points on the tradeoff curve with a higher coefficient leading to lower BPP (low rate) but also lower PSNR (high distortion). We additionally show the convergence effects of longer training for each group with increased number of iterations. Results are evaluated on the Jockey video of the UVG-HD dataset. While varying each parameter along with the entropy coefficient, we fix other parameters to their default values of patch size at 32×32 , group size at 3, and number of iterations at 2000. We sample values of the entropy coefficient λ_I between $1e-5$ and

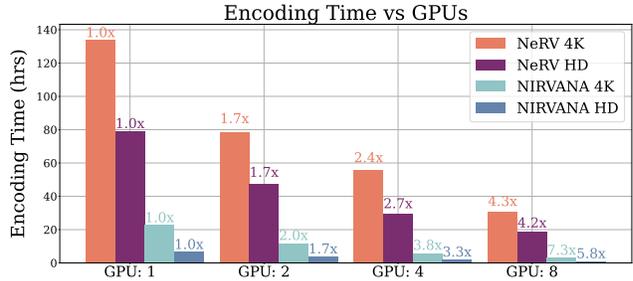


Figure 4. **GPU scalability of NIRVANA:** We compare scalability of our approach with NeRV in terms of encoding time with increasing number of GPUs at two video resolutions: 1080p and 4K. We scale close to linear for 4K and have much lower overhead compared to NeRV for both resolutions.

$5e-4$ to obtain various points on the tradeoff curve.

5.1. Effect of Entropy Regularization

We analyze the effect of varying the entropy coefficient λ_I and obtain a PSNR-BPP curve visualized in Fig. 5(a). In general, we see that increasing λ_I decreases the BPP at the cost of PSNR. This is to be expected as increasing the entropy regularization forces the quantized weights of each frame group’s networks to lie in fewer quantization bins. Consequently, more weight residual (difference between quantized weights of subsequent frame groups’ networks) values are 0 leading to lower entropy of the weights and subsequently lower BPP of the model. The entropy coefficient thus provides a natural way of controlling the PSNR-BPP tradeoff according to the required application.

5.2. Effect of Patch Size

We vary the patch prediction size of our network from 8×8 to 48×48 in steps of 8. We visualize the results in Fig. 5(b). In general, increasing patch size shifts the curve upwards and to the right corresponding to higher PSNR but also high BPP. A higher patch size results in an increase in number of network parameters (both in convolutional and SIREN layers) and hence its expressivity, leading to higher PSNR. However, as patches are less localized, the outputs of networks between subsequent frame-groups vary more significantly with dynamic scenes (such as Jockey), leading to larger residuals. This increases the entropy of the residuals and as a result, the BPP.

5.3. Effect of Frame Group Size

We vary the frame group size from 2 to 8 in steps of 1 and visualize the results in Fig. 5(c). Increasing the group size, in general, reduces BPP at the cost of PSNR. This is expected as a single model shares computation across a larger group of frames effectively leading to fewer parameters per frame and lower BPP. However, we notice that group size 3 obtains the best tradeoff curve in the < 0.8

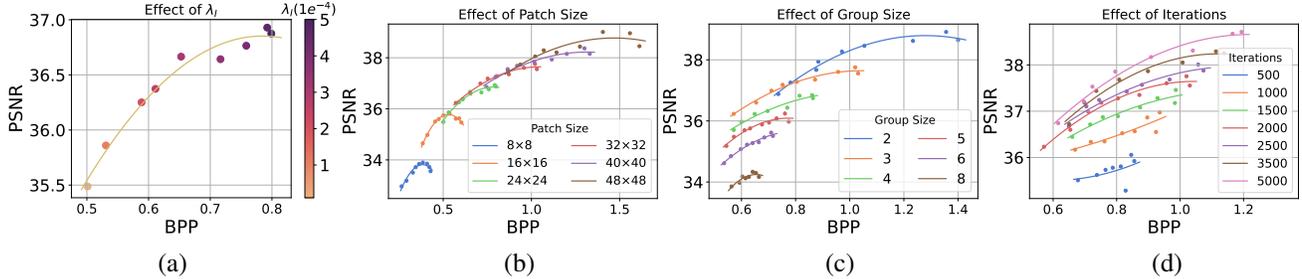


Figure 5. (a) Effect of entropy regularization: the larger λ_I is, the lower the entropy and the BPP. (b, c) As the patch size or group size increases, PSNR increases at the cost of higher BPP. (d) The longer the training iteration (encoding time) is, the higher the PSNR gains.

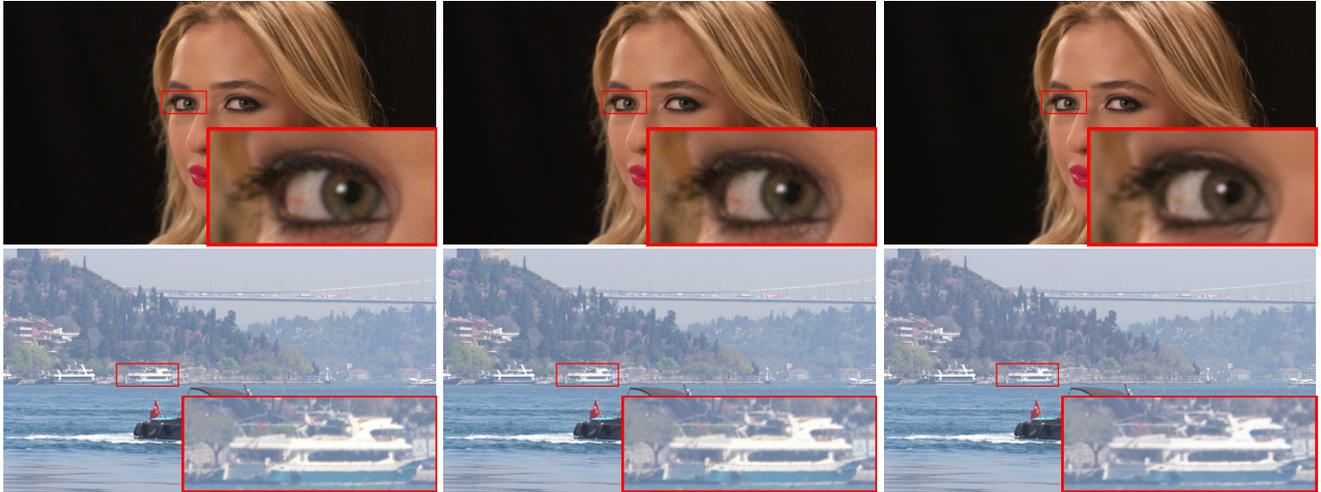


Figure 6. (Left) Ground truth video frames. (Center) Reconstruction from NIRVANA. (Right) Reconstruction from NeRV. We show that NIRVANA is able to preserve the image fidelity after reconstruction, capturing important details such as the veins in the eye of Beauty, and the color quality in the Bosphorus video.

BPP regime with higher group size detrimental to the performance. This is likely because of the fixed MLP representation capability which learns a shared global representation for all frames within a group. For a dynamic scene such as Jockey with significant pixel shift between frames (Section 4.5), a larger MLP is necessary for capturing the variations within a group.

5.4. Effect of Number of Training Iterations

To analyze the effect of longer training schedules, we vary the number of training iterations for the network for each frame group from 500 to 5000. Results are shown in Fig. 5(d). As expected, increasing the number of iterations improves the PSNR-BPP tradeoff with the curve shifting upwards. This shows that our network can obtain higher quality reconstructions for longer training times at no cost to the compression rate. This can be made feasible with higher number of GPUs as shown in Section 4.6. However, increasing iterations provides diminishing gains as we observe the curves approaching closer to each other with higher iterations.

6. Qualitative results

We visualize the reconstructions of frames in videos from the UVG-HD dataset in Fig. 6. Our approach (center) preserves image fidelity after reconstruction when compared to the ground-truth frame (left). It captures subtle details such as the veins in the eye (top row) or maintaining the color information of the boat (bottom row) when compared to the ground-truth. Further visualizations from videos in the UVG-4K dataset are provided in the Appendix.

7. Conclusion

In this work, we propose an autoregressive video INR framework, NIRVANA, which segments videos into groups of frames and fits separate neural networks to each group. Each network performs a patch-wise prediction across the group of frames thus exploiting both the spatial and temporal redundancy present in videos, improved from the previous works. Each network is initialized with the weights of the previous frame-group’s trained network. We additionally quantize weights during training, requiring no post-hoc pruning or quantization and store weight residuals between

subsequent frame group’s networks to obtain high levels of compression. NIRVANA achieves $12\times$ speedups on standard datasets compared to previous methods while maintaining similar levels of reconstruction quality and compression rate. NIRVANA adapts to varying video resolution and duration without large performance degradation and no architectural modifications. Additionally, our framework adaptively compresses videos based on their inter-frame variation. We achieve high levels of decoding speed compared to prior video INR approaches and also scale better with higher number of GPUs.

References

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016. **5**
- [2] Pontus Andersson, Jim Nilsson, Tomas Akenine-Möller, Magnus Oskarsson, Kalle Åström, and Mark D. Fairchild. Flip: A difference evaluator for alternating images. *Proc. ACM Comput. Graph. Interact. Tech.*, 3(2), aug 2020. **12**
- [3] Yunpeng Bai, Chao Dong, and Cairong Wang. Ps-nerv: Patch-wise stylized neural representations for videos. *arXiv preprint arXiv:2208.03742*, 2022. **3**
- [4] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016. **2**
- [5] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436*, 2018. **2, 11**
- [6] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. **11**
- [7] Linfeng Cao, Aofan Jiang, Wei Li, Huaying Wu, and Nanyang Ye. Oodhdr-codec: Out-of-distribution generalization for hdr image compression. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(1):158–166, Jun. 2022. **2**
- [8] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. Nerv: Neural representations for videos. *Advances in Neural Information Processing Systems*, 34:21557–21568, 2021. **2, 3, 5, 12**
- [9] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. *arXiv preprint arXiv:2012.09161*, 2020. **3**
- [10] Yue Chen, Debargha Murherjee, Jingning Han, Adrian Grange, Yaowu Xu, Zoe Liu, Sarah Parker, Cheng Chen, Hui Su, Urvang Joshi, et al. An overview of core coding tools in the av1 video codec. In *2018 Picture Coding Symposium (PCS)*, pages 41–45. IEEE, 2018. **1**
- [11] Yinbo Chen and Xiaolong Wang. Transformers as meta-learners for implicit neural representations. In *European Conference on Computer Vision*, 2022. **3**
- [12] Emilien Dupont, Adam Goliński, Milad Alizadeh, Yee Whye Teh, and Arnaud Doucet. Coin: Compression with implicit neural representations. *arXiv preprint arXiv:2103.03123*, 2021. **3**
- [13] Emilien Dupont, Hrushikesh Loya, Milad Alizadeh, Adam Goliński, Yee Whye Teh, and Arnaud Doucet. Coin++: Data agnostic neural compression. *arXiv preprint arXiv:2201.12904*, 2022. **3**
- [14] Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Rémi Gribonval, Herve Jegou, and Armand Joulin. Training with quantization noise for extreme model compression. *arXiv preprint arXiv:2004.07320*, 2020. **3**
- [15] Rizal Fathony, Anit Kumar Sahu, Devin Willmott, and J. Zico Kolter. Multiplicative filter networks. In *ICLR*, 2021. **3**
- [16] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. **3**
- [17] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Stabilizing the lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*, 2019. **3**
- [18] Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv preprint arXiv:1902.09574*, 2019. **3**
- [19] Sharath Girish, Kamal Gupta, Saurabh Singh, and Abhinav Shrivastava. Lilnetx: Lightweight networks with extreme model compression and structured sparsification. *ArXiv*, abs/2204.02965, 2022. **2, 3, 4, 11**
- [20] Sharath Girish, Shishira R Maiya, Kamal Gupta, Hao Chen, Larry S Davis, and Abhinav Shrivastava. The lottery ticket hypothesis for object recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 762–771, 2021. **3**
- [21] Andrei Ivanov, Nikoli Dryden, and Torsten Hoefer. Sten: An interface for efficient sparsity in pytorch. **3**
- [22] Jaeho Lee, Jihoon Tack, Namhoon Lee, and Jinwoo Shin. Meta-learning sparse implicit neural representations. In *Advances in Neural Information Processing Systems*, 2021. **3**
- [23] Xiaoqi Li, Jiaming Liu, Shizun Wang, Cheng Lyu, Ming Lu, Yurong Chen, Anbang Yao, Yandong Guo, and Shanghang Zhang. Efficient meta-tuning for content-aware neural video delivery. In *European Conference on Computer Vision*, pages 308–324. Springer, 2022. **2**
- [24] Zhi Li, Anne Aaron, Ioannis Katsavounidis, Anush Moorthy, and Megha Manohara. Toward A Practical Perceptual Video Quality Metric, 2016. **12**
- [25] Zizhang Li, Mengmeng Wang, Huaijin Pi, Kechun Xu, Jianbiao Mei, and Yong Liu. E-nerv: Expedite neural video representation with disentangled spatial-temporal context. *arXiv preprint arXiv:2207.08132*, 2022. **3**
- [26] Xiaofan Lin, Cong Zhao, and Wei Pan. Towards accurate binary convolutional neural network. *Advances in neural information processing systems*, 30, 2017. **3**
- [27] David B Lindell, Dave Van Veen, Jeong Joon Park, and Gordon Wetzstein. Bacon: Band-limited coordinate networks for multiscale scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16252–16262, 2022. **3**

- [28] Julien NP Martel, David B Lindell, Connor Z Lin, Eric R Chan, Marco Monteiro, and Gordon Wetzstein. Acorn: Adaptive coordinate networks for neural scene representation. *arXiv preprint arXiv:2105.02788*, 2021. 3
- [29] Alexandre Mercat, Marko Viitanen, and Jarno Vanne. Uvq dataset: 50/120fps 4k sequences for video codec analysis and development. In *Proceedings of the 11th ACM Multimedia Systems Conference, MMSys '20*, page 297–302, New York, NY, USA, 2020. Association for Computing Machinery. 2, 5
- [30] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 2
- [31] David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors for learned image compression. *Advances in neural information processing systems*, 31, 2018. 2
- [32] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 2
- [33] D. Oktay et al. Scalable model compression by entropy penalized reparameterization. In *ICLR*, 2020. 3, 4, 11
- [34] Daniel Rho, Junwoo Cho, Jong Hwan Ko, and Eunbyung Park. Neural residual flow fields for efficient video representations. *arXiv preprint arXiv:2201.04329*, 2022. 3
- [35] Vishwanath Saragadam, Jasper Tan, Guha Balakrishnan, Richard G Baraniuk, and Ashok Veeraraghavan. Miner: Multiscale implicit neural representations. *arXiv preprint arXiv:2202.03532*, 2022. 3
- [36] Shayan Shekarforoush, David B Lindell, David J Fleet, and Marcus A Brubaker. Residual multiplicative filter networks for multiscale reconstruction. *arXiv preprint arXiv:2206.00746*, 2022. 3
- [37] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020. 2, 3, 5
- [38] Pierre Stock, Armand Joulin, Rémi Gribonval, Benjamin Graham, and Hervé Jégou. And the bit goes down: Revisiting the quantization of neural networks. *arXiv preprint arXiv:1907.05686*, 2019. 3
- [39] Yannick Strümpfer, Janis Postels, Ren Yang, Luc Van Gool, and Federico Tombari. Implicit neural representations for image compression. *arXiv preprint arXiv:2112.04267*, 2021. 3
- [40] Vivienne Sze, Madhukar Budagavi, and Gary J Sullivan. High efficiency video coding (hevc). In *Integrated circuit and systems, algorithms and architectures*, volume 39, page 40. Springer, 2014. 1
- [41] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P. Srinivasan, Jonathan T. Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *CVPR*, 2021. 3
- [42] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 3
- [43] Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár. Lossy image compression with compressive autoencoders. *arXiv preprint arXiv:1703.00395*, 2017. 2
- [44] Frederick Tung and Greg Mori. Clip-q: Deep network compression learning by in-parallel pruning-quantization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7873–7882, 2018. 3
- [45] G.K. Wallace. The jpeg still picture compression standard. *IEEE Transactions on Consumer Electronics*, 38(1):xviii–xxxiv, 1992. 1
- [46] Ian H Witten, Radford M Neal, and John G Cleary. Arithmetic coding for data compression. *Communications of the ACM*, 30(6):520–540, 1987. 4
- [47] Mingtian Zhang, Andi Zhang, and Steven McDonagh. On the out-of-distribution generalization of probabilistic image modelling. *Advances in Neural Information Processing Systems*, 34:3811–3823, 2021. 2
- [48] Yunfan Zhang, Ties van Rozendaal, Johann Brehmer, Markus Nagel, and Taco Cohen. Implicit neural video compression. *arXiv preprint arXiv:2112.11312*, 2021. 3

Appendix

A. Additional Implementation Details

A.1. NeRV

We use the NeRV-L config from the original paper. The model takes in positional embedding of time coordinate as input. We set the number of sine levels to 80 and base value to 1.25. The hidden dimension of the 2-layer MLP in the beginning is set to 1024, with 128 output channels and an $8\times$ channel expansion for the first NeRV Block. We stack 5-NeRV blocks with upscale factors of $\{5, 3, 2, 2, 2\}$ for HD version and an additional block with $2\times$ upsampling for the 4K version. Following the standard training schedule, we use cosine learning rate schedule starting with $5e^{-4}$, with a warmup of 0.2. We use the combination of L1 and SSIM loss and train the model with a batch size of 1, as mentioned in the paper.

A.2. Entropy Loss and Weight Quantization

We follow the works of [19, 33] for performing model compression with small variations. We represent our MLP layer weights as $\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L$ for L layers with the l^{th} layer weight matrix $\mathbf{W}_l \in \mathbb{R}^{O_l \times I_l}$ having a shape $O_l \times I_l$ consisting of continuous values. For each weight matrix, we maintain a corresponding flattened latent quantized representation vector $\widetilde{\mathbf{W}}_l \in \mathbb{Z}^{O_l I_l}$. $\widetilde{\mathbf{W}}_l$ consists of integer values for each corresponding element in \mathbf{W}_l . For ease of notation, we drop the subscript l .

We then maintain decoders $f_\phi(\cdot)$ parameterized by parameters ϕ . The weight matrix \mathbf{W} is then obtained from the quantized weights $\widetilde{\mathbf{W}}$ as $\mathbf{W} = f_\phi(\widetilde{\mathbf{W}})$. Prior works use matrices or vectors to represent the weight decoders while we use a single scalar ϕ as a parameter of the decoder. The weight matrix is thus simply, $\mathbf{W} = \text{reshape}(\phi \widetilde{\mathbf{W}})$, where the scale parameter ϕ is multiplied with individual values of $\widetilde{\mathbf{W}}$. We make the scale parameter learnable by passing gradients. Thus, it effectively controls the bit width of the weight matrix. We maintain separate decoders for each layer of the MLP.

To make the network fully differentiable, we maintain continuous surrogates $\widehat{\mathbf{W}}$ for the discrete latents $\widetilde{\mathbf{W}}$. During training, we simply round the surrogates to their nearest integer to obtain the discrete latents which are then passed to the decoder. We make the rounding operation differentiable using the straight-through estimator [6] to pass the gradients from $\widehat{\mathbf{W}}$ to $\widetilde{\mathbf{W}}$.

To reduce the entropy of the quantized latents, we use probability models from [5]. For each continuous surrogate $\widehat{\mathbf{W}}$, we maintain probability models $c_\theta(\cdot)$ parameterized by θ which output the CDF of the latent distributions. Similar to prior works, we use uniform noise $n \sim \mathcal{U}(-\frac{1}{2}, \frac{1}{2})$ as a

Video	Frames	Link
Mario Kart	4000	http://bit.ly/3XjIvfR
Dota	4261	http://bit.ly/3Xf6Nru
Ride	4000	http://bit.ly/3TOEgWI
Submarine	3626	http://bit.ly/3EJKzGM
Water Scooter	4199	http://bit.ly/3OhF99h
Mortal Kombat	3239	http://bit.ly/3EdvNGU

Table 3. Details of videos from the Youtube-8M dataset.

substitute for quantization. The entropy of the model can now be minimized by minimizing the self-information I as follows:

$$I(\widehat{\mathbf{W}}) = -\log_2(c_\theta(\widehat{\mathbf{W}} + n)). \quad (3)$$

This serves as the entropy regularization loss which controls the rate-distortion tradeoff. A higher entropy coefficient λ_I leads to more compressed latents (lower rate), but usually at the cost of PSNR (higher distortion). The network latents, decoder parameter, and probability model parameters are learnable and jointly optimized. Following prior works, we use a learning rate of $1e-4$ for the probability model weights and the same learning rate for the decoder weights. We set the learning rate of the latents to $5e-4$. All the parameters are optimized in an end-to-end manner with an Adam optimizer during training, thus requiring no post-hoc approaches. After training, we discard the probability models and use the frequency of each quantized value in the latent vector to obtain the probability tables required for arithmetic entropy coding. Note that the continuous surrogates are discarded and only their rounded discrete latents are stored using entropy coding. These latents can then be decoded using the probability tables. The decoder parameters and probability tables have almost no overhead compared to the overall model latents.

B. Additional Dataset Details

B.1. UVG-4K

In addition to the datasets shown in Section 4 of the main paper, we show quantitative and qualitative results on 7 more videos from the UVG dataset at the 4K resolution: Twilight, Sunbath, CityAlley, FlowerFocus, FlowerKids, RiverBank, and RaceNight. We dub this dataset UVG-4K (Set 2).

B.2. Youtube-8M

We select 5 more videos from the Youtube-8M dataset, with varying video content to further test the ability of our model to encode longer videos. This is an extension of the experiments from Section 4.4 which consists of a single video (Mario Kart). We present the details of each video used in Table 3.

Video Name	NIRVANA (Ours)				NeRV			
	PSNR \uparrow	VMAF \uparrow	FLIP \downarrow	BPP \downarrow	PSNR \uparrow	VMAF \uparrow	FLIP \downarrow	BPP \downarrow
ReadySteadyGo	35.43	98.04	0.0862	1.26	34.59	96.95	0.0862	0.81
Bosphorus	40.53	90.97	0.0549	0.68	39.11	88.26	0.0621	0.81
Beauty	35.77	84.46	0.0524	0.96	34.57	86.35	0.0605	0.81
Honeybee	38.83	91.31	0.0505	0.51	39.71	95.08	0.0419	0.81
Jockey	37.56	93.07	0.0710	0.96	38.16	95.76	0.0653	0.81
Yachtride	37.94	91.31	0.0668	1.03	35.68	88.70	0.0786	0.81
ShakeNDry	37.82	88.81	0.0608	0.76	39.68	95.20	0.0468	1.61
Average	37.70	91.14	0.0632	0.86	37.35	92.33	0.0631	0.92

Table 4. **Video-wise performance on UVG-HD:** We show video-wise results of the 7 videos in UVG-HD and compare the reconstruction quality using the 3 metrics of PSNR, VMAF, FLIP along with compression rate measured by BPP. We see that we maintain similar performance as NeRV in all 3 metrics and BPP while having $12\times$ faster encoding speed (as also shown in Table 1 of the main paper).

Video Name	NIRVANA (Ours)		NeRV	
	PSNR	BPP	PSNR	BPP
ReadySteadyGo	33.85	0.41	33.22	0.24
Bosphorus	38.71	0.21	39.0	0.24
Beauty	31.96	0.28	31.05	0.24
Honeybee	35.64	0.14	36.36	0.24
Jockey	35.05	0.30	35.9	0.24
Yachtride	36.33	0.33	35.05	0.24
ShakeNDry	34.78	0.24	36.09	0.49
Average	35.18	0.27	35.23	0.28

(a) UVG-4K (Set 1)

Video Name	NIRVANA (Ours)		NeRV	
	PSNR	BPP	PSNR	BPP
FlowerFocus	36.50	0.12	37.08	0.24
CityAlley	37.43	0.17	38.39	0.24
Twilight	38.02	0.13	20.99	0.24
FlowerKids	34.62	0.26	33.77	0.24
RiverBank	33.83	0.26	32.35	0.24
RaceNight	32.72	0.27	32.92	0.24
Sunbath	37.75	0.24	44.17	0.49
Average	35.84	0.21	34.23	0.28

(b) UVG-4K (Set 2)

Table 5. **Video-wise comparison on different sets of UVG-4K:** We show video-wise results on 2 different sets of 7 UVG videos at 4K resolution. Set-1 consists of videos from UVG-HD at 4K resolution while Set 2 consists of additional 7 videos from the dataset. We maintain similar performance in terms of PSNR and BPP as NeRV while also being $\sim 6\times$ faster for both sets and being $6\times$ faster in terms of encoding time.

C. Video-wise comparison

We show additional quantitative results on UVG-4K (Set 2) and Youtube-8M mentioned above.

C.1. UVG-HD

We provide video-wise results of our approach along with that of NeRV [8]. We evaluate the video on the additional perceptual quality metrics of FLIP [2] and VMAF [24] as well, along with the standard PSNR. Results are summarized in Table 4. We see that we continue to obtain similar performance compared to NeRV in terms of these metrics while being $\sim 12\times$ faster. Also note the adaptive BPP of our method, which is based on the amount of motion in each video. In contrast, NeRV maintains a fixed BPP due to fixed model size (ShakeNDry shows twice the BPP due to half the number of frames). We observe a small drop in VMAF ($92.33 \rightarrow 91.14$) while maintaining similar value of FLIP (~ 0.0632) compared to NeRV.

C.2. UVG-4K

We provide video-wise results on the 2 sets of UVG at 4K resolution. For Set 1, we obtain comparable performance to NeRV while obtaining $\sim 6\times$ faster encoding speed. Similar to UVG-HD, we continue to show the benefits of adaptive compression, with static videos such as Honeybee showing lower levels of BPP (0.14) compared to the most dynamic video, ReadySteadyGo (0.41 BPP). For Set 2, we outperform NeRV by 1.5 PSNR while still obtaining 25% lower BPP $0.28 \rightarrow 0.21$. The PSNR drop of NeRV on the Twilight video is largely due to quantization at the fixed bit width of 20. Hand-tuning is necessary in order to maintain higher PSNR but at the cost of BPP. In contrast, our approach maintains the reconstruction quality for a variety of videos and adaptively quantizes for each video.

C.3. Youtube-8M

We now provide video-wise results of 5 videos picked from the Youtube-8M dataset at 1080p resolution. Details

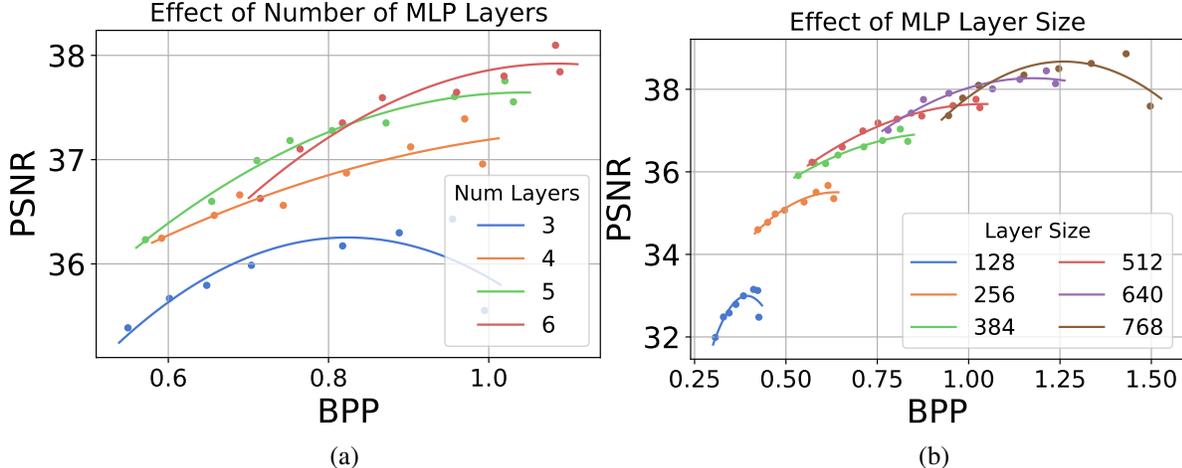


Figure 7. Increasing number of layers improves the PSNR/BPP curve upto 5 layers in the lower BPP regime (<0.8). Increasing layer size shifts the PSNR/BPP curve upwards and to the right as representation capacity increases along with more parameters.

Video Name	NIRVANA (Ours)		NeRV	
	PSNR	BPP	PSNR	BPP
Dota	38.03	0.62	35.53	0.34
Ride	36.65	1.09	29.74	0.36
Submarine	38.48	0.69	33.64	0.40
Water Scooter	37.79	0.84	30.46	0.34
Mortal Kombat	36.02	1.03	31.36	0.45
Average	37.39	0.85	32.14	0.38

Table 6. **Results on Youtube-8M videos with long duration.** We provide video-wise results on 5 videos picked from the Youtube-8M datasets with approximately 4000 frames compared to the typical 600 from UVG. Still we maintain PSNR/BPP with no change in hyperparameters, whereas NeRV shows large degradation in performance for the same network and similar encoding times.

of the videos are provided in Table 3. Results are summarized in Table 6. We see that our reconstruction quality does not degrade with longer videos. This behavior is different from NeRV, which obtains a significant drop in PSNR (about -5.2). This is in line with the observations in Section 4.4 of the main paper, where we see that increasing number of frames results in drop of NeRV’s reconstruction quality while we maintain similar levels of performance.

D. Additional Ablations

In addition to the ablations shown in Fig. 5, we analyze the effect of layer size and the number of layers of the MLP in our network when evaluating on the Jockey video of the UVG-HD dataset. Note that the default values of layer size is 512 and the number of layers is 5.

D.1. Effect of Layer Number

We increase the number of layers from 3 to 6, while keeping other parameters at their default values and varying the entropy coefficient λ_I for each curve as in Section 5.

Results are summarized in Figure 7(a). We see that increasing the number of layers from 3 to 5 improves the tradeoff curve (shifts upwards) in the low BPP regime (<0.8). However, increasing it further shifts the curve upwards and to the right. This might be because the MLP network requires higher levels of non-linearity to learn a global representation for a group of 3 frames which typically contain significant motion in the case of Jockey. However, for 6 layers the network shifts the curve upwards and to the right, and we no longer obtain increase in PSNR at no cost of BPP.

D.2. Effect of Layer Size

We vary the layer size from 128 to 768 progressively, in steps of 128 for each of the 5 layers. Results are summarized in Figure 7(b). We see that increasing the layer size simply shifts the curve upwards and to the right, which is expected as a higher number of parameters leads to more representation capability of the network at the cost of more parameters. While increasing the number of layers increases number of parameters as well, a similar tradeoff is not present in that case up to a certain level, suggesting that a minimum number of non-linearities/activation functions are important to achieve the optimal tradeoff.

E. Qualitative Results

We qualitatively visualize the reconstruction results for 3 videos from UVG-4K (Set 2) in Fig. 8. We obtain higher-quality, more faithful reconstructions while preserving details at similar/lower BPP compared to NeRV; *e.g.*, Twilight ($0.24 \rightarrow 0.13$), RiverBank ($0.24 \rightarrow 0.26$), CityAlley ($0.24 \rightarrow 0.17$). Notice the bird which is reconstructed by our approach in Twilight (top), or finer details of the branches in RiverBank (middle), or maintaining the right color information of the door and the people’s shirts in CityAlley (bottom).



Figure 8. **Qualitative results from UVG-4K Set-2:** (Left) Ground truth video frames. (Center) Reconstruction from NIRVANA. (Right) Reconstruction from NeRV. Top to bottom: We show additional examples where NIRVANA is able to preserve the image fidelity after reconstruction, such as the bird in Twilight, the tree in RiverBank, and the humans in CityAlley.