# Cross-Attention Transformer for Video Interpolation

Hannah Halin Kim[0000−0003−2588−0190], Shuzhi Yu[0000−0003−2514−381X], Shuai Yuan[0000−0003−4039−0464], and Carlo Tomasi[0000−0001−6104−6641]

Duke University, Durham NC 27708, USA
{hannah,shuzhiyu,shuai,tomasi}@cs.duke.edu

**Abstract.** We propose TAIN (Transformers and Attention for video INterpolation), a residual neural network for video interpolation, which aims to interpolate an intermediate frame given two consecutive image frames around it. We first present a novel vision transformer module, named Cross Similarity (CS), to globally aggregate input image features with similar appearance as those of the predicted interpolated frame. These CS features are then used to refine the interpolated prediction. To account for occlusions in the CS features, we propose an Image Attention (IA) module to allow the network to focus on CS features from one frame over those of the other. TAIN outperforms existing methods that do not require flow estimation and performs comparably to flow-based methods while being computationally efficient in terms of inference time on Vimeo90k, UCF101, and SNU-FILM benchmarks.

## 1 Introduction

Video interpolation [1,2,3,4,5,6,7] aims to generate new frames between consecutive image frames in a given video. This task has many practical applications, ranging from frame rate up-conversion [8] for human perception [9], video editing for object or color propagation [10], slow motion generation [11], and video compression [12].

Recent work on video interpolation starts by estimating optical flow in both temporal directions between the input frames. Some systems [3,4,5,13,14,15] use flow predictions output by an off-the-shelf pre-trained estimator such as FlowNet [16] or PWC-Net [17], while others estimate flow as part of their own pipeline [7,11,18,19,20,21,22,23,24]. The resulting bi-directional flow vectors are interpolated to infer the flow between the input frames and the intermediate frame to be generated. These inferred flows are then used to warp the input images towards the new one.

While these flow-based methods achieve promising results, they also come with various issues as follows. First, they are computationally expensive and rely heavily on the quality of the flow estimates (see Figure 1). Specifically, state-of-the art flow estimators [16,17,25] require to compute at least two four-dimensional cost volumes at all pixel positions. Second, they are known to suffer

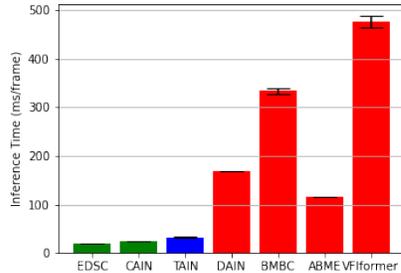| | Network | Inference Time | PSNR |
|---|---|---|---|
| No Flow | EDSC | **19.37 ± 0.06** | 34.84 |
| | CAIN | 25.21 ± 0.45 | 34.65 |
| | TAIN (ours) | 32.59 ± 0.82 | **35.02** |
| Flow | DAIN | 168.56 ± 0.33 | 34.71 |
| | BMBC | 333.24 ± 6.60 | 35.01 |
| | ABME | **116.66 ± 1.07** | 36.18 |
| | VFIformer | 476.06 ± 11.92 | **36.50** |



Fig. 1: Inference times (in milliseconds per frame prediction) of existing video interpolation methods and TAIN (ours) on a single P100 GPU. The inference times are computed as average and standard deviation over 300 inferences using two random images of size $256 \times 256$ as input. As a reference, we also list performance (PSNR) on Vimeo-90k [7]. Flow-based methods (red) have much higher inference times compared to TAIN (blue) and other non-flow-based methods (green). Bold values show the best performance in each panel, and underlined bold values show the best performance across both panels. Best viewed in color.

at occlusions, where flow is undefined; near motion boundaries, where flow is discontinuous; and in the presence of large motions [26,27]. For instance, RAFT [25] achieves an End-Point Error (EPE) of 1.4 pixels on Sintel [28], but the EPE rises to 6.5 within 5 pixels from a motion boundary and to 4.7 in occlusion regions. Warping input images or features with these flow estimates will lead to poor predictions in these regions. This is especially damaging for video interpolation as motion boundaries and occlusions result from motion [26], and are therefore important regions to consider for motion compensation. Further, these flow estimators are by-and-large trained on synthetic datasets [28,29] to avoid the cost and difficulties of annotating real video [30], and sometimes fail to capture some of the challenges observed in real data, to which video interpolation is typically applied.

Apart from flow inputs, many of these methods also utilize additional inputs and networks to improve performance, which adds to their computational complexity. These include depth maps [5], occlusion maps [5,7,11,14,31], multiple input frames [32] for better flow estimation, image classifiers (*e.g.*, VGG [33], ResNet [34]) pretrained on ImageNet [35] for contextual features [1,2,3,11,31], adversarial networks for realistic estimations [31], and event cameras [36,37,38,39] to detect local brightness changes across frames. Instead, we achieve competitive performance with a single network that uses two consecutive input frames captured with commonly available devices.

The proposed TAIN system (Transformers and Attention for video INterpolation) is a residual neural network that requires no estimation of optical flow. Inspired by the recent success of vision transformers [40,41,42,43], we employ a novel transformer-based module named Cross Similarity (CS) in TAIN. Video

interpolation requires the network to match corresponding points across frames, and our CS transformer achieves this through cross-attention. Vision Transformers typically use self-attention to correlate each feature of an image to every other feature *in the same image.* The resulting similarity maps are then used to collect relevant features from other image locations. Our CS module instead compares features *across frames*, namely, between an input frame and the current intermediate frame prediction. High values in this cross-frame similarity map indicate features with similar appearance. These maps are then used to aggregate features from the appropriate input image to yield CS features, and these are used in turn to refine the frame prediction. While there exists recent work on video interpolation that utilize vision transformers [22,32], our CS is specifically designed for video interpolation as it computes cross-frame similarity rather than within-frame self-similarity.

To account for occlusions or motion boundaries, we use CS scores in an Image Attention (IA) module based on spatial attention [44]. Given a feature in the tentative frame prediction, its maximum similarity score from our CS module will indicate whether or not a similar feature exists in either input frame. This maximum score will be high if such a feature exists and low if it does not. If the feature in the current frame prediction is occluded in one of the input images, the corresponding score will likely be low. Features at positions that straddle a motion boundary in the frame prediction often have a low maximum score as well. This is because boundary features contain information about pixels from both side of the boundary, and the particular mixture of information will change if the two sides move differently. The CS features aggregated from these low maximum similarity scores will not help in refining the current frame prediction, and our IA module learns to suppress them.

Thanks to CS transformer and IA module, TAIN improves or performs comparably to existing methods on various benchmarks, especially compared to those that do not require flow estimation. Our contributions are as follows:

- A novel Cross Similarity module based on vision transformer that aggregates features of the input image frames that have similar appearance to predicted-frame features. These aggregated features help refine the frame prediction.
- A novel Image Attention module that gives the predictor the ability to weigh features in one input frame over those in the other. This information is shown to be especially helpful near occlusions and motion boundaries.
- State-of-the-art performance on Vimeo-90k, UCF101, and SNU-FILM among methods that do not require optical flow estimation.

## 2    Related Work

### 2.1    Video Interpolation

Deep learning has rapidly improved the performance of video interpolation in recent work. Long *et al.* [45] are the first to use a CNN for video interpolation. In their system, an encoder-decoder network predicts the intermediate frame

directly from two image frames using inter-frame correspondences. Subsequent work can be categorized largely into kernel-, flow-, and attention-based.

**Kernel-based approaches** compute the new frame with convolutions over local patches, and use CNNs to estimate spatially-adaptive convolutional kernels [2,6,46]. These methods use large kernel sizes to accommodate large motion, and large amounts of memory are required as a result when frames have high resolution. Niklaus *et. al* propose to use separable convolutional kernels to reduce memory requirements and further improve results. Since kernel-based methods cannot handle motion larger than the pre-defined kernel size, EDSC [2] estimates not only adaptive kernels, but also offsets, masks, and biases to retrieve information from non-local neighborhoods. EDSC achieves the current state-of-the-art performance among methods that do not require optical flow.

**Flow-based approaches** to video interpolation rely on bi-directional flow estimates using off-the-shelf pre-trained flow estimators (e.g. FlowNetS, PWC-Net) [3,4,5,13,14,15], or estimate flow as part of their pipeline [7,11,18,19,20,21,22,23,24]. Most of these methods assume linear motion between frames and use the estimated flow to warp input images and their features to the intermediate time step for prediction. In order to avoid making this linear motion assumption, Gui *et al.* [18] do not predict flow between two input frames but instead attempt to produce flow directly between the intermediate frame being predicted and the two input frames. Park *et al.* [23] do compute a tentative intermediate frame from the given flow estimates, but they re-estimate flow between the tentative predicted frame and the input images. These new estimates are then used for a final estimation of the intermediate frame. Multiple optical flow maps have also been used to account for complex motion patterns and mitigate the resulting prediction artifacts [15,21]. While these flow-based methods show promising results, they are computationally expensive. We do not require flow estimates in our approach.

Some approaches [5,14,31] integrate **both flow-based and kernel-based methods** by combining optical flow warping with learned adaptive local kernels. These methods perform robustly in the presence of large motions and are not limited by the assumption of a fixed motion range. They use small kernels that require less memory but are still expensive.

Recently, work by Choi *et al.* [1] proposes a residual network called CAIN that interpolates video through **attention mechanism** [44] without explicit computation of kernels or optical flow to curb model complexity and computational cost. The main idea behind their design is to distribute the information in a feature map into multiple channels through PixelShuffle, and extract motion information by processing the channels through a Channel Attention module. In our work, we extend CAIN with a novel vision transformer module and spatial attention module, still without requiring flow estimates or adaptive kernels.

## 2.2   Vision Transformers

Transformers [42,43] have shown success in both computer vision [40,41,47] and natural language processing [42,48] thanks to their ability to model long-range

dependencies. Self-attention [42,48] has shown the most success among various modules in the transformer architecture, and derives query, key, and value vectors from the same image. Due to their content-adaptive nature, transformers have also been applied to video interpolation. Shi *et. al* [32] consider four frame inputs and propose to use a self-attention transformer based on SWIN [49] to capture long-range dependencies across both space and time. Lu *et. al* [22] propose a self-attention transformer along with a flow estimator to model long-range pixel correlation for video interpolation. Their work currently achieves the state-of-the-art performance on various video interpolation benchmarks. Different from existing work, we do not use self-attention in our work, but instead use cross-attention. In self-attention, query, key, and value are different projections of the same feature. We use cross-attention, where query and key are the same projections (shared weights) of different features. Specifically, our transformer module selects and refines features based on the similarities between the interpolated frame (query) and the two input frames (keys) while handling the occlusions that occur in the two input frames. On the other hand, existing transformer-based methods compare the features of the input frames without any explicit consideration of the interpolated frame or occlusions.

## 3   Method

The proposed TAIN method for video interpolation aims to predict frame $I_t \in \mathbb{R}^{h \times w \times 3}$ at time $t = 0.5$, given two consecutive images $I_0, I_1 \in \mathbb{R}^{h \times w \times 3}$ at times 0 and 1. We do not require any computation of flow, adaptive convolution kernel parameters, or warping, but instead utilize cross-similarity transformer and spatial attention mechanism. Specifically, a novel vision transformer module called the Cross Similarity (CS) module globally aggregates features from input images $I_0$ and $I_1$ that are similar in appearance to those in the current prediction $\hat{I}_t$ of frame $I_t$ (Section 3.2). These aggregated features are then used to refine the prediction $\hat{I}_t$, and the output from each residual group is a new refinement. To account for occlusions of the interpolated features in the aggregated CS features, we propose an Image Attention (IA) module to enable the network to prefer CS features from one frame over those of the other (Section 3.3). See Figure 2 for an overview of the network. Before describing our network we summarize CAIN [1], on which our work improves.

### 3.1   CAIN

CAIN [1] is one of the top performers for video interpolation and does not require estimation of flow, adaptive convolution kernels, or warping. Instead, CAIN utilizes PixelShuffle [50] and a channel attention module. PixelShuffle [50] rearranges the layout of an image or a feature map without any loss of information. To down-shuffle, activation values are merely rearranged by reducing each of the two spatial dimensions by a factor of $s$ and increasing the channel dimension by a factor of $s^2$. Up-shuffling refers to the inverse operation. This parameter-free
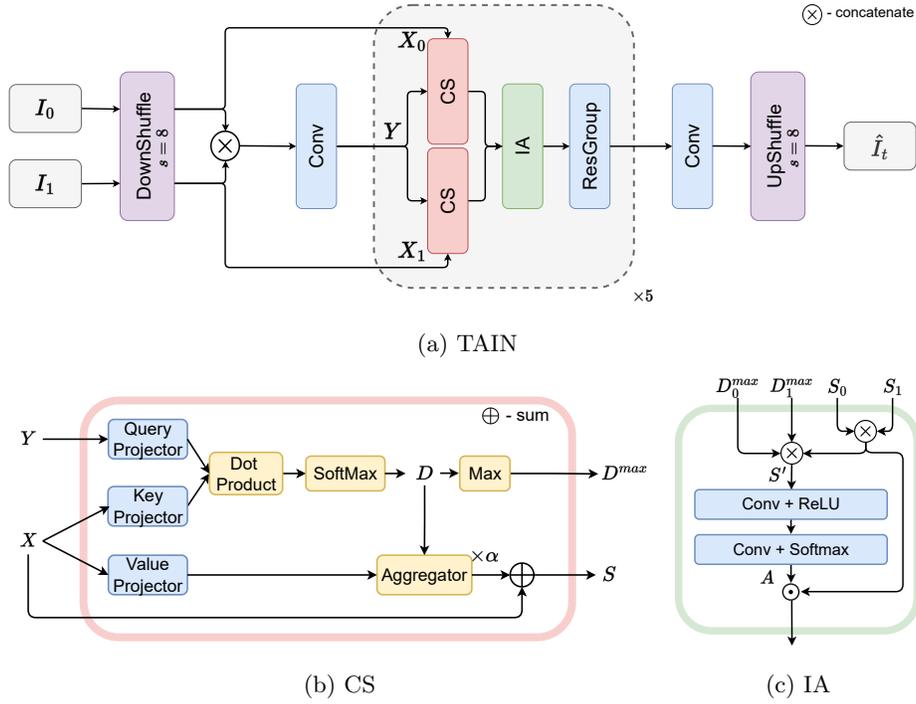
(a) TAIN



(b) CS



(c) IA

Fig. 2: (a) Overview of the proposed TAIN network for video interpolation. The two consecutive frames $I_0$ and $I_1$ are down-shuffled and concatenated along the channel dimension. They are then processed with five residual groups (Res-Groups) with a (b) CS transformer and an (c) IA module before being up-shuffled back to the original resolution to yield the final prediction $\hat{I}_t$. The output from each ResGroup is a refinement of the output from the previous block.

operation allows CAIN to increase the receptive field size of the network's convolutional layers without losing any information. CAIN first down-shuffles ($s = 8$) the two input images and concatenates them along the channel dimension before feeding them to a network with five ResGroups (groups of residual blocks). With the increased number of channels, each of the blocks includes a channel attention module that learns to pay attention to certain channels to gather motion information. The size of the features remains $h/8 \times w/8 \times 192$ throughout CAIN, and the final output map is up-shuffled back to the original resolution of $h \times w \times 3$.

## 3.2   Cross Similarity (CS) Module

All points in the predicted frame $I_t$ appear either in $I_0$ or $I_1$ or in both, except in rare cases where a point appears for a very short time between the two consecutive time frames. These ephemeral apparitions cannot be inferred from $I_0$ or $I_1$ and are ignored here. For the remaining points, we want to find features of $I_0$
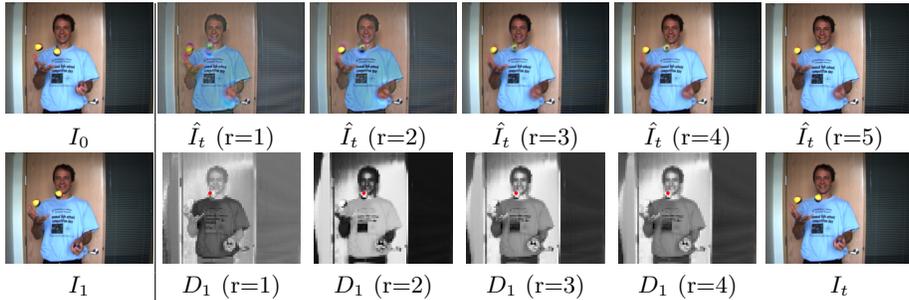
Fig. 3: Visualization of intermediate predictions and their cross similarity maps $D_1$ with $I_1$ across all five ResGroups ($r = 1 \dots 5$) using an example from Middlebury [51] dataset. The first column shows the two input frames, $I_0$ (top) and $I_1$ (bottom). The next four columns show the predicted intermediate frame $\hat{I}_t$ after each of the first four residual blocks with a query point highlighted in blue (top), and its corresponding similarity map $D_1$ from our CS module with the point of highest similarity highlighted in red (bottom). The values in the similarity maps $D_1$ show large scores (closer to white) whenever query and key features are similar in appearance. The last column shows the final output from TAIN (top), which is not used in any CS module, and the ground-truth intermediate frame $I_t$ (bottom) as a reference.

or $I_1$ that are similar in appearance to the features in the predicted intermediate frame $\hat{I}_t$, and use them to refine $\hat{I}_t$. We use a transformer to achieve this, where we compare each feature from $\hat{I}_t$ with features from $I_0$ and $I_1$, and use similar features to refine $\hat{I}_t$.

While transformers typically use self-similarity [42,48], wherein query, key, and value are based on similarities within the same feature array, we extend this notion through the concept of cross-image similarity. Specifically, we use features from the input images $I_0$ or $I_1$ (first column in Figure 3) as keys and values, and features from the current frame prediction $\hat{I}_t$ (the remaining images in the first row of Figure 3), as queries. Given a query feature normalized to have unit Euclidean norm (*e.g.*, features for the blue points on the ball in the first row of Figure 3), our CS module compares it to all the similarly-normalized key features through a dot product, which are then Softmax-normalized to obtain the corresponding similarity maps $D$ (*e.g.*, images in the second row of Figure 3). The similarity matrix $D$ is used to find the location of the largest similarity score (*e.g.*, red dots in the the second row of Figure 3), where we retrieve our value features (equation (3) later on), which are projections of the input image features, to yield aggregated input features $S$ based on similarity. These aggregated features are then used to refine the intermediate prediction of $\hat{I}_t$.

**Mathematical Formulation** Let $X \in \mathbb{R}^{h/s \times w/s \times d}$ denote the (down-shuffled) feature map from one of the input images $I_0, I_1 \in \mathbb{R}^{h \times w \times 3}$ and let $Y \in \mathbb{R}^{h/s \times w/s \times d}$

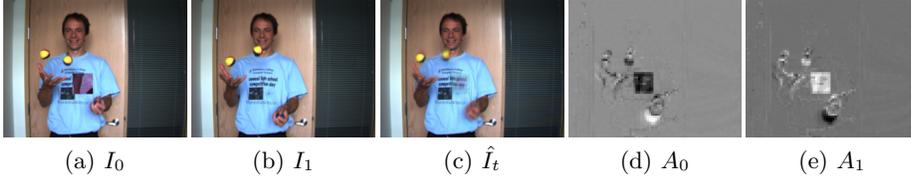(a) $I_0$     (b) $I_1$     (c) $\hat{I}_t$     (d) $A_0$     (e) $A_1$

Fig. 4: Visualization of Image Attention maps (d) $A_0$ and (e) $A_1$ using an example from Middlebury [51] (White is large and black is small). For visualization purposes only, we add a square patch to (a) $I_0$ on the person's chest, and show that the IA module assigns a higher weight to the CS features from $I_1$ on the person's chest in (c) $\hat{I}_t$ than to those from $I_0$, which the patch occludes.

denote the feature map of the predicted intermediate frame $\hat{I}_t \in \mathbb{R}^{h \times w \times 3}$. Let $M^{\mathbf{i}}$ represent feature at pixel location $\mathbf{i}$ in a given feature map $M$. Our CS module computes query feature map $Q \in \mathbb{R}^{h/s \times w/s \times d}$ from $Y$. It also computes key feature map $K \in \mathbb{R}^{h/s \times w/s \times d}$ and value feature map $V \in \mathbb{R}^{h/s \times w/s \times d}$ from $X$ as follows for all pixel locations $\mathbf{i}, \mathbf{j}$:

$$Q^{\mathbf{i}} = W_{qk}Y^{\mathbf{i}}, \quad K^{\mathbf{j}} = W_{qk}X^{\mathbf{j}}, \quad V^{\mathbf{j}} = W_v X^{\mathbf{j}} . \tag{1}$$

The matrices $W_{qk}, W_v \in \mathbb{R}^{d \times d}$ are learnable. Note that we use cross-attention where the query and key features are the same projections (shared weights $W_{qk}$) of different features. This is different from the self-attention transformers used in existing work [22,32], where the query, key, and value features are different projections (no shared weights) of the same features.

Each query feature $Q^{\mathbf{i}}$ is compared with all the key features $K^{\mathbf{j}}$ to compute a similarity matrix $D^{\mathbf{i}} \in \mathbb{R}^{h/s \times w/s}$ that captures their similarity:

$$D^{\mathbf{i}} = sim(Q^{\mathbf{i}}, K^{\mathbf{j}}) = \frac{\exp(Q^{\mathbf{i}T}K^{\mathbf{j}}/\sqrt{d})}{\sum_{\mathbf{j}} \exp(Q^{\mathbf{i}T}K^{\mathbf{j}}/\sqrt{d})} . \tag{2}$$

Using this similarity matrix $D^{\mathbf{i}}$, the CS module finally computes the aggregated similarity feature $S^{\mathbf{i}} \in \mathbb{R}^{h/s \times w/s \times d}$ by taking the value feature $V^{\mathbf{i}_{\max}}$ at location $\mathbf{i}_{\max}$ corresponding to the maximum similarity score for each query feature $Q^{\mathbf{i}}$ and adding the result to $Y^{\mathbf{i}}$

$$S^{\mathbf{i}} = Y^{\mathbf{i}} + \alpha V^{\mathbf{i}_{\max}} , \tag{3}$$

where $\alpha$ is a learnable scalar parameter initialized to zero.

### 3.3   Image Attention (IA) Module

We propose an Image Attention (IA) module based on spatial attention [44] to weigh our two CS features $S_0$ and $S_1$ computed from the input frames as shown above. The IA module enables TAIN to prioritize or suppress features

from one input image over those from the other for a given spatial location. This is useful especially at occlusions, where a feature from $I_t$ appears in one input image (likely yielding high similarity scores) but not in the other (likely yielding low similarity scores). This also helps on motion boundaries, where features encode information from both sides of the boundary. As the two sides move relative to each other, the specific mixture of features changes. The IA module compares the two CS feature maps $S_0$ and $S_1$ to construct two IA weight maps $A_0, A_1 \in [0,1]^{h/s \times w/s \times 1}$ where $A_0 + A_1 = 1^{h/s \times w/s \times 1}$ as shown below. These weight maps are multiplied with the corresponding CS maps $S_0$ and $S_1$ before they are concatenated and fed to the next ResGroup. See Figure 4.

**Mathematical Formulation** Let $S' \in \mathbb{R}^{h/s \times w/s \times 2(d+1)}$ be the concatenation of $S_0$, $S_1$, and the two maximum similarity maps $D_0$ and $D_1$ along the channel dimension. Our IA module first computes image attention weights $A \in [0,1]^{h/s \times w/s \times 2}$ by applying two $1 \times 1$ convolutional layers with ReLU and Softmax on $S'$ as shown in Figure 2c. More formally, we compute $A$ as:

$$A = \sigma(W_2 * (\rho(W_1 * S'))) \tag{4}$$

where $\sigma(\cdot)$ denotes the Softmax function, $\rho(\cdot)$ denotes the ReLU function, and $W_1$ and $W_2$ are the weights of the two $1 \times 1$ convolution layers. This image attention weight map $A$ can be seen as a concatenation of maps $A_0, A_1 \in [0,1]^{h/s \times w/s \times 1}$ where $A_0 + A_1 = 1$ which are used to weigh the aggregated similarity features $S_0$ and $S_1$ to obtain the weighted features $\tilde{S}_0$ and $\tilde{S}_1$:

$$\tilde{S}_0 = A_0 \odot S_0 \quad \text{and} \quad \tilde{S}_1 = A_1 \odot S_1 \ , \tag{5}$$

where $\odot$ is the element-wise product. Note that this is different from the channel attention module from CAIN [1]. The channel attention module computes a $(2d)$-dimensional vector weight over the channel dimension, while the IA module computes a $h/s \times w/s$ weight over the spatial dimension.

### 3.4   TAIN Architecture and Training Details

Figure 2a shows an overview of the TAIN network for video interpolation. TAIN extends CAIN by applying the proposed CS transformer and IA module after each of the intermediate ResGroups. Each CS transformer obtains query features $Q$ from the features $Y$ from the previous ResGroup, and key and value features from one of the two input image features $X_0$ or $X_1$. We remove the residual connections around each ResGroup so that the output from each ResGroup is a new refinement of the output from the previous ResGroup. Figure 3 shows sample predictions after each of the five ResGroups that the query features are based on. With the IA module, the CS features $S_0$ and $S_1$ are weighted based on their maximum similarity scores from $D_0$ and $D_1$. The two weighted CS features, $\tilde{S}_0$ and $\tilde{S}_1$, are then used to refine the prediction in the next ResGroup.

   In order to train the network for cases of occlusions and large motion, we add synthetic moving occlusion patches to the training data. Specifically, we

first randomly crop a square patch of size between $21 \times 21$ and $61 \times 61$ from a different sample in the training set. This cropped patch is pasted onto the input and label images and translated in a linear motion across the frames. We apply this occluder patch augmentation to first pre-train on TAIN and then fine-tune on the original training dataset. Following CAIN [1], we also augment the training data with random flips, crops, and color jitter.

Following the literature [1,3,6,7,11,31], we train our network using the $L_1$ loss on the difference between the predicted $\hat{I}_t$ and true $I_t$ intermediate frames: $\|\hat{I}_t - I_t\|_1$. As commonly done in the flow literature [52,53], we also include an $L_1$ loss on the difference between the gradients of the predicted and true intermediate frames: $\|\nabla \hat{I}_t - \nabla I_t\|_1$. We use the weighted sum of the two losses as our final training loss: $\mathcal{L} = \|\hat{I}_t - I_t\|_1 + \gamma \|\nabla \hat{I}_t - \nabla I_t\|_1$, where $\gamma = 0.1$.

Another common loss used in the literature [1,2,3,11,31] is perceptual loss: $\|\phi\left(\hat{I}_t\right) - \phi\left(I_t\right)\|_2^2$, where $\phi\left(\cdot\right)$ is a feature from a ImageNet pretrained VGG-19. As this loss depends on another network, adding to the computational complexity, we do not use this loss and still show performance improvements.

We implement TAIN in PyTorch [54] [1] and train our network with a learning rate of $10^{-4}$ through Adam [55] optimizer.

## 4    Datasets and Performance Metrics

As customary [1], we train our model on Vimeo90K [7], and evaluate it on four benchmark datasets for video interpolation, *i.e.*, Vimeo90K [7], UCF101 [56], SNU-FILM [1], and Middlebury [51]. Vimeo90K [7] consists of 51,312 triplets with a resolution of $256 \times 448$ partitioned into a training set and a testing set. UCF101 [56] contains human action videos of resolution $256 \times 256$. For the evaluation of video interpolation, Liu *et al.* [57] constructed a test set by selecting 379 triplets from UCF101. The SNU Frame Interpolation with Large Motion (SNU-FILM) [1] dataset contains videos with a wide range of motion sizes for evaluation of video interpolation methods. The dataset is stratified into four settings, Easy, Medium, Hard, and Extreme, based on the temporal gap between the frames. Middlebury [51] includes 12 sequences of images for evaluation. The images are combination of synthetic and real images that are often used as evaluation for video interpolation.

Following the literature, we use Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) [58] to measure performance. For the Middlebury dataset [51], we use Interpolation Error (IE).

## 5    Results

We first compare TAIN with the current state-of-the-art methods on video interpolation with two frame inputs in Table 1. Top panel of Table 1 lists kernel-based

---

[1] Code is available at https://github.com/hannahhalin/TAIN.

Table 1: Comparison to the existing methods across Vimeo90k, UCF101, SNU, and Middlebury (M.B.) datasets. Top panel shows kernel-based methods (K), second panel shows attention-based methods (Att.), third panel shows methods based on both kernels and flow (K+F), and bottom panel shows flow-based methods. Higher is better for PSNR and SSIM, and lower is better for IE. Bold values show the best performance in each panel, and underlined values show the best performance across all panels. TAIN (ours) outperforms existing methods based on kernel and attention, and performs comparably to those based on flow on Vimeo90k, UCF101, and SNU datasets.

|  | Method | Vimeo90k | | UCF101 | | SNU-easy | | SNU-extreme | | M.B. |
|  |  | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | IE |
|---|---|---|---|---|---|---|---|---|---|---|
| K | SepConv | 33.79 | 97.02 | 34.78 | 96.69 | 39.41 | 99.00 | 24.31 | **84.48** | 2.27 |
|  | EDSC | **34.84** | **97.47** | **35.13** | **96.84** | **40.01** | **99.04** | **24.39** | 84.26 | **2.02** |
| Att. | CAIN | 34.65 | 97.29 | 34.91 | 96.88 | 39.89 | 99.00 | 24.78 | 85.07 | **2.28** |
|  | TAIN (Ours) | **35.02** | **97.51** | **35.21** | **96.92** | <u>**40.21**</u> | **99.05** | **24.80** | **85.25** | 2.35 |
| K+F | AdaCoF | 34.47 | 97.30 | 34.90 | 96.80 | **39.80** | 99.00 | 24.31 | 84.39 | 2.24 |
|  | MEMC | 34.29 | 97.39 | 34.96 | 96.82 | - | - | - | - | 2.12 |
|  | DAIN | **34.71** | **97.56** | **34.99** | **96.83** | 39.73 | **99.02** | **25.09** | **85.84** | **2.04** |
| Flow | TOFlow | 33.73 | 96.82 | 34.58 | 96.67 | 39.08 | 98.90 | 23.39 | 83.10 | 2.15 |
|  | CyclicGen | 32.09 | 94.90 | 35.11 | 96.84 | 37.72 | 98.40 | 22.70 | 80.83 | - |
|  | BMBC | 35.01 | 97.64 | 35.15 | 96.89 | 39.90 | 99.03 | 23.92 | 84.33 | 2.04 |
|  | ABME | 36.18 | 98.05 | 35.38 | 96.98 | 39.59 | 99.01 | 25.42 | 86.39 | 2.01 |
|  | VFIformer | <u>**36.50**</u> | <u>**98.16**</u> | <u>**35.43**</u> | <u>**97.00**</u> | 40.13 | <u>**99.07**</u> | <u>**25.43**</u> | <u>**86.43**</u> | <u>**1.82**</u> |

methods, *i.e.*, SepConv [2] and EDSC [2], second panel lists attention-based methods, *i.e.*, CAIN [1] and TAIN (ours), third panel lists methods based on both kernels and flow estimations, *i.e.*, AdaCoF [31], MEMC [14], and DAIN [5], and bottom panel lists flow-based methods, *i.e.*, TOFlow [7], CyclicGen [59], BMBC [4], ABME [23], and VFIformer [22]. Bold values show the highest performance in each panel while underlined values show the highest performance across all panels. TAIN outperforms existing kernel-based methods across all benchmarks. Comparing to the flow-based methods, TAIN outperforms all listed methods except for ABME and VFIformer. However, as shown in Figure 1, the inference times of AMBE and VFIformer are about 4 and 15 times longer, respectively, than those of TAIN. DAIN and BMBC take around 5 times longer than TAIN to inference while performing comparably to TAIN.

Figure 5 visualizes examples of predictions from TAIN and the best performing methods in each panel of Table 1. Compared with kernel- and attention-based methods, TAIN is able to find the correct location of the moving object, *e.g.* shovel and hula hoop in red box, while keeping the fine details, *e.g.* shaft of the shovel, letters on the plane, and the hula hoop.

Fig. 5: Visualization of our proposed method and its comparison to the current state-of-the-art methods [1,2,5,22] on examples from Vimeo90k and UCF101.

**Inference Time** Figure 1 compares the inference time of TAIN and other state-of-the-art methods listed in Table 1. To measure time, we create two random images of size $256 \times 256$, the same size as those of UCF101 dataset, and evaluate it 300 times using a P100 GPU and report their average and standard deviation. As mentioned above, TAIN achieves competitive performance as the existing flow-based methods while taking a fraction of time for inference.
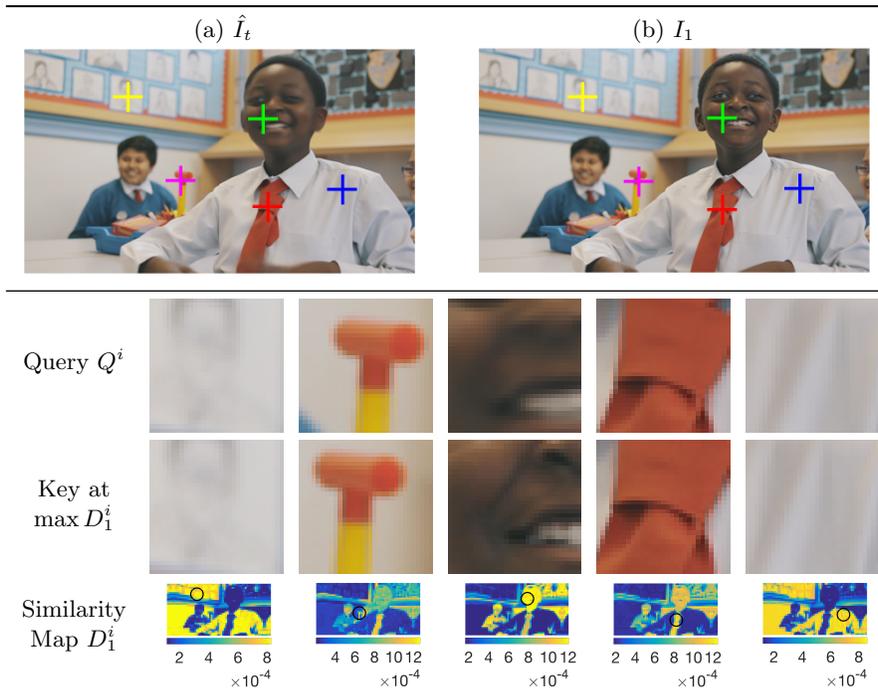
Fig. 6: Visualization of patches with the highest similarity scores from $D$ in our proposed CS transformer. Top panel shows an example of (a) $\hat{I}_t$ and (b) $I_1$ from the test set of Vimeo90k [7] dataset. Example query features in $\hat{I}_t$ are shown with '+' mark, and their key feature with the highest similarity score are shown in $I_1$ in corresponding colors. Bottom panel shows $31 \times 31$ patches extracted from the query points ('+' in $\hat{I}_t$), their corresponding key patches with the highest similarity score ('+' in $I_1$), and their corresponding similarity map $D_1^i$ with the highest score circled. Our CS module successfully aggregates similar appearance features when refining the interpolation prediction $\hat{I}_t$.

## 6   Ablation Study

**Visualization of the Components of CS Module** We visualize the components of our proposed CS module using an example from the test set of Vimeo90k [7] dataset in Figure 6. The top panel of this figure shows (a) $\hat{I}_t$ and (b) $I_1$, where example query features $Q^i$ are shown with '+' mark in (a) and their corresponding key features with the highest similarity score are shown with '+' mark in (b) with the corresponding colors. Bottom panel shows $31 \times 31$ patches extracted from the location '+' of each query $Q^i$ and key $K_1^i$ features from (a) $\hat{I}_t$ and (b) $I_1$, respectively. We also include visualization of the corresponding similarity maps $D_1^i$ and highlight their maximum score. As shown, our CS module successfully extracts similar appearance features when refining $\hat{I}_t$.

Table 2: Performance changes with our proposed modules: **IA** - Image Attention module; **CS** - Cross Similarity transformer; **#RG** - Number of ResGroups. Using all the components yields the best overall performance. While increasing the number of ResGroups yields consistently higher performance, the performance plateaus after 5 ResGroups, which we use for TAIN (top row).

| IA | CS | #RG | Vimeo90k PSNR | SSIM | UCF101 PSNR | SSIM | SNU-easy PSNR | SSIM | SNU-extreme PSNR | SSIM | M.B. IE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ | 5 | 35.02 | 97.51 | 35.21 | **96.92** | **40.21** | **99.05** | 24.80 | 85.25 | 2.35 |
| ✓ | ✓ | 4 | 34.82 | 97.38 | **35.22** | **96.92** | 40.20 | **99.05** | 24.76 | 84.97 | 2.39 |
| ✓ | ✓ | 6 | **35.06** | 97.52 | **35.22** | **96.92** | 40.20 | **99.05** | 24.74 | 85.09 | 2.33 |
| ✓ | ✓ | 7 | **35.06** | **97.53** | **35.22** | **96.92** | **40.21** | **99.05** | 24.75 | 85.09 | **2.32** |
|  | ✓ | 5 | 34.81 | 97.40 | 35.17 | 96.91 | 40.14 | 99.04 | 24.81 | 85.21 | 2.49 |
|  |  | 5 | 34.76 | 97.38 | 35.05 | 96.88 | 40.00 | 99.02 | **24.82** | **85.28** | 2.66 |

**Effect of Each Component**  Table 2 shows the performance changes with varying combinations of our proposed components, *i.e.*, Image Attention (IA), Cross Similarity transformers (CS), and the number of ResGroups (RG). Comparing the first four rows that list performances with the changing number of ResGroups from 4 to 7, we see that the performance increases with the number of ResGroups. However, it plateaus after 5 ResGroups, which we choose to use for TAIN (top row) for computational efficiency. In addition, each of the two main components, IA and CS module, contributes to the success of our method.

## 7    Conclusion

We propose TAIN, an extension of the CAIN network, for video interpolation. We utilize a novel vision transformer we call Cross Similarity module to aggregate input image features that are similar in appearance to those in the predicted frame to further refine the prediction. To account for occlusions in these aggregated features, we propose a spatial attention module we call Image Attention to suppress any features from occlusions. Combining both these components, TAIN outperforms existing methods that do not require flow estimation on multiple benchmarks. Compared to methods that utilize flow, TAIN performs comparably while taking a fraction of their time for inference.

# References

1. Choi, M., Kim, H., Han, B., Xu, N., Lee, K.M.: Channel attention is all you need for video frame interpolation. In: AAAI. (2020)
2. Niklaus, S., Mai, L., Liu, F.: Video frame interpolation via adaptive separable convolution. In: IEEE International Conference on Computer Vision. (2017)
3. Niklaus, S., Liu, F.: Softmax splatting for video frame interpolation. In: IEEE Conference on Computer Vision and Pattern Recognition. (2020)
4. Park, J., Ko, K., Lee, C., Kim, C.S.: Bmbc: Bilateral motion estimation with bilateral cost volume for video interpolation. In: European Conference on Computer Vision. (2020)
5. Bao, W., Lai, W.S., Ma, C., Zhang, X., Gao, Z., Yang, M.H.: Depth-aware video frame interpolation. In: IEEE Conference on Computer Vision and Pattern Recognition. (2019)
6. Niklaus, S., Mai, L., Wang, O.: Revisiting adaptive convolutions for video frame interpolation. In: IEEE Winter Conference on Applications of Computer Vision. (2021)
7. Xue, T., Chen, B., Wu, J., Wei, D., Freeman, W.T.: Video enhancement with task-oriented flow. International Journal of Computer Vision (IJCV) **127** (2019) 1106–1125
8. Bao, W., Zhang, X., Chen, L., Ding, L., Gao, Z.: High order model and dynamic filtering for frame rate up conversion. IEEE Transactions on Image Processing **PP** (2018) 1–1
9. Kuroki, Y., Nishi, T., Kobayashi, S., Oyaizu, H., Yoshimura, S.: A psychophysical study of improvements in motion-image quality by using high frame rate. Journal of The Society for Information Display - J SOC INF DISP **15** (2007)
10. Meyer, S., Cornillère, V., Djelouah, A., Schroers, C., Gross, M.H.: Deep video color propagation. In: BMVC. (2018)
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. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2018) 9000–9008
12. Wu, C., Singhal, N., Krähenbühl, P.: Video compression through image interpolation. In: European Conference on Computer Vision (ECCV). (2018)
13. Niklaus, S., Liu, F.: Context-aware synthesis for video frame interpolation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2018)
14. Bao, W., Lai, W.S., Zhang, X., Gao, Z., Yang, M.H.: Memc-net: Motion estimation and motion compensation driven neural network for video interpolation and enhancement. IEEE Transactions on Pattern Analysis and Machine Intelligence (2018)
15. Hu, P., Niklaus, S., Sclaroff, S., Saenko, K.: Many-to-many splatting for efficient video frame interpolation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2022) 3553–3562
16. Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., van der Smagt, P., Cremers, D., Brox, T.: Flownet: Learning optical flow with convolutional networks. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). (2015)
17. Sun, D., Yang, X., Liu, M.Y., Kautz, J.: Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In: Conference on Computer Vision and Pattern Recognition. (2018)

18. Gui, S., Wang, C., Chen, Q., Tao, D.: Featureflow: Robust video interpolation via structure-to-texture generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2020)
19. Liu, Z., Yeh, R., Tang, X., Liu, Y., Agarwala, A.: Video frame synthesis using deep voxel flow. In: Proceedings of International Conference on Computer Vision (ICCV). (2017)
20. Xiang, X., Tian, Y., Zhang, Y., Fu, Y., Allebach, J.P., Xu, C.: Zooming slow-mo: Fast and accurate one-stage space-time video super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2020)
21. Danier, D., Zhang, F., Bull, D.: St-mfnet: A spatio-temporal multi-flow network for frame interpolation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2022) 3521–3531
22. Lu, L., Wu, R., Lin, H., Lu, J., Jia, J.: Video frame interpolation with transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2022) 3532–3542
23. Park, J., Lee, C., Kim, C.S.: Asymmetric bilateral motion estimation for video frame interpolation. In: International Conference on Computer Vision. (2021)
24. Choi, M., Lee, S., Kim, H., Lee, K.M.: Motion-aware dynamic architecture for efficient frame interpolation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). (2021) 13839–13848
25. Teed, Z., Deng, J.: Raft: Recurrent all-pairs field transforms for optical flow. In: European Conference on Computer Vision, Springer (2020) 402–419
26. Kim, H.H., Yu, S., Tomasi, C.: Joint detection of motion boundaries and occlusions. In: British Machine Vision Conference (BMVC). (2021)
27. Yu, S., Kim, H.H., Yuan, S., Tomasi, C.: Unsupervised flow refinement near motion boundaries. In: British Machine Vision Conference (BMVC). (2022)
28. Butler, D.J., Wulff, J., Stanley, G.B., Black, M.J.: A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), ed.: European Conference on Computer Vision. Part IV, LNCS 7577, Springer-Verlag (2012) 611–625
29. Mayer, N., Ilg, E., Häusser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T.: A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In: IEEE International Conference on Computer Vision and Pattern Recognition. (2016) arXiv:1512.02134.
30. Yuan, S., Sun, X., Kim, H., Yu, S., Tomasi, C.: Optical flow training under limited label budget via active learning. In: European Conference on Computer Vision (ECCV). (2022)
31. Lee, H., Kim, T., Chung, T.y., Pak, D., Ban, Y., Lee, S.: Adacof: Adaptive collaboration of flows for video frame interpolation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2020)
32. Shi, Z., Xu, X., Liu, X., Chen, J., Yang, M.H.: Video frame interpolation transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2022) 17482–17491
33. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition (2014)
34. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2016) 770–778
35. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. (2009) 248–255

36. Yang, M., Liu, S.C., Delbruck, T.: A dynamic vision sensor with 1% temporal contrast sensitivity and in-pixel asynchronous delta modulator for event encoding. IEEE Journal of Solid-State Circuits **50** (2015) 2149–2160
37. Tulyakov, S., Gehrig, D., Georgoulis, S., Erbach, J., Gehrig, M., Li, Y., Scaramuzza, D.: Time lens: Event-based video frame interpolation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2021) 16155–16164
38. Zhang, X., Yu, L.: Unifying motion deblurring and frame interpolation with events. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2022) 17765–17774
39. Tulyakov, S., Bochicchio, A., Gehrig, D., Georgoulis, S., Li, Y., Scaramuzza, D.: Time lens++: Event-based frame interpolation with parametric non-linear flow and multi-scale fusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). (2022) 17755–17764
40. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale (2021)
41. Huang, Z., Wang, X., Wei, Y., Huang, L., Shi, H., Liu, W., Huang, T.S.: Ccnet: Criss-cross attention for semantic segmentation (2020)
42. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need (2017)
43. Jiang, S., Campbell, D., Lu, Y., Li, H., Hartley, R.: Learning to estimate hidden motions with global motion aggregation (2021)
44. Zhang, X., Wang, T., Qi, J., Lu, H., Wang, G.: Progressive attention guided recurrent network for salient object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2018)
45. Long, G., Kneip, L., Alvarez, J.M., Li, H., Zhang, X., Yu, Q.: Learning image matching by simply watching video. In Leibe, B., Matas, J., Sebe, N., Welling, M., eds.: Computer Vision – ECCV 2016, Cham, Springer International Publishing (2016) 434–450
46. Niklaus, S., Mai, L., Liu, F.: Video frame interpolation via adaptive convolution. In: IEEE Conference on Computer Vision and Pattern Recognition. (2017)
47. Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., Shlens, J.: Stand-alone self-attention in vision models (2019)
48. Galassi, A., Lippi, M., Torroni, P.: Attention in natural language processing. IEEE Transactions on Neural Networks and Learning Systems **32** (2021) 4291–4308
49. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 10012–10022
50. Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network (2016)
51. Baker, S., Roth, S., Scharstein, D., Black, M.J., Lewis, J., Szeliski, R.: A database and evaluation methodology for optical flow. In: 2007 IEEE 11th International Conference on Computer Vision. (2007) 1–8
52. Brox, T., Bregler, C., Malik, J.: Large displacement optical flow. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition, IEEE (2009) 41–48
53. Janai, J., Guney, F., Ranjan, A., Black, M., Geiger, A.: Unsupervised learning of multi-frame optical flow with occlusions. In: Proceedings of the European Conference on Computer Vision (ECCV). (2018)

54. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Advances in Neural Information Processing Systems 32. Curran Associates, Inc. (2019) 8024–8035
55. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. CoRR **abs/1412.6980** (2014)
56. Soomro, K., Zamir, A.R., Shah, M.: UCF101: A dataset of 101 human actions classes from videos in the wild. CoRR **abs/1212.0402** (2012)
57. Liu, Z., Yeh, R., Tang, X., Liu, Y., Agarwala, A.: Video frame synthesis using deep voxel flow. In: Proceedings of International Conference on Computer Vision (ICCV). (2017)
58. Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E.: Image quality assessment: from error visibility to structural similarity. IEEE Transactions on Image Processing **13** (2004) 600–612
59. Liu, Y., Liao, Y., Lin, Y.Y., Chuang, Y.Y.: Deep video frame interpolation using cyclic frame generation. In: AAAI. (2019)