

# STIF: Learning Continuous Video Representation for Space-Time Super-Resolution

Anonymous CVPR submission

Paper ID 2977

## Abstract

Videos typically record the streaming and continuous visual data as discrete consecutive frames. Since the storage cost is expensive for videos of high fidelity, most of them are stored in a relatively low resolution and frame rate. Recent works of Space-Time Video Super-Resolution (STVSR) are developed to incorporate temporal interpolation and spatial super-resolution in a unified framework. However, most of them only support a fixed up-sampling scale, which limits their flexibility and applications. In this work, instead of following the discrete representations, we propose a Space-Time Implicit Function (STIF) as a continuous representation for videos, and we show its applications for STVSR. The learned implicit neural representation can be decoded to videos of arbitrary spatial resolution and frame rate. We show that STIF achieves competitive performances with state-of-the-art STVSR methods on common up-sampling scales and significantly outperforms prior works on continuous and out-of-training-distribution scales.

## 1. Introduction

We observe the visual world in the form of streaming and continuous data. However, when we record such data with a video camera in a computer, it is often stored with limited spatial resolutions and temporal frame rates. Because of the high cost on recording and storing large time-scales of video data, oftentimes our computer vision system will need to process low-resolution and low frame rate videos. This introduces challenges in recognition systems such as video object detection [52], and we are still struggling at learning to recognize motion and actions from discrete frames [4, 12]. When presenting the video back to humans (e.g., on a TV), it is essential to visualize it in high resolution and high frame rate for user experience. How to recover the low resolution video back to high resolution in space and time becomes an important problem and the first step for many downstream applications.

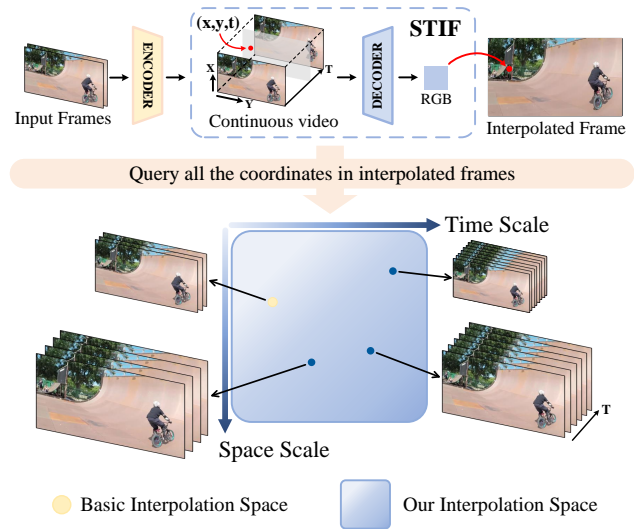


Figure 1. Our Space-Time Implicit Function (STIF) learns a continuous video representation, which maps any 3D space-time coordinate to an RGB value. This nature enables STIF to extend the latent interpolation space of STVSR from fixed space and time scales to arbitrary frame rate and spatial resolution.

Space-time video super-resolution (STVSR) approaches [15, 21, 28, 37, 38, 47, 48] are developed to increase the spatial resolution and frame rate at the same time given a low-resolution and low frame rate video as the input. Instead of performing super-resolution in space and time separately in two stages, researchers recently propose to simultaneously perform super-resolution in one stage [15, 21, 47, 48]. Intuitively, the aggregated information in time from multiple frames can reveal missing details for each frame when spatial scaling is applied, and the temporal interpolation can be more smooth and accurate given higher and richer spatial representation. The one-stage end-to-end training has shown to unify the benefits from both sides. While these results are encouraging, most approaches can only perform super-resolution to a fixed space and time scale ratio.

In this paper, instead of super resolution in a fixed scale,

we propose to learn a continuous video representation, which allows to sample and interpolate the video frames in arbitrary frame rate and spatial resolution at the same time. Our key idea is to learn a neural implicit function, which takes a space-time coordinate as input, and outputs the corresponding RGB value. Since we can sample the coordinate continuously, the video can be decoded in any spatial resolution and frame rate. Our work is inspired by recent progress on implicit functions for 3D shape representations [10, 13, 14, 26] and image representations with Local Implicit Image Functions (LIIF) using a ConvNet [7]. Different from images, where interpolation in space is based on the gradients between pixels, pixel gradients across frames with low frame rates are hard to compute. The network will need to understand the motion of the pixels and objects to perform interpolation, which is hard to model by 2D or 3D convolutions alone.

We propose a novel Space-Time Implicit Function (STIF) for continuous video representation. In the STVSR task, two low-resolution image frames are concatenated and forwarded to an encoder which generates a feature map with spatial dimensions. STIF then defines a continuous video representation over the generated feature map. It first uses a spatial implicit function module to learn a continuous spatial feature domain, from which a high-resolution image feature is sampled according to all query coordinates. Instead of using convolutional operations to perform temporal interpolation, we design the temporal implicit function module to first output a motion flow field given the high-resolution feature and the sampling time as inputs. This flow field will be applied back to warp the high-resolution feature which will be decoded to the target video frame. Since all the operations are differentiable, we can learn the motion in feature level end-to-end without any extra supervision besides the reconstruction error. To summarize, given the input frames, an encoder generates a feature map, which can be then decoded by STIF to arbitrary spatial resolution and frame rate.

In our experiments, we demonstrate that STIF can not only represent video in arbitrary space and time resolutions on the scales within the training distributions, but also extrapolate to out-of-distribution frame rates and spatial resolutions. Given the learned continuous function, instead of decoding the whole video each time, it allows the flexibility to decode only a certain region and time scale when needed. We conduct experiments with Vid4 [23], Go-Pro [29] and Adobe240 [41] datasets. We demonstrate that STIF achieves competitive performances with state-of-the-art STVSR methods on in-distribution spatial and temporal scales and significantly outperforms other methods on out-of-distribution scales.

We highlight our main contributions as follows:

- We propose a novel Space-Time Implicit Function as a

continuous video representation.

- The proposed approach allows for representing videos in arbitrary space and time resolution efficiently with one single network.
- STIF achieves out-of-distribution generalization and outperforms baselines by a large margin.

## 2. Related Work

**Implicit neural representation.** Implicit neural representations have been demonstrated as compact yet powerful continuous representations for various tasks, including 3D reconstruction [10, 13, 14, 26] and generation [5, 11, 36]. These representations typically represent signals as a neural function that maps coordinates to signed distance [34], occupancy [8, 24], or density and RGB values in a neural radiance field (NeRF [27]). Recent works also show promising results of applying this idea for modeling 2D images [1, 7, 20, 40]. Our continuous video representation is inspired by this rapidly growing field and has specific designs for videos, where a learnable flow can exploit the correspondences in video frames with inductive bias.

**Video frame interpolation.** Video frame interpolation (VFI) aims to synthesize unseen frames between the input video frames. Meyer *et al.* [25] proposed a phase-based method where information across levels of a multi-scale pyramid is combined for the synthesis of interpolated frames. Niklaus *et al.* [32, 33] introduced a series of kernel-based VFI algorithms in which they took pixel synthesis for the target frame as local convolution over input frames. Optical flow based VFI methods [2, 18, 30, 31, 49, 50] utilized optical flow prediction networks (e.g. PWC-Net [42]) to compute bidirectional flows between input frames, which served as the guidance for new frame synthesis. Additional information including occlusion masks [18, 50], depth maps [2], and cycle consistency [35] were also incorporated in the models for better performances.

**Video super-resolution.** Video super-resolution (VSR) aims at increasing the spatial resolutions of low-resolution videos. Earlier approaches [3, 43, 50] were typically built on the sliding-window framework, where they predicted optical flows between input frames and performed spatial warping for explicit feature alignment. Later on, implicit alignment started a new trend in this task [6, 17, 19, 44, 45]. For instance, TDAN [44] adopts deformable convolutions (DCNs) [9, 51] to align different input frames at feature levels. EDVR [45] further extends DCNs to a multi-scale fashion for more accurate alignment. Kelvin *et al.* introduced BasicVSR [6], in which they analyzed basic components for VSR models and suggested a bidirectional propagation scheme to maximize the gathered information from input frames.

**Space-time video super-resolution** The target of Space-time video super-resolution (STVSR) is to simultaneously

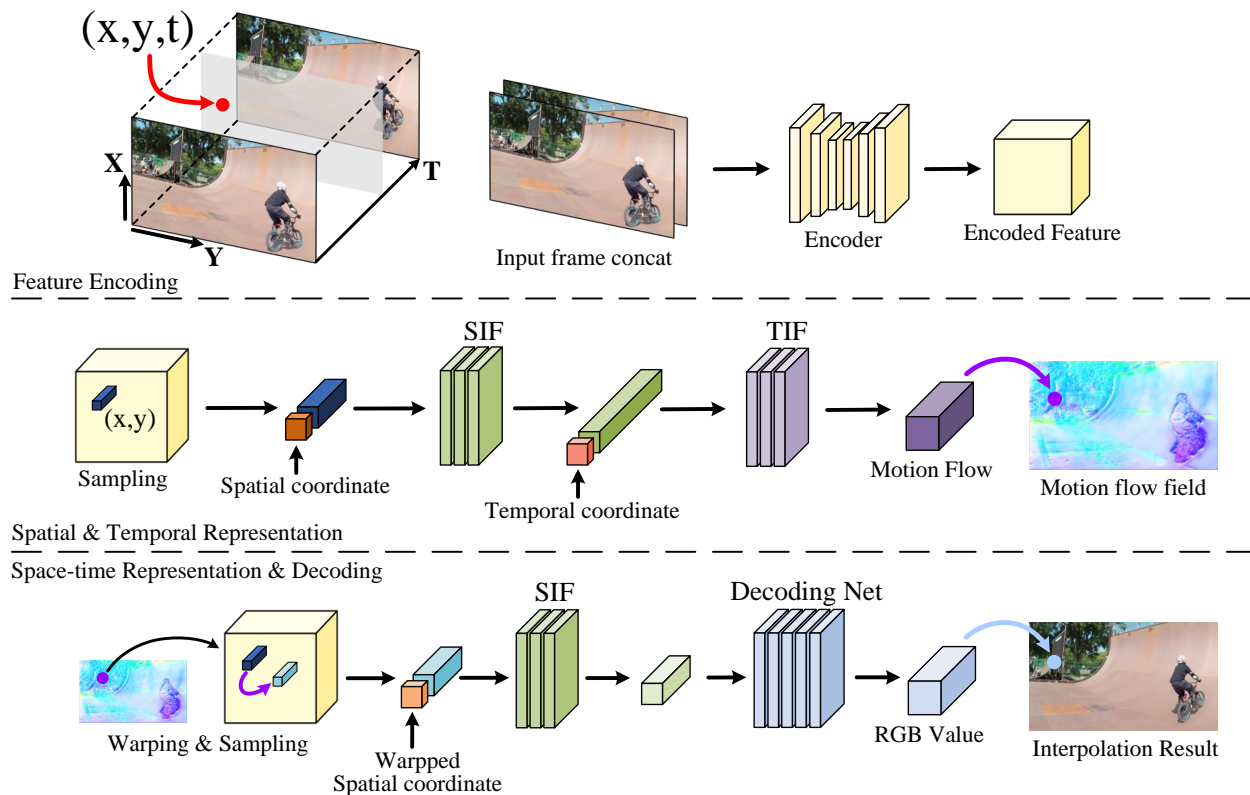


Figure 2. **A continuous video representation defined by Space-Time Implicit Function (STIF).** Two input frames are concatenated and encoded as a discrete feature map. Based on the feature, a 3D space-time coordinate is first decoded to a motion flow vector by a Spatial Implicit Function (SIF) and a Temporal Implicit Function (TIF). We then sample a new feature vector by warping according to the motion flow, and decode it as the RGB prediction of the query coordinate using a neural network. We omit the multi-scale feature aggregation part in this figure.

increase the spatial and temporal resolutions of the given low-resolution low frame rate videos. Shechtman *et al.* [38] tackled this problem by combining information from multiple input video sequences and applying a directional space-time regularization. Mudénagudi *et al.* [28] proposed a unified framework for STVSR in which videos are modeled as Markov random fields, and the maximum a posteriori estimates are taken as final solutions. Shahar *et al.* [37] introduced an effective space-time patch recurrence prior for STVSR. Recently, with the advances in deep learning, researchers started to employ powerful convolutional neural networks to address the task [15, 21, 47, 48]. Xi-ang *et al.* [47] proposed a unified neural network for synthesizing the feature of the missing frame and used a deformable ConvLSTM to align and aggregate extracted temporal information for reconstruction. STARNet [15] leveraged mutually informative relationships between time and space with the assistance of additional optical flow inputs. TMNet [48] proposed a temporal modulation block to modulate deformable convolution kernels for supporting frame interpolation at arbitrary time instances. All these STVSR

methods are designed to perform super-resolution on a specific up-sampling space scale defined before training, and some of them [15, 47] can only infer intermediate frames at pre-defined times. Therefore, the application scopes of these methods are limited. Our STIF is proposed to learn a continuous video representation that supports frame interpolation at arbitrary spatial resolution and frame rate. STIF is more flexible during the application and can be employed in more circumstances, such as non-uniform interpolation and video zoom-in in local regions.

### 3. Space-Time Implicit Function

Given a video with limited spatial resolution and frame rate, our goal is to find a continuous representation for the video. The representation interprets arbitrary space-time coordinate  $(x_s, x_t)$  into RGB values. To this end, we introduce the Space-Time Implicit Function (STIF), which produces continuous video representations of all videos. It is parameterized by multi-layer perceptrons (MLPs) and takes the form

$$s = f_{\theta}(x_s, x_t, \mathcal{V}) \quad (1)$$

where  $f_\theta$  is the proposed space-time implicit function,  $\mathcal{V}$  is the given video,  $x_s$  is the 2D spatial coordinate,  $x_t$  is the temporal coordinate, and  $s$  is the predicted RGB value. In order to learn such neural implicit representation, we propose to decouple space and time and adopt two implicit functions to represent them separately.

Figure 2 illustrates an overview of our model. Given a space-time coordinate  $(x_s, x_t)$  and the feature extracted from input frames by an encoder, a spatial implicit function (SIF) is first utilized to decode the spatial coordinate  $x_s$  and output a corresponding feature vector (Sec. 3.1). The feature is then forwarded to the temporal implicit function (TIF) for the motion flow at the query coordinate (Sec. 3.2). The flow is applied back to warp the continuous feature defined by SIF for a new feature vector (Sec. 3.3) which is finally decoded to the target RGB value (Sec. 3.4).

### 3.1. Continuous Spatial Representation

Inspired by LIIF [7], we use a neural implicit function for learning a continuous spatial representation. The implicit function converts the discrete encoded feature map to a continuous feature domain that decodes arbitrary 2D spatial coordinate into a corresponding feature vector. Specifically, the feature vectors generated by the encoder are evenly distributed in the 2D space. We sample the feature vector (the dark blue cuboid in Fig 2) nearest to the queried spatial coordinate  $x_s$ , concatenate it with the relative position information between query coordinate and feature vector, and input them into the Spatial Implicit Function (SIF)  $f_s$  to output the continuous feature at  $x_s$  (the green cuboid in Fig 2). This process could be expressed as

$$\mathcal{F}_s(x_s) = f_s(z^*, x_s - v^*) \quad (2)$$

where  $\mathcal{F}_s$  is the continuous feature domain defined by SIF,  $z^*$  is the feature vector nearest to the query coordinate  $x_s$  and  $v^*$  is the spatial coordinate of the feature vector  $z^*$ .

The main difference between LIIF and SIF is that LIIF is proposed for continuous image representation, while SIF defines a continuous feature domain, which is supposed to be further utilized for modeling temporal information in videos.

### 3.2. Continuous Temporal Representation

The proposed SIF defines a new continuous feature domain in 2D space. Our next step is to learn the continuous temporal representation and extend the feature domain from 2D space to 3D space and time, which can be achieved by decoding the temporal coordinate  $x_t$ . Directly generating the target decoded feature by a network can be fairly difficult, as the network has to learn not only the motion patterns between input frames but also the context information. Instead, we propose to learn a continuous motion flow field

for temporal representation. We introduce a temporal implicit function (TIF) to produce the motion flow.

Given a 3D space-time coordinate  $(x_s, x_t)$  and input frames  $I_0$  and  $I_1$ , the goal of TIF is to learn a mapping from the coordinate to a motion flow

$$\mathcal{M}(x_s, x_t) = f_t(x_s, x_t, I_0, I_1) \quad (3)$$

where  $\mathcal{M}$  is the continuous motion flow field and  $f_t$  is the temporal implicit function. Benefiting from the 2D continuous feature domain provided by SIF, we could replace the two input frames and the spatial coordinate  $x_s$  in input parameters of TIF with the continuous feature at  $x_s$ . Thus the equation could be written as

$$\mathcal{M}(x_s, x_t) = f_t(x_t, \mathcal{F}_s(x_s)) \quad (4)$$

where  $\mathcal{F}_s(x_s)$  is the feature domain defined in Eq 2.

In practice, we set the output of TIF as the combination of two motion flows. Based on our observation, TIF would implicitly learn bi-directional flows under such setting, which could be interpreted as correspondences between the target frame and two input frames.

### 3.3. Space-Time Continuous Representation

With two continuous representations for space and time, we aim at combining them into a unified space-time continuous representation. Starting from a space-time coordinate  $(x_s, x_t)$ , we first use SIF to predict the continuous feature at  $x_s$ . TIF is then utilized for calculating the motion flow of the query coordinate. Based on these outputs, we obtain the space-time feature by warping the continuous feature domain. The wrapped feature at  $x_s$  corresponds to the continuous feature at  $x'_s$ . The relationship between two coordinates can be written as

$$x'_s = x_s + \mathcal{M}(x_s, x_t) \quad (5)$$

where  $\mathcal{M}(x_s, x_t)$  is the motion flow vector at  $(x_s, x_t)$ .

We query this new spatial coordinate in the continuous 2D feature domain and obtain a new feature vector (the light green cuboid in Fig 2), which is treated as the feature of our continuous space-time representation at coordinate  $(x_s, x_t)$ . Accordingly, the continuous space-time feature  $\mathcal{F}_{st}$  can be formulated as

$$\mathcal{F}_{st}(x_s, x_t) = \mathcal{F}_s(x'_s) = \mathcal{F}_s(x_s + \mathcal{M}(x_s, x_t)) \quad (6)$$

### 3.4. Feature Decoding

Based on the continuous space-time representation, we can get the feature corresponding to any space-time coordinate. The final step is to decode the feature as an RGB value. A straightforward design is to take the obtained space-time feature for decoding directly. However, due to



the MLP-based network architecture, the RGB value of every predicted pixel depends on a single feature vector, leading to a limited size of the network receptive field. To alleviate the negative impact of this disadvantage, we enrich the input information of the decoding network by aggregating features of different scales. In detail, we incorporate the encoded feature as well as two input frames for decoding. Since these additional features are typically of low-resolution compared with the target resolution, we sample feature vectors corresponding to the query coordinate by bilinear interpolation. All features are then combined together for predicting the RGB output.

### 3.5. Frame synthesis

From Section 3.1 to 3.4, we focus on predicting the RGB value at a specific coordinate. To synthesize an entire frame, we need to query coordinates of all pixels in it. Given these coordinates, we can convert the continuous feature from SIF into a high-resolution feature map. We can also generate a complete motion flow field for the latent high-resolution interpolated frame. Therefore, we do not have to forward SIF twice before and after warping as in the situation of one input coordinate. Instead, we warp the whole high-resolution feature map based on the motion flow and input the warped feature into the decoding network to synthesize the target frame at one time.

## 4. Experiments

### 4.1. Experimental Setup

**Dataset.** We use Adobe240 dataset [41] as the training set, which includes 133 videos in 720P taken by hand-held cameras. We follow [48] to split these videos into the train, validation, and test subsets with 100, 16, and 17 videos. All videos are converted into image sequences for training and testing. Each sequence contains approximately 3000 frames which are treated as high-resolution frames in training. The low-resolution counterparts are then generated by imresize function in Matlab with the default setting of bicubic interpolation. We use a sliding window to select frames from the image sequences for training. The length of the sliding window is set to 9. We take the 1<sup>st</sup> and 9<sup>th</sup> frames as network inputs. The 2<sup>nd</sup> to 7<sup>th</sup> frames serve as ground-truth frames, and we randomly select three of them as the supervision of our network in every iteration. STIF is trained by two stages. In the first stage, we fixed the down-sampling space scale to  $\times 4$ . In the second stage, we randomly sample scales in a uniform distribution  $\mathcal{U}(1, 4)$ . We provide more discussion about this two-stage training strategy in Section 4.3.

Datasets including Vid4 [23], Adobe240 [41], and GoPro [29] are used for evaluation. On Vid4, we only conduct experiments on single frame interpolation of STVSR. For Adobe240 and GoPro, we evaluate on their test set. The im-

age sequences extracted from videos in the datasets are split into groups of 9-frame video clips. We feed the 1<sup>st</sup> and 9<sup>th</sup> frames down-sampled by scale  $\times 4$  in each clip into models to generate 9 high-resolution frames from 1<sup>st</sup> to 9<sup>th</sup>. We separately evaluate the average metrics of the *center* frames (i.e. the 1<sup>st</sup>, 4<sup>th</sup>, 9<sup>th</sup> frames) and all 9 output frames. They are denoted as *-Center* and *-Average* in Table 1.

**Implementation details.** We use Adam optimizer [22] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The learning rate is initialized as  $1 \times 10^{-4}$  and is decayed to  $1 \times 10^{-7}$  with a cosine annealing for every 150,000 iterations. The model is trained in a total of 600,000 iterations with batch size 24. The first training stage includes 450,000 iterations while the second stage includes 150,000 iterations. The input frames in one batch are down-sampled by the same space scale and randomly cropped into patches with size  $32 \times 32$ . We perform data augmentation by randomly rotating  $90^\circ$ ,  $180^\circ$  and  $270^\circ$ , and horizontal-flipping. We use Zooming SlowMo [47] as the encoder. For the spatial implicit function and temporal implicit function, we utilize two 3-layer SIRENs [39] with hidden dimensions of [64, 64, 256]. For the decoding network, we employ a 4-layer SIREN with hidden dimensions of [64, 64, 256, 256]. As suggested in [47, 48], we select the Charbonnier loss function for optimization.

**Evaluation.** Peak-Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) [46] are employed to evaluate model performances. We also compare the model size and inference time to measure the efficiency of models.

### 4.2. Comparison to State-of-the-arts

We compare the proposed STIF with state-of-the-art two-stage and one-stage STVSR methods. For two-stage methods, we employ SuperSloMo [18], QVI [49], and DAIN [2] for video frame interpolation (VFI); Bicubic Interpolation, EDVR [45], and BasicVSR [6] for video super-resolution (VSR). For one-stage methods, we compare STIF with recently developed Zooming SlowMo [47] and TMNet [48]. To perform fair comparisons, we train the three VFI methods and Zooming SlowMo from scratch on Adobe240 dataset. For TMNet, as mentioned in the original paper that a two-stage training scheme is needed for convergence, we pre-train the model on Vimeo90K [50] dataset and fine-tune it on Adobe240 dataset [41]. Therefore, TMNet is trained on more data compared with other methods, which may lead to some advantages in the comparison. To compare with Zooming SlowMo that only supports fixed frame interpolation, we train a new version of STIF named STIF-*fixed* of which the interpolation time is fixed to 0.5.

**Quantitative results.** We present in-distribution quantitative comparisons between STIF and other STVSR methods in Table 1. On single frame interpolation of STVSR including Vid4, GoPro-*Center*, and Adobe-*Center*, STIF-*Fixed* achieves competitive performance compared with

Table 1. **Quantitative comparison on benchmark datasets** including Vid4 [23], GoPro [29] and Adobe240 [41]. The best three results are highlighted in **red**, **blue**, and **bold**. We omit the results of Zooming SlowMo and STIF-Fixed on GoPro-Average and Adobe240-Average as the two models are trained for synthesizing frames only at fixed times.

VFI Method	SR Method	Vid4		GoPro-Center		GoPro-Average		Adobe-Center		Adobe-Average		Parameters (Million)
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	
SuperSloMo [18]	Bicubic	22.42	0.5645	27.04	0.7937	26.06	0.7720	26.09	0.7435	25.29	0.7279	19.8
SuperSloMo [18]	EDVR [45]	23.01	0.6136	28.24	0.8322	26.30	0.7960	27.25	0.7972	25.95	0.7682	19.8+20.7
SuperSloMo [18]	BasicVSR [6]	23.17	0.6159	28.23	0.8308	26.36	<b>0.7977</b>	27.28	0.7961	25.94	0.7679	19.8+6.3
QVI [18]	Bicubic	22.11	0.5498	26.50	0.7791	25.41	0.7554	25.57	0.7324	24.72	0.7114	29.2
QVI [18]	EDVR [45]	23.60	0.6471	27.43	0.8081	25.55	0.7739	26.40	0.7692	25.09	0.7406	29.2+20.7
QVI [18]	BasicVSR [6]	23.15	0.6428	27.44	0.8070	26.27	0.7955	26.43	0.7682	25.20	0.7421	29.2+6.3
DAIN [2]	Bicubic	22.57	0.5732	26.92	0.7911	26.11	0.7740	26.01	0.7461	25.40	0.7321	24.0
DAIN [2]	EDVR [45]	23.48	0.6547	28.01	0.8239	26.37	0.7964	27.06	0.7895	26.01	0.7703	24.0+20.7
DAIN [2]	BasicVSR [6]	23.43	0.6514	28.00	0.8227	<b>26.46</b>	0.7966	27.07	0.7890	<b>26.23</b>	<b>0.7725</b>	24.0+6.3
Zooming SlowMo [47]		<b>25.72</b>	<b>0.7717</b>	<b>30.69</b>	<b>0.8847</b>	-	-	<b>30.26</b>	<b>0.8821</b>	-	-	<b>11.10</b>
TMNet [48]		<b>25.96</b>	<b>0.7803</b>	30.14	0.8692	<b>28.83</b>	<b>0.8514</b>	29.41	0.8524	<b>28.30</b>	<b>0.8354</b>	<b>12.26</b>
STIF-fixed		<b>25.78</b>	<b>0.7730</b>	<b>30.73</b>	<b>0.8850</b>	-	-	<b>30.21</b>	<b>0.8805</b>	-	-	<b>11.31</b>
STIF		25.61	0.7709	<b>30.26</b>	<b>0.8792</b>	<b>29.41</b>	<b>0.8669</b>	<b>29.92</b>	<b>0.8746</b>	<b>29.27</b>	<b>0.8651</b>	<b>11.31</b>

Table 2. **Quantitative comparison for out-of-distribution scales** on GoPro dataset. Model performances are evaluated by PSNR and SSIM. Some results of TMNet are bolded as it does not support generalizing to out-of-training-distribution space scales.

Time Scale	Space Scale	SuperSloMo [18] + LIIF [7]	DAIN [2] + LIIF [7]	TMNet [48]	STIF
×6	×4	26.70 / 0.7988	26.71 / 0.7998	30.49 / 0.8861	<b>30.78 / 0.8954</b>
×6	×6	23.47 / 0.6931	23.36 / 0.6902	-	<b>25.56 / 0.7671</b>
×6	×12	21.92 / 0.6495	22.01 / 0.6499	-	<b>24.02 / 0.6900</b>
×12	×4	25.07 / 0.7491	25.14 / 0.7497	26.38 / 0.7931	<b>27.32 / 0.8141</b>
×12	×6	22.91 / 0.6783	22.92 / 0.6785	-	<b>24.68 / 0.7358</b>
×12	×12	21.61 / 0.6457	21.78 / 0.6473	-	<b>23.70 / 0.6830</b>
×16	×4	24.42 / 0.7296	24.20 / 0.7244	24.72 / 0.7526	<b>25.81 / 0.7739</b>
×16	×6	23.28 / 0.6883	22.80 / 0.6722	-	<b>23.86 / 0.7123</b>
×16	×12	21.80 / 0.6481	22.22 / 0.6420	-	<b>22.88 / 0.6659</b>

Table 3. **Quantitative comparison of out-of-distribution performance between STIF and the baseline Zooming SloMo model [47]**. Evaluated on GOPRO dataset. -×A×B refers to A up-sampling space scale and B up-sampling time scale.

Method	GoPro - ×4×2		GoPro - ×16×4	
	PSNR	SSIM	PSNR	SSIM
Zooming SloMo	30.69	0.8847	23.38	0.6708
STIF	30.26	0.8792	23.45	0.6710

other state-of-the-art models, while the performance of STIF slightly suffers. We attribute this observation to the difference of training targets between STIF and STIF-Fixed. The training settings of STIF-Fixed aim for synthesizing frames at pre-defined times. Therefore, it only learns fixed patterns between input frames instead of learning a continuous representation as STIF does, leading to advantages in performances. On Vid4, TMNet performs the best, and we assume this is because TMNet is trained with more data as we noted in Section 4.2. For multiple frame interpolation of STVSR including GoPro-Average and Adobe-Average, STIF achieves the best performance, which indicates that

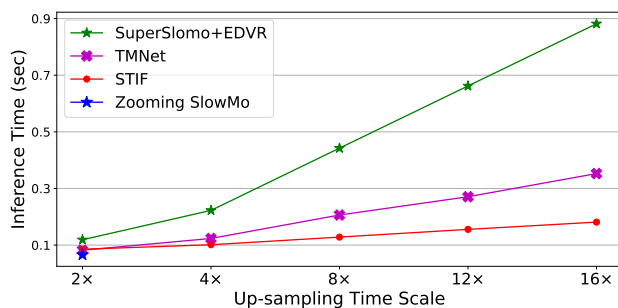


Figure 3. **Inference time of STVSR methods on different up-sampling time scales.** Note that We only select the most efficient two-stage method (SuperSloMo + EDVR) for comparison.

the proposed implicit neural representation provides advances on modeling the temporal information in videos.

In Table 2, we present comparisons of STVSR methods on out-of-distribution space and time scales. For two-stage STVSR methods, we select SuperSloMo and DAIN as VFI methods, and LIIF as the SR method since it can perform super-resolution on arbitrary up-sampling scales. We also take TMNet into the comparison as it could general-



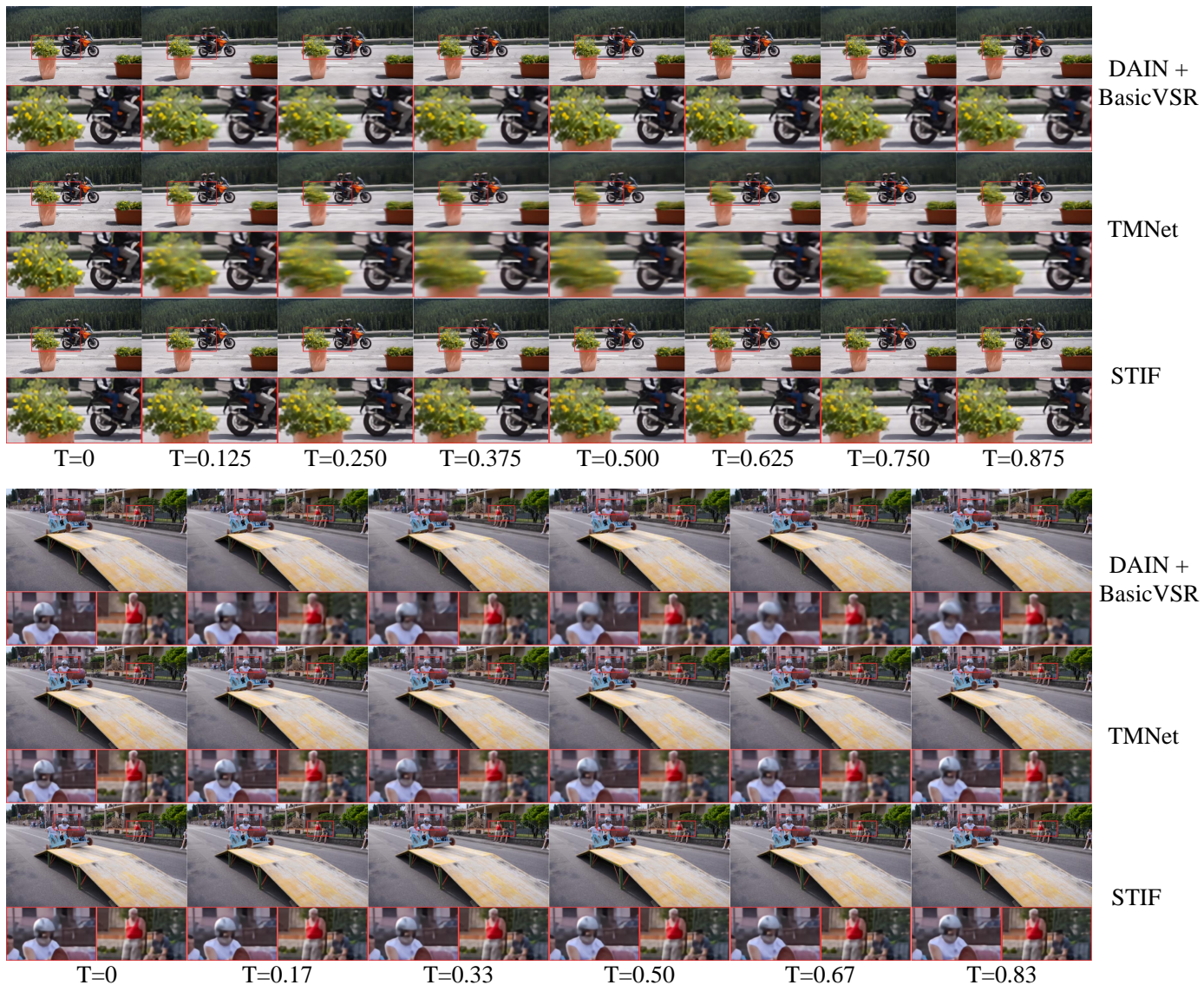


Figure 4. **Qualitative comparisons of different STVSR methods on arbitrary frame interpolation.** The interpolation times of the first example are in the training distribution and the times of the second example are out-of-distribution. Best zoom in for better visualization.

ize on time scales. We produce experiments on GoPro [29] dataset. We observe that STIF outperforms other methods by a large margin, demonstrating the advantage of our continuous video representation in out-of-distribution generalization. In addition, we further compare STIF with Zooming SlowMo (the encoder for STIF) in out-of-distribution scales. As Zooming SlowMo only supports interpolating fixed frames, we apply the model twice to achieve out-of-distribution inferences. In Table 3, we observe that while Zooming SlowMo performs slightly better on single frame interpolation ( $\times 4 \times 2$ ), STIF achieves better performance in out-of-distribution testing ( $\times 16 \times 4$ ).

We compare the inference time of STVSR methods in Figure 3. We observe that the efficiency of STIF is close to Zooming-SlowMo and TMNet at up-sampling time scale

$\times 2$ , and STIF inferences faster than other models on multi-frame interpolation. We attribute this feature to the design of STIF, where all the latent frames between two input frames can be synthesized by MLPs after encoding.

**Qualitative Results** We demonstrate a qualitative comparison in Figure 4. We compare STIF with two STVSR methods, DAIN + BasicVSR and TMNet. The selected temporal coordinates of the first sample are in the training distribution, while the coordinates of the second sample are out-of-distribution. We find that the performance of DAIN + BasicVSR degrades in out-of-distribution circumstances (see the rider’s head in the second sample). TMNet fails to recover objects with large motion between two input frames (see the flowers in the first sample). The performance of STIF is steady across both in-distribution and out-of-

Table 4. **Ablation study on architecture designs of STIF.** Evaluated on GOPRO and Adobe240 dataset. -f/m refers to removing flow correspondence and multi-scale feature aggregation. -s refers to decoding both time and space by a single implicit function.

Architecture Design	GoPro-Center		GoPro-Average		Adobe-Center		Adobe-Average	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
STIF	30.26	0.8792	29.41	0.8669	29.92	0.8746	29.27	0.8651
STIF (-f)	29.63	0.8719	28.76	0.8614	29.19	0.8641	28.50	0.8569
STIF (-m)	29.99	0.8751	29.28	0.8655	29.68	0.8690	29.04	0.8606
STIF (-s)	29.86	0.8741	29.20	0.8654	29.42	0.8678	28.95	0.8613

Table 5. **Ablation study on STIF trained with different data settings.** Evaluated on GOPRO-Average.  $\times 4$  refers to fixing the down-sampling space scale to  $\times 4$  throughout the training. *-continuous* refers to training STIF by continuous space scales from scratch.

Training Settings	Space $\times 2$		Space $\times 3$		Space $\times 4$		Space $\times 6$		Space $\times 12$	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
STIF	29.61	0.8734	29.14	0.8685	29.41	0.8669	25.40	0.7590	24.11	0.6913
STIF ( $\times 4$ )	28.25	0.8490	28.62	0.8626	29.50	0.8696	25.24	0.7567	23.82	0.6857
STIF ( <i>-continuous</i> )	27.46	0.8268	28.35	0.8507	28.82	0.8541	25.10	0.7533	23.62	0.6801

distribution temporal coordinates, indicating that learning continuous video representations helps to improve model generalization in STVSR task.

### 4.3. Ablation Study

**Motion Flow Field.** Motion flow is one critical component of STIF. Previous video interpolation methods [16, 18] have already demonstrated that such a learnable flow helps to interpolate frames with sharp edges and clear details. We propose that the motion flow field brings two main advantages. First, the flow field could capture non-local information and temporal contexts of large motions. Second, we explicitly apply spatial warping on features, which works as an inductive bias for the training. In Table 4 between STIF and STIF (-f), we show that the performance degrades when the motion flow is not incorporated in STIF.

**STIF trained with different data settings.** In Table 5, we compare the performances of STIF trained on different data settings. As noted before, STIF follows a two-stage training strategy: fixed down-sampling space scale for the first stage and continuous space scales sampled from a uniform distribution for the second stage. STIF- $\times 4$  indicates that the space scale is fixed to  $\times 4$  throughout the training of STIF. STIF-*continuous* represents STIF trained with continuous down-sampling space scales from scratch. We find that the performance suffers a significant drop when we train STIF only on continuous scales. We hypothesize this is because the network needs to learn spatial and temporal representations at the same time, and it becomes extremely difficult to learn such temporal representation when the scale of spatial features keeps varying. Besides, we observe that training STIF with a fixed space scale achieves slightly better performance for that specific scale. However, its generalization performance is competed by STIF trained by two stages, which is demonstrated by the comparisons between

STIF and STIF ( $\times 4$ ) on space scales other than  $\times 4$ .

**Other design choices.** We provide more ablation studies in Table 4. By comparing STIF with STIF (-m), we find that the proposed multi-scale feature aggregation contributes to performance improvement. We also try to combine SIF and TIF into a single implicit function, that is, we use one implicit function to generate the continuous motion flow and apply spatial warping only on the encoded feature and input frames. The results between STIF and STIF (-s) indicate that using two implicit functions for representing space and time outperforms only one implicit function for them all.

## 5. Discussion

**Conclusion.** In this paper, we present a Space-Time Implicit Function (STIF) for continuous video representation. STIF can represent videos in arbitrary spatial and temporal resolution, which brings natural advantages for solving space-time video super-resolution (STVSR) tasks. Extensive experiments show that STIF performs competitively with state-of-the-art STVSR methods on common up-sampling scales and outperforms prior works by a large margin on out-of-distribution scales.

**Limitations and Future Work.** We observe that there exist few cases for which STIF does not perform very well. These cases typically need to handle very large motions, which is still an open challenge for video interpolation.

**Ethical Concerns.** STIF interpolates frames based on learned statistics of the training dataset. Thus, it would reflect biases in those data, including ones with negative societal impacts. STIF may generate inexistent or fake contents. These issues warrant further consideration before processing videos by STIF. As researchers, we are committed to against misconduct behaviors and pursue research that is to the benefit of society.



## References

- [1] Ivan Anokhin, Kirill Demochkin, Taras Khakhulin, Gleb Sterkin, Victor Lempitsky, and Denis Korzhenkov. Image generators with conditionally-independent pixel synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14278–14287, 2021. 2
- [2] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3703–3712, 2019. 2, 5, 6
- [3] Jose Caballero, Christian Ledig, Andrew Aitken, Alejandro Acosta, Johannes Totz, Zehan Wang, and Wenzhe Shi. Real-time video super-resolution with spatio-temporal networks and motion compensation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4778–4787, 2017. 2
- [4] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 1
- [5] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5799–5809, 2021. 2
- [6] Kelvin CK Chan, Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Basicvsr: The search for essential components in video super-resolution and beyond. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4947–4956, 2021. 2, 5, 6
- [7] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8628–8638, 2021. 2, 4, 6
- [8] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 764–773, 2017. 2
- [10] Boyang Deng, John P Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Nasa neural articulated shape approximation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 612–628. Springer, 2020. 2
- [11] Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W Taylor, and Joshua M Susskind. Unconstrained scene generation with locally conditioned radiance fields. *arXiv preprint arXiv:2104.00670*, 2021. 2
- [12] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 1
- [13] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4857–4866, 2020. 2
- [14] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7154–7164, 2019. 2
- [15] Muhammad Haris, Greg Shakhnarovich, and Norimichi Ukita. Space-time-aware multi-resolution video enhancement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2859–2868, 2020. 1, 3
- [16] Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. Rife: Real-time intermediate flow estimation for video frame interpolation. *arXiv preprint arXiv:2011.06294*, 2020. 8
- [17] Takashi Isobe, Xu Jia, Shuhang Gu, Songjiang Li, Shengjin Wang, and Qi Tian. Video super-resolution with recurrent structure-detail network. In *European Conference on Computer Vision*, pages 645–660. Springer, 2020. 2
- [18] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9000–9008, 2018. 2, 5, 6, 8
- [19] Younghyun Jo, Seoung Wug Oh, Jaeyeon Kang, and Seon Joo Kim. Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3224–3232, 2018. 2
- [20] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *arXiv preprint arXiv:2106.12423*, 2021. 2
- [21] Soo Ye Kim, Jihyong Oh, and Munchurl Kim. Fsr: deep joint frame interpolation and super-resolution with a multi-scale temporal loss. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11278–11286, 2020. 1, 3
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [23] Ce Liu and Deqing Sun. A bayesian approach to adaptive video super resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 209–216. IEEE, 2011. 2, 5, 6
- [24] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks:

- Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [25] Simone Meyer, Oliver Wang, Henning Zimmer, Max Grosse, and Alexander Sorkine-Hornung. Phase-based frame interpolation for video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1410–1418, 2015. 2
- [26] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4743–4752, 2019. 2
- [27] 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
- [28] Uma Mudenagudi, Subhashis Banerjee, and Prem Kumar Kalra. Space-time super-resolution using graph-cut optimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):995–1008, 2010. 1, 3
- [29] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3883–3891, 2017. 2, 5, 6, 7
- [30] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1701–1710, 2018. 2
- [31] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5437–5446, 2020. 2
- [32] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 670–679, 2017. 2
- [33] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 261–270, 2017. 2
- [34] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [35] Fitsum A Reda, Deqing Sun, Aysegul Dundar, Mohammad Shoeybi, Guilin Liu, Kevin J Shih, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Unsupervised video interpolation using cycle consistency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 892–900, 2019. 2
- [36] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *arXiv preprint arXiv:2007.02442*, 2020. 2
- [37] Oded Shahrar, Alon Faktor, and Michal Irani. *Space-time super-resolution from a single video*. IEEE, 2011. 1, 3
- [38] Eli Shechtman, Yaron Caspi, and Michal Irani. Increasing space-time resolution in video. In *European Conference on Computer Vision*, pages 753–768. Springer, 2002. 1, 3
- [39] 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, 2020. 5
- [40] Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial generation of continuous images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10753–10764, 2021. 2
- [41] Shuo Chen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1279–1288, 2017. 2, 5, 6
- [42] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8934–8943, 2018. 2
- [43] Xin Tao, Hongyun Gao, Renjie Liao, Jue Wang, and Ji-aya Jia. Detail-revealing deep video super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4472–4480, 2017. 2
- [44] Yapeng Tian, Yulun Zhang, Yun Fu, and Chenliang Xu. Tdan: Temporally-deformable alignment network for video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3360–3369, 2020. 2
- [45] Xintao Wang, Kelvin CK Chan, Ke Yu, Chao Dong, and Chen Change Loy. Edvr: Video restoration with enhanced deformable convolutional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 2, 5, 6
- [46] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 5
- [47] Xiaoyu Xiang, Yapeng Tian, Yulun Zhang, Yun Fu, Jan P Allebach, and Chenliang Xu. 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*, pages 3370–3379, 2020. 1, 3, 5, 6
- [48] Gang Xu, Jun Xu, Zhen Li, Liang Wang, Xing Sun, and Ming-Ming Cheng. Temporal modulation network for controllable space-time video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6388–6397, 2021. 1, 3, 5, 6
- [49] Xiangyu Xu, Li Siyao, Wenxiu Sun, Qian Yin, and Ming-Hsuan Yang. Quadratic video interpolation. *Advances in Neural Information Processing Systems*, 32:1647–1656, 2019. 2, 5

1080		1134
1081		1135
1082		1136
1083		1137
1084	[50] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and	1138
1085	William T Freeman. Video enhancement with task-	1139
1086	oriented flow. <i>International Journal of Computer Vision</i> ,	1140
1087	127(8):1106–1125, 2019. 2, 5	1141
1088	[51] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. De-	1142
1089	formable convnets v2: More deformable, better results. In	1143
1090	<i>Proceedings of the IEEE/CVF Conference on Computer Vi-</i>	1144
1091	<i>sion and Pattern Recognition</i> , pages 9308–9316, 2019. 2	1145
1092	[52] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen	1146
1093	Wei. Flow-guided feature aggregation for video object detec-	1147
1094	tion. In <i>Proceedings of the IEEE International Conference</i>	1148
1095	<i>on Computer Vision</i> , pages 408–417, 2017. 1	1149
1096		1150
1097		1151
1098		1152
1099		1153
1100		1154
1101		1155
1102		1156
1103		1157
1104		1158
1105		1159
1106		1160
1107		1161
1108		1162
1109		1163
1110		1164
1111		1165
1112		1166
1113		1167
1114		1168
1115		1169
1116		1170
1117		1171
1118		1172
1119		1173
1120		1174
1121		1175
1122		1176
1123		1177
1124		1178
1125		1179
1126		1180
1127		1181
1128		1182
1129		1183
1130		1184
1131		1185
1132		1186
1133		1187