# Video Compression through Image Interpolation

Chao-Yuan Wu, Nayan Singhal, Philipp Krähenbühl

The University of Texas at Austin
{cywu, nayans, philkr}@cs.utexas.edu

**Abstract.** An ever increasing amount of our digital communication, media consumption, and content creation revolves around videos. We share, watch, and archive many aspects of our lives through them, all of which are powered by strong video compression. Traditional video compression is laboriously hand designed and hand optimized. This paper presents an alternative in an end-to-end deep learning codec. Our codec builds on one simple idea: Video compression is repeated image interpolation. It thus benefits from recent advances in deep image interpolation and generation. Our deep video codec outperforms today's prevailing codecs, such as H.261, MPEG-4 Part 2, and performs on par with H.264.

## 1    Introduction

Video commands the lion's share of internet data, and today makes up three-fourths of all internet traffic [18]. We capture moments, share memories, and entertain one another through moving pictures, all of which are powered by ever powerful digital camera and video compression. Strong compression significantly reduces internet traffic, saves storage space, and increases throughput. It drives applications like cloud gaming, real-time high-quality video streaming [21], or 3D and 360-videos. Video compression even helps better understand and parse videos using deep neural networks [32]. Despite these obvious benefits, video compression algorithms are still largely hand designed. The most competitive video codecs today rely on a sophisticated interplay between block motion estimation, residual color patterns, and their encoding using discrete cosine transform and entropy coding [24]. While each part is carefully designed to compress the video as much as possible, the overall system is not jointly optimized, and has largely been untouched by end-to-end deep learning.

This paper presents, to the best of our knowledge, the first end-to-end trained deep video codec. The main insight of our codec is a different view on video compression: We frame video compression as repeated image interpolation, and draw on recent advances in deep image generation and interpolation. We first encode a series of anchor frames (key frames), using standard deep image compression. Our codec then reconstructs all remaining frames by interpolating between neighboring anchor frames. However, this image interpolation is not unique. We additionally provide a small and compressible code to the interpolation network

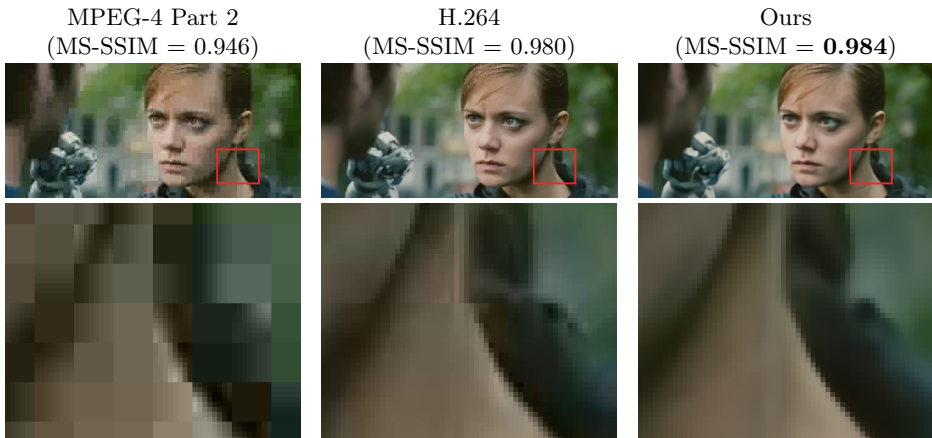| MPEG-4 Part 2 (MS-SSIM = 0.946) | H.264 (MS-SSIM = 0.980) | Ours (MS-SSIM = **0.984**) |
|---|---|---|



Fig. 1: Comparison of our end-to-end deep video compression algorithm to MPEG-4 Part 2 and H.264 on the Blender *Tears of Steel* movie. All methods use 0.080 BPP. Our model offers a visual quality better than MPEG-4 Part 2 and comparable to H.264. Unlike traditional methods, our method is free of block artifacts. The MS-SSIM [29] measures the image quality of the video clip compared to the raw uncompressed ground truth. (Best viewed on screen.)

to disambiguate different interpolations, and encode the original video frame as faithfully as possible. The main technical challenge is the design of a compressible image interpolation network.

We present a series of increasingly powerful and compressible encoder-decoder architectures for image interpolation. We start by using a vanilla U-net interpolation architecture [23] for reconstructing frames other than the key frames. This architecture makes good use of repeating static patterns through time, but it struggles to properly disambiguate the trajectories for moving patterns. We then directly incorporate an offline motion estimate from either block-motion estimation or optical flow into the network. The new architecture interpolates spatial U-net features using the pre-computed motion estimate, and improves compression rates by an order of magnitude over deep image compression. This model captures most, but not all of the information we need to reconstruct a frame. We additionally train an encoder that extracts the content not present in either of the source images, and represents it compactly. Finally, we reduce any remaining spatial redundancy, and compress them using a 3D PixelCNN [20] with adaptive arithmetic coding [31].

To further reduce bitrate, our video codec applies image interpolation in a hierarchical manner. Each consecutive level in the hierarchy interpolates between ever closer reference frames, and is hence more compressible. Each level in the hierarchy uses all previously decompressed images.

We compare our video compression algorithm to state-of-the-art video compression (HEVC, H.264, MPEG-4 Part 2, H.261), and various image interpolation baselines. We evaluate all algorithms on two standard datasets of uncompressed video: Video Trace Library (VTL) [2] and Ultra Video Group (UVG) [1].

We additionally collect a subset of the Kinetics dataset [7] for both training and testing. The Kinetics subset contains high resolution videos, which we downsample to remove compression artifacts introduced by prior codecs on YouTube. The final dataset contains 2.8M frames. Our deep video codec outperforms all deep learning baselines, MPEG-4 Part 2, and H.261 in both compression rate and visual quality measured by MS-SSIM [29] and PSNR. We are on par with the state-of-the-art H.264 codec. Figure 1 shows a visual comparison. All the data is publicly available, and we will publish our code upon acceptance.

## 2   Related Work

Video compression algorithms must specify an encoder for compressing the video, and a decoder for reconstructing the original video. The encoder and the decoder together constitute a codec. A codec has one primary goal: Encode a series of images in the fewest number of bits possible. Most compression algorithms find a delicate trade-off between compression rate and reconstruction error. The simplest codecs, such as motion JPEG or GIF, encode each frame independently, and heavily rely on image compression.

**Image compression.** For images, deep networks yield state-of-the-art compression ratios with impressive reconstruction quality [6, 12, 22, 25, 26]. Most of them train an autoencoder with a small binary bottleneck layer to directly minimize distortion [12, 22, 26]. A popular variant progressively encodes the image using a recurrent neural network [5, 12, 26]. This allows for variable compression rates with a single model. We extend this idea to variable rate video compression.

Deep image compression algorithms use fully convolutional networks to handle arbitrary image sizes. However, the bottleneck in fully convolutional networks still contains spatially redundant activations. Entropy coding further compresses this redundant information [6, 17, 22, 25, 26]. We follow Mentzer *et al.* [17] and use adaptive arithmetic coding on probability estimates of a Pixel-CNN [20].

Learning the binary representation is inherently non-differentiable, which complicates gradient based learning. Toderici *et al.* [26] use stochastic binarization and backpropagate the derivative of the expectation. Agustsson *et al.* [4] use soft assignment to approximate quantization. Balle *et al.* [6] replace the quantization by adding uniform noise. All of these methods work similarly and allow for gradients to flow through the discretization. In this paper, we use stochastic binarization [26].

Combining this bag of techniques, deep image compression algorithms offer a better compression rate than hand-designed algorithms, such as JPEG or WebP [3], at the same level of image quality [22]. Deep image compression algorithms heavily exploit the spatial structure of an image. However, they miss out on a crucial signal in videos: time. Videos are temporally highly redundant. No deep image compression can compete with state-of-the-art (shallow) video compression, which exploits this redundancy.

**Video compression.** Hand-designed video compression algorithms, such as H.263, H.264 or HEVC (H.265) [14] build on two simple ideas: They decompose each frame into blocks of pixels, known as macroblocks, and they divide frames into image (I) frames and referencing (P or B) frames. I-frames directly compress video frames using image compression. Most of the savings in video codecs come from referencing frames. P-frames borrow color values from preceding frames. They store a motion estimate and a highly compressible difference image for each macroblock. B-frames additionally allow bidirectional referencing, as long as there are no circular references. Both H.264 and HEVC encode a video in a hierarchical way. I-frames form the top of the hierarchy. In each consecutive level, P- or B-frames reference decoded frames at higher levels. The main disadvantages of traditional video compression is the intensive engineering efforts required and the difficulties in joint optimization. In this work, we build a hierarchical video codec using deep neural networks. We train it end-to-end without any hand-engineered heuristics or filters. Our key insight is that referencing (P or B) frames are a special case of image interpolation.

Learning-based video compression is largely unexplored, in part due to difficulties in modeling temporal redundancy. Tsai *et al.* propose a deep post-processing filter encoding errors of H.264 in domain specific videos [27]. However, it is unclear if and how the filter generalizes in an open domain. To the best of our knowledge, this paper proposes the first general deep network for video compression.

**Image interpolation and extrapolation.** Image interpolation seeks to hallucinate an unseen frame between two reference frames. Most image interpolation networks build on an encoder-decoder network architecture to move pixels through time [10, 11, 15, 19]. Jia *et al.* [10] and Niklaus *et al.* [19] estimate a spatially-varying convolution kernel. Liu *et al.* [15] produce a flow field. All three methods then combine two predictions, forward and backward in time, to form the final output.

Image extrapolation is more ambitious and predicts a future video from a few frames [16], or a still image [28, 33]. Both image interpolation and extrapolation works well for small timesteps, e.g. for creating slow-motion video [11] or predicting a fraction of a second into the future. However, current methods struggle for larger timesteps, where the interpolation or extrapolation is no longer unique, and additional side information is required. In this work, we extend image interpolation and incorporate few compressible bits of side information to reconstruct the original video.

## 3    Preliminary

Let $I^{(t)} \in \mathbb{R}^{W \times H \times 3}$ denote a series of frames for $t \in \{0, 1, \ldots\}$. Our goal is to compress each frame $I^{(t)}$ into a binary code $b^{(t)} \in \{0, 1\}^{N_t}$. An encoder $E : \{I^{(0)}, I^{(1)}, \ldots\} \to \{b^{(0)}, b^{(1)}, \ldots\}$ and decoder $D : \{b^{(0)}, b^{(1)}, \ldots\} \to \{\hat{I}^{(0)}, \hat{I}^{(1)}, \ldots\}$ compress and decompress the video respectively. $E$ and $D$ have two competing aims: Minimize the total bitrate $\sum_t N_t$, and reconstruct the original video

as faithfully as possible. We measure the reconstruction error with an L1 loss $\ell(\hat{I}, I) = \|\hat{I} - I\|_1$.

**Image compression.** The simplest encoders and decoders process each image independently: $E_I : I^{(t)} \rightarrow b^{(t)}$, $D_I : b^{(t)} \rightarrow \hat{I}^{(t)}$. Here, we build on the model of Toderici *et al.* [26], which encodes and reconstructs an image progressively over $K$ iterations. At each iteration, the model encodes a residual $r_k$ between the previously coded image and the original frame:

$$r_0 := I$$
$$b_k := E_I\left(r_{k-1}, g_{k-1}\right), \qquad r_k := r_{k-1} - D_I\left(b_k, h_{k-1}\right), \qquad \text{for } k = 1, 2, \ldots$$

where $g_k$ and $h_k$ are latent Conv-LSTM states that are updated at each iteration. All iterations share the same network architecture and parameters forming a recurrent structure. The training objective minimizes the distortion at all the steps $\sum_{k=1}^{K} \|r_k\|_1$. The reconstruction $\hat{I}_K = \sum_{k=1}^{K} D_I(b_k)$ allows for a variable bitrate encoding depending on the choice of $K$.

Both the encoder and the decoder consist of 4 Conv-LSTMs with stride 2. The bottleneck consists of a binary feature map with $L$ channels and 16 times smaller spatial resolution in both width and height. Toderici *et al.* use a stochastic binarization to allow a gradient signal through the bottleneck. Mathematically, this reduces to *REINFORCE* [30] on sigmoidal activations. At inference time, the most likely state is selected.

This architecture yields state-of-the-art image compression performance. However, it fails to exploit any temporal redundancy.

**Video compression.** Modern video codecs process I-frames using an image encoder $E_I$ and decoder $D_I$. P-frames store a block motion estimate $\mathcal{T} \in \mathbb{R}^{W \times H \times 2}$, similar to an optical flow field, and a residual image $\mathcal{R}$, capturing the appearance changes not explained by motion. Both motion estimate and residual are jointly compressed using entropy coding. The original color frame is then recovered by
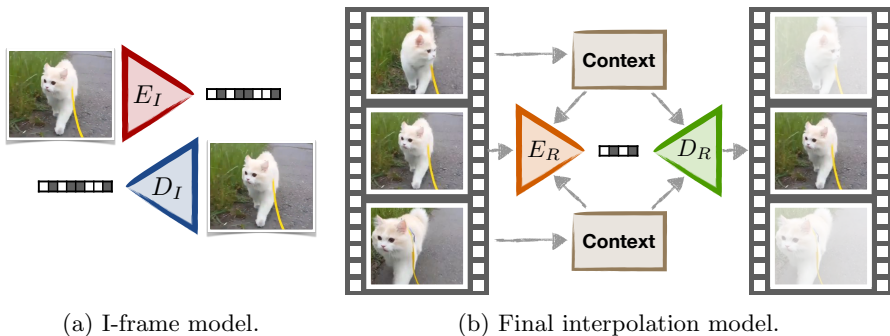
$$I_i^{(t)} = I_{i-\mathcal{T}_i^{(t)}}^{(t-1)} + \mathcal{R}_i^{(t)}, \tag{1}$$

for every pixel $i$ in the image. The compression is uniquely defined by a block structure and motion estimate $\mathcal{T}$. The residual is simply the difference between the motion interpolated image and the original.

In this paper, we propose a more general view on video compression through image interpolation. We augment image interpolation network with motion information and add a compressible bottleneck layer.

## 4   Video Compression through Interpolation

Our codec first encodes I-frames using the compression algorithm of Toderici *et al.*, see Figure 2a. We chose every $n$-th frame as an I-frame. The remaining $n - 1$ frames are interpolated. We call those frames R-frames, as they reference other frames. We choose $n = 12$ in practice, but also experimented with larger groups of pictures. We will first discuss our basic interpolation network, and then show a hierarchical interpolation setup, that further reduces the bitrate.

(a) I-frame model.                    (b) Final interpolation model.

Fig. 2: Our model is composed of an image compression model that compresses the key frames, and a conditional interpolation model that interpolates the remaining frames.

### 4.1  Interpolation network

In the simplest version of our codec, all R-frames use a blind interpolation network to interpolate between two key-frames $I_1$ and $I_2$. Specifically, we train a context network $C : I \rightarrow \{f^{(1)}, f^{(2)}, \ldots\}$ to extract a series of feature maps $f^{(l)}$ of various spatial resolutions. For notational simplicity let $f := \{f^{(1)}, f^{(2)}, \ldots\}$ be the collection of all context features. In our implementation, we use the upconvolutional feature maps of a U-net architecture with increasing spatial resolution $\frac{W}{8} \times \frac{H}{8}$, $\frac{W}{4} \times \frac{H}{4}$, $\frac{W}{2} \times \frac{H}{2}$, $W \times H$, in addition to the original image.

We extract context features $f_1$ and $f_2$ for key-frames $I_1$ and $I_2$ respectively, and train a network $D$ to interpolate the frame $\hat{I} := D(f_1, f_2)$. $C$ and $D$ are trained jointly. This simple model favors a high compression rate over image quality, as none of the R-frames capture any information not present in the I-frames.

Without any further information, it is impossible for the network to faithfully reconstruct a frame. What can we provide to the network to make interpolation easier?

**Motion compensated interpolation.** A great candidate is ground truth motion information. It defines where pixels move through time and greatly disambiguates interpolation. We tried both optical flow [8] and block motion estimation [21]. Block motion estimates are easier to compress, but optical flow retains finer details.

We use the motion information to warp each context feature map

$$\tilde{f}_i^{(l)} = f_{i-\mathcal{T}_i}^{(l)}, \tag{2}$$

at every spatial location $i$. We scale the motion estimation with the resolution of the feature map, and use bilinear interpolation for fractional locations. The decoder now uses the warped context features $\tilde{f}$ instead, which allows it to focus solely on image creation, and ignore motion estimation.

Motion compensation greatly improves the interpolation network, as we will show in Section 5. However, it still only produces content seen in either reference image. Variations beyond motion, such as change in lighting, deformation, occlusion, etc. are not captured by this model.

Our goal is to encode the remaining information in a highly compact from.

**Residual motion compensated interpolation.** Our final interpolation model combines motion compensated interpolation with a compressed residual information, capturing both the motion and appearance difference in the interpolated frames. Figure 2b show an overview of the model.

We jointly train an encoder $E_R$, context model $C$ and interpolation network $D_R$. The encoder sees the same information as the interpolation network, which allows it to compress just the missing information, and avoid a redundant encoding. Formally, we follow the progressive compression framework of Toderici *et al.* [26], and train a variable bitrate encoder and decoder conditioned on the warped context $\tilde{f}$:

$$r_0 := I$$
$$b_k := E_R(r_{k-1}, \tilde{f}_1, \tilde{f}_2, g_{k-1}), \quad r_k := r_{k-1} - D_R(b_k, \tilde{f}_1, \tilde{f}_2, h_{k-1}), \quad \text{for } k = 1, 2, \ldots$$

This framework allows the encoder to learn a variable rate compression at high reconstruction quality. The interpolation network generally requires fewer bits to encode temporally close images and more bits for images that are farther apart. In one extreme, when key frames do not provide any meaningful signal to the interpolated frame, our residual motion compensated interpolation reduces to image compression. In the other extreme, when the image content does not change, our algorithm reduces to a vanilla interpolation, and requires close to zero bits.

In the next section, we use this to our advantage, and design a hierarchical interpolation scheme, maximizing the number of temporally close interpolations.

## 4.2 Hierarchical interpolation

The basic idea of hierarchical interpolation is simple: We interpolate some frames first, and use them as key-frames for the next level of interpolations. See Figure 3
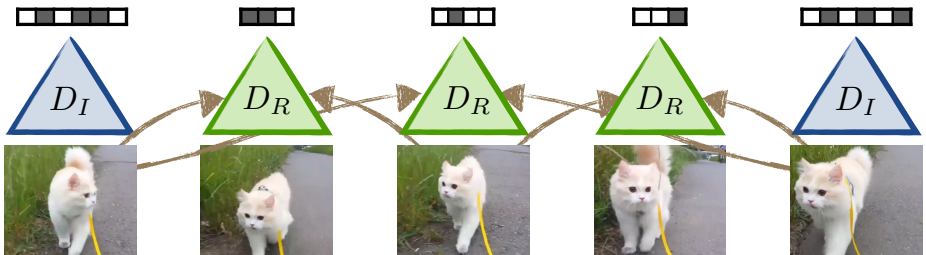


Fig. 3: We apply the interpolation in a hierarchical manner. Each level in the hierarchy uses previously decompressed images.

for example. Each interpolation model $\mathcal{M}_{a,b}$ references $a$ frames into the past and $b$ frames into the future. There are a few things we need to trade off. First, every level in our hierarchical interpolation compounds error. The shallower the hierarchy, the fewer errors compound. In practice, the error propagation for more than three levels in the hierarchy significantly reduces the performance of our codec. Second, we need to train a different interpolation network $\mathcal{M}_{a,b}$ for each temporal offset $(a, b)$, as different interpolations behave differently. To maximally use each trained model, we repeat the same temporal offsets as often as possible. Third, we need to minimize the sum of temporal offsets used in interpolation. The compression rate directly relates to the temporal offset, hence minimizing the temporal offset reduces the bitrate.

Considering just the bitrate and the number of interpolation networks, the optimal hierarchy is a binary tree cutting the interpolation range in two at each level. However, this cannot interpolate more than $n = 2^3 = 8$ consecutive frames, without significant error propagation. We extend this binary structure to $n = 12$ frames, by interpolating at a spacing of three frames in the last level of the hierarchy. For a sequence of four frames $I_1, \ldots, I_4$, we train an interpolation model $\mathcal{M}_{1,2}$ that predicts frame $I_2$, given $I_1$ and $I_4$. We use the exact same model $\mathcal{M}_{1,2}$ to predict $I_3$, but flip the conditioned images $I_4$ and $I_1$. This yields an equivalent model $\mathcal{M}_{2,1}$ predicting the third instead of the second image in the series. Combining this with an interpolation model $\mathcal{M}_{3,3}$ and $\mathcal{M}_{6,6}$ in a hierarchy, we extend the interpolation range from $n = 8$ frames to $n = 12$ frames while keeping the same number of models and levels. We tried applying the same trick to all levels in the hierarchy, extending the interpolation to $n = 27$ frames, but performance dropped, as we had more distant interpolations.

To apply this to a full video of $N$ frames, we divide them into $\lceil N/n \rceil$ groups of pictures (GOPs). Two consecutive groups share the same boundary I-frame. We apply our hierarchical interpolation to each group independently.

**Bitrate optimization.** Each interpolation model at a level $l$ of the hierarchy, can choose to spend $K_l$ bits to encode an image. Our goal is to minimize the overall bitrate, while maintaining a low distortion for all encoded frames. The challenge here is that each selection of $K_l$ affects all lower levels, as errors propagate. Selecting a globally optimal set of $\{K_l\}$ thus requires iterating through all possible combinations, which is infeasible in practice.

We instead propose a heuristic bitrate selection based on beam search. For each level we chose from $m$ different bitrates. We start by enumerating all $m$ possibilities for the I-frame model. Next, we expand the first interpolation model with all $m$ possible bitrates, leading to $m^2$ combinations. Out of these combinations, not all lead to a good MS-SSIM per bitrate, and we discard combinations not on the envelope of the MS-SSIM vs bitrate curve. In practice, only $O(m)$ combinations remain. We repeat this procedure for all levels of the hierarchy. This reduces the search space from $m^L$ to $O(Lm^2)$ for an $L$-level hierarchy. In practice, this yields sufficiently good bitrates.

### 4.3   Implementation

**Architecture.** Our encoder and decoder (interpolation network) architecture follows the image compression model in Toderici *et al.* [26]. While Toderici *et al.* use $L = 32$ latent bits to compress an image, we found that for interpolation, $L = 8$ bits suffice for distance 3 and $L = 16$ for distance 6 and 12. This yields a bitrate of 0.0625 bits per pixel (BPP) and 0.03125 BPP for each iteration respectively.

   We use the original U-net [23] as the context model. To speed-up training and save memory, we reduce the number of channels of all filters by half. We did not observe any significant performance degradation.

   To make it compatible with our architecture, we remove the final output layer and takes the feature maps at the resolutions that are $2\times$, $4\times$, $8\times$ smaller than the original input image.

**Conditional encoder and decoder.** To add the information of the context frames into the encoder and decoder, we fuse the U-net features with the individual Conv-LSTM layers. Specifically, we perform the fusion before each Conv-LSTM layer by concatenating the corresponding U-net features of the same spatial resolution. To increase the computational efficiency, we selectively turn some of the conditioning off in both encoder and decoder. This was tuned for each interpolation network; see supplementary material for details.

   To help the model compare context frames and the target frame side-by-side, we additionally stack the two context frames with target frame, resulting in a 9-channel image, and use that instead as the encoder input.

**Entropy coding.** Since the model is fully-convolutional, it uses the same number of bits for all locations of an image. This disregards the fact that information is not distributed uniformly in an image. Following Mentzer *et al.* [17], we train a 3D Pixel-CNN on the $\{0,1\}^{W/16 \times H/16 \times L}$ binary representations to obtain the probability of each bit sequentially. We then use this probability with adaptive arithmetic coding to encode the feature map. See supplementary material for more details.

**Motion compression.** We store forward and backward block motion estimates as a lossless 4-channel WebP [3] image. For optical flow we train a separate lossy deep compression model, as lossless WebP was unable to compress the flow field.

## 5   Experiments

In this section, we perform a detailed analysis on the series of interpolation models (Section 5.1), and present both quantitative and qualitative (Section 5.2) evaluation of our approach.
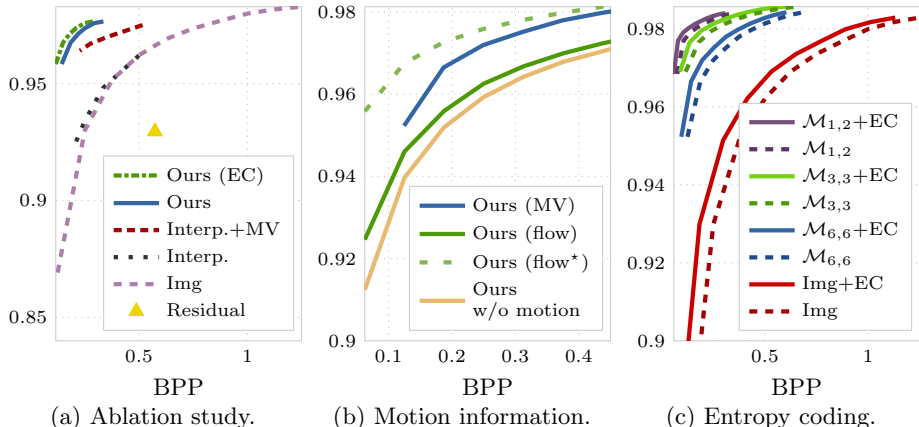
Fig. 4: MS-SSIM of different models evaluated on the VTL dataset.

**Datasets and Protocol.** We train our models using videos from the Kinetics dataset [7]. We only use videos with a width and height greater than 720px. To remove artifacts induced by previous compression, we downsample those high resolution videos to $352 \times 288$px. We allow the aspect ratio to change. The resulting dataset contains 2.8M frames in 75K videos. We train on 65K, use 5K for validation, and 5k for testing. For faster testing on Kinetics, we only use a single group of $n = 12$ pictures per video.

We additionally test our method on two raw video datasets, Video Trace Library (VTL) [2] and Ultra Video Group (UVG) [1]. The VTL dataset contains $\sim 40$K frames of resolution $352 \times 288$ in 20 videos. The UVG dataset contains $3,900$ frames of resolution $1920 \times 1080$ in 7 videos.

We evaluate our method based on the compression rate in bits per pixel (BPP), and the quality of compression in multi-scale structural similarity (MS-SSIM) [29] and peak signal-to-noise ratio (PSNR). We report the average performance of all videos, as opposed to the average of all frames, as our final performance. We use a GOP size of $n = 12$ frames, for all algorithms unless otherwise stated.

**Training Details.** All of our models are trained from scratch for 200K iterations using ADAM [13], with gradient norms clipped at 0.5. We use a batch size of 32 and a learning rate of 0.0005, which is divided by 2 when the validation MS-SSIM plateaus. We augment the data through horizontal flipping. For image models we train on $96 \times 96$ random crops, and for the interpolation models we train on $64 \times 64$ random crops. We train all models with 10 reconstruction iterations.

## 5.1   Ablation study

We first evaluate the series of image interpolation models in Section 4 on the VTL dataset. Figure 4a shows the results.

We can see that image compression model requires by far the highest BPP to achieve high visual quality and performs poorly in the low bitrate region. This is not surprising as it does not exploit any temporal redundancy and needs to encode everything from scratch. Vanilla interpolation does not work much better. We present results for interpolation from 1 to 4 frames, using the best image compression model. While it exploits the temporal redundancy, it fails to accurately reconstruct the image.

Motion-compensated interpolation works significantly better. The additional motion information disambiguates the interpolation, improving the accuracy. The presented BPP includes the size of motion vectors.

Our final model efficiently encodes residual information and makes good use of hierarchical referencing. It achieves the best performance when combined with entropy coding. Note the large performance gap between our method and the image compression model in the low bitrate regime – our model effectively uses context information and achieves a good performance even with very few bits per pixel.

As a sanity check, we further implemented a simple deep codec that uses image compression to encode the residual $\mathcal{R}$ in traditional codecs. This simple baseline stores the video as the encoded residuals, compressed motion vectors, in addition to key frames compressed by a separate deep image compression model. The residual model struggles to learn patterns from noisy residual images, and works worse than an image-only compression model. This suggests that trivially extending deep image compression to videos is not sufficient. Our end-to-end interpolation network performs considerably better.

**Motion.** Next, we analyze different motion estimation models, and compare optical flow to block motion vectors. For optical flow, we use the OpenCV implementation of the Farnebäck's algorithm [8]. For motion compensation, we use the same algorithm as H.264.

Figure 4b shows the results of the $\mathcal{M}_{6,6}$ model with both motion sources. Using motion information clearly helps improve the performance of the model, despite the overhead of motion compression. Block motion estimation (MV) works significantly better than the optical flow based model (flow). Almost all of this performance gain comes from better compressible motion information. The block motion estimates are smaller, easier to compress, and fit in a lossless compression scheme.

To understand whether the worse performance of optical flow is due to the errors in flow compression or the property of the flow itself, we further measure the *hypothetical* performance upper bound of an optical flow based model *assuming* a lossless flow compression at no additional cost (flow$^\star$). As shown in Figure 4b, this upper bound performs better than motion vectors, leaving room for improvement through compressible optical flow estimation. However, finding such a compressible flow estimate is beyond the scope of this paper. In the rest of this section, we use block motion estimates in all our experiments.

**Individual interpolation models and entropy coding.** Figure 4c shows the performance of different interpolation models with and without entropy coding. For all models, entropy coding saves up to 52% BPP, at low bitrate regions, and at least 10% BPP, at high bitrate region. More interestingly, the short time-frame interpolation is almost free, achieving the same visual quality as an image-based model at one or two orders of magnitude fewer bits per pixel. This shows that most of our bitrate saving comes from the interpolation models at lower levels in the hierarchy.

## 5.2   Comparison to prior work

We now quantitatively evaluate our method on all three datasets and compare our method with today's prevailing codecs, HEVC (H.265), H.264, MPEG-4 Part 2, and H.261. For consistent comparison, we use the same GOP size, 12, for H.264 and HEVC. We test H.261 on only VTL and Kinetics-5K, since it does not support high-resolution (1920 × 1080) videos of the UVG dataset.

Figures 5-7 present the results. Despite its simplicity, our model greatly outperforms MPEG-4 Part 2 and H.261, performs on par with H.264, and close to state-of-the-art HEVC. In particular, on the high-resolution UVG dataset, it outperforms H.264 by a good margin and matches HEVC in terms of PSNR.

Our testing datasets are not just large in scale (>100K frames of >5K videos), they also consist of videos in a wide range of sizes (from 352 × 288 to 1920 × 1080), time (from 1990s for most VTL videos to 2018 for Kinetics), quality (from professional UVG to user uploaded Kinetics), and contents (from scenes, animals, to the 400 human activities in Kinetics). Our model, trained on only one dataset, works well on all of them.

Finally, we present qualitative results of three of the best performing models, MPEG-4 Part 2, H.264 and ours in Figure 8. All models here use $0.12 \pm 0.01$ BPP. We can see that in all datasets, our method shows faithful images without any blocky artifacts. It greatly outperforms MPEG-4 Part 2 without bells and whistles, and matches state-of-the-art H.264.

## 6   Conclusion

This paper presents, to the best of our knowledge, the first end-to-end trained deep video codec. It relies on repeated deep image interpolation. To disambiguate the interpolation, we encode a few compressible bits of information representing information not inferred from the neighboring key frames. This yields a faithful reconstruction instead of pure hallucination. The network is directly trained to optimize reconstruction, without prior engineering knowledge.

Our deep codec is simple, and outperforms the prevailing codecs such as MPEG-4 Part 2 or H.261, matching state-of-the-art H.264. We have not considered the engineering aspects such as runtime or real-time compression. We think they are important directions for future research.

In short, video compression powered by deep image interpolation achieves state-of-the-art performance without sophisticated heuristics or excessive engineering.
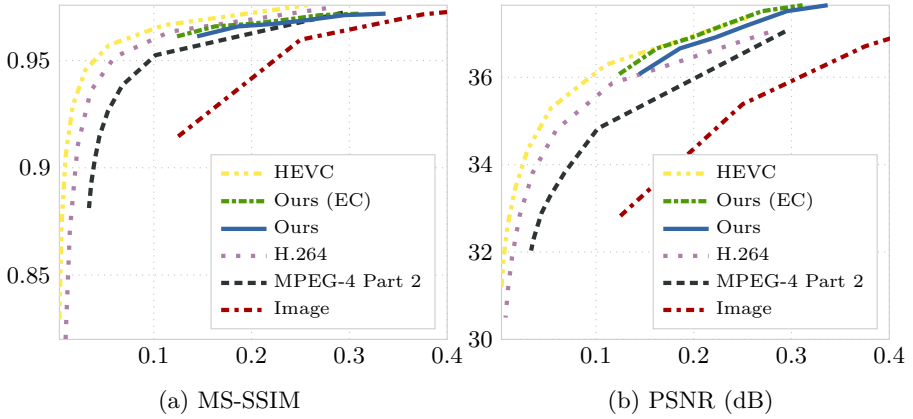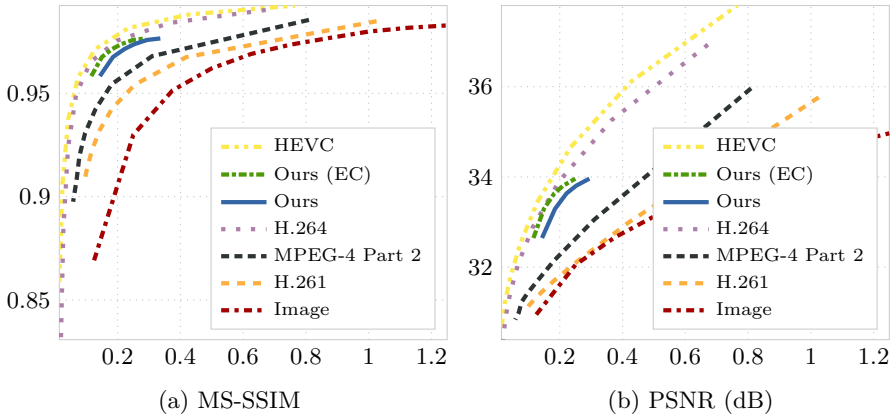
(a) MS-SSIM

(b) PSNR (dB)

Fig. 5: Performance on the UVG dataset.



(a) MS-SSIM

(b) PSNR (dB)

Fig. 6: Performance on the VTL dataset.



(a) MS-SSIM

(b) PSNR (dB)

Fig. 7: Performance on the Kinetics-5K dataset.

| Ground truth | MPEG-4 Part 2 | H.264 | Ours |



(a) Kinetics-5K



(b) VTL



(c) UVG

Fig. 8: Comparison of compression results at $0.12 \pm 0.01$ BPP. Our method shows faithful images without any blocky artifacts. (Best viewed on screen.) More examples and demo videos showing temporal coherence are available at https://chaoyuaw.github.io/vcii/.

## Acknowledgment

## References

1. Ultra video group test sequences. http://ultravideo.cs.tut.fi, accessed: 2018-03-11
2. Video trace library. http://trace.eas.asu.edu/index.html, accessed: 2018-03-11
3. WebP. https://developers.google.com/speed/webp/, accessed: 2018-03-11
4. Agustsson, E., Mentzer, F., Tschannen, M., Cavigelli, L., Timofte, R., Benini, L., Gool, L.V.: Soft-to-hard vector quantization for end-to-end learning compressible representations. In: NIPS (2017)
5. Baig, M.H., Koltun, V., Torresani, L.: Learning to inpaint for image compression. In: NIPS (2017)
6. Ballé, J., Laparra, V., Simoncelli, E.P.: End-to-end optimized image compression. In: ICLR (2017)
7. Carreira, J., Zisserman, A.: Quo vadis, action recognition? a new model and the kinetics dataset. In: CVPR (2017)
8. Farnebäck, G.: Two-frame motion estimation based on polynomial expansion. In: SCIA (2003)
9. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML (2015)
10. Jia, X., De Brabandere, B., Tuytelaars, T., Gool, L.V.: Dynamic filter networks. In: NIPS (2016)
11. Jiang, H., Sun, D., Jampani, V., Yang, M.H., Learned-Miller, E., Kautz, J.: Super slomo: High quality estimation of multiple intermediate frames for video interpolation. CVPR (2018)
12. Johnston, N., Vincent, D., Minnen, D., Covell, M., Singh, S., Chinen, T., Hwang, S.J., Shor, J., Toderici, G.: Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks. arXiv preprint arXiv:1703.10114 (2017)
13. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
14. Le Gall, D.: MPEG: A video compression standard for multimedia applications. Communications of the ACM (1991)
15. Liu, Z., Yeh, R., Tang, X., Liu, Y., Agarwala, A.: Video frame synthesis using deep voxel flow. In: ICCV (2017)
16. Mathieu, M., Couprie, C., LeCun, Y.: Deep multi-scale video prediction beyond mean square error. In: ICLR (2016)
17. Mentzer, F., Agustsson, E., Tschannen, M., Timofte, R., Van Gool, L.: Conditional probability models for deep image compression. arXiv preprint arXiv:1801.04260 (2018)
18. Networking Index, C.V.: Forecast and methodology, 2016-2021. CISCO White paper (2016)

19. Niklaus, S., Mai, L., Liu, F.: Video frame interpolation via adaptive separable convolution. In: ICCV (2017)
20. Oord, A.v.d., Kalchbrenner, N., Kavukcuoglu, K.: Pixel recurrent neural networks. In: ICML (2016)
21. Richardson, I.E.: Video codec design: developing image and video compression systems. John Wiley & Sons (2002)
22. Rippel, O., Bourdev, L.: Real-time adaptive image compression. In: ICML (2017)
23. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI (2015)
24. Schwarz, H., Marpe, D., Wiegand, T.: Overview of the scalable video coding extension of the H.264/AVC standard. TCSVT (2007)
25. Theis, L., Shi, W., Cunningham, A., Huszár, F.: Lossy image compression with compressive autoencoders. In: ICLR (2017)
26. Toderici, G., Vincent, D., Johnston, N., Jin Hwang, S., Minnen, D., Shor, J., Covell, M.: Full resolution image compression with recurrent neural networks. In: CVPR (2017)
27. Tsai, Y.H., Liu, M.Y., Sun, D., Yang, M.H., Kautz, J.: Learning binary residual representations for domain-specific video streaming. In: AAAI (2018)
28. Vondrick, C., Pirsiavash, H., Torralba, A.: Generating videos with scene dynamics. In: NIPS (2016)
29. Wang, Z., Simoncelli, E.P., Bovik, A.C.: Multiscale structural similarity for image quality assessment. In: ACSSC (2003)
30. Williams, R.J.: Simple statistical gradient-following algorithms for connectionist reinforcement learning. In: Reinforcement Learning. Springer (1992)
31. Witten, I.H., Neal, R.M., Cleary, J.G.: Arithmetic coding for data compression. Communications of the ACM (1987)
32. Wu, C.Y., Zaheer, M., Hu, H., Manmatha, R., Smola, A.J., Krähenbühl, P.: Compressed video action recognition. In: CVPR (2018)
33. Xue, T., Wu, J., Bouman, K., Freeman, B.: Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In: NIPS (2016)
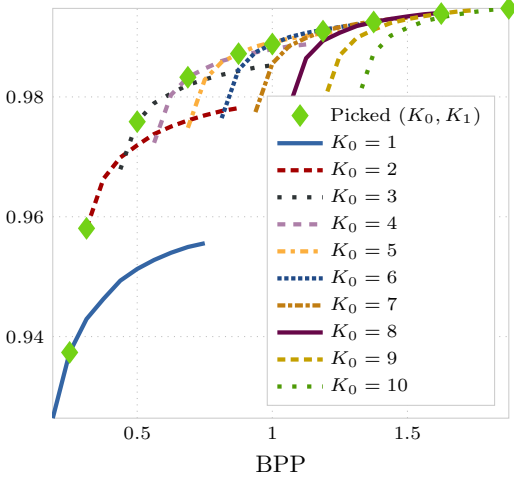
# Appendix

## A   Model Details

**Context feature fusion.** In experiments, we found that fusing the U-net features at resolution $\frac{W}{2} \times \frac{H}{2}$ yields good performance for all models. Fusing features at more resolutions improves the performance slightly, but requires more memory and computation. In this paper, we additionally use $\frac{W}{4} \times \frac{H}{4}$ features for the encoder of $\mathcal{M}_{1,2}$ and $\mathcal{M}_{3,3}$, and $\frac{W}{4} \times \frac{H}{4}$ and $\frac{W}{8} \times \frac{H}{8}$ features for the decoder of $\mathcal{M}_{3,3}$. Models are selected based on the performance on validation set.

**Probability estimation in AAC.** Adaptive arithmetic coding (AAC) relies on a good probability estimation for each bit, given previously decoded bits, to efficiently encode a binary representation. To estimate the probability, we follow Mentzer *et al.* and use a 3D-Pixel-CNN model [17, 20]. The model contains 11 layers of masked convolution. Each layer has 128 channels, and is followed by batch normalization [9] and relu. We train the models using Adam [13] with a learning rate of 0.0001 for 30K iterations.

## B   Bitrate Optimization

We present detailed results of bitrate optimization. Figure 9a shows the explored performance at the first level of the hierarchy. We pick good combinations from the envelope of the curves, and they proceed to the next level. Figure 9b presents the final combinations. We use these combinations for the experiments in our paper unless otherwise noted.

| $K_0$ | $K_1$ | $K_2$ | $K_3$ | Average bitrate |
|---|---|---|---|---|
| 5 | 3 | 2 | 1 | 0.109 |
| 7 | 3 | 2 | 2 | 0.151 |
| 10 | 4 | 2 | 2 | 0.188 |
| 10 | 10 | 2 | 2 | 0.219 |
| 10 | 10 | 6 | 2 | 0.260 |
| 10 | 10 | 10 | 2 | 0.302 |

(a) MS-SSIM. Beam search for $K_1$ given $K_0$.      (b) Final combinations.

Fig. 9: Bitrate optimization. $K_0$, $K_1$, $K_2$ and $K_3$ correspond to the number of iterations used in the I-frame model, $\mathcal{M}_{6,6}$, $\mathcal{M}_{3,3}$ and $\mathcal{M}_{1,2}$ respectively. (a) Beam search at the first level of hierarchy. Picked $(K_0, K_1)$-combinations proceed to the next level. (b) Final combinations returned by our algorithm.